



Acoustic cue variability affects eye movement behaviour during non-native speech perception: a GMM model

Jessie S. Nixon¹, Catherine T. Best²

¹New Zealand Institute of Language, Brain and Behaviour, New Zealand

²MARCS Institute, Western Sydney University, Australia

jess.s.nixon@gmail.com, C.Best@westernsydney.edu.au

Abstract

Participants in the ‘visual world’ paradigm simultaneously process both auditory and visual cues in order to match speech to target images. Previous research has shown that when native speakers listen to speech that has high within-category variability in the contrastive dimension, auditory perceptual uncertainty is reduced, resulting in increased looks to competitor objects. This suggests a cross-modal effect, where reduced reliability in the auditory domain leads to increased search for evidence in the visual domain.

The present study investigated the effects of within-category acoustic variability on eye movements during the acquisition of a new acoustic dimension not present in the native language, namely English speakers’ acquisition of lexical tone. All participants heard a bimodal distribution of stimuli, with distribution peaks at the prototypical pitch values for high and mid tones; however, presentation frequency differed between conditions: high-variance vs. low-variance. Based on previous research, we expected lower uncertainty and better learning in the low-variance condition.

GMM models of eye movement data showed that within-category acoustic variance increases perceptual uncertainty in the auditory domain and hinders acquisition of a cue dimension. The results also show a cross-modal effect: lower reliability in the auditory domain leads to increased search for cues in the visual domain, even when visual cues are held constant across conditions.

Index Terms: speech perception, cross-modal effects, statistical learning, second language acquisition, Cantonese, lexical tone, visual world eyetracking, generalised additive mixed models (GAMMs)

1. Introduction

In the ‘visual world’ eyetracking paradigm [1], participants simultaneously process auditory speech cues and visual cues in order to match auditory input to target images. Previous research suggests that, as the reliability of the acoustic signal decreases, the search for additional *visual* cues to support or reject activated candidates increases [2, 3, 4].

The organisation of acoustic cues varies substantially between languages. Cue dimensions that are lexically contrastive in one language may not be contrastive in another. Therefore, acquisition of a new language often involves learning to substantially adjust cue weights (i.e. to adjust the degree to which various cues in the signal are utilised, consciously or unconsciously) for lexical contrasts. In some cases, this can pose significant challenges. Expert knowledge of statistical regularities in one’s native language can lead to expectations that hinder non-native speech perception [5, 6]. Statistical properties that seem to play a role in shaping cue perception in-

clude the number of distribution peaks along a cue dimension [7, 8, 9, 10], acoustic distance between peaks in a bimodal distribution [11, 12] and within-category acoustic variance [2, 3, 4].

Many recent studies have emphasised the role of variability in shaping and reshaping native and non-native sound systems. For example, early first language acquisition [13, 14] and second language acquisition [15, 16, 17, 18] seem to benefit from multiple speakers, compared to a single speaker in the training input. When there are multiple speakers in the training input, this increases variability in *non-contrastive* indexical dimensions, which seems to have the effect of highlighting the relative invariance of the contrastive dimensions. This is consistent with learning models, which posit that learning not only involves acquisition of knowledge, but also learning to ignore irrelevant cues.

An aspect that has received less attention is variability within *contrastive* dimensions. Nixon and colleagues [4] investigated the temporal dynamics of perceptual uncertainty during native speech perception. Using the visual world paradigm, they found that effects of acoustic variance had a direct effect on visual processing. Effects emerged very early, in the first fixations of the trial. As auditory variability increased and speech cues became less reliable, listeners looked around more in search of visual cues to provide further support for partially activated candidates. The idea that listeners were seeking additional evidence in the high-variability condition seems to suggest the appropriate conditions for adjusting cue weights, and perhaps increasing weights of previously downweighted cues.

What is not yet known is whether such within-category acoustic variance also affects audio-visual processing during acquisition of a new acoustic dimension in a non-native language. The present study aimed to address this question by examining the effect of within-category acoustic variance on eye movements during native English speakers’ acquisition of a pitch cue (fundamental frequency; f_0) in a Cantonese lexical tone contrast. English does not use f_0 as a lexical contrast, and tone can be notoriously difficult for beginning learners of non-tonal languages. Based on previous studies [2, 3, 4], we expected greater weighting of the pitch cue over the course of the experiment - that is, better learning - in the low-variance, compared to the high-variance condition.

2. Method

2.1. Participants

Thirty-seven native English-speaking students from the University of Western Sydney who had not previously studied any tone language were recruited for the experiment for course

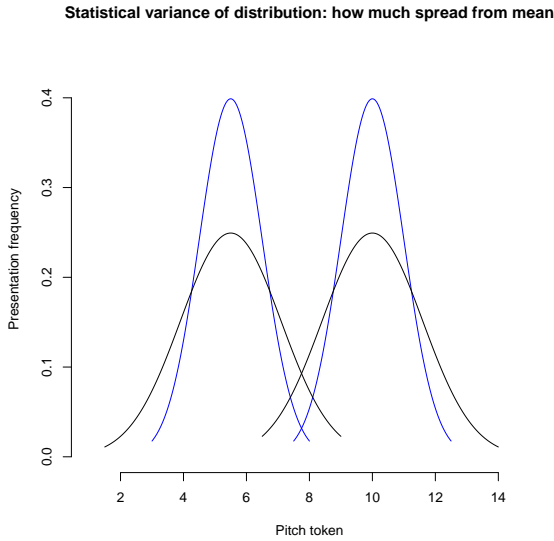


Figure 1: Illustration of the presentation frequency distributions in the high-variance (black lines) and low-variance conditions (blue lines).

credit¹. Participants were tested individually in a sound-attenuated booth.

2.2. Experiment design and stimuli

Visual stimuli were black-on-white line drawings of eight common objects. Auditory stimuli were four minimal pairs of mid- and high-tone words (e.g. *gon_mid* and *gon_high*). They were real Cantonese nouns; however, the matching images were not the true meaning in Cantonese. The purpose of using common, very high-frequency objects was to reduce cognitive load in the language-learning task. All auditory stimuli were recorded by a male native speaker of Hong Kong Cantonese. Stimuli were then resynthesised into a 14-step pitch continuum (e.g. *gon_mid* to *gon_high*) using PRAAT [19]. One half of the continuum corresponded to the mid tone and one half to the high tone.

The number of times participants heard each token of the continuum followed a bimodal distribution, with the two peaks of the distribution corresponding to the prototypical f0 for the mid- and high-tone stimuli, respectively. All participants heard the same number of tokens; but the number of times they heard each token differed between conditions, with greater spread from the mean (statistical variance) in the *high-variance* versus the *low-variance* distribution (see Figure 1). The experiment consisted of 240 experimental trials, divided into six blocks of 40 trials, with breaks between the blocks. The order of presentation was pseudo-randomised for each participant.

2.3. Procedure

Participants sat at a viewing distance of 60 cm from a computer screen equipped with an SR Research Eyelink 1000 remote eye-tracker. A chinrest and headrest were used to minimise movement. Stimulus presentation and data acquisition were con-

¹Participants were not explicitly asked whether they had studied a tone language, as this might influence the experiment results. Instead, they were asked to list all languages they spoke or had studied, and were screened if they did not meet this criterion.

ducted using SR Research Experiment Builder computer software with a sampling rate of 1000 Hz. The session began with ten practice trials to familiarize participants with the experimental procedure. None of the images or auditory stimuli from the experimental block appeared in the practice block. Each experimental trial began with a brief (1000 ms) presentation of four pictures, one in each quadrant of the screen. The purpose of the preview was to reduce noise in the data by reducing the time and likelihood of participants scanning the images at the beginning of the trial. The display always contained a target, a competitor and two distractor items. The target and competitor had the same segmental syllable, but differed in tone. The location of each picture condition on the screen and their location relative to each other were randomised to avoid strategic effects. The preview disappeared, followed by a gaze-contingent fixation cross to ensure participants were fixating the centre of the screen at the beginning of the critical trial period. The pictures then reappeared simultaneously with presentation of the auditory stimulus. Participants were instructed to select the picture corresponding to the word they heard by clicking on it with the mouse, and to guess if they did not know. They were given feedback ('correct'/'wrong') after each trial. Participants were told that this was a language-learning task, but were not informed about the pitch or tone manipulation or the target language.

3. Analysis

Eye movement data were analysed using *Generalised Additive Mixed Models* [20, 21, 22, GAMM] using the *mgcv* package (version 1.8.13) in R [23, version 3.3.0]. Generalised Additive Models (GAMs) are a type of Generalised Linear Models (GLM) that use smooth functions to model nonlinear effects of continuous predictors. The 'mixed' in GAMMs refers to the inclusion of random effects in addition to fixed effects.

GAM is a well-established method of analysis that is increasingly being used in the cognitive and language sciences, and has been applied to EEG data [24, 25, 26, 27], reaction times [28, 29], articulatory [30, 31], acoustic analysis [32], temporal clustering of sociolinguistic variants [33] and dialectology [34]. Recently, it has also been applied to eye movement data in the visual world paradigm [3, 4]. GAMMs have also been used to analyse single-image eye movement data [35] and pupilometry [36].

GAMMs are a valuable method for analysing visual world fixation data for several reasons, including their ability to capture nonlinear changes in eye movements over the course of the trial and/or over the course of the experiment, the inclusion of random effects to deal with taking repeated measures from the same participants and items, and methods for dealing with autocorrelation [4]. An important aspect of eyetracking data is how fixations change over time. In experimental data sets, and especially time series data, autocorrelation can occur between data points [37]. In the *mgcv* package, functions have been implemented to deal with autocorrelation in GAM models.

Eye movement data were modelled as a continuous predictor of Euclidean distance of fixations from the centre of the target image [4]. A sample trial which illustrates the *Target Distance* measure is shown in Figure ?? . All predictors of interest were entered into the initial model, and predictors that did not contribute to model fit were removed. Model comparison was conducted by means of χ^2 tests of fREML scores, using the *compareML* function in the *itsadug* package [38, version 2.2] in R. Because we were interested in the time course of fixations over the course of the trial, the predictor *time* was

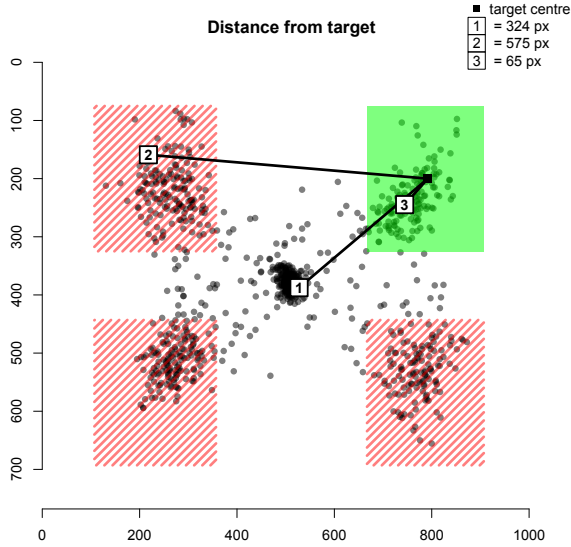


Figure 2: Illustration of the dependent measure Euclidean distance from the target. The x and y positions of the screen are on the x and y axes, respectively. Distance (in pixels) is calculated from the centre of the target picture.

included. Data was downsampled to 50 Hz to reduce autocorrelation between data points. A 3200 ms window was selected for analysis, from 200 ms prior to 3000 ms after auditory stimulus presentation. To test whether there was a learning effect over the course of the experiment, the model included a predictor of trial, centred around 0 (*centred trial*). To determine whether participants were using pitch as a cue to distinguish between target and competitor images, the model included a continuous predictor of pitch, also centred around zero (*centred pitch*). The centred values ranged from -5.5 to 5.5, with the distribution peaks at -3 and 3. Distribution variance was modelled as a two-level factor, low-variance and high-variance. Previous research with the visual world paradigm has shown that the location of the target object on the screen significantly affects eye movement behaviour [39, 4]. Therefore, a smooth for *target position* over time was included as a control variable, a factor with four levels: top-left, top-right, bottom-left and bottom-right. Random smooths for subject by item over time and subject by item over trial were included to account for differences in individual participants and items.

The initial model included intercepts for the two factor variables, variance condition and target position, and smooths (for each of the main effects) and nonlinear regression lines² (for each two- and three-way interaction) for each level of condition. A smooth was also included for each level of target position. Random effects were modelled with shrunk factor smooths. After running the model, the model residuals were examined to check for autocorrelation. An AR1 model was included to account for autocorrelation in the residuals.

4. Results

Model comparisons showed that model fit was improved by including smooths for centred trial by condition ($p < .001$); target

²The partial effects tensors are modelled with the `ti()` function in the `mgcv` package.

position over time ($p < .001$); and nonlinear regression smooths for time by trial ($p < .01$); time by pitch ($p < .05$); trial by pitch by condition ($p < .001$); time by pitch by condition ($p < .01$); trial by time by pitch ($p < .001$); and trial by time by pitch by condition ($p < .001$).

Figure 3 shows the difference between the high-variance and low-variance conditions (high minus low) over time for four representative centred pitch values, 3.5 (top left panel), 1.5 (top right panel), -1.5 (bottom left panel) and -3.5 (bottom right panel). To examine the results of exposure to the two distributions, a trial near the end of the experiment (centred trial 109) is selected. For both high and mid tones, fixations are significantly closer to the target in the low-variance condition, compared to the high-variance condition. For the high tone (pitch 3.5 and 1.5), this effect emerges early and increases over the remainder of the trial period. Interestingly, there is a nonlinear effect of pitch over time: the difference between conditions is greater and emerges earlier for the stronger cue, the higher pitch value (3.5), compared to the pitch value closer to the category boundary (1.5). This is similar to the nonlinear effect of pitch value that Nixon and colleagues [4] found for native Cantonese listeners, reflecting the ‘perceptual magnet effect’. For the mid tone (pitch -3.5 and -1.5), the low- and high-variance conditions diverge later than the high tones: around 1 second after auditory stimulus presentation for the mid tone, compared to around 500 ms for the high tones. This will be returned to below.

The difference between the high-variance and low-variance conditions over the course of the experiment is shown in Figure 4 for four representative pitch values (-3.5, -1.5, 1.5 and 3.5) at 2500 ms. Interestingly, at the beginning of the experiment, the distance from the target is greater in the low-variance condition than the high-variance condition for all pitch values. However, as the experiment progresses, the distance gets smaller, and by the end of the experiment, the distance from the target is significantly smaller in the low-variance condition.

The summed effects are shown in Figure 5. The figure shows a topographic plot of the Euclidean distance of fixations from the target object in the low-variance (left panels) and high-variance conditions (right panels). Time (ms) is on the horizontal axis. Centred pitch is on the vertical axis: positive pitch values correspond to the high tone and negative values to the mid tone. Distance from the target (in pixels) is on the z-axis and is colour-coded. Higher values (warmer colours) indicate fixations were further from the target image; lower values (cooler colours) indicate fixations were closer to the target image. The key in the top right corner indicates the corresponding values and z-limits. The panel rows show snapshots of trials throughout the experiment.

At the beginning of the experiment, fixations are further from the target in the low-variance condition, compared to the high-variance condition, for most of the trial. This effect lessens over the following trials, and has disappeared by trial 100. For the remainder of the experiment, fixations become gradually closer to the target in the low-variance condition. This becomes significant earlier in the low pitch values, as seen in trial 170. By the end of the experiment, fixations are significantly closer to the target in the low-variance condition for both the low and high pitch values.

5. Discussion

The present study investigated the effects of within-category acoustic variance on audio-visual processing during non-native acquisition of a new acoustic dimension, that is, pitch (f0) in

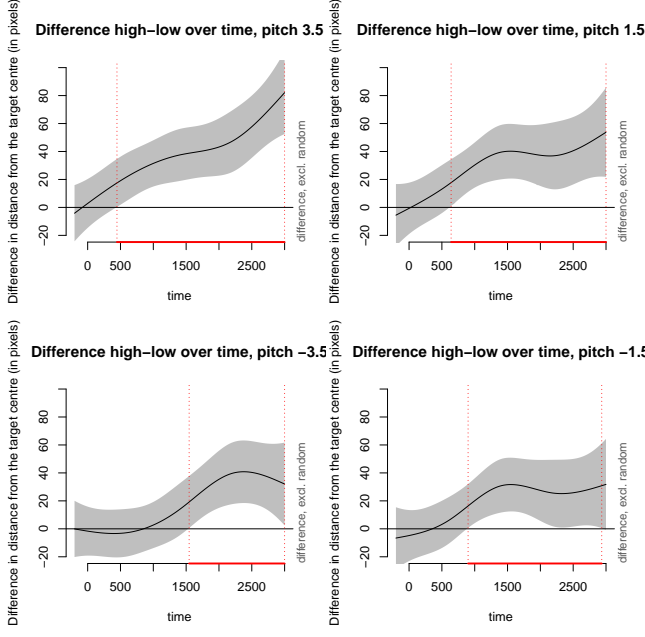


Figure 3: Smooth of the difference between the high-variance and low-variance conditions over time for centred pitch values 3.5 (top left), 1.5 (top right), -1.5 (bottom left) and -3.5 (bottom right). Time (ms) is on the x-axis. The difference (in pixels) between conditions (high minus low) is on the y-axis. Centred trial is set to 109.5. Random effects are removed from these plots.

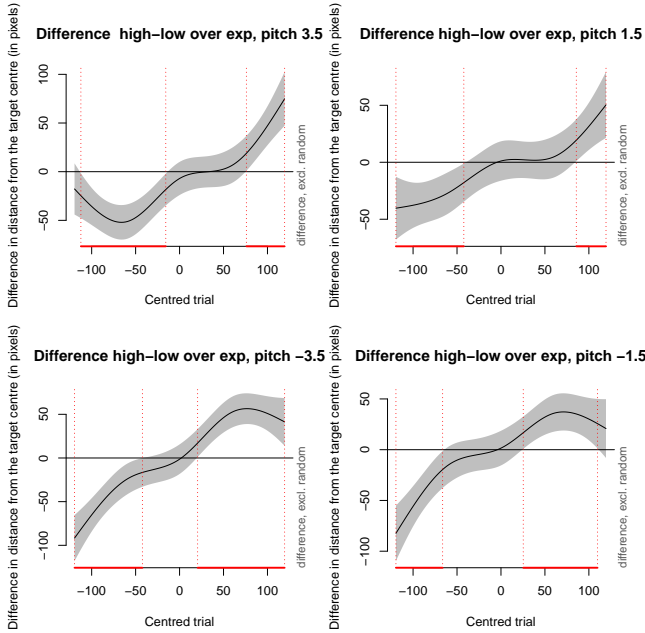


Figure 4: Smooth of the difference between the high-variance and low-variance conditions over the course of the experiment for centred pitch values 3.5 (top left), 1.5 (top right), and -3.5 (bottom left) and -1.5 (bottom right). Trial is on the x-axis and is centred around 0. The difference (in pixels) between conditions (high minus low) is on the y-axis. Time is set to 2500 ms. Random effects are removed from these plots.

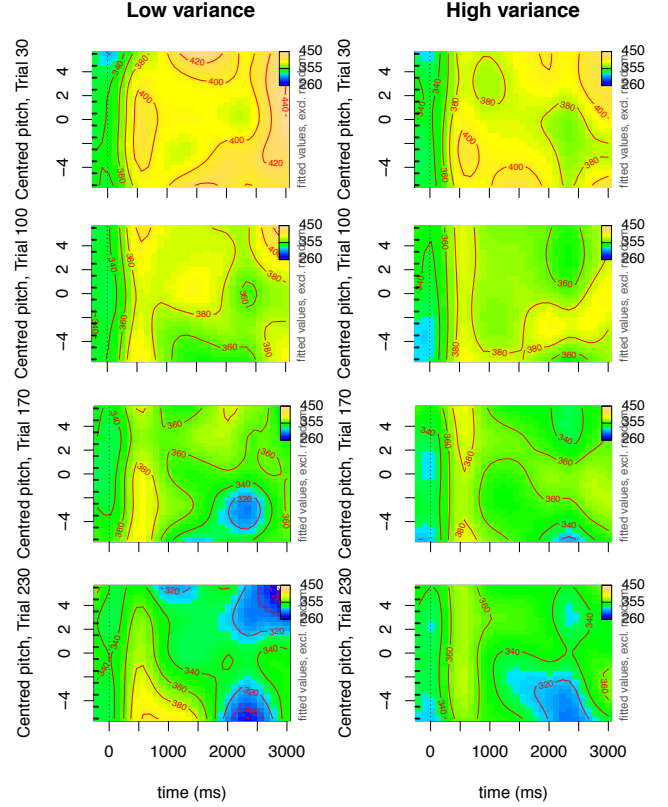


Figure 5: Topographic maps of the model fit for the best fit model of Euclidean distance from the centre of target object in the low-variance (left panels) and high-variance conditions (right panels). Time (in milliseconds) is represented on the x-axis. Pitch is on the y-axis. Pitch is centred around 0, the category boundary. Positive values correspond to the high tone, negative values to the mid tone. Distribution peaks were at 3 and -3. Distance of fixations from the target object is plotted on the z-axis and is colour coded. Higher values (warmer colours) indicate greater distance; lower values (cooler colours) indicate a smaller distance. The key in the top-left corner shows the corresponding distance (in pixels) and z-limits. Random effects are excluded from this plot.

a lexical tone contrast. Participants saw pictures of common objects and heard minimal word pairs, differing only in lexical tone. The tones were based on two Cantonese level tones, which are distinguished by pitch height. Auditory stimuli were sampled from pitch continua corresponding to the words. Stimuli were sampled according to a bimodal distribution. The critical manipulation was the statistical variance of the distribution, i.e. the amount of acoustic variability within the critical contrastive dimension, pitch. Participants heard either a *high-variance* or a *low-variance* distribution. Based on literature investigating effects of variance in native language processing [2, 4], we predicted that acquisition of the pitch cue would be better in the low-variance condition. GAMM models of eye movements showed that when variance was low, participants learned to use the pitch cue better over the course of the experiment. The Euclidean distance between fixations and the centre of the target picture reduced over the course of the experiment and, by the end of the experiment, was lower in the low-variance condition, compared to the high-variance condition.

The results demonstrate an interesting interplay between auditory and visual cues. As variance in the auditory stimuli increased, variability in the eye movements also increased. The fixations became further from the target, probably due to scanning the images and fixating the competitor. This suggests that as the reliability of the auditory cues decreased, participants increased their search for visual cues, looking around more for verification from the image stimuli. This was despite the fact that the image stimuli were kept constant between conditions. This suggests that as perceivers, we adapt to the current input, including cross-modal adjustments to make optimal use of multisensory input.

This result also provides new evidence that within-category acoustic variance shapes nonnative acoustic cue acquisition. Previous studies have shown that acoustic variance affects native speech perception, with increased variance leading to increased perceptual uncertainty [2, 4]. The present results show that the same mechanism can also help shape acquisition of a new acoustic dimension not present as a lexical contrast in the native language.

Cue variance has been investigated previously in native Japanese listeners' learning of the English /l/-/r/ contrast [40]. Many native Japanese listeners have trouble attending to the third formant (F3) cue - which native English listeners tend to use to distinguish /l/ and /r/ - and rely instead on the less reliable second formant (F2). Using a video game as training over several days, Lim and Holt found that by presenting stimuli with high variability in the F2 cue dimension and low variability in the F3 dimension, participants' categorisation accuracy significantly increased and cue weighting shifted towards F3. While this study was very interesting and informative, it differs from the present study in several respects. Firstly, the present study directly compared the effects of high- vs. low-variance; the Lim and Holt study compared the effects of training to a control condition that did not involve exposure to English words. Secondly, the participants in the Lim and Holt study were proficient English speakers. They had been studying English for at least 12 years and had lived in an English-speaking environment for up to 2.5 years. The present study investigated acquisition of a new cue dimension, not encountered before in a lexical contrast. Participants were not only improving an already partially acquired contrast, but instead experiencing for the first time both the language and the tonal contrast. Thirdly, the Lim and Holt study used flat distributions - four steps of F2 and two steps of F3 - presented at equal frequency, whereas the present study used ap-

proximately Gaussian distributions. Therefore the present study makes an important contribution by directly testing effects of distributional variance in a Gaussian distribution on acquisition of a new acoustic dimension.

While several recent studies have emphasised the facilitative effect of variability on learning [41, 17, 16, 18, 42, 13, 14], it is important to distinguish between within-category acoustic variance in the critical dimension and variability in non-contrastive dimensions. Variability can lower cue weighting. If that variability is in a contrastive dimension, it will hinder discrimination.

6. Acknowledgements

7. References

- [1] P. D. Allopenna, J. S. Magnuson, and M. K. Tanenhaus, "Tracking the time course of spoken word recognition using eye movements: Evidence for continuous mapping models," *Journal of Memory and Language*, vol. 38, no. 4, pp. 419–439, 1998.
- [2] M. Clayards, M. K. Tanenhaus, R. N. Aslin, and R. A. Jacobs, "Perception of speech reflects optimal use of probabilistic speech cues," *Cognition*, vol. 108, no. 3, pp. 804–809, 2008.
- [3] J. S. Nixon, J. van Rij, P. Mok, R. H. Baayen, and Y. Chen, "Eye movements reflect acoustic cue informativity and statistical noise," in *ExLing 2015: Proceedings of the International Conference of Experimental Linguistics*, A. Botonis, Ed., 2015, pp. 54–57.
- [4] —, "The temporal dynamics of perceptual uncertainty: eye movement evidence from Cantonese segment and tone perception," *Journal of Memory and Language*, vol. 90, pp. 103–125, 2016.
- [5] C. T. Best, "A direct realist view of cross-language speech perception," in *Speech perception and linguistic experience: Issues in cross-language research*, W. Strange, Ed. Timonium, MD: York Press, 1995, pp. 171–204.
- [6] J. E. Flege, "Second language speech learning: Theory, findings, and problems," *Speech perception and linguistic experience: Issues in cross-language research*, pp. 233–277, 1995.
- [7] J. Maye and L. Gerken, *Learning phonemes without minimal pairs*. Proceedings of the 24th Annual Boston University Conference on Language Development, 2000.
- [8] J. Maye, D. Weiss, and R. Aslin, "Statistical phonetic learning in infants: Facilitation and feature generalization," *Developmental Science*, vol. 11, no. 1, 2008.
- [9] J. Maye, J. F. Werker, and L. Gerken, "Infant sensitivity to distributional information can affect phonetic discrimination," *Cognition*, vol. 82, no. 3, 2002.
- [10] K. Wanrooij, P. Boersma, and T. L. van Zuijlen, "Fast phonetic learning occurs already in 2-to-3-month old infants: an ERP study," *Frontiers in Psychology*, vol. 5., 2014.
- [11] P. Escudero, T. Benders, and K. Wanrooij, "Enhanced bimodal distributions facilitate the learning of second language vowels," *The Journal of the Acoustical Society of America*, vol. 130, no. 4, 2011.
- [12] K. Wanrooij, P. Escudero, and M. E. Raijmakers, "What do listeners learn from exposure to a vowel distribution? An analysis of listening strategies in distributional learning," *Journal of Phonetics*, vol. 41, no. 5, pp. 307–319, 2013.
- [13] G. C. Rost and B. McMurray, "Speaker variability augments phonological processing in early word learning," *Developmental Science*, vol. 12, no. 2, pp. 339–349, 2009.
- [14] —, "Finding the signal by adding noise: The role of non-contrastive phonetic variability in early word learning," *Infancy*, vol. 15, no. 6, pp. 608–635, 2010.

- [15] A. R. Bradlow, D. B. Pisoni, R. Akahane-Yamada, and Y. Tohkura, "Training Japanese listeners to identify English /r/ and /l/: Iv. some effects of perceptual learning on speech production," *The Journal of the Acoustical Society of America*, vol. 101, no. 4, pp. 2299–2310, 1997.
- [16] J. S. Logan, S. E. Lively, and D. B. Pisoni, "Training Japanese listeners to identify English/r/and/l: A first report," *The Journal of the Acoustical Society of America*, vol. 89, no. 2, 1991.
- [17] S. E. Lively, J. S. Logan, and D. B. Pisoni, "Training Japanese listeners to identify English/r/and/l: II: The role of phonetic environment and talker variability in learning new perceptual categories," *The Journal of the Acoustical Society of America*, no. 1242., 1993.
- [18] R. A. Yamada, "Effect of extended training on /r/ and /l/ identification by native speakers of Japanese," *The Journal of the Acoustical Society of America*, vol. 93, no. 4, pp. 2391–2391, 1993.
- [19] P. Boersma and D. Weenink, "Praat (version 5.5)," 2012.
- [20] X. Lin and D. Zhang, "Inference in generalized additive mixed models using smoothing splines," *Journal of the Royal Statistical Society*, vol. 61, no. 7, p. 381, 1999.
- [21] S. Wood, *Generalized additive models: an introduction with R*. CRC press, 2006.
- [22] —, "Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models," *Journal of the Royal Statistical Society (B)*, vol. 73, no. 1, pp. 3–36, 2011.
- [23] R Core Team, *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2016. [Online]. Available: <https://www.R-project.org/>
- [24] C. de Cat, E. Klepousniotou, and H. Baayen, "Representational deficit or processing effect? An electrophysiological study of noun-noun compound processing by very advanced L2 speakers of English," *Frontiers in Psychology*, vol. 6, p. 77, 2015.
- [25] —, "Electrophysiological correlates of noun-noun compound processing by non-native speakers of English," *ComAComA 2014*, pp. 41–52, 2014.
- [26] J. S. Nixon, J. van Rij, X. Q. Li, and Y. Chen, "Cross-category phonological effects on ERP amplitude demonstrate context-specific processing during reading aloud," in *ExLing 2015: Proceedings of the International Conference of Experimental Linguistics*, A. Botonis, Ed., 2015, pp. 50–53.
- [27] A. Tremblay and A. Newman, "Modelling non-linear relationships in ERP data using mixed-effects Regression with R examples," *Psychophysiology*, vol. TBA, pp. 1–16, 2014.
- [28] L. B. Feldman, P. Milin, K. W. Cho, F. Moscoso del Prado Martín, and P. A. O'Connor, "Must analysis of meaning follow analysis of form? a time course analysis," *Frontiers in human neuroscience*, vol. 9, 2015.
- [29] Pham, Hien, and H. R. Baayen, "Semantic relations and compound transparency: A regression study in CARIN theory," *Psihologija*, vol. 46, no. 4, pp. 455–478, 2013.
- [30] D. Arnold, P. Wagner, and H. Baayen, "Using generalized additive models and random forests to model German prosodic prominence," *Proceedings of Interspeech 2013*, pp. 272–276, 2013.
- [31] F. Tomaschek, M. Wieling, D. Arnold, and R. H. Baayen, "Word frequency, vowel length and vowel quality in speech production: an EMA study of the importance of experience," in *Interspeech*, 2013, pp. 1302–1306.
- [32] S. Kawase, *Examination of the role of native speech rhythm in non-native speech production and its perception (Doctoral dissertation, University of Western Sydney)*. Australia, 2017.
- [33] M. Tamminga, C. Ahern, and A. Ecay, "Generalized additive mixed models for intraspeaker variation," *Linguistics Vanguard*, vol. 2, no. s1, 2016.
- [34] M. Wieling, S. Montemagni, J. Nerbonne, and R. H. Baayen, "Lexical differences between tuscan dialects and standard italian: Accounting for geographic and sociodemographic variation using generalized additive mixed modeling," *Language*, vol. 90, no. 3, pp. 669–692, 2014.
- [35] J. van Rij, B. Hollebrandse, and P. Hendriks, "Children's eye gaze reveals their use of discourse context in object pronoun resolution," in *Empirical perspectives on anaphora resolution: Information structural evidence in the race for salience*, A. Holler, C. Goeb, and K. Suckow, Eds. Berlin, Walter de Gruyter, 2016.
- [36] K. Loo, J. van Rij, J. Järvikivi, and H. Baayen, "Individual differences in pupil dilation during naming task," 2016.
- [37] H. Baayen, S. Vasisht, R. Kliegl, and D. Bates, "The cave of shadows: Addressing the human factor with generalized additive mixed models," *Journal of Memory and Language*, vol. 94, pp. 206–234, 2017.
- [38] J. van Rij, M. Wieling, R. H. Baayen, and H. van Rijn, "itsadug: Interpreting time series and autocorrelated data using gams," 2016, r package version 2.2.
- [39] D. Dahan, M. K. Tanenhaus, and A. P. Salverda, "The influence of visual processing on phonetically-driven saccades in the 'visual world' paradigm," in *Eye movements: A window on mind and brain*, 2007, pp. 471–486.
- [40] S.-j. Lim and L. L. Holt, "Learning foreign sounds in an alien world: Videogame training improves non-native speech categorization," *Cognitive science*, vol. 35, no. 7, pp. 1390–1405, 2011.
- [41] A. R. Bradlow and T. Bent, "Perceptual adaptation to non-native speech," *Cognition*, vol. 106, no. 2, pp. 707–729, 2008.
- [42] J. S. Allen, J. L. Miller, and D. DeSteno, "Individual talker differences in voice-onset-time," *Journal of the Acoustical Society of America*, 2003.