

1
2
3 Disentangling the role of biphone probability from neighborhood density in the perception of
4 nonwords
5

6 Jeremy Steffman & Megha Sundara
7 UCLA Department of Linguistics
8
9

10
11 Address for correspondence

12 Jeremy Steffman
13 UCLA Department of Linguistics
14 3125 Campbell Hall
15 Los Angeles, CA 90095–1543
16 Email: jeremysteffman@gmail.com
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
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1

2

3 *Abstract*

4

5

6

7

8 In four experiments we explored how biphone probability and lexical neighborhood density influence
9
10 listeners' categorization of vowels embedded in nonword sequences. We found independent effects of
11
12 each. Listeners shifted categorization of a phonetic continuum to create a higher probability sequence,
13
14 even when neighborhood density was controlled. Similarly, listeners shifted categorization to create a
15
16 non-word from a denser neighborhood, even when biphone probability was controlled. Next, using a
17
18 visual world eye-tracking task, we determined that biphone probability information is used rapidly by
19
20 listeners in perception. In contrast, task complexity and irrelevant variability in the stimuli interfere with
21
22 neighborhood density effects. These results support a model in which biphone probability, but not
23
24 neighborhood density information, is encoded prelexically.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1. Introduction

Listeners rely on knowledge about the phonological and lexical organization of their language when they process speech. Two such influences are biphone probability - the probability of two sounds occurring in sequence, and lexical neighborhood density - the number and frequency of similar sounding words in the lexicon (Vitevich & Luce, 1999). Both have been shown to influence phonetic categorization. In a series of experiments, we evaluated the independent contribution and time course of biphone probability and neighborhood density effects in a phonetic categorization paradigm.

Lexical neighborhood density is defined as the number of known words that are similar to a string by a given metric. Commonly, a neighbor is defined in terms of phoneme overlap: a word's (or non-word's) neighbors are words which can be created from substituting, adding or deleting a single phoneme. As captured in the Neighborhood Activation Model (Luce and Pisoni 1998), the central idea is that in speech processing, multiple lexical candidates are activated based on their similarity to the input, with various consequences for the processing of both words and non-words. As a result of competition, the recognition of sequences from high density neighborhoods is slower compared to sequences from low density neighborhoods (e.g., Luce & Pisoni, 1999; Vitevitch, 2002a), although the production and recall of sequences from high density neighborhoods is privileged in contrast to those from lower density neighborhoods (e.g., Vitevitch, 2002b; Roodenrys & Hinton, 2002). Children as well recognize words from high density neighborhoods more slowly than those from low density neighborhoods (e.g., Garlock, Walley, & Metsala, 2001; Munson, Swenson & Manthei, 2005). Sensitivity to neighborhood density emerges gradually only during the second year of life. Thus, 14-month-olds are sensitive to the details of pronunciation of familiar words from high as well as low density neighborhoods (Swingley & Aslin, 2002), but by 17-months infants are more likely to learn novel words from low density neighborhoods compared to those from high density neighborhoods (Hollich, Jusczyk & Luce, 2002).

1
2
3 It is typically challenging to distinguish effects of neighborhood density from those of biphone
4
5 probability because these measures are highly correlated, at least in English (Pitt & McQueen 1998,
6
7 Vitevitch & Luce 1998; Vitevitch, Luce, Pisoni & Auer, 1999; Landauer & Streeter 1973): words in
8
9 denser lexical neighborhoods tend to be comprised of higher probability sequences. Nonetheless,
10
11 whether neighborhood density and biphone probability independently affect speech perception is central
12
13 to the distinction between theories that do and do not advocate for the role of feedback in models of
14
15 spoken word recognition.
16
17
18
19

20 In this paper we focused on isolating the role of biphone probability and neighborhood density
21
22 using a phonetic categorization task. To do so, we built on an experiment by Newman et al. (1997).
23
24 Newman et al. tested phonetic processing using a 2AFC task in which listeners categorized a VOT
25
26 continua, with two non-word endpoints. They found that listeners categorization of the VOT continua
27
28 was biased towards non-words from denser neighborhoods. Newman et al., argue that their results can
29
30 only be captured by models of speech processing which are interactive, and allow for feedback from
31
32 lexical to prelexical levels of representation. TRACE (McClelland & Elman 1986) represents one such
33
34 interactive model which implements feedback from activated lexical entries to a lower, phonemic layer
35
36 of representation. In TRACE, an item with an ambiguous stop consonant activates both non-word
37
38 endpoints, which in turn activate lexical neighbors. Top-down activation from these neighbors then
39
40 boosts activation for the denser-neighborhood non-word to a greater extent, biasing categorization in its
41
42 direction. Thus, Newman et al., argue that their results are supportive of such a model where activation
43
44 of neighbors modulates the activation of prelexical nodes via feedback.
45
46
47
48
49
50

51 Norris, McQueen & Cutler (2000) argue that Newman et al.'s results can be explained without
52
53 recourse to feedback. One possibility they suggest is that Newman et al.'s neighborhood density effects
54
55 may be attributed to differences in biphone probability alone. It has been previously shown that listeners
56
57 tend to categorize an ambiguous stop consonant as one that results in a higher probability sequence given
58
59
60
61
62
63
64
65

1
2
3 the preceding segment (Pitt & McQueen, 1998). Crucially, if such differences can be explained by
4
5 differences in biphone probability alone, this obviates the need for feedback from the lexicon.
6

7
8 However, Norris et al.'s hypothesis has been only partially supported. As Newman et al., argue,
9
10 their results cannot be explained by differences in the probabilities between the initial consonant and the
11
12 following vowel because they controlled for it. Similarly, Brancazio & Fowler (2000) argue that at least
13
14 for some continua tested by Newman et al., the neighborhood effects cannot be explained by differences
15
16 in the probabilities of the non-adjacent consonants; although Norris et al. provide some evidence that
17
18 Newman et al.'s results could be attributed to higher order (triphone) probabilities.
19
20

21
22 Alternately, Norris et al., (2000) argue that Newman et al.'s results could be lexical and influence
23
24 categorization, but at the later decision stage, and thus be modeled as a response bias. Because lexical
25
26 effects at the decision stage do not alter the early recognition of non-words, feedback is not necessary to
27
28 explain them. Consistent with this hypothesis, Newman et al. report neighborhood density effects only
29
30 at intermediate and long (cf. Fox, 1984) but not short reaction times. These late effects of neighborhood
31
32 density could well emerge from the influence of lexical variables at the decision stage.
33
34
35

36
37 Based on these findings, Pitt & McQueen (1998) argue for autonomous models of speech
38
39 processing where listeners' expectations about sound sequences, as indexed by biphone probabilities,
40
41 alone feed-forward activate phonemic units, as in Shortlist A (Norris 1994). Further, effects of lexical
42
43 neighborhoods do not provide interactive feedback from lexical to prelexical levels of processing, but
44
45 instead influence decisions due to feed-forward activation of decision nodes.
46
47
48

49
50 Biphone probability effects are also well established in the literature. Adults are more likely to
51
52 recognize, name (e.g., Frisch, Large & Pisoni, 2000; Vitevitch, Armbruster & Chu, 2004), recall (Thorn
53
54 & Frankish, 2005) and accept as word-like (Pierrehumbert, Needle & Hay, 2018), high probability
55
56 sequences, compared to sequences with a lower probability. This advantage for high probability
57
58 sequences is evident in children as well who produce nonwords with high probability sequences more
59
60 accurately (e.g., Munson, Edwards & Beckman, 2005; Gathercole, Frankish, Pickering & Peaker, 1999).
61
62
63

1
2
3 Finally, these effects seem to be evident even in infancy. Whether infants are learning English (Jusczyk,
4
5 Luce & Charles-Luce, 1994; Mattys, Jusczyk, Luce & Morgan, 1999), Dutch (Freiderici & Wessels,
6
7 1993) or Catalan (Sebastian-Gallés & Bosch, 2002), 9-months listen longer to high probability sequences
8
9 compared to those with a low probability. Additionally, English learning 9-month-olds can use dips in
10
11 biphone probability sequences to segment words (Mattys & Jusczyk, 2001); they can also segment nonce
12
13 words beginning with high biphone probability sequences but not those with low biphone probabilities
14
15 (Archer & Curtin, 2016). In sum, biphone probability effects on speech perception and production are
16
17 evident early in acquisition and through adulthood.
18
19
20
21

22
23 These two sets of findings thus offer contrasting views of the variables implicated in phonetic
24
25 processing. In the view advocated by Newman et al (1997), neighborhood activations play a central role
26
27 in phonetic processing. As exemplified in TRACE, Newman et al attribute these neighborhood effects
28
29 to feedback from the lexicon. Critically in TRACE, feedback alters the prelexical activation of phones;
30
31 but there is no independent representation of biphone information. Given the high correlation between
32
33 biphone probability and neighborhood density, phonotactic probability effects in phonetic processing in
34
35 such models are simply a by-product of neighborhood activations. That is, a higher density neighborhood
36
37 increases activation for high probability words and non-words. Further, because feedback introduces a
38
39 delay, neighborhood density effects are not immediate. However, this account fails to capture how
40
41 biphone sensitivity in young infants might correspond to neighborhood effects seen in the second year
42
43 of life.
44
45
46
47
48

49
50 Alternatively, there are models where it is biphone probabilities that play a central role in
51
52 phonetic processing. As exemplified in Shortlist A (Norris, 1994) and Merge (Norris, 1999), such models
53
54 include architecture compatible with a prelexical, autonomous representation of sequential/phonotactic
55
56 information with no independent role for neighborhood density. Because biphone probability effects are
57
58 represented prelexically, they are expected to influence phonetic processing with little to no delay.
59
60
61
62
63
64
65

1
2
3 Finally, Norris et al (2000) outline a third possibility where both biphone probability and
4
5 neighborhood density independently influence phonetic processing. In this proposal as well, biphone
6
7 probability influences are prelexical. In addition, neighbors are activated and feed-forward activation to
8
9 decision nodes. Thus, unlike biphone probability, neighborhood density does not affect the prelexical
10
11 activation of phones. Instead, it acts as a bias at the decision stage. In this account, neighborhood density
12
13 effects are delayed relative to biphone probability effects, though not as late as might be expected from
14
15 a feedback account.
16
17
18
19

20 As is clear from the preceding discussion, answers to two questions are critical in teasing apart
21
22 these accounts. First, are biphone probability and neighborhood density effects independent? Second,
23
24 what is the time course for biphone probability and neighborhood density effects? In four experiments,
25
26 we used phonetic categorization of non-words to disentangle the contribution of biphone probability and
27
28 neighborhood density. Following Newman et al. (1997) and Pitt and McQueen (1998) we used non-
29
30 words because this allowed us to test listeners' use of information which does not directly depend on
31
32 word-hood, word frequency, semantic associations with words, and so on. First, we tested whether
33
34 biphone probability and neighborhood density independently influence phonetic categorization when the
35
36 other variable is controlled. Next, we used eye-tracking to determine the time course of each of these
37
38 effects. Together, these results address how and when listeners use lexical and phonological information
39
40 in speech processing, and thus inform models of spoken word recognition.
41
42
43
44
45
46
47
48

49 2. *Experiment 1: Biphone probability effects controlling for neighborhood density*

51 The goal of Experiment 1 was to test whether differences in biphone probability influenced
52
53 listeners' categorization of a continuum, when neighborhood density was controlled. For this, we created
54
55 a vowel continuum from English vowel /ε/ to /æ/ by manipulating F1, F2 and F3. The continuum was
56
57 presented in one of two CVC frames and listeners were asked to categorize the vowel as /ε/ to /æ/. The
58
59 two frames were: /mεb~/~ /mæb/ and /mεv~/~ /mæv/. The consonant frames were selected such that neither
60
61
62
63
64
65

endpoint was a word in English, and both coda consonants /b/ and /v/ involved labial constrictions so formant trajectories at the offset of the vowel could be expected to be similar, allowing for identical continuum steps to be used with each frame. Critically, continuum endpoints were selected to control for neighborhood density biases, while varying BP biases.

TABLE 1 HERE

Table 1 shows the biphone-probabilities and frequency-weighted neighborhood densities for the endpoints of all the continua used in this study. This includes BP information for both C₁V₂ and V₂C₃ in the CVC sequence. It also includes ND for C₁V₂ as well as for the CVC sequence as a whole. See section 2.1.1 for details on how these measures were calculated.

First consider neighborhood density for the full CVC sequence of the two continua used in Experiment 1: the non-word /mɛb/ has a frequency-weighted neighborhood density of 17.96. The other endpoint of the continuum, /mæb/ has a frequency-weighted neighborhood density of 29.54. Following Newman et al., listeners should be biased towards a denser-neighborhood non-word when exposed to an ambiguous stimulus. In the present case, a denser neighborhood for /mæb/ would bias listeners to respond /æ/ when exposed to ambiguous items on a /mɛb~/mæb/ continuum. We can index the magnitude of this bias by subtracting the frequency-weighted neighborhood density of /mɛb/ from that of /mæb/. The /mɛb~/mæb/ continuum therefore has a neighborhood density bias that is positive, i.e., towards /æ/. The bias is 11.94 (29.54-17.96). The /mɛv~/mæv/ continuum as well has a neighborhood density bias for /æ/ of 12.88. Comparing the biases for the two continua, we see that although both have an /æ/ bias, the /mɛv~/mæv/ continuum has a slightly larger one. This would predict that if listeners are sensitive to neighborhood density alone, they should show increased /æ/ responses to the /mɛv~/mæv/ continuum compared to /mɛb~/mæb/ continuum. However it should be noted that the difference in bias

across continua here is much smaller than reported for the continua used by Newman et al. (1997), suggesting its influence may be minimal.

Continuum biases were calculated in the same way for biphone probability, for both C₁V₂ and V₂C₃ in the CVC string, though note that the onset consonant is the same across continua so it cannot contribute to any sequential probability differences across conditions. As shown in Table 1, the V₂C₃ portion of the /mɛb/~mæb/ continuum exhibited an /æ/ bias (0.0019), while the V₂C₃ portion of the /mɛv/~mæv/ continuum exhibited a /ɛ/ bias (-0.0007). This differential predicts that a coda /b/ should bias listeners towards /æ/ responses, such that they prefer a relatively higher probability sequence /mæb/ (as compared to /mɛb/), and vice versa for coda /v/.

If listeners are sensitive to biphone probability information, they should thus show *increased* /æ/ responses for the /mɛb/~mæb/ continuum compared to the /mɛv/~mæv/ continuum, where /v/ should bias towards /ɛ/. Note that the bias based on biphone probability is in the opposite direction on the bias predicted by neighborhood density, making this a fairly conservative test for biphone probability effects (though density biases are minimally different). A finding in the predicted direction would therefore provide strong evidence that biphone probability exerts an independent influence on listeners' categorization of speech sounds.

2.1 Materials

Stimuli for Experiment 1 were created by resynthesizing the speech of an adult male speaker of American English. The stimuli were first recorded at 44.1 kHz (32 bit) in a sound-attenuated room, using an SM10A ShureTM microphone and headset. The starting point for the creation of stimuli was the speaker's natural production of two CVC nonwords: /mɛv/ and /mæv/. The vocalic portion of both of these nonwords was excised from the CVC frame. The continuum was then synthesized in Praat via LPC decomposition and resynthesis of F1, F2 and F3 (Winn, 2016). Resynthesis used /ɛ/ as a base and interpolated F1, F2, and F3 in evenly Bark-spaced steps to their respective values for the /æ/ token in 12

steps. The higher frequency energy and pitch contour were preserved during resynthesis such that they matched that of the original /ε/ token. The resulting 12 step continuum therefore varied only in the frequencies of the first three formants. The onset /m/ from the original production of /mεv/ was then re-spliced onto each continuum. The coda /b/ and /v/ were cross-spliced from productions of /mεb/ and /mæv/ respectively. This was done to remove any possible acoustic traces of co-articulatory information from the preceding vowel cuing these consonants; though note it is unlikely that the cross-spliced stop closure/release and fricative noise contained cues to identify the original preceding vowel. Specifically, given that we predicted a following /v/ should bias listeners towards /ε/ categorization, as outlined above, the cross-spliced /v/ came from a post -/æ/ context, ensuring any possible acoustic information from the preceding vowel would predict the opposite adjustment in categorization. For the same reason /b/ was cross-spliced from a post-/ε/ context. These manipulations created 24 unique stimuli (12 continuum steps × 2 consonant frames).

2.1.1 Calculation of biphone frequency and neighborhood density

All density and biphone probability measurements were made using the KU Phonotactic Probability Calculator and KU Neighborhood Density Calculator (Vitevich & Luce 2004), which provides frequency-weighted positional estimates for individual phones in a sequence, as well as biphone co-occurrence probabilities. The lexicon used in the calculators is based on the Merriam Webster Pocket Dictionary, with frequency measures from Kučera & Francis (1967). Neighborhood density was calculated using the same formula as in Newman et al. (1997), where each neighbor's contribution was frequency weighted. Each neighbor's frequency contribution was calculated by taking the logarithm (base 10) of the raw frequency times 10. This value was then summed for all neighbors for a given word, to provide a frequency-weighted neighborhood density. To ensure that the words entered into the calculation were likely known by our participants, we used only words that have previously been rated as familiar (Nusbaum, Pisoni & Davis 1984), using a familiarity index of 5.0 or higher as a cut off (on

a 7 point scale, see Nusbaum et al. 1984). We also made the same calculations including all (even less familiar) words, this did not change the direction of any predicted effects.

2.2 *Participants*

Thirty-five self-identified native speakers of American English with normal hearing participated in Experiment 1. Participants were students at North American University and received course credit for participation.

2.3 *Procedure*

Participants completed the task seated in front of a desktop computer, in a sound-attenuated booth in the lab. Stimuli were presented binaurally via a 3M™ Peltor™ listen-only headset. Participants were told that they would hear a speaker of English say nonce words, and that their task was simply to select which word they heard.

During a trial, participants were presented visually with two letters placed on either side of the computer screen: ‘E’ and ‘A’. Prior to the trials beginning, participants were instructed that they should select ‘E’ if they heard the sound /ε/, and ‘A’ if they heard the sound /æ/. This was conveyed by giving examples of real words that rhymed with of the non-word continuum endpoints in the written computer instructions. Participants indicated their response by keypress, where an ‘f’ key-press indicated the letter on the left side of the screen and a ‘j’ keypress indicated a letter on the right side of the screen. The side of the screen on which each letter appeared was counterbalanced across participants. Participants completed 8 practice trials in which they heard each continuum end point in each CVC frame two times. During test trials participants heard each unique stimulus 8 times for a total of 192 trials. Stimuli were completely randomized. Testing took about 15 minutes.

2.4 Results & Discussion

Results were analyzed using Bayesian mixed-effects logistic regression, with the *brms* package (Bürkner, 2018) in R. We predicted the log odds of selecting an /æ/ response as a function of the step of the continuum, the consonantal frame (manipulating BP), and the interaction of these two fixed effects. Continuum step was treated as a continuous variable and centered. Consonantal frame was contrast coded, with /mVb/ mapped to 0.5 and /mVv/ mapped to -0.5. Random effects in the model included by-participant intercepts with maximally specified random slopes including both fixed effects and their interaction. The default uniform prior distribution was employed in the model. In describing the results, we report the model estimates and 95% credible intervals for them. An effect is taken to have a reliable impact on responses when the 95% credible interval for an estimate excludes zero.

In the model, the main effect of step was credible as expected ($\beta = -2.47$, $CI = [-3.07, -1.85]$), confirming that listeners /æ/ responses decreased at the /ε/ end of both continua. The main effect of consonantal frame was also credible ($\beta = 0.62$, $CI = [0.05, 1.21]$), but its interaction with step was not ($\beta = -0.12$, $CI = [-0.34, 0.11]$). The effect of frame indicates that, consistent with biphone probability effects, participants showed an overall bias to categorize the target as /æ/ in the /mVb/ frame compared to the /mVv/ frame. This can be clearly seen in Figure 1.

FIG. 1 HERE

Notably, this effect involved a vertical shift in the categorization function not restricted to ambiguous regions of the continuum. Vertical shifts in categorization function have typically been attributed to decision biases (e.g. Massaro & Cowan, 1993, Norris et al., 2000). Contextual effects that directly modify input representations are predicted, instead, to only influence categorization at ambiguous steps on a continuum (e.g. Massaro, 1989, Massaro & Cowan, 1993). A main effect of consonantal frame, in the absence of an interaction with step is then consistent with the models where

biphone probability information feeds into a higher level representation, such as the decision nodes in MERGE (Norris et al. 2000), where explicit decisions about the stimuli are made. Whether listeners use biphone probability information online and early in processing though, remains unclear. We tested this in Experiment 3.

Irrespective of the nature of its influence, the results of Experiment 1 show that biphone probability can indeed modulate listeners' categorization of phonetic continua, in agreement with Pitt & McQueen (1998). Crucially, these results cannot be attributed to differences in neighborhood densities because we carefully controlled them during the stimulus selection.

3 Experiment 2: Neighborhood density effects controlling for biphone probability

Experiment 1 showed a clear effect of biphone probability in phonetic categorization of a non-word continuum, which was independent of neighborhood density. In Experiment 2 we tested if we could obtain evidence for an independent effect of neighborhood density. Two new continua were created: /bɛp/ ~ /bæp/ and /bɛb/ ~ /bæb/. V₂C₃ biphone probability was matched completely for these two pairs (see Table 1), such that they both exhibited an equally strong /æ/ bias (0.0019). Unlike Experiment 1 however, the neighborhood density bias for these continua was substantially different, both exhibited an /æ/ bias, though the bias for the /bɛp/ ~ /bæp/ continuum (29.35) was stronger than that for the /bɛb/ ~ /bæb/ continuum (19.85). As outlined above, a denser neighborhood should bias listeners towards /æ/ responses, predicting that more /æ/ responses should be observed for the /bɛp/ ~ /bæp/ continuum. Such a finding could not be explained by differences in biphone probability, which are identical (0.0046 summed CVC BP in both continua, see Table 1). The empirical prediction is thus that a coda /p/ frame should show *increased* /æ/ responses, due to the stronger /æ/ ND bias in this frame.

3.1 Materials

Experiment 2 used the same vowel continuum created in Experiment 1, however, we presented them in different frames. The new frame consonants were cross-spliced from the same speakers' productions. The initial /b/ was cross-spliced from a production of /beb/. The coda /b/ was cross-spliced from a production of /bæb/, and the coda /p/ was cross-spliced from a production of /bɛp/. As with Experiment 1, this method of cross splicing was chosen to remove any possible acoustic traces of the preceding vowel on cross-spliced coda consonants. Specifically, because we predicted that the /bæp/~bɛp/ continuum should bias categorization towards /æ/ (as compared to /bæb/~beb/), the coda /p/ was cross-spliced from an original /ɛp/ sequence. Likewise, the coda /b/ was cross-spliced from an original /æb/ sequence. Because the consonants used in Experiment 1 also involved labial constrictions, formant transitions at the onset and offset of the vowel continuum were judged to sound natural in these new frames.

3.2 Participants and procedure

Thirty-five self-identified native speakers of American English participated in Experiment 2. Participants were students at a North American University and received course credit for participation. The procedure for Experiment 2 was identical to that in Experiment 1.

3.4 Results & Discussion

The model specifications and model fitting procedure were identical to that in Experiment 1. Results are plotted in Figure 2. In contrast coding consonant frame, /bVp/ was mapped to -0.5 and /bVb/ was mapped to 0.5. As in Experiment 1, an expected main effect of step was found ($\beta = -3.84$, $CI = [-4.80, -2.91]$). The main effect of consonant frame was also credible ($\beta = -3.84$, $CI = [-0.67, -0.10]$). As in Experiment 1, the interaction of consonant frame and continuum step was not credible ($\beta = -0.40$, $CI = [-0.94, -0.01]$). In Figure 2 we can see the effect of consonant frame: consistent with predicted

neighborhood density effects, listeners showed *increased* /æ/ responses with the /bVp/ frame, shifting categorization in accordance with neighborhood density. This result replicates Newman et al.'s findings, and crucially precludes a biphone probability differences as a possible explanation. Like the BP effect, this adjustment in categorization was a vertical shift, consistent with a response bias, as indicated by the lack of a significant interaction between frame and continuum step.

FIG. 2 HERE

4 Experiment 3: Time course of biphone probability and neighborhood density effects

Taking Experiment 1 and 2 together, we have evidence for an independent influence of both biphone probability and neighborhood density, as indexed by listeners' categorization responses. However, both of these effects are consistent with a response bias given the vertical shifts in the categorization function observed in Experiments 1 and 2. Additionally, categorization performance only provides a measure of the endpoint of the speech recognition process. To obtain precise timing information about when BP and ND effect recognition, we need evidence from online tasks. Some previous research, outlined below, offers some relevant timecourse comparisons.

Using brain imaging, Pylkkänen, Stringfellow, & Marantz (2002) provide some evidence that biphone probability effects are consistently observed between 300 and 400ms post stimulus onset. In an MEG experiment, they administered a lexical decision task to listeners who were presented with CVC sequences that were either high probability and high density or low probability and low density. They investigated an MEG response component - M350 - which peaks between 300 and 400 ms post stimulus onset. Because the M350 was facilitated in response to the manipulated probability, and not inhibited as expected for a density manipulation, Pylkkänen et al argue that the M350 is sensitive to biphone probability. They did not find a clear correlate of the density effect in later MEG components. Thus, the

MEG results present an estimate of the timeline for probability effects, and indirect support that this may be different from the effect of neighborhood density (see also Pylkkänen & Marantz, 2003),

More recently, Kingston and colleagues (Kingston, Levy, Rysling, & Staub, 2016) report on two experiments where they evaluated the time course of lexical effects on phonetic processing. In Kingston et al.'s experiments, listeners were asked to categorize a word to nonword phonetic continuum. They reasoned that if lexical effects are driven by feedback, they should be delayed as demonstrated in TRACE simulations (McClelland & Elman 1986). However, a rapid use of lexical information in categorization would constitute as evidence against feedback, and be more consistent with a feed-forward account. Based on results from two eye-tracking experiments Kingston et al., claim that lexical effects influence phonetic processing between 300 and 400 ms after stimulus onset; and thus are too early to be consistent with feedback.

A closer look at Kingston et al.'s experiments, however, offers an alternative explanation for their findings. First, in Kingston et al.'s Experiment 2a – the lexical effect is confounded with a biphone probability effect. In this experiment, listeners were presented with a continuum ranging between the vowels /ε/ and /Λ/ in a CVC(C) frame; whether the end point was a word or non-word was determined by the final consonant. The continuum was placed in one of four frames: (1) /b _ ŋk/ forming the word “bunk” with /Λ/, (2) /d _ ŋk/ forming the word “dunk” with /Λ/, (3) /b _ f/ and (4) /d _ f/ (both resulting in nonwords). The initial consonant was varied to manipulate spectral context, and will not be discussed here; its inclusion does not alter the conclusions based on biphone probability differences discussed below. Because a coda /ŋk/ creates words with the vowel /Λ/, but not /ε/, the Kingston et al. predicted that /ŋk/ should increase looks to an orthographic representation of /Λ/ (“U”), as compared to a following /f/. This is what the authors found, with the influence of the coda consonant(s) emerging within 300-400 ms of stimulus onset. A different interpretation of these finding emerges if we compare the biphone probabilities for the vowel and following consonant sequence. In the /f/ context, the biphone probabilities are essentially matched with a very slight /Λ/ bias: 0.0002 for /Cεf/ and 0.0004 /CΛf/. However, the

biphone probability for the vowel and following consonant /ŋ/, reveals an asymmetry: a following /ŋ/ engenders a strong /ʌ/ bias: 0.0003 for /Cɛŋ/ and 0.0027 for /Cʌŋ/. The magnitude of this /ʌ/ bias is comparable to our own biphone probability manipulation in Experiment 1. Thus, an alternate explanation for Kingston et al.’s results is that the time course from Experiment 2a reflects a difference in biphone probability between the sequences, and therefore, like in Pykkänen et al.’s MEG experiment, occurs between 300 and 400 ms post stimulus onset.

In the other eye tracking experiment reported by Kingston et al. (Experiment 1a), listeners categorized a continuum of fricative noise that ranged from /s/ to /f/. The continuum was followed by one of three frames: (1) /_ aɪ /, creating a word with /f/ “file”, but not with /s/, (2) /_ aɪd /, creating a word with /s/ “side”, but not with /f/, and (3) control frame /_ aɪm / for which both continuum endpoints were non-words. The online effect was significant only in the /_ aɪ / frame, with increased looks to a visual ‘F’ target on the screen, in comparison to the control frame. This effect cannot be explained by biphone probability differences; the summed biphone probability of “file” (0.0043) is slightly lower than that of “sile” (0.0058). However, there was no significant difference in looks between the /_ aɪm / frame and the control frame /_ aɪd /, where we would expect to see more looks to a visual ‘S’ target when the lexical context “side” reinforces /s/. This asymmetry in online processing between the two experimental frames, makes it difficult to interpret the results from Kingston et al.’s Experiment 1a.

In Experiment 3 we used Kingston et al.’s experimental design with the carefully controlled stimuli used in Experiments 1 and 2, where the effects of biphone probability and neighborhood density were orthogonally manipulated. Specifically, we were interested in how these effects unfold online using a visual world eye-tracking task. Combining categorization with eye-tracking data allowed us to investigate the pre-decision stage integration of information as speech unfolds (unlike reaction times), as discussed in e.g. Norris et al. (2000). The eye movement response to the vowel spectra served as our baseline because it indexes a response to the signal. Given the independence of biphone probability and neighborhood density effects documented in Experiments 1 and 2 respectively, we expected to see an

1
2
3 independent influence of each variable in the online task as well. If biphone probability effects are
4
5 prelexical, we expected them to emerge soon after the spectral response (once listeners have heard the
6
7 coda consonant). In contrast, if biphone probability (and neighborhood density) effects are because of a
8
9 response bias, as indicated by a vertical shift in categorization, we expected minimal, or very late, effects
10
11 on looks. Of crucial interest was the relative timing of each effect. If neighborhood density effects
12
13 originate from a feedback loop between the lexicon and prelexical information, because feedback takes
14
15 time as shown in TRACE simulations (McClelland and Ellman 1986), the influence of ND should be
16
17 delayed in comparison to a spectral response. Recall that Newman et al. (1997) reported reliable ND
18
19 effects only at slow and intermediate reaction times, suggesting a later influence in processing.
20
21
22
23
24
25
26

27 *4.1 Materials*

28
29 The materials used in Experiment 3 were a subset of those used in Experiments 1 and 2. In order
30
31 to present listeners with relatively ambiguous stimulus tokens (following e.g. Mitterer & Reinisch 2013,
32
33 Reinisch & Sjerps 2013), we presented listeners the most ambiguous region of each continuum. This
34
35 was identified as the 4 step window centered around the 50% crossover points in the interpolated
36
37 categorization functions derived from Experiment 1 and 2. In both experiments, categorization was most
38
39 variable on steps 4 through 7. Participants heard all four continua (/mVb/, /mVv/, /bVb/, /bVp/) at these
40
41 four steps. There were thus 16 unique stimuli used in Experiment 3 (4 continuum steps \times 4 consonant
42
43 frames).
44
45
46
47
48
49
50
51

52 *4.2 Participants*

53
54 Seventy-two self-identified native speakers of American English with normal or corrected to
55
56 normal visions participated in Experiment 3. For five participants, a technical error caused the
57
58 experiment to terminate early. Out of the 128 trials in the experiment (see below), these five participants
59
60 completed 72, 120, 126, 126 and 127 trials respectively. These participants' data was retained and
61
62
63
64
65

analyzed. Participants were students at a North American University and received course credit for participation.

4.3 Procedure

In Experiment 3, we used a visual world eye-tracking task, with a similar design to that used by Kingston et al. (2016). Participants were seated in front of an arm-mounted SR Eyelink 1000 (SR Research, Mississauga, Canada), which was set to track the left eye remotely, at a sampling rate of 500 Hz, and at a distance of approximately 550 mm. The visual display was presented to participants on a 1920 × 1080 ASUS HDMI monitor. Participants were tested in a sound-attenuated room in the lab. Participants' gaze was calibrated using a 5-point calibration procedure at the start of each experiment.

During an experimental trial, participants were presented with orthographic E and A on the target screen (Kingston et al. 2016) and were instructed to click on the letter corresponding to sound they heard. As in Experiments 1 and 2, examples of real English words that rhymed with the nonwords were given to convey the intended letter-to-sound mapping. Participants' eye movements were monitored while they performed the task. The orthographic targets were arranged vertically in the visual display, with each letter centered horizontally, and positioned 270 pixels above and below the midpoint of the display. Each letter was presented in 60pt black Arial font. The location of each letter was counterbalanced across participants. Each trial began with the appearance of a black fixation cross in the center of the visual display (60 px by 60 px). Following Kingston et al. (2016), stimulus onset was look-contingent, such that the audio stimulus played only after a look was registered on the fixation cross. Eye-movements were recorded from the first appearance of the fixation cross until a click response was registered by participants. After a click response was provided, the location of the mouse cursor was re-centered on the screen. Each trial was separated by a 1 second interval.

During the experiment, participants heard 8 repetitions of the 16 unique stimuli in a random order, for a total of 128 trials. Participants additionally completed 8 training trials prior to test trials in

which they heard step 4 and step 7 for each frame, to give them practice with the experimental paradigm. The experiment took approximately 20 minutes to complete in total.

4.4. Analysis

We analyzed two measures in Experiment 3. First, we analyzed listeners' click responses; we used Bayesian mixed-effects to model the log odds of selecting an /æ/ response as a function of frames (/mVb/, /mVv/, /bVb/, /bVp/), like in Experiments 1 and 2. . The model was fit with the same fixed effects and random effect structure as previous models. We did this to replicate the offline effects from Experiments 1 and 2.

Second, we analyzed eye-movement data to understand the time course of BP and ND effects. Statistical analysis of the eye movement data was carried out using a Generalized Additive Mixed Model (GAMM), which offers a powerful tool for analyzing time-series and visual world data (Nixon, van Rij, Mok, Baayen & Chen, 2016; Zahner, Kutscheid & Braun, 2019). The model was implemented with the *mgcv* and *itsadug* packages in R (van Rij, Wieling, Baayen, & van Rijn, 2020; Wood 2011). As the numerical model output is fairly uninformative for understanding the timing questions asked here (Wood 2011; Zahner et al., 2019), the summary is included in the appendix. An AR1 error model predicted empirical-logit transformed looks to a target (Barr, 2008) binned in 20ms intervals. The model included parametric terms for continuum step and frame, smooths fit to model a non-linear interaction of continuum step by condition over time, and main effects for these interacting terms (using the *te()* function, cf. Nixon et al. 2016). By-participant random smooths over time (factor smooths) were additionally included with the *m* parameter set to 1 (Baayen, van Rij, de Cat, & Wood 2018). The analysis window was set to be 0 to 1200 ms from the onset of the target vowel; this was typically when listeners made a categorization decision, a point after which there was a substantial drop in recorded eye-movements. To assess the precise timing of effects of interest, the difference between two given smooths was visualized (Zahner et al. 2019). In this assessment, confidence intervals generated from the model

can inform us when the difference between two given smooths becomes reliably different from zero, i.e. the point in time at which smooths diverge.

4.5 Results & Discussion

4.5.1 Click responses

FIG. 3 HERE

Overall, in this experiment, the steps 4-7 were perceived as more /æ/-like; this is evident from the significant intercept ($\beta = 0.88$, CI = [0.68, 1.10]) and in Figure 3. This /æ/ bias was stronger than in Experiments 1 and 2. It is possible that despite sampling from around the 50% crossover point, listeners recalibrated categorization because of the absence of steps from the endpoint of the continua. Nevertheless continuum step still showed a credible effect as expected. Listeners decreased /æ/-responses ($\beta = -1.29$, CI = [-1.52, -1.07]) progressively from Step 4 towards Step 7 where formants were more /ε/-like. Thus, although listeners exhibited an overall /æ/ bias, listeners still used formant cues in the expected way.

The first comparison of interest was that of the frames manipulating biphone probability: /mVb/ versus /mVv/. We examined the biphone probability effect by extracting pairwise contrasts from the model with the *emmeans* package to compare the estimated effect and CI for given frames. As shown in Figure 3, there was a clear difference in listeners' categorization across these two frames, replicating the BP effect observed in Experiment 1. This difference was credible, and in the expected direction with the /mVb/ frame showing increased /æ/ responses relative to the /mVv/ frame ($\beta = 0.42$, CI = [0.03, 0.83]).

A comparison between the two frames manipulating neighborhood density: /bVb/ versus /bVp/, showed a different pattern. There was no credible difference between these two frames ($\beta = 0.03$, CI =

1
2
3 [-0.26, 0.29]). As we can see from Figure 3, these frames did not induce any reliable shift in
4
5 categorization, that is, we did not replicate the neighborhood density effect seen in Experiment 2.
6

7
8 However, an unexpected finding emerged: there were markedly fewer /æ/ responses to all /m/-
9
10 initial frames compared to /b/-initial frames. Differences in categorization between both /b/-initial
11
12 frames, as compared to both /m/ initial frames were credible in each pairwise comparison. Additionally,
13
14 as can be seen in the model output, with the reference level set to /mVb/, both /bVb/ ($\beta = 1.49$, CI =
15
16 [0.98, 2.04]) and /bVp/ ($\beta = 1.43$, CI = [0.86, 1.99]) showed credibly greater /æ/ responses. In fact, this
17
18 difference in /æ/ responses between the /m/- vs /b/-initial frames was even larger than the biphone
19
20 probability effect.
21
22
23

24
25 This effect emerged in Experiment 3 because we used a within-subject design in contrast to the
26
27 between-subjects design in Experiments 1 and 2 where the effects of the /m/-initial and /b/-initial frames
28
29 were investigated respectively. We can rule out that this effect was driven by the differences in biphone
30
31 probabilities of the /m/-initial and /b/-initial frames. From Table 1 we see that the biphone sequence
32
33 /mV/ has a stronger /æ/ bias (0.0042) compared to /bV/ (0.0027); this difference in biphone probability
34
35 would predict the opposite of the effect observed here. Even considering the summed biphone probability
36
37 of the whole CVC sequence, we see the following gradation in the strength of /æ/ biases, from largest to
38
39 smallest: /m_b/ (0.0061) > /b_p/ and /b_b/ (0.0046) > /m_v/ (0.0035). This too cannot explain the
40
41 difference we see between /m/- and /b/-initial frames, because on this basis alone we'd predict the most
42
43 /æ/ responses for /m_b/, which is clearly not the case.
44
45
46
47

48
49 The direction of difference in /æ/ responses between the frames is more consistent with a
50
51 difference in neighborhood density between the /m/- and /b/-initial frames with the latter having a
52
53 stronger /æ/ bias (Table 1). However, there is reason to be skeptical that neighborhood density
54
55 differences are driving the frame effect. The difference in neighborhood density between the two /b/-
56
57 initial frames was at least as large, if not larger in magnitude than the neighborhood density difference
58
59
60
61
62
63
64
65

between the /m/- and the /b/-initial frames. Yet, the frame effect was credible, but the neighborhood density difference between the two /b/-initial frames was not.

Instead, we speculate that by introducing different initial consonants in our frames, we may have introduced a new variable that influenced listeners' perception of the target vowel. A change in initial-consonant from /b/ to /m/ is a switch between an oral and a nasal onset. Although our vowel didn't vary in terms of nasality across frames (being originally produced in /m/ initial frames), listeners' perception of F1 and/or F2 is likely to have been modulated because they were compensating for the typical coarticulatory effects of nasals on vowel formants.

Nasalization of vowels adjacent to nasal consonants is well-attested in American English (e.g. Chen, Slifka & Stevens 2007; Cohn 1990). Nasalization typically lowers perceived F1 for low vowels (Diehl, Kleunder & Walsh, 1990), directly impacting listeners' perception of vowel height adjacent to nasal consonants (Beddor, 1993; Ohala, Beddor, Krakow & Goldstein 1986; Wright 1980). Ohala et al. (1986) present a test case that offers a close comparison to the present stimuli. They found that when a vowel on an /ε/ ~ /æ/ continuum was adjacent to a nasal consonant, but had only very weak nasalization (comparable to the present stimuli where vowels were originally produced in /m/-initial carryover contexts) listeners "overcompensated" for the expected effect of vowel nasalization. An adjacent nasal consonant accordingly lead to decreased /æ/ responses, i.e. perception of a higher vowel, /ε/. Thus, it is quite likely that the frame effect is due to the listener's compensation for the nasal context, and not attributable to either biphone probability or neighborhood density differences. We addressed this directly in Experiment 4.

4.4.2 Eye movement data

Given the strong /æ/ bias observed in categorization responses, we opted to model and plot *looks to* /ε/, effectively inverting the y axis from the categorization plot. Below we report effects of interest in turn.

4.4.2.1 Baseline continuum effects

FIG. 4 HERE

First, we explored the effect of the vowel spectrum – i.e., the changing formants along the continuum. In Figure 4, panel A we show raw eye movement data, with listeners' look to /ε/ over time split by continuum step. Here we can see that listeners' looks to the /ε/ target are indeed being impacted by the continuum, that is, as expected, more /ε/-like formants cause listeners to look more towards the /ε/ target.

In Figure 4 panel B, we show difference smooths from the model. Following Maslowski, Meyer & Bosker (2020), each trajectory shows the estimated *difference between steps 4 and 7 over time*, for each consonant frame. The divergence between the acoustically most distinct two steps allows an estimate of the time when listeners are using acoustic information to distinguish the most salient spectral differences in the vowel portion of the stimuli. The point in time at which a given difference smooth diverges from zero reliably (i.e. when the 95% CI exclude zero) in Figure 4, panel B can be taken as the time point at which listeners' looks are impacted by the continuum (i.e. different formant values in the target). The divergence time for each of the four difference smooths corresponding to each of the four frames is indexed by a dashed vertical line, matching that smooth in color. The average of these four divergences, which can be taken as a more general estimate for the effect of the continuum overall, is indexed by the blue dashed vertical line. We see in Figure 4B that listeners looking behavior was influenced by the continuum step about 373 ms following the onset of the target vowel. Given that it takes approximately 200 ms to initiate a saccade (Dahan, Magnuson, Tanenhaus & Hogan, 2001; Matin, Shao & Boff 1993), we can see that listeners waited to hear approximately 170 ms of the target vowel before using spectral information (on average across consonant frames). This effect is slightly more delayed than previous findings for the use of intrinsic spectral cues (cf. Reinisch & Sjerps, 2013; Bosker,

Reinisch & Sjerps 2017) likely because of the longer duration of the vowel used here (260 ms) compared to the ones used in these previous studies (120 -140 ms in the case of Reinisch & Sjerps, 2013).

Additionally, we can see that the timing of the effect of continuum varied by consonant frame. One clear difference here is manifested in the unexpected initial consonant effect, discussed above. Listeners were able to use the formant information earlier in both /m/-initial frames compared to the /b/-initial frames. This timing is consistent with the biasing effect of nasal consonants on the perception of vowel height. This early timing makes it unlikely that the frame effect is because of correlated neighborhood density effects; recall that at least in categorization tasks (Newman et al, 1997), neighborhood density effects are typically delayed.

4.4.2.1 Biphone probability and neighborhood density effects

FIG. 5 HERE

Figure 5 shows listeners' looks to the /ε/ target over time, split by consonant frame. As expected, we can see a separation based on initial consonant, lining up with the impact of initial consonant on the use of the continuum discussed above. The biphone probability effect of interest, comparing /m_b/ and /m_v/ frames, is also evident; listeners looked to /ε/ more in the /m_v/ frame compared to the /m_b/. In contrast, there was no robust separation between the two frames that differed in neighborhood density: /b_b/ and /b_p/.

FIG. 6 HERE

To assess these effects statistically and to establish the time course of each, we again inspected difference smooths from the model, this time comparing smooths for consonant frames of interest. First, to inspect

the influence of the biphone probability manipulation we plotted the divergence between /m_b/ and /m_v/ frames at each continuum step, as shown in panel A of Figure 6. Again, dotted vertical lines show divergences for continuum steps, indexed by color. The blue dashed vertical line shows the average of these four by-step divergence times. The averaged divergence estimate for the biphone probability manipulation obtained by the model was 508 ms from the onset of the target vowel, shown by the dashed blue vertical line. The onset of both coda consonants (i.e. the offset of the vowel) was 260 ms. Considering a 200 ms delay programing a saccade, the biphone probability effect emerged about 300 ms after vowel onset, that is within 40-50 ms of listeners hearing the coda consonant. That is, listeners used biphone probability with little to no delay from when they receive information about the coda consonants, consistent with findings from Kingston et al., (2016). Note that the absolute value of the timing of the effect in this experiment is about 100 ms longer than that reported in Kingston et al., (2016), which is consistent with the difference in vowel duration across the our stimuli and theirs (260ms here compared to 170ms in Kingston et al., Experiment 2a).

We can additionally see that the influence of biphone probability varied based on continuum step. In fact, model fit was substantially improved by including frame in the non-linear interaction between step and continuum, as compared to being a separate smooth that did not interact with step ($\chi^2(7) = 23.31$; $p < 0.001$) We used the `compareML()` function from *itsadug* to assess this. Specifically, the biphone probability effect was more rapidly integrated when vowel information was more /ε/-like (Step 6 and 7) facilitating looks to the /ε/ target and accordingly showing divergence between smooths earlier in time.

Finally, we turn to the neighborhood density manipulation, comparing /b_b/ and /b_p/ frames in Figure 6, panel B. Looks did not diverge at any point in the analysis window. That is, we did not observe an ND effect online, lining up with findings from listeners' click responses reported previously. In Experiment 4, we probed the unexpected difference between the /m/- and /b/-initial frames further.

5 Experiment 4

Recall that the stronger /æ/ bias for /b/-initial frames observed in Experiment 3 was consistent with the small neighborhood density difference favoring /b/-initial compared to /m/-initial frames. However, its early timing as well as the difference in magnitude of the effect compared to the neighborhood density effect observed in Experiment 2 led us to hypothesize that this effect was not driven by the neighborhood density differences. Instead, we conjectured that the frame effect was driven by perceptual adjustments related to nasal consonants and their effects on judgements of vowel height. Experiment 4 was designed to confirm that the difference between /m/- and /b/-initial frames seen in Experiment 3 was unrelated to neighborhood density and biphone probability. In Experiment 4 we presented listeners with another /m/-initial and /b/-initial frame where *both* biphone probability and neighborhood density predicted the opposite of the observed difference between /m/- and /b/-initial frames seen in Experiment 3. If we replicated the nasal vs oral frame effect from Experiment 3 here, we could be sure that it was not driven by either biphone probability or neighborhood density differences.

5.1 Materials

The frames used in Experiment 4 were /m_v/ (used in Experiment 1) and /b_v/. To create the new /b_v/ frames, the initial /b/ from the continua used in Experiment 2 was cross spliced, replacing the /m/ in the /m_v/ frames. As shown in Table 1, both biphone probability and neighborhood density predict that an /m_v/ should show increased /æ/ responses relative to the /b_v/ frame. This is the opposite of the effect seen in Experiment 3 (where the /b_v/ frame showed increased /æ/ responses), and accordingly, we can test if the effect observed there is independent of both biphone probability and neighborhood density.

5.2 Participants and procedure

Thirty-two self-identified monolingual English-speaking participants were recruited to participate in Experiment 4. Unlike previous experiments, these participants were recruited online, via the platform Prolific, and completed the experiment over the internet (due to lab closure driven by the COVID-19 pandemic). Participants were instructed to complete the experiment seated in a quiet room with a pair of headphones. Participants were paid \$4 for this experiment which took 15-20 minutes to complete. The experimental procedure was otherwise identical to that in Experiments 1 and 2.

FIG. 7 HERE

5.3 Results and discussion

Listeners' categorization responses were assessed by the same method and model structure as used in previous experiments. In contrast coding the frames, /m_v/ was mapped to -0.5 and /b_v/ was mapped to 0.5. Continuum step had a credible effect on responses, as seen in all previous experiments ($\beta = -3.85$, CI = [-4.55,-3.18]). Additionally, consonant frame had a credible effect ($\beta = 0.68$, CI = [-0.03,1.35]), replicating the effect observed in Experiment 3. Listeners showed increased /æ/ responses for the /b_v/ frame, as shown in Figure 6. The interaction between frame and step was not credible ($\beta = 0.19$, CI = [-0.30,0.73]). This replicates the effect of initial consonant seen in Experiment 3. Because the direction of the effect in this experiment, despite opposing neighborhood density and biphone probability effects confirms that the robust difference between /m/-initial and /b/-initial frames in Experiment 3 was not driven by differences in neighborhood density (or biphone probability).

7 General discussion

In four experiments, we tested how differences in biphone probability and neighborhood density influence listeners' categorization of a vowel continuum embedded in nonwords. In Experiment 1, we found that listeners shifted categorization to form a high probability sequence even when stimuli were controlled for neighborhood density. Likewise, in Experiment 2, we found that listeners shifted

1
2
3 categorization to favor a denser neighborhood even when stimuli were controlled for biphone
4
5 probability. Finally, in Experiment 3, we found evidence for a robust and early influence of biphone
6
7 probability. In contrast, density effects affected neither categorization nor looking behavior. In one
8
9 additional experiment, we probed an unexpected influence uncovered in Experiment 3. This effect
10
11 resulted from mixing the stimuli from Experiment 1 and 2, and was not driven by biphone probability or
12
13 neighborhood density; instead, it was due to context effects related to a preceding nasal consonant.
14
15

16
17 Our results provide both direct and indirect evidence for a dissociation between biphone
18
19 probability and neighborhood density effects. In Experiments 1 and 2 we showed that biphone
20
21 probability and neighborhood density effects exert an independent influence on offline categorization.
22
23 That is, despite the correlation between biphone probability and neighborhood density in English,
24
25 biphone probability effects in phonetic processing cannot be explained by differences in neighborhood
26
27 activation alone as we show in Experiment 1. Similarly, neighborhood density effects in phonetic
28
29 processing can also not be explained by differences in biphone probabilities alone, as we show in
30
31 Experiment 2.
32
33
34
35
36

37 Categorization data from Experiment 3 also provided evidence for a dissociation, albeit indirect.
38
39 In Experiment 3, the mixing of stimuli from Experiments 1 and 2 increased the variability of frames
40
41 (which had a clear effect on responses as confirmed in Experiment 4). Despite the inclusion of more
42
43 variable frames in Experiment 3, the biphone probability effect on categorization replicated the effect
44
45 observed in Experiment 1 with fewer frames. However, neighborhood density influences, which were
46
47 fairly small in magnitude in Experiment 2, disappeared when an irrelevant dimension of variation (in the
48
49 initial consonant) was introduced into Experiment 3. That is, biphone probability effects were robust
50
51 across online and offline tasks, and not affected by the increased variability in Experiment 3. In
52
53 comparison, the increased variability in frames and task complexity in Experiment 3 led listeners to
54
55 disregard neighborhood density differences in the stimuli. Together, these categorization results are
56
57
58
59
60
61
62
63
64
65

consistent only with accounts where both biphone probability and neighborhood density independently influence processing, albeit in qualitatively distinct ways.

Independent contributions of BP and ND effects seen here provide clear constraints on existing models of spoken word recognition. This is problematic for models like TRACE that do not independently represent biphone probability information (cf. Pitt & McQueen, 1996). This is also incompatible with Norris et al.'s (2000) proposal that neighborhood density effects are rooted in biphone probability differences, and therefore suggests that for models like Merge (Norris, 1999, Norris et al., 2000) and Shortlist (Norris, 1994), information from the lexicon must be used to model effects of ND on phonetic processing.

The categorization results from Experiment 3 also suggest that biphone probability and neighborhood density affect processing at different times. A general consensus in the literature is that early influences in processing are not impacted by task factors (e.g. Miller & Dexter 1988), including the presence of orthogonal variation in stimuli of the kind introduced in Experiment 3 (Green, Tomiak & Kuhl, 1997), as well as cognitive load (Bosker et al., 2017). Thus, based on robustness across tasks and stimulus variability, it is likely that biphone probability, but not neighborhood density, affects processing early.

The eye-tracking data from Experiment 3 confirmed that biphone probability effects are indeed early; biphone probability information is incorporated within 300 to 400 ms after the onset of the vowel, about 130 ms after the spectral response to the continuum itself. Further, when the vowel formants were more /ε/-like, biphone probability information was integrated earlier in processing. The time course of the biphone probability effect in our experiments is similar to the timing of the effect in Kingston et al.'s (2016) findings. In their experiments, as well as in Experiment 3, biphone probability effects emerged less than 50ms into the coda consonant. Given that our biphone probability results cannot be attributed to lexical factors because our continuum endpoints were non-words, and neighborhood density was

controlled, we take these aligning time course results to strengthen our argument that biphone probability differences are responsible for Kingston et al.'s findings in Experiment 2a.

The independence of the biphone probability effect and its early timing precludes biphone probability effects from being an epiphenomenon of lexical feedback (cf. Newman et al. 1997). Instead, the rapid, independent biphone probability effects observed here are consistent with proposals that biphone probability is encoded prelexically, and varies as a function of the robustness of the speech signal (Pitt & McQueen, 1998; Pytkänen, Stringfellow, & Marantz, 2002; Norris et al. 2000).

While we were able to disassociate BP and ND effects, we were unable to delineate the time course of the ND effect. By virtue of its fragility, the increased variability introduced in the eye-tracking experiment (Experiment 3) due to the inclusion of multiple frames erased the small effect of ND observed in the offline task in Experiment 2. In the absence of precise timing information about ND effects, we cannot conclude whether ND effects are a consequence of feedforward activation to decision nodes, or due to feedback effects from the lexicon. However, based on the lack of robustness of ND effects, we can rule out the possibility that ND affects processing as early as BP. Future eye tracking experiments will be required to confirm if they are as delayed as might be expected if they are a result of feedback (Newman et al., 1997), or only moderately so as expected if they feedforward to decision nodes (Norris et al., 2000).

In conclusion, we present new evidence for the dissociation of biphone probability and neighborhood density effects using a combination of categorization and online processing measured with eye tracking. Our results offer support for the claim that biphone probability influences in perception are independent from that of neighborhood density, prelexical in nature, and operate early in processing. Based on these results we argue in favor of models that encode both biphone probability and neighborhood density, albeit with asynchronous timing effects on early processing. Further research will be needed to establish a precise time course for neighborhood density effects, and from extending the

present findings to determine how they combine with other known influences, such as word-hood and word frequency.

Open access: The stimuli and data for all experiments and the analysis code are accessible from the Open Science Foundation website, at https://osf.io/eba2v/?view_only=b76743b2afb34c6c8f46dd7bd2050899

References:

- Archer, S. L., & Curtin, S. (2016). Nine-month-olds use frequency of onset clusters to segment novel words. *Journal of Experimental Child Psychology, 148*, 131-141.
- Barr, D. J. (2008). Analyzing ‘visual world’ eye tracking data using multilevel logistic regression. *Journal of memory and language, 59*(4), 457-474.
- Baayen, R. H., van Rij, J., de Cat, C., & Wood, S. (2018). Autocorrelated errors in experimental data in the language sciences: Some solutions offered by Generalized Additive Mixed Models. In *Mixed-effects regression models in linguistics* (pp. 49-69). Springer, Cham.
- Beddor, P. S. (1993). The perception of nasal vowels. In *Nasals, nasalization, and the velum* (pp. 171-196). Academic Press.
- Bürkner P. (2018). “Advanced Bayesian Multilevel Modeling with the R Package brms.” *The R Journal, 10*(1), 395–411.
- Bosker, H. R., Reinisch, E., & Sjerps, M. J. (2017). Cognitive load makes speech sound fast, but does not modulate acoustic context effects. *Journal of Memory and Language, 94*, 166–176.
- Chen, N. F., Slifka, J. L., & Stevens, K. N. (2007). Vowel nasalization in American English: acoustic variability due to phonetic context. *Speech Communication, 905-918*.

- 1
2
3 Dahan, D., Magnuson, J. S., Tanenhaus, M. K., & Hogan, E. M. (2001). Subcategorical mismatches and
4
5 the time course of lexical access: Evidence for lexical competition. *Language and Cognitive*
6
7 *Processes*, 16(5-6), 507-534.
8
9
10 Diehl, R. L., Kluender, K. R., & Walsh, M. A. (1990). Some auditory bases of speech perception and
11
12 production. *Advances in speech, hearing and language processing*, 1, 243-268.
13
14
15 Fox, R. A. (1984). Effect of lexical status on phonetic categorization. *Journal of Experimental*
16
17 *Psychology. Human Perception and Performance*, 10(4), 526-540.
18
19
20 Friederici, A. D., & Wessels, J. M. (1993). Phonotactic knowledge of word boundaries and its use in
21
22 infant speech perception. *Perception & psychophysics*, 54(3), 287-295.
23
24
25 Frisch, S. A., Large, N. R., & Pisoni, D. B. (2000). Perception of wordlikeness: Effects of segment
26
27 probability and length on the processing of nonwords. *Journal of memory and language*, 42(4),
28
29 481-496.
30
31
32 Garlock, V. M., Walley, A. C., & Metsala, J. L. (2001). Age-of-acquisition, word frequency, and
33
34 neighborhood density effects on spoken word recognition by children and adults. *Journal of*
35
36 *Memory and language*, 45(3), 468-492.
37
38
39 Gathercole, S. E., Frankish, C. R., Pickering, S. J., & Peaker, S. (1999). Phonotactic influences on short-
40
41 term memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 25(1),
42
43 84.
44
45
46 Green, K. P., Tomiak, G. R., & Kuhl, P. K. (1997). The encoding of rate and talker information during
47
48 phonetic perception. *Perception & Psychophysics*, 59(5), 675-692.
49
50
51 Hollich, G., Jusczyk, P. W., & Luce, P. A. (2002, November). Lexical neighborhood effects in 17-month-
52
53 old word learning. In *Proceedings of the 26th annual Boston University conference on language*
54
55 *development* (Vol. 1, pp. 314-23). Boston, MA: Cascadilla Press.
56
57
58 Jusczyk, P. W., Luce, P. A., & Charles-Luce, J. (1994). Infants' sensitivity to phonotactic patterns in the
59
60 native language. *Journal of Memory and Language*, 33(5), 630.
61
62
63
64
65

- Kingston, J., Levy, J., Rysling, A., & Staub, A. (2016). Eye movement evidence for an immediate Ganong effect. *Journal of Experimental Psychology: Human Perception and Performance*, 42(12), 1969–
- Landauer, T. K., & Streeter, L. A. (1973). Structural differences between common and rare words: Failure of equivalence assumptions for theories of word recognition. *Journal of Verbal Learning and Verbal Behavior*, 12(2), 119–131.
- Luce, P. A., & Pisoni, D. B. (1998). Recognizing spoken words: The neighborhood activation model. *Ear and hearing*, 19(1), 1.
- Maeda, S. (1993). Acoustics of vowel nasalization and articulatory shifts in French nasal vowels. In *Nasals, nasalization, and the velum* (pp. 147-167). Academic Press.
- Maslowski, M., Meyer, A. S., & Bosker, H. R. (2020). Eye-tracking the time course of distal and global speech rate effects. *Journal of Experimental Psychology: Human Perception and Performance*.
- Massaro, D. W. (1989). Testing between the TRACE model and the fuzzy logical model of speech perception. *Cognitive psychology*, 21(3), 398-421.
- Massaro, D. W. & Cowan, N. (1993) Information processing models: Microscopes of the mind. *Annual Review of Psychology* 44:383–425
- Matin, E., Shao, K. C., & Boff, K. R. (1993). Saccadic overhead: Information-processing time with and without saccades. *Perception & psychophysics*, 53(4), 372-380.
- Mattys, S. L., & Jusczyk, P. W. (2001). Phonotactic cues for segmentation of fluent speech by infants. *Cognition*, 78(2), 91-121.
- Mattys, S. L., Jusczyk, P. W., Luce, P. A., & Morgan, J. L. (1999). Phonotactic and prosodic effects on word segmentation in infants. *Cognitive psychology*, 38(4), 465-494.
- McClelland, J. L., & Elman, J. L. (1986). The TRACE model of speech perception. *Cognitive psychology*, 18(1), 1-86.

- 1
2
3 Miller, J. L., & Dexter, E. R. (1988). Effects of speaking rate and lexical status on phonetic
4 perception. *Journal of Experimental Psychology: Human Perception and Performance*, 14(3),
5
6 369.
7
8
9
10 Mitterer, H., & Reinisch, E. (2013). No delays in application of perceptual learning in speech
11 recognition: Evidence from eye tracking. *Journal of Memory and Language*, 69(4), 527-545.
12
13
14
15 Munson, B., Swenson, C. L., & Manthei, S. C. (2005). Lexical and phonological organization in
16 children. *Journal of Speech, Language, and Hearing Research*.
17
18
19
20 Munson, B., Edwards, J., & Beckman, M. E. (2005). Phonological knowledge in typical and atypical
21 speech–sound development. *Topics in language disorders*, 25(3), 190-206.
22
23
24
25 Newman, R. S., Sawusch, J. R., & Luce, P. A. (1997). Lexical neighborhood effects in phonetic
26 processing. *Journal of Experimental Psychology: Human Perception and Performance*, 3(23),
27
28 873–889.
29
30
31
32 Nixon, J. S., van Rij, J., Mok, P., Baayen, R. H., & Chen, Y. (2016). The temporal dynamics of perceptual
33 uncertainty: eye movement evidence from Cantonese segment and tone perception. *Journal of*
34
35 *Memory and Language*, 90, 103-125.
36
37
38
39 Norris, D. (1994). Shortlist: a connectionist model of continuous speech recognition. *Cognition*, 52, 189–
40
41
42 Norris, D., McQueen, J. M., & Cutler, A. (2000). Merging information in speech recognition: Feedback
43 is never necessary. *Behavioral and Brain Sciences*, 23(3), 299–325.
44
45
46
47 Ohala, J. J., Beddor, P. S., Krakow, R. A., & Goldstein, L. M. (1986). Perceptual constraints and
48 phonological change: a study of nasal vowel height. *Phonology*, 3, 197-217.
49
50
51
52 Pierrehumbert, J. B., Needle, J., & Hay, J. B. (2018). Phonological and morphological effects in the
53 acceptability of pseudowords (A. Sims & A. Ussishkin, Eds.). Cambridge University Press.
54
55
56
57 Pitt, M. A., & McQueen, J. M. (1998). Is Compensation for Coarticulation Mediated by the Lexicon?
58
59 *Journal of Memory and Language*, 39(3), 347–370.
60
61
62
63
64
65

- 1
2
3 Pykkänen, L., Stringfellow, A., & Marantz, A. (2002). Neuromagnetic evidence for the timing of lexical
4
5 activation: An MEG component sensitive to phonotactic probability but not to neighborhood
6
7 density. *Brain and language*, 81(1-3), 666-678.
8
9
- 10 Reinisch, E., & Sjerps, M. J. (2013). The uptake of spectral and temporal cues in vowel perception is
11
12 rapidly influenced by context. *Journal of Phonetics*, 41(2), 101-116.
13
14
- 15 Roodenrys, S., & Hinton, M. (2002). Sublexical or lexical effects on serial recall of nonwords?. *Journal*
16
17 *of Experimental Psychology: Learning, Memory, and Cognition*, 28(1), 29.
18
19
- 20 Sebastián-Gallés, N., & Bosch, L. (2002). Building phonotactic knowledge in bilinguals: Role of early
21
22 exposure. *Journal of Experimental Psychology: Human Perception and Performance*, 28(4),
23
24 974.
25
26
- 27 Stevens, K. (1998). Acoustic phonetics. Cambridge, MA: MIT Press.
28
29
- 30 Swingle, D., & Aslin, R. N. (2002). Lexical neighborhoods and the word-form representations of 14-
31
32 month-olds. *Psychological science*, 13(5), 480-484.
33
34
- 35 Thorn, A. S., & Frankish, C. R. (2005). Long-term knowledge effects on serial recall of nonwords are
36
37 not exclusively lexical. *Journal of Experimental Psychology: Learning, Memory, and*
38
39 *Cognition*, 31(4), 729.
40
41
- 42 van Rij J., Weling M., Baayen R., van Rijn H. (2020). itsadug: Interpreting Time Series and
43
44 Autocorrelated Data Using GAMMs. R package version 2.4.
45
46
- 47 Vitevitch, M. S. (2002a). Naturalistic and experimental analyses of word frequency and neighborhood
48
49 density effects in slips of the ear. *Language and speech*, 45(4), 407-434.
50
51
- 52 Vitevitch, M. S. (2002b). The influence of phonological similarity neighborhoods on speech
53
54 production. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 28(4), 735.
55
56
- 57 Vitevitch, M. S., Armbrüster, J., & Chu, S. (2004). Sublexical and lexical representations in speech
58
59 production: Effects of phonotactic probability and onset density. *Journal of Experimental*
60
61 *Psychology: Learning, Memory, and Cognition*, 30(2), 514.
62
63
64
65

- 1
2
3 Vitevitch, M. S., & Luce, P. A. (1998). When Words Compete: Levels of Processing in Perception of
4
5 Spoken Words. *Psychological Science*, 9(4), 325–329. Retrieved from JSTOR.
6
7
8 Vitevitch, M. S., & Luce, P. A. (1999). Probabilistic Phonotactics and Neighborhood Activation in
9
10 Spoken Word Recognition. *Journal of Memory and Language*, 40(3), 374–408.
11
12
13 Vitevitch, M. S., & Luce, P. A. (2004). A web-based interface to calculate phonotactic probability for
14
15 words and nonwords in English. *Behavior Research Methods, Instruments, & Computers*, 36(3),
16
17 481-487.
18
19
20 Vitevitch, M. S., Luce, P. A., Pisoni, D. B., & Auer, E. T. (1999). Phonotactics, Neighborhood
21
22 Activation, and Lexical Access for Spoken Words. *Brain and Language*, 68(1), 306–311.
23
24
25 Winn, M. (2016). Praat script: Make formant continuum [Computer software]. Retrieved 15 January,
26
27 2018, from <http://www.mattwinn.com/praat.html>.
28
29
30 Wood SN (2011). Fast stable restricted maximum likelihood and marginal likelihood estimation of
31
32 semiparametric generalized linear models. *Journal of the Royal Statistical Society (B)*, 73(1), 3-
33
34 36.
35
36
37 Wright, J. (1980). The behavior of nasalized vowels in the perceptual vowel space. *Report of the*
38
39 *Phonology Laboratory Berkeley, Cal*, (5), 127-163.
40
41
42 Zahner, K., Kutscheid, S., & Braun, B. (2019). Alignment of f0 peak in different pitch accent types
43
44 affects perception of metrical stress. *Journal of Phonetics*, 74, 75-95.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

TABLES AND FIGURES

Table 1: Lexical statistics and biases for the continuum endpoints used in the experiments. See section 2 for details on calculation of biphone probability (BP) and neighborhood density (ND). Endpoint biases are calculated with /æ/ as reference; thus, positive numbers favor greater /æ/ responses. The absolute difference between endpoint responses is given below each endpoint pair in bold.

Experiment 1	BP		ND	
	C ₁ V ₂	V ₂ C ₃	C ₁ V ₂	CVC
/mæb/	0.0101	0.0026	31.95	29.54
/mɛb/	0.0059	0.0007	31.06	17.96
Extent of /æ/ <i>bias</i>	<i>0.0042</i>	<i>0.0019</i>	<i>0.89</i>	<i>11.58</i>
/mæv/	0.0101	0.0019	31.95	30.25
/mɛv/	0.0059	0.0026	31.06	17.37
Extent of /æ/ <i>bias</i>	<i>0.0042</i>	<i>-0.007</i>	<i>0.89</i>	<i>12.88</i>
Relative /æ/ <i>bias</i>	identical	0.0026	identical	identical
/mæb/-/mɛb/ > /mæv/-/mɛv/				
Experiment 2	BP		ND	
	C ₁ V ₂	V ₂ C ₃	C ₁ V ₂	CVC
/bæb/	0.0059	0.0026	46.84	41.11
/bɛb/	0.0032	0.0007	33.47	21.26
Extent of /æ/ <i>bias</i>	<i>0.0027</i>	<i>0.0019</i>	<i>13.37</i>	<i>19.85</i>
/bæp/	0.0059	0.0048	46.84	44.42
/bɛp/	0.0032	0.0029	33.47	14.46
Extent of /æ/ <i>bias</i>	<i>0.0027</i>	<i>0.0019</i>	<i>13.37</i>	<i>29.96</i>
Relative /æ/ <i>bias</i>	identical	identical	identical	-10.11
/bæp/-/bɛp/ > /bæb/-/bɛb/				
Experiment 4	BP		ND	
	C ₁ V ₂	V ₂ C ₃	C ₁ V ₂	CVC
/bæv/	0.0059	0.0019	46.84	24.74
/bɛv/	0.0032	0.0026	33.47	15.19
Extent of /æ/ <i>bias</i>	<i>0.0027</i>	<i>-0.007</i>	<i>13.37</i>	<i>9.55</i>
/æ/ <i>bias</i> relative to	<i>-0.0015</i>	identical	12.48	-3.33
/mɛv/ ~ /mæv/	/mæv/-/mɛv/ > /bæv/-/bɛv/			

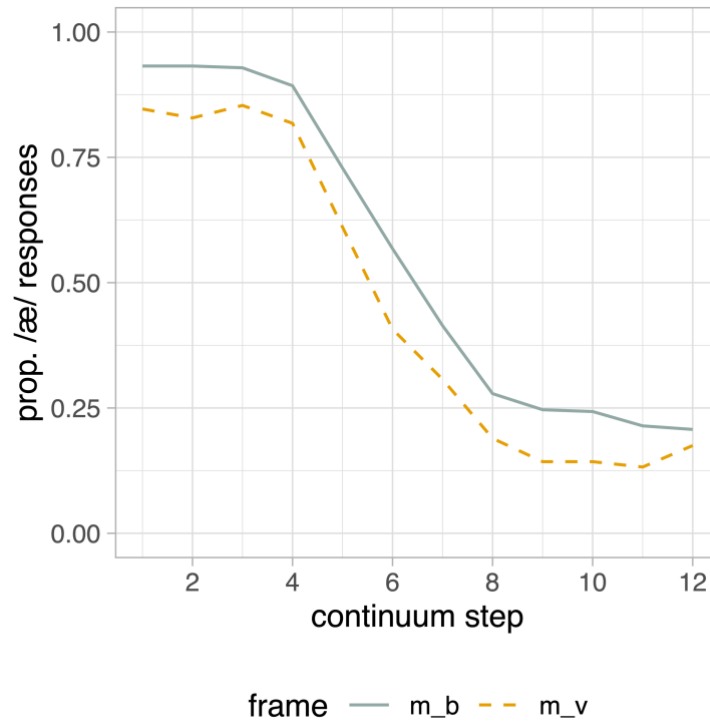


Figure 1: Experiment 1 categorization responses along the continuum (x axis, where step 1 is the most /æ/-like), split by consonant frame. The proportion of /æ/ responses is plotted on the y axis.

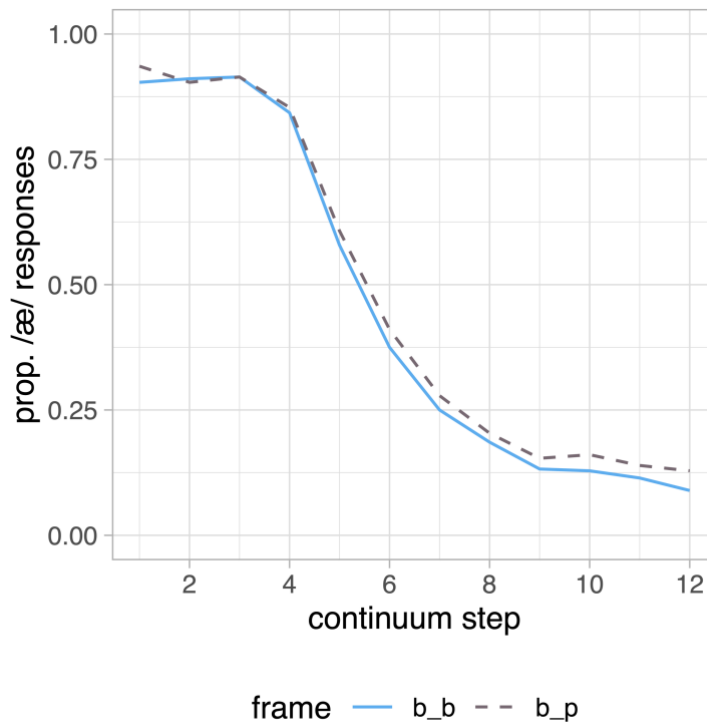


Figure 2: Experiment 2 categorization responses along the continuum, split by consonant frame.

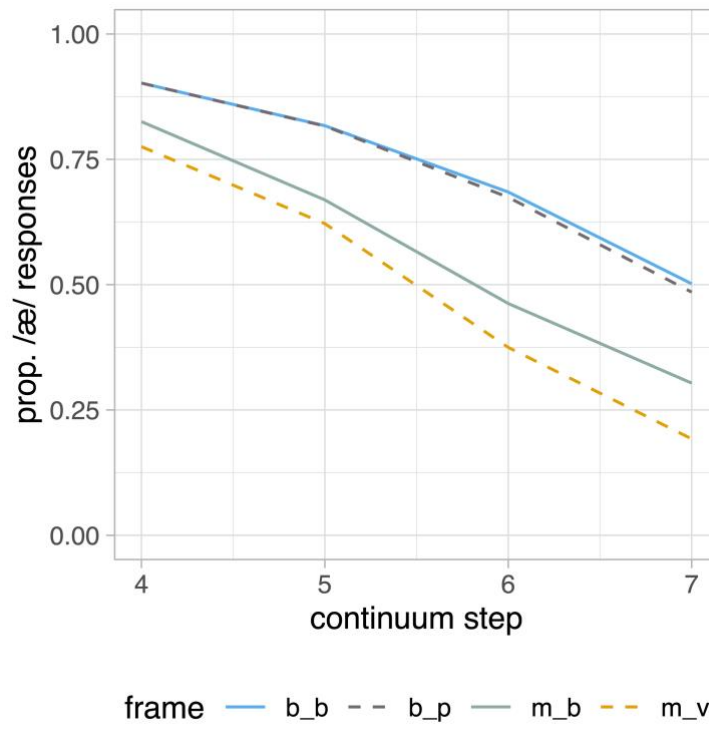


Figure 3: Experiment 3 categorization responses along the continuum, split by consonant frame.

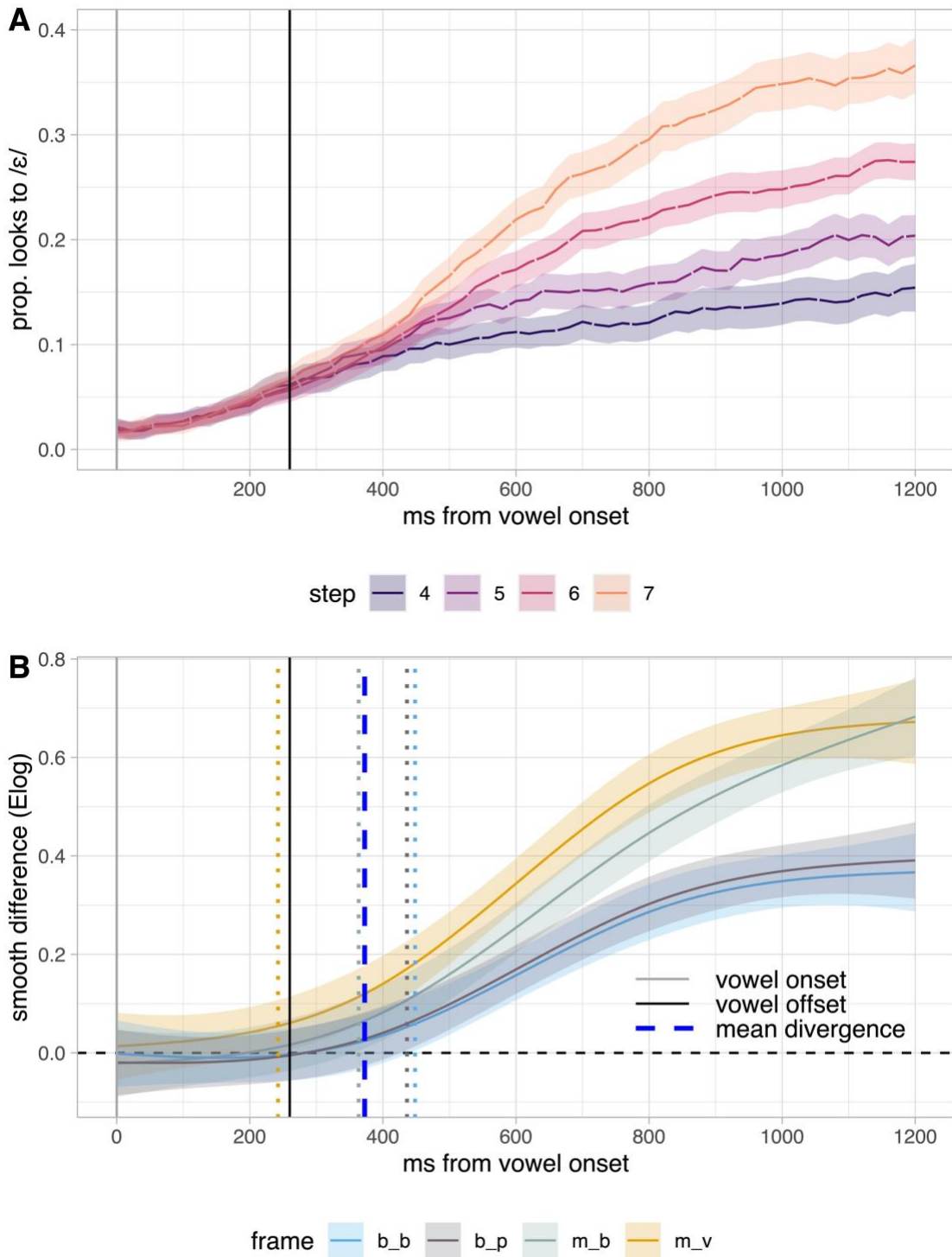


Figure 4: Eye movement data in Experiment 3 showing the effect of continuum step. Panel A plots listeners' looks to /ε/ as a function of continuum step. Temporal landmarks are indexed by vertical lines which are labeled in panel B. Panel B shows difference smooths for the effect of continuum step, split by consonant frame (see text). Dotted vertical lines show the estimated divergence times within each frame (indexed by matching color). The larger blue dashed vertical line shows the mean of these four times.

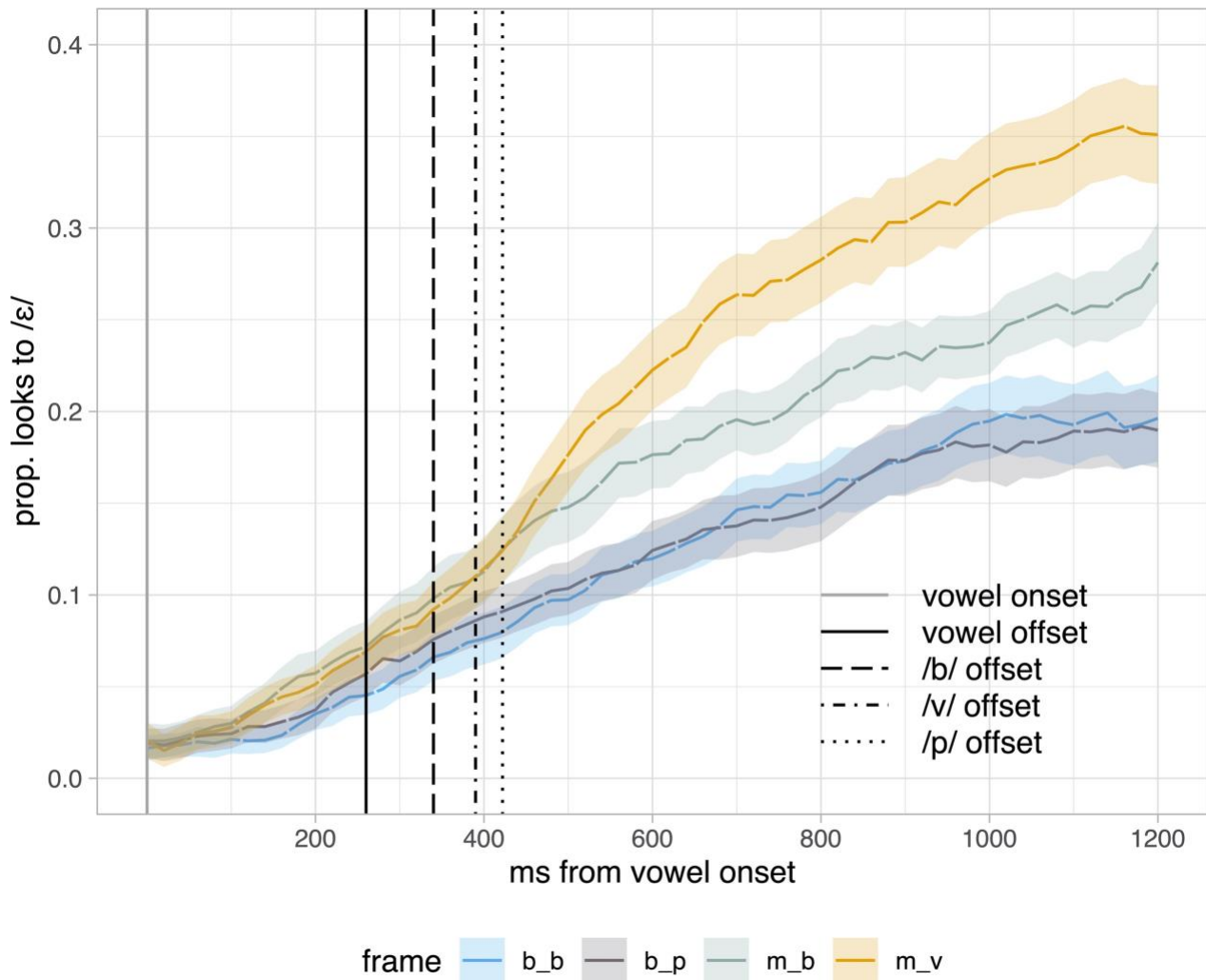
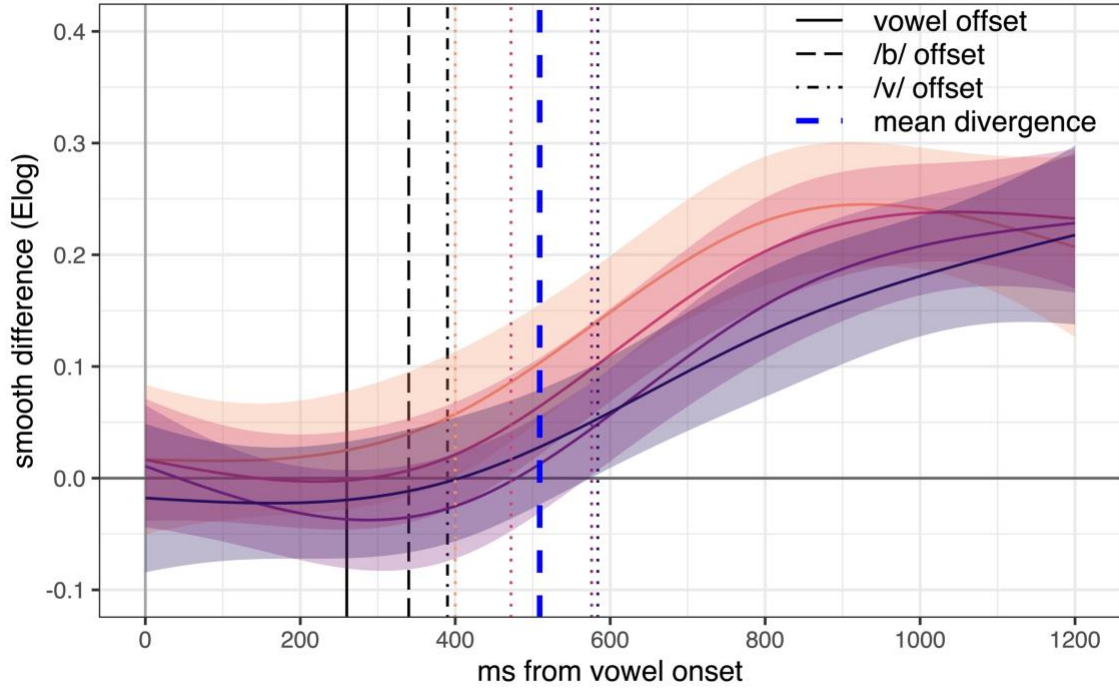


Figure 5: Eye movement data in Experiment 3 showing listeners' looks to /ε/ as a function of consonant frame. Temporal landmarks are indexed by vertical lines which are labeled on the plot.

A BP effect at each continuum step



B ND effect at each continuum step (/b_b/ versus /b_p/)

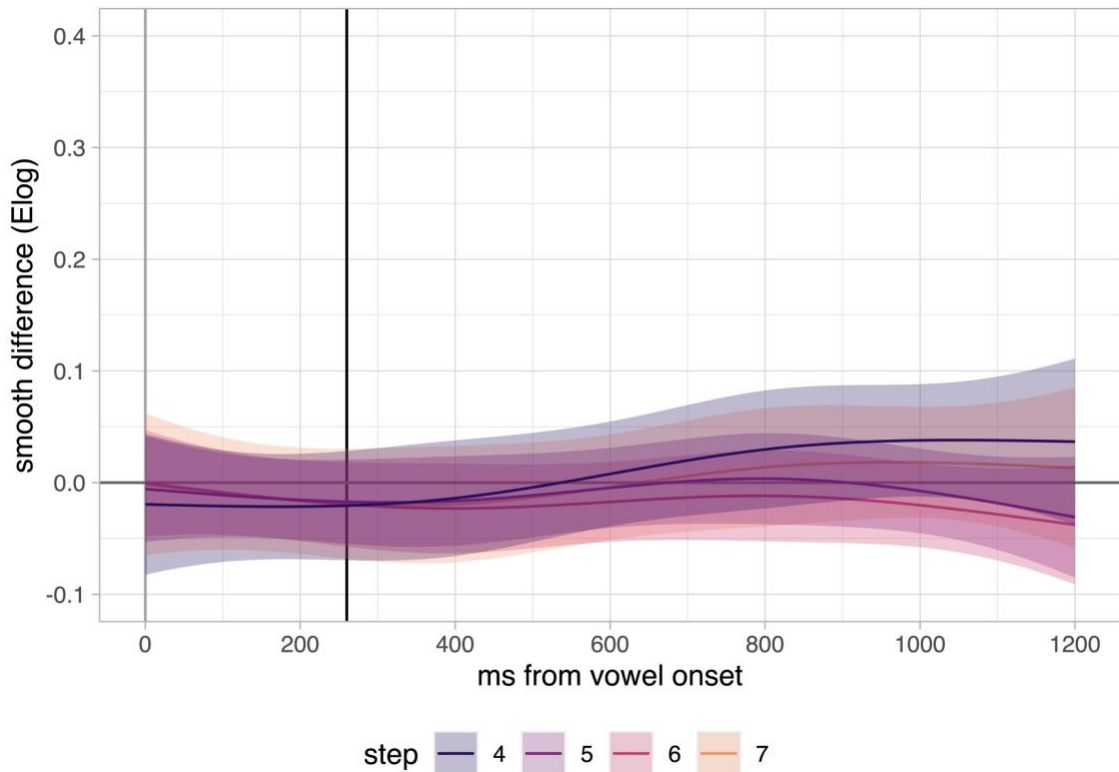


Figure 6: Difference smooths for the effect of consonant frame, split by continuum step (see text). Dotted vertical lines show the estimated divergence times within each frame (indexed by matching color). The larger blue dashed vertical line shows the mean of these four times. Panel A compares the /m_b/ and /m_v/ frames which manipulate biphone probability, while Panel B compares the /b_p/ and /b_b/ frames which manipulate neighborhood density.

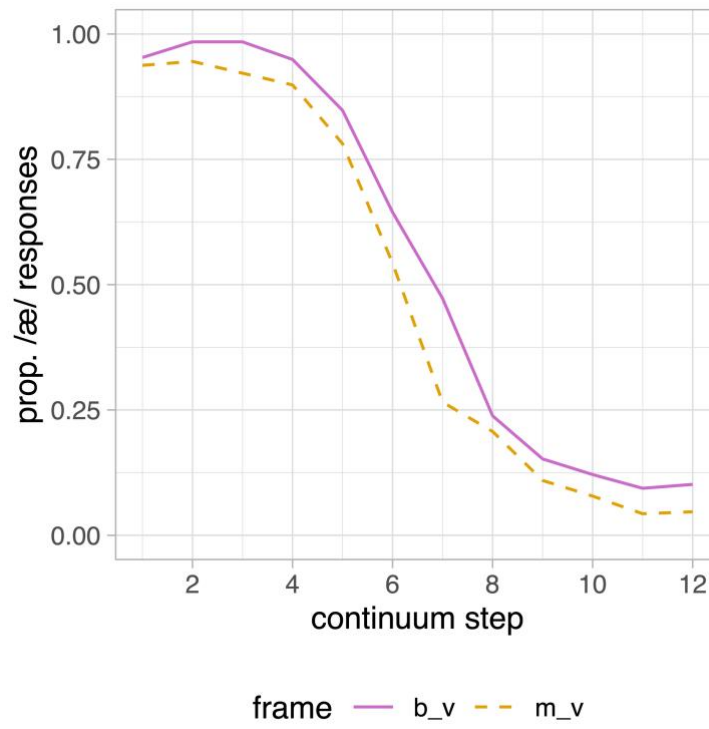


Figure 7: Experiment 4 categorization responses along the continuum, split by consonant frame.

Appendix

Table 2: Output for the GAMM, showing parametric and smooth terms. Note the reference level in the parametric terms in /m_b/. See section 4.4.2 for details.

<i>Parametric terms</i>	Est.	SE	t	p
Intercept	-0.80	0.02	-38.78	< 0.001
continuum (scaled)	0.10	0.02	6.52	< 0.001
frame /b_b/	-0.15	0.01	-12.04	< 0.001
frame /b_p/	-0.13	0.01	-11.09	< 0.001
frame /m_v/	0.05	0.01	4.07	< 0.001
<i>Smooth terms</i>	edf	Ref. df	F	p
te(time, continuum):frame /b_b/	10.20	12.10	19.49	< 0.001
te(time, continuum):frame /b_p/	6.80	7.66	29.87	< 0.001
te(time, continuum):frame /m_b/	9.64	11.59	42.07	< 0.001
te(time, continuum):frame /m_v/	10.50	12.14	62.16	< 0.001
s(time, participant)	498.67	647.00	9.38	< 0.001