

SUPPLEMENTARY INFORMATION FOR PERSSON & JAEGER (2023).
 EVALUATING NORMALIZATION ACCOUNTS AGAINST THE DENSE VOWEL
 SPACE OF CENTRAL SWEDISH

1090 Both the main text and these supplementary information (SI) are derived from the same R
 1091 markdown document available via OSF at <https://osf.io/zb8gx/>.

1 REQUIRED SOFTWARE

1092 The document was compiled using knitr (Xie, 2021) in RStudio with R:

```
1093
1094 platform      - aarch64-apple-darwin20
1095 arch          aarch64
1096 os            darwin20
1097 system        aarch64, darwin20
1098 status
1099 major         4
1100 minor         2.2
1101 year          2022
1102 month         10
1103 day           31
1104 svn rev       83211
1105 language      R
1106 version.string R version 4.2.2 (2022-10-31)
1107 nickname      Innocent and Trusting
```

1108 We used the following R packages to create this document: R (Version 4.2.2; R Core Team,
 1109 2021b) and the R-packages *assertthat* (Version 0.2.1; Wickham, 2019a), *brms* (Version 2.19.0;
 1110 Bürkner, 2017, 2018, 2021), *cowplot* (Version 1.1.1; Wilke, 2020), *data.table* (Version 1.14.8;
 1111 Dowle and Srinivasan, 2021), *dipptest* (Maechler, 2021), *dplyr* (Version 1.1.2; Wickham et al.,
 1112 2021a), *forcats* (Version 1.0.0; Wickham, 2021a), *ggridge* (Version 1.0.8; Pedersen and Robinson,
 1113 2020), *ggplot2* (Version 3.4.2; Wickham, 2016), *LaplaceDemon* (Version 16.1.6; Statisticat and
 1114 LLC., 2021), *latexpdf* (Version 0.1.0; Hugh-Jones, 2021), *linguisticsdown* (Version 1.2.0; Liao,
 1115 2019), *lme4* (Version 1.1.32; Bates et al., 2015), *magick* (Ooms, 2021), *magrittr* (Version 2.0.3;
 1116 Bache and Wickham, 2020), *Matrix* (Version 1.5.4; Bates and Maechler, 2021), *modelr* (Version
 1117 0.1.11; Wickham, 2020), *MVBeliefUpdatr* (Version 0.0.1.2; Kleinschmidt and Jaeger, 2015b), *papaja*
 1118 (Version 0.1.1.9001; Aust and Barth, 2020), *plotly* (Version 4.10.1; Sievert, 2020), *processx* (Version
 1119 3.8.1; Csárdi and Chang, 2021), *purrr* (Version 1.0.1; Henry and Wickham, 2020), *Rcpp* (Version
 1120 1.0.10; Eddelbuettel and François, 2011; Eddelbuettel and Balamuta, 2018), *readr* (Version 2.1.4;
 1121 Wickham et al., 2021b), *rlang* (Version 1.1.1; Henry and Wickham, 2021), *stringr* (Version 1.5.0;
 1122 Wickham, 2019b), *tibble* (Version 3.2.1; Müller and Wickham, 2021), *tidyR* (Version 1.3.0; Wickham,
 1123 2021b), *tidyverse* (Version 2.0.0; Wickham et al., 2019), *tinylabels* (Version 0.2.3; Barth, 2022), and
 1124 *tufte* (Xie and Allaire, 2022).



anna.persson@su.se

Figure S1. Example flyer for recruiting Stockholm Swedish talkers for recording of the SwehVd database.

2 ADDITIONAL INFORMATION ABOUT THE SWEHVD DATABASEE

1125 2.1 Participant recruitment

1126 Participants were recruited through word-of-mouth, flyers at Stockholm University Campus,
 1127 and online channels (accindi.se). Figure S1 is an example of flyers distributed at Stockholm
 1128 University Campus. The flyer gives information on criteria for participation, recording procedure,
 1129 Word list with all target and filler words, recorded by all talkers in the SwehVd database,
 reimbursement and contact information to experimenter (first author).

Table S1. Words recorded by the female talkers of Stockholm Swedish for the SwehVd database

Target words	Vowel IPA	Filler words	
hid	[i:]	titt	tand
hidd	[ɪ]	damm	dipp
hyd	[y:]	tå	buss
hydd	[ʏ]	bål	ding
hed	[e:]	dill	porr
hedd	[ɛ]	tugga	mitt
häd	[ɛ:]	mat	dopp
hädd	[ɛ]	norr	tal
härd	[æ:]	must	namn
härr	[æ]	pil	pall
höd	[ø:]	dina	bar
hödd	[ø]	biff	till
hörd	[œ:]	Tina	mål
hörر	[œ]	borr	Nina
hud	[u:]	dal	då
hudd	[θ]	Pål	nick
hod	[u:]	nunna	ditt
hodd	[ʊ]	mil	dugga
håd	[o:]	ting	mall
hådd	[ɔ:]	ball	bil
had	[a:]	piff	par
hadd	[a]	tipp	morr

1132 2.3 Unanticipated challenges during recording (and how they were addressed)

1133 In a small-scale pilot preceding recordings, the expected transparency of the orthography for
 1134 eliciting the long and short vowels was confirmed by three native talkers and one non-native talker
 1135 of Swedish (these talkers did not participate in the study). However, *hodd* [ʊ] and *hod* [u:] sometimes
 1136 elicited [ɔ].^{S1} We therefore decided to add instructions to the participants for these two words.
 1137 When *hod* or *hodd* appeared on screen, a written guide indicating the target vowel appeared below
 1138 the word in smaller font size: “*hod som i hot*”, “*hodd som i hosta*”, with *hot* and *hosta* being
 1139 real Swedish words containing [u:] and [ʊ], respectively.^{S2} Whenever the experimenter noticed that
 1140 the pronunciations clearly targeted another vowel, recordings were stopped and participants were
 1141 reminded to carefully read the guide. Despite our recording instructions, five of the talkers rarely
 1142 ever produced the targeted [ʊ] vowel. Instead, they often mispronounced the vowel, hence they are
 1143 not included in the subsetted SwehVd we use in this study.

1144 2.4 Neutralization of *hedd* and *hädd*

1145 For the vast majority of talkers, *hädd* productions elicited the same vowel as *hedd* (see Figure S2).
 1146 This confirms the common assumption that the short allophone to /e/ neutralizes with the short
 1147 allophone to [ɛ] in Central Swedish.

^{S1} The difficulty for some native talkers to produce [ʊ] when reading *hodd* might be due to frequency effects. Forms with stressed [ʊ] are rare in the Swedish language, and phonotactically similar words are most often pronounced as [ɔ] (see e.g., Riad, 2014).

^{S2} English translations: “*hod as in threat*”(phonologically [u:]), “*hodd as in cough*”(phonologically [ʊ]).

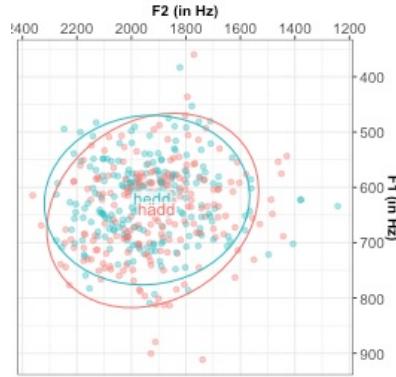


Figure S2. The *hadd* and *hedd* words in the SwehVd vowel data in unnormalized F1-F2 space. Points show recordings of the *hadd* and *hedd* words ([ɛ]) by the 24 female native talkers in the database, averaged across the five measurement points within each vowel segment. Word labels indicate word means across talkers. Since *hadd* and *hedd* resulted in the same allophone, we exclude *hadd* from this and all other visualizations below. This facilitates comparison of, for example, densities across vowels (see diagonal of Figure 4).

3 EVALUATION OF IMPLEMENTATIONS OF SYRDAL & GOPAL'S (1986) SECOND DIMENSION

1148 For the second dimension, distinguishing between front and back vowels, Syrdal and Gopal (1986)
 1149 evaluates two different bark-difference measures: F2-F1 and F3-F2. Previous studies had concluded
 1150 that F2-F1 distinguishes between all Swedish vowels (Fant, 1983), however, in Syrdal and Gopal
 1151 (1986)'s evaluation of American English, the F3-F2 dimension provided a better fit. Given that
 1152 there seems to be language specific effects concerning Syrdal and Gopal (1986)'s second dimension
 1153 (e.g., Adank, 2003), here we compare the two difference measures for the vowels in the SwehVd
 1154 database.

1155 Figure S3 displays the categorization accuracy for models trained on normalized data under the
 1156 two implementations of the Syrdal & Gopal account. The first version uses the F2-F1 bark-difference
 1157 metric for the second dimension, whereas the second version (labelled *SyrdalGopal2 (Bark)*)
 1158 implements the second dimension as suggested by Syrdal and Gopal (1986), F3-F2. Figure S3
 1159 indicates that the first implementation outperforms the implementation using F3-F2 bark-difference
 1160 measure, which replicates Fant (1983).

1161 We also compared the two implementations in terms of the separability index used in the auxiliary
 1162 study (7). Figure S4 mirrors the results from the modelling—the first implementation performs
 1163 better at separating categories in the SwehVd data. These results taken together indicate that the
 1164 F2-F1 implementation is more suitable for the materials used here, we therefore decided to use the
 1165 first implementation throughout this paper.

4 VISUALIZING THE DISTRIBUTION OF VOWEL PRODUCTIONS

1166 Figures S5 and S6 visualize the Central Swedish vowels in the test data, after applying the 15
 1167 different scale-transformations and normalization accounts for a visual inspection. For this purpose,
 1168 we focus on F1 and F2 only. In Section Correlation matrices for all normalization accounts below,
 1169 we plot pairwise correlation plots of all cues for all different normalization accounts we compare.

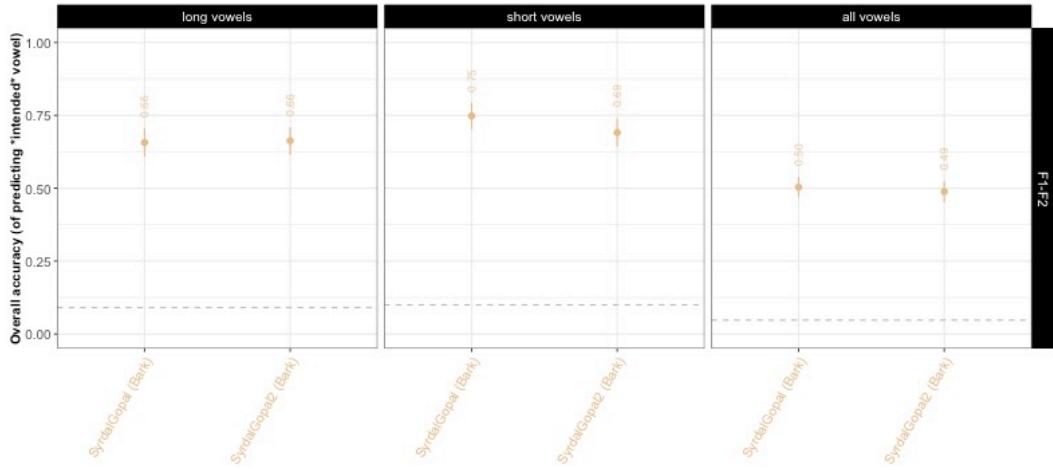


Figure S3. Predicted recognition accuracy of ideal observer under two versions of the Syrdal & Gopal (1986) account for long vowels, short vowels, and long and short vowels together, shown for the F1-F2 cue combination. Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. The dashed horizontal line indicates chance (different across columns because of the different number of long and short vowels).

1170 Visual inspection suggests a few initial observations. The most striking difference is perhaps
 1171 between intrinsic normalization accounts (Syrdal and Gopal, 1986; Miller, 1989) and all other
 1172 approaches, though it is not immediately visually obvious which type of approach achieves better
 1173 separability. Second, transforming the vowels to a different perceptual scale does not seem to
 1174 affect the vowel distributions much, besides a minor decrease in category variance for some of the
 1175 vowels. Some transformations bring the vowel categories closer together, towards the center of the
 1176 vowel space, e.g., ERB and semitones. Third, centering formants by subtracting each talkers' mean
 1177 (McMurray and Jongman, 2011; Nearey, 1978) reduces some of the category variance, and as a
 1178 result, increases the category separability. Transforming the vowel data into different scales prior
 1179 to centering also seems to further improve separability (compare e.g., C-CuRE (Hz) and C-CuRE
 1180 (semitones)). Overall, the top two performing accounts across the long and short vowels appear
 1181 to be Lobanov (1971) and Nearey (1978). However, even for the best performing normalization
 1182 accounts, there is still considerable category overlap. This involves some of the high long vowels,
 1183 and some of the mid-center short vowels. This highlights the need to more systematically quantify
 1184 the effects of normalization, as we do in this study.

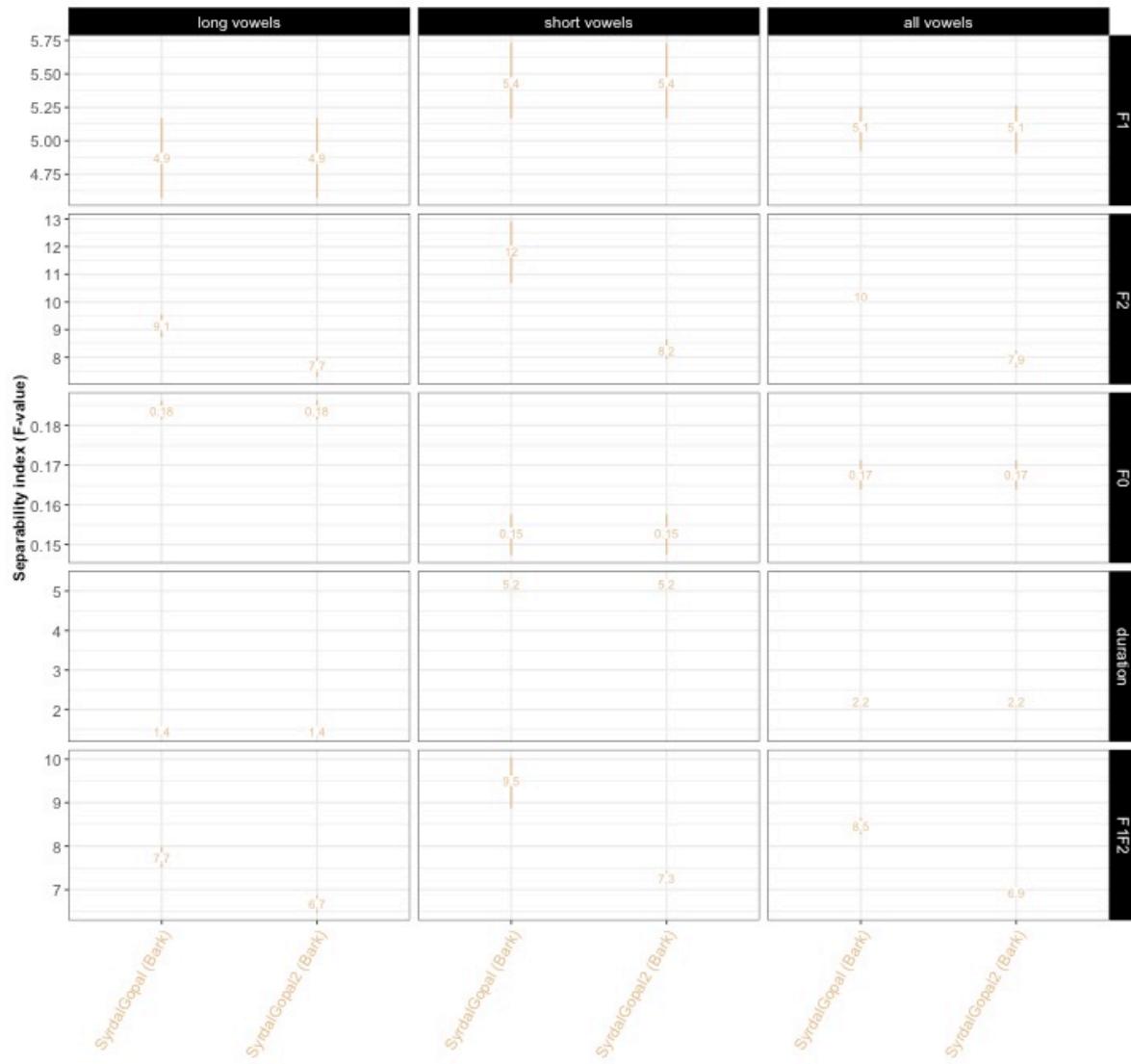


Figure S4. Separability indices of the two versions of the Syrdal & Gopal (1986) account for long vowels, short vowels, and long and short vowels together, shown for four of the five cues considered in this study and the combined F1-F2. Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. Note that the ranges of the y-axes varies across plots.

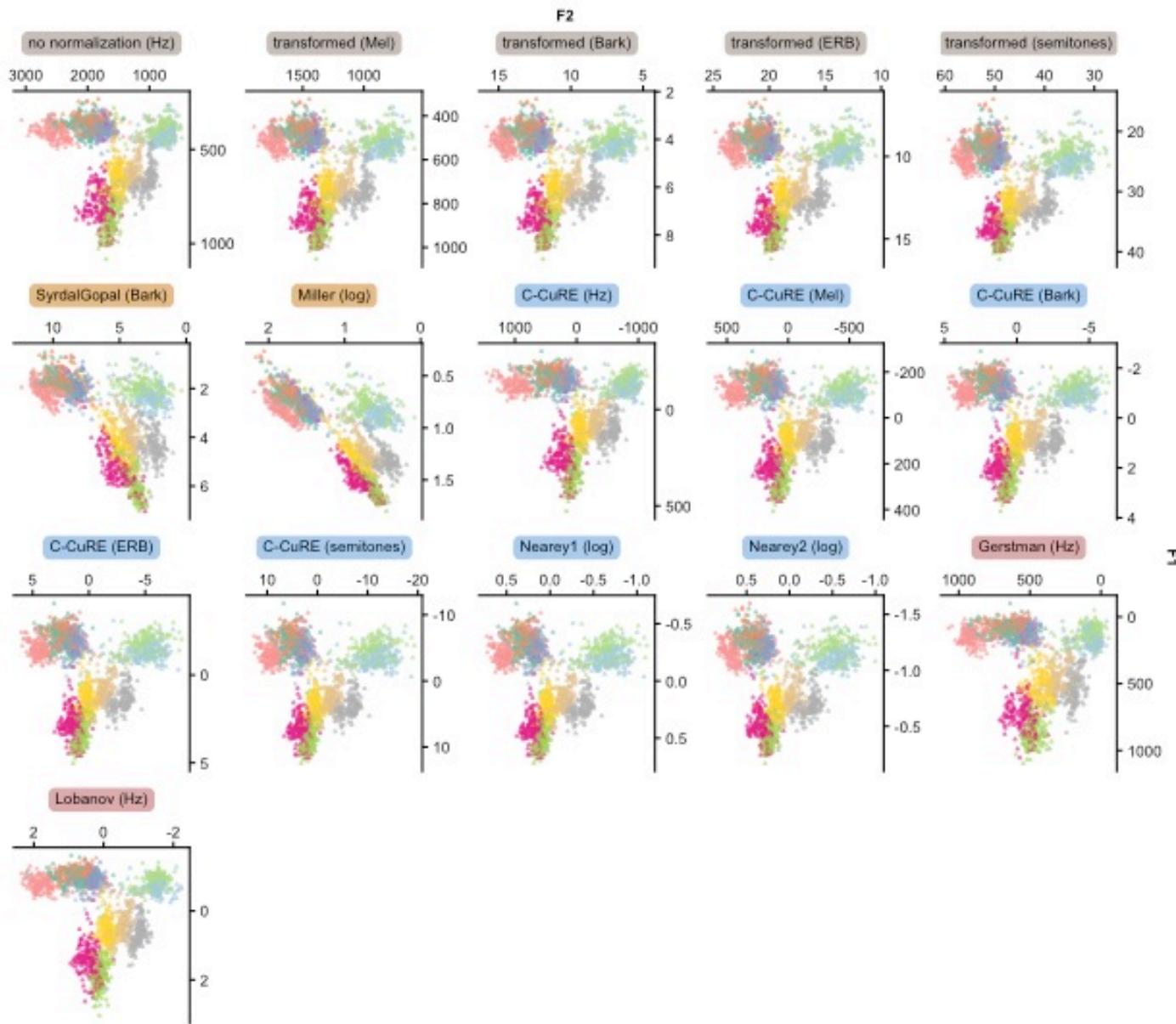


Figure S5. The 11 long vowels of Central Swedish when F1 and F2 are left unnormalized or transformed into a perceptual scales (grey), intrinsically normalized (yellow), or extrinsically normalized through centering (blue) or standardizing (purple). Each point corresponds to one recording, averaged across the five measurement points within each vowel segment. Each panel combines the data from all five test folds.

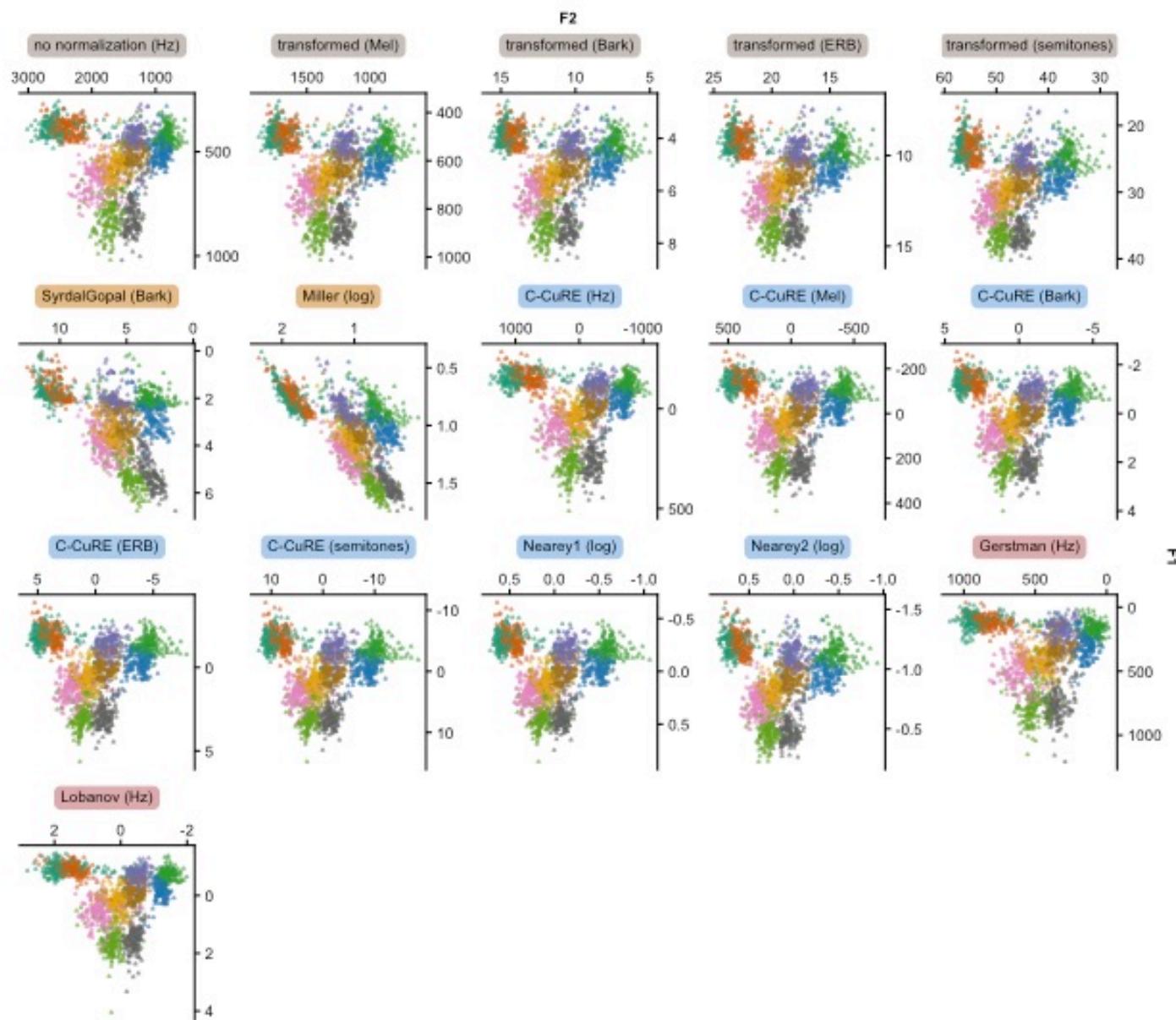


Figure S6. The 10 short vowels of Central Swedish when F1 and F2 are left unnormalized or transformed into a perceptual scales (grey), intrinsically normalized (yellow), or extrinsically normalized through centering (blue) or standardizing (purple). Each point corresponds to one recording, averaged across the five measurement points within each vowel segment. Each panel combines the data from all five test folds.

5 CUE CORRELATION MATRICES FOR ALL NORMALIZATION ACCOUNTS

1185 Here we include correlation matrices for the SwehVd vowel data, transformed into the 15 different
1186 normalization spaces.

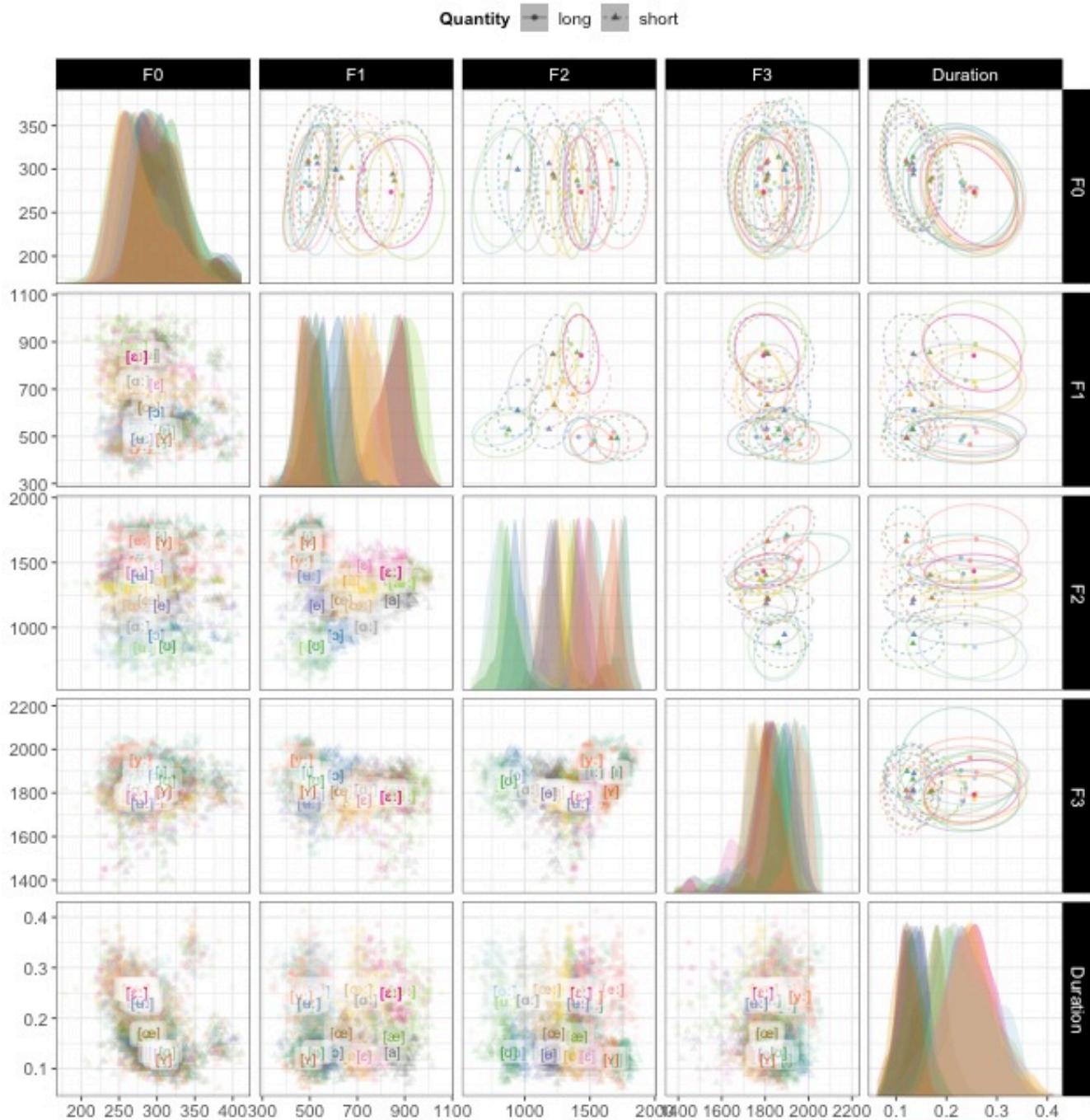


Figure S7. The SwehVd vowel data in Mel space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

6 VOWEL-SPECIFIC IDEAL OBSERVER ANALYSES

1187 The use of a perceptual model (here: ideal observers) also makes it straightforward to assess vowel-
 1188 specific effects of normalization. The next two subsections provide both the predicted categorization
 1189 accuracy per vowel in the different evaluations, as well as confusion matrices of the best and

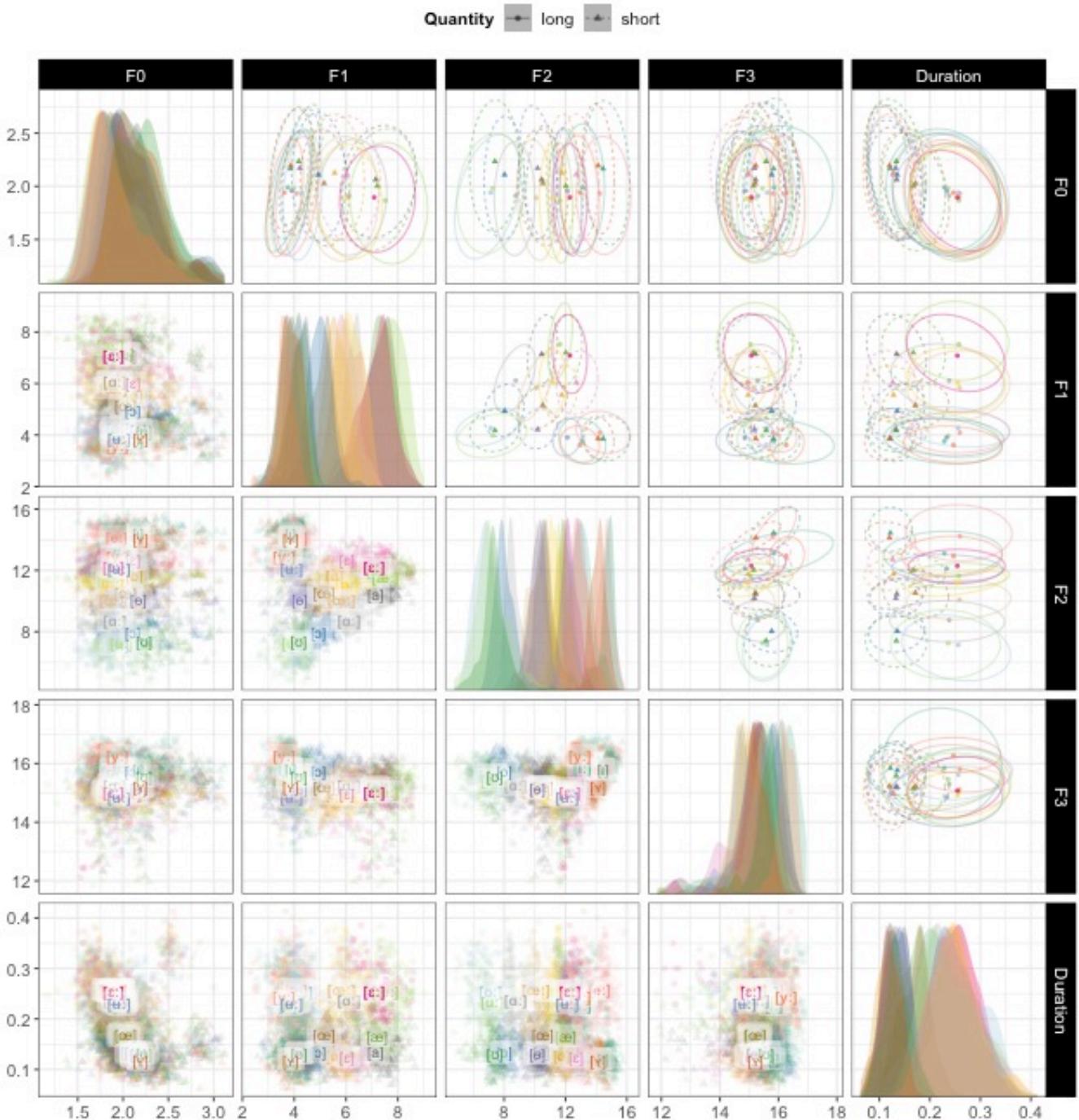


Figure S8. The SwehVd vowel data in Bark space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1190 the worst performing ideal observers, shedding light on *how* normalization improves recognition
 1191 accuracy.

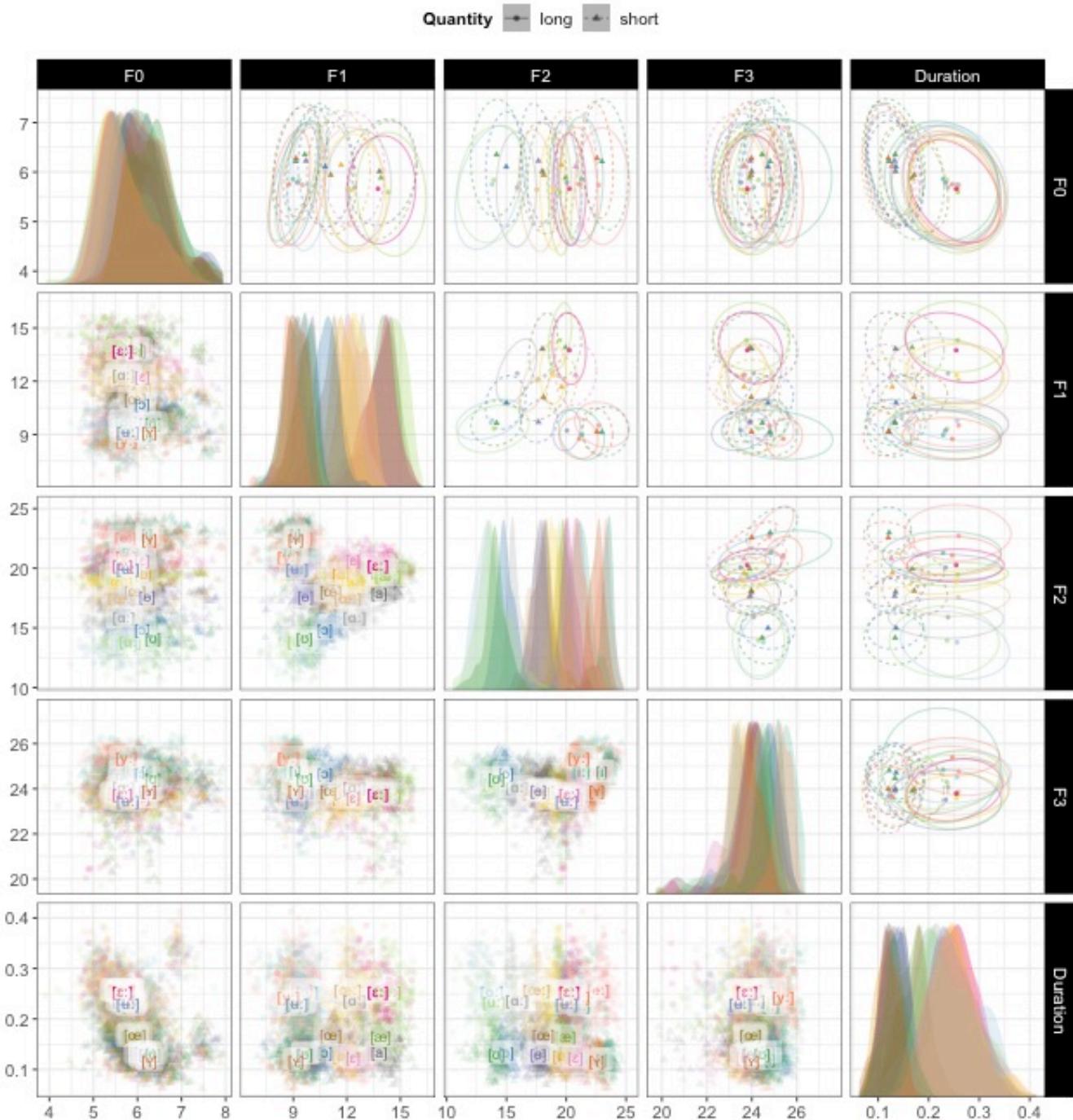


Figure S9. The SwehVd vowel data in ERB space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1192 6.1 Per-vowel categorization accuracy of models trained on long and short vowels separately

1193 Unsurprisingly, some vowels are recognized with much higher accuracy than others—at least when
 1194 uniform category priors are assumed, as we did here. This is a direct consequence of the position
 1195 of the vowel in the acoustic-phonetic space, relative to neighboring vowels: the more neighboring

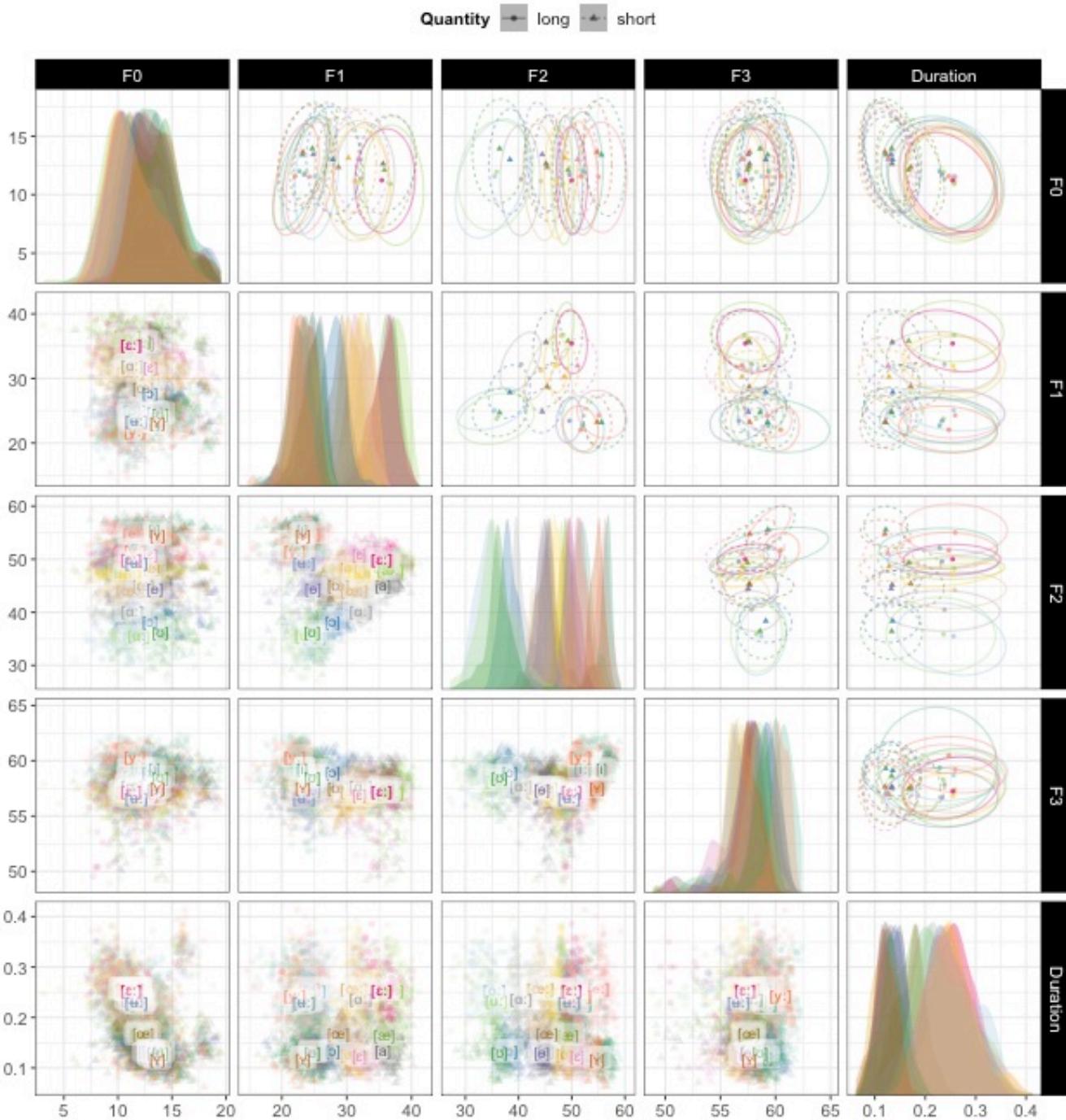


Figure S10. The SwehVd vowel data in semitones space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1196 vowels overlap with each other, the lower the accuracy with which they are recognized. Which
1197 vowels will benefit from normalization will thus naturally vary between languages, reflecting the
1198 language-specific properties of the vowel space. For instance, [i:] is often described as more easily
1199 recognized in previous work on other languages. This contrasts with our findings for Central Swedish:

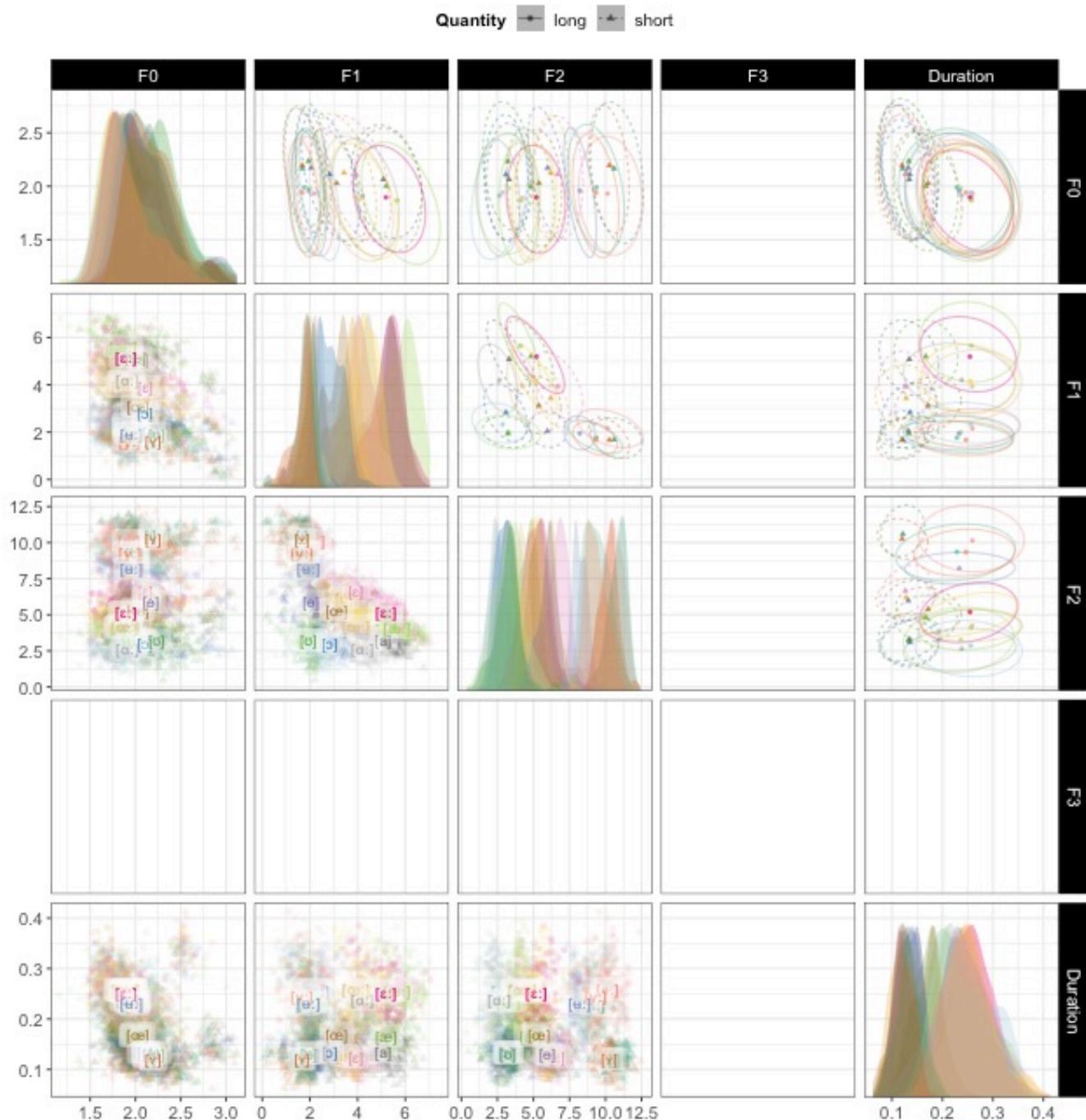


Figure S11. The SwehVd vowel data in SyrdalGopal (Bark) space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1200 here, [i:] is part of the dense clustering of vowels along the height dimension and so has many close
 1201 competitors. This highlights that recognition accuracy is due to the position of a vowel *relative* to
 1202 its competitors (e.g., Peterson and Barney, 1952; Kuhl, 1991; Polka and Bohn, 2003), rather than
 1203 its *absolute location* in the vowel space (e.g., [i:] being a peripheral vowel).

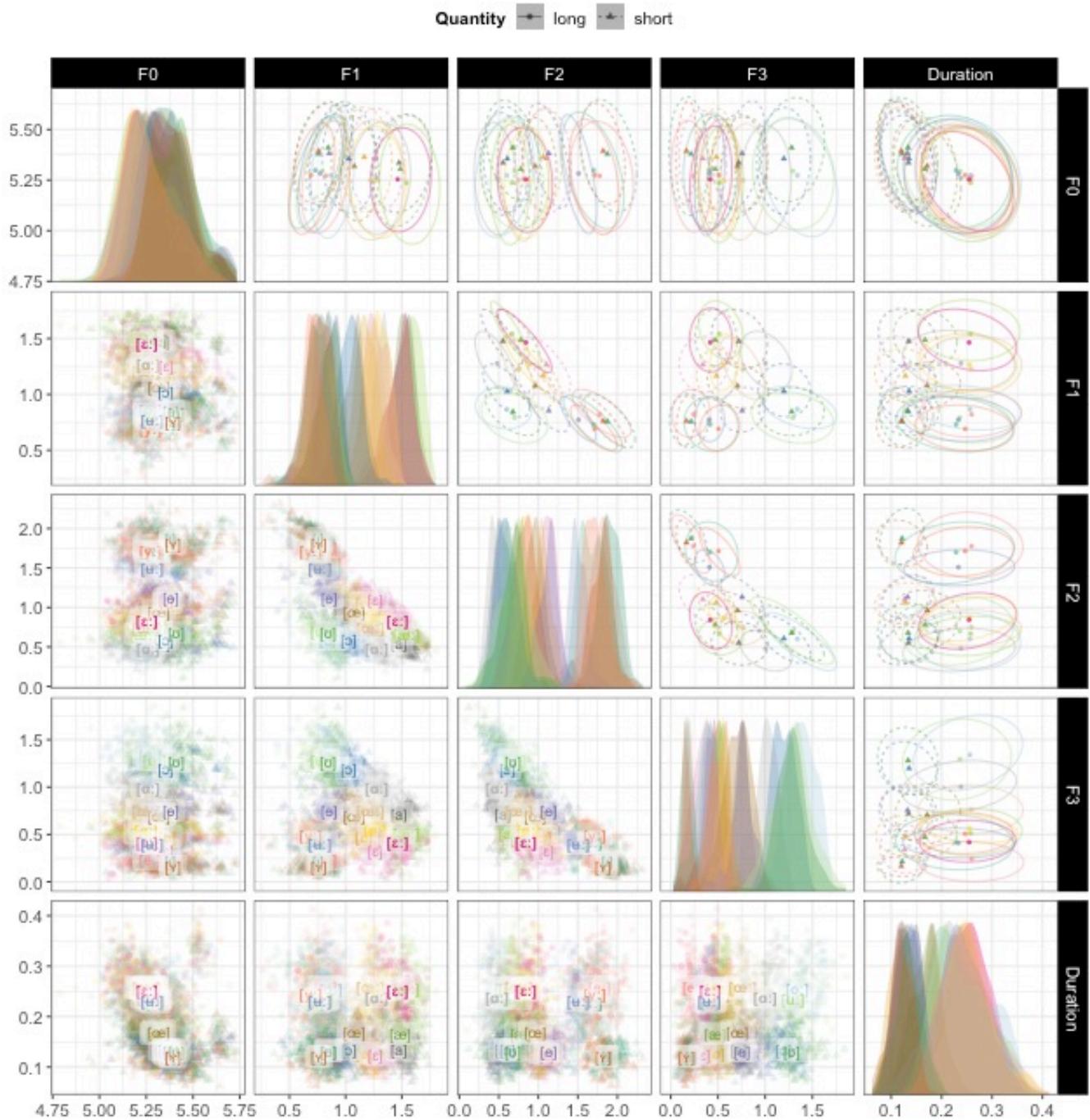


Figure S12. The SwehVd vowel data in Miller (log) space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1204 Also of interest is that not all vowels exhibit the benefit of normalization. In general, across
 1205 evaluations, it seems to be vowels that are already recognized with high accuracy that does not
 1206 benefit from normalization, which replicate previous studies that have included per-vowel accuracies
 1207 (e.g., Adank, 2003; Syrdal and Gopal, 1986). For one vowel in particular, normalization can actually

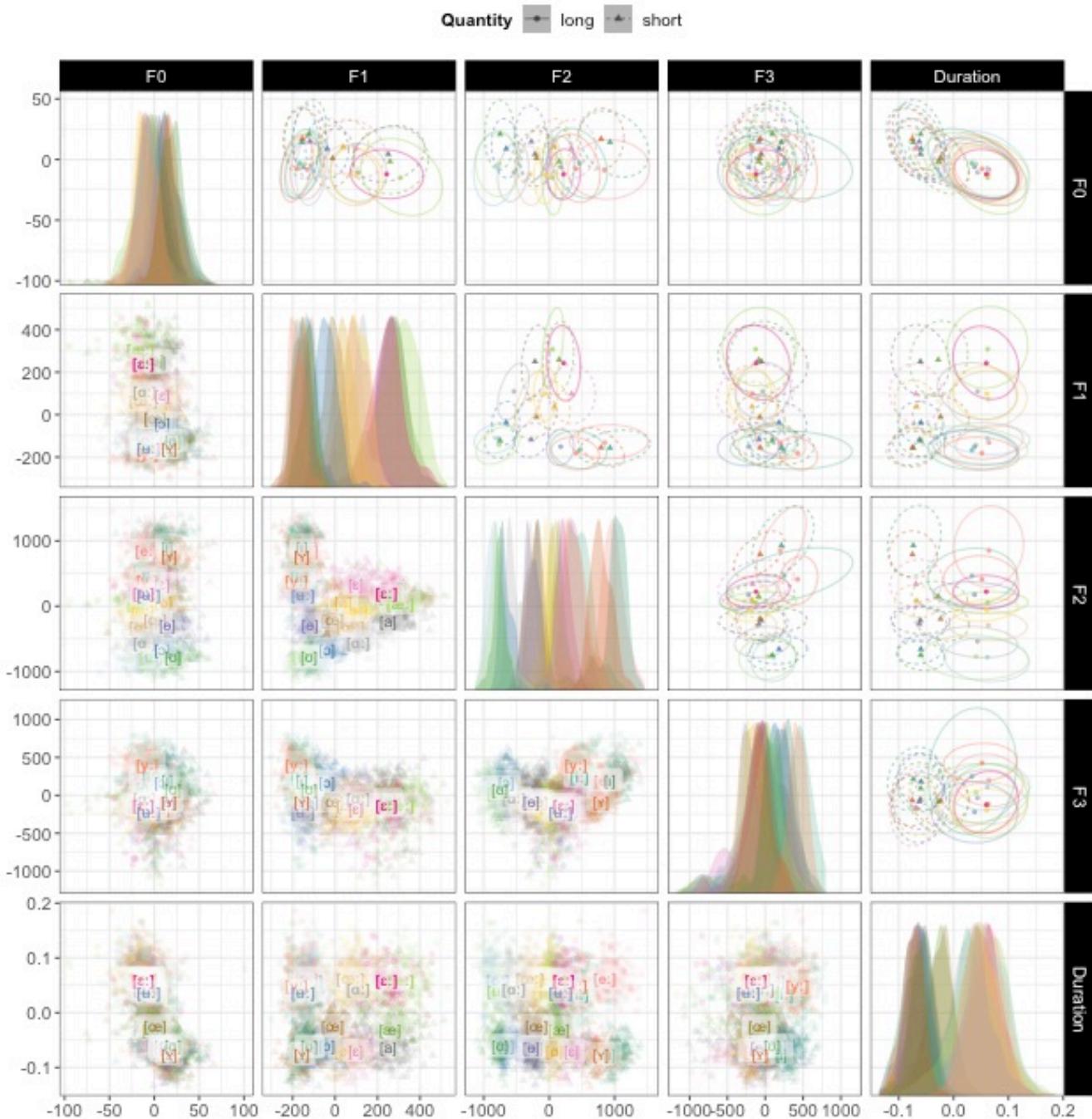


Figure S13. The SwehVd vowel data in C-CuRE Hz space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1208 be detrimental to recognition. The accuracy of some normalized models is reduced compared to
 1209 unnormalized models for [ɛ] when more cues than F1 and F2 are considered. Finally, while there
 1210 are minor differences across vowels in the relative goodness of different normalizations, the models

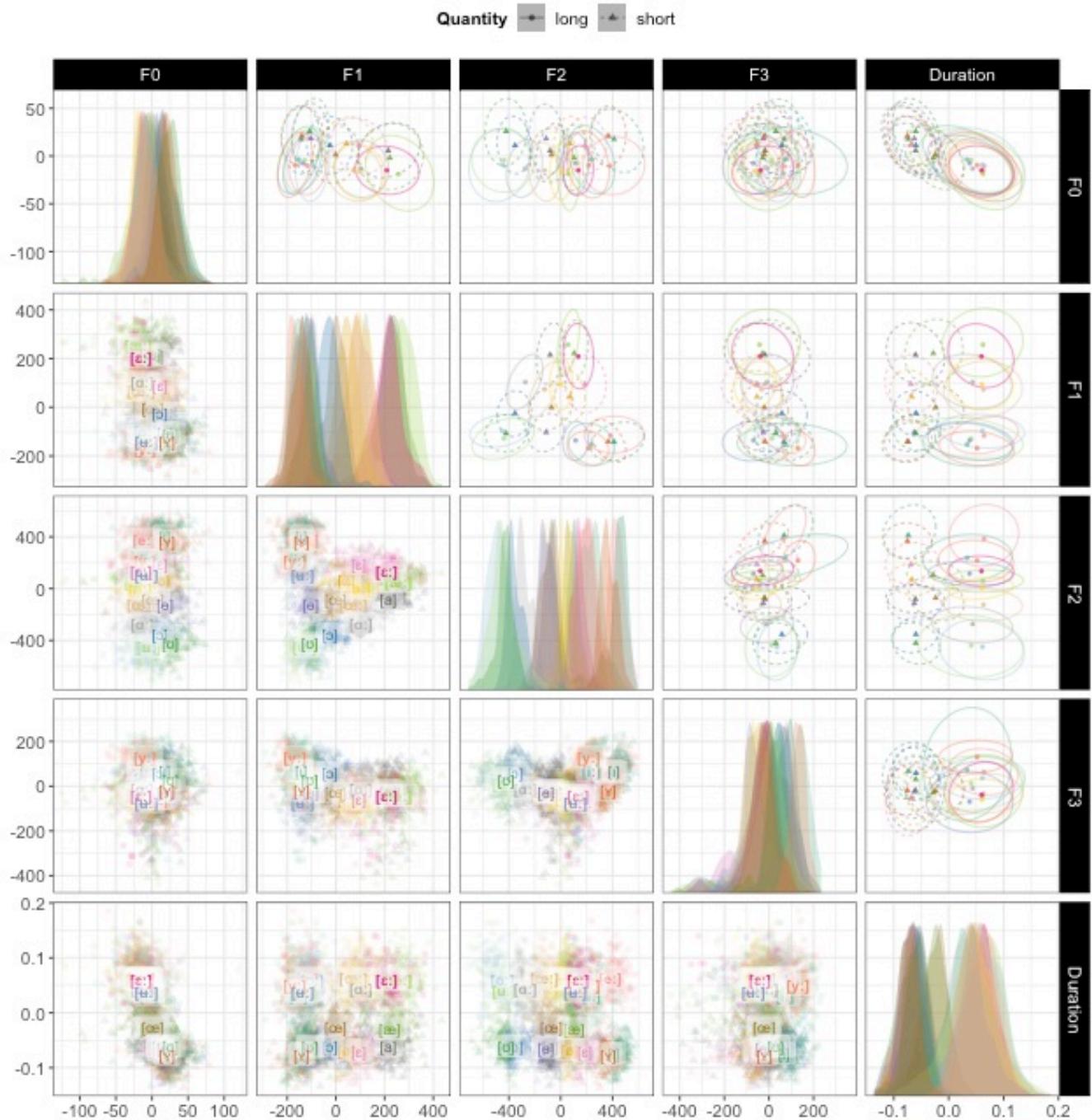


Figure S14. The SwehVd vowel data in C-CuRE Mel space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

1211 that perform overall best also perform best on each vowel (in line with Adank, 2003). This further
 1212 demonstrates the plausibility of these normalization accounts for perception.

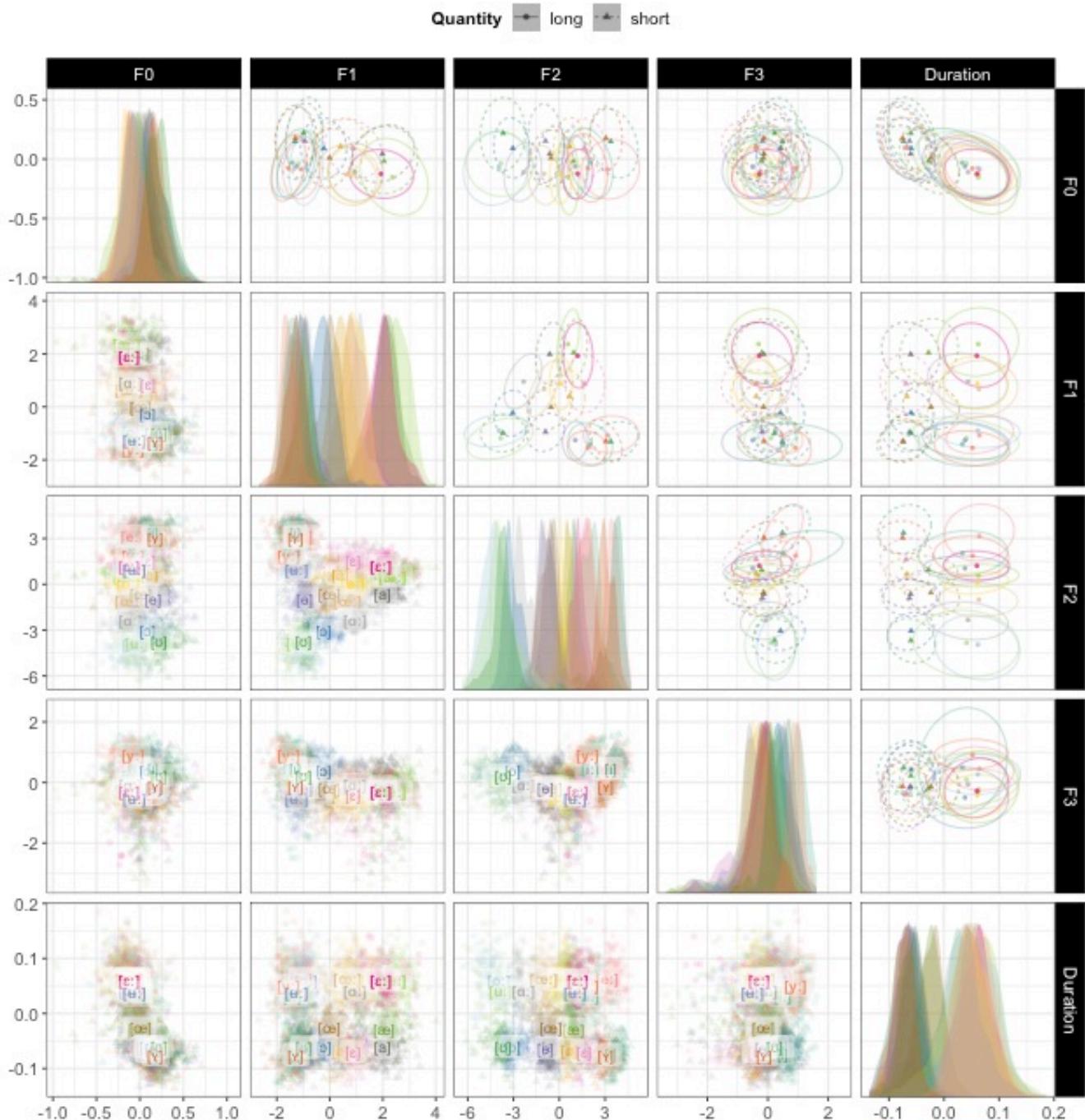


Figure S15. The SwehVd vowel data in C-CuRE Bark space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

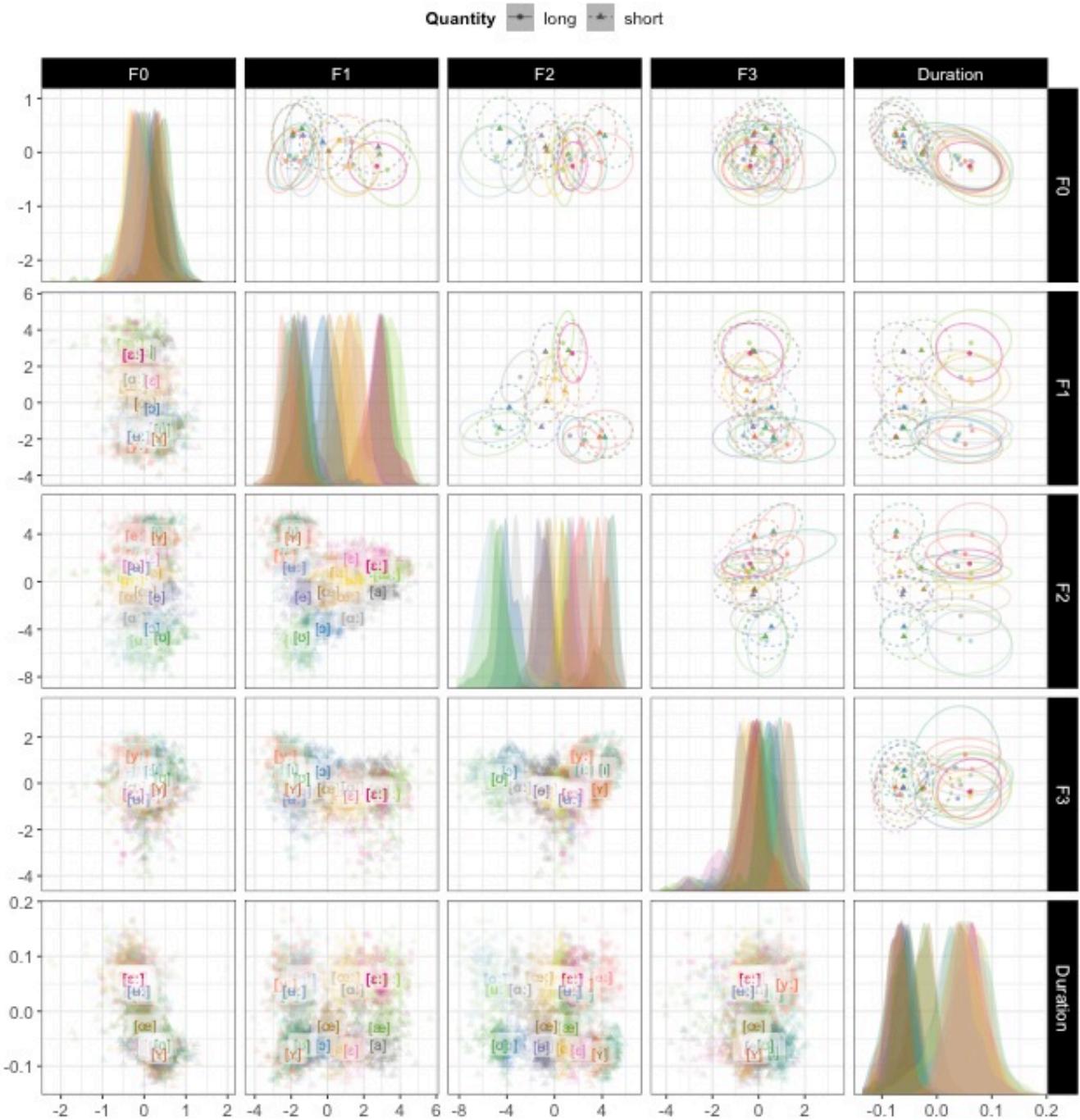


Figure S16. The SwehVd vowel data in C-CuRE ERB space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

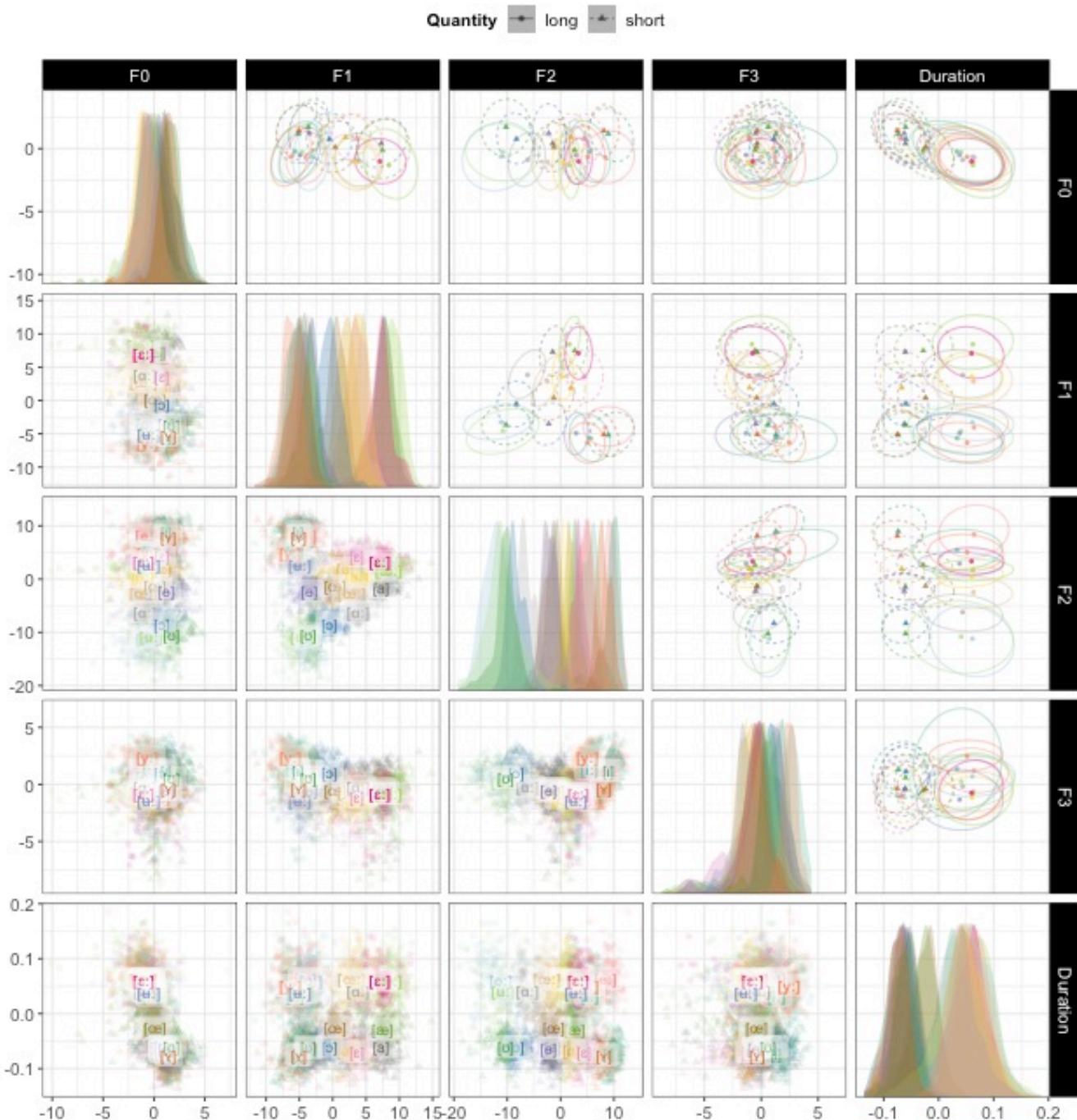


Figure S17. The SwehVd vowel data in C-CuRE semitones space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

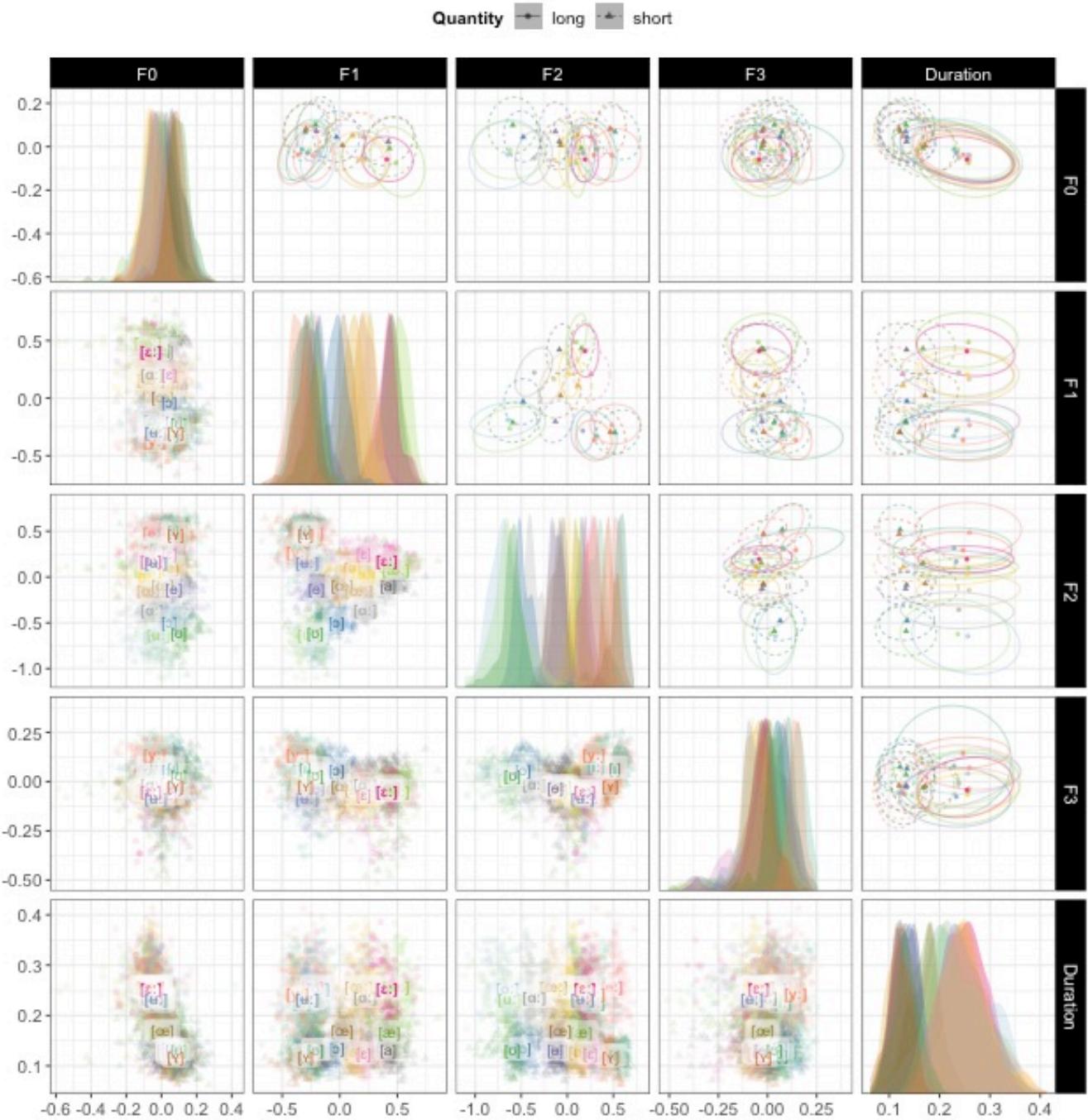


Figure S18. The SwehVd vowel data in Nearey1 (log) space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

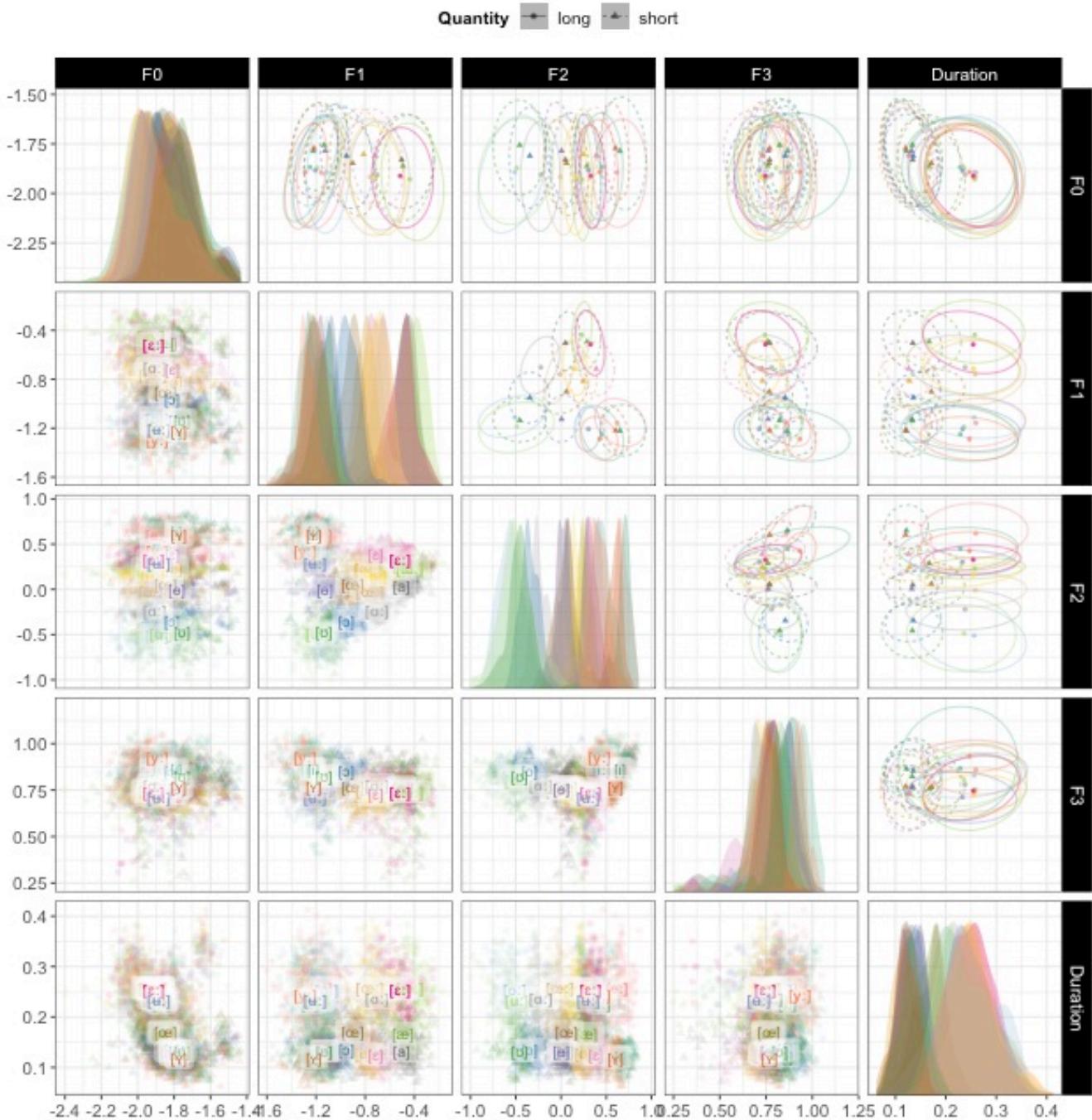


Figure S19. The SwehVd vowel data in Nearey2 (log) space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

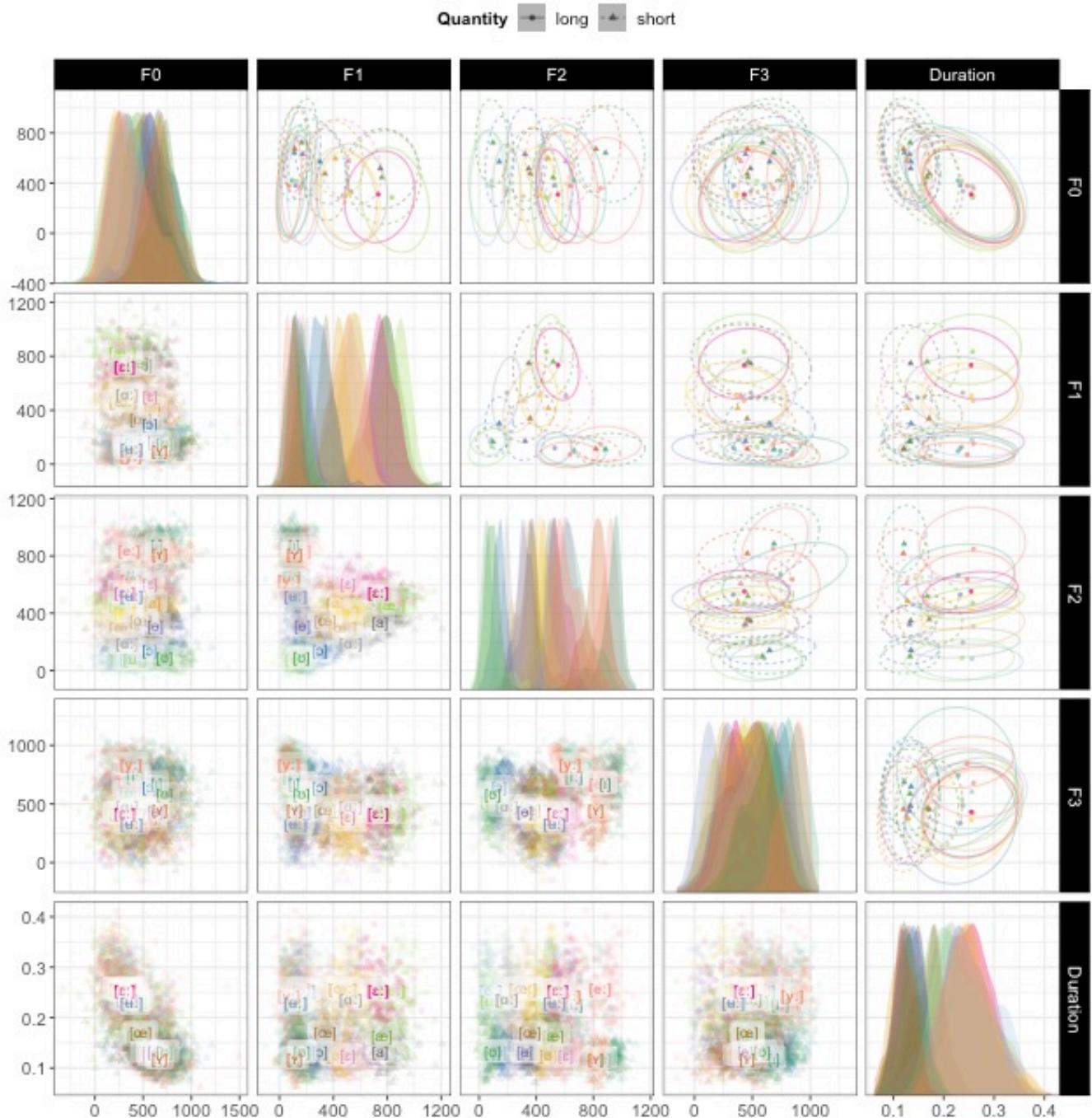


Figure S20. The SwehVd vowel data in Gerstman (Hz) space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

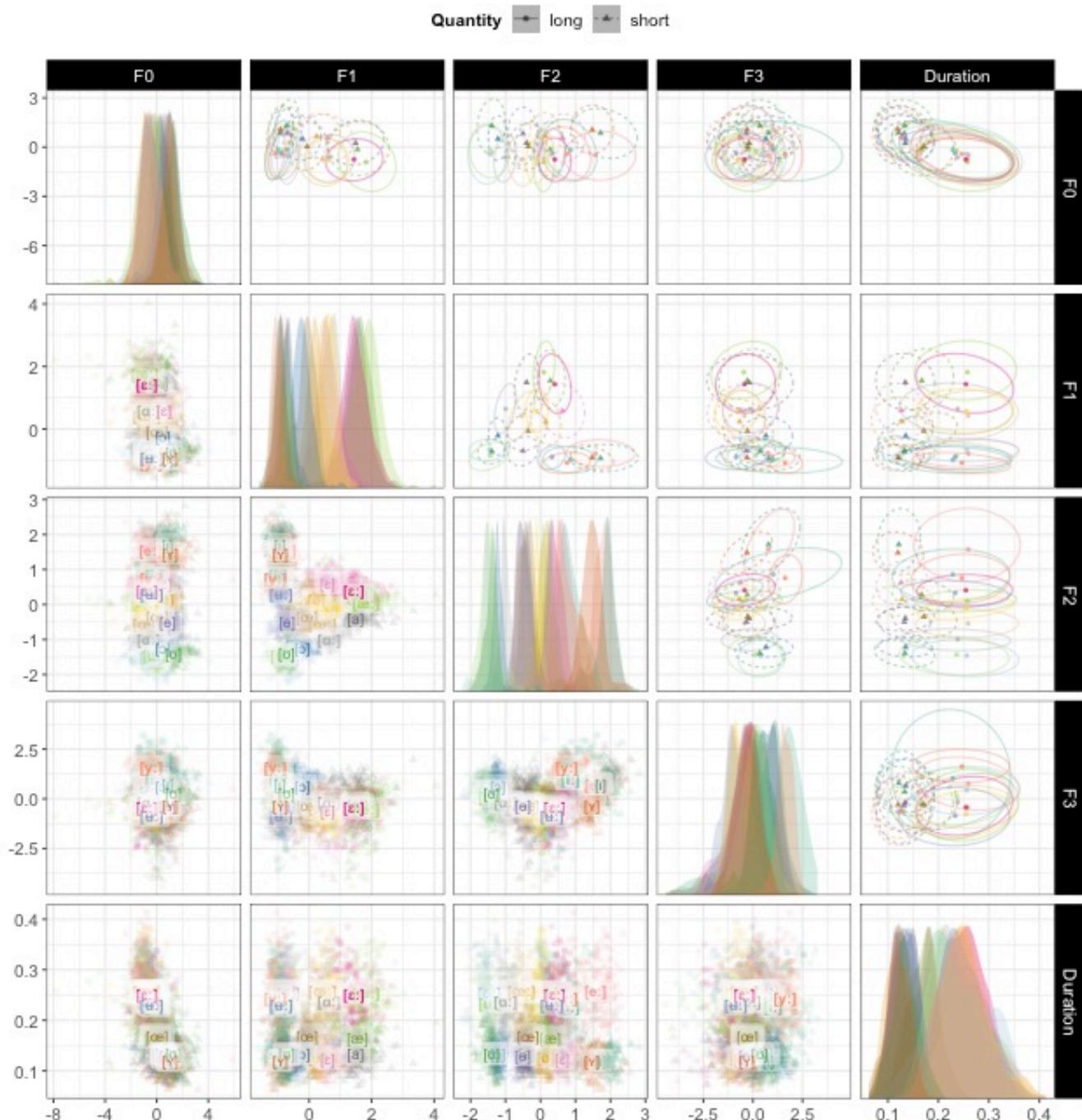


Figure S21. The SwehVd vowel data in Lobanov (Hz) space. Points show repetitions of each of the 21 Central Swedish vowels by 16 female native talkers in the database in F0-F3 and vowel duration cue space. Vowel labels indicate category means across talkers. Long vowels are boldfaced. Ellipses show bivariate Gaussian 95% confidence interval of category means.

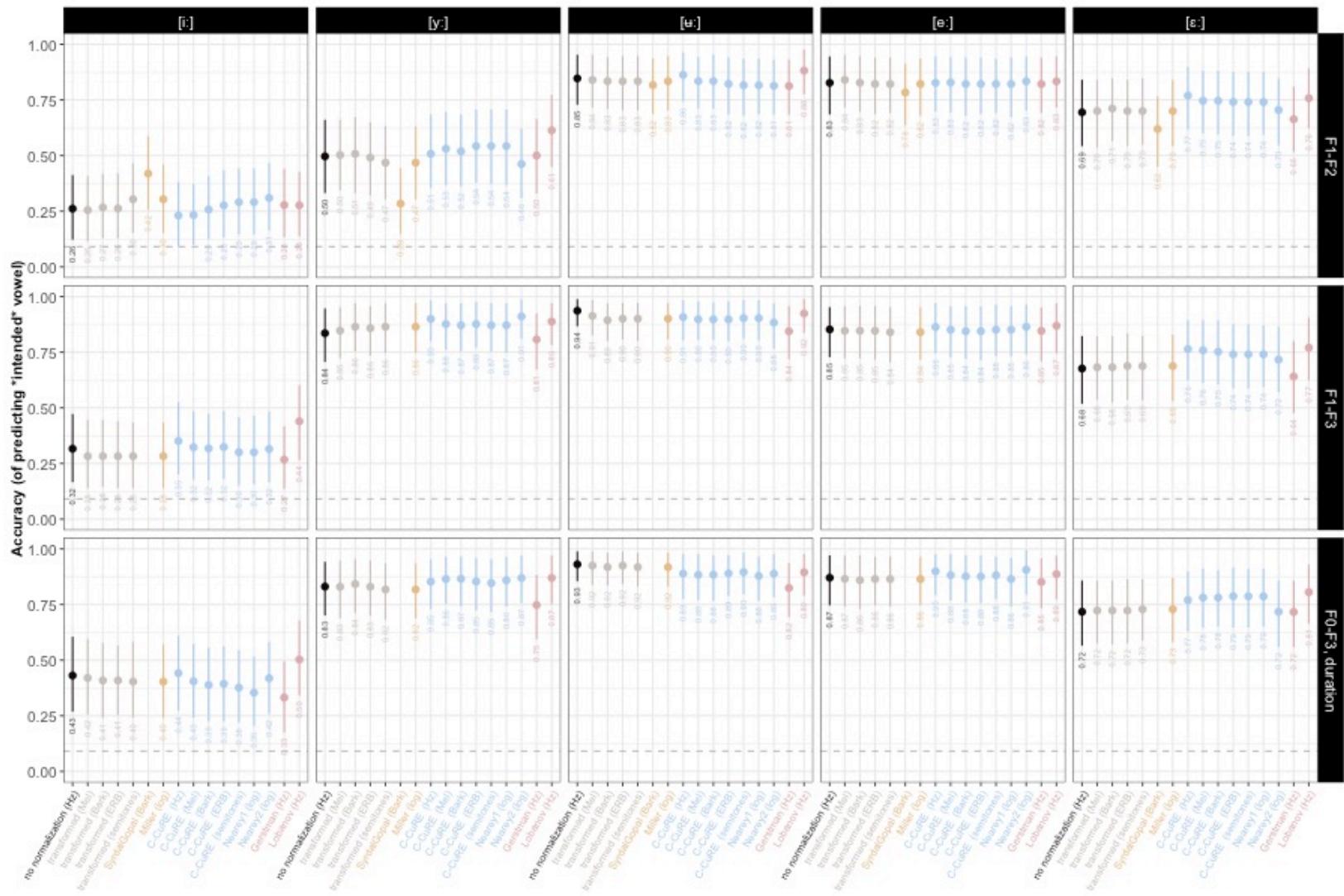


Figure S22. Per-vowel predicted categorization accuracy of the ideal observers trained on the **long** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

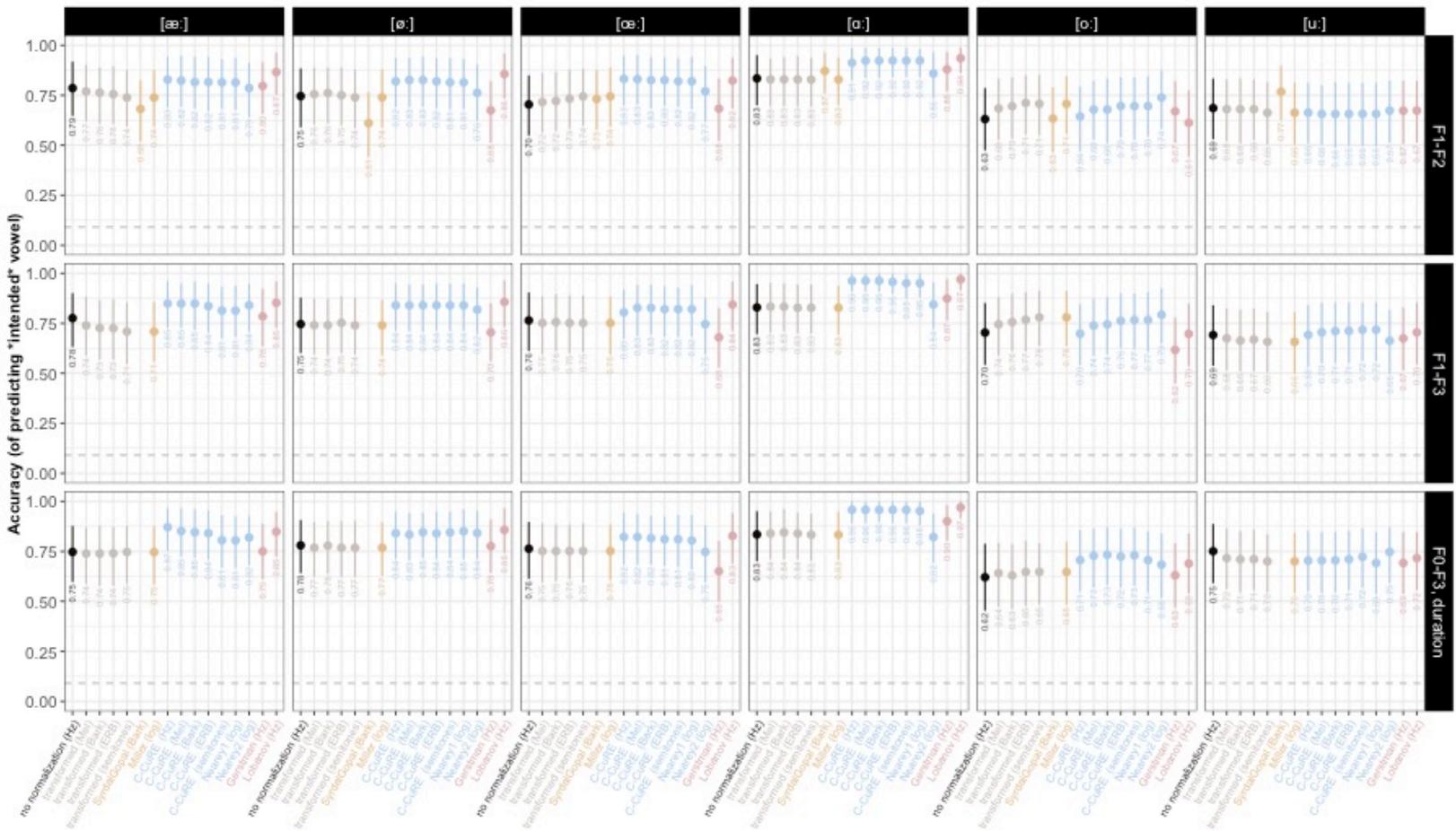


Figure S23. (Continued from last page) Per-vowel predicted categorization accuracy of the ideal observers trained on the **long** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

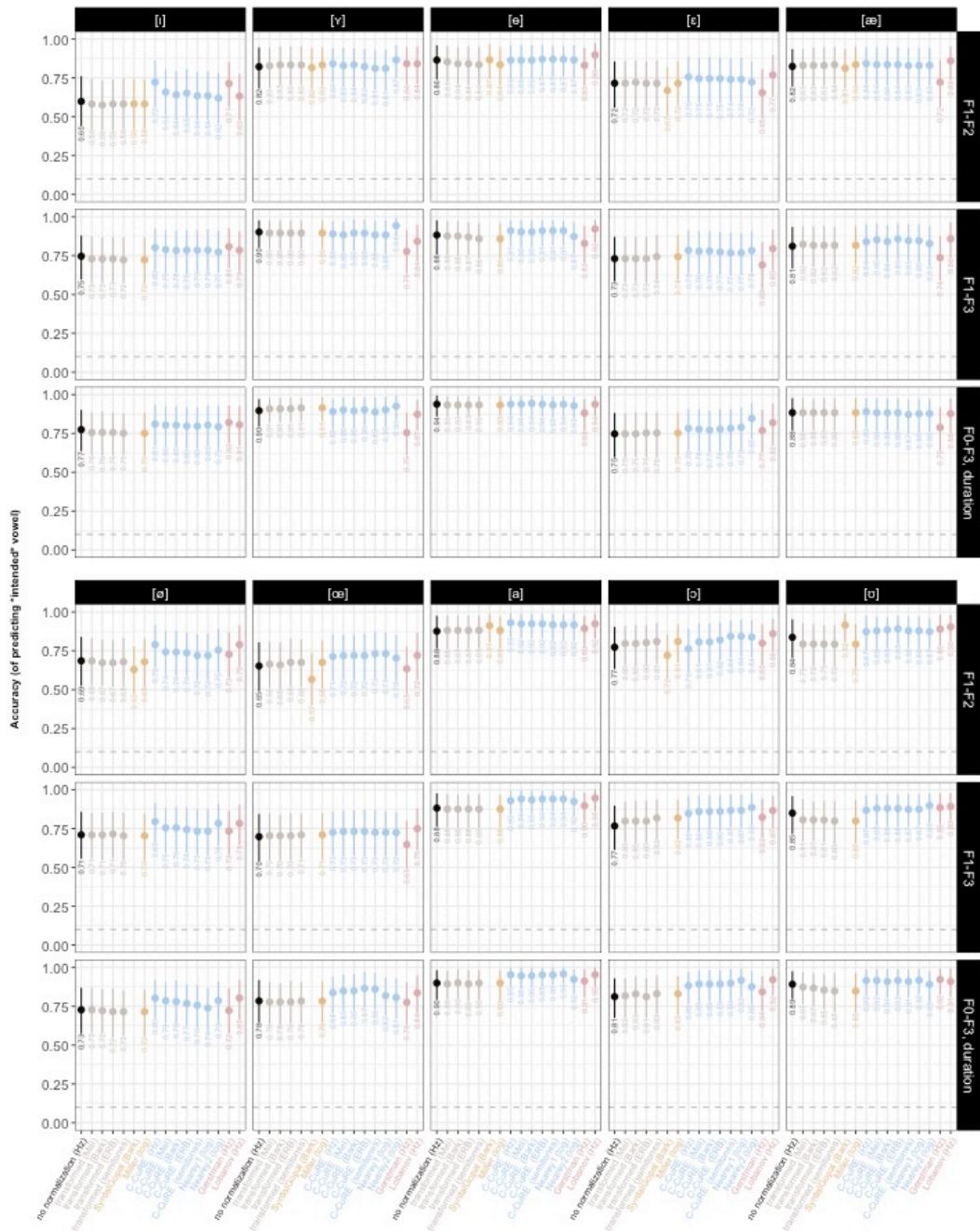


Figure S24. Per-vowel predicted categorization accuracy of the ideal observers trained on the **short** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

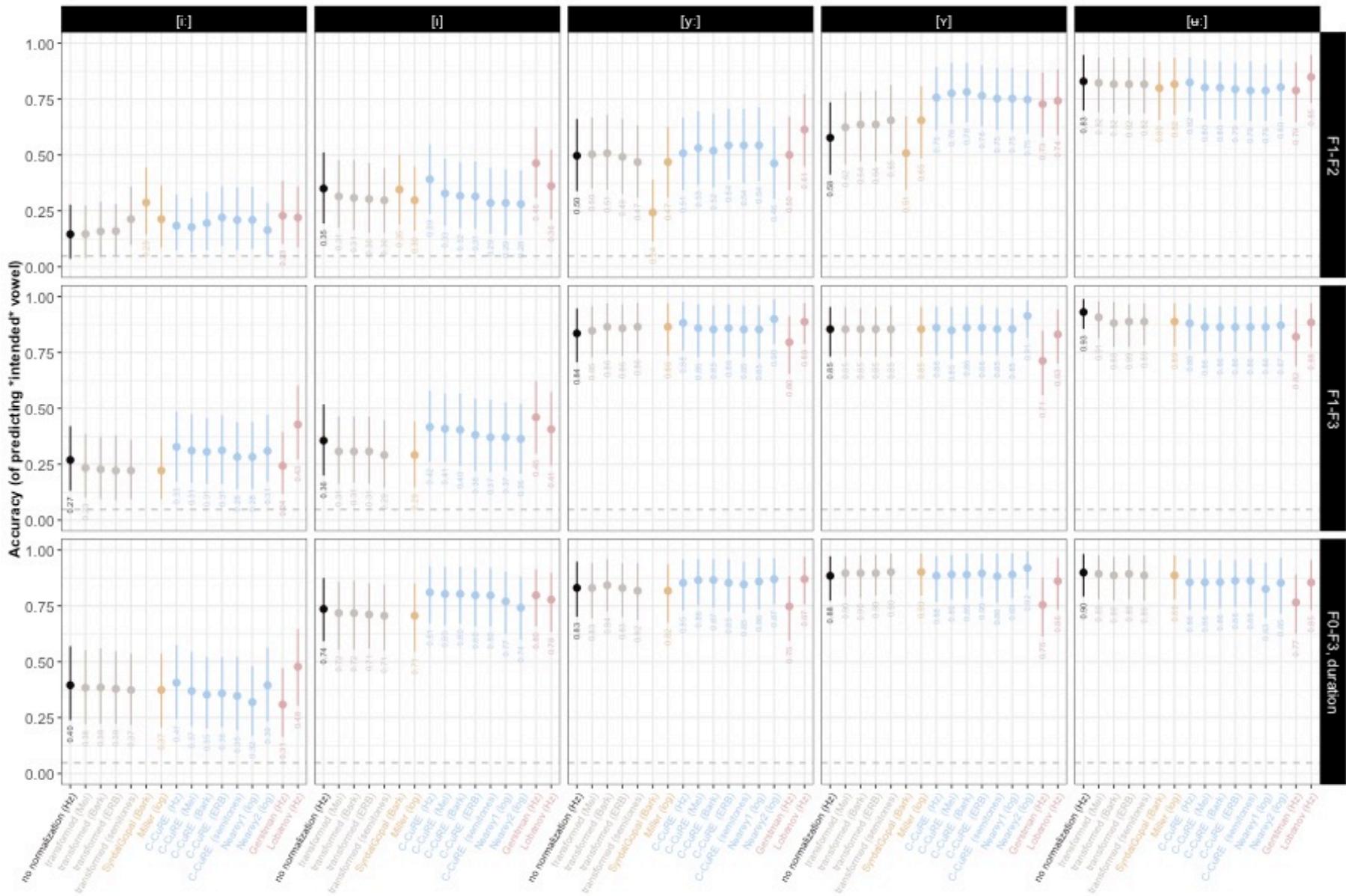


Figure S25. Per-vowel predicted categorization accuracy of the ideal observers trained on **all** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

Evaluating normalization accounts

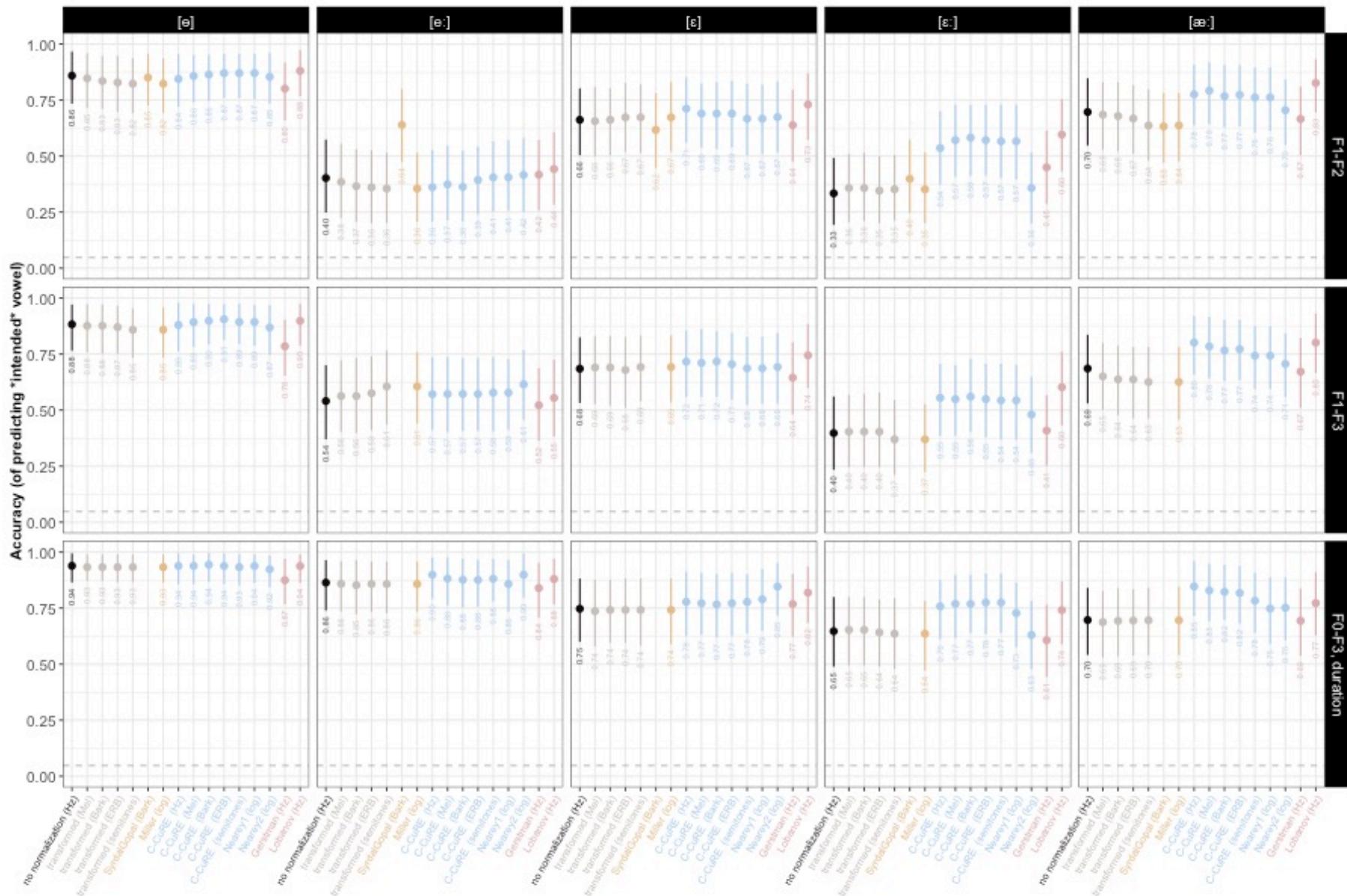


Figure S26. (Continued from last page) Per-vowel predicted categorization accuracy of the ideal observers trained on **all** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

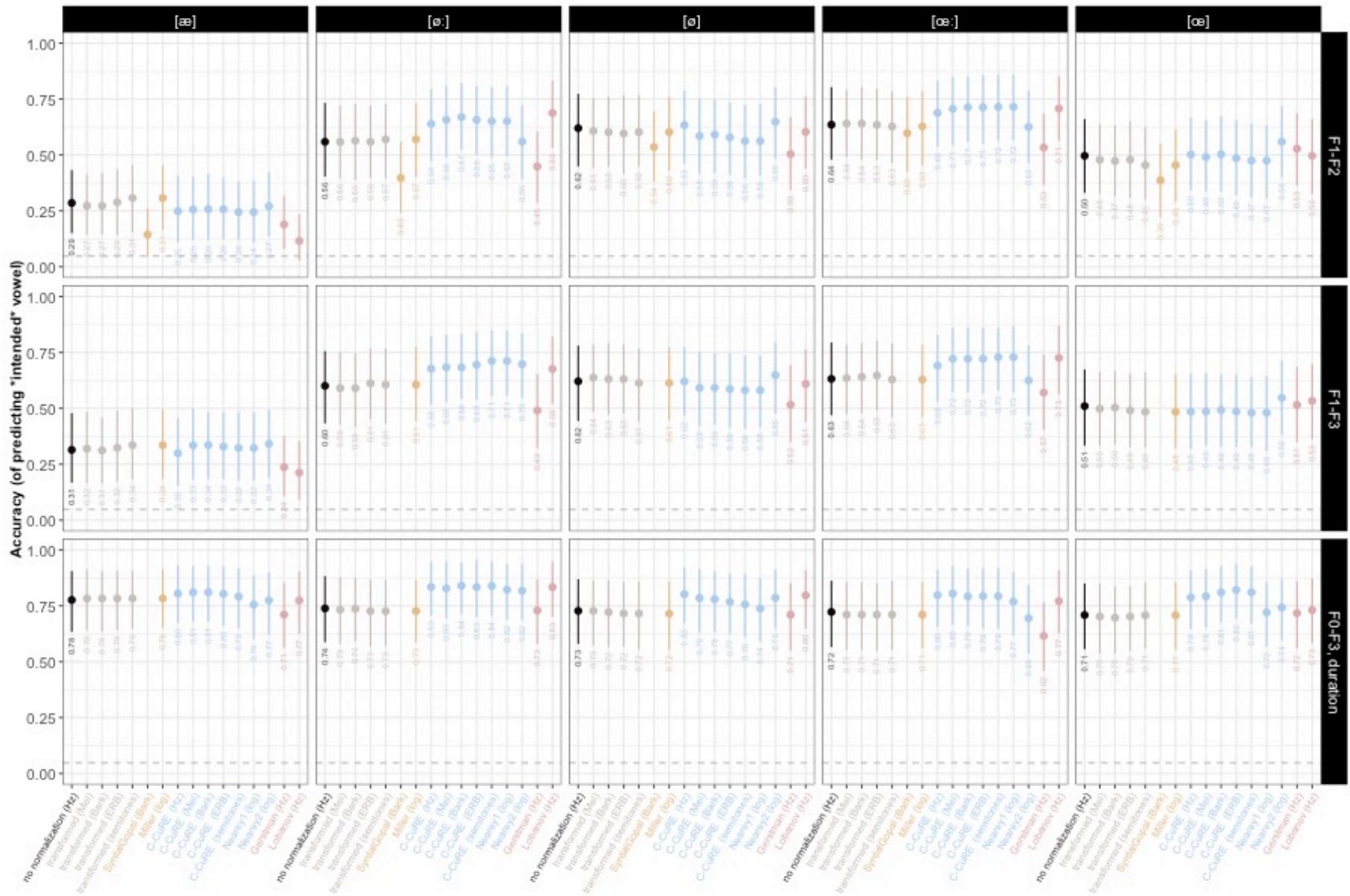


Figure S27. (Continued from last page) Per-vowel predicted categorization accuracy of the ideal observers trained on **all** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

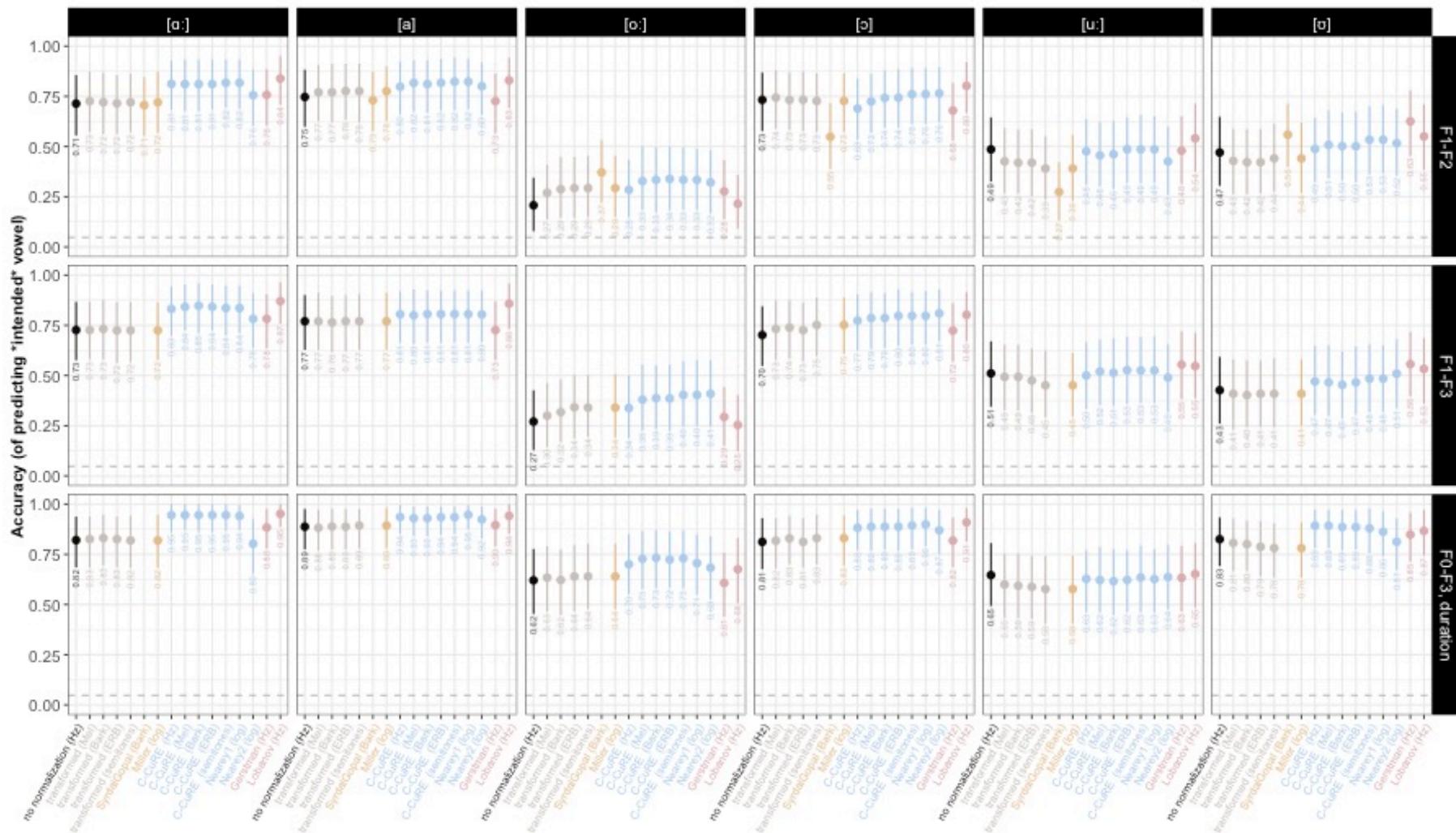


Figure S28. (Continued from last page) Per-vowel predicted categorization accuracy of the ideal observers trained on **all** vowels, under different assumptions about the relevant cues. Point ranges indicate the average mean accuracy and average 95% bootstrapped CI across the five folds. Chance level is indicated by grey line.

1213 6.2 Confusion and difference matrices of ideal observers

1214 To further explore effects of neighbouring categories, and which categories are more easily confused
1215 by the models and with what, we plot confusion matrices of the worst and best performing models
1216 trained on the long, short or all Central Swedish vowels, under the different assumptions about the
1217 relevant cues. Next to the confusion matrices, we plot difference matrices to facilitate comparison.

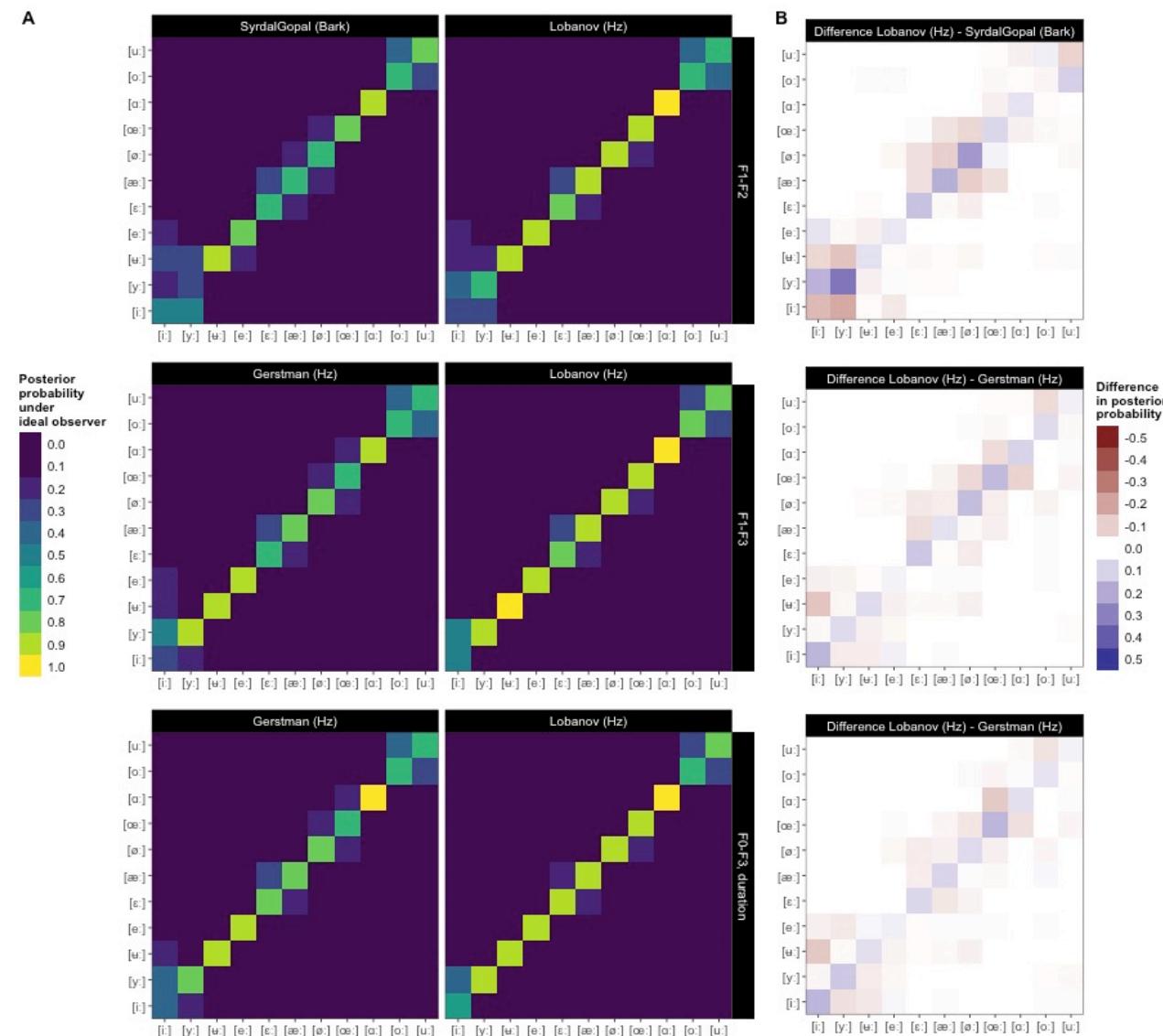


Figure S29. Illustration of the category-specific differences in predictions of the worst and best performing normalization models for each combination of cues (**rows**). The confusion matrices (**Panel A**) plot the predictions for the worst (**left**) and best (**right**) performing models in predicting the **long** vowels, under different assumptions about the relevant cues. Vowel intended by talker (x-axis) is plotted against vowel selected by ideal observer model (y-axis). Color fill indicates the posterior probability of the models predicting the intended vowel. The difference matrices (**Panel B**) illustrates the differences in predictions between the best and the worst performing models. Color fill indicates the difference in the posterior probability of the models predicting the intended vowel. More **purple** indicates an increase in posterior probability for the former over the latter model, more **red** indicates an advantage for the latter over the former.

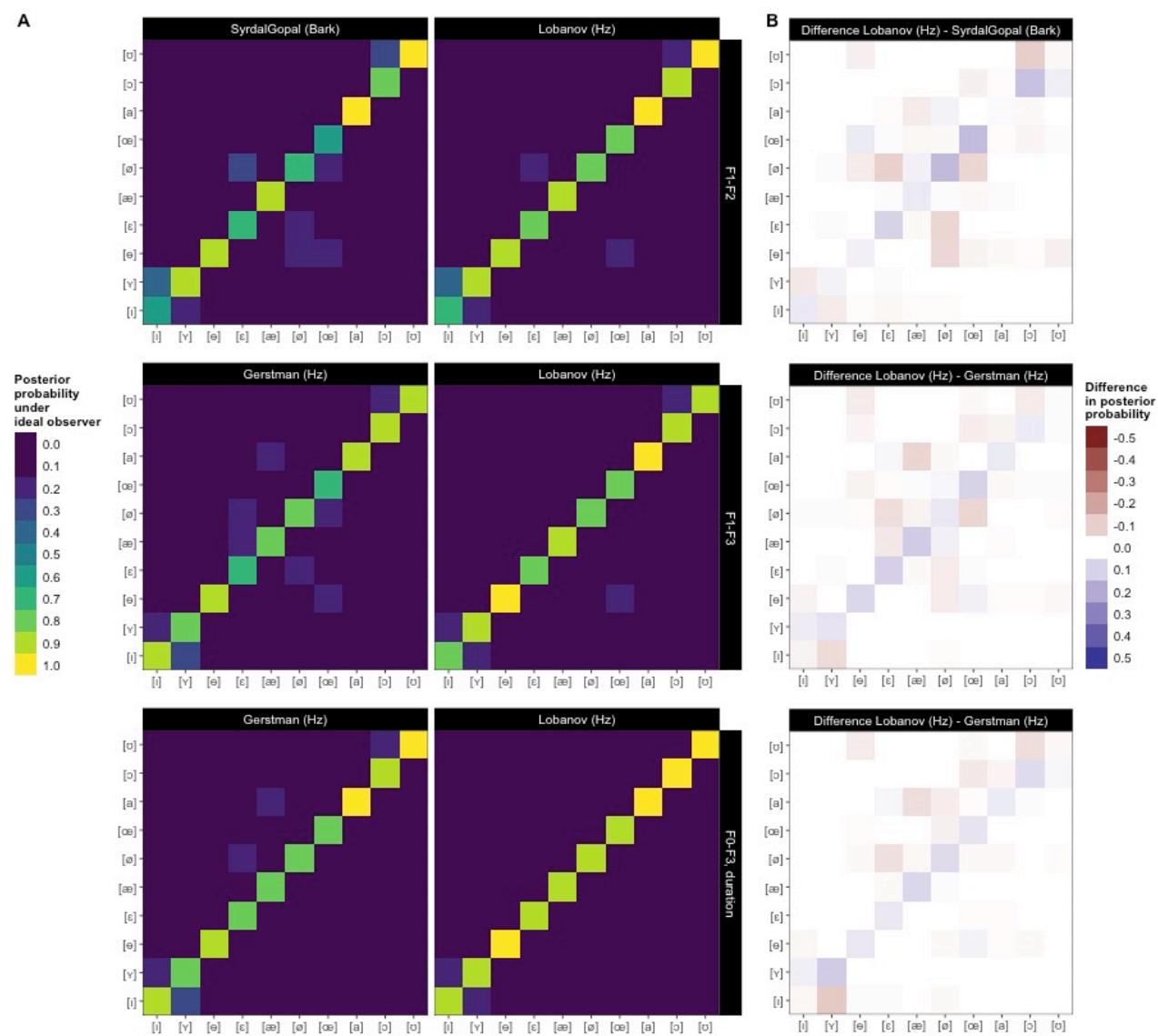


Figure S30. Illustration of the category-specific differences in predictions of the worst and best performing normalization models for each combination of cues (rows). The confusion matrices (Panel A) plot the predictions for the worst (left) and best (right) performing models in predicting the **short** vowels, under different assumptions about the relevant cues. Vowel intended by talker (x-axis) is plotted against vowel selected by ideal observer model (y-axis). Color fill indicates the posterior probability of the models predicting the intended vowel. The difference matrices (Panel B) illustrates the differences in predictions between the best and the worst performing models. Color fill indicates the difference in the posterior probability of the models predicting the intended vowel. More **purple** indicates an increase in posterior probability for the former over the latter model, more **red** indicates an advantage for the latter over the former.

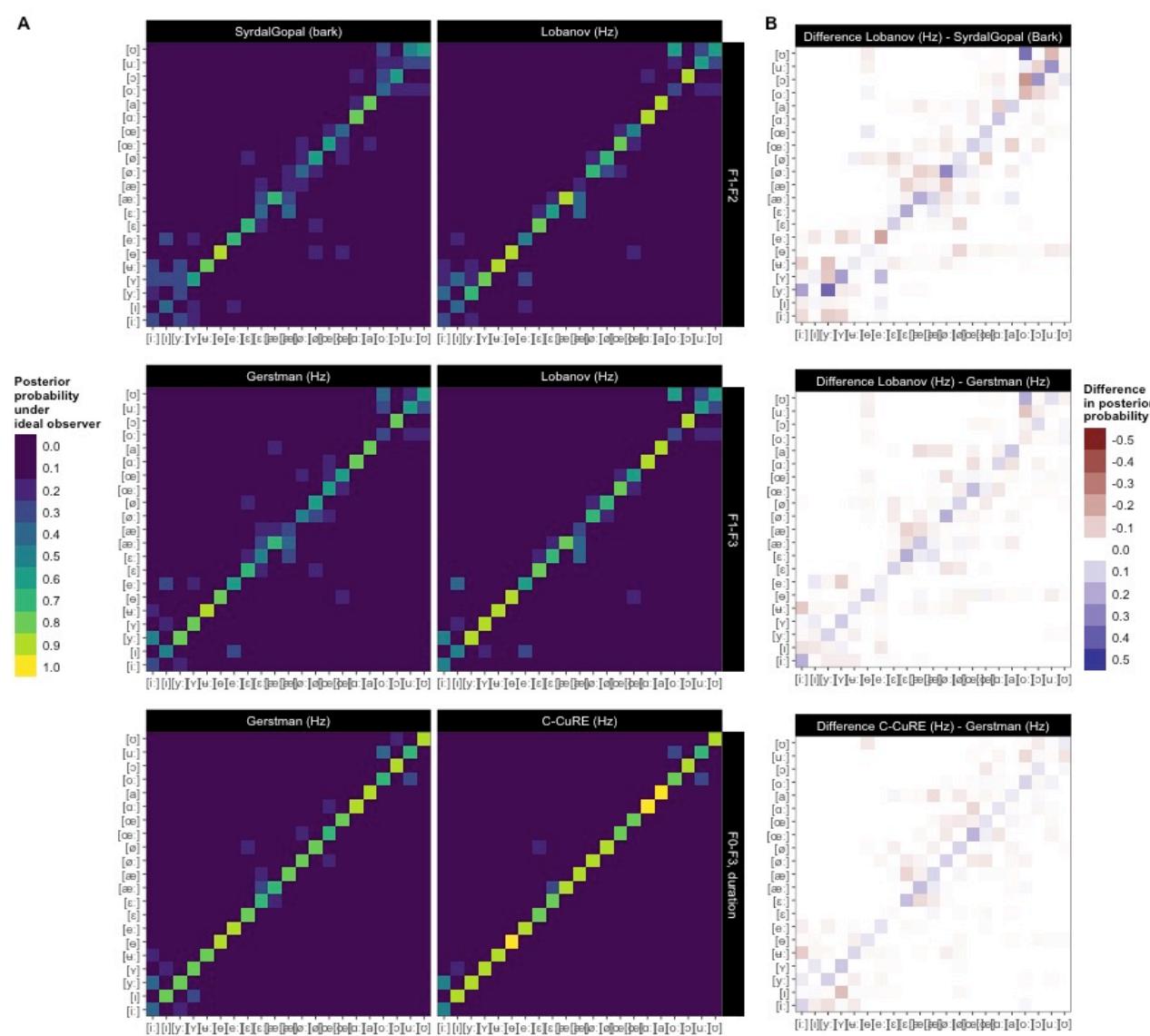


Figure S31. Illustration of the category-specific differences in predictions of the worst and best performing normalization models for each combination of cues (rows). The confusion matrices (Panel A) plot the predictions for the worst (left) and best (right) performing models in predicting the all vowels, under different assumptions about the relevant cues. Vowel intended by talker (x-axis) is plotted against vowel selected by ideal observer model (y-axis). Color fill indicates the posterior probability of the models predicting the intended vowel. The difference matrices (Panel B) illustrates the differences in predictions between the best and the worst performing models. Color fill indicates the difference in the posterior probability of the models predicting the intended vowel. More purple indicates an increase in posterior probability for the former over the latter model, more red indicates an advantage for the latter over the former.

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

1218 A large portion of previous studies evaluating normalization accounts against production data, has
 1219 compared approaches in terms of how they affect category variability. In this additional study,
 1220 we follow this traditional approach and evaluate how effectively different normalization accounts
 1221 reduce the within-category variability of Central Swedish vowels. We calculate a separability index
 1222 under different assumptions about the relevant cues and the size of the vowel space (the long and
 1223 short vowels separately, or the entire space) and assess the effects on vowel category variability. To
 1224 anticipate one take-home point, the results highlight important shortcomings of separability indices
 1225 in evaluating normalization accounts and underlines the benefits of using a perceptual model to
 1226 compare the effects of different normalization accounts.

1227 Before we evaluate how category separability is affected by normalization in F1-F2, F1-F3, and
 1228 F0-F3 and duration space, we look at how the normalization accounts affect the separability of
 1229 vowels along each cue separately (Figure S32). As we show below, this is helpful in understanding
 1230 the subsequently presented results for combinations of cues.

1231 7.1 Methods

1232 7.1.1 Speech materials

1233 This study employs the same speech materials as in the main study. Paralleling the main study, we
 1234 evaluated category separability for each combination of normalization account, cues, and training-
 1235 test fold. Specifically, we use the exact same cross-validation folds as in the main study.

1236 7.1.2 Separability index

1237 Previous studies have used different measures to assess the relative success of a normalization
 1238 procedure in reducing inter-talker variability (see Table 2 and Nearey, 1989 for an overview on
 1239 classification accounts). This includes assessing the reduction in variance or distance between
 1240 means by visual inspection (e.g., Clopper, 2009; Disner, 1980; Hindle, 1978), or by calculating
 1241 the reduction in within-category variance across talkers (e.g., Disner, 1980; Fabricius et al., 2009;
 1242 Flynn and Foulkes, 2011; Hindle, 1978), or comparing the degree of separation between category
 1243 means for unnormalized and normalized data, i.e., an F-ratio (e.g., Labov, 2010). We will assess how
 1244 distinguishable vowels become under different normalization accounts by calculating a separability
 1245 index, as described in Equation (S1). Following some previous studies (e.g., Labov, 2010), this
 1246 separability index is essentially an F statistics, where the F statistics is the ratio of the within- and
 1247 between-category variances:

$$\begin{aligned}
 \text{separability index} &= \frac{\text{between category } MS}{\text{within category } MS} \\
 &= \frac{\sum_{c=1,\dots,K} (N_c - 1)}{K - 1} \frac{\sum_{c=1,\dots,K} (\bar{x}_c - \bar{x})^2}{\sum_{c=1,\dots,K} \sum_{i=1,\dots,N_c} (x_{i,c} - \bar{x}_c)^2} \tag{S1}
 \end{aligned}$$

1248 where K is the number of categories, N_c is the number of observations for category c , $x_{i,c}$ is the
 1249 cue vector (for all cues considered in the calculation of the separability index) for observation i
 1250 of category c , \bar{x}_c is the cue mean vector for category c , and \bar{x} is the overall cue mean vector. We

1251 calculated this separability index separately for each combination of normalization account, cues,
1252 and training-test fold, as described next.

1253 7.2 Results

1254 For F1 (first row of Figure S32), we see a clear advantage for centering (in blue) and standardizing
1255 (in purple) compared to transformations (in grey) and intrinsic accounts (in yellow). In particular
1256 Lobanov normalization seems to maximize category separability along F1, at least for the long
1257 vowels and all vowels together. Notably, the accounts pattern differently along F2 (second row
1258 of Figure S32). Overall, differences between accounts are much smaller along F2, and the clear
1259 advantage of centering and standardizing accounts along F1 does not extend to F2.

1260 An altogether different picture is observed for F3. Compared to F1 and F2, the intrinsic account
1261 (Miller) performs substantially better in separating categories along F3, while all other accounts
1262 perform poorly. This result is surprising: one of the downsides of intrinsic approaches that has been
1263 noted in previous work is their sensitivity to measurement error (Thomas and Kendall, 2007). This
1264 sensitivity is caused by the fact that intrinsic accounts use a single measurement for normalization,
1265 rather than the less noisy estimates resulting from aggregating across segments that are used in
1266 extrinsic accounts. Since the third formant is often described as more difficult to reliably estimate
1267 than other formants (leading to more measurement error), F3 would be expected to be particularly
1268 affected by this weakness of intrinsic accounts.

1269 Yet, further visualization in Figure S33 confirms that F3 indeed separates categories particularly
1270 well when intrinsic normalization is applied. Compared to other accounts, Miller (1989) seems to be
1271 particularly successful in separating vowels that differ in lip rounding. For example, Miller (1989)
1272 separates two clusters among the high and mid-high vowels, one consisting of the back vowels [o:]
1273 and [u:], and the other one of the front [i:] and rounded [y:] and [ɯ:]. One possible explanation
1274 for this result is that intrinsic normalization is indeed particularly effective for F3, and that our
1275 correction of measurement errors—equally applied to all formants—effectively reduced the issue
1276 with F3 measurement errors (presumably the human brain, too, can do better than an uncorrected
1277 Praat algorithm without error correction). As we show below, this result for F3 carries over to any
1278 combination of cues that includes F3. It is, however, an artifact of using category separability to
1279 assess the effectiveness of normalization, as we show in the main study. We elaborate on this issue
1280 in the discussion further down.

1281 Returning to Figure S32, normalization does not increase category separability for F0. This is
1282 expected given that F0 is known to affect vowel separability primarily through its indirect influence
1283 on the interpretation of other formants (e.g., Barreda and Nearey, 2012; Barreda, 2020). Finally,
1284 for duration all of the C-CuRE accounts group together against the remaining accounts. This, too,
1285 is expected since all other accounts are formant-specific and thus do not normalize duration. In
1286 summary, the five cues contribute to category separability in different ways, and this is reflected
1287 in varying effectiveness of different normalization accounts. We also note that the best performing
1288 normalization account for any combination of cues and vowel qualities is typically never significantly
1289 better than the next best performing model (the 95% confidence intervals of the best model overlap
1290 with the mean of the next best model). In fact, for many combinations of cues and vowel qualities,
1291 many of the models perform similarly.

1292 Next, we summarize how normalization affects category separability when combinations of the
1293 five cues are considered. Figure S34 shows the separability index for the different normalization



Figure S32. Separability indices by normalization accounts for long vowels, short vowels, and all vowels together (columns), shown for each of the five cues considered in this study (rows). Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. Note that the ranges of the y-axes varies across plots.

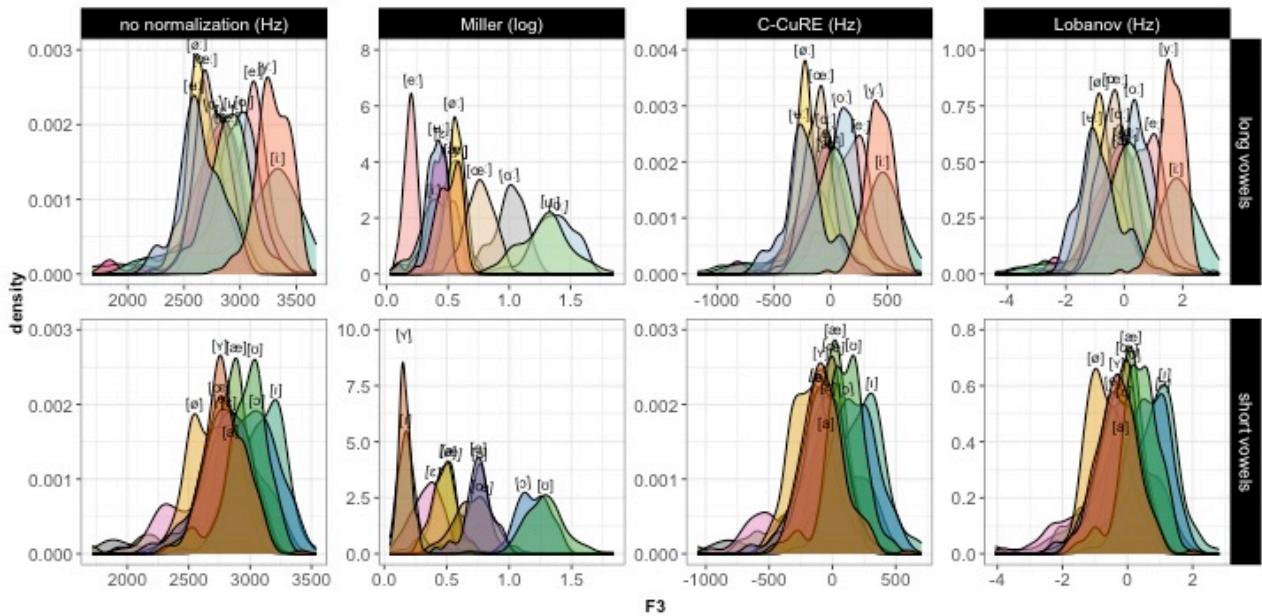


Figure S33. Category densities along F3 illustrates the effectiveness of vowel-intrinsic normalization for this cue. Here shown for Miller, compared to vowel-extrinsic accounts that center and/or standardize cues. For reference, densities in the absence of normalization are also shown.

1294 accounts for three different combinations of cues. For the first row of Figure S34, we followed
 1295 most previous research in assessing category separability for the combination of F1 and F2 (e.g.,
 1296 Disner, 1980; Fabricius et al., 2009; Flynn and Foulkes, 2011; Hindle, 1978; Labov, 2010).
 1297 Accounts that center against the talker's overall formant mean (in blue) are among the best
 1298 performing normalization accounts. No matter the assumed perceptual scale, centering always
 1299 improves category separability. Standardizing accounts (in purple), primarily Lobanov (1971), also
 1300 perform well at separating categories, more so for the long vowels. However, scale transformations
 1301 (in grey), and intrinsic accounts (in yellow), do not improve category separability compared to
 1302 unnormalized Hz, at least not when assessed on the long vowels or the entire vowel space.

1303 The remaining rows of Figure S34 compare normalization accounts when F3 (second row) or F0,
 1304 F3, and duration are included (third row). Overall, the category separability is now lower, a result
 1305 of how the accounts affect category separability along the cues added (see Figure S32). The most
 1306 drastic change in performance concerns the intrinsic Miller (1989) and the standardizing accounts.
 1307 When including F3, Miller (1989) performs as well or better, in absolute numbers, as when evaluated
 1308 on only the combination of F1 and F2, thereby increasing its performance relative to other accounts.
 1309 This increase in performance might be particularly pronounced for languages like Swedish, where F3
 1310 carries important information about lip rounding and thus vowel identity. In contrast, performance
 1311 of standardizing accounts drops substantially if F3 or any other cue besides F1 and F2 is included.^{S3}

^{S3} We confirmed this by conducting additional comparisons using only F1, F2 and F0, or only F1, F2 and duration. For both of these comparisons too, we found that standardizing accounts perform poorly.

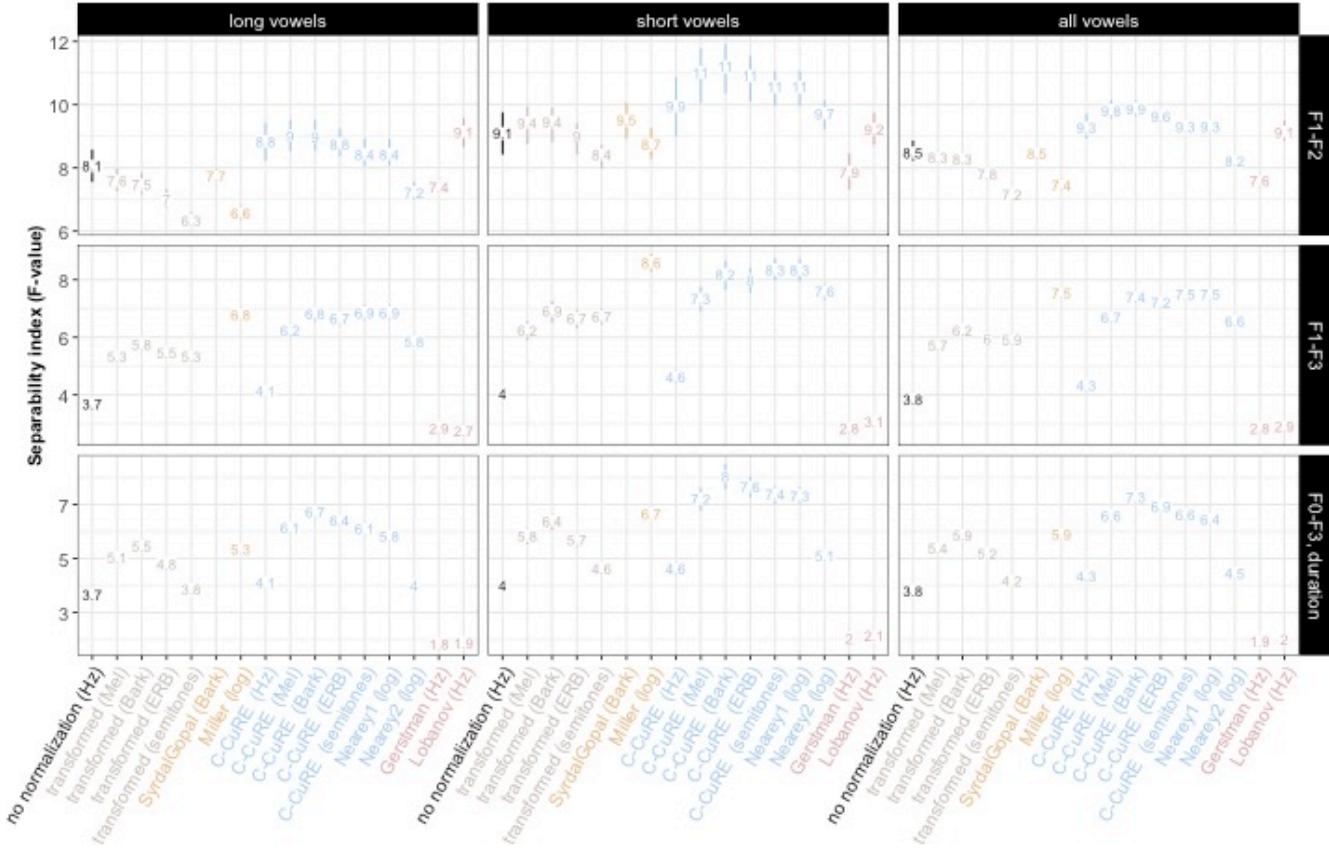


Figure S34. Separability indices by normalization accounts for long vowels, short vowels, and both long and short vowels together (columns) shown for three different combinations of cues (rows). Labels indicate mean across the five test folds. Intervals show average bootstrapped 95% confidence intervals across the test folds. Note that the ranges of the y-axes varies across plots.

1312 This mirrors what was found when assessing category separability separately for each cue (Figure
 1313 S32).

1314 Finally, looking across all three rows, category separability is consistently higher for short than
 1315 long vowels. The same pattern is evident for each cue separately in Figure S32. This result
 1316 conceptually replicates an initially surprising result of the main study: while short vowels are more
 1317 densely clustered in the center of the vowel space, and thus occupy a smaller perceptual space, they
 1318 also exhibit less category variability and less category overlap, making them overall more separable.
 1319

7.3 Discussion

1320 When only F1 and F2 are considered, as in most previous work on vowel normalization, we find
 1321 that extrinsic centering and standardizing accounts achieve the best category separability. Within
 1322 these two types of accounts, there is considerable variability. For example, among the intrinsic
 1323 accounts, Miller performs worse than Syrdal & Gopal, among the extrinsic accounts, versions of
 1324 C-CuRE seem to consistently perform best. It is also worth noting, however, that there is never a
 1325 single account that performs significantly better than all other normalizations. This points to the

1326 inherent similarities across normalization accounts, and perhaps limitations of the approach taken
1327 here (and in some previous work). This point is also raised in the general discussion in the main
1328 paper. Regardless of these caveats, the findings for F1 and F2 in this additional study, revise the
1329 results of Disner (1980) for Swedish, and instead replicates previous findings for the other Germanic
1330 languages in Disner's sample as well as the majority of previous studies on other languages (e.g.,
1331 Fabricius et al., 2009; Flynn and Foulkes, 2011; Labov, 2010).

1332 However, when F3 is considered along with F1 and F2, this result does no longer hold. Key
1333 to understanding this result and what it says about the suitability of category separability as a
1334 measure of normalization accounts is Figure S32: while extrinsic normalization performs better
1335 than other approaches for F1 and F2, the absolute differences in performance are small compared
1336 to the advantage of the intrinsic account observed for F3. Combined with a seemingly innocuous
1337 aspect of the separability index in Equation (S1), this allows separability along F3 to dominate
1338 separability along the other cue dimensions. Our separability index takes the *sum* of (squared)
1339 distances along each cue dimension, essentially assuming that the effect of all cues is simply a sum
1340 of each cue's effect considered separately. This means that the separability index cannot capture
1341 the *joint* effect of cues—whether, for example, one cue effectively separates one set of categories
1342 and another cue separates another set of categories, rather than both cues separating the same
1343 categories. The separability index thus cannot recognize, for example, that F1 and F2 capture
1344 largely complementary aspects of the vowel inventory (as evident in, for example, Figures S5 and
1345 S6).

1346 This is not the only deficiency of the separability index or similar measures of category variability.
1347 The use of *squared* distances means that even a small number of observations located far away from
1348 the category mean can disproportionately affect the index. Consider, for example, the F3 densities
1349 in Figure S33. For non-intrinsic normalizations, some categories have low but non-zero densities far
1350 away from the mode. Because of the use of squared distances, this results in low category separability
1351 for these normalization accounts despite the fact that observations with such cue values are rare
1352 and thus not expected to have a large effect on the *average* perceptual separability of vowels. For
1353 the same reason (the use of squared distances), category separability can be high even if a cue
1354 separates only a small subset of categories (as is the case for F3), compared to another cue that
1355 more gradually separates *all* categories (as is the case for F1 and F2; see Figure 4).

1356 In sum, indices of variability and category separability like that in Equation (S1) fail to adequately
1357 assess the expected consequences of normalization for perception, which is the primary interest of
1358 this paper, and addressed by the methodology we employed in the main study.