

## Citation for published version

Esteve-Gibert, N. [Núria], Borràs-Comes, J. [Joan], Asor, E. [Eli], Swerts, M. [Marc], & Prieto, P. [Pilar]. (2017). The timing of head movements: the role of prosodic heads and edges. *Journal of the Acoustical Society of America* 141(6), 4727-4739. doi: 10.1121/1.4986649

### DOI

<https://doi.org/10.1121/1.4986649>

### Handle

<http://hdl.handle.net/10609/150781>

### Document Version

This is the Accepted Manuscript version.

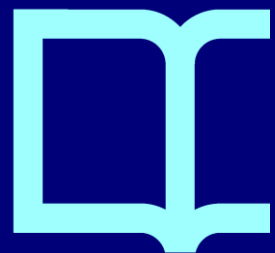
The version published on the UOC's O2 Repository may differ from the final published version.

### Copyright

© 2024 AIP Publishing LLC.

### Enquiries

If you believe this document infringes copyright, please contact the UOC's O2 Repository administrators: [repositori@uoc.edu](mailto:repositori@uoc.edu)



# The timing of head movements: The role of prosodic heads and edges

Núria Esteve-Gibert

Aix Marseille Université, CNRS, Laboratoire Parole et Langage, Aix-en-Provence, France

Joan Borràs-Comes<sup>a)</sup>

Universitat Autònoma de Barcelona, Bellaterra, Spain

Eli Asor

Department of Translation and Language Sciences, Universitat Pompeu Fabra, Barcelona, Spain

Marc Swerts

Department of Communication and Information Sciences, Tilburg University, Tilburg, the Netherlands

Pilar Prieto<sup>b)</sup>

ICREA - Department of Translation and Language Sciences, Universitat Pompeu Fabra, Barcelona, Spain

(Received 23 September 2016; revised 23 March 2017; accepted 4 June 2017; published online xx xx xxxx)

This study examines the influence of the position of prosodic heads (accented syllables) and prosodic edges (prosodic word and intonational phrase boundaries) on the timing of head movements. Gesture movements and prosodic events tend to be temporally aligned in the discourse, the most prominent part of gestures typically being aligned with prosodically prominent syllables in speech. However, little is known about the impact of the position of intonational phrase boundaries on gesture-speech alignment patterns. Twenty-four Catalan speakers produced spontaneous (experiment 1) and semi-spontaneous head gestures with a confirmatory function (experiment 2), along with phrase-final focused words in different prosodic conditions (stress-initial, stress-medial, and stress-final). Results showed (a) that the scope of head movements is the associated focused prosodic word, (b) that the left edge of the focused prosodic word determines where the interval of gesture prominence starts, and (c) that the speech-anchoring site for the gesture peak (or apex) depends both on the location of the accented syllable and the distance to the upcoming intonational phrase boundary. These results demonstrate that prosodic heads and edges have an impact on the timing of head movements, and therefore that prosodic structure plays a central role in the timing of co-speech gestures.

© 2017 Acoustical Society of America. [<http://dx.doi.org/10.1121/1.4986649>]

[BVT]

Pages: 1–13

## I. INTRODUCTION

Studies in the last few decades have shown that co-speech gestures are closely linked to speech in several ways. First, gestures and speech align in terms of semantic and pragmatic meaning (e.g., Bergmann *et al.*, 2014; Kelly *et al.*, 2010; Özyürek *et al.*, 2007). If you tell your friend that you just called your sister, it could well be that you produce a concomitant “calling” gesture in a way that the gesture represents what you also say in speech. Second, gesture and speech co-occur together, they are temporally aligned (e.g., Kendon, 1980; McNeill, 1992). When we speak, the timing of our gestures is not random but is determined by the accompanying speech. In this study, we will examine in detail the temporal alignment patterns between head gestures and speech.

Kendon (1980) and McNeill (1992) stated that the central part of a gesture movement tends to occur within the limits of the prominent prosodic elements of the speech stream. Depending on the gesture and the way it is produced, this prominent part of the gesture can be either an interval, called “gesture stroke,” or a peak in the gesture movement, called “gesture apex.” Many studies have further investigated the specifics of this temporal alignment, revealing that gesture strokes and gesture apexes are aligned with stressed syllables in the speech stream (see Wagner *et al.*, 2014, for a complete review). Interestingly, certain stressed syllables seem to attract more strongly the presence of co-speech gestures: gesture apexes (the peak of prominence in a gesture movement) are more frequently aligned with pitch-accented syllables and with focal pitch accents than with stressed syllables that have a lesser degree of prosodic emphasis (e.g., Alexanderson *et al.*, 2013; De Ruiter, 1998; Ferré, 2014; Yasinnik *et al.*, 2004).

Gesture-speech temporal patterns have been analysed in several contexts, from spontaneous conversations (e.g., Jannedy and Mendoza-Denton, 2005; Loehr, 2012; Yasinnik

<sup>a)</sup>Also at Department of Translation and Language Sciences, Universitat Pompeu Fabra, Barcelona, Spain.

<sup>b)</sup>Electronic mail: pilar.prieto@upf.edu

65 *et al.*, 2004) to controlled laboratory settings (e.g., De  
66 Ruiter, 1998; Esteve-Gibert and Prieto, 2013; Leonard and  
67 Cummins, 2011; Rochet-Capellan *et al.*, 2008; Rusiewicz  
68 *et al.*, 2013). Manual gestures are by far the most-studied  
69 gestures, beat and pointing manual movements traditionally  
70 receiving most of the researchers' attention (e.g., Kendon,  
71 1980; Leonard and Cummins, 2011; Treffner *et al.*, 2008, for  
72 beat gestures; De Ruiter, 1998; Levelt *et al.*, 1985; Rochet-  
73 Capellan *et al.*, 2008; Roustan and Dohen, 2010, for pointing  
74 gestures). Leonard and Cummins (2011) used a motion cap-  
75 tion system to track hand gestures while participants were  
76 reading a short fable. The authors correlated five movement  
77 points (the onset of the movement, the peak velocity of the  
78 extension phase, the point of maximum extension of the  
79 hand before retraction, the peak velocity of the retraction  
80 phase, and the termination of the gesture) with three speech  
81 landmarks (the vowel onset of the stressed syllable in each  
82 word, the estimated P-centre, and the pitch peak within the  
83 stressed syllable). They found that the point of maximum  
84 arm extension (the apex) occurred while the speaker pro-  
85 duced the stressed syllable, and that this pattern was very  
86 stable, meaning that this was the gesture landmark that  
87 showed less variability with respect to its speech anchoring.

88 Yet, another prosodic event might be influencing gesture  
89 timing as well, i.e., intonational phrase boundaries. There is  
90 evidence that the scope of gestural movements typically fin-  
91 ishes at the end of intonational phrases (Loehr, 2012;  
92 Shattuck-Hufnagel *et al.*, 2010; see Krivokapić, 2014, for a  
93 review) and that listeners can automatically extract prosodic  
94 structure by using the temporal scope of manual beat ges-  
95 tures and thus use these gestural features disambiguating the  
96 syntactic structure (Guellai *et al.*, 2014). Interestingly,  
97 phrase boundaries seem to impact not only the ending point  
98 of a gesture movement, but also the timing of the distinct  
99 gesture phases in relation to speech landmarks (De Ruiter,  
100 1998; Esteve-Gibert and Prieto, 2013; Krivokapić *et al.*,  
101 2015; Krivokapić *et al.*, 2016; Levelt *et al.*, 1985). Esteve-  
102 Gibert and Prieto (2013) observed that the movement pattern  
103 of the manual pointing gestures mimicked that of F0 move-  
104 ments. That is, both gesture peaks of pointing gestures and  
105 F0 peaks in rising pitch accents were retracted when the  
106 accented syllable was in phrase-final position; by contrast,  
107 they occurred at the end of the accented syllable when this  
108 syllable was non-phrase-final. Interestingly, Krivokapić *et al.*  
109 (2015) controlled the level of prosodic phrasing (no bound-  
110 ary, prosodic word, intermediate phrase, intonational phrase)  
111 and of prominence (de-accented, broad focus, narrow focus,  
112 contrastive focus) to see how these patterns affected the  
113 alignment of oral and manual pointing gestures with speech.  
114 The authors measured the duration of closing and opening  
115 oral movements and the duration of launching (the distance  
116 between the beginning of the pointing and its apex) and  
117 retraction (the distance between the apex and the end of the  
118 pointing) phases of the pointing gesture. The results showed  
119 that the pattern of manual gestures was very similar to that of  
120 oral gestures: oral movements were longer in trials with  
121 stronger phrase boundaries (just like the launching part of  
122 pointing gestures was), and oral movements were also longer

under prominence (just like the retraction part of the pointing 123  
gestures was). 124

125 Motion caption systems have been used to explore the  
126 timing of head gestures with the aim of creating virtual  
127 agents that can engage in synthesized dialogues that are as  
128 natural as possible. These studies take the position of the  
129 accented syllables as the key prosodic landmark with which  
130 gesture movements align, but they do not take into account  
131 intonational phrase boundaries. In general, they found a sim-  
132 ilar temporal alignment pattern as had been shown for hand  
133 gestures: accented syllables are the anchoring point in  
134 speech for the most prominent part of a head movement, the  
135 gesture apex (defined as the specific point in time when the  
136 head changes its direction in the vertical or lateral move-  
137 ment) (Alexanderson *et al.*, 2013; Ambrazaitis *et al.*, 2015;  
138 Fernández-Baena *et al.*, 2014; Goldenberg *et al.*, 2014; Graf  
139 *et al.*, 2002; Hadar *et al.*, 1983; Ishi *et al.*, 2014; Kim *et al.*,  
140 2014). However, these studies also reported variability in  
141 this alignment pattern. Alexanderson *et al.* (2013), for  
142 instance, analysed 54 head nods that co-occurred with target  
143 words in 20 min of spontaneous conversations, and found  
144 that the head gesture apexes occurred within the accented  
145 syllable, but that there was a great temporal variability in the  
146 precise anchoring point of the gesture apexes within that  
147 syllable. We hypothesize that this variability can be partly  
148 explained by the effects of upcoming intonational phrase  
149 boundaries.

150 The present study aims at investigating the role of the  
151 position of prosodic heads (accented syllables) and prosodic  
152 edges (prosodic word boundaries and intonational phrase  
153 boundaries) on the timing of head nod gestures. To our  
154 knowledge, only three studies have previously alluded at the  
155 combined effect of prosodic heads and edges but without  
156 testing it in a systematic way. Ishi *et al.* (2014) found that, in  
157 Japanese, head nods co-occur with the phrase-final syllables  
158 that are immediately followed by strong intonational phrase  
159 boundaries. Barkhuysen *et al.* (2008) observed that speakers  
160 use the visual information of head movements together with  
161 acoustic cues to mark the ends of utterances. Finally, Hadar  
162 *et al.* (1983) observed that some head gestures were associ-  
163 ated with stress and with junctures (ends of phrases). None  
164 of these previous studies on head nod timing, however, con-  
165 trolled the potential effect of the position of intonational  
166 phrase boundaries on the timing of head nod movements. In  
167 our study, we want to contribute to the previous literature by  
168 adding this factor to our analysis. On the one side, we  
169 hypothesize that accented syllables (prosodic heads) attract  
170 the peak of head movements (the gesture apex). On the other  
171 side, we hypothesize that the role of prosodic edges is crucial  
172 in determining the precise location of the head apex within  
173 the accented syllable. This would imply that speakers plan  
174 the timing of their co-speech gestures by taking into account  
175 the specific characteristics of the prosodic units of speech  
176 they are associating the gesture with, and, importantly, they  
177 take into account both its prominent bits and its ending  
178 edges. If this is the case, our results would help clarifying  
179 the nature of the temporal alignment between head move-  
180 ments and speech events.

To investigate these hypotheses, two experiments were designed. Experiment 1 elicited spontaneous head movements that co-occurred with end-of-utterance target words displaying different stress patterns (stress-initial, stress-final, stress-medial, or monosyllables). This enabled us to test how different positions of the accented syllable and of the phrase boundary influence the timing of head movements. Experiment 2 sought to confirm the findings from experiment 1 in a more controlled way by (a) narrowing down the pragmatic function of head gestures (e.g., a confirmatory function), (b) analysing a balanced number of cases per condition, and (c) varying systematically the position of prosodic heads and edges.

## II. EXPERIMENT 1

Experiment 1 examines the influence of the position of accented syllables and intonational phrase boundaries on the timing of head gestures that co-occur with spontaneous speech.

### A. Method

#### 1. Participants

Thirteen Catalan speakers (1 male and 12 females), between 19 and 24 years of age (mean age 20.9 years) participated in the experiment. All of them were undergraduates at the Universitat Pompeu Fabra in Barcelona, Spain. The participants signed a consent form and received 5 Euro as monetary compensation.

#### 2. Materials

Two digital variants of the Guess Who board game were presented (Ahmad *et al.*, 2011), each containing 24 coloured drawings of human faces. These faces differed regarding various parameters, such as gender or the colour of skin, hair, and eyes. Some faces were bald, some had beards or moustaches, and some were wearing hats, glasses, or earrings. As in the traditional version of Guess Who, the purpose of the game was to try to guess the opponent's mystery person before he or she could guess the participant's own.

The game was designed to naturally elicit sentences containing target words that had different metrical patterns and different distances to upcoming intonational phrase boundaries: stress-initial words (or strong-weak words, hereafter SW) such as *dona* "woman" or *barba* "beard," stress-final words (or weak-strong words, hereafter WS) such as *marrons* "brown" or *barret* "hat," monosyllables (hereafter S) such as *ros* "blond" or *blau* "blue," and stress-medial words (or weak-strong-weak words, hereafter WSW) such as *bigoti* "moustache" or *ulleres* "glasses." These patterns displayed variability in terms of the position of the accented syllable within the prosodic word and also in terms of the distance of the accented syllable from an upcoming intonational phrase boundary. More specifically, while in the WS and S words, the accented syllables were adjacent to the right-edge intonational phrase boundary, in the SW and WSW words, there was one unaccented syllable preceding the upcoming phrase boundary.

Two variants of the game were created, a question-eliciting version (the traditional version of the game) and a statement-eliciting version. In the statement-eliciting version, players produced statements about their own mystery person while the other player listened and eliminated all characters that did not exhibit a particular feature. In the question-eliciting version, players asked questions about the other player's mystery person by asking about specific features of this person. Note that in Catalan statements and yes-no questions have the same word order and they are only distinguished by intonation, rising for questions and falling for statements (unlike in English, for instance, where there is also subject/verb inversion).

All utterances and gestures were spontaneously produced as a result of the natural interaction between players. Crucially for our goals, participants spontaneously produced utterances that had target words in broad focus position and that were immediately followed by an intonational phrase boundary because they were produced at the end of the intonational phrase (see Table I for examples of a dialogue).

### 3. Procedure

While being paired up with another native speaker, all participants played the two versions of the game. The order was counter-balanced across pairs and both versions took place consecutively. During the game, participant A had to request information from participant B in order to find out the mystery person on B's board (question-eliciting version), or had to provide information to participant B so that participant B could guess the mystery person on A's board (statement-eliciting version). Players took turns asking questions or producing statements about the physical features of the "mystery persons." The winner was the player who first guessed the other's mystery person. No specific instructions were given to participants on the type of utterances they had to produce or on specific gestures they could use.

Participants sat facing each other across a table and in front of two laptop computers arranged so that they could not see each other's screen. Participants were audio-visually recorded using two Panasonic HD AVCCAMs at 50 frames per second. The cameras were placed on a tripod at a distance of approximately 1 m from the participants, each one facing a different member of the dyad. The cameras' height

TABLE I. Examples of a dialogue observed in the question-eliciting version of the game (dialogue 1) and in the statement-eliciting version of the game (dialogue 2). Words in bold are target prosodic words produced in broad focus position at the end of the prosodic phrase, and accented syllables are underlined.

Dialogue 1	Dialogue 2
Player A: És una <b>dona</b> ? 'Is it a woman'	Player A: És un <b>home</b> . 'Is it a man'
Player B: Sí. 'Yes'	Player B: D'acord. 'Ok'
Player A: Porta <b>barret</b> ? 'Does she wear a hat?'	Player A: Porta <b>bigoti</b> . 'He has got a moustache'



277 was adjusted to the participants' height in such a way that  
 278 the recording area included the participants' upper body and  
 279 head. Once the participants were seated, the experimenter  
 280 explained the game and gave instructions about the proce-  
 281 dure to be followed for each of the two variations.  
 282 Altogether each version of the game lasted approximately  
 283 20 min.

284 **4. Coding**

285 All utterances about the physical properties of the mys-  
 286 tery person were orthographically annotated and classified as  
 287 being accompanied by a head movement or not. Whenever  
 288 the annotator doubted on this classification, a conservative  
 289 criterion was used, meaning that utterances were coded as  
 290 not being accompanied by a head gesture. The types of head  
 291 movements that were included in the analyses were *head*  
 292 *nods* (following Poggi et al., 2010, a head nod was any verti-  
 293 cal head movement in which the head, after a slight tilt up,  
 294 bends downward and then goes back to its starting point),  
 295 *upward movements* (a head movement directed upward in  
 296 the opposite direction from nodding), and *head tilts* (a head  
 297 inclination or sideward movement) (see Wagner et al., 2014,  
 298 for a complete overview of the head gesture forms). All  
 299 selected sentences had the form of verb + article + noun/  
 300 adjective (the article being optional), as in the statement  
 301 *Porta barret* "(S)he has a hat."

302 From the total amount of utterances produced by partici-  
 303 pants ( $N = 492$ ), 111 utterances (22.6% of the total) were  
 304 spontaneously accompanied by a head gesture. This proportion  
 305 of gesture production per total amount of utterances is consis-  
 306 tent with previous studies (e.g., Alexanderson et al., 2013;  
 307 Ferré, 2014). All head gestures co-occurred with the target  
 308 word in the sentence (i.e., the content word featuring the phys-  
 309 ical property of the character, be it noun or adjective).

310 Table II displays the summary distribution of spontane-  
 311 ously produced utterances across participants, the amount of  
 312 head gestures accompanying the target word, and the stress  
 313 patterns of the target prosodic words. It illustrates that stress-

initial (SW) target words were the most frequently produced, 314  
 followed by monosyllabic words (S), and stress-medial 315  
 words (WSW). The least frequent pattern was the stress-final 316  
 (WS). 317

All utterances that were accompanied by a head gesture 318  
 were further coded in terms of speech and gesture features. 319  
 For gestures, we used ELAN annotation software, a tool that 320  
 allows precise, frame-by-frame navigation through the video 321  
 recording (Lausberg and Sloetjes, 2009). As Fig. 1 illus- 322  
 trates, head nods are characterized by a fall-rise movement 323  
 that is generally preceded by an upward motion (see Ishi 324  
 et al., 2014, for a detailed description of the head nod 325  
 shapes). For the gesture annotation we identified the follow- 326  
 ing three points within the gesture movement: the *onset of* 327  
*the gesture* (the point where the head starts moving from 328  
 its rest position), the *gesture apex* (the point where the 329  
 bi-directional fall-rise head movement changes its direction), 330  
 and the *end of the gesture* (the point where the gesture move- 331  
 ment returns to its rest position). 332

For speech, we manually annotated the beginning and 333  
 endpoints of the entire utterance, of the target prosodic 334  
 word, and of the accented syllable within that target prosodic 335  
 word (see Fig. 2). We used Praat (Boersma and Weenink, 336  
 2012) for speech coding, and Praat annotations were then 337  
 imported into ELAN. The following criteria were used for 338  
 speech segmentation: utterances were pause-bounded mean- 339  
 ingful semantic units; target prosodic words were end-of- 340  
 utterance content words (nouns or adjectives) forming a tone 341  
 group bearing one word stress; and the accented syllable 342  
 within the target prosodic word was the syllable within the 343  
 prosodic word that carried the stress (and consequently the 344  
 pitch accent of the entire utterance). 345

346 **B. Results**

347 For the analyses, the following dependent variables 348  
 were taken into account: (1) the distance in time between the 349  
 beginning of the gesture and the beginning of the prosodic 350  
 word, (2) the distance in time between the end of the gesture

TABLE II. Summary of all the utterances produced, classified as a function of the participant, the presence of a speech-accompanying gesture, and the stress pattern of the target prosodic word.

Participant	Target words without co-speech head gesture				Target words with co-speech head gesture					Total
	WSW	WS	SW	S	WSW	WS	SW	S		
1	14	5	18	10	1	0	6	4	58	
2	12	3	16	10	1	2	7	0	51	
3	15	16	20	17	0	0	1	1	70	
4	11	9	17	10	4	1	9	6	67	
5	11	8	28	10	5	4	13	3	82	
6	2	1	15	3	1	1	4	2	29	
7	3	9	12	6	1	0	2	0	33	
8	4	2	3	1	1	0	0	0	11	
9	1	0	0	1	0	0	1	1	4	
10	0	0	0	1	2	1	6	0	10	
11	3	2	3	1	0	3	6	2	20	
12	1	7	12	5	0	0	5	2	32	
13	0	1	12	5	1	0	1	0	20	
TOTAL	77	63	156	80	17	12	61	21	492	

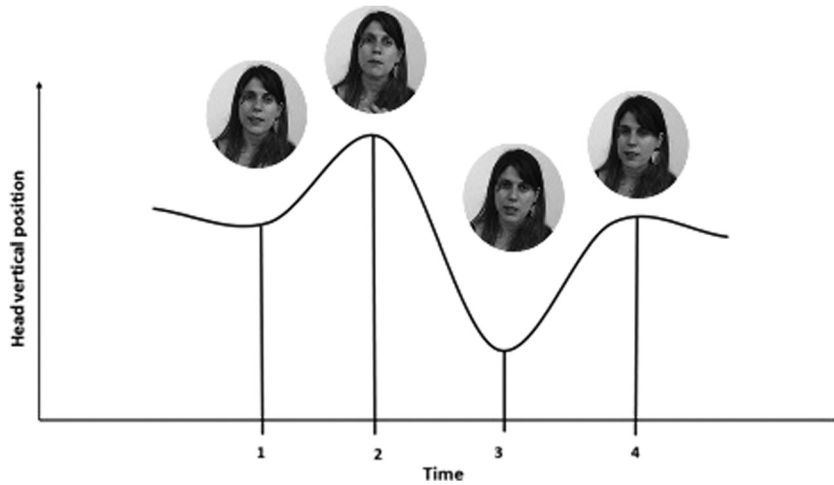


FIG. 1. Schematic representation of the relevant landmarks in a head nod gesture: the beginning of the gesture movement (1), the endpoint of the initial upward motion preceding the falling part of the movement (2), the gesture apex (3), and the end of the gesture (4). The preparation phase of the gesture corresponds to the temporal distance between points 1 and 2, the gesture stroke interval refers to the distance between 2 and 3, and the retraction phase interval is the distance between 3 and 4.

351 and the end of the prosodic word, and (3) the distance in  
 352 time between the gesture apex and the end of the accented  
 353 syllable. In all statistical analyses the fixed factor was the  
 354 metrical pattern of the target prosodic word (4 levels: SW,  
 355 WS, WSW, S), and the random factors were participant and  
 356 item (simple random effects structure). Variables were  
 357 assessed with linear mixed-effects models, using the *lmer*  
 358 function within the *lme4* package in R (Bates *et al.*, 2011).  
 359 The models predicting the first two dependent variables will  
 360 reveal what is the scope of the gesture movement, and  
 361 whether it varies as a function of the position of the accented  
 362 syllable and of the phrase boundary. The model predicting  
 363 the third dependent variable will show if the gesture apex is  
 364 produced within the temporal limits of the accented syllable,  
 365 and whether the position of the intonational phrase boundary  
 366 influences the precise location of the apex within this  
 367 accented syllable.

368 Table III summarizes the results of the mixed-effects  
 369 models. Results showed that the stress pattern of the pro-  
 370 sodic word did not influence the distance between the ges-  
 371 ture start and the start of the prosodic word or the distance  
 372 between the gesture end and the end of the prosodic word.  
 373 This means that, independently of the position of the pro-  
 374 sodic prominence and of the upcoming phrase boundary,  
 375 head movements started several milliseconds before the pro-  
 376 sodic word started, and ended several milliseconds after the  
 377 prosodic word ended (for descriptive values of all the

analyses, see the Appendix). Instead, the stress patterns sig- 378  
 nificantly impacted the temporal distance between the ges- 379  
 ture apex and the end of the accented syllable, in that the 380  
 stress-final patterns (S and WS) differed significantly from 381  
 non-final stress patterns (SW and WSW). As Fig. 3 shows, 382  
 the apex was aligned towards the middle of the accented syl- 383  
 lable when there was non-accented material preceding the 384  
 right-edge phrase boundary (SW and WSW), while it was 385  
 much more retracted when the end of the accented syllable 386  
 coincided with the presence of a right-edge phrase boundary 387  
 (S and WS). 388

389 Three additional linear mixed-effects analyses with the  
 390 same dependent variables and random factors were con-  
 391 ducted, but now with sentence type as fixed factor (2 levels:  
 392 question, statement). They revealed that the alignment pat-  
 393 terns did not vary significantly as a function of this param-  
 394 eter (temporal distance between word start and gesture start:  
 395  $\beta = 0.09$ ,  $t = 1.33$ ; temporal distance between word end and  
 396 gesture end:  $\beta = 0.02$ ,  $t = 0.14$ ; temporal distance between  
 397 apex and end of accented syllable:  $\beta = 0.07$ ,  $t = 1.04$ ).

**C. Discussion** 398

399 In experiment 1 participants took part in two variants of  
 400 the Guess Who game (one designed to elicit questions and  
 401 the other to elicit statements), while being audio-visually  
 402 recorded. Our aim was to see how speakers temporally

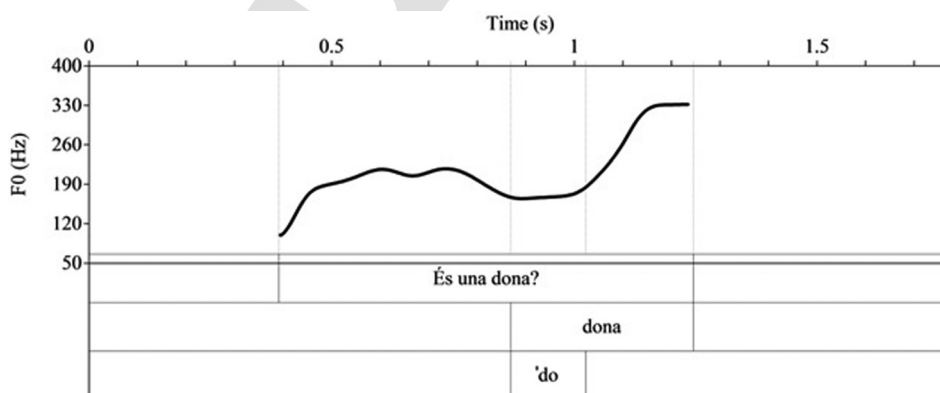


FIG. 2. Speech annotation of the utterances accompanied by a head gesture in Praat. First tier, temporal limits of the entire utterance. Second tier, temporal limits of the target prosodic word. Third tier, temporal limits of the accented syllable within that prosodic word.

TABLE III. Summary of the liner mixed-effects analyses for each dependent variable in experiment 1. Significant comparisons are in bold (we considered statistical significance to be  $p \leq 0.05$ ).

	$\beta$	SE	$t$
Gesture onset / word onset			
S vs WS	0.091	0.119	0.761
S vs SW	0.059	0.087	0.682
S vs WSW	0.099	0.113	0.881
WS vs SW	-0.031	0.104	-0.307
WS vs WSW	0.008	0.126	0.067
SW vs WSW	0.050	0.096	0.420
Gesture end / word end			
S vs WS	-0.039	0.183	-0.216
S vs SW	-0.194	0.133	-1.460
S vs WSW	-0.092	0.172	-0.535
WS vs SW	-0.154	0.157	-0.983
WS vs WSW	-0.052	0.192	-0.275
SW vs WSW	0.102	0.145	0.700
Gesture apex / end accented syllable			
S vs WS	-0.106	0.117	-0.905
S vs SW	0.257	0.085	<b>3.023</b>
S vs WSW	0.248	0.110	<b>2.245</b>
WS vs SW	0.363	0.101	<b>3.608</b>
WS vs WSW	0.354	0.123	<b>2.882</b>
SW vs WSW	-0.009	0.093	-0.102

403 aligned the head movements with speech while spontane-  
 404 ously interacting with an interlocutor. Specifically, we were  
 405 interested in the influence of the prosodic heads (accented  
 406 syllables) and phrase boundaries on the timing of head  
 407 gestures.

408 The first main result was that speakers spontaneously  
 409 produced head gestures together with the target prosodic

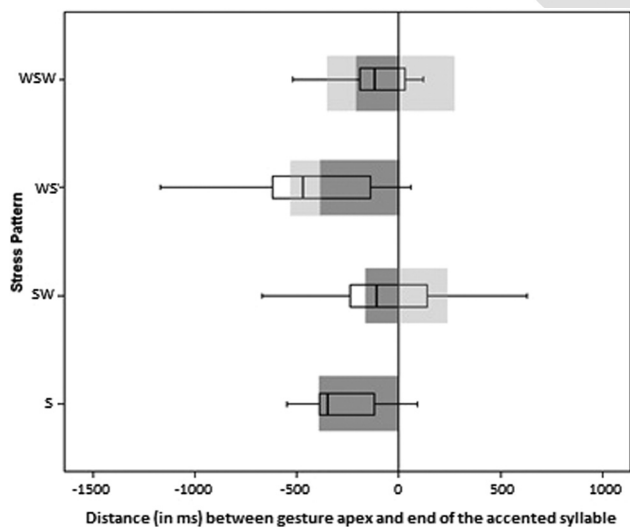


FIG. 3. Box plots displaying the temporal distance (in ms) between the gesture apex and the end of the accented syllable. The 0 represents the end of the accented syllable. Negative values show cases where the apex occurred before the end of the accented syllable. The dark grey shadow on top of box plots indicates the temporal limits of the accented syllable (means values) and the light grey shadows indicate the temporal limits of the non-accented syllables within the prosodic word (means values).

word. Participants were neither instructed regarding the type  
 of sentences to be produced and were not explicitly told to  
 gesture. Yet, all utterances included a phrase-final target  
 word in broad focus position, and almost one fourth of the  
 phrase-final target words were accompanied by a head ges-  
 ture (head nod, head tilt, or upward movement). Despite the  
 inter-individual variability in gestures production (also  
 observed in Graf *et al.*, 2002; Ishi *et al.*, 2014; Swerts and  
 Krahmer, 2010), the ratio of head gesture per utterance is  
 similar to what previous studies have found when examining  
 spontaneous interactions (Alexanderson *et al.*, 2013; Ferré,  
 2014) and indicates that the procedure was useful for the  
 purposes of our study. Spontaneous data are valuable  
 because they reveal the patterns of real-world interactions,  
 but at the same time they complicate the examination of  
 whether this variability is the result of different speaking  
 styles or maybe of different pragmatic functions served by  
 the head gesture (see experiment 2, and also the end of this  
 section for a discussion of this issue).

The second main result was that the scope of the head  
 gestures was the focused prosodic word. Irrespectively of the  
 position of the prosodic prominence within the prosodic  
 word, head gestures start close to the beginning of the corre-  
 sponding prosodic word and they end after prosodic words  
 are finished. This result contradicts those observed by Kim  
*et al.* (2014), who found that head movements occurred dur-  
 ing the critical focused word in narrow-focus conditions but  
 they occurred everywhere in broad-focus conditions. Yet, it  
 goes in line with previous studies on gesture-speech align-  
 ment, which observed that the onset and offset of gesture  
 movements are aligned with the onset and offsets of affilia-  
 ted target words (e.g., Butterworth and Beattie, 1978;  
 Kendon, 1980; Nobe, 2000; Roustan and Dohen, 2010;  
 Schegloff, 1984).

The third main result, and in our view the most interest-  
 ing one, refers to the temporal alignment of the gesture apex  
 with the accented syllable. We found that the position of the  
 head apex (the peak of gesture prominence) was influenced  
 by the position of the accented syllable and of the upcoming  
 phrase boundary. First, gesture apexes were produced within  
 the temporal limits of the accented syllable (except for the  
 WS case, in which the apex occurred during the pre-  
 accented interval). Second, the exact anchoring point of the  
 apex within the accented syllable depended on the position  
 of the upcoming phrase boundary: the gesture apex was  
 retracted if the prosodic word had the stress in phrase-final  
 position (as in S and WS, possibly due to the prosodic pres-  
 sure exerted by the upcoming prosodic boundary), and it was  
 lagged if the prosodic word did not have the stress in phrase-  
 final position (as in SW and WSW, where there is enough  
 post-accentual material where the retraction of the head  
 movement can be accommodated). The case of the phrase-  
 final WS stress pattern is interesting because the apex is so  
 retracted that it is produced out of the temporal limits of the  
 accented syllable, suggesting that the position of the upcom-  
 ing intonational phrase boundary has a stronger impact than  
 the position of the accented syllable.

In sum, results from experiment 1 reveal that focused  
 prosodic words determine the scope of head movements,



469 that accented syllables seem to attract the peak of the gesture  
470 movement, and that phrase boundaries seem to determine  
471 the position of the peak within the accented syllable. The  
472 results of the WS patterns might also suggest that the effect  
473 phrase boundary might be stronger than that of the accented  
474 syllable. Thus, the prosodic structure of the utterance seems  
475 to have a strong impact on the timing of the apexes of  
476 speech-accompanying head gestures. This effect is consis-  
477 tent with previous results on the alignment of pitch peaks in  
478 rise-fall intonation contours (Prieto and Ortega-Llebaria,  
479 2009), and of gesture peaks in manual pointing gestures (De  
480 Ruiter, 1998; Esteve-Gibert and Prieto, 2013).

481 However, some caveats in this experiment prevent us  
482 from drawing strong conclusions, mostly as a consequence  
483 of the spontaneous nature of the corpus. First, the spontane-  
484 ous corpus yielded an unbalanced number of cases within  
485 each stress pattern condition. The results for the SW pattern,  
486 for instance, were based on a substantial number of cases,  
487 but the other patterns were three to five times less frequent.  
488 Second, although we controlled for sentence type (yes-no  
489 question versus statement), the spontaneous elicitation pro-  
490 cedure did not allow us to finely control for the speakers'  
491 pragmatic intent. Previous studies have found that head nods  
492 can have different communicative functions: inclusivity,  
493 intensification, uncertainty, agreement, approval or emphasis  
494 (McClave, 2000; Poggi *et al.*, 2011; Poggi *et al.*, 2010). The  
495 emphatic function of head nods has also been observed in  
496 perception studies. It has been found that eyebrow move-  
497 ments and head nods help listeners to perceive prominent  
498 events in speech (House *et al.*, 2001; Kraemer and Swerts,  
499 2007) and facilitate the recognition of prosodic contrastive  
500 focus (Dohen and Loevenbruck, 2004; Prieto *et al.*, 2015). It  
501 has been proposed that the temporal patterns of the gesture-  
502 speech integration can be influenced by semantic and prag-  
503 matic reasons (e.g., Bergmann *et al.*, 2011; Esteve-Gibert  
504 *et al.*, 2014). It could well be that the participants in our  
505 game responded with different degrees of commitment to the  
506 proposition and with different pragmatic intentions in mind.  
507 Maybe in experiment 1 the speaker's pragmatic intention  
508 had influenced the temporal alignment of the gesture-speech  
509 landmarks. Third, we do not know if the "attraction effect"  
510 of the accented syllable over the gesture apex is still main-  
511 tained when there are larger distances between the accented  
512 syllable and the upcoming phrase boundary. It could be that  
513 this effect is reduced, maybe leading to gesture apexes that  
514 occur during the post-accented material. Experiment 2 was  
515 designed to remedy these concerns.

### 516 III. EXPERIMENT 2

517 The purpose of experiment 2 was to find additional sup-  
518 port for the findings obtained in experiment 1. We designed  
519 a more controlled setting that would allow us to elicit head  
520 nod gestures with a co-referential meaning of confirmation,  
521 accompanying target words with specific stress patterns, and  
522 a balanced number of cases per stress pattern. Furthermore,  
523 an additional measure was taken into account in order to dis-  
524 entangle whether phrase boundaries have a stronger impact  
525 than accented syllables in determining the alignment of head

gesture apexes with speech: the temporal distance between 526  
the beginning of the gesture stroke and the beginning of the 527  
accented syllable. This new measure will show us if the posi- 528  
tion of the prominent gesture interval (the gesture stroke) is 529  
determined by the position of the prosodic head (the 530  
accented syllable), by the upcoming phrase boundary, or by 531  
the entire prosodic word. Finally, in order to test whether the 532  
"attraction effect" of prosodic heads over gesture apexes 533  
is maintained when these heads are more distant to prosodic 534  
edges, a new stress pattern condition was included in the 535  
analyses (namely strong-weak-weak words, hereafter SWW). 536

## A. Method 537

### 1. Participants 538

Eleven Catalan speakers (4 male, 7 female), between 22 539  
and 54 years of age (mean age 30.5 years) participated in 540  
this experiment. All of them were students or staff at the 541  
Universitat Pompeu Fabra in Barcelona. They participated 542  
voluntarily and were not aware of the purpose of the experi- 543  
ment. None of them had participated in experiment 1. 544

### 2. Materials 545

Speakers were asked to participate in a Discourse 546  
Completion Task (DCT; Billmyer and Varghese, 2000; 547  
Blum-Kulka *et al.*, 1989) involving a set of 25 discourse con- 548  
texts. A set of 25 cards was created, each containing a situa- 549  
tion in which a hypothetical interlocutor is not sure whether a 550  
certain city (whose name appeared on the card) is the capital 551  
of a foreign country, a Spanish autonomous community, or a 552  
particular district in Catalonia. We chose to use names of 553  
world capital cities (and cities in Catalonia that would be 554  
well-known to all participants) as target words so that the sit- 555  
uations described in the DCTs would be as close as possible 556  
to natural conversational situations. 557

Example (1) shows an example of a DCT. In this 558  
instance the target word is *Roma* "Rome," as indicated by 559  
the boldface. 560

(1) *Esteu jugant al Trivial i tu i en Joan sou part del mateix 561*  
*equip. Surt una fitxa que demana la capital d'Itàlia. En 562*  
*Joan en aquell moment dubta si la capital d'Itàlia és 563*  
*Roma i t'ho diu dubtant. Tu li dius que és cert, que és 564*  
*Roma, la capital d'Itàlia. 565*

Expected answer: *Sí, sí, la capital d'Itàlia és **Roma**. 566*

"You and Joan are playing Trivial Pursuits and you are 567  
on the same team. The card you get asks you to name the 568  
capital of Italy. Joan is unsure and asks you whether it is 569  
Rome or not. You tell him that yes, it is Rome." 570

Expected answer: "Yes, yes, the capital of Italy is 571  
**Rome**." 572

All of the discourse contexts used for the DCT task 573  
were designed to elicit a declarative sentence expressing 574  
confirmation. The target words had one of five different 575  
stress patterns, as described in Table IV. There were five tar- 576  
get words for each pattern and they were expected to occur 577  
at the end of prosodic phrases. Each metrical pattern was 578  
chosen to represent a different position of prosodic promi- 579  
nence and prosodic edges, with stressed syllables in word 580



TABLE IV. The different stress patterns of the Catalan target words controlled for in experiment 2. In the examples column, stressed syllables are underlined.

Stress patterns of the target word	Position of the prosodic prominence	Examples
S	initial and final	<u>Vic</u> , <u>Valls</u>
WS	final	<u>París</u> , <u>Dakar</u>
SW	initial	<u>Roma</u> , <u>Lima</u>
SWW	initial	<u>Mònaco</u> , <u>Washington</u>
WSW	medial	<u>Figueres</u> , <u>Caracas</u>

581 initial, medial, or final position, and with unaccented syllables preceding, following, or surrounding the accented syllable.  
582  
583

### 584 3. Procedure

585 Participants were presented with one card at a time in random order, and were asked to read it carefully, to imagine themselves in the situation described in the discourse context, and, finally, to provide an appropriate verbal response. When 588 participants provided a response that did not include the target word (e.g., *Sí, sí, és veritat* “Yes, yes, that’s right”), the experimenter asked them to provide another response using the name of the capital city within the sentence. In order to elicit head nods as spontaneously as possible, participants were asked to produce spontaneous responses and were never prompted to gesture or produce utterances in an “expressive” manner.

596 Participants were audio-visually recorded using a Panasonic HD AVCCAM at 50 frames per second. The camcorder was placed on a tripod at a distance of approximately 1 m from the participant, and its height was adjusted to the participant’s height in such a way that the recording area included the participant’s upper body and head. The participants were recorded while standing up and were asked not to hold the DCT cards while providing a response. The entire procedure lasted approximately 15 min. A total of 275 trials (11 participants  $\times$  5 stress patterns  $\times$  5 items per pattern) were elicited.

### 607 4. Coding

608 We selected all utterances that were produced with a head nod gesture accompanying the target prosodic word, which occurred in focus position and was immediately followed by a prosodic boundary. The criterion for including head nods was the same as in experiment 1. From the total amount of trials ( $N=275$ ), 155 trials (56.4% of the total) were produced with a confirmation head nod gesture accompanying the target prosodic word. The remaining 120 trials were excluded from our analysis because speakers did not produce any head nod ( $N=48$ ), produced repetitive head nods associated with the adverb(s) *sí* “yes” and that continued during the entire utterance (called “hybrid” gestures in Yasinnik *et al.*, 2004) ( $N=39$ ), the target word was mispronounced ( $N=3$ ), or due to experimenter error ( $N=3$ ). We also excluded instances of head nods that co-occurred with the copular verb *és* “is” instead of with the target prosodic word ( $N=27$ ). Although these latter cases were

pragmatically appropriate in the context of the task, they would have been included in the group of head nods accompanying monosyllabic S words and thus they would have unbalanced the number of trials per stress pattern.

Responses analyzed in this study had one of the following two structures: in 96.2% of the trials ( $N=149$ ) the target name was produced in the main clause at the end of the prosodic phrase (e.g., *Sí, sí, la capital de França és París*. “Yes, yes, the capital of France is Paris”) and in 3.8% of the trials ( $N=6$ ) the target name appeared in a left-dislocated position, also at the end of the prosodic phrase (e.g., *Sí, sí, és París, la capital de França*. “Yes, yes, it is Paris, the capital of France”).

All 155 valid trials were annotated in terms of speech and gesture. The speech annotation was the same as in experiment 1. The gesture annotation was very similar to experiment 1 except with the addition of an extra temporal landmark: the onset of the gesture stroke (point 2 in Fig. 1). As a result, four points within the head movement were identified in experiment 2: the onset of the gesture (the point at which the head starts moving from its rest position, the onset of the gesture stroke (the start of the falling part of the head movement), the gesture apex (the point in which directions change), and the end of the gesture (the point in which the gesture movement returns to its rest position).

## B. Results

The following dependent variables were assessed using linear mixed-effects models (*lmer* function of the *lme4* package in R, Bates *et al.*, 2011): (1) the start of the head movement with respect to the start of the target prosodic word, (2) the end of the head movement with respect to the end of that prosodic word, (3) the start of the gesture stroke with respect to the start of the accented syllable, and (4) the position of the gesture apex with respect to the end of the accented syllable. The fixed factor in all the analyses was the metrical pattern of the target prosodic word (five levels: S, SW, SWW, WS, and WSW), and random factors were participant and item (simple random effects structure).

Table V summarizes the results of the analyses and Fig. 4 illustrates these results in a visually succinct way. First, results revealed that the gesture started before the onset of the target word, and that the temporal distance between the two landmarks was the same across conditions. Only the stress-medial WSW pattern differed: compared to the other patterns, the gesture start was slightly closer to the word start (for descriptive values of all the analyses, see the Appendix). All target words in the elicited sentences were preceded by the copular verb *és* “is,” hence gesture events that preceded the target prosodic word occurred during this preceding speech material.

Second, the temporal distance between the beginning of the gesture stroke and the beginning of the accented syllable varied significantly depending on whether there was pre-accented material within the prosodic word, as it occurred closer to the beginning of the accented syllable in stress-initial words (S, SW, and SWW) and further from it in stress-final and stress-medial patterns. Figure 5 illustrates

TABLE V. Summary of the linear mixed-effects analyses for each dependent variable in experiment 2. Significant comparisons are in bold (we considered statistical significance to be  $p \leq 0.05$ ).

	$\beta$	SE	$t$
Gesture onset / word onset			
S vs SW	10.01	29.00	0.345
S vs SWW	-11.63	28.58	-0.407
S vs WS	-10.68	30.27	-0.353
S vs WSW	69.70	29.94	<b>2.328</b>
SW vs SWW	-21.64	27.84	-0.777
SW vs WS	-20.69	29.74	-0.696
SW vs WSW	59.69	29.24	<b>2.041</b>
SWW vs WS	0.94	29.28	0.032
SWW vs WSW	81.32	28.80	<b>2.823</b>
WS vs WSW	80.373	30.56	<b>2.630</b>
Gesture end / word end			
S vs SW	-19.852	29.768	-0.667
S vs SWW	-88.267	29.295	<b>-3.013</b>
S vs WS	-8.208	31.106	-0.264
S vs WSW	-87.092	30.696	<b>-2.837</b>
SW vs SWW	-68.42	28.50	<b>-2.400</b>
SW vs WS	11.64	30.66	0.380
SW vs WSW	-67.24	30.00	<b>-2.241</b>
SWW vs WS	80.069	30.163	<b>2.654</b>
SWW vs WSW	1.175	29.487	0.040
WS vs WSW	-78.884	31.442	<b>-2.509</b>
Stroke onset / onset accented syllable			
S vs SW	-1.517	21.614	-0.070
S vs SWW	1.326	21.287	0.062
S vs WS	-102.790	22.580	<b>-4.552</b>
S vs WSW	-47.114	22.306	<b>-2.112</b>
SW vs SWW	2.843	20.721	0.137
SW vs WS	-101.272	22.226	<b>-4.556</b>
SW vs WSW	-45.597	21.785	<b>-2.093</b>
SWW vs WS	-104.116	21.875	<b>-4.760</b>
SWW vs WSW	-48.440	21.440	<b>-2.259</b>
WS vs WSW	55.68	22.81	<b>2.440</b>
Gesture apex / end accented syllable			
S vs SW	280.63	20.87	<b>13.449</b>
S vs SWW	285.94	20.53	<b>13.925</b>
S vs WS	-10.15	21.80	-0.465
S vs WSW	235.15	21.52	<b>10.929</b>
SW vs SWW	5.309	19.978	0.266
SW vs WS	-290.779	21.493	<b>-13.529</b>
SW vs WSW	-45.485	21.029	<b>-2.163</b>
SWW vs WS	-296.088	21.145	<b>-14.003</b>
SWW vs WSW	-50.794	20.668	<b>-2.458</b>
WS vs WSW	245.29	22.04	<b>11.129</b>

682 that this distance varied as a function of whether the onset of  
 683 the prosodic word coincided with the accented syllable or  
 684 not, since speakers always aligned the gesture stroke some  
 685 milliseconds before the onset of the prosodic word.

686 Third, regarding the temporal distance between the end  
 687 of the gesture and the end of the prosodic word, we found  
 688 that the gesture end was aligned significantly differently in  
 689 the trisyllabic words (SWW and WSW) compared to the  
 690 other patterns (S, WS, and SW): in trisyllabic words the ges-  
 691 ture end occurred a little before the end of the prosodic  
 692 word, while in the other patterns it occurred closer to it.

Finally, the position of the gesture apex with respect to  
 the end of the accented syllable differed depending on whether  
 there was unaccented material preceding the phrase boundary.  
 Stress-final (S and WS) patterns differed from stress-initial  
 (SW and SWW) and stress-medial WSW patterns. Figure 6  
 shows that the gesture apexes occurred during the temporal  
 limits of the accented syllable, but that their precise alignment  
 within that syllable varied depending on the presence of unac-  
 cented material preceding the phrase boundary. Thus, the ges-  
 ture apex was largely retracted when the accented syllable  
 occurred in phrase-final position (S and WS patterns), but was  
 produced towards the middle of the accented syllable when  
 there was post-accentual material preceding the right-edge  
 phrase boundary (SW, SWW, and WSW patterns).

### C. Discussion

Three main results can be observed from experiment 2.  
 First, we could confirm that the scope of a confirmatory head  
 nod gesture is the accompanying focused prosodic word, not  
 the accented syllable. This is evidenced by the fact that speak-  
 ers start head movements several milliseconds before the pro-  
 sodic word and end them several milliseconds before the  
 prosodic word is finished. Speakers maintain these patterns  
 even if there are strong edge constraints within the prosodic  
 word (i.e., the prosodic word being initiated or finished with  
 an accented syllable, as in the S, WS, SW, and SWW items).  
 Likewise, when speakers produce a gesture together with a  
 prosodic word that is less constrained in its edges (as in the  
 WSW condition), these general patterns are maintained  
 although with minor variations: the gesture onset is slightly  
 closer to the word onset and the end of the gesture is slightly  
 more distant to the end of the word.

Second, we found that the position of the peak of promi-  
 nence in the gesture (the gesture apex) is sensitive not only to  
 the position of the accented syllable (which had been found  
 in many previous studies; Fernández-Baena *et al.*, 2014; Graf  
*et al.*, 2002; Hadar *et al.*, 1983; Ishi *et al.*, 2014), but that it is  
 also highly sensitive to the distance to the upcoming intona-  
 tional phrase boundary. The position of the accented syllable  
 within the prosodic word determined where the gesture apex  
 will be produced (because gesture apexes tend to occur  
 within its limits). But the specific position of the apex within  
 the accented syllable depended on the upcoming prosodic  
 phrase boundary, because the position of the gesture apex is  
 adapted to the presence or absence of post-accentual material:  
 the gesture apex occurred closer to the end of the accented  
 syllable when there were one or more unaccented syllables  
 before the upcoming prosodic boundary; instead, the apex  
 was retracted if the upcoming prosodic boundary occurred  
 immediately after the accented syllable.

Third, complementary evidence regarding the important  
 role of the prosodic word as the domain of head nod move-  
 ments comes from the timing of the start of the gesture  
 stroke, which in our data is associated with the left-edge of  
 the prosodic word (e.g., where the word starts) rather than  
 with the accented syllable. In our data, speakers started the  
 gesture stroke before the beginning of the prosodic word,  
 and thus the gesture stroke was aligned further from the

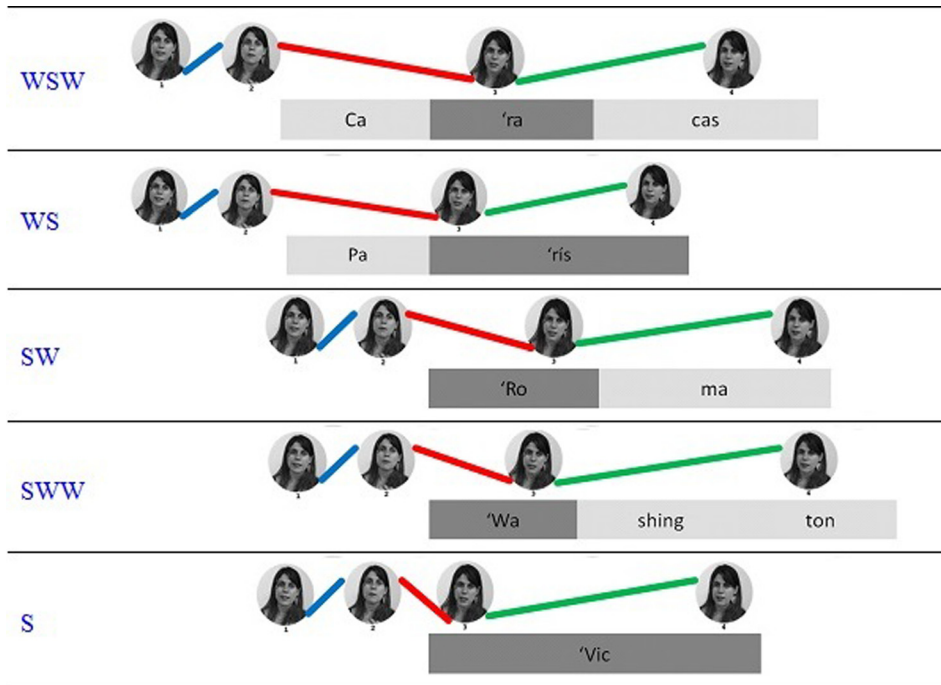


FIG. 4. (color online) Schematic representation of the alignment patterns of the head gesture and prosodic landmarks for each stress pattern. The dark grey cells represent the mean duration of the accented syllable within the prosodic word and the light grey cells the unaccented syllables. The lines connecting head images represent the gesture phases: the blue line from 1 to 2 is the preparation phase, the red line from 2 to 3 is the gesture stroke (the end of it being the gesture apex), and the green line from 3 to 4 is the retraction phase.

750 prosodic head in prosodic words with pre-accentual material  
 751 (WS and WSW patterns), and closer to the start of the pro-  
 752 sodic head when no pre-accentual material was available  
 753 (e.g., S, SW, and SWW).

754 **IV. GENERAL DISCUSSION AND CONCLUSION**

755 The aim of this study was to investigate the effects of  
 756 prosodic structure (i.e., the location of prosodic prominences  
 757 and prosodic phrase boundaries) on the timing of head nod  
 758 gestures. We designed two experiments, one that elicited  
 759 spontaneous head gestures through a *Guess Who* game and  
 760 another one that elicited semi-controlled head gestures in

which we could better control for the speakers' communi- 761  
 cative intent and the stress pattern of the target focused word. 762  
 The results of experiment 1 showed that the scope of head 763  
 movements is the whole prosodic word they accompany, and 764  
 that the peak of the head movement (the gesture apex) 765  
 occurs within the accented syllable of the prosodic word, its 766  
 exact position depending on the presence or absence of an 767  
 upcoming phrase boundary. A second experiment was 768  
 required in order to refine and confirm these results, now (1) 769  
 balancing the number of target prosodic words per stress pat- 770  
 tern, (2) analysing a more complete set of stress patterns, (3) 771  
 controlling for the speakers' communicative intent by eliciting 772  
 confirmatory sentences, and (4) measuring also the 773

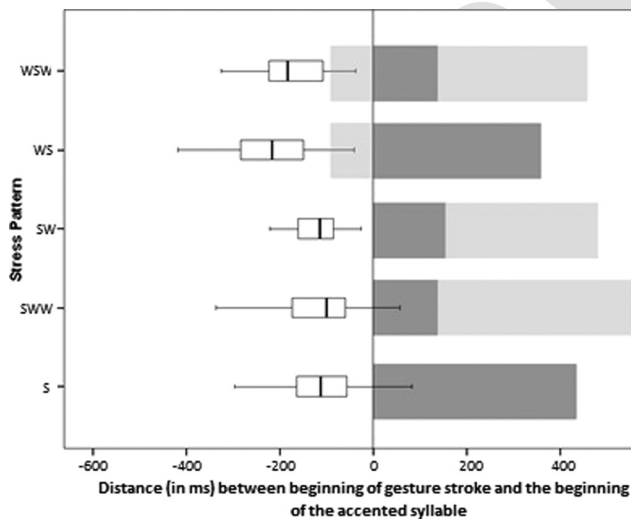


FIG. 5. Box plots displaying the temporal distance between the beginning of the gesture stroke and the beginning of the accented syllable. The 0 represents the beginning of the accented syllable, negative values showing cases where the gesture stroke started before the accented syllable and positive values the opposite. The dark grey boxes indicate the temporal limits of the accented syllable (mean values) and the light grey boxes indicate the temporal limits of the unaccented material within the prosodic word (mean values).

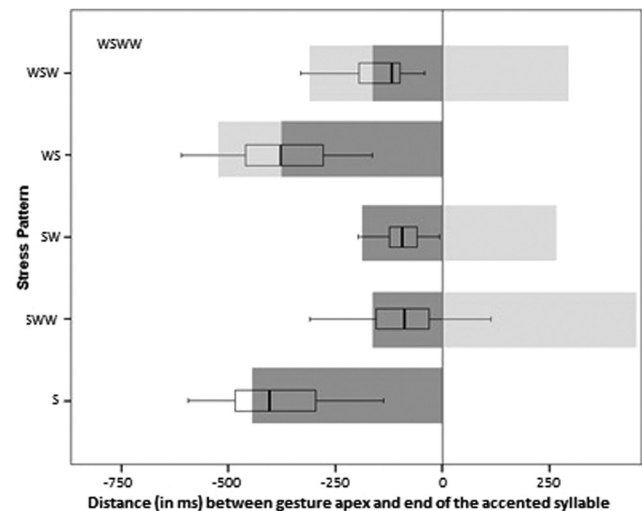


FIG. 6. Box plots displaying the temporal distance (in ms) between the gesture apex and the end of the accented syllable. The 0 represents the end of the accented syllable, negative values showing cases where the apex occurred before the end of the accented syllable and positive values the opposite. The dark grey boxes indicate the temporal limits of the accented syllable (mean values) and the light grey boxes indicate the temporal limits of the unaccented material within the prosodic word (mean values).



774 impact of the prosodic structure on the beginning of the  
 775 prominent gesture interval, the gesture stroke.  
 776 Experiment 2 confirmed that the scope of the head  
 777 movement is the accompanying focused prosodic word.  
 778 Likewise, we found that the beginning of the prosodic word  
 779 is the anchoring point for the start of the prominent interval  
 780 of the gesture movement (the gesture stroke), hence moving  
 781 it away from the accented syllable in prosodic words with  
 782 pre-accented material. Crucially, we confirmed that the peak  
 783 of the gesture movement, the apex, is timed as a function of  
 784 the prosodic heads and edges: it occurs within the accented  
 785 syllable independently of the metrical pattern of the target  
 786 word, but its exact anchoring point within that syllable is  
 787 retracted if there is an upcoming prosodic phrase boundary  
 788 and lagged if there is post-accentual material before the pro-  
 789 sodic phrase boundary occurs.  
 790 Previous research on the alignment of head gestures  
 791 with speech had shown that accented syllables were the  
 792 anchoring site for head apexes (Alexanderson *et al.*, 2013;  
 793 Fernández-Baena *et al.*, 2014; Goldenberg *et al.*, 2014; Graf  
 794 *et al.*, 2002; Hadar *et al.*, 1983; Ishi *et al.*, 2014). Yet, they  
 795 also reported variability in this pattern. Our results suggest  
 796 that an important source of variability is related to the po-  
 797 sition of prosodic edges, and specifically the distance between  
 798 the accented syllable and the upcoming prosodic phrase  
 799 boundary, a factor that none of these studies had controlled  
 800 for. Previous research on pointing gestures had shown that  
 801 the timing of pointing apexes resembles that of F0 move-  
 802 ments (because pointing apexes align with F0 peaks, and  
 803 these are retracted or lagged depending on the position of  
 804 phrase boundaries) and of oral gestures (because manual ges-  
 805 tures are lengthened at phrase boundaries) (Esteve-Gibert  
 806 and Prieto, 2013; Krivokapić *et al.*, 2015; Krivokapić *et al.*,  
 807 2016; Rochet-Capellan *et al.*, 2008). Our results reveal that  
 808 head movements are also affected by prosodic phrasing.  
 809 This seems to be due to the fact that speakers plan the timing  
 810 of their co-speech gestures by taking into account the pro-  
 811 sodic features of the interval that will accommodate their  
 812 associated gesture movements, and importantly the prosodic  
 813 head and edge positions.  
 814 These results have direct implications for applied  
 815 research. The temporal alignment of head gestures and  
 816 speech is relevant for those researchers interested in design-  
 817 ing virtual agents that interact in conversations as naturally  
 818 as possible, the so-called “talking heads.” Models of gesture-  
 819 speech temporal integration should incorporate the effects of  
 820 prosodic structure at several levels of speech planning.  
 821 Research studying the semantic integration of gesture and

speech has proposed that co-speech gestures refer to “lexical 822  
 affiliates” (Schegloff, 1984). Here we propose that the tem- 823  
 poral patterns of the gesture-speech alignment are explained 824  
 by the impact of the different levels of the prosodic hierar- 825  
 chy on the planning and execution of the gesture movement. 826

Future studies should further investigate this entrain- 827  
 ment between gesture and prosodic structure in speech. 828  
 More work is needed to investigate how prosodic domains 829  
 affect the temporal patterns in the realization of co-speech 830  
 gestures. In our materials, for instance, we cannot disentangle 831  
 whether the scope of the gesture movement is the lexical 832  
 word or the prosodic word. Also, if prosodic structure 833  
 strongly constrains the timing of head nod gestures (and co- 834  
 speech gestures in general), speakers should have fine- 835  
 grained perceptual expectations about gesture timing if a 836  
 specific prosodic structure is predicted in the discourse. 837  
 Finally, the influence of the semantic and pragmatic aspects 838  
 of a gesture on its temporal implementation deserves further 839  
 investigation, as recent studies examining spontaneously eli- 840  
 cited gestures suggest that this influence can induce different 841  
 types of gesture-speech temporal integration (e.g., 842  
 Bergmann *et al.*, 2011; Esteve-Gibert *et al.*, 2014). 843

What seems to be beyond question is that there is tight 844  
 temporal integration of gesture and speech, and that prosodic 845  
 structure is one of the main aspects controlling this temporal 846  
 coordination. Speakers use speech and gesture together to 847  
 transmit their message, and discourse prominence is commu- 848  
 nicated at both the visual and acoustic levels by integrating 849  
 the phases of gesture movements with the prosodic structure 850  
 of oral messages. 851

**ACKNOWLEDGMENTS** 852

Thanks to Igor Jauk for the gesture coding in 853  
 experiment 1, and Suleman Shahid and Constantijn Kaland 854  
 for helping us with the experimental setting and recordings. 855  
 This research has been funded by the Spanish MINECO 856  
 (grant FFI2015-66533-P), and by the Generalitat de 857  
 Catalunya to the Prosodic Studies Group (2014SGR-925), by 858  
 the 2010 BE1 00207 travelling grant awarded to the second 859  
 author of the study, and by the Labex BLRI (ANR-11- 860  
 LABX-0036) grant awarded to the first author of the study. 861

**APPENDIX** 862

Descriptive results of all the analyses in experiments 1 863  
 and 2 are given in Table VI (all duration and distance mea- 864  
 sures are in milliseconds). 866

TABLE VI. Descriptive results of all the analyses in Experiments 1 and 2 (all duration and distance measures are in milliseconds).

	S	WS	WSW	SW	SWW <sup>a</sup>
Experiment 1					
Duration accented syllable	M = 434.5 (SD = 116.3) <sup>b</sup>	M = 420 (SD = 92.4)	M = 169.3 (SD = 35.3)	M = 164.2 (SD = 54.3)	—
Distance onset word / onset accented syllable	M = 0 (SD = 0)	M = -149.8 (SD = 40.2)	M = -124.1 (SD = 48.8)	M = 0 (SD = 0)	—

TABLE VI. (Continued.)

	S	WS	WSW	SW	SWW <sup>a</sup>
Distance offset accented syllable / offset word	M = 0 (SD = 0)	M = 0 (SD = 0)	M = -291.5 (SD = 94.9)	M = -251.2 (SD = 79.3)	—
Distance onset gesture / onset word	M = -335.3 (SD = 326)	M = -245.6 (SD = 339)	M = -236.4 (SD = 392)	M = -277.1 (SD = 346)	—
Distance offset gesture / offset word	M = 286.2 (SD = 599)	M = 249.4 (SD = 674)	M = 224.8 (SD = 492)	M = .098 (SD = 491)	—
Distance apex / offset accented syllable	M = -371.4 (SD = 352.7)	M = -482.8 (SD = 368.7)	M = -118.2 (SD = 309.1)	M = -116.9 (SD = 345.5)	—
<i>Experiment 2</i>					
Duration of the accented syllable	M = 431.2 (SD = 116.7)	M = 378.7 (SD = 92.3)	M = 149.9 (SD = 28.2)	M = 177.2 (SD = 46.9)	M = 149.9 (SD = 38.1)
Distance onset word / onset accented syllable	M = 0 (SD = 0)	M = -134.2 (SD = 37.8)	M = -132.7 (SD = 30.1)	M = 0 (SD = 0)	M = 0 (SD = 0)
Distance offset accented syllable / offset word	M = 0 (SD = 0)	M = 0 (SD = 0)	M = -288.3 (SD = 81.3)	M = -273.5 (SD = 76.5)	M = -382.1 (SD = 109.6)
Duration preparation phrase of the gesture	M = 164.3 (SD = 88.8)	M = 211.5 (SD = 102.9)	M = 177.5 (SD = 61.2)	M = 148.4 (SD = 72.7)	M = 179.3 (SD = 113.8)
Duration of the gesture stroke	M = 170.8 (SD = 53.6)	M = 221.1 (SD = 83.3)	M = 184.4 (SD = 58.1)	M = 204.9 (SD = 65.7)	M = 181.1 (SD = 52.2)
Duration retraction phase of the gesture	M = 247.9 (SD = 118.5)	M = 248.9 (SD = 109.2)	M = 227.4 (SD = 105.8)	M = 234.8 (SD = 111.1)	M = 266.1 (SD = 122.1)
Distance onset gesture / onset word	M = -290.2 (SD = 129.7)	M = -300.7 (SD = 132.7)	M = -221.8 (SD = 94.3)	M = -277.9 (SD = 114.4)	M = -301.4 (SD = 112.3)
Distance offset gesture / offset word	M = -138.3 (SD = 158.1)	M = -131.2 (SD = 138.4)	M = -203.3 (SD = 109.1)	M = -140.4 (SD = 89.1)	M = -206.8 (SD = 161.5)
Distance onset stroke / onset accented syllable	M = -125.9 (SD = 102.8)	M = -223.5 (SD = 109.5)	M = -177 (SD = 79.8)	M = -129.5 (SD = 72.5)	M = -122.1 (SD = 91)
Distance apex / offset accented syllable	M = -386.2 (SD = 115.2)	M = -381.1 (SD = 120.5)	M = -142.5 (SD = 77.6)	M = -101.7 (SD = 63.6)	M = -90.9 (SD = 91.3)

<sup>a</sup>This column is empty in experiment 1 because this stress pattern was not observed in Experiment 1.

<sup>b</sup>Mean, M; standard deviation, SD.

867 Ahmad, M. I., Tariq, H., Saeed, M., Shahid, S., and Kraemer, E. (2011).  
868 "Guess who? An interactive and entertaining game-like platform for investigat-  
869 ing human emotions," in *Human Computer Interaction. Towards*  
870 *Mobile and Intelligent Interaction Environments*, Lecture Notes in  
871 Computer Science 6763, edited by J. A. Jacko (Springer, Berlin,  
872 Germany), Vol. 3, pp. 543–551.

873 Alexanderson, S., House, D., and Beskow, J. (2013). "Aspects of co-  
874 occurring syllables and head nods in spontaneous dialogue," in  
875 *Proceedings of 12th International Conference on Auditory-Visual Speech*  
876 *Processing (AVSP2013)*.

877 Ambrazaitis, G., Svensson Lundmark, M., and House, D. (2015). "Head  
878 movements, eyebrows, and phonological prosodic prominence levels in  
879 Stockholm Swedish news broadcasts," in *FAAVSP - The 1st Joint*  
880 *Conference on Facial Analysis, Animation, and Auditory-Visual Speech*  
881 *Processing*, Vienna, Austria, pp. 42–42.

882 Barkhuysen, P., Kraemer, E., and Swerts, M. (2008). "The interplay  
883 between the auditory and visual modality for end-of-utterance detection,"  
884 *J. Acoust. Soc. Am.* **123**, 354–365.

885 Bates, D., Maechler, M., and Bolker, B. (2011). "lme4: Linear mixed-effects  
886 models using S4 classes [R package version 0.99375-39]," [http://](http://CRAN.R-project.org/package=lme4)  
887 [CRAN.R-project.org/package=lme4](http://CRAN.R-project.org/package=lme4) (Last viewed January 27, 2017).

888 Bergmann, K., Aksu, V., and Kopp, S. (2011). "The relation of speech and  
889 gestures: Temporal synchrony follows semantic synchrony," in  
890 *Proceedings of the 2nd Workshop on Gesture and Speech in Interaction*,  
891 pp. 1–6.

892 Bergmann, K., Kahl, S., and Kopp, S. (2014). "How is information distrib-  
893 uted across speech and gesture? A cognitive modeling approach," *Cognit.*  
894 *Processing* **15**(1), S84–S87.

895 Billmyer, K., and Varghese, M. (2000). "Investigating instrument-based  
896 pragmatic variability: Effects of enhancing discourse completion tests,"  
897 *Appl. Linguist.* **21**(4), 517–552.

898 Blum-Kulka, S., House, J., and Kasper, G. (1989). "Investigating cross- cul-  
899 tural pragmatics: An introductory overview," in *Cross-Cultural Pragmatics: Requests and Apologies*, edited by S. Blum-Kulka, J. House, and G. Kasper (Ablex, Norwood, NJ), pp. 1–34.

Boersma, P., and Weenink, D. (2012). "Praat: Doing phonetics by com-  
puter," <http://www.praat.org/> (Last viewed July 25, 2016).

Butterworth, B., and Beattie, G. (1978). "Gesture and silence as indicators  
of planning in speech," in *Recent Advances in the Psychology of*  
*Language: Formal and Experimental Approaches*, edited by R. Campbell  
and G. T. Smith (Plenum Press, New York), pp. 347–360.

De Ruiter, J. P. (1998). "Gesture and speech production," doctoral disserta-  
tion, Katholieke Universiteit, Nijmegen, the Netherlands.

Dohen, M., and Loevenbruck, H. (2004). "Pre-focal rephrasing, focal  
enhancement and post-focal deaccentuation in French," in *Proceedings of*  
*the 8th International Conference on Spoken Language Processing*, pp.  
2–5.

Esteve-Gibert, N., Pons, F., Bosch, L., and Prieto, P. (2014). "Are gesture  
and prosodic prominences always coordinated? Evidence from perception  
and production," in *Proceedings of the Speech Prosody Conference*, edited  
by N. Campbell, D. Gibbon, and D. Hirst, pp. 222–226.

Esteve-Gibert, N., and Prieto, P. (2013). "Prosodic structure shapes the tem-  
poral realization of intonation and manual gesture movements," *J. Speech*  
*Language Hear. Res.* **56**, 850–864.

Fernández-Baena, A., Montaña, R., Antonijoan, M., Roversi, A., Miralles,  
D., and Alías, F. (2014). "Gesture synthesis adapted to speech emphasis,"  
*Speech Commun.* **57**, 331–350.

Ferré, G. (2014). "A multimodal approach to markedness in spoken  
French," *Speech Commun.* **57**, 268–282.

Goldenberg, D., Tiede, M., Honorof, D. N., and Mooshammer, C. (2014).  
"Temporal alignment between head gesture and prosodic prominence in  
naturally occurring conversation: An electromagnetic articulometry  
study," *J. Acoust. Soc. Am.* **135**, 2294.

Graf, H. P., Cosatto, E., Strom, V., and Huang, F. J. (2002). "Visual prosody:  
Facial movements accompanying speech," in *Proceedings of the 5th IEEE*  
*International Conference on Automatic Face Gesture Recognition*, pp. 396–401.

- 933 Guellai, B., Langus, A., and Nespov, M. (2014). "Prosody in the hands of  
934 the speaker," *Front. Psychol.* **5**, 1–8.
- 935 Hadar, U., Steiner, T. J., Grant, E. C., and Rose, F. C. (1983). "Kinematics  
936 of head movements accompanying speech during conversation," *Human  
937 Movement Sci.* **2**(1–2), 35–46.
- 938 House, D., Beskow, J., and Granström, B. (2001). "Timing and interaction  
939 of visual cues for prominence in audiovisual speech perception," in  
940 *Proceedings of Eurospeech*, pp. 387–390.
- 941 Ishi, C. T., Ishiguro, H., and Hagita, N. (2014). "Analysis of relationship  
942 between head motion events and speech in dialogue conversations,"  
943 *Speech Commun.* **57**, 233–243.
- 944 Jannedy, S., and Mendoza-Denton, N. (2005). "Structuring Information  
945 through Gesture and Intonation," *Interdisciplinary Stud. Inf. Struct.* **3**,  
946 199–244.
- 947 Kelly, S. D., Ozyürek, A., and Maris, E. (2010). "Two sides of the same  
948 coin: Speech and gesture mutually interact to enhance comprehension,"  
949 *Psychol. Sci.* **21**(2), 260–267.
- 950 Kendon, A. (1980). "Gesticulation and speech: Two aspects of the process  
951 of utterance," in *The Relationship of Verbal and Nonverbal  
952 Communication*, edited by M. R. Key (Mouton, the Hague, the  
953 Netherlands), pp. 207–227.
- 954 Kim, J., Cvejic, E., and Davis, C. (2014). "Tracking eyebrows and head ges-  
955 tures associated with spoken prosody," *Speech Commun.* **57**, 317–330.
- 956 Krahmer, E., and Swerts, M. (2007). "The effects of visual beats on prosodic  
957 prominence: Acoustic analyses, auditory perception and visual  
958 perception," *J. Mem. Language* **57**(3), 396–414.
- 959 Krivokapić, J. (2014). "Gestural coordination at prosodic boundaries and its  
960 role for prosodic structure and speech planning processes," *Philos. Trans.  
961 R. Soc. London Ser. B Biol. Sci.* **369**(1658), 20130397.
- 962 Krivokapić, J., Tiede, M. K., and Tyrone, M. E. (2015). "A kinematic analy-  
963 sis of prosodic structure in speech and manual gestures," in *Proceedings  
964 of the 18th International Congress of Phonetic Sciences*.
- 965 Krivokapić, J., Tiede, M. K., Tyrone, M. E., and Goldenberg, D. (2016).  
966 "Speech and manual gesture coordination in a pointing task," in  
967 *Proceedings of the 8th International Conference on Speech Prosody*, pp.  
968 1240–1244.
- 969 Lausberg, H., and Sloetjes, H. (2009). "Coding gestural behavior with the  
970 NEUROGES-ELAN system," *Behav. Res. Methods Instrum. Comput.*  
971 **41**(3), 841–849.
- 972 Leonard, T., and Cummins, F. (2011). "The temporal relation between beat  
973 gestures and speech," *Lang. Cognit. Processes* **26**(10), 1457–1471.
- 974 Levelt, W. J. M., Richardson, G., and La Heij, W. (1985). "Pointing and  
975 voicing in deictic expressions," *J. Mem. Language* **24**, 133–164.
- 976 Loehr, D. P. (2012). "Temporal, structural, and pragmatic synchrony  
977 between intonation and gesture," *Lab. Phonol.* **3**, 71–89.
- 978 McClave, E. Z. (2000). "Linguistic functions of head movements in the con-  
979 text of speech," *J. Pragmatics* **32**, 855–878.
- 980 McNeill, D. (1992). *Hand and Mind: What Gestures Reveal About Thought*  
981 (University of Chicago Press, Chicago, IL).
- Nobe, S. (2000). "Where to most spontaneous representational gestures  
982 actually occur with respect to speech?," in *Language and Gesture*, edited  
983 by D. McNeill (Cambridge University Press, Cambridge, UK), pp.  
984 186–198.
- Özyürek, A., Willems, R. M., Kita, S., and Hagoort, P. (2007). "On-line  
986 integration of semantic information from speech and gesture: Insights  
987 from event-related brain potentials," *J. Cognit. Neurosci.* **19**(4), 605–616.
- Poggi, I., D'Errico, F., and Vincze, L. (2011). "68 Nods. But not only of  
989 agreement," in *68 Zeichen Für Roland Posner. Ein Semiotisches Mosaik.  
990 (68 Signs for Roland Posner. A Semiotic Mosaic)* (Stauffenburg Verlag,  
991 Tübingen, Germany).
- Poggi, I., D'Errico, F., Vincze, L., and Milazzo, V. (2010). "Types of nods.  
993 The polysemy of a social signal," in *Proceedings of the Seventh conference  
994 on International Language Resources and Evaluation (LREC'10)*, Malta.
- Prieto, P., and Ortega-Llebaria, M. (2009). "Do complex pitch gestures  
996 induce syllable lengthening in Catalan and Spanish?," in *Phonetics and  
997 Phonology: Interactions and Interrelations*, edited by M. Vigário, S.  
998 Frota, and M. J. Freitas (John Benjamins, Philadelphia, PA), pp. 51–70.
- Prieto, P., Pugliesi, C., Borràs-Comes, J., Arroyo, E., and Blat, J. (2015).  
1000 "Exploring the contribution of prosody and gesture to the perception of  
1001 focus using an animated agent," *J. Phonetics* **49**, 41–54.
- Rochet-Capellan, A., Laboissière, R., Galván, A., and Schwartz, J. (2008).  
1003 "The speech focus position effect on jaw-finger coordination in a pointing  
1004 task," *J. Speech Language Hearing Res.* **51**(6), 1507–1521.
- Roustan, B., and Dohen, M. (2010). "Gesture and speech coordination: The  
1006 influence of the relationship between manual gesture and speech," in  
1007 *Proceedings of 11th Annual Conference of the International Speech  
1008 Communication Association (INTERSPEECH 2010)*, Makuhari, Japan.
- Rusiewicz, H. L., Shaiman, S., Iverson, J. M., and Szuminsky, N. (2013).  
1010 "Effects of prosody and position on the timing of deictic gestures,"  
1011 *J. Speech Language Hear. Res.* **56**(2), 458–470.
- Schegloff, E. A. (1984). "On some gestures' relation to talk," in *Structures  
1013 of Social Action*, edited by J. M. Atkinson and J. Heritage (Cambridge  
1014 University Press, Cambridge, UK), pp. 266–298.
- Shattuck-Hufnagel, S., Ren, P. L., and Tauscher, E. (2010). "Are torso  
1016 movements during speech timed with intonational phrases?," in  
1017 *Proceedings of the Speech Prosody 2010*, Chicago, IL.
- Swerts, M., and Krahmer, E. (2010). "Visual prosody of newsreaders:  
1019 Effects of information structure, emotional content and intended audience  
1020 on facial expressions," *J. Phonetics* **38**, 197–206.
- Treffner, P., Peter, M., and Kleidon, M. (2008). "Gestures and phases: The  
1022 dynamics of speech-hand communication," *Ecol. Psychol.* **20**(1), 32–64.
- Wagner, P., Malisz, Z., and Kopp, S. (2014). "Gesture and speech in interac-  
1024 tion: An overview," *Speech Commun.* **57**, 209–232.
- Yasinnik, Y., Renwick, M., and Shattuck-Hufnagel, S. (2004). "The timing  
1026 of speech-accompanying gestures with respect to prosody," in  
1027 *Proceedings From Sound to Sense: 50+ Years of Discoveries in Speech  
1028 Communication* (MIT, Cambridge, MA), pp. C97–C102.