

Comparing accounts of formant normalization against US English listeners' vowel perception

Anna Persson,¹ Santiago Barreda,² and T. Florian Jaeger³

¹*Swedish Language and Multilingualism, Stockholm University*^a

²*Linguistics, University of California, Davis*

³*Brain and Cognitive Sciences, Data Science, University of Rochester*

(Dated: 23 July 2024)

1 Human speech perception tends to achieve robust speech recognition, despite sub-
2 stantial cross-talker variability. Believed to be critical to this ability are auditory
3 normalization mechanisms whereby listeners adapt to individual differences in vocal
4 tract physiology in vowel perception. This study asks what types of computations
5 are involved in such normalization. Two 8-way alternative forced-choice experiments
6 assessed L1 listeners' categorizations across the entire US English vowel space—both
7 for unaltered and for synthesized stimuli. Listeners' responses in these experiments
8 were compared against the predictions of twenty influential normalization accounts
9 that differ starkly in the inference and memory capacities they imply for speech
10 perception. Listeners' responses were best explained by *extrinsic* normalization ac-
11 counts, suggesting that listeners learn and store distributional properties of talkers'
12 speech. Of the extrinsic accounts, it was the *computationally least complex* variants
13 that best fit listeners' responses, using a single parameter. These findings have conse-
14 quences for any research that aims to investigate the perceptual, social, and linguistic
15 information of vowel productions. This includes research in phonetics and phonol-
16 ogy, sociolinguistics, and language acquisition. In these fields, it remains common
17 to employ normalization accounts that the present study confirms to be inadequate
18 models of human perception (e.g., Lobanov normalization).

^aanna.persson@su.se

19 **I. INTRODUCTION**

20 One of the central challenges for speech perception originates in cross-talker variability:
 21 depending on the talker, the same acoustic signal can encode different sound categories (Allen
 22 *et al.*, 2003; Liberman *et al.*, 1967; Newman *et al.*, 2001). This results in ambiguity in the
 23 mapping from acoustics to words and meanings. Research has identified several mechanisms
 24 through which listeners resolve this ambiguity, ranging from early perceptual processes, to
 25 adaptation of phonetic categories, all the way to adjustments in post-linguistic decision
 26 processes (for review, see Xie *et al.*, 2023). The present study focuses on the first type of
 27 mechanism, early auditory processes that transform and normalize the acoustic input into
 28 the perceptual cues that constitute the input to linguistic processing (for reviews, Barreda,
 29 2020; Johnson and Sjerps, 2021; McMurray and Jongman, 2011; Stilp, 2020; Weatherholtz
 30 and Jaeger, 2016). We seek to respond, in particular, to recent calls to put theories of
 31 adaptive speech perception to stronger tests (Baese-Berk *et al.*, 2018; Schertz and Clare,
 32 2020; Xie *et al.*, 2023).

33 Evidence for the presence of early normalization mechanisms comes from neuroimaging
 34 and neurophysiological studies (e.g., Oganian *et al.*, 2023; Skoe *et al.*, 2021). These studies
 35 have decoded effects of talker identity from subcortical brain areas like the brain stem, and
 36 thus prior to the cortical regions believed to encode linguistic categories (e.g., Sjerps *et al.*,
 37 2019; Tang *et al.*, 2017). This includes brain responses that lag the acoustic signal by as
 38 little as 20-50 msec (Lee, 2009), suggesting very fast and highly automatic processes. By
 39 removing talker-specific variability from the phonetic signal early, auditory normalization

40 offers elegant and effective solutions to cross-talker variability, that might reduce the need
 41 for more complex adaptation of individual phonetic categories further upstream (Apfelbaum
 42 and McMurray, 2015; Xie *et al.*, 2023).¹

43 While it is relatively uncontroversial *that* normalization contributes to robust speech
 44 perception, it is still unclear what types of computations this implicates. We address this
 45 question for the perception of vowels, which cross-linguistically relies on peaks in the distri-
 46 bution of spectral energy over acoustic frequencies (formants). Vowel perception has long
 47 been a focus in research on normalization (e.g., Bladon *et al.*, 1984; Fant, 1975; Gerstman,
 48 1968; Johnson, 2020; Joos, 1948; Lobanov, 1971; Miller, 1989; Nearey, 1978; Nordström
 49 and Lindblom, 1975; Syrdal and Gopal, 1986; Traunmüller, 1981; Watt and Fabricius, 2002;
 50 Zahorian and Jagharghi, 1991; for review, see Barreda, 2020), with some reviews citing over
 51 100 competing proposals (Carpenter and Govindarajan, 1993). Importantly, these accounts
 52 differ in the types and complexity of computations they assume to take place during nor-
 53 malization. On the lower end of computational complexity, comparatively simple static
 54 transformations of the acoustic signal might suffice to achieve invariance in the mapping
 55 from cues to phonetic categories. For example, there is evidence that a transformation of
 56 acoustic frequencies (measured in Hz) into the psycho-acoustic Mel-space better describes
 57 how listeners perceive differences in the frequency of sine tones (e.g., Fastl and Zwicker,
 58 2007; Stevens and Volkmann, 1940; for a critique, see Greenwood, 1997; Siegel, 1965). It
 59 is thus possible that cross-talker variability in vowel pronunciations is effectively reduced
 60 when formants are represented in Mel, rather than Hz. Similar arguments have been made
 61 about other psycho-acoustic transformations (e.g., Bark, Traunmüller, 1990; ERB, Glasberg

62 and Moore, 1990; or semitones, Fant *et al.*, 2002). Most of these accounts share that they
63 log-transform acoustic frequencies—in line with neurophysiological evidence that the audi-
64 tory representations in the brain seem to follow a roughly logarithmic organization, so that
65 auditory perception is (up to a point) more sensitive to differences between lower frequen-
66 cies than to the same difference between higher frequencies (e.g., Merzenich *et al.*, 1975;
67 for review, see Saenz and Langers, 2014). If such static psycho-acoustic transformations are
68 sufficient for formant normalization, this would offer a particularly parsimonious account of
69 vowel perception as listeners would not have to infer talker-specific properties.

70 The parsimony of psycho-acoustic transformations contrasts with the majority of accounts
71 for vowel normalization, which introduce additional computations. This includes accounts
72 that normalize formants relative to other information that is available at the same point in
73 the acoustic signal (intrinsic normalization, e.g., Miller, 1989; Peterson, 1961; Syrdal and
74 Gopal, 1986). For example, according to one proposal, listeners normalize vowel formants
75 by the vowel’s fundamental frequency or other formants estimated at the same point in
76 time (Syrdal and Gopal, 1986). To the extent that the fundamental frequency is correlated
77 with the talkers’ vocal tract size (for review, see Vorperian and Kent, 2007), this allows
78 the removal of physiologically-conditioned cross-talker variability in formant realizations.
79 While such intrinsic accounts arguably entail more computational complexity than static
80 transformations, they do not require that listeners *maintain* talker-specific estimates over
81 time. This distinguishes intrinsic from extrinsic accounts, which introduce additional com-
82 putational complexity.

83 According to extrinsic accounts, normalization mechanisms infer and store estimates of
84 talker-specific properties that then are used to normalize subsequent speech from that talker
85 (Gerstman, 1968; Lobanov, 1971; Nearey, 1978; Nordström and Lindblom, 1975; Watt and
86 Fabricius, 2002; for review, see also Weatherholtz and Jaeger, 2016). At the upper end of
87 computational complexity, some accounts hold that listeners continuously infer and maintain
88 both talker-specific means for each formant and talker-specific estimates of each formant's
89 variability (Gerstman, 1968; Lobanov, 1971). These estimates are then used to normalize
90 formants, e.g., by centering and standardizing them (essentially z-scoring formants, Lobanov,
91 1971), removing cross-talker variability in the distribution of formant values. There are,
92 however, more parsimonious extrinsic accounts that require inference and maintenance of
93 fewer talker-specific properties. The most parsimonious of these is Nearey's *uniform scaling*
94 account, which assumes that listeners infer and maintain a single talker-specific parameter.
95 This parameter (Ψ) can be thought of as capturing the effects of the talkers' vocal tract
96 length on the spectral scaling applied to the formant pattern produced by a talker (Nearey,
97 1978).² Uniform scaling deserves particular mention here as it is arguably one of the most
98 developed normalization accounts, and rooted in principled considerations about the physics
99 of sound and the evolution of auditory systems (for review, see Barreda, 2020).

100 In summary, hypotheses about the computations implied by formant normalization differ
101 in the flexibility they afford as well as the inference and memory complexity they entail.
102 Considerations about the complexity of inferences—essentially the number of parameters
103 that listeners are assumed to estimate at any given moment in time—arguably gain in
104 importance in light of the speed at which normalization seems to unfold. In the present

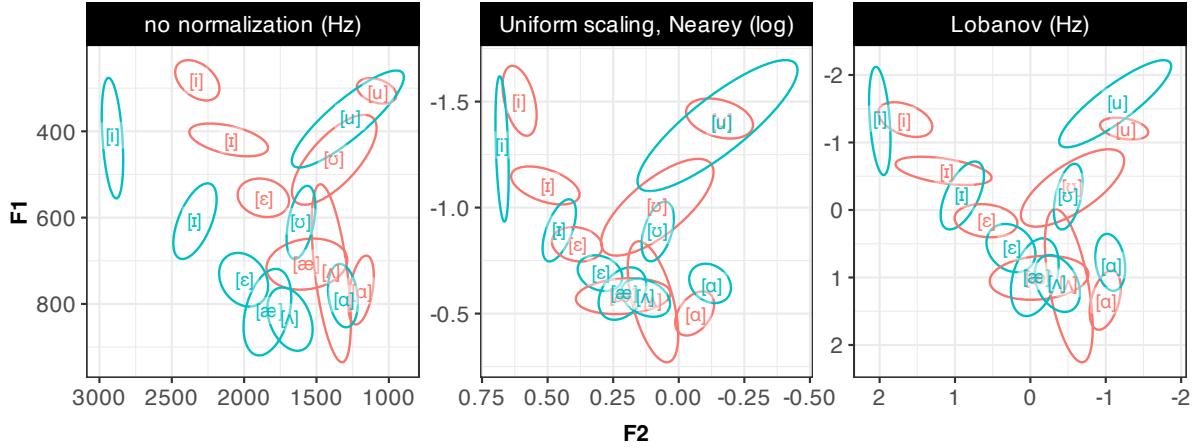


FIG. 1. Illustration of how height, which is positively correlated with vocal tract size, affects vowels' F1 and F2, and how normalization can partially remove this effect. Shown here are realizations of all 8 monophthong vowels of US English by a short (cyan) and a tall native talker (red). **Panel A:** In the acoustic space, prior to any normalization (Hz). **Panel B:** After uniform scaling (Nearey, 1978). **Panel C:** After Lobanov normalization (Lobanov, 1971). The present study compares three of these accounts, along with 17 other normalization accounts.

105 study, we thus ask whether computationally simple accounts are sufficient to explain human
 106 vowel perception.

107 While previous research has compared normalization accounts across languages, most of
 108 this work has evaluated proposals in terms of how well the normalized phonetic space sup-
 109 ports the separability of vowel categories (Adank *et al.*, 2004; Carpenter and Govindarajan,
 110 1993; Cole *et al.*, 2010; Escudero and Bion, 2007; Johnson and Sjerps, 2021; Syrdal, 1985).

111 This approach is illustrated in Figure 1. These studies have found that computationally
 112 more complex accounts—which also afford more flexibility—tend to achieve higher cate-
 113 gory separability and higher categorization accuracy (for review, see Persson and Jaeger,
 114 2023). This includes Lobanov normalization, which continues to be highly influential in,
 115 for example, variationist and sociolinguistic research because of its effectiveness in removing

¹¹⁶ cross-talker variability (for a critique, see Barreda, 2021). It is, however, by no means clear
¹¹⁷ that human speech perception employs the same computations that achieve the best cate-
¹¹⁸ gory separability or accuracy (see also discussion in Barreda, 2021; Nearey and Assmann,
¹¹⁹ 2007).

¹²⁰ A substantially smaller body of research has addressed this question by comparing nor-
¹²¹ malization accounts against *listeners' perception* (Barreda and Nearey, 2012; Barreda, 2021;
¹²² Nearey, 1989; Richter *et al.*, 2017; for a review, see Whalen, 2016). Interestingly, these
¹²³ works seems to suggest that computationally simpler accounts might provide a better fit
¹²⁴ against human speech perception than the influential Lobanov model (Barreda, 2021; Richter
¹²⁵ *et al.*, 2017). For example, Barreda (2021) compared the predictions of uniform scaling and
¹²⁶ Lobanov normalization against listeners' categorization responses in a forced-choice catego-
¹²⁷ rization task over parts of the US English vowel space. In his experiment, listeners' catego-
¹²⁸ rization responses were better predicted by uniform scaling than by Lobanov normalization.
¹²⁹ Findings like these suggest that comparatively simple corrections for vocal tract size—such
¹³⁰ as uniform scaling—might provide a better explanation of human perception than more
¹³¹ computationally complex accounts (see also Johnson, 2020; Richter *et al.*, 2017).

¹³² This motivates the present work. We take a broad-coverage approach by comparing
¹³³ the 20 normalization accounts in Table I against the perception of all 8 monophthongs
¹³⁴ of US English ([i] as in *heed*, [ɪ] in *hid*, [ɛ] in *head*, [æ] in *had*, [ʌ] in *hut*, [ʊ] in *hood*,
¹³⁵ [u] in *who'd*, [ɑ] in *odd*).³ We do so for the perception of both natural and synthesized
¹³⁶ speech. Our broad-coverage approach complements previous studies, which have typically
¹³⁷ compared a small number of accounts (up to 3) and focused on parts of the vowel inventory,

138 and thus parts of the formant space (typically 2-4 vowels, Barreda, 2021; Barreda and
139 Nearey, 2012; Nearey, 1989; Richter *et al.*, 2017). The accounts we consider include the
140 most influential examples of psycho-acoustic transformations (Fant *et al.*, 2002; Glasberg
141 and Moore, 1990; Stevens and Volkmann, 1940; Traunmüller, 1981), intrinsic (Syrdal and
142 Gopal, 1986), extrinsic (Gerstman, 1968; Johnson, 2020; Lobanov, 1971; McMurray and
143 Jongman, 2011; Nearey, 1978; Nordström and Lindblom, 1975), and hybrid accounts that
144 contain intrinsic and extrinsic components (Miller, 1989). This broad-coverage approach
145 allows us to assess, for example, whether the preference for computationally simple accounts
146 observed in Barreda (2021) replicates on new data that span the entire vowel space. It
147 also allows us to ask whether accounts even simpler than uniform scaling—such as psycho-
148 acoustic transformations—provide an even better fit to human perception.

₁₄₉ Next, we motivate and describe the two experiments we conducted. Then we compare
₁₅₀ the normalization accounts in Table I against listeners responses from these experiments.

₁₅₁ **A. Open Science Statement**

₁₅₂ All stimulus recordings, results, and the code for the experiment, data analysis, and
₁₅₃ computational modeling for this article can be downloaded from OSF at <https://osf.io/zemwn/>. The OSF repo also include extensive supplementary information (SI). Both the ar-
₁₅₅ ticle and SI are written in R markdown, allowing readers to replicate our analyses with the
₁₅₆ click of a button, using freely available software (R Core Team, 2023; RStudio Team, 2020).

₁₅₇ Readers can revisit the assumptions we committed to for the present project—for example,
₁₅₈ by substituting alternative normalization accounts or categorization models. Researchers
₁₅₉ can also substitute their own experiments on vowel normalization for our Experiments 1a
₁₆₀ and 1b, to see whether our findings generalize to novel data. We see this as an important
₁₆₁ contribution of the present work, as it should make it substantially easier to consider ad-
₁₆₂ ditional normalization accounts—including variants to the accounts we considered—and to
₁₆₃ assess the generalizability of the conclusions we reach based on the present data.

₁₆₄ **II. EXPERIMENTS 1A AND 1B**

₁₆₅ To compare the performance of different normalization accounts against listeners' percep-
₁₆₆ tion, we conducted two small web-based experiments on US English listeners' perception of
₁₆₇ US English vowels. The two experiments employ the same 8-alternative forced-choice vowel
₁₆₈ categorization task (Figure 2), and differ only in the whether they employed 'natural' (Ex-

TABLE I. Normalization accounts considered in the present study. Unless otherwise marked, formant variables (F s) in the right-hand side of normalization formulas are in Hz.

	Normalization procedure	Perceptual scale	Source	Formula
	n/a	Hz	n/a	n/a
transformation	—	log	—	$F_n^{log} = \ln(F_n)$
	Bark	—	Traunmüller (1990)	$F_n^{Bark} = \frac{26.81 \times F_n}{1960 + F_n} - 0.53$
	ERB	—	Glasberg & Moore (1990)	$F_n^{ERB} = 21.4 \times \log_{10}(1 + F_n) \times 0.00437$
	Mel	—	Stevens & Volkmann (1940)	$F_n^{Mel} = 2595 \times \log_{10}(1 + \frac{F_n}{700})$
	Semitones conversion	—	Fant et al. (2002)	$F_n^{ST} = 12 \times \ln(\frac{F_n}{100})$
	Bark	Syrdal & Gopal (1986)	—	$F1_{SyrdalGopal1} = F1_{Bark} - F0_{Bark}$
	Syrdal & Gopal 1 (Bark-distance model)	—	—	$F2_{SyrdalGopal1} = F2_{Bark} - F1_{Bark}$
	Syrdal & Gopal 2 (Bark-distance model)	—	—	$F1_{SyrdalGopal2} = F1_{Bark} - F0_{Bark}$
	Miller (formant-ratio)	log	Miller (1989)	$F2_{SyrdalGopal2} = F3_{Bark} - F2_{Bark}$
	Miller (formant-ratio)	—	—	$SR = k(\frac{G_{Miller}}{k})^{1/3}$
	—	—	—	$F1_{Miller} = \log(\frac{F1}{SR})$
	—	—	—	$F2_{Miller} = \log(\frac{F2}{F1})$
	—	—	—	$F3_{Miller} = \log(\frac{F3}{F2})$
	log	Nearey (1978)	Nearey (1978)	$F_n^{Nearey} = \ln(F_n) - \text{mean}(\ln(F))$
intrinsic	Uniform scaling, Nearey	Hz	Nordström & Lindblom (1975)	$F_n^{NordströmLindblom} = \frac{F_n}{\text{mean}(\frac{F_n}{F_n \geq 600})}$
	Uniform scaling, Nordström & Lindblom	—	—	$F_n^{Johnson} = \frac{F_n}{\text{mean}(\frac{F1}{0.5}, \frac{F2}{1.5}, \frac{F3}{2.5})}$
	Uniform scaling, Johnson	Hz	Johnson (2020)	$F_n^{Nearey} = \ln(F_n) - \text{mean}(\ln(F_n))$
	Nearey's formantwise log-mean	log	Nearey (1978)	$F_n^{C-CuRE} = F_n - \text{mean}(F_n)$
	C-CuRE	Hz	McMurray & Jongman (2011)	$F_n^{Gerstman} = 999 \times \frac{F_n - F_n^{min}}{F_n^{max} - F_n^{min}}$
	—	Bark	—	$F_n^{Lobanov} = \frac{F_n - \text{mean}(F_n)}{sd(F_n)}$
	—	ERB	—	
	—	Mel	—	
	Semitones conversion	—	Gerstman (1968)	
extrinsic	standardizing	Hz	—	
	Gerstman (range normalization)	—	—	
	Lobanov (z-score)	Hz	Lobanov (1971)	

heed who'd hood

hid  hud

head had hod

FIG. 2. Screen shot of the eight-alternative forced-choice (8-AFC) task used in both Experiment 1a and 1b.

¹⁶⁹ periment 1a) or synthesized stimuli (Experiment 1b). To the best of our knowledge, these
¹⁷⁰ two experiments are the first designed to compare normalization accounts against listeners'
¹⁷¹ perception over the entire monophthong inventory of a language.

¹⁷² Experiment 1a employs recordings of *hVd* word productions from a female talker of US
¹⁷³ English, these recordings are ‘natural’ in the sense that they were not synthesized or other-
¹⁷⁴ wise phonetically manipulated. One consequence of this is that the formant values of these
¹⁷⁵ recordings are clustered around the category means, and thus span only a comparatively
¹⁷⁶ small part of the phonetic space. This can limit the statistical power to distinguish between
¹⁷⁷ competing accounts. Natural recordings furthermore vary not only along the primary cues
¹⁷⁸ to vowel quality in US English (F1, F2) but also along potential secondary cues (e.g., F0,
¹⁷⁹ F3, and vowel duration) as well as other unknown properties, which can make it difficult to
¹⁸⁰ discern whether the performance of a normalization model is due to the normalization itself
¹⁸¹ or other reasons, e.g., because a normalized cue happens to correlate with another cue that
¹⁸² listeners are sensitive to but that is not included in the model.

183 Experiment 1b thus adopts an alternative approach and uses synthesized vowels. Unlike
 184 most previous work, which has used isolated vowels as stimuli (Barreda, 2021; Barreda and
 185 Nearey, 2012; Nearey, 1989; Richter *et al.*, 2017), Experiment 1b uses synthesized *hVd* words
 186 to facilitate comparison to Experiment 1a. This allowed us to sample larger parts of the F1-
 187 F2 space, which has two advantages. First, it allowed us to collect responses over parts of the
 188 formant space for which we expect listeners to have more uncertainty, and thus exhibit more
 189 variable responses. This can increase the statistical power to distinguish between competing
 190 accounts. Second, differences in the predictions of competing normalization account will
 191 tend to become more pronounced with increasing distance from the category centers. By
 192 collecting responses at those locations, we can thus increase the contrast between competing
 193 accounts.

194 The use of resynthesized stimuli does, however, also come with potential disadvantages.
 195 Synthesized stimuli can suffer in ecological validity, lacking correlations between cues, and
 196 across the speech signal (e.g., due to co-articulation) that are characteristic of human speech.
 197 This raises questions about the extent to which processing of such stimuli engages the same
 198 mechanisms as everyday speech perception. Additionally, it is possible that the use of robotic
 199 sounding synthesized speech affects listener engagement. This can lead to an increased rate
 200 of attentional lapses, and thus a decrease in the proportion of trials on which listeners'
 201 responses are based on the acoustics of the speech stimulus rather than random guessing
 202 (compare, e.g., Kleinschmidt, 2020; Tan and Jaeger, 2024). By comparing normalization
 203 accounts against both natural and synthesized stimuli, we investigate the extent to which

204 the accounts that best describe human perception depend on the type of stimuli used in the
205 experiment.

206 **A. Methods**

207 **1. Participants**

208 We recruited 24 (Experiment 1a) and 24 (Experiment 1b) participants from Amazon's
209 Mechanical Turk. Participants were paid \$6/hour prorated by the duration of the exper-
210 iments (15 minutes). Participants only saw the experiment advertised, and could only
211 participate in it, if (i) they were located within the US, (ii) had an approval rating of 99%
212 or higher, (iii) met the software requirements (a recent version of the Chrome browser en-
213 gine), and (iv) had not previously completed any other experiments on vowel perception in
214 our lab. Before the experiment could be accepted, participants had to confirm that they
215 were (i) native speakers of US English (defined as having spent their childhood until the
216 age of 10 speaking English and living in the United States), (ii) in a quiet room without
217 distractions, (ii) wearing over-the-ear headphones. Participants' responses were collected via
218 Javascript developed by the Human Language Processing Lab at the University of Rochester
219 ([Kleinschmidt *et al.*, 2021](#)).

220 An optional post-experiment survey recorded participant demographics using NIH pre-
221 scribed categories, including participant sex (Male: 27, Female: 20), age (mean = 35.5 years;
222 SD = 11.4; 95% quantiles = 24-63.25 years), race (White: 36, Asian: 3, Black: 6, multiple:

²²³ 1, declined to report: 1), and ethnicity (Non-Hispanic: 42, Hispanic: 4, declined to report:
²²⁴ 1).

²²⁵ **2. Materials**

²²⁶ Experiment 1a employed *hVd* word recordings by one adult female talker from a photo-
²²⁷ netically annotated database of L1-US English vowel productions (Xie and Jaeger, 2020).
²²⁸ Specifically, we used all 9 recordings of each of the eight *hVd*-words—*heed*, *hid*, *head*, *had*,
²²⁹ *hut*, *odd*, *who'd*, *hood* (the use of “hut” and “odd” rather than “hud” and “hod” follows
²³⁰ Assmann *et al.*, 2008; but see Hillenbrand *et al.*, 1995).

²³¹ The stimuli for Experiment 1b were synthesized from a single *had* recording used in
²³² Experiment 1a. Specifically, we used a script (based on descriptions in Wade *et al.*, 2007)
²³³ in Praat (Boersma and Weenink, 2022) to concatenate the original /h/ with a synthesized
²³⁴ vowel and the original /d/ recording. Unlike in Experiment 1a, all eight words thus had an
²³⁵ *hVd* context (including “hud” and “hod”, rather than “hut” and “odd”). The Praat script
²³⁶ first segmented the original *had* token into /h/, /ae/ and /d/ portions. It then filtered
²³⁷ the /h/ sound inversely with its LPC, and concatenated this neutral fricative sound with
²³⁸ a complex waveform generated from the pitch and intensity patterns of the original vowel,
²³⁹ to create a neutral hV-section that did not reflect any vocal tract resonances. The script
²⁴⁰ then created a formant grid that filtered the hV-section to create the intended vowel, and
²⁴¹ finally concatenated this segment to the final /d/ to create an *hVd* word. For each *hVd*
²⁴² word, the formant grid was populated with the F1, F2 and F3 values that we handed to the
²⁴³ script at five time-points transitioning from the /h/ to the vowel, to the final /d/ through

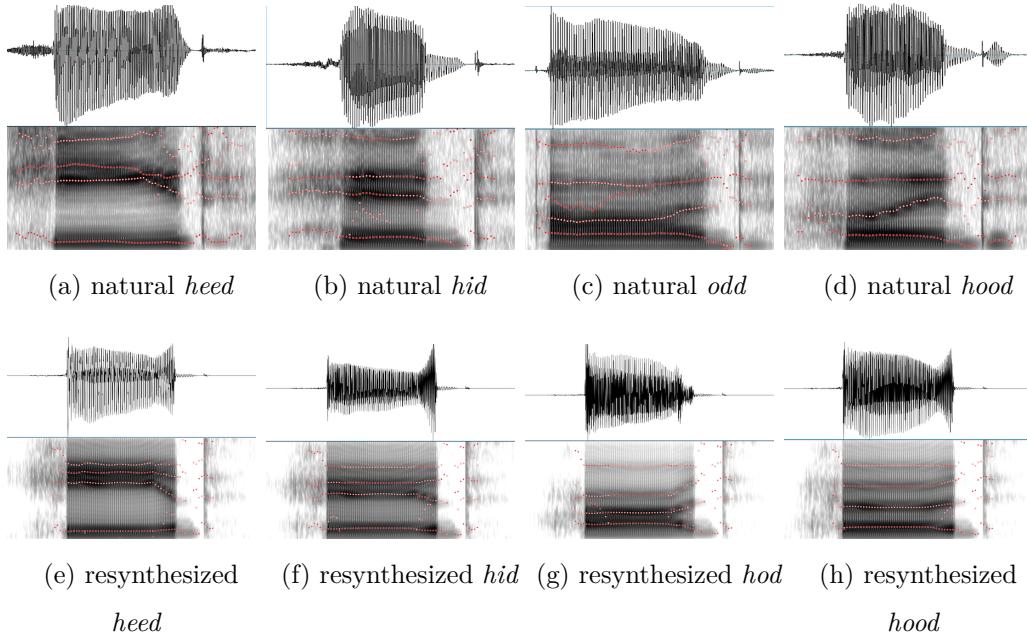


FIG. 3. **Top:** Spectrograms of four natural recordings from Experiment 1a. **Bottom:** Same for four synthesized tokens with similar formant values from Experiment 1b.

244 linear interpolation. Formant bandwidths were 500 Hz at the initial two time-points (the
 245 /h/ and beginning of transition to vowel), and then decreased linearly during vowel onset
 246 and throughout the final three time-points to 50 Hz (F1), 100 Hz (F2), 200 Hz (F3), 300
 247 Hz (F4), and 400 Hz (F5-F8, following Wade *et al.*, 2007). The bandwidth manipulation
 248 implied that formants became stronger as the vowel unfolded (see Figure 3). We used this
 249 approach to create synthesized vowels for arbitrary F1-F2 combinations. F3 was set based
 250 on those F1-F2 values. Specifically, we ran a linear regression over the natural productions
 251 of the talker from Experiment 1a, predicting F3 from F1, F2 and their interaction. We then
 252 used that regression to predict F3 values for any F1-F2 combination in Experiment 1b. F4
 253 to F8, as well as vowel duration, were held identical across all tokens (using the same values
 254 as Wade *et al.*, 2007).

255 We generated 146 synthesized *hVd* recordings that spanned the F1 and F2 space. The
 256 specific F1-F2 locations chosen were determined by a mix of modeling (using ideal observers
 257 described in the next section to predict listeners' categorization responses) and intuition.
 258 Specifically, we selected 64 recordings that we expected to fall within the bivariate 95%
 259 confidence intervals (CIs) of the eight US English monophthongs, and 82 recordings that we
 260 expected to fall between those CIs. Figure 4 under *Results* shows the distribution of stimuli
 261 for both experiments. Of note, our procedure also generated formant combinations that are
 262 physiologically unlikely to have all been produced by the same talker during 'normal' vowel
 263 production (also known as "off-template" instances, Nearey, 1978).

264 **3. Procedure**

265 The procedure for both experiments was identical. Live instances of each experiment
 266 can be found at <https://www.hlp.rochester.edu/experiments/DLPL2S/experiment-A/experiments.html>. At the start of the experiment, participants acknowledged that they
 267 met all requirements and provided consent, as per the Research Subjects Review Board of
 268 the University of Rochester. Before starting the experiment, participants performed a sound
 269 check and signed a consent form. Participants were then instructed to listen to a female talker
 270 saying words, and click on the word on screen to report what word they heard. On each trial,
 271 all eight *hVd*-words were displayed on screen. Half of the participants in each experiment
 272 saw the response options organized as in Figure 2 (resembling the IPA representation of a
 273 vowel space), half saw the response options in the opposite order (flipping top and bottom
 274 and left and right in Figure 2). Each trial started with the response grid on screen, together
 275 and left and right in Figure 2).

276 with a light green dot centered on screen. After 1000 ms, an *hVd* recording played, and
277 participants indicated their response by a mouse-click. After a 1000 ms intertrial interval,
278 the screen reset, and the next trial started.

279 In both experiments, participants heard two blocks of the materials described in the
280 previous sections, for a total of 144 trials in Experiment 1a and 292 trials in Experiment 1b.
281 Presentation within each block was randomized for each participant. Participants were not
282 informed about the block structure of the experiment.

283 After completing the experiment, participants filled out a language background question-
284 naire and the optional demographic survey. On average, participants took 10.3 minutes to
285 complete Experiment 1a ($SD = 6.6$) and 18.4 minutes for Experiment 1b ($SD = 7.3$).

286 **4. Exclusions**

287 We excluded participants who failed to follow instructions and did not wear over-the-ear
288 headphones (as indicated in the post-experiment survey). We also excluded participants
289 with mean (log-transformed) reaction times that were unusually slow or fast (absolute z-
290 score over by-participant means > 3), or if they clearly did not do the task (e.g., by answering
291 randomly). This excluded 6 participants from Experiment 1a and 2 from Experiment 1b
292 (for details, see [§2 A](#)).

293 We further excluded all trials that were unusually fast or slow. Specifically, we first z-
294 scored the log-transformed response times *within each participant* and then z-scored these
295 z-scores *within each trial* across participants. Trials with absolute z-scores > 3 were removed
296 from analysis. This double-scaling approach was necessary as participants' response times

decreased substantially over the first few trials and then continued to decrease less rapidly throughout the remainder of the experiment. The approach removes response times that are unusually fast or slow *for that participant at that trial*, while avoiding specific assumptions about the shape of the speed up in response times across trials. This excluded 1.2% of the trials in Experiment 1a and 1.1% in Experiment 1b. This left for analysis 2565 observations from 18 participants in Experiment 1a, and 6354 observations from 22 participants in Experiment 1b.

304 B. Results

Participants' categorization responses in Experiments 1a and 1b are shown in Figure 4, with larger labels indicating recordings that participants agreed on more.⁴ We make two observations. The first pertains to the degree of (dis)agreement between the two experiments. The second observation pertains to the degree of (dis)agreement across participants within each experiment.

310 1. *Similarities and differences between Experiments 1a and 1b*

Unsurprisingly, participants in both experiments divided the F1-F2 space into the eight vowel categories in ways that qualitatively resembled each other (after taking into account that Experiment 1b covers a larger range of F1-F2 values). Also unsurprisingly, there were some differences between participants' responses across the two experiments, at least when plotted in Hz. For example, [u] rarely was the most frequent response in Experiment 1b, even for stimuli that were predominantly categorized as [u] in Experiment 1a. There are at

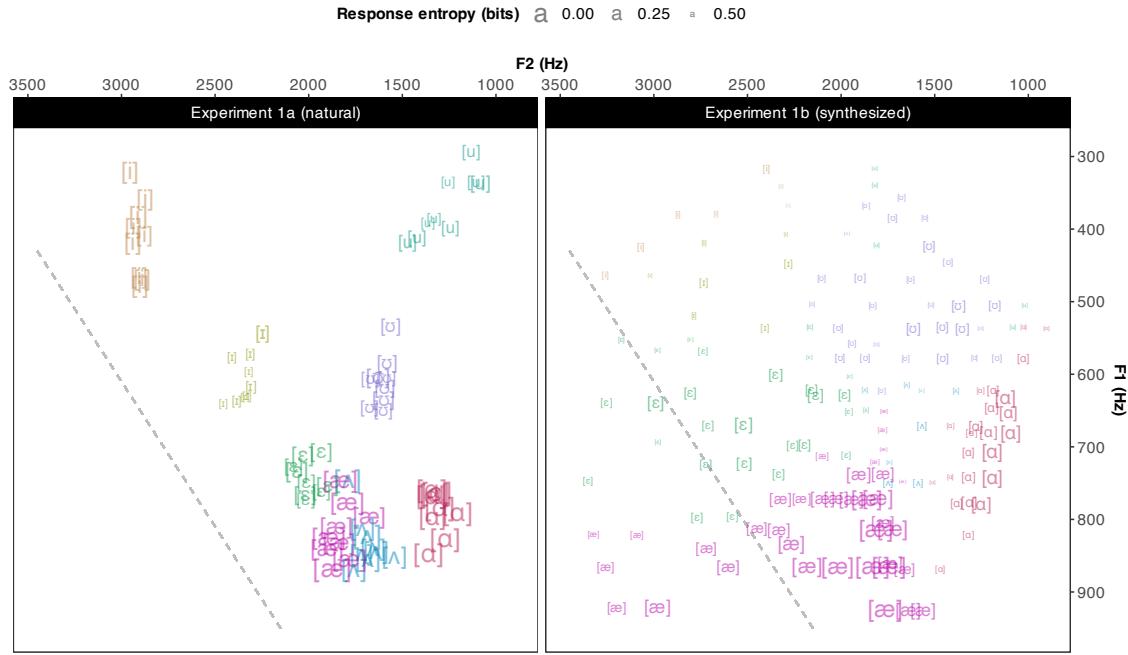


FIG. 4. Summary of listeners' categorization responses in Experiments 1a and 1b in F1-F2 space. The vowel label indicates the most frequent response provided across participants on each test location. Size indicates how consistent responses were across participants, with larger symbols indicating more consistent responses (lower entropy). F1-F2 combinations below the gray dashed line are articulatory unlikely to come from the same talker.

317 least two reasons to expect such differences. First, stimuli with similar F1-F2 values across
 318 the two experiments still differed in other acoustic properties (e.g. vowel duration or F3).
 319 These acoustic differences might have affected participants' responses. Second, it is possible
 320 that *formant normalization* affected participants' responses—i.e., the very mechanism we
 321 seek to investigate in the remainder of the paper. The two experiments differ in the means,
 322 variances, and other statistical properties that some normalization accounts predict to affect
 323 perception. As a consequence, Hz might not be the space in which we should expect identical
 324 responses across experiments.

325 Auxiliary analyses presented in the SI (§2C) suggest that *some but not all* of the dif-
 326 ferences in response entropy between the two experiments were caused by the placement of
 327 the stimuli in F1-F3 space: when comparing categorization responses for tokens from the
 328 two experiments with similar acoustic properties (differences of ≤ 30 Hz along F1 and F2),
 329 response entropies still differed substantially (for $N = 40$ acoustically similar tokens, mean
 330 by-item response entropy for Experiment 1a = 0.18 bits, SE = 0.03; Experiment 1b = 0.39
 331 bits, SE = 0.03). We see two mutually compatible explanations. One possibility is that
 332 the differences in listeners' responses across the two experiments originate in *normalization*.
 333 The two experiments differ in the means, variances, and other statistical properties of their
 334 formant distributions—i.e., in the statistical properties that some normalization accounts
 335 predict to affect perception. It is, however, also possible that the relation between formants
 336 in the synthesized stimuli or some other unknown acoustic-phonetic differences between the
 337 experiments explain the difference in response. For example, the absence of vowel inherent
 338 spectral change (VISC) or differences in tilt in the synthesized stimuli might have deprived
 339 listeners of information that is actually crucial for establishing phonemic identity (Hillen-
 340 brand and Nearey, 1999). This would result in increased uncertainty on each trial, leading
 341 to increased entropy of listeners' responses. The computational studies we present below
 342 shed some light on these two mutually compatible possibilities.

343 Similarly, the two experiments differed in the extent to which participants agreed with
 344 each other. Participants in Experiment 1b exhibited overall less agreement in their responses
 345 (mean by-item response entropy = 0.45 bits, SE = 0.01) than participants in Experiment
 346 1a (mean by-item response entropy = 0.23 bits, SE = 0.02). This was expected given that

347 Experiment 1b explored the entire F1-F2 space, including—by design—formant combina-
348 tions located *between* the centers of the natural vowel categories. Experiment 1b therefore
349 achieved its goal of eliciting less categorical response distributions, which is expected to
350 facilitate comparison of competing normalization accounts.⁵

351 Auxiliary analyses presented in the SI (§2C) suggest that *some but not all* of the dif-
352 ferences in response entropy between the two experiments were caused by the placement of
353 the stimuli in formant space: when comparing categorization responses for tokens from the
354 two experiments with similar acoustic properties (differences of ≤ 30 Hz along F1 and F2),
355 response entropies still differed substantially (for $N = 40$ acoustically similar tokens, mean
356 by-item response entropy for Experiment 1a = 0.18 bits, SE = 0.03; Experiment 1b = 0.39
357 bits, SE = 0.03). We see two mutually compatible explanations. First, similar to the dif-
358 ferences between experiments in the dominant response pattern discussed above, differences
359 in the degree of agreement between participants might originate in *normalization*. Second,
360 it is possible that the relation between formants in the synthesized stimuli or some other
361 unknown acoustic-phonetic differences between the experiments explain the difference in
362 response. For example, the absence of vowel inherent spectral change (VISC) or differences
363 in tilt in the synthesized stimuli might have made it difficult for listeners in Experiment
364 1b to reliably estimate the formants (Hillenbrand and Nearey, 1999). This would result in
365 increased uncertainty on each trial, leading to increased entropy of listeners' responses. The
366 computational studies we present below shed some light on these two mutually compatible
367 possibilities.

368 **2. Similarities and differences between participants**

369 Since the intended category was known for Experiment 1a, it was possible to calculate
 370 participants' recognition accuracy. As also evident in the left panel of Figure 4, participants'
 371 most frequent response *always* matched the intended vowel in Experiment 1a. Overall,
 372 participants' responses matched the intended vowel on 81.2% (SE = 4.8%) of all trials
 373 (Experiment 1b had no such ground truth). This is much higher than chance (12.5%). It is,
 374 however, also quite a bit lower than 100%. To better understand the reasons for this, Figure
 375 5A plots the confusion matrix. This suggests that participants' performance was largely
 376 affected by confusions between [i]-to-[ɛ] (*hid-to-head*), [ɛ]-to-[æ] (*head-to-had*), and [u]-to-[ʊ]
 377 (*who'd-to-hood*).

378 One plausible explanation for this pattern of vowel confusions lies in the substantial
 379 variation that exists across US English dialects (Labov *et al.*, 2006). Differences in the
 380 realization of vowel categories, and associated representations, across dialects will directly
 381 affect the expected classification for any given token. In addition, listeners might differ in
 382 terms of experience with different dialects, or in the dialect they attribute to the talker who
 383 produced the stimuli. To test this hypothesis, we calculated the [i]-to-[ɛ], [ɛ]-to-[æ], and
 384 [u]-to-[ʊ] confusion rates for each participant in Experiment 1a. These data are summarized
 385 in the left panel of Figure 5B. The data in the left panel suggest that most participants in
 386 Experiment 1a either heard [i] tokens consistently as the intended [i] (clustering on the left
 387 side of the panel) or as [ɛ] (clustering on the right side of the panel). Only a few participants
 388 exhibited mixed responses for items intended to be [i]. Tellingly, many of the participants

A

Experiment 1a (natural)											Experiment 1b (synthesized)										
	[i]	[ɪ]	[ɛ]	[æ]	[ʌ]	[ɑ]	[ɔ]	[ʊ]			[i]	[ɪ]	[ɛ]	[æ]	[ʌ]	[ɑ]	[ɔ]	[ʊ]			
Response vowel	[u]	0	0.9	0.3	0	0.6	0	5.1	67		3.9	4.9	0.8	0.6	3.7	4.6	24.2	37.4			
	[v]	0	0.6	0.6	0	0	0	75	30.5		2	5.9	2.9	1.6	17	7.5	45.9	35.2			
	[ɑ]	0	0.3	0.3	1.2	0.6	97.5	0.3	0		1	1.3	0.2	3.9	13.2	56.2	5.3	11.9			
	[ʌ]	0	0	1.9	0	93.8	0.3	4.7	0.6		3	2.3	3.7	6.4	42.2	19	13.2	8.7			
	[æ]	0.6	0.6	15.2	91	1.2	0.9	3.8	0.9		2.3	4.9	25.2	66.5	14.7	11.3	3.1	1.4			
	[ɛ]	2.5	38.1	78.3	6.5	1.2	0.3	7.6	0		10.2	20.9	51.2	19.2	8.6	0.7	3.5	1.4			
	[i]	7.2	57.9	2.5	0.3	2.2	0.6	3.2	0.3		33.1	44.4	11.7	1	0.6	0.4	3.7	1.4			
	[i]	89.7	1.5	0.9	0.9	0.3	0.3	0.3	0.6		44.6	15.4	4.4	0.9	0	0.3	1.2	2.7			

B

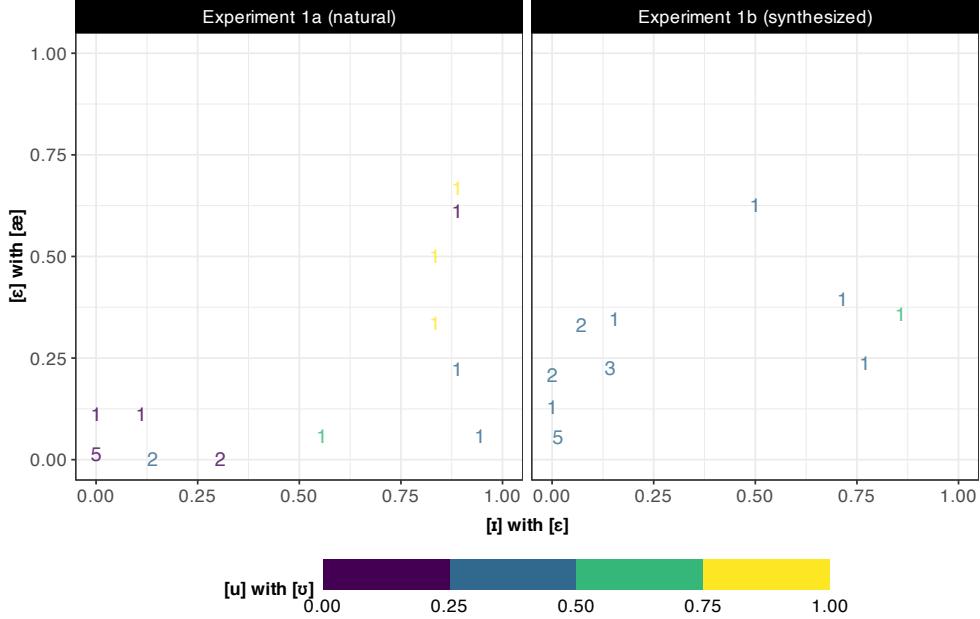


FIG. 5. Category confusability in Experiments 1a and 1b. **Panel A** summarizes the category confusability. Since correct responses were not defined for Experiment 1b, we grouped items along the x-axis based on most frequent response that listeners provided (for Experiment 1a, this was always identical to the intended response). Response percentages sum to 100 in each column, showing the response distribution depending on the most frequent response. **Panel B** summarizes individual differences across listeners, in terms of the listener-specific confusability of [i] with [ɛ] (x-axis), [ɛ] with [æ] (y-axis), and [u] with [ʊ] (color fill).

³⁸⁹ who exhibited increased [i]-to-[ɛ] confusion *also* exhibited increased [ɛ]-to-[æ] confusion. This

390 is precisely what would be expected by listeners who assume a dialect in which these vowels
 391 are articulated lower (with higher F1) than in the dialect of the talker in Experiment 1a.
 392 A similar, but less pronounced, pattern was also found with regard to [u]-to-[ʊ] confusions.⁶
 393 Finally, a qualitatively similar relation between [ɪ]-to-[ɛ], [ɛ]-to-[æ], and [ʊ]-to-[ʊ] confusions
 394 was also observed in Experiment 1b (right panel of Figure 5B), though the pattern was
 395 unsurprisingly less pronounced given that the stimuli in Experiment 1b by design often
 396 fell into the ambiguous region *between* vowels. Taken together, Experiments 1a and 1b
 397 thus suggest that systematic dialectal differences between participants may be a substantial
 398 contributor of the relatively low correct classification rate observed for experiment 1a.

399 This highlights two important points. First, the data from Experiment 1a demonstrate
 400 the perceptual challenges associated with an unfamiliar talker: in the absence of lexical or
 401 other context to distinguish between the eight available response options, listeners can only
 402 rely on the acoustic information in the input. In such a scenario, even listeners who are
 403 in principle familiar with the dialect spoken by the talker have comparatively little infor-
 404 mation to determine the talker’s dialect, making apparent what Matt Winn (2018) aptly
 405 summarizes as “speech [perception] is not as acoustic as [we] think”. Second, when dialect
 406 variability is taken into account, listeners’ recognition accuracy improved substantially. Af-
 407 ter removing 7 listeners who heard more than 50% of the [ɪ] items as [ɛ], *all* vowels were
 408 correctly recognized at least 88.3% of the time (overall accuracy = 95.9%). This suggests
 409 that dialect differences affected the recognition of all vowels. This aspect of our results serves
 410 as an important reminder that formant normalization is only expected to erase inter-talker
 411 variability associated with *physiological* differences: variation in dialect, sociolect, or other

412 non-physiologically-conditioned variation pose separate challenges to human perception, and
 413 require additional mechanisms (see discussion in [Barreda, 2021](#); [Weatherholtz and Jaeger,](#)
 414 [2016](#)). This introduces noise—variability in listeners’ responses that cannot be accounted
 415 for by normalization—to any comparison of normalization accounts, potentially reducing
 416 the power to detect differences between accounts.

417 III. COMPARISON OF NORMALIZATION ACCOUNTS

418 In order to evaluate normalization accounts against speech perception, it is necessary to
 419 map the phonetic properties of stimuli—under different hypotheses about normalization—
 420 onto listeners’ responses in Experiments 1a and 1b. Previous work has done so by directly
 421 predicting listeners’ responses from the raw or normalized phonetic properties of stimuli
 422 ([Apfelbaum and McMurray, 2015](#); [Barreda, 2021](#); [Crinnion *et al.*, 2020](#); [McMurray and](#)
 423 [Jongman, 2011](#); [Nearey, 1989](#)). For example, McMurray and Jongman used multinomial
 424 logistic regression to predict 8-way fricative categorization responses in US English (see also
 425 [Barreda, 2021](#)).

426 Here we pursued an alternative approach by committing to a core assumption common to
 427 contemporary theories of speech perception: that listeners acquire implicit knowledge about
 428 the probabilistic mapping from acoustic inputs to linguistic categories, and draw on this
 429 knowledge during speech recognition (e.g., TRACE, [McClelland and Elman, 1986](#); exem-
 430 plar theory, [Johnson, 1997](#); Bayesian accounts, [Luce and Pisoni, 1998](#); [Nearey, 1990](#); Norris
 431 and [McQueen, 2008](#); ASR-inspired models like DIANA or EARSHOT, [ten Bosch *et al.*,](#)
 432 [2015](#); [Magnuson *et al.*, 2020](#)). Using a general computational framework for adaptive speech

433 perception (ASP, [Xie *et al.*, 2023](#)) we trained Bayesian ideal observers to capture the expec-
434 tations that a ‘typical’ L1 adult listener might have about the formant-to-vowel mappings of
435 US English. We approximated these expectations using a database of L1-US English vowel
436 productions ([Xie and Jaeger, 2020](#))—transformed to reflect the different normalization ac-
437 counts. We then ask which of the different ideal observer models—corresponding to different
438 hypotheses about formant normalization—best predicts listeners’ responses in Experiments
439 1a and 1b.

440 A welcome side effect of this is that far fewer degrees of freedom (DFs) are required
441 to predict listeners’ responses. For example, using ordinary multinomial logistic regression
442 trained on our perceptual data to predict 8-way vowel categorization as a function of F1,
443 F2 and their interaction would require up to 28 DFs. This problem increases with the
444 number of cues considered. Because the model is trained on data that is independent of
445 our perceptual data, the ASP-based approach we employ instead uses only 2 DFs (i.e.,
446 parameters estimated based on our perceptual data) to mediate the mapping from stimuli
447 properties to listeners’ responses, regardless of the number of cues considered. Over the
448 next few sections, we describe how this parsimony is made possible through a commitment
449 to strong linking hypotheses motivated by theories of speech perception.

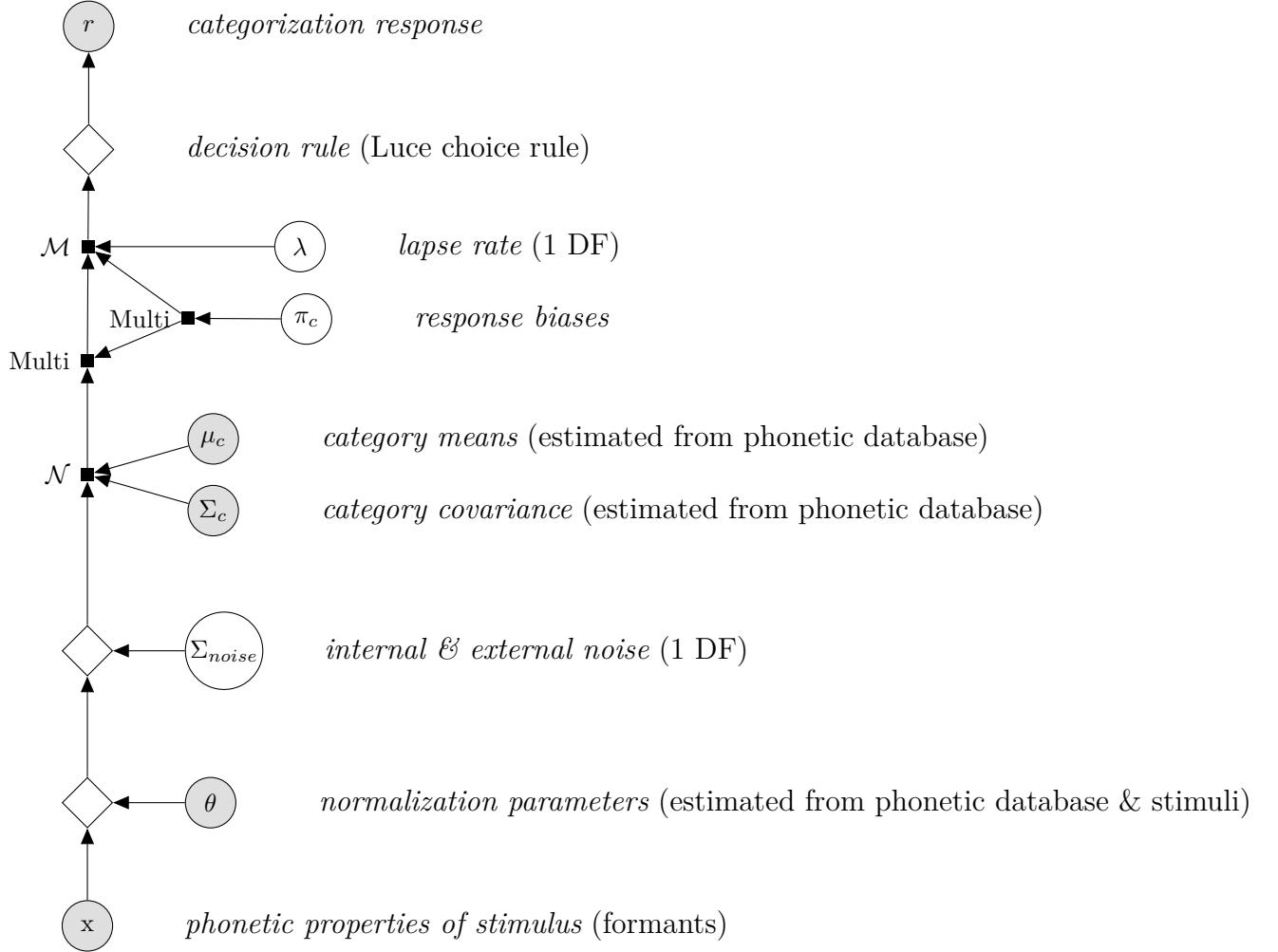


FIG. 6. Graphical model of ASP's general categorization framework (adapted for the current purpose from Xie et al., 2023, Figure 4). Here $J = 8$ (the eight vowel response options in Experiments 1a and 1b). We use this framework to compare normalization accounts against listeners' categorization responses from Experiments 1a and 1b. Filled gray circles represent variables that are known to the researcher. Empty circles represent latent variables that are not observable. Diamonds represent variable-free processes, annotated with the distributions resulting at that level of the model: \mathcal{N} (ormal), Multi (nomial), and \mathcal{M} (ixture) distributions.

450 **A. Methods**451 **1. A general-purpose categorization model for J-AFC categorization tasks**

452 Figure 6 summarises ASP’s categorization model for a J -alternative forced-choice task
 453 (for an in-depth description, we refer to [Xie et al., 2023](#)). The model combines Bayesian ideal
 454 observers (as used in e.g., [Clayards et al., 2008](#); [Feldman et al., 2009](#); [Norris and McQueen, 2008](#);
 455 [Xie et al., 2021](#); for a closely related approach, see also [Nearey and Hogan, 1986](#)) with
 456 psychometric lapsing models ([Wichmann and Hill, 2001](#)). To reduce researchers’ degrees of
 457 freedom, we adopt all assumptions made in [Xie et al. \(2023\)](#), and do not introduce additional
 458 assumptions.

459 Starting at the bottom of the figure, the phonetic input x is normalized. Here, $x =$
 460 the F1 and F2 of our stimuli (the SI, §3 E reports additional analyses that instead employ
 461 F1-F3; these analyses support the same conclusion presented here, and we mention them
 462 below where relevant). The specific computations applied to the input x depend on the
 463 normalization accounts (see Table I). We use θ to refer to the parameters required by the
 464 normalization account. For example, for the uniform scaling account ([Nearey, 1978](#)), θ is the
 465 overall mean of all log-transformed formants. For Lobanov normalization ([Lobanov, 1971](#)),
 466 θ is a vector of means and standard deviations for each formant (in Hz).

467 The normalized input is then perturbed by perceptual and environmental noise. Following
 468 [Feldman et al. \(2009\)](#), this noise is assumed to be Gaussian distributed centered around the
 469 transformed stimulus with noise variances that are independent and identical for all formants
 470 (i.e., Σ_{noise} is a diagonal matrix, and all diagonal entries have the same value). Next, the

471 likelihood of the normalized percept under each of the eight vowel categories is calculated,
 472 $p(F1, F2|vowel)$. This requires specifying listeners' expectations about the cue-to-category
 473 mapping (listeners' likelihood function). We followed [Xie et al. \(2023\)](#) and previous work and
 474 assume that each vowel maps onto a multivariate Gaussian distribution over the phonetic
 475 cues, here bivariate Gaussians over F1 and F2 (cf. [Clayards et al., 2008](#); [Feldman et al., 2009](#);
 476 [Kleinschmidt and Jaeger, 2015](#); [Norris and McQueen, 2008](#); [Xie et al., 2021](#)). The posterior
 477 probability of each vowel is obtained by combining its likelihood with its prior probability
 478 or response bias π_c , according to Bayes theorem:⁷

$$p(vowel = c|F1, F2) = \frac{\mathcal{N}(F1, F2|\mu_c, \Sigma_c + \Sigma_{noise}) \times \pi_c}{\sum_{c_i} \mathcal{N}(F1, F2|\mu_{c_i}, \Sigma_{c_i} + \Sigma_{noise}) \times \pi_{c_i}} \quad (1)$$

479 Up to this point, the model is identical to a standard Bayesian ideal observer over noisy
 480 input ([Feldman et al., 2009](#); [Kronrod et al., 2016](#)) for which the input has been transformed
 481 based on the normalization account. ASP's categorization model adds to this the potential
 482 that participants experience attentional lapses—or for other reasons do not respond based
 483 on the input—on some proportion of all trials (λ , as in standard psychometric lapsing
 484 models, [Wichmann and Hill, 2001](#)). On those trials, the posterior probability of a category
 485 is determined solely by participants' response bias, which we assume to be identical to the
 486 response bias on non-lapsing trials (following [Xie et al., 2023](#)). This results in a posterior
 487 that is described by weighted mixture of two components, describing participants' posterior
 488 on non-lapsing and lapsing trials, respectively:

$$p(vowel = v | F1, F2) = (1 - \lambda) \frac{\mathcal{N}(F1, F2 | \mu_c, \Sigma_c + \Sigma_{noise}) \times \pi_c}{\sum_{c_i} \mathcal{N}(F1, F2 | \mu_{c_i}, \Sigma_{c_i} + \Sigma_{noise}) \times \pi_{c_i}} + \lambda \frac{\pi_c}{\pi_{c_i}} \quad (2)$$

489 Finally, a decision rule is applied to the posterior to determine the response of the model,
 490 conditional on the input (one of the eight vowels in Experiments 1a and 1b). We followed
 491 the gross of research on speech perception and assume Luce's choice rule (Luce, 1959; for
 492 discussion, see Massaro and Friedman, 1990). Under this choice rule, the model can be
 493 seen as sampling from the posterior, responding with each category proportional to that
 494 category's posterior probability.

495 Next, we describe how we estimated the θ s, μ_c s and Σ_c s for each normalization account
 496 from a phonetic database. We use this database as a—very coarse-grained—approximation
 497 of a the speech input a ‘typical’ listener might have experienced previously. By fixing θ , μ_c
 498 and Σ_c based on the distribution of phonetic cues in the database, we substantially reduce
 499 the DFs that are allowed to mediate the mapping from stimulus properties to listeners’
 500 responses (following Xie *et al.*, 2023). In addition, this approach naturally penalizes overly
 501 complex models by validating these against out-of-sample data. Finally, we describe how
 502 we fit the remaining parameters as DFs to participants’ responses from Experiments 1a and
 503 1b.

504 2. *Modeling listeners' prior experience (and guarding against overfitting): θ , μ_c ,*
 505 *and Σ_c*

506 By fixing θ , μ_c , and Σ_c based on a database of vowel *productions*, we impose strong
 507 constraints on the functional flexibility of the model in predicting listeners' responses. This
 508 benefit is made possible by committing to a strong linking hypothesis—that listeners' cate-
 509 gories are learned from, and reflect, the distributional mapping from formants to vowels in
 510 previously experienced speech input (e.g., [Abramson and Lisker, 1973](#); [Massaro and Fried-](#)
 511 [man, 1990](#); [Nearey and Hogan, 1986](#)). The database we use to approximate listeners' prior
 512 experience was originally developed to compare the production of L1 and L2 speakers ([Xie](#)
 513 [and Jaeger, 2020](#)). It contains 9-10 recordings of the 8 *hVd* words from each of 17 (5 fe-
 514 male) L1 talkers of a Northeastern dialect of US English (ages 18 to 35 years old). Since
 515 Experiments 1a and 1b used recordings of one of these talkers, we excluded that talker prior
 516 to fitting training ideal observers on the data. In total, this yields 5842 recordings that are
 517 annotated for F0, F1-F3, and vowel duration. The SI ([§3 A 1](#)) summarizes the distribution
 518 of these cues, and how the different normalization accounts affect those distributions.

519 To avoid over-fitting the ASP model to the database, we used 5-fold cross-validation:
 520 we randomly split the [Xie and Jaeger \(2020\)](#) database into five approximately evenly-sized
 521 folds (following [Persson and Jaeger, 2023](#)). This split was performed within each vowel to
 522 guarantee that all five folds had the same relative amount of data for each vowel category.
 523 These splits were combined into five training sets, each containing one of the folds (20% of

524 the data). This way, each training set was different from the others, increasing the variability
 525 between sets.⁸

526 For each training set and for each normalization account, we then estimated the required
 527 normalization parameters θ for all talkers, and normalized all formants based on those talker-
 528 specific parameters. This yielded 5 (training sets) * 20 (accounts) = 100 normalized training
 529 sets. For each of these normalized training sets, we fit the category means, μ_c , and covariance
 530 matrices, Σ_c , of all eight vowels, using the R package `MVBeliefUpdatr` (Jaeger, 2024).⁹

531 This yielded 100 ideal observer models, 5 for each of the 20 normalization accounts in
 532 Table I. Of note, the 20 ideal observers fit on each fold differ *only* in the assumptions
 533 they make about the normalization that is applied to cues before they are mapped onto
 534 the eight vowel categories. Figure 7 visualizes the resulting bivariate Gaussian categories
 535 for four of the 20 normalization accounts. This illustrates one advantage of the cross-
 536 validation approach: it takes a modest step towards simulating differences across listeners'
 537 prior experience (represented by the five different folds).

538 **3. Transforming the stimuli from Experiments 1a and 1b into the normalized
 539 phonetic spaces**

540 Next, we transformed the stimuli of Experiments 1a and 1b into the formant space defined
 541 by the 20 normalization accounts in Table I. This requires estimating the required normal-
 542 ization parameters θ for each experiment and normalization account. We calculated these θ s
 543 over all stimuli (of each experiment and normalization account). For example, for the uni-
 544 form scaling account (Nearey, 1978), we calculated the overall mean of all log-transformed

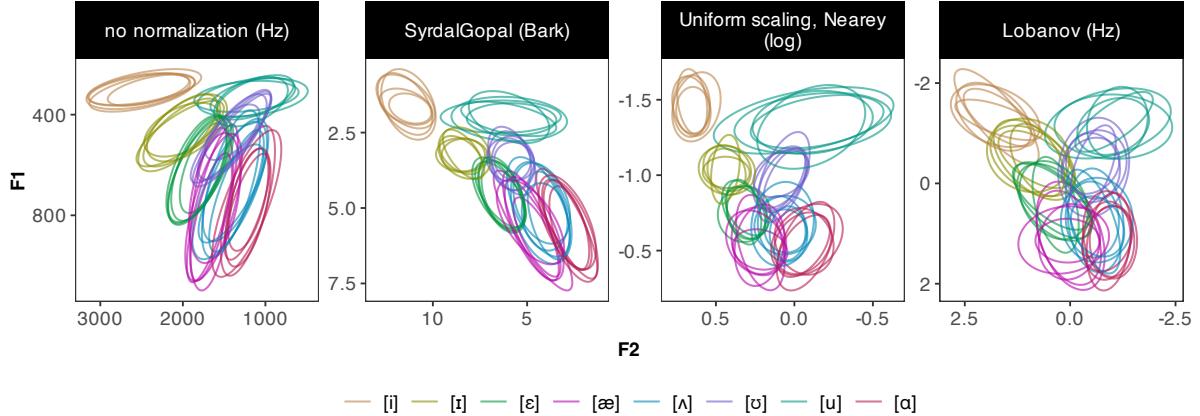


FIG. 7. Visualizing the bivariate Gaussian categories (prior to adding Σ_{noise}) of four example normalization accounts in F1-F2 space. Separate ellipses are shown for each of the five training sets (each set corresponds to one set of eight ellipses). The relative stability of the category ellipses across training sets indicates that the database is sufficiently large for the present purpose.

545 formants over all stimuli. For Lobanov normalization (Lobanov, 1971), we calculated the
 546 mean and standard deviation of each formant (in Hz) over all stimuli. For each combina-
 547 tion of experiment and normalization account, we then normalized the stimuli using those
 548 parameter estimates.

549 Combining the 100 normalized training sets described in the previous section with the
 550 matching normalized stimuli from each of the two experiments yielded 200 data sets.

551 **4. Noise (Σ_{noise}) and attentional lapses (λ)**

552 Finally, we describe the two parameters of the ASP model that we fit against listeners'
 553 responses in Experiments 1a and 1b. These two parameters constitute the only DFs that
 554 mediate the link from ideal observers' predictions to listeners' responses, and which are
 555 specifically tuned to these. The first DF (Σ_{noise}) models the effects of internal (perceptual)
 556 and external (environmental) noise on listeners' perception. While previous work provides

557 estimates of the internal noise in formant perception, these estimates were obtained under
 558 *assumptions* about the relevant formant space. For example, [Feldman *et al.* \(2009\)](#) estimated
 559 the internal noise variance to be about 15% of the average category variance along F1 and F2.
 560 This estimate was based on the assumption that human speech perception transforms vowel
 561 formants into Mel, without further normalization. Since we aim to *test* which normalization
 562 account best explains speech perception, we cannot rely on this or other internal noise
 563 estimates obtained under a single specific assumption. Additionally, internal noise can vary
 564 across individuals and external noise can vary across environments (a point particularly
 565 noteworthy, given that we conducted Experiments 1a and 1b over the web). We thus allowed
 566 the noise variance Σ_{noise} to vary in fitting participants' responses. Following [Feldman *et al.*](#)
 567 ([2009](#)), we assumed that perceptual noise had identical effects on all formants in the phonetic
 568 space defined by the normalization account (see also [Kronrod *et al.*, 2016](#)). This reduces
 569 Σ_{noise} to a single DF, regardless of the normalization account (for details, see SI §3 A 3).

570 The magnitude of Σ_{noise} affects the slope of the categorization functions that predict
 571 listeners' responses from stimulus properties (here, F1 and F2): higher Σ_{noise} imply more
 572 shallow categorization slopes. To facilitate comparison of Σ_{noise} values across normaliza-
 573 tion accounts, we report results in terms of the best-fitting *noise ratios* (τ^{-1}), rather than
 574 Σ_{noise} s. Specifically, Σ_{noise} is best understood *relative* to the inherent variability of the
 575 vowel categories (Σ_c). This variability in turn depends on the phonetic space defined by
 576 the normalization account. We thus divide Σ_{noise} by the mean of the diagonals of all Σ_c s to
 577 obtain the *noise ratio* τ^{-1} . For example, noise ratio of 0 corresponds to the absence of any
 578 noise, and a noise ratio of 1 corresponds to noise variance of the same magnitude as the av-

verage category variance along F1 and F2 in the phonetic space defined by the normalization account.¹⁰ Figure 8B illustrates the effects of this noise ratio for Nearey’s uniform scaling account.

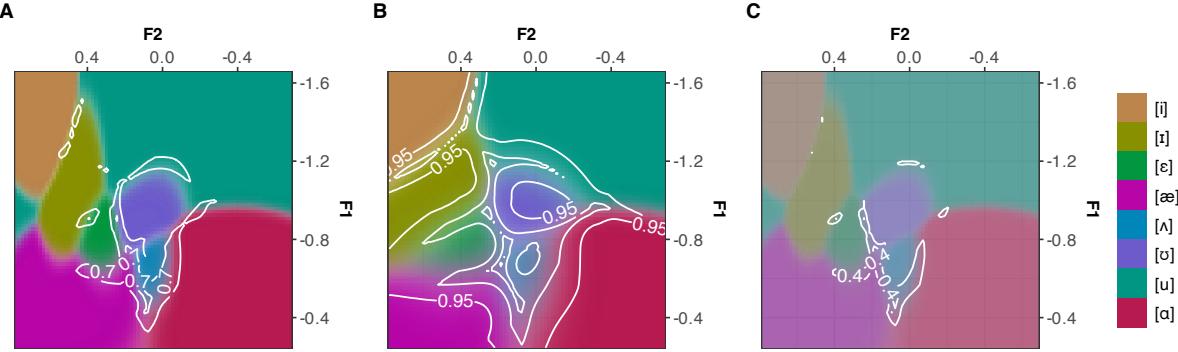


FIG. 8. Illustrating the consequences of perceptual and external noise (Σ_{noise}) and attentional lapse rates (λ) on the predicted posterior distribution of vowel categorizations. Shown are the average predicted posteriors across all five folds for Nearey’s uniform scaling account. **Panel A:** Predicted posterior distribution for noise ratio $\tau^{-1} = \lambda = 0$. **Panel B:** Same for $\tau^{-1} = 1$ and $\lambda = 0$. **Panel C:** Same for $\tau^{-1} = 0$ and $\lambda = 0.5$. Transparency of a color is determined by that vowel’s posterior probability. Contours indicate the highest posterior probability of any vowel (at .4, .5, .7, .95 probability level).

Second, participants can attentionally lapse or for other reasons reply without considering the speech input. We thus allowed lapse rates (λ) to vary while fitting human responses. This introduces a second DF, which we fit against listeners’ responses. Together, the inclusion of freely varying lapse rates and a uniform response bias allows the ASP models to capture that some unknown proportion of listeners’ responses might be more or less random, rather than reflecting properties of the vowel stimuli. This is illustrated in Figure 8C.

Finally, participants can have response biases that reflect their beliefs about the prior probability of each category. However, to reduce the DFs fit to participants’ responses, we

590 did *not* fit this response bias against listeners' responses (thus avoiding $J - 1 = 7$ additional
 591 DFs). Instead, we assumed uniform response biases—i.e., that listeners believed all eight
 592 response options in the experiments to be equally likely ($\forall c \pi_c = .125$). This decision implies
 593 that our models would not be able to capture any potential non-uniformity in listeners'
 594 response biases—including potential effects of additional acoustic differences (the absence
 595 of [h] in *odd* or the coda [t], rather than [d] in *hut*) and orthographically particular response
 596 options in Experiment 1a (“who’d”, “odd”, and “hut”). We do, however, see no reasons to
 597 expect this decision to bias the comparison of normalization accounts.

598 5. *Fitting normalization accounts to listeners' responses*

599 For each of the 200 combinations of experiment, normalization account, training set,
 600 we used constrained quasi-Newton optimization (Byrd *et al.*, 1995, as implemented in R's
 601 `optim()` function) to find the λ and τ^{-1} values that best described listener's responses.
 602 Specifically, we used the 100 ideal observers described in the previous sections, applied them
 603 to the normalized stimuli of the experiment, and determined which λ and τ^{-1} maximized
 604 the likelihood of listener's responses (for details, see SI [§3 A 3](#)). This procedure yielded five
 605 maximum likelihood estimates for both λ and τ^{-1} for each combination of experiment and
 606 normalization account—one for each training set. All result presented below were validated
 607 and confirmed by grid searches over the parameter spaces (SI, [§3 F](#)).

608 We compare normalization accounts in terms of the likelihood of listeners' responses
 609 under these maximum likelihood estimates of λ and τ^{-1} . Comparing accounts in terms
 610 of their data likelihood, rather than the accuracy of predicting intended productions (e.g.,

611 Johnson, 2020; Persson and Jaeger, 2023), or correlations with human response proportions
 612 (e.g., Hillenbrand and Nearey, 1999; Nearey and Assmann, 1986), follows more recent work
 613 (e.g., Barreda, 2021; McMurray and Jongman, 2011; Richter *et al.*, 2017; Xie *et al.*, 2023)
 614 and parallels standard approaches to model comparison in contemporary data analysis. We
 615 note that this approach puts normalization accounts to a stronger test. For example, a
 616 model can exhibit high correlations with listeners' responses even when its predictions are
 617 systematically 'off'. Similarly, a model can achieve high accuracy in predicting listeners'
 618 responses simply because it always predicts the most frequent response, and that response
 619 accounts for sufficiently much of the data. In contrast, the likelihood of listeners' responses
 620 under a model is a direct measure of how well the model captures the distribution of listeners'
 621 responses conditional on the stimulus properties. In particular, data likelihood will be
 622 maximized if, and only if, the model-predicted posterior probabilities of each vowel for each
 623 stimulus are identical to the proportion with which those vowels occur in listeners' responses.

624 **B. Results**

625 We begin by comparing the fit of different accounts against listeners' responses in Ex-
 626 periments 1a and 1b. Given the comparatively large number of accounts compared here, we
 627 provide initial conclusions based on the best-fitting accounts along with the description of
 628 the results (more in-depth discussion is provided in the general discussion). Following this
 629 comparison, we visualize how different normalization accounts predict the formant space to
 630 be divided into the eight vowel categories.

631 1. Comparing normalization accounts in terms of fit against human behavior

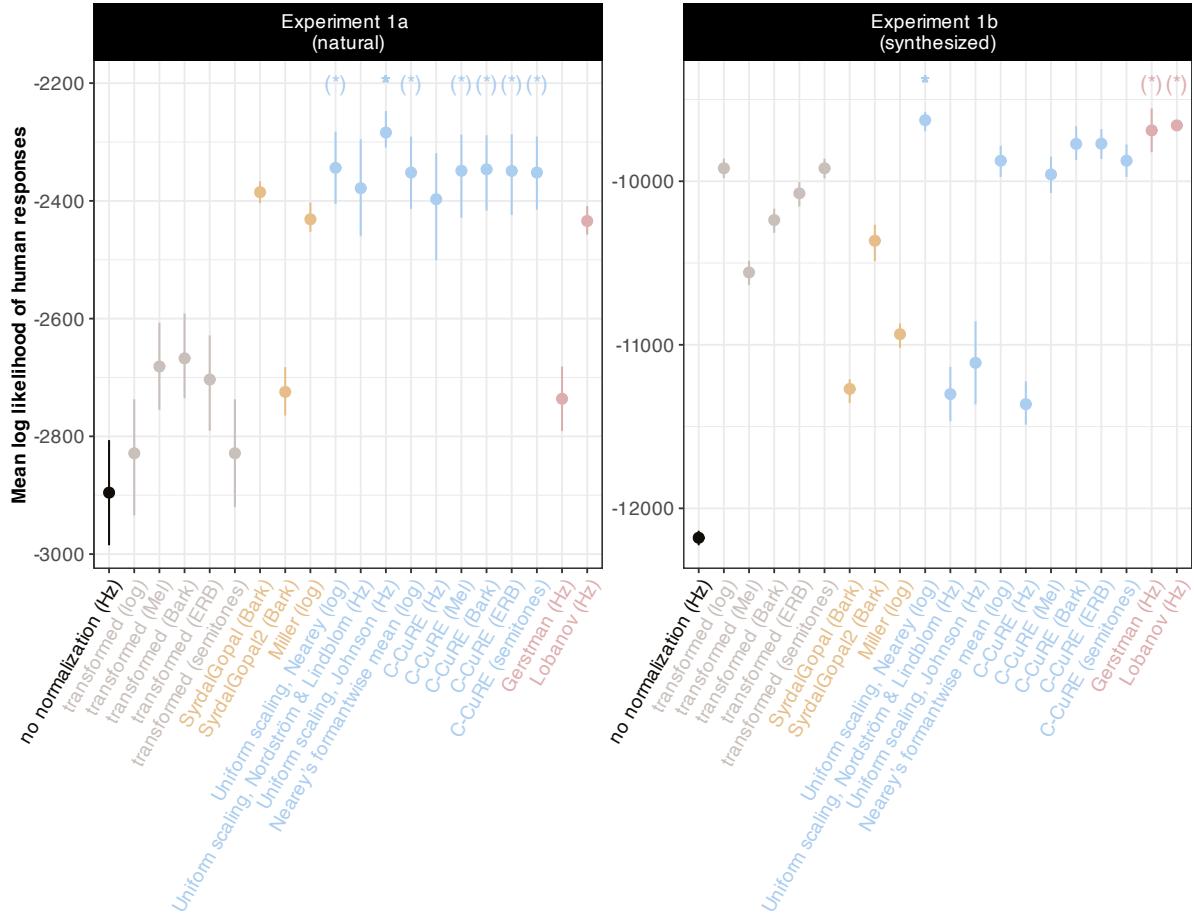


FIG. 9. Comparison of normalization accounts against listeners' responses. Pointranges indicate mean and 95% bootstrapped CIs of the log-likelihood summarized over the five training sets (higher is better). Accounts that fit listeners' responses to an extent that is statistically indistinguishable from the best-fitting account are marked by (*). Note that y-axis range differs across panels, and that it is *not* meaningful to compare the absolute log-likelihood values across the two experiments (just as it is not meaningful to compare the data likelihood of regressions that are fit on two different data sets).

632 Figure 9 compares how well the different normalization accounts fit listeners' responses
 633 in Experiments 1a and 1b. All accounts performed well above chance guessing (chance log

⁶³⁴ likelihood in Experiment 1a: -5334; Experiment 1b: -13213) but also well below the highest
⁶³⁵ possible performance (in Experiment 1a, log-likelihood = -1348, in Experiment 1b: -7225).

⁶³⁶ Normalization significantly improved the fit to listeners' responses relative to no normal-
⁶³⁷ ization. This was confirmed by paired one-sided *t*-tests comparing the maximum likelihood
⁶³⁸ values for each normalization account against those in the absence of normalization (all *ps*
⁶³⁹ < .05; see SI [§3 B 1](#)). Not all normalization accounts achieved equally good fits, however:
⁶⁴⁰ only some extrinsic accounts fit listeners' behavior well across both experiments. This sup-
⁶⁴¹ ports two conclusions. First, it suggests that the normalization mechanisms operating during
⁶⁴² human speech perception involve computations that go beyond static transformations into
⁶⁴³ psycho-acoustic spaces. Second, it suggests that the input to these computations is not
⁶⁴⁴ limited to intrinsic information—i.e., that the computations draw on information beyond
⁶⁴⁵ what is available in the acoustic signal *at that moment*. In particular, extrinsic normaliza-
⁶⁴⁶ tion requires the estimation and memory maintenance of talker-specific properties from the
⁶⁴⁷ speech signal.

⁶⁴⁸ While the accounts that achieved the best fit against listeners' responses differed be-
⁶⁴⁹ tween experiments, both were variants of uniform scaling. For Experiment 1a, Johnson
⁶⁵⁰ normalization account provided the best fit (log likelihood = -2284, SD = 41 across the
⁶⁵¹ five crossvalidation folds), while Nearey's uniform scaling account provided the best fit to
⁶⁵² Experiment 1b (log likelihood = -9626, SD = 78). Both accounts essentially slide the repre-
⁶⁵³ sentational 'template' of a dialect—here the eight bivariate Gaussian categories of an ideal
⁶⁵⁴ observer—along a single line in the formant space. They differ only in *which* space this linear
⁶⁵⁵ relation between formants is assumed. The same two accounts still fit listeners' responses

best when F3 was included in the analysis in addition to F1 and F2 (SI, §3E).¹¹ This suggests that formant normalization might involve comparatively parsimonious maintenance of talker-specific properties: in its simplest form, uniform scaling employs a single formant statistic to normalize all formants. In contrast, computationally more complex accounts like Lobanov normalization might require the estimation and maintenance of two formant statistics (mean and standard deviation) for each formant that is normalized (e.g., a total of four formant statistics for F1 and F2, or six statistics for F1-F3).

For both experiments, there were several accounts that fit listeners' responses similarly well as the best-fitting accounts ($ps > .065$). All of these were extrinsic accounts, though the specific accounts differed between experiments. Notably, only Nearey's uniform scaling either provided the best fit (Experiment 1b) or achieved performance statistically indistinguishable from the best fit *for both experiments* (for Experiment 1a: $p > 0.08$, log likelihood = -2344, $SD = 84$). Beyond the performance of Nearey's uniform scaling, there was little evidence of a correlation in relative ordering of accounts between experiments (Spearman rank $r = 0.09$, $p = 0.72$). Some accounts fit listeners' responses well for Experiment 1a, but not for Experiment 1b, and vice versa. Of note is the particularly variable performance of the centering accounts operating in Hertz space, i.e., C-CuRE Hz, Nordström & Lindblom and Johnson normalization. Similar variability across the two experiments is also observed for the two standardizing accounts, both of which operate in Hz space.

That an account operating over log-transformed formants—Nearey's uniform scaling—fits human behavior better should not be surprising. While questions remain about the exact organization of auditory formant representations, it is uncontroversial that the perceptual

678 sensitivity to acoustic frequency information is better approximated by a logarithmic scale
679 than by a linear scale (see [Moore, 2012](#)). As a result, a 30 Hz difference in an F1 of 300
680 Hz (a 10% change) is expected to be perceptually more salient than a 30 Hz change in an
681 F2 of 2500 Hz (a 1.2% change). In line with this reasoning, additional tests not reported
682 here found that Johnson normalization would provide the best fit to *both* experiments if
683 it was applied to log-transformed formants (instead of Hertz). In summary, variability in
684 how well different accounts predict human behavior across the two experiments highlights
685 the importance of psycho-acoustic transformations for human speech perception. This also
686 highlights the importance of comparing normalization accounts against multiple types of
687 data.

688 **2. *Visualizing the consequences of different normalization mechanisms***

689 Before we turn to the general discussion, we briefly visualize how different normalization
690 mechanisms affect vowel categorization. This sheds light on *why* the accounts differ in
691 how well they fit listeners' responses. Figure 10 visualizes the categorization functions
692 predicted by four different normalization accounts, using the best-fitting λ and τ^{-1} values
693 for each account (i.e., the values that lead to the fit shown in Figure 9). Figure 10 highlights
694 three points. First, a comparison across rows of Figure 10 shows how much the choice
695 of normalization can affect how the acoustic space gets carved up into vowel categories: a
696 comparison of the first (no normalization), third (Johnson), and fourth row (Lobanov) shows
697 that even normalization accounts operating over the same space can yield very different
698 categorization behavior.

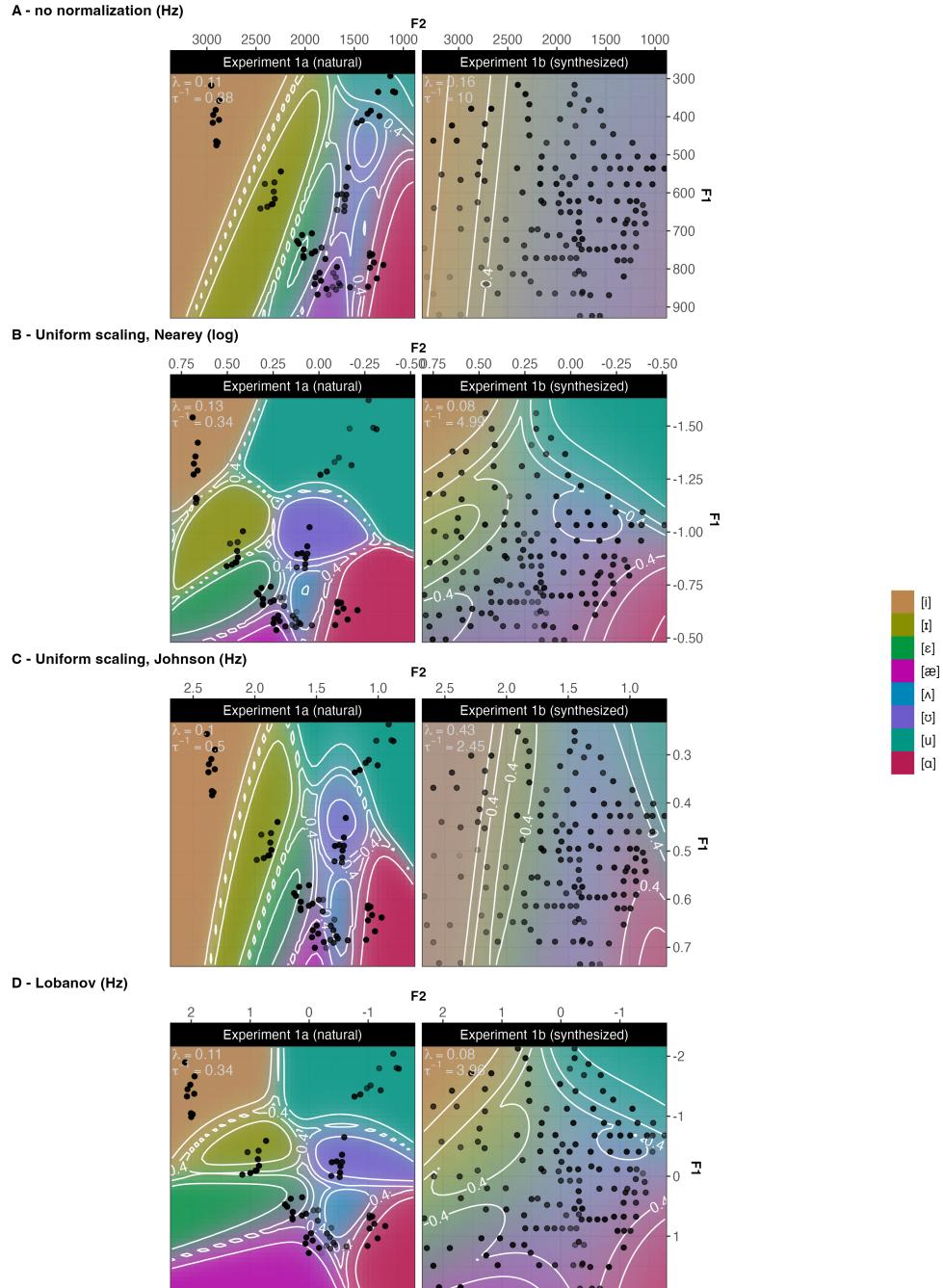


FIG. 10. Predicted categorization functions over the F1-F2 space under four different normalization accounts. For each account, we show the predicted posterior probabilities of all eight vowels obtained by averaging over the maximum likelihood parameterizations (of λ and τ^{-1}) for the five training sets (shown at top of each panel). **Top:** absence of normalization shown for reference. **2nd row:** the best-fitting account for Experiment 1a. **3rd row:** the best-fitting account for Experiment 1b (and second best for Experiment 1a). **Bottom:** the second best-fitting account in Experiment 1b. Contours indicate the highest posterior probability of any vowel. Points indicate location of test stimuli. The increasing opacity of points indicates a worse fit by the account (see text for detail).

699 Second, the best-fitting parameters (shown at the top of each panel) were relatively com-
 700 parable across accounts but differed more substantially across experiments. Specifically,
 701 the best-fitting estimates of lapse rates λ were generally comparable across the two exper-
 702 iments (with the exception of Nordström & Lindblom and Johnson normalization, which
 703 exhibited substantially higher lapse rates in Experiment 1b; SI [§3B2](#)). This suggests that
 704 participants in both experiments were about equally likely to pay attention to the stimulus.
 705 The best-fitting noise ratios τ^{-1} , however, differed substantially across experiments, and
 706 were 9 times larger for Experiment 1b (mean $\tau^{-1} = 4.74$, SD = 2.57 across normalization
 707 accounts) than for Experiment 1a (mean $\tau^{-1} = 0.52$, SD = 0.49). This difference most likely
 708 reflects the fact that the synthesized stimuli in Experiment 1b left listeners with substan-
 709 tially more uncertainty about the intended category, as discussed during the description of
 710 the experiments.

711 Since noise is assumed to be independent of category variability (see also [Feldman *et al.*,](#)
 712 [2009](#); [Kronrod *et al.*, 2016](#)), differences in noise ratios can substantially change the catego-
 713 rization function. This is particularly evident for the accounts that had more variable per-
 714 formance across the two experiments. For example, both Johnson (third row) and Lobanov
 715 normalization (fourth row) resulted in very different best-fitting categorization functions for
 716 Experiments 1a and 1b.

717 Third and finally, Figure 10 also shows how well accounts fit listeners' responses for each
 718 test stimulus (opaqueness of the black points). This begins to explain *why* some accounts
 719 fit listeners' responses in Experiment 1b less well. For example, the Johnson normalization
 720 account (third row) predicts the responses to the test stimuli in Experiment 1a well, but

721 fails to predict the responses to the test stimuli in Experiment 1b. This drop in performance
 722 seems to be primarily driven by stimuli that are unlikely to be articulated by the same talker
 723 (lower left, cf. dashed line in Figure 4). This might suggest that this account was over-
 724 engineered to explain naturally occurring productions—the type of data, it was originally
 725 tested on (Johnson, 2020). A plausible account of normalization, however, should be able
 726 to explain human perception to any type of stimulus, including synthesized stimuli. The SI
 727 (§3 B 3) presents more detailed by-item comparisons of normalization accounts that might
 728 be of interest to some readers.

729 **IV. GENERAL DISCUSSION**

730 Research on vowel normalization has an influential history. Cognitive scientists have
 731 long aimed to understand the organization of frequency information in the human brain
 732 (Siegel, 1965; Stevens and Volkmann, 1940), and how it helps listeners overcome cross-talker
 733 variability in the formant-to-vowel mapping (e.g., Fant, 1975; Joos, 1948; Nordström and
 734 Lindblom, 1975). Auditory processes that normalize speech inputs for differences in vocal
 735 tract physiology are now recognized to be an integral part of speech perception (Johnson
 736 and Sjerps, 2021; McMurray and Jongman, 2011; Xie *et al.*, 2023). Here, we set out to
 737 investigate what types of computations are implicated in the normalization of the frequency
 738 information that plays a critical role in the recognition of vowels.

739 Our results support three theoretical insights. First, human speech perception draws on
 740 more than psycho-acoustic transformations or intrinsic information, in line with previous
 741 research on normalization (Adank *et al.*, 2004; Ladefoged and Broadbent, 1957; Nearey,

742 1989). Rather, formant normalization seems to involve the estimation and storing of talker-
743 specific formant properties. Second, computationally simple uniform scaling accounts pro-
744 vide the best fit to listeners' responses, suggesting comparatively parsimonious maintenance
745 of talker-specific properties. This replicates and extends previous findings that uniform scal-
746 ing or similarly simple corrections for vocal tract size provide a better explanation for human
747 perception than more complex extrinsic accounts (Barreda, 2021; Richter *et al.*, 2017). It is
748 impossible to rule out more complex approaches to perceptual normalization given the large
749 number of possible alternatives. However, given that uniform scaling provides a parsimo-
750 nious explanation for human formant normalization, and the current absence of empirical
751 evidence for more complex computations, we submit that researchers ought to adapt uni-
752 form scaling as our working hypothesis. Third, the psycho-acoustic representation assumed
753 by different normalization accounts matter, as indicated by the comparison of otherwise
754 computationally similar accounts (e.g. Nearey's vs. Johnson's uniform scaling).

755 These results contribute to a still comparatively small body of work that has evaluated
756 competing normalization accounts against listeners' perception, whereas most previous work
757 evaluates accounts against intended productions. Complementing previous work, we took
758 a broad-coverage approach: the present study compared 20 of the most influential nor-
759 malization accounts against listeners' perception of *hVd* words with all eight US English
760 monophthongs in both natural and synthesized speech. This contrasts with previous work,
761 which has typically focused on subsets of the vowel system, either using natural *or* synthe-
762 sized speech, and considering a much smaller subset of accounts (typically 2-3 at a time). By
763 considering a wider range of accounts, a wider range of formant values and vowel categories,

764 and multiple types of speech, we aimed to contribute to a more comprehensive evaluation
 765 of competing accounts.

766 Next, we discuss the theoretical consequences of these findings for research beyond for-
 767 mant normalization. Following that, we discuss limitations of the present work, and how
 768 future research might overcome them.

769 **A. Consequences for theories of speech perception and beyond**

770 Understanding the perceptual space in which the human brain represents vowel categories—
 771 i.e., the normalized formant space—has obvious consequences for research on speech percep-
 772 tion. To illustrate how far reaching these consequences can be, we discuss a few examples.
 773 For instance, research on *categorical perception* has found that vowels seem to be per-
 774 ceived less categorically than some types of consonants. Recent work has offered an elegant
 775 explanation for this finding: the perception of formants—relevant to the recognition of
 776 vowels—might be more noisy than the perception of the acoustic cues that are critical to
 777 the recognition of less categorically perceived consonants (Kronrod *et al.*, 2016). This is
 778 a parsimonious explanation, potentially preempting the need for separate explanations for
 779 the perception of different types of phonemic contrasts. Kronrod and colleagues based their
 780 argument on estimates they obtained for the relative ratio of meaningful category variability
 781 to perceptual noise (τ , the inverse of our noise ratios, τ^{-1}). Critically, this ratio depends
 782 both on (i) the perceptual space in which formants are assumed to be represented (Kronrod
 783 et al used Mel-transformed formants), and on (ii) whether the meaningful category variabil-
 784 ity is calculated prior to, or following, normalization (Kronrod et al assumed the former,

785 which increases estimates of category variability). Our point here is not to cast doubt on
786 the results of [Kronrod *et al.* \(2016\)](#) —the fact that the best-fitting noise ratios in our study
787 were relatively similar across accounts (while varying across experiments) suggests that the
788 result of Kronrod and colleagues are likely to hold even under different assumptions about
789 (i) and (ii)—but rather to highlight how research on the perception and recognition of
790 vowels depends on assumptions about formant normalization. For example, similar points
791 could be raised about experiments on statistical learning that manipulate formant or other
792 frequency statistics (e.g., [Chládková *et al.*, 2017](#); [Colby *et al.*, 2018](#); [Wade *et al.*, 2007](#); [Xie
793 *et al.*, 2021](#)). Such experiments, too, need to make assumptions about the space in which
794 formants are represented. If these assumptions are incorrect, this can affect whether the
795 experimental manipulations have the intended effects, increasing the chance of null effects
796 or misinterpretation of observed effects.

797 Understanding the perceptual space in which the human brain represents vowel cate-
798 gories also has consequences for research beyond speech perception, perhaps more so than is
799 sometimes recognized. For instance, in sociolinguistics and related fields, Lobanov remains
800 the norm for representing vowels due to its efficiency in removing cross-talker variability
801 (for review, see [Adank *et al.*, 2004](#); [Barreda, 2021](#)). However, as shown in the present study,
802 removing cross-talker variability is not the same as representing vowels in the perceptual
803 space that listeners actually employ. Here, we do *not* find Lobanov to describe human
804 perception particularly well. On the contrary, we find no support for the hypothesis that
805 human speech perception employs these more complex computations that have been found
806 to perform best at reducing category variability. This should worry sociolinguists. In order

807 to understand how listeners infer a talker's background or social identity, it is important
808 to understand the perceptual space in which inferences are actually rooted. Critically, the
809 representations resulting from formant normalization presumably form an important part of
810 the information that listeners use to draw social and linguistic inferences. It should thus be
811 obvious that the use of normalization accounts that do not actually correspond to human
812 perception can both mask real markers of social identity, and hallucinate markers that are
813 not actually present. For example, in order to determine how a talker's social identity influ-
814 ences their vowel realizations, it is important to discount *all and only* effects that listeners'
815 will attribute to physiology, rather than social identity (Disner, 1980; Hindle, 1978).

816 Similar concerns apply to dialectology, research on language change, second language
817 acquisition research, etc. For example, the perceptual space in which vowels are represented
818 is critical to well-formed tests of hypotheses about the factors shaping the organization of
819 vowel inventories across languages of the world (Lindblom, 1986; Stevens, 1972, 1989). It is
820 essential in testing hypotheses about the extent to which the cross-linguistic realization of
821 those systems is affected by perceptual processes (Flemming, 2010; Steriade, 2001), or by
822 preferences for communicatively efficient linguistic systems (e.g., Hall *et al.*, 2018; Lindblom,
823 1990; Moulin-Frier *et al.*, 2015). Similarly, tests of the hypothesis that vowel *articulation*
824 during natural interactions is shaped by communicative efficiency do in obvious ways depend
825 on assumptions about the perceptual space in which talkers—by hypothesis—aim to reduce
826 perceptual confusion (cf. Buz and Jaeger, 2016; Gahl *et al.*, 2012; Scarborough, 2010; Wedel
827 *et al.*, 2018). The same applies to any other line of research that aims to understand the
828 perceptual consequences of formant variation across talkers, including research on infant- or

829 child-directed speech (Eaves Jr *et al.*, 2016; Kuhl *et al.*, 1997), and research on whether non-
 830 native talkers are inherently more variable than native talkers (Smith *et al.*, 2019; Vaughn
 831 *et al.*, 2019; Xie and Jaeger, 2020). In short, the perceptual space in which vowels are
 832 represented is a critical component of understanding the structure of vowel systems, the
 833 factors that shape them, and the ways in which they are used in natural language.

834 **B. Limitations and future directions**

835 The present work shares a few limitations with previous work. Here we focus on limita-
 836 tions that follow from the assumptions we made in our computational framework. While
 837 theories and hypotheses often contain substantial vagueness, *quantitative tests* of those
 838 theories—as we have done here—require assumptions about *every* aspect of the model.
 839 Here, this included all the steps necessary to link properties of the stimuli to listeners’ re-
 840 sponses. For this purpose, we adopted the ASP framework (Xie *et al.*, 2023), and visualized
 841 the graphical model that links stimuli (x) to responses (r) in Figure 6.

842 Many of the assumptions we made should be quite uncontroversial—e.g., the decision to
 843 include both external (environmental) and internal (perceptual) noise in our model. While
 844 these noise sources are often ignored in modeling human behavior, it is uncontroversial that
 845 they exist. Other assumptions we made were introduced as simplifying assumptions for
 846 the sake of feasibility—e.g., we expressed the effect of both types of noise through a single
 847 parameter that related the average within-category variability of formants to noise variability
 848 in the transformed and normalized formant space. In reality, however, environment noise
 849 can have effects that are independent of internal noise, and internal noise likely affects

850 information processing at multiple (or all) of the steps shown in Figure 6. Such simplifying
851 assumptions are both inevitable, and not necessarily problematic: as long as they do not
852 introduce systematic bias to the evaluation of normalization accounts, they should not limit
853 the generalizability of our results.

854 Some of our assumptions, however, might be more controversial. For example, we as-
855 sumed that category representations can be expressed as multivariate Gaussian distributions
856 in the formant space. This assumption, too, is a simplifying assumption—it simplified the
857 computation of likelihoods—rather than a critical feature of the ASP framework we em-
858 ployed. While human category representations are unlikely to be Gaussians, the alternative,
859 e.g., exemplar representations, would come with its own downsides, such as increased sen-
860 sitivity to the limited size of phonetic databases and substantial increases in computation
861 time (exemplar representations afford researchers with much larger degrees of freedom). For
862 researchers curious how this and other assumptions we made affect our results, our data
863 and code are shared on OSF. This includes the R markdown document that generates this
864 PDF, making it comparatively easy to revisit any of our assumptions to then regenerate the
865 entire study with a click of a button in RStudio.

866 Like previous work, we further assumed that all listeners in our experiments use the
867 same underlying vowel representations—the same dialect template(s). However, as already
868 discussed, it is rather likely that not all of our listeners employed the same dialect tem-
869 plate(s). An additional analysis reported in the SI (§3D) thus compared normalization
870 accounts against only the subset of listeners who employed the dialect template used by
871 the majority of participants (see lower-left of Figure 5B). This left only 11 participants for

872 Experiment 1a (61.1%) and 14 for Experiment 1b (77.8%), substantially reducing statistical
 873 power. Replicating the main analysis, uniform scaling accounts again fit listeners' behavior
 874 well across both experiments. The best-performing accounts did, however, differ from the
 875 ones obtained for the superset of data (see SI, §3D).

876 A related assumption was introduced by the use of a phonetic database to approximate
 877 listeners' vowel representations. This deviates from most previous evaluations of normal-
 878 ization accounts (McMurray and Jongman, 2011; Barreda, 2021; but see Richter *et al.*,
 879 2017), and reflects our commitment to a strong assumption made by most theories of speech
 880 perception: that listeners' representations reflect the formant statistics previously experi-
 881 enced speech input. By using a phonetic database to estimate listeners' representations, we
 882 *substantially* reduced the degrees of freedom in the evaluation of normalization accounts,
 883 reducing the chance of over-fitting to the data from our experiments. Our approach does,
 884 however, also introduce two new assumptions.

885 First, our approach assumes that the mixture of dialect template(s) used by talkers in the
 886 database sufficiently closely approximates those of the listeners in our experiments. Some
 887 validation for this assumption comes from the additional analysis reported in the preceding
 888 paragraph: when we subset listeners to only those who used the majority dialect template,
 889 this improved the fit of all normalization accounts—as expected, if the category representa-
 890 tions we trained on the phonetic database primarily reflect those listeners' representations
 891 (see SI, §3D). Future work could further address this assumption in a number of ways. One
 892 the one hand, dialect analyses like the ones we presented for our listeners (in Figure 5B)
 893 could compare listeners' templates against the templates used by talkers in the database.

894 Alternatively or additionally, researchers could see whether our results replicate if ideal
 895 observers are instead trained on other databases that have been hypothesized to reflect a
 896 ‘typical’ L1 listeners’ experience with US English.

897 Second, we made the simplifying assumption that listeners’ category representation—or
 898 at least the representations listeners’ drew on during the experiment—are talker-*independent*
 899 (we trained a single set of multivariate Gaussian categories, rather than, e.g., hierarchically
 900 organized set of multiple dialect templates). While this assumption is routinely made in
 901 research on normalization and beyond, it might well be wrong (see e.g., [Xie *et al.*, 2021](#)).

902 Finally, the evaluation of normalization accounts in the present work shares with all previ-
 903 ous work (e.g., [Apfelbaum and McMurray, 2015](#); [Barreda, 2021](#); [Cole *et al.*, 2010](#); [McMurray](#)
 904 and [Jongman, 2011](#); [Nearey, 1989](#); [Richter *et al.*, 2017](#)) another simplifying assumption that
 905 is clearly wrong: the assumption that listeners *know* the talker-specific formant properties
 906 required for normalization. Specifically, we normalized the input for each ideal observer
 907 using the maximum likelihood estimates of the normalization parameters for the respective
 908 experiment. For example, for the evaluation of the ideal observer trained on Lobanov nor-
 909 malized formants against listeners’ responses in Experiment 1a, we used the formant means
 910 and standard deviations of the stimuli used in Experiment 1a to normalize F1 and F2.

911 While this follows previous work, it constitutes a problematic assumption for the evaluation
 912 of extrinsic normalization accounts. For these accounts, the approach adopted essentially
 913 assumes the ability to predict the future: even on the first trial of the experiment, the input
 914 to the ideal observers were formants that were normalized based on the maximum likelihood
 915 estimate of the normalization parameters given the acoustic properties of *all* stimuli. Lis-

916 teners instead need to *incrementally infer* talker-specific properties from the speech input
 917 (Nearey and Assmann, 2007; Xie *et al.*, 2023). The development and testing of incremental
 918 variants of formant normalization strikes us an important avenue for future research.

919 **C. Concluding remarks**

920 We set out to compare how well competing accounts of formant normalization explain
 921 listeners' perception of vowels. We developed a computational framework that makes it
 922 possible to compare a large number of different accounts against multiple data sets. The
 923 code we share on OSF makes it possible to 'plug in' different accounts of vowel normalization,
 924 different phonetic databases, and different perception experiments. This, we hope, will
 925 substantially reduce the effort necessary to conduct similar evaluations on other datasets,
 926 dialects, and languages.

927 Comparing 20 of the most influential normalization accounts against L1 listeners' per-
 928 ception of US English monophthongs, we found that the normalization accounts that best
 929 describe listeners' perception share that they (1) learn and store talker-specific properties
 930 and (2) that they seem to be computationally very simple—taking advantage of the physics
 931 of sound generation to use as few as a single parameter to normalize inter-talker variability
 932 in vocal tract size. While the number of studies that have compared normalization accounts
 933 against *listeners'* behavior remains surprisingly small, these two results confirm the findings
 934 from more targeted comparisons that were focused on 2-3 accounts at a time (Barreda, 2021;
 935 Nearey, 1989; Richter *et al.*, 2017). Overall then, we submit that it is time for research in

936 speech perception and beyond to consider simple uniform scaling the most-likely candidate
937 for human formant normalization.

938 **ACKNOWLEDGMENTS**

939 Earlier versions of this work were presented at 2023 ASA meeting, ExLing 2022, at the
940 Department of Computational Linguistics at the University of Zürich and at the Depart-
941 ment of Swedish language and multilingualism at Stockholm University. We are grateful to
942 OMITTED FOR REVIEW.

943 **AUTHOR CONTRIBUTIONS**

944 AP designed the experiments and collected the data, with input from TFJ. TFJ pro-
945 grammed the experiments with input from AP. AP analyzed the experiments, with input
946 from TFJ. AP and TFJ wrote the code to implement and fit the normalization models, with
947 input from SB. AP developed the visualizations within input from SB and TFJ. AP wrote
948 the first draft of the manuscript with edits by SB and TFJ.

949 **SUPPLEMENTARY INFORMATION FOR *PERSSON, BARREDA & JAEGER***
 950 **(2024). COMPARING ACCOUNTS OF FORMANT NORMALIZATION AGAINST**
 951 **US ENGLISH LISTENERS' VOWEL PERCEPTION**

952 **§1. REQUIRED SOFTWARE**

953 Both the main text and these supplementary information (SI) are derived from the same
 954 R markdown document available via <https://osf.io/zemwn/>. It is best viewed using Acrobat
 955 Reader. The document was compiled using `knitr` in RStudio with R:

```
956 ## -  

957 ## platform      aarch64-apple-darwin20  

958 ## arch         aarch64  

959 ## os           darwin20  

960 ## system       aarch64, darwin20  

961 ## status  

962 ## major        4  

963 ## minor        3.2  

964 ## year         2023  

965 ## month        10  

966 ## day          31  

967 ## svn rev      85441  

968 ## language     R  

969 ## version.string R version 4.3.2 (2023-10-31)  

970 ## nickname     Eye Holes
```

971 Readers interested in working through the R markdown, and knitting it into a PDF will
 972 also need to download the IPA font **SIL Doulos** and a Latex environment like (e.g., **MacTex**
 973 or the R library **tinytex**).

974 We used the following R packages to create this document: R (Version 4.3.2; **R Core**
 975 **Team, 2023**) and the R-packages *assertthat* (Version 0.2.1; **Wickham, 2019**), *brms* (Ver-
 976 sion 2.21.0; **Bürkner, 2017, 2018, 2021**), *Cairo* (Version 1.6.2; **Urbanek and Horner, 2023**),
 977 *cmdstanr* (Version 0.5.3; **Gabry and Češnovar, 2022**), *dplyr* (Version 1.1.4; **Wickham *et al.***,

978 *2023*), *ellipse* (Version 0.5.0; Murdoch and Chow, 2023), *forcats* (Version 1.0.0; Wickham,
 979 2023a), *furrr* (Version 0.3.1; Vaughan and Dancho, 2022), *fuzzyjoin* (Version 0.1.6; Robinson, 2020), *ggforce* (Version 0.4.2; Pedersen, 2024a), *ggh4x* (Version 0.2.8; van den Brand, 2024), *ggnewscale* (Version 0.4.10; Campitelli, 2024), *ggplot2* (Version 3.5.1; Wickham, 2016),
 982 *ggtext* (Version 0.1.2; Wilke and Wiernik, 2022), *knitr* (Version 1.47; Xie, 2024), *linguistics-*
 983 *down* (Version 1.2.0; Liao, 2019), *lubridate* (Version 1.9.3; Grolemund and Wickham, 2011),
 984 *magrittr* (Version 2.0.3; Bache and Wickham, 2022), *mgcv* (Version 1.9.1; Wood *et al.*, 2016;
 985 Wood, 2003, 2004, 2011), *modelr* (Version 0.1.11; Wickham, 2023b), *MVBeliefUpdatr* (Ver-
 986 sion 0.0.1.10; Jaeger, 2024), *nlme* (Version 3.1.164; Pinheiro *et al.*, 2023), *patchwork* (Version
 987 1.2.0; Pedersen, 2024b), *phonR* (Version 1.0.7; McCloy, 2016), *phonTools* (Version 0.2.2.2;
 988 Barreda, 2023), *plotly* (Version 4.10.4; Sievert, 2020), *purrr* (Version 1.0.2; Wickham and
 989 Henry, 2023), *Rcpp* (Version 1.0.12; Eddelbuettel *et al.*, 2024), *readr* (Version 2.1.5; Wick-
 990 ham *et al.*, 2024a), *remotes* (Version 2.5.0; Csárdi *et al.*, 2024), *RJ-2021-048* (Bengtsson,
 991 2021), *rlang* (Version 1.1.4; Henry and Wickham, 2024), *stringr* (Version 1.5.1; Wickham,
 992 2023c), *tibble* (Version 3.2.1; Müller and Wickham, 2023), *tidybayes* (Version 3.0.6; Kay,
 993 2023), *tidyverse* (Version 2.0.0; Wickham *et al.*, 2019).

995 If opened in RStudio, the top of the R markdown document should alert you to any
 996 libraries you will need to download, if you have not already installed them. The full session
 997 information is provided at the end of this document.

998 **A. Interested in using R markdown do create APA formatted documents that
 999 integrate your code with your writing?**

1000 A project template, including R markdown files that result in APA-formatted PDFs, is
 1001 available at <https://github.com/hlplab/template-R-project>. Feedback welcome. We
 1002 aim to help others avoid the mistakes and detours we made when first deciding to embrace
 1003 literal coding to increase transparency in our projects.

1004 **§2. ADDITIONAL INFORMATION IN EXPERIMENTS 1A AND 1B**1005 **A. Participant exclusion**

1006 We adopted the following exclusion criteria: participants would get excluded if they failed
 1007 to pay attention to the instruction to wear over-the-ear-headphones, if they had unusually
 1008 slow or fast RT-means compared to other participants, or if they clearly did not do the task
 1009 (e.g., randomly clicking on different response options).

1010 N=1 participant in Experiment 1a was excluded based on the first criteria, as s/he used
 1011 external speakers instead of head set (based on response in post-experiment questionnaire).
 1012 This participant was also more than 3 standard deviations faster in her/his mean log-RTs
 1013 than other participants (second criteria), as were another participant in Experiment 1a. We
 1014 decided to exclude participants who were more than 3 standard deviations faster or slower
 1015 in their mean (log-transformed) RTs compared to other participants. We further excluded
 1016 *all trials* with RTs more than 3 standard deviations faster or slower than expected. This was
 1017 determined by first z-scoring the log-transformed RTs *within each participant* (by subtracting
 1018 the participants' mean from each observation and dividing through the participants standard
 1019 deviation) and then z-scoring these z-scores *within each trial* across participants. This
 1020 double-scaling approach was necessary as participants' RTs decreased substantially over
 1021 the first few trials and then continued to decrease less rapidly until converging against a
 1022 participant-specific minimum. This criterion did not remove just the first few trials but
 1023 rather removed RTs that were unusually fast or slow *for that participant at that trial*. And,
 1024 unlike more complicated methods (like developing a model of cross-trial decreases in RTs),
 1025 the approach employed here does not make any assumptions about the shape of the speed
 1026 up in RTs across trials. In total, N=117 trials were excluded, however, no participant was
 1027 excluded based on too high proportion of missing trials. Figure S1 summarizes participant
 1028 exclusions due to reaction times and not wearing headphones.

1029 The experiments did not contain independent catch trials. We therefore looked into par-
 1030 ticipants' individual responses in order to identify participants that seem to have randomly
 1031 answered, independent of stimulus. Figure S2 suggests that participants 13, 15, 21, 22, and
 1032 24, have not performed the task. Their responses are indicating that they have not payed
 1033 attention to the stimuli but rather randomly selected responses irrespective of the vowel they
 1034 heard, hence, e.g., *had* responses for all different locations in the vowel space (e.g., partici-

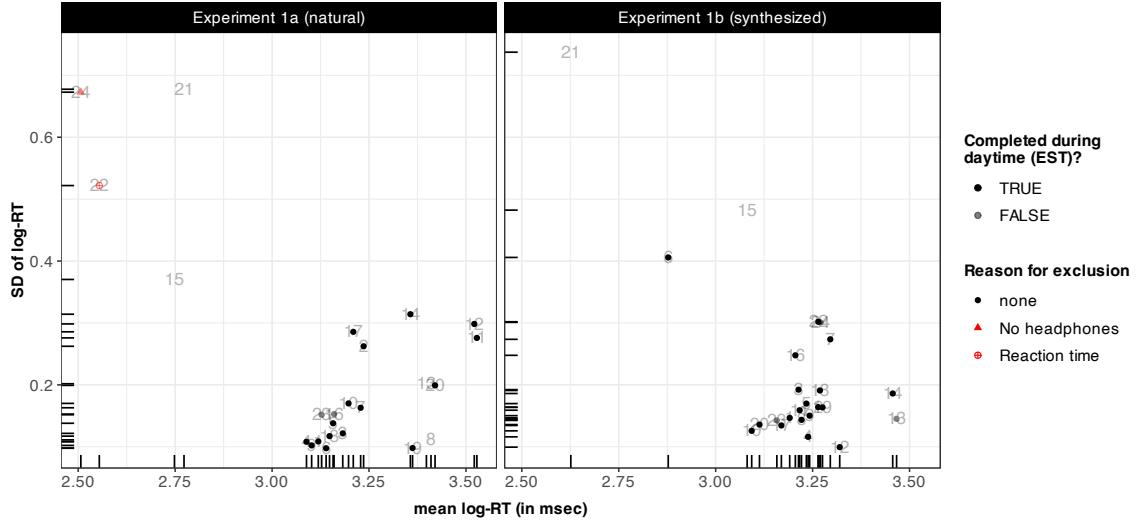


FIG. S1. Participant exclusions in Experiment 1a and Experiment 1b. Two participants in Experiment 1a (participants 22 and 24) were excluded based on their log-transformed RTs and/or for not wearing headphones.

1035 pants 15, 21, 22, 24), or *heed* responses for the low back vowels (e.g., 24), or *odd* responses
 1036 for high front vowels (e.g., participants 13, 21). Participant 8 seem to have performed the
 1037 task but clearly used the phonetic space in ways different from everyone else, as s/he seems
 1038 to have inverted the *who'd-hood* categories. This behaviour more likely indicates different
 1039 dialect patterns, we still however, decided to exclude participant 8 as well. Participants
 1040 22 and 24 had already been excluded based on unusual RT-patterns (22) and not wearing
 1041 over-the-ear headphones (24).

1042 Excluding participants because of unusual vowel responses is more complicated for Ex-
 1043 periment 1b as the stimuli is synthesized and more difficult to categorize in general. Figure
 1044 S3 nevertheless indicates that most participants made use of the phonetic space in similar
 1045 ways (and in line with where natural categories fall), besides participant 15 and 21. Partic-
 1046 ipant 15 often responds *who'd* for tokens in the high front part of the space, and *heed* for
 1047 tokens in the high center and back parts, while participant 21 is overall more random in
 1048 responses, but often selects *heed* for high back tokens and several times selects *hod* for front
 1049 tokens.

1050 Unusual vowel responses is not an objective criterion. There are many participants that
 1051 gave unexpected responses on some occasions, e.g., participants 10 and 12 in Experiment
 1052 1a, or participants 5 and 24 in Experiment 1b, however, we decided not to exclude them

1053 as they were not systematic in their response patterns, i.e., there were no indications of a
 1054 definite dialect shift (as with participant 8 in Experiment 1a), or systematic randomness in
 1055 selecting any kind of vowel for any kind of stimuli (as with participants 13, 15, 21, 22, 24,
 1056 in Experiment 1a, and to some extent, participants 15 and 22 in Experiment 1b).

1057 **B. Distribution of stimuli F1-F3 in Experiments 1a and 1b**

1058 Figure S4 visualizes the stimuli in Experiments 1a and 1b in F1-F3 space.

1059 **C. Auxiliary analysis of participant responses in Experiments 1a and 1b**

1060 Participants in Experiment 1b showed overall less agreement in their responses to the
 1061 stimuli than participants in Experiment 1a, as indicated by the higher response entropy in
 1062 Experiment 1b. In order to assess the extent to which this was a result of the placement of
 1063 the tokens in the F1-F3 space, we compared linear regression models that predicted response
 1064 entropy from experiment, to models that employed residuals from a general additive model
 1065 including response entropy and the tokens placement in the phonetic space (F1, F2) as
 1066 response variable, and experiment as predictor. We furthermore compared against models
 1067 based on the full data set to models that excluded all *hut* and *odd* responses from Experiment
 1068 1a in order to assess effects of lexical context.

1069 When adding effects of lexical context to the model, the difference between experiments is
 1070 reduced by 23.9%. Adding a nonlinear model with F1-F2 values of the tokens, the difference
 1071 is reduced by an additional 35.6%, while adding F3-values reduces the difference by 50.9%.
 1072 In sum, the result suggest that approximately two-thirds of the difference in response entropy
 1073 between experiments can be attributed to the placement of stimuli in the formant space,
 1074 while the remaining one-third is influenced by other factors, most likely the synthesized
 1075 stimuli sounding highly unnatural.

1076 As stated in the main paper, response entropies differed even for tokens that overlap in
 1077 Hertz space. Figure S5 visualizes differences in categorization behaviour for these tokens.
 1078 For many of these tokens, the most frequent response is the same category across exper-
 1079 iments, however, with substantially higher disagreement for tokens in Experiment 1b. In
 1080 the bottom part of the acoustic space, participants in Experiment 1b seem to respond *had*
 1081 disproportionately often.

Comparing normalization against perception



FIG. S2. Participants' categorization responses in Experiment 1a, shown in F1-F2 space. Color and vowel label indicate response provided by participants on each test location. Each vowel was repeated twice.

Comparing normalization against perception

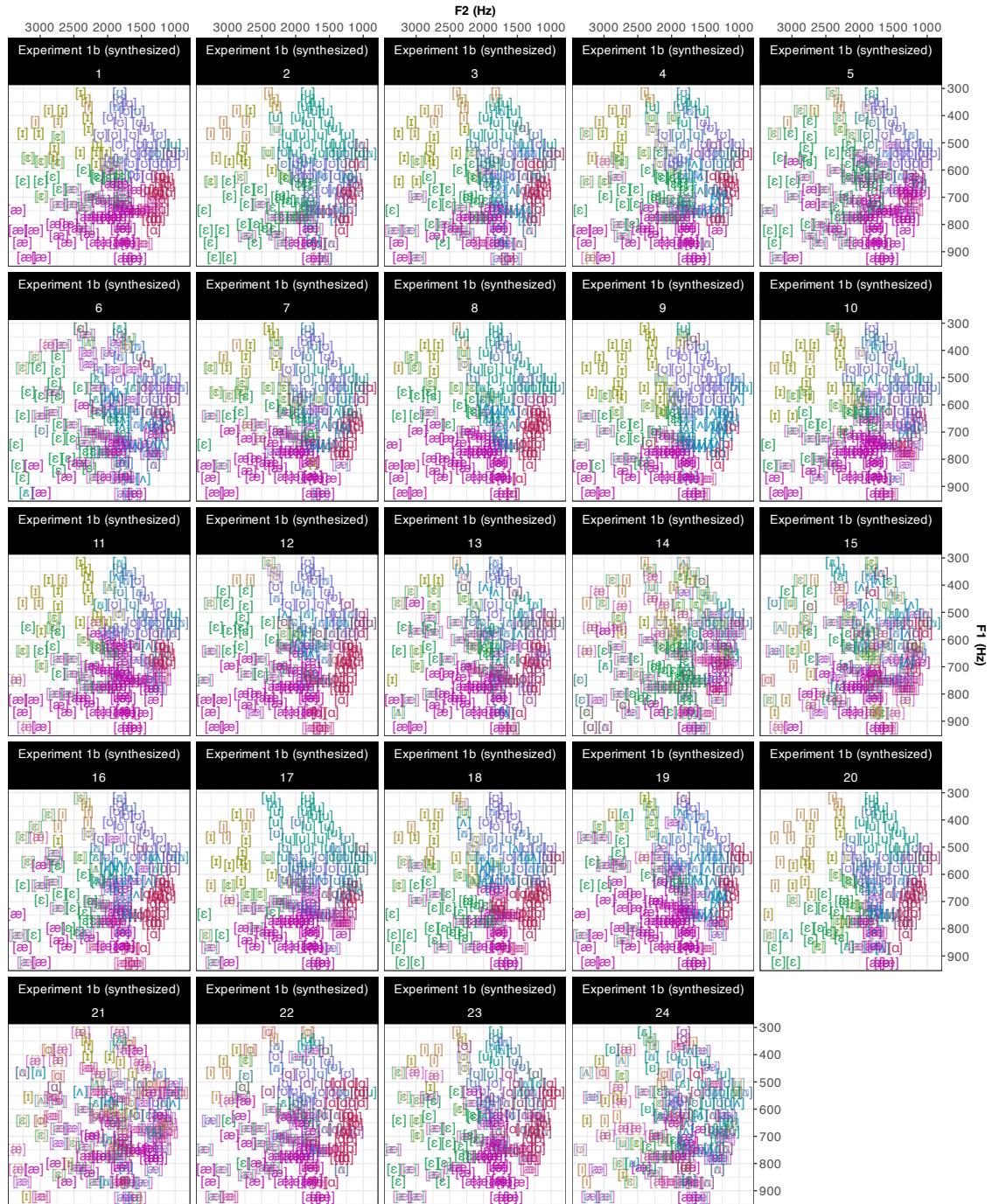


FIG. S3. Participants' categorization responses in Experiment 1b, shown in F1-F2 space. Color and vowel label indicate response provided by participants on each test location. Each vowel was repeated twice.

Comparing normalization against perception

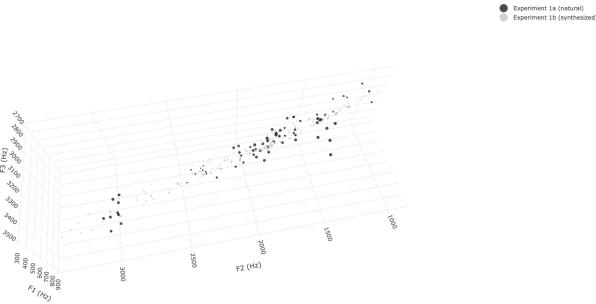


FIG. S4. Stimuli of Experiments 1a and 1b in F1-F3 space. Point size indicates response entropy: larger points represent higher listener agreement, and vice versa.

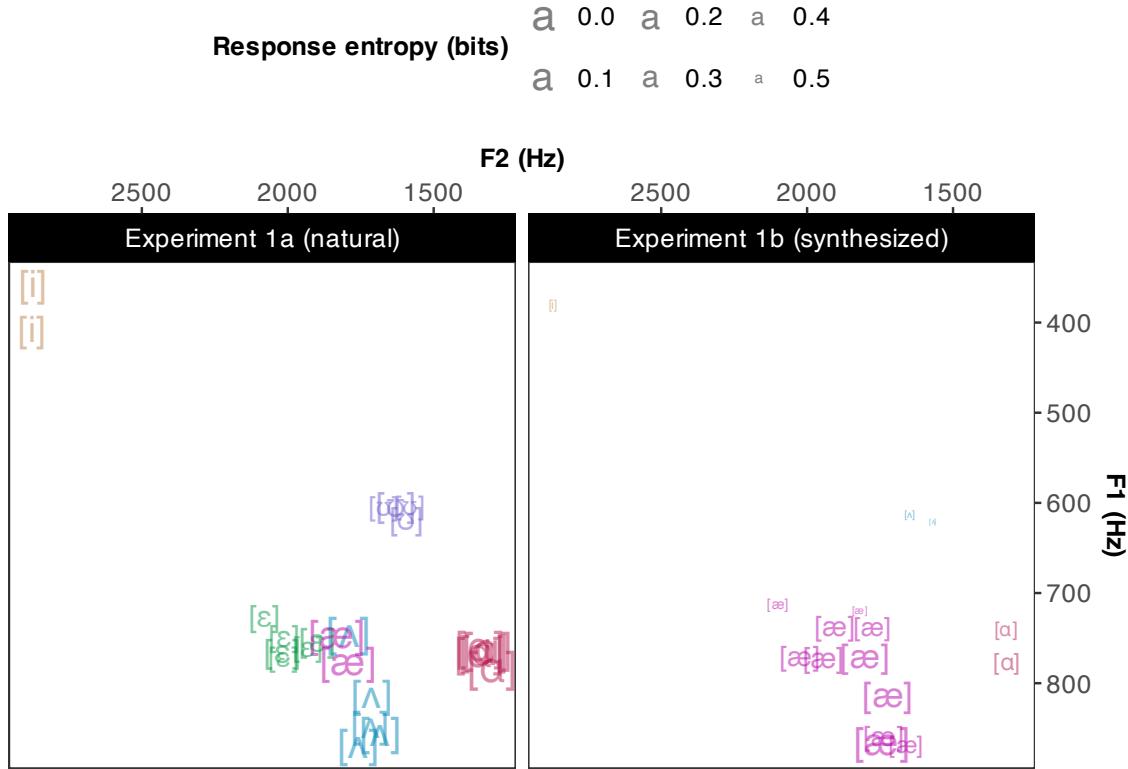


FIG. S5. Listeners' categorization responses in Experiments 1a and 1b, for comparable tokens in Hertz space. The vowel label indicates the most frequent response provided by participants on each test location. Size indicates how consistent responses were across participants, which larger symbols indicating more consistent responses (lower entropy).

1082 **§3. ADDITIONAL INFORMATION ON THE COMPUTATIONAL COMPAR-
1083 SON OF NORMALIZATION ACCOUNTS**

1084 **A. Methods**

1085 **1. Vowel data used to train ideal observers (Xie and Jaeger, 2020)**

1086 The Xie and Jaeger database consists of $N=1168$ hVd word recordings from 17 (5 female)
 1087 L1 talkers of a Northeastern dialect of US English (ages 18 to 35 years old). The talkers
 1088 were recorded reading a list of 180 English monosyllabic words, a list of short sentences, and
 1089 a list of ten hVd words—the eight US English monophthongs as well as *aid* and *owed* (for
 1090 further information, see [Xie and Jaeger, 2020](#)). For each talker, the database contains 9-10
 1091 recordings of each hVd word. An automatic aligner [Penn Phonetics Lab Forced Aligner;
 1092 [Yuan and Liberman \(2008\)](#)] was used to obtain estimates for word and segment boundaries.¹²

1093 The first author manually corrected the automatic alignments for all vowel segmentations.
 1094 We then used the Burg algorithm in Praat ([Boersma and Weenink, 2022](#)) to extract estimates
 1095 of the first three formants (F1-F3) at three points of the vowel (35, 50, and 65 percent into
 1096 the vowel). The following parameterization of the Burg algorithm was used:

- 1097 • Time step (s): 0.01
- 1098 • Max. number of formants: 5
- 1099 • Formant ceiling (Hz): 5500 (5000 for the male talkers)
- 1100 • Window length (s): 0.025
- 1101 • Pre-emphasis from (Hz): 50

1102 In addition to F1-F3, we automatically extracted vowel duration and the fundamental
 1103 frequency (F0) across the entire vowel. These are the data that we used in the cross-
 1104 validation procedure to train ideal observers, as described in the main text. Figure S6
 1105 visualizes the vowel data from the [Xie and Jaeger \(2020\)](#) for all pairwise combinations of
 1106 F0, F1, F2, F3 and vowel duration, in raw Hertz. Figure S7 shows the distribution of F1
 1107 and F2 in the different normalization spaces used in the main study.

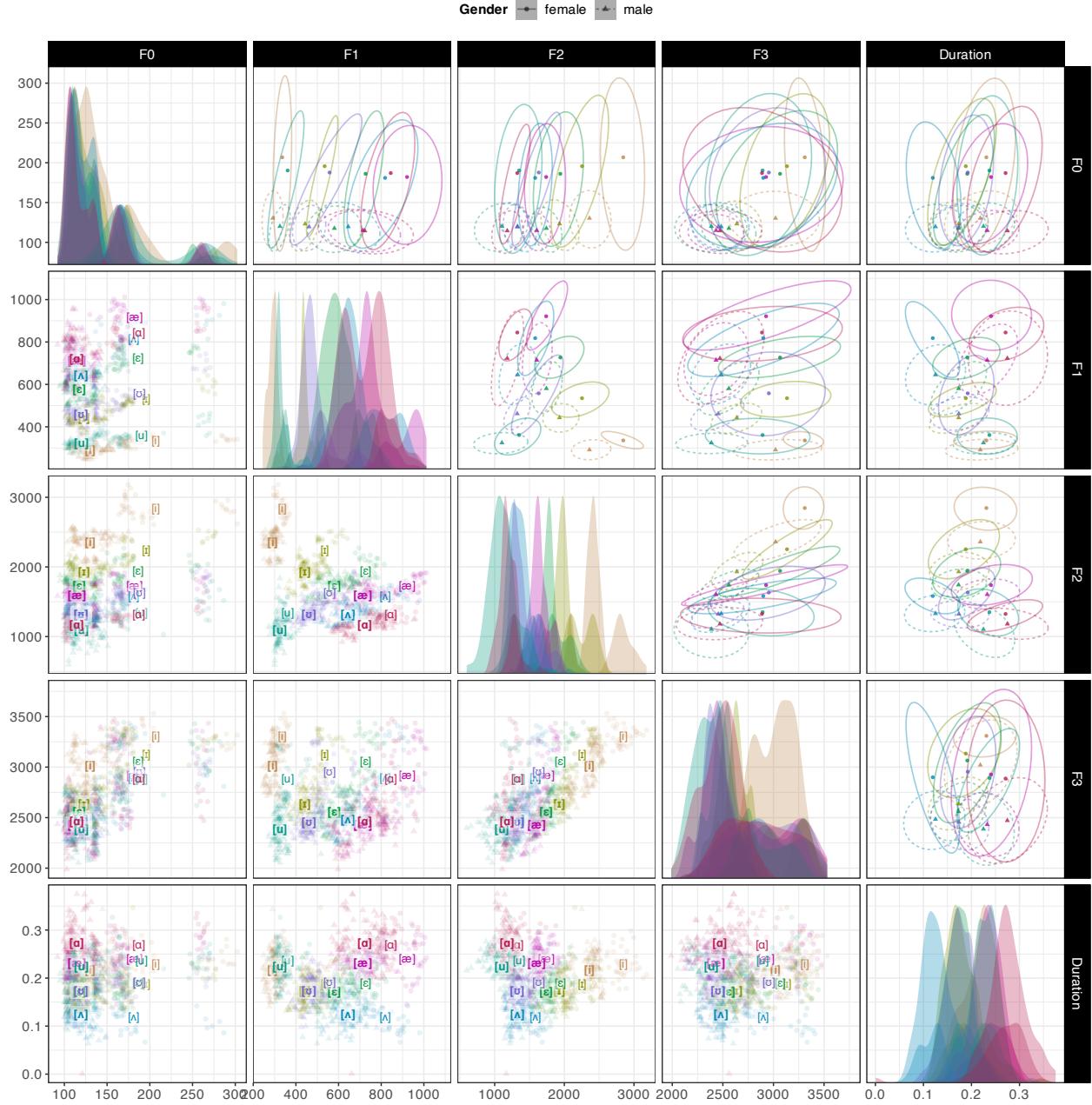


FIG. S6. The pairwise distributions of F0, F1, F2, F3, and duration for all 1168 recordings of the eight monophthong hVd words in [Xie and Jaeger \(2020\)](#). Note that axis directions are not reversed.

Panels on diagonal: marginal cue densities of all five cues. **Lower off-diagonal panels:** each point corresponds to a recording, averaged across the three measurement points within each vowel segment. Vowel labels indicate category means across talkers. Male talkers' vowels are boldfaced. **Upper off-diagonal panels:** Same data as in the lower off-diagonal panels but showing bivariate Gaussian 95% probability mass ellipses around category means.

Comparing normalization against perception

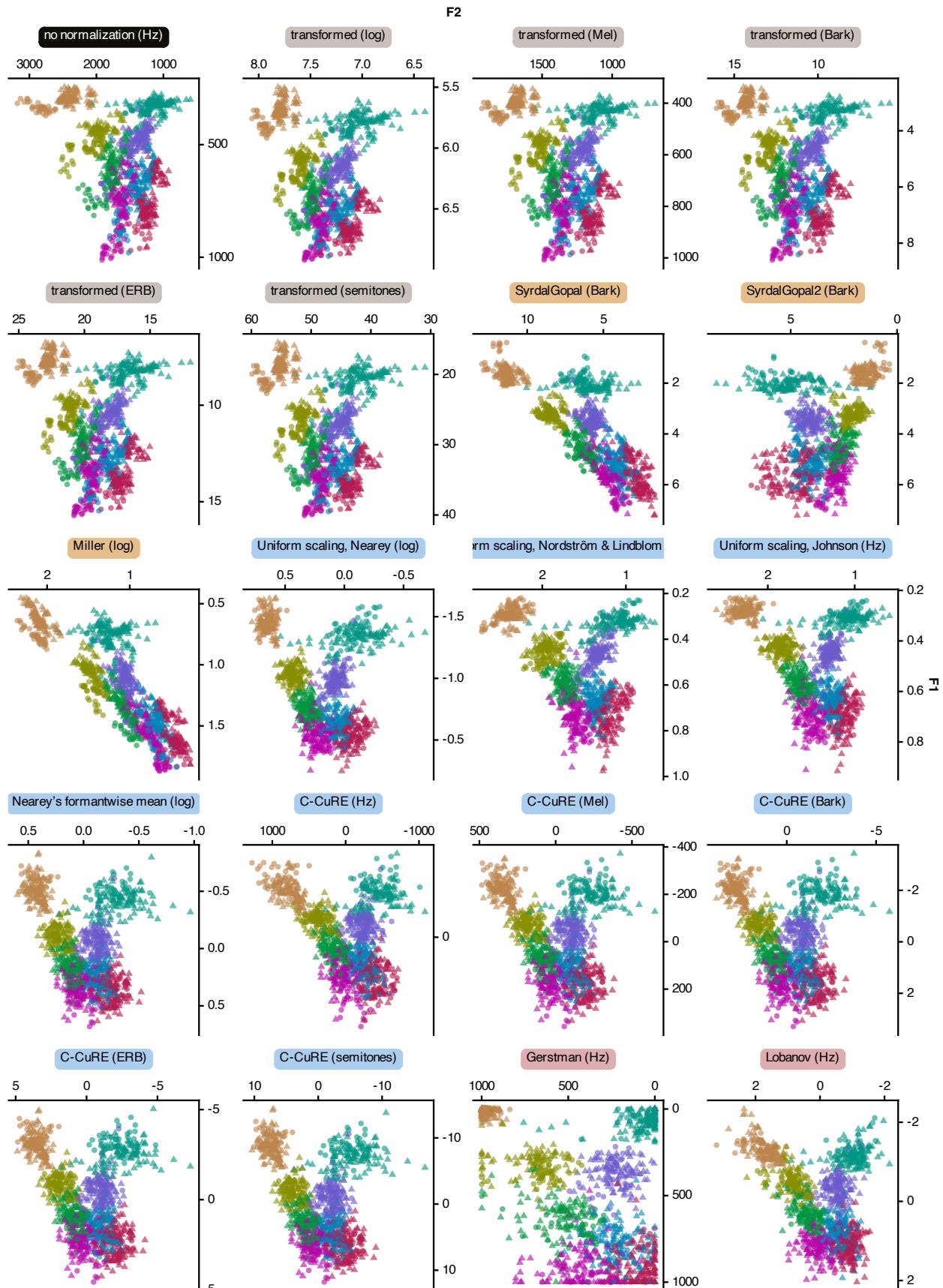


FIG. S7. The 8 monophthong vowels of US English from the [Xie and Jaeger \(2020\)](#) database when F1 and F2 are transformed into a perceptual scale (grey), intrinsically normalized (yellow), or extrinsically normalized through centering (blue) or standardizing (purple). Each point corre-

1108 **2. Normalization parameters θ**

1109 Figure S8 relates the normalization parameters θ obtained for each experiment to those
 1110 found for the five training sets of the [Xie and Jaeger \(2020\)](#) database. This servers two pur-
 1111 poses. First, by comparing the θ of Experiment 1a, which was based on natural productions,
 1112 to the θ obtained from [Xie and Jaeger \(2020\)](#), we can assess the extent to which the talker
 1113 used for Experiment 1a is ‘typical’ relative to the other talkers of that database. Second,
 1114 by comparing the range and variability of the θ across normalization accounts and experi-
 1115 ments, we can assess the volatility of different types of parameters, and assess the difference
 1116 between the beliefs the ideal observers have about the parameters and the parameters in
 1117 the experiment. How reliably the statistics of the input is established for the same amount
 1118 of data seems to depend on the space. For instance, parameters in Hertz space display
 1119 more variability. Within a given scale, we also note that some parameters are more difficult
 1120 to estimate than others, for instance, mean estimates display less variability than SD, and
 1121 range values (min and max).

1122 **3. Optimization process to fit models to human responses**

1123 We used constrained quasi-Newton optimization ([Byrd *et al.*, 1995](#)) to determine the best-
 1124 fitting values for the two degrees of freedom—lapse rate (λ) and noise ratio (τ^{-1}). Optimiza-
 1125 tion was performed separately for each of the 200 combinations of normalization account,
 1126 experiment, and cross-validation fold. Specifically, we maximized the *likelihood* of the hu-
 1127 man categorization responses in each experiment under the categorization model conditional
 1128 on the model’s lapse rate and perceptual noise, $\Sigma_i^N \log p(response_i | F1_{i,\theta}, F2_{i,\theta}, M_{\theta,\lambda,\Sigma_{noise}})$,
 1129 where $response_i$ is the i th categorization response, $F1_{i,\theta}, F2_{i,\theta}$ are the F1 and F2 values for
 1130 the i th observation after normalization (with parameters θ being estimated based on the
 1131 distribution of phonetic cues across the stimuli in the experiment). $M_{\theta,\lambda,\Sigma_{noise}}$ is the cate-
 1132 gorization model in Figure 6, with normalization parameters θ fixed based from the prior
 1133 cue distribution in the phonetic database ([Xie and Jaeger, 2020](#)), and λ and Σ_{noise} as the
 1134 only free parameters to maximize the likelihood. The best-fitting parameterizations were
 1135 determined by means of the `optim()` function in R’s `stats` package ([R Core Team, 2023](#)).
 1136 The starting value of lapse rates and perceptual noise were set to 0.1 and 0.15, respectively.
 1137 We set the lower and upper bounds to $10^{-10} \geq$ lapse rate ≥ 1 , and $10^{-10} \geq$ perceptual noise

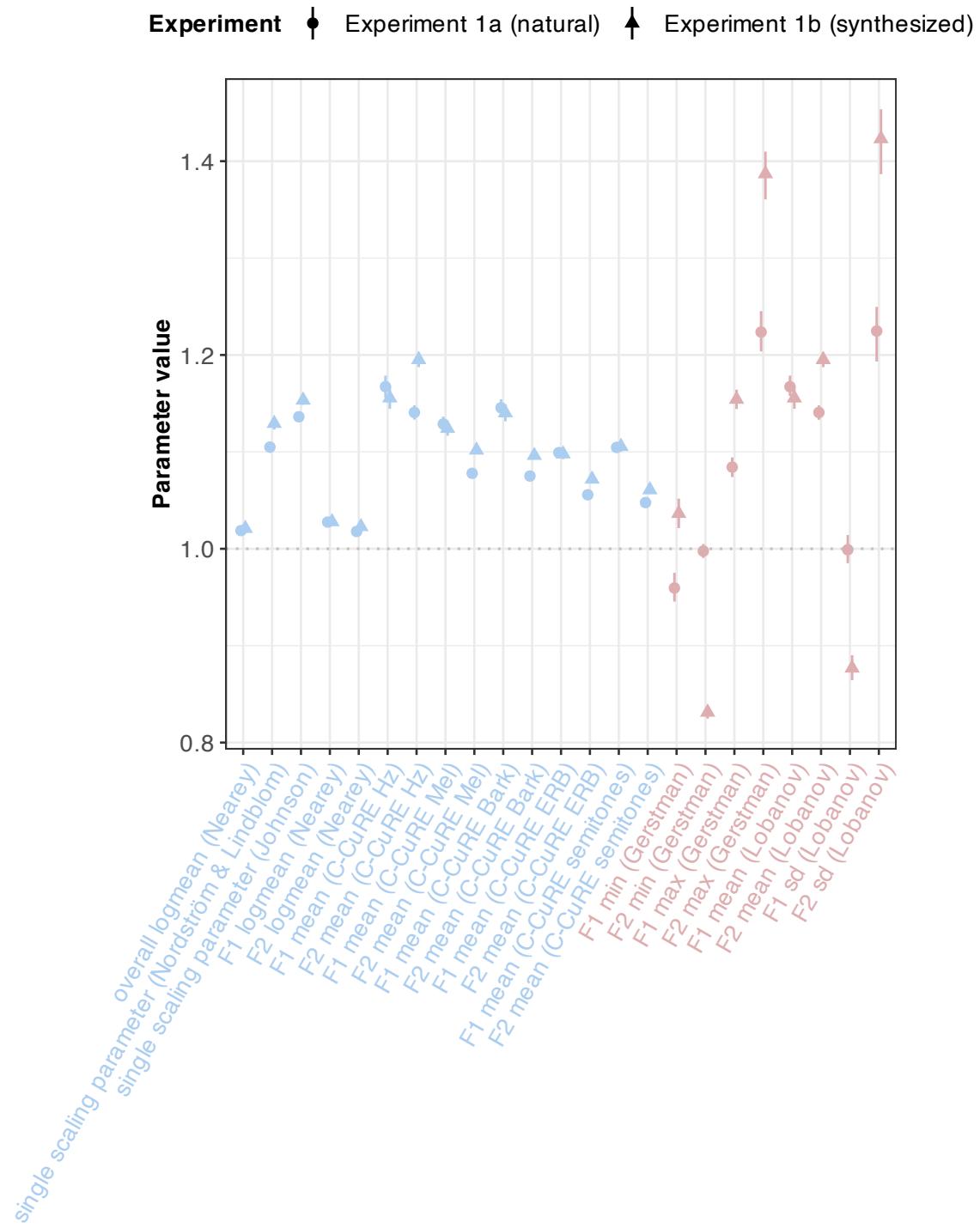


FIG. S8. Comparing normalization parameters θ across the phonetic database used to estimate listeners' prior experience (Xie and Jaeger, 2020) and Experiments 1a and 1b. Only accounts that assume talker-specific normalization parameters are shown.

₁₁₃₈ ≥ 10 (well above previously observed estimates for perceptual noise in [Kronrod *et al.*, 2016](#),
₁₁₃₉ p. 1698).

₁₁₄₀ Section [§3 F](#) presents additional analyses that instead used a grid search over the parameter
₁₁₄₁ space. These analyses confirm the results presented in the main paper.

₁₁₄₂ **B. Results for F1-F2**

₁₁₄₃ **1. Significance test of model performance**

₁₁₄₄ Tables [S1](#) and [S2](#) present the results from the paired one-sided t-tests conducted, pre-
₁₁₄₅ dicting model log likelihood as a function of normalization account for Experiment 1a and
₁₁₄₆ 1b (dummy coded with no normalization model as reference model). The log likelihoods are
₁₁₄₇ averaged across the five cross-validation folds and ordered by best-fitting models.

₁₁₄₈ **2. Parameter estimates for best-fitting models**

₁₁₄₉ In this section, we provide the estimates found for the two degrees of freedom—noise
₁₁₅₀ (Σ_{noise}) and attentional lapses (λ)—when fitting the models to human behaviour. This will
₁₁₅₁ provide insights into the relative contribution of these factors to explaining the variability
₁₁₅₂ found in the behavioral data between the two experiments, and to understanding the relative
₁₁₅₃ performance of the different normalization accounts as models of human behavior.

₁₁₅₄ Figure [S9](#) visualizes the parameter estimates for each account, averaged across the five
₁₁₅₅ training sets (see also Tables [S3](#), [S4](#) for summary of fitted values, and [S10](#) for an illustration
₁₁₅₆ of how the fitted noise affects the bivariate Gaussian categories). The Figure indicates that
₁₁₅₇ fitted λ s are very similar across experiments. In Experiment 1a, the mean λ across models
₁₁₅₈ was 0.1 ($sd = 0.04$), and for Experiment 1b, 0.12 ($sd = 0.14$). These estimates can in part
₁₁₅₉ confirm what was hypothesized with regard to listeners' categorization accuracy—that the
₁₁₆₀ performance of listeners in inferring the category intended by the talker in part reflected
₁₁₆₁ attentional lapses (mean accuracy in Experiment 1a = 81.2% ($SE = 4.8\%$); Experiment 1b
₁₁₆₂ had no such ground truth).

₁₁₆₃ What is perhaps more obvious from Figure [S9](#) is that Σ_{noise} estimates clearly differ
₁₁₆₄ between experiments. In Experiment 1a, the best-fitting Σ_{noise} estimates are comparable to
₁₁₆₅ what [Kronrod *et al.* \(2016\)](#) found (mean Σ_{noise} in Experiment 1a = 0.52 ($sd = 0.49$)). In

TABLE S1. T-test predicting the model log likelihood as a function of normalization account for Experiment 1a

Normalization account	Statistic	Estimate mean	Diff. in means	p_value
Uniform scaling, Johnson (Hz)	-15.085	-2523.406	611.644	0.000
Uniform scaling, Nearey (log)	-9.100	-2523.406	551.676	0.000
C-CuRE (Bark)	-13.229	-2523.406	549.192	0.000
C-CuRE (Mel)	-12.722	-2523.406	546.779	0.000
C-CuRE (ERB)	-11.847	-2523.406	546.472	0.000
Nearey's formantwise mean (log)	-10.428	-2523.406	543.718	0.000
C-CuRE (semitones)	-10.428	-2523.406	543.718	0.000
Uniform scaling, Nordström & Lindblom (Hz)	-7.791	-2523.406	517.075	0.001
SyrdalGopal (Bark)	-9.381	-2523.406	510.360	0.000
C-CuRE (Hz)	-10.169	-2523.406	498.443	0.000
Miller (log)	-7.466	-2523.406	464.305	0.001
Lobanov (Hz)	-8.970	-2523.406	461.405	0.000
transformed (Bark)	-13.415	-2523.406	228.190	0.000
transformed (Mel)	-12.472	-2523.406	214.317	0.000
transformed (ERB)	-10.291	-2523.406	192.156	0.000
SyrdalGopal2 (Bark)	-3.673	-2523.406	171.154	0.011
Gerstman (Hz)	-2.104	-2523.406	159.477	0.052
transformed (log)	-2.617	-2523.406	66.928	0.029
transformed (semitones)	-2.617	-2523.406	66.928	0.029

¹¹⁶⁶ Experiment 1b, this is not the case (mean $\Sigma_{noise} = 4.74$ (sd=2.57). However, there is no a priori reason to expect internal perceptual noise to differ between experiments, which is why these high noise ratios likely reflect external noise. Given what was shown for the human data (cf. discussion on differences in response entropy between experiments, Section II B),

TABLE S2. T-test predicting the model log likelihood as a function of normalization account for Experiment 1b

Normalization account	Statistic	Estimate mean	Diff. in means	p_value
Uniform scaling, Nearey (log)	-64.722	-10372.71	2553.594	0.000
Lobanov (Hz)	-82.707	-10372.71	2521.647	0.000
Gerstman (Hz)	-34.270	-10372.71	2491.092	0.000
C-CuRE (ERB)	-35.509	-10372.71	2409.974	0.000
C-CuRE (Bark)	-33.226	-10372.71	2408.381	0.000
Nearey's formantwise mean (log)	-35.495	-10372.71	2305.669	0.000
C-CuRE (semitones)	-35.495	-10372.71	2305.669	0.000
transformed (log)	-63.629	-10372.71	2259.438	0.000
transformed (semitones)	-63.629	-10372.71	2259.438	0.000
C-CuRE (Mel)	-28.933	-10372.71	2221.903	0.000
transformed (ERB)	-51.487	-10372.71	2106.368	0.000
transformed (Bark)	-46.651	-10372.71	1942.912	0.000
SyrdalGopal2 (Bark)	-29.866	-10372.71	1816.140	0.000
transformed (Mel)	-41.633	-10372.71	1622.458	0.000
Miller (log)	-26.334	-10372.71	1244.978	0.000
Uniform scaling, Johnson (Hz)	-8.574	-10372.71	1069.757	0.001
SyrdalGopal (Bark)	-20.657	-10372.71	910.361	0.000
Uniform scaling, Nordström & Lindblom (Hz)	-11.293	-10372.71	878.721	0.000
C-CuRE (Hz)	-10.441	-10372.71	817.140	0.000

1170 this is perhaps not surprising. The stimuli in Experiment 1b were clearly more noisy and
 1171 presumably left listeners with more uncertainty about the true value of the formants, and
 1172 how to best make use of previous experience. Even if the task itself was identical across
 1173 experiments, the nature of the stimuli in Experiment 1b likely contributed to making the

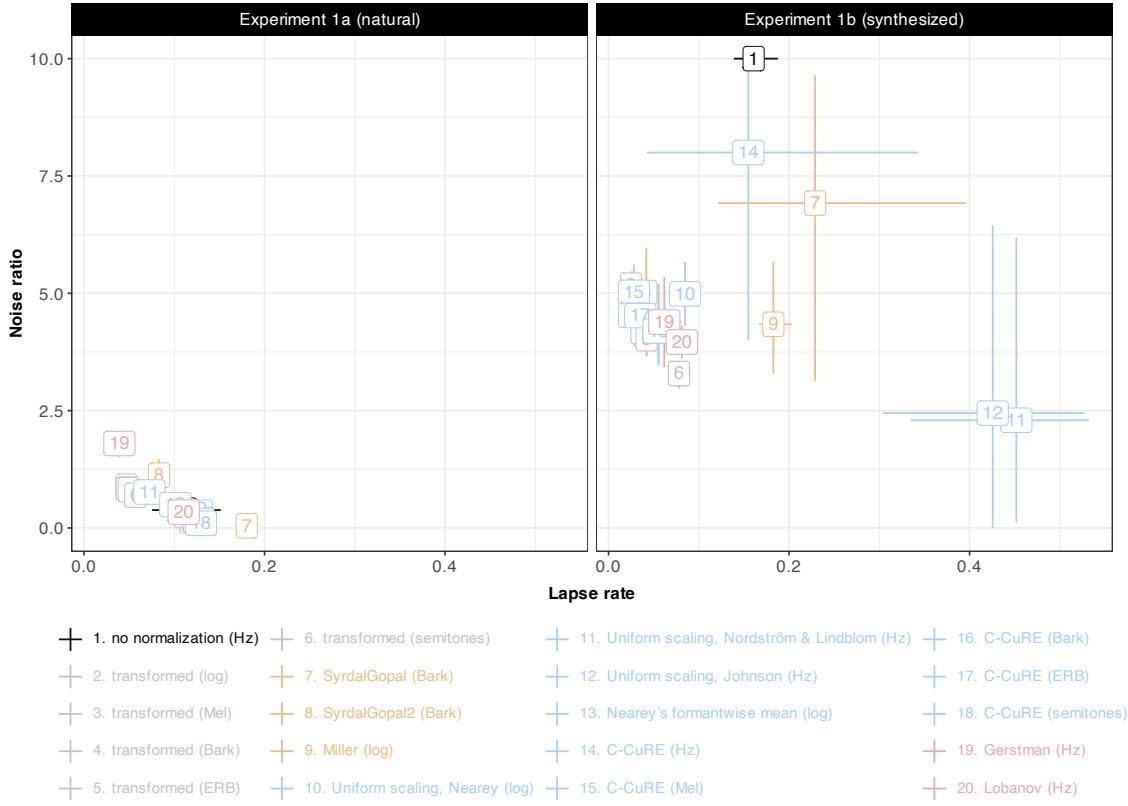


FIG. S9. Best-fitting estimates obtained for λ and Σ_{noise} . Numeric label is placed at the mean across the five folds, line ranges represent the 95% CIs.

1174 experiment overall more demanding. In addition, Experiment 1b was longer (N=74 more
 1175 trials, and took on average 8.1 more minutes to complete), both of these aspects might also
 1176 affect the amount of attentional lapsing.

1177 Finally, for the majority of accounts, there is little variability in parameter estimates—and
 1178 likelihoods—across training sets (Tables S3 and S4). This suggests that models achieved
 1179 their maximum likelihood fit to human data on similar estimates for the two degrees of
 1180 freedom, which provides a sanity check of the modelling approach adopted.

1181 3. By-item analysis

1182 To provide further insight into model performance, we visualize model fits against hu-
 1183 man behavior on a by-item level for three of the best-performing models across experiments,
 1184 Nearey's uniform scaling, Johnson's uniform scaling and Lobanov. This allows us to as-

TABLE S3. The best-fitting estimates obtained for noise ratios and lapse rates in Experiment 1a (averaged across the five cross-validation folds and ordered by best-performing models)

Normalization account	mean likelihood	log noise percentage	lapse rate
Uniform scaling, Johnson (Hz)	-2214.93	mean=0.5 (SD=0.28)	mean=0.1 (SD=0.01)
Uniform scaling, Nordström & Lindblom (Hz)	-2219.16	mean=0.77 (SD=0.39)	mean=0.07 (SD=0)
C-CuRE (Bark)	-2261.27	mean=0.28 (SD=0.26)	mean=0.12 (SD=0.02)
C-CuRE (Mel)	-2261.61	mean=0.29 (SD=0.27)	mean=0.12 (SD=0.02)
C-CuRE (ERB)	-2262.53	mean=0.16 (SD=0.23)	mean=0.13 (SD=0.02)
C-CuRE (semitones)	-2264.28	mean=0.1 (SD=0.17)	mean=0.13 (SD=0.01)
Nearey's formantwise mean (log)	-2264.28	mean=0.1 (SD=0.17)	mean=0.13 (SD=0.01)
Uniform scaling, Nearey (log)	-2273.68	mean=0.34 (SD=0.36)	mean=0.13 (SD=0.02)
C-CuRE (Hz)	-2305.16	mean=0.19 (SD=0.15)	mean=0.13 (SD=0.02)
SyrdalGopal (Bark)	-2357.42	mean=0.05 (SD=0.04)	mean=0.18 (SD=0.01)
Miller (log)	-2374.96	mean=0.15 (SD=0.06)	mean=0.12 (SD=0.01)
Lobanov (Hz)	-2390.23	mean=0.34 (SD=0.27)	mean=0.11 (SD=0.02)
transformed (Bark)	-2569.87	mean=0.83 (SD=0.2)	mean=0.05 (SD=0)
transformed (Mel)	-2587.65	mean=0.89 (SD=0.18)	mean=0.05 (SD=0)
transformed (ERB)	-2598.16	mean=0.81 (SD=0.24)	mean=0.05 (SD=0)
Gerstman (Hz)	-2647.38	mean=1.8 (SD=0.33)	mean=0.04 (SD=0)
SyrdalGopal2 (Bark)	-2661.45	mean=1.12 (SD=0.39)	mean=0.08 (SD=0.01)
transformed (semitones)	-2682.60	mean=0.7 (SD=0.34)	mean=0.06 (SD=0.01)
transformed (log)	-2682.60	mean=0.7 (SD=0.34)	mean=0.06 (SD=0.01)
no normalization (Hz)	-2779.13	mean=0.38 (SD=0.26)	mean=0.11 (SD=0.05)

TABLE S4. The best-fitting estimates obtained for noise ratios and lapse rates in Experiment 1b (averaged across the five cross-validation folds and ordered by best-performing models)

Normalization account	mean	log noise percentage	lapse rate
	likelihood		
Gerstman (Hz)	-9474.08	mean=4.39 (SD=1.24)	mean=0.06 (SD=0.03)
Uniform scaling, Nearey (log)	-9551.73	mean=4.99 (SD=0.87)	mean=0.08 (SD=0)
C-CuRE (Bark)	-9595.30	mean=4.53 (SD=0.9)	mean=0.03 (SD=0.02)
C-CuRE (ERB)	-9601.26	mean=4.52 (SD=0.93)	mean=0.04 (SD=0.02)
Lobanov (Hz)	-9608.45	mean=3.96 (SD=0.47)	mean=0.08 (SD=0.01)
Nearey's formantwise mean (log)	-9702.20	mean=4.2 (SD=1.08)	mean=0.06 (SD=0.03)
C-CuRE (semitones)	-9702.20	mean=4.2 (SD=1.08)	mean=0.06 (SD=0.03)
C-CuRE (Mel)	-9772.00	mean=5.02 (SD=0.71)	mean=0.03 (SD=0.01)
transformed (semitones)	-9815.97	mean=3.29 (SD=0.37)	mean=0.08 (SD=0.01)
transformed (log)	-9815.97	mean=3.29 (SD=0.37)	mean=0.08 (SD=0.01)
transformed (ERB)	-9956.82	mean=4.03 (SD=0.42)	mean=0.04 (SD=0.01)
transformed (Bark)	-10123.34	mean=4.18 (SD=0.42)	mean=0.04 (SD=0.01)
SyrdalGopal2 (Bark)	-10231.19	mean=5.02 (SD=1.14)	mean=0.04 (SD=0.01)
transformed (Mel)	-10431.06	mean=5.17 (SD=0.39)	mean=0.02 (SD=0.01)
Uniform scaling, Johnson (Hz)	-10791.39	mean=2.45 (SD=4.33)	mean=0.43 (SD=0.14)
Miller (log)	-10848.88	mean=4.34 (SD=1.54)	mean=0.18 (SD=0.02)
Uniform scaling, Nordström & Lindblom (Hz)	-11085.94	mean=2.29 (SD=4.32)	mean=0.45 (SD=0.13)
C-CuRE (Hz)	-11098.93	mean=8 (SD=4.47)	mean=0.15 (SD=0.21)
SyrdalGopal (Bark)	-11187.23	mean=6.92 (SD=4.17)	mean=0.23 (SD=0.19)
no normalization (Hz)	-12118.49	mean=10 (SD=0)	mean=0.16 (SD=0.03)

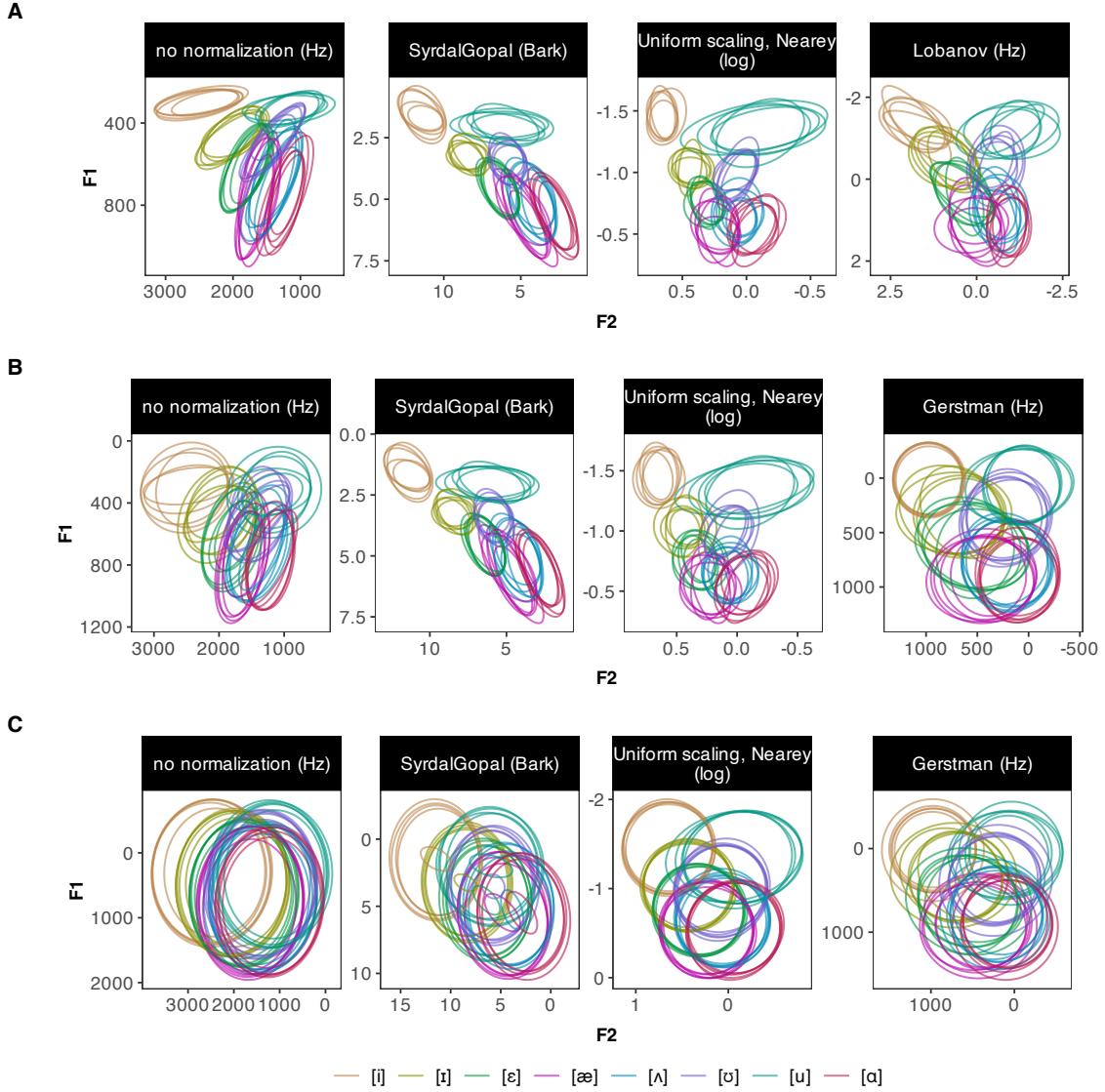


FIG. S10. Visualizing the bivariate Gaussian categories of four example normalization accounts for each of the five cross-validation folds (each fold corresponds to one set of eight ellipses). **Panel A** prior to adding Σ_{noise} , **Panel B** with added noise from best-fitting models in Experiment 1a, **Panel C** with added noise from best-fitting models in Experiment 1b. For most of the accounts in Panel B and C, noise ratios adds so much category variability that models could presumably only make correct predictions at the outer range of the ellipses. If allowing for separate noise estimates for F1 and F2, this might however not be the case.

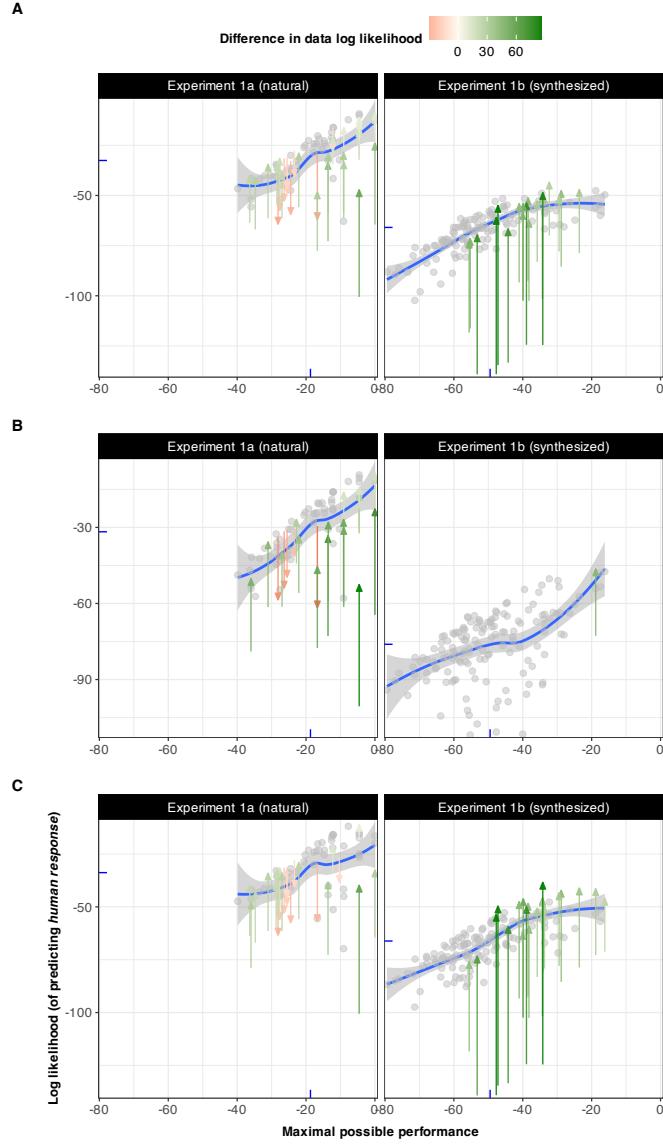


FIG. S11. By-item model improvement from no normalization, relative to the maximum possible performance (predicting human responses from human responses). Maximum log likelihood across items indicated by ticks on axis. Arrows indicate change from no normalization to Nearey's uniform scaling (**panel A**), Johnson (**panel B**), and Lobanov (**panel C**), for items with a change of more than 35%. Points represent items for which change is less than 35%. Color and arrow head indicate decrease or increase in log likelihood.

1185 sess whether normalization always improves model fit in absence of normalization, and if
 1186 normalization models perform equally well in different parts of the acoustic-phonetic space.

1187 Figure S11 indicates a general tendency for increased model performance as humans'
 1188 predictions about human behavior become stronger, even though models' improvements
 1189 are not limited to items for which humans have strong predictions. Normalization does
 1190 not, however, improve model fit across the board. Relative to no normalization, all three
 1191 accounts both increase and decrease in performance on a by-item level. The advantage
 1192 of Nearey's uniform scaling relative to no normalization seems to be driven by smaller
 1193 improvements (<35% change) on many items in Experiment 1a (proportion of items with
 1194 increase in performance = 76.4%, mean improvement in likelihood by item = 13.92 (sd
 1195 = 11.87), mean likelihood by item for items where there is *no* improvement = -12.59 (sd
 1196 = 10.81)), whereas for Experiment 1b, Nearey's uniform scaling improves substantially on
 1197 many items (proportion = 93.2, mean improvement in likelihood by item = 19.32 (sd =
 1198 15.38), mean likelihood by item for items where there is *no* improvement = -7.37 (sd =
 1199 4.62)). Johnson follows the same pattern for Experiment 1a only (proportion = 80.6%,
 1200 mean improvement in likelihood by item = 13.16 (sd = 11.36), mean likelihood by item
 1201 for items where there is *no* improvement = -10.84 (sd = 9.91)), while for Experiment 1b,
 1202 improvements are less pronounced (proportion = 75.3%, mean improvement in likelihood by
 1203 item = 11.9 (sd = 7.74), mean likelihood by item for items where there is *no* improvement
 1204 = -6.65 (sd = 4.47)). Lobanov seems to follow the same pattern as Nearey (for Experiment
 1205 1a, proportion = 73.6%, mean improvement in likelihood by item = 12.67 (sd = 12.08),
 1206 mean likelihood by item for items where there is *no* improvement = -11.06 (sd = 9.82); for
 1207 Experiment 1b, proportion = 87%, mean improvement in likelihood by item = 20.59 (sd
 1208 = 17.07), mean likelihood by item for items where there is *no* improvement = -4.92 (sd =
 1209 4.62)).

1210 To explore whether differences in model performance are related to where in the acoustic-
 1211 phonetic space items are located, we plot the likelihood of the unnormalized model in the
 1212 acoustic-phonetic space, along with likelihood differences between the best-performing mod-
 1213 els (see Figure S12).

1214 Figure S12 suggests that normalization does not improve things universally across the
 1215 acoustic-phonetic space. Overall, model performance is better for items for which human
 1216 predictions are stronger, that is, models perform better in parts of the acoustic space where
 1217 humans can easier predict human behavior (Figure S12, *Panel A*). To the extent that this
 1218 is not the case, it seems that normalization in general can adjust for this, improving model
 1219 performance on many tokens where the maximum performance is high but the unnormalized

Comparing normalization against perception

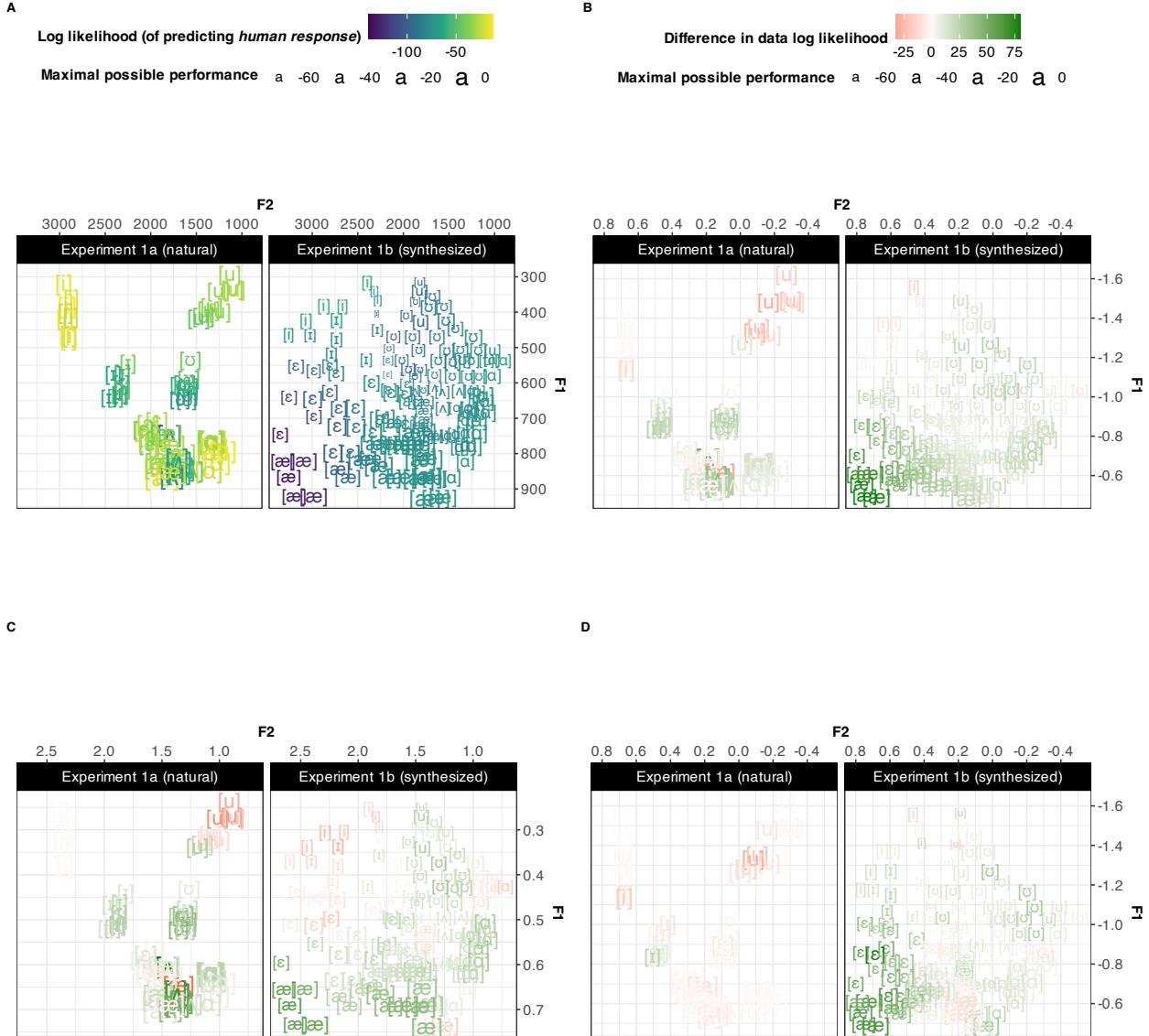


FIG. S12. In which part of the acoustic-phonetic space does normalization fail to improve fit against human responses? For each test location, the vowel label indicates the most frequent response provided by participants. Size of vowel label relates model performance to maximum performance (predicting human responses from human responses). **Panel A** shows the likelihood of the unnormalized model in predicting human responses to both experiments. **Panels B-D** shows difference in likelihood between models, Nearey's uniform scaling vs. no normalization (**panel B**), Johnson vs. no normalization (**panel C**), Nearey's uniform scaling vs. Johnson (**panel D**).

model's predictions are low, e.g., in the left bottom and center part of the acoustic space (*Panels B-C*). There is overall less improvement in Experiment 1a, presumably because models are already performing well predicting human behavior in the first experiment. Both Nearey's uniform scaling and Johnson clearly perform worse relative to the unnormalized model in the upper right part of the space, more specifically for the [u] category (*Panels B-C; left*), which could indicate that models are overly categorical in a part of the space where humans are less categorical. Possible reasons to this, could be 1) the stimuli sounding more like a neighbouring category to many listeners, or 2) potential effects of orthography, making humans less inclined to select the [u] category. The potential effect of the infrequent non-word response option *who'd* could have been checked against the synthesized stimuli in Experiment 1b. If there was indeed an effect of orthography, we should have observed a better model fit and larger between-account differences in predictions in this part of the acoustic space. Unfortunately, we under-sampled that part, which is an important caveat for Experiment 1b. For the items closest to the area in question, participants however often responded *hood*, which might indicate that items in this part of the space for this talker overall sounded more like *hood* and not *who'd* for many listeners (c.f., discussion on listeners' dialect templates in Section II B).

Comparing the two best-performing models across experiments (*Panel D*), there are no evident patterns of improvement in one model relative to the other. In Experiment 1a, Johnson provided the best fit to listeners' responses and appears to improve the fit relative to Nearey across the entire space (with the exception of one [i] token). For Experiment 1b, Nearey overall improves the fit relative to Johnson, with the exception of some locations in the mid part of the phonetic space (including high, center and low vowels).

1243 C. Results for F1-F2 (subsets of Experiments 1a and 1b)

1244 To evaluate two potential concerns with our stimuli, we decided to compare the 20 normalization accounts against a subset of the data from Experiment 1a and 1b. For Experiment 1245 1a, we excluded listeners' responses to the two *hVd* stimuli that differed in phonological 1246 context from all other words: *odd* and *hut*. For Experiment 1b, we excluded responses to 1247 stimuli that were presumed physiologically implausible under the assumption of a single 1248 talker (all stimuli below the diagonal dashed line in Figure 4).

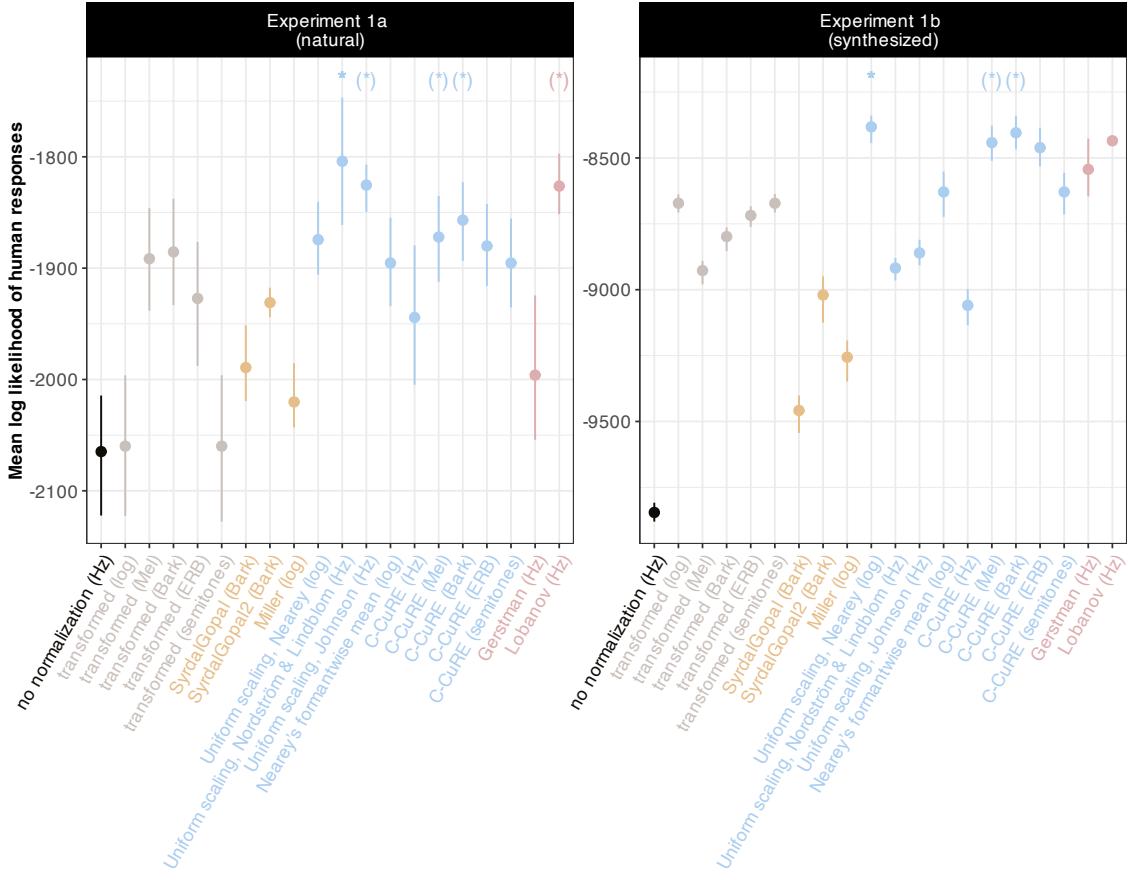


FIG. S13. Results of model fit to subset data. Pointranges indicate mean and 95% bootstrapped CIs of the log-likelihood summarized over the five training sets (higher is better). Accounts that fit listeners' responses to an extent that is statistically indistinguishable from the best-fitting account are marked by (*).

1250 This analysis overall replicates the results from the main analysis: uniform scaling ac-
 1251 counts again provide the best fit against listeners' responses in both experiments (Figure
 1252 S13). For Experiment 1a, Nordström & Lindblom achieved the best fit (log likelihood =
 1253 -1804, SD =), while Nearey's uniform scaling again provided the best fit to Experiment 1b
 1254 (log likelihood = -8382, SD = 71).

1255 **D. Results for F1-F2 (subset of listeners sharing dialect template)**

1256 Analyses in the main paper suggested that not all listeners in Experiment 1a and 1b
 1257 shared dialect template (Section II B). To investigate the effect of excluding listeners that
 1258 likely did not use the same underlying vowel representations for categorization, we compared

1259 the 20 normalization accounts against a subset of listeners who employed the dialect template
 1260 used by the majority of participants (see lower-left of both panels in Figure 5B). This left
 1261 11 participants for Experiment 1a (61.1%) and 14 for Experiment 1b (77.8%). Under the
 1262 assumptions that 1) our model of listeners is adequate, that 2) the subset group of listeners
 1263 now share dialect template and that 3) the priors, the phonetic database, can approximate
 1264 this template, we would expect all model to increase their likelihood fit to listeners' responses
 1265 (c.f., Section IV).

1266 As expected, Figure S14 suggests that all models overall provide higher likelihood fits
 1267 against human responses in both experiments compared to the main model. To increase
 1268 comparability to the results of the main model, we scaled the log likelihood of models in the
 1269 subset data to those of the main analysis by multiplying the model log likelihoods with the
 1270 ratio of the number of observations in the main model over the number of observations in
 1271 the subset model. This suggested that the improvement in likelihood for the dialect subset
 1272 model to the original dataset was 41.1%.

1273 Replicating the results from the main analysis, uniform scaling accounts again fit listeners'
 1274 behavior well across both experiments. While Nearey's uniform scaling provided the best-fit
 1275 in Experiment 1b (log likelihood = -5591, SD = 47), Syrdal & Gopal now achieved the best fit
 1276 to Experiment 1a (log likelihood = -618, SD = 19). Only one additional account performed
 1277 within the range of the best-performing accounts: for Experiment 1a, all $ps < .0341$; for
 1278 Experiment 1b, Lobanov achieved likelihood fits statistically indistinguishable from Nearey's
 1279 uniform scaling ($p > .15$, log likelihood = -5612, SD = 30).

1280 While Nearey's uniform scaling displayed relatively stable performance across experi-
 1281 ments, Syrdal & Gopal varied drastically, achieving one of the worst fits to listeners' re-
 1282 sponses in Experiment 1b (log likelihood = -6684, SD = 75). As mentioned in Section
 1283 III B 2, a possible explanation to large fluctuations in model fits between experiments, is
 1284 that this account has been over-engineered on specific types of natural vowel productions.
 1285 Given that formant normalization is a pre-linguistic mechanism, it ought to be able to ex-
 1286 plain listeners' responses to any type of data, including data that does not follow correlations
 1287 in natural data. This would suggest that Syrdal & Gopal might not be a plausible account
 1288 of normalization.

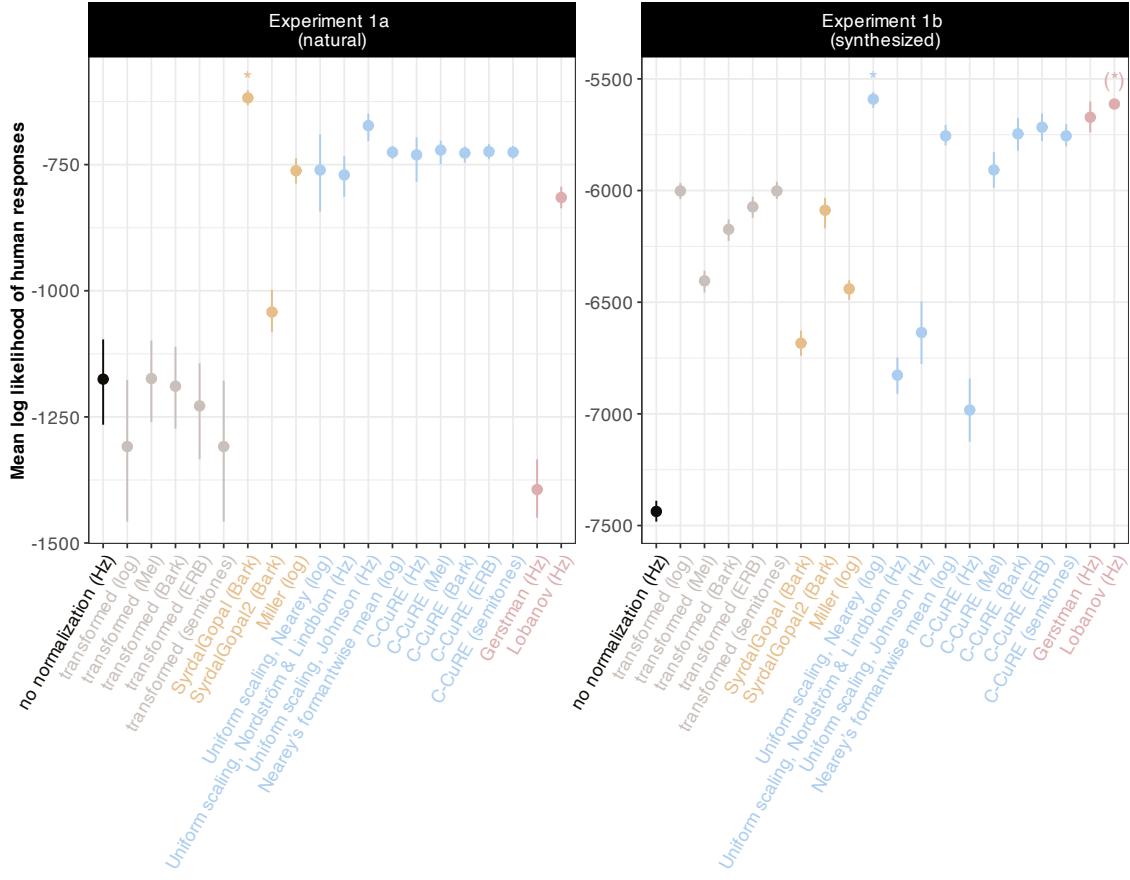


FIG. S14. Results of model fit to data excluding listeners that do not seem to share dialect template. Pointranges indicate mean and 95% bootstrapped CIs of the log-likelihood summarized over the five training sets (higher is better). Accounts that fit listeners' responses to an extent that is statistically indistinguishable from the best-fitting account are marked by (*).

1289

E. Results for F1-F3

1290

To investigate whether the inclusion of F3, a cue known to be important for vowel category distinctions, would improve the model fit to human behavior, we trained ideal observers on multivariate (F1-F2-F3) categories from the same database as in the main study. Here, we first report the results of the F1-F3 model and qualitatively compare them to the results in the main text for F1-F2. This will highlight that the results are largely similar and support the same conceptual conclusions, but there are some differences in model fit. To understand these differences better, we then also directly compare the results quantitatively to see for which accounts the inclusion of F3 improved the fit against listeners' responses and for which accounts it decreased the fit.

1299 Figure S15 summarizes how well the different accounts fit listeners' responses in Experiments 1a and 1b when assuming F1-F2-F3 multivariate category representations. Many 1300 aspects replicate the F1-F2 results reported in the main text. First, normalization sig- 1301 nificantly improved the fit relative to no normalization. Second, the same uniform scaling 1302 accounts again achieved the best fit against listeners' responses: for Experiment 1a, Johnson 1303 normalization account provided the best fit (log likelihood = -2345, SD = 23), while Nearey's 1304 uniform scaling account provided the best fit to Experiment 1b (log likelihood = -9610, SD 1305 = 76). However, we note that the inclusion of F3 does not improve the fit to listeners' 1306 responses for several accounts (compare *squares* and *circles* in Figure S15). In fact, with 1307 the exception of the raw Hertz, scale transformations, and intrinsic accounts, most extrinsic 1308 accounts seem to decrease their fit, more so in Experiment 1a than 1b. This includes the 1309 overall best-performing account in the main text, Nearey's uniform scaling, that no longer 1310 achieves a statistically indistinguishable fit from Johnson in Experiment 1a. At first blush, 1311 this is puzzling given that the model now has access to more information of a type that 1312 is broadly believed to be informative for US English vowel recognition (Hillenbrand *et al.*, 1313 1995; Nearey, 1989; Peterson and Barney, 1952). What might be underlying the apparent 1314 lack of improvement, and why does it appear as if some accounts actually achieve worse fits?

1315 One possible explanation is that listeners were only exposed to one talker in Experiment 1316 1a. According to some theories, F3 is expected to contribute to vowel recognition when 1317 there are multiple talkers, acting as a sort of normalizer for vocal tract length (Nearey, 1318 1989). In the absence of other talkers, this advantage might instead introduce noise to the 1319 models—an additional source of information that are not of use for listeners in this context. 1320 It is also possible that this particular talker has a pattern of F3 use across categories that 1321 is atypical given the other talkers in the database. This might explain why the raw Hertz 1322 model improves the fit with F3-inclusion. We checked for additional outliers along F3 for this 1323 talker, and also inspected the talker's categories in 3D-space (S4), but we could not find that 1324 outliers would be a likely explanation. To gain further knowledge into this talker's use of F3 1325 compared to other talkers in the database, we used the same models to predict the ground 1326 truth, i.e., the category the talker actually intended to produce. These models patterned 1327 with the other prediction results, again indicating that F3-inclusion did not improve model 1328 performance. We take this to suggest that the F1-F3 results is not about how our model 1329 uses F3, but rather about how this specific talker uses F3. The results might thus link back 1330 to the potential dialect differences between talkers in the database, reported in Section II B.

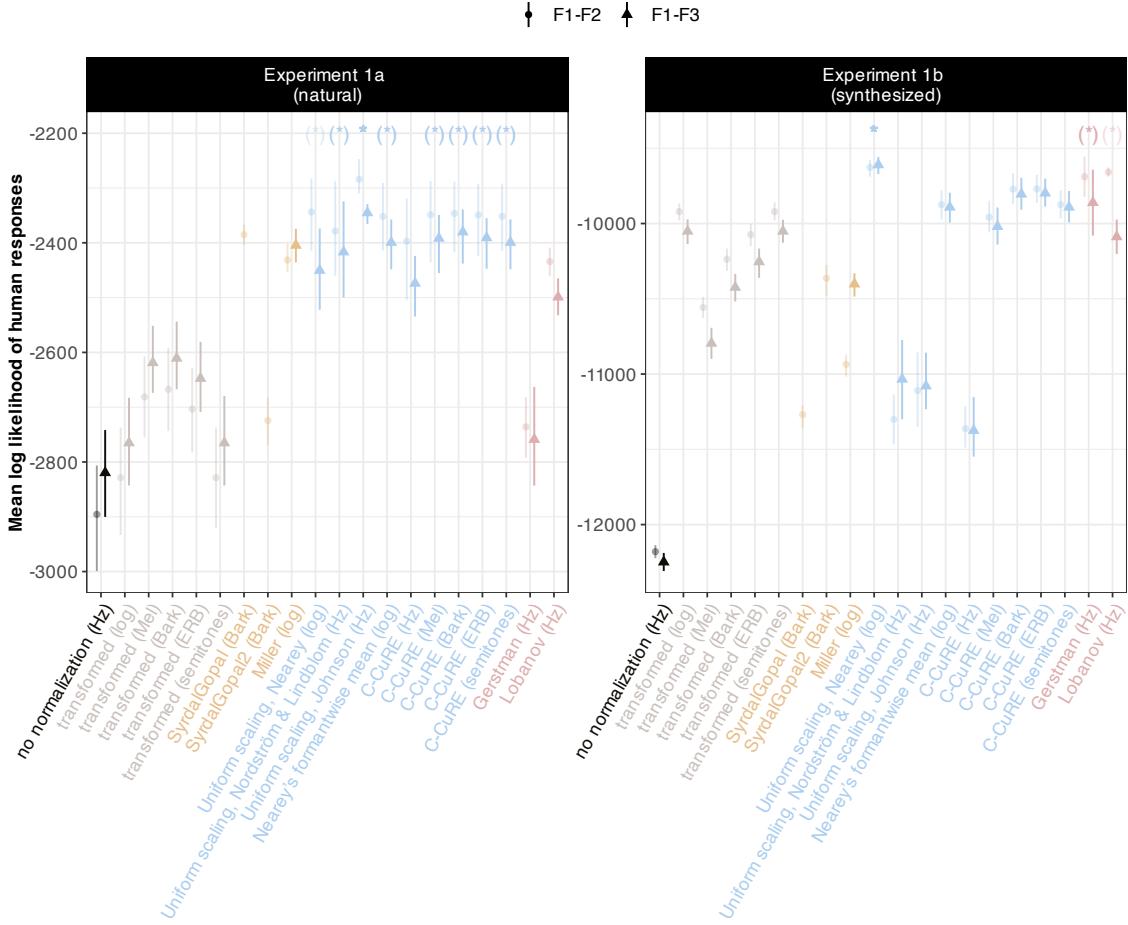


FIG. S15. Results of ideal observer models trained on F1, F2 and F3 as cues to vowel identity. As in Figure 9 in the main text, pointranges indicate mean and 95% bootstrapped CIs of the log-likelihood summarized over the five training sets (higher is better). For comparison, results from the F1-F2 models are included (more transparent circles).

1332 **F. Grid search over parameter space for F1-F2 and F1-F3**

1333 As an alternative to the quasi-Newton optimization presented in the main text, we also
 1334 conducted a grid search over the space defined by the two parameters lapse rate and noise
 1335 ratio. Figure S16 summarizes the results for a grid of lapse rates $\in 0, .02, .06, .18, .36, .72$
 1336 and noise ratios $\in 0, .3, .6, 1.25, 2.5, 5$ for Experiment 1a. For Experiment 1b (Figure S17),
 1337 the range of noise ratios explored was $\in 0, 1.5, 3, 6, 12.5, 25$.

1338 This search confirmed the pattern described in the main text. Additional grid searches
 1339 confirmed this pattern held beyond the values shown here. For all normalization accounts,
 1340 all combinations of cues, and both experiments, the goodness of fit of the ideal observers

Comparing normalization against perception

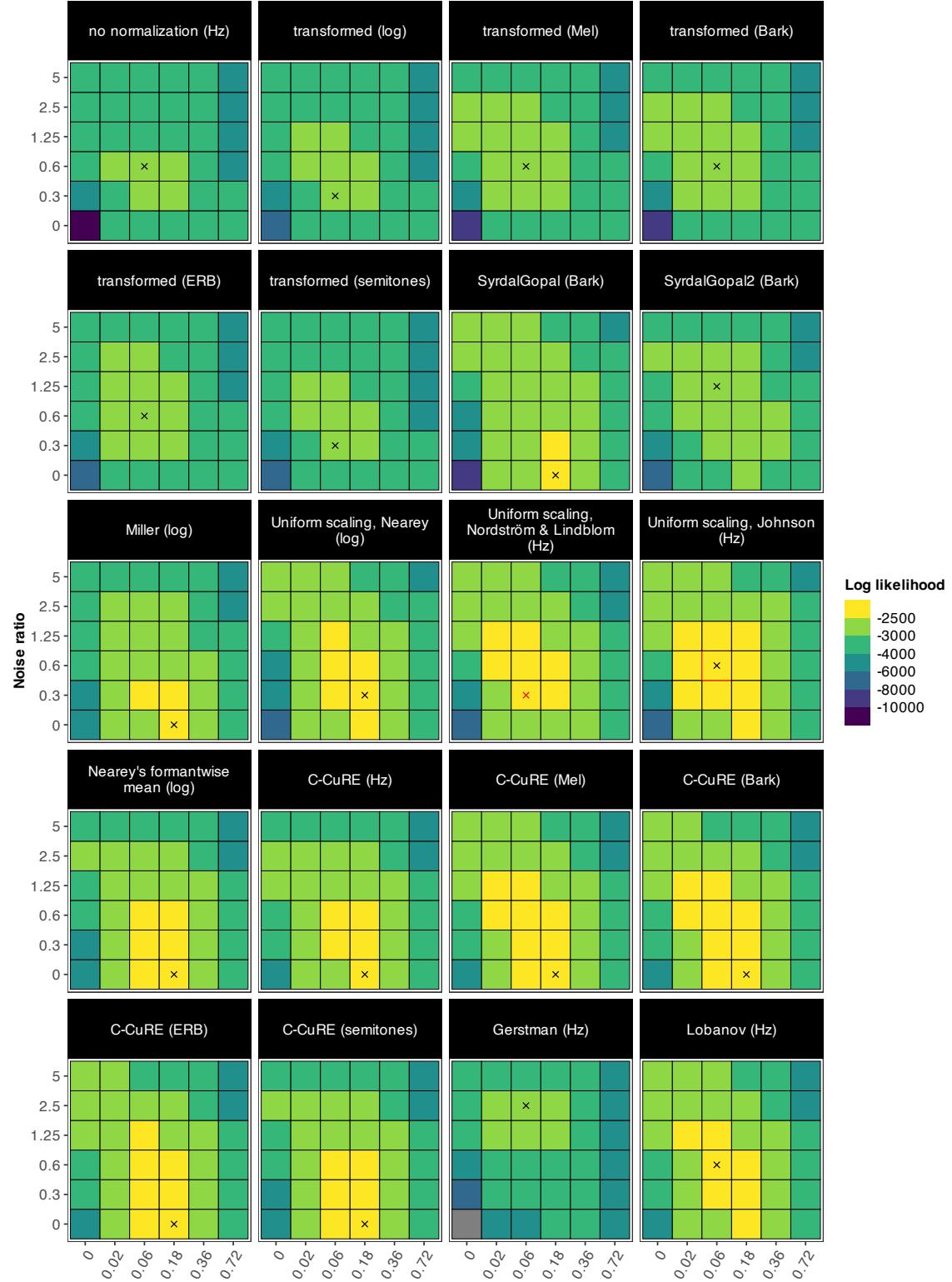


FIG. S16. Predicted likelihoods of ideal observer for human vowel responses in Experiment 1a, under different normalization accounts, different λ s and different τ^{-1} s. Likelihood is aggregated across vowels. Crosses are placed at the combination of parameters for which the maximum likelihood for an account and a cross-validation fold was found. The red cross indicates the maximum likelihood across all accounts and folds.

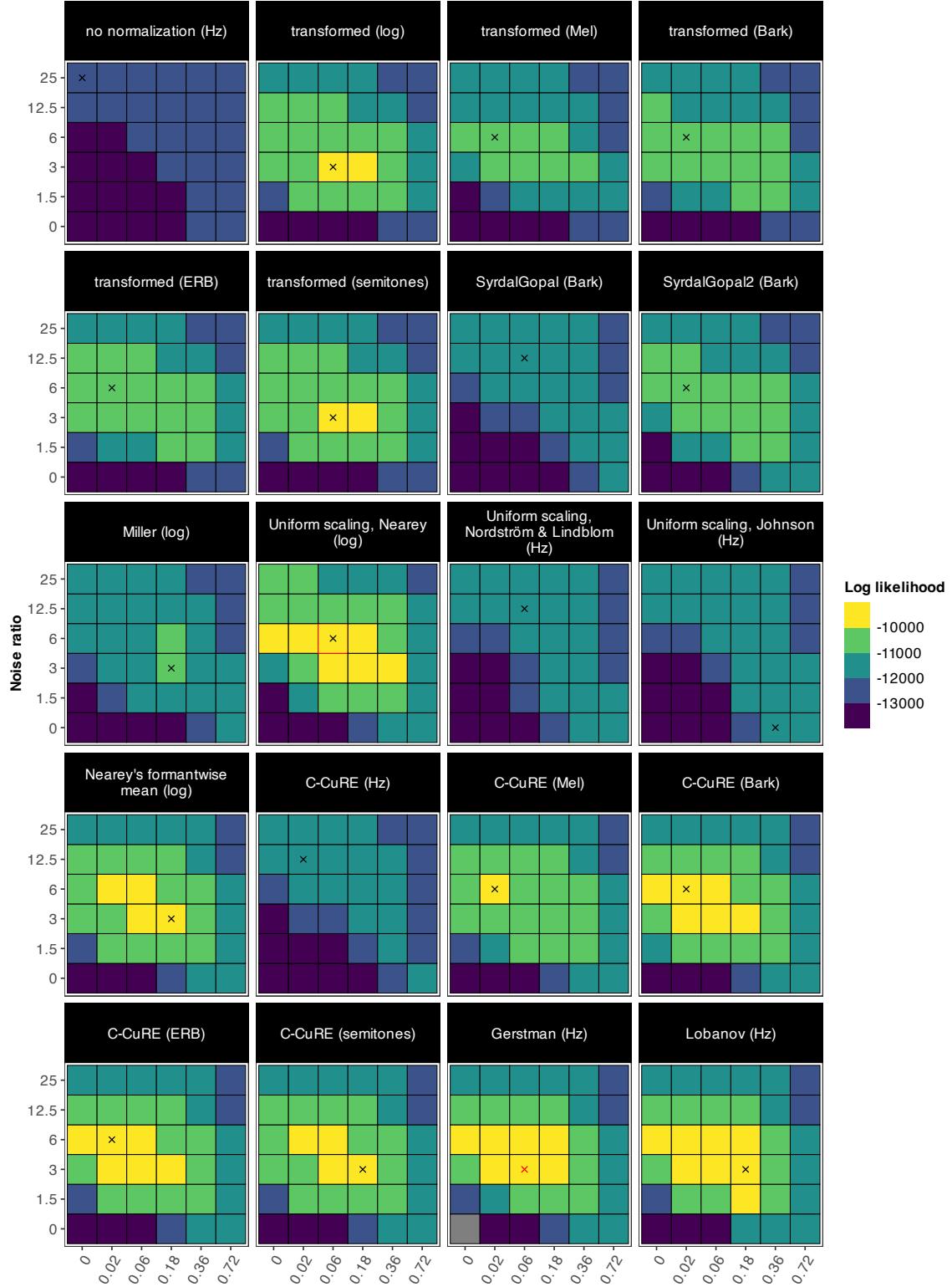


FIG. S17. Predicted likelihoods of ideal observer for human vowel responses in Experiment 1b, under different normalization accounts, different λ s and different τ^{-1} s. Likelihood is aggregated across vowels. Crosses are placed at the combination of parameters for which the maximum likelihood for an account and a cross-validation fold was found. The red cross indicates the maximum likelihood across all accounts and folds.

1341 initially improved with increasing lapse rate and increasing noise ratios, and then decreased
 1342 once lapse rates or noise ratios reached the best-fitting values (which depended on the
 1343 combination of normalization account, cues, and experiment). It further indicated that
 1344 Nearey’s uniform scaling, together with the other uniform scaling accounts and some of the
 1345 C-CuRE accounts (Experiment 1a) and Gerstman (Experiment 1b), improved faster and
 1346 performed consistently well for a good range of parameters, even for high τ^{-1} . Many of the
 1347 other models were less consistent and only performed well for a smaller range of estimates.

1348 **§4. REFERENCES**

1349 ¹Normalization does not necessarily imply that *only* talker-normalized auditory percepts are available to
 1350 subsequent processing. There is ample evidence that subcategorical information can enter listeners’ repre-
 1351 sentations of sound categories (e.g., [Hay et al., 2017, 2019](#); [Johnson et al., 1999](#); [McGowan, 2015](#); [Walker and Hay, 2011](#)), in line with episodic ([Goldinger, 1996](#)) and exemplar theory of speech perception ([Johnson, 1997](#); [Sumner, 2011](#)).

1354 ²Under uniform scaling accounts, listeners essentially ‘slide’ the center of their category representations
 1355 (e.g, the ‘template’ of vowel categories for a given dialect) along a single line in formant space, with Ψ
 1356 determining the target of this sliding. Later extensions of this account maintain its memory parsimony but
 1357 increased its inference complexity by allowing both intrinsic (the current F0) and extrinsic information (the
 1358 talker’s single mean of log-transformed formants) to influence the inference of Ψ ([Nearey and Assmann, 2007](#)). We return to this extension in the general discussion.

1360 ³We use Johnson’s (2020) implementation of [Nordström and Lindblom \(1975\)](#). We group both [Nordström](#)
 1361 and [Lindblom \(1975\)](#) and [Johnson \(2020\)](#) with the centering accounts, as they are essentially variants of
 1362 uniform scaling, differing in their estimation of Ψ . We also include both versions of Syrdal & Gopal’s
 1363 Bark-distance model. The two versions differ only in their normalization of F2, and have not previously
 1364 been compared against human perception.

1365 ⁴[Shannon \(1948\)](#) response entropy is defined as $H(x) = -\sum_{i=1}^n P(x_i) \log P(x_i)$. The maximum possible
 1366 response entropy for an 8-way response choice is 3 bits, which means that all eight vowels are responded
 1367 equally often. The minimum response entropy = 0 bits, which means that the same vowel is responded all
 1368 the time.

1369 ⁵Note that participants in Experiment 1a exhibited high agreement on [ʌ], [æ], and [ɑ], despite the close
 1370 proximity between, and partial overlap of, these vowels in F1-F2 space. To understand this pattern, it is

1371 important to keep in mind that the recordings for [ʌ] and [ɑ] differed from the recordings for other stimuli
 1372 in their word onset (“odd” for [ɑ]) or offset (“hut” for [ʌ]).

1373 ⁶[u] has been undergoing changes in many varieties of US English. Whereas the talker in Experiment 1a
 1374 produces [u] with low F1 and F2 (high and back), other L1 talkers of US English produce this vowel
 1375 considerably more forward (higher F2).

1376 ⁷For Gaussian noise and Gaussian category likelihoods, the resulting noise-convolved likelihood is a Gaussian
 1377 with variance equal to the sum of the noise and category variances (Kronrod *et al.*, 2016).

1378 ⁸We intentionally did *not* split the data within talkers since normalization accounts are meant to make
 1379 speech perception robust to cross-talker variability. Further, splitting the data by speaker rather than
 1380 by vowel category avoids the potential for biases in the normalization parameter estimates for different
 1381 speakers in the case of missing or unbalanced tokens across vowel categories, see (Barreda and Nearey,
 1382 2018). Additional analyses not reported here confirmed that the same results are obtained when splits are
 1383 performed within talkers and within vowels (except that this lead to smaller CIs, and thus *more* significant
 1384 differences, in Figure 9). These analyses can be replicated by downloading the R markdown document this
 1385 article is based on from our OSF (see comments in our code).

1386 ⁹Alternatively, it would be possible to treat these parameters as DFs in the link to listeners’ responses,
 1387 and infer them from the responses in Experiments 1a and 1b (cf., Kleinschmidt and Jaeger, 2016). This
 1388 approach would afford the model with a high degree of functional flexibility, regardless of which normal-
 1389 ization approach is applied (similar to previous approaches that have employed, e.g., multinomial logistic
 1390 regression).

1391 ¹⁰This ratio is a generalization of the inverse of the “meaningful-to-noise variance ratio (τ)” used in Kronrod
 1392 *et al.* (2016). However, whereas Kronrod and colleagues committed to the simplifying assumption that
 1393 all categories have identical variance (along all formants), we allowed category variances to differ between
 1394 vowels, and between F1 and F2 (matching the empirically facts). We merely assume that the *noise* variance
 1395 is identical across all formants (in the phonetic space defined by the normalization account, e.g., log-Hz for
 1396 uniform scaling and Hz for Lobanov).

1397 ¹¹Additional analyses reported in the SI (§3C) replicated this result for subsets of Experiments 1a and 1b.
 1398 For Experiment 1a, we excluded responses to the two *hVd* stimuli that differed from the other stimuli in
 1399 the preceding (*odd*) or following phonological context (*hut*). For Experiment 1b, we excluded responses to
 1400 any stimuli that were physiologically implausible for the talker (stimuli below the diagonal dashed line in
 1401 Figure 4).

1402 ¹²We thank Xin Xie and Leslie Li for providing us with the recordings and aligned Praat textgrids.

1403

1404 Abramson, A. S., and Lisker, L. (1973). “Voice-timing perception in spanish word-initial
 1405 stops,” Journal of Phonetics 1(1), 1–8, doi: [10.1016/S0095-4470\(19\)31372-5](https://doi.org/10.1016/S0095-4470(19)31372-5).

- 1406 Adank, P., Smits, R., and van Hout, R. (2004). “A comparison of vowel normalization pro-
 1407 cedures for language variation research,” The Journal of the Acoustical Society of America
 1408 **116**(5), 3099–3107, doi: [10.1121/1.1795335](https://doi.org/10.1121/1.1795335).
- 1409 Allen, J. S., Miller, J. L., and DeSteno, D. (2003). “Individual talker differences in voice-
 1410 onset-time,” Journal of the Acoustical Society of America **113**(1), 544–552, doi: [10.1121/1.1528172](https://doi.org/10.1121/1.1528172).
- 1412 Apfelbaum, K., and McMurray, B. (2015). “Relative cue encoding in the context of sophisti-
 1413 cated models of categorization: Separating information from categorization,” Psychonomic
 1414 Bulletin and Review **22**(4), 916–943, doi: [10.3758/s13423-014-0783-2](https://doi.org/10.3758/s13423-014-0783-2).
- 1415 Assmann, P. F., Nearey, T. M., and Bharadwaj, S. (2008). “Analysis of a vowel database,”
 1416 Canadian Acoustics **36**(3), 148–149.
- 1417 Bache, S. M., and Wickham, H. (2022). *magrittr: A Forward-Pipe Operator for R*, <https://CRAN.R-project.org/package=magrittr>, r package version 2.0.3.
- 1419 Baese-Berk, M. M., Walker, K., and Bradlow, A. (2018). “Variability in speaking rate of
 1420 native and non-native speakers,” The Journal of the Acoustical Society of America **144**(3),
 1421 1717–1717, doi: [10.1121/1.5067612](https://doi.org/10.1121/1.5067612).
- 1422 Barreda, S. (2020). “Vowel normalization as perceptual constancy,” Language **96**(2), 224–
 1423 254, doi: [10.1353/lan.2020.0018](https://doi.org/10.1353/lan.2020.0018).
- 1424 Barreda, S. (2021). “Perceptual validation of vowel normalization methods for vari-
 1425 ationist research,” Language Variation and Change **33**(1), 27–53, doi: [10.1017/S0954394521000016](https://doi.org/10.1017/S0954394521000016).
- 1427 Barreda, S. (2023). “phontools: Functions for phonetics in r” R package version 0.2-2.2.
- 1428 Barreda, S., and Nearey, T. M. (2012). “The direct and indirect roles of fundamental fre-
 1429 quency in vowel perception,” The Journal of the Acoustical Society of America **131**(1),
 1430 466–477, doi: [10.1121/1.3662068](https://doi.org/10.1121/1.3662068).
- 1431 Barreda, S., and Nearey, T. M. (2018). “A regression approach to vowel normalization for
 1432 missing and unbalanced data,” The Journal of the Acoustical Society of America **144**(1),
 1433 500–520, doi: [10.1121/1.5047742](https://doi.org/10.1121/1.5047742).
- 1434 Bengtsson, H. (2021). “A unifying framework for parallel and distributed processing in r us-
 1435 ing futures,” The R Journal **13**(2), 208–227, <https://doi.org/10.32614/RJ-2021-048>,
 1436 doi: [10.32614/RJ-2021-048](https://doi.org/10.32614/RJ-2021-048).
- 1437 Bladon, A., Henton, C., and Pickering, J. (1984). “Towards an auditory theory of speaker
 1438 normalization,” Language and Communication **4**, 59–69.

- 1439 Boersma, P., and Weenink, D. (2022). “Praat: Doing phonetics by computer [Computer
 1440 program]” .
- 1441 Buz, E., and Jaeger, T. F. (2016). “The (in) dependence of articulation and lexical planning
 1442 during isolated word production,” *Language, Cognition and Neuroscience* **31**(3), 404–424.
- 1443 Byrd, R. H., Lu, P., Nocedal, J., and Zhu, C. (1995). “A limited memory algorithm for
 1444 bound constrained optimization,” *SIAM Journal on Scientific Computing* **16**(5), 1190–
 1445 1208, doi: [10.1137/0916069](https://doi.org/10.1137/0916069).
- 1446 Bürkner, P.-C. (2017). “brms: An R package for Bayesian multilevel models using Stan,”
 1447 *Journal of Statistical Software* **80**(1), 1–28, doi: [10.18637/jss.v080.i01](https://doi.org/10.18637/jss.v080.i01).
- 1448 Bürkner, P.-C. (2018). “Advanced Bayesian multilevel modeling with the R package brms,”
 1449 *The R Journal* **10**(1), 395–411, doi: [10.32614/RJ-2018-017](https://doi.org/10.32614/RJ-2018-017).
- 1450 Bürkner, P.-C. (2021). “Bayesian item response modeling in R with brms and Stan,” *Journal*
 1451 *of Statistical Software* **100**(5), 1–54, doi: [10.18637/jss.v100.i05](https://doi.org/10.18637/jss.v100.i05).
- 1452 Campitelli, E. (2024). *ggnewscale: Multiple Fill and Colour Scales in 'ggplot2'*, <https://CRAN.R-project.org/package=ggnewscale>, r package version 0.4.10.
- 1453 Carpenter, G. A., and Govindarajan, K. K. (1993). “Neural Network and Nearest Neighbor
 1454 Comparison of Speaker Normalization Methods for Vowel Recognition,” in *ICANN '93*,
 1455 edited by S. Gielen and B. Kappen (Springer London, London), pp. 412–415, doi: [10.1007/978-1-4471-2063-6_98](https://doi.org/10.1007/978-1-4471-2063-6_98).
- 1456 Chládková, K., Podlipský, V. J., and Chionidou, A. (2017). “Perceptual adaptation of
 1457 vowels generalizes across the phonology and does not require local context.,” *Journal of*
 1458 *Experimental Psychology: Human Perception and Performance* **43**(2), 414.
- 1459 Clayards, M., Tanenhaus, M. K., Aslin, R. N., and Jacobs, R. A. (2008). “Perception of
 1460 speech reflects optimal use of probabilistic speech cues,” *Cognition* **108**(3), 804–809, doi:
 1461 [10.1016/j.cognition.2008.04.004](https://doi.org/10.1016/j.cognition.2008.04.004).
- 1462 Colby, S., Clayards, M., and Baum, S. (2018). “The role of lexical status and individual dif-
 1463 ferences for perceptual learning in younger and older adults,” *Journal of Speech, Language,*
 1464 and Hearing Research **61**(8), 1855–1874.
- 1465 Cole, J., Linebaugh, G., Munson, C., and McMurray, B. (2010). “Unmasking the acous-
 1466 tic effects of vowel-to-vowel coarticulation: A statistical modeling approach,” *Journal of*
 1467 *Phonetics* **38**(2), 167–184, doi: [10.1016/j.wocn.2009.08.004](https://doi.org/10.1016/j.wocn.2009.08.004).
- 1468 Crinnion, A. M., Malmskog, B., and Toscano, J. C. (2020). “A graph-theoretic approach to
 1469 identifying acoustic cues for speech sound categorization,” *Psychonomic Bulletin & Review*

- 1472 **27**(6), 1104–1125, doi: [10.3758/s13423-020-01748-1](https://doi.org/10.3758/s13423-020-01748-1).
- 1473 Csárdi, G., Hester, J., Wickham, H., Chang, W., Morgan, M., and Tenenbaum, D. (2024).
 1474 *remotes: R Package Installation from Remote Repositories, Including 'GitHub'*, <https://CRAN.R-project.org/package=remotes>, r package version 2.5.0.
- 1475
- 1476 Disner, S. F. (1980). “Evaluation of vowel normalization procedures,” The Journal of the
 1477 Acoustical Society of America **67**(1), 253–261, doi: [10.1121/1.383734](https://doi.org/10.1121/1.383734).
- 1478 Eaves Jr, B. S., Feldman, N. H., Griffiths, T. L., and Shafto, P. (2016). “Infant-directed
 1479 speech is consistent with teaching..,” Psychological Review **123**(6), 758.
- 1480 Eddelbuettel, D., Francois, R., Allaire, J., Ushey, K., Kou, Q., Russell, N., Ucar, I., Bates,
 1481 D., and Chambers, J. (2024). *Rcpp: Seamless R and C++ Integration*, <https://CRAN.R-project.org/package=Rcpp>, r package version 1.0.12.
- 1482
- 1483 Escudero, P., and Bion, R. A. H. (2007). “Modeling vowel normalization and sound per-
 1484 ception as sequential processes,” ICPHS **XVI**, 1413–1416.
- 1485 Fant, G. (1975). “Non-uniform vowel normalization,” STL-QPSR **16**(2-3), 001–019.
- 1486 Fant, G., Kruckenberg, A., Gustafson, K., and Liljencrants, J. (2002). “A New Approach
 1487 to Intonation Analysis and Synthesis of Swedish,” ISCA 2002 283–286.
- 1488 Fastl, H., and Zwicker, E. (2007). *Psychoacoustics* (Springer, Berlin, Heidelberg).
- 1489 Feldman, N. H., Griffiths, T. L., and Morgan, J. L. (2009). “The influence of categories
 1490 on perception: Explaining the perceptual magnet effect as optimal statistical inference,”
 1491 Psychological Review **116**(4), 752–782.
- 1492 Flemming, E. (2010). “Modeling listeners: Comments on pluymakers et al. and scarbor-
 1493 ough,” in *Laboratory Phonology*, edited by C. Fougeron, B. Kühnert, M. D’Imperio, and
 1494 N. Vallée, **10**, pp. 587–606.
- 1495 Gabry, J., and Češnovar, R. (2022). *cmdstanr: R Interface to 'CmdStan'*, <https://mc-stan.org/cmdstanr/>, r package version 0.5.3, <https://discourse.mc-stan.org>.
- 1496
- 1497 Gahl, S., Yao, Y., and Johnson, K. (2012). “Why reduce? phonological neighborhood
 1498 density and phonetic reduction in spontaneous speech,” Journal of Memory and Language
 1499 **66**(4), 789–806.
- 1500 Gerstman, L. (1968). “Classification of self-normalized vowels,” IEEE Transactions on Au-
 1501 dio and Electroacoustics **16**(1), 78–80, doi: [10.1109/TAU.1968.1161953](https://doi.org/10.1109/TAU.1968.1161953).
- 1502 Glasberg, B. R., and Moore, B. C. J. (1990). “Derivation of auditory filter shapes from
 1503 notched-noise data,” Hearing Research **47**(1), 103–138, doi: [10.1016/0378-5955\(90\)90170-T](https://doi.org/10.1016/0378-5955(90)90170-T).
- 1504

- 1505 Goldinger, S. D. (1996). "Words and voices: Episodic traces in spoken word identification
 1506 and recognition memory," *Journal of Experimental Psychology: Learning Memory and*
 1507 *Cognition* **22**(5), 1166–1183.
- 1508 Greenwood, D. D. (1997). "The mel scale's disqualifying bias and a consistency of pitch-
 1509 difference equisections in 1956 with equal cochlear distances and equal frequency ra-
 1510 tios," *Hearing Research* **103**(1), 199–224, <https://www.sciencedirect.com/science/article/pii/S037859559600175X>, doi: [https://doi.org/10.1016/S0378-5955\(96\)00175-X](https://doi.org/10.1016/S0378-5955(96)00175-X).
- 1513 Gromlund, G., and Wickham, H. (2011). "Dates and times made easy with lubridate,"
 1514 *Journal of Statistical Software* **40**(3), 1–25, <https://www.jstatsoft.org/v40/i03/>.
- 1515 Hall, K. C., Hume, E., Jaeger, T. F., and Wedel, A. (2018). "The role of predictability in
 1516 shaping phonological patterns," *Linguistics Vanguard* **4**(s2), 20170027.
- 1517 Hay, J., Podlubny, R., Drager, K., and McAuliffe, M. (2017). "Car-talk: Location-specific
 1518 speech production and perception," *Journal of Phonetics* **65**, 94–109, doi: [10.1016/j.wocn.2017.06.005](https://doi.org/10.1016/j.wocn.2017.06.005).
- 1520 Hay, J., Walker, A., Sanchez, K., and Thompson, K. (2019). "Abstract social categories
 1521 facilitate access to socially skewed words," *PLoS ONE* **14**(2), 1–29, doi: [10.1371/journal.pone.0210793](https://doi.org/10.1371/journal.pone.0210793).
- 1523 Henry, L., and Wickham, H. (2024). *rlang: Functions for Base Types and Core R and*
 1524 *'Tidyverse' Features*, <https://CRAN.R-project.org/package=rlang>, r package version
 1525 1.1.4.
- 1526 Hillenbrand, J. M., Getty, L. A., Clark, M. J., and Wheeler, K. (1995). "Acoustic charac-
 1527 teristics of american english vowels," *Journal of the Acoustical Society of America* **97**(5),
 1528 3099–3111.
- 1529 Hillenbrand, J. M., and Nearey, T. M. (1999). "Identification of resynthesized /hvd/ utter-
 1530 ances: Effects of formant contour," *Journal of the Acoustical Society of America* **105**(6),
 1531 3509–3523, doi: [10.1121/1.424676](https://doi.org/10.1121/1.424676).
- 1532 Hindle, D. (1978). "Approaches to Vowel Normalization in the Study of Natural Speech,"
 1533 in *Linguistic Variation: Models and Methods*, edited by D. Sankoff (Academic Press, New
 1534 York), pp. 161–171.
- 1535 Jaeger, T. F. (2024). *MVBeliefUpdatr: Fitting, Summarizing, and Visualizing of*
 1536 *Multivariate Gaussian Ideal Observers and Adaptors*, <https://github.com/hlplab/MVBeliefUpdatr>, r package version 0.0.1.0010.

- 1538 Johnson, K. (1997). "Speech perception without speaker normalization," in *Talker Variability in Speech Processing*, edited by K. Johnson and W. Mullennix (CA: Academic Press, San Diego), pp. 146–165.
- 1539
- 1540
- 1541 Johnson, K. (2020). "The Δf method of vocal tract length normalization for vowels," *Laboratory Phonology* 11(1), doi: [10.5334/labphon.196](https://doi.org/10.5334/labphon.196).
- 1542
- 1543 Johnson, K., and Sjerps, M. J. (2021). "Speaker normalization in speech perception," in *The Handbook of Speech Perception*, edited by J. S. Pardo, L. C. Nygaard, R. E. Remez, and D. B. Pisoni (John Wiley & Sons, Inc), Chap. 6, pp. 145–176, doi: [10.1002/9781119184096.ch6](https://doi.org/10.1002/9781119184096.ch6).
- 1544
- 1545
- 1546
- 1547 Johnson, K., Strand, E. A., and D'Imperio, M. (1999). "Auditory–visual integration of talker gender in vowel perception," *Journal of Phonetics* 27(4), 359–384, <https://www.sciencedirect.com/science/article/pii/S0095447099901006>, doi: <https://doi.org/10.1006/jpho.1999.0100>.
- 1548
- 1549
- 1550
- 1551 Joos, M. (1948). "Acoustic Phonetics," *Language* 24(2), 5–136, doi: [10.2307/522229](https://doi.org/10.2307/522229).
- 1552 Kay, M. (2023). *tidybayes: Tidy Data and Geoms for Bayesian Models*, <http://mjskay.github.io/tidybayes/>, doi: [10.5281/zenodo.1308151](https://doi.org/10.5281/zenodo.1308151), r package version 3.0.6.
- 1553
- 1554 Kleinschmidt, D. (2020). "What constrains distributional learning in adults?," .
- 1555 Kleinschmidt, D., and Jaeger, T. F. (2015). "Robust speech perception: Recognize the familiar, generalize to the similar, and adapt to the novel," *Psychological Review* 122(2), 148–203, doi: [10.1037/a0038695](https://doi.org/10.1037/a0038695).
- 1556
- 1557
- 1558 Kleinschmidt, D., and Jaeger, T. F. (2016). "What do you expect from an unfamiliar talker?," Proceedings of the 38th Annual Meeting of the Cognitive Science Society, CogSci 2016 2351–2356.
- 1559
- 1560
- 1561 Kleinschmidt, D., Liu, L., Bushong, W., Burchill, Z., Xie, X., Tan, M., Karboga, G., and Jaeger, F. (2021). "JSEXP" <https://github.com/hplab/JSEXP>.
- 1562
- 1563 Kronrod, Y., Coppess, E., and Feldman, N. H. (2016). "A unified model of categorical effects in consonant and vowel perception," *Psychological Bulletin and Review* 1681–1712, doi: [10.3758/s13423-016-1049-y](https://doi.org/10.3758/s13423-016-1049-y).
- 1564
- 1565
- 1566 Kuhl, P. K., Andruski, J. E., Chistovich, I. A., Chistovich, L. A., Kozhevnikova, E. V., Ryskina, V. L., Stolyarova, E. I., Sundberg, U., and Lacerda, F. (1997). "Cross-language analysis of phonetic units in language addressed to infants," *Science* 277(5326), 684–686.
- 1567
- 1568
- 1569 Labov, W., Ash, S., and Boberg, C. (2006). *The atlas of North American English: Phonetics, phonology, and sound change* (De Gruyter Mouton, Berlin • New York).
- 1570

- 1571 Ladefoged, P., and Broadbent, D. E. (1957). “Information conveyed by vowels,” Journal of
 1572 the Acoustical Society of America **29**, 98–104.
- 1573 Lee, C.-Y. (2009). “Identifying isolated, multispeaker mandarin tones from brief acoustic
 1574 input: A perceptual and acoustic study,” The Journal of the Acoustical Society of America
 1575 **125**(2), 0001–4966, doi: [10.1121/1.3050322](https://doi.org/10.1121/1.3050322).
- 1576 Liao, Y. (2019). *linguisticsdown: Easy Linguistics Document Writing with R Markdown*,
 1577 <https://CRAN.R-project.org/package=linguisticsdown>, r package version 1.2.0.
- 1578 Liberman, A. M., Cooper, F. S., Shankweiler, D. P., and Studdert-Kennedy, M. (1967).
 1579 “Perception of the speech code,” Psychological review **74**(6), 431–461.
- 1580 Lindblom, B. (1986). “Phonetic universals in vowel systems,” Experimental phonology 13–
 1581 44.
- 1582 Lindblom, B. (1990). “Explaining phonetic variation: A sketch of the H&H theory,” Speech
 1583 Production and Speech Modeling 403–439.
- 1584 Lobanov, B. M. (1971). “Classification of Russian vowels spoken by different speakers,” The
 1585 Journal of the Acoustical Society of America **49**(2B), 606–608, doi: [10.1121/1.1912396](https://doi.org/10.1121/1.1912396).
- 1586 Luce, P. A., and Pisoni, D. B. (1998). “Recognizing spoken words: The neighborhood acti-
 1587 vation model,” Ear and Hearing **19**(1), 1–36, doi: [10.1097/00003446-199802000-00001](https://doi.org/10.1097/00003446-199802000-00001).
- 1588 Luce, R. D. (1959). *Individual Choice Behavior* (John Wiley, Oxford).
- 1589 Magnuson, J. S., You, H., Luthra, S., Li, M., Nam, H., Escabí, M., Brown, K., Allopenna,
 1590 P. D., Theodore, R., Monto, N., and Rueckl, J. G. (2020). “EARSHOT: A minimal neural
 1591 network model of incremental human speech recognition,” Cognitive Science **44**(4), 1–17,
 1592 doi: [10.1111/cogs.12823](https://doi.org/10.1111/cogs.12823).
- 1593 Massaro, D. W., and Friedman, D. (1990). “Models of integration given multiple sources of
 1594 information,” Psychological Review **97**(2), 225–252, doi: [10.1037/0033-295X.97.2.225](https://doi.org/10.1037/0033-295X.97.2.225).
- 1595 McClelland, J. L., and Elman, J. L. (1986). “The TRACE model of speech perception,”
 1596 Cognitive Psychology **18**(1), 1–86.
- 1597 McCloy, D. R. (2016). *phonR: tools for phoneticians and phonologists*, r package version
 1598 1.0-7.
- 1599 McGowan, K. B. (2015). “Social expectation improves speech perception in noise,” Lan-
 1600 guage and Speech **58**(4), 502–521, doi: [10.1177/0023830914565191](https://doi.org/10.1177/0023830914565191).
- 1601 McMurray, B., and Jongman, A. (2011). “What information is necessary for speech catego-
 1602 rization?: Harnessing variability in the speech signal by integrating cues computed relative
 1603 to expectations,” Psychological Review **118**(2), 219–246, doi: [10.1037/a0022325](https://doi.org/10.1037/a0022325).

- 1604 Merzenich, M. M., Knight, P. L., and Roth, G. L. (1975). “Representation of cochlea
 1605 within primary auditory cortex in the cat,” *Journal of Neurophysiology* **38**(2), 231–249,
 1606 doi: [10.1152/jn.1975.38.2.231](https://doi.org/10.1152/jn.1975.38.2.231).
- 1607 Miller, J. D. (1989). “Auditory-perceptual interpretation of the vowel,” *The Journal of
 1608 Acoustical Society of America* **85**(5), 22.
- 1609 Moore, B. C. (2012). *An introduction to the psychology of hearing* (Brill).
- 1610 Moulin-Frier, C., Diard, J., Schwartz, J.-L., and Bessière, P. (2015). “Cosmo (“communi-
 1611 cating about objects using sensory–motor operations”): A bayesian modeling framework
 1612 for studying speech communication and the emergence of phonological systems,” *Journal
 1613 of Phonetics* **53**, 5–41.
- 1614 Murdoch, D., and Chow, E. D. (2023). *ellipse: Functions for Drawing Ellipses and Ellipse-
 1615 Like Confidence Regions*, <https://CRAN.R-project.org/package=ellipse>, r package
 1616 version 0.5.0.
- 1617 Müller, K., and Wickham, H. (2023). *tibble: Simple Data Frames*, <https://CRAN.R-project.org/package=tibble>, r package version 3.2.1.
- 1618 Nearey, T. M. (1978). Indiana University Linguistics Club *Phonetic Feature Systems for
 1619 Vowels*.
- 1620 Nearey, T. M. (1989). “Static, dynamic, and relational properties in vowel perception,” *The
 1621 Journal of the Acoustical Society of America* **85**(5), 2088–2113, doi: [10.1121/1.397861](https://doi.org/10.1121/1.397861).
- 1622 Nearey, T. M. (1990). “The segment as a unit of speech perception,” *Journal of Phonetics*
 1623 **18**(3), 347–373, doi: [10.1016/S0095-4470\(19\)30379-1](https://doi.org/10.1016/S0095-4470(19)30379-1).
- 1624 Nearey, T. M., and Assmann, P. F. (1986). “Modeling the role of inherent spectral change in
 1625 vowel identification,” *The Journal of the Acoustical Society of America* **80**(5), 1297–1308,
 1626 doi: [10.1121/1.394433](https://doi.org/10.1121/1.394433).
- 1627 Nearey, T. M., and Assmann, P. F. (2007). “Probabilistic ‘sliding template’ models for
 1628 indirect vowel normalization,” in *Experimental approaches to phonology*, edited by J.-J.
 1629 Solé, P. S. Beddor, and M. Ohala (Oxford University Press), pp. 246–270.
- 1630 Nearey, T. M., and Hogan, J. (1986). “Phonological contrast in experimental phonetics: Re-
 1631 lating distributions of measurements production data to perceptual categorization curves,”
 1632 in *Experimental Phonology*, edited by J. J. Ohala and J. Jaeger (Academic Press, New
 1633 York), pp. 141–161.
- 1634 Newman, R. S., Clouse, S. A., and Burnham, J. L. (2001). “The perceptual consequences
 1635 of within-talker variability in fricative production,” *The Journal of the Acoustical Society*

- 1637 of America **109**(3), 1181–1196.
- 1638 Nordström, P., and Lindblom, B. (1975). “A normalization procedure for vowel formant
1639 data,” Proceedings of ICPHS VIII, Leeds 212.
- 1640 Norris, D., and McQueen, J. M. (2008). “Shortlist B: A Bayesian model of continuous speech
1641 recognition.,” Psychological review **115**(2), 357–95, doi: [10.1037/0033-295X.115.2.357](https://doi.org/10.1037/0033-295X.115.2.357).
- 1642 Oganian, Y., Bhaya-Grossman, I., Johnson, K., and Chang, E. F. (2023). “Vowel and
1643 formant representation in the human auditory speech cortex,” Neuron **111**(13), 2105–2118.
- 1644 Pedersen, T. L. (2024a). *ggforce: Accelerating 'ggplot2'*, <https://CRAN.R-project.org/package=ggforce>, r package version 0.4.2.
- 1645 Pedersen, T. L. (2024b). *patchwork: The Composer of Plots*, <https://CRAN.R-project.org/package=patchwork>, r package version 1.2.0.
- 1646 Persson, A., and Jaeger, T. F. (2023). “Evaluating normalization accounts against the dense
1647 vowel space of central swedish,” Frontiers in Psychology **14**, [https://www.frontiersin.org/articles/10.3389/fpsyg.2023.1165742](https://doi.org/10.3389/fpsyg.2023.1165742), doi: [10.3389/fpsyg.2023.1165742](https://doi.org/10.3389/fpsyg.2023.1165742).
- 1648 Peterson, G. E. (1961). “Parameters of vowel quality,” Journal of Speech and Hearing
1649 Research **4**(1), 10–29, [https://pubs.asha.org/doi/abs/10.1044/jshr.0401.10](https://doi.org/10.1044/jshr.0401.10), doi:
1650 [10.1044/jshr.0401.10](https://doi.org/10.1044/jshr.0401.10).
- 1651 Peterson, G. E., and Barney, H. L. (1952). “Control methods used in a study of the vowels,”
1652 Journal of the Acoustical Society of America **24**(2), 175–184.
- 1653 Pinheiro, J., Bates, D., and R Core Team (2023). *nlme: Linear and Nonlinear Mixed Effects
1654 Models*, <https://CRAN.R-project.org/package=nlme>, r package version 3.1-164.
- 1655 R Core Team (2023). *R: A Language and Environment for Statistical Computing*, R Foun-
1656 dation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>.
- 1657 Richter, C., Feldman, N. H., Salgado, H., and Jansen, A. (2017). “Evaluating low-level
1658 speech features against human perceptual data,” Transactions of the Association for Com-
1659 putational Linguistics **5**, 425–440, doi: [10.1162/tacl_a_00071](https://doi.org/10.1162/tacl_a_00071).
- 1660 Robinson, D. (2020). *fuzzyjoin: Join Tables Together on Inexact Matching*, <https://CRAN.R-project.org/package=fuzzyjoin>, r package version 0.1.6.
- 1661 RStudio Team (2020). *RStudio: Integrated Development Environment for R*, RStudio,
1662 PBC., Boston, MA.
- 1663 Saenz, M., and Langers, D. R. (2014). “Tonotopic mapping of human auditory cortex,”
1664 Hearing Research **307**, 42–52, [https://www.sciencedirect.com/science/article/pii/S0378595513001871](https://doi.org/10.1016/j.heares.2013.07.016), doi: <https://doi.org/10.1016/j.heares.2013.07.016> hu-

- 1670 man Auditory NeuroImaging.
- 1671 Scarborough, R. (2010). "Lexical and contextual predictability: Confluent effects on the
1672 production of vowels," *Laboratory Phonology* **10**, 557–586.
- 1673 Schertz, J., and Clare, E. J. (2020). "Phonetic cue weighting in perception and production,"
1674 *Wiley Interdisciplinary Reviews: Cognitive Science* **11**(2), doi: [10.1002/wcs.1521](https://doi.org/10.1002/wcs.1521).
- 1675 Shannon, C. E. (1948). "A mathematical theory of communication," *The Bell System Technical
1676 Journal* **27**(3), 379–423, doi: [10.1002/j.1538-7305.1948.tb01338.x](https://doi.org/10.1002/j.1538-7305.1948.tb01338.x).
- 1677 Siegel, R. J. (1965). "A replication of the mel scale of pitch," *The American Journal of
1678 Psychology* **78**(4), 615–620, <http://www.jstor.org/stable/1420924>.
- 1679 Sievert, C. (2020). *Interactive Web-Based Data Visualization with R, plotly, and shiny*
1680 (Chapman and Hall/CRC), <https://plotly-r.com>.
- 1681 Sjerps, M. J., Fox, N. P., Johnson, K., and Chang, E. F. (2019). "Speaker-normalized
1682 sound representations in the human auditory cortex," *Nature Communications* **10**(1), doi:
1683 [10.1038/s41467-019-10365-z](https://doi.org/10.1038/s41467-019-10365-z).
- 1684 Skoe, E., Krizman, J., Spitzer, E. R., and Kraus, N. (2021). "Auditory cortical changes
1685 precede brainstem changes during rapid implicit learning: Evidence from human EEG,"
1686 *Frontiers in Neuroscience* 1007.
- 1687 Smith, B. L., Johnson, E., and Hayes-Harb, R. (2019). "Esl learners' intra-speaker vari-
1688 ability in producing american english tense and lax vowels," *Journal of Second Language
1689 Pronunciation* **5**(1), 139–164.
- 1690 Steriade, D. (2001). "The phonology of perceptibility effects: the P-map and its conse-
1691 quences for constraint organization" .
- 1692 Stevens, K. N. (1972). "The quantal nature of speech: Evidence from articulatory-acoustic
1693 data," in *Human communication: a unified view* (McGraHill, New York), pp. 51–66.
- 1694 Stevens, K. N. (1989). "On the quantal nature of speech," *Journal of phonetics* **17**(1-2),
1695 3–45.
- 1696 Stevens, S. S., and Volkmann, J. (1940). "The Relation of Pitch to Frequency: A Revised
1697 Scale," *The American Journal of Psychology* **53**(3), 329–353, doi: [10.2307/1417526](https://doi.org/10.2307/1417526).
- 1698 Stilp, C. (2020). "Acoustic context effects in speech perception," *WIREs Cognitive Science*
1699 **11**(1), 1–18, doi: [10.1002/wcs.1517](https://doi.org/10.1002/wcs.1517).
- 1700 Sumner, M. (2011). "The role of variation in the perception of accented speech," *Cognition*
1701 **119**(1), 131–136, doi: [10.1016/j.cognition.2010.10.018](https://doi.org/10.1016/j.cognition.2010.10.018).

- 1702 Syrdal, A. K. (1985). "Aspects of a model of the auditory representation of american english
 1703 vowels," *Speech Communication* 4(1-3), 121–135, doi: [10.1016/0167-6393\(85\)90040-8](https://doi.org/10.1016/0167-6393(85)90040-8).
- 1704 Syrdal, A. K., and Gopal, H. S. (1986). "A perceptual model of vowel recognition based on
 1705 the auditory representation of American English vowels," *The Journal of the Acoustical
 1706 Society of America* 79(4), 1086–1100, doi: [10.1121/1.393381](https://doi.org/10.1121/1.393381).
- 1707 Tan, M., and Jaeger, T. F. (2024). "Incremental adaptation to an unfamiliar talker,"
 1708 Manuscript, Stockholm University .
- 1709 Tang, C., Hamilton, L. S., and Chang, E. F. (2017). "Intonational speech prosody encod-
 1710 ing in the human auditory cortex," *Science* 357(6353), 797–801, doi: [10.1126/science.aam8577](https://doi.org/10.1126/science.aam8577).
- 1712 ten Bosch, L., Boves, L., Tucker, B., and Ernestus, M. (2015). "DIANA: towards compu-
 1713 tational modeling reaction times in lexical decision in north American English," in *Proc.
 1714 Interspeech 2015*, pp. 1576–1580, doi: [10.21437/Interspeech.2015-366](https://doi.org/10.21437/Interspeech.2015-366).
- 1715 Traunmüller, H. (1981). "Perceptual dimension of openness in vowels," *The Journal of the
 1716 Acoustical Society of America* 69(5), 1465–1475, doi: [10.1121/1.385780](https://doi.org/10.1121/1.385780).
- 1717 Traunmüller, H. (1990). "Analytical expressions for the tonotopic sensory scale," *The Jour-
 1718 nal of the Acoustical Society of America* 88(1), 97–100, doi: [10.1121/1.399849](https://doi.org/10.1121/1.399849).
- 1719 Urbanek, S., and Horner, J. (2023). *Cairo: R Graphics Device using Cairo Graphics Library
 1720 for Creating High-Quality Bitmap (PNG, JPEG, TIFF), Vector (PDF, SVG, PostScript)
 1721 and Display (X11 and Win32) Output*, <https://CRAN.R-project.org/package=Cairo>, r
 1722 package version 1.6-2.
- 1723 van den Brand, T. (2024). *ggh4x: Hacks for 'ggplot2'*, [https://CRAN.R-project.org/
 1724 package=ggh4x](https://CRAN.R-project.org/package=ggh4x), r package version 0.2.8.
- 1725 Vaughan, D., and Dancho, M. (2022). *furrr: Apply Mapping Functions in Parallel using
 1726 Futures*, <https://CRAN.R-project.org/package=furrr>, r package version 0.3.1.
- 1727 Vaughn, C., Baese-Berk, M., and Idemaru, K. (2019). "Re-examining phonetic variability
 1728 in native and non-native speech," *Phonetica* 76(5), 327–358.
- 1729 Vorperian, H. K., and Kent, R. D. (2007). "Vowel acoustic space development in children: A
 1730 synthesis of acoustic and anatomic data," *Journal of Speech, Language & Hearing Research
 1731* 50(6), 1510–1545, doi: [10.1044/1092-4388\(2007/104\)](https://doi.org/10.1044/1092-4388(2007/104)).
- 1732 Wade, T., Jongman, A., and Sereno, J. (2007). "Effects of acoustic variability in the per-
 1733 ceptual learning of non-native-accented speech sounds," *Phonetica* 64(2-3), 122–144, doi:
 1734 [10.1159/000107913](https://doi.org/10.1159/000107913).

- 1735 Walker, A., and Hay, J. (2011). “Congruence between ‘word age’ and ‘voice age’ facilitates
 1736 lexical access,” *Laboratory Phonology* **2**(1), 219–237, doi: [10.1515/labphon.2011.007](https://doi.org/10.1515/labphon.2011.007).
- 1737 Watt, D., and Fabricius, A. (2002). “Evaluation of a technique for improving the mapping
 1738 of multiple speakers’ vowel spaces in the F1 ~ F2 plane,” in *Leeds Working Papers in
 1739 Linguistics and Phonetics*, edited by D. Nelson, 9, pp. 159–173.
- 1740 Weatherholtz, K., and Jaeger, T. F. (2016). “Speech perception and generalization
 1741 across talkers and accents,” Oxford Research Encyclopedia of Linguistics doi: [10.1093/acrefore/9780199384655.013.95](https://doi.org/10.1093/acrefore/9780199384655.013.95).
- 1743 Wedel, A., Nelson, N., and Sharp, R. (2018). “The phonetic specificity of contrastive hy-
 1744 perarticulation in natural speech,” *Journal of Memory and Language* **100**, 61–88.
- 1745 Whalen, D. H. (2016). “A double-Nearey theory of vowel normalization: Approaching con-
 1746 sensus,” *The Journal of the Acoustical Society of America* **140**(4_Supplement), 3163–3164,
 1747 doi: [10.1121/1.4969932](https://doi.org/10.1121/1.4969932).
- 1748 Wichmann, F. A., and Hill, N. J. (2001). “The psychometric function: I. Fitting, sampling,
 1749 and goodness of fit,” *Perception & psychophysics* **63**(8), 1293–1313.
- 1750 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag New
 1751 York), <https://ggplot2.tidyverse.org>.
- 1752 Wickham, H. (2019). *assertthat: Easy Pre and Post Assertions*, <https://CRAN.R-project.org/package=assertthat>, r package version 0.2.1.
- 1754 Wickham, H. (2023a). *forcats: Tools for Working with Categorical Variables (Factors)*,
 1755 <https://CRAN.R-project.org/package=forcats>, r package version 1.0.0.
- 1756 Wickham, H. (2023b). *modelr: Modelling Functions that Work with the Pipe*, <https://CRAN.R-project.org/package=modelr>, r package version 0.1.11.
- 1758 Wickham, H. (2023c). *stringr: Simple, Consistent Wrappers for Common String Operations*,
 1759 <https://CRAN.R-project.org/package=stringr>, r package version 1.5.1.
- 1760 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund,
 1761 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M.,
 1762 Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., Takahashi, K., Vaughan, D.,
 1763 Wilke, C., Woo, K., and Yutani, H. (2019). “Welcome to the tidyverse,” *Journal of Open
 1764 Source Software* **4**(43), 1686, doi: [10.21105/joss.01686](https://doi.org/10.21105/joss.01686).
- 1765 Wickham, H., François, R., Henry, L., Müller, K., and Vaughan, D. (2023). *dplyr: A
 1766 Grammar of Data Manipulation*, <https://CRAN.R-project.org/package=dplyr>, r pack-
 1767 age version 1.1.4.

- 1768 Wickham, H., and Henry, L. (2023). *purrr: Functional Programming Tools*, <https://CRAN.R-project.org/package=purrr>, r package version 1.0.2.
- 1770 Wickham, H., Hester, J., and Bryan, J. (2024a). *readr: Read Rectangular Text Data*, <https://CRAN.R-project.org/package=readr>, r package version 2.1.5.
- 1772 Wickham, H., Vaughan, D., and Girlich, M. (2024b). *tidyr: Tidy Messy Data*, <https://CRAN.R-project.org/package=tidyr>, r package version 1.3.1.
- 1774 Wilke, C. O., and Wiernik, B. M. (2022). *ggtext: Improved Text Rendering Support for 'ggplot2'*, <https://CRAN.R-project.org/package=ggtext>, r package version 0.1.2.
- 1776 Winn, M. (2018). “Speech: It’s not as acoustic as you think,” *Acoustics Today* **12**(2), 43–49.
- 1777 Wood, S., Pya, and S”afken, B. (2016). “Smoothing parameter and model selection for 1778 general smooth models (with discussion),” *Journal of the American Statistical Association* **111**, 1548–1575.
- 1780 Wood, S. N. (2003). “Thin-plate regression splines,” *Journal of the Royal Statistical Society* **(B)** **65**(1), 95–114.
- 1782 Wood, S. N. (2004). “Stable and efficient multiple smoothing parameter estimation for 1783 generalized additive models,” *Journal of the American Statistical Association* **99**(467), 673–686.
- 1785 Wood, S. N. (2011). “Fast stable restricted maximum likelihood and marginal likelihood 1786 estimation of semiparametric generalized linear models,” *Journal of the Royal Statistical Society* **(B)** **73**(1), 3–36.
- 1788 Xie, X., Buxó-Lugo, A., and Kurumada, C. (2021). “Encoding and decoding of meaning 1789 through structured variability in speech prosody,” *Cognition* **211**, 1–27, doi: [10.1016/j.cognition.2021.104619](https://doi.org/10.1016/j.cognition.2021.104619).
- 1791 Xie, X., and Jaeger, T. F. (2020). “Comparing non-native and native speech: Are L2 1792 productions more variable?,” *The Journal of the Acoustical Society of America* **147**(5), 3322–3347, doi: [10.1121/10.0001141](https://doi.org/10.1121/10.0001141).
- 1794 Xie, X., Jaeger, T. F., and Kurumada, C. (2023). “What we do (not) know about the 1795 mechanisms underlying adaptive speech perception: A computational review,” *Cortex* .
- 1796 Xie, Y. (2024). *knitr: A General-Purpose Package for Dynamic Report Generation in R*, 1797 <https://yihui.org/knitr/>, r package version 1.47.
- 1798 Yuan, J., and Liberman, M. (2008). “Speaker identification on the SCOTUS corpus,” *The 1799 Journal of the Acoustical Society of America* **123**(5), 3878–3878, doi: [10.1121/1.2935783](https://doi.org/10.1121/1.2935783).

- 1800 Zahorian, S. A., and Jagharghi, A. J. (1991). “Speaker normalization of static and dynamic
1801 vowel spectral features,” The Journal of the Acoustical Society of America **90**(1), 67–75,
1802 doi: [10.1121/1.402350](https://doi.org/10.1121/1.402350).