



ELSEVIER

Available at

www.ElsevierComputerScience.com

POWERED BY SCIENCE @ DIRECT®

Signal Processing ■ (■■■■) ■■■-■■■

SIGNAL  
PROCESSING

www.elsevier.com/locate/sigpro

# Fusion of children's speech and 2D gestures when conversing with 3D characters

Jean-Claude Martin<sup>a,\*</sup>, Stéphanie Buisine<sup>a</sup>, Guillaume Pitel<sup>a</sup>, Niels Ole Bernsen<sup>b</sup><sup>a</sup>Laboratoire d'Informatique pour la Mécanique et les Sciences de l'Ingénieur (LIMSI-CNRS), BP 133, 91403 Orsay Cedex, France<sup>b</sup>Natural Interactive Systems Lab, Campusvej 55, DK 5230 Odense M, Denmark

Received 1 July 2005; received in revised form 5 December 2005; accepted 1 February 2006

## Abstract

Most existing multi-modal prototypes enabling users to combine 2D gestures and speech input are task-oriented. They help adult users solve particular information tasks often in 2D standard Graphical User Interfaces. This paper describes the NICE Andersen system, which aims at demonstrating multi-modal conversation between humans and embodied historical and literary characters. The target users are 10–18 years old children and teenagers. We discuss issues in 2D gesture recognition and interpretation as well as temporal and semantic dimensions of input fusion, ranging from systems and component design through technical evaluation and user evaluation with two different user groups. We observed that recognition and understanding of spoken deictics were quite robust and that spoken deictics were always used in multi-modal input. We identified the causes of the most frequent failures of input fusion and suggest possible improvements for removing these errors. The concluding discussion summarises the knowledge provided by the NICE Andersen system on how children gesture and combine their 2D gestures with speech when conversing with a 3D character, and looks at some of the challenges facing theoretical solutions aimed at supporting unconstrained speech/2D gesture fusion.

© 2006 Elsevier B.V. All rights reserved.

*Keywords:* Multi-modal interface; Design and evaluation; 2D gestures; Children; Conversational agent

## 1. Introduction

Since Bolt's seminal Put-that-there paper which heralded multi-modal interaction [1], several system prototypes have been developed that enable users to interact through combined speech-gesture input. It is widely recognised today that this form of multi-modal input might constitute a highly natural and intuitive multi-modal "compound" which all or

most humans use for many different communicative purposes. However, most of those prototypes are task-oriented, i.e., they help the user to solve particular information tasks in more or less standard Graphical User Interface (GUI) environments. Moreover, the target user group tends to be adults rather than children. This dominant paradigm of GUI-based task-oriented information systems for adults only addresses a fraction of the potentially relevant domains of application for using combined speech and gesture. Outside the paradigm we find, for instance, systems for children, non-task-oriented systems, systems for edutainment

\*Corresponding author. Tel.: +33 684 21 62 05.

E-mail addresses: martin@limsi.fr (J.-C. Martin), buisine@limsi.fr (S. Buisine), pitel@limsi.fr (G. Pitel), nob@nis.sdu.dk (N. Ole Bernsen).

1 and entertainment, and systems for making-friends  
 3 conversation with 3D embodied characters. The  
 5 challenges to combined speech–gesture input tech-  
 7 nologies posed by systems like those, including  
 9 systems, which include all of the extra-paradigm  
 11 properties mentioned, have not been addressed yet  
 13 to any substantial extent. No existing theory can  
 15 provide reliable predictions for questions, such as:  
 17 how do children combine speech and gesture?  
 19 Would they avoid using combined speech and  
 21 gesture if they can convey their communicative  
 intention in a single modality? Is their behaviour  
 dependent upon whether they use their mother  
 tongue or a second language? To what extent would  
 the system have to check for semantic consistency  
 between their speech and the perceptual features of  
 the object(s) they gestured at? How to manage  
 temporal relations between speech input, gesture  
 input and multi-modal output? How do we evaluate  
 the quality of such systems? What do the target  
 users think of them?

23 This paper addresses the questions and issues  
 25 mentioned above in the context of system prototype  
 27 development and evaluation. We discuss issues in  
 29 semantic input fusion of speech and 2D gesture,  
 31 ranging from systems and component design  
 33 through technical evaluation and user evaluation  
 35 to taking a look at the future challenges, which the  
 37 work reported has uncovered in a very concrete  
 39 manner. The work reported was carried out in the  
 41 EU project NICE on Natural Interactive Commu-  
 43 nication for Edutainment 2002–2005 ([www.niceproject.com](http://www.niceproject.com)). The NICE project has developed two  
 45 prototypes of each of two related systems, one for  
 47 conversation with fairytale author Hans Christian  
 49 Andersen and one for playful computer game–style  
 51 interaction with some of his fairytale characters in a  
 fairytale world. As we shall focus on the Andersen  
 system below, we would like to point out here that  
 both systems are the results of extensive European  
 collaboration, as follows. For both systems, Swed-  
 ish computer games company Liquid Media did the  
 graphics rendering; Scansoft, Germany, trained the  
 speech recognisers with children’s speech; and  
 LIMSI-CNRS, France, did the 2D gesture compo-  
 nents and the input fusion. What makes the two  
 systems different is that the Andersen system’s  
 natural language understanding, conversation man-  
 agement, and response generation components were  
 built by NISLab, Denmark, whereas the corre-  
 sponding components for the fairytale world system  
 were built by Telia-Sonera, Sweden.

### 1.1. Goals of the NICE Andersen project

53 The main goal of Andersen system development  
 55 is to demonstrate natural human–system interaction  
 57 for edutainment by developing natural, fun and  
 59 experientially rich communication between humans  
 61 and embodied historical and literary characters. The  
 63 target users are 10–18 years old children and  
 65 teenagers. The primary use setting for the system  
 67 is in museums and other public locations. Here,  
 69 users from many different countries are expected to  
 have English conversation with Andersen for an  
 average duration of, say, 5–20 min. The main goal  
 mentioned above subsumes a number of sub-goals,  
 none of which had been achieved, and some of  
 which had barely been addressed, at the start of  
 NICE, i.e. to:

- 71 • demonstrate domain-oriented spoken conversa-  
 73 tion as opposed to task-oriented spoken dialo-  
 75 gue, the difference being that, in domain-oriented  
 77 systems there are no tasks to be performed  
 79 through user-system interaction. Rather, the user  
 and the system can have free-style, fully mixed-  
 initiative conversation about any topic in one or  
 several semi-open domains of knowledge and  
 discourse;
- 81 • investigate the challenges involved in combining  
 domain-oriented spoken conversation input with  
 2D gesture input;
- 83 • investigate the use of spoken conversation  
 85 technologies for edutainment and entertainment  
 as opposed to their use in standard information  
 applications;
- 87 • demonstrate workable speech recognition for  
 89 children’s speech which is notoriously difficult  
 to recognise with standard speech recognisers  
 trained on adult speech-only;
- 91 • demonstrate spoken computer games, in a novel  
 93 and wider sense of this term, based on a  
 professional computer games platform; and
- 95 • create a system architecture which optimises re-  
 97 use, so that it is easy to replace Andersen by, e.g.,  
 Newton, Ghandi, or the 40-some past US  
 presidents.

99 The challenge of addressing domains of edutainment  
 101 and entertainment rather than information  
 103 systems was, in fact, chosen to make things slightly  
 easier. Our assumption was that users of the former  
 systems would be more tolerant to system error as  
 long as the conversation as a whole would be

perceived as entertaining. Furthermore, the museum context-of-use requirement mentioned earlier would reduce the performance requirements on the system to those needed for 5–20 min of fun and edutaining interaction. Based on the reasoning just outlined, we chose fairytale author Hans Christian Andersen for our embodied conversational agent because of yet another pragmatic consideration. Given the need to train the system’s speech recogniser with large amounts of speech data to be collected in the project, we needed a natural and convenient place to gather this data, such as the Andersen museum in his native city of Odense, Denmark, where partner NISLab is located.

### 1.2. Interacting with Andersen

The user meets Andersen in his study in Copenhagen (Fig. 1) and communicates with him in fully mixed-initiative conversation using spontaneous speech and 2D gesture. Thus, the user can change the topic of conversation, back-channel comments on what Andersen is saying, or point to objects in Andersen’s study at any time, and receive his response when appropriate. 3D animated Andersen communicates through audiovisual speech, gesture, facial expression, body movement and action. The high-level theory of conversation underlying Andersen’s conversational behaviour is derived from analyses of social conversations aimed at making new friends, emphasising common ground, expressive story-telling, rhapsodic topic shifts, balance of interlocutor “expertise” (stories to tell), etc. [2]. When Andersen is alone in his study, he goes about his work, thinking, meandering in



Fig. 1. Andersen gesturing in his study.

locomotion, looking out at the streets of Copenhagen, etc. When the user points at an object in his study, he looks at the object and then looks back at the user before telling a story about the object.

Andersen’s domains of knowledge and discourse are: his works, primarily his fairytales, his life, his physical and personal presence, his study, and his interest in the user, such as to know basic facts about the user and to know which games children like to play nowadays. The user is, of course, likely to notice that Andersen does not know everything about those domains, such as whether his father actually did see Napoleon when joining his army or whether Andersen’s visit to Dickens’ home in England was a pleasant one. The cover story, which Andersen tells his visitors on occasion, is that he is just back and that there is still much he is trying to remember from his past.

Visiting Andersen, the user can not only talk to him, but also gesture towards objects in his study, such as pictures on the wall or his travel bag on the floor, using a touch screen. Andersen encourages his visitors to do so and has stories to tell about those objects. Using a keyboard key, the user can choose between a dozen different virtual camera angles onto Andersen and his study. The user can also control Andersen’s locomotion using the arrow keys and assuming that Andersen is not presently in autonomous locomotion mode.

Some user input has emotional effects on Andersen, such as when they talk about his poor mother, the washerwoman who died early and had her bottle of aquavit to keep her company when washing other people’s clothes in the Odense River. Andersen is friendly by default but he can also turn sad, as illustrated in Fig. 2, angry, such as when a child tries to offend him by asking about his false teeth, or happy, such as when the self-indulgent author gets a chance to talk about how famous he has become.

### 1.3. Related work

The development of the NICE Andersen system relies on several research fields, in particular those of multi-modal input systems, Embodied Conversational Agents, and interactive systems for young users.

Regarding multi-modal input, numerous prototypes have been developed for combining speech and gesture input in, e.g., task-oriented spatial applications [3], crisis management [4], bathroom



Fig. 2. Close-up of a sad Andersen.

design [5], logistic planning [6,7], tourist maps [8,9], real estate [10], graphic design [11] or intelligent rooms [12,13]. Users' multi-modal behaviour was also investigated in order to ground system development on empirical data, e.g. for the temporal parameterisation of input fusion [14].

Some general requirements to multi-modal 2D gesture/speech input systems have been proposed in standardisation efforts [15]. Unification algorithms have been applied successfully to the interpretation of task-based applications [6]. Techniques have been proposed for managing ambiguity in both the speech and the gesture modality when each of them has limited complexity, such as in [16] where different spoken commands can be combined with different gestural commands for, e.g., mutual disambiguation. Different approaches were considered for multi-modal fusion, including early fusion, which integrates signals at the feature level (for example for simultaneously training lip-reading and speech recognition), and late fusion, which merges individual modalities based on temporal and semantic constraints.

One particular characteristic of the NICE Andersen system is that it offers multi-modal interaction with an animated character—a kind of interface also called Embodied Conversational Agent (ECA) [17] or Pedagogical Agent when applied to education [18]. Given the enormous challenges to achieving full human-style natural interactive communication, research on ECAs is a multi-dimensional endeavour, ranging from fine-tuning lip synchronisation details through adding computer vision to ECAs to theoretical papers on social conversation skills and multiple emotions which

ECAs might come to include in the future. So far, the ECA community has put less emphasis on advanced spoken interaction than has been done in the NICE Andersen system and ECA researchers are only now beginning to face the challenges of domain-oriented conversation. Moreover, few ECA researchers have ventured into the complex territory of conversational gesture/speech input fusion.

For these reasons, we know of few ECA research systems that come close to the Andersen system prototype in being a complete demonstrator of interactive spoken computer games for edutainment and entertainment. One of the research systems closest to the Andersen system may be the US Mission Rehearsal system [19]. By contrast with the Andersen system but similar to the NICE fairytale world system (Section 1), the Mission Rehearsal system is a multi-agent one, so that users can speak to several virtual agents. On the other hand, the sophisticated spoken dialogue with the Mission Rehearsal system is more task-oriented than is the conversation with Andersen; does not enable gesture and gesture/speech input; and does not target children. A few other prototypes involve bi-directional multi-modal communication and hence communication with an ECA via multi-modal input. The MAX agent [20] recognises and interprets combinations of speech and gesture, such as deictic and iconic gesture used for pointing, object manipulation, and object description in virtual reality assembly task. Combination of speech and 2D mouse gestures for interacting with a 3D ECA in a navigation task within a virtual theatre is presented in [21]. The CHIMP project had goals similar to NICE, i.e., to enable children to communicate with animated characters using speech and 2D gestures in a gaming application [22]. Similarly, some projects address fusion of users' gestures and speech when interacting with a robot. Combination of natural language and gesture to communicate commands involving directions (e.g., «turn left») and locomotion (e.g., «go over there») with a robot is described in [23]. Interaction with a humanoid robot in a kitchen scenario is described in [24]. Yet, for several of these bidirectional systems, the interaction still remains task-oriented or only addresses rather restricted conversational interaction experimentally evaluated with a children user group. The conversational dimension notably showed that turn-taking was a main issue, requiring proper output for notifying the user that the agent wants to take, keep, or give the turn.

Another domain likely to provide interesting data for the NICE Andersen project is the research on computer systems dedicated to cognitive development and child education. For example, using a simulated ECA system, Oviatt observed convergence between the spoken behaviour of children and the spoken behaviour of an animated character in a pedagogical application [25]. She also showed the differences in children's speech with this agent as compared to their speech with a human adult [26]. The effect of interacting with an agent was also observed in storytelling abilities of five-year-old girls [27]. However, neither gestural nor multi-modal children's behaviour has been studied to any great extent. Read et al. [28] studied handwritten text input from children but, to our knowledge, only (Xiao) analysed children's multi-modal behaviour with ECAs, primarily focusing on temporal integration of speech and pen input. In this context, the evaluation of the NICE Andersen system provides more data on children's interaction with ECAs, as well as a semantic analysis of their multi-modal constructions.

#### 1.4. Plan for the paper

In what follows, Section 2 describes the analytical steps performed prior to the design of gesture input processing as well as the specifications and algorithm of the Gesture Recogniser (GR) and the Gesture Interpreter (GI). Section 3 presents the design of the Input Fusion (IF) module. Technical and user test results on gesture-related conversation are presented in Section 4. Section 5 concludes the paper by taking a broad look at some of the challenges ahead, which have become increasingly

familiar to us in the course of the work presented in this paper. Throughout, we describe the design and evaluation of the 2nd Andersen prototype, which was in part grounded on observations made on the first Andersen prototype in which the speech recognition was simulated by human wizards [29,30].

## 2. Gesture recognition and interpretation

### 2.1. Requirements on gestural and multi-modal input

In view of the richness and complexity of spoken interaction in the Andersen system, we opted for having basic and robust gesture input. Thus, gesture input has the relatively simple generic semantics and pragmatics of getting information about objects in Andersen's study, which can then be combined with the expected, richer semantics of the spoken input. We did not consider strict unification as in the task-based systems described above, as such strict semantic checking did not appear relevant in an edutainment application for children. Furthermore, the graphical on-screen objects were designed so as to avoid possible overlaps between objects in order to facilitate gesture recognition.

Fig. 3 shows the Andersen system's overall architecture, including the modules involved in gestural and multi-modal input processing: GR, GI and IF. The modules communicate via a message broker, which is publicly available from KTH [31]. The broker is a server that routes function calls, results, and error codes between modules, using TCP/IP for communication. Input processing is distributed across two input "chains" which come together in IF. Speech recognition uses

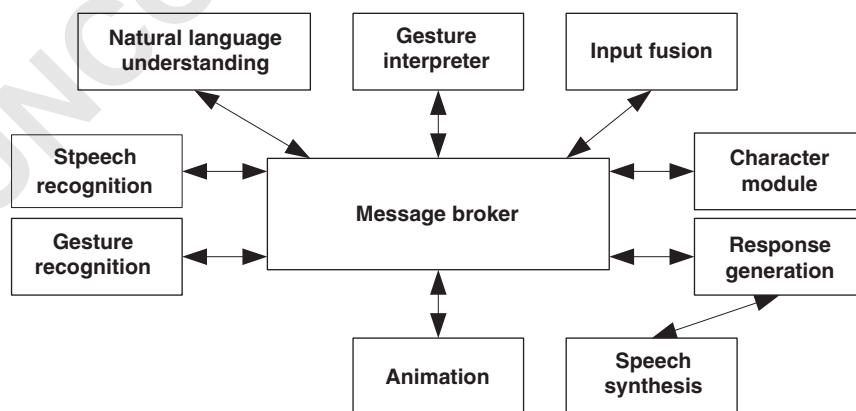


Fig. 3. General NICE Andersen system architecture.

1 a 1977 word vocabulary and a language model  
 2 developed on the basis of three Wizard of Oz  
 3 corpora and two domain-oriented training corpora  
 4 collected in the project. The recogniser's acoustic  
 5 models are tuned to children's voices, using  
 6 approximately 70 h of data most of which has been  
 7 collected in the project. A large part of this data was  
 8 collected in the Odense Andersen museum, using a  
 9 Wizard of Oz-simulated speech-only version of the  
 10 system. The recogniser does not have barge-in  
 11 (constant listening to spoken input) because of the  
 12 potentially noise-filled public use environment. This  
 13 restriction on the naturalness of conversation with  
 14 Andersen was decided upon in order to limit the  
 15 number of speech recognition errors that the system  
 16 would have to deal with. Some effects are that turn-  
 17 taking negotiation becomes curtailed and that the  
 18 user is not able to stop Andersen from completing  
 19 the story he is presently telling. It is also possible  
 20 that the system would miss some backchannelling  
 21 input produced by users while Andersen is speaking.  
 22 Natural language understanding uses the best-  
 23 recognised input string to generate a frame-based  
 24 attribute/value representation of the user's spoken  
 25 input, including dialogue act information. The  
 26 gesture input "chain" is described in detail in the  
 27 following sections.

28 The Andersen character module matches results  
 29 produced by the IF module to potential Andersen  
 30 output in context. Andersen keeps track of what he  
 31 has said already and changes domain when, having  
 32 the initiative, he has nothing more to tell about a  
 33 domain; takes into account certain long-range  
 34 implications of user input; remembers his latest  
 35 output; and keeps track of repeated generic user  
 36 input, including input which requires some form of  
 37 system-initiated meta-communication. The charac-  
 38 ter module's emotion calculator calculates a new  
 39 emotional state for each conversation turn. If the  
 40 input carries information which tends to change  
 41 Andersen's emotional state from its default friendly  
 42 state towards angry (e.g., "You are stupid"), sad  
 43 (e.g., "How was your mom?"), or happy ("Who are  
 44 you?"—I am the famous author Hans Christian  
 45 Andersen...)—the emotion calculator updates his  
 46 emotional state. If the user's input does not carry  
 47 any such information, Andersen's emotional state  
 48 returns stepwise towards default friendly.

49 Design-wise, Andersen is always in one of three  
 50 output states, i.e., non-communicative action when  
 51 he is alone in his study working, communicative  
 function when he pays attention to the user's input,

and communicative action when he actually re- 53  
 sponds to input. In the current system version, these 55  
 three output states are not fully integrated and can 57  
 only be demonstrated in isolation. The exception is 59  
 when the user gestures towards an object in 61  
 Andersen's study, making him turn towards the 63  
 object gestured at and then turn back to face the 65  
 user (the virtual camera). Response generation 67  
 generates a surface language string with animation 69  
 and control (e.g., camera view) tags. The string is 71  
 sent to the speech synthesiser, which synthesises the 73  
 verbal output and helps synchronise speech and 75  
 non-verbal output, including audio-visual speech. 77  
 Speech synthesis is off-the-shelf software from 79  
 AT&T. Andersen's voice was chosen partly for its 81  
 inherent intelligibility and naturalness, and partly 83  
 for matching the voice one would expect from a 85  
 55 years old man. Finally, animation renders Ander-  
 sen's study, animates Andersen, and enables the  
 user to change camera angle and control Andersen's  
 locomotion.

As described in the introduction, the part of the  
 scenario related to the graphical objects displayed in  
 Andersen's study is for the user to "indicate an  
 object to get information about it or express an  
 opinion about it". Table 1 lists the communicative  
 acts identified a priori, which were likely to lead to  
 gestural or multi-modal behaviours. The only  
 generic gesture semantics they feature is the gestural  
 selection of object(s) or location(s). Other possible  
 semantics, such as drawing to add or refer to an  
 object, or crossing an object to remove it, were not  
 considered compatible with the NICE scenario.

A 2D gestural input has several dimensions that  
 need to be considered by the GR/GI/IF modules:  
 shape (e.g., pointing, circle, line) including orienta-  
 tion (e.g., vertical, horizontal, diagonal); points of  
 interest (e.g., two points for a line); number of  
 strokes; location relative to objects; input device  
 (mouse or tactile screen); size (absolute size of  
 bounding box, size of bounding box relative to  
 objects); and timing between sequential gestures.  
 Gesture processing of these dimensions is a multi-  
 level process involving the GR, GI and IF modules.  
 The GR computes a "low-level" semantics from  
 geometrical features of the gesture without con-  
 sidering the objects in the study. The GI computes  
 a higher-level semantics by considering the list of  
 visible objects and their locations at the time of  
 gesturing as sent by the object tracker from the  
 rendering engine. Thus, the possibility that several  
 objects are selected simultaneously cannot be

1	Table 1	Table 2	53
	List of identified communicative acts	Definition of GR output classes	
3	<hr/>	<hr/>	55
	Communicative acts	GR output class	57
5	1. Ask for clarification on what to do with gesture	Pointer	59
	2. Ask for initial information about the study		
7	3. Select one referenceable object	Surrounder	61
	4. Select one non referenceable object		
9	5. Select several referenceable objects		63
	6. Select an area		
11	7. Explicitly ask information about selected object		65
	8. Negatively select an object (e.g. "I do not want to have information on this one")		67
13	9. Negatively select several objects		69
	10. Confirm the selection		
15	11. Reject the selection	Connect	71
	12. Correct the selection		
17	13. Interrupt Andersen	Unknown	73
	14. Ask Andersen to repeat the information on the currently selected object		
19	15. Ask Andersen to provide more information on the currently selected object		75
	16. Comment on information provided by Andersen		
21	17. Comment on another object than the one currently selected		77
	18. Select another object while referring to the previous one		
23	18. Select another object of the same type than the one currently selected		79
	20. Move an object (user may try to do that although not possible and not explicitly related to the user's communicative intention)		
25	21. Compare objects		81
	22. Thank		
29	<hr/>		83
31		2.2. <i>Gesture recognition</i>	85
33	detected by the GR and has to be detected by the GI. The IF computes a final interpretation of gesture by combining the GI output with the Natural Language Understanding (NLU) output.		
35	In the test of the 1st Andersen prototype, some users made several sequential gestures (e.g., parts of a circle) on the same object, which might be due to the fact that the gesture stroke was not highlighted on the screen (which might be due to insufficient finger pressure on the touch screen or a faulty touch screen setting), that Andersen would not give any feedback, such as gazing at the gestured object, or that their finger simply slipped on the tactile screen. This resulted in duplicated messages sent by the GI and thus to output repetitions by the system. In order to avoid this, we decided to have the GI group several sequential strokes on the same object as a single gesture on this object.	The gestural analysis described above resulted in the set of shapes described in Table 2.	87
37	Other difficulties include the facts that some objects have overlapping bounding boxes some of	As a result of gesture recognition, the GR sends to the GI a «grFrame» including the 1st best gesture shape recognised. The two-stroke "cross" shape is recognised when two crossing lines are drawn. It is recognised by the GI (instead of the GR) in order to avoid confusing the delay between the two strokes of the cross with the delays between different gestures. If the multi-stroke gestures were to be recognised by the GR, the GR would have to delay the sending of recognised lines to the GI as, e.g., the GR would wait for the second line of the cross. This delay would add to the delay in the GI for grouping sequential gestures of any type on the same object. In order to avoid this sum of delays, we decided to have multi-stroke gestures recognised by the GI since, there, the delay is used both for waiting for (1) a possible 2nd stroke of a multi stroke gesture and (2) another single-stroke gesture on the same object.	89
39			91
41			93
43			95
45			97
47			99
49			101
51			103

1 When a gesture is detected by the GR, a  
 2 «startOfGesture» message is sent by the GR to the  
 3 IF before launching shape recognition in order to  
 4 enable appropriate timing behaviour in the IF.  
 5 When the GR is not able to recognise the shape or  
 6 when the user makes noisy gestures, the GI can try  
 7 to recover, considering them as surrounder gestures,  
 8 and hopefully detect any associated object. The goal  
 9 is to reduce the non-detection of gestured objects.  
 10 Indeed, surrounder gestures logged during Proto-  
 11 type-1 evaluation were quite noisy and included  
 12 contours of objects. Another possibility would have  
 13 been to induce the user to gesture properly and not  
 14 to forward unknown shapes to the GI, but that was  
 15 considered inappropriate for a conversational ap-  
 16 plication for children. The GR also sends the  
 17 gesture-bounding box to the GI.

18 The GR uses a back-propagation neural network  
 19 trained with gestural data logged from Prototype-1.  
 20 Training involves several steps: manual labelling of  
 21 logged shapes, training of the neural network, and  
 22 testing and tuning its parameters. The general  
 23 algorithm of the GR is shown below.

#### 25 *Algorithm* GR

```

26 When a gesture is detected:
27   Send a ``startOfGesture`` mes-
28   sage to IF
29   If the bounding box of the gesture
30   is very small (10 × 10)
31     Then set shape = ``pointer``
32   Else
33     Convert the gesture points to a
34     slope features array.
35     Test the feature array with the
36     neural network.
37     set shape = result from the
38     neural network
39     (either ``surrounder`` |
40     ``connect`` | ``unknown``)
41     If the shape is ``connect``
42     Then compute start and end
43     points of the line
44     Build a grFrame for this newly
45     detected gesture
46     Send the grFrame to the GI
  
```

47 *End of Algorithm* GR

#### 53 2.3. Gesture interpretation

54 The GI module aims at detecting the object(s) the  
 55 user gestures at. It has been designed by considering  
 56 the properties of the graphical objects that are  
 57 displayed and which the user is able to refer to. The  
 58 properties are:

- 59 • spatial ambiguities due to objects that have  
 60 overlapping bounding boxes, or objects that are  
 61 in front of larger objects, such as the objects on  
 62 Andersen's desk;
- 63 • the singular/plural affordance of objects, e.g., a  
 64 picture showing a group of people might elicit  
 65 either singular spoken deictics, such as «this  
 66 picture», or plural spoken deictics («these peo-  
 67 ple»);
- 68 • perceptual groups which might elicit multiple-  
 69 object selection with a single gesture, or for which  
 70 a gesture on a single object might have to be  
 71 interpreted as a selection of the whole group,  
 72 such as the group of pictures on the wall [32].

73 Following gesture interpretation, the GI sends a  
 74 «giFrame» to the IF module. This frame includes  
 75 one of the three attributes "select" (a gesture on a  
 76 single object), "reference ambiguity" (several ob-  
 77 jects were gestured at), or "no object" (a gesture was  
 78 done, but no associated referenceable object could  
 79 be detected), as defined in Table 3. Gesture  
 80 recognition confidence scores are not considered  
 81 since a fast answer from the character is preferred  
 82 over an in-depth resolution of ambiguity in order to  
 83 enable fluent conversation. Moreover, due to the  
 84 challenging complexity in recognising children's  
 85 conversational speech, it was preferred to ensure  
 86 robust gesture interpretation by avoiding, as far as  
 87 possible, overlaps between graphical objects. Such  
 88 design choices wrt. to the graphical environment  
 89 enabled us to reach high-accuracy recognition of  
 90 gestured objects during monomodal tests held prior  
 91 to the test involving multi-modal fusion and  
 92 children users. Indeed, as it will be described in  
 93 the section on evaluation, assigning scores to results  
 94 of gesture interpretation would not have addressed  
 95 the problems observed in the management of multi-  
 96 modal behaviour.

97 The conversational context of the Andersen  
 98 system requires management of timing issues at  
 99 several levels (Fig. 4). In order to avoid endless  
 100 buffering of the user's input while Andersen is  
 101 speaking, gesture interpretation is inhibited during



1 Table 3 53  
 Definition of GI output classes

3 GI output semantic class	GR output class	Graphical context	55
5 Select	Pointer Cross Surrounder Connect	Gesture bounding box overlaps with bounding box of only one object.	57
7	Sequential: Pointer Cross Surrounder Connect	On the <i>same</i> object (close in time).	59 61 63
9 referenceAmbiguity	Surrounder Cross Connect	Bounding box of gesture overlaps with the bounding boxes of several objects.	65
11	Sequence of pointers or other shapes than unknown		67
13			69
15 noObject	Any except unknown	GI failed to detect any object although a gesture was made by the user (gesture on empty space; selection of non referenceable objects).	71
17			73
19			75
21			77
23			79
25			81
27			83
29			85
31			87
33			89
35			91
37			93
39			95
41			97
43			99
45			101
47			103

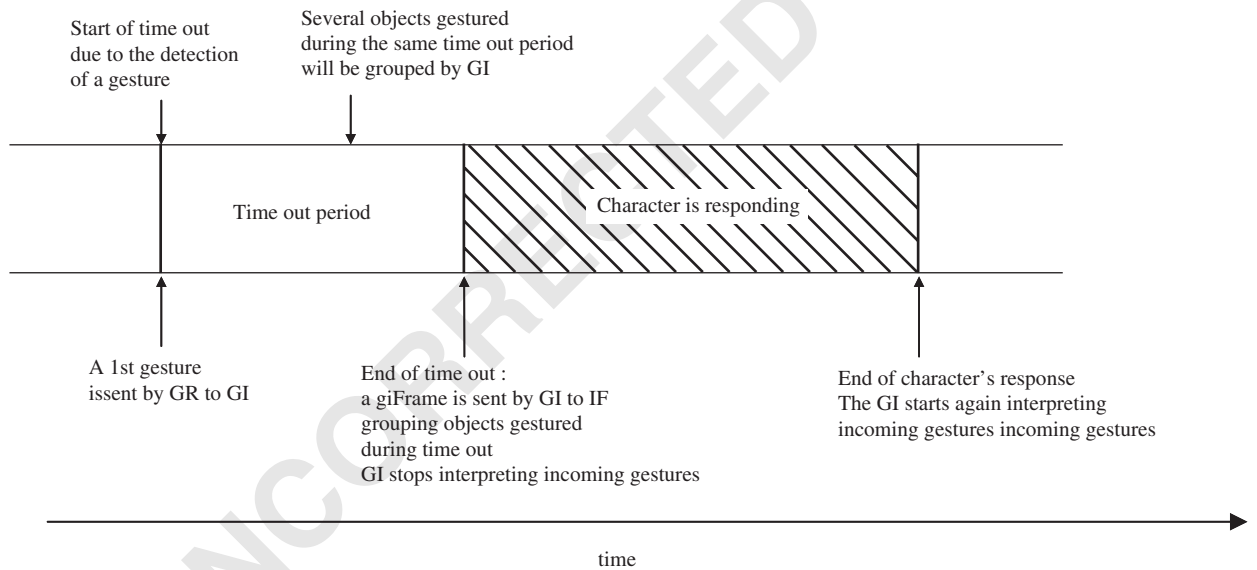


Fig. 4. Temporal management in the GI module.

preparation and execution of Andersen's verbal and non-verbal behaviour. In order to sequentially group objects gestured at, the GI has a relatively fast timeout. It collects what it gets before the timeout and then passes it on to the IF. The message sent by the GI to the IF may include reference to one or several objects. If several objects are

referenced, this may mean either that a single gesture was done on several objects or that sequential gestures were done on different objects. An object does not appear twice in the giFrame even in the case of multiple gestures on the same object. The GI collects references to one or several objects in the given time window and passes them to the IF

1 as a single gesture turn. The timeout period is reset  
each time a new gesture is recognised.

3 The 2nd Andersen prototype requires that once  
the timeout has been started, incoming gestures are  
5 ignored by the GI. The Character Module notifies  
the GI with an «EndOfBehavior» message that  
7 Andersen has finished his verbal and nonverbal  
output turn, so that the GI can start interpreting  
9 gestures again. The same notification is sent to the  
speech recogniser. The GI timeout is analogous to  
11 the lack of barge-in in the speech recogniser.  
However, the GI timeout may be less of a restriction  
13 on the naturalness of conversation since few users  
tend to do 2D touch screen gesture without speak-  
15 ing.

17 The following durations were selected as default  
values for the GI module:

- 19 ● timeout period duration: 1.5 s. This is compatible  
with observations made during the Prototype-1  
21 user tests;
- maximum duration of waiting for the character's  
23 response = 6 s. After this the GI starts interpret-  
ing gestures again.

25 These specifications resulted in the design of the  
27 following algorithm for time management in the GI:

29 **Algorithm** GI

**Input:** incoming messages from GR and CM

**Output:** messages sent by GI to IF

**Variable:** list of object names gestured  
33 during timeout

**{Processing of an incoming grFrame from GR}**

35 **If** a grFrame is received from GR

**Then**

37 **If** the character's response is cur-  
rently pending

**Then**

39 Ignore grFrame

**Else**

41 **If** gesture time out period is not  
43 started

**Then** start gesture time out period  
45 Call bounding box algorithm to  
detect objects

47 Store name of detected object(s)  
in the list of gestured objects (avoid  
49 duplicates)

**{Gesture time out period has finished}**

51 **If** end of timeout period

**Then**

**If** no object was detected during  
timeout 53

**Then** 55

Build a ``noObject`` giFrame

**If** a single object has been detected  
during timeout 57

**Then** 59

Build a ``select`` GIFrame with  
name of this object 61

**If** several objects have been de-  
tected during timeout 63

**Then**

Group objects names in a ``refer-  
enceAmbiguity`` GIFrame 65

Send the GIFrame to IF 67

Set characterResponsePending to  
true 69

**{Character's response is finished}**

**If** message is ``EndOfBehavior`` is  
received from the Character/Dialog  
Module OR 71

message ``EndOfBehavior`` has been  
waited for too long 75

**Then**

Set characterResponsePending to  
false 77

Set gesture detection period not  
started 79

Enable GI to start new timeout if a  
gesture is detected 81

**End of Algorithm** GI 83

85 In 3D graphics, some objects hide others, such as  
87 when a vase is hiding a table. Yet, the graphical  
application only delivers the coordinates of all the  
89 objects, which are partly in the camera viewpoint  
without informing the GI if these objects are hidden  
91 or not by some other visible objects. The objects,  
which are hidden, must not be selectable by gesture,  
93 even if the gesture is spatially relevant. In the  
bounding box algorithm, we used the depth (Z  
95 dimension) of the closest side of the bounding box  
of objects to compute hidden objects. The salience  
97 value computed for each object is weighted by a  
factor of the distance, which is maximal when the  
99 front of the object is near the camera and decreases  
quickly for objects, which are far from the camera.  
101 Yet, an object closer on its Z-dimension can  
actually be partially hidden by one further away,  
103 such as a vase on a table, which hides the part of the  
table, which is behind the vase. Thus, the size of the  
object is also considered in the algorithm. An object,

1 which better fits the size of the gesture is more likely  
to be selected.

3

### 3. Input fusion

5

#### 3.1. Requirements and specifications of input fusion

7

9 IF in the Andersen project aims at integrating  
children's speech and 2D gestures when conversing  
with virtual characters about 3D objects. In  
11 principle, IF is subject to some general requirements  
to multi-modal input systems, such as the need to  
13 manage and represent timestamps of input events,  
multi-level interpretation, composite input, and  
15 confidence scores [15]. Yet, the conversational goal  
of the system and the fact that it aims at being used  
17 by children make it different from current research  
on systems which use speech and gesture for task-  
19 oriented applications as described in the introduc-  
tion.

21 Both speech-only input and gesture-only input  
can be semantically and pragmatically independent.  
23 In other words, using either, the user can input a  
complete communicative intention to the system. As  
25 for combined gesture and speech in an input turn,  
their relationship regarding the semantics of object  
27 selection may be of several different kinds. Thus, the  
input speech may be either (i) redundant relative to  
29 the input gesture as in <pointing at the picture of  
Andersen's mother> "Tell me about your mother,"  
31 (ii) complementary to the input gesture as in <  
pointing at object> "What is this?," (iii) conflicting  
33 with the input gesture as in <pointing at the picture  
of Andersen's mother> "Tell me about your wife,"  
35 or (iv) independent of the input gesture as in <  
pointing at the feather pen> "Do you live here?".

37 Given the formal patterns of relationship between  
speech and gesture input just described, it would  
39 appear that speech-gesture IF is required in the two  
cases of redundancy and complementarity. Con-  
41 versely, IF is excluded in all cases of speech-gesture  
independence, i.e., speech-only input, gesture-only  
43 input, and independent, but simultaneous speech  
and gesture inputs. When independent gesture and  
45 speech occur at the same time, the system should  
not merge them. As for speech/gesture conflict, we  
47 decided to trust the gesture modality, as it is more  
robust than the speech recognition in the context.

49 The IF module integrates the messages sent by the  
NLU and the GI modules and sends the result to  
51 the character module. The IF parses the message  
sent by the NLU to find any explicit object reference

(e.g., "this picture") or implicit reference (e.g., 53  
"Jenny Lind?," "Do you like travelling?") which 55  
might be integrated with gestures on objects in the 57  
study. In order to do so, the IF parses the frame 59  
produced by the NLU and spots the following 61  
concepts: object in study, fairy tale, fairy tale 63  
character, family, work, friends, country, and 65  
location. It produces messages containing the 67  
"fusion status" which can be either "ok," i.e., the 69  
utterance and the gestured object were integrated 71  
because a reference was detected in the NLU 73  
message and in the GI; "none," i.e., the utterance 75  
and the gesture were not integrated either because 77  
there was only one of them, or because the IF could 79  
not decide if they were consistent or not regarding 81  
the number of references to objects in speech and 83  
gesture; or "inconsistent," i.e., the utterance and the 85  
gesture were inconsistent regarding the number of 87  
referenced objects. In case of successful integration, 89  
the semantic representation of gesture (the detected 91  
object(s)) is inserted into the semantic representa- 93  
tion sent by the NLU. The IF module also manages 95  
temporal delays between gesture and speech via 97  
several timeouts and messages signalling start of 99  
speech and start of gesture.

The IF specifications described above were driven 79  
by a conversation analysis that generated a set of 81  
233 multi-modal combinations which users might 83  
produce. This set includes the multi-modal beha- 85  
viours observed during the Prototype-1 user tests.

#### 3.2. Multi-modal behaviours in the Prototype-1 user tests

87 During the Prototype-1 user tests, 2h were 89  
videotaped (about 22% of the tests). Only 8 multi- 91  
modal behaviours were observed. These are shown 93  
in Table 4.

95 These examples provide illustrative semantic 97  
combinations of modalities:

- Deictic: "What's this?" + circling gesture on the 95  
picture of the Coliseum.
- Type of object mentioned in speech: 97
- "What's that picture?" + circling gesture on the 99  
picture of Andersen's mother;
- "I want to know something about your hat" + 101  
circling gesture on the hat.
- Linguistic reference to concepts related to the 103  
graphical object (e.g., "dad" and gesture on a  
picture) instead of direct reference to the object  
type or name ("picture");

Table 4  
Description of multimodal sequences observed in the Prototype-1 video corpus

Succession of modalities	Delay <sup>a</sup> between modalities (s)	Object gestured at	Shape of gesture	Spoken utterance + NLU frame	Cooperation between modalities
Gesture–speech	2	Picture of Coliseum	Circle	“What’s this?”	Complementarity
Simultaneous	0	Picture of Andersen’s mother	Circle	“What’s that picture?”	Complementarity
Simultaneous	0	Hat	Circle	“I want to know something about your hat.”	Redundancy
Gesture–speech	4	Statue of 2 people	Circle	“Do you have anything to tell me about these two?”	Complementarity
Simultaneous	0	Statue of 2 people	Point	“What are those statues?”	Complementarity
Gesture–speech	4	Picture above book-case	Circle	“Who is the family on the picture?”	Complementarity
Gesture–speech	3	Picture above book-case	Circle	“Who is in that picture?”	Complementarity
Simultaneous	0	Vase	Circle	“How old are you?”	Concurrency

<sup>a</sup>The delay between modalities was measured between end of first modality and end of second modality.

• Incompatibility between internal singular representation of objects and their plural/singular perceptual “affordance,” e.g., a single object is referred to in the user’s speech as a plurality of objects: “Do you have anything to tell me about these two?” (or “What are those statues?”) with a circling gesture on the statue of two characters which are internally represented as a single object.

Several objects might elicit such plural/singular incompatibility. They visually represent several entities of the same kind, but they are (system-) internally represented as a single object. They could be thus referred to as a single object or as several objects, their number being foreseeable for some of them: books (number > 2); boots (2); papers (> 2); pens (2); statue (2).

Conversely, although this was not observed as such in the Prototype-1 user test video, several objects of similar type and in the same area might be perceived as a single “perceptual group” [32] and might elicit a plural spoken reference combined with a singular gesture on only one of the items in the group: the group of pictures on the wall above the desk, the “clothes group” (coat–boots–hat–umbrella), the furniture (table and chairs), the small objects on the small shelf.

### 3.3. Temporal dimension of input fusion

A main issue for IF is to have a newly detected gesture wait for a possibly related spoken utterance.

How long should the gesture wait before the IF decides that it was indeed a mono-modal behaviour? We decided to use default values for delays to drive the IF to have gestures wait a little for speech (3 s) and have speech wait for gesture for a very short while only, since this is compatible with the literature [33] and the Prototype-1 user tests observations. We have also introduced management of “StartOfSpeech” and “StartOfGesture” messages sent to the IF in order to enable adequate waiting behaviour by the IF. Four temporal parameters of the IF have been defined to answer the following questions:

- How long should an NLU frame wait in the IF for a gesture when no “StartOfGesture” has been detected (Speech-waiting-for-gesture-short-delay)? The default value is 1 s.
- How long should an NLU frame wait in the IF for a gesture when a “StartOfGesture” has been detected (Speech-waiting-for-gesture-long-delay)? The default value is 6 s.
- How long should a GI frame wait in the IF for a NLU frame when no StartOfSpeech has been detected (Gesture-waiting-for-speech-short-delay)? The default value is 3 s.
- How long should a GI frame wait in the IF for an NLU frame when StartOfSpeech has been detected (Gesture-waiting-for-speech-long-delay)? The default value is 6 s.

The part of the IF algorithm that manages temporal behaviour is specified with the instructions

1 to be executed for each event that can be detected by  
 2 the IF: a new NLU frame is received by the IF, a  
 3 new GI frame is received by the IF, a “StartOf-  
 4 Speech” message is received by the IF, a “StartOf-  
 5 Gesture” message is received by the IF, a “Speech-  
 6 waiting-for-gesture” times out, and a “Gesture-  
 7 waiting-for-speech” times out.

8 The IF behaviour is described informally below  
 9 for each of these events.

#### 10 Init()

11 {Starts with ‘‘short’’ delays when no  
 12 start of speech or gesture has been  
 13 received. When start of speech/ges-  
 14 ture will be received, these will be  
 15 set to longer delays since there is a  
 16 very high probability that an asso-  
 17 ciated speech or gesture frame will be  
 18 received afterwards by the IF}  
 19 Speech-waiting-for-gesture-de-  
 20 lay = Speech-waiting-for-gesture-  
 21 *short*-delay  
 22 Gesture-waiting-for-speech-de-  
 23 lay = Gesture-waiting-for-speech-  
 24 *short*-delay  
 25

#### 26 When a new NLU frame is received by the IF

27 {Test if a gesture was already waiting  
 28 for this NLU frame}  
 29 If the timeout *Gesture-waiting-for-speech* is  
 30 running  
 31 Then  
 32 {A GI frame was already waiting for  
 33 this NLU frame}  
 34 Call semantic fusion on the NLU and  
 35 the GI frames  
 36 Stop-Timer (Gesture-waiting-for-  
 37 speech)  
 38 Else  
 39 {This new NLU frame will wait for  
 40 incoming gesture}  
 41 Start-Timer (Speech-waiting-for-  
 42 gesture)  
 43  
 44

#### 45 When a new GI frame is received by the IF

46 {Test if a NLU frame was already wait-  
 47 ing for this GI frame}  
 48 If the timeout *Speech-waiting-for-gesture* is  
 49 running  
 50 Then  
 51

A NLU frame was already waiting for 53  
 this GI frame}

Call semantic fusion on the NLU and 55  
 the GI frames

Stop-Timer (Speech-waiting-for- 57  
 gesture)

Else 59

{This new GI frame will wait for 61  
 incoming speech}

Start-Timer (Gesture-waiting-for- 63  
 speech)

#### 64 When a *startOfSpeech* message is received

65 {A new NLU frame will soon arrive. Ensure that 67  
 the GI frame that is already waiting waits longer 69  
 or that if a new GI frame arrives soon (since a 71  
 StartOfGesture was received) it will wait for the 71  
 NLU frame}

Gesture-waiting-for-speech-delay = Gesture- 73  
 waiting-for-speech-*long*-delay

If Gesture-waiting-for-speech is running 75

Then 75  
 Restart-Timer (Gesture-waiting-for-speech) 77

#### 78 When a *startOfGesture* message is received

79 {A new GI frame will soon arrive. 81  
 Ensure that the NLU frame that is 81  
 already waiting waits longer or that 83  
 if a new NLU frame arrives soon (since 83  
 a StartOfSpeech was received) it will 85  
 wait for the GI frame}

Speech-waiting-for-gesture-de- 87  
 lay = Speech-waiting-for-gesture- 87  
*long*-delay

If Speech-waiting-for-gesture is 89  
 running

Then 91  
 Restart-Timer (Speech-waiting- 93  
 for-gesture)

#### 94 When timeout *Speech-waiting-for-gesture* is over

95 {A NLU frame has waited for a GI frame 97  
 which did not arrive}

Build and send an IF frame containing 99  
 only the NLU frame

Stop-Timer (Speech-waiting-for-ges- 101  
 ture)

Init() 103

#### 96 When timeout *Gesture-waiting-for-speech* is over

```

1   {A GI frame has waited for a NLU frame
    which did not arrive.}
3   Build and send an IF frame containing
    only the GI frame
5   Stop-Timer (Gesture-waiting-for-
    speech)
7   Init()

```

3.4. Semantic dimension of input fusion

Regarding semantic IF we have decided to focus on (1) the semantic compatibility between gestured and spoken objects, and (2) the plural/singular property of these objects. We limited ourselves to one reference per NLU frame and identified 16 possible semantic combinations of speech and gesture (Table 5).

Only cases 11, 12, 15, and 16 can possibly lead to fusion in the IF, as described above. We systematically analysed each of the 16 cases. Below, we specify the instructions to be executed by the IF and the output it produces for each case. The instructions consider the following features of speech and gesture references: singular/plural, reference/no reference, semantic compatibility.

Semantic compatibility between gestured and spoken objects is evaluated by the IF via semantic distance computation which is less strict than object type unification and was expected to be more appropriate for conversational systems for children. Semantic distance computation makes use of a graph of concepts connected with an “is-related-to” relation. Each concept is represented by: a name (e.g., “feather Pen,” “\_Family”), a plural Boolean (e.g., “true” for the statue of two people), a singular Boolean (e.g., “true” for the feather Pen), a Boolean describing if it is an object in the study (“picture-

ColiseumRome”) or an abstract concept (“\_Mother”), and the set of semantically related concepts (generic relation “isRelatedTo”).

A reference detected by the NLU module is represented in the IF by: a Boolean stating if it is solved, a Boolean stating if it is plural/singular, and a Boolean stating if it is numbered (if yes, an attribute gives the number of referenced objects, e.g., “two” in the reference “these two pictures”).

A perceptual group is represented by the same attributes as a single concept, and by the set of concepts, which might be perceived as a group (e.g., the set of pictures above the desk).

The identified cases of semantic combination described above are integrated in a single algorithm for semantic fusion. The informal algorithm below only details cases for which one message has been sent by the NLU and one by the GI, i.e., cases 6–7–8, 10–11–12, 14–15–16 in our analysis.

After IF, when required, an IF frame is sent to the character module. An attribute called “fusion Status” is used in the IF frame to indicate if the input was mono-modal (“none”), successful (“ok”) or unsuccessful (“inconsistency”). Gestures towards objects that cannot be referenced are ignored and hence are not passed to the character module.

**Algorithm Semantic Fusion (NLU frame, GI frame)**

```

{Manage each multimodal combination
case. We suppose that one NLU frame
and one GI frame have been received by
the IF}
IF there is no explicit reference in
the NLU frame
THEN {CASES 6–7–8}
    Group both frames
    Send them to the Character Module
    with a fusion status set to none

```

Table 5  
Analysing 16 combinations of speech and gesture along the singular/plural dimension of references (only cases 11, 12, 15, and 16 can possibly lead to fusion in the IF)

GI/NLU	No message from GI	1 message from GI “noObject”	1 object detected by GI “select”	Several objects detected by GI
“referenceAmbiguity”				
No message from NLU	1	2	3	4
1 message from NLU but no explicit reference in NLU frame	5	6	7	8
1 message from NLU with 1 singular reference	9	10	11	12
1 message from NLU with 1 plural reference	13	14	15	16

1	ELSE	Character Module	53
3	IF there is only one reference in the NLU frame	ELSE	55
5	THEN	{Manage perceptual groups}	57
7	IF the reference is singular	IF there is only one object from GI	59
9	THEN call Semantic Fusion Singu- lar NLU (NLU frame, GI frame)	compatible with NLU reference and this object belongs to a percep- tual group	61
11	ELSE call Semantic Fusion Plural NLU (NLU frame, GI frame)	THEN	63
13		{Do semantic fusion}	65
15	<b>Semantic Fusion Singular NLU (NLU frame, GI frame)</b>	Resolve the plural NLU reference with the perceptual group of ob- jects	67
17	{The referential Expression in the NLU frame is singular:	Send the modified NLU frame to the Character Module	69
19	CASES 10-11-12 (not perceptual group) }	ELSE	71
21	IF there is at least one object se- lected by GI,	IF the GI object is compatible with the NLU reference	73
23	which is semantically compatible with the NLU reference	but does not belong to a percep- tual group	75
25	THEN	THEN {Do semantic fusion (not considering plural constraint) }	77
27	{Do semantic fusion (possibly not considering plural constraint	Resolve NLU reference with the compatible gestured object	79
29	if there was several gestured ob- jects) }	Send the modified NLU frame to the Character Module	81
31	Resolve the NLU reference with the compatible gestured object(s)	ELSE {No gestured object compa- tible with the NLU plural ref. }	83
33	Send the modified NLU frame to the Character Module	Signal inconsistency ; Send NLU frame and GI frame	85
35	ELSE		87
37	{No gestured object revealed com- patible with the NLU reference}		89
39	Signal inconsistency		91
41	Send NLU frame and GI frame to the Character Module		93
43			95
45	<b>Semantic Fusion Plural NLU (NLU frame, GI frame)</b>		97
47	{The Referential Expression is plur- al: CASES 14-15-16-12 (perceptual group) }		99
49	IF more than one object from GI is semantically compatible with the NLU reference		101
51	THEN		103
	{Do semantic fusion}		
	Resolve the plural NLU reference with the compatible gestured ob- ject(s)		
	Send the modified NLU frame to the		

The different feedforward and feedback mechanisms that have been implemented to enable proper coordination of multi-modal input with Andersen's behaviour are summarised in Fig. 5.

### 3.5. Character module processing

Given the many design-time uncertainties concerning how children would use combined speech and gesture input, we chose a simple processing scheme for gesture-related input in the character module. The IF frame goes to the character module's conversation mover, which tries to match the input to candidate system output. The conversation mover passes on its results to the conversation mover post-processor whose task it is to select among the conversation mover outputs a single output candidate to pass on to the move processor which analyses the candidate in the discourse history and domain knowledge contexts. The conversation mover does nothing about gesture-related input, i.e., gesture-only input and

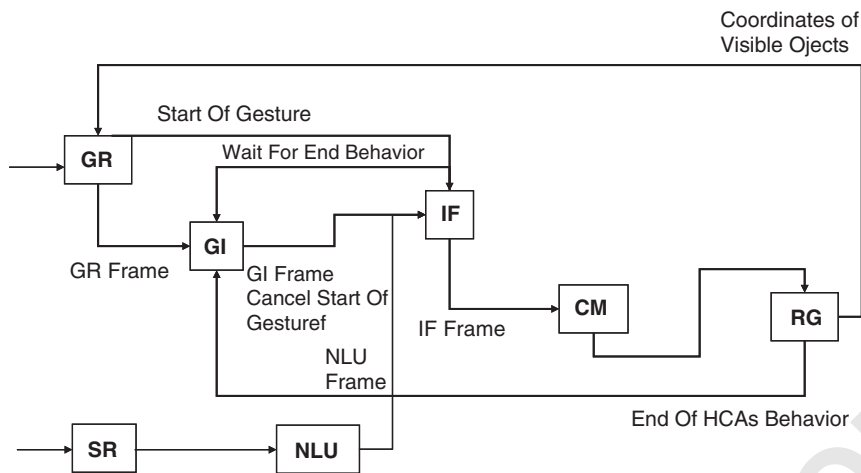


Fig. 5. Feedforward and feedback messages for managing multi-modal input conversation with Andersen (GR = Gesture Recogniser, GI = Gesture Interpreter, SR = Speech Recogniser, NLU = Natural Language Understanding, IF = Input Fusion, CM = Character Module, RG = Graphic Renderer). Messages “GRFrame,” “GIFrame,” “NLUFrame,” and “IF Frame” enable the transmission of processing results of modules. Messages “StartOfGesture,” and “CancelStartOfGesture” enable the proper management of temporal relations between speech and gestures. Messages “WaitForEndBehavior” and “EndOfHCABehavior” enable inhibition of gesture processing while the character is responding, hence regulating turn-taking. Message “GIFrame” is also used by the character to provide gaze feedback on the gestured object.

combined gesture-speech input, but simply passes them on to the conversation mover post-processor. Informally, the post-processor’s algorithm for gesture-related input is:

- check if multiple labels include label(s) prefixed by  $g\_$  [these are gesture object labels]
  - if yes, remove all labels **not** prefixed by  $g\_$
  - if only one label remains, send label to move processor **END**
  - if several labels remain, continue
- randomly select a label among the multiple labels left and send the selected label to move processor **END**

Thus, the character module ignores the “inconsistency” label from the IF and does not attempt to produce meta-communication output in an attempt to resolve the inconsistency claimed by the IF. We selected this solution because of the problems we have identified with singular vs. plural deictic expressions and what they might refer to (cf. Section 4). Furthermore, the character module does not process the spoken input in cases where the IF has deemed IF to be “ok”. Also, by not processing the spoken input in cases of independent concurrency, i.e., when the user points to some object(s), but speaks about something else entirely, the strategy

adopted means that Andersen at least manages to address one of the user’s concerns, i.e., that of getting a story about a referenceable object. What he does not do is keep in mind that the user had spoken about something else entirely whilst pointing to some object(s). Our design reasoning was that the user, when noticing this, might simply come back and repeat the spoken input in a subsequent turn. Arguably, this design decision is an acceptable one since the user (i) does get a reply wrt. to the object pointed to and (ii) has ample opportunity to come back to the unrelated issue posed in the spoken part of the input. Given the overall design of the Prototype-2 system, the only apparent flaw would seem to be the fact that the user’s spoken input might relate more closely to gesture input information randomly *discarded* by the post-processor than to the gesture input information randomly *chosen* by the post-processor. However, selecting wisely in this situation would either (i) require the conversation mover to have contextual knowledge which it does not possess or (ii) that the post-processor forward multiple output candidates to the move processor which does have contextual knowledge, and this is not possible in the Andersen Prototype-2 system.



## 1 4. Evaluation

### 3 4.1. Methodology

5 The Prototype-2 Andersen system was tested with  
 7 13 users (six boys and seven girls) from the target  
 9 user population of 10–18 years old children and  
 11 teenagers in February 2005. All users were Danish  
 13 school kids aged between 11 and 16 and with an  
 15 average age of 13 years. Their English skills were  
 17 not rated prior to the test as we wanted to test the  
 19 system with a random sample of target users. The  
 21 Prototype-1 test—following which the 18 children  
 23 users' English skills were rated for speech recogniser  
 25 training purposes—had shown that Danish kids are  
 27 generally able to conduct conversation with Ander-  
 29 sen even though, of course, their English proficiency  
 31 varies significantly depending upon factors, such as  
 33 age, individual differences, and temerity in address-  
 35 ing Andersen in the presence of unfamiliar adults.  
 37 As in the Prototype-1 user test, in the test of  
 39 Prototype-2 only a single child had significant  
 41 difficulties carrying out conversation with Ander-  
 43 sen. In the post-test structured interview, the users  
 45 were asked about their knowledge of Andersen's  
 47 fairytales. Their responses were all rated by two  
 49 independent raters at 2 on a 3-point scale, which  
 51 corresponds closely to the findings in the post-test  
 interview following the Prototype-1 test. Danish  
 children generally have substantial knowledge  
 about Andersen's fairytales. Only two of the  
 Prototype-2 users had had conversation with  
 Andersen before, i.e., in the Prototype-1 user test.

The test was a controlled laboratory test rather  
 than a field test in the Andersen museum. For the  
 first user test of a strongly modified second  
 prototype, it is often preferable to make use of the  
 laboratory environment in order to be able to fully  
 control the conditions of interaction, such as  
 advance notice of users in order for them to plan  
 for the entire (60–75 min) duration of the test which  
 included structured post-trial interviews, common  
 instructions to all users for each test phase, timing  
 of the two different test conditions that were used  
 for all users, etc. Admittedly, a field trial would have  
 provided more realistic data on system use, but this  
 data would also have been very different from the  
 data collected in the lab.

Users were wearing a microphone/loudspeaker  
 headset. They used a touch screen for gesture input  
 and a keyboard for controlling virtual camera  
 angles and for controlling Andersen's locomotion.



Fig. 6. A user talking to the 2nd Andersen system prototype.

Each user had a total of 35 min of multi-modal  
 interaction with Andersen, the conversation being  
 conducted in English. Each user interacted with  
 the system in two different test conditions. In the first  
 condition, they received basic instructions on how  
 to operate the system but not on how to speak to it,  
 and then spent approx. 15 min exploring the system  
 through conversation with Andersen. In the second  
 condition, in order to steer the users through a  
 cross-section of Andersen's domain knowledge and  
 put pressure on the system's ability to handle  
 substantial user initiative in conversation, they  
 received a handout with 11 issues they might wish  
 to address during conversation at their leisure for  
 20 min, such as "Try to offend Andersen" or "Tell  
 Andersen about the games you like to play". Fig. 6  
 shows a user in action.

Two cameras captured the user's behaviour  
 during interaction and all main module outputs  
 were logged. Following the test, each user was  
 interviewed separately about his/her experience  
 from interacting with Andersen, views on system  
 usability, proposals for system improvements, etc.

#### 4.2. Comparative analysis of video and log files

Eight hours of interaction were logged and  
 captured on video. In order to evaluate the GR,  
 GI and IF modules, the gesture-only and gesture-  
 combined-with-speech behaviours were analysed  
 based on the videos and the log files. The videos  
 were used to annotate the real behaviours displayed  
 by users in terms of: spoken utterances related to  
 gestural behaviour, the objects gestured at (includ-  
 ing each non-referenceable object, i.e., objects in

1 Andersen's study for which the animation does not  
 3 have an id to forward to the GI), and obvious or  
 5 possible misuse of the tactile screen in case the  
 7 corresponding gesture was not detected by the GR.  
 9 The log files were used to check the output of each  
 11 module, to compare the output to the observed  
 13 behaviour from the video, and to classify reasons  
 15 for, and cases of, failure.

17 We made a distinction between the success of the  
 19 interaction and the success of the processing done  
 21 by the gesture and multi-modal modules. *Multi-*  
 23 *modal interaction* was considered successful if the  
 25 system responded adequately to the user's beha-  
 27 viour, i.e., if the character provided information  
 29 about the object the user gestured at and/or spoke  
 31 about. *Module success* was evaluated by comparing  
 the user's behaviour and the output produced by the  
 modules in the log files. In some cases, the  
 interaction was successful although the output of  
 the module was incorrect, implying that the module  
 error was counter-balanced by other means or  
 modules. In some other cases, the interaction was  
 unsuccessful although the output of the module was  
 correct, implying that an error occurred in some  
 other module(s). Interaction success for multi-  
 modal input provides information on, among other  
 things, the use of inhibition and timing strategies,  
 which enable proper management of some redun-  
 dant multi-modal cases via the processing of only  
 one of the modalities.

#### 33 4.2.1. Gesture recognition

35 281 gesture shapes onto the tactile screen were  
 37 logged. The shapes were manually labelled without  
 39 displaying the result of GR processing (blind  
 41 labelling). To enable fine-grained analysis of gesture  
 43 shapes, the labelling made use of 25 categories of  
 45 shapes. We found that 87.2% (245) of the logged  
 47 gestures had been assigned the same category by the  
 GR and by the manual labelling process. The fine-  
 grained categories reveal a high number of diagonal  
 lines ( $90/281 = 32\%$ ) and explicitly noisy categories  
 (44/281 = 16%), such as garbage, noisy circle, and  
 open circle of various orientations. The distribution  
 of shapes in the GR and the manual labelling are  
 similar.

#### 49 4.2.2. Gesture interpretation

51 As observed in the videos, the users made 186  
 gesture-only turns. If we use the number of IF  
 frames (957) for counting the number of user  
 turns—this is not exact as sometimes a single

53 spoken turn might be divided into several recog-  
 55 nised utterances—gesture-only turns correspond to  
 19% of the user turns.

One hundred and eighty seven messages were  
 produced by the GI module. By comparing the log  
 files and the videos, we found that 54% of the user  
 gestures led to a GI frame, 30% were cancelled  
 because detected after GI timeout and during or  
 before the character's response, and 16% were  
 grouped because they were done on the same object.

The repartition of the gesture interpretation  
 categories is the following:  $125/187 = 67\%$  detected  
 a single referenceable gestured object,  $61/$   
 $187 = 33\%$  did not detect any referenceable object,  
 and only one detected several referenceable objects  
 in a single gesture. One multi-object gesture was  
 observed in the video, but this gesture included one  
 referenceable object and two non-referenceable  
 objects and was thus interpreted as selection of a  
 single object by the system.

73 Fifty one percent of the gesture-only behaviours  
 75 led to interaction success. The reasons for the 49%  
 77 cases of interaction failure were classified as follows:  
 79 gesture on non-referenceable objects (62%), gesture  
 during GI inhibition (17%), system crash (14%),  
 81 unexplained (4%), gestured object not detected  
 (2%), gesture not detected (1%). Most of the  
 83 interaction failures (76%) were thus due either to  
 gestures onto non-referenceable objects or to input  
 inhibition. On average, each user gestured at 11  
 referenceable objects and 4 non-referenceable ob-  
 85 jects.

#### 87 4.2.3. Input fusion

89 As observed in the videos, the users made 67  
 91 multi-modal turns combining gesture and spoken  
 93 input. If we use the number of IF frames as our  
 number of user turns, multi-modal turns correspond  
 to 7% of the user turns. Among the 957 messages  
 logged by the IF, only 21 (2%) were processed by  
 the system as multi-modal constructions.

95 Seventy percent of the multi-modal turns were  
 97 produced in the first test condition, cf. Section 4.1.  
 99 This is the same proportion as for gesture-only  
 behaviours. It is probable that, during the first test  
 phase, the users explored the 3D environment,  
 testing objects by gesturing and sometimes speaking  
 at the same time to find out if Andersen had stories  
 to tell about those objects. When the second test  
 101 condition started, the users had already received  
 103 information about a number of objects and  
 preferred to address topics other than the objects

1 in the study. In support of this interpretation it may  
2 be added that only one of the 11 issues in the  
3 second-condition handout concerned objects in  
Andersen's study (cf. Section 4.1).

5 Regarding the users' multi-modal behaviours, we  
6 also analysed interaction success and IF success. In  
7 24 multi-modal turns, the IF was unsuccessful, but  
8 interaction was successful. Sixty percent of the  
9 multi-modal behaviours led to interaction success.  
10 Analysis of the output of the IF module reveals that  
11 it worked well for 25% of the multi-modal cases. It  
12 is quite difficult to compare such results with the  
13 literature since there are very little experimental  
14 results on multi-modal fusion in conversational  
15 applications for children. For example, Kaiser et al.  
16 [16] observed an overall success in functional  
17 accuracy of 59.1% and 81.4% for multi-modal  
18 recognition but during adult's speech and 3D  
19 gestures multi-modal commands for manipulating  
3D objects.

21 The reasons for failure of processing multi-modal  
22 behaviours were collected from the video and log  
23 files and are listed in Table 6.

25 A closer analysis was done of the many "timer  
26 too small" cases, i.e., the cases in which the IF's  
27 1.5s waiting time for linguistic input after having  
28 received gesture input from the GI, was not long  
29 enough. The linguistic input did arrive and was  
30 temporally related to the gesture input, but it  
31 arrived too late for IF to take place, the gesture  
32 input already having been sent to the character  
33 module. In 85% of these 21 cases, the timestamp of  
34 the IF's "StartOfSpeech" message was evaluated as  
35 being incorrect compared to the start of speech  
36 observed in the video. It would have been inap-  
37 propriate to have the user wait for such a long  
38 period, e.g., 10s in several cases. For example, the  
39 "start of speech" would be logged as arriving in the

41 Table 6  
42 Reasons of failure in processing of multimodal behaviours

	NB	%
45 Timer too small	21	43
Speech recognition error	9	18
47 Input inhibited	6	12
Not a referenceable object	4	8
Gesture not detected	4	8
49 System crash	2	4
Unexplained	2	4
51 Gestured object not detected	1	2
Total	49	100

IF 14s after the "start of gesture" although, in the 53  
video, the user starts to speak only 1s after the start 55  
of gesture. Indeed, given the limited semantics of 57  
gesture involved, i.e., only selection of objects, and 59  
the frequent redundancy of speech and gesture in 61  
the conversational context, the strategy to take an 63  
early decision for gesture-only behaviour enabled us 65  
to obtain 60% of interaction success for multi- 67  
modal behaviour while avoiding the user waiting 69  
too long for the system's response. The IF would 71  
briefly wait for NLU input and then send its frame 73  
to the character module, ignoring any delayed NLU 75  
input. The explanation for the delayed "start of 77  
speech," as this is labelled by the IF, turned out to 79  
be a flaw in the speech recogniser's detection of *end* 81  
of speech, so that the recogniser would continue to 83  
listen until timeout even if the user had stopped 85  
speaking maybe 10s before. This flaw turned out 87  
not to be more complex to correct than expected 89  
because it was due to the fact, unknown to us at the 91  
time, that we should have used a different approach 93  
for implementing end of speech detection in the 95  
Scansoft recogniser.

In line with previous observations [34], 6% of the 97  
multi-modal input turns proved to be concurrent, 99  
i.e., speech and gesture were synchronised, but 101  
semantically unrelated. For example, one user said 103  
"Denmark" to answer the system's question about  
the user's country of origin while gesturing on the  
picture of the Coliseum. Another user said "Where  
do you live?" while gesturing on the feather pen on  
the desk.

The evaluation of the GR, GI and IF modules 85  
can be summarised as follows: 87

- GR failures represent 12.8% of gestural inputs, 89  
but had no impact on interaction success.
- Failures in processing gesture-only input for 91  
*referenceable* objects involved the GI module in  
only 4% of the cases.
- Fusion failures occurred for 40% of the multi- 93  
modal behaviours. Three-fourth of these cases  
correspond to missing fusions and 1/4 to 95  
irrelevant fusions.

Thus, our comparative analysis of the video and 97  
log files shows that the gestures done on non- 99  
referenceable objects and the gestures done while 101  
the character was speaking or preparing to speak, 103  
had a quite negative impact on gesture interpreta-  
tion. This is true both for the processing of gesture-  
only and multi-modal behaviours. Both might be

1 due to the graphical affordance of referenceable  
 3 objects and the lack of visibility of the non-verbal  
 5 cues shown by the character. Indeed, graphical  
 7 affordance could be improved in our system so that  
 9 (1) the users can visually detect the objects the  
 11 character can speak of, e.g., these referenceable  
 13 objects could be permanently highlighted, (2) the  
 15 users understand that the character is willing to take  
 17 or to keep the turn, e.g., the camera could be  
 19 directed towards the character's face in such cases,  
 21 thus enhancing the visibility of the non-verbal cues  
 23 for turn-taking management. Our analysis also  
 25 reveals how the dimensions of fusion were used by  
 27 the user and processed by our system. We observed  
 29 that the proper management of temporal informa-  
 31 tion, such as the reception of a start of speech  
 33 message at the right time has a huge impact on IF  
 35 success. Regarding the semantic dimension, users  
 37 only rarely did multi-object selection with a single  
 39 gesture or made implicit spoken references to  
 41 objects.

### 4.3. Interviews

53  
 55 Fig. 7 presents a summary of the users' answers in  
 57 the post-test interviews. For each interview ques-  
 59 tion, each user's answer was scored independently  
 61 by two scorers on a 3-point scale from (1) positive  
 63 with minor or no qualifications, over (2) positive  
 65 with qualifications, to (3) negative/with substantial  
 67 qualifications [35].

69 Six questions (Q(n)s) in the user interviews  
 71 address gesture-related issues. On the question  
 73 (Q3) *if Andersen was aware what the user pointed*  
 75 *to*, most users were quite positive although some  
 77 pointed out that Andersen ignored their gestures in  
 79 some cases. This was expected due to the large  
 81 number of non-referenceable objects in Andersen's  
 83 study and is confirmed by the analysis in Section  
 85 4.2. The kids were almost unanimously positive in  
 87 their comments on Q4, *how it was to use the touch*  
 89 *screen*, which they found easy and fun. Like in the  
 91 first prototype user interviews [2], the children were  
 93 divided in their opinions on Q5 as to *whether they*  
 95 *would like to do more with gesture*. Half of the users  
 97 were happy with the 2D gesture affordances while

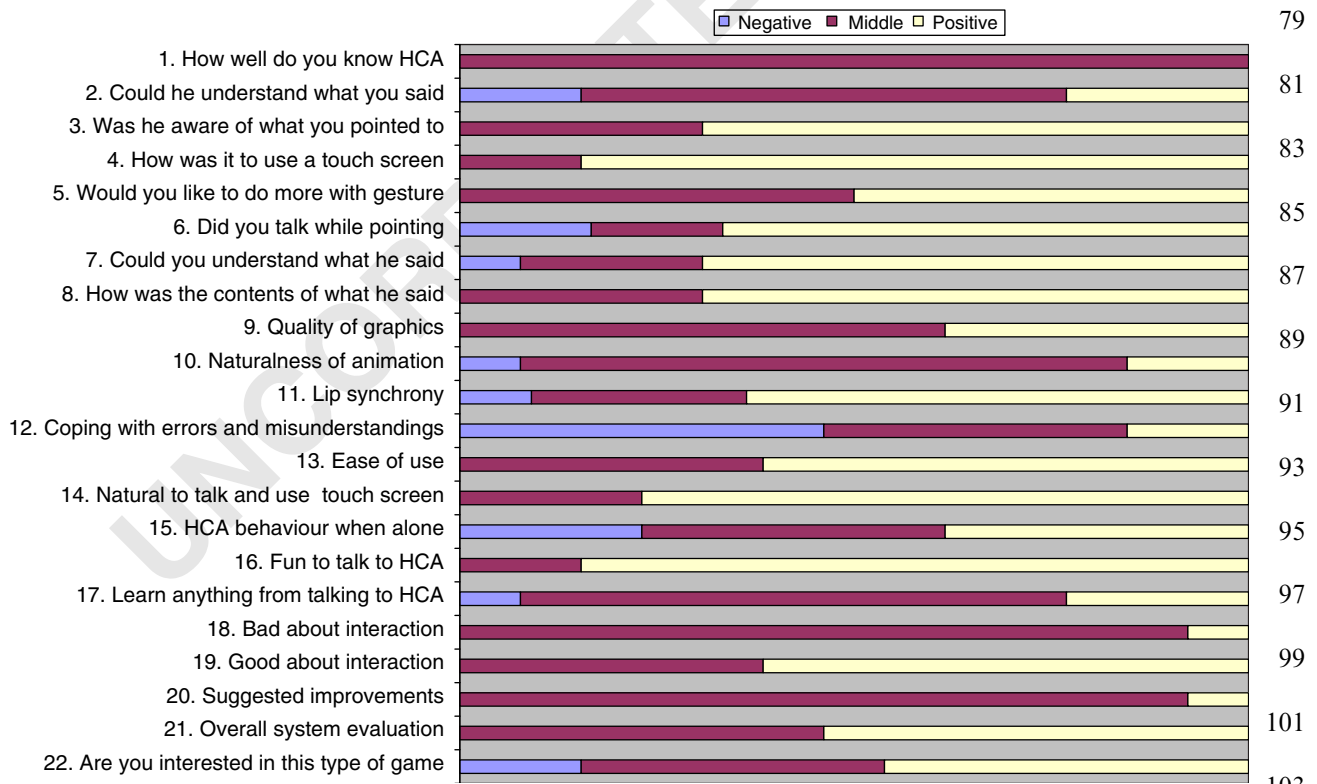


Fig. 7. Summary of interview results.

1 the other half wished to be able to gesture towards  
 2 more objects in Andersen's study. On the question  
 3 (Q6) *whether they talked while pointing*, only a  
 4 couple of users said that they never tried to talk and  
 5 point at the same time. We will return to this point  
 6 below. Finally, on the question (Q14) if the users *felt*  
 7 *it to be natural to talk and use the touch screen*, the  
 8 large majority of users were again quite positive.

9 In summary, the Danish users of the second  
 10 Andersen prototype were almost unanimously  
 11 happy about the available modality/device input  
 12 combinations, i.e., pointing gesture input via touch  
 13 screen and speech input via microphone headset  
 14 (Q4, Q14). Andersen sometimes ignored the users'  
 15 pointing gestures (Q3), which perhaps partly ex-  
 16 plains why half of the users wished to be able to  
 17 elicit more stories from Andersen through gesture  
 18 input (Q5). Finally, the majority of users claimed  
 19 that they, at least sometimes, talked while pointing  
 20 (Q6).

21 Globally, users were happy with gestural and  
 22 multi-modal input and many wished to do more  
 23 with gestures, which is congruent with previous  
 24 observation that gesture is a key modality for young  
 25 users to have fun and take initiative in the  
 26 interaction [36].

#### 27 4.4. Follow-up experiment with native English 28 speakers

31 Following the second prototype user test, de-  
 32 scribed above, with Danish children having English  
 33 as their second language, we did a small user test  
 34 with four children, two girls and two boys, 11–13  
 35 years old, all of whom had English as their first  
 36 language. The primary purpose of the test was to  
 37 explore the effects of (i) users' first language and (ii)  
 38 the amount of instruction received on how to speak  
 39 to the system. Thus, the English children were  
 40 provided with extensive instructions on how to  
 41 speak to the system during the first test condition,  
 42 whereupon they carried out the second test condi-  
 43 tion in the same way as the Danish kids did, cf.  
 44 Section 4.1. In what follows, we focus on a single  
 45 finding in the test related to the Danish kids'  
 46 response to Q6, i.e., that they often talked while  
 47 gesturing.

48 To compare the Danish children with the English  
 49 children, we randomly sampled four Danish chil-  
 50 dren from the Danish user population, two girls and  
 51 two boys. We then looked at the transcriptions from  
 the directly comparable 2nd-condition trials in

Table 7  
 Combined speech and gesture input in two user groups

	Danish children	English children
No. of input turns	201	267
No. of speech–gesture turns	0	30
No. of speech–gesture turns per user	0-0-0-0	12-2-4-12
No. of gesture-only turns	15	4

which all children were invited to address, at their  
 leisure, topics from a list of 11 topics in conversa-  
 tion with Andersen. Table 7 shows what we found  
 on the use of combined speech and gesture input in  
 the two test groups.

Table 7 shows that the randomly sampled Danish  
 users did *not* speak while gesturing at all. This is in  
 sharp contrast to Danish group's response to (Q6)  
*whether they talked while pointing*. Even if, by  
 (unlikely) chance, the sampled Danish group  
 includes the two Danish users who admittedly never  
 tried to talk and point at the same time, Table 7  
 includes four users who did not do that in the 2nd  
 test condition. They might, of course, have done so  
 in the first test condition. Whatever the explanation  
 might be, this contrasts markedly with the English  
 users, all of whom spoke when they gestured except  
 in 12% of the turns in which they used gesture  
 input. When the Danish kids in the sampled group  
 used gesture, they never spoke at the same time.

The hypothesis arising from Table 7 is that there  
 are significant behavioural differences between  
 children having English as their first language and  
 children having English as their second language, in  
 the way they use the speech and gesture input  
 affordances available. In order to obtain informa-  
 tion on objects that can be indicated through  
 gesture, the former naturally speak while gesturing  
 whereas the latter tend to choose gesture input-only.  
 The explanation for this hypothesis probably is that  
 the opportunity to complete a conversation act  
 without speaking a foreign language tends to be  
 favoured whereas, for users speaking their mother  
 tongue, it is more natural to speak and gesture at  
 the same time. It should be noted here that the  
 English users were very young, which speaks against  
 attributing their more frequent use of multi-modal  
 input to speaker maturity. This finding, hypothet-  
 ical as it remains due to the small user populations  
 involved, must be kept in mind when interpreting

1 the results presented in this paper, most of which  
 3 have been gathered with users having English as  
 their second language.

## 5. Discussion

7 In this paper, we have presented early results on  
 how 10–18 years old Danish children having English  
 9 as their second language use speech and 2D gesture  
 to express their communicative intentions in con-  
 11 versation with a famous 3D animated character  
 from the past. In a small control study with 11–13  
 13 years old children having English as their first  
 language, we found that the pattern of multi-modal  
 15 interactive input apparent in the Danish kids might  
 be significantly different in the English-speaking  
 17 children. In essence, the English-speaking kids  
 practice what the Danish children preach, lending  
 19 strong joint support for the conclusion that the  
 multi-modal input combination of speech and touch  
 21 screen-enabled conversational input is a highly  
 natural input combination for conveying users'  
 23 communicative intentions to embodied conversa-  
 tional characters.

25 From a technical point of view, the work reported  
 shows, first of all, that we are only at the very  
 27 beginning of addressing the enormous challenges  
 facing developers of natural interactive systems  
 29 capable of understanding combined speech and  
 2D gesture input. In the following, we describe some  
 31 of those challenges viewed from the standpoint of  
 having completed and tested the 2nd Andersen  
 33 system prototype.

### 5.1. Mouse vs. touch screen gesture input

37 It seems clear that gesture input via the touch  
 screen device is far more natural for conversational  
 39 purposes than gesture input via the mouse or similar  
 devices, such as controllers. The mouse (controller)  
 41 is a haptic input device, which a large user  
 population is used to employ for, among other  
 43 things, purposes of fast haptic control of computer  
 game characters and other computer game entities.  
 45 However, these input devices are far from being  
 natural in the context of natural interactive *con-*  
 47 *versation*. When offered these devices, as we  
 observed in the Prototype-1 user tests [30], the users  
 49 tend to “click like crazy,” following their—natural  
 or trained—tendency to gesture around in the  
 51 graphical output space without considering the  
 conversational context. Conversely, when offered

the more natural option of gesturing via the touch 53  
 screen in a speech-gesture conversational input 55  
 environment, no user seems to be missing the fast  
 interaction afforded by the mouse (controller). On 57  
 the contrary, given the interactive environment just  
 described, users seem perfectly happy with gesturing 59  
 via the touch screen, thereby emulating quite closely  
 their real-life-familiar 3D pointing gestures, cf. Fig. 61  
 7, Question 4.

### 5.2. Referential disambiguation through gesture

65 While the Danish users clearly seem to have  
 understood that they could achieve unambiguous 67  
 reference to objects without having to speak, they  
 also understood that spoken deictics require gesture 69  
 for referential disambiguation. Confirming the  
 users' claims about the intuitive naturalness of 71  
 using touch screen-mediated 2D gesture, the chil-  
 73 dren seem to be keenly aware of the need to point  
 while referring in speech to the object pointed  
 towards.

75 Another important point is that the users'  
 coordinated spoken references to pointed-to objects 77  
 were generally deictic in nature, making them  
 amenable to handling by the IF component we 79  
 had designed. Thus, in the large fraction of the 67  
 coordinated speech-gesture inputs in which the  
 81 speech part actually did refer to the object(s)  
 pointed towards, only one did not include deictics,  
 83 i.e., “Would you please tell me about the watch”.

### 5.3. Deictics fusion is only the tip of the iceberg

87 Essentially, the IF approach adopted for the  
 Andersen system aims at semantic fusion of singular 89  
 vs. plural spoken deictics with the number of named  
 objects identified through gesture interpretation. IF 91  
 also manages implicit or explicit references to  
 concepts related to (system-internally) named ob- 93  
 jects in Andersen's study. For instance, “Do you  
 like travelling” would be merged with a gesture on 95  
 one particular object, i.e., Andersen's travel bag.  
 What we found was that most users employed 97  
 spoken deictics, i.e., pure demonstratives, such as  
 ‘this’ in “What is this?” and only rarely used more 99  
 explicit referential phrases, such as noun phrases.

101 However, even this simple fusion domain is  
 subject to the fundamental ambiguity between, on  
 the one hand, how many physical objects the user 103  
 intends to refer to and, on the other, how many  
*within-object* entities the user intends to refer to,

1 such as several objects depicted in a single picture.  
 2 To resolve this ambiguity, the system would need  
 3 knowledge about the inherent structure and con-  
 4 tents of objects, such as pictures. Moreover, spoken  
 5 deictics do not necessarily refer to gestured-towards  
 6 objects. It is perfectly normal for spoken deictics to  
 7 anaphorically refer to the spoken discourse context  
 8 itself, as in “Are these your favourite fairytales?”  
 9 Given the fact that users sometimes perform  
 10 mutually independent (or concurrent) conversation  
 11 acts through speech and gesture, respectively, the  
 12 system would need quite sophisticated meta-com-  
 13 munication defences to pick up the fact that the user  
 14 is not performing a single to-be-fused conversa-  
 15 tion act but, rather, two quite independent con-  
 16 versation acts. Finally, requiring the system to be  
 17 able to manage, and hence to have knowledge  
 18 about, the internal structure and contents of objects,  
 19 such as pictures, is a demanding proposition. In the  
 20 foreseeable future, we would only expect highly  
 21 domain-specific applications to be able to handle  
 22 this problem, such as museum applications for users  
 23 to inquire about details in museum exhibit paint-  
 24 ings.

#### 25 5.4. Other chunks of the iceberg

26  
 27 As we saw in Section 4, users may, in principle,  
 28 point to anything in Andersen’s study and speak at  
 29 the same time. Furthermore, what they may  
 30 relevantly say when gesturing is open-ended, in-  
 31 cluding, for instance, the volunteered conversation  
 32 act <pointing to a chair> “My grandfather has a  
 33 chair like that”. This conversation act is relevant  
 34 simply because Andersen’s study is one of the  
 35 system’s domains of conversation. Users may also  
 36 explore relationships among objects, requiring the  
 37 character to have a model of these, as in <pointing  
 38 to picture of Coliseum> “Do you have other  
 39 pictures from your travels?”

40 We do not believe that the current Andersen  
 41 system architecture (Fig. 3) is the best solution for  
 42 handling the just illustrated, full-scale speech-  
 43 gesture IF for domain-oriented conversation. At  
 44 the very least, it seems, NLU must be made aware  
 45 that the currently processed spoken input is being  
 46 accompanied by gesture input. Otherwise, the  
 47 complexity to be handled by IF is likely to become  
 48 monstrous. An even better solution may be to  
 49 process speech and gesture input together, removing  
 50 the need for a subsequent late semantic IF  
 51 component. As regards conversation management

(in the character module) and response generation, 53  
 on the other hand, we see no evident obstacles for 55  
 the current architectures to process far more 57  
 complex IF than what is currently being processed 59  
 by the Andersen system.

In conjunction with Andersen’s injunctions to do 59  
 so, the design of Andersen’s study did lead the users 61  
 to gesture at the pictures on the walls. Inevitably, 63  
 however, these factors also made the users try to 65  
 find out which objects Andersen could actually tell 67  
 stories about. In the first Andersen prototype, we 69  
 had an additional class of “anonymous objects” 71  
 which were referenceable, but which, when gestured 73  
 upon, made Andersen say that he did not know 75  
 much about them at present. In the second 77  
 prototype, we dropped this class because it was felt 79  
 that Andersen’s response was not particularly 81  
 informative or interesting, and tended to be tedious 83  
 when frequently repeated. Since, for Prototype-2, 85  
 we did not increase the number of objects which 87  
 Andersen had stories to tell about, the consequence 89  
 was an increase in the number of failures in gesture 91  
 interpretation and IF since the users continued to 93  
 gesture at objects which were presented graphically, 95  
 but which the system did not know about (i.e., the 97  
 non-referenceable objects). There is no easy solution 99  
 to this problem. One solution is to increase the 101  
 number of objects which Andersen can tell stories 103  
 about until that number converges with the objects  
 which the majority of users want to know about.  
 Another solution is to make Andersen know about  
 all objects in his study, including the ceiling and the  
 carpet. A third, more heavy-handed and less  
 natural, solution might be to have specific rendering  
 for the objects the user can gesture at to get  
 Andersen to tell about them, such as by using some  
 form of permanent highlighting.

The user did not use the cross shape in their 91  
 gestures. This might be due to the fact that this 93  
 gesture shape is not that appropriate for the tactile 95  
 screen.

Selection of several objects in a single gesture, 95  
 using, e.g., encirclement or a connecting line, never 97  
 occurs in our data. Nor does the data show a single 99  
 case of plural spoken deictics, such as “these 101  
 books”. This may be due in part to the fact that 103  
 the placement of the individual objects on the walls  
 of Andersen’s study did not facilitate the making of  
 connections between them, and partly to the relative  
 scarcity of our data. Arguably, sooner or later, a  
 user might say, <pointing to the books on the  
 bookshelf> e.g., “Tell me about these books”. We

1 did not observe perceptual grouping behaviours,  
 3 e.g., using a deictic plural in speech, such as “these  
 5 pictures,” and selecting a single picture in a group of  
 7 pictures with a pointing gesture. This might be due  
 9 to several reasons. It was not demonstrated in the  
 11 simple multi-modal example the users were shown  
 at the start of the test. Another reason might be the  
 current layout of the graphical objects and the  
 richness of their perceptual properties (e.g., the  
 pictures) as compared to the 2D geometric shapes  
 investigated in [32].

13 As we explained in the analysis of the users’  
 15 multi-modal behaviour, users nearly always used  
 17 spoken deictics (pure demonstratives) rather than  
 19 actually naming the objects referred to, probably  
 21 because this was included in the short demonstra-  
 23 tion they had prior to the experiment and because  
 25 the recognition of deictics happened to work quite  
 27 well. They nevertheless also used a variety of  
 29 references that were not demonstrated (e.g., “who  
 31 is this woman?”), showing that they were able to  
 33 generalise to other kinds of references. This never-  
 35 theless raises the issue of natural vs. trained multi-  
 modality [37]. On the one hand, full natural multi-  
 modality (e.g., not showing any gesture or multi-  
 modal examples to the users prior to testing) will  
 probably lead to an even smaller proportion of  
 multi-modal behaviours than the one we observed.  
 On the other hand, trained multi-modality might  
 generate a larger variety of examples, such as  
 multiple-object gestures and implicit spoken refer-  
 ences without any deictics. We believe that the  
 approach we selected, i.e., that of demonstrating a  
 single example of a multi-modal input combination,  
 is a reasonable trade off between these two  
 extremes.

37 It follows that there are a serious number of  
 39 challenges ahead in order to be able to handle  
 41 natural interactive speech-gesture conversation,  
 including issues arising from the Andersen system,  
 such as:

- 43 1. the plural deictics/one object problem (the user  
 refers to several items in a single picture);
- 45 2. demonstratives may refer to spoken discourse as  
 well as to the visual environment;
- 47 3. addressing object details: a very demanding  
 proposition for developers;
- 49 4. addressing—potentially several—objects by a  
 (user-) stated criterion, such as “Can you show  
 51 me all the pictures from your fairytales?”
5. users may point at anything visible (and possibly

- ask as well); 53
6. users may meaningfully ask about, or comment  
 on, objects without pointing, as in “Who painted  
 the portrait of Jenny Lind?” 55
7. using visible objects as illustrations in spoken  
 discourse. 57

59 However, as regards the children who partici-  
 61 pated in the Prototype-2 user test, only Point 5  
 63 posed a significant problem, whereas Points 1 and 3  
 65 posed minor problems. Points 4, 6 and 7 never  
 occurs in the data whereas Point 2 occurs a few  
 times.

## 6. Conclusions 67

69 In this paper, we have described the modules that  
 71 we have developed for processing gesture and multi-  
 modal input in the Andersen system, as well as their  
 evaluation with two different groups of young users.  
 We have identified the causes of the most frequent  
 module failures, i.e., end of speech management in  
 the speech recogniser, gestures on non-referenceable  
 objects, and input gesturing while the character is  
 preparing to speak. We have suggested possible  
 improvements for removing these errors, such as  
 improvement of graphical and non-verbal afford-  
 79 dance, and proper management of end of speech  
 messages by the speech recogniser. 81

83 The Andersen project described in this paper has  
 85 provided data on how children gesture and combine  
 their gesture with speech when conversing with a 3D  
 character. Below, we revisit the issues that were  
 raised in the introduction.

87 How do children combine speech and gesture?  
 89 They do so more or less like adults do but (i)  
 probably in a slightly simpler fashion and (ii) only if  
 they are first-language speakers of the language  
 used for interaction with the ECA. 91

93 Would children avoid using combined speech and  
 gesture if they can convey their communicative  
 intention in a single modality? No, not if they are  
 first-language speakers of the language used in the  
 interaction; but yes, if the language of interaction is  
 their second language. 97

99 Is their behaviour dependent upon whether they  
 use their mother tongue or a second language? This  
 seems likely to be the case, but we need more data  
 analysis for confirmation. 101

103 To what extent would the system have to check  
 for semantic consistency between the speech and the  
 perceptual features of the object(s) gestured at? We



1 observed that the recognition and understanding of  
 3 spoken deictics was quite robust in the system and  
 5 that spoken deictics were nearly always used in  
 7 multi-modal input. We also observed behaviour in  
 9 which there was semantic inconsistency between the  
 11 speech and the perceptual features of the gestured  
 13 object. One user would ask “Who is this woman?”  
 15 when pointing to the picture of a man. This man is  
 17 wearing old-fashioned clothes and the picture,  
 19 which is in the corner of the room, might be less  
 21 visible than the other pictures. Another user would  
 23 say, “What is this?” when pointing to a picture  
 25 showing the picture of Andersen’s mother. We  
 27 might have expected “Who is this?” Finally, the  
 29 difficulties of speech recognition observed show that  
 31 it was better for the system to primarily trust the  
 33 gesture modality as it appeared, and was expected,  
 35 to be more robust than the speech. Since this paper  
 37 focused on gesture and combined speech-gesture in  
 39 the Prototype-2 user tests, we have not discussed the  
 41 speech processing findings made in those tests.  
 Suffice it here to say that the percentage of perfect  
 speech recognition was 23% for the Danish users  
 and 33% for the English users, whereas the  
 percentages for perfect gesture recognition and  
 interpretation were in the range of +90% for both  
 user groups. The system’s 2000 words speech  
 recogniser vocabulary was adequate for recognising  
 and understanding the spoken parts of the users’  
 multi-modal input despite the fact that the vocabu-  
 lary had been developed on the basis of spoken-  
 input-only corpora.

33 How do we evaluate the quality of such systems?  
 In this paper, we have used standard evaluation  
 35 methodologies, technical as well as usability-related,  
 37 for assessing the quality of the design solutions  
 adopted for gesture and combined speech-gesture  
 input processing. The solutions themselves represent  
 39 relatively complex trade-offs within the, still par-  
 41 tially uncharted, design space for multi-modal  
 speech/gesture input systems.

Some more specific evaluation methodologies  
 43 have also been considered in the literature. For  
 45 example, in their book dedicated to the evaluation  
 of ECAs, [38] point out the difference between  
 47 micro-level evaluation focused on a single feature of  
 the ECA and macro-level evaluation focused on the  
 global contribution of the ECA to an application. In  
 49 the same book, [38] provide a taxonomy of macro-  
 level dimensions to evaluate in an ECA, such as  
 51 believability or sociability, with corresponding  
 evaluation criteria. Another evaluation issue con-

cerns the target users of the Andersen system, i.e. 53  
 children and teenagers, who may require some 55  
 specific methods to optimise the data collection. In 57  
 this respect, [39] recommend methods, such as 59  
 thinking aloud, peer tutoring or user diaries in 61  
 order to access children’s mental model and 63  
 unbiased comments on a system. The authors also 65  
 point out the inadequacy of using some methods 67  
 with children, such as the use of focus groups. 69  
 Finally, the context of a game application raises 71  
 additional evaluation issues in the Andersen project, 73  
 because a game has to be usable and challenging at 75  
 the same time in order to be entertaining [40,41]. 77  
 Computer games can be evaluated by complemen- 79  
 tary means, such as classical usability methods, 81  
 psycho-physiological measures and behavioural 83  
 analysis [42–44]. However, among all these meth- 85  
 odologies—for evaluation of ECAs, doing tests with 87  
 children, and evaluating computer games—none 89  
 especially focus on investigating multi-modal input. 91  
 Therefore, we chose to rely on classical methods for 93  
 this particular topic, and we might draw on those 95  
 specific methods for evaluating other dimensions of 97  
 the Andersen system, e.g., Andersen’s believability 99  
 and entertainment qualities. 101

What do the users think of ECA systems 103  
 affording speech and gesture input? They clearly 105  
 like to use the touch screen and they very much 107  
 appreciate the idea of combined speech-gesture 109  
 input even if they do not massively practice 111  
 combined speech-gesture input when the language 113  
 of interaction is not their first language. Speech and 115  
 gesture input is, indeed, a “natural multi-modal 117  
 compound” for ECA systems. 119

How to manage temporal relations between 121  
 speech input, gesture input and multi-modal out- 123  
 put? We have proposed algorithms for managing 125  
 the temporal dimension and provided an illustration 127  
 of the multiple considerations involved when the 129  
 system is large and complex. According to our 131  
 evaluation, as reported above, the algorithms 133  
 proved suitable for the management of the users’ 135  
 behaviour. 137

The data we have collected clearly needs to be 139  
 complemented by data obtained with behaviours in 141  
 other multi-modal conversational contexts, possibly 143  
 more complex regarding graphical affordance for 145  
 multi-modal behaviour, such as many different 147  
 types of graphical objects, complex occlusion 149  
 patterns, etc. This might elicit more ambiguous 151  
 gesture semantics requiring the management of 153  
 gesture confidence scores, speech confidence scores

1 being notoriously unreliable for many important  
2 purposes.

3 In the current state of the art in the field of  
4 embodied conversational agents, Andersen is prob-  
5 ably one-of-a-kind. We know of no other running  
6 system, which integrates solutions to the challenges  
7 listed in Section 1.1. There is a sense in which the  
8 Andersen system is simply a computer game with  
9 spontaneous spoken interaction between the user  
10 and the character. This field of interactive spoken  
11 computer games was close to non-existent when the  
12 NICE project began. Spoken *output* in computer  
13 games was commonplace when the project began,  
14 however. Today, several computer games offer  
15 spoken input command words, which make a game  
16 character perform some action. So far, these  
17 products do not seem terribly popular with the  
18 games reviewers, probably because they typically  
19 assume that the game player is able to learn,  
20 sometimes quite large, numbers of spoken com-  
21 mands, and because their speech recognition and  
22 understanding is too fragile as well. We are not  
23 aware of any interactive spoken computer game  
24 products in the market. This is hardly surprising.  
25 Viewed from the perspective of the Andersen  
26 system, it may be too early to offer customers  
27 interactive spoken computer games in the standard  
28 sense of the term “computer game,” knowing that a  
29 computer game is being used, on average, for  
30 30–50 h of game-playing. By contrast, the Andersen  
31 system addresses the more modest challenge of  
32 providing edutaining conversation with a new user  
33 every 5–20 min.

## 37 7. Uncited references

39 [45]; [46].

## 43 Acknowledgements

45 We gratefully acknowledge the support for the  
46 NICE project by the European Commission’s Hu-  
47 man Language Technologies Programme, Grant  
48 IST-2001-35293. We would also like to thank all  
49 participants in the NICE project for the three  
50 productive years of collaboration that led to the  
51 running system prototypes presented in this paper.

## References

- 53
- [1] R.A. Bolt, “Put-that-there”: voice and gesture at the  
54 graphics interface, Seventh Annual International Conference  
55 on Computer Graphics and Interactive Techniques, ACM,  
56 Seattle, Washington, US, 1980, pp. 262–270.
- [2] N.O. Bernsen, L. Dybkjær, Evaluation of Spoken Multi-  
57 modal Conversation. Sixth International Conference on  
58 Multimodal Interaction (ICMI’2004), Association for Com-  
59 puting Machinery (ACM), New York, 2004, pp. 38–45.
- [3] S.L. Oviatt, Multimodal interfaces. Human–computer inter-  
60 action handbook: fundamentals, in: J. Jacko, A. Sears  
61 (Eds.), *Evolving Technologies and Emerging Applications*,  
62 vol. 14, Lawrence Erlbaum Assoc., Mahwah, NJ, 2003, pp.  
63 286–304.
- [4] R. Sharma, M. Yeasin, N. Krahnstoeber, I. Rauschert, G.  
64 Cai, I. Brewer, A. MacEachren, K. Sengupta, Speech–ges-  
65 ture driven multimodal interfaces for crisis management,  
66 Proc. IEEE VR2004 91 (9) (2003) 1327–1354 [http://](http://spatial.ist.psu.edu/cai/2003-Gesture-speech-interfaces-for%20crisis-management.pdf)  
67 [spatial.ist.psu.edu/cai/2003-Gesture-speech-interfaces-](http://spatial.ist.psu.edu/cai/2003-Gesture-speech-interfaces-for%20crisis-management.pdf)  
68 [for%20crisis-management.pdf](http://spatial.ist.psu.edu/cai/2003-Gesture-speech-interfaces-for%20crisis-management.pdf).
- [5] R. Catizone, A. Setzer, Y. Wilks, Multimodal Dialogue  
69 Management in the COMIC Project. EACL 2003 Workshop  
70 on Dialogue Systems: Interaction, Adaptation, and Styles of  
71 Management, 2003, [http://www.hcrd.ac.uk/comic/docu-](http://www.hcrd.ac.uk/comic/documents/publications/eaclCOMICFinal.pdf)  
72 [ments/publications/eaclCOMICFinal.pdf](http://www.hcrd.ac.uk/comic/documents/publications/eaclCOMICFinal.pdf).
- [6] M. Johnston, Unification-based multimodal parsing, in:  
73 17th International Joint Conference of the Association for  
74 Computational Linguistics, Montreal, Canada, August  
75 Association for Computational Linguistics Press, Morgan  
76 Kaufmann Publishers, Los Altos, CA, 1998, pp. 624–630.
- [7] M. Johnston, P. Cohen, D. McGee, S. Oviatt, J. Pittman, I.  
77 Smith, Unification-based Multimodal Integration, ACL’97,  
78 1997.
- [8] L. Almeida, I. Amdal, N. Beires, M. Boualem, L. Boves, E.  
79 Os, P. Filoche, R. Gomes, J.E. Knudsen, K. Kvale, J.  
80 Rugelbak, C. Tallec, N. Warakagoda, The MUST Guide to  
81 Paris; Implementation and expert evaluation of a multi-  
82 modal tourist guide to Paris. Multi-Modal Dialogue in  
83 Mobile Environments, ISCA Tutorial and Research Work-  
84 shop (IDS’2002), Kloster Irsee, Germany, June 17–19 [http://](http://www.isca-speech.org/archive/ids_02)  
85 [www.isca-speech.org/archive/ids\\_02](http://www.isca-speech.org/archive/ids_02), 2002.
- [9] M. Johnston, S. Bangalore, Multimodal Applications from  
86 Mobile to Kiosk. W3C Workshop on Multimodal Interac-  
87 tion, Sophia Antipolis, France, 19–20 July 2004, 2004 [http://](http://www.w3.org/2004/02/mmi-workshop/papers)  
88 [www.w3.org/2004/02/mmi-workshop/papers](http://www.w3.org/2004/02/mmi-workshop/papers)
- [10] S. Oviatt, Multimodal interactive maps: designing for  
89 human performance, *Hum. Comput. Interact.* 12 (1997)  
90 93–129.
- [11] A.D. Milota, Modality Fusion For Graphic Design Appli-  
91 cations, ICMI’2004, 2004.
- [12] P. Gieselmann, M. Denecke, Towards multimodal interac-  
92 tion with an intelligent room. Eighth European Conference  
93 On Speech Communication and Technology (Euro-  
94 speech’2003), Geneva, Switzerland, September 1–4, 2003,  
95 [http://isl.ira.uka.de/fame/publications/FAME-A-WP10-](http://isl.ira.uka.de/fame/publications/FAME-A-WP10-007.pdf)  
96 [007.pdf](http://isl.ira.uka.de/fame/publications/FAME-A-WP10-007.pdf).
- [13] J. Juster, D. Roy, Elvis: situated speech and gesture  
97 understanding for a robotic chandelier. Sixth International  
98 Conference on Multimodal Interfaces (ICMI’2004), October  
99 100

- 1 13–15, State College, Pennsylvania, USA, ACM, New York, 2004, pp. 90–96.
- 3 [14] S.L. Oviatt, R. Coulston, S. Tomko, B. Xiao, R. Lunsford, M. Wesson, L. Carmichael, Toward a theory of organized multimodal integration patterns during human–computer interaction, in: International Conference on Multimodal Interfaces (ICMI'2003), ACM Press, Vancouver, BC, 2003, pp. 44–51 [http://www.cse.ogi.edu/CHCC/Publications/toward\\_theory\\_organized\\_multimodal\\_integration\\_oviat.pdf](http://www.cse.ogi.edu/CHCC/Publications/toward_theory_organized_multimodal_integration_oviat.pdf).
- 5 [15] W.C. Avaya, D. Dahl, M. Johnston, R. Pieraccini, D. Ragget, EMMA: Extensible MultiModal Annotation markup language. W3C Working Draft 14 December 2004, W3C. <http://www.w3.org/TR/emma/>
- 7 [16] E. Kaiser, A. Olwal, D. McGee, H. Benko, A. Corradini, X. Li, P. Cohen, S. Feiner, Mutual disambiguation of 3D multimodal interaction in augmented and virtual reality, in: Fifth International Conference on Multimodal Interfaces (ICMI'03), ACM Press, Vancouver, British Columbia, Canada, 2003, pp. 12–19 <http://www1.cs.columbia.edu/~aolwal/projects/maven/maven.pdf>.
- 9 [17] J. Cassell, J. Sullivan, S. Prevost, E. Churchill, Embodied Conversational Agents, MIT Press, Cambridge, MA, 2000, 0-262-03278-3.
- 11 [18] W.L. Johnson, J.W. Rickel, J.C. Lester, Animated pedagogical agents: face-to-face interaction in interactive learning environments, *Int. J. Artif. Intell. Educ.* 11 (2000) 47–78 <http://www.csc.ncsu.edu/eos/users/l/lester/www/imedia/apajiaied-2000.html>.
- 13 [19] D. Traum, J. Rickel, Embodied Agents for Multi-party Dialogue in Immersive Virtual Worlds, First International Joint Conference on “Autonomous Agent and Multiagent Systems” (AAMAS'02), July 15–19, Bologna, Italy, ACM Press, New York, 2002, pp. 766–773.
- 15 [20] T. Sowa, S. Kopp, M.E. Latoschik, A Communicative Mediator in a Virtual Environment: Processing of Multimodal Input and Output, In: Proc. of the International Workshop on Information Presentation and Natural Multimodal Dialogue, Verona, Italy, 2001, pp. 71–74, <http://www.techfak.uni-bielefeld.de/~skopp/download/CommunicativeMediator.pdf>.
- 17 [21] D. Hofs, H.J.A. op den Akker, A. Nijholt, A generic architecture and dialogue model for multimodal interaction. 1st Nordic Symposium on Multimodal Communication, Copenhagen, Denmark, 25–26 September 2003, pp. 79–92.
- 19 [22] S. Narayanan, A. Potamianos, H. Wang, Multimodal systems for children: building a prototype, Sixth European Conference on Speech Communication and Technology (Eurospeech'99), Budapest, Hungary, September 5–9, 1999.
- 21 [23] D. Perzanowski, A.C. Schultz, W. Adams, E. Marsh, M. Bugajska, Building a multimodal human-robot interface, *IEEE Intell. Syst.* 16 (1) (2001) 16–21.
- 23 [24] H. Holzapfel, K. Nickel, R. Stiefelhagen, Implementation and Evaluation of a Constraint-Based Multimodal Fusion System for Speech and 3D Pointing Gestures. ICMI 2004, 2004, <http://isl.ira.uka.de/fame/publications/FAME-A-WP10-028.pdf>.
- 25 [25] S. Oviatt, C. Darves, R. Coulston, Toward Adaptive Conversational Interfaces: Modeling Speech Convergence with Animated Personas, 2004, [http://www.cse.ogi.edu/CHCC/Publications/TOCHI\\_Oviatt\\_MAI04-503.pdf](http://www.cse.ogi.edu/CHCC/Publications/TOCHI_Oviatt_MAI04-503.pdf)
- 27 [26] S.L. Oviatt, B. Adams, Designing and evaluating conversational interfaces with animated characters, in: J. Cassell, J. Sullivan, S. Prevost, E. Churchill (Eds.), Embodied Conversational Agents, MIT Press, Cambridge, MA, 2000, pp. 319–345. 53
- 29 [27] K. Ryokai, C. Vaucelle, J. Cassell, Virtual peers as partners in storytelling and literacy learning, *J. Comput. Assist. Learn.* 19 (2003) 195–208. 55
- 31 [28] J. Read, S. MacFarlane, C. Casey, Oops! silly me! errors in a handwriting recognition-based text entry interface for children. Nordic Conference on Human-Computer Interaction (NordCHI '02), 2002, pp. 35–40. 57
- 33 [29] N.O. Bernsen, M. Charfuelàn, A. Corradini, L. Dybkjær, T. Hansen, S. Küllerich, M. Kolodnytsky, D. Kupkin, M. Mehta, First prototype of conversational H.C. Andersen. International Working Conference on Advanced Visual Interfaces (AVI'2004), Gallipoli, Italy, May 2004, ACM, New York, 2004, pp. 458–461. 61
- 35 [30] S. Buisine, J.-C. Martin, N.O. Bernsen, Children's Gesture and Speech in Conversation with 3D Characters. HCI International 2005, Las Vegas, USA, 22–27 July 2005. 63
- 37 [31] E. Lewin, “KTH Broker, 1997, <http://www.speech.kth.se/broker/> 65
- 39 [32] F. Landragin, N. Bellalem, L. Romary, Visual salience and perceptual grouping in multimodal interactivity. First International Workshop on Information Presentation and Natural Multimodal Dialogue, Verona, Italy, 2001, pp. 151–155, <http://www.loria.fr/~landragi/publis/ipnmd.pdf>. 67
- 41 [33] S. Oviatt, A. De Angeli, K. Kuhn, Integration and synchronization of input modes during multimodal human–computer interaction, in: Human Factors in Computing Systems (CHI'97), ACM Press, New York, 1997, pp. 415–422. 69
- 43 [34] S. Buisine, J.-C. Martin, Children's and Adults' Multimodal Interaction with 2D Conversational Agents. CHI'2005, Portland, Oregon, 2–7 April 2005. 71
- 45 [35] N.O. Bernsen, L. Dybkjær, User evaluation of Conversational Agent H. C. Andersen, Ninth European Conference on Speech Communication and Technology (InterSpeech'2005), Lisboa, Portugal, 2005. 73
- 47 [36] S. Buisine, J.-C. Martin, Experimental evaluation of bidirectional multimodal interaction with conversational agents, Proceedings of the Ninth IFIP TC13 International Conference on Human–Computer Interaction (INTERACT'2003), Zürich, Switzerland, September 1–5, IOS Press, 2003, pp. 168–175, <http://www.interact2003.org/>. 75
- 49 [37] J. Rugelbak, K. Hamnes, Multimodal Interaction—Will Users Tap and Speak Simultaneously? *Elektronikk*, 2003, [http://www.eurescom.de/~ftproot/web-deliverables/public/P1100-series/P1104/Multimodal\\_Interaction\\_118\\_124.pdf](http://www.eurescom.de/~ftproot/web-deliverables/public/P1100-series/P1104/Multimodal_Interaction_118_124.pdf). 77
- 51 [38] K. Ibister, P. Doyle, The blind men and the elephant revisited, in: Z. Ruttkay, C. Pelachaud (Eds.), From Brows to Trust: Evaluating Embodied Conversational Agents, Kluwer Academic Publishers, Dordrecht, 2004, pp. 3–26. 79
- [39] S. MacFarlane, J. Read, J. Höysniemi, P. Markopoulos, Evaluating interactive products for and with children. Tutorial Notes, Interact'2003 Conference, 2003. 81
- [40] D. Johnson, J. Wiles, Effective affective user interface design in games, *Ergonomics* 46 (2003) 1332–1345. 83
- [41] K. Keeker, R. Pagulayan, J. Sykes, N. Lazzaro, The untapped world of video games, CHI'2004, 2004, 1610–1611, 85
- [42] S. Kaiser, T. Wehrle, S. Schmidt, Emotional episodes, facial expressions, and reported feelings in human-computer 87
- 91 89
- 93 91
- 95 97
- 97 99
- 101 103

- 1 interactions, Proceedings of Conference of the International  
2 Society for Research on Emotions, 1998, pp. 82–86.
- 3 [43] N. Lazzaro, K. Keeker, What’s my method? A game show  
4 on games, in: Proceedings of CHI’2004, 2004, pp.  
5 1093–1094.
- 6 [44] R. Pagulayan, K. Keeker, D. Wixon, R.L. Romero, T.  
7 Fuller, User-centered design in games, in: J.A. Jacko, A.  
8 Sears (Eds.), The Human–Computer Interaction Handbook:  
9 Fundamentals, Evolving Technologies and Emerging Appli-  
10 cations Archive, Lawrence Erlbaum Associates, Inc., Mah-  
11 wah, 2002, pp. 883–906.
- [45] Z. Ruttkay, C. Pelachaud, From Brows to Trust—Evaluat-  
ing Embodied Conversational Agents, Kluwer, 1-4020-2729-  
X, 2004, [http://wwwhome.cs.utwente.nl/~zsofi/Kluwer-  
Book.htm](http://wwwhome.cs.utwente.nl/~zsofi/Kluwer-Book.htm). 13
- [46] B. Xiao, C. Girand, S.L. Oviatt, Multimodal Integration  
Patterns in Children, in: J. Hansen, B. Pellom (Ed.),  
15 Proceedings of the Seventh International Conference on  
16 Spoken Language Processing (ICSLP’2002), Denver, CO,  
17 Casual Prod. Ltd., Sept. 2002, pp. 629–632. Abstract, [http://  
18 www.cse.ogi.edu/CHCC/Publications/multimodal\\_integra-  
19 tion\\_patterns\\_in\\_children\\_xiao.pdf](http://www.cse.ogi.edu/CHCC/Publications/multimodal_integration_patterns_in_children_xiao.pdf). 20  
21

UNCORRECTED PROOF