# Anaphora and Gestures in Multimodal Communication

Costanza Navarretta

Centre for Language Technology,
University of Copenhagen,
Njalsgade 140, build. 25,
2300 Copenhagen S, Denmark

**Abstract.** This paper describes a pilot study of gestures which are connected to anaphoric expressions and to their antecedents in video-recordings of spontaneous two- and three party interactions. The recordings have been transcribed, multimodally annotated and analysed. The results of this analysis and of machine learning experiments run on the annotated data show that gestures which are related to anaphora (or co-referring expression) and to their antecedents have many common shape attributes and values. They also show that shape attributes and values can be used for identifying gestures connected to referring expressions automatically. These results are promising for both anaphora resolution and for the generation of plausible conversational agents.

**Keywords:** anaphora, gestures, multimodal communication, spontaneous interactions, semi-supervised and supervised learning

## 1  Introduction

This paper deals with gestures which are connected with anaphoric and co-referring expressions in video-recordings of spontaneous interactions between two or more participants. Human communication is by nature multimodal because it involves both speech and gestures. Here, we use "gestures" as a general term comprising non-verbal communicative behaviours, such as facial expressions, head movements, hand gestures, and body postures.

Most studies on co-reference and anaphora have focused on their occurrences in texts and speech. An exception is the work by Eisenstein and Davis [6, 7] who extract via computer vision features of hand gestures co-occurring with co-referring nominal expressions in an English multimodal corpus. They base their study on a corpus collected for that task where the dialogue participants wore coloured hand gloves in order to facilitate the automatic recognition of the gestures. The results of Eisenstein and Davis' study show that the position of holds in the hand gestures can be useful to co-reference resolution. Eisenstein et al. [8] add these features into a reference system which they run on the same corpus. The added features improve co-reference. Chen et al. [21] replicate the experiment in a larger corpus.They report that the F-score of the resolution algorithm

without gestural features is 0.624, while it becomes 0.67 when information about deictic gestures is added to the system.

Differing from these studies, we investigate gestures that are semantically related to anaphora and co-referring expressions in a corpus of Danish video-recorded spontaneous two-party and three-party interactions between well-acquainted people. The recordings have been transcribed and the gestures have been annotated manually. In this study, we present a first analysis of the gestures related to anaphoric expressions. We also describe the results of machine learning experiments which we performed on the annotated data in order to investigate whether the gestures accompanying anaphoric expressions and their antecedents are related and whether these gestures can be automatically recognised on the basis of attributes describing their shape[1].

Investigating the relations between gestures and the referring expressions to which they are linked is not only important for understanding human communication, but it can also contribute to the resolution of anaphora in multimodal interactions and to the generation of plausible conversational agents.

The article is organised as follows. In section 2, we present the types of gestures that are related to anaphora, and we shortly introduce some relevant studies on these gestures. In section 3, we describe our corpus and the multimodal annotations, then we present a first analysis of the annotated data ( section 4). In section 5, we account for the machine learning experiments conducted on the data, and finally, we conclude and discuss future work (section 6).

## 2  Gestures and Anaphora

Various classifications of gesture types have been proposed, inter alia [5, 9, 23]. We use a classification based on Peirce [28] who recognises three main semiotic types: *symbolic*, *iconic*, and *indexical*.

*Indexical gestures* have a real and direct connection with the objects they denote and comprise *deictic* (pointing) and *non-deictic* gestures, e.g. beats and displays. *Iconic gestures*, also known as *emblems* or *illustrators*, denote their objects by similarity. They include metaphoric gestures. Finally, *symbolic gestures* are established by means of an arbitrary conventional relation.

The gestures that are related to anaphora are deictics that point towards objects which are not physically in the interaction room and iconic gestures These gestures often co-occur with speech, but they can also occur alone.

Pointing gestures have been extensively investigated from several points of view, inter alia in intercultural communication studies and in cognitive as well as in language acquisition studies [15, 18, 20]. In western European cultures, pointing is mainly done with hands, but can also involve the head, the body, and the gaze. In other cultures, pointing mainly involves other body parts, e.g. the mouth [18]. Furthermore, Kendon [15] presents differences in the way Italians point to individual objects (tokens) and to object types.

---

The main communicative function of pointing is to indicate specific objects in the interaction's physical space, but it can also be related to other communicative functions such as turn management, feedback, and focusing (information structure). As noticed above, gestures can point towards objects that are evoked in speech but are not present in the interaction room. These objects can also be abstract [15].

Iconic gestures resemble in their execution and manner of performance a concrete object, event, or action [23]. Some of the functions of iconic gestures which have been studied comprise their use to help recovering words, to facilitate the addressee's comprehension of some concepts, and to distinguish between more objects. Studies of iconic gestures also include aspects such as the relation between the represented objects and the way gestures depict them [31, 17], the analysis of the gestures with respect to the language type, [10], the creation of a gesture lexicon, i.a. [29, 16].

Symbolic gestures, which in some cases can be related to anaphora, have mainly been investigated in intercultural and cognitive studies, see inter alia [11, 15].

## 3 The Data

Our data consist of four video-recordings of spontaneous interactions between Danish native speakers (approx. 15 minutes each). The participants are sitting around a small table in private houses and discussing various subjects, such as economic crisis and family relations. The videos have been collected and CA transcribed with the CLAN tool[22] by researchers at University of Southern Denmark as part of the MOVIN database. A typical conversational settings is in Fig. 1.

We have re-transcribed the conversations orthographically in PRAAT [2] assigning time stamps to each word. Then, we have imported the PRAAT and the CLAN transcriptions into the ANVIL multimodal annotation tool [16]. The annotation is done following an extension of the MUMIN annotation scheme [1] which provides pre-defined attribute-value pairs describing the shape and the communicative functions of gestures[2]. The description of the gestures' shape is quite coarse-grained.

The MUMIN scheme has been applied to annotate video-recordings in several languages, comprising Danish [26, 27] Estonian [14], Finnish [12], Greek[19], and Swedish[1]. Some of these annotations have been evaluated in terms of inter-coder agreement with acceptable results given the type of task (Cohen's kappa [3] between 50-90% depending on the categories) [13, 26, 24].

Gestures are often multifunctional, and they can be related to one or more words if the annotators find that the gestures are semantically related to these

---

[2] Only gestures which are judged to have a communicative function are coded.

**Fig. 1.** An Interaction Setting

words. The functions annotated in this corpus are related to feedback, turn management, sequencing and information extraction. An emotion or attitude attribute can also be assigned to gestures. Examples of values connected to the emotion/attitude attribute are the following: *happy*, *sad*, *satisfied*, *certain*, *uncertain*, and *nervous*.

Gestures are divided into Peirce's three semiotic types [28] as described in section 2. In the present work, we have distinguished various kinds of deictics, indicating the object type to which gestures "point". The recognized deictic subtypes are the following: deictic-1person (the deictic points to the speaker), deictic-2person (the deictic points to the interlocutor), deictic-3person (the deictic points to an object in the interaction's room), deictic-3person-no (the deictic points to an individual object that is not in the interaction room), deictic-3person-abstract (the deictic points to an object that is not in the room, and it is related to an abstract anaphor[3]), deictic-space (the deictic is linked to a space adjunct), deictic-time (the deictic is linked to a time adjunct).

The shape of Head Movements is described with the name of the movement and information about whether the movement is single or repeated (Table 1). The features describing the shape of facial expressions and body postures are mostly as proposed in MUMIN. The description of the shape of hand gestures is a simplification of the scheme used at the McNeill Lab [4] and consists of eight dimensions, comprising the trajectory and amplitude of the gesture, the

---

[3] These anaphora have as antecedents constructions such as verbal phrases and discourse segments.

5

**Table 1.** Features for Head Movement

| Behaviour attribute | Behaviour value |
| --- | --- |
| Head_Movement | Nod, Jerk, HeadBackward HeadForward, TiltRight, TiltLeft, SideTurnRight, SideTurnLeft, Shake, Waggle, HeadOther |
| Head_Repetition | Single, Repeated |
| FaceInterlocutor | ToInterlocutor, AwayFromInterlocutor |
| GazeDirection | Up, Down, Forward, Left, Right, GazeDirectionOther |
| GazeInterlocutor | ToInterlocutor, AwayFromInterlocutor |

orientation of the palm, the extension of the fingers. The attributes and values describing hand gestures are in Table 2. A print-screen of the ANVIL tool with

**Table 2.** Features for Hand Gestures

| Behaviour attribute | Behaviour value |
| --- | --- |
| Handedness | SingleHand, BothHands |
| Hand-Repetition | Single, Repeated |
| Palm | Open, Close, PalmOther |
| PalmOrientation | Up, Down, Side, Vertical, OrientationOther |
| Fingers | IndexExtended, ThumbExtended, AllFingersExtended,FingersOther |
| Amplitude | Centre, Periphery,AmplitudeOther |
| TrajectoryRightHand | Forward, Backward, |
| or | Up, Down, SideRight, SideLeft, |
| TrajectoryLeftHand | HandComplex, HandOther |

the annotations of one of the interactions from the MOVIN database is in Fig 2.

## 4  Analysis

The annotations of the video recordings used in this study comprise 2619 gestures. The distribution of the annotations in the various gesture types are in Table 3. 110 of the 470 hand gestures involve both hands, and 282 of all gestures are repeated. 168 of the gestures have been classified as deictics (31 head movements, 31 gaze, and 106 hand gestures) and 61 hand gestures have been classified as iconic. In the following we only focus on head movements and hand gestures, because all gaze movements in these data co-occur with deictic head movements in the same direction.

127 of the deictic gestures (87 hand gestures and 30 head movements) are linked semantically to referring nominal expressions, and only five hand gestures
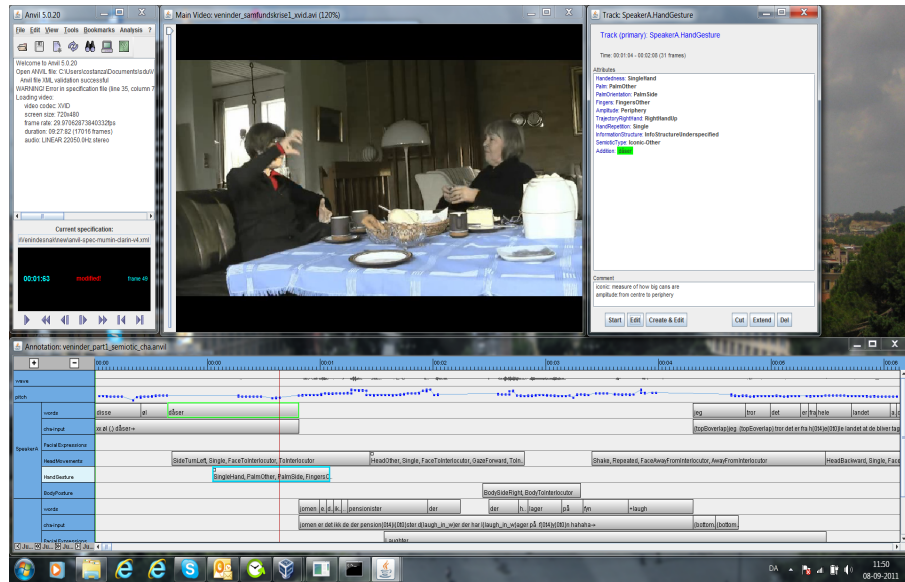
**Fig. 2.** Print-screen from the Anvil Tool

**Table 3.** Annotated Gestures

| Gesture | Total |
|---|---|
| Head Movements | 1069 |
| Gaze Direction | 868 |
| Hand Gesture | 470 |
| Facial Expressions | 98 |
| Eyebrows | 6 |
| Eyes | 11 |
| Mouth | 8 |
| Lips | 7 |
| Body Posture | 82 |
| Total | 2619 |

are linked to abstract anaphora. Two gestures are classified as first-person-deictic (the speaker points to herself) and 21 are coded as second-person-deictics (the speaker points to an interlocutor). Only in one fifth of the cases gestures co-occurred with anaphors and antecedents.

In 14 cases two or more referring expressions and their antecedents are accompanied by a gesture.

Only four iconic gestures are relevant to our study. In all these occurrences, the iconic gestures accompanying the anaphor or the co-referring expressions and their antecedents have similar shape attributes and values, although in one case there are 15 utterances (comprising expressions such as *Yes* and *Ok*) between the anaphor and the antecedent.

In all relevant cases (both deictic and iconic gestures) the gestures are performed by the same interlocutor. In these data there are no symbolic gestures connected to anaphora.

A first analysis of the data shows that if both anaphora and antecedents are related to gestures, the shape descriptions of the gestures have many common attributes and values. When only anaphora are linked to gestures, the hands and/or head of the speaker point towards a place in the room where there are no objects of the referent's type. Exceptions are cases in which the speaker talks about persons or objects somehow related to an interlocutor and points toward the interlocutor, or in which the deictic linked to the anaphor points to the place where the speaker previously made an iconic gesture to illustrate the antecedent.

## 5  Machine Learning Experiments

Although the data-set of gestures relevant to our study is limoted, we run machine learning experiments on the data in order to investigate whether the various attributes describing the gestures related to the anaphora and those related to the antecedents are the same and to which extent gestures related to referring expressions can be distinguished on the basis of their form. The experiments are run in WEKA [32].

In the first experiment we run a clustering algorithm (Expectation Maximisation) on the shape annotations of the relevant hand gestures. The aim of the experiment is to test the hypothesis that gestures which co-occur with expressions that are linked by reference relations (identity relation and other relations) should be clustered together because they have similar shape attributes. We only included a restricted number of the annotated shape attributes in the experiment, that is Palm, Palm Orientation, Fingers, Hand Repetition, and Trajectory of Right and Left Hand. We obtained eight clusters. Nearly all gestures which are connected with referring expressions and their antecedents in our data are grouped in the same cluster. The only two exceptions are cases in which the speaker makes an iconic gesture co-occurring with an action and then points to the place where she made the iconic gesture while referring to that action by the abstract pronoun *det* (it/this/that). Of course the results of this experiment do not imply necessarily that eventually competing antecedents also occur

in different clusters. They only show that when both anaphora (or other co-referring expressions) and their antecedents co-occur with gestures in our data, the gestures' shape attributes and values are similar.

These results are also in line with the results of the co-reference study performed by Eisenstein and Davis [6, 7].

In the second group of experiments we wanted to investigate to which extent the features describing the form of the hand gestures and their various functions help recognising their semiotic type. Thus, we classified the semiotic types of all hand gestures automatically. First, we used the annotated shape features, then we added the gesture's functional attributes (feedback, turn management, sequencing and information structure) to the dataset. The used classifier is WEKA's SMO, a support vector classifier. The baseline are the results obtained with the ZeroR classifier, which always chooses the most frequent nominal category.

Table 4 contains the results of both classifiers evaluated via ten-fold cross validation. These results show that the shape description of hand gestures which

**Table 4.** Results of Classification Experiments

| Algorithm | Dataset | F-score |
|-----------|---------|---------|
| ZeroR | | 28.6 |
| SMO | Shape | 58.3 |
| SMO | Shape+Feedback | 56.5 |
| SMO | Shape+Turn | 61.4 |
| SMO | Shape+Sequencing | 59.2 |
| SMO | Shape+IS | 47.8 |
| SMO | Shape+Functions | 67.1 |

is coded in our corpus helps to classify the semiotic type of hand gestures, even if this description is quite coarse-grained. The results of the classifier improve slightly adding turn management or sequencing, while they get worse adding the feedback or information structure attributes. Finally, the results improve significantly if all the functions of gestures are used. These results can be in part explained by the fact that the various semiotic types of gestures are connected to a limited number of functions and that some of these functions can co-occur. The number of gesture related to each type of function also influences the results. In particular hand gestures are only seldom related to feedback, and information structure nearly always co-occurs with other function types.

Concluding, our experiments confirm that the shape of hand gestures which accompany anaphora (or co-referring expressions) and that of gestures which accompany the anaphora's (or co-referring expressions') antecedents are similar. The shape of hand gestures might therefore contribute to the resolution of anaphora and co-referring expressions.

Furthermore, the experiments indicate that, to some extent, it is possible to identify automatically gestures which are relevant for reference from other types of gestures on the basis of their shape, even if the shape's description is not fine-grained.

## 6   Conclusion and Future Work

This paper contains a study of head movements and hand gestures which co-occur with anaphora and their antecedents in video-recordings of spontaneous dyadic and triadic interactions. Because gestures relevant to anaphora are deictic or iconic, we have focused on these types of gesture. The gestures' shape and dynamics, their function, their semiotic type and their relation to co-speech was manually annotated in the recordings. The analysis of the pairs of gestures which co-occur with anaphoric expressions and with their antecedents in the data indicates that their shape is described with many common attributes and values. The similarity of the shape of gestures accompanying the anaphora and their antecedents is confirmed by machine learning experiments in which the gestures annotated have been clustered on the basis of their shape annotation. As expected, almost all gesture pairs related to anaphora and their antecedents occur in the same clusters. These results are also in line with the study of hand gestures related to co-referring expressions made by Eisenstein and Davis [6, 7].

Classification experiments run on the annotated show that gestures related to referring expressions can, to some extent, be automatically recognised on the basis of their shape description. This is promising given that the shape description annotated in the corpus is not very fine-grained.

This study is only a pilot study and the utility of gestures for anaphora and co-reference resolution must be confirmed on more data and including gestures in resolution experiments as attempted by [8, 21]. Because the manual segmentation and the annotation of the gestures' shape is extremely time-consuming, it should be partially replaced by automatic annotations.

In our machine learning experiments we have only focused on hand gestures because they were the most frequent occurring gesture types relevant to this study. However, all behaviours related to referring expressions, that is hand gestures, head movements, gaze direction and, in some cases also body posture, should be analysed together.

## References

1. Allwood, J., Cerrato, L., Jokinen, K., Navarretta, C. and Paggio, P.: The MUMIN coding scheme for the annotation of feedback, turn management and sequencing phenomena. In J.-C. Martin et al. (Eds.), Multimodal Corpora for Modelling Human Multimodal Behaviour, Special issue of the International Journal of Language Resources and Evaluation, pp. 273-287. Springer, (2007).
2. Boersma, P. and Weenink D.: Praat: doing phonetics by computer (version 5.1.05). Retrieved May 1, 2009, from http://www.praat.org/.

3. Cohen, J. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, (1960).

4. Duncan, S.: McNeill Lab Coding Methods. Technical Report available from http://mcneilllab.uchicago.edu/topics/proc.html. (2004).

5. Efron, D. *Gesture and Environment.* King's Crown Press, N.Y. (1941).

6. Eisenstein, J., Davis, R.: Gesture Features for Coreference Resolution. In Renals,S., Bengio, S., Fiscus, J. (Eds.) *MLMI 2006*, pp. 154–165, (2006).

7. Eisenstein, J. and Davis, R.: Gesture Improve Coreference Resolution. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the ACL*, New York, June, pages 37–40, (2006).

8. Eisenstein, J., Barzilay, R. and Davis, R. Gesture Salience as a Hidden Variable for Coreference Resolution and Keyframe Extraction. In *Journal of Artificial Intelligence Research* 31:353-398 (2008).

9. Ekman, P. and Friesen, W. V. The repertoire of nonverbal behavior: Categories, origins, usage, and encoding. *Semiotica*, 1, 4998, (1969).

10. Goldin-Meadow, S., Chee So, W., Ozyurek, A., and Mylander, C. The natural order of events: how speakers of different languages represent events nonverbally. Proceedings of the National Academy of Sciences of the USA, 105(27), 9163-9168. (2008).

11. Goodwyn, S.W, and Acredolo, L.P and Brown, C.A. I.mpact of Symbolic Gesturing on Early Language Development. Journal of Nonverbal Behavior, 24, 81-103 (2000).

12. Jokinen, K. and Ragni, A. Clustering experiments on the communicative properties of gaze and gestures. In *Proceeding of the 3rd. Baltic Conference on Human Language Technologies*, Kaunas, October (2007).

13. Jokinen, K., Navarretta, C. and Paggio, P.: Distinguishing the communicative functions of gestures. In Proceedings of the 5th Joint Workshop on Machine Learning and Multimodal Interaction.Springer Verlag (2008).

14. Jokinen, K. and Vanhasalo, M. Stand-up Gestures  Annotation for Communication Management. In Navarretta et al. (Eds.) *Proceedings of the NODALIDA 2009 Workshop Multimodal Communication: from Human Behaviour to Computational Models*. Odense, Denmark, May, pp. 15-20 (2009).

15. Kendon, A.: *Gesture: Visible Action as Utterance.* Cambridge: Cambridge University Press (2004).

16. Kipp, M.: *Gesture Generation by Imitation - From Human Behavior to Computer Character Animation.* Ph.D. thesis, Saarland University, Saarbruecken, Germany, Boca Raton, Florida, dissertation.com. (2004).

17. Kita, S. How representational gestures help speaking. In McNeill D (ed.) *Language and Gesture.* Cambridge University Press, New York, pp. 379415.

18. Kita, S. (Ed.): Pointing: where language, culture, and cognition meet. Mahwah, NJ: Lawrence Erlbaum, (2002).

19. . Koutsombogera, M. and Touribaba L. and Papageorgiou, H. Multimodality in Conversation Analysis: A Case of Greek TV Interviews. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC 2008) Workshop on Multimodal Corpora from Models of Natural Interaction to Systems and Applications*, Marrakesh, May, pp. 12-15 (2008).

20. Liebal, K., Behne, T., Carpenter, M., and Tomasello, M.: Infants use shared experience to interpret pointing gestures. *Developmental Science*, 12, 264-71, (2009).

21. Chen, L., Wang, A. and Di Eugenio, B. Improving Pronominal and Deictic Co-Reference Resolution with Multi-Modal Features. In *Proceedings of the SIGDIAL 2011: the 12th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 307311, Portland, Oregon, June 17-18, (2011).

22. MacWhinney B.: The CHILDES Project: Tools for Analyzing Talk. Mahwah, NJ: Lawrence Erlbaum Associates (2000).

23. McNeill, D.: *Hand and mind: What gestures reveal about thought.* Chicago: University of Chicago Press. (2000).

24. Navarretta, C., Ahlsn,E., Allwood, J., Jokinen, K., Paggio, P: Creating Comparable Multimodal Corpora for Nordic Languages. In *Proceedings of the 18th Nordic Conference of Computational Linguistics (Nodalida 2011)*, Riga, Latvia, May 11-13, 2011, pp. 153–160 (2011).

25. Navarretta, C. Annotating Nonverbal Behaviours in Informal Interactions. In A. Esposito et al. (Eds.) *Analysis of Verbal and Nonverbal Communication and Enactment: The Processing Issues*, Springer, LNCS 6800, 317-324 (2011).

26. Navarretta, C. and Paggio, P.: Classification of Feedback Expressions in Multimodal Data. Proceedings of ACL 2010 Uppsala, Sweden, pp. 318-324 (2010).

27. P. Paggio and C. Navarretta. Head Movements, Facial Expressions and Feedback in Danish First Encounters Interactions: A Culture-Specific Analysis. In Constantine Stephanidis (Ed.) Universal Access in Human-Computer Interaction- Users Diversity. 6th International Conference. UAHCI 2011, Held as Part of HCI International 2011 Orlando Florida,July 9-14 2011, LNCS 6766, Springer Verlag, pp. 583-590 (2011).

28. Peirce, C. S.: Collected Papers of Charles Sanders Peirce, 1931-1958, 8 vols. Edited by C. Hartshorne, P. Weiss and A. Burks. Cambridge, MA: Harvard University Press.

29. Poggi, I.: The Lexicon and the Alphabet of Gesture, Gaze, and Touch. In A. de Antonio and R. Aylett and D. Ballin (Eds.) *Intelligent Virtual Agents* Third International Workshop, IVA 2001 Madrid, Spain, September 1011, 2001 Proceedings, Lecture Notes in Computer Science Volume 2190, 235–236 (2001).

30. Poggi, I.: *Mind, Hands, Face and Body - A Goal and Belief View of Multimodal Communication.* Weidler Buchverlag Berlin (2007).

31. Poggi, I.: Iconicity in different types of gestures. In *Gesture* 8:1, 45–61, John Benjamins Publishing Company (2008).

32. Witten, I.H.; Frank, E.: *Data Mining: Practical machine learning tools and techniques.* Morgan Kaufmann, San Francisco, second edition. (2005)