

Polysp: a polysystemic, phonetically-rich approach to speech understanding

Sarah Hawkins & Rachel Smith

ABSTRACT

We outline an approach to speech understanding, Polysp (for POLYsystemic SPEech understanding) that combines a richly-structured, polysystemic linguistic model derived from Firthian prosodic analysis and declarative phonology, with psychological and neuropsychological approaches to the organization of sensory experience into knowledge. We propose that the type of approach exemplified by Polysp promises a fruitful way of conceptualising how meaning is understood from spoken utterances, partly by ascribing an important role to all kinds of systematic fine phonetic detail available in the physical speech signal and by rejecting assumptions that the physical signal is analysed as early as possible into abstract linguistic units. Polysp provides a framework by which episodic multimodal sensory experience of speech can be simultaneously processed into different types of linguistic and non-linguistic knowledge at a variety of levels of abstraction, with the emphasis always on understanding meaning in order to interact with another person rather than on building a complete description of a given utterance at successive, obligatory stages of formal linguistic analysis. We discuss phonetic data consistent with these views.

1. Introduction

This paper explores the contribution of phonetic knowledge to how we understand words, and some implications for what makes a plausible model of spoken word understanding. We show that certain types of fine phonetic detail systematically reflect not just the phonemic content but the wider phonological and grammatical structure of the message; that while some systematic differences in phonetic fine detail are relatively localised in the speech signal, others stretch over several syllables; and that both types can make speech easier to understand. We make the case that one consequence of neglecting fine phonetic detail in models of spoken word recognition and understanding is that other processes and stages of analysis may be given inappropriate emphasis, and that this has happened in models which adopt the convenient fiction that the phoneme is the basic input unit to the lexicon. In consequence, no current phonetic or psycholinguistic theory accounts satisfactorily for how normal connected speech is understood.

We turn then to discuss central properties needed to model the process of speech understanding. First, we propose a model, derived from Firthian prosodic analysis and declarative phonology, that provides a clear linguistic-phonetic structure onto which incoming sensory speech information might be mapped. Then we set out critical characteristics of human listeners that may define the way they make sense of richly informative sensory signals. Amongst these, we emphasize various forms of learning, and some current neuropsychological views about the nature of memory and the organisation of mental categories, both linguistic and non-linguistic. We suggest that phonetic categories are like all other mental categories: self-organising (emerging from the distribution of incoming sensory information in combination with pre-existing relevant knowledge), multimodal and distributed within the brain, dynamic, and context-sensitive (or relational) and therefore plastic, or labile. These two threads, linguistic and neuropsychological/neurophysiological, together form the theoretical approach we name Polysp. We conclude by discussing central properties needed to model how normal speech is understood, identifying at each stage an existing (usually computational) model which includes properties that we consider promising.

2. *The dominance of the phoneme in models of speech perception and spoken word recognition*

2.1. *Overview*

It is well known that there are interdependencies between grammatical, prosodic and segmental parameters in speech. Yet the linguistic concept of the phoneme as the basic unit of sound contrast has dominated the thinking of phoneticians and psychologists over the last 50 years to such an extent that it is central to most theories. Either the phoneme (or an ill-defined 'phonetic segment' that is treated functionally as a phoneme) is taken as axiomatic, or it is effectively given a crucial role even while acknowledging that it has limitations and that other units may be more fundamental. Thus, phonetic models of speech perception usually focus on how phonemes are distinguished from one another, while psychological models of spoken word recognition use the phoneme as a crucial unit in the early stages of the process of recognition, and indeed usually as the input unit. One consequence is that phoneticians, speech scientists and psycholinguists all tend to assume that the sensory signal is transformed into an abstract form early in the process of word recognition, and certainly before lexical access.

Phoneme labels are valuable for certain types of linguistic and phonetic description, but we question their value in modelling speech understanding under normal conditions, as will become clear below. However, the purpose of this section is not to reiterate arguments against the phoneme as a necessary unit of speech perception (see for example Marslen-Wilson & Warren 1994; Hawkins 1995; Coleman 1998; Nguyen & Hawkins 1999; Warren 1999:169ff; Hawkins & Nguyen in press) but to argue that over-emphasis on phonemes has deflected attention from the presence of other types of linguistic information in the speech signal and, in consequence, distorted the relative importance of various processes in models of how speech is understood.

2.2. *Information conveyed by unstructured strings of phonemes or allophones*

The phoneme is an abstraction without physical reality. Its close relative, the allophone, can be loosely interpreted as having physical reality, but as is well known, a phoneme string cannot be uniquely related to a string of lexical items, and allophones cannot be related to the right phonemes, nor even to the right phoneme slots, independently of the linguistic structure in which they occur. For example, the phoneme string /katsaɪz/ signals *cat's eyes* (or *cats' eyes*) and *cat size* but in natural speech is distinguished by a number of durational and spectral differences predictable from the linguistic structure. The /s/ will be relatively longer when it is the onset of the second syllable and, depending on the accent, there can be other differences such as in the quality and degree of diphthongization of the second nucleus (in *eyes/size*). These allophonic differences presumably help the listener find the right meaning.

Similarly, in Standard Southern British English (SSBE), *Carter Knight* and *car tonight* have identical phoneme sequences, and in this case they can be pronounced very similarly. However, in at least one London accent, they are usually differentiated by the pattern of glottal stops: *Carter Knight* would be [kɑʔənaiʔ] whereas *car tonight* would be [kətənaiʔ], because in this accent glottal stops cannot substitute for [t] word-initially. All the glottal stops in this pair of phrases happen to be allophones of /t/, but glottal stops do not always signal /t/ in this accent. For instance, the phrase *hand it over*, said as [handɪtəʊvə] or [handɪʔəʊvə] in SSBE, also has (at least) two glottal stops in the London accent, [ʔandɪʔəʊvə]. But while the second one still signals /t/, the first is an allophone of neither /t/ nor /h/, since this accent lacks /h/; instead, it signifies (optionally) that the vowel is utterance-initial. In other words, allophones like glottal stop do not map simply onto phonemes, so there is still room for confusion unless a larger structure than the phoneme-sized segment is used to access meaning.

Similar arguments and illustrations can be made for the way an utterance's phonetic structure can indicate its grammatical structure. For example, in English and many other languages, the phoneme structure of function words is much more limited than that of content words, and function words are subject to rather different connected speech processes with surrounding words. For example, whereas many English function words begin with /ð/, no content words do, and word-initial

/ð/ conditions distinctive connected speech processes after word-final /n/ that are not found for other /n/-fricative sequences. Thus *ban that*, phonemically /bāðat/, is commonly pronounced [baŋːat] in SSBE and other accents of English, but other /n/-fricative sequences must retain their fricative, although they may lose or almost lose the /n/ completely as long as the preceding vowel is nasalized. So *ban thatch* can be [baŋθatʃ] or [bāⁿθatʃ] or similar, and *ban zips* can be [banzɪps], [bāⁿzɪps] or even [bāzɪps]; but neither would be understood if, following the rules for /nð/ in function words, they were [baŋːatʃ] and [banːɪps] respectively.

Note that these transcriptions oversimplify, and in particular do not do justice to the subtleties of timing, vowel quality, and other variation that tend to co-occur with the more obvious allophonic differences, creating a coherent sound that systematically reflects much of the linguistic structure of the intended utterance. Examples of such subtle systematic variation can be found in Section 3.2 below, while Kelly & Local (1989), Manuel *et al.* (1992), Manuel (1995), and Ogden (1999) offer more detailed treatments of particular examples.

What these examples do illustrate is that, by itself, an allophone string is effectively about as abstract and uninformative as a phoneme string unless given an explicit context which, for connected speech, includes grammar as well as syllable and word structure. That being so, although replacing phoneme strings with allophone strings appears to add desirable phonetic detail to the input to models of word recognition, it is unlikely to solve the problems introduced by the use of phoneme strings. What is needed is an input to the lexical identification process that preserves all the linguistic information inherent in the speech signal, and an output that is also sensitive to the linguistic structure—lexical and grammatical—so that speech can be understood correctly as soon as it is heard. Although experiments show that both meanings of homophonic single words can be simultaneously activated even when an appropriate context is provided (Swinney 1979), it seems unlikely that this will standardly happen in a normal spoken exchange. For example, it is unlikely that the apparent phoneme structure of [baŋːat] means it is understood first as *ban Nat* and only later corrected to *ban that*, or even that both meanings are simultaneously activated, for the sorts of phrases we are talking about are not homophonic if the fine phonetic detail is attended to. In other words, we suggest that systematic differences in fine phonetic detail provide one sort of disambiguating context that constrains which meanings are accessed, much as has been demonstrated in a number of experiments for other types of context (see Simpson 1994 for a review). Moreover, because fine phonetic detail often signals grammatical structure, it co-occurs with non-phonetic types of disambiguating context, and the two together can presumably provide even stronger constraints on which meaning is accessed. Information of this type has not normally been considered the domain of standard phonetic theories of speech perception, nor of the phonetic input to most psycholinguistic models of lexical access.

2.3. *Consequences for phonetic and psycholinguistic models of a focus on phonemes*

The focus of early phonetic and word recognition research on abstract, idealised and unstructured linguistic units like phonemes as the primary unit of perception is understandable and indeed defensible, but it has had at least two biasing consequences on the development of theory. It has encouraged (1) ‘short-domainism’ and (2) the introduction of separate and often arbitrary processes to explain how speech can be understood despite the impoverished information provided by unstructured phoneme strings.

‘Short-domainism’ is exemplified by much phonetic speech perception research, which has typically focussed on simple sound contrasts in highly controlled phonetic environments, if only to keep the size of investigations manageable. Experiments exploring the perceptual correlates of consonants, for example, often use only one vowel, while those on vowels normally restrict the consonantal context. There is surprisingly little literature on the perception (or production) of unstressed syllables, and relatively few experiments examine the perception of phoneme identity in more complex environments—when they do, results are often quite different (e.g. Kewley-Port & Zheng 1999). Most research on prosody (apart from stress) is conducted in isolation from that on segmental identity, and results of experiments that examine influences on segmental perception of domains of more than one or two syllables have tended not to be integrated into the most influential theories of perception. They may even be seen as part of some other process, such as vocal-tract

normalisation (e.g. Ladefoged & Broadbent 1957) or how we recognise words as opposed to phonemes (e.g. Warren 1970, 1984).

Thus, word recognition and sound recognition (of features or phonemes) are often axiomatically distinct processes in many of the most influential theories arising from the various branches of speech science. For example, the Motor Theory (Liberman & Mattingly 1985) distinguishes between acoustic-phonetic trading relations, phonetic context effects, and 'higher-order' effects of lexical knowledge (cf. Repp 1982; pers. comm. 1989), although they can alternatively be seen as due to the same underlying process of using multiple cues to make complex decisions, though often operating on different temporal domains. As the above discussion of information conveyed by allophones suggests, we maintain that sound recognition and word recognition should not be axiomatically separate, because knowing about words and larger structures makes it easier to interpret allophonic detail and *vice versa*. (This position provides a viable context for explaining the classical intelligibility experiment by Pickett & Pollack 1963). It will become clear later in this paper that, for us, the appeal of direct realist theories of speech perception (Fowler 1986; Best 1994, 1995) would gain significantly if the focus were on identifying all units of linguistic structure from all aspects of the highly complex vocal tract behaviour for speech, rather than on recognition of individual phonemes from particular details of movement.

The second problem of a principal focus on phonemes is exemplified by much of the debate on psychological models of spoken word recognition over the last 20 or more years, for they typically introduce a linguistically impoverished phoneme-like input to the lexicon. Even theoreticians who acknowledge the difficulties of identifying phonetic units of speech perception nevertheless usually do most of their actual work with discrete, abstract units corresponding to phonemes (or to distinctive features that have clear temporal boundaries and are grouped into units that correspond to phonemes), if only because to do otherwise is so complicated that it would effectively mean abandoning their own work and moving into another discipline.

Unfortunately, emphasis on an unstructured phoneme string as the basis for lexical activation or search brings with it the implicit assumptions that the input is structured into discrete units capable of distinguishing lexical meaning but is linguistically incomplete in other respects. Theoretical debate thus focuses on such issues as whether various hypothesized processes are autonomous or interactive (that is, on whether an earlier process is affected by feedback from later processes), and on whether inhibitory as well as excitatory processes are necessary to produce the desired output, as theorists struggle to build in enough supplementary sources of information to compensate for the impoverished information available from an unstructured phoneme string or its like. Comparison of the architecture and performance of computational models will not necessarily offer a resolution to these issues, since the same effects can result from quite different models, depending on the detailed assumptions and definitions adopted (e.g. Norris 1992). However, for a recent example of this type of debate, see Norris' *et al.* (2000) discussion of their model Merge, and the accompanying commentaries.

Merge, a development of Shortlist (Norris 1994), provides an interesting example of how early abstraction of linear strings of phonetic units demands invocation of particular, equally abstract processes to compensate for this simplification. Merge is presented as a model of phonemic decision making. It is unusual in having phonemes represented in two places and subject to quite different processes. The motivation for this is to distinguish initial sensory information (called prelexical information) from decisions about phoneme identity that are made from all the available information, lexical and sensory/prelexical. Sensory information is coded in terms of probabilities, with the probability for each prelexical phoneme being independent of that of the others, while the phoneme decision stage is categorical and subject to influence from other phoneme decisions.

The properties of the prelexical phoneme stage are introduced to preserve information about the sensory input so that a wrong decision about one phoneme need not jeopardise word recognition—in other words, to allow error correction. But neither simple feature bundles nor phonemes reflect all the systematic linguistic information available in the speech signal that could facilitate error correction or even avoid the need for it in the first place. Additionally, to account for certain experimental findings such as 'compensation for coarticulation', the prelexical stage includes abstract knowledge of transitional phoneme probabilities. Other types of linguistic knowledge are represented later, lexically

or post-lexically. Here we have an example of invocation of a distinction between types of abstract knowledge that is forced on the model because it assumes phonemic prelexical representations¹.

Although psycholinguists normally represent higher-order structure as linguistic knowledge, distinct from the sensory signal, some of that structure is directly available from the sensory signal, and we suggest it should not be considered as any more different from the signal than is an abstract prelexical phoneme or feature representation. Information about linguistic structure conveyed by connected speech processes such as those illustrated in Section 2.2 above seem to us to be no more ‘higher-order knowledge’ than is an attribution of feature or phoneme identity from spectro-temporal properties of the signal. Psycholinguistic processes that capture these higher-order regularities direct from the signal thus seem desirable, if only for reasons of model economy. Merge takes a welcome step in this direction, though we believe it would be more effective if it took a more radical approach to the prelexical representation.

In short, we see the incorporation of pre-lexical knowledge of transitional probabilities as a sort of metaphor for just one type of linguistic structure that is available in the speech signal but absent in unstructured phoneme strings, probabilistically encoded or not. Without a wider temporal focus to allow systematic covariation of fine phonetic detail to be detected, and even allowing for transitional probabilities, correcting errors by reference to probabilities of phonemes in undifferentiated strings is unlikely to reflect the way error-correction takes place in normal listening situations. We suggest that the claims made for Merge and indeed all similar models of spoken word recognition could be more convincing if the focus were widened from individual phonemes to a range of different linguistic units with varying temporal domains, and if detailed phonetic information were preserved as late as possible in the model. Furthermore, the output as well as the input must be defined in enough detail to provide a grammar of the mapping between speech and meaning, so that it takes into account the systematic relationship between phonetic detail and phonological and grammatical status.

2.4. *Summary*

In summary, because the phoneme has dominated thinking in both speech science and psycholinguistic research on spoken word recognition, at the expense of other types of phonological and grammatical structure, much of the systematic variation in speech that indicates linguistic structure has been ignored. Short-domain spectral-temporal events that relate most directly to phoneme identity have dominated perceptual research in speech science, together with a tendency to separate segmental and prosodic information in thinking and in research. Partly for practical reasons, this local focus has either been adopted in many computational models of spoken word recognition and lexical access, or else it has strongly influenced them. Despite long-standing acknowledgement of the potential importance of systematic fine phonetic detail (e.g. Elman & McClelland 1986; Pols 1986) the information conveyed by the speech stream has been assumed to be more impoverished than it really is, with consequent biases in the mental processes proposed to recognize and understand spoken words.

3. *Properties of the speech signal*

We have rejected the traditional view that speech is understood by being organised by listeners as well as linguists into independent prosodic and segmental strands and discrete, pure forms of abstract units. Instead, we suggest that the perceptual correlates of linguistic units are typically complex, often spread over relatively long sections of the signal, simultaneously contribute to more than one linguistic unit, and don’t cluster into discrete bundles in time. We believe that this complexity is a crucial determinant of how we understand speech. It is crucial because it reflects

¹ Norris *et al.* (2000:318) note that Merge would get the same results if it used features instead of phonemes prelexically. Since they do not say whether they would structure the features in terms of higher- and low-level nodes with different domains of application, one assumes that they have in mind discrete feature bundles that represent essentially the same thing as the phonemes they actually use. We would therefore expect the same results.

vocal-tract dynamics. And vocal-tract dynamics provide perceptual coherence, which is a central property of natural speech, and possibly the main thing that sets it apart from most synthetic speech.

These views form the springboard for a more detailed consideration of properties of the speech signal that we believe are fundamental, and thus that need to be included in a phonetically-valid model of spoken word understanding. Some of these properties are well understood and well accepted; others are less well understood and in some cases controversial. The section has seven parts. Section 3.1 deals with aspects of what makes a signal robust, natural-sounding, and interpretable as speech; we call this 'perceptual coherence'. Section 3.2 develops the point introduced in Section 2.2 that certain types of grammatical and phonological information are systematically distinguished by fine phonetic detail in the speech signal. Section 3.3 presents a polysystemic declarative linguistic-phonetic model capable of describing this type of information. We suggest that it encapsulates the type of information a listener looks for—and finds—in the speech stream, and that it may be a reasonable metaphor for how that information is represented in the brain during the process of understanding normal connected speech. Section 3.4 discusses the temporal distribution of information about phoneme-sized phonetic segments in the speech signal, distinguishing between short and long domains of influence, and exploring relationships between linguistic informativeness and the temporal extent of the domain of influence of a given phonetic segment. Sections 3.5 and 3.6 discuss the role of rhythm and time respectively in speech understanding. And Section 3.7 formally notes that phonetic categories are like other linguistic categories, and indeed all categories, in being fundamentally relational, or contrastive, in nature. Although the main focus is on systematic acoustic-phonetic properties of speech, results from experiments assessing the perceptual relevance of these properties are referred to when available.

3.1. *Perceptual coherence*

In natural speech, there is a tight relationship between vocal tract behaviour and acoustics, and this relationship provides the signal with perceptual coherence. A speech signal is perceptually coherent when it appears to come from a single talker because its properties reflect the detailed vocal-tract dynamics. Hawkins (1995) used the term acoustic coherence, but changed it to perceptual coherence to place the emphasis on the listener rather than the talker and to reflect the fact that multimodal information can contribute perceptual coherence (Fowler & Dekle 1991; Faulkner & Rosen 1999), or, conversely, incoherence: it is well known that when visual information conflicts with information in the acoustic signal, the resulting incoherence can radically change perception (McGurk & Macdonald 1976; Massaro 1998). However, the present discussion focuses on the acoustic signal because it is the most important medium for speech.

The concept of perceptual coherence is part hypothesis, part factually-based, in that we do not know exactly what properties make speech perceptually coherent, but we do know from many different types of work that small perturbations can change its perceived coherence (e.g. Huggins 1972a, b; Darwin & Gardner 1985). To be heard as speech, time-varying acoustic properties must bear the right relationships to one another. When they do, the perceptual system groups them together into an internally coherent auditory stream (Bregman 1990) or, in some views, into a more abstract entity (cf. Remez *et al.* 1994, Remez *et al.* 1998). A wide range of acoustic properties seem to contribute to perceptual coherence. The influence of some, such as patterns of formant frequencies, is widely acknowledged (cf. Remez *et al.* 1981). Others are known to be important but their contribution is not always well understood. Examples are the amplitude envelope which governs some segmental distinctions (e.g. Rosen & Howell 1987) and also perceptions of rhythm and of 'integration' between stop bursts and following vowels (van Tasell *et al.* 1987); and correlations between the mode of glottal excitation and the behaviour of the upper articulators, especially at abrupt segment boundaries (Gobl 1988; Pierrehumbert & Talkin 1992; Ní Chasaide & Gobl 1993; Stevens 1998). This partial list of properties that can make a signal perceptually coherent presumably includes both general acoustic consequences of vocal-tract behaviour, such as the complex events at the boundaries of vowels and obstruents, and coarticulatory patterns of the particular language and accent being spoken.

Perceptual coherence is a fundamentally dynamic concept, as the examples above show, and it interacts with knowledge. An example of the dynamic nature of perceptual coherence interacting with general knowledge comes from experiments using 'silent-center' syllables. Silent-center CVC

syllables leave intact formant transitions and the duration of the vocalic steady state, but replace the excitation during the steady state with a silence of the same duration, thus preserving only the vowels' durational and dynamic spectral information. Listeners identify the vowels of such syllables with remarkable accuracy (Strange *et al.* 1983; Jenkins *et al.* 1983; Parker & Diehl 1984; Strange 1989). More interesting is that listeners tend to hear such syllables as interrupted by a hiccup or glottal stop, suggesting that the drive to attain perceptual coherence introduces percepts of real-world experiences inherent to vocal-tract behaviour (even while speaking) but not necessarily to normal English speech behaviour. The relevance of direct realist theories (J. Gibson 1966, E. Gibson 1991; Fowler 1986; Best 1994, 1995) to the concept of perceptual coherence is obvious.

It is tempting to try to separate 'general vocal-tract dynamics' from 'knowledge' of language-specific vocal-tract dynamics like coarticulatory patterns (e.g. Perkell 1986; Manuel 1990; Magen 1997; Beddor & Krakow 1999), but this does not seem to be a valid distinction for listeners. How is a listener to know what is general and what is language-specific? Certainly, infants seem to know by 18 weeks or so of age which facial configuration is associated with particular vowel sounds (Kuhl & Meltzoff 1982) but in our view, a perceptually-coherent signal conveys far more subtle distinctions than the fairly obvious ones between the acoustic or visual properties of /i/ and /a/. A mature ability to detect and use perceptual coherence must depend on cumulative experience over several years.

If general and language-specific differences cannot be separated in practice, is it worth separating them on logical grounds? We believe that it is not. From the listener's point of view, the vocal-tract dynamics of his or her language(s) are systematic and must be known about; it is immaterial that they represent only a subset of the possible range (cf. Keating's (1985) similar conclusion for speech production). Furthermore, sounds the vocal tract can make, even ones that all or almost all humans do make, may not be perceived as speech when the relevant experience and consequent knowledge are missing, or when expectations bias the interpretation. For example, when native speakers of non-click languages hear a click language for the first time, they may not hear the clicks as part of the acoustic stream from the speaker's mouth, let alone as a structured part of his or her language. This is especially likely to happen if the context provides an opportunity for a more familiar interpretation. Anecdotally, on hearing Bushman spoken by a man striding rapidly across a desert at the outset of the movie *The Gods Must Be Crazy*, the first author, who was expecting to hear only English, first 'heard' speech plus twigs breaking as they were stepped on, then, realising there were no twigs, immediately 'heard' speech plus a disappointingly crackly sound track to the movie. It took a palpable amount of time, possibly several seconds, to realise that the clicks were integral to the speech, and thus that the language being spoken was probably genuine Bushman, even though she was aware of the existence of click languages and had heard them spoken before. (In her defence, she was very tired and not trying very hard.) Similarly, many Westerners hear dissociations between two parts of Mongolian chanting produced simultaneously by a single person. For example, one mode of Mongolian chant raises the amplitude of a single high-frequency harmonic significantly higher than that of surrounding harmonics. Westerners often hear this mode as having two sources: a deep chant from one person and a whistle from a second source. Though chanting may not be speech, it is closely related to it, and there are parallels with speech. For example, native speakers of English are typically amazed when they hear a baby making an aryepiglottic trill, because its extremely low pitch seems lower than any baby could make; but speakers of languages in which these trills are part of the speech repertoire (John Esling, pers. comm.) are presumably much less surprised.

We conclude that perceptual coherence is part properties of the physical signal and part mental construct. It is not useful to try to separate the two, because the coherent sensory experience of any physical signal is also a mental construct. In other words, just as phonemes are not in the signal, but may be constructed from it, so is all experience of words and their meanings; different expectations, different experiences or training, and different tasks, can influence how listeners respond to structured stimuli. These points have been made in detail by a number of people working in different fields of perception and taking a variety of approaches. Representative summaries can be found, for example, in Remez (2001), Moore (1997), and Warren (1999).

To sum up, perceptual coherence is the perceptual 'glue' of speech (cf. Remez & Rubin 1992). It is rooted in the sensory signal but relies on knowledge; the two are not distinct in this respect, but feed each other. It underlies the robustness of natural speech and determines why it sounds natural, and, conversely, it may be the key to understanding a number of the challenges in producing good

synthetic speech, such as why synthetic speech can be very intelligible in good conditions but tends not to be robust in difficult listening conditions, and why speech synthesized by concatenating chunks of natural speech that preserve natural segment boundaries can give an impression of sounding natural even if it is not very intelligible (see Ogden *et al.* 2000).

Although we can identify a number of good candidate contributors of perceptual coherence, we do not yet know exactly what it is about the phonetics of a particular unit or linguistic structure that causes it to be perceived as a coherent unit. Because coherence seems to result from complex relationships between physical and language-systemic constraints experienced in a particular context, we need a linguistic model that allows clear predictions about the phonetic properties associated with the different phonological and grammatical structures that exist in a particular language. It needs to systematize the complexity introduced by the diverse linguistic factors that influence speech production and acoustics. The next section gives examples of some of these linguistic factors, while the following one describes the type of model we use to try to systematise them.

3.2. Systematic phonetic variation

A speech signal will not sound as if the talker is using a consistent accent and style of speech unless all the systematic phonetic details are right. This requires producing often small distinctions that reflect different combinations of linguistic properties. Speech that lacks these properties will lack perceptual coherence to one degree or another, and in consequence will be harder to understand.

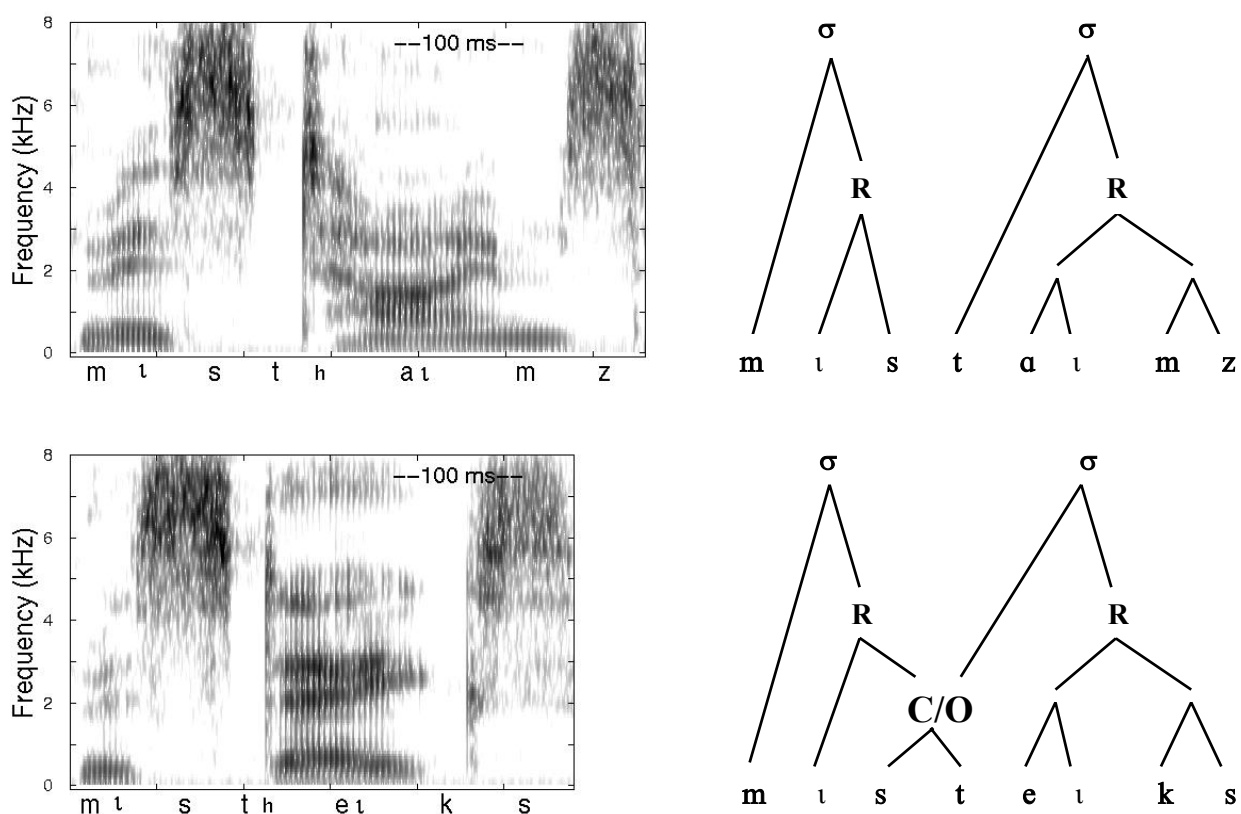


Figure 1. Left: spectrograms of the words *mistimes* (top) and *mistakes* (bottom) spoken by a British English woman in the sentence *I'd be surprised if Tess _____ it* with main stress on *Tess*. Right: syllabic structures of each word.

There are many well-known examples of systematic phonetic variation, some of which have been mentioned above. This section does not provide a complete review, but gives just two less-known examples of systematic variation in fine phonetic detail. Other cases are discussed in Section 3.4.2 below. The first example in this section reflects complex interplay between metrical-phonological and morphological influences and is found in both clear and casual speech. The second comes from casual speech. It can be described as a simple phonetic process of partial assimilation, but it nevertheless results in subtle phonetic marking of a grammatical distinction. Both these examples have potential perceptual salience but would be neglected in most models of speech understanding.

The first example compares the words *mistakes* and *mistimes*, whose spectrograms are shown at the left of Figure 1. The beginnings of these two words are phonetically different in a number of ways, although the first four phonemes are the same. The /t/ of *mistimes* is aspirated and has a longer closure, whereas the one in *mistakes* is not aspirated and has a shorter closure. The /s/ of *mistimes* is shorter, and its /m/ and /ɪ/ are longer, which is heard as a rhythmic difference: the first syllable of *mistimes* has a heavier beat than that of *mistakes*.

These acoustic-phonetic differences in the words' first four phonemes arise because their morphological structure differs. The word *mistimes* contains the morphemes *mis*+*time*, which each have a separate meaning; and the meaning of *mistimes* is straightforwardly related to the meaning of each of the two morphemes. But the meaning of *mistakes* is not obviously related to the meaning of its constituent morphemes, and the word is best regarded as monomorphemic. This difference in the productivity of *mis*- is reflected phonologically in the syllable structure, shown on the right of Figure 1. When *mis*- is non-productive, the /st/ cluster is ambisyllabic (*mistakes*), but the morpheme boundary in *mistimes* prevents /st/ from being ambisyllabic. Hence, in *mistimes*, /s/ is the coda of syllable 1, and /t/ is the onset of syllable 2. In contrast, the /s/ and /t/ in *mistakes* form both the coda of syllable 1 and the onset of syllable 2. In an onset /st/, the /t/ is always relatively unaspirated in English (cf. *step*, *stop*, *start*). The durational differences in the /m/ and the /ɪ/ arise because the morphologically-conditioned differences in syllable structure result in *mist* being a phonologically heavy syllable whereas *mis* is phonologically light, while both syllables are metrically weak. Thus the morphological differences between the words are reflected in structural phonological differences; and these in turn have implications for the phonetic detail of the utterances, despite the segmental similarities between the words. Complex though these influences are, listeners seem to keep track of them and use them, in that unless the various properties have the right relationships to one another, they will be heard as unnatural, and we predict they would be harder to understand.

The second example concerns vestigial acoustic indications of the presence of a /z/ in a heavily-assimilated /zʃ/ context. The two spectrograms at the top of Figure 2 show the first two syllables from *Who sharpened the meat cleaver?* (left) and *Who's sharpened the meat cleaver?* (right) respectively, spoken by the same person. Notice that the /z/ of *who's* is heavily coarticulated with the /ʃ/ of *sharpened*, so that there seems to be little or no change in the quality of the friction that represents the two fricatives. (It is, of course, longer in the *who's* than the *who* version, which provides additional phonetic evidence that there is an auxiliary verb in this structure). A conventional analysis would say that the /z/ is fully assimilated to the place of articulation of the following fricative. However, the assimilation is not complete, because the two /u/ vowels before the fricatives are very different, in ways consistent with alveolar versus palatal-alveolar articulations. The panel at the bottom of Figure 2 shows lpc spectra from 50-ms windows at the beginning of each vowel, as indicated by the connecting lines to the two spectra. Both F2 and F3 are considerably higher in frequency in *who's* than in *who*. That is, an 'underlying /z/' engenders higher F2 and F3 frequencies in the preceding /u/ and, of course, in the /h/ preceding the /u/.

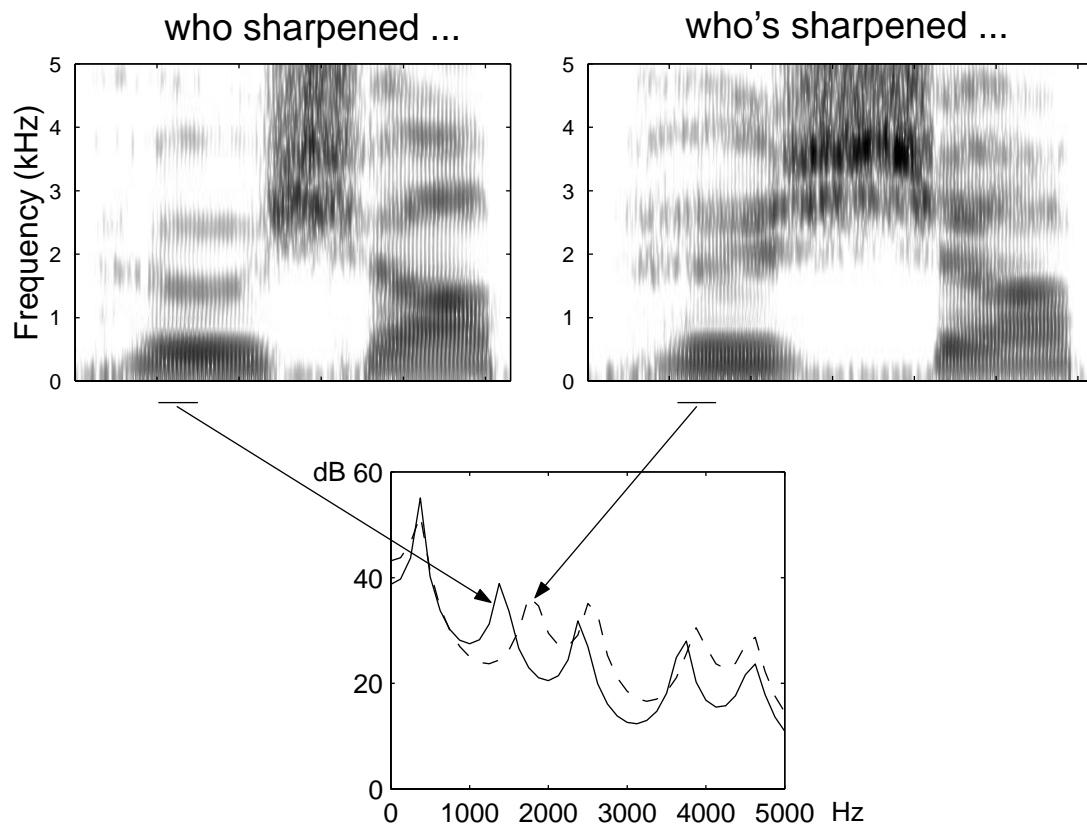


Figure 2. Top: spectrograms of the first two syllables from *Who sharpened the meat cleaver?* (left) and *Who's sharpened the meat cleaver?* (right). Bottom: 50-ms lpc spectra (18-pole autocorrelation, Hanning window) of the first part of the vowel in *who* and *who's*, as indicated by the arrows: solid line spectrum from *who*; dashed line spectrum from *who's*. The horizontal lines under the spectrograms indicate the 50-ms portions of the signal over which the spectra were made.

Although we have not tested the perceptual power of this particular example of systematic phonetic variation, we know the difference in the vowels is perceptible, and are confident that listeners would use it. Both classical and recent experiments (Repp 1982; Strange 1989; Duffy & Pisoni 1992; Pisoni 1997a; Kwong & Stevens 1999) suggest that most systematically varying properties of speech will enhance perception in at least some circumstances, perhaps especially in adverse listening conditions or when the cognitive load is high (Hawkins & Slater 1994; Tunley 1999; Heid & Hawkins 1999). Slightly raised F2 and F3 frequencies may not have a strong effect by themselves, and they are not always present in such sequences; but when they are there, and especially when they co-vary with other cues such as duration of the fricative, they could offer good information about the grammatical structure of the utterance.

3.3. *The linguistic model for Polysp: polysystemic and nonsegmental*

If we are right in suggesting that the properties of speech discussed in the previous two sections significantly affect our understanding of connected speech, then we need a model that can capture them. Hence we regard each 'phonetic segment' as best described in terms of all of its structural properties, rather than solely or mainly in traditional phonetic terms. Acknowledging these structural properties allows much phonetic variation to be seen as systematic, and therefore informative.

3.3.1. *Antecedents and principles from Firthian prosodic analysis*

The linguistic model we propose as most useful for this purpose was developed from Firthian prosodic analysis. Its principles are described by (amongst others) Local (1992) and Ogden & Local (1994), who are at pains to distinguish it from autosegmental phonology (Ogden & Local 1994); instantiations have been used as the theoretical model for synthesizing speech that sounds natural and is robust in adverse listening conditions, namely YorkTalk (Coleman 1994; Local & Ogden 1997) and ProSynth (Ogden *et al.* 2000). Related grammatical models have been used in psycholinguistics e.g. Jurafsky (1996). We outline here only the most important properties of this class of models for our present purposes; fuller treatments are available in the references cited and elsewhere.

The most important properties of this model for us are that it is (1) declarative, (2) nonsegmental and (3) polysystemic. The most important of these is the third—that it is polysystemic. The other two properties more or less follow from this, but we deal with them first because they provide a familiar basic structure from which to begin to discuss polysystemicity.

(1) Because the model is declarative, it attempts to provide a single invariable phonological and grammatical description of any given word or phrase. Differences in phonetic forms of words or phrases are accounted for by differences in the hierarchical structural description associated with each one, as partially illustrated in Figure 1 for *mistimes* and *mistakes*.

(2) The model is explicitly and primarily nonsegmental in that traditionally segmental distinctions are seen as differences in the features placed on particular nodes in the structural hierarchy, and they thus exist in paradigmatic and syntagmatic relationships with one another. Each feature is represented at the highest possible point in the hierarchical phonological structure. When a phonological node possesses a particular feature, that feature affects every phonetic segment within its domain. For example, lip rounding due to a particular rounded vowel can be represented high in the structural hierarchy of the phrase in question, so that an appropriate duration of anticipatory rounding is achieved. How lip rounding is actually realised in each acoustic segment in this domain depends on general principles specific to lip-rounding in the particular language, and on the other properties of the structure—that is, on the full specification of the segments in question.

(3) The model is polysystemic in that language is seen as a set of interacting systems rather than one single system. We have already discussed some examples of polysystemicity in English: distinctions between content and function words in the connected speech processes they take part in (Section 2.2), and distinctions due to morphological differences between words (Section 3.2). Likewise, the nominal and verbal systems of a language can be viewed as separate, and auxiliaries can be systematically distinguished from other verbs, thereby making explicit the fact that they obey different phonetic constraints (Ogden 1999). For example, auxiliaries can be much more phonetically reduced than other verbs: compare *I've seen it* with *I have it*. Another pervasive distinction in English is that between the Germanic and Latinate lexical stress systems. These can be viewed as having different phonological systems and structures, different accentual patterns and so on, which explains why superficially similar polysyllabic words like *unnatural* and *unknown* have different rhythmic patterns from *innate*: *un-* is a Germanic prefix, one of whose properties makes the /n/ in *unnatural* and *unknown* long; in contrast, *in-* is a Latinate prefix, whose properties make the /n/ short. Latinate and Germanic prefixes may both precede Latinate stems, but only Germanic prefixes can precede a Germanic stem. Thus a dual-system contrast combines to make three (but not four) systems. The picture is further complicated by the fact that some prefixes, such as *dis-*, can behave as if they are Latinate or Germanic (cf. *disperse*, *disburse*), while others are superficially similar but arise for different reasons. For example, *mistakes* and *mistimes* are both Germanic throughout, and the contrast between them, illustrated in Figure 1, is due to differences in morphological productivity resulting in differences in ambisyllabicity. These patterns are well known in lexical phonology (cf. Halle & Mohanan 1985) but accounted for there by using derivational rules, whereas our model captures the differences declaratively.

3.3.2. *Some attributes of the model*

The above properties mean that the model has different systems within linked grammatical and prosodic hierarchies. The prosodic hierarchy, for example, extends from intonational phrases to syllabic constituents, with phonological information distributed across the entire structure. Phonological contrasts can be expressed over any domain in the structure e.g. over the domain of the

phrase, syllable, syllable onset, syllable rhyme, and so on. Links with grammatical structure are made at appropriate nodes in the hierarchy through shared features. So, for example, the phonetic properties consistent with a voiced dental fricative in the onset of a weak syllable signal that the syllable is a function word, and provide a link between phonetic-prosodic structure and grammatical structure; they also indicate that the word is likely to be monosyllabic so that the next syllable onset probably signals a new word. The phonological categories and features are abstract: for instance, a contrast bearing the feature label [+ plosive] has different phonetic realisations depending on where it occurs in the phonological structure. The phonetic interpretation of the feature could be entirely different if it occurs in a syllable onset or a syllable coda. Within this system, the /p/ of *potato* has more properties in common with the /t/ of *toboggan* than it does with the /p/ of *cap*.

With reference to the discussion in Section 2.2, there is nothing to stop strings of phonetic segments being identified and referred to in the traditional way, but phonemes are not a part of the structure, and allophones bring with them their full systemic description. Thus, the structural differences between the two glottal stops of [ʔandiʔaʊvə] are immediately apparent and unproblematic. Similarly, the differing representations of /t/ in [kəʔənaiʔ] and [kətənaiʔ] bring with them their structural motivation.

Fundamental to this approach is that each unit of linguistic-phonetic analysis is seen IN RELATION TO its neighbours, and to other units above and below it in the prosodic hierarchy. The relational nature of speech sound categories is likewise fundamental to the long- and short-term phonetic properties discussed in Section 3.4. It is also relevant to memory, and hence to the brain's representation of linguistic categories, as discussed in Sections 4.2 and 4.4. In short, (1) a phonetic segment cannot be properly described independently of the larger structure(s) within which it occurs, and (2) each phonetic segment is definable only in relational terms.

For our purposes, the polysystemic attributes of the model are more important than the others. As we see below (Section 4.2 *ff*) there is neuropsychological evidence to encourage the view that some linguistic subsystems are functionally separate in speech processing. If some are functionally separate, it is worth pursuing the idea that they all are. Thus, although superficially the rhythmic system of English is chaotic, it can be analysed in terms of distinct subsystems, and experienced listeners are presumably able to relate new constructions to one or other of these subsystems.

We are less convinced that the model must be declarative but one reason to favour a declarative over a derivational model is that a polysystemic model requires a declarative structure, because otherwise the links between systems may be lost. Moreover, the declarative model seems broadly compatible with the types of functional groupings of neural activity described below. This does not mean to say that rule-derivations cannot be equally well represented, but what, after all, is a rule, if not the formal description of a general pattern?

Another advantage of a declarative structure is that it can be seen as providing an explicit way of modelling different degrees of casualness onto the same underlying form. This may be easier to envisage from the point of view of speech production or synthesis: when parameter values are changed to affect the tempo of speech synthesized by this model, there are concomitant spectral changes that affect perceived segmental quality (providing that the parameterization is properly implemented). A similar inverse mapping can be envisaged for understanding speech, if only as a metaphor for what really happens.

Lastly, this type of linguistic model in principle allows extension to discourse and interactions such as real conversations, which can be analysed as other systems subject to the same general constraints. To develop this point would go far beyond the scope of our present work, but the fact that it should be possible increases the appeal of the model.

3.3.3. *Relationship between Firthian linguistic analysis and speech understanding*

This model is linguistic. We do not claim that it encapsulates the conceptual structure that the brain constructs in order to understand spoken messages. However, we do suggest (1) that it can help guide the search for system within the acoustic-phonetic variation of speech, (2) that something like it must be used if the linguistic information available in speech is to be exploited by models of speech understanding, (3) that its principles are broadly compatible with current views of how the brain processes sensory information, and that it may therefore represent a more reasonable metaphor for

how the brain constructs meaning from speech than those offered by more standard models that assume that the information in the sensory signal is relatively impoverished.

For convenience, then, we call this model Polysp (for POLYsystemic SPeech understanding), and do not try to distinguish what is linguistic from what might be psychological. We stress that we prefer the term speech understanding to speech perception, because we are trying to address how speech is understood, rather than only how it is given phonological form, the latter being the end-point of most phonetic theories of speech perception.

3.4. *The temporal distribution of information about segmental identity*

We have proposed that properties of speech that make the signal perceptually coherent and that systematically reflect linguistic structure all contribute crucially to whether it is understood easily. We turn now to consider how acoustic cues to phonetic segments are distributed over stretches of speech. Although we have just proposed that a polysystemic nonsegmental linguistic model should underlie models of speech understanding, this is one of those times when it is nevertheless convenient to refer to phoneme-sized phonetic segments as entities in their own right, because phoneme names are familiar and simple.

We draw a fundamental distinction between information that happens relatively fast and information that extends over longer time intervals. This distinction is perhaps more conceptual than actual, inasmuch as it represents not a simple binary division, but rather two points spaced rather far apart along a continuum. However, we believe that the distinction has important implications for modelling speech understanding, and so is a fiction worth supporting.

3.4.1. *Short-domain segmental information*

In our view, the best systematisation of short-time events that we have at present is Stevens' and colleagues' work on relational acoustic or auditory invariants which developed from a search for invariant acoustic or auditory correlates of phonological distinctive features. Ideally, these correlates are invariant relational properties described in terms of differences between successive short-term spectra. Stevens (1998) offers the fullest treatment of this work, using a slightly different style of presentation from his earlier work, summaries of which can be found in Stevens (1983, 1989, 1995) and Stevens *et al.* (1992). Stevens distinguishes three types of features: landmarks, which describe degree of constriction and are hence articulator-free; and two types of articulator-specific feature, one representing place of articulation, and the other representing additional properties like voicing and nasality.

Landmarks in the signal mainly describe manner of articulation via degree of constriction, although in practice, 'pure' manner features may be inextricably associated with excitation type. Landmarks exploit quantal relationships, and most are found in short (10-40 ms) sections representing either local maxima or minima in degree of spectral change in that region. The successive spectra that represent landmarks for consonantal closures and releases fall on either side of an abrupt acoustic boundary. For vowels, they fall where the vocal tract is maximally open between two consonants, and hence approximately the centre of a vowel steady state, marked acoustically by a local maximum in F1 frequency or by relatively high amplitudes at low frequency. So landmarks distinguish between the syllable nucleus and its boundaries.

These short-term events in the acoustic signal thus reflect the way the signal changes as the vocal tract moves at critical points in the speech stream. They are relational, reflecting the fact that even short-term events are interpreted with reference to their context. Notice, however, that landmarks for vowels may extend over longer time domains than landmarks at abrupt segment boundaries involving consonants, at least for relatively slow speech. Longer durations will be necessary for identification of the phonological class of certain vowels (cf. Hillenbrand *et al.* 2000). This is one example of the relative nature of the distinction we draw between short-term and long-term information in the signal.

Articulator-specific features also tend to be short-term and occur in the vicinity of the landmarks, but, unlike landmarks, they can involve several different acoustic properties and can extend for longer, up to about 100 ms.

Since landmarks distinguish between the syllable nucleus and its margins, the phonological features they include (consonantal, sonorant, continuant, syllabic, lateral, and strident), (a) mark parts

of the signal that have status relatively high in our model of linguistic structure, including, vitally, those contributing to rhythm, and (b) singly or in combination with one another, correspond fairly closely to categories identified as robust features by Zue (1985). (The main difference between Stevens and Zue is that whereas robust features are seen as relatively steady state, Stevens' landmarks (especially the consonantal ones) are typically dynamic, and define points in time at which particular critical events occur.) Both these attributes of landmarks seem particularly satisfactory from the point of view of offering evidence that they are central to understanding speech. Although lateral and strident features may be considered as rather different from the others, and not applicable at high levels, in fact there are some independent reasons why they could reasonably be modelled relatively high in structure. Some of these are discussed in Section 3.4.2 below.

Opinions differ as to the validity of Stevens' arguments (cf. the commentaries following Stevens 1989). In our view, relational invariants for features tell part of the story, though not all of it since the documented spectral patterns work better in some contexts than in others. It may be that, when they work less well, other properties of the signal need to be considered. Some of these may be as local and short-domain as Stevens' proposed relational invariants, but in other cases, information over rather longer domains may be relatively more informative. For example, auditory enhancement theory represents a systematic attempt to show how a variety of acoustic properties, some local, and one of at least a syllable in duration, may combine to produce robust percepts of phonological voicing, not by enhancing the abstract feature of voicing itself, but by enhancing one or more intermediate (abstract) perceptual properties that together contribute to the overall perception of voicing (Kingston & Diehl 1994). The details are both controversial (cf. Nearey 1995, 1997) and not fully worked out, but it is sufficient for our present purpose simply to agree that, by and large, a great number of perceptual cues to segment identity provide their information in less than about 40 ms; and that, although that information is of variable quality, there is often much justification for regarding it as the 'primary' cue. That is, short-domain acoustic properties like those Stevens identifies typically convey reliable linguistic information.

Not all short-domain cues are as short as those Stevens identifies, but all appear to be relational, and the time-scale over which they operate is itself relational. An early indication of this is Miller & Liberman's (1979) demonstration that the same synthetic formant transitions can cue /b/ or /w/ depending on the duration of the following vowel, which is interpreted as indicating the rate of speech (cf. also Miller & Baer 1983; Schwab *et al.* 1981). There are, of course, limits to the range of variation possible, and by and large these limits are determined by the surrounding context: each acoustic property is interpreted in relation to other acoustic properties in the signal. Such limits parallel those in the spectral domain. For instance, coarticulated vowels in context are identified no worse than isolated vowels and sometimes better, although the steady states of different coarticulated vowels are not as distinctive in F1-F2 space as those of isolated vowels (Strange *et al.* 1976, Strange *et al.* 1979, Gottfried & Strange 1980, Macchi 1980, Assmann *et al.* 1982). Thus, understanding linguistic meaning appears to be very much dependent on the 'Gestalt' conveyed by the whole signal, rather than on the gradual accumulation of information from a sequence of quasi-independent cues.

3.4.2. Long-domain segmental information

Long-domain segmental information can be defined in terms of time, but is perhaps more sensibly defined in terms of syllables, since syllables vary so much in duration. We define long-domain perceptual information as extending for at least a syllable, or, somewhat arbitrarily, for about 100 ms or more. However, some long-domain information lasts much longer.

Some of the information conveyed over long domains is well established. For example, vowel-to-vowel coarticulation (Öhman 1966; Recasens 1989; Manuel 1990; Magen 1997; Beddor & Yavuz 1995), lip rounding (Benguerel & Cowan 1974), and nasalization (Clumeck 1976; Bell-Berti 1993; Krakow 1993, 1999; Solé 1995). However, with some notable exceptions (e.g. Alfonso & Baer 1982; Fowler & Smith 1986; Beddor *et al.* 1986; Warren & Marslen-Wilson 1987; Krakow *et al.* 1988; Marslen-Wilson & Warren 1994; Beddor & Krakow 1999) there has been comparatively little research on the perceptual power of such cues, and little attempt to integrate them explicitly into either psycholinguistic or even phonetic models of speech perception. Although activation models can in principle accommodate long-domain effects with no trouble, most psycholinguistic and computational models neglect them in practice, because they assume a discrete and abstract input unit,

usually phonemic. Amongst phonetic models, gestural models (Liberman & Mattingly 1985; Fowler 1986; Best 1994, 1995) could be said to be predicated on principles of long-domain influences. However, perhaps because gestures are axiomatic in this class of models, and related to 'linear', phoneme-like segments independent of other linguistic structure, the relative influence of long-domain and short-domain information has received no more attention in gestural than in auditory phonetic models. In particular, not even gestural models of speech perception can account for the class of long-domain perceptual cues known as resonance effects. As will become clear, we believe that resonance effects have great theoretical significance (possibly far greater than their contribution to speech understanding) and for this reason, and the fact that these data are comparatively new, we discuss them in some detail.

In many accents of British English, syllable-onset /l/ is realised as a clear palatalised segment that contrasts with a relatively dark syllable-onset /r/. Kelly & Local (1986) observed that these clear and dark resonances of /l/ and /r/ not only colour the entire syllable of which they are a part, so that the /i/ of *Henry* is darker than the /i/ of *Henley*, but can also affect vowels in neighbouring syllables. Subsequent work has supported their observations with acoustic (Hawkins & Slater 1994; Tunley 1999; Heid & Hawkins 2000) and EMA measurements (West 1999a) of a relatively wide range of carefully-controlled sentences, such as *We heard that it could be a mirror* and *We heard that it could be a miller*. Formant frequencies are lower in parts of the utterance when the last word is *mirror* than when it is *miller*.

Collectively, these studies show that the realisation of these liquid resonances is affected by a wide range of factors about which there is still much to be learned. They are fiendishly difficult to study in controlled experiments, partly because the effects are subtle, but mainly because, naturally, they interact with other linguistic-phonetic variables to produce complicated effects. For example, for independent reasons, Hawkins & Slater (1994) and Tunley (1999) both suggested that unstressed syllables will be more subject to resonance effects than will stressed syllables, but Tunley (1999) also found that vowel height affects them, and this is not straightforwardly separable from syllable stress. Moreover, Heid & Hawkins (2000) conjectured that acoustic consequences of severe vowel reduction might obscure liquid resonance effects.

Heid & Hawkins (2000) tried to disentangle some of the effects of metrical and segmental structure on the spread of anticipatory resonance effects in the speech of one British English male. They found resonance effects in some utterances for up to five syllables before the conditioning /r/ or /l/; this represents a temporal range of half to one second, depending on segmental and metrical structure. These durations are far longer than current models of speech understanding deal with. As expected, they found both segmental and metrical influences on the appearance of these effects, but not always in the expected patterns. For example, it was expected that stressed syllables and intervening lingual (especially velar) consonants would show weaker effects than unstressed syllables and labial consonants. The observed patterns showed that these expectations, though broadly correct, were too simplistic. Intriguingly, some unstressed syllables that showed anticipatory resonance effects preceded a stressed syllable that did not show the effect; in other words, the resonance appears to 'pass through' a stressed syllable to colour earlier unstressed syllables. For example, *it* showed anticipatory /r/-resonance effects in *We 'heard that it 'could be a \mirror* when *could* was stressed but not when it was unstressed (*We 'heard that it could be a \mirror*), even though stressed *could* did not itself show evidence of /r/-resonance effects. This stress-dependent difference may have been connected with the fact that the vowel in *it* was longer and less reduced before stressed than before unstressed *could*. Similarly, segment type (velars *versus* bilabials) appears to have a local influence on the appearance of resonance effects, but need not block their anticipatory spread to yet earlier syllables.

Though much more work is needed before we can fully describe the realisation of resonance effects, listeners know about them and can use them. In a forced choice test, West (1999b) showed that when regions of speech surrounding an /r/ or /l/ are replaced by noise, differences (resonance effects) in the remaining, more remote, regions of the utterance are sufficiently distinctive to allow listeners to tell whether the excised word contained /l/ or /r/. Perceptual tests in our laboratory have compared the intelligibility of synthetic speech with and without resonance effects when it is heard in

naturally-fluctuating noise from a cafeteria. When resonance effects are included, phoneme intelligibility increases by about 10-15% (Hawkins & Slater 1994; Tunley 1999).

Synthetic speech that includes these resonance effects does not usually sound very different from synthetic speech that omits them, and many listeners say they can hear no difference even though their intelligibility scores indicate that they can. We thus conjecture that resonance effects have a subtle effect on the perceptual coherence of the signal, sufficient to increase intelligibility in adverse listening conditions but not to change anything but the most detailed phonetic transcription of the utterance. We do not know whether they are perceptually salient in all listening conditions or only in adverse circumstances, but since speakers produce them and listeners use them to understand connected speech, models of speech understanding must be able to account for integration of subtle information relating to phonetic segments over several hundreds of milliseconds.

3.4.3. *Relationships between informativeness and domain of influence*

Every phonetic segment is probably cued to one extent or another by both long- and short-domain acoustic properties, but, generally speaking, short-domain events within the framework outlined above tend to be highly informative about some aspect of segmental identity, whereas information that takes longer to unfold is likely to carry weaker information, time unit for time unit. However, like the fiction of short-domain *versus* long-domain cues, the correlation between the duration of a cue and its informativeness is imperfect.

Standard examples of short-domain, highly informative events are acoustic cues associated with stop identity that are found in the vicinity of the acoustic segment boundaries, and traditionally described in terms of formant transition patterns and the shape of the burst spectrum, although alternative descriptions such as Stevens' may make more auditory sense. Examples of more long-domain, less informative acoustic cues include information about place of articulation of a coda obstruent that results from coarticulatory effects distributed throughout the rhyme of a syllable (e.g. Warren & Marslen-Wilson 1987), information about the voicing of a coda obstruent that is present in an onset /l/ in the same syllable (Hawkins & Nguyen 2000, 2001, *in press*), and the resonance effects discussed in Section 3.4.2, which may extend over several syllables.

In contrast, however, there are weak short-domain cues (an example is first formant cutback for the voicing of English stops) and very powerful, more long-domain cues such as vowel duration in cueing the phonological voicing of coda obstruents in English and many other languages. Other properties, such as those indicating vowel identity, also take time to build up, and it is not clear whether they should be classed as short- or long-domain. For example, at normal rates of speech, all diphthongs, and vowels like /a/ in some consonantal contexts, can take 100 ms or more before they are identified with certainty. Similarly, reliable information about manner of articulation, including excitation source and certain distinctions that rely on relative duration (e.g. between affricates and stop-fricative sequences as in *grey chip*, *great ship*, etc), by its nature extends over relatively long durations. Stevens' work is interesting in this context because his model looks for highly informative short-domain acoustic events to identify features such as nasality that can in fact spread (often weakly) over rather long domains.

These complex interrelationships between domain size, domain duration, and quality of linguistic information pose problems for theoretical descriptions that rely on phoneme-sized phonetic segments that are unstructured, or in which the various types of structure are separated into rather independent strands. They are more easily accommodated in a theory that accords less status to phoneme-sized segments, and more to time-varying relationships within a systematic, rich linguistic structure.

In summary, we adopt the position that the information the signal provides includes the immediate effect of clear, usually short-term, events (which may resemble Stevens' relational invariants), and also the cumulative effect of information that is distributed more widely across the signal. The distributed information can be very clear, as for strident fricatives and aspects of vowel quality and nasalization, or it can be rather weaker, as in the coarticulatory influences discussed with Figure 2 and in this section. As long as weak information points consistently to the same linguistic structure(s) across time, it seems reasonable to suppose that it could have a strong influence on speech understanding (Hawkins 1995; Hawkins & Warren 1994).

Taken together with the type of structured, nonsegmental, polysystemic model that is Polysp, it follows that we assume that there is no pre-defined order in which a listener must make decisions in order to identify words or understand meaning: decisions are made in parallel, and can be at any level of linguistic analysis. For example, a strong syllable can be identified with high certainty even when its component phones are only partially identified; a strident fricative can be identified with high certainty when there is only low certainty about which syllable it is a member of, and whether that syllable is stressed. These types of information offer islands of reliability at different levels within the linguistic structure. All available evidence is, as it were, fitted together until the best linguistic analysis is found. Hence the importance to perception of weak but systematically-varying information as well as of powerful information, because it adds perceptual coherence.

3.5. *Speech rhythm*

Although both segmental timing and phonological relationships must be organised in particular ways to produce the percept of natural-sounding rhythm, neither segmental durations nor relational phonological contrasts 'are' rhythm, nor do they fully explain its role in perception (e.g. Buxton 1983; Beckman 1986; Couper-Kuhlen 1986; Ladd 1996; Tajima & Port in press; Zellner Keller & Keller forthcoming). The potential units and/or focal points in the signal that are rhythmically important have been investigated for years and some of the most interesting experimental demonstrations relevant to our purposes are in fact some of the earliest (e.g. Huggins 1972b; Kozhevnikov & Chistovich 1965; Lindblom & Rapp 1973).

One important issue for models of speech understanding is that local changes in timing can have effects on perceived rhythm and perceptual grouping that are distributed over a domain of several syllables or the entire utterance. Thus Huggins (1972a), using natural speech to investigate just noticeable differences (JNDs) for segment durations, found that altering segmental durations sometimes caused the perceived distribution of stress in the experimental utterances to change. For instance, when the duration of initial /l/ was altered in the phrase *more lawful than it was*, primary stress was reported to shift between *lawful* (when /l/ was long) and *was* (when /l/ was short). The duration of a segment thus affected both the overall rhythm of the utterance and the perceived prominence of the word *was*, five syllables away from the segment that had been altered. Distributed effects of local changes can also be seen in cases where phonetic detail affects perceived syllabicity and, as a result, word identity. For example, Manuel *et al.* (1992) analysed casually-spoken natural tokens of the words *support* and *sport*, and found that the oral gestures for the consonants /s/ and /p/ in *support* may be timed so as to be contiguous, leading to potential confusion with *sport*. The acoustic consequences of glottal events, however, differed in their timing between *support* and *sport*. They report that these subtle timing differences affected word identification, and by implication also the overall rhythm in the sense of the number of syllables of the perceived word.

In the study cited above, Huggins (1972a) also found that JNDs for segment durations tended to be smaller when subjects reported that they had been attending to the rhythm of the sentence rather than the component phonemes. These observations suggest that the relationship between segmental detail and higher-order organisation is not one-way: the presence of higher-order patterning seems to refine listeners' ability to detect fine-grained properties at the level of segments.

These and similar data have a natural interpretation in the context of a declarative linguistic model, in which rhythm is not associated with any one particular linguistic unit, but is implicit throughout the hierarchy by nature of the relational contrasts that exist at all levels. These relationships collectively determine segmental-level timing (Zellner Keller & Keller forthcoming). So, for example, Ogden *et al.* (2000) argue that in speech synthesis, if all the phonological relationships are stated correctly and the entire structure receives a suitable phonetic interpretation, the resulting utterance is heard as having an appropriate rhythm.

3.6. *Representation of time in speech understanding*

With the exception of the simple durational parameter of VOT and effects due to rate of speech, the temporal properties of speech have received less attention than spectral properties in research on speech perception, and the time domain is relatively underplayed in theories of spoken word understanding. Models frequently assume, for example, that the signal is sampled at small, constant

intervals of 5-10 ms (e.g. McClelland & Elman 1986; Johnson 1997; Nearey 1997), and some of the more interesting neural models represent time as gross quantal changes (e.g. Abu-Bakar & Chater 1995).

Yet the temporal structure of speech offers a rich source of information for listeners. Low-frequency information may normally be processed in the time domain rather than the frequency domain, so speech frequencies that have rates slower than a few hundred Hertz may be perceptible from the temporal structure of the speech signal alone (Greenberg 1996; Moore 1997; Faulkner & Rosen 1999). At these relatively slow rates, amplitude variation reflects manner of articulation and hence syllabic units, and also aspects of segmental identity such as stop release *versus* affricate *versus* fricative onsets, and presumably sonorant consonants *versus* vowels. Other candidates for temporal processing are fundamental frequency and low first formant frequencies.

Nevertheless, when there is a choice, phoneticians typically emphasise the spectral in preference to the temporal aspects of perceptually important acoustic-phonetic properties, although temporal aspects are acknowledged to be important and heavy use of illustrative spectrograms does not hide them. Even exceptions like Stevens' use of frequency and amplitude changes in successive short-time spectra to describe perceptually crucial events (Section 3.4.1) make relatively restricted use of time.

One reason for this preference for the spectral domain may be the common assumption of early abstract representation of the physical signal in terms of a string of feature bundles or phonemes, which simplifies the modelling of the mental representation of speech time to that of sequencing abstract linguistic units. Such simple representations of phonological form may be encouraged by a further tendency to assume that lexical and even sentence meaning remains unchanged no matter how fast the speech: invariance of form and meaning if not of the physical signal. In reality, an important aspect of meaning in normal conversation includes so-called paralinguistic information like overall rate and rate changes, but even if we ignore this, acoustic-phonetic models of speech perception must be able to account for variations in rate of speech, if only because of the segmental reorganisation that typically accompanies them. Time as rate is therefore crucial.

Another obstacle to representation of time in speech perception is that rather little is known about how temporal and spectral properties are integrated in speech processing. Auditory models like AIM (Patterson <http://>) that address this issue are themselves in development, so, in addition to the significant practical challenges of combining work from different disciplines, it is tempting to pursue phonetic and psycholinguistic research independently of psychoacoustic research, and especially of particular models, for they may change.

However, no one doubts that the time domain is a key defining attribute of speech and its perception. Rhythm is clearly crucial, and systematic variation in the time domain appears to underlie the perception of speech as a coherent signal. As noted three decades ago by Huggins (1972b:1280), in speech synthesis “*spectral* detail can sometimes be dispensed with, provided that *temporal* detail is intact”. Huggins cites the example of Cohen *et al.* (1962), who synthesized “surprisingly good speech” from spectrally gross building blocks, by carefully reproducing the temporal properties within and between the blocks. More recently Shannon *et al.* (1995), using a technique that preserved amplitude information but virtually abolished spectral variation, showed that listeners can recognize short utterances and even nonsense syllables almost perfectly from the temporal envelopes of spoken utterances (cf. van Tasell *et al.* 1987). Sine-wave versions of utterances may be perceived as speech if they exhibit temporal patterning characteristic of normal spoken utterances; but if they consist only of sequences of isolated vowels, they are not perceived as speech (Remez *et al.* 1981).

The fundamental question, then, concerns what might be the appropriate window or windows of temporal analysis. Overall rhythm and intonation provide good evidence that we need a long time window in addition to a short one, but there has been little suggestion that we need a time window longer than two or three segments for the perception of segmental identity, even though models of speech production use windows longer than that to model, for example, spread of lip rounding. However, one consequence of our view that phoneme-sized segments should be interpreted in their appropriate structure is that we immediately have a long-domain window, because temporal factors contributing to overall rhythm are distributed throughout the structure and are not separate from segmental identity. Direct empirical support for this point comes from findings that lingual articulation varies systematically with prosodic structure, in a way that emphasizes the differences

between consonants and vowels at prosodic boundaries. Consonant articulations are often 'strengthened' at the edge of a new prosodic domain, and such strengthening is greater for larger domains than for smaller ones, i.e. for domains higher up in the prosodic hierarchy (Fougeron & Keating 1997; Keating *et al.* in press).

In addition to our views on perceptual coherence, two sources of empirical evidence indicate that segments need a long as well as a short time window. The most compelling empirical evidence is the long-domain resonance data discussed in Section 3.4.2. The other evidence comes from a word-spotting experiment by Smith & Hawkins (2000). In word-spotting experiments, the task is to respond as fast as possible when a real word is heard within a longer nonsense sequence, such as *oath* in *puzoath*. Smith & Hawkins (2000) independently varied the syllable stress and the presence or absence of a word boundary in utterances like *puzoath* to make four variants of each. These manipulations produced allophonic differences in the first (nonsense) syllables reminiscent of those in *mistimes* and *mistakes* (Figure 1). As predicted, listeners spotted the embedded words faster and more accurately when the sequence had been spoken as if it had a word boundary, and the segmental allophonic variation had a stronger influence than syllable stress. These data affirm the perceptual relevance of the linguistic structures used in Polysp, and suggest that systematic segmental variation in whole syllables can combine to facilitate word recognition. In other words, optimal perception of segmental information (which, after all, can be seen as the medium that 'carries' prosodic information) may require analysis windows that span at least one or two syllables and sometimes much more. On independent grounds, others have proposed processing buffers and attentive processes that last 100-300 ms (e.g. Silipo *et al* 1999; Grossberg 2000a) we think some such processes may last even longer.

3.7. *The relational status of linguistic categories*

This final part of the discussion of properties of the speech signal largely reiterates individual points made above, but they merit being drawn together in one place. One advantage of Polysp is that it emphasises the relational nature of linguistic categories at the phonetic as well as the phonological level: many linguistic categories can be identified within levels in the model, but none can be properly described independently of its place within the larger structure of the utterance. It is obviously possible and for some purposes entirely justifiable to discuss a single type of category—syllables, and their stress, for example—independently of other categories in linguistic structure, but a complete description of their phonetic realisation is forced to take account of the complete, rich linguistic structure in which they are realised. Phonetic variation is thus simultaneously highlighted and constrained—or systematized—by position in structure, and no one level of analysis can be taken as more important than any other. In particular, it is possible to include phoneme labels in a separate layer, but abstract, context-free phonemes cannot be represented because each node in linguistic structure can only be properly described with reference to its place in the entire structure—i.e. with respect to its context. Thus the gain in using phoneme labels is the provision of a convenient and familiar label, but no more. To be useful at the acoustic-phonetic level, one must know about their context. We believe that this system, in which 'units' are functionally inseparable from 'context', comes close to reflecting current understanding of acoustic-phonetic and perceptual reality, as discussed below.

4. *Properties of listeners*

One purpose of the previous section was to highlight some of the important properties of speech that have tended to be sidelined in models of speech understanding. Similarly, the purpose of this section is to highlight some aspects of cognitive behaviour that we believe should play an important part in models of speech understanding. We arrived at the ideas put forward here over a number of years, as a consequence of trying to explain relationships between phonetics, word forms and meaning, in ways that take into account a wide range of findings from language acquisition, language breakdown in the aphasias, and episodic memory, together with exemplar-based theories of speech perception. We were delighted that, far from being hopelessly speculative, there is good support for our conclusions in recent neuroscientific literature. We are therefore encouraged to make them a central part of our theoretical position.

4.1. *Learning: 'Adaptation' and the plasticity of phonetic categories*

Listeners typically adjust rapidly to different speakers, rates of speech, accents and other sources of variability in the speech signal. Their responses in laboratory tasks also evince a fast-acting sensitivity to statistical, temporal and other properties of stimulus sets. The available data suggest that listeners have detailed memory for specific perceptual experiences and are accordingly capable of continuous learning, and that speech sound categories are labile. These facts present a challenge for models of spoken word recognition and speech understanding in general.

4.1.1. *Adaptation to speakers, rates of speech, and regional accents*

For a long time, the widespread conception of adjustments to different speakers, rates of speech, and so on was that they involved 'normalization' which eliminated some of the undesirable inter- and intra-speaker variability in the signal, allowing an abstract linguistic code to be derived. The normalization approach was not, however, a monolith: Goldinger (1997), Pisoni (1997b), and Johnson *et al.* (1999) give examples of alternative earlier approaches, such as Richard Semon's theory of memory in the early 20th century. Since voice and rate properties have been shown to be retained in memory and to affect phonetic processing (see Pisoni 1997b and Goldinger 1997 for reviews), the term 'adaptation' seems more appropriate than normalization.

Multiple talkers and rates of speech can make some speech tasks harder, but they can also add useful information. The presence of many voices in training materials facilitates the learning of foreign segmental contrasts (Pisoni *et al.* 1994). Multiple-talker word lists are recalled more accurately than single-talker lists at slow rates of presentation (although less accurately at fast rates, perhaps due to distraction: Martin *et al.* 1989; Mullennix *et al.* 1989; Goldinger *et al.* 1991). Presumably the presence of many voices in conversations should enhance communication, since turn-taking between speakers is part of the dynamics of the conversation and changes in the speaking voice convey information.

Although adaptation to all aspects of a new regional accent presumably takes several months, adaptation to a new voice seems to be almost immediate. Listeners appear to use not only 'intrinsic' factors like formant spacing (Ladefoged & Broadbent 1957; Remez *et al.* 1997; Fellowes *et al.* 1997), but also diverse multimodal 'extrinsic' cues including breath sounds (Whalen & Sheffert 1997) and seen or imagined information about the gender of a speaker (Johnson *et al.* 1999). Perceptual adaptation also appears to be reflected in individuals' own speech production in a rather immediate way. Goldinger (2000) showed that when speakers read words aloud before and after hearing a list containing those words, their productions after exposure to the list were judged (by third parties) to be better imitations of the tokens that they had heard.

4.1.2. *Adaptation to properties of stimulus sets*

The extraordinary sensitivity of listeners to the detailed properties of stimulus sets is demonstrated directly in the vast literature on boundary shifts in categorical perception and other experiments involving phoneme identification, as well as indirectly by the type of experimental controls typically built into such experiments to minimise unwanted influences on establishing a boundary between two phonemic categories (cf. Repp 1984; Repp & Liberman 1987). Many experiments on categorical perception show category-boundary shifts as the distribution of the stimulus set changes, and shifts can be induced in the voicing boundary by randomising stimuli so that different talkers are heard in the same session, or in separate blocks (e.g. Johnson 1990; Green *et al.* 1997). These psychophysical context effects on category boundaries are well known, and are described in more detail in Section 5.4 below.

4.1.3. *Adaptation as a form of learning*

The type of training that listeners undergo can change their responses to acoustic stimuli. In the real world of phonetics teaching, we have long known that this is true, for much time in practical phonetics is spent on learning to perceive distinctions that are phonemic in someone else's language. Evidence from infant speech perception indicates that what we are probably doing is relearning what we apparently could do as babies (for a review, see Pickett 1999:ch 13). There are also a number of laboratory demonstrations of this process. Recent examples conducted in the context of the perceptual magnet effect and the internal structure of sound categories (Grieser & Kuhl 1989; Kuhl 1992;

Iverson & Kuhl 1995, 1996; Kuhl & Iverson 1995) include Guenther *et al.* (1999), who found that categorization training made listeners less sensitive to differences between non-speech stimuli, while discrimination learning increased sensitivity, and Barrett (1997), who found similar results for music and speech stimuli. This work is described in more detail in Section 4.4.2 below. Learning is also involved in more commonplace experiences such as adaptation to a new regional accent. We have recently observed that listeners learn to use very subtle acoustic properties of the stimulus set and change criteria rather fast when the situation demands it (Hawkins & Nguyen 2001).

4.1.4. *The theoretical significance of adaptation effects*

Adaptation phenomena challenge conceptions of speech sound categories in a number of ways. One of these concerns the specificity of information remembered: how does a stimulus in a familiar voice map on to similar productions by the same or similar speakers? Then, given this specificity, how do listeners generalize even from limited experience with a speaker's voice so that novel tokens in that voice can be identified more easily than novel tokens in an unfamiliar voice, as shown for example by Nygaard *et al.* (1994) and Nygaard & Pisoni (1998)?

In addition, adaptation effects force us to recognize the lability of speech sound categories: that boundaries between categories are plastic, changing with factors such as surrounding phonetic context, talker, speaking rate, and all kinds of linguistic experience such as lexicality, frequency, and recency, as well as the detailed acoustic properties of the particular speech sound. This context-sensitivity is not simply a 'performance' factor that influences how the category is identified from the sensory signal: it indicates that the mental representation of linguistic-phonetic categories is fundamentally relational and multidimensional. Finally, although many factors that cause category boundary shifts are well known, the temporal evolution of the perceptual process which results in one categorization or another is poorly understood. Adaptation phenomena, reflecting as they do real speech situations, suggest that this dynamic aspect of categorization is particularly important. These conclusions contribute to our conceptualisation of memory, language acquisition, and the nature of linguistic categories, as discussed in the following sections.

4.2. *Memory*

Memory is obviously crucially involved in understanding speech, although its prominence in the literature on models of speech perception has waxed and waned over the years. Early work distinguishing between short-term 'auditory' and longer-term 'phonetic' memory modes was introduced to explain categorical perception (Fujisaki & Kawashima 1970; Pisoni 1973) and has influenced work in that area ever since, often by manipulation of interstimulus intervals to gain insights, for example, into discriminatory vs. classificatory capabilities (e.g. Werker & Tees 1984, Werker & Logan 1985). This basic distinction has parallels in the psychoacoustic literature, for example in Durlach & Braida's (1969) distinction between a sensory-trace mode and a context-coding mode, applied respectively to basic auditory sensitivity and phonetic labelling in work on perception of consonants and vowels (Macmillan *et al.* 1988). Note the explicit connection in the psychoacoustic model between phonetic categories and context.

These memory modes are undoubtedly important, but they focus on explaining one particular type of (possibly not very natural) speech perception, and there are probably other 'memories' capable of holding other types of information, such as a 'working memory' for processing syntax, for example. When the focus of phonetic research moves from labelling phonemes to understanding a message, the importance of more wide-ranging connections between speech and memory becomes clearer.

At the risk of stating the obvious, memory encodes experiences. Hence the memory of an object or event is based on sensation: what it looked like, sounded like, what it felt like to touch or manipulate, what happened when it was manipulated in a particular way, what emotions were experienced at the time, and so on. Memories are thus complex, multimodal, and intensely personal: no two people have identical memories of the same event, not least because connections between the various modalities of a single memory must be deeply affected by what the individual is attending to at the time. Concepts related to words are formed from experiences in interaction with the individual in just the same way. MacWhinney (1999:218) eloquently describes multimodal, dynamic memories for the concept *banana*: the shape, colour, texture, smell, range of sensations while peeling a banana,

taste, and possibly other sensory experiences associated with bananas may all be activated when we hear or see the word *banana*. He terms these sensations ‘affordances’, and describes a theory of ‘embodiment’ of which affordances (neural responses reflecting features of the object relevant to interacting with it) are one of four perspectival (cognitive) systems necessary to understanding sentence meaning, interpreted in the rich and realistic sense of implied as well as denotative meaning.

This interesting polysystemic approach is intrinsic to behaviourism as well as to branches of philosophy, but it did not play a major part in mainstream linguistics and experimental phonetics of the last 40 years. However, it was not completely neglected. For example, MacWhinney takes the term affordance from Gibsonian philosophy (J. Gibson 1966, E. Gibson 1991), which underpins the direct realist class of phonetic theories of speech perception (e.g. Fowler 1986; Best 1994, 1995). The approach is also compatible with exemplar-based theories of speech perception (e.g. Goldinger *et al.* 1991; Goldinger 1997, 2000; Johnson 1997).

Most pertinently, there is now a variety of neuropsychological evidence that memories and concepts are linked to the modalities in which they were experienced, and thus are complex and multimodal and can include a temporal aspect. Evidence from aphasia, and from functional brain imaging (fMRI, PET, and ERP) and neurophysiological (EEG, MEG) studies of normal people converges to show that, while some particular parts of the brain are closely associated with specific behaviours, no one part can be said to govern a particular behaviour. Some of this evidence and associated argument is quite old (e.g. Warrington & Shallice 1984; Warrington & McCarthy 1987), but the surge of relevant work in functional brain imaging is recent. Excellent recent reviews and discussions are provided by Coleman (1998) from the phonologist’s perspective and Pulvermüller (1999 and associated commentaries) from the psycholinguist’s perspective, so the main points need only be summarised here.

Memory for language, and hence language processing, including processing individual words, is distributed in both cortical hemispheres, as well as more primitive parts of the brain such as the limbic system and possibly the cerebellum. The limbic system is basic to the experience and expression of emotion, while the cerebellum is concerned with coordinated movement. While certain parts of the brain such as the visual cortex are involved with many words, there are differences in the distribution of brain activation associated with different categories of word. The evidence that ‘vision’ and ‘action’ words activate different parts of the brain is particularly convincing. This includes the gross grammatical distinction between nouns and verbs, and more fine-grained distinctions such as words for animals versus for small man-made objects like tools. Broadly speaking, words that involve action are associated with parts of the brain that are concerned with action and motor control, while words that do not generally involve action are more strongly associated with visual areas. There is also evidence for differences in brain processing between so-called content and function words. While both hemispheres are strongly involved in processing content words with physical referents like concrete nouns, the more abstract function words are more localised to the left perisylvian region, i.e. the region around the Sylvian fissure which separates the temporal from the frontal lobe, and that is traditionally associated with speech and language; it includes Broca’s area and primary auditory cortex.

This brief summary indicates not only that memory is fundamental to models of speech understanding, but that the theory of memory adopted can deeply influence other aspects of the model. That memories are dynamic and structured according to the modes of experience they encompass lends particular force to the arguments in Section 4.1 above that the speech understanding process is fundamentally adaptive to the detailed circumstances in which speech is heard.

Further, this model of memory is compatible both with exemplar-based models of episodic memory and with our view of the mental organisation of speech represented by Polysp’s declarative, polysystemic properties (Section 3.3). By implication, if the linguistic structure of Polysp is relevant to speech perception, then the mental representation of linguistic structure is grounded in tangible phonetic memories. These phonetic memories are crucially linked to other sensory memories which may not be linguistic. Coleman has likewise argued that the mental representation of words is “essentially phonetic, rather than symbolic-phonological.” (ms: abstract), as discussed in Section 4.4 below.

Another consequence of this reasoning is that it offers the chance of laying to rest some of the inconsistencies of standard phonetic theory, especially those that attempt to separate ‘levels’ of

phonetic analysis. If the representation of speech and language is based on distributed, multimodal, dynamically-organised memory traces, then the same sensory information can be implicated in more than one class of memory, or more than one use to which memories are put. Thus, as described above, systematic fine phonetic detail simultaneously contributes segmental (allophonic) and prosodic information (e.g. Fougeron & Keating 1997; Smith & Hawkins 2000; Keating *et al.* in press). Perhaps most significantly, we can re-evaluate the traditional separation of speaker and message—of voice quality from segmental from prosodic information—as normally undesirable. (Equally, separating them is unproblematic if it is done for practical reasons.) Fundamental frequency, for example, can feed functional groupings of brain cells that themselves combine with other such functional groupings to produce complex memories of both speaker and message. Once this is acknowledged, ‘message’ can be interpreted as far more than simple phonological form and lexical item: it can encompass the whole range of sociophonetics, together with phonetic markers that facilitate successful conversation, such as slowing rate and adding creak to indicate the end of a conversational turn. In short, most information available from the spoken signal is truly polysystemic. We have long known that this must be the case, and that some languages (for example tone languages like Mandarin and Burmese) are probably better analysed in this way. We lacked the independent empirical support for its psychological reality, which the neuropsychological evidence now provides. In sum, we suggest that this approach allows phonetics to take a central place within a broad theory of the communication of meaning, rather than being seen as an arbitrary and somewhat inconsequential carrier of meaning. This argument for phonetics parallels MacWhinney’s (1999) approach to how we understand sentence meaning.

4.3. *Acquisition*

Our view of how babies learn to understand speech is more or less dictated by our position on broadly-defined adaptation effects and memory and the polysystemic, non-segmental principles of Polysp, although it is fairer to say that it was evidence from infant speech perception that influenced our thinking on the other areas, rather than *vice versa*. Following Jusczyk (1993, 1997) we assume that children exploit statistical properties of spoken language to ‘bootstrap’ their way into understanding speech. These distributional regularities are inherently context-sensitive and thus relational, and can arise at any level of formal linguistic analysis; Jusczyk’s work has identified prosody and phonotactics (i.e. rhythm/intonation, and sequential sound dependencies) as early influences on learning to segment words in the first year of life (Jusczyk *et al.* 1999a, b). The process is gradual, and relies on picking out salient properties of the signal that make discernible patterns in particular contexts, building up mental categories from these features, and using those early categories to further systematize the structure of sound and meaning with reference to categories that are already established. This model is broadly compatible with Plaut & Kello’s (1999) model in which phonetic input is mapped directly to meaning, and with a number of other models of how babies learn phonology (e.g. Suomi 1993). The approach is also compatible with models of how children learn grammar from distributional regularities in the input, without recourse to innate mental constructs that specifically concern grammatical relations (e.g. Rumelhart & McClelland 1986; Allen & Seidenberg 1999). Our general view has been beautifully put by Smith in the context of how children learn about nouns: “.... word learning biases that constrain and propel learning in certain directions are themselves made out of general associative and attentional processes. Each new word learned by a child changes what that child knows about learning words—adding to, strengthening, weakening associations among linguistic contexts and attention to object properties. In the end, word learning looks special and predestined. But specialness and destiny are themselves made out of more ordinary stuff.” (1999:301).

4.4. *Linguistic categories*

The foregoing discussion makes it clear that we see linguistic categories as (1) self-organising, (2) multimodal and distributed within the brain, (3) dynamic, and (4) context-sensitive (or relational) and therefore plastic, or labile. We suggest that speech sound categories have these same properties. A familiar way to conceptualise this is in terms of polysystemic hierarchies such as those Polysp provides.

4.4.1. Self-organising, multimodal and distributed

The assumption of self-organising categories is consistent with current computational models of acquisition such as those mentioned in Section 4.3 as well as those discussed in Section 5.5 below. Self-organising categories also allow continuous adaptation to new information, as proposed in Section 4.1. Categories ‘emerge’ as a result of the distributional regularities of the input experienced, modulated by attention and the task(s) at hand, because organization of sensory information into coherent groups of ‘like’ and ‘unlike’ is a fundamental property of how organisms interact with their environment.

Speech categories are multimodal and distributed within the brain because they rely on memory, which is inherently multimodal. The association of meaning with speech learned under normal communicative circumstances is likewise multimodal. Both Pulvermüller (1999) and Coleman (1998, ms) have proposed detailed arguments consistent with this idea.

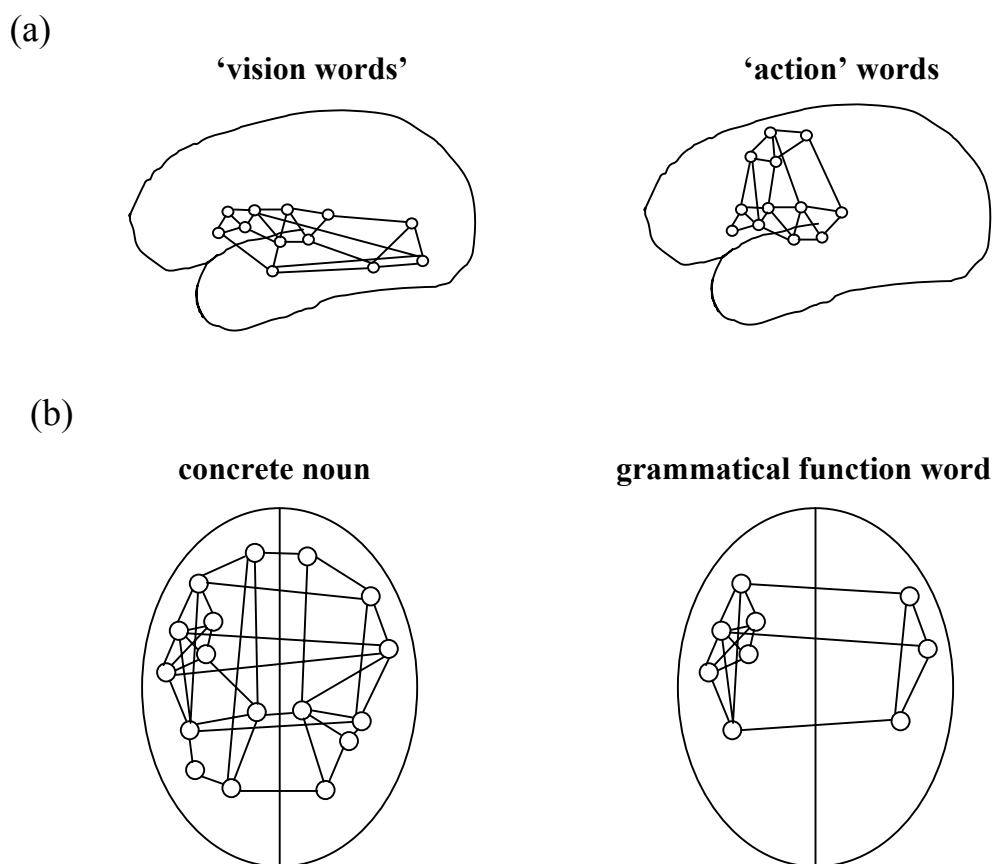


Figure 3. Schematic diagrams of the brain showing representations of the type of functional cell groupings that may represent different types of spoken words. (a) View of the left cerebral hemisphere, showing distributions of activity during processing of vision and action words. (b) View from above (left hemisphere on the left), showing different degrees of lateralisation in the two hemispheres for content and function words. Adapted from Pulvermüller (1999).

Pulvermüller’s (1999) model of how the brain encodes meaning for words postulates that word representations typically have two parts, a perisylvian part related to the word form, and thus to phonetics/phonology, and a part distributed mainly in other parts of the brain that represents its semantic word properties, and thus reflects the individual’s experience with that word, as described in Section 4.2 and illustrated schematically in Figure 3.

Pulvermüller's model is based on associative learning and in particular the principles of complex, functional cell assemblies proposed by Hebb (1949), modified by more recent research in neurophysiology and functional brain imaging. Hebb's basic principles of functional groupings of cell assemblies have parallels with the coordinative structures proposed for control of skilled movement (Kelso *et al.* 1979, Saltzman & Kelso 1987, Saltzman & Munhall 1989) that are axiomatic in articulatory phonology (Browman & Goldstein 1989, 1990, 1992) and may be more familiar to linguists and phoneticians. (Coordinative structures are functional synergies between groups of articulators.) Hebb's functional cell assemblies are also broadly compatible with the affordances of direct realism, and MacWhinney's (1999) general theory of embodiment (Section 4.2). Moreover, Hebb's formulation includes Gestalt principles, which have been incorporated into some models of hearing, notably by Bregman (1990). Remez and colleagues have argued that Gestalt principles do not fit phonetic evidence, and that phonetic categories are inherently abstract, but, while we acknowledge the value of many of their arguments, we are not convinced that wholesale rejection of Gestalt principles is necessary (cf. also Fixmer 2001). There is considerable appeal in finding parallels in the organisation of categories in phonetics, grammar, and meaning, if at the same time they can be externally motivated, in this case by being linked to neurolinguistic data. It is thus worth looking at Pulvermüller's arguments in more detail.

In Pulvermüller's model, a cell assembly is essentially the functional physical substrate of the computational scientist's emergent category. All mental categories are emergent because all are groups of neurons that develop functional unity due to persistent stimulation of the same groups of cells. That is, cell assemblies form when similar sensory episodes are repeated sufficiently often. Once part of a cell assembly, connections between the cells that comprise the assembly become more easily activated, and may continue for longer in a self-perpetuating process of excitation called reverberation. Initial activation (ignition) of all cells in the assembly is more or less simultaneous, and is followed by reverberation, which defines the dynamic activity pattern of the cells. Thus an activated cell assembly is a dynamic, spatiotemporal pattern of neural firing. The particular pattern of activity can in turn activate other cell assemblies.

Pulvermüller suggests that 'synfire chains' (Abeles 1982, 1991) provide support for the neurological reality of Hebbian cell assemblies. Although regarded with some scepticism when first proposed, synfire chains now seem to be accepted (cf. Arbib 1995) as spatio-temporal patterns of cortical neuronal firing that create wavelike patterned sequences of activity. What characterizes a synfire chain is that a subgroup of neurons fires synchronously, activating a second subgroup of (at least as many) neurons that in turn fire synchronously, and so on. Their timing is precise to 1-2 ms. They form networks that appear to reverberate, and amongst their interesting properties for speech and language understanding are that a given neuron can take part in more than one subgroup, and that connections can skip subgroups in the chain to directly link, say, subgroups 1 and 3 in the chain (Pulvermüller 1999:256). Abeles (1991:258) points out that groups of synfire chains are equivalent to multilayer perceptrons (essentially feedforward neural networks) in the sense that the first stage of all the synfire chains can be seen as equivalent to the input layer of the multilayer perceptron, the second stage as a hidden layer, and so on. Although each synfire chain acts only as a feedforward mechanism, synfire chains can be formed in a network that has feedback connections.

Pulvermüller suggests that entire words or morpheme strings might be laid down as synfire chains. Synfire chains might form the neural basis for Coleman's (ms) proposal that words are represented in the brain as trajectories through space-state networks that relate auditory states which vary through time. These networks include spectral information and deal well with rate changes, although more work is needed to establish how to capture the correlation between rate and spectral change.

The application of this work to speech and language is exciting, but the thinking behind it is highly speculative at present, and many details are under debate both by Pulvermüller and by others (see commentaries on his paper in the same journal volume). For example, instead of a separate cell assembly for each word, a number of words of the same semantic category, such as those with similar meanings, or action words vs. vision words, could have the same cell assembly, but would be distinguished by the firing patterns within it. Alternatively, different functional structures with overlapping sets of cells seems equally plausible: an example is the semantic attributes of *crocodile*

versus *alligator*, for which the colour neurons would differ somewhat (and presumably, for experts, some other cells, such as those reflecting aspects of shape).

However, there seems general agreement that there is good evidence that functional groupings of brain cells associate modality-specific sensory memories for words with concepts—that is, with meaning—and that some sort of word form representation is localised to the perisylvian area, particularly in the left hemisphere. For example, so-called function words are represented more strongly than content words in the left perisylvian area, and less strongly elsewhere, and it is argued that this could be because many of them are less intimately connected with sensory experience than, say, concrete nouns. Likewise, cells in the left perisylvian cortex are activated when meaningful words are processed, but not when pseudo-words are processed. Our interpretation of this last finding is that phonemes, as abstract entities, may not be represented in this region, for if they were, then phonotactically-acceptable pseudo-words should activate them. However, some aspects of phonological form do seem to be implicated in this area (cf. Coleman 1998).

There is less agreement about how these principles work for grammar and sentence meaning. Relatively simple cell assemblies may combine into more complex functional assemblies, by various types of connections between cell assemblies and groups of cell assemblies, and by their activation dynamics. Such links between functional groups of cells potentially offer the flexibility and possibility for nuance that characterise language and speech. Extrapolating from these views, we suggest that complex connections between assemblies could be involved even in recognising isolated words: as noted earlier, cross-modal priming studies indicate initial activation of all possible forms of an ambiguous word, with a relatively fast decay of those that are contextually inappropriate (Swinney 1979). Although much work remains to be done before these processes will be understood for either isolated words or words in context, the relevance to the hierarchical structures of standard linguistic theory is obvious.

The possibility of increasingly complex links between simple cell assemblies is obviously also relevant to the hierarchical structures of polysystemic, nonsegmental phonology and thus of Polysp. This interpretation is consistent with Coleman's views, and although it is not developed by Pulvermüller, to do so would require only a small shift in focus. Both Coleman and Pulvermüller see phonological representation as emergent categories associated with lexical representation and anatomically localised in the perisylvian area. In Pulvermüller's terms, these categories could be cell assemblies, or parts of cell assemblies for lexical items realised as synfire chains. Coleman suggests that semantic and all manner of phonetic information are bound together and accessed almost simultaneously when words are being processed; phonemes play no part. Pulvermüller has the same interpretation except that he favours strings of phoneme-sized units sensitive to their segmental context and neglects the relationship between the physical phonetic signal and prosodic and grammatical structure. If the context-sensitivity included information about rich (polysystemic) linguistic structure, then the two views would be in agreement.

Coleman (1998) further argues that there is just one phonological lexical representation for speech, and it is mainly auditory rather than motoric. Cross-modal integration of auditory and visual sensory information is thought to take place in the superior temporal lobes, but though the visual information is about movement, it is not of course the actual movement required to produce speech; sight of movement, and haptic or kinaesthetic sensations of one's own or others' movements, are represented in different, modality-specific regions of the brain. Coleman suggests that, to speak, an individual must translate the auditory lexical representation into motor gestures. This has a number of interesting implications, of which the two most important for our present purposes concern phoneme monitoring experiments and long-domain resonance effects.

Coleman (1998, ms) points out that phoneme monitoring is almost the only single-word perceptual task known to activate Broca's area, which is associated with articulatory encoding. He concludes that articulatory encoding is involved in phoneme monitoring but not in long-term memory for words. He also summarises work by Démonet and colleagues showing that activation of perisylvian cortex in phoneme monitoring is much later than in the normal course of lexical access. This work suggests that decomposing a nonsense word into constituent phonemes is a different and slower task from accessing a lexical entry from auditory input. Démonet suggests a sequential mode when decomposition into phonemes is necessary, but a probabilistic, nonexhaustive way of processing

lexical items in semantic tasks. These data lend independent support to our position that phoneme identification is neither a necessary nor a normal part of understanding speech.

Coleman's conclusions that lexical representations do not involve memory for articulation offer a pleasing way of describing how resonance effects and other forms of long-domain coarticulation can arise. If the word is translated into articulator movement only after the individual words have been slotted into place, consistent with many models of speech production—cf. Harley (1995) for a textbook account—then the necessary actions can be planned as a coherent, efficient sequence according to the demands of the particular accent and language. This does not explain what constrains the spread of resonance effects, nor why listeners are able to use them to facilitate perceptual tasks. If, moreover, Pulvermüller is right in his view (1999:321) that articulatory as well as acoustic and semantic aspects of each word are bound together into one functional unit, then this explanation of how long-domain resonances arise does not work. Much more work is needed before these questions can be answered.

However, the idea that long-domain 'units' of linguistic analysis could emerge from repeated exposure to consistent resonance effects fits comfortably within Pulvermüller's general Hebbian approach and can in principle be accommodated within the rich polysystemic structures of Polysp. Presumably, such coherent long-domain effects could operate in a similar way to longer prosodic units and might be represented high up in linguistic structure. In this case the complexity and flexibility offered by synfire chains or similar processes may be sufficient to account for both local phonetic information and long-domain resonance effects. Local information could be mediated by one type of detailed pattern-matching, while a more general picture is sought at a more global level of analysis in which local phonetic detail might be effectively treated as noise.

Another possibly contrary view to Coleman's (1998) is provided by Rizzolatti & Arbib (1998) and Arbib (2000), who show that monkey cortex includes 'mirror neurons' that discharge both when the monkey grasps or otherwise manipulates an object itself, and when it sees someone else grasp or manipulate the object. The mirror neurons appear to link the observer and the actor by matching similar actions regardless of who makes them. Rizzolatti & Arbib suggest that such a system could be fundamental to successful communication. Mirror neurons offer striking relevance to infants' early sensitivity to the connection between facial expression and speech sounds (Kuhl & Meltzoff 1982) and to motor theories of speech perception, especially Gibsonian approaches in which affordances are fundamental. Those mirror neurons found so far, however, appear to be for activities that can actually be seen; since most activity required for speech cannot be seen or otherwise directly experienced by the listener, it is a big—though not inconceivable—step to assume that mirror neurons or their like underlie perception of speech.

In summary, we suggest that all linguistic categories, including phonetic ones, are constructed by the individual by organising memories of his or her own sensory experiences into functional groups that reflect frequently associated factors. Memories are axiomatically modality-specific, but the linguistic categories they underlie are multimodal because they relate sensory linguistic-phonetic experience to experience of the referents of language. Successively more complex or more all-inclusive groupings represent abstraction. Consistent with Polysp's polysystemic properties, individual language-relevant memories can be part of a large number of different functional groups, so a given piece of phonetic information may contribute to several quite distinct percepts. According to this view, attributes of voice quality can contribute to phonetic information about prosody, segmental identity, and discourse structure while simultaneously contributing to the percept of the speaker's identity and affecting the listener's attitudes. For example, if the listener likes the speaker, then he or she will probably have a positive attitude towards other speakers who have similar voices, regional accents, rhythmic patterns and so on. Supportive data are provided, for example, by Johnson *et al.* (1999).

4.4.2. *Dynamic, context-sensitive (relational) and therefore plastic*

The case for dynamic and context-sensitive phonetic categories needs no reiteration, but comments in the recent literature suggest it is worth discussing the claim that plasticity (or lability) is a fundamental property of phonetic categories. We offer both a logical argument and empirical evidence. The logical argument for plasticity of phonetic categories is that all mental representations of categories are plastic and context-sensitive, and there is no reason to suppose phonetic categories

should be different. Birds are often identified by their silhouette, yet the prototypical blackbird is expected to have a plumper silhouette when you identify it on a lawn covered in snow compared with one bathed in sunshine, because it will have fluffed its feathers out to keep warm. Thus, in identifying the species of a particular bird, the absolute criteria for shape will change according to the perceived ambient temperature, and therefore the relative importance of shape and other criterial attributes such as colouring and beak size may also change. Just so with speech: the analogy with stimulus range and trading relations between multiple, apparently redundant, acoustic cues is obvious. In other words, in terms of perception, phonetic categories are no less contrastive than phonological ones. To have a contrast, there must be a context. If the context changes, then the relative perceptual salience of the identifying physical attributes of the phonetic category will probably also change.



Figure 4. Top left: A stick figure of a person. It has no distinguishing characteristics and might be interpreted as representing a man or a woman. Top right: when long hair is added to the same stick figure, the standard interpretation is that the one on the left is male and the one on the right is female. Hair length is criterial of gender for stick people. The photograph below shows that hair length is irrelevant to classifying the gender of real people, at least when they are hippies, whereas attributes such as bone and muscle structure are important. The stick figures are analogous to highly controlled synthetic stimuli, and the photograph to natural speech.

Figure 4 illustrates these principles. A ‘stick person’ without hair is normally taken to be male, as in the top left of the figure. Add hair, as at the top right, and the figure is conventionally taken to represent a female (especially if both figures appear together). But this convention only works in the context of stick people, for many men have long hair. Hair length is irrelevant in deciding the gender of the two people in the photograph at the bottom of the figure. Instead, one must focus on quite

different attributes, such as bone and muscle structure of the face, arms and hands, possibly the type and distribution of tattoos, and so on. Thus a property that completely defines class membership in one context is completely unhelpful in another. This figure has a message for speech perception methodology, in particular the issue of what constitutes appropriate experimental control: when synthetic stimuli are stylized and strongly over-simplified, or natural utterances are cross-spliced to produce unnatural combinations, there is a real danger that the relative perceptual salience of particular acoustic properties may be misinterpreted because their necessary context is missing or distorted.

The empirical evidence for plasticity comes from a number of sources. One is the adaptation effects discussed in Section 4.1. Another, and amongst the most interesting, is Coleman's (1998) interpretation from neurological and other evidence that lexical representations are phonetic, for lexical representations are by definition contrastive. This is supported by work we have reported elsewhere that suggests that systematic fine acoustic-phonetic detail in natural speech facilitates lexical decisions. Hawkins & Nguyen (2000, in press) have shown that in several varieties of British English, longer and phonetically darker word-onset /l/ tends to co-occur with voiced coda obstruents, and that listeners can exploit this knowledge in lexical decision and lexical identification tasks. As noted above (Section 4.1.3), work in progress suggests that listeners may be very quick to assess whether this type of acoustic information is sufficiently systematic to use as a perceptual cue (Hawkins & Nguyen 2001). These data are consistent with dynamic models of perception in which new exemplars have a large effect on category boundaries, as discussed in Nguyen & Hawkins (1999).

Our final example concerns the perceptual magnet effect (PME: Kuhl 1992; Iverson & Kuhl 1995, 1996; Kuhl & Iverson 1995). Kuhl coined this term to describe experimental findings that were interpreted as showing that listeners' experience with the sounds of a particular language causes poorer discrimination around the best exemplar of a phonetic category in that language than around a poor exemplar of the category (i.e. one that falls nearer the boundary with another phoneme). She interpreted this loss in discrimination as a reorganisation or distortion of psychoacoustic space, such that some sort of prototype develops, towards which physically similar stimuli are drawn as if by a magnet. For a number of reasons, theoretical and methodological, this interpretation has generated a lot of interest, with many researchers' views being strongly polarised for or against the PME. The aspect of this literature that is relevant to the present discussion is that much of the negative criticism of the PME, and of Kuhl's rather theory-neutral use of the term prototype, seems to be predicated on the assumption that the mental representation of phonemes (or phoneme-like phonetic segments) must be invariant. For example, one criticism of Kuhl's work made by Lotto *et al.* (1998) is that her proposed prototypes are pliable (their word, which we take to be the equivalent of our use of the word plastic). Lotto *et al.* (1998:3654) note that context can change the phonemic label a listener assigns to a given stimulus, and suggest that "category structure is a result of the input distribution itself and its relation to other categories", but they seem to expect that, for a phonetic prototype to exert a perceptual magnet effect, it must have a stable mental representation which cannot be easily manipulated. Neither we nor Kuhl dispute the first two points (e.g. Iverson & Kuhl 1995:560), but we do not think mental representations of phonetic categories must necessarily be acoustically stable. For example, Kuhl and Iverson (1995:146) suggest that prototypes might differ for gender.

One problem with interpreting these arguments is that the literature on the PME does not make it clear whether a so-called phonetic category is to be thought of as a phoneme. Our impression is that this issue has not been thought through: for example, the names of so-called phonetic categories are consistently given between slashes rather than square brackets, thus /i/ rather than [i], and most work has been done on isolated vowels. An isolated vowel, whose immediate phonetic context is silence, is arguably closest to the canonical quality of a phoneme, should one exist; but isolated vowels make a negligible contribution to real speech. In fact, experiments by Barrett (1997; Barrett Jones & Hawkins in preparation) confirm that Kuhl's proposed phonetic prototypes must be context-sensitive. Barrett showed that the best exemplars of synthetic /u/ vowels in /u/, /lu/ and /ju/ syllables must have different F2 frequencies if they are to sound natural, and that each best exemplar syllable produces a magnet effect that is specific to that syllable. The same F2 frequency in one of the other two contexts does not produce a magnet effect, but instead acts like a nonprototype. Notice that Barrett's findings

are consistent with our views on perceptual coherence, as well as plasticity of linguistic-phonetic categories.

In other experiments, both Barrett (1997) and Guenther *et al.* (1999) showed that the PME can be reversed: given an appropriate task, discrimination is increased around the best exemplar of a category and the magnet acts as a repeller rather than an attractor. For example, musicians who need to tune their instruments need to discriminate especially finely around the best (most in-tune) example of a particular note or chord. In contrast, speech understanding is presumably helped most by discriminating between but not within relevant categories.

In sum, no matter what theoretical status one wishes to give to the PME, the evidence is clear that the behavioural effect is demonstrable. As Guenther and Gjaja (1996) point out, logic, as well as neuroscientific and computational evidence, suggest that there is no need to postulate an actual prototypical representation in the brain, in that the same effect can be produced by representing memory of many similar instances and fewer more dissimilar ones, together with consideration of the demands of the task at hand. Equally, however, and central to our current argument, plasticity is fundamental to the concept.

4.4.3. *Summary: The nature of phonetic categories*

The neuropsychological research summarised above suggests that rich, context-sensitive structure may more nearly resemble the way words (and probably all aspects of language) are represented in the brain than the abstract, more independent strands that characterize much phonological-phonetic analysis. In other words, no category can be described independently of its context, and though the contexts can be combined to produce an apparently infinite range of finely-graded phonetic detail, the rules by which they are combined are relatively constrained. To extend the analogy of a blackbird seen in snow or sunshine: ornithologists learn a set of criteria necessary to identify it in either weather, from observation or explicit instruction; with experience, they develop a finely-tuned sense of how to weight the various criteria in different circumstances (contexts) so that the label *blackbird* is accurately used despite superficial differences. Although this knowledge may be gained faster with explicit instruction, what distinguishes the novice from the expert is that the expert has experience and knows what to attend to, which presumably allows rich and nuanced structuring of mental categories.

Phonetic categories, then, seem well conceptualised at present as part of a hierarchical structure of emergent categories, formed from repeated associations of meaning with patterns of experienced sound and other relevant sensations. To the extent that formal phonological and grammatical structure is represented in the brain, it is a by-product of this type of organisation, the union or intersection of related dynamic links. What is actually represented at the neural level depends on what tasks the individual habitually carries out; other aspects of linguistic analysis can be arrived at, but must be computed on the fly, possibly after some particular sensory input and its meaning have been linked. This view means that no two individuals will necessarily have identical phonetic categories.

5. *Attributes of Polysp in relation to some recent models of perception*

A number of recent models of speech understanding offer promising new approaches to some of the properties that we have identified as important, although no single extant model deals with them all. In this final section of the paper, we try to identify models whose strengths reflect views compatible with ours. This might help us, or others who agree with us, to combine the various strengths into a single, more comprehensive system. Some of the issues have been addressed in detail in the context of many different models, and we do not have space here for a full review of all the contributions that have been made. Instead, for several of the properties outlined in Sections 3 and 4, we discuss in some detail the one or two models that best address that property, making only brief reference to other related or contrasting approaches.

Most of the models we favour have in common that they are self-organising, in that the perceptually relevant units are not specified in advance but rather emerge dynamically in the course of processing information. In other respects, the various models we discuss differ computationally, and their processing assumptions may even contradict one another: something that is treated as a process in one model may be part of the architecture of another. These issues need not usually concern us, for the theoretical position we are developing here does not depend crucially on whether some important

property is modelled as a process or built in to the architecture of a model. Indeed this is an issue to which phonetics may not in general contribute a great deal, inasmuch as the status of memory and concepts is ultimately a question for neuroscience and molecular biology, or alternatively for philosophy.

5.1. *Modelling perception of temporally distributed information*

We have proposed that a model of speech understanding needs a mechanism for identifying both short-domain and long-domain events in the speech signal. It must be sensitive both to weak, subtle information as well as to clear information, bearing in mind that there are complex dependencies between the informational value of phonetic properties and the time domains over which they occur. In Polysp, the aim is not necessarily to arrive at abstract, time-free linguistic representations as soon as possible, but rather to allow the temporal structure of the speech signal to play a determining role in the course of processing. Rate variation, for instance, should not be modelled as time-warping of a sequence of underlying elements, but rather as a potentially meaningful restructuring of the temporal organisation of an utterance, which may result in an apparent reorganisation of segmental quality.

Most computational psycholinguistic models are neural networks, and developing any kind of internal representation of time in neural networks presents its own challenges (see Protopapas 1999 for a review). Perhaps as a consequence, comparatively little attention has yet been paid in neural net models of spoken word understanding to the specific constraints imposed by the temporal structure of phonetic events.

Short-domain events are comparatively easy to handle computationally given a sufficiently fine temporal resolution, such as spectra sampled at 10 ms intervals (cf. Klatt 1979; Johnson 1997; McClelland & Elman 1986), though 5 ms might be more realistic. The identification and integration of information that becomes available over longer time scales poses more of a problem, and certain kinds of model, such as Simple Recurrent Networks (e.g. Elman 1990; Norris 1992), have great difficulty integrating information over long time scales. The reason is that relationships between non-adjacent acoustic segments are subject to what can be termed a degrees-of-freedom problem. Local coarticulation and local aspects of rate variation can be modelled fairly easily because there are few degrees of freedom: that is, a property could be assigned to the current segment, the preceding segment or the following segment. But with greater distance between a segment containing some acoustic property and the segment that gives rise to that property (e.g. 6 phonemes if the first vowel in *it's a berry* exhibits /r/-colouring), the degrees of freedom required to correctly interpret that property increase substantially, at the expense of the model's explanatory power. Thus acoustic cues distributed over stretches of the signal as long or longer than a syllable may prove difficult to capture using recurrent neural networks (Nguyen & Hawkins 1999).

The difficulty with long-domain events can be attributed in part to the fact that word recognition is still often assumed to be a process which entails partitioning the signal into segments. A related problem is that most models take as the INPUT to speech processing either a series of static spectra (Klatt 1979; Johnson 1997) or auditory parameter values at successive time slices (e.g. Plaut & Kello 1999). Phonetic information is assumed simply to accumulate from one time slice to the next. Implicit in this view is the idea that speech processing proceeds uniformly (because the signal undergoes the same analysis process at each time step) and thus that processing responds only passively to the fact that the speech signal unfolds in time. Against this background, the integration of information forwards and backwards in time is bound to appear problematic.

In contrast, we suggest that processing itself is modulated or driven by the temporal nature of the speech signal, including its rhythm. That is, that the occurrence of particular acoustic properties, or the rate at which acoustic events occur, influences how subsequent information is processed. This is by no means an original idea: similar views are expressed by Lehiste (1972), by Whalen (2000), in Stevens' concept of landmarks (Section 3.4.1), and in Cutler & Norris's (1988) Metrical Segmentation Strategy, where strong syllables trigger lexical access attempts, and thus prosody is a major determinant of word segmentation.

The idea that temporal properties might drive speech processing has not been fully worked out, and hence is poorly implemented or unimplemented in computational models of speech

understanding. However, recent developments in the class of computational models developed in the context of Adaptive Resonance Theory (ART; Grossberg 1986) are promising in this respect. They are PHONET (Boardman *et al.* 1999), ARTPHONE (Grossberg *et al.* 1997), and ARTWORD (Grossberg & Myers 2000).

PHONET is designed to model phonetic context effects such as the effects of speech rate on consonant perception in a synthetic continuum between /ba/ and /wa/, and is applicable to speech rate variation in general. In outline, the model features separate, parallel auditory streams, one of which responds to transient and the other to sustained properties of the speech signal (e.g. formant transitions versus steady states); they store their inputs in parallel working memories. The model neurons in the transient channel operate together to detect rapid events like bursts and formant transitions, preserving information about the frequency, rate and extent of each transition: each is sensitive to frequency changes in its own particular narrow band of frequencies, and is subject to lateral inhibition so that their combined response reflects the relative degree of excitation in each frequency band. Activation in the transient channel is increased by a faster rate of spectral change, and decreased by a broader frequency extent.

The activation in the transient stream can modify the processing rate in the sustained stream, because greater activation in the transient stream (e.g. by faster formant transitions) increases the processing rate in the sustained channel. This cross-stream interaction ensures that the ratio of activation levels between the transient and sustained channels is maintained across syllables spoken at different rates. The logic is that, if there were no interaction between the channels, then when a CV syllable was spoken fast, it could produce a greater response in the transient channel, but no accompanying faster processing in the sustained channel, so that the ratio of transition-to-steady state durations would be different at different rates. Instead, the processing rate in the sustained channel is speeded up when the transient channel indicates that transitions are faster. The output is a ratio of transient-channel activation to sustained-channel activation, which is assumed to be mapped onto phonetic feature categories via a self-organising feature map.

The consequence is an invariant perceptual response for invariant relative durations; in this case, a relatively rate-invariant representation in working memory of the /b/-/w/ contrast. Since learning in the 'adaptive filters' of Adaptive Resonance Theory encodes the *ratio* of activations across working memory, phonetic categories in turn come to encode the ratios of sustained to transient information: a phonetic category boundary is enshrined in relative terms as relationships between distinct 'events', taking the whole context into account.

This appealing solution to an old problem has been tested on highly controlled synthetic speech. In natural speech, the ratios of transitions to steady states is not constant, and so some changes might be required. This might not prove too difficult. For example, even in simple CV syllables, rate changes are normally accompanied by changes in vowel quality, with concomitant changes in transition trajectories. Dynamic spectral differences of this type could presumably be related to the type of spectral difference Stevens uses to distinguish onset /b/ and /w/ (how abruptly the amplitude and spectral tilt change in the vicinity of the opening). Presumably the final decision should be based on sensitivity to a complex mix of phonetic detail.

ARTPHONE and ARTWORD do not explicitly consider the integration of acoustic information over time. The assumption is that this has already been performed in a prior stage, to produce what are termed 'phonemic item representations'. That is, in ART, speech input activates 'items' which are composed of feature clusters. Items in turn activate 'list chunks' in short-term memory, corresponding to possible groupings of features, such as segments, syllables, and words. Disappointingly from our point of view, it is only this latter process, with phonemic items as its starting point, which is modelled in ARTPHONE and ARTWORD. Even so, because ART models operate in real time, phonetic temporal structure still plays a part in these models, as we outline below.

Adaptive Resonance Theory distinguishes the rate of external input and various kinds of internal processing rate. Particularly important, as the name suggests, is the notion of resonance. Resonance is a positive feedback loop between items in working memory and list chunks. It involves non-specific top-down inhibition, and specific top-down confirmation of expected items. When listeners perceive speech, a wave of resonant activity plays across working memory, binding the phonemic items into larger language units and raising them into the listener's conscious perception.

The time scale of conscious speech is not equal to the time scale of bottom-up processing (the resonance evolves more slowly than working memory activation) nor to the rate of external input. The resonant dynamics include resonance ‘reset’, which prevents a resonance from continuing indefinitely, and may be triggered either by bottom-up mismatch, or by the resonance self-terminating. Self-termination is called ‘habituated collapse’, and represents the gradual cessation of resonance as synaptic neurotransmitters habituate to the same stimulus, so that they gradually stop transferring excitation between working memory and stored representations. An example is the decay of resonance at the end of an utterance because there is no new stimulation. Grossberg *et al.* (1997) suggest that geminates are perceived as two consonants because one resonance self-terminates and another begins, and that this is why the closure has to be longer for intervocalic geminates than for stop clusters with different places of articulation. Within a narrow temporal window, it may also be possible for a resonance to transfer seamlessly from one, smaller list chunk to another, larger chunk. This can happen if new information arrives while the support for the smaller chunk is beginning to weaken due to habituation (i.e. as the resonance is winding down), but before that chunk’s activation levels have fallen too low. Notice the congruence of these model processes with those of Pulvermüller (1999).

The nature of resonant dynamics in the ART models allows fine phonetic detail to play an important role in the grouping of items into list chunks that correspond to larger language units. The way this happens relates to a distinction drawn by Grossberg & Myers (2000) and Mattys (1997) between two kinds of phonetic evidence. One is informational, contributing to mapping of the acoustic stimulus into ‘phonemic items’, often over quite extended time periods; it is only implemented in PHONET. The second type, durational evidence, relates to the time of arrival of information such as silence, noise, and so on. Durational evidence can affect processing in a global way, because the time of arrival of a particular piece of information at a particular processing stage influences the competition between activated list chunks and the development of a resonance.

These ideas are illustrated in ARTPHONE by simulating the emergence of percepts of geminate stops versus two phonetically distinct stops. ARTWORD (Grossberg & Myers 2000) uses similar principles to simulate data on the contribution of silence and noise durations to percepts of *grey ship*, *great ship*, *grey chip* and *great chip* (Repp *et al.* 1978). Take, for example, the way in which increasing fricative noise duration produces a transition between percepts of *grey chip* and *great ship*. Both signals contain enough information for the perception of a stop-like sound, but the acoustic information groups differently, to yield a /tʃ/ percept in the former case, and a /t/ percept in the latter. The development of this grouping involves competition, which emerges at a slow enough rate to allow the initial competition between *grey* and *great* to be influenced by the later-occurring noise and the new competition it engenders between *great* and *chip*. When evidence for /tʃ/ is strong, at low noise durations, the chunks corresponding to *grey* and *chip* competitively team against *great*. At longer noise durations, the excitation of /ʃ/ has more time to develop, and is proportionally greater, so *ship* out-competes *chip*, thereby allowing *great* to receive greater levels of activation than *grey*.

In general, we are strongly sympathetic to the way in which ART models allow the temporal properties of the speech signal to affect the way the signal is processed. This contrasts with the approach taken in many models, where analysis of the signal is identical at each time step and phonetic information merely accumulates over time. Its principles bear some resemblance to those of the FLMP (Massaro 1998), although the two classes of model use different mathematical procedures. The fact that phonetic categories come to encode relational information is also very much in tune with our views. We do not have space here to discuss the other appealing properties of ART, which include its account of learning, its relative neural plausibility, and its very wide application in areas that include visual perception, memory and attention.

One less appealing property of ART models is that phonemic representations are obligatory, which at first sight seems incompatible with Polysp. We suspect, though, that the problem is fairly superficial. Because the range of contrasts examined so far is small, the objects of the simulations could as well be termed allophones as phonemes. More importantly, information is not thrown away when phonemic items are identified, because fine phonetic detail is passed up to subsequent stages as durational evidence. At the least, these models evince a strong sensitivity to syllable structure. In consequence, we do not view the ART approach as incompatible with nonsegmental phonological

representations, especially since the ART constructs of items and list chunks are not fixed responses corresponding to particular linguistic levels, but instead represent attractors of various sizes (Grossberg *et al.* 1997). Indeed, in typical ART resonance, longer chunks (e.g. words) mask smaller chunks (e.g. phonemes), so the largest coherent unit constitutes the focus of attention (Luce *et al.* 2000). This might be one way of dealing with the perceptual cueing function of long-domain phonetic resonance effects such as those described in Section 3.4.2 above.

We suspect, therefore, that the difference between Grossberg's and our views of phonemes is more terminological than substantive. This issue could become clearer if ART models investigated natural rather than synthetic speech, because in natural speech (unlike the synthetic stimuli in the studies simulated) spectral detail co-varies along with durational factors, as outlined above for rate.

We would also welcome a focus on long-domain properties other than speech rate, such as nasality, lip rounding, vowel-to-vowel coarticulation, and, at least for English, resonances due to liquids. If such longer-domain resonance effects had a status in the model equivalent to a Hebbian cell assembly, then a word containing /r/ might have its own cell assembly, and also an associated cell assembly that represented the resonance effect, which develops because the word has been associated with the long-domain resonance effect many times. When a putative /r/-resonance effect is detected, it would presumably raise activation levels of words containing /r/ in anticipation that one of them will be heard. So, when a particular word containing /r/ is heard, its activation level rises dramatically because it is congruent with the longer-domain resonance effect within which it occurs. In a broad sense, modelling these effects could be a first step towards making ART-type models explicitly polysystemic.

5.2. *Modelling the percept of rhythm*

We have suggested that a speech understanding model needs to allow local and sometimes subtle timing changes to influence the overall perceived rhythm of an utterance and the prominence of units within it, in ways which may convey grammatical, affective or other meaning. Such a mechanism might be naturally incorporated into future ART systems, given the attention paid to timing in these models. However, we also think a model should take into account the more general fact that people have a SENSE of rhythm in speech and other domains, especially since this may be an important bootstrapping tool for prelinguistic infants (for an overview, see Jusczyk 1997:ch. 6). No current speech understanding models can adequately explain the percept of rhythmic structure, but there do exist general accounts outside of speech research. One of the most appealing, for models of speech understanding, is Dynamic Attending Theory.

Dynamic Attending Theory (Jones 1976; Large & Jones 1999) addresses the puzzle that there exist all manner of dynamic events in which a clear temporal structure is apparent, yet the main beats are not isochronous. One example is the sound of a horse's hoof beats as it gallops faster and faster. People appear to apprehend stable rhythmic structures in such events, and they can perceive as meaningful the fluctuations in the periodicities that compose them, much as we have argued with regard to rate variation in speech. Dynamic Attending Theory explains this behaviour in a broadly Gibsonian framework, by postulating 'attending rhythms,' or internal oscillations which are capable of entraining to external events and targeting attentional energy to expected points in time.

The mathematical formulation of dynamic attending theory employs two entities: external rhythms and internal (attending) rhythms. External rhythms involve a sequence of related environmental events whose onsets can be localizable in time and which therefore define a sequence of time intervals. An example is the successive sounds of the horse's hoofs. An attending rhythm is internal to the perceiver. It is a 'self-sustaining oscillation' which can be thought of as generating an expectation.

When the individual attends to an environmental (external) rhythm, the phase and the period of the internal attending rhythm's oscillations become coupled, via entrainment, to those of the external rhythm, creating stable attractor states. The period of the oscillation reflects the rhythmic rate, or overall tempo, while the phase relationship between the coupled external and internal attending rhythms expresses the listener's expectation about when an (external) onset, or beat, should happen.

Changes in the beat are dealt with as follows. When the external rhythm is momentarily perturbed, then the phase relationship becomes desynchronised, but it can be adjusted to its original

phase when the perturbation is not large. So when the only perturbations in an ongoing external rhythm are random and relatively minor, entrainment can continue through small phase adjustments of the attending rhythm. However, when the external rhythm changes to a new overall rate, the attending rhythm must change its period rather than continually adjust its phase. If the period does not change to reflect the new rate, then all succeeding external onsets will be expected at the wrong time and heard as either late or early, because the phase relationship will be wrong. Thus, while period coupling captures the overall rate, phase coupling prevents synchrony between the external and attending rhythms from being lost when the external rhythm is momentarily perturbed, and can capture the disparity between an expected onset and an actual onset.

Within each cycle of an attentional rhythm, the allocation of attentional energy is also modelled with a variable termed 'focus of attention'. There is a pulse of attentional energy, whose locus is determined by phase and whose extent is determined by focus of attention. It contributes an expectancy region where attentional energy is nonzero: a narrow pulse reflects an expectation that an event will take place within a narrow time range. A broader pulse reflects greater uncertainty about when an external event will happen. Focus allows the model to make predictions about the noticeability of time-changes. For instance, a large deviation from expectancy is more likely to be noticed than a small one, and any deviation is more likely to be noticed if attention is focused narrowly rather than broadly. The distinguishing principle is the same as the distinction between narrow-band and wide-band spectrograms. Wide-band spectrograms are made with a filter that has a shorter time window than that used to make a narrow-band spectrogram. Because the window is short, it allows more precise temporal definition—which is why wide-band spectrograms have sharper 'edges' between acoustic segments than narrow-band spectrograms.

Since most external rhythms are complex, involving ratio relationships among the phases and periods of their different temporal components, the model likewise employs multiple attending rhythms which are coupled to one another to preserve phase and period relationships. These relationships offer the potential for dynamically expressing relational information about the complex rhythm's underlying temporal form. Evidence from maturational changes in motor tapping and tempo discrimination is interpreted in the context of the model to suggest that the coupling of multiple oscillators develops with age and experience (Drake *et al.* 2000).

Among the appealing properties of Dynamic Attending Theory for our purposes are its explicit rejection of the view that temporal events like speech are time-warped sequences of underlying discrete elements whose serial order is of paramount importance. Instead, like ART, Dynamic Attending Theory incorporates a notion of internal psychological timescales which attune to those in the external environment. This broadly Gibsonian view is compatible with contemporary understanding of cyclical behaviour in biology, such as adaptation of circadian cycles to daylight hours (e.g. Sawyer *et al.* 1997; Kyriacou *in press*). As a result, both Dynamic Attending Theory and ART allow a prominent role for expectations and for informed (i.e. systematic, knowledge-driven) fluctuations in degree and focus of attention. This contrasts with psycholinguistic models which minimise the role of feedback (e.g. Norris *et al.* 2000), because, while expectations about some rhythms might not involve knowledge about the world, it seems to us that expectations about speech and many other rhythms (including music) must involve knowledge. The percept of rhythm in music is at least partly culturally determined (e.g. Stobart & Cross 2000), and the percept of rhythm in speech is at least partly language-specific.

Another appealing property of Dynamic Attending Theory is its simultaneous focus on different time windows, each suitably precise for the domain. Links between focus and expectations suggest that the better the listener knows the speaker, for example, the better he or she can predict upcoming events, and hence the faster and easier it will be to understand the message. The model thus has particularly interesting implications for making maximally efficient use of cues to segmental identity that are found in rhythmically predictable locations, such as landmark features and cues to place of articulation at the edges of syllables. When you expect these events, you adjust the focus of attention to a suitable time window to allow you to find them. When you do not expect or need them, you can transfer attention elsewhere.

Given the complexity of speech rhythms and their interactions with linguistic units of various kinds, we do not know to what extent Dynamic Attending Theory and mathematical formulations thereof could be adopted wholesale to form the basis of a speech understanding model. Its insights

may be best viewed as a guiding principle showing what can be achieved when the idea of rhythmically-driven attention is explicitly built into a model. However, its similarities with principles used in other speech models such as Tuller *et al.*'s (1994) dynamical model of speech sound categorisation, described in Section 5.4 below, encourage us to suggest that principles used in Dynamic Attending Theory have a promising future in a polysystemic model of speech understanding.

5.3. *Modelling adaptation to new speakers and rates: episodic representation*

Adaptation to new speakers and rates is naturally accommodated in episodic theories of perception (Pisoni 1997b, Goldinger 1997) in which details of the speech signal that have traditionally been viewed as incidental are considered integral to the linguistic representation. Episodic theories see mental linguistic representation as an agglomeration of highly diverse language-related experiences, which is compatible with our ideas on multi-modal memory.

Exemplar models (e.g. Hintzman 1986, 1988; Nosofsky 1988) show the strongest commitment to episodic theories, retaining a separate memory trace for each stimulus experienced. When a new stimulus is matched against these multiple traces, it may excite anything from a highly specific response (as would be the case e.g. for an unusual word in a familiar or recently-heard voice), to an entirely generic one (e.g. to a common word in an unfamiliar voice). In these models, abstraction thus occurs dynamically during perception, rather than at storage.

There are growing numbers of exemplar models for speech, some of which have as a main objective the modelling of particular phonetic phenomena (e.g. Johnson 1997; Lacerda 1995), while others are not intended to model the details of the speech signal and use purely schematic, arbitrary representations of spoken words (e.g. Goldinger's adaptation of Hintzman's MINERVA 2: Goldinger 1997, 2000). Both approaches are valuable, though for the long term we stress the importance of an input that is phonetically realistic. In this paper, however, we concentrate on the work of Goldinger, because his simplifications have allowed a clear focus to develop on the global contours of adaptation to speakers over time.

Goldinger (1997) found support for MINERVA 2's prediction that some stimuli excite a generic response, others a specific response. Subsequent investigations (Goldinger 2000) have mapped some of the task-specific ways in which these responses can change with repeated exposure to isolated word tokens. For instance, the effect of word frequency diminishes over time in a recognition memory test, but increases when the measure is imitation of a particular speech token. The interpretation is that the word frequency effect decreases in recognition memory tests because more traces have accumulated, which reduces the discrepancy between high- and low-frequency words, presumably specifically in this task. Conversely, it increases in imitation tasks because all the memory traces are of the particular token of the low-frequency word, so that production comes to reflect that token more and more closely.

So far, there has been little suggestion of how an appropriate set of episodic traces might be activated; as we have suggested, this process is unlikely to be as simple as a direct mapping of raw acoustic information onto stored traces, time-slice by time-slice, and indeed Goldinger (1997) emphasises that context will play a key role. Moreover, it is difficult to see how a simple exemplar-based matching process could convey some types of linguistic information, such as that provided by long-domain resonance effects: this type of acoustic-phonetic information may require some degree of abstraction in order to integrate it over long time domains (see Section 4.4.1 above). In any case, Goldinger's rationale for working with MINERVA 2 is primarily that, since it is a strict exemplar model, predictions derived from it offer a very stringent test of the extent to which episodic memory is implicated in linguistic processing. Since the model does predict a number of experimental results, it does strongly suggest that episodic memory plays a role in speech understanding. Of course, models that learn more abstract representations than MINERVA 2 can capture episodic memory for speech, provided that they learn in a flexible enough way that when the system needs to be able to make extremely fine discriminations, the details of individual exemplars can be retained. For instance, ART formalisms are compatible with Goldinger's ideas, since abstract prototypes and 'exemplar prototypes' can co-exist in these systems (Grossberg 2000b). Accordingly, even if some properties of speech turn out to be best explained not by a strict exemplar model, but perhaps by something more

like ART, simulations with MINERVA 2 can make a useful contribution by showing where such a model should be constrained.

5.4. *Modelling the emergence of linguistic categories*

In the light of the overwhelming evidence that linguistic categories are plastic, a model should be able to produce appropriate changes in a particular phonetic percept as factors such as phonetic context and method of stimulus presentation vary. While many neural networks can simulate the influence of phonetic context, their success is usually due to training in which the network is told the correct categorisation for stimuli occurring in a variety of contexts. However, this type of training seems unlikely to be able to account for why repeated presentation of a stimulus may give different results depending on the order of presentation or the composition of the stimulus set². One model which seems promising in this regard is the dynamical model of Tuller *et al.* (1994).

Tuller *et al.* (1994) note that speech sound categorization is nonlinear in that when a control parameter is varied over a wide range, often no observable change in behaviour results, but when a critical value is reached, behaviour may change qualitatively or discontinuously. (Note the congruence with quantal theory.) To investigate the dynamics of this process, they varied the duration of a silent gap between a natural utterance of /s/ and a synthetic vowel /ei/ in a *say-stay* continuum like those typically used in categorical perception experiments. When the stimuli were presented in sequential order (the duration of the silent gap either increasing then decreasing, or first decreasing then increasing), it was relatively rare for listeners to hear *say* switch to *stay* (with increasing gap duration) at the same gap duration as *stay* switched to *say* (with decreasing gap duration). Instead, and as expected, listeners usually switched either relatively late in the series (response perseveration, or hysteresis) or early (enhanced contrast).

The dynamical model developed to account for these well-known psychophysical context effects uses equations of motion to describe the temporal evolution of the perceptual process, especially factors affecting stability *versus* plasticity of perceptual classes. In broad outline, each perceptual category is represented as an attractor. The attractor's range of influence is represented as a 'basin of attraction' whose breadth and depth can vary with stimulus conditions. Deeper and narrower basins are associated with more stable percepts. An acoustic parameter (in this case the duration of the silent gap) influences the perceptual dynamics. When the value of this parameter is low, only a single attractor exists, corresponding to the percept of *say*. This situation is maintained for increasing values of the parameter corresponding to gap duration (although the basin within which attraction occurs becomes progressively shallower and narrower), until a critical value is reached, at which point an additional attractor corresponding to the percept of *stay* is created. Both attractors coexist until a second critical value is reached at which the attractor corresponding to *say* disappears. Then, once again, only a single attractor is present, with the basin of attraction becoming deeper and wider with increasing gap duration. When both attractors coexist, random disturbances (corresponding to fatigue, attention, boredom, and so on) can produce spontaneous switches between the two percepts.

The effects of hysteresis and enhanced contrast are accounted for by including in the model a term corresponding to the number of perceived repetitions of a stimulus, so that when the listener

² Many of the experiments that demonstrate these properties of speech perception were conducted within the context of the early search for auditory neural feature detectors in the brain comparable to those found for vision. That approach was abandoned when it became clear that central, integrative processes seem to be at least as important to speech perception as relatively simple, automatic responses to the presence or absence of auditory properties (see Darwin (1976) and Eimas & Miller (1978) for reviews). At that time, demonstrations of phonetic trading relations were considered to provide one of the most damaging blows to the idea that auditory feature detectors would prove to underlie phonological features, especially if the trading relations were between complicated combinations of acoustic properties that tend to co-occur in natural speech because of the way the vocal tract works, rather than because of any acoustic similarities. However, the basic idea seems broadly compatible with more recent thinking about neuronal behaviour as characterised by complex interactions between different functional groupings of cell units, although the process is undeniably more complicated than had originally been hoped.

hears many repetitions of a stimulus (or of different stimuli that are interpreted as belonging to the same category) the location of the boundary for that category shifts. The effects of this term are, in turn, sensitive to cognitive factors (represented in the model as the combined effects of learning, attention, and experience). Among the resulting predictions are that factors such as learning and experience with the stimuli will make listeners less likely to cling to the initial perceptual state (hysteresis), and more likely to undergo an early switch in percept (enhanced contrast).

The appeal of this model is that it accounts for the emergence and plasticity of linguistic categories by making their inherent relational nature fundamental to perception, and allowing for multiple influences on the final, context-sensitive percept. This contrasts with the more traditional approach of listing the attributes of a particular category. We admire the focus on the temporal evolution of the percept, and on the influence of recent sensory experience on perceptual decisions (cf. also Ding *et al.* 1995). The model focuses mainly on properties of a restricted stimulus set—the way the stimuli are presented—since this offers an interesting window on the stability of percepts. Although we are more interested in the effects of a more varied phonetic context, as found in normal listening situations, we expect that modelling these effects must similarly take into account rich perceptual dynamics, as some of Tuller's later work suggests is the case (e.g. Rączaszek *et al.* 1999).

There are parallels in Tuller's model with, on the one hand, Goldinger's investigations of the effects of repeated exposure to identical or very similar tokens, and on the other hand Grossberg's modelling of the dynamic emergence of categorisation judgements. Where Tuller *et al.*'s model differs from Grossberg's is in its explicit focus on the plasticity of categories defined in terms of immediate relations in the sensory signal (as well as in terms of other factors such as fluctuations in attention). In contrast, the focus of an ART model such as PHONET is on achieving an invariant percept for any given ratio of activity between the model's two processing streams, which implies that the model consistently produces the same response given successive identical inputs. As we have made clear, ART systems do have complex learning dynamics. For example, each time a resonant state occurs in the system, learning takes place, which presumably means that multiple repetition of a stimulus modifies the structure of the category onto which that stimulus is mapped. It is certainly possible, then, that the important aspect of plasticity captured by Tuller *et al.*'s model could be simulated in the ART framework, but this remains to be explored.

5.5. *Modelling acquisition*

Language acquisition plays an important role in the kind of model of speech understanding represented by Polysp. In part this is because listeners' capacity to learn about speech does not stop once a language has been acquired, so that, for us, acquisition and adult behaviour need not be radically discontinuous. Polysp assumes that the basic processes are similar, although the adult's greater experience means that actual patterns of neural activation will often differ, so that the details of how a particular piece of knowledge is learned and structured may also differ between adults and children. Understanding of how babies and young children form speech sound categories should thus provide insights into the nature of their plasticity, and may provide insights into how adults use them.

Another reason why acquisition must feature in a phonetically-sensitive polysystemic model of speech understanding is that the processes presumed to underlie the development of speech sound categories and phonology are the same as those underlying the developing organisation of grammar and meaning. This view is supported by increasing evidence that sensitivity to the distribution of sound patterns in the speech signal provides a basis from which infants learn about properties of language, such as the nature of different word classes (e.g. Smith 1999). Thus the infant and young child constructs all his or her knowledge of the language in parallel from the primary speech input, which is usually connected speech. Just as speech perception or understanding does not involve an orderly sequence in which the signal is mapped onto phonological units which are then mapped onto grammatical units and then onto meaning, so there is no learning sequence of the sort in which phonology provides the input to grammar which in turn provides the input to meaning. Indeed, logically, it seems more likely that meaning and the phonetic signal are closely related, and that grammatical and phonological systems arise later, as frequent patterns of activation become associated into common structures, or functional groupings. Looked at this way, phonology and grammar are, as it were, by-products of the association of meaning with particular sound patterns.

One recent development which reflects well some of our views on acquisition is Plaut & Kello's (1999) distributed connectionist model of the emergence of phonology in infants from the interplay of speech production and comprehension. In the model acoustic input is mapped directly onto a semantic level, via a level of phonological representations which are not predefined, but are learned by the system under the pressure of understanding and producing speech. Accessing meaning is of central importance, therefore, and phonological representations emerge to support this goal by recoding time-varying acoustic input into a more stable format. The idea that phonological representations are emergent tallies well with the neuropsychological evidence discussed by Coleman (1998, ms), and Pulvermüller (1999) in Section 4.4.1 above. Unfortunately, in order to make the model computationally tractable, Plaut & Kello were forced to simplify the phonetic parameters so that much of the important detail is neglected.

Jusczyk's WRAPSA model (Word Recognition and Phonetic Structure Acquisition, Jusczyk 1993, 1997) includes much more fine phonetic detail than Plaut & Kello's model can, because it has not been computationally implemented. WRAPSA is potentially richly polysystemic, synthesizing work from a number of areas, including phonetic temporal structure, linguistic categories, attention and memory, and assuming that input patterns derived from weighted auditory properties are matched against stored traces in the manner proposed by Hintzman (1986, 1988). In these respects it is immensely appealing, although its focus on word identification as a final goal seems a limitation compared with Plaut & Kello's attempt to map the input directly onto meaning.

WRAPSA's treatment of phonetic temporal structure is particularly promising in that it combines attention to auditory fine detail with integration of information over a window larger than the segment. The auditory analyzers are assumed to carry out some temporal tagging of features that co-occur within the same syllable-sized unit. (A syllable-sized unit has a local maximum between two local minima in the amplitude contour.) We have emphasised, on the other hand, that when adults understand speech, they must integrate information over multiple time domains, some of which are longer than the syllable. It seems probable that the domain or range of domains used by listeners should change during maturation (cf. Drake *et al.* 2000). It is interesting in this regard that infants may use a smaller window for processing information than adults, due to limited attentional and memory capacities, and there is evidence that 'starting small' could help infants identify perceptually relevant units in the input (cf. Elman 1999).

WRAPSA also shares with Polysp the properties that potential lexical units are not initially represented in fully phonetically specified form, and are not decomposed into phoneme-like segments, although temporal relationships between the acoustic properties within each unit are represented. For instance, infants' early perceptual representations might encode enough information for some of the relational contrasts of the prosodic hierarchy to be abstracted, although the entire set of contrasts would not yet be in place. Jusczyk proposes that early representations are syllabic, that is, focused around the middle of a fully-specified prosodic hierarchy of an utterance. We would add to this proposal the speculation that highly reliable acoustic information is more likely to be represented than less clear information, and that this distinction can cut across that of syllabic identity. For example, a strong syllable can be identified with high certainty even when its component phones are only partially identified; on the other hand, it may not always be clear, at least to novice listeners, how many unstressed syllables are present, especially in a sequence of unaccented syllables. For example, the phrase (*of*) *course it is* often has no voiced /t/ for *it*, and instead that syllable is marked by durational changes and possibly some palatalisation of the other sounds in that vicinity e.g. [k^hɔs:tɪz]; these preserve some of the rhythmic and critical segmental features of the carefully-spoken phrase. Likewise, a strident fricative can be identified with high certainty when there is only low certainty about which syllable it is a member of, or, as in this example, about how many lexical syllables it represents. Here, the integration between rhythmic and segmental structures will help make the phrase readily understandable.

In the light of the evidence (DeCasper & Fifer 1980; Jusczyk *et al.* 1992; Jusczyk *et al.* 1993) that infants in the first year of life retain detailed memory for the speaker's voice, which is normally classed as nonlinguistic information, Jusczyk (1997:227) suggests that it is paradoxical that phonetic representations might be less than fully detailed. He proposes instead that infants' capacities for representing the internal structure of words and syllables may not yet be fully developed; this might

be the case if full phonetic representation required links between details for production as well as perception. As should be clear by now, our explanation is rather different. Nonlinguistic and linguistic representations of utterances are not fed by separate strands of information; instead, the same phonetic information is implicated in both. On this view, even incompletely specified phonetic representations are sure to contain information about the speaker, because voice information is bound up as part of the signal and is reflected in many, if not all, phonetic parameters.

In summary, we suggest that the processes represented by Polysp for an adult understanding speech are not fundamentally different from those of an infant who only partially understands the same speech, or even fails to understand it. The main difference between adults and babies in this respect is that the adult has a lot more experience—and hence available structures—within which to place the newly heard utterance. Polysp is compatible with aspects of Plaut & Kello's phonetic-semantic mapping in a self-organising connectionist framework and with Jusczyk's more clearly polysystemic, phonetically-detailed approach that builds on episodic memories. Polysp adds to these models a rather strong stand on the status of linguistic and non-linguistic sensory information: that it is only classifiable as linguistic or non-linguistic by virtue of its current functional role, which is determined not so much by the quality of the sensory information as by the listener's attention. This is one aspect of the plasticity fundamental to Polysp. One implication of this approach is that individuals are assumed to differ in how they structure the same information. In consequence, an individual's linguistic system need be neither fully specified nor correct in the linguist's sense: as long as it allows the individual to understand meaning correctly, the details of how that meaning is arrived at are immaterial, from the point of view of the success of the communication. That is why it is perfectly possible to believe throughout adulthood that the last syllables of *could've* and *would've* contain *of* rather than *have*. They sound as if they end in *of*, and there is no reason why the underlying message should be misunderstood because the construction is thought to contain *of*. It is also why most educated speakers of English do believe that these phrases end in *have*, because what they are taught to write in school (*have*) changes the way they structure this particular spoken unit.

5.6. *Summary: attributes of Polysp in existing models of perception*

Each of the models we have discussed addresses a different aspect of the ideas fundamental to Polysp: the complex temporal distribution of acoustic information relevant to linguistic distinctions, the role of rhythmically-governed expectations and knowledge in focusing listeners' attention on particular points in the sensory signal, the inherently dynamic nature of speech understanding, including the plasticity of linguistic-phonetic categories and their sensitivity to context, experience, and linguistic development. Most of them go much further than Polysp in that they are computationally implemented, and in general, we have commented more on their principles than their details. Importantly, though, our focus on phonetic detail does lead us to rule out certain approaches and favour others, especially those that have a realistic way of processing temporally distributed information.

Not only do these models offer a rather coherent set of principles compatible with Polysp's basic principles, but the principles themselves are usually thought to have some reality in how the brain works. For example, a recent handbook (Arbib 1995) includes articles that between them propose neural mechanisms for the general and detailed principles we have discussed, such as modified Hebbian cell assemblies, phase locking between external and internal rhythms, different responses to faster and slower events, coupled oscillators, attractors, feedback, and so on. Equally, it is clear from this same book that we cannot be certain that any of the proposed mechanisms functions as described. That we are beginning to have access to them is encouraging, however.

Of particular interest to us is the fact that the models we have discussed in this section rely less than many traditional models on using linguistic units such as phonemes to define their processing stages, and make more use of constructs with independent motivation from other disciplines such as neurobiology, dynamical systems theory or more general aspects of behaviour such as rhythm perception. It is because of their comparative agnosticism about the units of speech perception, combined with their rich processing dynamics, that models like these are relevant to Polysp and *vice versa*.

Why then is it worth developing Polysp? As a linguistic model, Polysp complements dynamic, self-organising and exemplar models such as these because it makes strong claims about the wealth of linguistic information that is systematically available from the sensory signal and the way this linguistic information is organised; it stresses more than most models that linguistic processing must draw upon fine-grained experience and expectations. Yet, on phonetic as well as cognitive grounds, Polysp also proposes that there is no one way to understand a speech signal: polysystemic linguistic structure can be identified by many routes, in many different orders or in parallel.

6. Concluding remarks

Polysp offers a framework that can potentially hold good for all the varieties of speech that speakers standardly produce and understand, and that is able to include all the linguistic and nonlinguistic knowledge that listeners use when they understand what another person is saying. Let us take an example to see how this might work. Most native speakers of English have an impressively wide variety of ways of conveying the meaning of *I do not know*. The most common forms probably range between *I don't know* and *dunno*, both of which can be pronounced in a number of different ways. However, there are many other variants, the most extreme forms of which can only be used in particular circumstances. For example, it is hard to say the fully expanded form *I do not know* without conveying some degree of exasperation. An even more extreme form has pauses between the words (*I...do...not...know*) and (in most cases) is so rude that it can only be used when the listener does not seem willing to accept that the speaker really does not know. At the other extreme, it is possible—again, only in the right circumstances—to convey one's meaning perfectly adequately by means of a rather stylized intonation and rhythm, with very weak segmental articulation, ranging between something like [ã³n:ə⁰] and [ã³ã³] (intonation not marked). This type of utterance could allow successful communication between relaxed family members, for example when A asks B where the newspaper is, and B does not know, but does not feel that she needs to stop reading her book in order to help find it. Notice that the intonation pattern alone is not enough: at least the vowels must be there ([m]s will not do), and the vowels must start more open (and probably more fronted) than they finish, just as in the more clearly-spoken utterance, so that, at least in this situational context, [ã³ã³] is nonsense whereas [ã³ã³] is not.

These minimalist utterances are understood with no difficulty in the type of situation we have described, but they are unlikely to be understood outside them. They can be explained within the Polysp framework by postulating a functional neuronal grouping (a cell assembly for instance) linking these speech variants to meaning and to the affective, socio-cognitive and other properties of the experienced situational sensations. That meaning, in turn, is part of another functional cell assembly that includes representations of other ways of expressing the same meaning in other circumstances. Directly, or indirectly via other links, these other assemblies will include representation of the actual lexical items *I do not know* or *I don't know* and their individual and sentence meaning. In essence, these cell assemblies describe relationships between phonetic forms and connotative meanings. One could speculate that denotative meaning might lie at the intersection of all these variants on connotative meaning. Another way to express this is that the phonetic form reflects the perceived pragmatic meaning of the utterance, which is represented at the neuronal level by nuances in the composition or functioning of particular cell assemblies.

In this type of system, no one type of information need be more important than any other type, and the same result can be arrived at from different starting points. When B says [ã³ã³], A learns a great deal more than simply that B does not know where the newspaper is. Equally, because A sees that B is deeply involved with her book, he will be less likely to interpret her [ã³ã³] as a sort of dysarthric grunt preparatory to a more helpful utterance. In other words, the meaning can be arrived at by linking the perceived multi-faceted (i.e. detailed) situation with perceived sound in a mutually reinforcing way, no matter which is attended to most at the outset, or which type of neuron fires first. What is crucial is that the experience is coherent, which is to say that its components must bear appropriate relationships to one another.

These arguments seem uncontroversial in that they are basically common sense. We suggest that the same type of argument can—and indeed should—be made for how we understand more

standard forms of speech. We have already made the case that components of a speech signal must be coherent, bearing the right relationships with one another. We have shown that the acoustic fine detail of the sound signal is a rich source of information about all aspects of the message, and argued that understanding will be most efficient when these details are attended to rather than when they are abstracted from early in the perceptual process. These views suggest that speech understanding is ultimately rooted in multi-modal episodic memory and subject to continuous adjustment (plasticity) throughout life.

We have also suggested (Section 5.5) that the perceptual system may access meaning from speech by using the most salient sensory information from any combination of levels of formal linguistic analysis, and from this create a representation that may be incomplete in some respects. By implication, (a) although some units might normally be more fundamental than others (stressed syllables are good candidates, for example—cf. Cutler & Norris 1988; Jusczyk 1997; Greenberg 1999) there is no one basic unit of speech analysis; and (b) there must be perceptual processes which are distinguished by the duration of the phonetic information they process: some process short-time events, while others process more long-lasting information (Section 3.6). This requirement in a model of speech understanding brings it closer to models of syntactic processing and contrasts sharply with the many computational models of speech perception which use a standard time-window and rather abstract input, thereby risking distortion of the relative importance of different types of processing. Nearey engagingly endorses this criticism of most models, including his own: “For example, in my models, instead of temporally-evolving acoustic waveforms, I start with neatly packaged, pre-parsed ‘cues’ which feed just the right nodes in my models in just the right way. Furthermore, my models are simple static-pattern recognizers where all evidence is presented instantaneously, rather than emerging through time.” (2000:343).

Another implication is that there is no obligatory order a listener must follow in order to understand speech, and the process of speech understanding can be circular. The physical signal will usually dominate and drive perceptual decisions, and processing itself is modulated or driven by the temporal nature of the speech signal, including its rhythm; but listeners must look for grammar as well as lexical identity in speech, and (normally) work out grammar and meaning simultaneously from the spoken signal. Phoneme or allophone strings will not allow this unless they include the type of polysystemic context-sensitivity intrinsic to Polysp, rather than just sensitivity to other phonemes in the immediate vicinity. And if the nature of their context-sensitivity is polysystemic as in Polysp, then phonemes become irrelevant to the process of understanding meaning from spoken utterances. Similarly, meaning might sometimes be understood before individual words and their grammatical relations are identified, and certainly before a complete phonological representation is arrived at (cf. Warren 1999; Grossberg 2000a). For this reason, one remembers the meaning (or perceived meaning) but not always the actual sounds or lexical items that were spoken, just as has been demonstrated for grammar (Sachs 1967). Moreover, individual differences stemming from ‘incomplete’ or idiosyncratic phonological and grammatical structures are likely to be the norm rather than unusual.

In sum, we have tried (a) to show that the search for a better model of speech understanding depends partly on acknowledging that the phonetic signal is a rich source of information about a large number of different types of linguistic and nonlinguistic information, (b) to outline the types of attributes a more comprehensive model should have and (c) to suggest processes or mechanisms by which these attributes can be built into a model. In common with a number of other recent theories, Polysp assumes that, rather than moving as fast as possible towards an abstract representation of speech, listeners retain the richness of acoustic information at least until the meaning has been identified. We are thus suggesting that the speech signal could be seen as an integral aspect of meaning rather than a relatively uninteresting carrier of meaning; and that phonetic categories behave just like other linguistic categories: they are emergent, dynamic, plastic throughout life, and not simply context-sensitive, but can only be discussed in relational terms.

Sarah Hawkins and Rachel Smith
Department of Linguistics
University of Cambridge
Sidgwick Avenue
Cambridge CB3 9DA
United Kingdom

sh110@cam.ac.uk rhs20@cam.ac.uk

Acknowledgements

We acknowledge with gratitude the intellectual debt we owe to John Local and Richard Ogden for sharing their knowledge and ideas about phonetics, and for many hours of patient exposition of linguistic-phonetic structure; and to Noël Nguyen for collaborating in experiments and the development of theory about the role of systematic acoustic-phonetic fine detail in spoken word understanding. We also thank Friedemann Pulvermüller for helpful and stimulating discussions. This work was supported in part by EPSRC grants GR/N19595 and GR/L53069 to the first author, by Swiss National Science Foundation grant 83BC-056471 to Noël Nguyen and the first author, and by an AHRB research studentship to the second author.

7. Bibliographical references

- Abeles, Moshe (1982), *Local Cortical Circuits: An Electrophysiological Study*, Berlin, Springer-Verlag.
- Abeles, Moshe (1991), *Corticonics: Neural Circuits of the Cerebral Cortex*, Cambridge, Cambridge University Press.
- Abu-Bakar, Mukhlis & Nick Chater (1995), "Time-warping tasks and recurrent neural networks", in Levy, Bairaktaris, Bullinaria & Cairns (1995:269-288).
- Alfonso, Peter J. & Thomas Baer (1982), "Dynamics of vowel articulation", *Language and Speech* 25:151-173.
- Allen, Joseph & Mark S. Seidenberg (1999), "The emergence of grammaticality in connectionist networks", in MacWhinney (1999:115-151).
- Arbib, Michael A. (1995), *The Handbook of Brain Theory and Neural Networks*, Cambridge, MA: MIT Press.
- Arbib, Michael A. (2000), "The mirror system, imitation, and the evolution of language", to appear in Chrystopher Nehaniv & Kerstin Dautenhahn (eds.) *Imitation in Animals and Artifacts*, MIT Press.
http://www-hbp.usc.edu/_Documentation/papers/arbib/LanguageEvolution1-00.doc.
- Assmann, Peter F., Terrance M. Nearey & John T. Hogan (1982), "Vowel identification: Orthographic, perceptual, and acoustic aspects", *Journal of the Acoustical Society of America* 71:975-989.
- Barrett, Sarah E. (1997), *Prototypes in Speech Perception*, unpublished Ph.D. dissertation, University of Cambridge.
- Barrett Jones, Sarah E. & Sarah Hawkins (in preparation), "The sensitivity of perceptual magnet effects to phonetic context".
- Beckman, Mary E. (1986), *Stress and Non-stress Accent*, Netherlands Phonetic Archives 7, Dordrecht, Foris.
- Beddor, Patrice Speeter & Rena Arens Krakow (1999), "Perception of coarticulatory nasalization by speakers of English and Thai: Evidence for partial compensation", *Journal of the Acoustical Society of America* 106:2868-2887.
- Beddor, Patrice Speeter & Handan Kopkalli Yavuz (1995), "The relation between vowel-to-vowel coarticulation and vowel harmony in Turkish", in Kjell Elenius & Peter Branderud (eds.), *Proceedings of the XIIIth International Congress of Phonetic Sciences*, Stockholm, KTH and Stockholm University, 2:44-51.
- Beddor, Patrice Speeter, Rena Arens Krakow & Louis M. Goldstein (1986), "Perceptual constraints and phonological change: A study of nasal vowel height", *Phonology Yearbook* 3:197-217.
- Bell-Berti, Fredericka (1993), "Understanding velic motor control: Studies of segmental context", in Huffman & Krakow (1993:63-85).
- Benguerel, André-Pierre & Helen A. Cowan (1974), "Coarticulation of upper lip protrusion in French", *Phonetica* 30:41-55.
- Best, Catherine T. (1994), "The emergence of native-language phonological influences in infants: A perceptual assimilation model", in Goodman & Nusbaum (1994:167-224).
- Best, Catherine T. (1995), "A direct realist view of cross-language speech perception", in Strange (1995:171-204).
- Boardman, Ian, Stephen Grossberg, Christopher Myers & Michael Cohen (1999), "Neural dynamics of perceptual order and context effects for variable-rate speech syllables", *Perception and Psychophysics* 61:1477-1500.
- Bregman, Arthur S. (1990), *Auditory Scene Analysis, The Perceptual Organization of Sound*, Cambridge, MA, MIT Press.
- Browman, Catherine P. & Louis Goldstein (1989), "Articulatory gestures as phonological units", *Phonology* 6:201-251.
- Browman, Catherine P. & Louis Goldstein (1990), "Gestural specification using dynamically-defined articulatory structures", *Journal of Phonetics* 18:299-320.
- Browman, Catherine P. & Louis Goldstein (1992), "Articulatory Phonology: An overview", *Phonetica* 49:155-180.

- Buxton, Hilary (1983), "Temporal predictability in the perception of English speech", in Cutler & Ladd (1983:111-121).
- Carterette, Edward C. & Morton P. Friedman, eds. (1976), *Handbook of Perception*, vol. 7, New York, Academic Press.
- Clumeck, Harold (1976), "Patterns of soft palate movements in six languages", *Journal of Phonetics* 4:337-351.
- Cohen, Antonie, J.F. Schouten & Johan t'Hart (1962), "Contribution of the time parameter to the perception of speech", in *Proceedings of the IVth International Congress of Phonetic Sciences*, The Hague, Mouton: 555-560.
- Coleman, John S. (1994), "Polysyllabic words in the YorkTalk synthesis system", in Keating (1994:293-324).
- Coleman, John S. (1998), "Cognitive reality and the phonological lexicon: A review", *Journal of Neurolinguistics* 11:295-320.
- Coleman, John S. (ms), "Phonetic representations in the mental lexicon", under review for publication in Jacques Durand & Bernard Laks (in preparation), *Phonology: From Phonetics to Cognition*.
- Connell, Bruce & Amalia Arvaniti, eds. (1995), *Phonology and Phonetic Evidence: Papers in Laboratory Phonology IV*, Cambridge, Cambridge University Press.
- Couper-Kuhlen, Elizabeth (1986), *An Introduction to English Prosody*, London, Edward Arnold.
- Cutler, Anne & D. Robert Ladd, eds. (1983), *Prosody: Models and Measurements*, Berlin / Heidelberg, Springer-Verlag.
- Cutler, Anne & Dennis Norris (1988), "The role of strong syllables in segmentation for lexical access", *Journal of Experimental Psychology: Human Perception and Performance* 14:113-121.
- Darwin, Christopher J. (1976), "The perception of speech", in Carterette & Friedman (1976:175-226).
- Darwin, Christopher J. & Roy B. Gardner (1985), "Which harmonics contribute to the estimation of first formant frequency?", *Speech Communication* 4:231-235.
- DeCasper, Anthony J. & William P. Fifer (1980), "Of human bonding: Newborns prefer their mothers' voices", *Science* 208:1174-1176.
- Ding, Mingzhou, Betty Tuller & J.A. Scott Kelso (1995), "Characterizing the dynamics of auditory perception", *Chaos* 5:70-75.
- Docherty, Gerard J. & D. Robert Ladd, eds. (1992), *Papers in Laboratory Phonology II: Gesture, Segment, Prosody*, Cambridge, Cambridge University Press.
- Drake, Carolyn, Mari Riess Jones & Clarisse Baruch (2000), "The development of rhythmic attending in auditory sequences: Attunement, referent period, focal attending", *Cognition* 77:251-288.
- Duffy, Susan A. & David B. Pisoni (1992), "Comprehension of synthetic speech produced by rule: A review and theoretical interpretation", *Language and Speech* 35:351-389.
- Durlach, Nathaniel I. & Louis D. Braida (1969), "Intensity perception. I. Preliminary theory of intensity resolution", *Journal of the Acoustical Society of America* 46:372-383.
- Eimas, Peter D. & Joanne L. Miller (1978), "Effects of selective adaptation on the perception of speech and visual patterns: Evidence for feature detectors", in Walk & Pick (1978:307-345).
- Elman, Jeffrey L. (1990), "Finding structure in time", *Cognitive Science* 14:179-211.
- Elman, Jeffrey L. (1999), "The emergence of language: A conspiracy theory", in MacWhinney (1999:1-27).
- Elman, Jeffrey L. & James L. McClelland (1986), "Exploiting lawful variability in the speech wave", in Perkell & Klatt (1986:360-380).
- Faulkner, Andrew & Stuart Rosen (1999), "Contributions of temporal encodings of voicing, voicelessness, fundamental frequency, and amplitude variation to audio-visual and auditory speech perception", *Journal of the Acoustical Society of America* 106:2063-2073.
- Fellowes, Jennifer M., Robert E. Remez & Philip E. Rubin (1997), "Perceiving the sex and identity of a talker without natural vocal timbre", *Perception and Psychophysics* 59:839-849.
- Fixmer, Eric (2001), *Grouping of Auditory and Visual Information in Speech*, dissertation submitted to fulfil the requirements of the Ph.D. degree, University of Cambridge.
- Fougeron, Cécile & Patricia A. Keating (1997), "Articulatory strengthening at edges of prosodic domains", *Journal of the Acoustical Society of America* 101:3728-3740.

- Fowler, Carol A. (1986), "An event approach to the study of speech perception from a direct-realist perspective", *Journal of Phonetics* 14:3-28.
- Fowler, Carol A. & Dawn J. Dekle (1991), "Listening with eye and hand: Cross-modal contributions to speech perception", *Journal of Experimental Psychology: Human Perception and Performance* 17:816-828.
- Fowler, Carol A. & Mary R. Smith (1986), "Speech perception as 'vector analysis': An approach to the problems of invariance and segmentation", in Perkell & Klatt (1986:123-139).
- Fromkin, Victoria A. (1985), *Phonetic Linguistics: Essays in Honor of Peter Ladefoged*, Orlando, Academic Press.
- Fujisaki, Hiroya & Takako Kawashima (1970), "A model of the mechanisms for speech perception: Quantitative analysis of categorical effects in discrimination", University of Tokyo *Electrical Engineering Research Institute, Division of Electrical Engineering, Annual Report* 3:59-68.
- Gernsbacher, Morton A., ed. (1994), *Handbook of Psycholinguistics*, San Diego, Academic Press.
- Gibson, Eleanor J. (1991), *An Odyssey in Learning and Perception*, Cambridge, MA / London, MIT Press.
- Gibson, James J. (1966), *The Senses Considered as Perceptual Systems*, Boston, Houghton Mifflin.
- Gobl, Christer (1988), "Voice source dynamics in connected speech", *Speech Transmission Laboratory-Quarterly Progress and Status Report (STL-QPSR)*, Stockholm, KTH, 1:123-159.
- Goldinger, Stephen D. (1997), "Words and voices: Perception and production in an episodic lexicon", in Johnson & Mullennix (1997:33-66).
- Goldinger, Stephen D. (2000), "The role of perceptual episodes in lexical processing", in Anne Cutler, James M. McQueen & Rian Zondervan (eds.), *Proceedings of the Workshop on Spoken Word Access Processes (SWAP)*, Nijmegen, Max-Planck Institute for Psycholinguistics, 155-158.
- Goldinger, Stephen D., David B. Pisoni & John S. Logan (1991), "On the nature of talker variability effects on serial recall of spoken word lists", *Journal of Experimental Psychology: Learning, Memory and Cognition* 17:152-162.
- Goodman, Judith C. & Howard C. Nusbaum, eds. (1994), *The Development of Speech Perception: The Transition from Speech Sounds to Spoken Words*, Cambridge, MA / London, MIT Press.
- Gottfried, Terry L. & Winifred Strange (1980), "Identification of coarticulated vowels", *Journal of the Acoustical Society of America* 68:1626-1635.
- Green, Kerry P., Gail R. Tomiak & Patricia K. Kuhl (1997), "The encoding of rate and talker information during phonetic perception", *Perception and Psychophysics* 59:675-692.
- Greenberg, Steven (1996), "Understanding speech understanding: Towards a unified theory of speech perception", in *Proceedings of the ESCA Workshop on the Auditory Basis of Speech Perception*, (1996:1-8).
- Greenberg, Steven (1999), "Speaking in shorthand – A syllable-centric perspective for understanding pronunciation variation", *Speech Communication* 29:159-176.
- Grieser, DiAnne & Patricia K. Kuhl (1989), "Categorization of speech by infants: Support for speech sound prototypes", *Developmental Psychology* 25:577-588.
- Grossberg, Stephen (1986), "The adaptive self-organization of serial order in behavior: Speech, language, and motor control", in Schwab & Nusbaum (1986:187-293).
- Grossberg, Stephen (2000a), "Brain feedback and adaptive resonance in speech perception", *Behavioral and Brain Sciences* 23:332-333.
- Grossberg, Stephen (2000b), "How hallucinations may arise from brain mechanisms of learning, attention, and volition", *Journal of the International Neuropsychological Society* 6:583-592.
- Grossberg, Stephen & Christopher W. Myers (2000), "The resonant dynamics of speech perception: Interword integration and duration-dependent backward effects", *Psychological Review* 107:735-767.
- Grossberg, Stephen, Ian Boardman & Michael Cohen (1997), "Neural dynamics of variable-rate speech categorization", *Journal of Experimental Psychology: Human Perception and Performance* 23:481-503.
- Guenther, Frank H. & Marin N. Gjaja (1996), "The perceptual magnet effect as an emergent property of neural map formation", *Journal of the Acoustical Society of America* 100:1111-1121.

- Guenther, Frank H., Fatima T. Husain, Michael A. Cohen & Barbara G. Shinn-Cunningham (1999), "Effects of categorization and discrimination training on auditory perceptual space", *Journal of the Acoustical Society of America* 106:2900-2912.
- Halle, Morris & K.P. Mohanan (1985), "Segmental phonology of modern English", *Linguistic Inquiry* 16:57-116.
- Harley, Trevor A. (1995), *The Psychology of Language: From Data to Theory*, Hove, Psychology Press.
- Harnad, Stevan, ed. (1987), *Categorical Perception: The Groundwork of Cognition*, Cambridge, Cambridge University Press.
- Hawkins, Sarah (1995), "Arguments for a nonsegmental view of speech perception", in Kjell Elenius & Peter Branderud (eds.), *Proceedings of the XIIIth International Congress of Phonetic Sciences*, Stockholm, KTH and Stockholm University, 3:18-25.
- Hawkins, Sarah & Noël Nguyen (2000), "Predicting syllable-coda voicing from the acoustic properties of syllable onsets", in Anne Cutler, James M. McQueen & Rian Zondervan (eds.), *Proceedings of the Workshop on Spoken Word Access Processes (SWAP)*, Nijmegen, Max-Planck Institute for Psycholinguistics: 167-170.
- Hawkins, Sarah & Noël Nguyen (2001), "Perception of coda voicing from properties of the onset and nucleus of *led* and *let*", paper submitted to Eurospeech 2001.
- Hawkins, Sarah & Noël Nguyen (in press), "Effects on word recognition of syllable-onset cues to syllable-coda voicing", in Local, Ogden & Temple (in press).
- Hawkins, Sarah & Andrew Slater (1994), "Spread of CV and V-to-V coarticulation in British English: Implications for the intelligibility of synthetic speech", in *Proceedings of the 1994 International Conference on Spoken Language Processing*, Yokohama: 57-60.
- Hawkins, Sarah & Paul Warren (1994), "Implications for lexical access of phonetic influences on the intelligibility of conversational speech", *Journal of Phonetics* 22:493-511.
- Hebb, Donald O. (1949), *The Organization of Behavior: A Neurophysiological Theory*, New York / London, Wiley.
- Heid, Sebastian & Sarah Hawkins (1999), "Synthesizing systematic variation at boundaries between vowels and obstruents", in John J. Ohala, Yoko Hasegawa, Manjari Ohala, Daniel Granville & Ashlee C. Bailey (eds.), *Proceedings of the XIVth International Congress of Phonetic Sciences*, University of California, Berkeley, 1:511-514.
- Heid, Sebastian & Sarah Hawkins (2000), "An acoustical study of long-domain /r/ and /l/ coarticulation", in *Proceedings of the 5th Seminar on Speech Production: Models and Data (ISCA)*, Kloster Seeon, Bavaria, Germany: 77-80.
- Hillenbrand, James M., Michael J. Clark & Robert A. Houde (2000), "Some effects of duration on vowel recognition", *Journal of the Acoustical Society of America* 108:3013-3022.
- Hintzman, Douglas L. (1986), "'Schema abstraction' in a multiple-trace memory model", *Psychological Review* 93:411-428.
- Hintzman, Douglas L. (1988), "Judgments of frequency and recognition memory in a multiple-trace memory model", *Psychological Review* 95:528-551.
- Huffman, Marie K. & Rena Arens Krakow (1993), *Nasals, Nasalization, and the Velum*, New York, Academic Press.
- Huggins, A. William F. (1972a), "Just noticeable differences for segment duration in natural speech", *Journal of the Acoustical Society of America* 51:1270-1278.
- Huggins, A. William F. (1972b), "On the perception of temporal phenomena in speech", *Journal of the Acoustical Society of America* 51:1279-1290.
- Hume, Elizabeth V. & Keith A. Johnson, eds. (2001), *The Role of Speech Perception in Phonology*, New York, Academic Press.
- Iverson, Paul & Patricia K. Kuhl (1995), "Mapping the perceptual magnet effect for speech using signal detection theory and multidimensional scaling", *Journal of the Acoustical Society of America* 97:553-562.
- Iverson, Paul & Patricia K. Kuhl (1996), "Influences of phonetic identification and category goodness on American listeners' perception of /r/ and /l/", *Journal of the Acoustical Society of America* 99:1130-1140.

- Jenkins, James J., Winifred Strange & Thomas R. Edman (1983), "Identification of vowels in 'vowelless' syllables", *Perception and Psychophysics* 34:441-450.
- Johnson, Keith (1990), "Contrast and normalization in vowel perception", *Journal of Phonetics* 18:229-254.
- Johnson, Keith (1997), "Speech perception without speaker normalization: An exemplar model", in Johnson & Mullennix (1997:145-165).
- Johnson, Keith & John W. Mullennix, eds. (1997), *Talker Variability in Speech Processing*, San Diego / London, Academic Press.
- Johnson, Keith, Elizabeth A. Strand & Mariapaola D'Imperio (1999), "Auditory-visual integration of talker gender in vowel perception", *Journal of Phonetics* 27:359-384.
- Jones, Mari Riess (1976), "Time, our lost dimension: Toward a new theory of perception, attention and memory", *Psychological Review* 83:323-355.
- Jurafsky, Daniel (1996), "A probabilistic model of lexical and syntactic access and disambiguation", *Cognitive Science* 20:137-194.
- Juszyk, Peter W. (1993), "From general to language specific capacities: The WRAPSA model of how speech perception develops", *Journal of Phonetics* 21:3-28.
- Juszyk, Peter W. (1997), *The Discovery of Spoken Language*, Cambridge, MA, MIT Press.
- Juszyk, Peter W., David B. Pisoni & John Mullennix (1992), "Some consequences of stimulus variability on speech processing by 2-month old infants", *Cognition* 43:253-291.
- Juszyk, Peter W., Elizabeth A. Hohne, Ann Marie Juszyk & Nancy J. Redanz (1993), "Do infants remember voices?", *Journal of the Acoustical Society of America* 93:2373.
- Juszyk, Peter W., Elizabeth A. Hohne & Angela Bauman (1999a), "Infants' sensitivity to allophonic cues for word segmentation", *Perception and Psychophysics* 61:1465-1476.
- Juszyk, Peter W., Derek M. Houston & Mary Newsome (1999b), "The beginnings of word segmentation in English-learning infants", *Cognitive Psychology* 39:159-207.
- Keating, Patricia A. (1985), "Universal phonetics and the organization of grammars", in Fromkin (1985:115-132).
- Keating, Patricia A., ed. (1994), *Phonological Structure and Phonetic Form: Papers in Laboratory Phonology III*, Cambridge, Cambridge University Press.
- Keating, Patricia A., Taehong Cho, Cécile Fougerson & Chai-Shune Hsu (in press), "Domain-specific articulatory strengthening in four languages", in Local, Ogden & Temple (in press).
- Kelly, John & John Local (1986), "Long-domain resonance patterns in English", in *International Conference on Speech Input/Output; Techniques and Applications*, London, Institution of Electrical Engineers, Conference publication no. 258:304-309.
- Kelly, John & John K. Local (1989), *Doing Phonology*, Manchester, Manchester University Press.
- Kelso, J.A. Scott, Dan L. Southard & David Goodman (1979), "On the nature of human interlimb coordination", *Science* 203:1029-1031.
- Kewley-Port, Diane & Yijian Zheng (1999), "Vowel formant discrimination: Towards more ordinary listening conditions", *Journal of the Acoustical Society of America* 106:2945-2958.
- Kingston, John & Randy L. Diehl (1994), "Phonetic knowledge", *Language* 70:419-454.
- Klatt, Dennis H. (1979), "Speech perception: A model of acoustic-phonetic analysis and lexical access", *Journal of Phonetics* 7:279-312.
- Kozhevnikov, Valerij A. & Ludmilla A. Chistovich (1965), *Speech: Articulation and Perception*, Moscow-Leningrad, English translation: J.P.R.S., Washington, D.C., No. JPRS 30543.
- Krakov, Rena Arens (1993), "Nonsegmental influences on velum movement patterns: Syllables, sentences, stress, and speaking rate", in Huffman & Krakow (1993:87-113).
- Krakov, Rena Arens (1999), "Articulatory organization of syllables: A review", *Journal of Phonetics* 27:23-54.
- Krakov, Rena Arens, Patrice S. Beddor, Louis M. Goldstein & Carol A. Fowler (1988), "Coarticulatory influences on the perceived height of nasal vowels", *Journal of the Acoustical Society of America* 83:1146-1158.
- Kuhl, Patricia K. (1992), "Infants' perception and representation of speech: Development of a new theory", in John J. Ohala, Terrance M. Nearey, Bruce L. Derwing, Megan M. Hodge & Grace E. Wiebe (eds.), *Proceedings of the 1992 International Conference on Spoken Language Processing*, Banff, University of Alberta, 1:449-456.

- Kuhl, Patricia K. & Paul Iverson (1995), "Linguistic experience and the 'perceptual magnet effect'", in Strange (1995:121-154).
- Kuhl, Patricia K. & Andrew N. Meltzoff (1982), "The bimodal perception of speech in infancy", *Science* 218:1138-1141.
- Kwong, Katherine & Kenneth N. Stevens (1999), "On the voiced-voiceless distinction for writer/rider", in *Speech Communication Group Working Papers—Research Laboratory of Electronics*, MIT, XI:1-20.
- Kyriacou, Charalambos P. (in press), "The genetics of time", in *Time*, The Darwin College Lectures 2000, Cambridge, Cambridge University Press.
- Lacerda, Francisco (1995), "The perceptual-magnet effect: An emergent consequence of exemplar-based phonetic memory", in Kjell Elenius & Peter Branderud (eds.), *Proceedings of the XIIIth International Congress of Phonetic Sciences*, Stockholm, KTH and Stockholm University, 2:140-147.
- Ladd, D. Robert (1996), *Intonational Phonology*, Cambridge, Cambridge University Press.
- Ladefoged, Peter & Donald E. Broadbent (1957), "Information conveyed by vowels", *Journal of the Acoustical Society of America* 29:98-104.
- Large, Edward W. & Mari Riess Jones (1999), "The dynamics of attending: How people track time-varying events", *Psychological Review* 106:119-159.
- Lass, Norman J., ed. (1984), *Speech and Language: Advances in Basic Research and Practice*, vol.10, San Diego, Academic Press.
- Lehiste, Ilse (1972), "Manner of articulation, parallel processing, and the perception of duration", *Computer and Information Science Research Center Working Papers in Linguistics*, Ohio State University, 12:33-52.
- Levy, Joseph P., Dimitris Bairaktaris, John A. Bullinaria & Paul Cairns, eds. (1995), *Connectionist Models of Memory and Language*, London, UCL Press.
- Lieberman, Alvin M. & Ignatius G. Mattingly (1985), "The motor theory of speech perception revised", *Cognition* 21:1-36.
- Lindblom, Björn & Karen Rapp (1973), "Some temporal regularities of spoken Swedish", *PILUS (Papers from the Institute of Linguistics)*, Stockholm, University of Sweden, 21:1-58.
- Local, John K. (1992), "Modelling assimilation in a non-segmental rule-free phonology", in Docherty & Ladd (1992:190-223).
- Local, John K. & Richard Ogden (1997), "A model of timing for non-segmental phonological structure", in van Santen, Sproat, Olive & Hirschberg (1997:109-122).
- Local, John K., Richard A. Ogden & Rosalind A.M. Temple, eds. (in press), *Papers in Laboratory Phonology VI*, Cambridge, Cambridge University Press.
- Lotto, Andrew J., Keith R. Kluender & Lori L. Holt (1998), "Depolarizing the perceptual magnet effect", *Journal of the Acoustical Society of America* 103:3648-3655.
- Luce, Paul A., Stephen D. Goldinger & Michael S. Vitevich (2000), "It's good... but is it ART?", *Behavioral and Brain Sciences* 23:336.
- Macmillan, Neil A., Rina F. Goldberg & Louis D. Braida (1988), "Resolution for speech sounds: Basic sensitivity and context memory on vowel and consonant continua", *Journal of the Acoustical Society of America* 84:1262-1280.
- MacNeilage, Peter F., ed. (1983), *The Production of Speech*, New York, Springer-Verlag.
- MacWhinney, Brian (1999), *The Emergence of Language*, Mahwah, Lawrence Erlbaum Associates.
- McClelland, James L. & Jeffrey L. Elman (1986), "The TRACE model of speech perception", *Cognitive Psychology* 18:1-86.
- McGurk, Harry & John MacDonald (1976), "Hearing lips and seeing voices", *Nature* 264:746-748.
- Macchi, Marian J. (1980), "Identification of vowels spoken in isolation versus vowels spoken in consonantal context", *Journal of the Acoustical Society of America* 68:1636-1642.
- Magen, Harriet S. (1997), "The extent of vowel-to-vowel coarticulation in English and Japanese", *Journal of Phonetics* 25:187-205.
- Manuel, Sharon Y. (1990), "The role of contrast in limiting vowel-to-vowel coarticulation in different languages", *Journal of the Acoustical Society of America* 88:1286-1298.

- Manuel, Sharon Y. (1995), "Speakers nasalize /ð/ after /n/, but listeners still hear /ð/", *Journal of Phonetics* 23:453-476.
- Manuel, Sharon Y., Stefanie Shattuck-Hufnagel, Marie Huffman, Kenneth N. Stevens, Rolf Carlson & Sheri Hunnicutt (1992), "Studies of vowel and consonant reduction", in John J. Ohala, Terrance M. Nearey, Bruce L. Derwing, Megan M. Hodge & Grace E. Wiebe (eds.), *Proceedings of the 1992 International Conference on Spoken Language Processing*, Banff, University of Alberta, 2:943-946.
- Marslen-Wilson, William & Paul Warren (1994), "Levels of perceptual representation and process in lexical access: Words, phonemes, and features", *Psychological Review* 101:653-675.
- Martin, Christopher S., John W. Mullennix, David B. Pisoni & W. Van Summers (1989), "Effects of talker variability on recall of spoken word lists", *Journal of Experimental Psychology: Learning, Memory, and Cognition* 15:676-684.
- Massaro, Dominic W. (1998), *Perceiving Talking Faces: From Speech Perception to a Behavioral Principle*, Cambridge, MA / London, MIT Press.
- Mattys, Sven L. (1997), "The use of time during lexical processing and segmentation: A review", *Psychonomic Bulletin and Review* 4:310-329.
- Miller, Joanne L. & Thomas Baer (1983), "Some effects of speaking rate on the production of /b/ and /w/", *Journal of the Acoustical Society of America* 73:1751-1755.
- Miller, Joanne L. & Alvin M. Liberman (1979), "Some effects of later-occurring information on the perception of stop consonant and semivowel", *Perception and Psychophysics* 46:505-512.
- Moore, Brian C.J. (1997), *An Introduction to the Psychology of Hearing*, 4th edition, San Diego / London, Academic Press.
- Mullennix, John W., David B. Pisoni & Christopher S. Martin (1989), "Some effects of talker variability on spoken word recognition", *Journal of the Acoustical Society of America* 85:365-378.
- Nearey, Terrance M. (1995), "A double-weak view of trading relations: Comments on Kingston and Diehl", in Connell & Arvaniti (1995:28-39).
- Nearey, Terrance M. (1997), "Speech perception as pattern recognition", *Journal of the Acoustical Society of America* 101:3241-3254.
- Nearey, Terrance M. (2000), "Some concerns about the phoneme-like inputs to Merge", *Behavioral and Brain Sciences* 23:342-343.
- Nguyen, Noël & Sarah Hawkins (1999), "Implications for word recognition of phonetic dependencies between syllable onsets and codas", in John J. Ohala, Yoko Hasegawa, Manjari Ohala, Daniel Granville & Ashlee C. Bailey (eds.), *Proceedings of the XIVth International Congress of Phonetic Sciences*, University of California, Berkeley, 1:647-650.
- Ní Chasaide, Ailbhe & Christer Gobl (1993), "Contextual variation of the vowel voice source as a function of adjacent consonants", *Language and Speech* 36:303-330.
- Norris, Dennis G. (1992), "Connectionism: A new breed of bottom-up model?", in Reilly & Sharkey (1992:351-371).
- Norris, Dennis G. (1994), "Shortlist: a connectionist model of continuous speech recognition", *Cognition* 52:189-234.
- Norris, Dennis G., James M. McQueen & Anne Cutler (2000), "Merging information in speech recognition: Feedback is never necessary", *Behavioral and Brain Sciences* 23:299-370.
- Nosofsky, Robert M. (1988), "Exemplar-based accounts of relations between classification, recognition, and typicality", *Journal of Experimental Psychology: Learning, Memory, and Cognition* 14:700-708.
- Nygaard, Lynne C. & David B. Pisoni (1998), "Talker-specific learning in speech perception", *Perception and Psychophysics* 60:355-376.
- Nygaard, Lynne C., Mitchell S. Sommers & David B. Pisoni (1994), "Speech perception as a talker-contingent process", *Psychological Science* 5:42-46.
- Ogden, Richard (1999), "A declarative account of strong and weak auxiliaries in English", *Phonology* 16:55-92.
- Ogden, Richard & John K. Local (1994), "Disentangling autosegments from prosodies: a note on the misrepresentation of a research tradition in phonology", *Journal of Linguistics* 30:477-498.

- Ogden, Richard, Sarah Hawkins, Jill House, Mark Huckvale, John Local, Paul Carter, Jana Dankovičová & Sebastian Heid (2000), "ProSynth: An integrated prosodic approach to device-independent, natural-sounding speech synthesis", *Computer Speech and Language*, 14:177-210.
- Öhman, S.E.G. (1966), "Coarticulation in VCV utterances: Spectrographic measurements", *Journal of the Acoustical Society of America* 39:151-168.
- Parker, Ellen M. & Randy L. Diehl (1984), "Identifying vowels in CVC syllables: Effects of inserting silence and noise", *Perception and Psychophysics* 36:369-380.
- Patterson, Roy, <http://www.mrc-cbu.cam.ac.uk/personal/roy.patterson/aim/>.
- Perkell, Joseph S. (1986), "On sources of invariance and variability in speech production", in Perkell & Klatt (1986:260-263).
- Perkell, Joseph S. & Dennis H. Klatt, eds. (1986), *Invariance and Variability in Speech Processes*, Hillsdale, Lawrence Erlbaum Associates.
- Pickett, J.M. (1999), *The Acoustics of Speech Communication: Fundamentals, Speech Perception Theory, and Technology*, Needham Heights, Allyn and Bacon.
- Pickett, J.M. & Irwin Pollack (1963), "Intelligibility of excerpts from fluent speech: Effects of rate of utterance and duration of excerpt", *Language and Speech* 6:151-164.
- Pierrehumbert, Janet & David Talkin (1992), "Lenition of /h/ and glottal stop", in Docherty & Ladd (1992:90-117).
- Pisoni, David B. (1973), "Auditory and phonetic memory codes in the discrimination of consonants and vowels", *Perception and Psychophysics* 13:253-260.
- Pisoni, David B. (1997a), "Perception of synthetic speech", in van Santen, Sproat, Olive & Hirschberg (1997:541-560).
- Pisoni, David B. (1997b), "Some thoughts on 'normalization' in speech perception", in Johnson & Mullennix (1997:9-32).
- Pisoni, David B., Scott E. Lively & John S. Logan (1994), "Perceptual learning of nonnative speech contrasts: Implications for theories of speech perception", in Goodman & Nusbaum (1994:121-166).
- Plaut, David C. & Christopher T. Kello (1999), "The emergence of phonology from the interplay of speech comprehension and production: A distributed connectionist approach", in MacWhinney (1999:381-415).
- Pols, Louis C.W. (1986), "Variation and interaction in speech", in Perkell & Klatt (1986:140-154).
- Protopapas, Athanasios (1999), "Connectionist modeling of speech perception", *Psychological Bulletin* 125:410-436.
- Pulvermüller, Friedemann (1999), "Words in the brain's language", *Behavioral and Brain Sciences* 22:253-336.
- Pulvermüller, Friedemann (forthcoming), *Neuroscience of Language: From Brain Principles to Mechanisms Organizing Words and Serial Order*, Cambridge, Cambridge University Press.
- Rączaszek, Joanna, Betty Tuller, Lewis P. Shapiro, Pamela Case & Scott Kelso (1999), "Categorization of ambiguous sentences as a function of a changing prosodic parameter: A dynamical approach", *Journal of Psycholinguistic Research* 28:367-393.
- Recasens, Daniel (1989), "Long-range coarticulation effects for tongue-dorsum contact in VCVCV sequences", *Speech Communication* 8:293-307.
- Reilly, Ronan G. & Noel E. Sharkey, eds. (1992), *Connectionist Approaches to Natural Language Processing*, Hove, Erlbaum.
- Remez, Robert E. (2001), "The interplay of phonology and perception considered from the perspective of perceptual organization", in Hume & Johnson (2001:27-52).
- Remez, Robert E. & Philip E. Rubin (1992), "Acoustic shards, perceptual glue", *Haskins Laboratories Status Report on Speech Research* SR-111/112:1-10.
- Remez, Robert E., Philip E. Rubin, David B. Pisoni & Thomas D. Carrell (1981), "Speech perception without traditional speech cues", *Science* 212:947-950.
- Remez, Robert E., Philip E. Rubin, Stefanie M. Berns, Jennifer S. Pardo, and Jessica M. Lang (1994), "On the perceptual organization of speech", *Psychological Review* 101:129-156.
- Remez, Robert E., Jennifer M. Fellowes & Philip E. Rubin (1997), "Talker identification based on phonetic information", *Journal of Experimental Psychology: Human Perception and Performance* 23:651-666.

- Remez, Robert E., Jennifer M. Fellowes, David B. Pisoni, Winston D. Goh & Philip E. Rubin (1998), "Multimodal perceptual organization of speech: Evidence from tone analogs of spoken utterances", *Speech Communication* 26:65-73.
- Repp, Bruno H. (1982), "Phonetic trading relations and context effects: New experimental evidence for a speech mode of perception", *Psychological Bulletin* 92:81-110.
- Repp, Bruno H. (1984), "Categorical perception: Issues, methods, findings", in Lass (1984:243-335).
- Repp, Bruno H. & Alvin M. Liberman (1987), "Phonetic category boundaries are flexible", in Harnad (1987:89-112).
- Repp, Bruno H., Alvin M. Liberman, Thomas Eccardt & David Pesetsky (1978), "Perceptual integration of acoustic cues for stop, fricative, and affricate manner", *Journal of Experimental Psychology: Human Perception and Performance* 4:621-637.
- Rizzolatti, Giacomo & Michael A. Arbib (1998), "Language within our grasp", *Trends in Neurosciences* 21:188-194.
- Rosen, Stuart & Peter Howell (1987), "Auditory, articulatory, and learning explanations of categorical perception in speech", in Harnad (1987:113-160).
- Rumelhart, David E. & James L. McClelland (1986), "On learning the past tense of English verbs", in Rumelhart, McClelland & the PDP Research Group (1986:216-271).
- Rumelhart, David E., James L. McClelland & the PDP Research Group (1986), *Parallel Distributed Processing: Vol. 2. Psychological and Biological Models*, Cambridge, MA, MIT Press.
- Sachs, Jacqueline S. (1967), "Recognition memory for syntactic and semantic aspects of connected discourse", *Perception and Psychophysics* 2(9):437-442.
- Saltzman, Elliot & J.A. Scott Kelso (1987), "Skilled actions: A task-dynamic approach", *Psychological Review* 94:84-106.
- Saltzman, Elliot L. & Kevin G. Munhall (1989), "A dynamical approach to gestural patterning in speech production", *Ecological Psychology* 1:333-382.
- Sawyer, L., Michael Hennessey, Alexandre A. Peixoto, E. Rosato, H. Parkinson, Rodolfo Costa and Charalambos P. Kyriacou (1997), "Natural variation in a *Drosophila* clock gene and temperature compensation", *Science* 278:2117-2120.
- Schwab, Eileen C. & Howard C. Nusbaum, eds. (1986), *Pattern Recognition by Humans and Machines: Volume 1, Speech Perception*, Orlando, Academic Press.
- Schwab, Eileen C., James R. Sawusch & Howard C. Nusbaum (1981), "The role of second formant transitions in the stop-semivowel distinction", *Perception and Psychophysics* 29:121-128.
- Shannon, Robert V., Fan-Gang Zeng, Vivek Kamath, John Wygonski & Michael Ekelid (1995), "Speech recognition with primarily temporal cues", *Science* 270:303-304.
- Silipo, Rosaria, Steven Greenberg & Takayuki Arai (1999), "Temporal constraints on speech intelligibility as deduced from exceedingly sparse spectral representations", *Proceedings of Eurospeech 1999*: 2687-2690.
- Simpson, Greg B. (1994), "Context and the processing of ambiguous words", in Gernsbacher (1994:359-374).
- Smith, Linda B. (1999), "Children's noun learning: How general learning processes make specialized learning mechanisms", in MacWhinney (1999:277-303).
- Smith, Rachel & Sarah Hawkins (2000), "Allophonic influences on word-spotting experiments", in Anne Cutler, James M. McQueen & Rian Zondervan (eds.), *Proceedings of the Workshop on Spoken Word Access Processes (SWAP)*, Nijmegen, Max-Planck Institute for Psycholinguistics: 139-142.
- Solé, Maria-Josep (1995), "Spatio-temporal patterns of velopharyngeal action in phonetic and phonological organization", *Language and Speech* 38:1-23.
- Stevens, Kenneth N. (1983), "Design features of speech sound systems", in MacNeilage (1983:247-261).
- Stevens, Kenneth N. (1989), "On the quantal nature of speech", *Journal of Phonetics* 17:3-45.
- Stevens, Kenneth N. (1995), "Applying phonetic knowledge to lexical access", *4th European Conference on Speech Communication and Technology* 1:3-11.
- Stevens, Kenneth N. (1998), *Acoustic Phonetics*, Cambridge, MA / London, MIT Press.
- Stevens, Kenneth, Sharon Y. Manuel, Stefanie Shattuck-Hufnagel & Sharlene Liu (1992), "Implementation of a model for lexical access based on features", in John J. Ohala, Terrance

- M. Nearey, Bruce L. Derwing, Megan M. Hodge & Grace E. Wiebe (eds.), *Proceedings of the 1992 International Conference on Spoken Language Processing*, Banff, University of Alberta, 1:499-502.
- Stobart, Henry & Ian Cross (2000), "The Andean anacrusis? Rhythmic structure and perception in Easter songs of Northern Potosí, Bolivia", *British Journal of Ethnomusicology* 9:63-94.
- Strange, Winifred (1989), "Dynamic specification of coarticulated vowels spoken in sentence context", *Journal of the Acoustical Society of America* 85:2135-2153.
- Strange, Winifred, ed. (1995), *Speech Perception and Linguistic Experience: Issues in Cross-Language Research*, Baltimore, York Press.
- Strange, Winifred, Robert R. Verbrugge, Donald P. Shankweiler & Thomas R. Edman (1976), "Consonant environment specifies vowel identity", *Journal of the Acoustical Society of America* 60:213-224.
- Strange, Winifred, Thomas R. Edman & James J. Jenkins (1979), "Acoustic and phonological factors in vowel identification", *Journal of Experimental Psychology: Human Perception and Performance* 5:643-656.
- Strange, Winifred, James J. Jenkins & Thomas L. Johnson (1983), "Dynamic specification of coarticulated vowels", *Journal of the Acoustical Society of America* 74:695-705.
- Suomi, Kari (1993), "An outline of a developmental model of adult phonological organization and behavior", *Journal of Phonetics* 21:29-60.
- Swinney, David A. (1979), "Lexical access during sentence comprehension: (Re)consideration of context effects", *Journal of Verbal Learning and Verbal Behavior* 18:545-569.
- Tajima, Keiichi & Robert F. Port (in press), "Speech rhythm in English and Japanese", in Local, Ogden & Temple (in press).
- Tuller, Betty, Pamela Case, Mingzhou Ding & J.A. Scott Kelso (1994), "The nonlinear dynamics of speech categorization", *Journal of Experimental Psychology: Human Perception and Performance* 20:3-16.
- Tunley, Alison (1999), *Coarticulatory Influences of Liquids on Vowels in English*, unpublished Ph.D. dissertation, University of Cambridge.
- van Santen, Jan P.H., Richard W. Sproat, Joseph P. Olive & Julia Hirschberg, eds. (1997), *Progress in Speech Synthesis*, New York, Springer.
- van Tasell, Dianne J., Sigfrid D. Soli, Virginia M. Kirby & Gregory P. Widin (1987), "Speech waveform envelope cues for consonant recognition", *Journal of the Acoustical Society of America* 82:1152-1161.
- Walk, Richard D. & Herbert L. Pick (1978), *Perception and Experience*, New York, Plenum Press.
- Warren, Paul & William Marslen-Wilson (1987), "Continuous uptake of acoustic cues in spoken word recognition", *Perception and Psychophysics* 41:262-275.
- Warren, Richard M. (1970), "Perceptual restoration of missing speech sounds", *Science* 167:392-393.
- Warren, Richard M. (1984), "Perceptual restoration of obliterated sounds", *Psychological Bulletin* 96:371-383.
- Warren, Richard M. (1999), *Auditory Perception: A New Analysis and Synthesis*, Cambridge, Cambridge University Press.
- Warrington, Elizabeth K. & Rosaleen A. McCarthy (1987), "Categories of knowledge: Further fractionations and an attempted integration", *Brain* 110:1273-1296.
- Warrington, Elizabeth K. & T. Shallice (1984), "Category specific semantic impairments", *Brain* 107:829-854.
- Werker, Janet F. & John S. Logan (1985), "Cross-language evidence for three factors in speech perception", *Perception and Psychophysics* 37:35-44.
- Werker, Janet F. & Richard C. Tees (1984), "Phonemic and phonetic factors in adult cross-language speech perception", *Journal of the Acoustical Society of America* 75:1866-1878.
- West, Paula (1999a), "The extent of coarticulation of English liquids: An acoustic and articulatory study", in John J. Ohala, Yoko Hasegawa, Manjari Ohala, Daniel Granville & Ashlee C. Bailey (eds.), *Proceedings of the XIVth International Congress of Phonetic Sciences*, University of California, Berkeley, 3:1901-1904.

- West, Paula (1999b), "Perception of distributed coarticulatory properties of English /l/ and /ɫ/", *Journal of Phonetics* 27:405-426.
- Whalen, Douglas H. (2000), "Occam's razor is a double-edged sword: Reduced interaction is not necessarily reduced power", *Behavioral and Brain Sciences* 23:351.
- Whalen, Douglas H. & Sonya Sheffert (1997), "Normalization of vowels by breath sounds", in Johnson & Mullennix (1997:133-144).
- Zellner Keller, Brigitte & Eric Keller (forthcoming), "Representing speech rhythm", to appear in *Working Papers of COST 258*, University of Lausanne.
- Zue, Victor W. (1985), "The use of speech knowledge in automatic speech recognition", in *Proceedings of the Institute of Electrical and Electronics Engineers* 73:1602-1615.