

Vision Systems with the Human in the Loop

C. Bauckhage, M. Hanheide, T. Käster, M. Pfeifer, G. Sagerer, and S. Wrede
Faculty of Technology, Bielefeld University,
P.O. Box 100131, 33501 Bielefeld, Germany

Abstract

The emerging cognitive vision paradigm deals with vision systems that apply machine learning and automatic reasoning in order to learn from what they perceive. Cognitive vision systems can rate the relevance and consistency of newly acquired knowledge, they can adapt to their environment and thus will exhibit high robustness.

This contribution presents vision systems that aim at flexibility and robustness. One is tailored for content-based image retrieval, the others are cognitive vision systems that constitute prototypes of *visual active memories* which evaluate, gather and integrate contextual knowledge for visual analysis. All three systems are designed to interact with human users. After we will have discussed adaptive content-based image retrieval and object and action recognition in an office environment, the issue of assessing cognitive systems will be raised. Experiences from psychologically evaluated human-machine interactions will be reported and the promising potential of psychologically based usability experiments will be stressed.

Key words: cognitive vision, adaption, learning, contextual reasoning, architecture, evaluation

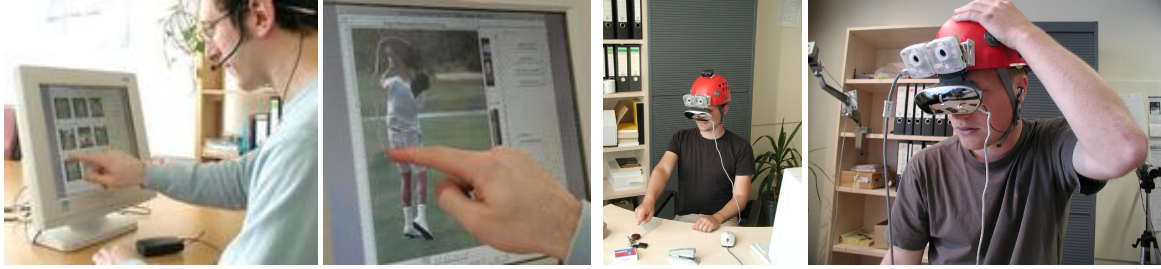
1 Introduction

Currently, the computer vision community is witnessing the emergence of a new paradigm. Even though its roots at least date back to work by Christensen, Crowley and Kittler [10] from the early 1990's, the idea of bringing together the achievements of 30 years of research in artificial intelligence, automatic perception, machine learning and robotics was termed *cognitive computer vision* just recently (cf. [13]).

Rather than trying to tackle the philosophical, psychological, or biological subtleties of the question what characterises cognition, we will adopt Christensen's point of view and restrict ourselves to a limited notion of cognition. Following his argument we will consider cognition as the generation of knowledge based on prior models, learning, reasoning and perception [8]. In this sense, cognition is an active process. Instead of just monitoring its surroundings, a cognitive vision system is able to communicate or interact with its environment. This underlines that the acquisition, storage, retrieval and use of knowledge is no end in itself but guides

the system's perception and (re)action. Simultaneously, the capabilities to perceive and act guide cognitive processes. Without perception and the possibility to manipulate or communicate perceived entities or events, knowledge cannot be acquired. Memory, however, is a limited resource. Besides mechanisms for learning, cognitive vision thus also implies attention control and a sense for relevance which comes along with the capability to forget irrelevant information. This requires flexible knowledge representation and techniques for top-down and bottom-up processing as well as functionalities for contextual reasoning and categorisation. Together with the biologically motivated principle of multiple computations [11], categorisation yields adaptability, flexibility and robustness.

Christensen even argues that embodiment is a prerequisite for cognitive vision systems. Only the capability to interfere with the environment can close the so called *perception action cycle*. However, even though there is considerable progress in the fields of mechatronics and robotics, machines that independently explore their environments are still in their infancy. In this contribution



(a) Multi-modal interactive CBIR

(b) Head mounted cameras and AR display

Figure 1: 1(a) Interactive content-based image retrieval using speech and haptics. 1(b) Head mounted cameras and display for augmented reality visualisation of recognised objects and events in an office environment.

we will thus argue that human-machine interaction can compensate embodiment. We will report results and experiences from two joint research projects on complex vision systems that make extensive use of the idea of the *human-in-the-loop*.

First, we shall present a system for interactive content-based image retrieval (CBIR). Although state of the art retrieval systems adapt to the preferences of their users, the involved learning processes only occur on the feature level of vision and there is no real knowledge acquisition. Claiming CBIR as a subfield of cognitive vision would therefore mean to overstretch the idea. However, CBIR systems are a perfect example of the benefits of bringing together pure computer vision and human-machine interaction. The retrieval system introduced in section 2 combines machine learning and adaption with intuitive multi-modal interfaces for image retrieval. While working with the system, the user may use natural language or a touch screen facility to indicate interesting image content (s. Fig. 1(a)).

Then, we will introduce systems which follow the cognitive vision paradigm. They are being developed in a research project dedicated to architectures and computational models for *visual active memories* (VAMs). Visual active memories are systems which evaluate given facts or gather and integrate contextual knowledge for visual analysis. VAMs can learn new concepts and categories as well as new spatio-temporal relations. They can adapt to unknown situations and may be scaled to different domains. Furthermore, the project investigates techniques and interfaces for advanced *interactive retrieval*. As an example, Fig. 1(b) shows impression

from experiments with a prototype of a mobile VAM. Working in a natural office environment, the user wears a head-mounted device which is equipped with cameras and a display. Information about recognised objects and results of user queries are visualised using augmented reality (AR). Likewise, by displaying status messages and prompts into the user's field of view the system can communicate with its user and thus close the perception-action cycle. Asking for manipulations of the environment in order to study their effects can accomplish interactive object and event learning.

The long-term perspective for interactive VAM research is to proceed towards *memory prosthetic devices*. The system in Fig. 1(b), for instance, can be seen as a first prototype of *memory spectacles* that may assist the memory challenged. But of course, expecting assistive technology to answer questions like "Where did I put my keys?" requires vision systems that will operate in everyday environments. The VAM demonstrators presented in section 3 are situated in unconstrained office environments. Applying multiple computations and contextual reasoning, the systems are able to identify different objects, actions and activities. They can be operated using speech and gesture; they cope with varying illumination as well as cluttered video signals and have capabilities in interactive object learning.

Given complex, interactive and adaptive vision systems, the problem of system evaluation arises. Obviously, the evaluation of an interactive system must not be restricted to a snapshotted performance testing. Rather, it has to take into account that failures that appear at a certain stage of an interactive session might be

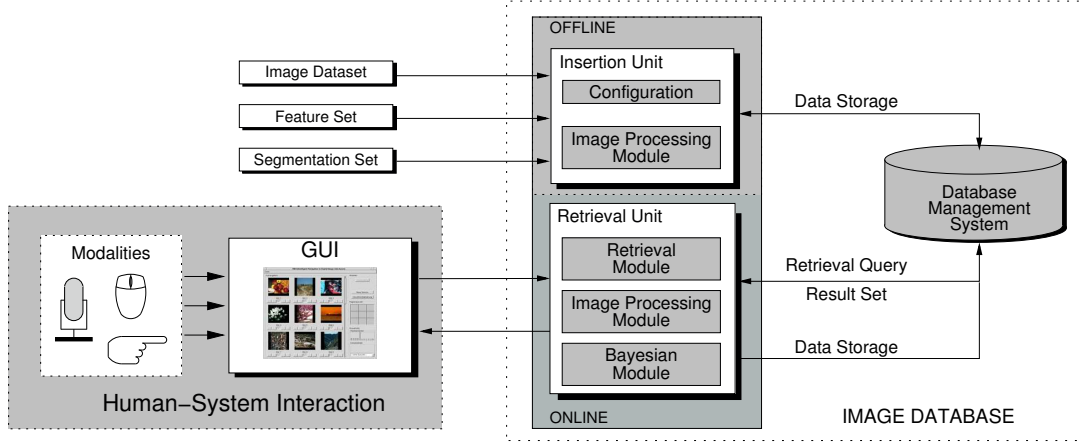


Figure 2: Components and conceptual architecture of the INDI system.

corrected later on. Also, learning and adaption might improve the system performance over time. However, up to now no commonly accepted evaluation framework that deals with these aspects has been established. In section 4, we will point out that usability experiments provide an promising avenue to solve this problem. We will report on a study designed by Psychologists we performed with naive users of our CBIR system. As we will see, this methodology can lead to surprising insight on how the human in the loop experiences his interaction with a cognitive system. Finally, a conclusion will close this contribution.

2 The INDI System

This section will present a system for content-based image retrieval (CBIR) that results from a project on **Intelligent Navigation in Digital Image** databases. Its characteristics are adaptability and multi-modal interaction. Adaption to the peculiarities of a certain retrieval task is guided by user feedback and happens on the feature level of computer vision. Multi-modal input devices are provided in order to facilitate intuitive handling. Figure 2 sketches the conceptual architecture of the INDI system. In the following, we will concentrate on the retrieval module displayed in the middle of the figure as well as on the user interface seen on the left.

2.1 A Hierarchical CBIR Approach

Image retrieval usually starts with low-level feature extraction either from an entire image or from certain image regions. The INDI system considers the following features: local moments in the LUV colour space as introduced by Stricker and Dimai [34], fuzzy histograms of the hue channel of the HSV colour space, and edge co-occurrence histograms which according to Brandt and Oja [7] are local shape descriptors.

Since local image signatures increase the precision in CBIR, our system automatically extracts regions of interest. In an initial keypoint detection process, the most salient points in a colour image are identified using the generalised Harris keypoint detector [28]. Afterwards, they are clustered using support vector clustering [4]. Pixels within the resulting clusters represent regions of interest which allow the computation of meaningful signatures and can be referenced during a retrieval process.

Following the approach of Rui and Huang [32], we assume an image object O_k , i.e. an image or parts of an image, to be characterised by several attributes: i) a set of pixels; ii) a set of feature classes, such as colour or texture; iii) for each feature class f_i there is a set of specific features. Examples of specific colour features could be histograms in different colour spaces or some sort of brightness information. All instances j of specific features are stored as sets of feature vectors $R = \{\vec{r}_{ij} \in R_{ij}\}$.

In the INDI system, we follow the common *query-by-example* approach and compute similarities between the database image objects O_k and a query object Q .

Using generalised Euclidian distances

$$m_{ij}(\vec{r}_{ij}, \vec{q}_{ij}) = (\vec{r}_{ij} - W_{ij}\vec{q}_{ij})^T \Omega_{ij} (\vec{r}_{ij} - W_{ij}\vec{q}_{ij})$$

where \vec{r}_{ij} and \vec{q}_{ij} are the feature vectors of the image object and the query object, respectively, similarities are computed separately for each feature class.

Again for each feature class, the image objects O_k are sorted yielding several ranked lists L_{ij} . Then, the ranks of the objects are linearly combined which produces an overall similarity ranking of the image objects $O_k, k = 1, \dots, n$, of the database with respect to the query object. Since the user of a content-based retrieval systems will only want to see reasonable matches, only the l most similar images (where $l \ll n$) will be selected from the database and displayed on the screen.

2.2 Adaption from Relevance Feedback

Iterative improvement during content-based image retrieval requires relating the user's high-level conception to low-level visual features. This is realised by means of relevance feedback. The user can rate objects in the current result list using scores $V \in \{2, 1, 0, -1, -2\}$ which represent ratings from *highly relevant* to *highly non-relevant*.

Preserving the information of previous search steps is accomplished by adapting the feature weights W_{ij} . Weights of features that allow the distinction of relevant and non-relevant images and thus allow to characterise the user's intention are increased, others are decreased:

$$W'_{ij} = W_{ij} + \epsilon \sum_{k=1}^l V(O_k) \cdot \lambda(\rho(O_k, L_{ij}))$$

Here, $V(O_k)$ is the score of image object O_k assigned by the user. ρ represents the rank of image object O_k in the feature dependent, ascendingly ordered result list L_{ij} . λ is a continuous descending function and ϵ is a learning rate.

Adopting another idea by Rui and Huang [33], the dissimilarity measures are refined as well. The matrix Ω_{ij} is adapted using the covariances of the feature vectors of to the image objects rated to be *relevant* or *highly relevant*. Finally, a query vector adaption is applied where the the query vectors in the feature spaces R_{ij} are slowly moved towards feature vectors of *relevant* and *highly relevant* image objects [2, 23].



Figure 3: Exemplary query images taken from a database of 1250 images from 10 different domains.

2.3 Evaluation of the CBIR Components

The adaptability of the INDI system was evaluated in different query tasks which were formulated as category searches like "show me images resembling Q ". Independence of the image domain was ensured by testing different categories, namely *autoracing*, *flowers* and *golfing* examples of which can be seen in Fig. 3.

Following the usual custom in information retrieval, a *precision* value was applied to evaluate efficiency and effectiveness. For the s th step of an interactive query, it is defined as

$$precision(s, t) = \frac{N_{s,t}}{t}$$

where $N_{s,t}$ represents the number of correct category images retrieved in session s within the first $t = 1, \dots, l$ retrieved images.

The adaptivity of our system is illustrated in Fig. 4. It shows the evolution of the precision values for $l = 27$ returned images over sequences of six interactive retrieval steps.

2.4 User Interface

In order to enable easy and intuitive handling, the INDI system provides different modalities for interaction. Except for the mouse there also is a touch screen facility. Both input devices enable the selection of images or image regions. They can be used to rate displayed database content and to initiate further selections from the database. Furthermore, a speech recognition component developed by Fink and colleagues was integrated whose core component is a statistical speech recogniser based on Hidden-Markov-Models [14].

Often, it is natural to use several input modalities simultaneously. For instance, users may point to the screen saying things like "this image". Therefore, a hierarchical event handling module was developed



Figure 4: Adaptation to the users intention expressed in terms of the evolution of the precision values in different category searches. Beginning with the second out of six search steps, positive feedback was provided. The depicted precision values are averaged over 10 experiments.

that can fuse asynchronous input events from different sources [24].

Given the all these input devices, the system must be able to relate verbally uttered commands to currently selected images or image regions in order to comprehend the user’s intentions. However, fusing results from speech and vision processing suffers from uncertainties like erroneous recognition or partial or unspecific descriptions. Consequently, we treat the task of speech and image integration as a *probabilistic decoding process* which is modelled using Bayesian networks (cf. e.g. [30]).

Adopting algorithms developed by Wachsmuth [38], each region description recognised in an utterance and each region detected in an image are represented as separate subnetworks. Matches between attributes obtained from speech recognition and those derived from image processing can be found by means of the relations in the network. After the relaxation of such a network, regions intended by the user will have the highest joint probability of being part of the image and also being referred to in an utterance [2].

3 VAMPIRE Systems

In this section, we will describe how the concept of human-machine interaction for computer vision can be extended to higher cognitive levels. While the previous section demonstrated how interaction can trigger adaption on the mere feature level of vision, this section will introduce *cognitive* vision systems that can learn new concepts and can adapt to a physical environ-

ment. We will present two systems that are being developed in a research project called **Visual Active Memory Processes for Interactive REtrieval** [36]. Both systems are able to recognise objects and activities in an unconstrained office environment. They can be operated using speech and gesture. Both make use of the principle of multiple computations and store results from different perceptual modules in a hierarchically organised memory. Processes registered in the memories apply contextual reasoning to verify the consistence and correctness of the incoming data. The memories themselves coordinate the registered processes and provide a notification mechanism to activate them if the memory content requires it. As such a memory is thus not a passive unit but rather is another active component of a system, we call it a visual active memory (VAM).

The VAM demonstrator shown in Figs. 5 and 6 analyses video signals from calibrated static cameras. Figure 5(a) depicts a human sitting in front of an office desk which is monitored by two cameras. One is observing the scene from above the other provides a side view of the desk. Figure 5(b) shows a snapshot recorded with the top view camera. In this example, the user is pointing to one of the objects on the desktop. In Fig. 5(c), the results of a view-based object recognition algorithm are cast into the image and Fig. 5(d) displays the results of skin colour segmentation and hand detection. As the index finger is stretched out, a gesture recognition algorithm identified a pointing gesture. Figure 5(e) visualises the angular probability distribution that indicates the most likely direction of this gesture.

Figure 6(a) exemplifies the side view on the scene. This viewpoint is used to recognise actions and activi-



(a) Office desk monitored by two cameras



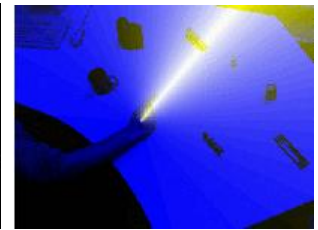
(b) Gesture seen from above



(c) Object recognition results



(d) Skin colour detection

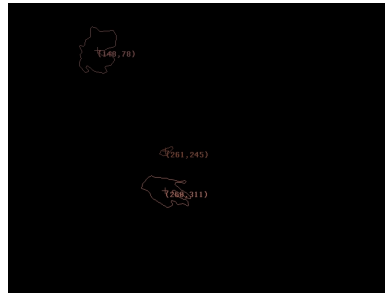


(e) Estimated pointing cone

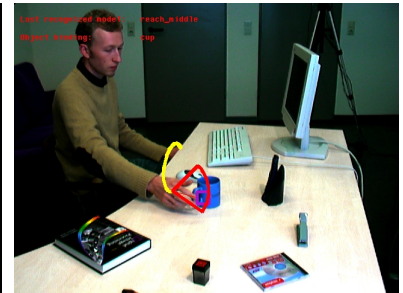
Figure 5: 5(a) VAM demonstrator with two static cameras monitoring a human sitting at an office desk. 5(b)–5(e) Exemplary results from processing top view images.



(a) Side view of the office scene



(b) Skin coloured regions



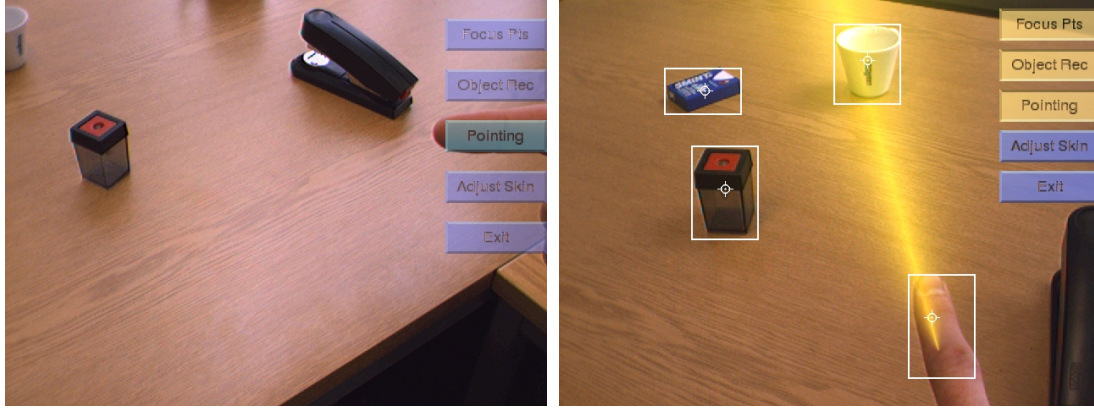
(c) Results from action recognition

Figure 6: Office scene and results obtained from the side view camera.

ties. Figure 6(b) shows a skin colour segmentation procedure for this example. While larger regions are assumed to depict faces, smaller ones are assumed to represent hands. In Fig. 6(c) the trajectory of one of the hands is cast into the image. Such trajectories are analysed by a module for action recognition. Furthermore, we see a fan projected into the middle of the image.

It indicates the image area near a moving hand where the system expects objects which might be manipulated next. According to the text displayed at the top of the image, the activity that was recognised last in this example, was 'reach middle' and the object that is currently expected to be manipulated is a cup.

Figure 7 shows the interaction with the mobile VAM



(a) Menu selection using pointing gestures

(b) Object referencing using pointing gestures

Figure 7: Office desk as seen through the mobile memory spectacles shown in Fig. 1(b).

demonstrator that was introduced in Fig. 1(b). By means of verbal commands or pointing gestures the user can browse through a command menu displayed on the right of his field of view. Selecting or deselecting menu buttons activates different operational modes of the system. Pointing gestures may also be used to reference objects or regions of interest in current the field of view. This resembles the use of the touch screen discussed in the last section. Here, however, space is becoming the interface; gestures are no longer bound to the operation of a physical input device.

3.1 Architecture and Components

Figure 8(a) sketches the conceptual architecture of our systems. In the centre, we recognise the memory component. It is organised hierarchically and is able to store image data (i.e. patches cropped from images) and feature based object descriptions as well as more abstract descriptions of observed events or categories. Several computational modules are grouped around the memory. Note that there is no direct communication between these modules but all data exchange is mediated through the memory. Also note that some of the building blocks represent several algorithms running in parallel.

All algorithms perform in real-time and run simultaneously. As we will detail below, the results they forward to the *active memory* are not considered as irrevocable facts but as hypotheses. Processes registered

on the database that provides the infrastructure for the memory continuously verify the consistency of incoming hypotheses and assign them a reliability. Corresponding hypotheses from different object recognition modules as well as from the action or gesture recognition components are fused into single abstract descriptions of the scene content. Moreover, since earlier results are stored in the memory, temporary occlusion or misinterpretations of the current scene can be filtered out using temporal context. Next, we shall outline the employed algorithms and technologies. For implementation details please refer to [1] and [18] for the static and mobile system, respectively.

3.1.1 Object Recognition

For object recognition, the VAMPIRE systems employ appearance based methods. On the one hand, VPL classifiers as introduced by Heidemann et al. [19] are applied. First, combining local entropy, symmetry and edge and corner detection, a saliency value is calculated for each image pixel. Where there is high saliency, patches are cropped from the image and classified in a three steps procedure using vector quantisation, PCA and LLM neural networks. On the other hand, we also use cascaded weak classifiers (cf. [25, 37]) for object recognition. For each object, windows of different sizes are shifted over the image. For each window, simple texture features are fed into the cascade. Already in

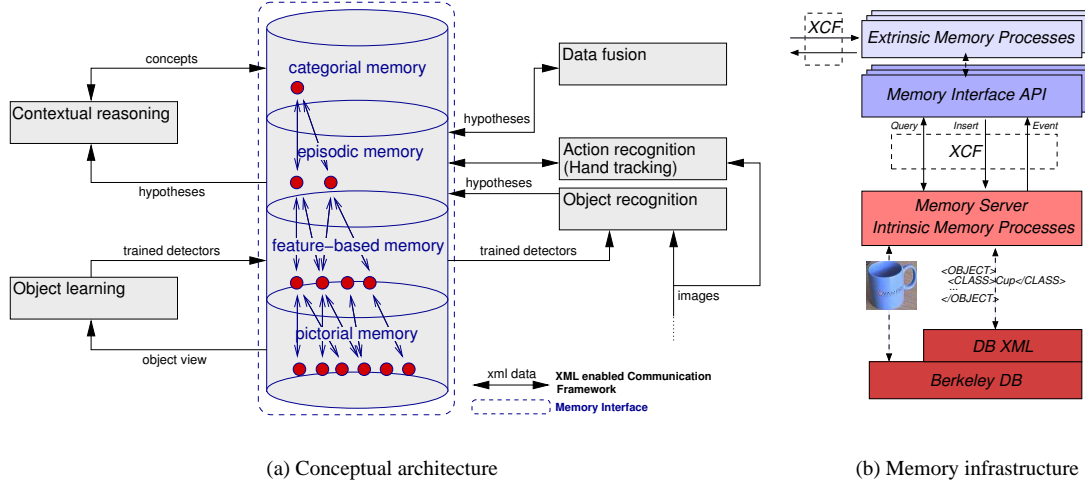


Figure 8: Conceptual architecture of the current VAMPIRE demonstrators and active memory infrastructure.

its first layer, most windows not depicting an object are rejected. Windows successfully passing through the whole cascade depict a known object. Either method is initially trained given manually labelled views of objects which were recorded in different positions and under varying illumination.

Both methods allow for interactive online object learning. Two techniques are being used. Either, the mobile AR-gear is used to focus on an unknown object. To acquire useful views of the object, template based image feature tracking as proposed by Gräßl et al. [16] compensates head movements. The second method incorporates the pointing mechanism described above. Introducing a rejection class label that is assigned to salient image regions which cannot be classified, these regions can be pointed to. If the user then moves the referred object to produce different views, the system can acquire a series of exemplary image patches. Randomly warping and distorting them yields artificial views which are then used to retrain the classifiers [3]. In either case, object labels are assigned verbally; to this end the systems are equipped with a speech recognition component [14] that was already mentioned in section 2.

3.1.2 Gesture and Action Recognition

Both, gesture as well as action recognition, rely on the detection of skin coloured image regions. To ensure ro-

business, we apply adaptive skin colour segmentation based on Gaussian mixtures models as described by Fritsch [15]. The mobile system provides yet another way for skin colour adjustment. After selecting a command for colour retraining from the interaction menu, moving the hand in front of the head mounted cameras produces data required for the adaption. Skin coloured image patches of a certain size are analysed by VPL classifier which decides whether they depict a hand or a even a pointing gesture.

Our action recognition framework is based on CONDENSATION particle filtering as introduced by Isard and Blake [21]. Black and Jepson [6] adapted this approach to the classification of hand trajectories. Using parameterised trajectory models, their techniques enable the recognition of activities solely on the basis of hand motions without incorporating any kind of context. For instance, 'pick' motions can be detected without information about *what* part has been taken.

In [15], Fritsch proposes an extension to the work of Black and Jepson in order to incorporate contextual knowledge. He distinguishes the *situational* context and the *spatial* context of a gesture.

The situational context of a gesture describes its necessary preconditions as well as the effect the gesture has on the scene. The spatial context of a gesture relates hand trajectories to objects being manipulated. Obviously, these objects must be close enough to a hand trajectory to be touched or picked for interaction. There-

fore, we define a *context area* to be the image area depicting objects potentially relevant for a specific gesture. The context area is given as a circle segment of a certain radius and angle. For interaction with objects that do not have an intrinsic 'handling direction' its orientation is defined relative to the moving direction of a hand. For objects that have an intrinsic 'handling direction', the context area has an absolute orientation. Besides defining *where* symbolic context is expected, we need to specify *what* context is expected. This includes the relevance of the context (irrelevant, necessary, or optional) as well as the type of the context object.

Actually incorporating context into recognition is done in two ways: The situational context is applied in the select step of the particle filter in order to initialise and select only those samples whose preconditions match the current situation. The spatial context is taken into account in the update step where it changes the weights of samples that match the observations. The calculation of sample weights is extended by a multiplicative *context factor* representing how well the observed scene fits the expected symbolic context.

3.1.3 Probabilistic Information Fusion

Due to flawed results of the perceptual modules or to a change in the environment, it might occur that hypotheses stored in the memory contradict one another. Consistency validation has to detect such conflicts and resolve them. As motivated in [39] and [40], elements in the memory are stored as XML fragments. Apart from information describing objects, these fragments also contain metadata like, for instance, the reliability of a hypothesis. An intrinsic memory process that lowers the reliability of stored data guides the removal, i.e. the *forgetting*, of conflicting hypotheses. The risk of conflicting results from object and action recognition is minimised by considering contextual and functional relations among incoming hypotheses. As they easily integrate different types of information, we apply Bayesian networks to model dependencies among the various facts our system gathers during runtime.

Consistency validation is realized as a memory process that uses *Functional Dependency Concepts* (FDCs) to rate stored hypotheses. FDCs basically consist of Bayesian networks that model expectations for the relations between specific types of hypotheses.

As an example, consider a situation where the user is sitting in front of a terminal and occasionally performs an action called 'typing'. Images of this situa-

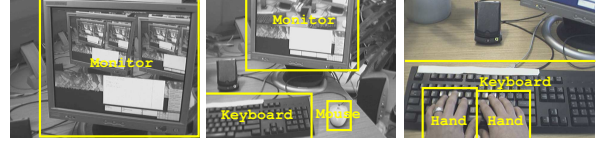


Figure 9: Three images of a sequence with annotated observations

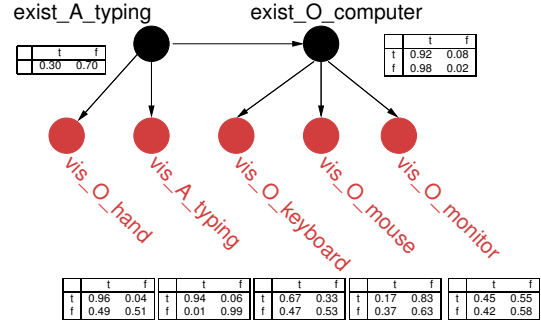


Figure 10: Bayesian network for a computer setup scenario

tion that were recorded with a head-mounted camera are shown in Fig 9. Recognising a 'typing' action is reasonable only under certain contextual prerequisites. For example, if there is no keyboard in the scene, 'typing' hypotheses have to be doubted. Figure 10 shows a Bayesian network and the corresponding conditional dependency tables used to represent contextual prerequisites for the 'typing'-action.

Nodes with the prefix *vis_* denote *observable* variables, whereas *exist_*-nodes are *hidden* and can only be inferred by the process. Inferring a computer, for instance, requires the observation of a keyboard, a mouse and a monitor. The object context required by a 'typing'-action is modelled as a directed arc from the action node *exist_A_typing* to the object node *exist_O_computer*.

The power of this approach lies in its applicability to any functional context. It allows for top-down as well as for bottom-up control and, as described in [29], this representation of contextual knowledge can guide object recognition and scene understanding. Conflicting memory content is detected as follows: For a given VAM content, the variables of an FDC are assigned *evidences* $e = \{e_1, e_2, \dots, e_m\}$. From evaluating the whole network, a conflict value *conf* can be calculated as a kind

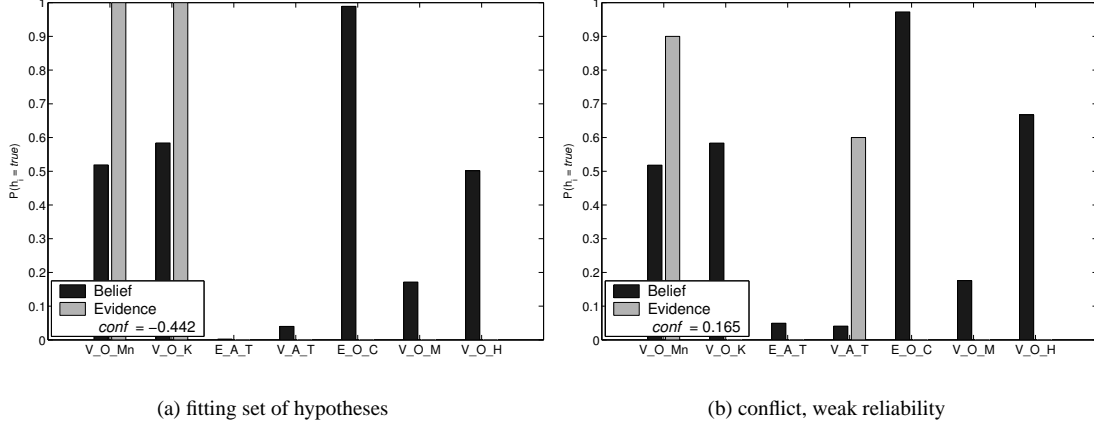


Figure 11: Two examples of beliefs and $conf$ -value for the FDC of the 'typing' action.

of emergence measure defined in [22]:

$$conf(\mathbf{e}) = \log_2 \frac{\prod_{i=1}^m P(e_i)}{P(\mathbf{e})}.$$

Here, $P(\mathbf{e})$ denotes the overall probability of the given evidences while the $P(e_i)$ are the marginal probabilities of the involved random variables of the Bayesian network. If there is a conflict, the probability $P(\mathbf{e})$ is expected to be small compared to the product of the probabilities $P(e_i)$ because in this case the evidences are not explained by the given FDC. Therefore, we will have $conf(\mathbf{e}) > 0$ which allows the detection of conflicts.

In order to cope with uncertainty of the underlying perception processes, *soft evidences* are used for the observable nodes. Their variables are assigned an evidence-vector $\vec{e} = (e_{true}, e_{false})^T$ with $0 \leq e_i \leq 1$ and $\sum e_i = 1$. A node's evidence is controlled by the reliability of the corresponding hypothesis. The more reliable the hypothesis is, the *harder* is its evidence. Evidences are set according to

$$(e_{true}, e_{false})^T = (0.5(1+r), 0.5(1-r))^T.$$

Thus, for a reliability $r = 1$, the evidence is set to $\vec{e} = (1, 0)^T$ while $r = 0$ will yield $\vec{e} = (0.5, 0.5)^T$ which is equivalent to an unobserved variable with no evidence. Details on the lowering of reliabilities in case of conflicts can be found in [17].

Probabilities for the conditional dependencies of the networks were estimated from manually annotated or correctly preprocessed video data. Figure 9 shows three out of 700 training images for the network in Fig. 10. If all nodes of the network are observable, parameter estimation simply means to counting the different configurations. Otherwise, with some nodes being not observed, an EM-algorithm is used (cf. [26]).

To evaluate our consistency validation approach, we defined FDCs for different constellations of objects and actions that are typical for an office scenario. Figure 11 displays prototypic results for the FDC of the 'typing'-action.

Diagram 11(a) depicts a situation corresponding to a consistent memory content. It shows highly reliable hypotheses *vis_O_monitor* and *vis_O_keyboard*, which mutually support each other. Note that $conf(\mathbf{e}) < 0$.

On the other hand, the configuration in 11(b) represents a conflict leading to $conf(\mathbf{e}) > 0$. In this example there are hypotheses of a monitor and a 'typing' action but no hypothesis for the keyboard which violates the expectation that a keyboard should be visible while typing.

3.2 System Integration

Developing complex vision systems is not only a matter of conceptual design but also a software engineering task. Concerning the development of a VAM, there

are two major issues: i) information storage and data organisation for the VAM and ii) a suitable communication framework allowing to distribute the different algorithms over several computers.

3.2.1 VAM Infrastructure

Since it is very flexible and suited for abstract concept descriptions, XML was chosen to describe content stored in the memory. Thus, a schema for symbolic data derived from vision algorithms (e.g. objects, actions, ...) was developed whose instance documents are composed of common and specific element structures (e.g. meta-data like reliability values). Beyond the simple and self describing nature of XML documents this has several other advantages. For example, the partition into common and specific elements is beneficial for the realisation of generic software modules where schema evolution allows for extensibility and XQuery/XPath techniques provide standardised access and selection mechanisms.

According to these consideration, a native XML database [12] provides the basic infrastructure for the VAM. On top of this embedded library, a server architecture as shown in Fig. 8(b) was implemented, that provides data management not only for XML but also for referenced binary data. Thus, pictorial data can also be used in the active memory and shared by several processes in parallel. Reference management is carried out using RDF information that links symbolic vision data to pictorial memory data. For both kinds of data, powerful standard DBMS methods like *insert*, *update*, *remove* and *query* are exposed. Node selection and referral is based on XPath statements.

Within this active memory server, for reasons of close coupling and performance, a runtime environment for intrinsic memory processes like *forgetting* or other, more generic, statistical processes was realised. A typical scenario for the use of this kind of processes are small, fast computations that work on large subsets of the system data. Furthermore, a subscription model for distributed event listeners was implemented, so that memory events can trigger registered processes and the memory indeed becomes *active*. Though realized in C++ there also is a MATLAB interface for rapid prototyping of further recognition or active memory components.

3.2.2 Communication Framework

Faced with the problem of distributing the algorithms discussed above over different machines in order to guarantee real time performance, a comparative study of existing framework technologies was carried out [41]. It yielded that by now there is no suitable integration framework tailored to the needs of cognitive vision. As most vision researchers are not middleware experts, the use of CORBA, for example, was ruled out due to its complexity and bloated standardisation. Rather, owing to the academic background of this work, an integration framework for an agile software process (cf. [9]) is needed.

This led to the development of an XML enabled Communication Framework (XCF) based on the Internet Communication Engine [20]. It provides an easy to use middleware for building distributed object oriented systems. Its architecture features a pattern based design and offers communication semantics like (a)synchronous streams, remote procedure calls and event channels. Similar to the data storage in the VAM component of our systems, data exchange between different modules is based on XML but wrapping and transport of binary data (e.g. images) is possible as well. Since interfaces are specified using XML schema, runtime type safety is ensured, rapid prototyping is possible, and interface programming is intuitive even for middleware novices.

Figure 12 presents a more technical sketch of the consistency validation example discussed above. After an extrinsic memory process, like object recognition, inserts a new hypothesis into the database, consistency validation is triggered. Related database content is queried using XPath and a conflict value is computed. Changes in the reliability values of stored hypotheses will trigger another intrinsic process. If they become too unreliable, hypotheses will be purged from the memory.

This example underlines that, in combination, the XML based memory infrastructure and the XCF framework enable to realize an architecture with low coupling between components. Furthermore, this decoupling and the capability of the memory to asynchronously gather and provide information yields a high robustness against component failure.

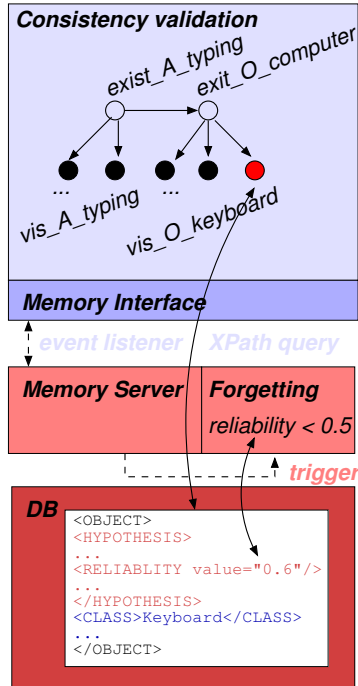


Figure 12: Example of interaction between extrinsic and intrinsic memory processes.

3.3 Technical Performance

Currently, the static system is running on five standard Linux PCs (Pentium 4, 2.4GHz, 512MB); images are captured using SONY DFW VL 500 firewire cameras providing a resolution of 640×480 pixels. The mobile demonstrator is running on a high performance DELL notebook (Pentium 4, 1.8GHz, 512MB); images are captured from fire-I firewire cameras with a resolution of 320×200 pixels.

Evaluating the core components of our systems as if they were stand alone modules yielded the following results: at a frame rate of 4 Hz, the VPL based recognition of gestures and objects yields an accuracy between 90% and 82% depending on the number of objects that have been trained [3]. The cascaded classifier approach to object recognition processes 6 images per second and yields 92% correctness. Trained with averaged trajectories from different videos and manually annotated information about object context, actions like 'drinking from a cup', 'reading a book', 'phoning' or 'typing on the keyboard' can reliably be recognised. A test with 420 sequences, yielded an accuracy of 93% [15]. Fi-

nally, local queries with low selectivity (approx. 1% of the whole dataset is returned) on a memory instance require an average of 0.57 seconds on a basis of 100,000 documents in a persistent memory (for an in-depth technical discussion of the evaluation of the XML enabled framework and the memory component please refer to [39] and [40]).

Having read all these figures, it appears that traditional performance assessment does not tell much about the overall performance of an integrated vision system. It is obvious that it does not take into account the continuous nature of human-machine interaction. Interaction with a flexible vision system is a process throughout which there will be mutual adaption. Learning and adaption may improve the system performance over time; recognition and interpretation errors that may appear during an interactive session might be corrected later on.

These considerations thus raise the problem of how to assess the long-term performance of an interactive vision system. Based on the experience reported in the next section, we are tempted to claim that asking the *human* in the loop may provide a solution.

4 Integrated System Evaluation

Modern evaluation of intelligent systems for advanced human-machine interaction has a history of about 10 years (cf. e.g. [27, 31]). Proposed approaches range from assessment by means of exemplary benchmarks [35] to the definition of measurable performance indices [5]. However, practical experience with performance measures was not reported. Moreover, neither do the methods known from literature consider situations of triadic interaction, i.e. situations where two agents coordinate their perception about a third person, thing, or event. Nor do they regard adaptive systems.

In the following, we will outline a holistic evaluation methodology that was applied to assess the capabilities of the INDI system [24]. Apart from collecting technical data like mentioned in the previous section, we also examined the *usability* of our system. To this end we carried out interactive experiments where we not only measured features like the average success rate in target search but also asked our subjects to fill out questionnaires in order to investigate human factors in interactive image retrieval. This focused on the following criteria adopted from Preece [31]:



Figure 13: Target images for query tasks.

- The *speed* of task execution.
- The *functionality* of the system, i.e. how many different tasks can be performed
- The *quality* of the results, i.e. how good is the average performance in different tasks
- The *speed of learning*, i.e. how quick can users learn to perform tasks with the system.
- the *mental load*, i.e. have users to think carefully while interacting with the system.
- *user satisfaction*, i.e. do users like working with the system.

4.1 Procedure and Design

We considered a database of 1250 images from 10 semantic categories which are taken the ArtExplosion image collection. A total of 20 computer experienced subjects (2 female and 18 male) who had never before operated a CBIR system were tested. They were divided into four groups of five people each and the input modalities

- *mouse* (M)
- *mouse and speech* (MS)
- *touch screen* (T)
- *touch screen and speech* (TS)

were evaluated. The modalities *mouse and touch screen* as well as *mouse, touch screen and speech* were not examined since initial tests revealed that people never used mouse and touch screen simultaneously.

Each subject took part in three interactive experiments. In each experiment, they were asked to retrieve an image from the database that was shown to them at the beginning (s. Fig. 13).

In every iteration of an interactive search, 27 images were displayed to the subjects which they could rate

in order to navigate through the database and find the query image. They could either score entire images or select certain regions from an image. The maximum amount of time for each experiment was limited to three minutes; if a subject was not able to retrieve the requested image within this time, the experiment was counted as a failure.

Besides the success rate S_E averaged over all experiments, the quality of interaction is characterised by the average time T_E the subjects needed to perform an experiment and by the mean number FB_E of user inputs, i.e. the amount of feedback provided in an experiment. Given the average number N_I of iterations of a query, it is possible to deduce the ratios T_I and FB_I describing the average time per iteration and number of feedbacks per iteration, respectively. The above mentioned aspects of learning, mental load, and user satisfaction were examined by means of the questionnaires the subjects were asked to fill out. Faced with statements like “It was fun to interact with the system” they ranked their sensation on a scale from 1 (no) to 5 (yes).

4.2 Results

Tables 1 and 2 and Fig. 14 summarise our findings. Looking at the figures in Tab. 1, it is noticeable that the three target searches were of increasing complexity. This is expressed in the increasing amount of time and feedback as well as in the growing number of interactions shown in the table.

Table 2 lists the figures we measured with respect to the different input modalities. We can see that subjects who only used the mouse provided most relevance feedback but did not achieve the best success rate. We also see that users of the touch screen device performed best and fastest while users of speech and touch screen were the slowest and least successful ones.

The latter observation is especially interesting if we regard Fig. 14. The diagrams in this figure depict the average ranking of the factors asked for in the questionnaires. In Fig. 14(a), for instance, we notice that the easiness of handling the mouse and of handling mouse and speech were both ranked 4.4, for the touch screen and speech modality it yielded 4.0 and the easiness of only using the touch screen reached 3.4 These figures accord with those in Fig. 14(e) which summarise our subjects notion regarding the patience their interaction required. Here, the touch screen users felt that they had to be most patient. Another interesting result becomes apparent from Fig. 14(d): users of multi-modal input

Target Image	$T_E = \frac{\text{time[s]}}{\text{experiment}}$	$FB_E = \frac{\text{feedbacks}}{\text{experiment}}$	$N_I = \frac{\text{iterations}}{\text{experiment}}$	$T_I = \frac{\text{time[s]}}{\text{iteration}}$	$FB_I = \frac{\text{feedbacks}}{\text{iteration}}$
RaceCar-78	73.0	9.2	2.1	33.95	4.28
Balloon-36	81.3	10.8	3.3	24.65	3.29
Flowers-32	96.5	15.3	4.2	22.98	3.65

Table 1: Experimental results w.r.t. target image.

Modality	S_E	$T_E = \frac{\text{time[s]}}{\text{experiment}}$	$FB_E = \frac{\text{feedbacks}}{\text{experiment}}$	$N_I = \frac{\text{iterations}}{\text{experiment}}$	$T_I = \frac{\text{time[s]}}{\text{iteration}}$
M	0.73	88.6	15.13	4.33	20.46
T	0.8	71.8	9.33	2.86	25.1
MS	0.73	79.66	11.8	2.93	27.18
TS	0.67	94.4	10.93	2.73	34.57

Table 2: Experimental results w.r.t. input modality.

devices rated their interaction with our CBIR system to be more efficient than those subjects who only worked with the mouse or touch screen.

4.3 Discussion

With respect to our six evaluation criteria our findings suggest: *Speed, functionality, and quality*: Concerning the time T_E , the number of iterations N_I as well as the number of user feedbacks FB_E , performances of mono-modal and multi-modal interaction diverge. While using mouse and speech is faster than only using the mouse, it is the other way round for using touch screen and speech. However, in any case, different target searches can be performed satisfyingly with regard to the average success as well as to average time need. *Learnability*: Regarding the tested input facilities, users did not sense a significant difference among modalities. *Mental load*: Measured results and user sensations are inconsistent. Even though the touch screen group performed best, their sensations concerning easiness and efficiency were worst. *User satisfaction*: Multi-modal input facilities are well appreciated by the users of our system. Even though their results in interactive image retrieval were not the best, the subjects who could use speech and another modality felt least annoyed and considered the interaction they had with the system to be efficient and fun.

5 Conclusion

This contribution reported on vision systems which make use of the concept of the human-in-the-loop. The first system is designed to enable efficient, intuitive, and easy content-based retrieval from image databases. On the one hand, it applies flexible techniques for image feature extraction and adaption on the lower levels of computer vision. On the other hand, it provides several input modalities. Understanding the problem of integrating the different modalities as a probabilistic decoding task enables to fuse the different types into consistent interpretations. As a consequence, natural and seamless interaction with the system becomes possible.

The two other systems we presented follow the cognitive vision paradigm. They are intended to demonstrate the idea of *visual active memory* (VAM). Situated in an unconstrained office environment, both systems recognise typical office objects as well as actions involving them. Information about recognised objects and events is stored in a memory and can be retrieved later on. Both systems are operated using speech or gesture; the mobile demonstrator uses AR technology to display memory content or control interfaces.

Robustness results from applying the principles of multiple computations and contextual reasoning. Different algorithms for object and gesture recognition process image sequences obtained from different views or from a set of head mounted cameras. The results of these computations are not seen as irrevocable facts but

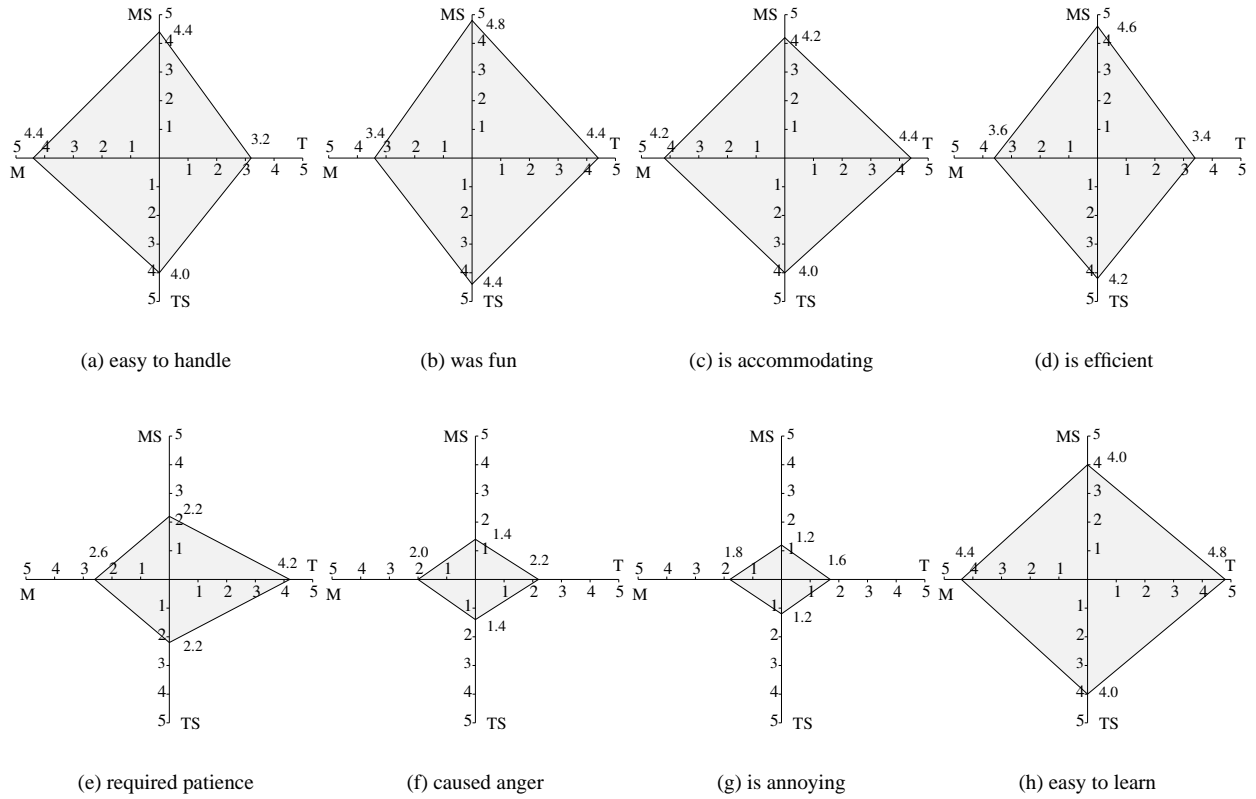


Figure 14: Averaged results of a questionnaire survey on usability aspects in interactive CBIR. For each interaction modality (mouse (M), touch screen (T), mouse and speech (MS), touch screen and speech (TS)). Each aspect had to be rated on a scale from 1 (no) to 5 (yes).

first of all as hypotheses. Hypotheses resulting from recognition processes applied to salient parts of the signal are forwarded to a memory component. There, processes that make use of probabilistic, top-down and bottom-up Bayesian reasoning verify their consistency. As processes like consistency verification and data deletion are triggered by the memory component, the memory indeed is an *active* module.

Basing the memory infrastructure on an XML database and realising the technical system integration using an XML enabled framework results in ease of use, extensibility and robustness against component failure. Moreover, the human-in-the-loop approach provides an avenue to even more flexibility. While for the image-retrieval system, adaption was only possible by weight adjustment on the feature level of visual processing, the presented VAMs can learn on higher cognitive levels. Through interaction with their users, they can extend

pre-acquired knowledge and learn representations and labels for new objects.

The systems introduced in this contribution thus demonstrate that the goals of the cognitive vision paradigm are not just illusory. Machine learning, contextual reasoning, relevance control and active system introspection can be brought together and human-machine interaction can compensate for embodiment. And indeed, in combination these techniques result in integrated systems of high robustness and flexibility.

However, dealing with the evaluation of complex integrated vision systems, human-machine interaction comes along with new challenges. Up to now, there is only scarce literature on how to characterise the mid- and long term performance of interactive systems. By means of our image retrieval system we thus exemplified how usability studies might help to assess the cognitive capabilities of artificial systems. As a matter of

fact, some of the results are surprising: even though they performed best users of simple interaction devices felt least content with the performance of the system. On the other hand, users of input devices of higher cognitive adequacy (natural language) experienced their interaction with the system to be very pleasant and efficient. Even though they practically obtained the worst retrieval results. Therefore, at least for now, it seems fair to conclude that research in cognitive vision must face the fact that cognition first of all lies in the eye of the beholder.

6 Acknowledgement

This work has been supported by the BMB+F under contract 01IB 001B and by the European Union IST 2001-34401 Project VAMPIRE. The authors would like to thank Silke Fischer for the valuable support and suggestions she provided for our usability experiments.

References

- [1] C. Bauckhage, M. Hanheide, S. Wrede, and G. Sagerer. A Cognitive Vision System for Action Recognition in Office Environments. In *Proc. Int. Conf. on Computer Vision and Pattern Recognition*, 2004.
- [2] C. Bauckhage, T. Käster, M. Pfeiffer, and G. Sagerer. Content-Based Image Retrieval by Multimodal Interaction. In *Proc. 29th Ann. Conf. of the IEEE Industrial Electronics Society*, 2003.
- [3] H. Bekel, I. Bax, G. Heidemann, and H. Ritter. Adaptive Computer Vision: Online Learning for Object Recognition. In *Proc. Pattern Recognition Symp.*, volume 3175 of *LNCS*. Springer, 2004.
- [4] A. Ben-Hur, D. Horn, H. Siegelmann, and V. Vapnik. Support vector clustering. *J. of Machine Learning Research*, 2:125–137, 2001.
- [5] N. Beringer, U. Kartal, K. Louka, F. Schiel, and U. Türk. Promise – a procedure for multimodal interactive system evaluation. In *Proc. Workshop 'Multimodal Ressources and Multimodal System Evaluation'*, Las Palmas, Spain, 2002.
- [6] M.J. Black and A.D. Jepson. A Probabilistic Framework for Matching Temporal Trajectories: CONDENSATION-Based Recognition of Gestures and Expressions. In *Proc. Eur. Conf. on Computer Vision*, 1998.
- [7] S. Brandt, J. Laaksonen, and E. Oja. Statistical Shape Features in Content-Based Image Retrieval. In *Proc. Int. Conf. on Pattern Recognition*, 2000.
- [8] H.I. Christensen. Cognitive (vision) systems. *ERCIM News*, pages 17–18, April 2003.
- [9] A. Cockburn. *Agile Software Development*. Addison-Wesley, 2001.
- [10] J.L. Crowley and H.I. Christensen, editors. *Vision as Process*. Springer, 1995.
- [11] H. Cruse. The evolution of cognition – a hypothesis. *Cognitive Science*, 27(1):135–155, 2003.
- [12] Berkely DB XML, Sleepycat Software. <http://www.sleepycat.com/products/xml.shtml>, 2004.
- [13] European research network for cognitive vision systems. <http://www.ecvision.info>, 2004.
- [14] G. A. Fink. Developing HMM-based recognizers with ESMERALDA. In V. Matoušek, J. Ocelíková P. Mautner, and P. Sojka, editors, *LNAI*, volume 1692. Springer, 1999.
- [15] J. Fritsch. *Vision-based Recognition of Gestures with Context*. PhD thesis, Bielefeld University, 2003.
- [16] C. Gräßl, T. Zinßer, and H. Niemann. Illumination insensitive template matching with hyperplanes. In *Proc. Pattern Recognition Symp.*, volume 2781 of *LNCS*. Springer, 2003.
- [17] M. Hanheide, C. Bauckhage, and G. Sagerer. Memory Consistency Validation in a Cognitive Vision System. In *Proc. Int. Conf. on Pattern Recognition*, 2004.
- [18] G. Heidemann, I. Bax, H. Bekel, C. Bauckhage, S. Wachsmuth, G. Fink, A. Pinz, H. Ritter, and G. Sagerer. Multimodal Interaction in an Augmented Reality Scenario. In *Proc. Int. Conf. on Multimodal Interaction*, 2004. to appear.
- [19] G. Heidemann, R. Rae, H. Bekel, I. Bax, and H. Ritter. Integrating context free and context-dependent attentional mechanisms for gestural object reference. In *Proc. Int. Conf. on Computer Vision Systems*, 2003.

- [20] The Internet Communications Engine. <http://www.zeroc.com/ice.html>, 2004.
- [21] M. Isard and A. Blake. Condensation – conditional density propagation for visual tracking. *Int. J. Computer Vision*, 29(1):5–28, 1998.
- [22] F.V. Jensen. *Bayesian Networks and Decision Graphs*. Springer, 2001.
- [23] T. Kämpfe, T. Käster, M. Pfeiffer, H. Ritter, and G. Sagerer. INDI – Intelligent Database Navigation by Interactive and Intuitive Content-Based Image Retrieval. In *Proc. Int. Conf. on Image Processing*, 2002.
- [24] T. Käster, M. Pfeiffer, C. Bauckhage, and G. Sagerer. Combining Speech and Haptics for Intuitive and Efficient Navigation through Image Databases. In *Proc. Int. Conf. on Multimodal Interfaces*, 2003.
- [25] J. Kittler, A. Ahmadyfard, and D. Windridge. Serial multiple classifier systems exploiting a coarse to fine output coding. In *Proc. Int. Workshop Multiple Classifier Systems*, volume 2709 of *LNCS*. Springer, 2003.
- [26] S. Lauritzen. The EM-algorithm for graphical association models with missing data. *Computational Statistics and Data Analysis*, 19(2):191–201, 1995.
- [27] G. Lindegaard. *Usability Testing and System Evaluation*. Chapman & Hall, London, 1994.
- [28] P. Montesinos, V. Gouet, and R. Deriche. Differential invariants for color images. In *Proc. Int. Conf. on Pattern Recognition*, 1998.
- [29] K. Murphy, A. Torralba, and W.T. Freeman. Using the forest to see the trees: A graphical model relating features, objects, and scenes. In *Proc. Conf. on Neural Information Processing Systems*, 2003.
- [30] J. Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, San Francisco, 1988.
- [31] J. Preece, Y. Rogers, and H.C. Sharp. *Beyond human-computer interaction*. Wiley & Sons, Chichester, 2002.
- [32] Y. Rui and T. Huang. Relevance feedback: A power tool for interactive content-based image retrieval. *IEEE Trans. on Circuits and Video Tech.*, 8(5), 1998.
- [33] Y. Rui and T. Huang. Optimizing Learning in Image Retrieval. In *Proc. Int. Conf. on Computer Vision and Pattern Recognition*, pages 1236–1243, 2000.
- [34] M. Stricker and A. Dimai. Spectral Covariance and Fuzzy Regions for Image Indexing. *Machine Vision and Applications*, 10:66–73, 1997.
- [35] C.-P. Tung and A.C. Kak. Integrating Sensing, Task Planning and Execution for Robotic Assembly. *IEEE Trans. on Robotics and Automation*, 12(2):187–201, 1996.
- [36] VAMPIRE. <http://www.vampire-project.org>.
- [37] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *Proc. Int. Conf. on Computer Vision and Pattern Recognition*, volume 1, 2001.
- [38] S. Wachsmuth and G. Sagerer. Bayesian Networks for Speech and Image Integration. In *Proc. AAAI*, 2002.
- [39] S. Wrede, J. Fritsch, C. Bauckhage, and G. Sagerer. An XML Based Framework for Cognitive Vision Architectures. In *Proc. Int. Conf. on Pattern Recognition*, 2004.
- [40] S. Wrede, M. Hanheide, C. Bauckhage, and G. Sagerer. An Active Memory as a Model for Information Fusion. In *Proc. Int. Conf. on Information Fusion*, 2004.
- [41] S. Wrede, W. Ponweiser, C. Bauckhage, G. Sagerer, and M. Vincze. Integration Frameworks for Large Scale Cognitive Vision Systems - An Evaluative Study. In *Proc. Int. Conf. on Pattern Recognition*, 2004.