

Attention and salience in lexically-guided perceptual learning

by

Michael McAuliffe

B.A., University of Washington, 2009

A THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS FOR THE DEGREE OF

Doctor of Philosophy

in

THE FACULTY OF ARTS

(Linguistics)

The University Of British Columbia

(Vancouver)

April 2015

© Michael McAuliffe, 2015

Preface

At University of British Columbia (UBC), a preface may be required. Be sure to check the Graduate and Postdoctoral Studies (GPS) guidelines as they may have specific content to be included.

Contents

Preface	i
List of Tables	iv
List of Figures	v
Glossary	vi
Acknowledgments	vii
1 Introduction	1
1.1 Perceptual learning	2
1.2 Linguistic expectations and perceptual learning	4
1.2.1 Lexical bias	4
1.2.2 Semantic predictability	5
1.3 Attention and perceptual learning	5
1.4 Category typicality and perceptual learning	7
1.5 Current contribution	8
2 Lexical decision	9
2.1 Motivation	9
2.2 Experiment 1	9
2.2.1 Methodology	10
2.2.2 Results	13
2.2.3 Discussion	15
2.3 Experiment 2	15
2.3.1 Methodology	15
2.3.2 Results	16
2.4 Grouped results across experiments	17
2.5 General discussion	18
3 Cross-modal word identification	20
3.1 Motivation	20
3.2 Methodology	21
3.2.1 Participants	21
3.2.2 Materials	21
3.2.3 Pretest	21
3.2.4 Procedure	22
3.3 Results	22
3.3.1 Exposure	22
3.3.2 Categorization	22

3.4	Discussion	22
4	Conclusions	24
4.1	Attentional control of perceptual learning	24
4.2	Specificity versus generalization in perceptual learning	24
4.3	Effect of increased linguistic expectations	25
4.4	Category atypicality	25
4.5	Conclusion	25
	Bibliography	26

List of Tables

2.1	Mean and standard deviations for frequencies (log frequency per million words in SUBTLEXus) and number of syllables of each item type	10
2.2	Frequencies (log frequency per million words in SUBTLEXus) of words used in categorization continua .	10

List of Figures

2.1	Proportion of word-responses for /s/-initial exposure words. Solid lines represent Experiment 1 selection criteria (50% word-response rate) and dashed lines represent Experiment 2 selection criteria (30% word-response rate). Dots are averaged word-response across subjects, and the blue line is a binomial model constructed from the responses.	11
2.2	Proportion of word-responses for /s/-final exposure words. Solid lines represent Experiment 1 selection criteria (50% word-response rate) and dashed lines represent Experiment 2 selection criteria (30% word-response rate). Dots are averaged word-response across subjects, and the blue line is a binomial model constructed from the responses	12
2.3	Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 1. In the S-Final condition, participants in the Attention condition showed a larger perceptual learning effect than those in the No Attention condition. In the S-Initial condition, there were no differences in perceptual learning between the Attention conditions. Error bars represent 95% confidence intervals.	14
2.4	Correlation of crossover point in categorization with the proportion of word responses to critical items containing an ambiguous /s/ token.	14
2.5	Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 2. Participants showed no significant differences across conditions. Error bars represent 95% confidence intervals.	16
2.6	Correlation of crossover point in categorization with the proportion of word responses to critical items containing an ambiguous /s/ token in Experiment 2.	17
2.7	Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 1 and Experiment 2. Error bars represent 95% confidence intervals.	18
3.1	Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 3. Error bars represent 95% confidence intervals.	23

Glossary

This glossary uses the handy `acroynym` package to automatically maintain the glossary. It uses the package's `printonlyused` option to include only those acronyms explicitly referenced in the \LaTeX source.

GPS Graduate and Postdoctoral Studies

Acknowledgments

Thank those people who helped you.

Molly!

Jobie!

Jamie!

Michelle!

Don't forget your parents or loved ones.

Mom!

Dad!

Laura!

You may wish to acknowledge your funding sources.

Chapter 1

Introduction

Listeners of a language are faced with a large degree of phonetic variability when interacting with their fellow language users. Speakers can have different sizes, different genders, and different backgrounds that make speech sound categories, at first blush, overlapping in distribution and hard to separate in acoustic dimensions. In addition to properties of the speaker varying, a listener's attention or goals in an interaction can vary, such as paying attention to a non-native speaker or being distracted by planning upcoming utterances or by another task entirely. Despite variability on the part of both the listener and the speaker, listeners can interpret disparate and variable productions as belonging to a single word type or sound category, a phenomenon referred to as perceptual constancy in categorization studies [31, 55] and as recognition equivalence in word recognition tasks [57]. One of the processes for achieving this constancy is perceptual learning, whereby perceivers update context-dependent categories. In the speech perception literature, perceptual learning or perceptual adaptation refers to the updating of distributions corresponding to a sound category, such as /s/ or /ʃ/, for a particular speaker after exposure to speech with a modification to the sound category's distribution [40]. Exposure to a speaker affects a listener's perceptual system for that speaker, even after only a few tokens [30, 60] or within the span of a single utterance [32]. Perceptual learning is not limited to single speakers, and exposure to multiple speakers sharing a non-native accent increases the intelligibility of new speakers with that accent [8].

However, perceptual learning, particularly in speech perception, relies on the linguistic system, which interacts with listener attention and signal properties. Cutler et al. [15] puts forth two attentional sets that are routinely used in speech perception tasks. The first is a comprehension-oriented attentional set, which is the most similar to normal language use. In this attentional set, listeners are focused on comprehending the intended message of the speech. The second attentional set is a perception-oriented attentional set, where a listener is focused more on the low-level, signal properties of the speech rather than the message. This dissertation investigates the interaction between attention and salience on the attentional sets used during perceptual learning when the /s/ category of a speaker is modified to be more /ʃ/-like. The salience manipulated in this dissertation is not that inherent the physical stimulus per se, but rather linguistic salience. The first way of manipulating salience relies on differences in lexical and sentential expectations. The second way relies on differences in degree of the modified category's typicality, with the modified /s/ category more /s/-like or more /ʃ/-like. This dissertation largely adopts a hierarchical predictive coding framework [11] for perceptual learning, where errors between expectations and perceptions are propagated hierarchically to modify future expectations. Incorporating findings from the psycholinguistic literature on attention, I propose a modification to how attention is handled in a predictive coding model. The key hypothesis being tested is that the higher errors are propagated, the more perceptual learning occurs and the higher the rate of generalization to novel forms will be. Error propagation is proposed to be blocked by attention, and expectations will be updated at the level that attention is oriented. For instance, in a more perception-oriented attentional set, the updating of expectations will be restricted to more stimuli-specific cognitive representations, preventing generalization. In a more comprehension-oriented attentional set, errors will be propagated higher levels of representation, allowing for more expectation updating, leading to generalization.

The adoption of perception- and comprehension-oriented attentional sets will be induced through selective attention and manipulations of linguistic salience. Attention will be manipulated through task demands focusing on word recognition and additional instructions to some listeners to pay attention to the specific sound category to be learned. I predict that focusing attention on the sound category will result in less perceptual learning through the adoption of a more perception-

oriented attentional set. Two linguistic expectations will be manipulated that have been found to affect the categorization and discrimination of speech sounds embedded in words, namely lexical bias and semantic predictability. According to the hypothesis above, increasing linguistic expectations for the words that contain the target sounds will result in greater perceptual learning through the maintenance of a more comprehension-oriented attentional set. Finally, the acoustic tokens that listeners are exposed to are manipulated to expose listeners to /s/ categories that are closer or farther from a typical /s/ category. I predict that exposure only to tokens farther from the typical category will result in less perceptual learning than exposure to tokens closer to the typical category, with the atypicality inducing a more perception-oriented attentional set. Selective attention directed to sound category should only show effects in conditions where participants would otherwise be biased towards using a more comprehension-oriented attentional set. For instance, lower lexical bias positions, such as the beginnings of words, are salient positions in formal linguistic theory [3], and productions farther from the typical production are more likely to be noticed, such as in phoneme restoration tasks where restorations are more likely when the replacing sound is similar to the restored sound [50]. In both of these cases, increasing external attention on the /s/ category will likely not produce an effect beyond what attention the signal itself draws.

The structure of the thesis is as follows. This chapter will review the recent literature on perceptual learning in speech perception (Section 1.1), as well as literature on linguistic expectations (Section 1.2), attention (Section 1.3), and category typicality (Section 1.4) as they relate to perceptual learning and to the three experiments of this dissertation. Chapter 2 will detail two experiments using a lexically-guided perceptual learning paradigm, each with different conditions for levels of lexical bias and attention. The two experiments differ in the acoustic properties of the exposure tokens, with the first experiment using a slightly atypical /s/ category that is halfway between /s/ and /ʃ/, and the second experiment using a more atypical /s/ category that is more /ʃ/-like than /s/-like. Chapter 3 details an experiment using a novel perceptual learning paradigm that manipulates an additional linguistic factor, namely semantic predictability, to increase the linguistic expectations during exposure. Finally, Chapter 4 will summarize the results and place them in a theoretical framework. The perceptual learning literature has generally used consistent processing conditions to elicit perceptual learning effects, and a goal of this dissertation is to examine the robustness and degree of perceptual learning across different processing conditions.

1.1 Perceptual learning

Perceptual learning is a well established phenomenon in the psychology and psychophysics literature. Training can improve a perceiver's ability to discriminate in many disparate modalities, such as visual acuity, somatosensory spatial resolution, weight estimation, and discrimination of hue and acoustic pitch [21, for review]. In this literature, perceptual learning is the improvement of a perceiver to judge the physical characteristics of objects in the world through training that assumes attention on the task, but doesn't require reinforcement, correction or reward. This definition of perceptual learning corresponds more to what is termed "selective adaptation" in the speech perception literature rather than what is termed "perceptual learning", "perceptual adaptation" or "perceptual recalibration." In speech perception, selective adaptation is the phenomenon where listeners that are exposed repeatedly to a narrow distribution for a sound category, narrow their own perceptual category, resulting in a change in variance of the category, not of the mean of the category (along some acoustic-phonetic dimension) [16, 51, 60]. Perceptual learning or recalibration in the speech perception literature is a more broad updating of perceptual categories, either in mean or variance [40, 60].

Norris et al. [40] began the recent set of investigations into lexically-guided perceptual learning in speech. Norris et al. [40] exposed one group of Dutch listeners to a fricative halfway between /s/ and /f/ at the ends of words like *olif* "olive" and *radijs* "radish", while exposing another group to the ambiguous fricative at the ends of nonwords, like *blif* and *blis*. Following exposure, both groups of listeners were tested on their categorization on a fricative continuum from 100% /s/ to 100% /f/. Listeners exposed to the ambiguous fricative at the end of words shifted their categorization behaviour, while those exposed to them at the end of nonwords did not. The exposure using words was further differentiated by the bias introduced by the words. Half the tokens ending in the ambiguous fricative formed a word if the fricative was interpreted as /s/ but not if it was interpreted as /f/, and the others were the reverse. Listeners exposed only to the /s/-biased tokens categorized more of the /f/-/s/ continuum as /s/, and listeners exposed to /f/-biased tokens categorized more of the continuum as /f/. The ambiguous fricative was associated with either /s/ or /f/ dependent on the bias of the word, which led to an expanded category for that fricative at the expense of the other category. These results crucially show

that perceptual categories in speech are malleable to new exposure, and that the linguistic system of the listener facilitates generalization to that category in new forms and contexts.

In addition to lexically-guided perceptual learning, unambiguous visual cues to sound identity can cause perceptual learning as well; this is referred to as perceptual recalibration in that literature. In Bertelson et al. [4], an auditory continuum from /aba/ to /ada/ was synthesized and paired with a video of a speaker producing /aba/ and a video of /ada/. Participants first completed a pretest that identified the maximally ambiguous step of the /aba/-/ada/ auditory continuum. In eight blocks, participants were randomly exposed to the ambiguous auditory token paired with video for /aba/ or with the video for /ada/. Following each block, they completed a short categorization test. Participants showed perceptual learning effects, such that they were more likely to respond with /aba/ if they had been exposed to video of /aba/ paired with the ambiguous token in the preceding block, and likewise for /ada/.

van Linden and Vroomen [58] compared the perceptual recalibration effects from the visually-guided perceptual learning paradigm [4] to the standard lexically-guided perceptual learning paradigm [40]. Lipread recalibration and perceptual learning effects had comparable size, lasted equally as long, were enhanced when presented with a contrasting sound, and were both unaffected by periods of silence between exposure and categorization. The effects did not last through prolonged testing in this study, unlike in other studies [18, 27]. In those studies, lexically-guided perceptual learning effects persisted through intermediate tasks, as long as the task did not involve any contradictory evidence of the trait learned [27], and also persisted across 12 hours [18].

Perceptual learning in speech perception can be captured well in terms of Bayesian belief updating [26] or as part of a predictive coding model of the brain [11]. In Bayesian belief updating, the model categorizes the incoming stimuli based on multimodal cues, and then updates the distribution to reflect that categorization. This updated conditional distribution is then used for future categorizations in an iterative process. Kleinschmidt and Jaeger [26] model the results of the behavioural study in Vroomen et al. [60] in a Bayesian framework, with models fit to each participant capturing the perceptual recalibration and selective adaptation shown by the participants over the course of the experiment. A similar, but more broad, framework is that of predictive coding [11]. This framework uses a hierarchical generative model that aims to minimize prediction error between bottom-up sensory inputs and top-down expectations. Mismatches between the top-down expectations and the bottom-up signals generate error signals that are used to modify future expectations. Perceptual learning then is the result of modifying expectations to match learned input.

Perceptual learning in the psychophysics literature has shown a large degree of exposure-specificity, where observers only show learning effects on the same or very similar stimuli as those they were trained on. As such, perceptual learning has been argued to reside or affect the early sensory pathways, where stimuli are represented with the greatest detail [22]. Perceptual learning in speech perception has also shown a large degree of exposure-specificity, where participants do not generalize cues across speech sounds [49] or across speakers unless the sounds are similar across exposure and testing [17, 27, 28, 47]. On the other hand, lexically-guided perceptual learning in speech has shown a greater degree of generalization than would be expected from a purely psychophysical standpoint. The testing stimuli are generally quite different from the exposure stimuli, with participants exposed to multisyllabic words ending in an ambiguous sound and tested on monosyllabic words [48] and nonwords [27, 40], though exposure-specificity is found when exposure and testing use different positional allophones [38].

Why is lexically-guided perceptual learning more context-general? The experiments performed in this dissertation will provide evidence that this context-generality is the result of a listener's attentional set, which can be influenced by linguistic, attentional and signal properties. A more comprehension-oriented attentional set, where a listener's goal is to understand the meaning of speech, promotes generalization and leads to greater perceptual learning. A more perception-oriented attentional set, where a listener's goal is to perceive specific qualities of a signal, does not promote generalization. In terms of a Bayesian framework with error propagation, a more perception-oriented attentional set may keep error propagation more local, resulting in the exposure-specificity seen more in the psychophysics literature. A more comprehension-oriented attentional set would propagate errors farther upward in the hierarchy of expectations. In both cases, errors would propagate to where attention is focused, but larger updating associated with farther error propagation would lead to larger observed perceptual learning effects. These attentional sets will be explored in more detail in Section 1.3 following an examination of the linguistic factors that will be manipulated in the experiments in Chapters 2 and 3, namely lexical bias and semantic predictability; both of which are hypothesized to aid the listener in adopting comprehension-oriented attentional sets.

1.2 Linguistic expectations and perceptual learning

The two linguistic factors manipulated in this dissertation are lexical bias and semantic predictability. Chapter 2 presents two experiments using a standard lexically-guided perceptual learning paradigm, which uses lexical bias as the means to link an ambiguous sound to an unambiguous category. In Chapter 3, a novel, sententially-guided perceptual learning paradigm is used to manipulate the listener’s linguistic expectations in a different manner than manipulating lexical bias.

1.2.1 Lexical bias

Lexical bias is the primary way through which perceptual learning is induced in the experimental speech perception literature. Lexical bias, also known as the Ganong Effect, refers to the tendency for listeners to interpret a speaker’s (noncanonical) production as a particular, meaningful word rather than a nonsense word. For instance, given a continuum from a nonword like *dask* to word like *task* that differs only in the initial sound, listeners in general are more likely to interpret any step along the continuum as the word endpoint rather than the nonword endpoint [20]. This bias is exploited in perceptual learning studies to allow for noncanonical, ambiguous productions of a sound to be linked to pre-existing sound categories. Given that ambiguous productions must be associated with a word to induce lexically-guided perceptual learning [40], differing degrees of lexical bias could lead to differing degrees of perceptual learning, a prediction which will be tested in Chapter 2. The stronger the lexical bias, the stronger the link will be between the ambiguous sound and the sound category, as mediated by the word.

Lexical bias has been found to vary in strength according to several factors. First, the length of the word in syllables has a large effect on lexical bias, with longer words showing stronger lexical bias than shorter words [46]. Continua formed using trisyllabic words, such as *establish* and *malpractice*, were found to show consistently large lexical bias effects than monosyllabic words, such as *kiss* and *fish*. Pitt and Samuel [46] also found that lexical bias from trisyllabic words was robust across experimental conditions, but lexical bias from monosyllabic words was more fragile and condition dependent. They argue that these effects arise from both the greater bottom-up information present in longer words and the greater lexical competition for shorter words.

Pitt and Szostak [43] used a lexical decision task with a continuum of fricatives from /s/ to /ʃ/ embedded in words differing in the position of a sibilant. They found that ambiguous fricatives earlier in the word, such as *serenade* or *chandelier*, lead to greater nonword responses than the same ambiguous fricatives embedded later in a word, such as *establish* or *embarass*. In the experiments and meta-analysis of phoneme identification results presented in Pitt and Samuel [45], they found that, for monosyllabic word frames, token-final targets produce more robust lexical bias effects than token-initial targets. Between word length and position in the word, lexical bias appears to be strengthened over the course of the word. As a listener hears more evidence for a particular word, their expectations for hearing the rest of that word increase.

Lexical bias has been found to affect phoneme restoration tasks as well [50]. In this paradigm, listeners heard words with noise added to or replacing sounds and were asked to identify which they had heard for each trial. Lower sensitivity to noise addition versus replacement and increased bias for responding that noise was added is indicative of phoneme restoration, where listeners are perceiving sounds not actually present in the signal. Several factors are identified in the study as increasing the likelihood of the phoneme restoration effect. In the lexical domain, words are more likely than nonwords to have phoneme restorations, and this discrepancy strengthens when listeners are primed with a form without noise before the trial. More frequent words were also more likely to exhibit phoneme restoration effects, and longer words also showed greater phoneme restoration effects. Position of the sound in the word also influenced the decision, with non-initial positions showing greater phoneme restoration effects. The other influences on phoneme restoration discussed in this paper, namely the signal properties and sentential context will be discussed in subsequent sections.

In the stimuli used in Norris et al. [40] and most other lexically-guided perceptual learning experiments, lexical bias tends to be maximized by using multisyllabic words with the ambiguous sound at the end for exposure stimuli. Perceptual learning has also been found when the ambiguous stimuli is embedded earlier in the word, such as the onset of the final syllable [27, 29, 30] or the even the onset of the first syllable [10]. In light of the findings in Pitt and Szostak [43], we would expect less lexical biases the earlier in the word those ambiguous sounds were heard, and therefore, I hypothesize, lower endorsement rates and smaller perceptual learning effect sizes. This prediction will be explicitly tested in Experiment 1 in Chapter 2.

1.2.2 Semantic predictability

The second type of linguistic expectation manipulation used in this dissertation is known as semantic predictability [23]. Sentences are semantically predictable when they contain words prior to the final word that points almost definitively to the identity of that final word. For instance, the sentence fragment *The cow gave birth to the...* from Kalikow et al. [23] is almost guaranteed to be completed with the word *calf*. On the other hand, a fragment like *She is glad Jane called about the...* is far from having a guaranteed completion, other than having the category of noun.

Despite the temporal and spectral reduction found in high predictability contexts [12, 52], high predictability sentences are generally more intelligible. Sentences that form a semantically coherent whole have higher word identification rates across varying signal-to-noise ratios [23], which has been found across children and adults [19], and across native monolingual and early bilingual listeners, but not late bilingual listeners [37]. However, when words at the ends of predictive sentences are excised from their context, they tend to be less intelligible than words excised from non-predictive contexts [35]. Highly predictable sentences are more intelligible to native listeners in noise, even when signal enhancements are not made, though non-native listeners require both signal enhancements and high predictability together to see any benefit [7].

Two studies looking at similarities between lexical bias and semantic predictability found similar effects on phoneme categorization, though differing effects in reaction time [6, 13]. Connine [13] found evidence that semantic predictability operated similar to the “postperceptual” processes in Connine and Clifton [14]. The perceptual process was lexical bias towards one endpoint of a continuum, and the postperceptual process was a monetary benefit for responding with one of the end points of a continuum [14]. The difference between “perceptual” and “postperceptual” processes in those studies was primarily one of reaction time; both kinds of bias produced comparable shifts in categorization. Borsky et al. [6] attempted to replicate Connine [13] while removing potential confounds from the methodological design. In contrast to Connine [13], they used only one voicing continuum (from *goat* to *coat*), embedded the acoustic target in the middle of the sentence rather than at the end, and only presented one instance of each sentence to each participant rather than all continuum steps for each sentence to each participant. The reaction time profile in their study more closely aligned to the profile of lexical bias in Connine and Clifton [14], but again the shift in how the continuum was categorized aligned with the sentential context.

In phoneme restoration, higher semantic predictability has been found to bias listeners toward interpreting the stimuli as intact with noise rather than replaced with noise [50]. This increased bias towards interpreting the stimuli as an intact word was also coupled with an increase in sensitivity between the two types of stimuli, which Samuel [50] suggests is the result of a lower cognitive load in predictable contexts, and therefore greater phonetic encoding is available [see also 36].

The previous literature on semantic predictability has shown largely similar effects as lexical bias in how sounds are categorized and restored. From this, I hypothesize that increasing the expectations for a word through semantic predictability will increase the link between sound categories and words in a similar fashion as increasing lexical bias. Listeners that are exposed to an /s/ category that is more /ʃ/-like only in words that are highly predictable from context should show larger perceptual learning effects than listeners exposed to the same category only in words that are unpredictable from context. However, there may be an upper limit for listener expectations when both semantic predictability and lexical bias are high, as committing too much to a particular expectation could lead to garden path phenomena [34]. The effect of semantic predictability on perceptual learning will be explicitly tested in Chapter 3.

1.3 Attention and perceptual learning

Attention is a large topic of research in its own right, and this section only reviews literature that is directly relevant to perceptual learning. Attention has been found to have a role on perceptual learning in the psychophysics literature. For instance, Ahissar and Hochstein [1] found that in general, attending to global features for detection (i.e., discriminating different orientations of arrays of lines) does not make participants better at using local features for detection (i.e., detection of a singleton that differs in angle in the same arrays of lines), and vice versa. However, there was a small degree of local detection learned from global detection, with the singleton popping out from its field as highly salient.

Attentional sets are a widely used term in the attention literature. In the visual domain, attentional sets can refer to the strategies that the perceiver uses to perform a task. For instance, in a visual search task, colour, orientation, motion and size are the predominant strategies [61]. Two broad categories of attentional sets are generally used. Focused sets direct

attention to components of the sensory input, and diffuse sets direct attention to global properties of the sensory input. In the visual search literature, a focused attentional set is the feature search mode, which gives priority to a single feature, such as the colour of the target, and a diffuse attentional set is singleton detection mode, which gives priority to any salient features [2]. In the auditory streaming literature, two attentional sets have been identified as “selective listening”, where the perceiver attempts to hear the components of two streams, and “comprehensive listening”, where the perceiver tries to hear all components as a single stream [59]. Finally, in a lexical decision tasks, the diffuse attentional set is where primary attention is on detecting words from nonwords, and the focused attentional set is where instructions direct participants attention to a potentially misleading sound [43]. Attentional set selection is not necessarily optimal on the part of the perceiver, and it has been shown to be biased based on experience, with the amount of training performed influencing the length of time that perceivers will continue to use non-optimal sets after the task has changed [33].

Listeners can selectively attend to many aspects of a signal, broadly falling into the categories of perception and comprehension. For instance, listeners can attend to particular syllables or sounds in syllable or phoneme monitoring tasks [39, and others], and even particular linguistically relevant temporal positions [44]. However, even in these low-level, signal based tasks, lexical properties of the signal can exhibit some influence, if the stimuli are not monotonous enough to disengage comprehension [15]. Variation in speech in general seems to lead towards a more diffuse, comprehension-oriented attentional set, which is likely the default attentional set employed in everyday use of language, where the goal is firmly more comprehension than perception.

Attention has not been manipulated in previous work on perceptual learning in speech perception, but some work has been done on how individual differences in attention control can impact perceptual learning. Scharenborg et al. [54] presents a perceptual learning study of older Dutch listeners in the model of Norris et al. [40]. In addition to the exposure and test phases, these older listeners completed tests for hearing loss, selective attention and attention-switching control. They found no evidence that perceptual learning was influenced by listeners’ hearing loss or selective attention abilities, but they did find a significant relationship between a listener’s attention-switching control and their perceptual learning. Listeners with worse attention-switching control showed greater perceptual learning effects, which the authors ascribed to an increased reliance on lexical information. Older listeners had previously been shown to have smaller perceptual learning effects as compared to younger listeners, but the differences were most prominent directly following exposure [53]. Younger listeners initially had a larger perceptual learning effect in the first block of testing, but the effect lessened over the subsequent blocks. Older listeners showed smaller initial perceptual learning effects, but no such decay. Scharenborg and Janse [53] also found that performance in the lexical decision task significantly affected the perceptual learning in the testing phase.

Lexical bias is affected by attentional processing conditions. Pitt and Szostak [43] additionally investigated the role of attention in modulating lexical bias. When listeners were told that the speaker’s /s/ and /ʃ/ were ambiguous and to listen carefully to ensure correct responses, they were less tolerant of noncanonical productions across all positions in the word. That is, participants attending to the speaker’s sibilants were less likely to accept the modified production as a word than participants given no particular instructions about the sibilants. Given that the task listeners performed was a lexical decision task, the default attentional set for the task, termed “diffuse” by Pitt and Szostak [43], would have attention distributed across both acoustic-phonetic and lexical domains. Listeners given the instructions about the speaker’s /s/ productions had a “focused” attentional set, with more weighting on the acoustic-phonetic domain than the “diffuse” attentional set. Under higher cognitive load, such as performing a more difficult concurrent task, listeners show an increased lexical bias, as a result of weaker encoding of the auditory details [36]. These results suggest that detailed encoding requires attentional resources. In Mattys and Wiget [36], the primary task was a phoneme identification task, where a “focused” attentional set would likely be the default for participants, but a more “diffuse” attentional set, or one that weights lexical information more heavily, seems to be employed in the higher cognitive load conditions.

Attentional sets are thought to be constant across a task, but the fact that attention switching in older adults was a significant predictor of the size of perceptual learning effects [54] suggests that listeners are indeed switching sets through an experiment. The task is oriented toward comprehension, so the primary attentional set is likely to be a diffuse one relying more on lexical information than acoustic. Participants with worse attention-switching control would have adopted this attentional set for more of the exposure than those with better attention-switching control, and it is precisely those with worse attention-switching control that showed the larger perceptual learning effects. The ability to switch attention to a more perception-oriented set could have allowed those listeners prevent over-generalization, leading to a

smaller perceptual learning effect for participants with better attention-switching control.

The predictive coding framework [11] provides a gain-based attentional mechanism. In this view, attention causes greater weight to be attached error units, increasing their weight and their effect on expectations. However, as noted in a response to Clark [11] by Block and Siegel [5], this view of attention does not capture the full range of experimental results. For instance, in a texture segregation task, spatial attention to the periphery improves performance where spatial resolution is poor, but attention to central locations, where spatial resolution is high, actually harms performance [62]. This detrimental effect is an instance of missing the forest for the trees, with spatial resolution is increased too much in the central locations to perceive the larger texture. The attentional mechanism that I propose for the predictive coding framework is one where attention is not simply gain, but affects where errors must be resolved and expectations updated. The more attention is oriented towards initial perception, the less general perceptual learning will be observed as a result of exposure. According to the mechanism proposed in Clark [11], any increases in attention, perception-oriented or otherwise, should lead to greater perceptual learning.

In this dissertation, focusing a listener’s attention on the signal through explicit instructions is hypothesized to lead to smaller perceptual learning effects. In a focused perception-oriented attentional set, comprehension will be de-emphasized and expectations will be updated at a lower level, leading to a more exposure-specific learning that will not generalize to novel tokens as readily. The work on phoneme restoration suggests that higher predictability sentences are lower cognitive load than unpredictable sentences [50], which may lead to greater encoding for the sounds to be learned, leading to larger perceptual learning effects. However, if fewer attentional resources are required for comprehending the word, there could be more attention on perception in the predictable sentences, leading to a reduced perceptual effect.

1.4 Category typicality and perceptual learning

A primary finding across the perceptual learning literature is that learning effects are only found on testing items that are similar to the exposure items. However, a less studied question is what properties of the exposure items cause different degrees of perceptual learning.

Variability is a fundamental property of the speech signal, so sound categories must have some variance associated with them, and certain contexts can have increased degrees of variability. Kraljic et al. [29] exposed participants to ambiguous sibilants between /s/ and /ʃ/ in two different contexts. In one, the ambiguous sibilants were intervocalic, and in the other, they occurred as part of a /str/ cluster in English words. Participants exposed to the ambiguous sound intervocalically showed a perceptual learning effect, while those exposed to the sibilants in /str/ environments did not. The sibilant in /str/ often surfaces closer to [ʃ] in many varieties of English, due to coarticulatory effects from the other consonants in the cluster, but the coarticulatory effects for merging /s/ and /ʃ/ are much weaker in intervocalic position. They argue that the interpretation of the ambiguous sound is done in context of the surrounding sounds, and only when the pronunciation variant is unexplainable from context is the variant learned and attributed to the speaker Kraljic et al. [see also 30]. A more /ʃ/-like /s/ category is typical in the context of the /str/ clusters, but is atypical in intervocalic position.

Sumner [56] investigated whether differences in presentation order in a perceptual learning experiment led to different learning effects of the /b/-p/ category boundary for a native French speaker of English. The presentation order that showed the greatest perceptual learning effects was the one where tokens started out close to what a listener would expect for the categories (English-like voice onset time for /b/ and /p/) and shifted over the course of the experiment to what the speaker’s actual categories were (French-like voice onset time for /b/ and /p/), despite the fact that this presentation order is not anything like what a listener would normally encounter when interacting with a non-native speaker of English. The condition that mirrored the more normal course of non-native speaker pronunciation changes, starting as more French-like and ending as more English-like, did not produce significantly different behaviour than control participants. One explanation for these results are that a small difference between the listener’s expectations and the input provides more robust learning, and the future input uses the updated expectations for future input, allowing for greater shifts over time through smaller shifts per trial, which aligns with some conceptualizations for exemplar model dynamics. For instance, in the exemplar model proposed by Pierrehumbert [42], only input similar to the learned distribution is used for updating that distribution, and input too ambiguous given learned distributions is discarded.

Additionally, in the phoneme restoration literature, sounds that are similar to the replacement sound are more likely to be restored [50]. When the replacement noise is white noise, fricatives and stops are more likely to be restored than

vowels and liquids. Acoustic signals that better match expectations are thus less likely to be noticed as atypical.

In this dissertation, the degree of typicality of the modified category is manipulated. In one case, the /s/ category for the speaker is maximally ambiguous between /s/ and /ʃ/, but in the other, the category is more like /ʃ/ than like /s/. The maximally ambiguous category is less likely to be noticed as an atypical production than the more /ʃ/-like /s/ category. I hypothesize that the more /ʃ/-like category will shift listeners' attentional sets to be more perception-oriented due to their greater atypicality, which will lead to lower perceptual learning. As such, perceptual learning effects are predicted to be greater when the modified /s/ category is more like the expected /s/ category than the expected /ʃ/ category.

1.5 Current contribution

Perceptual learning in speech perception generalizes to new forms and contexts far more than would be expected from a purely psychophysical perspective [22, 40]. This dissertation examines how different conditions for exposure affect generalization to subsequent categorization, leading to larger or smaller observed perceptual learning effects. The primary hypothesis being tested is that the more a listener is focused on comprehension rather than perception, the more the error between expectations and the actual signal will propagate and the larger the perceptual learning effect will be. The specific predictions being tested are that increased linguistic expectation, either through lexical bias or semantic predictability, will increase the reliance on comprehension-oriented attentional sets over perception-oriented ones. Directing selective attention to the specific sound category will shift the attentional set of participants to a more perception-oriented one, decreasing overall perceptual learning.

In this dissertation, I induce manipulations of lexical bias and selective attention in Experiments 1 and 2 (Chapter 2), and manipulations of sentence predictability and selective attention in Experiment 3 (Chapter 3). Lexical bias has been shown to affect phoneme categorization tasks [20], and can be manipulated by position of the ambiguous sound in the word and attention [43]. Sentence predictability has likewise been found to affect phoneme categorization and phoneme restoration tasks similar to lexical bias [6, 50], and can be manipulated by the preceding words in the sentence [23]. In all experiments, attention will be either comprehension-oriented, with a focus on lexical recognition or identity, or perception-oriented, with a focus on the /s/ category that is atypical.

Chapter 2

Lexical decision

2.1 Motivation

The experiments in this chapter implement a standard lexically-guided perceptual learning experiment with manipulations to lexical bias and selective attention. In the perceptual learning literature, word endorsement rates during exposure are reported as high, such as no lower than 85% in some studies [48]. However, word endorsement rate of words containing ambiguous sounds varies critically as a result of the position of that sound, and thus the lexical bias present when a listener encounters it, and whether a listener's attention is directed to that sound [43]. Experiment 1 contains manipulations to these two factors across participants.

Additionally, studies have reported greater perceptual learning when ambiguous stimuli are closer to the distribution expected by a listener than when the ambiguous stimuli are farther away from expected distributions [56]. Words containing stimuli farther away from the target production are in general less likely to be endorsed as words, but similar effects of attention were found across word position [43]. Experiment 2 contains the same manipulations to attention and lexical bias as Experiment 1, but with ambiguous stimuli farther from the target production than those used in Experiment 1.

The hypothesis tested in this dissertation predicts that greater perceptual learning effects should be observed when no attention is directed to the ambiguous sounds and when lexical bias is maximized. In that instance, participants should use a comprehension-oriented attentional set. If selective attention is directed to the ambiguous sounds, a more perception-oriented attentional set should be adopted. Likewise, if the ambiguous sound is in a linguistically salient position with little to no lexical bias, a listener's attention should be drawn to the ambiguous sound, causing an adoption of a more perception-oriented attentional set. Finally, if the ambiguous sounds are more atypical, they should be more salient to listeners regardless of lexical bias, leading to a more perceptual oriented attentional set. Regardless of the cause, adopting a perception-oriented attentional set is predicted to inhibit a generalized perceptual learning effect.

2.2 Experiment 1

In these two experiments, listeners will be exposed to ambiguous productions of words containing a single instance of /s/, where the /s/ has been modified to sound more like /ʃ/ in a lexical decision task. In one group, the S-Initial group, the critical words will have an /s/ in the onset of the first syllable, like in *cement*, with no /ʃ/ neighbour, like *shement*. In the other group, the S-Final group, the critical words will have an /s/ in the onset of the final syllable, like in *tassel*, with no /ʃ/ neighbour like *tashel*. In addition, half of each group will be given instructions that the speaker has an ambiguous /s/ and to listen carefully, following Pitt and Szostak [43].

Given the difference in word response rates depending on position in the word, we would predict that listeners exposed to ambiguous sounds earlier in words would be less likely to accept these productions as words as compared to listeners exposed to ambiguous sounds later in words. In addition, given the reliance of perceptual learning on lexical scaffolding, this lower acceptance rate for the former group should lead to a smaller perceptual learning effect as compared to the latter group.

Table 2.1: Mean and standard deviations for frequencies (log frequency per million words in SUBTLEXus) and number of syllables of each item type

Item type	Frequency	Number of syllables
Filler words	1.81 (1.05)	2.4 (0.55)
/s/-initial	1.69 (0.85)	2.4 (0.59)
/s/-final	1.75 (1.11)	2.3 (0.47)
/ʃ/-initial	2.01 (1.17)	2.3 (0.48)
/ʃ/-final	1.60 (1.12)	2.4 (0.69)

Table 2.2: Frequencies (log frequency per million words in SUBTLEXus) of words used in categorization continua

Continuum	/s/-word frequency	/ʃ/-word frequency
sack-shack	1.11	0.75
sigh-shy	0.53	1.26
sin-shin	1.20	0.48
sock-shock	0.95	1.46

2.2.1 Methodology

Participants

A total of 173 participants from the UBC population completed the experiment and were compensated with either \$10 CAD or course credit. Of these, 77 were non-native speakers of English and two native speakers of English reported speech or hearing disorders, and their results were not included in the analyses, resulting in data from 94 participants. Twenty additional native English speakers participated in a pretest to determine the most ambiguous sounds. Twenty five other native speakers of English participated for course credit in a control experiment.

Materials

One hundred and twenty English words and 100 nonwords that were phonologically legal in English were used as exposure materials. The set of words consisted of 40 critical items, 20 control items and 60 filler words. Half of the critical items had an /s/ in the onset of the first syllable and half had an /s/ in the onset of the final syllable. All critical tokens formed nonwords if their /s/ was replaced with /ʃ/. Half the control items had an /ʃ/ in the onset of the first syllable and half had an /ʃ/ in the onset of the final syllable. Each critical item and control item contained just the one sibilant, with no other /s z ʃ ʒ ʧ ʤ/. Filler words and nonwords did not contain any sibilants. Frequencies and number of syllables across item types are in Table 2.1

Four monosyllabic minimal pairs of voiceless sibilants were selected as test items for categorization (*sack-shack*, *sigh-shy*, *sin-shin*, and *sock-shock*). Two of the pairs had a higher log frequency per million words (LFPM) from SUBTLEXus [9] for the /s/ word, and two had higher LFPM for the /ʃ/ word, as shown in Table 2.2.

All words and nonwords were recorded by a male Vancouver English speaker in quiet room. Critical words for the exposure phase were recorded in pairs, once normally and once with the sibilant swapped forming a nonword. The speaker was instructed to produce both forms with comparable speech rate, speech style and prosody.

For each critical item, the word and nonword versions were morphed together in an 11-step continuum (0%-100% of the nonword /ʃ/ recording, in steps of 10%) using STRAIGHT [24] in Matlab. Prior to morphing, the word and nonword versions were time aligned based on acoustic landmarks, like stop bursts, onset of F2, nasalization or frication, etc. All control items and filler words were processed and resynthesized by STRAIGHT to ensure a consistent quality across stimulus items.

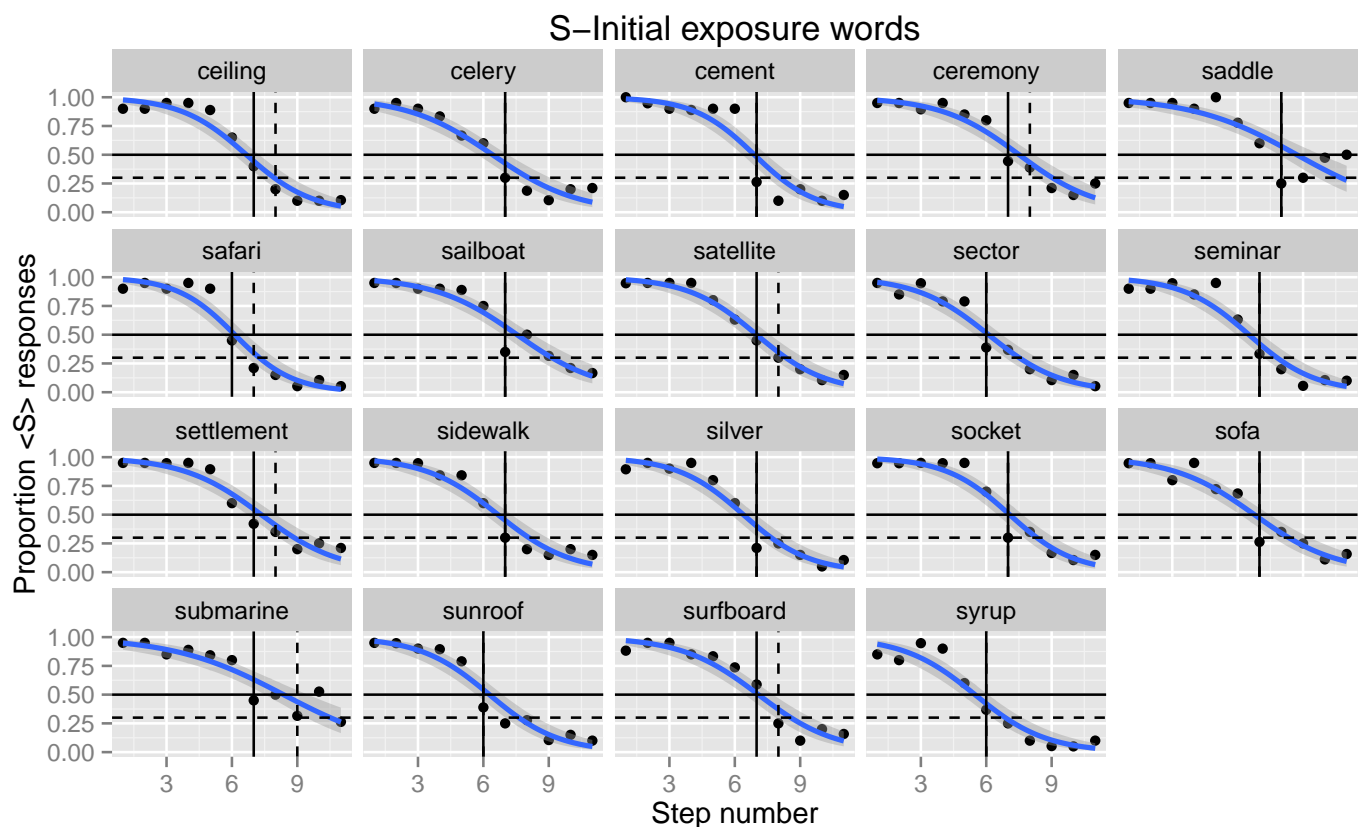


Figure 2.1: Proportion of word-responses for /s/-initial exposure words. Solid lines represent Experiment 1 selection criteria (50% word-response rate) and dashed lines represent Experiment 2 selection criteria (30% word-response rate). Dots are averaged word-response across subjects, and the blue line is a binomial model constructed from the responses.

Pretest

To determine which step of each continua would be used in exposure, a phonetic categorization experiment was conducted. Participants were presented with each step of each exposure word-nonword continua and each categorization minimal pair continua, resulting in 495 trials (40 exposure words plus five minimal pairs by 11 steps). The experiment was implemented in E-prime (cite). As half of the critical items had a sibilant in the middle of the word (onset of the final syllable), participants were asked to respond with word or non word rather than asking for the identity of the ambiguous sound, as in previous research [48].

The proportion of s-responses (or word responses for exposure items) at each step of each continua was calculated and the most ambiguous step chosen. The threshold for the ambiguous step for this experiment was when the percentage of s-response dropped near 50%. A full list of steps chosen for each stimulus item is in the appendix. For the minimal pairs, six steps surrounding the 50% cross over point were selected for use in the phonetic categorization task. Due to experimenter error, the continua for *seedling* was not included in the stimuli, so the chosen step was the average chosen step for the /s/-initial words. The average step chosen for /s/-initial words was 6.8 ($SD = 0.5$), and for /s/-final words the average step was 7.7 ($SD = 0.8$).

Procedure

Participants in the experimental conditions completed two tasks, an exposure task and a categorization task. The exposure task was a lexical decision task, where participants heard auditory stimuli and were instructed to respond with either

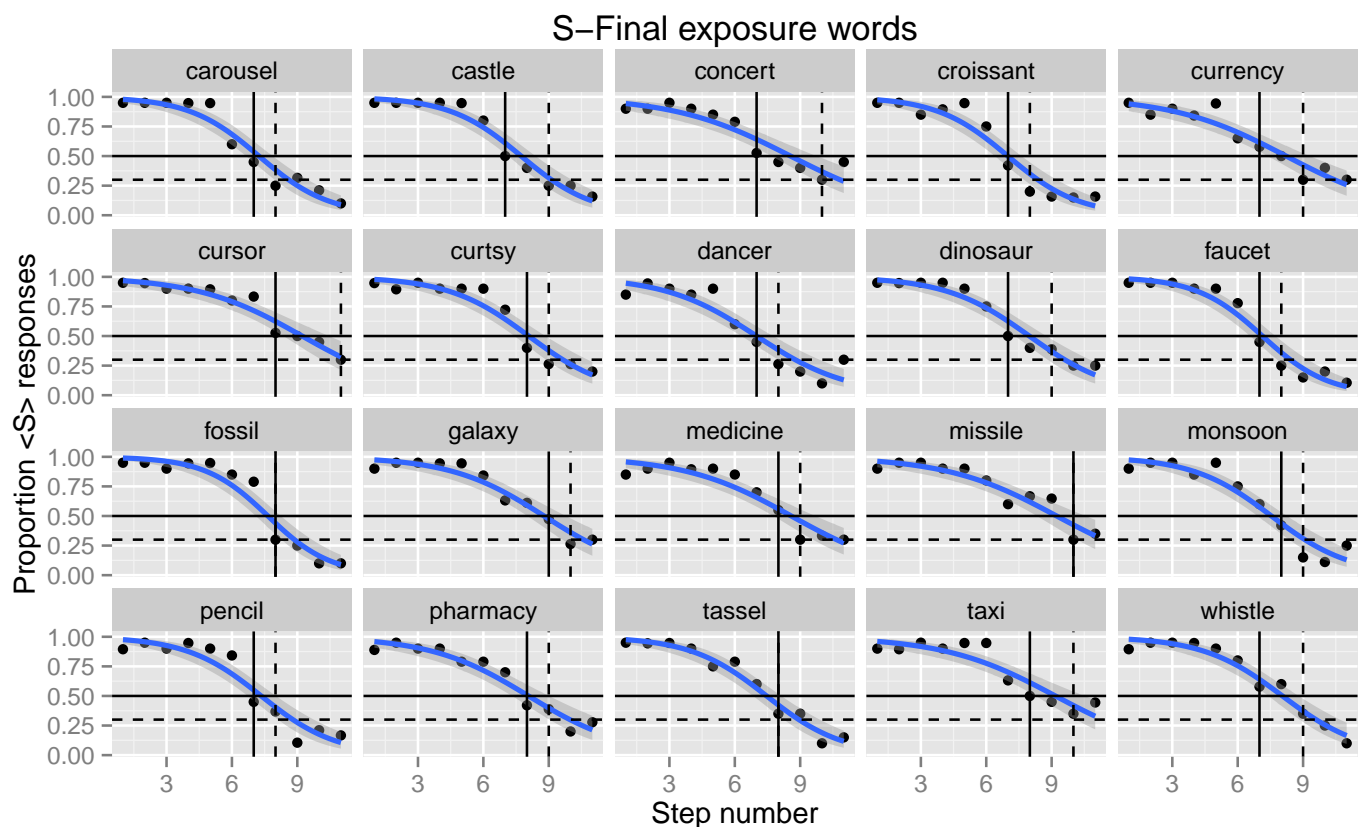


Figure 2.2: Proportion of word-responses for /s/-final exposure words. Solid lines represent Experiment 1 selection criteria (50% word-response rate) and dashed lines represent Experiment 2 selection criteria (30% word-response rate). Dots are averaged word-response across subjects, and the blue line is a binomial model constructed from the responses

”word” if they thought what they heard was a word or ”nonword” if they didn’t think it was a word. The buttons corresponding to ”word” and ”nonword” were counterbalanced across participants. Trial order was pseudorandom, with no critical or control items appearing in the first six trials, and no critical or control trials in a row, but random otherwise, following Reinisch et al. [48].

In the categorization task, participants heard an auditory stimulus and had to categorize it as one of two words, differing only in the onset sibilant (s vs sh). The buttons corresponding to the words were counterbalanced across participants. The six most ambiguous steps of the minimal pair continua were used with seven repetitions each, giving a total of 168 trials. Participants were instructed that there would be two tasks in the experiment, and both tasks were explained at the beginning to remove experimenter interaction between exposure and categorization.

Participants were assigned to one of four conditions. Two of the conditions exposed participants to only critical items that began with /s/, and the other two exposed them to only critical items that had an /s/ in the onset of the final syllable, giving a consistent 200 trials in all exposure phases with control and filler items shared across all participants. Additionally participants in half the conditions received additional instructions that the speaker’s ”s” sounds were sometimes ambiguous, and to listen carefully to ensure correct responses in the lexical decision.

2.2.2 Results

Control experiment

Responses with reaction times less than 200 ms or greater than 2500 ms were excluded from analyses. A logistic mixed effects models was fit with Subject and Continua as random effects and Step as a fixed effect with by-Subject and by-Item random slopes for Step. The intercept was not significant ($\beta = 0.43, SE = 0.29, z = 1.5, p = 0.13$), and Step was significant ($\beta = -2.61, SE = 0.28, z = -9.1, p < 0.01$).

Exposure

Trials with nonword stimuli and responses faster than 200 ms or slower than 2500 ms were excluded from analysis. Performance on the exposure task was high overall, with accuracy on filler trials averaging 92%. Word response rates for each of the four conditions did not differ significantly from each other, though S-Final/No Attention participants had a slightly higher average rate of 81% (SD= 17%) than the other conditions (S-Final/Attention: mean = 74%, SD = 18%; S-Initial/No Attention: mean = 74%, SD = 27%; S-Initial/Attention: mean = 76%, SD = 23%). A logistic mixed effects model with accuracy as the dependent variable was fit with fixed effects for trial type (Filler, S, SH), Attention (No Attention, Attention), Exposure Type (S-Initial, S-Final) and their interactions. The random effect structure was as maximally specified as possible with random effects for Subject and Word, and by-Subject random slopes for trial type and by-Word random slopes for Attention. The only fixed effects that were significant were a main effect of trial type for /s/ trials compared to filler trials ($\beta = -1.71, SE = 0.43, z = -3.97, p < 0.01$) and a main effect of Attention ($\beta = 0.76, SE = 0.38, z = 2.02, p = 0.04$). Trials containing an ambiguous /s/ were less likely to be responded to as a word, and participants instructed to pay attention to /s/ were more likely to correctly respond to words in general.

Categorization

Responses with reaction times less than 200 ms or greater than 2500 ms were excluded from analyses. Participants were excluded if their initial estimated cross over point for the continuum lay outside of the 6 steps presented (2 participants). A logistic mixed effects model was constructed with Subject and Continua as random effects and continua Step as random slopes, with 0 coded as a /f/ response and 1 as a /s/ response. Fixed effects for the model were Step, Exposure Type, Attention and their interactions.

There was a significant effect for the intercept ($\beta = 0.83, SE = 0.31, z = 2.6, p < 0.01$), indicating that participants categorized more of the continua as /s/ in general. There was also a significant main effect of Step ($\beta = -2.10, SE = 0.20, z = -10.3, p < 0.01$), and a significant interaction between Exposure Type and Attention ($\beta = -0.93, SE = 0.43, z = -2.14, p = 0.03$). There was a marginal main effect of Exposure Type ($\beta = 0.58, SE = 0.30, z = 1.8, p = 0.06$).

These results are shown in Figure 2.3. The solid lines show the control participants' categorization function across the 6 steps of the continua. The error bars show within-subject 95% confidence intervals at each step. When exposed to ambiguous /s/ tokens in the first syllables of words, participants show a general expansion of the /s/ category, but no differences in behaviour if they are warned about ambiguous /s/ productions. However, when the exposure is to ambiguous /s/ tokens later in the words, we can see differences in behaviour beyond the general /s/ category expansion. Participants not warned of the speaker's ambiguous tokens categorized more of the continua as /s/ than those who were warned of the speaker's ambiguous /s/ productions.

As an individual predictor of participants' performance we took the proportion critical word endorsements and compared these values to the estimated cross-over points. The crossover point was determined from the Subject random effect in the logistic mixed effects model [25]. There was a significant positive correlation between a participant's tolerance for the ambiguous exposure items and their crossover point on the continua ($r = 0.39, t(90) = 4, p < 0.01$), shown in Figure 2.4.

An ANOVA with cross-over point as the independent variable and word endorsement rate, Exposure Type, Attention and their interactions, found only a main effect of word endorsement rate ($F(1, 89) = 17.82, p < 0.01$), suggesting that listeners in different conditions were not affected differently from one another.

Figure 2.3: Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 1. In the S-Final condition, participants in the Attention condition showed a larger perceptual learning effect than those in the No Attention condition. In the S-Initial condition, there were no differences in perceptual learning between the Attention conditions. Error bars represent 95% confidence intervals.

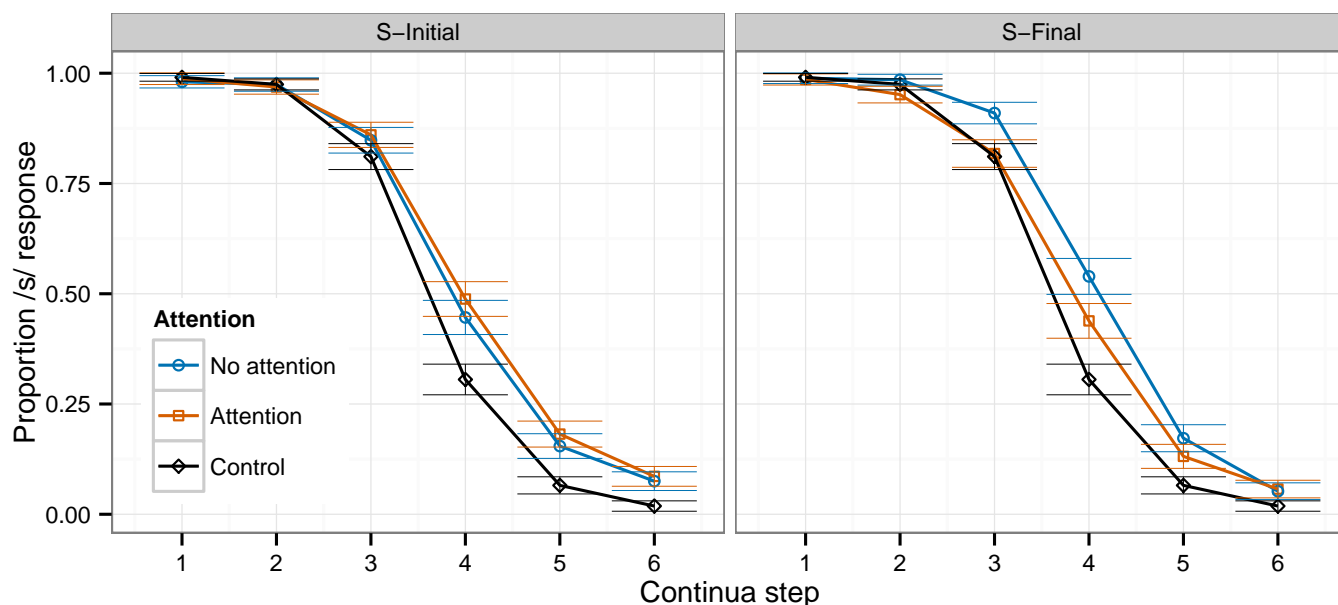
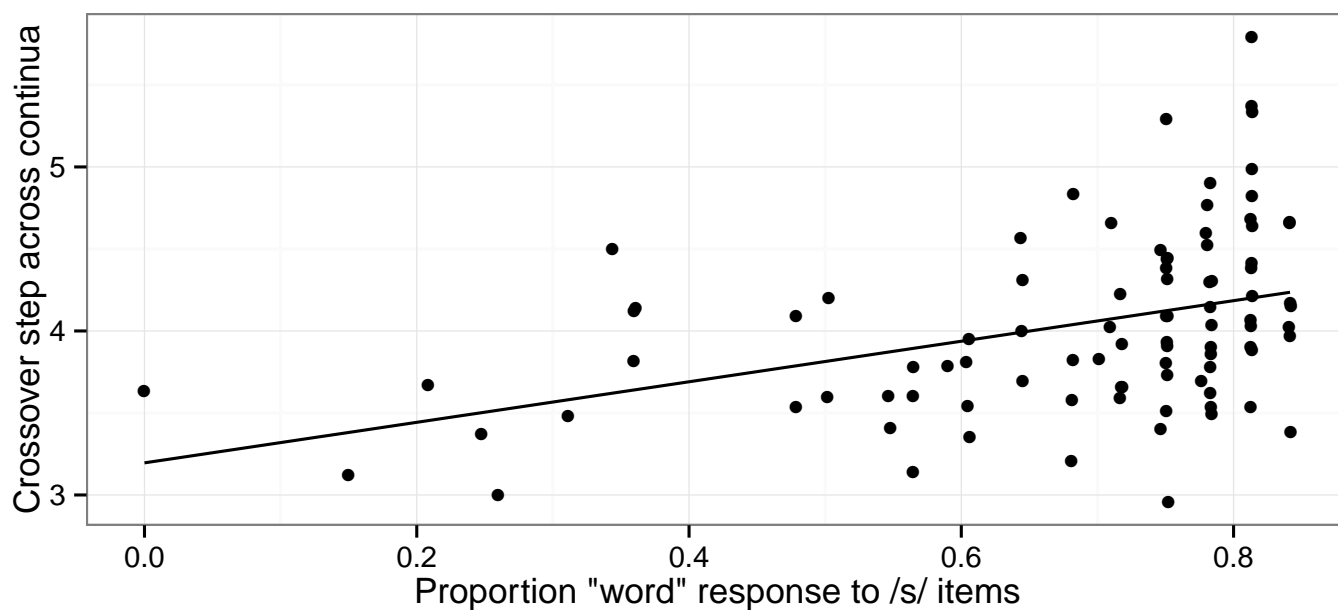


Figure 2.4: Correlation of crossover point in categorization with the proportion of word responses to critical items containing an ambiguous /s/ token.



2.2.3 Discussion

Perceptual learning effects in this experiment were robust across continua and experimental conditions, indicating the general automaticity of perceptual adaptation, but the degree of adaptation differed across conditions. Differences in attention only had an effect in the S-Final conditions, which replicates earlier findings that the effect of attention increased as the position of the ambiguous sound moved toward the end of the word [43]. The initial sound of a word is already a prominent position and listeners focus on the initial sounds to narrow the set of possible words they might be hearing. Later in a word, expectations for particular words given the preceding phonetic content would be greater, so directed attention on the phonetic detail shows a clear effect. That attention affected perceptual learning at all suggests that the adaptation is not wholly automatic and there is some degree of listener control.

The findings of this experiment do not support the attention mechanism of Clark [11]. If attention functioned as a gain mechanism, increasing the weight of error signals generated by mismatches between expectations and incoming signals, we should expect to see greater perceptual learning when listeners were instructed to attend to the speaker's /s/ sounds. Instead, the opposite occurred. Participants told to attend to the speaker's /s/ sounds showed smaller perceptual learning effects. This outcome may be determined by the nature of the instructions, which suggested that the ambiguity could harm accuracy, but different valences of attention are not predicted to behave any differently from one another in a gain mechanism.

The threshold for stimuli selection differed from that in previous studies. The closest study to this one used 70% of word responses in the pretest as the threshold for selection, creating a more typical, though still modified, category [48]. In that study, the lowest critical word endorsement rates in the exposure phase were 17 out of 20 (85%). In contrast, this experiment used 50% as the threshold and had correspondingly lower word endorsement rates (mean = 76%, $sd = 22\%$). Interestingly, despite the lower word endorsement rates and the less canonical stimuli used, perceptual learning effects remained robust, which raises the question, can perceptual learning occur from a modified category even more atypical than the one used in this experiment?

2.3 Experiment 2

Experiment 2 follows up on Experiment 1 by using stimuli that are farther from the canonical productions of the /s/ critical words. If the correlation from Experiment 1 holds, we expect to see lower word endorsement rates and lower (or perhaps non-existent) perceptual learning effects.

2.3.1 Methodology

Participants

A total of 127 participants from the UBC population completed the experiment and were compensated with either \$10 CAD or course credit. Of these, 31 were non-native speakers of English and their results were not included in the analyses, resulting in data from 96 participants.

Materials

Experiment 2 used the same items as Experiment 1, except that the step along the /s/-/ʃ/ continua chosen as the ambiguous sound had a different threshold. For this experiment, 30% identification as the /s/ word was used the threshold. The average step chosen for /s/-initial words was 7.3 ($SD = 0.8$), and for /s/-final words the average step was 8.9 ($SD = 0.9$).

Procedure

The procedure and instructions were identical between Experiment 1 and Experiment 2.

2.3.2 Results

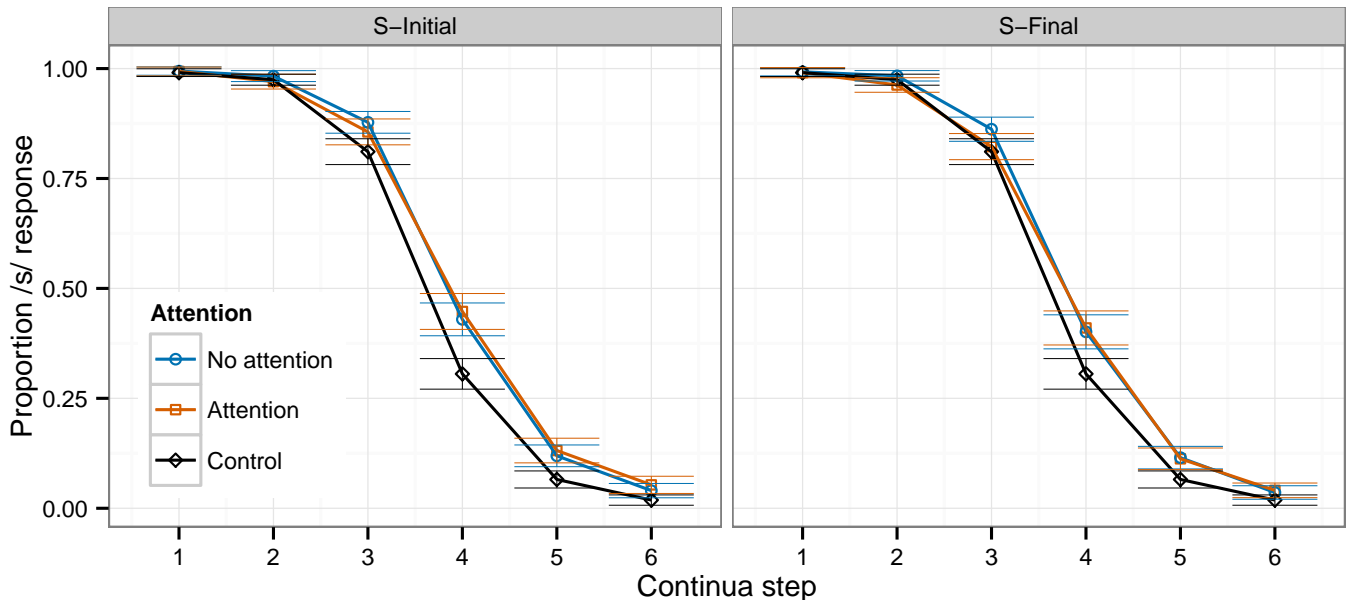
Exposure

Trials with nonword stimuli and responses faster than 200 ms or slower than 2500 ms were excluded from analysis. Performance on the exposure task was high overall, with accuracy on filler trials averaging 92%. An ANOVA of critical word endorsement rates revealed a marginal effect of Exposure Type ($F(1,92) = 3.86, p = 0.05$), with participants in the S-Final conditions having lower word endorsement rates (S-Final/Attention: mean = 56%, sd = 30%; S-Final/No Attention: mean = 52%, sd = 25%) than participants in the S-Initial conditions (S-Initial/Attention: mean = 68%, sd = 25%; S-Initial/No Attention: mean = 61%, sd = 23%). A logistic mixed effects model with accuracy as the dependent variable was fit with fixed effects for trial type (Filler, S, SH), Attention (No Attention, Attention), Exposure Type (S-Initial, S-Final) and their interactions. The random effect structure was as maximally specified as possible with random effects for Subject and Word, and by-Subject random slopes for trial type and by-Word random slopes for Attention. The only fixed effect that was significant were a main effect of trial type for /s/ trials compared to filler trials ($\beta = -2.51, SE = 0.46, z = -5.35, p < 0.01$).

Categorization

Responses with reaction times less than 200 ms or greater than 2500 ms were excluded from analyses. Participants were excluded if their initial estimated cross over point for the continuum lay outside of the 6 steps presented (2 participants). A logistic mixed effects model was constructed with Subject and Continua as random effects and continua Step as random slopes, with 0 coded as a /f/ response and 1 as a /s/ response. Fixed effects for the model were Step, Exposure Type, Attention and their interactions.

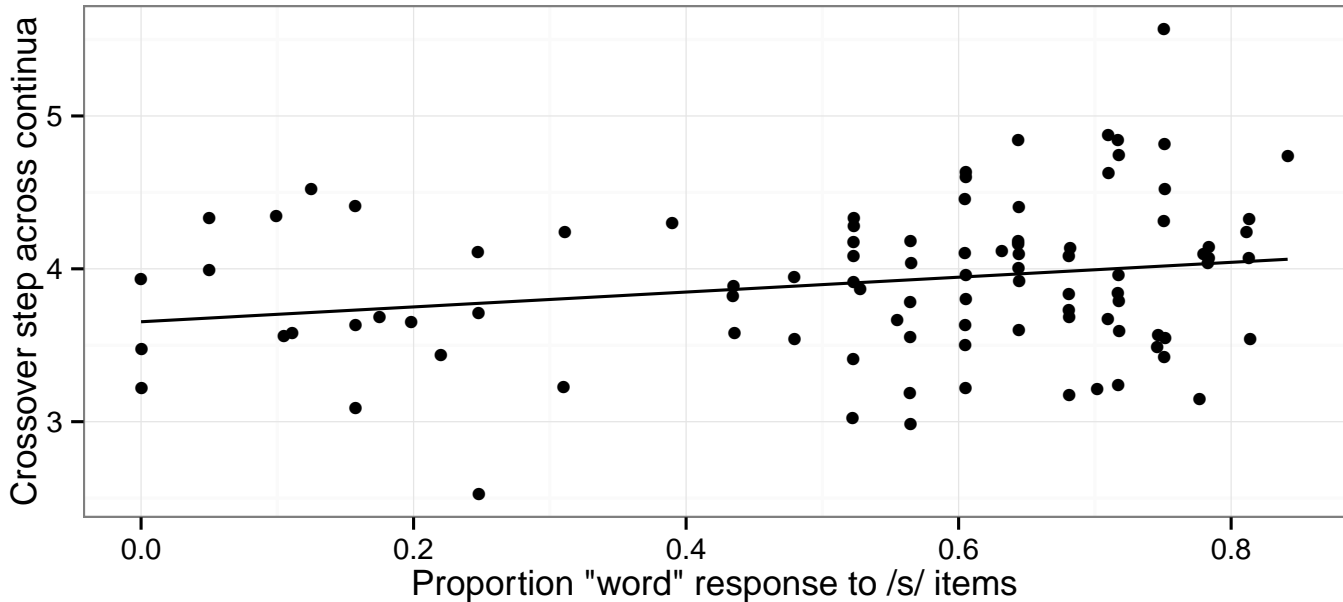
Figure 2.5: Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 2. Participants showed no significant differences across conditions. Error bars represent 95% confidence intervals.



There was a significant effect for the Intercept ($\beta = 1.01, SE = 0.38, z = 2.6, p < 0.01$), indicating that participants categorized more of the continua as /s/ in general. There was also a significant main effect of Step ($\beta = -2.67, SE =$

0.23, $z = -11.2$, $p < 0.01$). There were no other significant main effects or interactions, though an interaction between Step and Attention trended toward significant ($\beta = 0.35$, $SE = 0.21$, $z = 1.6$, $p = 0.09$).

Figure 2.6: Correlation of crossover point in categorization with the proportion of word responses to critical items containing an ambiguous /s/ token in Experiment 2.



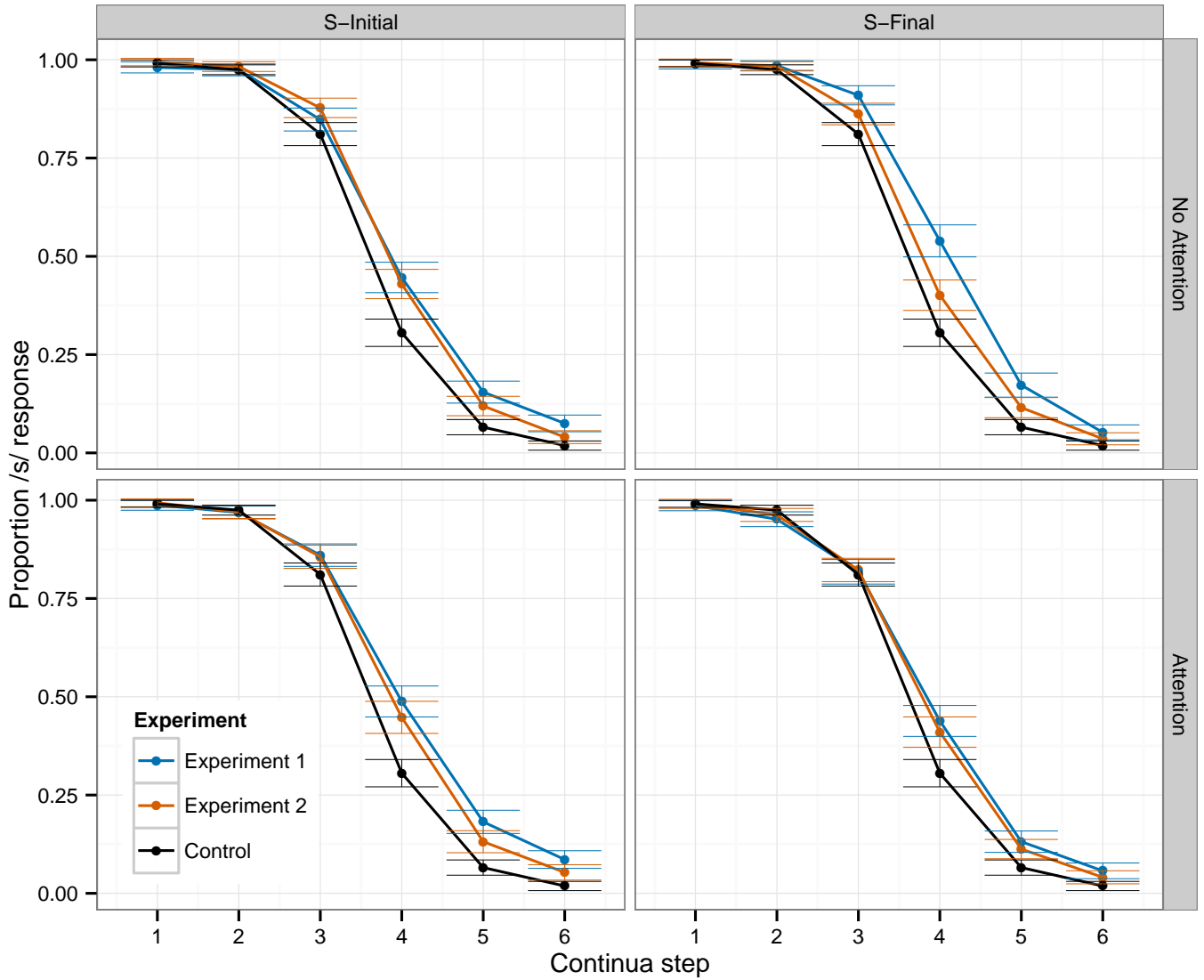
As in Experiment 1, the proportion critical word endorsements was calculated for each subject and assessed for correlation with participants' crossover points. There was a significant positive correlation between a participant's tolerance for the ambiguous exposure items and their crossover point on the continua ($r = 0.22$, $t(92) = 2.25$, $p = 0.02$), shown in Figure 2.6.

2.4 Grouped results across experiments

To see what degree the stimuli used had an effect on perceptual learning, the data from Experiment 1 and Experiment 2 were pooled and analyzed identically as above, but with Experiment and its interactions as fixed effects. In the logistic mixed effects model, there was significant main effects for Intercept ($\beta = 1.00$, $SE = 0.36$, $z = 2.7$, $p < 0.01$) and Step ($\beta = -2.64$, $SE = 0.21$, $z = -12.1$, $p < 0.01$), and a significant two-way interaction between Experiment and Step ($\beta = 0.51$, $SE = 0.20$, $z = 2.5$, $p = 0.01$), and a marginal four-way interaction between Step, Exposure Type, Attention and Experiment ($\beta = 0.73$, $SE = 0.42$, $z = 1.7$, $p = 0.08$). These results can be seen in Figure 2.7. The four-way interaction can be seen in S-Final/No Attention conditions across the two experiments, where Experiment 1 has a significant difference between the Attention and No Attention condition, but Experiment 2 does not. The two-way interaction between Experiment and Step and the lack of a main effect for Experiment potentially suggests that while the category boundary was not significantly different across experiments, the slope of the categorization function was.

To see if there was a difference with the Experiment 1 in how word endorsement rates affected crossover points, the data was pooled for the two experiments. An ANOVA with cross-over point as the independent variable and word endorsement rate, Exposure Type, Attention, Experiment and their interactions, found a main effect of word endorsement rate ($F(1, 185) = 21.82$, $p < 0.01$) and marginal interaction between word endorsement rates and Experiment ($F(1, 185) = 3.11$, $p = 0.07$).

Figure 2.7: Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 1 and Experiment 2. Error bars represent 95% confidence intervals.



2.5 General discussion

Perceptual learning was found across all experimental conditions. Attention only had effect on perceptual learning in the S-Final condition of Experiment 1.

Compared to Experiment 1, Experiment 2 had a weaker, yet still significant, correlation between critical word endorsement rates and crossover boundary points, suggesting that although the stimuli used in Experiment 2 were farther from the canonical production, they did not shift the category boundary as much as the stimuli in Experiment 1. While neither attention or position of the ambiguous sound in the word had an effect on the correlation, the distance from the canonical production did. This potentially suggests that the degree to which a category is shifted is inversely related to its distance. Or perhaps more likely, the distance is inversely related to its goodness or plausibility that the ambiguous sound was indeed meant to be an /s/.

The correlation between word response rate in the exposure phase and the category boundary in categorization phase

across both experiments raises two possible explanations. In a causal interpretation between exposure and categorization, as each ambiguous sound is processed and errors propagate, the distribution for that category (for that particular speaker) is updated. Participants who processed more of the ambiguous sound as an /s/ updated their perceptual category for /s/ more. This explanation fits within a larger neo-generative model of spoken language processing [41]. A non-causal story is also plausible: the correlation may reveal individual differences on the part of the participants, where some participants are more adaptable or tolerant of variability than others, leading to greater degrees of perceptual adaptation. Individual differences in attention-switching control have previously been found to affect perceptual learning [54], which supports a non-causal interpretation as well.

The perceptual learning effects found in Experiment 1 and 2 appear to fall into two categories, aligning with either comprehension-oriented or perception-oriented attentional sets. For most of exposure conditions, participants showed a small perceptual learning effect of similar magnitude. These conditions perception-oriented attentional sets were likely to be recruited, such as when attention was explicitly drawn to the ambiguous sounds, or the ambiguous sounds were at the beginnings of words, or when the ambiguous sounds were the least typical of /s/. The participants exposed to the ambiguous sounds with the most lexical bias and with no attention drawn to them, however, showed a significantly stronger perceptual learning effect. This condition is the one where comprehension-oriented attentional sets were predicted to dominate. The perception-oriented attentional sets exhibit less generalization, like what is seen in psychophysics literature and in visually-guided perceptual learning in speech perception [49].

As mentioned in the discussion for Experiment 1, the findings do not support a simple gain mechanism for attention as proposed in Clark [11]. In Yeshurun and Carrasco [62], attention to areas with finer spatial resolution caused observers to miss larger patterns. If attention simply boosted the error signal, attention should always be beneficial to perceptual learning. The findings of these two experiments supports a larger role for attention in a predictive coding framework, as advanced in Block and Siegel [5]. The propagation-limiting mechanism attention proposed earlier explains both the findings in the visual domain and the current findings. Attention to a level of representation causes errors between expectations and observed signals to be resolved and updated at that level. If attention is more oriented towards comprehension, either of language or other patterns in the world, errors can be propagated higher before updating expectations.

The manipulations in the two experiments in this chapter were ones that induced perception-oriented attentional sets. The task itself focuses on recognition, and a comprehension-oriented attentional set is likely more helpful for this task than a perception-oriented one. However, participants likely switch between the two throughout the course of the experiment. To fully examine the use of perception- and comprehension-oriented attentional sets in perceptual learning, manipulations that induce comprehension-oriented attentional sets are necessary. In the following chapter, such a manipulation is implemented through increasing linguistic expectations from semantic predictability.

Chapter 3

Cross-modal word identification

3.1 Motivation

In Chapter 2, listeners showed greater perceptual learning effects when the lexical bias was greater. However, this relationship was not found when attention was drawn to the ambiguous sound and when the ambiguous sound was farther from the canonical production. The experiment in this chapter seeks to increase the linguistic expectations for the words containing the ambiguous productions as a way to induce more comprehension-oriented attentional sets, even in the face of explicit instructions promoting perception-oriented attentional sets. To this end, semantic predictability is used to boost linguistic expectations in conjunction with lexical bias.

Semantic predictability in production studies has been found to have an effect on vowel realization and duration independent of lexical factors such as frequency and neighbourhood density [12, 52]. In speech perception work, semantic predictability has been found to have an effect on phoneme categorization similar to that of lexical bias [6], and increased intelligibility in noise, particularly for native speakers [7, 19, 23, 37, and others]. In phoneme restoration, semantic predictability both increased the bias to respond that a word intact and also increased the sensitivity to noise addition versus noise replacement [50].

As in previous work, two levels of semantic predictability are used in this experiment. Highly predictable sentences are ones where the final word is constrained to only a few words. Unpredictable sentences are ones where the possible completions are numerous, and any completions given are more guesses than anything else. All target words and the threshold for ambiguity are taken from Experiment 1 and have the highest lexical bias available. When a listener is exposed to these words in unpredictable sentences, they should show perceptual learning effects equivalent to those participants in the S-Final condition of Experiment 1.

The exposure task in this experiment differs from lexical decision task used in the previous chapter. Instead, participants identify the picture that matches the final word in the sentence, and all trials involve a real word. The attentional set for this task, where participants make a decision about the identity of a specific word, may be different from the previous experiment, where they make a decision about the lexicality of a word, but the decision is about the word in both instances.

If linguistic expectations are cumulative, we should expect to see a larger perceptual learning effect for participants exposed to ambiguous sounds in words that are highly predictable, due to greater use of comprehension-oriented attentional sets. Additionally, if the explicit attention does exert a categorical effect on perceptual learning, we should see a similar size of perceptual learning in this experiment in the attention conditions as in the previous experiments. Previous research has found an increased sensitivity to signal properties in semantically predictable sentences [50], which could lead to one of two outcomes. Given that bias to restore the phoneme was greater in these conditions as well, we might see greater perceptual learning effects, as the link between the word and the production would be enhanced, but the signal would be encoded more accurately from the lower cognitive load. On the other hand, the lower cognitive load could lead to more attention available to notice the deviant production, and thus adopt a perception-oriented attentional set.

3.2 Methodology

3.2.1 Participants

One hundred native speakers of English participated in the experiment and were compensated with either \$10 CAD or course credit. They were recruited from the UBC student population. Twenty additional native English speakers participated in a pretest to determine sentence predictability, and 10 other native English speakers participated in a picture naming pretest. Participants were native North American English speakers with no reported speech or hearing disorders.

3.2.2 Materials

One hundred and twenty sentences were used as exposure materials. The set of sentences consisted of 40 critical sentences, 20 control sentences and 60 filler sentences. The critical sentences ended in one of 20 of the critical words in Experiments 1 and 2 that had an /s/ in the onset of the final syllable. The 20 control sentences ended in the 20 control items used in Experiments 1 and 2, and the 60 filler sentences ended in the 60 filler words in Experiments 1 and 2. Half of all sentences were written to be predictive of the final word, and the other half were written to be unpredictable of the final word. Unlike previous studies using sentence or semantic predictability [23], Unpredictive sentences were written with the final word in mind with a variety of sentence structures, and the final words were plausible objects of lexical verbs and prepositions. A full list of words and their contexts can be found in the appendix. Aside from the sibilants in the critical and control words, the sentences contained no sibilants (/s z ʒ ʃ ʒ ʃ/). The same minimal pairs for phonetic categorization as in Experiments 1 and 2 were used.

Sentences were recorded by the same male Vancouver English speaker used in Experiments 1 and 2. Critical sentences were recorded in pairs, with one normal production and then a production of the same sentence with the /s/ in the final word replaced with an /ʃ/. The speaker was instructed to produce both sentences with comparable speech rate, speech style and prosody.

As in Experiments 1 and 2, the critical items were morphed together into an 11-step continuum using STRAIGHT [24]; however, only the final word in sentence was morphed. For all steps, the preceding words in the sentence were kept as the natural production to minimize artifacts of the morphing algorithm. The control and filler items were also processed and resynthesized to ensure consistent quality. The ambiguous point selection was based on the pretest performed for Experiment 1 and 2 exposure items. The ambiguous steps of the continua chosen corresponded to the 50% cross over point in Experiment 2.

Pictures of 200 words, with 100 pictures for the final word of the sentences and 100 for distractors, were selected in two steps. First, a research assistant selected five images from a Google image search of the word, and then a single image representing that word was selected from amongst the five by me. To ensure consistent behaviour in E-Prime, pictures were resized to fit within a 400x400 area with a resolution of 72x72 DPI and converted to bitmap format. Additionally, any transparent backgrounds in the pictures were converted to plain white backgrounds.

3.2.3 Pretest

The same twenty participants that completed the lexical decision continua pre-test also completed a sentence predictability task before the phonetic categorization task described in Experiment 1. Participants were compensated with \$10 CAD for both tasks, and were native North American English speakers with no reported speech, language or hearing disorders. In this task, participants were presented with sentence fragments that were lacking in the final word. They were instructed to type in the word that came to mind when reading the fragment, and to enter any additional words that came to mind that would also complete the sentence. There was no time limit for entry and participants were shown an example with the fragment “The boat sailed across the...” and the possible completions “bay, ocean, lake, river”. Responses were collected in E-Prime (cite), and were sanitized by removing miscellaneous keystrokes recorded by E-Prime, spell checking, and standardizing variant spellings and plural forms.

The measure used for determining rewriting of sentences was the proportion of participants that included the target word in their responses. For predictive sentences, the mean proportion was 0.49 (range 0-0.95) and for unpredictable

sentences, the mean proportion was 0.03 (range 0-0.45). Predictive sentences that had target response proportions of 20% or less were rewritten. The predictive sentences for *auction*, *brochure*, *carousel*, *cashier*, *cockpit*, *concert*, *cowboy*, *currency*, *cursor*, *cushion*, *dryer*, *graffiti*, and *missile* were rewritten to remove any ambiguities.

Five volunteers from the Speech in Context lab participated in another pretest to determine how suitable the pictures were at representing their associated word. All participants were native speakers of North American English, with reported corrected-to-normal vision. Participants were presented with a single image in the middle of the screen. Their task was to type the word that first came to mind, and any other words that described the picture equally well. There was no time limit and presentation of the pictures was self-paced. Responses were sanitized as in the first pretest.

Pictures were replaced if 20% or less of the participants (1 of 5) responded with the target word and the responses were semantically unrelated to the target word. Five pictures were replaced, *toothpick* and *falafel* with clearer pictures and *ukulele*, *earmuff* and *earplug* were replaced with *rollerblader*, *anchor* and *bedroom*. All five replacements were for distractor words.

3.2.4 Procedure

As in Experiments 1 and 2, participants completed two tasks, an exposure task and a categorization task. For the exposure task, participants heard a sentence via headphones for each trial. Immediately following the auditory presentation, they were presented with two pictures on the screen. Their task was to select the picture on the screen that corresponded to the final word in the sentence they heard. As in Experiments 1 and 2, the order was pseudorandom, with the same constraints.

Participants were assigned to one of four groups of 25 participants. In the exposure phase, half of the participants were exposed to a modified /s/ sound only in Predictive sentences and half were exposed to it only in Unpredictive sentences. Half of all participants were told that the speaker's production of "s" was sometimes ambiguous, and to listen carefully to ensure correct responses.

Following the exposure task, participants completed the same categorization task described in Experiments 1 and 2.

3.3 Results

Collection of data is currently still ongoing.

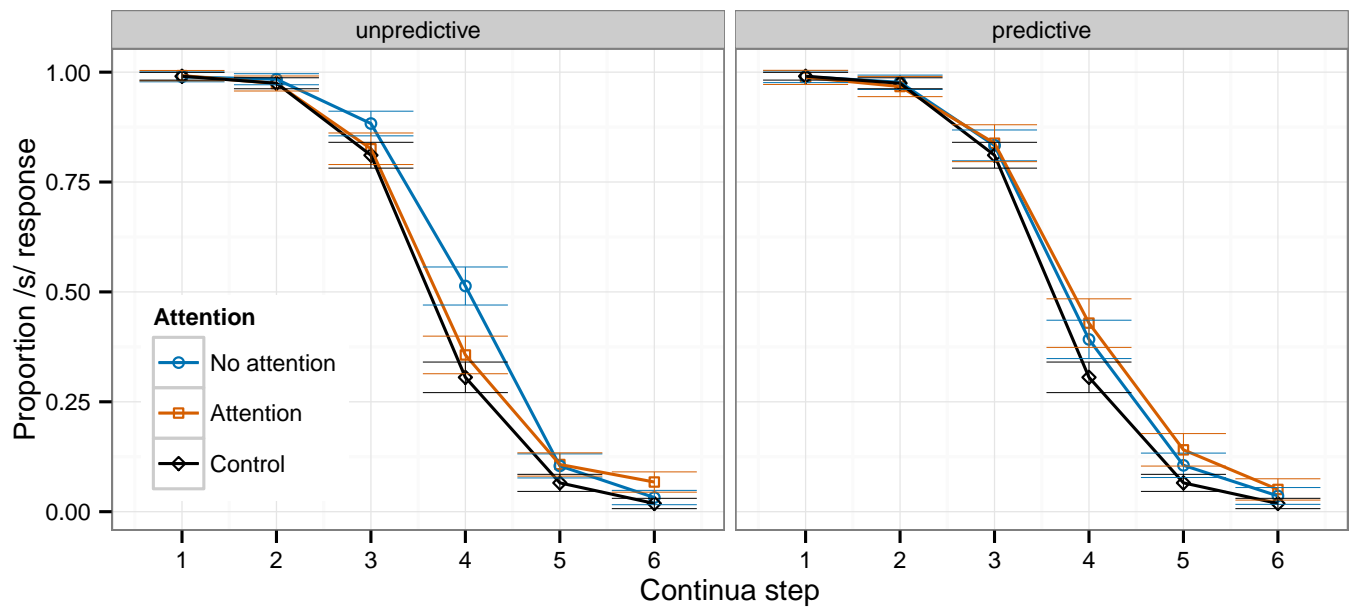
3.3.1 Exposure

Performance in the task was high, with accuracy at ceiling.

3.3.2 Categorization

3.4 Discussion

Figure 3.1: Proportion /s/ response along the 6 step continua as a function of Exposure Type and Attention in Experiment 3. Error bars represent 95% confidence intervals.



Chapter 4

Conclusions

This dissertation has examined the influence of attention and linguistic salience on perceptual learning in speech perception. Through explicit instructions and experimental manipulations, listeners were induced towards a more comprehension-oriented attentional set or a more perception-oriented attentional set. Those participants in conditions biased towards comprehension showed larger perceptual learning effects than those in conditions biased more towards perception. I adopted a predictive coding framework [11] that provides an intuitive model of perceptual learning and has been computationally implemented [26]. However, the attentional mechanism in the predictive coding framework does not work well for some instances of visual attention [5], or for the current results. I propose a new attentional mechanism for predictive coding, one in which attention blocks error propagation beyond the level to which attention is directed. Such a mechanism explains both the previous findings and the current results.

4.1 Attentional control of perceptual learning

The findings of Experiment 1 suggest also that perceptual learning is not a wholly automatic process. Although all participants showed perceptual learning effects, those exposed to the ambiguous sound with increased lexical bias only showed larger perceptual learning effects when the instructions about the speaker’s ambiguous sound were withheld. Attention on the ambiguous sound equalized the perceptual learning effects across lexical bias. However, in Experiment 2, there is no such effect of attention, suggesting that the ambiguous sounds that were farther away from the canonical production drew the listener’s attention to those sounds.

One question raised by the current results is whether attention always decreases perceptual learning. The instructions used to focus the listener’s attentional set framed the ambiguity in a negative way, with listeners being cautioned to listen carefully to ensure they made the correct decision. If the attention were directed to the ambiguous sound by framing the ambiguity in a positive way, would we still see the same pattern of results? The current mechanism would predict that attention of any kind to signal properties would block the propagation of errors, reducing perceptual learning. This prediction will be tested in future work.

4.2 Specificity versus generalization in perceptual learning

The results of this dissertation speak to dichotomy between specificity and generalization found in perceptual learning work. In Experiment 1, greater perceptual learning was shown by participants exposed to ambiguous sounds later in the words, not to those at the beginnings of words. And yet, the testing continua consisted of stimuli with the sibilant at the beginnings of words, which are more similar to the words beginning with the ambiguous sound.

The effect of attention in Experiments 1 and 2 suggest that when speakers are focused on the signal, or have their attention drawn to the signal, they are less likely to generalize across forms. When attention on the signal is not present, errors can propagate more widely causing greater updating. These different modes of perceptual learning may account for some differences that have been observed in the literature. While visually-guided perceptual learning shows comparable effects to lexically-guided perceptual learning [58], lexically-guided perceptual learning is more likely to be expanded to new contexts [29, 40], though with some restrictions [38], than visually-guided perceptual learning [49].

4.3 Effect of increased linguistic expectations

Comprehension-oriented attentional sets were promoted mainly through the task demands and the use of increased lexical and semantic expectations. The conditions that correspond most closely to those in other lexically-guided perceptual learning experiments [27, 40] are those participants not told of the ambiguous sound in Experiment 1. Under those conditions, there is a clear effect of lexical bias, such that listeners exposed to ambiguous stimuli with greater lexical bias update their speaker-specific category for that sound more.

4.4 Category atypicality

In Experiment 2, there was no effect of attention or lexical bias on perceptual learning, with a stable effect present for all listeners. There are two potential, non-exclusive explanations for the lack of effects. As stated above, the increased distance to the canonical production drew the listener’s attention to the ambiguous productions, resulting in a perception-oriented attentional set. The second potential explanation is that the productions farther from canonical produce a weaker effect on the updating of a listener’s categories, as predicted from the neo-generative model in [41]. This explanation is supported in part by the weaker correlation between word endorsement rate and crossover point found in Experiment 2, and the findings of Sumner [56] where the most perceptual learning was found when the categories begin like the listener expects and gradually shift toward the speaker’s actual boundaries over the course of exposure. This explanation could be tested straightforwardly by implementing the same gradual shift paradigm used in Sumner [56] with the manipulations used in this dissertation.

An interesting extension to the current findings would be to observe the perceptual learning effects in a cognitive load paradigm. Speech perception under cognitive load has shown to have greater reliance on lexical information due to weaker initial encoding of the signal [36]. Following exposure to a modified ambiguous category, we might expect to see less perceptual learning if the exposure was accompanied by high cognitive load. However, Scharenborg et al. [54] found that hearing loss of older participants did not significantly influence their perceptual learning. Therefore it may be that perceptual learning would likewise be similar across cognitive loads. Higher cognitive loads, however, might allow for more atypical ambiguous stimuli to be learned, due to the increased reliance on lexical information during initial encoding.

4.5 Conclusion

This dissertation investigated the influences of attention and linguistic salience on perceptual learning in speech perception. Perceptual learning was modulated by the attentional set of the listener. Decreased linguistic salience by increasing linguistic expectations induced more comprehension-oriented attentional sets, resulting in larger perceptual learning effects. Inducing perception-oriented attentional sets through explicit instructions or making the modified sound category more atypical resulted in smaller perceptual learning effects. These results support a greater role of attention than previously assumed in predictive coding frameworks, such as the proposed propagation-blocking mechanism. Finally, these results suggest that the degree to which listeners perceptually adapt to a new speaker is under their control to a degree, but given the robust perceptual learning effects across all conditions, perceptual learning is a largely automatic process.

Bibliography

- [1] M. Ahissar and S. Hochstein. Attentional control of early perceptual learning. *Proceedings of the National Academy of Sciences of the United States of America*, 90(12):5718–22, 1993. ISSN 0027-8424. doi:10.1073/pnas.90.12.5718. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=46793&tool=pmcentrez&rendertype=abstract>. → pages 5
- [2] W. F. Bacon and H. E. Egeth. Overriding stimulus-driven attentional capture. *Perception & psychophysics*, 55(5): 485–496, 1994. ISSN 0031-5117. → pages 6
- [3] J. N. Beckman. *Positional Faithfulness*. PhD thesis, University of Massachusetts Amherst, 1998. URL <http://onlinelibrary.wiley.com/doi/10.1002/9780470756171.ch16/summary>. → pages 2
- [4] P. Bertelson, J. Vroomen, and B. De Gelder. Visual Recalibration of Auditory Speech Identification: A McGurk Aftereffect. *Psychological Science*, 14(6):592–597, 2003. → pages 3
- [5] N. Block and S. Siegel. Attention and perceptual adaptation. *The Behavioral and brain sciences*, 36(3):205–6, 2013. ISSN 1469-1825. URL <http://www.ncbi.nlm.nih.gov/pubmed/23663745>. → pages 7, 19, 24
- [6] S. Borsky, B. Tuller, and L. P. Shapiro. "How to milk a coat:" the effects of semantic and acoustic information on phoneme categorization. *Journal of the Acoustical Society of America*, 103(5 Pt 1):2670–2676, 1998. URL <http://www.ncbi.nlm.nih.gov/pubmed/9604360>. → pages 5, 8, 20
- [7] A. R. Bradlow and J. a. Alexander. Semantic and phonetic enhancements for speech-in-noise recognition by native and non-native listeners. *Journal of the Acoustical Society of America*, 121(4):2339–2349, 2007. ISSN 00014966. doi:10.1121/1.2642103. URL <http://link.aip.org/link/JASMAN/v121/i4/p2339/s1&Agg=doi>. → pages 5, 20
- [8] A. R. Bradlow and T. Bent. Perceptual adaptation to non-native speech. *Cognition*, 106(2):707–729, 2008. → pages 1
- [9] M. Brysbaert and B. New. Moving beyond Kucera and Francis: a critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior research methods*, 41(4):977–990, 2009. ISSN 1554-3528. → pages 10
- [10] E. Clare. Applying phonological knowledge to phonetic accommodation, 2014. → pages 4
- [11] A. Clark. Whatever next? Predictive brains, situated agents, and the future of cognitive science. *The Behavioral and brain sciences*, 36(3):181–204, 2013. ISSN 1469-1825. URL <http://www.ncbi.nlm.nih.gov/pubmed/23663408>. → pages 1, 3, 7, 15, 19, 24
- [12] C. G. Clopper and J. B. Pierrehumbert. Effects of semantic predictability and regional dialect on vowel space reduction. *Journal of the Acoustical Society of America*, 124(3):1682–1688, Sept. 2008. ISSN 1520-8524. doi:10.1121/1.2953322. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2676620&tool=pmcentrez&rendertype=abstract>. → pages 5, 20

- [13] C. Connine. Constraints on interactive processes in auditory word recognition: The role of sentence context. *Journal of Memory and Language*, 538:527–538, 1987. URL <http://www.sciencedirect.com/science/article/pii/0749596X87901380>. → pages 5
- [14] C. M. Connine and C. Clifton. Interactive use of lexical information in speech perception. *Journal of experimental psychology. Human perception and performance*, 13(2):291–299, 1987. ISSN 0096-1523. → pages 5
- [15] A. Cutler, J. Mehler, D. Norris, and J. Segui. Phoneme identification and the lexicon, 1987. ISSN 00100285. → pages 1, 6
- [16] P. D. Eimas, W. E. Cooper, and J. D. Corbit. Some properties of linguistic feature detectors, 1973. ISSN 0031-5117. → pages 2
- [17] F. Eisner and J. M. McQueen. The specificity of perceptual learning in speech processing. *Perception & psychophysics*, 67(2):224–238, 2005. → pages 3
- [18] F. Eisner and J. M. McQueen. Perceptual learning in speech: Stability over time. *The Journal of the Acoustical Society of America*, 119(4):1950, 2006. ISSN 00014966. doi:10.1121/1.2178721. URL <http://scitation.aip.org/content/asa/journal/jasa/119/4/10.1121/1.2178721>. → pages 3
- [19] M. Fallon, S. E. Trehub, and B. A. Schneider. Children’s use of semantic cues in degraded listening environments. *The Journal of the Acoustical Society of America*, 111(5 Pt 1):2242–2249, 2002. ISSN 00014966. → pages 5, 20
- [20] W. F. Ganong. Phonetic categorization in auditory word perception. *Journal of Experimental Psychology: Human Perception and Performance*, 6(1):110–125, 1980. URL <http://www.ncbi.nlm.nih.gov/pubmed/6444985>. → pages 4, 8
- [21] E. J. GIBSON. Improvement in perceptual judgments as a function of controlled practice or training. *Psychological bulletin*, 50(6):401–431, 1953. → pages 2
- [22] C. Gilbert, M. Sigman, and R. Crist. The neural basis of perceptual learning. *Neuron*, 31:681–697, 2001. ISSN 0896-6273. doi:10.1016/S0896-6273(01)00424-X. URL <http://www.sciencedirect.com/science/article/pii/S089662730100424X>. → pages 3, 8
- [23] D. Kalikow, K. Stevens, and L. Elliott. Development of a test of speech intelligibility in noise using sentence materials with controlled word predictability. ... *Journal of the Acoustical Society of ...*, 61(5), 1977. URL <http://link.aip.org/link/?JASMAN/61/1337/1>. → pages 5, 8, 20, 21
- [24] H. Kawahara, M. Morise, T. Takahashi, R. Nisimura, T. Irino, and H. Banno. Tandem-straight: A temporally stable power spectral representation for periodic signals and applications to interference-free spectrum, F0, and aperiodicity estimation. In *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, pages 3933–3936, 2008. → pages 10, 21
- [25] F. Kleber, J. Harrington, and U. Reubold. The Relationship between the Perception and Production of Coarticulation during a Sound Change in Progress, 2012. ISSN 0023-8309. → pages 13
- [26] D. Kleinschmidt and T. F. Jaeger. A Bayesian belief updating model of phonetic recalibration and selective adaptation. *Science*, (June):10–19, 2011. URL <http://www.aclweb.org/anthology-new/W/W11/W11-06.pdf#page=20>. → pages 3, 24
- [27] T. Kraljic and A. G. Samuel. Perceptual learning for speech: Is there a return to normal? *Cognitive psychology*, 51(2):141–78, Sept. 2005. ISSN 0010-0285. doi:10.1016/j.cogpsych.2005.05.001. URL <http://www.ncbi.nlm.nih.gov/pubmed/16095588>. → pages 3, 4, 25

- [28] T. Kraljic and A. G. Samuel. Perceptual adjustments to multiple speakers. *Journal of Memory and Language*, 56(1):1–15, Jan. 2007. ISSN 0749596X. doi:10.1016/j.jml.2006.07.010. URL <http://linkinghub.elsevier.com/retrieve/pii/S0749596X06000842>. → pages 3
- [29] T. Kraljic, S. E. Brennan, and A. G. Samuel. Accommodating variation: dialects, idiolects, and speech processing. *Cognition*, 107(1):54–81, Apr. 2008. ISSN 0010-0277. doi:10.1016/j.cognition.2007.07.013. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2375975&tool=pmcentrez&rendertype=abstract>. → pages 4, 7, 24
- [30] T. Kraljic, A. G. Samuel, and S. E. Brennan. First impressions and last resorts: how listeners adjust to speaker variability. *Psychological science*, 19(4):332–8, Apr. 2008. ISSN 0956-7976. doi:10.1111/j.1467-9280.2008.02090.x. URL <http://www.ncbi.nlm.nih.gov/pubmed/18399885>. → pages 1, 4, 7
- [31] P. K. Kuhl. Speech perception in early infancy: perceptual constancy for spectrally dissimilar vowel categories. *The Journal of the Acoustical Society of America*, 66(6):1668–1679, 1979. ISSN 00014966. → pages 1
- [32] P. Ladefoged and D. E. Broadbent. Information conveyed by vowels. *The Journal of the Acoustical Society of America*, 29(1):98–104, 1957. ISSN 00014966. URL <http://scitation.aip.org/content/asa/journal/jasa/29/1/10.1121/1.1908694>. → pages 1
- [33] A. B. Leber and H. E. Egeth. Attention on autopilot: Past experience and attentional set, 2006. ISSN 1350-6285. → pages 6
- [34] R. Levy. Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177, 2008. → pages 5
- [35] P. Lieberman. Some Effects of Semantic and Grammatical Context on the Production and Perception of Speech. *Language and Speech*, 6(3):172–187, 1963. URL <http://las.sagepub.com/content/6/3/172.abstract>. → pages 5
- [36] S. L. Mattys and L. Wiget. Effects of cognitive load on speech recognition. *Journal of Memory and Language*, 65(2):145–160, 2011. → pages 5, 6, 25
- [37] L. H. Mayo, M. Florentine, and S. Buus. Age of second-language acquisition and perception of speech in noise. *Journal of speech, language, and hearing research : JSLHR*, 40(3):686–693, 1997. ISSN 1092-4388. → pages 5, 20
- [38] H. Mitterer, O. Scharenborg, and J. M. McQueen. Phonological abstraction without phonemes in speech perception. *Cognition*, 129(2):356–361, 2013. → pages 3, 24
- [39] D. Norris and A. Cutler. The relative accessibility of phonemes and syllables. *Perception & psychophysics*, 43(6):541–550, 1988. ISSN 0031-5117. → pages 6
- [40] D. Norris, J. M. McQueen, and A. Cutler. Perceptual learning in speech, 2003. → pages 1, 2, 3, 4, 6, 8, 24, 25
- [41] J. Pierrehumbert. Word-specific phonetics. *Laboratory phonology*, 2002. URL <http://books.google.com/books?hl=en&lr=&id=N5quZxfmpswC&oi=fnd&pg=PA101&dq=Word-specific+phonetics&ots=3kCkpgcbgl&sig=kMSAYjSD1UDqQr8VOpP0cst4LDE>. → pages 19, 25
- [42] J. B. Pierrehumbert. Exemplar dynamics: Word frequency, lenition and contrast. *Frequency and the emergence of linguistic structure*, pages 137–158, 2001. ISSN 01677373. → pages 7
- [43] M. Pitt and C. Szostak. A lexically biased attentional set compensates for variable speech quality caused by pronunciation variation. *Language and Cognitive Processes*, (April 2013):37–41, 2012. URL <http://www.tandfonline.com/doi/abs/10.1080/01690965.2011.619370>. → pages 4, 6, 8, 9, 15
- [44] M. A. Pitt and A. G. Samuel. Attentional allocation during speech perception: How fine is the focus? *Journal of Memory & Language*, 29(5):611–632, 1990. ISSN 0749596X. → pages 6

- [45] M. A. Pitt and A. G. Samuel. An empirical and meta-analytic evaluation of the phoneme identification task. *Journal of experimental psychology. Human perception and performance*, 19(4):699–725, 1993. ISSN 0096-1523. → pages 4
- [46] M. A. Pitt and A. G. Samuel. Word length and lexical activation: longer is better. *Journal of experimental psychology. Human perception and performance*, 32(5):1120–1135, 2006. ISSN 0096-1523. → pages 4
- [47] E. Reinisch and L. L. Holt. Lexically Guided Phonetic Retuning of Foreign-Accented Speech and Its Generalization. *Journal of experimental psychology. Human perception and performance*, 40(2):539–555, 2013. ISSN 1939-1277. URL <http://www.ncbi.nlm.nih.gov/pubmed/24059846>. → pages 3
- [48] E. Reinisch, A. Weber, and H. Mitterer. Listeners retune phoneme categories across languages. *Journal of experimental psychology. Human perception and performance*, 39(1):75–86, Feb. 2013. ISSN 1939-1277. doi:10.1037/a0027979. URL <http://www.ncbi.nlm.nih.gov/pubmed/22545600>. → pages 3, 9, 11, 12, 15
- [49] E. Reinisch, D. R. Wozny, H. Mitterer, and L. L. Holt. Phonetic category recalibration: What are the categories? *Journal of Phonetics*, 45:91–105, 2014. ISSN 00954470. doi:10.1016/j.wocn.2014.04.002. URL <http://linkinghub.elsevier.com/retrieve/pii/S009544701400045X>. → pages 3, 19, 24
- [50] A. G. Samuel. Phonemic restoration: insights from a new methodology. *Journal of experimental psychology. General*, 110(4):474–494, 1981. ISSN 0096-3445. → pages 2, 4, 5, 7, 8, 20
- [51] A. G. Samuel. Red herring detectors and speech perception: in defense of selective adaptation. *Cognitive psychology*, 18(4):452–499, 1986. ISSN 00100285. → pages 2
- [52] R. Scarborough. Lexical and contextual predictability: Confluent effects on the production of vowels. *Laboratory phonology*, pages 575–604, 2010. URL <http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Lexical+and+contextual+predictability:+confluent+effects+on+the+production+of+vowels#0>. → pages 5, 20
- [53] O. Scharenborg and E. Janse. Comparing lexically guided perceptual learning in younger and older listeners. *Attention, perception & psychophysics*, 75(3):525–36, 2013. ISSN 1943-393X. URL <http://www.ncbi.nlm.nih.gov/pubmed/23354594>. → pages 6
- [54] O. Scharenborg, A. Weber, and E. Janse. The role of attentional abilities in lexically guided perceptual learning by older listeners. *Attention, Perception, & Psychophysics*, pages 1–15, 2014. → pages 6, 19, 25
- [55] D. Shankweiler, W. Strange, and R. R. Verbrugge. Speech and the problem of perceptual constancy. *Perceiving, acting, and knowing: Toward an ecological psychology*, pages 315–345, 1977. → pages 1
- [56] M. Sumner. The role of variation in the perception of accented speech. *Cognition*, 119(1):131–136, 2011. ISSN 0010-0277. doi:10.1016/j.cognition.2010.10.018. URL <http://dx.doi.org/10.1016/j.cognition.2010.10.018>. → pages 7, 9, 25
- [57] M. Sumner and R. Kataoka. Effects of phonetically-cued talker variation on semantic encoding. *The Journal of the Acoustical Society of America*, 134(6):EL485, Dec. 2013. ISSN 1520-8524. doi:10.1121/1.4826151. URL <http://www.ncbi.nlm.nih.gov/pubmed/25669293>. → pages 1
- [58] S. van Linden and J. Vroomen. Recalibration of phonetic categories by lipread speech versus lexical information. *Journal of experimental psychology. Human perception and performance*, 33(6):1483–1494, 2007. ISSN 0096-1523. → pages 3, 24
- [59] L. van Noorden and J. F. Schouten. *Temporal Coherence in the Perception of Tone Sequences*. PhD thesis, 1975. → pages 6

- [60] J. Vroomen, S. van Linden, B. de Gelder, and P. Bertelson. Visual recalibration and selective adaptation in auditory-visual speech perception: Contrasting build-up courses. *Neuropsychologia*, 45(3):572–577, 2007. → pages 1, 2, 3
- [61] J. M. Wolfe and T. S. Horowitz. What attributes guide the deployment of visual attention and how do they do it? *Nature reviews. Neuroscience*, 5(6):495–501, 2004. → pages 5
- [62] Y. Yeshurun and M. Carrasco. Attention improves or impairs visual performance by enhancing spatial resolution. *Nature*, 396(6706):72–75, 1998. ISSN 0028-0836. → pages 7, 19