

1

The acoustic characteristics of Swedish vowels

2

Anna Persson<sup>1</sup>

3

<sup>1</sup> Department of Swedish Language and Multilingualism, Stockholm University

4

Author Note

5

Correspondence concerning this article should be addressed to Anna Persson,

6

Department of Swedish Language and Multilingualism, Stockholm University. E-mail:

7

[anna.persson@su.se](mailto:anna.persson@su.se)

8 Abstract

9 The Swedish vowel space is relatively densely populated with 21 categories that differ in  
10 quality and quantity. Existing descriptions of the entire space rest on recordings made in  
11 the late 1990s or earlier, while recent work in general has focused on sub-sets of the space.

12 The present paper reports two studies. The first presents static and dynamic acoustic  
13 analyses of the entire vowel space using a recently released database of *h-VOWEL-d* words  
14 (SwehVd). The second study compares the acoustic characteristics of SwehVd against  
15 formant data from 8 talkers of the same dialect recorded in 1999, to investigate the extent  
16 to which the phonetic realization of vowels has changed over the last generation. The  
17 results highlight the importance of static and dynamic spectral and temporal cues for  
18 Swedish vowel category distinction, and indicate possible ongoing shifts in the high front  
19 part of the space

20 *Keywords:* vowels, category separability, formant dynamics, vowel change

21 Word count: X

<sup>22</sup> The acoustic characteristics of Swedish vowels

## <sup>23</sup> 1 Introduction

<sup>24</sup> The Swedish vowel inventory consists of 21 categories that differ in spectral (formant  
<sup>25</sup> frequencies) and temporal cues (duration). It forms a typologically rather complex space,  
<sup>26</sup> characterized by a systematic quantity distinction resulting in 9 long and short vowel pairs,  
<sup>27</sup> 3 different levels of lip-rounding, and contextually conditioned allophones to /ɛ/ and /ø/ in  
<sup>28</sup> position before /r/ or any retroflex segments. Given the crowdedness of the space and  
<sup>29</sup> resulting category overlap for some vowels, previous work has reported on the need to look  
<sup>30</sup> beyond static point estimates of the two primary determinants to vowel identity  
<sup>31</sup> cross-linguistically, the first two formants (F1 and F2, e.g., Joos, 1948; Ladefoged &  
<sup>32</sup> Broadbent, 1957; Nearey & Assmann, 1986; Peterson, 1961), as well as the hypothesized  
<sup>33</sup> importance of additional cues besides F1 and F2, such as the third formant (F3) for  
<sup>34</sup> rounded vs. unrounded categories (e.g., Fant, 1959; Fant, Henningsson, & Stålhammar,  
<sup>35</sup> 1969; Fujimura, 1967; Kuronen, 2000), duration for certain long-short vowel pairs (Behne,  
<sup>36</sup> Czigler, & Sullivan, 1997), and formant dynamics for some front contrasts (Kuronen, 2000;  
<sup>37</sup> Pelzer & Boersma, 2019).

<sup>38</sup> This paper investigates the acoustic characteristics of modern-day Swedish vowels in  
<sup>39</sup> two studies that aim to contribute to our understanding of language-specific and  
<sup>40</sup> language-general patterns of vowel acoustics. The first study presents a comprehensive  
<sup>41</sup> description of the primary acoustic cues to vowel identity, using a recently released database  
<sup>42</sup> of *h-VOWEL-d* (short: hVd) words recorded by 44 male and female talkers of Swedish, the  
<sup>43</sup> SwehVd database (Persson & Jaeger, 2023). The second study assesses whether the  
<sup>44</sup> Swedish vowel space has been submitted to changes over the last generation, by comparing  
<sup>45</sup> the acoustics of SwehVd against that of a reference material from 1999 (Eriksson, 2004).  
<sup>46</sup> The variety investigated is Central Swedish, the regional standard variety of Swedish

47 spoken in an area around and beyond Stockholm (eastern Svealand) (Bruce, 2009; Elert,  
48 1994; Riad, 2014).<sup>1</sup> Existing descriptions of the entire space of 21 vowels rest on recordings  
49 made more than 25 years ago (reported in, e.g., Engstrand, 1999; Kuronen, 2000; T.  
50 Leinonen, 2010; Riad, 2014). Two of the most recent studies are T. Leinonen (2010) and  
51 Kuronen (2000) (Table 1). The former is based on recordings obtained around 1999 of all  
52 vowels, of which four short vowels were omitted from analysis. It covers 98 rural locations  
53 in Sweden and Swedish-speaking parts of Finland, including reference talkers of Standard  
54 Swedish. The latter covers the entire vowel space but is based on recordings from 1981 (K.  
55 Leinonen, Pitkänen, & Vihanta, 1981). More recent work over the last two decades has  
56 focused on parts of the phonological space, such as the long vowels for diphthongization  
57 studies (Pelzer & Boersma, 2019), two vowels for merger studies (e.g., [θ] - [œ] in Wenner,  
58 2010), or a single vowel, e.g., the damped [i:] (Schötz, Frid, & Löfqvist, 2011), see Table 1.  
59 These studies all provide detailed mappings of different parts of the space, and contribute  
60 important insights into the current state of as well as ongoing processes. However, given  
61 their focus on subsets of the space, a comprehensive acoustic mapping of the modern-day  
62 Central Swedish vowel space *in its entirety* is lacking. Given that there is some evidence  
63 that productions of minimal pairs can lead to enhanced contrasts (e.g., Schertz, 2013;  
64 Seyfarth, Buz, & Jaeger, 2016), how representative such subsets are for the vowel space as  
65 a whole, remains an open question. In addition, most previous studies differ in the  
66 materials used, in terms of the size of the database (e.g., number of talkers and repetitions  
67 per vowel), the demographics of talkers (e.g., male/female talkers, region of origin), and  
68 phonological contexts used for recording. For instance, the majority of previous work has  
69 either not held the phonetic context constant across vowels, or has investigated isolated  
70 vowel production out of context or in different CVC contexts (Table 1). This diversity  
71 restricts comparison across studies on Swedish, as well as cross-linguistically.

---

<sup>1</sup> For the reader unfamiliar with Central Swedish, Section 1.1 provides an overview of the acoustics of Central Swedish vowel space.

Table 1  
*A selection of previous studies on Central Swedish vowels*

Article	Speech materials	Participants	Approach	Focus
Eklund & Traunmüller, 1997	3 repetitions of all 9 long vowels in isolation	12 talkers (6 female) from the Greater Stockholm area, 20-58 years of age	Formant trajectories at 10 measurement points; F1-F2 means and SD; linear regression	Comparing the acoustics and perception of whispered and phonated vowels
Elert, 1964	1 repetition of all 21 vowels in sentence lists and word list (word list recorded by 2 talkers only; different V:C and VC: contexts)	11 talkers (5 female) from Stockholm, born in the 1930's	Duration means, SD and SE; t-statistics	Phonological quantity
Eriksson, 2004	3-4 repetitions of all 21 vowels, different phonological contexts for each vowel	12 talkers (6 female; 6 young) each from 107 locations across Sweden and Finland; 12 reference talkers (6 female; 6 young) of standard Swedish from the Greater Stockholm area	Project description	SweDia dialect database development
Fant, Henningsson & Stålhammar, 1969	1 repetition of all 9 long vowels in isolation	24 male talkers, students at KTH in Stockholm	F1 to F4 and duration means	Formant frequencies of the long vowels
Kuronen, 2000	5-7 repetitions of all 21 vowels, different phonological contexts for each vowel and each repetition	4 male talkers from Nyköping, 17-18 years of age (4 female reference talkers), recorded by Leinonen, Pitkänen & Vihanta, 1982	F0, F1 to F4 and duration means; formant trajectories at 4 measurement points with 30 ms intervals (long vowels)	Spectral and temporal acoustics of all Central Swedish vowels
Leinonen, 2010	a subset of 19 vowels from the SweDia database (see Eriksson, 2004 above)	see Eriksson, 2004	Principal component analysis of Bark-filtered vowel spectra; multidimensional scaling for dialectometric analysis	Dialectal variation in Swedish vowel pronunciation
Lindblom, 1963	5 repetitions of all 8 short vowels in bVb, dVd, gVg contexts, with four different stress patterns	1 male talker from Stockholm, 19 years of age	F1 to F3 means, duration; formants as a function of duration and consonantal context	The effect of vowel duration on formant frequencies
McAllister, Lubker & Carlson, 1974	10 repetitions of all 6 long rounded vowels in iV context	6 talkers of standard Swedish	EMG; F1 to F3 means	The articulation and acoustics of rounded vowels
Pelzer & Boersma, 2019	8 repetitions of all 9 long vowels, different phonological contexts for each vowel and each repetition	8 talkers (4 female) from Stockholm	F1-F2 median values at 20, 50, 80% into the vowel; linear mixed effects model	Diphthongization in the long vowels
Schötz, Frid & Löfqvist, 2011	6 repetitions of [i:] in two different contexts (bibel, papipa)	1 male talker from Stockholm	Articulography; F1 to F4 means; Bark-circles	The articulation and acoustics of the damped [i:]
Wenner, 2010	3 repetitions of [œ] and [ø] in different phonological contexts	78 talkers (40 female) from 4 locations in Uppland, 12-85 years of age	F1 to F3 means; linear regression; correlation analysis	The merger of [œ] and [ø]

72        The materials and methodological approach adopted in the current paper is motivated  
73    by the goal to complement previous work for a comprehensive picture of modern-day  
74    Central Swedish vowels. The first study provides an up-to-date acoustic description of the  
75    entire vowel space of 21 categories, using the SwehVd database (Persson & Jaeger, 2023).  
76    The use of an hVd database minimizes coarticulatory effects from the surrounding phonetic  
77    context, and increases cross-linguistic comparability (e.g., Hillenbrand, Getty, Clark, &  
78    Wheeler, 1995; Peterson & Barney, 1952). The main spectral and temporal cues to vowel  
79    identity are reported in static analyses (following e.g., Engstrand, 1999; Fant et al., 1969),  
80    as well as dynamic analyses, given the well documented importance of formant dynamics  
81    on vowel production and perception (e.g., Assmann & Katz, 2005; Eklund & Traunmüller,  
82    1997; Hillenbrand & Nearey, 1999; Kuronen, 2000; Nearey & Assmann, 1986). The static  
83    analysis assesses what cues contribute to vowel distinctions and evaluates some of the  
84    claims introduced in previous work, such as the hypothesized importance of F3 for rounded  
85    vs. unrounded high front contrasts (Fant, 1959; Fant et al., 1969; Fujimura, 1967; Kuronen,  
86    2000), and to what extent spectral and temporal cues contribute to long-short vowel pair  
87    distinctions (e.g., Behne et al., 1997; Kuronen, 2000). The dynamic analysis explores what  
88    part of the space seems more prone to diphthongization, and investigates how formant  
89    dynamics contribute to vowel distinctions. In contrast with previous work investigating the  
90    dynamics of Central Swedish vowels, the present study includes both long and short vowels,  
91    thus submitting the entire vowel space to the same analyses.

92        Anticipating some of the main results from Study 1, the acoustic analyses suggested  
93    differences in the phonetic realization of some vowels compared to other qualitative  
94    descriptions of Central Swedish (Eklund & Traunmüller, 1997; Engstrand, 1999; Fant et al.,  
95    1969; Kuronen, 2000; Riad, 2014). These included the fronting of [e:], the lowering of [ε:],  
96    and the centralization of [i:] and [y:] (c.f., T. Leinonen, 2010; Pelzer & Boersma, 2019;  
97    Schötz et al., 2011). This led me to conduct Study 2, which aimed to assess the scope and  
98    spread of these changes over the last generation, by comparing the acoustic characteristics

99 of SwehVd against vowel data from a database of 8 (female = 4) talkers of Central Swedish  
100 recorded in 1999 as part of the SweDia dialect project (Eriksson, 2004).

101 The paper is organized as follows. A background to the acoustics of Central Swedish  
102 vowels is provided by a review of previous work. This is followed by the two studies'  
103 methods and results, and finally, a discussion of the results and its consequences for the  
104 Central Swedish vowel system. All analyses and visualization code for this study can be  
105 found in an online repository (<https://osf.io/7uvj4/>). This article is written in R  
106 Markdown, which allows readers to easily replicate the analyses using freely available  
107 software (R Core Team, 2023; RStudio Team, 2020).

## 108 1.1 The acoustics of Central Swedish vowels

109 This section provides a description of the overall inventory of Central Swedish  
110 monophthongs, and discusses the role of cues beyond F1 and F2. It furthermore presents a  
111 review of previous studies on diphthongization and formant dynamics.

112 Central Swedish is most often described as having nine vowel phonemes: /i/, /y/,  
113 /u/, /e/, /ɛ/, /ø/, /ɑ/, /o/, /u/. The long allophones are [i:], [y:], [u:], [e:], [ɛ:], [ø:], [ɑ:], [o:],  
114 [u:], and the short allophones are [i], [y], [ø], [ɛ], [ø], [ɑ], [ɔ], [u]. The short allophones of /e/  
115 and /ɛ/ has been reported to neutralize as [ɛ] in Central Swedish, resulting in 17 vowels,  
116 rather than 18 (Riad, 2014). There is also evidence of neutralization of the short /ø/ and  
117 /u/ as [ø] among some talkers, primarily in position before a retroflex (Ståhlé, 1965;  
118 Wenner, 2010). In addition to these 17 vowels, there are 4 additional long and short  
119 allophones—[æ:], [æ], [œ:], and [œ], as /ɛ/ and /ø/ lower in position before /r/ or any  
120 retroflex segment (e.g., Kuronen, 2000; Riad, 2014). Traditionally, Central Swedish has  
121 been described using four height levels and three backness levels (Riad, 2014).

122 It has furthermore been suggested that Central Swedish is defined by three levels of  
123 lip-rounding, where the rounded vowels are most often referred to as either inrounded, with

<sup>124</sup> an extreme narrowing of the lips—[u:] and [u], or outrounded, with a lesser degree of  
<sup>125</sup> lip-narrowing and more protruded lips—[y:], [ø:], [œ:], [ɔ:], and the remaining vowels defined  
<sup>126</sup> as unrounded (e.g., Fant, 1971; McAllister, Lubker, & Carlson, 1974). Previous work has  
<sup>127</sup> claimed that lip-rounding is particularly important for some vowel distinctions. For  
<sup>128</sup> instance, [i:] and [y:] have been described as overlapping in F1-F2 space, but as more  
<sup>129</sup> separable when F3 is considered (Fant, 1959; Fant et al., 1969; Fujimura, 1967; Kuronen,  
<sup>130</sup> 2000).

<sup>131</sup> The vowels in each pair have been reported to differ systematically in duration, with  
<sup>132</sup> short-long vowel to vowel ratios on average .65-.67 for Central Swedish (Elert, 1964;  
<sup>133</sup> Kuronen, 2000; Strangert, 2001). Spectral differences have traditionally been interpreted as  
<sup>134</sup> a consequence of the durational distinction, hence assuming a trading relationship between  
<sup>135</sup> spectral and temporal cues (for a review, see Schaeffler, 2005). It has been hypothesized  
<sup>136</sup> that most of the durational variation is carried by F2 (e.g., Kuronen, 2000; Lindblom,  
<sup>137</sup> 1963). Previous work has found the largest spectral differences for the [u:] - [ø], and [a:] - [a]  
<sup>138</sup> vowel pairs, and the smallest differences for [ε:] - [ɛ], and [ø:] - [ø] (e.g., Kuronen, 2000). For  
<sup>139</sup> pairs with small spectral differences, duration is presumably more important for vowel  
<sup>140</sup> distinction. Perceptual studies on synthesized speech from talkers of Stockholm Swedish  
<sup>141</sup> have confirmed that duration is the primary cue for [i:] - [ɪ], and [o:] - [ɔ] (Behne et al.,  
<sup>142</sup> 1997; for results on Southern Swedish and additional vowel pairs, [ε:] - [ɛ], [ø:] - [ø], see  
<sup>143</sup> Hadding-Koch & Abramson, 1964). The extent to which *all* long-short vowel pairs rely on  
<sup>144</sup> spectral cues is less known, as studies have focused on subsets of pairs.

<sup>145</sup> According to previous work, several of the long vowels in Central Swedish tend to  
<sup>146</sup> diphthongize in their phonetic realization. Diphthongization is considered prosodically  
<sup>147</sup> conditioned and is the strongest in stressed vowels (Bleckert, 1987; T. Leinonen, 2010).<sup>2</sup>  
<sup>148</sup> Previous studies have characterized the diphthongal glide in the later part of the long

---

<sup>2</sup> In general, true phonological diphthongs are not considered part of the phonological inventory of Central Swedish (Eliasson, 2022).

149 vowels as either a centralization of the vowel segment towards [ø] or a more open quality, or  
150 as a consonantal offglide (e.g., Elert, 1981, 2000; Fant, 1971; Fant et al., 1969; Kuronen,  
151 2000; McAllister et al., 1974; Pelzer & Boersma, 2019; Riad, 2014). Results are  
152 inconclusive as to how widespread diphthongization is across the vowel space, and what  
153 direction it takes. Most work has however found substantial diphthongization towards a  
154 more open quality for the mid and mid-high vowels [e:], [ø:] and [o:] (Eklund &  
155 Traunmüller, 1997; Elert, 2000; Fant et al., 1969; Pelzer & Boersma, 2019). In addition,  
156 diphthongization has been hypothesized to cue vowel distinctions for certain high vowels  
157 ([i:], [y:], [œ], and [u:]) (e.g., Fant, 1971; Kuronen, 2000). For instance, Kuronen (2000)  
158 reported that [i:] - [y:] - [e:], and [u:] - [o:], differed in formant patterns only at later  
159 time-points of the vowel for some talkers and that the contrast between [ɛ:] and [æ:] was  
160 maintained solely by trajectory movements. Of importance for the present study, less is  
161 known about the formant dynamics in the short vowels, given the almost exclusive focus on  
162 the long vowels in diphthongization studies.

163 Some talkers of Central Swedish have been reported to realize [i:], [y:], [œ] and [u:]  
164 with a consonantal offglide, where the end-point of [i:] is described as a palatal  
165 approximant [j], the end-point of [y:] a voiced labio-palatal approximant [ɥ], the end-point  
166 of [œ] and [u:] a voiced bilabial fricative [β] (Elert, 1980; Hammarström & Norman, 1957;  
167 McAllister et al., 1974). Furthermore, both [i:] and [y:] can be damped and produced with  
168 a buzzing sound, phonetically realized as [i̯]. The damped [i̯] has been found in several  
169 dialects across Sweden, both in rural areas and in the cities of Gothenburg and Stockholm  
170 (Björsten & Engstrand, 1999; Elert, 1980; Engstrand, Björsten, Lindblom, Bruce, &  
171 Eriksson, 2000; Gross & Forsberg, 2020; Riad, 2014). In work on Swedish dialectology, it is  
172 often referred to as the Viby-*i*, and in the Stockholm area, as the Lidingö-*i*. Acoustically, it  
173 manifests primarily as a lowering of F2, thus occupying a more centralized position in the  
174 F1-F2 space. Schötz et al. (2011) describe it as a central palatal vowel, as the articulatory  
175 correlates involve a retracted and lower tongue position, the tip of the tongue being higher

<sup>176</sup> than blade and dorsum.

## <sup>177</sup> **2 Study 1: Static and dynamic analyses of the 178 spectral and temporal properties of Central Swedish**

<sup>179</sup> Study 1 describes the spectral acoustics of Central Swedish in static and dynamic analyses  
<sup>180</sup> of F0, F1, F2, F3, and duration. It aims to provide a detailed description of the acoustics  
<sup>181</sup> of all Central Swedish vowels, and to evaluate the relative importance of certain cues for  
<sup>182</sup> specific vowel contrasts, as hypothesized in previous work. These include the importance of  
<sup>183</sup> lip-rounding (F3) for high vowel distinctions (Fant, 1959; Fant et al., 1969; Fujimura, 1967;  
<sup>184</sup> Kuronen, 2000), to what extent all long-short vowel pairs differ in quality (formants) and  
<sup>185</sup> quantity (duration) (e.g., Behne et al., 1997; Kuronen, 2000), and what vowels seem to  
<sup>186</sup> undergo diphthongization (Kuronen, 2000; Pelzer & Boersma, 2019). The dynamic analysis  
<sup>187</sup> furthermore explores which cues carry information about neighboring vowel distinctions  
<sup>188</sup> once dynamic information is considered.

<sup>189</sup> The methodology employed in Study 1 is presented next, beginning with a  
<sup>190</sup> description of the materials used.

### <sup>191</sup> **2.1 Methods**

#### <sup>192</sup> **2.1.1 Materials**

<sup>193</sup> The materials used in Study 1 is a corpus of Swedish hVd word recordings, collected by  
<sup>194</sup> Anna Persson and Maryann Tan (Stockholm University) in 2020-2024, the SwehVd. An  
<sup>195</sup> initial version of the corpus with 24 female talkers is described in Persson and Jaeger  
<sup>196</sup> (2023). For this paper, an updated release is presented, including 20 additional male talkers  
<sup>197</sup> (targeted number of male talkers = 24). All recordings, annotations, and acoustic  
<sup>198</sup> measurements are available at <https://osf.io/ruxnb/>. SwehVd covers the entire

<sup>199</sup> monophthong inventory of Central Swedish, including all nine long vowels, eight short  
<sup>200</sup> vowels, and the four allophones to /ɛ/ and /ø/.<sup>3</sup> SwehVd focuses on a single regional  
<sup>201</sup> variety, providing high resolution within and across talkers for this variety with N = 10  
<sup>202</sup> recordings of each hVd word from each of the N = 44 talkers (N = 24 female), for a total N  
<sup>203</sup> of tokens = 9103. All talkers in the database were L1 talkers of Swedish, born and raised in  
<sup>204</sup> the Greater Stockholm area or surroundings, of 18-44 years of age (mean age = 30; SD =  
<sup>205</sup> 6.82). For more details on the recruitment, recording, pre-processing, segmentation and  
<sup>206</sup> annotation procedure, see Persson and Jaeger (2023).

<sup>207</sup> For the vast majority of talkers in the SwehVd, *hädd* productions elicited the same  
<sup>208</sup> vowel as *hedd* (see Supplementary Information—SI, Figure S1), which confirms the  
<sup>209</sup> commonly held assumption that the short allophone of /e/ neutralizes with the short  
<sup>210</sup> allophone of /ɛ/ in Central Swedish. In order to have a balanced number of tokens for each  
<sup>211</sup> vowel, all *hädd* words were excluded from the subsetted SwehVd materials used in this  
<sup>212</sup> study (following Persson & Jaeger, 2023). Recordings on which the talker did not produce  
<sup>213</sup> the targeted vowel were also excluded.<sup>4</sup> Furthermore, outliers were identified and removed  
<sup>214</sup> by estimating the relative probability of each token's F1-F2 values given the joint  
<sup>215</sup> distribution of F1-F2 for that vowel and talker. Tokens outside of the 2.50th to 97.50th  
<sup>216</sup> quantile of the bivariate Gaussian distribution were filtered out. To facilitate empirical  
<sup>217</sup> analyses and statistical models, all talkers (N = 7) with fewer than 4 remaining recordings  
<sup>218</sup> for at least one of the vowels were removed. This left data from 37 L1 talkers (N=20  
<sup>219</sup> female talkers), with on average 322 (SD = 20) tokens per vowel (range = 277 to 345), for  
<sup>220</sup> a total of 6759 observations.

---

<sup>3</sup> The words used to elicit the 21 vowels were: *hid-[i]*, *hyd-[y]*, *hud-[ɯ]*, *hed-[ε]*, *häd-[ɛ]*, *höd-[ø]*, *had-[ɑ]*, *håd-[ɔ]*, *hod-[u]*, *hidd-[ɪ]*, *hydd-[ʏ]*, *hudd-[ø]*, *hedd-[ɛ]*, *hädd-[ɛ]*, *hödd-[ø]*, *hadd-[ɑ]*, *hådd-[ɔ]*, *hodd-[ʊ]*, *härd-[æ]*, *härr-[æ]*, *hörd-[œ]*, *hörr-[œ]*.

<sup>4</sup> The SwehVd database contains information on both targeted vowel and what vowel was actually produced.

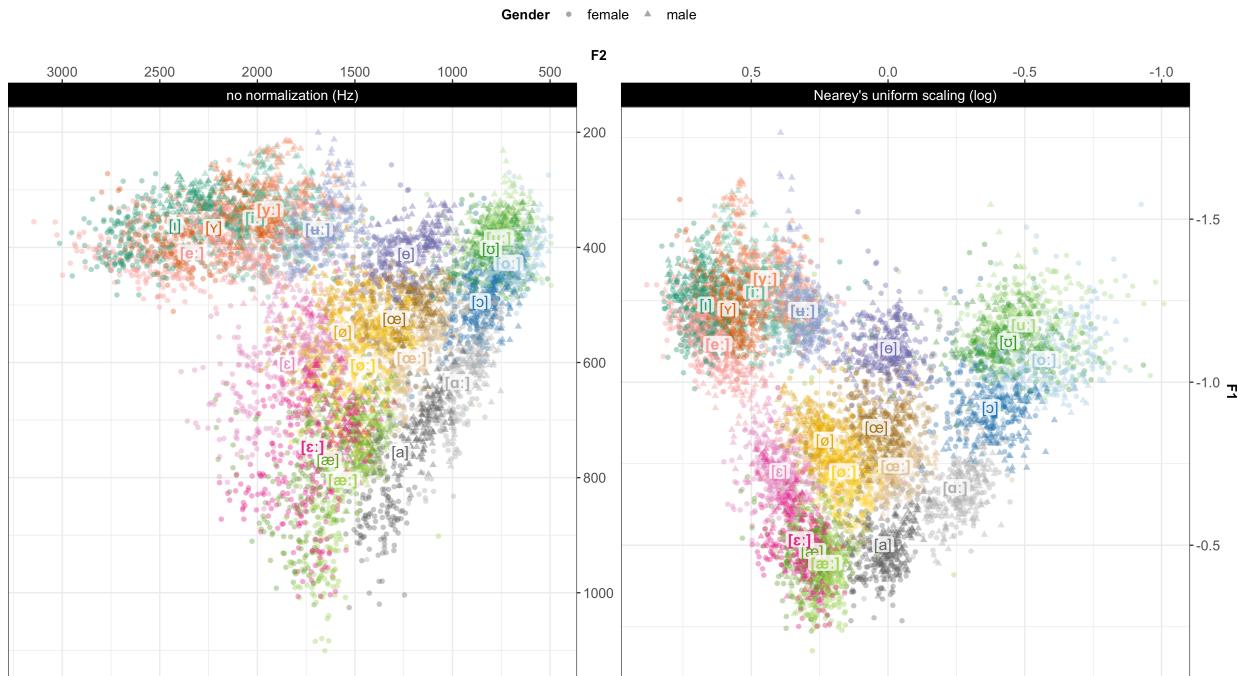
221 **2.1.2 Acoustic analyses**

222 **2.1.2.1 Measuring acoustic cues to vowel identity** The Swedish version of the  
223 Montreal Forced Aligner developed by Young and McGarrah (2021) was used to obtain  
224 estimates of word and segment boundaries. The boundaries were then manually corrected  
225 by the author (an L1-talker of Swedish). The formant analysis was carried out in Praat  
226 (Boersma & Weenink, 2022), using the Burg algorithm to extract estimates of the first three  
227 formants (F1-F3) at five time-points of the vowel (20, 35, 50, 65, 80% into the vowel), while  
228 vowel duration and F0 were extracted across the entire vowel segment. The Burg algorithm  
229 was parameterized with a time step of 0.01 seconds, a window length of 0.025 seconds, and  
230 pre-emphasis was applied from 50 Hz. The maximum number of formants was set to 5,  
231 with a formant ceiling of 5500 Hz for the female talkers, and 5000 Hz for the male talkers.

232 **2.1.2.2 Vowel normalization** The raw formant values were transformed into a vowel  
233 normalized space using Nearey's uniform scaling account (Nearey, 1978). Formant  
234 measurements in Hertz are reported in the SI, Section S1.2.4. Vowel normalization is used  
235 in studies on vowel production and perception to account for acoustically irrelevant  
236 inter-talker variation, as caused by differences in anatomical structure, e.g., vocal tract size  
237 (for reviews see e.g., Barreda & Nearey, 2018; Johnson & Sjerps, 2021; Stilp, 2020). In  
238 vowel production studies such as the present, normalization is primarily used as a  
239 methodological tool. Transforming the formant data into a normalized space reduces  
240 differences in F1 and F2 due to physiology, which can reduce between-talker variability and  
241 increase category separability, as visualized in Figure 1 (compare left and right panel).

242 Previous work on Swedish has primarily analyzed vowel data in raw Hertz (Björsten  
243 & Engstrand, 1999; Fant et al., 1969; Pelzer & Boersma, 2019), or transformed into Bark  
244 (Fant, 1983; Kuronen, 2000; Schötz et al., 2011; Wenner, 2010), Mel (Lindblom, 1963), or  
245 Lobanov (Gross & Forsberg, 2020). The choice of Nearey's uniform scaling in the present  
246 study was motivated by its previous use in socio-phonetic research to describe and compare

<sup>247</sup> languages and varieties (e.g., Barreda, 2021; Labov, 2001; Labov, Ash, & Boberg, 2005),  
<sup>248</sup> and by its plausibility as perceptual model of how we come to achieve robust cross-talker  
<sup>249</sup> perception, as it has provided a good fit against both production (e.g., Persson & Jaeger,  
<sup>250</sup> 2023; Syrdal, 1985) and perception data (e.g., Barreda, 2021; Persson, Barreda, & Jaeger,  
<sup>251</sup> 2024).



*Figure 1.* The SwehVd vowel data in unnormalized Hertz (*left*) and Nearey's uniform scaling space (*right*), along the first two formants, F1 and F2. Points show recordings of each of the 21 Central Swedish vowels by 44 (24 female) L1 talkers in the database, averaged across the three middle time-points (at 35, 50, 65% into the vowel). Vowel labels are placed at the vowel mean across talkers. Long vowels are boldfaced. Vowels that mismatched intended label are excluded (1.33% of all recordings).

<sup>252</sup> **2.1.2.3 Static acoustic analysis** The static analysis of SwehVd presents formant  
<sup>253</sup> measurements at the steady state of the vowel, by averaging across the three mid-points.<sup>5</sup>

<sup>5</sup> The choice of time-point for extracting formants, or whether to average across several time-points, affects the acoustic characterizations given how formants move across the vowel segment. The SI (S1.2.2) presents evaluations of the effect of different measurement points. For increased comparability across studies in this paper, the steady state of the vowel was selected for both studies.

254 It maps the entire vowel space of 21 categories and evaluates the relative contribution of  
 255 F0, F1, F2, F3 and duration to vowel distinctions, using visualizations of cues and cue  
 256 correlations. While fundamental frequency (F0) is not considered an important cue to  
 257 vowel identity in itself, it is known to vary between languages, dialects and speech styles  
 258 (e.g., Henton, 2005; Jacewicz & Fox, 2018; Johnson, 2005; Mennen, Schaeffler, & Docherty,  
 259 2012; Weirich, Simpson, Öjbro, & Ericsdotter Nordgren, 2019) and is therefore reported.

260 In order to evaluate the hypothesized importance of lip-rounding (F3) for neighboring  
 261 unrounded and rounded categories, a category separability index was employed. Following  
 262 work by Wedel, Nelson, and Sharp (2018) and X. Xie and Jaeger (2020), each vowel's  
 263 separability from the neighboring vowel was calculated as the average distance of vowel  
 264 tokens to the centroid of the neighboring vowel, operationalized as (1).

$$\text{separability of } /y:/ \text{ from } /i:/ = \frac{\sum_{k=1}^n \sqrt{(F1_{\text{token } k \text{ of } /y:/} - F1_{\text{Center of } /i:/})^2 + (F2_{\text{token } k \text{ of } /y:/} - F2_{\text{Center of } /i:/})^2}}{n} \quad (1)$$

265 For instance, for the [y:] - [i:] contrast, first, each talker's [i:] center was calculated for  
 266 F1-F2. Next, the distances between each [y:] token to the neighboring [i:] center from the  
 267 same talker were calculated for F1-F2. Finally, the distances were averaged across all [y:]  
 268 tokens from a talker, resulting in a separability measure for that vowel and talker. The  
 269 higher the index, the greater the separation between categories. The same was  
 270 subsequently done for F1-F2-F3. These two measures of separability for each contrast  
 271 (F1-F2, F1-F2-F3) were then compared to assess whether including F3 would lead to  
 272 increased category separability. The contrasts investigated were [y:] - [i:], [e:] - [y:], [ɪ] - [ʏ]  
 273 for comparing unrounded vs. outrounded vowels, and [y:] - [œ], [o:] - [u:], [ɔ] - [ʊ] for  
 274 outrounded vs. inrounded vowels.

275 To quantify the effect of including F3 on category separability, separate linear  
 276 mixed-effects model (LMM) were fit for each contrast, predicting separability from cue

277 combination (F1-F2-F3 vs. F1-F2) while including by talker random intercepts.<sup>6</sup> The  
278 model was formulated as follows:  $\text{separability} \sim \text{cuecombination} + (1|\text{Talker})$ . Cue  
279 combination was sum-coded ( $F1-F2 = 1$ ,  $F1-F2-F3 = -1$ ).

280 The same process was applied to investigate to what extent long-short vowel pairs  
281 differ in spectral cues, by assessing what combination of spectral cues could provide the  
282 largest separability between the two vowels in each pair. For this evaluation of quantity  
283 contrasts, the category separability index was calculated for each pair and four different  
284 cue combinations: F1-F2, F1-F3, F2-F3 or F1-F2-F3. The models were the same as the  
285 previous sets, however, cue combination was treatment-coded with F1-F2 as reference  
286 category, thus comparing each cue combination against the F1-F2 combination.

287 The results of the static analysis are presented in Section 2.2.1.

288 **2.1.2.4 Dynamic acoustic analysis** Formant measurements at all five time-points  
289 were used in the dynamic analysis to assess the importance of formant dynamics for vowel  
290 distinctions. The dynamic analysis is divided into two main sections. In the first section,  
291 formant trajectory plots were used to assess the scope and direction of formant movements,  
292 to what extent vowels seemed to diphthongize, and to evaluate the hypothesized  
293 importance of formant trajectories for the [i]-[y]-[e], [o]-[u] and [ɛ]-[æ] contrasts reported  
294 in previous work (e.g., Kuronen, 2000; Pelzer & Boersma, 2019). Lastly, trajectories of  
295 short vowels were also visualized as they have not been typically explored in the past.

296 In the second part of the dynamic analysis, the hypothesized contribution of formant  
297 dynamics to category information was modeled using generalized additive mixed-effects  
298 models (GAMMs) (Baayen, Vasishth, Kliegl, & Bates, 2017). GAMMs were employed to  
299 assess what cues carry information about vowel quality once formant dynamics were  
300 inspected. GAMMs are increasingly used in phonetic research, due to their suitability in  
301 modeling the non-monotonic complex phonetic patterns found in formants without

---

<sup>6</sup> By-talker random intercepts was the maximum random effect structure that converged.

302 assuming linearity or having to rely on the simplifying assumption that vowels can be  
 303 reduced to a single F1-F2 point estimate (e.g., Chuang, Fon, Papakyritsis, & Baayen, 2021;  
 304 Sóskuthy, 2021; Wieling, 2018). GAMMs have been used in studies on vowels in different  
 305 English varieties, e.g., on /u/-fronting in Derby English (Sóskuthy, Foulkes, Hughes, &  
 306 Haddican, 2018) and on the front vowel system of Southern American English (Renwick &  
 307 Stanley, 2020) but to the best of my knowledge, they have not been implemented in studies  
 308 of Swedish vowels. The use of GAMMs thus complements previous work on Central  
 309 Swedish that has primarily used visual inspection, formant measurements and linear  
 310 models (Table 1).

311 Two main groups of GAMMs were fit in the dynamic analysis. In the first group,  
 312 GAMMs were fit to 4 subsets of neighboring contrasts, hypothesized to differ primarily in  
 313 formant dynamics: [i:] - [y:] - [ɛ̄] - [e:], [o:] - [u:] - [ɛ̄] - [æ:] (Fant, 1971; Kuronen, 2000; Pelzer  
 314 & Boersma, 2019). Given the directionality in formant trajectories found for [ø:]-[œ:]  
 315 (Figure 7), this contrast was also included. To explore potential effects of dynamics in the  
 316 corresponding short vowels, an additional 4 contrasts were modeled. These were not  
 317 entirely identical to the long subsets, for reasons of evident separability in F1-F2 space: [ɪ] -  
 318 [ʏ], [ɔ] - [ʊ], [ɛ] - [æ] and [ø]-[œ]. The general model formulation was as follows  
 319  $formant \sim category + Gender + s(timepoint, by = category, k = 5) + s(Talker, bs =$   
 320 "re") +  $s(Talker, category, bs = "re")$ . The GAMMs were treatment coded with [i:], [o:],  
 321 [ɛ:], [ɪ], [ɔ], [ɛ], and [ø] as reference categories in respective set.

322 The second group consisted of GAMMs fit to all 21 categories, aiming for an  
 323 evaluation of differences between categories within vowel pairs. Vowel was backwards  
 324 difference coded, with [i:] as reference vowel. The vowel variable was ordered alternating  
 325 long and short vowels by the pair, so the short vowels were compared against their  
 326 preceding long counterpart.

327 All GAMMs were fit separately for each of the formants, which necessarily meant  
 328 committing to the simplifying assumption of cue independence. Previous work has shown

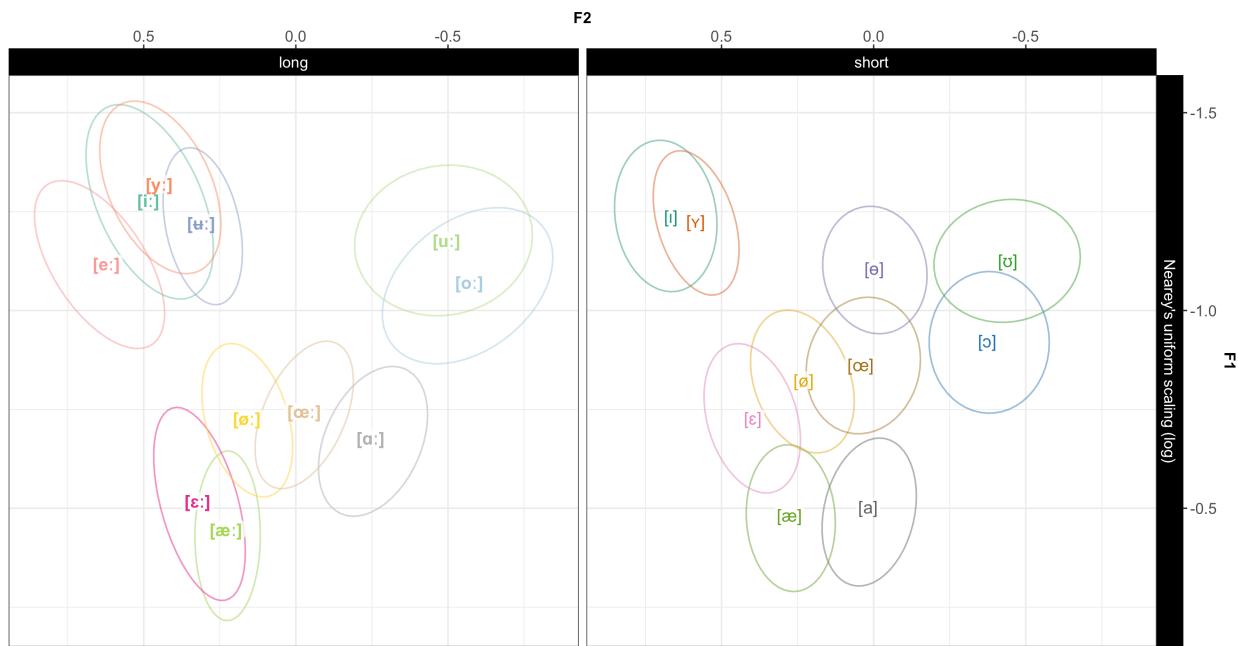
329 that acoustic cues tend to co-vary (for a review, see Schertz & Clare, 2020). For vowels,  
 330 this is the case for F1 and F2, as shown by the shape and orientation of ellipses in Figure 2.

331 The results of the dynamic analysis are presented in Section 2.2.2.

## 332 2.2 Results

### 333 2.2.1 Static spectral and temporal cues to vowel identity

334 The static analysis begins with a mapping of the entire 21 category space along F1-F2.  
 335 Next, the relative contribution of additional cues beyond F1 and F2 is assessed, as well as  
 336 the extent to which all long-short vowel pairs are qualitatively and quantitatively different.



*Figure 2.* The SwehVd vowel data separated by quantity. Ellipses show bivariate Gaussian 95% confidence interval of vowel means. Vowel labels indicate vowel means across female and male talkers.

337 Figure 2, left panel, visualizes the long vowels along the two primary cues to vowel  
 338 identity, F1-F2.<sup>7</sup> Four vowels cluster in the high front part of the space. The mid-high [e:]

<sup>7</sup> The SI presents the mean cue values for the male and female talkers, Tables S1 and S2. As expected, the

339 occupies a substantially higher position than in many previous descriptions, and is also the  
 340 most fronted vowel (c.f., Fant et al., 1969; Kuronen, 2000; but see Engstrand et al., 2000;  
 341 Pelzer & Boersma, 2019). The high [i:] and [y:] are rather mid-central, and exhibit  
 342 substantial overlap with [œ:]. The [u:] - [o:], and [ɛ:] - [æ:] categories are also partly  
 343 overlapping.

344 The short vowels (right panel), present a slightly more compact space, however with  
 345 increased category separability (c.f., Riad, 2014).<sup>8</sup> For some vowel pairs, overlap is clearly  
 346 reduced for the short vowels, e.g., for [ɪ] - [ʏ], [ɛ] - [æ], and [ɔ] - [ʊ]. Of note, the high vowels  
 347 [ɪ] and [ʏ] are more fronted than their long counterparts, which does not replicate previous  
 348 studies on Central Swedish (e.g., Fant, 1971; Kuronen, 2000).

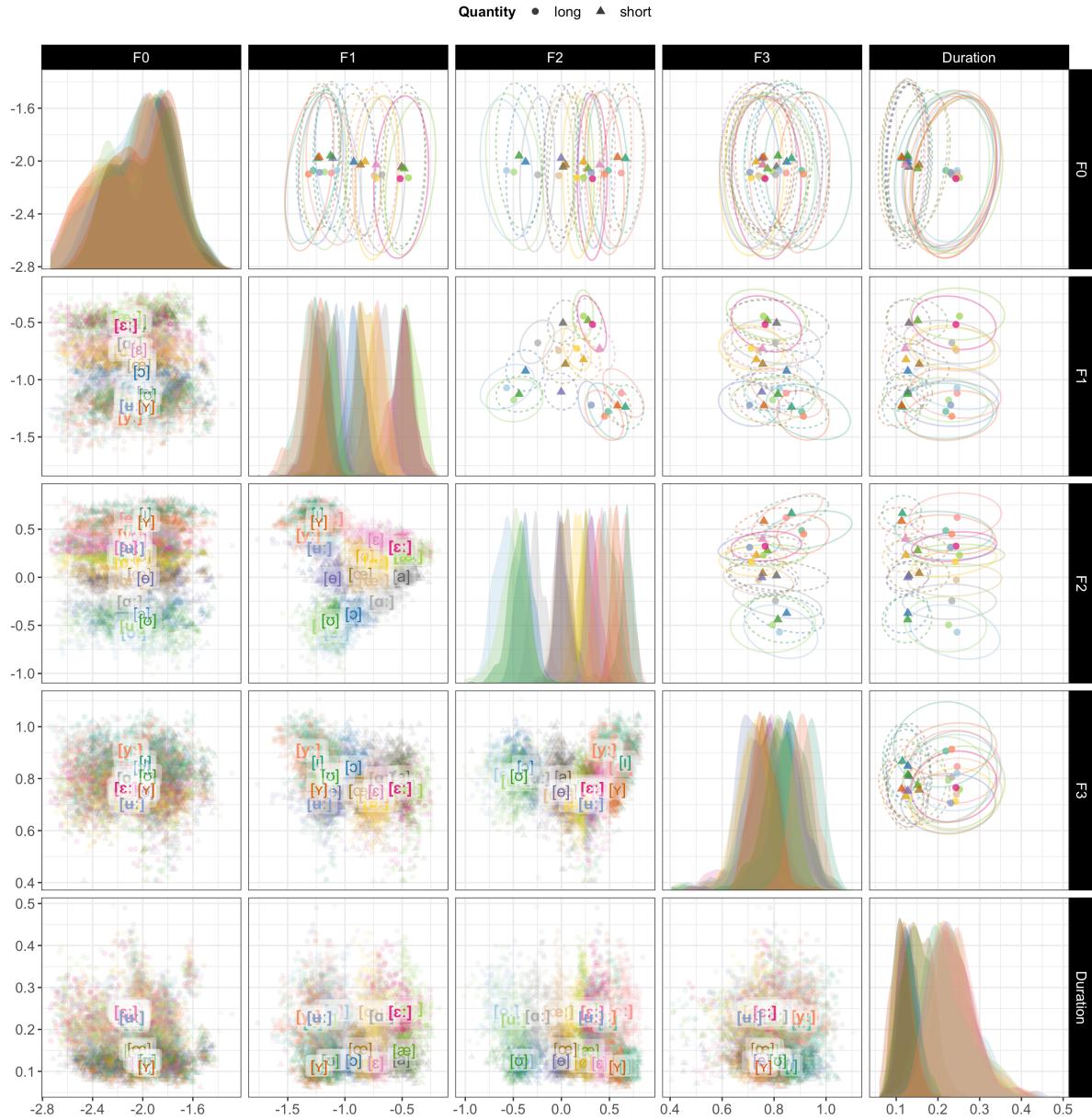
349 **2.2.1.1 Cues and cue correlations** For the pairwise combinations of the five spectral  
 350 and temporal cues—F0, F1, F2, F3 and duration, see Figure 3 from Persson and Jaeger  
 351 (2023) updated to include data from the 17 male talkers. Unsurprisingly, the densities  
 352 along the diagonal suggest that F0 carries the least information about vowel identity,  
 353 exhibiting less between-category separation than all other cues.

354 As is to be expected, vowels differing in quality are most separated in the F1-F2  
 355 panels. The F1-F3 and F3-F2 panels both display increased separation between the  
 356 neighboring outrounded [y:] and inrounded [œ:], and unrounded [ɪ] and outrounded [ʏ],  
 357 compared to when plotted along F1-F2, which points to the importance of F3 for these  
 358 vowels. Interestingly, the almost complete overlap between [i:] and [y:] in F1-F2 space  
 359 overall remains when F3 is considered, even if some individual differences in the amount of  
 360 overlap exist. Most talkers produce these two vowels very close in F1-F2 space, and only

---

male talkers have lower formant values and lower F0s than the female talkers (average F0 across long and short categories for female talkers = 203, for male talkers = 121).

<sup>8</sup> There is a possibility that the increased separability found for the short vowels is partly an artifact of how time-points for cue measurements were selected. Time-points based on percentage of vowel duration will necessarily render measurement points that are closer in time for the shorter vowels, potentially providing a better estimate of the formant value that is most distinctive, i.e., the steady state in the center of the vowel.



*Figure 3.* The SwehVd vowel data shown for all pairwise combinations of five cues: F0, F1, F2, F3 and duration. Panels on the diagonal show marginal cue densities of all five cues. The off-diagonal panels show vowel means across talkers, represented by points and with bivariate Gaussian 95% probability mass ellipses in the upper panels, and represented by vowel labels and with points for each recording in the lower panels. Note that, unlike in Figure 1, axis directions are not reversed.

<sup>361</sup> slightly separated in F2-F3 space, while others display a continued overlap when  
<sup>362</sup> considering F3 (for reference, one talker of each type are displayed in SI Figure S4). This  
<sup>363</sup> would seem to suggest that F3 might carry less importance as distinctive feature for [i:] -  
<sup>364</sup> [y:] than previously established (c.f., Fant, 1959; Fant et al., 1969).

<sup>365</sup> In order to quantitatively asses whether the distinction between closely neighboring  
<sup>366</sup> unrounded and rounded categories increased when F3 was considered, the category  
<sup>367</sup> separability of these vowels was calculated based on F1 and F2, and subsequently compared  
<sup>368</sup> against the separability calculated when including F3. If separability were to increase when  
<sup>369</sup> F3 was added, it would suggest that F3 does contribute to category distinctions.<sup>9</sup>

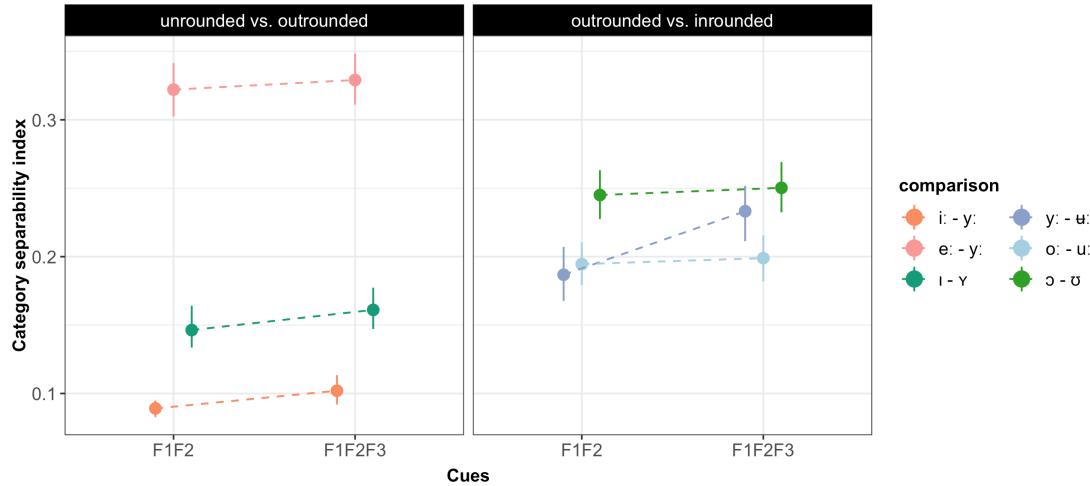
<sup>370</sup> Two general observations can be made from Figure 4. Category separability was  
<sup>371</sup> overall lower for some contrasts when only F1 and F2 were considered, e.g., the [i:] - [y:],  
<sup>372</sup> and [i] - [y] contrasts, presumably indicating their overlap in F1-F2 space. Second,  
<sup>373</sup> including F3 increased overall category separability, but only marginally for most contrasts.  
<sup>374</sup> The contrast that seem to benefit most from the inclusion of F3 is the [y:] - [u:] contrast.

<sup>375</sup> The LMMs fit to the data (presented in Section 2.1.2.3) indicated that including F3  
<sup>376</sup> improved category separability for all contrasts (all  $p < .001$ ). This suggests that the  
<sup>377</sup> subtle differences observed by visual inspection were nevertheless significant (Summary  
<sup>378</sup> tables in SI S1.2.6).

<sup>379</sup> **2.2.1.2 Quantity vs. quality in long and short vowel pairs** To gain more insight  
<sup>380</sup> into the extent to which there are spectral and temporal differences between long and short  
<sup>381</sup> vowels, the acoustics of categories within vowel pairs were evaluated. This allows for an  
<sup>382</sup> assessment of whether quantity and quality distinctions seem to be separate from each  
<sup>383</sup> other.

---

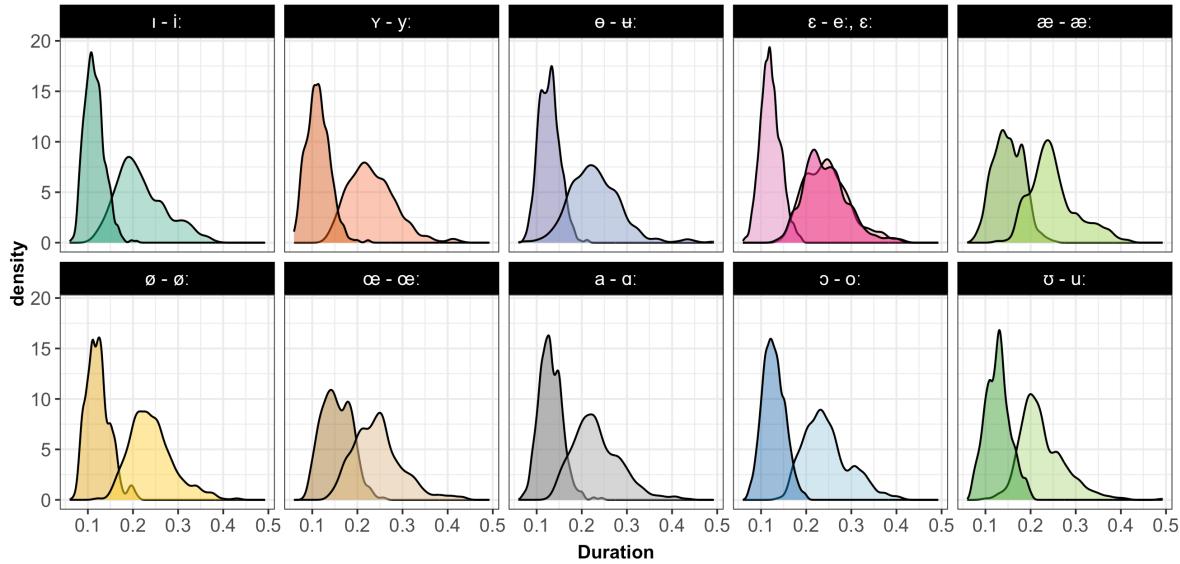
<sup>9</sup> To capture possible differences in separability between vowels in a contrast, the separability index was first calculated in both directions, that is, both assessing the separability of [i:] from [y:], and [y:] from [i:], following X. Xie and Jaeger (2020). Given that no significant differences in directionality were found, the separability index reports the average separability across the two categories in each contrast.



*Figure 4.* The effect of including F3 in measures of category separability for the distinction between neighboring unrounded vs. outrounded vowels (**left panel**), and outrounded vs. inrounded vowels (**right panel**). Pointranges indicate mean and 95% bootstrapped CIs of the category separability summarized across talkers for each cue combination.

384 As expected, long-short vowel pairs differ systematically in duration (Figure 5). For  
 385 each vowel pair, the duration densities in Figure 5 are overlapping but with two clearly  
 386 separable peaks (mean duration for the long vowels = 0.19 ms, SD = 0.10; mean duration  
 387 for the short vowels = 0.08 ms, SD = 0.09). Overall, the short vowels display less variability  
 388 in duration than the long vowels, a common pattern for measures with a lower bound.

389 All long-short vowel pairs furthermore display spectral differences in F1-F2. In fact,  
 390 as indicated in Figure 1, formant differences are apparent for *all* vowel pairs, even for vowel  
 391 distinctions for which duration has been found to be the primary cue—[ε:] - [ɛ], [ø:] - [ø], [i:]  
 392 - [ɪ], and [ɔ:] - [ɔ] (e.g., Behne et al., 1997; Hadding-Koch & Abramson, 1964; Kuronen,  
 393 2000). The vowel pairs that display larger spectral differences along F1-F2 seem to be [ɛ:] -  
 394 [ɛ] and [ɑ:] - [a] (in line with e.g., Fant, 1983; Kuronen, 2000), but also [ε:] - [ɛ], which  
 395 contrasts with previous studies. The large spectral differences in [ε:] - [ɛ] are presumably  
 396 due to [ε:] being produced very low in the SwehVd materials, which increases the distance

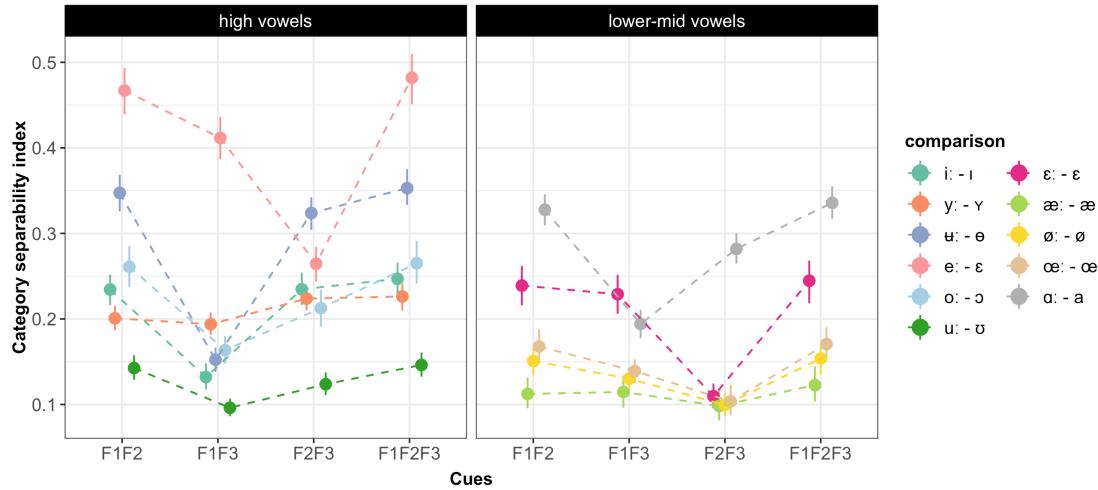


*Figure 5.* Illustrating the systematic differences in duration between the long and short vowel pairs in SwehVd.

397 to [ɛ] and in addition leads to a gap along the top-left to bottom-left diagonal between [e:]  
 398 and [ɛ]. Overall, F2 appears to carry more of the spectral variation between the long and  
 399 short vowel phonemes, as categories display increased separability in the pairwise  
 400 combination of F2 and duration (Figure 3, rightmost column, third row).

401 In order to evaluate what cue combination would provide the largest separability  
 402 between vowels in long-short contrasts, the category separability index was calculated for  
 403 each pair and four different cue combinations: F1-F2, F1-F3, F2-F3 or F1-F2-F3.

404 Figure 6 indicates that category separability is generally higher for some pairs, e.g.,  
 405 [e:] - [ɛ], [ɯ] - [ø], and [œ] - [ɑ], and that the F1-F2-F3 combination overall maximizes  
 406 separability between the pairs, although the difference to the F1-F2 space appears marginal  
 407 for most pairs. This would seem to suggest that the first two formants are the most  
 408 important cues to long-short vowel pair distinctions, while the inclusion of an additional  
 409 cue unsurprisingly does not punish the separability. For the F1-F3 and F2-F3 cue spaces,  
 410 there seems to be two main groups for which either cue combination generates the lowest



*Figure 6.* Category separability for long-short vowel pair distinctions depending on the cue combination considered. For visualization purposes, the pairs are split into high vowels (**left panel**), and lower-mid vowels (**right panel**). Pointranges indicate mean and 95% bootstrapped CIs of the category separability summarized across talkers for each cue dimension.

411 separability index. For the [u:] - [ø], [æ:] - [a], [o:] - [ɔ], [i:] - [ɪ], and [u] - [ø] vowel pairs, the  
 412 F1-F3 cue combination generates the lowest separability, which would seem to suggest the  
 413 informativity of F2 for distinguishing between these pairs. For the second group, mainly  
 414 consisting of allophones to /ε/ and /ø/: [e:] - [ɛ], [ɛ:] - [ɛ], [æ:] - [æ], [ø:] - [ø] and [œ:] - [œ],  
 415 the F2-F3 cue combination produces the lowest separability, presumably highlighting the  
 416 importance of F1. Interestingly, the [y:] - [y] pair generates almost identical separability  
 417 indices in F1-F2 and F1-F3 spaces, and increased separability for both F2-F3 and  
 418 F1-F2-F3, which seems to signal the importance of all three cues.

419 These findings were confirmed by the statistical analysis (Summary tables in SI  
 420 S1.2.6): the LMMs indicated that the F1-F2-F3 cue combination generated the highest  
 421 separability, however, for all but the [y:] - [y] ( $\hat{\beta} = .025, SE = .006, p < .0001$ ) and [e:] - [ɛ]  
 422 ( $\hat{\beta} = .015, SE = .007, p < .041$ ) vowel pairs, this change was significantly indistinguishable  
 423 from the F1-F2 cue combination (all other  $ps > .06$ ). For the other comparisons, F1-F2

424 vs. F1-F3, and F1-F2 vs. F2-F3, the LMMs indicated significant negative changes for both  
 425 comparisons for all but three vowel pairs: for the [i:] - [ɪ], there was a significant effect of  
 426 cue combination in the F1-F2 vs. F1-F3 comparison only  
 427 ( $\hat{\beta} = -1.02, SE = .007, p < .0001$ ), for the [y:] - [ʏ] and [æ:] - [æ] vowel pairs, there was a  
 428 significant effect in the F1-F2 vs. F2-F3 comparison only (for [y:] - [ʏ],  
 429  $\hat{\beta} = .023, SE = .006, p < .0002$ ; for [æ:] - [æ],  $\hat{\beta} = -.014, SE = .005, p < .011$ ).

430 To sum up, the results of the static analysis suggests that F1 and F2 are the most  
 431 important cues to vowel distinctions in Central Swedish. While visual inspection suggested  
 432 that including F3 did not substantially increase category separability for neighboring  
 433 rounded vs. unrounded contrasts, the statistical analysis found significant improvements in  
 434 separability for all contrasts. This highlights subtle but significant differences, and the  
 435 advantages of expanding empirical analyses to modeling approaches.

436 In addition, even though all long-short vowel pairs differed systematically in duration,  
 437 they also displayed considerable spectral differences, suggesting that quantity  
 438 distinctions—long vs. short vowels—are not separate from quality distinctions—high, low,  
 439 front, back vowels. The comparison of how the category separability within each pair  
 440 changed as a function of cue combination, furthermore highlighted the importance of F1  
 441 and F2, with F2 carrying much of the informativity for several pairs. The increase in  
 442 separability found when including F3 was numerical for all pairs but the [y:] - [ʏ] and [e:] -  
 443 [ɛ] pairs, indicating the importance of F3 for these two contrasts. Given that /y/ is a  
 444 rounded category, the lower F3 values found in [ʏ] relative to [y:] would seem to indicate a  
 445 relaxation of lip-rounding in [y:].

446 Given that the category separability index assigns equal weight to all cues included,  
 447 there is no direct way of knowing which cue contributes more to the separability index.  
 448 Furthermore, similar to other evaluations presented in this subsection, the separability  
 449 index cannot account for the fact that formants are not static but rather fluctuate across  
 450 the signal. A more holistic mapping of the acoustics should therefore aim to assess how

451 formant dynamics contribute to vowel distinctions. The next section investigates how  
 452 formants move across the segment and how much information is gained by accounting for  
 453 this dynamics.

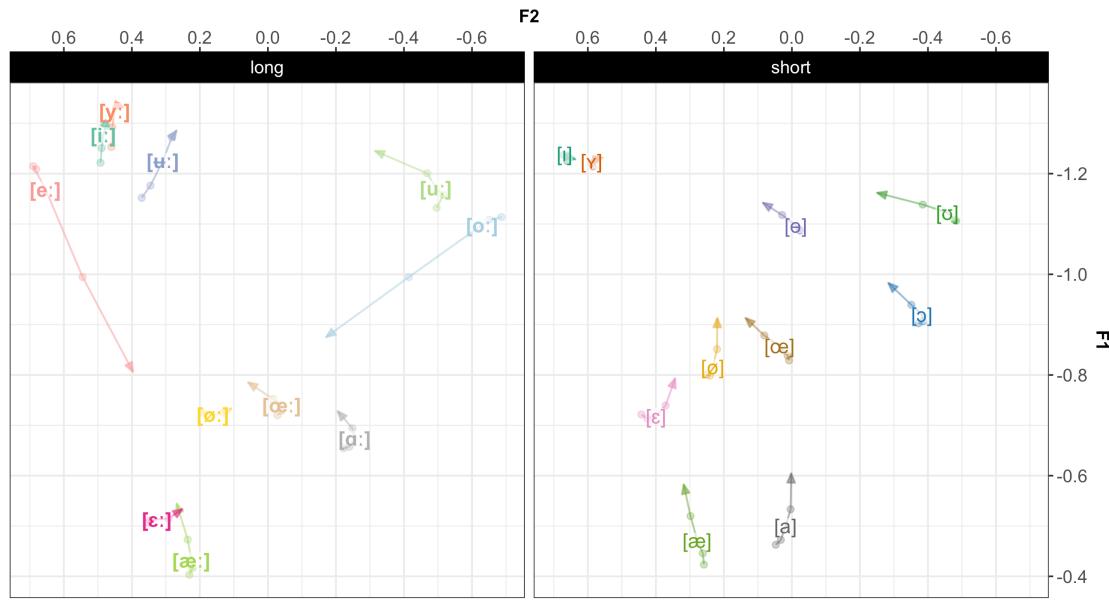
454 **2.2.2 Dynamic spectral analysis**

455 This subsection begins with visualizations of the empirical data in formant trajectory plots.  
 456 Next, the results of the GAMMs fit to the data are presented, first focusing on the  
 457 movements in eight sets of neighboring vowel contrasts, then on spectral differences  
 458 between the long and short vowel pairs.

459 **2.2.2.1 Formant movements across the space** Figure 7 displays the formant  
 460 trajectories across all 5 time-points for the long and the short vowels. In almost all vowels,  
 461 long or short, formants showed a dynamic pattern. Only [ø:], [i], and [y] showed little  
 462 movement over the measurement points. The scope and direction, however, vary. Across  
 463 vowels, the scope of movements appear to be larger moving from vowel mid-point to 80%  
 464 into the vowel, as indicated by the length of the line from vowel label to end of arrow.

465 Most of the formant dynamics thus take place *after* vowel mid-point. The long high front  
 466 vowels are important exceptions—the dynamics in [i:] and [y:] mostly occur at the  
 467 beginning of the vowel segment (between 20 and 50% into the vowel), whereas [u:] displays  
 468 movements of almost equal magnitude across the first four time-points. The largest  
 469 movements overall seem to concern [e:] and [o:].

470 In terms of directionality, there is a general tendency to move towards the centre in  
 471 most of the vowels, both long and short. According to previous studies (e.g., Bleckert, 1987;  
 472 Elert, 2000), the high vowels [i:], [y:], [u:] and [u:] tend to be realized with an offglide, which  
 473 would generate a falling F1 for all four vowels, a rising F2 for [i:] and [y:] and a falling F2  
 474 for [u:] and [u:]. These predictions were borne out for F1 in all cases, but for F2, only for  
 475 [u:]. Both [i:] and [y:] display very little movement along F2, whereas [u:] moves towards a



*Figure 7.* The trajectory of all vowels across the five time-points, along F1-F2. The arrow indicates the direction of the trajectory and ends at the final time-point, at 80% into the vowel. The vowel label is placed at the third time-point, at vowel mid-point (50%). The first (20%), second (35%) and forth (65%) time-points are represented by points.

476 more central quality, possibly indicating diphthongization ending in [ə] rather than a  
 477 consonantal offglide.

478 Parts of the movements could be due to coarticulatory effects in anticipation of the  
 479 upcoming coda ([d], [d̄], [r]). If so, one would expect F2 to centralize in the later part of  
 480 the segment, as tongue movements mark transitions into the alveolar (e.g., Hillenbrand,  
 481 Clark, & Nearey, 2001; K. N. Stevens & House, 1963). The formant movements along F2  
 482 from the last point (65%) to arrow tip (80%) in e.g., [ə], [œ:], [œ], [ɑ], [ɔ], [u:], and [i:],  
 483 might at least partly be caused by such coarticulation. Given the scope and direction of  
 484 movements, Figure 7 suggests diphthongization in primarily [e:], [u:] and [o:], replicating  
 485 previous work (e.g., Eklund & Traunmüller, 1997; Elert, 2000; Pelzer & Boersma, 2019),  
 486 while the other vowels appear to merely display formant movement, that partly could be  
 487 caused by e.g., coarticulation. The previously reported diphthongization in [ø:], however,  
 488 does not seem to be particularly pronounced in these data.

489 Figure 7 further demonstrates that some neighbouring categories either converge at

490 end points or diverge at end points. For instance, [u:] - [o:] are fairly closely located at

491 earlier time-points, but differ substantially towards the end of the vowel segment, while

492 [ɛ:]-[æ:], [ø:]-[œ:] and [ø]-[œ] start at different locations but end up in approximately the

493 same (c.f., Kuronen, 2000). Finally, the formant trajectories suggest that the empty spots

494 identified in the vowel space under a static analysis (Figure 1), may indeed be occupied

495 when vowel dynamics are considered. This is especially true for [e:] that travels from the

496 mid-high front to the mid center of the space as the signal unfolds, down to a position

497 closer to its short counterpart, [ɛ]. Given the amount of overlap when static spectral cues

498 are considered (Figure 2), formant dynamics are likely highly informative for several of

499 these distinctions. Figure 8 parallels Figure 2 and illustrates the effect of considering

500 formant movements for neighboring categories. As visualized in Figure 8, the overlap

501 between [u:] - [o:] is substantially reduced at the later time-points, while [ɛ:]-[æ:], [ø:]-[œ:]

502 and [ø]-[œ] are most distinguishable at earlier time-points.

### 503 2.2.2.2 Models of formant dynamics

#### 504 2.2.2.2.1 The effect of modeling formant dynamics for neighboring

505 **contrasts** The first set of GAMMs were modeled separately for each cue (F1, F2, F3)

506 and each of the sets of neighboring vowels hypothesized to (at least for some talkers) rely

507 on formant dynamics (Fant, 1971; Kuronen, 2000; Pelzer & Boersma, 2019): the high front

508 vowel contrasts [i:] - [y:] - [ɛ:] - [e:] and its short counterpart [i] - [y], the high back vowel

509 contrast [o:] - [u:] and its short counterpart [ɔ] - [ʊ], the lower-mid front contrast [ɛ:] - [æ:]

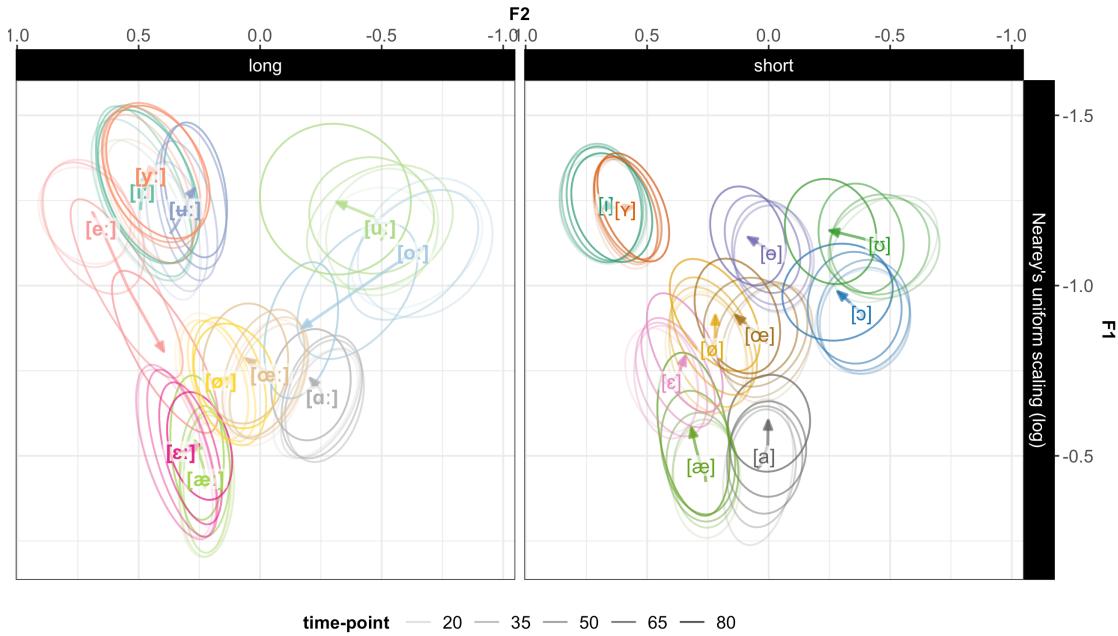
510 and short [ɛ] - [æ], and the mid center [ø:] - [œ:] and [ø] - [œ]. Summary tables of models are

511 included in the SI, Section S1.2.7.

512 The GAMMs fit to the high front vowel contrasts suggested an effect of vowel on all

513 cues for all vowels (all  $ps < .0025$ ), except for [y:] predicting F3

514 ( $\hat{\beta} = .0007, SE = .013, p > .95$ ), and for [y] predicting F1 ( $\hat{\beta} = .01, SE = .01, p > .32$ ).



*Figure 8.* Vowel placement in F1-F2 space at each of the five time-points. Ellipses show bivariate Gaussian 95% confidence interval of vowel means at each of the five time-points. Transparency indicates time-point, more transparent ellipses for earlier times. The vowel label is placed at vowel mid-point (50%). The arrow indicates the direction of the formant trajectory and ends at the final time-point, at 80% into the vowel.

515 The long [y:] predicting F2 was also close to insignificant ( $p > .046$ ). This indicates that  
 516 F3-dynamics does not aid in distinguishing between [i:] and [y:], and neither does  
 517 F1-dynamics for [i] and [y]. These results would thus overall seem to suggest that while F3  
 518 increases separability between [i:] and [y:] under static analysis, the *dynamics* in F3 does  
 519 not seem to add anything to the contrast. For the [i] - [y] contrast, F3 is likely more  
 520 informative, while the F1-dynamics does not contribute to distinguishing between the two.  
 521 Figure 9 panels A and E demonstrate the importance of formant dynamics for the high  
 522 front long vowels in F1 and F3, as many of the vowels are overlapping at some time-points,  
 523 but never along the entire segment.

524 An effect of category on all three cues was found in all GAMMs fit to the high back

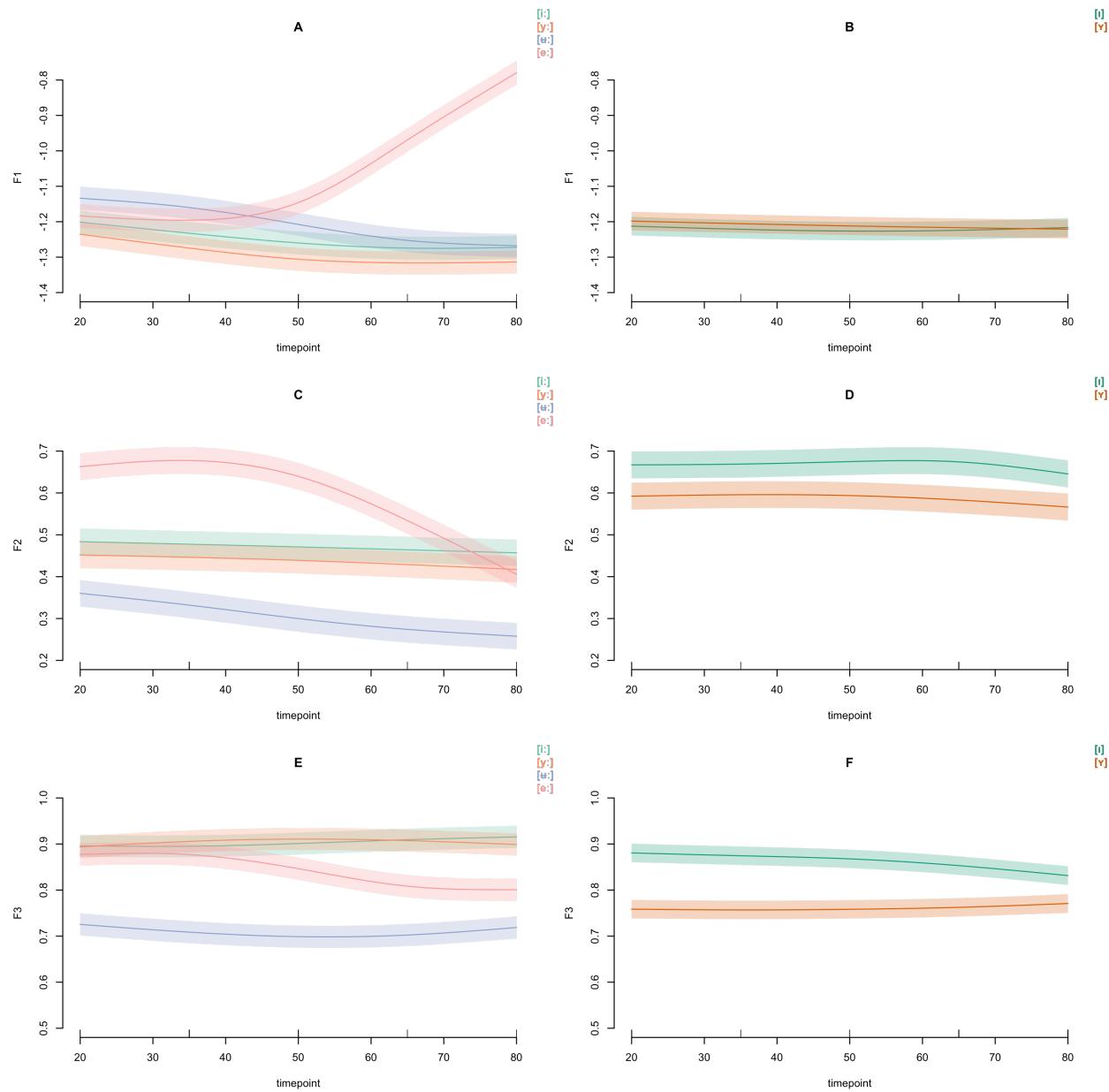


Figure 9. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the long vowels [i:], [y:], [u:], and [e:] (**left**), and the short vowels [i] and [y] (**right**).

525 long and short contrasts (all  $ps < .0006$ ). These results suggest that formant dynamics are  
 526 likely informative for these contrasts, and that all three cues contribute to distinguishing  
 527 between these vowels also when formant dynamics are considered (Figure 10).<sup>10</sup> Figures 9  
 528 and 10 further highlight the presumable diphthongization in [e:] and [o:], as indicated by  
 529 the steepness of the fitted curve.

530 The GAMMs fit to the lower-mid front long and short vowel contrasts suggested an  
 531 effect of vowel on F1 and F2 (all  $ps < .0001$ ), but not on F3 (for [æ:],  
 532  $\hat{\beta} = -.012, SE = .01, p > .25$ ; for [æ],  $\hat{\beta} = .012, SE = .008, p > .14$ ). Figure 11  
 533 demonstrates how these vowels overlap in F3-dynamics, but are distinguished for larger  
 534 parts of the segment along F1 and F2.

535 Finally, for the GAMMs fitted to the mid center vowels, there was an effect of  
 536 category on all cue evaluations for both long and short vowels, with the exception of the  
 537 long vowels fit to F1 (for [ø:] - [œ:],  $\hat{\beta} = -.017, SE = .01, p > .12$ ). For several of the other  
 538 evaluations, the effect was relatively small, however significant (all  $ps < .044$ ; Figure 12).  
 539 These results would seem to suggest that when formant dynamics are considered, the [ø:] -  
 540 [œ:] contrast is primarily upheld by F2 and F3.

541 The sets of neighboring contrasts investigated here all exhibited varying degrees of  
 542 category overlap in static analysis. However, when formant dynamics was considered, the  
 543 vowels in each contrast were all significantly different from each other along at least two  
 544 cues (c.f., Fant, 1971; Kuronen, 2000; Pelzer & Boersma, 2019). For all contrasts, the  
 545 vowels overlapped in parts of the segment, but importantly never along the entire segment  
 546 (Figures 9A-E, 10A-D-E, 11A, 12A-B-E-F). This indicates that category overlap found  
 547 in static analysis is mitigated once temporal analysis is included, which suggests that  
 548 category distinctions unfold over time.

---

<sup>10</sup> An effect of gender was found for the high back vowels predicting F1 for the long and short vowels (both  $ps < .022$ ), which suggests that the normalization approach likely reduced some talker-specific differences related to anatomical differences, but not all.

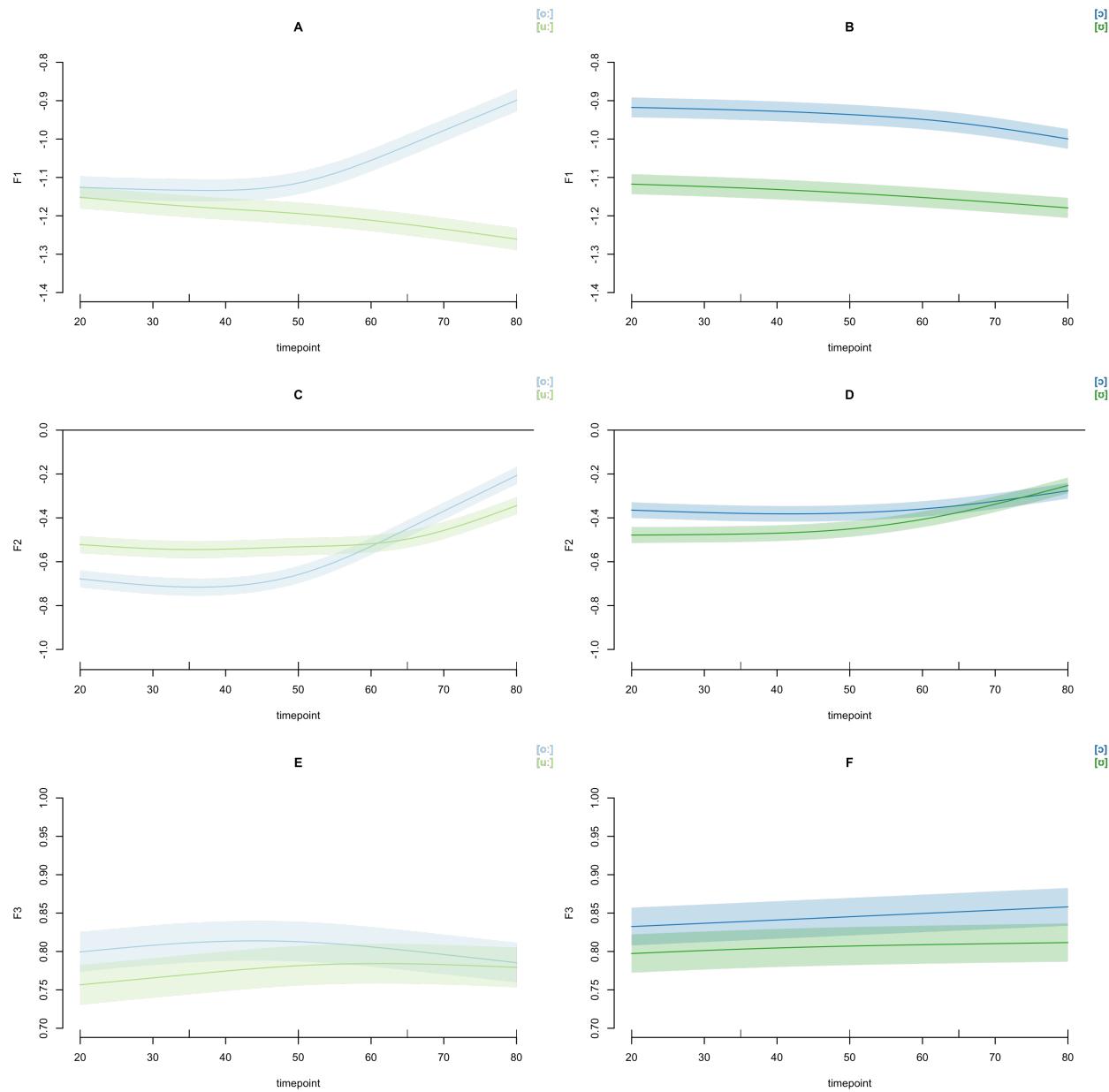


Figure 10. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the long vowels [o:] - [u:] (**left**), and the short vowels [ɔ] - [ʊ] (**right**).

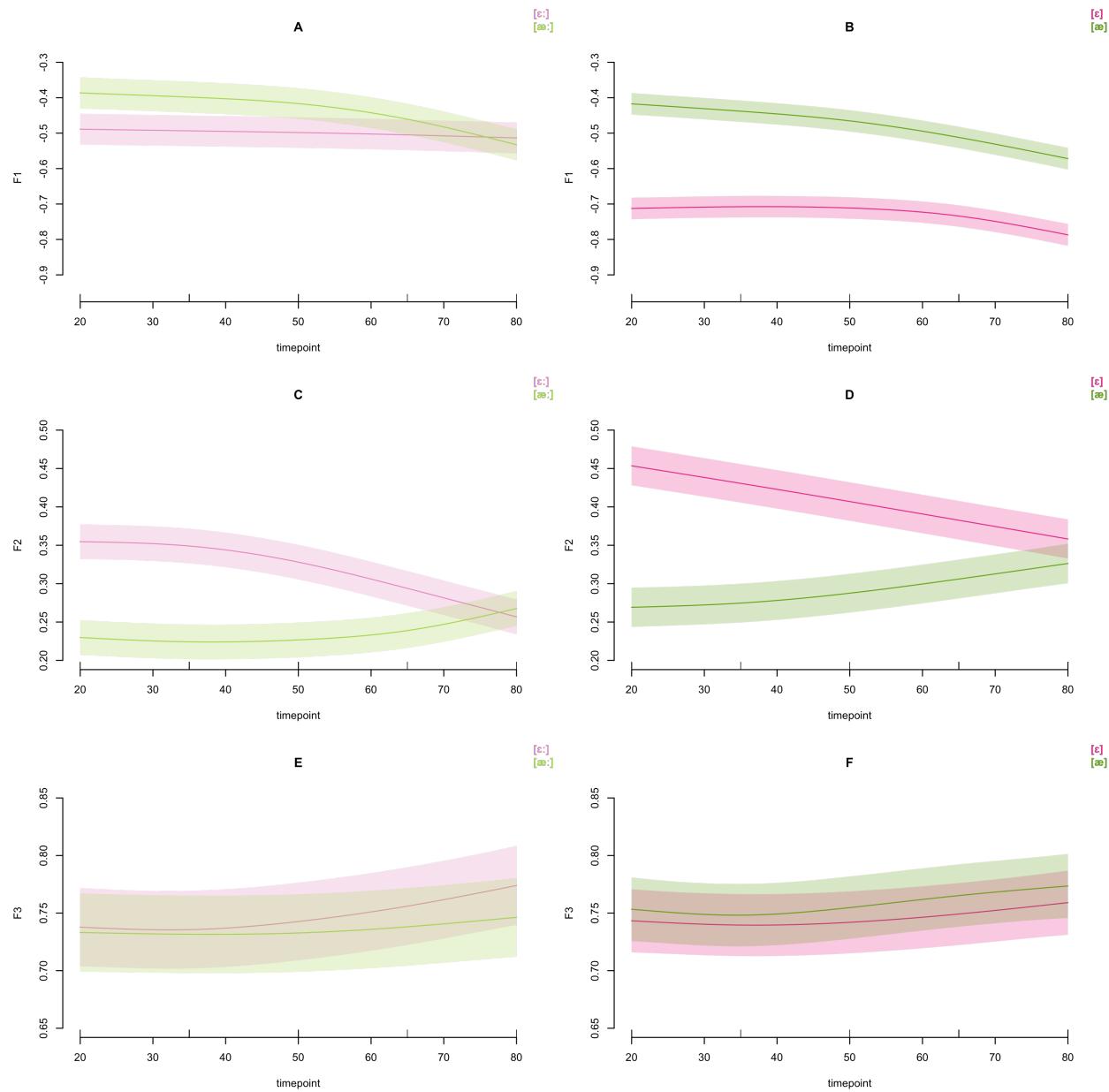


Figure 11. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the long vowels [ɛ:] - [æ:] (**left**), and the short vowels [ɛ] - [æ] (**right**).

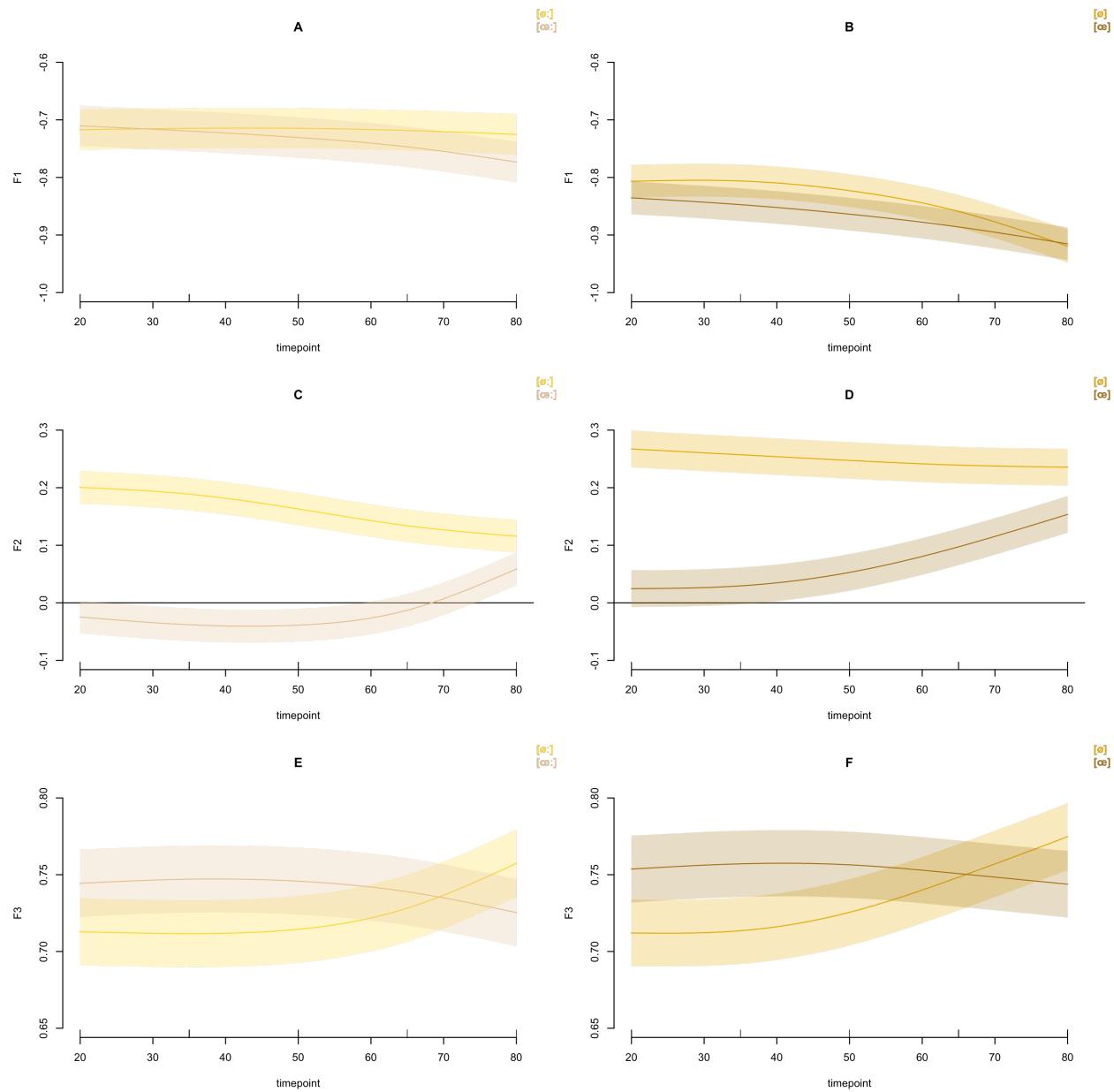


Figure 12. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the long vowels [ø:] - [œ:] (**left**), and the short vowels [ø] - [œ] (**right**).

549 The results further indicate that F3-dynamics carry little information for the [i:] - [y:]

550 contrast. For the short [i] - [y], however, it is F1 rather than F3 that does not distinguish,

551 suggesting overlap in F1 (height), but not F3 (lip-rounding).

#### 552 2.2.2.2.2 Formant dynamics in long-short vowel pairs The second set of

553 GAMMs modeled the effect of F1, F2, F3 and duration on all long and short vowel pairs.

554 Summary tables and visualizations of these GAMMs are included in the SI (Section S1.2.7).

555 There was a treatment effect of vowel on all spectral and temporal cues for all vowel pairs

556 (for F1, all  $ps < .04$ ; for F2 and duration, all  $ps < .0001$ , for F3, all  $ps < .0005$ ). With the

557 exception of the [y:] - [y] vowel pair, all rounded vowels displayed lower F3 values, hence

558 more lip-rounding, in their long allophones (e.g., Hadding, Hirose, & Harris, 1976;

559 Stålhammar, Karlsson, & Fant, 1973). The results suggest that when formant dynamics

560 are considered, *all* long-short Central Swedish vowel pairs differ in spectral and temporal

561 cues. Among the pairs that displayed smaller, albeit statistically significant, differences in

562 spectral cues are the [i:] - [i] pair predicting F1 and F3, possibly indicating a tendency for

563 stronger duration dependency, in line with previous perceptual work (Behne et al., 1997).

564 While the static analysis suggested that F1-F2-F3 only numerically increased separability

565 for long-short pairs, with the exception of [y:] - [y] and [e:] - [ɛ], the dynamic analysis

566 suggests that all cues, when considered separately, carry information about vowel quality in

567 quantity contrasts under the assumption of formant dynamics.

#### 568 2.2.3 Results summary

569 The results from the static and dynamic analyses in Study 1 suggest that F1, F2, F3 and

570 duration all contribute to the distinction of Central Swedish vowels. While F1 and F2 are

571 or primary importance for most contrasts, F3 does contribute to increasing the separability

572 between some neighboring vowels differing in lip-rounding. The static analysis furthermore

573 suggested that the inclusion of F3 had less importance for the long-short vowel contrasts,

574 as statistically indistinguishable performance was found for the F1-F2 model.

575 Some vowels displayed overlap in the static analysis but increased separability when  
 576 formant dynamics was considered, as indicated by formant trajectory analysis and  
 577 GAMMs. The dynamic analyses highlighted that the short vowels also display formant  
 578 movements, and that for most of both long and short categories, a larger portion of the  
 579 dynamics resides in the later part of the segment. Given the increased separability of  
 580 neighboring contrasts found in dynamic analysis, it is reasonable to assume formant  
 581 movements as an auxiliary cue to vowel identity, more so for some contrasts than others.  
 582 For instance, the [i:] - [y:] - [ɯ:] - [e:], [ɪ] - [ʏ], [o:] - [u:], [ɔ] - [ʊ], [ɛ:] - [æ:], [ɛ] - [æ], [ø:] - [œ:],  
 583 and [ø] - [œ] contrasts displayed considerable overlap in static analyses but increased  
 584 distinguishability when analysed dynamically.

585 While the static analysis suggested increased separability for the [i:] - [y:] contrast  
 586 when F3 was included, this was not found under the assumption of formant dynamics. The  
 587 GAMM fit to [i:] - [y:] - [ɯ:] - [e:] contrast found no effect of [y:] relative to [i:] predicting F3.  
 588 However, the GAMM fit to [y:] - [ʏ] suggested statistically significant differences between  
 589 the two categories predicting F3. These two analyses taken together suggest more effects  
 590 on F3 in the short vowel compared to the long vowel.

Table 2

*The phonetic characterization of long (left) and short (right) Central Swedish vowels (as represented in the SwehVd database). Rounded vowels are shaded.*

	front	central			back		front	central	back
high		[i:]	[y:]	[ɯ:]	[u:]		[ɪ]	[ʏ]	[ʊ]
mid-high	[e:]				[o:]		[ɛ]	[ø]	[œ]
lower-mid	[æ:]		[ø:]	[œ:]	[ɑ:]		[æ]	[ɑ]	

591 The resulting phonetic characteristics of the long and short Central Swedish vowels  
 592 presented in Study 1, is summarized in Table 2. Beginning with the long vowels, there are  
 593 4 high vowels. The current acoustic description suggests that none of them are front.  
 594 Instead, [i:] and [y:] group with [ɯ:] as central vowels, and [u:] is back (c.f., Pelzer &

<sup>595</sup> Boersma, 2019; Schötz et al., 2011).<sup>11</sup> There are 2 mid-high vowels—[e:] (front) and [o:] (back)—and 4 lower-mid vowels, [æ:] (front), [ø:], [œ:] (both central) and [ɑ:] (back). Given the substantial lowering of [ɛ:] and its overlap with [æ], it is reasonable to assume one long allophone for /ɛ/, which is [æ:] (c.f., Pelzer & Boersma, 2019).

<sup>599</sup> The short vowel space contains 4 high vowels, two of which are front, [i] and [y], one <sup>600</sup> is central [ø] and one back [ʊ]. There are 4 mid-high vowels, [ɛ] (front), [ø], [œ] (both <sup>601</sup> central), and [ɔ] (back), and 2 low vowels, [æ] (front), and [a] (central). The analysis of this <sup>602</sup> database supports what Riad (2014) anticipated and Pelzer and Boersma (2019) suggested, <sup>603</sup> namely, a vowel system consisting of three height levels only, in contrast to the traditional <sup>604</sup> four height levels system (e.g., Engstrand, 1999, 2004; Riad, 2014).

<sup>605</sup> The motivation of the summarized acoustics presented in Table 2 rests on a pairwise <sup>606</sup> grouping of the long and short vowels, similar to phonological analyses of Central Swedish <sup>607</sup> (Riad, 2014). For instance, despite its high position, [e:] is defined as mid-high, on par with <sup>608</sup> [o:], as both vowels share diphthongizational patterns, and their short versions are both <sup>609</sup> lower than their long counterparts. Furthermore, [æ:] is front rather than central as its <sup>610</sup> short version is clearly more front than [ø]. Because of their overall centralized positions in <sup>611</sup> SwehVd, [i:] and [y:] groups with [œ:], while their short versions are still clearly front. One <sup>612</sup> could thus argue that both /i/ and /y/ are under-specified for the front-back dimension, <sup>613</sup> which would also be the case for /ɑ/, as [ɑ:] is back and [a] is central. It is important to <sup>614</sup> note that while Table 2 may point to possible updates of Central Swedish vowel phonology, <sup>615</sup> this is only tentative as more evidence is required for a definite update. These include, e.g., <sup>616</sup> more investigations in different contexts.

<sup>617</sup> The phonetic characterization in Table 2 suggests that some categories have shifted <sup>618</sup> in comparison to previous work. In order to assess the scope and spread of these changes, <sup>619</sup> Study 2 was conducted.

---

<sup>11</sup> Whether [y:] can still be considered rounded given the high F3 values, is a question for future research.

## 620 3 Study 2: Vowel category movements over time -

### 621 from SweDia to SwehVd

622 The aim of Study 2 was to investigate possible vowel shifts over the last generation by  
 623 comparing the acoustic characteristics of SwehVd against that of 8 reference talkers in the  
 624 SweDia materials, along the two primary cues to vowel identity, F1 and F2. Next, the  
 625 reference materials used for comparison is presented, alongside the distance metric used for  
 626 assessing potential vowel shifts.

### 627 3.1 Methods

#### 628 3.1.1 Materials

629 The materials used for comparison is a database of  $N = 8$  L1 talkers ( $N = 4$  female) of  
 630 Central Swedish that were recorded as reference talkers of Standard Swedish for the  
 631 SweDia dialect database (Eriksson, 2004). All talkers were L1 talkers of Swedish, born and  
 632 raised in the Greater Stockholm area or surroundings (provinces of Södermanland,  
 633 Uppland), and of 22-35 years of age (mean age = 27; SD = 4.41). The talkers were  
 634 recorded reading a list of 41 words containing all 21 vowels of Central Swedish (for the  
 635 complete wordlist used for recording, see Table S33 in SI). Thirty-nine words on the list  
 636 were monosyllabic, and two were bi-syllabic (*leta*, *flytta*). The target vowels extracted ([e:],  
 637 [y]) from the bisyllabic words had primary stress and were thus included in this study.

638 All [æ] vowels were excluded from the subsetted materials used here, as this vowel  
 639 was recorded by only two of the talkers. All repetitions recorded by a talker were included,  
 640 even though the variability in the number of repetitions for each vowel across talkers was  
 641 substantial.<sup>12</sup> The phonological contexts used were the same across talkers (with two  
 642 exceptions, see Table S33 in SI) but differed across categories. Importantly, none of the

---

<sup>12</sup> Three to 8 repetitions of [i:], [y:], [ɯ], [ø], [e:], [ɛ:], [œ:], [ø], [œ], [ɑ:], [ɑ], [ɔ:], [ʊ], [ʊ], 6 repetitions of [ɪ], [ʏ], [ɛ], 3 to 14 repetitions of [ø:], [ɔ], [ɛ] per talker.

643 contexts were hVd. For details on the recruitment, recording, pre-processing and  
644 annotation procedure in the SweDia materials, see Eriksson (2004).

645 The two databases were of unequal sample sizes overall: N = 9103 datapoints for  
646 SwehVd, and N = 669 datapoints for SweDia. The regional origin of the talkers matched  
647 across databases, age at the time of recording (mean age in SweDia = 27, mean age in  
648 SwehVd = 29) and the relative proportion of male and female talkers were also highly  
649 similar (50% female talkers in SweDia2000, 55% female talkers in SwehVd). Any potential  
650 effect of differences in gender distributions on the comparison were presumably reduced by  
651 the use of normalized vowel data.

652 Cue measurements from SweDia were extracted in the same way as for SwehVd  
653 (Section 2.1.2.1). To correct for measurement errors in the automatic extraction of cues, 5  
654 separate univariate distributions of the five extracted cues (F0, F1, F2, F3 and duration)  
655 was estimated for each distinct combination of talker and vowel. Points that fell outside of  
656 the 0.25th to the 97.5th quantile of the distributions for each vowel were identified,  
657 examined for measurement errors, and subsequently corrected. This followed the approach  
658 employed for the SwehVd corpus (Persson & Jaeger, 2023), and strikes a middle-ground  
659 between the ideal (manual correction of all tokens) and feasibility. Outliers were identified  
660 and removed following the approach for the subsetted SwehVd dataset used in Study 1,  
661 leaving N = 664 datapoints for analysis.

### 662 3.1.2 Assessing vowel shifts

663 In order to assess whether the relative placement of the categories in the acoustic space was  
664 different between datasets, formant data was visualized in F1-F2 space and distances  
665 assessed with the *orthogonal projection ratio*, henceforth *op* (M. Stevens, Harrington, &  
666 Schiel, 2019).<sup>13</sup> Both visualizations and *op* calculations were based on measurements at the

<sup>13</sup> Other distance measures are possible, such as Mahalanobi's distance. The primary advantages of using the *op* is that it considers the overall structure of vowel spaces rather than individual distances, and that it

667 steady-state portion of the vowel in both datasets—taking the average F1-F2 values across  
 668 the three mid-points (at 35, 50, and 65% into the vowel). This approach presumably  
 669 reduced the risk of coarticulatory effects resulting from the use of different phonological  
 670 contexts in the datasets.

671 **3.1.2.1 Orthogonal projection ratio** The *op*, formalized in (2), compares the  
 672 position of a single token of a  $V_{vowel}$  from a talker, relative to the means of two anchor  
 673 categories,  $U_{anchor}$  and  $W_{anchor}$ , from the same talker. The anchor categories are selected  
 674 based on their position in the space relatively to  $V_{vowel}$  (for an illustration, see SI Figure  
 675 ??). After calculating the *op* for each token of  $V_{vowel}$  from each talker relative to the  
 676 anchors, the *op* is summarized across all tokens for that talker. This procedure is repeated  
 677 for each of the vowels and their respective anchors for each talker in each dataset and then  
 678 compared across datasets, in order to assess vowel shifts over time.

$$op(V_{vowel}) = (-2 \cdot \frac{((V_{vowel} - U_{anchor})\% * \% (U_{anchor} - W_{anchor}))}{((U_{anchor} - W_{anchor})\% * \% (U_{anchor} - W_{anchor}))}) - 1 \quad (2)$$

679 An *op* = 0 indicates that the vowel  $V_{vowel}$  is positioned exactly in the middle between  
 680  $U_{anchor}$  and  $W_{anchor}$ , an *op* = ±1 indicates that  $V_{vowel}$  is positioned exactly on the mean  
 681 of one of the two anchors, whereas a  $V_{vowel}$  positioned somewhere between  $U_{anchor}$  and  
 682  $W_{anchor}$  has an *op* that varies between ±1. The *op* can thus be used as an index of e.g.,  
 683 lowering, fronting, or centralizing.

684 Table 3 lists all the vowels  $V_{vowel}$  tested along with their corresponding anchor points,  
 685  $U_{anchor}$  and  $W_{anchor}$ . These were all vowels for which the overall distance in F1 and F2  
 686 vowel means in SwehVd to SweDia,  $\geq .85$  standard deviations from the vowel means in

---

offers a standardized way to compare the degree of category membership across different datasets.

687 SweDia.<sup>14</sup>

Table 3

*The anchors used for calculating orthogonal projection ratio across the two datasets*

$U_{anchor}$	$V_{vowel}$	$W_{anchor}$
[ɑ:]	[y]	[ɪ]
[ɪ]	[ɯ]	[o:]
[e:]	[ɛ:]	[æ:]
[ε:]	[ø:]	[a:]
[ɯ:]	[ø]	[a]
[ø]	[œ]	[œ:]
[ɑ:]	[ɔ]	[o:]

688 To assess the magnitude of differences in  $op$  between datasets, separate LMMs for  
 689 each of the categories investigated were fit, predicting  $op$  from database, with talkers as  
 690 random effect:  $op \sim Database + (1|Talker)$ . Database was treatment coded with SweDia  
 691 as reference level.<sup>15</sup>

## 692 3.2 Results

693 Figure 13 displays the vowel movements in F1-F2 space across talkers over time. The figure  
 694 suggests that a few vowels display considerable stability over time, e.g., the back vowels [ɑ:]  
 695 and [o:], the center vowels [œ:] and [a], and the front vowels [e:] and [ɛ]. The remaining  
 696 vowels have all moved to a lesser or greater extent. Two general tendencies can be noted.  
 697 Firstly, several of the vowels that form the edges of the space slightly shift further  
 698 outwards, causing the overall space to widen. This is the case for [ɪ], [æ:], [o:], and [ɯ].  
 699 Most of the inner vowels either also shift outwards—[y], [ø], [ɯ], [ʊ], [ɔ]—or centralize—[ø],

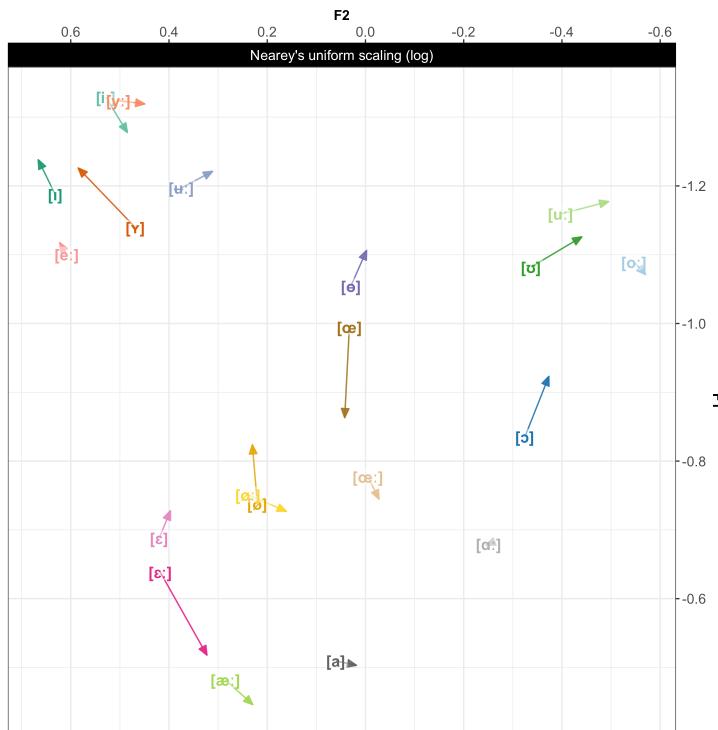
---

<sup>14</sup> The front [æ:] was also included in this subset. Given that it occupies the lowest position in the space, the calculation of  $op$  was not possible.

<sup>15</sup> The two databases employed different phonological context. One way to control for this difference would be to use random effects by word. However, this was not possible since this random effect would be highly collinear with the database, as SwehVd only used one word context. Indeed, most effects were no longer significant when word was included. Furthermore, it is reasonable to assume that coarticulatory effects from using different phonological contexts would be minimal given that the measurements were extracted from the steady-state of the vowel. Therefore, excluding word as random effect likely had very little effect on the results.

700 [ø], [œ]. Second, these movements lead to some neighboring categories closing in on each  
 701 other, e.g., [ɪ] - [Y], [i:] - [ɪ:], [ø] - [œ], [o:] - [ɔ] - [u:] - [u], or increasing their separability—[e]  
 702 - [œ], [ø] - [ø:], [Y] - [ɪ:].

Figure 13 further suggests a series of consecutive movements in the front part of the space, where [e:], [i], [y] have traveled further up and front in SwehVd, while [i:], [y:] and [u:] have centralized. Finally, Figure 13 shows that [ɛ:] have lowered substantially compared to SweDia, down to a position close to [æ:] (and to [æ]; recall Figure 1). The already low [æ:] has in turn moved further down, presumably constrained by the lower bound and cannot move further.



*Figure 13.* Vowel category movements in the Central Swedish vowel space from 1999 to 2023, in F1-F2 space. The arrow indicates the direction of movement over time, with base of arrow representing SweDia (vowel label) and tip of arrow representing SvehVd. Vowel data is taken at the steady state of the vowel (averaged across the three mid-points) and averaged across talkers in the two datasets.

To quantify these potential vowel shifts, the distances between the category locations in the two datasets were assessed with the *op*. Figure 14 visualizes the *op* of all vowels in

711 Table 3 in both datasets. By this way of assessing the relative movement between anchor  
 712 points across datasets, Figure 14 indicates that all selected categories have moved over  
 713 time in the directions suggested by the qualitative analysis in Figure 13. The largest  
 714 differences appear to concern [y], [ɛ:], [œ], and [ɔ]. The front [y] is relatively closer to [i] in  
 715 SwehVd than in SweDia, which suggest fronting. The [ɛ:] is relatively closer to [æ:] in  
 716 SwehVd than in SweDia, as is [œ] to [œ:], indicating the lowering of both categories. The  
 717 back [ɔ] is relatively closer to [o:] in SwehVd, indicating back raising.

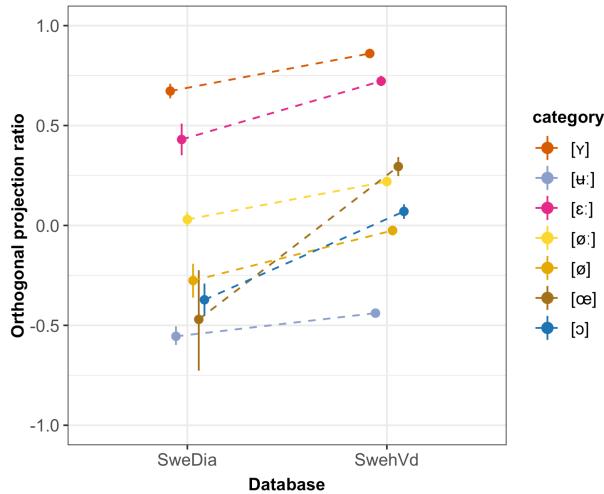


Figure 14. The relative distance between categories in SweDia and SwehVd.  $op = 0$  indicates that the vowel is equidistant between the two anchor categories in Table 3, positive numbers indicate a position closer to the second anchor vowel, negative numbers indicate a position closer to the first anchor. Pointranges indicate mean and 95% bootstrapped CIs of the  $op$  summarized across talkers in each database. The large CIs for the [œ] vowel in the SweDia materials likely reflect the tendency of some talkers to neutralize the short /ø/ and /u:/ as [œ].

718 The LMMs fit to the data confirmed the directionality of the vowel shifts, and found  
 719 a main effect of Database on  $op$  for all categories (all  $p < .015$ ), with the exception of [u:]  
 720 ( $\hat{\beta} = 0.11, SE = 0.06, p > .07$ ). Summary tables of all 6 models are included in the SI,  
 721 Section S1.3.4.

722 A discussion of these results, and the results from Study 1, follows next.

## 723 4 General discussion

724 The objective of the present paper was to provide an extensive acoustic analysis of  
725 contemporary Central Swedish vowels and to evaluate the extent of vowel shifts over the  
726 last generation. The main findings from both studies are next discussed, alongside  
727 methodological considerations and future directions.

### 728 4.1 Static and dynamic analyses of the 21 vowels of modern-day 729 Central Swedish

730 The purpose of the first study was to present up-to-date static and dynamic acoustic  
731 analyses of Central Swedish vowels that included both empirical formant data and models  
732 of formant dynamics. The first study thus aimed to expand on and complement previous  
733 work by 1) the scope of the analysis, performing the same type of analyses on all 21  
734 categories, 2) the materials chosen, using a recently collected hVd corpus with high  
735 resolution within and across talkers for a single variety, and 3) the methodological  
736 approach employed, with the use of traditional formant analysis, category separability  
737 index, trajectory visualizations and models of formant dynamics (GAMMs). Study 1  
738 furthermore aimed to evaluate the hypothesized importance of F3 for contrasting rounded  
739 vs. unrounded categories, the extent to which *all* long-short vowel pairs display spectral  
740 differences, and what part of the vowel space is more susceptible to diphthongization.

741 Beginning with the static analysis, the results seem to suggest that the most  
742 important cues to vowel identity were F1, F2 and duration, which replicated previous work  
743 (e.g., Kuronen, 2000; Lindblom, 1963). For some contrasts, F3 increased separability,  
744 highlighting the F3 dependency for certain rounding contrasts. The category separability  
745 index calculated for the long-short vowel pairs suggested that categories in each quantity  
746 contrasts were, with the exception of /y/ and /e/, just as separable in F1-F2 as F1-F2-F3  
747 space. This would suggest differences in lip-rounding for the long and short /y/, that were

748 not found for the other rounded vowels. The results furthermore suggested that static  
749 measurements across talkers, while informative, were insufficient to accurately capture the  
750 spectral acoustics of some vowels, as indicated by the amount of overlap when static  
751 formants were considered. A more exhaustive description can be achieved by including  
752 dynamic analyses, as shown in this paper.

753 This was especially true for the [ɛ] - [æ:], [u:] - [o:] and [e:] - [i:] - [y:] distinctions, given  
754 the direction and scope of formant trajectory movements across the vowel segment. The  
755 analysis of the empirical data furthermore suggested that a larger portion of the movement  
756 took place at the later part of the segment for most vowels, irrespective of the magnitude  
757 of change in F1 and F2. However, [i:], [y:] and [ɯ] constituted important exceptions as they  
758 displayed more movement in the first three to four time-points. Of note, some of the short  
759 vowels showed formant movements of equal or larger magnitude as certain long vowels,  
760 which seems to signal the importance of vowel dynamics for long and short vowels alike.

761 This has been investigated in work on other languages (e.g., Hillenbrand et al., 1995;  
762 Watson & Harrington, 1999), but has largely been lacking in studies on Swedish. In terms  
763 of distinguishing between diphthongization and merely formant movement, the trajectory  
764 plots suggested diphthongization towards an open quality in primarily [e:] and [o:].

765 GAMMs fit to the data contributed with further insights into the formant dynamics  
766 of individual contrasts as well as for vowels differing in quantity, and allowed for an  
767 assessment of the relative contribution of formant dynamics to cue dependencies for  
768 neighboring and more distant contrasts. For instance, GAMMs fitted to the long-short  
769 vowel pairs indicated that all pairs were reliably distinguished by all included spectral cues,  
770 highlighting the informativity carried by formant dynamics in quantity distinctions. The  
771 static and dynamics analyses thus both suggest that quantity distinctions are not separate  
772 from quality distinctions in Central Swedish.

773 With regard to the neighboring contrasts hypothesized to rely on formant  
774 dynamics—the [i:] - [y:] - [ɯ] - [e:], [o:] - [u:], [ɛ:] - [æ:], and [ø]-[œ] contrasts, and their short

775 counterparts—the results indicated significant differences for all comparisons and all cues  
776 with the exception of [i:] - [y:], [ɛ:] - [æ:], and [ɛ] - [æ] predicting F3, and [ɪ] - [ʏ], [ø:]-[œ:]  
777 predicting F1. This would seem to suggest that the movements in these vowels along these  
778 cues, are not contributing to vowel quality information. Presumably, the dynamics in  
779 remaining cues is sufficient for distinguishability.

780 Since the GAMMs were fit to each cue separately, they allowed for an assessment of  
781 the relative weight of each cue, compared to the separability index where the by-cue  
782 contributions can only be evaluated indirectly. An inherent limitation in the separability  
783 index, as implemented in this paper, is the simplifying assumption that all dimensions  
784 within a cue combination carry equal weight. It is therefore not possible to assess whether  
785 the relation between the cues in each space is symmetrical or not, that is, in the F1-F2  
786 space, we do not know whether F1 carries as much information as F2 for separability, and  
787 vice versa. In addition, given that the comparisons are pairwise, they are limited to  
788 explaining the relation between two vowels in a contrast. As such, they cannot inform us of  
789 the separability of a given vowel from other neighboring vowels, or the overall category  
790 separability in the entire space. This is, however, a limitation that the separability index  
791 shares with the GAMMs. Neither of these methods are able to assess the distinguishability  
792 or confusability of all vowels under different cue combinations in one analysis. Nor can they  
793 inform us of the *perceived* distinguishability, even if it is reasonable to assume that reduced  
794 overlap between tokens of neighboring categories would increase intelligibility (e.g.,  
795 Bradlow, 1995; Wright, Local, Ogden, & Temple, 2004). For this, one could either fit a  
796 multinomial model predicting vowel from cue combinations, or a perceptual model that can  
797 assess the predicted consequences for perception by the use of categorization accuracies and  
798 confusion matrices. Such models can be an important avenue for future research, a starting  
799 point for the design of perception studies that can shed more light on the consequences of  
800 the present results for the perception of Central Swedish vowels. For instance, in a  
801 language with a systematic quantity distinction such as Swedish, the role of spectral cues

802 in long-short vowel pair distinctions could be assessed by exposing listeners to synthesized  
803 versions where long and short vowel duration is crossed with the allophones' spectral  
804 information for any given phoneme. Furthermore, as the results seem to support claims of  
805 the hypothesized importance of formant dynamics for vowel distinctions, more insight into  
806 the effect of formant dynamics for vowel perception could be gained from having listeners  
807 categorize tokens extracted from different segments of the long vowels, e.g., the first three  
808 time-points vs. the three final time-points (c.f., Jenkins, Strange, & Miranda, 1994;  
809 Strange, 1989). The design of such experiments can be informed by modeling the predicted  
810 perceptual consequences of different cue spaces, and of considering different vowel segments.

## 811 4.2 Vowel shifts

812 In Study 2, the acoustic characteristics of SwehVd was subsequently compared to that of 8  
813 reference talkers of Central Swedish from the SweDia materials recorded approximately one  
814 generation ago, in order to investigate the scope and spread of vowel changes noted in  
815 Study 1 (some of which, previously reported in Pelzer & Boersma, 2019). The aim was to  
816 present a broad-scale comparison by evaluating changes in all but one of the 21 categories,  
817 using empirical formant analysis and the *op*.

818 Several vowels had shifted quite considerably compared to earlier work on Central  
819 Swedish (c.f., Fant et al., 1969; Kuronen, 2000). Two of the most important shifts  
820 concerned the fronting of [i] and [y] and the centralization of [i:] and [y:], where [i] and [y]  
821 appeared to maintain their positions as high front vowels, while [i:] and [y:] had centralized  
822 to the mid-center part of the space. Compared to the SweDia reference materials, the shifts  
823 in [i], [i:] and [y:] were not as substantial as for other vowels, which suggests that these  
824 changes had begun already in 1999, and then continued throughout the first two decades of  
825 the 2000s (see also, Pelzer & Boersma, 2019).<sup>16</sup>

---

<sup>16</sup> Unfortunately, Pelzer and Boersma (2019)'s study on diphthongization only included the long vowels, it is therefore difficult to know whether the fronting of the short vowels was as pronounced in 2019.

826        The high front vowels form a particularly interesting part of the space given that  
 827    several possibly related processes might be ongoing. The visualizations and GAMMs  
 828    suggested that F3 dynamics did not contribute to the distinction between [i:] and [y:], and  
 829    that [y:] appeared to be less rounded than [y]. This finding conflicted with previous work  
 830    on Central Swedish (e.g., Fant et al., 1969; Fujimura, 1967; Kuronen, 2000), and would  
 831    seem to indicate that the [i:] - [y:] contrast might primarily be supported by F1-F2  
 832    dynamics, or by additional acoustic cues not investigated in this study. There is of course  
 833    also the possibility that listeners might disambiguate the two categories using primarily  
 834    visual cues or linguistic information, or perhaps these two categories are not distinguished,  
 835    which would suggest a merger in process among these talkers. A merger might be driven by  
 836    relaxation of lip-rounding, as supported by the higher F3 values found for [y:].

837        Both of these processes—centralization and relatively sustained overlap—could  
 838    however be explained by [i:] and [y:] being produced as damped versions. The presence of a  
 839    damped [i:] would be supported by the lower F2 values, as the consonantal offglide in [i:]  
 840    lowers F2 (Engstrand et al., 2000). A merger of [i:] and [y:] into [i:] has been observed  
 841    among younger talkers in other regions in Sweden, e.g., for Gothenburg Swedish, as  
 842    reported by Gross and Forsberg (2020). If centralization of [i:] is a prerequisite for such a  
 843    merger, as suggested by Gross and Forsberg (2020), the present results might indicate the  
 844    beginning of a merger. A future scenario might involve the loss of the defining feature  
 845    lip-rounding in [y:] if a merger continues, given its weaker perceptual salience (as  
 846    hypothesized by Gross & Forsberg, 2020).<sup>17</sup>

847        Impressionistic listening by the author did support the presence of a buzzing sound,  
 848    similar to [i:], among the majority of talkers in SwehVd, both male and female. The

---

<sup>17</sup> A parallel exist in iotaism in Norwegian and Swedish dialects (Eliasson, 2000), characterized by the loss of lip-rounding resulting in vowel shifts such as [y:] to [i:] for monophthongs, and [øy] to [ei] for diphthongs. Iotaism describes a feature loss for a fairly unusual phonological characteristic—rounded front vowels. Talkers that merge because of iotaism would not necessarily produce [i:] as the damped [i:], which could potentially explain the vowel patterns of some talkers in SwehVd that merge the two vowels but have a rather fronted [i:] without an offglide consonant element.

strength and scope of [i:] varied across talkers, from a relatively strong [i:] (17 talkers), to a weaker buzzing sound (12 talkers), or more of a consonant offglide element similar to [j] following [i:] (3 talkers). Interestingly, several of the talkers that did not produce any apparent final consonant glide or buzz (11 talkers), seemed to have overall less retracted [i:] and [y:], hence supporting the hypothesized link between centralization and consonantal offglide. Further insights into these individual differences can be gained by studying the SwehVd materials on a talker-specific level.

The changed placement of [i:] and [y:] might have affected the surrounding front vowels. For instance, [e:] was *not* the most fronted vowel for some of the talkers who did not centralize [i:] and [y:], which would seem to suggest that the fronted placement of [e:] might be conditioned by a more centralized position of [i:] and [y:]. The order and cause of these events is of course unclear, however, they might suggest underlying structural co-variation, such as chain-shifts, causing several categories to move (e.g., Brand, Hay, Clark, Watson, & Sóskuthy, 2021). Chain-shifts are hypothesized to be causally related, in that the shift in one vowel sets off changes in other vowels (Gordon, 2013; Hay, Pierrehumbert, Walker, & LaShell, 2015; Maclagan & Hay, 2007; Martinet, 1952). Push-chain-shifts are characterized by a moving vowel encroaching on another vowel's space, causing that vowel to move, whereas pull-chain-shifts describe a situation where a travelling vowel leaves a vacuum for another vowel to enter (Gordon, 2013; Łubowicz, 2011). Underlyingly, this structural co-variation might be driven by system-wide evolutionary pressure, i.e. the drive towards symmetry (Boersma, 1998) or vowel spaces that maintain or maximize contrast or dispersion (Liljencrants & Lindblom, 1972; Martinet, 1952; Schwartz, Boë, Vallée, & Abry, 1997). The centralization of [i:] and [y:] could thus have pushed [i] and [y] to travel front in order to maintain a contrast between the long-short pairs. Alternatively, the fronting of [e:] might have pushed back [i:] and [y:], however, what might be causing [e:] to move remains unclear, given that the neighbouring vowel [ɛ:] substantially lowers. A future scenario might involve [e:] once more descending to a position closer to its shorter counterpart, and

876 closer to the position it occupies in the latter part of the segment (see Figure 7). Such a  
877 change would potentially result in a sparser front space, possibly setting the stage for a  
878 more compact system overall. Future studies should aim to assess this co-variation in a  
879 principled way for further insight into which of these factors might be driving the shifts.

880 Other noticeable changes outside of the high front area concern [ɛ:] and [œ], that have  
881 both lowered since 1999. The lowering of [ɛ:] was anticipated by e.g., Riad (2014), and  
882 confirmed in T. Leinonen (2010) and Pelzer and Boersma (2019). Even if [œ] has lowered  
883 compared to SweDia, [ø:] and [ø] are still located beneath [œ], which suggest a backing  
884 rather than lowering of /ø/ in position before retroflex segments (c.f., Riad, 2014).  
885 Furthermore, the difference in placement of [œ] might result from a presence (SweDia) or  
886 lack of (SwehVd) neutralization of [ø] into [ø]. In both materials, [ø] precedes /r/ (SwehVd -  
887 *hörr*, SweDia - *dörr*), and in SweDia, [ø] is considerably closer to [ø] than in the SwehVd. If  
888 this result holds, it would suggest that the previously documented merger of short [ø] and  
889 [ø] might have lost its importance in Central Swedish over the last generation (Ståhle, 1965;  
890 Wenner, 2010; but see Kotsinas, 1995 for younger Stockholm talkers).

### 891 4.3 Methodological considerations

892 Methodological considerations for the present study that goes beyond those already  
893 discussed concern 1) the choice of how to measure some of the acoustic cues, 2) methods  
894 for assessing ongoing changes in the vowel space, and 3) materials used for comparing  
895 vowel changes over time. The first consideration concerns two of the cue measurements, F0  
896 and duration, and how they were measured and evaluated. Here, the average F0 across the  
897 vowel segment was reported. However, pitch contours have been known to influence  
898 perceived duration and prominence of vowels (e.g., Gussenhoven & Zhou, 2013). One could  
899 therefore claim that a more accurate acoustic representation of vowels should have included  
900 pitch contours. Similarly, the temporal analysis of the effect of duration on long-short  
901 vowel pairs could be supplemented with measures of consonant ratios (e.g., Pelzer &

902 Boersma, 2019; Schaeffler, 2005). These limitations can be addressed in future studies as  
903 the SwehVd database is publicly available in an online repository (<https://osf.io/ruxnb/>).

904 Second, many acoustic-phonetic studies often include systematic listening by trained  
905 phoneticians, validated through measures of inter-rater reliability (for a review, see  
906 Cucchiarini, 1995; Gross & Forsberg, 2020; Kuronen, 2000; Pelzer & Boersma, 2019). In  
907 this study, proposals of the presence of a possible ongoing merger and/or centralization of  
908 [i:] - [y:], and the absence of merger of [œ] - [ø], would all receive stronger support if  
909 supplemented by systematic (and not impressionistic) listening. Furthermore, the amount  
910 and scope of vowel changes were assessed using the *op*. This non-parametric measure  
911 calculates each vowel's movement between two anchor vowels separately, as such, it does  
912 not take the entire system into account and can therefore miss patterns of co-variation  
913 within and across talkers. Employing a modeling approach similar to what Brand et al.  
914 (2021) suggest, could provide a more principled way of assessing vowel change system-wide,  
915 while accounting for by-talker variability.

916 Third and finally, different phonological contexts were employed for recording of the  
917 two databases compared, which can constrain claims made about vowel changes given  
918 differences in potential effects of coarticulation. This limitation was partly mitigated by  
919 selecting the steady-state portion of the vowel for comparison. A more comparable  
920 database would of course strengthen claims of vowel shifts. Unfortunately, given that  
921 existing databases on Central Swedish display substantial variation in composition (Table  
922 1), there would inevitably be some loss in comparability irrespective of what reference  
923 materials chosen. The SwehVd sought to address this issue by recording hVd words, with  
924 the intention of enhancing comparability in studies on Central Swedish, and with the  
925 welcome side-effect of facilitating cross-linguistic investigations.

## 926 5 Conclusions

927 The present study has reported on the acoustic properties of Central Swedish vowels and  
928 how they have changed over the last generation. The spectral and temporal cues  
929 investigated all contributed to distinguishing between the 21 vowels in the Central Swedish  
930 vowel space, with varying weight. More insight into formant dynamics within and between  
931 quantities have been gained by the dynamic analysis presented, which is also of value for  
932 cross-linguistic research. What has been gained with the broad-scale approach of  
933 characterizing the *entire* vowel space adopted here, is of course lost in terms of detailed  
934 investigations of individual vowel contrasts. There is certainly a lot more to say about the  
935 centralization of [i:] - [y:], the potential loss of lip-rounding in [y:], the direction of vowel  
936 shifts, and lowering of [ɛ:], among other things. The acoustic descriptions outlined in this  
937 paper, together with the publicly available SwehVd database, can provide a reference point  
938 for future investigations into these acoustic events and beyond.

## 939 **Ethics statement**

940 This study on human participants was granted an exemption from requiring ethics  
941 approval in accordance with the local legislation and institutional requirements  
942 (Etikprövningsmyndigheten, Uppsala, Sweden). The participants provided their written  
943 informed consent to participate in this study.

## 944 **Data availability statement**

945 The SwehVd dataset presented in this study can be found in an online repository (SwehVd:  
946 <https://osf.io/ruxnb/>). All analyses and visualization code can be found in a separate  
947 online repository (<https://osf.io/7uvj4/>).

## 948 **Funding**

949 The work presented in this study was partially funded by a grant from the Kinander's  
950 foundation (2021), a grant from Kungliga Vetenskapsakademien (2023), and by the  
951 Department of Swedish Language and Multilingualism at Stockholm University.

## 952 **Acknowledgments**

953 Omitted for review

## 954 **Conflict of Interest**

955 The author declares that the research was conducted in the absence of any commercial or  
956 financial relationships that could be construed as a potential conflict of interest.

## 957 6 References

- 958 Assmann, P. F., & Katz, W. F. (2005). Synthesis fidelity and time-varying spectral  
959 change in vowels. *Journal of the Acoustical Society of America*, 117(2), 886–895.  
960 <https://doi.org/10.1121/1.1852549>
- 961 Baayen, H., Vasishth, S., Kliegl, R., & Bates, D. (2017). The cave of shadows:  
962 Addressing the human factor with generalized additive mixed models. *Journal of  
963 Memory and Language*, 94, 206–234.
- 964 Barreda, S. (2021). Perceptual validation of vowel normalization methods for  
965 variationist research. *Language Variation and Change*, 33(1), 27–53.  
966 <https://doi.org/10.1017/S0954394521000016>
- 967 Barreda, S., & Nearey, T. M. (2018). A regression approach to vowel normalization  
968 for missing and unbalanced data. *The Journal of the Acoustical Society of  
969 America*, 144(1), 500–520. <https://doi.org/10.1121/1.5047742>
- 970 Behne, D. M., Czigler, P. E., & Sullivan, K. P. H. (1997). Swedish Quantity and  
971 Quality: A Traditional Issue Revisited. *Reports from the Department of  
972 Phonetics, Umeå University*, 4, 81–83.
- 973 Björsten, S., & Engstrand, O. (1999). Swedish “damped” /i/ and /y/:  
974 Experimental and typological observations. *Proc. 14th ICPPhS*. San Francisco.
- 975 Bleckert, L. (1987). *Centralsvensk diftongering som satsfonetiskt problem*  
976 [*Diphthongization in Central Swedish as a problem of sentence phonetics*].  
977 Uppsala: Skrifter Utgivna Av Institutionen För Nordiska Språk Vid Uppsala  
978 Universitet.
- 979 Boersma, P. (1998). *Functional phonology: Formalizing the interactions between  
980 articulatory and perceptual drives*. Amsterdam: Den HaagHolland Academic  
981 Graphics/IFOTT.
- 982 Boersma, P., & Weenink, D. (2022). *Praat: Doing phonetics by computer [Computer  
983 program]*.

- 984 Bradlow, A. R. (1995). A comparative acoustic study of english and spanish vowels.  
985 *The Journal of the Acoustical Society of America*, 97(3), 1916–1924.
- 986 Brand, J., Hay, J., Clark, L., Watson, K., & Sóskuthy, M. (2021). Systematic  
987 co-variation of monophthongs across speakers of new zealand english. *Journal of*  
988 *Phonetics*, 88, 101096.
- 989 Bruce, G. (2009). Components of a prosodic typology of Swedish intonation. In  
990 *Components of a prosodic typology of Swedish intonation* (pp. 113–146). De  
991 Gruyter Mouton. <https://doi.org/10.1515/9783110207569.113>
- 992 Chuang, Y.-Y., Fon, J., Papakyritsis, I., & Baayen, H. (2021). Analyzing phonetic  
993 data with generalized additive mixed models. In *Manual of clinical phonetics*  
994 (pp. 108–138). Routledge.
- 995 Cucchiarini, C. (1995). Assessing transcription agreement: Methodological aspects.  
996 *Clinical Linguistics and Phonetics*, 10(2), 131–155.  
997 <https://doi.org/10.3109/026992096089851670269-9206>
- 998 Eklund, I., & Traunmüller, H. (1997). Comparative Study of Male and Female  
999 Whispered and Phonated Versions of the Long Vowels of Swedish. *Phonetica*,  
1000 54(1), 1–21. <https://doi.org/10.1159/000262207>
- 1001 Elert, C.-C. (1964). *Phonologic studies of quantity in swedish: Based on material*  
1002 *from Stockholm speakers*. Uppsala: Almqvist & Wiksell.
- 1003 Elert, C.-C. (1980). Diftongeringar och konsonantinslag: Drag i uttalet av långa  
1004 vokaler i svenska av i dag [Diphthongizations and consonantal offglides: Traits  
1005 in the pronunciation of long vowels in the Swedish of today]. *Språken i vårt*  
1006 *Språk*, 168–181.
- 1007 Elert, C.-C. (1981). *Ljud och ord i svenska [Sounds and words in Swedish*  
1008 *language]*. Umeå: Universitetet i Umeå, Almqvist & Wiksell international.
- 1009 Elert, C.-C. (1994). Indelning och gränser inom området för den talade svenska:  
1010 En aktuell dialektografi [Distribution and boundaries within the area of spoken

- 1011 Swedish: An up-to-date dialectography]. In *Diabas: Vol. 4. Kulturgränser - myt*  
1012 *eller verklighet?* (pp. 215–228). Institutionen för nordiska språk vid Umeå  
1013 Universitet.
- 1014 Elert, C.-C. (2000). *Allmän och svensk fonetik [General and Swedish phonetics]* (8.,  
1015 omarb. uppl). Stockholm: Norsteds.
- 1016 Eliasson, S. (2000). Typologiska och areallingvistiska aspekter på de nordeuropeiska  
1017 språkens fonologi [Typological and area linguistic aspects of the phonology of  
1018 the Northern European languages]. In: *Ernst Håkon Jahr (Ed.), Språkkontakt —*  
1019 *Innverknaden Frå Nedertysk På Andre Nordeuropeiske Språk, 21–70. (Nord*  
1020 *2000:19.) København: Nordisk Ministerråd, 2000, 21–70.*
- 1021 Eliasson, S. (2022). The phonological status of Swedish au and eu: Proposals,  
1022 evidence, evaluation. *Nordic Journal of Linguistics*, 1–42.  
1023 <https://doi.org/10.1017/S0332586522000233>
- 1024 Engstrand, O. (1999). Swedish. In *Handbook of the International Phonetic  
1025 Association: A guide to the usage of the International Phonetic Alphabet*.  
1026 Cambridge: Cambridge University Press.
- 1027 Engstrand, O. (2004). *Fonetikens grunder [The basics of phonetics]*. Lund:  
1028 Studentlitteratur.
- 1029 Engstrand, O., Björsten, S., Lindblom, B., Bruce, G., & Eriksson, A. (2000). Hur  
1030 udda är Viby-i? Experimentella och typologiska observationer [How peculiar is  
1031 Viby-i? Experimental and typological observations]. *Folkmålsstudier*, 39, 83–95.
- 1032 Eriksson, A. (2004). SweDia 2000: A swedish dialect database. In P. J. Henrichsen  
1033 (Ed.), *Babylonian confusion resolved. Proceedings of the nordic symposium on*  
1034 *the comparison of languages, no. 1 in copenhagen working papers in LSP* (pp.  
1035 33–48).
- 1036 Fant, G. (1959). Acoustic analysis and synthesis of speech with applications to  
1037 swedish. *Ericsson Technics*, 15, 3–108.

- 1038 Fant, G. (1971). Notes on the Swedish Vowel System. In L. Hammerich, R.  
1039 Jakobson, E. Zwirner, & E. Fischer-Jørgensen (Eds.), *Form and substance:*  
1040 *Phonetic and linguistic papers*. Odense: Andelsbogtrykkeriet.
- 1041 Fant, G. (1983). Feature analysis of Swedish vowels - a revisit. *STL-QPSR*, 24(2-3),  
1042 001–019.
- 1043 Fant, G., Hennigsson, G., & Stålhammar, U. (1969). Formant frequencies of  
1044 Swedish vowels. *STL-QPSR*, 10(4), 026–031.
- 1045 Fujimura, O. (1967). On the Second Spectral Peak of Front Vowels: A Perceptual  
1046 Study of the Role of the Second and Third Formants. *Language and Speech*,  
1047 10(3), 181–193. <https://doi.org/10.1177/002383096701000304>
- 1048 Gordon, M. J. (2013). Investigating chain shifts and mergers. *The Handbook of*  
1049 *Language Variation and Change*, 203–219.
- 1050 Gross, J., & Forsberg, J. (2020). Weak Lips? A Possible Merger of /i:/ and /y:/ in  
1051 Gothenburg. *Phonetica*, 77(4), 268–288. <https://doi.org/10.1159/000499107>
- 1052 Gussenhoven, C., & Zhou, W. (2013). *Revisiting pitch slope and height effects on*  
1053 *perceived duration*. 1365–1369. <https://doi.org/10.21437/Interspeech.2013-360>
- 1054 Hadding, K., Hirose, H., & Harris, K. S. (1976). Facial muscle activity in the  
1055 production of swedish vowels: An electromyographic study. *Journal of Phonetics*,  
1056 4(3), 233–245. [https://doi.org/10.1016/S0095-4470\(19\)31246-X](https://doi.org/10.1016/S0095-4470(19)31246-X)
- 1057 Hadding-Koch, K., & Abramson, A. S. (1964). Duration Versus Spectrum in  
1058 Swedish Vowels: Some Perceptual Experiments2. *Studia Linguistica*, 18(2),  
1059 94–107. <https://doi.org/10.1111/j.1467-9582.1964.tb00451.x>
- 1060 Hammarström, G., & Norman, L. (1957). Om den frikativa slutfasen vid de svenska  
1061 långa vokalerna i, y, u, w.[On the final fricative phase in the Swedish long vowels  
1062 i, y, u, w.]. *Nordisk Tidsskrift for Tale Og Stemme*, 17(3).
- 1063 Hay, J. B., Pierrehumbert, J. B., Walker, A. J., & LaShell, P. (2015). Tracking word  
1064 frequency effects through 130 years of sound change. *Cognition*, 139, 83–91.

- 1065 Henton, C. G. (2005). Creak as a sociophonetic marker. *The Journal of the  
1066 Acoustical Society of America*, 80(S1), S50–S50.  
1067 <https://doi.org/10.1121/1.2023837>
- 1068 Hillenbrand, J. M., Clark, M. J., & Nearey, T. M. (2001). Effects of consonant  
1069 environment on vowel formant patterns. *The Journal of the Acoustical Society of  
1070 America*, 109(2), 748–763. <https://doi.org/10.1121/1.1337959>
- 1071 Hillenbrand, J. M., Getty, L. A., Clark, M. J., & Wheeler, K. (1995). Acoustic  
1072 characteristics of american english vowels. *Journal of the Acoustical Society of  
1073 America*, 97(5), 3099–3111.
- 1074 Hillenbrand, J. M., & Nearey, T. M. (1999). Identification of resynthesized /hVd/  
1075 utterances: Effects of formant contour. *The Journal of the Acoustical Society of  
1076 America*, 105(6), 3509–3523. <https://doi.org/10.1121/1.424676>
- 1077 Jacewicz, E., & Fox, R. A. (2018). Regional variation in fundamental frequency of  
1078 American English vowels. *Phonetica*, 75(4), 273–309.  
1079 <https://doi.org/10.1159/000484610>
- 1080 Jenkins, J. J., Strange, W., & Miranda, S. (1994). Vowel identification in  
1081 mixed-speaker silent-center syllables. *The Journal of the Acoustical Society of  
1082 America*, 95(2), 1030–1043. <https://doi.org/10.1121/1.410014>
- 1083 Johnson, K. (2005). Speaker normalization in speech perception. In D. B. Pisoni &  
1084 R. E. Remez (Eds.), *The Handbook of Speech Perception* (pp. 363–389). John  
1085 Wiley & Sons, Inc.
- 1086 Johnson, K., & Sjerps, M. J. (2021). Speaker normalization in speech perception. In  
1087 J. S. Pardo, L. C. Nygaard, R. E. Remez, & D. B. Pisoni (Eds.), *The handbook  
1088 of speech perception* (pp. 145–176). John Wiley & Sons, Inc.  
1089 <https://doi.org/10.1002/9781119184096.ch6>
- 1090 Joos, M. (1948). Acoustic Phonetics. *Language*, 24(2), 5–136.  
1091 <https://doi.org/10.2307/522229>

- 1092 Kotsinas, U.-B. (1995). *Dialekt [Dialect]* (pp. 267–269). Höganäs:  
1093 Nationalencyklopedin.
- 1094 Kuronen, M. (2000). *Vokaluttalets akustik i sverigesvenska, finlandssvenska och*  
1095 *finska [The acoustics of vowel pronunciation in Sweden Swedish, Finland*  
1096 *Swedish and Finnish]*. Jyväskylä: University of Jyväskylä.
- 1097 Labov, W. (2001). *Principles of linguistic change. 2: Social factors*. Oxford:  
1098 Wiley-Blackwell.
- 1099 Labov, W., Ash, S., & Boberg, C. (2005). *The atlas of north american english:*  
1100 *Phonetics, phonology, and sound change*. De Gruyter Mouton.
- 1101 Ladefoged, P., & Broadbent, D. E. (1957). Information conveyed by vowels. *Journal*  
1102 *of the Acoustical Society of America*, 29, 98–104.
- 1103 Leinonen, K., Pitkänen, Antti J., & Vihanta, V. V. (1981). Rikssvenskt och  
1104 finlandssvenskt ljudsystem ur perceptionssynpunkt [Perceptual perspectives on  
1105 the sound systems of Standard Swedish and Finland Swedish]. *X Fonetikan*  
1106 *päivät*, TaYSYLJ 7, 163–218. Tampere, Finland: Department of Finnish  
1107 language and general linguistics University of Tampere.
- 1108 Leinonen, T. (2010). An Acoustic Analysis of Vowel Pronunciation in Swedish  
1109 Dialects. *Groningen Dissertations in Linguistics*, 83.
- 1110 Liljencrants, J., & Lindblom, B. (1972). Numerical simulation of vowel quality  
1111 systems: The role of perceptual contrast. *Language*, 839–862.
- 1112 Lindblom, B. (1963). *On vowel reduction*. Uppsala: Uppsala University.
- 1113 Łubowicz, A. (2011). Chain shifts. *The Blackwell Companion to Phonology*, 1–19.
- 1114 Maclagan, M., & Hay, J. (2007). Getting fed up with our feet: Contrast  
1115 maintenance and the new zealand english “short” front vowel shift. *Language*  
1116 *Variation and Change*, 19(1), 1–25.
- 1117 Martinet, A. (1952). Function, structure, and sound change. *WORD*, 8(1), 1–32.  
1118 <https://doi.org/10.1080/00437956.1952.11659416>

- 1119                   McAllister, R., Lubker, J., & Carlson, J. (1974). An EMG study of some  
1120                   characteristics of the Swedish rounded vowels. *Journal of Phonetics*, 2(4),  
1121                   267–278. [https://doi.org/10.1016/S0095-4470\(19\)31297-5](https://doi.org/10.1016/S0095-4470(19)31297-5)
- 1122                   Mennen, I., Schaeffler, F., & Docherty, G. (2012). Cross-language differences in  
1123                   fundamental frequency range: A comparison of English and German. *The  
1124                   Journal of the Acoustical Society of America*, 131(3), 2249–2260.  
1125                   <https://doi.org/10.1121/1.3681950>
- 1126                   Nearey, T. M. (1978). *Phonetic Feature Systems for Vowels*. Indiana.
- 1127                   Nearey, T. M., & Assmann, P. F. (1986). Modeling the role of inherent spectral  
1128                   change in vowel identification. *The Journal of the Acoustical Society of America*,  
1129                   80(5), 1297–1308. <https://doi.org/10.1121/1.394433>
- 1130                   Pelzer, J. A., & Boersma, P. (2019). Diphthongization in three regional varieties of  
1131                   Swedish. *Proceedings of the 19th International Congress of Phonetic Sciences*,  
1132                   1144–1148. Canberra, Australia: Australian Speech Science and Technology  
1133                   Association.
- 1134                   Persson, A., Barreda, S., & Jaeger, T. F. (2024). *Comparing accounts of formant  
1135                   normalization against US english listeners' vowel perception*. manuscript;  
1136                   Stockholm University.
- 1137                   Persson, A., & Jaeger, T. F. (2023). Evaluating normalization accounts against the  
1138                   dense vowel space of central swedish. *Frontiers in Psychology*, 14, 01–21.  
1139                   <https://doi.org/10.3389/fpsyg.2023.1165742>
- 1140                   Peterson, G. E. (1961). Parameters of Vowel Quality. *Journal of Speech and  
1141                   Hearing Research*, 4(1), 10–29. <https://doi.org/10.1044/jshr.0401.10>
- 1142                   Peterson, G. E., & Barney, H. L. (1952). Control methods used in a study of the  
1143                   vowels. *Journal of the Acoustical Society of America*, 24(2), 175–184.
- 1144                   R Core Team. (2023). *R: A language and environment for statistical computing*.  
1145                   Vienna, Austria: R Foundation for Statistical Computing. Retrieved from

- 1146 https://www.R-project.org/
- 1147 Renwick, M. E. L., & Stanley, J. A. (2020). Modeling dynamic trajectories of front  
1148 vowels in the american south. *The Journal of the Acoustical Society of America*,  
1149 147(1), 579–595. <https://doi.org/10.1121/10.0000549>
- 1150 Riad, T. (2014). *The phonology of Swedish*. Oxford: Oxford University Press.
- 1151 RStudio Team. (2020). *RStudio: Integrated development environment for r*. Boston,  
1152 MA: RStudio, PBC. Retrieved from <http://www.rstudio.com/>
- 1153 Schaeffler, F. (2005). *Phonological quantity in Swedish dialects: Typological aspects,*  
1154 *phonetic variation and diachronic change*. Umeå: Umeå University, Dep. of  
1155 philosophy and linguistics.
- 1156 Schertz, J. (2013). Exaggeration of featural contrasts in clarifications of misheard  
1157 speech in english. *Journal of Phonetics*, 41(3-4), 249–263.
- 1158 Schertz, J., & Clare, E. J. (2020). Phonetic cue weighting in perception and  
1159 production. *WIREs Cognitive Science*, 11(2), e1521.  
1160 <https://doi.org/https://doi.org/10.1002/wcs.1521>
- 1161 Schötz, S., Frid, J., & Löfqvist, A. (2011). Exotic vowels in swedish – an  
1162 articulographic and acoustic pilot study of /i:/. *Proc. 17th ICPHS*. Hong Kong.
- 1163 Schwartz, J.-L., Boë, L.-J., Vallée, N., & Abry, C. (1997). Major trends in vowel  
1164 system inventories. *Journal of Phonetics*, 25(3), 233–253.
- 1165 Seyfarth, S., Buz, E., & Jaeger, T. F. (2016). Dynamic hyperarticulation of coda  
1166 voicing contrasts. *The Journal of the Acoustical Society of America*, 139(2),  
1167 EL31–EL37.
- 1168 Sóskuthy, M. (2021). Evaluating generalised additive mixed modelling strategies for  
1169 dynamic speech analysis. *Journal of Phonetics*, 84.  
1170 <https://doi.org/10.1016/j.wocn.2020.101017>
- 1171 Sóskuthy, M., Foulkes, P., Hughes, V., & Haddican, B. (2018). Changing words and  
1172 sounds: The roles of different cognitive units in sound change. *Topics in*

- 1173           *Cognitive Science*, 10(4), 787–802.
- 1174           <https://doi.org/https://doi.org/10.1111/tops.12346>
- 1175           Ståhle, C. I. (1965). 'Mötet uppnas på sundag'[The meeting uppnas on sundag].
- 1176           *Språkvård*, 3(3-8), 1–15.
- 1177           Stålhammar, U., Karlsson, I., & Fant, G. (1973). Contextual effects on vowel nuclei.
- 1178           *STL-QPSR*, 14(4), 001–018.
- 1179           Stevens, K. N., & House, A. S. (1963). Perturbation of vowel articulations by
- 1180           consonantal context: An acoustical study. *Journal of Speech and Hearing*
- 1181           *Research*, 6(2), 111–128. <https://doi.org/10.1044/jshr.0602.111>
- 1182           Stevens, M., Harrington, J., & Schiel, F. (2019). Associating the origin and spread
- 1183           of sound change using agent-based modelling applied to /s/-retraction in
- 1184           English. *Glossa: A Journal of General Linguistics*, 4(1), 01–30.
- 1185           <https://doi.org/10.5334/gjgl.620>
- 1186           Stilp, C. (2020). Acoustic context effects in speech perception. *WIREs Cognitive*
- 1187           *Science*, 11(1), 1–18. <https://doi.org/10.1002/wcs.1517>
- 1188           Strange, W. (1989). Evolving theories of vowel perception. *The Journal of the*
- 1189           *Acoustical Society of America*, 85(5), 2081–2087.
- 1190           <https://doi.org/10.1121/1.397860>
- 1191           Strangert, E. (2001). Quantity in ten swedish dialects in northern sweden and
- 1192           Österbotten in finland. *Working Papers/Lund University, Department of*
- 1193           *Linguistics and Phonetics*, 49, 144–147.
- 1194           Syrdal, A. K. (1985). Aspects of a model of the auditory representation of american
- 1195           english vowels. *Speech Communication*, 4(1-3), 121–135.
- 1196           [https://doi.org/10.1016/0167-6393\(85\)90040-8](https://doi.org/10.1016/0167-6393(85)90040-8)
- 1197           Watson, C. I., & Harrington, J. (1999). Acoustic evidence for dynamic formant
- 1198           trajectories in Australian English vowels. *The Journal of the Acoustical Society*
- 1199           *of America*, 106(1), 458–468. <https://doi.org/10.1121/1.427069>

- 1200 Wedel, A., Nelson, N., & Sharp, R. (2018). The phonetic specificity of contrastive  
1201 hyperarticulation in natural speech. *Journal of Memory and Language*, 100,  
1202 61–88.
- 1203 Weirich, M., Simpson, A. P., Öjbro, J., & Ericsdotter Nordgren, C. (2019). The  
1204 phonetics of gender in Swedish and German. *Fonetik 2019, Stockholm, Sweden*,  
1205 10-12 June, 2019, 49–53.
- 1206 Wenner, L. (2010). *När lögnare blir lugnare. En sociofonetisk studie av*  
1207 *sammanfallet mellan kort ö och kort u i uppländskan [When lögnare become*  
1208 *lugnare. A sociophonetic study of the merger between short ö and short u in*  
1209 *Uppland Swedish]*. Uppsala: Uppsala Universitet.
- 1210 Wieling, M. (2018). Analyzing dynamic phonetic data using generalized additive  
1211 mixed modeling: A tutorial focusing on articulatory differences between L1 and  
1212 L2 speakers of english. *Journal of Phonetics*, 70, 86–116.
- 1213 Wright, R., Local, J., Ogden, R., & Temple, R. (2004). Factors of lexical  
1214 competition in vowel articulation. *Papers in Laboratory Phonology VI*, 75–87.
- 1215 Xie, X., & Jaeger, T. F. (2020). Comparing non-native and native speech: Are L2  
1216 productions more variable? *The Journal of the Acoustical Society of America*,  
1217 147(5), 3322–3347. <https://doi.org/10.1121/10.0001141>
- 1218 Young, N. J., & McGarrah, M. (2021). Forced alignment for Nordic languages:  
1219 Rapidly constructing a high-quality prototype. *Nordic Journal of Linguistics*,  
1220 1–27. <https://doi.org/10.1017/S033258652100024X>