

The effect of pre-linguistic normalization in vowel perception

Anna Persson<sup>1</sup> & T. Florian Jaeger<sup>2,3</sup>

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

<sup>2</sup> Brain and Cognitive Sciences, University of Rochester

<sup>3</sup> Computer Science, University of Rochester

Author Note

We are grateful to ### omitted for review ###

Correspondence concerning this article should be addressed to Anna Persson, Department of Swedish Language and Multilingualism, Stockholm University, SE-106 91 Stockholm, Sweden.

E-mail: anna.persson@su.se

11 Abstract

12 XXX. All data and code for this study are shared via OSF, including the R markdown document  
13 that this article is generated from, and an R library that implements the models we present.

14 *Keywords:* speech perception; vowels; distributional learning; computational model

15 Word count: XXX

The effect of pre-linguistic normalization in vowel perception

## 1 To Do

- check that treatment coding for glm is correct (7 variables instead of 3)

## 2 Abstract

One of the central challenges for speech perception is that talkers differ in pronunciation—i.e., how they map linguistic categories and meanings onto the acoustic signal (Lieberman et al., 1967). While this challenge is always present, it is most evident when listeners first encounter talkers with unfamiliar pronunciations. Yet, listeners typically overcome even these difficulties within minutes (e.g., Clarke & Garrett, 2004; Xie et al., 2017). What mechanisms underlie these adaptive abilities remains unclear. One highly influential general hypothesis holds that inter-talker differences are removed via low-level pre-linguistic auditory normalization of acoustic cues. There are now at least a dozen of competing normalization proposals (e.g., Lobanov, 1971; Nearey, 1978; McMurray & Jongman, 2011). Despite the fact that these proposals differ in the cognitive capacities they entail, comparisons of different normalization models against human perception remain largely lacking (but see Bion & Escudero, 2007; McMurray & Jongman, 2011). Here, we seek to address this gap by comparing normalization accounts against the perception of American English vowels.

We first trained 7 ideal observer (IO) models that differed in whether they were trained on unnormalized or normalized acoustic cues. The following influential normalization models was employed: C-CuRE, Lobanov, Miller, Gerstman, Nearey1 and Nearey2. All models were trained against the same phonetic database of productions of all 8 h-VOWEL-d words of American English (heed, hid, head, had, odd, hut, hood, who'd, N=9 tokens per vowel from 17 talkers each), previously described in (REF to Li & Xie). We then compared the predictions of all IOs against L1 American English listeners' 8-way categorization responses for productions of the 8 h-VOWEL-d words in a web-based experiment (N=22 participants).

The best performing IO only centered cues relative to talkers' cue mean (as in C-CuRE normalization, mean accuracy 64.5%, SE 1.0%). This is significantly above chance (12.5%) but

also significantly below the best possible performance (always guessing the response most frequently given by human listeners, 72.0%). Statistically indistinguishable performance was achieved when cues were both centered and scaled (as in, e.g., Lobanov normalization, 63.5%, SE 1.0%). This suggests that simple normalization operations might be sufficient to explain perception. Either type of normalization model performed significantly better than IOs based on unnormalized cues (53.2%, SE 0.7%). Additional comparisons showed that the benefit of normalization over unnormalized cues held regardless of whether cues were first transformed from acoustic (Hz) into perceptual spaces (Mel, Bark, and similar). In conclusion, these results indicate that pre-linguistic normalization (or computationally similar algorithms) contribute to the remarkable adaptive abilities of human speech perception. Our results further suggest that human perception only employs simple normalization operations—such as centering cues relative to a talker’s mean. Finally, we find that simple ideal observers achieve performance far above chance in predicting perception, although our results also indicate that human perception might employ additional adaptive algorithms (given the 7.5% accuracy difference between the best performing models and human performance).

### 3 Introduction

One of the central challenges for speech perception is that talkers differ in pronunciation—i.e., how they map linguistic categories and meanings onto the acoustic signal (Liberman et al., 1967). While this challenge is always present, it is most evident when listeners first encounter talkers with unfamiliar pronunciations. Yet, listeners typically overcome even these difficulties within minutes (e.g., Clarke & Garrett, 2004; Xie et al., 2017). What mechanisms underlie these adaptive abilities remains unclear. One highly influential general hypothesis holds that inter-talker differences are removed via low-level pre-linguistic auditory normalization of acoustic cues. There are now at least a dozen of competing normalization proposals (e.g., Lobanov, 1971; Nearey, 1978; McMurray & Jongman, 2011), that have indeed been shown to reduce irrelevant inter-talkers variability due to e.g. anatomical or physiological factors. Despite the fact that these proposals differ in the cognitive capacities they entail, comparisons of different normalization models against human perception remain largely lacking (but see Bion & Escudero, 2007;

McMurray & Jongman, 2011). Here, we seek to address this gap by comparing normalization accounts against the perception of American English vowels.

Normalization procedures that reduce the category variance and increase category separability should generally reduce perceptual difficulties, hence improve speech perception/categorization. Reducing category overlap is known to increase category separability in perception (Feldman, Griffiths, and Morgan (2009); Kleinschmidt and Jaeger (2015); Kronrod, Coppess, and Feldman (2016)). Following previous research, we hypothesize that normalization procedures that have previously been shown to reduce inter-talker variability, would better explain human perception of vowels... The following influential normalization models was employed: C-CuRE, Lobanov, Miller, Gerstman, Nearey1 and Nearey2.

## 4 Methods

In order to assess the effect of normalization procedures, we use a computational model based on Bayesian probability theory, ideal observers (see e.g. Feldman et al., 2009; Kleinschmidt & Jaeger, 2015; Kronrod et al., 2016). This allows us to simulate the effects of normalization procedures and to ask what the predicted consequences are for perception. We trained 7 ideal observer (IO) models on unnormalized or normalized acoustic cues from a phonetic database of American English, previously described in (REF to Li & Xie). The Li & Xie corpus consists of recordings from 15 native (five female) and 15 non-native (five female) talkers of a Northeastern dialect of American English (ages 18 to 35 years old). For this study, we selected the male and female native talkers (N=17). The IOs were trained on all 8 h-VOWEL-d words of English in the database (heed, hid, head, had, odd, hut, hood, who'd, N=9 tokens per vowel from 17 talkers each). The IOs' categorization performances were compared against each other, as a measurement of the efficiency of each model in reducing irrelevant inter-talker variability in the data. The predictions of the IOs were then evaluated against human categorization data from a web-based experiment on vowel categorization of English vowels.

## 4.1 Vowel categorization experiment

In order to evaluate the IOs’ predictions against human perceptual data, we exposed native English listeners (N=22) to the all vowel productions from one of the female native English talkers in the phonetic database (9 repetitions \* 8 vowels = 72 words) in an 8-way categorization experiment. The experiment was administrated on Amazon Mechanical Turk and consisted of one test block with 144 trials (each word repeated twice over the experiment). Participants were instructed to listen to a female talker saying words, and click on a word on screen to report what word they heard. At each trial, all eight hVd-words were displayed on screen in a response grid. Eligibility requirements, besides being a native speaker of American English, were to complete the experiment in a quiet room, wearing over-the-ear headphones of good sound quality. Before taking the experiment, participants performed a sound check and signed a consent form. After completing the experiment, participants filled out a language background questionnaire and were reimbursed.

## 4.2 Normalization procedures

The normalization procedures used in this study are listed in Table 1. We selected procedures that are commonly used in speech production and perception research and that have been compared and evaluated in several studies (see e.g. (adankComparisonVowelNormalization2004a?), XXX). All procedures take formant frequencies (in Hertz) for the first two formants, F1 and F2, as input and outputs normalized versions of the same formants.

Table 1

*Normalization procedures selected*

Normalization procedure	Source
Miller (formant-ratio)	(Miller, 1989)
Nearey1 (logmean)	(Nearey, 1977)
Nearey2 (shared logmean)	(Nearey, 1977)
C-CuRE (centering; subtracting from expected)	(McMurray et al., 2011)
Lobanov (z-score)	(Lobanov, 1971)
Gerstman (range normalization)	(Gerstman, 1968)

#### 4.2.1 Training-test data split

Prior to modelling the effects of different normalization procedures on the perception of English vowels, we split the vowel data into training and testing portions in a K-fold cross-validation approach in order to minimize the risk of overfitting the models. Overfitting would in this case mean that the ideal observers learn overly specific vowel categorization functions that perform well on the vowel data they have been trained on, but cannot generalize to unseen data. We set  $K = 5$ , and split the vowel data randomly into five equally sized bins. We then trained five ideal observers on four of the bins in a latin square design, and tested each ideal observers' predictions for vowel perception on the test data from the vowel categorization experiment.

```
Warning: Using an external vector in selections was deprecated in tidyselect 1.1.0.
i Please use `all_of()` or `any_of()` instead.

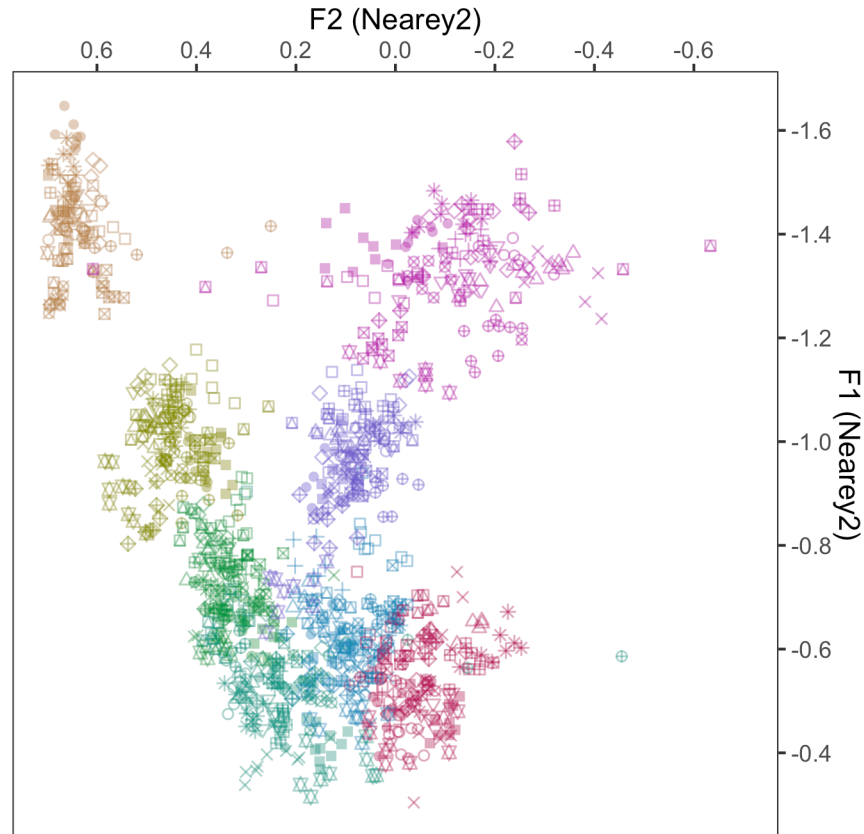
# Was:
data %>% select(cues)

# Now:
data %>% select(all_of(cues))

See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
```

The test data from the perceptual experiment is then subsequently normalized with the same normalization procedures, but using the parameters generated from running the procedures on the training data. This means that we use the means and variances obtained from the training data to normalize the test data for all normalization procedures that make use of formant means and variances. This applies to Lobanov normalization, Nearey1, Nearey2, Miller, C-CuRE and Gerstman????? (min,max). MOTIVATION ————— simulating that learning has finished when exposure is over???

```
Warning: Removed 36 rows containing missing values (`geom_point()`).
```



## 5 Results

The best performing IO only centered cues relative to talkers' cue mean (as in C-CuRE normalization, mean accuracy 64.5%, SE 1.0%). This is significantly above chance (12.5%) but also significantly below the best possible performance (always guessing the response most frequently given by human listeners, 72.0%). Statistically indistinguishable performance was achieved when cues were both centered and scaled (as in, e.g., Lobanov normalization, 63.5%, SE 1.0%). This suggests that simple normalization operations might be sufficient to explain perception. Either type of normalization model performed significantly better than IOs based on unnormalized cues (53.2%, SE 0.7%). Additional comparisons showed that the benefit of normalization over unnormalized cues held regardless of whether cues were first transformed from acoustic (Hz) into perceptual spaces (Mel, Bark, and similar).



## 5.1 Predicting the category heard by listeners

We tested the IOs in predicting the category heard by listeners in the categorization experiment.

```
Warning: Using an external vector in selections was deprecated in tidyselect 1.1.0.
```

```
i Please use `all_of()` or `any_of()` instead.
```

```
# Was:
```

```
data %>% select(levels.vowel.IPA)
```

```
# Now:
```

```
data %>% select(all_of(levels.vowel.IPA))
```

```
See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
```

Figure ?? visualizes the accuracy of the seven IOs in predicting the vowel responses from the 22 native English participants in the vowel categorization experiment.

```
Warning: `cols` is now required when using `unnest()`.
```

```
i Please use `cols = c(ci)`.
```

```
`cols` is now required when using `unnest()`.
```

```
i Please use `cols = c(ci)`.
```

```
$predicates
```

```
<list_of<quosure>>
```

```
[[1]]
```

```
<quosure>
```

```
expr: ^IO.cue_normalization == "C-CuRE (Hz)"
```

```
env: global
```

```
178 $n
179 NULL
180
181 $max_highlight
182 [1] 5
183
184 $unhighlighted_params
185 list()
186
187 $use_group_by
188 NULL
189
190 $use_direct_label
191 NULL
192
193 $line_label_type
194 [1] "ggrepel_label"
195
196 $label_key
197 <quosure>
198 expr: ^NULL
199 env: empty
200
201 $label_params
202 $label_params$fill
203 [1] "white"
204
205
206 $keep_scales
207 [1] FALSE
```

208

209 `$calculate_per_facet`210 `[1] FALSE`

211

212 `attr("class")`213 `[1] "gg_highlighter"`

214       Response accuracy if human responses are to predict human responses. For this, we are  
 215 assuming that someone knows the distribution of human responses for all items and uses the  
 216 criterion choice rule to predict what the response would be. This gives us an upper limit—no  
 217 model could ever predict human responses more accurately.

218 `$predicates`219 `<list_of<quosure>>`

220

221 `[[1]]`222 `<quosure>`223 `expr: ^IO.cue_normalization == "C-CuRE (Hz)"`224 `env: global`

225

226

227 `$n`228 `NULL`

229

230 `$max_highlight`231 `[1] 5`

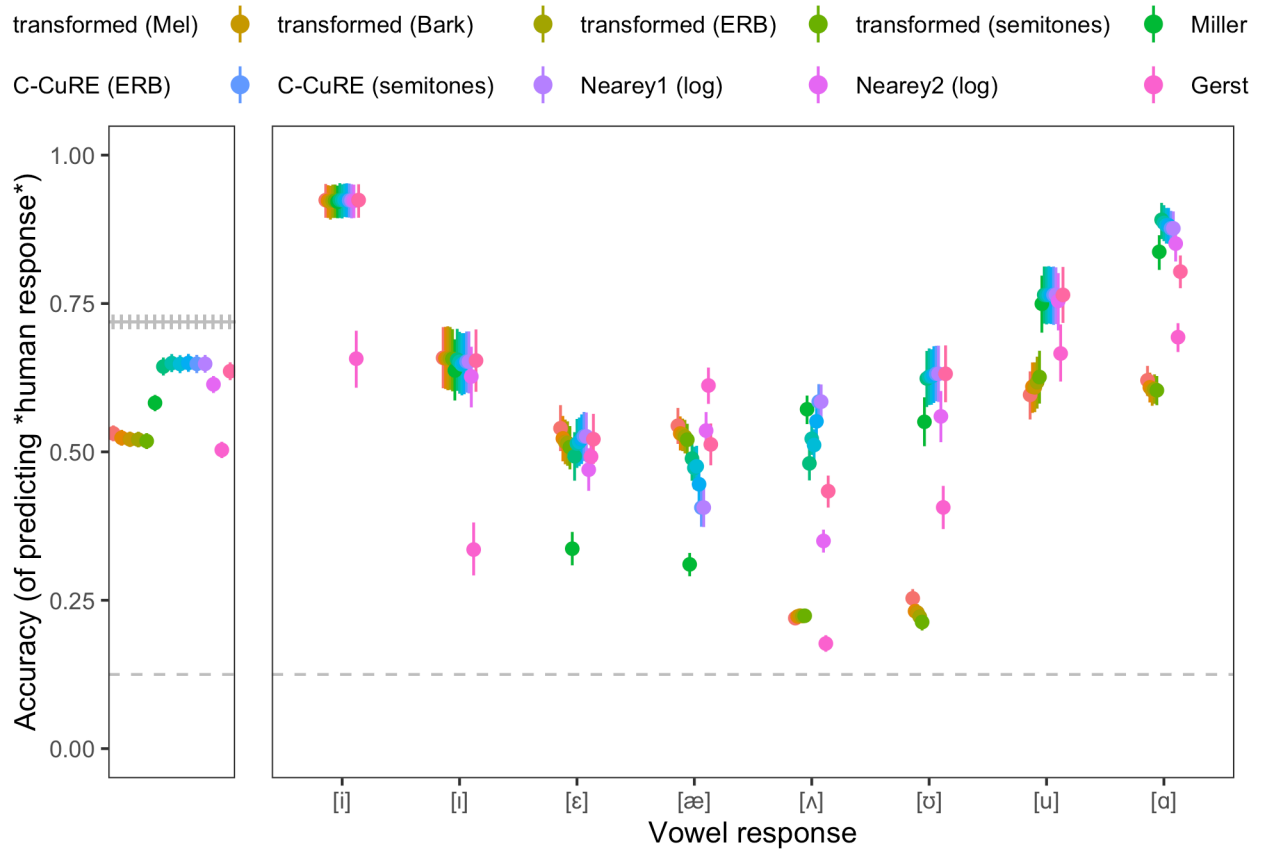
232

233 `$unhighlighted_params`234 `list()`

235

236 `$use_group_by`

```
237 NULL
238
239 $use_direct_label
240 NULL
241
242 $line_label_type
243 [1] "ggrepel_label"
244
245 $label_key
246 <quosure>
247 expr: ^NULL
248 env: empty
249
250 $label_params
251 $label_params$fill
252 [1] "white"
253
254
255 $keep_scales
256 [1] FALSE
257
258 $calculate_per_facet
259 [1] FALSE
260
261 attr(,"class")
262 [1] "gg_highlighter"
```



Fit log model to compare predictions of the ios (test of sign)+ evaluate against max accuracy and chance

## 6 Conclusions

In conclusion, these results indicate that pre-linguistic normalization (or computationally similar algorithms) contribute to the remarkable adaptive abilities of human speech perception. Our results further suggest that human perception only employs simple normalization operations—such as centering cues relative to a talker’s mean. Finally, we find that simple ideal observers achieve performance far above chance in predicting perception, although our results also indicate that human perception might employ additional adaptive algorithms (given the 7.5% accuracy difference between the best performing models and human performance).

## References

- Feldman, N. H., Griffiths, T. L., & Morgan, J. L. (2009). The influence of categories on perception: Explaining the perceptual magnet effect as optimal statistical inference. *Psychological Review*, 116, 752–782.
- Kleinschmidt, D. F., & Jaeger, T. F. (2015). Robust speech perception: Recognize the familiar, generalize to the similar, and adapt to the novel. *Psychological Review*, 122, 148–203. <https://doi.org/10.1037/a0038695>
- Kronrod, Y., Coppess, E., & Feldman, N. H. (2016). A unified model of categorical effects in consonant and vowel perception. *Psychological Bulletin and Review*, 1681–1712. <https://doi.org/10.3758/s13423-016-1049-y>