# Evaluating automatic speech recognition systems as quantitative models of cross-lingual phonetic category perception

Thomas Schatz, Francis Bach, Emmanuel Dupoux

**HAL Id: hal-01888735**
**https://hal.science/hal-01888735**

Submitted on 7 Dec 2018

# Evaluating automatic speech recognition systems as quantitative models of cross-lingual phonetic category perception

**Thomas Schatz**

*Department of Linguistics & UMIACS*
*University of Maryland*
*College Park, USA*

*Department of Linguistics*
*Massachusetts Institute of Technology*
*Cambridge, USA*

*LSCP, Département d'études cognitives de l'ENS*
*École normale supérieure, EHESS, CNRS, PSL Research University*
*29, rue d'Ulm 75005 Paris, France*
*thomas.schatz@laposte.net*


**Francis Bach**

*SIERRA project-team, Département d'informatique de l'ENS*
*École normale supérieure, INRIA, CNRS, PSL Research University*
*45, rue d'Ulm 75005 Paris, France*
*francis.bach@ens.fr*


**Emmanuel Dupoux**

*LSCP, Département d'études cognitives de l'ENS*
*École normale supérieure, EHESS, CNRS, PSL Research University*
*29, rue d'Ulm 75005 Paris, France*
*emmanuel.dupoux@gmail.com*

**Quantitative models of phonetic category perception**

Quantitative models of
phonetic category perception
Page 2
Schatz, JASA-EL

1      **Abstract**

2        Theories of cross-linguistic phonetic category perception posit that listeners per-
3    ceive foreign sounds by mapping them onto their native phonetic categories, but,
4    until now, no way to effectively implement this mapping has been proposed. In this
5    paper, Automatic Speech Recognition (ASR) systems trained on continuous speech
6    corpora are used to provide a fully specified mapping between foreign sounds and
7    native categories. We show how the *machine ABX* evaluation method can be used
8    to compare predictions from the resulting quantitative models with empirically at-
9    tested effects in human cross-linguistic phonetic category perception.

10

## 1. Introduction

The way we perceive phonetic categories (i.e. basic speech sounds such as consonants and vowels) is largely determined by the language(s) to which we were exposed as a child. For example, native speakers of Japanese have a hard time discriminating between American English (AE) /ɹ/ and /l/, a phonetic contrast that has no equivalent in Japanese (Goto, 1971; Miyawaki et al., 1975). Perceptual specialization to the phonological properties of the native language has been extensively investigated using a varieties of techniques (see Strange 1995 and Cutler 2012 for reviews). Many of the proposed theoretical accounts of this phenomenon concur that foreign sounds are not perceived faithfully, but rather, are 'mapped' onto one's pre-existing (native) phonetic categories, which act as a kind of 'filter' resulting in the degradation of some non-native contrasts (Best, 1995; Flege, 1995; Kuhl and Iverson, 1995; Werker and Curtin, 2005). In none of these theories, however, is the mapping specified in enough detail to allow a concrete implementation. In addition, in most of the existing theories[1], even if a fully specified mapping was available, it remains unclear how predictions on patterns of error rates could be derived from it (the filtering operation). These theories remain therefore mainly descriptive.

In this paper, we propose to leverage ASR technology to obtain fully specified mappings between foreign sounds and native categories and then use the *machine ABX* evaluation task (Schatz et al., 2013; Schatz, 2016) to derive quantitative predictions from these mappings regarding cross-linguistic phonetic category perception. More specifically, our approach can be broken down into three steps. First, train a *phoneme recognizer* in a 'native' language using annotated continuous speech recordings. Second, use the trained system to derive *perceptual representations* for test stimuli in a foreign language. In this paper, these will be vectors of posterior probabilities over each of the native phonemes. Third, obtain predictions for perceptual errors by running a *psychophysical test* over these representations for each foreign contrast. *Machine ABX* discrimination tasks will be used for this.

To showcase the possibilities offered by the approach, we look at predictions obtained for three empirically-attested effects in cross-linguistic phonetic category perception. The first two effects are *global* effects that apply to the set of phonetic contrasts in a language as a whole. First: native contrasts tend to be easier to distinguish than non-native ones (Gottfried, 1984). Second: patterns of perceptual confusions are function of the native language(s): two persons with the same native language tend to confuse the same foreign sounds, which can be different from sounds confused by persons with another native language (Strange, 1995). Thanks to the quantitative and systematic nature of the proposed approach, these effects are straightforward to study. We show that ASR models can account for both of them. Most effects documented in the empirical literature on cross-linguistic phonetic category perception are more *local* however. They describe patterns of confusion observed for very specific choices of languages and contrasts. We illustrate how such effects can be studied with our method through the classical example of AE /ɹ/-/l/ perception by native Japanese listeners (Goto, 1971; Miyawaki et al., 1975). We show that ASR models correctly predict the difficulty of perceiving this distinction for Japanese listeners.

Previous attempts at specifying mappings between foreign and native categories relied on phonological descriptions of the languages involved. Analyses at the level of abstract (context-independent) phonemes, however, were found not to be sufficient to fully account for perceptual data (Kohler, 1981; Strange et al., 2004). For example, the French [u-y] contrast can be either easy or hard to perceive for native AE listeners, depending on the specific phonetic context in which it is realized (Levy and Strange, 2002). Attempting to specify mappings *explicitly* through finer-grain phonetic analyses certainly remains an option, but involves a formidable amount of work. An

66 attractive and potentially less costly alternative consists in specifying mappings *implic-*
67 *itly*, through quantitative models of native speech perception. By this, we mean models
68 that map any input sound to a perceptual representation adapted to the model's 'native
69 language'. This representation can take the form of a phonetic category label, a vector
70 of posterior probabilities over possible phones or some other, possibly richer, form of
71 representation. Predictions regarding human perception of foreign speech sounds are
72 then derived by analyzing the 'native representations' produced by the model when
73 exposed to these foreign sounds.

74       Let us now explain the rationale for turning toward ASR technology, when the
75 goal is to model *human* speech perception. This approach is best understood in the
76 context of a top-down effort, where the focus is on developing models first at the *in-*
77 *formation processing* level, before considering issues at the algorithmic and biological
78 implementation levels (Marr, 1982). Native speech perception is thought to arise pri-
79 marily from a need to reliably identify the linguistic content in the language-specific
80 speech signal to which we are exposed, despite extensive para-linguistic variations.
81 ASR systems, whose goal is to map input speech to corresponding sequences of words,
82 face the same problem. ASR systems seek optimal performance, and can thus be inter-
83 esting as potential normative models of human behavior from an *efficient coding* point
84 of view (Barlow, 1961), even though biological plausibility is not taken into account in
85 their development.

86       We found two previous studies taking steps in the proposed direction. In the
87 first one (Strange et al., 2004), a Linear Discriminant Analysis model was trained to
88 classify AE vowels from F1/F2/F3 formant plus duration representations. The classi-
89 fication of North German vowels by this model was then compared to assimilation
90 patterns from a phoneme classification task performed by native AE speakers exposed
91 to North German vowels. The model's predictions only partially matched observed hu-
92 man behavior. In the second study (Gong et al., 2010), Hidden-Markov-Models (HMM)
93 with a structure inspired from ASR technology were trained to classify Mandarin con-
94 sonants from Mel-Frequency Cepstral Coefficients[2] (MFCC). The classification of AE
95 consonants by this model was then compared to assimilation patterns from a phoneme
96 classification task performed by native Mandarin speakers exposed to AE consonants.
97 There was a good consistency between model's predictions and human assimilation
98 patterns in most cases, although the model provided more variable answers overall
99 and differed markedly from humans in its preferred Mandarin classification of certain
100 AE fricatives.

101       The present work expands over these previous studies in several respects. First,
102 we replace ad hoc speech processing models trained on restricted stimuli[3] with general-
103 purpose ASR systems trained on natural continuous speech. This has both conceptual
104 and practical benefits. Conceptually, the information processing problem our models
105 attempt to solve is closer to the one solved by humans, who have to deal with the full
106 variability of natural speech. From a practical point of view, this allows us to capital-
107 ize on existing corpora of annotated speech recordings developed for ASR. A second
108 difference with previous studies is that we improve on the evaluation methodology,
109 by replacing informal analysis of assimilation patterns with quantitative evaluations
110 based on a simple model of an ABX discrimination task, leading to clean and clearly
111 interpretable results. Finally, we conduct more systematic evaluations, testing for two
112 *global* and one *local* effect in cross-linguistic phonetic category perception.

### 2. Methods

113

114 *2.1. Speech recordings*

115 To train and evaluate ASR models, 5 corpora of recorded speech in different languages
116 were used: a subset of the Wall Street Journal corpus (WSJ) (Paul and Baker, 1992),

117  the Buckeye corpus (BUC) (Pitt et al., 2005), a subset of the Corpus of Spontaneous
118  Japanese (CSJ) (Maekawa, 2003), the Global Phone Mandarin (GPM) corpus (Schultz,
119  2002) and the Global Phone Vietnamese (GPV) corpus (Vu and Schultz, 2009). Impor-
120  tant characteristics of the corpora are summarized in Table 1. Two corpora in American
121  English were included to dissociate *language-mismatch* effects, in which we are inter-
122  ested, from *channel-mismatch* effects due to differences across corpora in recording
123  conditions, microphones, speech register, etc. Phonetic transcriptions were obtained
124  by combining word-level transcriptions with a phonetic dictionary for the WSJ, BUC,
125  GPM and GPV corpora. For the CSJ corpus, manual phonetic transcriptions were used.
126  For all corpora, timestamps for the phonetic transcriptions were obtained by forced
127  alignment using an ASR system similar to those described in the next section, but
128  trained on the whole corpus.

129  *2.2. ASR models*

130  State-of-the-art ASR systems are built from deep recurrent neural networks. These sys-
131  tems, however, typically require hundreds of hours of data to be reliably trained and
132  we decided to focus in this study on using older, but more stable, Gaussian-Mixture
133  based Hidden-Markov Models (GMM-HMM) to ensure reasonable performance across
134  all corpora. Each corpus was randomly split into a training and a test set of approx-
135  imately the same size, each containing an equal number of speakers. There was no
136  overlap between training and test speakers. Models were trained with the Kaldi toolkit
137  (Povey et al., 2011) using the same recipe with the same parameters and input fea-
138  tures to train all models[4]. The Word-Error Rate[5] (WER) on the test set for each of the
139  resulting models is reported in Table 1.

140          We will not attempt to describe the inner workings of the models beyond men-
141  tioning that a generative model is trained for each phone, with explicit mechanisms for
142  handling variability due to changes in speaker, phonetic context or word-position. We
143  refer to the Kaldi documentation for further detail [6]. Input to the models takes the form
144  of 39 MFCC coefficients[7] plus 9 pitch-related features[8] extracted every 10ms of signal.
145  These 48-dimensional input features can be seen as a *universal* auditory-like baseline
146  representation that is not tuned to any particular 'native language'. The model pro-
147  duces 'native' representations under the form of output vectors produced every 10ms,
148  which list the posterior probabilities, according to the model, that the corresponding
149  stretch of speech signal belongs to each of the segment in the phonemic inventory of
150  the model's 'native language'[9]. The test set of each corpus is decoded with each of the
151  5 ASR models and we also use the input features directly, without any GMM-HMM
152  decoding, as a language-independent control, yielding a total of 6 different represen-
153  tations of each corpus to be evaluated.

Table 1. Word-Error-Rates obtained by the ASR systems trained on each corpus as
well as the language, total duration, speech register and number of speakers for
each corpus. AE stands for American English, Spont. stands for Spontaneous.

| Corpus | Language | Time | Type | Spk | WER |
|--------|----------|------|------|-----|------|
| WSJ | AE | 143h | Read | 338 | 8.5% |
| BUC | AE | 19h | Spont. | 40 | 48.0% |
| CSJ | Japanese | 15h | Spont. | 75 | 30.0% |
| GPM | Mandarin | 30h | Read | 132 | 31.0% |
| GPV | Vietnamese | 20h | Read | 129 | 23.5% |

154 *2.3. Machine ABX evaluation*

155 We evaluate our ASR models with a machine version of an ABX discrimination task
156 (Schatz et al., 2013; Schatz, 2016) that allows us to quantify how easy it is to distin-
157 guish two phonetic categories based on representations produced by one of our models.
158 The basic idea is to take two acoustic realizations $A$ and $X$ from one of the phonetic
159 categories and one acoustic realization $B$ from the other category and to test whether
160 the model representation for $X$ is closer to the model representation for $A$ than to
161 the model representation for $B$. The probability for this to be false for $A$, $B$ and $X$
162 randomly chosen in a corpus is defined as the *ABX error rate* for the two phonetic
163 categories according to the model. If it is equal to $0$, the two categories are perfectly
164 discriminated. If it is equal to $.5$, discrimination is at chance level.

165       For each $A$, $B$ and $X$ triplet, we use the phone-level time alignments to select
166 corresponding model representations. Because the stimuli have variable durations, the
167 resulting representations can have different lengths. To find a good alignment and
168 obtain a quantitative measure of dissimilarity between $A$ and $X$ and $B$ and $X$, we use
169 Dynamic Time Warping based on a frame-wise symmetric Kullback-Leibler divergence
170 for posterior probability vectors and a frame-wise cosine distance for the input features
171 control. In the specific ABX task considered here, we select only triplets such that $A$, $B$
172 and $X$ occur in the same phonetic context (same preceding phone and same following
173 phone) and are uttered by the same speaker. For each phonetic contrast an aggregated
174 ABX error rate is obtained by averaging over stimulus order, context and speaker. Let
175 us illustrate this through the example of the /u/-/i/ contrast. First, we average error
176 rates obtained when A and X are chosen to be /u/ and B is chosen to be /i/ and vice-
177 versa, then we average over all possible choices of speaker and finally we average over
178 all possible choices of preceding and following phones. We either report directly the
179 scores obtained for individual phonetic contrasts or we average them over interesting
180 classes of contrasts, such as consonant contrasts or vowel contrasts.

181       Note that, because we are studying very robust empirical effects that reflect
182 what subjects learn outside the lab and that are expected to be observed in any well-
183 designed experimental task, our evaluation method focus on simplicity of application
184 rather than detailed modeling of human performance in a specific experimental setting.

185 **3. Results**

186 See supplementary material for the raw (unanalyzed) confusion matrices obtained for
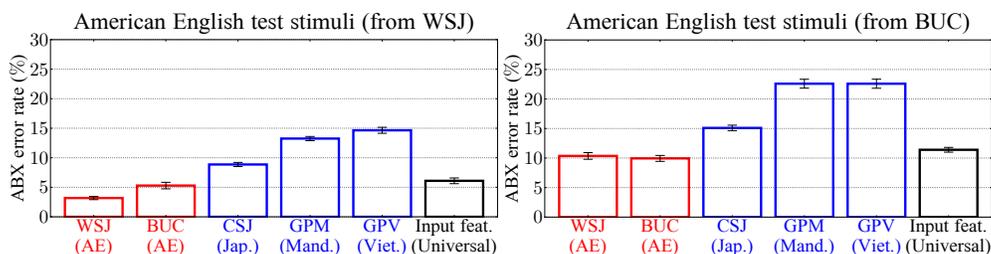187 each model on each test corpus.

188 *3.1. Native vs. non-native contrasts*



Fig. 1. (color online) ABX error-rates averaged over all consonant contrasts of AE.
Left: using stimuli from the WSJ corpus test set. Right: using stimuli from the BUC
corpus test set.

189       Native phonetic categories are easier to distinguish than non-native categories
190 (Gottfried, 1984). This is consistent with the predictions of our models shown in Figure

1. The AE models (in red) separate AE phonetic categories better than other models (in blue). This is true even when they are tested with AE stimuli from a corpus different from the one on which they were trained, showing that the differences observed cannot be explained simply by *channel-mismatch* effects and reflect a true *language-specificity* of the representations learned by the models. Another interesting observation is that, while a moderate improvement in phone separability is observed when comparing 'native' AE models to the 'universal' input features control, the most salient effect is a large decrease in performance for 'non-native' models. A possible interpretation is that, while ASR models can provide categorical representations of 'native' speech that are much more compact than the input features, they do it at the expense of a loss of representation power for coding speech in other languages[10].

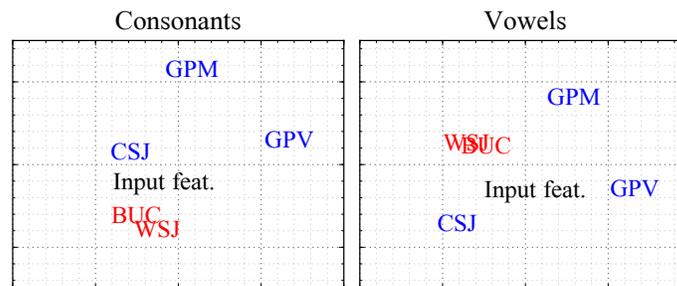*3.2. Native-language-specific confusion patterns*



Fig. 2. (color online) Two-dimensional embeddings of the different models based on the average cosine similarity between their patterns of ABX errors across the five test corpora. The distance between models in the embedding space directly reflects whether they make the same type of confusions or not. Left: for consonant contrasts. Right: for vowel contrasts. Text labels are centered horizontally and vertically on the point they represent.

The specific confusions we make between sounds of a foreign language differ according to our native language (Strange, 1995). Consistent with this effect, Figure 2 shows that, for both consonant and vowel contrasts, the confusion patterns obtained with the two AE models over the different corpora are more similar to each other than to the confusion patterns obtained with models trained on other languages. Confusion patterns for input features occupy a somewhat central role. In this figure, the distance between two points is proportional to the observed similarity between confusion patterns obtained from the associated models[11]. Confusion patterns on a given corpus consist of vectors listing the ABX errors for either all consonant contrasts or all vowel contrasts in this corpus. For example for a language with $n$ consonants, $n(n-1)/2$ consonant contrasts can be formed and the corresponding ABX errors are listed in a vector of size $n(n-1)/2$. The similarity between confusion patterns of two models is defined as the average of the cosine similarity between the confusion patterns obtained with these models on each of the five corpora[12]. Importantly, the rescaling invariance of the cosine similarity ensures that our analysis of confusion patterns is independent from the average ABX error rates studied in Section 3.1.

*3.3. Japanese listeners and American English /ɹ/-/l/*

AE /ɹ/ and /l/ are much harder to perceive for Japanese than for AE native speakers (Goto, 1971; Miyawaki et al., 1975). Figure 3 shows that our models' predictions are fully consistent with this effect: when comparing the Japanese model to both AE models and to the input features, the /ɹ/-/l/ discriminability drops spectacularly, much
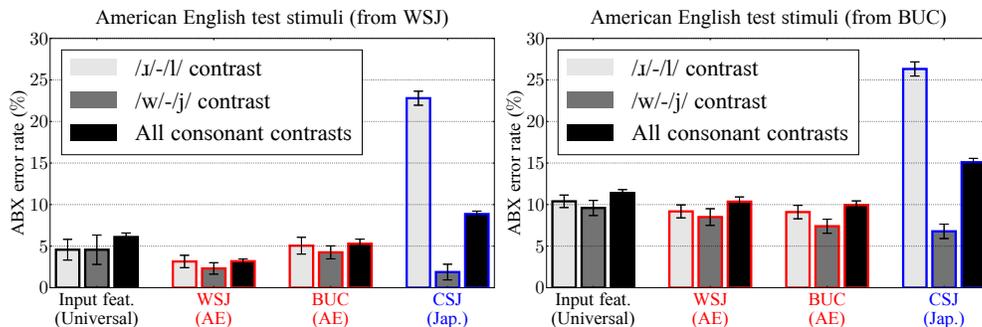
Fig. 3. (color online) Comparison of the ABX error-rates obtained with the input features, with the two AE models and with the Japanese model on the AE /ɹ/-/l/ contrast. ABX Error-rates for the /w/-/j/ contrast and ABX Error-rates averaged over all consonant contrasts of AE are also shown as controls. Left: using stimuli from the WSJ corpus test set. Right: using stimuli from the BUC corpus test set.

more than the discriminability of two controls. This is observed both when using test stimuli from the WSJ and from the BUC corpora.The first control is the AE /w/-/j/ contrast. Like /ɹ/ and /l/, /w/ and /j/ are liquid consonants, but unlike those, they have a clear counterpart in Japanese. The second control is the average ABX error rate from Section 3.1. This control allows to check that there is a specific deficit of the Japanese model on AE /ɹ/-/l/ discrimination, that cannot be explained by an overall weakness of this model.

## 4. Discussion

Fully specified mappings between foreign sounds and native phonetic categories were obtained for several language pairs through GMM-HMM ASR systems. Coupled with a simple model of a discrimination task, they successfully accounted for several empirically attested effects in cross-linguistic phonetic category perception by monolingual listeners. This includes two types of *global* effects: first, that the phonetic categories of a language are overall harder to discriminate for non-native speakers than for native speakers and second, that the pattern of confusions between phonetic categories for non-native speakers is specific to their native language (e.g. native speakers of Japanese do not make the same confusions between phonetic categories of American English than native speakers of French). We also showed that the proposed model can account for a well-known *local* effect: American English /ɹ/ and /l/ are very hard to discriminate for native speakers of Japanese.

These results provide a proof-of-concept for the proposed approach to evaluating ASR systems as quantitative models of phonetic category perception. They also show promise regarding the possibility of modeling human phonetic category perception with ASR systems. Yet we do not claim, at this point, to have provided definitive evidence that the particular GMM-HMM ASR systems considered provide the best, or even a particularly 'good', such model. A host of *local* effects have been documented in the empirical literature on phonetic category perception beyond the one investigated here (Strange, 1995; Cutler, 2012) and the empirical adequacy of the proposed models with respect to more of these effects will need to be determined before any conclusion can be reached. Effects that are hard to predict from conventional phonological analyses, such as how the phonetic or prosodic context can modulate the difficulty of perceiving certain foreign contrasts (Levy and Strange, 2002; Kohler, 1981; Strange et al., 2004), should be of particular interest. Finally, let us underline that we only investigated predictions obtained with one particular ASR architecture. There

258 are multiple ways of instantiating ASR systems, which might yield different predic-
259 tions. For example, modeling variability in the signal due to the phonetic context
260 explicitly with context-dependent phone models, as in this article, or implicitly with
261 context-independent phone models, might affect predictions regarding the aforemen-
262 tioned context-dependent effects. Another example of a potentially significant decision
263 is whether to use HMM-GMM or neural-network systems. HMM models have known
264 structural limitations for modeling segment duration (Pylkkönen and Kurimo, 2004),
265 from which neural-network models do not suffer. Thus, neural-network ASR systems
266 may provide better models of native perception in languages like Japanese, where du-
267 ration is contrastive. The multiplicity of documented empirical effects and available
268 computational models calls for an extensive investigation, which could in turn trigger
269 a more systematic *experimental* investigation of non-native perception and result in
270 applications in foreign language education.

## Acknowledgments

## Notes

283 [1] Best 1995 being a possible exception.
284 [2] MFCC (Mermelstein, 1976) are speech features commonly used as a front-end to ASR systems. They
285 can be thought of as moderate-dimensional descriptor ($d = 13$) of the whole shape of regularly-spaced
286 spectral-slices in a mel-scale log-spectrogram. They are usually taken every 10ms and augmented with their
287 first and second time derivatives to incorporate dynamic information, leading to 100 vector descriptors of
288 dimension $d = 39$ per second of signal.
289 [3] Previous studies used as training stimuli a limited sample of 264 AE vowels occurring either in
290 [hVba] context or within a unique carrier sentence (Strange et al., 2004) and 3331 Chinese consonants
291 occurring in isolated VCV context (Gong et al., 2010).
292 [4] See https://goo.gl/RsKMA3.
293 [5] Error-rate obtained in a word recognition task using the trained acoustic model with a language
294 model (in our case a word-level bigram estimated from the training set).
295 [6] See http://kaldi-asr.org/.
296 [7] See footnote 1.
297 [8] Pitch features were added because two of the languages considered (Mandarin and Vietnamese) are
298 tonal languages.
299 [9] More specifically, we use Viterbi-smoothed phone-level posteriorgrams obtained with a phone-level
300 bigram language model estimated on the training set of each corpus.
301 [10] Note that Renshaw et al. (2015) observed a different pattern when testing a neural-network-based
302 ASR system trained on AE on the Xitsonga language: the 'AE-native' model improved Xitsonga phone sep-
303 arability relative to the input features control. There are, at least, two possible interpretations for this dis-
304 crepancy: it could be due to general differences between GMM-HMM and neural-network architectures or
305 it could be due to differences in the representation format chosen (they used 'bottleneck features' extracted
306 from a middle layer of the neural network, which are not constrained to represent phonetic categories, while
307 our posterior features are)
308 [11] Two-dimensional embeddings are obtained with scikit-learn's non-metric multi-dimensional-scaling.
309 [12] Observed range of cosine similarities: [0.90-0.96] for consonants and [0.85-0.94] for vowels.

## References and links

311 H. B. Barlow. *Possible principles underlying the transformations of sensory messages*. MIT press, 1961.
312 C. T. Best. A direct realist view of cross-language speech perception. *Speech Perception and Linguistic
313 Experience: Issues in Cross-Language Research*, pages 171–204, 1995.

Quantitative models of
phonetic category perception
Page 10
Schatz, JASA-EL

314   A. Cutler. *Native listening: Language experience and the recognition of spoken words*. Mit Press, 2012.

315   J. E. Flege. Second language speech learning: Theory, findings, and problems. *Speech perception and*
316   *linguistic experience: Issues in cross-language research*, pages 233–277, 1995.

317   J. Gong, M. Cooke, and M. Garcia Lecumberri. Towards a quantitative model of mandarin chinese
318   perception of english consonants. *Proc. NewSounds 2010*, 2010.

319   H. Goto. Auditory perception by normal japanese adults of the sounds l and r. *Neuropsychologia*, 9
320   (3):317–323, 1971.

321   T. L. Gottfried. Effects of consonant context on the perception of french vowels. *Journal of Phonetics*,
322   12(2):91–114, 1984.

323   K. Kohler. Contrastive phonology and the acquisition of phonetic skills. *Phonetica*, 38(4):213–226,
324   1981.

325   P. K. Kuhl and P. Iverson. Linguistic experience and the perceptual magnet effect. *Speech perception*
326   *and linguistic experience: Issues in cross-language research*, pages 121–154, 1995.

327   E. S. Levy and W. Strange. Effects of consonantal context on perception of french rounded vowels by
328   american english adults with and without french language experience. *The Journal of the Acoustical*
329   *Society of America*, 111(5):2361–2362, 2002.

330   K. Maekawa. Corpus of spontaneous japanese: Its design and evaluation. In *Proc. ISCA & IEEE*
331   *Workshop on Spontaneous Speech Processing and Recognition*, 2003.

332   D. Marr. *Vision: A computational approach*. Freeman.[aAC], 1982.

333   P. Mermelstein. Distance measures for speech recognition, psychological and instrumental. *Pattern*
334   *recognition and artificial intelligence*, 116:91–103, 1976.

335   K. Miyawaki, J. J. Jenkins, W. Strange, A. M. Liberman, R. Verbrugge, and O. Fujimura. An effect of
336   linguistic experience: The discrimination of [r] and [l] by native speakers of japanese and english.
337   *Perception & Psychophysics*, 18(5):331–340, 1975.

338   D. B. Paul and J. M. Baker. The design for the wall street journal-based csr corpus. In *Proc. Workshop*
339   *on Speech and Natural Language*, pages 357–362, 1992.

340   M. A. Pitt, K. Johnson, E. Hume, S. Kiesling, and W. Raymond. The buckeye corpus of conversational
341   speech: Labeling conventions and a test of transcriber reliability. *Speech Communication*, 45(1):89–
342   95, 2005.

343   D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlček,
344   Y. Qian, P. Schwarz, et al. The kaldi speech recognition toolkit. In *Proc. Workshop on Automatic*
345   *Speech Recognition and Understanding*, 2011.

346   J. Pylkkönen and M. Kurimo. Duration modeling techniques for continuous speech recognition. In
347   *Proc. INTERSPEECH*, 2004.

348   D. Renshaw, H. Kamper, A. Jansen, and S. Goldwater. A comparison of neural network methods for
349   unsupervised representation learning on the zero resource speech challenge. In *Proc. INTERSPEECH*,
350   2015.

351   T. Schatz. *ABX-Discriminability Measures and Applications*. Doctoral dissertation, Université Paris 6
352   (UPMC), 2016.

353   T. Schatz, V. Peddinti, F. Bach, A. Jansen, H. Hermansky, and E. Dupoux. Evaluating speech features
354   with the minimal-pair ABX task: Analysis of the classical MFC/PLP pipeline. In *Proc. INTERSPEECH*,
355   2013.

356   T. Schultz. Globalphone: a multilingual speech and text database developed at karlsruhe university.
357   In *Proc. INTERSPEECH*, 2002.

358   W. Strange. *Speech perception and linguistic experience: Issues in cross-language research*. York Press,
359   1995.

360   W. Strange, O.-S. Bohn, S. A. Trent, and K. Nishi. Acoustic and perceptual similarity of north german
361   and american english vowels. *The Journal of the Acoustical Society of America*, 115(4):1791–1807,
362   2004.

363   N. T. Vu and T. Schultz. Vietnamese large vocabulary continuous speech recognition. In *Proc. ASRU*,
364   2009.

365   J. F. Werker and S. Curtin. Primir: A developmental framework of infant speech processing. *Language*
366   *learning and development*, 1(2):197–234, 2005.