

# On the Annotation of Multimodal Behavior and Computation of Cooperation Between Modalities

Jean-Claude MARTIN

L.I.M.S.I.-C.N.R.S., BP 133, 91403  
Orsay cedex, FRANCE

LINC, IUT de Montreuil, 140 rue de la  
Nouvelle France, 93100 Montreuil,  
FRANCE

+33.1.69.85.81.04

[martin@limsi.fr](mailto:martin@limsi.fr)

Sarah GRIMARD

L.I.M.S.I.-C.N.R.S., BP 133, 91403  
Orsay cedex, FRANCE

+33.1.69.85.81.04

Katerina ALEXANDRI

L.I.M.S.I.-C.N.R.S., BP 133, 91403  
Orsay cedex, FRANCE

+33.1.69.85.81.04

## ABSTRACT

With the success of multimedia and mobile devices, human-computer interfaces combining several communication modalities such as speech and gesture may lead to more "natural" human-computer interaction. Yet, developing multimodal interfaces requires an understanding (and thus the observation and analysis) of human multimodal behavior. In the field of annotation of multimodal corpus, there is no standardized coding scheme. In this paper, we describe a coding scheme we have developed. We give examples on how we applied it to a multimodal corpus by producing descriptions. We also provide details about the software we have developed for parsing such descriptions and for computing metrics measuring the cooperation between modalities. Although this paper is concerned with the input side (human towards machine) and thus deals with the annotation of human behavior observed in multimodal corpora, we also provide some ideas on how it might be of use for specifying cooperation between output modalities in multimodal agents.

## Categories and Subject Descriptors

[**multimedia tools, end-systems and application**]: multi-modal interaction and integration, coding of multi-modal video corpus.

## General Terms

Design, Experimentation, Human Factors, Standardization.

## Keywords

Multi-modal coding scheme.

## 1. INTRODUCTION

With the success of multimedia and mobile devices, human-computer interfaces combining several communication modalities such as speech and gesture may lead to more "natural" human-computer interaction.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

*Autonomous Agent '01*, May 29, 2001, Montreal, Canada.  
Copyright 2000 ACM 1-58113-000-0/00/0000...\$5.00.

Yet, developing multimodal interfaces requires an understanding (and thus the observation and analysis) of human multimodal behavior. To develop multimodal human-computer interfaces, researchers have been producing and analyzing corpora of multimodal behavior made of video tapes and electronic or textual annotations [7, 8, 11, 9, 12]. A survey of such experiments can be found in [6]. Some corpora of multimodal behavior have also been built in other domains such as Sociology, for example the analysis of interactions during a meeting [10].

During multimodal analysis, researchers in this field use their own coding scheme for annotating multimodal behavior and for computing the metrics they are interested in for measuring multimodal behavior (i.e. the temporal relationships between speech and gesture). Since no common standard is used, corpora of multimodal behavior can not be shared (i.e. one researcher can not compute the statistics she is interested in using a corpus constructed by someone else). Targeting a standard (or at least guidelines) for the encoding of multimodal behavior could make easier the sharing of multimodal corpora on a large scale. This could also be fruitful to future developers and providers of multimedia search engines where annotation of relationships between media is a key issue.

The work described in this paper has been done as part of the ISLE European project [3]. We present a coding scheme for multimodal corpus that we have implemented with a Document Type Definition (DTD). We give examples on how we applied it to a multimodal corpus by producing Extensible Markup Language (XML) descriptions. We have also provide details about Java software we have developed for parsing such descriptions and for computing metrics measuring cooperation among modalities.

## 2. THEORETICAL FRAMEWORK FOR MULTIMODALITY

TYCOON is a framework for observing, evaluating and specifying cooperation among modalities during multimodal human-computer interaction. TYCOON stands for Types of COOperation [5]. This typology is composed of six basic types of cooperation: equivalence, specialization, transfer, redundancy, complementarity, concurrency.

**Modality.** In TYCOON, a *modality* is seen as a process for analyzing and producing *chunks of information*. We consider both the modalities used by humans (speech, hand gesture,...) and the modalities used by computers (graphics, speech synthesis,...).

**Referenceable object.** A *referenceable object* embeds an object of the application with knowledge on how to refer to this object (with linguistic or non-linguistic means). In a classical map application, such objects are hotels, restaurants or streets the user is able to refer to in commands.

**Saliency value.** In TYCOON, the reference resolution process is based on the computation of *saliency values* [2]. The saliency value of a referenceable object in a modality gives an idea of how much this object is explicitly referred to in this modality. A global saliency value is computed across several modalities to find the best candidate for the reference resolution. In case of ambiguity, two referenceable objects may have the same saliency in one modality (e.g., graphics). Yet, this ambiguity might be removed when considering the saliency of these objects in another modality (e.g., gesture). The saliency value associated to an object within a reference can take any value between 0 and 1.

**Equivalence.** A cooperation by *equivalence* is defined by a set of modalities, a set of chunks of information, which can be produced by either of the modalities and a criterion, which is used to select one of the modalities. When several modalities cooperate by equivalence, this means that a chunk of information may be produced as an alternative, by either of them.

**Redundancy.** Several modalities, a set of chunks of information and three functions define cooperation by *redundancy*. The first function checks that there are some common attributes in chunks produced by the modalities, the second function computes a new chunk out of them, and the third function is used as a fusion criterion. If modalities cooperate by redundancy, this means that these modalities produce the same information.

**Complementarity.** Cooperation by *complementarity* is similar to cooperation by redundancy except that there are several non-common attributes between the chunks produced by the two processes. The common value of some attributes might be used to drive the fusion process. When modalities cooperate by complementarity, different chunks of information are produced by each modality and have to be merged.

**Specialization.** Cooperation by *specialization* is defined by a modality, a set of modalities *A* and a set of chunks of information this modality is specialized in when compared to the modalities of the set *A*. When modalities cooperate by specialization, this means that a specific kind of information is always produced by the same modality.

**Transfer.** Cooperation by *transfer* is defined by two modalities and a function mapping the output of the first modality into the input of the second modality. When several modalities cooperate by transfer, this means that a chunk of information produced by one modality is used by another modality.

**Concurrency.** Cooperation by concurrency means that several modalities produce independent chunks of information at the same time. These chunks must not be merged.

The TYCOON framework has already been applied to the analysis of the multimodal behavior of subjects in a Wizard of Oz experiment at the Stanford Research Institute [1, 6]. During this experiment, subjects were asked to interact with a simulated system using speech and pen to get touristic information about Toronto. Sessions were videotaped. During the analysis of the video corpus, saliency of the reference to objects was computed manually. Some results of the computation of such statistics have

been described in [4]. In this previous work no structured coding scheme for annotating multimodal behavior had been proposed. Videos were transcribed manually and statistics were computed manually from the transcriptions. On a wide scale this methodology is not possible as it requires extra work and may lead to errors during the manual computation of statistics. Furthermore, as more and more multimedia resources become available on the Internet, standardized formats for corpus annotation may help in achieving a better usage and exchange of corpus data.

### 3. CODING SCHEME FOR THE ANNOTATION OF MULTIMODAL BEHAVIOR

Our goal is to ease the computation of metrics of multimodal behavior from video corpora. The metrics we are interested in are based on the theoretical framework described in the previous section.

Thus, the coding scheme we propose features the annotation of available referenceable objects and the annotation of references to such objects in each modality. Pieces of information enabling the computation of saliency values associated to referred objects are also included in the coding scheme.

The logical structure of the coding scheme we propose is the following one:

- A corpus of multimodal behavior is annotated as a multimodal session
- A multimodal session includes one referenceable objects section and one or more multimodal segments
- A *multimodal segment* is made of a speech segment, a gesture segment, the annotation of temporal relation between these two segments and a graphics segment

We have implemented this coding scheme as a Document Type Definition (DTD) for defining the generic structure of multimodal behavior annotations. The DTD is provided in appendix. Such annotations are done in the eXtensible Markup Language (XML). We will take the example of the XML annotation of a sample multimodal command observed in the SRI corpus [1]. Such an annotation is composed of a *ReferenceableObjects* section describing the graphical objects the user is able to refer to, and a *MultimodalSegment* section composed of four sub-sections: speech, gesture, synchrony, and graphics (Figure 1).

The first section contains annotation about the referenceable objects the user may refer to such as restaurants, hotels (Figure 2). This section about the referenceable objects is followed by one or several multimodal segment sections.

Each multimodal segment section may contain annotations about speech, gesture, synchrony or the state of the graphics modality. Both speech and gesture annotations may contain annotation of references to objects (Figure 3 and 4). Details about the annotation of synchrony and graphics can be found in the DTD provided in appendix.



Figure 1: Example of the XML annotation of a sample command observed in the SRI corpus [1].



Figure 2: The “referenceable objects” section of a multimodal annotation.



Figure 3: A speech segment (“Senator dinner ... ? can I eat a hamburger there ?” which contains two references.

## 4. COMPUTING METRICS MEASURING MULTIMODAL BEHAVIOR

In this section, we describe the algorithm for computing metrics from the XML annotation files and we provide some examples of the execution of the Java software we have developed using the SUN Java optional package for XML parsing (JAXP).

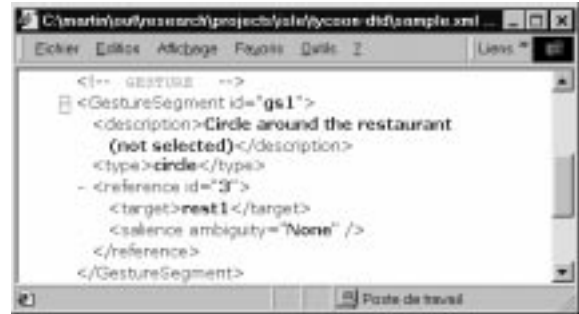


Figure 4: A gesture segment including a reference to the object *rest1*.

### 4.1 STEP 1: Parse the XML file

- Build the document tree out of the XML file.
- Build Java representation of referenceable objects (Figure 5) and references (Figure 6).
- Build the table associating each couple (objects, reference) with a salience value (Figure 7) ; these values are computed according to pre-defined salience rules such as “if the reference contains the full name of this object, set the salience of this object in this reference to 1.0” ; these rules are expected to be dependent on the corpus at hand.
- Build the table computing the average salience values for all the references in the different modalities within the same multimodal segment (Figure 8).

### 4.2 STEP 2: Compute statistics

- Number of referenceable objects, of multimodal segments, of references in each modality.
- Number of multimodal references (i.e. segments including several references in different modalities).
- Average salience of references to an object (or to a type of object) across all modalities with different weight assigned to different modalities (values of weight are fixed a priori).
- Average salience of references in a modality across all objects (or type of object).

id	name	type
dir	a direction	direction
anyd	any direction	direction
hot1	Novotel	hotel
rest1	Senator dinner	restaurant
rest2	Atlas	restaurant
rest3	Tiger Lily	restaurant
site1	Campbell House	site
site2	Eaton Centre	site
site3	Royal Ontario Museum	museum
win1	window	window

Figure 5: Referenceable objects extracted from the XML file and displayed by the Java program.

id	multi...	modality	target	content	ambiguity
11	ms6	Gesture...	site2		None
12	ms7	Speech...	site3	here	
13	ms7	Gesture...	site3		None
14	ms8	Speech...	site3	the museum	
15	ms9	Speech...	site3	it	
16	ms10	Speech...	anyd	there	
17	ms10	Gesture...	adir		None
18	ms11	Speech...	rest1	here	
19	ms11	Gesture...	rest1		LittleAmbiguity
20	ms12	Speech...	rest1	this place	
21	ms12	Gesture...	rest1		None
22	ms13	Speech...	rest3	a Chinese rest...	
23	ms14	Speech...	rest3	this	
24	ms14	Gesture...	rest3		None
25	ms15	Speech...	win1	the window	
26	ms16	Speech...	rest3	there	

Figure 6: Annotated references extracted from the XML file and displayed by the Java program.

### 4.3 STEP 3: Compute metrics

Thanks to these statistics, two metrics are computed for measuring the multimodal behavior. The first metric is the rate of redundancy / complementarity of user's multimodal behavior which is computed as the average salience value assigned to objects when they are referred to:

```

For each object o
  For each reference r
    If the target of r is o
      Then sum += salienceOf o in r
Ratered-compl = sum / nbReferences

```

The value of this metric in the example is provided in Table 1.

Table 1: Average complementarity / redundancy rate of the reference to each object as well as the average across the whole session (0.64).

Object	Object type	Complementarity / Redundancy rate
Atlas	Rest	1
Novotel	Hotel	1
A direction	Dir	1
Eaton Centre	Site	0.7
Senator dinner	Rest	0.6
Tiger Lily	Rest	0.6
Royal Ontario	Site	0.6
Window	Misc	0.6
any direction	Dir	0.4
Rate red-compl		<b>0.64</b>

The second metric measuring the multimodal behavior is the equivalence rate computed in the following way when considering the references to a given object  $o$ :

$$x = \frac{nbRefInSpeech - nbRefInGesture}{MAX}$$

$MAX$  is the total number of references to the object  $o$

```

Rateeq: R+ => [0, 1]
If 0 <= x <= MAX
Then Rateeq(x) = 1 - x / MAX

If MAX < x
Then Rateeq(x) = 0

```

This is the average equivalence rate for one object, the average can be computed for all objects giving an idea of how much the subject switches between different modalities.

## 5. DISCUSSION

We have proposed some preliminary steps towards the proposal of a standard for sharing multimodal analyses, as well as automated tools that gather statistics from corpora expressed in this format. Although multimedia annotation schemes were also considered, XML was selected because of its adequacy to our goals such as the possibility to compute statistics on annotations thanks to a Java API enabling the parsing of XML data. Yet, the current version of the TYCOON-DTD is limited considering the annotation of each single modality such as speech, gesture, body and face, when compared to other systems such as SignStream. Some similar systems include a query language for searching annotated corpora. We will investigate the possibility to define ourselves a multimodal query language based on TYCOON operators. We will also investigate the possibility to store results of multimodal statistics such as equivalence, redundancy and complementarity in the metadata section associated to a corpus for making easier the searching of multimodal patterns in multimedia databases. We are also looking at XML schemas which seem to be more interesting than DTDs when considering expressiveness.

The TYCOON-DTD has already been applied to the annotation of 40 multimodal segments coming from 5 different corpora. We were interested in testing the applicability of our DTD on several corpora. In the near future, we plan to apply the multimodal DTD to more examples to evaluate the tycoon metrics on a significant number of annotations. We believe that the metrics we propose will enable the computation of multimodal behavior features as a function of objects type. When multimodal behavior is too low, graphical attributes, which seem to have an impact on multimodal behavior as observed in [7], could be modified to induce a more redundant behavior which would make recognition easier. Thus, the annotation of multimodal behavior could lead to specifications to be used by a multimodal recognition engine. Finally, as multimedia resources become available on the Internet, one needs to have a better understanding of the potential users of multimedia search engines, but also of the technical requirements on coding schemes for annotating multimedia resources.

Although we have worked only on the input side (human towards machine), our work might also be of use for specifying cooperation between output modalities in multimodal agents. The DTD could be used for specifying at an abstract level the multimodal cross-references including the salience of referents. Such a XML description could itself be the output of a module specifying the intended cooperations between modalities (i.e. either a redundant or complementary behavior of the agent).

object	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	
Senator dinner (rest1)	1	0,4	1	0	0	0	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0,5	0	1	0,6	0,4	0	0	0,4	
Atlas (rest2)	0	0,4	0	0	0	1	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,6	0,4	0	0	0,4	
Tiger Lily (rest3)	0	0,4	0	0	0	0	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,6	0,4	1	0	0,4	
Novotel (hot1)	0	0,4	0	0	0	0	1	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,4	0	0	0	0,4	
Campbell House (sit...)	0	0,4	0	1	1	0	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,4	0	0	0	0,4	
Eaton Centre (site2)	0	0,4	0	0	0	0	0	0,4	0	0,4	1	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,4	0	0	0	0,4	
Royal Ontario Museu...	0	0,4	0	0	0	0	0	0,4	0	0,4	0	0,4	1	0,6	0,4	0,4	0	0,4	0	0	0	0,4	0	0	0	0,4	
window (win1)	0	0,4	0	0	0	0	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,4	0	0	0	0,6	0,4
any direction (anyd)	0	0,4	0	0	0	0	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	0	0,4	0	0	0	0,4	0	0	0	0,4	
a direction (adir)	0	0,4	0	0	0	0	0	0,4	0	0,4	0	0,4	0	0	0,4	0,4	1	0,4	0	0	0	0,4	0	0	0	0,4	

Figure 7: A 2D table is used for storing the saliency value computed for each object (line) in each reference (column). Let's consider the saliency value assigned to Senator dinner restaurant in reference #1 (upper left corner of the table) ; according to Figure 3, this reference is a spoken reference containing the full name of this object ("Senator dinner"), hence the saliency value is set to 1 ; reference # 2 is a reference with a deictic ("there"), hence the saliency value of the Senator dinner in reference #2 is lower : 0.4.

object name	ms1	ms2	ms3	ms4	ms5	ms6	ms7	ms8	ms9	ms10	ms11	ms12	ms13	ms14	ms15	ms16
Senator dinner (rest1)	0,79	0	0	0	0,14	0,28	0,28	0	0,28	0,28	0,43	0,3	0,42	0,28	0	0,28
Atlas (rest2)	0,14	0	0,3	0	0,14	0,28	0,28	0	0,28	0,28	0,28	0	0,42	0,28	0	0,28
Tiger Lily (rest3)	0,14	0	0	0	0,14	0,28	0,28	0	0,28	0,28	0,28	0	0,42	0,58	0	0,28
Novotel (hot1)	0,14	0	0	0,7	0,14	0,28	0,28	0	0,28	0,28	0,28	0	0	0,28	0	0,28
Campbell House (site1)	0,14	1	0	0	0,14	0,28	0,28	0	0,28	0,28	0,28	0	0	0,28	0	0,28
Eaton Centre (site2)	0,14	0	0	0	0,14	0,58	0,28	0	0,28	0,28	0,28	0	0	0,28	0	0,28
Royal Ontario Museum...	0,14	0	0	0	0,14	0,28	0,58	0,42	0,28	0,28	0,28	0	0	0,28	0	0,28
window (win1)	0,14	0	0	0	0,14	0,28	0,28	0	0,28	0,28	0,28	0	0	0,28	0,42	0,28
any direction (anyd)	0,14	0	0	0	0,14	0,28	0,28	0	0,28	0,28	0,28	0	0	0,28	0	0,28
a direction (adir)	0,14	0	0	0	0,14	0,28	0,28	0	0,28	0,58	0,28	0	0	0,28	0	0,28

Figure 8: A 2D table is used for storing the average saliency values of all references of the same multimodal segment.

## 6. ACKNOWLEDGEMENTS

Part of this work was financed by a NSF-STIMULATE grant and the IST-ISLE project. The authors wish to thank A. Cheyer, L. Julia, J. Hobbs, A. Kehler for their work during the Wizard of Oz at SRI ; A. Cheyer for his comments on this paper.

## 7. REFERENCES

- [1] Cheyer, A., Julia, L. & Martin, J.C. (1998), A Unified Framework for Constructing Multimodal Experiments and Applications, Proceedings of the Second International Conference on Cooperative Multimodal Communication, Theory and Applications (CMC'98), 28-30 January 1998, Tilburg, pp.63-69, The Netherlands. <http://cwis.kub.nl/~fdl/research/ti/Docs/CMC/>
- [2] Huls, C., Bos, E., Claasen, W. (1995), Automatic Referent Resolution of Deictic and Anaphoric Expressions. Computational Linguistics, 21, 59-79.
- [3] ISLE (2000) International Standard for Language Engineering. European IST project. [http://www.ilc.pi.cnr.it/EAGLES96/isle/ISLE\\_Home\\_Page.htm](http://www.ilc.pi.cnr.it/EAGLES96/isle/ISLE_Home_Page.htm)
- [4] Kehler, A., Martin, J.C., Cheyer, A., Julia, L., Hobbs, J. & Bear, J. (1998), On Representing Saliency and Reference in Multimodal Human-Computer Interaction, Proceedings of the AAAI'98 workshop on Representations for Multimodal Human-Computer Interaction, July 26-27, 1998, Madison, Wisconsin, USA <http://tiger.cs.uwm.edu/~syali/AAAI-98-Workshop/aaai-wrkshp.html>
- [5] Martin, J.C., Veldman, R., & Béroule, D. (1998a), Developing multimodal interfaces: a theoretical framework and guided propagation networks, Multimodal Human-Computer Communication. Bunt, H., Beun, R.J. & Borghuis, T. (Eds.). Lecture notes in Artificial intelligence 1374. Springer, pp.158-187.
- [6] Martin, J. C., Julia, L. & Cheyer, A. (1998b), A Theoretical Framework for Multimodal User Studies. Proceedings of the Second International Conference on Cooperative Multimodal Communication, Theory and Applications (CMC'98), 28-30 January 1998, Tilburg, The Netherlands, pp.104-110, <http://cwis.kub.nl/~fdl/research/ti/Docs/CMC/>
- [7] Oviatt, S., De Angeli, A., Kuhn, K. (1997), Integration and synchronization of input nodes during multimodal human-computer interaction, Proceedings of the workshop "Referring phenomena in a multimedia context and their computational treatment, ACL/EACL'97, July 11th, Madrid, Spain, pp.1-13, <http://www.dfki.uni-sb.de/imedia/workshops/mm-references.html>
- [8] Petrelli, D., De Angeli, Gerbino, W., Cassano, G. (1997) Referring in multimodal systems: The importance of user expertise and system features. Proceedings of the workshop

"Referring phenomena in a multimedia context and their computational treatment", ACL / EACL'97, July 11th, Madrid, Spain, pp.14-19, <http://www.dfki.uni-sb.de/imedia/workshops/mm-references.html>

- [9] Salmon Alt, S., Romary, L., Schaaf, A. (2000), Increasing the genericity of the MATE annotation framework : the case of reference, Proceedings of the Workshop "Meta-Descriptions and Annotation Schemes for Multimodal Language Resources" - Second International Conference on Language Resources and Evaluation (LREC'2000), 29-30 may, Athens, Greece, pp.57-62.
- [10] Sannino, A. (1998), L'accomplissement interlocutoire et intergestuel d'une interaction en situation de travail. Communications interactives dans les groupes de travail. K.

A.Trognon, Collection langage, cognition, interaction - Presses Universitaires de Nancy, pp.123-155.

- [11] Trafton, J.G., Wauchope, K., Stroup, J. (1997) Errors and usability of natural language in a multimodal system, Proceedings of the IJCAI'97 workshop on "Intelligent multimodal systems", Nagoya, Japan, august 24, pp.49-53.
- [12] Villaseñor, L., Massé, A., Pineda, L. A. (2000), A Multimodal Dialogue Contribution Coding Scheme. Proceedings of the Workshop "Meta-Descriptions and Annotation Schemes for Multimodal Language Resources" - Second International Conference on Language Resources and Evaluation (LREC'2000), 29-30 may, Athens, Greece, pp.52-56.

## 8. APPENDIX: Part of the DTD ([www.limsi.fr/Individu/martin/research/projects/isle](http://www.limsi.fr/Individu/martin/research/projects/isle))

```

<!ELEMENT MultimodalSession (Info?, ReferenceableObjects, MultimodalSegment*)>
...
<!-- REFERENCEABLE OBJECTS -->
<!ELEMENT ReferenceableObjects (object)*>
<!ELEMENT object (type, id, name?, position?, address?)>
<!ELEMENT type (#PCDATA)>
<!ELEMENT id (#PCDATA)>
<!ELEMENT name (#PCDATA)>
<!ELEMENT position (#PCDATA)>
<!ELEMENT address (#PCDATA)>
<!-- ***** -->
<!-- MULTIMODAL SEGMENT -->
<!ELEMENT MultimodalSegment (SpeechSegment*, GestureSegment*, HeadSegment?, BodyMvtSegment?,
Synchrony?, GraphicSegment?)>
<!ATTLIST MultimodalSegment id ID #REQUIRED start CDATA #IMPLIED end CDATA #IMPLIED >
<!-- ***** -->
<!-- SPEECH SEGMENT -->
<!ELEMENT SpeechSegment (content, reference*) >
<!ATTLIST SpeechSegment id CDATA #REQUIRED>
<!ELEMENT content (#PCDATA)>
<!ELEMENT reference (target*, salience*)>
<!ATTLIST reference id CDATA #REQUIRED>
<!ELEMENT target (#PCDATA)>
<!ELEMENT salience (#PCDATA)>
<!ATTLIST salience content CDATA #IMPLIED >
<!-- ***** -->
<!-- GESTURE SEGMENT -->
<!ELEMENT GestureSegment (desc, type, reference*, direction?)>
<!ATTLIST GestureSegment id CDATA #REQUIRED>
<!ELEMENT desc (#PCDATA)>
<!ATTLIST salience ambiguity (None | LittleAmbiguity | Ambiguous) #IMPLIED >
<!ELEMENT direction (#PCDATA)>
<!-- Defines the type of hand movement ; terms are taken from "Hand Gestures for HCI" -->
<!ENTITY % description " changing-position | changing-orientation | changing-shape | contact-objects | join-objects |
indirect-manipulation | empty-handed | haptic-exploration ">
<!ATTLIST type hand (%description;) #IMPLIED>
...

```