

TOWARDS AUTOMATIC EVALUATION OF MULTIMODAL USER INTERFACES

Joëlle Coutaz, Daniel Salber, Sandrine Balbo

Laboratoire de Génie Informatique, IMAG
B.P. 53X, 38041 Grenoble Cedex
Tel. +33 76 51 48 54, e-mail: balbo@imag.fr, joelle@imag.fr, salber@imag.fr

ABSTRACT

The evaluation of the usability and the learnability of a computer system may be performed with predictive models during the design phase. It may be done on the executable code as well as by observing the user in action. In the latter case, data collected in vivo must be processed. Our goal is to provide software supports for performing this difficult and time consuming task.

This article presents an early analysis and experience towards the automatic evaluation of multimodal user interfaces. With this end in view, a generic Wizard of Oz platform has been designed to allow the observation and automatic recording of subjects' behavior while interacting with a multimodal interface. We then show how recorded data can be analyzed to detect behavioral patterns, and how deviations of such patterns from a data flow-oriented task model can be exploited by a software usability critic.

KEYWORDS

Capture of behavioral data, multimodal user interface, Wizard of Oz, user interface evaluation techniques.

INTRODUCTION

The development of interactive systems is an iterative process composed of three steps: design, construction and evaluation. Software tools such as interaction toolkits and UIMS technology, or software architecture models such as PAC⁷ and the Abstraction-Link-View paradigm¹⁴, have been developed to facilitate the construction of graphical user interfaces (GUI). Although the construction of user interfaces has been widely addressed by the software engineering community, little attention has been paid to support user interface design and evaluation.

Parallel to the development of graphical user interfaces, natural language processing, computer vision, and gesture analysis have made significant progress. Clearly, the combination of medias and modalities

open a complete new world of experience but our current understanding on how to design, build and evaluate such interactive systems is still primitive.

This article presents our early analysis and experience with the automatic evaluation of multimodal interfaces. Our goal is to provide designers with a Wizard of Oz multimodal software platform flexible enough to support the evaluation of multiple interactive systems. The next two sections define the problem space: first, the main streams for user interface evaluation are presented and our own approach is situated in this framework. Then, a taxonomy for the study of multimodal interfaces is presented. (A full description of this classification is available in ⁸.) The last two sections are dedicated to our own solution to the problem: the description of the platform and a first experience with automatic evaluation.

AN OVERVIEW OF EVALUATION TECHNIQUES

As shown in figure 1, evaluation techniques for interactive systems may be divided into two broad categories. Predictive methods are applicable during the design phase. They do not require any system implementation, nor do they need effective users. At the opposite, experimental techniques rely on the existence of a physical apparatus ranging from mock-ups of the real system up to the full implementation of a running prototype.

Predictive techniques

In general, predictive techniques are theory-based. For example, GOMS⁵ and its related models such as the Cognitive Complexity Theory (CCT)¹⁵, rely on an explicit hierarchical decomposition of the user's tasks. This static representation is supposed to model the user's plan for accomplishing a particular task. The genuine GOMS is a pure analytic model of errorless performance: it is able to predict the time required to accomplish a task without errors. Recently, GOMS has been extended to predict errors due to cognitive overload¹⁶. CCT, on the other hand, is useful for comparing several designs in terms of learnability and knowledge transfer.

At the opposite of GOMS and CCT, PUM builds a dynamic model of the user's plan. It predicts errors through a programmable cognitive architecture²⁸. The designer specifies the knowledge that the user needs to accomplish a particular task. This description, which includes domain knowledge as well as knowledge about the user interface, is compiled in terms of rules. These rules represent the user's ability to accomplish this particular task with this particular user interface. Based on this knowledge, the PUM cognitive architecture tries to elaborate a plan. If no plan can be generated, then the designer is notified of a potential usability or learnability problem.

Predictive evaluation techniques may also be based on HCI heuristics. Typically, the assessor looks for properties in the user interface design that, he knows from experience, lead to usability or learnability

problems. Such knowledge, exemplified by the Smith & Mosier's work, may be embedded in an expert system such as KRI¹⁸. KRI is able to detect anomalies from a formal description of the user interface. However, only the lexical and syntactical levels of the interaction are covered by the critic. Task modelling and any high level cognitive activity are discarded.

Assessing a design through HCI heuristics is a difficult task. The task-based "cognitive walkthrough" method proposed by Lewis et al. provides a useful framework for extracting evaluation guidance from a formal theory of human-computer interaction¹⁷. It consists of "a list of theoretically motivated questions about the system" such as "how will the user access description of action?" or "how will the user associate description to action?".

In summary, the main benefit from predictive models and techniques is that they allow the evaluation of user interfaces at the design stage. A design can be improved before a costly implementation takes place. On the other hand, specifying data to a predictive model may be as time consuming as the implementation per se. In addition, predictions made by theoretical models are based on hypotheses, not on real data. As demonstrated by Pollier²⁴ as well as by Nielsen et al.²², heuristic evaluation is difficult to achieve. At least three assessors are necessary to discover a reasonable number of usability problems (i.e., half of the problems at best!)

Experimental techniques

Experimental techniques and methods deal with real data observed from real users accomplishing real tasks with a physical artefact. This artefact may include paper scenarios, mock-ups, computer system prototypes, or Wizard of Oz (WOz) platforms.

With a WOz setting, "designers can illustrate how users will interact with yet-to-build software"¹⁹. In general, WOz experiments have been applied to natural language interfaces only. Corpus collected in vivo would be used to tune the linguistic parameters of the system, and thus would improve the robustness of the interaction⁹.

In general, behavioral data from WOz experiments, are tape-recorded. As a result, they must be retrieved and interpreted by hand. This is a time consuming task which requires expertise and patience. However, recent WOz platforms are able to capture and mix digitized and analogical behavioral data¹². By doing so, automatic tools can be developed to support the evaluation process. However, to our knowledge, none of the WOz platforms has tackled the problem of multimodal interfaces.

In summary, analysis from experimental methods are performed on real data, not on uncertain hypothetical values. This benefit is counterbalanced by the volume of behavioral data to process and by

the difficulty to identify the appropriate parameters for a particular experiment. We believe that a WOz computer platform which automatically captures selected behavioral data, provides a good basis for the development of evaluation and design tools. In addition, a WOz computer platform to study multimodal interaction is certainly a promising enterprise. Our approach to this problem is presented in the rest of the article.

A TAXONOMY FOR MULTIMODAL USER INTERFACES

In psychology, a modality refers to a human sensory channel such as vision, audition and touch. In the theoretical framework of the Model Human Processor⁵ as well as in ICS³, these channels are modelled as specialized processors. Whereas a modality denotes a type of human communication channel, a media such as a computer system, is an artefact that conveys information by triggering one or several human communication channels. According to these definitions, how can a user interface be qualified as being multimodal? The ICS model can provide us with a useful starting point.

Multimodal User Interfaces

In ICS, the human information processing system is subdivided into a set of specialized subsystems. The sensory subsystems transform sense data into specific mental codes that represent the structure and content of incoming data. These representations are then handled by subsystems specialized in the processing of higher-level representations: the morphonolexical subsystem for processing the surface structure of language, the object subsystem for processing visuospatial structures, and the propositionnal and implicational subsystems for more abstract and conceptual representations. The output of these higher-level subsystems are directed to the effector subsystems (articulatory and limb).

Using a similar process, a multimodal system is able to represent and manipulate information at multiple levels of abstraction along multiple input or output channels. A channel covers a set of physical sensory (or effector) means through which particular types of information can be received (or transmitted)⁸ and processed. One observes that both multimodal and multimedia systems are characterized by communicating information either through multiple input channels or through multiple output channels or both. The multiplicity of communication channels along one direction (whether it be input or output) provides the basis for multimedia-lity and multimodality.

The distinction between multimedia-lity and multimodality lies in the degree of built-in cognitive sophistication of the system along the axis “level of abstraction”. Multimodality is characterized by the capacity of the system to interpret raw inputs up to high levels of abstraction (e.g., that of the task domain) or to render information starting from high level representations. Although multimedia-lity

includes interpretation and rendering, it is not capable of handling the highest task-domain level representations.

By extension, an interactive system may be both multimedia and multimodal. For example, an hypermedia system would illustrate task-domain concepts using images and sound replayed from a CD-ROM, and it would be controlled by the user in a multimodal way using both speech and mouse to navigate through the hyperspace. Note that current multimedia systems are all able to handle the highest task-domain level representations but they do so for commands only and through a unique channel. Thus any multimedia system is at least monomodal in order to recognize input commands.

If “striving for meaning along multiple channels” denotes multimodal user interfaces, the fusion/fission phenomenon and the granularity of concurrency define a classification for multimodal interaction per se.

Two dimensions for classifying multimodal interfaces: Fusion/Fission, granularity of parallelism

Fusion and Fission

Fusion refers to the combination of several chunks of information to form new chunks. Fission refers to the decomposition phenomenon. Fusion and fission are part of the interpretation and rendering functions, i.e., the sequence of transformations applied to input and output data respectively.

Considering fusion for the interpretation function, information chunks may (or may not) originate from distinct digital input channels. For example, the sequence of events “mouse-down, mouse-up” that occurs in the palette of a graphics editor are two information chunks that originate from the same input channel. They are combined within the context of the palette to form a higher information chunk (i.e., the selection of a geometric class). At the opposite the “put that there” paradigm as in Cubricon²¹ offers an example of fusion between chunks originating from distinct input digital channels. In this example, fusion is required to solve the coreferences expressed through distinct channels.

The interpretation function may also perform fission. It may be the case that information coming from a single input channel need to be decomposed in order to be understood at a higher level of abstraction. For example, consider the utterance “show me the red circle in a new window”. This sentence, received through a single digital channel, references two domains of discourse: that of the graphics task (i.e., “the red circle”) and that of the user interface (i.e., “a new window”). In order to satisfy the request, the system has to decompose the sentence into two high level functions: “create a window” and “draw a red circle” in the newly created window.

Similarly, the rendering function can perform fusion and fission. As an example of fusion, the picture of a town may be combined to a graphical representation of the population growth. The notions of town and population, which are modelled as two distinct data structures within the internal processes of the system, are combined at the lowest level and presented through a single output channel. Fission occurs when an information chunk gives birth to multiple representations whether it be through a single or multiple output channels. For example, the spoken message “watch this wall!” along with a blinking red line on the screen uses two distinct output channels to denote the same wall.

Parallelism

Representation and usage of time is a complex issue. In our discussion, we are concerned with the role of time within the interpretation and rendering functions. How does time relate to levels of abstraction? How does it interfere with fusion and fission? Parallelism at the user interface may appear at multiple grains: at the physical level and at the task level.

At the physical level, input corresponds to the user actions that can be sensed by input channels as an information chunk (e.g., an event). For example, a mouse click, a spoken utterance are information chunks. For output, the physical level denotes output primitives, that is the information chunks that can be produced by output channels in one burst. For example, a spoken message or the reverse video of an icon. The fusion example “Put that there” and the fission example “watch this wall” both require parallelism at the physical level using multiple input and output channels respectively.

From the system’s perspective, a task (i.e, an elementary task) cannot be decomposed further but in terms of physical actions. For input, an elementary task is usually called a command, that is, the smallest fusion/fission of physical user’s actions that changes the system state. For output, an elementary task is the set of output physical primitives used to express a system state change.

True parallelism at the command level allows the user to issue multiple commands simultaneously. It necessarily relies on the availability of parallelism at the physical level. Pseudo-parallelism at the command level as in Matis²³, allows the user to build several commands in an interleaved way as in multithread dialogues. Then, parallelism at the physical level is not required.

Figure 2 shows a possible classification for multimodal systems based on the parallelism and fusion/fission dimensions. A multimodal user interface may be:

- exclusive if input (or output) expressions are built up from one channel only and no parallelism is permitted at the interface,
- alternate if input (or output) expressions are built up from multiple channels but no parallelism is supported,

- concurrent if input (or output) expressions are built up from one channel only and parallelism is permitted,
- synergistic if input (or output) expressions are built up from multiple channels and parallelism is permitted.

As an example of exclusive multimodal user interface, we can imagine the situation where, to open a window, the user can choose among double-clicking an icon or say "open window". One can observe the redundancy of the ways for specifying input expressions but, at a given time, an input expression uses one channel only. In a concurrent system, the user would be able to express both commands simultaneously but using distinct independent modalities (no fusion/fission would be supported).

As an example of synergistic multimodal system, the user of a graphics editor can say *put that there* while pointing at the object to be moved and showing the location of the destination with the mouse or a data glove. Within an alternate system, fusion/fission is supported but sequentiality is imposed. As a result, if we consider the *put that there* example, sequentiality would require the user to say *put that* followed by a mouse click to denote *that*. He would then say *there* and click a second time to indicate the destination. Although not desirable, sequentiality may be imposed for technical reasons.

In summary, the two dimensions, fusion/fission and temporal constraints on the usage of modalities, define a problem space which, as shown in the next paragraph, provides a useful framework for reasoning about multimodality.

Benefits from the classification

As discussed in^{8,10}, one can study the implications of fusion/fission and temporal constraints on software architectures. Bourguet and Caelen⁴ exploit the framework for the interpretation of multimodal expressions in a dialogue model.

For example, in a graphic editor supporting concurrency, consider the vocal command *rotate the triangles* combined with the selection of a set of triangles with the mouse. Depending on the presence of fusion or not, the interpretation of the vocal command may have different effects:

- In an exclusive multimodal user interface, the vocal and gesture commands are independent. Then, *rotate the triangles* can be interpreted in two ways: a) rotation of the triangles selected in a previous command (i.e., the pronoun *the* acts as an anaphoric reference); b) rotation of all of the triangles in the picture if no triangle has been previously selected.
- In a synergic multimodal user interface, the vocal and gesture commands can be coupled. Thus, in addition to interpretations a) and b), *rotate the triangles* allows a third interpretation c) where *the* acts as a deictic related to the concurrent gesture. The choice between the 3 interpretations relies heavily on the dialogue model used to drive the

interface. In particular, if parallelism at the interface prevails, then solution c is selected. If sequentiality dominates, then solutions a and b are good candidates. In addition, if vocal modality is privileged, then a prevails over b.

A second interest of our framework is to provide the basis for human factors experiments and cognitive psychology studies of multimodality. For example, for a given task, one can identify the usage of modalities, or determine the constraints imposed on the user by sequentiality, or even consider the usefulness of synergism, etc. Although psychological theories and human factor principles provide useful recommendations for the GUI technology, they are not directly extensible to multimodal user interfaces. In the absence of hard-core theories, we opt for the Wizard of Oz experimental approach.

NEIMO, A MULTIMODAL WIZARD OF OZ PLATFORM

Objectives

The goal of Neimo is to provide designers with a “Wizard of Oz” environment to observe and evaluate how users interact with multimodal interfaces. Wizards are used to supplement missing functions such as recognizers or generators for a particular modality. Modalities studied with the platform include graphics, oral communication, and gesture. Gesture covers facial expression and 2-D mousing. We consider that the redundancy provided by facial expression can be used to disambiguate the interpretation of end-user’s behavior.

An essential requirement for such an environment is flexibility. Neimo must support any number of wizards, any number of modalities, and any type of application whose source code is available (An application is the functional core that implements task domain concepts.) In addition, for a given application and a given set of modalities, it should allow multiple experiments with fusion/fission and time constraints.

Hardware configuration

Figure 3 illustrates a typical hardware configuration for the Neimo platform. It includes one workstation for the user (i.e., the subject) and several workstations for wizards. All of the workstations are Apple Quadras.

Voice Navigator, a word-Markov network-based pattern matching system, is used as the speech recognizer. Because of its poor performance with regard to “naturalness”, we intend to replace Voice Navigator with a wizard and evaluate the behavioral differences.

A CCD camera is focussed on the user's face and connected to a video acquisition board installed in one of the wizards' workstation. Thus, it is possible for that wizard to observe an image of the user's face digitized in real time (2 to 3 images per second).

Wizards are classified in two categories: the functional wizards and the modal wizards. A functional wizard accomplishes the task domain dependent services that are not implemented in the application. A modal wizard replaces or complements software components specialized in the processing of a particular modality or in the fusion of multiple modalities. For example, when Voice Navigator is turned off, a speech modal wizard comes into play. Similarly, a modal vision wizard is used to interpret facial expressions²⁷.

Software organisation and services

In order to satisfy the flexibility requirement, the Neimo platform is organized around a minimal communication kernel, NeimoCom, interfaced by libraries. As shown in Figure 4, NeimoCom acts as a communication server for transferring messages between workstations, and serves as a message recorder in history files. Each workstation runs a set of client functions linked to a NeimoCom library. These functions are not provided by the kernel but are developed for specific purpose. For example, the client functions of the speech wizard workstation support the wizard's task. In particular, a list of prerecorded answers is proposed to the wizard in order to alleviate his cognitive load as well as to guarantee a consistent behavior with regard to the user.

A Neimo library provides client functions with the following main services:

- open and close a connection. A wizard workstation can dynamically open or close a connection during a session.
- declare new types of message or redefine previously defined types with additional fields and/or suppression of obsolete fields. A message type is declared as a data structure with named typed fields. Since data fields are named, their order is not significant. By so doing, old clients can be modified or new clients added to the environment without jeopardizing previous settings. For example, a color field may be added to the original type *face* which was first designed for black and white pictures of faces.
- dynamically subscribe to a set of message types. This service allows clients to express their interest for a category of messages. For example, the vision wizard client subscribes to the *face* type to receive images of the user's face in order to be able to process and display them onto the wizard's screen. In addition, the dynamicity of the subscription allows wizards to change roles on the fly.

- send and receive messages synchronously. Messages are time-stamped by NeimoCom which maintains the universal time.
- open, close history files and record messages in an opened history file. In addition to navigation functions such as *go to the last record of history file h*, or *select record i from history file h*, clients can define views. A view acts as a filter on a history file. It includes a start and a stop date to identify the temporal window of interest, the origin (i.e., wizard VS user) as a selector of the source of the recorded messages, and a list of the message types of interest (for example face, mouse and speech).

A more detailed description of the run time kernel is available in¹. To summarize, the originality of the Neimo kernel is three-fold:

- 1) messages from multiple medias are processed in a uniform and integrated way,
- 2) message types are not imposed by the system. Instead, their level of abstraction is client defined. It is then possible to record information from low level events such as mouse clicks up to high level commands;
- 3) messages are recorded on request in dedicated files. Messages can be subsequently retrieved either directly according to their record sequencing number, or indirectly through the notion of views. Views allow client programs to extract messages through temporal windows on a set of message types. A user interface critic is one of such potential clients.

In addition to the run time kernel, a minimal user interface common to all of the wizards has been designed. This user interface allows a wizard to:

- set up the configuration for the session: which user's workstation to observe, which message types to receive, etc.;
- control message recording and message load during the session;
- observe the user's behavior: reception of the user's utterances, message sequence received from the user, replication of the user's screen as well as the user's mouse movements.

The Neimo platform is under development using MacApp in the MPW environment. Although we have not yet been able to make full-fledged experiments with Neimo, an early experiment of message recording and interpretation under X window² will be re-implemented in the Neimo environment. The model developed in this experiment is described in the next section.

OUR APPROACH TO AUTOMATIC EVALUATION

Our automatic evaluation of user interfaces is a four step process: 1) definition of a task model, 2) acquisition of behavioral data, 3) identification of behavioral patterns, and 4) critic per se.

A task model defines the optimal way of performing the task in a particular context for a particular domain. It is a behavioral reference model. Recorded data, as those captured by Neimo, reflect the effective behavior of the user performing a task in a quasi-realistic setting. As in the MRP technique²⁶, behavioral patterns are repeated user actions that may reveal usability problems. The critic per se combines general heuristic HCI knowledge with data specific to the case at hand: it detects deviations of the behavioral patterns from the reference task model. These four aspects are developed in the following paragraphs.

The task model

The definition of a task model depends on its end use.

In the design stage, the task model expresses the logical use of the system. It structures the work space in terms of tasks and subtasks showing the relationships between clusters of logically connected tasks. As in GOMS and CLG²⁰, the representation is a task hierarchy whose leaves denote the tasks that are conceptually indivisible. This is the conceptual task model.

When considering the running prototype, the task model aims at specifying the way the system functions. At the opposite of the conceptual task model, it specifies the way the user should perform the task with the real system. It does not necessarily describe a manner that is convenient for the user. Ideally, the conceptual task model and the effective task model should be isomorphic. In particular, the elementary tasks of the conceptual model should correspond to commands in the effective task model.

Although mapping elementary tasks to commands is straightforward, the correspondence for compound tasks is difficult to formalize. This situation results mainly from the discrepancy between the points of view adopted at the design and implementation steps. In the design phase, attention is focussed on the domain. As a result, syntactic tasks such as window manipulations which, by definition are domain independent, are not considered. Clearly, our evaluation technique, based on the acquisition of real data, speaks in favor of an effective task model: the sequence tree.

A sequence tree represents the sequence of possible elementary tasks. An elementary task is indivisible and modifies the system state. A node denotes an elementary task. As shown in figure 5, a node may be decorated with preconditions P. For example, to open a file from the Macintosh desktop, the folder which contains the file must be opened. Edges express ordering between elementary tasks. If A and B denote

two elementary tasks, then $A \dashrightarrow B$ means that A must be executed before B, and A cannot be executed again. In $A \rightarrow B$, A must be executed before B but A may be executed again in the future. If a node has multiple siblings, then the user may choose any one of them and switch freely between the subtrees. This property allows for the expression of interleaving between tasks.

In summary, a sequence tree shows the possible migration of the user through an organized and constraint space of elementary tasks. In turn, elementary tasks are expressed in terms of the physical actions that the user can perform with input medias. A notation like UAN¹³ can be used to specify the correspondence.

Collecting behavioral data

Recorded data correspond to the physical actions of the user.

In an X Window environment, actions such as mouse clicks and key presses, are modelled as events. We have recorded them in a history file with additional information such as the name of the elementary task to which it belongs. In a Neimo environment, actions are conveyed through messages and additional actions such as facial expressions and meta-comments, can be interpreted by a wizard and recorded.

Data has been recorded for a set of simple tasks selected from a complex problem space: the design of user interfaces for everyday cars. This design environment, Compo, is similar to an authoring system like HyperCard, where scripts are expressed with an iconic language⁶.

Behavioral patterns

The analysis of the recorded data within Compo led us to identify three classes of interesting behavioral patterns: direction shift, action repetition and action cancellation as part of syntactic tasks.

A direction shift occurs when the user stops following a downward path in the sequence tree. For example, although it is useless to do so, he may systematically select an icon file in the Macintosh desktop before invoking the *page setup* item in the *file* menu.

Action repetition as part of syntactic tasks cover actions that do not modify the functional core but concern the user interface portion only. Typically, we have observed systematic resizing, or iconification, or scrolling tasks on newly opened windows.

Action cancellation as part of syntactic tasks refer to closing newly opened windows or navigating through menus without selecting any item.

The critic

Behavioral pattern types are then related to usability problems based on human factor and psychological principles. These are numerous, inflationary or even contradictory. As a first step experience, we have used the simple taxonomy provided by Scapin²⁵ with the notions of compatibility, homogeneity, concision, feedback pertinence, explicit control, cognitive load, and error management.

Rules encoded in the critic point out user interface flaws through behavioral patterns. For example, the systematic occurrence of a direction shift at a particular node in the sequence tree expresses a cognitive dead-end. The interface does not lead the user to build the appropriate model. There is an incompatibility between the effective task model and the user's mental model. In the same way, systematic iconification of newly opened window may correspond to a messy screen or expresses the irrelevance of the information contained in the window.

CONCLUSION

In summary, we have developed an early experience with the automatic evaluation of interfaces using a simple apparatus. From low level captured data (i.e., user's mouse clicks and key presses), from a task model, and from general HCI heuristic knowledge, we have been able to detect anomalies for a graphical user interface. This technique needs to be extended to multimodal interfaces. With this end in view, we have designed Neimo, a generic platform, able to register a wide range of behavioral data in an integrated way. These data may result either from the automatic interpretation of the user's behavior by the system, or from on the fly interpretation by human wizards. With such a tool, experimenters should be able to build rich but focussed history files and conduct experiments about the usage of modalities along the dimensions of our framework. As an additional benefit from history files, we expect to devise user models applicable to the design of multimodal interfaces.

ACKNOWLEDGEMENTS

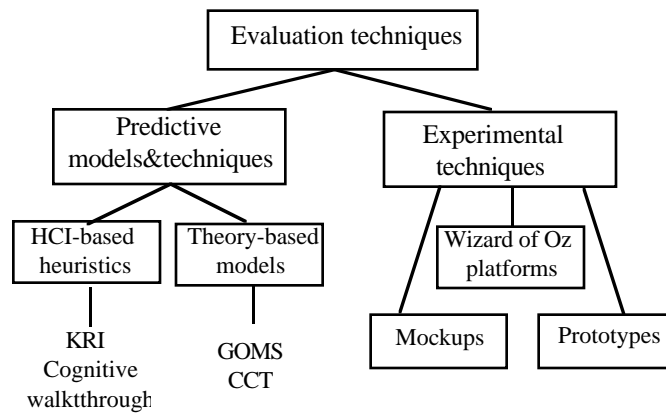
This work has been supported by project ESPRIT BR 7040 AMODEUS2 and by PRC Communication Homme-Machine.

REFERENCES

- 1 Ambone, G Noz, B & Salber, D 'Projet Neimo, Spécifications externes', *rapport équipe IHM LGI-IMAG* (1992).
- 2 Balbo, S Coutaz, J 'Automatic Evaluation in Human Computer Interaction' *The Ergonomics Society 1993 Annual Conference* Edinburgh 13-16 April (1993).
- 3 Barnard, P J 'Cognitive Resources and the Learning of Human-Computer Dialogs', in *Interfacing Thought, Cognitive Aspects of Human-Computer Interaction*, Carroll Ed. MIT Press Publ.(1987) pp112-158.
- 4 Bourguet, M L & Caelen, J 'Interfaces Homme-Machine Multimodales: Gestion des Evénements et Représentation des Informations', *ERGO-IA '92 proceedings* Biarritz (1992).

- 5 Card, S K, Moran, T P & Newell, A *The Psychology of Human Computer Interaction* Lawrence Erlbaum Associates (1983).
- 6 Chabert, A 'La programmation visuelle', *Rapport de DEA d'Informatique* Institut National Polytechnique de Grenoble (1991).
- 7 Coutaz, J 'PAC: an Implementation Model for Dialog Design' *Proceedings of the Interact'87 conference* Stuttgart, Bullinger & Shackel (ed.) North Holland (1987) pp 431-436.
- 8 Coutaz, J, Nigay, L. & Salber, D 'The MSM Framework: A Design Space for Multi-Sensory-Motor Systems', *EWHCI'93 proceedings* Moscow (1993).
- 9 Dahlbäck, N & Jönsson, A 'Empirical studies of discourse representations for natural language interfaces' *Fourth Conference of the European Chapter of the ACL Proceedings* (1989) pp 291-298.
- 10 Gourdol, A Nigay, L Salber, D & Coutaz, J 'Two cases studies of software architecture for multimodal interactive systems: VoicePaint and a Voice enabled graphical notebook' *IFIP'92 Congress* (1992)
- 11 Diaper, D 'The Wizard's Apprentice: A Program to Help Analyse Natural Languages Dialogues' *proceedings of the fifth conference of the British Computer Society Human-Computer Interaction Specialist Group*, University of Nottingham, (1989) pp 231-243.
- 12 Hammontree, M L Hendrickson, J J & Hensley, B W 'Integrated Data Capture and Analysis Tools for Research and Testing on Graphical User Interfaces' *in the CHI'92 Conference Proceeding*, ACM Press Publ. (1992) pp 431-432.
- 13 Hartson, R & Gray, P D 'Temporal Aspects of Tasks in the User Action Notation' *Human-Computer Interactio* Laurence Erlbaum vol. 7 No 1 (1992) pp 1-45.
- 14 Hill, R 'The Abstraction-Link-View paradigm: using constraints to connect user interfaces to applications' *in Proceedings of CHI'92* ACM Press (1992) pp 335-342.
- 15 Kieras, D & Polson P G 'An Approach to the Formal Analysis of User Complexity' *International Journal of Man-Machine Studies* 22 (1985) pp 365-394.
- 16 Lerch, F L Mantei, M M & Olson, J R 'Skilled financial planning: The Cost of Translating Ideas into Actions' *in Proceedings of CHI'89* ACM Press (1989) pp 121-126.
- 17 Lewis, C Polson, P Wharton, C & Rieman, J 'Testing a Walkthrough Methodology for Theory-Based Design of Walk-Up-and-Use Interfaces' *in Proceedings of CHI'90* ACM Press (1990) pp 235-241.
- 18 Löwgren, J & Nordqvist, T 'A Knowledge-Based Tool for User Interface Evaluation and its Integration in a UIMS' *Human-Computer Interaction-INTERACT'90* (1990) pp 395-400.
- 19 Mackay, W 'Video: Data for Studying Human-Computer Interaction' *CHI '88* (1988) pp 133-137.
- 20 Moran, T 'The Command Langage Grammar : a representation for the user interface of interactive computer systems' *International Journal of Man-Machine Studies* vol. 15 (1981) pp 3-50.
- 21 Neal, J Thielman, C Bettinger, K & Byoun, J 'Multi-modal References in Human-Computer Dialogue' *Proceedings of AAAI-88* (1988) pp 819-823

- 22 Nielsen, J & Molich, R 'Heuristic Evaluation of User Interfaces' in *CHI'90 Proceedings* ACM Press Publ. (1990) pp 249-256.
- 23 Nigay, L. & Coutaz, J. 'A Design Space for Multimodal Systems: Concurrent Processing and Data Fusion' *INTERCHI'93 Proceedings* Amsterdam (1993) pp 172-178
- 24 Pollier, A 'Evaluation d'une interface par des ergonomes : diagnostics et stratégies' *Rapport de recherche INRIA no 1391*(1991).
- 25 Scapin, D 'Guide ergonomique de conception des interfaces homme-machine' *Rapport Technique INRIA no 77* (1986).
- 26 Siochi, A & Hix, D 'A study of Computer-Supported User interface Evaluation Using Maximal Repeating Patern Analysis' in *Proceedings of the CHI'91 Conference* ACM Press (1991) pp 301-305.
- 27 Turk, M & Pentland, A 'Eigenfaces for recognition' *Journal of Cognitive Neuroscience* Vol. 3 No. 1 (1991) pp 71-86.
- 28 Young, R M &Whittington, J 'Interim Report on the Instruction Language' *AMODEUS Project Document: Deliverable D5, ESPRIT Basic Research Action 3066* (1990).



		USE OF MODALITIES	
		Sequential	Parallel
FUSION/FISSION	Combined	ALTERNATE	SYNERGISTIC
	Independent	EXCLUSIVE	CONCURRENT



Speech Wizard



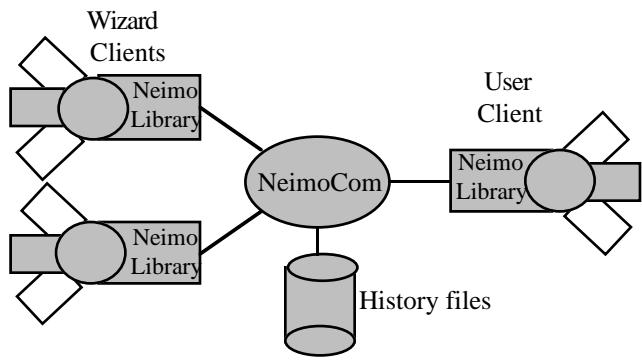
Face Wizard



Mouse Wizard



Observed User



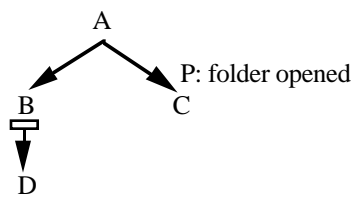


Figure 1: An overview of evaluation techniques for interactive systems.

Figure 2. A taxonomy of multimodal interactive systems.

Figure 3: Hardware configuration of the Neimo platform.

Figure 4: The software organisation of the Neimo platform. Dimmed areas denote common services and white areas, specific components.

Figure 5: The sequence tree model.