

Active listening

Active Listening

1

2

3 Karl J. Friston, Noor Sajid, David Ricardo Quiroga-Martinez, Thomas Parr, Cathy J. Price, Emma

4

Holmes

5

6 The Wellcome Centre for Human Neuroimaging, UCL Queen Square Institute of Neurology, London, UK

7

WC1N 3AR.

8 Emails: k.friston@ucl.ac.uk, noor.sajid.18@ucl.ac.uk, dquiroga@clin.au.dk, thomas.parr.12@ucl.ac.uk,

9

c.j.price@ucl.ac.uk, emma.holmes@ucl.ac.uk

10

11 Address for correspondence:

12 Emma Holmes, emma.holmes@ucl.ac.uk

13 The Wellcome Centre for Human Neuroimaging,

14 UCL Queen Square Institute of Neurology,

15 London, UK WC1N 3AR.

Active listening

16

Abstract

17 This paper introduces active listening, as a unified framework for synthesising and recognising speech. The
18 notion of *active listening* inherits from active inference, which considers perception and action under one
19 universal imperative: to maximise the evidence for our (generative) models of the world. First, we describe
20 a generative model of spoken words that simulates (i) how discrete lexical, prosodic, and speaker attributes
21 give rise to continuous acoustic signals; and conversely (ii) how continuous acoustic signals are recognised
22 as words. The ‘active’ aspect involves (covertly) segmenting spoken sentences and borrows ideas from
23 active vision. It casts speech segmentation as the selection of internal actions, corresponding to the
24 placement of word boundaries. Practically, word boundaries are selected that maximise the evidence for an
25 internal model of how individual words are generated. We establish face validity by simulating speech
26 recognition and showing how the inferred content of a sentence depends on prior beliefs and background
27 noise. Finally, we consider predictive validity by associating neuronal or physiological responses, such as
28 the mismatch negativity and P300, with belief updating under active listening, which is greatest in the
29 absence of accurate prior beliefs about what will be heard next.

30

31

32 **Key words:** speech recognition, voice, active inference, active listening, segmentation, variational Bayes,
33 audition.

Active listening

34

Introduction

35 This paper could be read at three complementary levels: it could be regarded as a foundational paper
36 introducing a *generative model* of spoken word sequences and an accompanying inversion (i.e., word
37 recognition) scheme that has some biological plausibility; e.g., (Kleinschmidt and Jaeger 2015).
38 Alternatively, one could read this article as a proposal for a speech recognition scheme based upon first
39 (Bayesian) principles; e.g., (Rosenfeld 2000). Finally, one could regard this work as computational
40 neuroscience, which makes some predictions about the functional brain architectures that mediate
41 hierarchical auditory perception, when listening or repeating spoken words; e.g., (Hickok and Poeppel
42 2007, Houde and Nagarajan 2011, Tourville and Guenther 2011, Ueno, Saito et al. 2011). In the latter
43 setting, the generative model can be used to predict the effects of synthetic lesions, i.e., as the basis for
44 computational neuropsychology. In other words, one could optimise the parameters of the active listening
45 scheme described below to best explain empirical (electrophysiological or behavioural) responses of
46 individual subjects. We hope to pursue this in subsequent work. The current paper focuses on the form of
47 the generative model, the accompanying recognition or inference scheme, and the kinds of behavioural and
48 neuronal responses it predicts.

49 Speech recognition is not a simple problem. The auditory system receives a continuous acoustic signal and,
50 in order to understand the words that are spoken, must parse a continuous signal into discrete words. To a
51 naïve listener, the acoustic signal provides few cues to indicate where words begin and end (Altenberg
52 2005, Thiessen and Erickson 2013). Furthermore, even when word boundaries are made clear, there exists
53 a many-to-many mapping between lexical content and the acoustic signal. This is because speech is not
54 ‘invariant’ (Liberman, Cooper et al. 1967)—words are never heard out of a particular context. When
55 considering how words are generated, there is wide variability in the pronunciation of the same word among
56 different speakers (Hillenbrand, Getty et al. 1995, Remez 2010)—and even when spoken by the same
57 speaker, pronunciation depends on prosody (Bänziger and Scherer, 2005). From the perspective of
58 recognition, two signals that are acoustically identical can be perceived as different words or phonemes by
59 human listeners, depending on their context—for example, the preceding words or phonemes (Mann 1980,
60 Miller, Green et al. 1984), preceding spectral content (Holt, Lotto et al. 2000), or the duration of a vowel
61 that follows a consonant (Miller and Liberman 1979). The current approach considers the processes
62 involved in segmenting speech—and inferring the words that were spoken—as complementary.

Active listening

63 The idea that speech segmentation and lexical inference operate together did not figure in early accounts of
64 speech recognition. For example, the Fuzzy Logic Model of Perception (FLMP) (Oden and Massaro 1978,
65 Massaro 1987, Massaro 1989) matches acoustic features with prototype representations to recognise
66 phonemes, even when considered in the context of words and sentences. Similarly, the Neighbourhood
67 Activation Model (NAM) (Luce 1986, Luce and Pisoni 1998) considers individual word recognition; it
68 accounts for effects of word frequency, but does not address the segmentation problem. Later connectionist
69 accounts, such as TRACE (McClelland and Elman 1986), assumed that competition between lexical nodes
70 drives recognition, where competition is mediated by inhibitory connections between nodes: bottom-up
71 cues determine recognition of phonemes and top-down cues take into account the plausible words in the
72 lexicon. Shortlist B (Norris and McQueen 2008) reformulates this problem as one of an optimal Bayesian
73 observer and incorporates word frequency effects.

74 Implicit in these connectionist and Bayesian accounts is the idea that speech segmentation depends on
75 words in the listener’s lexicon. For example, word recognition under TRACE assumes that speech will be
76 segmented into words rather than combinations of words and non-words. However, it does not explain how
77 alternative segmentations leading to valid word combinations are reconciled—for example, distinguishing
78 “Grade A” from “grey day”. This example is problematic for the above accounts, because the two
79 segmentations are phonetically identical, acoustically similar, and are both valid word combinations in
80 English. Early accounts also ignored the problem of converting the acoustic signal into words or phonemes.
81 Specifically, they assume that phonetic features (McClelland and Elman 1986) or acoustic features that
82 underlies perceptual confusions in human listeners (NAM; Shortlist B) have already been successfully
83 extracted from the signal. In short, accounts of inputs that are not continuous acoustic signals cannot explain
84 findings that acoustically identical signals are perceived as different words or phonemes depending on their
85 context (Miller and Liberman 1979, Mann 1980, Holt, Lotto et al. 2000).

86 Here, we consider speech recognition as a Bayesian inference problem. We introduce a simplified
87 generative model that maps from the continuous acoustic signal (i.e., a time varying auditory signal or
88 spectral fluctuations containing particular formant frequencies) to discrete words using lexical, speaker,
89 and prosodic information. Generating continuous states from a succession of discrete states is a non-trivial
90 issue for a first principle (i.e., ideal Bayesian observer) approach. However, the requisite neuronal message
91 passing can be solved by combining variational (marginal) message passing and predictive coding (a.k.a.
92 Bayesian filtering). This allows one to simulate perception using generative models that entertain mixtures

Active listening

93 of continuous and discrete states (Friston, Parr et al. 2017, Friston, Rosch et al. 2017).

94 Previous Bayesian accounts (e.g., Shortlist B: Norris and McQueen 2008) have assumed that listeners use
95 exact Bayesian inference. However, performing the calculations required for exact inference would be
96 difficult for biological systems like ourselves, given the complexity of the speech generation process; see
97 (Friston 2010, Bogacz 2017, Friston, FitzGerald et al. 2017). Appealing to variational inference (Beal 2003)
98 affords a much simpler implementation, which has been applied to a variety of other domains in human
99 perception and cognition (Brown, Friston et al. 2011, Brown, Adams et al. 2013, Parr and Friston 2017).
100 Consequently, speech recognition becomes an optimisation problem that corresponds to minimising
101 variational free energy—or, equivalently, maximising the evidence for a particular generative model.

102 In this paper, we provide a computational perspective on the segmentation problem—addressing the
103 challenge that there are often several ways in which a sentence can be parsed, and multiple segmentations
104 engender valid word combinations. We therefore treat speech recognition as a problem of selecting the
105 most appropriate segmentation among several alternatives. We assume that the listener selects the
106 segmentation that is least surprising from the perspective of their generative model. In doing so, we cast
107 segmentation as an internal action that selects among competing hypotheses for the most likely causes of
108 the acoustic signal. Although this is a novel computational implementation of speech segmentation, it aligns
109 with the basic idea that competing segmentations are held in working memory before a listener decides on
110 the most appropriate segmentation, as supported by behavioural studies of word recognition in human
111 listeners (Shillcock 1990, Davis, Marslen-Wilson et al. 2002). This idea is similar to that used in previous
112 accounts such as TRACE and Shortlist B. Here, we address the problem of selecting among multiple
113 segmentations of valid word combinations. Our approach accounts for contextual effects using priors; we
114 show that alternative segmentations—such as “Grade A” and “grey day”—can be accounted for by
115 appealing to these (e.g., semantic or contextual) priors.

116 Conceptualising speech segmentation as an internal (covert) action appeals to the ‘active’ aspect of
117 listening. It is distinct from ‘passive’ listening, which—if truly passive—would not require mental or covert
118 actions. This conceptualisation is grounded in active inference, which has previously been applied to active
119 vision (Grossberg, Roberts et al. 1997, Davison and Murray 2002, Ulanovsky and Moss 2008,
120 Andreopoulos and Tsotsos 2013, Ognibene and Baldassarre 2014, Mirza, Adams et al. 2016, Parr and
121 Friston 2017, Veale, Hafed et al. 2017). Here, we consider the covert placement of word boundaries from

Active listening

122 the same computational perspective as has been used to model an observer whose task is to decide where
123 to sample the visual scene by making overt saccades (Mirza, Adams et al. 2016, Parr and Friston 2017).
124 The types of computations in this framework therefore appeal to general principles that the brain may use
125 to solve a variety of problems.

126 This paper comprises four sections, which each describe different elements of active listening. The first
127 section reviews active inference and then describes a simplified but plausible generative model of how
128 (continuous) sound waves are generated from a discrete word with particular (discrete) attributes. The
129 attributes include lexical content, prosody, and speaker characteristics. The division of attributes into
130 lexical, prosodic, and speaker attributes is logical from a generative perspective—and is consistent with
131 neuropsychological studies showing selective deficits in the processing of these attributes (Miller and
132 Liberman 1979, Peretz, Kolinsky et al. 1994). Indeed, these attributes have been considered fundamental
133 characteristics in qualitative models of speech perception such as the ‘auditory face’ model (Belin, Fecteau
134 et al. 2004)—and are known to interact to affect human speech perception (Nygaard, Sommers et al. 1994,
135 Johnsrude, Mackey et al. 2013, Holmes, Domingo et al. 2018). We, therefore, assume these are the types
136 of attributes that human listeners infer when trying to explain the (hidden) causes of an acoustic (speech)
137 signal. This section describes how the generative model can be inverted to determine the most likely lexical,
138 prosodic, and speaker attributes of a word, given a continuous sound wave.

139 The second section deals with the speech segmentation problem, which becomes important when
140 recognising words within sentences, rather than individual words. It considers the question: how do we
141 determine the most likely onsets and offsets of words within a sentence? For example, how do we parse
142 auditory input to disambiguate "Grade A" from "grey day"? To address this question, we use simple acoustic
143 properties to identify plausible word boundaries. We then appeal to the ‘active’ element of active inference,
144 considering the (implicit) placement of word boundaries as a covert ‘action’. This allows us to use
145 established inference schemes to select among competing segmentations (i.e., hypotheses about different
146 word boundaries). These inference schemes essentially ask: which of the possible segmentations minimise
147 free energy or, equivalently, provide the greatest evidence for the listener’s (internal) model of how words
148 are generated? It is at this point that the relationship between the generative model from the first section
149 and ‘active’ speech segmentation becomes clear: these different elements work in unison when inferring
150 words within a sentence. The generative model operates at the individual word level, whereas speech
151 segmentation operates at the sentence level: the best speech segmentation will maximise the combined

Active listening

152 evidence for attributes of constituent words. This section concludes with an illustration of the face validity
153 of the active listening scheme by comparing speech recognition (i.e., lexical inference) with and without
154 prior beliefs about the sequence of plausible words that could be encountered—demonstrating how different
155 segmentations that contain valid English words can be disambiguated.

156 The third section highlights an aspect of speech recognition that has not been simulated under previous
157 accounts. We show that a quantity within active listening can predict neurophysiological responses of the
158 sort measured by electromagnetic recordings (Hasson, Yang et al.) or functional magnetic resonance
159 imaging (fMRI). In particular, the magnitude of belief updating in active listening appears to capture the
160 fluctuations in evoked (or induced) responses that have been demonstrated empirically; e.g., the mismatch
161 negativity (Garrido, Kilner et al. 2009, Morlet and Fischer 2014), P300 (Donchin and Coles 1988, Morlet
162 and Fischer 2014), and N400 (Kutas and Hillyard 1980). Broadly speaking, this suggests that elements of
163 speech perception are consistent with predictive coding (see (Poeppel and Monahan 2011) for a review).
164 Formally, belief updating is related to the difference between *prior* beliefs about states in the generative
165 model to *posterior* beliefs. In other words, the amount that beliefs change after sampling sensory evidence.
166 This is variously known as *Bayesian surprise*, salience, information gain, or complexity. In this section, we
167 illustrate the similarity between belief updates and violation responses, showing that the magnitude of belief
168 updating depends upon prior expectations about particular words in the lexicon (Cole, Jakimik et al. 1980,
169 Mattys and Melhorn 2007, Mattys, Melhorn et al. 2007, Kim, Stephens et al. 2012) and the quality of
170 sensory evidence; e.g., when speech is acoustically masked by background noise (“speech-in-noise”)
171 (Sams, Paavilainen et al. 1985, Winkler, Denham et al. 2009). We conclude by discussing how the model
172 could be developed for future applications, and its potential utility in the cognitive neuroscience (and
173 neuropsychology) of auditory perception and language.

174 **A generative model of spoken words**

175 Active inference is a first principle account of action and perception in sentient creatures (Friston,
176 FitzGerald et al. 2017). It is based upon the idea that synaptic activity, efficacy and connectivity all change
177 to maximise the evidence for a model of how our sensations are generated. Formally, this means treating
178 neuronal dynamics as a *gradient flow* on a quantity that is always greater than (negative) log evidence
179 (Friston, Parr et al. 2017). This quantity is known as variational free energy in physics and statistics

Active listening

180 (Feynman 1972, Hinton and Zemel 1993). The complement (i.e., negative) of this quantity is known as an
181 evidence lower bound (ELBO) in machine learning (Winn and Bishop 2005). A gradient flow is simply a
182 way of writing down dynamics in terms of equations of motion that ensure a certain function is minimised—
183 in this case, variational free energy. The resulting dynamics furnish a model of neuronal fluctuations (and
184 changes in synaptic efficacy and connectivity) that necessarily minimise free energy or maximise model
185 evidence. In short, if one simulates speech recognition using active inference, one automatically provides
186 an account of the accompanying neuronal dynamics.

187 This approach to understanding and modelling (active) inference in the brain has been applied in many
188 settings, using exactly the same schemes and principles. The only thing that distinguishes one application
189 from another is the form of the generative model. In other words, if one can write down a probabilistic
190 model of how some sensory input was generated, one can invert the model—using standard model inversion
191 schemes—to simulate neuronal dynamics and implicit belief updating in the brain: See (Friston, Parr et al.
192 2017) for a detailed summary of these schemes that cover models of both discrete and continuous states
193 generating sensations. See also (Bastos, Usrey et al. 2012, Friston, FitzGerald et al. 2017) for a discussion
194 of neurobiological implementation, in terms of attending process theories, for continuous and discrete state
195 space models, respectively.

196 In this section, we focus on the form of a (simplified) generative model that can be used to generate
197 continuous acoustic signals associated with a particular word. A benefit of this active inference approach
198 is that the generative model can be used to both generate synthetic speech (by applying the forward model)
199 and recognise speech (by inverting the model). The goal is not to provide a state-of-the-art speech synthesis
200 system, but rather to use the generative model and accompanying inference scheme to simulate listening
201 behaviour and neural responses. The work reported in this paper is a prelude to a model of natural language
202 processing, in which the current generative model is equipped with higher levels to enable dyadic
203 exchanges; namely, conversations that entail questions and answers that resolve uncertainty about shared
204 narratives or beliefs. In this paper, we restrict ourselves to inference about sequences of words—and assume
205 that simulated subjects are equipped with prior beliefs about which words are more or less likely in a short
206 sentence or phrase. In a more complex (i.e., deep hierarchical) model, these beliefs would be available from
207 a higher level. These prior beliefs are about the likely semantic content of spoken words; for example, based
208 on previous words in a sentence (Dubno, Ahlstrom et al. 2000) or the topic of conversation (Holmes,
209 Folkeard et al. 2018). Note that previous accounts of speech recognition, such as Shortlist B (Norris and

Active listening

210 McQueen 2008), assume that priors reflect only word frequency, rather than priors that can be flexibly
211 updated based on context. Technically, these kinds of context-sensitive priors are known as empirical
212 priors—and are an integral part of hierarchical generative models.

213 In this paper, we deal with the lowest level of the generative model; namely, given a particular lexical
214 content, prosody and speaker identity, how would one generate a spoken word in terms of its acoustic
215 timeseries. In the next section of this paper, we turn to the problem of segmentation (i.e., identifying word
216 boundaries) and the enactive aspects of the current scheme. It will become apparent later on that these two
217 (perceptual and enactive) aspects of active listening go hand-in-hand.

218 Figure 1 summarises the modelling of a spoken word, from the perspectives of generation and recognition.
219 The model considers: how is an acoustic signal generated given the causes of a spoken word, in terms of
220 ‘what’ word is spoken (*lexical*), ‘how’ it is spoken (*prosody*), and ‘who’ is speaking (*speaker identity*)?
221 From the perspective of word generation, it takes *lexical*, *speaker*, and *prosody* parameters and generates
222 an expected acoustic signal. The *lexical* state consists of frequency and temporal coefficients corresponding
223 to words in the lexicon. The model includes two *speaker* states: fundamental frequency and formant scaling.
224 It includes four *prosody* states: amplitude, duration, timbre, and inflection. Within each of these states,
225 different factors correspond to different lexical items, or the fundamental frequency associated with
226 different speakers, for example.

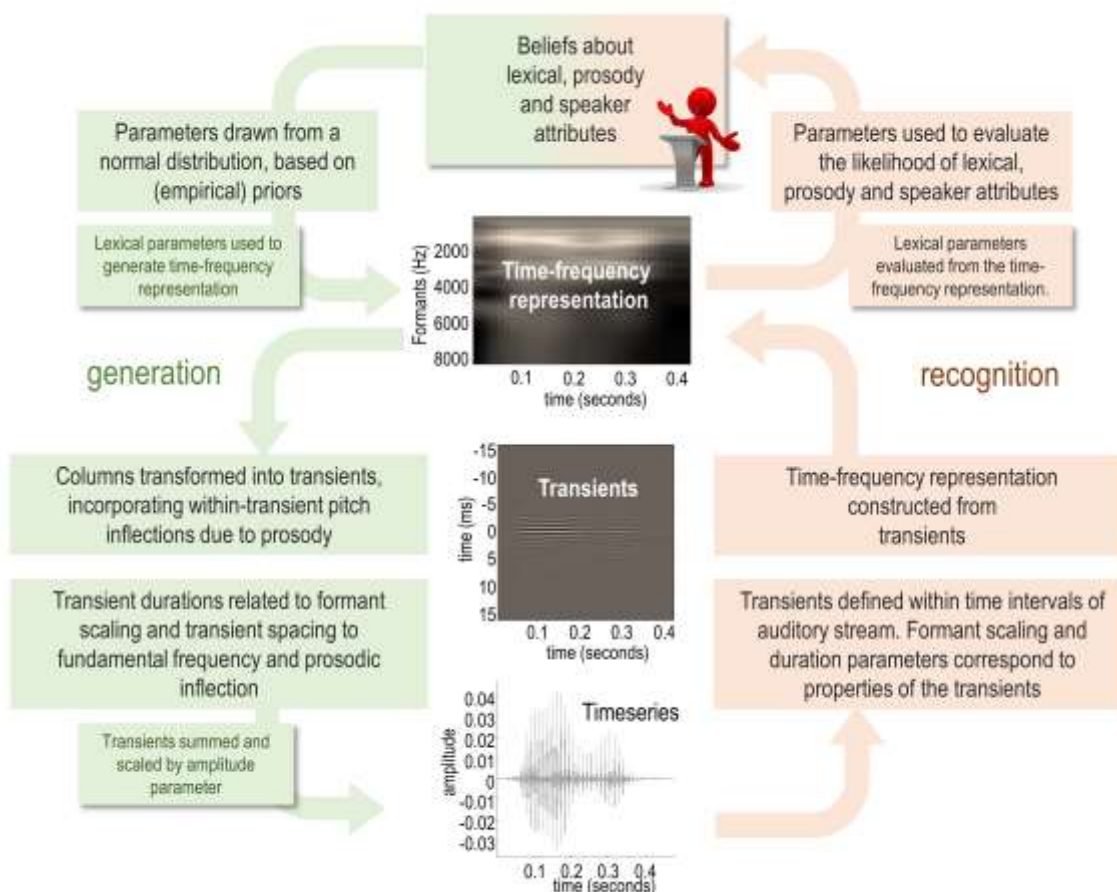
227 The model starts by sampling parameters from a set of probability distributions, which are modelled as
228 separate Gaussians. The means and covariances of the Gaussians have been specified in advance; they can
229 be entered into the model explicitly (by hand) or they can be estimated empirically based on training
230 samples of speech. Sampling parameters from distributions with particular means and variances accounts
231 for the fact that the same lexical item spoken by the same speaker with the same prosody does not always
232 produce an identical acoustic signal, and—conversely—because the distributions are allowed to overlap, a
233 similar acoustic signal can be generated by different combinations of factors. The (discrete) lexical content
234 of a word is sampled from a (categorical) probability distribution over words in a lexicon. This is based on
235 how likely particular words are to be spoken. Ultimately, the selected parameters are combined, in a
236 nonlinear way, to generate an acoustic timeseries corresponding to the articulated word.

237 The acoustic timeseries is generated from a sequence of transients, whose properties are determined by the

Active listening

238 selected parameters. Each word (i.e., *lexical* item) is associated with a matrix of frequency and temporal
239 coefficients (for a discrete cosine transform) that can be used to generate a time-frequency representation
240 of the spoken word (i.e., the spectrogram) when combined with *speaker* and *prosody* information. Each
241 column of the time-frequency representation is used to generate a transient. These transients can be thought
242 of as pulses or ‘shockwaves’ at the glottal pulse rate, which are modulated by the shape of the vocal tract.
243 The instantaneous fundamental frequency is related to the average fundamental frequency of a particular
244 speaker, but also varies smoothly over time based on inflections due to prosody. The prosodic inflection
245 parameters encode: (1) the average fundamental frequency relative to the speaker average, (2) increases or
246 decreases in fundamental frequency over time, and (3) the acceleration or deceleration of changes in
247 fundamental frequency. The instantaneous fundamental frequency determines the spacing of the transients.
248 The durations of the transients are determined by the formant frequencies, which depend on the lexical
249 parameters and the speaker formant scaling parameter. The formant frequencies correspond to the
250 frequency bins in the time-frequency representation. The number of transients that are aggregated to
251 construct the timeseries is determined by the time intervals in the time-frequency representation. Figure 2
252 provides an illustration of how a sequence of transients is generated. In the final step, the transients are
253 summed together and scaled by an amplitude parameter. For mathematical detail, the equations
254 corresponding to the generative model are shown in Figure 11 and are described in Appendix 1. For an
255 algorithmic description, please see the demonstration (annotated Matlab) code—that reproduces the
256 simulations below—which can be read as pseudocode (see Software note).

Active listening



257

258

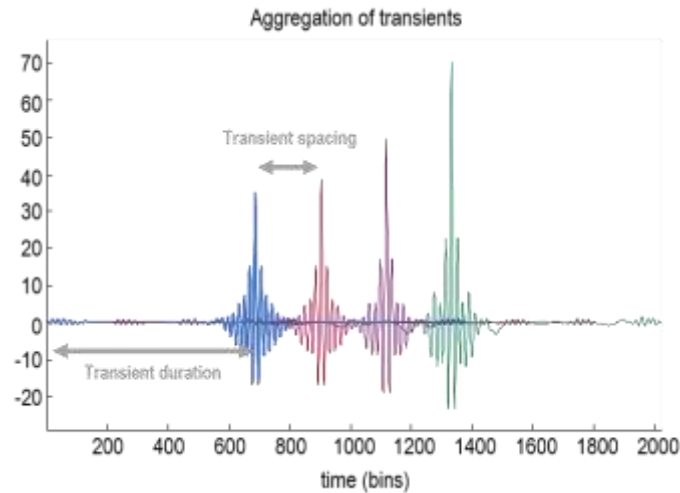
FIGURE 1

259 *A generative model of a word.* This figure illustrates the generative model from the perspective of word generation
260 (green panels) and accompanying inversion (orange panels), which corresponds to word recognition. For the equations
261 describing these probabilistic transformations, please see Appendix 1.

262

263

Active listening



264

265

FIGURE 2

266 *Fundamental and formant intervals.* This figure illustrates the way in which an acoustic timeseries is generated by
267 assembling a succession of transients separated by an interval that is inversely proportional to the (instantaneous)
268 fundamental frequency. The duration of each transient places an upper bound on the wavelength of the formant
269 frequencies—and corresponds to the minimum frequency, which we take to be the first formant frequency.

270

271

272 In effect, the lexical parameters—which, under this generative model, determine the formant frequencies—
273 parameterise a trajectory through high-dimensional formant frequency space, which becomes apparent as
274 the word unfolds. The prosody of the word determines the duration and inflection of the fundamental
275 interval function, while speaker identity determines the average fundamental frequency—which relates to
276 the interval between transients—and a formant scaling parameter that determines the duration of each
277 transient. With such a model in place, one can, in principle, generate any word, spoken with any prosody
278 by any speaker, by sampling the correct parameters from their appropriate distributions. In what follows,
279 we briefly review the inversion of this model given an acoustic timeseries.

Active listening

280 **Model inversion or word recognition**

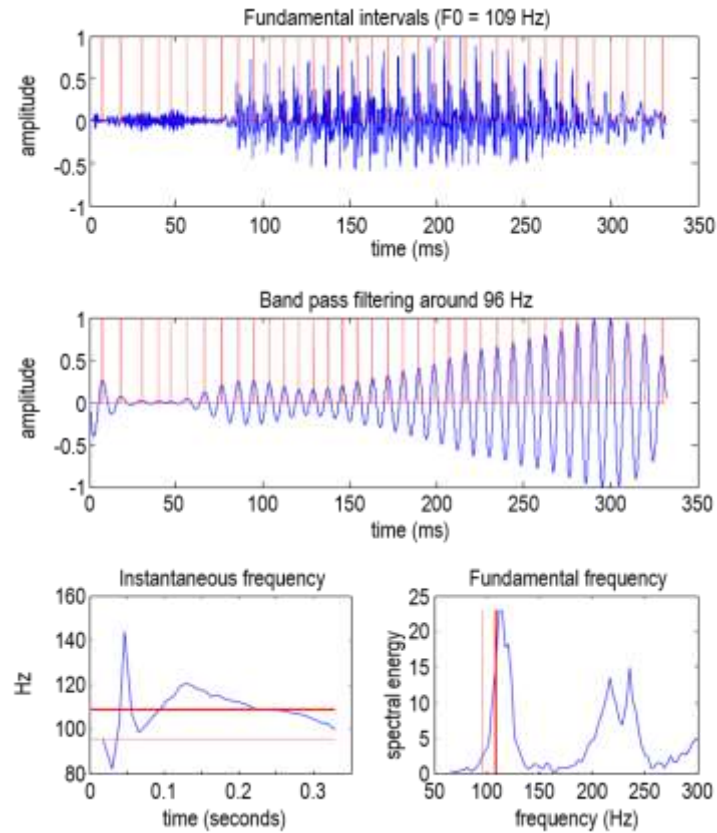
281 Now we have established a generative model that is capable of producing a spoken word, word recognition
282 can be achieved by inverting the model. This section describes a plausible inversion scheme in the context
283 of our particular generative model of spoken words. In principle, given any generative model it should be
284 possible to use Bayesian model inversion to invert the timeseries, using generalised (variational or
285 Bayesian) filtering; also known as predictive coding (Norris, McQueen et al. 2016). However, given we
286 have assumed a deterministic generation of acoustic signals from parameters, we know that the posterior
287 beliefs about parameters will take the form of Dirac delta functions, whose only parameter is a mode. This
288 means that in practice, it is simpler to cache an epoch of the timeseries and use *maximum a posteriori* (Kim,
289 Frisina et al.) estimates of the parameters, based upon least squares. One can then evaluate the posterior
290 probability of discrete *lexical*, *prosody* and *speaker* states, using the respective likelihood of the (Kim,
291 Frisina et al.) parameter estimates (and any priors over discrete states should they be available). This MAP
292 scheme can be read in the spirit of predictive coding that has been *amortised* (Zhang, Butepage et al. 2018).
293 In other words, the inversion scheme reduces to a nonlinear recognition function—a series of equations that
294 map from epochs of the acoustic signal to parameters encoding lexical content, prosody and identity.

295 Model inversion rests on the assumption that we have isolated the acoustic timeseries corresponding to an
296 individual word. The next section deals with the segmentation problem, which involves enactive processes.
297 For now, we will assume that we have identified an epoch of the acoustic signal that might plausibly contain
298 one word—and that we wish to evaluate the probabilities of *lexical*, *prosody*, and *speaker* states within this
299 epoch.

300 In brief, the recognition scheme comprises the following steps (see Figure 1). The instantaneous frequency
301 is estimated by first calculating ‘fundamental intervals’, which are the reciprocal of the instantaneous
302 frequency. The fundamental intervals are calculated by bandpass filtering the acoustic signal around the
303 prior value for the speaker fundamental frequency parameter; the positions of peaks in the filtered signal
304 correspond to the fundamental intervals. Please see Figure 3 for an illustration of how the fundamental
305 intervals are estimated and Figure 4 to see the fundamental frequency and formant frequencies projected
306 onto the spectrum of a speech sample.

307

Active listening



308

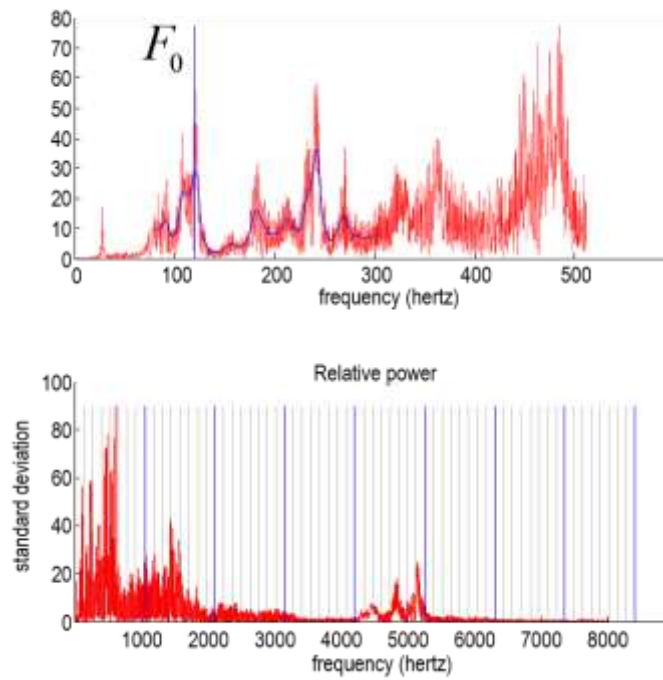
309

FIGURE 3

310 *Fundamental frequencies and intervals.* This figure illustrates the estimation of fluctuations around the fundamental
311 frequency during the articulation of (the first part of) a word. These fluctuations correspond to changes in the
312 fundamental interval; namely, the reciprocal of the instantaneous frequency. The upper panel shows the original
313 timeseries, while the middle panel shows the same timeseries after bandpass filtering. The peaks (i.e., phase crossings)
314 then determine the intervals, which are plotted in terms of instantaneous frequencies on the lower left (as a blue line).
315 The solid red line corresponds to the mean frequency (here, 109 Hz), while the broken red line corresponds to the
316 centre frequency of the bandpass filtering (here, 96 Hz) which is centred on the prior for the speaker average
317 fundamental frequency. The same frequencies are shown on the lower right panel, superimposed on the spectral energy
318 (the absolute values of the accompanying Fourier coefficients of the timeseries in the upper panel). The ensuing
319 fundamental intervals are depicted as red lines in the upper two panels.

320

Active listening



321

322

FIGURE 4

323 *Fundamental and formant frequencies:* Both plots show the root mean square power (i.e., absolute value of Fourier
324 coefficients) following the Fourier transform of a short segment of speech. The frequency range in the upper plot
325 covers the first 500 Hz. The first peak in power (illustrated by the blue vertical line) corresponds to the *fundamental*
326 *frequency*, which is typically between 80 and 150 Hz for adult men and up to 350 Hz for adult women. The lower
327 panel shows the same spectral decomposition but covers 8000 Hz to illustrate formant frequencies. The solid blue
328 lines show the calculated formant frequency and its multiples, while the grey lines arbitrarily divide the frequency
329 intervals into eight bins. These frequencies define the frequencies used for the spectral decomposition.

330

331 Next, the inversion scheme essentially deconstructs transients (i.e., segments) from the epoch. The formant
332 frequencies are estimated by evaluating the cross-covariance function over short segments; the length of
333 the segments is the inverse of the first formant frequency and the segments are centred on each fundamental
334 interval. This is based on the simplifying assumption that the spectral content of each transient, within each
335 segment, is sufficient to generate the word. The formant frequencies are then used to project back to a time-
336 frequency representation.

Active listening

337 To infer the lexical content, prosody and speaker, the parameter estimates from the nonlinear
338 transformations above can be used to evaluate the likelihood of each discrete attribute. This likelihood is
339 then combined with a prior to produce a posterior categorical distribution over the attributes in question.
340 For the *lexical* content of the word, this just corresponds to an index in the lexicon. Here, the lexicon is
341 assumed to be small for simplicity, although it would be trivial to extend the model to accommodate more
342 comprehensive lexicons. The likelihood is based upon the mean and precision (i.e., inverse covariance) of
343 the lexical parameters in the usual way, where the sufficient statistics of this (likelihood) model—for each
344 word—are evaluated using some exemplar or training set of words. This completes the description of word
345 recognition based upon the generative model above. For details of the equations used in model inversion,
346 please see Appendix 2.

347 In summary, the above transformations simply reverse the operations used for word generation in the
348 previous section. The combination of prior expectations with the likelihoods of each attribute is a key
349 feature of this inversion scheme that will allow the model to accommodate contextual effects on speech
350 recognition. In other words, we are more likely to interpret speech consistent with our prior expectations.
351 This will become evident in the simulations later in this paper.

352 After the discrete parameters have been inferred from a continuous timeseries through model inversion,
353 they could be entered back into the generative model to synthesise a new timeseries that would share some
354 properties with the timeseries that was used to infer the discrete parameters. This simply involves projecting
355 the lexical coefficients back into a time frequency representation, implementing the inverse discrete cosine
356 transform to produce (after scaling with the timbre parameter and exponentiation) a series of (time
357 symmetric) transients, which are aggregated to form the acoustic timeseries. This is essentially what is
358 illustrated in Figure 1. Indeed, the processes of inversion and generation can be iterated (see below) to
359 check the fidelity of the forward and inverse transformations that map between the acoustic timeseries and
360 formant representation.

361 **Speech segmentation as an active process**

362 So far, we have a generative model (and amortised elements of a predictive coding scheme) that generates
363 an appropriate time series, given discrete *lexical* (i.e., what), *prosody* (i.e., how) and *speaker* (i.e., who)

Active listening

364 states (i.e., latent causes of the word). It can also be inverted to infer the attributes of a word given an
365 acoustic timeseries. However, in our everyday lives, we usually hear series of words rather than words in
366 isolation. In this section, we combine the generative model with an active segmentation process, to infer
367 the most likely *sequence* of words given a continuous timeseries.

368 This requires us to address the following problem: we have not specified how the onsets and offsets of the
369 interval containing the word are generated (i.e., when). Clearly, there are some prior constraints on the
370 generation of these intervals. For example, the offset of one word should precede the onset of the subsequent
371 word. Furthermore, the intervals contained between the onset and offset must lie in some plausible time
372 range. We also know that segmentations are more likely to contain words than non-words (Ganong 1980,
373 Billig, Davis et al. 2013), and listeners have prior knowledge of the words that are possible in a language
374 ('possible word constraint') (Norris, McQueen et al. 1997). In the current segmentations, we account for
375 these simple constraints and, effectively, offload inference about word boundaries to the *active* part of active
376 inference. The only acoustic cue we use is the contour of the amplitude envelope, which has previously
377 been identified as a cue that human listeners use for speech segmentation (Lehiste 1960).

378 In brief, we assume that boundary segmentations are not entirely specified by the acoustic signal, and
379 conceptualise the segmentation problem as a problem of choosing which boundaries to select given several
380 possible segmentations; in a similar way as we would select visual actions (e.g., saccadic eye movements
381 or oculomotor pursuit) to fixate or track a visual object given multiple possible actions. In the current
382 setting, this simply means identifying a number of plausible boundary intervals and finding the interval that
383 provides the greatest evidence for our prior beliefs about the words we hear. This is the same principle used
384 to explain motor and autonomic action under active inference (Friston, Mattout et al. 2011). For example,
385 classical motor reflexes can be construed as minimising proprioceptive prediction error (i.e. minimising
386 variational free energy or maximising model evidence) as described in (Adams, Shipp et al. 2013). Formally
387 identical arguments have been applied in the setting of interoceptive inference where motor reflexes are
388 replaced by autonomic reflexes that realise autonomic set-points or homoeostasis (Seth 2014).

389 In the current context, we essentially treat the decision about speech segmentation as a covert action from
390 a computational perspective, which shares similarities with the overt actions used in other settings. This
391 can be implemented in a straightforward fashion by selecting boundary pairs (i.e., offsets and onsets) and
392 evaluating their free energy under some prior beliefs about the next word. Ultimately, we want to select the

Active listening

393 boundary pairs with the smallest free energy—which effectively selects the interval with the greatest
394 evidence (a.k.a., marginal likelihood) of auditory outcomes contained in that interval. This follows because
395 the variational free energy, by construction, represents an upper bound on log evidence (see Appendix 3
396 for more details and the corresponding equations). Importantly, both posterior beliefs about latent states
397 (i.e., *lexical*, *prosody*, and *speaker*) and the active selection of acoustic intervals optimise free energy. This
398 is the signature of active inference. In this instance, the posterior beliefs obtain from the likelihood of the
399 lexical, prosody and identity parameters, given the associated states.

400 For words spoken in isolation, one can identify candidate boundaries using threshold crossings of the
401 amplitude envelope (where the threshold is a low value, roughly corresponding to the noise floor).
402 However, it is well known that a continuous stream of words does not always contain ‘silent’ (i.e., below-
403 threshold) gaps between words and, conversely, silence can occur between two syllables of the same word.
404 We therefore include local minima of the amplitude envelope as candidate boundaries. It is important to
405 note that these are only *candidate* boundaries—in other words, plausible hypotheses for segmentations of
406 the acoustic signal. We will turn to the question of which interval is *selected* later, during which candidate
407 segmentations are combined with (lexical) priors. In practice, this means that two syllables separated by a
408 silent gap are not always classified as separate words—consistent with the knowledge that naturally spoken
409 words often contain silent gaps that—to a naïve listener—could be confused with word boundaries. An
410 example of the candidate boundary points is illustrated in Figure 5. Please see figure legend for details.

411 Using this procedure to identify candidate intervals, one can select the interval that minimises free energy
412 (or has the greatest evidence under prior beliefs about the next word). In other words, for each candidate
413 interval, the likelihood of the lexical parameters is evaluated—for all plausible words—to create a belief
414 over lexical content, in terms of a probability distribution. This posterior belief is then used to evaluate the
415 log evidence (i.e., free energy) of each interval. The interval (and associated posterior beliefs) with the
416 greatest evidence is selected. The offset of this interval specifies the onset of the next segment and the
417 process starts again.

418 Treating speech segmentation as a problem of (covertly) sampling among plausible intervals is interesting
419 from a mathematical perspective. The free energy associated with a particular action is a trade-off between
420 the accuracy of sensory observations under the generative model and the complexity of belief updating on
421 the basis of those observations (see Appendix 3 for the equations). In the current setting, these quantities

Active listening

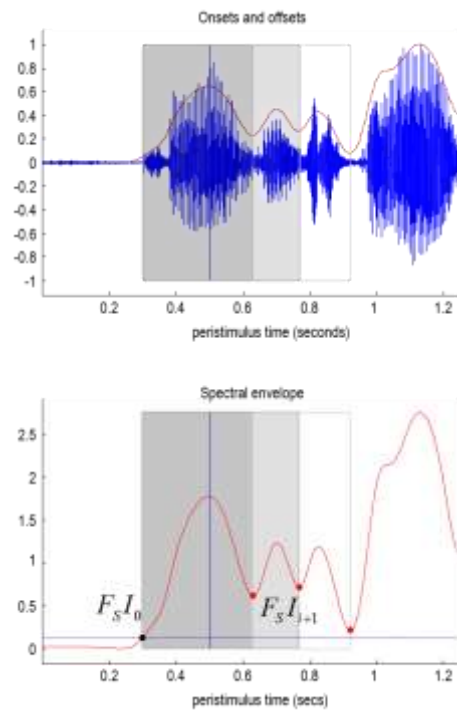
422 can be evaluated explicitly, because the evidence has already been accumulated. Thus, the accuracy term
423 simply scores the expected log likelihood of the auditory observations under posterior beliefs about the
424 lexical categories that generated them. The complexity term scores the difference between the prior beliefs
425 and the new beliefs based on auditory observations. This will become an important quantity later and,
426 essentially, reflects the degree of belief updating associated with selecting one lexical parsing over another.
427 Phrased another way, the goal of segmentation under active listening is to sample data in a way that requires
428 the most parsimonious degree of belief updating, in accord with Ockham's principle (Maisto, Donnarumma
429 et al. 2015).

430 Figure 6 shows the consequence of this form of active listening by comparing segmentation and recognition
431 with and without appropriate prior beliefs (please see the figure legend for details). The input to this
432 simulation is a continuous acoustic signal that has alternative parsings, leading to different lexical
433 segmentations. The timeseries in Figures 6A and 6E are identical, but the segmentation (as indicated by the
434 colours) differs. The point of this simulation is to show that the selected segmentation depends on the
435 distribution of the priors. When the artificial listener has no particular prior beliefs about which words will
436 be heard (left panel), the priors are uniform, and recognition goes awry after the first two words (“triangle
437 square”). The scheme inferred that the best possible explanation for the subsequent words was a series of
438 shorter words (“a is red a is red”; Figure 6B). From Figure 6C, we can tell that the artificial listener was
439 uncertain about the correct parsing—reflecting the fact that this signal was difficult to segment because
440 there were several parsings that would be plausible in English (displayed as grey shaded regions). However,
441 when the artificial listener was equipped with strong prior beliefs that the words they would hear would be
442 shape words (the words “triangle” and “square”), it recovered the correct parsing (“triangle square triangle
443 square triangle square”; Figure 6F). Note that the acoustic boundaries for these two lexical segmentations
444 differ—highlighting that speech segmentation and lexical inference go hand-in-hand, under this framework.

445

446

Active listening



447

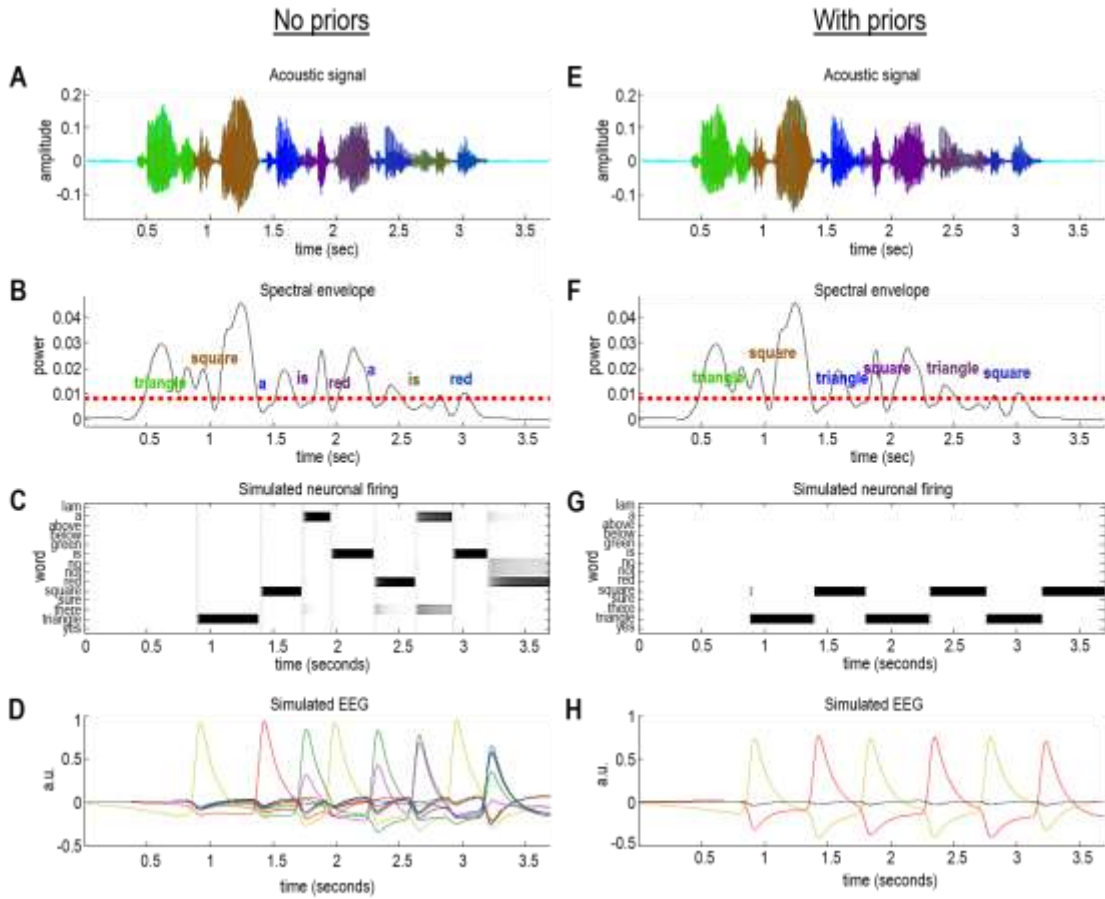
448

FIGURE 5

449 *Spectral envelopes and segment boundaries.* This figure provides an example of how candidate intervals containing
450 words are identified using the spectral envelope. The upper panel shows a timeseries produced by saying "triangle,
451 square". The timeseries is high pass filtered and smoothed using a Gaussian kernel. The dotted red line in the upper
452 panel shows the resulting spectral envelope, after subtracting the minimum. The broken line corresponds to a
453 threshold: $1/16^{\text{th}}$ of the maximum encountered during the (1250 ms) epoch. This envelope is reproduced in the lower
454 panel (red line). Boundaries are then identified as the first crossing (black dot) of the threshold (horizontal blue line)
455 before the spectral peak and the last crossing after the peak. These boundaries are then supplemented with the internal
456 minima between the peak and offset (red dots). These boundaries then generate a set of intervals for subsequent
457 selection during the recognition or inference process. Here, there are three such intervals. The first contains the first
458 two syllables of triangle, the second contains the word "triangle". The third additionally includes the first phoneme of
459 "square". In this example, the second interval was selected as the most plausible (i.e., free energy reducing) candidate
460 to correctly infer that this segment contained the word "triangle". The vertical blue line corresponds to the first spectral
461 peak following the offset of the last word, which provides a lower bound on the onset.

462

Active listening



463

464

FIGURE 6

465 *Speech recognition and segmentation. Left panel:* This panel shows the results of active listening to a sequence of
 466 words: a succession of “triangle, square, triangle, square...”. Its format will be used in subsequent figures and is
 467 described in detail here. Panel A shows the acoustic timeseries as a function of time in seconds. The different colours
 468 correspond to the segmentation selected by the active listening scheme, with each colour corresponding to an inferred
 469 word. Regions of cyan denote parts of the timeseries that were not contained within a word boundary. Panel B shows
 470 the accompanying spectral envelope (back line) and the threshold (red dashed line) used to identify subsequent peaks.
 471 The first peak of each successive word centres the boundary identification scheme of Panel A. The words that have
 472 been inferred are shown in the same colours as the upper panel at their (inferred) onset. Panels C–D show the results
 473 of simulated neuronal firing patterns and local field potentials or electroencephalographic responses. These are based
 474 upon a simple form of belief updating cast as a neurally plausible gradient descent on variational free energy (please
 475 see main text). Panel C shows the activity of neuronal populations encoding each potential word (here, 14 alternatives
 476 listed on the Y axis). These are portrayed as starting at the offset of each word. Effectively, these reflect a competition
 477 between lexical representations that record the selection of the most likely explanation. Sometimes this selection is
 478 definitive: for example, the first word (“triangle”) supervenes almost immediately. Conversely, some words induce a
 479 belief updating that is more uncertain. For example, the last word (“red”) has at least three competing explanations
 480 (i.e., “no”, “not” and “a”). Even after convergence to a particular posterior belief, there is still some residual

Active listening

481 uncertainty about whether “red” was heard. Note that the amplitude of the spectral envelope is only just above
482 threshold. In other words, this word was spoken rather softly. Panel D shows the same data after taking the temporal
483 derivative and filtering between 1 and 16 Hz. This reveals fluctuations in (simulated) depolarisation that drives the
484 increases or decreases in neuronal firing of the panels above. In this example, the sequence of words was falsely
485 inferred to be a mixture of several words not actually spoken. This failure to recognise the words reflects the fact that
486 the sequence was difficult to parse or segment. Once segmentation fails, it is difficult to pick up the correct sequence
487 of segmentations that will, in turn, support veridical inference. These results can be compared with the equivalent
488 results when appropriate priors are supplied to enable a more veridical segmentation and subsequent recognition.
489 **Right panel:** This panel shows the results of active listening using the same auditory stream as in the left panel. The
490 only difference here is that the (synthetic) subject was equipped with strong prior beliefs that the only words in play
491 were either “triangle” or “square”. This meant that the agent could properly identify the succession of words, by
492 selecting the veridical word boundaries and, by implication, the boundaries of subsequent words. If one compares the
493 ensuing segmentation with corresponding segmentation in the absence of informative priors, one can see clearly where
494 segmentation failed in the previous example. For example, the last word (i.e., “square”) is correctly identified in dark
495 blue in Panel F. Whereas, in Panel B (without prior constraints), the last phoneme of the word “square” was inferred
496 as “red” and the first phoneme was assigned to a different word (“is”). The comparative analysis of these segmentations
497 highlights the ‘handshake’ between inferring the boundaries in a spectral envelope and correctly inferring the lexical
498 content on the basis of fluctuations in formant frequencies.

499

500

501 These two examples are analogous to the “Grade A” versus “grey day” example that we considered in the
502 introduction. As in our simulated example, there is no consistent acoustic cue that differentiates “Grade A”
503 from “grey day”—and, therefore, priors play an essential disambiguating role. The active segmentation
504 would identify these two (and perhaps additional) possible segmentations, and the percept would be the one
505 that was most similar to the priors. In other words, these two segmentations would be distinguished by
506 different prior beliefs, which could originate from a higher (semantic or contextual) level—for example,
507 whether the topic of conversation was about the weather or a student’s exam results. In a comprehensive
508 treatment, these would be empirical prior beliefs generated by deep temporal models of the sort described
509 in (Kiebel, Daunizeau et al. 2009, Friston, Rosch et al. 2017). For simplicity and focus, we assume here
510 that priors about sequential lexical content—of the sort that could be formed by lexical and semantic
511 predictions—are available to a subject in the form of categorical probability distributions.

512 **Belief updating and neuronal dynamics**

513 Figure 6 includes a characterisation of simulated word recognition in terms of neuronal responses (Figure

Active listening

514 6C–D, G–H). These (simulated) neuronal responses inherit from the neuronal (marginal) message passing
515 scheme described in (Friston, Parr et al. 2017, Parr, Markovic et al. 2019). They reflect belief updating
516 about the lexical category for each word; the simulated neuronal responses are simply the gradient flow on
517 free energy that is associated with belief updating in active listening. The prediction error is the (negative)
518 free energy gradient that drives neuronal dynamics. Mathematically, the prediction error is the difference
519 between the optimal log posterior and current estimate of this. As detailed in Appendix 3, log expectations
520 about hidden states can be associated with depolarisation of neurons or neuronal populations encoding
521 expectations about hidden states, while firing rates encode expectations *per se*.

522 Figure 6 reproduces these simulated neuronal responses following the processing of each word. These
523 responses are shown in terms of spike rates, as would be recorded with single unit electrodes (Figure 6C,
524 G) and depolarisation that would be measured with EEG (Figure 6D, H). Under this formulation, neuronal
525 activity starts off from some prior expectations and evolves, via a gradient flow on free energy (i.e.,
526 prediction error) to encode posterior expectations. Because depolarisation corresponds to the rate of change
527 of these beliefs (expressed as log expectations) they show peak responses during the greatest degree of
528 belief updating from priors to posterior expectations. After filtering, the simulated depolarisations look like
529 evoked responses that are typically observed in human studies (as discussed in more detail below).

530 **Summary**

531 The message from the simulations in Figure 6 is that proper segmentation and subsequent inference about
532 lexical content obtain only with particular priors. If we remove prior constraints entirely, the synthetic
533 listener failed to identify the correct intervals; it falsely inferred the presence of words that were not uttered
534 and ‘missed’ words that were spoken. It is worth mentioning that the absence of priors would be extremely
535 unlikely in realistic contexts, because our knowledge of language generates expectations about plausible
536 words in any given sentence (e.g., due to syntactic and semantic constraints, as well as simple effects of
537 word frequency) and contextual knowledge (e.g., knowing the topic of conversation, or being in a particular
538 setting) will also supply empirical priors. Indeed, the effect of priors on speech segmentation is well-
539 established in human speech perception. The common observation that word boundaries are difficult to
540 ascertain in an unknown language is an intuitive example that priors based on lexical knowledge help to
541 determine speech segmentation. In addition, the way that humans segment speech depends on previous

Active listening

542 words in a sentence (Cole, Jakimik et al. 1980, Mattys and Melhorn 2007, Mattys, Melhorn et al. 2007,
543 Kim, Stephens et al. 2012)—a simple demonstration that priors are flexibly applied in different contexts.
544 The aim of this simulation was to demonstrate the role of priors in speech recognition under active listening.

545 This simulation also shows that active listening goes beyond simply inferring the best explanation for a
546 particular sensory signal: active listening also infers which signals to ‘sample’. By this, we mean that
547 different segments (corresponding to plausible word boundaries) of the speech signal are evaluated, with
548 the goal of ‘sampling’ or selecting one set of intervals. The action (here, covert placement of word
549 boundaries, which can be considered more generally as active sampling) therefore goes hand-in-hand with
550 perception. This is demonstrated in the left panel of Figure 6: Although the words recognised provide the
551 best (Kim, Frisina et al.) explanation for acoustic sensations, both the words themselves and the placement
552 of word boundaries are categorically different from the right panel of Figure 6, in which the model was
553 equipped with different (uniform) prior beliefs. This ability to integrate different levels of beliefs and
554 inference is consistent with a hierarchical architecture, as suggested by (i) experimental studies that have
555 measured brain responses during speech perception (Davis and Johnsrude 2003, Vinckier, Dehaene et al.
556 2007, DeWitt and Rauschecker 2012), (ii) studies that examine the weights participants assign to different
557 cue types during speech segmentation; e.g., (Mattys, White et al. 2005), and (iii) cognitive accounts of
558 speech processing (McClelland and Elman 1986, Gaskell and Marslen-Wilson 1997). In the next section,
559 we turn to the electrophysiological correlates of this belief updating and ask what predictions this model of
560 auditory inference can offer.

561 **Face validity: Simulating sentence recognition**

562 Here, we use the generative model and inversion scheme described above, under simple prior beliefs about
563 a sentence, to illustrate the circular causality implicit in Bayesian belief updating. In brief, we will examine
564 how prior beliefs underwrite word segmentation and how segmentation changes in the absence of
565 appropriate priors. We then look at how the selected speech segmentation updates subsequent prior beliefs
566 and how the ensuing Bayesian surprise may manifest electrophysiologically. To illustrate the effect of
567 priors, we chose the following sentence: “Is there a square above?” This is a completely arbitrary sentence
568 but is interesting because the formant frequencies in the word “square” have a bimodal (biphone) structure
569 (Bashford, Warren et al. 2008), which means there is a fairly severe segmentation problem at hand. Will a

Active listening

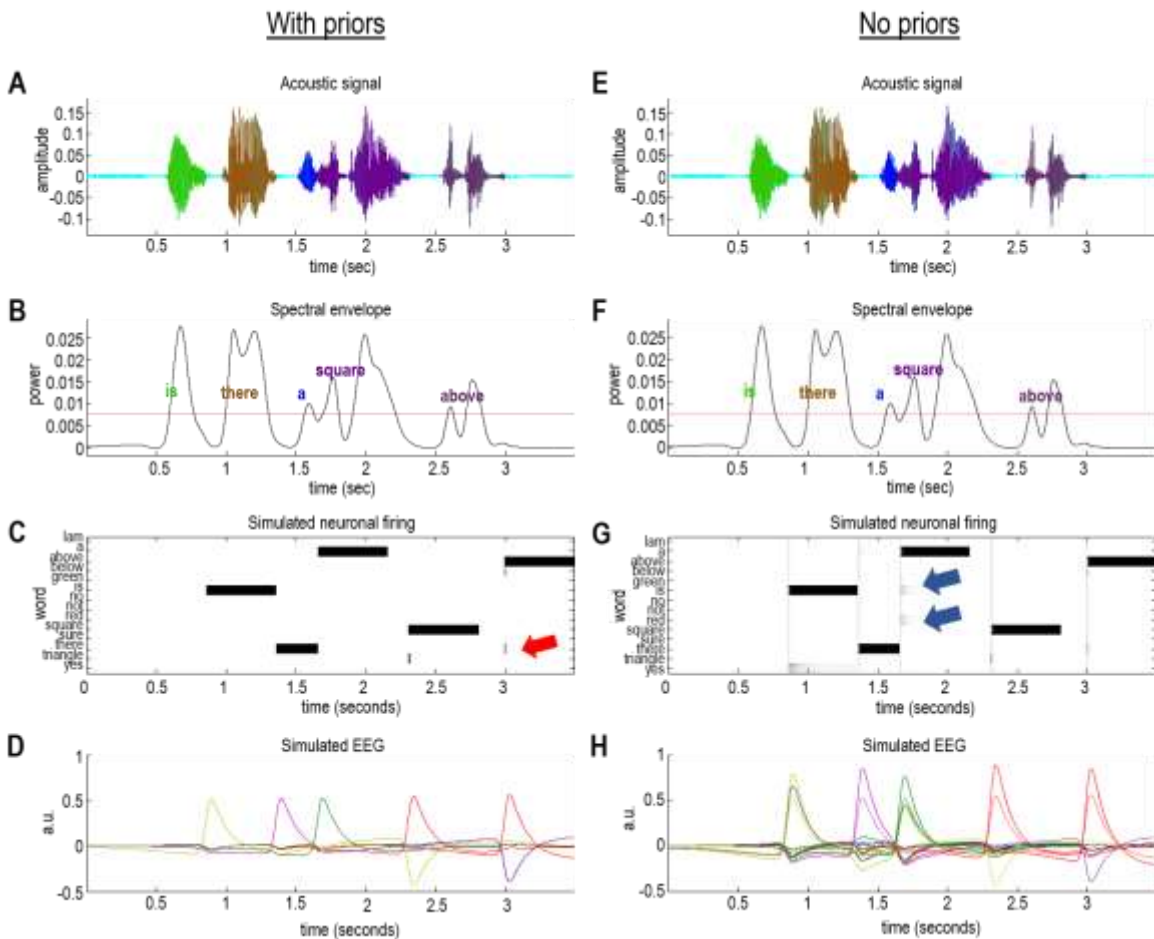
570 simulated subject segment “square” properly or—as in Figure 6—append the first phone to the previous
571 word? If they do infer the words correctly, how do priors manifest in terms of belief updating?

572 Figure 7 shows the results of integrating the active inference scheme above with strong (left panels) or
573 uniform (right panels) prior beliefs. In this example, prior beliefs were definitive for the first three words
574 (“is there a”) with more ambiguous prior for the last two words: for the fourth word, the possibilities
575 included “square” and “triangle”. For the final word, the possibilities included “above”, “below” and
576 “there”). These priors were selected because they are lexically congruent and represent a plausible belief
577 that a listener might have about the content of a sentence. Please see the figure legend for technical details.
578 The message from this simulation is that priors play a key role in resolving uncertainty and subsequent
579 competition among neuronal representation.

580 In the absence of precise prior constraints, the uncertainty associated with speech recognition is expressed
581 as an increased amplitude of simulated electrophysiological responses. This can be seen most clearly by
582 comparing the simulated electrophysiological responses in the lower right panel: the dotted lines reflect
583 belief updating in the absence of specific priors, while the dashed lines are the same responses under
584 informative priors. Figure 8 drills down on these differences by focusing on the responses to the third word.
585 In so doing, the simulated waveform looks very much like a P300 that is frequently observed in
586 electrophysiological studies (Donchin and Coles 1988, Morlet and Fischer 2014, Ylinen, Huuskonen et al.
587 2016). To understand this more formally, the next section explains how these simulated
588 electrophysiological responses were derived and how they can be interpreted in terms of belief updating
589 and Bayesian surprise.

590 To conclude this section, we will use this example to illustrate the fidelity of recursively generating and
591 recognising words, under this generative model. Figure 9 shows the segmentation and word recognition
592 following the presentation of the sentence above (“is there a square above”), without priors. The sentence
593 was then generated using the recognised lexical, prosodic and speaker attributes. The synthetic speech was
594 then presented to the active listening scheme, to recover the original utterance. This shows that the scheme
595 can understand itself and perform rudimentary speech repetition. More formally, it illustrates the validity
596 of the amortised inversion scheme.

Active listening



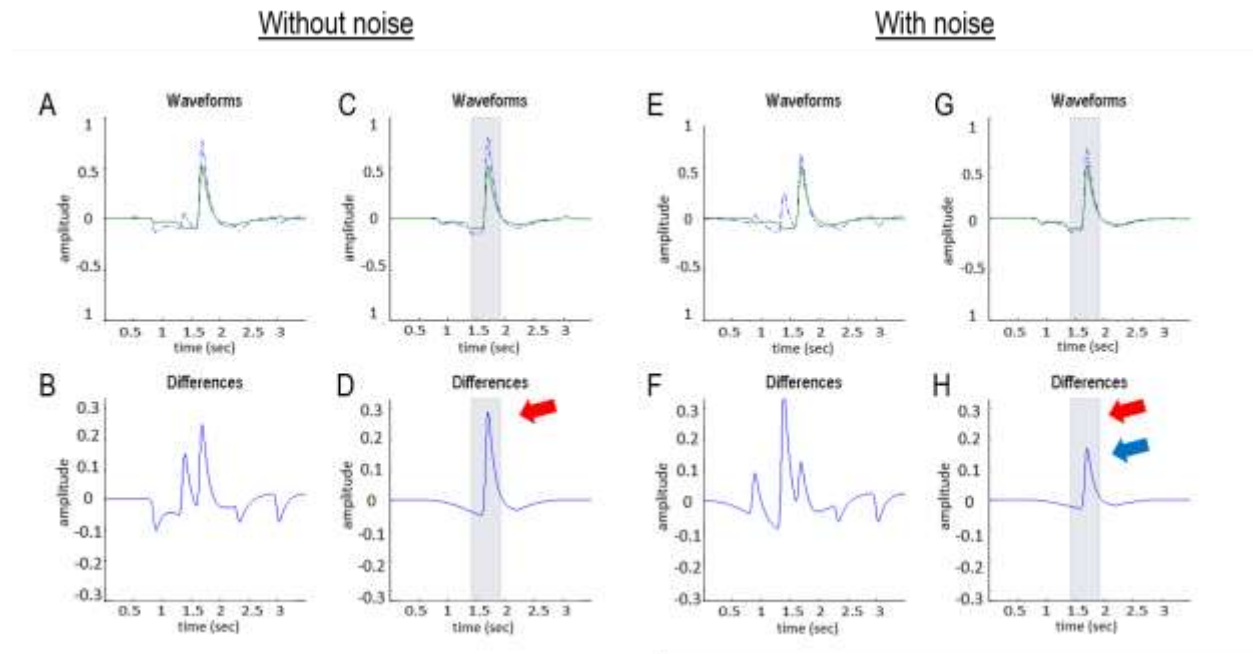
597

598

FIGURE 7

599 *The role of priors in a word recognition:* This figure uses the same format as Figure 6. In this example, the spoken
 600 sentence was “Is there a square above?” The left panel (A–D) shows the results of segmentation and word recognition
 601 under informative priors about the possible words. In other words, for each word in the sequence, a small number of
 602 plausible options were retained for inference. For example, the word “above” could have been “below” or “there”, as
 603 shown by the initial neuronal firing in Panel C at the end of the last word (red arrow). The right panel (E–H) shows
 604 exactly the same results but in the absence of any prior beliefs. The inference is unchanged; however, one can see in
 605 the neuronal firing (Panel G) that other candidates are competing to explain the acoustic signal (e.g., blue arrows).
 606 The key observation is that the resulting uncertainty—and competition among neuronal representations—is expressed
 607 in terms of an increased amplitude of simulated electrophysiological responses. This can be seen by comparing the
 608 simulated EEG trace in Panel H—in the absence of priors (solid lines)—with the equivalent EEG response under
 609 strong priors (solid lines in Panel D, reproduced as dashed lines in Panel H). In this example, there has been about a
 610 50% increase in the amplitude of evoked responses. A more detailed analysis of the differences in simulated EEG
 611 responses is provided in Figure 8.

Active listening



612

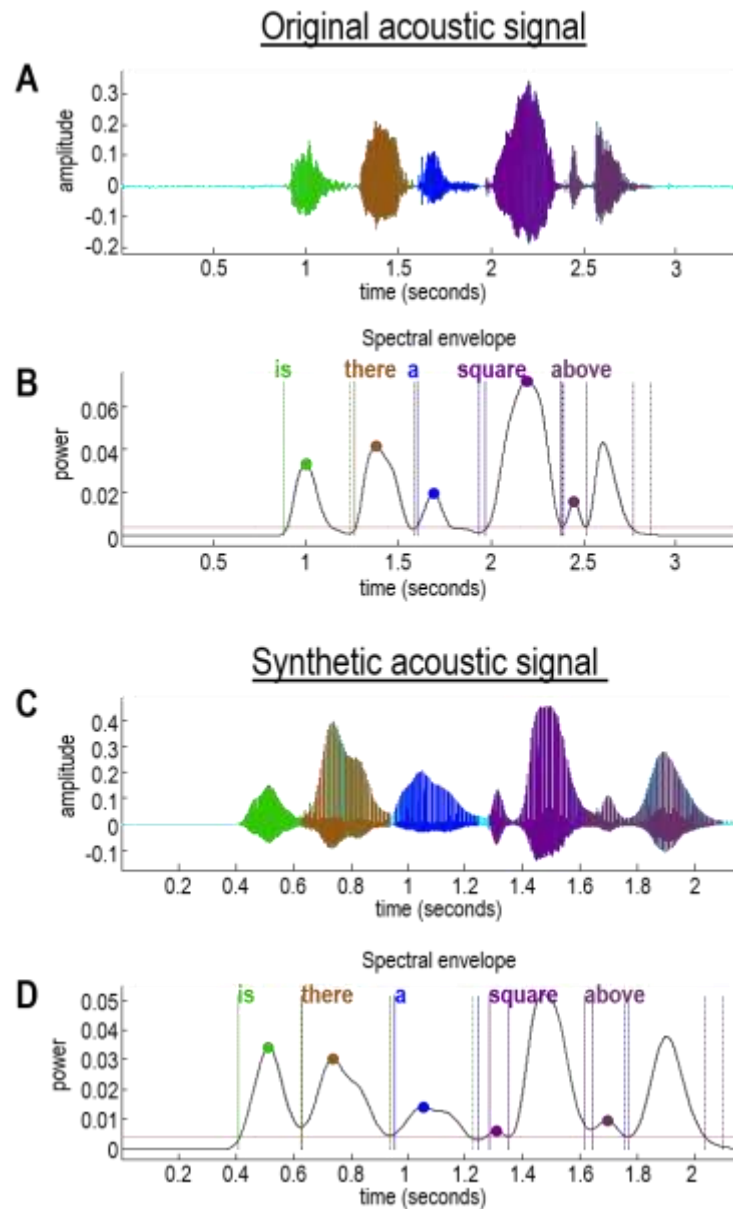
613

FIGURE 8

614 *Mismatch responses and speech-in-noise:* Panel A reproduces the results of Figure 7H, but focuses on the simulated
615 electrophysiological responses of a single neuronal population responding to the third word (“a”). The upper row
616 reports simulated responses evoked with (green lines) and without (blue dashed lines) priors (as in Figure 7), while
617 the lower row shows the differences between these two responses. These differences can be construed in the spirit of
618 a mismatch negativity or P300 waveform difference. Removing the priors over the third word (Panels C–D) isolates
619 the evoked responses and their differences more clearly. The grey shaded area corresponds to a peristimulus time of
620 500 ms, starting 250 ms before the offset of the word in question. Assuming update time bins of around 16 ms means
621 that we can associate this differential response with a P300. In other words, when the word is more surprising—in
622 relation to prior beliefs about what will be heard—they evoke a more exuberant response some 300 ms after its offset.
623 Panels E–H reports the same analysis with one simple manipulation; namely, the introduction of noise to simulate
624 speech-in-noise. In this example, we doubled the amount of noise; thereby shrinking the coefficients by about a factor
625 of half. This attenuates the violation (i.e., surprise) response by roughly a factor of two (compare difference waveform
626 in Panel D without noise—red arrows—with the difference waveform in Panel H without noise—blue arrow).
627 Interestingly, in this example, speech-in-noise accentuates the differences evoked in this simulated population when
628 the word is not selected (i.e., on the previous word). The underlying role of surprise and prior beliefs in determining
629 the amplitude of these responses is addressed in greater detail in the final figure.

630

Active listening



631

632

FIGURE 9

633 *Recursive recognition and generation:* The upper part of this figure shows the recognition of words (Panel B)
634 contained within an acoustic signal (Panel A). Here, the acoustic signal is parsed into the words “is there a square
635 above”. The corresponding lexical states can be used to synthesise a new acoustic signal (Panel C) containing the
636 same words. Here, we inverted the model a second time, to recover the words contained within the synthetic acoustic
637 signal (Panel D). Happily, the recovered words from the synthetic signal (Panel D) match those from the original
638 signal (Panel B).

Active listening

639 **Predictive validity: Belief updating and neurophysiology**

640 Figure 8 suggests that belief updating during word recognition depends sensitively on prior beliefs and
641 implicit differences in the confidence with which a particular word is inferred. Here, we pursue the
642 predictive validity of this active listening formulation, by looking in greater detail at belief updating under
643 the model. In doing so, we highlight qualitative similarities to canonical violation responses measured with
644 EEG and MEG that are well-established in the empirical literature (as discussed in more detail below). In
645 brief, the message of this section is that evoked or induced responses in the brain will increase in proportion
646 to the degree of belief updating following sensory input.

647 Generally speaking, the idea that belief updating may underpin vigorous neuronal responses to surprising
648 sensations is broadly consistent with experimental observations. Under predictive coding models of
649 auditory perception, the mismatch negativity has been considered in light of precision weighted prediction
650 error responses (Garrido, Kilner et al. 2009, Wacongne, Changeux et al. 2012, Heilbron and Chait 2018).
651 In this literature, the mismatch negativity is related to deviants in elementary acoustic events, such as
652 frequency (Näätänen, Gaillard et al. 1978, Giard, Lavikahen et al. 1995, Jacobsen, Schröger et al. 2003),
653 intensity (Näätänen, Gaillard et al. 1978, Giard, Lavikahen et al. 1995, Jacobsen, Horenkamp et al. 2003),
654 or timbre (Tervaniemi, Ilvonen et al. 1997, Tervaniemi, Winkler et al. 1997, Toiviainen, Tervaniemi et al.
655 1998)—and its amplitude covaries with the probability of a deviant (Picton, Alain et al. 2000, Sato, Yabe
656 et al. 2000, Sato, Yabe et al. 2003). Mismatch negativity responses have also been recorded in the context
657 of spoken phonemes (Dehaene-Lambertz 1997, Näätänen, Lehtokoski et al. 1997). In the current
658 framework, precision weighted prediction errors induced by acoustic deviations reflect the surprise and
659 concomitant belief updating induced by heard (spoken) words. At a slightly longer latency, reorientation
660 responses could also be construed as a reflection of belief updating at higher levels of hierarchical inference.
661 For example, the P300 has been proposed to reflect contextual violations (Donchin and Coles 1988) and
662 the N400 has been proposed to reflect semantic violations (Kutas and Hillyard 1980, Kutas and Hillyard
663 1984, Van Petten, Coulson et al. 1999, Kutas and Federmeier 2000). The whole field of repetition
664 suppression and adaptation in functional magnetic resonance imaging rests upon exactly the same notion;
665 namely, an attenuation of neuronal responses that induce less belief updating, in virtue of being predictable
666 or repetitious (Larsson and Smith 2012, Grotheer and Kovács 2014).

667 In the current simulations, our agenda is to identify generic principles that may underpin neuronal responses

Active listening

668 to surprising sensations under active listening. Our goal was not to simulate any particular type of ERP
669 component, but merely to observe belief updating in the current framework. In the discussion section, we
670 visit the finer details of the mismatch negativity and later endogenous (e.g., P300, N400) responses, which
671 would be interesting avenues for future work. An advantage of the current setup is that we can expand upon
672 the qualitative explanation for violation or surprise related responses using explicit, quantitative
673 simulations.

674 If we take the average change in depolarisation under expected firing rates (after belief updating), we
675 recover a quantity that scores the degree of belief updating (see Appendix 4 for details)—a quantity that
676 emerges in many guises in different disciplines. For example, in statistics, it is known as the *complexity*
677 (see equation A.18), which scores the departure from prior beliefs required to provide an accurate account
678 of some data (Penny 2012). In the visual neurosciences, this quantity is known as *Bayesian surprise*
679 (Schmidhuber 1991, Itti and Baldi 2009) that underwrites the *salience* or epistemic affordance of locations
680 in the visual scene that attract saccadic eye movements (Parr and Friston 2017). In robotics, this quantity is
681 known as *intrinsic motivation*; namely the *information gain* associated with a particular move or action
682 (Ryan and Deci 1985, Oudeyer and Kaplan 2007). In short, we have a link between the information theoretic
683 quantity that reflects the degree of Bayesian belief updating and the average neuronal responses that
684 perform belief updating.

685 There are a number of reasons that one might consider this a sensible predictor of evoked responses in the
686 brain, above and beyond the idealised dynamics described above. These reasons rest upon the statistical
687 physics of belief updating in any sentient system making inferences about external states of affairs. The
688 technical back story to active inference—that is, the free energy principle—allows one to associate the
689 degree of belief updating and implicit changes in variational free energy in terms of a thermodynamic
690 potential (Landauer 1961, Bennett 2003, Friston 2013). This means that for an ensemble of neurons (or
691 neuronal processes) belief updating can be translated directly into thermodynamic free energy. The
692 corresponding thermodynamic cost of belief updating may be reflected in nearly every sort of
693 electrophysiological neuroimaging measurement. For example, the excursions of transmembrane potentials
694 from their Nernst equilibrium in EEG (c.f., a mismatch negativity amplitude). Similarly, in fMRI,
695 activations may reflect the metabolic costs of belief updating (Attwell and Iadecola 2002).

696 The second line of argument is based upon the common sense observation that, in the absence of an

Active listening

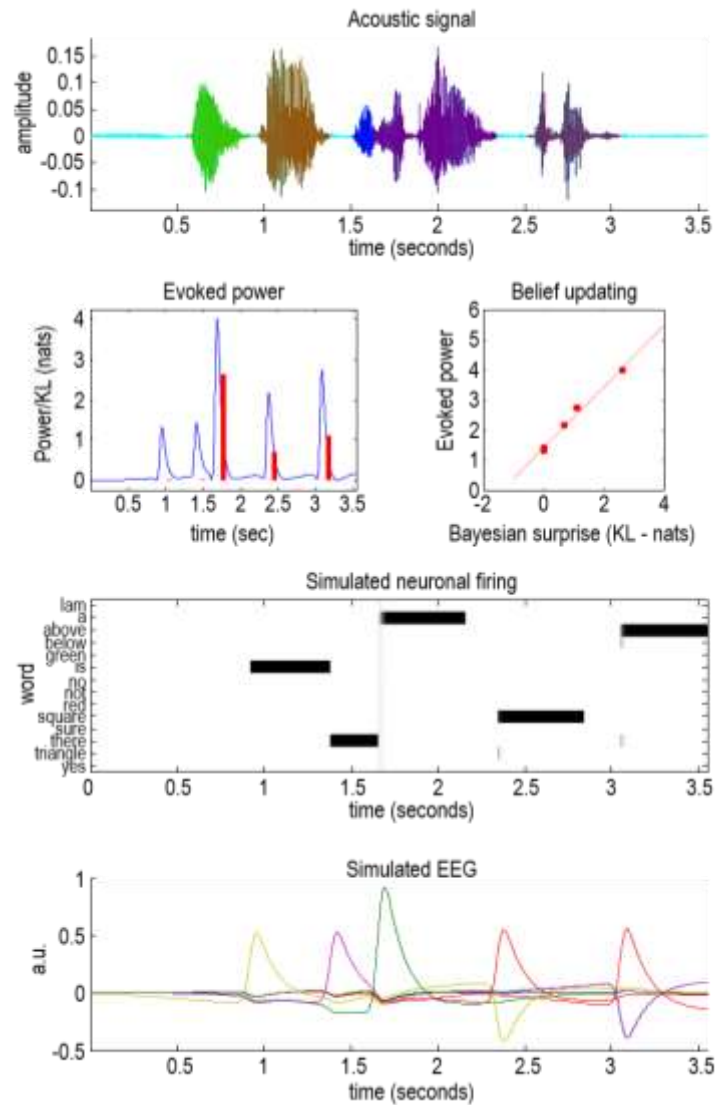
697 informative sensory cue, there can be no belief updating and no complexity cost or accompanying
698 thermodynamic cost (Sengupta, Tozzi et al. 2016). In this instance, there will be, clearly, no evoked or
699 induced response. This argument further suggests that the precision of continuous sensory (e.g., auditory)
700 signals will determine the degree of belief updating and related violation responses, such as the mismatch
701 negativity. In speech perception, reduced precision could correspond to speech-in-noise, for which this
702 model predicts an attenuation of mismatch responses as noise levels increase. The basis of this effect rests
703 upon the estimation of random fluctuations in sensory cues that, under predictive coding, shrink the
704 posterior expectations of the lexical coefficients towards their prior mean.

705 If we revisit the results in Figure 6 and Figure 7, and compare responses evoked with and without priors, it
706 is immediately obvious that, on average, evoked responses in the absence of (accurate) priors have a larger
707 amplitude. This is sensible because priors that are congruent with the words presented mean that the belief
708 updating has a smaller complexity cost because the prior is closer to the posterior. In other words, there is
709 less information gain because the (synthetic) subject already had accurate prior beliefs about the lexical
710 content of the spoken words.

711 To illustrate the sort of effect more quantitatively, we repeated the simulations reported in Figure 7 but
712 introduced uncertainty about the third word by relaxing its priors. This allowed us to introduce differences
713 in belief updating, from word to word, and show that simulated neuronal responses vary monotonically
714 with information gain or Bayesian surprise. Figure 10 reports the results of this numerical analysis in terms
715 of the variance of depolarisation over neurons encoding lexical expectations (blue line in the second panel)
716 and the corresponding Kullback-Leibler divergence (red bars). Their monotonic relationship is apparent
717 (see the third panel), although the relationship is not perfect due to filtering the simulated EEG data and our
718 *ad hoc* measure of neuronal responses. At the (coarse-grained) level of the current treatment, this can be
719 regarded as a simulation of neuronal responses to Bayesian surprise at a fairly high level in the auditory
720 hierarchy (encoding the lexical content of a word).

721

Active listening



722

723

FIGURE 10

724 *Bayesian surprise and evoked responses*: this shows the same results as in Figure 7 but after removing priors from the
725 third word (“a” in blue). The result is a more vigorous simulated event related response after the onset of the third
726 word (green line in the bottom panel). A simple measure of these surprise-related responses can be obtained by taking
727 the variance of the (simulated) responses over all populations as a function of time (c.f., evoked power). This is shown
728 in the second panel as a solid blue line (normalised to a maximum of four arbitrary units). The red bars correspond to
729 the degree of belief updating or Bayesian surprise, as measured by the KL divergence between prior and posterior
730 beliefs after updating. The key conclusion from these numerical analyses is that there is a monotonic relationship
731 between the evoked power and Bayesian surprise, as shown by the nearly linear relationship between Bayesian surprise
732 and the maxima of evoked power in the third panel. In short, the greater the Bayesian surprise, the greater the belief
733 updating and the larger the fluctuations in neuronal activity.

Active listening

734 With this characterisation of mismatch responses, we can now return to the effect of noise, which highlights
735 a key feature of active listening—that the quality of sensory evidence affects the magnitude of belief
736 updating. In Figure 8, noise was simulated by decreasing the prior precision associated with the lexical
737 coefficients at the auditory level of inference (namely, the prior precision in Equation A.20). This
738 manipulation attenuates the mismatch or surprise response because the degree of belief updating has been
739 reduced. The attenuation arises because there is less confidence placed in the evidence ascending from
740 lower (sensory) levels of auditory processing. In other words, the attenuation of belief updating (and
741 mismatch responses) in Figure 8 arises because the posteriors have been moved closer to the priors. This
742 contrasts Figure 7, in which belief updating and mismatch responses were attenuated by one moving the
743 priors closer to the posteriors. In subsequent work, we will revisit the effects of manipulating speech-in-
744 noise—and prior beliefs—to demonstrate their effects empirically and, crucially, how they interact in the
745 genesis of difference waveforms. For the purposes of this paper, the basic phenomenology illustrated above
746 will be taken as a validation of the belief updating scheme by appealing to the literature on the canonical
747 mismatch and violation responses of this sort.

748 Discussion

749 Active listening considers the enactive synthesis or inference that might underwrite the recognition—and
750 generation—of spoken sentences. The notion of *active listening* inherits from active inference, which
751 considers perception and action under a universal imperative—to maximise the evidence for our
752 (generative) models of the world. Here, the ‘active’ component is the (covert) parsing of words from a
753 continuous auditory signal. Active listening entails the selection of internal actions (i.e., placement of word
754 boundaries) that minimise variational free energy. Practically, word boundaries are selected so as to
755 minimise surprise or maximise the evidence for an internal model of word generation. We have described
756 the formal basis of this kind of active listening, using simulations of speech recognition to establish its face
757 validity in behavioural terms. We then considered predictive validity, in terms of neuronal or physiological
758 responses to violations and surprise, of the sort associated with the mismatch negativity, P300, and N400.

759 In treating the segmentation of a continuous sensory stream into meaningful words as an active sensing
760 problem, we imagine that several segmentation operations are applied by the auditory system in parallel
761 and the interval that maximises model evidence or marginal likelihood (i.e., minimises variational free

Active listening

762 energy) is selected for further hierarchical processing. From the perspective of hierarchical Bayesian
763 inference, this follows the usual way of mapping from posterior density estimates, based upon continuous
764 signals, to posterior beliefs about the discrete causes of those signals. This is generally cast in terms of
765 Bayesian model selection. In other words, selecting some discrete explanation or hypothesis for the data
766 that is most consistent with the estimated parameters of a generative model at the lower (sensory) level
767 (Friston, Parr et al. 2017). The twist here is that this model selection has been framed in terms of action
768 selection by treating the selection of word boundaries as an active process.

769 The generative model of word production that we considered has been stripped down to its bare essentials.
770 More complex models could be conceived that synthesise more natural speech. Expanding the parameter
771 space would not only allow it to produce more natural speech, but also allow the model to explain more
772 domains of auditory production and perception. We discuss some of these possibilities in the discussion
773 that follows. Nevertheless, we have demonstrated with this simplified generative model that inversion of
774 the model—which corresponds to speech recognition—is associated with belief updating that makes
775 plausible predictions for neuronal dynamics. In this paper, we produced quantitative simulations of
776 electrophysiological responses and showed that they depend on the prior knowledge of the listener—a
777 phenomenon that has commonly been observed in human speech perception (Marslen-Wilson 1975,
778 Marslen-Wilson and Welsh 1978, Cole, Jakimik et al. 1980, Mattys and Melhorn 2007, Mattys, Melhorn et
779 al. 2007, Kim, Stephens et al. 2012).

780 In borrowing ideas from active vision, we highlight parallels by which the brain could plausibly accumulate
781 evidence among sensory modalities. The covert actions considered in this paper (i.e., the placement of word
782 boundaries) follow in the spirit of overt (motor or autonomic) actions that have been used to simulate
783 saccadic searches of the visual scene (Mirza, Adams et al. 2016, Parr and Friston 2017). We discuss the
784 relationship between covert and overt actions in greater depth below. Intuitively, sensory observations in
785 the auditory and visual modalities may appear to differ because speech unfolds over time, whereas visual
786 experiments frequently use static stimuli that are spatially distributed. However, many parallels can be
787 drawn between cortical processing in these modalities (O'Leary 1989), consistent with findings that sensory
788 cortices can reorganise and subsequently process inputs from a different sensory modality (Sur, Garraghty
789 et al. 1988, Shnell, Champoux et al. 2015). Shamma and colleagues (Shamma 2001, Shamma, Elhilali et al.
790 2011) propose a unified computational framework for auditory and visual perception, suggesting that the
791 neural processes proposed for vision could also operate in auditory cortex. In short, this is based on the idea

Active listening

792 that the cochlea transforms temporally unfolding sound into spatiotemporal response patterns early in
793 auditory processing. In other words, this is a ‘spatial’ view of auditory processing. Under this view, the
794 computations for analysing auditory signals in time could be similar to the computations used for analysing
795 visual signals in space; e.g., (Bar, Kassam et al. 2006).

796 **Active listening and Bayesian surprise**

797 Selecting intervals containing auditory cues that minimise free energy (i.e., maximise marginal likelihood
798 or model evidence) follows from the basic premise of the free energy principle; namely, both action and
799 perception are in the game of self-evidencing (Hohwy 2016). Having said this, there is something unique
800 about the particular selective process (which are implicit in Equation A.19) that distinguishes it from overt
801 actions, such as moving one’s head or making visual saccades to a location in a visual scene. This is because
802 the corresponding selection of ‘where to look next’ is based upon anticipated data that would be sampled
803 if one looked ‘over there’. However, predictive coding (in some amortised form) of speech segmentation
804 here is based on evidence *that has already accumulated* under different interval or segmentation schemes.
805 In other words, there is a distinction between overt actions—such as moving one’s eyes or moving one’s
806 head—which changes observations in the future, and covert actions—such as covert visual attention, or
807 selecting a particular segmentation of speech—which is based on sampling current observations. In the case
808 of these covert actions, the sensory evidence (and subsequent posterior) can be computed explicitly to
809 evaluate the free energy expected under a particular interval choice. In contrast, expected free energy based
810 on overt actions has to be averaged under predicted sensory outcomes—known technically as a posterior
811 predictive density. This means that evaluating the *free energy* for particular speech segmentation intervals
812 is much simpler than evaluating the *expected free energy* under a posterior predictive density, conditioned
813 upon a particular overt action. It is useful to bear this distinction in mind because it can resolve some
814 apparent paradoxes.

815 These paradoxes pertain largely to the question: does active inference minimise or maximise Bayesian
816 surprise? In the current setting, covert actions associated with speech segmentation minimise Bayesian
817 surprise, because Bayesian surprise relates to the complexity (i.e., cost) associated with belief updating
818 based on current observations. In other words, because the free energy associated with covert actions can
819 be evaluated explicitly, a listener can choose the covert action that requires the least belief updating (i.e.,

Active listening

820 that is closest to their priors), but still provides an accurate explanation for the auditory observations. This
821 leads to a conceptualisation in which neuronal dynamics and implicit message passing aim to explain
822 sensory input with minimal complexity and, therefore, minimum accompanying thermodynamic cost
823 (Sengupta, Stemmler et al. 2013). On this view, large mismatch or violation responses indicate that an
824 accurate explanation for sensory inputs required a costly update to posterior beliefs.

825 The situation flips for overt actions, for which action selection depends on *expected* free energy—which is
826 evaluated on the basis of predicted (i.e., unknown) outcomes in the future. Future sensory outcomes are
827 random (i.e., unknown or hidden) variables and active inference maximises expected Bayesian surprise,
828 which corresponds to expected information gain. In other words, it reflects the reduction in uncertainty in
829 how the world is sampled. Actions that maximise Bayesian surprise will lead to the greatest reduction in
830 uncertainty. This is why *expected* Bayesian surprise has to be maximised when selecting actions, where it
831 plays the role of epistemic affordance (Parr and Friston 2017). As noted above, this is an important
832 imperative that underwrites uncertainty reducing, exploratory behaviour; known as intrinsic motivation in
833 neurorobotics (Schmidhuber 2006) or salience when ‘planning to be surprised’ (Sun, Gomez et al. 2011,
834 Barto, Mirolli et al. 2013). An intuitive way of thinking about whether surprise should be maximised or
835 minimised is to appeal to the analogy of scientific experiment. We may attempt to analyse empirical data
836 that we have collected in a way that minimises how surprising it appears; for example, by giving greater
837 weight to hypotheses consistent with our measurements. Having done so, we may want to design a future
838 experiment, which would aim is to collect data that will tell us something new; in this case, we should
839 design an experiment that we expect to maximise our (Bayesian) surprise (a.k.a., information gain).

840 In future work, we will expand upon this distinction by using the current model to simulate conversations.
841 The act of speaking is an overt action, and the basic principle of conversational turn taking has been
842 simulated using active inference in the setting of bird song (Friston and Frith 2015). We hope to combine
843 the current active listening implementation with an agent who is able to ask questions. In brief, the agent
844 will actively listen to speech by *minimising* Bayesian surprise at the level of word recognition considered
845 in this paper, and select words to speak (i.e., overt actions, here in the form of questions) that *maximise*
846 expected Bayesian surprise to maximise information gain (i.e., resolve uncertainty). This leads to a first
847 principle account of language ‘understanding’ that can be described in terms of self-evidencing: namely,
848 minimising free energy through belief updating, and planning to take actions that minimise expected free
849 energy.

Active listening

850 Although evaluating the free energy of alternative data features (i.e., segments) that have already been
851 sampled is more straightforward than evaluating the expected free energy when planning how to sample
852 data, it is not as straightforward as reflexive action; e.g., (Adams, Shipp et al. 2013). Reflexive or
853 elementary action, under active inference, changes the sensory data solicited, e.g., the stretch receptor
854 signals that are attenuated by classical motor reflexes. However, this kind of reflexive action does not
855 change internal brain states or the posterior beliefs that they parameterise. This means that the only part of
856 free energy that can be minimised directly is the accuracy term (Equation A.18). This is why it is sufficient
857 to minimise interoceptive and proprioceptive prediction errors when accounting for autonomic and motor
858 action; very much along the lines of the equilibrium point hypothesis (Feldman and Levin 1995) and the
859 passive movement paradigm (Mohan and Morasso 2011). However, in the active listening framework
860 proposed here, the situation is a little more involved. This is because hierarchical inference means that
861 committing to one data feature (i.e., interval) or another will change posterior beliefs. This means that to
862 comply with the free energy principle, it is necessary to select data features (i.e., intervals) that not only
863 maximise accuracy but also minimise complexity. This entails a more nuanced form of action selection, in
864 virtue of the fact that it requires the (covert) selection of data features that have been (overtly) acquired.
865 Even though the data have already been acquired, and selecting different data features does not change the
866 auditory outcomes (acoustic timeseries), these processes are nevertheless ‘active’ from our perspective,
867 because the agent has an epistemic imperative to sample auditory outcomes in a way that reduces
868 uncertainty. In other words, the agent is in charge of the *data features* (i.e., segmentation). Thus, we can
869 think of speech segmentation as a kind of action that is internal or attentional, related to how the acoustic
870 timeseries is covertly sampled. The framework we have introduced in this paper highlights that—
871 mathematically—these covert actions can be considered in a similar way as overt actions.

872 **Acoustic envelope and spectral fluctuations**

873 Under active listening, the implicit generative model of an envelope, which is used to create a repertoire of
874 intervals from which to select, is distinct from the spectral fluctuations (i.e., formant frequencies) generated
875 by latent states (i.e., lexical and prosody). This formulation of speech recognition may explain why there
876 are ‘envelope following responses’ in distinct parts of the auditory system, whose functional architecture
877 can be distinguished from the tonotopic mapping of auditory cortex per se (Easwar, Purcell et al. 2015,
878 Braiman, Fridman et al. 2018). This leads to an interesting picture of how the brain thinks words are
879 generated that echoes the distinction between ‘what’ and ‘where’ in the visual hierarchy (Ungerleider and

Active listening

880 Haxby 1994). In other words, there may be a homologous distinction between ‘what’ and ‘when’ in the
881 auditory system that manifests as an anatomical separation of the pathways inferring ‘what’ is being spoken
882 (i.e., tonotopic predictions and representations) and when this content is deployed (i.e., envelope following
883 responses) (Romanski, Tian et al. 1999, Alain, Arnott et al. 2001). From the point of view of word
884 generation, these two streams converge to generate the correct formants at the correct time. From the point
885 of view of recognition or generative model inversion; this would imply a functional segregation of the sort
886 seen in other modalities (Ungerleider and Haxby 1994, Friston and Buzsaki 2016); for example, the
887 segregation into dorsal and ventral streams – or, indeed, parvocellular and magnocellular streams (Zeki and
888 Shipp 1988, Nealey and Maunsell 1994). Interestingly, this sort of segregation into ‘what’ and ‘how’
889 pathways has already been proposed for the auditory system (Kaas and Hackett 1999, Belin and Zatorre
890 2000).

891 **Active listening and electrophysiological responses**

892 In a general sense, we have shown that belief updating under active listening qualitatively resembles
893 physiological responses to violations and surprise that are already in the literature. Our goal was not to
894 simulate any particular type of ERP component or the empirical results from any particular study, but rather
895 to explore belief updating in an artificial agent whose goal is to generate and/or recognise speech. So, can
896 we interpret this belief updating in light of particular ERP responses?

897 One canonical violation response is the mismatch negativity. The mismatch negativity is observed in classic
898 ‘oddball’ paradigms (Garrido, Kilner et al. 2009), in which a deviant sound follows a sequence of sounds
899 that all share a particular acoustic property. Mismatch negativity responses have been observed when a
900 sound deviates in frequency (Näätänen, Gaillard et al. 1978, Giard, Lavikahen et al. 1995, Jacobsen,
901 Schröger et al. 2003), intensity (Näätänen, Gaillard et al. 1978, Giard, Lavikahen et al. 1995, Jacobsen,
902 Horenkamp et al. 2003), or timbre (Tervaniemi, Ilvonen et al. 1997, Tervaniemi, Winkler et al. 1997,
903 Toiviainen, Tervaniemi et al. 1998) from preceding stimuli. Crucially, the mismatch negativity has recently
904 been interpreted in terms of predictive coding—specifically, it has been assumed to reflect precision
905 weighted prediction errors (Garrido, Kilner et al. 2009, Wacongne, Changeux et al. 2012, Heilbron and
906 Chait 2018)—which relates nicely to the current framework. The finding that the amplitude of the mismatch
907 negativity covaries with the probability of a deviant (Picton, Alain et al. 2000, Sato, Yabe et al. 2000, Sato,

Active listening

908 Yabe et al. 2003) is consistent with the idea that it reflects belief updating. Most previous studies of the
909 mismatch negativity have used basic auditory stimuli, such as artificial pure or complex tones; it is therefore
910 assumed to reflect deviations to low-level acoustic properties, rather than processes that are specific to
911 speech. Nevertheless, observations of the mismatch negativity during phoneme perception (Dehaene-
912 Lambertz 1997, Näätänen, Lehtokoski et al. 1997) can be interpreted as reflecting acoustic violations that
913 occur within speech.

914 The P300 is often observed in similar ‘oddball’ settings as the mismatch negativity (Polich 2007). It has a
915 longer latency than the mismatch negativity and has been related to higher-level context violations
916 (Donchin and Coles 1988). It could, therefore, be interpreted as reflecting belief updating when the
917 listener’s context changes. In the domain of speech, the P300 has been associated with word frequency
918 (Polich and Donchin 1988).

919 The N400 is commonly observed in response to meaningful speech, and has also been associated with word
920 frequency (Kutas and Hillyard 1984, Van Petten and Kutas 1990, Van Petten, Coulson et al. 1999). Kutas
921 and Hillyard (Kutas and Hillyard 1984) found that the amplitude of the N400 was inversely correlated with
922 a word’s cloze probability—that is, participants’ ratings of the probability that a particular word would
923 come at the end of the sentence in question. They found that the same effect transferred to words that were
924 semantically related to high-probability words. They, therefore, concluded that the N400 relates to semantic
925 activation. Modulations of N400 responses have been reported in a variety of semantic contexts (reviewed
926 by (Kutas and Federmeier 2000))—including sentence-final words, the semantic congruency of words that
927 occur mid-sentence, and the semantic relatedness of word pairs—and has been shown to build up as the
928 semantic context becomes increasingly constrained throughout a sentence. Syntactic violations do not elicit
929 an N400 response (Kutas and Federmeier 2009), but instead evoke a P600 (Osterhout and Holcomb 1992,
930 Friederici, Hahne et al. 1996, Kuperberg, Sitnikova et al. 2003).

931 An N400-like negativity, termed the frontocentral negativity (‘FN400’) has been related to speech
932 segmentation by transitional probabilities (Balaguer, Toro et al. 2007, Cunillera, Càmara et al. 2009,
933 François, Cunillera et al. 2017). For example, stronger FN400 responses were elicited from acoustic signals
934 that comprised strong statistical relationships between syllables than syllables that were selected randomly
935 (François, Cunillera et al. 2017). The FN400 also appears to increase in amplitude as the segmentation
936 process becomes more prominent as new words are learned (Balaguer, Toro et al. 2007, Cunillera, Càmara

Active listening

937 et al. 2009).

938 Speech segmentation by prosodic cues has been associated with a different ERP: the closure positive shift
939 (CPS) (Steinhauer, Alter et al. 1999). The closure positive shift is evoked around the time of a prosodic
940 boundary, and has been reported to last until the onset of the next word (Bögels, Schriefers et al. 2011). It
941 has been found in several different languages (see (Bögels, Schriefers et al. 2011) for a review) and even
942 in hummed speech (Pannekamp, Toepel et al. 2005), which has no lexical content.

943 So, which level of processing does belief updating in the current scheme reflect? This level could be
944 intermediate between lower acoustic levels at which a mismatch negativity is generated, and the kind of
945 violation responses associated with a change in context or semantics. Possibly, this could be something like
946 the phonological mismatch negativity, which has been interpreted as reflecting acoustic-phonetic
947 processing in response to the initial phoneme of a spoken word, occurring 270–300 ms after onset
948 (Connolly, Phillips et al. 1992). Connolly and Phillips (Connolly and Phillips 1994) observed the
949 phonological mismatch negativity when the final word of a sentence was semantically congruent, but the
950 word (and the initial phoneme) differed from the word with the highest Cloze probability. An N400 was
951 not observed in this condition and was instead observed when the word was semantically incongruent.
952 Interestingly, the phonological mismatch negativity was not observed when a word was semantically
953 incongruent, but the initial phoneme matched the word with the highest Cloze probability. These
954 observations are consistent with the idea that the phonological mismatch negativity reflects acoustic-
955 phonetic processing.

956 One advantage of the current framework is that it generates quantitative predictions that can be explicitly
957 tested in future electrophysiological studies. The predictive validity we have considered here is a first step:
958 the next step is to scrutinise the particular parameters of the simulation using empirical data. To study this
959 in more detail, specific sequences of words and/or acoustic features could be posed to the model that
960 generate particular violations. Belief updating in active listening—and, for comparison, parameters of other
961 models (Aitchison and Lengyel 2017)—could be quantitatively compared to empirical electrophysiological
962 results. This speaks again to future directions, in which the current framework will be extended to a
963 hierarchical model that can simulate conversations. Speech has a deep temporal structure, with phrases
964 evolving over longer time intervals than words or phonemes—and a more complete generative model of
965 speech will have to incorporate this temporal hierarchy (Friston, Rosch et al. 2017). The idea of an

Active listening

966 interlocutor asking questions to resolve uncertainty relates to a higher-level semantic processing of
967 speech—and violations of semantic expectations might be associated with later electrophysiological
968 responses, such as the N400. Consistent with the types of hierarchies that have often been suggested based
969 on empirical data (Kumar, Stephan et al. 2007, Ding, Melloni et al. 2015), a deep generative model implies
970 that belief updating occurs at multiple time scales, and we anticipate that this will give rise to more
971 structured ERPs that include contributions from later components.

972 **Background noise during active listening**

973 In this paper, we simulated a simple case of speech-in-noise, in which we imposed random fluctuations (of
974 constant amplitude) on the speech signal. We showed that noisier signals attenuate belief updating. We plan
975 to extend this model to incorporate other types of noise, including fluctuating-amplitude maskers such as
976 multi-speaker environments. This should allow one to investigate which aspects of the signal are most
977 informative for minimising Bayesian surprise, when some parts of the signal (but not others) undergo
978 energetic masking (Brungart 2001, Brungart, Simpson et al. 2001, Durlach 2006) or when informational
979 masking (Durlach, Mason et al. 2003, Durlach, Mason et al. 2003, Kidd, R. Mason et al. 2007) comes into
980 play. In other words, in the presence of noise, a listener needs to reduce their uncertainty about the words
981 that were spoken by deciding which attributes of the acoustic signal they should attend to.

982 One problem that the current segmentation algorithm would face—when adding background noise to
983 speech—is that envelope minima may not always be present at word boundaries. In human listeners,
984 segmentation at envelope minima could be achieved based on envelope following responses. Indeed, the
985 magnitude of envelope following responses (i) has been linked to speech intelligibility in humans (Drullman
986 1995, Muralimanohar, Kates et al. 2017, Vanthornhout, Decruy et al. 2018), (ii) is greater for attended than
987 unattended speakers (Ding and Simon 2012, O'Sullivan, Power et al. 2014), and (iii) can be reconstructed
988 from measurements of brain activity (Pasley, David et al. 2012, O'Sullivan, Power et al. 2014). These
989 envelope responses could, therefore, reflect the success of speech segmentation. Other cues to segmentation
990 have been reported in the literature—and may be particularly important when background noise is present.
991 These cues include durations: a lengthening of syllables at the end of words (Klatt 1975, Beckman and
992 Edwards 1990), and possibly also the beginning (Lehiste 1960, Lehiste 1972, Oller 1973, Klatt 1976,
993 Nakatani and Dukes 1977, Gow Jr and Gordon 1995). They also include a shortening of the middle portion

Active listening

994 of words (Lehiste 1973, Oller 1973, Harris and Umeda 1974, Klatt 1976). Other work has also reported
995 metrical (stress) cues (Cutler and Norris 1988), allophonic variation (Christie Jr 1974, Nakatani and Dukes
996 1977, Gow Jr and Gordon 1995), and fundamental frequency contour (Ladd and Schepman 2003) as
997 segmentation cues. Although the current algorithm of finding envelope minima was sufficient for the
998 current simulations, these other cues could be implemented into active listening in other contexts in which
999 segmentation may be particularly challenging. While the current implementation retrospectively places
1000 word boundaries, future work could also consider that word boundaries are somewhat predictable from the
1001 lexical statistics of the preceding sequences (Marslen-Wilson 1984)—for example, the offset of “trombone”
1002 may be predicted upon hearing “trom”, given it is the only valid ending to the word in English.

1003 **Active listening and language production and perception**

1004 The active listening scheme can also be used as a foundation to gain a neuronal-level understanding of
1005 language production and perception behaviours. For example, engaging in a two-way dialogue (Kuhlen,
1006 Bogler et al. 2017), verbal fluency (Paulesu, Goldacre et al. 1997) and reading (Fiez and Petersen 1998,
1007 Landi, Frost et al. 2013, Taylor, Rastle et al. 2013); see (Price 2012) for a detailed overview. Previous
1008 investigations of these behaviours have been motivated by the desire to better understand the underlying
1009 neuropsychology (Aring 1963, Hodges, Patterson et al. 1992, Warburton, Price et al. 1999, Thiel, Habedank
1010 et al. 2005, Nardo, Holland et al. 2017, Hope, Leff et al. 2018). In other words, what are the causal
1011 mechanisms associated with (language) behavioural modifications following neurological disorders?
1012 Despite valiant efforts, none of the current computational accounts of language can fully explain these
1013 behaviours (Rueschemeyer, Gaskell et al. 2018): examples include Directions Into Velocities of
1014 Articulators model (Tourville and Guenther 2011), State Feedback Control model (Houde and Nagarajan
1015 2011), and Hierarchical State Feedback Control model (Hickok 2014). Crucially, these approaches do not
1016 simultaneously account for higher-order language processing (semantic, syntactic, *etc.*) and lower level
1017 articulatory control (prosody, *etc.*); however, human language processing requires both. The active listening
1018 scheme presented here departs from previous approaches: it explicitly considers the segmentation of
1019 continuous signals (which come into play through the accuracy term in Equation (A.18) and relate to lower-
1020 level processing) and beliefs about the lexical content of those signals (key to the complexity term in
1021 Equation (A.18) and relating to higher-level language processing). Not only do these two aspects exist in
1022 the model, but they go hand-in-hand during word recognition. This makes the generative model described
1023 here a prime candidate for developing a mechanistic and neurobiologically plausible account of (healthy

Active listening

1024 and impaired) language behaviour.

1025 The idea that a generative model for speech generation can be inverted for the purpose of recognising speech
1026 touches upon a longstanding debate in the literature—are similar neural processes used to recognise speech,
1027 as those that are used to produce speech? This is an interesting question, and one that the current formulation
1028 does not address. Of relevance, the properties of spoken sentences that active listening uses to produce and
1029 recognise speech are acoustic (e.g., fundamental and formant frequencies) rather than biological (e.g., vocal
1030 chords and vocal tract) attributes (Guenther and Vladusich 2012). Thus, it does not necessarily follow from
1031 this framework that an individual who is unable to speak is unable to comprehend speech. On the contrary,
1032 we expect that an individual who is unable to speak could still generate an internal model that specifies the
1033 causes of spoken words, which they have learnt by perceiving speech. Whether the experience of producing
1034 speech contributes to the same model is an interesting question. In short, there may be an opportunity to
1035 examine how computational lesions to the model impair speech perception and production.

1036 **Active listening and voice recognition**

1037 One strength of the current scheme is that it deals with both speech generation and recognition, and can be
1038 iteratively applied to recognise the lexical content of simulated speech (see Figure 9). The simulated speech
1039 that the model produces is discernibly artificial, but the key message here is that the model reduces the
1040 problems of speech generation and recognition to their necessary parameters. The generative model
1041 introduced in this paper lays the groundwork for a complete model of voice recognition. In other words, a
1042 model that infers *who* is speaking. The current model includes states for the speaker attributes of their
1043 average fundamental frequency and formant spacing. From a speech production perspective, a speaker's
1044 fundamental frequency relates to the rate of vocal fold vibration (known as glottal pulse rate), and formant
1045 spacing is affected by the length and shape of the vocal tract—which are relatively fixed for a speaker,
1046 although can be modified slightly by changing the positions of the articulators, such as the tongue and lips.
1047 Previous research demonstrates that listeners use both fundamental frequency and speech formants to judge
1048 the identity of people who are familiar (LaRiviere 1975, Abberton and Fourcin 1978, Van Dommelen 1987,
1049 Van Dommelen 1990, Lavner, Gath et al. 2000, Lavner, Rosenhouse et al. 2001, Holmes, Domingo et al.
1050 2018) and unfamiliar (Matsumoto, Hiki et al. 1973, Walden, Montgomery et al. 1978, Murry and Singh
1051 1980, Baumann and Belin 2009, Gaudrain, Li et al. 2009). To extend the current model to recognise voices,

Active listening

1052 the next step is to specify how combinations of fundamental and formant frequencies are used to infer
1053 speaker identity. From the perspective of the generative model, fundamental and formant frequencies are
1054 generated from hidden states that correspond to particular speakers. This approach differs from that
1055 proposed by Kleinschmidt and Jaeger (Kleinschmidt and Jaeger 2015), who assume that listeners construct
1056 a separate generative model for each talker they encounter. In the current implementation, we have focused
1057 on fundamental and formant frequencies, because these attributes are most prevalent in the voice
1058 recognition literature. However, they are not the only relevant speaker attributes (Cai, Gilbert et al. 2017,
1059 Holmes, Domingo et al. 2018). More complex models of voice recognition could incorporate additional
1060 speaker parameters, for example, relating to speaker-specific accent, stress, and intonation.

1061 **Active listening and music**

1062 Finally, the generative and inversion schemes presented here could also form the basis for models of other
1063 complex auditory signals. Music, for example, shares several features with language (Patel 2010) and relies
1064 on partly overlapping brain networks (Musso, Weiller et al. 2015), which makes it a natural choice for
1065 future work. It is not difficult to imagine how the generative model in Figure 1 could be adapted to simulate
1066 music in an active listening framework. For example, somewhat akin to determining the correct onsets and
1067 offsets of word boundaries, we need to decide where a musical phrase—or longer section of music—begins
1068 and ends.

1069 Recent empirical findings have shown that mismatch responses to unexpected musical sounds are larger in
1070 contexts with low than high uncertainty (Quiroga-Martinez, Hansen et al. 2019). This fits comfortably with
1071 the proposed explanation of evoked responses as reflecting Bayesian surprise or salience, which would be
1072 reduced when sensory signals are unreliable or imprecise. Since music is rich and multifaceted and relies
1073 greatly on statistical learning (Pearce 2018), it would be an ideal means to understand how neuronal
1074 dynamics change with uncertainty.

1075 **Summary**

1076 In summary, this paper introduces active listening—a unified framework for generating and recognising

Active listening

1077 speech. The generative model specifies how discrete *lexical*, *prosodic*, and *speaker* attributes give rise to a
1078 continuous acoustic timeseries. As the name implies, the framework also includes an active component, in
1079 which plausible segmentations of the acoustic timeseries—corresponding to the placement of word
1080 boundaries—are considered, and segmentation that minimises Bayesian surprise is selected. In the
1081 simulations presented here, we demonstrate that speech can be iteratively recognised and generated under
1082 this model. We show that the words that the model recognises depend on prior expectations about the
1083 content of the words, as is the case for human listeners, and that simulated neuronal responses resemble
1084 human electrophysiological responses. This work establishes a foundation for future work that will simulate
1085 human conversations, voice recognition, speech-in-noise, and music—and which we anticipate will provide
1086 key insights into neuropsychological impairments to language processing.

1087

1088

Software note

1089 The routines described in this paper are available as Matlab code in the SPM academic software:
1090 <http://www.fil.ion.ucl.ac.uk/spm/>. The simulations reported in the figures can be reproduced (and
1091 customised) via a graphical user interface by typing (in the Matlab command window) **DEM** and selecting
1092 appropriate (speech recognition) demonstration routines. The accompanying Matlab scripts are called
1093 **spm_voice_*.m**.

1094

1095

Acknowledgements

1096 The Wellcome Trust funded K.J.F. (Ref: 088130/Z/09/Z), E.H. (Ref: WT091681MA), and the Wellcome
1097 Centre for Human Neuroimaging (Ref: 203147/Z/16/Z), where this work was conducted. N.S. is funded by
1098 the Medical Research Council (Ref: MR/S02522/1). D.R.Q. is funded by the Danish National Research
1099 Foundation (Project number: DNRF117). T.P. is supported by the Rosetrees Trust (Award number:
1100 173346).

Active listening

1101

Disclosure statement

1102 The authors have no disclosures or conflict of interest.

1103

1104

Appendices

1105 **Appendix 1: The generative model**

1106 This appendix covers technical details of the generative model introduced in Figure 1. Figure 11 is designed
1107 to supplement Figure 1, and includes the equations corresponding to word generation (left column) and
1108 word recognition (right column). This section first provides a summary of the technical details of the
1109 generative model, then goes on to unpack each of the equations of the generative model in Figure 11.
1110 Although these may seem complicated for a non-technical reader, they are simply a sequence of non-linear
1111 transforms that specify the mapping from lexical, speaker, and prosody parameters to an acoustic timeseries.

1112 In brief, each word (i.e., lexical item) is associated with a matrix of a discrete cosine transform coefficients
1113 ($\theta^{\mathcal{L}}$) that generate a time-frequency representation (W) of the spoken word (i.e., the spectrogram), when
1114 combined with speaker and prosody information. In this scheme, the lexical form and structure comprise a
1115 discrete cosine transform with 8 basis functions over time and 32 over formant frequencies (see Figure
1116 11C). The number of basis functions was selected as a compromise between the quality of the generated
1117 acoustic timeseries and computational efficiency. Each column of the time-frequency representation
1118 generates a transient: thus, the number of transients corresponds to the number of columns in the time-
1119 frequency representation.

1120 The transients are emitted at an instantaneous fundamental frequency, which is inversely proportional to
1121 the time intervals between successive transients (Δ_i). These time intervals are stored in a fundamental
1122 interval variable (I). The instantaneous fundamental frequency is affected by the average fundamental
1123 frequency of the speaker (θ^0), corresponding to their average *glottal pulse rate*. It also depends on a discrete

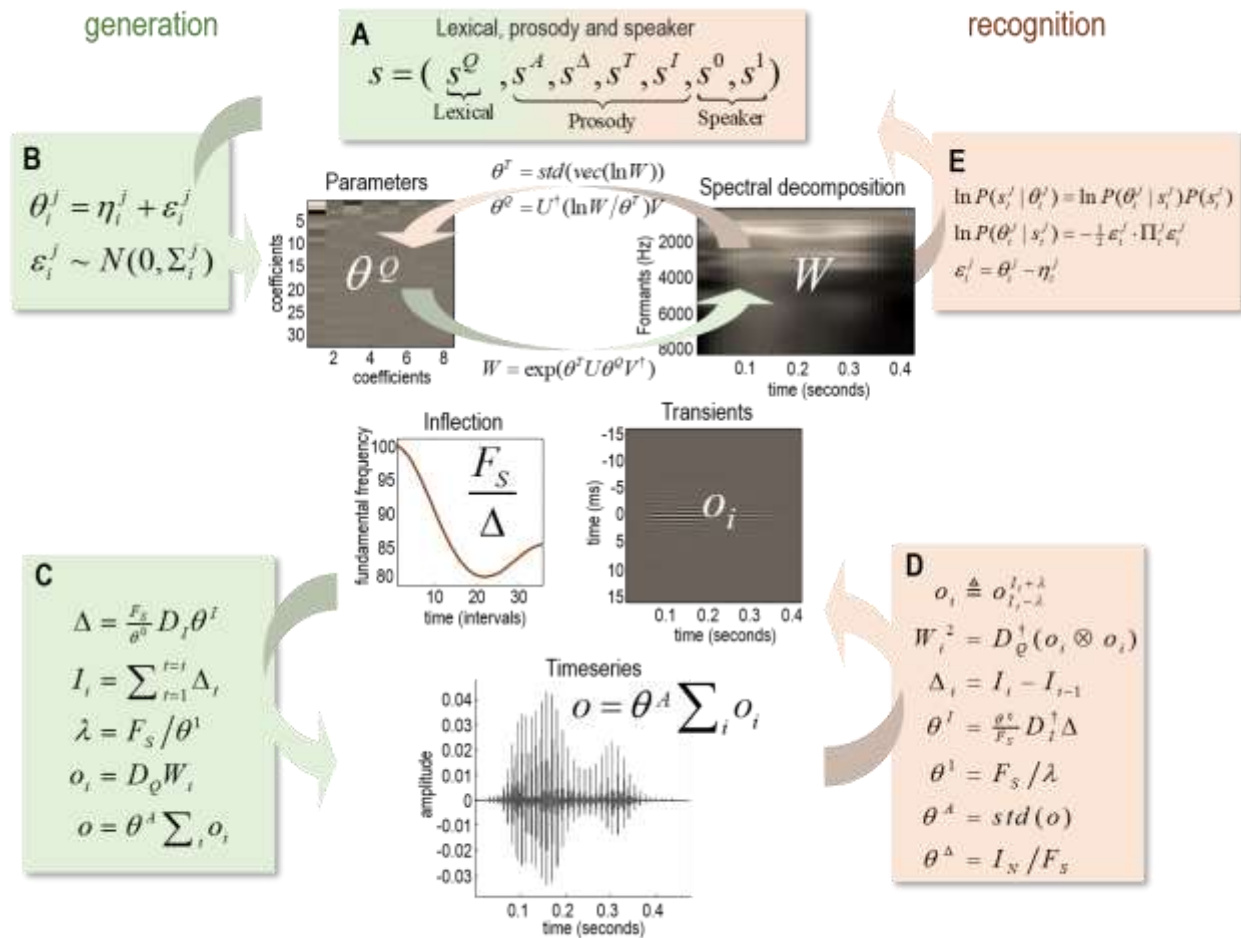
Active listening

1124 cosine transform (D) based upon (three) coefficients (θ^i) that encode inflection around the speaker's
1125 average fundamental frequency (θ^0): (1) the average fundamental frequency relative to the speaker average,
1126 (2) increases or decreases in fundamental frequency over time, and (3) the acceleration or deceleration of
1127 changes in fundamental frequency. The ensuing time-frequency representation is then multiplied by an
1128 inverse temperature (θ^T) parameter, which affects the quality of the sound and can be thought of as a timbre
1129 parameter. Its exponential is, effectively, Fourier transformed to create a succession of transients that are
1130 deployed over fundamental intervals. The resulting timeseries is then scaled by an amplitude parameter (θ
1131 ^A) to furnish the final (continuous) acoustic timeseries.

1132

1133

Active listening



1134

1135

FIGURE 11

1136 *A generative model of a word.* This figure illustrates the generative model from the perspective of word generation
 1137 (green panels) and accompanying inversion (orange panels), which corresponds to word recognition. This model maps
 1138 from hidden states (s ; shown in box A), which denote the attributes of a spoken word (in this case lexical content,
 1139 prosody, and speaker identity), to outcomes (o ; shown in box C), which corresponds to the continuous acoustic
 1140 timeseries. Box B shows how parameters are sampled for word generation. The centre panels illustrate the non-linear
 1141 mappings between model parameters and the acoustic spectrum (i.e., time-frequency representation). Box C specifies
 1142 how the transients are then aggregated to form a timeseries. Recognition (boxes D–E) corresponds to the inversion of
 1143 the generative model: a given time series is transformed to parameterise the time-frequency representation (box D) by
 1144 simply inverting or ‘undoing’ the generative operations. These parameters are used to evaluate the likelihood of
 1145 lexical, prosody and speaker states (box E). The equations displayed in this figure are unpacked in the text.

Active listening

1146 In what follows, we unpack each of the equations in Figure 11, from the perspective of word generation
1147 (left column of Figure 11). Note that word generation simply involves a sequence of non-linear
1148 transformations, which specify the relationship between parameters and the acoustic timeseries.

1149 Each discrete state generates a parameter that is sampled from a Gaussian distribution (Figure 11B) with a
1150 mean η and covariance Σ . The subscript notation indicates hidden state j and its i -th possible value:

$$\begin{aligned} 1151 \quad \theta_i^j &= \eta_i^j + \varepsilon_i^j \\ \varepsilon_i^j &\sim N(0, \Sigma_i^j) \end{aligned} \tag{A.1}$$

1152 The spectrum is constructed from frequency (U) and temporal (V) basis functions, which are combined with
1153 a matrix of coefficients (θ^ℓ) corresponding to lexical parameters. The spectrum is scaled with an inverse
1154 temperature (i.e., precision; θ^T) parameter, which is then exponentiated to create a matrix of fluctuations
1155 W of (formant) frequencies over time:

$$1156 \quad W = \exp(\theta^T U \theta^\ell V^\dagger) \tag{A.2}$$

1157 Each column of W is transformed into a transient as a function of time (using discrete cosine transform
1158 matrix D):

$$1159 \quad o_i = D_\ell W_i \tag{A.3}$$

1160 The duration of the transients (λ) is determined by the speaker formant spacing (θ^1)—such that a high
1161 formant spacing value squashes (shortens) the transients, rendering the frequencies higher when placed in
1162 the timeseries. F_s indicates the sampling rate of the audio timeseries:

$$1163 \quad \lambda = F_s / \theta^1 \tag{A.4}$$

1164 The spacing (Δ) of the transients is inversely proportional to the speaker fundamental frequency parameter

Active listening

1165 (θ^0) , and is also affected by inflections due to prosody (θ^l) :

$$1166 \quad \Delta = \frac{F_s}{\theta^0} D_l \theta^l \quad (\text{A.5})$$

1167 A fundamental interval (I) variable stores the absolute positions of all of the transients:

$$1168 \quad I_i = \sum_{t=1}^{t=i} \Delta_t \quad (\text{A.6})$$

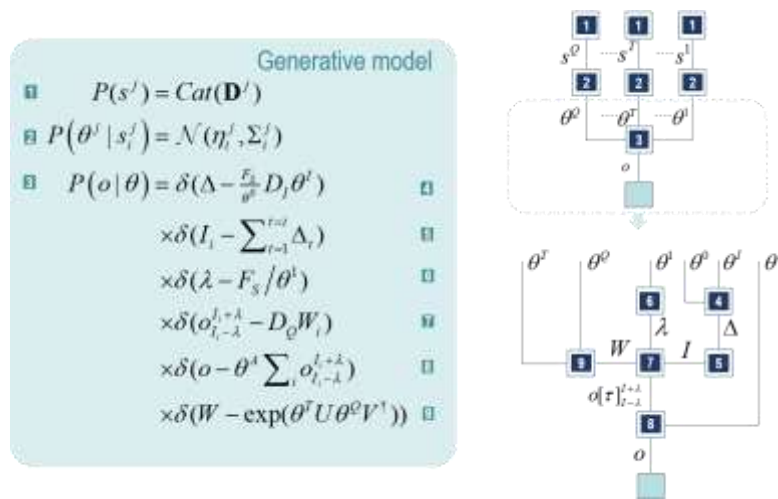
1169 The timeseries (o) is constructed by summing the transients and multiplying this by the amplitude
1170 parameter:

$$1171 \quad o = \theta^A \sum_i o_i \quad (\text{A.7})$$

1172 For readers familiar with graphical formulations of generative models, Figure 12 illustrates the same model
1173 in factor graph form (Forney 2001). This provides an alternative visual representation of the generative
1174 model, and highlights inferences based on message passing. This perspective is used below to describe the
1175 form of local (neuronal) message passing that underwrites simulated electrophysiological responses.

1176

Active listening



1177

1178

FIGURE 12

1179 *A graphical formulation of the generative model.* This figure illustrates the same model as described in Figure 11, but
 1180 uses a normal (Forney) factor graph form. This graphical notation relies upon the factorisation of the probability
 1181 density that underwrites the generative model. Each factor is specified in the panel on the left. Factor 1 is the prior
 1182 probability associated with the hidden states and takes a categorical form. Factor 2 is a normal distribution that
 1183 specifies the dependence of parameters on states. Each discrete state is associated with a different expectation and
 1184 covariance for the parameters. Factor 3 describes how the observed timeseries is generated from the parameters, and
 1185 this is decomposed into factors 4–9. These are Dirac delta functions that may be thought of as normal distributions,
 1186 centred on zero, with infinite precision (i.e., zero covariance). In the graphs on the right, factors are indicated by
 1187 numbered squares, and these are connected by edges (Hasson, Yang et al.), which represent the variables common to
 1188 the factors they connect. The upper right graph shows factors 1–3, and the lower graph unpacks factor 3 in terms of
 1189 factors 4–9. The process of generating data may be thought of in terms of a series of local operations taking place at
 1190 each factor from top to bottom (i.e., sample states from factor 1, then parameters from factor 2, then perform the series
 1191 of operations in factor 3 to get the timeseries). The recognition process can be thought of as bidirectional message
 1192 passing across each factor node, such that empirical priors and likelihoods are combined at each edge to form posterior
 1193 beliefs about the associated variable. Factor 5 is of particular interest here, as it determines the internal ‘action’ that
 1194 selects the interval for segmentation.

1195

Active listening

1196 **Appendix 2: Model inversion or word recognition**

1197 Next, we turn our attention to word recognition (right column of Figure 11). Inversion of the generative
1198 model simply requires ‘undoing’ the sequence of events that we used for word generation. Like word
1199 generation, word recognition simply requires a series of non-linear transforms—except, for word
1200 recognition, we map from epochs of the acoustic signal to discrete *lexical*, *speaker*, and *prosody* parameters.

1201 In brief, the recognition scheme comprises the following steps. The peak energy of the auditory timeseries
1202 is identified by convolving its absolute values with a Gaussian kernel. A one second epoch, centred on the
1203 peak, is selected as a signal to search for the onset and offset of the word (although in principle this epoch
1204 could be any length). Onsets and offsets are identified based on threshold crossings of the amplitude
1205 envelope. Here, the amplitude envelope is calculated from the absolute values of the timeseries convolved
1206 with a Gaussian kernel. This is, for all practical purposes, equivalent to the absolute values of the Hilbert
1207 transform, but is computationally more efficient. The threshold we use here is $1/16^{\text{th}}$ of the maximum
1208 envelope value across the window, after subtracting the minimum; this value was selected to be above the
1209 noise floor.

1210 The fundamental interval function is estimated using a discrete cosine transform (with three coefficients)
1211 of the fundamental intervals. The fundamental intervals are defined as phase crossings following a Hilbert
1212 transform and bandpass filtering around the prior for the speaker average fundamental frequency (e.g., 100
1213 Hz, with a standard deviation of 8 Hz).

1214 Equipped with the fundamental interval function, the formant frequencies are then estimated by evaluating
1215 the cross-covariance function over short segments centred on each fundamental interval. The duration of
1216 these segments corresponds to the inverse of the first formant frequency. The formant frequencies *per se*
1217 are evaluated using a modified (by retaining even terms) discrete cosine transform at each slice, to evaluate
1218 the spectral density over the acoustic range (in 256 frequency bins, where each bin is determined by the
1219 formant spacing; for example, with a formant spacing of 32 Hz, the highest spectral density is 8000 Hz).
1220 Following a log transform and normalisation, fluctuations in (log) spectral density are recovered with a
1221 discrete cosine transform with 32 basis functions over (formant) frequencies and eight basis functions over
1222 intervals. The inverse temperature (timbre) parameter corresponds to the standard deviation of these lexical
1223 (formant frequency) parameters, which is used to normalise the lexical (32x8) parameter matrix.

Active listening

1224 To infer the lexical content, prosody and speaker, the MAP parameter estimates above can be used to
1225 evaluate the likelihood of each discrete attribute. As described in the main text, the likelihoods are combined
1226 with a prior to produce a posterior categorical distribution over the attributes in question. For the prosody
1227 parameters, each parameter is divided into eight bins and the likelihood of belonging to any particular bin
1228 is evaluated under Gaussian assumptions as above; using *a priori* means and precisions of the discrete levels
1229 of each prosody attribute (i.e., amplitude, duration, timbre, inflection). Similarly, the categorical speaker
1230 identity is determined by a 16 x 16 discrete states space, covering fundamental and formant frequencies.

1231 In what follows, we unpack each of the equations in Figure 11—this time, from the perspective of word
1232 recognition (right column of Figure 11).

1233 The amplitude parameter is the standard deviation of the timeseries (o):

$$1234 \quad \theta^A = std(o) \tag{A.8}$$

1235 Each transient (o_i) is defined as an interval of the timeseries, based on the positions of fundamental intervals
1236 (I) and transient durations (λ):

$$1237 \quad o_i \triangleq o_{I_i - \lambda}^{I_i + \lambda} \tag{A.9}$$

1238 The spacing (Δ) of the transients corresponds to the difference between successive fundamental intervals
1239 (I):

$$1240 \quad \Delta_i = I_i - I_{i-1} \tag{A.10}$$

1241 Inflection parameters are proportional to the speaker fundamental frequency (θ^0) and are constructed using
1242 discrete cosine transform matrix D . F_s indicates the sampling rate of the audio timeseries:

$$1243 \quad \theta^I = \frac{\theta^0}{F_s} D_i^\dagger \Delta \tag{A.11}$$

Active listening

1244 The formant scaling parameter (θ^l) is inversely proportional to the transient duration (λ):

$$1245 \quad \theta^l = F_s / \lambda \quad (\text{A.12})$$

1246 The duration parameter (θ^Δ) is proportional to the fundamental interval (I):

$$1247 \quad \theta^\Delta = I_N / F_s \quad (\text{A.13})$$

1248 The (squared) matrix of fluctuations of (formant) frequencies over time (W) is constructed from the
1249 transients using discrete cosine transform matrix D :

$$1250 \quad W_i^2 = D_Q^\dagger (o_i \otimes o_i) \quad (\text{A.14})$$

1251 The timbre parameter (θ^T) is the standard deviation of the log spectral decomposition:

$$1252 \quad \theta^T = \text{std}(\text{vec}(\ln W)) \quad (\text{A.15})$$

1253 Lexical parameters (θ^ϱ) are a matrix of coefficients that control the joint expression of formant frequency
1254 and temporal basis functions. These are calculated from the frequency (U) and temporal (V) basis functions
1255 and the log spectral decomposition, scaled by the timbre parameter:

$$1256 \quad \theta^\varrho = U^\dagger (\ln W / \theta^T) V \quad (\text{A.16})$$

1257 The parameters are used to evaluate the likelihood of lexical, prosody and speaker states, as shown in the
1258 following equations:

Active listening

$$\begin{aligned}
 & \ln P(s_i^j | \theta_i^j) = \ln P(\theta_i^j | s_i^j) P(s_i^j) \\
 1259 \quad & \ln P(\theta_i^j | s_i^j) = -\frac{1}{2} \varepsilon_i^j \cdot \Pi_i^j \varepsilon_i^j \tag{A.17} \\
 & \varepsilon_i^j = \theta_i^j - \eta_i^j
 \end{aligned}$$

1260

1261 **Appendix 3: Speech segmentation as an active process**

1262 In the current framework, speech segmentation is treated as a covert action from a computational
 1263 perspective: We select boundary pairs (I_0 and I_T) and evaluate their free energy under prior beliefs about
 1264 the word. Formally, this can be expressed as minimising free energy both with respect to (approximate)
 1265 posterior beliefs about the attributes of the word (Q) and the intervals selected (I_0, I_T):

$$\begin{aligned}
 Q &= \arg \min_Q F(Q, o_{I_0}^{I_T}) \\
 (I_0, I_T) &= \arg \min_I F(Q, o_{I_0}^{I_T})
 \end{aligned}$$

$$\begin{aligned}
 1266 \quad F(Q, o) &= E_Q[\ln Q(s) - \ln P(o, s)] \tag{A.18} \\
 &= \underbrace{E_Q[\ln Q(s) - \ln P(s | o)]}_{\text{Evidence bound}} - \underbrace{\ln P(o)}_{\text{Log evidence}} \\
 &= \underbrace{E_Q[\ln Q(s) - \ln P(s)]}_{\text{Complexity}} - \underbrace{E_Q[\ln P(o | s)]}_{\text{Accuracy}} \geq -\ln P(o)
 \end{aligned}$$

1267 Choosing the interval with the smallest free energy effectively selects the interval that maximises the
 1268 evidence or marginal likelihood of auditory outcomes contained in that interval; namely, $P(o)$. This follows
 1269 because the variational free energy, by construction, represents an upper bound on log evidence. In (A.18),
 1270 the free energy is expressed in terms of *log evidence* and an *evidence bound*. It is also expressed as the
 1271 difference between *complexity* and *accuracy* by rearranging the equation. Complexity is the Kullback-
 1272 Leibler divergence between a posterior over latent states $Q(s)$, and prior beliefs $P(s)$, while accuracy is the
 1273 expected log likelihood of auditory signals contained in the interval in question. Importantly, both posterior
 1274 beliefs about latent states (i.e., *lexical*, *prosody*, and *speaker*) and the active selection of acoustic intervals

Active listening

1275 optimise free energy. This is the signature of active inference. In this instance, the posterior beliefs obtain
1276 from the likelihood of the lexical, prosody and identity parameters, given the associated states. From Figure
1277 11, the optimal posterior beliefs satisfy (A.18) when (ignoring constants):

$$\begin{aligned} \ln Q(s_i^j) &= \ln P(s_i^j | o) = \ln P(s_i^j | \theta_i^j) \\ &= \ln P(s_i^j) + \ln P(\theta_i^j | s_i^j) \\ 1278 \quad &= \ln P(s_i^j) - \frac{1}{2} \varepsilon_i^j \cdot \Pi_i^j \varepsilon_i^j && \text{(A.19)} \\ &\Rightarrow \\ &F(Q, o) = -\ln P(o) \end{aligned}$$

1279 Here, Π is the prior precision of lexical parameters from Figure 11. The second equality on the first line
1280 may seem a little counterintuitive, but rests upon the assumed relationship between the parameters and the
1281 timeseries. The equality holds in virtue of the absence of random fluctuations in this mapping, such that a
1282 given parameter deterministically generates time-series data. In other words, the implicit conditional
1283 probability density describing the generation of the timeseries from the parameters (and the associated
1284 posterior distribution over parameters) takes the form of a Dirac delta function. The last equality reflects
1285 the fact that when the evidence bound in **Error! Reference source not found.** collapses to zero, free
1286 energy becomes negative log evidence. The subscript notation indicates the value that a discrete state might
1287 take (i.e. $P(s_i)$ should be read as ‘the probability that the hidden state j takes its i -th possible value’).

1288 From the equations above, it should be clear that we can identify a variety of *candidate* boundaries for
1289 words and evaluate their free energy to select the final parsing of the acoustic signal. But where should
1290 these candidate boundaries be placed? In an extreme case, we could place boundaries at every combination
1291 of time points within the acoustic signal—but that would be computationally inefficient given that we can
1292 reduce the scope of possibilities by using sensible priors. Here, we use the simple prior that word boundaries
1293 are more likely to occur at local minima of the amplitude envelope—so these are the boundaries that we
1294 choose to evaluate.

1295 Practically, based upon the spectral content of speech, we estimate the amplitude envelope by removing
1296 low frequencies up to about 512 Hz. The envelope is then simply the average of the ensuing absolute values,
1297 smoothed with a Gaussian kernel (with a standard deviation of $F_S/16$). This method is less computationally

Active listening

1298 demanding than using the absolute values of the Hilbert transform, yet practically gives the same result in
1299 this setting.

1300

1301 **Appendix 4: Belief updating and neuronal dynamics**

1302 The form of neuronal dynamics is calculated by constructing ordinary differential equations whose solution
1303 satisfies Equation (A.18). Using $\mathbf{v} = \ln \mathbf{s}$ to denote the log of the approximate posterior expectation about
1304 hidden states and introducing a prediction error ($\boldsymbol{\varepsilon}$) one obtains the following update scheme (Friston,
1305 FitzGerald et al. 2017) (dropping the superscript j for clarity):

$$\mathbf{v}_i \triangleq \ln Q(s_i)$$

$$\mathbf{s}_i \triangleq Q(s_i)$$

$$1306 \quad \boldsymbol{\varepsilon}_i = \ln P(s_i) - \frac{1}{2} \boldsymbol{\varepsilon}_i \cdot \Pi_i \boldsymbol{\varepsilon}_i - \mathbf{v}_i \tag{A.20}$$

$$\dot{\mathbf{v}}_i = \boldsymbol{\varepsilon}_i$$

$$\mathbf{s} = \sigma(\mathbf{v})$$

1307 Here, σ denotes the softmax (normalised exponential) function and Π is the prior precision of lexical
1308 parameters from Figure 11. The prediction error ($\boldsymbol{\varepsilon}$) is the difference between the optimal log posterior and
1309 current estimate of this (\mathbf{v}). The log posterior, via Bayes theorem, is equal to the sum of the log prior and
1310 the log likelihood (minus a normalisation constant). As the likelihood is assumed to be normally distributed,
1311 its log is quadratic in the difference ($\boldsymbol{\varepsilon}$) between the mode and lexical parameters. The mode of this
1312 distribution is different under each state, so the likelihood of a given parameter value varies with states. For
1313 readers familiar with clustering procedures, this is like having a series of clusters (states) with different
1314 centroids (i.e., modes of the likelihood).

1315 The prediction error ($\boldsymbol{\varepsilon}$) is the (negative) free energy gradient that drives neuronal dynamics. Intuitively, the
1316 fourth line of Equation A.20 drives \mathbf{v} to change until it is equal to the Bayes optimal posterior, at which

Active listening

1317 point $\boldsymbol{\varepsilon}$ is zero. To account for the normalisation constant that would have appeared in Bayes theorem, the
 1318 conversion from \mathbf{v} to \mathbf{s} requires not only that we exponentiate (i.e., convert a log probability into a
 1319 probability), but that we normalise the result. This ensures that \mathbf{s} comes to encode a vector of posterior
 1320 probabilities for each hidden state.

1321 The sigmoid (softmax) function in Equation A.20 can be thought of as a sigmoid (voltage–firing rate)
 1322 activation function, which mediates competition among posterior expectations. Equation A.20 therefore,
 1323 provides a process theory for neuronal dynamics. Based on this equation, log expectations about hidden
 1324 states can be associated with depolarisation of neurons or neuronal populations encoding expectations about
 1325 hidden states (\mathbf{v}_i), while firing rates (\mathbf{s}_i) encode expectations *per se*. The simulated responses in Figure 6
 1326 use a finite difference scheme that has the same solution as A.20:

$$\begin{aligned}
 \mathbf{v}(\tau)_i &= \ln \mathbf{s}(\tau)_i \\
 \boldsymbol{\varepsilon}(\tau)_i &= \ln P(s_i) - \frac{1}{2} \boldsymbol{\varepsilon}_i \cdot \Pi_i \boldsymbol{\varepsilon}_i - \mathbf{v}_i \\
 \mathbf{s}(\tau + d\tau)_i &= \sigma(\mathbf{v}_i + \kappa \cdot \boldsymbol{\varepsilon}_i)
 \end{aligned}
 \tag{A.21}$$

1328 where κ is chosen to reproduce dynamics at a plausible, neuronal timescale.

1329 When considering electrophysiological responses in terms of belief updating, our formal interpretation
 1330 relates to Equation (A.20), which suggests that depolarisation corresponds to the log posterior. The change
 1331 in depolarisation is the difference between the log posterior and prior expectations. The average of these
 1332 differences is the Kullback-Leibler divergence between the posterior and prior:

$$\begin{aligned}
 \mathbf{v}_i &= \ln Q(s_i) \\
 \mathbf{v}(\tau)_i - \mathbf{v}(0)_i &= \ln Q(s_i) - \ln P(s_i) \\
 &\Rightarrow \\
 E_Q[\mathbf{v}(\tau) - \mathbf{v}(0)] &= E_Q[\ln Q(s) - \ln P(s)] = D[Q(s) \parallel P(s)]
 \end{aligned}
 \tag{A.22}$$

1334

Active listening

References

1335

- 1336 Abberton, E. and A. J. Fourcin (1978). "Intonation and Speaker Identification." *Language and Speech* **21**(4): 305-318.
- 1337 Adams, R. A., S. Shipp and K. J. Friston (2013). "Predictions not commands: active inference in the motor system."
- 1338 *Brain Struct Funct.* **218**(3): 611-643.
- 1339 Aitchison, L. and M. Lengyel (2017). "With or without you: predictive coding and Bayesian inference in the brain."
- 1340 *Current opinion in neurobiology* **46**: 219-227.
- 1341 Alain, C., S. R. Arnott, S. Hevenor, S. Graham and C. L. Grady (2001). "'What' and 'where' in the human auditory
- 1342 system." *Proceedings of the National Academy of Sciences* **98**(21): 12301-12306.
- 1343 Altenberg, E. P. (2005). "The perception of word boundaries in a second language." *Second Language Research* **21**(4):
- 1344 325-358.
- 1345 Andreopoulos, A. and J. Tsotsos (2013). "A computational learning theory of active object recognition under
- 1346 uncertainty." *International journal of computer vision* **101**(1): 95-142.
- 1347 Aring, C. D. (1963). "Traumatic Aphasia: A Study of Aphasia in War Wounds of the Brain." *JAMA Neurology* **8**(5):
- 1348 579-580.
- 1349 Attwell, D. and C. Iadecola (2002). "The neural basis of functional brain imaging signals." *Trends in Neurosciences*
- 1350 **25**(12): 621-625.
- 1351 Balaguer, R. D. D., J. M. Toro, A. Rodriguez-Fornells and A.-C. Bachoud-Lévi (2007). "Different neurophysiological
- 1352 mechanisms underlying word and rule extraction from speech." *PLoS One* **2**(11): e1175.
- 1353 Bar, M., K. S. Kassam, A. S. Ghuman, J. Boshyan, A. M. Schmid, A. M. Dale, M. S. Hämäläinen, K. Marinkovic, D.
- 1354 L. Schacter and B. R. Rosen (2006). "Top-down facilitation of visual recognition." *Proceedings of the national*
- 1355 *academy of sciences* **103**(2): 449-454.
- 1356 Barto, A., M. Mirolli and G. Baldassarre (2013). "Novelty or Surprise?" *Frontiers in Psychology* **4**.
- 1357 Bashford, J. A., Jr., R. M. Warren and P. W. Lenz (2008). "Evoking biphone neighborhoods with verbal
- 1358 transformations: illusory changes demonstrate both lexical competition and inhibition." *J Acoust Soc Am* **123**(3):
- 1359 E132.
- 1360 Bastos, A. M., W. M. Usrey, R. A. Adams, G. R. Mangun, P. Fries and K. J. Friston (2012). "Canonical microcircuits
- 1361 for predictive coding." *Neuron* **76**(4): 695-711.
- 1362 Baumann, O. and P. Belin (2009). "Perceptual scaling of voice identity: Common dimensions for different vowels and
- 1363 speakers." *Psychological Research* **74**(1): 110--120.
- 1364 Beal, M. J. (2003). "Variational Algorithms for Approximate Bayesian Inference." *PhD. Thesis, University College*
- 1365 *London*.
- 1366 Beckman, M. E. and J. Edwards (1990). "of prosodic constituency." *Between the grammar and physics of speech:*
- 1367 152.
- 1368 Belin, P., S. Fecteau and C. Bzdard (2004). "Thinking the voice: Neural correlates of voice perception." *Trends in*
- 1369 *Cognitive Sciences* **8**(3): 129--135.
- 1370 Belin, P. and R. J. Zatorre (2000). "'What', 'where' and 'how' in auditory cortex." *Nature Neuroscience* **3**(10): 965--
- 1371 966.
- 1372 Bennett, C. H. (2003). "Notes on Landauer's principle, reversible computation, and Maxwell's Demon." *Studies in*
- 1373 *History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics* **34**(3): 501-510.
- 1374 Billig, A. J., M. H. Davis, J. M. Deeks, J. Monstrey and R. P. Carlyon (2013). "Lexical influences on auditory
- 1375 streaming." *Current Biology* **23**(16): 1585--1589.
- 1376 Bogacz, R. (2017). "A tutorial on the free-energy framework for modelling perception and learning." *Journal of*
- 1377 *Mathematical Psychology* **76**: 198--211.
- 1378 Bögels, S., H. Schriefers, W. Vonk and D. J. Chwilla (2011). "Prosodic Breaks in Sentence Processing Investigated
- 1379 by Event-Related Potentials." *Language and Linguistics Compass* **5**(7): 424-440.
- 1380 Bögels, S., H. Schriefers, W. Vonk and D. J. Chwilla (2011). "The role of prosodic breaks and pitch accents in
- 1381 grouping words during on-line sentence processing." *Journal of Cognitive Neuroscience* **23**(9): 2447-2467.
- 1382 Braiman, C., E. A. Fridman, M. M. Conte, H. U. Voss, C. S. Reichenbach, T. Reichenbach and N. D. Schiff (2018).
- 1383 "Cortical Response to the Natural Speech Envelope Correlates with Neuroimaging Evidence of Cognition in Severe
- 1384 Brain Injury." *Curr Biol* **28**(23): 3833-3839.e3833.
- 1385 Brown, H., R. A. Adams, I. Parees, M. Edwards and K. J. Friston (2013). "Active inference, sensory attenuation and
- 1386 illusions." *Cognitive Processing* **14**(4): 411--427.

Active listening

- 1387 Brown, H., K. J. Friston and S. Bestmann (2011). "Active inference, attention, and motor preparation." Frontiers in
1388 psychology **2**: 218.
- 1389 Brungart, D. S. (2001). "Evaluation of speech intelligibility with the coordinate response measure." The Journal of the
1390 Acoustical Society of America **109**(5 Pt 1): 2276--2279.
- 1391 Brungart, D. S., B. D. Simpson, M. A. Ericson and K. R. Scott (2001). "Informational and energetic masking effects
1392 in the perception of multiple simultaneous talkers." The Journal of the Acoustical Society of America **110**(5): 2527--
1393 2538.
- 1394 Cai, Z. G., R. A. Gilbert, M. H. Davis, M. G. Gaskell, L. Farrar, S. Adler and J. M. Rodd (2017). "Accent modulates
1395 access to word meaning: Evidence for a speaker-model account of spoken word recognition." Cognitive Psychology
1396 **98**: 73-101.
- 1397 Christie Jr, W. M. (1974). "Some cues for syllable juncture perception in English." the Journal of the Acoustical
1398 Society of America **55**(4): 819-821.
- 1399 Cole, R. A., J. Jakimik and W. E. Cooper (1980). "Segmenting speech into words." The Journal of the Acoustical
1400 Society of America **67**(4): 1323-1332.
- 1401 Connolly, J. F. and N. A. Phillips (1994). "Event-related potential components reflect phonological and semantic
1402 processing of the terminal word of spoken sentences." Journal of cognitive neuroscience **6**(3): 256-266.
- 1403 Connolly, J. F., N. A. Phillips, S. H. Stewart and W. Brake (1992). "Event-related potential sensitivity to acoustic and
1404 semantic properties of terminal words in sentences." Brain and language **43**(1): 1-18.
- 1405 Cunillera, T., E. Càmara, J. M. Toro, J. Marco-Pallares, N. Sebastián-Galles, H. Ortiz, J. Pujol and A. Rodríguez-
1406 Fornells (2009). "Time course and functional neuroanatomy of speech segmentation in adults." Neuroimage **48**(3):
1407 541-553.
- 1408 Cutler, A. and D. Norris (1988). "The role of strong syllables in segmentation for lexical access." Journal of
1409 Experimental Psychology: Human perception and performance **14**(1): 113.
- 1410 Davis, M. H. and I. S. Johnsrude (2003). "Hierarchical processing in spoken language comprehension." Journal of
1411 Neuroscience **23**(8): 3423-3431.
- 1412 Davis, M. H., W. D. Marslen-Wilson and M. G. Gaskell (2002). "Leading up the lexical garden path: Segmentation
1413 and ambiguity in spoken word recognition." Journal of Experimental Psychology: Human Perception and Performance
1414 **28**(1): 218.
- 1415 Davison, A. J. and D. W. Murray (2002). "Simultaneous localization and map-building using active vision." Ieee
1416 Transactions on Pattern Analysis and Machine Intelligence **24**(7): 865-880.
- 1417 Dehaene-Lambertz, G. (1997). "Electrophysiological correlates of categorical phoneme perception in adults." NeuroReport
1418 **8**(4): 919-924.
- 1419 DeWitt, I. and J. P. Rauschecker (2012). "Phoneme and word recognition in the auditory ventral stream." Proceedings
1420 of the National Academy of Sciences of the United States of America **109**(8): E505-E514.
- 1421 Ding, N., L. Melloni, H. Zhang, X. Tian and D. Poeppel (2015). "Cortical tracking of hierarchical linguistic structures
1422 in connected speech." Nature Neuroscience **19**(1): 158--164.
- 1423 Ding, N. and J. Z. Simon (2012). "Neural coding of continuous speech in auditory cortex during monaural and dichotic
1424 listening." Journal of neurophysiology **107**(1): 78--89.
- 1425 Donchin, E. and M. G. H. Coles (1988). "Is the P300 component a manifestation of context updating?" Behavioral
1426 and Brain Sciences **11**(3): 357.
- 1427 Drullman, R. (1995). "Temporal envelope and fine structure cues for speech intelligibility." Journal of the Acoustical
1428 Society of America **97**(1): 585-592.
- 1429 Dubno, J. R., J. B. Ahlstrom and a. R. Horwitz (2000). "Use of context by young and aged adults with normal hearing." The
1430 Journal of the Acoustical Society of America **107**(1): 538--546.
- 1431 Durlach, N. (2006). "Auditory masking: Need for improved conceptual structure." The Journal of the Acoustical
1432 Society of America **120**(4): 1787-1790.
- 1433 Durlach, N. I., C. R. Mason, G. K. Jr., T. L. Arbogast, H. S. Colburn and B. G. Shinn-Cunningham (2003). "Note on
1434 informational masking (L)." The Journal of the Acoustical Society of America **113**(6): 2984-2987.
- 1435 Durlach, N. I., C. R. Mason, B. G. Shinn-Cunningham, T. L. Arbogast, H. S. Colburn and G. Kidd (2003).
1436 "Informational masking: Counteracting the effects of stimulus uncertainty by decreasing target-masker similarity."
1437 The Journal of the Acoustical Society of America **114**(1): 368.
- 1438 Easwar, V., D. W. Purcell, S. J. Aiken, V. Parsa and S. D. Scollie (2015). "Evaluation of Speech-Evoked Envelope
1439 Following Responses as an Objective Aided Outcome Measure: Effect of Stimulus Level, Bandwidth, and
1440 Amplification in Adults With Hearing Loss." Ear Hear **36**(6): 635-652.

Active listening

- 1441 Feldman, A. G. and M. F. Levin (1995). "The origin and use of positional frames of reference in motor control." Behav
1442 Brain Sci. **18**: 723-806.
- 1443 Feynman, R. P. (1972). Statistical mechanics. Reading MA, Benjamin.
- 1444 Fiez, J. A. and S. E. Petersen (1998). "Neuroimaging studies of word reading." Proc Natl Acad Sci U S A **95**(3): 914-
1445 921.
- 1446 Fitch, W. T. (1997). "Vocal tract length and formant frequency dispersion correlate with body size in rhesus
1447 macaques." The Journal of the Acoustical Society of America **102**(2): 1213-1222.
- 1448 Forney, G. D. (2001). "Codes on graphs: Normal realizations." IEEE Transactions on Information Theory **47**(2): 520-
1449 548.
- 1450 François, C., T. Cunillera, E. Garcia, M. Laine and A. Rodriguez-Fornells (2017). "Neurophysiological evidence for
1451 the interplay of speech segmentation and word-referent mapping during novel word learning." Neuropsychologia **98**:
1452 56-67.
- 1453 Friederici, A. D., A. Hahne and A. Mecklinger (1996). "Temporal structure of syntactic parsing: early and late event-
1454 related brain potential effects." Journal of Experimental Psychology: Learning, Memory, and Cognition **22**(5): 1219.
- 1455 Friston, K. (2013). "Life as we know it." J R Soc Interface **10**(86): 20130475.
- 1456 Friston, K. and G. Buzsaki (2016). "The Functional Anatomy of Time: What and When in the Brain." Trends Cogn
1457 Sci.
- 1458 Friston, K., T. FitzGerald, F. Rigoli, P. Schwartenbeck and G. Pezzulo (2017). "Active Inference: A Process Theory."
1459 Neural Comput **29**(1): 1-49.
- 1460 Friston, K. and C. Frith (2015). "A duet for one." Consciousness and cognition **36**: 390-405.
- 1461 Friston, K., J. Mattout and J. Kilner (2011). "Action understanding and active inference." Biol Cybern. **104**: 137-160.
- 1462 Friston, K. J. (2010). "The free-energy principle: A unified brain theory?" Nature Reviews Neuroscience **11**(2): 127-
1463 -138.
- 1464 Friston, K. J., T. FitzGerald, F. Rigoli, P. Schwartenbeck and G. Pezzulo (2017). "Active Inference: A Process
1465 Theory." Neural computation **29**: 1--49.
- 1466 Friston, K. J., T. Parr and B. de Vries (2017). "The graphical brain: belief propagation and active inference." Network
1467 Neuroscience: 1--78.
- 1468 Friston, K. J., T. Parr and B. de Vries (2017). "The graphical brain: Belief propagation and active inference." Netw
1469 Neurosci **1**(4): 381-414.
- 1470 Friston, K. J., R. Rosch, T. Parr, C. Price and H. Bowman (2017). "Deep temporal models and active inference."
1471 Neurosci Biobehav Rev **77**: 388-402.
- 1472 Ganong, W. F. (1980). "Phonetic categorization in auditory word perception." Journal of experimental psychology:
1473 Human perception and performance **6**(1): 110.
- 1474 Garrido, M. I., J. M. Kilner, K. E. Stephan and K. J. Friston (2009). "The mismatch negativity: a review of underlying
1475 mechanisms." Clin Neurophysiol **120**(3): 453-463.
- 1476 Gaskell, M. G. and W. D. Marslen-Wilson (1997). "Integrating form and meaning: A distributed model of speech
1477 perception." Language and cognitive Processes **12**(5-6): 613-656.
- 1478 Gaudrain, E., S. Li, V. S. Ban and R. D. Patterson (2009). "The role of glottal pulse rate and vocal tract length in the
1479 perception of speaker identity." Proceedings of the Annual Conference of the International Speech Communication
1480 Association, INTERSPEECH(January 2009): 148--151.
- 1481 Giard, M., J. Lavikahen, K. Reinikainen, F. Perrin, O. Bertrand, J. Pernier and R. Näätänen (1995). "Separate
1482 representation of stimulus frequency, intensity, and duration in auditory sensory memory: an event-related potential
1483 and dipole-model analysis." Journal of cognitive neuroscience **7**(2): 133-143.
- 1484 Gow Jr, D. W. and P. C. Gordon (1995). "Lexical and prelexical influences on word segmentation: Evidence from
1485 priming." Journal of Experimental Psychology: Human perception and performance **21**(2): 344.
- 1486 Grossberg, S., K. Roberts, M. Aguilar and D. Bullock (1997). "A neural model of multimodal adaptive saccadic eye
1487 movement control by superior colliculus." J Neurosci. **17**(24): 9706-9725.
- 1488 Grotheer, M. and G. Kovács (2014). "Repetition probability effects depend on prior experiences." The Journal of
1489 neuroscience : the official journal of the Society for Neuroscience **34** **19**: 6640-6646.
- 1490 Guenther, F. H. and T. Vladusich (2012). "A Neural Theory of Speech Acquisition and Production." J Neurolinguistics
1491 **25**(5): 408-422.
- 1492 Harris, M. and N. Umeda (1974). "Effect of speaking mode on temporal factors in speech: Vowel duration." The
1493 Journal of the Acoustical Society of America **56**(3): 1016-1018.
- 1494 Hasson, U., E. Yang, I. Vallines, D. J. Heeger and N. Rubin (2008). "A hierarchy of temporal receptive windows in

Active listening

- 1495 human cortex." *J Neurosci* **28**(10): 2539-2550.
- 1496 Heilbron, M. and M. Chait (2018). "Great Expectations: Is there Evidence for Predictive Coding in Auditory Cortex?"
- 1497 *Neuroscience* **389**: 54-73.
- 1498 Hickok, G. (2014). "The architecture of speech production and the role of the phoneme in speech processing." *Lang*
- 1499 *Cogn Process* **29**(1): 2-20.
- 1500 Hickok, G. and D. Poeppel (2007). "Opinion - The cortical organization of speech processing." *Nature Reviews*
- 1501 *Neuroscience* **8**(5): 393-402.
- 1502 Hillenbrand, J. M., L. A. Getty, M. J. Clark and K. Wheeler (1995). "Acoustic characteristics of American English
- 1503 vowels." *Journal of the Acoustical Society of America* **97**(5): 3099--3111.
- 1504 Hinton, G. E. and R. S. Zemel (1993). Autoencoders, minimum description length and Helmholtz free energy.
- 1505 *Proceedings of the 6th International Conference on Neural Information Processing Systems*. Denver, Colorado,
- 1506 Morgan Kaufmann Publishers Inc.: 3-10.
- 1507 Hodges, J. R., K. Patterson, S. Oxbury and E. Funnell (1992). "Semantic dementia. Progressive fluent aphasia with
- 1508 temporal lobe atrophy." *Brain* **115 (Pt 6)**: 1783-1806.
- 1509 Hohwy, J. (2016). "The Self-Evidencing Brain." *Noûs* **50**(2): 259-285.
- 1510 Holmes, E., Y. Domingo and I. S. Johnsrude (2018). "Familiar voices are more intelligible, even if they are not
- 1511 recognized as familiar." *Psychological Science* **29**(10): 1575--1583.
- 1512 Holmes, E., P. Folkeard, I. S. Johnsrude and S. Scollie (2018). "Semantic context improves speech intelligibility and
- 1513 reduces listening effort for listeners with hearing impairment." *Int J Audiol* **57**(7): 483-492.
- 1514 Holt, L. L., A. J. Lotto and K. R. Kluender (2000). "Neighboring spectral content influences vowel identification."
- 1515 *Journal of the Acoustical Society of America* **108**(2): 710-722.
- 1516 Hope, T. M. H., A. P. Leff and C. J. Price (2018). "Predicting language outcomes after stroke: Is structural
- 1517 disconnection a useful predictor?" *NeuroImage. Clinical* **19**: 22-29.
- 1518 Houde, J. and S. Nagarajan (2011). "Speech Production as State Feedback Control." *Frontiers in Human Neuroscience*
- 1519 **5**(82).
- 1520 Itti, L. and P. Baldi (2009). "Bayesian Surprise Attracts Human Attention." *Vision Res.* **49**(10): 1295-1306.
- 1521 Jacobsen, T., T. Horenkamp and E. Schröger (2003). "Preattentive memory-based comparison of sound intensity."
- 1522 *Audiology and Neurotology* **8**(6): 338-346.
- 1523 Jacobsen, T., E. Schröger, T. Horenkamp and I. Winkler (2003). "Mismatch negativity to pitch change: varied stimulus
- 1524 proportions in controlling effects of neural refractoriness on human auditory event-related brain potentials."
- 1525 *Neuroscience letters* **344**(2): 79-82.
- 1526 Johnsrude, I. S., A. Mackey, H. Hakyemez, E. Alexander, H. P. Trang and R. P. Carlyon (2013). "Swinging at a
- 1527 cocktail party: voice familiarity aids speech perception in the presence of a competing voice." *Psychological science*
- 1528 **24**(10): 1995--2004.
- 1529 Kaas, J. H. and T. A. Hackett (1999). "'What' and 'where' processing in auditory cortex." *Nat Neurosci* **2**(12): 1045--
- 1530 1047.
- 1531 Kidd, G., C. R. Mason, V. M. Richards, F. Gallun and N. Durlach (2007). Informational Masking. **29**: 143-189.
- 1532 Kiebel, S. J., J. Daunizeau and K. J. Friston (2009). "Perception and hierarchical dynamics." *Front Neuroinform* **3**:
- 1533 20.
- 1534 Kim, D., J. D. Stephens and M. A. Pitt (2012). "How does context play a part in splitting words apart? Production and
- 1535 perception of word boundaries in casual speech." *Journal of memory and language* **66**(4): 509-529.
- 1536 Kim, S., R. D. Frisina, F. M. Mapes, E. D. Hickman and D. R. Frisina (2006). "Effect of age on binaural speech
- 1537 intelligibility in normal hearing adults." *Speech Communication* **48**(6): 591--597.
- 1538 Klatt, D. H. (1975). "Vowel lengthening is syntactically determined in a connected discourse." *Journal of phonetics*
- 1539 **3**(3): 129-140.
- 1540 Klatt, D. H. (1976). "Linguistic uses of segmental duration in English: Acoustic and perceptual evidence." *The Journal*
- 1541 *of the Acoustical Society of America* **59**(5): 1208-1221.
- 1542 Kleinschmidt, D. F. and T. F. Jaeger (2015). "Robust Speech Perception: Recognize the Familiar, Generalize to the
- 1543 Similar, and Adapt to the Novel." *Psychological Review* **122**(2): 148-203.
- 1544 Kuhlen, A. K., C. Bogler, S. E. Brennan and J.-D. Haynes (2017). "Brains in dialogue: decoding neural preparation
- 1545 of speaking to a conversational partner." *Social cognitive and affective neuroscience* **12**(6): 871-880.
- 1546 Kumar, S., K. E. Stephan, J. D. Warren, K. J. Friston and T. D. Griffiths (2007). "Hierarchical processing of auditory
- 1547 objects in humans." *PLoS computational biology* **3**(6): e100.
- 1548 Kuperberg, G. R., T. Sitnikova, D. Caplan and P. J. Holcomb (2003). "Electrophysiological distinctions in processing

Active listening

- 1549 conceptual relationships within simple sentences." *Cognitive brain research* **17**(1): 117-129.
- 1550 Kutas, M. and K. D. Federmeier (2000). "Electrophysiology reveals semantic memory use in language
1551 comprehension." *Trends in cognitive sciences* **4**(12): 463-470.
- 1552 Kutas, M. and K. D. Federmeier (2009). "N400." *Scholarpedia* **4**(10): 7790.
- 1553 Kutas, M. and S. A. Hillyard (1980). "Reading senseless sentences: Brain potentials reflect semantic incongruity."
1554 *Science* **207**(4427): 203-205.
- 1555 Kutas, M. and S. A. Hillyard (1984). "Brain potentials during reading reflect word expectancy and semantic
1556 association." *Nature* **307**(5947): 161.
- 1557 Ladd, D. R. and A. Schepman (2003). "'Sagging transitions" between high pitch accents in English: Experimental
1558 evidence." *Journal of phonetics* **31**(1): 81-112.
- 1559 Landauer, R. (1961). "Irreversibility and Heat Generation in the Computing Process." *IBM Journal of Research and
1560 Development* **5**(3): 183-191.
- 1561 Landi, N., S. J. Frost, W. E. Menci, R. Sandak and K. R. Pugh (2013). "Neurobiological bases of reading
1562 comprehension: Insights from neuroimaging studies of word level and text level processing in skilled and impaired
1563 readers." *Read Writ Q* **29**(2): 145-167.
- 1564 LaRivière, C. (1975). "Contributions of Fundamental Frequency and Formant Frequencies to Speaker Identification."
1565 *Phonetica* **31**(3-4): 185-197.
- 1566 Larsson, J. and A. T. Smith (2012). "fMRI repetition suppression: neuronal adaptation or stimulus expectation?" *Cereb
1567 Cortex* **22**(3): 567-576.
- 1568 Lavner, Y., I. Gath and J. Rosenhouse (2000). "Effects of acoustic modifications on the identification of familiar
1569 voices speaking isolated vowels." *Speech Communication* **30**(1): 9--26.
- 1570 Lavner, Y., J. Rosenhouse and I. Gath (2001). "The prototype model in speaker identification by human listeners."
1571 *International Journal of Speech Technology* **4**(1): 63--74.
- 1572 Lehiste, I. (1960). "An acoustic-phonetic study of internal open juncture." *Phonetica* **5**(Suppl. 1): 5-54.
- 1573 Lehiste, I. (1972). "The timing of utterances and linguistic boundaries." *The Journal of the Acoustical Society of
1574 America* **51**(6B): 2018-2024.
- 1575 Lehiste, I. (1973). "Rhythmic units and syntactic units in production and perception." *The Journal of the Acoustical
1576 Society of America* **54**(5): 1228-1234.
- 1577 Liberman, A. M., F. S. Cooper, D. P. Shankweiler and M. Studdert-Kennedy (1967). "Perception of the speech code."
1578 *Psychological review* **74**(6): 431.
- 1579 Luce, P. A. (1986). "Neighborhoods of words in the mental lexicon." *Research on speech perception, Technical Report
1580 6*: 1-91.
- 1581 Luce, P. A. and D. B. Pisoni (1998). "Recognizing spoken words: the neighborhood activation model." *Ear and hearing
1582 19*(1): 1-36.
- 1583 Maisto, D., F. Donnarumma and G. Pezzulo (2015). "Divide et impera: subgoaling reduces the complexity of
1584 probabilistic inference and problem solving." *12*(104): 20141335.
- 1585 Mann, V. A. (1980). "Influence of preceding liquid on stop-consonant perception." *Perception & Psychophysics* **28**(5):
1586 407-412.
- 1587 Marslen-Wilson, W. D. (1975). "Sentence perception as an interactive parallel process." *Science* **189**(4198): 226-228.
- 1588 Marslen-Wilson, W. D. (1984). Function and process in spoken word recognition: A tutorial review. *Attention and
1589 performance: Control of language processes*, Erlbaum: 125-150.
- 1590 Marslen-Wilson, W. D. and A. Welsh (1978). "Processing interactions and lexical access during word recognition in
1591 continuous speech." *Cognitive psychology* **10**(1): 29-63.
- 1592 Massaro, D. W. (1987). Categorical partition: A fuzzy-logical model of categorization behavior. *Categorical
1593 perception: The groundwork of cognition*. New York, NY, US, Cambridge University Press: 254-283.
- 1594 Massaro, D. W. (1989). "Testing between the TRACE model and the fuzzy logical model of speech perception."
1595 *Cognitive psychology* **21**(3): 398-421.
- 1596 Matsumoto, H., S. Hiki, T. Sone and T. Nimura (1973). "Multidimensional representation of personal quality of
1597 vowels and its acoustical correlates." *IEEE Transactions on Audio and Electroacoustics* **21**(5): 428--436.
- 1598 Mattys, S. L. and J. F. Melhorn (2007). "Sentential, lexical, and acoustic effects on the perception of word boundaries."
1599 *The Journal of the Acoustical Society of America* **122**(1): 554-567.
- 1600 Mattys, S. L., J. F. Melhorn and L. White (2007). "Effects of syntactic expectations on speech segmentation." *Journal
1601 of Experimental Psychology: Human Perception and Performance* **33**(4): 960.
- 1602 Mattys, S. L., L. White and J. F. Melhorn (2005). "Integration of multiple speech segmentation cues: A hierarchical

Active listening

- 1603 framework." *Journal of Experimental Psychology-General* **134**(4): 477-500.
- 1604 McClelland, J. L. and J. L. Elman (1986). "The TRACE model of speech perception." *Cognitive Psychology* **18**(1):
- 1605 1-86.
- 1606 Mermelstein, P. (1967). "Determination of the Vocal-Tract Shape from Measured Formant Frequencies." *The Journal*
- 1607 *of the Acoustical Society of America* **41**(5): 1283-1294.
- 1608 Miller, J. L., K. Green and T. M. Schermer (1984). "A distinction between the effects of sentential speaking rate and
- 1609 semantic congruity on word identification." *Perception & Psychophysics* **36**(4): 329-337.
- 1610 Miller, J. L. and A. M. Liberman (1979). "Some effects of later-occurring information on the perception of stop
- 1611 consonant and semivowel." *Perception & Psychophysics* **25**(6): 457-465.
- 1612 Mirza, M. B., R. A. Adams, C. D. Mathys and K. J. Friston (2016). "Scene Construction, Visual Foraging, and Active
- 1613 Inference." *Frontiers in Computational Neuroscience* **10**(56).
- 1614 Mirza, M. B., R. A. Adams, C. D. Mathys and K. J. Friston (2016). "Scene Construction, Visual Foraging, and Active
- 1615 Inference." *Front Comput Neurosci* **10**: 56.
- 1616 Mohan, V. and P. Morasso (2011). "Passive motion paradigm: an alternative to optimal control." *Front Neurobot* **5**:
- 1617 4.
- 1618 Morlet, D. and C. Fischer (2014). "MMN and novelty P3 in coma and other altered states of consciousness: a review." *Brain Topogr*
- 1619 **27**(4): 467-479.
- 1620 Muralimanohar, R. K., J. M. Kates and K. H. Arehart (2017). "Using envelope modulation to explain speech
- 1621 intelligibility in the presence of a single reflection." *J Acoust Soc Am* **141**(5): E1482.
- 1622 Murry, T. and S. Singh (1980). "Multidimensional analysis of male and female voices." *The Journal of the Acoustical*
- 1623 *Society of America* **68**(5): 1294--1300.
- 1624 Musso, M., C. Weiller, A. Horn, V. Glauche, R. Umarova, J. Hennig, A. Schneider and M. Rijntjes (2015). "A single
- 1625 dual-stream framework for syntactic computations in music and language." *Neuroimage* **117**: 267-283.
- 1626 Näätänen, R., A. W. Gaillard and S. Mäntysalo (1978). "Early selective-attention effect on evoked potential
- 1627 reinterpreted." *Acta psychologica* **42**(4): 313-329.
- 1628 Näätänen, R., A. Lehtokoski, M. Lennes, M. Cheour, M. Huotilainen, A. Iivonen, M. Vainio, P. Alku, R. J. Ilmoniemi
- 1629 and A. Luuk (1997). "Language-specific phoneme representations revealed by electric and magnetic brain responses." *Nature*
- 1630 **385**(6615): 432.
- 1631 Nakatani, L. H. and K. D. Dukes (1977). "Locus of segmental cues for word juncture." *The Journal of the Acoustical*
- 1632 *Society of America* **62**(3): 714-719.
- 1633 Nardo, D., R. Holland, A. P. Leff, C. J. Price and J. T. Crinion (2017). "Less is more: neural mechanisms underlying
- 1634 anomia treatment in chronic aphasic patients." *Brain* **140**(11): 3039-3054.
- 1635 Nealey, T. A. and J. H. Maunsell (1994). "Magnocellular and parvocellular contributions to the responses of neurons
- 1636 in macaque striate cortex." *The Journal of Neuroscience* **14**(4): 2069.
- 1637 Norris, D. and J. M. McQueen (2008). "Shortlist B: A Bayesian model of continuous speech recognition." *Psychological review*
- 1638 **115**(2): 357--395.
- 1639 Norris, D., J. M. McQueen and A. Cutler (2016). "Prediction, Bayesian inference and feedback in speech recognition." *Lang Cogn Neurosci*
- 1640 **31**(1): 4-18.
- 1641 Norris, D., J. M. McQueen, A. Cutler and S. Butterfield (1997). "The possible-word constraint in the segmentation of
- 1642 continuous speech." *Cognitive Psychology* **34**(3): 191-243.
- 1643 Nygaard, L. C., M. S. Sommers and D. B. Pisoni (1994). "SPEECH PERCEPTION AS A TALKER-CONTINGENT
- 1644 PROCESS." *Psychol Sci* **5**(1): 42-46.
- 1645 O'Leary, D. D. M. (1989). "Do cortical areas emerge from a protocortex?" *Trends in Neurosciences* **12**(10): 400-406.
- 1646 O'Sullivan, J. A., A. J. Power, N. Mesgarani, S. Rajaram, J. J. Foxe, B. G. Shinn-Cunningham, M. Slaney, S. a.
- 1647 Shamma and E. Lalor (2014). "Attentional selection in a cocktail party environment can be decoded from single-trial
- 1648 EEG." *Cerebral Cortex*: 1--10.
- 1649 Oden, G. C. and D. W. Massaro (1978). "Integration of featural information in speech perception." *Psychological*
- 1650 *review* **85**(3): 172.
- 1651 Ognibene, D. and G. Baldassarre (2014). Ecological Active Vision: Four Bio-Inspired Principles to Integrate Bottom-
- 1652 Up and Adaptive Top-Down Attention Tested With a Simple Camera-Arm Robot. *IEEE Transactions on Autonomous*
- 1653 *Mental Development, IEEE*.
- 1654 Oller, D. K. (1973). "The effect of position in utterance on speech segment duration in English." *The journal of the*
- 1655 *Acoustical Society of America* **54**(5): 1235-1247.

Active listening

- 1656 Osterhout, L. and P. J. Holcomb (1992). "Event-related brain potentials elicited by syntactic anomaly." Journal of
1657 memory and language **31**(6): 785-806.
- 1658 Oudeyer, P.-Y. and F. Kaplan (2007). "What is intrinsic motivation? a typology of computational approaches."
1659 Frontiers in Neurobotics **1**: 6.
- 1660 Pannekamp, A., U. Toepel, K. Alter, A. Hahne and A. D. Friederici (2005). "Prosody-driven sentence processing: An
1661 event-related brain potential study." Journal of cognitive neuroscience **17**(3): 407-421.
- 1662 Parr, T. and K. J. Friston (2017). "The active construction of the visual world." Neuropsychologia **104**: 92-101.
- 1663 Parr, T. and K. J. Friston (2017). "Working memory, attention, and salience in active inference." Scientific Reports
1664 **7**(1): 14678.
- 1665 Parr, T., D. Markovic, S. J. Kiebel and K. J. Friston (2019). "Neuronal message passing using Mean-field, Bethe, and
1666 Marginal approximations." Scientific Reports **9**(1): 1889.
- 1667 Pasley, B. N., S. V. David, N. Mesgarani, A. Flinker, S. A. Shamma, N. E. Crone, R. T. Knight and E. F. Chang
1668 (2012). "Reconstructing speech from human auditory cortex." PLoS biology **10**(1): e1001251.
- 1669 Patel, A. D. (2010). Music, language, and the brain. Oxford, UK, Oxford Univ. Press.
- 1670 Paulesu, E., B. Goldacre, P. Scifo, S. F. Cappa, M. C. Gilardi, I. Castiglioni, D. Perani and F. Fazio (1997). "Functional
1671 heterogeneity of left inferior frontal cortex as revealed by fMRI." Neuroreport **8**(8): 2011-2017.
- 1672 Pearce, M. T. (2018). "Statistical learning and probabilistic prediction in music cognition: mechanisms of stylistic
1673 enculturation." Ann N Y Acad Sci.
- 1674 Penny, W. D. (2012). "Comparing dynamic causal models using AIC, BIC and free energy." Neuroimage **59**(1): 319-
1675 330.
- 1676 Peretz, I., R. Kolinsky, M. Tramo, R. Labrecque, C. Hublet, G. Demeurisse and S. Belleville (1994). "Functional
1677 dissociations following bilateral lesions of auditory cortex." Brain **117**(6): 1283-1301.
- 1678 Picton, T. W., C. Alain, L. Otten, W. Ritter and A. Achim (2000). "Mismatch negativity: different water in the same
1679 river." Audiology and Neurotology **5**(3-4): 111-139.
- 1680 Poeppel, D. and P. J. Monahan (2011). "Feedforward and feedback in speech perception: Revisiting analysis by
1681 synthesis." Language and Cognitive Processes **26**(7): 935-951.
- 1682 Polich, J. (2007). "Updating P300: an integrative theory of P3a and P3b." Clinical neurophysiology **118**(10): 2128-
1683 2148.
- 1684 Polich, J. and E. Donchin (1988). "P300 and the word frequency effect." Electroencephalography and clinical
1685 neurophysiology **70**(1): 33-45.
- 1686 Price, C. J. (2012). "A review and synthesis of the first 20 years of PET and fMRI studies of heard speech, spoken
1687 language and reading." NeuroImage **62**(2): 816-847.
- 1688 Quiroga-Martinez, D. R., N. C. Hansen, A. Højlund, M. Pearce, E. Brattico and P. Vuust (2019). "Reduced prediction
1689 error responses in high- as compared to low-uncertainty musical contexts." bioRxiv: 422949.
- 1690 Remez, R. E. (2010). "Spoken expression of individual identity and the listener." Expressing oneself/expressing one's
1691 self: Communication, cognition, language, and identity: 167--181.
- 1692 Romanski, L. M., B. Tian, J. Fritz, M. Mishkin, P. S. Goldman-Rakic and J. P. Rauschecker (1999). "Dual streams of
1693 auditory afferents target multiple domains in the primate prefrontal cortex." Nat Neurosci **2**(12): 1131-1136.
- 1694 Rosenfeld, R. (2000). "Two decades of statistical language modeling: Where do we go from here?" Proceedings of
1695 the Ieee **88**(8): 1270-1278.
- 1696 Rueschemeyer, S.-A., M. G. Gaskell, G. Walker and G. Hickok (2018). Speech Production Integrating
1697 psycholinguistic, neuroscience, and motor control perspectives, Oxford University Press.
- 1698 Ryan, R. and E. Deci (1985). Intrinsic motivation and self-determination in human behavior. New York, Plenum.
- 1699 Sams, M., P. Paavilainen and K. Alho (1985). "Auditory frequency discrimination and event-related potentials."
1700 Electroencephalography and Clinical Neurophysiology **62**: 437--448.
- 1701 Sato, Y., H. Yabe, T. Hiruma, T. Sutoh, N. Shinozaki, T. Nashida and S. Kaneko (2000). "The effect of deviant
1702 stimulus probability on the human mismatch process." Neuroreport **11**(17): 3703-3708.
- 1703 Sato, Y., H. Yabe, J. Todd, P. Michie, N. Shinozaki, T. Sutoh, T. Hiruma, T. Nashida, T. Matsuoaka and S. Kaneko
1704 (2003). "Impairment in activation of a frontal attention-switch mechanism in schizophrenic patients." Biological
1705 psychology **62**(1): 49-63.
- 1706 Schmidhuber, J. (1991). "Curious model-building control systems." In Proc. International Joint Conference on Neural
1707 Networks, Singapore. IEEE **2**: 1458-1463.
- 1708 Schmidhuber, J. (2006). "Developmental robotics, optimal artificial curiosity, creativity, music, and the fine arts."
1709 Connection Science **18**(2): 173-187.

Active listening

- 1710 Sengupta, B., M. B. Stemmler and K. J. Friston (2013). "Information and efficiency in the nervous system—a
1711 synthesis." PLoS computational biology **9**(7): e1003157.
- 1712 Sengupta, B., A. Tozzi, G. K. Cooray, P. K. Douglas and K. J. Friston (2016). "Towards a Neuronal Gauge Theory."
1713 PLoS Biol **14**(3): e1002400.
- 1714 Seth, A. (2014). The cybernetic brain: from interoceptive inference to sensorimotor contingencies. MINDS project.
1715 Metzinger, T; Windt, JM, MINDS.
- 1716 Shamma, S. (2001). "On the role of space and time in auditory processing." Trends in cognitive sciences **5**(8): 340-
1717 348.
- 1718 Shamma, S. A., M. Elhilali and C. Micheyl (2011). "Temporal coherence and attention in auditory scene analysis."
1719 Trends in neurosciences **34**(3): 114--123.
- 1720 Shiell, M. M., F. Champoux and R. J. Zatorre (2015). "Reorganization of auditory cortex in early-deaf people:
1721 Functional connectivity and relationship to hearing aid use." Journal of Cognitive Neuroscience **27**(1): 150-163.
- 1722 Shillcock, R. (1990). "Lexical hypotheses in continuous speech."
1723 Steinhauer, K., K. Alter and A. D. Friederici (1999). "Brain potentials indicate immediate use of prosodic cues in
1724 natural speech processing." Nature neuroscience **2**(2): 191.
- 1725 Sun, Y., F. Gomez and J. Schmidhuber (2011). Planning to Be Surprised: Optimal Bayesian Exploration in Dynamic
1726 Environments. Artificial General Intelligence: 4th International Conference, AGI 2011, Mountain View, CA, USA,
1727 August 3-6, 2011. Proceedings. J. Schmidhuber, K. R. Thórisson and M. Looks. Berlin, Heidelberg, Springer Berlin
1728 Heidelberg: 41-51.
- 1729 Sur, M., P. E. Garraghty and A. W. Roe (1988). "Experimentally induced visual projections into auditory thalamus
1730 and cortex." Science **242**(4884): 1437-1441.
- 1731 Taylor, J. S., K. Rastle and M. H. Davis (2013). "Can cognitive models explain brain activation during word and
1732 pseudoword reading? A meta-analysis of 36 neuroimaging studies." Psychol Bull **139**(4): 766-791.
- 1733 Tervaniemi, M., T. Ilvonen, K. Karma, K. Alho and R. Näätänen (1997). "The musical brain: brain waves reveal the
1734 neurophysiological basis of musicality in human subjects." Neuroscience letters **226**(1): 1-4.
- 1735 Tervaniemi, M., I. Winkler and R. Näätänen (1997). "Pre-attentive categorization of sounds by timbre as revealed by
1736 event-related potentials." NeuroReport **8**(11): 2571-2574.
- 1737 Thiel, A., B. Habedank, L. Winhuisen, K. Herholz, J. Kessler, W. F. Haupt and W. D. Heiss (2005). "Essential
1738 language function of the right hemisphere in brain tumor patients." Ann Neurol **57**(1): 128-131.
- 1739 Thiessen, E. and L. Erickson (2013). "Discovering Words in Fluent Speech: The Contribution of Two Kinds of
1740 Statistical Information." Frontiers in Psychology **3**(590).
- 1741 Toiviainen, P., M. Tervaniemi, J. Louhivuori, M. Saher, M. Huotilainen and R. Näätänen (1998). "Timbre similarity:
1742 Convergence of neural, behavioral, and computational approaches." Music Perception: An Interdisciplinary Journal
1743 **16**(2): 223-241.
- 1744 Tourville, J. A. and F. H. Guenther (2011). "The DIVA model: A neural theory of speech acquisition and production."
1745 Lang Cogn Process **26**(7): 952-981.
- 1746 Ueno, T., S. Saito, T. T. Rogers and M. A. Lambon Ralph (2011). "Lichtheim 2: synthesizing aphasia and the neural
1747 basis of language in a neurocomputational model of the dual dorsal-ventral language pathways." Neuron **72**(2): 385-
1748 396.
- 1749 Ulanovsky, N. and C. F. Moss (2008). "What the bat's voice tells the bat's brain." Proceedings of the National Academy
1750 of Sciences of the United States of America **105**(25): 8491-8498.
- 1751 Ungerleider, L. G. and J. V. Haxby (1994). "'What' and 'where' in the human brain." Current Opinion in Neurobiology
1752 **4**(2): 157-165.
- 1753 Van Dommelen, W. A. (1987). "The Contribution of Speech Rhythm and Pitch to Speaker Recognition." Language
1754 and Speech **30**(4): 325-338.
- 1755 Van Dommelen, W. A. (1990). "Acoustic parameters in human speaker recognition." Language and Speech **33**(3):
1756 259-272.
- 1757 Van Petten, C., S. Coulson, S. Rubin, E. Plante and M. Parks (1999). "Time course of word identification and semantic
1758 integration in spoken language." Journal of Experimental Psychology: Learning, Memory, and Cognition **25**(2): 394.
- 1759 Van Petten, C. and M. Kutas (1990). "Interactions between sentence context and word frequency in event-related
1760 brain potentials." Memory & cognition **18**(4): 380-393.
- 1761 Vanthornhout, J., L. Decruy, J. Wouters, J. Simon and T. Francart (2018). "Speech intelligibility predicted from neural
1762 entrainment of the speech envelope." bioRxiv(637424): 246660.
- 1763 Veale, R., Z. M. Hafed and M. Yoshida (2017). "How is visual salience computed in the brain? Insights from

Active listening

- 1764 behaviour, neurobiology and modelling." **372**(1714).
- 1765 Vinckier, F., S. Dehaene, A. Jobert, J. P. Dubus, M. Sigman and L. Cohen (2007). "Hierarchical coding of letter strings
- 1766 in the ventral stream: Dissecting the inner organization of the visual word-form system." Neuron **55**(1): 143-156.
- 1767 Wacongne, C., J. P. Changeux and S. Dehaene (2012). "A neuronal model of predictive coding accounting for the
- 1768 mismatch negativity." J Neurosci **32**(11): 3665-3678.
- 1769 Walden, B. E., A. A. Montgomery, G. J. Gibeily, R. A. Prosek and D. M. Schwartz (1978). "Correlates of
- 1770 psychological dimensions in talker similarity." Journal of speech, language, and hearing research **21**: 265--275.
- 1771 Warburton, E., C. J. Price, K. Swinburn and R. J. S. Wise (1999). "Mechanisms of recovery from aphasia: evidence
- 1772 from positron emission tomography studies." Journal of Neurology, Neurosurgery & Psychiatry **66**(2): 155-161.
- 1773 Winkler, I., S. L. Denham and I. Nelken (2009). "Modeling the auditory scene: predictive regularity representations
- 1774 and perceptual objects." Trends in Cognitive Sciences **13**(12): 532--540.
- 1775 Winn, J. and C. M. Bishop (2005). "Variational message passing." Journal of Machine Learning Research **6**: 661-694.
- 1776 Ylinen, S., M. Huuskonen, K. Mikkola, E. Saure, T. Sinkkonen and P. Paavilainen (2016). "Predictive coding of
- 1777 phonological rules in auditory cortex: A mismatch negativity study." Brain Lang **162**: 72-80.
- 1778 Zeki, S. and S. Shipp (1988). "The functional logic of cortical connections." Nature **335**: 311-317.
- 1779 Zhang, C., J. Butepage, H. Kjellstrom and S. Mandt (2018). "Advances in Variational Inference." IEEE Trans Pattern
- 1780 Anal Mach Intell.
- 1781