



## 2 A Tangible Multimodal Dialog Scenario

Our experimental scenario attempts to combine the benefits of both physical and digital worlds in a mixed-reality setting by targeting an in-store scene, but augmented by instrumented devices like a Personal Digital Assistant (PDA) and a shopping trolley with a mounted display. Whereas the PDA is used as a communication channel through which users can associate directly with the products rather than through a sales assistant, the shopping trolley is capable of offering shopping advice based on its current contents. An in-store setting encompasses the down-to-earth basics that only a traditional store in a real world and with real physical products can provide such as the sense of touch. When instrumented, it further provides the convenience inherent in digital worlds such as ubiquitous information access. The unification of these two worlds is achieved through a Tangible MultiModal interface (TMM) that is seamlessly integrated into existing shopping practices. TMM's are now being incorporated into a wide-range of fields, up to an including safety-critical applications [Cohen and McGee 2004].



Fig. 1. Anthropomorphized object initiating a dialog.

Tangible User Interfaces (TUI's) [Ullmer and Ishi 2001] couple physical representations (e.g., spatially manipulable physical objects) with digital representations (e.g., graphics and audio), yielding interactive systems that are computationally mediated. In our scenario, we use an intuitive "one-to-one" mapping between physical shopping items on the shelf and elements of digital information. The spa-

tial relation of a physical token partially embodies the dialog state, which can be seen in our example in that a product can be either in a product shelf, in the shopping trolley or outside of these containers. The position of the product is mapped to a physical action of the user, where the physical movements of the artifacts serve as a means to controlling the dialog state.

The Mobile ShopAssist (MSA) is a demonstrator that aids users in product queries and comparisons. The goal is to provide rich symmetric multimodal interaction and the ability for users to converse directly with the products. Using the MSA, a shopper interested in buying a digital camera would for example walk up to a shelf and synchronize its contents with their PDA. After synchronization, they may ask a product about its attributes (e.g. “*What is your optical zoom?*”), or even compare multiple products together (e.g. “*<gesture> Compare yourself with this camera <gesture>*”). Comparisons may be made between products from the physical world, digital world, or a mixture of both (i.e. mixed-reality).

When interacting with the digital cameras, the user may decide to communicate indirectly with the object “*What is the price of this camera <gesture>*”, or directly “*What is your price?*”. The input modalities available to the user include speech, handwriting, gesture, and combinations thereof. It is direct interaction and the concept of anthropomorphization (i.e. assigning inanimate objects human-like characteristics, see section 5) that we focus on. Assuming the user has chosen to interact directly with the objects, the objects will in return communicate directly with the user and may also initiate mini-dialogs when picked up or put down on a shelf or in a shopping trolley, similar to (1) in Fig. 1. Once the user has finished conversing with the products, they may decide to buy the product, or to simply take the information that they have downloaded back home with them to think about later on. The objects (not limited to digital cameras), can then be placed into the shopping trolley and taken to the cashier. On request, the user’s interactions are logged and summarized in a personal shopping diary [Kröner et al. 2004].

The MSA is a mixed-initiative dialog system, which means that both a product and the user can start a dialog or take the initiative in a sub-dialog. For instance, when the product is picked up – and no accompanying user query is issued – the product will introduce itself. Another system-initiated dialog phase is that of cross-selling, which occurs when a product is placed into the shopping trolley. Such a dialog might give advice on accessories available for the product, for example: “*You may also find the NB-2LH batteries in the accessories shelf to be useful*”.

Instrumented environments containing RFID tagged products have till now primarily benefited the retailer through improved inventory management and tracking. Our scenario also highlights user benefits in the form of comparison shopping, cross-selling recommendations, and product information retrieval based on real physical indexes.

### 3 Instrumented Environment Infrastructure

The main infrastructure components that exist in our shopping environment include the mobile device which is used as a communication channel, the containers (e.g. shelves, trolley) and the objects (i.e. shopping products) belonging to a shop (see Fig. 3). Each shelf is identified by an infrared beacon that is required when a user synchronizes the shelf's product data. The products are identified through the use of passive RFID tags, which allow a product to be classified as being in or out of a container. Each container has an RFID antenna and a reader connected to it, and this allows the shelves and shopping trolley to recognize when products are put in or taken out of them.

The instrumented shelves may be scattered over several rooms, and communicate via a WLAN connection with the ambient intelligence server, as shown in Fig. 2 (similar instrumented shopping environments without a tangible multimodal interface are the Metro Future Store and MyGrocer [Kourouthanasis et al. 2002] . It is this server that maintains the product database, and the event heap [Fox et al. 2000], which is used for recording extra-gesture events. As described in [Butz et al. 2004], a searchlight in the form of a steerable projector further allows the system to find and highlight products based on optical markers. This is important in establishing a link between the physical products and their digital counterparts (and vice-versa), which do not need to be sorted in the same way. Such a situation could for example arise when digital objects are re-sorted based on specific product features such as price or manufacturer, instead of their physical location.

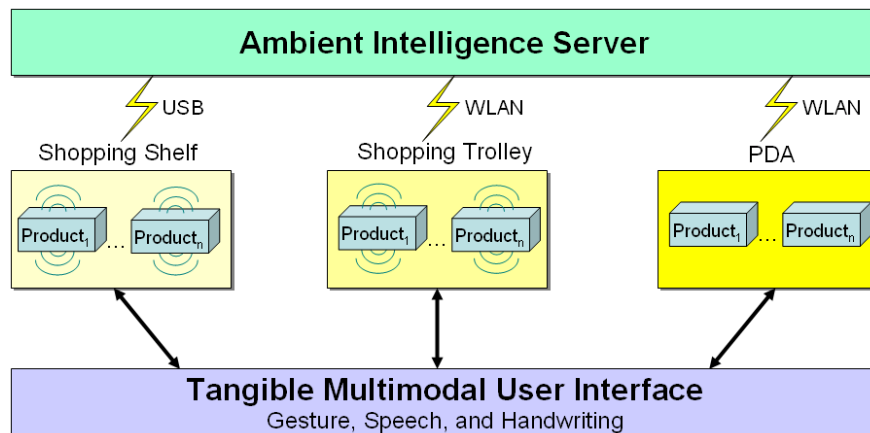


Fig. 2. Distributed architecture of the MSA.

After a client device such as a PDA has been synchronized with a shelf, it will maintain its own blackboard of events, on which it stores not only the extra-gesture interactions broadcast by the server, but also speech, handwriting, and intra-gesture interactions that the PDA is capable of recognizing and interpreting lo-

cally. When a shopping trolley is added to the scenario [Schneider, 2004], the contained products are listed on a trolley mounted display, and the trolley will offer advice on additional products that may be relevant to the user.

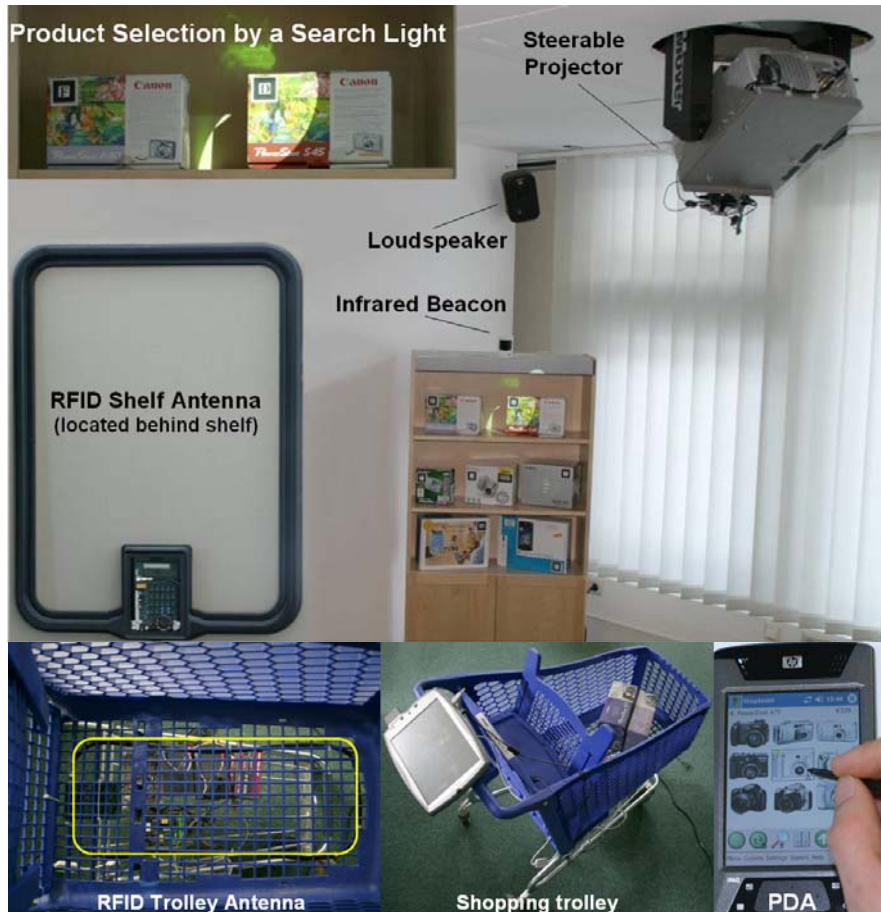


Fig. 3. Instrumented shopping environment.

The data downloaded upon shelf synchronization is contained within the product database. This is located on the ambient intelligence server, and contains product feature-value lists for attributes like “price”, “optical zoom” and “mega pixels”. The database also contains images, links to URL manufacturer sites, RFID and optical marker values for the products, and a reference to the associated grammar file used for input recognition. This data is retrieved by SQL queries and transferred from the server to the mobile device in XML format. The input grammar files contain a similar feature-value list, in which grammar entries for each feature are defined for the different modalities like speech and handwriting. The

input grammars are assigned to a group of products based on their product type, which allows multiple products to share a single grammar file, as is the case for the product type “digital camera”.

After a user has synchronized with a shelf and starts to browse through the products, internal data representations are created for both the objects, and feature keywords displayed on the PocketPC display. This representation is for example used by the modality of intra-gesture, first during input interaction as screen coordinates are mapped to underlying graphical objects and visual WCIS keywords currently on the screen, and later during output presentation through the use of object and keyword lookup functions, which locate a particular reference and display it on the screen as a selected referent.

## 4 Symmetric Multimodal Interaction

As defined in [Wahlster, 2003], symmetric multimodality refers to the ability of a system to use all input modes as output modes, and vice-versa. Empirical studies have shown that the robustness of multimodal interfaces increase substantially as the number and heterogeneity of modalities expand [Oviatt, 2002]. Information provided by one or more sources, can be used to resolve ambiguities or manage recognition and sensor uncertainties in another modality, thereby reducing errors both in the system’s interpretation of the user’s input, or the user’s understanding of the system’s output. Whereas modality fusion maps multimodal input to a semantic representation language, the modality fission component provides the inverse functionality of the modality fusion component, since it maps a communicative intention of the system onto a coordinated multimodal presentation.

Most of the previous multimodal interfaces (see [Oviatt and Wahlster 1997] and [Maybury and Wahlster 1998]) do not support symmetric multimodality, since they focus either on multimodal fusion (e.g. QuickSet [Cohen et al. 1997], and MATCH [Johnston et al. 2002]) or multimodal fission (e.g. WIP [Wahlster et al. 1993]). Symmetric multimodal dialog systems like SmartKom and the MSA create a natural experience for the user in the form of daily human-to-human communication, by allowing both the user and the system to combine the same spectrum of modalities for the input as for the output. The MSA represents a new generation of multimodal dialog systems that deal not only with simple modality integration and synchronization, but covers the full spectrum of dialog phenomena that are associated with symmetric multimodality. Symmetric multimodality supports the mutual disambiguation of modalities, as well as multimodal or crossmodal deixis and anaphora resolution.

#### 4.1 Base Modalities

Multimodal interaction in the MSA is based on the modalities: speech, handwriting and gesture, whereby gesture can be further grouped into the types intra and extra. Intra-gestures refer to product and feature selections on the display of the PDA (“intra\_point”), while extra-gestures refer to actions in the physical real world such as picking an object up from a shelf (“extra\_pick\_up”), or putting an object back onto a shelf (“extra\_put\_down”).

From this limited number of base modalities, a wide range of mixed and overlapped input combinations can be formed. [Wasinger and Krüger 2004] outline a total of 23 input modality combinations that were tested within a laboratory setting for use with the system. The modalities included both unimodal (e.g. speech-only) and multimodal (e.g. speech-gesture) combinations, as well as overlapped and non-overlapped modality combinations. Overlapped modality combinations are ones in which (possibly conflicting) information is provided multiple times in potentially different modalities, as seen in the following non-conflicting speech-gesture overlapped feature interaction: “*What is the price <intra\_gesture=price> of the EOS10D?*”. Such redundant information is useful for reference resolution.

All of these input modality combinations are however only one side of the interaction equation. The flip-side encompasses the output modalities used by the anthropomorphized objects when replying to the user. Speech output for example is presented to the user via an embedded synthesizer. We currently use two synthesizers, one is a formant synthesizer which requires a small memory footprint (around 2MB per language), while the other is a high quality concatenative synthesizer that has a much larger footprint (between 7 and 15MB per language for a single voice). Although the formant synthesizer sounds robotic, it provides far greater flexibility in manipulating voice characteristics such as age and gender, which is important in providing the anthropomorphized objects with their own personality (see section 5.2). The output equivalent to handwriting is the use of system fonts that are displayed in a predefined location on the PDA’s display. Intra-gesture output for object selection is achieved by drawing a border around the selected object, while intra-gesture output for feature selection is achieved by highlighting the active keyword within the visual What-Can-I-Say (WCIS) text bar, which scrolls across the bottom of the PDA’s display. Extra-gesture output is made possible through the use of a steerable projector, which provides for real-world product selection by placing the product under a spot light. Fig. 4 shows the use of the primary modalities within our system, for both input interaction and output presentation. This graphic also shows how objects and features can be referred to within the modality types. The output for intra-gesture for example shows a selected feature and below it a selected object.









	User Input	System Output
Speech:		
Handwriting:		
Intra Gesture:		<p>⏪   (next   previous) page.&lt;root&gt; = what can i si price   mega pixels   optical zoom   focal length  </p> 
Extra Gesture:	<p>pick_up, and put_down</p> 	<p>Searchlight:</p> 

Fig. 4. The modalities used for both input and output.

## 4.2 Symmetric Modality Combinations

Systems that support multimodal interaction such as speech, handwriting, and gesture, require an efficient means of fusing the interactions together to form a single unambiguous dialog result, which can then be passed onto subsequent modules in the system such as a retrieval component. Multimodal user input interaction within our system generally consists of a single feature and one or more object references, for example: “*What is your price <gesture=PowerShot S70>?*”. Valid values for the feature tag include (in reference to digital cameras) ‘price’, ‘optical zoom’ and ‘mega pixels’, while valid values for the object tag include ‘PowerShot S70’ and ‘CoolPix 4300’. Before such interactions can be parsed however, they must first be converted into a modality-free language, where the maximum number of objects is currently limited to two due to the limited space available on the PDA display. This language is formatted in XML and closely resembles the W3C EMMA standard (see [www.w3.org/TR/emma/](http://www.w3.org/TR/emma/)) in that each tag (i.e. FEATURE and OBJECT) contains a number of attributes like the modality type, timestamp, confidence value, and N-best list values.

On the flip-side, multimodal output from our anthropomorphized products must provide the resulting value information alongside reproducing the feature and object information, and be flexible enough to cater for both direct interaction: “*My price is €500*”, and indirect interaction: “*PowerShotS50, price, €500*”.

Fig. 5 summarizes the potential range of modality combinations that exist for user input and anthropomorphized object output, when the modalities Speech (S), Handwriting (H), Intra-Gesture (I) and Extra-Gesture (E) are available. For input alone, possible multimodal combinations can be seen to include: SS, SH, SI, SE, HS, HH, HI, HE, IS, IH, II and IE. This figure does not consider multiple object referents, or overlaid input, which would create an even larger number of modality combinations to choose from. In this diagram, the interaction manager is responsible for recognizing and interpreting user interactions with the system, while the presentation planner is responsible for coordinating output for presentation back to the user. This output must be consistent not only in providing the correct information in response to user queries, but also in the choice of modality combinations that are used to present the information.

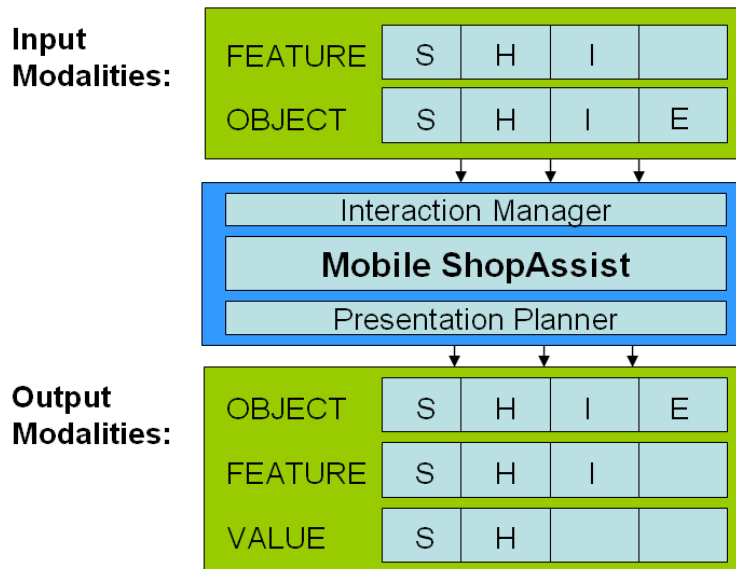


Fig. 5. Symmetric modality matching in the MSA.

### 4.3 Output Modality Allocation Strategy

The output strategy within the MSA uses speech as a base modality to present Object (O), Feature (F), and Value (V) information. Complementing speech is the modality of handwriting, which is used first to present the user with the transient information (O) and (F), and then a short time later with the non-transient information (V). Gesture is additionally used to show the selected (O) to the user as non-transient information, either solely via an intra-gesture on the PDA display,

but also as an extra-gesture via the Searchlight, if it is available. Intra-gesture output for the feature (i.e. highlighting an active keyword from the scrolling WCIS text) only occurs if the scrolling text is currently visible. At the end of an interaction dialog, a user will have been presented with the same information in two complete modalities (speech and handwriting), and part of the modality gesture.

In comparison to recent usability tests in which our subjects stated that they preferred non-overlapped input modalities to overlapped ones ( $\chi^2(2, N=27) > 24.889, p < 0.000$ ), our users were keen to be provided with overlapped modalities for the output. Redundancy in the output modalities as used above compensated for the names of objects such as “PowerShot S11S” and “EOS 10D” being pronounced incorrectly by the speech synthesizers, and for the transience required in presenting the written language as two separate events on the limited display space.

The current output strategy is just one of several possibilities. Other output allocation strategies include for example the exact replication of modalities used for input as for output (mimicking), user defined profiles, or profiles that limit the media to types that 3rd person parties can not observe (e.g. handwriting, intra-gesture, and speech output through a PDA-based headset), or that do not require a PDA (e.g. server-sided speech, and extra-gestures).

## 5 Anthropomorphized Products

In this section, we outline the concept of anthropomorphization. We describe the difference between direct and indirect interaction, and also outline how we account for anthropomorphized objects in the MSA, with particular focus on the language grammars, the product personalities, and the state-based object models that define when our objects may initiate dialog interaction with the user.

### 5.1 The Role of Anthropomorphization

Anthropomorphism is the tendency for people to think of inanimate objects as having human-like characteristics. Many early cultures made no distinction between animate and inanimate objects [Todd 2002]. Animism is looking at all nature as if it were alive. It’s one of the oldest ways of explaining how things work, when people have no good functional model. When users interact with ambient intelligence environments rather than with a desktop screen, there is a need for communication with a multitude of embedded computational devices in mass-marketed products. For human-environment interaction with thousands of networked smart objects, a limited animistic design metaphor seems to be appropriate [Nijholt et al. 2004, see also the chapter by Nijholt, de Ruyter, Heylen, and Priver in this book].

Although there are various product designs that use an anthropomorphic form (like the Gaultier perfume bottles that have the shape of a female torso), in the work presented here we stimulate anthropomorphization solely by the pretended conversational abilities of the products. Since the shopper's hands are often busy with picking up and comparing products, in many situations the most natural mode to ask for additional information about the product is the use of speech. When a product talks and answers the shopper's questions with its own voice, the product is being anthropomorphized.

There is a longstanding tradition among some HCI researchers against the use of anthropomorphism [Don 1992], because it may create wrong user expectations. This has led to taboos like "Don't use the first person in error messages". People are however used to dealing with disembodied voices on the telephone, and our empirical user studies also provide evidence that most shoppers have little concern about speaking with shopping items such as digital cameras (see section 6). In addition, through the world of TV commercials, shoppers are used to anthropomorphized products like "Mr. Proper", a liquid cleaning product that is morphed into an animated cleaning Superman or the animated "M&M" round chocolates.

Of course, anthropomorphized interaction can be irritating or misleading, but our system is designed in such a way that it presents its limitations frankly. The What-Can-I-Say (WCIS) mechanism in the MSA guides the user in their decision-oriented dialog and makes it clear that it has only restricted, but very useful communication capabilities. We contend that anthropomorphism can be a useful framework for interaction design in ambient intelligence environments, if its strengths and weaknesses are understood.

## 5.2 Adding Human-Like Characteristics

Distinct from the assortment of modality combinations available in the MSA, users may choose to interact either directly or indirectly with the shopping products. These products will in return also need to respond correspondingly. We derive the terms direct and indirect interaction from the mode of reference being made to the "person" segment of a dialog. In English for example, there exist the tenses: first person (the person speaking), second person (the person being spoken to), and third person (the person being spoken about). From an input perspective, direct interaction refers to the 2<sup>nd</sup> person (e.g. "*What is your price?*"), while indirect interaction refers to the 3<sup>rd</sup> person (e.g. "*What is the price of this/that camera?*"). From an output perspective, direct interaction (as used by the anthropomorphized objects) takes the 1<sup>st</sup> person (e.g. "*My price is €599*"), while indirect interaction takes the 3<sup>rd</sup> person (e.g. "*The price of this/that camera is €599*").

Within the MSA, grammar files exist for each product type, such as "digital camera", and for both English and German. These grammar files define the recognizable input (e.g. product and feature information) for the modalities handwriting, intra-gesture, and speech, whereby the handwriting and intra-gesture gram-

grams are identical for feature resolution. Although the individual modalities may be used to communicate complete dialog acts (i.e. product and feature information), speech is the only modality in which complete sentences may be used. Three forms of speech input are accepted by the system, namely “keyword” (i.e. speaking only the keyword, e.g. “*price*”), “indirect” (e.g. “*What is the price of <product>?*”), and “direct” (e.g. “*What is your price?*”). The grammar files for each of the product types are downloaded onto the PDA together with the product information, each time the user synchronizes with a particular shelf container. These files are then parsed by the PDA to create the individual grammars required for each of the recognizers.

Objects within the MSA are further personalized by one of five different formant synthesizer voice profiles (3 male, 2 female, and all adult), which are based on parameters such as gender, head size, pitch, roughness, breathiness, speed and volume. A limitation of our approach is that 5 different voices can not provide each product in a shelf (let alone an entire store) with a unique voice. An alternative would be to use pre-recorded audio samples for each product, but this would require different magnitudes of storage space. A different approach might be to allow the PDA to assign the voices to products, which would allow at least the first 5 products interacted with to have a unique voice. Such an approach would also allow the use of personality matching strategies to better market products to specific user groups. Dynamic voice assignment would however also create the need for storing voice to product mappings for future use, so that returning users are not faced with anthropomorphized objects with multiple personalities.

### 5.3 State-based Object Model

A further feature of our anthropomorphized objects is their ability to initiate interaction with the user when in a particular state (see Fig. 6). These states are based on variables such as a products location, a recent extra-gesture action, and an elapsed period of time. The location of a product may be either “in a shelf”, “out of a shelf”, or “in a shopping trolley”, and extra-gesture events include: “pick\_up” and “put\_down”. Thus, the physical acts of the user like “Pick\_Up (product007, shelf02)” and “Put\_Down (product007, trolley01)” are mapped onto dialog acts like “Activate\_Dialog\_With (product007)” and “Finish\_Dialog\_With (product007)” respectively. In this case, the Put\_Down action reflects a positive buying decision as the product was placed inside the trolley, but the product could just as equally have been put down on the shelf instead, thus reflecting a negative buying decision.

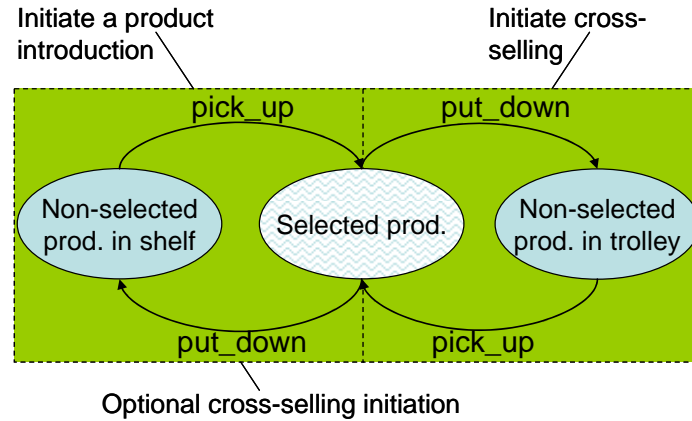


Fig. 6. Base product states used for object-initiated interaction.

As an example, an object will initiate a dialog interaction if it is picked up from the shelf for the first time and no further user interaction is observed within a 5 second time frame. Silence as a powerful form of communication is well documented [Knapp 2000], and in our case such silence forces the product to introduce itself (see (1) in Fig. 1). A product might also initiate an interaction when for example placed inside the shopping trolley, in order to alert the user of any further products (e.g. accessories) that they might be interested in purchasing (i.e. cross-selling).

## 6 Usability Study

Ben Shneiderman, as a prominent critic of anthropomorphized user interfaces, stated at a panel discussion documented by [Don, 1992] “I call on those who believe in the anthropomorphic scenarios to build something useful and conduct usability studies and controlled experiments”. That is exactly what we have done in the described research.

In this section, we describe an empirical field study on user interaction with anthropomorphized objects. The goal of the study, which was conducted at an electronics store of the “Conrad” chain, was to identify how accepting people would be to conversing with shopping products such as digital cameras. This study was part of a larger experiment designed to test modality preference and modality intuition. A total of 1489 interactions were logged over the two week test period, averaging 55 interactions per subject. Each test session generally took between 45 and 60 minutes to complete, during which time an average of 13.8 users could be seen from the shelf’s location.

## 6.1 Method

Our sample of test persons consisted of 27 people, 16 female and 12 male, and ranging in age from 19 to 55 (mean: 28.3 years). We advertised the study by posting notices around the University of Saarland, and setting up a registration desk at the main cafeteria. Only two subjects were from the faculty of computer science. Our setup consisted of a shelf of digital cameras located in a prominent part of a local electronics store. Each participant was allocated a PDA and headset, and asked to stand in front of the shelf containing real-world camera boxes. The subjects were briefed on how to use the system and the individual modality combinations. They were then instructed to interact indirectly using the 3<sup>rd</sup> person tense (e.g. “*what is the price of this camera?*”), and then later on directly by using the 2<sup>nd</sup> person tense (e.g. “*what is your price?*”). In each case, the products responded in an aligned manner (i.e. 3<sup>rd</sup> and 1<sup>st</sup> person tenses respectively). To ensure that our subjects spent enough time within both modes, they were given a series of smaller sub tasks to complete, such as to find the cheapest camera on the shelf, or to find the camera with the largest number of mega pixels. During the test, system output was limited to a single female concatenative synthesizer voice. This configuration was chosen to minimize the effect that voice quality and limited number of voice types might have on the study. After having completed the practical component, the participants were given a small questionnaire.

## 6.2 Results

The first question that we asked our subjects was which of the two interaction modes they preferred best. The proportion of users that preferred direct interaction over indirect interaction (18 from 27, 66%) signifies a distinct trend for anthropomorphization,  $\text{Chi}^2(1, N=27)=3.00$ ,  $p=0.083$ . This result is seen clearer in men than in women, in which 10 from 12 men (83%) stated that they preferred direct interaction:  $\text{Chi}^2(1, N=12)=5.22$ ,  $p=0.021$ , which is significant. An advantage seen by several subjects with direct interaction was that the dialog interactions were shorter and simpler (e.g. “what is your price” compared to “what is the price of the PowerShot S50”).

Following this question, we asked our subjects if they would reciprocate with direct interaction if the objects only spoke directly to them. 22 from 27 users (81%) stated that they would allow themselves to be coerced into communicating directly:  $\text{Chi}^2(1, N=27)=10.70$ ,  $p=0.001$ , which is significant. Courtesy and conformity were cited reasons for this allowed coercion. Note that a “no” response to this question would result in incoherent language similar to the following:

U: “What is the price of this <gesture> camera?”

O: “My price is €99”.

U: “How many megapixels does this camera have? <gesture>”.

We then asked our subjects whether they would interact directly with a given range of products (soap, digital camera, personal computer, and a car), first as a Buyer (B), and then as the Owner (O) of the product. For brevity, we report only the resulting significance values obtained from our non-parametric chi-square tests, where  $df=1$ , and  $N=27$ . Whereas only around 30% of people would interact directly with a bar of soap (as B:  $p=0.034$ , as O:  $p=0.201$ ), around 70% of people said that they would interact directly with digital cameras (as B:  $p=0.034$ , as O:  $p=0.033$ ), personal computers (as B:  $p=0.012$ , as O:  $p=0.003$ ) and cars (as B:  $p=0.336$ , as O:  $p=0.003$ ). Our subjects were more inclined to interact directly with the products as the owner rather than as a buyer, and this difference is best seen for the product type “car”, in which a Wilcoxon signed rank test bordered on statistical significance ( $z=-1.890$ ,  $p=0.059$ ). As the owner of the products “personal computer” and “car”, men were more inclined than women to talk directly with the objects, with a Mann-Whitney U-test showing this trend in gender difference to be:  $U(16,12)=40.5$ , equating to  $p=0.072$  for both product types. Other objects that our subjects said they would consider talking directly with included plants, soft toys, computer games, and a variety of electronic devices like TVs and refrigerators.

Finally, we tested which modalities people would be comfortable using in a public environment (e.g. when surrounded by other shoppers), compared to a private environment (e.g. when no shoppers are around). Given the choice of “comfortable”, “hesitant”, and “embarrassed”, the results showed that our subjects would feel comfortable using all modalities except speech when in a public environment ( $\text{Chi}^2(2, N=27) > 12.667$ ,  $p < 0.002$ ), and would feel comfortable using all modalities in a private environment ( $\text{Chi}^2(2, N=27) > 10.889$ ,  $p < 0.004$ ).

### 6.3 Lessons Learnt

From this empirical study, our hypothesis that users would not simply reject the concept of anthropomorphized objects was confirmed, and indeed many of our users actually enjoyed the concept. The study has also shown that product type (e.g. toiletries, electronics, automobile), relationship to a product (e.g. buyer, or owner), and gender (male, female) all have an effect on a persons preference for direct interaction with anthropomorphized objects. Future tests on the benefits of anthropomorphization could focus on a broader set of product types, the acceptance of cross-selling, and richer product personalities including distinct voices.

## 7 Conclusions and Future Work

This chapter has described a new interaction paradigm for instrumented environments based on tangible multimodal dialogs with anthropomorphized objects. For this purpose, we introduced the concept of symmetric multimodality and applied it

to speech, handwriting, and gesture. Finally, we showed via a usability field study that direct interaction with anthropomorphized objects is accepted and indeed preferred by the majority of users. Such findings have already been exploited in two other projects of our research group in which interactive installations for museums and theme parks are being developed.

Future work will now focus on scalability aspects of our approach, which will be particularly important if the system is to provide a shop full of differing products with rich forms of communication and personalities. The underlying grammars of this mobile system have currently been handcrafted for each product-type. This is acceptable when many products all have the same attributes, such as with digital cameras, but is less acceptable when many different product types exist, as would be the case when modeling the products of an entire store. We are currently developing a module to automatically generate the direct and indirect grammars based on keyword information available in the product database, and the type of question to be associated with the keyword, for e.g. a wh-question (who, what, when, where, why, and how), or a yn-question (yes, and no), and perhaps later also alternate and tag questions.

## **Acknowledgements**

This work was partially funded by the German Federal Ministry for Education and Research (BMBF) under the contract no. 01 IN C02, as part of COLLATE II.

---

## References

- Butz A, Schneider M, Spassova M (2004) SearchLight - A Lightweight Search Function for Pervasive Environments. In *Proc. of the 2<sup>nd</sup> Intern. Conf. on Pervasive Computing*, Springer LNCS, pp. 351-356.
- Cohen PR, McGee DR (2004) Tangible multimodal interfaces for safety-critical applications, In *Commun. ACM*, 47(1), ISSN 0001-0782, pp. 41-46.
- Cohen PR, Johnston M, McGee D, Oviatt S, Pittman J, Smith I, Chen L, Clow J (1997). QuickSet: Multimodal interaction for distributed applications, In: *Proc. of the Fifth International Multimedia Conference*, pp. 31-40.
- Don A (1992) Anthropomorphism: from ELIZA to Terminator 2. Panel Session in: *Proceedings of CHI'92*, ACM, pp. 67-70.
- Fox A, Johanson B, Hanrahan P, Winograd T (2000) Integrating information appliances into an interactive workspace. In *IEEE Computer Graphics and Applications*, 20(3), pp. 54-65.
- Jeremijenko N (2001) Dialogue with a Monologue: Voice Chips and the Products of Abstract Speech. In *SIGGRAPH*.
- Johnston M, Bangalore S, Vasireddy G, Stent A, Ehlen P, Walker M, Whittaker S, Maloor P (2002) MATCH: An Architecture for Multimodal Dialogue Systems. In *Proc. of Association for Computational Linguistics (ACL)*, pp. 376-383.
- Knapp K (2000). Metaphorical and interactional uses of silence. In *EESE: Erfurt electronic studies in English*.
- Kourouthanasis P, Spinellis D, Roussos G, Giaglis, G (200) Intelligent cokes and diapers: MyGrocer ubiquitous computing environment. In: *Proc. First International Mobile Business Conference*, pages 150–172, July 2002
- Kröner A, Baldes S, Jameson A, Bauer M (2004) Using an Extended Episodic Memory Within a Mobile Companion. In *Proc. of the Pervasive 2004 Workshop on Memory and Sharing of Experiences*, pp 59-66.
- Maybury M, Wahlster W (eds.) (1998): *Readings in Intelligent User Interfaces*. San Francisco: Morgan Kaufmann, 1998.
- Nijholt A, Rist T, Tuijnbreijer K (2004). Lost in Ambient Intelligence?. Panel Session in: *Proc. of CHI'04*, ACM, pp. 1725-1726.
- Oviatt, S W, Wahlster W (eds): (1997): *Multimodal Interfaces*. Human-Computer Interaction Journal, Vol. 12, No. 1-2.
- Oviatt SL (2002) Breaking the Robustness Barrier: Recent Progress on the Design of Robust Multimodal Systems, In *M. Zelkowitz (ed.): Advances in Computers*, vol. 56, 305-341.
- Schneider M (2004) Towards a Transparent Proactive User Interface for a Shopping Assistant. In *Workshop on Multi-User and Ubiquitous User Interfaces (MU3I 2004)*, pp. 31-35.
- Todd J (2002): The Hopi Environmental Ethos. In: Ferrero P (ed): *Hopi: Songs of the Fourth World. Resource Handbook*, Ferrero Films.
- Ullmer B, Ishi H (2001) Emerging Frameworks for Tangible User Interfaces. In *John M. Carroll(ed.) Human-Computer Interaction in the New Millenium*, Addison-Wesley, pp. 579-601.

- Wahlster W, André E, Finkler W, Profitlich HJ, Rist T (1993) Plan-Based Integration of Natural Language and Graphics Generation. In: *Artificial Intelligence*, 63, pp. 387-427.
- Wahlster W (2003) Towards Symmetric Multimodality: Fusion and Fission of Speech, Gesture, and Facial Expression. In Günter A, Kruse R, Neumann B (eds.): KI 2003: Advances in Artificial Intelligence. Proc. of the 26th German Conference on Artificial Intelligence, Berlin, Heidelberg: Springer, LNAI 2821, 2003, pp. 1-18.
- Wasinger R, Krüger A (2004) Multi-modal Interaction with Mobile Navigation Systems. In: W. Wahlster (ed.): *Special Journal Issue "Conversational User Interfaces"*, *it - Information Technology*, 46(6), ISSN 1611-2776, pp. 322-331.