

Evaluating normalization accounts against the dense vowel space of Central Swedish

Anna S. Persson^{1*}, T. Florian Jaeger^{2,3}

¹ Department of Swedish Language and Multilingualism, Stockholm University, Stockholm, Sweden

² Brain and Cognitive Sciences, University of Rochester, Rochester, NY, USA

³ Computer Science, University of Rochester, Rochester, NY, USA

Correspondence*:

Anna S. Persson

anna.persson@su.se

2 ABSTRACT

3 Talkers vary in the phonetic realization of their vowels. One influential hypothesis holds that listeners
4 overcome this inter-talker variability through pre-linguistic auditory mechanisms that normalize the
5 acoustic or phonetic cues that form the input to speech recognition. Dozens of competing normalization
6 accounts exist—including both accounts specific to vowel perception and general purpose accounts that
7 can be applied to any type of cue. We add to the cross-linguistic literature on this matter by comparing
8 normalization accounts against a new phonetically annotated vowel database of Swedish, a language
9 with a particularly dense vowel inventory of 21 vowels differing in quality and quantity. All data and code
10 for this study are shared via OSF, including the R markdown document that this article is generated
11 from, and an R library that implements the models we present.

12 Keywords: vowel normalization; category separability; ideal observers; speech production; speech perception

INTRODUCTION

13 Talkers differ in their pronunciation of individual speech sounds due to both physiological differences
 14 and socio-cultural factors, including style, regional dialect, and second language accents. For
 15 listeners, this means that the mapping from acoustic cues to linguistic categories—phonemes,
 16 syllables, words, and ultimately word meanings—varies depending on the talker. How listeners
 17 manage to typically understand talkers despite this “lack of invariance” (Liberman et al., 1967)
 18 has remained one of the central questions for research on speech perception. Hypotheses about
 19 the mechanisms underlying this ability can be grouped into three, mutually compatible and
 20 complementary, accounts: (1) low-level, pre-linguistic auditory transformation of the acoustic signal,
 21 (2) learning of changes in the linguistic representations, and (3) post-linguistic changes in decision-
 22 making biases (see e.g., Johnson, 2006; Pardo and Remez, 2006; Xie et al., 2022). The present
 23 study focuses on the first type of account, that the acoustic signal is transformed and normalized
 24 early on during auditory processing (for recent reviews, Johnson and Sjerps, 2021; Stilp, 2020).

25 The existence of some form of pre-linguistic normalization is motivated by *a priori* considerations
 26 about the physics of sounds (cf. the discussion of uniform scaling in Barreda, 2020), evolutionary
 27 arguments (e.g., even non-human animals exhibit similar abilities, as reviewed in Barreda, 2020),
 28 as well as brain imaging evidence: talker-normalized information about the speech signal can be
 29 decoded from areas as early as the brain stem (e.g., Sjerps et al., 2019; Sussman, 1986), and thus
 30 prior to even the earliest cortical areas typically associated with linguistic category representations
 31 or decision-making. While it is rather uncontroversial that normalization is part of adaptive speech
 32 perception, questions remain about the specific nature of the operations involved in normalization.
 33 We contribute to this line of research by comparing different types of normalization accounts against
 34 data from the production of short and long vowels of Central Swedish.

35 Normalization accounts were originally proposed as a theory of how the brain removes
 36 *physiologically*-conditioned variation from the speech signal, reducing variability in, for example,
 37 category means between talkers, and thus reducing the overlap of phonological categories in
 38 the acoustic-phonetic space (e.g., Bladon et al., 1984; Gerstman, 1968; Lobanov, 1971; Miller,
 39 1989; Nearey, 1978; Nordström and Lindblom, 1975; Peterson, 1961; Syrdal and Gopal, 1986;
 40 Sussman, 1986). Most of this early work focused specifically on differences in formants, which
 41 cross-linguistically are the primary cues to distinctions in vowel quality, and known to be affected
 42 by the vocal tract size of the talker (e.g., Fox et al., 1995; Peterson and Barney, 1952; Verbrugge
 43 and Shankweiler, 1977; Yang and Fox, 2014). Figure 1 illustrates the effect of one of the most
 44 commonly applied normalization accounts (Lobanov, 1971) for the vowel spaces of three L1
 45 speakers of US English (from the database reported in Xie and Jaeger, 2020). By normalizing
 46 cues prior to categorization, physiological differences between talkers can be reduced, resulting in
 47 reduced between-talker variability (compare Figure 1B to 1A). If listeners’ category representations
 48 pool experiences across talkers into a single talker-independent model, this reduced inter-talker
 49 variability results in reduced category overlap (compare Figure 1D to 1C).¹

¹ Talker-independent category representations are assumed in many influential models of spoken word recognition (e.g., Luce and Pisoni, 1998; McClelland and Elman, 1986; Norris and McQueen, 2008). While talker-independent representations might be a simplifying assumption for some of these theories, this assumption has persisted for decades (e.g. Magnuson et al., 2020; ten Bosch et al., 2022). Exceptions include exemplar accounts (e.g., Johnson, 1997; Pierrehumbert, 2001) and the Bayesian ideal adaptor account (Kleinschmidt and Jaeger, 2015). Importantly, it is an unresolved question whether—or for which cues and phonetic contrasts—listeners maintain talker-specific category representations (for findings and discussion, see Kraljic and Samuel, 2007; Kleinschmidt and Jaeger, 2015; Kleinschmidt, 2019; Xie et al., 2021). In the present paper, we follow previous work and compare the effectiveness of normalization under the assumption of talker-independent category representations.

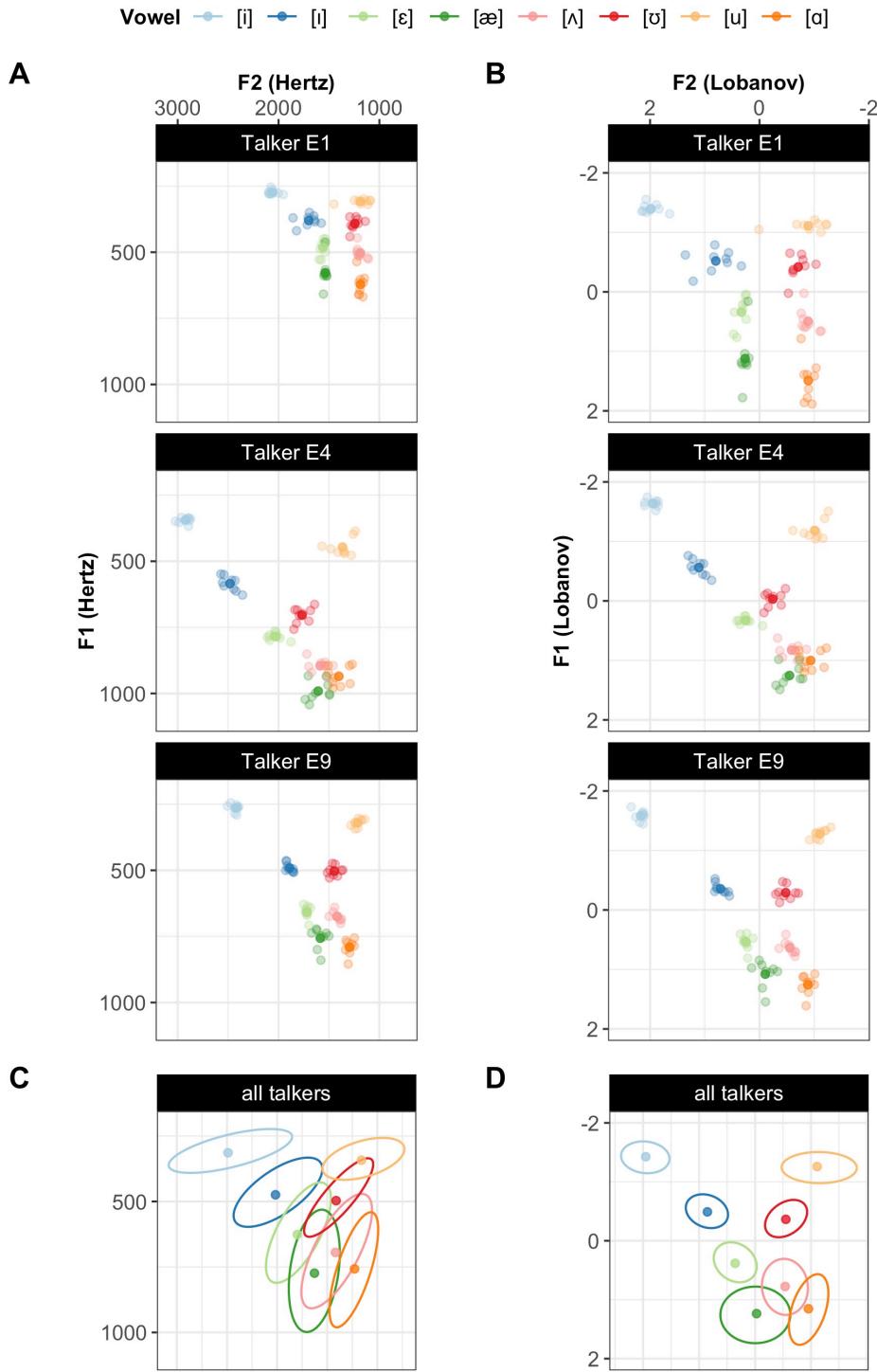


Figure 1. Illustrating how normalization reduces category overlap for the 8 monophthongs of L1 US English. Three talkers from the Xie and Jaeger (2020) database are shown before (**Panel A**) and after Lobanov normalization (**Panel B**). Lobanov normalization reduces inter-talker variability in the category means and, to some extent, in the category variances. The bottom two panels aggregate the data from all 17 talkers in the database (5 female, 12 male), showing the means and 95% probability mass bivariate Gaussian densities for each vowel before (**Panel C**) and after Lobanov normalization (**Panel D**).

50 Dozens of different accounts of vowel normalization have been proposed over the years (e.g.,
 51 Bladon et al., 1984; Fant, 1975; Gerstman, 1968; Joos, 1948; Lobanov, 1971; Miller, 1989; Nearey,
 52 1978; Nordström and Lindblom, 1975; Syrdal and Gopal, 1986; Traunmüller, 1981; Watt and
 53 Fabricius, 2002; Zahorian and Jagharghi, 1991; for reviews, see Barreda, 2020; Weatherholtz
 54 and Jaeger, 2016). Carpenter and Govindarajan (1993) summarize over 100 different vowel-
 55 specific normalization accounts, many of them closely related to each other. More recently,
 56 additional *general* normalization accounts have emerged that can be applied to *any* type of cue
 57 and phonological contrast (e.g., Cole et al., 2010; McMurray and Jongman, 2011). The most widely
 58 used of these proposals, C-CuRE, has since been successfully applied to the categorization of US
 59 English fricatives (Apfelbaum et al., 2014; Crinnion et al., 2020; McMurray and Jongman, 2011),
 60 stop voicing (Kulikov, 2022; Toscano and McMurray, 2015; Xie et al., 2022), vowels (Kleinschmidt,
 61 2019; McMurray and Jongman, 2016), and sentence-final rising question vs. statement intonation
 62 (Xie et al., 2021). In each of these studies, C-CuRE reduced inter-talker variability and improved
 63 categorization. C-CuRE stands for “computing **cues** relative to **expectations**”, capturing the
 64 motivation behind many of the earlier normalization accounts that cue values should be interpreted
 65 relative to the distribution they are expected to have in the present context. In addition to being
 66 more general than earlier accounts, C-CuRE also addresses one of the main arguments against
 67 normalization as an account of human speech perception (e.g., Johnson, 2006, which continues to
 68 be frequently cited in reviews of exemplar theory): by focusing on listeners’ expectations rather
 69 than talkers’ physiology, accounts like C-CuRE capture that inter-talker variability is not limited
 70 to physiology.

71 Table 1 lists the normalization accounts investigated in the present study. This includes both
 72 the most influential vowel-specific normalization accounts that have been found to perform well
 73 in previous works (e.g., Lobanov and Nearey² normalization) and several variants of the general
 74 purpose normalization C-CuRE. As indicated through shading in the table, the accounts can be
 75 grouped into four types based on the computational assumptions they make. *Transformations* are
 76 meant to transform the formant data from acoustic (Hz) into a perceptual space that approximates
 77 the perceptual organization of auditory information in the human brain. All other accounts instead
 78 or additionally adjust each formant value based on either the values of other formants on the same
 79 segment (*vowel-intrinsic* approaches) or summary statistics of the formant across segments (*vowel-*
 80 *extrinsic* approaches).² We further distinguish two types of vowel-extrinsic approaches that differ
 81 in their computational complexity and tractability: approaches that *center* each cue relative to its
 82 mean across all vowel segments, and approaches that instead/additionally *standardize* cues relative
 83 to the overall variability or range of the cue across all vowel segments (for reviews, see also e.g.,
 84 Johnson, 2005; Kohn and Farrington, 2012; Weatherholtz and Jaeger, 2016).³ The former type
 85 includes C-CuRE, and we consider different variants of this approach, one for each transformation
 86 approach in Table 1.

87 Our selection of accounts to consider in the present study is primarily based on their influence
 88 and performance in previous evaluations against other data sets. Additionally, we only consider

² Miller’s formant-ratio account (Miller, 1989) is technically a hybrid approach: the first formant (F1) is normalized with regard to an extrinsic sensory reference (based on the average F0 across segments); subsequent formants are (intrinsically) normalized using the normalized lower formants on the same vowel segment.

³ Here we group accounts based on their computational complexity (the number of parameters listeners are assumed to estimate). For example, we group Nearey1 and Nearey2 with the centering accounts because they require estimation of only cue means. However, since these accounts perform centering over log-transformed Hz, they can also be considered as a form of functionally constrained standardization in non-log space (Barreda and Nearey, 2018).

89 accounts that are sufficiently general in nature to be applied across languages. This decision stems
90 from our goal to understand the mechanisms underlying *human* speech perception. This means that
91 we for instance do not include Watt & Fabricius (Watt and Fabricius, 2002; Fabricius et al., 2009),
92 another frequently used normalization account, as it requires making specific assumptions of vowel
93 inventories of the language. Finally, we do not consider *combinations* of accounts. This follows the
94 majority of previous work but is an important limitation that we return to in the general discussion.

Table 1. Normalization accounts considered in the present study. Unless otherwise marked, formant variables (F 's) in the right-hand side of normalization formulas are in Hz.

	Normalization procedure	Perceptual scale	Source	Formula	
transformation	n/a	Hz	n/a	n/a	
	n/a	Bark	Traunmüller (1990)	$F_n^{Bark} = \frac{26.81 \times F_n}{1960 + F_n} - 0.53$	
	—	ERB	Glasberg & Moore (1990)	$F_n^{ERB} = 21.4 \times \log_{10}(1 + F_n \times 0.00437)$	
	—	Mel	Stevens & Volkmann (1940)	$F_n^{Mel} = 2595 \times \log_{10}(1 + \frac{F_n}{700})$	
	—	Semitones conversion	Fant et al. (2002)	$F_n^{ST} = 12 \times \frac{\ln(\frac{F_n}{100})}{\ln}$	
intrinsic	Syrdal & Gopal's Bark-distance model ⁴	Bark	Syrdal & Gopal (1986)	$F1^{SyrdalGopal} = F1^{Bark} - F0^{Bark}$	
	Miller (formant-ratio)	log	Miller (1989)	$F2^{SyrdalGopal} = F2^{Bark} - F1^{Bark}$ $SR = k(\frac{GMf_0}{k})^{1/3}$ $F1^{Miller} = \log(\frac{F1}{SR})$ $F2^{Miller} = \log(\frac{F2}{F1})$ $F3^{Miller} = \log(\frac{F3}{F2})$	
extrinsic	centering	C-CuRE	Hz	McMurray & Jongman (2011)	$F_n^{C-CuRE} = F_n - \text{mean}(F_n)$
		—	Bark		
		—	ERB		
		—	Mel		
		—	Semitones conversion		
	Nearey1 (log-mean)	log	Nearey (1978)	$F_n^{Nearey1} = \ln(F_n) - \text{mean}(\ln(F_n))$	
	Nearey2 (single parameter log-mean)	log	Nearey (1978)	$F_n^{Nearey2} = \ln(F_n) - \text{mean}(\ln(F))$	
standardizing	Gerstman (range normalization)	Hz	Gerstman (1968)	$F_n^{Gerstman} = 999 \times \frac{F_n - F_{min}}{F_{max} - F_{min}}$	
	Lobanov (z-score)	Hz	Lobanov (1971)	$F_n^{Lobanov} = \frac{F_n - \text{mean}(F_n)}{sd(F_n)}$	

⁴ Previous work has considered two different implementations of Syrdal & Gopal's Bark-distance model for F2, depending on the language (Adank, 2003; Fant, 1983; Syrdal & Gopal, 1986). In the SI (see Section Evaluation of implementations of Syrdal & Gopal), we compare these two implementations, and find that the F2-F1 implementation performs better for the present data. We thus present that version of Syrdal & Gopal's model in the main text of Studies 1 and 2.

95 Previous comparisons of normalization accounts have primarily focused on English (e.g., Adank
96 et al., 2004; Barreda and Nearey, 2018; Carpenter and Govindarajan, 1993; Clopper, 2009; Disner,
97 1980; Escudero and Bion, 2007; Fabricius et al., 2009; Flynn and Foulkes, 2011; Hindle, 1978;
98 Kohn and Farrington, 2012; Labov, 2010; Richter et al., 2017; Syrdal, 1985). Additional studies
99 have investigated, for example, Dutch (Adank et al., 2004; Disner, 1980), Russian (Lobanov, 1971),
100 and Brazilian Portuguese (Escudero and Bion, 2007). The complexity of the vowel inventories
101 (7-11 monophthongs) and the number of these vowels included in the comparison (2-11) varied
102 across these studies. We add to this literature by comparing normalization accounts against a
103 new phonetically annotated database of Central Swedish (SwehVd, introduced below). With a
104 total of 21 monophthong allophones that vary in quantity (long vs. short vowels) and quality,
105 the vowel inventory of Swedish is crowded compared to most languages previously studied in the
106 normalization literature. This allows us to test whether the same normalization accounts that work
107 well for simpler vowel inventories generalize well to more crowded vowel spaces. Additionally, the
108 presence of quantity contrasts between long and short allophones means that Swedish provides a
109 suitable case study to bridge the literature between vowel-specific normalization accounts (which
110 focus on formants, and thus only quality contrasts) and general normalization accounts that can
111 be applied to any type of cue (and thus also vowel duration, which is expected to be the primary
112 cue to vowel quantity). Relatedly, previous studies have found both F3 and vowel duration to
113 be important cues to vowel categorization in Swedish (e.g., Behne et al., 1997; Fujimura, 1967;
114 Hadding-Koch and Abramson, 1964). These two cues have never (duration) or rarely (F3, but see
115 e.g. Adank et al., 2004; Barreda and Nearey, 2018; Carpenter and Govindarajan, 1993; Nearey,
116 1989; Syrdal, 1985) been included in comparisons of normalization accounts.

117 We compare the normalization accounts in Table 1 both in terms of their effectiveness in reducing
118 within-category variability relative to between-category variability (Study 1) and in terms of the
119 predicted consequences for recognition accuracy (Study 2). We originally only conducted Study 2,
120 as it more adequately addresses our goal of assessing the predicted consequences of normalization for
121 perception. However, we later added Study 1 because—despite serious shortcomings—measures of
122 between- vs. within-category separability/variability continue to be commonly used in research on
123 normalization. The primary purpose of including Study 1 is to illustrate some of these shortcomings,
124 highlighting the benefits of using models of speech perception when evaluating normalization
125 accounts (as we do in Study 2, for related discussion, see Barreda, 2021).

126 Following the bulk of previous work (but see e.g., Barreda, 2021; McMurray and Jongman, 2016;
127 Nearey, 1989; Richter et al., 2017), the present studies do not compare the predicted consequences
128 of normalization against categorization responses from human listeners, though we plan to do so
129 in future work. We discuss this and other limitations—most shared with previous work—after
130 presenting our studies. As part of that discussion, we also raise conceptual considerations about
131 the *a priori* plausibility of different normalization accounts.

132 Both Study 1 and 2 compare normalization accounts applied to (1) only F1 and F2, as in the
133 majority of previous studies, (2) F1-F3, as in, e.g., Adank et al. (2004), and (3) F0-F3 as well as
134 vowel duration. This allows us to assess whether differences in the effectiveness of normalization
135 accounts depend on the number and types of cues that are considered. Since listeners integrate cues
136 beyond F1 and F2 (e.g., Assmann et al., 1982; Hillenbrand and Nearey, 1999; Nearey and Assmann,
137 1986), this is an important gap in evaluating the plausibility of different normalization accounts
138 as models of adaptive speech perception. All three comparisons are evaluated both separately for

139 short and long vowels, and for the entire space of the 21 vowels. This allows us to assess whether
 140 the same types of normalization perform well across the entire vowel inventory.

141 To the best of our knowledge, only one previous study has compared normalization accounts
 142 against Swedish, as part of a cross-linguistic comparison across six Germanic languages (Disner,
 143 1980). Disner (1980) compared 4 normalization accounts, using F1 and F2 means of the nine
 144 long Swedish vowels spoken by 24 male Swedish talkers (from a database presented in Fant et al.,
 145 1969). Of interest to the present study, the results for Swedish differed from the other Germanic
 146 languages in two unexpected ways. Whereas Lobanov normalization—which involves centering and
 147 standardizing (cf. Table 1)—performed best for Swedish, Nearey2 normalization—which involves
 148 only centering—performed best for the other four languages. And, while normalization effectively
 149 reduced inter-talker variability in category variances for the other four languages by 61%-71%, it
 150 was substantially less effective for Swedish (41%). As discussed by Disner (1980), this raises the
 151 question as to whether these findings reflect an inherent property of Swedish or merely differences
 152 in the phonetically annotated databases available for each language. In particular, the Swedish
 153 data consisted of *vowels* produced in isolation without any lexical or phonetic context, whereas the
 154 data for the five other languages consisted of isolated *word* productions (paralleling the majority
 155 of research on normalization). The present study addresses this difference: the new database we
 156 introduce consists of *h-VOWEL-d* word recordings, which makes our stimuli directly comparable
 157 to those used in previous work on normalization. Additionally, we complement Disner’s study by
 158 focusing on female, rather than male talkers, by considering both long and short vowels (separately
 159 and together), and by including the general normalization account C-CuRE. This lets us revisit
 160 whether simple *centering* accounts perform best for Swedish—like for the other languages in Disner
 161 (1980).

162 All data and code for this article can be downloaded from OSF at <https://osf.io/zb8gx/>. This
 163 article is written in R markdown, allowing readers to replicate our analyses using freely available
 164 software (R Core Team, 2021; RStudio Team, 2020), while changing any of the parameters of our
 165 models. Readers can revisit the assumptions we make—for example, by substituting alternative
 166 normalization models. The supplementary information (SI) lists the software/libraries required
 167 to compile this document. Next we introduce the new phonetically annotated corpus of Central
 168 Swedish vowel productions employed in the present studies.

THE SWEHVD DATABASE

169 The SwehVd database is a new phonetically annotated corpus of Swedish *h-VOWEL-d* (short: hVd)
 170 word recordings. All recordings, annotations, and acoustic measurements are available on OSF at
 171 <https://osf.io/ruxnb/>. SwehVd was collected with the goal to characterize the Central Swedish
 172 vowel space—specifically, the Stockholm variety—within and across talkers. It covers the entire
 173 monophthong inventory of Central Swedish, including all nine long vowels (*hid, hyd, hud, hed, häd,*
 174 *höd, had, håd, hod*), eight short vowels (*hidd, hydd, hudd, hedd, hädd, hödd, hadd, hådd, hodd*),
 175 and four allophones (*härd, härr, hörd, hörr*). To our knowledge, there are few publicly available
 176 databases of Swedish vowel productions that are phonetically annotated (e.g., Bruce et al., 1999;
 177 Eklund and Traunmüller, 1997; Fant et al., 1969; Kuronen, 2000). The largest and perhaps best-
 178 known is SweDia 2000 (Bruce et al., 1999). SweDia 2000 was developed to characterize differences
 179 in vowel pronunciations *across* regional varieties of Swedish. It consists of recordings of spontaneous
 180 speech, isolated words in varying phonological contexts, and phrases in isolation from approximately

181 1300 talkers of 107 regional backgrounds, with 10-12 recorded talkers per region and 5-15 recordings
 182 per vowel for each talker.

183 Unlike most existing databases, SwehVd focuses on a single regional variety, providing high
 184 resolution within and across talkers for this variety: SwehVd consists of $N=10$ recordings of each
 185 hVd word (for a total of 210 recordings for the 21 different hVd words) per talker. We note that this
 186 makes the demographic composition of SwehVd relatively homogenous, compared both to some
 187 other vowel corpora of Swedish (e.g., Bruce et al., 1999; Kuronen, 2000) and to some previous
 188 studies on normalization—a point we return to in the general discussion. Specifically, we target N
 189 = 24 male and female talkers each (current $N = 23$, all female) for a total targeted N of tokens
 190 = 10,080 (current $N = 4511$ tokens). The accompanying database contains first to third formant
 191 (F1-F3) measurements for each talker at five time points across each vowel, together with vowel
 192 duration and mean F0 over the entire vowel.

193 SwehVd follows the gross of research on normalization and uses hVd words for recording in
 194 order to minimise coarticulatory effects from the surrounding phonetic context. The hVd context
 195 was originally chosen for studies on English because the glottal /h/ in onset position minimizes
 196 supraglottal articulations (confirmed in e.g., Chesworth et al., 2003; Robb and Chen, 2009). Since
 197 then hVd words have played a central role in research on vowel production (e.g., Hillenbrand et
 198 al., 1995; Peterson and Barney, 1952) and perception (e.g., Malinsky et al., 2020; Peterson and
 199 Barney, 1952). Since Swedish onset /h/ is a glottal approximant (Riad, 2014) similar to English,
 200 the use of this context in SwehVd facilitates comparison to similar databases from other languages.
 201 It deviates, however, from the majority of previous studies on Swedish vowels, which have either
 202 not held phonetic context constant across vowels (e.g., Bruce et al., 1999), or have investigated
 203 vowel production out of context (Eklund and Traunmüller, 1997; Fant et al., 1969; Disner, 1980)
 204 or in different CVC contexts (e.g., kVp and pV k in Nordstrand et al., 2004; vV t , vV t t , fV t , fV t t ,
 205 in Behne et al., 1997).

206 The Swedish vowel inventory

207 The Central Swedish vowel inventory contains 21 monophthong vowels. Seventeen of these vowels
 208 form nine pairs distinguished by quantity (long and short): in Central Swedish, the two long vowels
 209 [ɛ:] and [e:] both neutralize to the same short vowel [ɛ] (resulting in a total of 17, rather than
 210 18, distinct vowels). The two variants of a pair are considered allophones, the selection of which
 211 is determined primarily by stress and syllable complexity. Quantity is neutralized in unstressed
 212 positions (Riad, 2014).⁵ Vowels lengthen in open word-final syllables, before morpheme-final single
 213 consonants, and in non-final syllables.

214 Additionally, there are four contextually conditioned allophones to [ɛ] and [ø]. Before /r/ (or any
 215 retroflex segment), both the long and short versions of these vowels lower to long and short [æ] and
 216 [œ], respectively. As shown in Table 2 (adapted from Riad, 2014), some long-short vowel pairs are
 217 described to differ not only in quantity but also in quality: generally, short vowels are described
 218 as more open and also more centralized, forming a more condensed vowel space. In ongoing work
 219 (Persson, 2023), we found this to be confirmed for SwehVd.

220 Several of the long vowels have been claimed to be diphthongized in Central Swedish (e.g., Elert,
 221 1981; Fant et al., 1969; Fant, 1971; Kuronen, 2000) and/or with consonantal elements (McAllister

⁵ This reflects the mainstream analytical position in present-day Swedish phonology. The opposite position, distinctive vowel quantity, has also been proposed (e.g., Linell, 1978, 1979; Schaeffer, 2005). This theoretical debate does not affect the interpretation of our results.

Table 2. The phonological characterization of long (left) and short (right) Central Swedish vowels (based on Riad, 2014)

	front	rounded	central	back		front	rounded	central	back
high	[i:]	[y:]		[u:]	high	[I]	[Y]	[e]	[ø]
mid-high	[e:]	[ø:]		[o:]	mid	[ɛ]	[ø]	[œ]	[ɔ:]
mid	[ɛ:]	[ø:]			low	[æ]		[a]	
low	[æ:]	[œ:]		[ɑ:]					

et al., 1974), though empirical evaluations of this claim have returned mixed results (Eklund and Traunmüller, 1997; Fant et al., 1969; Leinonen, 2010). In separate work (Persson, 2023), we have found evidence for diphthongization of some vowels in SwehVd, including some for which it has not previously been reported. Here we do not discuss this issue further since it has no direct consequences for the present study: normalization can be applied to both monophthongs and diphthongs, and listeners presumably use the full information contained in the dynamics of formant trajectories—rather than just aggregate point estimates (the simplifying information we make below, following other work on normalization).

Participants

Stockholm Swedish is a variety subsumed under Central Swedish—the regional standard variety of Swedish spoken in an area around and beyond Stockholm (eastern Svealand), including Mälardalssvenska, Sveamål, Uppsvenska, Mellansvenska (see e.g., Bruce, 2009; Elert, 1994; Riad, 2014). SwehVd targets the contemporary Central Swedish spoken in the larger Stockholm area.

Native talkers of Stockholm Swedish were recruited through word-of-mouth, flyers at Stockholm University Campus (see example flyer in SI, Participant recruitment), and online channels (accindi.se). Participants were selected based on the following criteria: L1 talkers of Swedish, born and raised in the greater Stockholm area or its surroundings, 20-40 years old (mean age=28; SD=5.35). Four of the participants were bilingual from birth, the L1s of each talker are provided in the database. All participants were reimbursed with a voucher to the value of SEK 100 after completing the recordings.

Recording procedure

Recording for the SwehVd database began in 2020 and is ongoing. The data were collected by the first author and Maryann Tan (Stockholm University). The hVd words were recorded together with another set of recordings targeting the production of Swedish word-initial stop voicing. Recording took place in a sound-attenuated room at the Multilingualism Laboratory, Department of Swedish Language and Multilingualism, Stockholm University.

Prior to recording, participants were informed about the study and given the possibility to ask questions before signing a consent form. They were then given instructions and seated at approximately 10 cm distance from an Audio Technica AT3035 microphone facing a computer screen. Words were presented one at a time, centered on screen, using PsychoPy software (Peirce et al., 2019). Participants were instructed to read the words with their natural voice as they appeared on screen. Each talker read the same 21 target words, with 48 mono- and bi-syllabic filler words interspersed. Each target word was repeated 10 times and each filler word was repeated five times, generating a total of 450 productions per talker, 210 target productions and 240 filler productions. We generated two pseudo-randomized lists of the words, each list divided into four different blocks. Words were blocked across block lists and randomized within block lists, with the constraint that

258 the same word would not appear more than twice in succession. Each participant was randomly
 259 assigned to one of the two lists. The pace of the presentation of the words was controlled by the
 260 experimenter, who was listening over Sennheiser HD215 headphones in the next room. A Yamaha
 261 MG102c mixing console with a built-in preamplifier was used together with a high-end ground
 262 isolator for preventing signal interference (Monacor FGA-40HQ). The speech was recorded at 44.1
 263 kHz in Audacity (Audacity, 2021). Each long sound file was split into individual short sound files
 264 of one word each. The boundaries of each file were slightly trimmed and the files were labelled with
 265 the target word. All sound files from the same talker were concatenated into one long file before
 266 further processing.

267 The complete list of target hVd words is provided in Table 5 in the SI. It consists of four
 268 real Swedish words, *hed*, *härd*, *hörd*, *hud* (English translations: *heath*, *hearth*, *heard*, and
 269 *skin*, respectively) and 17 phonotactically legal pseudowords. Following Swedish orthographical
 270 conventions for quantity, we used orthographic *hVdd* to elicit the short vowel allophone (e.g., *hudd*
 271 for [ø]) and orthographic *hVd* to elicit the long vowel allophone (e.g., *hud* for [u:]). This orthography
 272 reflects systematic phonological process of complementary quantity in Swedish (Riad, 2014). In
 273 order to elicit the contextual allophones to [ɛ] and [ø], we added the supradental [d̪] to elicit the
 274 long allophones (*härd*, *hörd*), and [r] to elicit the short allophones (*härr*, *hörr*). In a small-scale pilot
 275 preceding recordings, the expected transparency of the orthography for eliciting the long and short
 276 vowels was confirmed by three native talkers and one non-native talker of Swedish (these talkers did
 277 not participate in the study). However, *hodd* [u] and *hod* [u:] sometimes elicited [ɔ].⁶ We therefore
 278 decided to add instructions to the participants for these two words. When *hod* or *hodd* appeared
 279 on screen, a written guide indicating the target vowel appeared below the word in smaller font size:
 280 “*hod som i hot*”, “*hodd som i hosta*”, with *hot* and *hosta* being real Swedish words containing [u:]
 281 and [u], respectively.⁷ Whenever the experimenter noticed that the pronunciations clearly targeted
 282 another vowel, recordings were stopped and participants were reminded to carefully read the guide.

283 The recordings were divided into five blocks: one practice block and four recording blocks, with
 284 breaks in between. The purpose of the practice block was threefold: to familiarize the participants
 285 with the recording procedure, to adjust the recording level, and if necessary, to further instruct the
 286 participant (e.g., if the participant used inappropriate or inconsistent intonation or stress pattern).
 287 Each recording block consisted of either 110 (N=2 blocks) or 120 (N=2 blocks) trials. The length
 288 of each block was approximately eight minutes, for a total of roughly 30 minutes recording time
 289 per talker. After the recording, participants filled out a language background questionnaire and
 290 received their reimbursement.

291 Word and vowel segmentation

292 SweFA, a Swedish version of the Montreal Forced Aligner developed by Young and McGarrah
 293 (2021), was used to obtain estimates for word and segment boundaries. The boundaries were
 294 manually corrected by the first author (an L1 talker of Central Swedish). Following standard
 295 segmentation protocol and guidelines in Engstrand et al. (2001), segment boundaries were adjusted
 296 using spectrogram, waveforms and pitch and intensity tracks. The boundaries between /h/ and the
 297 vowel were adjusted to align with clear appearance of an F1, and the boundaries between the vowel
 298 and the coda consonant were aligned to a simultaneous rapid cessation of most or all formants.

⁶ The difficulty for some native talkers to produce [u] when reading *hodd* might be due to frequency effects. Forms with stressed [u] are few in the Swedish language, and phonotactically similar words are most often pronounced as [ɔ] (see e.g., Riad, 2014).

⁷ English translations: “*hod* as in threat”(phonologically [u:]), “*hodd* as in cough”(phonologically [u]).

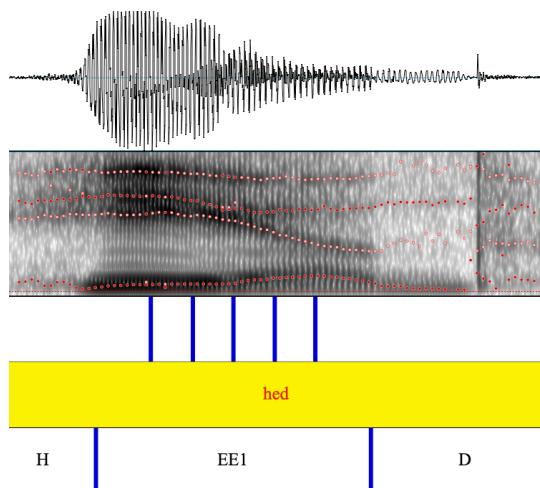


Figure 2. Example of Praat textgrid with annotated segment boundaries and measurement points for the automatic extraction of F1-F3 formant frequencies.

299 Extraction of phonetic cues

300 We used the Burg algorithm in Praat (Boersma and Weenink, 2022) to extract estimates of the
 301 first three formants (F1-F3) at five points of the vowel (20, 35, 50, 65, and 80 percent into the
 302 vowel). The following parameterization of the Burg algorithm was used:

- 303 • Time step (s): 0.01
 304 • Max. number of formants: 5
 305 • Formant ceiling (Hz): 5500
 306 • Window length (s): 0.025
 307 • Pre-emphasis from (Hz): 50

308 In addition to F1-F3, we automatically extracted vowel duration and the fundamental frequency
 309 (F0) across the entire vowel. The Praat scripts that extract this information are shared as part of
 310 the SWehVd OSF repository, allowing researchers to choose additional or alternative time points at
 311 which to extract formants.

312 In order to correct for measurement errors in the automatic extraction of cues, we estimated the
 313 joint multivariate distribution along all five extracted cues (F0, F1, F2, F3, and vowel duration)
 314 for each unique combination of vowel and talker. This approach allowed us to detect outliers
 315 relative to the joint distribution of the five cues for that vowel and talker. Points outside of the
 316 0.5th to 99.5th quantile of the multivariate Gaussian distribution of each vowel were identified,
 317 checked for measurement errors, and corrected. For measurements of the first three formants, we
 318 first checked the segmentation boundaries in the Praat textgrid and then manually measured new
 319 formant values using visual approximation of time points and Praat's function *Formant: Formant*
 320 *listing* or manually reading off the spectrogram. Segmentation boundaries were also checked for the
 321 identified vowel duration outliers. For measurements of F0, we extracted new estimated F0s across
 322 the vowel, after changing the pitch range settings. Given that there were still instances of pitch
 323 halving after measurement correction, in order to be conservative, we also checked all F0 values
 324 below the point of intersection between the two halves. The database available via OSF reports

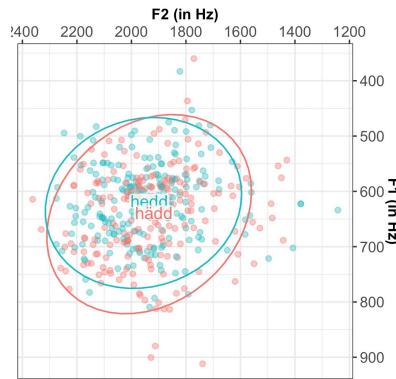


Figure 3. The *hedd* and *hädd* words in the SwehVd vowel data in unnormalized F1-F2 space. Points show recordings of the *hedd* and *hädd* words ([ɛ]) by the 23 female native talkers in the database, averaged across the five measurement points within each vowel segment. Word labels indicate word means across talkers. Since *hädd* and *hedd* resulted in the same allophone, we exclude *hädd* from this and all other visualizations below. This facilitates comparison of, for example, densities across vowels (see diagonal of Figure 5).

325 these corrected values (since the recordings and annotation grids are also available via OSF, other
 326 researchers can easily derive alternative measurements).⁸

327 The procedure of adding written guides to *hod* and *hodd* to facilitate vowel identification was
 328 mostly successful, however not for all talkers. Some talkers corrected themselves after one trial,
 329 others failed to produce the intended vowel altogether. The SwehVd database contains columns for
 330 both the targeted vowel category, and the vowel category that the talker actually produced (as
 331 annotated by the first author).

332 For the vast majority of talkers, *hädd* productions elicited the same vowel as *hedd* (see Figure 3).
 333 This confirms the common assumption that the short allophone to /e/ neutralizes with the short
 334 allophone to [ɛ] in Central Swedish.

335 Characterizing vowel productions in SwehVd

336 Figure 4 visualizes the vowel data from the SwehVd in F1-F2 space. The plot highlights the density
 337 of the Central Swedish vowel space, the categories are numerous and closely located. Category
 338 overlap is especially large among some of the high vowels (e.g., [i:] & [y:]; [u:], [o:] & [ø]). The
 339 contextually conditioned allophone [æ], almost completely overlaps with the long [ɛ:], whereas
 340 the contextual allophones to [ø] are more separated. Not all contextual allophones are articulated
 341 lower (higher F1) in relation to their phonemes (compare e.g., Riad, 2014). They are, however, all
 342 articulated further back (lower F2).

343 In line with Riad (2014, cf. Table 2 above), the short vowels are overall more centralized and
 344 form a more condensed space, whereas the long vowels are more dispersed.⁹ Differences in formant
 345 patterns between long and short vowels have been found to be smallest for the allophones to [ɛ] and
 346 [ø], and largest for /u/ ([u:] and [ø]), and /a/ ([a:] and [a]) (e.g., Fant, 2001; Kuronen, 2000). This

⁸ We note that outlier detection and correction was based on raw, rather than transformed/normalized, cue values. For the studies we report below, this potentially introduces a bias *against* normalization. If anything, the present study is thus likely to under-estimate the effects of normalization.

⁹ [i] and [y] are, however, more fronted than their long counterparts, which does not replicate previous descriptions of Central Swedish. We elaborate on this result in other ongoing work (Persson, 2023).

347 pattern does not entirely replicate here. We do see large differences in F1-F2 for the allophones to
 348 /u/ and /a/, but also substantial differences between [ɛ] and [ɛ:]. Formant differences are apparent
 349 even for some category distinctions for which quantity has been found to be the primary cue (as
 350 shown in perceptual studies, see e.g., Behne et al., 1997; Hadding-Koch and Abramson, 1964), ([ɛ]
 351 - [ɛ], [ø:] - [ø], [i:] - [i], and [o:] - [ɔ]).

352 Figure 5 visualizes the vowel data from the SwehVd database for all pairwise combinations of five
 353 cues: F0, F1, F2, F3 and vowel duration. As is to be expected, vowels differing in quality are most
 354 separated in the F1-F2 plot, indicating the two cues most important for vowel category distinction.
 355 However, the F1-F3 and F3-F2 plots both display less overlap between the high vowels [i:], [y:] and
 356 [u:], comparing to when plotted along F1-F2. The increased separation of these categories along F3
 357 in vowel production data could point to the importance of F3 for some category distinctions, as
 358 found in previous studies (see e.g., Fant et al., 1969; Fujimura, 1967; Kuronen, 2000, for [i:] and
 359 [y:] categorization).

360 Also as expected, duration is the primary cue that distinguishes vowel quantity: in the last column
 361 of Figure 5, the short vowels cluster on the left, and the long vowels on the right. They are separable,
 362 but overlapping. Overall, the short vowels seem to display less variability in duration than the long
 363 vowels, a common pattern for measures with a lower bound. In addition to duration, F1-F3 can
 364 also carry information about vowels differing in quantity. This is evident, for example, for [i:] vs. [i],
 365 [y:] vs. [y], [u:] vs. [ø], [a:] vs. [a], [ɛ:] vs. [ɛ] in F1-F2 space, and for [i:] vs. [i], [y:] vs. [y], [u:] vs. [ø]
 366 in F2-F3 space.

367 Finally, the densities along the diagonal of Figure 5 suggest that F0 carries the least information
 368 about vowel identity, exhibiting the least between-category separation, followed by F3. This, too, is
 369 not surprising: while some accounts use F0 to *normalize* F1 and F2 (e.g., Miller, 1989; Syrdal and
 370 Gopal, 1986), F0 is not considered an important cue to vowel identity by itself (for demonstrations
 371 that F0 can, however, have strong *indirect* effects on vowel categorization, see Barreda and Nearey,
 372 2012; Barreda, 2020).

STUDY 1: COMPARING THE EFFECTS OF NORMALIZATION ACCOUNTS ON BETWEEN- VS. WITHIN-CATEGORY VARIABILITY

373 In Study 1, we follow previous research on the effects of different normalization accounts on
 374 category variability, by evaluating how effectively different normalization accounts reduce the
 375 within-category variability of Central Swedish vowels. We visualize the vowel space under different
 376 normalization accounts, and assess the effects on vowel category variability by calculating a measure
 377 of category separability. To anticipate one take-home point of Study 1, the results highlight
 378 important shortcomings of separability indices in evaluating normalization accounts.

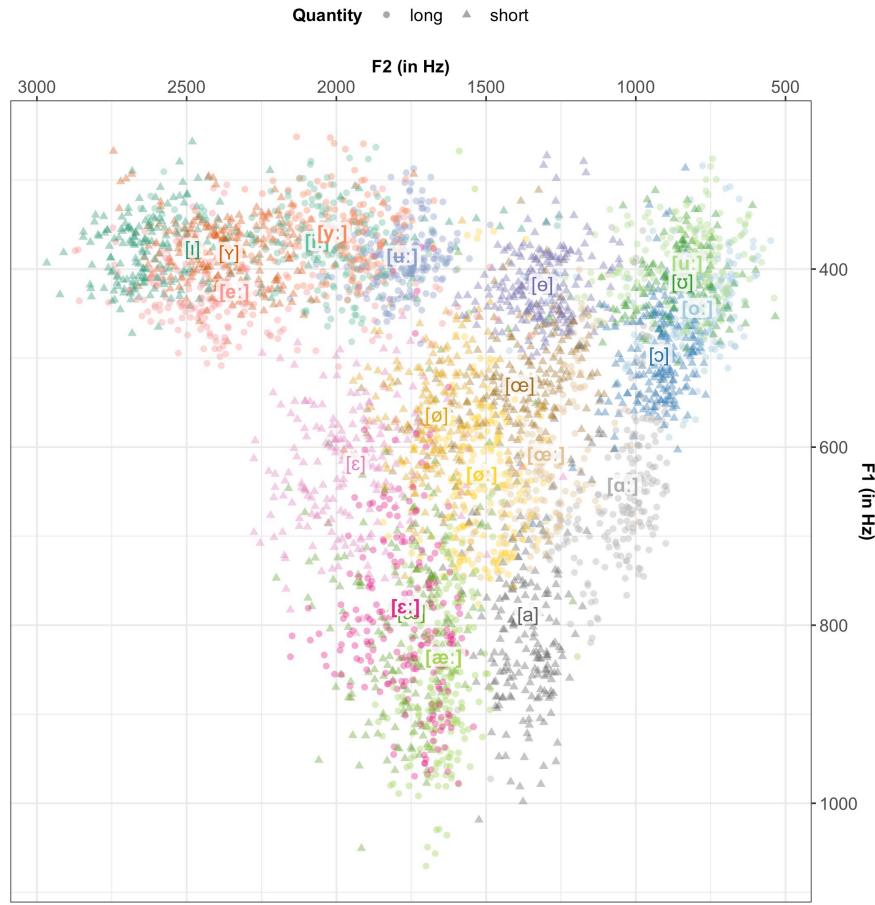


Figure 4. The SwehVd vowel data in unnormalized F1-F2 space. Points show recordings of each of the 21 Central Swedish vowels by the 23 female native talkers in the database, averaged across the five measurement points within each vowel segment. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Vowels that mismatched intended label are excluded (1.23976% of all recordings). Note that the F1 and F2 axes are reversed. We follow this convention whenever plotting vowels in the F1-F2 space.

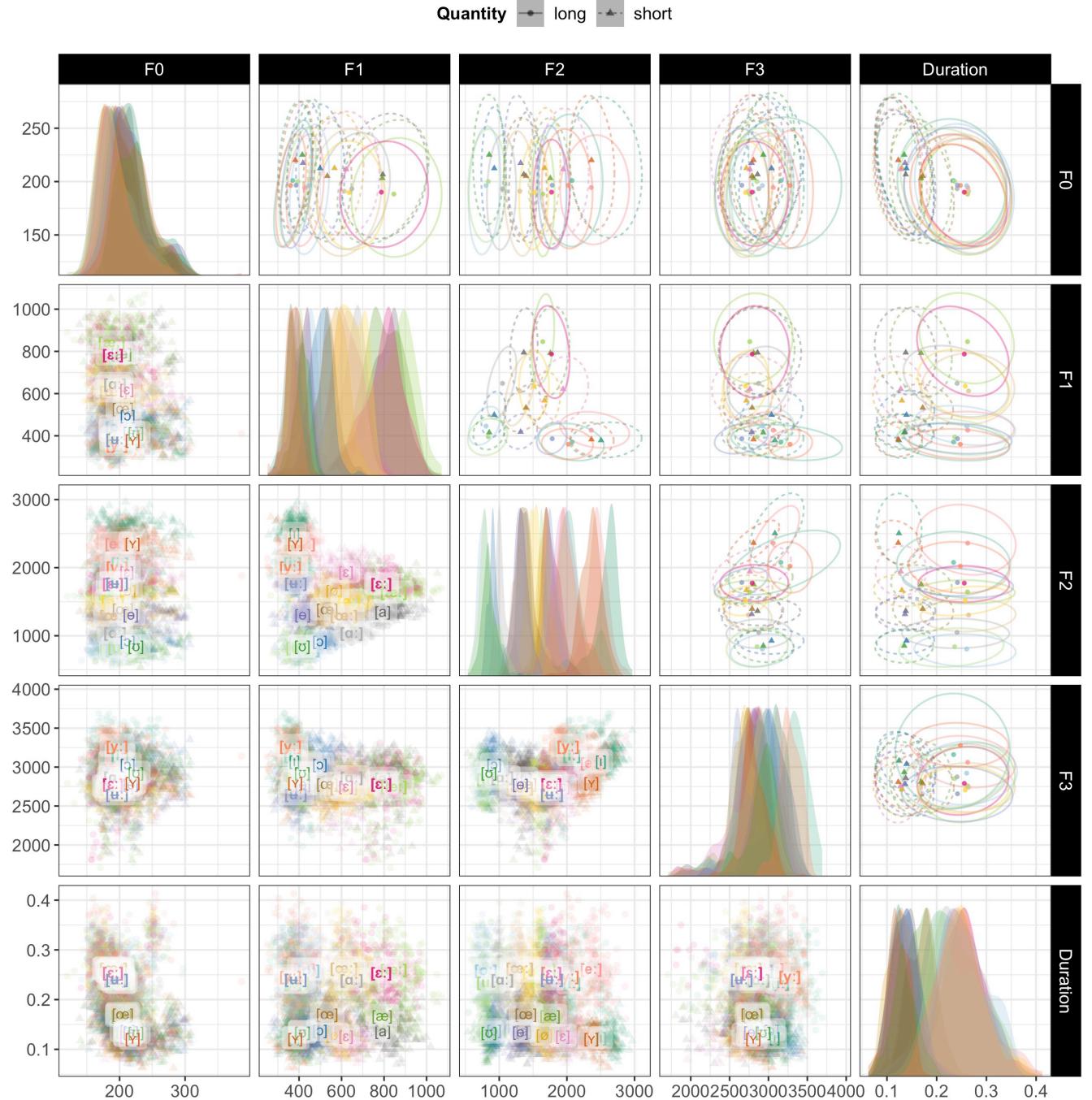


Figure 5. The same data as in Figure 4 but for all pairwise combinations of five cues: F0, F1, F2, F3, and vowel duration. The primary purpose of this figure is to provide an overview of the SwehVd data. Additionally, comparisons across the panels sheds light on which cues carry information about vowel quality and vowel quantity, respectively. Note that, unlike in Figure 4, axis directions are not reversed. **Panels on diagonal:** marginal cue densities of all five cues. **Lower off-diagonal panels:** each point corresponds to a recording, averaged across the five measurement points within each vowel segment. Vowel labels indicate category means across talkers. Long vowels are boldfaced. **Upper off-diagonal panels:** Same data as in the lower off-diagonal panels but showing bivariate Gaussian 95% probability mass ellipses around category means. This makes it more obvious, for example, that long and short vowels are primarily distinguished by vowel duration (top right panel).

Table 3. Previous studies comparing the effectiveness of normalization accounts in reducing within-category cue variability

Language investigated	Article	Speech materials	Normalization accounts	Approach	Best two performing
US English	Barreda & Nearey, 2018	120,000 simulated languages (of 5 or 9 vowels) modeled on Hillenbrand et al.'s (1995) data (98 female/male child/adult talkers * 12 vowels)	Nearey2, Lobanov, log-mean in linear regression framework	distance between means (Euclidian distance)	log-mean in linear regression framework (1), Nearey2 (2)
	Clopper, 2009	2 female/male talkers from Ohio (1 token * 10 vowels)	Bladon et al.'s scale factor of 1 Bark (1994), Syrdal & Gopal, Nordström & Lindblom, Nearey1, Nearey2, Watt & Fabricius, Gerstman, Lobanov, Miller	variance reduction (visual inspection)	Nearey, Watt & Fabricius, Gerstman, Lobanov (no order)
	Hindle, 1978	Peterson & Barney's (1952) database; 19 female/male talkers from Philadelphia + 60 telephone informants (minimum 3 tokens per category; analysis focus on /ay/)	Nearey2, Nordström-Lindblom, Sankoff-Shorrock-McKay	distance between means, variance reduction (regression)	Sankoff (1)
	Kohn & Farrington, 2012	Longitudinal data from 10 female/male African American talkers from North Carolina (approx. 10 tokens * 10 vowels * 5 ages)	Lobanov, Gerstman, Nearey1, Nordström & Lindblom, Syrdal & Gopal/Thomas, Watt & Fabricius	variance reduction (regression)	Lobanov (1), Gerstman, Watt & Fabricius (2)
	Labov, 2010	Peterson & Barney's (1952) database; Philadelphia/Linguistic Change and Variation project (120 female/male talkers, stratified for age, sociolinguistic factors)	Nearey2, Nordström-Lindblom, Sankoff-Shorrock-McKay	distance between means (F-statistics)	Sankoff (1), Nearey2 (2)
US English, Norwegian, Swedish, German, Danish, Dutch	Disner, 1980	Differing number of tokens, vowels, and phonetic contexts across the six languages	Gerstman, Lobanov, Nearey2, Harshman's PARAFAC model	variance reduction (visual inspection)	Nearey2 (1), Lobanov (2)
UK English	Fabricius, Watt & Johnson, 2009	20 old/young female/male talkers of Received pronunciation (11 vowels); 6 old/young female/male talkers of Aberdeen English (8 vowels in different phonetic contexts)	Watt & Fabricius, Lobanov, Nearey1		Lobanov (1), Watt & Fabricius (2)
	Flynn & Foulkes, 2011	20 old/young female/male Nottingham talkers (mean 180 recordings per talker; categories not reported)	log-transformation (base 10), log-transformation (natural), Mel, ERB, Bark (*2 gender-specific versions), Syrdal & Gopal, Nordström (*2 gender-specific versions), LCE, Gerstman, Lobanov, Watt & Fabricius (* 4 versions), lettER, Nearey (*4 versions)	variance reduction (SCV in talker-means)	Gerstman (1), LCE (2)
Russian	Lobanov, 1971	5 female/male talkers (9 vowels in different phonetic contexts)	linear compression or expansion (Fant, 1960), Gerstman, Lobanov	distance between means	Lobanov (1), Gerstman (2)

379 Table 3 lists previous studies that have compared normalization accounts in terms of their
 380 effectiveness in reducing inter-talker variability. These studies varied in the specific metrics they
 381 employed to assess the effects of normalization, the languages they studied, and the conclusions they
 382 arrived at. Two generalizations emerge from Table 3. First, transformation to a perceptual scale
 383 alone does not seem to be sufficient to reduce inter-talker variability (see also Adank et al., 2004;
 384 Carpenter and Govindarajan, 1993; Clopper, 2009; Escudero and Bion, 2007; Flynn and Foulkes,
 385 2011; Kohn and Farrington, 2012). Second, normalization accounts that include centering and/or
 386 standardizing seem to perform best in reducing inter-talker variability (see e.g. Barreda and Nearey,
 387 2018; Disner, 1980; Fabricius et al., 2009; Kohn and Farrington, 2012; Labov, 2010; Lobanov, 1971).
 388 When Lobanov and Gerstman normalization—both involving standardizing—were included in a
 389 study, they often rank among the top two performing accounts. Of note, Nearey normalization
 390 (Nearey, 1978) seems to perform well even though it does not involve the computationally more
 391 complex operation of standardizing. This suggests that simple centering of formants relative to the
 392 talker’s mean *might* be sufficient to achieve significant variance reduction (but see Disner, 1980 for
 393 Swedish, which is revisited in this study).

394 Methods

395 Speech materials

396 We use the SwehVd database with some exclusions. Since we are interested in assessing the effects
 397 of normalization, we excluded any productions on which the talker did not produce the targeted
 398 vowel. We then excluded all talkers ($N = 7$) with fewer than 5 remaining recordings for at least one
 399 of the vowels. Five of these talkers rarely ever produced the targeted [u] vowel despite our recording
 400 instructions (see Extraction of phonetic cues). Instead, they often mispronounced the vowel. This
 401 left data from 16 female native talkers, with on average 797 ($se = 2.5$) tokens per vowel (range =
 402 765 to 815), for a total of 16730 observations.

403 Since our goal is to obtain a reliable estimate of the formant values during the steady state of
 404 the vowel, we use only the three formant measurements extracted from the middle of the vowel (at
 405 35%, 50%, and 65% into the vowel).¹⁰

406 We also exclude all *hädd* productions, as they elicited the same vowel as *hedd* (see Extraction
 407 of phonetic cues). This way, we have about equally many tokens from all vowels, simplifying the
 408 cross-validation procedure presented below and facilitating visual comparisons across vowels in our
 409 figures.

410 Cues included in the normalization

411 We compare the effect of different normalization accounts on the variability of Central Swedish
 412 vowels under three different assumptions about the relevant cues. The first comparison follows most
 413 previous research and focuses on the two primary cues to vowel perception, F1 and F2. The second
 414 comparison considers F3 in addition to F1 and F2, following Adank et al. (2004), Barreda and
 415 Nearey (2018), Nearey (1989) and Syrdal (1985).¹¹ Finally, the third comparison includes F0 and
 416 duration in addition to F1-F3. Since Syrdal and Gopal (1986)’s bark-difference model only considers

¹⁰ While this is the approach most commonly employed in the literature, it has the potential downside that co-articulation might affect formant values at the measurement points differently for long and short vowels (since the long and short vowels differ in overall duration). An alternative approach would be to extract formants at fixed durations (e.g., 30 ms) after the vowel onset and before the vowel offset. Since the findings we present below do not indicate any systematic differences in the performance of normalization accounts between long and short vowels, we do not consider this issue further here.

¹¹ Some of these studies additionally included F0 (Adank et al., 2004; Nearey, 1989; Syrdal, 1985). However, since F0 is a cue that can display substantial cross-talker variability without directly contributing much information to vowel categorization (recall Figure 5), its inclusion can reduce the informativeness of the separability index. We therefore decided to add only F3 to F1-F2 in the second evaluation.

417 normalization along two dimensions—height, implemented as F1-F0, and backness, implemented
 418 as F2-F1—this account will only be included in the first comparison. Furthermore, given that C-
 419 CuRE is the only account that applies to any type of cue, we will consider duration as centered to
 420 each talker's mean (for the C-CuRE accounts), or as raw input (in ms; for all other accounts). To
 421 our knowledge, duration has not been considered in previous studies on normalization, presumably
 422 because the present study is the first to evaluate normalization accounts against a vowel system
 423 with a systematic quantity distinction (long vs. short vowels). We evaluate the effect on variability
 424 both separately for long and short vowels, and on all 21 vowels together.

425 Separability index

426 Previous studies have used different measures to assess the relative success of a normalization
 427 procedure in reducing inter-talker variability (see Table 3 and Nearey, 1989 for an overview on
 428 classification accounts). This includes assessing the reduction in variance or distance between
 429 means by visual inspection (e.g., Clopper, 2009; Disner, 1980; Hindle, 1978), or by calculating
 430 the reduction in within-category variance across talkers (e.g., Disner, 1980; Fabricius et al., 2009;
 431 Flynn and Foulkes, 2011; Hindle, 1978), or comparing the degree of separation between category
 432 means for unnormalized and normalized data, i.e., an F-ratio (e.g., Labov, 2010).

433 We will assess how distinguishable vowels become under different normalization accounts by
 434 calculating a separability index, as described in Equation (1). Following some previous studies
 435 (e.g., Labov, 2010), this separability index is essentially an F statistics, where the F statistics is
 436 the ratio of the within- and between-category variances:

$$\text{separability index} = \frac{\text{between category } MS}{\text{within category } MS} = \frac{\sum_{c=1, \dots, K} (N_c - 1)}{K - 1} \frac{\sum_{c=1, \dots, K} (\bar{x}_c - \bar{x})^2}{\sum_{c=1, \dots, K} \sum_{i=1, \dots, N_c} (x_{i,c} - \bar{x}_c)^2} \quad (1)$$

437 where K is the number of categories, N_c is the number of observations for category c , $x_{i,c}$ is the
 438 cue vector (for all cues considered in the calculation of the separability index) for observation i
 439 of category c , \bar{x}_c is the cue mean vector for category c , and \bar{x} is the overall cue mean vector. We
 440 calculated this separability index separately for each combination of normalization account, cues,
 441 and training-test fold, as described next.

442 Guarding against over-fitting: cross-validation

443 As shown in Table 1, many of the normalization accounts involve parameters that are set
 444 based on the data (e.g., Gerstman, 1968; Lobanov, 1971; McMurray and Jongman, 2011; Miller,
 445 1989; Nearey, 1978). This raises the question of how much these parameters can be affected by
 446 outliers, or other issues such as over-fitting to the sample. Unlike previous work, we thus use 5-
 447 fold cross-validation to obtain 5 separate estimates of the separability index for each combination
 448 of normalization procedure and cues. Specifically, we randomly split the data for each unique
 449 combination of talker and vowel into 5 even parts (folds). On each of the five folds, we then fit the
 450 normalization parameters based on four of the folds (the training data) and evaluated the effects
 451 of the normalization on the fifth fold (the test data). This resulted in five separability indices for
 452 each combination of normalization procedure and cues.

453 Results

454 Visualizing the distribution of vowel productions

455 Figures 6 and 7 visualize the Central Swedish vowels in the test data, after applying the 15
456 different scale-transformations and normalization accounts for a visual inspection. For this purpose,
457 we focus on F1 and F2 only. The SI includes similar pairwise correlation plots of all cues for all
458 different normalization accounts we compare (see SI, Correlation matrices).

459 Visual inspection suggests a few initial observations. The most striking difference is perhaps
460 between intrinsic normalization accounts (Syrdal and Gopal, 1986; Miller, 1989) and all other
461 approaches, though it is not immediately visually obvious which type of approach achieves better
462 separability. Second, transforming the vowels to a different perceptual scale does not seem to
463 affect the vowel distributions much, besides a minor decrease in category variance for some of the
464 vowels. Some transformations bring the vowel categories closer together, towards the center of the
465 vowel space, e.g., ERB and semitones. Third, centering formants by subtracting each talkers' mean
466 (McMurray and Jongman, 2011; Nearey, 1978) reduces some of the category variance, and as a
467 result, increases the category separability. Transforming the vowel data into different scales prior
468 to centering also seems to further improve separability (compare e.g., C-CuRE (Hz) and C-CuRE
469 (semitones)). Overall, the top two performing accounts across the long and short vowels appear
470 to be Lobanov (1971) and Nearey (1978). However, even for the best performing normalization
471 accounts, there is still considerable category overlap. This involves some of the high long vowels,
472 and some of the mid-center short vowels. This highlights the need to more systematically quantify
473 the effects of normalization, as we do next for category separability and then in more depth in
474 Study 2.

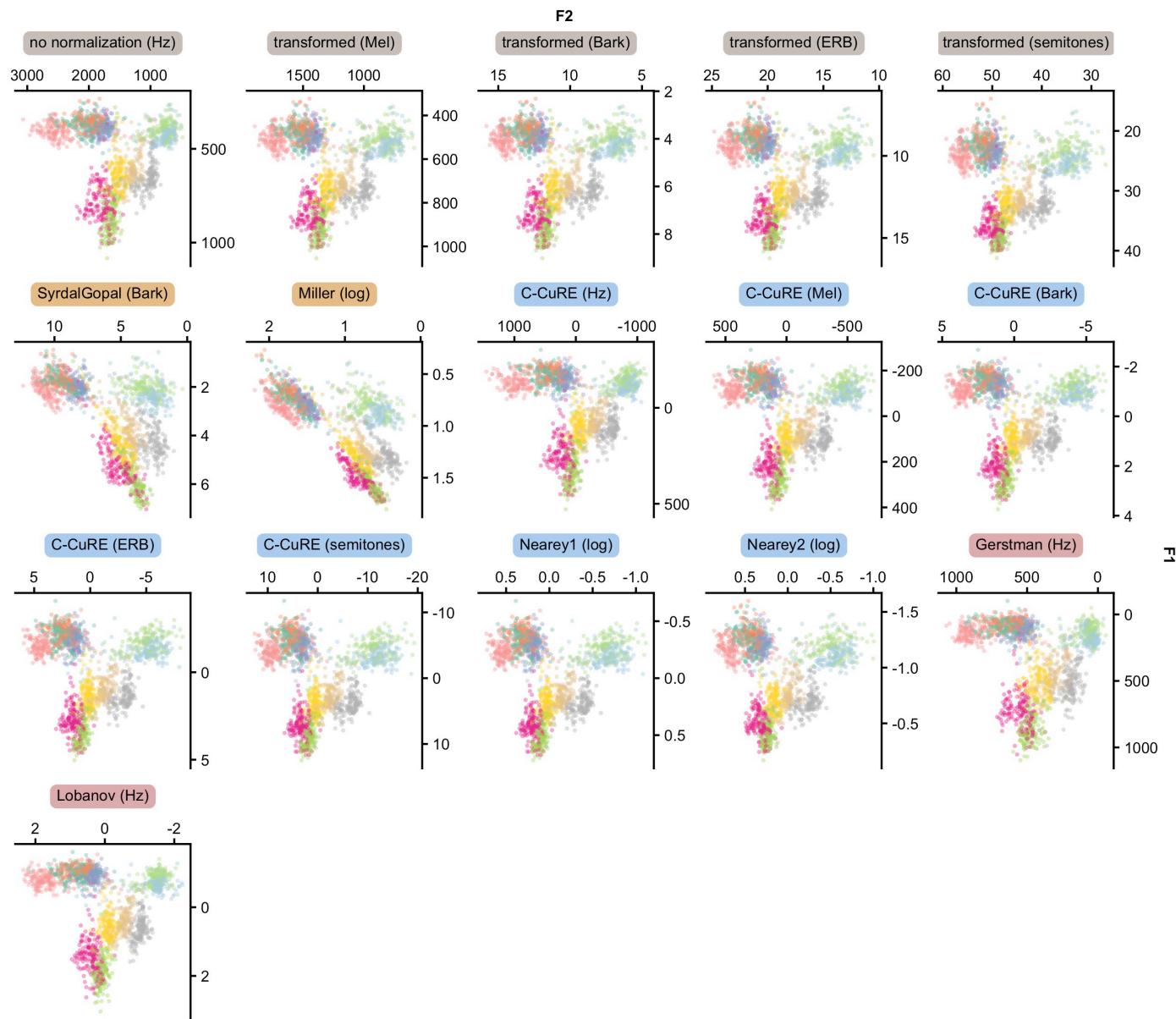


Figure 6. The 11 long vowels of Central Swedish when F1 and F2 are left unnormalized or transformed into a perceptual scales (grey), intrinsically normalized (yellow), or extrinsically normalized through centering (blue) or standardizing (purple). Each point corresponds to one recording, averaged across the five measurement points within each vowel segment. Each panel combines the data from all five test folds.

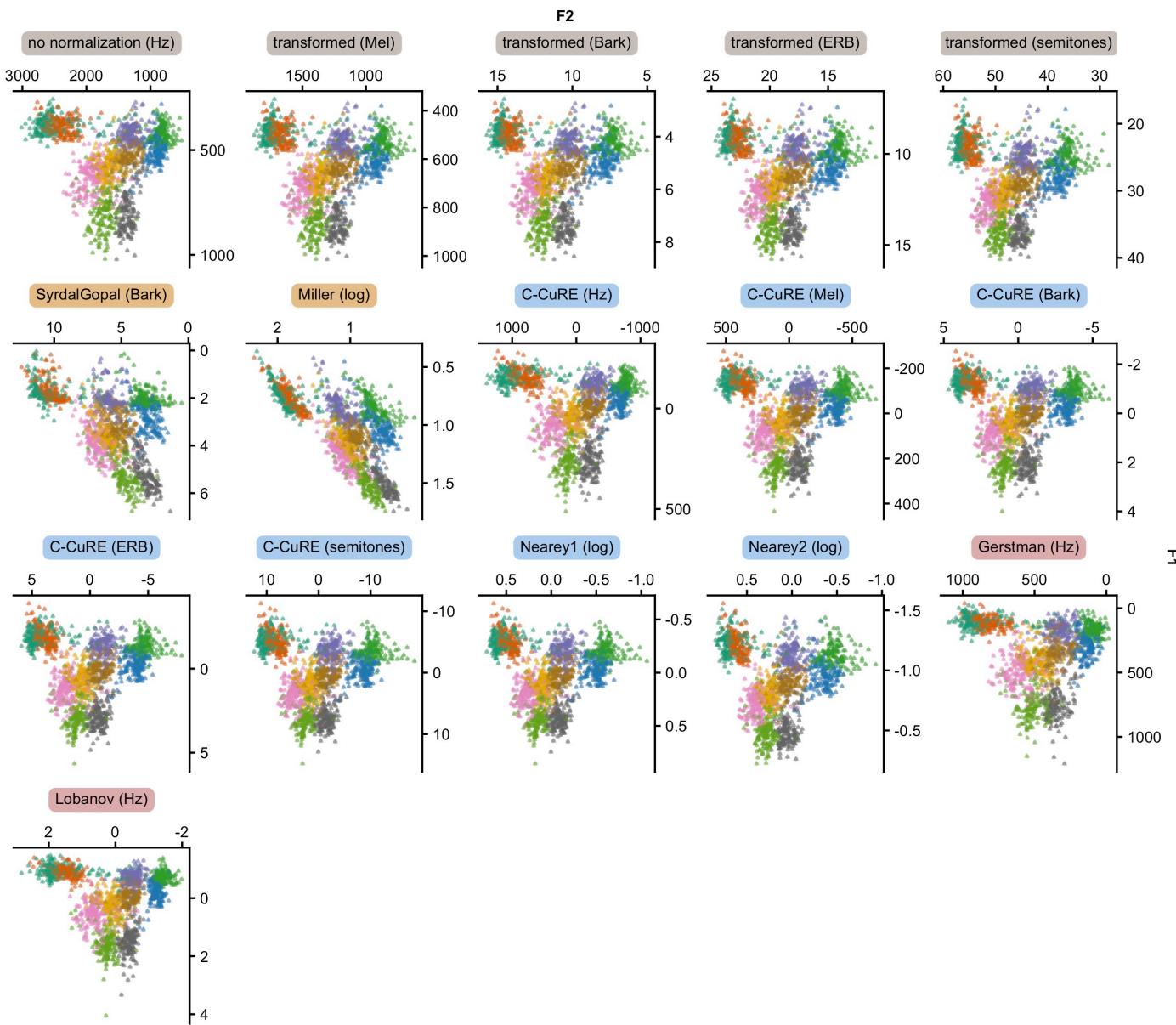


Figure 7. The 10 short vowels of Central Swedish when F1 and F2 are left unnormalized or transformed into a perceptual scales (grey), intrinsically normalized (yellow), or extrinsically normalized through centering (blue) or standardizing (purple). Each point corresponds to one recording, averaged across the five measurement points within each vowel segment. Each panel combines the data from all five test folds.

475 The effect of normalization accounts on the variability of Central Swedish vowels

476 In order to assess the effects of normalization accounts on category separability, we calculate a
477 separability index under different assumptions about the relevant cues and the size of the vowel
478 space (the long and short vowels separately, or the entire space). Before we evaluate how category
479 separability is affected by normalization in F1-F2, F1-F3, and F0-F3 and duration space, we look at
480 how the normalization accounts affect the separability of vowels along each cue separately (Figure
481 8). As we show below, this is helpful in understanding the subsequently presented results for
482 combinations of cues.

483 For F1 (first row of Figure 8), we see a clear advantage for centering (in blue) and standardizing
484 (in purple) compared to transformations (in grey) and intrinsic accounts (in yellow). In particular
485 Lobanov normalization seems to maximize category separability along F1, at least for the long
486 vowels and all vowels together. Notably, the accounts pattern differently along F2 (second row of
487 Figure 8). Overall, differences between accounts are much smaller along F2, and the clear advantage
488 of centering and standardizing accounts along F1 does not extend to F2.

489 An altogether different picture is observed for F3. Compared to F1 and F2, the intrinsic account
490 (Miller) performs substantially better in separating categories along F3, while all other accounts
491 perform poorly. This result is surprising: one of the downsides of intrinsic approaches that has been
492 noted in previous work is their sensitivity to measurement error (Thomas and Kendall, 2007). This
493 sensitivity is caused by the fact that intrinsic accounts use a single measurement for normalization,
494 rather than the less noisy estimates resulting from aggregating across segments that are used in
495 extrinsic accounts. Since the third formant is often described as more difficult to reliably estimate
496 than other formants (leading to more measurement error), F3 would be expected to be particularly
497 affected by this weakness of intrinsic accounts.

498 Yet, further visualization in Figure 9 confirms that F3 indeed separates categories particularly
499 well when intrinsic normalization is applied. Compared to other accounts, Miller (1989) seems to be
500 particularly successful in separating vowels that differ in lip rounding. For example, Miller (1989)
501 separates two clusters among the high and mid-high vowels, one consisting of the back vowels [o:]
502 and [u:], and the other one of the front [i:] and rounded [y:] and [ɯ:]. One possible explanation
503 for this result is that intrinsic normalization is indeed particularly effective for F3, and that our
504 correction of measurement errors—equally applied to all formants—effectively reduced the issue
505 with F3 measurement errors (presumably the human brain, too, can do better than an uncorrected
506 Praat algorithm without error correction). As we show below, this result for F3 carries over to any
507 combination of cues that includes F3. It is, however, an artifact of using category separability to
508 assess the effectiveness of normalization, as we show in Study 2. We elaborate on this issue in the
509 discussion of Study 1.

510 Returning to Figure 8, normalization does not increase category separability for F0. This is
511 expected given that F0 is known to affect vowel separability primarily through its indirect influence
512 on the interpretation of other formants (e.g., Barreda and Nearey, 2012; Barreda, 2020). Finally,
513 for duration all of the C-CuRE accounts group together against the remaining accounts. This, too,
514 is expected since all other accounts are formant-specific and thus do not normalize duration. In
515 summary, the five cues contribute to category separability in different ways, and this is reflected
516 in varying effectiveness of different normalization accounts. We also note that the best performing
517 normalization account for any combination of cues and vowel qualities is typically never significantly
518 better than the next best performing model (the 95% confidence intervals of the best model overlap

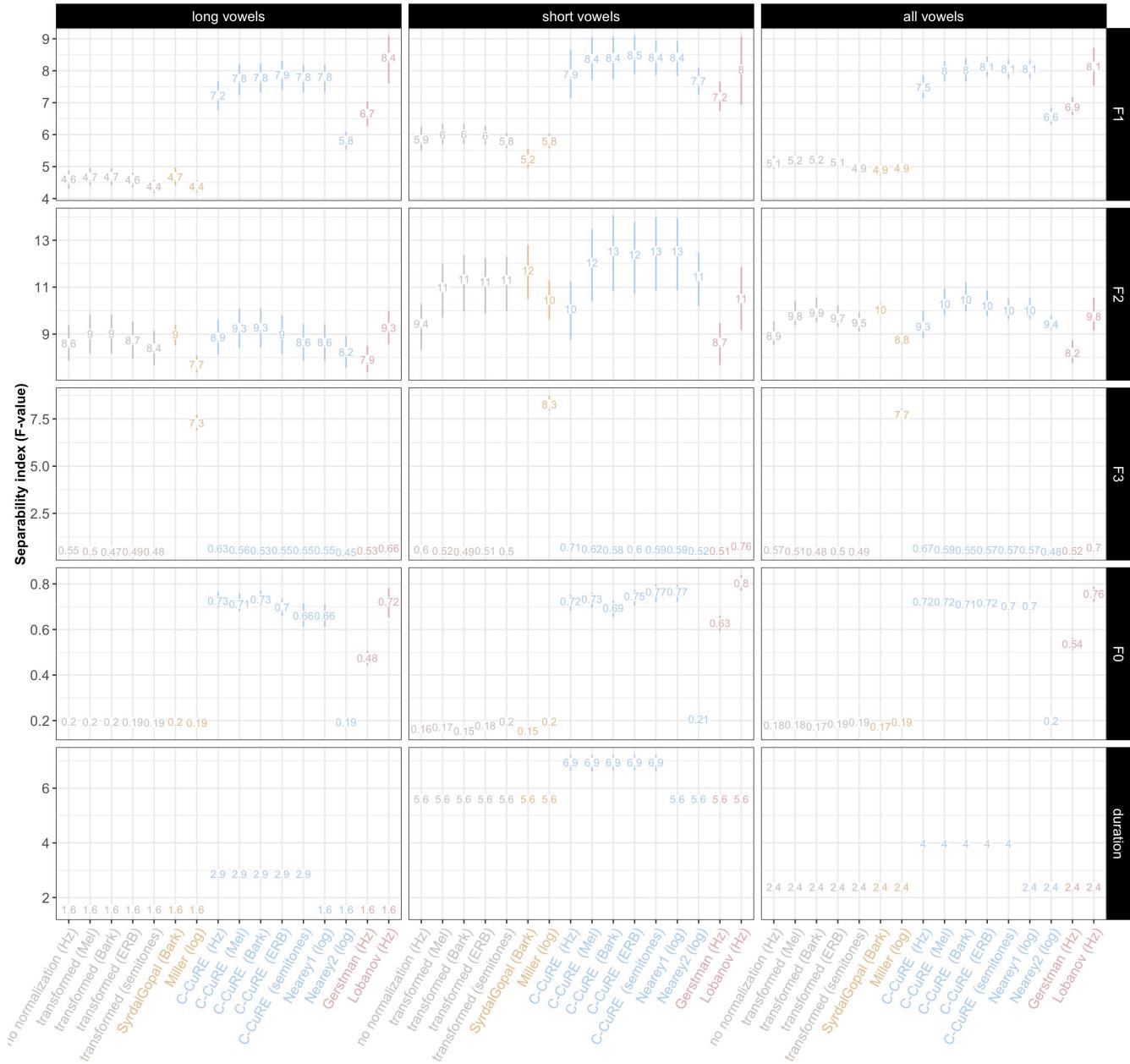


Figure 8. Separability indices by normalization accounts for long vowels, short vowels, and long and short vowels together (columns), shown for each of the five cues considered in this study (rows). Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. Note that the ranges of the y-axes varies across plots.

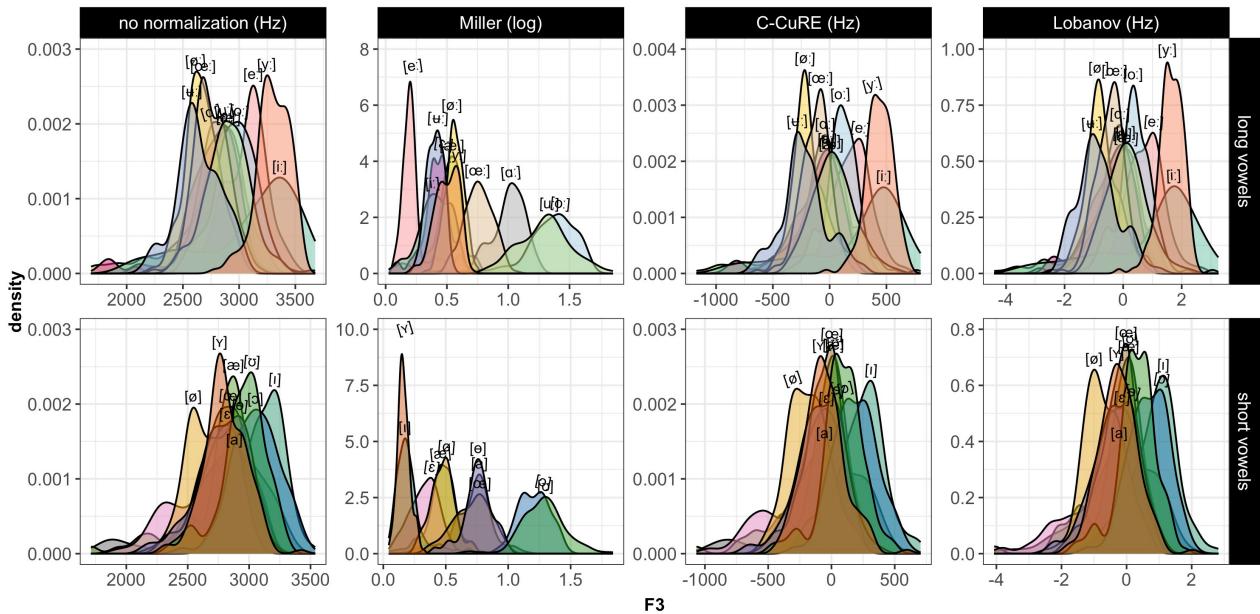


Figure 9. Category densities along F3 illustrates the effectiveness of vowel-intrinsic normalization for this cue. Here shown for Miller, compared to vowel-extrinsic accounts that center and/or standardize cues. For reference, densities in the absence of normalization are also shown.

519 with the mean of the next best model). In fact, for many combinations of cues and vowel qualities,
520 many of the models perform similarly.

521 Next, we summarize how normalization affects category separability when combinations of the
 522 fives cues are considered. Figure 10 shows the separability index for the different normalization
 523 accounts for three different combinations of cues. For the first row of Figure 10, we followed most
 524 previous research in assessing category separability for the combination of F1 and F2 (e.g., Disner,
 525 1980; Fabricius et al., 2009; Flynn and Foulkes, 2011; Hindle, 1978; Labov, 2010). Accounts
 526 that center against the talker’s overall formant mean (in blue) are among the best performing
 527 normalization accounts. No matter the assumed perceptual scale, centering always improves
 528 category separability. Standardizing accounts (in purple), primarily Lobanov (1971), also perform
 529 well at separating categories, more so for the long vowels. However, scale transformations (in grey),
 530 and intrinsic accounts (in yellow), do not improve category separability compared to unnormalized
 531 Hz, at least not when assessed on the long vowels or the entire vowel space.

532 The remaining rows of Figure 10 compare normalization accounts when F3 (second row) or F0, F3,
 533 and duration are included (third row). Overall, the category separability is now lower, a result of
 534 how the accounts affect category separability along the cues added (see Figure 8). The most drastic
 535 change in performance concerns the intrinsic Miller (1989) and the standardizing accounts. When
 536 including F3, Miller (1989) performs as well or better, in absolute numbers, as when evaluated on
 537 only the combination of F1 and F2, thereby increasing its performance relative to other accounts.
 538 This increase in performance might be particularly pronounced for languages like Swedish, where F3

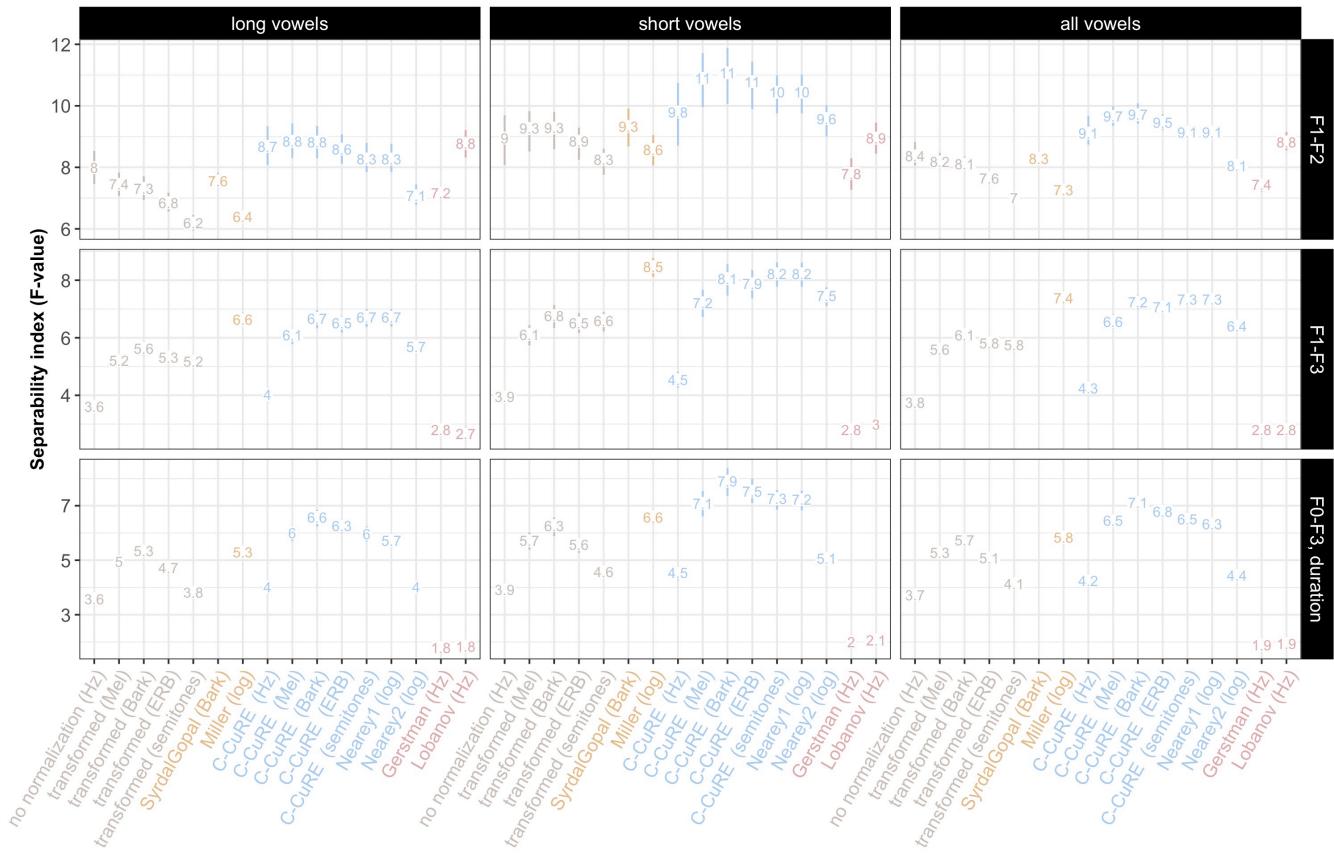


Figure 10. Separability indices by normalization accounts for long vowels, short vowels, and both long and short vowels together (columns) shown for three different combinations of cues (rows). Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. Note that the ranges of the y-axes varies across plots.

539 carries important information about lip rounding and thus vowel identity. In contrast, performance
 540 of standardizing accounts drops substantially if F3 or any other cue besides F1 and F2 is included.¹²
 541 This mirrors what was found when assessing category separability separately for each cue (Figure
 542 8).

543 Finally, looking across all three rows, category separability is consistently higher for short
 544 than long vowels. The same pattern is evident for each cue separately in Figure 8. This result
 545 might initially be puzzling, given that previous descriptions of Central Swedish vowel inventories
 546 characterize the inventory of short vowels as being more centralized and more densely clustered
 547 (e.g., Kuronen, 2000; Riad, 2014). Indeed, this claim seems to hold for SwehVd—compare Figures 6
 548 and 7. However, the short vowels exhibit less category variability and less category overlap, making
 549 them overall more separable.

¹² We confirmed this by conducting additional comparisons using only F1, F2 and F0, or only F1, F2 and duration. For both of these comparisons too, we found that standardizing accounts perform poorly.

550 Discussion

551 When only F1 and F2 are considered, as in most previous work on vowel normalization, we find
552 that extrinsic centering and standardizing accounts achieve the best category separability. Within
553 these two types of accounts, there is considerable variability. For example, among the intrinsic
554 accounts, Miller performs worse than Syrdal & Gopal, among the extrinsic accounts, versions of
555 C-CuRE seem to consistently perform best. It is also worth noting, however, that there is never a
556 single account that performs significantly better than all other normalizations. This points to the
557 inherent similarities across normalization accounts, and perhaps limitations of the approach taken
558 here (and in some previous work). We return to this point in the general discussion. Regardless of
559 these caveats, the findings for F1 and F2 of Study I, revise the results of Disner (1980) for Swedish,
560 and instead replicates previous findings for the other Germanic languages in Disner's sample as
561 well as the majority of previous studies on other languages (e.g., Fabricius et al., 2009; Flynn and
562 Foulkes, 2011; Labov, 2010).

563 However, when F3 is considered along with F1 and F2, this result does no longer hold. Key
564 to understanding this result and what it says about the suitability of category separability as a
565 measure of normalization accounts is Figure 8: while extrinsic normalization performs better than
566 other approaches for F1 and F2, the absolute differences in performance are small compared to
567 the advantage of the intrinsic account observed for F3. Combined with a seemingly innocuous
568 aspect of the separability index in Equation (1), this allows separability along F3 to dominate
569 separability along the other cue dimensions. Our separability index takes the *sum* of (squared)
570 distances along each cue dimension, essentially assuming that the effect of all cues is simply a sum
571 of each cue's effect considered separately. This means that the separability index cannot capture
572 the *joint* effect of cues—whether, for example, one cue effectively separates one set of categories
573 and another cue separates another set of categories, rather than both cues separating the same
574 categories. The separability index thus cannot recognize, for example, that F1 and F2 capture
575 largely complementary aspects of the vowel inventory (as evident in, for example, Figures 6 and 7).

576 This is not the only deficiency of the separability index or similar measures of category variability.
577 The use of *squared* distances means that even a small number of observations located far away from
578 the category mean can disproportionately affect the index. Consider, for example, the F3 densities in
579 Figure 9. For non-intrinsic normalizations, some categories have low but non-zero densities far away
580 from the mode. Because of the use of squared distances, this results in low category separability
581 for these normalization accounts despite the fact that observations with such cue values are rare
582 and thus not expected to have a large effect on the *average* perceptual separability of vowels. For
583 the same reason (the use of squared distances), category separability can be high even if a cue
584 separates only a small subset of categories (as is the case for F3), compared to another cue that
585 more gradually separates *all* categories (as is the case for F1 and F2; see Figure 5).

586 Finally, measures like the separability index suffer from a conceptual issue: the goal of speech
587 perception is presumably not to reduce cue variability around the category mean but rather to
588 increase the probability of correct recognition (of both linguistic and social information, where we
589 focus on the former here). These two goals are not the same (see also discussion in Barreda, 2020).
590 In sum, indices of variability and category separability like that in Equation (1) fail to adequately
591 assess the expected consequences of normalization for perception, which is the primary interest of
592 this study. Study 2 addresses this issue.

STUDY 2: COMPARING THE EXPECTED EFFECTS OF NORMALIZATION ON PERCEPTION

593 Study 2 takes a more direct approach to evaluating the expected effects of normalization for the
594 perception of vowels. Specifically, we use a general model of speech perception, Bayesian ideal
595 observers (e.g., Clayards et al., 2008; Nearey and Hogan, 1986; Norris and McQueen, 2008), to
596 predict the vowel identities in the SwehVd database under different normalization accounts. We
597 then compare normalization accounts based on the recognition accuracy that they achieve when
598 the (un)normalized cues are fed into the otherwise identical categorization model. We evaluate the
599 same normalization accounts as those investigated in Study 1. And, as in Study 1, we repeat this
600 comparisons for different combinations of cues, and while categorizing just long vowels, just short
601 vowels, or both long and short vowels together.

¹³ Barreda & Nearey (2018) identify a mistake in the implementation of the Nearey2 account in Adank et al. (2004), so that the relative performance of Nearey2 reported by Adank and colleagues should be interpreted with caution.

Table 4. Previous studies comparing normalization accounts in terms of their predicted consequences for perception

Language(s) investigated	Article	Speech materials	Normalization accounts	Approach	Accuracy assessed	Best two performing
US English	Barreda, 2021	synthesized stimuli representing 6 talker types (based on data from 30 female/male talkers of California English (15 tokens * 11 vowels))	Nearey2, Watt & Fabricius, Lobanov	regression	against perceived category	Nearey2 (1), Watt & Fabricius (2)
	Carpenter & Govindarajan, 1993	Peterson & Barney's (1952) database, 75 female/male child/adult talkers (2 tokens * 10 vowels)	Bark, Mel, ERB, 2 log-transformations, Syrdal & Gopal, Miller, Nearey1, Nearey2, Gerstman, linear transformation (Watrous, 1993)	fuzzy ARTMAP, K-nearest neighbour		linear transformation (1), Nearey1 (2)
	Cole, Linebaugh, Munson & McMurray, 2010	10 female/male talkers (3 tokens * 2 target vowels * 4 context vowels * 6 consonants)	C-CuRE	regression		C-CuRE (1)
	Johnson & Sjerps, 2021	Peterson & Barney's (1952) database, 75 female/male child/adult talkers (2 tokens * 10 vowels); Hillenbrand et al.'s (1995) database, 138 female/male child/adult talkers (1-3 tokens * 12 vowels)	Mean λ , F3 anchor, F1 anchor, Mean F* anchor (Sussman, 1986), Nordström, VTLN (Lammert & Narayanan, 2015), Nearey2, Gerstman, VTLN (ΔF), Nearey1, Watt & Fabricius, Lobanov, Miller, Syrdal & Gopal	support vector machine classification models	against intended category	Lobanov (1), Watt & Fabricius (2)
	McMurray, Cole & Munson, 2011	Cole et al. (2010) database, 10 female/male talkers (1 token * 2 target vowels * 4 context vowels * 6 consonants)				
	McMurray & Jongman, 2016	Jongman et al. (2000) database, 10 female/male talkers (1 token * 4 vowels * 8 fricatives)	C-CuRE	regression		C-CuRE (1)
	Nearey, 1989	synthesized stimuli of male child/adult talker (based on male talker data from Fant, 1973, and Peterson & Barney, 1952)	intrinsic normalization, extrinsic normalization	response patterns (F-ratio)	against perceived category	extrinsic effects (1), intrinsic effects (2)
	Richter, Feldman, Salgado & Jansen, 2017	models based on Clopper & Pisoni's (2006) NSP vowel corpus, 60 female/male talkers, 6 varieties (5 tokens * 10 vowels); perceptual data from Feldman et al., 2009 (synthesized stimuli of male talker)	Vocal Tract Length Normalization (VTLN), Lobanov	discrimination model likelihoods		VTLN (1), Lobanov (2)
	Syrdal, 1985	Peterson & Barney's (1952) database, 75 female/male child/adult talkers (2 tokens * 10 vowels)	log-transformation, Bark, Syrdal's bark-difference model, Miller (2 accounts), Nearey1, Nearey2, Gerstman	linear discriminant analysis		Nearey1 (1), Nearey2 (2)
Brazilian Portuguese & US English	Escudero & Hoffman Bion, 2007	models trained on 400,000 F1-F2 combinations generated on recordings of 8 female/male talkers (20 tokens * 7 vowels and 15 tokens * 11 vowels)	Nearey1, Lobanov, Gerstman	constraint rankings	against intended category	
Dutch	Adank, Smits & van Hout, 2004	160 female/male talkers, 8 varieties (2 tokens * 9 vowels)	log-transformation, Bark, Mel, ERB, Syrdal & Gopal, Lobanov, Nearey1, Nearey2 ¹³ , Gerstman, Nordström, Miller	linear discriminant analysis		Lobanov (1), Nearey1 (2)

602 Table 4 lists previous studies that have compared normalization accounts in terms of their
 603 expected consequences for perception. While some previous studies have evaluated normalization
 604 accounts against listeners' responses in perception experiments (Barreda, 2021; McMurray and
 605 Jongman, 2016; Nearey, 1989; Richter et al., 2017), the majority of previous work has evaluated
 606 accounts against the category intended by the talker—i.e., against production data. Here, we follow
 607 this majority (though we plan to conduct perception studies on Central Swedish in the future).

608 Across languages and approaches, the studies in Table 4 yield similar results as those summarized
 609 in Table 3 for category separability. Normalization improves recognition accuracy compared
 610 to unnormalized cues. Perceptual transformations perform worse than intrinsic or extrinsic
 611 normalization at predicting perception from production data. Furthermore, the same normalization
 612 accounts that have been found to be most successful at reducing inter-talker variability (e.g.,
 613 Lobanov, 1971; Nearey, 1978) are also often found to achieve higher recognition accuracy (e.g.,
 614 Adank et al., 2004; Escudero and Bion, 2007; Johnson and Sjerps, 2021; Syrdal, 1985). Despite
 615 these overall similarities in results, it is important to keep in mind that the two approaches—
 616 category separability and models of perception—do not *have* to yield the same results. In the
 617 present study, we go beyond previous work by modeling perception of both vowel quality and
 618 vowel quantity over a particularly dense vowel space, and by considering additional cues.

619 As also shown in Table 4, previous work has employed a number of model types to compare the
 620 expected effects of normalization on perception, ranging from models based on phonological theory
 621 (e.g., optimality theory, Escudero and Bion, 2007), to more general models of categorization (e.g.,
 622 linear discriminant analysis, Adank et al., 2004; Syrdal, 1985; k-nearest neighbors as in exemplar
 623 theory or ARTMAP, Carpenter and Govindarajan, 1993; Bayesian inference, Kleinschmidt et al.,
 624 2018; Richter et al., 2017; support vector machine classification models, Johnson and Sjerps, 2021),
 625 to general frameworks for data analysis (e.g., regression, Cole et al., 2010; McMurray and Jongman,
 626 2016). We use ideal observers, rather than other approaches, because *all* of their degrees of freedom
 627 can be estimated from SwehVd. In contrast, k-nearest neighbor categorization requires the choice
 628 of a similarity metric, which can introduce one or more degrees of freedom into the modeling,
 629 and requires a choice for k . Similarly, linear discriminant analysis or regression introduce at least
 630 one degree of freedom for each cue considered. This means that any comparison of normalization
 631 accounts needs to be conducted over the entire range of possible values for these degrees of
 632 freedom, making comparisons computationally more demanding and interpretation of the results
 633 more difficult.¹⁴ Bayesian ideal observers avoid this issue because of their assumption that listeners
 634 use and integrate cues *optimally*. As a consequence, the predicted posterior probabilities of all
 635 categories for a given acoustic input are determined by the combination of (1) the category-specific
 636 distribution of cues in the previous input and (2) the cue values of the input. Next, we describe
 637 ideal observers in more detail. After introducing ideal observers, we detail how this approach avoids
 638 the pitfalls of the separability index employed in Study 1.

639 Methods

640 Speech materials

641 Study 2 employs the same speech materials as in Study 1.

¹⁴ The consequences of additional degrees of freedom that we describe here hold for comparisons of normalization accounts against *production* data (as done here and in the majority of research on normalization).

642 Ideal Observers to predict the consequences of normalization for perception

643 Ideal observers provide an analytical framework for estimating how a rational listener would
 644 optimally behave in response to input (here: *n*-way alternative forced-choice categorization). Ideal
 645 observer models have been found to provide a good qualitative and quantitative fit against human
 646 speech perception (e.g., Clayards et al., 2008; Feldman et al., 2009; Kleinschmidt and Jaeger,
 647 2015; Kronrod et al., 2016; Norris and McQueen, 2008; Xie et al., 2021). Unlike most other
 648 models of speech perception, ideal observers in their simplest form—as employed here—have zero
 649 degrees of freedom in the link from production to perception: once the ideal observer is trained on
 650 phonetic data from a database of *productions*, its predictions about *perception* are not mediated by
 651 additional parameters (unlike, e.g., exemplar models, connectionist accounts, or neural networks).
 652 This removes researchers' degrees of freedom in the evaluation of normalization accounts, which is
 653 the reason we chose ideal observers for Study 2. We emphasize, however, that other researchers can
 654 download the R markdown document for this article (which contains the R code for our models)
 655 from OSF and substitute any other perceptual model for the ideal observers to assess the extent
 656 to which our choice of computational framework affects the findings reported below.

657 In line with influential theories of speech perception (e.g., exemplar theory, Johnson, 1997;
 658 Bayesian accounts, Luce and Pisoni, 1998; Nearey, 1990; Norris and McQueen, 2008; interactive-
 659 activation accounts and their offsprings, Magnuson et al., 2020; McClelland and Elman, 1986),
 660 ideal observers describe the posterior probability of a category as dependent both on the prior
 661 probability of the category in the current context, $p(\text{category})$, and the likelihood of the acoustic
 662 input under the hypothesis that it originates from the category, $p(\text{cues}|\text{category})$:

$$p(\text{category}|\text{cues}) = \frac{p(\text{cues}|\text{category}) \times p(\text{category})}{\sum_c p(\text{cues}|\text{category}_c) \times p(\text{category}_c)} \quad (2)$$

663 The category prior, $p(\text{category})$, describes how much the surrounding context favors each category.
 664 For Study 2, the choice of category prior cannot affect the qualitative results since category priors
 665 are independent of the cues and held identical across all normalization accounts.¹⁵ We arbitrarily
 666 assume uniform category priors. Specifically, for ideal observers trained and tested on the long and
 667 short vowels separately, we model categorization as an 11- and 10-alternatives-forced-choice task,
 668 respectively, resulting in $p(\text{category}) = .091$ for the former and $p(\text{category}) = .1$ for the latter.
 669 For ideal observers trained and tested on the entire vowel space, we model categorization as a
 670 21-alternatives-forced-choice task, resulting in $p(\text{category}) = .048$.

671 The likelihood, $p(\text{cues}|\text{category})$, describes the distribution of cues for each category. Here, we
 672 follow previous work and assume multivariate Gaussian distributions to describe the cue likelihood
 673 (e.g., Clayards et al., 2008; Kleinschmidt and Jaeger, 2015; Kronrod et al., 2016; Xie et al., 2021).
 674 That is, we use the model in equation (3), where μ and Σ refer to the category mean and variance-
 675 covariance matrix of the category's multivariate normal distribution. In terms of representational
 676 complexity, the assumption of multivariate Gaussian categories strikes a compromise between
 677 exemplar storage (less representationally parsimonious, Johnson, 1997; Pierrehumbert, 2001) and
 678 cue integration over multiple separate univariate Gaussians (more parsimonious, Toscano and
 679 McMurray, 2010).¹⁶

¹⁵ Specifically, category priors have a constant additive effect on the posterior log-odds of categories.

¹⁶ Human perception is affected by an additional source of uncertainty beyond category variability: perceptual noise (for review, see Feldman et al., 2009). There is now evidence that including such perceptual noise in models provides a better fit against human behavior

$$p(\text{category}|\text{cues}) = \frac{\mathcal{N}(\text{cues}|\mu, \Sigma) \times p(\text{category})}{\sum_c \mathcal{N}(\text{cues}|\mu_c, \Sigma_c) \times p(\text{category}_c)} \quad (3)$$

680 Paralleling Study 1, we trained and tested separate ideal observers for each combination of
 681 normalization account, cues, and training-test fold. Specifically, we use the exact same cross-
 682 validation folds as in Study 1. Each ideal observer was trained on the training portion of the folded
 683 unnormalized and normalized data, and subsequently evaluated on the held-out test fold. This
 684 means that the parameters of each normalization account (e.g., the cue means in C-CuRE) and the
 685 resulting category parameters (the μ_c s and Σ_c s for all categories) were set on the training data, and
 686 not changed for the test data. This reflects the realities of speech perception: although this is often
 687 ignored in evaluations of normalization accounts (e.g., Barreda, 2021; McMurray and Jongman,
 688 2011), listeners do not *a priori* know the cue means, cue variance, etc. of an unfamiliar talker.
 689 Rather, listeners need to incrementally *infer* those statistical properties from the talker's speech
 690 input (for discussion and a model, see Xie et al., 2022). An additional advantage of cross-validation
 691 is that it gives us an estimate of the uncertainty about the model predictions. The performance
 692 of each ideal observer during test is assessed by calculating the ideal observer's predicted posterior
 693 probability of the *intended* category for each test token. This is identical to assuming Luce's choice
 694 rule for categorization (Luce, 1959), as in most models of speech perception and spoken word
 695 recognition.

696 Advantages compared to separability index and similar measures of category variability

697 By using a categorization model to estimate the consequence of normalization for perception,
 698 we avoid the pitfalls of the separability index discussed in Study 1. First, by using the density of
 699 an acoustic input under the multivariate distribution of all cues, the ideal observers we employ
 700 capture the *joint* effect of all cues. This captures that an input can be an improbable instance of a
 701 category based on one of its cue values but a probable instance given the values for the other cues.
 702 In particular, since we assume multivariate likelihoods, the model can account for category-specific
 703 covariances between cues. In the presence of strong correlations between cues, an acoustic input
 704 can be an improbable instance of a category based on the marginal distribution of each cue—i.e., if
 705 each cue is considered separately—but be a highly probable instance of that category if considered
 706 relative to the joint distribution of all cues (and vice versa). This is illustrated by Figure 11A.

707 Second, by normalizing the support for a category by the support for all other categories (the
 708 denominator in Equations (2) and (3)), ideal observers consider the perceptual consequences of
 709 an acoustic input *relative to all possible categories*. This means that a token that is relatively far
 710 away from its category mean does not necessarily result in low recognition accuracy. Rather, low
 711 recognition accuracy is only predicted if the relative position of the acoustic input in the acoustic-
 712 phonetic space makes it more probable that the input originated from another (unintended) category.
 713 This is illustrated in Figure 11B, and parallels human perception: e.g., while a more mid-fronted
 714 [a:] with high F1- and F2-values is atypical, human listeners are more likely to recognize it as
 715 a [a:] compared to a more high-back articulated, but equally atypical, [a:], presumably because
 716 the observed phonetics would be equally likely to occur if the talker intended a [o:]. Measures of

(e.g., Kronrod et al., 2016; Tan and Jaeger, 2022). Since Study 2 assesses the *relative* recognition accuracy of different normalization accounts, it is not immediately obvious how the inclusion of noise could affect our results. To avoid additional researchers degrees of freedom—such as the decision as to which acoustic or perceptual space (Hz, Mel, Bark, etc.) perceptual noise is additive in—we do not model the perceptual consequences of noise. To the best of our knowledge, no previous comparisons of normalization accounts has considered perceptual noise.

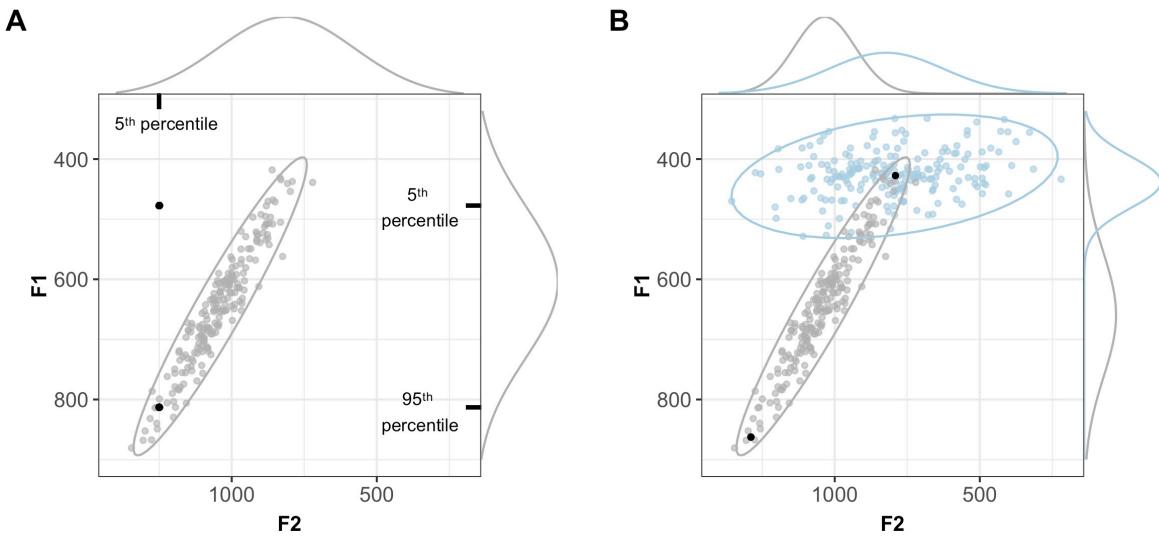


Figure 11. Using a perceptual model to evaluate normalization accounts avoids the pitfalls of separability/variability indices. **Panel A:** an acoustic token can be an improbable instance of a category if each cue is considered separately (the marginal densities along the sides of the plot), but highly probable if considered relative to the joint distribution of cues (the bivariate distribution indicated by the ellipse). **Panel B:** two acoustic tokens that are equally far from the mean of one category can have radically different consequences for perception, depending on where the tokens fall relative to other categories. Under the hypothesis that the two black points are instances of the gray category, they would be attributed the same separability index but radically different probabilities given the joint distribution of cues relative to the other category in the space. Both points are on the 90% highest density interval isoline.

717 between- vs. within-category variability like the separability index in Study 1, however, have no
 718 means of directly capturing this.

719 Finally, when an acoustic input is indeed improbable under its intended category relative to other
 720 unintended categories, the ideal observer model will correctly predict low recognition accuracy.
 721 Unlike the separability index in Study 1, however, the ideal observer will not disproportionately
 722 weight this single observation: it will simply be one of many observations—much like human speech
 723 perception does not collapse simply because one word (or one vowel) has been misunderstood. As
 724 we show next, these differences between the approaches taken in Studies 1 and 2 can, and sometimes
 725 do, cause differences in results.

726 Results aggregating across vowels

727 Figure 12 visualizes the unnormalized and normalized models' predictions for perception of
 728 Central Swedish vowels, under different assumptions about the relevant cues. As in Study 1, this
 729 figure aggregates results across vowels of a given type (long, short, both). We first discuss these
 730 results, and then briefly summarize how different vowels are affected by normalization.

731 The first observation is that—unlike for the separability index in Study 1—the relative
 732 performance of the different normalization accounts within each panel is remarkably constant
 733 across all panels (compare to Figure 10). Regardless of the combination of cues or the vowel
 734 types considered (long, short, both), transformation into a perceptual space does little to improve
 735 recognition accuracy, compared to unnormalized cues. Intrinsic normalization, too, does not
 736 improve recognition accuracy. And, unlike in Study 1, this result holds even when F3 is included as a

cue. This replicates previous work on Dutch (Adank et al., 2004) but conflicts with some evaluations of English (e.g., Syrdal, 1985). Adank et al. (2004) discussed whether the discrepancy in results might be attributed to implementations of the Bark-transformation, or to what Syrdal (1985) describes as language-specificity of the second dimension of Syrdal and Gopal (1986) normalization. The present results would seem to confirm this vulnerability of intrinsic normalizations. Extrinsic normalization, however, tends to substantially improve recognition accuracy (with the exception of Gerstman normalization). Depending on the specific combination of cues and the vowel qualities considered, the best-performing normalization model increases recognition accuracy by at least 46.8% (from 40.8% for unnormalized cues for all vowels when only F1-F2 are considered) to 81.2% (from 75.6% for short vowels when all cues are considered). The benefit of extrinsic normalization models, as well as the lower performance of perceptual transformations, replicates previous findings on other languages (e.g., Adank et al., 2004; Escudero and Bion, 2007; Nearey, 1989 found effects of both intrinsic and extrinsic accounts, but larger effects for extrinsic).

We also see that all models—even for unnormalized cues—perform substantially above chance. When long and short vowels are considered separately, the best ideal observers achieve recognition accuracies of 73.9% for long vowels and 81.2% for short vowels. For reference, in a recent perception experiment we conducted on the eight monophthongs of US English, L1-US English listeners achieved 71.1% accuracy in categorizing isolated hVd words (chance = 12.5%, Persson and Jaeger, 2023). The ideal observers for the Central Swedish vowel system thus achieve performance that is comparable to that of human listeners, at least when cues are normalized.

Looking across columns of Figure 12, short vowels are always recognized with higher accuracy compared to long vowels. This increase in performance cannot be explained by the small increase in the chance baseline alone (10% for the 10 short vowels, compared to 9.1% for the 11 long vowels). It conceptually replicates an initially surprising result of Study 1: while short vowels are more densely clustered in the center of the vowel space, and thus occupy a smaller perceptual space, they also exhibit less variability. Overall, this makes those vowels *easier* to recognize.

When long and short vowels are categorized together, performance of the ideal observers is comparatively poor unless vowel duration is included as a cue. This is expected given that vowel duration is the primary cue to vowel quantity. Of interest, however, is that even the inclusion of only F3 (second row) yields a substantial improvement in recognition accuracy, in line with Johnson and Siersp (2021). Remarkably, once vowel duration is included, the best-performing ideal observer achieves 75% recognition accuracy across the 21 long and short vowels (compared to chance = 4.8%).

Looking across rows, we note that Lobanov normalization clearly performs best when only the first two formants are considered. Indeed, when only F1 and F2 are considered and all vowels are categorized (top right panel of Figure 12), Lobanov performs *significantly* better than any of the other normalization models (the 95% confidence intervals of Lobanov normalization does not include the mean of the second best performing model). However, this advantage of Lobanov normalization decreases and is no longer significant when additional cues are considered.¹⁷

¹⁷ Indeed, when all five cues are considered for the categorization of all 21 short and long vowels, the best centering account performs numerically better (74%) than Lobanov normalization (72.7%). This is, however, an artifact of our decision to only center vowel duration—the primary cue to vowel quantity—for the C-CuRE model (paralleling Study 1). Separate modeling not shown here confirmed that Lobanov normalization achieves the same recognition accuracy as the C-CuRE models when duration is centered and combined with Lobanov-normalized formants (74.9%, 95%-CI: 73.9-75.9%).

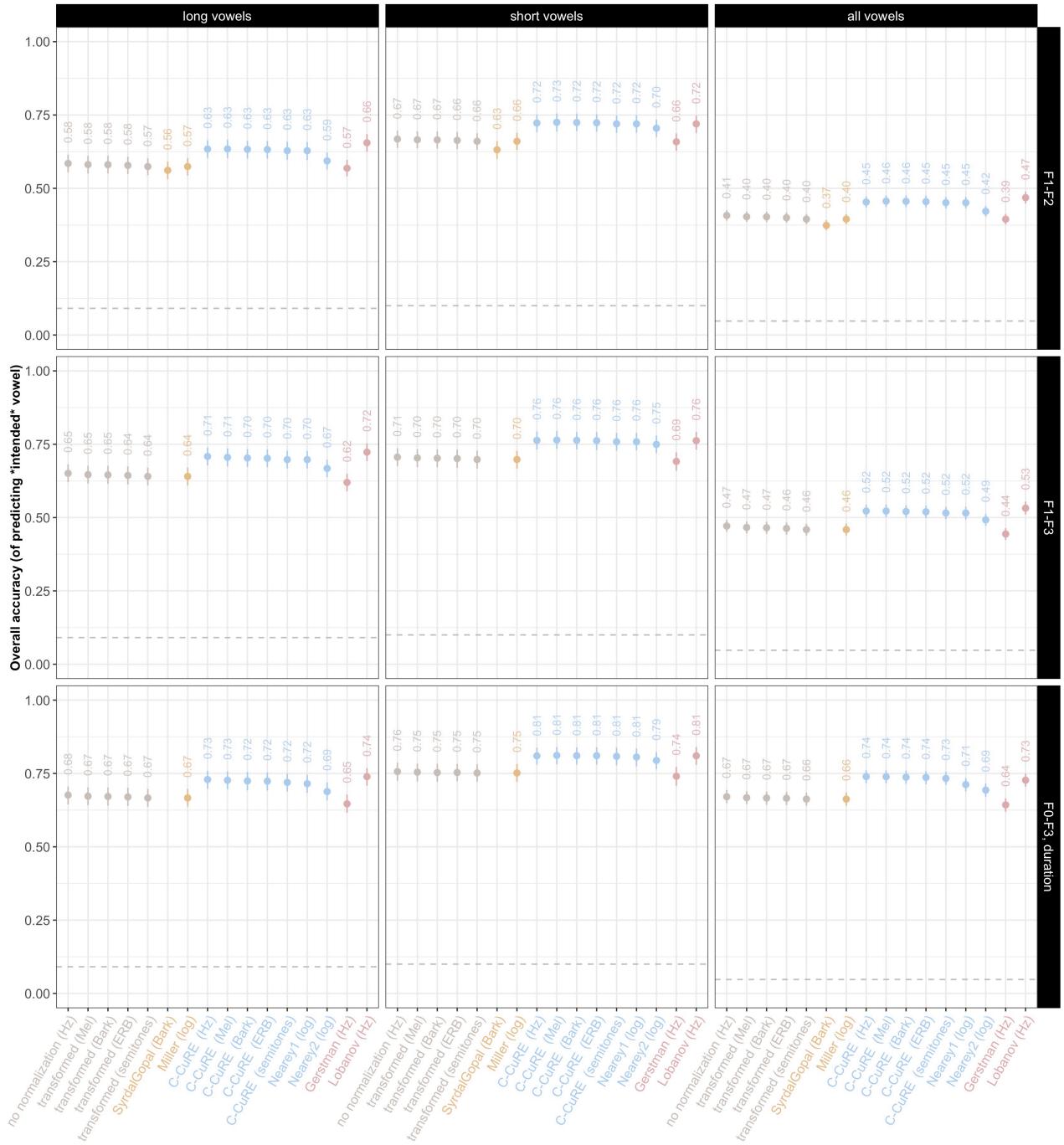


Figure 12. Predicted recognition accuracy of ideal observer under different normalization accounts for long vowels, short vowels, and long and short vowels together (columns), shown for three different combinations of cues (rows). Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. The dashed horizontal line indicates chance (different across columns because of the different number of long and short vowels).

776 Results for specific vowels

777 Finally, the use of a perceptual model lets us assess vowel-specific effects on the predictions of
778 the ideal observers. While this is not the main focus of this paper, it might be of interest to other
779 researchers. In the SI, we therefore plot both the predicted categorization accuracy per vowel in
780 the different evaluations (see SI, Per-vowel categorization accuracy), as well as confusion matrices
781 of the best and the worst performing ideal observers (see SI, Confusion and difference matrices),
782 to further investigate *how* normalization improves recognition accuracy. Here, we briefly mention
783 some of the main findings that emerge from these additional visualizations. For reference, Figure 13
784 visualizes the predicted recognition accuracy for five vowels that illustrate some of the vowel-specific
785 effects across different types of normalization.

786 Unsurprisingly, some vowels are recognized with much higher accuracy than others—at least when
787 uniform category priors are assumed, as we did here. This is a direct consequence of the position
788 of the vowel in the acoustic-phonetic space, relative to neighboring vowels: the more neighboring
789 vowels overlap with each other, the lower the accuracy with which they are recognized. Which
790 vowels will benefit from normalization will thus naturally vary between languages, reflecting the
791 language-specific properties of the vowel space. For instance, [i:] is often described as more easily
792 recognized in previous work on other languages. This contrasts with our findings for Central Swedish:
793 here, [i:] is part of the dense clustering of vowels along the height dimension and so has many close
794 competitors. This highlights that recognition accuracy is due to the position of a vowel *relative* to
795 its competitors (e.g., Peterson and Barney, 1952; Kuhl, 1991; Polka and Bohn, 2003), rather than
796 its *absolute location* in the vowel space (e.g., [i:] being a peripheral vowel).

797 Also of interest is that not all vowels exhibit the benefit of normalization. In general, across
798 evaluations, it seems to be vowels that are already recognized with relatively high accuracy that
799 particularly benefit from normalization, which does not replicate previous studies that have included
800 per-vowel accuracies (e.g., Adank, 2003; Syrdal and Gopal, 1986). In fact, Adank (2003) reports
801 larger improvements in error rates after normalization for vowels with higher error rates prior to
802 normalization. Finally, while there are minor differences across vowels in the relative goodness of
803 different normalizations, the models that perform overall best also perform best on each vowel (in
804 line with Adank, 2003). This further demonstrates the plausibility of these normalization accounts
805 for perception.

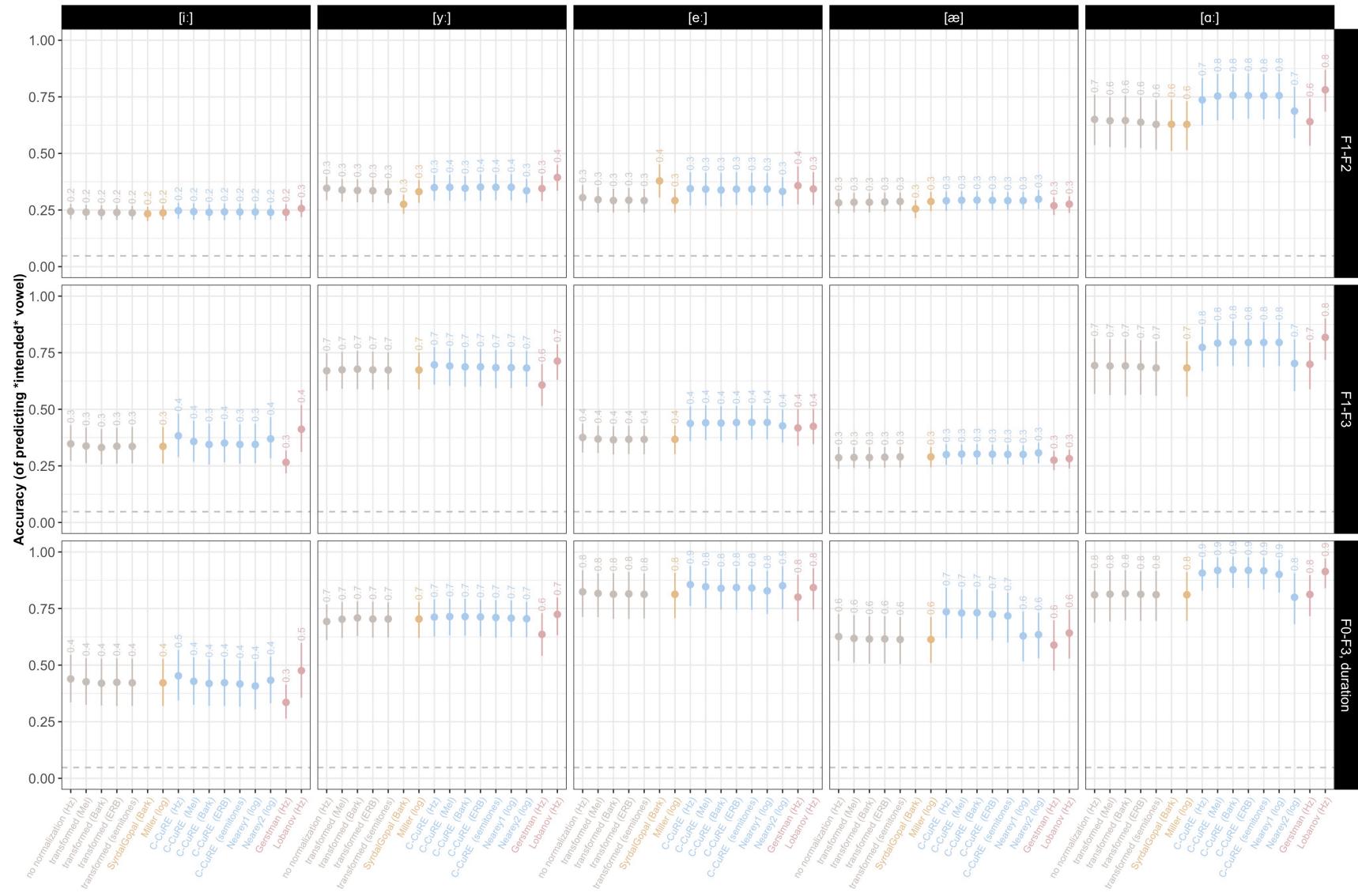


Figure 13. Predicted recognition accuracy of ideal observer under different normalization accounts for five of the 21 vowels. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Labels indicate mean across the five test folds. Chance level is indicated by grey line.

GENERAL DISCUSSION

806 We have compared low-level pre-linguistic normalization accounts against a new phonetically
807 annotated database of Central Swedish vowels. We set out to evaluate how the different accounts
808 affect the category separability (Study 1), and how they differ in predicted consequences for
809 perception (Study 2). As we have already discussed the methodological shortcomings of separability
810 and variability indices like that employed in Study 1, we focus here on the discussion of Study 2.

811 Previous work found that the types of normalization accounts that performed well on other
812 languages did not seem to perform well on Swedish vowel data (Disner, 1980). However, as pointed
813 out by Disner, the Swedish data differed from the data for other languages in that study, and the
814 majority of studies on other languages. Here, we followed the majority of previous work on vowel
815 productions and analyzed productions of hVd recordings. We find—both in Study 2 and in Study
816 1 (as long as additional cues are not included)—that the same accounts found in previous work to
817 perform well on other languages also perform well for the dense vowel space of Swedish. Specifically,
818 Lobanov and centering approaches—incl. Nearey normalization and C-CuRE normalization—were
819 the top-performing accounts, replicating the pattern found in previous studies on other languages
820 (e.g., Adank et al., 2004; Carpenter and Govindarajan, 1993; Escudero and Bion, 2007; Syrdal,
821 1985). This result suggests that the (somewhat) diverging results for Swedish in Disner (1980)'s
822 study, were not caused by properties inherent to Swedish, but more likely were an artifact of
823 the dataset employed by Disner. It also suggests that languages with dense vowel spaces do not
824 necessarily require more complex normalization mechanisms.

825 Additionally, we find that the best-performing centering accounts (C-CuRE) often achieve
826 performance that is statistically indistinguishable from the best-performing standardization
827 accounts (Lobanov). This is the case, in particular, when all five cues were considered and all
828 21 vowels were included in the categorization (see footnote 17). Together with similar findings from
829 research on consonants and supra-segmental categories (e.g., Apfelbaum et al., 2014; Crinnion
830 et al., 2020; Kleinschmidt, 2019; Kulikov, 2022; McMurray and Jongman, 2011, 2016; Toscano
831 and McMurray, 2015; Xie et al., 2021, 2022), this suggests that simple centering operations
832 might be sufficient to maximize the benefits achievable by normalization. In the remainder, we
833 first summarize some methodological considerations based on the present study, and then discuss
834 limitations of our work, and how they can be addressed in future work.

835 Methodological considerations

836 We considered two methods of evaluations. In Study 1, we compared the reduction of within-
837 category variability relative to between-category variability by calculating a separability index (i.e.,
838 an F-statistics). This methodology facilitated comparison to previous studies, that has most often
839 assessed the efficacy of normalization accounts by means of increase in between-talker category
840 overlap or category separability, or decrease in category variability (Disner, 1980; Fabricius et
841 al., 2009; Flynn and Foulkes, 2011; Hindle, 1978; Labov, 2010). However, given the limitations
842 inherent in this index, in the second study, we used Bayesian ideal observers to investigate
843 how normalization accounts differ in predictions for perception. Previous studies have employed
844 linear discriminant analysis (e.g., Adank et al., 2004; Syrdal, 1985), fuzzy ARTMAP or K-nearest
845 neighbour classification (e.g., Carpenter and Govindarajan, 1993), or constraint ranking grammars
846 (e.g., Escudero and Bion, 2007). While these different approaches will generally return similar
847 results, we see two advantages with ideal observers. First, ideal observers remove researchers'
848 degrees of freedom in the evaluation of normalization accounts (see Tan et al., 2021): the statistics

849 of the cue distributions in the training data fully determine the ideal observers' predictions.
 850 Second, representing the cue likelihood for each category as multivariate Gaussian distributions
 851 (e.g., Kronrod et al., 2016; Norris and McQueen, 2008; Xie et al., 2021), strikes a middle ground
 852 between less parsimonious models such as exemplar models (e.g., Johnson, 1997; Pierrehumbert,
 853 2001), and more parsimonious models, such as models of cue integration over multiple separate
 854 univariate Gaussians (e.g., Toscano and McMurray, 2010).¹⁸ However, like any other model of
 855 speech perception, the approach adopted here comes with a set of simplifying assumptions. All
 856 of the assumptions in the model, e.g., uniform priors, models of linguistic representations, the
 857 normalization accounts selected, can, however, be revisited and altered by the reader, as this paper
 858 is written in R markdown and all data and code is provided on OSF.

859 As a final note on methodology, we find the 5-fold cross-validation approach used for normalization
 860 of the acoustic data adopted here, advantageous for several reasons, the two most important
 861 being that: (1) it allows us to avoid over-fitting to the sample, while also (2) providing a more
 862 realistic reflection of how parameters used for normalization are incrementally inferred from the
 863 talker's speech input. Even though many of the commonly adopted normalization accounts involve
 864 parameters that are set based on the data, previous studies have rarely considered how these
 865 parameters might be affected by the specific dataset used for normalization (for exceptions, see e.g.,
 866 Barreda and Nearey, 2018). In the present study, we have seen that accounting for researchers'
 867 uncertainty about the effects of normalization, highlights that many of the normalization accounts
 868 exhibit statistically indistinguishable performance—at least under the approach taken here and in
 869 the majority of previous work. This means that future studies should aim to increase statistical
 870 power, in order to determine the normalization mechanism that best describes human behavior.
 871 This will require even larger datasets, and/or by targeted sampling of vowel tokens for which
 872 predictions of different normalization accounts are maximally contrasted (see e.g., Barreda, 2021;
 873 and see Xie et al., 2022 for more general discussion of how to increase the statistical power to
 874 determine what mechanism underlie adaptive speech perception). Next, we close by discussing
 875 additional limitations of the present work and future directions.

876 Limitations and future directions

877 Four limitations of the present study, three of which are shared with most previous work, deserve
 878 discussion. First, the present study compared normalization accounts against speech from only
 879 female talkers of one regional variety of Central Swedish (Stockholm Swedish). In contrast, many
 880 previous studies included data from talkers of different genders (e.g., Barreda, 2021; Clopper, 2009;
 881 Cole et al., 2010; McMurray et al., 2011; McMurray and Jongman, 2016), and sometimes from
 882 talkers of different ages (e.g., Barreda and Nearey, 2018; Carpenter and Govindarajan, 1993; Flynn
 883 and Foulkes, 2011; Hindle, 1978; Johnson and Sjerps, 2021; Kohn and Farrington, 2012; Syrdal,
 884 1985) and/or language backgrounds (e.g., Adank et al., 2004; Disner, 1980; Escudero and Bion,
 885 2007; Fabricius et al., 2009; Labov, 2010; Richter et al., 2017). Given that age, gender, etc. tends
 886 to affect formants (and other cues) beyond talker-variability, it is likely that the inclusion of more
 887 diverse talkers would increase the lack of invariance problem. For example, we would expect the
 888 ideal observers over unnormalized cues (Study 2) to achieve lower recognition performance if vowel
 889 productions from male talkers would be included in the data. In short, the models in Study 2

¹⁸ Parsimony here refers both to the number of degrees of freedom these models afford to the researcher *and* to the amount of information that listeners are assumed to store (for discussion, see Xie et al., 2022).

890 likely over-estimates the recognition accuracy that can be achieved for unnormalized cues if a more
891 diverse range of talkers is considered.

892 What does that mean for the relative effect of normalization? To the extent that normalization
893 successfully overcomes inter-talker variability that is caused by gender, age, and other social or
894 physiological factors, we expect that the benefit of normalization accounts should show more clearly,
895 relative to unnormalized cues. In this sense, the present study might *under-estimate* the relative
896 benefits of normalization. Whether the *relative* performance of normalization accounts—i.e., the
897 finding of primary interest to us—would differ if a more diverse range of talkers was considered is
898 unclear. To the extent that vowel-specific accounts were originally developed specifically to eliminate
899 physiological differences that are correlated with gender (as reviewed in, e.g., Johnson and Sjerps,
900 2021), it is theoretically possible that the high performance of general normalization accounts (e.g.,
901 C-CuRE, McMurray and Jongman, 2011) might not replicate when talkers of different genders are
902 included. Future releases of the SwehVd database will contain data from male talkers, which will
903 allow us or other researchers to revisit these questions.

904 Second, the present study aggregated acoustic-phonetic measurements taken at different points of
905 the vowel (at 35%, 50%, and 65% into the vowel) into a single formant measurement. This follows
906 previous comparisons of normalization accounts but is a simplifying assumption that should be
907 revisited in future work. Formant dynamics carry important information for category distinctions
908 (e.g., Assmann and Katz, 2005; Hillenbrand and Nearey, 1999; Nearey and Assmann, 1986), and
909 are hypothesized to be of particular importance for some vowel distinctions in other varieties of
910 Central Swedish (e.g., Kuronen, 2000). Prior to other consideration, this means that Study 2 likely
911 under-estimates the recognition accuracy that could be achieved even from unnormalized cues
912 alone. It is an open question whether the findings of primary interest—the relative performance
913 of different normalization accounts—would be affected if formant dynamics were considered. Some
914 normalization accounts, for example, consider normalization of such formant dynamics to take place
915 *after* basic formant normalization (but before the mapping of cues to category representations, S.
916 Barreda, personal communication, 01/06/2023). Future work could employ SwehVd to compare
917 ideal observers or other classification models while taking into consideration formant measurements
918 throughout the vowel.

919 Third, we only considered competing *normalization* accounts. This, too, follows previous research
920 on normalization but is potentially problematic. As mentioned in the introduction, it is now believed
921 that at least three different mechanisms contribute to adaptive speech perception, including not
922 only normalization but also changes in category representations and decision-making (for review,
923 see Xie et al., 2022). This has consequences for research on normalization. For example, Xie et
924 al. (2021) compared normalization accounts against the production of prosodic phrasing in L1-
925 US English, while also considering alternative hypotheses about listeners' ability to adapt category
926 representations. Xie and colleagues found that the effectiveness of cue normalization is substantially
927 reduced if listeners can learn and maintain talker- or group-specific category representations (as
928 assumed in some influential theories of speech perception, exemplar models, e.g., Johnson, 1997;
929 Pierrehumbert, 2001; Bayesian ideal adaptors, Kleinschmidt and Jaeger, 2015). Xie and colleagues
930 only considered two general types of normalization, and their focus was on the interpretation of
931 prosodic signals. But their results call for caution in interpreting studies like the present that do
932 not consider the possibility of talker-specific representations—an assumption shared with basically
933 all previous work on vowel normalization.

934 Similarly, as mentioned in the introduction, we limited our evaluation to a single level of
935 normalization (and combinations of perceptual transformations and a single level of normalization).
936 Some proposals, however, assume multiple separate normalization steps. For example, some
937 accounts hold that evolutionarily early mechanisms first transform spectral percepts into a phonetic
938 space (e.g., uniform scaling accounts, Barreda, 2020; Nearey, 1983), on which additional subsequent
939 normalization might operate. There is also evidence that speech perception combines aspects of
940 intrinsic and extrinsic normalization (Johnson and Sjerps, 2021) review relevant evidence from
941 brain imaging; early behavioral evidence is found in Nearey, 1989). The present study—like most
942 existing evaluations—did not consider these possibilities (for exceptions, see e.g., Barreda, 2021;
943 Nearey and Assmann, 2007).

944 Fourth and finally, we followed the majority of previous work and evaluated normalization
945 accounts against *production* data. This is potentially problematic, especially when measures like
946 category separability or reduced cross-talker variability in category means are used to evaluated
947 normalization accounts (as in our Study 1 and many previous studies). These evaluations essentially
948 assume that the goal of speech perception is to make the perceptual realizations of the same category
949 by different talkers as similar as possible in the normalized space (for an in-depth critique, see
950 Barreda, 2021). However, the goal of speech perception is presumably to reliably infer the category
951 intended by the talker,¹⁹ and this aim does not necessarily entail perfect removal of cross-talker
952 variability (as evidenced, for example, by the different findings of Studies 1 and 2).

953 To some extent, Study 2 addresses this potential issue by evaluating normalization accounts in
954 terms of how well they predict the vowel category intended by the talker. However, if the goal is to
955 explain human perception, the most informative evaluations of normalization accounts are arguably
956 those that compare their predictions against *listeners'* behavior (for examples, see Barreda, 2020,
957 2021; McMurray and Jongman, 2016; Nearey, 1989; Richter et al., 2017; Xie et al., 2021). In short,
958 approaches like that employed in Study 2 take an important step away from the most misleading
959 evaluation of normalization accounts in terms of reduced category variability/increased category
960 separability. Ultimately, however, normalization accounts should be evaluated in terms of how well
961 they predict listeners' perception, not talker's intention.

DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

962 The authors declare that the research was conducted in the absence of any commercial or financial
963 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

964 AP proposed project idea, and both authors jointly developed the conceptual approach. AP designed
965 SwehVd materials, recorded and annotated vowel productions, and coded cue extraction. AP coded
966 data analyses and visualization with guidance from TFJ. Both authors contributed to the writing
967 of the manuscript.

ACKNOWLEDGMENTS

968 We are grateful to OMITTED FOR REVIEW.

¹⁹ Or somewhat more precisely, cooperative listeners aim to understand the meaning intended by the talker, and this inference is generally believed to benefit from the correct recognition of phonological categories, such as phonemes, syllables, or word forms. For discussion, see also Hume et al. (2016).

REFERENCES

- 969 Adank, P. (2003). *Vowel normalization: A perceptual-acoustic study of Dutch vowels*.
970 Wageningen: Ponsen & Looijen.
- 971 Adank, P., Smits, R., and van Hout, R. (2004). A comparison of vowel normalization
972 procedures for language variation research. *The Journal of the Acoustical Society of*
973 *America* 116, 3099–3107. doi:[10.1121/1.1795335](https://doi.org/10.1121/1.1795335).
- 974 Apfelbaum, K., Bullock-Rest, N., Rhone, A., Jongman, A., and McMurray, B. (2014).
975 Contingent categorization in speech perception. *Language, Cognition and Neuroscience*
976 29, 1070. doi:[10.1080/01690965.2013.824995](https://doi.org/10.1080/01690965.2013.824995).
- 977 Assmann, P. F., and Katz, G. S. (2005). Synthesis fidelity and time-varying spectral change
978 in vowels. *Journal of the Acoustical Society of America* 117, 886–895.
- 979 Assmann, P. F., Nearey, T. M., and Hogan, J. (1982). Vowel identification: Orthographic,
980 perceptual, and acoustic aspects. *Journal of the Acoustical Society of America* 71, 975–
981 989.
- 982 Audacity, T. (2021). Audacity(R): Free Audio Editor and Recorder [Computer Application].
- 983 Barreda, S. (2021). Perceptual validation of vowel normalization methods for variationist
984 research. *Language Variation and Change* 33, 27–53. doi:[10.1017/S0954394521000016](https://doi.org/10.1017/S0954394521000016).
- 985 Barreda, S. (2020). Vowel normalization as perceptual constancy. *Language* 96, 224–254.
986 doi:[10.1353/lan.2020.0018](https://doi.org/10.1353/lan.2020.0018).
- 987 Barreda, S., and Nearey, T. M. (2018). A regression approach to vowel normalization for
988 missing and unbalanced data. *The Journal of the Acoustical Society of America* 144,
989 500–520. doi:[10.1121/1.5047742](https://doi.org/10.1121/1.5047742).
- 990 Barreda, S., and Nearey, T. M. (2012). The direct and indirect roles of fundamental
991 frequency in vowel perception. *The Journal of the Acoustical Society of America* 131,
992 466–477. doi:[10.1121/1.3662068](https://doi.org/10.1121/1.3662068).
- 993 Behne, D. M., Czigler, P. E., and Sullivan, K. P. H. (1997). Swedish Quantity and Quality: A
994 Traditional Issue Revisited. *Reports from the Department of Phonetics, Umeå University*
995 4, 81–83.
- 996 Bladon, A., Henton, C. G., and Pickering, J. B. (1984). Towards an auditory theory of
997 speaker normalization. *Language and Communication* 4, 59–69.
- 998 Boersma, P., and Weenink, D. (2022). Praat: Doing phonetics by computer [Computer
999 program].
- 1000 Bruce, G. (2009). “Components of a prosodic typology of Swedish intonation,” in
1001 *Components of a prosodic typology of Swedish intonation* (De Gruyter Mouton), 113–146.
1002 doi:[10.1515/9783110207569.113](https://doi.org/10.1515/9783110207569.113).
- 1003 Bruce, G., Elert, C.-C., Engstrand, O., and Wretling, P. (1999). Phonetics and phonology of
1004 the Swedish dialects - a project presentation and a database demonstrator. *Proceedings of*
1005 *the 14th International Congress of Phonetic Sciences, University of California.*, 321–324.
- 1006 Carpenter, G. A., and Govindarajan, K. K. (1993). “Neural Network and Nearest Neighbor
1007 Comparison of Speaker Normalization Methods for Vowel Recognition,” in *ICANN '93*,
1008 eds. S. Gielen and B. Kappen (London: Springer London), 412–415. doi:[10.1007/978-1-4471-2063-6_98](https://doi.org/10.1007/978-1-4471-2063-6_98).
- 1009 Chesworth, J., Coté, K., Shaw, C., Williams, S., and Hodge, W. (2003). Effect of phonetic
1010 context on women’s vowel area. *Canadian Acoustics/Acoustique Canadienne* 31, 20–21.

- 1012 Clayards, M., Tanenhaus, M. K., Aslin, R. N., and Jacobs, R. A. (2008). Perception
1013 of speech reflects optimal use of probabilistic speech cues. *Cognition* 108, 804–809.
1014 doi:[10.1016/j.cognition.2008.04.004](https://doi.org/10.1016/j.cognition.2008.04.004).
- 1015 Clopper, C. G. (2009). Computational Methods for Normalizing Acoustic Vowel Data
1016 for Talker Differences: Computational Methods for Normalizing Acoustic Vowel Data.
1017 *Language and Linguistics Compass* 3, 1430–1442. doi:[10.1111/j.1749-818X.2009.00165.x](https://doi.org/10.1111/j.1749-818X.2009.00165.x).
- 1018 Cole, J., Linebaugh, G., Munson, C., and McMurray, B. (2010). Unmasking the acoustic
1019 effects of vowel-to-vowel coarticulation: A statistical modeling approach. *Journal of*
1020 *Phonetics* 38, 167–184. doi:[10.1016/j.wocn.2009.08.004](https://doi.org/10.1016/j.wocn.2009.08.004).
- 1021 Crinnion, A. M., Malmskog, B., and Toscano, J. C. (2020). A graph-theoretic approach
1022 to identifying acoustic cues for speech sound categorization. *Psychonomic Bulletin &*
1023 *Review* 27, 1104–1125. doi:[10.3758/s13423-020-01748-1](https://doi.org/10.3758/s13423-020-01748-1).
- 1024 Disner, S. F. (1980). Evaluation of vowel normalization procedures. *The Journal of the*
1025 *Acoustical Society of America* 67, 253–261. doi:[10.1121/1.383734](https://doi.org/10.1121/1.383734).
- 1026 Eklund, I., and Traunmüller, H. (1997). Comparative Study of Male and Female
1027 Whispered and Phonated Versions of the Long Vowels of Swedish. *Phonetica* 54, 1–21.
1028 doi:[10.1159/000262207](https://doi.org/10.1159/000262207).
- 1029 Elert, C.-C. (1994). “Indelning och gränser inom området för den talade svenska: En
1030 aktuell dialektografi,” in *Kulturgränser - myt eller verklighet?* Diabas. (Institutionen för
1031 nordiska språk vid Umeå Universitet), 215–228.
- 1032 Elert, C.-C. (1981). *Ljud och ord i svenska*. Umeå: Universitetet i Umeå, Almqvist &
1033 Wiksell international.
- 1034 Engstrand, O., Bruce, G., Elert, C.-C., Eriksson, A., and Strangert, E. (2001).
1035 Databearbetning i SweDia 2000: Segmentering, transkription och taggning. Version 2.2.
- 1036 Escudero, P., and Bion, R. A. H. (2007). Modeling vowel normalization and sound perception
1037 as sequential processes. *ICPhS XVI*, 1413–1416.
- 1038 Fabricius, A., Watt, D., and Johnson, D. E. (2009). A comparison of three speaker-
1039 intrinsic vowel formant frequency normalization algorithms for sociophonetics. *Language*
1040 *Variation and Change* 21, 413–435. doi:[10.1017/S0954394509990160](https://doi.org/10.1017/S0954394509990160).
- 1041 Fant, G. (1983). Feature analysis of Swedish vowels - a revisit. *STL-QPSR* 24, 001–019.
- 1042 Fant, G. (1975). Non-uniform vowel normalization. *STL-QPSR* 16, 001–019.
- 1043 Fant, G. (1971). “Notes on the Swedish Vowel System,” in *Form and substance: Phonetic and*
1044 *linguistic papers.*, eds. L. Hammerich, R. Jakobson, E. Zwirner, and E. Fischer-Jørgensen
1045 (Odense: Andelsbogtrykkeriet).
- 1046 Fant, G. (2001). Swedish vowels and a new three-parameter model. *TMH-QPSR* 42, 043–049.
- 1047 Fant, G., Henningsson, G., and Stålhammar, U. (1969). Formant frequencies of Swedish
1048 vowels. *STL-QPSR* 10, 026–031.
- 1049 Feldman, N. H., Griffiths, T. L., and Morgan, J. L. (2009). The influence of categories
1050 on perception: Explaining the perceptual magnet effect as optimal statistical inference.
1051 *Psychological Review* 116, 752–782.
- 1052 Flynn, N., and Foulkes, P. (2011). Comparing vowel formant normalization methods.
1053 *Proceedings of ICPHS XVII*, 683–686.
- 1054 Fox, R. A., Flege, J. E., and Munro, M. J. (1995). The perception of English and Spanish
1055 vowels by native English and Spanish listeners: A multidimensional scaling analysis. *The*
1056 *Journal of the Acoustical Society of America* 97, 2540–2551. doi:[10.1121/1.411974](https://doi.org/10.1121/1.411974).

- 1057 Fujimura, O. (1967). On the Second Spectral Peak of Front Vowels: A Perceptual Study
1058 of the Role of the Second and Third Formants. *Language and Speech* 10, 181–193.
1059 doi:[10.1177/002383096701000304](https://doi.org/10.1177/002383096701000304).
- 1060 Gerstman, L. (1968). Classification of self-normalized vowels. *IEEE Transactions on Audio*
1061 and *Electroacoustics* 16, 78–80. doi:[10.1109/TAU.1968.1161953](https://doi.org/10.1109/TAU.1968.1161953).
- 1062 Hadding-Koch, K., and Abramson, A. S. (1964). Duration Versus Spectrum in
1063 Swedish Vowels: Some Perceptual Experiments2. *Studia Linguistica* 18, 94–107.
1064 doi:[10.1111/j.1467-9582.1964.tb00451.x](https://doi.org/10.1111/j.1467-9582.1964.tb00451.x).
- 1065 Hillenbrand, J. M., Getty, L. A., Clark, M. J., and Wheeler, K. (1995). Acoustic
1066 characteristics of american english vowels. *Journal of the Acoustical Society of America*
1067 97, 3099–3111.
- 1068 Hillenbrand, J. M., and Nearey, T. M. (1999). Identification of resynthesized /hVd/
1069 utterances: Effects of formant contour. *The Journal of the Acoustical Society of America*
1070 105, 3509–3523. doi:[10.1121/1.424676](https://doi.org/10.1121/1.424676).
- 1071 Hindle, D. (1978). “Approaches to Vowel Normalization in the Study of Natural Speech,” in
1072 *Linguistic variation: Models and methods*, ed. D. Sankoff (New York: Academic Press),
1073 161–171.
- 1074 Hume, E. V., Jaeger, T. F., and Hall, K. (2016). The Message Shapes Phonology.
- 1075 Johnson, K. (2006). Resonance in an exemplar-based lexicon: The emergence of social
1076 identity and phonology. *Journal of Phonetics* 34, 485–499.
- 1077 Johnson, K. (2005). Speaker normalization in speech perception. *The Handbook of Speech*
1078 *Perception*. doi:[10.1002/9780470757024.ch15](https://doi.org/10.1002/9780470757024.ch15).
- 1079 Johnson, K. (1997). “Speech perception without speaker normalization,” in *Talker variability*
1080 in *speech processing*, eds. K. Johnson and W. Mullennix (San Diego: CA: Academic
1081 Press), 146–165.
- 1082 Johnson, K., and Sjerps, M. J. (2021). “Speaker normalization in speech perception,”
1083 in *The handbook of speech perception* (John Wiley & Sons, Ltd), 145–176.
1084 doi:[10.1002/9781119184096.ch6](https://doi.org/10.1002/9781119184096.ch6).
- 1085 Joos, M. (1948). Acoustic Phonetics. *Language* 24, 5–136. doi:[10.2307/522229](https://doi.org/10.2307/522229).
- 1086 Kleinschmidt, D. F. (2019). Structure in talker variability: How much is there
1087 and how much can it help? *Language, Cognition and Neuroscience* 34, 43–68.
1088 doi:[10.1080/23273798.2018.1500698](https://doi.org/10.1080/23273798.2018.1500698).
- 1089 Kleinschmidt, D. F., and Jaeger, T. F. (2015). Robust speech perception: Recognize the
1090 familiar, generalize to the similar, and adapt to the novel. *Psychological Review* 122,
1091 148–203. doi:[10.1037/a0038695](https://doi.org/10.1037/a0038695).
- 1092 Kleinschmidt, D. F., Weatherholtz, K., and Jaeger, T. F. (2018). Sociolinguistic
1093 perception as inference under uncertainty. *Topics in Cognitive Science* 10, 818–834.
1094 doi:[10.1111/tops.12331](https://doi.org/10.1111/tops.12331).
- 1095 Kohn, M. E., and Farrington, C. (2012). Evaluating acoustic speaker normalization
1096 algorithms: Evidence from longitudinal child data. *The Journal of the Acoustical Society*
1097 of *America* 131, 2237–2248. doi:[10.1121/1.3682061](https://doi.org/10.1121/1.3682061).
- 1098 Kraljic, T., and Samuel, A. G. (2007). Perceptual adjustments to multiple speakers. *Journal*
1099 of *Memory and Language* 56, 1–15.

- 1100 Kronrod, Y., Coppess, E., and Feldman, N. H. (2016). A unified model of categorical
1101 effects in consonant and vowel perception. *Psychological Bulletin and Review*, 1681–1712.
1102 doi:10.3758/s13423-016-1049-y.
- 1103 Kuhl, P. K. (1991). Human adults and human infants show a "perceptual magnet effect"
1104 for the prototypes of speech categories, monkeys do not. *Perception & psychophysics* 50,
1105 93–107. doi:10.3758/BF03212211.
- 1106 Kulikov, V. (2022). Voice and Emphasis in Arabic Coronal Stops: Evidence for Phonological
1107 Compensation. *Language and Speech* 65, 73–104. doi:10.1177/0023830920986821.
- 1108 Kuronen, M. (2000). *Vokaluttalets akustik i sverigesvenska, finlandssvenska och finska*.
1109 Jyväskylä: University of Jyväskylä.
- 1110 Labov, W. (2010). *Principles of linguistic change. 2: Social factors*. repr. Chichester: Wiley-
1111 Blackwell.
- 1112 Leinonen, T. (2010). An Acoustic Analysis of Vowel Pronunciation in Swedish Dialects.
1113 *Groningen Dissertations in Linguistics* 83.
- 1114 Liberman, A. M., Cooper, F. S., Shankweiler, D. P., and Studdert-Kennedy, M. (1967).
1115 Perception of the speech code. *Psychological review* 74, 431–461.
- 1116 Linell, P. (1979). *Psychological reality in phonology: A theoretical study*. Cambridge:
1117 Cambridge University Press.
- 1118 Linell, P. (1978). "Vowel length and consonant length in Swedish word level phonology,"
1119 in *Nordic prosody: Papers from a symposium* Travaux de l'Institut de Linguistique de
1120 Lund., eds. E. Gårding, G. Bruce, and R. Bannert, 123–136.
- 1121 Lobanov, B. M. (1971). Classification of Russian vowels spoken by different speakers. *The
1122 Journal of the Acoustical Society of America* 49, 606–608. doi:10.1121/1.1912396.
- 1123 Luce, P. A., and Pisoni, D. B. (1998). Recognizing spoken words: The neighborhood
1124 activation model. *Ear and Hearing* 19, 1–36. doi:10.1097/00003446-199802000-00001.
- 1125 Luce, R. D. (1959). *Individual choice behavior*. Oxford: John Wiley.
- 1126 Magnuson, J. S., You, H., Luthra, S., Li, M., Nam, H., Escabí, M., Brown, K.,
1127 Alloppenna, P. D., Theodore, R., Monto, N., et al. (2020). EARSHOT: A minimal neural
1128 network model of incremental human speech recognition. *Cognitive Science* 44, 1–17.
1129 doi:10.1111/cogs.12823.
- 1130 Malinasky, M., Shafiro, V., Moberly, A. C., and Vasil, K. J. (2020). Perception of vowels and
1131 consonants in cochlear implant users. *The Journal of the Acoustical Society of America*
1132 148, 2711–2711. doi:10.1121/1.5147511.
- 1133 McAllister, R., Lubker, J., and Carlson, J. (1974). An EMG study of some characteristics
1134 of the Swedish rounded vowels. *Journal of Phonetics* 2, 267–278. doi:10.1016/S0095-
1135 4470(19)31297-5.
- 1136 McClelland, J. L., and Elman, J. L. (1986). The TRACE model of speech perception.
1137 *Cognitive Psychology* 18, 1–86.
- 1138 McMurray, B., Cole, J., and Munson, C. (2011). *Features as an emergent product of
1139 computing perceptual cues relative to expectations*., eds. G. N. Clements and R. Ridouane
1140 John Benjamins Publishing Company doi:10.1075/lfab.6.08mcm.
- 1141 McMurray, B., and Jongman, A. (2016). What comes after /f/?: Prediction in
1142 speech derives from data-explanatory processes. *Psychological Science* 27, 43–52.
1143 doi:10.1177/0956797615609578.

- 1144 McMurray, B., and Jongman, A. (2011). What information is necessary for speech
1145 categorization?: Harnessing variability in the speech signal by integrating cues computed
1146 relative to expectations. *Psychological Review* 118, 219–246. doi:10.1037/a0022325.What.
- 1147 Miller, J. D. (1989). Auditory-perceptual interpretation of the vowel. *The Journal of
1148 Acoustical Society of America* 85, 22.
- 1149 Nearey, T. M. (1978). *Phonetic Feature Systems for Vowels*. Indiana.
- 1150 Nearey, T. M. (1989). Static, dynamic, and relational properties in vowel perception. *The
1151 Journal of the Acoustical Society of America* 85, 2088–2113. doi:10.1121/1.397861.
- 1152 Nearey, T. M. (1990). The segment as a unit of speech perception. *Journal of Phonetics* 18,
1153 347–373. doi:10.1016/S0095-4470(19)30379-1.
- 1154 Nearey, T. M. (1983). Vowel-space normalization procedures and phone-preserving
1155 transformations of synthetic vowels. *The Journal of the Acoustical Society of America*
1156 74, S17–S17. doi:10.1121/1.2020835.
- 1157 Nearey, T. M., and Assmann, P. F. (1986). Modeling the role of inherent spectral change
1158 in vowel identification. *The Journal of the Acoustical Society of America* 80, 1297–1308.
1159 doi:10.1121/1.394433.
- 1160 Nearey, T. M., and Assmann, P. F. (2007). “Probabilistic ‘sliding template’ models for
1161 indirect vowel normalization.” in *Experimental approaches to phonology*, eds. M.-J. Solé,
1162 P. S. Beddor, and M. Ohala (Oxford: Oxford University Press).
- 1163 Nearey, T. M., and Hogan, J. (1986). “Phonological contrast in experimental phonetics:
1164 Relating distributions of measurements production data to perceptual categorization
1165 curves,” in *Experimental Phonology*, eds. J. J. Ohala and J. Jaeger (New York: Academic
1166 Press), 141–161.
- 1167 Nordstrand, M., Svanfeldt, G., Granström, B., and House, D. (2004). Measurements
1168 of articulatory variation in expressive speech for a set of Swedish vowels. *Speech
1169 Communication* 44, 187–196. doi:10.1016/j.specom.2004.09.003.
- 1170 Nordström, P. E., and Lindblom, B. (1975). A normalization procedure for vowel formant
1171 data. *Proceedings of ICPHS VIII*.
- 1172 Norris, D., and McQueen, J. M. (2008). Shortlist B: A Bayesian model of continuous speech
1173 recognition. *Psychological review* 115, 357–95. doi:10.1037/0033-295X.115.2.357.
- 1174 Pardo, J. S., and Remez, R. E. (2006). “Chapter 7 - The Perception of Speech,” in *Handbook
1175 of Psycholinguistics (Second Edition)*, eds. M. J. Traxler and M. A. Gernsbacher (London:
1176 Academic Press), 201–248. doi:10.1016/B978-012369374-7/50008-0.
- 1177 Peirce, J., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., Kastman,
1178 E., and Lindeløv, J. K. (2019). PsychoPy2: Experiments in behavior made easy. *Behavior
1179 Research Methods* 51, 195–203. doi:10.3758/s13428-018-01193-y.
- 1180 Persson, A. (2023). Acoustic-perceptual cues to vowel identity in Stockholm Swedish.
1181 *Manuscript, Stockholm University*.
- 1182 Persson, A., and Jaeger, T. F. (2023). The effect of pre-linguistic normalization in vowel
1183 perception. *Manuscript, Stockholm University*.
- 1184 Peterson, G. E. (1961). Parameters of Vowel Quality. *Journal of Speech and Hearing Research*
1185 4, 10–29. doi:10.1044/jshr.0401.10.
- 1186 Peterson, G. E., and Barney, H. L. (1952). Control methods used in a study of the vowels.
1187 *Journal of the Acoustical Society of America* 24, 175–184.

- 1188 Pierrehumbert, J. B. (2001). "Exemplar dynamics: Word frequency, lenition and contrast,"
1189 in *Frequency and the Emergence of Linguistic Structure*, eds. J. Bybee and P. Hopper
1190 (John Benjamins), 137–157.
- 1191 Polka, L., and Bohn, O.-S. (2003). Asymmetries in vowel perception. *Speech Communication*
1192 41, 221–231.
- 1193 R Core Team (2021). *R: A language and environment for statistical computing*. Vienna,
1194 Austria: R Foundation for Statistical Computing.
- 1195 Riad, T. (2014). *The phonology of Swedish*. Oxford: Oxford University Press.
- 1196 Richter, C., Feldman, N. H., Salgado, H., and Jansen, A. (2017). Evaluating low-level
1197 speech features against human perceptual data. *Transactions of the Association for
1198 Computational Linguistics* 5, 425–440. doi:[10.1162/tacl_a_00071](https://doi.org/10.1162/tacl_a_00071).
- 1199 Robb, M. P., and Chen, Y. (2009). Is /h/ phonetically neutral? *Clinical Linguistics &
1200 Phonetics* 23, 842–855. doi:[10.3109/02699200903247896](https://doi.org/10.3109/02699200903247896).
- 1201 RStudio Team (2020). *RStudio: Integrated development environment for R*. Boston, MA:
1202 RStudio, PBC.
- 1203 Schaeffler, F. (2005). *Phonological quantity in swedish dialects: Typological aspects, phonetic
1204 variation and diachronic change*. Umeå: Umeå University, Dep. of philosophy and
1205 linguistics.
- 1206 Sjerps, M. J., Fox, N. P., Johnson, K., and Chang, E. F. (2019). Speaker-normalized
1207 sound representations in the human auditory cortex. *Nature Communications* 10, 1–9.
1208 doi:[10.1038/s41467-019-10365-2](https://doi.org/10.1038/s41467-019-10365-2).
- 1209 Stilp, C. (2020). Acoustic context effects in speech perception. *WIREs Cognitive Science* 11.
1210 doi:[10.1002/wcs.1517](https://doi.org/10.1002/wcs.1517).
- 1211 Sussman, H. M. (1986). A neuronal model of vowel normalization and representation. *Brain
1212 and Language* 28, 12–23. doi:[10.1016/0093-934X\(86\)90087-8](https://doi.org/10.1016/0093-934X(86)90087-8).
- 1213 Syrdal, A. K. (1985). Aspects of a model of the auditory representation of american english
1214 vowels. *Speech Communication* 4, 121–135. doi:[10.1016/0167-6393\(85\)90040-8](https://doi.org/10.1016/0167-6393(85)90040-8).
- 1215 Syrdal, A. K., and Gopal, H. S. (1986). A perceptual model of vowel recognition based on
1216 the auditory representation of American English vowels. *The Journal of the Acoustical
1217 Society of America* 79, 1086–1100. doi:[10.1121/1.393381](https://doi.org/10.1121/1.393381).
- 1218 Tan, M., and Jaeger, T. F. (2022). Listeners adjust prior expectations. *Manuscript,
1219 Stockholm University*.
- 1220 Tan, M., Xie, X., and Jaeger, T. F. (2021). Using rational models to understand experiments
1221 on accent adaptation. *Frontiers in Psychology* 12, 1–19. doi:[10.3389/fpsyg.2021.676271](https://doi.org/10.3389/fpsyg.2021.676271).
- 1222 ten Bosch, L., Boves, L., and Ernestus, M. (2022). DIANA, a Process-
1223 Oriented Model of Human Auditory Word Recognition. *Brain Sciences* 12, 681.
1224 doi:[10.3390/brainsci12050681](https://doi.org/10.3390/brainsci12050681).
- 1225 Thomas, E. R., and Kendall, T. (2007). NORM: The vowel normalization and plotting suite.
- 1226 Toscano, J. C., and McMurray, B. (2010). Cue integration with categories: Weighting
1227 acoustic cues in speech using unsupervised learning and distributional statistics.
1228 *Cognitive Science* 34, 434–464.
- 1229 Toscano, J. C., and McMurray, B. (2015). The time-course of speaking rate compensation:
1230 Effects of sentential rate and vowel length on voicing judgments. *Language, Cognition
1231 and Neuroscience* 30, 529–543. doi:[10.1080/23273798.2014.946427](https://doi.org/10.1080/23273798.2014.946427).

- 1232 Traunmüller, H. (1981). Perceptual dimension of openness in vowels. *The Journal of the*
1233 *Acoustical Society of America* 69, 1465–1475. doi:[10.1121/1.385780](https://doi.org/10.1121/1.385780).
- 1234 Verbrugge, R. R., and Shankweiler, D. P. (1977). Prosodic information for vowel identity.
1235 *Journal of the Acoustical Society of America* 61.
- 1236 Watt, D., and Fabricius, A. (2002). “Evaluation of a technique for improving the mapping
1237 of multiple speakers’ vowel spaces in the F1 ~ F2 plane,” in *Leeds Working Papers in*
1238 *Linguistics and Phonetics*, ed. D. Nelson, 159–173.
- 1239 Weatherholtz, K., and Jaeger, T. F. (2016). Speech perception and generalization across
1240 talkers and accents. *Oxford Research Encyclopedia of Linguistics*. doi:[10.1093/acrefore/9780199384](https://doi.org/10.1093/acrefore/9780199384)
- 1241 Xie, X., Buxó-Lugo, A., and Kurumada, C. (2021). Encoding and decoding
1242 of meaning through structured variability in speech prosody. *Cognition* 211.
1243 doi:[10.1016/j.cognition.2021.104619](https://doi.org/10.1016/j.cognition.2021.104619).
- 1244 Xie, X., and Jaeger, T. F. (2020). Comparing non-native and native speech: Are L2
1245 productions more variable? *The Journal of the Acoustical Society of America* 147,
1246 3322–3347. doi:[10.1121/10.0001141](https://doi.org/10.1121/10.0001141).
- 1247 Xie, X., Jaeger, T. F., and Kurumada, C. (2022). What we do (not) know about
1248 the mechanisms underlying adaptive speech perception: A computational review.
1249 doi:[10.17605/OSF.IO/Q7GJP](https://doi.org/10.17605/OSF.IO/Q7GJP).
- 1250 Yang, J., and Fox, R. A. (2014). Perception of English Vowels by Bilingual Chinese–
1251 English and Corresponding Monolingual Listeners. *Language and Speech* 57, 215–237.
1252 doi:[10.1177/0023830913502774](https://doi.org/10.1177/0023830913502774).
- 1253 Young, N. J., and McGarrah, M. (2021). Forced alignment for Nordic languages:
1254 Rapidly constructing a high-quality prototype. *Nordic Journal of Linguistics*, 1–27.
1255 doi:[10.1017/S033258652100024X](https://doi.org/10.1017/S033258652100024X).
- 1256 Zahorian, S. A., and Jagharghi, A. J. (1991). Speaker normalization of static and dynamic
1257 vowel spectral features. *The Journal of the Acoustical Society of America* 90, 67–75.
1258 doi:[10.1121/1.402350](https://doi.org/10.1121/1.402350).