

# Journal of Speech, Language, and Hearing Research

## Eye Gaze and Perceptual Adaptation to Audiovisual Degraded Speech

--Manuscript Draft--

<b>Manuscript Number:</b>	JSLHR-21-00106R1	
<b>Full Title:</b>	Eye Gaze and Perceptual Adaptation to Audiovisual Degraded Speech	
<b>Article Type:</b>	Research Article	
<b>Section/Category:</b>	Speech	
<b>Corresponding Author:</b>	Briony Banks Lancaster University Lancaster, Lancashire UNITED KINGDOM	
<b>Other Authors:</b>	Emma Gowen	
	Kevin J Munro	
	Patti Adank	
<b>Funding Information:</b>	Biotechnology and Biological Sciences Research Council	Dr Briony Banks
<b>Keywords:</b>	Speech Perception; audiovisual speech; eye tracking	
<b>Manuscript Classifications:</b>	Cognition; Speech perception; Speech recognition	
<b>Abstract:</b>	<p><b>Purpose :</b> Visual cues from a speaker’s face may benefit perceptual adaptation to degraded speech, but current evidence is limited. We aimed to replicate results from previous studies to establish the extent to which visual speech cues can lead to greater adaptation over time, extending existing results to a real-time adaptation paradigm (i.e., without a separate training period). A second aim was to investigate whether eye gaze patterns towards the speaker’s mouth were related to better perception, hypothesising that listeners who looked more at the speaker’s mouth would show greater adaptation.</p> <p><b>Method:</b> A group of listeners ( N =30) were presented with 90 noise-vocoded sentences in audiovisual format while a control group ( N =29) were presented with the audio signal only. Recognition accuracy was measured throughout and eye tracking was used to measure fixations towards the speaker’s eyes and mouth in the audiovisual group.</p> <p><b>Results:</b> Previous studies were partially replicated: the audiovisual group had better recognition throughout and adapted slightly more rapidly, but both groups showed an equal amount of improvement overall. Longer fixations on the speaker’s mouth in the audiovisual group were related to better overall accuracy. An exploratory analysis further demonstrated that the duration of fixations to the speaker’s mouth decreased over time.</p> <p><b>Conclusions:</b> The results suggest that visual cues may not benefit adaptation to degraded speech as much as previously thought. Longer fixations on a speaker’s mouth may play a role in successfully decoding visual speech cues, however this will need to be confirmed in future research to fully understand how patterns of eye gaze are related to audiovisual speech recognition. All materials, data, and code are available at <a href="https://osf.io/2wqkf/">https://osf.io/2wqkf/</a> .</p>	
<b>Response to Reviewers:</b>	<p>Thank you for reviewing our manuscript “Eye Gaze and Perceptual Adaptation to Audiovisual Degraded Speech”, manuscript no. JSLHR-21-00106. We appreciate the helpful and informative comments from the editor and the reviewers. We have addressed all of these in our revision of the manuscript; particularly, we have toned down some of our conclusions, provided extra details as requested, and revised the introduction and discussion accordingly, including additional literature where appropriate. We have also ensured that all data and materials are now publicly available on the OSF. We attach a response to reviewers addressing each point in detail. All changes are highlighted in yellow in the manuscript. We have also addressed the minor points highlighted by the editorial team, again highlighted in yellow. Thank you for considering the revised version of our paper for publication in the Journal of Speech, Language, and Hearing Research, and we look forward to hearing from you.</p>	

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

## Eye Gaze and Perceptual Adaptation to Audiovisual Degraded Speech

Briony Banks<sup>1</sup>, Emma Gowen<sup>1</sup>, Kevin J Munro<sup>2,3</sup> and Patti Adank<sup>4</sup>

<sup>1</sup> Faculty of Biology, Medicine and Health, The University of Manchester

<sup>2</sup> Manchester Centre for Audiology and Deafness, Faculty of Biology, Medicine and Health, The University of Manchester

<sup>3</sup> Manchester University Hospitals NHS Foundation Trust, Manchester Academic Health Science Centre

<sup>4</sup> Speech, Hearing and Phonetic Sciences, University College London

### Author Note

Briony Banks is now at Department of Psychology, Lancaster University.

We have no known conflict of interest to disclose.

BB was funded by a BBSRC research studentship and by The University of Manchester. KJM is supported by the NIHR Manchester Biomedical Research Centre.

Correspondence concerning this article should be addressed to Briony Banks, Department of Psychology, Fylde College, Lancaster University, Lancaster, LA1 4YF, UK, email:

[b.banks@lancaster.ac.uk](mailto:b.banks@lancaster.ac.uk)

23

**Abstract**

24 **Purpose:** Visual cues from a speaker's face may benefit perceptual adaptation to degraded speech,  
25 but current evidence is limited. We aimed to replicate results from previous studies to establish the  
26 extent to which visual speech cues can lead to greater adaptation over time, extending existing  
27 results to a real-time adaptation paradigm (i.e., without a separate training period). A second aim  
28 was to investigate whether eye gaze patterns towards the speaker's mouth were related to better  
29 perception, hypothesising that listeners who looked more at the speaker's mouth would show  
30 greater adaptation.

31 **Method:** A group of listeners ( $N=30$ ) were presented with 90 noise-vocoded sentences in audiovisual  
32 format while a control group ( $N=29$ ) were presented with the audio signal only. Recognition  
33 accuracy was measured throughout and eye tracking was used to measure fixations towards the  
34 speaker's eyes and mouth in the audiovisual group.

35 **Results:** Previous studies were partially replicated: the audiovisual group had better recognition  
36 throughout and adapted slightly more rapidly, but both groups showed an equal amount of  
37 improvement overall. Longer fixations on the speaker's mouth in the audiovisual group were related  
38 to better overall accuracy. An exploratory analysis further demonstrated that the duration of  
39 fixations to the speaker's mouth decreased over time.

40 **Conclusions:** The results suggest that visual cues may not benefit adaptation to degraded speech as  
41 much as previously thought. Longer fixations on a speaker's mouth may play a role in successfully  
42 decoding visual speech cues, however this will need to be confirmed in future research to fully  
43 understand how patterns of eye gaze are related to audiovisual speech recognition. All materials,  
44 data, and code are available at <https://osf.io/2wqkf/>.

45 **Key words:** Speech perception, audiovisual speech, perceptual adaptation, eye tracking

46

47 **Eye Gaze and Perceptual Adaptation to Audiovisual Degraded Speech**

48 Human communication often takes place in suboptimal listening conditions such as in noisy  
49 environments, listening to a distorted phone or video signal, or encountering unfamiliar speech such  
50 as a foreign accent. Most listeners are adept at dealing with such difficult conditions by rapidly  
51 adapting to them – that is, undergoing a period where they learn and ‘tune in’ to the acoustic and  
52 perceptual differences in the particular listening condition. This perceptual adaptation to degraded  
53 or unfamiliar speech has been consistently and empirically demonstrated for a variety of adverse  
54 conditions, such as noise-vocoded (M. H. Davis et al., 2005; Hervais-Adelman et al., 2008), accented  
55 (Adank & Janse, 2010; Banks et al., 2015a, 2015b), and time-compressed speech (Pelle & Wingfield,  
56 2005; Sebastian-Galles & Mehler, 2000). Artificially degrading the speech through noise-vocoding  
57 (Shannon et al., 1995) is particularly useful in such experiments due to the level of control that it  
58 offers the experimenter, particularly with regards to intelligibility (e.g., Dorman et al., 1997; Faulkner  
59 et al., 2000). Noise-vocoding distorts the spectral structure of speech while preserving the temporal  
60 structure, creating a speech signal that contains enough detail to be intelligible but with significantly  
61 less spectral, specifically harmonic, detail than the original (M. H. Davis et al., 2005). The relative  
62 intelligibility of the signal is associated with the number of channels initially used to divide the  
63 acoustic signal, with more channels resulting in higher levels of intelligibility (Loizou et al., 1999).  
64 Listeners can adapt to noise-vocoded sentences after relatively short exposure; for example, Davis et  
65 al. (2005) report a steady linear increase in recognition performance after listening to 30 sentences  
66 noise-vocoded into six channels, with participants improving from ~20% of words correctly reported  
67 to ~60%. Distortions such as noise-vocoding can reflect particularly challenging conditions that we  
68 might encounter in modern digital communication. However, the processes and individual strategies  
69 used during perceptual adaptation are still not fully understood, particularly the role of visual speech  
70 cues, as although we often communicate face-to-face with a speaker, the majority of research into  
71 perceptual adaptation of degraded speech has only examined auditory perception.

72 It is well established that access to visual cues from a speaker's face substantially improves  
73 speech recognition in difficult listening conditions; this *audiovisual benefit* has been demonstrated,  
74 for example, in the presence of background noise or with a distorted speech signal (Erber, 1975;  
75 MacLeod & Summerfield, 1987; Sommers et al., 2005; Sumbly & Pollack, 1954). Listeners benefit  
76 from viewing articulatory cues, particularly from a speaker's mouth, integrating them with auditory  
77 cues and thus enhancing the overall speech signal and improving recognition (Summerfield, 1987).  
78 Attending to visual speech cues may thus improve or speed up the adaptation process required to  
79 adapt to unfamiliar or degraded speech, leading to greater improvements in speech recognition.

80 A handful of studies have investigated the benefits of visual speech cues in perceptual  
81 adaptation to degraded (noise-vocoded) speech, but with varying types of linguistic stimuli. At the  
82 syllable level, Bernstein, Auer, Eberhardt & Jiang (2013) found that the presence of visual speech  
83 cues leads to greater perceptual adaptation of noise-vocoded syllables. Kawase et al., (2009)  
84 extended this finding to individual noise-vocoded words, comparing perceptual adaptation with and  
85 without audiovisual speech cues (i.e., with and without the speaker's face visible), finding that  
86 listeners adapted a greater amount when visual speech cues were available to listeners compared to  
87 when they were not. However, listening to individual syllables or words, without any additional  
88 linguistic context, is not representative of everyday communication. Pilling and Thomas (2011)  
89 therefore tested auditory recognition of degraded sentences. Participants listened to 3 blocks of 76  
90 noise-vocoded sentences, whereby the middle block was a training condition with either  
91 audiovisual, audio-only or non-degraded sentences. They observed a greater improvement in  
92 performance after training with visual cues compared to without (i.e., after exposure to audiovisual  
93 compared to audio-only sentences during training). Wayne & Johnsrude (2012) also assessed the  
94 contribution of training with visual speech information, comparing several training conditions during  
95 adaptation to noise-vocoded sentences. They found that training with audiovisual cues resulted in  
96 no more adaptation than training with non-degraded feedback – i.e., training where the listener  
97 heard the sentences both with and without noise-vocoding. However, the paradigm did not directly

98 compare adaptation to noise-vocoded speech with and without visual speech cues as in Pilling &  
99 Thomas (2011), and it is therefore impossible to ascertain the amount of improvement that visual  
100 cues contributed to adaptation over and above the auditory signal alone. Moreover, if one is  
101 listening to speech in adverse conditions (e.g., a degraded phone or video signal) it is not always  
102 possible to obtain the type of clear (i.e. non-degraded) feedback as used in the training conditions by  
103 Wayne & Johnsrude, and visual cues may thus provide a more readily accessible source of  
104 perceptual information that can help listeners adapt to difficult listening conditions.

105 Both Pilling & Thomas (2011), and Wayne & Johnsrude (2012), used a training paradigm  
106 whereby adaptation was measured by testing participants *after* being exposed to audiovisual  
107 speech; however, adaptation to unfamiliar or degraded speech most likely occurs in real time – that  
108 is, we adapt to the listening conditions we are exposed to at the time, integrating useful visual cues  
109 as we adapt. Furthermore, the sentences used in both Pilling & Thomas (2011) and Wayne &  
110 Johnsrude (2012) were relatively simple in terms of vocabulary and structure. Such sentences may  
111 be relatively easy to perceive and adapt to compared to more challenging and less predictable  
112 sentences; for example, the more challenging IEEE sentences (e.g., ‘Sickness kept him home the  
113 third week’, ‘The hog crawled under the high fence’; Rothauser et al., 1969) result in poorer  
114 recognition than the BKB sentences (e.g., ‘A cat sits on the bed’, ‘The ice cream was pink’; Bench et  
115 al., 1979) used by Pilling & Thomas (2011), when presented in fluctuating masking (Schoof & Rosen,  
116 2015). It is therefore possible that an equivalent audiovisual benefit to perceptual adaptation may  
117 not be present for different linguistic stimuli.

118 The benefit gained from visual speech cues has potential applications for listeners adapting  
119 to a variety of difficult listening conditions – whether these originate from the environment (for  
120 example background noise or a distorted phone line) or from listeners themselves in the form of a  
121 hearing impairment (Mattys et al., 2012). Nevertheless, current evidence of an audiovisual benefit to  
122 adaptation using naturalistic stimuli (i.e., sentences) comes essentially from a single study (Pilling &  
123 Thomas, 2011). The first aim of the present study was thus to replicate and extend the finding by

124 Pilling and Thomas (2011) that visual speech cues improve perceptual adaptation to degraded  
125 sentences, using a more naturalistic and real-time (i.e., continuous) adaptation paradigm whereby  
126 participants were continually exposed to noise-vocoded sentences with and without visual speech  
127 cues, and where recognition was measured throughout the task, rather than after a period of  
128 training. Additionally, we used the IEEE sentences (Rothausser et al., 1969), which are more complex  
129 than the BKB sentences, and thus potentially more challenging for listeners to integrate the auditory  
130 and visual signals, to more strongly test the effects of visual speech cues.

131         A second aim of the present study was to examine the role of eye gaze in comprehending  
132 and adapting to audiovisual degraded speech. Interest in listeners' eye gaze during speech  
133 perception has seen a recent increase (e.g., Barenholtz et al., 2016; Birulés et al., 2020; Lusk &  
134 Mitchel, 2016; Morin-Lessard et al., 2019; Wang Jianrong et al., 2020; Worster et al., 2018), with  
135 some studies suggesting a link between where and how listeners view a speaker's face and their  
136 resulting comprehension (Lusk & Mitchel, 2016; Worster et al., 2018). Adult listeners normally show  
137 a preference for looking at a speaker's eyes during communication (Morin-Lessard et al., 2019;  
138 Yarbus, 1967), which is likely for social reasons (Birmingham & Kingstone, 2009). Indeed, speech  
139 recognition studies employing eye-tracking have shown that in optimal listening conditions (i.e., in  
140 quiet and with a clear auditory signal), adults look more towards a speaker's eyes than the mouth  
141 (Buchan et al., 2007, 2008; Vatikiotis-Bateson et al., 1998). However, when listening conditions are  
142 challenging, e.g., when background noise is present, listeners look more often at a speaker's mouth  
143 (Buchan et al., 2007, 2008; Lansing & McConkie, 2003; Vatikiotis-Bateson et al., 1998). This pattern  
144 has also been found for artificial (Lusk & Mitchel, 2016) and non-native language (Barenholtz et al.,  
145 2016; Birulés et al., 2020). Indeed, the more challenging the condition (e.g., as background noise  
146 increases), the more frequently listeners look towards a speaker's mouth (Vatikiotis-Bateson et al.,  
147 1998) and the more attentional weighting is given to visual over auditory cues (Hazan et al., 2010).  
148 Although some useful speech cues can be gained from extra-oral areas such as the upper face and  
149 eye region (e.g., Preminger et al, 1998; Scheinberg, 1980), visible mouth movements are

150 considerably more important for successful audiovisual speech comprehension in challenging  
151 listening conditions (Thomas & Jordan, 2004). Thus, in such conditions, listeners likely shift their  
152 attention (and thus their eye gaze) more frequently towards the speaker's mouth to benefit from  
153 the most useful visual cues (i.e., articulatory mouth movements), potentially to improve lexical  
154 segmentation (Lusk & Mitchel, 2016; Mitchel & Weiss, 2014). These observations fit well with the  
155 cognitive relevance framework of visual attention (Henderson et al., 2009), which stipulates that the  
156 weight allocated to a particular visual feature is dependent on the cognitive needs of the perceiver.  
157 Accordingly, gaze patterns towards facial features during audiovisual speech perception have been  
158 shown to vary depending on the task (Buchan et al., 2007; Malcolm et al., 2008) and the type of  
159 stimuli presented (Lansing & McConkie, 2003; Vo et al., 2012).

160 Observations that listeners look more towards the speaker's mouth in adverse listening  
161 conditions would suggest a direct relationship between listeners' patterns of eye gaze and successful  
162 recognition of audiovisual degraded speech – i.e., listeners' performance. Indeed, in both deaf and  
163 hearing children, the amount of time spent looking at a speaker's mouth has been related to better  
164 speech-reading (i.e., lip-reading) accuracy (Worster et al., 2018), although the same relationship was  
165 not observed in normal-hearing adults (Lansing & McConkie, 2003; Wilson, Alsius, Pare, & Munhall,  
166 2016). Perception of the McGurk effect has also been related to listeners' patterns of eye gaze,  
167 whereby significantly more time is spent looking at a speaker's mouth in trials when it is perceived  
168 (Stacey et al., 2020), and stronger perceivers of the effect spend overall more time looking at the  
169 speaker's mouth than their eyes (Gurler et al., 2015). Nevertheless, the relevance of the McGurk  
170 illusion to audiovisual speech recognition is unclear (Alsius et al., 2018), and an equivalent  
171 relationship between patterns of eye gaze and audiovisual speech recognition has still not been  
172 found.

173 Two studies have reported correlational analyses between measurements of eye gaze and  
174 audiovisual speech recognition (Buchan et al., 2007; Everdell et al., 2007), but no significant  
175 correlations were observed. However, these analyses were not the main aim of the above studies,



176 and certain aspects of their methodology may explain the lack of observed correlations, namely  
177 ceiling effects in recognition accuracy which likely reduced variability in the measure. Furthermore,  
178 different measures of eye gaze have been used between studies; while some have focused on the  
179 length of time spent fixating on the eyes and mouth (Worster et al., 2018), others have measured  
180 the number of fixations (Lansing & McConkie, 2003) or trials (Buchan et al., 2007) spent looking at  
181 the speaker's mouth, or even left-right asymmetry of eye gaze on the eyes and mouth (Everdell et  
182 al., 2007), so it is unclear if one particular pattern of eye movements is particularly important during  
183 speech perception.

184           More recently, Lusk & Mitchell (2016) demonstrated that, after a period of familiarisation,  
185 better speech segmentation of an artificial language (i.e., strings of non-words) was related to  
186 greater shifts in attention between the eyes and mouth during familiarisation – however, these  
187 shifts took place in either direction (i.e., participants looked more or less at the mouth over time), so  
188 it is unclear if a particular eye gaze strategy was directly related to learning the new language.  
189 Lewkowicz & Hansen-Tift (2012) demonstrated that infants shift their eye gaze more towards a  
190 speaker's mouth when learning to speak, but look more at the eyes at a later stage of development  
191 when they have become more proficient, indicating that looking at a speaker's mouth is important  
192 during language acquisition. Conversely, Birulés, Bosch, Pons & Lewkowicz (2020) demonstrated that  
193 non-native adult listeners look more at a speaker's mouth than native speakers regardless of their  
194 language proficiency, suggesting that eye gaze towards the mouth is not necessarily linked to  
195 learning or performance. In summary, evidence in support of a relationship between eye gaze  
196 patterns and language learning are mixed, and nevertheless, the mechanisms of learning a language  
197 (as investigated in the above studies), may differ from the mechanisms of adapting to unfamiliar  
198 speech in one's native language.

199           The following questions therefore remain unanswered with regards to eye gaze and  
200 perception of audiovisual degraded speech: first, are measures of eye gaze on a speaker's mouth  
201 related to i) listeners' speech recognition accuracy, and ii) amount of adaptation to the unfamiliar

202 speech? Secondly, if such a relationship exists, is there a particular pattern of eye gaze on the  
203 speaker's mouth (for example, longer or more frequent fixations) that is related to better speech  
204 recognition and adaptation? Using eye tracking to investigate patterns of eye gaze towards a  
205 speaker's eyes and mouth during a relatively challenging speech recognition task, that avoids ceiling  
206 effects and where performance has room to improve over time, may reveal a direct relationship  
207 between eye gaze towards a speaker's mouth and audiovisual speech recognition.

208         The current study therefore had two aims: 1) To replicate and extend previous findings that  
209 the presence of visual speech cues improves perceptual adaptation to degraded speech, and 2) to  
210 examine the relationship between eye gaze on a speaker's mouth and speech recognition, as well as  
211 amount of adaptation (i.e., improvements in speech recognition over time). To address these aims,  
212 we measured recognition of degraded sentences in a real-time adaptation paradigm (i.e., where  
213 adaptation occurs during continuous exposure rather than after a training period), with and without  
214 visual speech cues. We recorded audiovisual sentences spoken from a single speaker and degraded  
215 these sentences using noise-vocoding; thus, we could create a relatively challenging speech  
216 recognition task that would avoid the ceiling and floor effects found in previous studies.

217         In a between-subjects design, we exposed a test group to audiovisual degraded speech  
218 stimuli, and a control group to audio-only degraded speech stimuli, using eye-tracking to measure  
219 participants' eye gaze. The control group was included to allow for direct comparison of speech  
220 recognition with and without visual speech cues. For consistency in our methods, we carried out eye  
221 tracking in both conditions, but presented the audio-only group with a static image of the speaker's  
222 face, therefore offering no dynamic visual cues that could be used to benefit speech recognition (see  
223 Methods for full details). To analyse eye gaze patterns during audiovisual speech recognition, we  
224 selected two commonly used eye-tracking variables in line with previous studies of audiovisual  
225 speech recognition: fixation duration and percentage fixations (Buchan et al., 2007; Everdell et al.,  
226 2007; Lansing & McConkie, 2003). Fixations (i.e., any period of time when eye gaze is relatively still;  
227 see Methods for full details) reflect the perceiver's foveal field of vision and thus the area of greatest

228 visual acuity. The frequency and duration of fixations can indicate where and to what extent a  
229 perceiver's visual attention is primarily directed at any given time (Christianson et al., 1991), and so  
230 are a good indicator of when listeners are attending to visual speech cues.

231 We predicted that perceptual adaptation would be greater when visual speech cues were  
232 visible – that is, recognition of the noise-vocoded speech would improve more in the audiovisual  
233 group compared to the audio-only group. Secondly, we predicted that recognition accuracy and  
234 adaptation in the audiovisual group would be related to the percentage and duration of fixations to  
235 the speaker's mouth, with more and longer fixations on the mouth relating to better performance  
236 (i.e., higher accuracy and a greater amount of improvement over time).

### 237 Method

#### 238 Participants

239 Seventy young adults (10 male, *Mdn* = 23 years, age range 19-30 years) were initially  
240 recruited from the University of Manchester to participate in the study, which was approved by the  
241 university ethics committee. All participants were native British English speakers with no history of  
242 neurological, speech or language problems (self-declared), and gave their written informed consent.  
243 Participants were included if their corrected binocular vision was 6/6 or better using a reduced  
244 Snellen chart, and their stereoacuity was at least 60 seconds of arc using a TNO test. Participants'  
245 hearing was measured using pure-tone audiometry for the main audiometric frequencies of speech  
246 (0.5, 1, 2, and 4 kHz) in each ear separately. Any participant with a hearing threshold level greater  
247 than 20dB for more than one frequency in either ear was excluded from participation. Eleven  
248 participants in total (one male) were excluded; two based on the hearing criteria, two based on the  
249 visual criteria, five due to data loss during the eye tracking procedure (see Data Analysis for full  
250 details), one due to poor eye tracking calibration, and one due to technical failure. 59 participants  
251 (nine male, *Mdn* = 23 years, age range 19-30 years) were thus included in the final analyses reported  
252 here. Our sample size was based on the expected effect size for the audiovisual benefit to  
253 adaptation. Pilling & Thomas (2011) observed a 'benefit' of 12% accuracy for adaptation to

254 audiovisual compared to audio-only degraded sentences using a similar measure of keywords to the  
255 present study, although insufficient statistics were reported to obtain an effect size. Bernstein et al.  
256 (2013) observed a large effect size of  $d = 1.21$  for adaptation to degraded syllables; as our task was  
257 more challenging, we predicted a medium-sized effect. Brysbaert & Stevens (2018) recommend a  
258 minimum of 1600 observations per cell for linear mixed effect models detecting medium-sized  
259 effects, which we achieved with 60 keywords per testing block, and at least 29 participants per  
260 group (i.e., we had at least 1740 observations per cell).

### 261 **Materials**

262 Experimental materials are available at <https://osf.io/2wqkf/>. Our stimuli consisted of 91 randomly  
263 selected Institute of Electrical and Electronics Engineers Harvard sentences (IEEE; Rothauser et al.,  
264 1969). As we wanted to compare our adaptation results as far as possible to Pilling & Thomas (2011),  
265 we selected 4 keywords per sentence to score participant accuracy. These were content and  
266 function words, selected by the experimenters, that were considered important to the meaning of  
267 each sentence. A list of the sentences and keywords used is available as supplemental materials at  
268 the above link. Recordings were carried out in a soundproofed laboratory using a Shure SM58  
269 microphone and a High Definition Canon HV30 camera. A 26-year-old female native British English  
270 speaker recited the sentences, and was asked to look directly at the camera, to remain still, and to  
271 maintain a neutral facial expression throughout the recordings to minimise head movement. Video  
272 recordings were imported into iMovie 11 running on an Apple MacBook Pro, as large (960 x 540)  
273 high-definition digital video (.dv) files. Recordings were edited to create individual video clips for  
274 each sentence. These were checked by the experimenter and any that were not deemed suitable  
275 (for example due to mispronunciation) were re-recorded. The audio tracks for each clip were  
276 extracted as audio (.wav) files, then normalised by equating the root mean square amplitude,  
277 resampled at 22 kHz in stereo, cropped at the nearest zero crossings at voice onset and offset, and  
278 vocoded using Praat speech processing software (Boersma & Weenink, 2018). Speech recordings  
279 were noise-vocoded (Shannon et al., 1995) using four frequency bands (cut-offs: 50 Hz → 369 Hz →

280 1160 Hz → 3124 Hz → 8000 Hz), selected to represent equal spacing along the basilar membrane  
281 (Greenwood, 1990). In the audio-only (control) condition, a static image of the speaker's face with  
282 the mouth in different "speaking" positions was displayed congruently with the audio files so that a  
283 visual component was also present in this condition, but with no useful linguistic information. Static  
284 faces have previously been used as a control condition for analysing speech perception in dynamic  
285 faces (e.g., Calvert & Campbell, 2003; C. Davis & Kim, 2004; Jerger Susan et al., 2018). Using a static  
286 face as a control allowed us to assess the contribution of visible articulatory cues to speech  
287 recognition, whilst controlling for visual attention towards any salient features of the speaker's face,  
288 and also allowing for eye tracking to be conducted in both groups for consistency. To create the still  
289 images (one image per trial), screen shots saved as TIFF files were taken from the videos of the  
290 speaker displaying a variety of mouth positions, to make the mouth visually salient and to make it  
291 evident that she was speaking. The still images, video files and the noise-vocoded audio files were  
292 imported into Experiment Builder software (SR Research, Ontario, Canada) to create the  
293 experimental stimuli. In the audio-only condition, the still images of the speaker were displayed for  
294 the exact length of each audio file, and for the audiovisual condition the audio and video files were  
295 played congruently.

## 296 **Procedure and apparatus**

297 Data were collected in a soundproofed booth in a single test lasting approximately 40  
298 minutes. Participants were randomly allocated into either the audiovisual ( $N=30$ ) or audio-only  
299 ( $N=29$ ) control group. In both conditions, participants sat facing the screen approximately 50 cm  
300 from the monitor, with their chin on a chin-rest. They were asked not to move their head during the  
301 experiment and to look continuously at the screen. Before starting the experiment, the eye-tracker  
302 was calibrated for each participant (see 'Data analysis' for details). Participants first listened to one  
303 practice sentence (a clear version and a noise-vocoded version) that was not included in the  
304 experiment, to prepare them for hearing the unusual distortion. They then completed 90 trials with  
305 the remaining noise-vocoded sentences. Participants triggered the start of the experiment and each

306 subsequent trial by pressing the space bar on the keyboard; there were no structured breaks and all  
307 90 trials were presented in a single continuous session. All stimuli were presented through  
308 Sennheiser HD 25-SP II headphones. The experimenter set the volume for all stimuli at a  
309 comfortable level for the first participant, and kept it at the same level for all participants thereafter.  
310 A Panasonic lapel microphone attached to the chin-rest recorded their verbal responses.

311 To measure speech recognition, we asked participants to repeat out loud as much of each  
312 sentence as they could. The experimenter retrospectively scored participants' responses according  
313 to how many keywords they correctly repeated out of a maximum of four. Responses were scored  
314 as correct despite incorrect suffixes (such as -s, -ed, -ing) or verb endings; however if only part of a  
315 word (including compound words) was repeated this was scored as incorrect (Dupoux & Green,  
316 1997; Golomb et al., 2007).

317 We used a desktop-mounted EYELINK 1000 eye-tracker with Experiment Builder software (SR  
318 Research, Ontario, Canada) to present all stimuli, and to record participants' eye movements. The  
319 pupil and corneal reflection of each participant's right eye were tracked at a sample rate of 1000 Hz,  
320 with a spatial resolution of 0.01° RMS and average accuracy of 0.25°–0.5°. Calibration was carried  
321 out for each participant before the experiment using a standard nine-point configuration, and again  
322 five minutes after the experiment began. Each calibration was validated for accuracy, and accepted  
323 if the average error was <1° and the maximum error was <1.5°. A drift check preceded each trial  
324 using a fixation point presented in the centre of the screen, and if the error between the computed  
325 fixation position and the on-screen target was >1.5°, calibration was repeated to correct this drift.

### 326 **Data analysis**

327 The dependent variables were recognition accuracy, fixation duration, and percentage  
328 fixations. Recognition accuracy was calculated as the percentage of keywords correctly repeated in  
329 each trial. To analyse recognition accuracy over time, we divided all consecutive trials into six blocks  
330 of 15 trials, and calculated mean percentage accuracy per testing block based on the number of

## EYE GAZE AND ADAPTATION

331 correctly repeated keywords<sup>1</sup>. Fixations were defined as any period that was not a saccade (saccades  
332 were defined as eye movements with velocity  $>30^\circ/\text{sec}$ , acceleration  $>8000^\circ/\text{sec}^2$ , and motion  
333  $>0.1^\circ$ ). Fixations were evaluated in relation to one of two regions of interest (ROIs). For each video  
334 clip, we created two elliptical ROIs (see Figure 1) based on the first video frame. These comprised  
335 the eye area (extending from just below the speaker's eyebrows to the tip of the nose) and the  
336 mouth area (from the septum to just below the bottom lip). Fixation duration and percentage  
337 fixations in these regions were then analysed to compare patterns of eye gaze between the two  
338 ROIs. We also created a third interest area that surrounded the speaker's face that was used to  
339 verify the proportion of eye gaze directed to the speaker's face rather than peripheral areas of the  
340 screen. *Fixation duration* was calculated as the mean duration of fixations in milliseconds.  
341 *Percentage fixations* was calculated as the percentage of all fixations in a trial falling in the current  
342 ROI. We selected these variables to indicate where listeners were allocating their attention at  
343 particular time points. Measurements of eye gaze were computed using Data Viewer (SR Research,  
344 Ontario, Canada), and we calculated the mean of each variable per testing block, and per interest  
345 area.

346 Data were analysed using linear mixed effects hierarchical regression models in the lmerTest  
347 package (Kuznetsova et al., 2017), which uses the lme4 package, running in R v3.4.1. All models  
348 included the random effect of participant to account for individual differences in baseline speech  
349 recognition. Fixed effects of group, ROI and testing block (i.e., time) were tested by comparing  
350 models pairwise using likelihood ratio tests and Bayes Factors calculated using the BIC (e.g.,  
351 Wagenmakers, 2007). For effects of individual predictors within the model, beta ( $B$ ) coefficients and  
352 estimated  $p$ -values are reported. The variable of fixation duration was rescaled (ms/1000) to make  
353 the coefficient more interpretable; estimates of this variable are therefore expressed in seconds.  
354 --Include Figure 1 about here--

---

<sup>1</sup> Trials were only divided into testing blocks during data analysis – i.e., participants were not aware of the testing blocks during the procedure.

355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380

## Results

### Perceptual adaptation to noise-vocoded speech

Figure 2 shows mean recognition accuracy for the noise-vocoded speech across the six testing blocks for each group. We first tested for group effects against the baseline random effect of participant. Recognition was overall significantly better in the audiovisual group ( $M = 54\%$ ,  $SD = 2.0\%$ ) compared to the audio-only group ( $M = 35\%$ ,  $SD = 1.6\%$ ),  $B = 19.48$ ,  $SE = 2.53$ ,  $p < 0.001$ ;  $\chi^2 = 41.02$ ,  $p < .001$ ,  $BF_{10} = 42952865$ . We then added a group \* testing block interaction to the model to test whether the audiovisual group improved more over the six testing blocks than the audio-only group. The comparison was significant,  $\chi^2 = 145.45$ ,  $p < .001$ , and the large Bayes Factor indicated strong evidence in favour of including the interaction in the model,  $BF_{10} = 6.911289e+18$ . However, across the whole experiment (i.e., between block 1 and block 6), recognition accuracy increased equally in both groups by approximately 19%,  $B = 18.68$ ,  $SE = 1.56$ ,  $p < .001$ .

### Exploratory Analysis: Rate of Adaptation

Although we observed a group\*testing block interaction, results of the mixed effects model described above indicated that the only significant difference in adaptation occurred between block 1 and block 5, where the audiovisual group adapted by 18.47% compared to 12.51% in the audio-only group,  $B = 6.69$ ,  $SE = 3.08$ ,  $p = 0.031$ . This suggested that listeners adapted more rapidly in the audiovisual group. To examine the rate of adaptation across the experiment in more detail, we conducted exploratory analyses of the amount of adaptation between groups for each consecutive pair of testing blocks. Figure 3 shows that the rate of adaptation was not consistent between blocks or groups. Most adaptation occurred during exposure to the first 30 sentences, when both groups showed ~9% improvement in recognition accuracy. Between blocks 2-5 adaptation slowed in both groups, but the audiovisual group consistently adapted slightly faster, improving by approximately 9% compared to only 2% in the audio-only group. However, between blocks 5 and 6 the audio-only group adapted more than the audiovisual group, improving by 6.4% compared to <1% in the audiovisual group.



381 We conducted exploratory Bayesian hierarchical regression analyses of adaptation to  
382 quantify the evidence for group differences in adaptation rate between consecutive testing blocks.  
383 We used forward difference coding whereby a contrast variable was calculated for each pair of  
384 consecutive blocks (e.g., B1-B2, B2-B3 etc.), representing differences in recognition accuracy  
385 between each pair of blocks. The resulting five coded variables were added as fixed effects to a  
386 baseline model that also included group as a main fixed effect, and participant as a random effect.  
387 The interaction between each coded variable and group (e.g., B1-B2\*group, which represents group  
388 differences in adaptation between blocks 1 and 2) was added individually and compared to the  
389 baseline model to test for group differences in adaptation at different time points. As these were  
390 exploratory analyses we report Bayes Factors and effect sizes only (see Table 1). The baseline model  
391 of adaptation between each consecutive pair of testing blocks, and a main effect of group,  
392 accounted for approximately 46% variance in recognition accuracy. Bayes factors indicated that  
393 there was either no evidence ( $BF < 0.3$ ), or inconclusive evidence ( $BF > 0.3 < 1$ ), of a difference in  
394 adaptation between groups for each consecutive pair of testing blocks, and indeed, adding the  
395 interaction variables increased the explained variance by a maximum of just 0.3% (for the B2-  
396 B3\*Group interaction).

397 --Include Figures 2 and 3, and Table 1 about here--

### 398 **Patterns of Eye Gaze**

399 We first examined overall patterns of eye gaze in both groups, to establish whether our eye  
400 tracking methods and stimuli had successfully replicated the patterns of eye gaze frequently seen in  
401 studies of audiovisual speech perception and when viewing static faces; particularly, to confirm that  
402 there were no unusually salient features in our stimuli that attracted viewer's visual attention. In the  
403 audiovisual group, 99% of all fixations fell on the speaker's face and 98% fell on the eyes and mouth.  
404 In line with previous studies of audiovisual speech recognition in difficult listening conditions  
405 (Buchan et al., 2007, 2008; Lansing & McConkie, 2003; Vatikiotis-Bateson et al., 1998), fixations on  
406 the speaker's mouth ( $M = 984.32\text{ms}$ ,  $SD = 405\text{ms}$ ) were significantly longer than fixations on the

407 eyes ( $M = 363.37\text{ms}$ ,  $SD = 164\text{ms}$ ),  $\chi^2 = 350.83$ ,  $p < .001$ ,  $BF_{10} = 8.024141\text{e}+74$ ,  $B = 0.621$ ,  $SE = 0.02$ ,  
408 confirming that, as expected, listeners attended more to the speaker's mouth than the eyes.  
409 However, there was no difference in percentage fixations on the mouth ( $M = 49\%$ ,  $SD = 18\%$ ) and  
410 eyes ( $M = 49\%$ ,  $SD = 18\%$ ),  $\chi^2 = 0$ ,  $p = .988$ ,  $BF_{10} = 0.05$ .

411 In the audio-only group, 83% of fixations were located on the speaker's face, with 74% on  
412 the eyes and mouth. The duration of fixations on the eyes ( $M = 443.46\text{ms}$ ,  $SD = 179\text{ms}$ ) and mouth  
413 ( $M = 443.30\text{ms}$ ,  $SD = 189\text{ms}$ ) did not differ,  $\chi^2 = 0$ ,  $p = .980$ ,  $BF_{10} = 0.05$ . However, a higher  
414 percentage of fixations fell on the eyes ( $M = 65\%$ ,  $SD = 21\%$ ) than on the mouth ( $M = 18\%$ ,  $SD = 17\%$ ),  
415  $\chi^2 = 315.59$ ,  $p < .001$ ,  $BF_{10} = 1.818774\text{e}+67$ ,  $B = -0.47$ ,  $SE = 0.02$ ,  $p < 0.001$ , in line with previous  
416 results from viewing static faces (e.g., Birmingham & Kingstone, 2009). As there were no useful  
417 visual cues available in the audio-only group that could benefit speech recognition, and the stimuli  
418 was not dynamic, we did not analyse this data in relation to speech recognition; however all data is  
419 available as supplemental material here: <https://osf.io/2wqkf/>.

#### 420 ***Are audiovisual speech recognition and perceptual adaptation related to patterns of eye gaze?***

421 To test this hypothesis, we analysed speech recognition data from the audiovisual group,  
422 first establishing a baseline model of adaptation with testing block as a predictor; compared to a  
423 random effects model of participants' baseline accuracy, there was strong evidence for the baseline  
424 model of adaptation to the noise-vocoded speech:  $\chi^2 = 84.53$ ,  $p < .001$ ,  $BF_{10} = 5.22229\text{e}+12$ . We then  
425 compared this baseline model to four experimental models, each of which included one of the  
426 following eye tracking measures as a predictor variable: 1) duration of fixations on the mouth; 2)  
427 duration of fixations on the eyes; 3) percentage fixations on the mouth, and 4) percentage fixations  
428 on the eyes (see Table 2 for models and corresponding  $R^2$  values). Only the model including duration  
429 of fixations on the mouth was significantly different to the baseline model,  $\chi^2 = 5.47$ ,  $p = 0.019$ ;  
430 longer fixations on the speaker's mouth were related to better recognition of the noise-vocoded  
431 sentences,  $B = 7.68$ ,  $SE = 3.21$ ,  $p = 0.018$ , however, evidence in support of this relationship was  
432 relatively weak ( $BF_{10} = 1.15$ ). We then tested for an interaction between testing block and the

433 duration of fixations on the mouth to ascertain whether the duration of fixations could predict  
434 adaptation. The results did not support the presence of an interaction,  $\chi^2 = 9.17$ ,  $p = 0.102$ ,  $BF_{10} =$   
435  $0.0002$ , indicating that there was no overall relationship between eye gaze and adaptation over the  
436 course of the experiment.

437 --Include Table 2 about here--

#### 438 ***Exploratory Analyses: Changes in Eye Gaze Over Time***

439 As speech recognition and adaptation rate varied across the time course of the experiment,  
440 we conducted exploratory analyses to examine whether patterns of eye gaze in the audiovisual  
441 group, as well as their relationship with speech recognition, varied over time. As before, we used  
442 Bayesian hierarchical linear mixed effects models, comparing the inclusion of each experimental  
443 predictor to a baseline model with participant as a random effect. As these were exploratory  
444 analyses we report descriptive statistics, effect sizes and Bayes Factors only. Figure 4 shows the  
445 mean duration of fixations and percentage fixations over the time course of the experiment. There  
446 was strong evidence that the duration of fixations on the mouth decreased over time by an average  
447 of 268.77ms between block 1 and block 6 ( $BF_{10} = 7522.16$ ,  $B = -0.26877$ ,  $SE = 0.04256$ , marginal  $R^2 =$   
448  $0.05$ ). There was no evidence that the duration of fixations on the eyes changed over time ( $BF_{10} =$   
449  $0.0002$ , marginal  $R^2 = 0.01$ ), nor percentage fixations on the mouth ( $BF_{10} = 0.0003$ , marginal  $R^2 =$   
450  $0.01$ ) or the eyes ( $BF_{10} = 0.0004$ , marginal  $R^2 = 0.01$ ).

451 Based on the variability in speech recognition, amount of adaptation and the duration of  
452 fixations on the speaker's mouth over time, it was possible that longer fixations on the speaker's  
453 mouth were more useful at particular time points of the experiment than others, for example during  
454 earlier testing blocks. We therefore explored whether the duration of fixations on the speaker's  
455 mouth were related to speech recognition in early (blocks 1-2), middle (blocks 3-4) or late (blocks 5-  
456 6) testing blocks. For each time period, we compared a model including the duration of fixations on  
457 the mouth to the baseline random effects model. We found evidence for a relationship between  
458 speech recognition and the duration of fixations on the mouth for middle testing blocks (blocks 3-4)

459 only,  $BF_{10} = 19.90$ ,  $B = 18.08$ ,  $SE = 5.42$ , marginal  $R^2 = 0.21$ ; conversely, we found evidence *against* a  
460 relationship between speech recognition and the duration of fixations on the mouth in early (blocks  
461 1-2:  $BF_{10} = 0.13$ , marginal  $R^2 = 0.001$ ), and late blocks (blocks 5-6:  $BF_{10} = 0.18$ , marginal  $R^2 = 0.02$ ).  
462 --Figure 4 about here--

### 463 Discussion

464 We investigated perceptual adaptation to noise-vocoded speech with and without visual speech  
465 cues, aiming to replicate and extend previous findings (Bernstein et al., 2013; Kawase et al., 2009;  
466 Pilling & Thomas, 2011) that being able to view a speaker's face can lead to greater improvement in  
467 recognition over time. We used a real-time (i.e., continuous) adaptation paradigm to better reflect  
468 real-life adaptation, and eye tracking to investigate eye gaze patterns during audiovisual speech  
469 recognition. We tested the relationship between performance and the duration and percentage of  
470 fixations on the speaker's eyes and mouth, predicting that looking more at the speaker's mouth  
471 would be related to better recognition accuracy and greater adaptation.

472 We partially replicated previous studies which found an audiovisual benefit to perceptual  
473 adaptation, but our observations are somewhat more complex. There was a clear overall benefit to  
474 speech recognition from the visual speech cues, with accuracy in the audiovisual group consistently  
475 ~20% better than in the audio-only group. However, we found no overall difference in the amount of  
476 adaptation between groups as expected – by the final testing block (i.e., after exposure to all 90  
477 sentences), both groups had improved by ~19% accuracy overall. Instead, we only observed a  
478 difference between blocks 1 and 5 (after exposure to 75 sentences). Exploratory analyses suggested  
479 that the rate of adaptation between blocks varied across the experiment, with the greatest amount  
480 of adaptation within the first 30 trials in both groups, who initially adapted at an equal rate despite  
481 different baseline levels of accuracy. After this point, the audiovisual group adapted slightly faster  
482 until testing blocks 5 and 6, when the audio-only group improved more quickly. However, in  
483 Bayesian terms, there was no evidence for group differences in adaptation rate between most

484 blocks, although evidence was inconclusive between testing blocks 2 and 3. Overall, the benefit from  
485 visual speech cues to adaptation to degraded speech in our data is smaller and less clear than  
486 expected; particularly, we expected the audiovisual group to adapt more overall than the audio-only  
487 group.

488 Our findings are in contrast to studies which found a clear audiovisual benefit to adaptation  
489 for noise-vocoded syllables (Bernstein et al., 2013), words (Kawase et al., 2009), and sentences  
490 (Pilling & Thomas, 2011); these studies all found greater overall adaptation when the speaker's face  
491 was visible compared to when it was not. However, there is some similarity between our findings  
492 and those of Pilling & Thomas (2011); we found that adaptation was greater in our audiovisual group  
493 following exposure to 75 sentences (testing block 5), while Pilling & Thomas observed the same  
494 effect after a similar amount of exposure (76 sentences) during an audiovisual training period.  
495 Nevertheless, we did not predict that the audiovisual benefit to adaptation would only be limited to  
496 the fifth testing block, and this finding could therefore be due to chance.

497 There are several possible conclusions from our data. First, that providing a specific  
498 audiovisual training period (as in Pilling & Thomas, 2011) is more effective than real-time  
499 adaptation; this may, for example, be due to participants attending more to audiovisual speech cues  
500 during a separate period of training, in comparison to continuous exposure which may result in  
501 lessened attention or fatigue; indeed, the rate of adaptation slowed considerably for our audiovisual  
502 group between the final two testing blocks. Second, the amount of benefit to adaptation gained  
503 from visual speech cues may depend on the type of stimuli, whereby a greater benefit is possible  
504 with simpler and more predictable linguistic items, or from particular speakers (Blackburn et al.,  
505 2019). Indeed, using the linguistically more complex IEEE sentences, we observed less improvement  
506 in our audiovisual condition (19%) than with the BKB sentences used by Pilling & Thomas (26%) even  
507 after greater exposure, although this difference could also be explained by the different speakers  
508 used in each study. Lastly, visual speech cues may in fact lead to *faster* adaptation rather than

509 greater overall improvement; that is, without visual cues listeners can still adapt equally well but  
510 require more exposure to do so, as was the case for our audio-only group. Our exploratory analyses  
511 of adaptation rate seem to support this, as speech recognition rapidly improved in both groups  
512 initially, but then slowed in the audio-only condition; however, this group difference was small, and  
513 the Bayesian evidence from our data didn't support a clear difference in adaptation rate. The  
514 amount of adaptation observed may thus depend on exactly *when* it is measured, and how much  
515 exposure participants have had to the degraded speech.

516 Overall, our results indicate that the benefits of visual speech cues to adaptation are not as  
517 **great or clearcut** as results from previous studies suggest. Instead, the benefits potentially depend  
518 on factors such as the linguistic items used (i.e., the specific linguistic characteristics of the stimuli  
519 such as length, syntactic complexity or semantic predictability), speaker, and amount of exposure,  
520 and the contribution of these factors will need to be confirmed in future studies. The small  
521 advantage to adaptation in the audiovisual group during middle testing blocks suggests that benefits  
522 from visual cues could further be related to participants' attention or energy levels, whereby visual  
523 cues are particularly beneficial to learning at points where attention and motivation are low – such  
524 as in the middle of a challenging laboratory experiment. The benefits of visual cues in real-life  
525 contexts may thus depend on the type of communication taking place; while these cues do not  
526 necessarily lead to greater adaptation early on, they may be particularly useful in contexts where  
527 longer periods of sustained adaptation are required, for instance, listening to a lecture or when  
528 participating in a longer conversation. The interaction between use of visual speech cues and  
529 attention or fatigue may thus be an interesting line for future research into speech recognition in  
530 adverse listening conditions. **Nevertheless, the small audiovisual benefit that we observed during**  
531 **middle testing blocks could just have been an anomaly – i.e., it could have occurred by chance.**

532 **It should be noted that recognition of noise-vocoded sentences (with or without visual cues)**  
533 **varies considerably between studies. We observed mean performance of 35% accuracy in our audio-**

534 only condition, but similar studies have found differing levels of performance. For example, using 4-  
535 band noise-vocoding and the IEEE sentences (as in the present study), McGettigan et al., (2014)  
536 observed approximately 40% mean accuracy for recognition of only 10 sentences; however, this was  
537 following exposure to 70 noise-vocoded BKB sentences, perhaps accounting for the higher level of  
538 accuracy than in the present study. In comparison, using 6-band noise-vocoding, Paulus et al. (2020)  
539 observed approximately 60% accuracy after exposure to 48 IEEE sentences. Using the simpler BKB  
540 sentences, Scott et al., (2006) observed approximately 40% accuracy using 4-band noise-vocoding,  
541 but after exposure to only 16 sentences, while Rosen et al., (1999) observed 64% mean accuracy  
542 after exposure to 112 sentences also vocoded with 4 channels. Thus, recognition of noise-vocoded  
543 speech can vary greatly depending on the amount of exposure, the type of linguistic stimuli, and the  
544 exact vocoding transformation. In the present study, we specifically chose to use the IEEE sentences  
545 and 4-band noise-vocoding to create a more challenging task (and particularly to prevent ceiling  
546 effects in the audiovisual condition). Nevertheless, the intelligibility of our stimuli may also have  
547 been affected by the speaker we used (e.g., Bradlow & Bent, 2008). Indeed, specific acoustic-  
548 phonetic features (namely vowel space dispersion and mean energy in mid-range frequencies) can  
549 account for differing levels of intelligibility between speakers for noise-vocoded speech, although  
550 these features do not necessarily impact listeners' amount of adaptation (Paulus et al., 2020).  
551 Furthermore, the amount of benefit that visual cues can provide also varies between speakers  
552 (Blackburn et al., 2019). As changing speakers can interfere with adaptation (e.g., Dupoux & Green,  
553 1997), we used the same speaker throughout our study. However, we note that a limitation of the  
554 current findings is that we cannot confirm whether mean levels of performance in either condition,  
555 or indeed the benefit that listeners obtained from the speaker's visual cues, would be the same for  
556 other speakers.

557           The second aim of our study was to examine patterns of eye gaze during adaptation to  
558 audiovisual degraded speech, and specifically to test whether there is a direct relationship between  
559 eye gaze towards a speaker's mouth movements and speech recognition. We found that longer

560 fixations on the speaker's mouth were related to better recognition, but not to the amount of  
561 adaptation. This supports findings from speechreading (Worster et al., 2018) which found that  
562 longer time spent fixating the speaker's mouth was related to better speechreading in both deaf and  
563 normal-hearing children. Two previous studies have also directly tested the relationship between  
564 eye gaze patterns and speech recognition (Buchan et al., 2007; Everdell et al., 2007), but found no  
565 significant relationship. However, methodological differences can potentially account for the  
566 different results reported here. First, audiovisual speech recognition was at ceiling in both studies,  
567 i.e., 86% (Buchan et al., 2007) and 90% (Everdell et al., 2007), compared to 41-61% in the present  
568 study. Second, neither study analysed the duration of fixations (as in the present study), or time  
569 spent fixating the speaker's mouth (as in Worster et al., 2018). Everdell et al. (2007) analysed an  
570 index of left-right asymmetry of eye gaze on the eyes and mouth, while Buchan et al. (2007)  
571 analysed percentage trials spent looking at the speaker's mouth, but neither observed correlations  
572 between these measures and speech recognition. Current evidence thus suggests that  
573 measurements of the *time* spent fixating a speaker's mouth is indicative of effective use of visual  
574 speech cues, rather than the frequency or proportion of fixations; indeed, we found no correlation  
575 between percentage fixations on the speaker's mouth and speech recognition, similar to Lansing &  
576 McConkie (2003) who found no relationship between the number of fixations on the mouth and  
577 speechreading. More recently, Lusk & Mitchell (2016) observed a positive relationship between  
578 changes in the amount of eye gaze on a speaker's mouth during passive listening to an artificial  
579 language, and subsequent segmentation of non-words from this language. However, note that Lusk  
580 & Mitchell's finding only partially supports the current findings, as the relationship was irrespective  
581 of direction – i.e., the shift could involve looking more or less at the mouth. Thus, to our knowledge,  
582 ours is the first study to observe a direct relationship between looking more at a speaker's mouth  
583 and audiovisual speech recognition.

584           The results add to a growing body of literature indicating that patterns of eye gaze – that is,  
585 where and how listeners look at a speaker's face – are important for successfully understanding



586 unfamiliar or degraded audiovisual speech. Thus, it is not merely the presence of visual speech cues,  
587 but also the particular visual strategies employed by listeners, that relate to successful speech  
588 recognition. As we compared two measures of eye gaze commonly used in eye tracking studies, we  
589 can further conclude that the *duration* of fixations on a speaker's mouth are likely more important  
590 than the proportion of fixations. Longer fixations on the mouth likely reflected a greater focus of  
591 attention on this region, particularly as visual perception is reduced during eye movements (Matin &  
592 Ethel, 1974). Thus, with longer fixations and less eye movement, listeners could better or more  
593 efficiently decode articulatory cues from a speaker's mouth, improving recognition. The duration of  
594 fixations on a speaker's mouth is thus potentially a useful measure when assessing the use or  
595 relevance of visual speech cues. Indeed, longer fixations on a speaker's mouth have indicated  
596 increased use of visual cues in other studies of adverse listening conditions (Buchan et al., 2007,  
597 2008), although the measure has not previously been related to performance. The importance of  
598 this measure was indirectly supported by our exploratory observation that the duration of fixations  
599 on the speaker's mouth decreased over time, as performance improved (while no such change was  
600 observed for percentage fixations). This decrease would suggest that participants' use of visual cues  
601 from the speaker's mouth decreased as they adapted to the degraded speech. A similar observation  
602 was made by Lusk & Mitchel (2016) who noted a decrease in overall gaze time on a speaker's  
603 mouth, but not on the eyes or nose, during a period of familiarisation to an artificial language (i.e.,  
604 passive listening/viewing), prior to listeners being tested on non-word recognition. The duration of  
605 fixations on a speaker's mouth may thus be an important indicator of effective use of visual speech  
606 cues when learning or adapting to unfamiliar speech – for example helping word segmentation  
607 (Mitchel & Weiss, 2014); however, we did not observe a correlation between the duration of  
608 fixations and amount of adaptation.

609 Another interpretation of our finding is that the decrease in fixation durations indicates  
610 changes in attention or effort. After the period of rapid adaptation between testing blocks 1 and 2,  
611 decoding the noise-vocoded speech perhaps no longer required as much cognitive effort, or

612 attention, from participants. Listening effort (as measured by relative pupil size) is greater during  
613 perception of noise-vocoded speech compared to undegraded speech in quiet (Paulus et al., 2020);  
614 furthermore, it has been shown to decrease during a period of adaptation to unfamiliar accented  
615 speech (Brown et al., 2020), just as the duration of fixations decreased in our study. An  
616 interpretation of our results related to cognitive effort is compatible with those of Birulés et al.  
617 (2020), who found that listeners looked more towards a speaker's mouth (measured as proportion  
618 of total gaze time) during recognition of non-native speech than native, regardless of linguistic  
619 ability; that is, the cognitive demands required to understand non-native speech were consistently  
620 greater – indicated by more time spent looking at the speaker's mouth (and potentially greater  
621 reliance on visual speech cues). Outside of the speech perception literature, changes in eye gaze  
622 patterns have also been associated with cognitive load; for example, fewer and longer fixations  
623 during scene viewing are observed with greater memory loads (Cronin et al., 2020), again suggesting  
624 that greater cognitive demands can influence patterns of eye gaze. Although the present results  
625 cannot confirm this interpretation of our data, they nonetheless offer an interesting avenue for  
626 future research.

627       Some limitations to the current findings should be noted. First, the evidence for a  
628 relationship between eye gaze and speech recognition was relatively weak in Bayesian terms.  
629 Exploratory analyses suggested that the relationship was in fact only present in middle testing  
630 blocks, but why this would be the case is unclear; the pattern somewhat matches our observation  
631 that audiovisual cues were most beneficial to adaptation during middle testing blocks, rather than in  
632 early or later blocks, and so could indicate a particular reliance on visual cues during this time. Visual  
633 cues from the speaker's mouth could potentially serve to compensate for decreasing attention or  
634 motivation, resulting in a stronger relationship between longer fixations and performance during  
635 this period. Nevertheless, the results require further testing. A second limitation is that the result  
636 was correlational, and we therefore cannot ascertain whether longer fixations on the speaker's  
637 mouth resulted in better recognition, or whether participants who performed better looked more

638 steadily at the speaker's mouth. Again, this correlational result would benefit from further testing  
639 whereby particular eye gaze strategies are manipulated to observe the effects on performance.  
640 Finally, we note that using a static face as a control condition for the audio-only condition is less  
641 naturalistic than, for example, providing no visual information at all, and thus does not have an exact  
642 'real-world' equivalent (except, perhaps, a frozen screen during a video call). Our motivation in  
643 including this condition was to equate the procedure for both groups as far as possible, including  
644 visual information and eye tracking in both. However, we are confident that performance in this  
645 condition was not significantly worse than would be expected without a visible static face (for an  
646 online replication see Trotter et al., 2020), and thus that it was a valid comparison for speech  
647 adaptation.

648 We report several exploratory analyses in the current paper to support interpretation of the  
649 findings, and these are intended as hypothesis-generating observations rather than hypothesis-  
650 testing, whereby our aim is to open up further lines of enquiry regarding adaptation to unfamiliar  
651 speech and related patterns of eye gaze. For example, the decrease in the duration of fixations  
652 during adaptation may be further investigated by comparing eye gaze during audiovisual speech  
653 recognition to a control condition with non-informative mouth movements, or compared to  
654 measures of listening effort. Furthermore, differences in the rate of adaptation to unfamiliar speech  
655 with and without visual cues should be investigated in more detail to establish the exact parameters  
656 that determine when visual cues offer a clear benefit to listeners. The analyses and observations  
657 presented here will thus be beneficial to the research fields of audiovisual speech perception and,  
658 more broadly, communication in difficult listening conditions.

## 659 **Conclusion**

660 We have demonstrated that the benefit of visual speech cues to adaptation to degraded (noise-  
661 vocoded) speech is more limited than previously thought – potentially resulting in slightly faster  
662 adaptation only after a period of initial exposure and rapid adaptation, but not resulting in an overall

## EYE GAZE AND ADAPTATION

663 greater amount of improvement after a longer period of exposure. Longer fixations on the speaker's  
664 mouth were related to better overall recognition accuracy of the audiovisual speech, adding to a  
665 growing body of evidence that patterns of eye gaze are related to effective use of visual speech  
666 cues. Nevertheless, evidence for this relationship was relatively weak and will need further testing to  
667 be fully confirmed and understood. We further observed that the duration of fixations on the  
668 speaker's mouth decreased over time; future research will need to determine the relevance of this  
669 finding, as well as whether particular patterns of eye gaze can intentionally bring benefits to  
670 listeners in adverse listening conditions.

671

672

**Bibliography**

- 673 Adank, P., & Janse, E. (2010). Comprehension of a novel accent by young and older listeners.  
674 *Psychology and Aging, 25*(3), 736–740. <https://doi.org/10.1037/a0020054>
- 675 Alsius, A., Paré, M., & Munhall, K. G. (2018). Forty Years After Hearing Lips and Seeing Voices: The  
676 McGurk Effect Revisited. *Multisensory Research, 31*(1–2), 111–144.  
677 <https://doi.org/10.1163/22134808-00002565>
- 678 Banks, B., Gowen, E., Munro, K. J., & Adank, P. (2015a). Audiovisual cues benefit recognition of  
679 accented speech in noise but not perceptual adaptation. *Frontiers in Human Neuroscience,*  
680 *9*(AUGUST). <https://doi.org/10.3389/fnhum.2015.00422>
- 681 Banks, B., Gowen, E., Munro, K. J., & Adank, P. (2015b). Cognitive predictors of perceptual  
682 adaptation to accented speech. *Journal of the Acoustical Society of America, 137*(4).  
683 <https://doi.org/10.1121/1.4916265>
- 684 Barenholtz, E., Mavica, L., & Lewkowicz, D. J. (2016). Language familiarity modulates relative  
685 attention to the eyes and mouth of a talker. *Cognition.*  
686 <https://doi.org/10.1016/j.cognition.2015.11.013>
- 687 Bench, J., Kowal, Å., & Bamford, J. (1979). The Bkb (Bamford-Kowal-Bench) Sentence Lists for  
688 Partially-Hearing Children. *British Journal of Audiology, 13*(3), 108–112.  
689 <https://doi.org/10.3109/03005367909078884>
- 690 Bernstein, L. E., Auer, E. T., Eberhardt, S. P., & Jiang, J. (2013). Auditory Perceptual Learning for  
691 Speech Perception Can be Enhanced by Audiovisual Training. *Frontiers in Neuroscience, 7,*  
692 *34.* <https://doi.org/10.3389/fnins.2013.00034>
- 693 Birmingham, E., & Kingstone, A. (2009). Human Social Attention. *Annals of the New York Academy of*  
694 *Sciences, 1156*(1), 118–140. <https://doi.org/10.1111/j.1749-6632.2009.04468.x>
- 695 Birulés, J., Bosch, L., Pons, F., & Lewkowicz, D. J. (2020). Highly proficient L2 speakers still need to  
696 attend to a talker’s mouth when processing L2 speech. *Language, Cognition and*  
697 *Neuroscience, 0*(0), 1–12. <https://doi.org/10.1080/23273798.2020.1762905>

- 698 Blackburn, C. L., Kitterick, P. T., Jones, G., Sumner, C. J., & Stacey, P. C. (2019). Visual Speech Benefit  
699 in Clear and Degraded Speech Depends on the Auditory Intelligibility of the Talker and the  
700 Number of Background Talkers. *Trends in Hearing*, *23*, 2331216519837866.  
701 <https://doi.org/10.1177/2331216519837866>
- 702 Boersma, P., & Weenink, D. (2018). Praat: Doing phonetics by computer. In *Retrieved October 2017*  
703 *from <http://www.praat.org/>* (6.0.31).
- 704 Bradlow, A. R., & Bent, T. (2008). Perceptual adaptation to non-native speech. *Cognition*, *106*(2),  
705 707–729. <https://doi.org/10.1016/j.cognition.2007.04.005>
- 706 Brown, V. A., McLaughlin, D. J., Strand, J. F., & Van Engen, K. J. (2020). Rapid adaptation to fully  
707 intelligible nonnative-accented speech reduces listening effort. *Quarterly Journal of*  
708 *Experimental Psychology*, 1747021820916726. <https://doi.org/10.1177/1747021820916726>
- 709 Brysbaert, M., & Stevens, M. (2018). Power Analysis and Effect Size in Mixed Effects Models: A  
710 Tutorial. *Journal of Cognition*, *1*(1). <https://doi.org/10.5334/joc.10>
- 711 Buchan, J. N., Paré, M., & Munhall, K. G. (2007). Spatial statistics of gaze fixations during dynamic  
712 face processing. *Social Neuroscience*, *2*(1), 1–13.  
713 <https://doi.org/10.1080/17470910601043644>
- 714 Buchan, J. N., Paré, M., & Munhall, K. G. (2008). The effect of varying talker identity and listening  
715 conditions on gaze behavior during audiovisual speech perception. *Brain Research*, *1242*,  
716 162–171. <https://doi.org/10.1016/j.brainres.2008.06.083>
- 717 Calvert, G. A., & Campbell, R. (2003). Reading Speech from Still and Moving Faces: The Neural  
718 Substrates of Visible Speech. *Journal of Cognitive Neuroscience*, *15*(1), 57–70.  
719 <https://doi.org/10.1162/089892903321107828>
- 720 Christianson, S. a, Loftus, E. F., Hoffman, H., & Loftus, G. R. (1991). Eye fixations and memory for  
721 emotional events. *Journal of Experimental Psychology. Learning, Memory, and Cognition*,  
722 *17*(4), 693–701. <https://doi.org/10.1037/0278-7393.17.4.693>

- 723 Cronin, D. A., Peacock, C. E., & Henderson, J. M. (2020). Visual and verbal working memory loads  
724 interfere with scene-viewing. *Attention, Perception, & Psychophysics*, 82(6), 2814–2820.  
725 <https://doi.org/10.3758/s13414-020-02076-1>
- 726 Davis, C., & Kim, J. (2004). Audio–Visual Interactions with Intact Clearly Audible Speech. *The*  
727 *Quarterly Journal of Experimental Psychology Section A*, 57(6), 1103–1121.  
728 <https://doi.org/10.1080/02724980343000701>
- 729 Davis, M. H., Johnsrude, I. S., Hervais-Adelman, A., Taylor, K., & McGettigan, C. (2005). Lexical  
730 Information Drives Perceptual Learning of Distorted Speech: Evidence From the  
731 Comprehension of Noise-Vocoded Sentences. *Journal of Experimental Psychology: General*,  
732 134(2), 222–241. <https://doi.org/10.1037/0096-3445.134.2.222>
- 733 Dorman, M. F., Loizou, P. C., & Rainey, D. (1997). Speech intelligibility as a function of the number of  
734 channels of stimulation for signal processors using sine-wave and noise-band outputs. *The*  
735 *Journal of the Acoustical Society of America*, 102(4), 2403–2411.  
736 <https://doi.org/10.1121/1.419603>
- 737 Dupoux, E., & Green, K. (1997). Perceptual adjustment to highly compressed speech: Effects of talker  
738 and rate changes. *Journal of Experimental Psychology. Human Perception and Performance*,  
739 23(3), 914–927. <https://doi.org/10.1037/0096-1523.23.3.914>
- 740 Erber, N. P. (1975). Auditory-visual perception of speech. *Journal of Speech and Hearing Disorders*,  
741 40(4), 481–492. <https://doi.org/10.1044/jshd.4004.481>
- 742 Everdell, I. T., Marsh, H., Yurick, M. D., Munhall, K. G., & Paré, M. (2007). Gaze behaviour in  
743 audiovisual speech perception: Asymmetrical distribution of face-directed fixations.  
744 *Perception*, 36(10), 1535–1545. <https://doi.org/10.1068/p5852>
- 745 Faulkner, A., Rosen, S., & Smith, C. (2000). Effects of the salience of pitch and periodicity information  
746 on the intelligibility of four-channel vocoded speech: Implications for cochlear implants. *The*  
747 *Journal of the Acoustical Society of America*, 108(4), 1877–1887.  
748 <https://doi.org/10.1121/1.1310667>

- 749 Golomb, J. D., Peelle, J. E., & Wingfield, A. (2007). Effects of stimulus variability and adult aging on  
750 adaptation to time-compressed speech. *The Journal of the Acoustical Society of America*,  
751 *121*(3), 1701–1708. <https://doi.org/10.1121/1.2436635>
- 752 Greenwood, D. D. (1990). A cochlear frequency-position function for several species—29 years later.  
753 *The Journal of the Acoustical Society of America*, *87*(6), 2592–2605.  
754 <https://doi.org/10.1121/1.399052>
- 755 Gurler, D., Doyle, N., Walker, E., Magnotti, J., & Beauchamp, M. (2015). A link between individual  
756 differences in multisensory speech perception and eye movements. *Attention, Perception, &*  
757 *Psychophysics*, *77*(4), 1333–1341. <https://doi.org/10.3758/s13414-014-0821-1>
- 758 Hazan, V., Kim, J., & Chen, Y. (2010). Audiovisual perception in adverse conditions: Language,  
759 speaker and listener effects. *Speech Communication*, *52*(11–12), 996–1009.  
760 <https://doi.org/10.1016/j.specom.2010.05.003>
- 761 Henderson, J. M., Malcolm, G. L., & Schandl, C. (2009). Searching in the dark: Cognitive relevance  
762 drives attention in real-world scenes. *Psychonomic Bulletin & Review*, *16*(5), 850–856.  
763 <https://doi.org/10.3758/PBR.16.5.850>
- 764 Hervais-Adelman, A., Davis, M. H., Johnsrude, I. S., & Carlyon, R. P. (2008). Perceptual learning of  
765 noise vocoded words: Effects of feedback and lexicality. *Journal of Experimental Psychology:*  
766 *Human Perception and Performance*, *34*(2), 460–474. [https://doi.org/10.1037/0096-](https://doi.org/10.1037/0096-1523.34.2.460)  
767 [1523.34.2.460](https://doi.org/10.1037/0096-1523.34.2.460)
- 768 Jerger Susan, Damian Markus F., Karl Cassandra, & Abdi Hervé. (2018). Developmental Shifts in  
769 Detection and Attention for Auditory, Visual, and Audiovisual Speech. *Journal of Speech,*  
770 *Language, and Hearing Research*, *61*(12), 3095–3112. [https://doi.org/10.1044/2018\\_JSLHR-](https://doi.org/10.1044/2018_JSLHR-H-17-0343)  
771 [H-17-0343](https://doi.org/10.1044/2018_JSLHR-H-17-0343)
- 772 Kawase, T., Sakamoto, S., Hori, Y., Maki, A., Suzuki, Y., & Kobayashi, T. (2009). Bimodal audio-visual  
773 training enhances auditory adaptation process. *Neuroreport*, *20*(14), 1231–1234.  
774 <https://doi.org/10.1097/WNR.0b013e32832fbef8>



## EYE GAZE AND ADAPTATION

- 775 Kuznetsova, A., Brockhoff, B., & Christensen, R. H. B. (2017). lmerTest Package: Tests in Linear Mixed  
776 Effects Models. *Journal of Statistical Software*, 82(13). <https://doi.org/10.18637/jss.v082.i13>
- 777 Lansing, C. R., & McConkie, G. W. (2003). Word identification and eye fixation locations in visual and  
778 visual-plus-auditory presentations of spoken sentences. *Perception & Psychophysics*, 65(4),  
779 536–552. <https://doi.org/10.3758/BF03194581>
- 780 Lewkowicz, D. J., & Hansen-Tift, A. M. (2012). Infants deploy selective attention to the mouth of a  
781 talking face when learning speech. *Proceedings of the National Academy of Sciences of the*  
782 *United States of America*, 109(5), 1431–1436. <https://doi.org/10.1073/pnas.1114783109>
- 783 Loizou, P. C., Dorman, M., & Tu, Z. (1999). On the number of channels needed to understand speech.  
784 *The Journal of the Acoustical Society of America*, 106(4), 2097–2103.  
785 <https://doi.org/10.1121/1.427954>
- 786 Lusk, L. G., & Mitchel, A. D. (2016). Differential Gaze Patterns on Eyes and Mouth During Audiovisual  
787 Speech Segmentation. *Frontiers in Psychology*, 7, 52.  
788 <https://doi.org/10.3389/fpsyg.2016.00052>
- 789 MacLeod, a, & Summerfield, Q. (1987). Quantifying the contribution of vision to speech perception  
790 in noise. *British Journal of Audiology*, 21(October), 131–141.  
791 <https://doi.org/10.3109/03005368709077786>
- 792 Malcolm, G. L., Lanyon, L. J., Fugard, A. J. B., & Barton, J. J. S. (2008). Scan patterns during the  
793 processing of facial expression versus identity: An exploration of task-driven and stimulus-  
794 driven effects. *Journal of Vision*, 8(8), 2–2. <https://doi.org/10.1167/8.8.2>
- 795 Matin, E., & Ethel. (1974). Saccadic suppression: A review and an analysis. *Psychological Bulletin*,  
796 81(12), 899–917. <https://doi.org/10.1037/h0037368>
- 797 Mattys, S. L., Davis, M. H., Bradlow, A. R., & Scott, S. K. (2012). Speech recognition in adverse  
798 conditions: A review. *Language and Cognitive Processes*, 27(7–8), 953–978.  
799 <https://doi.org/10.1080/01690965.2012.705006>

- 800 McGettigan, C., Rosen, S., & Scott, S. K. (2014). Lexico-semantic and acoustic-phonetic processes in  
801 the perception of noise-vocoded speech: Implications for cochlear implantation. *Frontiers in*  
802 *Systems Neuroscience*, 8. <https://doi.org/10.3389/fnsys.2014.00018>
- 803 Mitchel, A. D., & Weiss, D. J. (2014). Visual speech segmentation: Using facial cues to locate word  
804 boundaries in continuous speech. *Language, Cognition and Neuroscience*, 29(7), 771–780.  
805 <https://doi.org/10.1080/01690965.2013.791703>
- 806 Morin-Lessard, E., Poulin-Dubois, D., Segalowitz, N., & Byers-Heinlein, K. (2019). Selective attention  
807 to the mouth of talking faces in monolinguals and bilinguals aged 5 months to 5 years.  
808 *Developmental Psychology*, 55(8), 1640–1655. <https://doi.org/10.1037/dev0000750>
- 809 Paulus, M., Hazan, V., & Adank, P. (2020). The relationship between talker acoustics, intelligibility,  
810 and effort in degraded listening conditions. *The Journal of the Acoustical Society of America*,  
811 147(5), 3348–3359. <https://doi.org/10.1121/10.0001212>
- 812 Peelle, J. E., & Wingfield, A. (2005). *Dissociations in Perceptual Learning Revealed by Adult Age*  
813 *Differences in Adaptation to Time-Compressed Speech*. [https://doi.org/10.1037/0096-](https://doi.org/10.1037/0096-1523.31.6.1315)  
814 1523.31.6.1315
- 815 Pilling, M., & Thomas, S. (2011). Audiovisual Cues and Perceptual Learning of Spectrally Distorted  
816 Speech. *Language and Speech*, 54(4), 487–497. <https://doi.org/10.1177/0023830911404958>
- 817 Rosen, S., Faulkner, A., & Wilkinson, L. (1999). Adaptation by normal listeners to upward spectral  
818 shifts of speech: Implications for cochlear implants. *The Journal of the Acoustical Society of*  
819 *America*, 106(6), 3629–3636. <https://doi.org/10.1121/1.428215>
- 820 Rothausser, E. H., Chapman, N. D., Guttman, N., Nordby, K. S., Silbiger, H. R., Urbanek, G. E., &  
821 Weinstock, M. (1969). IEEE recommended practice for speech quality measurements. *IEEE*  
822 *Transactions on Audio and Electroacoustics*, 17(3), 225–246.  
823 <https://doi.org/10.1109/TAU.1969.1162058>

- 824 Schoof, T., & Rosen, S. (2015). High sentence predictability increases the fluctuating masker benefit.  
825 *The Journal of the Acoustical Society of America*, *138*(3), EL181–EL186.  
826 <https://doi.org/10.1121/1.4929627>
- 827 Scott, S. K., Rosen, S., Lang, H., & Wise, R. J. S. (2006). Neural correlates of intelligibility in speech  
828 investigated with noise vocoded speech—A positron emission tomography study. *The*  
829 *Journal of the Acoustical Society of America*, *120*(2), 1075–1083.  
830 <https://doi.org/10.1121/1.2216725>
- 831 Sebastian-Galles, N., & Mehler, J. (2000). Adaptation to time-compressed speech: Phonological  
832 determinants. In *Perception & Psychophysics* (Vol. 62, Issue 4).
- 833 Shannon, R. V., Zeng, F.-G. G., Kamath, V., Wygonski, J., & Ekelid, M. (1995). Speech recognition with  
834 primarily temporal cues. *Science (New York, N.Y.)*, *270*(5234), 303–304.  
835 <https://doi.org/10.1126/science.270.5234.303>
- 836 Sommers, M. S., Tye-Murray, N., & Spehar, B. (2005). Auditory-visual speech perception and  
837 auditory-visual enhancement in normal-hearing younger and older adults. *Ear and Hearing*,  
838 *26*(3), 263–275. <https://doi.org/10.1097/00003446-200506000-00003>
- 839 Stacey, J. E., Howard, C. J., Mitra, S., & Stacey, P. C. (2020). Audio-visual integration in noise:  
840 Influence of auditory and visual stimulus degradation on eye movements and perception of  
841 the McGurk effect. *Attention, Perception, & Psychophysics*, *82*(7), 3544–3557.  
842 <https://doi.org/10.3758/s13414-020-02042-x>
- 843 Sumbly, W. H., & Pollack, I. (1954). Visual Contribution to Speech Intelligibility in Noise. *The Journal of*  
844 *the Acoustical Society of America*, *26*(2), 212–215. <https://doi.org/10.1121/1.1907309>
- 845 Summerfield, Q. A. (1987). Some preliminaries to a comprehensive account of audio-visual speech  
846 perception. In *Hearing by Eye: The Psychology of Lip-reading* (pp. 3–51).  
847 <https://doi.org/citeulike-article-id:795404>

- 848 Thomas, S. M., & Jordan, T. R. (2004). Contributions of Oral and Extraoral Facial Movement to Visual  
849 and Audiovisual Speech Perception. *Journal of Experimental Psychology: Human Perception*  
850 *and Performance*, 30(5), 873–888. <https://doi.org/10.1037/0096-1523.30.5.873>
- 851 Trotter, A., Banks, B., & Adank, P. (2020). *Effects of the availability of visual cues during adaptation*  
852 *to noise—Vocoded speech*. PsyArXiv. <https://doi.org/10.31234/osf.io/jxaeb>
- 853 Vatikiotis-Bateson, E., Eigsti, I. M., Yano, S., & Munhall, K. G. (1998). Eye movement of perceivers  
854 during audiovisual speech perception. *Perception & Psychophysics*, 60(6), 926–940.  
855 <https://doi.org/10.3758/BF03211929>
- 856 Vo, M. L.-H., Smith, T. J., Mital, P. K., & Henderson, J. M. (2012). Do the eyes really have it? Dynamic  
857 allocation of attention when viewing moving faces. *Journal of Vision*, 12(13), 3–3.  
858 <https://doi.org/10.1167/12.13.3>
- 859 Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic*  
860 *Bulletin & Review*, 14(5), 779–804. <https://doi.org/10.3758/BF03194105>
- 861 Wang Jianrong, Zhu Yumeng, Chen Yu, Mamat Abdilbar, Yu Mei, Zhang Ju, & Dang Jianwu. (2020). An  
862 Eye-Tracking Study on Audiovisual Speech Perception Strategies Adopted by Normal-Hearing  
863 and Deaf Adults Under Different Language Familiarities. *Journal of Speech, Language, and*  
864 *Hearing Research*, 63(7), 2245–2254. [https://doi.org/10.1044/2020\\_JSLHR-19-00223](https://doi.org/10.1044/2020_JSLHR-19-00223)
- 865 Wayne, R. V., & Johnsrude, I. S. (2012). The role of visual speech information in supporting  
866 perceptual learning of degraded speech. *Journal of Experimental Psychology: Applied*, 18(4),  
867 419–435. <https://doi.org/10.1037/a0031042>
- 868 Worster, E., Pimperton, H., Ralph-Lewis, A., Monroy, L., Hulme, C., & MacSweeney, M. (2018). Eye  
869 movements during visual speech perception in deaf and hearing children. *Language*  
870 *Learning*, 68(Suppl 1), 159–179. <https://doi.org/10.1111/lang.12264>
- 871 Yarbus, A. L. (1967). *Eye Movements and Vision*. Springer US. <https://doi.org/10.1007/978-1-4899->  
872 5379-7
- 873



875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891

**Figure Captions**

**Figure 1**

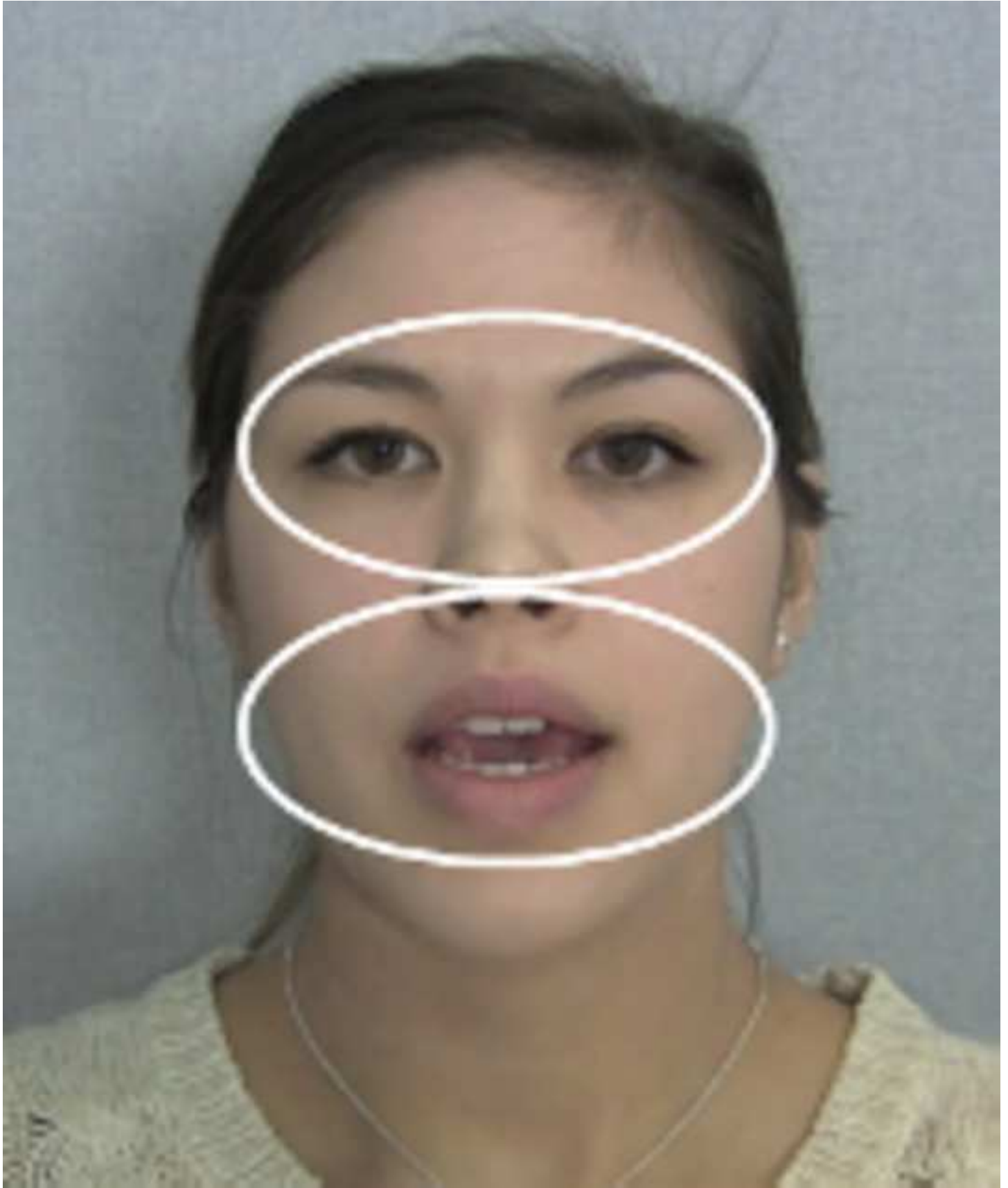
Image of the speaker with regions of interest ('mouth' and 'eyes').

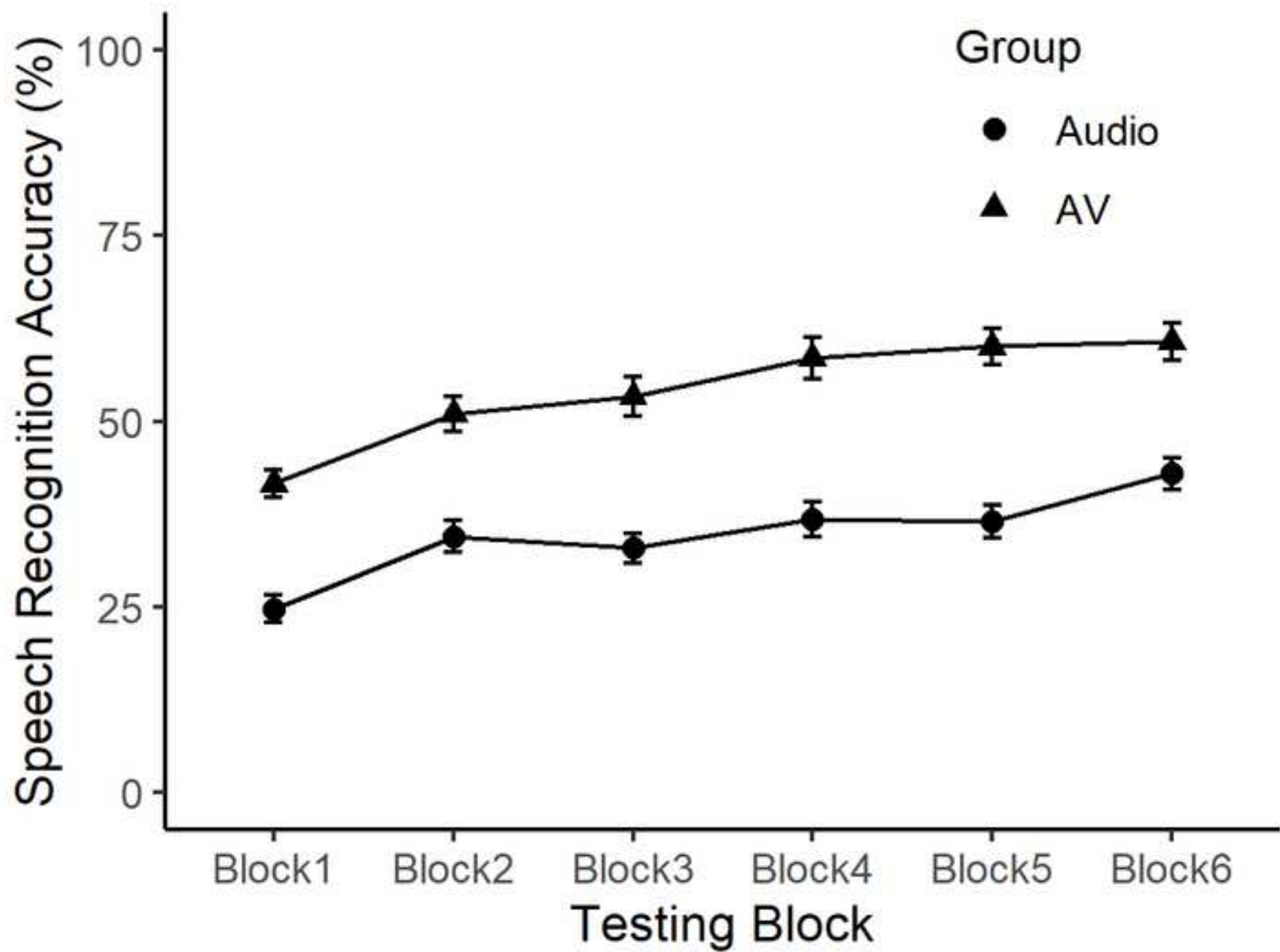
**Figure 2.**

Mean recognition accuracy per testing block, per group. Error bars show  $\pm 1SE$ .

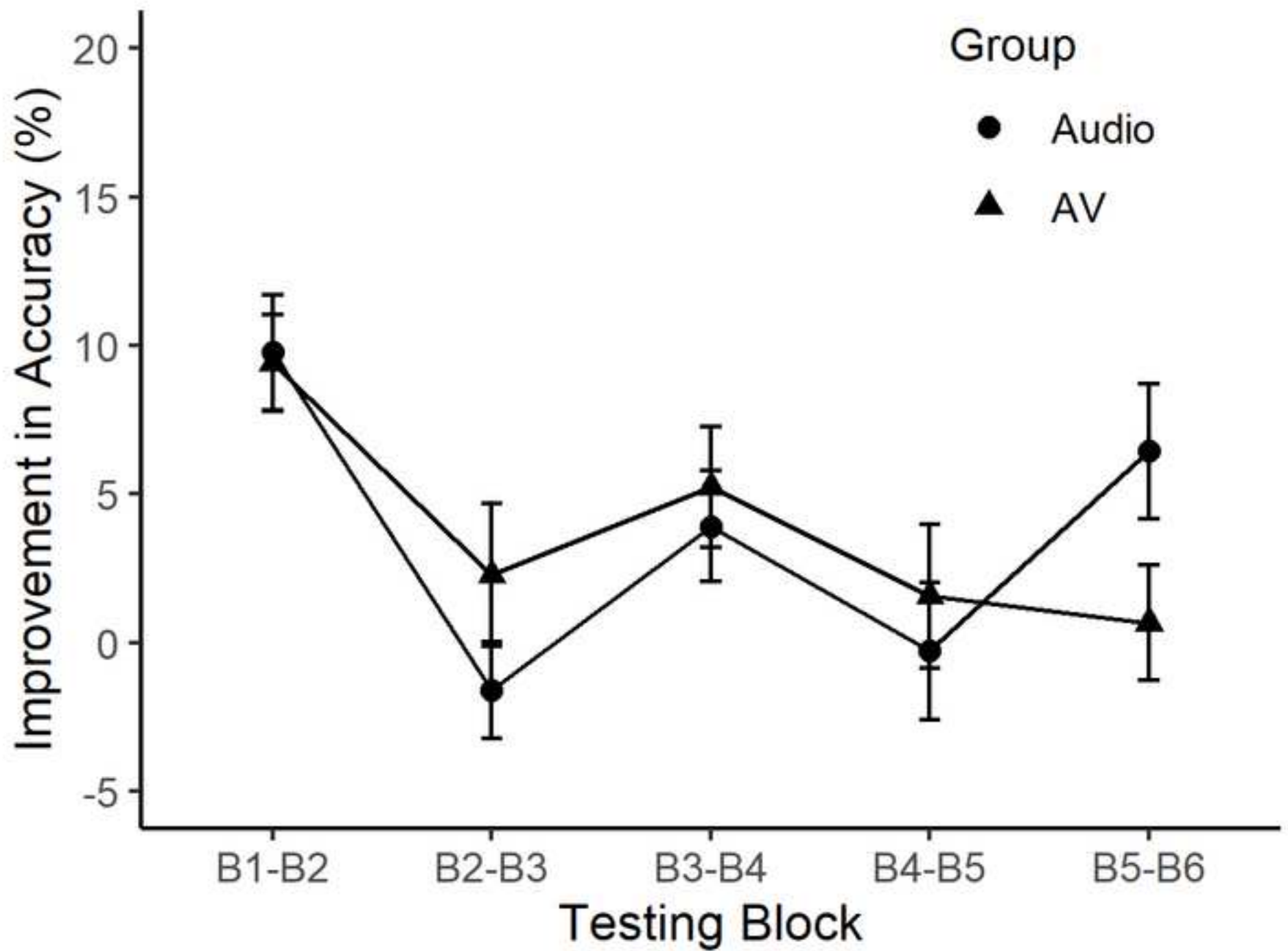
**Figure 3.** Adaptation (amount of improvement) between consecutive testing blocks per group. Error bars show  $\pm 1SE$ .

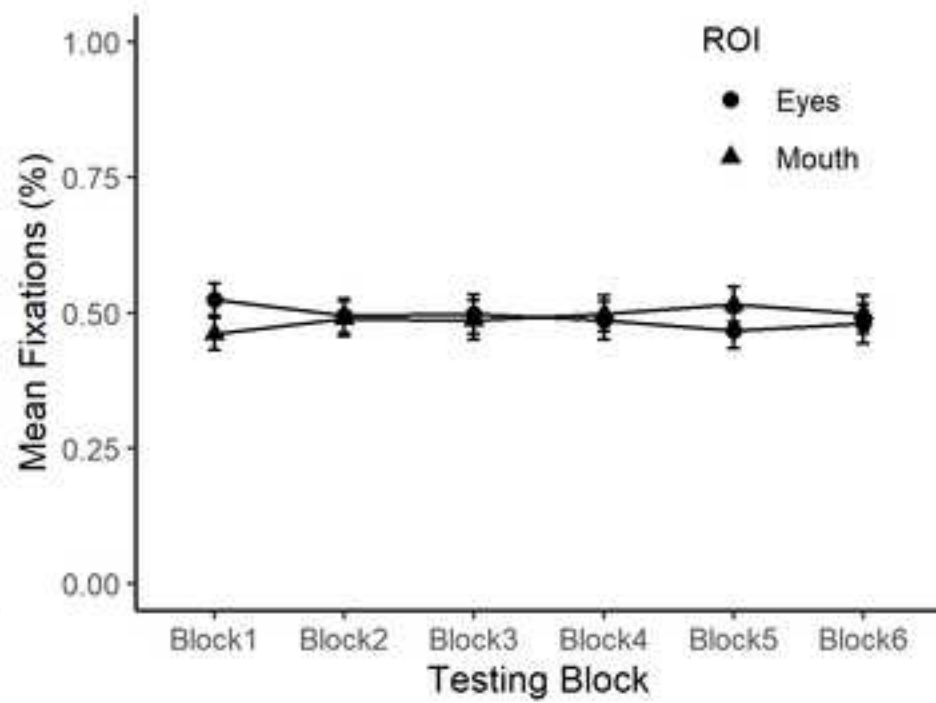
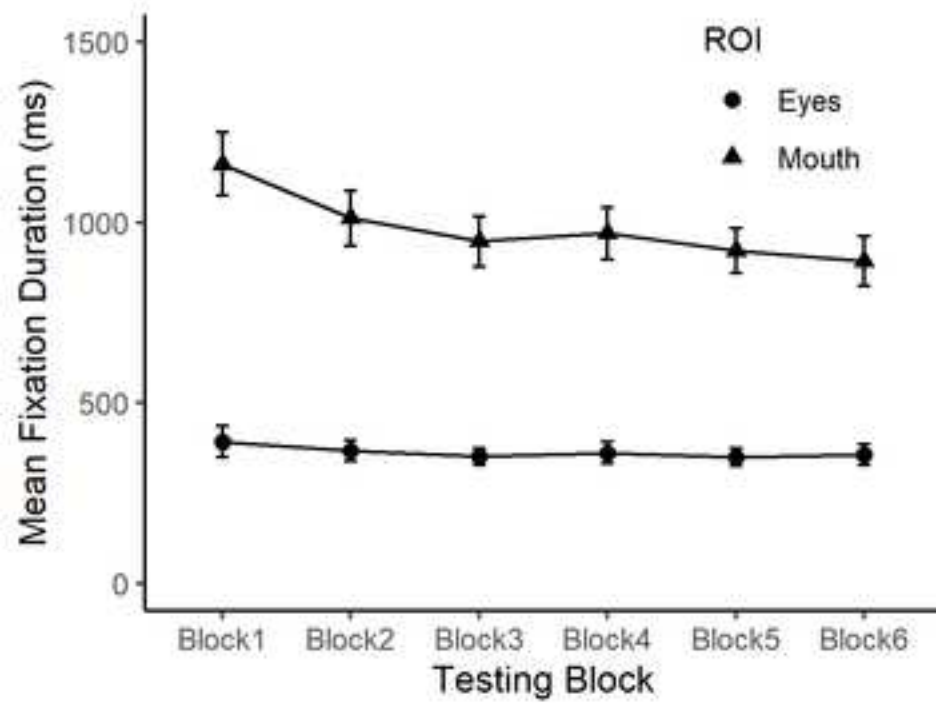
**Figure 4.** Duration of fixations (left panel) and percentage fixations (right panel) on the mouth and eyes, per testing block in the audiovisual group. Error bars  $\pm 1SE$ .











**Table 1.** Exploratory Bayesian hierarchical regression analyses of group differences in adaptation rate.

Model	BF <sub>10</sub>	$R^2m$	$\Delta R^2$	Interpretation
Baseline model	-	.458	-	-
B1-B2*Group	0.12	.459	.001	No group difference
B2-B3*Group	0.59	.462	.003	Inconclusive
B3-B4*Group	0.23	.460	.001	No group difference
B4-B5*Group	0.08	.459	.000	No group difference
B5-B6*Group	0.08	.459	.000	No group difference

*Note.*  $R^2m$  = marginal  $R^2$  (fixed effects only);  $\Delta R^2$  indicates change in marginal  $R^2$  based on difference between baseline model and the addition of the model interaction. BF<sub>10</sub> = Bayes Factor indicating evidence of a difference between groups in the amount of adaptation between each consecutive pair of testing blocks.

**Table 2.** Hierarchical mixed model comparisons for the audiovisual group predicting overall speech recognition by each measure of eye gaze.

Model	R <sup>2</sup>	<i>p</i> -value	BF <sub>10</sub>
Testing Block (baseline model of adaptation)	0.20	<.001**	5.22229e+12
Testing Block + Duration of Fixations on Mouth	0.25	.019*	1.15
Testing Block + Duration of Fixations on Eyes	0.20	.805	0.08
Testing Block + Percentage Fixations on Mouth	0.20	.496	0.09
Testing Block + Percentage Fixations on Eyes	0.20	.613	0.08
Testing Block * Duration of Fixations on Mouth (interaction)	0.20	1.00	0.07

*Note:* All models contain the random effect of participant. We report marginal  $R^2$  representing the variance explained by fixed effects only.

\*  $p < .05$ ; \*\*  $p < .001$