

1 **Vowel recognition at fundamental frequencies up to 1 kHz**
2 **reveals point vowels as acoustic landmarks**

3 Daniel Friedrichs^{a)}

4 Department of Speech, Hearing and Phonetic Sciences, UCL
5 2 Wakefield Street, London WC1N 1PF, United Kingdom

6 Dieter Maurer

7 Institute of the Performing Arts and Film, Zurich University of the Arts (ZHdK)
8 Toni-Areal, Pfingstweidstrasse 96, CH-8031 Zurich, Switzerland

9 Stuart Rosen

10 Department of Speech, Hearing and Phonetic Sciences, UCL
11 2 Wakefield Street, London WC1N 1PF, United Kingdom

12 Volker Dellwo

13 Phonetics Group, Department of Computational Linguistics, University of Zurich
14 Andreasstrasse 15, CH-8050 Zurich, Switzerland

15 **4 Figures**
16 **11 Tables (Appendix)**

^{a)} Author to whom correspondence should be addressed: daniel.friedrichs@ucl.ac.uk

Abstract

18 The phonological function of vowels can be maintained at fundamental frequencies (f_o) up
19 to 880 Hz [Friedrichs et al. (2015). J. Acoust. Soc. Am. **138**, EL36–EL42]. Here, we
20 test the influence of talker variability and multiple response options on vowel recognition at
21 high f_o s. The stimuli (n=264) consisted of eight isolated vowels (/i y e ø ε a o u/) produced
22 by three female native German talkers at eleven f_o s within a range of 220–1046 Hz. In a
23 closed-set identification task, 21 listeners were presented excised 700-ms vowel nuclei with
24 quasi-flat f_o contours and resonance trajectories. The results show that listeners can identify
25 the point vowels /i a u/ at f_o s up to almost 1 kHz, with a significant decrease for the vowels
26 /y ε/ and a drop to chance level for the vowels /e ø o/ towards the upper f_o s. Auditory
27 excitation patterns reveal highly differentiable representations for /i a u/ that can be used
28 as landmarks for vowel category perception at high f_o s. These results suggest that theories
29 of vowel perception based on overall spectral shape will provide a fuller account of vowel
30 perception than those based solely on formant frequency patterns.

31

32 © The Journal of the Acoustical Society of America

33 **PACS number:** 43.71.-k

34 **I. INTRODUCTION**

35 Patterns of formant frequencies are commonly assumed to be the most salient cues to
36 vowel perception. The assumption that the vowel identification process is mainly driven by
37 such an underlying acoustic representation contributes largely to the pervasive idea that
38 listeners' ability to recognize vowels has to be poor at very high fundamental frequencies
39 (f_o) due to a sparse sampling of the vocal tract transfer function. This holds true, in
40 particular, when the normal range of the first formant frequency (F_1) is exceeded by f_o ,
41 and the higher formants are poorly specified due to a wide spacing of the harmonics.

42 Support for this view is mainly provided by studies on Western operatic singing.
43 Howie and Delattre (1962), for example, found in a study on the perception of high-pitched
44 vowels (f_o range 132–1056 Hz) sung by a baritone and a soprano that vowels lose their
45 identity increasingly with increasing f_o . This degradation starts with the categories usually
46 characterized by a low F_1 (i.e., high vowels such as /i/ and /u/) and leaving only those
47 with the highest F_1 (i.e., low vowels such as /a/ and /ɑ/) identifiable at very high f_o s. Ever
48 since, numerous studies have reported that only /a/-like vowels can remain identifiable at
49 the highest musical notes near 1 kHz (see Sundberg, 2013, p. 87, for an overview). It seems
50 plausible, however, that this loss of vowel contrast is primarily due to articulatory changes
51 applied by Western operatic singers when they perform at higher pitches. In experimental
52 studies such as Joliveau et al. (2004) it has been shown, for example, that sopranos shift

53 the first resonant frequency (f_{R1}) of their vocal tract – and thus F_1 – to the vicinity of f_o
54 as soon as f_o drastically exceeds the normal range of f_{R1} of an intended vowel. This tuning
55 of f_{R1} is achieved by increasing the jaw opening and reducing the maximum constriction of
56 the vocal tract (Sundberg, 1975; Sundberg, 2013). As f_o gains considerable amplitude
57 when being closer to a resonant frequency, these maneuvers may help a singer to maintain
58 vocal power and timbral homogeneity (Smith and Wolfe, 2009). However, the acoustic
59 modifications associated with shifting a resonant frequency may lead to ambiguous formant
60 frequency patterns and consequently to a confusion of vowel categories.

61 Given this situation, it is surprising that few studies have investigated vowel
62 recognition outside Western operatic singing at very high f_o s as there is evidence that even
63 a sparsely sampled vocal tract transfer function still carries information, which can be used
64 by listeners to recognize different vowels, despite a likely absence of the supposed F_1 and
65 an undersampling of the higher formants. Smith and Scott (1980), for example, reported
66 listeners' identification performance significantly above chance level (mean of 70% correct)
67 for the four front vowels /i ɪ ε æ/, which were produced by a soprano in isolation at an f_o
68 of about 880 Hz (i.e., the musical note A5) with a raised larynx (i.e., a shortened vocal
69 tract), and thus not in an articulation mode typical for Western operatic singers. When
70 asked to produce the same vowels in her operatic singing style, identification dropped to a
71 mean of 4% correct at the same f_o . Maurer and Landis (1996) showed that infant and

72 adult talkers can produce identifiable versions of the vowels /i a o u/ but not of /e/ at an
73 f_o between about 500–870 Hz that was individually chosen by the talker. In a more recent
74 study, Maurer et al. (2014) investigated the high-pitched vowels /i y œ a ɔ u/ produced by
75 a female Cantonese opera singer in isolation and monosyllabic consonant-vowel utterances
76 and found that /i a ɔ u/ could be identified by more than 80% of the listeners within an f_o
77 range of 820–860 Hz. In a study using a two-alternative forced choice task, Friedrichs et al.
78 (2015a) provided evidence that the phonological function of the eight vowels /i y e ø ε a o
79 u/ (i.e., the function they fulfil in linguistic contrastive position to help listeners
80 distinguish between words) can be maintained at f_o s up to at least 880 Hz when they were
81 produced in minimal pairs. These judgments were made on excised steady-state vowel
82 nuclei (250 ms) excluding consonantal context phenomena such as co-articulation and
83 formant transitions. This is particularly surprising for vowels that typically have a low F_1
84 that were tested in combination with adjacent vowels with similar F_2 (e.g., /i/ vs. /e/ and
85 /u/ vs. /o/), because an absent F_1 has been argued to make vowels with a similar F_2
86 indistinguishable (Smith and Wolfe, 2009, p. E196; see Ito et al., 2001, for contradictory
87 results). In a follow-up study (Friedrichs et al., 2015b), a female talker produced the same
88 vowels except /u/ in the German word context /l-V-gən/ (/u/ was excluded as it would
89 have resulted in a meaningless utterance), and a multiple-choice identification task was
90 used. It was found that the words including /i y a o/ remained identifiable – and thus the

91 vowels' phonological function could be maintained – throughout the investigated f_o range
92 from 220 to 880 Hz. For the vowels /e ø ε/, however, a significant decrease was observed in
93 listeners' identification performance within this range (for /ø/ from about 587 Hz and for
94 /e ε/ from about 784 Hz). At the highest f_o used (880 Hz), listeners could recognize the
95 vowel /ε/ again.

96 The acoustic features and perceptual mechanisms underlying accurate vowel category
97 perception at such high f_o s remain unclear. As some of these studies found high
98 identification rates even when excluding cues that play an important secondary role in
99 vowel perception (e.g., vowel duration and formant frequency movement, see Lehiste and
100 Peterson, 1961), it seems possible that spectral information apart from formant frequencies
101 allowed listeners to identify vowels at very high f_o s. Besides vowel identification models
102 that are based on formant frequency distribution, speech scientists (in particular, from the
103 automatic speech recognition community) have long recognized that overall spectral shape
104 as reflected by, for example, Mel Frequency Cepstral Coefficients (MFCCs) (Davis and
105 Mermelstein, 1980), are a more robust feature set than formants. Pols et al. (1969) and
106 Klein et al. (1970) showed that a simple filter bank analysis (essentially an auditory
107 excitation pattern approach which encodes the overall shape of the spectrum) matched
108 perceptual vowel spaces well. Zahorian and Jagharghi (1993) found in an automatic vowel
109 classification experiment that spectral-shape features (the discrete cosine transform

110 coefficients of a bark frequency scaled spectrum) are superior acoustic cues for vowel
111 identity classification compared to formants. Ito et al. (2001) showed that also the
112 amplitude ratio of high- to low-frequency components (i.e., the spectral tilt) affects the
113 perceived vowel category and is at least equally effective as F_2 as a cue for vowel
114 identification. Several overall-spectral-shape models have been advocated over the last
115 decades (see Kiefte et al., 2013, for a more comprehensive review of this approach). Most
116 of them do not pay special attention to the distribution of formants, but are based on the
117 assumption that the gross shape of a smoothed spectral envelope underlies the
118 identification process. As it is very unlikely to find common formant frequency patterns at
119 f_o s of about 880 Hz, it seems possible that the overall spectral shape – despite a severe
120 undersampling of the spectral envelope (see de Cheveigné and Kawahara, 1999, and
121 Hillenbrand and Houde, 2003, for more details on this problem) – might have conveyed the
122 information that allowed listeners to identify different vowel categories (but see Maurer,
123 2016, for an argument that perceived vowel categories are more a result of a complex
124 systematic interaction between spectral shapes and f_o than has generally been assumed in
125 phonetic theory).

126 However, it is also possible that the lack of between-talker acoustic vowel variation
127 facilitated identification of the vowels (excepting Maurer and Landis, 1996, who used
128 vowels of infant and adult talkers, all of the above-mentioned studies showing accurate

129 vowel category perception at high f_o s were single-talker studies). In that situation, listeners
130 may have adapted to the talker's individual articulatory behavior (i.e., the within-talker
131 acoustic vowel variation). Thus, it is not clear whether the results can be generalized to
132 other talkers and whether an experimental design including more than one talker would
133 lead to similar results. In addition, it seems likely that the number of response options
134 (i.e., binary and multiple-choice tasks were used) had an effect on the identification
135 performance as listeners perform better when fewer response options are provided.

136 The present study addresses these issues. Here, we asked three female talkers to
137 produce the eight vowels /i y e ø ε a o u/ in isolation (thus eliminating possible
138 confounding effects due to co-articulation with adjacent consonants) at eleven f_o s within a
139 range of 220–1046 Hz. In a multiple-choice task (mixed-talker condition) with all possible
140 vowels as response options, listeners had to identify single 700-ms nuclei with quasi
141 steady-state acoustic characteristics. These center portions of the vowels were used to
142 exclude possible secondary cues, in particular, sweeping harmonics in the on- and off-sets,
143 which might sample the vocal tract transfer function more continuously and thus provide
144 information about the position of the formants.

145 To investigate possible spectral properties underlying listeners' identification process
146 at high f_o s, we calculated simple versions of the excitation patterns that these vowels
147 would be expected to generate in the auditory periphery and discuss them with respect to

148 the results of the identification test.

149 **II. METHODS**

150 **A. Subjects**

151 21 native German listeners (10 female, 11 male; mean age = 23.2, s.d. = 2.25)
152 participated in a multiple-choice vowel identification task. All were students at the
153 University of Zurich and none of them reported any hearing impairments when asked
154 before the experiment.

155 **B. Stimuli and apparatus**

156 Three female native German talkers with professional voice training (one soprano,
157 age: 33; one Musical-Theatre singer, age: 34; one actress, age: 34) were recorded with a
158 cardioid condenser microphone (Sennheiser MKH 40 P48 with pop shield,
159 Wedemark-Wennebostel, Germany) on a PC via an audio interface (RME Fireface UCX,
160 RME, Halmhausen, Germany) in a noise-controlled room at Zurich University of the Arts
161 (ZHdK) (Switzerland). The sampling frequency of the recordings was 44.1 kHz. Subjects
162 were recorded keeping a constant distance of about 30 cm to the microphone when
163 standing on a drawn position reference on the floor. They were selected based on samples
164 from a corpus of recordings of 60 talkers because of their extended vocal range and
165 noticeable skill of maintaining vowel categories at high f_o s. As part of the standard

166 procedure as implemented in an associated project (see Maurer et al., 2016, for more
167 details), the latter was assessed in a listening test using a blocked-talker condition and a
168 multiple-choice identification task carried out by five phonetically trained listeners. The
169 other 57 talkers (both female and male) had more limited vocal ranges and were not
170 capable of producing vowels throughout the designated f_o range from 220 to 1046 Hz.

171 The three subjects were then asked to produce the eight long vowels /i y e ø ε a o u/
172 in isolation at eleven f_o s (220, 330, 440, 523, 587, 659, 698, 784, 880, 988, 1046 Hz) with a
173 monotone pitch contour resulting in 264 recordings (11 frequencies * 8 vowels * 3 talkers).
174 Piano notes were presented as reference sounds to the subjects via loudspeaker
175 immediately preceding the production. The talkers were asked to focus on producing
176 recognizable vowels and to ignore typical voice aesthetics that might be important in their
177 respective artistic style. The lowest f_o (220 Hz) corresponds to the female average f_o in
178 citation-form words (Hillenbrand et al., 1995). The highest f_o (1046 Hz) corresponds to the
179 high C (the musical note C6) in soprano singing and exceeds the normal range of F_1 of all
180 German vowels produced by female talkers (see Pätzold and Simpson, 1997). The average
181 f_o of each vowel was measured in Praat (Boersma and Weenink, 2016) using it's
182 autocorrelation method (Boersma, 1993) and later checked manually. All vowels used in
183 this study were recorded several times to ensure that at least one had an actual f_o close to
184 the target f_o and a minimum duration of 1 second. All vowels that met these criteria were

185 then evaluated again in the same listening test carried out by the five phonetically trained
186 listeners, and the vowels with the highest identification scores were selected as stimuli. The
187 mean duration of the final recordings was 1.49 s (range from on- to offset of voicing: 1.18 –
188 2.83 s).

189 Only vowel centers of 700 ms (\pm 350 ms from the vowel midpoint) with quasi-flat f_o
190 contours and steady-state spectral characteristics were used as stimuli. On- and offsets of
191 the excised sounds were faded over 5 ms by amplitude modulating the waveform with
192 raised cosines. All stimuli were normalized to an arbitrary intensity. The overall output
193 level was chosen by listeners individually to be comfortable.

194 C. Procedure

195 A mixed-talker listening test was carried out in a small and noise-controlled room at
196 the University of Zurich (Switzerland) using closed dynamic headphones (Beyerdynamic
197 DT 770 Pro, 250 Ω). The experiment consisted of a multiple-choice identification task with
198 all 8 vowels as response options. Listeners (n=21) were presented the excised 700-ms vowel
199 nuclei while they saw a screen that contained eight circularly arranged buttons, each button
200 labeled with one category (randomly arranged). Above the response buttons listeners could
201 read the question *Welchen Vokal hörst Du?* (*Which vowel do you hear?*). The listener's
202 task was to identify the vowel presented from the eight response options provided. After
203 listeners made their choice they heard the next stimulus automatically with a delay of one

204 second. Listeners could not repeat a stimulus. Each listener heard each token only once
205 which means that any particular vowel at each f_o was responded to 63 times.

206 D. Data analysis

207 We performed a set of statistical analyses on correct/incorrect responses using
208 mixed-effects logistic regression models in R (version 3.3.1; R Development Core Team,
209 2016, lmerTest package; Kuznetsova et al., 2014), in which listeners and items were entered
210 as random variables (Baayen et al., 2008). The predictors were vowel category, f_o , talker,
211 and all their interaction. The significance of the main effects and interactions was assessed
212 with likelihood ratio tests that compared the model with the main effect or interaction to a
213 model without it. For clarity's sake, the results and figures are presented in percentages,
214 although all statistical analyses were performed on raw data (correct/incorrect responses).
215 The estimates (β) that are reported in the results section are expressed in logit units and
216 were computed taking "incorrect response" as the reference level for the dependent variable.

217 To investigate possible shifts towards other than the intended vowel categories, 11
218 confusion matrices (one for each f_o , each based on a total of 504 samples, i.e., 8 vowels x 3
219 talkers x 21 listeners' responses) with the two dimensions *intended vowel* (actual class) and
220 *response vowel* (predicted class) were calculated.

221 E. Excitation patterns

222 Simple auditory excitation patterns were generated for each vowel using a 200-channel
 223 linear gammatone filter bank, whose bandwidths and centre frequencies were calculated
 224 according to the ERB formulae given by Glasberg and Moore (1990). The rms level of the
 225 output wave was calculated for each filter channel, and converted to dB. In addition, a
 226 frequency weighting was applied to account for the transmission properties of the middle
 227 ear, as based on measurements made by Puria et al. (1997).

228 III. RESULTS

229 Results obtained from the logistic regression revealed a highly significant effect of f_o
 230 ($\chi^2(10) = 30.8$, $p < .001$), a highly significant effect of vowel category ($\chi^2(7) = 28.21$, $p <$
 231 $.001$), no main effect of talker ($\chi^2(2) = 2.24$, $p = .33$), and a highly significant interaction
 232 between the three ($\chi^2(244) = 627.91$, $p < .001$). For the ease of interpretation, and as a
 233 complex three-way interaction makes it impossible to ignore any one of them in accounting
 234 for the effects of the other two, we decided to break down the data into three sets to test
 235 for a two-way interaction between vowel category and f_o for the individual talkers. The
 236 results of the three analyses showed consistently a highly significant interaction between
 237 vowel category and f_o (talker 1: $\chi^2(70) = 188.42$, $p < .001$; talker 2: $\chi^2(70) = 182.74$, $p <$
 238 $.001$; talker 3: $\chi^2(70) = 209.5$, $p < .001$). Significant effects of vowel category were found
 239 for all talkers (talker 1: $\chi^2(7) = 28.19$, $p < .001$; talker 2: $\chi^2(7) = 22.01$, $p < .01$; talker 3:
 240 $\chi^2(7) = 35.77$, $p < .001$), and f_o (talker 1: $\chi^2(10) = 30.79$, $p < .001$; talker 2: $\chi^2(10) =$

241 32.61, $p < .001$; talker 3: $\chi^2(10) = 30.2$, $p < .001$). Taken together, these effects suggest
 242 that listeners' identification performance showed high variability between vowel categories
 243 and across f_o s generally.

244 Figure 1 shows the distribution of the percentage of correct identification for each f_o
 245 and talker across vowels. Throughout the f_o range the overall performance declined more
 246 or less continuously for all talkers.

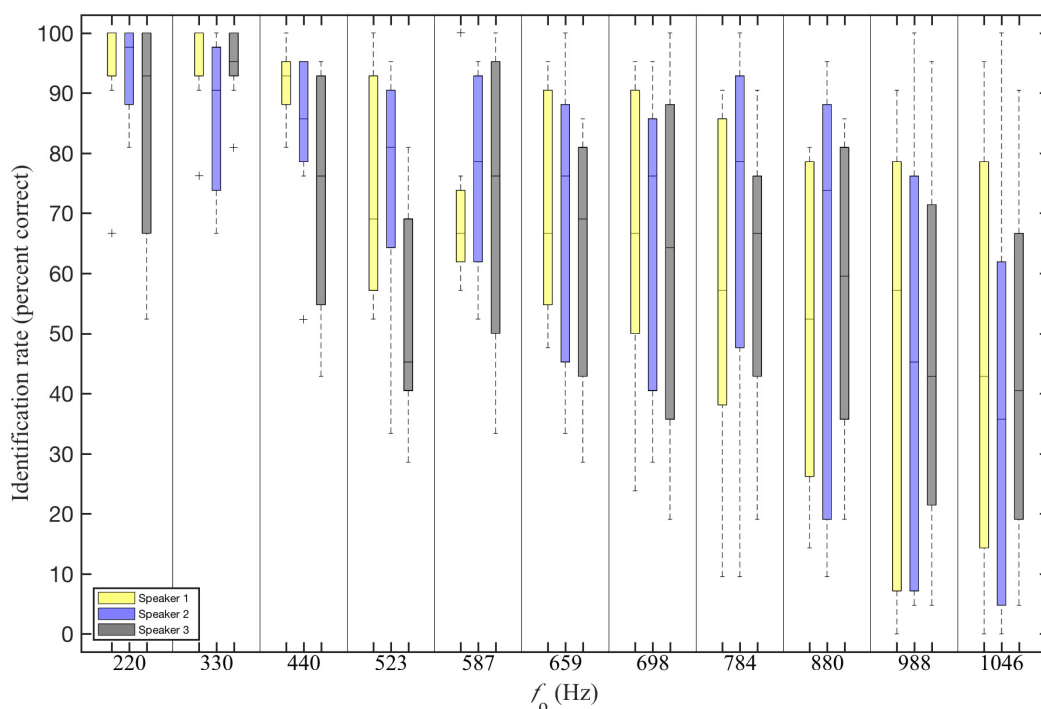


Figure 1: (Color online) Box plots showing the distribution of percent correct for the identification of all investigated vowels at the eleven f_o s for the individual talkers.

247 The increasing variability toward the higher f_o s can be explained by an increasing

248 inter-vowel variability, as the identification rate of individual vowel categories differed
 249 greatly between low and high f_o s. This can be seen in Figure 2 showing the mean percent
 250 correct scores for each individual vowel at the different f_o s. Listeners' identification
 251 performance for the vowels /i ε a u/ is surprisingly stable up to at least 880 Hz, and
 252 percent correct values can typically be found in the range above 70%. At the two highest
 253 f_o s (988 and 1046 Hz), the identification rate for /ε/ drops to intermediate ranges between
 254 40 and 50% correct. Only the point vowels /i a u/ remain in the upper third of the percent
 255 correct scale. On the contrary, for the vowels /e ø o/ an extensive decrease in listeners'
 256 identification performance can be found throughout the f_o s from 220 to 1046 Hz. While
 257 identification scores range between 90–100% at the two lowest f_o s (220 and 330 Hz), they
 258 drop fairly continuously toward chance level for these three vowels, which is reached at 988
 259 Hz. The identification rate of /y/ drops substantially at an f_o of 523 Hz (from about 85 to
 260 60% correct) and decreases despite some variability towards upper f_o s. From 988 Hz
 261 identification scores are similar to those of /ε/ (i.e., within the 35–50% correct range).

262 Confusion matrices (see Figure 3, for a graphical illustration; the raw data can be
 263 found in Appendix A) reveal dominant shifts toward the vowel categories /i a u/ in cases of
 264 false identifications at the highest f_o s. For /ε/, strong confusions at the highest two f_o s
 265 (988 and 1046 Hz) were found with /a/, which also showed the highest response
 266 proportions of all vowels at these f_o s (28% and 24.4%). The drop in identification

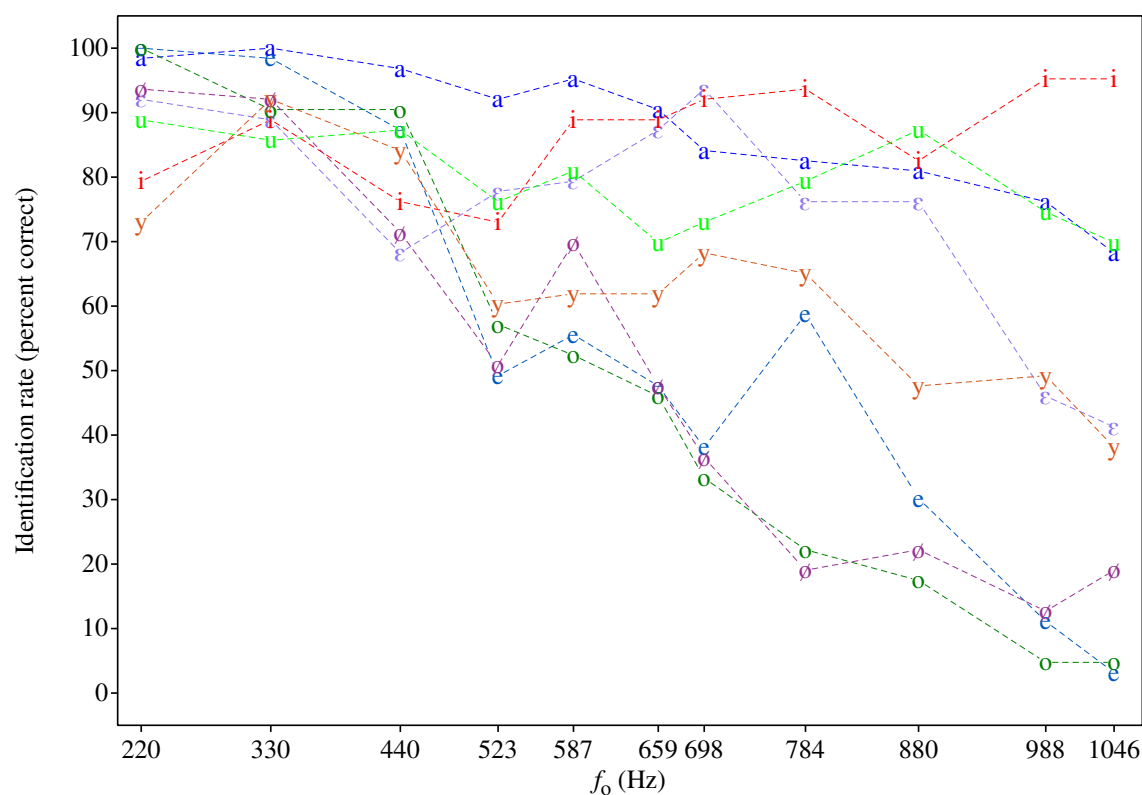


Figure 2: (Color online) Line graphs showing percent correct values, summed over all talkers, for the identification of each of the eight vowels over the investigated f_o range.

267 performance for the vowel /y/ in the range from 523 Hz on upwards is due to a confusion
 268 with other front vowels and from 784 Hz upwards mainly due to a confusion with /i/. A
 269 confusion between these two vowels also explains the relatively poor performance for /i/ at
 270 the lowest f_o 220 Hz (15.9% of the listeners responded /i/ when /y/ was presented to
 271 them). In case of /ø/, shifts in perception were generally found to be widely spread, that
 272 is, toward all the investigated vowel categories except /i/. The majority of false

273 identification of /o/ shifted from a perceived /a/ at 523 and 587 Hz to /u/ at all higher
274 f_o s. Within the range 523–784 Hz, the vowel /e/ was often confused with /i/. At higher
275 f_o s the perceived vowel category shifted toward / ϵ / and /a/.

276 Figure 4 shows the auditory excitation patterns for the eight vowels used in this study
277 produced at an f_o of about 988 Hz. Both the patterns calculated for individual talkers and
278 those averaged across talkers reveal that the point vowels /i a u/ show maximally distinct
279 spectral shapes, which can be easily distinguished by the overall excitation level in the
280 higher frequency region above about 1.5 kHz. The obtained confusions of the vowel
281 categories /y e \emptyset ϵ o/ at this f_o show a high degree of correspondence to the excitation
282 patterns of the respective point vowels they were confused with most often. For example,
283 the pattern calculated for /o/ shows high similarity with the pattern of the point vowel
284 /u/, that is, a relatively low excitation level in the high frequency region. The excitation
285 pattern of /y/ exhibits a relatively high excitation level in the high frequency region, which
286 is also the case for the point vowel /i/. The patterns of the vowels /e \emptyset ϵ / show
287 intermediate levels of excitation in the high frequency region, which is also the case for /a/,
288 the vowel which was most often responded by the listeners when these vowels were
289 presented to them at 988 Hz.

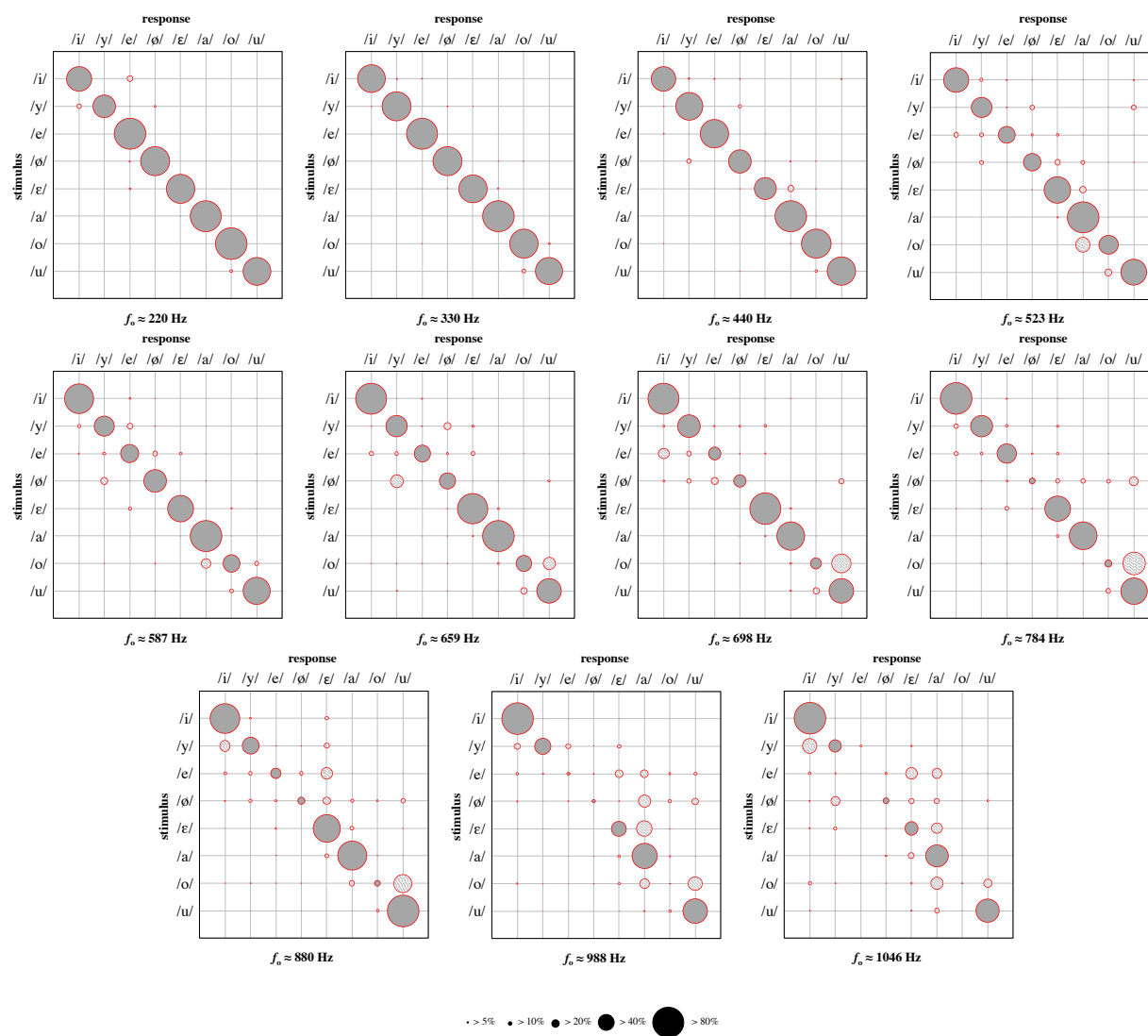


Figure 3: (Color online) Graphical confusion matrices showing the intended and response vowel categories for each f_o . The radius of each circle is proportional to the number of times that a particular stimulus (given by the row) was identified as the column response. Correct responses (down the diagonal) are solid gray, whereas identification errors (confusions) are indicated by diagonal lines through the circles.

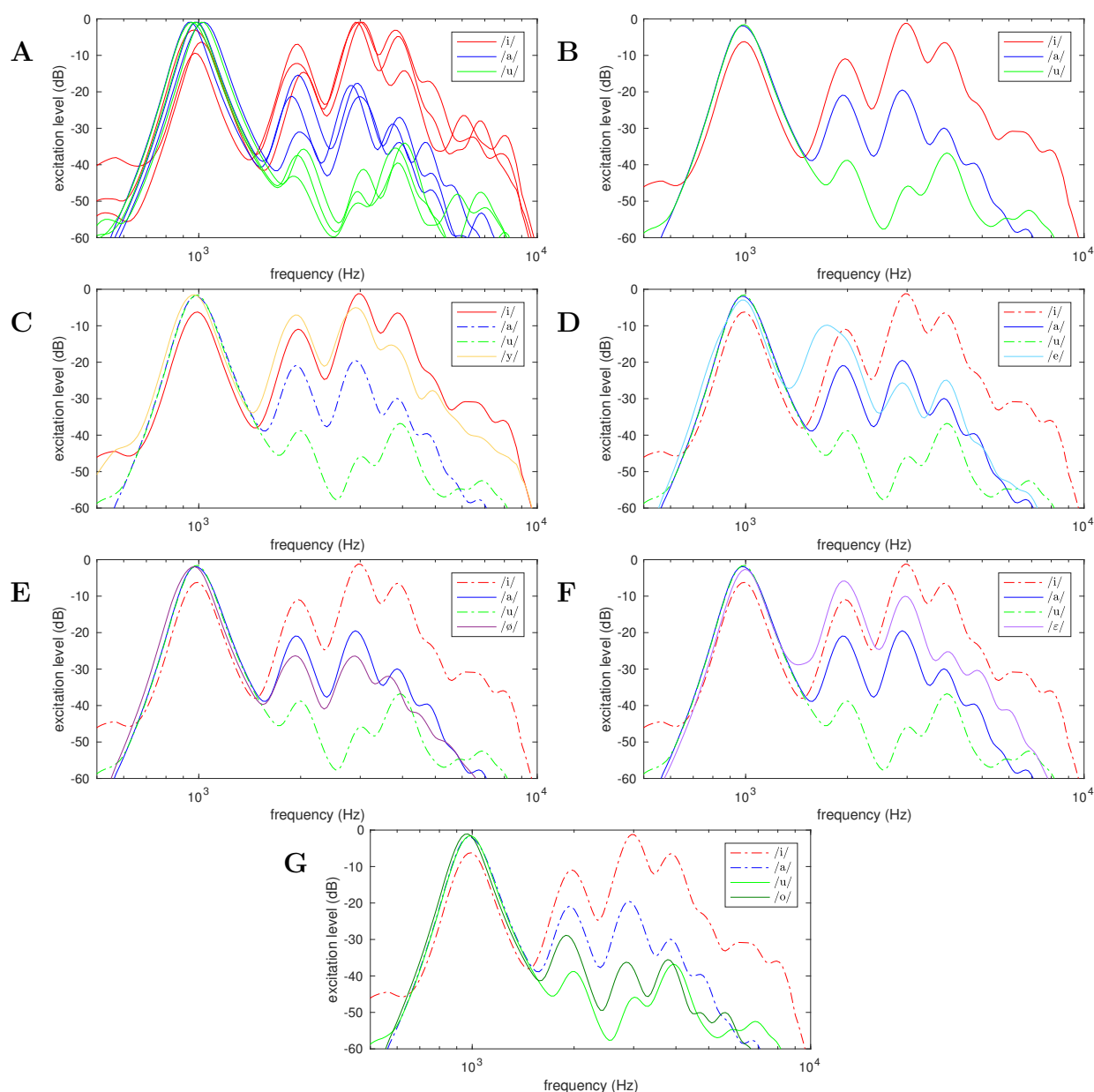


Figure 4: (Color online) Excitation patterns for the vowels used in this study that had an f_o of about 988 Hz. Part (A) shows the excitation patterns for the individual point vowels /i a u/ produced by all talkers. Part (B) shows the excitation patterns of the same vowels averaged across talkers. All other parts (C–G) show each of the other investigated vowels together with the point vowels. In these graphs, solid lines are used to indicate the strongest confusion of a respective vowel with one of the point vowels. (The information in this figure may not be properly conveyed in black and white.)

290 **IV. DISCUSSION**

291 The results have shown that listeners' abilities to recognize vowels within a
292 fundamental frequency range from 220 to 1046 Hz differ greatly across vowel categories and
293 the range of f_o s. Listeners could perform well even with a variety of talkers, which means
294 that good performance at high f_o s is not being done through some odd mechanism or
295 sensitivity which would be idiosyncratic for each talker. It is not surprising that all vowels
296 could be identified accurately at the lowest f_o s used here (220 and 330 Hz), but it is
297 striking that only the performance for the vowels /y e ø o/, but not for /i a ε u/ decreased
298 drastically within the f_o range from around 523 to 880 Hz. The results also revealed that
299 the point vowels /i a u/ remain identifiable at an f_o close to 1 kHz or even above (in the
300 case of /i/).

301 Thus, the results differ substantially from those provided by numerous studies on
302 vowel identification in Western classical singing, which have reported consistently that high
303 vowels such as /i/ and /u/ are the first vowels to lose their identity when f_o is
304 progressively increased. This means that findings from the field of operatic singing cannot
305 be generalized to other forms of speech production. In addition, the findings reported here
306 support the hypothesis that articulatory changes which have been found in Western
307 classical singers like resonance tuning (e.g., shifting f_{R1} to the vicinity of a higher f_o), must
308 indeed have a strong effect on the identifiability of vowels.

309 Given the degree to which the vocal tract transfer function is undersampled at an f_o
310 around 1 kHz a significant loss of formant information has to be considered as very likely
311 (e.g., here, the vowels' typical medians of F_1 are exceeded by about 220–660 Hz, and there
312 is only one harmonic every 1 kHz). Although it is possible that the loss of formant
313 information can explain the decreasing identification performance, it seems likely that
314 formants cannot be the primary acoustic correlates for vowel category perception at very
315 high f_o s.

316 Calculations of auditory excitation patterns for the eight vowels at an f_o of 988 Hz,
317 revealed maximally distinct excitation levels in the frequency region above roughly 1.5 kHz
318 for the point vowels /i a u/. Excitation patterns of the other vowels have been found to
319 exhibit very similar spectral shapes as those of the point vowels they have been confused
320 with most often. Both the excitation patterns of /u/ and /o/, for example, show relatively
321 low excitation in the frequency region above 1.5 kHz, but the identification rate of /u/
322 (about 75% correct) was considerably higher than that of /o/ (about 10% correct), while a
323 substantial proportion of responses (about 43%) were /u/ when /o/ was presented. As
324 similar observations were found for other non-point and point vowel combinations, it seems
325 likely that distinctive excitation patterns can be used by listeners as landmarks (in terms of
326 reference points) for vowel category perception at high f_o s.

327 Using distinctive excitation patterns as landmarks for vowel identification could also

328 explain most of the findings reported in earlier studies on vowel identification at high f_o s.
329 Regarding the vowels used by Smith and Scott (1980) in their perception experiment (i.e.,
330 /i ɪ ε æ/), it is possible that the information conveyed by the distinct spectral shapes
331 might have been sufficient for the listeners to distinguish at least between the two pairs /i
332 ɪ/ and /ε æ/. However, it is difficult to draw conclusions from this as vowel duration
333 differed substantially in this study, and not enough detail about performance with the
334 different vowels and the instructions given to the listeners were provided.

335 Comparing the results of the present study to those reported by Friedrichs et al.
336 (2015b), the diverging identification performance for the vowel /o/ is surprising. While a
337 perfect identification rate (100% correct) was found at an f_o of 880 Hz by Friedrichs et al.
338 (2015b), a performance near chance (17.5% correct) was observed in the present study.
339 Although the lack of between-talker acoustic vowel variation (as being a single talker
340 study) and secondary cues to vowel identity (vowels were presented in word context) in the
341 former study might have helped listeners to perform better it seems possible that this
342 difference is also due to the importance of perceptual and acoustic landmarks. The
343 strongest support for this hypothesis is the fact that the vowel /u/ was not included in the
344 study of Friedrichs et al. (2015b), and thus, a confusion of /o/ and /u/ like the one found
345 in the present study was not possible (e.g., /u/ received more than 50% of the responses
346 for the intended vowel /o/ at an f_o of 880 Hz). It seems, therefore, likely that listeners

347 used the vowel /o/ as a substitute because /u/ was not presented to them as a response
348 option. The results by Friedrichs et al. (2015a), who found the same eight vowels used in
349 the present study identifiable up to an f_o of 880 Hz when recorded in minimal pairs and
350 tested in a two-alternative forced choice task, could also be explained within this context.
351 As a single talker was asked to produce several different two-word combinations containing
352 a vowel in contrastive position (e.g., the German words *Buden* vs. *Boden*), it is possible
353 that the talker produced vowels with acoustic features alike or different from those of a
354 point vowel at higher f_o s to make them distinguishable (e.g., producing an /o/ more
355 toward /a/ to distinguish it from /u/). This way the phonological function of vowels in
356 linguistic contrastive positions could be maintained for all vowels even at very high f_o s.
357 Given this, it is plausible that the number of response options has a strong effect on
358 listeners' identification performance, and obviously, a better performance should be
359 expected when fewer responses options are provided.

360 It is possible that the results presented here may have been driven in part by the
361 relative frequency of German vowels. For example, in German, /i/ is more frequent than
362 /y/, and /u/ is more frequent than /o/ (Pätzold and Simpson, 1997). Forced to choose
363 between two vowels that otherwise match the spectral characteristics of the stimulus
364 equally well, listeners are most likely to pick the one with the higher a priori probability.
365 However, it is unlikely that this can explain listeners' identification performance entirely as,

366 for example, the long /e/ is more frequent than the long /a/, with which it has been
367 confused most often in this study at an f_o of 988 Hz. In addition, relative frequency may
368 be the driving force behind which vowel label is applied to a cluster of similar vowels, but
369 it cannot explain the fact that vowels were categorized into three distinct groups.

370 In summary, the results presented here make it clear that a theory of vowel perception
371 based solely on formant peak patterns cannot account for the relatively preserved
372 performance listeners demonstrate in identifying vowels at high f_o s. Formal modelling of
373 the relationship between the perceptual and physical spaces of vowels at high and low f_o s
374 are required for a convincing demonstration, but it seems likely that overall spectral shape
375 features will play an important role in a coherent account of vowel perception generally.

376 **Acknowledgements**

377 This study was supported by the Forschungskredit of the University of Zurich, Grant
378 No. FK-14-062, and the Swiss National Science Foundation (SNSF), Grants No.
379 P2ZHP1_168375 and 100016_143943/1. Thanks to Nick Clark, whose software was used to
380 perform the gammatone filtering, and Sandra Schwab for her helpful contributions and
381 comments on an earlier draft of this paper.

382 **REFERENCES**

383 Baayen, R. H., Davidson, D. J., and Bates, D. M. (2008). "Mixed-effects modeling
384 with crossed random effects for subjects and items," *J. Mem. Lang.* **59**(4),
385 390–412.

386 Boersma, P., and Weenink, D. (2016). "Praat: Doing phonetics by computer
387 [Computer program]," Version 6.0.15, retrieved March 23, 2016 from
388 <http://www.praat.org/> (Last viewed April 30, 2016).

389 Boersma, P. (1993). "Accurate short-term analysis of the fundamental frequency
390 and the harmonics-to-noise ratio of a sampled sound," *Proceedings of the Institute
391 of Phonetic Sciences* (17), University of Amsterdam, 97–110.

392 de Cheveigné, and Kawahara, H. (1999). "Missing-data model of vowel
393 identification," *J. Acoust. Soc. Am.* **105**, 3497–3508.

394 Davis, S., and Mermelstein, P. (1980). Comparison of parametric representations of
395 monosyllabic word recognition in continuously spoken sentences. *Proceedings of
396 the IEEE Interational Conference on Acoustics, Speech and Signal Processing* **28**,
397 357–366.

398 Friedrichs, D., Maurer, D., and Dellwo, V. (2015a). "The phonological function of

- 399 vowels is maintained at fundamental frequencies up to 880Hz,” *J. Acoust. Soc.*
400 *Am.* **138** , EL36–EL42.
- 401 Friedrichs, D., Maurer, D., Suter, H., and Dellwo, V. (**2015b**). ”Vowel identification
402 at high fundamental frequencies in minimal pairs,” *Proc. 18th Int. Congr.*
403 *Phonetic Sci.*, paper number 0438, 1–5.
- 404 Glasberg, B. R., and Moore, B. C. J. (**1990**). ”Derivation of auditory filter shapes
405 from notched-noise data,” *Hearing Research* **47**, 103–138.
- 406 Hillenbrand, J., Getty, L. A., Clark, M. J., and Wheeler, K. (**1995**). ”Acoustic
407 characteristics of American English vowels,” *J. Acoust. Soc. Am.* **97**, 3099–3111.
- 408 Hillenbrand, J. M., and Houde, R. A. (**2003**). ”A narrow band pattern-matching
409 model of vowel perception,” *J. Acoust. Soc. Am.* **113**, 1044–1055.
- 410 Howie, J., and Delattre, P. (**1962**). ”An experimental study of the effect of pitch on
411 the intelligibility of vowels,” *Natl. Assoc. Teachers Singing Bull.* **18**(4), 6–9.
- 412 Ito, M., Tsuchida, J., and Yano, M. (**2001**). ”On the effectiveness of whole spectral
413 shape for vowel perception,” *J. Acoust. Soc. Am.* **110**(2), 1141–1149.
- 414 Joliveau, E., Smith, J., and Wolfe, J. (**2004**). ”Vocal tract resonances in singing: the
415 soprano voice,” *J. Acoust. Soc. Am.* **116**, 2434–2439.

416 Kiefte, M., Neary, T. M., and Assmann, P. F. (2013). "Vowel Perception in Normal
417 Speakers," Handbook of Vowels and Vowel Disorders, edited by M. J. Ball, and
418 F. E. Gibbon (Taylor & Francis, LLC, New York), pp. 161–185.

419 Klein, W., Plomp, R., and Pols, L. C. (1970). "Vowel spectra, vowel spaces, and
420 vowel identification," J. Acoust. Soc. Am. **48**(4B), 999–1009.

421 Kuznetsova, A., Brockhoff, P. B., and Christensen, R. H. B. (2014). lmerTest: Tests
422 in Linear Mixed Effects Models. R package version 2.0–20. Retrieved from
423 <http://CRAN.R-project.org/package=lmerTest>. (Accessed on June 30, 2016)

424 Lehiste, I., and Peterson, G. E. (1961). "Transitions, glides, and diphthongs," J.
425 Acoust. Soc. Am. **33**, 268–277.

426 Maurer, D., and Landis, T. (1996). "Intelligibility and spectral differences in
427 high-pitched vowels," Folia Phoniatr. Logop. **48**, 1–10.

428 Maurer, D., Mok, P., Friedrichs, D., and Dellwo, V. (2014). "Intelligibility of
429 high-pitched vowel sounds in the singing and speaking of a female Cantonese
430 Opera singer," Proceedings of the Fifteenth Annu. Conf. Int. Speech Commun.
431 Assoc. Singapore, 2132–2133.

432 Maurer, D., Suter, H., Friedrichs, D., and Dellwo, V. (2016). "Acoustic

433 characteristics of voice in music and straight theatre: topics, conceptions,
434 questions,” Trends in Phonetics and Phonology. Studies from German speaking
435 Europe, edited by A. Leemann, M. J. Kolly, S. Schmid, and V. Dellwo (Peter
436 Lang, Bern/Frankfurt), pp. 256–265.

437 Maurer, D. (2016). *Acoustics of the Vowel - Preliminaries* (Peter Lang AG,
438 International Academic Publishers, Bern).

439 Pätzold, M., and Simpson, A. (1997). ”Acoustic analysis of German vowels in the
440 Kiel Corpus of read speech,” Arbeitsberichte des Instituts für Phonetik und Digit.
441 Sprachverarbeitung Univ. Kiel **32**, 215–247.

442 Pols, L. C., Van der Kamp, L. T., and Plomp, R. (1969). ”Perceptual and physical
443 space of vowel sounds,” J. Acoust. Soc. Am. **46**(2B), 458–467.

444 Puria, S., Peake, W. T., and Rosowski, J. J. (1997) Sound-pressure measurements in
445 the cochlear vestibule of human cadaver ears, Journal Of the Acoustical Society Of
446 America **101**, 2754–2770.

447 R Core Team. (2016). R: A Language and Environment for Statistical Computing.
448 Version 3.1.3. [Computer software] Vienna: R Foundation for Statistical
449 Computing. Retrieved from <https://www.r-project.org> (Accessed on June 30,
450 2016)

- 451 Smith, J., and Wolfe, J. (2009). "Vowel-pitch matching in Wagner's operas:
452 Implications for intelligibility and ease of singing," *J. Acoust. Soc. Am.* **125**,
453 EL196–EL201.
- 454 Smith, L. A., and Scott, B. L. (1980). "Increasing the intelligibility of sung vowels,"
455 *J. Acoust. Soc. Am.* **67**, 1795–1797.
- 456 Sundberg, J. (1975). "Formant technique in a professional female singer," *Acustica*
457 **32**, 89–96.
- 458 Sundberg, J. (2013). "Perception of singing," in *Psychology of Music*, 3rd ed., edited
459 by D. Deutsch (Academic Press, London), pp. 69–106.
- 460 Zahorian, S., and Jagharghi, A. (1993). "Spectral-shape features versus formants as
461 acoustic correlates for vowels," *J. Acoust. Soc. Am.* **94**, 1966–82.

462 Appendices

463 **A.** Confusion matrices for each f_o containing the raw data of the identification test in
 464 percentages.

$f_o \approx 220$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	79.4	0	20.6	0	0	0	0	0
/y/	15.9	73	3.2	7.90	0	0	0	0
/e/	0	0	100	0	0	0	0	0
/ø/	0	0	6.3	93.7	0	0	0	0
/ɛ/	0	0	7.9	0	92.1	0	0	0
/a/	0	0	0	0	1.6	98.4	0	0
/o/	0	0	0	0	0	0	100	0
/u/	0	0	0	0	0	0	11.1	88.9
response proportions	11.9	9.10	17.3	12.7	11.7	12.3	13.9	11.1

$f_o \approx 330$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	88.9	6.3	4.8	0	0	0	0	0
/y/	4.8	92.1	0	1.6	1.6	0	0	0
/e/	1.6	0	98.4	0	0	0	0	0
/ø/	0	0	0	92.1	0	4.8	3.2	0
/ɛ/	0	0	3.2	1.6	88.9	6.3	0	0
/a/	0	0	0	0	0	100	0	0
/o/	0	0	1.6	0	0	0	90.5	7.9
/u/	0	0	0	0	0	0	14.3	85.7
response proportions	11.9	12.3	13.5	11.9	11.3	13.9	13.5	11.7

$f_o \approx 440$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	76.2	7.9	6.3	4.8	0	0	0	4.8
/y/	4.8	84.1	0	11.1	0	0	0	0
/e/	4.8	1.6	87.3	3.2	3.2	0	0	0
/ø/	0	15.9	0	71.4	3.2	6.3	3.2	0
/ɛ/	0	0	1.6	4.8	68.3	20.6	3.2	1.6
/a/	0	0	0	0	1.6	96.8	1.6	0
/o/	1.6	0	0	0	0	4.8	90.5	3.2
/u/	0	1.6	0	1.6	0	0	9.5	87.3
response proportions	10.9	13.9	11.9	12.1	9.5	16.1	13.5	12.1

$f_o \approx 523$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	73	11.1	6.3	1.6	0	1.6	0	6.3
/y/	1.6	60.3	4.8	15.9	0	0	1.6	15.9
/e/	15.9	12.7	49.2	7.9	9.5	3.2	0	1.6
/ø/	0	12.7	1.6	50.8	17.5	12.7	1.6	3.2
/ɛ/	0	0	0	1.6	77.8	20.6	0	0
/a/	0	0	0	0	4.8	92.1	3.2	0
/o/	0	0	0	0	0	42.9	57.1	0
/u/	0	0	0	0	0	1.6	22.2	76.2
response proportions	11.3	12.1	7.7	9.7	13.7	21.8	10.7	12.9

$f_o \approx 587$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	88.9	0	7.9	1.6	0	0	0	1.6
/y/	12.7	61.9	19	4.8	0	1.6	0	0
/e/	6.3	11.1	55.6	15.9	7.9	1.6	0	1.6
/ø/	0	22.2	1.6	69.8	0	4.8	0	1.6
/ɛ/	0	0	11.1	0	79.4	0	6.3	3.2
/a/	0	0	0	0	1.6	95.2	3.2	0
/o/	0	1.6	0	1.6	0	30.2	52.4	14.3
/u/	0	0	0	1.6	0	3.2	14.3	81
response proportions	13.5	12.1	11.9	11.9	11.1	17.1	9.5	12.9

$f_o \approx 659$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	88.9	1.6	4.8	0	3.2	0	0	1.6
/y/	3.2	61.9	4.8	20.6	7.9	0	0	1.6
/e/	14.3	11.1	47.6	7.9	14.3	0	3.2	1.6
/ø/	0	38.1	1.6	47.6	1.6	1.6	1.6	7.9
/ɛ/	0	0	0	3.2	87.3	7.9	1.6	0
/a/	0	0	1.6	1.6	6.3	90.5	0	0
/o/	1.6	3.2	3.2	3.2	0	6.3	46	36.5
/u/	0	4.8	0	1.6	1.6	1.6	20.6	69.8
response proportions	13.5	15.1	8	10.7	15.3	13.5	9.1	14.9

$f_o \approx 698$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	92.1	0	3.2	0	1.6	3.2	0	0
/y/	6.3	68.3	6.3	7.9	9.5	0	0	1.6
/e/	33.3	15.9	38.1	4.8	6.3	0	0	1.6
/ø/	7.9	14.3	22.2	36.5	0	0	1.6	17.5
/ɛ/	0	0	0	0	93.7	6.3	0	0
/a/	0	1.6	3.2	3.2	6.3	84.1	1.6	0
/o/	0	0	1.6	1.6	0	6.3	33.3	57.1
/u/	0	0	0	0	0	6.3	20.6	73
response proportions	17.5	12.5	9.3	6.8	14.7	13.3	7.1	18.9

$f_o \approx 784$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	93.7	0	4.8	0	1.6	0	0	0
/y/	15.9	65.1	9.5	1.6	7.9	0	0	0
/e/	14.3	9.5	58.7	6.3	9.5	0	1.6	0
/ø/	0	3.2	7.9	19	14.3	14.3	12.7	28.6
/ɛ/	4.8	3.2	12.7	3.2	76.2	0	0	0
/a/	0	1.6	1.6	0	9.5	82.5	3.2	1.6
/o/	0	3.2	1.6	0	0	4.8	22.2	68.3
/u/	0	0	0	0	1.6	3.2	15.9	79.4
response proportions	16.1	10.7	12.1	3.8	15.1	13.1	7	22.2

$f_o \approx 880$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	82.5	6.3	0	0	11.1	0	0	0
/y/	30.2	47.6	3.2	3.2	15.9	0	0	0
/e/	9.5	11.1	30.2	11.1	33.3	3.2	0	1.6
/ø/	4.8	11.1	7.9	22.2	22.2	11.1	6.3	14.3
/ɛ/	1.6	0	6.3	0	76.2	12.7	0	3.2
/a/	0	0	3.2	0	11.1	81	3.2	1.6
/o/	3.2	4.8	3.2	4.8	0	15.9	17.5	50.8
/u/	0	1.6	0	1.6	0	1.6	7.9	87.3
response proportions	16.5	10.3	6.8	5.4	21.2	15.7	4.4	19.9

$f_o \approx 988$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	95.2	1.6	1.6	0	1.6	0	0	0
/y/	20.6	49.2	15.9	1.6	12.7	0	0	0
/e/	9.5	6.3	11.1	4.8	23.8	25.4	7.9	11.1
/ø/	6.3	1.6	4.8	12.7	4.8	38.1	11.1	20.6
/ɛ/	1.6	1.6	0	0	46	47.6	3.2	0
/a/	0	0	3.2	1.6	9.5	76.2	6.3	3.2
/o/	6.3	1.6	3.2	3.2	7.9	30.2	4.8	42.9
/u/	3.2	3.2	1.6	0	1.6	6.3	9.5	74.6
response proportions	17.8	8.1	5.2	3	13.5	28	5.4	19.1

$f_o \approx 1046$ Hz	/i/	/y/	/e/	/ø/	/ɛ/	/a/	/o/	/u/
/i/	95.2	1.6	0	0	3.2	0	0	0
/y/	44.4	38.1	7.9	0	6.3	1.6	1.6	0
/e/	9.5	6.3	3.2	7.9	36.5	31.7	3.2	1.6
/ø/	6.3	28.6	1.6	19	17.5	17.5	1.6	7.9
/ɛ/	6.3	11.1	0	4.8	41.3	33.3	0	3.2
/a/	0	3.2	1.6	6.3	19	68.3	1.6	0
/o/	11.1	4.8	3.2	4.8	6.3	38.1	4.8	27
/u/	4.8	1.6	1.6	0	4.8	15.9	1.6	69.8
response proportions	22.2	11.9	2.4	5.4	16.9	25.8	1.8	13.7