

1 **34. Prosody and spoken-word recognition**

2 **James M. McQueen**

3 Radboud University & Max Planck Institute for Psycholinguistics, Nijmegen, The Netherlands

4 **&**

5 **Laura C. Dilley**

6 Michigan State University, East Lansing, USA

7
8 **Abstract**

9 This chapter outlines a Bayesian model of spoken-word recognition and reviews how
10 prosody is part of that model. The review focusses on the information which assists the
11 listener in recognizing the prosodic structure of an utterance and on how spoken-word
12 recognition is also constrained by prior knowledge about prosodic structure. Recognition is
13 argued to be a process of perceptual inference which ensures that listening is robust to
14 variability in the speech signal. In essence, the listener makes inferences about the
15 segmental content of each utterance, about its prosodic structure (simultaneously at
16 different levels in the prosodic hierarchy) and about the words it contains, and uses these
17 inferences to form an utterance interpretation. Four characteristics of the proposed
18 prosody-enriched recognition model are discussed: parallel uptake of different information
19 types, high contextual dependency, adaptive processing and phonological abstraction. The
20 next steps that should be taken to develop the model are also discussed.

21
22
23 **Keywords**

24 Spoken-word recognition; suprasegmental information; prosodic structure; prosody
25 hierarchy; phonological abstraction; perceptual learning; variability problem; Bayesian
26 model

27

1 **34.1 Introduction**

2 Each spoken utterance is potentially unique and is one of an infinite range of possible
 3 utterances. But each is made from words that usually have been heard before, sampled from
 4 the finite set of words the speaker/listener knows. To understand the speaker's intended
 5 message in any utterance, therefore, the listener must recognize the utterance's words. We
 6 argue here that she, the listener, achieves spoken-word recognition through Bayesian
 7 perceptual inference. Her task, over and over again for each word, is to infer the identity of
 8 the current word, and build an interpretation, integrating current acoustic information with
 9 prior knowledge. In this chapter, we consider the role of 'prosody' in this process of
 10 perceptual recovery of spoken words.

11 **34.2 Defining prosody in spoken-word recognition**

12 We begin with a definition of 'prosody'. This is not only because it can mean different things
 13 to different people, but also because one of our goals is to highlight the utility of an abstract
 14 definition of prosody that has to do with structures built in the mind of the perceiver.
 15 Critically, this definition is tied to the cognition in question: the process of spoken-word
 16 recognition. Our definition therefore does not start from linguistic material (words,
 17 sentences) or from the acoustic properties of speech (e.g., spectral and durational features)
 18 but instead from a psychological perspective, focussing on the representations and
 19 processes the listener uses as she understands speech.

20
 21 The basis of our definition is that, during word recognition, two types of structure are
 22 built in the listener's mind. The former structures are 'segmental' in that they are based on
 23 abstractions about segments – the traditional combinatorial 'building blocks' of words. The
 24 latter structures are 'suprasegmental' and relate to abstractions about the prominence,
 25 accentuation, grouping, expressive tone of voice, etc., of syllables relative to each other and
 26 also of words relative to each other. The latter structures are prosodic, and hence to
 27 understand the role of prosody in word recognition is to have an adequate account of how
 28 these structures are built, but also how the segmental structures are built, and how these
 29 two types of structure jointly support speech understanding.

30 This definition thus highlights the interdependency, during processing, of signal
 31 characteristics often classified as 'segmental' and 'suprasegmental'. For example, pitch
 32 characteristics (i.e., perceptual indices of fundamental frequency variations) – often
 33 considered to be suprasegmental in the spoken-word recognition literature – may frequently
 34 contribute simultaneously to extracting both segmental and suprasegmental structures, as
 35 well as other kinds of structures (e.g., syntactic). In the same vein, acoustic characteristics
 36 relating to distributions of periodic (i.e., vocal fold vibration) or aperiodic energy – often
 37 considered to be segmental in the spoken-word recognition literature – contribute to
 38 extracting both segmental structures (e.g., words) and suprasegmental structures (e.g.,
 39 prosodic phrase-level structures through domain-initial strengthening of segments, see
 40 below), as well as other kinds of structures (e.g., syntactic). Again, this happens in an
 41 interdependent fashion across levels of structure. That such interdependences among
 42 different levels of structure exist in spoken-word recognition is consistent with the
 43 observation that lexical entries are defined in part by the constructs of *syllable* and *stress* –
 44 each of which has both a 'segmental' and a 'suprasegmental' interpretation. That a given
 45 acoustic attribute (e.g., fundamental frequency in speech, which gives rise to a harmonic
 46 spectrum) contributes simultaneously to perception of both segmental and suprasegmental
 47 structures, has long been recognized (e.g., Lehiste 1970).

1 Consideration of this interdependence across different levels of the linguistic hierarchy
 2 during structure extraction is also motivated by our perspective on speech recognition. In
 3 our view, a core challenge to be explained is how words are extracted from the speech
 4 stream in spite of considerable variability. That is, a spoken-word recognizer needs to be
 5 robust in the face of acoustic variability of various kinds, for example, differences among
 6 phonological contexts, speakers, speaking styles and listening conditions. We argue that
 7 redundancy in encoding multi-leveled tiers of structure across different kinds of acoustic
 8 information means that the system is more robust to any one kind of acoustic degradation.
 9 That is, listeners build interlocking segmental and suprasegmental phonological structures as
 10 a means to solving the variability problem.

11 We believe that our cognitive definition of prosody allows us to avoid several
 12 problems. In particular, we do not need to define particular types of acoustic cue as strictly
 13 either ‘segmental’ or ‘suprasegmental’. Such attempts come with the implication that
 14 whatever phonetic properties are taken to define ‘suprasegmental’ – usually timing and
 15 pitch – are via logical opposition ‘not segmental’, and thus that these do not cue segmental
 16 contrasts. Indeed, such a view is highly problematic, as has been noted by many researchers
 17 (e.g., Lehiste 1970). Much work has documented the role of timing in cueing of segmental
 18 contrasts, including both consonants (Lieberman, Cooper, Shankweiler & Studdert-Kennedy
 19 1967; Lisker & Abramson 1964; Wade & Holt 2005) and vowels (cf. vowel length or
 20 tenseness; Ainsworth 1972; Miller 1981).

21 Under our proposal, acoustic information can nevertheless still be categorized as that
 22 which assists the listener in recognizing either the segments of an utterance (‘segmental
 23 information’) or its prosodic structure (‘suprasegmental information’). Our definition is in
 24 service of the view that spoken-word recognition involves simultaneously recognizing the
 25 words being said, the prosodic (e.g., grouping, prominence) structures associated with those
 26 words, and the larger structures (e.g., syntactic ones) in which the words are embedded. On
 27 this view, it becomes easier to see how diverse acoustic cues – ranging from pitch to timing
 28 to allophonic phonetic variation – could be employed to help extract structure (lexical and
 29 otherwise) at various hierarchical levels.

30 The same acoustic information can therefore help the listener simultaneously identify
 31 segmental and prosodic structures. Take the case of domain-initial strengthening, in which
 32 acoustic cues for consonants and vowels tend to be strengthened (e.g., become longer, or
 33 louder, or add glottal stops or other fortification) at the beginnings of prosodic domains
 34 (Beňuš & Šimko 2014; Cho 2016; Cho & Keating 2001; Dilley, Shattuck-Hufnagel & Ostendorf
 35 1996; Fougeron & Keating 1997; Garellek 2014; Krivokapić & Byrd 2012; Tabain 2003; Turk &
 36 Shattuck-Hufnagel 2000). Domain-initial strengthening affects pitch, timing and spectral
 37 details, but also concerns systematic variation at the lexical level, such that it can help with
 38 lexical disambiguation (Cho, McQueen & Cox 2007) and at the utterance level (such that it
 39 helps the listener with sentential parsing and interpretation building). That is, domain-initial
 40 strengthening concerns variation simultaneously at (at least) two levels of structure.

41 Domain-initial strengthening is an example of cross-talk between segmental and
 42 suprasegmental domains. Another example relates to the widespread usage of pitch in the
 43 world’s languages to convey lexical contrast. Not only is pitch used throughout the lexicon to
 44 convey lexical contrasts in lexical tone languages (e.g., Mandarin, Thai, Igbo), but pitch also
 45 plays a role in distinguishing words in languages such as Japanese and Swedish (Beckman
 46 1986; Bruce 1977; Heldner & Strangert 2001). Even intonation languages (e.g., English,
 47 Spanish, German, and Dutch) include lexical contrasts based on stress (e.g., *IMPact* (noun)

1 vs. *imPACT* (verb)) which may be signalled by a difference in pitch in many structural and
 2 communicative contexts, but certainly not all (Fry 1958; Gussenhoven 2004). Indeed, the
 3 acoustic cues that signal lexical stress contrasts are many and varied and include not only
 4 segmental vowel-quality differences but also differences in timing, amplitude, and/or
 5 spectral balance as well as pitch (Banzina, Dilley & Hewitt 2016; Beckman & Edwards 1994;
 6 Mattys 2000; Morrill 2012; Sluijter & van Heuven 1996; Turk & White 1999).

7 Our definition also highlights how prosody can assist in the perceptual recovery of
 8 spoken words when the speech signal is degraded. For example, fine spectral details in
 9 signals usually associated with segmental information can be replaced with a few frequency
 10 bands of noise, producing noise-vocoded speech, or the dynamic formants can be replaced
 11 with sine waves, producing sinewave speech. Such degraded speech is often highly
 12 intelligible, especially with practice (Davis, Johnsrude, Hervais-Adelman, Taylor &
 13 McGettigan 2005; Dorman, Loizou & Rainey 1997; Shannon, Zeng, Kamath, Wygonski &
 14 Ekelid 1995). Such perceptual recovery of spoken words is possible partly because the
 15 listener is able to make contact with her prior experiences of timing and frequency
 16 properties of spoken words experienced over her lifetime. That is, this ability indicates that
 17 stored knowledge about word forms may include timing, pitch and amplitude information.
 18 A critical feature of our fundamentally cognitive definition is thus that it refers not only to
 19 relevant acoustic information but also to relevant prior knowledge. To explore prosody in
 20 spoken-word recognition is thus to ask how suprasegmental information and prior
 21 knowledge about prosodic structures, together with segmental information and prior
 22 knowledge about segments, jointly support speech comprehension. We propose that the
 23 answer to this question is that speech recognition involves Bayesian inference.

24 **34.3 The Bayesian Prosody Recognizer: Robustness under variability**

25 A growing body of evidence supports a Bayesian account of spoken-word recognition in
 26 which simultaneous multiple interdependent hypotheses are considered about the words
 27 being said, their component segments, and aspects of expressiveness that are heard to
 28 accompany those words. According to this view, the linguistic structures which are perceived
 29 are those that ultimately best explain experienced sensory information. Our proposal is that
 30 a Bayesian Prosody Recogniser (BPR) supports this inferential process by extracting prosodic
 31 structures (syllables, phrases) and words while deriving utterance interpretations. The BPR
 32 draws inspiration from other Bayesian models of speech recognition and understanding and
 33 analysis-by-synthesis approaches (Gibson, Bergen & Piantadosi 2013; Halle & Stevens 1962;
 34 Kleinschmidt & Jaeger 2015; Norris & McQueen 2008; Poeppel, Idsardi & van Wassenhove
 35 2008) which envision the inferential, predictive process of spoken language understanding as
 36 involving simultaneous determination of multiple levels of linguistic structures, including
 37 hierarchical prosodic structures. In essence, as guaranteed by Bayes' rule, the listener
 38 combines prior knowledge with signal-driven likelihoods to obtain an optimal interpretation
 39 of current input. The BPR also draws inspiration from previous accounts arguing that speech
 40 recognition requires parallel evaluation of segmental and suprasegmental interpretations (in
 41 particular the Prosody Analyser of Cho et al., 2007). Evidence for predictive and inferential
 42 processes in speech recognition is reviewed in multiple sources (Kuperberg & Jaeger 2016;
 43 Norris, McQueen & Cutler 2016; Pickering & Garrod 2013; Tavano & Scharinger 2015).

44 A central motivation for the BPR is the variability problem, as already introduced:
 45 Structure extraction needs to be robust in spite of variability in speech. Bayesian inference is
 46 a response to this challenge because it ensures optimal interpretation of the current input.
 47

1 The BPR instantiates four key characteristics about prosodic processing in spoken-word
 2 recognition. All are further specifications of how the BPR offers ways to ensure robustness of
 3 recognition under acoustic variability.

5 **34.3.1 Parallel uptake of information**

6 As we review below, considerable evidence from studies examining the temporal dynamics
 7 of the recognition process supports our contention that timing and pitch characteristics
 8 constrain word identification, and that they do so at the same time as segmental
 9 information. In our view, parallel uptake of information has at least two important
 10 consequences. First, it makes it possible that structures can be extracted at different
 11 representational levels simultaneously. This can readily be instantiated in the BPR. Just like
 12 there can be, in a Bayesian framework, a hierarchy of segments (Kleinschmidt & Jaeger
 13 2015), words (Norris & McQueen 2008), and sentences (Gibson et al. 2013), there can also
 14 be a Bayesian prosodic hierarchy, potentially from syllables up to intonational phrases.
 15 Second, it means that the same acoustic information can contribute simultaneously to
 16 construction of different levels of linguistic representation, including the prosodic,
 17 phonological, lexical, and higher (syntactic, semantic, pragmatic) levels. In order to
 18 accomplish the above, the BPR must analyse information across windows of varying sizes
 19 simultaneously (some quite long, such as recognizing a tune or determining turn-taking
 20 structures in discourse). As an example of both of the above, consider that as durational
 21 information for a prosodic word (i.e., a single lexical item) accumulates, it can also provide
 22 the basis of evidence for a phrase which contains that word. Evidence about that word
 23 influences the interpretation of syntactic information, and so forth. Suprasegmental
 24 information (as acoustically defined) has been shown to influence recognition in at least four
 25 different ways.

27 *34.3.1.1 Influences on processing segmental information*

28 Segments belonging to stressed syllables in sentences are processed more quickly than those
 29 belonging to unstressed syllables (Cutler & Foss 1977; Shields, McHugh & Martin 1974).
 30 Segmental content in stressed syllables is more accurately perceived than that in unstressed
 31 syllables (Bond & Garnes 1980), and mispronounced segments are more easily detected in
 32 stressed syllables than in unstressed syllables (Cole & Jakimik 1978). Distortion of normal
 33 word stress information also impairs word processing and recognition (Bond & Small 1983;
 34 Cutler & Clifton 1984; Slowiaczek 1990, 1991). Recent findings indicate that categorisation of
 35 speech segments is modulated by the type of prosodic boundary preceding those segments
 36 (Kim & Cho 2013; Mitterer, Cho & Kim 2016). All of the above evidence supports the view
 37 that suprasegmental and segmental sources of acoustic information in words are the basis of
 38 parallel inference processes at multiple levels of linguistic structure. In keeping with this
 39 view, it has been shown that the same information (durational cues; Tagliapietra &
 40 McQueen 2010) can simultaneously help listeners determine which segments they are
 41 hearing and the location of word boundaries.

43 *34.3.1.2 Influences on lexical segmentation*

44 Consistent with the BPR, the metrical properties of a given syllable affect the likelihood with
 45 which listeners infer the syllable to be word-initial (Cutler, Dahan & van Donselaar 1997;
 46 Cutler & Norris 1988). For instance, strong syllables are more likely heard as word-initial in
 47 errors in perception (Cutler & Butterfield 1992). There is evidence that listeners use multiple

1 cues (some lexical and some signal-driven, based on segmental and suprasegmental acoustic
 2 properties) to segment continuous speech into words (Norris, McQueen, Cutler & Butterfield
 3 1997). Suprasegmental cues appear to play a more important role under more difficult
 4 listening conditions. Thus, for example, the tendency to assume that strong syllables are
 5 word-initial is stronger when stimuli are presented in background noise than when there is
 6 no noise (Mattys 2004; Mattys, White & Melhorn 2005).

7 8 *34.3.1.3 Influences on lexical selection*

9 Suprasegmental pronunciation modifications modulate which words the listener considers
 10 and which words she eventually recognizes. For example, subtle differences in segment
 11 durations or whole syllables can help her determine the location of syllable boundaries
 12 (Tabossi, Collina, Mazzetti & Zoppello 2000), word boundaries (Gow & Gordon 1995) and
 13 prosodic boundaries (e.g., in making the distinction between a monosyllabic word such as
 14 *cap* and the initial syllable of a longer word such as *captain* (Blazej & Cohen-Goldberg 2015;
 15 Davis, Marslen-Wilson & Gaskell 2002; Salverda, Dahan & McQueen 2003). Additional kinds
 16 of suprasegmental acoustic-phonetic information, including pitch and intensity, also
 17 modulate perception of syllable boundaries (Garellek 2014; Heffner, Dilley, McAuley & Pitt
 18 2013; Hillenbrand & Houde 1996). The rapidity with which these kinds of lexical
 19 disambiguation take place (as measured, e.g., with eye tracking; Salverda et al., 2003)
 20 indicates that suprasegmental processing is not delayed relative to segmental processing.
 21 Variation in pronunciation associated with distinct positions of words in prosodic phrases
 22 (e.g., whether the two words in the phrase ‘bus tickets’ span an intonation phrase boundary
 23 or not) has also been shown to modulate lexical selection (Cho et al. 2007; Christophe,
 24 Peperkamp, Pallier, Block & Mehler 2004; see also Tremblay, Broersma & Coughlin 2018;
 25 Tremblay, Broersma, Coughlin & Choi 2016 for similar non-native language effects).

26 Some earlier studies (Cutler 1986; Cutler & Clifton 1984) suggested that stress
 27 differences cued by suprasegmental information (e.g., the distinction between the ‘ancestor’
 28 and ‘tolerate’ senses of *forbear*, which is not due to a difference in the segments of the
 29 words; Cutler, 1986) did not constrain lexical access substantially. Subsequent experiments,
 30 however, indicated that stress does constrain lexical access, albeit to different extents in
 31 different languages, as a function of the informational value of suprasegmental stress cues in
 32 the language in question (Cooper, Cutler & Wales 2002; Cutler & van Donselaar 2001; Soto-
 33 Faraco, Sebastian-Galles & Cutler 2001). For example, the influence of suprasegmental stress
 34 cues on word recognition is stronger in Dutch, where such cues tell listeners more about
 35 which words have been spoken, than in English, where segmental differences are more
 36 informative (Cooper et al., 2002). Eye-tracking studies indicate that suprasegmental cues to
 37 stress are taken up without delay and can thus support lexical disambiguation before any
 38 segmental cues could disambiguate the input (Brown, Salverda, Dilley & Tanenhaus 2015;
 39 Reinisch, Jesse & McQueen 2010). Relatedly, work on word recognition in tone languages
 40 has shown how pitch characteristics of the input constrain word identification in parallel
 41 with the uptake of segmental information (Lee 2009; Sjerps, Zhang & Peng 2018).

42 43 *34.3.1.4 Influences on inferences about other structures*

44 Consistent with the BPR, there is considerable evidence that suprasegmental information
 45 influences the listener’s inferences about various levels of structure beyond the word level,
 46 simultaneously, in real time. The focus of this chapter is on spoken-word recognition, but
 47 since perception of lexical forms influences higher levels of linguistic structure and inference,

1 it is important to note that there is evidence that prosody and other higher levels of
 2 linguistic knowledge are extracted in parallel. That is, perception of prosodic information and
 3 of syntactic structure are interdependent (Buxó-Lugo & Watson 2016; Carlson, Clifton &
 4 Frazier 2001) and prosody influences semantic and pragmatic inference (Ito & Speer 2008;
 5 Rohde & Kurumada 2018).

7 **34.3.2 High contextual dependency**

8 Another characteristic of prosodic processing in spoken-word recognition is its high
 9 contextual dependency. That is, the interpretation of the current prosodic event depends on
 10 the context which occurs before and/or after that event. Context can be imagined as a
 11 timeline, where ‘left context’ temporally precedes an event and ‘right context’ follows it.

13 **34.3.2.1 Left-context effects**

14 Under the BPR account, regularities in context which are statistically predictive of properties
 15 of upcoming words will be used to infer lexical properties of upcoming words, giving rise to
 16 left-context effects. It is well-attested that suprasegmental aspects of sentential context
 17 affect the speed of processing of elements. For example, suprasegmental cues in a sequence
 18 of words preceding a given word affect processing speed on that word (Cutler 1976; Pitt &
 19 Samuel 1990) and accuracy of word identification (Slowiaczek 1991). The rhythm of stressed
 20 and unstressed syllables is an important cue for word segmentation in continuous speech
 21 (Nakatani & Schaffer 1978). Further, a metrically regular speech context has also been
 22 shown to promote spoken-word recognition (Quené & Port 2005). Our BPR proposal
 23 accounts for these findings in terms of statistical inference on the basis of regularities in the
 24 speech signal. Structures in utterances formed by prosodic (e.g., rhythmic) patterning in
 25 production engenders predictability of structure and timing of upcoming sentential elements
 26 (Jones 1976; Martin 1972) at multiple, hierarchical levels and points (Lieberman & Prince
 27 1977). Statistical regularities in stress alternation and timing are attested in speech
 28 production experiments, corpus studies and theoretical linguistics (Farmer, Christiansen &
 29 Monaghan 2006; Hayes 1995; Kelly & Bock 1988; Selkirk 1984). Changes in the priors in a
 30 Bayesian model can account easily for the effects of left prosodic context (and other types of
 31 preceding context) on recognition of the current word.

32 Contextual influences of suprasegmental cues on perception of segmental
 33 information (e.g., VOT) are well known, particularly for timing (Kidd 1989; Miller & Liberman
 34 1979; Repp 1982) but also for pitch (Dilley 2010; Dilley & Brown 2007; Holt 2006; Sjerps et al.
 35 2018). However, such effects have by and large been found to involve fairly proximal speech
 36 context within about 300 ms of a target segment (Kidd 1989; Newman & Sawusch 1996;
 37 Sawusch & Newman 2000; Summerfield 1981; but see Wade & Holt 2005).

38 More recent work has shown that suprasegmental information from the more distant
 39 (‘distal’) left context can also influence which words are heard – including how syllables are
 40 grouped into words, and even whether certain words (and hence certain phonemes) are
 41 heard at all. For example, the rate of distal context speech influences whether listeners hear
 42 reduced words such as *are* spoken as ‘err’ (Dilley & McAuley 2008; Pitt, Szostak & Dilley
 43 2016). Statistical distributions of distal contextual speech rates influence listeners’ word
 44 perception over the course of ~1 hour (Baese-Berk et al. 2014). Further, the patterns of pitch
 45 and timing on prominent and nonprominent syllables in the left context influences where
 46 listeners hear word boundaries in lexically-ambiguous sequences such as *crisis turnip* vs. *cry*
 47 *sister nip* (Dilley, Mattys & Vinke 2010; Dilley & McAuley 2008; Morrill, Dilley & McAuley

1 2014). These patterns also influence the extent to which listeners hear reduced words or
 2 syllables (Baese-Berk, Dilley, Henry, Vinke & Banzina 2019; Morrill, Dilley, McAuley & Pitt
 3 2014). Distal rate and rhythm influence lexical processing early in perception and modulate
 4 the extent to which lexically-stressed syllables are heard to be word-initial (Breen, Dilley,
 5 McAuley & Sanders 2014; Brown, Salverda, Dilley & Tanenhaus 2011; Brown et al. 2015).
 6 Consistent with BPR, whether a listener hears a word depends in gradient, probabilistic
 7 fashion on the joint influence of distal rate cues and proximal information signalling a word
 8 boundary (Heffner et al. 2013).

9 10 34.3.2.2 *Right-context effects*

11 Information which follows can be informative about lexical content that may have already
 12 elapsed. A growing body of evidence that listeners often commit to an interpretation of
 13 lexical content only *after* the temporal offset of that content (Bard, Shillcock & Altmann
 14 1988; Connine, Blasko & Hall 1991; Grossberg & Myers 2000; McMurray 2007). In segmental
 15 perception, temporal information to the right of a given segment can influence listeners'
 16 judgments of segmental perception (e.g., Miller and Liberman 1979). Eyetracking studies
 17 show that later-occurring distal temporal information (e.g., relative duration of a subsequent
 18 phoneme sequence that includes the morpheme /s/) can influence whether listeners hear a
 19 prior reduced function word (Brown, Dilley & Tanenhaus 2014). All of these findings indicate
 20 that acoustic information must be held in some kind of memory buffer and hence that
 21 perceptual decisions can be delayed until after the acoustic offset of that information. The
 22 extent to which listeners hold alternative parses in mind after a given portion of signal
 23 consistent with a given word has elapsed, as opposed to abandoning them, is an active area
 24 of research and debate (Christiansen & Chater 2016).

25 While effects of right context might at first glance appear to be more problematic,
 26 they too can be explained in a Bayesian framework. The key notion here is that different
 27 hierarchical levels of structure and constituency (e.g., segments, syllables, words, prosodic
 28 phrases) entail different time windows over which relevant evidence is collected and applied
 29 to generation of inferences about representations at that level. This implies that acoustic
 30 evidence at a given moment might be taken as highly informative for structure at one level,
 31 while simultaneously being taken as only weakly informative (or indeed uninformative)
 32 about structure at another level. Depending on the imputed reliability of evidence as it
 33 appertains to each level, inferences about structure at different levels may be made at
 34 different rates (i.e., are staggered in time). Because evidence bearing on the structure of a
 35 larger constituent (e.g., a prosodic phrase) typically will appear in the signal over a longer
 36 time span than evidence bearing on the structure of a smaller one (e.g., a syllable),
 37 completion of the inferences about the larger constituent may often entail consideration of
 38 evidence from some amount of subsequent 'right-context' material. This apparent delay
 39 with respect to inferences about the structure of the larger constituent does not imply that
 40 the BPR does not always attempt to use all information simultaneously nor that it does not
 41 attempt to draw inferences at different levels simultaneously. Rather, it implies only that in
 42 some cases the current information is insufficient for inferences at a given level of structure
 43 to be made with confidence, and hence that the BPR may wait for further information in the
 44 upcoming context before committing to an interpretation of structure at that level. This view
 45 also entails that later-occurring information might provide evidence that an earlier
 46 assumption about structure was not well-supported and hence the possibility of revision of
 47 inferences drawn earlier.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

34.3.2.3 Syntagmatic representation of pitch

Phonological interpretation of pitch cues in spoken language comprehension requires consideration of both left and right pitch context (Francis, Ciocca, Wong, Leung & Chu 2006; Sjerps et al. 2018). Left and right context is also important for listeners to draw abstractions about the tonal properties of a given syllable, including that which is relevant to perceiving distinct lexical items (Dilley & Brown 2007; Dilley & McAuley 2008; Wong & Diehl 2003). Such findings support a view in which the representation of linguistically-relevant pitch information is fundamentally syntagmatic (i.e., relational), and that paradigmatic aspects of tonal information involve inferences driven by abstract knowledge about a typical speaker's pitch range in relation to incoming pitch information (Dilley 2005, 2008; Dilley & Breen to appear; Lai 2018). This view is adopted in the BPR.

34.3.3 Adaptive processing

The perceptual apparatus must dynamically adapt to variation in order to remain robust in understanding intended messages. The available evidence suggests that prosodic processing is indeed very flexible. For instance, listeners adapt rapidly to the rate of compressed speech (Dupoux & Green 1997). The evidence just reviewed on context effects shows that listeners track characteristics of the current speech (e.g., distributional properties of speaking rate variation and the metrical properties of utterances) and flexibly adjust to that context (Baese-Berk et al. 2014; Dilley & McAuley 2008; Dilley & Pitt 2010; Morrill, Baese-Berk, Heffner & Dilley 2015).

Another way in which prosodic processing has been shown to be adaptive is that it involves perceptual learning. It has been established that listeners can adapt to variation in the realization of segments (Norris, McQueen & Cutler 2003; Samuel & Kraljic 2009): they tune in, as it were, to the segmental characteristics of the speech of the current talker. It is thus plausible to expect that there are similar adjustments with respect to suprasegmental characteristics. There is indirect evidence that this may be the case. Listeners adapt to the characteristics of accented as well as distorted speech (Baese-Berk, Bradlow & Wright 2013; Borrie et al. 2012; Bradlow & Bent 2008; Mitterer & McQueen 2009), which presumably includes adjustments to suprasegmental features. But there is also more direct evidence. Dutch listeners in a perceptual-learning paradigm can adjust the way they interpret the reduced syllables of a particular Dutch speaker (Poellmann, Bosker, McQueen & Mitterer 2014), and Mandarin listeners adjust the way they interpret the tonal characteristics of syllables through exposure to stimuli with ambiguous pitch contours in contexts which encouraged a particular tonal interpretation (Mitterer, Chen & Zhou 2011).

The BPR therefore needs to be flexible. Detailed computational work on perceptual learning in a Bayesian model with respect to speech segments has already been performed (Kleinschmidt & Jaeger 2015). The argument, in a nutshell, is that learning is required for the listener to be able to recognise speech optimally, in the context of an input that is noisy and highly variable due, for instance, to differences among talkers (Kleinschmidt & Jaeger 2015; Norris et al. 2003). That is, the ideal observer needs to be an ideal adapter. Exactly the same arguments apply to prosodic variability. Learning processes, for example based on changes in the probability density function of a given prosodic constituent for a given idiosyncratic talker, should be instantiated in the BPR in a similar way to those already implemented for segments.

1 **34.3.4 Phonological abstraction**

2 The final characteristic of prosodic processing in spoken-word recognition is that it is based
 3 on phonological abstraction. The listener must be able to form abstractions so as to remain
 4 optimally robust and capable of handling not-yet-encountered variation. Phonological
 5 abstraction is thus also a feature of the BPR. As in the previous Bayesian accounts focussing
 6 on segmental recognition (Kleinschmidt & Jaeger 2015; Norris & McQueen 2008), the
 7 representations that inferences are drawn about are abstract categories so that (as the
 8 adaptability of the BPR also guarantees) the recognition process is robust to variation due to
 9 differences across talkers and listening situations. Evidence suggests that the abstractions
 10 about categories entail generalizations about segmental structures and allophonic variation
 11 (Mitterer, Reinisch & McQueen 2018), lexical stress and tone (Ramachers 2018; Sjerps et al.
 12 2018; Sulpizio and McQueen 2012), pitch accent, pitch range and boundary tone types
 13 (Cutler & Otake 1999; Dilley & Brown 2007; Dilley & Heffner 2013), and relationships
 14 between phonological elements and other aspects of the linguistic structure of information,
 15 such as grammatical categories (Farmer et al. 2006; Kelly 1992; Söderström, Horne,
 16 Mannfolk, van Westen & Roll 2017).

17 Prosodic processing in speech recognition appears to involve phonological abstraction.
 18 One line of evidence for this comes from the learning studies just reviewed. If perceptual
 19 learning generalizes to the recognition of words that have not been heard during the
 20 exposure phase, then some type of abstraction must have taken place – the listener must
 21 know which entities to apply the learning to (cf. McQueen, Cutler & Norris 2006). The
 22 studies on learning about syllables (Poellmann et al. 2014) and tones (Mitterer et al. 2011)
 23 both show generalization of learning to the recognition of previously unheard words.

24 Experiments on learning novel words also provide evidence that listeners have abstract
 25 knowledge about prosody. In these experiments (on prosodic words in Dutch, Shatzman &
 26 McQueen 2006; , and on lexical stress in Italian, Sulpizio & McQueen 2012) listeners learned
 27 new minimal pairs of words; the new words were acoustically altered to remove
 28 suprasegmental cues that distinguished between the pairs. In a final test phase, the listeners
 29 heard the altered (training) words and their unaltered (original) variants. Eye-tracking
 30 measures revealed that the listeners had knowledge about the suprasegmental cues that
 31 they could apply to the on-line recognition of the novel words, even though they had never
 32 heard those words with those cues (for the Dutch listeners, durational cues distinguishing
 33 monosyllabic words from the initial syllables of disyllabic words; for the Italian listeners,
 34 durational and amplitude cues to antepenultimate stress in trisyllabic words). These findings
 35 suggest that processing of prosody in spoken-word recognition involves not only the uptake
 36 of fine-grained acoustic-phonetic cues to prosodic structure, but also the storage of abstract
 37 knowledge about those cues. That is, while the fine phonetic details about the prosody in
 38 the current utterance are key determinants of word recognition and speech comprehension,
 39 the listener abstracts over those details in order to be able to understand future utterances.

40 Speakers also form phonological abstractions based on long-term knowledge of
 41 phonetic properties of talker attributes, such as gender (Johnson, Strand & D'Imperio 1999;
 42 Lai 2018), that contribute to Bayesian inferences about spoken words and other aspects of
 43 linguistic meaning. Phonological abstractions are also formed based on simultaneous or
 44 sequential statistical correspondences among phonetic properties, such as pitch and non-
 45 modal voice quality, which are phonetic properties that co-vary in many lexical tone
 46 languages (Garellek & Keating 2011; Garellek, Keating, Esposito & Kreiman 2013; Gerfen &
 47 Baker 2005; Gordon & Ladefoged 2001). Such phonological abstraction – formed from long-

1 term statistical knowledge of correspondences – is essential for drawing correct inferences
 2 based on otherwise highly ambiguous suprasegmental cues (including those for pitch and
 3 duration) about, for example, intended words, meaning, and structure (Bishop & Keating
 4 2012; Gerfen & Baker 2005; Lai 2018). For instance, knowledge about co-occurrences of
 5 pitch and spectral (e.g., formant frequency) information for male vs. female voices can be
 6 used to infer a typical or mean pitch of a talker’s voice and/or pitch span, from which
 7 Bayesian inferences can be drawn about phonological structures (such as those for pitch
 8 accents and lexical tones) and associated meanings (Dilley 2005; Dilley & Breen to appear).
 9 The BPR assumes that such long-term abstracted statistical knowledge about talkers and the
 10 simultaneous and sequential distributional properties of the phonetic cues they produce is,
 11 along with talker-independent abstract phonological knowledge, the basis of the Bayesian
 12 probabilistic inferences which enable optimal decoding of spoken signals.

13

14 **34.4 Conclusions and future directions**

15 We have argued that spoken-word recognition is robust under speech variability because it
 16 is based on Bayesian perceptual inference and that a vital component of this process is the
 17 BPR. As a spoken utterance unfolds over time, the BPR, based on prior knowledge about
 18 correspondences between acoustic variables and meanings and structures, makes Bayesian
 19 inferences about the prosodic organization, lexical content, and semantic and pragmatic
 20 information in the utterance, among other inferences. These inferences are both signal- and
 21 knowledge-driven and concern abstract structures at different levels in the prosodic
 22 hierarchy which are computed in parallel, informed by statistical distributions of
 23 relationships among acoustic cues often considered segmental or suprasegmental.
 24 Inferences about a given stretch of input are influenced by earlier input and by inferences
 25 about it, and can be revised based on later input. Importantly, the BPR adapts to current
 26 input to optimize its inferences.

27 We have suggested that the goal of the BPR is to derive the metrical and grouping
 28 structures in each utterance at different levels in the prosodic hierarchy. Especially for
 29 utterance-level inferences, the representation must include a sparse set of tones, including
 30 pitch accents, boundary tones and/or lexical tones, which are autosegmentally associated
 31 with particular positions in metrical and grouping structures indexed to the lexicon (Dilley &
 32 Breen to appear; Gussenhoven 2004; Ladd 2008). Establishing how listeners recover this
 33 prosodic hierarchy, and the number of levels that need to be built, are important challenges
 34 for future research.

35 The BPR will need to be implemented as part of a full Bayesian model of speech
 36 recognition, which includes, but is not limited to, prosodic inferences. Our view is that
 37 segmental and suprasegmental structures are built in parallel, based on information that
 38 may inform inferences about either or both types of structure. Over time, inferences about
 39 prosodic structure feed into (and are in turn influenced by) inferences made about segments
 40 and words of the unfolding utterance and its current interpretation. The model will need to
 41 specify how interacting processes determine spoken-word recognition and how inferences
 42 drawn about the speech signal change over time. It will also need to be tested, through
 43 simulations and experimentation.

44 One way to evaluate and develop the BPR would be to compare it to other models on
 45 the role of prosody in spoken-word recognition. Unfortunately, no such alternative models
 46 currently exist. Shuai and Malins (2017) have recently proposed TRACE-T, and
 47 implementation of TRACE (McClelland & Elman 1986) which seeks to account for the

1 processing of tonal information in Mandarin monosyllabic words. While this is a very
2 welcome addition to the literature, TRACE-T is much more limited in scope than the BPR.
3 Comparisons could potentially also be made to the Prosody Analyzer (Cho et al. 2007; but
4 the BPR can be seen as a development of that model) and to Shortlist B (Norris & McQueen
5 2008; but Shortlist B is limited, with respect to prosody, to the role of metrical structure in
6 lexical segmentation, and again the BPR is largely inspired by the earlier model). Detailed
7 comparisons of the BPR to other models (e.g., Kurumada, Brown & Tanenhaus 2018) will
8 have to wait for the implementation of the BPR and for the development of competitor
9 models of equivalent scope.

10 Another important aspect of future work will be cross-linguistic comparison. Most
11 work on prosody in spoken-word recognition has been done on English or a small set of
12 related European languages. There are some exceptions to this Eurocentric bias (Cutler &
13 Otake 1999; Lee 2007; Ye & Connine 1999), and there has been an upsurge of work on, for
14 example, pitch cues in conveying lexical and other meanings in typologically diverse
15 languages (Genzel & Kügler in press; Kula & Braun 2015; Ramachers 2018; Sjerps et al. 2018;
16 Wang, Xu & Ding 2018; Yamamoto & Haryu 2018). Much research nevertheless still is
17 needed to explore how the full set of prosodic phenomena in the world's languages
18 modulates the recognition process. We do not expect that experiments on non-European
19 languages will lead to falsification of the Bayesian model. For example, pitch conveys
20 different kinds of structure simultaneously in a given language, and is used to convey lexical
21 information to different degrees in different languages. Pitch is simply less informative about
22 lexical structure in a Bayesian statistical sense in intonation languages than in tone
23 languages and thus will be relied on less in discriminating among and recognizing words in
24 intonation languages. While such cross-linguistic differences can thus readily be captured in
25 a Bayesian model, it will be important to explore how pitch information can simultaneously
26 inform inferences about words and inferences about intonational structures in a tone
27 language, and how this weighting changes in intonation vs. lexical tone languages.

28 The Bayesian model will need to be developed in the direction of neurobiological
29 implementation. As in psycholinguistic research (including computational modelling), much
30 work in cognitive neuroscience focusses on how segments (e.g. individual consonants or
31 vowels) are recognized, and how that contributes to word recognition. Prosody had tended
32 to be ignored. There are some interesting new approaches, for example, evidence of neural
33 entrainment to the 4 Hz oscillations at which speech tends to be spoken (i.e., the 'syllable
34 rate') (Ding et al. 2017; Giraud & Poeppel 2012). Nevertheless, much work still needs to be
35 done to specify the brain mechanisms which support spoken-word recognition as a process
36 that depends on parallel inferences about segmental content and prosodic structures (e.g.,
37 whether entrainment is modulated by information arriving in the speech signal at faster or
38 slower rates than 4 Hz).

39 It will also be necessary to specify how the proposed model relates to other aspects of
40 language processing, speech production in particular. Knowledge that is needed to support
41 recognition (e.g., the acoustic characteristics of words with penultimate stress) may not be
42 relevant in speech production. It remains to be determined if and to what extent the
43 processes and representations involved in recognition are shared with those involved in
44 production. It is already clear, however, that there is an intimate relationship between input
45 and output operations. For example, the Bayesian recognition process depends on the ability
46 of the recogniser to track production statistics. There are undoubtedly constraints on which
47 statistics are tracked (e.g. with respect to the size of the structures which are tracked), but

1 future work will need to establish what those constraints are and why certain statistics are
2 tracked and not others.

3 There is also a need to evaluate the model not only relative to other domains of
4 cognitive psychology (such as speech production, language acquisition and second-language
5 processing) but also to other domains of linguistics. The representations of prosodic
6 structure that a listener needs for efficient speech recognition may or may not have a one-
7 to-one correspondence with those that are most relevant (for example) to language
8 typology. It is theoretically possible, for example, that a structure such as the prosodic word
9 may have an essential role in typological work and yet have no role in processes relating to
10 the cognitive construction of prosodic structures during spoken-word recognition. It is
11 another important challenge for future research to establish the extent to which
12 representations of prosody indeed vary across different domains of linguistic enquiry.

13 We have here reviewed the state-of-the-art of research on prosody in spoken-word
14 recognition. Rather than being theoretically neutral, we have advocated a specific model.
15 We look forward to future research testing our central claim that prosody influences speech
16 recognition through Bayesian perceptual inference.

17
18

1 References

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

- Ainsworth, W. A. (1972). Duration as a cue in the recognition of synthetic vowels. *Journal of the Acoustical Society of America* 51: 648-651.
- Baese-Berk, M., L. Dilley, M. Henry, L. Vinke, and E. Banzina (2019). Not just a function of function words: Distal speech rate affects perception of prosodically weak syllables. *Attention, Perception, & Psychophysics* 81: 571-589.
- Baese-Berk, M., C. Heffner, L. Dilley, M. Pitt, T. Morrill, and J. D. McAuley (2014). Long-term temporal tracking of speech rate affects spoken-word recognition. *Psychological Science* 25: 1546-1553.
- Baese-Berk, M. M., A. R. Bradlow, and B. A. Wright (2013). Accent-independent adaptation to foreign-accented speech. *Journal of the Acoustical Society of America* 133: 174-180.
- Banzina, E., L. Dilley, and L. Hewitt (2016). The role of secondary stressed and unstressed unreduced syllables in word recognition: acoustic and perceptual studies with Russian learners of English. *Journal of Psycholinguistic Research* 45: 813-831.
- Bard, E. G., R. C. Shillcock, and G. T. Altmann (1988). The recognition of words after their acoustic offsets in spontaneous speech: Effects of subsequent context. *Perception & Psychophysics* 44: 395-408.
- Beckman, M. (1986). *Stress and Non-stress Accent*. Dordrecht: Foris.
- Beckman, M.E., and J. Edwards (1994). Articulatory evidence for differentiating stress categories. In P. Keating (Ed.), *Papers in Laboratory Phonology III: Phonological Structure and Phonetic Form* (pp. 7-33). Cambridge University Press, Cambridge.
- Beňuš, Š., and J. Šimko (2014). Emergence of prosodic boundary: continuous effects of temporal affordance on inter-gestural timing. *Journal of Phonetics* 44: 110-129.
- Bishop, J., and P. A. Keating (2012). Perception of pitch location within a speaker's range: Fundamental frequency, voice quality, and speaker sex. *Journal of the Acoustical Society of America* 132: 1100-1112.
- Blazej, L. J., and A. M. Cohen-Goldberg (2015). Can we hear morphological complexity before words are complex? *Journal of Experimental Psychology: Human Perception and Performance* 41: 50-68.
- Bond, Z. S., and S. Garnes (1980). Misperceptions of fluent speech. In R. A. Cole (Ed.), *Perception and production of fluent speech*. Hillsdale, NJ: Erlbaum.
- Bond, Z. S., and L. H. Small (1983). Voicing, vowel, and stress mispronunciations in continuous speech. *Perception & Psychophysics* 34: 470-474.
- Borrie, S. A., M. J. McAuliffe, J. M. Liss, C. Kirk, G. A. O'Beirne, and T. Anderson (2012). Familiarisation conditions and the mechanisms that underlie improved recognition of dysarthric speech. *Language and Cognitive Processes* 27: 1039-1055.
- Bradlow, A. R., and T. Bent (2008). Perceptual adaptation to non-native speech. *Cognition* 106: 707.
- Breen, M., L. Dilley, J. D. McAuley, and L. Sanders (2014). Auditory evoked potentials reveal early perceptual effects of distal prosody on speech segmentation. *Language, Cognition and Neuroscience* 29: 1132-1146.
- Brown, M., L. Dilley, and M. Tanenhaus (2014). Probabilistic prosody: Effects of relative speech rate on perception of (a) word(s) several syllables earlier. In N. Campbell, D. Gibbon, & D. Hirst (Eds.), *Proceedings of the 7th International Conference on Speech Prosody* (pp. 1154-1158). Dublin, Ireland.

- 1 Brown, M., A. P. Salverda, L. Dilley, and M. Tanenhaus (2011). Expectations from preceding
2 prosody influence segmentation in online sentence processing. *Psychonomic Bulletin*
3 *and Review* 18: 1189-1196.
- 4 Brown, M., A. P. Salverda, L. Dilley, and M. Tanenhaus (2015). Metrical expectations from
5 preceding prosody influence perception of lexical stress. *Journal of Experimental*
6 *Psychology: Human Perception and Performance* 41: 306-323.
- 7 Bruce, G. (1977). *Swedish word accents in sentence perspective*. Lund: Gleerups.
- 8 Buxó-Lugo, A., and D. Watson (2016). Evidence for the influence of syntax on prosodic
9 parsing. *Journal of Memory and Language* 90: 1-13.
- 10 Carlson, K., C. J. Clifton, and L. Frazier (2001). Prosodic boundaries in adjunct attachment.
11 *Journal of Memory and Language* 45: 58-81.
- 12 Cho, T. (2016). Prosodic boundary strengthening in the phonetics–prosody interface.
13 *Language and Linguistics Compass* 10: 120-141.
- 14 Cho, T., and P. Keating (2001). Articulatory and acoustic studies on domain-initial
15 strengthening in Korean. *Journal of Phonetics* 29: 155-190.
- 16 Cho, T., J. M. McQueen, and E. A. Cox (2007). Prosodically driven phonetic detail in speech
17 processing: The case of domain-initial strengthening in English. *Journal of Phonetics*
18 35: 210-243.
- 19 Christiansen, M. H., and N. Chater (2016). The Now-or-Never bottleneck: A fundamental
20 constrain on language. *Behavioral and Brain Sciences* 39: e62.
- 21 Christophe, A., S. Peperkamp, C. Pallier, E. Block, and J. Mehler (2004). Phonological phrase
22 boundaries constrain lexical access: I. Adult data. *Journal of Memory and Language*
23 51: 523-547.
- 24 Cole, R. A., and J. Jakimik (1978). Understanding speech: How words are heard. In G.
25 Underwood (Ed.), *Strategies of information processing* (pp. 67-116). London:
26 Academic Press.
- 27 Connine, C. M., D. Blasko, and M. Hall (1991). Effects of subsequent sentence context in
28 auditory word recognition: Temporal and linguistic constraints. *Journal of Memory*
29 *and Language* 30: 234-250.
- 30 Cooper, N., A. Cutler, and R. Wales (2002). Constraints of lexical stress on lexical access in
31 English: Evidence from native and non-native listeners. *Language and Speech* 45, No.
32 3: 207-228.
- 33 Cutler, A. (1976). Phoneme-monitoring reaction time as a function of preceding intonation
34 contour. *Perception and Psychophysics* 20: 55-60.
- 35 Cutler, A. (1986). *Forbear* is a homophone: Lexical prosody does not constrain lexical access.
36 *Language and Speech* 29: 201-220.
- 37 Cutler, A., and S. Butterfield (1992). Rhythmic cues to speech segmentation: Evidence from
38 juncture misperception. *Journal of Memory and Language* 31: 218-236.
- 39 Cutler, A., and C. E. Clifton (1984). The use of prosodic information in word recognition. In H.
40 Bouma & D. G. Bouwhuis (Eds.), *Attention and Performance X* (pp. 183-196).
41 Hillsdale, NJ: Erlbaum.
- 42 Cutler, A., D. Dahan, and W. van Donselaar (1997). Prosody in the comprehension of spoken
43 language: A literature review. *Language and Speech* 40: 141-201.
- 44 Cutler, A., and D. Foss (1977). On the role of sentence stress in sentence processing.
45 *Language and Speech* 20: 1-10.

- 1 Cutler, A., and D. G. Norris (1988). The role of strong syllables in segmentation for lexical
2 access. *Journal of Experimental Psychology: Human Perception and Performance* 14:
3 113-121.
- 4 Cutler, A., and T. Otake (1999). Pitch accent in spoken-word recognition in Japanese. *Journal*
5 *of the Acoustical Society of America* 105: 1877-1888.
- 6 Cutler, A., and W. van Donselaar (2001). *Voornaam* is not (really) a homophone: Lexical
7 prosody and lexical access in Dutch. *Language and Speech* 44: 171-195.
- 8 Davis, M. H., I. S. Johnsrude, A. Hervais-Adelman, K. Taylor, and C. McGettigan (2005). Lexical
9 information drives perceptual learning of distorted speech: Evidence from the
10 comprehension of noise-vocoded sentences. *Journal of Experimental Psychology:*
11 *General* 134: 222-241.
- 12 Davis, M. H., W. D. Marslen-Wilson, and M. G. Gaskell (2002). Leading up the lexical garden
13 path: Segmentation and ambiguity in spoken word recognition. *Journal of*
14 *Experimental Psychology: Human Perception and Performance* 28: 218-244.
- 15 Dilley, L. 2005. *The phonetics and phonology of tonal systems*. (Ph.D. dissertation), MIT,
16 Cambridge, MA.
- 17 Dilley, L. (2008). On the dual relativity of tone. In *Proceedings from the Annual Meeting of*
18 *the Chicago Linguistics Society* (Vol. 41, pp. 129-144). Chicago, IL.
- 19 Dilley, L. (2010). Pitch range variation in English tonal contrasts is continuous, not
20 categorical. *Phonetica* 67: 63-81.
- 21 Dilley, L., and M. Breen (to appear). An *enhanced* autosegmental-metrical theory (AM⁺)
22 facilitates phonetically transparent prosodic annotation: A reply to Jun. In J. Barnes &
23 S. Shattuck-Hufnagel (Eds.), *Prosodic Theory and Practice*. Cambridge, MA: MIT Press.
- 24 Dilley, L., and M. Brown (2007). Effects of pitch range variation on F0 extrema in an imitation
25 task. *Journal of Phonetics* 35: 523-551.
- 26 Dilley, L., and C. Heffner (2013). The role of f0 alignment in distinguishing intonation
27 categories: Evidence from American English. *Journal of Speech Sciences* 3: 3-67.
- 28 Dilley, L., S. Mattys, and L. Vinke (2010). Potent prosody: comparing the effects of distal
29 prosody, proximal prosody, and semantic context on word segmentation. *Journal of*
30 *Memory and Language* 63: 274-294.
- 31 Dilley, L., and J. D. McAuley (2008). Distal prosodic context affects word segmentation and
32 lexical processing. *Journal of Memory and Language* 59: 294-311.
- 33 Dilley, L., and M. Pitt (2010). Altering context speech rate can cause words to appear or
34 disappear. *Psychological Science* 21: 1664-1670.
- 35 Dilley, L., S. Shattuck-Hufnagel, and M. Ostendorf (1996). Glottalization of word-initial
36 vowels as a function of prosodic structure. *Journal of Phonetics* 24: 423-444.
- 37 Ding, N., A. D. Patel, L. Chen, H. Butler, C. Luo, and D. Poeppel (2017). Temporal modulations
38 in speech and music. *Neuroscience and Biobehavioral Reviews* 81: 181-187.
- 39 Dorman, M. F., P. C. Loizou, and D. Rainey (1997). Speech understanding as a function of the
40 number of channels of stimulation for processors using sine-wave and noise-band
41 outputs. *Journal of the Acoustical Society of America* 102: 2403-2411.
- 42 Dupoux, E., and K. Green (1997). Perceptual adjustment to highly compressed speech:
43 Effects of talker and rate changes. *Journal of Experimental Psychology: Human*
44 *Perception and Performance* 23: 914-927.
- 45 Farmer, T. A., M. H. Christiansen, and P. Monaghan (2006). Phonological typicality influences
46 on-line sentence comprehension. *Proceedings of the National Academy of Sciences*
47 103: 12203-12208.

- 1 Fougeron, C., and P. A. Keating (1997). Articulatory strengthening at edges of prosodic
2 domains. *Journal of the Acoustical Society of America* 101: 3728-3740.
- 3 Francis, A. L., V. Ciocca, N. K. Y. Wong, W. H. Y. Leung, and P. C. Y. Chu (2006). Extrinsic
4 context affects perceptual normalization of lexical tone. *Journal of the Acoustical*
5 *Society of America* 119: 1712-1726.
- 6 Fry, D. B. (1958). Experiments in the perception of stress. *Language and Speech* 1: 126-152.
- 7 Garellek, M. (2014). Voice quality strengthening and glottalization. *Journal of Phonetics* 45:
8 106-113.
- 9 Garellek, M., and P. A. Keating (2011). The acoustic consequences of phonation and tone
10 interactions in Jalapa Mazatec. *Journal of the International Phonetic Association* 41:
11 185-205.
- 12 Garellek, M., P. A. Keating, C. Esposito, and J. Kreiman (2013). Voice quality and tone
13 identification in White Hmong. *The Journal of the Acoustical Society of America* 133:
14 1078-1089.
- 15 Genzel, S., and F. Kügler (in press). Production and perception of question prosody in Akan.
16 *Journal of the International Phonetic Association*.
- 17 Gerfen, C., and K. Baker (2005). The production and perception of laryngealized vowels in
18 Coatzospan Mixtec. *Journal of Phonetics* 33: 311-334.
- 19 Gibson, E., L. Bergen, and S. T. Piantadosi (2013). The rational integration of noisy evidence
20 and prior semantic expectations in sentence interpretation. *Proceedings of the*
21 *National Academy of Sciences* 110: 8051-8056.
- 22 Giraud, A.-L., and D. Poeppel (2012). Cortical oscillations and speech processing: Emerging
23 computational principles and operations. *Nature Neuroscience* 15: 511-517.
- 24 Gordon, M., and P. Ladefoged (2001). Phonation types: A cross-linguistic overview. *Journal of*
25 *Phonetics* 29: 383-406.
- 26 Gow, D., and P. C. Gordon (1995). Lexical and prelexical influences on word segmentation:
27 Evidence from priming. *Journal of Experimental Psychology: Human Perception and*
28 *Performance* 21: 344-359.
- 29 Grossberg, S., and C. Myers (2000). The resonant dynamics of speech perception: Interword
30 integration and duration-dependent backward effects. *Psychological Review* 107:
31 735-767.
- 32 Gussenhoven, C. (2004). *The phonology of tone and intonation*. Cambridge: Cambridge
33 University Press.
- 34 Halle, M., and K. N. Stevens (1962). Speech recognition: A model and a program for research.
35 *IEEE Transactions on Information Theory* 8: 155-159.
- 36 Hayes, B. (1995). *Metrical Stress Theory: Principles and Case Studies*. Chicago: University of
37 Chicago Press.
- 38 Heffner, C., L. Dilley, J. D. McAuley, and M. Pitt (2013). When cues combine: how distal and
39 proximal acoustic cues are integrated in word segmentation. *Language and Cognitive*
40 *Processes* 28: 1275-1302.
- 41 Heldner, M., and E. Strangert (2001). Temporal effects of focus in Swedish. *Journal of*
42 *Phonetics* 29: 329-361.
- 43 Hillenbrand, J., and R. A. Houde (1996). Role of F0 and amplitude in the perception of
44 intervocalic glottal stops. *Journal of Speech and Hearing Research* 39: 1182-1190.
- 45 Holt, L. L. (2006). The mean matters: Effects of statistically defined nonspeech spectral
46 distributions on speech categorization. *Journal of the Acoustical Society of America*
47 120: 2801-2817.

- 1 Ito, K., and S. R. Speer (2008). Anticipatory effects of intonation: Eye movements during
2 instructed visual search. *Journal of Memory and Language* 58: 541-573.
- 3 Johnson, K., E. A. Strand, and M. D'Imperio (1999). Auditory-visual integration of talker
4 gender in vowel perception. *Journal of Phonetics* 27: 359-384.
- 5 Jones, M. R. (1976). Time, our lost dimension: Toward a new theory of perception, attention,
6 and memory. *Psychological Review* 83: 323-355.
- 7 Kelly, M. H. (1992). Using sound to solve syntactic problems: The role of phonology in
8 grammatical category assignments. *Psychological Review* 99: 349-364.
- 9 Kelly, M. H., and J. K. Bock (1988). Stress in time. *Journal of Experimental Psychology: Human*
10 *Perception and Performance* 14: 389-403.
- 11 Kidd, G. R. (1989). Articulatory rate-context effects in phoneme identification. *Journal of*
12 *Experimental Psychology: Human Perception and Performance* 15: 736-748.
- 13 Kim, S., and T. Cho (2013). Prosodic boundary information modulates phonetic
14 categorization. *Journal of the Acoustical Society of America* 134: EL19-EL25.
- 15 Kleinschmidt, D. F., and F. T. Jaeger (2015). Robust speech perception: Recognize the
16 familiar, generalize to the similar, and adapt to the novel. *Psychological Review* 122:
17 148-203.
- 18 Krivokapić, J., and D. Byrd (2012). Prosodic boundary strength: An articulatory and
19 perceptual study. *Journal of Phonetics* 40: 430-442.
- 20 Kula, N. C., and B. Braun (2015). Mental representation of tonal spreading in Bemba:
21 Evidence from elicited production and perception. *Southern African Linguistics and*
22 *Applied Language Studies* 33: 307-323.
- 23 Kuperberg, G. R., and F. T. Jaeger (2016). What do we mean by prediction in language
24 comprehension? *Language, Cognition and Neuroscience* 31: 32-59.
- 25 Kurumada, C., M. Brown, and M. K. Tanenhaus (2018). Effects of distributional information
26 on categorization of prosodic contours. *Psychonomic Bulletin and Review* 25: 1153-
27 1160.
- 28 Ladd, D. R. (2008). *Intonational Phonology* (2nd ed.). Cambridge: Cambridge University Press.
- 29 Lai, W. (2018). *Voice gender effect on tone categorization and pitch perception*. Paper
30 presented at the Sixth International Symposium on Tonal Aspects of Language,
31 Berlin.
- 32 Lee, C.-Y. (2007). Does horse activate mother? Processing lexical tone in form priming.
33 *Language and Speech* 50: 101-123.
- 34 Lee, C.-Y. (2009). Identifying isolated, multispeaker Mandarin tones from brief acoustic
35 input: A perceptual and acoustic study. *Journal of the Acoustical Society of America*
36 125: 1125-1137.
- 37 Lehiste, I. (1970). *Suprasegmentals*. Cambridge, MA: MIT Press.
- 38 Liberman, A. M., F. S. Cooper, D. P. Shankweiler, and M. Studdert-Kennedy (1967).
39 Perception of the speech code. *Psychological Review* 74: 431-461.
- 40 Liberman, M., and A. Prince (1977). On stress and linguistic rhythm. *Linguistic Inquiry* 8: 249-
41 336.
- 42 Lisker, L., and A. S. Abramson (1964). A cross-language study of voicing in initial stops:
43 Acoustical measurements. *WORD: Journal of the International Linguistic Association*
44 20: 384-422.
- 45 Martin, J. G. (1972). Rhythmic (hierarchical) versus serial structure in speech and other
46 behavior. *Psychological Review* 79: 487-509.

- 1 Mattys, S. L. (2000). The perception of primary and secondary stress in English. *Perception*
2 *and Psychophysics* 62: 253-265.
- 3 Mattys, S. L. (2004). Stress versus coarticulation: Toward an integrated approach to explicit
4 speech segmentation. *Journal of Experimental Psychology: Human Perception and*
5 *Performance* 30: 397-408.
- 6 Mattys, S. L., L. White, and J. F. Melhorn (2005). Integration of multiple speech segmentation
7 cues: a hierarchical framework. *Journal of Experimental Psychology: General* 134:
8 477-500.
- 9 McClelland, J. L., and J. L. Elman (1986). The TRACE model of speech perception. *Cognitive*
10 *Psychology* 18: 1-86.
- 11 McMurray, B. (2007). Defusing the childhood vocabulary explosion. *Science* 317: 631-631.
- 12 McQueen, J. M., A. Cutler, and D. Norris (2006). Phonological abstraction in the mental
13 lexicon. *Cogn Sci* 30: 1113-1126.
- 14 Miller, J. L. (1981). Phonetic perception: evidence for context-dependent and context-
15 independent processing. *Journal of the Acoustical Society of America* 69: 822-831.
- 16 Miller, J. L., and A. M. Liberman (1979). Some effects of later-occurring information on the
17 perception of stop consonant and semivowel. *Perception & Psychophysics* 25: 457-
18 465.
- 19 Mitterer, H., Y. Chen, and X. Zhou (2011). Phonological abstraction in processing lexical-tone
20 variation: evidence from a learning paradigm. *Cogn Sci* 35: 184-197.
- 21 Mitterer, H., T. Cho, and S. Kim (2016). How does prosody influence speech categorization?
22 *Journal of Phonetics* 54: 68-79.
- 23 Mitterer, H., and J. M. McQueen (2009). Foreign subtitles help but native-language subtitles
24 harm foreign speech perception. *PloS One* 4: e7785.
- 25 Mitterer, H., E. Reinisch, and J. M. McQueen (2018). Allophones, not phonemes in spoken-
26 word recognition. *Journal of Memory and Language* 98: 77-92.
- 27 Morrill, T. (2012). Acoustic correlates of stress in English adjective-noun compounds.
28 *Language and Speech* 55: 167-201.
- 29 Morrill, T., M. Baese-Berk, C. Heffner, and L. Dilley (2015). Interactions between distal
30 speech rate, linguistic knowledge, and speech environment. *Psychonomic Bulletin*
31 *and Review* 22: 1451-1457.
- 32 Morrill, T., L. Dilley, and J. D. McAuley (2014). Prosodic patterning in distal speech context:
33 Effects of list intonation and f0 downtrend on perception of proximal prosodic
34 structure. *Journal of Phonetics* 46: 68-85.
- 35 Morrill, T., L. Dilley, J. D. McAuley, and M. Pitt (2014). Distal rhythm influences whether or
36 not listeners hear a word in continuous speech: support for a perceptual grouping
37 hypothesis. *Cognition* 131: 69-74.
- 38 Nakatani, L. H., and J. A. Schaffer (1978). Hearing "words" without words: Prosodic cues for
39 word perception. *Journal of the Acoustical Society of America* 63: 234-245.
- 40 Newman, R. S., and J. R. Sawusch (1996). Perceptual normalization for speaking rate: effects
41 of temporal distance. *Perception & Psychophysics* 58: 540-560.
- 42 Norris, D., and J. M. McQueen (2008). Shortlist B: A Bayesian model of continuous speech
43 recognition. *Psychological Review* 115: 357-395.
- 44 Norris, D., J. M. McQueen, and A. Cutler (2003). Perceptual learning in speech. *Cognitive*
45 *Psychology* 47: 204-238.
- 46 Norris, D., J. M. McQueen, and A. Cutler (2016). Prediction, Bayesian inference and feedback
47 in speech recognition. *Language, Cognition and Neuroscience* 31: 4-18.

- 1 Norris, D., J. M. McQueen, A. Cutler, and S. Butterfield (1997). The possible-word constraint
2 in the segmentation of continuous speech. *Cognitive Psychology* 34: 191-243.
- 3 Pickering, M. J., and S. Garrod (2013). An integrated theory of language production and
4 comprehension. *Behavioral and Brain Sciences* 36: 329-392.
- 5 Pitt, M., C. Szostak, and L. Dilley (2016). Rate dependent speech processing can be speech
6 specific: Evidence from the perceptual disappearance of words under changes in
7 context speech rate. *Attention, Perception, and Psychophysics* 78: 334-345.
- 8 Pitt, M. A., and A. G. Samuel (1990). The use of rhythm in attending to speech. *Journal of*
9 *Experimental Psychology: Human Perception and Performance* 16: 564-573.
- 10 Poellmann, K., H. R. Bosker, J. M. McQueen, and H. Mitterer (2014). Perceptual adaptation to
11 segmental and syllabic reductions in continuous spoken Dutch. *Journal of Phonetics*
12 46: 101-127.
- 13 Poeppel, D., W. J. Idsardi, and V. van Wassenhove (2008). Speech perception at the interface
14 of neurobiology and linguistics. *Philosophical Transactions of the Royal Society B* 363:
15 1071-1086.
- 16 Quené, H., and R. F. Port (2005). Effects of timing regularity and metrical expectancy on
17 spoken-word perception. *Phonetica* 62: 1-13.
- 18 Ramachers, S. T. M. R. (2018). *Setting the tone: Acquisition and processing of lexical tone in*
19 *East-Limburgian dialects of Dutch*. Utrecht: LOT.
- 20 Reinisch, E., A. Jesse, and J. M. McQueen (2010). Early use of phonetic information in spoken
21 word recognition: Lexical stress drives eye movements immediately. *Quarterly*
22 *Journal of Experimental Psychology* 63: 772-783.
- 23 Repp, B. H. (1982). Phonetic trading relations and context effects: New experimental
24 evidence for a speech mode of perception. *Psychological Bulletin* 92: 81-110.
- 25 Rohde, H., and C. Kurumada (2018). Alternatives and inferences in the communication of
26 meaning. In C. Fedemeier & D. Watson (Eds.), *Psychology of Learning and Motivation*
27 (Vol. 68, pp. 215-252).
- 28 Salverda, A. P., D. Dahan, and J. M. McQueen (2003). The role of prosodic boundaries in the
29 resolution of lexical embedding in speech comprehension. *Cognition* 90: 51-89.
- 30 Samuel, A. G., and T. Kraljic (2009). Perceptual learning for speech. *Attention, Perception, &*
31 *Psychophysics* 71: 1207-1218.
- 32 Sawusch, J. R., and R. S. Newman (2000). Perceptual normalization for speaking rate II:
33 Effects of signal discontinuities. *Perception and Psychophysics* 62: 285-300.
- 34 Selkirk, E. O. (1984). *Phonology and Syntax: The Relation Between Sound and Structure*.
35 Cambridge, MA: MIT Press.
- 36 Shannon, R. V., F. G. Zeng, V. Kamath, J. Wygonski, and M. Ekelid (1995). Speech recognition
37 with primarily temporal cues. *Science* 270: 303-304.
- 38 Shatzman, K. B., and J. M. McQueen (2006). Prosodic knowledge affects the recognition of
39 newly acquired words. *Psychological Science* 17: 372-377.
- 40 Shields, J. L., A. McHugh, and J. G. Martin (1974). Reaction time to phoneme targets as a
41 function of rhythmic cues in continuous speech. *Journal of Experimental Psychology*
42 102: 250-255.
- 43 Shuai, L., and J. G. Malins (2017). Encoding lexical tones in jTRACE: a simulation of
44 monosyllabic spoken word recognition in Mandarin Chinese. *Behavior Research*
45 *Methods* 49: 230-241.

- 1 Sjerps, M. J., C. Zhang, and G. Peng (2018). Lexical tone is perceived relative to locally
2 surrounding context, vowel quality to preceding context. *Journal of Experimental*
3 *Psychology: Human Perception and Performance* 44: 914-924.
- 4 Slowiaczek, L. M. (1990). Effects of lexical stress in auditory word recognition. *Language and*
5 *Speech* 33: 46-68.
- 6 Slowiaczek, L. M. (1991). Stress and context in auditory word recognition. *Journal of*
7 *Psycholinguistic Research* 20: 465-481.
- 8 Sluijter, A. M. C., and V. J. van Heuven (1996). Spectral balance as an acoustic correlate of
9 linguistic stress. *Journal of the Acoustical Society of America* 100: 2417-2485.
- 10 Söderström, P., M. Horne, P. Mannfolk, D. van Westen, and M. Roll (2017). Tone-grammar
11 association within words: Concurrent ERP and fMRI show rapid neural pre-activation
12 and involvement of left inferior frontal gyrus in pseudowords. *Brain and Language*
13 174: 119-126.
- 14 Soto-Faraco, S., N. Sebastian-Galles, and A. Cutler (2001). Segmental and suprasegmental
15 mismatch in lexical access. *Journal of Memory and Language* 45: 412-432.
- 16 Sulpizio, S., and J. M. McQueen (2012). Italians use abstract knowledge about lexical stress
17 during spoken-word recognition. *Journal of Memory and Language* 66: 177-193.
- 18 Summerfield, Q. (1981). Articulatory rate and perceptual constancy in phonetic perception.
19 *Journal of Experimental Psychology: Human Perception and Performance* 7: 1074-
20 1095.
- 21 Tabain, M. (2003). Effects of prosodic boundary on /aC/ sequences: Articulatory results.
22 *Journal of the Acoustical Society of America* 113: 2834-2849.
- 23 Tabossi, P., S. Collina, M. Mazzetti, and M. Zoppello (2000). Syllables in the processing of
24 spoken Italian. *Journal of Experimental Psychology: Human Perception and*
25 *Performance* 26: 758-775.
- 26 Tagliapietra, L., and J. M. McQueen (2010). What and where in speech recognition:
27 Geminates and singletons in spoken Italian. *Journal of Memory and Language* 63:
28 306-323.
- 29 Tavano, A., and M. Scharinger (2015). Prediction in speech and language processing. *Cortex*
30 68: 1-7.
- 31 Tremblay, A., M. Broersma, and C. E. Coughlin (2018). The functional weight of a prosodic
32 cue in the native language predicts the learning of speech segmentation in a second
33 language. *Bilingualism: Language and Cognition* 21: 640-652.
- 34 Tremblay, A., M. Broersma, C. E. Coughlin, and J. Choi (2016). Effects of the native language
35 on the learning of fundamental frequency in second-language speech segmentation.
36 *Frontiers in Psychology* 7: 985.
- 37 Turk, A. E., and White, L. (1999). Structural influences on accentual lengthening in English.
38 *Journal of Phonetics* 27: 171-206.
- 39 Turk, A. E., and S. Shattuck-Hufnagel (2000). Word-boundary-related duration patterns in
40 English. *Journal of Phonetics* 28: 397-440.
- 41 Wade, T., and L. L. Holt (2005). Perceptual effects of preceding nonspeech rate on temporal
42 properties of speech categories. *Perception & Psychophysics* 67: 939-950.
- 43 Wang, B., Y. Xu, and Q. Ding (2018). Interactive prosodic marking of focus, boundary and
44 newness in Mandarin. *Phonetica* 75: 24-56.
- 45 Wong, P. C. M., and R. L. Diehl (2003). Perceptual normalization for inter- and intra-talker
46 variation in Cantonese level tones. *Journal of Speech, Language, and Hearing*
47 *Research* 46: 413-421.

- 1 Yamamoto, H., and E. Haryu (2018). The role of pitch patterns in Japanese 24-month olds'
- 2 word recognition. *Journal of Memory and Language* 99: 90-98.
- 3 Ye, Y., and C. M. Connine (1999). Processing spoken Chinese: The role of tone information.
- 4 *Language and Cognitive Processes* 14: 609-630.

5