

Extracting and analyzing head movements accompanying spontaneous dialogue

Simon Alexanderson (simonal@kth.se)

Department of Speech, Music and Hearing, KTH, Lindstedtsvägen 24
10044 Stockholm, Sweden

David House (davidh@speech.kth.se)

Department of Speech, Music and Hearing, KTH, Lindstedtsvägen 24
10044 Stockholm, Sweden

Jonas Beskow (beskow@speech.kth.se)

Department of Speech, Music and Hearing, KTH, Lindstedtsvägen 24
10044 Stockholm, Sweden

Abstract

This paper reports on a method developed for extracting and analyzing head gestures taken from motion capture data of spontaneous dialogue in Swedish. Candidate head gestures with beat function were extracted automatically and then manually classified using a 3D player which displays time-synced audio and 3D point data of the motion capture markers together with animated characters. Prosodic features were extracted from syllables co-occurring with a subset of the classified gestures. The beat gestures show considerable variation in temporal synchronization with the syllables, while the syllables generally show greater intensity, higher F0, and greater F0 range when compared to the mean across the entire dialogue. Additional features for further analysis and automatic classification of the head gestures are discussed.

Keywords: Gestures; prosody; motion capture; beats; head nods; stressed syllable

Introduction

Prosody and especially intonation, or the melody of speech, plays a crucial role in regulating face-to-face spoken interaction. The meaning of a spoken utterance is shaped and reinforced by intonation whereby important words and phrases are highlighted and made prominent by changing intonation. This type of prominence is often referred to as focal or sentence accent and consists of a marked tonal movement synchronized with the word in the sentence which conveys new information (Ladd, 1996; Hirst & Di Cristo 1998). Old or given information can be backgrounded or deaccented (Bruce & Touati 1992) resulting in an intonational structure which reflects the information structure of the utterance and reveals how the information is organized in the speaker's consciousness (Chafe 1994).

Along with intonation there are a wealth of visual gestures, including head, facial and body movements, which co-occur with speech adding emphasis and prominence to portions of the utterance and contributing to the flow of the dialogue. There is a growing body of research concerning the type and frequency of gestures co-occurring with speech (e.g. Allwood & Cerrato, 2003; Bergmann et al., 2011; Beskow et al., 2006; Gullberg & Kita, 2009; Kendon, 1980; McNeill, 2005). Beat gestures such as rapid hand movement

as described by McNeill (2005) are particularly interesting in this regard as they coincide and appear to be synchronized with prosodic and intonational peaks related to prominence. They also share the same prominence-lending function as the focal or sentence accent, but they can be repetitive marking the stress or rhythmical structure of an utterance.

Studying synchronization between speech and gesture has played an important role in building theories of human communication which approach speech and gesture production as arising from a common generation process (Kendon, 2004; McNeill, 1992). In terms of timing and synchronization, there are many similarities between intonational gestures and visual gestures produced in accompaniment with speech. In terms of the speech production process, both intonation and visual gestures are free to vary across the vowels and consonants of the segments. In intonation, however, this variation is restricted by the specific patterns used by a language to signal meaning in spoken interaction. Intonation research carried out within the past several decades has led to the development of models of intonational meaning for a wide variety of languages in which timing differences can change the lexical identity of a word as well as adding prominence and signaling type of feedback such as confirmation or questioning (Hirst & Di Cristo 1998).

If we wish to study the timing of gestures in the same way as we approach timing in intonation, we currently lack an established methodology to extract and analyze gestures, especially gestures occurring in spontaneous dialogue. The Spontal corpus of Swedish dialogue provides a rich database as a point of departure for testing gesture extraction and analysis methodology. The database, containing more than 60 hours of unrestricted conversation in over 120 dialogues between pairs of speakers is comprised of high-quality audio and video recordings (high definition) and motion capture for body and head movements for all recordings (Edlund et al., 2010).

The progression and timing of the motion of a head nod can be described in much the same terms as an intonational excursion. However, many speakers move their heads extensively while speaking, and manual annotation of head-

gestures in spontaneous dialogue involves a number of difficulties. Among the sources of disagreement are segment boundaries and location of maximum extent. Gestures may be multifunctional and involve simultaneous rotations around several axes. In this study we present a semi-automatic approach to head-gesture annotation, in which the main goal is to test its viability and potential for annotation of gesture data on a large scale, such as is represented by the Spontal database.

Method

Motion capture offers many possibilities in gesture research. Compared to video-based head pose estimators, it offers increased accuracy and higher frame-rate (100 fps in our current study).

To overcome some of the difficulties of head-gesture annotation we developed a semi-automatic annotation procedure consisting of two steps. First an automatic head-gesture segmentation algorithm is applied to the motion capture data and then the segments are manually classified by the annotators. In addition to gesture annotation we also processed the audio files of the speakers and generated pitch and intensity data, talk vs. no-talk segmentation and syllable segmentation. This was done to be able to investigate the relationship between the head-gestures and prosodic features. In this exploratory stage we chose one of the dialogues from the Spontal database where the speakers demonstrated a relatively large number and variety of head-movements. The participants were a male and female who did not know each other.

Automatic segmentation of head-nods

A simplistic segmentation approach was used for head-nod segmentation. The head orientations were calculated from three markers attached to the headbands of the speakers and expressed in an Euler angle form. We then calculated the angular velocity of the pitch component as the basis for segmentation. During a head nod the angular velocity follows an oscillatory movement during a limited time period, and we use its local extreme values as segment boundaries. A segment is defined as a maximum velocity followed by a minimum or vice versa. Two thresholds are used in this process. The first enforces the peak velocity to be over a specified value, thus prohibiting small movements caused by noise to be interpreted as nods, and the second enforces the nod segments to be shorter than a specified duration. Using the velocity peaks as segment boundaries has some desirable features. During head-nods there are rapid changes in angular velocity causing clear detectable spikes in the data. It also naturally splits repeated head-nods into a consecutive sequence of down-up (nod) and up-down (jerk) segments, which fits well with the MUMIN multimodal annotation scheme proposed by Allwood et al., (2007).

For our data, the minimum peak velocity threshold was empirically set to a value of 0.0015 radians/s, which was the lowest value before noise in the data would be manifested as

segments. The maximum segment duration threshold was set to 1000 ms. The thresholds were set to detect a maximum number of valid candidates. In future work, we plan to hand-annotate part of the data in order to allow for an objective verification of the quality of the segmentation.

In the current study we were interested in gestures with beat-function produced in companion with stressed syllables. We therefore discarded all segments occurring while the subject was not speaking and further all nods in the up-down order, leaving all down-up nods occurring during speech as our candidate gestures. The segmentation of talk vs. no-talk was performed with an automatic speech activity detection algorithm (Heldner et al., 2011).

Manual classification

After running the automatic processing, the resulting segments were examined and manually classified by two annotators. To carry out the classification, the annotators viewed each segment in a specially designed 3D player which plays time-synchronized audio and displays 3D point data of the motion capture markers together with animated characters following the 3D marker motion, see figure 1. This was done to enable the annotators to work on the same signal as the algorithm. In addition, the 3D-player allows for alternate view-points of the scene and maintains the high frame rate of the motion data. In some uncertain cases the annotators also used the video data as a reference. In future releases we plan to include this feature in the player.



Figure 1: 3D viewer for manual annotation.

As expected, the segments from the automatic process did not only contain unambiguous beat gestures, but also gestures with other functions co-occurring with speech (McClave, 2000). Such other functions were feedback, confirmation, word or phrase intensification and listing of lexical items. Moreover, some of the extracted gestures did not appear to co-occur with a stressed syllable.

Therefore, an annotation scheme was devised with two main queries: Q1, “Is there a clear nod in synchrony with a stressed syllable?” and Q2, “Is the nod multifunctional?” If the answer to the first query is positive the second query is also answered. This scheme resulted in three categories: 1. No clear nod in synchrony (no sync), 2. A clear nod with a beat function (beat-function), and 3. A clear nod which is multifunctional (multi-function). The annotation time was approximately two nods per minute.

Prosodic features

The pitch and intensity curves were extracted from the audio signals of the speakers using the SNACK toolkit (Sjölander & Beskow, 2000). Also syllable boundaries and nuclei were derived by applying Mermelstien’s convex-hull algorithm (Mermelstien, 1975).

After gesture- and prosodic feature extraction was performed, we determined which syllable was closest in time to the maximum rotation of the nod. The time difference between each gesture and the start and nucleus of its closest syllable was then calculated. Also pitch and intensity properties of the closest syllable were compared with mean values across the total dialogue.

Results

Gesture extraction

The automatic segmentation algorithm was applied on the 20 minute dialog, extracting 64 nod-segments for speaker 1 (male) and 150 segments from speaker 2 (female). The manual classification by the two annotators resulted in a 69% and 65% agreement for the head-nods of speaker 1 and speaker 2 respectively. As is displayed in table 1 the annotators showed greatest agreement for the category with no syllable-synchronization. Less agreement was obtained from the categories beat-function and multi-function. Annotator 1 perceived more of the nods having beat-function while annotator 2 perceived more having multi-function.

Table 1: Results of the manual annotation showing class and agreement. The first number in each pair is the male speaker, the second is the female.

| Class | Annotator 1 | Annotator 2 | Agreement |
|----------------|-------------|-------------|-----------|
| No sync | 23/44 | 15/41 | 13/29 |
| Beat-function | 9/74 | 11/66 | 4/50 |
| Multi-function | 32/32 | 38/43 | 27/18 |
| Total | 64/150 | 64/150 | 44/97 |

Gesture duration

Figure 2 shows the durations of the segments in the different categories for speaker 1 and speaker 2 for those gestures for which the annotators agreed.

Note that the segment length is the part of the nod between the peak velocities of the downwards and upwards phase as described earlier. The results show a tendency for the multi-functional nods to be shorter than those with a beat function. The nods classified as non-synchronous showed greater temporal variation than the other categories.

Gesture timing related to syllables

The subset of gestures annotated as having a beat function for the female speaker was analyzed in terms of timing related to its closest syllable. Only the gestures from the

female speaker were analyzed due to the small number of beat gestures annotated for the male speaker. Figure 3 shows the time difference between two different anchor-points of the syllable (onset and nucleus) and three different phases of the nod: peak velocity of the downward phase (p1), max rotation (p2) and peak velocity of the upward phase (p3). The timing relationship between the gesture and the syllable does not seem to be influenced by the choice of syllable anchor-point. The timing relationships show a considerable amount of variation regarding the question of gesture synchronization with the syllable.

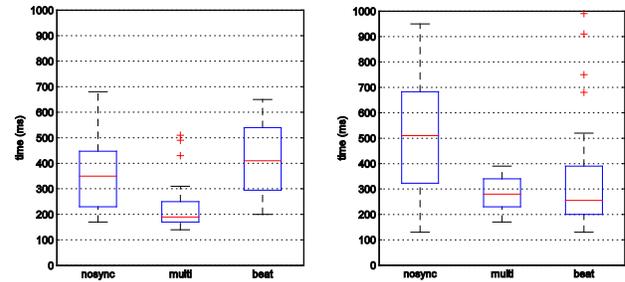


Figure 2: Durations of the agreed nods in the classes for speaker 1 (left) and speaker 2 (right).

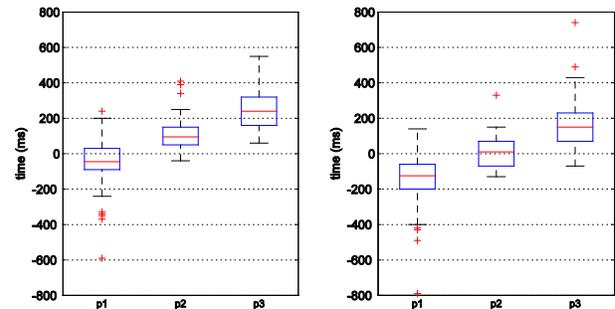


Figure 3: Timing of syllable anchor-points, onset (left) and nucleus (right), with respect to three different phases of the nod: peak velocity of the downward phase (p1), max rotation (p2) and peak velocity of the upward phase (p3).

Syllable comparison

When compared with mean values across the total dialogue the syllables closest to the annotated beat nods generally showed greater integrated intensity, higher F0 at the nucleus, and greater F0 range. This is displayed in table 2.

Table 2: Comparison of syllable features.

| | integrated intensity | F0 at nucleus | F0 range |
|-----------------------|----------------------|---------------|----------|
| Mean closest syllable | 13.2 dBs | 217 Hz | 77 Hz |
| Mean across dialogue | 11.1 dBs | 180 Hz | 67 Hz |

Discussion

The goal of this study was primarily to develop and test a new method for extracting and annotating gesture data. While head gestures, and in particular relatively small gestures, have been problematic for manual annotation schemes, the semi-automatic method tested here shows promise for the selection and fast annotation of multimodal data.

This process is a starting point for further work in the field of automatic recognition and classification of multimodal communication. Given the fact that non-verbal communication and verbal communication are tightly coupled, the motion data may provide important and robust features for machine-learning techniques. In this study we started an investigation along this path by analyzing features for prominence detection and their coupling to beat gestures. This method may also prove useful in analyzing features related to other communicative functions such as feedback and turn-taking.

While the results concerning the analysis of the characteristics of the head nods and their timing must be seen as quite preliminary due to the small sample and very limited classification categories, the timing results are consistent with those results reported in Leonard & Cummins (2011) for hand and arm beat gestures. More available data and the development of automatic methods and tools should better enable us to compare and evaluate results such as these.

Acknowledgments

The work reported here is carried out within the projects: "Timing of intonation and gestures in spoken communication," (P12-0634:1) funded by the Bank of Sweden Tercentenary Foundation, and "Large-scale massively multimodal modelling of non-verbal behaviour in spontaneous dialogue," (VR 2010-4646) funded by the Swedish Research Council.

References

- Allwood, J. & Cerrato, L. (2003). A study of gestural feedback expressions. In Paggio, P., Jokinen, K. & Jönsson, A. (Eds.) *Proc. of the First Nordic Symposium on Multi-modal Communication*. (pp. 7-20). Copenhagen
- Allwood, J., Cerrato, L., Jokinen, K., Navarretta, C., Paggio, P. (2007). The MUMIN coding scheme for the annotation of feedback, turn management and sequencing phenomena. *Language Resources and Evaluation*, 41, 273-287.
- Bergmann, K., Aksu, V. & Kopp, S. 2011. The Relation of Speech and Gestures: Temporal Synchrony Follows Semantic Synchrony, *Proceedings of the 2nd Workshop on Gesture and Speech in Interaction*. Bielefeld, Germany.
- Beskow, J., Granström, B., & House, D. (2006). Visual correlates to prominence in several expressive modes. In *Proceedings of Interspeech 2006* (pp. 1272–1275). Pittsburgh, PA.
- Bruce, G. & Touati, P. (1992). On the analysis of prosody in spontaneous speech with exemplification from Swedish and French. *Speech Communication*, 11, 453–458.
- Chafe, W. (1994). *Discourse, Consciousness, and Time: The Flow and Displacement of Conscious Experience in Speaking and Writing*. Chicago: The University of Chicago Press.
- Edlund, J., Beskow, J., Elenius, K., Hellmer, K., Strömbergsson, S., & House, D. (2010). Spontal: a Swedish spontaneous dialogue corpus of audio, video and motion capture. In Calzolari, N., Choukri, K., Maegaard, B., Mariani, J., Odjik, J., Piperidis, S., Rosner, M., & Tapias, D. (Eds.), *Proc. of the Seventh conference on International Language Resources and Evaluation (LREC'10)* (pp. 2992 - 2995). Valetta, Malta.
- Gullberg, M. & Kita, S. (2009). Attention to speech-accompanying gestures: Eye movements and information uptake. *Journal of Nonverbal Behavior*, 33, 251-277.
- Heldner, M., Edlund, J., Hjalmarsson, A., & Laskowski, K. (2011). Very short utterances and timing in turn-taking. *Proceedings of Interspeech 2011* (pp. 2837-2840). Florence, Italy.
- Hirst, D. & Di Cristo, A. (1998). A survey of intonation systems. In: Hirst, D. & Di Cristo, A. (Eds.), *Intonation Systems*. Cambridge: Cambridge University Press.
- Leonard, T. and Cummins, F. (2011). The temporal relation between beat gestures and speech. *Language and Cognitive Processes*, 26, 1457-1471.
- Kendon, A. (1980). Gesticulation and speech: Two aspects of the process of utterance. In M. R. Key (Ed.), *The relationship of verbal and nonverbal communication*. The Hague: Mouton.
- Kendon, A. (2004). *Gesture: Visible action as utterance*. Cambridge: Cambridge University Press.
- Ladd, D.R. (1996). *Intonation Phonology*. Cambridge: Cambridge University Press.
- McClave, E. (2000). Linguistic functions of head movements in the context of speech. *Journal of Pragmatics*, 32, 855-878.
- McNeill, D. (1992). *Hand and mind: What gestures reveal about thought*. Chicago: The University of Chicago Press.
- McNeil, D. (2005). *Gesture and thought*. Chicago: The University of Chicago Press.
- Mermelstien, P. (1975). Automatic segmentation of speech into syllabic units. *Journal of the Acoustical Society of America*, 58, 880-883.
- Sjölander, K., & Beskow, J. (2000). WaveSurfer - an open source speech tool. In Yuan, B., Huang, T., & Tang, X. (Eds.), *Proceedings of ICSLP 2000, 6th Intl Conf on Spoken Language Processing* (pp. 464-467). Beijing.