

1 | Ahn, et al.

1           **Evidence for a Causal Dissociation of the McGurk Effect and**  
2                           **Congruent Audiovisual Speech Perception via TMS**

3

4 EunSeon Ahn<sup>1</sup>, Areti Majumdar<sup>1</sup>, Taraz Lee<sup>\*1</sup>, David Brang<sup>\*1</sup>

5 <sup>1</sup>Department of Psychology, University of Michigan, Ann Arbor, MI 48109

6 \*Joint Senior Corresponding Authors: [djbrang@umich.edu](mailto:djbrang@umich.edu), [tarazlee@umich.edu](mailto:tarazlee@umich.edu)

7

8           **Conflict of Interest Statement**

9 The authors declare no competing financial interests.

10

11           **Acknowledgements**

12 This study was supported by NIH Grants R00DC013828 and R01DC020717. The authors report  
13 no conflicts of interest.

14

15           **Keywords**

16 Multisensory; Audiovisual; Speech; Language; TMS

17

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

## Abstract

Congruent visual speech improves speech perception accuracy, particularly in noisy environments. Conversely, mismatched visual speech can alter what is heard, leading to an illusory percept known as the McGurk effect. This illusion has been widely used to study audiovisual speech integration, illustrating that auditory and visual cues are combined in the brain to generate a single coherent percept. While prior transcranial magnetic stimulation (TMS) and neuroimaging studies have identified the left posterior superior temporal sulcus (pSTS) as a causal region involved in the generation of the McGurk effect, it remains unclear whether this region is critical only for this illusion or also for the more general benefits of congruent visual speech (e.g., increased accuracy and faster reaction times). Indeed, recent correlative research suggests that the benefits of congruent visual speech and the McGurk effect reflect largely independent mechanisms. To better understand how these different features of audiovisual integration are causally generated by the left pSTS, we used single-pulse TMS to temporarily impair processing while subjects were presented with either incongruent (McGurk) or congruent audiovisual combinations. Consistent with past research, we observed that TMS to the left pSTS significantly reduced the strength of the McGurk effect. Importantly, however, left pSTS stimulation did not affect the positive benefits of congruent audiovisual speech (increased accuracy and faster reaction times), demonstrating a causal dissociation between the two processes. Our results are consistent with models proposing that the pSTS is but one of multiple critical areas supporting audiovisual speech interactions. Moreover, these data add to a growing body of evidence suggesting that the McGurk effect is an imperfect surrogate measure for more general and ecologically valid audiovisual speech behaviors.

## 41 **Introduction**

42           While speech perception is predominantly an auditory process, face-to-face conversations  
43 benefit from concurrent visual cues that provide both complementary and redundant information  
44 about an auditory stimulus, particularly in noisy environments (Campbell, 2008; MacLeod &  
45 Summerfield, 1987; Sumbly & Pollack, 1954). For example, listeners can derive both timing and  
46 phonemic information about spoken words simply by watching a speaker's mouth movements  
47 (Luo et al., 2010; Plass et al., 2020; Schroeder et al., 2008). This visual augmentation can  
48 significantly improve speech perception accuracy, especially when acoustics are compromised.  
49 Given the inherent multisensory nature of speech, it is important to examine how the brain enables  
50 vision to support language to understand speech processing in naturalistic contexts.

51           Traditionally, studies of audiovisual processing have relied on the use of the McGurk  
52 effect, in which an auditory phoneme (e.g., /ba/) is paired with the visual movie from a different  
53 phoneme (e.g., /ga/), resulting in the perception of a fused or unique sound (e.g., /da/) (McGurk  
54 & MacDonald, 1976). Using McGurk stimuli, prior research has identified the left posterior superior  
55 temporal sulcus (pSTS) as a crucial region that facilitates the integration of auditory and visual  
56 information (Benoit et al., 2010; Bernstein et al., 2008; Irwin et al., 2011; Nath et al., 2011;  
57 Sekiyama et al., 2003; Szyckik et al., 2012). For example, individual differences in the strength of  
58 the McGurk effect are correlated with fMRI activity in the pSTS during the experience of the illusion  
59 (Nath & Beauchamp, 2012), and both transcranial magnetic stimulation (TMS) and damage  
60 following a stroke in this region are associated with reduced McGurk effect percepts (Beauchamp  
61 et al., 2010; Hickok et al., 2018). However, researchers have recently questioned whether the  
62 findings from research using McGurk combinations can generalize to more natural audiovisual  
63 integration processes (for review see Alsius et al., 2018). For example, numerous studies (Brown  
64 & Braver, 2005; Hickok et al., 2018; Van Engen et al., 2017) have reported weak correlations  
65 between the McGurk effect and other measures of audiovisual speech processing in individuals,

4 | Ahn, et al.

66 raising doubts about whether the mechanisms that enable the McGurk effect are the same as  
67 those that process natural audiovisual speech.

68 Whether the McGurk effect and more general audiovisual processes rely on the same  
69 mechanisms is a significant issue in the field because the McGurk effect had long been accepted  
70 as a standard measure for quantifying audiovisual speech integration (Alsius et al., 2018; Van  
71 Engen et al., 2022). This holds especially true for clinical populations including those with autism  
72 spectrum disorder (ASD). Individuals with ASD often exhibit difficulties in communication and  
73 social interaction and many researchers believe that this could be in part attributed to impairments  
74 in multisensory processing. To study this, researchers have repeatedly used the likelihood of  
75 McGurk percepts to demonstrate that individuals with ASD show altered multisensory processing  
76 (Feldman et al., 2022; Gelder et al., 1991; Stevenson et al., 2014; Williams et al., 2004; Zhang et  
77 al., 2019). These studies have consistently shown that individuals with ASD experience  
78 significantly weaker McGurk effects than the neurotypical population. However, it is important to  
79 note that the McGurk effect focuses on the cognitive cost of processing conflicting auditory and  
80 visual information. In contrast, in normal speech contexts, listeners avoid integrating conflicting  
81 audiovisual speech information (Brang, 2019; Seijdel et al., 2023). Specifically, the incongruent  
82 pairing that is necessary to elicit McGurk fusion responses has been regarded as artificial and  
83 unnatural, showing limited features present in everyday speech (Van Engen et al., 2022) because  
84 face-to-face conversations yield congruent combinations of auditory and visual speech.  
85 Moreover, McGurk studies tend to examine audiovisual processing using isolated phonemes  
86 rather than complete words, casting further doubt on their applicability to natural speech.  
87 Consequently, tasks that use more naturalistic stimuli, like complete words and congruent  
88 audiovisual pairings, may be better able to clarify the role of visual information in everyday speech  
89 perception, thus advancing beyond the limited context of the McGurk effect.

90           While strong correlative research has identified the left pSTS as a region associated with  
91 the generation of the McGurk effect, only two studies to date have used causal methods  
92 (Beauchamp et al., 2010; Hickok et al., 2018). In a 2010 study, Beauchamp et al. reported that  
93 single pulse transcranial magnetic stimulation (TMS) applied to the left pSTS significantly reduced  
94 perception of the McGurk effect, providing strong evidence for the role of the pSTS in the  
95 generation of this illusion. TMS is a noninvasive brain stimulation method that involves the  
96 application of magnetic pulses to targeted brain areas that causes neurons to immediately  
97 depolarize thus injecting noise into ongoing processes (Hallett, 2000). This method has been  
98 effectively used to study the causal mechanisms underlying numerous cognitive and perceptual  
99 processes (Hallett, 2000; Rossini & Rossi, 2007; Walsh & Cowey, 2000). In Beauchamp et al.'s  
100 2010 study, the authors conducted two separate experiments, each with sample sizes of 9, using  
101 similar task designs but two different speakers (experiment 1 used a female speaker and  
102 experiment 2 used a male speaker) and two different phonemes (experiment 1 used auditory /BA/  
103 with visual /GA/ and experiment 2 used auditory PA with visual /NA/ or /KA/). The authors used  
104 two separate experiments with different phonemes and speakers to ensure the robustness of their  
105 findings across different speakers and stimuli combinations. Results showed that single-pulse  
106 stimulation to the left pSTS reduced the average frequency of fused percepts in McGurk  
107 conditions by 54% (experiment 1) and 21% (experiment 2) compared to stimulation applied to the  
108 control site. The authors also showed that only stimulation applied within 100 ms of the onset of  
109 the auditory stimulus (100 ms before until 100 ms after) reduced the frequency of McGurk  
110 percepts, while stimulation applied outside of this time range did not influence frequency.

111           Extending this research to clinical populations, Hickok et al. (2018) tested patients with a  
112 recent stroke using a McGurk effect paradigm with the goal of identifying lesioned areas of the  
113 brain that reduced McGurk effect percepts. Partially replicating the TMS work, Hickok et al. (2018)  
114 showed that stroke lesions in the broad superior temporal lobe (as well as in auditory and visual

115 areas in the superior temporal and lateral occipital regions) resulted in the greatest deficits in  
116 McGurk perception, adding support to the model that the left pSTS enables the generation of the  
117 McGurk effect.

118 While both prior TMS and stroke lesion mapping studies identified the left pSTS as being  
119 causally relevant to the generation of the McGurk effect, those studies were not designed to test  
120 the relevance of this structure on more natural, congruent audiovisual speech perception  
121 behaviors. Building upon Beauchamp et al. (2010)'s findings, here we sought to address the  
122 involvement of left pSTS in other aspects of audiovisual speech processing beyond the generation  
123 of the McGurk percept and further investigate whether the McGurk effect is a good proxy for  
124 audiovisual processing. Based on past literature, two clear predictions emerged: 1) transient  
125 disruption of left pSTS activity will impair both the McGurk effect and the normal benefits from  
126 audiovisual speech. Such a finding would indicate that left pSTS is a critical hub for audiovisual  
127 speech generation in general and that the McGurk effect is a good proxy for natural audiovisual  
128 speech behaviors. Or 2) transient disruption of left pSTS activity impairs the McGurk effect with  
129 minimal impact on the normal benefits from audiovisual speech. Such findings would indicate that  
130 this region is likely only one of many critical structures necessary for audiovisual speech  
131 generation in general, reflecting only a subset of the information relayed from visual to auditory  
132 speech regions and laying the groundwork for research to understand which information is relayed  
133 through this hub.

134 To test these predictions, we assessed the impact of TMS on both the frequency of  
135 subjects' McGurk effect percepts (which captures how visual information can change, or  
136 modulate, the perception of auditory information) and measures of audiovisual facilitation (which  
137 focuses on how visual information aids and facilitates the processing of the auditory information).  
138 By distinguishing between these two measures, we can better understand how mismatched visual  
139 information can modulate auditory perception and whether this is dissociable from the ability of

140 congruent visual information to improve and facilitate the processing of concurrent auditory  
141 information. Towards this goal, we used an audiovisual task with real word stimuli, rather than  
142 phonemes that are typically utilized in most McGurk-type designs, as well as a larger sample ( $n$   
143 = 21) than the prior TMS study. We hypothesized that pSTS stimulation would affect the McGurk  
144 effect more than congruent audiovisual benefits, providing evidence that the visual modulation of  
145 auditory speech relies on different neural mechanisms and consequently brain regions from  
146 audiovisual facilitation.

147

## 148 **Methods**

### 149 **Subjects**

150 25 healthy subjects (10 males, mean age = 24.7, 4 left-handed) with self-reported normal  
151 hearing and vision without a history of neurological disorder participated in the study. Four of the  
152 25 total subjects either voluntarily withdrew from the study or experienced data errors resulting in  
153 21 total subjects who completed the study (7 male, mean age = 24.2). Beauchamp et al. (2010)  
154 prescreened their subjects to include only those who reported strong McGurk effects. They  
155 justified this pre-selection process because there are large individual differences in susceptibility  
156 to the McGurk effect (Nath & Beauchamp, 2012) and variable reliance on lip movements during  
157 speech perception (Gurler et al., 2015). However, we did not exclude any subjects based on their  
158 McGurk susceptibility. While previous research has shown that the specific audiovisual stimulus  
159 used affects the strength and frequency of McGurk effects experienced by subjects (Basu Mallick  
160 et al., 2015), the stimuli used in the current study successfully evoked McGurk percepts in the  
161 majority of individuals tested in a prior study by our lab (Brang et al., 2020) and the same stimulus  
162 set has been shown to produce robust congruent audiovisual benefits during speech recognition  
163 (Ross et al., 2007). Therefore, to maximize the generalizability of our study to everyday speech  
164 perception, we did not exclude subjects based on their susceptibility to the McGurk effect.

165 To estimate the necessary sample size, we conducted an a priori power analysis using  
166 G\*Power (Faul et al., 2007) based on Beauchamp et al. (2010)'s data. In comparing the frequency  
167 of reported fusion responses during pSTS stimulation versus control site stimulation, their results  
168 yielded Cohen's D values of 3.22 and 8.43 across two experiments. Given the effect size of 3.22  
169 (the smaller of the two estimates) is considered extremely large using Cohen's criteria (1988), we  
170 would need a minimum sample size of 4 to replicate their effects with a significance criterion of  
171  $\alpha = .05$  and power = .95. However, as our goal was to examine whether congruent audiovisual  
172 behaviors were affected as well, we made the more conservative assumption that if present, TMS  
173 effects on congruent audiovisual behaviors would be at least 25% the magnitude of effect of TMS  
174 on McGurk percepts. Repeating the power analysis with an alpha of 0.05, power of .80, and effect  
175 size of .805 in a two-tailed paired t-test design yielded a minimum sample size of 15. We sought  
176 to exceed this number and collect as much data as possible before the end of April 2023.

177 All participants gave informed consent prior to the experiment. This study was approved  
178 by the institutional review board at the University of Michigan.

179

### 180 ***Experimental Task***

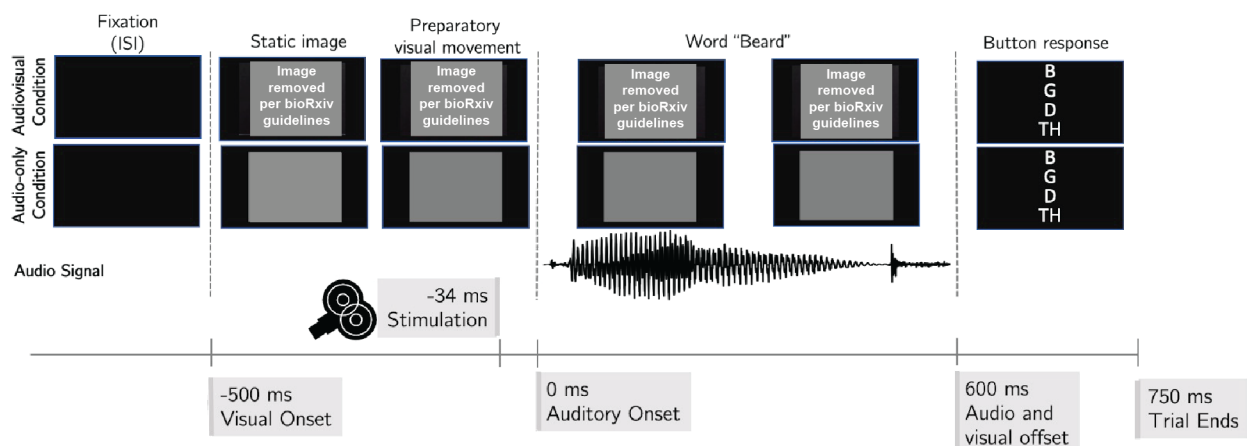
181 Our audiovisual speech task used stimuli adapted from a prior study by Ross et al. (2007).  
182 From their larger set, we selected 32 single syllable words starting with the consonants 'b', 'f', 'g',  
183 or 'd', with the initial vowel sound approximately balanced across the consonant groups (e.g.,  
184 'bag', 'gag', 'dad', 'fad').

185 The schematic of the trials is shown in Figure 1. On each trial, a female speaker produced  
186 a high frequency monosyllabic word starting with the consonant 'b', 'f', 'g', or 'd'. The trials were  
187 either audio-only, visual-only, audiovisual incongruent, or audiovisual congruent. Pink noise (SNR  
188 of -4.6 dB) was applied to all auditory stimuli (visual-only trials included pink noise but no speech)  
189 to increase the relative difficulty of the task, to avoid ceiling effects, and because the addition of



190 noise tends to increase the reliance on visual speech information (Ross et al., 2007). The noise  
191 level was set based on piloting to lower accuracy in the audio-only condition away from the ceiling.  
192 In trials with a visual stimulus, the speaker's video appeared 500 ms prior to the onset of the  
193 auditory stimulus. In audio-only trials, a gray box appeared 500 ms prior to the auditory onset to  
194 provide an equivalent temporal cue. Visual stimuli were recorded at 29.97 frames per second,  
195 trimmed to 1100 ms in length, and adjusted so that the first consonantal burst of sound occurred  
196 at 500 ms during each video.

197 In Beauchamp et al. (2010), the authors found maximally diminished fusion responses  
198 when the pSTS was stimulated 34 ms before auditory onset. Accordingly, here we applied a single  
199 TMS pulse 34 ms prior to audio onset (see TMS parameters below). Six hundred milliseconds  
200 following auditory onset, subjects were prompted to report the initial consonant of the seen (in  
201 visual trials) or heard word using a gamepad (Logitech F310) from the 4 options displayed on the  
202 screen. Subjects were asked to choose the option that sounded closest to what they heard if they  
203 were unsure. They were informed that the word that they heard may not be real words. The 4  
204 response options always included the initial consonant of the spoken word, the initial consonant  
205 of the video (i.e., a viseme, which is the visual equivalent of a phoneme in spoken language), as



**Figure 1.** Trial schematic for the word 'beard'. All trials started out with a black screen lasting between 125 - 375 ms. 500 ms prior to the auditory onset, either a gray box (audio-only condition) or the video of the speaker (AV and lipreading conditions) appeared. For blocks in which TMS was applied, stimulation was applied 34 ms prior to the onset of the audio. Following the audio/visual offset, participants were given 1.5 seconds to identify via button press which initial consonantal sound the word they heard. Note that each of the three conditions had pink noise mixed in with the audio signal (the visual-only condition contained only the pink noise at time 0).

206 well as the two common McGurk fusion percepts. These four response options remained  
207 consistent across all stimuli and conditions, even if some options were not relevant to certain  
208 conditions. The task was completed via a desktop computer using Psychtoolbox-3 (Brainard &  
209 Vision, 1997; Kleiner et al., 2007; Pelli & Vision, 1997) with participants seated approximately 60  
210 cm away from the screen at eye level.

211 The task consisted of 3 blocks in total. Two blocks included TMS stimulation of one  
212 anatomical region per block (the left pSTS or vertex; defined below) and one block included no  
213 stimulation. The order of stimulation was counterbalanced across participants. In total, the study  
214 consisted of 384 trials, with 128 trials in each block with 4 audiovisual conditions (audio-only,  
215 visual-only, audiovisual congruent, audiovisual incongruent) divided equally within each block (32  
216 trials per condition per block). Each block took approximately 7 minutes to complete, and  
217 participants were given the option to take breaks between each block.

218 Based on prior literature and internal piloting in which subjects provided free-response  
219 reports of what they heard, McGurk fusions were expected for two sets of audiovisual  
220 combinations: auditory words starting with either a B or F and visual words starting with either  
221 G or D, respectively. For example, auditory 'buy' + visual 'guy' typically resulted in the percept  
222 'die' or 'thigh', and auditory 'fad' + visual 'dad' typically resulted in the percept 'tad' or 'thad'. To  
223 ensure that the auditory and visual components of each word were presented the same number  
224 of times throughout the experiment, half of the incongruent audiovisual trials had the modality of  
225 these pairings flipped (auditory words starting with either a G or D and visual words starting with  
226 either or B or F, respectively). These flipped pairings were not expected to generate fused  
227 responses, although they were still expected to reduce accuracy and slow reaction times; during  
228 piloting subjects were more likely to report 'hearing' the visual percept than a fused percept.

229

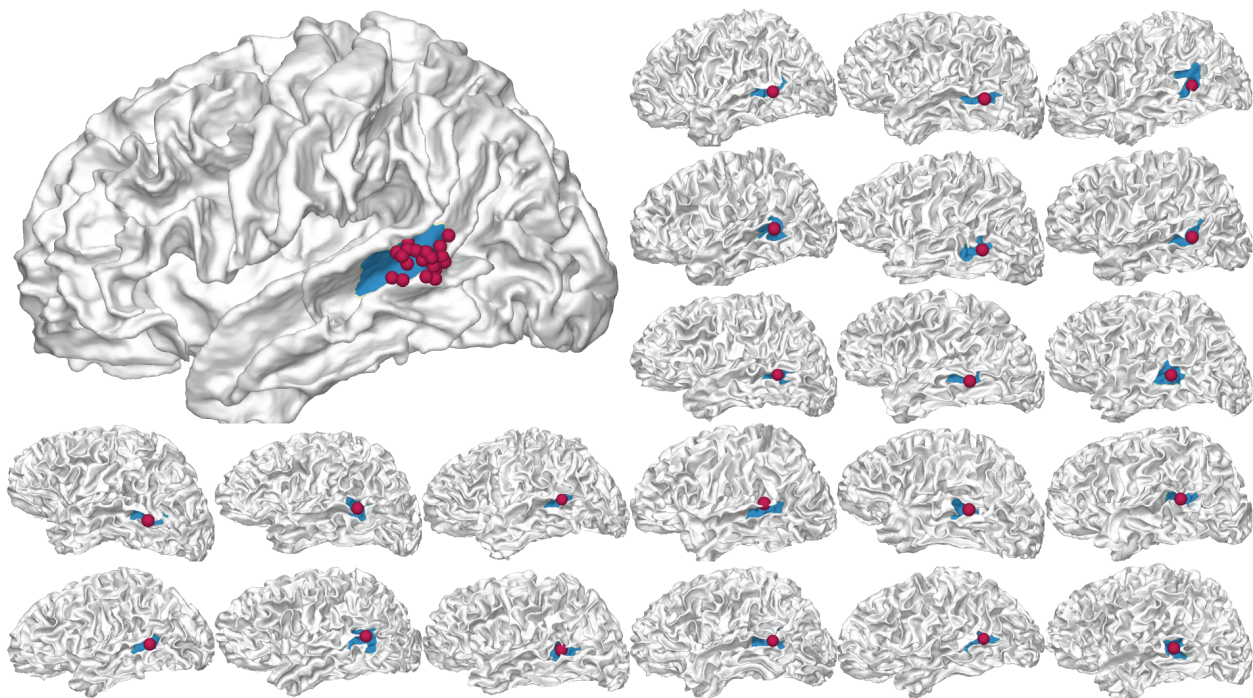
230 ***MRI and TMS Procedure***

231 Prior to the behavioral testing, on a previous day all participants underwent a structural  
232 MRI scan using a 3T GE MR 750 scanner to acquire T1-weighted images to be used for TMS  
233 guidance; three of 21 participants' MRI data was acquired as part of a different experimental  
234 paradigm from our lab (Ganesan et al., 2022). Individual participants' T1 scans were processed  
235 using Freesurfer (<http://surfer.nmr.mgh.harvard.edu/>) for cortical reconstruction and volumetric  
236 segmentation. The Freesurfer-generated pial and white matter reconstructions were then used to  
237 localize the left pSTS target coordinates according to the labels of the automatic cortical  
238 parcellation and automatic segmentation volumes. Specifically, we used the individual subject  
239 coordinates from the center of 'lh-bankssts', which is the bank of the left hemisphere superior  
240 temporal sulcus, as our pSTS target, and the vertex was defined as the midline of the postcentral  
241 gyrus using the subjects' structural MRI scan.

242 TMS was applied through a MagPro X100 using a 70 mm figure-8 shaped TMS coil (MCF-  
243 B70, MagVenture Inc.). For each subject, their respective stimulation intensity was determined by  
244 obtaining their resting motor threshold and multiplying it by 1.1 (110%) as is the standard for  
245 single-pulse TMS thresholding (Kallioniemi & Julkunen, 2016; Sondergaard et al., 2021). The  
246 resting motor threshold is the lowest stimulation intensity necessary to evoke a consistent motor  
247 response while targeting the motor cortex. Specifically, this is the threshold at which an  
248 electromyographic motor response that is greater than 50  $\mu$ V from the first dorsal interosseus  
249 muscle measuring is observed 5 out of 10 times (Mills & Nithi, 1997; Rossini & Rossi, 2007). In  
250 our motor thresholding, our study recorded from the right first dorsal interosseous muscle while  
251 stimulating the left primary motor cortex. The mean resting motor threshold used in our study was  
252 51.7% of the maximum stimulator output. Once the target intensity was determined for each  
253 subject, the same intensity (110% of the resting motor threshold) was used for the entirety of their  
254 TMS session.

255           Following motor thresholding, participants completed the audiovisual task with each block  
256 targeting a different TMS site: pSTS, vertex, and no stimulation. The vertex stimulation block  
257 served as a control condition for the non-specific effects of TMS on behavior (e.g., scalp  
258 sensation, auditory stimulation, induced current in the brain, etc.) (Jung et al., 2016). Brainsight's  
259 neuro-navigation system (Brainsight; Rogue Research) was used to target the stimulation sites in  
260 real time by registering participants' facial landmarks to participants' structural T1s using  
261 headbands containing trackers. The target coordinates for the pSTS that were obtained via  
262 Freesurfer reconstructions were used to guide TMS stimulation. Figure 2 shows the location of  
263 the left pSTS stimulation sites across subjects along with the anatomical label. The average  
264 coordinates of the left pSTS stimulation site across all subjects were ( $x = -65.1 \pm 3.3$ ,  $y = -48.5$   
265  $\pm 4.8$ ,  $z = 7.3 \pm 3.5$ ) on a standard MNI-152 brain.

266           To ensure precise delivery of the stimulation pulse relative to the onset of the auditory  
267 stimuli, the TMS pulse was auto-triggered through our MATLAB script using the MAGIC toolbox



**Figure 2.** (top left) White matter rendering of the *cv\_s\_avg35\_inMNI152* brain showing MNI transformed pSTS stimulation sites from all 21 subjects. Each sphere denotes the specific region where TMS was applied for a single subject. The blue region denotes the pSTS (Freesurfer label banksts) which was used for targeting. The smaller brains show the stimulation site for each of the 21 subjects on their own Freesurfer reconstructed brain.

268 (Saatlou et al., 2018). Due to the position of the coil required to target the pSTS, earbuds, instead  
269 of headphones, were used to deliver the audio. Headband positions were validated both before  
270 and after the motor thresholding prior to the audiovisual task. If significant deviations from the  
271 original facial landmarks were observed, the landmarks were re-registered.

272

### 273 **Data Analyses**

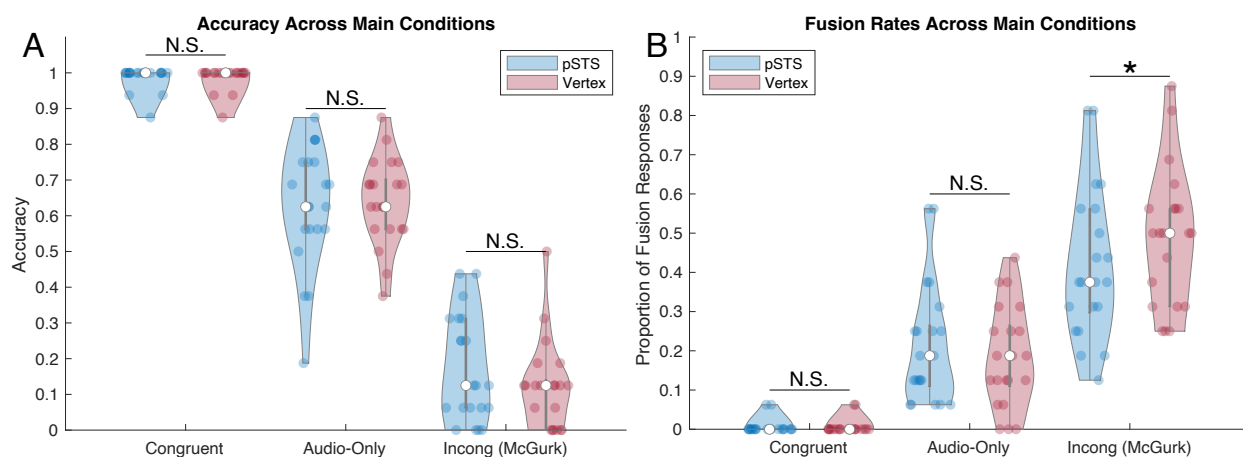
274 Our primary comparisons of interest examined changes in performance following  
275 stimulation of the left pSTS relative to Vertex stimulation for McGurk fusion frequency (the  
276 frequency at which subjects reported a fused percept on McGurk pairings) and for accuracy during  
277 congruent audiovisual trials. These measures were selected to be comparable with prior causal  
278 audiovisual research such as Beauchamp et al. (2010) and Hickok et al. (2018). As noted above,  
279 the incongruent audiovisual condition included 50% of stimuli combinations whose phonemic  
280 pairings typically evoke McGurk fusions (e.g., auditory 'bet' and visual 'get' yield the percept of  
281 'debt') and the reversed pairings (e.g., auditory 'get' and visual 'bet') that do not typically result in  
282 McGurk fusions, for the purpose of counterbalancing stimuli. McGurk fusion frequency was only  
283 estimated from the 50% of incongruent audiovisual trials that contained auditory words starting  
284 with either a B or F and visual words starting with either G or D, respectively. To ensure  
285 comparability with the McGurk fusion frequency measure, performance on congruent audiovisual  
286 trials was restricted to the same set of auditory words, unless noted otherwise. For example, the  
287 auditory word 'bet' was included in both the McGurk fusion frequency analysis (auditory 'bet' +  
288 visual 'get') and the congruent accuracy and RT analyses ( auditory 'bet' + visual 'bet'), but the  
289 auditory word 'get' was excluded (as auditory 'get' + visual 'bet' typically fails to yield fusion  
290 percepts). By restricting our analyses of congruent audiovisual trials, we ensured a similar base  
291 rate for the relevant comparisons across the conditions. Secondary analyses examined the overall  
292 pattern of results using a two-way ANOVA to test the impact of stimulation conditions (pSTS,

293 vertex, or no stimulation) across the four conditions. These secondary analyses were conducted  
294 on all trials (McGurk and Non-McGurk stimuli combinations). Although the original degrees of  
295 freedom are reported here for clarity,  $p$  values were subjected to Greenhouse–Geisser correction  
296 where appropriate (Greenhouse & Geisser, 1959).

297

## 298 Results

299 Figure 3 shows the accuracy and fusion response rates across the congruent, audio-only,  
300 and incongruent McGurk trials. As noted in the methods, McGurk fusion frequency was only  
301 estimated from the 50% of incongruent audiovisual trials that contained auditory words starting  
302 with either a B or F, and visual words starting with either G or D, respectively. To ensure  
303 comparability with the McGurk fusion frequency measure, we first restricted our analyses to the  
304 same set of auditory words across all conditions (later analyses and Figure 4 reflect the data from  
305 all trials). Collapsing across stimulation site, congruent audiovisual trials showed higher accuracy  
306 relative to audio-only trials ( $t(20) = 14.4$ ,  $p < .001$ ,  $d = 3.14$ ), and audio-only trials showed higher  
307 accuracy relative to incongruent McGurk trials ( $t(20) = 16.6$ ,  $p < .001$ ,  $d = 3.61$ ) validating the  
308 positive and negative influence of vision depending on congruent and incongruent contexts.

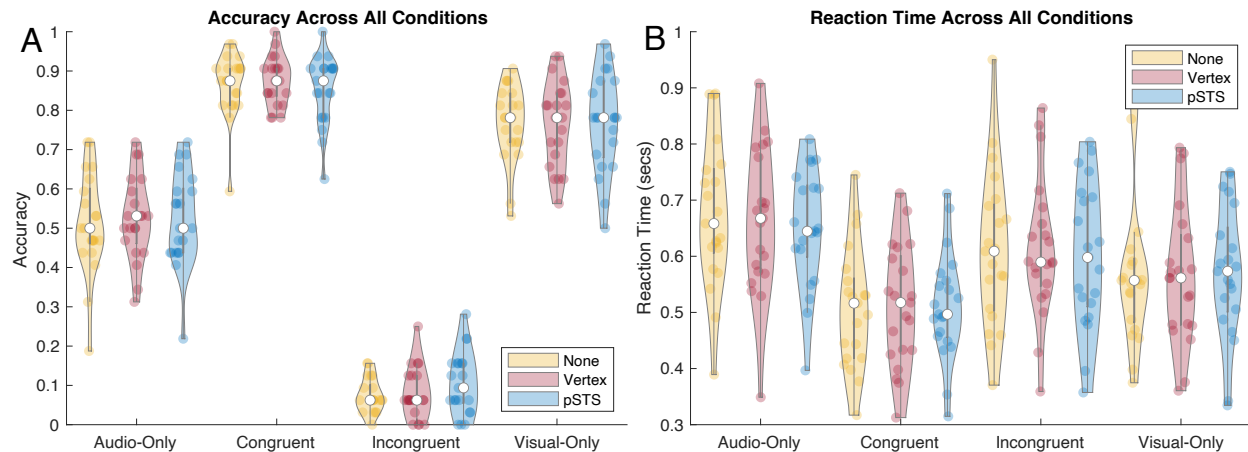


**Figure 3.** Violin plots showing accuracy (A) and fusion response rates (B) for pSTS and vertex stimulation sites, across the three main experimental conditions. Data were restricted to 'McGurk' stimuli to ensure comparability across analyses. Center circles indicate the median, gray boxes reflect the upper and lower quartile ranges, whiskers the min and max excluding outliers, and colored points are individual subject responses. TMS stimulation of the pSTS lowered the rate of fusion responses made by subjects on incongruent McGurk trials but did not impact congruent audiovisual trial accuracy. \* $p < .05$ .

309 To directly replicate the comparison made by Beauchamp et al. (2010), we examined the  
310 proportion of fusion responses made by subjects on Incongruent McGurk trials. Consistent with  
311 their data, single-pulse TMS to the left pSTS significantly reduced subjects' likelihood of  
312 perceiving a McGurk fusion percept compared to the TMS to the vertex ( $t(20) = 2.16$ ,  $p = 0.043$ ,  
313  $d = 0.472$ ; Fig 3). Next, we examined the effect of TMS on congruent audiovisual accuracy. In  
314 contrast to McGurk trials, behavior on congruent audiovisual trials did not differ between  
315 stimulation of the pSTS and Vertex ( $t(20) = 0.00$ ,  $p = 1.00$ ,  $d = 0.00$ ; Fig 3A); note that the statistics  
316 are at the absolute minimum because the average means across stimulation sites were identical.  
317 Comparing McGurk fusion rates and congruent audiovisual accuracy in a 2x2 repeated measure  
318 ANOVA additionally demonstrated a significant interaction between the two [ $F(1,20) = 5.15$ ,  $p =$   
319  $.034$ ,  $\eta_p^2 = 0.205$ ] indicating that pSTS stimulation affected McGurk fusion to a greater degree  
320 than congruent audiovisual accuracy.

321 Following these planned comparisons, we calculated repeated measures ANOVAs for  
322 accuracy and reaction time (RT) data across all conditions. Violin plots showing the distribution  
323 of accuracy and reaction time measures are shown in Figure 4. Repeated measures ANOVA  
324 applied to accuracy data revealed a main effect of visual type [ $F(3,60) = 561.9$ ,  $p = 7.4E-39$ ,  $\eta_p^2$   
325  $= 0.913$ ], but no effect of stimulation site [ $F(2,40) = 0.242$ ,  $p = .737$ ,  $\eta_p^2 = 0.002$ ], nor an interaction  
326 between visual type and stimulation site [ $F(6,120)$ ,  $p = .725$ ,  $\eta_p^2 = 0.007$ ]. These results  
327 demonstrated a strong influence of visual content on task performance (such that congruent visual  
328 speech improves speech recognition and incongruent visual speech impairs it) and that there was  
329 no systematic effect of stimulation across conditions. Notably, visual only performance was higher  
330 than audio-only performance likely due to the high level of auditory-noise.

331 RT data mirrored those of the accuracy data with a main effect of visual type [ $F(2.3,46.08)$   
332  $= 54.379$ ,  $p = 1.19E-13$ ,  $\eta_p^2 = 0.185$ ], but no effect of stimulation site [ $F(2,40) = 0.152$ ,  $p = 0.860$ ,  
333  $\eta_p^2 = 0.001$ ], nor an interaction between visual type and stimulation site [ $F(6,120) = 0.820$ ,  $p =$



**Figure 4.** Violin plots showing accuracy (A) and reaction time (B) for each stimulation site and condition. Center circles indicate the median, gray boxes reflect the upper and lower quartile ranges, whiskers, the min and max excluding outliers, and colored points are individual subject responses.

334 0.556,  $\eta_p^2 = 0.002$ ]. Moreover, regardless of stimulation site, congruent visual speech was  
335 responded to more quickly than incongruent speech ( $t(20) = 8.44, p < .001, d = 1.84$ ), which in  
336 turn was faster than audio-only trials ( $t(20) = 3.86, p < .001, d = 0.842$ ) showing the strong benefit  
337 of visual information regardless of congruity. Stimulation of the pSTS (relative to vertex  
338 stimulation) affected neither the RTs of congruent trials ( $t(20) = 0.00, p = 1.00, d = 0.00$ ) nor  
339 incongruent trials ( $t(20) = 0.552, p = .587, d = 0.121$ ) showing the resistance of this audiovisual  
340 speeding to pSTS stimulation.

341

## 342 Discussion

343 This study used real word stimuli to investigate the role of the left pSTS on audiovisual  
344 speech processes, including McGurk effect percepts, audiovisual facilitation, audiovisual reaction  
345 times, and lipreading. Using TMS in healthy subjects, our data demonstrate that stimulation of the  
346 left pSTS significantly disrupts the experience of the McGurk effect, reducing the frequency of  
347 reported fusion percepts, while leaving the other audiovisual processes intact. Our study extends  
348 the work of Beauchamp et al. (2010) demonstrating that TMS applied over the left pSTS, but not  
349 a control site (vertex), reduces the strength of the McGurk effect. However, as their behavioral  
350 task could not robustly capture the benefit of congruent visual stimuli on auditory speech



351 perception, it remained unclear whether this region contributes specifically to McGurk processes,  
352 or to the general audiovisual speech integration processes at large. Our study extends  
353 Beauchamp et al. (2010)'s finding by confirming that the neural mechanism facilitating the McGurk  
354 perception, which many consider unnatural and artificial, can be dissociated from the other  
355 audiovisual processes like audiovisual facilitation and lipreading, which occur in everyday speech.  
356 Similarly, our results are consistent with models proposing that the pSTS is only one of the  
357 multiple critical areas supporting audiovisual speech interactions. This work adds to the growing  
358 evidence that McGurk processing relies on additional neural mechanisms beyond our everyday  
359 audiovisual speech.

360         One way to explain the involvement of pSTS in McGurk processing but not in audiovisual  
361 facilitation is that direct projections from visual motion area MT/V5 to the auditory cortex (Besle  
362 et al., 2008) allow the visual facilitation of auditory process to bypass the pSTS. Similarly, it is  
363 also possible that the neural processes underlying audiovisual facilitation alternatively recruit  
364 frontal structures to recover speech information through sensorimotor integration (Du et al., 2014,  
365 2016; Hickok & Poeppel, 2007) without involving the pSTS. Along with Beauchamp et al. (2010)'s  
366 findings of reduced McGurk effect when targeting the left pSTS that we replicated here, a similar  
367 report of a weakened McGurk effect has been reported in patients following strokes near the left  
368 pSTS (Hickok et al., 2018). Adding onto these findings, our work provides compelling evidence  
369 that McGurk processing is a specific form of audiovisual speech integration that's independent of  
370 other audiovisual speech enhancement. While McGurk processing may share some naturalistic  
371 features of everyday audiovisual speech processing, it also contains additional less ecologically  
372 valid properties like mismatched auditory and visual information. Consequently, it is possible that  
373 the left pSTS is more responsible for detecting or reconciling minor incongruities across modalities  
374 and re-evaluating the transformation of visemes to phonemes.

375           The work of Hickok et al. (2018) and Van Engen et al. (2017) similarly aimed to investigate  
376 the relationship between McGurk susceptibility and the use of visual information to facilitate  
377 speech perception. Indeed, both studies reported minimal correlations between the two  
378 measures, providing evidence against the widespread use of McGurk susceptibility as an index  
379 for audiovisual speech integration (Alsius et al., 2007; Jones & Callan, 2003; Paré et al., 2003;  
380 Van Wassenhove et al., 2007).

381           While we replicated Beauchamp et al. (2010)'s main finding, such that single pulse TMS  
382 to left pSTS diminishes the McGurk effect, we observed a much smaller change in behavior. In  
383 comparison to the large effect size reported by Beauchamp et al. (2010), with a Cohen's D of 3.22  
384 and 8.43 (across two separate experiments using different speakers and phonemes), our effect  
385 size for the difference in the frequency of fusion responses was much more moderate, yielding a  
386 Cohen's D of 0.472. Such disparity in the effect size may be accounted for by the differences in  
387 the two study designs or the common trend for effect sizes to lower with larger sample sizes (e.g.,  
388 Slavin and Smith (2009).

389           Our study differed from Beauchamp et al. (2010) in multiple ways. First, Beauchamp et al.  
390 (2010) only had 2 conditions: congruent and incongruent McGurk conditions. All trials used the  
391 same single auditory phoneme either matched with its congruent viseme or an incongruent  
392 viseme that is known to create a fusion percept. Conversely, we included multiple additional  
393 conditions to provide context for the subjects' performance and to enable counter-balancing  
394 stimuli. Second, we used real monosyllabic words rather than phonemes for more naturalistic  
395 speech stimuli. This was done to address a common criticism of McGurk studies which argues  
396 that phonemes are highly artificial and do not reflect natural speech. Third, Beauchamp et al.  
397 (2010) determined the location of left STS using both anatomical (5 subjects based on landmarks)  
398 and functional (7 subjects based on individual subjects' fMRI activation patterns) approaches.  
399 However, our study relied only on anatomical landmarks to determine the intended stimulation

400 site. Fourth, we used the vertex as a control site rather than “a control TMS site dorsal and  
401 posterior to the STS” as reported by Beauchamp et al. (2010). This was to ensure we were using  
402 a consistent control site across subjects. Fifth, we did not exclude subjects based on whether  
403 they experienced strong McGurk effects. Our prior work (Brang et al., 2020) using similar stimuli  
404 set showed that most individuals report some level of fusion responses when presented with our  
405 word stimuli. Therefore, we wanted to ascertain that the results of our TMS were generalizable  
406 and not restricted to only those who experience strong McGurk percepts. Indeed, in the no-  
407 stimulation condition for the current dataset, all participants experienced a decrease in accuracy  
408 in the McGurk audiovisual incongruent condition relative to the audio-only condition (range 12.5 -  
409 75% decrease in accuracy, mean = 45.5%). Sixth, we added pink noise to all our auditory stimuli  
410 whereas no noises were added to Beauchamp et al. (2010)’s auditory stimuli. Our decision to add  
411 noise was based on the prior literature showing that dependence on visual speech information  
412 increases with introduction of noise (Alsius et al., 2016; Buchan et al., 2008; Stacey et al., 2020).  
413 The addition of pink noise may in part have aided in generating McGurk percepts in our  
414 participants. Lastly, the two studies differed slightly in the TMS threshold used to apply single  
415 pulse stimulation. Whereas Beauchamp et al. (2010) used 100% of the resting motor threshold  
416 (RMT) for pulses, we used a slightly higher threshold at 110% of the RMT. We opted for the higher  
417 threshold as 110% or 120% of RMT is the more widely reported approach found in TMS literature  
418 (Cuypers et al., 2014; Kallioniemi & Julkunen, 2016; Sondergaard et al., 2021) which would  
419 naturally be expected to produce greater disruption of the involved region.

420         Given these differences in the study designs, it is possible that the disparity in the effect  
421 sizes of the pSTS stimulation on the frequency of McGurk effect may have been driven by a few  
422 of these factors. Specifically, we speculate that the largest driver of the disparity is the difference  
423 in the exclusion criteria (which were accompanied by other task designs to ensure that subjects  
424 still get fusion percepts). By broadening the subject pool, we similarly broaden our inference

425 beyond the individuals who almost always experience the McGurk effect with a particular stimulus  
426 pairing. Because our pool of subjects is less likely to experience the McGurk percept compared  
427 to those from Beauchamp et al. (2010), it is possible that the effect of pSTS stimulation is less  
428 pronounced as the integration of the incongruent stimuli does not always occur as is the case for  
429 Beauchamp et al. (2010)'s subjects. In addition, it is also widely accepted that the likelihood of  
430 the McGurk effect varies largely across the stimuli used (Basu Mallick et al., 2015; Beauchamp  
431 et al., 2010). It is plausible that the audiovisual phoneme pairing used in Beauchamp et al. (2010)  
432 elicits a stronger McGurk percept compared to the audiovisual word pairings used in our study.

433 Taken together, our data points to evidence that audiovisual speech integration is not  
434 exclusively dependent on a single major hub in the left pSTS; instead, the left pSTS is more  
435 important for the generation of McGurk perception, resolving the conflict between auditory and  
436 visual information so that the information can be perceived as a single percept, rather than two  
437 mismatching percepts. While pSTS has been dubbed the multisensory hub of the brain and is  
438 indeed necessary for certain facets of multisensory perception, the importance of this region has  
439 likely been inflated due to the field's heavy reliance on McGurk stimuli in the study of audiovisual  
440 integration. Indeed, our data provides converging and complementary evidence with the growing  
441 number of both behavioral and electrophysiological studies (Arnal et al., 2009; Arnal et al., 2011;  
442 Eskelund et al., 2011; Faivre et al., 2014; Fingelkurts et al., 2003; Lange et al., 2013; Palmer &  
443 Ramsey, 2012; Roa Romero et al., 2015) that point to the existence of multiple processing  
444 pathways and question the generalizability of McGurk perception to more naturalistic audiovisual  
445 speech processing. Additionally, the right pSTS or other contextual feedback mechanisms could  
446 have supported intact congruent audiovisual benefits after temporary disruption of the left pSTS.

447 Collectively, this data is consistent with the emerging viewpoint that two distinct neural  
448 pathways underlie congruent audiovisual processes responsible for speech enhancements and  
449 incongruent audiovisual processing responsible for McGurk processing. Specifically, one that

450 enhances the initial encoding of auditory information based on the information passed by the  
451 visual cues and another that modifies auditory representation based on the integrated audiovisual  
452 information. The first early feedforward process may occur early in the processing stream and  
453 align auditory encoding with the temporal and acoustic features of the accompanying visual input  
454 (Arnal et al., 2009; Arnal et al., 2011; Van Wassenhove et al., 2005). The later feedback process  
455 is engaged following the detection of mismatch between the auditory and visual cues, with the  
456 brain subsequently altering and adjusting the processing of the unisensory speech based on the  
457 combined audiovisual information (Arnal et al., 2011; Kayser & Logothetis, 2009; Olasagasti et  
458 al., 2015). Given that this late feedback process is facilitated by higher order areas like pSTS, the  
459 limited reliance of this pathway during congruent audiovisual processing can explain why  
460 stimulation to pSTS shows limited disruption on the audiovisual speech enhancement benefits.  
461 Importantly, however, future research should aim to identify a double dissociation of sites  
462 responsible for congruent facilitation versus McGurk effects.

463         Importantly, while this study revealed a significant interaction between stimulation site and  
464 congruity, such that pSTS stimulation affected the McGurk effect but not congruent audiovisual  
465 benefits, the overall pattern of results was weaker than expected and warrants future replications.  
466 In particular, accuracy in the congruent condition approached ceiling which makes it more difficult  
467 to detect TMS stimulation effects. Nevertheless, there was no effect of stimulation on the strong  
468 reaction time benefits present in the audiovisual conditions, emphasizing that the pSTS does not  
469 appear responsible for audiovisual speeding. Moreover, as this was a within-subject, within-  
470 session study in which we observed effects on McGurk fusion rates, we would expect to observe  
471 some changes in congruent audiovisual condition if they were present.

472         In summary, our data demonstrate that while TMS to the left pSTS can limit audiovisual  
473 speech integration and result in a weaker McGurk effect, it does not universally reduce the  
474 ecologically important benefits of congruent visual information on speech perception. This

22 | Ahn, et al.

475 suggests a dissociation in neural mechanisms such that the pSTS reflects only one of multiple

476 critical areas necessary for audiovisual speech interactions.

477

478

## 479 References

- 480 Alsius, A., Navarra, J., & Soto-Faraco, S. (2007). Attention to touch weakens audiovisual  
481 speech integration. *Experimental Brain Research*, 183, 399-404.
- 482 Alsius, A., Paré, M., & Munhall, K. G. (2018). Forty years after hearing lips and seeing voices:  
483 the McGurk effect revisited. *Multisensory Research*, 31(1-2), 111-144.
- 484 Alsius, A., Wayne, R. V., Paré, M., & Munhall, K. G. (2016). High visual resolution matters in  
485 audiovisual speech perception, but only for some. *Attention, Perception, &*  
486 *Psychophysics*, 78, 1472-1487.
- 487 Arnal, L. H., Morillon, B., Kell, C. A., & Giraud, A.-L. (2009). Dual neural routing of visual  
488 facilitation in speech processing. *Journal of Neuroscience*, 29(43), 13445-13453.
- 489 Arnal, L. H., Wyart, V., & Giraud, A.-L. (2011). Transitions in neural oscillations reflect prediction  
490 errors generated in audiovisual speech. *Nature neuroscience*, 14(6), 797-801.
- 491 Basu Mallick, D., F Magnotti, J., & S Beauchamp, M. (2015). Variability and stability in the  
492 McGurk effect: contributions of participants, stimuli, time, and response type.  
493 *Psychonomic bulletin & review*, 22, 1299-1307.
- 494 Beauchamp, M. S., Nath, A. R., & Pasalar, S. (2010). fMRI-guided transcranial magnetic  
495 stimulation reveals that the superior temporal sulcus is a cortical locus of the McGurk  
496 effect. *Journal of Neuroscience*, 30(7), 2414-2417.
- 497 Benoit, M. M., Raij, T., Lin, F. H., Jääskeläinen, I. P., & Stufflebeam, S. (2010). Primary and  
498 multisensory cortical activity is correlated with audiovisual percepts. *Human brain*  
499 *mapping*, 31(4), 526-538.
- 500 Bernstein, L. E., Lu, Z.-L., & Jiang, J. (2008). Quantified acoustic–optical speech signal  
501 incongruity identifies cortical sites of audiovisual speech processing. *Brain research*,  
502 1242, 172-184.
- 503 Besle, J., Fischer, C., Bidet-Caulet, A., Lecaiguard, F., Bertrand, O., & Giard, M.-H. (2008).  
504 Visual activation and audiovisual interactions in the auditory cortex during speech  
505 perception: intracranial recordings in humans. *Journal of Neuroscience*, 28(52), 14301-  
506 14310.
- 507 Brainard, D. H., & Vision, S. (1997). The psychophysics toolbox. *Spatial vision*, 10(4), 433-436.
- 508 Brang, D. (2019). The stolen voice illusion. *Perception*, 48(8), 649-667.
- 509 Brang, D., Plass, J., Kakaizada, S., & Hervey-Jumper, S. L. (2020). Auditory-Visual Speech  
510 Behaviors are Resilient to Left pSTS Damage. *bioRxiv*, 2020.2009. 2026.314799.
- 511 Brown, J. W., & Braver, T. S. (2005). Learned predictions of error likelihood in the anterior  
512 cingulate cortex. *Science*, 307(5712), 1118-1121.
- 513 Buchan, J. N., Paré, M., & Munhall, K. G. (2008). The effect of varying talker identity and  
514 listening conditions on gaze behavior during audiovisual speech perception. *Brain*  
515 *research*, 1242, 162-171.
- 516 Campbell, R. (2008). The processing of audio-visual speech: empirical and neural bases.  
517 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 363(1493),  
518 1001-1010.
- 519 Cohen, J. (1988). edition 2. Statistical power analysis for the behavioral sciences. In: Hillsdale.  
520 Erlbaum.
- 521 Cuyppers, K., Thijs, H., & Meesen, R. L. (2014). Optimization of the transcranial magnetic  
522 stimulation protocol by defining a reliable estimate for corticospinal excitability. *PLoS one*,  
523 9(1), e86380.
- 524 Du, Y., Buchsbaum, B. R., Grady, C. L., & Alain, C. (2014). Noise differentially impacts  
525 phoneme representations in the auditory and speech motor systems. *Proceedings of the*  
526 *National Academy of Sciences*, 111(19), 7126-7131.

- 527 Du, Y., Buchsbaum, B. R., Grady, C. L., & Alain, C. (2016). Increased activity in frontal motor  
528 cortex compensates impaired speech perception in older adults. *Nature*  
529 *communications*, 7(1), 12241.
- 530 Eskelund, K., Tuomainen, J., & Andersen, T. S. (2011). Multistage audiovisual integration of  
531 speech: dissociating identification and detection. *Experimental Brain Research*, 208,  
532 447-457.
- 533 Faivre, N., Mudrik, L., Schwartz, N., & Koch, C. (2014). Multisensory integration in complete  
534 unawareness: Evidence from audiovisual congruency priming. *Psychological Science*,  
535 25(11), 2006-2016.
- 536 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G\* Power 3: A flexible statistical  
537 power analysis program for the social, behavioral, and biomedical sciences. *Behavior*  
538 *research methods*, 39(2), 175-191.
- 539 Feldman, J. I., Conrad, J. G., Kuang, W., Tu, A., Liu, Y., Simon, D. M., Wallace, M. T., &  
540 Woynaroski, T. G. (2022). Relations between the McGurk effect, social and  
541 communication skill, and autistic features in children with and without autism. *Journal of*  
542 *Autism and Developmental Disorders*, 52(5), 1920-1928.
- 543 Fingelkurts, A. A., Fingelkurts, A. A., Krause, C. M., Möttönen, R., & Sams, M. (2003). Cortical  
544 operational synchrony during audio–visual speech integration. *Brain and language*,  
545 85(2), 297-312.
- 546 Ganesan, K., Cao, C. Z., Demidenko, M. I., Jahn, A., Stacey, W. C., Wasade, V. S., & Brang, D.  
547 (2022). Auditory cortex encodes lipreading information through spatially distributed  
548 activity. *bioRxiv*, 2022.2011.2011.516209. <https://doi.org/10.1101/2022.11.11.516209>
- 549 Gelder, B. d., Vroomen, J., & Van der Heide, L. (1991). Face recognition and lip-reading in  
550 autism. *European Journal of Cognitive Psychology*, 3(1), 69-86.
- 551 Greenhouse, S. W., & Geisser, S. (1959). On methods in the analysis of profile data.  
552 *Psychometrika*, 24(2), 95-112.
- 553 Gurler, D., Doyle, N., Walker, E., Magnotti, J., & Beauchamp, M. (2015). A link between  
554 individual differences in multisensory speech perception and eye movements. *Attention*,  
555 *Perception, & Psychophysics*, 77, 1333-1341.
- 556 Hallett, M. (2000). Transcranial magnetic stimulation and the human brain. *Nature*, 406(6792),  
557 147-150.
- 558 Hickok, G., & Poeppel, D. (2007). The cortical organization of speech processing. *Nature*  
559 *Reviews Neuroscience*, 8(5), 393-402.
- 560 Hickok, G., Rogalsky, C., Matchin, W., Basilakos, A., Cai, J., Pillay, S., Ferrill, M., Mickelsen, S.,  
561 Anderson, S. W., & Love, T. (2018). Neural networks supporting audiovisual integration  
562 for speech: A large-scale lesion study. *Cortex*, 103, 360-371.
- 563 Irwin, J. R., Frost, S. J., Mencl, W. E., Chen, H., & Fowler, C. A. (2011). Functional activation for  
564 imitation of seen and heard speech. *Journal of neurolinguistics*, 24(6), 611-618.
- 565 Jones, J. A., & Callan, D. E. (2003). Brain activity during audiovisual speech perception: an  
566 fMRI study of the McGurk effect. *Neuroreport*, 14(8), 1129-1133.
- 567 Jung, J., Bungert, A., Bowtell, R., & Jackson, S. R. (2016). Vertex stimulation as a control site  
568 for transcranial magnetic stimulation: a concurrent TMS/fMRI study. *Brain stimulation*,  
569 9(1), 58-64.
- 570 Kallioniemi, E., & Julkunen, P. (2016). Alternative stimulation intensities for mapping cortical  
571 motor area with navigated TMS. *Brain topography*, 29, 395-404.
- 572 Kayser, C., & Logothetis, N. K. (2009). Directed interactions between auditory and superior  
573 temporal cortices and their role in sensory integration. *Frontiers in integrative*  
574 *neuroscience*, 7.
- 575 Kleiner, M., Brainard, D., & Pelli, D. (2007). What's new in Psychtoolbox-3?



- 576 Lange, J., Christian, N., & Schnitzler, A. (2013). Audio–visual congruency alters power and  
577 coherence of oscillatory activity within and between cortical areas. *Neuroimage*, *79*, 111-  
578 120.
- 579 Luo, H., Liu, Z., & Poeppel, D. (2010). Auditory cortex tracks both auditory and visual stimulus  
580 dynamics using low-frequency neuronal phase modulation. *PLoS biology*, *8*(8),  
581 e1000445.
- 582 MacLeod, A., & Summerfield, Q. (1987). Quantifying the contribution of vision to speech  
583 perception in noise. *British journal of audiology*, *21*(2), 131-141.
- 584 McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature*, *264*(5588), 746-  
585 748.
- 586 Mills, K. R., & Nithi, K. A. (1997). Corticomotor threshold to magnetic stimulation: normal values  
587 and repeatability. *Muscle & Nerve: Official Journal of the American Association of*  
588 *Electrodiagnostic Medicine*, *20*(5), 570-576.
- 589 Nath, A. R., & Beauchamp, M. S. (2012). A neural basis for interindividual differences in the  
590 McGurk effect, a multisensory speech illusion. *Neuroimage*, *59*(1), 781-787.
- 591 Nath, A. R., Fava, E. E., & Beauchamp, M. S. (2011). Neural correlates of interindividual  
592 differences in children's audiovisual speech perception. *Journal of Neuroscience*, *31*(39),  
593 13963-13971.
- 594 Olasagasti, I., Bouton, S., & Giraud, A.-L. (2015). Prediction across sensory modalities: A  
595 neurocomputational model of the McGurk effect. *Cortex*, *68*, 61-75.
- 596 Palmer, T. D., & Ramsey, A. K. (2012). The function of consciousness in multisensory  
597 integration. *Cognition*, *125*(3), 353-364.
- 598 Paré, M., Richler, R. C., ten Hove, M., & Munhall, K. (2003). Gaze behavior in audiovisual  
599 speech perception: The influence of ocular fixations on the McGurk effect. *Perception &*  
600 *psychophysics*, *65*, 553-567.
- 601 Pelli, D. G., & Vision, S. (1997). The VideoToolbox software for visual psychophysics:  
602 Transforming numbers into movies. *Spatial vision*, *10*, 437-442.
- 603 Plass, J., Brang, D., Suzuki, S., & Grabowecy, M. (2020). Vision perceptually restores auditory  
604 spectral dynamics in speech. *Proceedings of the National Academy of Sciences*,  
605 *117*(29), 16920-16927.
- 606 Roa Romero, Y., Senkowski, D., & Keil, J. (2015). Early and late beta-band power reflect  
607 audiovisual perception in the McGurk illusion. *Journal of neurophysiology*, *113*(7), 2342-  
608 2350.
- 609 Ross, L. A., Saint-Amour, D., Leavitt, V. M., Molholm, S., Javitt, D. C., & Foxe, J. J. (2007).  
610 Impaired multisensory processing in schizophrenia: deficits in the visual enhancement of  
611 speech comprehension under noisy environmental conditions. *Schizophrenia research*,  
612 *97*(1-3), 173-183.
- 613 Rossini, P. M., & Rossi, S. (2007). Transcranial magnetic stimulation: diagnostic, therapeutic,  
614 and research potential. *Neurology*, *68*(7), 484-488.
- 615 Saatlou, F. H., Rogasch, N. C., McNair, N. A., Biabani, M., Pillen, S. D., Marshall, T. R., &  
616 Bergmann, T. O. (2018). MAGIC: An open-source MATLAB toolbox for external control  
617 of transcranial magnetic stimulation devices. *Brain Stimulation: Basic, Translational, and*  
618 *Clinical Research in Neuromodulation*, *11*(5), 1189-1191.
- 619 Schroeder, C. E., Lakatos, P., Kajikawa, Y., Partan, S., & Puce, A. (2008). Neuronal oscillations  
620 and visual amplification of speech. *Trends in cognitive sciences*, *12*(3), 106-113.
- 621 Seijdel, N., Schoffelen, J.-M., Hagoort, P., & Drijvers, L. (2023). Attention drives visual  
622 processing and audiovisual integration during multimodal communication. *bioRxiv*,  
623 2023.2005. 2011.540320.
- 624 Sekiyama, K., Kanno, I., Miura, S., & Sugita, Y. (2003). Auditory-visual speech perception  
625 examined by fMRI and PET. *Neuroscience research*, *47*(3), 277-287.

- 626 Slavin, R., & Smith, D. (2009). The relationship between sample sizes and effect sizes in  
627 systematic reviews in education. *Educational evaluation and policy analysis*, 31(4), 500-  
628 506.
- 629 Sondergaard, R. E., Martino, D., Kiss, Z. H., & Condliffe, E. G. (2021). TMS motor mapping  
630 methodology and reliability: a structured review. *Frontiers in Neuroscience*, 15, 709368.
- 631 Stacey, J. E., Howard, C. J., Mitra, S., & Stacey, P. C. (2020). Audio-visual integration in noise:  
632 Influence of auditory and visual stimulus degradation on eye movements and perception  
633 of the McGurk effect. *Attention, Perception, & Psychophysics*, 82, 3544-3557.
- 634 Stevenson, R. A., Segers, M., Ferber, S., Barense, M. D., & Wallace, M. T. (2014). The impact  
635 of multisensory integration deficits on speech perception in children with autism  
636 spectrum disorders. In (Vol. 5, pp. 379): Frontiers Media SA.
- 637 Sumbly, W. H., & Pollack, I. (1954). Visual contribution to speech intelligibility in noise. *The*  
638 *journal of the acoustical society of america*, 26(2), 212-215.
- 639 Szycik, G. R., Stadler, J., Tempelmann, C., & Münte, T. F. (2012). Examining the McGurk  
640 illusion using high-field 7 Tesla functional MRI. *Frontiers in Human Neuroscience*, 6, 95.
- 641 Van Engen, K. J., Dey, A., Sommers, M. S., & Peelle, J. E. (2022). Audiovisual speech  
642 perception: Moving beyond McGurk. *The journal of the acoustical society of america*,  
643 152(6), 3216-3225.
- 644 Van Engen, K. J., Xie, Z., & Chandrasekaran, B. (2017). Audiovisual sentence recognition not  
645 predicted by susceptibility to the McGurk effect. *Attention, Perception, & Psychophysics*,  
646 79, 396-403.
- 647 Van Wassenhove, V., Grant, K. W., & Poeppel, D. (2005). Visual speech speeds up the neural  
648 processing of auditory speech. *Proceedings of the National Academy of Sciences*,  
649 102(4), 1181-1186.
- 650 Van Wassenhove, V., Grant, K. W., & Poeppel, D. (2007). Temporal window of integration in  
651 auditory-visual speech perception. *Neuropsychologia*, 45(3), 598-607.
- 652 Walsh, V., & Cowey, A. (2000). Transcranial magnetic stimulation and cognitive neuroscience.  
653 *Nature Reviews Neuroscience*, 1(1), 73-80.
- 654 Williams, J. H., Massaro, D. W., Peel, N. J., Bosseler, A., & Suddendorf, T. (2004). Visual-  
655 auditory integration during speech imitation in autism. *Research in developmental*  
656 *disabilities*, 25(6), 559-575.
- 657 Zhang, J., Meng, Y., He, J., Xiang, Y., Wu, C., Wang, S., & Yuan, Z. (2019). McGurk effect by  
658 individuals with autism spectrum disorder and typically developing controls: A systematic  
659 review and meta-analysis. *Journal of Autism and Developmental Disorders*, 49(1), 34-  
660 43.
- 661