

## **In spoken word recognition the future predicts the past**

\*Laura Gwilliams<sup>1,2</sup>, Tal Linzen<sup>4,5</sup>, David Poeppel<sup>1,3</sup> & Alec Marantz<sup>1,2</sup>

*<sup>1</sup>New York University*

*<sup>2</sup>NYUAD Institute*

*<sup>3</sup>Max-Planck Institute*

*<sup>4</sup>LSCP & IJN, CNRS, EHESS and ENS, PSL Research University*

*<sup>5</sup>Johns Hopkins University*

\* Corresponding Author:

Laura Gwilliams  
10 Washington Place  
New York,  
NY, 10003  
United States

Email: [laura.gwilliams@nyu.edu](mailto:laura.gwilliams@nyu.edu)

Tel: 1-347-725-5635

## Abstract

1 Speech is an inherently noisy and ambiguous signal. In order to fluently derive meaning, a  
2 listener must integrate contextual information to guide interpretations of the sensory  
3 input. While many studies have demonstrated the influence of *prior* context, the neural  
4 mechanisms supporting the integration of *subsequent* information remain unknown.  
5 Using magnetoencephalography, we analysed responses to spoken words with a varyingly  
6 ambiguous onset phoneme, the identity of which is later disambiguated at the lexical  
7 uniqueness point<sup>1</sup>. Our results uncover a three-level processing network. Subphonemic  
8 detail is preserved in primary auditory cortex over long timescales, and re-evoked at  
9 subsequent phoneme positions. Commitments to phonological categories occur in  
10 parallel, resolving on the shorter time-scale of ~450 ms. Finally, predictions are formed  
11 over likely lexical items. These findings provide evidence that future input determines the  
12 perception of earlier speech sounds by maintaining sensory features until they can be  
13 optimally integrated with top-down information.

14  
Keywords: *Speech; MEG; Lexical access; Auditory processing*

---

<sup>1</sup> For a demonstration of this retroactive phenomenon, see [here](#)

1 Perceptual interpretation of transient signals presents a fascinating puzzle. Because the  
2 system cannot choose to re-receive earlier input, it must either prioritise *accuracy* by  
3 allowing sensory evidence to accumulate over time, or prioritise *speed* of processing by  
4 committing to an interpretation before it is certain that interpretation is correct. This  
5 issue becomes even more problematic when the signal is hierarchically structured, as is  
6 true for speech: Comprehension of lower-level units is required to perceive the whole (e.g.  
7 phonemes are assembled into words, words into sentences); understanding the whole aids  
8 comprehension of its constituents. Therefore, in order to derive the correct interpretation,  
9 the speech-comprehension system must balance integration over short and long time-  
10 scales, and allow low and high-order representations to interact.

11  
12 Phonemes (e.g. /b/, /p/, /t/), and their phonetic features (e.g. -voiced, +plosive), are  
13 considered the smallest meaningful units of speech. Previous research indicates that  
14 phonemes are first processed in terms of low-level spectrotemporal extraction in primary  
15 auditory cortex ~50 ms after onset (A1 / Heschl's gyrus), followed by analysis of phonetic  
16 feature representations in superior temporal gyrus (STG) ~100 ms (Simos, Diehl et al.  
17 1998, Ackermann, Lutzenberger et al. 1999, Obleser, Lahiri et al. 2003, Papanicolaou,  
18 Castillo et al. 2003, Obleser, Lahiri et al. 2004, Mesgarani, Cheung et al. 2014, Di Liberto,  
19 O'Sullivan et al. 2015). This is in line with models of auditory processing that assume  
20 increasing abstraction as information moves along the auditory pathways (Scott and  
21 Johnsrupe 2003, Hickok and Poeppel 2004, Liebenthal, Binder et al. 2005, Rauschecker and  
22 Scott 2009). In many cases an acoustic signal will be consistent with more than one  
23 phoneme, and the system will need to decide which categorisation is the correct one. But  
24 it is currently unknown where the recognition and resolution of phonological ambiguity  
25 fits relative to these stages of processing. While ambiguity has been shown to modulate  
26 fMRI-recorded activation in posterior Heschl's gyrus (Kilian-Hutten, Valente et al. 2011),  
27 suggesting there are abstract perceptual sound representations in primary auditory  
28 cortex, without timing information it cannot be established whether this region is  
29 recruited before, during or after phonological categorisation in higher level areas.

30  
31 The interpretation of phonemes is not a purely bottom-up process; surrounding context,  
32 in the form of words and sentences, also has an influence. So far, integration of these low  
33 (phonetic/phonological) and high (lexical/sentential) order representations over long  
34 time-scales has primarily been investigated in terms of the influence of preceding context  
35 on following phoneme perception. For example, a sound that is acoustically ambiguous  
36 between [s] and [f] will be perceived as /f/ in the word "gift" but as /s/ in the word "kiss".  
37 The Cohort model accounts for this phenomenon by suggesting that a phonological  
38 sequence activates word candidates, which help form predictions about upcoming  
39 phonemes (Marslen-Wilson and Tyler 1980, Tyler 1984, Marslen-Wilson 1987).  
40 Phonological perception is then biased to be consistent with the sound that was predicted  
41 given the context (Warren 1970, Cole 1973, Samuel 1981). Investigations into the neural  
42 correlates of this process associate the (dis-)confirmation of predictions with activation in  
43 the left superior and transverse temporal gyri (Gagnepain, Henson et al. 2012, Ettinger,  
44 Linzen et al. 2014, Gwilliams and Marantz 2015), suggesting that these regions compare  
45 the expected with the received information, and any discrepancy between the two is  
46 reflected as an error-prediction signal.

1 One defining characteristic of the Cohort model and other feedforward models of speech  
2 processing is that information never flows backwards in time. Consequently, because  
3 phonemes unfold more quickly than words, information contained within a word cannot  
4 feed back to alter analysis of its previously heard phonemes. However, relevant context  
5 often does occur after the sensory input to be interpreted, and still serves to bias  
6 perception. For example, a speech sound that is ambiguous between [g] and [k] will be  
7 interpreted as /g/ if occurring at the onset of “gift” and /k/ at the onset of “kiss”, even  
8 though the acoustic signal at the onset of each word is identical (Ganong 1980, McQueen  
9 1991, Gordon, Eberhardt et al. 1993). Functional connectivity analysis of MEG data  
10 supports the existence of feedback in such situations, finding evidence that activity linked  
11 to lexical processing in supra marginal gyrus affects phonetic processing in STG at a  
12 word’s point of disambiguation (POD) (Gow, Segawa et al. 2008). The STG has also been  
13 implicated in fMRI studies of phoneme ambiguity (Blumstein, Myers et al. 2005, Myers  
14 and Blumstein 2008), and with sensitivity to post-assimilation context (Gow and Segawa  
15 2009) (also see (Gow and McMurray 2007) for an overview of related behavioural results).  
16 The inferior frontal gyrus is also recruited when ambiguous phonemes distinguish word  
17 pairs (e.g. *blade vs. glade*), but not pseudo-words, or pseudo-word/word pairs (Rogers and  
18 Davis 2017), further suggesting there is a lexical influence in resolving phonological  
19 identity.

20  
21 If speech processing does not advance in a purely feedforward manner (though see  
22 (Norris, McQueen et al. 2000)), what computations license retroactive processing? Some  
23 models do permit bi-directional flow of information between phonological and lexical  
24 levels of analysis, such as the TRACE model (McClelland and Elman 1986). McMurray and  
25 colleagues (McMurray, Tanenhaus et al. 2009) tested the accuracy of the TRACE model in  
26 an eye-tracking study, and found that while the model was not able to account for  
27 responses in its native form, it was able to accurately predict responses when phoneme-  
28 level inhibition was removed. This finding is consistent with psycholinguistic research  
29 suggesting that, in the face of ambiguity, sub-phonemic detail is maintained to delay  
30 commitments to discrete phonological categories. Exact estimates of how long the system  
31 can delay categorisation has varied across studies, although it appears to be on the order  
32 of ~1 second (Connine, Blasko et al. 1991, Samuel 1991, McMurray, Tanenhaus et al. 2009,  
33 Szostak and Pitt 2013). And eventual commitments are affected by both word frequency  
34 and sub-phonemic information (Dahan and Gareth Gaskell 2007). However, here we make  
35 the claim that operationalising delayed commitment as sustained sensitivity to sub-  
36 phonemic detail may be inaccurate, as these two processes are not necessarily equivalent:  
37 Sub-phonemic detail may be maintained in parallel to, or independent from, committing  
38 to discrete phonological categories. No study has yet sought to piece these two  
39 computations apart.

40  
41 Three fundamental questions remain regarding the ability of subsequent information to  
42 influence phonological perception. First, does the recognition of phonological ambiguity  
43 manifest as an early perceptual or higher-order post-perceptual process? Second, how  
44 does sub-phonemic maintenance prolong the malleability of phoneme perception? Third,  
45 what temporal constraints are placed on the system — what is the limit on how late  
46 subsequent context can be received and still be optimally integrated? The current study  
47 aims to address these questions.

1 Whole-head magnetoencephalography (MEG) was recorded across two experiments. In  
 2 the first, participants performed two-alternative forced choice (2AFC) classification on  
 3 11-step syllable continua (e.g., /ba/ <-> /pa/), which provide sensory information about  
 4 onset phoneme identity but no contextual information. In the second, participants  
 5 performed a non-metalinguistic auditory-to-visual word-matching task on word to non-  
 6 word continua (“parakeet” <-> “barakeet”), where the onset of the words contained the  
 7 syllables used in Experiment 1 (see *Materials and Methods*). These stimuli provided both  
 8 sensory *and* subsequent contextual information about onset phoneme identity. MEG was  
 9 used to time-lock analyses to temporally distinct positions during the time course of  
 10 processing, and to determine the cortical regions demonstrating sensitivity to  
 11 phonological ambiguity and ambiguity resolution.

## 12 13 2. Results

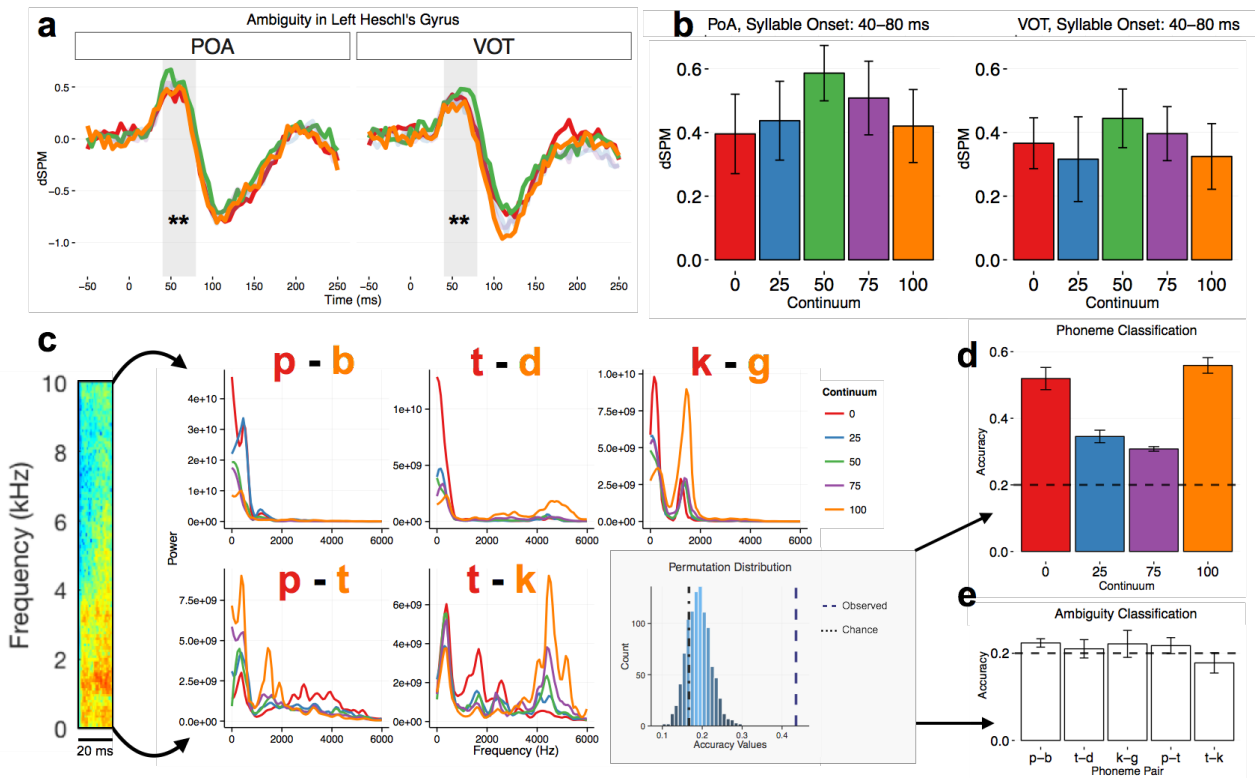
14  
15 Analyses on the behavioural responses of Experiment 1 are reported in Supplementary  
 16 materials (Section 1 & Figure 1). Analysis scripts are [available online](#); data is available on  
 17 request by contacting the first author.  
 18

19 In order to identify the neural processes that are sensitive to phonemic ambiguity, a  
 20 mixed effects regression-based temporal-spatial cluster analysis was applied to source  
 21 estimates in Heschl’s gyrus (HG) and superior temporal gyrus (STG) bilaterally, searching  
 22 a time-window of 0-200 ms after phoneme onset (see (Gwilliams, Lewis et al. 2016) for  
 23 more details concerning this analysis technique). The same analysis was conducted on  
 24 both Experiment 1 & 2 datasets. The 11 acoustically defined steps were re-sampled based  
 25 on participants’ behavioural responses, forming a perceptually defined 5-step continuum  
 26 (see *Materials & Methods* and Fig. 4). There were three primary variables of interest:  
 27 Ambiguity (distance from the behaviourally-defined perceptual boundary); Acoustics  
 28 (position on the 5-step re-sampled continuum); and Feature Type (binary variable coding  
 29 for place of articulation (POA) or voice onset time (VOT)). Trials were grouped into  
 30 phoneme categories based on participants’ behavioural responses in Experiment 1.  
 31 Number of trials into the experiment and block number were included in all models as  
 32 confound variables. We used a cluster forming threshold of  $p < .05$ , with a minimum of 10  
 33 neighbouring spatial samples, and 25 temporal samples. Bonferroni correction for  
 34 multiple comparisons over time and space was applied following (Maris and Oostenveld  
 35 2007).  
 36

37 The results for Experiment 1 are summarised in Table 1 (left) and displayed in Fig. 1A &  
 38 1B. In sum, early (~50 ms) Ambiguity effects were left lateralised, and early Acoustic

	Experiment 1: Syllable Onset				Experiment 2: Word Onset				Experiment 2: POD Onset			
	Ambiguity	Acoustics	VOT	POA	Ambiguity	Acoustics	VOT	POA	Ambiguity	Acoustics	VOT	POA
<i>p</i> -value	< .005*** .029*	0.125	< .001***	< .001***	.034 * .063 †	.019*	< .005**	< .005** .028*	.011*	.043*	< .01 **	< .001 ***
Time window (ms)	45-110 105-145	40-75	85-200	90-150	150-182 144-172	106-152	92-138	86-126 88-126	50-84	110-136	98-140	26-96
Hemisphere & Location	LH (HG) RH (STG)	RH (HG)	RH (STG)	LH (STG)	LH (HG) LH (HG)	RH (HG)	RH (STG)	LH (STG) LH (STG)	LH (HG)	LH (HG)	RH (STG)	RH (STG)

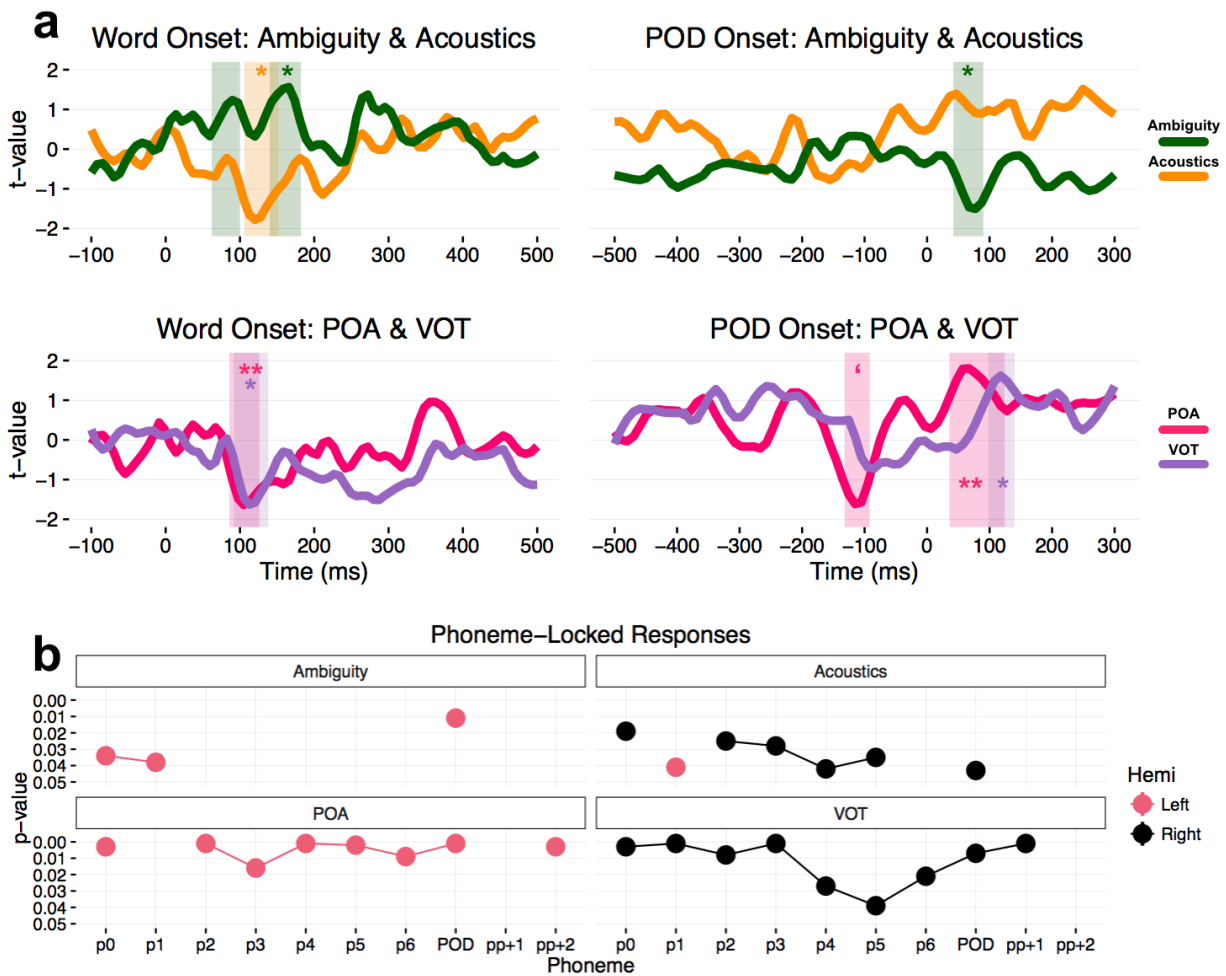
39  
40  
41  
42  
43  
44  
45  
46  
47 **Figure 1.** Summary of MEG results. VOT = “voice onset time”; POA = “place of articulation”; POD = “point of  
 48 disambiguation”.



**Figure 1.** Early responses to ambiguity in left Heschl's Gyrus (LHG) in Experiment 1. (a) Timecourse of responses for each ambiguity level averaged over source-localised responses in LHG. (b) Averaged responses in LHG over the p50m peak, plotted separately for place of articulation (POA) and voice onset time (VOT) continua. (c) FFT decomposition of first 20 ms of the auditory stimuli, plotted for each phoneme continuum. The histogram represents the 1000 permutations used to determine the significance of classification accuracy. (d) Accuracy of the logistic regression classifier in identifying the correct phoneme, based on leave-one-out cross validation — accuracy drops off for more ambiguous tokens. (e) Chance-level accuracy in classifying steps along the continuum.

effects were right lateralised. Later responses (~100 ms) interacted with the phonetic feature being manipulated (i.e., VOT vs. POA). POA modulated left hemisphere responses, and VOT modulated right hemisphere responses.

The results for Experiment 2 are summarised in Table 1 (middle). The lateralisation of effects observed in Experiment 1 was replicated: Sensitivity to Ambiguity and POA in the left hemisphere, and to Acoustics and VOT in the right hemisphere. The Ambiguity cluster was identified at ~150 ms in the lexical context, which is later than the effect found for syllable context. However, when looking at the cluster level *t*-values across time (Fig. 2A), there was a clear peak in sensitivity to Ambiguity at ~50 ms, too. To test if lexical items also elicit early sensitivity to Ambiguity, we ran a post-hoc mixed-effects regression analysis, averaging just in left Heschl's gyrus (the locus of the effect in Experiment 1) at 50 ms post word-onset (the peak of the effect in Experiment 1). Ambiguity, Acoustics, Feature Type and their interaction were coded as a fixed effect and random slopes over items. This revealed a significant interaction between Ambiguity and Feature Type ( $X = 5.9, p = .015$ ), and a significant effect of Feature Type ( $X = 13.14, p < .001$ ). When breaking the results down at each level of Feature Type, Ambiguity was a significant factor for POA contrasts ( $X = 4.84, p = .027$ ) and was approaching significance for VOT contrasts ( $X = 3.09, p = .078$ ). This analysis confirms that the early ambiguity effect is replicated in lexical contexts. Interestingly, the direction of the effect was reversed for POA contrasts, whereby more ambiguous tokens elicited less rather than more activity (SM, Fig. 3). This



**Figure 2.** Timecourse of regression analysis for the four primary variables of interest for Experiment 2. (a) Responses time-locked to word onset (left) and point of disambiguation (POD, right). Above:  $t$ -values of the Ambiguity and Acoustics variables when put into the same regression model. Below:  $t$ -values of place of articulation (POA) and voice onset time (VOT). Note that  $t$ -values are extracted from the most significant cluster formed for that variable. (b) Responses time-locked to the onset of each phoneme along the length of the word. p0 = word onset; pp = “post POD”. Statistical  $p$ -values are based on the most significant cluster. If no cluster is formed, the cell is empty for that entry.

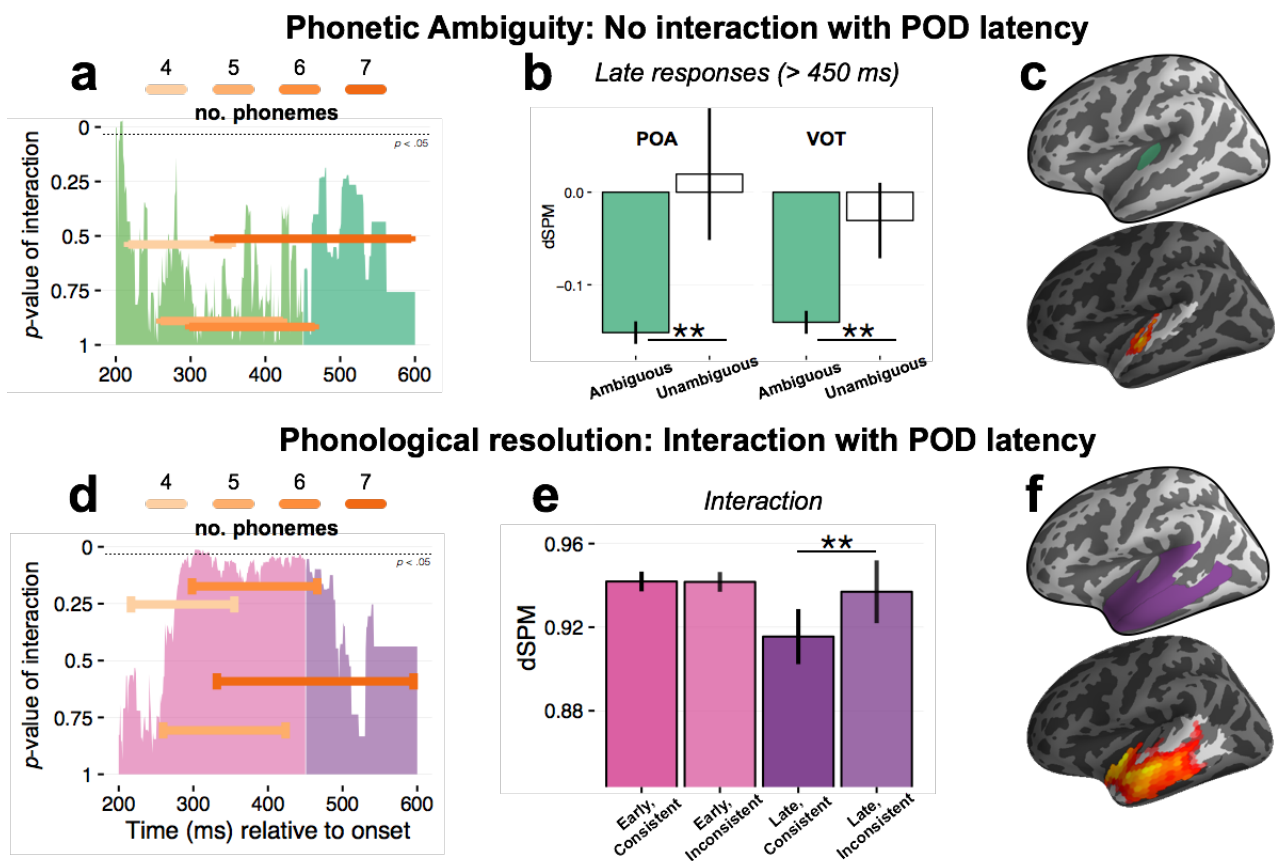
interaction may be due to differences in the task or due to processing syllables versus words. More research would need to be conducted to piece these apart.

Integration of top-down lexical information was analysed by time-locking MEG responses to the onset of a word’s point of disambiguation (POD). Main results are summarised in Table 1 (right). In sum, sensitivity to Ambiguity, Acoustics, POA and VOT re-emerges at point of disambiguation, with similar lateralisation to that observed at onset (see Fig. 2). The results for Ambiguity and Acoustics tested at each level of Feature Type are in Supplementary Materials 2.1-2.3.

In order to assess whether the results observed at POD were specific to disambiguation or were a general re-activation process, we ran the same analysis time-locked to each phoneme along the length of the word, up to two phonemes post-disambiguation (Fig. 2B). Sensitivity to Ambiguity was observed just at the initial phonemes and at disambiguation; however, sensitivity to Acoustics, VOT and POA was maintained throughout the word ( $p$ ’s < .05, corrected for multiple comparisons across phonemes).

1 To determine whether the system commits to a phonological category when  
 2 disambiguation occurs “too late”, we tested for an interaction between disambiguation  
 3 latency and whether the word resolves to the more or less likely word of the pair given  
 4 acoustics at onset. The rationale is that if the system commits to a /b/, for example, but  
 5 then the word resolves to a p-onset word, more effort is required to comprehend the  
 6 lexical item that was thrown away during the commitment process; however, if no  
 7 commitment has occurred, there should be a minimal difference between word and non-  
 8 word resolution because both the cohort of p-onset and b-onset words are still active.

10 First, we applied the spatio-temporal regression to responses 0-300 ms after to point of  
 11 disambiguation. Because this test assumes higher-level processing, the search area was  
 12 expanded to include middle temporal gyrus (Fig. 3F). The variables included in the model  
 13 were lexical resolution (word versus non-word) and its interaction with POD latency,  
 14 where latency was defined in terms of ms for one test, and in terms of phonemes in the  
 15 other. No interaction was found between phoneme-defined latency and lexical resolution.  
 16 However, defining latency in terms of ms did reveal a significant interaction in the left



41 **Figure 3.** Interaction between point of disambiguation (POD) latency and ambiguity level (above), and with lexical  
 42 resolution (below). For panels (a) and (d), POD latency defined in milliseconds (200-600 ms) is represented by filled  
 43 area, and in terms of phonemes (4-7) by horizontal orange bars. The span of the orange bars reflects the range of onset  
 44 latencies for each phoneme. (a) Non-significant interaction between latency and phoneme ambiguity when splitting  
 45 latency into “early” vs. “late” incrementally from 200-600 ms post word onset. (b) More activity for more ambiguous  
 46 sounds at disambiguation, even past 450 ms after word onset. (c) Above: Location of Heschl's Gyrus; Below: Location of  
 47 significant cluster sensitive to ambiguity. (d) Timecourse of interaction between lexical resolution and latency, again  
 48 defining latency as “early” and “late” at incremental distances from word onset. (e) Significant interaction when splitting  
 responses at 450 ms. Resolution is “consistent” when the acoustics at onset more closely match the onset of the  
 revealed word of the pair (e.g., 75% /b/ and hear “barricade”), and “inconsistent” when the acoustics better match the  
 onset of the other word of the pair (e.g., 75% /b/ and hear “parakeet”). (f) Above: Location of superior temporal and  
 middle temporal gyrus; Below: Location of cluster sensitive to the interaction between lexical resolution and latency.



1 hemisphere between 196-266 ms after POD ( $p = .02$ , Bonferroni corrected for multiple  
2 comparisons over time). Second, in order to identify the optimal split between “early” and  
3 “late”, we averaged activity over the spatio-temporal dimensions of the interaction  
4 cluster, and ran a linear mixed-effects regression analysis, testing for an interaction with  
5 latency, where latency was shifted incrementally by 1 ms from 200-600 ms after word  
6 onset. As can be seen in Fig. 3D, the interaction was maximised when setting the  
7 boundary between “early” and “late” between 292-447 ms. When running the same  
8 analysis for the Ambiguity variable, no interactions were observed with latency – words  
9 that had an ambiguous onset elicited a stronger response at POD regardless of how many  
10 ms or phonemes elapsed before disambiguation (Fig. 3A-C).

### 11 **3. Discussion**

12  
13  
14 The aim of the present study was to determine the neural computations that allow the  
15 perception of speech sounds to be affected by subsequent top-down information. Using  
16 MEG, we analysed responses to two critical points within a word: 1) varyingly ambiguous  
17 phoneme at word onset, and 2) point of disambiguation: onset of the phoneme that  
18 reveals the identity of the word, and consequently the identity of the word-initial  
19 phoneme. Our results uncover a three-stage processing network: *Maintain subphonemic*  
20 *information* in parallel to the *recognition of phonological categories*, in order to *predict*  
21 *lexical items*. These findings demonstrate how phonetic representations can be moulded  
22 by incoming higher-order information while being carried forward through the processing  
23 stream.

#### 24 **3.1 Early sensitivity to ambiguity and acoustics**

25  
26  
27 We found evidence for sensitivity to phonological ambiguity very early during processing,  
28 at just 50 ms after onset, in left Heschl’s Gyrus. This was orthogonal to sensitivity to  
29 position on the continuum, i.e., linear acoustic differences, which was right-lateralised at  
30 the same latency. While previous studies have found the p50m to be modulated by VOT  
31 (Steinschneider, Volkov et al. 1999, Hertrich, Mathiak et al. 2000) and POA (Tavabi,  
32 Obleser et al. 2007), and fMRI studies have found sensitivity to ambiguity in primary  
33 auditory cortex (Kilian-Hutten, Valente et al. 2011) (see *Introduction*), this is the first  
34 evidence of such early responses tracking proximity to perceptual boundaries. This  
35 finding supports a hierarchical over reverse-hierarchical processing model cf. (Kilian-  
36 Hutten, Valente et al. 2011) because sensitivity comes online before any top-down higher-  
37 order influence, and illustrates that early stages of processing are tuned to strikingly  
38 complex features of the acoustic signal.

39  
40 Because of the time it takes the acoustic signal to reach primary auditory cortex, the early  
41 ambiguity effect must be reflecting a response to (at most) the first 20 ms of the stimulus.  
42 As we were able to decode phoneme category from the spectrotemporal properties of the  
43 first 20 ms of the acoustic stimuli (Fig. 1C-E, Supplementary Materials 1.2), it is clear that  
44 phoneme category information is present in the signal (also see (Blumstein, Stevens et al.  
45 1977, Stevens and Blumstein 1978) for a similar conclusion in voiced POA contrasts). In  
46 light of this, a very plausible explanation for these observations is consistent with an  
47 analysis by synthesis model (Halle and Stevens 1962, Poeppel and Monahan 2011).  
48 Concretely, responses reflect the number of candidate phonemic representations

1 generated by the first ~20 ms of acoustic signal. Neurons fire more when the search space  
2 over phonemic hypotheses is large, and less when there are fewer possibilities.

3  
4 Furthering the stimuli decoding analysis, we applied the same logistic classifier to the  
5 first 60 ms of acoustic input – the likely amount of information driving the N100m  
6 response (see *Introduction*). The classifier was trained either on a single 60 ms spectral  
7 segment of the signal, or three sequential 20 ms spectral chunks. The former provides  
8 reasonable spectral resolution but poor temporal resolution; the latter provides the  
9 opposite. This novel analysis revealed intuitive results: The classifier more accurately  
10 distinguished VOT contrasts (a temporal cue) when trained on three 20 ms chunks, and  
11 POA contrasts (a spectral cue) when trained on a single 60 ms chunk (see Section 1.2 and  
12 Figure 2 in Supplementary Materials). It may be the case that the N100m response is  
13 driven by neuronal populations that sample both at fast (~20 ms) and slower (~60 ms)  
14 frequencies in order to accurately identify phonemes that vary across each phonetic  
15 dimension.

16  
17 Although not the focus of our study, the data revealed striking asymmetries across  
18 hemispheres. In early ~50 ms responses Ambiguity localised to the left hemisphere, and  
19 Acoustics to the right. In later ~100 ms responses and along the length of the word, POA  
20 was predominantly supported by the left hemisphere, while VOT was supported by the  
21 right (see Fig. 2B). Both p50m and N100m asymmetries were replicated across the two  
22 studies, suggesting that the right hemisphere does not only contribute in lexical contexts,  
23 or only in the presence of metalinguistic judgements cf. (Wolmetz, Poeppel et al. 2011).  
24 Rather it appears that, at least in the initial stages of processing, each hemisphere makes  
25 a functionally distinct contribution to the discretisation of a speech sound (Gage, Roberts  
26 et al. 2002, Poeppel 2003). Concretely, our results suggest that the left hemisphere  
27 analyses input relative to abstract phonetic features, whereas the right hemisphere  
28 extracts lower level subphonetic information, thus illustrating sensitivity to linearly  
29 varying spectral information.

### 30 31 *3.2 Re-emergence of subphonemic detail*

32  
33 We observed a re-emergence of sensitivity to the acoustics, POA and VOT of the phoneme  
34 heard at onset at each phoneme along the length of the word, at disambiguation point,  
35 and at the two phonemes *after* disambiguation. This was specifically time-locked to the  
36 onset of each incoming phoneme and was not apparent when analysing based on the time  
37 elapsed from word onset (contrast Fig. 2A with Fig. 2B). Sensitivity to onset phoneme  
38 ambiguity revealed a different pattern. It faded at intervening positions, resurfaced at the  
39 disambiguation point, and decayed after ambiguity resolution. While this may reflect a  
40 more subtle response that we have insufficient power to detect at intermediate positions,  
41 it is possible that sensitivity to this variable only resurfaces once phoneme identity can be  
42 fully resolved. This novel finding is critically important because it supports the hypothesis  
43 that the sub-phonemic representation of a speech sound is maintained in superior  
44 temporal regions throughout the duration of a word, even while subsequent phonemes are  
45 being received; perhaps suggesting that the percept of a speech sound is reassessed at  
46 each increment based on the provision of additional input. This finding is also consistent  
47 with a recent study using EEG (Khalighinejad, da Silva et al. 2017), which found evidence  
48 for continued maintenance of phoneme-category distinctions.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

Further, it appears that phonemic reactivation is a general feature of speech comprehension, rather than a specific mechanism recruited in the presence of ambiguity. Specifically, our results indicate that subphonemic information is maintained even when uncertainty about phoneme identity is low, in two capacities. First, re-emergence of phonetic properties was not specific to the ambiguous tokens – it also occurred for the unambiguous phonemes. Second, information about phonetic features continues to be conserved after disambiguating information became available. Overall, these observations are the first to reveal that subphonemic information is maintained, not just in terms of uncertainty about categorisation, but *also* in terms of fine-grained phonetic and acoustic detail of the phoneme under scrutiny. And both sources of information continue to be revisited over long timescales.

### 3.3 Commitment to phonological categories

Finally, we do see evidence for phonological commitment, resolving on a time-scale of ~300-450 ms (see Fig. 3). The superiority of defining latency in terms of elapsed ms rather than phonemes may indicate that commitment is based on the amount of time or number of completed processing cycles rather than intervening information (Fig. 3C). This process is supported by higher auditory processing regions in anterior STG, a location consistent with a recent meta-analysis of auditory word recognition (DeWitt and Rauschecker 2012). Critically, this seems to be computed in parallel to the maintenance of subphonemic detail in primary auditory regions. Before ~300 ms there is no cost associated with resolution to a lexical item less consistent with word onset: listeners do not get temporarily misled (garden-pathed) provided resolution comes early enough (Fig. 3D). This suggests that the cohort of words consistent with either phonological interpretation is considered together (e.g., in the presence of b/p ambiguity, both the p-onset and b-onset words are activated). This is fully consistent with previous behavioural studies (Martin and Bunnell 1981, Gow 2001, Gow and McMurray 2007), and a previous eye-tracking study (McMurray, Tanenhaus et al. 2009), which used similar materials and found look-contingent responses to be dependent upon phonetic information at lexical onset until at least ~300 ms (the longest disambiguation delay they tested). However, after ~450 ms a cost begins to emerge when there is a mismatch between the more likely word given word-onset and the resolving lexical information (e.g., “barricade” is more likely if the onset phoneme was more b-like than p-like, so hearing “parakeet” is a mismatch). This plausibly reflects the recruitment of a repair mechanism, a prediction-error response or re-analysis of the input from making an incorrect commitment.

Finding maintained sensitivity to subphonemic detail in parallel to phonological commitment is very important for the interpretation of psychophysical research, which has implicitly equated insensitivity to within-category variation with phonological commitment (Connine, Blasko et al. 1991, Szostak and Pitt 2013, Bicknell, Tanenhaus et al. 2015). This previous work has largely converged on a processing model whereby phonological commitment can be delayed for over one second after onset. Our results indicate, in contrast, that while subphonemic detail is indeed maintained over large time-scales, this does not implicate that commitment is also put off for this length of time.

### 3.4 Relationship to models of speech processing

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

It is unclear which model of speech processing can account for these data. While Shortlist (Norris 1994) and Shortlist B (Norris and McQueen 2008) may be able to model recovery from lexical garden-paths, they do not explicitly model processing of subphonemic detail. While the MERGE model (Norris, McQueen et al. 2000) is capable of modelling such detail, it proposes no feedback from the lexical to phoneme levels of analysis, which is inconsistent with our observation that (sub-)phonemic representations are reactivated when top-down lexical information becomes available. Although it has been demonstrated that TRACE (McClelland and Elman 1986) can be modified to simulate recovery by removing phoneme-level inhibition (McMurray, Tanenhaus et al. 2009), it does not provide the architecture to model initial sensitivity to phoneme ambiguity, or account for how the percept of speech sounds is modulated by past and future linguistic information (see (Grossberg and Kazerounian 2011) for an overview of TRACE limitations). It is also unclear whether this modification would interfere with TRACE's success in accounting for a range of observations in spoken word recognition (see (Gaskell 2007) for a review). One model proposed to deal with TRACE's shortcoming is Adaptive Resonance Theory (ART): each speech sound produces a resonance wave that is influenced by top-down information until it reaches equilibrium and surfaces to consciousness (Grossberg 2003). While this theory is consistent with the idea that there is a critical time-limit to receive top-down information, it suggests that there is a linear decay in subphonemic information as temporal distance from the phoneme increases. Our results do not support that conjecture. Instead, they suggest that subphonemic information is re-evoked later in processing, with a similar magnitude as that experienced at onset. In light of the present results, one shortcoming of these models is their attempt to explain spoken word recognition with a single mechanism, built on the assumption that acoustic-phonetic information is lost once a phonological categorisation is derived. Instead, our results suggest that a three-element processing model is more appropriate, allowing for a dynamic interaction between phonetic, phonological and lexical levels of analysis.

### 3.5 Conclusion

Later sounds determine the perception of earlier speech sounds through the simultaneous recruitment of acoustic-phonetic and phonological computational pathways. This facilitates contact with lexical items in order to derive the message of the utterance, as well as continued revisitation to the phonetic level of analysis to reduce parsing errors. In this manner, lexical selection can be achieved rapidly, while also reducing the likelihood of mistakes in phonological segmentation. The human brain therefore solves the issue of processing a transient hierarchically structured signal by recruiting complementary computations in parallel, rather than conceding to the trade-off between speed and accuracy.

## 4. Materials & Methods

### 4.1 Experiment 1

#### 4.1.1 Participants

1 Twenty-four right handed native English participants took part in the study (11 female;  
2 age:  $M=25.44$ ,  $SD=8.44$ ). They were recruited from the NYUAD community and were  
3 compensated for their time. All had normal or corrected vision, normal hearing and no  
4 history of neurological disorders.

#### 5 6 4.1.2 Materials

8 Phonological transcriptions and by-phoneme frequency counts were extracted from the  
9 English Lexicon Project (ELP) (using scripts [available online](#)). We selected all word pairs  
10 whose first phoneme was a plosive, and differed either in terms of voicing {t-d, p-b, k-g},  
11 or in terms of one place of articulation feature {t-k, p-t}. Any item that occurred in more  
12 than one pair was eliminated, resulting in 53 pairs: 31 VOT and 22 POA. Words ranged  
13 from 4-10 phonemes in length ( $M=6.8$ ;  $SD=1.33$ ), and from 291-780 ms in length ( $M=528$ ;  
14  $SD=97$ ). Latency of disambiguation ranged from 3-8 phonemes ( $M=5.1$ ;  $SD=0.97$ ) and  
15 142-708 ms ( $M=351$ ;  $SD=92$ ).

16  
17 A native English speaker was recorded saying the selected words in isolation. Each item  
18 pair was exported into TANDEM-STRAIGHT for the morphing procedure (Kawahara,  
19 Morise et al. 2008, Kawahara and Morise 2011). In short, the morphing works by taking the  
20 follows steps: 1) position anchor points to mark the onset of each paired phoneme; 2)  
21 place weights on each anchor point to determine the % contribution of each word at the  
22 end-points of the continuum; 3) specify the number of continuum steps to generate.

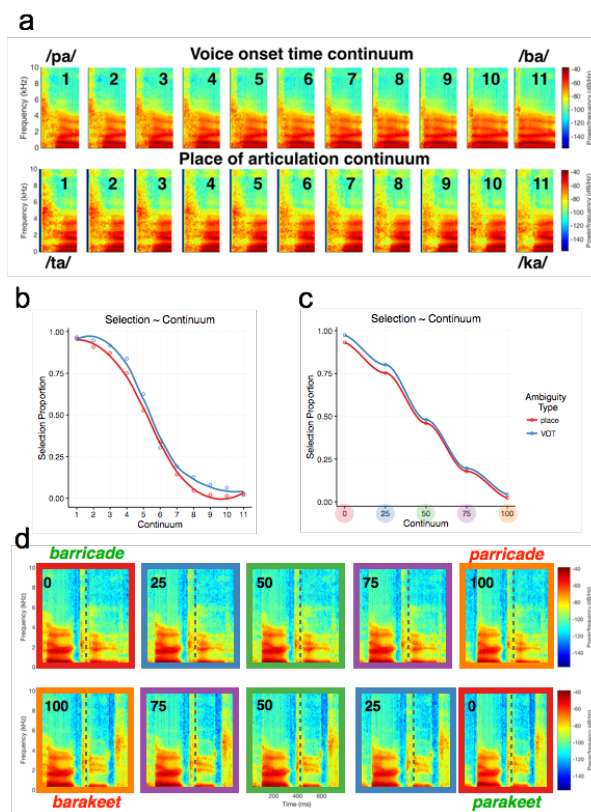
23  
24 In long, the process involves decomposing the auditory signal into its spectrotemporal  
25 components to make a confusion matrix between each auditory sample of the lexical pair.  
26 Then, anchor points are placed along this confusion matrix to mark temporal  
27 correspondence between the onset of phonemes in each word. For example, the onset  
28 occurs at different latencies for each of the seven phonemes of the words “barricade” and  
29 “parakeet”, but a single anchor point is placed to mark the correspondence in timing  
30 between the phonemes of the two words. This forms a morphing substrate that can be  
31 used to adjust the weighted contribution of the acoustics of each phoneme of each word.  
32 Then, the user specifies the anchor position weights at the two ends of the continuum. For  
33 example, to generate a “barakeet” <-> “parakeet” continuum, the onset anchor points will  
34 be placed at 100% “barricade” at one end of the continuum, and 100% “parakeet” at the  
35 other end of the continuum. Until and including disambiguation point (“arak-”), all  
36 anchors will be weighted at 50-50 of each word, and after disambiguation (“-eet”), all  
37 anchors are placed to be 100% “parakeet”, for both sides of the continuum. In general, for  
38 each pair, all anchor points before the POD phoneme were placed at the 50% position, and  
39 the first anchor point was positioned either in the congruent position, creating a word  
40 (“parakeet”) or with the competitor, creating a non-word (“barakeet”). This ensured that  
41 apart from the first phoneme, the acoustic signal remained identical across the two word  
42 pairs until the disambiguation point. Eleven continua steps were created for each item.  
43 The resulting 1166 auditory files were annotated in Praat for each phoneme onset using  
44 the Penn Forced Aligner, making manual adjustments as necessary. The first consonant-  
45 vowel sequence was then extracted for each item, and amplitude-normed to 70 dB.

#### 46 47 4.1.3 Procedure

48

1 Stimuli were separated into eleven blocks using a latin-square design. Each block  
2 contained two morphed tokens from the stimuli pair, with the constraint that each item  
3 had to be at least three morphed steps away from its paired counterpart. This resulted in  
4 106 trials per block. Items were randomised within each block.

5  
6 Participants performed a single interval 2AFC classification on each item, while MEG was  
7 recorded concurrently. First, two options appeared (e.g., “B or P”) on a screen 85 cm away  
8 from the participant’s face. The auditory stimulus began 500 ms later, and the options  
9 remained on the screen until a left or right selection was made via button press using the  
10 left hand. After the response, a blank screen was displayed for 500 ms until the onset of  
11 the next trial. At each block interval, participants had a self-terminated break. The  
12 recording session lasted ~50 minutes.



13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
**Figure 4.** Stimuli production. (a) Example 11-step continuum for voice onset time (above) and place of articulation (below) contrasts used in Experiment 1. (b) Behavioural classification across each continuum step of Experiment 1. (c) Re-sampled continuum at five perception-defined thresholds. (d) Example word (green) - nonword (red) continuum; point of disambiguation is represented with a dashed line.

## 4.2 Experiment 2

### 4.2.1 Participants

Twenty-five right handed native English participants took part in the study (15 female; age:  $M=24.84$ ,  $SD=7.3$ ). Six had taken part in Experiment 1 two months earlier. All had normal or corrected vision, normal hearing, no history of neurological disorders, and were recruited from the NYUAD community.

## 1 4.2.2 Materials

2  
3 Five out of the initial 11-step continuum were selected for each word pair based on the  
4 psychometric responses from Experiment 1. Placing a selection threshold at 0.05, 0.25,  
5 0.5, 0.75 and 0.95, the closest morph position was extracted (see Fig. 4). Items were  
6 rejected if the unambiguous endpoints of the continuum were not categorised with at  
7 least 80% accuracy for all subjects, or if the position of the ambiguous token was not at  
8 least three points away from either end-point of the continuum. This resulted in 49 pairs.  
9 All remaining words were amplitude-normed to 70 dB.

## 10 11 4.2.3 Procedure

12  
13 Each participant's head scan and digitised fiducial points were first collected, as described  
14 for Experiment 1.

15  
16 Participants performed an auditory-to-visual word matching task on 1/5 of the auditory  
17 items. The visual word was either the target ("match" response) or the competitor ("mis-  
18 match" response). One item of the five morphed tokens would be a "match" trial and  
19 another would be a "non-match" trial; randomly assigned using a latin-square procedure.  
20 The experiment was split into 5 blocks, and only one of the morphed tokens from each  
21 item appeared in each block. The assignment of item-block was also pseudo-randomised  
22 in a latin-square fashion. This resulted in 25 unique experimental orders, across which,  
23 items were matched for block order, and match-nonmatch assignment.

24  
25 A fixation cross was displayed for 500 ms, followed by the auditory word. If it was a task  
26 trial, the visual word would appear 500 ms after auditory word offset, and remain on  
27 screen until participants made a match (left button) or non-match (right button) decision  
28 with their left hand. If there was no task, a blank screen was presented and participants  
29 could move forward by pressing either button. The recording lasted ~40 minutes.

30  
31 All participants' head shapes were digitised using a hand-held FastSCAN laser scanner  
32 (Polhemus, VT, USA) to allow for co-registration during data preprocessing. Five points on  
33 each participant's head were also digitised: just anterior of the left and right auditory  
34 canal, and three points on the forehead. Marker coils were later placed at the same five  
35 positions to localise each participant's skull relative to the sensors. These marker  
36 measurements were recorded just before and after the experiment in order to track the  
37 degree of movement during the recording.

38  
39 MEG data were recorded continuously using a 208 channel axial gradiometer system  
40 (Kanazawa Institute of Technology, Kanazawa, Japan), with a sampling rate of 1000 Hz  
41 and applying an online low-pass filter of 200 Hz.

42  
43 MEG data from the two experiments underwent the same pre-processing steps. First, the  
44 continuous recording was noise reduced using Continuously Adjusted Least Squares  
45 Method (CALM: [53]), with MEG160 software (Yokohawa Electric Corporation and Eagle  
46 Technology Corporation, Tokyo, Japan). The noise-reduced data, digital scan and fiducials,  
47 and marker measurements were exported into MNE-Python [54]. Bad channels were  
48 removed through visual inspection. Independent Component Analysis (ICA) was

1 computed over the noise-reduced data using FastICA in MNE-Python. Components were  
2 removed from the raw recording if they contained ocular or cardiac artefacts, which were  
3 identified based on the topography of magnetic activity and time-course response. The  
4 data were then epoched from 500 ms pre-stimulus to 1000 ms post-stimulus onset. Any  
5 trials whose amplitude exceeded +/- 2000 femto-tesla absolute or peak-to-peak threshold  
6 were removed.

7  
8 In order to perform source localisation, the location of the subject's head was co-  
9 registered with respect to the sensory array in the MEG helmet. For subjects with  
10 anatomical MRI scans (n=4), this involved rotating and translating the digital scan to  
11 minimise the distance between the fiducial points of the MRI and the head scan. For  
12 participants without anatomical scans, the FreeSurfer "fsaverage" brain was used, which  
13 involved first rotation and translation, and then scaling the average brain to match the  
14 size of the head scan.

15  
16 Next, a source space created, consisting of 2562 potential electrical sources per  
17 hemisphere. At each source, activity was computed for the forward solution with the  
18 Boundary Element Model (BEM) method, which provides an estimate of each MEG  
19 sensor's magnetic field in response to a current dipole at that source. The inverse solution  
20 was computed from the forward solution and the grand average activity across all trials.  
21 Data were converted into noise-normalised Dynamic Statistical Parameter Map (dSPM)  
22 units (see [55]), employing an SNR value of 2. The inverse solution was applied to each  
23 trial at every source, for each millisecond defined in the epoch, employing a fixed  
24 orientation of the dipole current that estimates the source normal to the cortical surface  
25 and retains dipole orientation.

## Acknowledgements

This research was supported by ERC-2011-AdG 295810 BOOTPHON, ANR-10-IDEX-0001-02 PSL and ANR-10-LABX-0087 IEC to TL; NIH 2R01DC05660 to DP; NYU Abu Dhabi Institute under grant G1001 to AM. We would like to thank Kyriaki Neophytou for her help with data collection.



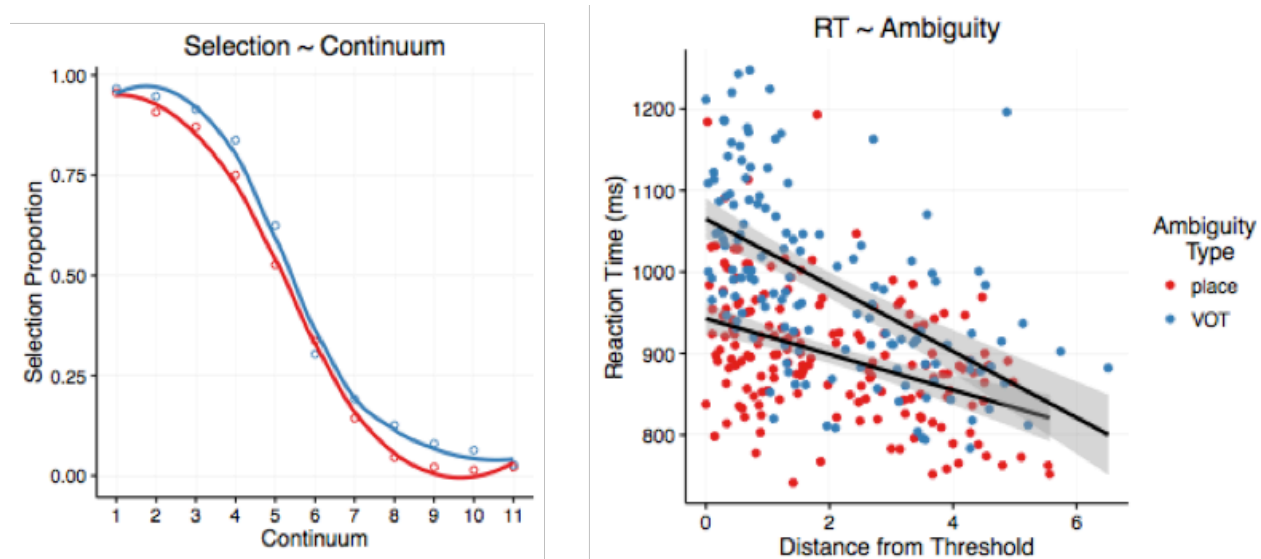
## Bibliography

- Ackermann, H., W. Lutzenberger and I. Hertrich (1999). "Hemispheric lateralization of the neural encoding of temporal speech features: a whole-head magnetencephalography study." Cognitive Brain Research **7**(4): 511-518.
- Bicknell, K., M. Tanenhaus and T. Jaeger (2015). "Listeners can maintain and rationally update uncertainty about prior words." Manuscript submitted for publication.[KB].
- Blumstein, S. E., E. B. Myers and J. Rissman (2005). "The perception of voice onset time: an fMRI investigation of phonetic category structure." Journal of Cognitive Neuroscience **17**(9): 1353-1366.
- Blumstein, S. E., K. N. Stevens and G. N. Nigro (1977). "Property detectors for bursts and transitions in speech perception." The Journal of the Acoustical Society of America **61**(5): 1301-1313.
- Cole, R. A. (1973). "Listening for mispronunciations: A measure of what we hear during speech." Perception & Psychophysics **13**(1): 153-156.
- Connine, C. M., D. G. Blasko and M. Hall (1991). "Effects of subsequent sentence context in auditory word recognition: Temporal and linguistic constraint." Journal of Memory and Language **30**(2): 234-250.
- Dahan, D. and M. Gareth Gaskell (2007). "The temporal dynamics of ambiguity resolution: Evidence from spoken-word recognition." J Mem Lang **57**(4): 483-501.
- DeWitt, I. and J. P. Rauschecker (2012). "Phoneme and word recognition in the auditory ventral stream." Proceedings of the National Academy of Sciences **109**(8): E505-E514.
- Di Liberto, G. M., J. A. O'Sullivan and E. C. Lalor (2015). "Low-Frequency Cortical Entrainment to Speech Reflects Phoneme-Level Processing." Curr Biol **25**(19): 2457-2465.
- Ettinger, A., T. Linzen and A. Marantz (2014). "The role of morphology in phoneme prediction: evidence from MEG." Brain Lang **129**: 14-23.
- Gage, N. M., T. P. Roberts and G. Hickok (2002). "Hemispheric asymmetries in auditory evoked neuromagnetic fields in response to place of articulation contrasts." Cognitive brain research **14**(2): 303-306.
- Gagnepain, P., R. N. Henson and M. H. Davis (2012). "Temporal predictive codes for spoken words in auditory cortex." Curr Biol **22**(7): 615-621.
- Ganong, W. F. (1980). "Phonetic categorization in auditory word perception." Journal of Experimental Psychology: Human Perception and Performance **6**(1): 110.
- Gaskell, M. G. (2007). "Statistical and connectionist models of speech perception and word recognition." The Oxford handbook of psycholinguistics: 55-69.
- Gordon, P. C., J. L. Eberhardt and J. G. Rueckl (1993). "Attentional modulation of the phonetic significance of acoustic cues." Cognitive Psychology **25**(1): 1-42.
- Gow, D. W. (2001). "Assimilation and anticipation in continuous spoken word recognition." Journal of Memory and Language **45**(1): 133-159.
- Gow, D. W., Jr. and J. A. Segawa (2009). "Articulatory mediation of speech perception: a causal analysis of multi-modal imaging data." Cognition **110**(2): 222-236.
- Gow, D. W., Jr., J. A. Segawa, S. P. Ahlfors and F. H. Lin (2008). "Lexical influences on speech perception: a Granger causality analysis of MEG and EEG source estimates." Neuroimage **43**(3): 614-623.
- Gow, D. W. and B. McMurray (2007). "Word recognition and phonology: The case of English coronal place assimilation." Papers in laboratory phonology **9**: 173-200.
- Grossberg, S. (2003). Adaptive resonance theory, Wiley Online Library.
- Grossberg, S. and S. Kazerounian (2011). "Laminar cortical dynamics of conscious speech perception: Neural model of phonemic restoration using subsequent context in noise." The Journal of the Acoustical Society of America **130**(1): 440-460.

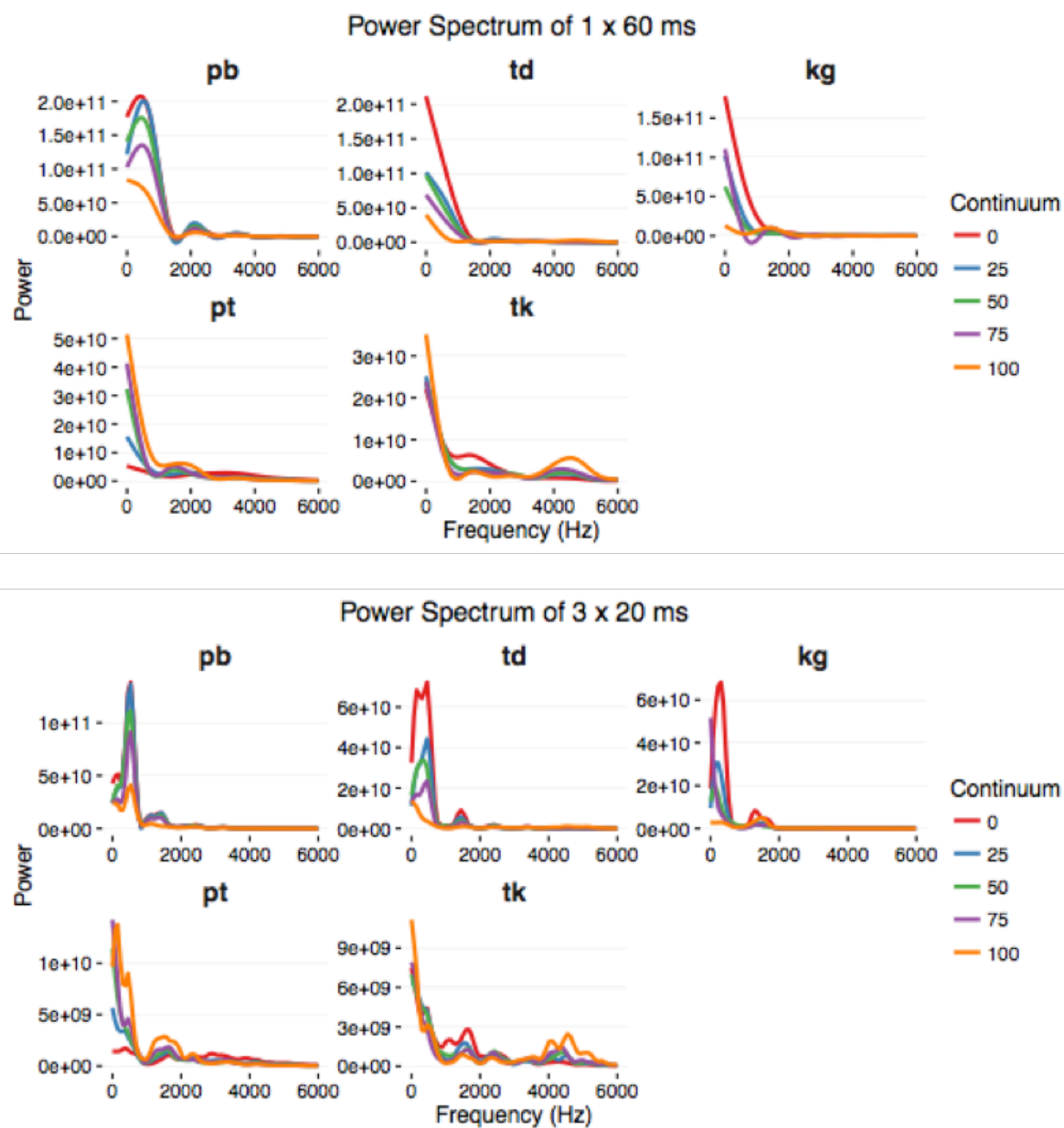
- Gwilliams, L., G. Lewis and A. Marantz (2016). "Functional characterisation of letter-specific responses in time, space and current polarity using magnetoencephalography." NeuroImage **132**: 320-333.
- Gwilliams, L. and A. Marantz (2015). "Non-linear processing of a linear speech stream: The influence of morphological structure on the recognition of spoken Arabic words." Brain and language **147**: 1-13.
- Halle, M. and K. Stevens (1962). "Speech recognition: A model and a program for research." IRE Transactions on Information Theory **8**(2): 155-159.
- Hertrich, I., K. Mathiak, W. Lutzenberger and H. Ackermann (2000). "Differential impact of periodic and aperiodic speech-like acoustic signals on magnetic M50/M100 fields." NeuroReport **11**(18): 4017-4020.
- Hickok, G. and D. Poeppel (2004). "Dorsal and ventral streams: a framework for understanding aspects of the functional anatomy of language." Cognition **92**(1-2): 67-99.
- Kawahara, H. and M. Morise (2011). "Technical foundations of TANDEM-STRAIGHT, a speech analysis, modification and synthesis framework." Sadhana **36**(5): 713-727.
- Kawahara, H., M. Morise, T. Takahashi, R. Nisimura, T. Irino and H. Banno (2008). Tandem-STRAIGHT: A temporally stable power spectral representation for periodic signals and applications to interference-free spectrum, F0, and aperiodicity estimation. 2008 IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE.
- Khalighinejad, B., G. C. da Silva and N. Mesgarani (2017). "Dynamic Encoding of Acoustic Features in Neural Responses to Continuous Speech." Journal of Neuroscience **37**(8): 2176-2185.
- Kilian-Hutten, N., G. Valente, J. Vroomen and E. Formisano (2011). "Auditory cortex encodes the perceptual interpretation of ambiguous sound." J Neurosci **31**(5): 1715-1720.
- Liebenthal, E., J. R. Binder, S. M. Spitzer, E. T. Possing and D. A. Medler (2005). "Neural substrates of phonemic perception." Cereb Cortex **15**(10): 1621-1631.
- Maris, E. and R. Oostenveld (2007). "Nonparametric statistical testing of EEG- and MEG-data." J Neurosci Methods **164**(1): 177-190.
- Marslen-Wilson, W. and L. K. Tyler (1980). "The temporal structure of spoken language understanding." Cognition **8**(1): 1-71.
- Marslen-Wilson, W. D. (1987). "Functional parallelism in spoken word-recognition." Cognition **25**(1): 71-102.
- Martin, J. G. and H. T. Bunnell (1981). "Perception of anticipatory coarticulation effects." The Journal of the Acoustical Society of America **69**(2): 559-567.
- McClelland, J. L. and J. L. Elman (1986). "The TRACE model of speech perception." Cognitive psychology **18**(1): 1-86.
- McMurray, B., M. K. Tanenhaus and R. N. Aslin (2009). "Within-category VOT affects recovery from "lexical" garden paths: Evidence against phoneme-level inhibition." J Mem Lang **60**(1): 65-91.
- McQueen, J. M. (1991). "The influence of the lexicon on phonetic categorization: stimulus quality in word-final ambiguity." Journal of Experimental Psychology: Human Perception and Performance **17**(2): 433.
- Mesgarani, N., C. Cheung, K. Johnson and E. F. Chang (2014). "Phonetic feature encoding in human superior temporal gyrus." Science **343**(6174): 1006-1010.
- Myers, E. B. and S. E. Blumstein (2008). "The neural bases of the lexical effect: an fMRI investigation." Cereb Cortex **18**(2): 278-288.
- Norris, D. (1994). "Shortlist: A connectionist model of continuous speech recognition." Cognition **52**(3): 189-234.

- Norris, D. and J. M. McQueen (2008). "Shortlist B: a Bayesian model of continuous speech recognition." *Psychol Rev* **115**(2): 357-395.
- Norris, D., J. M. McQueen and A. Cutler (2000). "Merging information in speech recognition: Feedback is never necessary." *Behavioral and Brain Sciences* **23**(03): 299-325.
- Obleser, J., A. Lahiri and C. Eulitz (2003). "Auditory-evoked magnetic field codes place of articulation in timing and topography around 100 milliseconds post syllable onset." *Neuroimage* **20**(3): 1839-1847.
- Obleser, J., A. Lahiri and C. Eulitz (2004). "Magnetic brain response mirrors extraction of phonological features from spoken vowels." *Journal of Cognitive Neuroscience* **16**(1): 31-39.
- Papanicolaou, A. C., E. Castillo, J. I. Breier, R. N. Davis, P. G. Simos and R. L. Diehl (2003). "Differential brain activation patterns during perception of voice and tone onset time series: a MEG study." *Neuroimage* **18**(2): 448-459.
- Poeppel, D. (2003). "The analysis of speech in different temporal integration windows: cerebral lateralization as 'asymmetric sampling in time'." *Speech Communication* **41**(1): 245-255.
- Poeppel, D. and P. J. Monahan (2011). "Feedforward and feedback in speech perception: Revisiting analysis by synthesis." *Language and Cognitive Processes* **26**(7): 935-951.
- Rauschecker, J. P. and S. K. Scott (2009). "Maps and streams in the auditory cortex: nonhuman primates illuminate human speech processing." *Nat Neurosci* **12**(6): 718-724.
- Rogers, J. C. and M. H. Davis (2017). "Inferior Frontal Cortex Contributions to the Recognition of Spoken Words and Their Constituent Speech Sounds." *J Cogn Neurosci*: 1-18.
- Samuel, A. G. (1981). "The role of bottom-up confirmation in the phonemic restoration illusion." *Journal of Experimental Psychology: Human Perception and Performance* **7**(5): 1124.
- Samuel, A. G. (1991). "A further examination of attentional effects in the phonemic restoration illusion." *The Quarterly Journal of Experimental Psychology* **43**(3): 679-699.
- Scott, S. K. and I. S. Johnsrude (2003). "The neuroanatomical and functional organization of speech perception." *Trends in neurosciences* **26**(2): 100-107.
- Simos, P. G., R. L. Diehl, J. I. Breier, M. R. Molis, G. Zouridakis and A. C. Papanicolaou (1998). "MEG correlates of categorical perception of a voice onset time continuum in humans." *Cognitive Brain Research* **7**(2): 215-219.
- Steinschneider, M., I. O. Volkov, M. D. Noh, P. C. Garell and M. A. Howard (1999). "Temporal encoding of the voice onset time phonetic parameter by field potentials recorded directly from human auditory cortex." *Journal of neurophysiology* **82**(5): 2346-2357.
- Stevens, K. N. and S. E. Blumstein (1978). "Invariant cues for place of articulation in stop consonants." *The Journal of the Acoustical Society of America* **64**(5): 1358-1368.
- Szostak, C. M. and M. A. Pitt (2013). "The prolonged influence of subsequent context on spoken word recognition." *Atten Percept Psychophys* **75**(7): 1533-1546.
- Tavabi, K., J. Obleser, C. Dobel and C. Pantev (2007). "Auditory evoked fields differentially encode speech features: an MEG investigation of the P50m and N100m time courses during syllable processing." *Eur J Neurosci* **25**(10): 3155-3162.
- Tyler, L. K. (1984). "The structure of the initial cohort: Evidence from gating." *Perception & Psychophysics* **36**(5): 417-427.
- Warren, R. M. (1970). "Perceptual restoration of missing speech sounds." *Science* **167**(3917): 392-393.
- Wolmetz, M., D. Poeppel and B. Rapp (2011). "What does the right hemisphere know about phoneme categories?" *Journal of cognitive neuroscience* **23**(3): 552-569.

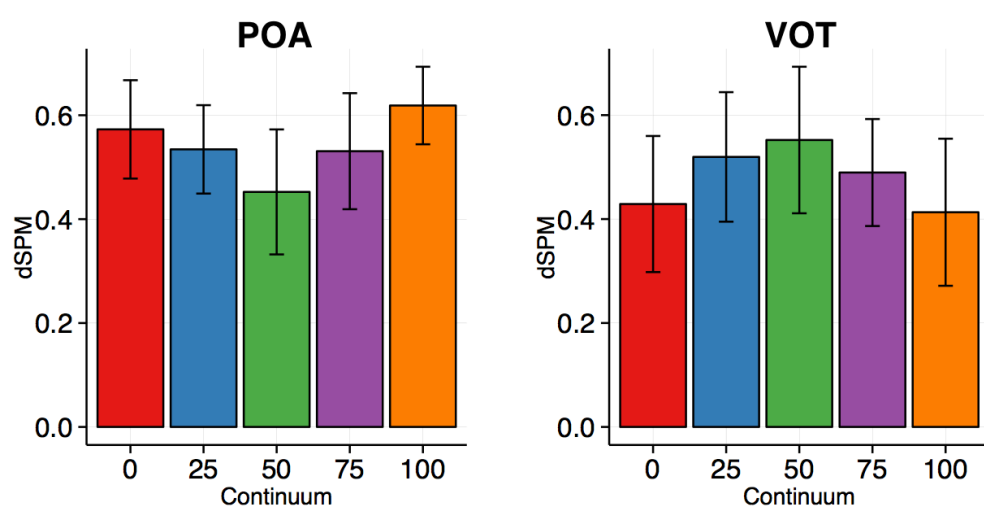
## Supplementary Materials



**Figure 1.** Behavioural phoneme selection and reaction time in experiment 1. Left: Proportion of positive selections as a function of the 11-step continuum. Right: Reaction time as a function of each trial's distance from the 50% selection threshold.



**Figure 2.** Frequency decomposition of first 60 ms of stimuli items. Above: Signal decomposed over a single 60 ms chunk. Below: Signal decomposed and averaged over three sequential 20 ms chunks.



**Figure 3.** Early ambiguity effects in Heschl's gyrus in Experiment 2. Effects are shown separately for place of articulation continua (left) and voice onset time continua (right).

## 1 Behavioural results: Experiment 1

To analyse behavioural responses in Experiment 1, we applied a mixed effects regression analysis, coding for Feature Type (place of articulation [POA] vs. voice onset time [VOT]), position on the morphed continuum (henceforth “Acoustics”), Ambiguity (distance from 50% selection point), and their interaction terms (see *Materials & Methods* for detailed information on these variables). Reaction times were significantly slower for VOT continua than POA continua ( $p < .001$ ). Responses significantly slowed with increased Ambiguity ( $p < .001$ ). Ambiguity and Feature Type revealed a significant interaction ( $p < .01$ ). Continuum had a marginal effect on reaction time ( $p = .068$ ). A logistic regression was applied to behavioural selection with the same model structure. Feature Type was a significant predictor ( $p < .001$ ), as well as Continuum ( $p < .001$ ). There were no significant interactions. See Supplementary Figure 1.

## 2. Ambiguity and Acoustics at each level of Feature Type

### 2.1 Experiment 1: Syllable Onset

There was a significant interaction between Ambiguity and Feature Type, forming two clusters in the right hemisphere: 110-150 ms ( $p = .033$ ); 155-200 ms ( $p = .049$ ). The interaction between Acoustics and Feature Type also formed two significant clusters, one in the left hemisphere: 100-145 ms ( $p = .045$ ) and one in the right: 105-170 ms ( $p = .05$ ). For POA continua, one marginally significant cluster was found for Ambiguity in the right hemisphere: 80-130 ms ( $p = .054$ ). For VOT, there was a marginally significant effect of Ambiguity 70-110 ms ( $p = .057$ ) and Acoustics 85-115 ms ( $p = .06$ ) both in the left hemisphere.

### 2.2 Experiment 2: Word Onset

A marginally significant interaction was found between Ambiguity and Feature Type in the left hemisphere: 176-200 ms ( $p = .087$ ), and between Acoustics and Feature Type in the right hemisphere: 84-152 ms ( $p < .005$ ); 84-148 ms ( $p = .028$ ). For POA, no clusters were formed for Ambiguity but Acoustics modulated later responses in the left hemisphere: 178-200 ms ( $p = .02$ ). Responses to VOT continua displayed a significant cluster for Ambiguity in the left hemisphere: 132-180 ms ( $p = .021$ ), and Acoustics revealed two clusters in the right hemisphere: 88-152 ms ( $p = .011$ ); 90-148 ms ( $p = .058$ ).

### 2.3 Experiment 2: POD

There was a significant interaction between Ambiguity and Feature Type in the left hemisphere: 126-200 ms ( $p < .001$ ), and a significant interaction between Acoustics and Feature Type in the left 60-152 ms ( $p < .001$ ) and right hemisphere 72-168 ms ( $p < .01$ ). Breaking down responses to each level of Feature Type revealed a significant effect of Ambiguity for POA in the left hemisphere: 132-182 ms ( $p = .05$ ). Responses to VOT continua revealed a significant effect of both Ambiguity in the left hemisphere: 44-108 ms ( $p < .001$ ), and Acoustics in the right hemisphere: 116-158 ms ( $p = .06$ ).

## 3. Classifying phonemes based on first 20 ms of auditory stimuli

1

2

3

4

5

6

7

8

9

10

11

12

13

Sensitivity to Ambiguity at 50 ms after onset must be reflecting a response to no more than the first 20 ms of the acoustic signal — just the noise burst of the voiceless items, and the initial voicing of the voiced items. To assess what information is available at this latency, we decomposed the first 20 ms of each stimulus into its frequency power spectra using Fast Fourier Transform (FFT). Power at each frequency band from 0-10 KHz, for all stimuli except the fully ambiguous items, was used to train a logistic regression classifier to decode the phonological category (Fig. 1C). Accuracy was significantly above chance level, as determined by 1000 random permutations of phoneme labels ( $p < .001$ ). Accuracy of classification decreased as a function of ambiguity (Fig. 1D), but all continua steps performed greater than chance. Importantly, continua steps themselves could not be decoded from this signal (Fig. 1E), suggesting that this early response indeed scales with distance from the perceptual boundary, and not acoustic properties *per se*.