

# Attention, Perception, & Psychophysics

## Effects of Stimulus Response Compatibility on Covert Imitation of Vowels

Journal:	<i>Attention, Perception, &amp; Psychophysics</i>
Manuscript ID	PP-ORIG-17-263.R2
Manuscript Type:	Original Manuscript
Date Submitted by the Author:	n/a
Complete List of Authors:	Adank, Patti; UCL, Speech, Hearing and Phonetic Sciences Nuttall, Helen; University of Lancaster, Department of Psychology Bekkering, Harold; Radboud University Nijmegen, Donders Institute for Brain, Cognition and Behaviour Maegherman, Gwijde; UCL, Speech, Hearing and Phonetic Sciences
Keywords:	speech perception, Multisensory Processing

SCHOLARONE™  
Manuscripts

Only

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1 RUNNING HEAD: Covert Vowel Imitation

2

3

4

5

6

7 Effects of Stimulus Response Compatibility on Covert Imitation of Vowels

8

9 Patti Adank<sup>1</sup>, Helen Nuttall<sup>2,1</sup>, Harold Bekkering<sup>3</sup>, Gwijde Maegherman<sup>1</sup>

10 <sup>1</sup>Department of Speech, Hearing and Phonetic Sciences, University College London,

11 Chandler House, 2 Wakefield Street, London, UK, WC1N 1PF

12 <sup>2</sup> Department of Psychology, Lancaster University, Lancaster, UK, LA1 4YF

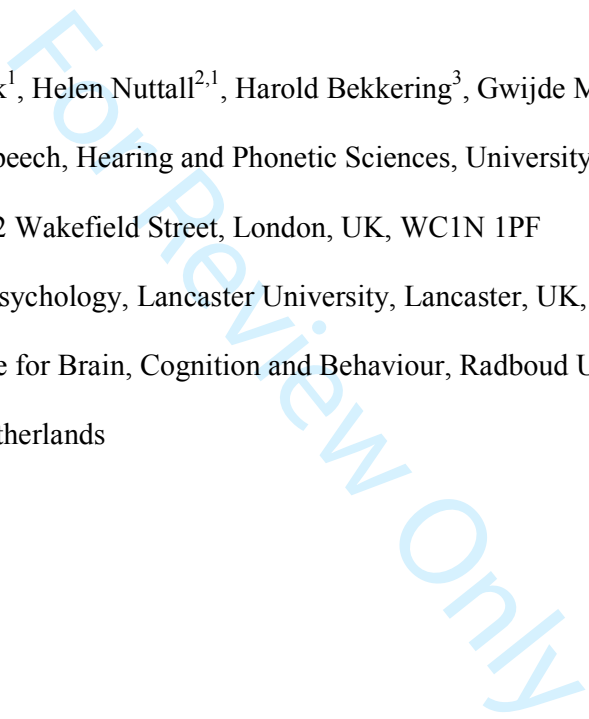
13 <sup>3</sup> Donders Institute for Brain, Cognition and Behaviour, Radboud University,

14 Nijmegen, the Netherlands

15

16

17



## 18 Abstract

19 When we observe someone else speaking, we tend to automatically activate the  
20 corresponding speech motor patterns. When listening we therefore covertly imitate  
21 the observed speech. Simulation theories of speech perception propose that covert  
22 imitation of speech motor patterns supports speech perception. Covert imitation of  
23 speech has been studied with interference paradigms including the Stimulus Response  
24 Compatibility paradigm (SRC). The SRC paradigm measures covert imitation by  
25 comparing articulation of a prompt following exposure to a distracter. Responses tend  
26 to be faster for congruent than incongruent distracters; thus showing evidence of  
27 covert imitation. Simulation accounts propose a key role for covert imitation in  
28 speech perception. However, covert imitation has thus far only been demonstrated for  
29 a select class of speech sounds, namely consonants, and it is unclear whether covert  
30 imitation extends to vowels. We aimed to demonstrate that covert imitation effects as  
31 measured with the SRC paradigm extend to vowels, in two experiments. We  
32 examined whether covert imitation occurs for vowels in a consonant-vowel-consonant  
33 context in Visual, Audio, and Audiovisual modalities. We presented the prompt at  
34 four time points to examine how covert imitation varied over the distracter's duration.  
35 The results of both experiments clearly demonstrated covert imitation effects for  
36 vowels, thus supporting simulation theories of speech perception. Covert imitation  
37 was not affected by stimulus modality and was maximal for later time points.

38

## 39 Keywords

40 Speech perception, speech production, multisensory processing

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 41 Effects of Stimulus Response Compatibility on Covert Imitation of Vowels

42

43 Observing someone else perform an action has been shown to activate neural  
44 mechanisms required to perform that action (Buccino et al., 2004; Fadiga, Craighero,  
45 Buccino, & Rizzolatti, 2002). For speech, this type of *covert imitation* occurs  
46 whenever we hear and/or see someone speaking and involves activation of speech  
47 production mechanisms (Nuttall, Kennedy-Higgins, Devlin, & Adank, 2017; Nuttall,  
48 Kennedy-Higgins, Hogan, Devlin, & Adank, 2016; Watkins, Strafella, & Paus, 2003).  
49 Covert imitation processes are proposed to play a key role in current speech  
50 perception theories, commonly referred to as *simulation accounts* (Pickering &  
51 Garrod, 2013; Wilson & Knoblich, 2005). Simulation accounts propose that listening  
52 to speech results in automatic activation of the articulatory motor plans for producing  
53 speech. These motor plans consist of simulations of the movements of articulators that  
54 are generated while the listener is processing the incoming speech signal. The  
55 generated motor plans then inform forward models of the heard speech that run in  
56 parallel with the unfolding speech signal (Kawato, 1999). Forward models are thought  
57 to use implicit knowledge of the perceiver's articulatory mechanics as a real-time  
58 mental simulation to track others' speech that support speech perception. These  
59 mental simulations generate top-down predictions of incoming speech, serving as a  
60 prediction signal supporting perception and thereby streamlining interaction.

61 Covert imitation in speech can be demonstrated using neuroimaging methods  
62 including functional Magnetic Resonance Imaging (fMRI), neurostimulation methods  
63 such as Transcranial Magnetic Stimulation (TMS), or using behavioural paradigms.  
64 Using fMRI, it was demonstrated that passively listening to speech broadly activates  
65 speech production regions, including motor and pre-motor areas (Wilson, Saygin,

1  
2  
3 66 Sereno, & Iacoboni, 2004). Areas in primary motor cortex (M1) have been found to  
4  
5 67 respond in a somatotopic manner during speech perception: Areas of M1 show  
6  
7 68 activation congruent with the primary articulator producing the perceived speech  
8  
9 69 stimulus. (Pulvermüller et al., 2006) used fMRI to demonstrate that lip and tongue  
10  
11 70 areas of M1 responded in a somatotopic manner when participants listened to sounds  
12  
13 71 produced with the lips (/p/) and the tongue (/t/).

14  
15  
16 72 Using TMS, a causal link has been demonstrated between articulatory M1 and  
17  
18 73 the efficacy of perception of sounds articulated using the congruent articulator  
19  
20 74 (D'Ausilio et al., 2009; Möttönen & Watkins, 2009). D'Ausilio et al. administered  
21  
22 75 TMS pulses to lip or tongue M1 while participants performed a discrimination task  
23  
24 76 for sounds produced with the lips (/p/ and /b/) or tongue (/t/ and /d/) as active  
25  
26 77 articulators. D'Ausilio et al. report a double dissociation in speech sound  
27  
28 78 discrimination: Participants showed poorer discrimination for lips sounds, but not for  
29  
30 79 tongue sounds, after a TMS pulse to the lips, and vice versa. Möttönen & Watkins  
31  
32 80 (2009) asked participants to perform a categorical perception task of spoken syllables  
33  
34 81 before administering 15 minutes of offline repetitive TMS to lip M1. After receiving  
35  
36 82 TMS, participants repeated the task and showed impaired categorical perception of  
37  
38 83 syllables involving lip sounds (/pa/-/ba/ and /pa/-/ta/) but not tongue sounds (/ka/-/ga/  
39  
40 84 and /da/-/ga/).

41  
42  
43 85 Besides establishing causal links between a brain area and behaviour, TMS has  
44  
45 86 also been used to estimate the relative excitability of the corticobulbar tract  
46  
47 87 innervating speech muscles (Adank, Nuttall, & Kennedy-Higgins, 2016) while  
48  
49 88 listening to speech. Following a TMS pulse to an area in articulatory M1, it is possible  
50  
51 89 to record the resulting action potentials, Motor Evoked Potentials (MEPs), in the  
52  
53 90 corresponding muscle. Increased MEPs while perceiving speech can be regarded to  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 91 imply covert imitation. This covert imitation response is also somatotopic in nature  
4  
5 92 and, for instance, also reflects the clarity with which the speech stimulus was  
6  
7 93 produced. (Nuttall et al., 2016) measured MEPs from lip M1 while participants were  
8  
9 94 listening to clearly spoken syllables (/apa/, /aba/, /ata/, and /ada/) and distorted  
10  
11 95 syllables (produced with a tongue depressor in the speaker's mouth). As in Möttönen  
12  
13 96 & Watkins and D'Ausilio et al., participants showed somatotopic effects: Lip M1 was  
14  
15 97 facilitated for lip sounds, and further facilitation was measured for distorted lip  
16  
17 98 sounds. Moreover, (Sato, Buccino, Gentilucci, & Cattaneo, 2009) demonstrated that  
18  
19 99 somatotopic effects extend to visual speech processing; they applied TMS to left  
20  
21 100 tongue M1 and recorded MEPs from participants' tongue muscles during perception  
22  
23 101 of congruent and incongruent audiovisual syllables incorporating tongue- and/or lip-  
24  
25 102 related phonemes (visual and acoustic /ba/, /ga/, and /da/, visual /ba/ and acoustic /ga/,  
26  
27 103 and visual /ga/ and acoustic /ba/). Greater excitability of tongue M1 was measured for  
28  
29 104 syllables incorporating visual and/or acoustic tongue-related speech sounds, compared  
30  
31 105 to the presentation of lip-related speech sounds.  
32  
33  
34

35 106 Behaviourally, covert imitation can be measured using interference paradigms,  
36  
37 107 such as the Stimulus Response Compatibility (SRC) paradigm. SRC tasks were  
38  
39 108 originally mostly used to study covert imitation of manual actions (Brass, Wohlschläger,  
40  
41 109 Bekkering, & Prinz, 2000), but have also been used for speech stimuli. In a manual  
42  
43 110 SRC task, participants are instructed to perform a manual action in response to a  
44  
45 111 prompt (e.g., lift index finger when a written '1' appears, lift middle finger when '2'  
46  
47 112 appears). The prompt is presented superimposed on a distracter: An image or video of  
48  
49 113 a hand lifting the index or middle finger. When the prompt is presented in the  
50  
51 114 presence of a congruent distracter ('1' with a video of a lifting index finger),  
52  
53 115 participants are faster to perform the correct response than when the prompt is  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 116 presented together with an incongruent distracter ('l' with a video of a lifting middle  
4  
5 117 finger). For congruent distracters, it is assumed that action observation invokes motor  
6  
7 118 patterns for performing the prompted action, thus reducing response times (RTs). In  
8  
9 119 contrast, incongruent distracters result in competition between the activated motor  
10  
11 120 patterns and those required to produce the prompted response, leading to slower RTs.  
12  
13 121 A larger SRC effect, i.e., a larger RT difference between incongruent and congruent  
14  
15 122 pairs, indicates that motor mechanisms were more activated for the distracter. SRC  
16  
17 123 paradigms are thought to provide a fairly direct measure of the relative activation of  
18  
19 124 motor mechanisms and of covert imitation (Heyes, 2011).

20  
21  
22 125 In speech SRC paradigms (Galantucci, Fowler, & Goldstein, 2009; Jarick &  
23  
24 126 Jones, 2009; Kerzel & Bekkering, 2000; Roon & Gafos, 2015), the participant  
25  
26 127 produces a speech response following a prompt (e.g., *ba*) while ignoring a distracter  
27  
28 128 (e.g., a video of someone saying *da*). As reported for manual SRC studies, responses  
29  
30 129 to the prompt are slower for incongruent (*da*) than congruent (*ba*) distracters (Kerzel  
31  
32 130 & Bekkering, 2000). Kerzel & Bekkering used video-only distracter stimuli, and later  
33  
34 131 studies extended the use of the SRC paradigm to audio and audiovisual modalities.  
35  
36 132 Jarick & Jones ran the SRC task with video-only, audio-only and audiovisual  
37  
38 133 distracters. Participants were required to respond by either pressing a button or  
39  
40 134 speaking when seeing the prompt *ba* or *da*, in separate tasks. They measured the  
41  
42 135 largest covert imitation effects for their video-only condition, and the smallest effect  
43  
44 136 for the audio-only condition for the speech response condition. They also report no  
45  
46 137 covert imitation effects for manual responses (a pattern also reported in Galantucci et  
47  
48 138 al.), thus demonstrating that covert imitation is effector-specific.

49  
50  
51  
52 139 Converging evidence from fMRI, TMS and behavioural studies thus indicates  
53  
54 140 that observing visual, auditory, or audiovisual speech sounds results in covert  
55  
56  
57  
58  
59  
60

1  
2  
3 141 imitation. However, covert imitation effects for speech sounds have only been  
4  
5 142 demonstrated for a select class of speech sounds, i.e., for stop consonants, either in a  
6  
7 143 CV syllable or in isolation. It is not clear if observing vowels also invokes covert  
8  
9 144 imitation, and if these effects would be comparable in size with covert imitation  
10  
11 145 effects reported for consonants. A single fMRI study examined whether vowels are  
12  
13 146 somatotopically represented in articulatory M1 (Grabski et al., 2013). Grabski et al.  
14  
15 147 presented listeners with recordings of participants' own monophthongal French  
16  
17 148 vowels (/i y u e ø o ε œ ə/). These vowels varied in vowel height (close, mid-close and  
18  
19 149 mid-open), tongue position (front or back), and lip rounding (rounded or unrounded).  
20  
21 150 If vowel articulation is represented somatotopically as is the case for stop consonants,  
22  
23 151 it could be expected that tongue position and rounding could be linked to tongue and  
24  
25 152 lip M1 respectively, and vowel height to the jaw muscle M1 representation. However,  
26  
27 153 Grabski et al report no activation in M1 related to vowel perception and neural  
28  
29 154 responses linked to vowel perception were diffusely distributed across a network of  
30  
31 155 bilateral temporal, left prefrontal, and left parietal areas. Thus, to our knowledge, no  
32  
33 156 fMRI, TMS, or behavioural SRC study has demonstrated that observers covertly  
34  
35 157 imitate vowel stimuli.

36  
37  
38  
39 158 There is evidence that consonants and vowel are processed differently at neural  
40  
41 159 levels. Brain damage has been shown to impair consonant processing while  
42  
43 160 preserving vowel processing and vice versa (Caramazza, Chialant, Capasso, & Miceli,  
44  
45 161 2000). Moreover, electrical stimulation of the temporal cortex in patients with aphasia  
46  
47 162 impaired consonant discrimination but not vowel discrimination (Boatman, Hall,  
48  
49 163 Goldstein, Lesser, & Gordon, 1997; Boatman, Lesser, Hall, & Gordon, 1994). Results  
50  
51 164 from fMRI studies also suggest a difference in the neural processing of consonant and  
52  
53 165 vowel sounds (Seifritz et al., 2002). Using behavioural studies, further evidence was  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 166 provided for a dissociation in the roles vowels and consonants play, in speech  
4  
5 167 perception specifically. Several perceptual phenomena occurring for stop consonants,  
6  
7 168 such as categorical perception (Liberman, Harris, Hoffman, & Griffith, 1957) and  
8  
9 169 duplex perception (Liberman, Isenberg, & Rakerd, 1981), were found to not extend to  
10  
11 170 vowels (Gerrits & Schouten, 2004; Whalen & Liberman, 1996). Results from patient  
12  
13 171 studies, electrical stimulation experiments, fMRI studies, and behavioural studies thus  
14  
15 172 converge on the notion that consonants and vowels may be treated differently by the  
16  
17 173 speech processing system. It is important to establish whether covert imitation occurs  
18  
19 174 for stop consonants and for vowels, and if it does, whether there is a difference in the  
20  
21 175 size of covert imitation effects. If it is the case covert imitation only occurs for (stop)  
22  
23 176 consonants, and not for vowels, then this implies that listening to vowel sounds may  
24  
25 177 not result in automatic activation of articulatory motor plans required for generating  
26  
27 178 simulations during speech perception.  
28  
29  
30

31 179 The present study tested whether listeners covertly imitate vowels. Past studies  
32  
33 180 used CV syllables where place of articulation or voicing was contrasted between the  
34  
35 181 initial consonants, and the following vowel remained the same (Galantucci et al.,  
36  
37 182 2009; Jarick & Jones, 2009; Kerzel & Bekkering, 2000; Roon & Gafos, 2015). In our  
38  
39 183 CVC (consonant-vowel-consonant) stimuli the consonants remained the same (/h/ and  
40  
41 184 /d/), while the vowel was either /i/ as in *heed*) or /ʊ/ (as in *hood*). The vowels in *heed*  
42  
43 185 and *hood* were selected as they are produced with either spread (*heed*) or rounded lips  
44  
45 186 (*hood*) and can thus be distinguished visually.  
46  
47

48 187 Using vowels allows also for more detailed scrutiny of variation over time in the  
49  
50 188 covert imitation effect, as vowels are less transient than consonants. We therefore  
51  
52 189 presented the prompt at four time points (Stimulus Onset Asynchronies, SOAs) during  
53  
54 190 articulation. SOA manipulations were also used in Roon & Gafos, Kerzel &  
55  
56  
57  
58  
59  
60

1  
2  
3 191 Bekkering, and Galantucci et al. However, all three studies used CV stimuli, and  
4  
5 192 SOAs were restricted to a short time-span, i.e., between 100-300ms for Roon & Gafos  
6  
7 193 (100, 200, 300ms), between 0-500ms for Kerzel & Bekkering (0, 167, 333, 500 ms),  
8  
9 194 and between 0-495ms (0, 165, 330, 495ms) for Galantucci et al. The SOAs used in  
10  
11 195 past studies were spaced apart in equal intervals of the distracter video duration and  
12  
13 196 not linked to specific articulatory features, such as the onset or offset of articulation.  
14  
15 197 In the present study, we presented the prompts at four SOAs coinciding with the start  
16  
17 198 of the distracter (0ms, SOA1), the onset of visible articulation (335ms, SOA2), the  
18  
19 199 point where the auditory signal started and where the visual articulatory difference  
20  
21 200 between the two vowels was maximal (670ms, SOA3), and the point at which visible  
22  
23 201 articulation ceased for both vowels (1700ms, SOA4). We expected smaller covert  
24  
25 202 imitation effects for SOA1 compared to later SOAs, as no distracting articulatory  
26  
27 203 information was present at 0ms. Previous studies found smaller or no interference  
28  
29 204 effects when the SOA was set to the start of the trial. We included SOA2 and SOA4  
30  
31 205 to establish whether the covert imitation effect is larger at the beginning or the end of  
32  
33 206 the articulatory sequence, and SOA3 to establish if the covert imitation effect is  
34  
35 207 maximal when the visual difference between the two distracters is also maximal.

36  
37  
38  
39 208 Finally, it is currently unclear how distracter modality affects covert imitation of  
40  
41 209 vowels. A single previous study examined the effect of video, audio, and audiovisual  
42  
43 210 distracter stimuli on covert imitation for consonants (Jarick & Jones, 2009). However,  
44  
45 211 as Jarick & Jones presented the prompt at a single time point (100ms from the start of  
46  
47 212 the distracter stimulus), it remains unclear how modality affects covert imitation over  
48  
49 213 time. The four SOAs will thus also serve to establish if and how distracter modality  
50  
51 214 interacts with covert imitation over time.

215

## Experiment 1

## 216 Methods

217 An *a priori* power analysis (G\*Power 3.1.9.2, (Faul, Erdfelder, Lang, & Buchner,  
218 2007) for a between-group design with three groups and 240 observations per  
219 participant suggested a sample size of 66 participants (22 per group) with an type I  
220 error of  $p < 0.05$  and observed power of 80% for an expected effect size of 0.25. Sixty-  
221 six participants, 22 per group, (46F, 20M, mean 22.4y, SD 4.8y, range: 18-40y) took  
222 part. One male participant from the Audio group was excluded for not following task  
223 instructions. Participants were randomly assigned to three groups: Video (16F, 6M,  
224 mean 23.6y, SD 4.8y, range: 18-40y), Audio (12F, 11M, mean 23.1y, SD 3.7y, range:  
225 19-31y), and Audiovisual (18F, 4M, mean 20.6y, SD 4.1y, range: 18-28y). All were  
226 native speakers of British English, who reported normal or corrected to normal vision,  
227 normal hearing, and no (history of) dyslexia. The study was approved by UCL's  
228 Research Ethics Committee (#0599.001). Participants gave informed consent and  
229 received course credit or payment.

230 The distracter stimuli consisted of two videos of a female speaker saying *heed*  
231 or *hood* (Figure 1). The video stimuli were recorded by a 29-year-old female speaker  
232 of British English, with a Canon Lagria HF G30 video camera on a tripod. The video  
233 recordings were edited using iMovie on an Apple iMac, and scaled down in resolution  
234 from 1920×1090 to 1280×720 in .avi format. The prompt was a jpeg image with a  
235 resolution of 300dpi, 0.38×0.16cm (45×19 pixels), was presented on-screen at a size  
236 of 1.1×0.5cm, and consisted of either *heed* or *hood* printed in boldfaced Arial font on  
237 a black background. Font size was adjusted so that the lip movements remained  
238 highly visible while the prompt appeared centred on the mouth (Figure 1). The audio  
239 stimuli were recorded simultaneously with the video recordings, using a RODE NO1-  
240 A Condenser Microphone, a Focusrite Scarlett 2i4 USB Computer Audio Interface

1  
2  
3 241 pre-amplifier plugged into the sound card input of a Dell PC in a sound-attenuated  
4  
5 242 room at 44.1kHz with 16 bits. Audio recordings were amplitude normalized offline,  
6  
7 243 down-sampled to 22.050kHz, and scaled to 70dB SPL (Sound Pressure Level) using  
8  
9 244 Praat (Boersma & Weenink, 2003). The audio file for *hood* had a total duration of  
10  
11 245 977ms (/h/ segment: 137ms, /o/ 732ms, /d/ 108ms) and the audio file for *heed* also  
12  
13 246 had a total duration of 977ms (/h/ segment: 133ms, /o/ 734ms, /d/ 110ms). The video  
14  
15 247 files were muted using iMovie (9.0.9), and the video and audio files were combined in  
16  
17 248 Presentation when the trial was presented.

18  
19  
20 249 The experiment was conducted in a sound-attenuated and light-controlled booth.  
21  
22 250 The stimuli appeared on a PC monitor located 70cm away from the participant.  
23  
24 251 Stimuli were presented using Presentation (Neurobehavioral Systems). Audio was  
25  
26 252 played through Sennheiser HD25 SP-II headphones. Instructions were provided on-  
27  
28 253 screen. Participants were instructed to look out for the prompt and speak the prompt  
29  
30 254 aloud as fast as possible, ignoring the video in the background. Participants completed  
31  
32 255 16 familiarisation trials to ensure they performed the task as instructed and spoke at  
33  
34 256 appropriate loudness levels, while avoiding making any other sounds. The  
35  
36 257 experimenter left the room after the familiarisation session.

37  
38  
39 258 --- Figure 1 about here ---  
40

41  
42 259 Trials in the main experiment proceeded as follows. First, a black screen with a  
43  
44 260 fixation cross was presented for either 500, 750 or 1000ms (jitter, following Kerzel  
45  
46 261 and Bekkering). Next, a tone (500Hz, 200ms) was presented to signal the start of the  
47  
48 262 trial. In the Video condition, subsequently the video was presented with the sound  
49  
50 263 muted. In the Audio condition, a still image of the speaker with her mouth closed was  
51  
52 264 presented in the background, and the sound file started 670ms from the start of the  
53  
54 265 trial. In the Audiovisual condition, the video started playing at 0ms and the sound file  
55  
56  
57  
58  
59  
60

## Covert Vowel Imitation 12

1  
2  
3 266 started playing 670 after the start of the video. Note that audible articulation of vowels  
4  
5 267 in an /hVd/ context tends to follow visible articulation. The start time of the audio was  
6  
7 268 selected as initial pilot testing revealed this time point optimal for a natural effect and  
8  
9 269 this time point was placed approximately in between the points in time when the  
10  
11 270 audio started for the original *heed* and *hood* audiovisual recordings.

13 271 In all conditions, the prompt appeared superimposed over the lips of the speaker  
14  
15 272 for a duration of 200ms (Figure 1). The prompt was presented at four Stimulus Onset  
16  
17 273 Asynchronies (SOA); chosen to coincide with key points in the stimulus: 0ms (start of  
18  
19 274 the trial), 335ms (onset of visible articulation in the Video and Audiovisual  
20  
21 275 conditions), 670ms, (the start of the auditory signal in all three conditions), 1700ms  
22  
23 276 (end of visible articulation). The video started and ended with the speaker's lips  
24  
25 277 closed and no eye-blinks were present.

28  
29 278 Responses were recorded via a voice key in Presentation, using a Rode  
30  
31 279 microphone plugged into a Scarlett pre-amplifier connected to the PC's USB input,  
32  
33 280 from voice onset for 2500ms. Responses could be made from the start of the trial (i.e.,  
34  
35 281 the start of the video). RTs were measured from the onset of the prompt across for all  
36  
37 282 three groups. When no response had been detected after 2500ms from the start of the  
38  
39 283 video, participants received a *no response* warning. Stimulus lists were randomised  
40  
41 284 for each individual participant, and the same randomised stimulus lists were used  
42  
43 285 across successive participants in the three groups. The experiment lasted  
44  
45 286 approximately 40 minutes. Data, stimulus materials and program code can be found  
46  
47 287 on the Open Science Network, under the name *SRC\_Vowels* (<https://osf.io/sn396/>).

50 288 We first converted the raw error percentages per participant to rationalized  
51  
52 289 arcsine units, or RAUs, (Studebaker, 1985), as this procedure is customary for  
53  
54 290 proportional scales (e.g., (Adank, Evans, Stuart-Smith, & Scott, 2009). Transforming  
55  
56  
57  
58  
59  
60

1  
2  
3 291 the raw proportions to RAU ensures that the mean and variance of the data are  
4  
5 292 relatively uncorrelated and that the data are on a linear and additive scale (Studebaker,  
6  
7 293 1985). After transforming the error percentages data to RAUs, we performed a three-  
8  
9 294 factor repeated-measures ANOVA with the transformed error rates as the dependent  
10  
11 295 variable and with Prompt (Heed or Hood), Congruence (Congruent or Incongruent),  
12  
13 296 SOA (SOA1-4) as within-subject factors and listener group as a between-subject  
14  
15 297 factor for experiment 1 and Modality (Video, Audiovisual, Audio) as an additional  
16  
17 298 within-subject factor for experiment 2.

19  
20 299 The factors Congruence (Congruent, Incongruent), Prompt (*heed, hood*), SOA  
21  
22 300 (1-4), and Modality (Video, Audio, Audiovisual) were manipulated to explore  
23  
24 301 changes in the response times in milliseconds (RT), and analysed in a repeated-  
25  
26 302 measures ANOVA, controlled for non-sphericity (Huynh-Feldt), and post-hoc tests  
27  
28 303 were corrected for multiple comparisons (Bonferroni). RTs were log-transformed  
29  
30 304 before entered into the statistical analyses (Baayen, 2008). Only correct responses  
31  
32 305 were analysed. Errors were responses that were too early (<200ms) or late (>1000ms),  
33  
34 306 following Jarick & Jones, absent or partial responses, plus trials in which participants  
35  
36 307 produced incorrect or multiple prompts. It was determined whether a participant had  
37  
38 308 produced a correct or incorrect response by two phonetically trained listeners. Sound  
39  
40 309 file editing was conducted by a research assistant blind to the Congruence condition.

#### 41 42 43 44 310 Results

45  
46 311 Participants made 9.4% errors on average. Of the 15600 responses in total, 1460 were  
47  
48 312 classed as errors and excluded: 228 (1.5%) were missed responses, 1042 (6.7%) were  
49  
50 313 too early or too late, and in 190 (1.2%) cases participants produced the wrong prompt.  
51  
52 314 The analysis of the errors showed main effects of Prompt and SOA, and significant  
53  
54 315 interactions for Prompt×Congruence, Prompt×SOA (see Table A in Supplementary  
55  
56  
57  
58  
59  
60

## Covert Vowel Imitation 14

1  
2  
3 316 Materials). Analysis of the errors showed that participants made more errors for *heed*  
4  
5 317 (10%) than *hood* (8%). Participants made more errors for SOA1 (19%) than for the  
6  
7 318 other three SOAs (SOA2: 8%, SOA3: 7%, SOA4: 4%). Participants also made  
8  
9 319 significantly more errors for congruent (12%) than incongruent (9%) pairs for *heed*,  
10  
11 320 but not *hood* (8% congruent and 9% incongruent). Participants also made more errors  
12  
13 321 for SOA1 for *heed* (22%) than *hood* (16%). No Congruence effects were found.

14  
15 322 The analysis of the RTs included only correct responses. Main effects were  
16  
17 323 found for Prompt, Congruence, SOA, and the following interactions: SOA×Modality,  
18  
19 324 Prompt×Congruence, Prompt×SOA, and Congruence×SOA. Participants responded  
20  
21 325 overall slower for *heed* than for *hood* prompts. The RTs showed an overall covert  
22  
23 326 imitation effect, as RT were faster for congruent than incongruent trials (Figure 2,  
24  
25 327 Table I). As predicted, covert imitation effects differed per SOA and were largest for  
26  
27 328 SOA3, and no covert imitation effect was found for SOA1. RTs were faster for later  
28  
29 329 consecutive time points, except between SOA2 and SOA3. The SOA×Modality  
30  
31 330 interaction was linked to slower responses for the Video than for the Audiovisual  
32  
33 331 group, for SOA4 only. The Prompt×Congruence interaction was related to larger  
34  
35 332 covert imitation effects for *heed* than *hood*. *Heed* responses were slower than *hood*  
36  
37 333 responses at SOAs 2 and 4. An analysis of difference scores (incongruent minus  
38  
39 334 congruent RTs) showed that covert imitation effects were found for heed across all  
40  
41 335 three groups, but for hood these effects were found for Video and Audio groups only.

42  
43  
44 336 --- Insert Table I and Figure 2 about here ---

45  
46  
47 337 In conclusion, the results of Experiment 1 showed a clear main covert imitation effect  
48  
49 338 for the response times only. Congruent trials were associated with faster responses  
50  
51 339 than incongruent trials across all three modalities. These results replicated earlier  
52  
53 340 work showing effects of congruence for consonants in CV syllables (Jarick & Jones,  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 341 2009, Kerzel & Bekkering, 2000) and extended these effects to vowels in CVC  
4  
5 342 syllables. However, the effects measured here were smaller than those for CV  
6  
7 343 syllables (13ms across all SOAs versus ~35ms for Experiment 1 in Kerzel and  
8  
9 344 Bekkering, averaged across both prompts). Jarick and Jones report smaller covert  
10  
11 345 imitation effects for Audio than their Video and Audiovisual conditions. However,  
12  
13 346 due to the between-group design, employed in Experiment 1, it was not feasible to  
14  
15 347 directly establish the extent to which participants changed their responses under  
16  
17 348 different modalities, as was done in Jarick and Jones (2009), who used a within-  
18  
19 349 subject design. Note that we chose to use a between-group design in Experiment 1 to  
20  
21 350 reduce the experimental duration (40 minutes) while optimising the number of trials  
22  
23 351 per participant (240 per modality), and to avoid potential order effects from switching  
24  
25 352 from one modality to the next. Experiment 2 used a within-group design, in which all  
26  
27 353 participants completed the task for all three modalities in separate blocks to further  
28  
29 354 explore the effect of modality on covert imitation.  
30  
31

### 32 33 355 Experiment 2

34  
35 356 Experiment 2 aimed to independently replicate effects found in Experiment 1 using a  
36  
37 357 within-group design in which all participants completed the task for the three  
38  
39 358 modalities in separate blocks.  
40

### 41 359 Methods

42  
43  
44 360 An *a priori* power analysis for a within-group design with 360 observations per  
45  
46 361 participant suggested a sample size of 24 with a type I error of  $p < 0.05$  and observed  
47  
48 362 power of 80%, for an expected effect size of 0.25. Twenty-four female participants  
49  
50 363 (19.0y, SD 1.4y, range: 18-23y) took part in Experiment 2. None of these participants  
51  
52 364 took part in Experiment 1. All participants were native speakers of British English,  
53  
54 365 who reported normal or corrected-to-normal vision, normal hearing, and no (history  
55  
56  
57  
58  
59  
60



1  
2  
3 366 of) dyslexia. Video data for one participant was missing due to a technical error.  
4  
5 367 Materials, task, and general procedure were similar to Experiment 1, except that  
6  
7 368 participants completed the three conditions Video, Audio, and Audiovisual (120 trials  
8  
9 369 each) in a counterbalanced order: participant 1 first completed the Video condition,  
10  
11 370 followed by the Audio and Audiovisual conditions. The order for the next participant  
12  
13 371 was Audiovisual, Video, Audio, and the next participant completed the experiment in  
14  
15 372 the order: Audio, Audiovisual, Video, in a single session lasting 60 minutes. The  
16  
17 373 procedure was the same for all other participants. Stimulus lists were randomised per  
18  
19 374 participant per condition, and the same randomised list was used across the three  
20  
21 375 conditions per participant, per the procedure used in Experiment 1.  
22  
23

#### 376 Results

24  
25  
26 377 Participants made 8.5% errors overall. Of the 8520 responses, 728 were classed as  
27  
28 378 errors and excluded: 164 (1.9%) were missed responses, 417 (4.9%) were too early or  
29  
30 379 too late, and in 147 (1.7%) cases participants produced the wrong prompt. Main  
31  
32 380 effects were found for Prompt, Congruence, SOA, plus the Prompt×SOA interaction  
33  
34 381 (see Table B in Supplementary Materials). Participants made more errors for *heed*  
35  
36 382 (10%) than *hood* (7%). Participants made more errors for SOA1 (19%) than for the  
37  
38 383 other SOAs (SOA2: 5%, SOA3: 5%, SOA4: 4%). Participants made fewer errors for  
39  
40 384 congruent (8%) than incongruent (9%) pairs. Participants also made more errors for  
41  
42 385 SOA1 for *heed* (22%) than for *hood* (16%).  
43  
44  
45

46 386 The analysis of the RTs included only correct responses. Main effects were  
47  
48 387 found for Congruence and SOA, plus the interactions Modality×SOA, Prompt×SOA,  
49  
50 388 Congruence×SOA, and Prompt×Congruence×SOA interactions. An overall covert  
51  
52 389 imitation effect was again found, as participants responded faster for congruent than  
53  
54 390 for incongruent pairs. However, covert imitation effects were only found for SOA2  
55  
56  
57  
58  
59  
60

1  
2  
3 391 and SOA3, as the difference between incongruent and congruent trials was not  
4  
5 392 significantly different for SOA1 and SOA4. Participants again responded overall  
6  
7 393 faster for later consecutive SOAs. Modality×SOA interactions were rather  
8  
9 394 inconsistent. Faster responses were recorded for Audio SOA2 than Audiovisual  
10  
11 395 SOA2, faster responses were found for Video SOA3 than Audio SOA3, and faster  
12  
13 396 responses were found for Audio SOA4 than Video SOA4. Slower *heed* responses  
14  
15 397 were reported for SOA1 and SOA2, but not for SOA3 and SOA4. No follow-up tests  
16  
17 398 survived correction for the Prompt×Congruence×SOA interaction.

19  
20 399 In conclusion, the results of experiment 2 replicated the covert imitation effect  
21  
22 400 for vowels reported for Experiment 1 for the response times and also reported a small  
23  
24 401 covert imitation effect for the errors, which was not reported for experiment 1. The  
25  
26 402 results did not reveal an effect of distracter modality on covert imitation, even when  
27  
28 403 participants performed the SRC task for all three modalities. Experiment 2 further  
29  
30 404 showed a replication of the interaction between SOA and congruence, covert imitation  
31  
32 405 was most prominent at SOA2 and SOA3.

33  
34  
35 406 --- *Insert Table II about here* ---

#### 36 37 407 General discussion

38  
39 408 This study aimed to establish whether observers covertly imitate vowel stimuli, how  
40  
41 409 covert imitation varies over time, and how distracter modality affects covert imitation.  
42  
43 410 We conducted two experiments in which participants produced vocal responses to a  
44  
45 411 CVC prompt in the presence of a background distracter in Video, Audio, or  
46  
47 412 Audiovisual modalities. A clear covert imitation effect was found on the response  
48  
49 413 times in both experiments; participants showed faster responses for congruent than  
50  
51 414 incongruent trials. Our study thus replicated earlier work that showed covert imitation  
52  
53 415 effects on consonants (Galantucci et al., 2009; Jarick & Jones, 2009; Kerzel &  
54  
55  
56  
57  
58  
59  
60

## Covert Vowel Imitation 18

1  
2  
3 416 Bekkering, 2000; Roon & Gafos, 2015) and extended these effects to vowels. We  
4  
5 417 found covert imitation effects of 13ms for Experiment 1 and 7ms for Experiment 2,  
6  
7 418 collapsed over the four SOAs. Kerzel & Bekkering report covert imitation effects of  
8  
9 419 35ms for their Experiment 1 and Galantucci et al. report an effect of 28ms for their  
10  
11 420 Experiment 2. Covert imitation effects for vowels seem to be overall smaller than  
12  
13 421 those reported for consonants. Observing incongruent vowel articulation may lead to  
14  
15 422 less activation of articulatory motor patterns compared to observing incongruent stop  
16  
17 423 consonant articulation. In the visual domain, the stop consonants generally used in  
18  
19 424 SRC paradigms differ in the active articulator, namely lips or tongue, while our vowel  
20  
21 425 stimuli differed only in the use of the primary articulator (lips rounded or unrounded).  
22  
23 426 A distracter employing a different effector could result in greater, more widespread,  
24  
25 427 activation of articulatory patterns than a distracter changing the use of a single  
26  
27 428 effector. Alternatively, observing a congruent vowel distracter may not facilitate the  
28  
29 429 production of the correct response as much as is the case for stop consonants, again  
30  
31 430 due to differences in articulation between the two classes of speech sounds. Follow-up  
32  
33 431 studies could address the issue of articulatory complexity, for instance, by exploring  
34  
35 432 somatotopy of perceived vowel stimuli using TMS, specifically by measuring MEPs  
36  
37 433 from lip and tongue muscles. Previous work has demonstrated somatotopy in tongue  
38  
39 434 M1 (Sato et al., 2009) and lip M1 (Nuttall et al., 2017; Nuttall et al., 2016) congruent  
40  
41 435 with the primary articulator of the observed speech sound. Somatotopy in TMS  
42  
43 436 speech perception studies refers to the notion that specific parts of articulatory M1  
44  
45 437 become active, or show relative facilitation, when listening to speech sounds  
46  
47 438 articulated using a congruent articulator (so lip M1 becomes relatively facilitated for  
48  
49 439 lip-produced sounds such as /t/ or /d/). By comparing relative facilitation of lip M1  
50  
51 440 and tongue M1 while observing lip-articulated (/p/), tongue-articulated (/t/) sounds  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 441 with unrounded (/i:/) and rounded vowels (/o/, or /y/ for languages other than British  
4  
5 442 English, e.g., Dutch), it could be established if greater differences in facilitation occur  
6  
7 443 for lip or tongue sounds.

8  
9 444 Modality did not directly affect covert imitation, as no evidence was found of an  
10  
11 445 interaction between congruence, modality, and SOA in either experiment. It must be  
12  
13 446 concluded that Modality effects on covert imitation seem to be moderate or small for  
14  
15 447 vowels, replicating and extending past findings by Jarick & Jones for consonants.

16  
17 448 Covert imitation effects were largest for SOA3 (26ms) in Experiment 1, and  
18  
19 449 SOA2 (20ms) and SOA3 (23ms) in Experiment 2. These results illustrate that covert  
20  
21 450 imitation is maximal for the time point (670ms) at which the difference between the  
22  
23 451 two distracters is maximal visually (in the Video and Audiovisual conditions) and/or  
24  
25 452 when the audio starts playing (in Audio and Audiovisual conditions). The absence of  
26  
27 453 a covert imitation effect at SOA1 (0ms) in either experiment shows that distracting  
28  
29 454 audio and/or visual distracter information was required to elicit covert imitation  
30  
31 455 effects. Participants also responded faster for later onsets in both experiments; a result  
32  
33 456 also reported by Kerzel & Bekkering and Galantucci et al. Interference effects also  
34  
35 457 differed across SOAs. For Experiment 1, interference effects were largest for SOA3  
36  
37 458 (26ms), while for Experiment 2 these were largest for SOA2 (22ms) and SOA3  
38  
39 459 (14ms) and no interference effect was found at SOA1 in either experiment. Note that  
40  
41 460 SOA3 (670ms) was chosen to coincide with the moment at which the audio signal  
42  
43 461 started in the Audio and Audiovisual modalities and also the point at which the visual  
44  
45 462 difference between the two distracters was maximal (spread vs. rounded lips).

46  
47 463 Covert imitation effects differed depending on the stimulus prompt; larger  
48  
49 464 effects were found for *heed* than *hood*, in analogy with Kerzel & Bekkering, who  
50  
51 465 report a trend towards smaller effects for /ba/ than /da/ prompts. Larger interference  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Covert Vowel Imitation 20

1  
2  
3 466 effects for *heed* imply more interference from *hood* and vice versa. Larger effects for  
4  
5 467 *heed* (with *hood* distracter) showed that a distracter with rounded lips results in more  
6  
7 468 covert imitation than the other way around. Alternatively, lip rounding might be more  
8  
9 469 visually salient than lip spreading, and as a result might subsequently lead to more  
10  
11 470 activation of motor substrates. Alternatively, it seems possible that the conflict  
12  
13 471 between prompt and distracter resulted in a perceived fusion between the distracter  
14  
15 472 and prompt. Results from previous work has shown that observing conflicting  
16  
17 473 audiovisual information can lead to perceived vowel fusions (Traunmüller &  
18  
19 474 Öhrström, 2007). Traunmüller & Öhrström found that acoustic /geg/ dubbed onto  
20  
21 475 visually presented /gyg/ was predominantly perceived as /gøg/. In Traunmüller &  
22  
23 476 Öhrström's study visual lip-rounding affected the auditory perception of spreading  
24  
25 477 more than the degree to which visual perception of lip-spreading affected the auditory  
26  
27 478 perception of lip rounding. It seems possible that similar asymmetric partial fusions  
28  
29 479 occur for conflicts between speech production and simultaneously presented  
30  
31 480 distracters and that such asymmetric partial fusions can explain the difference in how  
32  
33 481 participants perceived our incongruent prompt-distracter pairings. Finally, participants  
34  
35 482 could have found the video that involved lip-spreading (heed) more visually salient  
36  
37 483 than the lip rounding video (hood). Potential effects of the relative salience of lip-  
38  
39 484 spreading versus lip-rounding warrants further investigation in future studies.

40  
41  
42  
43 485 For both experiments, on average 9% errors were found. Participants made  
44  
45 486 more errors for *heed* than for *hood* prompts in both experiments. Error percentages  
46  
47 487 were higher than those reported in previous work (Galantucci et al., 2009; Jarick &  
48  
49 488 Jones, 2009; Kerzel & Bekkering, 2000) (~1-3% for across all three studies). Close  
50  
51 489 inspection of the results showed that, for both experiments, most errors were due to  
52  
53 490 participants failing to respond, or failing to respond on time, for SOA1 (0ms),  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 491 possibly as a result of missing the prompt altogether for this SOA. Jarick & Jones did  
4  
5 492 not include trials in which the prompt was presented at the very start of the trial; the  
6  
7 493 prompt was presented around 100ms into the trial duration, so participants were more  
8  
9 494 likely to not miss the prompt. Kerzel & Bekkering and Galantucci et al. showed the  
10  
11 495 prompt at 0ms, but do not provide detailed information on how errors were distributed  
12  
13 496 across SOAs. Finally, it is unclear whether error percentages in previous work  
14  
15 497 included incorrect responses (i.e., the wrong prompt) or whether they only included  
16  
17 498 early or late or missed responses (e.g., Experiments 2 and 3 in Galantucci et al.).

19  
20 499 In conclusion, our study provides the first experimental evidence of covert  
21  
22 500 imitation for vowels. Covert imitation effect for vowels were smaller than those  
23  
24 501 previously reported for stop consonants, which may be due to less activation of  
25  
26 502 articulatory motor plans during perception of vowel stimuli. Future studies could  
27  
28 503 explore the possibility raised by our results that the dampened covert imitation effects  
29  
30 504 for vowels compared to previously reported effects for consonants could be due to  
31  
32 505 greater similarity between vowel stimuli than between contrastive stop consonants.  
33  
34 506 Covert imitation of vowels is not modulated by stimulus modality, and appears linked  
35  
36 507 to differences between distracter and prompt. We replicated this finding in two  
37  
38 508 experiments. Our study thus supports simulation theories of speech perception, by  
39  
40 509 clearly showing that perceiving vowels links to activation of speech motor  
41  
42 510 mechanisms. Current theories (Pickering & Garrod, 2013; Wilson & Knoblich, 2005)  
43  
44 511 predict that observing an action activates articulatory plans congruent with the  
45  
46 512 observed action in a somatotopic fashion, based on the results of studies mostly using  
47  
48 513 stop consonants. Past work has so far not demonstrated that vowel stimuli are  
49  
50 514 processed in a similar somatotopic manner (Grabski et al., 2013). The lack of  
51  
52 515 evidence of somatotopic processing for vowels in combination with our reported  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 516 smaller covert activation effects imply that the type of articulatory plan activated  
4  
5 517 during perception differs for different classes of speech sounds.  
6

7 518 Acknowledgements

9 519 This work was supported by the BIAL Foundation under grant number 267/14 to PA.

11 520 We thank Dan Kennedy-Higgins and Flavia Bojescu for assistance in data collection.  
12  
13  
14 521

15 522 References

17  
18 523 Adank, P., Evans, B. G., Stuart-Smith, J., & Scott, S. K. (2009). Comprehension of

20 524 familiar and unfamiliar native accents under adverse listening conditions.

22 525 *Journal of Experimental Psychology Human Perception and Performance*,

24 526 35(2), 520-529. doi:10.1037/a0013552

26 527 Adank, P., Nuttall, H. E., & Kennedy-Higgins, D. (2016). Transcranial Magnetic

28 528 Stimulation (TMS) and Motor Evoked Potentials (MEPs) in Speech

30 529 Perception Research. *Language, Cognition & Neuroscience*, 1-10.

32 530 doi:10.1080/23273798.2016.1257816

34 531 Baayen, R. H. (2008). Data sets and functions with "Analyzing Linguistic Data: A

36 532 practical introduction to statistics". (Version R package version 0.953.).

38 533 Boatman, D., Hall, C., Goldstein, M. H., Lesser, R., & Gordon, B. (1997).

40 534 Neuroperceptual differences in consonant and vowel discrimination: as

42 535 revealed by direct cortical electrical interference. *Cortex*, 33(10.1016/S0010-

44 536 9452(97)80006-8), 83-98.

46 537 Boatman, D., Lesser, R., Hall, C., & Gordon, B. (1994). Auditory perception of

48 538 segmental features: a functional neuroanatomic study. *Journal of*

50 539 *Neurolinguistics*, 8(225-234). doi:10.1016/0911-6044(94)90028-0  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 540 Boersma, P., & Weenink, D. (Producer). (2003). Praat: doing phonetics by computer.  
4  
5 541 Retrieved from <http://www.praat.org>  
6  
7 542 Brass, M., Wohlsläger, A., Bekkering, H., & Prinz, W. (2000). Compatibility between  
8  
9 543 observed and executed finger movements: comparing symbolic, spatial and  
10  
11 544 imitative cues. *Brain and Cognition*, *44*, 124-143. doi:10.1006/brcg.2000.1225  
12  
13 545 Buccino, G., Lui, F., Canessa, N., Patteri, I., Lagravinese, G., Benuzzi, F., . . .  
14  
15 546 Rizzolatti, G. (2004). Neural circuits involved in the recognition of actions  
16  
17 547 performed by nonconspecifics: An fMRI study. *Journal of Cognitive*  
18  
19 548 *Neuroscience*, *16*(1), 114-126. doi:10.1162/089892904322755601  
20  
21  
22 549 Caramazza, A., Chialant, D., Capasso, R., & Miceli, G. (2000). Separable processing  
23  
24 550 of consonants and vowels. *Nature*, *403*(6768), 428-430.  
25  
26 551 doi:10.1038/35000206  
27  
28  
29 552 D'Ausilio, A., Pulvermuller, F., Salmas, P., Bufalari, I., Begliomini, C., & Fadiga, L.  
30  
31 553 (2009). The motor somatotopy of speech perception. *Current Biology*, *19*(5),  
32  
33 554 381-385. doi:10.1016/j.cub.2009.01.017  
34  
35 555 Fadiga, L., Craighero, L., Buccino, G., & Rizzolatti, G. (2002). Speech listening  
36  
37 556 specifically modulates the excitability of tongue muscles: a TMS study.  
38  
39 557 *European Journal of Neuroscience*, *15*(2), 399-402  
40  
41  
42 558 Faul, F., Erdfelder, E., Lang, A. G., & Buchner, A. (2007). G\*Power 3: A flexible  
43  
44 559 statistical power analysis program for the social, behavioral, and biomedical  
45  
46 560 sciences. *Behavior Research Methods*, *39*, 175-191. doi:10.3758/BF03193146  
47  
48 561 Galantucci, B., Fowler, C. A., & Goldstein, L. (2009). Perceptuomotor compatibility  
49  
50 562 effects in speech. *Attention, Perception & Psychophysics*, *71*(5), 1138–1149.  
51  
52 563 doi:10.3758/APP.71.5.1138  
53  
54  
55  
56  
57  
58  
59  
60



- 1  
2  
3 564 Gerrits, E., & Schouten, M. E. H. (2004). Categorical perception depends on the  
4  
5 565 discrimination task. *Perception and Psychophysics*, 66(3), 363-376.  
6  
7 566 doi:10.3758/BF03194885  
8  
9 567 Grabski, K., Schwartz, J. L. K., Lamalle, L., Vilain, C., Vallée, N., Baciú, M., . . .  
10  
11 568 Sato, M. (2013). Shared and distinct neural correlates of vowel perception and  
12  
13 569 production. *Journal of Neurolinguistics*, 26, 384-408.  
14  
15 570 doi:10.1016/j.jneuroling.2012.11.003  
16  
17  
18 571 Heyes, C. (2011). Automatic Imitation. *Psychological Bulletin*, 137(3), 463-483.  
19  
20 572 doi:10.1037/a0022288  
21  
22 573 Jarick, M., & Jones, J. A. (2009). Effects of seeing and hearing speech on speech  
23  
24 574 production: a response time study. *Experimental Brain Research*, 195, 175-  
25  
26 575 182. doi:10.1007/s00221-009-1765-x  
27  
28  
29 576 Kawato, M. (1999). Internal models for motor control and trajectory planning.  
30  
31 577 *Opinion in Neurobiology*, 9, 718-727. doi:10.1016/S0959-4388(99)00028-8  
32  
33 578 Kerzel, D., & Bekkering, H. (2000). Motor activation from visible speech: Evidence  
34  
35 579 from stimulus response compatibility. *Journal of Experimental Psychology:*  
36  
37 580 *Human Perception and Performance*, 26, 634-647. doi:10.10371/0096-  
38  
39 581 1523.26.2.634  
40  
41  
42 582 Liberman, A. M., Harris, K., Hoffman, H. S., & Griffith, B. (1957). The  
43  
44 583 discrimination of speech sounds within and across phoneme boundaries.  
45  
46 584 *Journal of Experimental Psychology*, 54, 358-368  
47  
48 585 Liberman, A. M., Isenberg, D., & Rakerd, B. (1981). Duplex perception of cues for  
49  
50 586 stop consonants: Evidence for a phonetic mode. *Attention, Perception, &*  
51  
52 587 *Psychophysics*, 30(2), 133-143  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 588 Möttönen, R., & Watkins, K. E. (2009). Motor representations of articulators  
4  
5 589 contribute to categorical perception of speech sounds. *Journal of*  
6  
7 590 *Neuroscience*, 5(29), 9819-9825. doi:10.1523/JNEUROSCI.6018-08.2009  
8  
9 591 Nuttall, H. E., Kennedy-Higgins, D., Devlin, J. T., & Adank, P. (2017). The role of  
10  
11 592 hearing ability and speech distortion in the facilitation of articulatory motor  
12  
13 593 cortex. *Neuropsychologia*, 94(8), 13-22.  
14  
15 594 doi:10.1016/j.neuropsychologia.2016.11.016  
16  
17 595 Nuttall, H. E., Kennedy-Higgins, D., Hogan, J., Devlin, J. T., & Adank, P. (2016).  
18  
19 596 The effect of speech distortion on the excitability of articulatory motor cortex  
20  
21 597 *NeuroImage*, 128, 218-226. doi:10.1016/j.neuroimage.2015.12.038  
22  
23 598 Pickering, M. J., & Garrod, S. (2013). An integrated theory of language production  
24  
25 599 and comprehension. *Behavioral and Brain Sciences*, 36(4), 329-347.  
26  
27 600 doi:10.1017/S0140525X12001495  
28  
29 601 Pulvermüller, F., Huss, M., Kherif, F., Moscoso del Prado Martin, F., Hauk, O., &  
30  
31 602 Shtyrov, Y. (2006). Motor cortex maps articulatory features of speech sounds.  
32  
33 603 *Proceedings of the National Academy of Sciences of the United States of*  
34  
35 604 *America*, 103(20), 7865-7870. doi:10.1073/pnas.0509989103  
36  
37 605 Roon, K. D., & Gafos, A. I. (2015). Perceptuo-motor effects of response-distractor  
38  
39 606 compatibility in speech: beyond phonemic identity. *Psychonomic Bulletin &*  
40  
41 607 *Review*, 22(1), 242-250. doi:10.3758/s13423-014-0666-6  
42  
43 608 Sato, M., Buccino, G., Gentilucci, M., & Cattaneo, M. (2009). On the tip of the  
44  
45 609 tongue: modulation of the primary motor cortex during audio - visual speech  
46  
47 610 perception. *Speech Communication*, 52(6), 533-541.  
48  
49 611 doi:10.1016/j.bandl.2009.03.002  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Covert Vowel Imitation 26

- 1  
2  
3 612 Seifritz, E., Esposito, F., Hennel, F., Mustovic, H., Neuhoff, J. G., Bilecen, D., . . . Di  
4  
5 613 Salle, F. (2002). Spatiotemporal pattern of neural processing in the human  
6  
7 614 auditory cortex. *Science*, *297*, 1706-1708. doi:10.1126/science.1074355  
8  
9 615 Studebaker, G. A. (1985). A "rationalized" arcsine transform. *Journal of Speech,*  
10  
11 616 *Language, and Hearing Research*, *28*, 455-462  
12  
13 617 Traunmüller, H., & Öhrström, N. (2007). Audiovisual perception of openness and lip  
14  
15 618 rounding in front vowels. *Journal of Phonetics*, *35*, 244–258.  
16  
17 619 doi:10.1016/j.wocn.2006.03.002  
18  
19 620 Watkins, K. E., Strafella, A. P., & Paus, T. (2003). Seeing and hearing speech excites  
20  
21 621 the motor system involved in speech production. *Neuropsychologia*, *41*(8),  
22  
23 622 989-994. doi:10.1016/S0028-3932(02)00316-0  
24  
25 623 Whalen, D., & Liberman, A. M. (1996). Limits on phonetic integration in duplex  
26  
27 624 perception. *Perception and Psychophysics*, *58*(6), 857-870.  
28  
29 625 Wilson, M., & Knoblich, G. (2005). The case for motor involvement in perceiving  
30  
31 626 conspecifics. *Psychological Bulletin*, *131*, 460-473. doi:10.1037/0033-  
32  
33 627 2909.131.3.460  
34  
35 628 Wilson, S. M., Saygin, A. P., Sereno, M. I., & Iacoboni, M. (2004). Listening to  
36  
37 629 speech activates motor areas involved in speech production. *Nature*  
38  
39 630 *Neuroscience*, *7*, 701-702. doi:doi:10.1038/nm1263  
40  
41  
42  
43  
44 631  
45  
46 632  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 633 Figures and Tables

634

635 Table I. Averages plus standard deviations “( )” for % error and response times in ms

636 for congruent and incongruent stimulus pairs, per prompt, per Stimulus Onset

637 Asynchrony (SOA), and modality, for Experiment 1.

		ERRORS				RESPONSE TIMES		
		Video	Audio	Audiovisual	Video	Audio	Audiovisual	
<i>Heed</i>	<i>Congruent</i>	SOA1	21 (41)	26 (44)	25 (44)	648 (123)	648 (123)	606 (140)
		SOA2	5 (23)	13 (33)	10 (30)	590 (115)	590 (115)	537 (130)
		SOA3	4 (20)	12 (33)	5 (22)	534 (103)	534 (103)	558 (144)
		SOA4	4 (19)	6 (24)	7 (26)	535 (85)	535 (85)	498 (109)
	<i>Incongruent</i>	SOA1	19 (39)	21 (41)	22 (42)	660 (131)	660 (131)	604 (136)
		SOA2	6 (23)	10 (30)	4 (20)	608 (106)	608 (106)	537 (145)
		SOA3	2 (15)	10 (29)	5 (23)	579 (111)	579 (111)	578 (130)
		SOA4	4 (19)	4 (19)	6 (24)	555 (75)	555 (75)	501 (100)
<i>Hood</i>	<i>Congruent</i>	SOA1	14 (34)	17 (38)	15 (36)	655 (134)	655 (134)	600 (119)
		SOA2	5 (21)	10 (30)	5 (23)	575 (110)	575 (110)	534 (129)
		SOA3	4 (19)	10 (31)	5 (23)	553 (107)	553 (107)	553 (139)
		SOA4	2 (15)	4 (20)	3 (17)	524 (72)	524 (72)	492 (139)
	<i>Incongruent</i>	SOA1	16 (37)	19 (39)	15 (36)	635 (125)	635 (125)	607 (130)
		SOA2	5 (23)	9 (29)	9 (28)	590 (119)	590 (119)	529 (138)
		SOA3	3 (16)	11 (31)	8 (28)	567 (102)	567 (102)	587 (149)
		SOA4	2 (13)	4 (19)	3 (18)	537 (80)	537 (80)	498 (104)

638

639

## Covert Vowel Imitation 28

640 Table II. Averages plus standard deviations “( )” for response times in milliseconds  
 641 for congruent and incongruent stimulus pairs, per prompt, per Stimulus Onset  
 642 Asynchrony (SOA), and modality, for Experiment 2.

		ERRORS				RESPONSE TIMES		
		Video	Audio	Audiovisual	Video	Audio	Audiovisual	
<i>Heed</i>	<i>Congruent</i>	SOA1	20 (40)	24 (43)	20 (40)	20 (40)	671 (130)	670 (131)
		SOA2	4 (21)	6 (24)	7 (26)	4 (21)	582 (118)	608 (130)
		SOA3	3 (18)	5 (23)	9 (29)	3 (18)	537 (107)	561 (133)
		SOA4	6 (24)	2 (15)	5 (23)	6 (24)	529 (86)	513 (96)
	<i>Incongruent</i>	SOA1	28 (45)	25 (44)	28 (45)	28 (45)	661 (138)	685 (132)
		SOA2	5 (23)	6 (24)	4 (20)	5 (23)	618 (123)	632 (140)
		SOA3	2 (15)	7 (26)	8 (27)	2 (15)	575 (112)	593 (120)
		SOA4	5 (23)	6 (23)	4 (20)	5 (23)	529 (83)	524 (97)
<i>Hood</i>	<i>Congruent</i>	SOA1	14 (35)	11 (31)	15 (36)	14 (35)	639 (126)	652 (120)
		SOA2	5 (22)	4 (19)	7 (25)	5 (22)	590 (138)	604 (137)
		SOA3	2 (15)	2 (12)	7 (25)	2 (15)	569 (142)	581 (125)
		SOA4	3 (17)	4 (20)	4 (19)	3 (17)	499 (95)	519 (84)
	<i>Incongruent</i>	SOA1	16 (37)	17 (38)	16 (37)	16 (37)	635 (129)	659 (132)
		SOA2	2 (15)	4 (19)	10 (30)	2 (15)	573 (135)	612 (133)
		SOA3	2 (13)	7 (26)	9 (28)	2 (13)	589 (136)	586 (126)
		SOA4	2 (15)	2 (14)	6 (23)	2 (15)	507 (87)	511 (91)

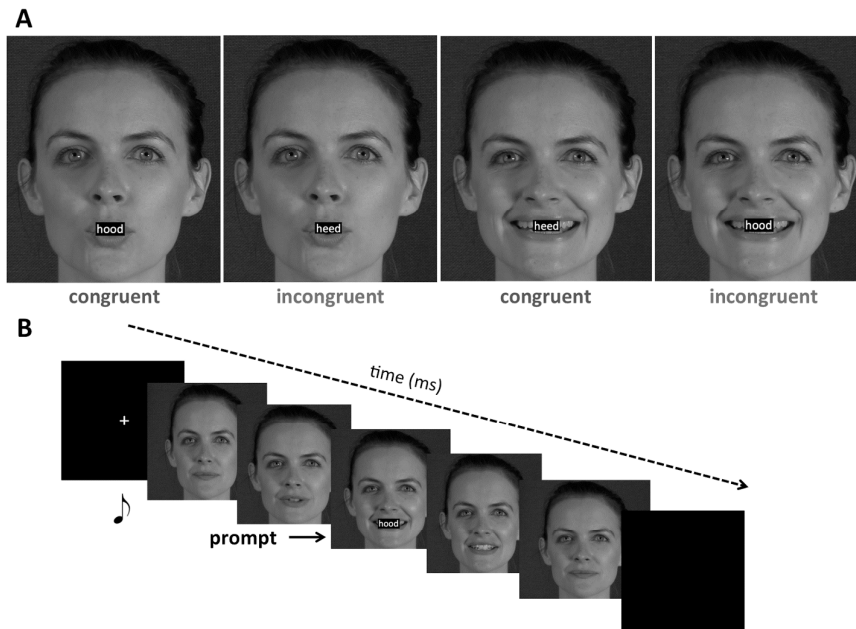
643

644

645

646

647



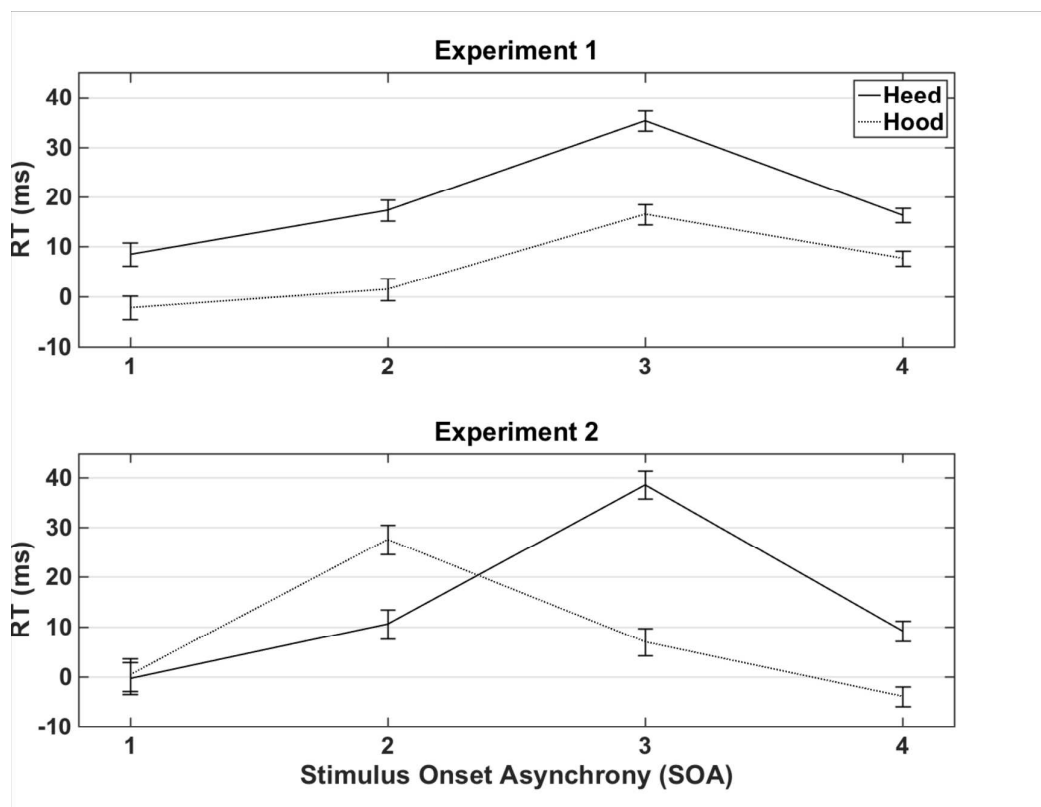
648

649 *Figure 1.* A: congruent trial for *hood* prompt, incongruent trial for *hood* prompt,650 congruent trial for *heed* prompt, incongruent trial for *heed* prompt. B. Example of the651 timeline of an incongruent stimulus pair with *hood* prompt and *heed* distracter.

652

653

654



655  
 656 *Figure 2.* Difference scores in milliseconds (incongruent minus congruent pairs)  
 657 pooled across the Video, Audio, and Audiovisual conditions, for each SOA and  
 658 separated by prompt, error bars represent one standard error. Top: Experiment 1, B:  
 659 Experiment 2.

660

661

662

663

664

665

666

667

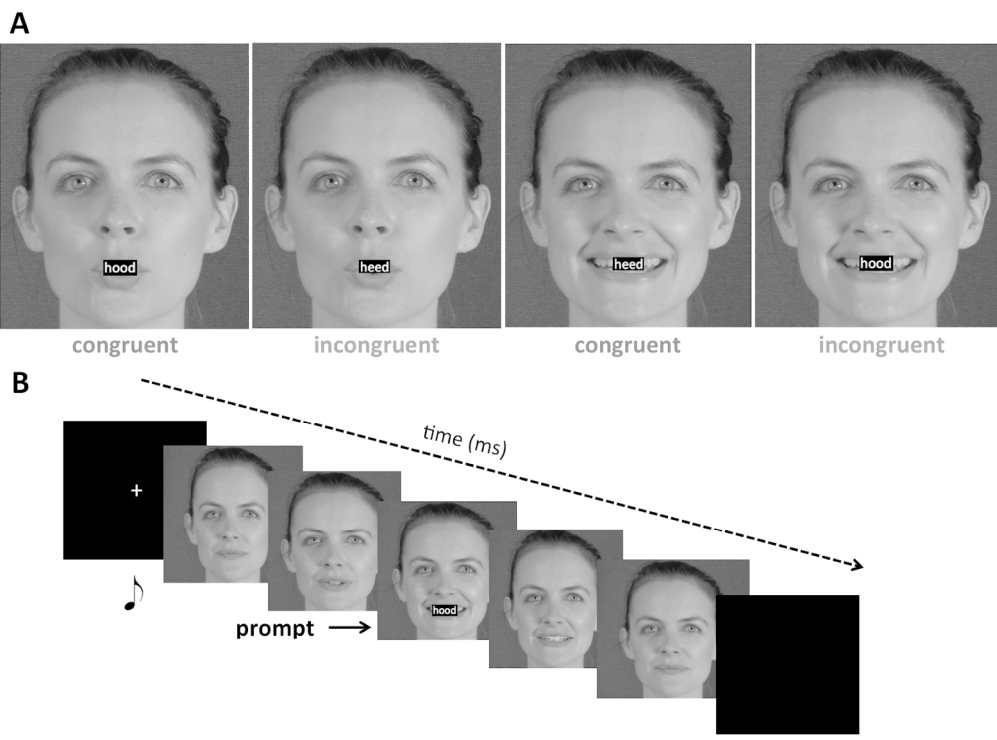
668

669

670

671

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

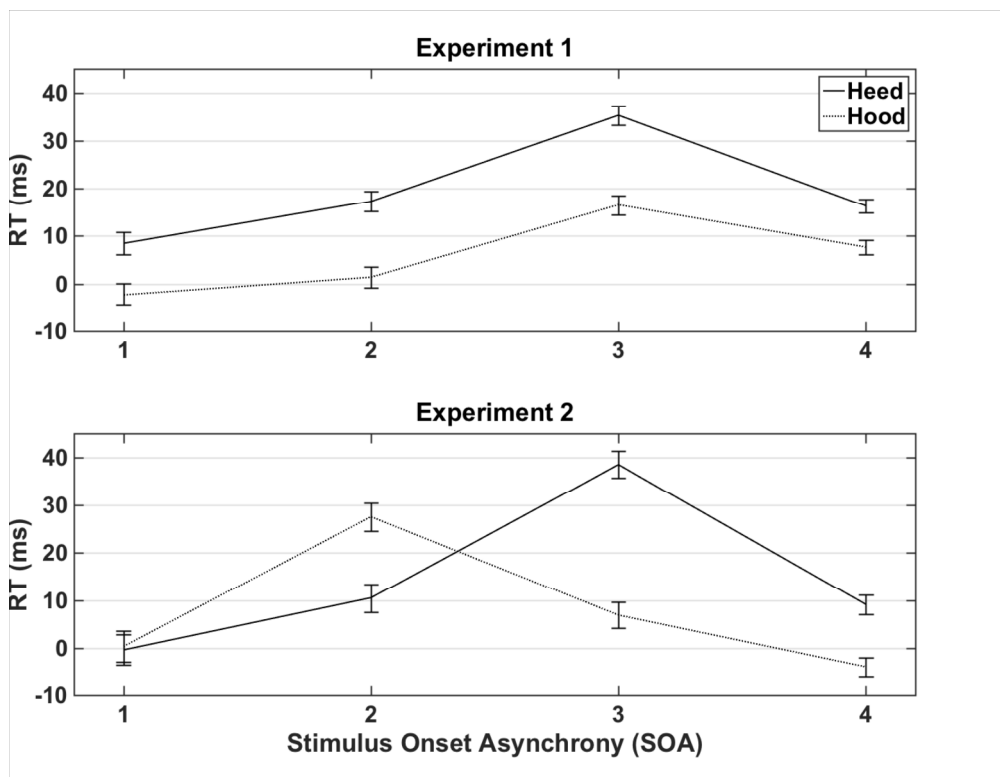


A: congruent trial for hood prompt, incongruent trial for hood prompt, congruent trial for heed prompt, incongruent trial for heed prompt. B. Example of the timeline of an incongruent stimulus pair with hood prompt and heed distracter.

180x134mm (300 x 300 DPI)

Only





Difference scores in milliseconds (incongruent minus congruent pairs) pooled across the Video, Audio, and Audiovisual conditions, for each SOA and separated by prompt, error bars represent one standard error. Top: Experiment 1, B: Experiment 2.

489x377mm (300 x 300 DPI)

only

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Table I. Averages plus standard deviations “( )” for % error and response times in ms for congruent and incongruent stimulus pairs, per prompt, per Stimulus Onset Asynchrony (SOA), and modality, for Experiment 1.

		ERRORS			RESPONSE TIMES			
		Video	Audio	Audiovisual	Video	Audio	Audiovisual	
<i>Heed</i>	<i>Congruent</i>	SOA1	21 (41)	26 (44)	25 (44)	648 (123)	648 (123)	606 (140)
		SOA2	5 (23)	13 (33)	10 (30)	590 (115)	590 (115)	537 (130)
		SOA3	4 (20)	12 (33)	5 (22)	534 (103)	534 (103)	558 (144)
		SOA4	4 (19)	6 (24)	7 (26)	535 (85)	535 (85)	498 (109)
	<i>Incongruent</i>	SOA1	19 (39)	21 (41)	22 (42)	660 (131)	660 (131)	604 (136)
		SOA2	6 (23)	10 (30)	4 (20)	608 (106)	608 (106)	537 (145)
		SOA3	2 (15)	10 (29)	5 (23)	579 (111)	579 (111)	578 (130)
		SOA4	4 (19)	4 (19)	6 (24)	555 (75)	555 (75)	501 (100)
<i>Hood</i>	<i>Congruent</i>	SOA1	14 (34)	17 (38)	15 (36)	655 (134)	655 (134)	600 (119)
		SOA2	5 (21)	10 (30)	5 (23)	575 (110)	575 (110)	534 (129)
		SOA3	4 (19)	10 (31)	5 (23)	553 (107)	553 (107)	553 (139)
		SOA4	2 (15)	4 (20)	3 (17)	524 (72)	524 (72)	492 (139)
	<i>Incongruent</i>	SOA1	16 (37)	19 (39)	15 (36)	635 (125)	635 (125)	607 (130)
		SOA2	5 (23)	9 (29)	9 (28)	590 (119)	590 (119)	529 (138)
		SOA3	3 (16)	11 (31)	8 (28)	567 (102)	567 (102)	587 (149)
		SOA4	2 (13)	4 (19)	3 (18)	537 (80)	537 (80)	498 (104)

Table II. Averages plus standard deviations “( )” for response times in milliseconds for congruent and incongruent stimulus pairs, per prompt, per Stimulus Onset Asynchrony (SOA), and modality, for Experiment 2.

		ERRORS			RESPONSE TIMES			
			Video	Audio	Audiovisual	Video	Audio	Audiovisual
<i>Heed</i>	<i>Congruent</i>	SOA1	20 (40)	24 (43)	20 (40)	20 (40)	671 (130)	670 (131)
		SOA2	4 (21)	6 (24)	7 (26)	4 (21)	582 (118)	608 (130)
		SOA3	3 (18)	5 (23)	9 (29)	3 (18)	537 (107)	561 (133)
		SOA4	6 (24)	2 (15)	5 (23)	6 (24)	529 (86)	513 (96)
	<i>Incongruent</i>	SOA1	28 (45)	25 (44)	28 (45)	28 (45)	661 (138)	685 (132)
		SOA2	5 (23)	6 (24)	4 (20)	5 (23)	618 (123)	632 (140)
		SOA3	2 (15)	7 (26)	8 (27)	2 (15)	575 (112)	593 (120)
		SOA4	5 (23)	6 (23)	4 (20)	5 (23)	529 (83)	524 (97)
<i>Hood</i>	<i>Congruent</i>	SOA1	14 (35)	11 (31)	15 (36)	14 (35)	639 (126)	652 (120)
		SOA2	5 (22)	4 (19)	7 (25)	5 (22)	590 (138)	604 (137)
		SOA3	2 (15)	2 (12)	7 (25)	2 (15)	569 (142)	581 (125)
		SOA4	3 (17)	4 (20)	4 (19)	3 (17)	499 (95)	519 (84)
	<i>Incongruent</i>	SOA1	16 (37)	17 (38)	16 (37)	16 (37)	635 (129)	659 (132)
		SOA2	2 (15)	4 (19)	10 (30)	2 (15)	573 (135)	612 (133)
		SOA3	2 (13)	7 (26)	9 (28)	2 (13)	589 (136)	586 (126)
		SOA4	2 (15)	2 (14)	6 (23)	2 (15)	507 (87)	511 (91)

## Supplementary materials

Table B. Results of the repeated measures ANOVAs on the errors transformed to Rationalised Arcsine Units (RAU) and log-transformed (LogRT) Response Times from Experiment 1. Significant results are indicated with ‘\*’.

Factor	RAU				LogRT			
	df	F	<i>p</i>	$\eta_{2pa}$	df	F	<i>p</i>	$\eta_{2par}$
<i>Prompt</i>	<b>1, 62</b>	<b>7.94</b>	<b>0.006*</b>	<b>0.11</b>	<b>1, 62</b>	<b>8.67</b>	<b>0.005*</b>	<b>0.12</b>
<i>Prompt</i> × <i>Modality</i>	2, 62	0.09	0.917	0	2, 62	0.34	0.712	0.01
<i>Congruence</i>	1, 62	0.07	0.796	0	<b>1, 62</b>	<b>42.45</b>	<b>&lt;0.001*</b>	<b>0.41</b>
<i>Congruence</i> × <i>Modality</i>	2, 62	1.48	0.236	0.05	2, 62	2.16	0.12	0.07
<i>SOA</i>	<b>2.56,</b> <b>159.57</b>	<b>93.75</b>	<b>&lt;0.001*</b>	<b>0.60</b>	<b>2.84, 176</b>	<b>250.61</b>	<b>&lt;0.001*</b>	<b>0.80</b>
<i>SOA</i> × <i>Modality</i>	6, 186	1.43	0.206	0.04	<b>6, 186</b>	<b>14.28</b>	<b>&lt;0.001*</b>	0.32
<i>Prompt</i> × <i>Congruence</i>	<b>1, 62</b>	<b>4.39</b>	<b>0.04*</b>	<b>0.07</b>	<b>1, 62</b>	<b>11.76</b>	<b>0.001*</b>	0.16
<i>Prompt</i> × <i>Congruence</i> × <i>Modality</i>	2, 62	0.63	0.536	0.02	2, 62	4.57	0.01*	0.13
<i>Prompt</i> × <i>SOA</i>	<b>3,</b> <b>185.73</b>	<b>6.27</b>	<b>0.001*</b>	<b>0.09</b>	<b>3, 186</b>	<b>4.01</b>	<b>0.01*</b>	<b>0.06</b>
<i>Prompt</i> × <i>SOA</i> × <i>Modality</i>	6, 186	1.50	0.184	0.05	6, 186	0.54	0.77	0.02
<i>Congruence</i> × <i>SOA</i>	3, 186	0.60	0.615	0	<b>3, 186</b>	<b>9.5</b>	<b>&lt;0.001*</b>	<b>0.13</b>
<i>Congruence</i> × <i>SOA</i> × <i>Modality</i>	6, 186	0.99	0.434	0.03	6, 186	1.44	0.20	0.04
<i>Prompt</i> × <i>Congruence</i> × <i>SOA</i>	3, 184.98	0.27	0.844	0	3, 186	0.78	0.50	0.01
<i>Prompt</i> × <i>Congruence</i> × <i>SOA</i> × <i>Modality</i>	6, 184.98	0.28	0.945	0.01	6, 86	0.86	0.53	0.03

Table B. Results of the repeated measures ANOVAs on the errors transformed to Rationalised Arcsine Units (RAU) and log-transformed (LogRT) Response Times from Experiment 2. Significant results are indicated with ‘\*’.

<i>Factor</i>	RAU				LogRT			
	df	F	<i>p</i>	$\eta_{2pa}$	df	F	<i>p</i>	$\eta_{2par}$
<i>Modality</i>	2, 44	1.33	0.287	0.06	2, 44	2.93	0.06	0.12
<i>Prompt</i>	<b>1, 22</b>	<b>8.38</b>	<b>0.008*</b>	<b>0.28</b>	1, 22	1.31	0.26	0.06
<i>Congruence</i>	<b>1, 22</b>	<b>5.80</b>	<b>0.025*</b>	<b>0.21</b>	<b>1, 22</b>	<b>23.41</b>	<b>&lt;0.001*</b>	<b>0.52</b>
<i>SOA</i>	<b>3, 66</b>	<b>34.22</b>	<b>&lt;0.001*</b>	<b>0.61</b>	<b>2.54, 55.8</b>	<b>130.19</b>	<b>&lt;0.001*</b>	<b>0.86</b>
<i>Modality</i> × <i>Prompt</i>	2, 44	.623	0.541	0.03	2, 44	0.58	0.56	0.03
<i>Modality</i> × <i>Congruence</i>	2, 44	1.293	0.541	0.06	2, 44	1.07	0.35	0.05
<i>Prompt</i> × <i>Congruence</i>	1, 22	0.03	0.955	0	1, 22	1.51	0.23	0.06
<i>Modality</i> × <i>Prompt</i> × <i>Congruence</i>	2, 44	0.85	0.435	0.04	2, 44	2.78	0.07	0.11
<i>Modality</i> × <i>SOA</i>	4.62, 101.59	.215	0.808	0.01	<b>6, 132</b>	<b>12.91</b>	<b>&lt;0.001*</b>	<b>0.37</b>
<i>Prompt</i> × <i>SOA</i>	<b>3, 66</b>	<b>5.795</b>	<b>0.001</b>	<b>0.21</b>	<b>2.54, 5.91</b>	<b>5.5</b>	<b>0.004*</b>	<b>0.2</b>
<i>Modality</i> × <i>Prompt</i> × <i>SOA</i>	6, 132	0.912	0.488	0.04	6, 132	0.8	0.57	0.04
<i>Congruence</i> × <i>SOA</i>	3, 66	1.519	0.218	0.02	<b>3, 66</b>	<b>4.95</b>	<b>0.004*</b>	<b>0.18</b>
<i>Modality</i> × <i>Congruence</i> × <i>SOA</i>	4.38, 96.296	0.51	0.745	0.02	6, 132	1.88	0.09	0.08
<i>Prompt</i> × <i>Congruence</i> × <i>SOA</i>	3, 66	1.207	0.314	0.05	<b>3, 66</b>	<b>4.31</b>	<b>0.008*</b>	<b>0.16</b>
<i>Modality</i> × <i>Prompt</i> × <i>Congruence</i> × <i>SOA</i>	6, 132	0.767	0.597	0.03	6, 132	0.85	0.53	0.04