

Auditory word comprehension is less incremental in isolated words

Phoebe Gaston,^{1,^,*} Christian Brodbeck,^{2,^} Colin Phillips,¹ and Ellen Lau¹

¹Department of Linguistics, University of Maryland, College Park, MD 20742

²Institute for Systems Research, University of Maryland, College Park, MD 20742

[^]Denotes equal contribution. Phoebe Gaston and Christian Brodbeck are now at Department of Psychological Sciences, University of Connecticut, Storrs, CT 06269.

* Corresponding author: phoebe.gaston@uconn.edu

Acknowledgements

We thank Daphne Amir, Fen Ingram, and Stephanie Pomrenke for assistance with stimulus selection, and Aura Cruz Heredia for assistance with some of the data collection.

Conflict of Interest

Authors report no conflict of interest.

Funding Sources

This material is based upon work supported by the National Science Foundation under Grants BCS-1749407 (E. Lau, PI) and DGE-1449815 (C. Phillips, PI) at the University of Maryland. Phoebe Gaston was also supported by a Flagship Fellowship from the University of Maryland and by NIH T32 DC017703 (I-M Eigsti & E. Myers, PIs) at the University of Connecticut.

Auditory word comprehension is less incremental in isolated words

Abstract

Speech input is often understood to trigger rapid and automatic activation of successively higher-level representations for comprehension of words. Here we show evidence from magnetoencephalography that incremental processing of speech input is limited when words are heard in isolation as compared to continuous speech. This suggests a less unified and automatic process than is often assumed. We present evidence that neural effects of phoneme-by-phoneme lexical uncertainty, quantified by cohort entropy, occur in connected speech but not isolated words. In contrast, we find robust effects of phoneme probability, quantified by phoneme surprisal, during perception of both connected speech and isolated words. This dissociation rules out models of word recognition in which phoneme surprisal and cohort entropy are common indicators of a uniform process, even though these closely related information-theoretic measures both arise from the probability distribution of wordforms consistent with the input. We propose that phoneme surprisal effects reflect automatic access of a lower level of representation of the auditory input (e.g., wordforms) while cohort entropy effects are task-sensitive, driven by a competition process or a higher-level representation that is engaged late (or not at all) during the processing of single words.

1 Introduction

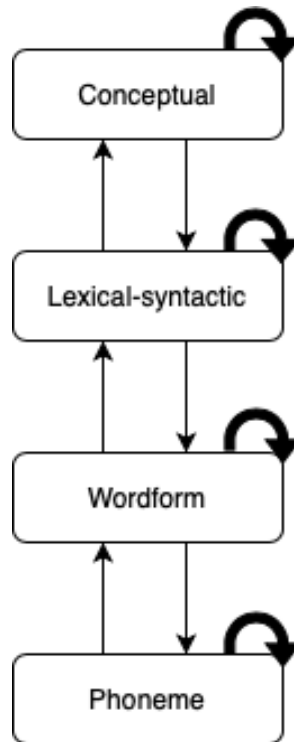
2 Speech recognition necessarily involves the access of multiple levels of representation in
3 response to auditory input, from phonemes to wordforms to higher-level lexical-syntactic
4 representations that link wordforms to meaning. While much about this process remains to be
5 elucidated, research on spoken word recognition has reached broad consensus on several points.
6 The contributions of a vast behavioral literature (reviewed by, e.g., Dahan & Magnuson (2006);
7 McQueen (2007); Magnuson, Mirman, & Myers (2013); Magnuson (2016)) indicate an
8 incremental, phoneme-by-phoneme process of winnowing down the phonological wordforms
9 that are consistent with the unfolding auditory input (e.g., Grosjean (1980); Zwitserlood (1989);
10 Allopenna et al. (1998); and following). Conceptual information associated with those
11 wordforms can be incrementally activated (e.g., Zwitserlood (1989); Yee & Sedivy (2006); and
12 following), and syntactic information is rapidly invoked (e.g., Marslen-Wilson & Tyler (1980);
13 McAllister (1988); and following). This process is highly sensitive to distributional statistics,
14 captured by word frequency (e.g., Connine et al. (1990); Dahan et al. (2001)).
15

16 The evidence leading to this consensus comes from a broad array of experimental approaches
17 that vary in which aspects of word recognition they can most effectively probe. These
18 approaches use stimuli that vary from sublexical phoneme sequences to natural, connected
19 speech. Combining evidence from these different paradigms is usually guided by an assumption
20 that there is a uniform, automatic progression of processing triggered by speech input, such that
21 we can expect datapoints from different points in that progression to cohere. Under this
22 assumption, simpler or single-word paradigms will straightforwardly capture the fundamental
23 word recognition sequence in isolation, while presenting more complex input allows us to
24 investigate how contextual information influences, for example, the speed of processing or the
25 set of lexical candidates under consideration.
26

27 In **Figure 1**, we sketch a representative sequence of processing proposed to occur in response to
28 each phoneme of speech input. TRACE (McClelland & Elman, 1986) is an example of a model
29 that is consistent with the illustrated principles. Each level of representation automatically
30 determines the most likely interpretation of the input through local competition and broadcasts
31 this interpretation through feed-forward and feed-back connections. The assumption of
32 automaticity implies that any speech input engages this processing hierarchy in the same manner.
33 The task context might change the information available at different levels, but not the basic
34 sequence of processing. However, if the assumption of automaticity is incorrect, then the basic
35 process of word recognition could deviate significantly according to the demands of different
36 comprehension scenarios. This deviation could occur because of variation in, for instance, the
37 relevance of different types of information to different experimental tasks, the ease of word
38 segmentation, and the degree to which word-to-word dependencies occur in the input.
39

40 In this paper we present neural evidence that word recognition in isolation may proceed in a
41 qualitatively different way from word recognition in continuous speech. Behavioral measures or
42 paradigms requiring an explicit response to each stimulus make comparison between isolated
43 words and continuous speech difficult, and a single trial generally reflects the status of just a
44 single item in the lexicon. Instead, we turn to a neural measure-- temporal response function
45 (TRF) analysis of magnetoencephalography (MEG) responses—that can be applied in exactly
46 the same way to single-word and continuous-speech listening, and that reflects distributional

1 properties of the entire class of word candidates consistent with each presented phoneme. We
2 show that the effects of two measures that have both been understood to reflect automatic
3 wordform-level processing in fact dissociate robustly according to the nature of the experiment.
4 This dissociation implicates a break in the automaticity of the sequence of activation and
5 indicates a difference between the processing of words presented in isolation and words
6 presented in continuous speech. Our findings have implications for the architecture of word
7 recognition models as well as for experimental approaches to studying speech perception.
8



9
10 **Figure 1.** Automatic sequence of processing assumed to occur in response to each phoneme of
11 speech input. Straight arrows indicate connections between levels of representation. Curved
12 arrows indicate a within-level competition/selection process.

13 **Phoneme surprisal and cohort entropy**

14 The neural response to speech has been shown to be modulated by information-theoretic
15 properties of the set of wordforms that match the auditory input at any given phoneme (Brodbeck
16 et al., 2018, 2021; Di Liberto et al., 2019; Donhauser & Baillet, 2020; Ettinger et al., 2014;
17 Gagnepain et al., 2012; Gaston & Marantz, 2018; Gillis et al., 2021; Gwilliams et al., 2020;
18 Gwilliams & Marantz, 2015; Kocagoncu et al., 2017). Two of these properties in particular –
19 cohort entropy and phoneme surprisal – have emerged as promising means of investigating the
20 time course of auditory word recognition.

21
22 Phoneme surprisal at a given phoneme reflects the conditional probability of that phoneme given
23 the preceding sequence of phonemes in the current word. Phoneme surprisal at position i in a
24 wordform is defined as $-\log_2 p(k_i | k_1, \dots, k_{i-1})$ where k_i is the phoneme at position i and $i = 1$
25 for the first phoneme in the wordform. Cohort entropy at that same phoneme, in contrast, is
26 determined by the probability distribution over wordforms that might complete that phoneme
27 sequence, reflecting how much uncertainty there is among those candidates. Cohort entropy at

1 position i in a wordform is defined as $-\sum_w^{C_i} (p(w | k_1, \dots, k_i) \times \log_2 p(w | k_1, \dots, k_i))$ where w is
2 each wordform in the cohort C_i of wordforms consistent with the sequence of phonemes k_1, \dots, k_i .
3 One of the critical differences between these formulations is that cohort entropy is forward-
4 looking in a way that phoneme surprisal is not. A cohort entropy effect reflects expectations for
5 potential candidates that would be consistent with the current input, while phoneme surprisal
6 provides evidence of the degree to which a phoneme could previously have been expected.

7
8 More neural activity is generally observed in response to higher surprisal, or lower probability,
9 phonemes, consistent with many cognitive domains in which predictable or higher probability
10 stimuli elicit reduced neural responses (see Aitchison & Lengyel (2017)). Exactly how cohort
11 entropy should be expected to drive neural activity is less clear. A larger set of candidates has a
12 higher cohort entropy than a smaller set of candidates, and a set of candidates in which
13 probability is equally distributed has a higher cohort entropy than a set of candidates in which
14 probability is concentrated on a single candidate. Greater uncertainty could be associated with
15 more neural activity due to an intensified process of lexical competition (Gagnepain et al., 2012),
16 or due to increased attentional gain on bottom-up input (Donhauser & Baillet, 2020), or it could
17 be that lower uncertainty is a precondition for other processes to be engaged (Ettinger et al.,
18 2014).

19
20 Despite these differences, phoneme surprisal and cohort entropy are often investigated and
21 presented in tandem as interchangeable indicators of wordform-level processing. One likely
22 reason for this approach in the literature is that the conditional phoneme probabilities underlying
23 both measures are calculated from the probabilities of wordforms consistent with the input. The
24 two variables are also often correlated, and their effects in neural data frequently co-occur.
25 Finally, in a hypothesized model of word recognition that includes automatic engagement of
26 successive representational levels regardless of task or context, phoneme surprisal and cohort
27 entropy effects are simply two different windows into the same automatic flow of activation
28 through the system.

29 **Variation in neural effects of cohort entropy and phoneme surprisal**

30 Despite frequently being treated interchangeably, a careful look at the prior literature reveals
31 considerable variation in whether, where, and when phoneme surprisal and cohort entropy effects
32 manifest across experiments. This variation has not previously been examined systematically.
33 Thus, before we proceed to our own study, we review this literature and consider whether there
34 are properties of the stimulus or experimental context that can help explain when cohort entropy
35 and phoneme surprisal effects do or do not occur, and what this might mean for the processes
36 and levels of representation they describe. An account of this variability is important for
37 improving the utility of phoneme surprisal and cohort entropy as measures for investigating
38 speech perception and specifically the class of active items in competition for recognition at any
39 given point in a word. However, understanding this variability also has the potential to illuminate
40 dissociable sub-processes in word recognition.

41
42
43 We begin by trying to characterize why these effects occur at all in some experiments and not in
44 others, though further efforts to understand variation in the localization and time course of these
45 effects will also be important. In Table 1, we summarize existing electrophysiology (primarily
46 MEG) studies that have tested for effects of cohort entropy and phoneme surprisal on neural

1 activity. Effects of both cohort entropy and phoneme surprisal have been reported in behavioral
 2 measures of auditory word recognition (Baayen et al., 2007; Balling & Baayen, 2012; Bien et al.,
 3 2011; Kemps et al., 2005; Wurm et al., 2006). However, testing for such effects in behavioral
 4 data generally requires constructing a cumulative measure of a phoneme-level variable across the
 5 course of the word or selecting the variable's value at just one phoneme position as the predictor.
 6 Therefore we restrict our focus here to neural measures that have the temporal resolution to
 7 examine cohort entropy and phoneme surprisal effects on a phoneme-by-phoneme basis. We
 8 exclude one additional study on the processing of continuous speech (Di Liberto et al., 2019),
 9 which did not report effects of cohort entropy and phoneme surprisal separately.

10
 11 **Table 1.** Properties of the Stimulus and Experimental Task for Existing Electrophysiology
 12 Studies Reporting Phoneme Surprisal or Cohort Entropy Effects

Study	Phoneme surprisal effect?	Cohort entropy effect?	Stimulus	Experimental task	Multimorphemic words included?
Gagnepain et al. (2012)	yes	no	single words	pause detection	no
Ettinger et al. (2014)	yes	yes [^]	single words	lexical decision	yes
Brennan et al. (2014)		no	single words	lexical decision	no
Lewis & Poeppel (2014)		no	single words	lexical decision	no
Gwilliams & Marantz (2015)	yes		single words	lexical decision	yes
Kocagoncu et al. (2017)		yes [^]	single words	nonword detection	not specified
Gaston & Marantz (2018)	yes	yes [*]	three-word phrases	phrase acceptability	yes
Brodbeck, Hong, et al. (2018)	yes	yes	continuous speech	comprehension questions	yes
Donhauser & Baillet (2020)	yes		continuous speech	comprehension questions	yes
Gwilliams et al. (2020)	yes	yes	continuous speech	comprehension questions	yes
Gillis et al. (2021)	yes	yes	continuous speech	comprehension questions	yes
Brodbeck et al. (2021)	yes	yes	continuous speech	comprehension questions	yes

13
 14 *Note.* Grey cells indicate studies that did not test for the specified effect. Superscripts indicate
 15 that a reported cohort entropy effect did not survive when phoneme surprisal was controlled for
 16 (*), or that such a test was not performed ([^]).
 17

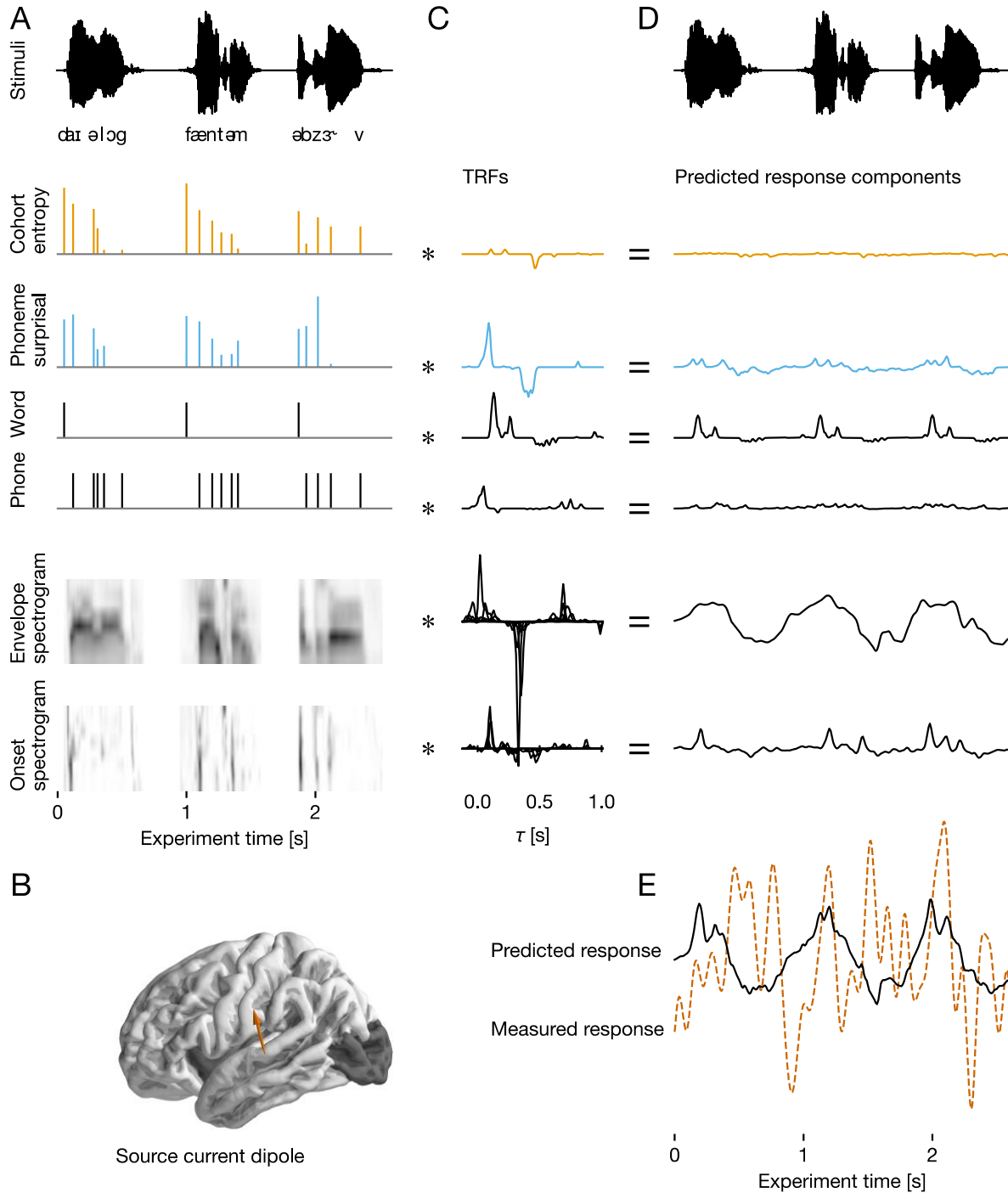
1 Table 1 demonstrates that phoneme surprisal and cohort entropy effects have very different
2 profiles across studies. Phoneme surprisal effects were reported in all studies that tested for them
3 (Brodbeck et al., 2018, 2021; Donhauser & Baillet, 2020; Ettinger et al., 2014; Gagnepain et al.,
4 2012; Gaston & Marantz, 2018; Gillis et al., 2021; Gwilliams et al., 2020; Gwilliams & Marantz,
5 2015), and thus appear to be robust to variation in stimulus and experimental task. Cohort
6 entropy, in contrast, produces mixed results. Among studies that presented single words and
7 short phrases, three reported cohort entropy effects (Ettinger et al., 2014; Gaston & Marantz,
8 2018; Kocagoncu et al., 2017) and three tested for but failed to find them (Brennan et al., 2014;
9 Gagnepain et al., 2012; Lewis & Poeppel, 2014). The presence or absence of multimorphemic
10 words in the study is potentially relevant, as the three studies that failed to find cohort entropy
11 effects included only monomorphemic words. However, more important in our view is that the
12 three single-word studies that reported cohort entropy effects did not exclude the possibility that
13 these effects were due to the highly correlated phoneme surprisal measure. Gaston and Marantz
14 (2018) in fact found that their significant cohort entropy effect was no longer significant in a
15 model that controlled for phoneme surprisal, and the other two studies (Ettinger et al., 2014;
16 Kocagoncu et al., 2017) did not conduct such a test. In continuous speech, cohort entropy effects
17 were reported in all studies that tested for them (Brodbeck et al., 2018, 2021; Gillis et al., 2021;
18 Gwilliams et al., 2020), with methods that controlled for effects of phoneme surprisal. We
19 conclude that, in the existing electrophysiology literature, there is strong evidence for phoneme
20 surprisal effects across the board, but for cohort entropy effects only in continuous speech.

21
22 A true dissociation between cohort entropy and phoneme surprisal effects would indicate not
23 only that these measures do not index the same level of representation or process, but also that
24 whatever drives cohort entropy effects does not occur during the processing of single words, or
25 at least does not occur incrementally (i.e., phoneme by phoneme). This is not consistent with all
26 processing steps being engaged in a fully automatic sequence during speech recognition.
27 However, this interpretation of the prior literature is complicated by the fact that many of these
28 studies did not control for potential confounds, such as acoustic variables and overlapping
29 responses to different phonemes. Differences in statistical power or analysis methods (which
30 vary widely) may also have contributed to the apparent influence of stimulus on cohort entropy
31 effects.

32 33 **The current study**

34 Hypothesizing that cohort entropy and phoneme surprisal do, indeed, dissociate, and that cohort
35 entropy effects do not occur for single words, we evaluated cohort entropy and phoneme
36 surprisal effects on the neural response to speech in a simple single-word paradigm and then
37 directly compared this data to an existing continuous-speech dataset (Brodbeck et al., 2021).
38 Comparing single-word and continuous-speech data requires that the two types of responses be
39 evaluated with the same method. Analysis techniques traditionally applied to single-word studies
40 are not suitable for responses to continuous speech, and generally fail to account for acoustic and
41 other confounding variables, as well as the overlapping nature of phoneme responses. Instead,
42 we modeled source-localized MEG data with temporal response functions (**Figure 2**), a method
43 that deals with acoustic confounds and was originally developed for continuous speech. This
44 allowed for novel comparison between single words and continuous speech as well as a more
45 accurate characterization of the single-word response relative to previous analyses.

46



1
 2 **Figure 2.** Temporal response function analysis. Brain activity was modeled as continuous
 3 response to multiple variables describing the sequence of words. (A) Predictor variables used to
 4 describe the stimuli were all represented as time series. Cohort entropy and phoneme surprisal
 5 were modeled as impulses at phoneme onset, scaled by the relevant quantity. Covariates
 6 included word and phoneme onsets, an 8-band auditory spectrogram, and an 8-band auditory
 7 onset spectrogram. (B) Neural activity was quantified as distributed minimum norm current
 8 estimates, i.e., estimated current at a grid of dipoles covering the cortical surface. The analysis
 9 was restricted to the temporal, frontal, and parietal lobes (the dark shading indicates regions

1 *excluded from the analysis*). *One dipole from one representative subject is used in this figure for*
2 *illustration. (C) Temporal response functions (TRFs) were estimated using a coordinate descent*
3 *algorithm to predict the neural signal from the predictor variables. (D) TRFs were estimated*
4 *jointly, i.e., each TRF, convolved with its corresponding predictor variable time series, predicted*
5 *a component of the neural activity. The sum of these component responses is the predicted brain*
6 *response (E). Model performance was evaluated by the proportion of the variability in the*
7 *measured response that was explained by the predicted response.*

8 Participants heard a list of 1000 monomorphemic words with an inter-stimulus interval of 267
9 ms, and responded to randomly occurring semantic relatedness probes. Models were fit using 5-
10 fold cross-validation in each subject separately. We evaluated the models by the proportion of
11 variability they explained in the source-localized MEG recordings, correcting for multiple
12 comparisons using threshold-free cluster enhancement (Smith & Nichols, 2009). Unless noted
13 specifically, analyses were performed on the surface of the temporal, frontal and parietal lobes
14 combined (see shaded area in **Figure 2B**).

15 **Materials & Methods**

16 **Participants**

17 We collected MEG data from 24 people. Sample size was chosen in accordance with the
18 previous studies cited in Table 1. All participants were right-handed, native speakers of English,
19 and seven were also native speakers of additional languages. None reported a history of
20 neurological or linguistic impairment, brain injury, or hearing loss. All reported normal or
21 corrected-to-normal vision. The procedure was approved by the University of Maryland
22 Institutional Review Board and all participants provided written informed consent. Participants
23 were compensated with their choice of \$15 or 1 course credit per hour of participation. The full
24 session (including another, unrelated study) lasted 2 hours.

25
26 One dataset was excluded before data processing because of participant fatigue and an earbud
27 falling out during the experiment. After this exclusion, we computed accuracy on the semantic
28 relatedness task and excluded any participant with accuracy lower than a cutoff one standard
29 deviation below the mean. This excluded three of 23 participants. After preprocessing, two
30 additional datasets were excluded due to excessive magnetic noise. 18 datasets are therefore
31 included in our analysis.

32 **Stimuli**

33
34 Our stimuli were word recordings from the Massive Auditory Lexical Decision (MALD)
35 database (Tucker et al., 2019), which includes the timing of phoneme boundaries from a forced
36 aligner. The set of 1000 words we selected had no missing variables in the database and were
37 monomorphemic per MALD, CELEX (Baayen et al., 1995), and first author judgment. We
38 excluded all items with the following labels in MALD: Preposition, Interjection, Name,
39 Unclassified, Conjunction, Pronoun, Determiner, Letter, Not, Ex, Article, To. We also removed
40 items with the 10% lowest frequency values, and excluded homophones, inappropriate and
41 particularly evocative words, and any item for which the pronunciation in the recording was
42 noticeably divergent from American English. The full lists of stimuli and semantic relatedness
43 probes (see below), as well as associated stimulus variables from MALD, are available on OSF
44 (<https://osf.io/u56ea/>).

1 **Procedure**

2 The study was always the second of two experiments in a session. Before the MEG recording,
3 we used a Polhemus 3SPACE FASTRAK to digitize participant head shapes as well as the
4 positions of five affixed marker coils. These marker coils were used to record head position
5 relative to the MEG sensors before and after each study in the session. We recorded continuous
6 MEG data, inside a magnetically shielded room, with a 160-channel axial gradiometer whole-
7 head system (Kanazawa Institute of Technology, Kanazawa, Japan). Our sampling rate was 1000
8 Hz, and we used an online 60 Hz notch filter and 200 Hz low-pass filter.
9

10 Participants lay supine and looked at a screen overhead, while holding a button box in each hand.
11 They wore foam earbuds and volume was adjusted to their comfort level. We instructed
12 participants that they would hear a long series of random words, and that they should simply
13 listen to the words while watching for probe words that would randomly appear on the screen
14 with a question mark. They were instructed to press a button (with left hand for No and right
15 hand for Yes) to indicate whether the word on the screen was related in any way to the word they
16 had heard just before it.
17

18 We used Presentation (Neurobehavioral Systems, Inc., www.neurobs.com) to present the
19 experiment. Our parameter and scenario files are available on OSF (<https://osf.io/u56ea/>). There
20 were 1000 auditory trials interspersed pseudo-randomly with 97 semantic relatedness probe
21 trials. The amount of time between trials was 267 ms. A visual fixation cross was on screen
22 continuously during auditory trials and during the inter-trial interval. Each auditory trial simply
23 consisted of presentation of the auditory stimulus and lasted the length of the auditory stimulus.
24 Visual probe trials were pseudo-randomly distributed with a maximum interlude of 20 trials
25 between probes. The probe (e.g., “podium?”) stayed on the screen until the participant pressed a
26 button to answer.
27

28 We selected this task so that it would apply equally well to all types of words, and because we
29 did not want button presses to occur on critical trials (as would happen in, e.g., lexical decision).
30 The probe trials for which we expected participants to answer “No” were selected randomly from
31 the list of eligible words that we did not end up using for auditory trials. Probe trials for which
32 we expected participants to answer “Yes” were synonyms taken from the WordNet
33 (<https://wordnet.princeton.edu>) page of the preceding auditory item and were also
34 monomorphemic so as not to be trivially distinguishable from “No” trials. There was no overlap
35 between probe words and words used in auditory trials. Which auditory trials would be followed
36 with a probe were randomly selected. “Yes” and “No” probes were equally distributed.
37

38 The experiment lasted roughly 17 minutes. There was no built-in break, but participants were
39 instructed that if they wished to take a break, they should simply delay their button press on a
40 probe trial.
41

42 **Data preprocessing**

43 We processed the data using mne-python version 0.22 (Gramfort et al., 2013, 2014) and Eelbrain
44 0.34 (Brodbeck et al., 2019).
45

1 During file conversion with mne-python’s kit2fiff GUI, we excluded any faulty marker
2 measurements. We co-registered each digitized head shape with the Freesurfer (Fischl, 2012)
3 “fsaverage” brain, using mne-python’s co-registration GUI. We first used rotation and translation
4 to align the digitized head shape and average MRI by the three fiducial points. We then used
5 rotation, translation, and 3-axis scaling to minimize the distance between digitized head shape
6 and average MRI points using the iterative closest point (ICP) algorithm. Convergence was
7 always achieved within 40 iterations. For one participant, outlying points on the digitized head
8 shape were removed between fitting to the fiducials and applying ICP.

9
10 Flat channels were automatically removed, and we used temporal signal space separation (Taulu
11 & Simola, 2006) for removal of extraneous artifacts, with a buffer duration of 10 seconds. We
12 then band-pass filtered the recordings between 1 – 40 Hz (mne-python default settings) and used
13 ICA (independent components analysis, with extended infomax method) for removal of ocular,
14 cardiac, and other extraneous artifacts. Components were selected manually based on their
15 topography and time-course. After removing artifactual ICA components, we further low-pass
16 filtered the data at 20 Hz, cropped it from 1 s before the first word to 2 s after the last word, and
17 down-sampled it to 100 Hz.

18
19 To compute a noise covariance matrix, we used two minutes of empty room data recorded before
20 or after each session. We defined the source space on the white matter surface with a four-fold
21 icosahedral subdivision, with 2562 sources per hemisphere. Orientation of the source dipoles was
22 fixed perpendicular to the white matter surface. Continuous data were source localized with the
23 regularized minimum norm estimator ($\lambda = 1/6$).

24 25 **Analysis**

26 ***Behavioral data***

27 Mean accuracy was computed after the exclusion of one participant a priori. The mean number
28 of correct probe responses was 73.6 (out of 97) with a standard deviation of 18.4. The number of
29 correct probe responses was lower than one standard deviation below the mean for three
30 participants, so they were excluded from further analysis. One participant answered 13 of 97
31 probes correctly. We kept this participant in the dataset because this was so far below chance that
32 the only plausible explanation seemed to be that they had reversed which hand they were
33 supposed to use to make Yes and No responses.

34 35 ***Predictors for neural data***

36 For each stimulus variable of interest, a time series was created indicating the value of the
37 predictor at each time point in the stimulus and aligned to the MEG data. For the acoustic
38 predictors, the value of the predictor can vary continuously at each time point. For linguistic
39 predictors, values are non-zero only at time points labeled as phoneme onsets, as determined by
40 the forced alignments (i.e., linguistic predictors consist of impulses at phoneme onsets). Of these
41 lexical predictors, the phoneme onset and word onset predictors each consist of binary impulses,
42 while the predictors for cohort entropy and phoneme surprisal consist of impulses that are scaled
43 continuously according to the variable value at that phoneme. Our study did not actually present
44 a continuous stimulus (rather, individual words with short intervening pauses), but a single time
45 series reflecting predictor values (or pauses) throughout the entire experiment could still be

1 created. Probe trials were modeled simply as silence. The timing of phoneme onsets was taken
2 from the forced aligner information made available with the MALD recordings.

3
4 *Acoustic spectrogram*: a gammatone spectrogram (Heeris, 2018) was computed for each stimulus
5 waveform with 256 channels regularly spaced in ERB space between 20 and 5000 Hz. These
6 spectrograms were resampled to 100 Hz to match the MEG data and binned into eight equally
7 spaced frequency bands.

8
9 *Acoustic onset spectrograms*: The high resolution gammatone spectrograms were processed with
10 an algorithm for acoustic edge extraction (Brodbeck et al., 2020; Fishbach et al., 2001). The
11 onset spectrograms were also resampled to 100 Hz and binned into eight bands.

12
13 *Word onsets*: Word onsets were represented as a single, equally valued impulse at the onset of
14 every word, as determined from the forced alignments.

15
16 *Phoneme onsets*: Phoneme onsets, including only phonemes that were not also word onsets, were
17 represented as equally valued impulses on a single predictor time series.

18
19 *Phoneme surprisal and cohort entropy*: these variables were calculated based on an
20 implementation of the cohort model of word perception (Marslen-Wilson, 1987), as in Brodbeck,
21 Hong, and Simon (2018). Initially, a dictionary was created combining frequency information
22 from SUBTLEX (Brysbaert & New, 2009) with pronunciations (phoneme sequences) from the
23 CMU pronouncing dictionary (Weide, 1994), adding any pronunciations from the stimuli that
24 were missing from the CMU dictionary. This dictionary was then used to compute the set of
25 words compatible with the input so far for each word at each phoneme position. These cohorts,
26 together with the SUBTLEX frequencies, were used to compute a probability distribution over
27 possible words for each phoneme position. The *cohort entropy* predictor contained an impulse at
28 each phoneme onset, scaled by the entropy of that cohort. The *phoneme surprisal* predictor
29 contained an impulse at each phoneme onset scaled by the surprisal of that phoneme, based on
30 the posterior probability of that phoneme given the preceding phoneme's cohort.

31 **TRF analysis**

32
33 A multivariate temporal response function (mTRF) maps a set of predictor variables to a single
34 outcome time series. Here, independent mTRFs were estimated for each subject and for each
35 virtual current source of source-localized MEG data (see **Figure 2**). The neural response at time
36 t , y_t is predicted jointly from N predictor time series, represented as $x_{i,t}$, convolved with a
37 corresponding mTRF $h_{i,\tau}$ of length T :

38

$$\hat{y}_t = \sum_i^N \sum_{\tau}^T h_{i,\tau} \cdot x_{i,t-\tau}$$

39 mTRFs were generated from a basis of 50 ms wide Hamming windows centered at $T=[-$
40 $100, \dots, 1000)$ ms. All responses and predictors were standardized by centering and dividing by
41 the mean absolute value.

42
43 For a given set of predictors, the predictive power was estimated through 5-fold cross-validation.
44 The continuous data and corresponding predictors were split into five contiguous partitions of

1 even length. The neural responses of each partition were predicted with an mTRF trained on the
2 remaining four partitions to minimize ℓ_1 error. Within each set of four training partitions, each
3 partition in turn served as validation data once as four mTRFs were estimated based on
4 coordinate descent with early stopping based on the validation data (David et al., 2007). The
5 validation data was used to selectively stop training predictors when they caused an increase in
6 error in the validation set. Those four mTRFs were then averaged to predict the responses to the
7 unseen (fifth) test segment.

8
9 For evaluating the predictive power of the relevant predictors, phoneme surprisal and cohort
10 entropy, we compared the predictive power of the full model with that of a model that was
11 identical except for not including the predictor under investigation. Together with the cross-
12 validation, this assures a conservative estimate of the unique predictive power of the predictor
13 under investigation, while controlling for the predictive power of all the other variables. The
14 anatomical maps of explanatory power of the two models were compared with a mass-univariate
15 related measures t -test based on threshold-free cluster enhancement (TFCE) (Smith & Nichols,
16 2009), with a null distribution based on 10,000 random permutations of condition (model) labels.

17
18 For analysis of the TRFs, the five estimates of the TRFs from the five different test partitions
19 were averaged in each subject. In order to visualize the TRF current over time, the TRF was
20 restricted to the anatomical area in which the surprisal predictor significantly improved
21 predictions ($p \leq .05$ based on TFCE). Within this area, and separately for each hemisphere and
22 each participant, principal component analysis was applied to the virtual current dipoles, and
23 only the first principal component was analyzed, i.e., a single spatial topography and
24 corresponding time course for each participant. The advantage of this approach over others, such
25 as root mean squared activity, is that the signed current direction can be visualized. Because the
26 sign of a principal components is arbitrary, the components were aligned across subjects such
27 that the average current vector was pointing upward. For components whose average current
28 vector pointed downward, both component and time-course were multiplied by -1.

29
30 TRF time-course was then evaluated in each hemisphere using a mass-univariate one-sample t -
31 test with TFCE, with the null hypothesis that the average current direction is random (i.e., not
32 different from 0). The null distribution was based on the maximum statistic in 10,000 random
33 permutations of the signs. To test for hemispheric differences, a mass-univariate repeated
34 measures t -test with the same parameters was used.

35 ***Comparison with connected speech***

36
37 For this comparison, data from 12 participants listening to 47 minutes of a non-fiction audiobook
38 were used (for more details see Brodbeck et al. (2021)). Data were acquired on the same MEG
39 equipment and with analogous procedures, with one exception: For estimation of the mTRF
40 models, data were split into four partitions instead of five. This was done to speed up
41 computations (requiring training of fewer models) and because the longer recording resulted in
42 more training data per participant. Audiobook stimuli were labeled using the Montreal Forced
43 Aligner (McAuliffe et al., 2017), and predictor variables were generated as for the single-word
44 data.

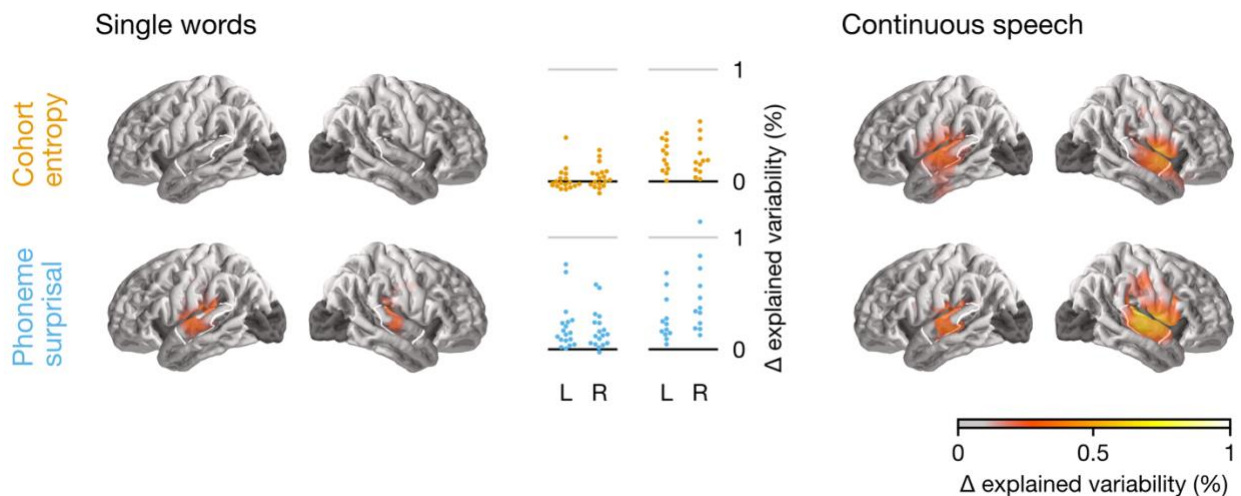
1 Results

2 To ensure that responses reflect attentive lexical processing, we applied behavioral exclusion
3 criteria (see Methods). Subjects included in the analyses presented here responded accurately to
4 at least 69% of relatedness probes (group mean 82.9%).

5 Our first question was whether first phonemes should be excluded from the phoneme surprisal
6 and cohort entropy estimates, as suggested by a lack of first phoneme cohort effects reported by
7 Brodbeck, Hong, and Simon (2018) and Gaston and Marantz (2018). To answer this question, we
8 compared the model treating all phonemes uniformly (**Figure 2**) to a model in which surprisal
9 and entropy at the first phoneme are modeled as separate predictors from surprisal and entropy at
10 non-initial phonemes. The more complex model, in which they are modeled separately, was not
11 significantly better ($t_{max} = 2.74, p = .341$, multiple comparison correction in temporal lobes
12 only). We therefore proceeded with the simpler model in which initial phonemes are not
13 modeled separately (as shown in **Figure 2**). Acoustic and segmentation variables were always
14 controlled for (see Methods).

15
16 We then tested our primary question: do phoneme surprisal and cohort entropy improve the
17 estimated neural response in a single-word design? We found that indeed, a model with phoneme
18 surprisal was significantly better than a model without it ($p < 0.001$). However, comparison with
19 a model lacking cohort entropy led to no significant difference ($p = .260$, see **Figure 3**). The
20 model improvement due to surprisal (i.e., the explanatory power of surprisal) was significantly
21 larger than that due to entropy ($p = .007$).

22



23
24 **Figure 3.** Model evaluation and comparison to continuous speech. The anatomical plots at left
25 and right show regions where the given predictor significantly improved the model fit. The white
26 outline indicates an anatomical region of interest (ROI) defined as the posterior 2/3 of the
27 superior temporal gyrus. The swarm plots (middle) show average proportion of variability in
28 that ROI that is uniquely explained by entropy or surprisal, respectively. Each dot represents one
29 participant. While surprisal improves the model fit in both experiments in almost all
30 participants, entropy does so only in the continuous-speech data. Explained variability
31 (explanatory power) is expressed as percentage of the maximum variability explained by the full
32 model in the single-word data.

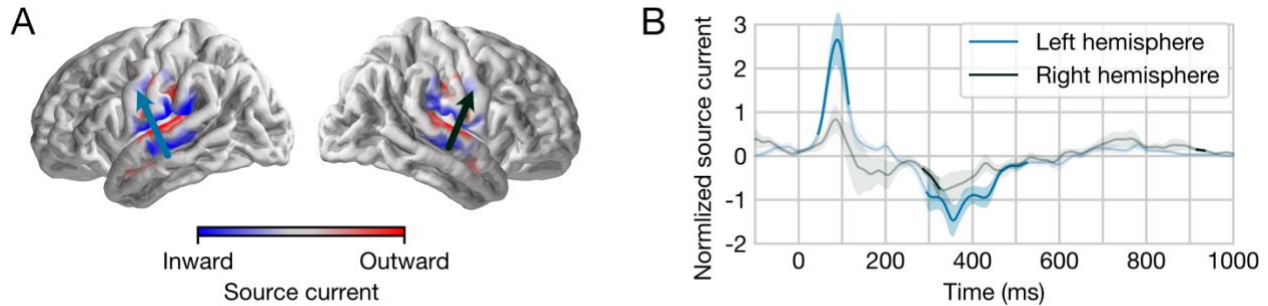
33

1 This finding contrasts with previously reported results in continuous speech (see Table 1). To
2 address this apparent difference, we compared our single-word data to an existing continuous-
3 speech dataset acquired with the same MEG scanner (Brodbeck et al., 2021), consisting of
4 recordings from 12 participants listening to 47 minutes of an audiobook, using closely matched
5 analysis methods (**Figure 3**). For the continuous-speech data, phoneme surprisal significantly
6 improved the model ($p < .001$) and cohort entropy did as well ($p < .001$). In the whole brain
7 analysis, the explanatory power of phoneme surprisal and cohort entropy did not differ
8 significantly ($p = .720$).

9
10 To confirm this difference between experiments statistically, we extracted the mean of the model
11 fit metric in a region of interest (ROI) defined as the posterior two thirds of the superior temporal
12 gyrus of each hemisphere. This value did not differ between the left and right hemisphere ROIs
13 in any of the four categories (surprisal/entropy, single words/continuous speech; all $t \leq 1.74$, $p \geq$
14 $.110$), so we averaged the values for the two hemispheres. We then calculated the ratio between
15 the predictive power of entropy and surprisal and compared this ratio for continuous speech and
16 single words. Such a test is unlikely to be affected by differences in the size of the datasets. This
17 ratio was significantly higher in continuous speech than for single words (continuous speech $M =$
18 0.68 , $SD = 0.45$; single words $M = 0.10$, $SD = 0.59$; $t_{28} = 2.80$, $p = .009$). Consistent with this,
19 effect sizes for predictive power in the ROI were large for surprisal in both paradigms (single
20 words: $d = 1.62$; connected speech: $d = 2.14$) but for entropy only in connected speech ($d = 1.72$)
21 and not in single words ($d = 0.39$).

22
23 Finally, we examined the nature of the estimated response functions for phoneme surprisal in the
24 single-word dataset (**Figure 4**). The TRF analysis was restricted to a mirror-symmetric
25 anatomical region, based on the area in which surprisal significantly improved the model fit in at
26 least one hemisphere. Because the TRF was relatively well captured by a single topography, we
27 extracted only the first principal component of the TRF for each participant and each hemisphere
28 (a parallel analysis using the response magnitude led to the same conclusions). **Figure 4A** shows
29 the average of the first principal component for each subject. The result in both hemispheres was
30 consistent with a current dipole in auditory cortex, indicated by the arrows in **Figure 4A**. The
31 time-course in the two hemispheres (**Figure 4B**) was analyzed with mass-univariate t -tests,
32 correcting for the time range from 0 to 1000 ms. In both hemispheres, an early peak around 90
33 ms was followed by more extended activity of opposite current direction, starting around 280 ms.
34 Even though activity in the early peak did not reach significance in the right hemisphere, the
35 difference between hemispheres, based on a mass-univariate related measures t -test, was not
36 significant ($p = .063$, at 70 ms).

37
38
39
40



1
2 **Figure 4.** TRF results for phoneme surprisal. The TRF is analyzed using the first principal
3 component in each subject. (A) The average first principal component across subjects. The
4 average current direction (indicated by arrows) is consistent with auditory cortex activity. (B)
5 The time-course of the component separately for the left and the right hemisphere. Solid line
6 segments indicate time ranges in which the respective TRF is significantly different from zero.

7 Discussion

8 This study examined cohort entropy and phoneme surprisal effects in a single-word paradigm
9 using a temporal response function analysis, modeling both acoustic and linguistic predictors of
10 neural activity. We found that phoneme surprisal is a significant predictor of neural activity
11 during speech recognition, as have many previous studies (Brodbeck et al., 2018, 2021;
12 Donhauser & Baillet, 2020; Ettinger et al., 2014; Gagnepain et al., 2012; Gaston & Marantz,
13 2018; Gillis et al., 2021; Gwilliams et al., 2020; Gwilliams & Marantz, 2015). The spatial
14 distribution of the effect along the superior temporal gyrus is also consistent with previous work.
15 The TRF for phoneme surprisal in our study appears to peak twice, in line with Gwilliams and
16 Marantz (2015), Gaston and Marantz (2018), and Brodbeck et al. (2021).

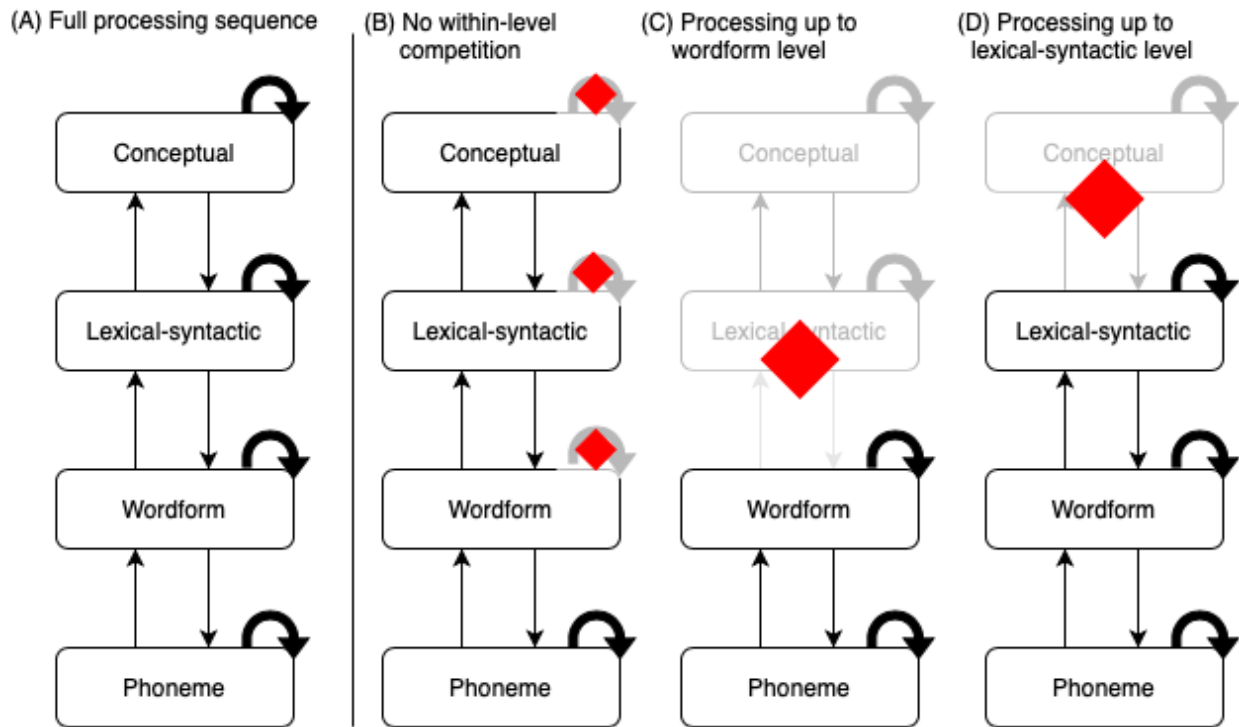
17
18 In contrast to the robust effect of phoneme surprisal, we did not observe a significant effect of
19 cohort entropy. In a direct comparison to our single-word dataset, we analyzed a continuous-
20 speech dataset (Brodbeck et al., 2021) in the same manner, and found effects of both phoneme
21 surprisal and cohort entropy. The direct comparison of these two datasets substantiates our
22 generalization about the existing literature, that cohort entropy effects are weak or non-existent
23 in studies that use single words or short phrases, while they are robust in studies that use
24 continuous, naturalistic speech as stimuli.

25
26 How could this dissociation between phoneme surprisal and cohort entropy occur? As reviewed
27 in the Introduction, it is frequently assumed that speech input triggers relatively automatic and
28 incremental activation of phoneme, wordform, lexical-syntactic, and conceptual units, but this
29 would predict cohort entropy effects for any task involving word recognition. In the following
30 sections, we hypothesize (1) that brain responses related to phoneme surprisal and cohort entropy
31 arise from different levels of representation or different sub-processes and (2) that their
32 dissociation therefore implies a break in the automatic sequence of processing involved in word
33 recognition.

34 35 **Non-automaticity in the lexical access sequence**

36 The pattern of dissociation that we observed could have several different explanations,
37 contingent on the precise neural processes indexed by cohort entropy and phoneme surprisal. In

1 **Figure 5A**, we reproduce our illustration in **Figure 1** of a fully automatic processing sequence in
2 response to each phoneme of speech input. In **Figure 5B-D**, we illustrate alternative partial
3 versions of this sequence that might better represent what occurs incrementally in single-word
4 paradigms that do not elicit cohort entropy effects. It is possible that the decoupled processes do
5 not occur at all in single-word processing; alternatively, they could be engaged much later or
6 engaged in a less strictly incremental, time-locked manner rather than on a phoneme-by-
7 phoneme basis.
8



9
10 **Figure 5.** (A) Fully automatic processing sequence in which both phoneme surprisal and cohort
11 entropy effects arise. (B)-(D) Proposed partial processing sequences in which phoneme surprisal
12 but not cohort entropy effects occur. Red diamonds indicate processes or levels of representation
13 that might be delayed or suspended from incremental (phoneme-locked) processing during
14 recognition of single words. As in **Figure 1**, straight arrows indicate connections between levels
15 of representation. Curved arrows indicate a within-level competition/selection process.

16
17 One possible explanation is based on the reasoning that cohort entropy is specifically a measure
18 of the amount of lexical competition occurring (Gagnepain et al., 2012). We can imagine a
19 scenario in which initial activation of multiple lexical candidates is automatic, but in which the
20 competitive process of winnowing out the weaker ones is only applied when rapid selection of a
21 single best candidate is particularly helpful or necessary for the task at hand. Accordingly,
22 phoneme surprisal effects might require only *activation* of, e.g., the wordform level of
23 representation, rather than the competition process that occurs within that level (a scenario
24 illustrated by **Figure 5B**, in which within-level competition processes are not occurring above
25 the phoneme level). In contrast, cohort entropy effects would reflect the competitive selection
26 process that allows a single best candidate to be identified as early as possible, and this process
27 might only be engaged when faced with the time pressure of processing connected speech.

1 Another possibility is that phoneme surprisal and cohort entropy effects reflect different levels of
2 representation which are not all automatically accessed to the same degree. Access to 'lower'
3 levels of representation like phoneme or wordform representations might be more automatic,
4 whereas access to 'higher' levels of representation like lexical-syntactic or conceptual units
5 might be more dependent on context and task demands. For instance, surprisal effects might
6 require only wordform-level activation, while cohort entropy effects might require lexical-
7 syntactic activation or higher. Similarly, phoneme surprisal effects could implicate up to lexical-
8 syntactic activation while cohort entropy effects require conceptual activation or higher. These
9 two possibilities are illustrated in **Figure 5C** and **Figure 5D**, respectively. Consistent with such
10 an explanation, semantic priming from partial wordforms seems to be more reliable in connected
11 speech (Zwitserslood, 1989) than in a single-word lexical-decision paradigm, where form-based
12 priming predominates (Gaskell & Marslen-Wilson, 2002). Even within a single-word paradigm,
13 Bentin et al. (1993) argue that the extent to which a task requires semantic processing can
14 influence the degree of semantic priming that occurs, as indexed by the N400 response.

15
16 Though less likely, we can acknowledge two alternative explanations in which phoneme
17 surprisal effects are not actually related to wordform representations analogous to the ones used
18 to calculate phoneme surprisal. We consider these less likely because they would imply an
19 absence of incremental wordform-level processing in the single-word tasks, despite behavioral
20 evidence to the contrary. One possibility is that apparent phoneme surprisal effects arise due to
21 prelexical phonotactic processing, involving representations sensitive to the probability of
22 phoneme sequences in the language independent of wordform representations. The second
23 possibility arises from the proposal of Norris & McQueen (2008) that 'off-line' perceptual
24 learning could lead to wordform frequency effects on phoneme probability without concurrent
25 wordform activation causing online top-down effects. In either of these scenarios, a phoneme
26 surprisal effect does not necessarily imply wordform activation; cohort entropy could reflect
27 anything at the wordform level or above. Similarly, it remains possible that correlations between
28 neural activity and cohort entropy are not driven by lexical competition or uncertainty per se but
29 by a secondary process that is sensitive to lexical competition or uncertainty. If that process is
30 only engaged by continuous speech, cohort entropy effects would also appear to be modulated.

31 32 **Single words vs. continuous speech**

33 What are the differences between single-word paradigms and continuous speech that would
34 make any of the distinctions described in the previous section possible? First, the reliable
35 presence of pauses between words in single-word studies may constitute a key change in task
36 demands, by leaving sufficient time for full lexical access to occur after wordform offset and
37 before the next wordform begins and thus reducing the necessity for incremental processing.
38 Early competitive selection might be unnecessary, and/or access to higher-level syntactic and
39 conceptual units could be deferred until the pause makes the auditory wordform uniquely
40 identifiable. Among the single-word studies we have reviewed, the pause detection (Gagnepain
41 et al., 2012) and nonword detection (Kocagoncu et al., 2017) tasks incorporate lengthy inter-
42 stimulus intervals averaging 2000 ms, and the lexical decision studies (Brennan et al., 2014;
43 Ettinger et al., 2014; Gwilliams & Marantz, 2015; Lewis & Poeppel, 2014) wait for a participant
44 response after each word. Our study used a shorter but still considerable inter-stimulus interval
45 of 267 ms with only occasional semantic relatedness probes and also did not find a cohort
46 entropy effect.

1 Second, the syntactic and semantic structure in continuous speech provides another motivation
2 for incremental processing: rapid access to lexical and conceptual content for the current word
3 provides information that might aid recognition of the subsequent word. This rationale for rapid
4 processing is absent in single-word paradigms that lack structure. Even beyond not requiring
5 speed in lexical or conceptual access, the tasks employed in single-word paradigms may in some
6 cases not require lexical or conceptual access at all. For instance, our task involved semantic
7 relatedness judgements with written probes. It is conceivable that this task might be solved
8 successfully by temporarily ‘buffering’ the input from each word as a form-based representation,
9 and only accessing conceptual representations if a probe occurs. By contrast, the speed of
10 continuous speech, its many between-word dependencies, and the imperative to build sentence-
11 level and message-level interpretations could be what drive competition or incremental higher-
12 level activation (and therefore cohort entropy effects) in naturalistic paradigms.

13
14 We might expect that cohort entropy effects could be observed for single words if a task were
15 designed such that earlier identification of the word is encouraged and incremental higher-level
16 activation becomes more advantageous, whether via the elimination of pauses or the addition of
17 some higher-level structure. Likewise, pauses could be added to continuous speech. The three-
18 word phrases used by Gaston and Marantz (2018) (e.g., “to chew gum,” “the shredder broke”)
19 are an interesting test of these hypotheses, as they lack within-phrase pauses and have syntactic
20 and semantic structure. Nevertheless, Gaston and Marantz did not find cohort entropy effects
21 when their cohort entropy variables were evaluated in the same model as phoneme surprisal.
22 This suggests that only longer sequences of continuous speech elicit cohort entropy effects, and
23 therefore that a buffering process may play a mediating role here.

24
25 Another line of investigation for understanding what drives neural cohort effects might involve
26 the contrast between monomorphemic and multimorphemic words. The two types of stimuli can
27 be closely matched in length, but only multimorphemic words can be viewed as structured
28 sequences of units of meaning of the kind that might encourage more incremental processing.
29 The inclusion of multimorphemic words in a single-word study could thus motivate earlier
30 selection and higher-level activation so that initial morphemes can be recognized in time to begin
31 processing any potential subsequent morphemes. In Table 1, we noted that all single-word
32 studies that do not include multimorphemic words also do not report cohort entropy effects
33 (Brennan et al., 2014; Gagnepain et al., 2012; Lewis & Poeppel, 2014). This is true of our study
34 as well. Among the single-word studies that do report cohort entropy effects, albeit without
35 controlling for phoneme surprisal, Ettinger et al. (2014) include multimorphemic words and
36 Kocagoncu et al. (2017) do not indicate whether multimorphemic words are included in their
37 stimuli or not. This factor deserves further investigation. In particular, it is unclear whether
38 hypothesized cohorts should be constructed on a morpheme-by-morpheme or wordform-by-
39 wordform basis.

40 41 **Implications**

42 If auditory word recognition in most single-word studies proceeds in the manner we have
43 proposed, with candidate selection or higher-than-wordform-level processing delayed or
44 suspended entirely, there are two major implications. The first is that the cascading, incremental
45 access process is not automatic but rather is motivated by time pressure and modulable with the
46 extent of that time pressure. The second is that auditory word recognition in many single-word

1 studies may differ fundamentally from the process most researchers assume they are studying
2 (that is, speech recognition in natural contexts). This would invite re-interpretation of existing
3 neural and behavioral data and would motivate increased use of more naturalistic designs in
4 future work, or identification of changes to single-word paradigms that would drive cohort
5 entropy effects so that these paradigms can be used with more confidence that they are
6 representative of the processing of natural, connected speech.

7 8 **Conclusion**

9 Our goal in this study was to evaluate whether an assumption of parallelism is warranted
10 between recognition of single words and word recognition in natural connected speech. We also
11 intended to establish a better understanding of what drives phoneme surprisal and cohort entropy
12 effects, while modeling the speech stimulus as thoroughly as current methods allow. We directly
13 compared single-word and continuous-speech data from MEG and demonstrated the occurrence
14 of phoneme surprisal effects in both paradigms but a cohort entropy effect only in continuous
15 speech, consistent with patterns in the existing literature. We proposed that this is because
16 phoneme surprisal effects arise from the activation of a lower level of representation while
17 cohort entropy effects arise from a competition process or higher level of representation whose
18 engagement is delayed or does not occur in single-word paradigms. This dissociation suggests
19 that the sequence of processing triggered by speech input is not automatic and the extent to
20 which competition processes or higher levels of representation are engaged depends on the
21 nature of the stimulus or experimental task. This study has also helped validate the TRF
22 approach as a promising method for future work in single-word paradigms.

23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

1 References

- 2 Aitchison, L., & Lengyel, M. (2017). With or without you: Predictive coding and Bayesian
3 inference in the brain. *Current Opinion in Neurobiology*, *46*, 219–227.
4 <https://doi.org/10.1016/j.conb.2017.08.010>
- 5 Allopenna, P. D., Magnuson, J. S., & Tanenhaus, M. K. (1998). Tracking the Time Course of
6 Spoken Word Recognition Using Eye Movements: Evidence for Continuous Mapping
7 Models. *Journal of Memory and Language*, *38*(4), 419–439.
8 <https://doi.org/10.1006/jmla.1997.2558>
- 9 Baayen, H., Piepenbrock, R., & Gulikers, L. (1995). *CELEX2 LDC96L14 [Web Download]*.
10 Linguistic Data Consortium.
- 11 Baayen, H., Wurm, L. H., & Aycocck, J. (2007). Lexical dynamics for low-frequency complex
12 words: A regression study across tasks and modalities. *The Mental Lexicon*, *2*(3), 419–
13 463. <https://doi.org/10.1075/ml.2.3.06baa>
- 14 Balling, L. W., & Baayen, H. (2012). Probability and surprisal in auditory comprehension of
15 morphologically complex words. *Cognition*, *125*(1), 80–106.
16 <https://doi.org/10.1016/j.cognition.2012.06.003>
- 17 Bentin, S., Kutas, M., & Hillyard, S. A. (1993). Electrophysiological evidence for task effects on
18 semantic priming in auditory word processing. *Psychophysiology*, *30*(2), 161–169.
19 <https://doi.org/10.1111/j.1469-8986.1993.tb01729.x>
- 20 Bien, H., Baayen, R. H., & Levelt, W. J. M. (2011). Frequency effects in the production of Dutch
21 deverbal adjectives and inflected verbs. *Language and Cognitive Processes*, *26*(4–6),
22 683–715. <https://doi.org/10.1080/01690965.2010.511475>
- 23 Brennan, J., Lignos, C., Embick, D., & Roberts, T. P. L. (2014). Spectro-temporal correlates of
24 lexical access during auditory lexical decision. *Brain and Language*, *133*, 39–46.
25 <https://doi.org/10.1016/j.bandl.2014.03.006>
- 26 Brodbeck, C., Bhattasali, S., Cruz Heredia, A., Resnik, P., Simon, J. Z., & Lau, E. (2021).
27 *Parallel processing in speech perception: Local and global representations of linguistic*
28 *context*. <https://doi.org/10.1101/2021.07.03.450698>
- 29 Brodbeck, C., Das, P., Brooks, T., & Reddigari, S. (2019). *Eelbrain 0.31* (v0.31) [Computer
30 software]. Zenodo. <https://doi.org/10.5281/ZENODO.3564850>
- 31 Brodbeck, C., Hong, L. E., & Simon, J. Z. (2018). Rapid Transformation from Auditory to
32 Linguistic Representations of Continuous Speech. *Current Biology*, *28*(24), 3976–
33 3983.e5. <https://doi.org/10.1016/j.cub.2018.10.042>
- 34 Brodbeck, C., Jiao, A., Hong, L. E., & Simon, J. Z. (2020). Neural speech restoration at the
35 cocktail party: Auditory cortex recovers masked speech of both attended and ignored
36 speakers. *PLOS Biology*, *18*(10), e3000883. <https://doi.org/10.1371/journal.pbio.3000883>
- 37 Brysbaert, M., & New, B. (2009). Moving beyond Kučera and Francis: A critical evaluation of
38 current word frequency norms and the introduction of a new and improved word
39 frequency measure for American English. *Behavior Research Methods*, *41*(4), 977–990.
40 <https://doi.org/10.3758/BRM.41.4.977>
- 41 Connine, C. M., Mullennix, J., Shernoff, E., & Yelen, J. (1990). Word Familiarity and Frequency
42 in Visual and Auditory Word Recognition. *Journal of Experimental Psychology:*
43 *Learning, Memory, and Cognition*, *16*(6), 1084–1096. [https://doi.org/10.1037/0278-](https://doi.org/10.1037/0278-7393.16.6.1084)
44 [7393.16.6.1084](https://doi.org/10.1037/0278-7393.16.6.1084)

- 1 Dahan, D., & Magnuson, J. S. (2006). Spoken Word Recognition. In M. J. Traxler & M. A.
2 Gernsbacher (Eds.), *Handbook of Psycholinguistics* (2nd ed., pp. 249–283). Elsevier.
3 <https://doi.org/10.1016/B978-012369374-7/50009-2>
- 4 Dahan, D., Magnuson, J. S., & Tanenhaus, M. K. (2001). Time Course of Frequency Effects in
5 Spoken-Word Recognition: Evidence from Eye Movements. *Cognitive Psychology*,
6 42(4), 317–367. <https://doi.org/10.1006/cogp.2001.0750>
- 7 David, S. V., Mesgarani, N., & Shamma, S. A. (2007). Estimating sparse spectro-temporal
8 receptive fields with natural stimuli. *Network: Computation in Neural Systems*, 18(3),
9 191–212. <https://doi.org/10.1080/09548980701609235>
- 10 Di Liberto, G. M., Wong, D., Melnik, G. A., & de Cheveigné, A. (2019). Low-frequency cortical
11 responses to natural speech reflect probabilistic phonotactics. *NeuroImage*, 196, 237–
12 247. <https://doi.org/10.1016/j.neuroimage.2019.04.037>
- 13 Donhauser, P. W., & Baillet, S. (2020). Two Distinct Neural Timescales for Predictive Speech
14 Processing. *Neuron*, 105(2), 385–393.e9. <https://doi.org/10.1016/j.neuron.2019.10.019>
- 15 Ettinger, A., Linzen, T., & Marantz, A. (2014). The role of morphology in phoneme prediction:
16 Evidence from MEG. *Brain and Language*, 129, 14–23.
17 <https://doi.org/10.1016/j.bandl.2013.11.004>
- 18 Fischl, B. (2012). FreeSurfer. *NeuroImage*, 62(2), 774–781.
19 <https://doi.org/10.1016/j.neuroimage.2012.01.021>
- 20 Fishbach, A., Nelken, I., & Yeshurun, Y. (2001). Auditory Edge Detection: A Neural Model for
21 Physiological and Psychoacoustical Responses to Amplitude Transients. *Journal of*
22 *Neurophysiology*, 85(6), 2303–2323. <https://doi.org/10.1152/jn.2001.85.6.2303>
- 23 Gagnepain, P., Henson, R. N., & Davis, M. H. (2012). Temporal Predictive Codes for Spoken
24 Words in Auditory Cortex. *Current Biology*, 22(7), 615–621.
25 <https://doi.org/10.1016/j.cub.2012.02.015>
- 26 Gaskell, M. G., & Marslen-Wilson, W. D. (2002). Representation and competition in the
27 perception of spoken words. *Cognitive Psychology*, 45(2), 220–266.
28 [https://doi.org/10.1016/S0010-0285\(02\)00003-8](https://doi.org/10.1016/S0010-0285(02)00003-8)
- 29 Gaston, P., & Marantz, A. (2018). The time course of contextual cohort effects in auditory
30 processing of category-ambiguous words: MEG evidence for a single “clash” as noun or
31 verb. *Language, Cognition and Neuroscience*, 33(4), 402–423.
32 <https://doi.org/10.1080/23273798.2017.1395466>
- 33 Gillis, M., Vanthornhout, J., Simon, J. Z., Francart, T., & Brodbeck, C. (2021). *Neural markers*
34 *of speech comprehension: Measuring EEG tracking of linguistic speech representations,*
35 *controlling the speech acoustics* [Preprint]. Neuroscience.
36 <https://doi.org/10.1101/2021.03.24.436758>
- 37 Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., Goj, R.,
38 Jas, M., Brooks, T., Parkkonen, L., & Hämäläinen, M. (2013). MEG and EEG data
39 analysis with MNE-Python. *Frontiers in Neuroscience*, 7.
40 <https://doi.org/10.3389/fnins.2013.00267>
- 41 Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C.,
42 Parkkonen, L., & Hämäläinen, M. S. (2014). MNE software for processing MEG and
43 EEG data. *NeuroImage*, 86(Supplement C), 446–460.
44 <https://doi.org/10.1016/j.neuroimage.2013.10.027>
- 45 Grosjean, F. (1980). Spoken word recognition processes and the gating paradigm. *Perception &*
46 *Psychophysics*, 28(4), 267–283. <https://doi.org/10.3758/BF03204386>

- 1 Gwilliams, L., King, J.-R., Marantz, A., & Poeppel, D. (2020). *Neural dynamics of phoneme*
2 *sequencing in real speech jointly encode order and invariant content* [Preprint].
3 Neuroscience. <https://doi.org/10.1101/2020.04.04.025684>
- 4 Gwilliams, L., & Marantz, A. (2015). Non-linear processing of a linear speech stream: The
5 influence of morphological structure on the recognition of spoken Arabic words. *Brain*
6 *and Language*, *147*, 1–13. <https://doi.org/10.1016/j.bandl.2015.04.006>
- 7 Heeris, J. (2018). *Gammatone Filterbank Toolkit*
8 (0626328ef7c31d3b33214db2fdcd52e8601eb4c5) [Computer software].
9 <https://github.com/detly/gammatone>
- 10 Kemps, R. J. J. K., Wurm, L. H., Ernestus, M., Schreuder, R., & Baayen, H. (2005). Prosodic
11 cues for morphological complexity in Dutch and English. *Language and Cognitive*
12 *Processes*, *20*(1–2), 43–73. <https://doi.org/10.1080/01690960444000223>
- 13 Kocagoncu, E., Clarke, A., Devereux, B. J., & Tyler, L. K. (2017). Decoding the Cortical
14 Dynamics of Sound-Meaning Mapping. *The Journal of Neuroscience*, *37*(5), 1312–1319.
15 <https://doi.org/10.1523/JNEUROSCI.2858-16.2016>
- 16 Lewis, G., & Poeppel, D. (2014). The role of visual representations during the lexical access of
17 spoken words. *Brain and Language*, *134*, 1–10.
18 <https://doi.org/10.1016/j.bandl.2014.03.008>
- 19 Magnuson, J. S. (2016). Mapping spoken words to meaning. In M. G. Gaskell & J. Mirkovic
20 (Eds.), *Speech Perception and Spoken Word Recognition* (pp. 76–96). Routledge.
- 21 Magnuson, J. S., Mirman, D., & Myers, E. (2013). Spoken Word Recognition. In D. Reisberg
22 (Ed.), *Oxford Handbook of Cognitive Psychology* (pp. 412–441). Oxford University
23 Press.
- 24 Marslen-Wilson, W. D. (1987). Functional parallelism in spoken word-recognition. *Cognition*,
25 *25*(1), 71–102. [https://doi.org/10.1016/0010-0277\(87\)90005-9](https://doi.org/10.1016/0010-0277(87)90005-9)
- 26 Marslen-Wilson, W. D., & Tyler, L. K. (1980). The temporal structure of spoken language
27 understanding. *Cognition*, *8*, 1–71.
- 28 McAllister, J. M. (1988). The use of context in auditory word recognition. *Perception &*
29 *Psychophysics*, *44*(1), 94–97. <https://doi.org/10.3758/BF03207482>
- 30 McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017, August).
31 Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. *Proceedings of*
32 *the 18th Conference of the International Speech Communication Association*.
- 33 McClelland, J. L., & Elman, J. L. (1986). The TRACE model of speech perception. *Cognitive*
34 *Psychology*, *18*(1), 1–86. [https://doi.org/10.1016/0010-0285\(86\)90015-0](https://doi.org/10.1016/0010-0285(86)90015-0)
- 35 McQueen, J. M. (2007). Eight questions about spoken word recognition. In M. G. Gaskell (Ed.),
36 *The Oxford Handbook of Psycholinguistics* (pp. 36–54). Oxford University Press.
37 <https://doi.org/10.1093/oxfordhb/9780198568971.013.0003>
- 38 Norris, D. (1994). Shortlist: A connectionist model of continuous speech recognition. *Cognition*,
39 *52*(3), 189–234. [https://doi.org/10.1016/0010-0277\(94\)90043-4](https://doi.org/10.1016/0010-0277(94)90043-4)
- 40 Norris, D., & McQueen, J. M. (2008). Shortlist B: A Bayesian model of continuous speech
41 recognition. *Psychological Review*, *115*(2), 357–395. [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-295X.115.2.357)
42 [295X.115.2.357](https://doi.org/10.1037/0033-295X.115.2.357)
- 43 Smith, S., & Nichols, T. (2009). Threshold-free cluster enhancement: Addressing problems of
44 smoothing, threshold dependence and localisation in cluster inference. *NeuroImage*,
45 *44*(1), 83–98. <https://doi.org/10.1016/j.neuroimage.2008.03.061>

- 1 Taulu, S., & Simola, J. (2006). Spatiotemporal signal space separation method for rejecting
2 nearby interference in MEG measurements. *Physics in Medicine and Biology*, *51*(7),
3 1759–1768. <https://doi.org/10.1088/0031-9155/51/7/008>
- 4 Tucker, B. V., Brenner, D., Danielson, D. K., Kelley, M. C., Nenadić, F., & Sims, M. (2019).
5 The Massive Auditory Lexical Decision (MALD) database. *Behavior Research Methods*,
6 *51*(3), 1187–1204. <https://doi.org/10.3758/s13428-018-1056-1>
- 7 Weide, R. (1994). *CMU pronouncing dictionary*. <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>
- 8 Wurm, L. H., Ernestus, M. T. C., Schreuder, R., & Baayen, H. (2006). Dynamics of the auditory
9 comprehension of prefixed words: Cohort entropies and Conditional Root Uniqueness
10 Points. *The Mental Lexicon*, *1*(1), 125–146. <https://doi.org/10.1075/ml.1.1.08wur>
- 11 Yee, E., & Sedivy, J. C. (2006). Eye movements to pictures reveal transient semantic activation
12 during spoken word recognition. *Journal of Experimental Psychology: Learning,*
13 *Memory, and Cognition*, *32*(1), 1–14. <https://doi.org/10.1037/0278-7393.32.1.1>
- 14 Zwitserlood, P. (1989). The locus of the effects of sentential-semantic context in spoken-word
15 processing. *Cognition*, *32*(1), 25–64. [https://doi.org/10.1016/0010-0277\(89\)90013-9](https://doi.org/10.1016/0010-0277(89)90013-9)
16