

AUTHOR QUERY FORM

	<p>Journal: J. Acoust. Soc. Am.</p> <p>Article Number: JASA-11243</p>	<p>Please provide your responses and any corrections by annotating this PDF and uploading it according to the instructions provided in the proof notification email.</p>
---	---	--

Dear Author,

Below are the queries associated with your article; please answer all of these queries before sending the proof back to AIP. Please list figures that are to appear as color in print _____ (a fee of \$325 per figure will apply). If you do not wish to have color page charges, the figures will be color online only.

Article checklist: In order to ensure greater accuracy, please check the following and make all necessary corrections before returning your proof.

1. Is the title of your article accurate and spelled correctly?
2. Please check affiliations including spelling, completeness, and correct linking to authors.
3. Did you remember to include acknowledgment of funding, if required, and is it accurate?

Location in article	Query / Remark: click on the Q link to navigate to the appropriate spot in the proof. There, insert your comments as a PDF annotation.
AQ1	Payment of page charges for papers over 12 pages is mandatory. If you will not honor the page charges then you must reduce your paper to 12 pages.
AQ2	Please check that the author names are in the proper order and spelled correctly. Also, please ensure that each author's given and surnames have been correctly identified (given names are highlighted in red and surnames appear in blue).
AQ3	Please define ERB at first occurrence.
AQ4	Please define IPA at first occurrence.
AQ5	In the sentence beginning "This excluded five..." please check the section citation and update as necessary.
AQ6	In the sentence beginning "Compared to participants in Experiment 1a..." please check the in text citation of what appears to be the supplementary material. Please confirm this citation and carefully check all mentions of supplementary material in this article, ensuring the correct section numbers are noted.
AQ7	Please define ASR at first occurrence.
AQ8	Please define AFC at first occurrence.
AQ9	Please update reference Barreda and Jaeger, submitted, if possible. If unable to update, then this reference must be removed from the text and references list or replaced with a published reference.
AQ10	Please specify which section or subsection "in the previous section" refers to here.
AQ11	Please provide a URL and last viewed date for the URL in Boersma and Weenink (2022).
AQ12	Goldinger (1996) was not cited in text. Please provide an in text citation for the reference Goldinger (1996) or remove from your references list.
AQ13	Please provide a last viewed date for the URL in reference Jaeger (2024).
AQ14	Please provide a last viewed date for the URL in reference Kleinschmidt <i>et al.</i> (2021).
AQ15	Please provide a last viewed date for the URL in R Core Team (2024).
AQ16	Sumner (2011) was not cited in text. Please provide and in text citation for Sumner (2011) or remove from your references list.
AQ17	We were unable to locate a digital object identifier (doi) for Ref(s). Assmann <i>et al.</i> (2008), Fant (1975), Fant <i>et al.</i> (2002), and Winn (2018). Please verify and correct author names and journal details (journal title, volume number, page number, and year) as needed and provide the doi. If a doi is not available, no other information is needed from you. For additional information on doi's, please select this link: http://www.doi.org/ .

AQ18

Please verify the single-page article in Ref(s). Chladkova *et al.* (2017), Eaves *et al.* (2016), and Johnson (2020); otherwise, please provide the full page range.

AQ19

Please specify which section or subsection “in the next section” refers to here.

Thank you for your assistance.

Author Proof



AQ1 1 Comparing accounts of formant normalization against US 2 English listeners' vowel perception

AQ2 8 **Anna Persson**,^{1,a)} **Santiago Barreda**,² and **T. Florian Jaeger**,³

7 ¹*Swedish Language and Multilingualism, Stockholm University, Stockholm, SE-106 91, Sweden*

8 ²*Linguistics, University of California, Davis, California 95616, USA*

9 ³*Brain and Cognitive Sciences, Data Science, University of Rochester, Rochester, New York 14627, USA*

ABSTRACT:

10 Human speech recognition tends to be robust, despite substantial cross-talker variability. Believed to be critical to
11 this ability are auditory normalization mechanisms whereby listeners adapt to individual differences in vocal tract
12 physiology. This study investigates the computations involved in such normalization. Two 8-way alternative forced-
13 choice experiments assessed L1 listeners' categorizations across the entire US English vowel space—both for unal-
14 tered and synthesized stimuli. Listeners' responses in these experiments were compared against the predictions of 20
15 influential normalization accounts that differ starkly in the inference and memory capacities they imply for speech
16 perception. This includes variants of *estimation-free* transformations into psycho-acoustic spaces, *intrinsic* normaliza-
17 tions relative to concurrent acoustic properties, and *extrinsic* normalizations relative to talker-specific statistics.
18 Listeners' responses were best explained by extrinsic normalization, suggesting that listeners learn and store distribu-
19 tional properties of talkers' speech. Specifically, *computationally simple* (single-parameter) extrinsic normalization
20 best fit listeners' responses. This simple extrinsic normalization also clearly outperformed Lobanov normalization—
21 a computationally more complex account that remains popular in research on phonetics and phonology, sociolinguistics,
22 typology, and language acquisition. © 2025 Author(s). All article content, except where otherwise noted, is
licensed under a Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).
<https://doi.org/10.1121/10.0035476>

(Received 25 July 2024; revised 23 December 2024; accepted 7 January 2025; published online xx xx xxxx)

[Editor: Li Xu]

Pages: 1–25

23 I. INTRODUCTION

24 One of the central challenges for speech perception
25 originates in cross-talker variability: depending on the
26 talker, the same acoustic signal can encode different sound
27 categories (Allen *et al.*, 2003; Liberman *et al.*, 1967;
28 Newman *et al.*, 2001). This results in ambiguity in the map-
29 ping from acoustics to words and meanings. Research has
30 identified several mechanisms through which listeners
31 resolve this ambiguity, ranging from early perceptual pro-
32 cesses to adaptation of phonetic categories, all the way to
33 adjustments in post-linguistic decision processes (for
34 review, see Xie *et al.*, 2023). The present study focuses on
35 the first type of mechanism, early auditory processes that
36 transform and normalize the acoustic input into the percept-
37ual cues that constitute the input to linguistic processing
38 (for reviews, Barreda, 2020; Johnson and Sjerps, 2021;
39 McMurray and Jongman, 2011; Stilp, 2020; Weatherholtz
40 and Jaeger, 2016). We seek to respond, in particular, to
41 recent calls to put theories of adaptive speech perception to
42 stronger tests (Baese-Berk *et al.*, 2018; Schertz and Clare,
43 2020; Xie *et al.*, 2023).

44 Evidence for the presence of early normalization mech-
45 anisms comes from neuroimaging and neurophysiological

studies (e.g., Oganian *et al.*, 2023; Skoe *et al.*, 2021), as
46 well as research on the peripheral auditory system suggest-
47 ing automatic transformations of the acoustic signal into
48 scale-invariant spectral patterns (e.g., Patterson and Irino,
49 2014; Smith *et al.*, 2005). Neurophysiological studies have
50 further decoded the effects of talker identity from subcorti-
51 cal brain areas like the brain stem, and thus prior to the cor-
52 tical regions believed to encode linguistic categories (e.g.,
53 Sjerps *et al.*, 2019; Tang *et al.*, 2017). This includes brain
54 responses that lag the acoustic signal by as little as 20–50
55 ms (Lee, 2009), suggesting very fast and highly automatic
56 processes. While this does not mean that *only* talker-
57 normalized auditory percepts are available to subsequent
58 processing—there is now convincing evidence that subcate-
59 gorical information can enter listeners' phonetic representa-
60 tions (e.g., Hay *et al.*, 2017, 2019; Johnson *et al.*, 1999;
61 McGowan, 2015; Walker and Hay, 2011)—it does suggest
62 that normalized auditory percepts are available to subse-
63 quent processing. By removing (some) cross-talker variabil-
64 ity early during auditory processing, normalization offers an
65 elegant and effective solution that can reduce the need for
66 more complex adaptive processes further upstream
67 (Apfelbaum and McMurray, 2015; Xie *et al.*, 2023).

68 While it is relatively uncontroversial that normalization
69 contributes to robust speech perception, it is still unclear
70

^{a)}Email: anna.persson@su.se

what types of computations this implicates. We address this question for the perception of vowels, which cross-linguistically rely on peaks in the distribution of spectral energy over acoustic frequencies (formants).¹ Vowel perception has long been a focus in research on normalization (e.g., Bladon *et al.*, 1984; Fant, 1975; Gerstman, 1968; Johnson, 2020; Joos, 1948; Lobanov, 1971; Miller, 1989; Nearey, 1978; Nordström and Lindblom, 1975; Syrdal and Gopal, 1986; Traunmüller, 1981; Watt and Fabricius, 2002; Zahorian and Jagharghi, 1991; for review, see Barreda, 2020), with some reviews citing over 100 competing proposals (Carpenter and Govindarajan, 1993). Importantly, these accounts differ in the types and complexity of computations they assume to take place during normalization.

On the lower end of computational complexity, *estimation-free* psycho-acoustic transformations involve zero degrees of freedom that listeners would need to estimate from the acoustic input. For example, there is evidence that a transformation of acoustic frequencies (measured in Hz) into the psycho-acoustic Bark-space better describes how listeners perceive differences along the frequency spectrum (in terms of critical bands, e.g., Traunmüller, 1990; Zwicker, 1961; Zwicker *et al.*, 1957; Zwicker and Terhardt, 1980). It is thus possible that cross-talker variability in vowel pronunciations is reduced when formants are represented in Bark, rather than Hz. Similar arguments have been made about other psycho-acoustic transformations (e.g., ERB, Glasberg and Moore, 1990; Mel, Stevens and Volkmann, 1940; or semitones, Fant *et al.*, 2002) most of which share that they log-transform acoustic frequencies—in line with neurophysiological evidence that the auditory representations in the brain seem to follow a roughly logarithmic organization so that auditory perception is (up to a point) more sensitive to differences between lower frequencies than to the same difference between higher frequencies (e.g., Merzenich *et al.*, 1975; for review, see Saenz and Langers, 2014). While each of these transformations was developed with different applications in mind (e.g., ERB and Bark to explain frequency selectivity, Glasberg and Moore, 1990; or semitones for the perception of musical pitch, Balzano, 1982), psycho-acoustic transformations might suffice for effective formant normalization. If so, this would offer a particularly parsimonious account of vowel perception as listeners would not have to infer talker-specific properties.

The parsimony of psycho-acoustic transformations contrasts with the majority of accounts for vowel normalization, which introduces additional computations. This includes accounts that normalize formants relative to other information that is available at the same point in the acoustic signal (intrinsic normalization, e.g., Miller, 1989; Peterson, 1961; Syrdal and Gopal, 1986). For example, according to one proposal, listeners normalize vowel formants by the vowel's fundamental frequency or other formants estimated at the same point in time (Syrdal and Gopal, 1986). To the extent that the fundamental frequency is correlated with the talkers' vocal tract size (for review, see Vorperian and Kent,

2007), this allows the removal of physiologically-conditioned cross-talker variability in formant realizations. While such intrinsic accounts arguably entail more computational complexity than estimation-free transformations, they do not require that listeners *maintain* talker-specific estimates over time. This distinguishes intrinsic from extrinsic accounts, which introduce additional computational complexity.

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

In summary, hypotheses about the computations implied by formant normalization differ in the flexibility they afford as well as the inference and memory complexity they entail. Considerations about the complexity of inferences—essentially the number of parameters that listeners are assumed to estimate at any given moment in time—arguably gain importance in light of the speed at which normalization seems to unfold. In the present study, we thus ask whether computationally simple accounts are sufficient to explain human vowel perception.

While previous research has compared normalization accounts across languages, most of this work has evaluated proposals in terms of how well the normalized phonetic space supports the separability of vowel categories (Adank *et al.*, 2004; Carpenter and Govindarajan, 1993; Cole *et al.*, 2010; Escudero and Bion, 2007; Johnson and Sjerps, 2021; Syrdal, 1985). This approach is illustrated in Fig. 1. These studies have found that computationally more complex accounts—which also afford more flexibility—tend to achieve higher category separability and higher categorization accuracy (for review, see Persson and Jaeger, 2023). This includes Lobanov normalization, which continues to be

AQ3

185 highly influential in, for example, variationist and sociolinguistic research because of its effectiveness in removing
 186 cross-talker variability (for a critique, see Barreda, 2021). It
 187 is, however, by no means clear that human speech perception
 188 employs the same computations that achieve the best
 189 category separability or accuracy (see also discussion in
 190 Barreda, 2021; Nearey and Assmann, 2007).

192 A substantially smaller body of research has addressed
 193 this question by comparing normalization accounts against
 194 listeners' perception (Barreda and Nearey, 2012; Barreda,
 195 2021; Nearey, 1989; Richter *et al.*, 2017; for a review, see
 196 Whalen, 2016). Interestingly, these works seem to suggest
 197 that computationally simpler accounts might provide a better
 198 fit against human speech perception than the influential
 199 Lobanov model (Barreda, 2021; Richter *et al.*, 2017). For
 200 example, Barreda (2021) compared the predictions of uniform
 201 scaling and Lobanov normalization against listeners'
 202 categorization responses in a forced-choice categorization
 203 task over parts of the US English vowel space. In his experiment,
 204 listeners' categorization responses were better predicted by uniform scaling than by Lobanov normalization.
 205 Findings like these suggest that comparatively simple corrections
 206 for vocal tract size—such as uniform scaling—might provide a better explanation of human perception
 207 than more computationally complex accounts (see also
 208 Johnson, 2020; Richter *et al.*, 2017).

211 This motivates the present work. We take a broad-
 212 coverage approach by comparing the 20 normalization
 213 accounts in Table I against the perception of eight monophthongs of US English [i] as in *heed*, [ɪ] as in *hid*, [ɛ] as in *head*,
 214 [æ] as in *had*, [ʌ] as in *hut*, [ʊ] as in *hood*, [u] as in *who would*, [ɑ] as in
 215 [ɒ] as in *odd*.³ We do so for the perception of both natural and synthesized speech. Our broad-coverage approach complements
 216 previous studies, which have typically compared a small
 217 number of accounts (up to 3) and focused on parts of the
 218 vowel inventory, and thus parts of the formant space (typi-
 219 cally 2–4 vowels, Barreda, 2021; Barreda and Nearey, 2012;
 220

221 Nearey, 1989; Richter *et al.*, 2017). The accounts we consider include the most influential examples of psycho-
 222 acoustic transformations (Fant *et al.*, 2002; Glasberg and
 223 Moore, 1990; Stevens and Volkmann, 1940; Traunmüller,
 224 1981), intrinsic (Syrdal and Gopal, 1986), extrinsic
 225 (Gerstman, 1968; Johnson, 2020; Lobanov, 1971; McMurray and Jongman, 2011; Nearey, 1978; Nordström
 226 and Lindblom, 1975), and hybrid accounts that contain
 227 intrinsic and extrinsic components (Miller, 1989). This
 228 broad-coverage approach allows us to assess, for example,
 229 whether the preference for computationally simple accounts
 230 observed in Barreda (2021) replicates new data that span the
 231 entire vowel space. It also allows us to ask whether accounts
 232 even simpler than uniform scaling—such as psycho-acoustic
 233 transformations—provide an even better fit to human
 234 perception.

235 Next, we motivate and describe the two experiments we
 236 conducted. Then we compare the normalization accounts in
 237 Table I against listeners' responses from these experiments.

A. Open science statement

238 All stimulus recordings, results, and the code for the
 239 experiment, data analysis, and computational modeling for
 240 this article can be downloaded from the Open Science
 241 Framework (OSF) at <https://osf.io/zemwn/>. The OSF repository
 242 also includes extensive supplementary information (SI).
 243 Both the article and SI are written in R markdown, allowing
 244 readers to replicate our analyses with the click of a button,
 245 using freely available software (R Core Team, 2024;
 246 RStudio Team, 2020). Readers can revisit the assumptions
 247 we committed to for the present project—for example, by
 248 substituting alternative normalization accounts or categorization
 249 models. Researchers can also substitute their own
 250 experiments on vowel normalization for our Experiments 1a
 251 and 1b, to see whether our findings generalize to novel data.
 252 We see this as an important contribution of the present
 253 work, as it should make it substantially easier to consider
 254

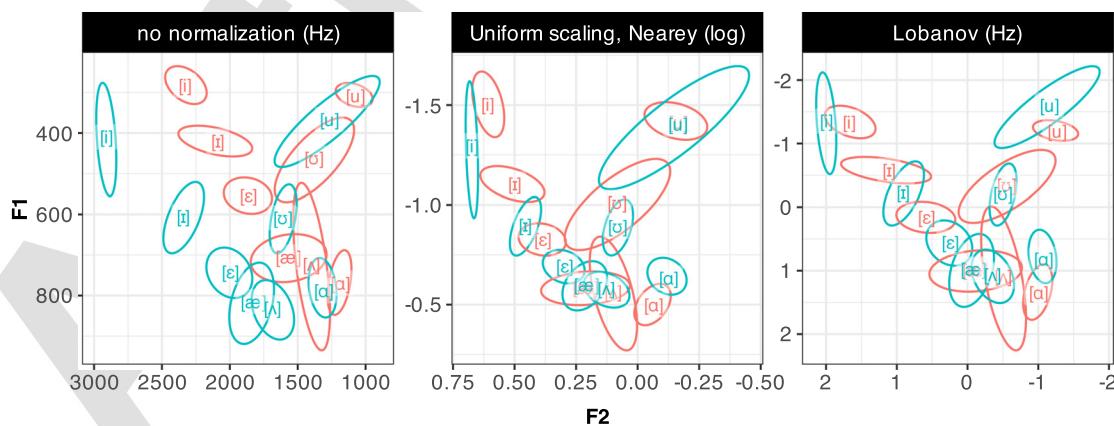


FIG. 1. Illustration of how height, which is positively correlated with vocal tract size, affects vowels F1 and F2, and how normalization can partially remove this effect. Shown here are realizations of eight monophthong vowels of US English by a short (cyan) and a tall native talker (red). (A) In the acoustic space, prior to any normalization (Hz). (B) After uniform scaling (Nearey, 1978). (C) After Lobanov normalization (Lobanov, 1971). The present study compares these three accounts, along with 17 other normalization accounts. Here, and throughout the paper, panel captions indicate the phonetic space in which normalization takes place in parenthesis. Note that this is not necessarily identical to the units of F1 and F2 *after* normalization (e.g., Lobanov normalization results in scale-free z-scores along the formant axes).

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

Normalization procedure		Perceptual scale	Source	Formula
Transformation	No normalization	Hz	n/a	n/a
	—	log		$F_n^{\log} = \ln(F_n)$
	—	Bark	Traunmüller (1990)	$F_n^{\text{Bark}} = \frac{26.81 \times F_n}{1960 + F_n} - 0.53$
	—	ERB	Glasberg and Moore (1990)	$F_n^{\text{ERB}} = 21.4 \times \log_{10}(1 + F_n \times 0.00437)$
	—	Mel	Stevens and Volkmann (1940)	$F_n^{\text{Mel}} = 2595 \times \log_{10}\left(1 + \frac{F_n}{700}\right)$
	—	Semitones conversion	Fant <i>et al.</i> (2002)	$F_n^{\text{ST}} = 12 \times \frac{\ln\left(\frac{F_n}{100}\right)}{\ln}$
Intrinsic	Syrdal and Gopal 1 (Bark-distance model)	Bark	Syrdal and Gopal (1986)	$F1_{\text{SyrdalGopal1}} = F1^{\text{Bark}} - F0^{\text{Bark}}$
	Syrdal and Gopal 2 (Bark-distance model)			$F2_{\text{SyrdalGopal1}} = F2^{\text{Bark}} - F1^{\text{Bark}}$
	Miller (formant-ratio)	log	Miller (1989)	$F1_{\text{SyrdalGopal2}} = F1^{\text{Bark}} - F0^{\text{Bark}}$
				$F2_{\text{SyrdalGopal2}} = F3^{\text{Bark}} - F2^{\text{Bark}}$
				$SR = k \left(\frac{GMf0}{k} \right)^{1/3}$
				$F1_{\text{Miller}} = \log\left(\frac{F1}{SR}\right)$
Extrinsic centering	Nearey's uniform scaling	log	Nearey (1978)	$F2_{\text{Miller}} = \log\left(\frac{F2}{F1}\right)$
	Nordström and Lindblom (vocal tract scaling)	Hz	Nordström and Lindblom (1975)	$F3_{\text{Miller}} = \log\left(\frac{F3}{F2}\right)$
	Johnson (average formant spacing)	Hz	Johnson (2020)	$F1_{\text{Nearey}} = \ln(F_n) - \text{mean}(\ln(F))$
	Nearey's formantwise log-mean	log	Nearey (1978)	$F2_{\text{Nearey}} = \ln(F_n) - \text{mean}(\ln(F_n))$
	C-CuRE	Hz	McMurray and Jongman (2011)	$F3_{\text{Nearey}} = F_n - \text{mean}(F_n)$
	—	Bark		
Extrinsic standardizing	—	ERB		
	—	Mel		
	—	Semitones conversion		
	Gerstman (range normalization)	Hz	Gerstman (1968)	$F_n^{\text{Gerstman}} = 999 \times \frac{F_n - F_n^{\min}}{F_n^{\max} - F_n^{\min}}$
	Lobanov (z-score)	Hz	Lobanov (1971)	$F_n^{\text{Lobanov}} = \frac{F_n - \text{mean}(F_n)}{\text{sd}(F_n)}$

258 additional normalization accounts—including variants to the
 259 accounts we considered—and to assess the generalizability
 260 of the conclusions we reach based on the present data.

261 II. EXPERIMENTS 1A AND 1B

262 To compare the performance of different normalization
 263 accounts against listeners' perceptions, we conducted two
 264 small web-based experiments on US English listeners' per-
 265 ception of US English vowels. Both experiments investigate
 266 listeners' perceptions of a single talker. This choice was
 267 made so as to not confound questions about formant normal-
 268 ization with questions about talker recognition, and infer-
 269 ences about talker switches (Magnuson and Nusbaum,

270 2007). The two experiments employ the same eight-
 271 alternative forced-choice vowel categorization task (Fig. 2),
 272 and differ only in the whether they employed "natural"
 273 (Experiment 1a) or synthesized stimuli (Experiment 1b). To
 274 the best of our knowledge, these two experiments are the
 275 first designed to compare normalization accounts against lis-
 276 teners' perception over a larger portion of the monophthong
 277 inventory of a language.

278 Experiment 1a employs recordings of *hVd* word produc-
 279 tions from a female talker of US English, these recordings
 280 are "natural" in the sense that they were not synthesized or
 281 otherwise phonetically manipulated. One consequence of
 282 this is that the formant values of these recordings are clus-
 283 tered around the talker's category means, and thus span only

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 Experiments 1a and 1b.

284 a comparatively small part of the phonetic space. This can
 285 limit the statistical power to distinguish between competing
 286 accounts. Natural recordings furthermore vary not only
 287 along the primary cues to vowel quality in US English (F1,
 288 F2) but also along secondary cues (e.g., F0, F3, vowel dura-
 289 tion, and vowel inherent spectral change—VISC) as well as
 290 other unknown properties, which can make it difficult to dis-
 291 cern whether the performance of a normalization model is
 292 due to the normalization itself or other reasons, e.g., because
 293 a normalized cue happens to correlate with another cue that
 294 listeners are sensitive to but that is not included in the
 295 model.

296 Experiment 1b thus adopts an alternative approach and
 297 uses synthesized vowels. Unlike most previous work, which
 298 has used isolated vowels as stimuli (Barreda, 2021; Barreda
 299 and Nearey, 2012; Nearey, 1989; Richter *et al.*, 2017),
 300 Experiment 1b uses synthesized *hVd* words to facilitate
 301 comparison to Experiment 1a. This allowed us to sample
 302 larger parts of the F1–F2 space, which has two advantages.
 303 First, it allowed us to collect responses over parts of the for-
 304 mant space for which we expect listeners to have more
 305 uncertainty, and thus exhibit more variable responses. This
 306 can increase the statistical power to distinguish between
 307 competing accounts. Second, differences in the predictions
 308 of competing normalization accounts will tend to become
 309 more pronounced with increasing distance from the category
 310 centers. By collecting responses at those locations, we can
 311 thus increase the contrast between competing accounts.
 312 Critically, an adequate model of formant normalization
 313 needs to capture human perception not only for prototypical
 314 vowel instances but also for instances of vowels that fall
 315 between category means.

316 The use of synthesized stimuli does, however, also
 317 come with potential disadvantages. Synthesized stimuli can
 318 suffer in ecological validity, lacking correlations between
 319 cues, and across the speech signal (e.g., due to co-articula-
 320 tion) that are characteristic of human speech. This raises
 321 questions about the extent to which the processing of such
 322 stimuli engages the same mechanisms as everyday speech
 323 perception. Additionally, it is possible that the use of robotic
 324 sounding synthesized speech affects listener engagement.
 325 This can lead to an increased rate of attentional lapses, and
 326 thus a decrease in the proportion of trials in which listeners'
 327 responses are based on the acoustics of the speech stimulus

rather than random guessing (compare, e.g., Kleinschmidt, 328
 2020; Tan and Jaeger, 2024). By comparing normalization 329
 accounts against both natural and synthesized stimuli, we 330
 investigate the extent to which the accounts that best 331
 describe human perception depend on the type of stimuli 332
 used in the experiment. 333

A. Methods 334

1. Participants 335

We recruited 33 (Experiment 1a) and 33 (Experiment 336
 1b) participants. The majority of these (24 for each experi- 337
 ment) were recruited from Amazon's Mechanical Turk. 338
 However, after exclusions, we were left with a relatively 339
 low number of participants (for Experiment 1a, 19, and for 340
 Experiment 1b, 22). We therefore decided to recruit an addi- 341
 tional 18 participants from Prolific (9 for each experiment; 342
 October 2024). Exclusions described in the following left 28 343
 and 31 participants for analysis in Experiments 1a and 1b, 344
 respectively. Results did not change after the inclusion of 345
 the new participants from Prolific. 346

Participants were paid \$6/h (\$12/h on Prolific) prorated 347
 by the duration of the experiments (15 min). Participants 348
 only saw the experiment advertised, and could only partici- 349
 pate in it, if (i) they were located within the US, (ii) had an 350
 approval rating of 99% or higher, (iii) met the software 351
 requirements (a recent version of the Chrome browser 352
 engine), and (iv) had not previously completed any other 353
 experiments on vowel perception in our lab. Before the 354
 experiment could be accepted, participants had to confirm 355
 that they were (i) native speakers of US English (defined as 356
 having spent their childhood until the age of ten speaking 357
 English and living in the United States), (ii) in a quiet room 358
 without distractions, (iii) wearing over-the-ear headphones. 359
 Participants' responses were collected via JavaScript devel- 360
 oped by the Human Language Processing Lab at the 361
 University of Rochester (Kleinschmidt *et al.*, 2021). 362

An optional post-experiment survey recorded partici- 363
 pant demographics using National Institutes of Health (NIH) 364
 prescribed categories, including participant sex (Male: 36, 365
 Female: 29), age [mean = 36.9 years; standard deviation 366
 (SD) = 12.2; 95% quantiles = 22.6–66 years], race (White: 367
 48, multiple: 3, Black: 10, Asian: 3, declined to report: 1), 368
 and ethnicity (Non-Hispanic: 60, Hispanic: 4, declined to 369
 report: 1). All but one participant completed the survey. 370

2. Materials 371

Experiment 1a employed *hVd* word recordings by one 372
 adult female talker of a Northeastern dialect (spoken in cen- 373
 tral Connecticut) from a phonetically annotated database of 374
 L1-US English vowel productions (Xie and Jaeger, 2020). 375
 Specifically, we used all nine recordings of each of the eight 376
hVd-words—*heed*, *hid*, *head*, *had*, *hut*, *odd*, *hood*, *who* 377
would (the use of “*hut*” and “*odd*” rather than “*hud*” and 378
 “*hod*” follows Assmann *et al.*, 2008; but see Hillenbrand 379
et al., 1995). 380

381 The stimuli for Experiment 1b were synthesized from a
 382 single *had* recording used in Experiment 1a (see Fig. 3 for
 383 example spectrograms). Specifically, we used a script (based
 384 on descriptions in Wade *et al.*, 2007) in Praat (Boersma and
 385 Weenink, 2022) to concatenate the original /h/ with a syn-
 386 thesized vowel and the original /d/ recording. Unlike in
 387 Experiment 1a, all eight words thus had an *hVd* context
 388 (including “*hud*” and “*hod*,” rather than “*hut*” and “*odd*”).
 389 The Praat script first segmented the original *had* token into
 390 the three segments /h/, /æ/, and /d/, with the /d/ segment
 391 consisting of the voiced closure and burst. The script then
 392 estimated the spectral envelope of the /h/ sound by linear
 393 predictive coding (LPC; autocorrelation method), and used
 394 the resulting coefficients to inversely filter the /h/. This
 395 resulted in an /h/ sound with the effects of the vocal tract
 396 removed, leaving the source signal. Next, a glottal wave-
 397 form was generated at each point in the pitch contour from
 398 the original /æ/ sound using the point process to phonation
 399 functionality in Praat. This waveform was multiplied with
 400 the intensity pattern from the same original /æ/ sound. The
 401 resulting sound was concatenated with the neutral fricative
 402 /h/ sound, to create a neutral hV-section that did not reflect
 403 any vocal tract resonances. The script then created a formant
 404 grid that filtered the hV-section to create the intended vowel,
 405 and finally concatenated this segment to the final /d/ to cre-
 406 ate an *hVd* word. For each *hVd* word, the formant grid was
 407 populated with the F1, F2, and F3 values that we handed to
 408 the script at five time-points transitioning from the /h/ to the
 409 steady-state vowel, to the first portion of the voiced closure
 410 of the final /d/ segment through linear interpolation, thus
 411 holding formants steady until transitioning into the final
 412 consonant. Formant bandwidths were 500 Hz at the initial

two time-points (the /h/ and beginning of the transition to vowel), and then decreased linearly during vowel onset and throughout the final three time-points to 50 Hz (F1), 100 Hz (F2), 200 Hz (F3), 300 Hz (F4), and 400 Hz (F5–F8, follow- 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436 AQ19
 437
 438
 439
 440
 441
 442
 443
 444

413 vowel), and then decreased linearly during vowel onset and 414 throughout the final three time-points to 50 Hz (F1), 100 Hz 415 (F2), 200 Hz (F3), 300 Hz (F4), and 400 Hz (F5–F8, follow- 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436 AQ19
 437
 438
 439
 440
 441
 442
 443
 444

413 vowel), and then decreased linearly during vowel onset and 414 throughout the final three time-points to 50 Hz (F1), 100 Hz 415 (F2), 200 Hz (F3), 300 Hz (F4), and 400 Hz (F5–F8, follow- 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436 AQ19
 437
 438
 439
 440
 441
 442
 443
 444

413 We generated 146 synthesized *hVd* recordings that 414 spanned the F1 and F2 space. The specific F1–F2 locations 415 chosen were determined by a mix of modeling (using ideal 416 observers described in the next section to predict listeners’ 417 categorization responses) and intuition. Specifically, we 418 selected 64 recordings that we expected to fall within the 419 bivariate 95% confidence intervals (CIs) of the eight US 420 English monophthongs, and 82 recordings that we expected 421 to fall between those CIs. Figure 4 shows the distribution of 422 stimuli for both experiments. Of note, our procedure also 423 generated formant combinations that are physiologically 424 unlikely to have all been produced by the same talker during 425 ‘normal’ vowel production (also known as “off-template” 426 instances, Nearey, 1978). 427

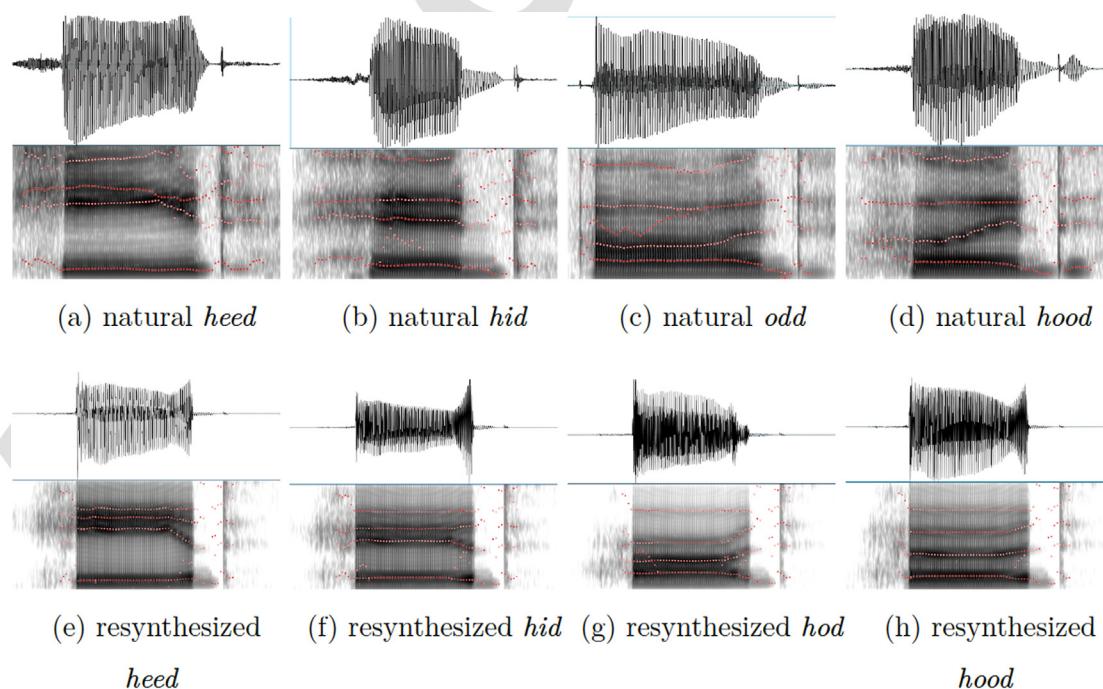


FIG. 3. Top: Spectrograms of four natural recordings from Experiment 1a. Bottom: Same for four synthesized tokens with similar formant values from Experiment 1b. Additional spectrograms are provided in the supplementary material, Sec. 2C.

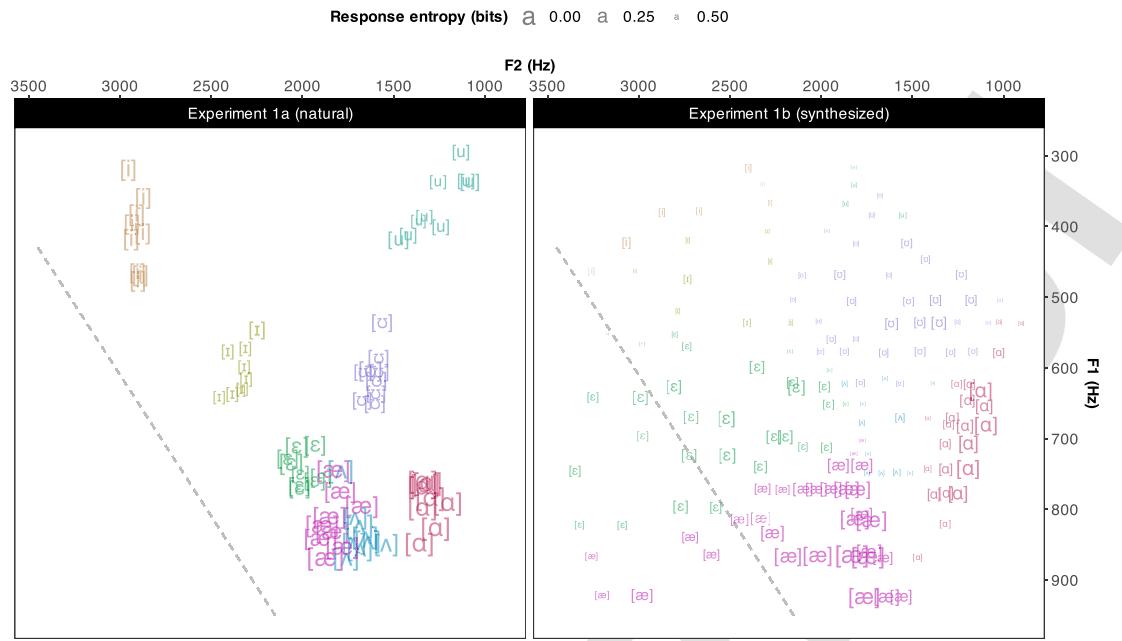


FIG. 4. Summary of listeners' categorization responses in Experiments 1a and 1b in F1–F2 space. The vowel label indicates the most frequent response provided by participants at each test location. Size indicates how consistent responses were across participants, with larger symbols indicating more consistent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

445 3. Procedure

446 The procedure for both experiments was identical. Live
 447 instances of each experiment can be found at <https://www.hlp.rochester.edu/experiments/DLPL2S/experiment-A/experiments.html>. At the start of the experiment, participants
 448 acknowledged that they met all requirements and provided
 449 consent, as per the Research Subjects Review Board of the
 450 University of Rochester. Before starting the experiment, participants
 451 performed a sound check. Participants were then
 452 instructed to listen to a female talker saying words and click
 453 on the word on screen to report what word they heard. On
 454 each trial, all eight *hVd*-words were displayed on screen.
 455 Half of the participants in each experiment saw the response
 456 options organized as in Fig. 2 (resembling the IPA represen-
 457 tation of a vowel space), and half saw the response options
 458 in the opposite order (flipping top and bottom and left and
 459 right in Fig. 2). Each trial started with the response grid on
 460 screen, together with a light green dot centered on screen.
 461 After 1000 ms, an *hVd* recording played, and participants
 462 indicated their response by a mouse-click. After a 1000 ms
 463 intertrial interval, the screen reset, and the next trial started.

464 In both experiments, participants heard two blocks of
 465 the materials described in the previous sections, for a total
 466 of 144 trials in Experiment 1a and 292 trials in Experiment
 467 1b. Presentation within each block was randomized for each
 468 participant in order to reduce confounds due to stimulus
 469 order (known to affect vowel perception, [Repp and](#)
 470 [Crowder, 1990](#), and references therein). Participants were
 471 not informed about the block structure of the experiment.

472 After completing the experiment, participants filled out
 473 a language background questionnaire and the optional
 474 demographic survey. On average, participants took 9.3 min
 475 in Experiment 1a, and 9.7 min in Experiment 1b.

476 to complete Experiment 1a ($SD = 5.5$) and 17.9 min for Experiment 1b ($SD = 6.5$). 477 478

479 4. Exclusions

480 We excluded participants who failed to follow instruc-
 481 tions and did not wear over-the-ear headphones (as indicated
 482 in the post-experiment survey). We also excluded partici-
 483 pants with mean (log-transformed) reaction times that were
 484 unusually slow or fast (absolute z-score over by-participant
 485 means > 3), or if they clearly did not do the task (e.g., by
 486 answering randomly). This excluded five participants from
 487 Experiment 1a and 2 from Experiment 1b (for details, see
 488 Sec. II A). 489 AQ5

490 We further excluded all trials that were unusually fast or slow. Specifically, we first z-scored the log-transformed
 491 response times *within each participant* and then z-scored
 492 these z-scores *within each trial* across participants. Trials
 493 with absolute z-scores > 3 were removed from analysis.
 494 This double-scaling approach was necessary as participants' response times decreased substantially over the first
 495 few trials and then continued to decrease less rapidly
 496 throughout the remainder of the experiment. The approach
 497 removes response times that are unusually fast or slow *for*
 498 *that participant at that trial*, while avoiding specific
 499 assumptions about the shape of the speed up in response
 500 times across trials. This excluded 1.3% of the trials in
 501 Experiment 1a and 0.9% in Experiment 1b. This left for
 502 analysis 3983 observations from 28 participants in
 503 Experiment 1a, and 8970 observations from 31 participants
 504 in Experiment 1b. 505

506 **B. Results**

507 Participants' categorization responses in Experiments
 508 1a and 1b are shown in Fig. 4, with larger labels indicating
 509 recordings that participants agreed on more.⁴ We make two
 510 observations. The first pertains to the degree of (dis)agree-
 511 ment between the two experiments. The second observation
 512 pertains to the degree of (dis)agreement across participants
 513 within each experiment.

514 **1. Similarities and differences between Experiments**
 515 **1a and 1b**

516 Unsurprisingly, participants in both experiments
 517 divided the F1–F2 space into the eight vowel categories in
 518 ways that qualitatively resembled each other (after taking
 519 into account that Experiment 1b covers a larger range of
 520 F1–F2 values). Also, unsurprisingly, there were some differ-
 521 ences between participants' responses across the two experi-
 522 ments, at least when plotted in Hz. For example, [u] rarely
 523 was the most frequent response in Experiment 1b, even for
 524 stimuli with similar F1–F2 values that were predominantly
 525 categorized as [u] in Experiment 1a. There are at least two
 526 reasons to expect such differences. First, stimuli with similar
 527 F1–F2 values across the two experiments still differed in
 528 other acoustic properties (e.g., vowel duration or F3). These
 529 acoustic differences might have affected participants'
 530 responses. Second, it is possible that *formant normalization*
 531 affected participants' responses—i.e., the very mechanism
 532 we seek to investigate in the remainder of the paper. The
 533 two experiments differ in the means, variances, and other
 534 statistical properties that some normalization accounts pre-
 535 dict to affect perception. As a consequence, Hz might not be
 536 the space in which we should expect identical responses
 537 across experiments.

538 Similarly, the two experiments differed in the extent to
 539 which participants agreed with each other. Participants in
 540 Experiment 1b exhibited overall less agreement in their
 541 responses [mean by-item response entropy = 0.45 bits, stan-
 542 dard error (SE) = 0.01] than participants in Experiment 1a
 543 (mean by-item response entropy = 0.19 bits, SE = 0.02).
 544 This was also confirmed by participants' responses during
 545 the post-experiment survey. Compared to participants in
 546 Experiment 1a, participants in Experiment 1b reported
 547 increased uncertainty about their responses, and that the
 548 stimuli were less distinguishable and more robotic-sounding
 549 (see supplementary material Sec. 2 B).

550 This increased uncertainty in Experiment 1b was
 551 expected—and, indeed, intended by the design: Experiment
 552 1b explored the entire F1–F2 space, including formant com-
 553 binations located *between* the centers of the natural vowel
 554 categories. Experiment 1b therefore achieved its goal of
 555 eliciting less categorical response distributions, which is
 556 expected to facilitate comparison of competing normaliza-
 557 tion accounts.⁵

558 Auxiliary analyses presented in the supplementary
 559 material Sec. 2 B suggest that *some but not all* of the differ-
 560 ences in response entropy between the two experiments

were caused by the placement of the stimuli in formant
 561 space: when comparing categorization responses for tokens
 562 from the two experiments with similar acoustic properties
 563 (differences of ≤ 30 Hz along F1 and F2), response entro-
 564 pies still differed substantially (for $N = 40$ acoustically simi-
 565 lar tokens, mean by-item response entropy for Experiment
 566 1a = 0.14 bits, SE = 0.02; Experiment 1b = 0.4 bits,
 567 SE = 0.03). The same section of the supplementary materi-
 568 als, Sec. 2 E, presents additional analyses grouping acousti-
 569 cally similar tokens in the phonetic space defined by the
 570 normalization account we find to best fit listeners'
 571 responses. These analyses support the same conclusion.
 572

We see two mutually compatible explanations to this
 573 difference in listener agreement between experiments. First,
 574 similar to the differences between experiments in the domi-
 575 nant response pattern discussed previously, differences in
 576 the degree of agreement between participants might origi-
 577 nate in *normalization*. Second, it is possible that the relation
 578 between formants in the synthesized stimuli or some other
 579 unknown acoustic-phonetic differences between the experi-
 580 ments explain the difference in response. For example, the
 581 absence of VISC or differences in spectral tilt in the synthe-
 582 sized stimuli might have deprived listeners of information
 583 that is actually crucial for establishing phonemic identity
 584 (Hillenbrand and Nearey, 1999). This would result in
 585 increased uncertainty on each trial, leading to increased
 586 entropy of listeners' responses. The computational study we
 587 present in the following sheds some light on these two mutu-
 588 ally compatible possibilities.
 589

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

591 Since the intended category was known for Experiment
 592 1a, it was possible to calculate participants' recognition
 593 accuracy. As also evident in the left panel of Fig. 4, partici-
 594 pants' most frequent response *always* matched the intended
 595 vowel in Experiment 1a. Overall, participants' responses
 596 matched the intended vowel on 84.7% (SE = 3.5%) of all
 597 trials (Experiment 1b had no such ground truth). This is
 598 much higher than chance (12.5%). It is, however, also quite
 599 a bit lower than 100%. To better understand the reasons for
 600 this, Fig. 5(A) plots the confusion matrix. This suggests that
 601 participants' performance was largely affected by confu-
 602 sions between [i]-to-[ɛ] (*hid-to-head*), [ɛ]-to-[æ] (*head-to-
 603 had*), and [u]-to-[u] (*who would-to-hood*).
 604

One plausible explanation for this pattern of vowel con-
 605 fusion lies in the substantial variation that exists across US
 606 English dialects (Labov *et al.*, 2006). Differences in the real-
 607 ization of vowel categories, and associated representations,
 608 across dialects will directly affect the expected classification
 609 for any given token. In addition, listeners might differ in
 610 terms of experience with different dialects, or in the dialect
 611 they attribute to the talker who produced the stimuli. To test
 612 this hypothesis, we calculated the [i]-to-[ɛ], [ɛ]-to-[æ], and
 613 [u]-to-[u] confusion rates for each participant in Experiment
 614 1a. These data are summarized in the left panel of Fig. 5(B).
 615 The data in the left panel suggest that most participants in

AQ6

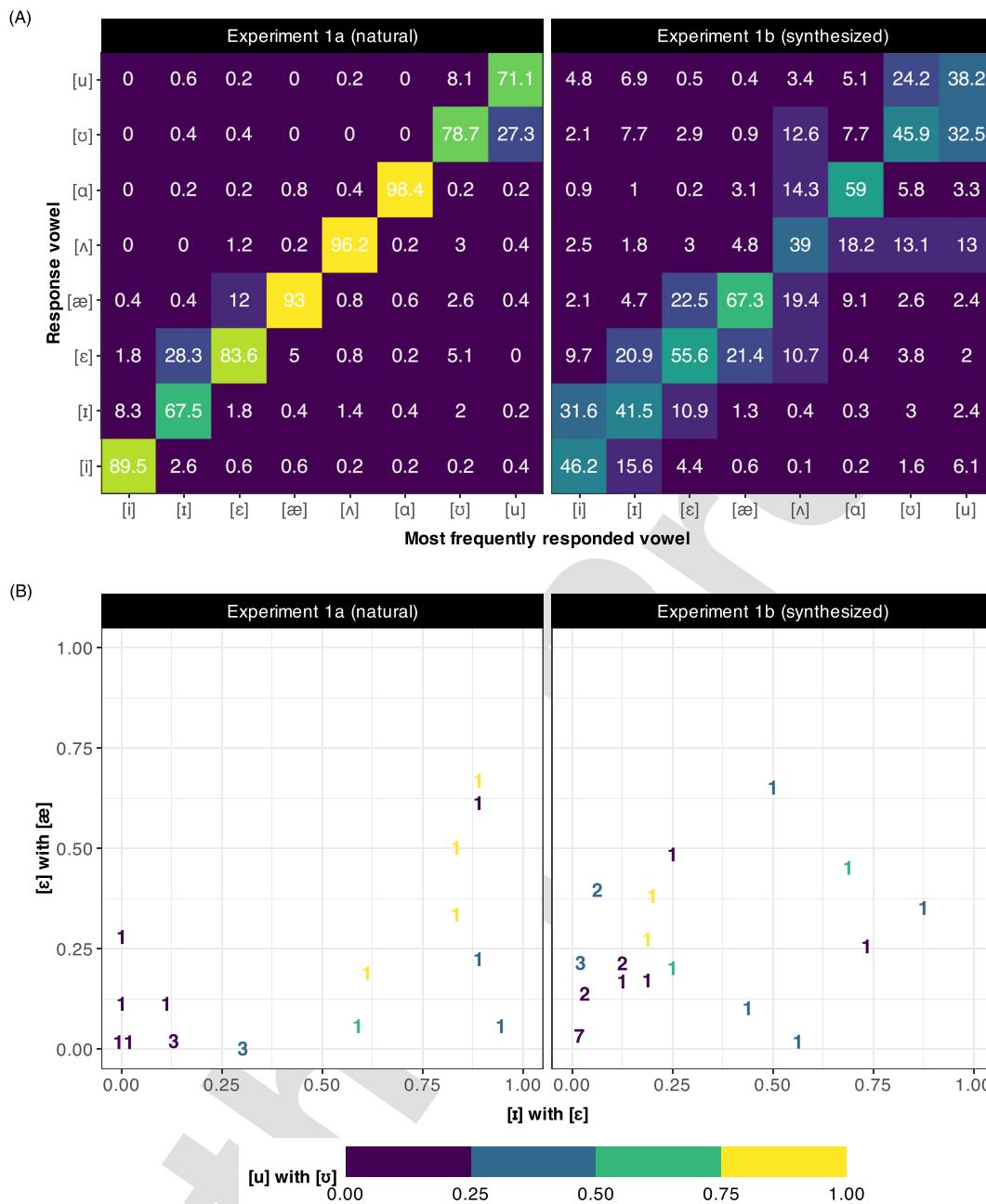


FIG. 5. Category confusability in Experiments 1a and 1b. (A) summarizes the category confusability. Since correct responses were not defined for Experiment 1b, we grouped items along the x -axis based on the most frequent response that listeners provided (for Experiment 1a, this was always identical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

616 Experiment 1a either heard [ɪ] tokens consistently as the
 617 intended [ɪ] (clustering on the left side of the panel) or as [ɛ]
 618 (clustering on the right side of the panel). Only a few partici-
 619 pants exhibited mixed responses for items intended to be [ɪ].
 620 Tellingly, many of the participants who exhibited increased
 621 [ɪ]-to-[ɛ] confusion *also* exhibited increased [ɛ]-to-[æ] con-
 622 fusion. This is precisely what would be expected by listeners
 623 who assume a dialect in which these vowels are articulated
 624 lower (with higher F1) than in the dialect of the talker in
 625 Experiment 1a. A similar, but less pronounced, pattern was

also found with regard to [u]-to-[ʊ] confusions.⁶ Finally, a qualitatively similar relation between [ɪ]-to-[ɛ], [ɛ]-to-[æ], and [u]-to-[ʊ] confusions was also observed in Experiment 1b [right panel of Fig. 5(B)], though the pattern was unsurprisingly less pronounced given that the stimuli in Experiment 1b by design often fell into the ambiguous region *between* vowels. Taken together, vowel-to-vowel confusion rates in Experiments 1a and 1b suggest that systematic dialectal differences contributed to the relatively low categorization accuracy.

636 This highlights two important points. First, the data
 637 from Experiment 1a demonstrate the perceptual challenges
 638 associated with an unfamiliar talker: in the absence of lexical
 639 or other context to distinguish between the eight available
 640 response options, listeners can only rely on the acoustic
 641 information in the input. In such a scenario, even listeners
 642 who are in principle familiar with the dialect spoken by the
 643 talker have comparatively little information to determine the
 644 talker's dialect, making apparent what Winn (2018) aptly
 645 summarizes as "speech [perception] is not as acoustic as
 646 [we] think." Second, when dialect variability is taken into
 647 account, listeners' recognition accuracy improved substantially.
 648 After removing eight listeners who heard more than
 649 50% of the [ɪ] items as [ɛ], all vowels were correctly recognized
 650 at least 87.1% of the time (overall accuracy = 94.8%).
 651 This suggests that dialect differences affected the recognition
 652 of all vowels. This aspect of our results serves as an
 653 important reminder that formant normalization is only
 654 expected to erase inter-talker variability associated with
 655 physiological differences: variation in dialect, sociolect, or
 656 other non-physiologically-conditioned variation pose separate
 657 challenges to human perception, and require additional
 658 mechanisms (see discussion in Barreda, 2021; Weatherholtz
 659 and Jaeger, 2016). This introduces noise—variability in listeners' responses that cannot be accounted for by normalization—to any comparison of normalization accounts, potentially reducing the power to detect differences between accounts.

664 III. COMPARISON OF NORMALIZATION ACCOUNTS

665 In order to evaluate normalization accounts against
 666 speech perception, it is necessary to map the phonetic properties
 667 of stimuli—under different hypotheses about normalization—onto listeners' responses in Experiments 1a and 1b.
 668 Previous work has done so by directly predicting listeners' responses from the raw or normalized phonetic properties of
 669 stimuli (Apfelbaum and McMurray, 2015; Barreda, 2021;
 670 Crinnion *et al.*, 2020; McMurray and Jongman, 2011;
 671 Nearey, 1989). For example, McMurray and Jongman used
 672 multinomial logistic regression to predict eight-way fricative
 673 categorization responses in US English (see also Barreda, 2021).

674 Here, we pursued an alternative approach by committing to a core assumption common to contemporary theories
 675 of speech perception: that listeners acquire implicit knowledge about the probabilistic mapping from acoustic inputs to
 676 linguistic categories, and draw on this knowledge during speech recognition (e.g., TRACE, McClelland and Elman,
 677 1986; exemplar theory, Johnson, 1997; Bayesian accounts,
 678 Luce and Pisoni, 1998; Nearey, 1990; Norris and McQueen,
 679 2008; ASR-inspired models like DIANA or EARSHOT, ten
 680 Bosch *et al.*, 2015; Magnuson *et al.*, 2020). Using a general
 681 computational framework for adaptive speech perception
 682 (ASP, Xie *et al.*, 2023) we trained Bayesian ideal observers
 683 to capture the expectations that a "typical" L1 adult listener
 684 might have about the formant-to-vowel mappings of US
 685 686 687 688 689 690

691 English. We approximated these expectations using a database of L1-US English vowel productions (Xie and Jaeger, 692 2020)—transformed to reflect the different normalization 693 accounts. We then ask which of the different ideal observer 694 models—corresponding to different hypotheses about formant 695 normalization—best predicts listeners' responses in Experiments 696 1a and 1b. 697

698 Training ideal observers on a database of vowel productions 699 have the advantage that it reduces the degrees of freedom (DFs) 700 used to predict listeners' responses. For example, using ordinary 701 multinomial logistic regression trained on our perceptual data to 702 predict eight-way vowel categorization as a function of F1, F2 and their 703 interaction would require up to 28 DFs. This problem increases with the 704 number of cues considered. By instead training ideal observers 705 on phonetic data that are independent of listeners' 706 responses, the ASP-based approach we employ uses only 707 two DFs to mediate the mapping from stimuli properties to 708 listeners' responses, regardless of the number of cues considered. 709 Over the next few sections, we describe how this parsimony 710 is made possible through a commitment to strong 711 linking hypotheses motivated by theories of speech 712 perception. 713

714 A. Methods

715 1. A general-purpose categorization model for J-AFC 716 categorization tasks AQ8

717 Figure 6 summarizes ASP's categorization model for a 718 J-alternative forced-choice task (for an in-depth description, 719 we refer to Xie *et al.*, 2023). The model combines Bayesian 720 ideal observers (as used in e.g., Clayards *et al.*, 2008; 721 Feldman *et al.*, 2009; Norris and McQueen, 2008; Xie *et al.*, 722 2021; for a closely related approach, see also Nearey and 723 Hogan, 1986) with psychometric lapsing models 724 (Wichmann and Hill, 2001). To reduce researchers' degrees 725 of freedom, we adopt all assumptions made in Xie *et al.* 726 (2023) and do not introduce additional assumptions. 727

728 Starting at the bottom of the figure, the acoustic input x 729 is normalized. Here, we follow most previous evaluations of 730 normalization accounts and focus on the point estimates of 731 formants at the center of the vowel as the inputs to normalization. 732 This leaves open the question of how considerations of additional 733 cues to vowel identity (e.g., VISC) or formant dynamics 734 might affect the findings we report in the following (a point to 735 which we return in the general discussion). Specifically, the main 736 analysis we present here focuses on $x = F1$ and $F2$. As one anonymous 737 reviewer pointed out, this focus on $F1$ – $F2$ might underestimate the potential of 738 *intrinsic* normalization accounts, which might perform better 739 when more acoustic-phonetic features are considered. The 740 supplementary material, Sec. 3 E, thus reports additional 741 analyses that instead employ $F1$ – $F3$. These analyses indeed 742 find that the fit of intrinsic normalization accounts improves 743 more than that of extrinsic accounts when $F3$ is included in 744 the analysis. However, the best-fitting accounts were still 745

AQ7

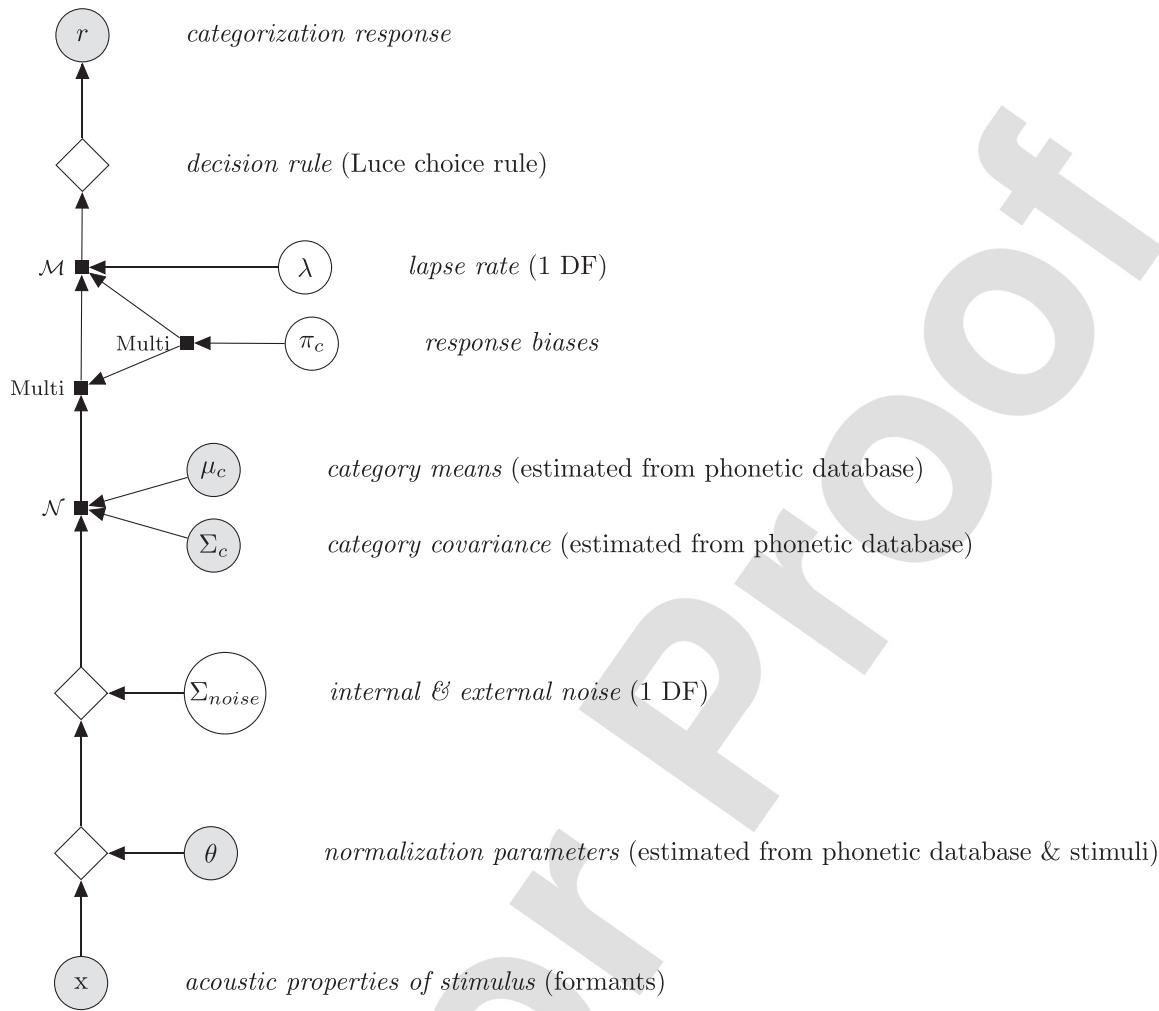


FIG. 6. Graphical model of ASP's general categorization framework (adapted for the current purpose from Xie *et al.*, 2023, Fig. 4). Here, $J = 8$ (the eight vowel response options in Experiments 1a and 1b). We use this framework to compare normalization accounts against listeners' categorization responses 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.

745 the same extrinsic accounts we found to best fit listeners' 745 responses when only F1 and F2 were considered.

746 The specific computations applied to the input x depend 746 on the normalization accounts (see Table I). We use θ to refer 747 to the parameters required by the normalization account. For 748 example, for Nearey's uniform scaling account (Nearey, 748 1978), θ is the overall mean of all log-transformed formants. 749 For Lobanov normalization (Lobanov, 1971), θ is a vector of 750 means and standard deviations for each formant (in Hz). The 751 normalized input is then perturbed by perceptual and environ- 752 mental noise. Following Feldman *et al.* (2009), this noise is 753 assumed to be Gaussian distributed centered around the trans- 754 formed stimulus with noise variances that are independent and 755 identical for all formants (i.e., Σ_{noise} is a diagonal matrix, and 756 all diagonal entries have the same value).

757 Next, the likelihood of the normalized percept under 757 each of the eight vowel categories is calculated, 758 $p(F1, F2|vowel)$. This requires specifying listeners' expecta- 759 tions about the cue-to-category mapping (listeners' likeli- 760 hood function). We followed Xie *et al.* (2023) and previous

761 work and assume that each vowel maps onto a multivariate 761 Gaussian distribution over the phonetic cues, here bivariate 762 Gaussians over F1 and F2 (cf. Clayards *et al.*, 2008; 763 Feldman *et al.*, 2009; Kleinschmidt and Jaeger, 2015; Norris 764 and McQueen, 2008; Xie *et al.*, 2021). We also followed 765 previous models in assuming a single dialect template—i.e., 766 a single set of bivariate Gaussian vowel categories (Nearey 767 and Assmann, 2007). The analyses of participants' 768 responses we provided previously in the description of 769 Experiments 1a and 1b suggest that this assumption is 770 wrong. However, more appropriate alternatives—such as 771 hierarchical or mixture models with multiple dialect tem- 772 plates—will require substantial additional research as well 773 as larger databases of vowel recordings that have high reso- 774 lution both within and across dialects. We return to this 775 issue in the general discussion.

776 Once the likelihood function for each vowel is speci- 776 fied, the posterior probability of each vowel is obtained by 777 combining its likelihood with its prior probability or 778 response bias π_c , according to Bayes theorem,⁷ 779

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

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

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

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

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

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

By fixing θ , μ_c , and Σ_c based on a database of vowel productions, we impose strong constraints on the functional

flexibility of the model in predicting listeners' responses. This benefit is made possible by committing to a strong linking hypothesis—that listeners' categories are learned from, and reflect, the distributional mapping from formants to vowels in previously experienced speech input (e.g., Abramson and Lisker, 1973; Massaro and Friedman, 1990; Nearey and Hogan, 1986). The database we use to approximate listeners' prior experience was originally developed to compare the production of L1 and L2 speakers (Xie and Jaeger, 2020). It contains 9–10 recordings of the eight *hVd* words from each of 17 (five female) L1 talkers of a Northeastern dialect of US English (ages 18 to 35 years old). Since Experiments 1a and 1b used recordings of one of these talkers, we excluded that talker prior to fitting training ideal observers on the data. In total, this yields 5842 recordings that are annotated for F0, F1–F3, and vowel duration. The supplementary material, Sec. 3 A 1, summarizes the distribution of these cues, and how the different normalization accounts affect those distributions.

To avoid over-fitting the ASP model to the database, we used fivefold cross-validation: we randomly split the Xie and Jaeger (2020) database into five approximately evenly-sized folds (following Persson and Jaeger, 2023). This split was performed within each vowel to guarantee that all five folds had the same relative amount of data for each vowel category. These splits were combined into five training sets, each containing one of the folds (20% of the data). This way, each training set was different from the others, increasing the variability between sets.⁸

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

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

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

Next, we transformed the stimuli of Experiments 1a and 1b into the formant space defined by the 20 normalization accounts in Table I. This requires estimating the required normalization parameters θ for each experiment and normalization account. We calculated these θ s over all stimuli

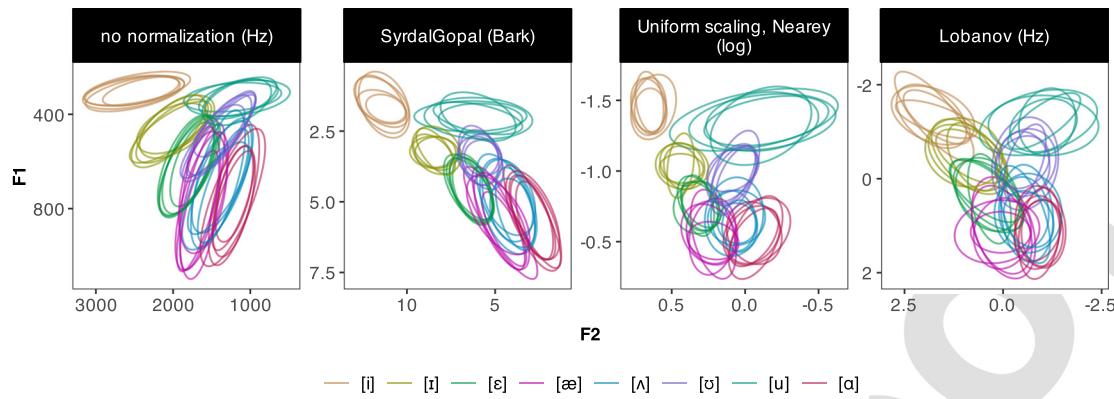


FIG. 7. Visualizing the bivariate Gaussian categories (prior to adding Σ_{noise}) of four example normalization accounts in F1–F2 space. Separate ellipses are shown for each of the five training sets (each set corresponds to one set of eight ellipses). The relative stability of the category ellipses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

(of each experiment and normalization account). For example, for Nearey's uniform scaling account (Nearey, 1978), we calculated the overall mean of all log-transformed formants over all stimuli. For Lobanov normalization (Lobanov, 1971), we calculated the mean and standard deviation of each formant (in Hz) over all stimuli. For each combination of experiment and normalization account, we then normalized the stimuli using those parameter estimates. The supplementary material, Sec. 3 A 2, summarizes the θ parameters of all normalization accounts for each experiment and how they relate to the values obtained from the training sets. For reasons outlined in that same section, we did not expect a clear relation between an account's ability to predict listeners' responses for an experiment, and the degree to which the account's normalization parameters differed between the experiment and the training database (and, indeed, no such relation was found).

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

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

Finally, we describe the two parameters of the ASP model that we fit against listeners' responses in Experiments 1a and 1b. These two parameters constitute the only DFs that mediate the link from ideal observers' predictions to listeners' responses, and which are fit to listeners' responses. The first DF (Σ_{noise}) models the effects of internal (perceptual) and external (environmental) noise on listeners' perception. While previous work provides estimates of the internal noise in formant perception, these estimates were obtained under *assumptions* about the relevant formant space. For example, Feldman *et al.* (2009) estimated the internal noise variance to be about 15% of the average category variance along F1 and F2. This estimate was based on the assumption that human speech perception transforms vowel formants into Mel, without further normalization. Since we aim to *test* which normalization account best explains speech perception, we cannot rely on this or other

internal noise estimates obtained under a single specific assumption. Additionally, internal noise can vary across individuals and external noise can vary across environments (a point particularly noteworthy, given that we conducted Experiments 1a and 1b over the web). We thus allowed the noise variance Σ_{noise} to vary in fitting participants' responses. Following Feldman *et al.* (2009), we assumed that perceptual noise had identical effects on all formants in the phonetic space defined by the normalization account (see also Kronrod *et al.*, 2016). This reduces Σ_{noise} to a single DF, regardless of the normalization account (for details, see supplementary material, Sec. 3 A 3).

The magnitude of Σ_{noise} affects the slope of the categorization functions that predict listeners' responses from stimulus properties (here, F1 and F2): higher Σ_{noise} imply more shallow categorization slopes. To facilitate comparison of Σ_{noise} values across normalization accounts, we report results in terms of the best-fitting *noise ratios* (τ^{-1}), rather than Σ_{noise} s. Specifically, Σ_{noise} is best understood relative to the inherent variability of the vowel categories (Σ_c). This variability in turn depends on the phonetic space defined by the normalization account. We thus divide Σ_{noise} by the mean of the diagonals of all Σ_c s to obtain the *noise ratio* τ^{-1} . For example, noise ratio of 0 corresponds to the absence of any noise, and a noise ratio of 1 corresponds to noise variance of the same magnitude as the average category variance along F1 and F2 in the phonetic space defined by the normalization account.¹⁰ Figure 8(B) illustrates the effects of this noise ratio for Nearey's uniform scaling account.

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 Fig. 8(C).

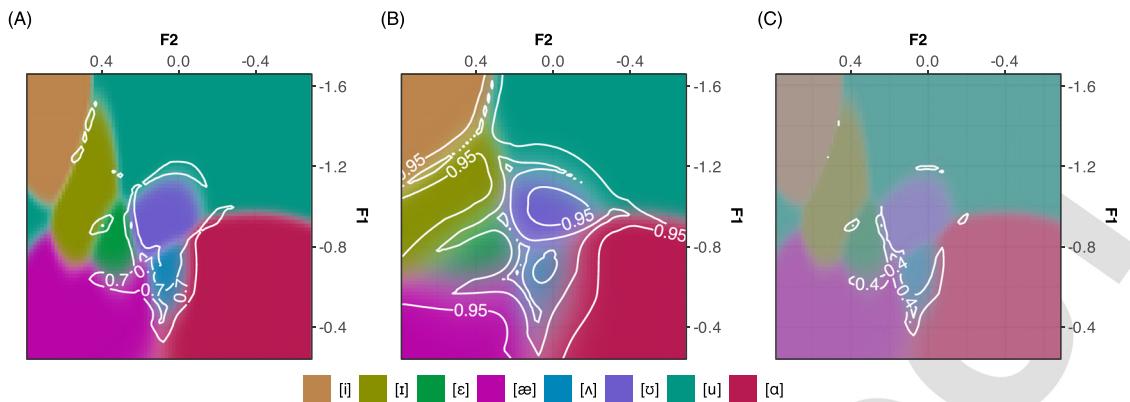


FIG. 8. Illustrating the consequences of perceptual and external noise (Σ_{noise}) and attentional lapse rates (λ) on the predicted posterior distribution of vowel categorizations. Shown are the average predicted posteriors across all five folds for Nearey's uniform scaling account. (A) Predicted posterior distribution for noise ratio $\tau^{-1} = \lambda = 0$. (B) Same for $\tau^{-1} = 1$ and $\lambda = 0$. (C) Same for $\tau^{-1} = 0$ and $\lambda = 0.5$. The transparency of a color is determined by that vowel's posterior probability. Contours indicate the highest posterior probability of any vowel (at 0.4, 0.5, 0.7, 0.95 probability level).

960 Finally, participants can have response biases that reflect
 961 their beliefs about the prior probability of each category.
 962 However, to reduce the DFs fit to participants' responses, we
 963 did *not* fit this response bias against listeners' responses (thus
 964 avoiding $J - 1 = 7$ additional DFs). Instead, we assumed
 965 uniform response biases—i.e., that listeners believed all eight
 966 response options in the experiments to be equally likely
 967 ($\forall c \pi_c = 0.125$). This decision implies that our models would
 968 not be able to capture any potential non-uniformity in listeners'
 969 response biases—including potential effects of additional
 970 acoustic differences (the absence of [h] in *odd* or the coda [t],
 971 rather than [d] in *hut*) and orthographically particular
 972 response options in Experiment 1a (“who would,” “odd,” and
 973 “hut”). We do, however, see no reason to expect this decision
 974 to bias the comparison of normalization accounts.

975 5. Fitting normalization accounts to listeners' 976 responses

977 For each of the 200 combinations of experiment, nor-
 978 malization account, training set, we used constrained quasi-
 979 Newton optimization [Byrd *et al.*, 1995, as implemented in
 980 R's `optim()` function] to find the λ and τ^{-1} values that
 981 best described listener's responses. Specifically, we used the
 982 100 ideal observers described in the previous sections,
 983 applied them to the normalized stimuli of the experiment,
 984 and determined which λ and τ^{-1} maximized the likelihood
 985 of the listener's responses (for details, see supplementary
 986 material, Sec. 3 A 3). This procedure yielded five maximum
 987 likelihood estimates for both λ and τ^{-1} for each combina-
 988 tion of experiment and normalization account—one for each
 989 training set. All results presented in the following were vali-
 990 dated and confirmed by grid searches over the parameter
 991 spaces (see supplementary material, Sec. 3 F).

992 We compare normalization accounts in terms of the like-
 993 lihood of listeners' responses under these maximum likelihood
 994 estimates of λ and τ^{-1} . Comparing accounts in terms of their
 995 data likelihood follows more recent work (e.g., Barreda, 2021;
 996 McMurray and Jongman, 2011; Richter *et al.*, 2017; Xie
 997 *et al.*, 2023). Previous work has instead compared

normalization accounts in terms of their accuracy (e.g., 998
 Johnson, 2020; Nearey and Assmann, 2007; Persson and 999
 Jaeger, 2023), or correlations with human response propor- 1000
 tions (e.g., Hillenbrand and Nearey, 1999; Nearey and 1001
 Assmann, 1986). Both of these approaches are problematic. 1002
 Correlations between the predictions of a model and human 1003
 responses can be high even when the model's predictions are 1004
 systematically “off.” Imagine three items for which listeners 1005
 respond [I] 10%, 30%, and 50% of the time. If a model pre- 1006
 dicts 30%, 50%, and 70% [I] responses, respectively, for the 1007
 same items, its predictions will perfectly correlate with listen- 1008
 ers' response proportions, and yet be systematically wrong. 1009
 Similarly, a model can achieve the highest possible accuracy 1010
 in predicting listeners' responses simply because it always 1011
 predicts the most frequent response (see discussion of criterion 1012
 choice rule in Massaro and Friedman, 1990). In contrast, the 1013
 likelihood of listeners' responses under a model is a direct 1014
 measure of how well the model captures the distribution of lis- 1015
 teners' responses conditional on the stimulus properties. In 1016
 particular, data likelihood will be maximized if, and only if, 1017
 the model-predicted posterior probabilities of each vowel for 1018
 each stimulus are identical to the proportion with which those 1019
 vowels occur in listeners' responses. 1020

B. Results

We begin by comparing the fit of different accounts 1022
 against listeners' responses in Experiments 1a and 1b. Given 1023
 the comparatively large number of accounts compared here, 1024
 we provide initial conclusions based on the best-fitting 1025
 accounts along with the description of the results (more in- 1026
 depth discussion is provided in the general discussion). 1027
 Following this comparison, we visualize how different nor- 1028
 malization accounts predict the formant space to be divided 1029
 into the eight vowel categories. 1030

1. Comparing normalization accounts in terms of fit 1031 against human behavior 1032

Figure 9 compares how well the different normalization 1033
 accounts fit listeners' responses in Experiments 1a and 1b. 1034

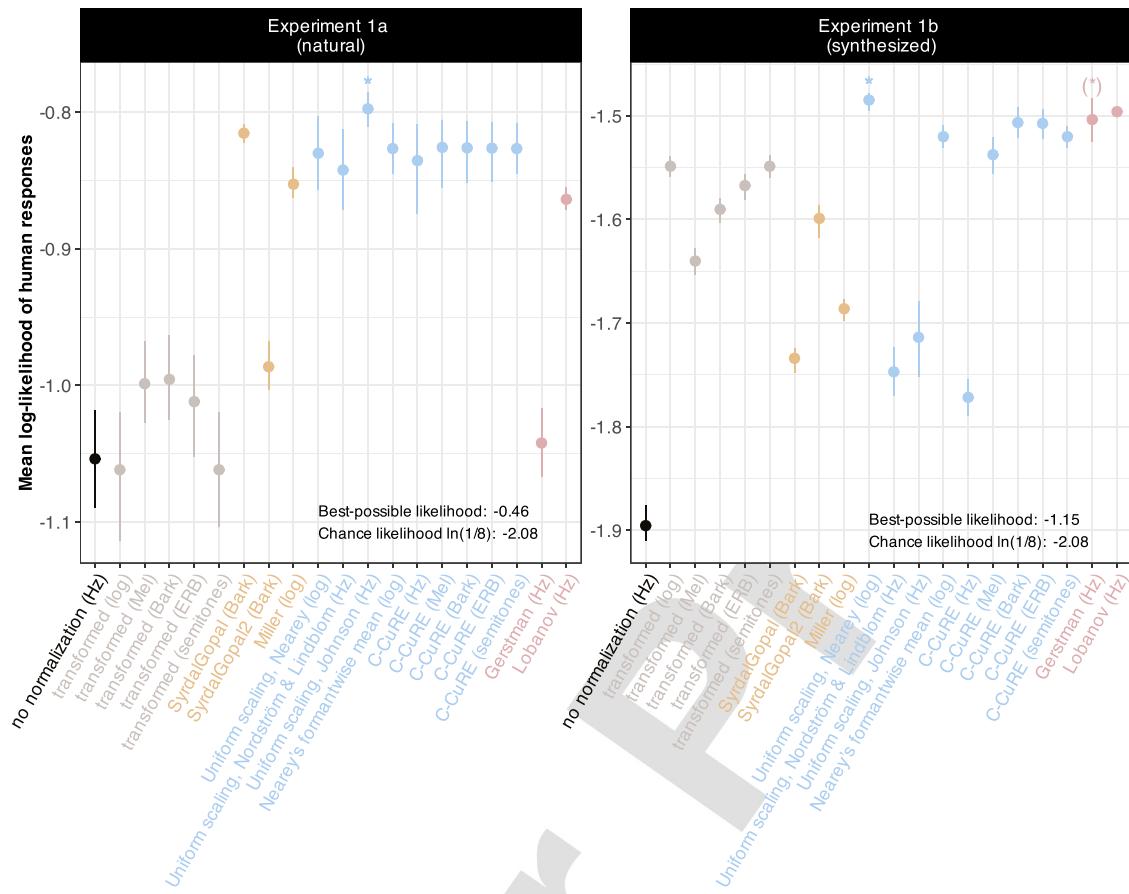


FIG. 9. Comparison of normalization accounts against listeners' responses. Point ranges indicate mean and 95% bootstrapped CIs of the per-token log-likelihoods summarized over the five training sets (higher is better), normalized by the number of listener responses in each experiment. Accounts that fit listeners' responses to an extent that is statistically indistinguishable from the best-fitting account are marked by (*). Note that per-token likelihoods cannot be directly compared across experiments because the best-possible likelihoods differ across experiments (due to differences in stimulus placement and other factors).

1035 All accounts performed well above chance guessing [chance
 1036 per-token log-likelihood in both experiments: $\ln(1/8)$
 1037 $= -2.08$] but also well below the highest possible performance
 1038 (in Experiment 1a, per-token log-likelihood $= -0.46$,
 1039 in Experiment 1b: -1.15).

1040 Normalization significantly improved the fit to listeners'
 1041 responses relative to no normalization. This was con-
 1042 firmed by paired one-sided *t*-tests comparing the maximum
 1043 likelihood values for each normalization account against
 1044 those in the absence of normalization (all $p < 0.05$ except
 1045 for Gerstman normalization, log-transformation, and
 1046 semitones-transformation and Experiment 1a; see supple-
 1047 mentary material Sec. 3 B 1). However, not all normalization
 1048 accounts achieved equally good fits: only some extrinsic
 1049 accounts fit listeners' behavior well across both experi-
 1050 ments. This supports two conclusions. First, it suggests that
 1051 the normalization mechanisms operating during human
 1052 speech perception involve computations that go beyond
 1053 estimation-free transformations into psycho-acoustic spaces.
 1054 Second, it suggests that the input to these computations is
 1055 not limited to intrinsic information—i.e., that the computa-
 1056 tions draw on information beyond what is available in the
 1057 acoustic signal *at that moment*. In particular, extrinsic

1058 normalization requires the estimation and maintenance of talker-specific properties from the speech signal. 1059

1060 While the accounts that achieved the best fit against lis-
 1061 teners' responses differed between experiments, both were
 1062 variants of uniform scaling. For Experiment 1a, Johnson 1063 normalization account provided the best fit (per-token log-
 1064 likelihood $= -0.8$, $SD = 0.02$ across the five cross- 1064 validation folds), while Nearey's uniform scaling account 1065 provided the best fit to Experiment 1b (per-token log- 1066 likelihood $= -1.48$, $SD = 0.01$). Both accounts essentially 1067 slide the representational “template” of a dialect—here the 1068 eight bivariate Gaussian categories of an ideal observer— 1069 along a single line in the formant space. They differ only in 1070 which space this linear relation between formants is 1071 assumed. The same two accounts still fit listeners' responses 1072 best when F3 was included in the analysis in addition to F1 1073 and F2 (see supplementary material, Sec. 3 E).¹¹ This sug- 1074 gests that formant normalization might involve compara- 1075 tively parsimonious maintenance of talker-specific 1076 properties: in its simplest form, uniform scaling employs a 1077 single formant statistic to normalize all formants. In con- 1078 trast, computationally more complex accounts like Lobanov 1079 normalization might require the estimation and maintenance 1080

1081 of two formant statistics (mean and standard deviation) for
 1082 each formant that is normalized (e.g., a total of four formant
 1083 statistics for F1 and F2, or six statistics for F1–F3).

1084 Also, of note is that accounts that were particularly sta-
 1085 ble across experiments operate in log space, whereas
 1086 accounts that operate in Hz space seemed to display a more
 1087 volatile performance (e.g., both standardizing accounts but
 1088 also C-CuRE Hz, Nordström and Lindblom and Johnson
 1089 normalization). That accounts operating over log-
 1090 transformed formants fit human behavior better should not
 1091 be surprising. While questions remain about the exact orga-
 1092 nization of auditory formant representations, it is uncontro-
 1093 versial that the perceptual sensitivity to acoustic frequency
 1094 information is better approximated by a logarithmic scale
 1095 than by a linear scale (see [Moore, 2012](#)). As a result, a
 1096 30 Hz difference in an F1 of 300 Hz (a 10% change) is
 1097 expected to be perceptually more salient than a 30 Hz
 1098 change in an F2 of 2500 Hz (a 1.2% change).¹² In summary,
 1099 variability in how well different accounts predict human
 1100 behavior across the two experiments highlights the impor-
 1101 tance of psycho-acoustic transformations for human speech
 1102 perception. This also highlights the importance of compar-
 1103 ing normalization accounts against multiple types of data.

1104 **2. Visualizing the consequences of different 1105 normalization mechanisms**

1106 Before we turn to the general discussion, we briefly
 1107 visualize how different normalization mechanisms affect
 1108 vowel categorization. This sheds light on *why* the accounts
 1109 differ in how well they fit listeners' responses. Figure 10 vis-
 1110 ualizes the categorization functions predicted by four differ-
 1111 ent normalization accounts, using the best-fitting λ and τ^{-1}
 1112 values for each account (i.e., the values that lead to the fit
 1113 shown in Fig. 9). Figure 10 highlights three points. First, a
 1114 comparison across Figs. 10(A)–10(C) shows different nor-
 1115 malization accounts can result in very different predictions
 1116 about how the acoustic space is carved into categories.

1117 Second, the best-fitting parameters (shown at the top of
 1118 each panel) were relatively comparable across accounts but
 1119 differed more substantially across experiments. Specifically,
 1120 the best-fitting estimates of lapse rates λ were generally
 1121 comparable across the two experiments (with the exception
 1122 of Nordström & Lindblom and Johnson normalization,
 1123 which exhibited substantially higher lapse rates in
 1124 Experiment 1b; see supplementary material, Sec. 3 B 2).
 1125 This suggests that participants in both experiments were
 1126 about equally likely to pay attention to the stimulus. The
 1127 best-fitting noise ratios τ^{-1} , however, differed substantially
 1128 across experiments and were ten times larger for
 1129 Experiment 1b (mean $\tau^{-1} = 4.32$, SD = 2.52 across normali-
 1130 zation accounts) than for Experiment 1a (mean $\tau^{-1} = 0.42$,
 1131 SD = 0.46). This difference most likely reflects the fact that
 1132 the synthesized stimuli in Experiment 1b left listeners with
 1133 substantially more uncertainty about the intended category,
 1134 as discussed during the description of the experiments.

1135 Since noise is assumed to be independent of category
 1136 variability (see also [Feldman et al., 2009](#); [Kronrod et al.,](#)

1137 differences in noise ratios can substantially change 1137
 1138 the categorization function. This is particularly evident for 1138
 1139 the accounts that had more variable performance across the 1139
 1140 two experiments. For example, Johnson normalization 1140
 1141 [Fig. 10(B)] resulted in very different best-fitting categoriza- 1141
 1142 tion functions for Experiments 1a and 1b. 1142

1143 Third and finally, Fig. 10 also shows how well accounts 1143
 1144 fit listeners' responses for each test stimulus (opaqueness 1144
 1145 of the points). This begins to explain *why* some accounts 1145
 1146 fit listeners' responses in Experiment 1b less well. For 1146
 1147 example, the Johnson normalization account [Fig. 10(B)] 1147
 1148 predicts the responses to the test stimuli in Experiment 1a 1148
 1149 well but fails to predict the responses to the test stimuli 1149
 1150 in Experiment 1b. This drop in performance seems to be 1150
 1151 primarily driven by stimuli that are unlikely to be 1151
 1152 articulated by the same talker (lower left, cf. dashed line in 1152
 1153 Fig. 4). This might suggest that this account was over- 1153
 1154 engineered to explain naturally occurring productions—the 1154
 1155 type of data, it was originally tested on ([Johnson, 2020](#)). A 1155
 1156 plausible account of normalization, however, should be able 1156
 1157 to explain human perception to any type of stimulus, includ- 1157
 1158 ing synthesized stimuli. The supplementary material, Sec. 1158
 1159 3 B 3, presents more detailed by-item comparisons of nor- 1159
 1160 malization accounts that might be of interest to some 1160
 1161 readers. 1161

1162 **IV. GENERAL DISCUSSION**

1163 Research on vowel normalization has an influential his- 1163
 1164 tory. Cognitive scientists have long aimed to understand the 1164
 1165 organization of frequency information in the human brain 1165
 1166 ([Siegel, 1965](#); [Stevens and Volkmann, 1940](#)), and how it 1166
 1167 helps listeners overcome cross-talker variability in the for- 1167
 1168 mant-to-vowel mapping (e.g., [Fant, 1975](#); [Joos, 1948](#); [Nordström and Lindblom, 1975](#)). Auditory processes that 1168
 1169 normalize speech inputs for differences in vocal tract physi- 1169
 1170 ology are now recognized to be an integral part of speech 1171
 1171 perception ([Johnson and Sjerps, 2021](#); [McMurray and Jongman, 2011](#); [Xie et al., 2023](#)). Here, we set out to inves- 1172
 1173 tigate what types of computations are implicated in the nor- 1174
 1175 malization of the frequency information that plays a critical 1175
 1176 role in the recognition of vowels. 1176

1177 Our results support three theoretical insights. First, 1177
 1178 human speech perception draws on more than psycho- 1178
 1179 acoustic transformations or intrinsic information, in line 1179
 1180 with previous research on normalization ([Adank et al., 2004](#); 1180
 1181 [Ladefoged and Broadbent, 1957](#); [Nearey, 1989](#)). 1181
 1182 Rather, formant normalization seems to involve the estima- 1182
 1183 tion and storing of talker-specific formant properties. 1183
 1184 Second, computationally simple uniform scaling accounts 1184
 1185 provide the best fit to listeners' responses, suggesting com- 1185
 1186 paratively parsimonious maintenance of talker-specific 1186
 1187 properties. This replicates and extends previous findings that 1187
 1188 uniform scaling or similarly simple corrections for vocal 1188
 1189 tract size provide a better explanation for human perception 1189
 1190 than more complex extrinsic accounts ([Barreda, 2021](#); [Richter et al., 2017](#)). It is impossible to rule out more 1191

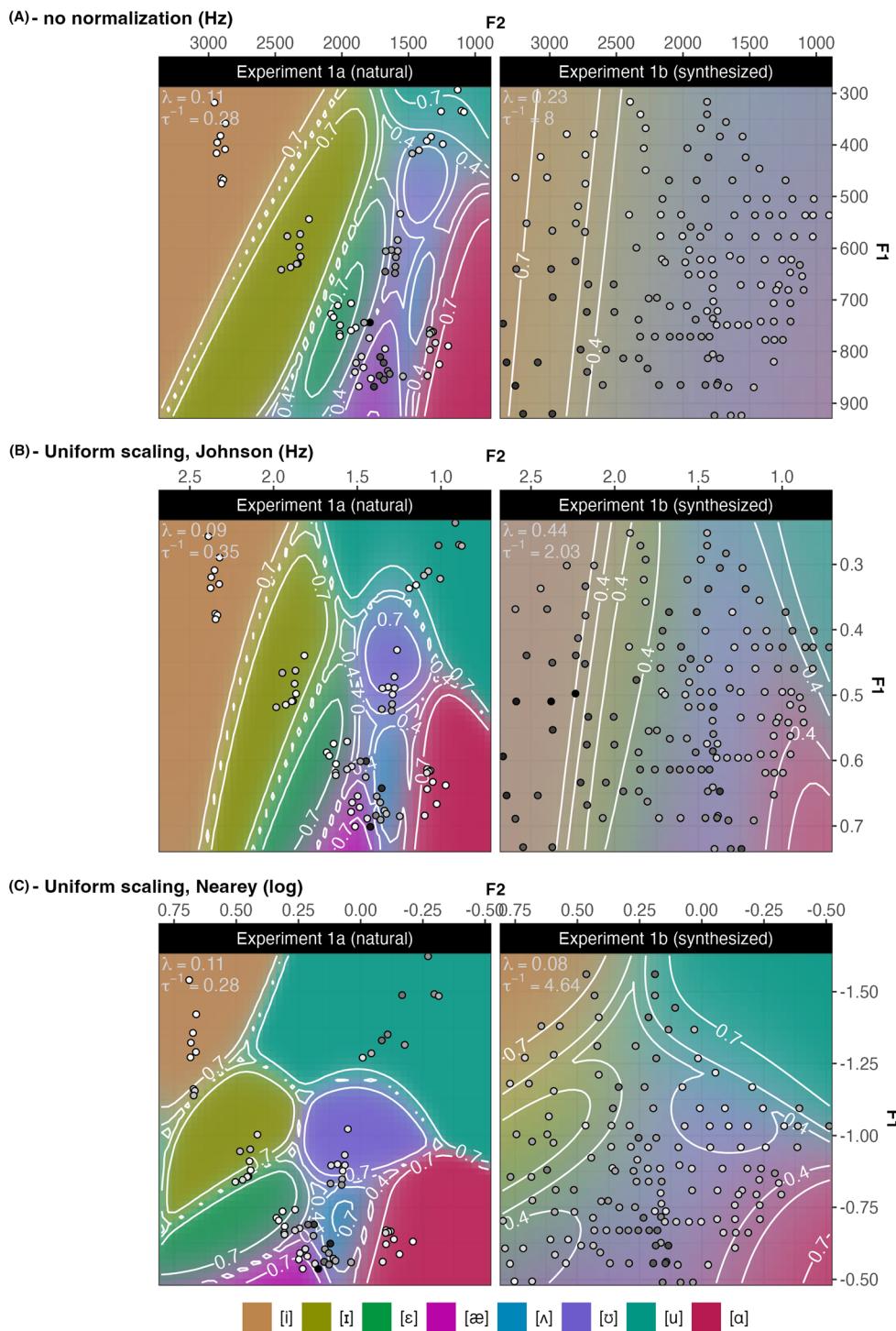


FIG. 10. Predicted categorization functions over the F1–F2 space under three different normalization accounts. For each account, we show the predicted posterior probabilities of all eight vowels obtained by averaging over the maximum likelihood parameterizations (of λ and τ^{-1}) for the five training sets (shown at the top of each panel). (A) Absence of normalization shown for reference. (B) The best-fitting account for Experiment 1a. (C) the best-fitting account for Experiment 1b. Contours indicate the highest posterior probability of any vowel. Points indicate the location of test stimuli. The increasing brightness of points indicates a better match between the account's prediction and listeners' responses (higher log-likelihood; see text for detail).

1192 complex approaches to perceptual normalization given the
 1193 large number of possible alternatives. However, given that
 1194 uniform scaling provides a parsimonious explanation for
 1195 human formant normalization, and the current absence of
 1196 empirical evidence for more complex computations, we sub-
 1197 mit that researchers ought to adapt uniform scaling as the

1198 working hypothesis. Third, the psycho-acoustic representation
 1199 assumed by different normalization accounts matters, as indi-
 1200 cated by the comparison of otherwise computationally similar
 1201 accounts (e.g., Nearey's vs Johnson's uniform scaling).

1202 The results contribute to a still comparatively small
 1203 body of work that has evaluated competing normalization

1204 accounts against listeners' perception, whereas most previous
 1205 work evaluates accounts against intended productions.
 1206 Complementing previous work, we took a broad-coverage
 1207 approach: the present study compared 20 of the most influential
 1208 normalization accounts against listeners' perception of *hVd* words with eight US English monophthongs in both
 1209 natural and synthesized speech. This contrasts with previous
 1210 work, which has typically focused on subsets of the vowel
 1211 system, either using natural *or* synthesized speech and consider-
 1212 ing a much smaller subset of accounts (typically 2–3 at a time). By considering a wider range of accounts, a wider
 1213 range of formant values and vowel categories, and multiple
 1214 types of speech, we aimed to contribute to a more comprehensive
 1215 evaluation of competing accounts.

1216 Next, we discuss the theoretical consequences of these
 1217 findings for research beyond formant normalization.
 1218 Following that, we discuss the limitations of the present
 1219 work, and how future research might overcome them.

1220 A. Consequences for theories of speech perception 1221 and beyond

1222 Understanding the perceptual space in which the human
 1223 brain represents vowel categories—i.e., the normalized for-
 1224 mant space—has obvious consequences for research on
 1225 speech perception. To illustrate how far-reaching these
 1226 consequences can be, we discuss a few examples. For instance,
 1227 research on *categorical perception* has found that vowels
 1228 seem to be perceived less categorically than some types of
 1229 consonants. Recent work has offered an elegant explanation
 1230 for this finding: the perception of formants—relevant to the
 1231 recognition of vowels—might be more noisy than the per-
 1232 ception of the acoustic cues that are critical to the recogni-
 1233 tion of more categorically perceived consonants (Kronrod
 1234 *et al.*, 2016). This is a parsimonious explanation, potentially
 1235 preempting the need for separate explanations for the per-
 1236 ception of different types of phonemic contrasts. Kronrod
 1237 *et al.* based their argument on estimates they obtained for
 1238 the relative ratio of meaningful category variability to per-
 1239 ceptual noise (τ , the inverse of our noise ratios, τ^{-1}).
 1240 Critically, this ratio depends both on (i) the perceptual space
 1241 in which formants are assumed to be represented (Kronrod
 1242 *et al.* used Mel-transformed formant frequencies), and on
 1243 (ii) whether the meaningful category variability is calculated
 1244 prior to, or following, normalization (Kronrod *et al.*
 1245 assumed the former, which increases estimates of category
 1246 variability). Our point here is not to cast doubt on the results
 1247 of Kronrod *et al.* (2016)—the fact that the best-fitting noise
 1248 ratios in our study were relatively similar across accounts
 1249 (while varying across experiments) suggests that the result
 1250 of Kronrod and colleagues are likely to hold even under dif-
 1251 ferent assumptions about (i) and (ii)—but rather to highlight
 1252 how research on the perception and recognition of vowels
 1253 depends on assumptions about formant normalization. For
 1254 example, similar points could be raised about experiments
 1255 on statistical learning that manipulate formant or other fre-
 1256 quency statistics (e.g., Chládková *et al.*, 2017; Colby *et al.*,
 1257 2018; Wade *et al.*, 2007; Xie *et al.*, 2021). Such

1258 experiments, too, need to make assumptions about the space 1259 in which formants are represented. If these assumptions are 1260 incorrect, this can affect whether the experimental manipu- 1261 lations have the intended effects, increasing the chance of 1262 null effects or misinterpretation of observed effects. 1263

1264 Understanding the perceptual space in which the human 1265 brain represents vowel categories also has consequences for 1266 research beyond speech perception, perhaps more so than is 1267 sometimes recognized. For instance, in sociolinguistics and 1268 related fields, Lobanov remains the norm for representing 1269 vowels due to its efficiency in removing cross-talker vari- 1270 ability (for review, see Adank *et al.*, 2004; Barreda, 2021). 1271 However, as shown in the present study, removing cross- 1272 talker variability is not the same as representing vowels in 1273 the perceptual space that listeners actually employ. Here, we 1274 do *not* find Lobanov to describe human perception particu- 1275 larly well. On the contrary, we find no support for the 1276 hypothesis that human speech perception employs these 1277 more complex computations that have been found to per- 1278 form best at reducing category variability. This should 1279 worry sociolinguists. In order to understand how listeners 1280 infer a talker's background or social identity, it is important 1281 to understand the perceptual space in which inferences are 1282 actually rooted. Critically, the representations resulting from 1283 formant normalization presumably form an important part 1284 of the information that listeners use to draw social and lin- 1285 guistic inferences. It should thus be obvious that the use of 1286 normalization accounts that do not actually correspond to 1287 human perception can both mask real markers of social 1288 identity and "hallucinate" markers that are not actually pre- 1289 sent. For example, in order to determine how a talker's 1290 social identity influences their vowel realizations, it is 1291 important to discount *all and only* effects that listeners will 1292 attribute to physiology, rather than social identity (Disner, 1293 1980; Hindle, 1978). 1294

1295 Similar concerns apply to dialectology, research on lan- 1296 guage change, second language acquisition research, etc. 1297 For example, the perceptual space in which vowels are rep- 1298 resented is critical to well-formed tests of hypotheses about 1299 the factors shaping the organization of vowel inventories 1300 across languages of the world (Lindblom, 1986; Stevens, 1301 1972, 1989). It is essential in testing hypotheses about the 1302 extent to which the cross-linguistic realization of those sys- 1303 tems is affected by perceptual processes (Flemming, 2010; 1304 Steriade, 2008), or by preferences for communicatively effi- 1305 cient linguistic systems (e.g., Hall *et al.*, 2018; Lindblom, 1306 1990; Moulin-Frier *et al.*, 2015). Similarly, tests of the 1307 hypothesis that vowel *articulation* during natural interac- 1308 tions is shaped by communicative efficiency do in obvious 1309 ways depend on assumptions about the perceptual space in 1310 which talkers—by hypothesis—aim to reduce perceptual 1311 confusion (cf. Buz and Jaeger, 2016; Gahl *et al.*, 2012; 1312 Scarborough, 2010; Wedel *et al.*, 2018). The same applies 1313 to any other line of research that aims to understand the per- 1314 ceptual consequences of formant variation across talkers, 1315 including research on infant- or child-directed speech 1316 (Eaves *et al.*, 2016; Kuhl *et al.*, 1997), and research on 1317

1317 whether non-native talkers are inherently more variable than
 1318 native talkers (Smith *et al.*, 2019; Vaughn *et al.*, 2019; Xie
 1319 and Jaeger, 2020). In short, the perceptual space in which
 1320 vowels are represented is a critical component of under-
 1321 standing the structure of vowel systems, the factors that
 1322 shape them, and the ways in which they are used in natural
 1323 language.

1324 B. Limitations and future directions

1325 As mentioned in Sec. I, we take it as relatively uncon-
 1326 troversial *that* normalization is part of human speech per-
 1327 ception. Independent of any benefits that such normalization
 1328 conveys for speech perception, its existence is supported by
 1329 evidence from cross-species comparisons and neuro-
 1330 physiologically studies (for review, see Barreda, 2020). There
 1331 are, however, important questions as to how decisions we
 1332 made in comparing normalization accounts against each
 1333 other might have affected their fit against listeners'
 1334 responses.

1335 For instance, we followed previous work in focusing on
 1336 formants, and specifically estimates of the formants in the
 1337 *center* of the vowel. There is, however, ample evidence that
 1338 formant dynamics throughout the vowel can strongly affect
 1339 perception (Assmann and Katz, 2005; Hillenbrand and
 1340 Nearey, 1999; Nearey and Assmann, 1986). In addition,
 1341 there are proposals that entirely give up the assumption that
 1342 formants are the primary cues to vowel identity (e.g., whole-
 1343 spectrum accounts, Hillenbrand *et al.*, 2006). While these
 1344 proposals might provide a more informative representation
 1345 of vowels, we consider it unlikely that they would entirely
 1346 remove the problem of cross-talker variability. For instance,
 1347 Richter *et al.* (2017) still found benefits of normalization
 1348 even when the entire frequency spectrum throughout vowels
 1349 was considered (in the form of Mel-frequency cepstral coef-
 1350 ficients and their derivatives). For the present work, auxil-
 1351 iary analyses in the supplementary material, Sec. 3 E
 1352 replicated our core findings when F3 was included in the
 1353 model. Still, it remains unclear whether the inclusion of
 1354 additional cues, such as VISC, or additional formant dyna-
 1355 mics, would alter the results of the present study.

1356 As is the case of any computational work, the present
 1357 study committed to a number of assumptions that are not
 1358 critical, but were necessary in order to deliver clear quantita-
 1359 tive predictions. Quantitative tests of theories—as we have
 1360 done here—require assumptions about *every* aspect of the
 1361 model. Here, this included all the steps necessary to link
 1362 properties of the stimuli to listeners' responses. For this pur-
 1363 pose, we adopted the ASP framework (Xie *et al.*, 2023), and
 1364 visualized the graphical model that links stimuli (x) to
 1365 responses (r) in Fig. 6.

1366 Many of the assumptions we made should be relatively
 1367 uncontroversial—e.g., the decision to include both external
 1368 (environmental) and internal (perceptual) noise in our
 1369 model. While these noise sources are often ignored in
 1370 modeling human behavior, it is uncontroversial that they
 1371 exist. Other assumptions we made were introduced as

1372 simplifying assumptions for the sake of feasibility—e.g., we
 1373 expressed the effect of both types of noise through a single
 1374 parameter that related the average within-category variabil-
 1375 ity of formants to noise variability in the transformed and
 1376 normalized formant space. In reality, however, environment
 1377 noise can have effects that are independent of internal noise,
 1378 and internal noise likely affects information processing at
 1379 multiple (or all) of the steps shown in Fig. 6. Such simplify-
 1379 ing assumptions are both inevitable and not necessarily
 1380 problematic: as long as they do not introduce systematic
 1381 bias to the evaluation of normalization accounts, they should
 1382 not limit the generalizability of our results.
 1383

1384 Some of our assumptions, however, might be more con-
 1385 troversial. For example, we assumed that category represen-
 1386 tations can be expressed as multivariate Gaussian
 1387 distributions in the formant space. This assumption, too, is a
 1388 simplifying assumption—it simplified the computation of
 1389 likelihoods—rather than a critical feature of the ASP frame-
 1390 work we employed. While human category representations
 1391 are unlikely to be Gaussians, the alternative, e.g., exemplar
 1392 representations, would come with its own downsides, such
 1393 as increased sensitivity to the limited size of phonetic data-
 1394 bases and substantial increases in computation time (exem-
 1395 plar representations afford researchers with much larger
 1396 degrees of freedom). For researchers curious how this and
 1397 other assumptions we made affect our results, our data, and
 1398 code are shared on OSF.

1399 Like previous work, we further assumed that all listen-
 1400 ers in our experiments use the same underlying vowel repre-
 1401 sentations—the same dialect template(s). However, as
 1402 already discussed, it is rather likely that not all of our listen-
 1403 ers employed the same dialect template(s). An additional
 1404 analysis reported in the supplementary material, Sec. 3 D,
 1405 thus compared normalization accounts against only the sub-
 1406 set of listeners who employed the dialect template used by
 1407 the majority of participants [see lower-left of Fig. 5(B)].
 1408 This left only 20 participants for Experiment 1a (71.4%) and
 1409 23 for Experiment 1b (82.1%), substantially reducing statis-
 1410 tical power. Replicating the main analysis, uniform scaling
 1411 accounts again fit listeners' behavior well across both
 1412 experiments. The best-performing account for Experiment
 1413 1a did, however, differ from the one obtained for the super-
 1414 set of data (the intrinsic Syrdal and Gopal achieved the best
 1415 fit to listeners' responses in Experiment 1a for the shared
 1416 dialect subset; see supplementary material, Sec. 3 D).
 1416

1417 A related assumption was introduced by the use of a
 1418 phonetic database to approximate listeners' vowel represen-
 1419 tations. This deviates from most previous evaluations of nor-
 1420 malization accounts (McMurray and Jongman, 2011; 1420
 1421 Barreda, 2021; but see Richter *et al.*, 2017), and reflects our
 1422 commitment to a strong assumption made by most theories
 1423 of speech perception: that listeners' representations reflect
 1424 the formant statistics previously experienced speech input.
 1424 By using a phonetic database to estimate listeners' represen-
 1425 tations, we *substantially* reduced the degrees of freedom in
 1426 the evaluation of normalization accounts, reducing the
 1427 chance of over-fitting the data from our experiments. Our
 1428

1429 approach does, however, also introduce two new
1430 assumptions.

1431 First, our approach assumes that the mixture of dialect
1432 template(s) used by talkers in the database sufficiently
1433 closely approximates those of the listeners in our experi-
1434 ments. Some validation for this assumption comes from the
1435 additional analysis reported in the preceding paragraph:
1436 when we subset listeners to only those who used the major-
1437 ity dialect template, this improved the fit of all normaliza-
1438 tion accounts—as expected, if the category representations
1439 we trained on the phonetic database primarily reflect those
1440 listeners’ representations (see supplementary material, Sec.
1441 3 D). Future work could further address this assumption in a
1442 number of ways. On the one hand, dialect analyses like the
1443 ones we presented for our listeners [in Fig. 5(B)] could com-
1444 pare listeners’ templates against the templates used by talk-
1445 ers in the database. Alternatively or additionally, researchers
1446 could see whether our results replicate if ideal observers are
1447 instead trained on other databases that have been hypothe-
1448 sized to reflect a “typical” L1 listeners’ experience with US
1449 English. Finally, it might be possible in future work to use
1450 larger databases of vowel recordings to train separate ideal
1451 observers for all major dialects of US English, and to try to
1452 estimate for each listener which mixture of dialects their
1453 responses are based on.

1454 Second, we made the simplifying assumption that
1455 listeners’ category representation—or at least the represen-
1456 tations listeners’ drew on during the experiment—are talker-
1457 independent (we trained a single set of multivariate
1458 Gaussian categories, rather than, e.g., a hierarchically
1459 organized set of multiple dialect templates). While this
1460 assumption is routinely made in research on normalization
1461 and beyond, it might well be wrong (see, e.g., Xie *et al.*,
1462 2021).

1463 Finally, the evaluation of normalization accounts in the
1464 present study shares with all previous work (e.g.,
1465 Apfelbaum and McMurray, 2015; Barreda, 2021; Cole

et al., 2010; McMurray and Jongman, 2011; Nearey, 1989; Richter *et al.*, 2017) another simplifying assumption that is clearly wrong: the assumption that listeners *know* the talker-specific formant properties required for normalization. Specifically, we normalized the input for each ideal observer using the maximum likelihood estimates of the normalization parameters over all stimuli for the respective experiment. For example, for the evaluation of the ideal observer trained on Lobanov normalized formants against listeners’ responses in Experiment 1a, we used the formant means and standard deviations of the stimuli used in Experiment 1a to normalize F1 and F2. While this follows previous work, it constitutes a problematic assumption for the evaluation of extrinsic normalization accounts. For extrinsic accounts, the approach adopted here would seem to entail the ability to predict the future: even on the first trial of the experiment, the input to the ideal observers were formants that were normalized based on the normalization parameters estimated over the acoustic properties of *all* stimuli. Listeners instead need to *incrementally infer* talker-specific properties from the speech input (Barreda and Jaeger, submitted; Nearey and Assmann, 2007; Xie *et al.*, 2023). An important avenue for future research is thus the development and evaluation of incremental normalization accounts.

AQ9

The present data only allow an initial, rather tentative, look at this question. For example, for Experiment 1a, for which each trial had a known correct answer (the vowel intended by the talker), we can assess whether participants’ recognition accuracy improved across trials, as would be expected if listeners need to incrementally infer the talker-specific normalization parameters. Figure 11(A) suggests that this was indeed the case: the non-parametric listeners’ average recognition accuracy improved over the course of the experiment from about 65% to 88%, with most of the improvements occurring during the first ten trials. To address potential confounds due to differences in the distribution of stimuli across trials, we used a generalized

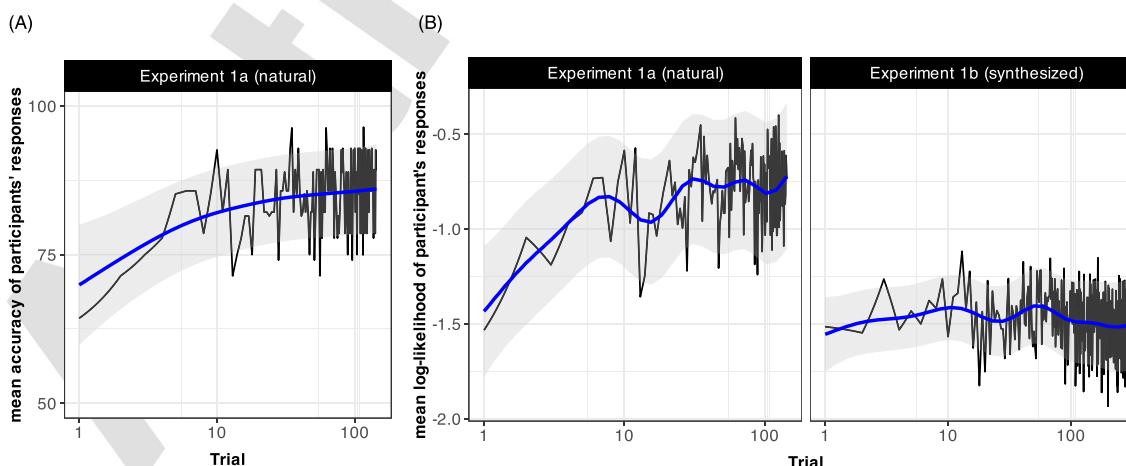


FIG. 11. (A) Changes across trials in listeners’ average accuracy in recognizing the vowel intended by the talker in Experiment 1a, averaged across items and participants (black line). The blue line shows a generalized additive mixed-effects model predicting accuracy from log-transformed trial order, with 95% CIs. (B) Log-likelihood of listeners’ responses under the best-fitting normalization account at each trial, averaged across items and participants (Johnson’s uniform scaling for Experiment 1a and Nearey’s uniform scaling for Experiment 1b). Blue lines show generalized additive mixed-effects models predicting log-likelihood from log-transformed trial order, with 95% CIs.

1503 additive mixed-effect model to predict listeners' accuracy
 1504 from log-transformed trial order while accounting for ran-
 1505 dom by-participant and by-item intercepts and slopes for the
 1506 log-transformed trial order (blue lines). Still, this result
 1507 should be interpreted with caution, as Experiment 1a was
 1508 not designed to reliably address questions about incremental
 1509 changes across the experiment.

1510 Figure 11(B) shows how the fit of the best-fitting
 1511 normalization model changes across trials. We used a gener-
 1512 alized additive mixed-effect model to predict the log-
 1513 likelihood of listeners' responses from log-transformed trial
 1514 order while accounting for random by-participant and by-
 1515 item intercepts and slopes for the log-transformed trial order
 1516 (blue lines). Given that our evaluation of normalization
 1517 accounts assumed that the normalization parameters were
 1518 already known on the first trial of the experiment, we would
 1519 expect that the likelihood of listeners' responses under a
 1520 normalization model would improve the more input listeners
 1521 have received (i.e., as the simplifying assumptions of our
 1522 evaluation become increasingly more plausible). For
 1523 Experiment 1a, this indeed appears to be the case. However,
 1524 no clear evidence for such incremental improvements in the
 1525 fit of the normalization model is observed for Experiment
 1526 1 b. In short, the present data does not support decisive con-
 1527 clusions about the extent to which normalization proceeds
 1528 incrementally.

1529 C. Concluding remarks

1530 We set out to compare how well competing accounts of
 1531 formant normalization explain listeners' perception of vowels.
 1532 We developed a computational framework that makes it
 1533 possible to compare a large number of different accounts
 1534 against multiple data sets. The code we share on OSF makes
 1535 it possible to "plug in" different accounts of vowel normali-
 1536 zation, different phonetic databases, and different perception
 1537 experiments. This, we hope, will substantially reduce the
 1538 effort necessary to conduct similar evaluations on other
 1539 datasets, dialects, and languages.

1540 Comparing 20 of the most influential normalization
 1541 accounts against L1 listeners' perception of US English
 1542 monophthongs, we found that the normalization accounts
 1543 that best describe listeners' perception share that they (1)
 1544 learn and store talker-specific properties and (2) seem to be
 1545 computationally very simple—taking advantage of the phys-
 1546 ics of sound generation to use as few as a single parameter
 1547 to normalize inter-talker variability in vocal tract size.
 1548 While the number of studies that have compared normaliza-
 1549 tion accounts against *listeners'* behavior remains surpris-
 1550 ingly small, these two results confirm the findings from
 1551 more targeted comparisons that were focused on 2–3
 1552 accounts at a time (Barreda, 2021; Nearey, 1989; Richter
 1553 et al., 2017). Overall then, we submit that it is time for
 1554 research in speech perception and beyond to consider simple
 1555 uniform scaling the most-likely candidate for human for-
 1556 mant normalization.

SUPPLEMENTARY MATERIAL

1557

See the supplementary material for more details on 1558 participant and experiment data, on the vowel database 1559 used, and on the computational models presented in 1560 the paper. The supplementary material also contains addi- 1561 tional auxiliary analyses, including models trained on 1562 different subsets of the data, and on additional cues besides 1563 F1–F2. 1564

ACKNOWLEDGMENTS

1565

Earlier versions of this work were presented at the 2023 1566 ASA Meeting, ExLing 2022, at the Department of 1567 Computational Linguistics at the University of Zürich and 1568 the Department of Swedish Language and Multilingualism 1569 at Stockholm University. We are grateful to Maryann Tan, 1570 Chigusa Kurumada, and Xin Xie for their feedback on this 1571 work. We thank Travis Wade for clarifications on the 1572 synthesis procedure used in his study. We thank Leslie Li 1573 and Xin Xie for sharing their database of L1-US English 1574 *hVd* productions, and the JASA copy editing staff for 1575 help with the Latex formatting. This work was partially 1576 funded by grants to A.P. from Kungliga 1577 Vetenskapsakademien, Kungliga Vitterhetsakademien, and 1578 the Department of Swedish Language and Multilingualism 1579 at Stockholm University, as well as Grants to T.F.J. by the 1580 Helge Ax:son Johnson foundation, the Stockholm 1581 University Board of Human Science (Funding for Strategic 1582 Investments), and the Stockholm University Faculty of 1583 Humanities' Research School (Kvalitetssäkrande medel 1584 grant). A.P. designed the experiments and collected the data, 1585 with input from T.F.J. T.F.J. programmed the experiments 1586 with input from A.P. A.P. analyzed the experiments, with 1587 input from T.F.J. A.P. and T.F.J. wrote the code to 1588 implement and fit the normalization models, with input from 1589 S.B. A.P. developed the visualizations within input from 1590 S.B. and T.F.J. A.P. wrote the first draft of the manuscript 1591 with edits by S.B. and T.F.J. 1593

1594

AUTHOR DECLARATIONS

1595

Conflict of Interest

1596

The authors have no conflicts to disclose. 1597

Ethics approval

1598

This study was reviewed and approved Research 1599 Subjects Review Board (RSRB) of the University of 1600 Rochester (STUDY00000417) under the OHSP and UR 1601 policies, and in accordance with Federal regulation 45 CFR 1602 46 under the university's Federal-Wide Assurance 1603 (FWA00009386). 1604

¹Some hypotheses hold that robust speech perception does not require 1605 normalization, and that research on normalization has over-estimated 1606 its effectiveness because studies tend to consider only a fraction of 1607 the phonetic information available to listeners (for review, see Strange 1608 and Jenkins, 2012). For vowel recognition, for example, listeners might 1609 use cues other than just formants (Hillenbrand et al., 2006; Nearey and 1610

1611 **Assmann, 1986**), and/or might use information about the dynamic development of formant trajectories over the entire vowel rather than just point estimates of formants at the vowel center (e.g., **Shankweiler et al., 1978**). We return to this in the general discussion but note that even studies that use much richer inputs have found that normalization provides a better fit for listeners' perceptions (**Richter et al., 2017**).

1612 ²Under uniform scaling accounts, listeners essentially "slide" the center of their category representations (e.g., the "template" of vowel categories for a given dialect) along a single line in formant space, with Ψ determining the target of this sliding. Later extensions of this account maintain its memory parsimony but increase its inference complexity by allowing both intrinsic (the current F0) and extrinsic information (the talker's single mean of log-transformed formants) to influence the inference of Ψ (**Nearey and Assmann, 2007**).

1613 ³We use the **Johnson (2020)** implementation of **Nordström and Lindblom (1975)**. We group both **Nordström and Lindblom (1975)** and **Johnson (2020)** with the centering accounts, as they are essentially variants of uniform scaling, differing in their estimation of Ψ . We also include both versions of Syrdal and Gopal's Bark-distance model. The two versions differ only in their normalization of F2 and have not previously been compared against human perception.

1614 ⁴Shannon (1948) response entropy is defined as $H(x) = -\sum_{i=1}^n P(x_i) \log P(x_i)$. The maximum possible response entropy for an eight-way response choice is three bits, which means that all eight vowels are responded to equally often. The minimum response entropy = 0 bits, which means that the same vowel is responded to all the time.

1615 ⁵Note that participants in Experiment 1a exhibited high agreement on [ʌ], [æ], and [ɑ], despite the close proximity between, and partial overlap of, these vowels in F1–F2 space. To understand this pattern, it is important to keep in mind that the recordings for [ʌ] and [ɑ] differed from the recordings for other stimuli in their word onset ("odd" for [ɑ]) or offset ("hut" for [ʌ]).

1616 ⁶[u] has been undergoing changes in many varieties of US English. Whereas the talker in Experiment 1a produces [u] with low F1 and F2 (high and back), other L1 talkers of US English produce this vowel considerably more forward (higher F2).

1617 ⁷For Gaussian noise and Gaussian category likelihoods, the resulting noise-convolved likelihood is a Gaussian with variance equal to the sum of the noise and category variances (**Kronrod et al., 2016**).

1618 ⁸We intentionally did *not* split the data among talkers since normalization accounts are meant to make speech perception robust to cross-talker variability. Further, splitting the data by speaker rather than by vowel category avoids the potential for biases in the normalization parameter estimates for different speakers in the case of missing or unbalanced tokens across vowel categories, see (**Barreda and Nearey, 2018**). Additional analyses not reported here confirmed that the same results are obtained when splits are performed within talkers and within vowels (except that this leads to smaller CIs, and thus *more* significant differences, in Fig. 9). These analyses can be replicated by downloading the R markdown document this article is based on from our OSF (see comments in our code).

1619 ⁹Alternatively, it would be possible to treat these parameters as DFs in the link to listeners' responses, and infer them from the responses in Experiments 1a and 1b (cf. **Kleinschmidt and Jaeger, 2016**). This approach would afford the model a high degree of functional flexibility, regardless of which normalization approach is applied (similar to previous approaches that have been employed, e.g., multinomial logistic regression).

1620 ¹⁰This ratio is a generalization of the inverse of the "meaningful-to-noise variance ratio (τ)" used in **Kronrod et al. (2016)**. However, whereas Kronrod and colleagues committed to the simplifying assumption that all categories have identical variance (along all formants), we allowed category variances to differ between vowels, and between F1 and F2 (matching the empirical facts). We merely assume that the *noise* variance is identical across all formants (in the phonetic space defined by the normalization account, e.g., log-Hz for uniform scaling and Hz for Lobanov).

1621 ¹¹Additional analyses reported in the supplementary material, Sec. 3C, overall replicated this result for subsets of Experiments 1a and 1b, with Nearey's uniform scaling achieving the best fit to listeners' responses in both experiments. For Experiment 1a, we excluded responses to the two *hVd* stimuli that differed from the other stimuli in the preceding (*odd*) or

1622 following phonological context (*hut*). For Experiment 1b, we excluded responses to any stimuli that were physiologically implausible for the talker (stimuli below the diagonal dashed line in Fig. 4). As requested by a reviewer, the supplementary material, Sec. 3B 4 also reports the accuracy of predicting listeners' responses for all normalization accounts. The best-performing accounts achieved 61.8% for Experiment 1a (Johnson normalization), and 29.2% for Experiment 1b (Nearey's uniform scaling), compared to 52.3% and 16.9%, respectively, without normalization.

1623 ¹²In line with this reasoning, additional tests found that Johnson normalization would provide the best fit to Experiment 1b if it was applied to log-transformed formants (instead of Hertz).

1624 Abramson, A. S., and Lisker, L. (1973). "Voice-timing perception in Spanish word-initial stops," *J. Phon.* **1**(1), 1–8.

1625 Adank, P., Smits, R., and van Hout, R. (2004). "A comparison of vowel normalization procedures for language variation research," *J. Acoust. Soc. Am.* **116**(5), 3099–3107.

1626 Allen, J. S., Miller, J. L., and DeSteno, D. (2003). "Individual talker differences in voice-onset-time," *J. Acoust. Soc. Am.* **113**(1), 544–552.

1627 Apfelbaum, K., and McMurray, B. (2015). "Relative cue encoding in the context of sophisticated models of categorization: Separating information from categorization," *Psychon. Bull. Rev.* **22**(4), 916–943.

1628 Assmann, P. F., and Katz, W. F. (2005). "Synthesis fidelity and time-varying spectral change in vowels," *J. Acoust. Soc. Am.* **117**(2), 886–895.

1629 Assmann, P. F., Nearey, T. M., and Bharadwaj, S. (2008). "Analysis of a vowel database," *Can. Acoust.* **36**(3), 148–149.

1630 Baese-Berk, M. M., Walker, K., and Bradlow, A. (2018). "Variability in speaking rate of native and non-native speakers," *J. Acoust. Soc. Am.* **144**(3), 1717–1717.

1631 Balzano, G. J. (1982). "The pitch set as a level of description for studying musical pitch perception," in *Music, Mind, and Brain: The Neuropsychology of Music*, edited by M. Clynes (Springer, New York), pp. 321–351.

1632 Barreda, S. (2020). "Vowel normalization as perceptual constancy," *Language* **96**(2), 224–254.

1633 Barreda, S. (2021). "Perceptual validation of vowel normalization methods for variationist research," *Lang. Var. Change* **33**(1), 27–53.

1634 Barreda, S., and Jaeger, T. F. (submitted). *Re-Introducing the Probabilistic Sliding Template Model of Vowel Perception* (Linguistic Vanguard, Berlin).

1635 Barreda, S., and Nearey, T. M. (2012). "The direct and indirect roles of fundamental frequency in vowel perception," *J. Acoust. Soc. Am.* **131**(1), 466–477.

1636 Barreda, S., and Nearey, T. M. (2018). "A regression approach to vowel normalization for missing and unbalanced data," *J. Acoust. Soc. Am.* **144**(1), 500–520.

1637 Bladon, A., Henton, C., and Pickering, J. (1984). "Towards an auditory theory of speaker normalization," *Lang. Commun.* **4**, 59–69.

1638 Boersma, P., and Weenink, D. (2022). "Praat: Doing phonetics by computer [computer program]."

1639 Buz, E., and Jaeger, T. F. (2016). "The (in) dependence of articulation and lexical planning during isolated word production," *Lang. Cogn. Neurosci.* **31**(3), 404–424.

1640 Byrd, R. H., Lu, P., Nocedal, J., and Zhu, C. (1995). "A limited memory algorithm for bound constrained optimization," *SIAM J. Sci. Comput.* **16**(5), 1190–1208.

1641 Carpenter, G. A., and Govindarajan, K. K. (1993). "Neural network and nearest neighbor comparison of speaker normalization methods for vowel recognition," in *Proceedings of the International Conference on Artificial Neural Networks*, edited by S. Gielen and B. Kappen (Springer London, London), pp. 412–415.

1642 Chládková, K., Podlipský, V. J., and Chionidou, A. (2017). "Perceptual adaptation of vowels generalizes across the phonology and does not require local context," *J. Exp. Psychol. Hum. Percept. Perform.* **43**(2), 414.

1643 Clayards, M., Tanenhaus, M. K., Aslin, R. N., and Jacobs, R. A. (2008). "Perception of speech reflects optimal use of probabilistic speech cues," *Cognition* **108**(3), 804–809.

1644 Colby, S., Clayards, M., and Baum, S. (2018). "The role of lexical status and individual differences for perceptual learning in younger and older adults," *J. Speech. Lang. Hear. Res.* **61**(8), 1855–1874.

- 1755 Cole, J., Linebaugh, G., Munson, C., and McMurray, B. (2010).
 1756 "Unmasking the acoustic effects of vowel-to-vowel coarticulation: A statistical modeling approach," *J. Phon.* **38**(2), 167–184.
 1757
 1758 Crinnion, A. M., Malmkog, B., and Toscano, J. C. (2020). "A graph-theoretic approach to identifying acoustic cues for speech sound categorization," *Psychon. Bull. Rev.* **27**(6), 1104–1125.
 1759
 1760 Disner, S. F. (1980). "Evaluation of vowel normalization procedures," *J. Acoust. Soc. Am.* **67**(1), 253–261.
 1761
 1762 Eaves, B. S., Jr., Feldman, N. H., Griffiths, T. L., and Shafto, P. (2016). "Infant-directed speech is consistent with teaching," *Psychol. Rev.* **123**(6), 758.
 1763
 1764 Escudero, P., and Bion, R. A. H. (2007). "Modeling vowel normalization and sound perception as sequential processes," in *Proceedings of the 16th International Congress of Phonetic Sciences*, August 6–10, Saarbrücken, Germany, 1413–1416.
 1765
 1766 Fant, G. (1975). "Non-uniform vowel normalization," *STL-QPSR* **16**(2–3), 001–019.
 1767
 1768 Fant, G., Kruckenberg, A., Gustafson, K., and Liljencrants, J. (2002). "A new approach to intonation analysis and synthesis of Swedish," *Proc. Fonetik, TMH-QPSR* **44**(1), 161–164.
 1769
 1770 Feldman, N. H., Griffiths, T. L., and Morgan, J. L. (2009). "The influence of categories on perception: Explaining the perceptual magnet effect as optimal statistical inference," *Psychol. Rev.* **116**(4), 752–782.
 1771
 1772 Flemming, E. (2010). "Modeling listeners: Comments on Pluymaekers *et al.* and Scarborough," in *Laboratory Phonology*, edited by C. Fougeron, B. Kühnert, M. D'Imperio, and N. Vallée (De Gruyter Mouton, Berlin), pp. 587–606.
 1773
 1774 Gahl, S., Yao, Y., and Johnson, K. (2012). "Why reduce? Phonological neighborhood density and phonetic reduction in spontaneous speech," *J. Memory Lang.* **66**(4), 789–806.
 1775
 1776 Gerstman, L. (1968). "Classification of self-normalized vowels," *IEEE Trans. Audio Electroacoust.* **16**(1), 78–80.
 1777
 1778 Glasberg, B. R., and Moore, B. C. J. (1990). "Derivation of auditory filter shapes from notched-noise data," *Hear. Res.* **47**(1), 103–138.
 1779
 1780 Goldinger, S. D. (1996). "Words and voices: Episodic traces in spoken word identification and recognition memory," *J. Exp. Psychology: Learn. Memory Cogn.* **22**(5), 1166–1183.
 1781
 1782 Hall, K. C., Hume, E., Jaeger, T. F., and Wedel, A. (2018). "The role of predictability in shaping phonological patterns," *Linguistics Vanguard* **4**(s2), 20170027.
 1783
 1784 Hay, J., Podlubny, R., Drager, K., and McAuliffe, M. (2017). "Car-talk: Location-specific speech production and perception," *J. Phon.* **65**, 94–109.
 1785
 1786 Hay, J., Walker, A., Sanchez, K., and Thompson, K. (2019). "Abstract social categories facilitate access to socially skewed words," *PLoS One* **14**(2), e0210793–29.
 1787
 1788 Hillenbrand, J. M., Getty, L. A., Clark, M. J., and Wheeler, K. (1995). "Acoustic characteristics of American English vowels," *J. Acoust. Soc. Am.* **97**(5), 3099–3111.
 1789
 1790 Hillenbrand, J. M., Houde, R. A., and Gayvert, R. T. (2006). "Speech perception based on spectral peaks versus spectral shape," *J. Acoust. Soc. Am.* **119**(6), 4041–4054.
 1791
 1792 Hillenbrand, J. M., and Nearey, T. M. (1999). "Identification of resynthesized /h/d/ utterances: Effects of formant contour," *J. Acoust. Soc. Am.* **105**(6), 3509–3523.
 1793
 1794 Hindle, D. (1978). "Approaches to vowel normalization in the study of natural speech," in *Linguistic Variation: Models and Methods*, edited by D. Sankoff (Academic Press, New York), pp. 161–171.
 1795
 1796 Jaeger, T. F. (2024). "MVBeliefUpdatr: Fitting, Summarizing, and visualizing of multivariate gaussian ideal observers and adaptors, R package version 0.0.1.0010," <https://github.com/hlplab/MVBeliefUpdatr> (Last viewed ■■■).
 1797
 1798 Johnson, K. (1997). "Speech perception without speaker normalization," in *Talker Variability in Speech Processing*, edited by K. Johnson and W. Mullenix (CA: Academic Press, San Diego), pp. 146–165.
 1799
 1800 Johnson, K. (2020). "The ΔF method of vocal tract length normalization for vowels," *Lab. Phonology* **11**(1), 10.
 1801
 1802 Johnson, K., and Sjerp, 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, New York), pp. 145–176..
 1803
 1804 Johnson, K., Strand, E. A., and D'Imperio, M. (1999). "Auditory–visual integration of talker gender in vowel perception," *J. Phon.* **27**(4), 359–384.
 1805
 1806 Joos, M. (1948). "Acoustic phonetics," *Language* **24**(2), 5–136.
 1807 Kleinschmidt, D. (2020). "What constrains distributional learning in adults?," *Psychol. Rev.* **122**(2), 148–203.
 1808
 1809 Kleinschmidt, D., and Jaeger, T. F. (2015). "Robust speech perception: Recognize the familiar, generalize to the similar, and adapt to the novel," *Psychol. Rev.* **122**(2), 148–203.
 1810 Kleinschmidt, D., and Jaeger, T. F. (2016). "What do you expect from an unfamiliar talker?," in *Proceedings of the 38th Annual Meeting of the Cognitive Science Society*, August 10–13, Philadelphia, PA, pp. 2351–2356.
 1811 Kleinschmidt, D., Liu, L., Bushong, W., Burchill, Z., Xie, X., Tan, M., Karboga, G., and Jaeger, F. (2021). "JSEXP," <https://github.com/hlplab/JSEXP> (Last viewed ■■■).
 1812
 1813 Kronrod, Y., Coppess, E., and Feldman, N. H. (2016). "A unified model of categorical effects in consonant and vowel perception," *Psychol. Bull.* **143**(2), 1681–1712.
 1814
 1815 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.
 1816
 1817 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).
 1818
 1819 Ladefoged, P., and Broadbent, D. E. (1957). "Information conveyed by vowels," *J. Acoust. Soc. Am.* **29**, 98–104.
 1820
 1821 Lee, C.-Y. (2009). "Identifying isolated, multispeaker mandarin tones from brief acoustic input: A perceptual and acoustic study," *J. Acoust. Soc. Am.* **125**(2), 1125–4966.
 1822
 1823 Liberman, A. M., Cooper, F. S., Shankweiler, D. P., and Studdert-Kennedy, M. (1967). "Perception of the speech code," *Psychol. Rev.* **74**(6), 431–461.
 1824
 1825 Lindblom, B. (1986). "Phonetic universals in vowel systems," in *Experimental Phonology*, edited by J. J. Ohala and J. J. Jaeger (Academic Press, Orlando, FL), pp. 13–44.
 1826
 1827 Lindblom, B. (1990). "Explaining phonetic variation: A sketch of the H&H theory," in *Speech Production and Speech Modeling*, edited by W. J. Hardcastle and A. Marchal (Kluwer, Dordrecht, the Netherlands), pp. 403–439.
 1828
 1829 Lobanov, B. M. (1971). "Classification of Russian vowels spoken by different speakers," *J. Acoust. Soc. Am.* **49**(2B), 606–608.
 1830
 1831 Luce, P. A., and Pisoni, D. B. (1998). "Recognizing spoken words: The neighborhood activation model," *Ear Hear.* **19**(1), 1–36.
 1832
 1833 Luce, R. D. (1959). *Individual Choice Behavior* (John Wiley, Oxford, UK).
 1834
 1835 Magnuson, J. S., and Nusbaum, H. C. (2007). "Acoustic differences, listener expectations, and the perceptual accommodation of talker variability," *J. Exp. Psychol. Human Percept. Perform.* **33**(2), 391–409.
 1836
 1837 Magnuson, J. S., You, H., Luthra, S., Li, M., Nam, H., Escabi, M., Brown, K., Allopenna, P. D., Theodore, R., Monto, N., and Rueckl, J. G. (2020). "EARSHOT: A minimal neural network model of incremental human speech recognition," *Cogn. Sci.* **44**(4), 1–17.
 1838
 1839 Massaro, D. W., and Friedman, D. (1990). "Models of integration given multiple sources of information," *Psychol. Rev.* **97**(2), 225–252.
 1840
 1841 McClelland, J. L., and Elman, J. L. (1986). "The TRACE model of speech perception," *Cogn. Psychol.* **18**(1), 1–86.
 1842
 1843 McGowan, K. B. (2015). "Social expectation improves speech perception in noise," *Lang. Speech* **58**(4), 502–521.
 1844
 1845 McMurray, B., and Jongman, A. (2011). "What information is necessary for speech categorization?: Harnessing variability in the speech signal by integrating cues computed relative to expectations," *Psychol. Rev.* **118**(2), 219–246.
 1846
 1847 Merzenich, M. M., Knight, P. L., and Roth, G. L. (1975). "Representation of cochlea within primary auditory cortex in the cat," *J. Neurophysiol.* **38**(2), 231–249.
 1848
 1849 Miller, J. D. (1989). "Auditory-perceptual interpretation of the vowel," *J. Acoust. Soc. Am.* **85**(5), 2114–2134.
 1850
 1851 Moore, B. C. (2012). *An Introduction to the Psychology of Hearing* (Brill, Berlin).
 1852
 1853 Moulin-Frier, C., Diard, J., Schwartz, J.-L., and Bessière, P. (2015). "Cosmo ('Communicating about objects using sensory-motor

- 1898 operations'): A Bayesian modeling framework for studying speech communication and the emergence of phonological systems," *J. Phon.* **53**, 1899 5–41.
- 1900 Nearey, T. M. (1978). *Phonetic Feature Systems for Vowels* (Indiana University Linguistics Club, Bloomington, IN).
- 1901 Nearey, T. M. (1989). "Static, dynamic, and relational properties in vowel 1902 perception," *J. Acoust. Soc. Am.* **85**(5), 2088–2113.
- 1903 Nearey, T. M. (1990). "The segment as a unit of speech perception," 1904 *J. Phon.* **18**(3), 347–373.
- 1905 Nearey, T. M., and Assmann, P. F. (1986). "Modeling the role of inherent 1906 spectral change in vowel identification," *J. Acoust. Soc. Am.* **80**(5), 1907 1297–1308.
- 1908 Nearey, T. M., and Assmann, P. F. (2007). "Probabilistic 'sliding template' 1909 models for indirect vowel normalization," in *Experimental Approaches to 1910 Phonology*, edited by J.-J. Solé, P. S. Beddor, and M. Ohala (Oxford University Press, Oxford, UK), pp. 246–270.
- 1911 Nearey, T. M., and Hogan, J. (1986). "Phonological contrast in 1912 experimental phonetics: Relating distributions of measurements production 1913 data to perceptual categorization curves," in *Experimental Phonology*, edited by J. J. Ohala and J. Jaeger (Academic Press, New 1914 York), pp. 141–161.
- 1915 Newman, R. S., Clouse, S. A., and Burnham, J. L. (2001). "The perceptual 1916 consequences of within-talker variability in fricative production," *J. Acoust. 1917 Soc. Am.* **109**(3), 1181–1196.
- 1918 Nordström, P., and Lindblom, B. (1975). "A normalization procedure for 1919 vowel formant data," in *Proceedings of the 8th International Congress of 1920 Phonetic Sciences*, August 17–23, Leeds, UK.
- 1921 Norris, D., and McQueen, J. M. (2008). "Shortlist B: A Bayesian model of 1922 continuous speech recognition," *Psychol. Rev.* **115**(2), 357–395.
- 1923 Oganian, Y., Bhaya-Grossman, I., Johnson, K., and Chang, E. F. (2023). 1924 "Vowel and formant representation in the human auditory speech cortex," *Neuron* **111**(13), 2105–2118.
- 1925 Patterson, R. D., and Irino, T. (2014). "Size matters in hearing: How the 1926 auditory system normalizes the sounds of speech and music for source 1927 size," in *Perspectives on Auditory Research* (Springer, New York), pp. 417–440.
- 1928 Persson, A., and Jaeger, T. F. (2023). "Evaluating normalization accounts 1929 against the dense vowel space of Central Swedish," *Front. Psychol.* **14**, 1165742.
- 1930 Peterson, G. E. (1961). "Parameters of vowel quality," *J. Speech Hear. Res.* 1931 **4**(1), 10–29.
- 1932 R Core Team (2024). "R: A language and environment for statistical 1933 computing," <https://www.R-project.org/> (Last viewed ■■■).
- 1934 Repp, B. H., and Crowder, R. G. (1990). "Stimulus order effects in vowel 1935 discrimination," *J. Acoust. Soc. Am.* **88**(5), 2080–2090.
- 1936 Richter, C., Feldman, N. H., Salgado, H., and Jansen, A. (2017). 1937 "Evaluating low-level speech features against human perceptual data," *Trans. Assoc. Comput. Ling.* **5**, 425–440.
- 1938 RStudio Team (2020). *RStudio: Integrated Development Environment for R*, RStudio (PBC, Boston, MA).
- 1939 Saenz, M., and Langers, D. R. (2014). "Tonotopic mapping of human 1940 auditory cortex," *Hear. Res.* **307**, 42–52.
- 1941 Scarborough, R. (2010). "Lexical and contextual predictability: Confluent 1942 effects on the production of vowels," in *Laboratory Phonology*, edited by C. Fougeron, B. Kühnert, M. D'Imperio, and N. Vallée (De Gruyter Mouton, Berlin), pp. 557–586.
- 1943 Schertz, J., and Clare, E. J. (2020). "Phonetic cue weighting in perception 1944 and production," *Wiley Interdiscip. Rev. Cognit. Sci.* **11**(2), e1521.
- 1945 Shankweiler, D., Verbrugge, R. R., and Studdert-Kennedy, M. (1978). 1946 "Insufficiency of the target for vowel perception," *J. Acoust. Soc. Am.* **63**(S1), S4–S4.
- 1947 Shannon, C. E. (1948). "A mathematical theory of communication," *Bell 1948 Syst. Tech. J.* **27**(3), 379–423.
- 1949 Siegel, R. J. (1965). "A replication of the mel scale of pitch," *Am. J. 1950 Psychol.* **78**(4), 615–620.
- 1951 Sjerps, M. J., Fox, N. P., Johnson, K., and Chang, E. F. (2019). "Speaker- 1952 normalized sound representations in the human auditory cortex," *Nat. 1953 Commun.* **10**(1), 01–09.
- 1954 Skoe, E., Krizman, J., Spitzer, E. R., and Kraus, N. (2021). "Auditory 1955 cortical changes precede brainstem changes during rapid implicit learning: 1956 Evidence from human EEG," *Front. Neurosci.* **15**, 01–09.
- 1957 Smith, B. L., Johnson, E., and Hayes-Harb, R. (2019). "ESL learners' intra- 1958 speaker variability in producing American English tense and lax vowels," *JSLP* **5**(1), 139–164.
- 1959 Smith, D. R., Patterson, R. D., Turner, R., Kawahara, H., and Irino, T. (1970). 1960 "The processing and perception of size information in speech sounds," *J. Acoust. Soc. Am.* **117**(1), 305–318.
- 1961 Steriade, D. (2008). "The phonology of perceptibility effects: The P-map 1962 and its consequences for constraint organization," in *The Nature of the 1963 Word: Studies in Honor of Paul Kiparsky*, edited by K. Hanson and S. 1964 Inkelas (MIT Press, Cambridge, MA).
- 1965 Stevens, K. N. (1972). "The quantal nature of speech: Evidence from 1966 articulatory-acoustic data," in *Human Communication: A Unified View* 1967 (McGraw Hill, New York), pp. 51–66.
- 1968 Stevens, K. N. (1989). "On the quantal nature of speech," *J. Phon.* **17**(1–2), 1969 3–45.
- 1970 Stevens, S. S., and Volkmann, J. (1940). "The relation of pitch to frequency: 1971 A revised scale," *Am. J. Psychol.* **53**(3), 329–353.
- 1972 Stilp, C. (2020). "Acoustic context effects in speech perception," *WIREs 1973 Cognit. Sci.* **11**(1), 1–18.
- 1974 Strange, W., and Jenkins, J. J. (2012). "Dynamic specification of coarticulated 1975 vowels: Research chronology, theory, and hypotheses," in *Vowel Inherent 1976 Spectral Change* (Springer, New York), pp. 87–115.
- 1977 Sumner, M. (2011). "The role of variation in the perception of accented 1978 speech," *Cognition* **119**(1), 131–136.
- 1979 Syrdal, A. K. (1985). "Aspects of a model of the auditory representation of 1980 American English vowels," *Speech Commun.* **4**(1–3), 121–135.
- 1981 Syrdal, A. K., and Gopal, H. S. (1986). "A perceptual model of vowel 1982 recognition based on the auditory representation of American English 1983 vowels," *J. Acoust. Soc. Am.* **79**(4), 1086–1100.
- 1984 Tan, M., and Jaeger, T. F. (2024). *Incremental Adaptation to an Unfamiliar 1985 Talker* (Stockholm University, Stockholm, Sweden).
- 1986 Tang, C., Hamilton, L. S., and Chang, E. F. (2017). "Intonational speech 1987 prosody encoding in the human auditory cortex," *Science* **357**(6353), 1988 797–801.
- 1989 ten Bosch, L., Boves, L., Tucker, B., and Ernestus, M. (2015). "DIANA: 1990 Towards computational modeling reaction times in lexical decision in 1991 north American English," in *Proceedings of Interspeech 2015*, September 1992 6–10, Dresden, Germany, pp. 1576–1580.
- 1993 Traunmüller, H. (1981). "Perceptual dimension of openness in vowels," *J. 1994 Acoust. Soc. Am.* **69**(5), 1465–1475.
- 1995 Traunmüller, H. (1990). "Analytical expressions for the tonotopic sensory 1996 scale," *J. Acoust. Soc. Am.* **88**(1), 97–100.
- 1997 Vaughn, C., Baese-Berk, M., and Idemaru, K. (2019). "Re-examining 1998 phonetic variability in native and non-native speech," *Phonetica* **76**(5), 1999 327–358.
- 2000 Vorperian, H. K., and Kent, R. D. (2007). "Vowel acoustic space development 2001 in children: A synthesis of acoustic and anatomic data," *J. Speech. 2002 Lang. Hear. Res.* **50**(6), 1510–1545.
- 2003 Wade, T., Jongman, A., and Sereno, J. (2007). "Effects of acoustic variability 2004 in the perceptual learning of non-native-accented speech sounds," *2005 Phonetica* **64**(2–3), 122–144.
- 2006 Walker, A., and Hay, J. (2011). "Congruence between 'word age' and 2007 'voice age' facilitates lexical access," *Lab. Phonol.* **2**(1), 219–237.
- 2008 Watt, D., and Fabricius, A. (2002). "Evaluation of a technique for improving 2009 the mapping of multiple speakers' vowel spaces in the F1 ~ F2 2010 plane," in *Leeds Working Papers in Linguistics and Phonetics*, edited by 2011 D. Nelson (University of Leeds, Leeds, UK), pp. 159–173.
- 2012 Weatherholtz, K., and Jaeger, T. F. (2016). *Speech Perception and 2013 Generalization across Talkers and Accents* (Oxford University Press, 2014 Oxford, UK).
- 2015 Wedel, A., Nelson, N., and Sharp, R. (2018). "The phonetic specificity of 2016 contrastive hyperarticulation in natural speech," *J. Memory Lang.* **100**, 2017 61–88.
- 2018 Whalen, D. H. (2016). "A double-Nearey theory of vowel normalization: 2019 Approaching consensus," *J. Acoust. Soc. Am.* **140**(4_Supplement), 2020 3163–3164.
- 2021 Wichmann, F. A., and Hill, N. J. (2001). "The psychometric function: I. 2022 Fitting, sampling, and goodness of fit," *Percept. Psychophys.* **63**(8), 2023 1293–1313.
- 2024 Winn, M. (2018). "Speech: It's not as acoustic as you think," *Acoust. 2025 Today* **12**(2), 43–49.

- 2040 Xie, X., Buxó-Lugo, A., and Kurumada, C. (2021). “Encoding and decoding of meaning through structured variability in speech prosody,”
2041 *Cognition* **211**, 104619–104627.
- 2042 Xie, X., and Jaeger, T. F. (2020). “Comparing non-native and native
2043 speech: Are L2 productions more variable?,” *J. Acoust. Soc. Am.* **147**(5),
2044 3322–3347.
- 2045 Xie, X., Jaeger, T. F., and Kurumada, C. (2023). “What we do (not) know
2046 about the mechanisms underlying adaptive speech perception: A compu-
2047 tational review,” *Cortex* **166**, 377–424.
- 2048 Zahorian, S. A., and Jagharghi, A. J. (1991). “Speaker normalization of static
2049 and dynamic vowel spectral features,” *J. Acoust. Soc. Am.* **90**(1), 67–75.
- 2050 Zwicker, E. (1961). “Subdivision of the audible frequency range into criti-
2051 cal bands (frequenzgruppen),” *J. Acoust. Soc. Am.* **33**(2), 248–248.
- 2052 Zwicker, E., Flottorp, G., and Stevens, S. S. (1957). “Critical band width in
2053 loudness summation,” *J. Acoust. Soc. Am.* **29**(5), 548–557.
- 2054 Zwicker, E., and Terhardt, E. (1980). “Analytical expressions for critical-
2055 band rate and critical bandwidth as a function of frequency,” *J. Acoust. Soc. Am.* **68**(5), 1523–1525.
- 2056
- 2057