

1 **Auditory detection is modulated by theta phase of silent lip movements**

2 Emmanuel Biau<sup>1, 2,\*</sup>, Danying Wang<sup>1, 2</sup>, Hyojin Park<sup>1, 2</sup>, Ole Jensen<sup>1, 2</sup>, Simon Hanslmayr<sup>1, 2</sup>.

3 <sup>1</sup> School of Psychology, University of Birmingham, Edgbaston, Birmingham, UK.

4 <sup>2</sup> Centre for Human Brain Health, University of Birmingham, Birmingham, UK

5 **CONTACT INFO:** e.biau@bham.ac.uk

6 **LEAD CONTACT:** Emmanuel Biau

7 **KEYWORDS:** Lip movements; theta oscillations; entrainment; auditory processing.

8

9

10

11

12

13

14

15

16

17

18

19

20

## 21 **SUMMARY**

22 Audiovisual speech perception relies on our expertise to map a speaker's lip movements  
23 with speech sounds. This multimodal matching is facilitated by salient syllable features that  
24 align lip movements and acoustic envelope signals in the 4 - 8 Hz theta band  
25 (Chandrasekaran et al., 2009). The predominance of theta rhythms in speech processing has  
26 been firmly established by studies showing that neural oscillations track the acoustic  
27 envelope in the primary auditory cortex (Giraud & Poeppel, 2012). Equivalently, theta  
28 oscillations in the visual cortex entrain to lip movements (Park et al., 2016), and the auditory  
29 cortex is recruited during silent speech perception (Bourguignon et al., 2020; Cross et al.,  
30 2019; Calvert et al., 1997). These findings suggest that neuronal theta oscillations play a  
31 functional role in organising information flow across visual and auditory sensory areas. We  
32 presented silent speech movies while participants performed a pure tone detection task to  
33 test whether entrainment to lip movements enslaves the auditory system and drives  
34 behavioural outcomes. We showed that auditory detection varied depending on the  
35 ongoing theta phase conveyed by lip movements in the movies. In a complementary  
36 experiment presenting the same movies while recording participants' electro-  
37 encephalogram (EEG), we found that silent lip movements entrained neural oscillations in  
38 the visual and auditory cortices with the visual phase leading the auditory phase. These  
39 results support the idea that the visual cortex entrained by lip movements increases the  
40 sensitivity of the auditory cortex at relevant time-windows for speech comprehension as a  
41 filtering modulator relying on theta phase synchronisation.

## 42 **RESULTS**

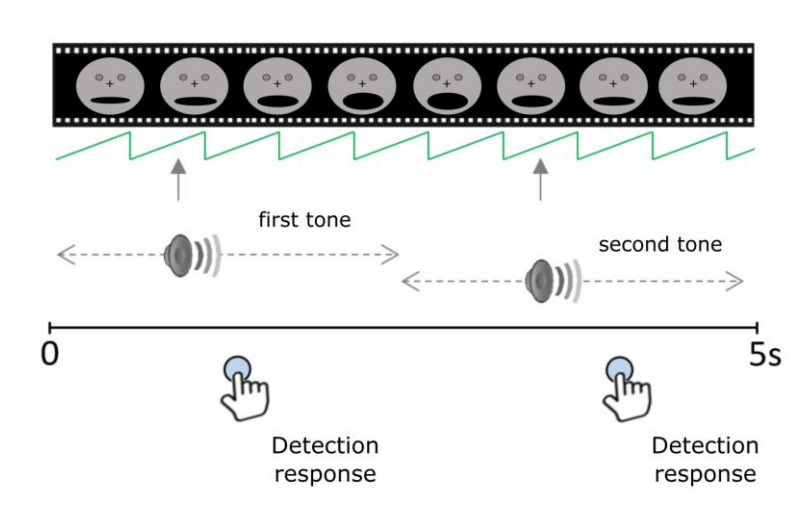
43 When hearing gets difficult, people often visually focus on their interlocutors' mouth to  
44 match lip movements with sounds and to improve speech perception. Mouth opening  
45 indeed shares common features with auditory speech envelope, which temporally  
46 synchronize on dominant 4 - 8 Hz theta rhythms imposed by syllables (Park et al., 2016;  
47 Chandrasekaran et al., 2009; Luo & Poeppel, 2007). Neural oscillations in the auditory cortex  
48 track the auditory envelope structure during speech perception, suggesting that this  
49 "entrainment" reflects signal analysis (Keitel et al., 2018; Pelle & Davis, 2012; Gross et al.,  
50 2013; Giraud & Poeppel, 2012). Although the term entrainment is currently under debate

51 (Meyer, Sun & Martin, 2019; Obleser & Keyser, 2019; Haegens & Zion Golumbic, 2018;  
52 Rimmele et al., 2018), here we use it to describe neural patterns tracking salient features  
53 conveyed in speech signals which occur at theta frequency (4-8 Hz). Previous studies  
54 demonstrated that the visual perception of silent moving lips entrains theta oscillations in  
55 the visual cortex and recruits auditory processing regions (Bourguignon et al., 2020; Cross et  
56 al., 2019), even in the absence of sound (Cross et al., 2015; Calvert et al., 1997). Further,  
57 information specific to lip movements is represented not only in the visual cortex but also in  
58 the auditory cortex (Park et al., 2018). These results beg the question of whether visual  
59 perception of lip movements modulates the auditory cortex in a functional way. In other  
60 words, do purely visually induced theta speech rhythms impose time windows that render  
61 the auditory cortex more sensitive to input in a phasic manner? If the answer to this  
62 question is yes, then visually focussing on your interlocutor's mouth when you have trouble  
63 understanding them would indeed be an effective filter modulator to increase auditory  
64 sensitivity.

### 65 **Entrainment to lip movements during silent speech drives behavioural performance**

66 To address this question, we adapted an auditory tone detection paradigm in which a  
67 continuous white noise was presented simultaneously with silent movies displaying  
68 speakers engaged in conversations (Figure 1 and STAR Methods). Participants were  
69 instructed to press a key as fast and accurate as possible every time when they detected a  
70 pure tone (1 kHz, 100 ms) embedded in the white noise at individual threshold (determined  
71 with a calibration task). In the condition of interest, there were two target tones: the first  
72 tone occurred randomly in the first half of the trial (0 to 2.5 s after trial onset; early window)  
73 and the second tone occurred randomly in the second half of the trial (2.5 to 5 s; late  
74 window). Two additional conditions containing zero or one single tone were introduced to  
75 estimate the false alarm rates (FA) and to reduce the predictability of the second tone by  
76 the occurrence of the first one. The three conditions were counterbalanced and randomised  
77 across six blocks of 50 trials (100 trials per condition). To test the first hypothesis of visual  
78 entrainment affecting auditory processing, participants were asked to attend carefully to  
79 the silent movies centred on the speakers' nose displayed with sound albeit non-  
80 informative. Crucially, the videos were preselected such that lip movements occurred in the

81 4 - 8 Hz theta range. We determined at which theta frequency the vertical mouth's  
82 apertures and auditory speech envelope showed significant dependencies in the original  
83 clips by using mutual information method (see STAR Methods). This paradigm allowed us to  
84 link directly the onset of detected tones with the phase of the ongoing theta activity  
85 conveyed by the lip movements. As neural entrainment increases over time (Thut et al.,  
86 2011; Hanslmayr, Axmacher & Inman, 2019), we compared behavioural performance  
87 between the early and late time-windows (containing respectively the first and second  
88 tones).

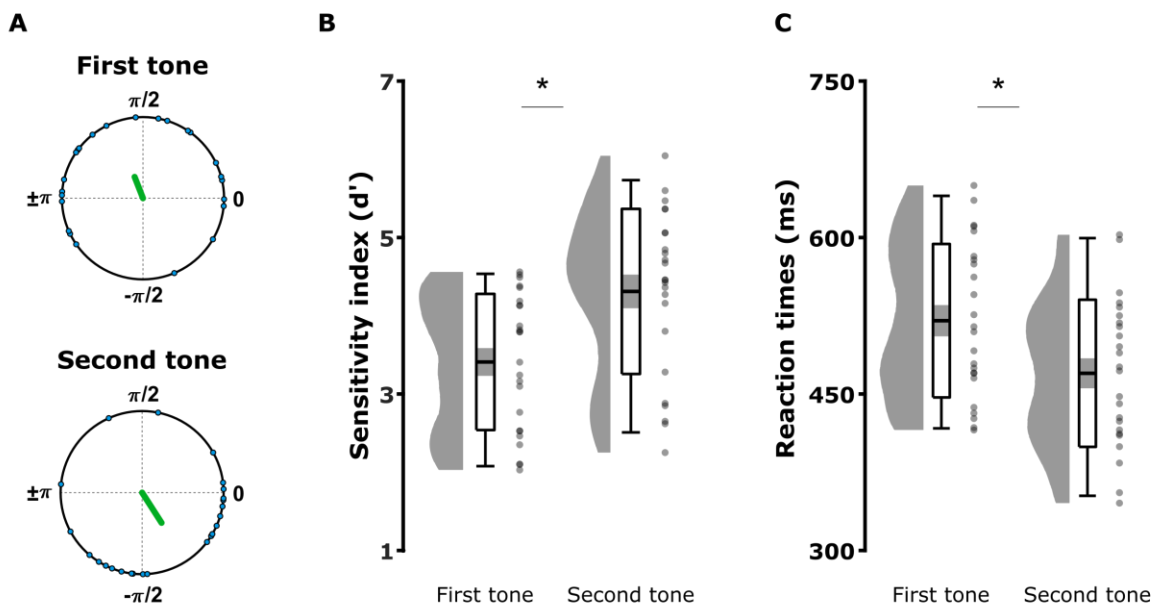


89

90 **Figure 1: Experimental Paradigm of the Tone Detection Task (TDT).** For each trial,  
91 continuous white noise and a silent movie were presented together during five seconds  
92 (drawing of a speaking face substituted to illustrate the relevant concept here). A first pure  
93 tone occurred randomly in the first half of the trial while a second tone occurred randomly  
94 in the second half of the trial. Participants were instructed to respond as fast and accurately  
95 as possible whenever they detected a tone. In the one tone condition, the white noise track  
96 contained only one tone that occurred randomly between the two halves of the trial. In the  
97 zero tone condition, the sound of the trial contained only white noise. The green line  
98 represents the ongoing theta phase conveyed by lip movements.

99 We compared the mean theta phase distributions between first and second tones' onsets  
100 across participants (Figure 2A; see STAR Methods). For each participant, the corresponding  
101 theta phases in ongoing lip activity at detected first and second tone onsets were averaged  
102 across hit trials. Individual mean theta phases were then averaged across subjects to

103 estimate phase locking of hits to the theta signal conveyed visually in the first and second  
104 tone time-windows. Two Rayleigh's uniformity tests were performed on the first and second  
105 grand average theta phase distributions separately. For the first tone window, the Rayleigh's  
106 test did not reject the hypothesis of uniform distribution ( $n = 24$ ;  $\mu = 1.944$  rad or  $111.384^\circ$ ;  
107  $r = 0.282$ ;  $p = 0.148$ , Bonferroni-corrected). In contrast, the Rayleigh's test revealed that  
108 mean phases were not uniformly distributed in the second tone window ( $n = 24$ ;  $\mu = -0.999$   
109 rad or  $302.763^\circ$ ;  $r = 0.44$ ;  $p < 0.01$ , Bonferroni-corrected). Further, a permutation test was  
110 performed on the resultant vector length ( $r$ ) difference between the first and second tones  
111 to test whether the strength of visual entrainment in the second tone window was  
112 significantly stronger than in the first tone window, which indeed was the case  
113 (permutations: 10,000; effect size = 0.158;  $p = 0.015$ ; Figure 2A; see STAR Methods). We  
114 performed two additional permutation tests on resultant vector length difference between  
115 hits and misses in the first and second tone windows separately to test that visual  
116 entrainment related to successful auditory processing. No significant difference of vector  
117 length was found in the first tone window (permutations: 10,000; effect size = 0.041;  $p =$   
118 0.405), whereas in the second tone window the resultant hits vector strongly tended to be  
119 longer than the misses vector (permutations: 10,000; effect size = 0.228;  $p = 0.052$ ).



120

121 **Figure 2: Visual entrainment and tone detection performance in the two tones condition.**

122 (A) Resultant  $r$  vector length (green line) from grand average phase at the onsets of

123 correctly detected first and second tones across participants. The individual mean theta  
124 phases are depicted in polar coordinates (blue circles). (B) Mean sensitivity index ( $d' = Z_{\text{Hit rate}}$   
125  $- Z_{\text{FA rate}}$ ) and (C) reaction times of first and second tone hits. The graphs depict the density,  
126 the grand average (mean  $\pm$  standard error of the mean; errors bars: 95 % confidence  
127 interval), and individual means (grey dots) for first/second tones. Significant contrasts are  
128 evidenced with stars ( $p < 0.05$ ).

129 Following up, we investigated whether tone detection differed between the first and  
130 second tones windows. Such a difference might reflect an auditory bias by visual inputs  
131 (Figure 2B). First, two independent one-sample t-tests established that participants  
132 detected the first and second tones in the two tones condition, as the  $d'$  scores were greater  
133 than zero (first tone:  $T(1,23) = 20.014$ ;  $p < 0.001$ , two-tailed; second tone:  $T(1,23) = 21.124$ ;  
134  $p < 0.001$ , two-tailed). Second, a paired-samples t-test confirmed that the second tones  
135 were better detected than the first ones ( $T(1, 23) = -4.488$ ;  $p < 0.001$ ; two-tailed; Fig. 2B).  
136 Third, a paired-sample t-test applied on the hit reaction times showed that participants  
137 responded faster to second compared to first tones ( $T(1, 23) = 5.486$ ;  $p < 0.001$ ; two-tailed;  
138 Figure 2C). Importantly, the improvement of the second tone detection could not be  
139 attributed to a simple attentional effect due to the presence of the preceding first one, as  
140 the single tone condition replicated the two tones condition performances (i.e. by sorting  
141 the single tones as first/second tones according to their onsets; see Figure S1B). Finally,  
142 a paired-samples t-test performed on the FA rates in the no tone condition confirmed that  
143 detection performance modulations did not reflect a change in response bias between the  
144 two windows ( $T(1, 23) = 0.627$ ;  $p = 0.537$ ; two-tailed). Altogether, these results established  
145 that entrainment to theta lips activity increased in time and coincided temporally with  
146 increases in auditory detection. In the next step, we aimed at establishing whether a  
147 potential audiovisual communication relying on the critical theta organisation of  
148 information flows was reflected in the brain.

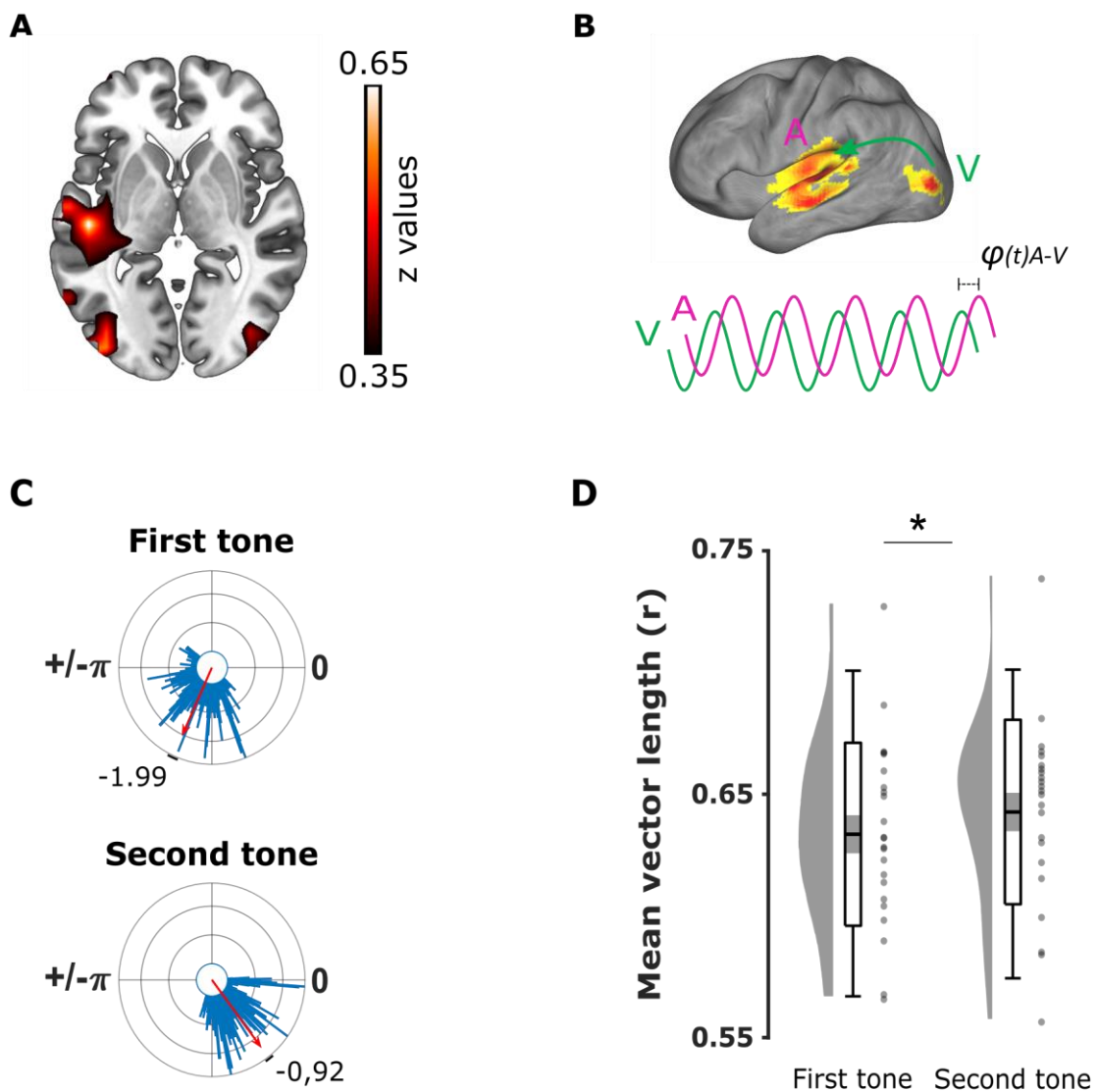
### 149 **Visual cortex leads synchronization to left auditory cortex via theta oscillations during** 150 **silent lips perception**

151 The above results suggest that visual speech stimuli may recruit the auditory regions via  
152 entrainment to render some time-windows more sensible to auditory detection than others.

153 To test this hypothesis on a neural level, we recorded the EEG signal of 23 new participants  
154 during the perception of the same 60 silent movies used in the previous tone detection task.  
155 Participants were instructed to attend to each movie and rate its emotional content based  
156 on the speaker's face. The movies were presented in a single block and randomised across  
157 participants. First, the sources of interest responding to speakers' lip movements were  
158 identified applying a linearly constrained minimum variance beamforming method. Neural  
159 entrainment to lip movements was estimated by computing mutual information (MI) on the  
160 theta phase between the EEG epochs and corresponding lip signals in the equivalent first (0  
161 to 2.5 s) and second (2.5 to 5 s) tone windows. Just as in the behavioural data, we assessed  
162 whether entrainment increased over time by contrasting the difference of MI between the  
163 first and second time window. Second, the EEG data at the identified visual and auditory  
164 sources were reconstructed to perform single-trial phase coupling analysis. The synchrony  
165 between visual and auditory sources was reflected by the distribution of theta phase angle  
166 differences  $\phi_{A-V} = \phi_{\text{audio}} - \phi_{\text{visual}}$  at each time-point within the first and second tone windows,  
167 and the directionality of the coupling was evidenced with the sign of  $\phi_{A-V}$  (i.e. a mean  
168 distribution of  $\phi_{A-V} = 0$  would mean perfect phase alignment, while  $\phi_{A-V} < 0$  would mean  
169 that the visual phase leads the auditory phase; see STAR Methods).

170 Source localisation analysis revealed that the maximum increases in  $MI_{\text{second}}$  as compared  
171 to  $MI_{\text{first}}$  were localised in the expected left visual and auditory cortices, as well as in the  
172 right visual cortex to a lesser extent (Figure 3A). This result supports the recruitment of  
173 both sensory areas during the perception of speakers' lip movements even in the absence of  
174 speech sound. Two separate Rayleigh tests confirmed non-uniform distributions of  $\phi_{A-V}$  in  
175 the first ( $n = 23$ ;  $\mu = -1.99$  rad or  $-114.45^\circ$ ;  $r = 0.768$ ;  $p < 0.001$ , Bonferroni-corrected) and  
176 second tone windows ( $n = 23$ ;  $\mu = -0.92$  rad or  $-52.79^\circ$ ;  $r = 0.875$ ;  $p < 0.001$ , Bonferroni-  
177 corrected). An additional Kuiper two-sample test confirmed that the mean  $\phi_{A-V}$  distributions  
178 between the first and second tone windows converged towards two different preferred  
179 angles ( $k = 3.614 \times 10^5$ ;  $p < 0.001$ ). Further, an increase in theta phase synchrony between  
180 visual and auditory areas would be reflected by a more consistent distribution of  $\phi_{A-V}$   
181 towards zero degree. To quantify the modulation of phase coupling with entrainment, we  
182 computed the resultant vector length  $r$  of the distance between the observed  $\phi_{A-V}$  in the  
183 data and a fixed zero  $\phi_{A-V}$  in the first and second windows separately (zero  $\phi_{A-V} = 0$ ; meaning

184 that visual and auditory theta phases are perfectly aligned with a constant offset of zero at  
185 each time-point of the time-window). A paired-samples t-test showed that the resultant  
186 vector length  $r$  of the distance between the observed  $\phi_{A-V}$  and the zero  $\phi_{A-V}$  was significantly  
187 greater in the second tone window ( $T(1, 22) = -2.135$ ;  $p = 0.044$ ; two-tailed), confirming that  
188 synchrony between auditory and visual sources improved with time (Figure 3B, C and D).  
189 The negative theta phase angle differences  $\phi_{A-V}$  in both the first and second tone windows  
190 confirmed that the visual phase led the auditory phase, in line with the idea of visual  
191 oscillations responding first to the lips inputs and then enslaving theta oscillations in the  
192 auditory cortex. Altogether, these results support our hypothesis that visual cortex led  
193 synchronization to left auditory cortex via theta oscillations during silent lips perception.



194



195 **Figure 3: Theta phase coupling analysis between visual and auditory areas during lips**  
196 **perception.** (A) Difference of mutual information between the second (2.5 - 5s) and first (0 -  
197 2.5s) time-windows ( $MI_{\text{second}} > MI_{\text{first}}$  contrast; z values). Auditory (MNI coordinates of  
198 maximum voxel: [-50 -21 0]; Temporal middle left) and visual (MNI coordinates of maximum  
199 voxel: [-40 -89 0]; Occipital middle Left) sources were localized in the left hemisphere. (B)  
200  $MI_{\text{first}} > MI_{\text{second}}$  contrast projected on brain's surface for illustrative purpose:  
201 Synchronisation was estimated through  $\phi_{A-V}$  theta phase offset between theta oscillations at  
202 identified auditory (pink line) and visual sources (green line) by mean of phase coupling  
203 analysis. (C) Audio-visual phase coupling in the first and second time-windows. The mean  
204  $\phi_{A-V}$  offset between auditory and visual theta phases (red arrows) confirmed that  
205 oscillations entrained by lip movements in the visual cortex preceded oscillations in the  
206 auditory cortex by 1.99 rad ( $\sim 56.10$  ms) and 0.92 rad ( $\sim 25.88$  ms), respectively in the first  
207 and second windows. (D) Theta synchronisation between visual and auditory areas improves  
208 with entrainment. The resultant vector length  $r$  of the distance between the observed  $\phi_{A-V}$   
209 and the theoretical  $\phi_{A-V} = 0$  was greater in the second than the first window, reflecting a  
210 more consistent distribution of  $\phi_{A-V}$ . The graphs depict the density, the grand average (mean  
211  $\pm$  standard error of the mean; errors bars: 95 % confidence interval), and individual resultant  
212 vector length  $r$  (grey dots). Significance evidenced with a star ( $p < 0.05$ ).

## 213 **DISCUSSION**

214 In two complementary experiments, we first established that visual entrainment to theta  
215 lip phase modulated auditory detection, even if information from silent movies was  
216 irrelevant to perform the task. Second, the perception of silent lip movements entrained  
217 theta oscillations in the visual cortex, which in turn synchronized with the auditory cortex.  
218 Together, these results suggest that the brain's natural reaction to visual speech stimuli  
219 might be to align the excitability of the auditory cortex with sharp mouth-openings because  
220 that is when one expects to hear corresponding acoustic syllable edges (Hickock & Poeppel,  
221 2007; Giraud & Poeppel, 2012; Peelle & Sommers, 2015 Park et al., 2016; Chandrasekaran  
222 et al., 2009). Such a neural process could be a very effective filtering method to increase the  
223 sensitivity of the auditory cortex in these relevant time windows for speech comprehension.

224 Our EEG results suggest that theta oscillations in the left visual cortex encoded the lips'  
225 activity first. Then information travelled to the left auditory cortex via phase coupling to  
226 shape its activity. Previous findings reported that the auditory cortex tracks both auditory  
227 and visual stimulus dynamics using low-frequency neuronal phase modulation during  
228 audiovisual movie perception (Luo, Liu and Poeppel, 2010). Other studies reported that the  
229 perception of silent lips also recruited the auditory regions (Bourguignon et al., 2020; Cross  
230 et al., 2015; Calvert et al., 1997). Our findings go beyond and establish how theta  
231 oscillations orchestrate visual and auditory cortices through phase coupling to ensure cross-  
232 region communication even in a unimodal condition. Furthermore, it is commonly agreed  
233 that entrainment takes several cycles from rhythmic inputs to build up (Doelling et al. 2014;  
234 Lakatos et al., 2008; Thut et al., 2011; Zoefel et al. 2018). Behavioural and neural indicators  
235 of entrainment reported here consistently increased from the first to the second half time-  
236 window of the trial in both experiments. This supports the idea that we indeed observed  
237 neural entrainment to lip movements and sheds light on the functional relevance of visual  
238 inputs modulating auditory theta rhythms.

239 As visual onsets naturally lead corresponding auditory onsets by 100-to-300 ms in  
240 audiovisual speech (Chandrasekaran et al., 2009; van Wassenhove et al., 2005; Pilling,  
241 2009), visual entrainment to lips may act as a filter by increasing excitability in the auditory  
242 cortex to windows containing relevant acoustic features. This hypothesis is corroborated by  
243 our phase coupling results where visual theta phase systematically led auditory theta phase  
244 during silent movie presentation. However, whether such filtering reflected direct  
245 enslavement of the auditory cortex or involved top-down modulations remains unclear.  
246 Indeed, higher-level sensorimotor areas also activate during speech perception (Park et al.,  
247 2016; 2018; Cognan & Poeppel, 2011; Pulvermüller et al., 2005; Wilson et al., 2004).  
248 Assaneo & Poeppel (2018) demonstrated recently that activity in the motor and auditory  
249 cortices couple at theta rate during syllable perception, correlating with the strength of  
250 coupling between speech signal and EEG in the auditory cortex. On the other hand, motor  
251 areas play a role in temporal analysis of rhythmic sensory stimulation (Biau & Kotz, 2018;  
252 Arnal et al., 2015; Fujioka et al., 2015; Morillon et al., 2019). Entrainment to lip movements  
253 may provide the temporal theta structure of speech signal to motor cortex, which in turn  
254 adjusts auditory excitability at critical windows containing the corresponding acoustic

255 features in a top-down fashion (in line with Park et al. 2015). Alternatively, mouth-opening  
256 perception may target internal articulatory representations and help to identify the  
257 corresponding sounds in the auditory signal. Theta activity in the auditory cortex would  
258 reflect the contribution of endogenous oscillations bearing linguistic inferences generated  
259 from motor representations activation. Although speculative, this could partially explain  
260 why the increase of entrainment in the auditory cortex was left-lateralized, i.e. by recruiting  
261 language-related representations classically associated with the left hemisphere. This  
262 hypothesis fits in recent debates on whether neural tracking during speech processing  
263 reflects online cooperation between pure entrainment to external salient features and  
264 endogenous rhythms providing abstract representations (Meyer, Sun & Martin, 2019;  
265 Obleser & Keyser, 2019; Haegens & Zion Golumbic, 2018; Rimmele et al., 2018). However,  
266 this would not explain why visual speech information improved the detection of unrelated  
267 pure tones here, which will be addressed in future experiments. Additional data-driven  
268 analysis suggested that two subpopulations of participants showed distinct visual theta  
269 phases shaping auditory perception in the TDT (Figure S2). One could hypothesize that the  
270 “good” subpopulation (i.e. group 2) were fine-tuned to a preferred visual theta phase that  
271 represents an optimal time-window. This optimal window allowed information to travel to  
272 the auditory cortex (either directly or via top-down modulations), and reset auditory activity  
273 at “perfect” moments when a tone occurred. Back to our filtering hypothesis, visual theta  
274 entrainment would increase auditory excitability coinciding with more time windows  
275 containing a tone in this “good” subpopulation regardless of the nature of sounds.

276 As a final note, although the auditory signal alone often provides enough structural  
277 information for the early analytic steps of continuous speech, e.g. telephone conversations,  
278 a visual filter may be especially helpful to sharpen auditory perception when hearing is  
279 impaired or in elders (Grant et al., 1998). Our results provide an important step toward  
280 understanding how visual information functionally drives auditory speech perception, and  
281 suggest future directions to investigate hearing loss compensation, i.e. to improve lip-  
282 reading along with hearing correction.

## 283 **ACKNOWLEDGMENTS**

284 This work was supported by a sir Henry Wellcome Postdoctoral Fellowship awarded to E.B  
285 (Grant reference number: 210924/Z/18/Z), as well as grants from the ERC (Consolidator  
286 Grant 647954) and ESRC (ES/R010072/1) awarded to S.H, who is further supported by the  
287 Wolfson Foundation and Royal Society. The authors would like to thank people from the  
288 Memory and Attention Lab and David Poeppel for their valuable comments and inputs  
289 during the preparation of this manuscript.

## 290 **AUTHOR CONTRIBUTION**

291 E.B, H.P and S.H designed the experiments and paradigms. E.B and D.W collected and  
292 analysed the data. E.B, D.W, H.P, O.J and S.H wrote the paper. All the authors discussed the  
293 results and commented on the manuscript.

## 294 **DECLARATION OF INTEREST**

295 The authors of this manuscript declare to have no conflicts of interest.

## 296 **STAR METHODS**

297 Key resources table:

REAGENT or	SOURCE	IDENTIFIER
Software and Algorithms		
MATLAB	The MathWorks	R2018a
Psychophysics Toolbox	<a href="http://psychtoolbox.org">http://psychtoolbox.org</a>	3
FieldTrip	<a href="http://www.fieldtriptoolbox.org">http://www.fieldtriptoolbox.org</a>	v.20161231
SPM8	Wellcome Trust Centre for Neuroimaging	8
ASIO4All	Steinberg Media Technologies	2.12
ActiView	BioSemi B.V. Amsterdam, Netherlands	7
Shotcut	Meltytech, LLC	v.18.06.02
Brainstorm Toolbox	<a href="https://neuroimage.usc.edu/brainstorm/">https://neuroimage.usc.edu/brainstorm/</a>	
CARET	Washington University School of Medicine	5.65
Circular Statistics Toolbox	<a href="https://uk.mathworks.com/matlabcentral/fileexchange/10676-circular-statistics-toolbox-directional-statistics">https://uk.mathworks.com/matlabcentral/fileexchange/10676-circular-statistics-toolbox-directional-statistics</a>	v.1.21.0.0
Other		
BioSemi ActiveTwo	BioSemi B.V. Amsterdam, Netherlands	EEG system

system		
ER-3C system	Etymotic Research, Elk Grove Village, IL	EEG compatible earphones
Fastrak	Polhemus, Colchester, VT, USA	Electromagnetic digitiser
ThorLabs DET36A	<a href="https://thorlabs.de">https://thorlabs.de</a>	Photodetector

298 **Contact for reagent and resource sharing:**

299 Further information and requests for resources and data should be directed to and will be  
300 fulfilled by the Lead Contact, Emmanuel Biau (e.biau@bham.ac.uk). Summarized data (cell  
301 means) are available; data for individual participants are available, as consent for sharing  
302 data at the level of the individual participant was received.

303 **Experimental model and subject details:**

304 Tone detection experiment: Twenty-eight healthy English native speakers (mean age = 19  
305 years  $\pm$  0.69; 21 females) took part in the first behavioural experiment. Five participants  
306 were left-handed. All of them reported normal or corrected-to-normal vision and hearing.  
307 All participants were granted experimental participation credit. The data from four  
308 participants were excluded because of extreme overall performances and the final analysis  
309 were applied on twenty-four data sets.

310 **Silent movie perception-EEG experiment:**

311 Twenty-five healthy English native speakers (mean age = 21.52 years  $\pm$  3.86; 17 females)  
312 took part in the first behavioural experiment. All of them reported normal or corrected-to-  
313 normal vision and hearing, and were right-handed. Twenty-one participants were granted  
314 credits and five participants received financial compensation for their participation (£20).  
315 The data from two participants were excluded from the final analyses due to too noisy EEG  
316 data. In the two experiments, all the participants signed informed consent and ethical  
317 approval was granted by the University of Birmingham Research Ethics Committee,  
318 complying with the Declaration of Helsinki.

319 **Method details:**

320 **Apparatus:**

321 The two tasks were programmed with Matlab (R2018a; The MathWorks, Natick, MA, USA)  
322 and presented with Psychophysics Toolbox (Brainard, 1997; Pelli, 1997; Kleiner et al., 2007).  
323 In both tasks, the silent videos were presented on a 21-inch CRT display with a screen  
324 refresh rate of 75 Hz (nVidia Quadro K600 graphics card: 875 MHz graphics clock, 1024 MB  
325 dedicated graphics memory; Nvidia, Santa Clara, CA, USA). The auditory stimuli in the tone  
326 detection task were presented through EEG-compatible insert earphones (ER-3C; Etymotic  
327 Research, Elk Grove Village, IL). In the tone detection task, the accuracy of movie and sound  
328 presentation timing was optimised by detecting a small white square displayed on the left of  
329 the first frame of each visual stimulus with a photodiode (ThorLabs DET36A, thorlabs.de),  
330 and Psychophysics Toolbox (PsychPort Audio and ASIO4ALL extensions for Matlab).  
331 Additionally, a parallel audio port was used to record the online audio signal of each trial  
332 during presentation. Continuous photodiode and audio data during trials were recorded  
333 through a BioSemi Analog Input Box (AIB) adding two separate channel inputs into BioSemi  
334 ActiveTwo system: the BioSemi AD-box was connected with the AIB through optical fibres.  
335 The input from the photodiode was connected through a BNC connector and the input from  
336 the microphone was connected through the 3.5 mm audio. Those two inputs were  
337 connected to the AIB through a 37 pin Sub-D connector. Data were digitized using the  
338 BioSemi ActiView software, with a sampling rate of 2048 Hz. Offline analysis were  
339 performed to calculate the real delay between visual and audio stimuli offset using in-house  
340 Matlab codes. Any lag between visual and auditory stimuli onsets was later compensated in  
341 the data analyses when computing the corresponding visual theta phase to the tones  
342 onsets. The experiments were run from a solid-state hard drive on a Windows 7-based PC  
343 (3.40 GHz processor, 16 Gb RAM). Participants used a standard computer keyboard to  
344 respond to the tasks.

345 Stimuli of the Tone detection Task:

346 Movies:

347 Sixty five-second movies were extracted from natural face-to-face interviews published on  
348 YouTube ([www.youtube.com](http://www.youtube.com)) by various universities channels and downloaded via free  
349 online application. Satisfying movies containing meaningful content (i.e. one complete  
350 sentence, speaker facing toward the camera) were edited using Shotcut (Melttytech, LLC).

351 For each movie, the video and the sound were exported separately (Video: .mp4 format,  
352 1280 x 720 resolution, 25 frame per second, 200 ms linear ramp fade in/out; Audio: .wav  
353 format, 44100 Hz sampling rate, mono).

354 Lip movements' detection:

355 Lips contour signal was extracted for each video using in-house Matlab codes. We computed  
356 the area information (area contained within the lips contour), the major axis information  
357 (horizontal axis within lip contour) and minor axis information (vertical axis within lip  
358 contour) as described in Park et al. (2016). In the present study, we used vertical aperture  
359 information of the lips contour to establish the theta correspondence between lips and  
360 auditory speech (i.e. aperture between the superior and inferior lips) but using area  
361 information gave very similar results, as also reported in Park et al. (2016). The lips time-  
362 series was resampled at 250 Hz for further analyses with corresponding auditory speech  
363 envelope.

364 Auditory speech signal:

365 The amplitude envelope of each movie sound was computed using in-house Matlab codes  
366 (Park et al., 2018; 2016; Chandrasekaran et al., 2009). First, eight frequency bands  
367 equidistant on the cochlear map in the range 100–10,000 Hz were constructed (Smith et al.,  
368 2002). Then, sound signals were then band-pass filtered in these bands with a fourth-order  
369 Butterworth filter (forward and reverse). Hilbert transform was applied to obtain amplitude  
370 envelopes for each band. These signals were then averaged across bands and resulted in a  
371 unique wideband amplitude envelope per sound signal. Each final signal was resampled to  
372 250 Hz for further theta correspondence analyses.

373 Mutual information between lip movements and corresponding auditory speech signal:

374 To identify the main oscillatory activity conveyed by the lip movements in each visual  
375 stimulus, we determined at which theta frequency the auditory and visual speech signals  
376 showed significant dependencies. To do so, we examined the audiovisual speech frequency  
377 spectrum (1 to 20 Hz) and computed the mutual information (MI) between the minor axis  
378 information and speech envelope signals sampled at 250 Hz. MI measures the statistical

379 dependence between two variables with no prior hypothesis, and with a meaningful effect  
380 size measured in bits (Ince et al., 2017; Shannon, 1948). We applied the Gaussian Copula  
381 Mutual Information approach described in Ince et al. (2017) in which the MI between two  
382 signals corresponds to the negative entropy of their joint copula transformed distribution.  
383 This method provides a robust, semiparametric lower bound estimator of MI by combining  
384 the statistical theory of copulas together with the closed-form solution for the entropy of  
385 Gaussian variables, allowing good estimation over circular variables, like phase as well as  
386 power. For each movie, the complex spectrum is normalized by its amplitude to obtain a 2D  
387 representation of the phase as points lying on the unit circle for both the lip movements and  
388 auditory envelope time-series. The real and imaginary parts of the normalized spectrums  
389 are rank-normalized separately and the phase dependence for each frequency between the  
390 two 2D signals is estimated using the multivariate GCM estimator giving a lower bound  
391 estimate of the MI between the phases of the two signals. Here, we applied the GCM  
392 analyses in two conditions to determine the frequency of interest in each movie: first, we  
393 computed MI between corresponding lips and envelope signals as well as non-matching  
394 signals (i.e. lips time-series paired with random auditory envelope signals). For the matching  
395 signals, the averaged MI spectrum revealed a greater peak in the expected 4 - 8 Hz theta  
396 frequencies, reflected by a bump in the band of interest. In contrast, there was no  
397 relationship between random auditory and visual signal pairs, which depicts a flat line  
398 profile along the whole spectrum (see Supplementary Information Figure S3 A). These  
399 results are well in line with previous studies using coherence or MI measures, and confirm  
400 the temporal coupling between lips and auditory speech streams at the expected syllable  
401 rate in our videos (Park et al., 2016; 2018; Chandrasekaran et al., 2009). Second, for each  
402 movie, we performed a peak detection on the MI spectrum to determine which specific  
403 frequency carried most theta information to maximize entrainment in the tone detection  
404 and silent movie perception tasks (4Hz frequency peak: 16 videos; 5Hz frequency peak: 15  
405 videos; 6Hz frequency peak: 9 videos; 7Hz frequency peak: 13 videos; 8Hz frequency peak: 7  
406 videos. See Supplementary Information Figure S3 B).

407 Audio tones and white noise:



408 Pure auditory tones and white noise stimuli were generated using in-house Matlab codes.  
409 The target tone consisted in a sinusoidal signal of 100 ms at one kHz (sampling rate: 44100  
410 Hz). The same noise consisted in a Gaussian white noise lasting two seconds for the  
411 calibration task and five seconds in the tone detection task (the white noise has been  
412 generated only once and loaded during each procedure to ensure that all the participants  
413 were tested with the same noise; sampling rate: 44100 Hz). Both the tone and the white  
414 noise signals were normalized between - 1 to 1 (arbitrary units).

415 Tones onsets:

416 For each trial, the target tones were embodied in the white noise at predetermined pseudo-  
417 random onsets counterbalanced across conditions (zero, one or two tones per trial, 100  
418 trials per condition). In the calibration task serving to determine the individual threshold of  
419 target tones detection (see below for the general procedure), there could be only zero or  
420 one tone maximum per trial. For the one tone condition, the onset of the target tone always  
421 randomly occurred between 300 and 1400 ms after the trial onset to allow participants to  
422 detect it properly and have time to respond before the end of the trial. In the zero tone  
423 condition, the auditory track consisted in two seconds of white noise only. In the tone  
424 detection task, there could be zero, one or two tones per trial. In the one tone condition,  
425 the onset of the target always occurred randomly between 300 and 4500 ms after the trial  
426 onset. In the two tones condition, the first tone randomly occurred in a time-window  
427 centred on the first half of the trial length, between 300 and 3000 ms (mean first tone  
428 onsets =  $1.68 \pm 0.78$  s). The second tone occurred in a time-window centred on the second  
429 half of the trial length, between 1000 ms after the first tone onset and 4500 ms after the  
430 trial onset (mean first tone onsets =  $3.45 \pm 0.62$  s). This design provided participants with  
431 enough time to detect and respond to both tones, and kept the two tones temporally  
432 unrelated from each other. In the zero tone condition, the auditory track consisted in five-  
433 second of white noise only. The signal-to-noise ratio between target tones and white noise  
434 was determined for each participant individually with the calibration task performances and  
435 adjusted consequently in the following tone detection task (see below).

436 Procedure of the calibration task and tone detection task (TDT):

437 The experiment began after the completion of a safety-screening questionnaire and the  
438 provision of informed consent. Participants sat in a well-lit testing room at approximately  
439 60 cm from the centre of the screen and wore the insert earphones for sound presentation.  
440 Participants performed first a short pure tone detection task with no visual stimuli (i.e.  
441 calibration task). This task served to determine the individual threshold at which each  
442 participant detected  $\sim 70 - 80$  % of the target tones in auditory modality only, and the  
443 signal-to-noise ratio (SNR) to be implemented between the amplitude of the target tones  
444 and the white noise in the following tone detection task (TDT). The calibration task was  
445 composed of a four-trial practice to identify the target tone itself, followed by five blocks  
446 containing 20 trials each. Each trial began with a black fixation cross (500 - 1000 ms  
447 duration, jittered) followed by the presentation of a red cross over a grey background during  
448 two seconds to indicate the period of possible target tones occurrence. A continuous white  
449 noise was displayed during the red cross presentation. In 50 % of the trials, a unique audio  
450 tone was embedded in the white noise at unpredictable onset, and participants had to press  
451 "1" key as fast and accurately as possible only when they perceived a target tone. The  
452 pseudo-random sequence of the procedure ensured that there were never more than two  
453 consecutive trials of the same condition. The participants received no feedback and the  
454 procedure continued to the next trial after the end of the two-second white noise. The  
455 signal-to-noise ratio was adjusted following an adapted *two-down, one-up* staircase  
456 procedure (see Leek, 2001): For the first five trials, the SNR was fixed (mean white noise  
457 power of 0.981) and served as a starting point across participants. After each trial, the  
458 keypress response of the participant was stored to adjust the SNR for the next trial as  
459 following: for two successive hits, the SNR was decreased by 2 % of the starting signal  
460 energy in the next trial. For two successive correct rejections (i.e. no response when no tone  
461 occurred) or one correct rejection following a hit, the SNR was kept identical for the next  
462 trial. After a miss or a false alarm, the SNR was always increased by 2 % of the starting signal  
463 energy. At the end of the calibration task, the individual SNR was averaged over the last 30  
464 trials and stored for the following real tone detection task (mean calibration accuracy rate:  
465  $0.75 \pm 0.05$ ). The participants took a short break and were recalled the instructions before  
466 starting the proper tone detection task. The calibration task lasted approximately seven  
467 minutes.

468 The main structure of the TDT was the same as in the precedent calibration task. The TDT  
469 was composed of a short four-trial practiced followed by 300 trials divided in 6 blocks of 50  
470 trials each and separated by breaks (the sixty silent movies were repeated five times each to  
471 generate the total 300 trials). Each trial began with a red fixation cross presentation (500-  
472 1250 ms duration, jittered). Then, a random five-second silent movie was presented with a  
473 black fixation cross in the centre of the screen to give the participants a point to gaze at and  
474 reduce saccades. The continuous white noise was displayed together with the silent movie  
475 according to the three random conditions: no tone (100 trials), one single tone (100 trials) or  
476 two tones (100 trials) hidden in the white noise. Participants were instructed to press “1”  
477 key as fast and accurately as possible only when they perceived a target tone. The  
478 participants received no feedback on their responses and the procedure continued with the  
479 next trial after the end of the silent movie. The SNR between the tones and the white noise  
480 was determined in the previous calibration task as explained above. The TDT lasted  
481 approximatively 50 minutes.

482 TDT conditions:

483 The condition of interest containing the two tones (i.e. first and second tone) served to  
484 assess our main hypothesis that entrainment increases in time with the perception of visual  
485 information conveyed by the speakers’ lip movements. According to this, the second tones  
486 should be better detected and associated to a greater theta entrainment as compared to  
487 the first tones to reflect enslavement of the auditory system by the entrained visual system  
488 to lip movements. The zero and single tone conditions were additional control conditions:  
489 the zero tone condition served to determine the false alarm rates (i.e. participants’  
490 keypresses in the absence of tone) and controlled whether participants tended to press  
491 more together with the tone onset delays (i.e. time-dependent response bias). The single  
492 tone condition served to counterbalance the number of trials containing two tones and  
493 control for the predictability of the second tone. The replication of the performances  
494 observed in the two tones condition by sorting the single tones according to their onsets  
495 equivalent to either first or second tone onsets would confirm that the detection of the  
496 second tone is not due to its predictability from a preceding tone but its position in time

497 only. The pseudo-random sequence of the procedure ensured that there were never more  
498 than three consecutive trials of the same condition.

499 Perception of silent movies-EEG task:

500 Movies:

501 The movies presented during the silent lips perception task were the exact same 60 movies  
502 used in the previous tone detection task. The order of movies was randomized across  
503 participants.

504 Procedure:

505 Participants sat in a well-lit testing room at approximately 60 cm from the centre of the  
506 screen to complete a safety-screening questionnaire and the provision of informed consent  
507 first. After the correct preparation of the EEG cap, the participants were instructed to attend  
508 to all the movies quietly and to avoid movements during the presentation. Each trial was  
509 preceded by a central fixation cross (500 - 1250 ms duration, jittered) followed by the  
510 presentation of a random five-second movie. A central fixation cross was displayed during  
511 the movie presentation to give participants a point to gaze at and reduce excessive  
512 saccades. Participants were instructed to attend to each movie carefully and rate its  
513 emotional content based on speaker's facial gestures by using the number keys on the  
514 keyboard after the presentation (i.e. 1 for neutral through 5 for very emotional; results not  
515 reported). The total presentation of the sixty movies lasted approximately 10 minutes.

516 Online EEG recordings: Continuous EEG signal was recorded using a 128 channel BioSemi  
517 ActiveTwo system (BioSemi, Amsterdam, Netherlands). Vertical and horizontal eye  
518 movements were recorded from additional electrodes placed approximately one cm to  
519 the left of the left eye, one cm to the right of the right eye, and one cm below the left eye.  
520 Online EEG signals were digitalized using BioSemi ActiView software at a sampling rate of  
521 2048 Hz. For each participant, the position of the electrodes on the scalp were tracked using  
522 a Polhemus FASTRAK device (Colchester) and recorded with Brainstorm (Tadel et al., 2011)  
523 implemented in MATLAB.

524 Offline EEG preprocessing: EEG data were preprocessed offline using Fieldtrip (Oostenveld  
525 et al., 2011) and SPM8 toolboxes (Wellcome Trust Centre for Neuroimaging). Continuous  
526 EEG signals were bandpass filtered between one and 100 Hz and bandstop filtered (48–52  
527 Hz and 98–102 Hz) to remove line noise at 50 and 100 Hz. Data were epoched from 2000 ms  
528 before stimulus onset to 7000 ms after stimulus onset, and downsampled to 512 Hz. Bad  
529 trials and channels with artefacts were excluded by visual inspection and numerical criteria  
530 (e.g., variance as well as kurtosis) before applying an independent component analysis (ICA)  
531 to remove components related to ocular artefacts. Bad channels were then interpolated  
532 using the method of triangulation of nearest. After re-referencing the data to average  
533 reference, trials with artefacts were manually rejected by a last visual inspection. On  
534 average,  $4.48 \pm 2.48$  trials were removed and  $4.04 \pm 1.82$  channels were interpolated per  
535 participants.

536 Head models:

537 For the 22 participants without individual MRI scans, the MNI-MRI and the volume  
538 conduction templates provided by Fieldtrip were used to construct the head models.  
539 Electrode positions of each participant were aligned to the template head model. Source  
540 models were prepared with the template volume conduction model and the aligned  
541 individuals' electrode positions following standard procedures. One participant provided his  
542 own MRI scans and his head model was built using his structural scans (Michelmann et al.,  
543 2016): the MRI scans were segmented into four layers (i.e. brain, CSF, skull and scalp) using  
544 the Statistical Parametric Mapping 8 (SPM8; <http://www.fil.ion.ucl.ac.uk/spm>) and Huang  
545 toolboxes (Huang et al., 2013). The volume conduction model was constructed using the  
546 dipoli method implemented in Fieldtrip. Participant's electrode positions were aligned to his  
547 individual head model. Finally, his MRI was warped into the same MNI template MRI of  
548 Fieldtrip and the inverse of the warp was applied to a template dipole grid to have each grid  
549 point position in the same normalized MNI space as the other participants for further group  
550 analyses.

551 Source localization during silent movie perception:

552 Source analyses on EEG data recorded during silent movies presentation were run using  
553 individual electrode positions, grid positions and template volume conduction model. For  
554 the participant who had his MRI scans, source analyses were calculated using normalized  
555 grid positions instead. Source activity was reconstructed using a linearly constrained  
556 minimum variance beamforming method implemented in Fieldtrip (LCMV; see Van Veen et  
557 al., 1997). The neural entrainment to lip movements at source level was determined by  
558 computing mutual information between EEG epochs and the lip movements during silent  
559 movie presentation (i.e. lips time-series from the silent movie presented during the trial). To  
560 test our hypothesis that entrainment builds up in time with perceived theta lips activity, we  
561 contrasted the difference of MI between the equivalent time-window to the second tone  
562 window ( $MI_{\text{second}}$ ), and the equivalent time-window to the first tone window ( $MI_{\text{first}}$ ) in the  
563 previous TDT. Accordingly, we expected first to observe an increase of theta activity in the  
564 visual cortex reflecting entrainment to lip movements. Second, we expected an equivalent  
565 pattern in the auditory correlates reflecting a tuning from visual activity. For each single  
566 trial, MI was first computed separately in the equivalent first (0 to 2.5 seconds after trial  
567 onset;  $MI_{\text{first}}$ ) and second time-windows (2.5 to 5 seconds after trial onset;  $MI_{\text{second}}$ ) at the  
568 2020 virtual electrodes by using the same approach described in the stimuli analysis section  
569 (i.e. where we established which frequency carried most correspondence between lips and  
570 envelope signals for each video; using a wavelet transform to compute the phase). Second,  
571 for each single trial, the MI spectrum was realigned respect to the frequency bin ( $\pm 2$  Hz)  
572 corresponding to the peak of MI between lips and envelope signal established in the movie  
573 analyses. This step was done to be able to average all the trials together taking into account  
574 the main theta activity carried in each individual movie. For instance, if the peak of MI  
575 between lip movements and auditory envelope was found at 4 Hz in the video number 1,  
576 the realigned MI spectrum between EEG and lips signals from the trials presenting video  
577 number 1 was now  $4 \pm 2$  Hz (2 to 6 Hz; 1 Hz bin) to insure that the central bin of each single  
578 trial corresponds to the objectively determined frequency peak of theta activity. Third, the  
579 realigned MIs of single trials were averaged across trials within each participant for further  
580 group analyses. For each participant, we calculated the normalized difference of MI at the  
581 frequency bin of interest (MI normalization:  $(MI_{\text{second}} - MI_{\text{first}}) / MI_{\text{first}}$ ; third bin in the  
582 realigned MI spectrum) in the equivalent second tone window minus equivalent first tone

583 window at all the 2020 virtual electrodes. Finally, the normalized difference of MI between  
584 second and first tone time-windows was grand averaged across participants and the grand  
585 average was interpolated to the MNI MRI template. The coordinates for auditory and visual  
586 sources of interest were determined by finding the maximum of  $MI_{\text{second}} - MI_{\text{first}}$  differences  
587 in regions corresponding to the auditory and visual areas, and defined using the automated  
588 anatomical labelling atlas (AAL).

589 Source reconstruction:

590 We performed time-series reconstruction analysis to investigate the synchronization at  
591 theta activity between the two sources of interest during silent movie presentation. The  
592 time series data were reconstructed and extracted at the visual and auditory coordinates  
593 determined by source localization analysis. LCMV beamformer reconstruction can cause  
594 random direction of source dipoles and eventually affect phase analysis results. To get  
595 around this issue, the event-related potentials (ERP) time-locked to movie onsets at visual  
596 and auditory sources were plotted to identify the visual component, i.e. N1-P2-N2  
597 waveform (Wang et al., 2018). After visual inspection, the sign of the reconstructed data  
598 were flipped in direction by multiplying the time-series by -1 if any visual or auditory source  
599 ERPs showed the opposite of the expected direction of a visual component (i.e. negative-  
600 positive-negative polarity). This “flipping” correction was applied consistently across all trials  
601 before sorting data between first and second time-windows, thus it did not bias results  
602 towards our hypothesis. The same phase coupling analyses were computed with unflipped  
603 source data as a control. Phase angle differences between visual and auditory theta  
604 activities in the first and second windows were also non-uniformly distributed according to  
605 Rayleigh tests with significantly different mean angles according to a Kuiper’s test,  
606 confirming that the flipping procedure only better reflected phase coupling modulation with  
607 entrainment.

608 Theta phase coupling between auditory and visual sources:

609 First, auditory signal was projected orthogonally onto the visual signal applying a Gram-  
610 Schmidt process (GSP; Hipp et al., 2012) for single trials before computing phase  
611 information. This was done to reduce the noise correlation patterns reflecting activity from

612 a common source (i.e. volume conduction) estimate captured at different electrodes (in that  
613 case, the phase alignment reflects the same source activity and not the phase coupling  
614 between two distinct source activities). The GSP increases the signal-to-noise ratio by  
615 leaving intact the proper activities conveyed at the two distinct electrodes while reducing  
616 noise correlation weight (see Hipp et al., 2012). Second, for each trial the instantaneous  
617 theta phase of the auditory and visual orthogonalized time-series were computed by  
618 applying a Hilbert transform with a bandpass filter centred on the frequency bin of MI peak  
619  $\pm 2$  Hz, accordingly to the mean theta frequency of the video presented during the trial.  
620 Third, the difference of unwrapped instantaneous phase between auditory and visual  
621 sources was computed for each single trial at each time-point in two windows  
622 corresponding to the first (0.5 to 2 seconds after trial onset) and second tone windows (3 to  
623 4.5 seconds after trial onset). The first and last 500 ms at the edges of the epoch were not  
624 included into phase coupling analyses to avoid the trial onset and offset responses.  
625 Additionally, the phase-slope index (PSI) was calculated in the first and second time-  
626 windows between the left auditory and visual sources to estimate the directionality of  
627 information flow between our two sources of interest, using Fieltrip procedure (Nolte et al.,  
628 2008). PSI analysis revealed negative values in both time-windows (respectively  $\text{psi}_{\text{first tone}} = -$   
629  $0.035 \pm 0.055$  and  $\text{psi}_{\text{second tone}} = -0.041 \pm 0.035$ ), confirming that that the left visual source  
630 led the left auditory source during silent movie perception.

631 Quantification and statistical analysis:

632 The Tone Detection and the Silent lips perception tasks were within-subject design.

633 Tone detection performance:

634 Tone detection performances: The hits (i.e. correctly detected tones) and false alarms (i.e.  
635 keypress responses during the zero tone condition allocated to the first or second tone  
636 window depending on their onsets) rates were computed to calculate the individual mean  
637 sensitivity index (i.e.  $d'$ ) in the two conditions for each participant (i.e. single tone and two  
638 tones conditions). The reaction times of the hits were computed to calculate the individual  
639 mean reaction times in the two conditions for each participant (i.e. single tone and two  
640 tones conditions). Additionally, we calculated the mean correct response rates and reaction



641 times of the two conditions concatenated together of each individual to exclude blindly  
642 potential outliers without favouring the results towards our hypothesis and performing as  
643 following: below chance level (correct response rate < 0.5) or perfectly (correct response  
644 rate = 1), or with mean reaction times outside the grand averaged reaction times  $\pm$  two  
645 standard deviations range. Accordingly, four participants were excluded from analyses (two  
646 participants performed below chance level, one participant performed perfectly and one  
647 participant's reaction times were slower than the grand average mean + two standard  
648 deviations). A paired-samples t-test was conducted on the averaged  $d'$  scores and hit  
649 reaction times between the first and second tones in the two tones condition and single  
650 tone condition separately. Additionally, a paired-samples t-test was performed on false  
651 alarm rates from the first and second windows in the zero tone condition to control for any  
652 response bias with time.

653 Visual entrainment to theta activity conveyed by lip movements:

654 To bridge visual entrainment to auditory processing together, we related the tone target  
655 onsets to the theta activity conveyed by the lip movements during silent movies perception:  
656 First, for each movie the theta phase of the lip movements' time-series was computed by  
657 applying a Hilbert transform with a bandpass filter centred on the frequency bin of MI peak  
658  $\pm$  2 Hz, accordingly to the mean theta frequency determined in MI stimuli analyses. Second,  
659 we computed the instantaneous theta phase of the lips signal corresponding to the onset of  
660 the tones occurring during each trial. All further circular statistics on angular scale were  
661 performed using the CircStat toolbox on Matlab (Berens, 2009). The circular uniformity in  
662 the first and second tones windows within and across participants were estimated  
663 separately by applying Rayleigh tests to calculate the mean direction and resultant vector  
664 length from hits/miss trials. To assess statistically the strength of visual entrainment  
665 between the first and second tone windows in the two tones condition (hits only), we  
666 performed a permutation test on the resultant vector length difference (z-value) second  
667 tone minus first tone reflecting the effect size. For each participant, we generated 10000  
668 iterations as following: first, the hit trial labels were shuffled between the first and second  
669 tones in the two tones condition. Second, two balanced subsamples of shuffled trials were  
670 selected, with a number matching the smallest number of trials between the first and

671 second tone hits. Third, the mean phase of the first and second tone shuffled trials were  
672 computed for each iteration and per participant. Fourth, a Rayleigh's test of uniformity was  
673 applied on the mean phases to determine a resultant vector length at the first and second  
674 tones per iteration (i.e. z-value). For each iteration, we computed the difference of z-  
675 value<sub>second tone</sub> - z-value<sub>first tone</sub> to quantify its effect size, and the resultant 10000 z-value  
676 differences were sorted in descending order. To estimate the final p-value and test the null  
677 hypothesis, the difference of z-value between the original first and second tone data was  
678 ranked in the sorted permuted z-value differences and divided by the total number of  
679 permutations+1. If the p-value was smaller than  $\alpha = 0.05$ , we rejected the null hypothesis  $H_0$   
680 = there is no difference of resultant vector length between the first and second tones (i.e.  
681 the visual entrainment is significantly greater in the second tone window). The exact same  
682 approach was applied on the single tone condition, as well as on the hit versus miss  
683 entrainment comparisons.

684 Subpopulations of group 1 and group 2:

685 Participants were sorted in two subpopulations according to their preferred theta phase in  
686 the second tone window, where visual entrainment supposedly took place after enough lip  
687 movements inputs in the condition of interest (i.e. two tones condition). The Rayleigh tests  
688 revealed non-uniform distributions of preferred phase at the second tones for group 1 ( $n =$   
689  $11$ ;  $\mu = 348.83^\circ$ ;  $p < 0.001$ ) and group 2 ( $n = 10$ ;  $\mu = 249.63^\circ$ ;  $p < 0.001$ ). A Kuiper two-sample  
690 test confirmed that the mean preferred phases were different between group 1 and 2 ( $k =$   
691  $121$ ;  $p < 0.01$ ).

692 Audio-visual theta synchrony in the silent lips perception-EEG task:

693 The phase coupling between auditory and visual sources was estimated through their theta  
694 phase angle difference in time-windows equivalent to the first and second tone windows  
695 from the TDT. To quantify the improvement of phase coupling with entrainment, we  
696 computed the resultant vector length  $r$  of the distance between the observed  $\phi_{A-V}$  in the  
697 data and a fixed  $\phi_{A-V} = 0$  in the first and second windows separately ( $\phi_{A-V} = 0$  meaning that  
698 visual and auditory theta phases are perfectly aligned with a constant offset of zero at each  
699 time-point of the time-window). A better synchrony between visual and auditory areas

700 would be reflected by a more consistent distribution of  $\phi_{A-V}$  towards a particular angle (i.e.  
701 0). For each trial, we calculated the resultant vector length of the distance between the real  
702 auditory-visual phase offset and the theoretical phase offset  $0^\circ$  at each time point in the two  
703 time-windows. The resultant vector length was collapsed across time in the first and second  
704 windows separately, resulting in two values per trial. Single-trial values in the first and  
705 second windows were then averaged across trials for each participant, and the difference of  
706 phase entrainment values was assessed with a paired samples t-test.

## 707 SUPPLEMENTARY INFORMATION

708 Supplemental Information includes four figures, one silent movie (Video\_1) and one white  
709 noise containing two tones (Audio\_1) used in the TDT and EEG tasks.

## 710 REFERENCES

711 Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., & Ghazanfar, A. A. (2009). The  
712 natural statistics of audiovisual speech. *PLoS Computational Biology*, *5*(7), e1000436.  
713 <https://doi.org/10.1371/journal.pcbi.1000436>

714 Giraud, A.-L., & Poeppel, D. (2012). Cortical oscillations and speech processing: emerging  
715 computational principles and operations. *Nature Neuroscience*, *15*(4), 511–517.  
716 <https://doi.org/10.1038/nn.3063>

717 Park, H., Ince, R. A. A., Schyns, P. G., Thut, G., & Gross, J. (2015). Frontal Top-Down Signals  
718 Increase Coupling of Auditory Low-Frequency Oscillations to Continuous Speech in  
719 Human Listeners. *Current Biology*, *25*(12), 1649–1653.  
720 <https://doi.org/10.1016/j.cub.2015.04.049>

721 Bourguignon, M., Baart, M., Kapnoula, E. C., & Molinaro, N. (2020). Lip-reading enables the  
722 brain to synthesize auditory features of unknown silent speech. *Journal of*  
723 *Neuroscience*, *40*(5), 1053–1065. <https://doi.org/10.1523/jneurosci.1101-19.2019>

724 Crosse, M. J., Butler, J. S., & Lalor, E. C. (2019). Congruent Visual Speech Enhances Cortical  
725 Entrainment to Continuous Auditory Speech in Noise-Free Conditions. *The Journal of*

- 726 *Neuroscience : The Official Journal of the Society for Neuroscience*, 35(42), 14195–  
727 14204. <https://doi.org/10.1523/jneurosci.1829-15.2015>
- 728 Calvert, G. A., Bullmore, E. T., Brammer, M. J., Campbell, R., Williams, S. C. R., McGuire, P. K.,  
729 Woodruff, P.W., Iversen, S.D., & David, A. S. (1997). Activation of auditory cortex during  
730 silent lipreading. *Science*, 276(5312), 593–596.  
731 <https://doi.org/10.1126/science.276.5312.593>
- 732 Park, H., Kayser, C., Thut, G., & Gross, J. (2016). Lip movements entrain the observers' low-  
733 frequency brain oscillations to facilitate speech intelligibility. *ELife*, 5.  
734 <https://doi.org/10.7554/eLife.14521>
- 735 Luo, H., & Poeppel, D. (2007). Phase patterns of neuronal responses reliably discriminate  
736 speech in human auditory cortex. *Neuron*, 54(6), 1001–1010.  
737 <https://doi.org/10.1016/j.neuron.2007.06.004>
- 738 Keitel, A., Gross, J., & Kayser, C. (2018). Perceptually relevant speech tracking in auditory  
739 and motor cortex reflects distinct linguistic features. *PLoS Biology*, 16(3), e2004473.  
740 <https://doi.org/10.1371/journal.pbio.2004473>
- 741 Peelle, J. E., & Davis, M. H. (2012). Neural Oscillations Carry Speech Rhythm through to  
742 Comprehension. *Frontiers in Psychology*, 3, 320.  
743 <https://doi.org/10.3389/fpsyg.2012.00320>
- 744 Gross, J., Hoogenboom, N., Thut, G., Schyns, P., Panzeri, S., Belin, P., & Garrod, S. (2013).  
745 Speech Rhythms and Multiplexed Oscillatory Sensory Coding in the Human Brain. *PLoS*  
746 *Biology*, 11(12), e1001752. <https://doi.org/10.1371/journal.pbio.1001752>
- 747 Meyer, L., Sun, Y., & Martin, A. E. (2019). Synchronous, but not entrained: exogenous and  
748 endogenous cortical rhythms of speech and language processing. *Language, Cognition*  
749 *and Neuroscience*, 1–11. <https://doi.org/10.1080/23273798.2019.1693050>

- 750 Obleser, J., & Kayser, C. (2019). Neural Entrainment and Attentional Selection in the  
751 Listening Brain. *Trends in Cognitive Sciences*, Vol. 23, pp. 913–926.  
752 <https://doi.org/10.1016/j.tics.2019.08.004>
- 753 Haegens S, Zion Golumbic E. (2018). Rhythmic facilitation of sensory processing: A critical  
754 review. *Neurosci Biobehav Rev*, 86:150-165. doi:10.1016/j.neubiorev.2017.12.002
- 755 Rimmele, J. M., Morillon, B., Poeppel, D., & Arnal, L. H. (2018). Proactive Sensing of Periodic and  
756 Aperiodic Auditory Patterns. *Trends in Cognitive Sciences*, 22, 870–882.  
757 <https://doi.org/10.1016/j.tics.2018.08.003>
- 758 Park, H., Ince, R. A. A., Schyns, P. G., Thut, G., & Gross, J. (2018). Representational  
759 interactions during audiovisual speech entrainment: Redundancy in left posterior  
760 superior temporal gyrus and synergy in left motor cortex. *PLoS Biology*, 16(8).
- 761 Crosse, M.J, ElShafei, H.A, Foxe, J.J, and Lalor E.C. (2015). Investigating the Temporal  
762 Dynamics of Auditory Cortical Activation to Silent Lipreading. 7th Annual International  
763 IEEE EMBS Conference on Neural Engineering
- 764 Hanslmayr, S., Axmacher, N., & Inman, C. S. (2019, July 1). Modulating Human Memory via  
765 Entrainment of Brain Oscillations. *Trends in Neurosciences*, Vol. 42, pp. 485–499.  
766 <https://doi.org/10.1016/j.tins.2019.04.004>
- 767 Thut, G., Veniero, D., Romei, V., Miniussi, C., Schyns, P., & Gross, J. (2011). Rhythmic TMS  
768 causes local entrainment of natural oscillatory signatures. *Current Biology*, 21(14),  
769 1176–1185. <https://doi.org/10.1016/j.cub.2011.05.049>
- 770 Hickok, G., & Poeppel, D. (2007, May). The cortical organization of speech processing.  
771 *Nature Reviews Neuroscience*, Vol. 8, pp. 393–402. <https://doi.org/10.1038/nrn2113>
- 772 Peelle, J. E., & Sommers, M. S. (2015). Prediction and constraint in audiovisual speech  
773 perception. *Cortex; a Journal Devoted to the Study of the Nervous System and Behavior*,  
774 68, 169–181. <https://doi.org/10.1016/j.cortex.2015.03.006>

- 775 Luo, H., Liu, Z., & Poeppel, D. (2010). Auditory cortex tracks both auditory and visual  
776 stimulus dynamics using low-frequency neuronal phase modulation. *PLoS Biology*, 8(8),  
777 25–26. <https://doi.org/10.1371/journal.pbio.1000445>
- 778 Doelling, K. B., Arnal, L. H., Ghitza, O., & Poeppel, D. (2014). Acoustic landmarks drive delta-  
779 theta oscillations to enable speech comprehension by facilitating perceptual parsing.  
780 *NeuroImage*, 85 Pt 2, 761–768. <https://doi.org/10.1016/j.neuroimage.2013.06.035>
- 781 Lakatos, P., Karmos, G., Mehta, A. D., Ulbert, I., & Schroeder, C. E. (2008). Entrainment of  
782 neuronal oscillations as a mechanism of attentional selection. *Science (New York, N.Y.)*,  
783 320(5872), 110–113. <https://doi.org/10.1126/science.1154735>
- 784 Zoefel, B., Archer-Boyd, A., & Davis, M. H. (2018). Phase Entrainment of Brain Oscillations  
785 Causally Modulates Neural Responses to Intelligible Speech. *Current Biology*, 28(3),  
786 401-408.e5. <https://doi.org/10.1016/j.cub.2017.11.071>
- 787 van Wassenhove, V., Grant, K. W., & Poeppel, D. (2005). Visual speech speeds up the neural  
788 processing of auditory speech. *Proceedings of the National Academy of Sciences of the*  
789 *United States of America*, 102(4), 1181–1186.  
790 <https://doi.org/10.1073/pnas.0408949102>
- 791 Pilling, M. (2009). Auditory event-related potentials (ERPs) in audiovisual speech perception.  
792 *Journal of Speech, Language, and Hearing Research: JSLHR*, 52(4), 1073–1081.  
793 [https://doi.org/10.1044/1092-4388\(2009/07-0276\)](https://doi.org/10.1044/1092-4388(2009/07-0276))
- 794 Cogan, G. B., & Poeppel, D. (2011). A mutual information analysis of neural coding of speech  
795 by low-frequency MEG phase information. *Journal of Neurophysiology*, 106(2), 554–  
796 563. <https://doi.org/10.1152/jn.00075.2011>
- 797 Pulvermüller, F., Huss, M., Kherif, F., Del Prado Martin, F. M., Hauk, O., & Shtyrov, Y. (2006).  
798 Motor cortex maps articulatory features of speech sounds. *Proceedings of the National*  
799 *Academy of Sciences of the United States of America*, 103(20), 7865–7870.  
800 <https://doi.org/10.1073/pnas.0509989103>

- 801 Wilson, S. M., Saygin, A. P., Sereno, M. I., & Iacoboni, M. (2004). Listening to speech  
802 activates motor areas involved in speech production. *Nature Neuroscience*, 7(7), 701–  
803 702. <https://doi.org/10.1038/nn1263>
- 804 Assaneo, F., M., & Poeppel, D. (2018). The coupling between auditory and motor cortices is  
805 rate-restricted: Evidence for an intrinsic speech-motor rhythm. *Sci Adv*, 4(2). DOI:  
806 10.1126/sciadv.aao3842
- 807 Biau, E., & Kotz, S. A. (2018). Lower beta: A central coordinator of temporal prediction in  
808 multimodal speech. *Frontiers in Human Neuroscience*, Vol. 12.  
809 <https://doi.org/10.3389/fnhum.2018.00434>
- 810 Arnal, L. H., Doelling, K. B., & Poeppel, D. (2015). Delta-Beta Coupled Oscillations Underlie  
811 Temporal Prediction Accuracy. *Cerebral Cortex*, 25(9), 3077–3085.  
812 <https://doi.org/10.1093/cercor/bhu103>
- 813 Fujioka, T., Ross, B., & Trainor, L. J. (2015). Beta-Band Oscillations Represent Auditory Beat  
814 and Its Metrical Hierarchy in Perception and Imagery. *Journal of Neuroscience*, 35(45),  
815 15187–15198. <https://doi.org/10.1523/jneurosci.2397-15.2015>
- 816 Morillon, B., Arnal, L. H., Schroeder, C. E., & Keitel, A. (2019, December 1). Prominence of  
817 delta oscillatory rhythms in the motor cortex and their relevance for auditory and  
818 speech perception. *Neuroscience and Biobehavioral Reviews*, Vol. 107, pp. 136–142.  
819 <https://doi.org/10.1016/j.neubiorev.2019.09.012>
- 820 Grant, K. W., Walden, B. E., & Seitz, P. F. (1998). Auditory-visual speech recognition by  
821 hearing-impaired subjects: consonant recognition, sentence recognition, and auditory-  
822 visual integration. *The Journal of the Acoustical Society of America*, 103(5 Pt 1), 2677–  
823 2690. <https://doi.org/10.1121/1.422788>
- 824 Brainard, D.H. (1997). The psychophysics toolbox. *Spat. Vis.* 10, 433–436
- 825 Pelli, D.G. (1997). The VideoToolbox software for visual psychophysics: transforming  
826 numbers into movies. *Spat. Vis.* 10, 437–442

- 827 Kleiner, M., Brainard, D., Pelli, D., Ingling, A., Murray, R., and Broussard, C (2007). What's  
828 new in psychtoolbox-3. *Perception*, 36, 1–16
- 829 Smith, Z. M., Delgutte, B., & Oxenham, A. J. (2002). Chimaeric sounds reveal dichotomies in  
830 auditory perception. *Nature*, 416(6876), 87–90. <https://doi.org/10.1038/416087a>
- 831 Ince, R. A. A., Giordano, B. L., Kayser, C., Rousselet, G. A., Gross, J., & Schyns, P. G. (2017). A  
832 statistical framework for neuroimaging data analysis based on mutual information  
833 estimated via a gaussian copula. *Human Brain Mapping*, 38(3), 1541–1573.  
834 <https://doi.org/10.1002/hbm.23471>
- 835 Shannon C. E. (1948). The mathematical theory of communication. *Bell Syst Tech J*, 27:379
- 836 Leek, M. R. (2001). Adaptive procedures in psychophysical research. *Perception and*  
837 *Psychophysics*, 63(8), 1279–1292. <https://doi.org/10.3758/BF03194543>
- 838 Tadel F, Baillet S, Mosher JC, Pantazis D, Leahy RM (2011) Brainstorm: a user-friendly  
839 application for MEG/EEG analysis. *Comput Intell Neurosci*, 2011:879716
- 840 Oostenveld R, Fries P, Maris E, Schoffelen JM (2011) FieldTrip: open source software for  
841 advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput Intell*  
842 *Neurosci*, 2011:156869
- 843 Michelmann S, Bowman H, Hanslmayr S (2016) The temporal signature of memories:  
844 identification of a general mechanism for dynamic memory replay in humans. *PLOS*  
845 *Biol*, 14:e1002528
- 846 Huang Y, Dmochowski JP, Su Y, Datta A, Rorden C, Parra LC (2013) Automated MRI  
847 segmentation for individualized modeling of current flow in the human head. *J Neural*  
848 *Eng*, 10:066004
- 849 Van Veen BD, van Drongelen W, Yuchtman M, Suzuki A (1997) Localization of brain electrical  
850 activity via linearly constrained minimum variance spatial filtering. *IEEE Trans Biomed*  
851 *Eng*, 44:867– 880



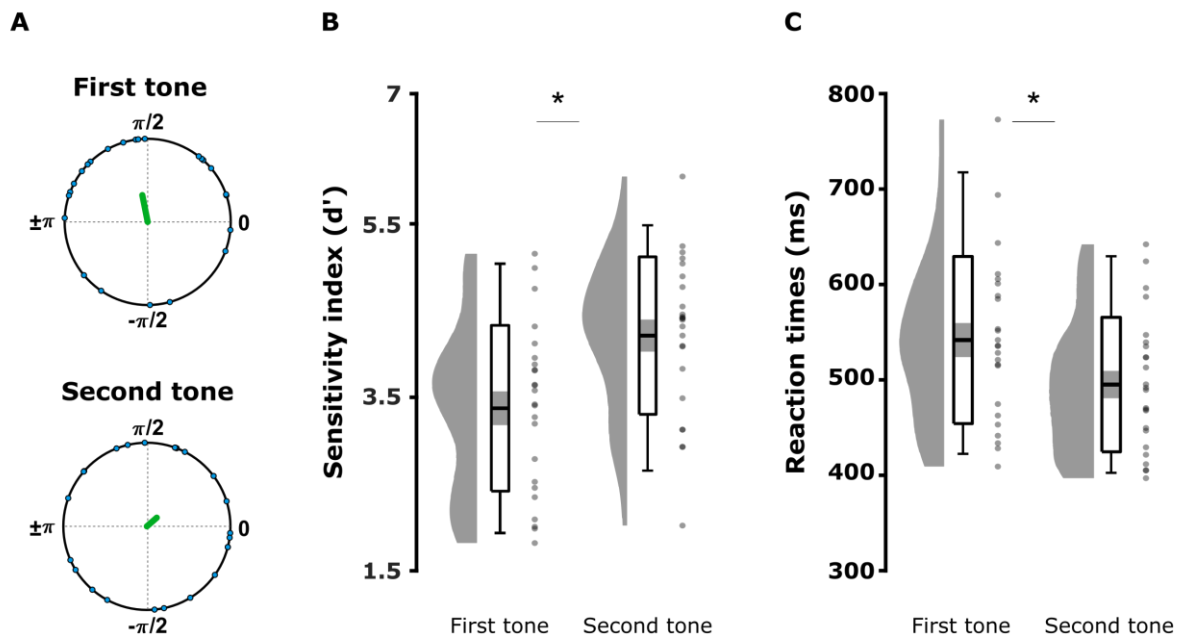
852 Wang, D., Clouter, A., Chen, Q., Shapiro, K. L., & Hanslmayr, S. (2018). Single-Trial Phase  
853 Entrainment of Theta Oscillations in Sensory Regions Predicts Human Associative  
854 Memory Performance. *The Journal of Neuroscience : The Official Journal of the Society*  
855 *for Neuroscience*, 38(28), 6299–6309. [https://doi.org/10.1523/JNEUROSCI.0349-](https://doi.org/10.1523/JNEUROSCI.0349-18.2018)  
856 18.2018

857 Hipp, J. F., Hawellek, D. J., Corbetta, M., Siegel, M., & Engel, A. K. (2012). Large-scale cortical  
858 correlation structure of spontaneous oscillatory activity. *Nature Neuroscience*, 15(6),  
859 884–890. <https://doi.org/10.1038/nn.3101>

860 Berens, P. (2009). CircStat: A MATLAB Toolbox for Circular Statistics. *Journal of Statistical*  
861 *Software*, 31(10), 1 - 21. doi:<http://dx.doi.org/10.18637/jss.v031.i10>.

## 862 SUPPLEMENTARY INFORMATION

863



864

865 **Figure S1: Visual entrainment and tone detection performance in the single tone**  
866 **condition.** Each single tone was sorted as a first or second tone according to its onset, i.e.  
867 between 0 and 2.5 s or 2.5 and 5s after trial onset. (A) Resultant  $r$  vector length (green line)  
868 from grand average phase at the onsets of first tones hits ( $n = 24$ ;  $\mu = 1.769$  rad or  $101.336^\circ$ ;  
869  $r_{\text{first}} = 0.329$ ;  $p = 0.074$ ) and second tones hits ( $n = 24$ ;  $\mu = 0.719$  rad or  $41.165^\circ$ ;  $r_{\text{second}} =$

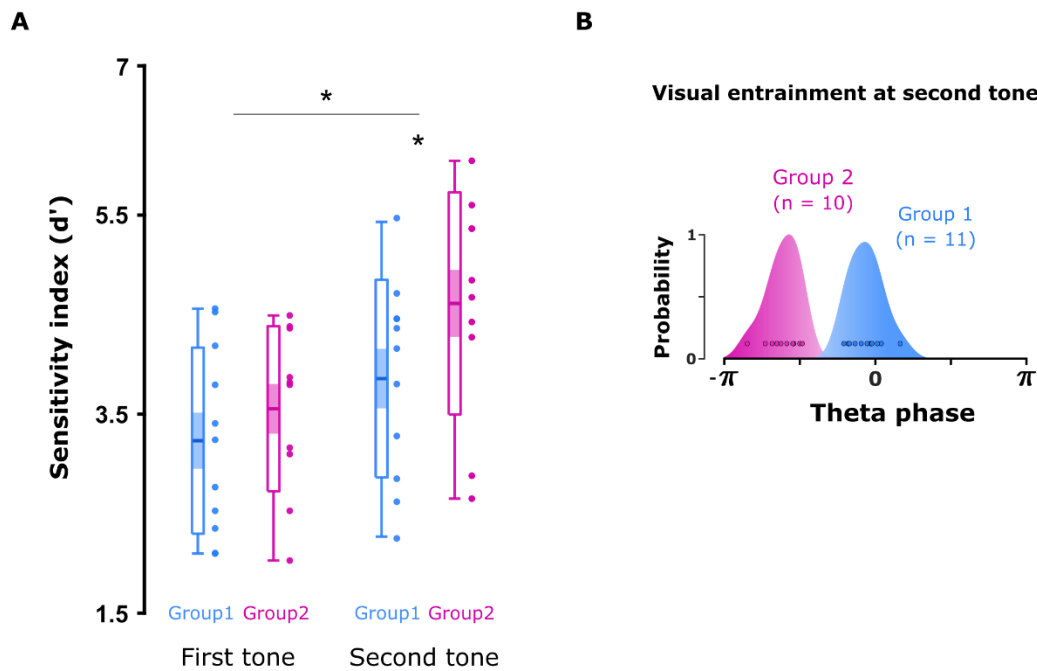
870 0.157;  $p = 0.558$ ). The individual mean theta phases are depicted in polar coordinates (blue  
871 circles). (B) Mean sensitivity index ( $d'$ ) and (C) reaction times. The graphs depict the density,  
872 the grand average (mean  $\pm$  standard error of the mean; errors bars: 95 % confidence  
873 interval), and individual means (grey dots) for tones sorted as first or second ones.  
874 Significant contrasts are evidenced with stars. A permutation test on the resultant vector  
875 length difference between the tones sorted as first and second tones did not reveal  
876 significant difference (figure S1 A; permutations: 10000; effect size = - 0.172;  $p = 0.835$ ). Two  
877 paired-samples t-test performed on  $d'$  scores and reaction times confirmed that the second  
878 tones were better detected (figure S1 B;  $T(1, 23) = -4.114$ ;  $p < 0.001$ ; two-tailed), and faster  
879 as compared to the first tones (figure S1 C;  $T(1, 23) = 4.778$ ;  $p < 0.001$ ; two-tailed). Finally,  
880 we compared the tone detection performances between the single tone and two tones  
881 conditions by mean of 2-by-2 repeated-measures ANOVAs (factors condition and tone  
882 position). A main effect of position showed that the second tones were better detected than  
883 the first tones in both conditions ( $F(1, 23) = 19.174$ ;  $p < 0.001$ ). No main effect of condition  
884 ( $F(1, 23) = 1.017$ ;  $p = 0.324$ ) or interaction between condition and tone position on  $d'$  were  
885 found ( $F(1, 23) = 0.567$ ;  $p = 0.459$ ). A repeated-measures ANOVA on reaction times showed  
886 a significant effect of tone position with faster responses to second tones than first tones  
887 ( $F(1, 23) = 33.797$ ;  $p < 0.001$ ). Additionally, a significant effect of condition showed overall  
888 faster reaction times in the two tones condition as compared to the single tone condition  
889 ( $F(1, 23) = 14.047$ ;  $p = 0.001$ ), but no interaction between tone position and condition ( $F(1,$   
890  $23) = 0.173$ ;  $p = 0.682$ ).

891

### 892 **Distinct preferred phases showed performance differences between two subpopulations** 893 **of listeners.**

894 The behavioural data of the TDT task suggests that two separate subpopulations  
895 entrained to different preferred theta phases in the second tone time-window (see Figure  
896 2A lower panel and Figure S2B below). In a post-hoc analysis, we assessed whether these  
897 apparently distinct populations also showed differences in tone detection performances.  
898 Arguably, any difference should be most pronounced only when visual entrainment  
899 eventually took place (second tone window) but not early in the trial. Participants were

900 sorted in two groups based on their mean theta phase in the second tone window (i.e.  
901  $n_{\text{group1}} = 11$  and  $n_{\text{group2}} = 10$ ; see STAR Methods) and we compared detection performance  
902 ( $d'$ ) in the two tones condition by means of a repeated-measures ANOVA (with factors tone  
903 position and group).



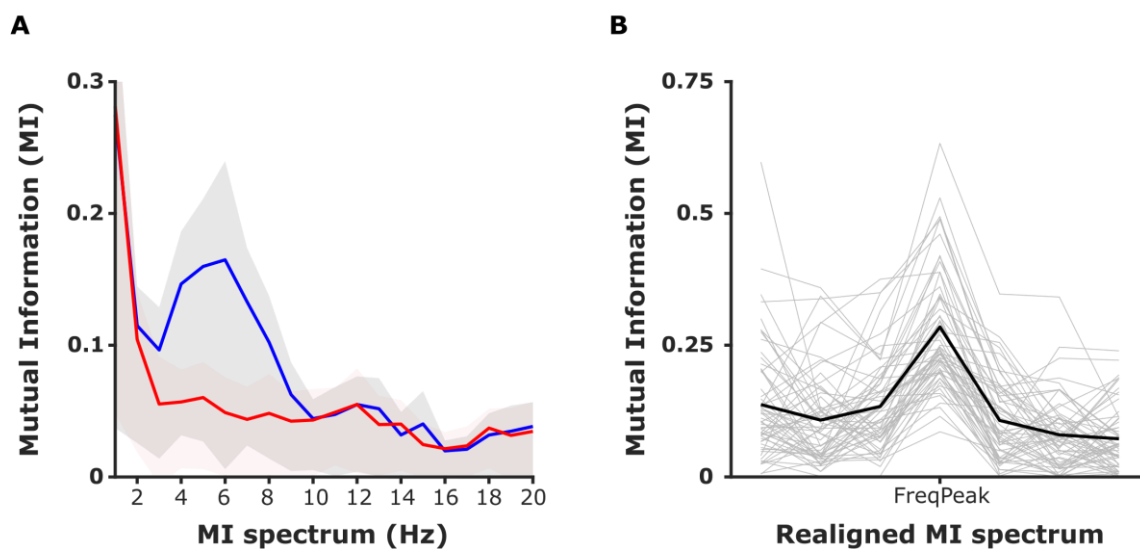
904

905 **Figure S2: Preferred theta phase predicted tone detection task performance.** (A) Tone  
906 detection sensitivity ( $d'$ ) of first and second tones between the group 1 and 2. The graphs  
907 depict the density, the grand average (mean  $\pm$  standard error of the mean; errors bars: 95 %  
908 confidence interval), and individual means for the first/second tones. Results reveal that  
909 participants from the group 2 were significantly better than group 1, only in the late second  
910 time-window when visual cortex supposedly entrained to lip movements' activity.  
911 Significant contrasts are evidenced with stars ( $p < 0.05$ ). (B) Mean phase distributions of the  
912 group 1 (blue) and group 2 (pink). The two separate populations were sorted based on their  
913 individual preferred phase at the second tones (blue and pink dots), where visual  
914 entrainment supposedly took place.

915

916 The ANOVA on  $d'$  scores revealed a significant interaction between tone position and group  
917 ( $F(1, 9) = 5.893$ ;  $p = 0.038$ ). Bonferroni-corrected pairwise t-tests showed that participants  
918 from the group 2 were better than the group 1 to detect the second tones ( $T(1,9) = -0.786$ ;  $p$

919 = 0.028) but not the first tones ( $T(1,9) = -0.208$ ;  $p = 0.629$ ). Results also replicated the effect  
920 of tone position ( $F(1, 9) = 7.715$ ;  $p = 0.021$ ) with greater  $d'$  for the tones sorted as second  
921 than first. Finally, no main effect of group was found ( $F(1, 9) = 2.103$ ;  $p = 0.181$ ). No  
922 difference between threshold ( $SNR_{group1} = 1.39e10^{-3} \pm 2.31e10^{-3}$ ;  $SNR_{group2} = 1.43e10^{-3} \pm$   
923  $3.47e10^{-3}$ ;  $T(1, 24) = -0.66$ ;  $p = 0.95$ ; two-tailed), nor hit rates ( $hit_{group1} = 0.761 \pm 0.003$ ;  
924  $hit_{group2} = 0.763 \pm 0.004$ ;  $T(1, 24) = -0.23$ ;  $p = 0.84$ ; two-tailed) were found in the calibration  
925 task, ruling out any hearing difference.  
926



927  
928 **Figure S3: Mutual information between lips movements and auditory envelope in the**  
929 **movies.** (A) Mean mutual information spectrum ( $\pm$  standard deviation) between the vertical  
930 aperture of the lips and the corresponding (blue line) or random (red line) speech envelope  
931 from movies. The greater dependency between the two signals is reflected by the bump  
932 localised in the theta frequency band of interest (4 - 8 Hz). (B) Realigned spectrum on the  
933 frequency with the greater MI peak ( $\pm$  1-3 Hz) of each movie (grey lines) and averaged (black  
934 line). For each movie, we applied a peak detection and selected the stimuli with a greater  
935 MI between the vertical aperture of the lips and auditory envelope situated in the  
936 frequencies of interest only (i.e. 4, 5, 6, 7 and 8 Hz).