

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 on static and dynamic acoustic analyses of the entire vowel  
13 space using a recently released database of *h-VOWEL-d* words (SwehVd). The results  
14 highlight the importance of static and dynamic spectral and temporal cues for Swedish  
15 vowel category distinction. The first two formants and vowel duration are the primary  
16 acoustic cues to vowel identity, however, the third formant contributes to increased  
17 category separability for neighboring contrasts presumed to differ in lip-rounding. In  
18 addition, even though all long-short vowel pairs differed systematically in duration, they  
19 also display considerable spectral differences, suggesting that quantity distinctions are not  
20 separate from quality distinctions in Swedish. The dynamic analysis further suggests  
21 formant movements in both long and short vowels, with [e:] and [o:] displaying clearer  
22 patterns of diphthongization.

23 *Keywords:* vowels, category separability, formant dynamics

24 Word count: X

<sup>25</sup> The acoustic characteristics of Swedish vowels

## <sup>26</sup> 1 Introduction

<sup>27</sup> The Swedish vowel inventory consists of 21 categories that differ in spectral (formant  
<sup>28</sup> frequencies) and temporal cues (duration). It forms a typologically rather complex space,  
<sup>29</sup> characterized by a systematic quantity distinction resulting in 9 long and short vowel pairs,  
<sup>30</sup> 3 different levels of lip-rounding, and contextually conditioned allophones to /ɛ/ and /ø/ in  
<sup>31</sup> position before /r/ or any retroflex segments. Given the crowdedness of the space and  
<sup>32</sup> resulting category overlap for some vowels, previous work has reported on 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, Hennigsson, & Stålhammar,  
<sup>35</sup> 1969; Fujimura, 1967; Kuronen, 2000), duration for certain long-short vowel pairs (Behne,  
<sup>36</sup> Czigler, & Sullivan, 1997), formant movements for some front contrasts (Kuronen, 2000;  
<sup>37</sup> Pelzer & Boersma, 2019), as well as the need to look beyond static point estimates of the  
<sup>38</sup> two primary determinants to vowel identity cross-linguistically, the first two formants (F1  
<sup>39</sup> and F2, e.g., Joos, 1948; Ladefoged & Broadbent, 1957; Nearey & Assmann, 1986;  
<sup>40</sup> Peterson, 1961). Single-point estimates extracted from the steady-state of the vowel where  
<sup>41</sup> the formant pattern is presumed to be static, continues to be widely used to represent  
<sup>42</sup> vowels in an F1-F2 plane (for a review, see Kent & Vorperian, 2018). However, static  
<sup>43</sup> estimates cannot capture how formants move as the signal unfolds; information that has  
<sup>44</sup> been shown to influence listeners' vowel perception (e.g., Assmann & Katz, 2005; Eklund &  
<sup>45</sup> Traunmüller, 1997; Hillenbrand & Nearey, 1999; Kuronen, 2000; Nearey & Assmann, 1986).

<sup>46</sup> This paper investigates the acoustic characteristics of modern-day Swedish vowels in  
<sup>47</sup> analyses that aim to contribute to our understanding of language-specific and  
<sup>48</sup> language-general patterns of vowel acoustics. The paper presents a comprehensive  
<sup>49</sup> description of the primary acoustic cues to vowel identity, using a recently released

50 database of *h-VOWEL-d* (short: hVd) words, the SwehVd database (Persson & Jaeger,  
51 2023). The variety investigated is Central Swedish, the regional standard variety of  
52 Swedish spoken in an area around and beyond Stockholm (eastern Svealand) (Bruce, 2009;  
53 Elert, 1994; Riad, 2014).<sup>1</sup> Existing descriptions of the entire space of 21 vowels rest on  
54 recordings made more than 25 years ago (reported in, e.g., Engstrand, 1999; Kuronen,  
55 2000; Leinonen, 2010; Riad, 2014). Two of the most recent studies are Leinonen (2010) and  
56 Kuronen (2000) (Table 1). The former is based on recordings obtained around 1999 of all  
57 vowels, of which four short vowels were omitted from analysis. It covers 98 rural locations  
58 in Sweden and Swedish-speaking parts of Finland, including reference talkers of Standard  
59 Swedish. The latter covers the entire vowel space but is based on recordings from 1981  
60 (Leinonen, Pitkänen, & Vihanta, 1981). More recent work over the last two decades has  
61 focused on parts of the phonological space, such as the long vowels for diphthongization  
62 studies (Pelzer & Boersma, 2019), two vowels for merger studies (e.g., [θ] - [œ] in Wenner,  
63 2010), allophonic variation in /ɛ/ (Gross, Boyd, Leinonen, & Walker, 2016), or a single  
64 allophone, e.g., the damped [ɪ̯] (Schötz, Frid, & Löfqvist, 2011), see Table 1. These studies  
65 all provide detailed mappings of different parts of the space, and contribute important  
66 insights into the current state of as well as ongoing processes. However, given their focus  
67 on subsets of the space, a comprehensive acoustic mapping of the modern-day Central  
68 Swedish vowel space *in its entirety* is lacking. Given that there is some evidence that  
69 productions of minimal pairs can lead to enhanced contrasts (e.g., Schertz, 2013; Seyfarth,  
70 Buz, & Jaeger, 2016), how representative such subsets are for the vowel space as a whole,  
71 remains an open question. In addition, most previous studies differ in the materials used,  
72 in terms of the size of the database (e.g., number of talkers and repetitions per vowel), the  
73 demographics of talkers (e.g., male/female talkers, region of origin), and phonological  
74 contexts used for recording. For instance, the majority of previous work has either not held

---

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

- <sup>75</sup> the phonetic context constant across vowels, or has investigated isolated vowel production  
<sup>76</sup> out of context or in different CVC contexts (Table 1). This diversity restricts comparison  
<sup>77</sup> across studies on Swedish, as well as cross-linguistically.

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
Gross, Boyd, Leinonen & Walker, 2016	3-15 repetitions of /ɛ/ before retroflex consonants or in other word contexts extracted from running speech (alongside corner vowel productions of [i], [a], [u])	57 female and male talkers of 16-17 years of age (28 from Stockholm; 13 with foreign-born mother)	Euclidean distance; Wilcoxon rank sum test; t-test	Sociolinguistic variation in allophones to /ɛ/ in Gothenburg and Stockholm Swedish
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 itV 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 [ɛ] 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 [ɛ]
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]

78        The materials and methodological approach adopted in the current paper is  
79    motivated by the goal to complement previous work for a comprehensive picture of  
80    modern-day Central Swedish vowels. The paper provides an up-to-date acoustic description  
81    of the entire vowel space of 21 categories, using the SwehVd database (Persson & Jaeger,  
82    2023). The hVd context continues to be widely used in vowel production and perception  
83    studies on languages where the glottal /h/ in onset minimizes supraglottal articulations and  
84    thereby reduces the risk of coarticulatory effects from the surrounding phonetic context, as  
85    in e.g., English and Swedish (as confirmed in e.g., Chesworth, Coté, Shaw, Williams, &  
86    Hodge, 2003; Robb & Chen, 2009). The use of an hVd database hence increases  
87    comparability to studies on other languages (e.g., Hillenbrand, Getty, Clark, & Wheeler,  
88    1995; Peterson & Barney, 1952). The main spectral (F1-F2-F3) and temporal cues to vowel  
89    identity are reported in static analyses (following e.g., Engstrand, 1999; Fant et al., 1969),  
90    as well as dynamic analyses, given the well documented importance of formant dynamics  
91    on vowel production and perception (e.g., Assmann & Katz, 2005; Eklund & Traunmüller,  
92    1997; Hillenbrand & Nearey, 1999; Kuronen, 2000; Nearey & Assmann, 1986). While  
93    fundamental frequency (F0) is not considered an important cue to vowel identity in itself,  
94    it is known to vary between languages, dialects and speech styles and is therefore reported  
95    in the static analysis for a comprehensive picture of the acoustics (e.g., Henton, 2005;  
96    Jacewicz & Fox, 2018; Johnson, 2005; Leung, Jongman, Wang, & Sereno, 2016; Mennen,  
97    Schaeffler, & Docherty, 2012; Weirich, Simpson, Öjbro, & Ericsdotter Nordgren, 2019).

98        The static analysis assesses what cues contribute to vowel distinctions and evaluates  
99    some of the claims introduced in previous work, such as the hypothesized importance of F3  
100   for rounded vs. unrounded high front contrasts (Fant, 1959; Fant et al., 1969; Fujimura,  
101   1967; Kuronen, 2000; Persson & Jaeger, 2024), and to what extent spectral and temporal  
102   cues contribute to long-short vowel pair distinctions (e.g., Behne et al., 1997; Kuronen,  
103   2000). The dynamic analysis explores what part of the space seems more prone to  
104   diphthongization, and investigates how formant dynamics contribute to vowel distinctions.

<sup>105</sup> In contrast with previous work investigating the dynamics of Central Swedish vowels, the  
<sup>106</sup> present study includes both long and short vowels, thus submitting the entire vowel space  
<sup>107</sup> to the same analyses.

<sup>108</sup> The paper is organized as follows. A background to the acoustics of Central Swedish  
<sup>109</sup> vowels is provided by a review of previous work. This is followed by methods and results,  
<sup>110</sup> and finally, a discussion of the results and its consequences for the Central Swedish vowel  
<sup>111</sup> system. All analyses and visualization code for this study can be found in an online  
<sup>112</sup> repository (<https://osf.io/7uvj4/>). This article is written in R Markdown, which allows  
<sup>113</sup> readers to easily replicate the analyses using freely available software (R Core Team, 2023;  
<sup>114</sup> RStudio Team, 2020).

## <sup>115</sup> 1.1 The acoustics of Central Swedish vowels

<sup>116</sup> This section provides a description of the overall inventory of Central Swedish  
<sup>117</sup> monophthongs, and discusses the role of cues beyond F1 and F2. It furthermore presents a  
<sup>118</sup> review of previous studies on diphthongization and formant dynamics.

<sup>119</sup> Central Swedish is most often described as having nine vowel phonemes: /i/, /y/,  
<sup>120</sup> /ɯ/, /e/, /ɛ/, /ø/, /ɑ/, /o/, /u/. The long allophones are [i:], [y:], [ɯ:], [e:], [ɛ:], [ø:], [ɑ:], [o:],  
<sup>121</sup> [u:], and the short allophones are [i], [y], [ɯ], [e], [ɛ], [ø], [ɑ], [ɔ], [u]. The short allophones of /e/  
<sup>122</sup> and /ɛ/ has been reported to neutralize as [ɛ] in Central Swedish, resulting in 17 vowels,  
<sup>123</sup> rather than 18 (Riad, 2014). There is also evidence of neutralization of the short /ø/ and  
<sup>124</sup> /ɯ/ as [ø] among some talkers, primarily in position before a retroflex (Ståhle, 1965;  
<sup>125</sup> Wenner, 2010). In addition to these 17 vowels, there are 4 additional long and short  
<sup>126</sup> allophones—[æ:], [æ], [œ:], and [œ], as /ɛ/ and /ø/ lower in position before /r/ or any  
<sup>127</sup> retroflex segment (e.g., Kuronen, 2000; Riad, 2014). Traditionally, Central Swedish has  
<sup>128</sup> been described using four height levels and three backness levels (Riad, 2014).

<sup>129</sup> It has furthermore been suggested that Central Swedish is defined by three levels of

<sup>130</sup> lip-rounding, where the rounded vowels are most often referred to as either inrounded, with  
<sup>131</sup> an extreme narrowing of the lips—[u:] and [u:], or outrounded, with a lesser degree of  
<sup>132</sup> lip-narrowing and more protruded lips—[y:], [ø:], [œ:], [o:], and the remaining vowels defined  
<sup>133</sup> as unrounded (e.g., Fant, 1971; McAllister, Lubker, & Carlson, 1974). Previous work has  
<sup>134</sup> claimed that lip-rounding is particularly important for some vowel distinctions. For  
<sup>135</sup> instance, [i:] and [y:] have been described as overlapping in F1-F2 space, but as more  
<sup>136</sup> separable when F3 is considered (Fant, 1959; Fant et al., 1969; Fujimura, 1967; Kuronen,  
<sup>137</sup> 2000).

<sup>138</sup> The vowels in each pair have been reported to differ systematically in duration, with  
<sup>139</sup> short-long vowel to vowel ratios on average .65-.67 for Central Swedish (Elert, 1964;  
<sup>140</sup> Kuronen, 2000; Strangert, 2001). Spectral differences have traditionally been interpreted as  
<sup>141</sup> a consequence of the durational distinction, hence assuming a trading relationship between  
<sup>142</sup> spectral and temporal cues (for a review, see Schaeffler, 2005). It has been hypothesized  
<sup>143</sup> that most of the durational variation is carried by F2 (e.g., Kuronen, 2000; Lindblom,  
<sup>144</sup> 1963). Previous work has found the largest spectral differences for the [u:] - [ø], and [u:] - [a]  
<sup>145</sup> vowel pairs, and the smallest differences for [ε:] - [ɛ], and [ø:] - [ø] (e.g., Kuronen, 2000). For  
<sup>146</sup> pairs with small spectral differences, duration is presumably more important for vowel  
<sup>147</sup> distinction. Perceptual studies on synthesized speech from talkers of Stockholm Swedish  
<sup>148</sup> have confirmed that duration is the primary cue for [i:] - [ɪ], and [o:] - [ɔ] (Behne et al.,  
<sup>149</sup> 1997; for results on Southern Swedish and additional vowel pairs, [ε:] - [ɛ], [ø:] - [ø], see  
<sup>150</sup> Hadding-Koch & Abramson, 1964). The extent to which *all* long-short vowel pairs rely on  
<sup>151</sup> spectral cues is less known, as studies have focused on subsets of pairs.

<sup>152</sup> According to previous work, several of the long vowels in Central Swedish tend to  
<sup>153</sup> diphthongize in their phonetic realization. Diphthongization is considered prosodically  
<sup>154</sup> conditioned and is the strongest in stressed vowels (Bleckert, 1987; Leinonen, 2010).<sup>2</sup>

---

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

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

170 Some talkers of Central Swedish have been reported to realize [i:], [y:], [œ] and [u:]  
171 with a consonantal offglide, where the end-point of [i:] is described as a palatal  
172 approximant [j], the end-point of [y:] a voiced labio-palatal approximant [ɥ], the end-point  
173 of [œ] and [u:] a voiced bilabial fricative [β] (Elert, 1980; Hammarström & Norman, 1957;  
174 McAllister et al., 1974). Furthermore, both [i:] and [y:] can be damped and produced with  
175 a buzzing sound, phonetically realized as [i̯]. The buzzing sound is presumably generated  
176 by the co-articulation of the vowel and a voiced fricative sound similar to [z] (Elert, 1980;  
177 Engstrand, Björsten, Lindblom, Bruce, & Eriksson, 2000). The damped [i:] has been found  
178 in several dialects across Sweden, both in rural areas and in the cities of Gothenburg and  
179 Stockholm (Björsten & Engstrand, 1999; Elert, 1980; Engstrand et al., 2000; Gross &  
180 Forsberg, 2020; Riad, 2014). It has been claimed to carry strong socio-indexical meaning  
181 across locations; indexing place in rural areas, and class and gender in urban areas (e.g.,

<sup>182</sup> Bruce, 2010; Gross, 2018; Kotsinas, 1994; Nilsson, Wenner, Leinonen, & Thorselius, 2021).

<sup>183</sup> In work on Swedish dialectology, it is often referred to as the Viby-*i*, and in the Stockholm  
<sup>184</sup> area, as the Lidingö-*i*. Acoustically, it manifests primarily as a lowering of F2, thus  
<sup>185</sup> occupying a more centralized position in the F1-F2 space. Schötz et al. (2011) describe it  
<sup>186</sup> as a central palatal vowel, as the articulatory correlates involve a retracted and lower  
<sup>187</sup> tongue position, the tip of the tongue being higher than blade and dorsum.

<sup>188</sup> The methodology employed is presented next, beginning with a description of the  
<sup>189</sup> materials used.

## <sup>190</sup> 2 Methods

### <sup>191</sup> 2.1 Materials

<sup>192</sup> The materials used is a corpus of Swedish hVd word recordings, collected by Anna Persson  
<sup>193</sup> and Maryann Tan (Stockholm University) in 2020-2024, the SwehVd. An initial version of  
<sup>194</sup> the corpus with 24 female talkers is described in Persson and Jaeger (2023). For this paper,  
<sup>195</sup> the final release is presented, including 24 additional male talkers. All recordings,  
<sup>196</sup> annotations, and acoustic measurements are available at <https://osf.io/ruxnb/>. SwehVd  
<sup>197</sup> covers the entire monophthong inventory of Central Swedish, including all nine long vowels,  
<sup>198</sup> eight short vowels, and the four allophones to /ɛ/ and /ø/.<sup>3</sup> SwehVd focuses on a single  
<sup>199</sup> regional variety, providing high resolution within and across talkers for this variety with 10  
<sup>200</sup> recordings of each hVd word from each of the 48 talkers (N = 24 female), for a total N of

---

<sup>3</sup> The words used to elicit the 21 vowels were: *hid-[i]*, *hyd-[y]*, *hud-[ɯ]*, *hed-[ε]*, *häd-[ɛ]*, *höd-[ɔ]*, *had-[ɑ]*, *håd-[o:]*, *hod-[u:]*, *hidd-[ɪ]*, *hydd-[ʏ]*, *hudd-[ø]*, *hedd-[ɛ]*, *hädd-[ɛ]*, *hödd-[ø]*, *hadd-[ɑ]*, *hådd-[ɔ]*, *hodd-[ʊ]*, *härd-[æ:]*, *härr-[æ]*, *hörđ-[œ:]*, *hörr-[œ]*. The mix of 4 real Swedish words – *hed*, *härd*, *hörđ*, *hud* (English translations: *heath*, *hearth*, *heard* and *skin*, respectively), and 18 phonotactically legal pseudowords in the word list, might have affected talkers' pronunciations. For instance, there are studies indicating frequency and neighborhood density effects on vowel productions, with low-frequency words and words with high neighborhood density being produced more distinctly (e.g., Munson & Solomon, 2004; Wright, Local, Ogden, & Temple, 2004). Similar effects of hyperarticulation interacting with neighborhood density have been found for pseudowords as well, however, with substantial cross-talker variability (e.g., Scarborough, 2012).

201 tokens = 9979. All talkers in the database were L1 talkers of Swedish, born and raised in  
202 the Greater Stockholm area or surroundings, of 18-44 years of age (mean age = 30; SD =  
203 6.82). For more details on the recruitment, recording, pre-processing, segmentation and  
204 annotation procedure, see Persson and Jaeger (2023).

205 For the vast majority of talkers in the SwehVd, *hädd* productions elicited the same  
206 vowel as *hedd* (see Supplementary Information—SI, Figure S1), which confirms the  
207 commonly held assumption that the short allophone of /e/ neutralizes with the short  
208 allophone of /ɛ/ in Central Swedish. In order to have a balanced number of tokens for each  
209 vowel, all *hädd* words were excluded from the subsetted SwehVd materials used in this  
210 study (following Persson & Jaeger, 2023). Recordings on which the talker did not produce  
211 the targeted vowel were also excluded.<sup>4</sup> Furthermore, outliers were identified and removed  
212 by estimating the relative probability of each token's F1-F2 values given the joint  
213 distribution of F1-F2 for that vowel and talker. Tokens outside of the 2.50th to 97.50th  
214 quantile of the bivariate Gaussian distribution were filtered out. To facilitate empirical  
215 analyses and statistical models, all talkers ( $N = 7$ ) with fewer than 4 remaining recordings  
216 for at least one of the vowels were removed. This left data from 41 L1 talkers ( $N=20$   
217 female talkers), with on average 359 ( $SD = 21.50$ ) tokens per vowel (range = 304 to 383),  
218 for a total of 7529 observations.

## 219 2.2 Acoustic analyses

### 220 2.2.1 Measuring acoustic cues to vowel identity

221 The Swedish version of the Montreal Forced Aligner developed by Young and McGarrah  
222 (2021) was used to obtain estimates of word and segment boundaries. The boundaries were  
223 then manually corrected by the author (an L1-talker of Swedish). The formant analysis was  
224 carried out in Praat (Boersma & Weenink, 1992-2022), using the Burg algorithm to extract

---

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

estimates of the first three formants (F1-F3) at five time-points of the vowel (20, 35, 50, 65, 80% into the vowel). The five points were selected to capture formant trajectories and potential diphthongization in the later time-points for the dynamic analysis (following, e.g., Holbrook & Fairbanks, 1962; Jacewicz, Fox, & Salmons, 2011; Kuronen, 2000; Lehiste & Peterson, 1961; Yang, 2019), as well as provide stable measures of the steady-state of the vowel for the static analysis. The Burg algorithm was parameterized with a time step of 0.01 seconds, a window length of 0.025 seconds, and pre-emphasis was applied from 50 Hz. The maximum number of formants was set to 5, with a formant ceiling of 5500 Hz for the female talkers, and 5000 Hz for the male talkers. Vowel duration and F0 was extracted across the entire vowel segment. The Praat script that extracts these cues is shared as part of the SwehVd OSF repository, allowing researchers to choose additional or alternative time points at which to extract formants and F0.

To correct for measurement errors in the automatic extraction of cues, 5 separate univariate distributions of the five extracted cues (F0, F1, F2, F3 and duration) was estimated for each distinct combination of talker and vowel. Points that fell outside of the 2.50th to the 97.50th quantile of the distributions for each vowel were identified, examined for measurement errors, and subsequently corrected. This followed the approach employed for the SwehVd corpus (Persson & Jaeger, 2023), and strikes a middle-ground between the ideal (manual correction of all tokens) and feasibility.

### 2.2.2 Vowel normalization

The raw formant values were transformed into a vowel normalized space using Nearey's uniform scaling account (Nearey, 1978). Formant measurements in Hertz are reported in the SI, Section S1.5. Vowel normalization is used in studies on vowel production and perception to account for acoustically irrelevant cross-talker variation, as caused by differences in anatomical structure, e.g., vocal tract size (for reviews see e.g., Barreda & Nearey, 2018; Johnson & Sjerps, 2021; Stilp, 2020). In vowel production studies such as the

251 present, normalization is primarily used as a methodological tool. Transforming the  
252 formant data into a normalized space reduces differences in F1 and F2 due to physiology,  
253 which can reduce between-talker variability and increase category separability, as visualized  
254 in Figure 1 (compare left and right panel).

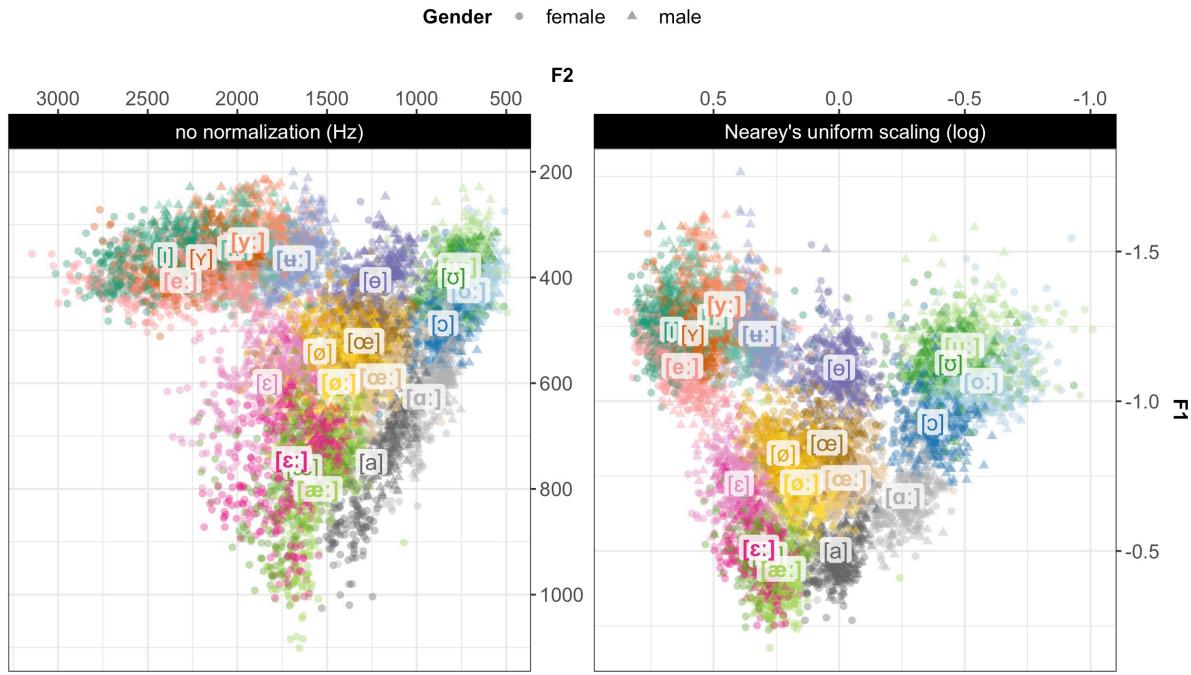
255 Previous work on Swedish has primarily analyzed vowel data in raw Hertz (Björsten  
256 & Engstrand, 1999; Fant et al., 1969; Pelzer & Boersma, 2019), or transformed into Bark  
257 (Fant, 1983; Kuronen, 2000; Schötz et al., 2011; Wenner, 2010), Mel (Lindblom, 1963), or  
258 Lobanov (Gross & Forsberg, 2020). The choice of Nearey's uniform scaling in the present  
259 study was motivated by its previous use in socio-phonetic research to describe and compare  
260 languages and varieties (e.g., Barreda, 2021; Labov, 2001; Labov, Ash, & Boberg, 2005),  
261 and by its plausibility as perceptual model of how we come to achieve robust cross-talker  
262 perception, as it has provided a good fit against both production (e.g., Persson & Jaeger,  
263 2023; Syrdal, 1985) and perception data (e.g., Barreda, 2021; Persson, Barreda, & Jaeger,  
264 2024).

### 265 2.2.3 Static acoustic analysis

266 The static analysis of SwehVd presents formant measurements at the steady state of the  
267 vowel, by averaging across the three mid-points.<sup>5</sup> It maps the entire vowel space of 21  
268 categories and evaluates the relative contribution of F0, F1, F2, F3 and duration to vowel  
269 distinctions, using visualizations of cues and cue correlations.

270 In order to evaluate the hypothesized importance of lip-rounding (F3) for neighboring  
271 unrounded and rounded categories, a category separability index was employed. Category  
separability index and similar measures of reduction in variance or distance between means  
continue to be frequently used in phonetic research (see e.g., Nycz & Hall-Lew, 2013 for a  
review; Fabricius, Watt, & Johnson, 2009; Flynn & Foulkes, 2011; Labov, 2010). While the

<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 (S1.3) presents evaluations of the effect of different measurement points.



*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. Vowel recordings on which the talker produced a different vowel than the intended are excluded (1.21% of all recordings).

275 separability index is a simple and straightforward measure of the relative separability of  
 276 vowel categories, it nevertheless comes with certain limitations to which I return in the  
 277 General discussion (4). Following work by Wedel, Nelson, and Sharp (2018) and Xie and  
 278 Jaeger (2020), each vowel’s separability from the neighboring vowel was calculated as the  
 279 average distance of vowel tokens to the centroid of the neighboring vowel, operationalized  
 280 as (1).

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

281 For instance, for the [i:] - [y:] contrast, first, each talker’s [i:] center was calculated for

282 F1-F2. Next, the distances between each [y] token to the neighboring [i] center from the  
 283 same talker were calculated for F1-F2. Finally, the distances were averaged across all [y]:  
 284 tokens from a talker, resulting in a separability measure for that vowel and talker. The  
 285 same was done in the opposite direction for the same contrast, thus calculating the  
 286 separability of [i] tokens from the [y] center, following Xie and Jaeger (2020). The  
 287 separability index reports the average separability across the two categories in each  
 288 contrast. The higher the index, the greater the separation between categories. The same  
 289 was subsequently done for F1-F2-F3. These two measures of separability for each contrast  
 290 (F1-F2, F1-F2-F3) were then compared to assess whether including F3 would lead to  
 291 increased category separability. The contrasts investigated were [i] - [y], [e] - [y], [ɪ] - [ʏ]  
 292 for comparing unrounded vs. outrounded vowels, and [y] - [œ], [o] - [u], [ɔ] - [ʊ] for  
 293 outrounded vs. inrounded vowels.<sup>6</sup>

294 To quantify the effect of including F3 on category separability, separate linear  
 295 mixed-effects models (LMM) were fit for each contrast, predicting separability from cue  
 296 combination (F1-F2-F3 vs. F1-F2) while including by talker random intercepts.<sup>7</sup> The  
 297 model was formulated as follows: *separability ~ cuecombination + (1|Talker)*. Cue  
 298 combination was treatment-coded with F1-F2 as reference category, thus comparing  
 299 F1-F2-F3 against the F1-F2 baseline.

300 The same process was applied to investigate to what extent long-short vowel pairs  
 301 differ in spectral and temporal cues, by assessing what combination of cues would provide  
 302 the largest separability between the two vowels in each pair. For this evaluation of quantity  
 303 contrasts, the category separability index was calculated for each pair and three different  
 304 cue combinations: F1-F2, F1-F2-F3, or F1-F2-duration. The models were the same as the  
 305 previous sets, with F1-F2 as reference category, comparing each cue combination against

---

<sup>6</sup> I note that this way of calculating category separability assumes talker-specific category representations. The SI S1.7.3 reports separability indices that instead assume talker-independent representations.

<sup>7</sup> By-talker random intercepts was the maximum random effect structure that converged.

306 the F1-F2 baseline.

307 The results of the static analysis are presented in Section 3.1.

308 **2.2.4 Dynamic acoustic analysis**

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

317 In the second part of the dynamic analysis, the hypothesized contribution of formant  
318 dynamics to category information was modeled using generalized additive mixed-effects  
319 models (GAMMs) (Baayen, Vasishth, Kliegl, & Bates, 2017). GAMMs were employed to  
320 assess what cues carry information about vowel quality once formant dynamics were  
321 inspected. GAMMs are increasingly used in phonetic research, due to their suitability in  
322 modeling the non-monotonic complex phonetic patterns found in formants without  
323 assuming linearity or having to rely on the simplifying assumption that vowels can be  
324 reduced to a single F1-F2 point estimate (e.g., Chuang, Fon, Papakyritsis, & Baayen, 2021;  
325 Sóskuthy, 2021; Wieling, 2018). GAMMs have been used in studies on vowels in different  
326 English varieties, e.g., on /u/-fronting in Derby English (Sóskuthy, Foulkes, Hughes, &  
327 Haddican, 2018) and on the front vowel system of Southern American English (Renwick &  
328 Stanley, 2020) but to the best of my knowledge, they have not been implemented in studies  
329 of Swedish vowels. The use of GAMMs thus complements previous work on Central  
330 Swedish that has primarily used visual inspection, formant measurements and linear  
331 models (Table 1).

332 Two main groups of GAMMs were fit in the dynamic analysis. In the first group,  
 333 GAMMs were fit to 6 subsets of neighboring contrasts, hypothesized to differ primarily in  
 334 formant dynamics: [i:] - [y:], [i:] - [ɯ], [i:] - [e:], [o:] - [ɯ], [ɛ:] - [æ:] (Fant, 1971; Kuronen,  
 335 2000; Pelzer & Boersma, 2019). Given the directionality in formant trajectories found for  
 336 [ø]-[œ] (Figure 7), this contrast was also included. To explore potential effects of dynamics  
 337 in the corresponding short vowels, an additional 4 contrasts were modeled. These were not  
 338 entirely identical to the long subsets, for reasons of evident separability in F1-F2 space: [ɪ] -  
 339 [ʏ], [ɔ] - [ʊ], [ɛ] - [æ] and [ø]-[œ]. The general model formulation employed ordered factor  
 340 difference smooths as follows:

341 formant ~ category + Gender + s(timepoint) + s(timepoint, by = category, k =  
 342 5) + s(Talker, bs = "re") + s(Talker, category, bs = "re"). The ordered factor predictor (=  
 343 category) was treatment coded with [i:], [o:], [ɛ:], [ø:], [ɪ], [ɔ], [ɛ], and [ø] as reference  
 344 categories in respective set.

345 The second group consisted of 11 sets of GAMMs fit to each of the long-short vowel  
 346 pairs, aiming for an evaluation of differences between categories within pairs. In each set,  
 347 vowel was treatment coded with the long vowel in each pair as reference vowel.

348 All GAMMs were fit separately for each of the three formants, which necessarily  
 349 meant committing to the simplifying assumption of cue independence. Previous work has  
 350 shown that acoustic cues tend to co-vary (for a review, see Schertz & Clare, 2020). For  
 351 vowels, this is the case for F1 and F2, as shown by the shape and orientation of ellipses in  
 352 Figure 2.

353 The results of the dynamic analysis are presented in Section 3.2.

### 3 Results

354 This study aims to provide a detailed description of the acoustics of all Central Swedish  
 355 vowels and to evaluate the relative importance of certain cues for specific vowel contrasts.

357 as hypothesized in previous work. These include the importance of lip-rounding (F3) for  
 358 high vowel distinctions (Fant, 1959; Fant et al., 1969; Fujimura, 1967; Kuronen, 2000), to  
 359 what extent all long-short vowel pairs differ in quality (formants) and quantity (duration)  
 360 (e.g., Behne et al., 1997; Kuronen, 2000), and what vowels seem to undergo  
 361 diphthongization (Kuronen, 2000; Pelzer & Boersma, 2019). The dynamic analysis  
 362 furthermore explores which cues carry information about neighboring vowel distinctions  
 363 once dynamic information is considered. The results section is divided into two main  
 364 sections, following the static and dynamic analyses.

### 365 3.1 Static spectral and temporal cues to vowel identity

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

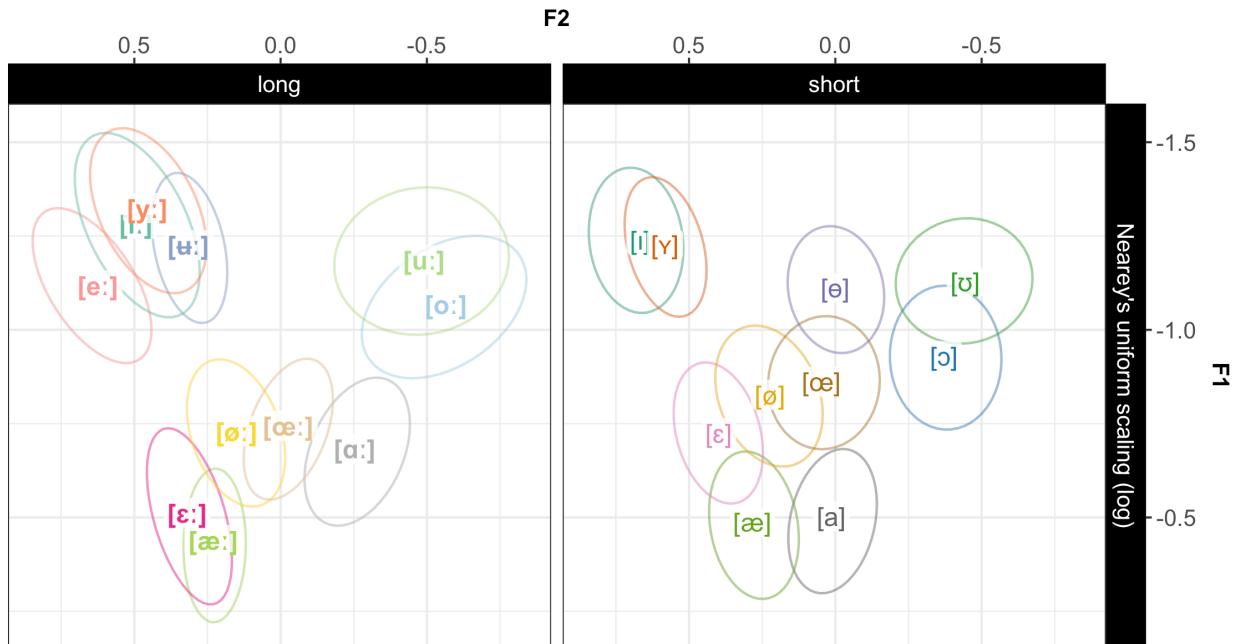


Figure 2. The SwebVd vowel data separated by quantity. Ellipses show bivariate Gaussian 95% confidence interval of vowel means. Vowel labels indicate vowel means across female and male talkers.

Figure 2, left panel, visualizes the long vowels along the two primary cues to vowel identity, F1-F2.<sup>8</sup> Four vowels cluster in the high front part of the space. The mid-high [e:] occupies a substantially higher position than in many previous descriptions, and is also the most fronted vowel (c.f., Fant et al., 1969; Kuronen, 2000; but see Engstrand et al., 2000; Pelzer & Boersma, 2019). The high [i:] and [y:] are rather mid-central, and exhibit substantial overlap with each other and with the neighboring [u:]. The [u:] - [o:], and [ε] - [æ] categories are also partly overlapping.

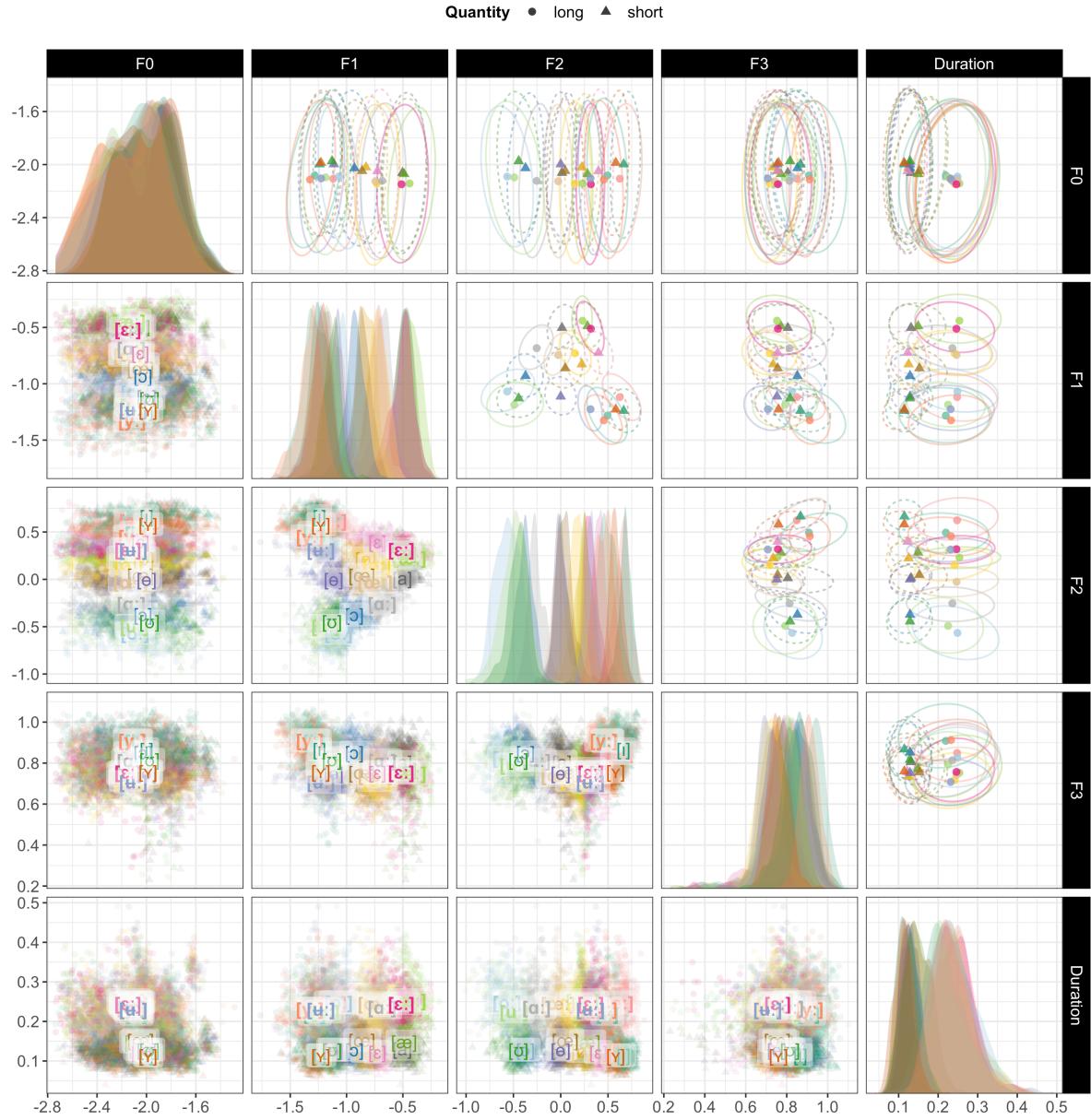
The short vowels (right panel), present a slightly more compact space, however with increased category separability (c.f., Riad, 2014).<sup>9</sup> For some vowel pairs, overlap is clearly reduced for the short vowels, e.g., for [ɪ] - [ʏ], [ɛ] - [æ], and [ɔ] - [ʊ]. Of note, the high vowels [ɪ] and [ʏ] are more fronted and more peripheral than their long counterparts, which does not replicate previous studies (for Central Swedish, see e.g., Fant, 1971; Kuronen, 2000; for short-long contrasts in other languages, see e.g., Clopper, Pisoni, & De Jong, 2005; Hillenbrand et al., 1995). There are further no indications of neutralization of the short /ø/ and /u/ as [ø] among these talkers (c.f., Ståhle, 1965; Wenner, 2010).

### 3.1.1 Cues and cue correlations

For the pairwise combinations of the five spectral and temporal cues—F0, F1, F2, F3 and duration, see Figure 3 from Persson and Jaeger (2023) updated to include data from the 21 male talkers. Unsurprisingly, the densities along the diagonal suggest that F0 carries the least information about vowel identity, exhibiting less between-category separation than all other cues.

<sup>8</sup> The SI presents the mean cue values for the male and female talkers, Tables S1 and S2. As expected, the male talkers have lower formant values and lower F0s than the female talkers (average F0 across long and short categories for female talkers = 204, for male talkers = 119).

<sup>9</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.

390 As is to be expected, vowels differing in quality are most separated in the F1 and F2

391 panels. Inspecting the off-diagonals in Figure 3, the F1-F3 and F3-F2 panels both display

392 increased separation between the neighboring outrounded [y:] and inrounded [ɯ:], and

393 unrounded [i] and outrounded [y], compared to when plotted along F1-F2, which points to

394 the importance of F3 for these vowels. Interestingly, the almost complete overlap between

395 [i:] and [y:] in F1-F2 space overall remains when F3 is considered, even if individual

396 differences in the amount of overlap exist. Most talkers produce these two vowels very close

397 in F1-F2 space, and only slightly separated in F2-F3 space, while others display a continued

398 overlap when considering F3 (for reference, one talker of each type are displayed in SI

399 Figure S4). This would seem to suggest that F3 might carry less importance as distinctive

400 feature for [i:] - [y:] than previously established (c.f., Fant, 1959; Fant et al., 1969).

401 In order to quantitatively assess whether the distinction between closely neighboring

402 unrounded and rounded categories increased when F3 was considered, the category

403 separability of these vowels was calculated based on F1 and F2, and subsequently compared

404 against the separability calculated when including F3. If separability were to increase when

405 F3 was added, it would suggest that F3 does contribute to category distinctions.

406 Three general observations can be made from Figure 4. First, category separability is

407 overall lower for some contrasts when only F1 and F2 were considered, e.g., the [i:] - [y:],

408 and [i] - [y] contrasts, indicating their overlap in F1-F2 space (Figure 4:**column A**).

409 Second, including F3 overall increases category separability, more so for some contrasts

410 than others. The proportional increase in category separability relative to F1-F2 baseline

411 (Figure 4:**column B**) is largest for the [i:] - [y:], [i] - [y], and [y:] - [ɯ:] contrasts. How much

412 separability increases by adding F3 varies quite substantially between talkers for the [i:] -

413 [y:] and [y:] - [ɯ:] contrasts, as indicated by the large confidence intervals. Indeed, when

414 assuming talker-independent representations, the relative increase in separability is less

415 pronounced for these contrasts (see additional analyses in SI S1.7.3). Third, the [y:] - [ɯ:]

416 contrast seems to benefit most from the inclusion of F3, resulting in an overall larger

<sup>417</sup> increase in separability for the outrounded vs. inrounded contrasts over the unrounded  
<sup>418</sup> vs. outrounded contrasts.

<sup>419</sup> The LMMs fit to the data (presented in Section 2.2.3) indicated that including F3  
<sup>420</sup> improved category separability for all contrasts (all  $ps < .003$ ). This suggests that the  
<sup>421</sup> subtle differences observed by visual inspection for the [e:] - [y:], [o:] - [u:] and [ɔ] - [ʊ]  
<sup>422</sup> contrasts were nevertheless significant (Summary tables in SI S1.7.1).<sup>10</sup>

### <sup>423</sup> 3.1.2 Quantity vs. quality in long and short vowel pairs

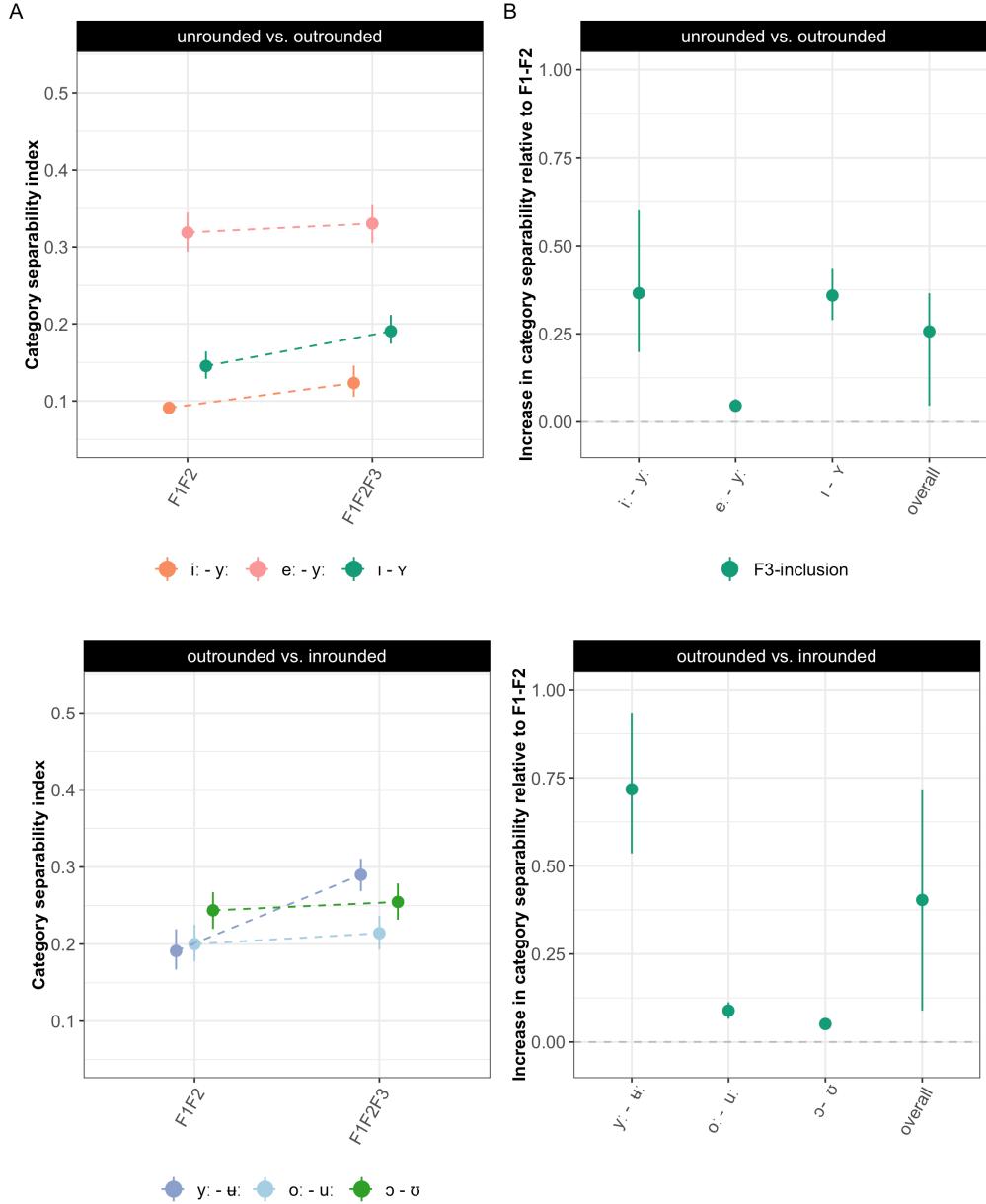
<sup>424</sup> To gain more insight into the extent to which there are spectral and temporal differences  
<sup>425</sup> between long and short vowels, the acoustics of categories within vowel pairs were  
<sup>426</sup> evaluated. This allows for an assessment of whether quantity and quality distinctions seem  
<sup>427</sup> to be separate from each other.

<sup>428</sup> As expected, long-short vowel pairs differ systematically in duration (Figure 5). For  
<sup>429</sup> each vowel pair, the duration densities in Figure 5 are overlapping but with two clearly  
<sup>430</sup> separable peaks (mean duration for the long vowels = 0.19 ms, SD = 0.10; mean duration  
<sup>431</sup> for the short vowels = 0.08 ms, SD = 0.09). Overall, the short vowels display less variability  
<sup>432</sup> in duration than the long vowels, a common pattern for measures with a lower bound.

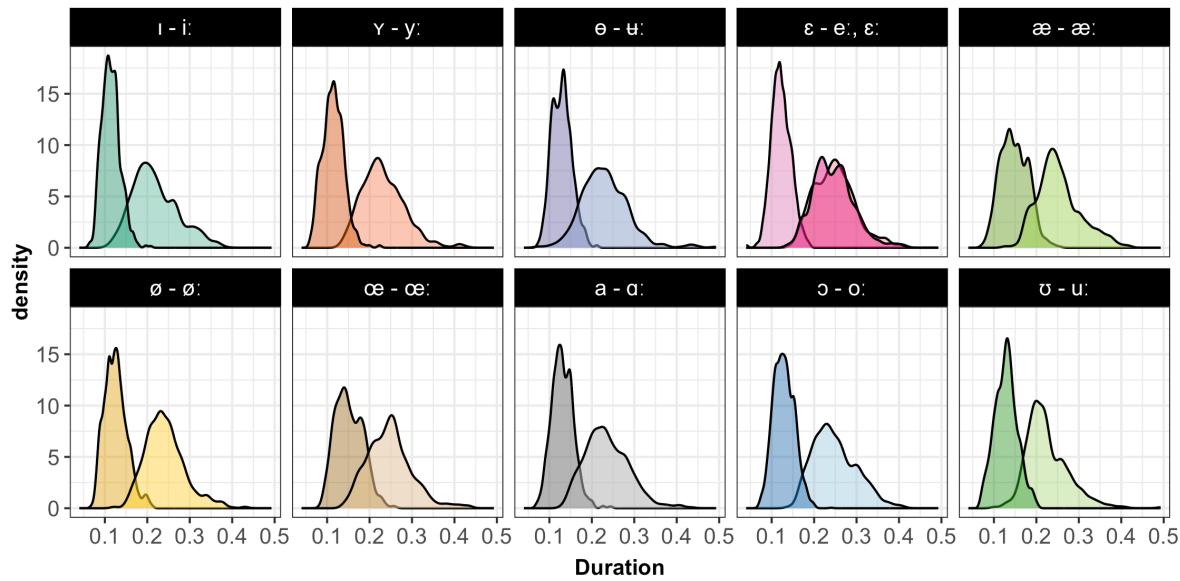
<sup>433</sup> All long-short vowel pairs furthermore display spectral differences in F1-F2. In fact,  
<sup>434</sup> as indicated in Figure 1, formant differences are apparent for *all* vowel pairs, even for vowel  
<sup>435</sup> distinctions for which duration has been found to be the primary cue—[ɛ:] - [ɛ], [ø:] - [ø], [i:]  
<sup>436</sup> - [ɪ], and [o:] - [ɔ] (e.g., Behne et al., 1997; Hadding-Koch & Abramson, 1964; Kuronen,  
<sup>437</sup> 2000). The vowel pairs that display larger spectral differences along F1-F2 seem to be [ɯ:] -  
<sup>438</sup> [ɵ] and [ɑ:] - [a] (in line with e.g., Fant, 1983; Kuronen, 2000), but also [ɛ:] - [ɛ], which  
<sup>439</sup> contrasts with previous studies. The large spectral differences in [ɛ:] - [ɛ] are presumably  
<sup>440</sup> due to [ɛ:] being produced very low in the SwehVd database, which increases the distance

---

<sup>10</sup> The alpha level for statistical significance used throughout the paper is  $p < .05$ .



*Figure 4.* The effect of including F3 in measures of category separability for the distinction between neighboring unrounded vs. outrounded vowels (**top row**), and outrounded vs. inrounded vowels (**bottom row**). **Left panels** plot the category separability for F1-F2 and F1-F2-F3 cue combinations. **Right panels** plot the proportional increase in category separability relative to F1-F2-baseline. Pointranges indicate mean and 95% bootstrapped CIs of the category separability summarized across talkers for each cue combination. Axis ranges are held constant across columns.



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

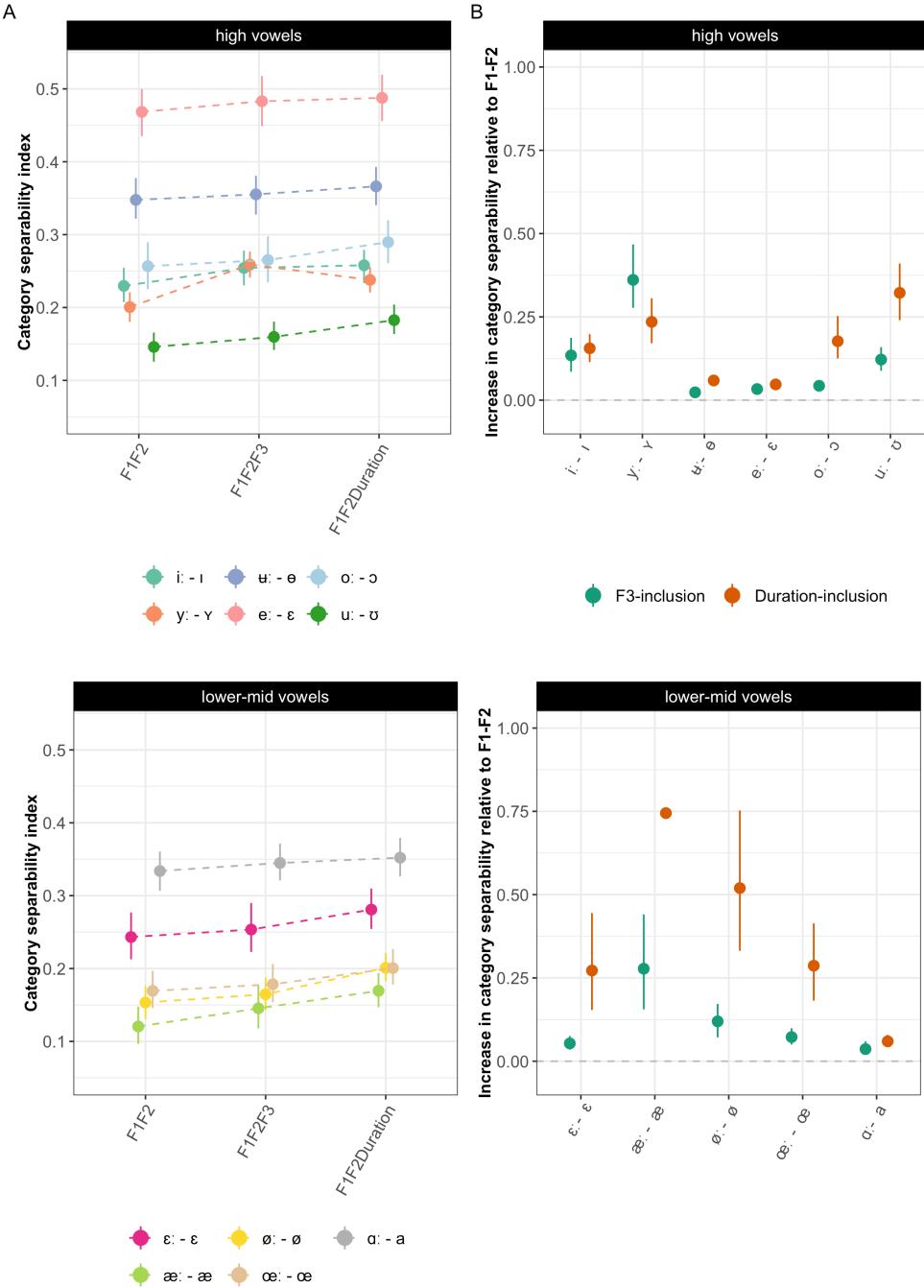
441 to [ɛ] and in addition leads to a gap along the top-left to bottom-left diagonal between [e:]  
 442 and [ɛ]. Overall, F2 appears to carry more of the spectral variation between the long and  
 443 short vowel phonemes, as categories display increased separability in the pairwise  
 444 combination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

445 In order to evaluate what cue combination would provide the largest separability  
 446 between vowels in long-short contrasts, the category separability index was calculated for  
 447 each pair and three different cue combinations: F1-F2, F1-F2-F3, or F1-F2-Duration.<sup>11</sup>

448 NULL

449 Figure 6 indicates that category separability is generally higher for some pairs. For  
 450 instance, the [e:] - [ɛ], [u:] - [ø], and [ɑ:] - [a] pairs display the largest separability when only

<sup>11</sup> For an approximation of the relative separability by cue, evaluations of F1-F3, F2-F3, F1-Duration and F2-Duration are included in the SI S1.7.2.



*Figure 6.* The effect of including F3 and duration in measures of category separability for long-short vowel pair distinctions. For visualization purposes, the pairs are split into high vowels (**top row**), and lower-mid vowels (**bottom row**). **Left panels** plot the category separability for F1-F2, F1-F2-F3 and F1-F2-duration cue combinations. **Right panels** plot the increase in category separability relative to F1-F2-baseline. Pointranges indicate mean and 95% bootstrapped CIs of the category separability summarized across talkers for each cue combination. Axis ranges are held constant across columns.

451 F1-F2 is considered (Figure 6:**column A**). The inclusion of additional cues unsurprisingly  
 452 increases category separability for all pairs, more so for pairs that are less separable in  
 453 F1-F2 space (e.g., [u:] - [u], [y:] - [y], [æ:] - [æ], [ø:] - [ø], and [œ:] - [œ]). The inclusion of  
 454 duration overall maximizes the increase in separability relative to baseline (Figure  
 455 6:**column B**). Interestingly, this is not the case for [y:] - [y] that achieves the highest  
 456 separability for the F1-F2-F3 combination. This would overall seem to suggest that the  
 457 first two formants and duration are the most important cues to long-short vowel pair  
 458 distinctions, while the inclusion of F3 unsurprisingly does not punish the separability<sup>12</sup>

459 These findings were confirmed by the statistical analysis (Summary tables in SI  
 460 S1.7.1): the LMMs indicated that the F1-F2-Duration cue combination generated the  
 461 highest separability for all pairs (all  $ps > .0001$ ), with the exception of the [y:] - [y]  
 462 contrast, for which the F1-F2-F3 combination achieved the highest separability  
 463 ( $\hat{\beta} = .058, SE = .004, p < .0001$ ). The inclusion of F3 nevertheless increased separability  
 464 relative to F1-F2 for all pairs (all  $ps > .01$ ), with the [ɛ:] - [ɛ]  
 465 ( $\hat{\beta} = .01, SE = .004, p > .009$ ) and [œ:] - [œ] ( $\hat{\beta} = .009, SE = .0034, p < .01$ ) pairs  
 466 displaying overall smaller but statistically significant differences.

467 To sum up, the results of the static analysis suggest that F1, F2 and duration are the  
 468 most important cues to vowel distinctions in Central Swedish. While visual inspection  
 469 suggested that including F3 did not substantially increase category separability for some  
 470 neighboring rounded vs. unrounded contrasts, the statistical analysis found significant  
 471 improvements in separability for all contrasts. This highlights subtle but significant  
 472 differences, and the advantages of expanding empirical analyses to modeling approaches.

473 In addition, even though all long-short vowel pairs differed systematically in duration,

---

<sup>12</sup> Analyses of additional cue combinations in the SI S1.7.2 (F1-F3, F2-F3, F1-Duration, F2-Duration) suggested that duration contributed more to separability than F3 for the [ɛ:] - [ɛ], [æ:] - [æ], [ø:] - [ø], [œ:] - [œ], and [u:] - [u] vowel pairs, while for the [u:] - [ø], [ɛ:] - [ɛ], [œ:] - [ɔ], [ɑ:] - [a] contrasts, combining one spectral cue with duration decreased separability relative to baseline, hence highlighting the reliance on both F1-F2 over duration.

474 they also displayed considerable spectral differences, suggesting that quantity  
 475 distinctions—long vs. short vowels—are not separate from quality distinctions—high, low,  
 476 front, back vowels. The comparison of how the category separability within each pair  
 477 changed as a function of cue combination furthermore highlighted the importance of F1  
 478 and F2, with F2 carrying much of the informativity for several pairs. The F1-F2-duration  
 479 combination generated the highest separability for all pairs but the [y:] - [y], where  
 480 F3-inclusion maximized separability of the cue combinations considered, highlighting the  
 481 importance of F3 for this contrast.

482 Given that the category separability index assigns equal weight to all cues included,  
 483 there is no direct way of knowing which cue contributes more to separability. Furthermore,  
 484 similar to other evaluations presented in this subsection, the separability index cannot  
 485 account for the fact that formants are not static but rather fluctuate across the signal. A  
 486 more holistic mapping of the acoustics should therefore aim to assess how formant  
 487 dynamics contribute to vowel distinctions. The next section investigates how formants  
 488 move across the segment and how much information is gained by accounting for this  
 489 dynamics.

## 490 3.2 Dynamic spectral analysis

491 This section begins with visualizations of the empirical data in formant trajectory plots.  
 492 Next, the results of the GAMMs fit to the data are presented, first focusing on the  
 493 dynamics in eight sets of neighboring vowel contrasts, and subsequently on the long and  
 494 short vowel pairs.

### 495 3.2.1 Formant movements across the space

496 Figure 7 displays the formant trajectories across all 5 time-points for the long and the short  
 497 vowels. In almost all vowels, long or short, formants showed a dynamic pattern. Only [ø:],  
 498 [i], and [y] showed very little movement over the measurement points. The scope and

499 direction, however, vary. Across vowels, the scope of movements appear to be larger  
 500 moving from vowel mid-point to 80% into the vowel, as indicated by the length of the line  
 501 from vowel label to end of arrow. Most of the formant dynamics thus take place *after*  
 502 vowel mid-point. The long high front vowels are important exceptions—the dynamics in [i:]  
 503 and [y:] mostly occur at the beginning of the vowel segment (between 20 and 50% into the  
 504 vowel), whereas [u:] displays movements of almost equal magnitude across the first four  
 505 time-points. The largest movements overall seem to concern [e:] and [o:].

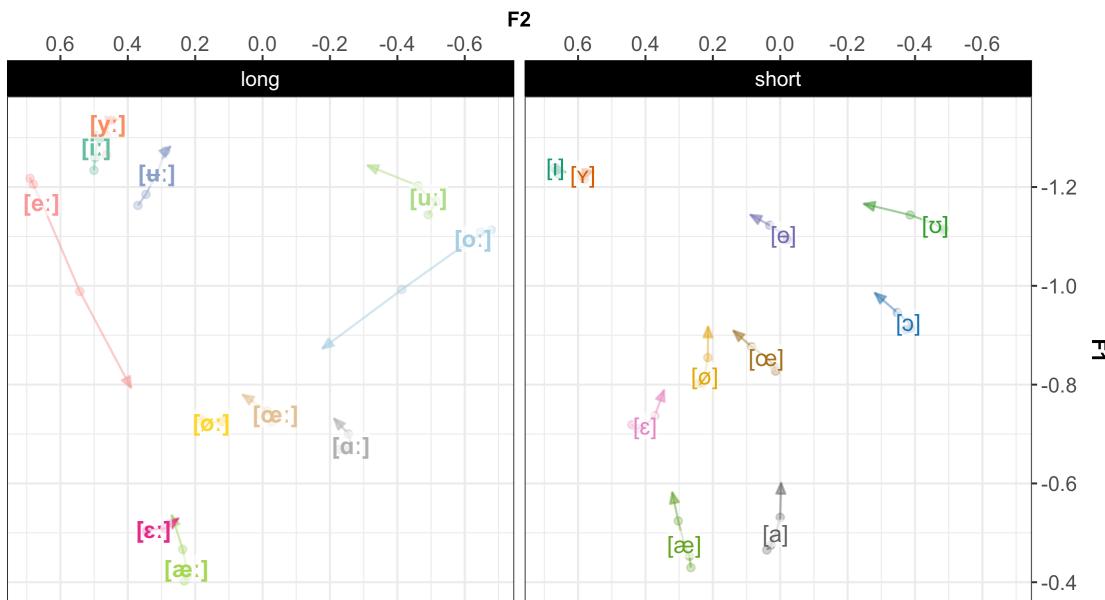
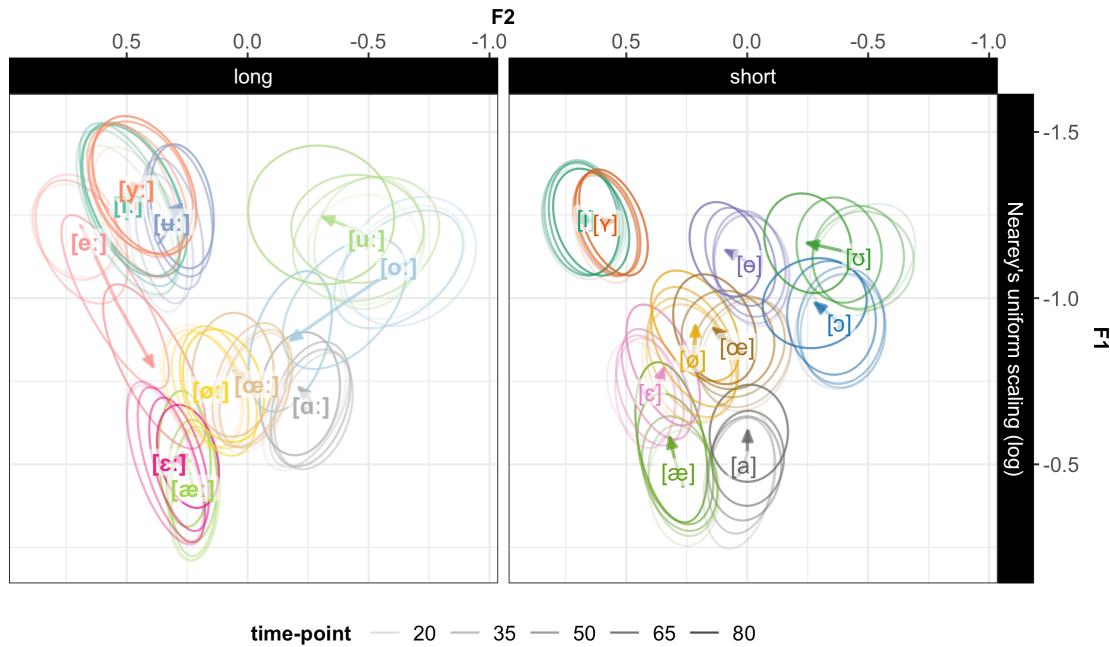


Figure 7. The trajectory of all vowels across the five time-points, along F1-F2. The arrow indicates the direction of the trajectory and ends at the final time-point, at 80% into the vowel. The vowel label is placed at the third time-point, at vowel mid-point (50%). The first (20%), second (35%) and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

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

512 more central quality, possibly indicating diphthongization ending in [ø] rather than a  
513 consonantal offglide. Parts of the movements could be due to coarticulatory effects in  
514 anticipation of the upcoming coda ([d], [d̥], [r]). If so, one would expect F2 to centralize in  
515 the later part of the segment, as tongue movements mark transitions into the alveolar (e.g.,  
516 Hillenbrand, Clark, & Nearey, 2001; Stevens & House, 1963). The formant movements  
517 along F2 from the last point (65%) to arrow tip (80%) in e.g., [e], [æ:], [œ], [ɑ:], [ɔ], [u:], and  
518 [v], might at least partly be caused by such coarticulation. Given the scope and direction of  
519 movements, Figure 7 suggests diphthongization in primarily [e:], [u:] and [o:], replicating  
520 previous work (e.g., Eklund & Traunmüller, 1997; Elert, 2000; Pelzer & Boersma, 2019),  
521 while the other vowels appear to merely display formant movement, that partly could be  
522 caused by e.g., coarticulation. The previously reported diphthongization in [ø:], however,  
523 does not seem to be particularly pronounced in these data.

524 Figure 7 further demonstrates that some neighbouring categories either converge at  
525 end points or diverge at end points. For instance, [u:] - [o:] are fairly closely located at  
526 earlier time-points, but differ substantially towards the end of the vowel segment, while  
527 [ɛ]-[æ:], [ø]-[œ] and [ø]-[œ] start at different locations but end up in approximately the  
528 same (c.f., Kuronen, 2000). Finally, the formant trajectories suggest that the empty spots  
529 identified in the vowel space under a static analysis (Figure 1), may indeed be occupied  
530 when vowel dynamics are considered. This is especially true for [e:] that travels from the  
531 mid-high front to the mid center of the space as the signal unfolds, down to a position  
532 closer to its short counterpart, [ɛ]. Given the amount of overlap when static spectral cues  
533 are considered (Figure 2), formant dynamics are likely highly informative for several of  
534 these distinctions. Figure 8 parallels Figure 2 and illustrates the effect of considering  
535 formant movements for neighboring categories. As visualized in Figure 8, the overlap  
536 between [u:] - [o:] is substantially reduced at the later time-points, while [ɛ]-[æ:], [ø]-[œ]  
537 and [ø]-[œ] are most distinguishable at earlier time-points.



*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.

### 538 3.2.2 Models of formant dynamics

#### 539 3.2.2.1 The effect of modeling formant dynamics for neighboring contrasts

540 The first set of GAMMs were modeled separately for each cue (F1, F2, F3) and each of the  
 541 sets of neighboring vowels hypothesized to (at least for some talkers) rely on formant  
 542 dynamics (Fant, 1971; Kuronen, 2000; Pelzer & Boersma, 2019): the high front vowel  
 543 contrasts [i:] - [y:] - [u:] - [e:] and the short counterpart [i] - [y], the high back vowel contrast  
 544 [o:] - [u:] and its short counterpart [ø] - [u], the lower-mid front contrast [ɛ:] - [æ:] and short  
 545 [ɛ] - [æ], and the mid center [ø:] - [œ:] and [ø] - [œ]. Summary tables of models are included  
 546 in the SI, Section S1.8.1.

547 The GAMMs fit to the high front vowel contrasts suggested significant constant and  
 548 non-linear differences over time between vowels in each contrast for all cues and contrasts

549 (for constant differences, all  $ps < .0001$ ; for non-linear differences, all  $ps < .039$ ), except for  
 550 [i:] - [y:] predicting F3 ( $\hat{\beta} = -.002, SE = .01, p > .84$ ), and [i] - [y] predicting F1  
 551 ( $\hat{\beta} = .007, SE = .007, p > .27$ ) (Figures 9 and 10). This suggests that there are no  
 552 differences in F3 for [i:] and [y:] or in F1 for [i] and [y] when formant dynamics are  
 553 considered. While F3 increases separability between [i:] and [y:] under static analysis, the  
 554 dynamics in F3 does not seem to add information about vowel quality for this contrast.

555 Significant constant as well as non-linear effects of category on all three cues was  
 556 found in all GAMMs fit to the high back long and short contrasts (for constant differences,  
 557 all  $ps < .0001$ ; for non-linear differences, all  $ps < .0197$ ). These results suggest that all  
 558 three cues contribute to distinguishing between these vowels also when formant dynamics  
 559 are considered (Figure 11).

560 The GAMMs fit to the lower-mid front long and short vowel contrasts suggested  
 561 constant and non-linear effects of vowel on F1 and F2 (for constant differences, all  
 562  $ps < .0001$ ; for non-linear differences, all  $ps < .0002$ ). For F3, there were no significant  
 563 constant differences between [ɛ:] and [æ:] ( $\hat{\beta} = -.006, SE = .007, p > .37$ ), and no constant  
 564 or non-linear differences between [ɛ] and [æ] ( $\hat{\beta} = .0099, SE = .005, p > .073$ ;  
 565  $EDF = 1, F = .879, p > .348$ ). Figure 12 demonstrates how these vowels overlap in  
 566 F3-dynamics, but are distinguished for most of the segment along F1 and F2.

567 For the GAMMs fit to the mid center vowels, there was an effect of category on all  
 568 cue evaluations for both long and short vowels (for constant differences, all  $ps < .048$ ; for  
 569 non-linear differences, all  $ps < .0001$ ; Figure 13), with the exception of both long and short  
 570 vowels fit to F1, for which no constant difference was found (for [ø:] - [œ:],  
 571  $\hat{\beta} = -.0098, SE = .007, p > .18$ ; for [ø] - [œ],  $\hat{\beta} = -.016, SE = .0097, p > .102$ ). These  
 572 contrasts nevertheless displayed non-linear differences over time (Figure 13). These results  
 573 would seem to suggest that when formant dynamics are considered, F1 does not carry  
 574 information about vowel quality for the [ø:] - [œ] and [ø] - [œ] contrasts. However, the  
 575 vowels within each pair display different non-linear patterns.

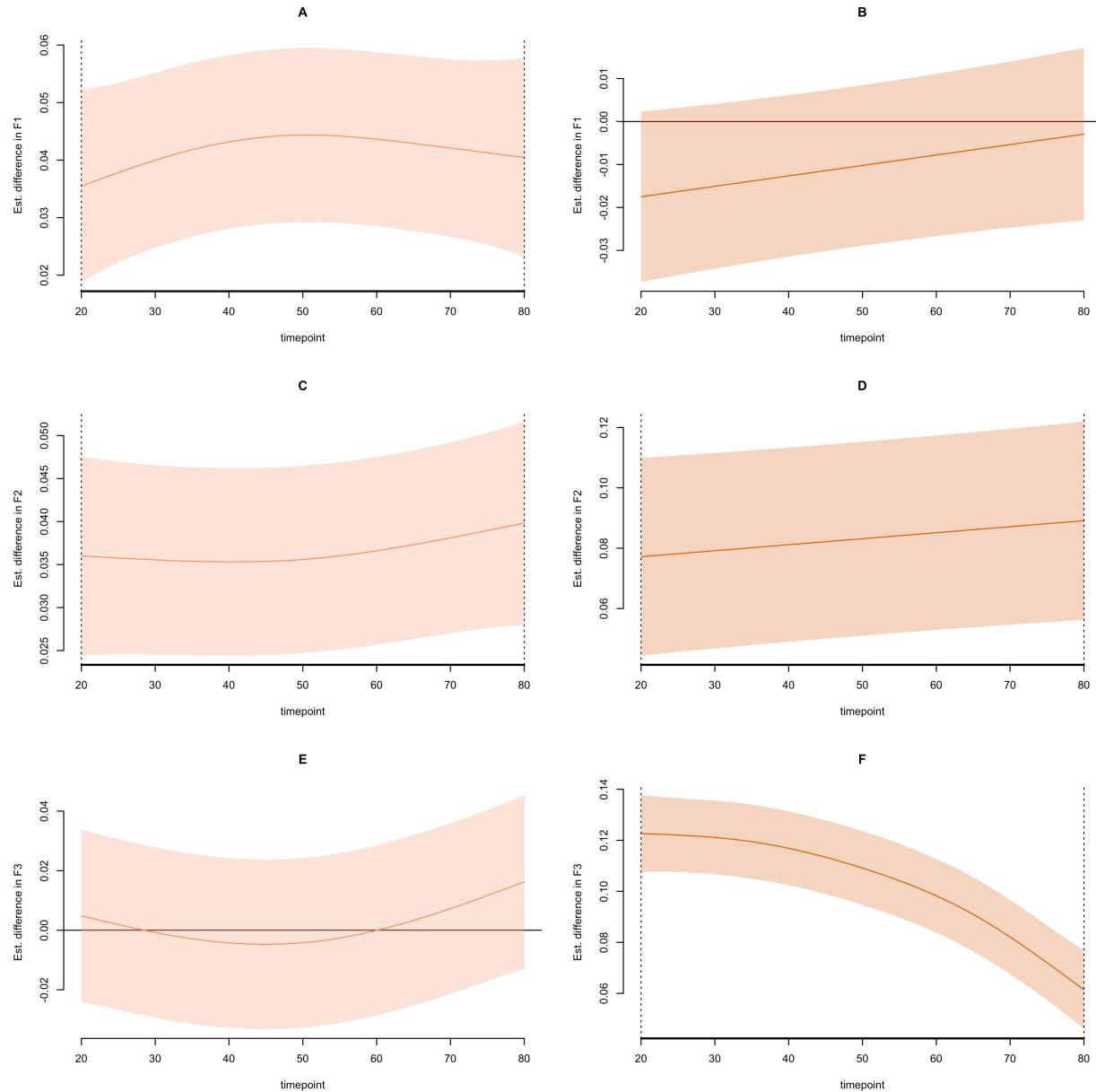


Figure 9. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the [i:] - [y:] contrast (**left**), and the [i] - [y] contrast (**right**). Differences significantly different from 0 are marked by black dotted vertical lines.

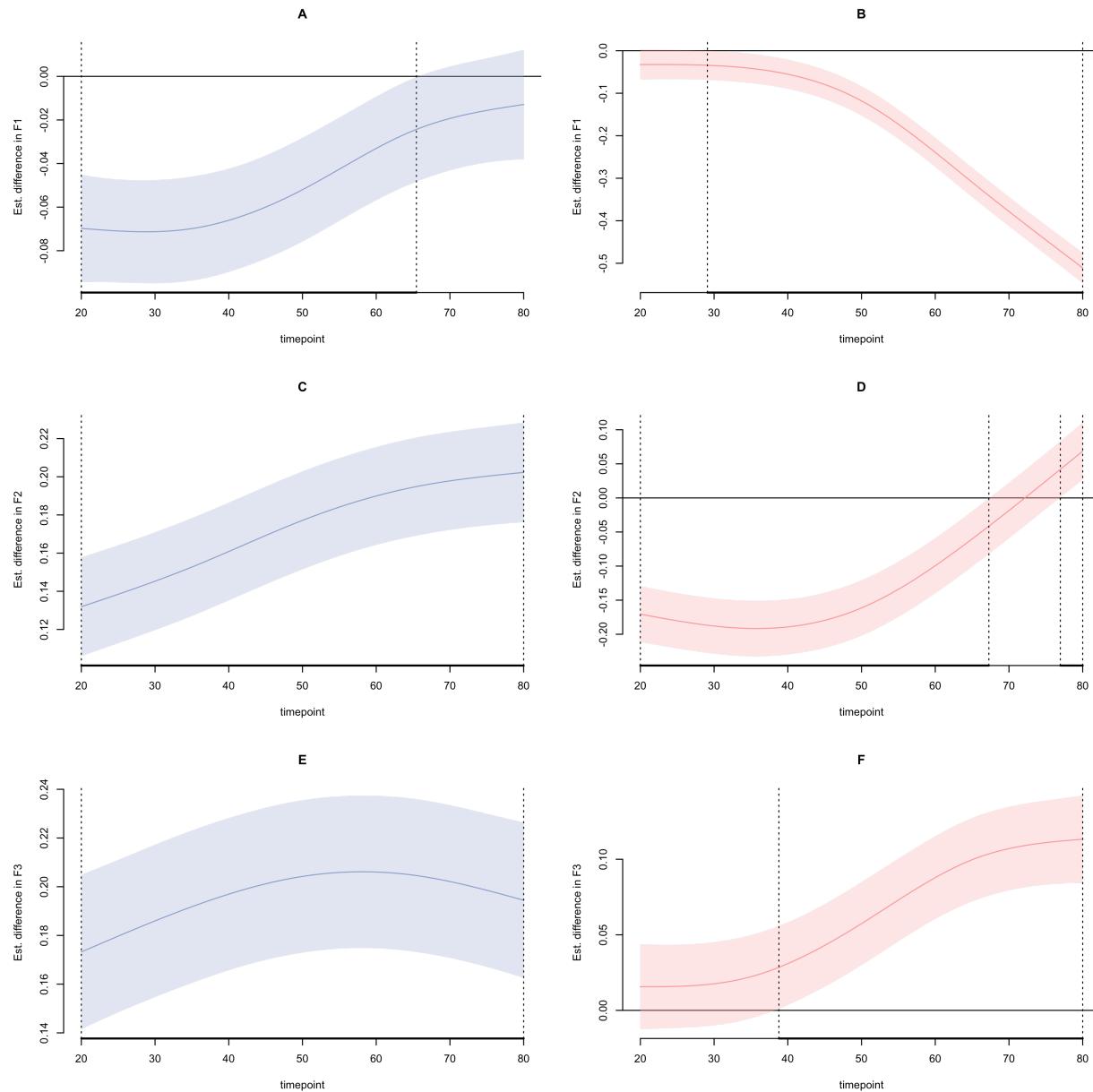


Figure 10. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the [i:] - [ɛ] contrast (**left**), and the [i:] - [e] contrast (**right**). Differences significantly different from 0 are marked by black dotted vertical lines.

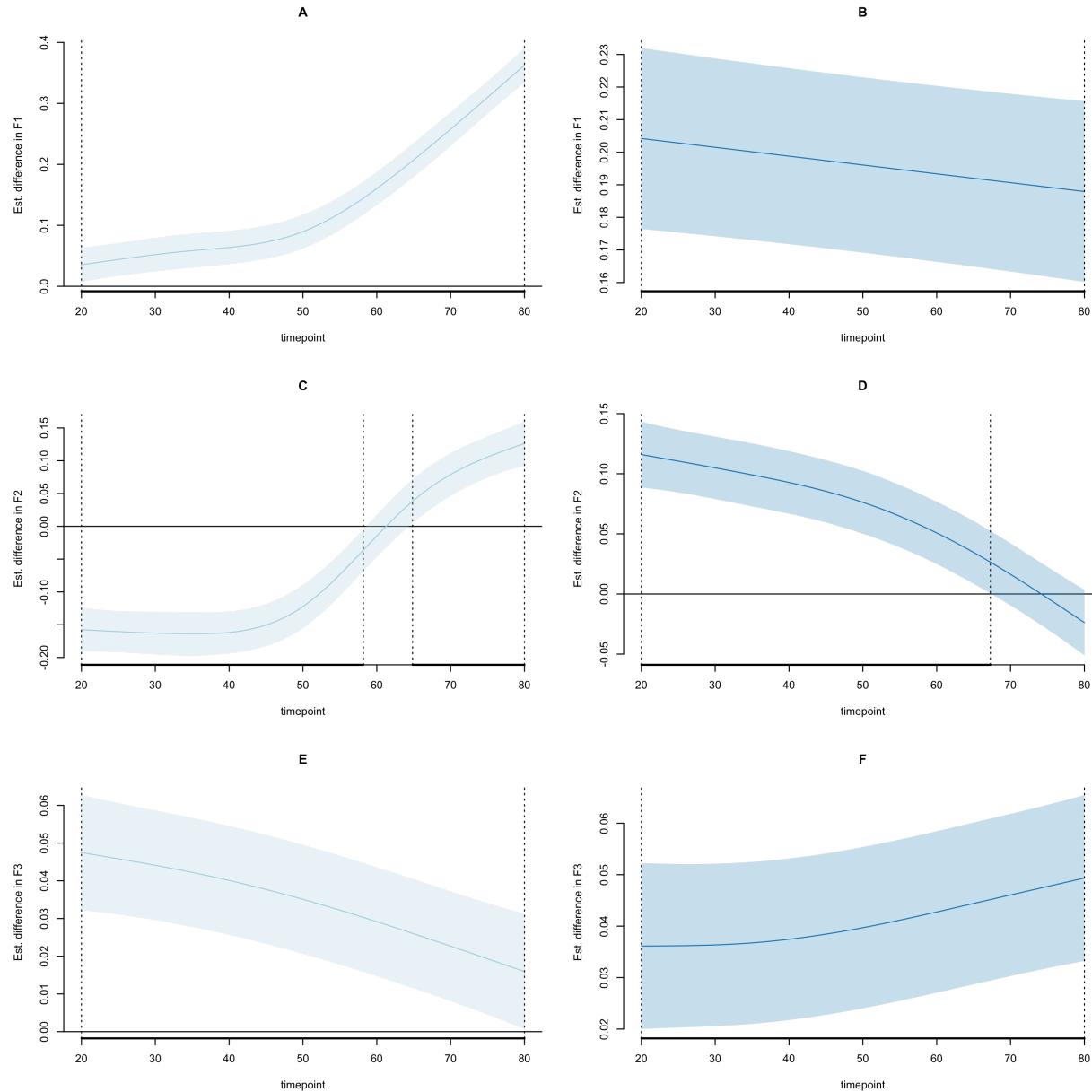
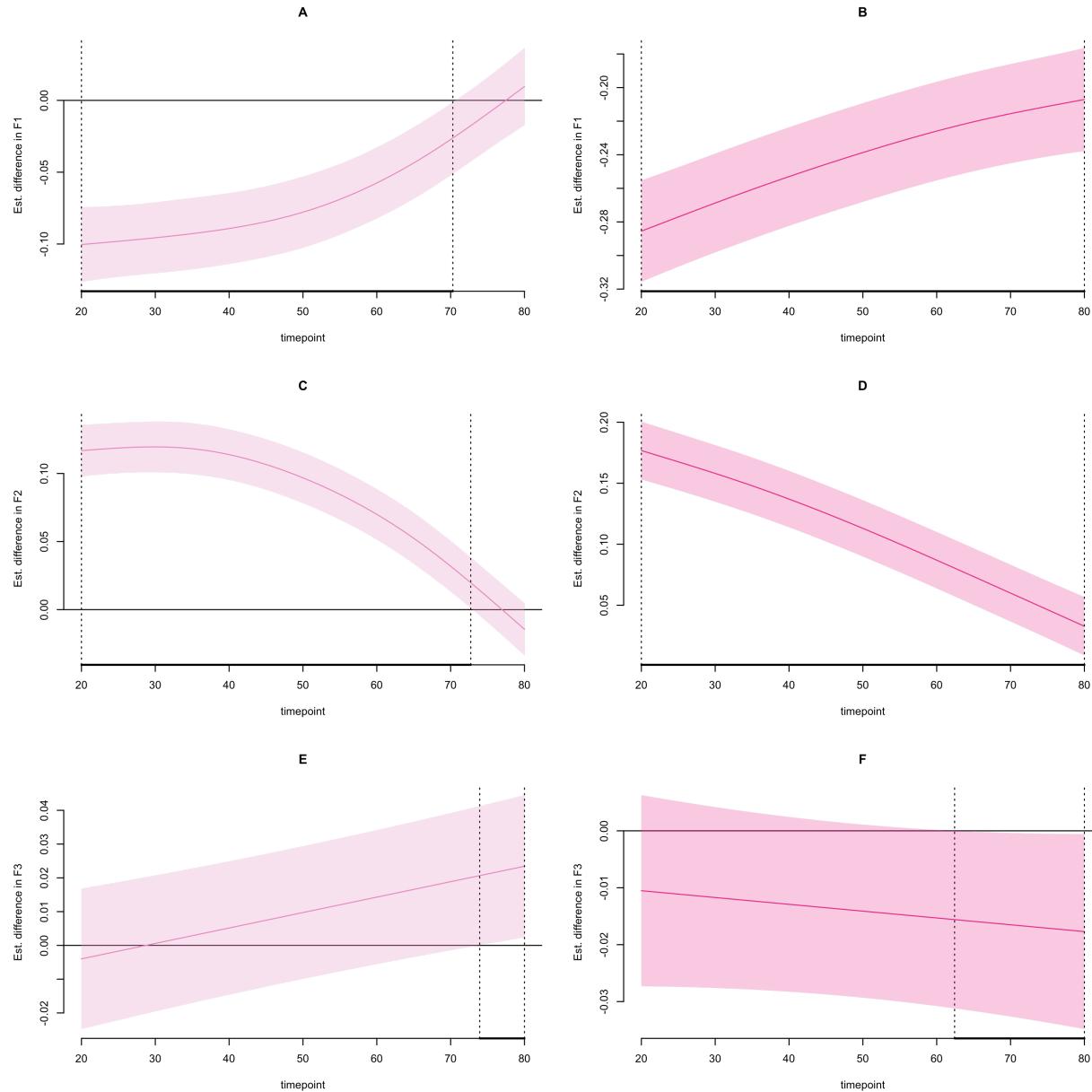
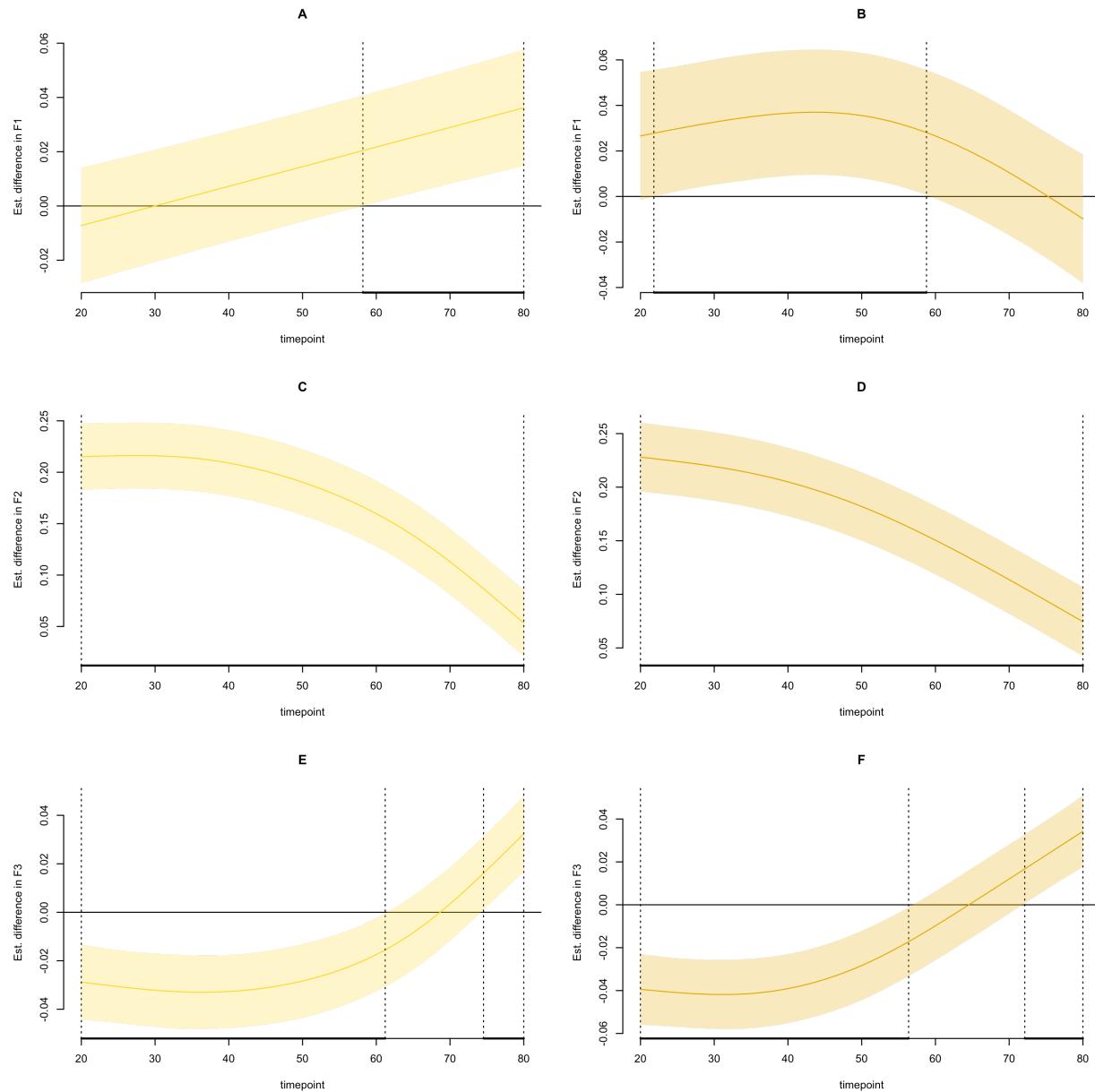


Figure 11. Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the [o:] - [u:] contrast (**left**), and the [ɔ] - [u] contrast (**right**). Differences significantly different from 0 are marked by black dotted vertical lines.



*Figure 12.* Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the [ε:] - [æ:] contrast (**left**), and the [ɛ] - [æ] contrast (**right**). Differences significantly different from 0 are marked by black dotted vertical lines.



*Figure 13.* Fitted smooths of GAMM for predicting F1 (**upper row**), F2 (**mid row**), F3 (**bottom row**) and 95% confidence intervals for the [ø:] - [œ:] contrast (**left**), and the [ø] - [œ] contrast (**right**). Differences significantly different from 0 are marked by black dotted vertical lines.

Finally, an effect of gender was found in two contrasts predicting F1: [i] - [y]

( $\hat{\beta} = -.075, SE = .026, p < .004$ ), and [i] - [u] ( $\hat{\beta} = -.06, SE = .024, p < .01$ ), indicating  
a difference in vowel height between female and male talkers. This suggests that the  
normalization approach likely reduced some talker-specificity related to anatomical  
differences, but not all. An alternative explanation is that the overall difference in height  
for these vowels when formant dynamics is considered could be attributed to a real  
pronunciation difference, driven by sociolinguistic factors (for research on young female  
talkers driving language change, see e.g., Boberg, 2019; Labov, 1990; for a review, see e.g..  
Woods, 1997).

The sets of neighboring contrasts investigated here all exhibited varying degrees of

category overlap in static analysis. However, when formant dynamics was considered, the  
vowels in each contrast were all significantly different from each other along at least two  
cues (c.f., Fant, 1971; Kuronen, 2000; Pelzer & Boersma, 2019). For some contrasts, the  
vowels overlapped only in parts of the segment, as indicated by the gaps in significant  
differences in Figures 10A-B-D-F, 11A-C-D-E, 12A-C-E, 13A-B-E-F. This indicates  
that category overlap found in static analysis is mitigated once temporal analysis is  
included, which suggests that category distinctions unfold over time. The results further  
indicated that vowel differences in dynamics were driven by both constant as well as  
non-linear differences in most cases. The contrasts for which no non-linear differences were  
found were [ɛ] - [æ] along F3 ( $\hat{\beta} = .0099, SE = .005, p > .073; EDF = 1, F = .88, p > .34$ ),  
[i] - [y] along F1 ( $EDF = 1.8, F = 1.06, p > .4$ ) and F2 ( $EDF = 1.59, F = 1.0, p > .43$ ),  
and [ɪ] - [ʏ] along F2 ( $EDF = 1, F = 3.14, p > .076$ ), which would seem to suggest their  
similarity in formant movements across the segment.

**3.2.2.2 Formant dynamics in long-short vowel pairs** The second set of GAMMs  
modeled the effect of category on F1, F2, and F3 for all long and short vowel pairs.  
Summary tables and visualizations of these GAMMs are included in the SI (Section S1.8.1).

602 There was a treatment effect of vowel on all spectral cues for all vowel pairs, driven by both  
 603 constant (for F1, all  $ps < .0067$ ; for F2, all  $ps < .002$ , for F3, all  $ps < .022$ ) and non-linear  
 604 differences (for F1, all  $ps < .0004$ ; for F2, all  $ps < .016$ , for F3, all  $ps < .05$ ), except for [ɛ:  
 605] - [ɛ̂] ( $\hat{\beta} = .00003, SE = .008, p > .96$ ) and [ɑ:] - [a] ( $\hat{\beta} = -.0043, SE = .01, p > .67$ ;  
 606  $EDF = 2.09, F = 1.67, p > .19$ ) predicting F3. This would thus seem to suggest that F3  
 607 does not reliably distinguish between long-short vowels in these pairs when dynamics is  
 608 considered. Differences in dynamics were overall driven by both constant and non-linear  
 609 differences. However, some vowel pairs displayed insignificant smooth differences between  
 610 vowels within pairs, suggesting similarity in formant movements across the segment (for the  
 611 [u:] - [ʊ] contrast predicting F2, the [æ:] - [æ] contrast predicting F1 and F3, and the [œ:] -  
 612 [œ] contrast predicting F3, all  $ps > .06$ ). With the exception of the [y:] - [ʏ] vowel pair, all  
 613 rounded vowels displayed lower F3 values in their long allophones, presumably indicating  
 614 greater lip-rounding as caused by more vigorous activity of the lips (e.g., Hadding, Hirose,  
 615 & Harris, 1976; Stålhammar, Karlsson, & Fant, 1973). Again, significant gender differences  
 616 were found in height for some high vowels (for the [i:] - [ɪ], [y:] - [ʏ] and [o:] - [ɔ] contrasts  
 617 predicting F1, all  $ps < .0086$ ), but also in backness (F2) for the [u:] - [ʊ] contrast  
 618 ( $p < .0185$ ).

619 The results suggest that when formant dynamics are considered, *all* long-short  
 620 Central Swedish vowel pairs differ in spectral cues. Among the pairs that displayed smaller,  
 621 albeit statistically significant, differences in spectral cues is the [i:] - [ɪ] pair predicting F1,  
 622 possibly indicating a tendency for stronger duration dependency, in line with previous  
 623 perceptual work (Behne et al., 1997). Overall, larger effects were found for F2 which would  
 624 support the hypothesis that F2 carries more of the durational variation in quantity  
 625 contrasts in Swedish. Non-significant differences within pairs were found only for GAMMs  
 626 predicting F3, which highlights the primary importance of F1-F2 as spectral cues to  
 627 quantity contrasts under the assumption of formant dynamics.

### 628 3.3 Results summary

629 The results from the static analysis suggest that F1, F2, F3 and duration all contribute to  
630 the distinction of Central Swedish vowels. While F1 and F2 are of primary importance for  
631 most contrasts, F3 contributes to increasing the separability between neighboring vowels  
632 differing in lip-rounding. The static analysis further suggested that while the  
633 F1-F2-duration cue combination maximized separability for long-short vowel pairs, the  
634 inclusion of F3 also increased separability relative to F1-F2. Importantly, including F3  
635 increased separability more than duration did for the [y:] - [y] contrast.

636 Some vowels displayed overlap in the static analysis but increased separability when  
637 formant dynamics was considered, as indicated by formant trajectory analysis and  
638 GAMMs. The dynamic analyses highlighted that the short vowels also display formant  
639 movements, and that for most of both long and short categories, a larger portion of the  
640 dynamics resides in the later part of the segment. Given the increased separability of  
641 neighboring contrasts found in dynamic analysis, it is reasonable to assume formant  
642 movements as an auxiliary cue to vowel identity, more so for some contrasts than others.  
643 For instance, the [i:] - [y:] - [u:] - [e:], [i] - [y], [o:] - [u], [ɔ] - [ʊ], [ɛ:] - [æ:], [ɛ] - [æ], [ø:] - [œ:],  
644 and [ø] - [œ] contrasts displayed considerable overlap in static analyses but increased  
645 distinguishability when analysed dynamically.

646 While the static analysis suggested increased separability for the [i:] - [y:] contrast  
647 when F3 was included, the GAMM fit to the same contrast found no significant differences  
648 between the two vowels predicting F3, suggesting less distinguishability under the  
649 assumption of formant dynamics. However, the GAMM fit to [y:] - [y] suggested  
650 statistically significant differences between the two categories predicting F3. These two  
651 analyses taken together suggest more effects on F3 in the short vowel compared to the long  
652 vowel.

653 The resulting phonetic characteristics of the long and short Central Swedish vowels

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]	[a]	

presented here, is summarized in Table 2. Beginning with the long vowels, there are 4 high vowels. The current acoustic description suggests that none of them are front. Instead, [i:] and [y:] group with [ɯ:] as central vowels, and [u:] is back (c.f., Pelzer & Boersma, 2019; Schötz et al., 2011).<sup>13</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).

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

The motivation of the summarized acoustics presented in Table 2 rests on a pairwise grouping of the long and short vowels, similar to phonological analyses of Central Swedish (Riad, 2014). For instance, despite its high position, [e:] is defined as mid-high, on par with [o:], as both vowels share diphthongizational patterns, and their short versions are both lower than their long counterparts. Furthermore, [æ:] is front rather than central as its short version is clearly more front than [ø]. Because of their overall centralized positions in SwehVd, [i:] and [y:] groups with [ɯ:], while their short versions are still clearly front. One

<sup>13</sup> Whether [y:] can still be considered rounded given the high F3 values, is a question for future research.

674 could thus argue that both /i/ and /y/ are under-specified for the front-back dimension,  
675 which would also be the case for /ɑ/, as [ɑ:] is back and [a] is central. It is important to  
676 note that while Table 2 may point to possible updates of Central Swedish vowel phonology,  
677 this is only tentative as more evidence is required for a definite update. These include, e.g.,  
678 more investigations in different contexts.

## 679 4 General discussion

680 The purpose of this paper was to present up-to-date static and dynamic acoustic analyses  
681 of Central Swedish vowels that included both empirical formant data and models of  
682 formant dynamics. The study thus aimed to expand on and complement previous work by  
683 1) the scope of the analysis, performing the same type of analyses on all 21 categories, 2)  
684 the materials chosen, using a recently collected hVd corpus with high resolution within and  
685 across talkers for a single variety, and 3) the methodological approach employed, with the  
686 use of traditional formant analysis, category separability index, trajectory visualizations  
687 and models of formant dynamics (GAMMs). The study further aimed to evaluate the  
688 hypothesized importance of F3 for contrasting rounded vs. unrounded categories, the  
689 extent to which *all* long-short vowel pairs display spectral differences, and what part of the  
690 vowel space is more susceptible to diphthongization. Next, the main findings are discussed,  
691 alongside methodological considerations and future directions.

692 Beginning with the static analysis, the results seem to suggest that the most  
693 important cues to vowel identity are F1, F2 and duration, which replicates previous work  
694 (e.g., Kuronen, 2000; Lindblom, 1963). Including F3 increased separability for all rounding  
695 contrasts which suggests that F3 adds information for rounding vs. unrounded neighboring  
696 contrasts, given the way category separability is calculated here. The category separability  
697 index for the long-short vowel pairs suggested that the F1-F2-duration cue combination  
698 maximized separability between vowels within pairs, highlighting the role of both spectral

and temporal cues for quantity contrasts and that *all* long-short vowel pairs display spectral differences. The results further indicated that F2 seemed to carry more of the spectral variation in quantity contrasts (Kuronen, 2000; Lindblom, 1963). Of note, however, one long-short vowel pair, [y:] - [y], achieved higher separability with F3-inclusion over Duration-inclusion, despite the fact that the two vowels clearly display systematic differences in duration similar to all other pairs (Figure 5). Across talkers, [y] is clearly more front than [y:] (mean F2 for [y] = 2201 Hz, SD=226; mean F2 for [y:] = 1939 Hz, SD=187), yet F3 is overall higher for [y:] (mean F3 for [y:] = 3058 Hz, SD=259; mean F3 for [y] = 2625 Hz, SD=226). Under the assumption of a lower F3 in high front rounded vowels being indicative of increased lip-rounding (compare, mean F3 for [i:] = 3054 Hz, SD=313; mean F3 for [i] = 2924 Hz, SD=237), this would seem to suggest more rounding in [y] relative to [y:].

The results furthermore suggested that static measurements across talkers, while informative, were insufficient to accurately capture the spectral acoustics of some vowels, as indicated by the amount of overlap when static formants were considered. A more exhaustive description can be achieved by including dynamic analyses, as shown in this paper. This was especially true for the [ɛ:] - [æ:], [u:] - [o:] and [e:] - [i:] - [y:] distinctions, given the direction and scope of formant trajectory movements across the vowel segment. The analysis of the empirical data furthermore suggested that a larger portion of the movement took place at the later part of the segment for most vowels, irrespective of the magnitude of change in F1 and F2. However, [i:], [y:] and [ɛ:] constituted important exceptions as they displayed more movement in the first three to four time-points. Of note, some of the short vowels showed formant movements of equal or larger magnitude as certain long vowels, which seems to signal the importance of vowel dynamics for long and short vowels alike. This has been investigated in work on other languages (e.g., Hillenbrand et al., 1995; Watson & Harrington, 1999), but has largely been lacking in studies on Swedish. In terms of distinguishing between diphthongization and merely formant

726 movement, the trajectory plots suggested diphthongization towards an open quality in  
727 primarily [e:] and [o:].

728 GAMMs fit to the data contributed with further insights into the formant dynamics  
729 of individual contrasts as well as for vowels differing in quantity, and allowed for an  
730 assessment of the relative contribution of formant dynamics to cue dependencies for  
731 neighboring and more distant contrasts. For instance, GAMMs fit to the long-short vowel  
732 pairs indicated that within all pairs, the two categories were significantly different in  
733 predicting the spectral cues in all cue evaluations, with the exception of [ɛ]-[ε] and [ɔ]-[ɑ]  
734 predicting F3. This highlights the informativity carried by formant dynamics for quantity  
735 distinctions. The static and dynamics analyses thus both suggest that quantity distinctions  
736 are not separate from quality distinctions in Central Swedish.

737 With regard to the neighboring contrasts hypothesized to rely on formant  
738 dynamics—the [i:] - [y:] - [ɯ] - [e:], [o:] - [u:], [ɛ:] - [æ:], and [ɔ:]-[œ:] contrasts, and their short  
739 counterparts—the results indicated significant differences for all comparisons and all cues  
740 with the exception of [i:] - [y:], [ɛ:] - [æ:], and [ε] - [æ] predicting F3, and [i] - [y], [ɔ:]-[œ:]  
741 predicting F1. This would seem to suggest that the movements in these vowels along these  
742 cues, are not contributing to vowel quality information. For the [i:] - [y:] contrast, including  
743 F3 under static measures either did not appear to change separability (visualizations), or  
744 increased the category separability index relative to F1-F2. However, the F1-F2-F3 index  
745 for [i:] - [y:] displayed substantial between-talker variability, and was comparatively low.  
746 These results overall seem to indicate that the [i:] - [y:] contrast might primarily be  
747 supported by F1-F2 dynamics, or by additional acoustic cues not investigated in this study.  
748 There is of course also the possibility that listeners might disambiguate the two categories  
749 using primarily visual cues or linguistic information, or perhaps these two categories are  
750 not distinguished, which would suggest a merger in process among some of these talkers. A  
751 merger might be driven by relaxation of lip-rounding, as supported by the higher F3 values  
752 found for [y:].

753 The vowel space as summarized in Table 2 suggests shifts in some vowels compared to

754 previous characterizations of Central Swedish. For instance, the [ɛ:] has lowered

755 substantially compared to earlier mappings of the space (e.g., Engstrand, 1999; Fant et al.,

756 1969; Kuronen, 2000). This lowering was anticipated by Riad (2014) and supported in

757 Leinonen (2010), Gross et al. (2016), and Pelzer and Boersma (2019). Another important

758 shift concerns the fronting of [i] and [y] and the centralization of [i:] and [y:], where [i] and

759 [y] appear to maintain their positions as high front vowels, while [i:] and [y:] have

760 centralized to the mid-center part of the space (see also, Pelzer & Boersma, 2019).<sup>14</sup> This

761 conflicts with previous work on Swedish as well as other languages demonstrating that

762 short vowels are generally less peripheral and more centralized than their long counterparts

763 (e.g., Clopper et al., 2005; Hillenbrand et al., 1995). Besides the already mentioned

764 hypothesis of a possible [i:] - [y:] merger due to relaxation of lip-rounding in [y:], an

765 alternative, yet related, hypothesis is that both, or either, vowels are produced as damped

766 versions, as previously found for talkers of Stockholm Swedish (e.g., Kotsinas, 1994; Schötz

767 et al., 2011). The presence of a damped [i:] would be supported by the lower F2 values, as

768 the consonantal offglide in [i:] lowers F2 (Engstrand et al., 2000). A merger of [i:] and [y:]

769 into [i] has been observed among younger talkers in other regions in Sweden, e.g., for

770 Gothenburg Swedish, as reported by Gross and Forsberg (2020). If centralization of [i:] is a

771 prerequisite for such a merger, as suggested by Gross and Forsberg (2020), the present

772 results might indicate the beginning of a merger. Impressionistic listening by the author

773 did support the presence of a final consonantal glide, similar to [i:], among the majority of

774 talkers in SwehVd, both male and female. The strength and scope of [i:] varied across

775 talkers, from a relatively strong [i:] (18 talkers), to a weaker voiced fricative offglide buzzing

776 [z] (14 talkers), or more of a consonant offglide element similar to [j] following [i:] (3

777 talkers). Interestingly, several of the talkers that did not produce any apparent final

<sup>14</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.

778 consonant glide or buzz (12 talkers), seemed to have overall less retracted [i:] and [y:], hence  
779 supporting the hypothesized link between centralization and consonantal offglide.

## 780 4.1 Methodological considerations and future directions

781 Two sets of methodological considerations for the present work not mentioned elsewhere,  
782 deserve further discussion. The first set concerns the measuring and evaluation of two of  
783 the acoustic cues, F0 and vowel duration. F0 is not considered an important cue to vowel  
784 identity in itself and therefore often not included in previous work (e.g., Fant et al., 1969;  
785 Gross et al., 2016; Leinonen, 2010; Pelzer & Boersma, 2019; Wenner, 2010). However, it is  
786 known to vary across languages and dialects (e.g., Henton, 2005; Jacewicz & Fox, 2018;  
787 Johnson, 2005; Leung et al., 2016; Mennen et al., 2012; Weirich et al., 2019), and has been  
788 shown to have strong indirect effects on vowel categorization (Barreda & Nearey, 2012; see  
789 also work on vowel-intrinsic F0, e.g., Whalen & Levitt, 1995; and the hypothesized use of  
790 F0 for tense-lax distinctions in German, Pape & Mooshammer, 2006). For the above  
791 reasons, analysis of F0 was included in this study but limited to reports of mean F0 across  
792 the segment for static visualizations. Figure 3 confirmed that F0 carried the least  
793 information about vowel identity of all cues included in this material. However, given  
794 evidence of pitch contours influencing perceived duration and prominence of vowels (e.g.,  
795 Gussenhoven & Zhou, 2013), one could claim that a more comprehensive acoustic  
796 representation of vowels should have included F0 in the dynamic analysis. It is not  
797 unreasonable to assume a role for F0 in crowded systems such as the Central Swedish vowel  
798 inventory with category overlap and potential mergers (c.f., research on adaptive dispersion,  
799 e.g., Liljencrants & Lindblom, 1972; Lindblom, 1998). For instance, talkers might combine  
800 tone distinction (F0) with voice quality (e.g., creak, or buzzing, as in the damped [i]) to  
801 increase distinctions between neighboring contrasts (for a review, see e.g., Davidson, 2021).

802 For duration, the temporal analysis of the effect of duration on long-short vowel pairs  
803 could have been supplemented with measures of consonant ratios, following, e.g., Pelzer

804 and Boersma (2019), and Schaeffler (2005). Since the SwehVd database is publicly  
805 available in an online repository (<https://osf.io/ruxnb/>), both the question of consonant  
806 ratios and the role of F0 can be addressed in future studies.

807 The second methodological consideration concerns methods used for assessing  
808 category distinguishability. Since the GAMMs were fit to each cue separately, they allowed  
809 for separate assessments of the relative weight of each cue, compared to the separability  
810 index where the by-cue contributions could only be evaluated indirectly. An inherent  
811 limitation in the separability index, as implemented in this paper, is the simplifying  
812 assumption that all dimensions within a cue combination carry equal weight. It is therefore  
813 not possible to assess whether the relation between the cues in each space is symmetrical or  
814 not, that is, in the F1-F2 space, we do not know whether F1 carries as much information as  
815 F2 for separability, and vice versa. In addition, given that the comparisons are pairwise,  
816 they are limited to explaining the relation between two vowels in a contrast. As such, they  
817 cannot inform us of the separability of a given vowel from other neighboring vowels, or the  
818 overall category separability in the entire space. This is, however, a limitation that the  
819 separability index shares with the GAMMs. Neither of these methods are able to assess the  
820 distinguishability or confusability of all vowels under different cue combinations in one  
821 analysis. Nor can they inform us of the *perceived* distinguishability, even if it is reasonable  
822 to assume that reduced overlap between tokens of neighboring categories would increase  
823 intelligibility (e.g., Bradlow, 1995; Wright et al., 2004). For this, one could fit other types  
824 of models, such as a multinomial logistic regression model predicting vowel category from  
825 cue combinations, or perceptual models assessing the predicted consequences for  
826 perception. In separate work conducted in parallel with this study, we have pursued a  
827 similar approach, evaluating the predicted consequences of F3-inclusion for high-front  
828 vowel contrasts in Swedish using a perceptual model based on Bayesian inference, ideal  
829 observers (Persson & Jaeger, 2024). The results of including F3 qualitatively replicated the  
830 results presented here for the investigated vowels in that it overall improved the predicted

831 recognition accuracy. Compared to measures of category separability, the ideal observers  
832 allowed for an assessment of the effect of F3-inclusion on category confusability among all  
833 vowels considered. This analysis suggested that while F3-inclusion overall decreased  
834 category confusability, especially for [y:] and [y], it increased the probability of confusing [i:]  
835 with [y:].<sup>15</sup>

836 Models predicting perception from production data can further inform the design of  
837 perception studies that can shed more light on the consequences of the present results for  
838 the perception of Central Swedish vowels. For instance, in a language with a systematic  
839 quantity distinction such as Swedish, the role of spectral cues in long-short vowel pair  
840 distinctions could be assessed by exposing listeners to synthesized versions where long and  
841 short vowel duration is crossed with the allophones' spectral information for any given  
842 phoneme. Furthermore, as the results seem to support claims of the hypothesized  
843 importance of formant dynamics for vowel distinctions, more insight into the effect of  
844 formant dynamics for vowel perception could be gained from having listeners categorize  
845 tokens extracted from different segments of the long vowels, e.g., the first three time-points  
846 vs. the three final time-points (c.f., Jenkins, Strange, & Miranda, 1994; Strange, 1989).  
847 The design of such experiments can be informed by modeling the predicted perceptual  
848 consequences of different cue spaces, and of considering different vowel segments.

849 Another avenue for future research to explore, are the preliminary hypotheses  
850 mapped out in this section concerning vowel change. Since this study's primary focus is  
851 mapping the acoustics of modern-day Central Swedish vowels, systematic investigations of  
852 the underlying reasons to the potential shifts in the vowel space compared to previous  
853 work, are left for future studies. These hypotheses, and others, could be investigated in  
854 perceptual studies or by using systematic listening by trained phoneticians, further

<sup>15</sup> In comparison to other model-based approaches such as logistic regression and linear discriminant analysis, Bayesian ideal observers have the advantage of reducing the number of degrees of freedom in the fit from production data to predicting perception, which substantially reduces the risk of over-fitting to the data (see discussion in e.g., Persson & Jaeger, 2023).

855 validated through measures of inter-rater reliability (for a review, see Cucchiarini, 1995;  
856 Gross & Forsberg, 2020; Kuronen, 2000; Pelzer & Boersma, 2019). Evaluations by trained  
857 phoneticians can provide assessments of both auditory distinguishability between categories  
858 and the correspondence between assigned IPA label and auditory impression. Perceptual  
859 studies with naïve listeners can complement evaluations by phoneticians, as training can  
860 differ across phoneticians and introduce biases (for reviews, see e.g., Heselwood & Howard,  
861 2008; Kerswill & Wright, 1990; Stemberger & Bernhardt, 2020). Relatedly, the extent to  
862 which individual talkers are driving these changes could be investigated by assessing the  
863 amount of cross-talker variability in SwehVd. This study suggests that even when keeping  
864 the background variables age and region of origin constant and normalizing for  
865 talker-specificity in physiology, there are still differences in the phonetic realization of some  
866 vowels that could be related to sociolinguistic differences. Further insights into these  
867 individual differences can be gained by studying the SwehVd materials on a talker-specific  
868 level.

## 869 5 Conclusions

870 The present study has reported on the acoustic properties of Central Swedish vowels. The  
871 spectral and temporal cues investigated all contributed to distinguishing between the 21  
872 vowels in the Central Swedish vowel space, with varying weight. More insight into formant  
873 dynamics within and between quantities have been gained by the dynamic analysis  
874 presented, which is also of value for cross-linguistic research. What has been gained with  
875 the broad-scale approach of characterizing the *entire* vowel space adopted here, is of course  
876 lost in terms of detailed investigations of individual vowel contrasts. There is certainly a  
877 lot more to say about the centralization of [i:] - [y:], the potential relaxation of lip-rounding  
878 in [y:], the lowering of [ε:], and the potential role of additional cues beyond those  
879 investigated here, among other things. The acoustic descriptions outlined in this paper,

880 together with the publicly available SwehVd database, can provide a reference point for  
881 future investigations into these acoustic events and beyond.

## 900 6 References

- 901 Assmann, Peter F., & William F. Katz. (2005). Synthesis fidelity and time-varying  
902 spectral change in vowels. *The Journal of the Acoustical Society of America*, 117(2),  
903 886–895. <https://doi.org/10.1121/1.1852549>
- 904 Baayen, Harald, Shravan Vasishth, Reinhold Kliegl, & Douglas Bates. (2017). The cave of  
905 shadows: Addressing the human factor with generalized additive mixed models.  
906 *Journal of Memory and Language*, 94, 206–234.
- 907 Barreda, Santiago. (2021). Perceptual validation of vowel normalization methods for  
908 variationist research. *Language Variation and Change*, 33(1), 27–53.  
909 <https://doi.org/10.1017/S0954394521000016>
- 910 Barreda, Santiago, & Terrance M. Nearey. (2012). The direct and indirect roles of  
911 fundamental frequency in vowel perception. *The Journal of the Acoustical Society of  
912 America*, 131(1), 466–477. <https://doi.org/10.1121/1.3662068>
- 913 Barreda, Santiago, & Terrance M. Nearey. (2018). A regression approach to vowel  
914 normalization for missing and unbalanced data. *The Journal of the Acoustical Society  
915 of America*, 144(1), 500–520. <https://doi.org/10.1121/1.5047742>
- 916 Behne, Dawn M., Peter E. Czigler, & Kirk P. H. Sullivan. (1997). Swedish Quantity and  
917 Quality: A Traditional Issue Revisited. *Reports from the Department of Phonetics,  
918 Umeå University*, 4, 81–83.
- 919 Björsten, Sven, & Olle Engstrand. (1999). Swedish “damped” /i/ and /y/: Experimental  
920 and typological observations. *Proceedings of the 14th International Congress of  
921 Phonetic Sciences*, 1957–1960. San Francisco.
- 922 Bleckert, Lars. (1987). *Centralsvensk diftongering som satsfonetiskt problem*  
923 [*Diphthongization in Central Swedish as a problem of sentence phonetics*]. Uppsala:  
924 Skrifter Utgivna Av Institutionen För Nordiska Språk Vid Uppsala Universitet.
- 925 Boberg, Charles. (2019). A closer look at the short front vowel shift in Canada. *Journal of  
926 English Linguistics*, 47(2), 91–119. <https://doi.org/10.1177/0075424219831353>

- 927 Boersma, Paul, & David Weenink. (1992–2022). *Praat: Doing phonetics by computer*  
928 [Computer program].
- 929 Bradlow, Ann R. (1995). A comparative acoustic study of English and Spanish vowels.  
930 *The Journal of the Acoustical Society of America*, 97(3), 1916–1924.
- 931 Bruce, Gösta. (2009). Components of a prosodic typology of Swedish intonation. In  
932 *Components of a prosodic typology of Swedish intonation* (pp. 113–146). De Gruyter  
933 Mouton. <https://doi.org/10.1515/9783110207569.113>
- 934 Bruce, Gösta. (2010). *Vår fonetiska geografi : Om svenska accenterna, melodi och uttal*  
935 [*Our phonetic geography: about the accents, melody and pronunciation of Swedish*]  
936 (första upplagan). Lund: Studentlitteratur.
- 937 Chesworth, Janine, Kim Coté, Colleen Shaw, Sandra Williams, & Megan Hodge. (2003).  
938 Effect of phonetic context on women's vowel area. *Canadian Acoustics/Acoustique*  
939 *Canadienne*, 31, 20–21.
- 940 Chuang, Yu-Ying, Janice Fon, Ioannis Papakyritsis, & Harald Baayen. (2021). Analyzing  
941 phonetic data with generalized additive mixed models. In *Manual of clinical phonetics*  
942 (pp. 108–138). Routledge.
- 943 Clopper, Cynthia G., David B. Pisoni, & Kenneth De Jong. (2005). Acoustic  
944 characteristics of the vowel systems of six regional varieties of American English. *The*  
945 *Journal of the Acoustical Society of America*, 118(3), 1661–1676.
- 946 Cuccharini, Catia. (1995). Assessing transcription agreement: Methodological aspects.  
947 *Clinical Linguistics and Phonetics*, 10(2), 131–155.  
948 <https://doi.org/10.3109/026992096089851670269-9206>
- 949 Davidson, Lisa. (2021). The versatility of creaky phonation: Segmental, prosodic, and  
950 sociolinguistic uses in the world's languages. *Wiley Interdisciplinary Reviews: Cognitive*  
951 *Science*, 12(3), e1547.
- 952 Eklund, Ingegerd, & Hartmut Traunmüller. (1997). Comparative Study of Male and  
953 Female Whispered and Phonated Versions of the Long Vowels of Swedish. *Phonetica*,

- 954 54(1), 1–21. <https://doi.org/10.1159/000262207>
- 955 Elert, Claes-Christian. (1964). *Phonologic studies of quantity in Swedish: Based on*  
956 *material from Stockholm speakers*. Uppsala: Almqvist & Wiksell.
- 957 Elert, Claes-Christian. (1980). Diftongeringar och konsonantinslag: Drag i uttalet av långa  
958 vokaler i svenska av i dag [Diphthongizations and consonantal offglides: Traits in the  
959 pronunciation of long vowels in the Swedish of today]. *Språken i vårt Språk*, 168–181.
- 960 Elert, Claes-Christian. (1981). *Ljud och ord i svenska [Sounds and words in Swedish*  
961 *language]*. Umeå: Universitetet i Umeå, Almqvist & Wiksell international.
- 962 Elert, Claes-Christian. (1994). Indelning och gränser inom området för den talade  
963 svenska: En aktuell dialektografi [Distribution and boundaries within the area of  
964 spoken Swedish: An up-to-date dialectography]. In *Diabas: Vol. 4. Kulturgränser - myt*  
965 *eller verklighet?* (pp. 215–228). Institutionen för nordiska språk vid Umeå Universitet.
- 966 Elert, Claes-Christian. (2000). *Allmän och svensk fonetik [General and Swedish phonetics]*  
967 (8., omarb. uppl). Stockholm: Norstedt.
- 968 Eliasson, Stig. (2022). The phonological status of Swedish au and eu: Proposals, evidence,  
969 evaluation. *Nordic Journal of Linguistics*, 1–42.  
970 <https://doi.org/10.1017/S0332586522000233>
- 971 Engstrand, Olle. (1999). Swedish. In *Handbook of the International Phonetic Association: A guide to the usage of the International Phonetic Alphabet*. Cambridge: Cambridge  
972 University Press.
- 973 Engstrand, Olle. (2004). *Fonetikens grunder [The basics of phonetics]*. Lund:  
975 Studentlitteratur.
- 976 Engstrand, Olle, Sven Björsten, Björn Lindblom, Gösta Bruce, & Anders Eriksson. (2000).  
977 Hur udda är Viby-i? Experimentella och typologiska observationer [How peculiar is  
978 Viby-i? Experimental and typological observations]. *Folkmålsstudier*, 39, 83–95.
- 979 Fabricius, Anne, Dominic Watt, & Daniel Ezra Johnson. (2009). A comparison of three  
980 speaker-intrinsic vowel formant frequency normalization algorithms for sociophonetics.

- 981        *Language Variation and Change*, 21(3), 413–435.
- 982        <https://doi.org/10.1017/S0954394509990160>
- 983        Fant, Gunnar. (1959). Acoustic analysis and synthesis of speech with applications to  
984        Swedish. *Ericsson Technics*, 15, 3–108.
- 985        Fant, Gunnar. (1971). Notes on the Swedish Vowel System. In L. Hammerich, R.  
986        Jakobson, E. Zwirner, & E. Fischer-Jørgensen (Eds.), *Form and substance: Phonetic*  
987        and linguistic papers. Odense: Andelsbogtrykkeriet.
- 988        Fant, Gunnar. (1983). Feature analysis of Swedish vowels - a revisit. *STL-QPSR*, 24(2-3),  
989        001–019.
- 990        Fant, Gunnar, G. Henningsson, & U. Stålhammar. (1969). Formant frequencies of Swedish  
991        vowels. *STL-QPSR*, 10(4), 026–031.
- 992        Flynn, Nicholas, & Paul Foulkes. (2011). Comparing vowel formant normalization  
993        methods. *Proceedings of ICPHS XVII, Hong Kong*, (August), 683–686.
- 994        Fujimura, Osamu. (1967). On the Second Spectral Peak of Front Vowels: A Perceptual  
995        Study of the Role of the Second and Third Formants. *Language and Speech*, 10(3),  
996        181–193. <https://doi.org/10.1177/002383096701000304>
- 997        Gross, Johan. (2018). *Mapping vowels: Variation and change in the speech of Gothenburg*  
998        adolescents. Gothenburg: University of Gothenburg.
- 999        Gross, Johan, Sally Boyd, Therese Leinonen, & James A. Walker. (2016). A tale of two  
1000        cities (and one vowel): Sociolinguistic variation in swedish. *Language Variation and*  
1001        *Change*, 28(2), 225–247.
- 1002        Gross, Johan, & Julia Forsberg. (2020). Weak Lips? A Possible Merger of /i:/ and /y:/ in  
1003        Gothenburg. *Phonetica*, 77(4), 268–288. <https://doi.org/10.1159/000499107>
- 1004        Gussenhoven, Carlos, & Wencui Zhou. (2013). *Revisiting pitch slope and height effects on*  
1005        *perceived duration*. 1365–1369. <https://doi.org/10.21437/Interspeech.2013-360>
- 1006        Hadding, Kerstin, Hajime Hirose, & Katherine S. Harris. (1976). Facial muscle activity in  
1007        the production of Swedish vowels: An electromyographic study. *Journal of Phonetics*,

- 1008 4(3), 233–245. [https://doi.org/https://doi.org/10.1016/S0095-4470\(19\)31246-X](https://doi.org/https://doi.org/10.1016/S0095-4470(19)31246-X)
- 1009 1 Hadding-Koch, Kerstin, & Arthur S. Abramson. (1964). Duration Versus Spectrum in  
1010 Swedish Vowels: Some Perceptual Experiments2. *Studia Linguistica*, 18(2), 94–107.  
1011 <https://doi.org/10.1111/j.1467-9582.1964.tb00451.x>
- 1012 1 Hammarström, G., & L. Norman. (1957). Om den friktiva slutfasen vid de svenska långa  
1013 vokalerna i, y, u, w.[On the final fricative phase in the Swedish long vowels i, y, u, w].  
1014 *Nordisk Tidsskrift for Tale Og Stemme*, 17(3).
- 1015 1 Henton, Caroline G. (2005). Creak as a sociophonetic marker. *The Journal of the  
1016 Acoustical Society of America*, 80(S1), S50–S50. <https://doi.org/10.1121/1.2023837>
- 1017 1 Heselwood, Barry, & Sara Howard. (2008). Clinical phonetic transcription. *The Handbook  
1018 of Clinical Linguistics*, 381–399.
- 1019 1 Hillenbrand, James M., Michael J. Clark, & Terrance M. Nearey. (2001). Effects of  
1020 consonant environment on vowel formant patterns. *The Journal of the Acoustical  
1021 Society of America*, 109(2), 748–763. <https://doi.org/10.1121/1.1337959>
- 1022 1 Hillenbrand, James M., Laura A. Getty, Michael J. Clark, & Kimberlee Wheeler. (1995).  
1023 Acoustic characteristics of American English vowels. *The Journal of the Acoustical  
1024 Society of America*, 97(5), 3099–3111.
- 1025 1 Hillenbrand, James M., & Terrance M. Nearey. (1999). Identification of resynthesized  
1026 /hVd/ utterances: Effects of formant contour. *The Journal of the Acoustical Society of  
1027 America*, 105(6), 3509–3523. <https://doi.org/10.1121/1.424676>
- 1028 1 Holbrook, Anthony, & Grant Fairbanks. (1962). Diphthong formants and their movements.  
1029 *Journal of Speech and Hearing Research*, 5(1), 38–58.
- 1030 1 Jacewicz, Ewa, & Robert Allen Fox. (2018). Regional variation in fundamental frequency  
1031 of American English vowels. *Phonetica*, 75(4), 273–309.  
1032 <https://doi.org/10.1159/000484610>
- 1033 1 Jacewicz, Ewa, Robert Allen Fox, & Joseph Salmons. (2011). Regional dialect variation in  
1034 the vowel systems of typically developing children. *Journal of Speech, Language, and*

- 1035        *Hearing Research*, 54(2), 448–470. [https://doi.org/10.1044/1092-4388\(2010/10-0161\)](https://doi.org/10.1044/1092-4388(2010/10-0161))
- 1036        Jenkins, James J., Winifred Strange, & Salvatore Miranda. (1994). Vowel identification in  
1037        mixed-speaker silent-center syllables. *The Journal of the Acoustical Society of America*,  
1038        95(2), 1030–1043. <https://doi.org/10.1121/1.410014>
- 1039        Johnson, Keith. (2005). Speaker normalization in speech perception. In D. B. Pisoni & R.  
1040        E. Remez (Eds.), *The Handbook of Speech Perception* (pp. 363–389). John Wiley &  
1041        Sons, Inc.
- 1042        Johnson, Keith, & Matthias J. Sjerps. (2021). Speaker normalization in speech perception.  
1043        In J. S. Pardo, L. C. Nygaard, R. E. Remez, & D. B. Pisoni (Eds.), *The handbook of  
1044        speech perception* (pp. 145–176). John Wiley & Sons, Inc.  
1045        <https://doi.org/10.1002/9781119184096.ch6>
- 1046        Joos, Martin. (1948). Acoustic Phonetics. *Language*, 24(2), 5–136.  
1047        <https://doi.org/10.2307/522229>
- 1048        Kent, Raymond D., & Houri K. Vorperian. (2018). Static measurements of vowel formant  
1049        frequencies and bandwidths: A review. *Journal of Communication Disorders*, 74, 74–97.  
1050        <https://doi.org/https://doi.org/10.1016/j.jcomdis.2018.05.004>
- 1051        Kerswill, Paul, & Susan Wright. (1990). The validity of phonetic transcription:  
1052        Limitations of a sociolinguistic research tool. *Language Variation and Change*, 2(3),  
1053        255–275. <https://doi.org/10.1017/S0954394500000363>
- 1054        Kotsinas, Ulla-Britt. (1994). *Ungdomsspråk [Youth language]*. Uppsala: Hallgren &  
1055        Fallgren.
- 1056        Kuronen, Mikko. (2000). *Vokaluttalets akustik i sverigesvenska, finlandssvenska och finska*  
1057        [*The acoustics of vowel pronunciation in Sweden Swedish, Finland Swedish and  
1058        Finnish*]. Jyväskylä: University of Jyväskylä.
- 1059        Labov, William. (1990). The intersection of sex and social class in the course of linguistic  
1060        change. *Language Variation and Change*, 2(2), 205–254.  
1061        <https://doi.org/10.1017/S0954394500000338>

- 1062 Labov, William. (2001). *Principles of linguistic change. 2: Social factors*. Oxford:  
1063 Wiley-Blackwell.
- 1064 Labov, William. (2010). *Principles of linguistic change. 2: Social factors* (repr).  
1065 Chichester: Wiley-Blackwell.
- 1066 Labov, William, Sharon Ash, & Charles Boberg. (2005). *The atlas of North American  
1067 English: Phonetics, phonology, and sound change*. Berlin • New York: De Gruyter  
1068 Mouton. <https://doi.org/doi:10.1515/9783110167467>
- 1069 Ladefoged, Peter, & Donald E. Broadbent. (1957). Information conveyed by vowels. *The  
1070 Journal of the Acoustical Society of America*, 29, 98–104.
- 1071 Lehiste, Ilse, & Gordon E. Peterson. (1961). Transitions, glides, and diphthongs. *The  
1072 Journal of The Acoustical Society of America*, 33(3), 268–277.
- 1073 Leinonen, Kari, Pitkänen Antti J., & Veijo V. Vihanta. (1981). Rikssvenskt och  
1074 finlandssvenskt ljudsystem ur perceptionssynpunkt [Perceptual perspectives on the  
1075 sound systems of Standard Swedish and Finland Swedish]. *X Fonetikan päivät,*  
1076 *TaYSYLJ* 7, 163–218. Tampere, Finland: Department of Finnish language and general  
1077 linguistics University of Tampere.
- 1078 Leinonen, Therese. (2010). *An Acoustic Analysis of Vowel Pronunciation in Swedish  
1079 Dialects*. Groningen: University of Groningen.
- 1080 Leung, Keith KW, Allard Jongman, Yue Wang, & Joan A. Sereno. (2016). Acoustic  
1081 characteristics of clearly spoken English tense and lax vowels. *The Journal of the  
1082 Acoustical Society of America*, 140(1), 45–58.
- 1083 Liljencrants, Johan, & Björn Lindblom. (1972). Numerical simulation of vowel quality  
1084 systems: The role of perceptual contrast. *Language*, 839–862.
- 1085 Lindblom, Björn. (1963). *On vowel reduction*. Uppsala: Uppsala University.
- 1086 Lindblom, Björn. (1998). Systemic constraints and adaptive change in the formation of  
1087 sound structure. In H. J. R., S.-K. M., & K. C. (Eds.), *Approaches to the evolution of  
1088 language: Social and cognitive bases* (pp. 242–264). Cambridge University Press

- 1089 Cambridge, United Kingdom.
- 1090 McAllister, Robert, James Lubker, & Johann Carlson. (1974). An EMG study of some  
1091 characteristics of the Swedish rounded vowels. *Journal of Phonetics*, 2(4), 267–278.  
1092 [https://doi.org/10.1016/S0095-4470\(19\)31297-5](https://doi.org/10.1016/S0095-4470(19)31297-5)
- 1093 Mennen, Ineke, Felix Schaeffler, & Gerard Docherty. (2012). Cross-language differences in  
1094 fundamental frequency range: A comparison of English and German. *The Journal of the  
1095 Acoustical Society of America*, 131(3), 2249–2260. <https://doi.org/10.1121/1.3681950>
- 1096 Munson, Benjamin, & Nancy Pearl Solomon. (2004). The effect of phonological  
1097 neighborhood density on vowel articulation. *Journal of Speech, Language, and Hearing  
1098 Research*, 47(5), 1048–1058. [https://doi.org/10.1044/1092-4388\(2004/078\)](https://doi.org/10.1044/1092-4388(2004/078))
- 1099 Nearey, Terrance M. (1978). *Phonetic Feature Systems for Vowels*. Indiana.
- 1100 Nearey, Terrance M., & Peter F. Assmann. (1986). Modeling the role of inherent spectral  
1101 change in vowel identification. *The Journal of the Acoustical Society of America*, 80(5),  
1102 1297–1308. <https://doi.org/10.1121/1.394433>
- 1103 Nilsson, Jenny, Lena Wenner, Therese Leinonen, & Eva Thorselius. (2021). New and old  
1104 social meanings in urban and rural Sweden. *Urban Matters*, 179.
- 1105 Nycz, Jennifer, & Lauren Hall-Lew. (2013). Best practices in measuring vowel merger.  
1106 *Proceedings of Meetings on Acoustics*, 20. AIP Publishing.
- 1107 Pape, Daniel, & Christine Mooshammer. (2006). Intrinsic F0 differences for German tense  
1108 and lax vowels. *Proceedings of the 7th International Seminar on Speech Production in  
1109 Ubatuba, Brazil, December 13th to 15th*, 271–278.
- 1110 Pelzer, Joppe Anna, & Paul Boersma. (2019). Diphthongization in three regional varieties  
1111 of Swedish. *Proceedings of the 19th International Congress of Phonetic Sciences*,  
1112 1144–1148. Canberra, Australia: Australian Speech Science and Technology  
1113 Association.
- 1114 Persson, Anna, Santiago Barreda, & T. Florian Jaeger. (2024). *Comparing accounts of  
1115 formant normalization against US English listeners' vowel perception*. Manuscript;

- 1116 Stockholm University.
- 1117 Persson, Anna, & T. Florian Jaeger. (2023). Evaluating normalization accounts against  
1118 the dense vowel space of Central Swedish. *Frontiers in Psychology*, 14, 01–21.  
1119 <https://doi.org/10.3389/fpsyg.2023.1165742>
- 1120 Persson, Anna, & T. Florian Jaeger. (2024). Measuring the informativity of F3 for  
1121 rounded and unrounded high-front vowels in Central Swedish. *Proceedings from*  
1122 *FONETIK 2024, Stockholm, June 3–5, 2024*, 13–18. Stockholm University.  
1123 <https://doi.org/10.5281/zenodo.11396050>
- 1124 Peterson, Gordon E. (1961). Parameters of Vowel Quality. *Journal of Speech and Hearing*  
1125 *Research*, 4(1), 10–29. <https://doi.org/10.1044/jshr.0401.10>
- 1126 Peterson, Gordon E., & Harold L. Barney. (1952). Control methods used in a study of the  
1127 vowels. *The Journal of the Acoustical Society of America*, 24(2), 175–184.
- 1128 R Core Team. (2023). *R: A language and environment for statistical computing*. Vienna,  
1129 Austria: R Foundation for Statistical Computing. Retrieved from  
1130 <https://www.R-project.org/>
- 1131 Renwick, Margaret E. L., & Joseph A. Stanley. (2020). Modeling dynamic trajectories of  
1132 front vowels in the American South. *The Journal of the Acoustical Society of America*,  
1133 147(1), 579–595. <https://doi.org/10.1121/10.0000549>
- 1134 Riad, Tomas. (2014). *The phonology of Swedish*. Oxford: Oxford University Press.
- 1135 Robb, Michael P., & Yang Chen. (2009). Is /h/ phonetically neutral? *Clinical Linguistics*  
1136 & *Phonetics*, 23(11), 842–855. <https://doi.org/10.3109/02699200903247896>
- 1137 RStudio Team. (2020). *RStudio: Integrated development environment for r*. Boston, MA:  
1138 RStudio, PBC. Retrieved from <http://www.rstudio.com/>
- 1139 Scarborough, Rebecca. (2012). Lexical similarity and speech production: Neighborhoods  
1140 for nonwords. *Lingua*, 122(2), 164–176.  
1141 <https://doi.org/https://doi.org/10.1016/j.lingua.2011.06.006>
- 1142 Schaeffler, Felix. (2005). *Phonological quantity in Swedish dialects: Typological aspects*,

- 1143        *phonetic variation and diachronic change*. Umeå: Umeå University, Dep. of philosophy  
1144        and linguistics.
- 1145        Schertz, Jessamyn. (2013). Exaggeration of featural contrasts in clarifications of misheard  
1146        speech in English. *Journal of Phonetics*, 41(3-4), 249–263.
- 1147        Schertz, Jessamyn, & Emily J. Clare. (2020). Phonetic cue weighting in perception and  
1148        production. *WIREs Cognitive Science*, 11(2), e1521.  
1149        <https://doi.org/https://doi.org/10.1002/wcs.1521>
- 1150        Schötz, Susanne, Johan Frid, & Anders Löfqvist. (2011). Exotic vowels in Swedish – an  
1151        articulographic and acoustic pilot study of /i:/. *Proceedings of the 17th International*  
1152        *Congress of Phonetic Sciences*. Hong Kong.
- 1153        Seyfarth, Scott, Esteban Buz, & T. Florian Jaeger. (2016). Dynamic hyperarticulation of  
1154        coda voicing contrasts. *The Journal of the Acoustical Society of America*, 139(2),  
1155        EL31–EL37.
- 1156        Sóskuthy, Márton. (2021). Evaluating generalised additive mixed modelling strategies for  
1157        dynamic speech analysis. *Journal of Phonetics*, 84.  
1158        <https://doi.org/10.1016/j.wocn.2020.101017>
- 1159        Sóskuthy, Márton, Paul Foulkes, Vincent Hughes, & Bill Haddican. (2018). Changing  
1160        words and sounds: The roles of different cognitive units in sound change. *Topics in*  
1161        *Cognitive Science*, 10(4), 787–802. <https://doi.org/https://doi.org/10.1111/tops.12346>
- 1162        Ståhle, Carl Ivar. (1965). 'Mötet uppnas på sundag'[The meeting uppnas on sundag].  
1163        *Språkvård*, 3(3-8), 1–15.
- 1164        Stålhammar, U., I. Karlsson, & Gunnar Fant. (1973). Contextual effects on vowel nuclei.  
1165        *STL-QPSR*, 14(4), 001–018.
- 1166        Stemberger, Joseph Paul, & Barbara May Bernhardt. (2020). Phonetic transcription for  
1167        speech-language pathology in the 21st century. *Folia Phoniatrica Et Logopaedica*, 72(2),  
1168        75–83.
- 1169        Stevens, Kenneth N., & Arthur S. House. (1963). Perturbation of vowel articulations by

- 1170 consonantal context: An acoustical study. *Journal of Speech and Hearing Research*,  
1171 6(2), 111–128. <https://doi.org/10.1044/jshr.0602.111>
- 1172 Stilp, Christian. (2020). Acoustic context effects in speech perception. *WIREs Cognitive  
1173 Science*, 11(1), 1–18. <https://doi.org/10.1002/wcs.1517>
- 1174 Strange, Winifred. (1989). Evolving theories of vowel perception. *The Journal of the  
1175 Acoustical Society of America*, 85(5), 2081–2087. <https://doi.org/10.1121/1.397860>
- 1176 Strangert, Eva. (2001). Quantity in ten Swedish dialects in Northern Sweden and  
1177 Österbotten in Finland. *Working Papers/Lund University, Department of Linguistics  
1178 and Phonetics*, 49, 144–147.
- 1179 Syrdal, Ann K. (1985). Aspects of a model of the auditory representation of American  
1180 English vowels. *Speech Communication*, 4(1-3), 121–135.  
1181 [https://doi.org/10.1016/0167-6393\(85\)90040-8](https://doi.org/10.1016/0167-6393(85)90040-8)
- 1182 Watson, Catherine I., & Jonathan Harrington. (1999). Acoustic evidence for dynamic  
1183 formant trajectories in Australian English vowels. *The Journal of the Acoustical Society  
1184 of America*, 106(1), 458–468. <https://doi.org/10.1121/1.427069>
- 1185 Wedel, Andrew, Noah Nelson, & Rebecca Sharp. (2018). The phonetic specificity of  
1186 contrastive hyperarticulation in natural speech. *Journal of Memory and Language*, 100,  
1187 61–88.
- 1188 Weirich, Melanie, Adrian P. Simpson, Jasmine Öjbro, & Christine Ericsdotter Nordgren.  
1189 (2019). The phonetics of gender in Swedish and German. *Fonetik 2019, Stockholm,  
1190 Sweden, 10-12 June, 2019*, 49–53.
- 1191 Wenner, Lena. (2010). *När lögnare blir lugnare. En sociofonetisk studie av sammanfallet  
1192 mellan kort ö och kort u i uppländskan [When lögnare become lugnare. A sociophonetic  
1193 study of the merger between short ö and short u in Uppland Swedish]*. Uppsala: Uppsala  
1194 Universitet.
- 1195 Whalen, D. H., & Andrea G. Levitt. (1995). The universality of intrinsic F0 of vowels.  
1196 *Journal of Phonetics*, 23(3), 349–366.

- 1197 https://doi.org/https://doi.org/10.1016/S0095-4470(95)80165-0
- 1198 Wieling, Martijn. (2018). Analyzing dynamic phonetic data using generalized additive  
1199 mixed modeling: A tutorial focusing on articulatory differences between L1 and L2  
1200 speakers of english. *Journal of Phonetics*, 70, 86–116.
- 1201 Woods, Nicola J. (1997). The formation and development of New Zealand English:  
1202 Interaction of gender-related variation and linguistic change. *Journal of Sociolinguistics*,  
1203 1(1), 95–125. https://doi.org/https://doi.org/10.1111/1467-9481.00005
- 1204 Wright, Richard, John Local, Richard Ogden, & Rosalind Temple. (2004). Factors of  
1205 lexical competition in vowel articulation. *Papers in Laboratory Phonology VI*, 75–87.
- 1206 Xie, Xin, & T. Florian Jaeger. (2020). Comparing non-native and native speech: Are L2  
1207 productions more variable? *The Journal of the Acoustical Society of America*, 147(5),  
1208 3322–3347. https://doi.org/10.1121/10.0001141
- 1209 Yang, Byunggon. (2019). A comparison of normalized formant trajectories of English  
1210 vowels produced by American men and women. *Phonetics and Speech Sciences*, 11(1),  
1211 1–8.
- 1212 Young, Nathan J., & Michael McGarrah. (2021). Forced alignment for Nordic languages:  
1213 Rapidly constructing a high-quality prototype. *Nordic Journal of Linguistics*, 1–27.  
1214 https://doi.org/10.1017/S033258652100024X