

1 **Hand gestures as visual prosody: BOLD responses to audio-**
2 **visual alignment are modulated by the communicative nature of**
3 **the stimuli**

4
5
6 Emmanuel Biau ^a

7 Luis Moris Fernandez ^a

8 Henning Holle ^c

9 César Avila ^d

10 Salvador Soto-Faraco ^{a, b}

11
12 ^a Multisensory Research Group, Center for Brain and Cognition, Universitat
13 Pompeu Fabra, Barcelona, Spain.

14 ^b Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain.

15 ^c Department of Psychology, University of Hull, UK.

16 ^d Department of Psychology, Universitat Jaume I, Castelló de la Plana, Spain.

17
18 **Manuscript accepted for publication in Neuroimage**

19
20 Corresponding author: Emmanuel Biau

21
22 Dept. de Tecnologies de la Informació i les Comunicacions

23 Universitat Pompeu Fabra

24 Roc Boronat, 138

25 08018 Barcelona

26 Spain

27 +34 691 752 040

28 emmanuel.biau@free.fr

30 **ABSTRACT**

31

32 During public addresses, speakers accompany their discourse with
33 spontaneous hand gestures (beats) that are tightly synchronized with the
34 prosodic contour of the discourse. It has been proposed that speech and beat
35 gestures originate from a common underlying linguistic process whereby both
36 speech prosody and beats serve to emphasize relevant information. We
37 hypothesized that breaking the consistency between beats and prosody by
38 temporal desynchronization, would modulate activity of brain areas sensitive to
39 speech-gesture integration. To this aim, we measured BOLD responses as
40 participants watched a natural discourse where the speaker used beat gestures.
41 In order to identify brain areas specifically involved in processing hand gestures
42 with communicative intention, beat synchrony was evaluated against arbitrary
43 visual cues bearing equivalent rhythmic and spatial properties as the gestures.
44 Our results revealed that left MTG and IFG were specifically sensitive to speech
45 synchronized with beats, compared to the arbitrary vision-speech pairing. Our
46 results suggest that listeners confer beats a function of visual prosody,
47 complementary to the prosodic structure of speech. We conclude that the
48 emphasizing function of beat gestures in speech perception is instantiated
49 through a specialized brain network sensitive to the communicative intent
50 conveyed by a speaker with his/her hands.

51

52 Speech perception; Gestures; Audiovisual speech; Multisensory Integration;
53 MTG; fMRI.

54

55

56

57 1. INTRODUCTION

58

59 In everyday life, most communicative interactions between humans involve
60 auditory and visual information. Indeed, in addition to auditory speech, listeners
61 often have visual access to the speaker's lips, head, body posture and hand
62 gestures. Here we concentrate on the communicative impact of the cospeech
63 gestures that speakers produce with their hand movements while talking to
64 someone (McNeill, 1992). By combining behavioral and physiological measures
65 like event-related potentials (ERPs), prior studies have demonstrated that, for
66 example, gestures describing an object or an action (i.e. iconic gestures) alter
67 semantic processing of the spoken message (Kelly et al., 2004; Kelly et al.,
68 2009; Wu & Coulson, 2010) or help disambiguate semantically complex
69 sentences (Holle et al., 2007). These studies suggest that gestures provide
70 information not present in the verbal modality alone, and support the idea that
71 both streams of information are in fact components of a common integrated
72 language system (McNeill, 1992; Kelly, Creigh & Bartolotti, 2009).

73

74 Many fMRI studies have investigated the degree to which gestures and speech
75 recruit common brain areas. For example, a recent study by Dick et al. (2014)
76 established the implication of a fronto-temporal network of language-related
77 areas when iconic gestures provide complementary information to speech. The
78 Superior Temporal Sulcus (STS) and the Middle and Superior Temporal Gyri
79 (MTG/STG), which are well known to respond to audiovisual (AV) speech (Nath
80 and Beauchamp, 2012; Calvert et al., 2000; Callan et al., 2004; Macaluso et al.,
81 2004; Meyer et al., 2004; Campbell, 2008), have been found to be sensitive to
82 the semantic relationship and congruency between gestures and the spoken
83 message (Marstaller & Burianova, 2014). Greater BOLD responses in the STS,
84 inferior parietal lobule and precentral sulcus were found for the perception of
85 spoken sentences accompanied by semantically corresponding iconic gestures,
86 as compared to meaningless movements or auditory-only versions (Holle et al.,
87 2010; Holle et al., 2008). Willems et al, (2009) also found greater activations in
88 the left STS/MTG when spoken sentences were presented with simultaneous
89 pantomimes (i.e. speech-independent gestures) whose shape matched the verb
90 of the utterance in meaning, as compared to incongruent ones. Additionally, the

91 left Inferior Frontal Gyrus (IFG) has been often found to respond to the
92 manipulation of the semantic relationship between gesture and speech
93 (Marstaller & Burianova, 2014; Willems et al., 2009; Willems et al., 2007),
94 suggesting a role in the integration of both streams of information to support
95 sentence comprehension (Glaser et al., 2013; Uchiyama et al., 2008; Willems et
96 al., 2007; Hagoort, 2005).

97 Although very relevant, these past studies have focused mostly on the
98 neural correlates of hand gestures conveying semantic content, leaving aside
99 other important functions of gestures, like their role as prosodic markers of
100 speech (Guellaï, Langus & Nespors, 2014). Additionally, in these prior studies,
101 participants were typically presented with single sentences where gesture-
102 speech interactions happen in an impoverished context (i.e., short speech
103 fragments containing an isolated gesture corresponding to a critical word). If
104 one considers gestures and speech as two complementary sides of a common
105 underlying language system, a natural continuous flow of visual (gestural) and
106 audio (speech) streams might be essential for the system to remain fully
107 functional (Hubbard et al., 2009; Biau & Soto-Faraco, 2013; Biau et al., 2015).

108
109 In the present study, we address the neural correlates of spontaneous beat
110 gestures. As compared to the more commonly studied iconic gestures, beats
111 are much less sophisticated in semantic content. Generally, beats are rapid
112 biphasic flicks of the hand with no semantic content, serving to highlight
113 relevant information and structure the narrative discourse (McNeill, 1992; So et
114 al., 2012). These kinds of gestures are, by far, the most frequent class of co-
115 speech gesture, and their use is very evident in public addresses, such as
116 political discourses. Based on several evidences, it is now widely hypothesized
117 that beat gestures may also play a role in prosodic processing (Guellaï, Langus
118 & Nespors, 2014). First, beats seem to be very precisely aligned with speech
119 envelope. The functional phase of beats - the moment of maximum extension of
120 the movement, called the “apex” – is temporally aligned with the pitch accent of
121 its affiliate spoken word, increasing its prominence by modulating the acoustic
122 properties of the accentuated syllable (Yasinnik, Renwick & Shattuck-Hufnagel,
123 2004; Krahmer & Swerts, 2007; Treffner and al., 2008; Leonard & Cummins,
124 2010). Second, the speakers use the timing of gesture’s apexes to pack related

125 information together, possibly playing a role in the syntactic organization of
126 sentences supported by prosody (Holle et al., 2012; Guellai, Langus & Nespor,
127 2014). The few studies that have investigated the neural correlates of beat
128 gestures support the prosodic hypothesis too. For instance, Biau & Soto-Faraco
129 (2013) found that beats modulate early ERPs time-locked to the affiliate words
130 onset, within the latency window corresponding to phonological processing.
131 Holle et al. (2012) also found that beats in complex sentences modulated the
132 P600 ERP component, associated to syntactic analysis. Finally, in an fMRI
133 study, observers watched a speaker producing beats while spontaneously
134 speaking (Hubbard et al., 2009). The authors reported greater activations in the
135 left STG/S in response to speech when it was accompanied by beats as
136 compared to unrelated sign language gestures. They also reported greater
137 BOLD responses in the bilateral posterior STG/S, including the Planum
138 Temporale (PT) for speech accompanied by beats compared to a still body.
139 Using beats from an actual fragment of continuous discourse ensured that
140 gestures were produced in a legitimate context and frequency. In addition,
141 spontaneous speech production ensured that the temporal relationship between
142 the continuous beats stream and the rhythm of speech was maintained as in
143 natural language conversation (Biau et al., 2015).

144
145 Scope of the present study

146
147 We hypothesize that beat gestures are produced as an integral part of
148 the language system, providing the listener with visual prosodic information that
149 is aligned with the prosodic contour of the speech message. For this reason, we
150 advance that precise temporal alignment is essential to engage brain processes
151 related to the integration of beats and speech. If this is true, brain activations in
152 relevant integration areas may be sensitive to a breach in the temporal
153 synchrony of beats with respect to their speech affiliates (Marstaller &
154 Burianova, 2014; Hubbard et al., 2009). To test this hypothesis, we used fMRI
155 while participants were presented with video clips in which the video was either
156 synchronized with the audio track or lagged behind 800 milliseconds. With this
157 manipulation, we assumed that when beat's apexes fall out of synchrony with
158 their affiliated speech accentuations, their highlighting function would falter. Yet,

159 please note that desynchronization between beats and speech involves
160 temporal misalignment at many levels, from mere spatio-temporal correlations
161 of low level features to the misalignment in linguistic functions. Therefore, an
162 integral question in this framework is whether the putative prosodic function of
163 beats relates to a generic mechanism of visual emphasis or, alternatively,
164 whether beats engage a specialized mechanism. Revealing such specialization
165 is essential to attribute any beat-speech interaction effects to a common
166 underlying language system. For instance, it is relevant that in the study by
167 Holle et al. (2012), mentioned above, the authors did not find the same effects
168 on the P600 ERP component when speaker's moving hands (producing the
169 beats) were replaced with discs following equivalent spatio-temporal trajectories
170 in the visual display. The authors concluded that beats bear additional
171 communicative intentions above and beyond simple visual emphasis (e.g.
172 intentions and postures that come along with the prosodic variations, which
173 might not be the case for an isolated disc).

174 Following Holle et al.'s logic, we wanted to single out brain areas that
175 play a relevant and specific role in beat-speech integration by looking at the
176 effect of beats-speech (de)synchronization, compared to the same effect when
177 the speaker's hands are replaced by arbitrary visual cues (i.e., moving discs).
178 We hypothesized that the visual emphasis from arbitrary stimuli may differ from
179 the linguistic function that gestures have when combined with speech (i.e. when
180 beat emphasis is synchronized with the speech prosody). If beat gestures
181 effectively confer a special communicative value to the spoken message, then
182 one should expect disparate effects of audio-visual synchrony for beat gestures
183 as compare to visual cues. We set up a 2x2 design with the factors AV
184 synchrony (synchronous or asynchronous) and visual information (beats or
185 discs) to test how the temporal alignment affects the integration of speech with
186 either type of visual information. The interaction between synchrony and visual
187 information is of essential interest because it allows isolating brain areas in
188 which the impact of synchrony depends on which kind of visual information
189 (beats or discs) accompanies audio speech prosody. Please note that a simple
190 comparison between synchronous-asynchronous would conflate brain areas
191 that are sensitive to generic, low level features as well as more specific
192 linguistic related attributes of the stimuli. Thus, in this study we will mainly

193 concentrate on brain areas where such an interaction arises. According to prior
194 literature, these areas might (though not exclusively) correspond to the ones
195 previously shown to be sensitive to gesture-speech integration, such as the left
196 STS/G but also the left IFG (Holle et al., 2007; Willems et al., 2007; Hubbard et
197 al., 2009; Holle et al., 2010; Marstaller & Burianova, 2014).

198

199 **2. MATERIAL AND METHODS**

200

201 2.1 Participants

202

203 Nineteen native speakers of Spanish (12 female, age range 19-29) took part in
204 the current study. All participants were right-handed with normal auditory acuity
205 as well as normal or corrected-to-normal vision. Participants gave informed
206 consent prior to participation in the experiment and the study was approved by
207 the University's ethics committee. Due to a technical problem, two participants
208 could not listen to the speech stream during fMRI data acquisition and were
209 therefore excluded from the statistical analysis. Thus, data from 17 participants
210 (12 females, age range: 22.4 ± 2.4 years old) were included in the imaging
211 analysis.

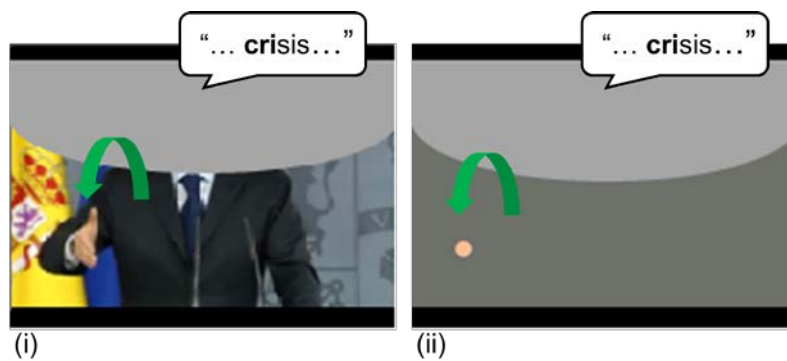
212

213 2.2 Material and stimuli

214

215 We extracted 44 video clips (18 s duration each) from a political discourse of
216 the former Spanish President Luis Rodríguez Zapatero, recorded at the palace
217 of La Moncloa and available on the official website (*Balance de la acción de*
218 *Gobierno en 2010*, 12-30-2010; <http://www.lamoncloa.gob.es>). During the whole
219 public address, the speaker stood behind a lectern, with the upper part of the
220 body in full sight. The video clips were edited using Adobe Premiere Pro CS3.
221 We visually inspected the entire discourse to select relevant segments of
222 speech, containing only beats and cohesive gestures (series of beats that link
223 successive points to a common concept) according to McNeill's definition. Clear
224 iconic gestures were not found but as gesture categories sit along a continuum
225 with fuzzy boundaries, some gestures may fall into multiple categories. Therefore
226 one cannot be absolutely certain that our stimuli never included a minimum of

227 semantic content in the hand shape. However, hand movements always conformed
 228 to McNeill's definition of beat gestures. To avoid abrupt onsets and offsets, we
 229 introduced 1 second audio-visual fade-in and -out at the beginning and end of
 230 each clip (respectively). In all the AV clips, the head of the speaker was masked
 231 with a superimposed ellipse-shaped patch in order to remove any facial
 232 information, such as lips or eyebrow movements, as well as head movements.
 233 After editing, videos were exported using the following parameters: video
 234 resolution 960x720, 25 fps compressor Indeo video 5.10, AVI format; audio
 235 sample rate 48 kHz 16 bits Mono. As explained below, we created four different
 236 versions for each video, corresponding to the four conditions of our
 237 experimental design: Beat Synchronous (Bs), Beat Asynchronous (Ba), Disc
 238 Synchronous (Ds) and Disc Asynchronous (Ds) (Fig. 1).
 239



240
 241 **Figure 1.** Screenshots from (i) Beat and (ii) Disc conditions. Audio and video streams were
 242 either synchronized (Bs and Ds conditions) or desynchronized (audio lagged video by 32
 243 frames, corresponding to 800 ms) with respect to audio in the Ba and Da conditions). Green
 244 arrow illustrates the trajectory of a beat gesture and the corresponding disc. The apex of the
 245 movement coincided in this case with the Spanish word 'crisis'.

246 *Beat conditions:* We selected 44 segments (18s each, 450 frames) of the
 247 discourse in which the speaker naturally produced spontaneous beats (McNeill,
 248 1992). For each clip, the speaker produced a minimum of 8 beats within the 18
 249 s (mean number of gestures per clip: 12.8 ± 4.2). To create the Beat-
 250 Synchronous condition, audio and visual information remained synchronized as
 251 in the original discourse, with the speaker's hands fully visible (beat synchrony,
 252 Bs). For the beat asynchrony (Ba) condition, audio and visual information were
 253 desynchronized by inserting a lag of 800 ms (32 frames), leading to speech
 254 preceding beat gestures.

255

256 *Disc conditions:* To create the disc conditions, the video was removed and the
257 hands were replaced by two discs that followed the hand trajectories of the
258 original clips. We defined the junction between the index and the thumb as the
259 reference point of both hands. We used *Skin Color Estimation Application* and
260 *ELAN* software to detect pixel coordinates of hands frame-by-frame in each
261 Beat video (<http://tla.mpi.nl/tools/tla-tools/elan>; Max Planck Institute for
262 Psycholinguistics, The Language Archive, Nijmegen, The Netherlands;
263 Wittenburg et al., 2006). Reference point coordinates were reviewed and
264 corrected were necessary for both hands using custom-made scripts for Matlab
265 (MATLAB Release 2012b, The MathWorks, Inc., Natick, Massachusetts, United
266 States). The two discs representing the hands had a 40 pixel diameter size and
267 were flesh-colored (Red, Green, Blue color values: 246, 187 and 146) at their
268 corresponding reference point. The background color was set to the average
269 value of a still frame of the speaker (Red Green Blue Value: 110, 114, and 104).
270 We then created a synchronized (Disc Synchrony, Ds) and a desynchronized
271 (Disc Asynchrony, Da) condition following the same process as in the beat
272 condition.

273

274 *Target videos:* To ensure that stimuli were attended, participants performed an
275 auditory detection task. For this, we used two clips from each experimental
276 condition to create 8 targets. For each target video, the fundamental pitch of the
277 original audio tracks was artificially shifted up three semitones (high pitch) for
278 one syllable using Adobe's PitchShift filter while the intensity remained the
279 same. In total, each participant was presented with 36 experimental and 8
280 target videos.

281

282 2.3 Procedure and Instructions

283

284 Participants were presented with 44 trials using E-Prime2 software. The order of
285 trials was pseudo-randomized to avoid direct repetition of experimental
286 conditions. Each trial consisted of a fixation cross with variable duration (from
287 7.5 to 8.5 seconds in steps of 0.25 seconds, uniformly distributed) followed by a
288 video clip. The next trial began automatically after the end of the preceding

289 video. A total of four experimental lists were created, counterbalanced for the
290 four experimental conditions. Each participant saw one of the four lists.

291

292 Participants were instructed to perform an auditory detection task and press a
293 button of the fMRI-compatible controller as soon as they detected an artificial
294 pitch change in the voice of the speaker. The hand holding the controller (left or
295 right hand) was counterbalanced across participants (even though target trials
296 were not included in the statistical analysis). Participants were also instructed to
297 always look at the screen during the whole experiment as if they were watching
298 television. Before the fMRI acquisition, participants performed a rapid training
299 with an extra target video presented in both Bs and Ds conditions as an
300 example of artificial pitch change. After the scanning session, participants were
301 given a questionnaire, asking 1) Did you perceive any asynchrony between
302 video and speech during the experiment? 2) What could the moving discs
303 represent? This questionnaire served to ensure that participants correctly
304 attended to all videos. More importantly, it allowed us to evaluate if they could
305 perceive the asynchrony between video and speech.

306

307 2.4 fMRI acquisition

308

309 Imaging was performed in a single session on a 1.5 T Siemens scanner. We
310 first acquired a high-resolution T1-weighted structural image (GR\IR
311 TR=2200ms, TE=3.79ms, FA=15°, 256 x 256 x 160, 1mm isotropic voxel size).
312 Functional data was acquired in a single run consisting of 610 Gradient Echo
313 EPI functional volumes (TE = 50 ms, TR = 2000 ms) not specifically co-planar
314 with the Anterior Commissure – Posterior Commissure line, acquired in an
315 interleaved ascending order using a 64x 64 acquisition matrix with a FOV =
316 224. Voxel size was 3.5 x 3.5 x 3.5 mm with a 0.6 mm gap between slices,
317 covering 94.3 mm in the Z axis.. The functional volumes were placed attempting
318 to cover the whole brain in 23 axial slices. The first four volumes were discarded
319 to allow for stabilization of longitudinal magnetization.

320

321 2.5 Imaging data analysing

322

323 FMRI data were analyzed using SPM12b (www.fil.ion.ucl.ac.uk/spm) and
324 Matlab R2013b (MathWorks).

325

326 2.5.1. Preprocessing

327

328 Standard spatial preprocessing was performed for all participants using the
329 following steps: Horizontal AC-PC reorientation; realignment and unwarp using
330 the first functional volume as reference, a least squares cost function, a rigid
331 body transformation (6 degrees of freedom) and a 2nd degree B-spline for
332 interpolation, creating in the process the estimated translations and rotations
333 occurred during the acquisition; slice timing correction using the middle slice as
334 reference using SPM8's Fourier phase shift interpolation; coregistration of the
335 structural image to the mean functional image using a normalized mutual
336 information cost function and a rigid body transformation. The image was then
337 normalized into the Montreal Neurological Institute (MNI) space (Voxel size was
338 changed during normalization to isotropic 3.5 × 3.5 × 3.5 mm and interpolation
339 was done using a 4th B-spline degree). Functional data was smoothed using an
340 8-mm full width half-maximum Gaussian kernel to increase signal to noise ratio
341 and reduce inter subject localization variability. To add an extra quality control
342 to the movement in participants, we used the Artifact Detection tools (ART)
343 (http://www.nitrc.org/projects/artifact_detect/) with which the composite
344 movement was calculated. This provides a single measure that comprises the
345 movement due to rotation and translation between volumes. All volumes with a
346 composite movement of more than 0.5 mm or more than 9 standard deviations
347 away from the global mean signal of the session were considered as outliers
348 (On average, 1.4% of the volumes per participant were detected as outliers).
349 One regressor per outlier was added at the first level to discard any possible
350 influence of these volumes in the final analysis.

351

352 2.5.2. fMRI analysis

353

354 The time series for each participant were high-pass filtered at 128 s and pre-
355 whitened by means of an autoregressive model AR(1). At the first level (subject-
356 specific) analysis, box-car regressors modelling the occurrence of the four

357 conditions of interest (Bs, Ba, Ds and Da) and a fifth regressor for trials
358 containing a target, all modelled as 18s blocks, were convolved with the
359 standard SPM12b hemodynamic response function. Additionally, several
360 regressors of no interest were included, including the six movement regressors
361 provided by SPM during the realign process, the extra composite movement
362 regressor calculated with ART and one regressor for each of the volumes
363 considered as outliers. The resulting general linear model produced an image
364 estimating the effect size of the response induced by each of the conditions of
365 interest. The images from the first level were used for the planned critical
366 contrasts in a second level analysis (inter-subject). At the second (inter-subject)
367 level, these images were entered into a random effects factorial design with five
368 levels, corresponding to the four critical conditions, plus an additional subject
369 constant to account for non-condition-specific inter-subject variance. Correction
370 for non-sphericity (Friston et al., 2002) was used to account for possible
371 differences in error variance across conditions and any non-independent error
372 terms for the repeated measures. Statistical images were assessed for cluster-
373 wise significance using a cluster-defining threshold of $p < 0.001$. The 0.05
374 Family-wise error correction critical cluster size was 31 voxels and was
375 determined using random field theory (Data smoothing FWHM: 11.4mm,
376 11.2mm, 11.3 mm. Resel Count: 749.2), considering the whole brain as a
377 volume of interest. Contrasts vectors assessing the two main effects and the
378 interaction were used. Although the whole interaction statistical parametric map
379 is presented, the discussion is limited to the clusters that showed an effect of
380 Beat gestures compared to Discs (Bs+Ba > Ds+Da), as our main interest is
381 focused on the parts of the brain that are involved in beat processing (for
382 unmasked results and additional contrasts, please see supplementary online
383 materials). To achieve this, we masked the interaction contrast, corrected as
384 explained above, with the Beat > Discs contrast (p -threshold (unc.) < 0.05). MNI
385 coordinates were classified as belonging to a particular anatomical region using
386 the SPM Anatomy Toolbox (Eickhoff et al., 2005).

387

388 **3. RESULTS**

389

390 **3.1 Behavioral results**

391

392 Participants correctly detected pitch deviation targets on $65.4\% \pm 31.7\%$ of the
393 target trials and gave False Alarm (FA) responses only on $7.0\% \pm 13.6\%$ of the
394 non-target trials.

395

396 3.2 Post-scanning questionnaire

397

398 When asked, after the scanning session, whether they perceived any
399 asynchrony between video and speech during the experiment, 12 participants
400 responded “yes”; 3 participants responded “yes, but not in the disc condition”
401 and 2 participants responded “no”. With respect to the second question (“What
402 could the moving discs represent?”), all participants responded “the hand of the
403 speaker. This suggests that the asynchrony between beats and speech was
404 noticeable, even though facial information was removed from videos.
405 Furthermore, this consistent response confirmed that the spatiotemporal
406 characteristics of disc movements successfully mimicked the hand trajectories
407 in the Disc conditions. Both the behavioural and post-scanning questionnaire
408 results suggest that participants were attentive to the AV stimuli.

409

410 3.3 fMRI results

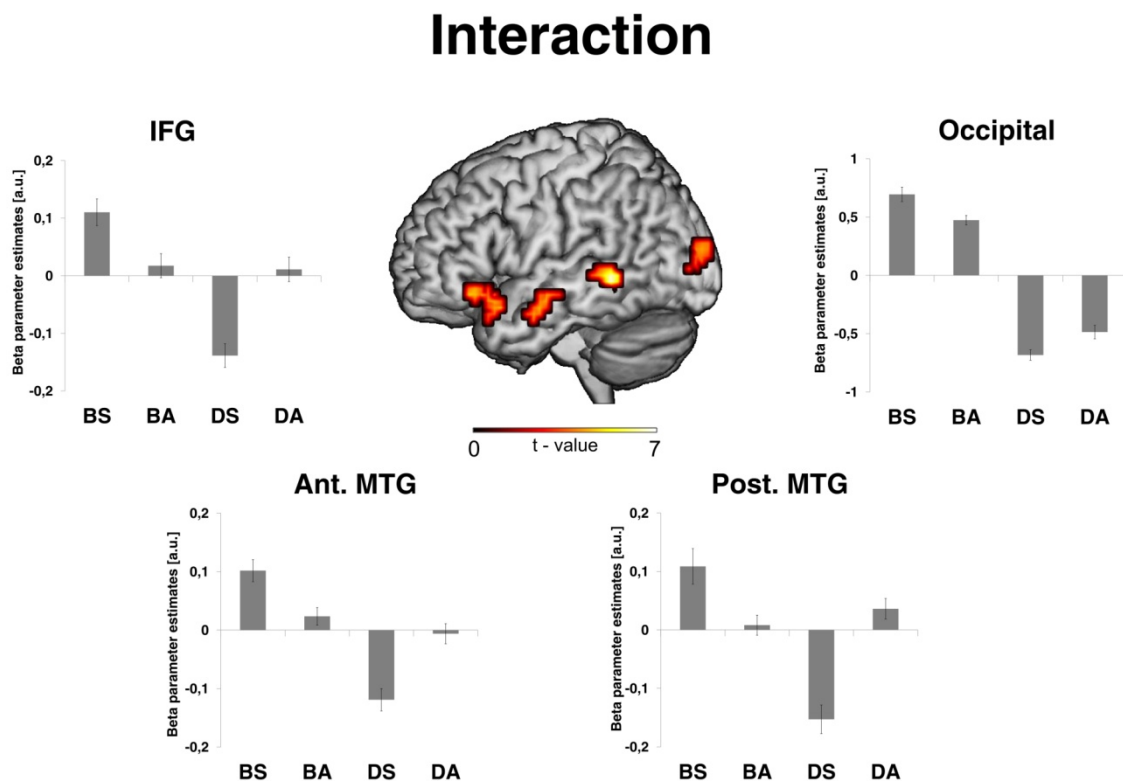
411

412 3.3.1 Differential effect of AV synchrony depending on visual information

413

414 The first contrast of interest concerns the interaction between synchrony and
415 visual information [(Bs-Ba) – (Ds-Da)]. This contrast is of particular interest as it
416 highlights the brain areas where the impact of synchrony depends on which
417 kind of visual information (beats or discs) accompanies speech. We studied this
418 interaction in the areas that showed an effect of Beat > Disc (uncorrected mask
419 $p < 0.05$), as explained in the methods section (see Table 1). This restricts our
420 analysis to areas that are related to beat processing. The results revealed a
421 significant interaction in BOLD responses in two different clusters of the left
422 Middle Temporal Gyrus and Superior Temporal Sulcus (MTG/STS), one more
423 posterior and one more anterior (respectively, pMTG and aMTG/STS).

424 Additionally, significant interactions in left IFG and left occipital cortex
425 (Brodmann area 18) were observed.



426

427

428 **Figure 2.** Interaction contrast [(Bs- Ba) – (Ds – Da)] inclusively masked with the main effect of
429 Beat (Bs+Da) compared to Disc (Ds+Da) using a $p < 0.05$ cluster-corrected threshold with a
430 minimum cluster size $k = 31$ and rendered on a 3D brain surface in MNI space (Left
431 hemisphere). Error bars show 1 S.E.M of parameter estimates. IFG: Inferior frontal gyrus (-41
432 32 -11); Ant.MTG: anterior Middle temporal gyrus (-52 -7 -18); Post. MTG: posterior MTG (-59 -
433 46 -4); Occipital (-20 -95 14).

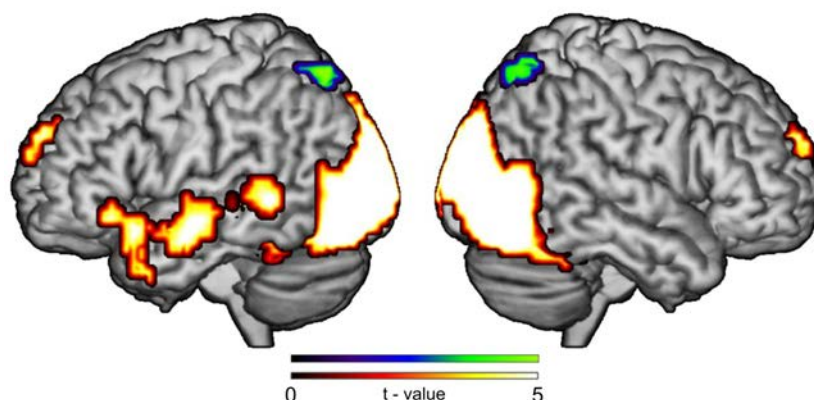
434

435 These results suggest that synchrony differentially affects speech integration,
436 depending on the content of visual information. In particular, speech-gesture
437 synchrony seems to recruit left-hemisphere brain areas preferentially, as
438 compared to other visual cues which share the same spatio temporal properties
439 but are arbitrary. Post-hoc analysis in the four significant clusters revealed that
440 activations were significantly greater when beats and audio were synchronized
441 (Bs) than asynchronous (Ba). Furthermore, the effect of synchrony on brain's
442 activations was exactly the opposite when beats were replaced by simple discs
443 (see Figure 2; see the significance of post-hoc simple main effects in the

444 Supplementary Material). It is worth noting that the areas which display this
445 pattern (MTG, IFG and Occipital cortex in the left hemisphere) and the
446 directionality of the numerical effects of beat synchrony are well in line with
447 previous studies investigating gesture perception (Hubbard et al., 2009; Willems
448 et al., 2009; Skipper et al., 2007; Holle et al., 2008, 2010), which further
449 reassures the interpretation of these activations. Yet, despite this is the pattern
450 expected from prior results and support our hypothesis, one should be careful
451 from putting too much weight on it, given the post-hoc nature of the test.

454 3.3.2 Effect of type of visual information within temporal synchrony

456 Looking at the main effect of type of visual cue within the synchronous
457 conditions can reveal differences arising from the type of visual stimulus. The
458 contrast Beat Synchronous > Disc Synchronous revealed a greater BOLD
459 response in various brain areas when speech was accompanied by
460 synchronized beats (Bs), relative to synchronized discs (Ds) (see figure 3 and
461 table 1). Not surprisingly, the greatest difference was observed in the occipital
462 cortex likely due to a pure difference in visual information between conditions.
463 The contrast also revealed differences in beyond visual brain areas, such as a
464 significantly greater BOLD activity in the left MTG/STS, as well as in the left
465 Inferior frontal Gyrus (left IFG) and left hippocampus. The contrast Ds>Bs
466 revealed greater BOLD activity when speech was accompanied by synchronous
467 discs rather than synchronous hand beats in the Superior Parietal areas
468 bilaterally and right Angular Gyrus (see figure 3 and table 1).



471

472 **Figure 3.** Main effect of Beat Synchronous (Bs) compared to Disc Synchronous (Ds). Statistical
473 maps are thresholded at P -uncorrected <0.001 with a minimum cluster size $k = 31$ and rendered
474 on a 3D brain surface in MNI space. From left to right: left hemisphere, right hemisphere and an
475 axial cut at $z=0$. Hot colors indicate $Bs > Ds$. Cold colors indicate $Ds > Bs$.

476

477 **3.3.3** Effect of synchrony between beat gestures and speech

478

479 The contrasts involving the comparisons $Bs > Ba$ and $Ba > Bs$, restricted within
480 the beat gesture conditions, revealed no main effect of synchrony, when
481 performed at the whole brain level. Note that this particular result deviates from
482 Hubbard et al. (2009), who reported an effect of synchrony in the left STS/G
483 area. However, it must be mentioned that in Hubbard’s study not only the actual
484 synchrony, but also the nature of the gestures themselves was substantially
485 changed between the synchronous and asynchronous condition (beats vs. ASL
486 gestures in the control condition, respectively). In any case, our result implies
487 that despite the BOLD responses for synchronous gestures tend to be larger
488 than the BOLD responses for asynchronous gestures in the areas of significant
489 interaction (as revealed in the interaction analysis). However, as discussed in
490 the introduction, this effect cannot be fully interpreted without factoring in the
491 responses of these areas to the disc synchrony/asynchrony conditions. This is
492 because several low-level generic, as well as language-specific responses to
493 synchrony are conflated in this contrast.

494

495

Hemisphere	Region	Corrected Cluster P-Value	Number of Voxels ^a	Z Score	Coordinates (mm) ^b x y z		
<i>Interaction [(Bs-Ba) – (Ds-Da)] masked with Beat > Disc (mask p-value <0.05)</i>							
L	Middle Temporal Gyrus	0,043	32	5,93	-59	-46	-4
L	Inferior frontal gyrus	0,048	31	4,36	-41	32	-11
L	Temporal Pole			4,35	-45	14	-18
L	Middle Temporal Gyrus	0,048	31	4,20	-52	-7	-18
L	Middle Temporal Gyrus			4,10	-59	-11	-14
L	Middle Temporal Gyrus			4,09	-59	-4	-21
L	Middle Occipital	0,039	33	4,04	-20	-95	14
L	Inferior Occipital			3,38	-31	-88	4

Beat Synchronous > Disc Synchronous

R	Lingual Gyrus	0,000	3080	Inf	8	-88	4
L	Cuneus			Inf	-10	-98	18
L	Calcarine			Inf	-3	-88	-4
L	Middle Temporal Gyrus	0,000	151	5,22	-62	-11	-14
L	Temporal Pole			4,75	-48	18	-14
L	Inferior Frontal Gyrus			4,33	-41	28	-11
L	Thalamus	0,006	52	5,20	-24	-28	0
L	Middle Temporal Gyrus	0,001	75	4,90	-55	-46	0
L	Middle Temporal Gyrus			3,93	-48	-32	0

Disc Synchronous > Beat Synchronous

L	Superior Parietal	0,006	50	4,75	-16	-70	56
R	Superior Parietal	0,009	47	3,73	22	-66	59
	Angular Gyrus			3,49	22	-56	49
	Superior Parietal			3,40	15	-59	63

Beat Synchronous > Beat Asynchronous

No significantly activate regions

Beat Asynchronous > Beat Synchronous

No significantly activate regions

496

497

498 **Table 1.**^a Number of voxels exceeding a voxel-height threshold of $p < 0.001$ using a $p < 0.05$
 499 cluster-extend FWE correction. ^b First three maximum peaks more than 8 mm apart are reported
 500 for each cluster.

501

502

503 4. DISCUSSION

504

505 In the present study, we investigated the neural correlates of spontaneous beat
 506 gestures accompanying continuous, natural spoken discourses. Based on
 507 previous reports (McNeill, 1992; Yasinnik et al., 2004; Guellaï et al., 2014; Biau
 508 et al., [2015](#)), we hypothesized that beats act as a visual counterpart of prosody.
 509 If this is the case, then breaking up the consistency between beat apexes and
 510 speech prosody may affect speech processing. In terms of neural expression,
 511 we hypothesized that if beats are integrated as linguistically relevant
 512 information, brain activity in relevant integration areas may be modulated by an
 513 asynchrony between visual and audio streams. As an integral aspect of this
 514 question, we addressed whether beats convey additional communicative
 515 aspects above and beyond arbitrary visual cues (discs) sharing the same
 516 spatiotemporal properties (Holle et al., 2012). Beats are thought to translate
 517 speaker intentions, extending body posture accompanying speaker's prosody to
 518 emphasize relevant segments of speech, which are available for listeners

519 during speech perception (So et al., 2012; Casasanto & Jasmin, 2009). If this is
520 the case, and beats play a linguistically relevant role above and beyond mere
521 emphasis acting at low-level stages of stimulus processing, then the effect of
522 synchrony for beats should be different as compared to visual discs, in the
523 relevant brain areas. Indeed, this question was answered with the interaction
524 term in our analysis, that indicates that the temporal synchrony of beats with
525 speech prosody has a differential impact on BOLD responses, as compared to
526 other kinds of visual information (here, discs that replaced the speaker's hands).
527 The tendencies in the pattern of the interaction simple contrasts suggest greater
528 activations when beats and speech were presented in synchrony as compared
529 to asynchrony. Instead, the opposite pattern was observed when discs
530 accompanied speech. Based on this significant interaction pattern, we interpret
531 that, in addition to their emphasizing trajectory, beats also convey
532 communicative aspects that simple discs are arguably lacking.

533

534 One surprising finding of our study is that the effect of synchrony for beats (i.e.,
535 greater activity for synchronous as compared to asynchronous beats in left IFG
536 and MTG) was not simply absent for the moving discs, but actually tended to be
537 reversed. When interpreting this cross-over interaction, it is also useful to take
538 into account whether the neural response in these areas represents an
539 activation or deactivation, relative to the implicit fixation cross baseline (see
540 parameter estimates in Fig. 2). Relative to this fixation cross baseline, only
541 speech accompanied by synchronous beats elicited activation in IFG, aMTG
542 and pMTG. This is consistent with the idea that IFG and posterior temporal lobe
543 are crucially involved in comprehending co-speech gestures (Holle et al., 2008,
544 2010, Willems et al., 2007, 2009). In contrast, a visual emphasis cue presented
545 in asynchrony with speech (regardless of whether emphasis consisted of beats
546 or moving discs) did not activate these areas, which may reflect that temporally
547 incongruent AV stimuli are less likely to be integrated and may even cause
548 suppression in multisensory areas (Noesselt et al., 2007). Interestingly,
549 processing speech accompanied by temporally congruent discs elicited a
550 reduction of activity in IFG, aMTG and pMTG, relative to fixation baseline. Such
551 a deactivation could possibly reflect a phasic inhibitory influence onto IFG,
552 aMTG and pMTG whenever speech is accompanied by temporally congruous

553 but unfamiliar visual emphasis cues, such as moving discs. An influence of
554 stimulus familiarity on AV integration in the temporal lobe has been
555 demonstrated before (Hein et al., 2007) and may extend to unfamiliar speech-
556 accompanying visual emphasis cues, such as moving discs.

557

558 Our results are in line with previous fMRI [studies that](#) investigated neural
559 correlates of iconic gestures (Holle et al., 2010; Holle et al., 2008; Willems et al.,
560 2009; Willems et al., 2007). Particularly, one previous fMRI addressed natural
561 hand beats co-occurring with continuous speech (Hubbard et al., 2009) and
562 reported a greater engagement of the STS compared to speech alone, an area
563 comparable to the one found in the present study. The authors also reported
564 greater BOLD activation in the left STS/G when speech was presented with the
565 corresponding beat as compared to when presented with unrelated hand
566 movements. Please note that this comparison does not allow one to infer
567 whether the difference in left STS activation was produced by the lack of
568 synchrony between control gestures and speech, the lack of communicative
569 value of control gestures, or an unknown combination of the two. When
570 Hubbard et al. compared speech accompanying beats to beats presented
571 without speech, no difference was observed, suggesting that the modulations in
572 the left STS/G reflect not only processing of biological movement but also
573 integration of speech with the synchronized beat gestures. Indeed, the STS is
574 sensitive to various types of cross-modal correspondence including AV speech
575 (sound-lip correspondence) in various previous studies (Nath and Beauchamp,
576 2012; Calvert et al., 2000; Callan et al., 2004; Macaluso et al., 2004; Meyer et
577 al., 2004).

578

579 In the present study, the interaction contrast suggests that BOLD response in
580 the left MTG was greater when speech was accompanied by beats as
581 compared to discs (regardless of whether they were synchronized or not with
582 speech). At first glance, the greater response to stimuli containing beats in
583 occipital areas compared to those with discs may reflect a pure bottom-up effect
584 of richness of visual information (Figure 3). However, the interaction (Figure 2)
585 revealed also that the significant difference of BOLD activity in the visual areas
586 between beat and disc were dramatically reduced under asynchronous

587 presentations. This suggests that mere physical differences between beats and
588 discs conditions were not sufficient to explain their respective impact of
589 synchrony in the identified areas. The difference between beats and discs
590 might bring about more profound consequences. For example, in a previous
591 ERP study, Holle et al. (2012) showed that a beat modulated the P600
592 component reflecting syntactic parsing, whereas a disc following the equivalent
593 trajectory did not. The authors suggested that the lack of communicative
594 intention may explain the failure of simple discs to affect the neural correlates of
595 syntactic parsing. Here, the significant simple contrast Bs>Ds supports this
596 claim as it revealed greater activations not only in the occipital areas (although
597 certainly due to differences of visual information, the results are only
598 orientative), but also in the left MTG and left IFG areas. Indirectly, this result
599 also converges toward the idea a differential response to synchrony for using
600 discs that are not functionally associated with speech as part of a common
601 language system.

602

603 According to the effect of interaction on the neural activations, it seems that the
604 MTG responded to some additional language-related aspects associated with
605 beat gestures during speech perception. Previous behavioral studies suggested
606 that some implicit pragmatic and intentional information from the speaker could
607 be extracted from beats, and influence speech encoding. For example, So et al,
608 (2012) showed that adult observers managed to remember more words from a
609 spoken list when the words had previously been accompanied by a beat
610 gesture. As this memory improvement was not found in children, the authors
611 concluded that beat gestures conveyed communicative information but the
612 effect was functionally dependent on experiencing social interactions during
613 development (McNeill, 1992). For example, listeners learn to interpret the
614 speaker's intention to underline relevant information with a beat through social
615 experience. This association of communicative aspects between beats and
616 pitch accentuations was highlighted by Krahmer and Swerts (2007) who
617 showed that listeners perceived words as more salient when accompanied with
618 a beat gesture compared to same words presented in isolation. What is often
619 missing in these studies is whether the value of gestures and their integration of
620 speech simply depended on the general salience of the stimulus, or whether co-

621 speech gestures engaged a more specialized system. Although the listeners in
622 the present study could associate moving discs with movements of the hands
623 and participants were able to detect an asynchrony between discs and speech,
624 synchronized gestures and synchronized discs elicited qualitatively distinct
625 patterns of brain activation (see contrast Bs>Ds). This suggests that during
626 perception listeners distinguished visual information functional related to some
627 aspect of speech (beats) from arbitrary visual cues (discs). Here, this
628 information may require additional processes reflected by the differences of
629 activations in the MTG between beats and discs conditions.

630 In addition to the above explanation, the possible linguistic aspects engaged
631 when beats are present may be directly related to human movement
632 understanding and body postures, over and above to their interaction with
633 speech. The STS was found to respond to point-light representations of
634 biological movements (Grossman et al., 2004; Pelphrey et al., 2004), actions
635 executed by humans (Thioux et al., 2008) and social visual cues (for reviews,
636 see Nummenmaa & Calder, 2009; Allison, Puce & McCarthy, 2000). Herrington
637 et al, (2009) showed that the posterior STS was significantly more activated for
638 trials in which participants perceived human point-light representations of
639 actions compared to non-human movements. In the present study, the discs did
640 not clearly represent a human form but clearly mimicked the trajectories
641 described by hands during speech. In reference to the present study, listeners
642 could have associated discs trajectories with hands (as they identified in the
643 post-task questionnaire). Yet, whatever aspect of biological motion engaged by
644 left MTG activations in the disc conditions, it was more strongly expressed
645 during beat conditions. Please note, however, that this possible perceptual
646 difference between beat gestures and discs in biological motion cannot explain
647 the whole pattern of results we found in the left MTG, because the interaction
648 term $[(Bs - Ba) - (Ds - Da)]$ effectively controls for the different amounts of
649 biological movement in the beat and disc conditions.

650

651 The present results also revealed an interaction between synchrony and visual
652 information effects in the left IFG. Several fMRI studies have showed that the
653 left IFG is sensitive to the semantic relationship between gesture and
654 corresponding speech (Skipper et al., 2007; Willems et al., 2007; Willems et al.,

655 2009; Dick et al., 2009) and may be engaged in the unification of visual
656 (gestures) and audio (speech) complementary streams to facilitate
657 comprehension (Willems et al., 2007; Hagoort, 2005). Recently, a meta-analysis
658 investigating the neural correlates shared between different types of gestures
659 reported a common engagement of the left IFG during the perception of speech
660 accompanied with gestures as compared to a still body (Marstaller & Burianova,
661 2014). However, beat gestures do not convey semantic content, therefore the
662 IFG responses observed in the present study cannot be explained in terms of
663 semantic integration. Beyond meaning integration, the left IFG was also shown
664 to be involved in the process of syntactic analysis during sentence
665 comprehension (Glaser et al., 2013; Meyer et al., 2012; Obleser et al., 2011;
666 Uchiyama et al., 2008). As beats play a role in syntactic parsing (Holle et al.,
667 2012), our results might correspond to an engagement of this area in the
668 integration of beat information toward the parsing of the spoken stream, as
669 compared to moving discs. When beats were delayed (Ba condition), their
670 apexes fell out from synchrony with pitch accents and likely out of the time
671 window of gesture-speech integration, potentially affecting the AV speech
672 processing load (Habets et al., 2011; Obermeier et al., 2011; Obermeier &
673 Gunter, 2014).

674

675 It is worth noting that the simple main effect of synchrony for beat stimuli
676 (contrast Bs vs Ba) in left MTG, IFG and occipital cortex did not reach
677 significance in the whole brain analysis, but it is only revealed by the patterns of
678 activations in the interaction contrasts following up on the interaction. Yet, the
679 post-hoc results obtained for the simple main effects restricted to the interaction
680 areas have to be often interpreted with caution (see Supplementary Materials).

681 In consequence, the interpretation of synchrony effects for beat gestures must
682 be linked to its effects relative to the disc condition. In other words, the disc
683 synchrony manipulation can be seen as a baseline for the beat-synchrony
684 manipulation. However, this is indeed a theoretically relevant type of
685 comparison as discussed Holle et al. (2012). In addition, if we go by the results
686 of previous studies, and extant knowledge the neural correlates of speech, we
687 feel safe in interpreting this pattern in line with the results of the interaction that
688 suggested a difference between synchronous and asynchronous beat

689 conditions (see Figure 2). Note, for example that a similar effect of AV
690 synchrony involving gestures in the left STG/S was reported in Hubbard et al.
691 (2009). In their study, however, as mentioned earlier, Hubbard et al. used
692 unrelated sign language movements as a control condition, which not only
693 constitute a more dramatic asynchrony manipulation altogether (as speech and
694 gestures had completely different rhythms), but also changed the very nature of
695 the visual stimuli from the synchronous to the asynchronous condition. Here, we
696 have looked at these two effects (confounded in Hubbard) separately, and
697 therefore it is not surprising that their individual neural correlates are more
698 subtle. That is, in the present study, although delayed with respect to speech,
699 the rhythm of beats was maintained and might still be associable with the global
700 speech envelope. This may have diminished the detrimental impact of
701 desynchronized gestures on a listener's perception. This may also explain why
702 we did not observe any effect of synchrony in the right auditory cortex related to
703 auditory processing and prosody, as it was reported in Hubbard et al.'s results.
704 A further relevant aspect in our study is that participants were asked to simply
705 focus on an auditory detection task. [This is interesting because our results](#)
706 [cannot be attributed to an explicit monitoring of speech-gesture synchrony. On](#)
707 [the contrary, our auditory detection task may have decreased attention on](#)
708 [visual information and effectively weakened the expression of beat synchrony](#)
709 [on speech processing networks.](#)

710

711 Taken together, the present results provide new insights about the specificity of
712 left MTG and IFG in the processing of multimodal language (for a review, see
713 Campbell, 2008; Özürek, 2014). As participants were not explicitly asked to pay
714 attention to the speaker's hands, this suggests that the temporal
715 correspondence between beats and speech prosody may be picked up
716 automatically. This is in line with previous proposals considering speech and
717 gestures as two side of a same underlying language system (McNeill, 1992;
718 Kelly, Creigh and Bartolotti, 2009). Beats appear to convey additional
719 communicative value such as speakers' intentions, which are not available (or
720 at least, not extracted) from simple visual stimuli (Holle et al., 2012; So et al.,
721 2012; Casasanto & Jasmin, 2009; McNeill, 1992). The access to concurrent
722 gestures during speech perception may engage the listeners and provide a

723 better alignment between listener and speaker, improving speech processing
724 and information encoding. Finally, the fact that the speaker was a well-known
725 former Spanish president may have engaged some political sensitivity from
726 listeners. However, such a possible bias is unlikely to influence our results,
727 since participants viewed the same speaker across all four experimental
728 conditions.

729

730 **5. CONCLUSION**

731

732 We investigated the neural correlates of spontaneous beat gestures
733 produced in continuous speech. Our results revealed that the synchrony
734 affected [brain's](#) activations differently according to the visual information
735 accompanying speech during perception. We concluded that beats [are linguistic](#)
736 [information](#) by their trajectories aligned with speech prosody, but also
737 communicative intentions of the speaker.

738

739 **ACKNOWLEDGMENTS**

740

741 This research was supported by the Ministerio de Economía y Competitividad
742 (PSI2013-42626-P), AGAUR Generalitat de Catalunya (2014SGR856), and the
743 European Research Council (StG-2010 263145).

744

745 **REFERENCES**

746

- 747 Allison, T., Puce, A., & McCarthy, G. (2000). Social perception from visual cues: role of the STS
748 region. *Trends in Cognitive Sciences*, 4(7), 267–278.
- 749 Biau, E., & Soto-Faraco, S. (2013). Beat gestures modulate auditory integration in speech
750 perception. *Brain and Language*, 124(2), 143–52.
- 751 Biau, E., Torralba, M., Fuentemilla, L., de Diego Balaguer, R., & Soto-Faraco, S. (2015).
752 Speaker's hand gestures modulate speech perception through phase resetting of ongoing
753 neural oscillations. *Cortex*, 68, 76-85.
- 754 Brett, M., Anton, J-L., Valabregue, R., & Poline, J-B. Region of interest analysis using an SPM
755 toolbox [abstract] Presented at the 8th International Conference on Functional Mapping of
756 the Human Brain, June 2-6, 2002, Sendai, Japan. Available on CD-ROM in NeuroImage,
757 Vol 16, No 2.

758 Callan, D. E., Jones, J. A., Callan, A. M., & Akahane-Yamada, R. (2004). Phonetic perceptual
759 identification by native- and second-language speakers differentially activates brain
760 regions involved with acoustic phonetic processing and those involved with articulatory-
761 auditory/orosensory internal models. *NeuroImage*, 22(3), 1182–94.

762 Calvert, G. A., Campbell, R., & Brammer, M. J. (2000). Evidence from functional magnetic
763 resonance imaging of crossmodal binding in the human heteromodal cortex. *Current*
764 *Biology: CB*, 10(11), 649–57.

765 Campbell, R. (2008). The processing of audio-visual speech: empirical and neural bases.
766 *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*,
767 363(1493), 1001–10.

768 Casasanto, D., & Jasmin, K. (2010). Good and bad in the hands of politicians: spontaneous
769 gestures during positive and negative speech. *PloS One*, 5(7), e11805.

770 Dick, A. S., Mok, E. H., Raja Beharelle, A., Goldin-Meadow, S., & Small, S. L. (2014). Frontal
771 and temporal contributions to understanding the iconic co-speech gestures that
772 accompany speech. *Human Brain Mapping*, 35(3), 900–17.

773 Dick, A. S., Goldin-Meadow, S., Hasson, U., Skipper, J. I., & Small, S. L. (2009). Co-speech
774 gestures influence neural activity in brain regions associated with processing semantic
775 information. *Human Brain Mapping*, 30(11), 3509–26.

776 [Friston, K. J., Glaser, D. E., Henson, R. N. A., Kiebel, S., Phillips, C., & Ashburner, J. \(2002\).
777 Classical and Bayesian inference in neuroimaging: applications. *NeuroImage*, 16\(2\), 484–
778 512.](#)

779 Glaser, Y. G., Martin, R. C., Van Dyke, J. A., Hamilton, A. C., & Tan, Y. (2013). Neural basis of
780 semantic and syntactic interference in sentence comprehension. *Brain and Language*,
781 126(3), 314–26.

782 Grossman, E. D., Blake, R., & Kim, C.-Y. (2004). Learning to see biological motion: brain activity
783 parallels behavior. *Journal of Cognitive Neuroscience*, 16(9), 1669–79.

784 Guellaï, B., Langus, A., & Nespors, M. (2014). Prosody in the hands of the speaker. *Frontiers in*
785 *Psychology*, 5, 700.

786 Habets, B., Kita, S., Shao, Z., Ozyurek, A., & Hagoort, P. (2011). The role of synchrony and
787 ambiguity in speech-gesture integration during comprehension. *Journal of Cognitive*
788 *Neuroscience*, 23(8), 1845–54.

789 Hagoort, P. (2005). On Broca, brain, and binding: a new framework. *Trends in Cognitive*
790 *Sciences*, 9(9), 416–23.

791 Hein, G., Doehrmann, O., Müller, N. G., Kaiser, J., Muckli, L., & Naumer, M. J. (2007). Object
792 familiarity and semantic congruency modulate responses in cortical audiovisual integration
793 areas. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*,
794 27(30), 7881–7.

795 Herrington, J. D., Nymberg, C., & Schultz, R. T. (2011). Biological motion task performance
796 predicts superior temporal sulcus activity. *Brain and Cognition*, 77(3), 372–81.

797 Holle, H., & Gunter, T. C. (2007). The role of iconic gestures in speech disambiguation: ERP
798 evidence. *Journal of Cognitive Neuroscience*, 19(7), 1175–92.

799 Holle, H., Gunter, T. C., Rueschemeyer, S. A., Hennenlotter, A., & Iacoboni, M. (2008). Neural
800 correlates of the processing of co-speech gestures. *Neuroimage*, 39(4), 2010–2024.

801 Holle, H., Obermeier, C., Schmidt-Kassow, M., Friederici, A. D., Ward, J., & Gunter, T. C.
802 (2012). Gesture facilitates the syntactic analysis of speech. *Frontiers in Psychology*, 3, 74.

803 Holle, H., Obleser, J., Rueschemeyer, S.-A., & Gunter, T. C. (2010). Integration of iconic
804 gestures and speech in left superior temporal areas boosts speech comprehension under
805 adverse listening conditions. *NeuroImage*, 49(1), 875–84.

806 Hubbard, A. L., Wilson, S. M., Callan, D. E., & Dapretto, M. (2009). Giving speech a hand:
807 gesture modulates activity in auditory cortex during speech perception. *Human Brain*
808 *Mapping*, 30(3), 1028–37.

809 Kelly, S. D., Kravitz, C., & Hopkins, M. (2004). Neural correlates of bimodal speech and gesture
810 comprehension. *Brain and Language*, 89(1), 253–60.

811 Kelly, S. D., Ozyürek, A., & Maris, E. (2010). Two sides of the same coin: speech and gesture
812 mutually interact to enhance comprehension. *Psychological Science*, 21(2), 260–7.

813 Kelly, S. D., Ward, S., Creigh, P., & Bartolotti, J. (2007). An intentional stance modulates the
814 integration of gesture and speech during comprehension. *Brain and Language*, 101(3),
815 222–33.

816 Krahmer, E., & Swerts, M. (2007). The effects of visual beats on prosodic prominence: Acoustic
817 analyses, auditory perception and visual perception. *Journal of Memory and Language*,
818 57(3), 396–414.

819 Leonard, T., & Cummins, F. (2011). The temporal relation between beat gestures and speech.
820 *Language and Cognitive Processes*, 26(10), 1457–1471.

821 Macaluso, E., George, N., Dolan, R., Spence, C., & Driver, J. (2004). Spatial and temporal
822 factors during processing of audiovisual speech: a PET study. *NeuroImage*, 21(2), 725–
823 32.

824 Marstaller, L., & Burianová, H. (2014). The multisensory perception of co-speech gestures – A
825 review and meta-analysis of neuroimaging studies. *Journal of Neurolinguistics*, 30, 69–77.

826 Meyer, M., Steinhauer, K., Alter, K., Friederici, A. D., & von Cramon, D. Y. (2004). Brain activity
827 varies with modulation of dynamic pitch variance in sentence melody. *Brain and*
828 *Language*, 89(2), 277–89.

829 Noesselt, T., Rieger, J. W., Schoenfeld, M. A., Kanowski, M., Hinrichs, H., Heinze, H.-J., &
830 Driver, J. (2007). Audiovisual temporal correspondence modulates human multisensory
831 superior temporal sulcus plus primary sensory cortices. *The Journal of Neuroscience: The*
832 *Official Journal of the Society for Neuroscience*, 27(42), 11431–41.

833 Nummenmaa, L., & Calder, A. J. (2009). Neural mechanisms of social attention. *Trends in*
834 *Cognitive Sciences*, 13(3), 135–43.

- 835 Obermeier, C., Holle, H., & Gunter, T. C. (2011). What iconic gesture fragments reveal about
836 gesture-speech integration: when synchrony is lost, memory can help. *Journal of Cognitive*
837 *Neuroscience*, 23(7), 1648–63.
- 838 Obermeier, C., & Gunter, T. C. (2014). Multisensory Integration: The Case of a Time Window of
839 Gesture-Speech Integration. *Journal of Cognitive Neuroscience*, 1–16.
- 840 Obleser, J., Meyer, L., & Friederici, A. D. (2011). Dynamic assignment of neural resources in
841 auditory comprehension of complex sentences. *NeuroImage*, 56(4), 2310–20.
- 842 Pelphrey, K. A., Morris, J. P., & McCarthy, G. (2004). Grasping the intentions of others: the
843 perceived intentionality of an action influences activity in the superior temporal sulcus
844 during social perception. *Journal of Cognitive Neuroscience*, 16(10), 1706–16.
- 845
- 846 Skipper, J. I., Goldin-Meadow, S., Nusbaum, H. C., & Small, S. L. (2007). Speech-associated
847 gestures, Broca's area, and the human mirror system. *Brain and Language*, 101(3), 260–
848 77.
- 849 So, W. C., Sim Chen-Hui, C., & Low Wei-Shan, J. (2012). Mnemonic effect of iconic gesture and
850 beat gesture in adults and children: Is meaning in gesture important for memory recall?
851 *Language and Cognitive Processes*, 27(5), 665–681.
- 852 Thioux, M., Gazzola, V., & Keysers, C. (2008). Action understanding: how, what and why.
853 *Current Biology: CB*, 18(10), R431–4.
- 854 Treffner, P., Peter, M., & Kleidon, M. (2008). Gestures and Phases: The Dynamics of Speech-
855 Hand Communication. *Ecological Psychology*, 20(1), 32–64.
- 856 Uchiyama, Y., Toyoda, H., Honda, M., Yoshida, H., Kochiyama, T., Ebe, K., & Sadato, N.
857 (2008). Functional segregation of the inferior frontal gyrus for syntactic processes: a
858 functional magnetic-resonance imaging study. *Neuroscience Research*, 61(3), 309–18.
- 859 Willems, R. M., Ozyürek, A., & Hagoort, P. (2007). When language meets action: the neural
860 integration of gesture and speech. *Cerebral Cortex (New York, N.Y. : 1991)*, 17(10), 2322–
861 33.
- 862 Willems, R. M., Ozyürek, A., & Hagoort, P. (2009). Differential roles for left inferior frontal and
863 superior temporal cortex in multimodal integration of action and language. *NeuroImage*,
864 47(4), 1992–2004.
- 865 Wu, Y. C., & Coulson, S. (2010). Gestures modulate speech processing early in utterances.
866 *Neuroreport*, 21(7), 522–6.