

# Adaptive Plasticity in Perceiving Speech Sounds

Citation for published version (APA):

Ullas, S., Bonte, M., Formisano, E., & Vroomen, J. (2022). Adaptive Plasticity in Perceiving Speech Sounds. In L. L. Holt, J. E. Peelle, A. B. Coffin, A. N. Popper, & R. R. Fay (Eds.), *Speech Perception, Springer Handbook of Auditory Research* (pp. 173-199). Springer. [https://doi.org/10.1007/978-3-030-81542-4\\_7](https://doi.org/10.1007/978-3-030-81542-4_7)

## Document status and date:

Published: 01/01/2022

## DOI:

[10.1007/978-3-030-81542-4\\_7](https://doi.org/10.1007/978-3-030-81542-4_7)

## Document Version:

Version created as part of publication process; publisher's layout

## Document license:

Taverne

## Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

## General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.umlib.nl/taverne-license](http://www.umlib.nl/taverne-license)

## Take down policy

If you believe that this document breaches copyright please contact us at:

[repository@maastrichtuniversity.nl](mailto:repository@maastrichtuniversity.nl)

providing details and we will investigate your claim.

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/358787845>

# Adaptive Plasticity in Perceiving Speech Sounds

Chapter · February 2022

DOI: 10.1007/978-3-030-81542-4\_7

CITATIONS

0

READS

143

4 authors, including:



**Milene Bonte**

Maastricht University

94 PUBLICATIONS 2,721 CITATIONS

[SEE PROFILE](#)



**Elia Formisano**

Maastricht University

256 PUBLICATIONS 14,655 CITATIONS

[SEE PROFILE](#)



**Jean Vroomen**

Tilburg University

254 PUBLICATIONS 9,268 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Cognitive Neuroscience of Speech production [View project](#)



Human auditory pathway [View project](#)

# Chapter 7 1

## Adaptive Plasticity in Perceiving Speech 2

### Sounds 3

Shruti Ullas, Milene Bonte, Elia Formisano, and Jean Vroomen 4

**Abstract** Listeners can rely on perceptual learning and recalibration in order to make reliable interpretations during speech perception. Lexical and audiovisual (or speech-read) information can disambiguate the incoming auditory signal when it is unclear, due to speaker-related characteristics, such as an unfamiliar accent, or due to environmental factors, such as noise. With experience, listeners can learn to adjust boundaries between phoneme categories as a means of adaptation to such inconsistencies. Recalibration experiments tend to use a targeted approach by embedding ambiguous phonemes into speech or speechlike items, and with continuous exposure, a learning effect can be induced in listeners, wherein disambiguating contextual information shifts the perceived identity of the same ambiguous sound. The following chapter will review current and past literature regarding lexical and audiovisual influences on phoneme boundary recalibration, as well as theories and neuroimaging data that potentially reveal what facilitates this perceptual plasticity. 5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18

**Keywords** Recalibration · Perceptual learning · Speech perception · Phonetic processing · Lexical processing · Audiovisual speech · Speech-reading 19  
20

### 7.1 Introduction 21

Speech perception is seemingly easy and automatic to the listener, and for healthy young listeners, it requires little to no effort to accomplish in most circumstances. While it may appear straightforward, a great deal of variability exists in the quality 22  
23  
24

---

S. Ullas (✉) · M. Bonte · E. Formisano  
Department of Cognitive Neuroscience, Faculty of Psychology and Neuroscience, Maastricht University, Maastricht, The Netherlands

Maastricht Brain Imaging Center, Maastricht University, Maastricht, The Netherlands  
e-mail: [shruti.ullas@maastrichtuniversity.nl](mailto:shruti.ullas@maastrichtuniversity.nl); [m.bonte@maastrichtuniversity.nl](mailto:m.bonte@maastrichtuniversity.nl);  
[e.formisano@maastrichtuniversity.nl](mailto:e.formisano@maastrichtuniversity.nl)

J. Vroomen  
Department of Cognitive Neuropsychology, Tilburg University, Tilburg, The Netherlands  
e-mail: [j.vroomen@tilburguniversity.edu](mailto:j.vroomen@tilburguniversity.edu)

© The Author(s), under exclusive license to Springer Nature  
Switzerland AG 2022

L. Holt et al. (eds.), *Speech Perception*, Springer Handbook of Auditory  
Research 74, [https://doi.org/10.1007/978-3-030-81542-4\\_7](https://doi.org/10.1007/978-3-030-81542-4_7)

25 of the speech signal, which requires the listener to adapt to the novel characteristics  
26 of the encountered speech. The acoustic signal can differ significantly across speak-  
27 ers, often due to unfamiliar accents, the presence of noise, or speech rate. The lis-  
28 tener is able to easily resolve these inconsistencies and understand what is spoken.  
29 No two speakers will pronounce a phoneme in the exact same way, and even the  
30 same speaker may not produce a phoneme identically across multiple instances, yet  
31 listeners are effortlessly able to recognize what speakers are saying. Auditory qual-  
32 ity can also vary within speakers, perhaps due to a cold or while speaking over the  
33 phone. Still, the listener is usually able to easily resolve these inconsistencies and  
34 understand what is spoken. In order to adapt to these irregularities, listeners can  
35 learn to reshape existing representations of speech sounds and categories to accom-  
36 modate any possible variability.

37 Acoustics are not the only source of information capable of changing speech  
38 sound representations, as other contextual cues are also highly influential. Contextual  
39 features may be just as useful as auditory information, and possibly even more so.  
40 Winn (2018) introduces some non-acoustic cues that impact what listeners perceive  
41 to hear, including visual cues, such as the lip movements of a speaker, as well as the  
42 listener's own lexical knowledge. These non-acoustic sources can also enable pro-  
43 cesses known as recalibration and lexically guided perceptual learning. Contextual  
44 information can guide the retuning process of phoneme category boundaries, after  
45 continuous exposure to speech or videos of speechlike tokens, edited to contain  
46 ambiguous versions of a phoneme. Listeners can learn to incorporate these ambigu-  
47 ous sounds into the phoneme category itself, particularly when the sounds resemble  
48 already familiar phonemes.

49 Norris et al. (2003) termed this effect lexically guided perceptual learning, and  
50 observed that with the help of lexical knowledge, listeners could learn to adjust a  
51 perceptual boundary between two phonemes by hearing ambiguous phonemes  
52 embedded into words. Similarly, Bertelson et al. (2003) identified a comparable  
53 effect as recalibration, where listeners utilized visual or speech-reading information  
54 to adjust the perceptual boundary. The two discoveries were made close in time, and  
55 while Norris et al. (2003) used recordings of words as stimuli, Bertelson et al.  
56 (2003) relied on video recordings of syllables. Still, while the types of contextual  
57 information differed between the two studies, the experimental designs and stimuli  
58 constructions were remarkably similar. Since then, in the literature on lexical influ-  
59 ences, the resulting aftereffect is often referred to as perceptual retuning or pho-  
60 neme adaptation, while the studies on visual/speech-reading influences refer to the  
61 analogous effect as audiovisual recalibration.

62 In laboratory settings, recalibration and perceptual learning are typically mea-  
63 sured in two phases, starting with an exposure phase and followed by a test phase  
64 (see Kraljic and Samuel 2009, for an overview). In the approach used to measure  
65 lexically guided perceptual learning, exposure stimuli are composed of audio  
66 recordings of words, whereas in audiovisual recalibration experiments, exposure  
67 stimuli comprise videos of a speaker's lip movements while pronouncing a syllable.  
68 Both types of stimuli contain edited audio, where one particular phoneme is replaced  
69 with an ambiguous sound halfway between two clear phonemes. For instance,

speech stimuli containing /f/ sounds are replaced with a token halfway between /f/ and /s/. Listeners are presented with many examples of such edited stimuli in the exposure phase, with words such as “half” and “paragraph” edited to remove the clear /f/ and replaced with the ambiguous version. Because “half” and “paragraph” are real words in English, whereas “halss” and “paragrass” are not, listeners tend to perceive the ambiguous token as an /f/. During subsequent test phases, listeners hear the ambiguous sounds again, but without any lexical or visual context available, and respond with the phoneme they perceive to be hearing. Consequently, listeners become more likely to respond hearing the same phoneme that was replaced in the previously presented words or videos. In the case of the aforementioned example, the listener would now report hearing the ambiguous token as /f/ as well. This response pattern is understood to reflect recalibration or perceptual retuning, and is a result of the listeners learning to include the ambiguous sound as a part of that particular phoneme category.

Listeners in such experiments can also learn to perceive the same ambiguous phoneme, with no change in acoustic features, in opposing ways, depending on the bias of the surrounding context. A 50–50 /f/-/s/ blend can be learned as either /f/ or /s/ depending on the type of exposure the listener has undergone. Again, in the same example, if listeners were instead presented with speech stimuli that replaced all /s/ sounds with the same ambiguous token (the 50–50 blend of /f/ and /s/), listeners would be more likely to perceive the ambiguous sound as /s/ as well. With this approach, the contributions of visual and lexical information on speech perception can be disentangled from the auditory signal itself, as the exact same ambiguous tokens can be learned as different phonemes depending on the contextual cues. Perceptual retuning and recalibration studies (Bertelson et al. 2003; Norris et al. 2003; Krajlic and Samuel 2009) also reveal how flexible the units of speech are, and how they can be adapted depending on the surroundings of the listener. These experiments illuminate non-acoustic contributions to speech perception, and what listeners rely on in addition to the acoustic signal itself, which, again, tends to fluctuate greatly both within and across speakers.

With the advancement of neuroimaging technologies, the ways in which the brain incorporates these perceptual shifts have been explored with greater detail and have revealed the areas of the brain likely to be involved in these processes. Techniques such as functional MRI (fMRI; see Table 7.1 for abbreviations) and electrocorticography (ECoG) recordings have proven especially useful in elucidating the potential neural mechanisms (Hickok and Poeppel 2007; Mesgarani et al. 2014). These findings, combined with existing theories of speech perception, are useful for understanding how the brain adapts to unclear speech and how the necessary changes may be implemented at the neural level.

This chapter will present an overview of the current literature regarding lexical (Sect. 7.2.1) and audiovisual influences (Sect. 7.3.1) on phoneme boundary recalibration, as well as some related works on selective speech adaptation (Sect. 7.3.2). Changes over time (Sect. 7.2.2), generalization over speakers and sounds (Sects. 7.2.3 and 7.3.3), and other features (Sect. 7.2.4) will also be discussed, as well as a comparison between lexical and audiovisual perceptual learning (Sect. 7.4).

t1.1 **Table 7.1** Table of abbreviations

Abbreviation	Full name	
ECoG	Electrocorticography	t1.2
EEG	Electroencephalogram	t1.3
fMRI	Functional MRI	t1.4
IFS	Inferior frontal sulcus	t1.5
IPL	Inferior parietal lobe	t1.6
ITS	Inferior temporal sulcus	t1.7
MEG	Magnetoencephalogram	t1.8
MTG	Medial temporal gyrus	t1.9
PT	Planum temporale	t1.10
STG	Superior temporal gyrus	t1.11
STS	Superior temporal sulcus	t1.12
SWS	Sine-wave speech	t1.13
		t1.14

115 Theories and neuroimaging studies that may explain the underlying mechanisms of  
 116 recalibration will also be reviewed (Sect. 7.5), followed by a final conclusion and  
 117 summary (Sect. 7.6).

## 118 **7.2 Lexical Knowledge and Auditory Perception**

### 119 **7.2.1 Introduction to Lexically Guided Perceptual Learning**

120 As mentioned earlier in the introduction (Sect. 7.1), top-down lexical knowledge  
 121 can assist listeners in interpreting unclear speech. To investigate this, some research-  
 122 ers have used noise-vocoded or degraded speech stimuli that systematically distort  
 123 frequency and amplitude components of the speech (Davis et al. 2005). Other  
 124 researchers have studied how listeners adapt to accented speech (Clarke and Garrett  
 125 2004; Bradlow and Bent 2008), how listeners adapt to non-native speech in noise  
 126 (Lecumberri et al. 2010), as well as how lexical knowledge supports understanding  
 127 accented speech (Maye et al. 2008). A review by Holt and Lotto (2008) describes  
 128 the various ways in which listeners can build links between acoustic information  
 129 and linguistic representations. Prior to many of these studies, the discovery of what  
 130 is now known as the Ganong effect (Ganong 1980) established a specific influence  
 131 of lexical information on speech sound perception. Ganong (1980) showed that lis-  
 132 teners were likely to report hearing words even when exposed to auditory stimuli  
 133 that were edited to begin with ambiguous sounds. Listeners who heard the word  
 134 “?eep,” where the /?/ sound was acoustically halfway between /d/ and /t/, were  
 135 likely to interpret the stimulus in the form of a word, such as “deep,” rather than  
 136 “teep.” The same held true in the opposite direction, when the same ambiguous  
 137 token replaced /t/ in recordings of words beginning with /t/, such as “?each.” Again,

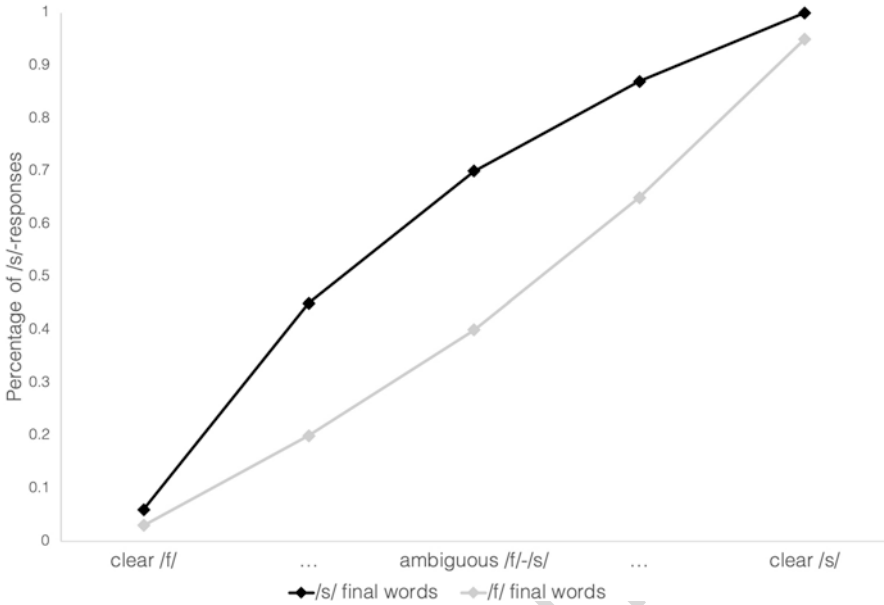
listeners were likely to report hearing a word, such as “teach,” rather than the non- 138  
word version, “deach.” In essence, listeners were not hindered by the unclear audi- 139  
tory information and were still able to infer the intended words. 140

Extending further from the Ganong effect, the findings of Norris et al. (2003) 141  
revealed how lexical information could not only affect perception of speech stimuli 142  
but could also reshape speech sound representations. Native Dutch speakers per- 143  
formed a lexical decision task while listening to audio recordings of Dutch words, 144  
some of which typically ended in /f/, such as “witlo??” (*witlof*, meaning chicory) 145  
and “druif??” (*druif*, meaning grape), where all /f/ sounds were replaced with an 146  
ambiguous token halfway between /f/ and /s/. During the following test phase, 147  
where listeners responded to a continuum of sounds ranging from more /f/-like to 148  
more /s/-like, they were likely to report a significantly greater number of tokens as 149  
/f/ sounding. Another group of participants conducted the same lexical decision task 150  
while hearing words, but in contrast, these words typically contained /s/ (such as 151  
*radijs* and *relaas*, meaning radish and account) and were spliced with the same 152  
ambiguous token in the place of /s/, and the opposite pattern of results was found. 153  
These listeners responded to the same continuum of /f/ to /s/ sounds during the test 154  
phase, and were more likely to report hearing the sounds as /s/. A third control group 155  
heard pseudo-words containing the ambiguous phoneme to test whether the absence 156  
of any lexical information could impact subsequent categorization. This group 157  
showed no bias toward either phoneme during the test phase. An example of the 158  
pattern of results is shown in Fig. 7.1. 159

Together, these results built further upon the lexical effect first described by 160  
Ganong and illustrated how lexical knowledge impacted the participants’ percep- 161  
tion in two ways. First, during the exposure phase, the words containing the ambigu- 162  
ous sounds were still perceived as words and nearly indistinguishable from 163  
unedited words, and replicated the Ganong effect. Then, in the test phase, listeners 164  
categorized ambiguous sounds of a continuum and were prone to hearing the conti- 165  
nuum sounds resembling the phoneme replaced in the prior exposure phase. That 166  
is, listeners were likely to perceive the ambiguous token as /f/ after exposure to 167  
f-final words containing the said token. Thus, phoneme category boundaries were 168  
found to be flexible, as listeners adjusted the boundary between two phonemes 169  
using their lexical knowledge. The authors proposed that the results mirrored what 170  
listeners may be doing in response to an unfamiliar accent, by shifting a category 171  
boundary to make room for the pronunciation of the newly encountered speaker 172  
(this will be discussed more in Sect. 7.2.3). 173

## 7.2.2 Perceptual Learning Over Time 174

Since Norris et al. (2003), later studies of perceptual learning explored the other 175  
attributes of this effect, such as the duration of time for which the retuning effects 176  
could last in the listener, as well as if these changes were permanent or if the catego- 177  
ries returned to their previous state. Kraljic and Samuel (2005) used nearly the same 178



**Fig. 7.1** Example graph of perceptual retuning results. After exposure to edited words, participants are presented with a continuum of sounds ranging from clear /f/ to clear /s/ in a test phase. Participants who hear words typically containing /f/ replaced with an ambiguous /f/-/s/ blend are likely to report hearing /f/ during the test phase (shown in gray), while participants who heard the same sound replacing /s/ in /s/-final words are likely to report hearing more /s/ (shown in black)

179 approach as Norris et al. (2003), testing native English speakers using words con-  
 180 taining either /s/ or /ʃ/ (the “sh” sound in shoe), with items such as *eraser* and *pub-*  
 181 *lisher*. After a 25-minute delay, participants were tested on a continuum from /s/ to  
 182 /ʃ/, and their responses reflected the shift induced by the preceding exposure phase  
 183 (i.e., more /s/ responses after /s/-final words, or more /ʃ/ after /ʃ/-final words). Despite  
 184 the delay, the listeners could still retain the newly learned phoneme boundary position.  
 185 Eisner and McQueen (2006) also measured perceptual learning effects in sub-  
 186 jects after a longer delay, where participants completed one test immediately after  
 187 exposure, and also returned 12 hours after the exposure to complete the test phase  
 188 again. The exposure phase was slightly altered from the original version by Norris  
 189 et al. (2003) and consisted of words with ambiguous segments, all embedded into a  
 190 short story. The potential confound of sleep was also accounted for, as one group  
 191 waited 12 hours during the day to be retested, while another group waited 12 hours  
 192 overnight, and returned for the second test phase after they had slept. Both groups  
 193 still maintained retuning effects after the 12-hour delay, with or without sleeping.  
 194 Perceptual learning is seemingly unaffected by long gaps between exposure and  
 195 test, which suggests that lexically guided perceptual learning is largely stable over  
 196 the order of hours.



### 7.2.3 Generalization of Perceptual Retuning

197

Although lexically driven perceptual learning appears to be quite robust, other investigators have identified the limitations of such learning. For example, perceptual learning tends to be restricted by the stimuli, particularly by the speakers of the tokens. The shift in perception resulting from experience with one phoneme pair by one speaker may not apply to the same pair produced by a new speaker. Eisner and McQueen (2005) had two groups of participants undergo exposure to Dutch words containing either an ambiguous /f/ or /s/ spoken by one speaker, but were tested on a continuum of /f/-/s/ sounds by a different speaker. Participants did not show the retuning effect when tested with the continuum by the novel speaker, so responses to the items on the continuum did not show a shift toward any particular phoneme. Thus, the authors concluded that the participants treated the sounds contained in the exposure stimuli as an idiosyncrasy, so it was tied specifically to the speaker of the ambiguous sounds and did not generalize to ambiguous sounds by a different speaker.

Kraljic and Samuel (2007) also addressed a possible discrepancy in generalization to new speakers based on phoneme types. Listeners who were exposed to words containing ambiguous /d/ or /t/ (plosives or stop consonants) sounds could generalize retuning to the same tokens of a new speaker during the test phase, translating to a shift in categorization responses toward the phoneme replaced in the prior exposure phase (i.e., more /d/ responses after exposure to /d/ words replaced with /d/-/t/ blend). However, those who were exposed to words spliced with ambiguous /s/ or /ʃ/ (fricatives) could not generalize any retuning to a new speaker, so no shift was found in categorization responses during the test phase. Evidently, perceptual learning may not always be constrained by the speaker, and depending on the type of phoneme pair used, it may also be token-specific.

Similarly, generalization to new speakers may also be dependent on the accent of the speaker. Kraljic et al. (2008a) compared effects of speaker characteristics on perceptual learning, with an idiosyncratic pronunciation versus an accent commonly known to the participants. The idiosyncrasy, or speaker-specific version, was designed by placing an ambiguous /s/-/ʃ/ sound before any consonants in the word stimuli, whereas the accented version only placed the ambiguous sound before an occurrence of /tr/ (such as /s/ in *string*), as is typical of some regional American accents. Phoneme boundary retuning was not successful in the latter group that was exposed to the tokens typical of the accented speech, but was detected in the non-accented group. Knowledge of reasonable and unrealistic deviations, which may be implicit or explicit, also seem to impact perceptual learning. In contrast, native English participants who heard exposure stimuli in English by a speaker with a Mandarin accent were more likely to generalize retuning to another acoustically similar Mandarin-accented speaker (Xie and Myers 2017), and to a lesser extent to speakers whose voices were acoustically more distant. The discrepancy in findings between Xie and Myers (2017) and Kraljic and Samuel (2008a) may once again reflect differences in learning effects due to the phoneme pair used.

238  
239

240 Just as speaker specificity of perceptual learning is tied to the type of phoneme  
 241 pairs, the same applies to generalization across phoneme pairs within a single  
 242 speaker. Kraljic and Samuel (2006) saw that perceptual learning could generalize  
 243 between pairs of plosives or stop consonants, particularly between /d/-/t/ and /b/-/p/.  
 244 During the exposure phase, listeners heard words containing either an ambiguous  
 245 /d/ or /t/, but during the test phase, they responded to both a /d/-/t/ continuum and a  
 246 /b/-/p/ continuum. Participants were able to extend retuning to the /b/-/p/ continuum  
 247 in the same direction of voicing, or the point in time at which the vocal folds vibrate,  
 248 where /b/ and /d/ are voiced, whereas /d/ and /t/ are unvoiced. Participants who  
 249 heard words with an ambiguous /b/ were more likely to report a greater amount of  
 250 both /b/ along the /b/-/p/ continuum, as well as more /d/ during an additional test  
 251 phase on a continuum of /d/-/t/. Mitterer et al. (2013) also explored phoneme speci-  
 252 ficity by creating exposure stimuli using Dutch words ending in an approximant /r/  
 253 (the /r/ in red) or dark /l/ (the /l/ in pool). Participants showed retuning effects during  
 254 a test phase with a continuum of the versions of /r/ or /l/ they previously heard dur-  
 255 ing exposure, but could not generalize to other allophones, or phonetic neighbors of  
 256 /r/ and /l/, such as a trill /r/ (not in American English phonology but similar to the  
 257 t-sound in better) or a light /l/ (the /l/ in leaf). Once again, the specificity of recal-  
 258 ibration seems to be dependent on the acoustic features of the phoneme pair being  
 259 learned.

260 Overall, it appears that retuning is often phoneme- and speaker-specific, but con-  
 261 tingent on the specific phoneme pair used. Generalization to a new speaker is more  
 262 likely to occur if the phoneme boundary is adjusted between two plosives and not  
 263 between fricatives. Perceptual retuning effects upon plosives or stop consonants are  
 264 also more likely to extend to other plosives, but, again, are unlikely to do so for  
 265 fricatives or approximants. Acoustic similarity also plays an important role as to  
 266 whether retuning effects can be applied to new sounds.

#### 267 **7.2.4 Other Attributes of Perceptual Learning**

268 Most studies of the lexically guided perceptual learning studies described through-  
 269 out Sect. 7.2 are twofold. They typically start with an exposure phase, with words  
 270 containing one particular ambiguous phoneme, presented along with other filler  
 271 words and pseudo-words. Listeners are also often asked to perform a lexical deci-  
 272 sion task during this exposure phase, in order to maintain their attention. This is  
 273 followed by a categorization task, or the test phase, on a continuum between two  
 274 clear phonemes with the aforementioned ambiguous phoneme in between. However,  
 275 this design is not always used, and other similar designs can still lead to measurable  
 276 retuning effects. McQueen et al. (2006b) concluded that perceptual learning is not  
 277 dependent on a lexical decision task during the exposure phase. Instead, the lexical  
 278 decision task was replaced with a simple counting task, and learning effects  
 279 remained intact. However, a more recent study by Samuel (2016) suggested that  
 280 targeted distractions during exposure that can prevent access to the lexicon are

detrimental to perceptual retuning. In this study, listeners heard two voices only separated by 200 ms during exposure, of words containing an ambiguous /s/-/ʃ/ phoneme by a male speaker, and irrelevant words by a female speaker, and were asked to perform a lexical decision task on the male speaker, or to count the number of syllables spoken by the female speaker. Listeners who attended to the female speaker showed no recalibration during subsequent testing; however, when the voices were separated by 1200 ms, recalibration effects were reinstated. Similarly, listeners were also unable to undergo learning in the presence of background noise (Zhang and Samuel 2015), suggesting that recalibration cannot be performed automatically and requires attentional resources. But attention alone is also not enough to induce retuning, as can listeners still account for potentially transient characteristics of a speaker. In a creative design by Kraljic et al. (2008b), listeners viewed stimuli of a speaker with a pen in their mouth while pronouncing words dubbed with an ambiguous phoneme. These listeners did not show retuning during the subsequent test phase, implying that listeners also acknowledge temporary atypical pronunciations of a speaker before adjusting phoneme representations.

Attention aside, the prototypical test phase, most often a continuum of sounds between two phonemes, is also not a requisite to detect perceptual retuning effects. Effects were still preserved when test phase items were replaced with minimal word pairs ending in an ambiguous phoneme (McQueen et al. 2006a). Participants were then more likely to hear one of the two words of the pair, predicated by the prior exposure phase. For instance, after exposure to words with an ambiguous /f/ (such as *paragraph*, ending with an /f/-/s/ blend), participants were likely to hear “knife” rather than “nice” when presented with “kni-,” ending in the same /f/-/s/ blend. The effect was observed in the opposite direction when listeners were presented with /s/ words ending in the ambiguous token during the exposure. In the same example, listeners were more likely to hear “nice.”

Even fully intact lexical information is not a necessity for retuning to occur, and implicit knowledge of phonotactic information, or the rules within a language regarding allowable phoneme combinations, can be sufficient (Cutler et al. 2008). Here, exposure stimuli were phonotactically valid pseudo-words containing an ambiguous phoneme. Perceptual retuning can also be observed with other known phonemes that are acoustically related, such as /θ/ (represented as theta, or the “th” sound in thing) in place of /s/ or /f/, in place of the oft-mentioned ambiguous phoneme (Sjerps and McQueen 2010). Again, the acoustic or perceptual similarity can determine whether retuning is induced or not.

Thus, the exposure and test phases do not necessarily have to follow one particular procedure for phoneme boundary retuning, but all of the studies discussed within Sect. 7.2, as well as most of the classical studies of lexically driven perceptual retuning, have focused on native listeners. More recent works have also studied non-native listeners, and retuning can take place in these listeners as well. Native Dutch speakers with high proficiency in English also showed perceptual learning effects in response to English stimuli spoken by a British English speaker (Drozdova et al. 2015). Native German speakers of Dutch were also observed to undergo retuning effects in response to Dutch stimuli, at levels comparable to native Dutch speakers

326 (Reinisch et al. 2013). However, proficiency in the second language can also deter-  
327 mine whether recalibration can occur, as a group of native Arabic speakers with  
328 lower English proficiency than another group of native Hebrew speakers showed no  
329 retuning effects with English phonemes, while the latter group did (Samuel and  
330 Frost 2015).

## 331 **7.2.5 Summary of Lexically Driven Perceptual Learning**

332 Section 7.2 summarized the seminal studies as well as some more recent findings  
333 about lexically guided perceptual learning. These effects are potentially long-lasting  
334 but may not generalize to new speakers. Non-native speakers are also capable of  
335 demonstrating learning effects, but this may be mitigated by the listener's profi-  
336 ciency in the second language. Generalization to new speakers and to other pho-  
337 nemes is mitigated by the type of phoneme category being adjusted. Retuning  
338 effects may be applied from stop consonants or plosives to other phonemes within  
339 this classification, but this is less likely for fricatives or approximants. While lexical  
340 knowledge is primarily driving the subsequent learning, acoustic features still place  
341 constraints on what can and cannot be extended to other speech sounds.

## 342 **7.3 Audiovisual Information and Speech**

### 343 **7.3.1 Overview of Audiovisual Recalibration**

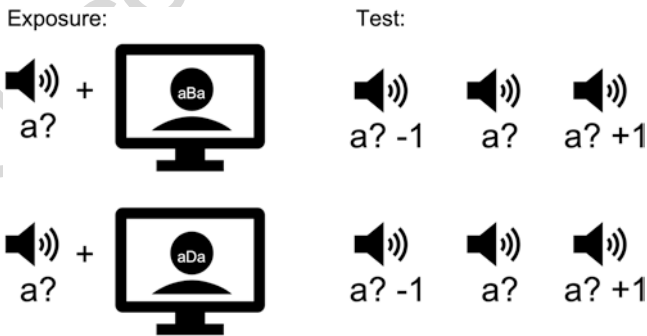
344 Visual or speech-read information, much like lexical information, can also provide  
345 clarity when the available acoustics are unclear. Speech-reading can be relied upon  
346 if noise is present (Sumbly and Pollack 1954), and also significantly alter what lis-  
347 teners perceive to hear. McGurk and MacDonald (1976) made the groundbreaking  
348 discovery that participants who viewed videos of a speaker pronouncing the syllable  
349 /gaga/, dubbed with audio of the syllable /baba/, perceived an entirely new percept,  
350 and reported hearing /dada/. Bertelson et al. (2003) extended this finding, and  
351 detected aftereffects on categorization responses following exposure to McGurk-  
352 like stimuli. Again, not only did speech-reading influence the perception of incon-  
353 gruent audiovisual tokens, but continuous exposure led to responses biased by the  
354 visual/speech-reading information. Much like the approach used by Norris et al.  
355 (2003) described in Sect. 7.2, participants first underwent an exposure phase, where  
356 they viewed audiovisual stimuli of a speaker's lip movements while pronouncing /  
357 aba/, dubbed with audio of an ambiguous phoneme halfway between /aba/ and /  
358 ada/. During a subsequent test phase, participants only heard the audio token of the  
359 ambiguous phoneme and its two neighbors from a continuum, and were more likely  
360 to report them as /aba/ sounding. Unlike Norris et al. (2003), a within-subjects

design was used, and the same group of participants also viewed videos of the speaker pronouncing /ada/, but dubbed with the same ambiguous token. In this case, participants were more likely to report hearing the token as /ada/ during the test phase (Fig. 7.2).

In a follow-up experiment, listeners were exposed to congruent stimuli, or clear audio of /aba/ combined with lip movements of /aba/, and the same for an audio and video combination of /ada/. These unambiguous stimuli showed the reverse effect of the recalibration experiment and led to selective speech adaptation (Eimas and Corbit 1973). As a result of said selective speech adaptation, participants made fewer /aba/ responses to the ambiguous sounds if exposed to clear /aba/ tokens, and similarly gave fewer /ada/ responses after exposure to clear /ada/ tokens. This response is unlike recalibration, where participants who listen to ambiguous sounds during the exposure phase then become more likely to report hearing the phoneme being biased for by the lip movements of the speakers (i.e., ambiguous audio coupled with video of /aba/ leading to more /aba/ responses during the test phase). Selective speech adaptation will be discussed in more detail in Sect. 7.3.2.

### 7.3.2 Audiovisual Recalibration and Selective Speech Adaptation

Prior to studies of audiovisual recalibration, a perceptual learning effect known as selective speech adaptation was discovered (Eimas and Corbit 1973) and has also been helpful for understanding the building blocks of speech perception. Recalibration and selective speech adaptation share considerable overlap, especially in terms of their experimental design, but are also distinct in their interpretations. Both styles of experiments use a similar two-part procedure with an exposure and



**Fig. 7.2** A typical audiovisual recalibration procedure. Exposure phases pair ambiguous phoneme blends (such as an /aba/-/ada/ blend) with video of a speaker pronouncing one of the two phonemes (/aba/ or /ada/). Following exposure to these videos, listeners are then presented with the auditory items (the ambiguous /aba/-/ada/ blend, along with other similar sounds) and asked to respond with what they hear

385 test phase. Unlike recalibration, which typically uses ambiguous sounds, selective  
386 speech adaptation relies on exposure to clear sounds. While recalibration experi-  
387 ments lead to an increase in responses of the phoneme indicated by the videos dur-  
388 ing exposure, selective adaptation results in a reduction. For example, listeners  
389 repeatedly exposed to tokens of a clear /ba/ become less likely to perceiving /ba/  
390 when given a categorization task on a /ba-/da/ continuum. Selective speech adapta-  
391 tion is thought to reflect a fatigue effect, where listeners become desensitized to the  
392 auditory token during the exposure phase. The listener then becomes more sensitive  
393 to the acoustic differences in other similar sounds, thereby reports hearing the  
394 ambiguous tokens as the phoneme opposing the preceding exposure phase. The  
395 original study of selective speech adaptation (Eimas and Corbitt 1973) relied on  
396 solely auditory stimuli, but later studies measured the same effects when exposure  
397 stimuli were coupled with videos of a speaker’s lip movements, as Bertelson et al.  
398 (2003) reported. These unambiguous, or congruent, audiovisual stimuli also led to  
399 fewer responses of the phoneme presented in the test phase, as described in  
400 Sect. 7.3.1.

401 Selective speech adaptation and recalibration are often discussed together, as  
402 they both reflect a change in auditory perception, following an exposure phase to  
403 syllables or speech sounds. Just as the response patterns of the two phenomena go  
404 in opposite directions, the two differ in numerous other ways as well. Vroomen and  
405 colleagues have compared an audiovisual form of selective speech adaptation to  
406 recalibration and have found that the overall buildup and dissipation also tend to  
407 differ (Vroomen et al. 2006). The number of exposure trials has been found to share  
408 a log-linear relationship with selective speech adaptation, as the effect was observed  
409 to increase as exposure trials accumulate, whereas recalibration was found to have  
410 a curvilinear relationship in relation to the number of exposure trials, as it steadily  
411 increased until eight exposure trials, but reduced with additional exposure.  
412 Recalibration and selective speech adaptation are also differentially affected by the  
413 number of test trials, as audiovisual recalibration effects are short-lived and can be  
414 present only up until approximately 6 test trials, while selective speech adaptation  
415 effect can be continuously sustained for up to 60 test items (Vroomen et al. 2004).

416 Sine-wave speech (SWS) is constructed by starting from clear speech but stripped  
417 down until approximately three sinusoids that follow the central frequency and  
418 amplitude of the first three formants remain (Remez et al. 1981). These stimuli are  
419 often unintelligible unless listeners are explicitly told that the sounds have been  
420 extracted from actual speech. Vroomen and Baart (2009b) also compared recalibra-  
421 tion and selective speech adaptation in groups that viewed audiovisual SWS tokens  
422 as speechlike versus non-speechlike. In this experiment, all of the ambiguous and  
423 clear sounds typical of recalibration and selective speech adaptation studies were  
424 replaced with SWS versions, so a continuum including and between two clear pho-  
425 nemes was converted into SWS. For exposure phases, these SWS sounds were still  
426 paired with videos of a speaker’s corresponding lip movements, but were presented  
427 without video for test phases. One “speech-mode” group viewed ambiguous SWS  
428 tokens paired with videos, which identified the tokens as /onso/ or /omso/, and  
429 showed recalibration effects. A “non-speech-mode” group viewed the same stimuli

but categorized the ambiguous SWS tokens as “1” or “2,” and did not show a recalibration effect, so a “speech mode” did impact any possible recalibration. In contrast, for selective speech adaptation, participants viewed videos coupled with endpoint SWS tokens (rather than ambiguous), and adaptation effects were observed. In this instance, listeners who performed a categorization test on SWS versions of the ambiguous tokens heard them as the opposite phoneme to the one biased for by the preceding exposure (i.e., hearing more /omso/ after exposure to SWS versions of a clear /onso/ paired with video). Selective speech adaptation was still measurable in another non-speech-mode group, who underwent the same types of exposure, but categorized the subsequent test phase ambiguous sounds as 1 or 2. Essentially, selective speech adaptation was unaffected by either set of labels, so speech mode had no impact on perception and listeners still adapted accordingly. The awareness of speechlike qualities was crucial for successful recalibration, but selective speech adaptation was not hindered by this lack of this awareness. While recalibration and selective speech adaptation can reshape speech sound representations, based on these comparisons, it appears the two may be controlled by distinct but related substrates. The authors concluded that audiovisual recalibration may emerge from speech and language networks, while selective speech adaptation is purely a bottom-up process that does not require higher-level feedback. Potential neural mechanisms will be discussed in more detail in Sect. 7.5.

### 7.3.3 Specificity of Audiovisual Recalibration

Whether recalibration can be generalized has been addressed with regard to audiovisual information as well, just as it has with lexical context. While recalibration is robust enough to not depend on working memory (Baart and Vroomen 2010), audiovisual recalibration tends to be token-specific (Reinisch et al. 2014), as exposure to either visual /aba/ or /ada/ tokens dubbed with ambiguous audio had no effect on listeners’ categorization of continua of either /ibi-/idi/ or /ama-/ana/ sounds during test. Therefore, audiovisual recalibration appears to be constrained by the acoustics features, as learning could not extend to other phonemes, or even to the same phonemes paired with different vowels. The ear itself can limit recalibration (Keetels et al. 2016a, b), as the effect was optimal if exposure and test stimuli were presented into the same ear, but was diminished for test stimuli presented into the opposite ear, and locations in between resulted in a gradient of responses as the presentations moved further away from the original ear. The authors argue that this is further evidence that recalibration is strongly tied to the token and context, and the encoding process even accounts for the exact location of the presented sound (neural mechanisms will be addressed further in Sect. 7.5). Notably, listeners also have the capacity to recalibrate each ear in opposite directions using the same ambiguous sounds, e.g., one ear recalibrated toward /aba/, the other toward /ada/, with test sounds presented into the corresponding ears of the exposure phase (Keetels et al. 2015). Thus, phoneme representations may not be completely

471 abstracted from the input received and can retain speaker- and context-specific  
 472 details. Keetels et al. (2015) argue that this could be due to the perceptual system  
 473 striking a balance between generalizing too often and too rarely. If recalibration is  
 474 employed when speech is unclear, then it is may be only necessary to apply the  
 475 newly learned boundary position to other instances that are similar both in acoustic  
 476 and contextual features, so as to not unnecessarily overgeneralize.

477 While audiovisual recalibration may be restricted in some respects, it is not nec-  
 478 essarily specific to the speaker, as listeners can recalibrate to another speaker's pro-  
 479 nunciation of the same phoneme, although to a substantially lesser extent compared  
 480 to the speaker during exposure (van der Zande et al. 2014). Recalibration is gener-  
 481 ally maximal in response to the sound used during exposure, which suggests that it  
 482 generally tends to be constrained by the acoustic features of the exposure sound.  
 483 Similarly, audiovisual recalibration is most often tested with consonant contrasts,  
 484 but Franken et al. (2017) have found that recalibration is possible with a vowel con-  
 485 trast pair of /e/-/ø/. In addition, recalibration with a vowel pair and multiple speakers  
 486 has also been observed, wherein the gender identity of the speakers combined with  
 487 the visual cue indicated by the speech-reading information influenced listeners' cat-  
 488 egorization responses (Burgering et al. 2020).

489 The majority of the studies described have also been centered on adults, but  
 490 audiovisual recalibration can also be adopted early in life and has been observed in  
 491 children as young as 8 years old. Van Linden and Vroomen (2008) measured recal-  
 492 ibration effects in two groups of children and determined that children at 8 years old  
 493 could recalibrate with audiovisual stimuli, but children at 5 years old could not, so  
 494 the ability may be developed within this window of 3 years. Dyslexia does not  
 495 restrict the effect either (Baart et al. 2012), as adults with dyslexia were compared  
 496 with fluently reading adults, and the dyslexic group showed no deficit in their ability  
 497 to recalibrate. Even children with dyslexia are capable of undergoing recalibration  
 498 driven by text (Romanovska et al. 2019), even though children with dyslexia often  
 499 experience difficulties in speech-reading and letter-speech sound mappings  
 500 (Snowling 1980; van Laarhoven et al. 2018).

### 501 **7.3.4 Summary of Audiovisual Recalibration**

502 Section 7.3 described audiovisual recalibration, originally described by Bertelson  
 503 et al. (2003), and its various attributes. Later studies by Vroomen and colleagues  
 504 have established the general buildup and dissipation, as well as similarities and dif-  
 505 ferences with another perceptual learning effect, called selective speech adaptation.  
 506 Audiovisual recalibration tends to both build up following a few exemplars during  
 507 exposure and diminish with increasing numbers of test items as well. In contrast,  
 508 selective speech adaptation requires much longer exposure phases, but subsequent  
 509 effects can last for longer durations. Recalibration also tends to be token- and  
 510 context-specific, even to the extent that listeners can recalibrate each ear in opposite  
 511 directions. It also does not easily generalize to other speakers, phonemes, or other



similar instances of the same phoneme, so it is considerably restricted by the acoustic features present during exposure. Nevertheless, it has shown to be utilized by a variety of listeners, including children and adults with dyslexia, and remains to be a helpful tool for listeners when the auditory signal is inadequate.

## 7.4 Comparison of Audiovisual Recalibration and Lexical Retuning

Sections 7.2 and 7.3 have discussed audiovisual recalibration and lexical retuning separately, but the two processes also share many common attributes. In realistic situations, listeners are likely to encounter lexical and visual information simultaneously, so it is possible that these two sources may interact while influencing speech perception. The designs of the two types of experiments share overlap in many respects, with exposure phases consisting of stimuli embedded with ambiguous phonemes, followed by forced-choice test phases where the ambiguous sounds are presented without lexical or speech-reading contextual cues. Even the response patterns between the two original studies by Bertelson et al. (2003) and Norris et al. (2003) paralleled each other, so it may appear that phoneme categories are affected comparably by both audiovisual and lexical information. Brancazio (2004) probed the influence of lexical and speech-reading information in audiovisual speech perception but found that speech-reading exerted a stronger influence on phoneme categorization. Audiovisual effects were similar irrespective of faster and slower response times, while lexical information showed a weaker effect overall and was associated with slower responses.

Based on this, van Linden and Vroomen (2007) proposed that audiovisual information may induce recalibration more effectively than lexical cues, and conducted a study comparing lexical and audiovisual recalibration to test this hypothesis. Two forms of recalibration were compared in native Dutch speakers using a /p/-/t/ phoneme contrast. One group was exposed to lexical stimuli, which consisted of audio Dutch words typically ending in either /op/ or /ot/ (such *bioscoop*, or movie theater, and *idoot*, or idiot), with all endings replaced by an ambiguous token halfway between /op/ and /ot/. Another group was exposed to audiovisual stimuli, comprised of videos of pseudo-words, where lip movements indicated a /op/ or /ot/ ending, and were also dubbed with audio of the ambiguous phoneme at the end of the token. Participants were also exposed to both /op/- and /ot/-biased stimuli, to explore whether they could recalibrate in both directions of the phoneme pair, such that half of the exposure blocks would induce a bias toward /p/, and the remaining half were biased toward /t/. Test phase judgments indicated that recalibration was indeed successful in both groups and in response to both phonemes as well. As the authors originally proposed, audiovisual information was largely more effective in producing recalibration than lexical information. The discrepancy may have resulted from the inherent differences in the stimuli and the processing levels affected, as lexical

552 information might only induce a phoneme preference with the help of top-down  
553 influences, whereas the incoming audiovisual information already contained a  
554 visual bias toward one phoneme. Theories of top-down and bottom-up processing  
555 will be discussed in more depth in Sect. 7.5.

556 In contrast to previous studies on lexical retuning, both audiovisual and lexical  
557 recalibration dissipated at the same rate. Although audiovisual recalibration has  
558 been known to dissipate relatively quickly (Vroomen et al. 2007b), other studies  
559 have found that lexically guided perceptual learning can be long-lasting (Eisner and  
560 McQueen 2006). Participants in the van Linden and Vroomen (2007) study were  
561 flexibly adjusting the phoneme boundary back and forth between the two phonemes,  
562 throughout the duration of the experiment, so the faster dissipation of lexical recal-  
563 ibration may have resulted from constant switching between the two phonemes.  
564 However, this was refuted in a follow-up experiment with a between-subjects  
565 design, where each group of participants were only exposed to one phoneme-  
566 modality combination, and no improvements to recalibration were found. Still, the  
567 chosen phoneme pair is also worth noting, as plosives or stop consonants such as /p/  
568 and /t/ may be more amenable to adjustment than fricatives (as mentioned in Sect.  
569 7.2), such as /f/ and /s/ (Kraljic and Samuel 2007). Overall, lexical and audiovisual  
570 recalibrations seem to be markedly similar, although the pathways supporting them  
571 may not be identical, and may only overlap.

572 The two types of retuning also tend to differ in their stability, as lexical retuning  
573 has been shown to be stable over time, but audiovisual recalibration can be more  
574 susceptible to decay with the passage of time. After a standard exposure phase,  
575 participants were tested after a 24-hour gap and effects had dissipated (Vroomen  
576 et al. 2007a), even if participants were tested both immediately after the exposure  
577 phase and again 24 hours later (Vroomen and Baart 2009b). Audiovisual recalibra-  
578 tion effects have also been shown to diminish within the test phase, as responses that  
579 corresponded with the preceding visual exposure (such as /b/ responses after view-  
580 ing /aba/ videos) were maximal at the start of the test phase, but consistently  
581 decreased as the test phase progressed (Vroomen and Baart 2009b). In contrast,  
582 lexical retuning effects can be preserved throughout longer testing sessions, often  
583 containing approximately 30 test items (Kraljic and Samuel 2009), or up to 12 hours  
584 later (Eisner and McQueen 2006). As mentioned earlier in Sect. 7.2, lexical retuning  
585 is capable of generalizing to new speakers and certain phonemes, while audiovisual  
586 recalibration is most often token-specific and may generalize if the critical pho-  
587 nemes are plosives/stop consonants.

588 More recently, studies comparing audiovisual recalibration and lexical retuning  
589 within both a single session and the same participants have found that the resulting  
590 effects were similar between the two, with similar patterns of dissipation as well  
591 (Ullas et al. 2020a). The simultaneous presentation of both audiovisual and lexical  
592 information within exposure (i.e., listeners presented with videos of words edited to  
593 contain an ambiguous final phoneme) also showed effects comparable to audiovi-  
594 sual recalibration alone, suggesting that the combination leads to no benefit in sub-  
595 sequent phoneme boundary retuning as a result of differences in the pathways  
596 involved in the two forms of perceptual learning (Ullas et al. 2020b). Overall,

lexical retuning and audiovisual recalibration share many similarities in terms of 597  
how the subsequent effects are exhibited, how the experiments measuring them are 598  
designed, as well as the resulting response patterns to presentations of ambiguous 599  
sounds. Both approaches are useful for adapting to speech in noise, even if their 600  
origins and functions may differ. 601

## 7.5 Theoretical and Neural Explanations of Recalibration 602

### 7.5.1 Theories of Speech Perception 603

The mechanisms that enable the auditory system to adjust phoneme boundaries are 604  
often debated. Numerous theories of speech perception have been invoked in expla- 605  
nations of recalibration and perceptual retuning as well. Cutler, McQueen, Norris, 606  
and colleagues (Norris et al. 2000) originally proposed a feed-forward model of 607  
speech perception called Merge and argued that listeners can retune phoneme cate- 608  
gories through a bottom-up abstraction process, which does not rely upon online 609  
feedback from the lexicon, not unlike the COHORT model which also states that 610  
word recognition primarily relies on bottom-up processes (Gaskell and Marslen- 611<sup>AU3</sup>  
Wilson 1997). COHORT presents a modular, unidirectional explanation, where 612  
word recognition is initiated first by acoustic information, triggering a possible 613  
“cohort” of matches, and later, other features such as context and semantics allow 614  
the listener to narrow down the possibilities. Similarly, according to the Merge 615  
model, top-down feedback during speech recognition and phoneme categorization 616  
is not essential, so recognition and categorization operate at a pre-lexical level. 617  
Feedback during categorization could be time-consuming or lead to misinterpreta- 618  
tions of the input, so interactions between lexical and pre-lexical processing would 619  
not be beneficial. Phonemic decisions can be made based on both lexical and pre- 620  
lexical information but do not necessitate interactions between the processes. Cutler 621  
et al. (2010) also emphasized that perceptual retuning cannot be explained purely by 622  
episodic information and that abstraction from such events must be involved as 623  
well. A more recent model by Norris et al. (2016) has been updated to include pre- 624  
dictions of perception based on Bayesian inference, but still does not rely upon 625  
online feedback during phoneme processing. Acoustic information and lexical 626  
knowledge are combined to calculate probable phonemes, but again, the two pro- 627  
cesses are not proposed to interact. 628

Others have described top-down (Davis et al. 2005; Davis and Johnsruide 2007) 629  
and bidirectional influences on speech perception (McClelland and Elman 1986; 630  
McClelland et al. 2006). A classical, interactive model of speech perception, 631  
TRACE (McClelland and Elman 1986), derives its name from a structure called 632  
“The Trace,” a perceptual processing tool. McClelland and Elman proposed that 633  
top-down feedback modulates connections between three layers, from words, to 634  
phonemes, down to features. Phoneme identification can be influenced by lexical 635

636 and speech-reading contexts, and can also be improved through experience.  
 637 According to TRACE, this influence is due to feedback from higher levels of pro-  
 638 cessing. Similarly, McClelland et al. (2006) contend that both top-down and bot-  
 639 tom-up information streams are essential for speech perception. Phoneme  
 640 representations can be influenced by both lexical and acoustic features, and  
 641 vice versa.

642 While most classical theories of speech perception have not accounted for the  
 643 role of visual information, more recently, Kleinschmidt and Jaeger (2011) have put  
 644 forth a belief-updating model based on Bayesian inference, by using data from  
 645 previous studies of recalibration and selective speech adaptation to calculate probabili-  
 646 ties of outcomes. This model, called the Ideal Adaptor Framework, is tailored to  
 647 explain audiovisual recalibration and selective speech adaptation. As described in  
 648 Sect. 7.3.2, audiovisual recalibration and selective speech adaptation are two forms  
 649 of perceptual learning, but their response profiles are in direct contrast. According  
 650 to the Ideal Adaptor Framework, both recalibration and selective speech adaptation  
 651 are described as forms of statistical learning, as a result of exposure to various dis-  
 652 tributions of phonemes. Listeners can create speaker-specific models of phoneme  
 653 categories which allow for initial speaker-level adaptation, but can eventually gen-  
 654 eralize to more speakers with additional experience and if they are also acoustically  
 655 close. The authors also posit recalibration and selective speech adaptation as two  
 656 response patterns along a continuum ranging from ambiguous to prototypical  
 657 sounds. As mentioned earlier in Sect. 7.2.2, recalibration effects tend to peak after  
 658 approximately eight exposure tokens and slowly diminish with additional expo-  
 659 sures, while selective speech adaptation tends to continuously build in a linear man-  
 660 ner with increasing exposure. According to the model, recalibration reflects a  
 661 response to ambiguous sounds, but with increasing amounts of exposure tokens and  
 662 as speech sounds become more prototypical, selective adaptation effects can be  
 663 observed.

AU4

## 664 7.5.2 *Neural Basis of Recalibration and Perceptual Learning*

665 While theoretical frameworks and models have been useful in understanding recal-  
 666 ibration and retuning, neuroimaging studies have shed additional light on areas of  
 667 the brain where these changes occur and how they might explain the levels of pro-  
 668 cessing involved. More general models of speech perception drawn from neuroim-  
 669 aging data and primate studies (Scott and Johnsrude 2003; Rauschecker and Scott  
 670 2009) have described the hierarchical and topographic nature of processing in the  
 671 auditory cortex and surrounding areas.

672 Hickok and Poeppel (2007) proposed the dual-stream processing model of  
 673 speech, with certain features equivalent to those found in visual-processing models.  
 674 According to the model, areas of the brain along a ventral pathway, including medial  
 675 temporal gyrus (MTG) and inferior temporal sulcus (ITS), are geared toward con-  
 676 necting phonological and lexical representations, while regions along a dorsal

pathway, including parietal-temporal, (pre)motor, and inferior frontal regions, are geared toward connecting phonological with sensorimotor and articulatory representations. Adank and Devlin (2010) also explored how listeners adjust to recordings of unclear sentences and found activation patterns consistent with the Hickok and Poeppel (2007) model. Jäncke et al. (2002) also identified structures of the brain in the planum temporale (PT) and middle superior temporal gyrus (STG) that are specific to phoneme perception. STG and the primary auditory cortex can also encode fine-tuned phonetic information (Mesgarani et al. 2008, 2014), with evidence for speaker-invariant phoneme representations distributed across both of these regions (Formisano et al. 2008; Bonte et al. 2014). Other regions implicated in categorical perception of speech sounds include the inferior frontal gyrus (Rogers and Davis 2017) and the supramarginal gyrus (Raizada and Poldrack 2007; see Davis and Johnsruide 2007 for a review).

While these studies paved the way toward delineating a network of regions possibly implicated in recalibration, they may still be insufficient, as this process relies on the integration of both acoustic and contextual information, which are often lexical or visual. In light of this, Obleser and Eisner (2009) proposed a model of pre-lexical abstraction based on prior neuroimaging studies of speech perception, reminiscent of the Merge model (with similarities to TRACE as well, but this model focuses on word recognition and not on abstraction). Pre-lexical abstraction may appear to resemble recalibration, but it also implies that the phoneme representation can be fully disentangled from the acoustic input and thereby abstracted. Pre-lexical abstraction could be implemented probabilistically, primarily along the STG, resulting in phoneme likelihoods rather than definitive phoneme identification. Likelihoods could be calculated by weighing various acoustic features, first processed by primary auditory cortex, and could be updated with talker and context-specific information. Similarly, Holdgraf et al. (2016) have found evidence for acoustic updating, using spectro-temporal receptive field mapping on ECoG recordings of the auditory cortex. Responses of cortical populations were observed to have increased sensitivity to speechlike spectro-temporal features of degraded speech, after exposure to intact speech. This sensitivity could reflect how listeners encode rudimentary acoustic features that also allow the listener to interpret less intelligible speech, or how listeners “fill in the gaps.”

The merits of these models of speech perception can be reexamined in light of fMRI studies of recalibration and retuning. Kilian-Hütten et al. (2011b) had participants undergo audiovisual recalibration using the classic /aba/-/ada/ stimuli while fMRI data was collected. It was discovered that a higher-order network of areas in and around the auditory cortex, including bilateral inferior parietal lobe (IPL), inferior frontal sulcus (IFS), superior temporal sulcus and superior temporal gyrus (STS/STG), and posterior MTG, were all active in recalibration. These areas showed overlapping activation during both the exposure phase and the subsequent test phase. These regions are also known to be involved in audiovisual integration and constructive processes, which would account for their increased activation during recalibration. Kilian-Hütten et al. (2011a) were also able to investigate audiovisual recalibration using MVPA, or multivariate pattern analysis, a technique using fMRI

722 data to train an algorithm to recognize differences in patterns of brain activity. They  
723 were successfully able to decode whether a participant perceived /aba/ or /ada/  
724 while presented with the ambiguous sounds during the test phase of the same audio-  
725 visual recalibration experiment, solely using the activation patterns. Active clusters  
726 were found in and around left PT and left Heschl's gyrus and sulcus, which are typi-  
727 cally viewed as low-level auditory areas, but they may have been influenced by  
728 information other than rudimentary acoustics features as they effectively predicted  
729 the percepts that were driven by the visual cue and not the auditory informa-  
730 tion alone.

731 More recently, Lüttke et al. (2016) investigated a form of adaptation induced by  
732 McGurk-style adaptors with fMRI. Exposure to McGurk adaptors, or clear auditory  
733 /aba/ paired with video of /aga/, resulted in the percept of /ada/. These stimuli led to  
734 an effect much like selective speech adaptation, where follow-up presentations of  
735 clear auditory /aba/ were incorrectly perceived as /ada/ as a result. This mistaken /  
736 ada/ percept showed closely related neural patterns to those elicited by correctly  
737 perceived auditory /ada/, and more so than to patterns associated with correct per-  
738 ception of clear /aba/ tokens. Again, neural activations echoed a shift in auditory  
739 perception due to adaptation through contextual cues.

740 fMRI has also been used to explore lexically driven perceptual learning and other  
741 related phenomena. Activation in posterior left STG and STS has been recorded in  
742 listeners receiving instructions to switch from an acoustic mode to speech mode  
743 while listening to SWS stimuli (Dehaene-Lambertz et al. 2005). While stimuli  
744 remained the same, instructions alone could induce a shift in both perception and  
745 the resulting activation patterns. Similarly, activity in left pSTS has also been asso-  
746 ciated with identification of nonphonemic, short-term sound categories, while left  
747 mSTS may store long-term representation of phoneme patterns already known to  
748 the listener (Liebenthal et al. 2010). Myers and Blumstein (2008) investigated the  
749 Ganong effect (described in Sect. 1.1), or the impact of lexical knowledge on per-  
750 ception of ambiguous speech tokens. Participants heard auditory items with ranging  
751 voice onset time (VOT) from *gift* to *kift* (i.e., word to nonword) and another con-  
752 tinuum ranging from *giss* to *kiss* (from nonword to word). Activity in STG was  
753 modulated by the lexical effect, such that boundary tokens that were perceived as  
754 words showed higher activations compared to acoustically similar tokens from the  
755 other continuum that were not perceived as words. As STG was engaged in both  
756 phonological and lexical processing, the authors suggested that this was evidence in  
757 support of top-down models similar to TRACE that accommodate higher-level  
758 information during processing (Liebenthal et al. 2010).

759 Similarly, Myers and Mesite (2014) tested participants in a classic lexically  
760 guided perceptual retuning experiment with the addition of fMRI, alternating  
761 between exposure phases containing edited words ending in an ambiguous pho-  
762 neme, followed by a forced-choice test phase on a continuum of the same ambigu-  
763 ous sounds. Participants were separated into two groups with the stimuli biased  
764 toward /s/ for one group, and toward /ʃ/ (the “sh” in shop) for the other. Behavioral  
765 results indicated a boundary shift, so over the course of the successive test phases,  
766 participants' perception of the ambiguous /s/-/ʃ/ phoneme had changed. Increased

activity in left IFG and STG was measured with boundary shifted items. These items reflected the perceptual shift, and were categorized as the biasing phoneme in test blocks following the exposure, but not during the earlier blocks at the start of the experiment. Activity both within the auditory cortex and in higher-level cognitive areas suggests that top-down information may have influenced the learning process and may also have been responsible for creating connections between phonetic information and the speaker. Together, the results of these two studies of lexical context imply that perceptual learning involves areas responsible for both lower and higher levels of information processing in resolving the perception of these sounds. However, it remains unclear as to whether the flow of information is simply feed-forward or not, as the exact timing as to when each region is engaged is not yet understood. The authors suggest that initial processing of the unclear sounds relies on higher-level executive regions, but once the listener undergoes sufficient training and has shifted the perceptual boundary, then regions responsible for lower levels of processing, such as STG, can be activated in response to the ambiguous sound.

Combined magnetoencephalogram (MEG) and electroencephalogram (EEG) data have also confirmed that activity in STG reduced over time, as participants learned to improve in identification of degraded speech sounds combined with matching text (Sohoglu and Davis 2016). Furthermore, the results were framed within a model of predictive coding, not unlike Bayesian inference, such that the listener learns to reduce prediction errors as a consequence of learning. STG is proposed to process acoustic features and receives predictions of phonological categories from higher-level frontal areas, and predictions are continuously updated with experience.

While many of the studies discussed thus far have identified STG to be involved in perceptual learning or recalibration, a recent study has also found evidence from the cerebellum (Guediche et al. 2015). Listeners learned to identify words distorted by noise vocoding, and consequently, cerebellar regions showed changes, as well as functional connections to cortical language and auditory regions. Stemming in part from this finding, another model of speech adaptation has been proposed, also relying on a predictive coding mechanism, but supervised by the cerebellum (see Guediche et al. 2014, for a complete review). In contrast, some areas of the brain may be uniquely engaged by either recalibration or retuning. When compared directly using fMRI within the same participants, audiovisual recalibration and lexical retuning showed largely similar areas of activation, over temporal, parietal, and motor cortex areas, although audiovisual recalibration specifically seems to retrigger activation within areas of the visual cortex, despite the lack of visual stimuli during the recalibration test trials (Ullas et al. 2020).

### 806 7.5.3 Summary of Theories of Speech Perception

807 Section 7.5 detailed various theories of speech perception as well as supporting  
808 neuroimaging data that propose the channels through which recalibration and per-  
809 ceptual retuning may operate. Proponents of these speech perception theories have  
810 debated the nature of how phoneme categories can be reshaped, as some argue that  
811 this is a unidirectional, bottom-up abstraction process (Merge, COHORT), while  
812 others postulate that both top-down and bottom-up processes contribute (TRACE).  
813 Theories incorporating distributional and statistical learning, such as the Ideal  
814 Adaptor Framework (Kleinschmidt and Jaeger 2011), have also been useful for  
815 understanding how listeners adapt to variability. Neuroimaging data suggest that  
816 both top-down and bottom-up influences are involved, based on the areas of the  
817 brain that tend to be active during perception of ambiguous tokens, such as STS/  
818 STG and IFS/IFG. Sophisticated analysis techniques such as MVPA have also been  
819 useful for pinpointing specific patterns of neural activity associated with the shifts  
820 in perception, but the directionality of influences upon these percepts remains  
821 unclear and may require more advanced neuroscientific methods.

## 822 7.6 Conclusion and Future Directions

823 The literature described throughout this chapter has focused on lexical and audiovi-  
824 sual information as contextual influences on speech perception, as well as their  
825 dimensions and limitations. Section 7.2 highlighted the seminal findings regarding  
826 lexical retuning, starting from Norris et al. (2003) and the studies since then that  
827 have illuminated the strengths and drawbacks. Section 7.3 discussed audiovisual  
828 recalibration, first described by Bertelson et al. (2003) and expanded upon by others.

829 These two contextual sources can differ in terms of their impact on perception,  
830 as lexical information can potentially lead to more stable and longer-lasting shifts in  
831 perception, while audiovisual information results in adjustments in shorter dura-  
832 tions that are not easily generalizable and are often either (or both) context- and  
833 token-dependent. The phoneme categories themselves can also impose restrictions,  
834 as plosives (also known as stop consonants) may allow for generalization to other  
835 speakers more so than other types of phonemes, such as fricatives or liquids.  
836 Evidently, contextual cues alone do not drive these phoneme boundary shifts, and  
837 acoustic information still modulates learning effects to a great extent. Theories of  
838 speech perception have also been helpful for understanding the basis of phoneme  
839 boundary adjustments, but disagreements exist with regard to the stages of process-  
840 ing that are thought to be involved.

841 Although questions remain in the field as to the precise details of retuning,  
842 researchers continue to pursue the answers with behavioral and neuroimaging stud-  
843 ies. Related works may also shed light upon how exactly these perceptual shifts may  
844 occur. Recent studies have investigated another related form of text-based



recalibration. Reading text of syllables while listening to ambiguous phonemes can also contribute to changes in phoneme categorization (Keetels et al. 2016a, b), and this has also been tested using fMRI (Bonte et al. 2017). Just as in audiovisual and lexical experiments, participants viewed either /aba/ or /ada/ written in text, while hearing an ambiguous blend of the two, and participants were able to effectively recalibrate depending on the text they viewed (Keetels et al. 2016a, b). In addition, fMRI results showed that text-based recalibration was linked to activity in posterior superior temporal cortex, and percepts of /aba/ and /ada/ during test could also be decoded with MVPA, primarily based on patterns of activity in left posterior STG and PT and right STS (Bonte et al. 2017). Functional connectivity was observed between IPL and left STG during exposure and may be indicative of higher-order influences leading to eventual retuning. While lexical and audiovisual recalibration studies have been useful for understanding how listeners adapt to ambiguity in speech, this new paradigm illuminates how mappings are acquired between auditory and written representations, and may also have the potential to detect disruptions of reading networks during development, particularly in individuals with dyslexia.

Together, these approaches using lexical and audiovisual information, and more recently with text, have proven useful in understanding the plasticity of speech sounds. These non-acoustic sources of information can not only sway how speech tokens are perceived but, moreover, can restructure the units of speech. Evidently, these units are malleable and are continuously updated with experience; they are susceptible to change even within short windows of time and with relatively little input required to do so. This adaptive tool is beneficial for adjusting to speakers, noise, or other obstacles that could impede successful speech comprehension, although the acoustic features of the input may restrict the extent to which recalibration can be generalized. Still, stimulus specificity may be advantageous, as a complete overhaul of speech sounds in response to deviations from the norm would be impractical. Speech perception theories and neuroimaging studies have highlighted the possible processing streams involved, and both lexical and speech-reading influences appear to share significant similarities in terms of the brain areas being recruited. The relative contributions of top-down and bottom-up information in processing the acoustic input are still hotly debated, but the continued application of advanced neuroimaging techniques, as well as statistical modeling, may aid in building a more cohesive picture of perceptual retuning.

Compliance with Ethics Requirements Shruti Ullas declares no conflict of interest. 880

Milene Bonte declares no conflict of interest. 881

Elia Formisano declares no conflict of interest. 882

Jean Vroomen declares no conflict of interest. 883

- 885 Adank P, Devlin JT (2010) On-line plasticity in spoken sentence comprehension: adapt-  
 886 ing to time-compressed speech. *NeuroImage* 49(1):1124–1132. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.neuroimage.2009.07.032)  
 887 [neuroimage.2009.07.032](https://doi.org/10.1016/j.neuroimage.2009.07.032)
- 888 Baart M, Samuel AG (2015) Turning a blind eye to the lexicon: ERPs show no cross-talk between  
 889 lip-read and lexical context during speech sound processing. *J Mem Lang* 85:42–59. [https://](https://doi.org/10.1016/j.jml.2015.06.008)  
 890 [doi.org/10.1016/j.jml.2015.06.008](https://doi.org/10.1016/j.jml.2015.06.008)
- 891 Baart M, Vroomen J (2010) Phonetic recalibration does not depend on working memory. *Exp*  
 892 *Brain Res* 203:575–582. <https://doi.org/10.1007/s00221-010-2264-9>
- 893 Baart M, de Boer-Schellekens L, Vroomen J (2012) Lipread-induced phonetic recalibration in  
 894 dyslexia. *Acta Psychol* 140(1):91–95. <https://doi.org/10.1016/j.actpsy.2012.03.003>
- 895 Bertelson P, Vroomen J, De Gelder B (2003) Visual recalibration of auditory speech identification:  
 896 a McGurk aftereffect. *Psychol Sci* 14(6):592–597. [https://doi.org/10.1046/j.0956-7976.2003.](https://doi.org/10.1046/j.0956-7976.2003.psci_1470.x)  
 897 [psci\\_1470.x](https://doi.org/10.1046/j.0956-7976.2003.psci_1470.x)
- 898 Bonte M, Hausfeld L, Scharke W, Valente G, Formisano E (2014) Task-dependent decod-  
 899 ing of speaker and vowel identity from auditory cortical response patterns. *J Neurosci*  
 900 34(13):4548–4557. <https://doi.org/10.1523/JNEUROSCI.4339-13.2014>
- 901 Bonte M, Correia JM, Keetels M, Vroomen J, Formisano E (2017) Reading-induced shifts of  
 902 perceptual speech representations in auditory cortex. *Sci Rep* 7:1–11. [https://doi.org/10.1038/](https://doi.org/10.1038/s41598-017-05356-3)  
 903 [s41598-017-05356-3](https://doi.org/10.1038/s41598-017-05356-3)
- 904 Bradlow AR, Bent T (2008) Perceptual adaptation to non-native speech. *Cognition*  
 905 106(2):707–729. <https://doi.org/10.1016/j.cognition.2007.04.005>
- 906 Brancazio L (2004) Lexical influences in audiovisual speech perception. *J Exp Psychol Hum*  
 907 *Percept Perform* 30(3):445–463. <https://doi.org/10.1037/0096-1523.30.3.445>
- 908 Burgering M, van Laarhoven T, Baart M, Vroomen J (2020) Fluidity in the perception of auditory  
 909 speech: cross-modal recalibration of voice gender and vowel identity by a talking face. *Q J Exp*  
 910 *Psychol (Hove)* 73(6):957–967. <https://doi.org/10.1177/1747021819900884>
- 911 Clarke CM, Garrett MF (2004) Rapid adaptation to foreign-accented English. *J Acoust Soc Am*  
 912 116(6):3647–3658. <https://doi.org/10.1121/1.1815131>
- 913 Cutler A, McQueen JM, Butterfield S, Norris D (2008) Prelexically-driven perceptual retuning  
 914 of phoneme boundaries. In: Fletcher J, Loakes D, Goecke R, Burnham D, Wagner M (eds)  
 915 *Proceedings of Interspeech, Brisbane, 2008*
- 916 Cutler A, Eisner F, McQueen JM, Norris D (2010) How abstract phonemic categories are nec-  
 917 cessary for coping with speaker-related variation. In: Fougeron C, Kühnert B, D’Imperio M,  
 918 Vallée N (eds) *Laboratory phonology, vol 10*. de Gruyter, Berlin, pp 91–111
- 919 Davis MH, Johnsruide IS (2007) Hearing speech sounds: top-down influences on the interface  
 920 between audition and speech perception. *Hear Res* 229(1–2):132–147. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.heares.2007.01.014)  
 921 [heares.2007.01.014](https://doi.org/10.1016/j.heares.2007.01.014)
- 922 Davis MH, Johnsruide IS, Hervais-Adelman AG, Taylor K, McGettigan C (2005) Lexical  
 923 information drives perceptual learning of distorted speech: evidence from the compre-  
 924 hension of noise-vocoded sentences. *J Exp Psychol Gen* 134(2):222–241. [https://doi.](https://doi.org/10.1037/0096-3445.134.2.222)  
 925 [org/10.1037/0096-3445.134.2.222](https://doi.org/10.1037/0096-3445.134.2.222)
- 926 Dehaene-Lambertz G, Pallier C, Serniclaes W, Sprenger-Charolles L, Jobert A, Dehaene S (2005)  
 927 Neural correlates of switching from auditory to speech perception. *NeuroImage* 24(1):21–33.  
 928 <https://doi.org/10.1016/j.neuroimage.2004.09.039>
- 929 Drozdova P, van Hout R, Scharenborg O (2015) Lexically-guided perceptual learning in  
 930 non-native listening. *Biling (Camb Engl)* 19(5):914–920. doi: [https://doi.org/10.1017/](https://doi.org/10.1017/S136672891600002X)  
 931 [S136672891600002X](https://doi.org/10.1017/S136672891600002X)
- 932 Eimas PD, Corbit JD (1973) Selective adaptation of linguistic feature detectors. *Cogn Psychol*  
 933 4:99–109. [https://doi.org/10.1016/0010-0285\(73\)90006-6](https://doi.org/10.1016/0010-0285(73)90006-6)
- 934 Eisner F, McQueen JM (2005) The specificity of perceptual learning in speech processing. *Atten*  
 935 *Percept Psychophys* 67:224–238. <https://doi.org/10.3758/BF03206487>

- Eisner F, McQueen JM (2006) Perceptual learning in speech: stability over time. *J Acoust Soc Am* 119:1950–1953. <https://doi.org/10.1121/1.2178721> 936  
937
- Formisano E, De Martino F, Bonte M, Goebel R (2008) “Who” is saying “what”? Brain based decoding of human voice and speech. *Science* 322(5903):970–973. <https://doi.org/10.1126/science.1164318> 938  
939  
940
- Franken MK, Eisner F, Schoffelen JM, Acheson DJ, Hagoort P, McQueen JM (2017) Audiovisual recalibration of vowel categories. In: *Proceedings of Interspeech, Stockholm*, pp 655–658. <https://doi.org/10.21437/Interspeech.2017-122> 941  
942  
943
- Ganong WF (1980) Phonetic categorization in auditory word perception. *J Exp Psychol Hum Percept Perform* 6(1):110–125. <https://doi.org/10.1037/0096-1523.6.1.110> 944  
945
- Gaskell MG, Marslen-Wilson WD (1997) Integrating form and meaning: a distributed model of speech perception. *Lang Cogn Process* 12(5–6):613–656. <https://doi.org/10.1080/016909697386646> 946  
947  
948
- Guediche S, Blumstein SE, Fiez JA, Holt LL (2014) Speech perception under adverse conditions: insights from behavioral, computational, and neuroscience research. *Front Syst Neurosci* 7:1–16. <https://doi.org/10.3389/fnsys.2013.00126> 949  
950  
951
- Guediche S, Holt LL, Laurent P, Lim S, Fiez JA (2015) Evidence for cerebellar contributions to adaptive plasticity in speech perception. *Cereb Cortex* 25:1867–1877. <https://doi.org/10.1093/cercor/bht428> 952  
953  
954
- Hickok G, Poeppel D (2007) The cortical organization of speech processing. *Nat Rev Neurosci* 8:393–402. <https://doi.org/10.1038/nrn2113> 955  
956
- Holdgraf CR, de Heer W, Pasley B, Rieger J, Crone N, Lin JJ, Knight RT, Theunissen FE (2016) Rapid tuning shifts in human auditory cortex enhance speech intelligibility. *Nat Commun* 7:13654. <https://doi.org/10.1038/ncomms13654> 957  
958  
959
- Holt LL, Lotto AJ (2008) Speech perception within an auditory cognitive science framework. *Curr Dir Psychol Sci* 17(1):42–46. <https://doi.org/10.1111/j.1467-8721.2008.00545.x> 960  
961
- Jäncke L, Wüstenberg T, Scheich H, Heinze HJ (2002) Phonetic perception and the auditory cortex. *NeuroImage* 15(4):733–746. <https://doi.org/10.1006/nimg.2001.1027> 962  
963
- Keetels MN, Pecoraro M, Vroomen J (2015) Recalibration of auditory phonemes by lipread speech is ear-specific. *Cognition* 141:121–126. <https://doi.org/10.1016/j.cognition.2015.04.019> 964  
965
- Keetels MN, Schakel L, Bonte M, Vroomen J (2016a) Phonetic recalibration of speech by text. *Atten Percept Psychophys* 78:938–945. <https://doi.org/10.3758/s13414-015-1034-y> 966  
967
- Keetels MN, Stekelenburg JJ, Vroomen J (2016b) A spatial gradient in phonetic recalibration by lipread speech. *J Phon* 56:124–130. <https://doi.org/10.1016/j.wocn.2016.02.005> 968  
969
- Kilian-Hütten N, Valente G, Vroomen J, Formisano E (2011a) Auditory cortex encodes the perceptual interpretation of ambiguous sound. *J Neurosci* 31(5):1715–1720. <https://doi.org/10.1523/JNEUROSCI.4572-10.2011> 970  
971  
972
- Kilian-Hütten N, Vroomen J, Formisano E (2011b) Brain activation during audiovisual exposure anticipates future perception of ambiguous speech. *NeuroImage* 57(4):1601–1607. <https://doi.org/10.1016/j.neuroimage.2011.05.043> 973  
974  
975
- Kleinschmidt DF, Jaeger TF (2011) Robust speech perception: recognize the familiar, generalize to the similar, and adapt to the novel. *Psychol Rev* 122(2):148–203. <https://doi.org/10.1037/a0038695> 976  
977  
978
- Kraljic T, Samuel AG (2005) Perceptual learning for speech: is there a return to normal? *Cogn Psychol* 51:141–178. <https://doi.org/10.1016/j.cogpsych.2005.05.001> 979  
980
- Kraljic T, Samuel AG (2006) Generalization in perceptual learning for speech. *Psychon Bull Rev* 13:262–268. <https://doi.org/10.3758/BF03193841> 981  
982
- Kraljic T, Samuel AG (2007) Perceptual adjustments to multiple speakers. *J Mem Lang* 56:1–15. <https://doi.org/10.1016/j.jml.2006.07.010> 983  
984
- Kraljic T, Samuel AG (2009) Perceptual learning for speech. *Atten Percept Psychophys* 71(3):1207–1218. <https://doi.org/10.3758/APP.71.6.1207> 985  
986
- Kraljic T, Brennan SE, Samuel AG (2008a) Accommodating variation: dialects, idiolects, and speech processing. *Cognition* 107:51–81. <https://doi.org/10.1016/j.cognition.2007.07.013> 987  
988

- 989 Kraljic T, Samuel AG, Brennan SE (2008b) First impressions and last resorts: how listeners adjust to  
 990 speaker variability. *Psychol Sci* 19:332–338. <https://doi.org/10.1111/j.1467-9280.2008.02090.x>
- 991 Lecumberri MLG, Cooke M, Cutler A (2010) Non-native speech perception in adverse conditions:  
 992 a review. *Speech Commun* 52(11–12):864–886. <https://doi.org/10.1016/j.specom.2010.08.014>.
- 993 Lieberthal E, Desai R, Ellingson MM, Ramachandran B, Desai A, Binder JR (2010)  
 994 Specialization along the left superior temporal sulcus for auditory categorization. *Cereb Cortex*  
 995 20(12):2958–2970. <https://doi.org/10.1093/cercor/bhq045>
- 996 Lüttke C, Ekman M, van Gerven M, de Lange FP (2016) McGurk illusion recalibrates subsequent  
 997 auditory perception. *Sci Rep* 6:32891. <https://doi.org/10.1038/srep32891>
- 998 Maye J, Aslin RN, Tanenhaus MK (2008) The Weckud Wetch of the Wast: Lexical adaptation to a  
 999 novel accent. *Cogn Sci* 32(3):543–562. <https://doi.org/10.1080/03640210802035357>
- 1000 McClelland JL, Elman JL (1986) The TRACE model of speech perception. *Cogn Psychol* 18:1–86.  
 1001 [https://doi.org/10.1016/0010-0285\(86\)90015-0](https://doi.org/10.1016/0010-0285(86)90015-0)
- 1002 McClelland JL, Mirman D, Holt LL (2006) Are there interactive processes in speech perception?  
 1003 *Trends Cogn Sci* 10(8):363–369. <https://doi.org/10.1016/j.tics.2006.06.007>
- 1004 McGurk H, MacDonald J (1976) Hearing lips and seeing voices. *Nature* 264:746–748. <https://doi.org/10.1038/264746a0>
- 1005
- 1006 McQueen JM, Cutler A, Norris D (2006a) Phonological abstraction in the mental lexicon. *Cogn*  
 1007 *Sci* 30:1113–1126. [https://doi.org/10.1207/s15516709cog0000\\_79](https://doi.org/10.1207/s15516709cog0000_79)
- 1008 McQueen JM, Norris D, Cutler A (2006b) The dynamic nature of speech perception. *Lang Speech*  
 1009 49(1):101–112. <https://doi.org/10.1177/00238309060490010601>
- 1010 Mesgarani N, David SV, Fritz JB, Shamma SA (2008) Phoneme representation and classification in  
 1011 primary auditory cortex. *J Acoust Soc Am* 123(2):899–909. <https://doi.org/10.1121/1.2816572>
- 1012 Mesgarani N, Cheung C, Johnson K, Chang EF (2014) Phonetic feature encoding in human superior  
 1013 temporal gyrus. *Science* 343(6174):1006–1010. <https://doi.org/10.1126/science.1245994>
- 1014 Mitterer H, Scharenborg O, McQueen JM (2013) Phonological abstraction without phonemes in  
 1015 speech perception. *Cognition* 129:356–261. <https://doi.org/10.1016/j.cognition.2013.07.011>
- 1016 Myers EB, Blumstein SE (2008) The neural basis of the lexical effect: an fMRI investigation.  
 1017 *Cereb Cortex* 18:278–288. <https://doi.org/10.1093/cercor/bhm053>
- 1018 Myers EB, Mesite LM (2014) Neural systems underlying perceptual adjustment to non-standard  
 1019 speech tokens. *J Mem Lang* 76:80–93. <https://doi.org/10.1093/cercor/bhm053>
- 1020 Norris D, McQueen JM, Cutler A (2000) Merging information in speech recognition: feedback is  
 1021 never necessary. *Behav Brain Sci* 23:299–325. <https://doi.org/10.1017/S0140525X00003241>
- 1022 Norris D, McQueen JM, Cutler A (2003) Perceptual learning in speech. *Cogn Psychol* 47:204–238.  
 1023 [https://doi.org/10.1016/S0010-0285\(03\)00006-9](https://doi.org/10.1016/S0010-0285(03)00006-9)
- 1024 Norris D, Cutler A, McQueen JM, Butterfield S (2006) Phonological and conceptual activation  
 1025 in speech comprehension. *Cogn Psychol* 53(2):146–193. <https://doi.org/10.1016/j.cogpsych.2006.03.001>
- 1026
- 1027 Norris D, McQueen JM, Cutler A (2016) Prediction, Bayesian inference and feedback in speech  
 1028 recognition. *Lang Cogn Neurosci* 31(1):4–18. <https://doi.org/10.1080/23273798.2015.1081703>
- 1029 3
- 1030 Obleser J, Eisner F (2009) Pre-lexical abstraction of speech in the auditory cortex. *Trends Cogn*  
 1031 *Sci* 13(1):14–19. <https://doi.org/10.1016/j.tics.2008.09.005>
- 1032 Raizada RD, Poldrack RA (2007) Selective amplification of stimulus differences during categori-  
 1033 cal processing of speech. *Neuron* 56(4):726–740. <https://doi.org/10.1016/j.neuron.2007.11.001>
- 1034 Reinisch E, Holt LL (2014) Lexically-guided phonetic retuning of foreign-accented speech and its  
 1035 generalization. *J Exp Psychol Hum Percept Perform* 40(2):539–555. <https://doi.org/10.1037/a0034409>
- 1036
- 1037 Reinisch E, Weber A, Mitterer H (2013) Listeners retune phoneme categories across languages. *J*  
 1038 *Exp Psychol Hum Percept Perform* 39:75–86. <https://doi.org/10.1037/a0027979>
- 1039 Reinisch E, Wozny D, Mitterer H, Holt LL (2014) Phonetic category recalibration: what are the  
 1040 categories? *J Phon* 45:91–105. <https://doi.org/10.1016/j.wocn.2014.04.002>

- Remez RE, Rubin PE, Pisoni DB, Carell TD (1981) Speech perception without traditional speech cues. *Science* 212:947–950 1041  
1042
- Roberts M, Summerfield Q (1981) Audiovisual presentation demonstrates that selective adaptation in speech perception is purely auditory. *Atten Percept Psychophys* 30(4):309–314. <https://doi.org/10.3758/BF03206144> 1043  
1044  
1045
- Rogers JC, Davis MH (2017) Inferior frontal cortex contributions to the recognition of spoken words and their constituent speech sounds. *J Cogn Neurosci* 29(5):919–936. [https://doi.org/10.1162/jocn\\_a\\_01096](https://doi.org/10.1162/jocn_a_01096) 1046  
1047  
1048
- Romanovska L, Janssen R, Bonte M (2019) Reading-induced shifts in speech perception in dyslexic and typically reading children. *Front Psychol* 10:221. <https://doi.org/10.3389/fpsyg.2019.00221> 1049  
1050  
1051
- Samuel AG, Frost R (2015) Lexical support for phonetic perception during non-native spoken word recognition. *Psychon Bull Rev* 22(6):1746–1752. <https://doi.org/10.3758/s13423-015-0847-y> 1052  
1053
- Sjerps MJ, McQueen JM (2010) The bounds on flexibility in speech perception. *J Exp Psychol Hum Percept Perform* 36:195–211. <https://doi.org/10.1037/a0016803> 1054  
1055
- Snowling MJ (1980) The development of grapheme-phoneme correspondence in normal and dyslexic readers. *J Exp Child Psychol* 29:294–305. [https://doi.org/10.1016/0022-0965\(80\)90021-1](https://doi.org/10.1016/0022-0965(80)90021-1) 1056  
1057
- Sohoglu E, Davis MH (2016) Perceptual learning of degraded speech by minimizing prediction error. *Proc Natl Acad Sci USA* 113(12):1747–1756. <https://doi.org/10.1073/pnas.1523266113> 1058  
1059
- Sumbly WH, Pollack I (1954) Visual contribution to speech intelligibility in noise. *J Acoust Soc Am* 26:212–215. <https://doi.org/10.1121/1.1907309> 1060  
1061
- Ullas S, Hausfeld L, Cutler A, Eisner F, Formisano E (2020) Neural correlates of phonetic adaptation as induced by lexical and audiovisual context. *J Cogn Neurosci*:1–14. [https://doi.org/10.1162/jocn\\_a\\_01608](https://doi.org/10.1162/jocn_a_01608) 1062  
1063  
1064
- Ullas S, Formisano E, Eisner F, Cutler A (2020a) Interleaved lexical and audiovisual information can retune phoneme boundaries. *Atten Percept Psychophys* 82:2018–2026. <https://doi.org/10.3758/s13414-019-01961-8> 1065  
1066  
1067
- Ullas S, Formisano E, Eisner F, Cutler A (2020b) Audiovisual and lexical cues do not additively enhance perceptual adaptation. *Psychon Bull Rev* 27:707–715. <https://doi.org/10.3758/s13423-020-01728-5> 1068  
1069  
1070
- Van der Zande P, Jesse A, Cutler A (2014) Hearing words helps seeing words: a cross-modal word repetition effect. *Speech Commun* 59:31–43. <https://doi.org/10.1016/j.specom.2014.01.001> 1071  
1072
- Van Laarhoven T, Keetels M, Schakel L, Vroomen J (2018) Audio-visual speech in noise perception in dyslexia. *Dev Sci* 21(1):e12504. <https://doi.org/10.1111/desc.12504> 1073  
1074
- Van Linden S, Vroomen J (2007) Recalibration of phonetic categories by lipread speech versus lexical information. *J Exp Psychol Hum Percept Perform* 33(6):1483–1494. <https://doi.org/10.1037/0096-1523.33.6.1483> 1075  
1076  
1077
- Van Linden S, Vroomen J (2008) Audiovisual speech recalibration in children. *J Child Lang* 35(4):809–822. <https://doi.org/10.1017/S0305000908008817> 1078  
1079
- Vroomen J, Baart M (2009a) Phonetic recalibration only occurs in speech mode. *Cognition* 110(2):254–259. <https://doi.org/10.1016/j.cognition.2008.10.015> 1080  
1081
- Vroomen J, Baart M (2009b) Recalibration of phonetic categories by lipread speech: measuring aftereffects after a twenty-four hours delay. *Lang Speech* 52:341–350. <https://doi.org/10.1177/0023830909103178> 1082  
1083  
1084
- Vroomen J, van Linden S, Keetels M, de Gelder B, Bertelson P (2004) Selective adaptation and recalibration of auditory speech by lipread information: dissipation. *Speech Commun* 44:55–61. <https://doi.org/10.1016/j.specom.2004.03.009> 1085  
1086  
1087
- Vroomen J, van Linden S, Baart M (2007a) Lipread aftereffects in auditory speech perception: measuring aftereffects after a twenty-four hours delay. In: Vroomen J, Swerts M, Krahmer E (eds) *Auditory-visual speech processing*, Hilvarenbeek, p P05 1088  
1089  
1090
- Vroomen J, van Linden S, de Gelder B, Bertelson P (2007b) Visual recalibration and selective adaptation in auditory-visual speech perception: contrasting build-up courses. *Neuropsychologia* 45(3):572–577. <https://doi.org/10.1016/j.neuropsychologia.2006.01.031> 1091  
1092  
1093

- 1094 Winn M (2018) Speech: it's not as acoustic as you think. *Acoust Today* 14(2):43–49
- 1095 Xie X, Myers EB (2017) Learning a talker or learning an accent: acoustic similarity constrains  
1096 generalization of foreign accent adaptation to new talkers. *J Mem Lang* 97:30–46. [https://doi.](https://doi.org/10.1016/j.jml.2017.07.005)  
1097 [org/10.1016/j.jml.2017.07.005](https://doi.org/10.1016/j.jml.2017.07.005)
- 1098 Zhang X, Samuel AG (2015) Perceptual learning of speech under optimal and adverse condition. *J*  
1099 *Exp Psychol Hum Percept Perform* 40(1):200–217. <https://doi.org/10.1037/a0033182>

Uncorrected Proof

# Author Queries

Chapter No.: 7      0005197830

Queries	Details Required	Author's Response
AU1	“Krajlic and Samuel (2008a), Samuel (2016), Eimas and Corbitt (1973), Vroomen et al. (2006), Scott and Johnsrude (2003), Rauschecker and Scott (2009)” is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information. Otherwise, please delete it from the text/body.	
AU2	The citation “Keetels et al. (2016)” has been changed to “Keetels et al. (2016a, b)” to match the author name/date in the reference list. Please check if the change is fine in this occurrence and modify the subsequent occurrences, if necessary.	
AU3	The citation “Gaskell and Marslen-Wilson 1987” has been changed to “Gaskell and Marslen-Wilson 1997” to match the author name/date in the reference list. Please check if the change is fine in this occurrence and modify the subsequent occurrences, if necessary.	
AU4	The citation “Kleinschmidt and Jaeger (2015)” has been changed to “Kleinschmidt and Jaeger (2011)” to match the author name/date in the reference list. Please check if the change is fine in this occurrence and modify the subsequent occurrences, if necessary.	
AU5	“Sect. 1.1” is not available in this chapter. Please check and provide alternate citation	
AU6	References “Baart & Samuel (2015), Norris et al. (2006), Reinisch & Holt (2014), Roberts & Summerfield (1981), Vroomen & Baart (2009a)” were not cited anywhere in the text. Please provide in text citation or delete the reference from the reference list.	