

1 **Full title:** Lip-reading enables the brain to synthesize auditory features of unknown silent  
2 speech.

3 **Running title:** Synthesizing auditory features of silent speech.  
4

5 **Authors:** Mathieu Bourguignon<sup>1,2,3,\*</sup>, Martijn Baart<sup>1,4</sup>, Efthymia C. Kapnoula<sup>1</sup>, Nicola  
6 Molinaro<sup>1,5</sup>  
7

### 8 **Affiliations**

9 <sup>1</sup>BCBL. Basque Center on Cognition, Brain and Language, 20009 San Sebastian, Spain.

10 <sup>2</sup>Laboratoire de Cartographie fonctionnelle du Cerveau, UNI – ULB Neuroscience Institute, Université libre de  
11 Bruxelles (ULB), Brussels, Belgium.

12 <sup>3</sup>Laboratoire Cognition Langage et Développement, UNI – ULB Neuroscience Institute, Université libre de  
13 Bruxelles (ULB), Brussels, Belgium.

14 <sup>4</sup>Department of Cognitive Neuropsychology, Tilburg University, Tilburg, the Netherlands.

15 <sup>5</sup>Ikerbasque, Basque Foundation for Science, Bilbao, Spain.

16 \*Corresponding author. E-mail: mabourgu@ulb.ac.be.  
17

18 **Number of pages:** 47, **number of figures:** 9, **number of tables:** 1, **number of words in**

19 **Abstract:** 250, **number of words in Introduction:** 624, **number of words in Discussion:**  
20 1490.  
21

### 22 **Conflict of interest**

23 The authors have no conflict of interest to declare.  
24

### 25 **Acknowledgment**

26 Mathieu Bourguignon was supported by the Innoviris Attract program (grant 2015-BB2B-  
27 10), by the Spanish Ministry of Economy and Competitiveness (grant PSI2016-77175-P), and  
28 by the Marie Skłodowska-Curie Action of the European Commission (grant 743562). Martijn  
29 Baart was supported by the Netherlands Organization for Scientific Research (NWO, VENI

30 grant 275-89-027). Efthymia C. Kapnoula was supported by the Spanish Ministry of  
31 Economy and Competitiveness, through the Juan de la Cierva-Formación fellowship, and by  
32 the Spanish Ministry of Economy and Competitiveness (grant PSI2017-82563-P). Nicola  
33 Molinaro was supported by the Spanish Ministry of Science, Innovation and Universities  
34 (grant RTI2018-096311-B-I00), the Agencia Estatal de Investigación (AEI), the Fondo  
35 Europeo de Desarrollo Regional (FEDER) and by the Basque government (grant  
36 PI\_2016\_1\_0014). The authors acknowledge financial support from the Spanish Ministry of  
37 Economy and Competitiveness, through the “Severo Ochoa” Programme for Centres/Units of  
38 Excellence in R&D” (SEV-2015-490) awarded to the BCBL.  
39 We thank Riitta Hari at Department of Art (Aalto University School of Arts, Design, and  
40 Architecture, Espoo, Finland) for helpful comments on the manuscript.

41

42 **Abstract**

43 Lip-reading is crucial for understanding speech in challenging conditions. But how  
44 the brain extracts meaning from—silent—visual speech is still under debate. Lip-reading in  
45 silence activates the auditory cortices, but it is not known whether such activation reflects  
46 immediate synthesis of the corresponding auditory stimulus or imagery of unrelated sounds.

47 To disentangle these possibilities, we used magnetoencephalography to evaluate how  
48 cortical activity in 28 healthy adults humans (17 females) entrained to the auditory speech  
49 envelope and lip movements (mouth opening) when listening to a spoken story without visual  
50 input (*audio-only*), and when seeing a silent video of a speaker articulating another story  
51 (*video-only*).

52 In *video-only*, auditory cortical activity entrained to the absent auditory signal at  
53 frequencies below 1 Hz more than to the seen lip movements. This entrainment process was  
54 characterized by an auditory-speech-to-brain delay of ~70 ms in the left hemisphere,  
55 compared to ~20 ms in *audio-only*. Entrainment to mouth opening was found in the right  
56 angular gyrus at below 1 Hz, and in early visual cortices at 1–8 Hz.

57 These findings demonstrate that the brain can use a silent lip-read signal to synthesize  
58 a coarse-grained auditory speech representation in early auditory cortices. Our data indicate  
59 the following underlying oscillatory mechanism: Seeing lip movements first modulates  
60 neuronal activity in early visual cortices at frequencies that match articulatory lip movements;  
61 the right angular gyrus then extracts slower features of lip movements, mapping them onto  
62 the corresponding speech sound features; this information is fed to auditory cortices, most  
63 likely facilitating speech parsing.

64

65 **Significance statement**

66 Lip-reading consists in decoding speech based on visual information derived from  
67 observation of a speaker's articulatory facial gestures. Lip reading is known to improve  
68 auditory speech understanding, especially when speech is degraded. Interestingly, lip-reading  
69 in silence still activates the auditory cortices, even when participants do not know what the  
70 absent auditory signal should be. However, it was uncertain what such activation reflected.  
71 Here, using magnetoencephalographic recordings, we demonstrate it reflects fast synthesis of  
72 the auditory stimulus rather than mental imagery of unrelated—speech or non-speech—  
73 sounds. Our results also shed light on the oscillatory dynamics underlying lip-reading.

74

75

76 **Keywords**

77 Lip reading; silent speech; audiovisual integration; speech entrainment;  
78 magnetoencephalography

79

## 80 **Introduction**

81 In everyday situations, seeing a speaker's articulatory mouth gestures, here referred to  
82 as lip-reading or visual speech, can help us decode the auditory speech signal (Sumbly and  
83 Pollack, 1954). In fact, lip movements are intelligible even without an auditory signal, likely  
84 because there is a strong connection between auditory and visual speech (Munhall and  
85 Vatikiotis-Bateson, 2004; Chandrasekaran et al., 2009). It is however not clear how the brain  
86 extracts meaning from visual speech.

87 Some evidence points to the possibility that visual speech is recoded into acoustic  
88 information. For example, seeing silent visual speech clips of simple speech sounds such as  
89 vowels or elementary words activates auditory cortical areas (Calvert et al., 1997; Pekkola et  
90 al., 2005), even when participants are not aware of what the absent auditory input should be  
91 (Calvert et al., 1997; Bernstein et al., 2002; Paulesu et al., 2003). However, recoding visual  
92 speech into an acoustic representation (here referred to as synthesis) is computationally  
93 demanding. It has therefore been suggested that meaning is directly extracted from visual  
94 speech within visual areas and heteromodal association cortices (Bernstein and Liebenthal,  
95 2014; O'Sullivan et al., 2016; Lazard and Giraud, 2017; Hauswald et al., 2018). According to  
96 this view, activation in early auditory cortices driven by lip reading might reflect imagery of  
97 unrelated—speech—sounds (Bernstein and Liebenthal, 2014), but not a direct recoding of  
98 visual speech into its corresponding acoustic representation. As previous work has relied on  
99 time-insensitive neuroimaging techniques (Calvert et al., 1997; Bernstein et al., 2002;  
100 Paulesu et al., 2003; Pekkola et al., 2005), there was no empirical evidence to disentangle  
101 these two alternatives. Here, we took advantage of auditory cortical entrainment to look for  
102 decisive evidence to support the existence of a synthesis mechanism whereby visual speech is  
103 recoded into its corresponding auditory information.

104           When people listen to continuous natural speech, oscillatory cortical activity  
105 synchronises with the auditory temporal speech envelope (Luo and Poeppel, 2007;  
106 Bourguignon et al., 2012; Gross et al., 2013; Peelle et al., 2013; Molinaro et al., 2016; Vander  
107 Ghinst et al., 2016; Meyer et al., 2017; Meyer and Gumbert, 2018). Such “speech-brain  
108 entrainment” originates mainly in auditory cortices at frequencies matching phrase (below 1  
109 Hz) and syllable rates (4–8 Hz), and is thought to be essential for speech comprehension  
110 (Ahissar et al., 2001; Luo and Poeppel, 2007; Peelle et al., 2013; Ding et al., 2016; Meyer et  
111 al., 2017). An electroencephalography study suggested that silent lip-read information  
112 entrains cortical activity at syllable rate when participants are highly familiar with speech  
113 content (Crosse et al., 2015). However, since participants knew what the absent speech sound  
114 should be in this study, it remains unclear whether entrainment is driven by the (i) lip-read  
115 information, (ii) covert production or repetition of the speech segment, (iii) top-down lexical  
116 and semantic processes, or (iv) some combination of these factors.

117           Here, we address the following critical question: does the brain use lip-read input to  
118 bring auditory cortices to entrain to the audio speech signal even when there is no physical  
119 speech sound and participants do not know the content of the absent auditory signal? To do  
120 so, we evaluated entrainment to a spoken story without visual input (*audio-only*), and  
121 compared these data to a silent condition with a video of a speaker articulating another story  
122 (*video-only*). To determine the ‘lip-read specificity’ of these entrainment patterns, we also  
123 included a condition in which the mouth configuration of the speaker telling another story  
124 was transduced into a dynamic luminance contrast (*control-video-only*). If the brain can  
125 synthesize properties of missing speech based on concurrent lip-reading in a timely manner,  
126 auditory cortical entrainment with the envelope of the audio signal should be similar in  
127 *audio-only* and *video-only*, even if the speech sound was not physically present in the latter  
128 condition.

129

## 130 **Materials and Methods**

### 131 *Participants*

132 Twenty-eight healthy human adults (17 females) aged  $24.1 \pm 4.0$  years (mean  $\pm$  SD)  
133 were included in the study. All reported being native speakers of Spanish and right-handed.  
134 They had normal or corrected-to-normal vision and normal hearing, had no prior history of  
135 neurological or psychiatric disorders, and were not taking any medication or substance that  
136 could influence the nervous system.

137 The experiment was approved by the BCBL Ethics Review Board and complied with  
138 the guidelines of the Helsinki Declaration. Written informed consent was obtained from all  
139 participants prior to testing.

### 140 *Experimental paradigm*

141 Figure 1 presents stimulus examples and excerpts. The stimuli were derived from 8  
142 audio-visual recordings of a female native Spanish speaker talking for 5 min about a given  
143 topic (animals, books, food, holidays, movies, music, social media, and sports). Video and  
144 audio were simultaneously recorded using a digital camera (Canon Legria HF G10) with an  
145 internal microphone. Video recordings were framed as head shots, and recorded at the PAL  
146 standard of 25 frames per second (videos were  $1920 \times 1080$  pixels in size, 24 bits/pixel, with  
147 an auditory sampling rate of 44100 Hz). The camera was placed  $\sim 70$  cm away from the  
148 speaker, and the face spanned about half of the vertical field of view. Final images were  
149 resized to a resolution of  $1024 \times 768$  pixels.

150 For each video, a “control” video was created in which mouth movements were  
151 transduced into luminance changes (Fig. 1C). To achieve this we extracted lip contours from  
152 each individual frame of the video recordings with an in-house Matlab code based on the

153 approach of Eveno et al. (2004). In the control video, the luminance of a Greek cross changed  
154 according to mouth configuration (Fig. 1C). Its size ( $300 \times 300$  pixels) was roughly matched  
155 with the extent of the eyes and mouth, which are the parts of the face people tend to look at  
156 when watching a speaker's face (Vatikiotis-Bateson et al., 1998). Mouth configuration  
157 variables (mouth opening, width, and surface) were rescaled so that their 1<sup>st</sup> and 99<sup>th</sup>  
158 percentiles corresponded to the minimum and maximum luminance levels. The center of the  
159 cross encoded the mouth surface area, its top and bottom portions encoded mouth opening,  
160 and its left- and rightmost portions encoded mouth width. In this configuration, the three  
161 represented parameters were spatially and temporally congruent with the portion of the mouth  
162 they parametrized. All portions were smoothly connected by buffers along which the weight  
163 of the encoded parameters varied as a squared cosine. These control videos were designed to  
164 determine if effects were specific to lip-reading. The transduced format was preferred to other  
165 classical controls such as meaningless lip movements or gum-chewing motions because  
166 preserved the temporal relation between the visual input and underlying speech sounds.

167 For each sound recording, we derived a non-speech "control" audio consisting of  
168 white noise modulated by the auditory speech envelope. These control sounds were designed  
169 to determine whether uncovered effects were specific to speech. However, conditions that  
170 included these control sounds were not analyzed because they were uninformative about lip-  
171 reading driven oscillatory entrainment.

172 In total, participants completed 10 experimental conditions while sitting with their  
173 head in a MEG helmet. This included all 9 possible combinations of 3 types of visual stimuli  
174 (original, control, no video) and 3 types of audio stimuli (original, control, no audio). The test  
175 condition with no audio and no video was trivially labeled as the *rest* condition and lasted 5  
176 min. Each of the other 8 conditions was assigned to 1 of the 8 stories (condition-story  
177 assignment counterbalanced across participants). In this way, we ensured that each condition



178 was presented continuously for 5 min, and that the same story was never presented twice. The  
179 tenth condition was a localizer condition in which participants attended 400-Hz pure tones  
180 and checkerboard pattern reversals lasting 10 min. This condition is not analyzed in this  
181 paper. All conditions were presented in random order, separated by short breaks. Videos were  
182 shown on a back-projection screen (videos were 41 cm × 35 cm in size) placed in front of the  
183 participants at a distance of ~1 m. Sounds were delivered at 60 dB (measured at ear-level)  
184 through a front-facing speaker (Panphonics Oy, Espoo, Finland) placed ~1 m behind the  
185 screen. Participants were instructed to watch the videos and listen to the sounds attentively.

186 To investigate our research hypotheses, we focussed on the following conditions: 1)  
187 the original speech audio with no video, referred to as *audio-only*, 2) the original video with  
188 no audio, referred to as *video-only*, 3) the control video with no audio, referred to as the  
189 *control-video-only*, and 4) the *rest*.

## 190 ***Data acquisition***

191 Neuromagnetic signals were acquired with a whole-scalp-covering  
192 neuromagnetometer (Vectorview; Elekta Oy, Helsinki, Finland) in a magnetically shielded  
193 room. The recording pass-band was 0.1–330 Hz and the signals were sampled at 1 kHz. The  
194 head position inside the MEG helmet was continuously monitored by feeding current to 4  
195 head-tracking coils located on the scalp. Head position indicator coils, three anatomical  
196 fiducials, and at least 150 head-surface points (covering the whole scalp and the nose surface)  
197 were localized in a common coordinate system using an electromagnetic tracker (Fastrak,  
198 Polhemus, Colchester, VT, USA).

199 Eye movements were tracked with an MEG-compatible eye tracker (EyeLink 1000  
200 Plus, SR Research). Participants were calibrated using the standard 9-point display and  
201 monocular eye movements were recorded at a sampling rate of 1 kHz. Eye-movements were  
202 recorded for the duration of all experimental conditions.

203 High-resolution 3D-T1 cerebral magnetic resonance images (MRI) were acquired on a  
204 3 Tesla MRI scan (Siemens Medical System, Erlangen, Germany) facility available at the  
205 BCBL.

### 206 *MEG preprocessing*

207 Continuous MEG data were first preprocessed off-line using the temporal signal space  
208 separation method (correlation coefficient, 0.9; segment length, 10 s) to suppress external  
209 sources of interference and to correct for head movements (Taulu et al., 2005; Taulu and  
210 Simola, 2006). To further suppress heartbeat, eye-blink, and eye-movement artifacts, 30  
211 independent components (Vigário et al., 2000; Hyvärinen et al., 2004) were evaluated from  
212 the MEG data low-pass filtered at 25 Hz using FastICA algorithm (dimension reduction, 30;  
213 non-linearity, tanh). Independent components corresponding to such artifacts were identified  
214 based on their topography and time course and were removed from the full-rank MEG  
215 signals.

### 216 *Coherence analysis*

217 Coherence was estimated between MEG signals and 1) the auditory speech temporal  
218 envelope, 2) mouth opening, 3) mouth width, and 4) mouth surface. The auditory speech  
219 temporal envelope was obtained by summing the Hilbert envelope of the auditory speech  
220 signal filtered through a third octave filter bank (central frequency ranging linearly on a log-  
221 scale from 250 Hz to 1600 Hz; 19 frequency bands), and was further resampled to 1000 Hz  
222 time-locked to the MEG signals (Fig. 1B). Continuous data from each condition were split  
223 into 2-s epochs with 1.6-s epoch overlaps, affording a spectral resolution of 0.5 Hz while  
224 decreasing noise on coherence estimates (Bortel and Sovka, 2014). MEG epochs exceeding 5  
225 pT (magnetometers) or 1 pT/cm (gradiometers) were excluded from further analyses to avoid  
226 data contamination by artifact sources that had not been suppressed by the temporal signal

227 space separation or removed with independent component analysis. These steps led to an  
228 average of 732 artifact-free epochs across participants and conditions ( $SD = 36$ ). A one-way  
229 repeated measures ANOVA revealed no differences between conditions ( $F_{2,54} = 1.07$ ,  $p =$   
230  $0.35$ ). Next, we estimated sensor-level coherence (Halliday, 1995) and combined gradiometer  
231 pairs based on the direction of maximum coherence (Bourguignon et al., 2015). Only values  
232 from these gradiometer pairs are presented in the results.

233 In coherence analyses, we focused on four frequency ranges (0.5 Hz, 1–3 Hz, 2–5 Hz,  
234 and 4–8 Hz) by averaging coherence across the frequency bins they encompassed. The 2–5-  
235 Hz, and 4–8-Hz frequency ranges were well matched to the count rate of words ( $3.34 \pm 0.12$   
236 Hz; mean  $\pm$  SD across the 8 videos) and syllables ( $5.91 \pm 0.12$  Hz), while the count rate of  
237 phrases ( $1.01 \pm 0.20$  Hz) fell in between the two lowest ranges. As in a previous study  
238 (Vander Ghinst et al., 2019), rates were assessed as the number of phrases, words, or  
239 syllables manually extracted from audio recordings divided by the corrected duration of the  
240 audio recording. For phrases, the corrected duration was trivially the total duration of the  
241 audio recording. For words and syllables, the corrected duration was the total time during  
242 which the talker was actually talking, that is the total duration of the audio recording (here 5  
243 min) minus the sum of all silent periods when the auditory speech envelope was below a  
244 tenth of its mean for at least 100 ms. Note that setting the threshold for the duration defining  
245 a silent period to a value obviously too low (10 ms) or too high (500 ms) changed the  
246 estimates of word and syllable count rates by only  $\sim 10\%$ . These frequency ranges were  
247 selected also because auditory speech entrainment dominates at 0.5 Hz and 4–8 Hz (Luo and  
248 Poeppel, 2007; Bourguignon et al., 2012; Gross et al., 2013; Peelle et al., 2013; Molinaro et  
249 al., 2016; Vander Ghinst et al., 2016; Meyer et al., 2017; Meyer and Gumbert, 2018) but is  
250 also present at intermediate frequencies (Keitel et al., 2018), and because lip entrainment has

251 previously been identified at 2–5 Hz (Park et al., 2016; Giordano et al., 2017). Coherence  
252 maps were also averaged across participants for illustration purposes.

253 We only report coherence estimated between MEG signals and 1) the auditory speech  
254 envelope and 2) mouth opening. Although tightly related, the two latter signals displayed  
255 only a moderate degree of coupling, that peaked at 0.5 Hz, and 4–8 Hz (Fig. 2A<sub>i</sub>), with a  
256 visual–to–auditory speech delay of ~120-ms (maximum cross-correlation between auditory  
257 speech envelope and mouth opening; Fig. 2B<sub>i</sub>). Mouth opening and mouth surface were  
258 coherent at > 0.7 across the 0–10-Hz range (Fig. 2A<sub>ii</sub>) and yielded similar results. Mouth  
259 width displayed a moderate level of coherence with mouth opening (Fig. 2A<sub>iii</sub>) and an unclear  
260 visual–to–auditory speech delay (Fig. 2B<sub>ii</sub>). Mouth width was not included in the main  
261 analyses because it led to lower coherence values with MEG signals than mouth opening, but  
262 was retained as a nuisance factor in the partial coherence analyses (see below).

263 It is worth noting that the magnitude of the coupling between the auditory speech  
264 envelope and mouth opening (as assessed by coherence) we report for our audio-visual  
265 stimuli was 2–3 times lower than that reported elsewhere (Park et al., 2016; Hauswald et al.,  
266 2018). To ensure that this discrepancy was not due to the inadequacy of our lip-extraction  
267 procedure, we compared our time-series of mouth parameters to those extracted using a deep-  
268 learning-based solution (Visage Technology; face tracking and analysis). This revealed a  
269 good correspondence between the estimated time-series for mouth opening ( $r = 0.95 \pm 0.01$ ;  
270 mean  $\pm$  SD across the 8 videos), mouth width ( $r = 0.88 \pm 0.01$ ), and mouth surface ( $r = 0.95$   
271  $\pm 0.01$ ). The genuine difference between the level of audio-visual speech coupling found in  
272 our study compared to others might be due to the language used (Spanish here vs. English  
273 elsewhere), or to the idiosyncrasies of our talker. Nevertheless, this relative decoupling  
274 between audio- and visual speech signals provided an opportunity to separate their respective  
275 cortical representations more efficiently.

276 Coherence was also estimated at the source level. To do so, individual MRIs were  
277 first segmented using the Freesurfer software (Reuter et al., 2012; RRID:SCR\_001847).  
278 Then, the MEG forward model was computed using the Boundary Element Method  
279 implemented in the MNE software suite (Gramfort et al., 2014; RRID:SCR\_005972) for three  
280 orthogonal tangential current dipoles (corresponding to the 3 spatial dimensions) placed on a  
281 homogeneous 5-mm grid source space covering the whole brain. At each source, the forward  
282 model was further reduced to its two first principal components, which closely corresponded  
283 to sources tangential to the skull; the discarded component corresponded to the radial source  
284 which is close to magnetically silent. Coherence maps were produced within the computed  
285 source space at 0.5 Hz, 1–3 Hz, 2–5 Hz, and 4–8 Hz using a linearly constrained minimum  
286 variance beamformer built based on the *rest* data covariance matrix (Van Veen et al., 1997;  
287 Hillebrand and Barnes, 2005). Source maps were then interpolated to a 1-mm homogenous  
288 grid and smoothed with a Gaussian kernel of 5 mm full-width-at-half-maximum. Both planar  
289 gradiometers and magnetometers were used for inverse modeling after dividing each sensor  
290 signal (and the corresponding forward-model coefficients) by the standard deviation of its  
291 noise. The noise variance was estimated from the continuous *rest* MEG data band-passed  
292 through 1–195 Hz, for each sensor separately.

293 Coherence maps were also produced at the group level. A non-linear transformation  
294 from individual MRIs to the MNI brain was first computed using the spatial normalization  
295 algorithm implemented in Statistical Parametric Mapping (SPM8; Ashburner et al., 1997;  
296 Ashburner and Friston, 1999; RRID:SCR\_007037) and then applied to individual MRIs and  
297 coherence maps. This procedure generated a normalized coherence map in the MNI space for  
298 each subject and frequency range. Coherence maps were then averaged across participants.

299 Individual and group-level coherence maps for the auditory speech envelope (mouth  
300 opening, respectively) were also estimated after controlling for mouth opening and mouth

301 width (the auditory speech envelope, respectively) using partial coherence (Halliday, 1995).  
302 Partial coherence is the direct generalization of partial correlation (Kendall and Stuart, 1968)  
303 to the frequency domain (Halliday, 1995).

304 The same approach was used to estimate coherence between MEG (in the sensor and  
305 source space) and global changes (or edges) in the visual stimulus, and to partial out such  
306 “global visual change” from coherence maps for the auditory speech envelope. The global  
307 visual change signal was computed at every video frame as the sum of squares of the  
308 difference between that frame and the previous frame, divided by the sum of squares of the  
309 previous frame. This signal predominantly identified edges corresponding to periods when  
310 the speaker moved her head, eyebrows and jaw (see Fig. 3). The rationale being that these  
311 periods may tend to co-occur with the onset of phrases and sentences (Munhall et al., 2004)  
312 and could modulate oscillatory activity in auditory cortices (Schroeder et al., 2008).

313 Finally, individual and group-level coherence maps for the auditory speech envelope  
314 in *video-only* were estimated after shifting the auditory speech envelope by ~30 s, ~60 s, ...  
315 ~240 s, and ~270 s. For each subject and time-shift, the exact time-shift applied was selected  
316 within a  $\pm 10$  s window around the target time-shift, at the silent period for which the auditory  
317 speech envelope smoothed with a 1-s square kernel was at the minimum. Ensuing values of  
318 coherence were used to rule out the possibility that coherence with the genuine auditory  
319 speech envelope results from general temporal characteristics of auditory speech.

### 320 ***Estimation of temporal response functions***

321 We used temporal response functions (TRFs) to model how the auditory speech  
322 envelope affected the temporal dynamics of auditory cortical activity. Based on our results,  
323 TRFs were estimated only for the 0.2–1.5-Hz frequency range, in the *audio-only* and *video-*  
324 *only* conditions. A similar approach has been used to model brain responses to speech at 1–8  
325 Hz (Lalor and Foxe, 2010; Zion Golumbic et al., 2013), and to model brain responses to

326 natural force fluctuations occurring during maintenance of constant hand grip contraction  
327 (Bourguignon et al., 2017b). TRFs are the direct analogue of evoked responses in the context  
328 of continuous stimulation.

329 We used the mTRF toolbox (Crosse et al., 2016) to estimate the TRF of auditory  
330 cortical activity associated with the auditory speech envelope. In all conditions, source  
331 signals were reconstructed at individual coordinates of maximum 0.5-Hz coherence with the  
332 auditory speech envelope in *audio-only*. These two-dimensional source signals were  
333 projected onto the orientation that maximized the coherence with the auditory speech  
334 envelope at 0.5 Hz. Then, the source signal was filtered at 0.2–1.5 Hz, the auditory speech  
335 envelope was convolved with a 50-ms square smoothing kernel and both were down-sampled  
336 to 20 Hz (note that for auditory speech envelope, this procedure is equivalent to taking the  
337 mean over samples 25 ms around sampling points). For each subject, the TRFs were modeled  
338 from –1.5 s to +2.5 s, for a fixed set of ridge values ( $\lambda = 2^0, 2^1, 2^2 \dots 2^{20}$ ). We adopted the  
339 following 10-fold cross-validation procedure to determine the optimal ridge value: For each  
340 subject, TRFs were estimated based on 90% of the data, and used to predict the 10% of data  
341 left out and the Pearson correlation was then estimated between predicted and measured  
342 signals. The square of the mean correlation value across the 10 runs provided an estimate of  
343 the proportion of variance explained by entrainment to the auditory speech envelope. TRFs  
344 were recomputed based on all the available data for the ridge value maximizing the mean  
345 explained variance. To deal with sign ambiguity, the polarity of each TRF was adapted so  
346 that correlation with the first singular vector of all subjects' TRF in the range –0.5 s to 1.0 s is  
347 positive.

348 Based on our results, the TRF framework was also used to model brain responses to  
349 mouth opening and the global visual change signal at 0.2–1.5 Hz and mouth opening at 2–5  
350 Hz, and to model the evolution of the auditory speech envelope at 0.2–1.5 Hz associated with

351 the time course of (i) mouth opening, (ii) global visual change, and (iii) the Hilbert envelope  
352 of mouth opening in the 2–5-Hz band. Note that the last TRF seeks phase–amplitude  
353 coupling between auditory speech envelope at 0.2–1.5 Hz (phase) and mouth opening at 2–5  
354 Hz (amplitude), with the—perhaps not that common—perspective that the amplitude signal  
355 drives the phase signal. We used exactly the same parameters as reported above, except the  
356 data for the brain response to mouth opening at 2–5 Hz where were downsampled to 50 Hz  
357 and modeled from –0.7 to 1.2 s.

### 358 *Eye-tracking data*

359 As in previous studies using eye-tracking (McMurray et al., 2002; Kapnoula et al.,  
360 2015), eye-movements were automatically parsed into saccades and fixations using default  
361 psychophysical parameters. Adjacent saccades and fixations were combined into a single  
362 “look” that started at the onset of the saccade and ended at the offset of the fixation.

363 A region of interest was identified for each of the three critical objects: mouth and  
364 eyes in *video-only* and flickering cross in *control-video-only* (Fig. 4). In converting the  
365 coordinates of each look to the object being fixated, the boundaries of the regions of interest  
366 were extended by 50 pixels in order to account for noise and/or head-drift in the eye-tracking  
367 record. This did not result in any overlap between the eye and mouth regions.

368 Based on these regions of interest, we estimated the proportion of eye fixation to the  
369 combined regions of interest encompassing eyes and mouth in *video-only* and flickering cross  
370 in *control-video-only*. Eyes and mouth regions were combined because these are the parts of  
371 the face people tend to look at when watching a talking face (Vatikiotis-Bateson et al., 1998).  
372 Importantly, even when people are looking at the eyes, lip movements—in the periphery of  
373 the field of view—still benefit speech perception (Paré et al., 2003; Kaplan and Jesse, 2019).  
374 The two resulting areas were of comparable size: 100,800 pixels for the flickering cross vs.  
375 77,300 pixels for the eyes and mouth. Data from one participant were excluded due to



376 technical issues during acquisition, and eye fixation analyses were thus based on data from 27  
377 participants.

### 378 ***Experimental design and statistical analyses***

379 Sample size was based on previous studies reporting entrainment to lip movements,  
380 which included 46 (Park et al., 2016) and 19 (Giordano et al., 2017) healthy adults.

381 The statistical significance of the local coherence maxima observed in group-level  
382 maps was assessed with a non-parametric permutation test that intrinsically corrects for  
383 multiple spatial comparisons (Nichols and Holmes, 2002). Subject- and group-level *rest*  
384 coherence maps were computed in a similar way to the *genuine* maps; MEG signals were  
385 replaced by *rest* MEG signals while auditory/visual speech signals were identical. Group-  
386 level difference maps were obtained by subtracting *genuine* and *rest* group-level coherence  
387 maps. Under the null hypothesis that coherence maps are the same irrespective of the  
388 experimental condition, *genuine* and *rest* labels should be exchangeable at the subject-level  
389 prior to computing the group-level difference map (Nichols and Holmes, 2002). To reject this  
390 hypothesis and to compute a threshold of statistical significance for the correctly labeled  
391 difference map, the permutation distribution of the maximum of the difference map's  
392 absolute value was computed for a subset of 1000 permutations. The threshold at  $p < 0.05$   
393 was computed as the 95<sup>th</sup> percentile of the permutation distribution (Nichols and Holmes,  
394 2002). Permutation tests can be too conservative for voxels other than the one with the  
395 maximum observed statistic (Nichols and Holmes, 2002). For example, dominant coherence  
396 values in the right auditory cortex could bias the permutation distribution and overshadow  
397 weaker coherence values in the left auditory cortex, even if these were highly consistent  
398 across subjects. Therefore, the permutation test described above was conducted separately for  
399 left- and right-hemisphere voxels. All supra-threshold local coherence maxima were

400 interpreted as indicative of brain regions showing statistically significant coupling with the  
401 auditory or visual signal.

402 A confidence volume was estimated for all significant local maxima, using the  
403 bootstrap-based method described in Bourguignon et al. (2017a). The location of the maxima  
404 was also compared between conditions using the same bootstrap framework (Bourguignon et  
405 al., 2017a).

406 For each local maximum, individual maximum coherence values were extracted  
407 within a 10-mm sphere centered on the group level coordinates, or on the coordinates of  
408 maxima for *audio-only*. Coherence values were compared between conditions or signals of  
409 reference with two-sided paired *t*-tests.

410 The bootstrap method was used to assess the timing of peak TRFs (Efron and  
411 Tibshirani, 1993). As a preliminary step, TRFs were upsampled by spline interpolation to  
412 1000 Hz. A bootstrap distribution based on 10000 random drawings of subjects (or videos)  
413 was then built for the timing of peak TFR, from which we extracted the mean and standard  
414 deviation. Also the bias-corrected and accelerated bootstrap (Efron and Tibshirani, 1993) was  
415 used to compare the timing of peak TRF between conditions.

416 For the eye-tracking data, individual proportions of fixations were transformed using  
417 the empirical-logit transformation (Collins et al., 1992). Fixations to eyes and mouth in *video-*  
418 *only* were compared to fixations to the flickering cross in *control-video-only* using a two-  
419 sided paired *t*-test across participants.

#### 420 ***Data and software availability***

421 MEG and eye-tracking data as well as video stimuli are available on request from the  
422 corresponding author.

423

## 424 **Results**

425 Table 1 provides the coordinates and significance level of the loci of statistically  
426 significant coherence with the auditory speech envelope (henceforth, speech entrainment) and  
427 mouth opening (henceforth, lip entrainment) in all conditions (*audio-only*, *video-only*, and  
428 *control-video-only*) at all the selected frequency ranges (0.5 Hz, 1–3 Hz, 2–5 Hz, and 4–8  
429 Hz).

### 430 ***Entrainment to heard speech***

431 In *audio-only*, significant speech entrainment peaked at sensors covering bilateral  
432 auditory regions in all the explored frequency ranges: 0.5-Hz (Fig. 5A), 1–3 Hz (Fig. 6A), 2–  
433 5 Hz (Fig. 6B), and 4–8 Hz (Fig. 6C). Underlying sources were located in bilateral auditory  
434 cortices (Fig. 5A, 6, and Table 1).

### 435 ***Auditory cortices entrain to absent speech at frequencies below 1 Hz***

436 In *visual-only*, there was significant 0.5-Hz entrainment to the speech sound that was  
437 actually produced by the speaker, but not heard by participants (see Fig. 5B and Table 1). The  
438 significant loci for speech entrainment were the bilateral auditory cortices, the left inferior  
439 frontal gyrus, and the inferior part of the left precentral sulcus (Fig. 5B and Table 1).  
440 Critically, the location of the auditory sources where we observed maximum 0.5-Hz  
441 entrainment did not differ significantly between *audio-only* and *video-only* (left,  $F_{3,998} = 1.62$ ,  
442  $p = 0.18$ ; right,  $F_{3,998} = 0.85$ ,  $p = 0.47$ ). Not surprisingly, the magnitude of 0.5-Hz speech  
443 entrainment was higher in *audio-only* than in *video-only* (left,  $t_{27} = 6.36$ ,  $p < 0.0001$ ; right,  $t_{27}$   
444  $= 6.07$ ,  $p < 0.0001$ ). Nevertheless, brain responses associated with speech entrainment at ~0.5  
445 Hz displayed a similar time-course in *audio-only* and *video-only* (see Fig. 5A and 5B). In the  
446 left hemisphere, brain response peaked after the auditory speech envelope with a delay that  
447 did not differ significantly between the two conditions (*audio-only*,  $18 \pm 19$  ms, *video-only*,

448 73 ± 47 ms;  $p = 0.27$ ); in the right hemisphere this delay was significantly shorter for *audio-*  
449 *only* (43 ± 38 ms) than *video-only* (216 ± 54 ms;  $p = 0.019$ ). These results demonstrate that  
450 within the auditory cortices, neuronal activity at ~0.5 Hz is modulated similarly by heard  
451 speech sounds and absent speech when lip-read information is available, but incurs an  
452 additional delay in the right hemisphere. Next, we address four critical questions related to  
453 this effect: 1) Can it be explained by the general temporal characteristics of auditory speech?  
454 2) Is it unspecific to seeing the speaker's face? 3) Is it a direct result of lip-reading induced  
455 visual activity simply being fed to auditory areas? 4) Is it mediated by edges in the visual  
456 stimuli (predominantly reflecting head, eyebrows and jaw movements) that would prime  
457 phrase/sentence onset and modulate auditory cortical activity. A negative answer to these 4  
458 questions would support the view that auditory speech envelope is “synthesized” through  
459 internal models that map visual speech onto sound features.

460 ***Below 1-Hz entrainment to absent speech is not explained by the general temporal***  
461 ***characteristics of auditory speech***

462 In *video-only*, auditory sources (coordinates identified in *audio-only*) entrained  
463 significantly more to the corresponding—though absent—auditory speech than to unrelated  
464 auditory speech, here taken as the corresponding speech shifted in time (left,  $t_{27} = 3.08$ ,  $p =$   
465  $0.0047$ ; right,  $t_{27} = 3.78$ ,  $p = 0.0008$ ; see Fig. 7A). In this analysis, individual subject values  
466 were computed as the mean value across all considered time shifts. In addition, inspection of  
467 the maps of entrainment to unrelated speech did not reveal any special tendency to peak in  
468 auditory regions. This demonstrates that entrainment to absent speech in auditory cortices is  
469 not a consequence of the general temporal characteristics of auditory speech.

470 ***Below 1 Hz entrainment to absent speech is specific to seeing speaker's face***

471 Analysis of a *control-visual-only* condition revealed that entrainment to unheard  
472 speech at auditory cortices was specific to seeing the speaker's face. In the control condition,  
473 participants were looking at a silent video of a flickering Greek cross whose luminance  
474 pattern dynamically encoded the speaker's mouth configuration. We observed luminance-  
475 driven entrainment at 0.5 Hz at occipital cortices (Table 1), but no significant entrainment  
476 with unheard speech ( $p > 0.1$ , Fig. 7B). Importantly, speech entrainment at auditory sources  
477 (coordinates identified in *audio-only*) was significantly higher in *video-only* than in *control-*  
478 *video-only* (left,  $t_{27} = 3.44$ ,  $p = 0.0019$ ; right,  $t_{27} = 4.44$ ,  $p = 0.00014$ , see Fig. 7A). These  
479 differences in auditory speech entrainment cannot be explained by differences in attention as  
480 participants attended the flickering cross in *control-video-only* approximately as much as  
481 speaker's eyes and mouth in *video-only* ( $81.0 \pm 20.9\%$  vs.  $87.5 \pm 17.1\%$ ;  $t_{26} = 1.30$ ,  $p = 0.20$ :  
482 fixation data derived from eye-tracking recordings). This demonstrates that auditory cortical  
483 entrainment to unheard speech is specific to seeing the speaker's face.

484 ***Below 1-Hz entrainment to absent speech does not result from a direct feeding of lip***  
485 ***movements to auditory cortices***

486 Although driven by lip-read information, auditory cortical activity at  $\sim 0.5$  Hz in  
487 *visual-only* entrained more to unheard speech than to seen lip movements. Indeed, speech  
488 entrainment was stronger than lip entrainment at the left auditory source coordinates  
489 identified in *audio-only* ( $t_{27} = 2.52$ ,  $p = 0.018$ , see Fig. 7A). The same trend was observed at  
490 the right auditory source ( $t_{27} = 1.98$ ,  $p = 0.058$ , see Fig. 7A). However, at 0.5 Hz, lip  
491 movements entrained brain activity in the right angular gyrus (Fig. 7C and Table 1), a visual  
492 integration hub implicated in biological motion perception (Allison et al., 2000; Puce and  
493 Perrett, 2003). Such entrainment entailed a visual-speech-to-brain delay of  $40 \pm 127$  ms. Note

494 that the dominant source of lip and speech entrainment were  $\sim 4$  cm apart ( $F_{3,998} = 4.68$ ,  $p =$   
495 0.0030). Still, despite being distinct, their relative proximity might be the reason why speech  
496 entrainment was only marginally higher than lip entrainment in the right auditory cortex.  
497 Indeed, due to issues inherent to reconstructing brain signals based on extracranial signals  
498 (known as source leakage), lip entrainment estimated at the auditory cortex was artificially  
499 enhanced by the source in the angular gyrus. This leads us to conclude that entrainment in  
500 bilateral auditory cortices occurred with unheard speech rather than with seen lip movements.  
501 As further support for this claim, speech entrainment was still significant bilaterally in  
502 auditory cortices after partialling out lip movements (mouth opening and width; see Fig. 7D).  
503 In the right hemisphere, it peaked 2.2 mm away from sources observed without partialling out  
504 lip movements. In the left hemisphere, the peak in the partial coherence map was displaced  
505 towards the middle temporal gyrus (MNI coordinates:  $[-64 -21 -9]$ ). Although it did not  
506 peak in the left auditory cortex, the source distribution of the partial coherence was clearly  
507 pulled towards that brain region.

508 ***Below 1-Hz entrainment to absent speech is not explained by modulation of auditory***  
509 ***activity by edges in the visual stimulus***

510 Speech entrainment did not differ significantly from entrainment to the global visual  
511 change signal at the coordinates of bilateral auditory sources identified in *audio-only* (left,  $t_{27}$   
512 = 1.17,  $p = 0.25$ ; right,  $t_{27} = 1.10$ ,  $p = 0.28$ , see Fig. 7A). However, entrainment to the global  
513 visual change signal at  $\sim 0.5$  Hz was significant only in the posterior part of the right superior  
514 temporal gyrus (MNI coordinates:  $[62 -32 21]$ ), with a visual-change-to-brain delay of  $149 \pm$   
515 33 ms (See Fig. 7E). Most importantly, speech entrainment corrected for the global visual  
516 change signal still peaked and was significant in three left hemisphere sources that were less  
517 than 2.5 mm away from those of uncorrected speech entrainment (see Fig. 7F). Corrected  
518 speech entrainment in the right hemisphere peaked 1 mm away from the right auditory source

519 of uncorrected speech entrainment and was only marginally significant ( $p = 0.085$ ). In sum,  
520 global changes in the visual stimulus modulated oscillatory brain activity at  $\sim 0.5$  Hz in the  
521 right posterior superior temporal gyrus, but such modulation did not mediate the entrainment  
522 to absent speech.

523         Altogether, our results support the view that auditory speech envelope is synthesized  
524 through lip-reading.

### 525 ***Entrainment to absent speech at other frequencies***

526         At 1–3 Hz, there was significant entrainment to the absent speech in *visual-only* but  
527 not in *control-visual-only* (see Table 1). Significant entrainment to absent speech in *visual-*  
528 *only* peaked in the posterior part of the left inferior temporal gyrus, and in the central part of  
529 the middle temporal gyrus (see Table 1).

530         Entrainment in the posterior part of the left inferior temporal gyrus was specific to  
531 seeing the speaker's face (comparison *visual-only* vs. *control-visual-only*:  $t_{27} = 2.72$ ,  $p =$   
532  $0.011$ ) but did not entail a synthesis process since speech entrainment at this location was not  
533 significantly different from lip entrainment ( $t_{27} = 1.30$ ,  $p = 0.20$ ). It did not reach significance  
534 after partialling out mouth movements (see Table 1).

535         Entrainment in the central part of the middle temporal gyrus was not specific to seeing  
536 the speaker's face (comparison *visual-only* vs. *control-visual-only*:  $t_{27} = 1.48$ ,  $p = 0.15$ ) and  
537 did not entail a synthesis process since speech entrainment at this location was not  
538 significantly different from lip entrainment ( $t_{27} = -0.10$ ,  $p = 0.92$ ) despite surviving  
539 partialling out of mouth movements.

540         At 2–5 Hz, there was no significant entrainment to the absent speech in *visual-only*  
541 nor in *control-visual-only*.

542 At 4–8 Hz, there was significant entrainment to the absent speech in *video-only* and  
543 *control-video-only*, but only in occipital areas, and it vanished after partialling out the  
544 contribution of lip movements.

#### 545 ***Entrainment to lip movements***

546 Lip entrainment at 1–3 Hz, 2–5 Hz, and 4–8 Hz trivially occurred in occipital cortices  
547 in *video-only* and *control-video-only* (Table 1). Figure 8 illustrates entrainment at 2–5 Hz  
548 which we had planned to focus on based on previous reports (Park et al., 2016; Giordano et  
549 al., 2017). Brain responses associated with lip entrainment at 2–5 Hz peaked with a delay of  
550  $115 \pm 8$  ms (first source) and  $159 \pm 8$  ms (second source).

551 Our data do not suggest the presence of entrainment to unseen lip movements in  
552 visual cortices in *audio-only*. Indeed, in that condition, significant lip entrainment at 0.5 Hz  
553 occurred only in auditory cortices, and disappeared when we partialled out entrainment to the  
554 auditory speech envelope. No significant lip entrainment in this condition was found at any of  
555 the other tested frequency ranges: 1–3 Hz, 2–5 Hz and 4–8 Hz.

#### 556 ***Delays between auditory and visual speech***

557 Time-efficient synthesis of the auditory speech envelope might rely on the visual-to-  
558 auditory lag inherent to natural speech. Indeed, in our audio-visual stimuli, the  $\sim 0.5$ -Hz  
559 auditory speech envelope peaked  $87 \pm 9$  ms after the  $\sim 0.5$ -Hz mouth-opening time-course  
560 (see Fig. 9 left). But our results indicate that in *visual-only*, visual activity entrains to 2–5-Hz  
561 mouth movements while auditory activity entrains to an  $\sim 0.5$ -Hz absent auditory speech  
562 envelope. The simplest way to connect these oscillations is through phase–amplitude  
563 coupling, whereby the amplitude of 2–5-Hz visual activity modulates the phase of  $\sim 0.5$ -Hz  
564 auditory activity. Accordingly, we also estimated the delay from the envelope of 2–5-Hz



565 mouth opening time-course to  $\sim 0.5$ -Hz auditory speech envelope, and found it was  $170 \pm 7$   
566 ms (see Fig. 9 middle).

567         Also important is the interplay between global changes in the visual stimulus (mainly  
568 driven by head, eyebrows and jaw movements) and auditory speech envelope. This is because  
569 global visual changes could in principle modulate auditory cortical activity and hence  
570 mediate entrainment to absent speech. And indeed, in our audio-visual stimuli, the  $\sim 0.5$ -Hz  
571 auditory speech envelope peaked  $73 \pm 22$  ms after the  $\sim 0.5$ -Hz global visual change signal  
572 (see Fig. 9 right), meaning that low-level visual changes can cue slow changes in speech  
573 envelope (indicating phrase/sentence boundaries). However, the global visual change signal  
574 and the auditory speech envelope were only weakly coupled at  $\sim 0.5$  Hz (mean  $\pm$  SD  
575 coherence across the 8 video stimuli:  $0.051 \pm 0.022$ ) and in the other frequency ranges we  
576 explored. For a comparison, this degree of coupling was significantly lower than that between  
577 mouth opening and the auditory speech envelope ( $t_7 = 5.63$ ,  $p = 0.0008$ ; paired t-test on the  
578 coherence values for the 8 videos). In other words, lip movements provide more information  
579 about speech envelope than global changes in the visual stimulus, and similar temporal lead  
580 on auditory speech envelope (see Fig. 9). This further supports the view that auditory cortical  
581 entrainment to silent speech results from a fast synthesis process driven by lip reading rather  
582 than from modulation of auditory activity driven by the identification of low-level cross-  
583 sensory correspondences.

584

## 585 **Discussion**

586         We have demonstrated that the brain synthesises the slow (below 1 Hz) temporal  
587 dynamics of unheard speech from lip-reading. Specifically, watching silent lip-read videos  
588 without prior knowledge of what the speaker is saying leaves a trace of the auditory speech  
589 envelope in auditory cortices that closely resembles that left by the actual speech sound.

590

591 *Entrainment to unheard speech in auditory cortices*

592 Our most striking finding was that lip-reading induced entrainment in auditory  
593 cortices to the absent auditory speech at frequencies below 1 Hz. This entrainment 1) was  
594 specific to lip-reading, 2) was not a consequence of the general temporal characteristics of  
595 auditory speech, 3) was not a mere byproduct of entrainment to lip movements, and 4) was  
596 not mediated by low-level changes in the visual stimulus (at least in the left hemisphere).  
597 Instead, this genuine entrainment is similar to the entrainment induced by actual auditory  
598 speech: both are rooted in bilateral auditory cortices and are characterized by similar time-  
599 courses, though with an additional delay of ~200 ms in the right hemisphere. This suggests  
600 the existence of a time-efficient synthesis mechanism that maps facial articulatory mouth  
601 gestures onto corresponding speech sound features. Such a mechanism would likely leverage  
602 the natural visual-to-auditory speech delay (90–170 ms) and could be explained by visually-  
603 driven predictive coding (Friston and Kiebel, 2009). Likewise, auditory-driven predictive  
604 coding could account for the short (below-50-ms) latencies observed here in *audio-only* (Park  
605 et al., 2015).

606 Importantly, such auditory entrainment is unlikely to be driven by auditory imagery.  
607 Auditory imagery reflects perceptual auditory processing not triggered by external auditory  
608 stimulation (Nanay, 2018). In principle, observation of lip movements could lead to auditory  
609 imagery of related or unrelated speech or non-speech sounds. **Clearly**, auditory imagery of  
610 the actual speech sounds **was never an option** since participants were not professional lip-  
611 readers and were not cued about speech content. **Furthermore, our results demonstrate that**  
612 the auditory entrainment we observed cannot be linked to auditory imagery of unrelated  
613 sounds **either** since it was stronger for the corresponding but absent sound than for either seen

614 lip movements or unrelated speech. Accordingly, the fast synthesis hypothesis we have  
615 suggested seems to be the most likely interpretation of the observed entrainment.

616         The synthesis mechanism we have uncovered is likely grounded in the fact that lip-  
617 read information is coupled to the auditory signal in space and time (Munhall and Vatikiotis-  
618 Bateson, 2004; Chandrasekaran et al., 2009). In addition, the phonetic identity of each  
619 phoneme is supported by sound as well by the configuration of the lips. Even young infants  
620 are sensitive to this type of correspondence (Kuhl and Meltzoff, 1982), and phonetic  
621 integration continues to develop into adulthood, where the first traces of speech-specific  
622 phonetic integration are observed within ~250 ms after sound onset (Stekelenburg and  
623 Vroomen, 2012; Baart et al., 2014). Presumably, the tight audiovisual coupling in speech lies  
624 at the foundation of lip-read-induced entrainment to absent auditory speech in the brain, and  
625 there is indeed much evidence for entrainment to auditory speech at phrase and syllable rates  
626 (Luo and Poeppel, 2007; Bourguignon et al., 2012; Gross et al., 2013; Peelle et al., 2013;  
627 Molinaro et al., 2016; Vander Ghinst et al., 2016; Meyer et al., 2017; Meyer and Gumbert,  
628 2018).

629         Frequencies below 1 Hz match with phrasal, stress and sentential rhythmicity.  
630 Accordingly, corresponding entrainment to heard speech sounds has been hypothesised to  
631 subserve parsing or chunking of phrases and sentences (Ding et al., 2016; Meyer et al., 2017),  
632 or to help align neural excitability with syntactic information to optimize language  
633 comprehension (Meyer and Gumbert, 2018). Hence, our data suggest that such  
634 entrainment/alignment can be obtained through lip-reading, thereby facilitating speech  
635 chunking, parsing, and extraction of syntactic information.

636         As 4–8 Hz frequencies match with syllable rate, corresponding entrainment has been  
637 hypothesised to reflect parsing or chunking of syllables. Supporting this view, 4–8-Hz  
638 entrainment is enhanced when listening to intelligible speech compared to non-intelligible

639 speech (Ahissar et al., 2001; Luo and Poeppel, 2007; Peelle et al., 2013). However, we did  
640 not observe such entrainment during silent lip-reading, which may suggest that the brain does  
641 not synthesise the detailed phonology of unfamiliar silent syllabic structures based on lip-read  
642 information only. After all, lip-reading is a very difficult task, even for professional lip-  
643 readers (Chung et al., 2017). This is because different phonemes correspond to very similar  
644 lip configurations (*e.g.*, /ba/, /pa/ and /ma/). However, when the auditory signal is known, this  
645 ambiguity in the mapping between lip-reading and the corresponding phonemes disappears.  
646 Indeed, it has been suggested that lip-reading can induce entrainment in auditory cortices at  
647 frequencies above 1 Hz when participants are aware of the content of the visual-only speech  
648 stimuli (Crosse et al., 2015).

649

### 650 *Entrainment to lip movements*

651 During silent lip-reading, activity in early visual cortices entrained to lip movements  
652 mainly at frequencies above 1 Hz, in line with previous studies (Park et al., 2016; Giordano et  
653 al., 2017). Such occipital lip entrainment was reported to be modulated by audio-visual  
654 congruence (Park et al., 2016). This is probably the first necessary step for the brain to  
655 synthesize features of the absent auditory speech. Our results suggest that corresponding  
656 signals are forwarded to the right angular gyrus (Hauswald et al., 2018).

657 The right angular gyrus was the dominant source of lip entrainment at frequencies  
658 below 1 Hz. It is the convergence area for the dorsal and ventral visual streams and is  
659 specialised for processing visual biological motion (Perrett et al., 1989; Allison et al., 2000;  
660 Puce and Perrett, 2003; Marty et al., 2015). The right angular gyrus—or more precisely an  
661 area close to it termed the temporal visual speech area (Bernstein et al., 2011; Bernstein and  
662 Liebenthal, 2014)—activates during lip-reading (Calvert et al., 1997; Allison et al., 2000;  
663 Campbell et al., 2001) and observation of mouth movements (Puce et al., 1998). It has also

664 been suggested that it maps visual input onto linguistic representation during reading  
665 (Démonet et al., 1992), and lipreading (Hauswald et al., 2018). Our results shed light on the  
666 oscillatory dynamics underpinning such mapping during lip-reading: based on visual input at  
667 dominant lip movement frequencies (above 1 Hz), the angular gyrus presumably extracts  
668 features of lip movements below 1 Hz, which can then serve as an intermediate step to  
669 synthesise speech sound features. Given the short lip-to-brain delay observed in this brain  
670 area (~40 ms), such extraction might rely on the prediction of mouth movements.

671

### 672 *Entrainment to unheard speech in visual cortices*

673 Previous studies that have examined the brain dynamics underlying lipreading of  
674 silent connected visual speech have essentially focused on visuo-phonological mapping in  
675 occipital cortices (O’Sullivan et al., 2016; Lazard and Giraud, 2017; Hauswald et al., 2018).  
676 For example, it was shown that occipital 0.3–15-Hz EEG signals are better predicted by a  
677 combination of motion changes, visual speech features and the unheard auditory speech  
678 envelope than by motion changes alone (O’Sullivan et al., 2016). Also, visual activity has  
679 been reported to entrain more to absent speech at 4–7 Hz when a video is played forward  
680 rather than backward (Hauswald et al., 2018). Importantly, this effect was not driven by  
681 entrainment to lip movements since lip entrainment was similar for videos played forwards  
682 and backwards. Instead, it came with increased top-down drive from left sensorimotor  
683 cortices to visual cortices, indicating that visuo-phonological mapping had already taken  
684 place in early visual cortices through top-down mechanisms (O’Sullivan et al., 2016;  
685 Hauswald et al., 2018). Our study complements these results by showing that auditory  
686 cortices also entrain to unheard speech, but at frequencies below 1 Hz, probably based on  
687 earlier processes taking place in the occipital regions and the right angular gyrus.

688

689 ***Limitations and future perspectives***

690 We did not collect behavioral data from our participants. Further studies should  
691 clarify how the synthesis mechanism we have uncovered relates to individual lip-reading  
692 abilities, or susceptibility to the McGurk effect.

693 It also remains to be clarified what features of speech are synthesised, and under  
694 which circumstances auditory cortices can entrain to absent speech at higher frequencies  
695 (especially 4–8-Hz).

696 Finally, it will be important to specify which features of the articulatory mouth  
697 gestures lead to below-1-Hz auditory entrainment to absent speech. This would require visual  
698 control conditions in which, for example, lip movements are shown in isolation, or replaced  
699 by point-light stimuli.

700

701 ***Conclusion***

702 Our results demonstrate that the brain can quickly synthesize a representation of  
703 coarse-grained auditory speech features in early auditory cortices and shed light on the  
704 underlying oscillatory dynamics. Seeing lip movements first modulates neuronal activity in  
705 early visual cortices at frequencies that match articulatory lip movements (above 1 Hz).  
706 Based on this activity, the right angular gyrus, putatively the temporal visual speech area,  
707 extracts and possibly predicts the slower features of lip movements. Finally, these slower lip  
708 movement dynamics are mapped onto their corresponding speech sound features and this  
709 information is fed to auditory cortices. Receiving this information likely facilitates speech  
710 parsing, in line with the hypothesised role of entrainment to heard speech at frequencies  
711 below 1 Hz.

712 **References**

- 713 Ahissar E, Nagarajan S, Ahissar M, Protopapas A, Mahncke H, Merzenich MM (2001)  
714 Speech comprehension is correlated with temporal response patterns recorded from  
715 auditory cortex. *Proc Natl Acad Sci U S A* 98:13367–13372.
- 716 Allison T, Puce A, McCarthy G (2000) Social perception from visual cues: role of the STS  
717 region. *Trends in Cognitive Sciences* 4:267–278 .
- 718 Ashburner J, Friston KJ (1999) Nonlinear spatial normalization using basis functions. *Hum*  
719 *Brain Mapp* 7:254–266.
- 720 Ashburner J, Neelin P, Collins DL, Evans A, Friston K (1997) Incorporating prior knowledge  
721 into image registration. *Neuroimage* 6:344–352.
- 722 Baart M, Stekelenburg JJ, Vroomen J (2014) Electrophysiological evidence for speech-  
723 specific audiovisual integration. *Neuropsychologia* 53:115–121.
- 724 Bernstein LE, Auer ET, Moore JK, Ponton CW, Don M, Singh M (2002) Visual speech  
725 perception without primary auditory cortex activation. *Neuroreport* 13:311–315.
- 726 Bernstein LE, Jiang J, Pantazis D, Lu Z-L, Joshi A (2011) Visual phonetic processing  
727 localized using speech and nonspeech face gestures in video and point-light displays.  
728 *Human Brain Mapping* 32:1660–1676.
- 729 Bernstein LE, Liebenthal E (2014) Neural pathways for visual speech perception. *Front*  
730 *Neurosci* 1:386.
- 731 Bortel R, Sovka P (2014) Approximation of the null distribution of the multiple coherence  
732 estimated with segment overlapping. *Signal Processing* 96:310–314.
- 733 Bourguignon M, De Tiège X, Op de Beeck M, Ligot N, Paquier P, Van Bogaert P, Goldman  
734 S, Hari R, Jousmäki V (2012) The pace of prosodic phrasing couples the listener's  
735 cortex to the reader's voice. *Hum Brain Mapp* 34:314–326.
- 736 Bourguignon M, Molinaro N, Wens V (2017a) Contrasting functional imaging parametric

737 maps: The mislocation problem and alternative solutions. *Neuroimage* 169:200–211.

738 Bourguignon M, Piitulainen H, De Tiège X, Jousmäki V, Hari R (2015) Corticokinematic  
739 coherence mainly reflects movement-induced proprioceptive feedback. *Neuroimage*  
740 106:382–390.

741 Bourguignon M, Piitulainen H, Smeds E, Zhou G, Jousmäki V, Hari R (2017b) MEG Insight  
742 into the Spectral Dynamics Underlying Steady Isometric Muscle Contraction. *J*  
743 *Neurosci* 37:10421–10437.

744 Calvert GA, Bullmore ET, Brammer MJ, Campbell R, Williams SC, McGuire PK, Woodruff  
745 PW, Iversen SD, David AS (1997) Activation of auditory cortex during silent  
746 lipreading. *Science* 276:593–596.

747 Campbell R, MacSweeney M, Surguladze S, Calvert G, McGuire P, Suckling J, Brammer  
748 MJ, David AS (2001) Cortical substrates for the perception of face actions: an fMRI  
749 study of the specificity of activation for seen speech and for meaningless lower-face  
750 acts (gurning). *Brain Res Cogn Brain Res* 12:233–243.

751 Chandrasekaran C, Trubanova A, Stillitano S, Caplier A, Ghazanfar AA (2009) The natural  
752 statistics of audiovisual speech. *PLoS Comput Biol* 5:e1000436.

753 Chung JS, Senior A, Vinyals O, Zisserman A (2017) Lip Reading Sentences in the Wild. In:  
754 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

755 Collins JJ, Fanciulli M, Hohlfeld RG, Finch DC, Sandri G v. H, Shtatland ES (1992) A  
756 random number generator based on the logit transform of the logistic variable.  
757 *Computers in Physics* 6:630.

758 Crosse MJ, Di Liberto GM, Bednar A, Lalor EC (2016) The Multivariate Temporal Response  
759 Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to  
760 Continuous Stimuli. *Front Hum Neurosci* 10:604.

761 Crosse MJ, ElShafei HA, Foxe JJ, Lalor EC (2015) Investigating the temporal dynamics of



762 auditory cortical activation to silent lipreading. In: 2015 7th International IEEE/EMBS  
763 Conference on Neural Engineering (NER).

764 Démonet JF, Chollet F, Ramsay S, Cardebat D, Nespoulous JL, Wise R, Rascol A,  
765 Frackowiak R (1992) The anatomy of phonological and semantic processing in normal  
766 subjects. *Brain* 115:1753–1768.

767 Ding N, Melloni L, Zhang H, Tian X, Poeppel D (2016) Cortical tracking of hierarchical  
768 linguistic structures in connected speech. *Nat Neurosci* 19:158–164.

769 Efron B, Tibshirani RJ (1993) *An Introduction to the Bootstrap*.

770 Eveno N, Caplier A, Coulon P-Y (2004) Accurate and Quasi-Automatic Lip Tracking. *IEEE*  
771 *Trans Circuits Syst Video Technol* 14:706–715.

772 Friston K, Kiebel S (2009) Predictive coding under the free-energy principle. *Philos Trans R*  
773 *Soc Lond B Biol Sci* 364:1211–1221.

774 Giordano BL, Ince RAA, Gross J, Schyns PG, Panzeri S, Kayser C (2017) Contributions of  
775 local speech encoding and functional connectivity to audio-visual speech perception.  
776 *Elife* 6:e24763.

777 Gramfort A, Luessi M, Larson E, Engemann DA, Strohmeier D, Brodbeck C, Parkkonen L,  
778 Hämäläinen MS (2014) MNE software for processing MEG and EEG data.  
779 *Neuroimage* 86:446–460.

780 Gross J, Hoogenboom N, Thut G, Schyns P, Panzeri S, Belin P, Garrod S (2013) Speech  
781 rhythms and multiplexed oscillatory sensory coding in the human brain. *PLoS Biol*  
782 11:e1001752.

783 Halliday D (1995) A framework for the analysis of mixed time series/point process data—  
784 Theory and application to the study of physiological tremor, single motor unit  
785 discharges and electromyograms. *Prog Biophys Mol Biol* 64:237–278.

786 Hauswald A, Lithari C, Collignon O, Leonardelli E, Weisz N (2018) A Visual Cortical

787 Network for Deriving Phonological Information from Intelligible Lip Movements. *Curr*  
788 *Biol* 28:1453–1459.e3.

789 Hillebrand A, Barnes GR (2005) Beamformer Analysis of MEG Data. *International Review*  
790 *of Neurobiology* 68:149–171.

791 Hyvärinen A, Karhunen J, Oja E (2004) *Independent Component Analysis*. John Wiley &  
792 Sons.

793 Kaplan E, Jesse A (2019) Fixating the eyes of a speaker provides sufficient visual  
794 information to modulate early auditory processing. *Biol Psychol* 146:107724.

795 Kapnoula EC, Packard S, Gupta P, McMurray B (2015) Immediate lexical integration of  
796 novel word forms. *Cognition* 134:85–99.

797 Keitel A, Gross J, Kayser C (2018) Perceptually relevant speech tracking in auditory and  
798 motor cortex reflects distinct linguistic features. *PLoS Biol* 16:e2004473.

799 Kendall MG, Stuart A (1968) *The Advanced Theory of Statistics*. *The Statistician* 18:163.

800 Kuhl PK, Meltzoff AN (1982) The bimodal perception of speech in infancy. *Science*  
801 218:1138–1141.

802 Lalor EC, Foxe JJ (2010) Neural responses to uninterrupted natural speech can be extracted  
803 with precise temporal resolution. *Eur J Neurosci* 31:189–193.

804 Lazard DS, Giraud A-L (2017) Faster phonological processing and right occipito-temporal  
805 coupling in deaf adults signal poor cochlear implant outcome. *Nat Commun* 8:14872.

806 Luo H, Poeppel D (2007) Phase patterns of neuronal responses reliably discriminate speech  
807 in human auditory cortex. *Neuron* 54:1001–1010.

808 Marty B, Bourguignon M, Jousmäki V, Wens V, Op de Beeck M, Van Bogaert P, Goldman  
809 S, Hari R, De Tiège X (2015) Cortical kinematic processing of executed and observed  
810 goal-directed hand actions. *Neuroimage* 119:221–228.

811 McMurray B, Tanenhaus MK, Aslin RN (2002) Gradient effects of within-category phonetic

812 variation on lexical access. *Cognition* 86:B33–42.

813 Meyer L, Gumbert M (2018) Synchronization of Electrophysiological Responses with  
814 Speech Benefits Syntactic Information Processing. *J Cogn Neurosci*:1–10.

815 Meyer L, Henry MJ, Gaston P, Schmuck N, Friederici AD (2017) Linguistic Bias Modulates  
816 Interpretation of Speech via Neural Delta-Band Oscillations. *Cereb Cortex* 27:4293–  
817 4302.

818 Molinaro N, Lizarazu M, Lallier M, Bourguignon M, Carreiras M (2016) Out-of-synchrony  
819 speech entrainment in developmental dyslexia. *Hum Brain Mapp* 37:2767–2783.

820 Munhall KG, Jones JA, Callan DE, Kuratate T, Vatikiotis-Bateson E (2004) Visual prosody  
821 and speech intelligibility: head movement improves auditory speech perception.  
822 *Psychol Sci* 15:133–137.

823 Munhall KG, Vatikiotis-Bateson E (2004) Spatial and Temporal Constraints on Audiovisual  
824 Speech Perception. In: *The handbook of multisensory processes* (Calvert GA, Spence  
825 C, Stein BE, eds), pp 177–188. Cambridge, MA, US: MIT Press.

826 Nanay B (2018) Multimodal mental imagery. *Cortex* 105:125–134.

827 Nichols TE, Holmes AP (2002) Nonparametric permutation tests for functional  
828 neuroimaging: a primer with examples. *Hum Brain Mapp* 15:1–25.

829 O’Sullivan AE, Crosse MJ, Di Liberto GM, Lalor EC (2016) Visual Cortical Entrainment to  
830 Motion and Categorical Speech Features during Silent Lipreading. *Front Hum Neurosci*  
831 10:679.

832 Paré M, Richler RC, ten Hove M, Munhall KG (2003) Gaze behavior in audiovisual speech  
833 perception: the influence of ocular fixations on the McGurk effect. *Percept Psychophys*  
834 65:553–567.

835 Park H, Ince RAA, Schyns PG, Thut G, Gross J (2015) Frontal top-down signals increase  
836 coupling of auditory low-frequency oscillations to continuous speech in human

837 listeners. *Curr Biol* 25:1649–1653.

838 Park H, Kayser C, Thut G, Gross J (2016) Lip movements entrain the observers' low-  
839 frequency brain oscillations to facilitate speech intelligibility. *Elife* 5:e14521.

840 Paulesu E, Perani D, Blasi V, Silani G, Borghese NA, De Giovanni U, Sensolo S, Fazio F  
841 (2003) A functional-anatomical model for lipreading. *J Neurophysiol* 90:2005–2013.

842 Peelle JE, Gross J, Davis MH (2013) Phase-locked responses to speech in human auditory  
843 cortex are enhanced during comprehension. *Cereb Cortex* 23:1378–1387.

844 Pekkola J, Ojanen V, Autti T, Jääskeläinen IP, Möttönen R, Tarkiainen A, Sams M (2005)  
845 Primary auditory cortex activation by visual speech: an fMRI study at 3 T. *Neuroreport*  
846 16:125–128.

847 Perrett DI, Harries MH, Bevan R, Thomas S, Benson PJ, Mistlin AJ, Chitty AJ, Hietanen JK,  
848 Ortega JE (1989) Frameworks of analysis for the neural representation of animate  
849 objects and actions. *J Exp Biol* 146:87–113.

850 Puce A, Allison T, Bentin S, Gore JC, McCarthy G (1998) Temporal cortex activation in  
851 humans viewing eye and mouth movements. *J Neurosci* 18:2188–2199.

852 Puce A, Perrett D (2003) Electrophysiology and brain imaging of biological motion. *Philos*  
853 *Trans R Soc Lond B Biol Sci* 358:435–445.

854 Reuter M, Schmansky NJ, Rosas HD, Fischl B (2012) Within-subject template estimation for  
855 unbiased longitudinal image analysis. *Neuroimage* 61:1402–1418.

856 Schroeder CE, Lakatos P, Kajikawa Y, Partan S, Puce A (2008) Neuronal oscillations and  
857 visual amplification of speech. *Trends Cogn Sci* 12:106–113.

858 Stekelenburg JJ, Vroomen J (2012) Electrophysiological evidence for a multisensory speech-  
859 specific mode of perception. *Neuropsychologia* 50:1425–1431.

860 Sumbly WH, Pollack I (1954) Visual Contribution to Speech Intelligibility in Noise. *J Acoust*  
861 *Soc Am* 26:212–215.

862 Taulu S, Simola J (2006) Spatiotemporal signal space separation method for rejecting nearby  
863 interference in MEG measurements. *Phys Med Biol* 51:1759–1768.

864 Taulu S, Simola J, Kajola M (2005) Applications of the signal space separation method.  
865 *IEEE Trans Signal Process* 53:3359–3372.

866 Vander Ghinst M, Bourguignon M, Niesen M, Wens V, Hassid S, Choufani G, Jousmäki V,  
867 Hari R, Goldman S, De Tiège X (2019) Cortical Tracking of Speech-in-Noise Develops  
868 from Childhood to Adulthood. *J Neurosci* 39:2938–2950.

869 Vander Ghinst M, Ghinst MV, Bourguignon M, Op de Beeck M, Wens V, Marty B, Hassid  
870 S, Choufani G, Jousmäki V, Hari R, Van Bogaert P, Goldman S, De Tiège X (2016)  
871 Left Superior Temporal Gyrus Is Coupled to Attended Speech in a Cocktail-Party  
872 Auditory Scene. *J Neurosci* 36:1596–1606.

873 Van Veen BD, van Drongelen W, Yuchtman M, Suzuki A (1997) Localization of brain  
874 electrical activity via linearly constrained minimum variance spatial filtering. *IEEE*  
875 *Trans Biomed Eng* 44:867–880.

876 Vatikiotis-Bateson E, Eigsti I-M, Yano S, Munhall KG (1998) Eye movement of perceivers  
877 during audiovisual speech perception. *Perception & Psychophysics* 60:926–940.

878 Vigário R, Särelä J, Jousmäki V, Hämäläinen M, Oja E (2000) Independent component  
879 approach to the analysis of EEG and MEG recordings. *IEEE Trans Biomed Eng*  
880 47:589–593.

881 Zion Golumbic EM, Ding N, Bickel S, Lakatos P, Schevon CA, McKhann GM, Goodman  
882 RR, Emerson R, Mehta AD, Simon JZ, Poeppel D, Schroeder CE (2013) Mechanisms  
883 underlying selective neuronal tracking of attended speech at a “cocktail party.” *Neuron*  
884 77:980–991.

885 **Figures and Tables:**

886

887 **Table 1.**

888 Significant peak of speech- and lip entrainment: peak MNI coordinates, significance level,  
 889 confidence volume, and anatomical location. Only significant peaks of speech-entrainment  
 890 that survived partialling out lip movements (exceptions marked with \*) and global visual  
 891 changes (exceptions marked with \*\*) are presented here. Likewise, only peaks of significant  
 892 lip entrainment that survived partialling out the auditory speech envelope are presented here.  
 893 For the exceptions, *ps* are displayed in between parentheses.

	Peak coordinates [mm]	<i>p</i>	Mean ± SD values	Confidence volume [cm <sup>3</sup> ]	Anatomical location
<b>Speech entrainment at 0.5 Hz</b>					
<i>Audio-only</i>	[-64 -19 8]	<10 <sup>-3</sup>	0.076 ± 0.045	2.6	Left auditory cortex
	[64 -21 6]	<10 <sup>-3</sup>	0.075 ± 0.046	5.5	Right auditory cortex
<i>Video-only</i>	[-46 -30 11]	0.003	0.025 ± 0.017	35.5	Left auditory cortex
	[68 -14 -2]**	0.029 (0.085)	0.024 ± 0.015	5.6	Right auditory cortex
	[-57 25 15]	0.005	0.021 ± 0.013	9.6	Left inferior frontal gyrus
	[-58 -15 41]*	0.018 (0.063)	0.023 ± 0.012	21.3	Left inferior precentral sulcus
<b>Lip entrainment at 0.5 Hz</b>					
<i>Video-only</i>	[49 -46 10]	0.002	0.022 ± 0.014	30.9	Right angular gyrus
<i>Control-video-only</i>	[10 -89 -21]	<10 <sup>-3</sup>	0.028 ± 0.023	6.3	Inferior occipital area
	[25 -96 -1]	0.008	0.027 ± 0.023	11.7	Right lateral occipital cortex
	[-23 -97 -4]	0.046	0.023 ± 0.014	39.1	Left lateral occipital cortex
<b>Speech entrainment at 1-3 Hz</b>					
<i>Audio-only</i>	[-62 -15 11]	<10 <sup>-3</sup>	0.031 ± 0.017	0.17	Left auditory cortex
	[66 -10 9]	<10 <sup>-3</sup>	0.036 ± 0.022	0.22	Right auditory cortex
<i>Video-only</i>	[-51 -65 -16]	0.020	0.012 ± 0.004	58.8	Left inferior temporal gyrus
	[-67 -20 -12]*	0.005 (0.22)	0.012 ± 0.004	2.8	Left middle temporal gyrus
<b>Lip entrainment at 1-3 Hz</b>					
<i>Video-only</i>	[5 -92 -13]	<10 <sup>-3</sup>	0.015 ± 0.007	22.6	Calcarine cortex
	[33 -92 6]	0.001	0.014 ± 0.007	5.2	Right lateral occipital sulcus
	[-15 -96 12]	<10 <sup>-3</sup>	0.015 ± 0.005	18.3	Left calcarine cortex
<i>Control-video-only</i>	[1 -98 10]	<10 <sup>-3</sup>	0.029 ± 0.016	0.3	Calcarine cortex
	[34 -92 -3]	<10 <sup>-3</sup>	0.028 ± 0.016	0.9	Right lateral occipital cortex
	[-28 -94 -10]	<10 <sup>-3</sup>	0.023 ± 0.012	3.4	Left lateral occipital cortex
<b>Speech entrainment at 2-5 Hz</b>					
<i>Audio-only</i>	[67 -11 10]	<10 <sup>-3</sup>	0.020 ± 0.008	0.3	Left auditory cortex
	[-62 -14 13]	<10 <sup>-3</sup>	0.016 ± 0.007	0.4	Right auditory cortex
<b>Lip entrainment at 2-5 Hz</b>					
<i>Video-only</i>	[-14 -97 11]	<10 <sup>-3</sup>	0.018 ± 0.007	8.3	Left calcarine cortex
	[2 -93 -2]	<10 <sup>-3</sup>	0.018 ± 0.008	15.1	Calcarine cortex
<i>Control-video-only</i>	[-1 -98 11]	<10 <sup>-3</sup>	0.026 ± 0.016	1.8	Calcarine cortex
	[25 -97 -8]	<10 <sup>-3</sup>	0.025 ± 0.012	1.9	Right lateral occipital cortex
	[-28 -94 -10]	<10 <sup>-3</sup>	0.024 ± 0.012	0.3	Left lateral occipital cortex
<b>Speech entrainment at 4-8 Hz</b>					
<i>Audio-only</i>	[-64 -18 7]	<10 <sup>-3</sup>	0.013 ± 0.005	1.4	Left auditory cortex

	[67 -13 5]	<10 <sup>-3</sup>	0.020 ± 0.009	0.3	Right auditory cortex
<b>Lip entrainment at 4–8 Hz</b>					
<i>Video-only</i>	[10 -94 -4]	<10 <sup>-3</sup>	0.013 ± 0.005	19.3	Right calcarine cortex
	[-11 -95 9]	0.001	0.013 ± 0.004	18.8	Left calcarine cortex
<i>Control-video-only</i>	[-4 -88 -18]	<10 <sup>-3</sup>	0.016 ± 0.008	5.3	Inferior occipital cortex
	[27 -94 -5]	<10 <sup>-3</sup>	0.016 ± 0.007	22.5	Right lateral occipital cortex
	[-5 -97 15]	0.011	0.015 ± 0.006	30.0	Calcarine cortex

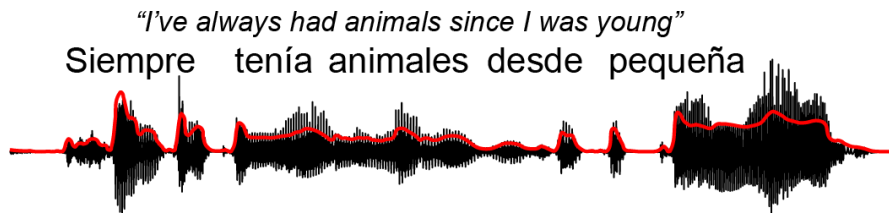
894

895

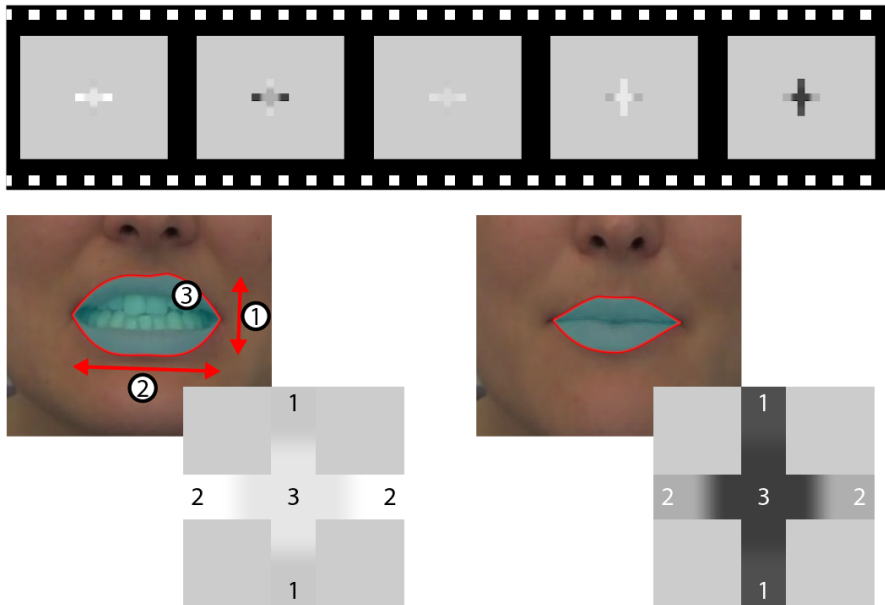
A. video signal



B. audio signal



C. control video



896

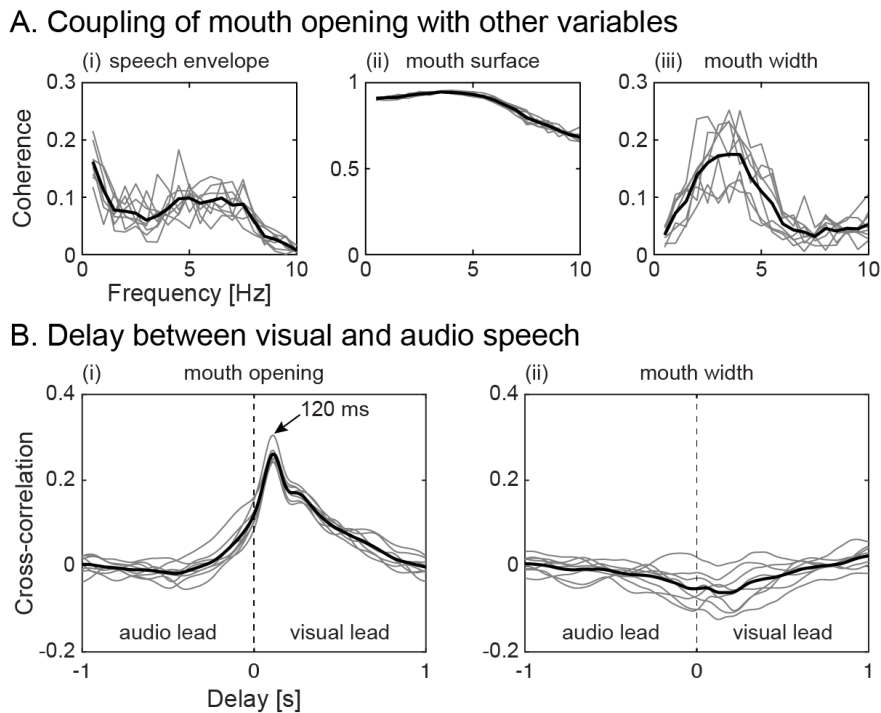
897 **Figure 1.** Experimental material. **A** and **B** — Two-second excerpt of video (**A**) and audio (**B**;

898 auditory speech envelope in red) of the speaker telling a 5-min story about a given topic.

899 There were 8 different videos. Video without sound was presented in *video-only*, and sound

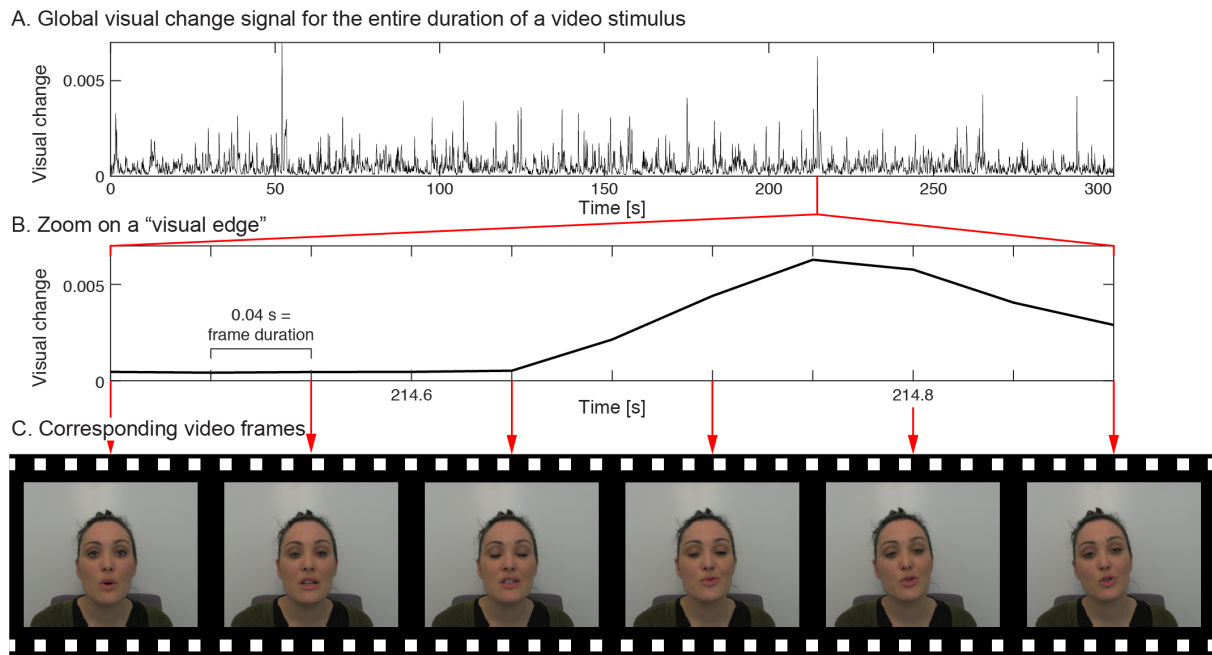
900 without video was presented in *audio-only*. **C** — Corresponding control video in which a

901 flickering Greek cross encoded speaker's mouth configuration. Based on a segmentation of  
 902 mouth contours, the cross encoded mouth opening (1), mouth width (2), and mouth surface  
 903 (3). The resulting video was presented in *control-video-only*.  
 904



905  
 906 **Figure 2.** Relation between audio and visual speech signals. **A** — Frequency-dependent  
 907 coupling (coherence) of mouth opening with auditory speech envelope (i), mouth surface (ii),  
 908 and mouth width (iii). Coupling is quantified with coherence. There is one gray trace per  
 909 video (8 in total), and thick black traces are the average across them all. **B** — Delay between  
 910 visual and audio speech assessed with cross-correlation of auditory speech envelope with  
 911 mouth opening (i) and mouth width (ii).  
 912





913

914 **Figure 3.** Global visual changes in the visual stimuli. **A** — The global visual change signal as  
 915 a function of time for the entire duration of a video stimulus. **B** — Zoom on one of the most  
 916 prominent edges (peaks) of the global visual change signal. **C** — Video frames corresponding  
 917 to this visual edge, showing that it was due to head movements.

918

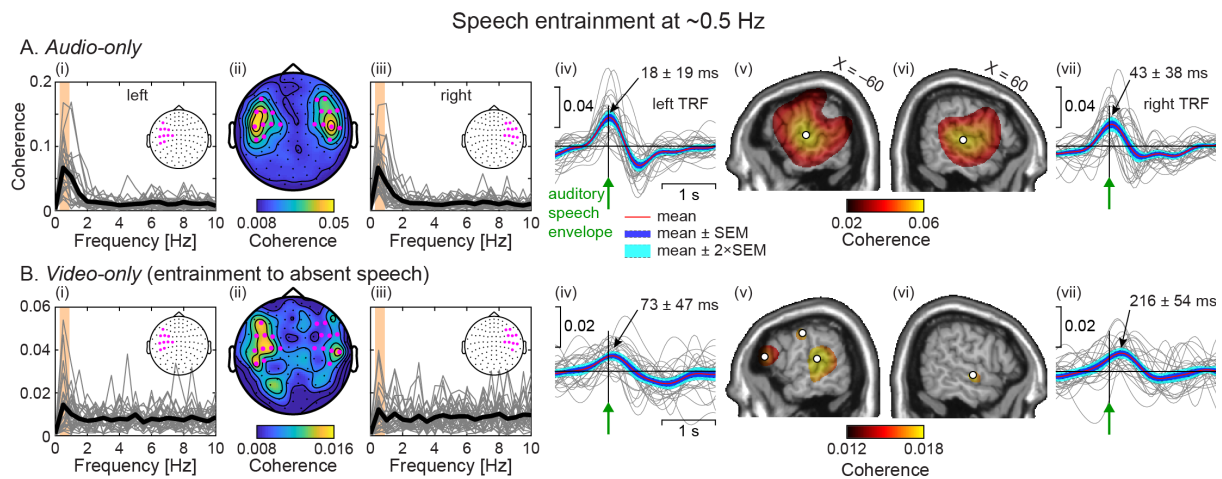


919

920 **Figure 4.** Regions of interest for eye fixation. The initial regions of interest are delineated in  
 921 yellow, and the extended ones in white. Eye fixation analyses were based on extended  
 922 regions. In *video-only (left)*, the final region of interest comprised the mouth and the eyes. In  
 923 *control-video-only (right)*, it encompassed the flickering cross.

924

925



926

927 **Figure 5.** Speech entrainment at 0.5 Hz. **A** — Speech entrainment in *audio-only*. (i–iii)

928 Sensor distribution of speech entrainment at 0.5 Hz quantified with coherence (ii) and its

929 spectral distribution at a selection of 10 sensors in the left (i) and right hemisphere (iii) of

930 maximal 0.5 Hz coherence (highlighted in magenta). Gray traces represent individual

931 subject’s spectra at the sensor of maximum 0.5 Hz coherence within the preselection, and the

932 thick black trace is their group average. (iv–vii) Brain distribution of significant speech

933 entrainment quantified with coherence in the left (v) and right hemispheres (vi) and the

934 temporal response function (TRF) associated with auditory speech envelope at coordinates of

935 peak coherence (marked with white discs) in the left (iv) and right hemispheres (vii). In brain

936 images, significant coherence values at MNI coordinates  $|X| > 40$  mm were projected

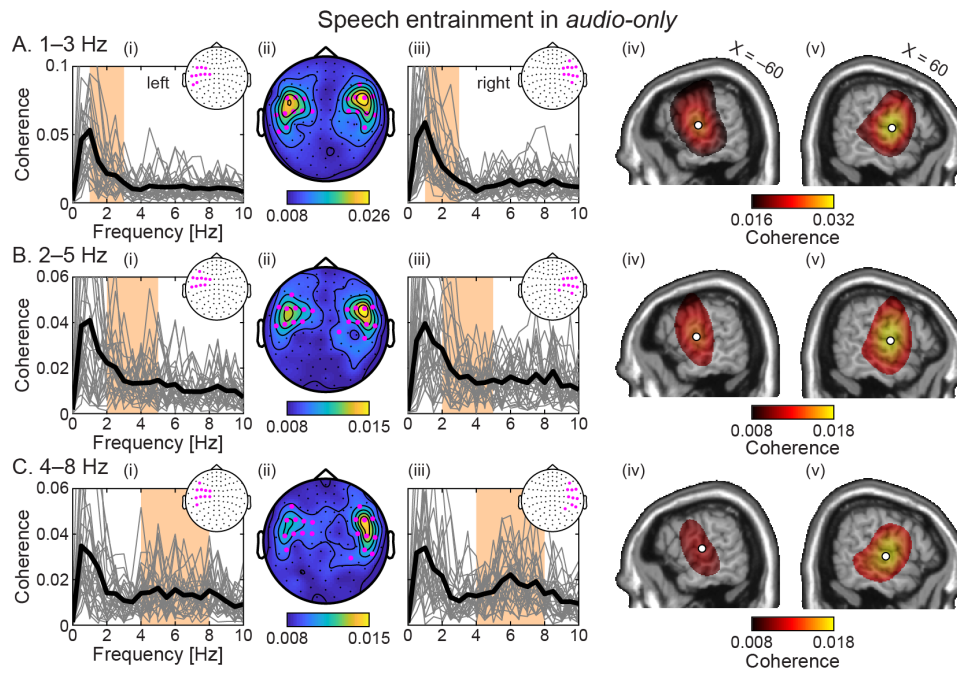
937 orthogonally onto the parasagittal slice of coordinates  $|X| = 60$  mm. **B** — Same as in **A** for

938 *video-only*, illustrating that seeing speaker’s face was enough to elicit significant speech

939 entrainment at auditory cortices. Note that coherence spectra were estimated at the subject-

940 specific sensor selected based on coherence in *audio-only*.

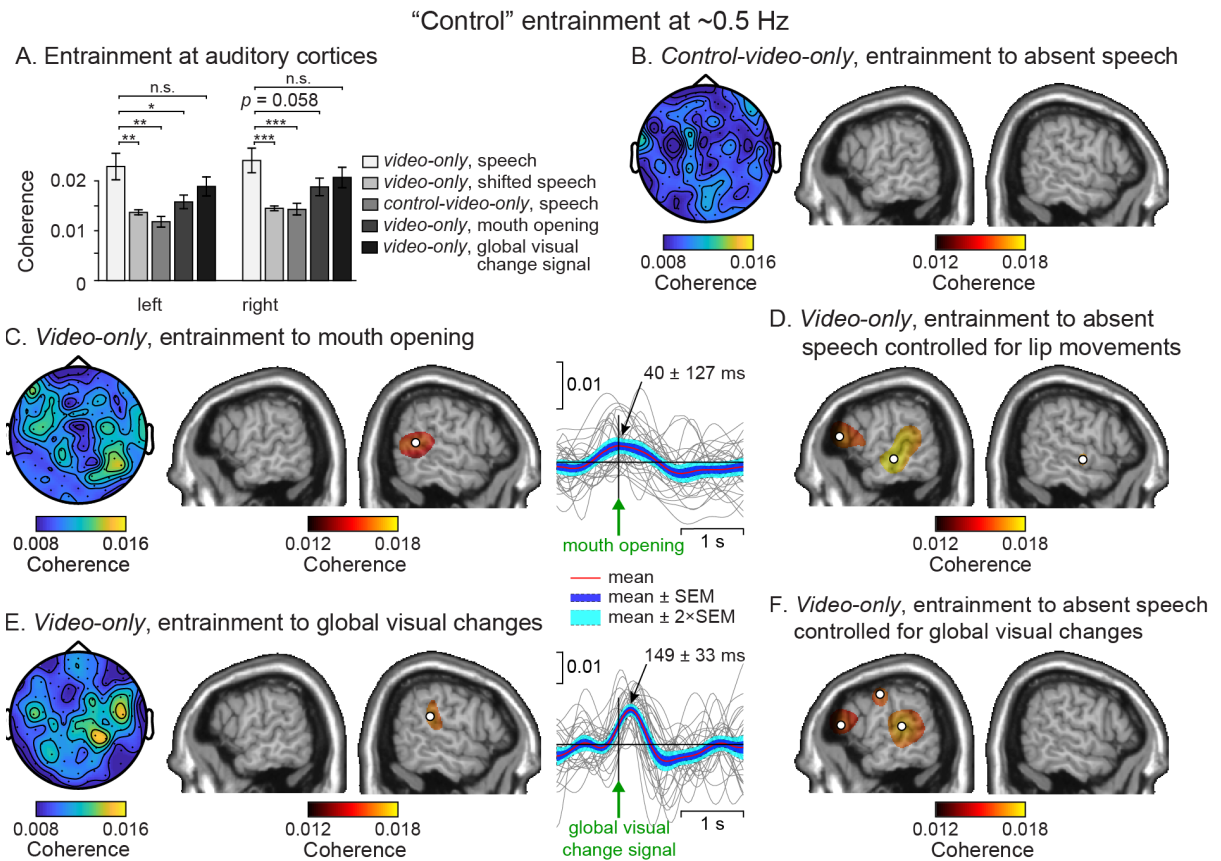
941



942

943 **Figure 6.** Speech entrainment quantified with coherence in *audio-only* at 1–3 Hz (A), 2–5 Hz  
 944 (B) and 4–8 Hz (C). (i–iii) Sensor distribution of speech entrainment (ii) and its spectral  
 945 distribution at a selection of 10 left- (i) and right-hemisphere (iii) sensors of maximum  
 946 coherence (highlighted in magenta). Gray traces represent individual subject's spectra at the  
 947 sensor of maximum coherence across the considered frequency range and within the  
 948 preselection, and the thick black trace is their group average. (iv & v) Brain distribution of  
 949 significant speech entrainment in the left (iv) and right hemispheres (v) produced as  
 950 described in Fig. 5.

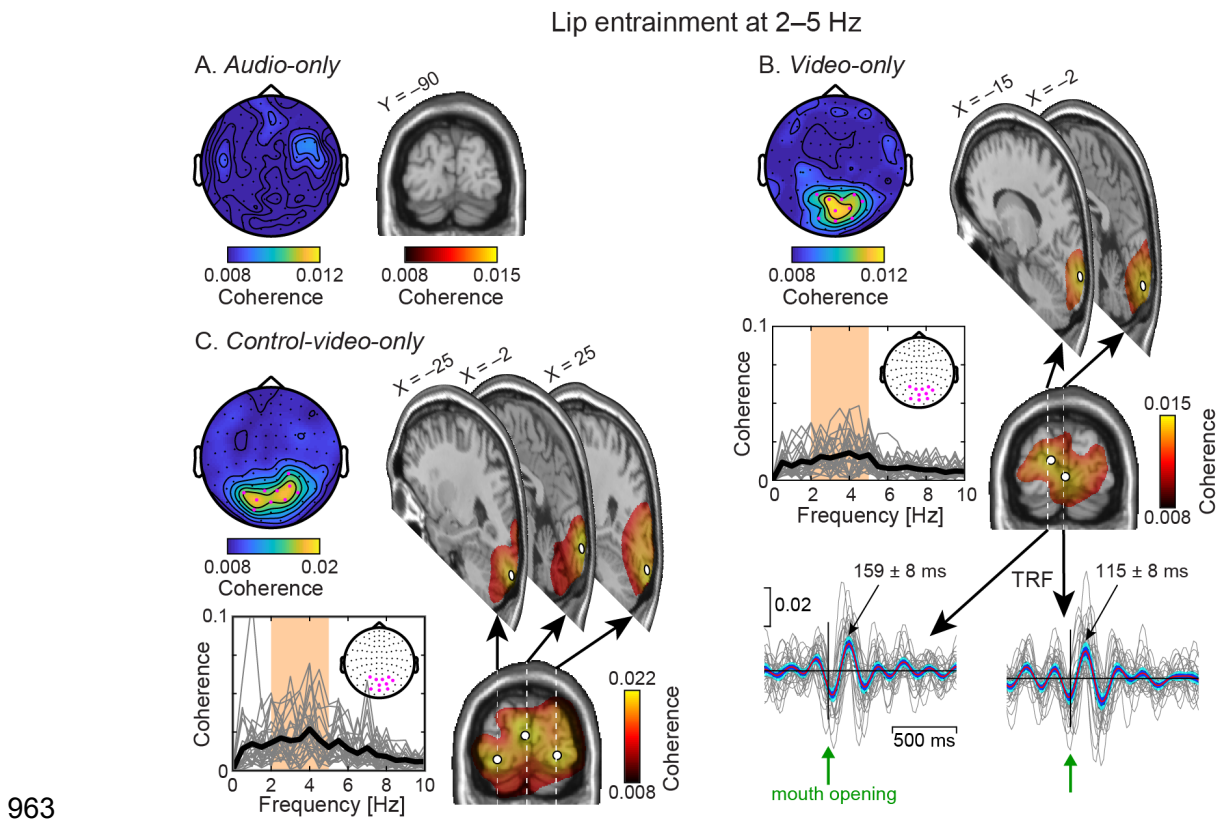
951



952

953 **Figure 7.** Control for the entrainment to absent speech at 0.5 Hz. **A** — Entrainment values  
 954 quantified with coherence at coordinates identified in *audio-only* (mean  $\pm$  SD across  
 955 participants). **B** — Sensor and brain distribution of auditory speech entrainment in *control-*  
 956 *video-only* wherein speech entrainment was not significant. **C** — Sensor and brain  
 957 distribution of lip entrainment in *video-only* and associated temporal evolution. Lip  
 958 entrainment was significant only in the right angular gyrus. **D** — Brain distribution of  
 959 significant speech entrainment at 0.5 Hz after partialling out lip movements (mouth opening  
 960 and mouth width). **E & F** — As in C & D but for the global visual change signal instead of  
 961 mouth opening. Brain images were produced as described in Fig. 5.

962

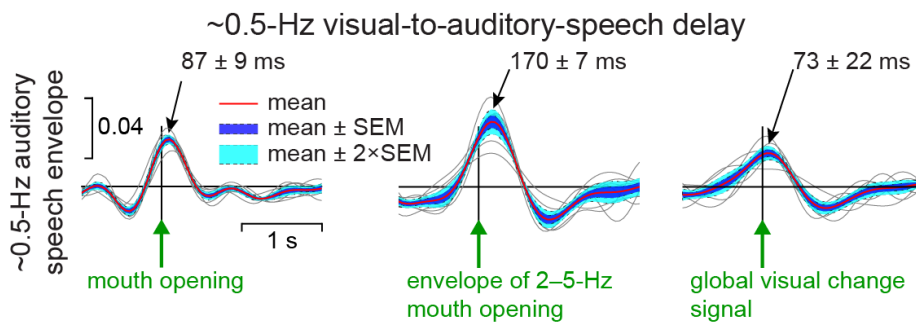


963

964 **Figure 8.** Lip entrainment at 2–5 Hz in *audio-only* (A), *video-only* (B) and *control-video-only*  
 965 (C). Lip entrainment is presented both in the sensor space and on the brain in all conditions  
 966 (*audio-only*, *video-only*, *control-video-only*). In brain maps, significant coherence values at  
 967 MNI coordinates  $Y < -70$  mm were projected orthogonally onto the coronal slice of  
 968 coordinates  $|Y| = 90$  mm. Locations of peak coherence are marked with white discs. Note that  
 969 coherence was not significant in *audio-only*. Additional parasagittal maps are presented for  
 970 all significant peaks of coherence. In these maps, the orthogonal projection was performed  
 971 for significant coherence values at Y coordinates less than 5 mm away from the selected slice  
 972 Y coordinate. The figure also presents a spectral distribution of coherence at a selection of 10  
 973 sensors of maximum 2–5 Hz coherence (highlighted in magenta) in *video-only* and *control-*  
 974 *video-only*. Gray traces represent individual subject's spectra at the sensor of maximum 2–5  
 975 Hz coherence within the preselection, and the thick black trace is their group average.  
 976 Finally, temporal response functions (TRF) to mouth opening are presented for the two  
 977 significant sources of peak entrainment to mouth opening in *video-only*.

978

979



980

981 **Figure 9.** Visual-to-auditory-speech delays at ~0.5-Hz. Temporal response function of  
982 auditory speech envelope filtered through 0.2–1.5 Hz associated with the time course of  
983 mouth opening (*left*), 2–5-Hz envelope of mouth opening (*middle*), and global visual changes  
984 in video stimuli. There is one gray trace per video (8 in total), and thick red traces are the  
985 average across them all.

986