

1

2

3

4

5 Reassessing the Benefits of Audio-Visual Integration to Speech Perception and Intelligibility

6

7 Brandon O'Hanlon<sup>1</sup>, Professor Christopher J. Plack<sup>1,2</sup>, Dr Helen E Nuttall<sup>1</sup>

8

9 1. Department of Psychology, Lancaster University, UK

10 2. Manchester Centre for Audiology & Deafness, The University of Manchester, UK

11

12 Corresponding Author: Dr Helen E Nuttall ([h.nuttall1@lancaster.ac.uk](mailto:h.nuttall1@lancaster.ac.uk))

13

14 Conflict of Interest: The authors declare no conflicts of interest.

15

16

## 17 Abstract

18 **Purpose:** In difficult listening conditions, the visual system assists with speech perception  
19 through lipreading. Stimulus onset asynchrony (SOA) is used to investigate the interaction  
20 between the two modalities in speech perception. Previous estimates of audiovisual benefit  
21 and SOA integration period differ widely. A limitation of previous research is a lack of  
22 consideration of visemes - categories of phonemes defined by similar lip movements when  
23 produced by a speaker - to ensure that selected phonemes are visually distinct. This study  
24 aimed to reassess the benefits of audiovisual lipreading to speech perception when different  
25 viseme categories are selected as stimuli and presented in noise. The study also aimed to  
26 investigate the effects of SOA on these stimuli.

27 **Method:** Sixty participants were tested online and presented with audio-only and audiovisual  
28 stimuli containing the speaker's lip movements. The speech was presented either with or  
29 without noise and had six different SOAs (0, 200, 216.6, 233.3, 250, and 266.6 ms).  
30 Participants discriminated between speech syllables with button presses.

31 **Results:** The benefit of visual information was weaker than that in previous studies. There  
32 was a significant increase in reaction times as SOA was introduced, but no significant effects  
33 of SOA on accuracy. Furthermore, exploratory analyses suggest that the effect was not equal  
34 across viseme categories: 'Ba' was more difficult to recognise than 'Ka' in noise.

35 **Conclusion:** In summary, the findings suggest that the contributions of audiovisual  
36 integration to speech processing are weaker when considering visemes but are not sufficient  
37 to identify a full integration period.

38 Keywords: audiovisual speech, speech perception, multisensory integration, visemes, vision

39

40

## 41 Reassessing the Benefits of Audio-Visual Integration to Speech Perception and Intelligibility

42 Intelligible speech is built up from speech phonemes. Phonemes are small linguistic  
43 units - such as the /b/ phoneme that begins ‘boy’ in the English language - and play a large  
44 role in the identification of speech (Ewen & Van der Hulst, 2001; Bowers *et al.*, 2016).  
45 Processing of speech can be made more difficult with the introduction of noise in the  
46 environment, which reduces the ability to discriminate successfully between phonemes  
47 (Summerfield, 1992). In many cases, information from the visual sense that is relevant to the  
48 speech – such as from lipreading – can be integrated into speech processing systems to  
49 improve comprehension. In background noise, this assisting sense is recruited further (Yuan  
50 *et al.*, 2021). Viewing the lip movements when an individual is speaking can help to improve  
51 the intelligibility of speech-in-noise versus when the lips are not visible (Sumbly & Pollack,  
52 1954; Maier *et al.*, 2011). The inverse can also occur, wherein incongruent lip movements  
53 influence our ability to discriminate between speech sounds. An example of this is the  
54 McGurk Effect (McGurk & MacDonald, 1976), where presenting the speech phoneme ‘Ba’  
55 with visual lip movements associated with ‘Ga’ leads to perceptions of the sound ‘Da’  
56 instead. A more recent example comes from face mask-wearing due to the COVID-19  
57 pandemic. Brown *et al.* (2021) found that if the speaker wore a facemask that either fully or  
58 partially covered lip movements, performance on speech discrimination tasks decreased  
59 dramatically. These data indicate that the visual and auditory systems interact to influence  
60 how we perceive speech.

61 However, estimates of audiovisual benefit vary widely in the literature, likely due to  
62 stimulus-dependent effects (Ma *et al.*, 2009), in that, how the stimuli are created for lab  
63 experimentation drastically affects how participants respond to speech discrimination tasks.  
64 For example, whilst it is important for research on audio-visual processing to consider how  
65 auditorily distinctive sounds are, visual distinctiveness is equally important. A way to

66 examine the effect of visual distinctiveness is to select phonemes from separate viseme  
67 categories for testing. A viseme category is a group of phonemes from the English language  
68 that share the lip movements and visual information portrayed by each phoneme when  
69 spoken (Massaro *et al.*, 2012). Fisher (1968) identified five viseme categories based purely  
70 on visual distinguishability for English phonemes. Examples of phonemes that belong to the  
71 same viseme category are /b/, /p/, and /m/ which in syllable form can correspond to 'Ba',  
72 'Pa', and 'Ma'. If two speech tokens share the same viseme, then it is impossible to discern  
73 which was spoken through lip-reading alone (Van Engen *et al.*, 2022) and are only  
74 distinguishable through sound. This means that any measure of audiovisual benefit derived  
75 from discriminating within a viseme category will be lessened. It is therefore important to  
76 select stimuli from separate viseme categories when investigating how auditory and visual  
77 systems work together during speech syllable discrimination.

78         When speaking with others, we typically see lip movements before we hear the  
79 spoken words (Chandrasekaran *et al.*, 2009) as a form of natural stimulus onset asynchrony  
80 (SOA). SOA is when two different modalities of information in cross-modal stimuli are  
81 presented at different onsets. The window of integration is the term given to the period in  
82 which visual information can lead or lag speech sounds before the visual information is no  
83 longer perceived as part of the same stimulus (Stein & Meredith, 1993). If the lip movements  
84 are desynchronised from the speech sounds within a specific period, then we still perceive the  
85 lip movements and the speech we hear to be congruent. If the SOA is large enough that the  
86 auditory and visual information do not fall within the same window of integration, we may  
87 perceive the two modalities as separate, and therefore not process the visual information as  
88 helpful extra information to discern and comprehend the speech. For speech signals, syllables  
89 have a window with an upper limit of about 240 ms and short words of about 300 ms  
90 (Navarra *et al.*, 2005). Although it is important to note that this window of integration can be

91 highly stimulus-dependent, and ranges in the literature between 150 and 800 ms (Colonius &  
92 Diederich, 2010; Schwartz & Savariaux, 2014, Ren *et al.*, 2017), and even differs between  
93 age groups (Ren *et al.*, 2017). This mixed range in the literature could also be due to a  
94 mismatch with reported display refresh rates (typically 60 Hz), video framerates (typically 30  
95 – 60 frames per second) and levels of SOA used in research if reported at all (Ren *et al.*,  
96 2017). For example, in 10 ms increments, a 60 Hz monitor can't display separate visual  
97 streams of information that refresh every 10 ms, as it is only capable of doing so every 16.6  
98 ms, assuming the video plays at a full 60 frames per second as well.

99         The present study aimed to reassess the benefits of visual information to speech-in-  
100 noise perception using stimuli with visual distinctiveness. We also aimed to determine the  
101 effect of SOA on audiovisual speech perception. We tested the following hypotheses:

- 102         (i) purely audio speech discrimination accuracy will be decreased when speech is  
103 presented in noise compared to without noise.
- 104         (ii) reaction time to correctly discriminated purely audio speech will be increased  
105 when speech is presented in noise compared to without noise.
- 106         (iii) speech-in-noise discrimination accuracy will be increased when speech is  
107 presented with congruent visual information of the speaker's lip movements  
108 (audiovisual stimuli) compared to when no visual information is present (purely  
109 audio stimuli).
- 110         (iv) reaction time to correctly discriminated speech-in-noise will be decreased when  
111 speech is presented with congruent visual information of the speaker's lip  
112 movements (audiovisual stimuli) compared to when no visual information is  
113 present (purely audio stimuli).
- 114         (v) as the visual information precedes the auditory information by larger SOAs (0 -  
115 266 ms), speech-in-noise discrimination accuracy will decrease.

116 (vi) as the visual information precedes the auditory information by larger SOAs,  
117 reaction time to correctly discriminated speech-in-noise will increase.

118 Further exploratory analysis also investigated the window of integration for these  
119 audiovisual stimuli, as well as differences in visual benefit between each syllable used.

120

121

**Method****Design**

123 To address hypotheses (i) and (ii), a single within factor (noise type: speech without  
124 noise, and speech-in-noise) design was used for purely audio trials with no stimulus  
125 asynchrony. For hypotheses (iii) and (iv), a single within factor (stimulus type: audiovisual,  
126 and purely audio) design was used for speech-in-noise trials with no stimulus asynchrony.  
127 Finally, for hypotheses (v) and (vi), a single within factor (SOA; 0, 200, 216.6, 233.3, 250,  
128 and 266.6 ms) design was used for audiovisual, speech-in-noise trials only. In total,  
129 participants took part in all 14 unique conditions (see Table 1), and both the accuracy of  
130 speech discrimination and reaction time in the discrimination task were recorded.

131 Ethical approval was granted by the Faculty of Science and Technology Research  
132 Ethics Committee at Lancaster University (approval reference: FST-2022-2122-RECR-2,  
133 project ID: 2122). The study was pre-registered on AsPredicted.org before commencing data  
134 collection. The pre-registration can be found at <https://aspredicted.org/aq98a.pdf>. All  
135 deviations from this pre-registration are listed in the section below. The collected data have  
136 been archived on the Open Science Framework (OSF: <https://osf.io/kcbzs>).

**Deviations from pre-registration**

138 In the original study pre-registration, there were three set hypotheses listed:

- 139 • There will be a decrease in the accuracy of speech discrimination (measured  
140 by correct responses in trials) or an increase in response times in the auditory-  
141 only condition when the speech is in noise compared to speech without noise.
- 142 • When visual information is present (audiovisual), the accuracy of speech  
143 discrimination and response times for each trial will not be as obstructed in

144 speech-in-noise conditions compared to audio-only conditions (when no visual  
145 information is present).

- 146 • As the visual information precedes the auditory information by larger margins  
147 (200 ms, 216 ms, 233 ms, 250 ms, 266 ms), the accuracy of speech  
148 discrimination in the speech-in-noise conditions will decrease - or response  
149 times will increase - in audiovisual conditions compared to when the  
150 audiovisual information is congruent (0ms).

151 These were changed to the hypotheses listed in the introduction by splitting the  
152 dependent measures into separate hypotheses and improving readability. This was done to  
153 make interpretations of results more clearly defined when referring to the hypotheses. To  
154 accompany this, the models used to test the hypothesis were also adjusted, giving six separate  
155 models of analysis - one for each hypothesis - instead of four. Generalised linear mixed-  
156 effects models (GLMER) were used to test all six hypotheses, instead of the mixture of  
157 GLMER models for accuracy data and LMER models for reaction time data that was listed in  
158 the pre-registration. This was done as GLMER models are more appropriate than LMER  
159 models for reaction time data, which is generally positively skewed (Lo & Andrews, 2015).  
160 These GLMER models were preferable still over repeated measures generalised linear  
161 models for considering random effects that may be present on a participant-by-participant  
162 basis. Finally, in our sample size calculation using data simulation (see section ‘Sample size  
163 calculation’), it was determined that 60 participants were needed to sufficiently power the  
164 study. In the pre-registration, we then added a further 10% after a priori calculations (another  
165 6 participants) to make a sample size estimate of 66. Due to the availability of resources, this  
166 extra 10% was not collected, leaving the sample of the study at the original number of 60  
167 participants.

## 168 **Participants**



169 All data were collected online, with 81 participants recruited for the study. Of these, a  
170 total of 60 participants completed the study (mean age = 25.66, 28 male, 30 female, two non-  
171 binary). The other 21 participants completed the eligibility questionnaire but were either not  
172 eligible or did not proceed to the study task and provide study data. Participants were  
173 recruited via online advertisements or through Prolific and were compensated for their time.  
174 All participants were monolingual, native speakers of British English to control any potential  
175 speech perception differences across languages and in bilingualism and multilingualism  
176 (Lotfi *et al.*, 2019). Participants reported no hearing disorders and had either normal or  
177 corrected-to-normal vision. Only those between the ages of 18 and 35 were tested, as the  
178 window of integration for audiovisual information increases significantly with age, which can  
179 make speech discrimination more difficult (Ganesh *et al.*, 2018; Sekiyama *et al.*, 2014).  
180 Participants reported no developmental disorders, such as dyslexia, or history of  
181 developmental disorders. This was important as the window of integration for audiovisual  
182 stimuli is wider in individuals with learning difficulties such as autism spectrum disorder and  
183 developmental dyslexia (Smith & Bennetto, 2007; Megnin-Viggars & Goswami, 2013;  
184 Michalek *et al.*, 2014; Noel *et al.*, 2018). All participants were right-handed. Finally,  
185 participants had no musical expertise, as previous research suggests that individuals with  
186 continuous experience as musicians can detect smaller SOAs, even for speech syllables (Lee  
187 & Noppeney, 2014; Sorati & Behne, 2019). Musical expertise was defined as training with a  
188 single musical instrument or voice for more than 7 years (Varnet *et al.*, 2015; Lee *et al.*,  
189 2020) and for at least 3.5 hours a week (Lee & Noppeney, 2014). Participants were screened  
190 for the experiment using Qualtrics (see section ‘Procedure’).

### 191 **Sample size calculation**

192 Before testing, data simulation was conducted using R studio for power and sample  
193 size analysis. Lme4 (vers. 1.1-27.1; Bates *et al.*, 2015), afex (vers. 1.0-1; Singmann *et al.*,

2024) and *simr* (vers. 1.0.5; Peter *et al.*, 2019) were the core packages utilised in this process. Firstly, means and standard deviations of accuracy were gathered from studies that used syllable or bi-syllable phonetic speech tokens to investigate visual integration in speech perception. These studies typically used either multiple signal-to-noise ratios (SNRs; between -12 and -18 dB: Altieri *et al.*, 2014; Grant & Seitz, 1998; Sekiyama *et al.*, 2014) or individualised ratios (Ten Oever *et al.* 2013). For those studies that used multiple speech-to-noise ratios, we took data from – or closest to – -16 dB SNR. -16 dB was selected for our speech-shaped noise as this was the average SNR at which there was a notable difference between perceiving speech with or without visual aid (Bernstein *et al.*, 2004). An average estimated mean and standard deviation were then calculated for each condition. A dataset was produced using the *rtruncnorm* function (*truncnorm* package; vers. 1.0-8; Mersmann *et al.*, 2018) - to randomly generate data for each condition that had a mean and standard deviation close to the ones calculated. This was repeated for each speech token ('Ba', 'Fa', and 'Ka') and all trials of each condition, providing a full dataset of expected results.

The dataset was then analysed using our planned experimental analyses (see below) to generate predicted results. Simulations were repeated 1000 times. An aggregation of power was then calculated. If the power was insufficient (below .80 at an alpha level of .05), the sample size of the dataset was manually adjusted, and the data simulation was conducted again. This was done until a minimal sample size with sufficient power was found. A total of 60 participants were calculated to be needed for sufficient power. The code for data simulation is available on OSF (<https://osf.io/kcbzs>).

## 215 **Materials**

216 The experiment was created using PsychoPy 3's builder tools (vers. 2021.2.3; Peirce  
217 *et al.*, 2019) and hosted online through Pavlovia. A consent form and a screening form were

218 created and hosted on Qualtrics (Qualtrics, 2005). Three single-syllable speech tokens were  
219 used: ‘Ba’, ‘Fa’, and ‘Ka’. These were chosen as they belong to three distinct viseme  
220 categories, did not rely on any tongue movements to distinguish that would have been  
221 obscured from sight (such as labiodental phonemes), and could be easily distinguished  
222 without visual aid when not in noise. These speech tokens were spoken by a native British  
223 English-speaking male speaker and were recorded using personal home equipment. An  
224 external USB 3.0 condenser microphone was used to record audio (HyperX Quadcast with  
225 default windshield, set to the cardioid position). The initial video footage was recorded at  
226 1920 x 1080 resolution and 60 frames per second using a mobile device (OnePlus 7 Pro).  
227 Both devices were connected to a single desktop machine, which recorded the audio and  
228 video in tandem using open-source OBS Studio software (Open Broadcaster Software,  
229 version 29.1.3). After the initial recording, the speech tokens were edited in length and  
230 converted to mp4 files at a resolution of 1280 x 720 and a frame rate of 60 frames per second.  
231 As the study would be completed on participants’ laptops or desktop systems and using their  
232 internet connection, we could not ensure that all participants were using a device with a 1920  
233 x 1080 resolution screen. By reducing the resolution of files to 1280 x 720, all likely  
234 participant resolution sizes could be accommodated whilst ensuring that all participants  
235 viewed the files at the same resolution. Sixty frames per second was chosen as the frame rate  
236 as home device monitors and laptop screens are typically to a standard 60 Hz or higher. By  
237 using the lower boundary and not a higher frame rate, we can be sure that all SOAs  
238 implemented in the stimuli were visually relayed to the participant. For audiovisual  
239 conditions, the video footage contained only the speaker’s lower face in view, containing  
240 mouth and lips. This meant that participants were only provided with visual information  
241 regarding the lip movements made when speaking, and not any other visual information  
242 relevant to other actions the speaker may have made during recordings. For audio-only

243 conditions, the video of the lips was overlaid with a plain black PNG image file. This kept  
244 the audio-only stimuli in a consistent video format rather than exporting the file as an mp3.  
245 All video files were the same length of 2 s.

246 Audacity software (Audacity Team, 2021) was then used to rip the audio from the  
247 MKV files to be edited as WAV files in Praat software (Boersma & Weenink, 2021) for the  
248 creation of speech-shaped noise. First, a sentence using English words – ‘*His plan meant*  
249 *taking a big risk*’ - was recorded to provide a base for the speech-shaped noise. White noise  
250 was then produced using Praat’s white noise generator. The noise was brought down to an  
251 intensity tier, then an amplitude tier. This was then multiplied with the sentence above to  
252 create speech-shaped noise (Van Engen *et al.*, 2017). Praat was then used to combine the  
253 speech-shaped noise with the speech-in-noise conditions at a speech-to-noise ratio of –16 dB.  
254 This was done using a Praat script developed by McCloy (2021). Finally, Audacity was used  
255 again to ramp up the start and ramp down the ends of all audio files for every condition. The  
256 audio was then stitched back onto the MP4 files.

257 For the conditions where the stimuli were asynchronous, Lightworks was again used  
258 to desynchronise the onset of the audio ahead of the onset of the lip movements using exact  
259 frames of the video footage (12, 13, 14, 15, and 16 frames per second) which corresponded  
260 with the SOAs of the relevant conditions (audio starting after the visual lip information by  
261 200, 216.6, 233.3, 250, and 266.6 ms). The result was 42 stimuli in MP4 format, representing  
262 three speech tokens (‘Ba’, ‘Fa’, and ‘Ka’) for each of the 14 condition levels presented to the  
263 participant.

## 264 **Procedure**

265 Participants were linked to Qualtrics once they had consented to the study.  
266 Participants were also reminded at this stage to ensure that they were in a quiet room with no

267 background noise, as well as to load the experiment on either Microsoft Edge, Google  
268 Chrome, or Mozilla Firefox internet browsers on a laptop or desktop computer. They were  
269 explicitly told not to open the experiment on any other browser, such as Safari, nor a mobile  
270 or tablet device as these were incompatible. Participants were also instructed to use  
271 headphones for the experiment, rather than to play the stimuli through their device's  
272 speakers.

273 A volume check began, in which a constant pure tone played (440 Hz frequency), and  
274 participants were asked to adjust the volume of their device as necessary for a comfortable  
275 auditory experience and to ensure that the audio was playing correctly at a sufficient volume  
276 level. This tone would play for as long as the participant wished to alter the volume levels of  
277 their device. Once complete, the spacebar would be pressed, and the tone stopped.  
278 Participants were informed that a video would play either showing no visual information or  
279 visual information of lips moving. Meanwhile, speech would be played. Participants were  
280 told to listen carefully to the speech sound spoken, and after hearing the sound to press one of  
281 three buttons on their keyboards that corresponded with the three available speech tokens.  
282 They were instructed to respond to each trial as quickly as possible. They were reminded  
283 before and after each trial to press 'z' on their keyboard if they heard 'Ba', 'x' for 'Fa', or 'c'  
284 for 'Ka'. Participants were told to answer as quickly as possible. If they were unsure, they  
285 were told to make a guess.

286 Participants were given six practice trials before data were collected. This was using  
287 the speech without noise, 0 ms, and audiovisual condition stimuli, with two trials for each of  
288 the three speech tokens (Ba, Fa, and Ka). A white crosshair would be displayed on the screen  
289 for 1000 ms before the trial began to bring attention to the centre of the screen where the  
290 video trials would be displayed. Stimuli were shown for 2500 ms, then the response screen  
291 would display. On this screen, the participants were reminded of the buttons to press for each

292 of the three speech sounds. Only the three buttons could be pressed and pressing the buttons  
293 whilst the stimuli were still playing would not record a response or stop the trial. A total of  
294 546 trials (not including the practice trials) were completed. The order of the trials and  
295 conditions was completely random to avoid any potential order bias. After every 42 trials, a  
296 break screen would appear. This screen told the participant to take a short break before  
297 continuing with a press of the spacebar. If the participant did not wish to take a break, they  
298 were permitted to continue with a spacebar press immediately. There was a total of 12 breaks  
299 in the experiment, each with a short attention check question to ensure participants remained  
300 attentive to the experiment. Upon completing the study, participants could close the browser  
301 tab or window down and all data would remain recorded on the Pavlovia system.

## 302 **Analysis**

303 Descriptive statistics were first gathered from each condition for both the accuracy  
304 ratings and the reaction times. Reaction times were taken from the offset of the stimuli to the  
305 participant response. The average accuracy and reaction time of accurately responded trials  
306 for each condition and each participant was calculated, with reaction times winsorised over  
307 the 95<sup>th</sup> percentile only. This was done to replace any large, outlying reaction times to trials  
308 that may be due to a distraction at home during testing or the participant taking a short break  
309 before the break period. The assumptions of linear and generalised linear mixed-effects  
310 models were tested, including residual plots to check for linearity, quantile-quantile plots for  
311 normality, assessing the levels of multicollinearity between stimulus type, noise, and SOA  
312 using variance inflation factors, and ensuring the assumption of homoscedasticity was met.  
313 All the above tests were conducted on the dataset and all assumptions were met. As we were  
314 testing six separate hypotheses, the experiment-wise error rate was controlled using the  
315 Bonferroni-Holm method (Holm, 1979).

316           With further regards to stimulus variability, previous studies often employ analyses  
317 such as repeated measures analysis of variance (ANOVA) tests which do not consider  
318 random effects (Bates *et al.*, 2015). Including random effects is important for ensuring that  
319 any effects found in the model are not influenced by differences in participant ability or by  
320 the stimuli themselves, as some stimuli may be easier to recognise and comprehend in noise  
321 than others. To counter this issue, mixed-effects models can be used that consider the random  
322 effects, such as participant number and stimuli number, across intercepts and slopes within  
323 the model to provide a more valid interpretation of the integration between visual and  
324 auditory systems in speech perception.

325           Using the lme4 package (Bates *et al.*, 2015), generalized linear mixed-effects  
326 regression model (GLMER) analyses were conducted for the accuracy scores to test  
327 hypotheses (i), (iii), and (v) and for reaction time scores to test hypotheses (ii), (iv), and (vi).  
328 GLMERs were chosen instead of repeated measures generalised linear models such as  
329 ANOVA tests because they consider random effects that may be present across all 546 trials  
330 on a participant-by-participant basis. GLMER was chosen over LMER for analysis with  
331 reaction times as these scores are typically positively skewed. As noted by Lo and Andrews  
332 (2015), generalised linear mixed models are more appropriate for skewed datasets in this  
333 context. Furthermore, accuracy in a trial is a binary outcome variable that can either be  
334 correct (1) or incorrect (0). Therefore, GLMERs were used to ensure that assumptions of  
335 categorical dependent variables in mixed-effects models were met. GLMERs were conducted  
336 using the lme4 package still, as this package supported a generalised approach. Due to the  
337 generalised nature of the model and package restrictions, no suitable p-values were provided  
338 with the GLMER analyses. Instead, significance was interpreted using 99.2% confidence  
339 intervals (CIs), chosen to reflect our lowest criterion of significance in the Bonferroni-Holm  
340 approach being  $p < .008$  for six comparisons. If the resulting confidence intervals showed

341 insignificance, the next boundary of Bonferroni-Holm ( $p < .01$ ) was checked using 99%  
342 confidence intervals. This kept going until either significance was found or no significance  
343 was found at a significance level of  $p < .05$ . Once detected or classed as insignificant, the test  
344 was ranked with the other p-values in our analyses as the lowest boundary of significance and  
345 Bonferroni-Holm was conducted as normal on our six ranked comparisons.

346 To test hypothesis (i), a GLMER analysis was conducted using the accuracy of  
347 responses on the speech discrimination task as the dependent variable and using noise type  
348 (no noise or speech-shaped noise) as the independent variable in the model. As we  
349 hypothesised that presenting speech in noise would significantly decrease accuracy compared  
350 to without noise, we expected to find a significant effect of noise type from this GLMER  
351 analysis. Hypothesis (ii) was the same as the first but looked at reaction times to correctly  
352 discriminated speech-in-noise on the same task. A GLMER was used to test this hypothesis,  
353 using reaction times as the dependent variable and noise type as the independent variable.  
354 Similarly, we expected to find a significant effect of noise type, increasing reaction times.

355 To test hypothesis (iii), a GLMER analysis was conducted using the accuracy of  
356 responses on the speech discrimination task as the dependent variable and using stimulus type  
357 (purely audio or audiovisual) as the independent variable in the model. As we hypothesised  
358 that presenting audiovisual stimuli in noise would significantly increase accuracy compared  
359 to purely audio stimuli in noise, we expected to find a significant effect of stimulus type from  
360 this GLMER analysis. Hypothesis (iv) was the same as the third but looked at reaction times  
361 to correctly discriminated speech-in-noise on the same task. A GLMER was used to test this  
362 hypothesis, using reaction times as the dependent variable and stimulus type as the  
363 independent variable. Similarly, we expected to find a significant effect of stimulus type,  
364 decreasing reaction times.



365 To test hypothesis (v), we conducted a GLMER analysis using accuracy as a  
366 dependent variable and SOA as the independent variable. SOA was treated as a categorical  
367 variable in this model and the model for hypothesis (vi) below. We expected to find a  
368 significant effect of SOA, with accuracy decreasing when more asynchrony was introduced  
369 to the stimuli. This would reflect that the window of integration for audiovisual speech is  
370 important for visual information to be beneficial to understanding speech in noise. Finally, in  
371 a similar manner, hypothesis (vi) was tested using a GLMER analysis with reaction times as  
372 the dependent variable and with SOA levels as the independent variable in the model. Again,  
373 we expected a significant effect of SOA on reaction times, with reaction times increasing  
374 with the introduction of asynchrony.

375 For all six GLMER models listed above, the speech sound token used (Ba, Fa, or Ka),  
376 participant age and the participant ID were all included as random effects. No further model  
377 selection of these random and fixed effects was undergone, as we wanted a conservative  
378 model that included a full random effects structure to account for the expected larger  
379 individual differences of an online experiment. All model equations and structures can be  
380 found in the supplementary materials (Table 2).

381 Furthermore, we also conducted exploratory analyses to assess the effect of noise on  
382 speech discrimination accuracy between the three visually distinct, chosen phonemes ('Ba',  
383 'Fa', and 'Ka'). To do this, a GLMER analysis was conducted using accuracy as the  
384 dependent variable and speech token as the independent variable. Purely audio trials in noise  
385 were used for this analysis. Furthermore, we also conducted pairwise comparisons within the  
386 GLMER models used to test hypotheses (v) and (vi) as another exploratory analysis,  
387 comparing between each level of our SOA independent variable. We expect that not all the  
388 SOA interactions will show significance. As we expected the benefits of visual stimuli to  
389 only be present during the window of integration, there would only be a significant decrease

390 in accuracy and an increase in reaction times at SOAs outside this window. Therefore, this  
391 exploratory analysis can be used to better understand the window of integration for our  
392 stimuli. All exploratory analyses will use an inference criterion of  $p < .008$  as this was the  
393 strictest threshold for significance included in our Bonferroni-Holm correction.

394

395

## Results

### 396 Descriptive statistics

397 The means and standard deviations of the accuracy of responses and reaction times of  
398 responses can be seen in Table 1. Descriptive statistics were also calculated for each speech  
399 token (Ba, Fa, and Ka). Figure 1 shows the mean reaction times and mean accuracy rates for  
400 both audio-only and audiovisual stimuli when no SOA is considered (0 ms SOA), whilst  
401 Figure 2 shows these data for all SOAs when audiovisual stimuli are used for speech-in-noise  
402 conditions. Figure 3 shows the mean reaction times and accuracy rates for all SOAs when  
403 audiovisual stimuli are presented without noise. Furthermore, Figure 4 shows accuracy rates  
404 and reaction times in purely audio and audiovisual stimuli in noise between each of the three  
405 speech tokens. Violin plots were used for all figures to highlight the distribution of accuracies  
406 and reaction times across participants for each condition, as individual differences were large  
407 in this dataset likely due to online experimentation.

### 408 Effect of noise on speech perception

409 The first planned GLMER analysis was conducted to test hypothesis (i). There was a  
410 significant effect of noise type (with or without noise), showing a decrease in accuracy in  
411 speech-in-noise discrimination when noise was introduced versus clear speech ( $\beta = -.29$ ,  $t = -$   
412  $12.95$ , 99.2%  $CI = [-.35, -.23]$ ,  $p < .008$ ). This model supports hypothesis (i), as we expected  
413 to find that the introduction of noise to speech would decrease performance. For testing  
414 hypothesis (ii), the planned GLMER analysis was conducted. There was a significant effect  
415 of noise type on reaction times ( $\beta = .06$ ,  $t = 3.10$ , 99.2%  $CI = [.01, .11]$ ,  $p < .008$ ). This model  
416 supports hypothesis (ii), as we expected to find that introducing noise would increase reaction  
417 times to correctly discriminated speech.

### 418 Effect of congruent, distinguishable visual information on speech perception

419 Our next planned GLMER analysis was conducted to test hypothesis (iii). There was a  
420 significant effect of stimulus type (purely audio or audiovisual), as there was an increase in  
421 accuracy in speech-in-noise discrimination when stimulus type was audiovisual versus purely  
422 audio ( $\beta = .26, t = 11.36, 99.2\% CI = [.20, .32], p < .008$ ). This model supports hypothesis  
423 (iii), as we expected to find that introducing relevant visual information would improve  
424 speech perception in noise. For testing hypothesis (iv), the planned GLMER analysis was  
425 conducted. There was a significant effect of stimulus type on reaction times ( $\beta = -.08, t = -$   
426  $4.15, 99.2\% CI = [-.13, -.03], p < .008$ ). This model supports hypothesis (iv), as we expected  
427 to find that introducing relevant visual information would decrease reaction times and  
428 improve speech perception in noise.

#### 429 **Effect of stimulus onset asynchrony on audiovisual speech perception**

430 When testing hypothesis (v), the planned GLMER analysis was done for data across  
431 all SOA levels for audiovisual speech-in-noise stimuli only. There was no significant effect  
432 of SOA on accuracy at any interval, even at a 95% confidence interval, showing no support  
433 for hypothesis (v). Finally, our planned GLMER analysis was run to test hypothesis (vi).  
434 There was a significant main effect of SOA ( $\beta = .04, t = 3.31, p < .008$ ) on reaction times,  
435 indicating reaction times increased with SOA. This supports hypothesis (vi).

#### 436 **Exploratory analyses**

437 As a further, exploratory analysis, a GLMER model was used to investigate phoneme  
438 differences in speech-in-noise discrimination. Looking at pairwise comparisons, there was a  
439 significant difference between accuracy rates of the 'Ba' and 'Fa' tokens ( $\beta = -.17, t = -5.03,$   
440  $p < .008$ ), 'Ba' and 'Ka' tokens ( $\beta = -.53, t = -15.50, p < .001$ ), and 'Fa' and 'Ka' tokens ( $\beta =$   
441  $-.36, t = -10.47, p < .008$ ) for purely audio stimuli. For audiovisual stimuli, however, there  
442 was no significant change in accuracy rate between the three tokens. A GLMER model for

443 reaction times showed similar patterns, although only ‘Ba’ and ‘Ka’ were significantly  
444 different for purely audio stimuli, with ‘Ba’ having increased reaction times in comparison to  
445 ‘Ka’ ( $\beta = .14, t = 4.21, p < .008$ ).

446 Finally, to explore differences between SOA intervals to see if a window of  
447 integration could be determined, pairwise comparisons were made on the GLMER analyses  
448 used to test hypothesis (vi). Pairwise comparisons were not made on the GLMER used to test  
449 hypothesis (v) as no significant effect of SOA on accuracy was observed. Pairwise  
450 comparisons made on the GLMER to test hypothesis (vi) indicated that reaction times were  
451 significantly reduced compared to 0 ms at 250 ( $\beta = -.05, t = -3.94, p = .001$ ) and 266.6 ms ( $\beta$   
452  $= -.05, t = -3.88, p = .002$ ). However, no other comparisons between levels of SOA were  
453 significantly different. Whilst this implies that a minimal end of the window of integration  
454 could lie above 233.3 ms (as SOAs between 233.3 and 250 ms were not tested), no accurate  
455 window of integration can be determined from the data.

456

457

458

459

460

**Discussion**

461

462

463

464

465

466

467

468

469

This study aimed to reassess the contribution of audiovisual integration to speech perception in noise when stimuli belonged to different viseme categories. As speech perception can differ wildly with stimuli sets, it was important to first reassess the detriment of noise on speech discrimination, as well as the benefits of speech-relevant visual integration. The study incorporated the visual distinguishability of each speech phoneme used in the speech discrimination task by selecting phonemes from separate viseme categories. Furthermore, the study also aimed to examine the effects of stimulus onset asynchrony (SOA) on audiovisual speech perception. This may assist in determining a window of integration for these stimuli, which was explored in further analyses.

470

**Reassessing the detriment of noise on speech perception**

471

472

473

474

475

476

477

478

479

GLMERs were used to investigate the influence of the predictor variables on accuracy ratings on the speech discrimination task. The first model, using noise type as the predictor, supported our first hypothesis, showing that there was a decrease in accuracy for purely audio stimuli when the speech was presented in noise and not without. Additionally, the introduction of noise to the speech signal increased reaction times significantly. These results support our second hypothesis. As both the accuracy and reaction time to trials with noise differed significantly from those without, it can be said that the detriment of noise on speech perception was present with our created stimuli and chosen SNR ratio of -16 dB using speech-shaped white noise.

480

**Reassessing the contribution of audiovisual information on speech processing in noise**

481

482

483

There was a significant increase in accuracy when relevant, congruent visual information was present with the stimuli versus purely audio stimuli in noise. This supports hypothesis (iii) and confirms previous findings regarding the contribution of audiovisual

484 information to speech-in-noise processing. However, it should be noted that whilst the effect  
485 is prominent, it is not as great as previous literature findings which used a similar speech-to-  
486 noise ratio (Van de Rijt *et al.*, 2019). Here, the effectiveness of audiovisual enhancement of  
487 speech recognition was assessed with SNR ratios as low as -21 dB SNR, where the  
488 introduction of relevant visual cues provided an increase in accuracy of up to 50% for some  
489 stimuli, with greater enhancements for words like ‘Pieter’. Even at -16 dB SNR, Van de Rijt  
490 *et al.*’s data suggests that greater audiovisual enhancement should have been seen, though  
491 reaction time data was not reported in the study.

492         This could also be explained using results from our exploratory analysis. When the  
493 speech was in noise and the stimuli contained auditory information only, the token ‘Ba’  
494 displayed much lower mean accuracy scores than the other tokens. This suggests that there  
495 are specific differences in the acoustic properties of the tokens used that are influencing the  
496 perception of speech-in-noise. In previous literature, ‘Ba’ and other tokens within the same  
497 viseme (such as ‘Pa’) are frequently used, which could suggest why results in previous  
498 literature show a larger speech discrimination effect in noise. It is therefore important for  
499 future research to determine if there are differences in speech perception between other  
500 viseme categories that were not used in this study (Fisher, 1968). In our LMER model for  
501 hypothesis (iii), the token used was loaded as a random factor. This variance between tokens  
502 was removed from the variance found in fixed effects in the outputs of the model. This mixed  
503 effect modelling also considered participant differences and age, unlike previous literature  
504 that did not investigate speech discrimination effects using more complex models (Bernstein  
505 *et al.*, 2004; Sekiyama *et al.*, 2014). As the tokens appear to be largely variant, this could  
506 further account for the weaker overall patterns of change seen between fixed effects.

507         Next, there was a significant decrease in reaction times when audiovisual stimuli were  
508 used over purely audio, supporting hypothesis (iv). Interestingly, there was a decrease in

509 reaction time in audiovisual conditions with noise over without noise as well. When  
510 processing multisensory stimuli that are not beneficial to us, reaction times likely increase  
511 due to extra unnecessary processing (Brown & Strand, 2019). In this case, the audiovisual  
512 information is only beneficial to us in noise. Therefore, in this model where no comparisons  
513 to clear speech are made, reaction times significantly decrease with the introduction of noise  
514 as the extra processing of visual information becomes beneficial. Comparatively, when  
515 audiovisual information is present without noise, reaction times increase as the added visual  
516 information is no longer beneficial to speech recognition as it is already clear to understand.

517 **Investigating the effects of stimulus onset asynchrony on the speech processing benefits**  
518 **of audiovisual information**

519 Our GLMER model testing hypothesis (v) uncovered no meaningful change in  
520 accuracy between any SOA value. In previous research, the maximal window of integration  
521 was around 250 to 260 ms for syllables (Dixon & Spitz, 1980). Here, SOAs up to 266.6 ms  
522 did not affect speech discrimination accuracy, implying that the stimuli were still inside the  
523 window of integration and that the maximal end of the window lies beyond 266.6 ms. Our  
524 final LMER model testing hypothesis (vi) found significant increases in reaction time when  
525 an SOA was introduced. Alternatively, this implies that the range of SOAs used does cover  
526 the maximal end of the window concerning processing speed, as there was a gradual increase  
527 in reaction times as SOA was further increased reducing the benefit of audiovisual  
528 information. When looking at exploratory pairwise comparisons between SOA levels, there  
529 was a distinct decrease in reaction times at 250 and 266.6 ms compared to no SOA. This  
530 implies that the ability to discriminate the speech was made less taxing past 250 ms  
531 asynchrony. It could be, based on these findings, that the minimal end of the window of  
532 integration for our created stimuli lies between 233.3 and 250 ms. Given that the stimuli were  
533 simple syllables, an alternative interpretation may be that the processing of the auditory and



534 the visual information was completed before integration had finished, although this would not  
535 explain the differences in reaction times between the SOA levels. Furthermore, as  
536 participants could only respond after the stimuli had played in full with visual cues preceding  
537 the auditory cues, we would expect integration to have occurred as long as the SOA remained  
538 within the window of integration. As these comparisons are exploratory, however, and there  
539 is no account of accuracy changing with SOAs, further research would be needed to  
540 determine the full window of integration.

#### 541 **Limitations of the study and future directions**

542 One explanation for the audiovisual benefit in our data not being as large as in  
543 previous studies could be the lack of ecological validity and the artificial nature of the online  
544 experimentation. Speech-shaped white noise was utilised for speech-in-noise conditions.  
545 Despite this noise modulating speech, it is still unlike that in a real environment. This may  
546 mean that the speech-shaped noise was too distinct from the speech itself, especially  
547 considering that we used syllables for recognition rather than words or sentences. Speaker  
548 babble or background noise such as light vocal music would be much more akin to that in  
549 everyday life, making it perhaps more suitable and valid for investigating audiovisual speech  
550 perception when speech is in noise (Krishnamurthy & Hansen, 2009). Furthermore, the  
551 stimuli used were single syllable speech tokens, which do not reflect typical communicative  
552 speech in a real-world environment. Given their simplicity, other aspects of speech  
553 perception, such as prediction of oncoming words in larger sentences, would not be used as a  
554 method of speech processing here (Solberg Økland *et al.*, 2019). The overall simplicity and  
555 artificial design of these stimuli may be obscuring other benefits of audiovisual integration in  
556 speech perception when applied to realistic speech settings. To better reassess audiovisual  
557 integration in speech, further research with more ecologically valid speech stimuli (e.g., full  
558 sentences) would be of benefit.

559           The SNR used for our study was -16 dB. This was selected based on previous  
560 research investigating audiovisual syllable perception in noise, for which there was a notable  
561 difference between perceiving speech with or without visual aid (Bernstein *et al.*, 2004).  
562 However, whilst this may have been true for speech token ‘Ba’, this did not seem to translate  
563 to ‘Ka’, indicating that different speech viseme categories were affected by speech-shaped  
564 noise at the SNR -16 dB. Furthermore, initial data collection for this study was conducted  
565 from 2021 to 2022 after multiple lockdowns in the UK due to the COVID-19 pandemic.  
566 Many adults in the UK during this time had been socially distancing and wearing facemasks  
567 to prevent contamination. These facemasks would obscure the lip and mouth area of the  
568 wearer, meaning that social interactions between many people in this period would have  
569 lacked visual information to assist with speech perception. In many cases, the facemasks  
570 obscured sound, making it more difficult to understand speech and imitating difficult  
571 listening conditions (Yi *et al.*, 2021; Smiljanic *et al.*, 2021). It is possible that due to  
572 facemask wearing for a year, participants had adapted to listening to speech in difficult  
573 conditions without visual aids. Furthermore, only three phonemes from three viseme  
574 categories were used in this study. As there was an apparent difference between these  
575 phonemes selected, with ‘Ba’ being more impacted by added noise than ‘Ka’, future studies  
576 may wish to investigate the differences between more viseme categories and the phonemes  
577 within them. It may also be beneficial to further apply this to more than single-syllable units  
578 of speech. This would provide a broader view of the contributions of visual information to  
579 speech processing.

580           Finally, this experiment did use home equipment to record stimuli as well as the home  
581 equipment of participants to play the stimuli through online experimentation. Whilst the  
582 recording equipment was of laboratory standard and the recording procedure rigorous, there  
583 will still be discrepancies between these stimuli and other lab-created stimuli which might

584 make replications difficult. Furthermore, the environments that participants were in whilst  
585 taking part in the study may be different between participants. We do not have measures of  
586 how well participants understood the task, how noisy their environment was during listening,  
587 the hardware they used to run the study, and if they followed pre-experiment instructions  
588 such as to wear headphones. These are likely to contribute to the large individual differences  
589 seen in the dataset. Whilst GLMER models can consider the participant differences, further  
590 in-person lab testing with similar methodologies may be needed to fully control these  
591 confounds.

## 592 **Conclusion**

593         A set of purely audio and viseme-controlled audiovisual stimuli was created to  
594 investigate the contributions of audiovisual information to speech-in-noise processing.  
595 Introducing visual information increased accuracy and decreased reaction times in speech-in-  
596 noise conditions relative to audio-only stimuli. When looking at accuracy and reaction times  
597 at varying SOA intervals in our audiovisual stimuli, introducing SOAs influenced reaction  
598 times, but not accuracy. In the future, more syllables from more viseme categories could be  
599 tested to investigate a full range of speech sounds in audio-only and audio-visual contexts, as  
600 well as with further SOA intervals to ensure that a window of integration can be determined  
601 with accuracy.

602

## Acknowledgements

603

This work was supported by the Economic and Social Research Council (ESRC)

604

Training Grant (O'Hanlon, ES/P000665/1), the Manchester Biomedical Research Centre and

605

the National Institute for Health and Care Research (NIHR) (Plack, NIHR203308), and the

606

Biotechnology and Biological Sciences Research Council (BBSRC) New Investigator Grant

607

(Nuttall, BB/S008527/1). We would like to thank all participants who expressed interest and

608

participated in this research. We also thank Kyle Stonehouse for their contributions to study

609

material creation.

610

## 611 Data Availability Statement

612 Upon publication, all collected data are available to view online through the Open  
613 Science Framework (OSF: <https://osf.io/kcbzs>), as well as all stimuli used in the experiment  
614 code relevant to data analysis.

615 **References**

- 616 Altieri, N., Townsend, J. T., & Wenger, M. J. (2014). A measure for assessing the effects of  
617 audiovisual speech integration. *Behavior Research Methods*, *46*(2), 406-415.
- 618 Audacity Team. (2021). *Audacity(R): Free Audio Editor and Recorder* [Computer  
619 application]. Version 3.0.0 retrieved March 17th, 2021, from <https://audacityteam.org/>  
620 . *Copyright statement: Audacity® software is copyright © 1999-2021 Audacity Team.*  
621 *The name Audacity® is a registered trademark.*
- 622 Bates, D., Maechler, M., & Bolker, B. (2015). Walker., S. Fitting linear mixed-effects models  
623 using lme4. *J Stat Softw*, *67*(1), 1-48.
- 624 Bernstein, L. E., Auer Jr, E. T., & Takayanagi, S. (2004). Auditory speech detection in noise  
625 enhanced by lipreading. *Speech Communication*, *44*(1-4), 5-18.
- 626 Boersma, P., & Weenink, D. (2021). *Praat: doing phonetics by computer* [Computer  
627 program]. Version 6.1.53, retrieved August 2021 from <http://www.praat.org/>
- 628 Bowers, J. S., Kazanina, N., & Andermane, N. (2016). Spoken word identification involves  
629 accessing position invariant phoneme representations. *Journal of Memory and*  
630 *Language*, *87*, 71-83.
- 631 Brown, V. A., & Strand, J. F. (2019). About face: Seeing the talker improves spoken word  
632 recognition but increases listening effort. *Journal of Cognition*, *2*(1).
- 633 Brown, V. A., Van Engen, K. J., & Peelle, J. E. (2021). Face mask type affects audiovisual  
634 speech intelligibility and subjective listening effort in young and older adults.  
635 *Cognitive Research: Principles and Implications*, *6*(1), 1-12.

- 636 Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., & Ghazanfar, A. A. (2009).  
637 The natural statistics of audiovisual speech. *PLoS computational biology*, *5*(7),  
638 e1000436.
- 639 Colonus, H., & Diederich, A. (2010). The optimal time window of visual-auditory  
640 integration: a reaction time analysis. *Frontiers in integrative neuroscience*, *4*, 1316.
- 641 Dixon, N. F., & Spitz, L. (1980). The detection of auditory visual desynchrony. *Perception*,  
642 *9*(6), 719-721.
- 643 Ewen, C. J., & Van der Hulst, H. (2001). *The phonological structure of words: an*  
644 *introduction*. Cambridge University Press.
- 645 Fisher, C. G. (1968). Confusions among visually perceived consonants. *Journal of speech*  
646 *and hearing research*, *11*(4), 796-804.
- 647 Ganesh, A. C., Berthommier, F., & Schwartz, J. L. (2018). Audiovisual binding for speech  
648 perception in noise and in aging. *Language Learning*, *68*, 193-220.
- 649 Grant, K. W., & Seitz, P. F. (1998). Measures of auditory–visual integration in nonsense  
650 syllables and sentences. *The Journal of the Acoustical Society of America*, *104*(4),  
651 2438-2450.
- 652 Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian*  
653 *journal of statistics*, 65-70.
- 654 Krishnamurthy, N., & Hansen, J. H. (2009). Babble noise: modeling, analysis, and  
655 applications. *IEEE transactions on audio, speech, and language processing*, *17*(7),  
656 1394-1407.

- 657 Kuznetsova, A., Brockhoff, P. B., Christensen, R. H. B., & Jensen, S. P. (2020). *lmerTest:*  
658 *Tests in linear mixed effects models* [computer manual]. Retrieved from:  
659 <https://cran.r-project.org/web/packages/lmerTest/index.html>
- 660 Lee, H., & Noppeney, U. (2014). Music expertise shapes audiovisual temporal integration  
661 windows for speech, sinewave speech, and music. *Frontiers in psychology, 5*, 868.
- 662 Lee, J., Han, J., & Lee, H. (2020). Long-Term Musical Training Alters Auditory Cortical  
663 Activity to the Frequency Change. *Frontiers in Human Neuroscience, 14*, 329.
- 664 Lo, S., & Andrews, S. (2015). To transform or not to transform: Using generalized linear  
665 mixed models to analyse reaction time data. *Frontiers in psychology, 6*, 148545.
- 666 Lotfi, Y., Chupani, J., Javanbakht, M., & Bakhshi, E. (2019). Evaluation of speech perception  
667 in noise in Kurd-Persian bilinguals. *Auditory and Vestibular Research, 28*(1), 36-41.
- 668 Luke, S. G. (2017). Evaluating significance in linear mixed-effects models in R. *Behavior*  
669 *research methods, 49*(4), 1494-1502.
- 670 Ma, W. J., Zhou, X., Ross, L. A., Foxe, J. J., & Parra, L. C. (2009). Lip-reading aids word  
671 recognition most in moderate noise: a Bayesian explanation using high-dimensional  
672 feature space. *PloS one, 4*(3), e4638.
- 673 Maier, J. X., Di Luca, M., & Noppeney, U. (2011). Audiovisual asynchrony detection in  
674 human speech. *Journal of Experimental Psychology: Human Perception and*  
675 *Performance, 37*(1), 245.
- 676 Massaro, D. W., Cohen, M. M., Tabain, M., & Beskow, J. (2012). Animated speech:  
677 Research progress and applications In Clark RB, Perrier J, P, & Vatikiotis-Bateson E  
678 (Eds.), *Audiovisual Speech Processing* (pp. 246–272). *Cambridge: Cambridge*  
679 *University.*



- 680 McCloy, D. (2021). *Praat Script: 'Mix speech with noise'* [Praat script]. LICENSED  
681 UNDER THE GNU GENERAL PUBLIC LICENSE v3.0:  
682 <http://www.gnu.org/licenses/gpl.html>
- 683 McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature*, 264(5588),  
684 746-748.
- 685 Megnin-Viggars, O., & Goswami, U. (2013). Audiovisual perception of noise vocoded  
686 speech in dyslexic and non-dyslexic adults: the role of low-frequency visual  
687 modulations. *Brain and language*, 124(2), 165-173.
- 688 Mersmann, O., Trautmann, H., Steuer, D., & Bornkamp, B. (2018). *Truncated normal*  
689 *distribution* [computer manual]. Retrieved from: [https://cran.r-](https://cran.r-project.org/web/packages/truncnorm/)  
690 [project.org/web/packages/truncnorm/](https://cran.r-project.org/web/packages/truncnorm/)
- 691 Michalek, A. M., Watson, S. M., Ash, I., Ringleb, S., & Raymer, A. (2014). Effects of noise  
692 and audiovisual cues on speech processing in adults with and without ADHD.  
693 *International journal of audiology*, 53(3), 145-152.
- 694 Navarra, J., Vatakis, A., Zampini, M., Soto-Faraco, S., Humphreys, W., & Spence, C. (2005).  
695 Exposure to asynchronous audiovisual speech extends the temporal window for  
696 audiovisual integration. *Cognitive Brain Research*, 25(2), 499-507.
- 697 Noel, J. P., Stevenson, R. A., & Wallace, M. T. (2018). Atypical audiovisual temporal  
698 function in autism and schizophrenia: similar phenotype, different cause. *European*  
699 *Journal of Neuroscience*, 47(10), 1230-1241.
- 700 Open Broadcaster Software. (2024). OBS Studio (Version 29.1.3) [Computer software].  
701 <https://obsproject.com/>

- 702 Peirce, J. W., Gray, J. R., Simpson, S., MacAskill, M. R., Höchenberger, R., Sogo, H.,  
703 Kastman, E., & Lindeløv, J. (2019). PsychoPy2: experiments in behavior made easy.  
704 *Behavior Research Methods*. 10.3758/s13428-018-01193-y
- 705 Peter, G., Catriona, M., & Phillip, A. (2019). *Power analysis for generalised linear mixed*  
706 *models by simulation* [computer manual]. Retrieved from: [https://cran.r-](https://cran.r-project.org/web/packages/simr/index.html)  
707 [project.org/web/packages/simr/index.html](https://cran.r-project.org/web/packages/simr/index.html)
- 708 Qualtrics. (2005). *Qualtrics software*, Provo, Utah, USA. Copyright@2021, Current version:  
709 09-21. Retrieved from: <https://www.qualtrics.com>
- 710 Ren, Y., Yang, W., Nakahashi, K., Takahashi, S., & Wu, J. (2017). Audiovisual integration  
711 delayed by stimulus onset asynchrony between auditory and visual stimuli in older  
712 adults. *Perception*, 46(2), 205-218.
- 713 Satterthwaite, F. E. (1941). *Synthesis of variance*. *Psychometrika*, 6(5), 309–316.
- 714 Schwartz, J. L., & Savariaux, C. (2014). No, there is no 150 ms lead of visual speech on  
715 auditory speech, but a range of audiovisual asynchronies varying from small audio  
716 lead to large audio lag. *PLoS Computational Biology*, 10(7), e1003743.
- 717 Sekiyama, K., Soshi, T., & Sakamoto, S. (2014). Enhanced audiovisual integration with  
718 aging in speech perception: a heightened McGurk effect in older adults. *Frontiers in*  
719 *Psychology*, 5, 323.
- 720 Singmann, H., Bolker, B., Westfall, J., Aust, F., Ben-Shachar, M. S., Højsgaard, S., Fox, J.,  
721 Lawrence, M. A., Mertens, U., Love, J., Lenth, R., & Christensen, R. H. B. (2024).  
722 *Analysis of factorial experiments* [computer manual]. Retrieved from: [https://cran.r-](https://cran.r-project.org/web/packages/afex/index.html)  
723 [project.org/web/packages/afex/index.html](https://cran.r-project.org/web/packages/afex/index.html)

- 724 Smiljanic, R., Keerstock, S., Meemann, K., & Ransom, S. M. (2021). Face masks and  
725 speaking style affect audio-visual word recognition and memory of native and non-  
726 native speech. *The Journal of the Acoustical Society of America*, *149*(6), 4013-4023.
- 727 Smith, E. G., & Bennetto, L. (2007). Audiovisual speech integration and lipreading in autism.  
728 *Journal of Child Psychology and Psychiatry*, *48*(8), 813-821.
- 729 Solberg Økland, H., Todorović, A., Lüttke, C. S., McQueen, J. M., & De Lange, F. P. (2019).  
730 Combined predictive effects of sentential and visual constraints in early audiovisual  
731 speech processing. *Scientific Reports*, *9*(1), 7870.
- 732 Sorati, M., & Behne, D. M. (2019). Musical Expertise Affects Audiovisual Speech  
733 Perception: Findings From Event-Related Potentials and Inter-trial Phase Coherence.  
734 *Frontiers in Psychology*, *10*, 2562.
- 735 Stein, B. E., & Meredith, M. A. (1993). *The merging of the senses*. The MIT Press.
- 736 Sumbly, W., & Pollack, I. (1954). Visual Contribution to Speech Intelligibility in Noise. *The*  
737 *Journal of the Acoustical Society of America*, *26*, 212-215.
- 738 Summerfield, Q. (1992). Lipreading and audio-visual speech perception. Philosophical  
739 Transactions of the Royal Society of London. *Series B: Biological Sciences*,  
740 *335*(1273), 71-78.
- 741 Ten Oever, S., Sack, A. T., Wheat, K. L., Bien, N., & Van Atteveldt, N. (2013). Audio-visual  
742 onset differences are used to determine syllable identity for ambiguous audio-visual  
743 stimulus pairs. *Frontiers in psychology*, *4*, 331.
- 744 Van de Rijt, L. P., Roye, A., Mylanus, E. A., Van Opstal, A. J., & Van Wanrooij, M. M.  
745 (2019). The principle of inverse effectiveness in audiovisual speech perception.  
746 *Frontiers in human neuroscience*, *13*, 335.

- 747 Van Engen, K. J., Dey, A., Sommers, M. S., & Peelle, J. E. (2022). Audiovisual speech  
748 perception: Moving beyond McGurk. *The Journal of the Acoustical Society of*  
749 *America*, 152(6), 3216-3225.
- 750 Van Engen, K. J., Xie, Z., & Chandrasekaran, B. (2017). Audiovisual sentence recognition  
751 not predicted by susceptibility to the McGurk effect. *Attention, Perception, &*  
752 *Psychophysics*, 79, 396-403.
- 753 Varnet, L., Wang, T., Peter, C., Meunier, F., & Hoen, M. (2015). How musical expertise  
754 shapes speech perception: Evidence from auditory classification images. *Scientific*  
755 *Reports*, 5(1), 14489.
- 756 Yi, H., Pingsterhaus, A., & Song, W. (2021). Effects of wearing face masks while using  
757 different speaking styles in noise on speech intelligibility during the COVID-19  
758 pandemic. *Frontiers in psychology*, 12.
- 759 Yuan, Y., Lleo, Y., Daniel, R., White, A., & Oh, Y. (2021). The impact of temporally  
760 coherent visual cues on speech perception in complex auditory environments.  
761 *Frontiers in neuroscience*, 15, 678029.

762 Supplementary Materials

763 The supplementary materials contain the following file:

764 A caption and display of Table 2, showing the model equations used for each of the  
765 six GLMERs.

766

767

## Tables and Figures

768

Speech	Stimuli	SOA (ms)	Accuracy Rate (%)		Reaction Time (ms)	
			Mean	Std. Dev	Mean	Std. Dev
Clear	AO	0	96.11	10.24	538	216
Clear	AV	0	96.60	13.55	564	257
Clear	AV	200	95.95	14.09	551	241
Clear	AV	216.6	96.56	13.82	573	232
Clear	AV	233.3	96.58	12.61	575	249
Clear	AV	250	96.54	13.90	568	236
Clear	AV	266.6	96.84	13.23	575	255
Noise	AO	0	67.33	21.91	597	285
Noise	AV	0	93.10	15.21	518	239
Noise	AV	200	92.87	16.08	553	223
Noise	AV	216.6	93.62	15.75	554	218
Noise	AV	233.3	93.35	17.52	562	227
Noise	AV	250	93.11	16.30	570	224
Noise	AV	266.6	93.26	15.01	569	237

769

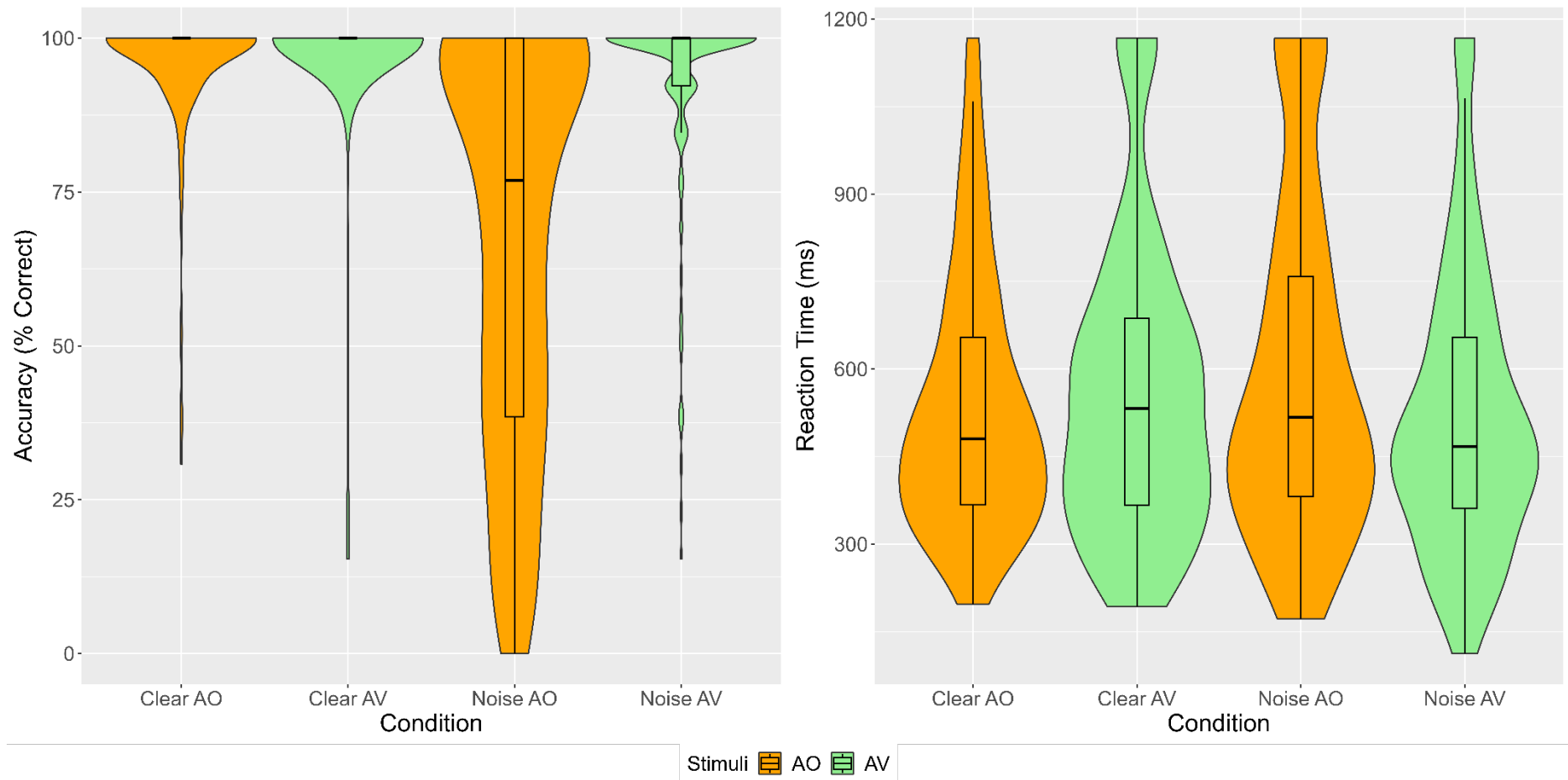
770 **Table 1:**

771 *Means and Standard Deviations (Std. Dev) of accuracy rates and reaction times for speech*  
772 *with and without noise, audio-only (AO) or audiovisual (AV) stimuli, and different stimulus*  
773 *onset asynchronies (SOAs), with each speech token and participant aggregated into a single*  
774 *mean.*

775

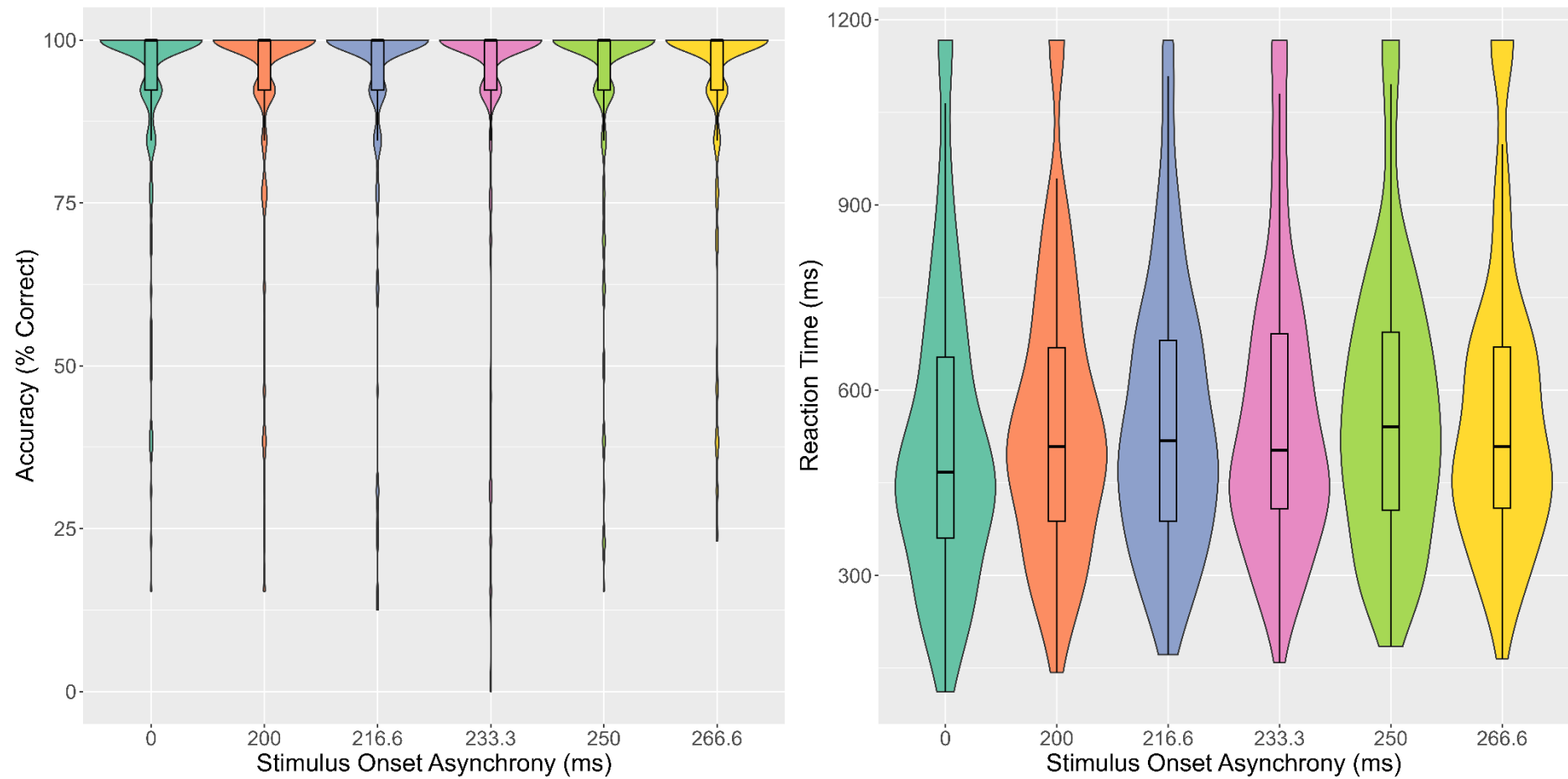
776 *Figure 1:* Violin plots showing the accuracy rates and reaction times of participants when speech was presented either with or without noise, for  
 777 both audio-only (AO) and audiovisual (AV) stimuli. Boxplots show the median and interquartile ranges for each condition.

778



779 *Figure 2.* Violin plots showing the accuracy rates and reaction times of participants when audiovisual stimuli were presented in noise at different  
780 SOAs. Boxplots show the median and interquartile ranges for each condition.

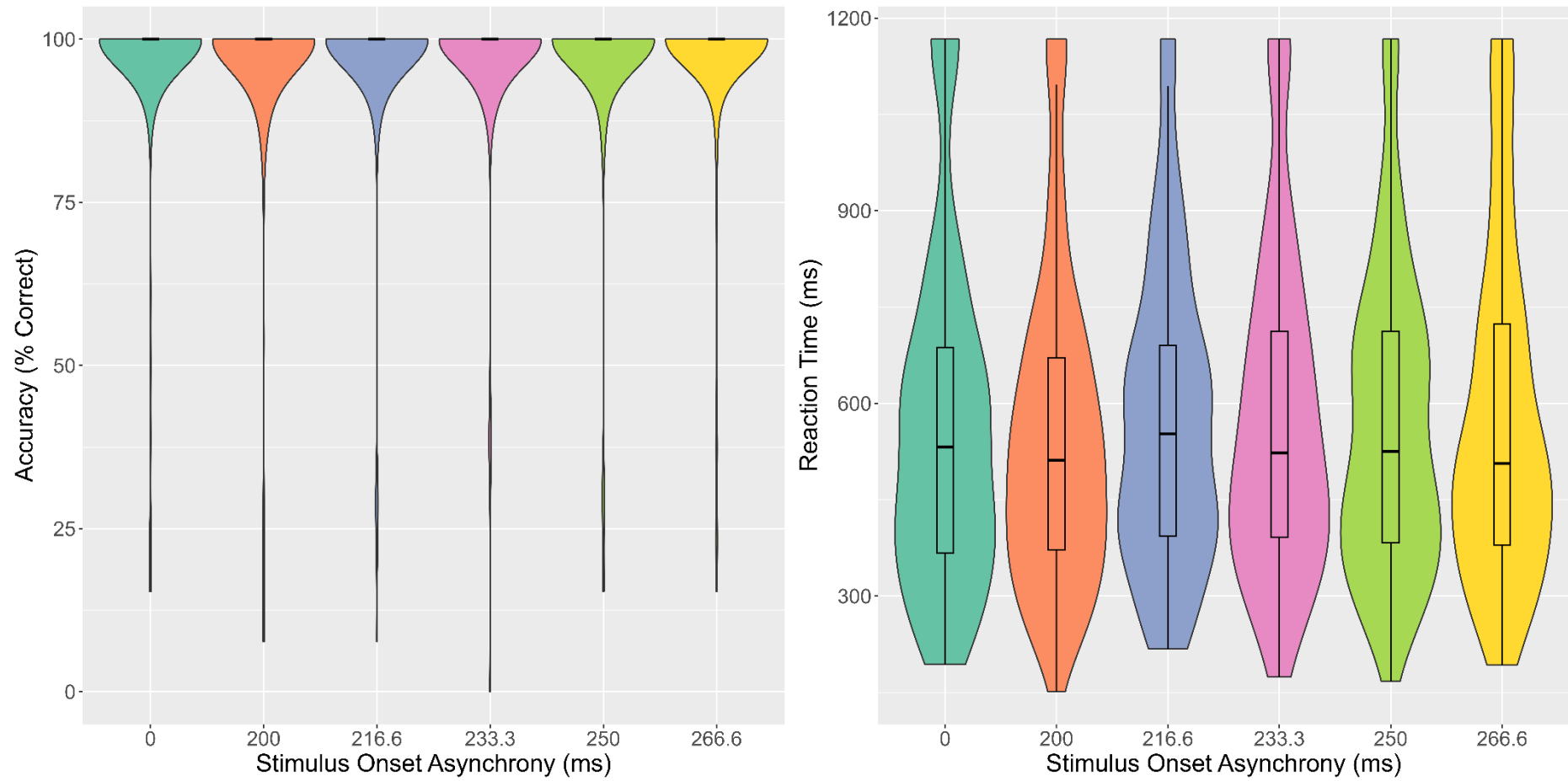
781





782 *Figure 3.* Violin plots showing the accuracy rates and reaction times of participants when audiovisual stimuli were presented without noise at  
783 different SOAs. Boxplots show the median and interquartile ranges for each condition.

784



785 *Figure 4.* Violin plots showing the accuracy rates and reaction times of participants when speech tokens were investigated individually in noise  
 786 for both Audio-Only (AO) and Audiovisual (AV) stimuli. Boxplots show the median and interquartile ranges for each condition.

787

