

Integrating Perceptual and Cognitive Modeling for Adaptive and Intelligent Human–Computer Interaction

ZORAN DURIC, WAYNE D. GRAY, RIC HEISHMAN, STUDENT MEMBER, IEEE, FAYIN LI, AZRIEL ROSENFELD, MICHAEL J. SCHOELLES, CHRISTIAN SCHUNN, AND HARRY WECHSLER, FELLOW, IEEE

Invited Paper

This paper describes technology and tools for intelligent human–computer interaction (IHCI) where human cognitive, perceptual, motor, and affective factors are modeled and used to adapt the H–C interface. IHCI emphasizes that human behavior encompasses both apparent human behavior and the hidden mental state behind behavioral performance. IHCI expands on the interpretation of human activities, known as W4 (what, where, when, who). While W4 only addresses the apparent perceptual aspect of human behavior, the W5+ technology for IHCI described in this paper addresses also the why and how questions, whose solution requires recognizing specific cognitive states. IHCI integrates parsing and interpretation of nonverbal information with a computational cognitive model of the user, which, in turn, feeds into processes that adapt the interface to enhance operator performance and provide for rational decision-making. The technology proposed is based on a general four-stage interactive framework, which moves from parsing the raw sensory-motor input, to interpreting the user’s motions and emotions, to building an understanding of the user’s current cognitive state. It then diagnoses various problems in the situation and adapts the interface appropriately. The interactive component of the system improves processing at each stage. Examples of perceptual, behavioral, and cognitive tools are described throughout the paper. Adaptive and intelligent HCI are important for novel applications of computing, including ubiquitous and human-centered computing.

Manuscript received May 31, 2001; revised February 15, 2002.

Z. Duric, R. Heishman, F. Li, and H. Wechsler are with the Department of Computer Science, George Mason University, Fairfax, VA 22030-4444 USA (e-mail: zduric@cs.gmu.edu; rheishman@cs.gmu.edu; fli@cs.gmu.edu; wechsler@cs.gmu.edu).

W. D. Gray is with the Cognitive Science Department, Rensselaer Polytechnic Institute, Troy, NY 12180-3590 USA (e-mail: gray@rpi.edu).

A. Rosenfeld is with the Center for Automation Research, University of Maryland, College Park, MD 20742-3275 USA (e-mail: ar@cfar.umd.edu).

M. J. Schoelles is with the Department of Psychology, ARCH Laboratory/HFAC Program, George Mason University, Fairfax, VA 22030-4444 USA (e-mail: mschoell@gmu.edu).

C. Schunn is with the Department of Psychology, Learning Research & Development Center, LRDC, University of Pittsburgh, Pittsburgh, PA 15260 USA (e-mail: schunn@pitt.edu).

Publisher Item Identifier 10.1109/JPROC.2002.801449.

Keywords—Adaptation, behavioral performance, cognitive modeling, decision-making, feedback, human-centered computing, human–computer interaction (HCI), intelligent interfaces, interpretation of human behavior, nonverbal information, perceptual modeling, ubiquitous computing.

I. INTRODUCTION

Imagine a computer interface that could predict and diagnose whether the user was fatigued, confused, frustrated, or momentarily distracted by gathering a variety of *nonverbal* information (e.g., pupillary responses, eye fixations, facial expressions, upper-body posture, arm movements, and keystroke force). Further imagine that the interface could adapt itself—simplify, highlight, or tutor—to improve the human–computer interaction (HCI) using these diagnoses and predictions. Nonverbal information facilitates a special type of communication where the goal is to probe the inner (cognitive and affective) states of the mind before any verbal communication has been contemplated and/or expressed. This paper addresses the technology and tools required to develop novel computer interfaces suitable for handling such nonverbal information.

Assume now that a private, single-engine plane wanders into a commercial flight sector. The air traffic controller does nothing. Has she noticed the plane, evaluated its flight path, and concluded that it will shortly leave the sector without posing a threat to commercial aviation? Or has the plane slipped in unnoticed, and the controller has not yet considered the need to alert and reroute the five commercial flights in her sector? From the simple data that the computer gets from its operator (i.e., decisions made), it is impossible to know whether the busy controller’s attention should be directed to the intruder or left alone to focus on more urgent matters. However, at the time when the intruder entered the sector, the controller’s upper body was erect and tilted

forward in the seat, her point of gaze was within 1° of visual angle of the intruder, her pupils were dilated, her facial expression indicated surprise, and the force of her mouse click (on a commercial jetliner) was much less intense than normal. Imagine a computer interface that gathered such information about the operator and correctly diagnosed the operator's current cognitive state. With the above data it might decide to do nothing. With another combination of nonverbal information, it might decide to make the intruder's icon blink on and off. With a third combination, it might zoom the screen so that the intruder's icon was at the controller's point of gaze.

Less dramatic types of human-computer problems could also benefit by processing nonverbal information for adaptive and intelligent interfaces. For example, an operator repeatedly uses the mouse to gesture at (i.e., point at, circle, or otherwise indicate) an already classified target. Is he confused, frustrated, or simply fatigued? If he is confused, the interface could be automatically simplified (since optimal display complexity is relative to the expertise of the operator), or a tutorial could be offered during the next work lull. Again, this diagnosis and remedial action could be carried out if the computer had access to nonverbal information about the operator. Arm movements can indicate cognitive states like surprise and fatigue in addition to being a substitute for verbal communication suitable for deaf and/or mute people, and gestures can be appropriate for noisy environments.

Yet another example of using nonverbal information to infer the cognitive state of the user comes from pupillometry, the psychology of the pupillary response. There is general agreement that the pupils dilate during increased cognitive activity and constrict (or return to some previous baseline) when the activity decreases (e.g., when a particular problem has been solved and relaxation sets in). There is also evidence to support the assertion that the constant motion of the pupil (referred to as "pupillary unrest" or "hippus") is more accentuated under conditions of fatigue or drowsiness.

Knapp and Hall provide further evidence regarding nonverbal information in the context of expressions of emotions and their locations. In particular, they note [50] "rarely is the eye area tested separately from the entire face in judging emotions. Sometimes, however, a glance at the [brow and] eye area may provide us with a good deal of information about the emotion being expressed. For example, if we see tears we certainly conclude that the person is emotionally aroused, though without other cues we may not know whether the tears reflect grief, physical pain, joy, anger, or some complex blend of emotions. Similarly, downcast or averted eyes are often associated with feelings of sadness, shame, or embarrassment."

Nonverbal information as a new communication medium is most suitable for behavior interpretation. For example, the existing work on facial processing can now be extended to task-relevant expressions rather than the typical arbitrary set of expressions identified in face processing research. Moreover, the technology and tools proposed will have the added benefit of developing a framework by which one can improve on predictions of the consequences of various interface decisions on behavior—an important goal in the science of HCI. In particular, this paper emphasizes that human behavior en-

compasses both apparent performance and the hidden mental state behind performance. Toward that end we suggest an integrated system approach that can measure the corresponding perceptual and cognitive states of the user, and then can adapt the HCI in an intelligent fashion for enhanced human performance and satisfaction.

The outline of this paper is as follows. Section II provides the conceptual and intellectual framework needed to address issues related to adaptive and intelligent nonverbal interfaces. Section III describes recent research related to the interpretation of human activities, also known as W4 (*what, where, when, who*). The shortcomings of W4, since it is dealing only with the apparent perceptual aspect, are discussed and provide the motivation for our novel proposed methodology, W5+ (*what, where, when, who, why, how*); the *why* and *how* questions are directly related to recognizing specific cognitive states. Section IV describes in detail the W5+ methodology and motivates the choices made. In addition to migration from W4 to W5+, emphasis is placed on the fact that performance needs to be monitored in terms of both apparent (external) and internal behavior. (Contrast our framework with the more traditional use of Bayesian networks [65] and dynamic belief networks, which "meter" only the external behavior). Section V describes the tools required to implement the perceptual processing module, focusing on the interpretation of lower arm movements, facial expressions, pupil size, and eye-gaze location. Section VI describes the technology required for behavioral processing, focusing on the novel area of mouse gesture interpretation. Section VII overviews the components of embodied models of cognition and how they can be extended to include affect. Section VIII elaborates on how user interfaces can be adapted dynamically using the embodied model of cognition, and what additional issues need to be considered. We conclude the paper in Section IX with a summary of the novel W5+ technology and recommendations for further research and tool development.

II. BACKGROUND

HCI has developed using two competing methodologies [78]: direct manipulation and intelligent agents (also known as delegation). These approaches can be contrasted as the computer sitting passively waiting for input from the human versus the computer taking over from the human. Another dimension for HCI is that of affective computing [67]. Affective computing is concerned with the means to recognize "emotional intelligence." Whereas emotional intelligence includes both bodily (physical) and mental (cognitive) events, affective computing presently focuses mainly on the apparent characteristics of verbal and nonverbal communication, as most HCI studies elicit emotions in relatively simple settings [67]. Specifically, recognition of affective states focuses on their physical form (e.g., blinking or face distortions underlying human emotions) rather than implicit behavior and function (their impact on how the user employs the interface). In contrast to the established paradigms of direct manipulation and intelligent agents, intelligent human-computer interaction (IHCI) uses computer intelligence to increase the bandwidth through

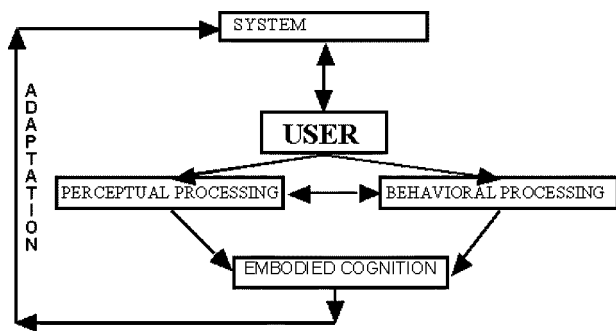


Fig. 1. System architecture for adaptive and intelligent HCI.

which humans interact with computers [63]. Nonverbal information such as facial expressions, posture, point of gaze, and the speed or force with which a mouse is moved or clicked can be parsed and interpreted by the computer to iteratively construct and refine a model of the human's cognitive and affective states. The availability of such users' models can be then used in an adaptive fashion to enhance HCIs and to make them appear intelligent, i.e., causal, to an outside observer.

It is not only computer technology that needs to change to make such novel interfaces a reality. People have to change as well and adapt to the interface that the computer presents them with. In the end, both people and the computer have to understand each other's intentions and/or motivations, provide feedback to each other as necessary, and eventually adapt to each other. W5+ systems are examples of novel intelligent interfaces, and they make the transition from HCI to IHCI where people and computers can augment each other's capabilities and display characteristics of team behavior. The methodology we propose for IHCI (see Fig. 1) integrates parsing and interpretation of nonverbal information with a computational cognitive model of the user that, in turn, feeds into processes that adapt the interface to enhance operator performance and provide for rational decision-making. Adaptive and intelligent HCI combines advanced work in perceptual recognition, machine learning, affective computing, and computational modeling of embodied cognition. Our methodology is based on a general four-stage, interactive framework. The system moves from parsing the raw sensory-motor input to interpreting the user's motions and emotions to building an understanding of the user's current cognitive state. It then diagnoses various problems in the situation and adapts the interface appropriately. The interactive component of the system improves processing at each stage. For example, knowledge of the user's current cognitive state helps predict changes in eye and head location, which in turn improves image parsing. We expect that our approach will have potential benefits for a broad class of HCIs. Moreover, our integrated methodology will also advance many areas of basic research (e.g., computer vision and facial processing, perception, cognition, human learning, and adaptation). We view our approach as a necessary step in developing IHCI systems where human cognitive, perceptual, motor, and affective factors are (fully) modeled and used to adapt the interface.

IHCI also promotes human activity and creativity. As part of emerging intelligent synthesis environments, IHCI supports human-centered and immersive computing, the infrastructure for distributed collaboration, rapid synthesis and simulation tools, and life-cycle system integration and validation. IHCI combines the (computational) ability to perceive mixed affordance (input) patterns, reasoning and abstraction, learning and adaptation, and finally, communication, language and visualization. These concepts echo those of Kant, for whom perception without abstraction is blind, while abstraction without perception is empty, and of Confucius, for whom learning without thought is useless, and thought without learning is dangerous.

Ubiquitous or pervasive computing, a new metaphor for computing, provides users with constant access to information and computation in the form of mobile and wireless computing devices—personal digital assistants (PDAs), such as cell phones, wearable computers, and appliances—that are endowed with intuitive user interfaces. Ubiquitous computing maintains computing at its core but removes it from our central focus. Computing is embedded into a surrounding, but almost invisible and friendly, world in order to facilitate collaborative work and the creation and dissemination of knowledge. Ubiquitous computing is more than virtual reality, which puts people inside a computer-generated world, and is more than PDAs. Ubiquitous computing involves explicit representations of oneself and other humans, possibly using avatars, and representation of cognitive, affective, social, and organizational aspects of the human behind the avatar (e.g., natural and expressive faces and gestures, representing and reasoning about others' places in organizational systems, and social relationships). The boundaries between the “real world,” augmented reality, and virtual environments are blurred to create a mixed reality. Unlike virtual reality, mixed reality seeks to enhance the real environment, not to replace it. That is, while an interface agent or PDA will alert a pilot of an impending collision, ubiquitous computing will display for the pilot airspace information that provides continuous spatial awareness of surrounding objects. The mixed reality metaphor avoids making systems appear too human-like in cases where they have very limited intelligence and are brittle in their interaction.

IHCI makes ubiquitous computing possible by continuously adapting the interface medium to meet specific user needs and demands. The emergence of human-centered interaction with intelligent systems buttresses the utilization of both verbal and nonverbal communication to create a richer, more versatile and effective environment for human activity. Human-centered design is problem-driven, activity-centered, and context-bound and employs computing technology as a tool for the user, not as a substitute. Thus, the emphasis is on supporting human activity using adaptive and intelligent interfaces rather than on building (fully) autonomous systems that mimic humans. One approach to a human-centered use of intelligent system technology seeks to make such systems “team players” in the context of human activity, where people and computer technology

interact to achieve a common purpose. Another possible approach focuses on building effective computational tools for modeling, interpreting, fusing, and analyzing cognitive and social interactions such as speech, vision, gesture, haptic inputs, and/or affective state expressed using body language. The goal of IHCI is to expand on human perceptual, intellectual, and motor activities.

People tend to misjudge the bounds of systems' capabilities, ranging from over-reliance on system performance—in cases where it is inappropriate to do so—to loss of trust and lack of acceptance in situations where the system performs well. Ubiquitous computing seeks to derive benefits from a truly complementary relationship with the human partners, forcing computing to thrive in an integrated (virtual and physical) world with people and to become predictable, comprehensible, and informative for the human partners. To forge a trusted partnership and expand human intellect and abilities, computing has to become compliant with our demands, communicative regarding our processes, and cooperative with our endeavors.

Ubiquitous computing emphasizes distributed affordances instead of focused expertise. People are most effective when they are fully engaged, mind and body, in the world. What becomes apparent is situated (grounded) computing, leading to practical intelligence facilitated by shared context acquired at the intersection of perception (senses, affective and physiological), language and communication, thought and reason, and action (purposive and functional). Shared context leads to recoordination of human behavior and subsequent reconceptualization. We will soon reach the point when computing power and storage are cheap commodities relative to the challenge of making them readily and effectively available to everyone, everywhere. The vision of ubiquitous computing can become reality only with the emergence of several components, and IHCI is one of them. Wearable devices need to supply correct and just-in-time information and reduce the wearer's cognitive load. Such devices combine communication, computation, and context sensitivity, while also supporting increased collaboration. Users must be able to leave their desktops and continue their daily tasks, possibly engaging in hands-free tasks, while still remaining connected to computing and communication resources. With wearable computers, interaction with the real world is the primary task, making current windows–icons–menu–pointers (WIMP) interfaces obsolete. An emerging augmented system that interfaces people and computing devices promotes mutual understanding—including background, goals, motivations and plans—and optimal sharing of the computational load.

III. REVIEW OF W4: WHAT, WHERE, WHEN, WHO

We briefly review here recent research related to the interpretation of human activities (W4). As will become apparent from our discussion, W4 deals only with the apparent perceptual aspect and makes no reference to the cognitive element

responsible for that aspect. Extensive work on analyzing images of humans' activities began in the 1990s. Many references are listed in [69]. A large number of papers on face and gesture recognition were presented in the four international conferences on the subject [89]–[92]. A very good review of early work on motion understanding approaches and applications was done by Cedras and Shah [17]. A review of research papers on hand gesture recognition for HCI was done by Pavlovic *et al.* [63], and a broader review of research papers on visual analysis of human motion was done by Gavrilu [28]. A review of papers on nonrigid motion analysis, in particular on articulated and elastic motion, was presented by Aggarwal *et al.* [1]. A comprehensive review of various methods for computer vision-based capture of human motions was recently done by Moeslund and Granum [59]. In the remainder of this section, we review recent work on motion analysis and understanding because it is the primary force behind human activities. Three main criteria can be used to classify research on human motion analysis. First, the research can be classified in terms of the tasks that it focuses on: detection, tracking, or recognition. Second, it can be classified in terms of the models used to represent objects and humans. Third, it can be classified in terms of the control mechanisms used.

Detection of humans in static or video images has mostly been addressed through background subtraction and matching. A background subtraction method that uses colors and edges was described in [42]. Some authors have used background subtraction as part of a system combining detection, body labeling, and tracking of humans [21], [29], [37], [85]. In some cases, cues such as skin color have been used to detect humans in images [41]. Other authors have used motion from single or multiple cameras to detect, label, and track humans or their body parts in video images [46], [48], [68], [72], [74], [83], [85]. Some authors have approached this problem as one of matching. Humans or their parts have been detected and tracked as configurations of points such as light displays, markers, and image features [81], as configurations of edges [30], [51]–[53], [62], [66], [73], and as collections of particularly shaped strips [27], cylinders, or superquadrics [22], [29], [72], [79], [83]. For tracking, some authors have focused on using motions of image points and edges. Human models have been initialized by hand in the first frame of each sequence [14], [18]. Some authors have considered the problem of action/activity/gesture recognition for humans using shape and/or motion information [8], [12], [20], [25], [26], [32], [36], [44], [45], [71], [82], [86], [87]. Dynamic recognition, most appropriate for interpreting video sequences, is done using recursive and neural networks, deformable templates, spatio-temporal templates [60], and graphical models [12] because they offer dynamic time warping and a clear Bayesian semantics for both individual (HMM) and interacting or coupled (CHMM) generative processes [61]. Finally, some authors have implemented systems that combine detection, tracking, and recognition [2], [15], [20], [37], [39], [57], [58], [68].

A second set of criteria that can be used for classifying research on human motions is based on how to model humans. Humans have been modeled as elongated, blob-like shapes either implicitly [20], [37], [39], [57], [58], [68] or

explicitly [30], [62], [66]. Deformable models have been utilized for body part (hand) and facial feature tracking [12], [73], [86]. Some authors have modeled humans as articulated stick figures [2], [51]–[53], [71], [81]; this approach has been particularly effective for moving light display analysis. Finally, humans have been modeled as articulated objects, where parts correspond to blobs [85], strips [14], [27], [46], tapered superquadrics [29], [48], or cylinders [22], [72], [74], [79], [83].

A third set of criteria that can be used for classifying research on human motions is based on the mechanisms used to control search in detection, tracking, and recognition. Kalman filtering has been used frequently; examples include [18], [48], [72], [74], and [83]. More recently, Bayesian inference has been used [22], [66], [73], [79]; these methods are also known as Condensation. Other strategies that have been used include search algorithms such as best-first [29] and/or “winner take all” [18], [57], [87].

Bobick [11] recently proposed a taxonomy of movement, activity, and action. In his taxonomy, movements are primitives, requiring no contextual or sequence knowledge in order to be recognized. Activities are sequences of movements or states, where the only knowledge required to recognize them involves statistics of the sequence. According to Bobick, most of the recent work on gesture understanding falls in this category. Actions are larger scale events which typically include interactions with the environment and causal relationships. An important distinction between these levels is the degree to which time must be explicitly represented and manipulated, ranging from simple linear scaling of speed to constraint-based reasoning about temporal intervals.

Other related work includes biomechanics and human modeling in computer graphics and movement notations in choreography. In biomechanics, researchers are interested in modeling forces and torques applied to human bodies and tissues during various physical activities [88]; these models provide tools for analyzing the relationship between movements and actions. In computer graphics, researchers are interested in producing realistic images for virtual reality applications, human factor analysis and computer animation [7], [8], [10], [38]. Formalisms for describing the motions of humans include movement notations [7], [8] such as Labanotation [35], which is mostly used for dance, and Eshkol–Wachmann notation [24], [35], which is also used for sign languages [19].

IV. ADAPTIVE AND INTELLIGENT HCI—METHODOLOGY

Our methodology is quite general and has been outlined in Fig. 1. The main modules are perceptual processing, behavioral processing, embodied cognition, and adaptive system interface. The user is interacting with an adaptive system interface, which changes as a function of the current task state and the cognitive or mental state of the user. The nonverbal front end includes the perceptual and behavioral processing modules, and its input consists of raw sensory information about the user. The perceptual module processes images of the face, the eye (gaze location and pupil size), and the upper body and analyzes their relative motions; the

behavioral module processes information about actions done to the computer interface directly, such as keystroke choices, the strength of keystrokes, and mouse gestures. Both the perceptual and behavioral modules provide streams of elementary features that are then grouped, parsed, tracked, and converted eventually to *subsymbolic, summative affective representations* of the information in each processing modality. In other words, one output of the perceptual and behavioral processing modules is a stream of affective states at each point in time. States that could be recognized include confusion, fatigue, stress, and other task-relevant affective states. The quest for subsymbolic and summative affective representations is motivated by abstraction and generalization, communication, and reasoning, in particular, and by the perception-control-action cycle [84]. Furthermore, “Signal and symbol integration and transformation is an old but difficult problem. It comes about because the world surrounding us is a mixture of continuous space–time functions with discontinuities. Recognition of these discontinuities in the world leads to representations of different states of the world, which in turn place demands on behavioral strategies. Similarly, an agent’s (biological or artificial) closed-loop interactions with the world/environment can be modeled as a continuous process, whereas switching between behaviors is naturally discrete. Furthermore, the tasks that are either externally given to the agents or internally self-imposed prespecify and, hence, discretize an otherwise continuous behavior. Thus, we have three sources for discretization of the agent-world behavioral space: 1) natural space–time discontinuities of the world; 2) the model of agent-world dynamics during execution of a given task; and 3) the task. Furthermore, in computer vision, symbols served mainly as a data reduction mechanism, while in AI the following were missing: 1) explicit acknowledgment that the transformation from signals to symbols results in the loss of information; 2) self-correction and updating mechanisms of the obtained symbolic information; and 3) explicit models of the dynamic interaction between an agent and its world. Symbols not only provide nice abstractions for low-level strategies, but also allow us to move one level up the modeling hierarchy and observe the properties of the systems and their interactions among each other and their environment at a more macroscopic level. Symbolic representation mediates reasoning about the sequential and repetitive nature of various tasks” [9]. The adaptive and intelligent HCI methodology proposed in this paper addresses the problems raised above using embodied cognition to connect the apparent perceptual and behavioral subsymbolic affective representations and symbolic mental states, and in the process adaptively derive the summative subsymbolic states from raw signals and also adapt the user/system interface for enhanced performance and human satisfaction. The task description language chosen for manipulation tasks using such an adaptive and Intelligent HCI is that of the ACT-R/PM cognitive architecture (see Section VII).

These affective subsymbols are fed into the embodied cognition module and mediate fusion and reasoning about possible cognitive states. While the subsymbols correspond to external manifestations of affective states, the cognitive

states are hidden and not directly observable. The embodied cognition module generates hypotheses about possible task-relevant cognitive states, resolves any resulting ambiguity by drawing from contextual and temporal information, and optimally adapts the interface in order to enhance human performance. Knowledge of the user's state and the system's state are used to diagnose potential problems, and these diagnoses trigger adaptive and compensatory changes in the computer interface. While this process has been described as a linear process, in fact it is an interactive process in which information from later phases can feed back to augment processing in earlier phases. For example, knowledge of current cognitive activities can be used to improve recognition of affective states.

The development of the perceptual and behavioral processing and embodied cognition modules and their interactions is the key to making progress on constructing adaptive intelligent interfaces. The embodied cognition module is a novel addition to the more traditional HCI approaches as: 1) it bridges between direct manipulation and intelligent agents through physical modeling of the users and 2) it augments emotional intelligence through cognitive modeling of user behavior. The capability of modeling the effects of the affective state on cognitive performance will have an impact on the choice of models as well as the computational techniques employed.

The experimental platform envisions continuous operation of the system over extended periods of time with a single user in a particular task environment. While it is possible to personalize the user interface and be able to detect characteristics of the user automatically, one can also assume that the system will be properly initialized for a particular user. This would comprise initialization of hardware and software by acquiring both a physical (perceptual and behavioral) and cognitive model of the user, calibration of the video cameras and eye tracker, and detection of the face and upper body. Initial localization of facial landmarks is crucial for successfully tracking features over time. For a first-time user, generic information concerning the user's expertise and experience is used to initialize the cognitive model. For return users, the cognitive model from the user's last session is retrieved and instantiated in the current setting.

The raw video data for the perceptual processing module include color and intensity. Standard low-level computer vision techniques are used to extract additional features such as corners, edge information and motion feature maps. The feature extraction process is guided by the initialized model, for which the spatial locations of head, upper body, eyes, and mouth have been determined. Our previous research on perceptual processing includes detecting faces and the upper body [40], [42], [57], [58] and automatic detection of facial landmarks [80]. The measurements acquired (color, intensity, edge information, and motion) can be considered as first-order features. These features are combined together in order to acquire reliable estimates of the shape (contour) and motion fields of the eye and mouth facial regions, and/or arms. In the next immediate level of the hierarchy, lower order parametric descriptions of the shapes and motion fields associated with smaller spatial regions corresponding to the

eyes and mouth are sought. The modes of these parametric descriptions accumulated over time are processed using learned vector quantization (LVQ) and yield subsymbolic indicators contributing to the assessment of the affective state. A variety of eye movements involving pupils, irises, eyelids, and eyebrows are captured as different modes. For example, the movement of the eyelids can reveal information about the stress level and tension of the user. At this level we also capture lower arm and hand movements; Section V provides an additional description of this module. The last level of the hierarchy uses upper body shape and motion information. One can estimate independently the 3-D poses of the head and shoulders, which can undergo independent or combined motions. Information on the two 3-D poses is then abstracted using modal analysis and then fed into the embodied cognition module.

The processing of raw eye data (pupil location and size) requires additional computations and is currently performed most effectively using special-purpose hardware. The eye tracker data includes time stamped x - y coordinates of the point of gaze (POG) and the diameter of the pupil (at 60 samples or more per second). The particular state information corresponds to a spatial location where a fixation occurred and the transitions between the states (events) correspond to saccadic eye movements. Changes in recorded eye positions are parsed as fixations or saccades according to their magnitudes and directions. Eye blinks are calculated from the raw data by detecting consecutive samples with zero pupil dilation. An eye blink is indicated if these samples span a time of 30–300 ms. Like fixation data, the number of eye blinks and rate of eye blinks between mouse clicks is calculated.

The *behavioral processing* module processes keystroke (choice and rate) and mouse data (clicks and movements). The keystroke data include key choice and timing of keystrokes. The mouse data include the time-stamped x - y coordinates of the pointer, the force and the frequency of mouse clicks, and the acceleration applied to mouse movements. The keystroke data are the primary means by which the cognitive model of the user is updated, through the process of model tracking (see below). Raw mouse data are collected at the same time that a raw eye data sample is collected. For the mouse, motion could be a movement from one position on the screen to another, and the dynamics would describe the applied force and the duration of the movement.

Parsing and interpreting the mouse data deserve additional notes, as they represent very novel uses of mouse data. The mouse data provide more than the obvious performance data about how fast and how accurately users make choices. We analyze the force data (how hard individuals click the mouse) and the trajectory data (how users move the mouse). The force data are divided into two dimensions: average force of a click and duration of the click. The trajectory data are treated as a form of gesture data. In other domains, we have found that people will gesture at various aspects of the screen using the mouse and these mouse gestures are indicators of preliminary cognitions [77]. For example, people sometimes circle objects, trace trajectories, or move rapidly between objects. Informally, we have seen the same behavior in the Argus domain (to be described below). To recognize mouse gestures,

we will use a technique that we developed in the context of sign language recognition [36]. There, hand gestures corresponding to American Sign Language are first located using projection analysis and then normalized in size, while recognition takes place using a hybrid mixture of (connectionist and symbolic) experts consisting of ensembles of radial basis functions (ERBFs) and decision trees (DTs).

The mouse data will be used in two different ways. First, the subsymbols corresponding to mouse gestures will be passed on to the embodied cognition module, with the goal of providing additional information about the aspects of the screen that are being attended to. Second, the addition of subsymbols corresponding to combined mouse gestures and force data will allow for the possibility of recognizing combinations of affective states in the user. That is, a person's face may reflect fatigue and the person's mouse gestures may reflect confusion. Research on hand gestures has often found that people can reflect different information about their internal cognition in speech than they do in co-occurring gestures [31]. Similarly, real affective states are often a combination of basic states, and our hypothesis is that the components of the combinations may be externalized simultaneously but in different external forms (e.g., fatigue in mouse movements, disgust in facial expressions).

One simple alternative to our approach would be to try to go directly from these diagnoses of affective states to adaptations of the interface (e.g., confusion = simplify interface). However, such a simple method is not likely to work because it does not take into account the cognitive state of the individual with respect to the task being performed. How to best adapt the interface will usually depend upon what cognitive operations are currently being performed. For example, simplifying an interface by removing information will only work if that information is not needed in the computations currently being performed by the individual. Moreover, affective states are often directed at particular aspects of the current task/interface. For example, a particular object on the screen or aspect of an interface is often a source of confusion, and it is better to clarify or simplify the offending object/aspect than to simplify random aspects of the interface or the entire interface, which would cause more confusion. The embodied cognition module uses the process of model tracing to understand the user's behavior, thereby making intelligent interface adaptation possible. In model tracing, a cognitive model is built that is capable of solving the human tasks in the same way as they are solved by the humans. The model is then aligned with the task and behavioral choice data (i.e., what the state of the world is and what the human chose to do) such that one can see which internal cognitive steps the human must have taken in order to produce the observed behavioral actions. Toward that end, the embodied cognition model also uses the affective subsymbols and their degrees of belief, derived earlier by the perceptual and behavioral processing modules. The embodied cognition module is described further in Section VII.

The *adaptive system interface* module adapts the interface to the current needs of the human participant. Different affective and cognitive diagnoses include confusion, fatigue, stress, momentary lapses of attention, and misunderstanding

of procedures. Different adaptations include simplifying the interface, highlighting critical information, and tutoring on selected misunderstandings. For instance, in one of the examples described earlier, if the force of the controller's mouse click and parsing of facial expressions concur in suggesting that the participant's visual attention is totally consumed by a commercial airliner, the system will intervene to alert the controller to the intruder's presence. Similarly, if later in her work shift, the controller's facial expressions and a wandering POG indicate a waning of attention, and the cognitive model interprets this as resulting from a decrease in cognitive resources (due to fatigue), steps may be taken to off-load parts of the tasks, to increase the salience of the most safety-critical components, or to relieve the controller. The types of interface adaptations that one can consider include: 1) addition and deletion of task details; 2) addition and deletion of help/feedback windows; 3) changing the formatting/organization of information; and 4) addition and removal of automation of simple subtasks. Further details will be given in Section VIII.

V. PERCEPTUAL PROCESSING

We describe here tools for perceptual processing, including lower arm movements, facial data processing, eye-gaze tracking, and mouse gestures. Additional tools are possible, including upper body posture (head and shoulders).

A. Interpretation of Lower Arm Movements

We describe next a method of detecting, tracking, and interpreting lower arm and hand movements from color video sequences (for details, see [23]). This method is relevant to parsing the raw sensory-motor input and in particular to interpreting the user's hand motions. It corresponds to perceptual processing (see Fig. 1) and its role is to transform signals to subsymbols expressing an affective state and suitable for the embodied cognition component. The method works as follows. The moving arm is detected automatically, without manual initialization, foreground, or background modeling. The dominant motion region is detected using normal flow. Expectation maximization (EM), uniform sampling, and Dijkstra's shortest path algorithm are used to find the bounding contour of the moving arm. An affine motion model is fit to the arm region; residual analysis and outlier rejection are used for robust parameter estimation. The estimated parameters are used for both prediction of the location of the moving arm and motion representation. In analogy to linguistic analysis, the processed sensory information is made compact and suitable for interpretation using LVQ, whose task is to abstract motion information. LVQ maps the affine motion parameters into a discrete set of codes {A, B, G, J, C, E, D, I, H}. The final transition from signals to a hierarchical and subsymbolic representation is enabled by clustering. In particular, clustering will map the discrete set of codes generated by LVQ into more abstract "subactivity" codes {up, down, circle} first, and finally into specific "activity" {pounding, swirling} subsymbols. Each activity or the expression of some cognitive state now corresponds to its own sequence of subsymbols and can

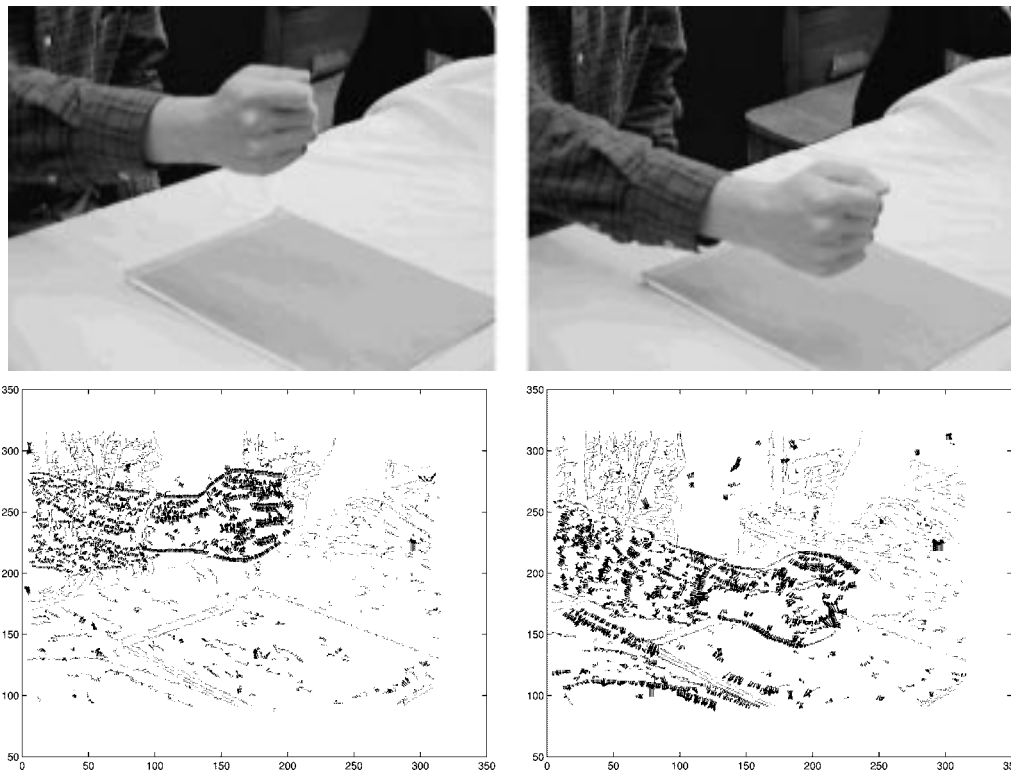


Fig. 2. The first frames from two image sequences and their corresponding normal flow. Left: the first frame from a “pounding” sequence of 400 frames. Right: the first frame from a “swirling” sequence of 100 frames. Upper row: images. Lower row: normal flows. These images were collected using a Sony DFW-VL500 progressive scan camera; the frame rate was 30 frames/s and the resolution was 320×240 pixels per frame.

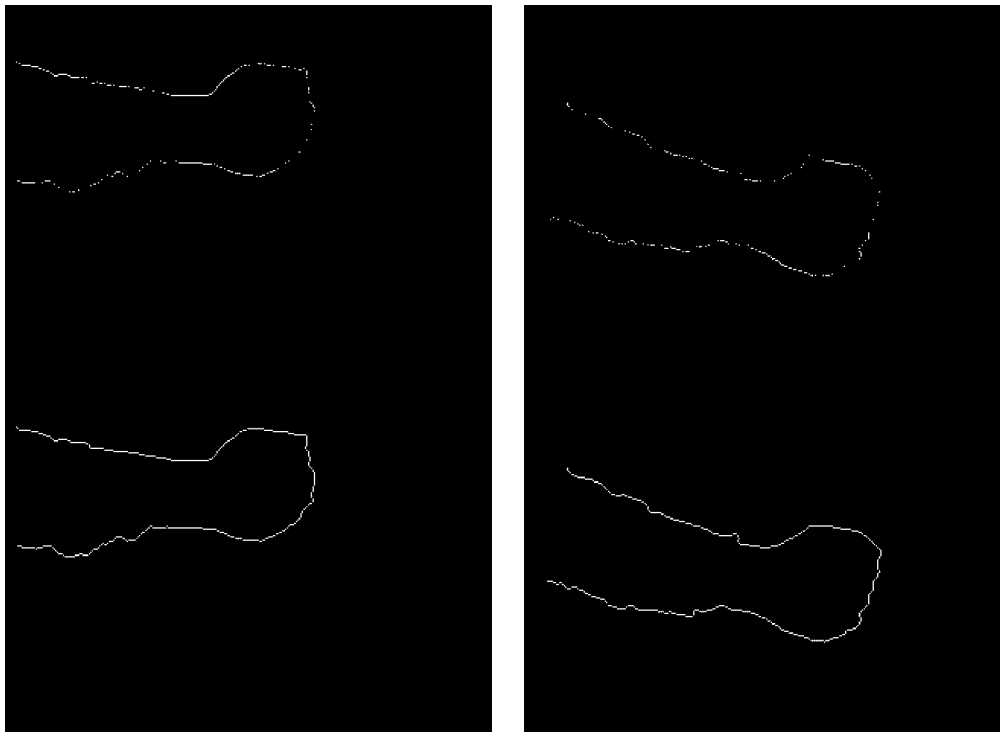


Fig. 3. Delineating the foreground objects for the images in Fig. 2. Top: points with high normal flow values and high gradient magnitudes. Bottom: foreground object outlines.

be properly distinguished from other activities or affective states of mind. Figs. 2–4 show some of the steps respon-

sible for parsing the raw signal to generate a subsymbolic description.

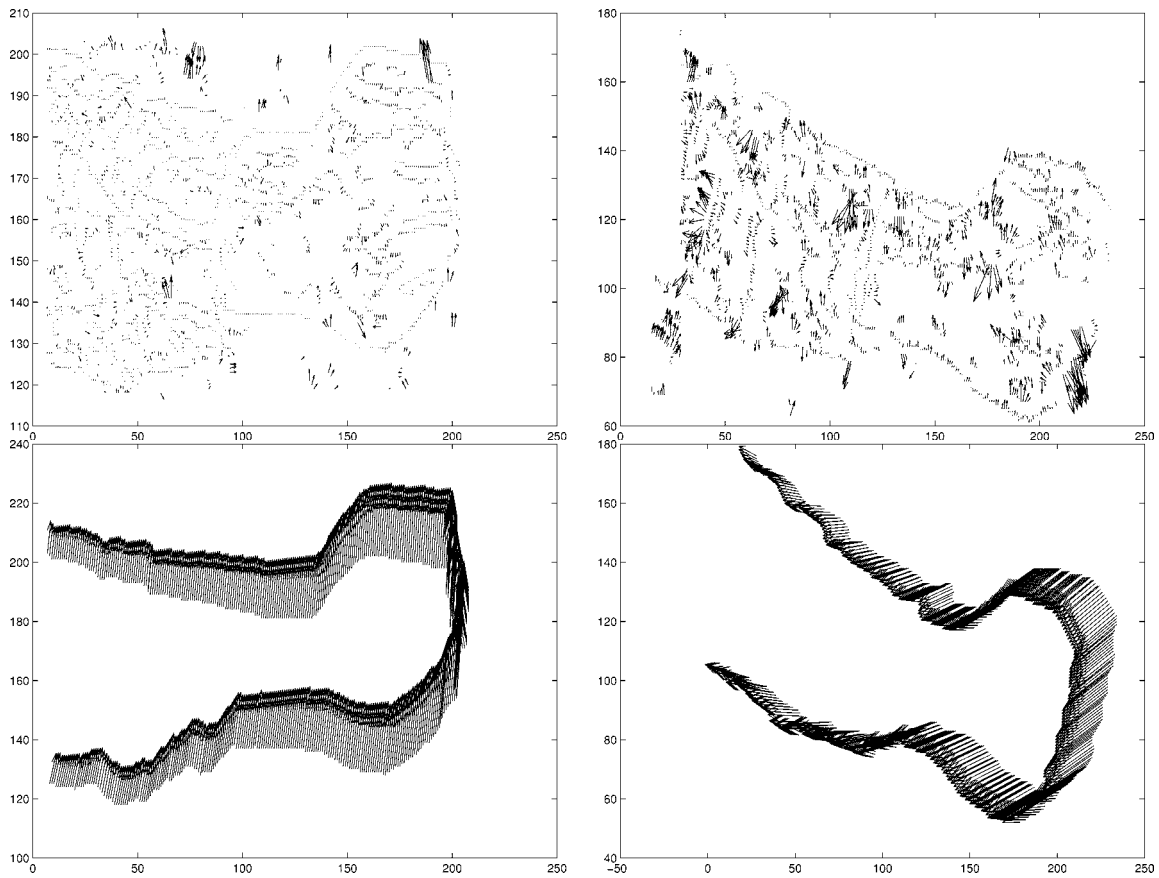


Fig. 4. Residual and reestimated flows for the detected arms in Fig. 2. Top: residual flow computed as the difference between the computed normal flow and the estimated affine normal motion field. Bottom: reestimated affine flow after outlier rejection.



Fig. 5. Facial expressions of an individual under (a) baseline, (b) easy, and (c) difficult load conditions.

B. Processing of Facial Data

The cognitive and emotional states of a person can be correlated with visual features derived from images of the mouth and eye regions [50]. Fig. 5 illustrates pilot data of facial expressions in a complex simulated radar classification task—expressions of the same subject in a baseline condition (9 targets to track), a low load condition (5 targets), and

a high load condition (30 targets). There are detectable differences of the type that one would expect, but the differences are subtle; in particular, the mouth and eye regions display an increase of tension for the difficult task.

For the eye region, the visual features related to cognitive states of a person include gaze direction (position of the irises relative to the eyes), pupil dilation, and the degree of occlusion of the iris and the pupil by the eyelids. For ex-

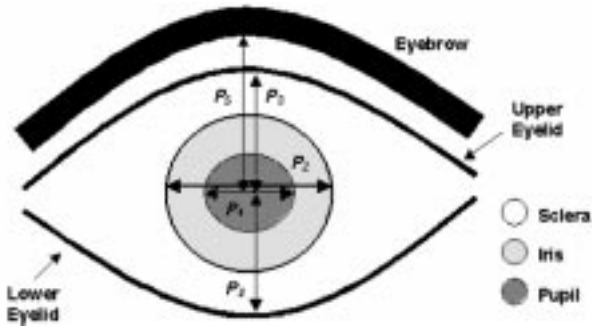


Fig. 6. Eye parameters: pupil diameter (P_1), iris diameter (P_2), distance from center of the iris to the upper eyelid (P_3), distance from center of the iris to the lower eyelid (P_4), and distance from center of the iris to the eyebrow (P_5).

ample, the pupil dilation indicates arousal or heightened cognitive activity, while an averted gaze may indicate increased mental activity associated with the processing of data. The visual features related to emotional states include the degree of openness of the eyes, the positions of the eyelids relative to the irises, the positions and shapes of the eyebrows (arched, raised, or drawn together), and the existence and shape of lines in particular eye regions (eye corners, between the eyebrows, below the lower eyelids). For example, surprise is indicated by wide-open eyes with the lower eyelids drawn down, and raised and arched eyebrows; fear is indicated by wide-open eyes with the upper eyelids raised (exposing the white of the eye) and the lower eyelids tensed and drawn up, and with the eyebrows raised and drawn together; happiness is indicated primarily in the lower eyelids, which have wrinkles below them and may be raised but are not tense.

Fig. 6 illustrates some of the parameters measured on the eyes. Note that we expect that the parameters P_1 – P_5 will always be positive; P_3 and/or P_4 can only become small when the eyes are almost closed. Note that other parameters, such as density of lines around the eyes and curvature of the eyelids and the eyebrows, can be added to complement these parameters.

The eye parameters are acquired from high-resolution color images. The color edges are detected by combining the gradients in red, green, and blue color bands [43]; the edges are thinned using nonmaxima suppression. Irises and pupils are located via generalized Hough transform using multiple size templates; it is assumed that their shapes are always circular. Eyelids and eyebrows are detected using a variation of the method described in [80]; the edges in the vicinity of irises are labeled as the candidates for the eyelids and the eyebrows. The left column of Fig. 7 shows eye regions corresponding to anger, surprise, and happiness; the right column of the figure shows the results of detecting the irises, eyelids, and eyebrows for the corresponding images on the left.

Numerical results for the examples shown in Fig. 7 are as follows. For the images showing anger (top row), the irises were detected at positions (114, 127) and (499, 103) with radii $P_2 = 32$. For the left eye (in the image) the distance from the upper eyelid to the iris was $P_3 = 20$, the distance from the lower eyelid to the iris center was $P_4 = 39$, and the

distance from the iris center to the eyebrow was $P_5 = 38$. For the right eye (in the image) the distance from the upper eyelid to the iris was $P_3 = 11$, the distance from the lower eyelid to the iris center was $P_4 = 35$, and the distance from the iris center to the eyebrow was $P_5 = 26$. For the images showing surprise (middle row) the irises were detected at positions (100, 173) and (483, 149) with radii $P_2 = 32$. For the left eye (in the image) the distance from the upper eyelid to the iris was $P_3 = 36$, the distance from the lower eyelid to the iris center was $P_4 = 34$, and the distance from the iris center to the eyebrow was $P_5 = 71$. For the right eye (in the image) the distance from the upper eyelid to the iris was $P_3 = 36$, the distance from the lower eyelid to the iris center was $P_4 = 39$, and the distance from the iris center to the eyebrow was $P_5 = 77$. For the images showing happiness (bottom row) the irises were detected at positions (148, 139) and (497, 115) with radii $P_2 = 30$; for the left eye (in the image) the distance from the upper eyelid to the iris was $P_3 = 14$, the distance from the lower eyelid to the iris center was $P_4 = 27$, and the distance from the iris center to the eyebrow was $P_5 = 49$. For the right eye (in the image) the distance from the upper eyelid to the iris was $P_3 = 17$, the distance from the lower eyelid to the iris center was $P_4 = 33$, and the distance from the iris center to the eyebrow was $P_5 = 48$. Computing the ratios $P_3/P_2, P_4/P_2, P_5/P_2$ we obtain the following results for the left and right eye pairs: for anger (0.34, 1.09, 0.81) and (0.62, 1.22, 1.19); for surprise (1.13, 1.06, 2.22) and (1.13, 1.22, 2.4); and for happiness (0.47, 0.9, 1.63) and (0.57, 1.1, 1.6).

C. Eye-Gaze Tracking

Because people can consider and discard various aspects of a task rather quickly (in less than 200 ms), eye movements can provide detailed estimates of what information an individual is considering. Eye tracking is becoming an increasingly popular online measure of high-level cognitive processing (e.g., [55]). By gathering data on the location and duration of eye fixations, psychologists are able to make many inferences about the microstructure of cognition. The use of eye tracking in estimating cognitive states rests on the immediacy assumption (people process information as it is seen) and the eye-mind assumption (the eye remains fixated on an object while the object is being processed). As long as the visual information requires fine discrimination, these assumptions are generally considered valid, but when the visual information is very coarse-scale, people can process the information without fixating on it.

In order to reliably separate eye fixations from saccades, one needs to sample gaze data at least 60 times per second with an accuracy of at least 2° of visual angle. A variety of eye-tracking methods exist. In terms of the data collected from the eye, two popular methods are: 1) shining a light on the eye and detecting corneal reflection and 2) simply taking visual images of the eye and then locating the dark iris area. Which method is best depends upon the external lighting conditions.

To compute where in the world a person is fixating, there are three popular methods. The first method simplifies the



Fig. 7. Left column: eye regions displaying (from top to bottom) anger, surprise, and happiness. Right column: processed eye regions.

calculations by having fixed geometries by forcing the person to hold still by biting on a bar or putting the head in a restraint. The second method has the person wear a head sensor that tracks the head orientation and location in three dimensions and then combines this information with eye-direction information. The third method places the eye-tracking apparatus on the person's head along with a scene camera so that a visual image is displayed showing what the person is currently looking at, with a point on the image indicating the object being fixated. To achieve the high levels of accuracy required, all three methods require recalibration for each individual being tracked in a given session; however, methods exist for automatic recalibration. While the first method of computing the point of fixation is the most accurate, it is purely a research method with no applicability to IHCI for obvious practical reasons. If one wants to track upper-body and facial gestures at the same time, a head-mounted camera is not practical either. The remote camera is the least accurate method, but extreme levels of precision are probably not needed for IHCI.

To separate out fixations from the raw point-of-regard data, the most popular method is to use a movement/time threshold: whenever the distance between consecutive points of regard is below a threshold for a sufficient length of time, a fixation is assumed. A more sophisticated and accurate approach uses a centroid submodel tracing methodology developed by Salvucci and Anderson [75]. The methodology involves categorizing eye movements using hidden Markov models (HMMs) and model tracing. Raw eye data are first categorized into fixations and saccades using a two-state HMM given velocity information alone. The centroid of each fixation is then determined. One could examine which object on the screen is closest to this centroid and simply assume the person was looking at that object. However,

because of noise in the eye data, this would frequently miscategorize the fixation. Instead, another HMM fitting process is used which takes into account the closeness of each fixation to objects on the screen AND the context of which other objects were just fixated. This model fitting process is done by comparing the sequence of fixations with all plausible sequences of fixations and selecting the sequence with the highest probability (best overall fit).

Fig. 8 presents an example of point-of-regard data extracted while a user interacts with a complex computer display. From the locations of the fixations, one can determine which objects were likely encoded at a detailed level. From the durations of the fixations, one can determine which objects were most likely involved in more detailed computations.

VI. BEHAVIORAL PROCESSING

Behavioral processing focuses on two kinds of data input: keyboard and mouse. Both the keyboard and mouse data are first used as primary inputs into the computer interface. We are not proposing that perceptual processing sources of information replace the mouse and keyboard. In addition to serving as direct interaction with the interface, keyboard, and mouse input will have additional functions of providing insights into the cognitive and affective states of the user. Keystroke data will provide information about the cognitive state of the user through the process of model tracing, which will be described in Section VII. Mouse data will provide information about user cognition and user affect; this process is described next.

Mouse data can be divided into two primary types: mouse gestures and mouse move and clicks. Mouse gestures are movements of the mouse that do not result in mouse clicks.

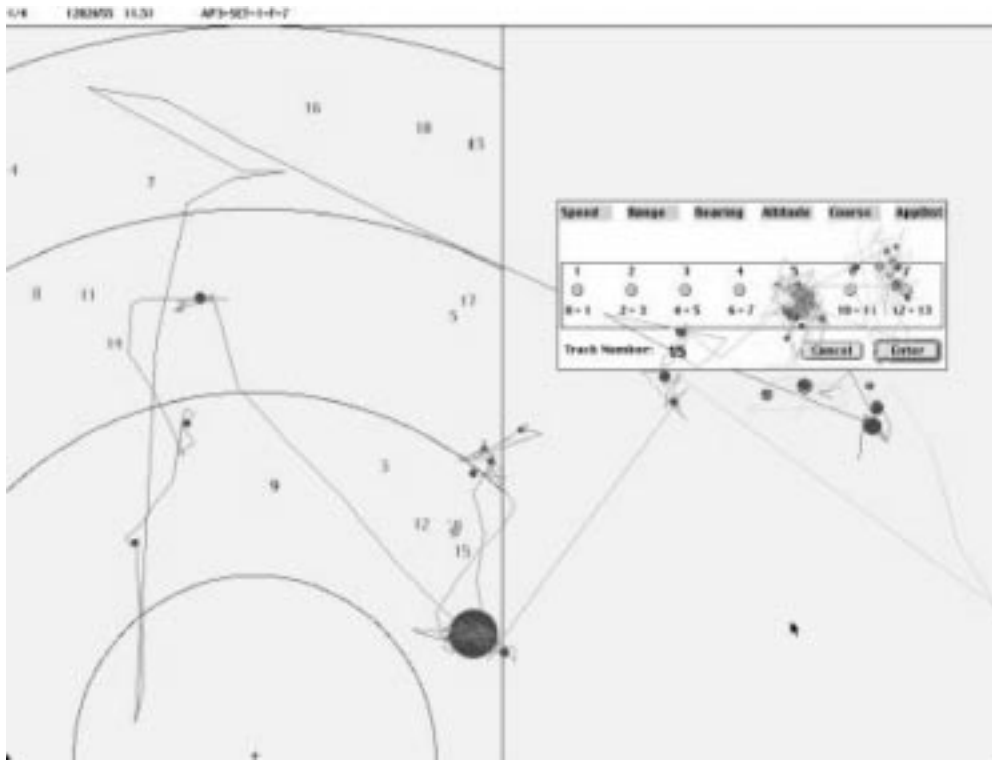


Fig. 8. Example of eye-tracking data from a complex simulated radar classification task. The moving targets to be identified are represented by the numbers in the left half of the screen. The details of target track 15 are in the table on the right half of the screen. Overlaid on the interface image is 10 s of POG data (changing from dark to light over time). Identified eye fixations (minimum = 100 ms) are indicated with disks (larger disks for longer fixations).

In most HCI environments, these movements are nonfunctional, although some systems provide rollover information. For the moment, we treat rollover cases as if they involve a mouse click. The mouse gestures, the nonfunctional mouse movements, can be viewed as windows into the mind. That is, they indicate what objects are currently being reasoned about [77]. Common mouse gestures include circling objects, linking objects, circling groups of objects, and shape tracing of large objects or groups. Object circling represents indecision about an object. Repetitive movement between objects is a linking gesture, and like circling groups of objects, represents a categorical decision that the linked objects are similar along an important dimension. Shape tracing of larger objects represents reasoning about the shape of the object.

Move-and-click movements have three important dimensions: speed of movement, force of click, and directness of movement to the clicked object. These dimensions are indicators of the general level of arousal and of indecision and confusion. Slower movements and weaker clicks combined with direct movements indicate low levels of arousal (i.e., fatigue). Slower movements and weaker clicks combined with indirect movement indicate confusion. Fast movements with strong clicks combined with indirect movements indicate frustration.

To recognize mouse gestures, one can use a technique that we developed in the context of sign language recognition [36]. There, hand gestures corresponding to American Sign Language are first located using projection analysis and then

normalized in size, while recognition takes place using a mixture of (connectionist and symbolic) experts consisting of ERBFs and DTs. ERBFs display robustness when facing the variability of the data acquisition process, using for training both original data and distortions caused by geometrical changes and blurring, while preserving the topology that is characteristic of raw data. In addition, ERBFs are similar to boosting, as each of their RBF components is trained to capture a subregion of the perceptual or behavioral landscape and can then properly recognize it. Inductive learning, using DTs, calls for numeric to subsymbolic data conversion, suitable for embodied cognition. The ERBF output vectors chosen for training are tagged as “CORRECT” (positive example) or “INCORRECT” (negative examples) and properly quantized. The input to the DT implemented using C4.5 [70], consists of a string of learning (positive and negative) events, each of them described as a vector of discrete attribute values. To further parse the move-and-click movements, one can use HMMs because the user is likely to stay in a given affective state for several movements. The states of the Markov model are the diagnosed affective states (alertness, fatigue, and confusion).

VII. EMBODIED COGNITION

The embodied cognition module has at its core an embodied cognitive model and a model tracing function. A cognitive model is capable of solving tasks using the same cog-

nitive steps as humans use to solve the tasks. An embodied cognitive model also has affective states to match the user's states and the ability to perceive and interact with an external world as the user does. In model tracing, the model is aligned with the task and behavioral choice data (i.e., what the state of the world is and what the human chose to do), so that one can see which internal cognitive steps the human must have taken in order to produce the observed behavioral actions. Toward that end, the embodied cognition model also uses the affective subsymbols and their degrees of belief, derived earlier by the perceptual and behavioral processing modules.

Currently the best way to build models of embodied cognition is to use a cognitive architecture (e.g., ACT-R, Soar, EPIC, 3CAPS) that has a relatively complete and well-validated framework for describing basic cognitive activities at a fine grain size. The currently most developed framework that works well for building models of embodied cognition is ACT-R/PM [16], a system that combines the ACT-R cognitive architecture [5] with a modal theory of visual attention [6] and motor movements [49]. ACT-R is a hybrid production system architecture, representing knowledge at both a symbolic level (declarative memory elements and productions) and subsymbolic level (the activation of memory elements, the degree of association among elements, the probability of firing productions, etc.). ACT-R/PM contains precise (and successful) methods for predicting reaction times and probabilities of responses that take into account the details of and regularities in motor movements, shifts of visual attention, and capabilities of human vision. The task for the embodied cognition module is to build a detailed mapping of the interpretations (i.e., motion/affective state) of the parsed sensory-motor data onto the ACT-R/PM model.

One can extend ACT-R/PM to make it a true model of embodied cognition by incorporating the effects of affect on performance. For example, in addition to handling the interactions among memory, vision, and motor movements, the model becomes fatigued over time and distracted when there is too much to attend to. Better than merely *becoming* fatigued and distracted, such an extended ACT-R/PM can *model the effects of fatigue and distraction* on memory, vision, and motor behavior and thereby on performance. Like people, as the model becomes fatigued, several changes may occur. First, the model slows down (increasing the interval between physical actions and shifts in visual attention, as well as increasing the time needed to store or retrieve information from memory). Second, the accuracy of its responses decreases (this includes the physical accuracy due to increased noise in eye-hand coordination and mental accuracy due to increased noise in memory retrieval, e.g., retrieving the target's old, rather than current, flight information). Third, the model becomes distracted, losing its focus of attention (running the risk of applying the "right" response to the "wrong" object or the "wrong" response to the "right" object). Fourth, it becomes narrower in what it chooses to encode [56].

This incorporation of affective subsymbols into models of embodied cognition is a product in its own right. Such a capability can be applied to other task environments (simu-

lated or prototypes of real systems) to determine changes in human performance over time. As such, models of embodied cognition could become an important tool for designers of real-time safety-critical systems (see, e.g., [34]). One novelty here lies in using a broader range of nonverbal data in guiding the model tracing process. Recent work in cognitive science suggests that nonverbal information, such as gestures, provides important insights into an individual's cognition [3]. Mouse gestures, the eye data, and affective states are important tools to improve this model tracing process.

One can explore three qualitatively different methods of incorporating affect into the cognitive model. First, affect can be thought of as directly modifying parameters in the cognitive model to produce relatively simple changes in behavior. For example, fatigue may affect processing speed (i.e., how fast someone thinks) as well as working memory capacity (i.e., how much information can be kept in mind). Similarly, confusion or frustration may influence the noise parameter in the decision process (i.e., the likelihood of making nonoptimal choices) or the threshold amount of effort a person will expend on a task (which influences the probability of giving up). Parameters controlling processing speed, working memory capacity, noise, and effort expended are formally defined within the ACT-R architecture. Second, affect can also change more structural or strategic aspects of the cognitive model. For example, when people become confused, fatigued, or frustrated, they may adopt an entirely different way of thinking about the task and making choices (i.e., alternative strategies). Thus, the performance parameters of the cognitive model may be held constant, but the action and decision rules themselves may change as the affect changes. A third possibility is some combination of these two types of changes in the model with changing affect. Individuals use a wide variety of qualitatively different strategies to solve any given type of problem [76], and changes in model performance parameters are likely to produce changes in strategy choice. For example, in a classic decision making study, Payne *et al.* [64] showed that, as the cognitive effort required for task performance increased (thereby placing greater demands on a limited-capacity working memory system), the decision-making strategies that people adopted changed. Lohse and Johnson [55] showed that changes in decision-making strategies were also induced by tradeoffs between perceptual-motor versus cognitive effort. Hence, it may well be that changes in strategies induced by changes in affective state are mediated by changes in underlying cognitive parameters. ACT-R contains clear predictions of how certain parameter changes will influence strategy choice, assuming a good characterization of the features of each strategy.

The process of model tracing keeps the model aligned with the user. It takes as primary input the behavioral interactions with the interface (i.e., the keystroke and mouse click data) and tries to match symbolic steps in the model. In a production system model, this amounts to matching to sequences of production firings. There are three factors that shape the model tracing process. First, any realistic model of human cognition acknowledges some stochastic variability

in human choices. That is, at many points in time, the model on its own must choose randomly (although typically with biases) among a set of alternative actions. Model tracing examines the behavioral data and identifies which of all the possible alternative paths the model could have taken best fits the observed behavioral data. The second factor shaping model tracing is that typically there are several internal steps for every external step. Thus, the model must be run for several steps, with each step potentially having alternative choices, to produce the full set of possible matches to the behavioral data. If there are many internal steps between external behaviors, then a large set of internal paths may need to be generated. The third factor is that the behavioral data may not uniquely distinguish among different model paths. In such circumstances, one must select the currently most probable path. With the addition of eye-tracking data, the density of observable data points goes up significantly, making it easier to match models to data.

VIII. ADAPTATION OF USER/SYSTEM INTERFACE

A system based on IHCI can adapt the interface based on current needs of the human participant as found in the embodied model of cognition. As stated earlier, different affective and cognitive diagnoses include confusion, fatigue, stress, momentary lapses of attention, and misunderstanding of procedures. Different adaptations include simplifying the interface, highlighting critical information, and tutoring on selected misunderstandings. The types of interface adaptations that one can consider include: 1) addition and deletion of task details; 2) addition and deletion of help/feedback windows; 3) changing the formatting/organization of information; and 4) addition and removal of automation of simple subtasks. These changes are described generically here.

With respect to the addition and deletion of task details, the important insight is that modern interfaces contain details relevant to many subtasks. When an operator becomes confused or distracted, it may well be because details relevant to subtask A interfere with the attention to details needed to accomplish subtask B. One general strategy is to identify the currently critical subtask with the goal of eliminating details relevant to other subtasks or enhancing details relevant to the critical subtask. Interface details relevant to other subtasks can be restored when the user appears able to handle them. The combination of the POG data (via eye tracking) and the affective response data (via facial expressions) provides important information regarding which aspects of the interface to change and in what manner to change them. For example, if an important aspect of the screen is not attended to and the individual appears fatigued, then that aspect should be highlighted. By contrast, if an aspect of the screen is attended to for an unusually long period of time and is coupled with a look of confusion, then a situation-relevant help window will be displayed. All of the possible interface structures that are possible will have advantages and disadvantages that are likely to vary with the cognitive and affective state of the user. Thus, different interface structures will be optimal at different points in time, and if a particular structure is gen-

erally suboptimal, there is no reason to ever use it. For example, having a help window display help messages may be useful for a confused individual, but may be distracting for a nonconfused individual. Alternatively, having less information on the screen may be helpful to a fatigued individual, but harmful to a fully attentive individual (who could make appropriate use of the extra information to handle more subtasks).

Because no one particular interface structure is better than others across all situations, one can avoid strange feedback loops in which a user becomes trained (either implicitly or explicitly) to always look frustrated because that makes the task easier. Instead, users will be trained to correctly externalize their internal states. For example, when frustrated, look frustrated, because that will produce a change that is useful for dealing with this particular source of frustration; but when not frustrated, do not look frustrated, because that will produce changes that reduce optimal performance (of someone who is not frustrated).

A model of embodied cognition that is continuously being updated to reflect the individual's cognitive, perceptual, motor, and affective states makes it possible to have two different methods of adapting the interface: reactive and proactive. In reactive adaptation, the system waits for external evidence of some cognitive or affective change before adapting the interface. For example, the user becomes confused and this confusion is manifested by a confused look, longer choice latencies, and longer fixations across a broader range of entities. A reactive system adapts the interface only after the confusion is manifested. Alternatively, a proactive system applies the model of embodied cognition (which is capable of performing the task) to the correct task state and predicts what kinds of problems the user is likely to encounter. These predictions are used to adapt the interface, that is, the interface changes before the user becomes confused, frustrated, or bored (or at least before this can be diagnosed from outward performance changes). Once the model tracing has approached a high level of accuracy, so that we believe it can use its broadened set of inputs, then one can explore including proactive interface adaptation in the system.

For either proactive or reactive adaptation, the adaptation will have to be conservative (i.e., relatively infrequent with relatively small changes at any one time). An interface that is constantly changing is a source of frustration in itself. Moreover, there should be a relatively small set of possible changes to the interface, and the set needs to be introduced during initial training. The embodied model provides insights into how conservative to be (i.e., to predict how disruptive various interface changes will be), in addition to providing insights into what interface adaptations are likely to be helpful.

IX. CONCLUSION

This paper has described a W5+ methodology for IHCI that extends current methods of interpreting human activities. Our approach to IHCI has four central pieces. First,

the behavioral interactions between the user and the interface are processed in a very rich fashion, including the novel use of mouse gestures, to provide a richer understanding of the user's cognitive and affective states. Second, additional nonverbal information is gathered through perceptual processing of eye gaze, pupil size, facial expressions, and arm movements to further enrich the understanding of the user's cognitive and affective states. Third, an embodied model of cognition is synchronized with the behavioral and perceptual data to produce a deep understanding of the user's state. Fourth, the computer interface is adapted in reaction to problems diagnosed in the user's cognitive and affective states in a task-sensitive way.

While our complete methodology has not yet been implemented in a running system, we have described the tools that would be required to implement such a system. Further, these tools appear to be well within current computational capabilities. We are currently engaging in research to further flesh out the computer science and cognitive psychology underlying these tools.

Some subtle components in our methodology require further comment. In particular, the process of building an embodied cognitive model has three separate advantages for intelligent adaptation of an interface. First, it forces one to develop highly detailed models of the precise cognitive and affective problems that a user might experience because the model must be designed to perform the task in the same way that the user does. This enforced detail allows for more precise diagnosis of sources of problems. Second, the embodied cognitive model allows one to test the consequences of different changes to the interface so that one can have a good understanding of which interface changes are likely to help and why they will help. Without the embodied cognitive model, one must rely on simple rules of thumb and past experiences to determine which changes will be effective. Third, the embodied cognitive model has a predictive component that allows one to predict what problems a user is likely to have in the future and warn the user about them in advance (e.g., when fatigue is likely to begin to occur given the recent load and current arousal level).

In developing a running system that implements our methodology, we recommend a strategy of using a simulated task environment. In field research, there is often too much complexity to allow for definite conclusions, and in laboratory research, there is usually too little complexity to allow for any interesting conclusions [13]. Those who study complex situations as well as those who wish to generalize their results to complex situations have often faced the dilemma so succinctly framed by Brehmer and Dörner. Simulated task environments are the solution to this dilemma. The term "simulated task environment" is meant to be both restrictive and inclusive. There are many types of simulations; however, the term is restricted to those that are intended as simulations of task environments. At the same time, the term includes the range of task simulations from high-fidelity ones that are intended as a substitute for the real thing, all the way to microworlds that enable the performance of existing tasks [33]. The common denominator

in these simulated task environments is the researcher's desire to study complex behavior. The task environment must be complex enough to challenge the current state of the art, but malleable enough so that task complexity and interface adaptivity can be controlled and increased as the research progresses. These requirements can be met by using simulated task environments. We are working on the IHCI approach in the context of human operators interacting with ARGUS, a simulated task environment for radar operator tasks [33]. These tasks represent a real-time safety-critical environment in which improving HCI is of utmost importance. Similar issues relating to image processing, cognitive modeling, and intelligent interface adaptation can be found in the HCI of a wide variety of domains (e.g., medical, educational, business). We can collect and time stamp every mouse click made by the subject, every system response, and every mouse movement with an accuracy of 17 ms and interleave this record with POG data collected 60 times per second. ACT-R/PM models currently interact with Argus. In addition, because we own the ARGUS code and it is written in Lisp, the simulated task environment is easy to modify.

Perception in general, and form and behavior analysis in particular, are important not only because they can describe things but, as Aristotle realized long ago, because they make us know and bring to light many differences between things so we can categorize them and properly respond to their affordances. Once both forms and behaviors are represented, their most important functionality is to serve for discrimination and classification. Recognition is thus based on both perceptual form and behavior, together with their associated functionalities. Form and behavior analysis considers things like average prototypes and/or the similarity between them, while functionality carves the perceptual and behavioral layout according to innate physical and geometrical constraints, sensor-motor affordances, and their corresponding cognitive mental states. According to this view, functional and purposive recognition takes precedence over perceptual and behavioral reconstruction. As things are always changing and constancy is an illusion, form and behavior recognition require generalization and motivate learning and adaptation. What form and behavior analysis do not require, however, is taking them to bits in ways that destroy the very relations that may be of the essence; as Lewontin [54] would say, "one murders to dissect." Embodied cognition, as described and advocated in this paper, provides the glue connecting the apparent visual form and behavior with hidden mental models, which bear on both functionality and performance. To further emphasize the important role functionality plays in perception, it is instructive to recall Oliver Sacks' well-known book, *The Man Who Mistook his Wife for a Hat*. The book describes someone who can see, but not interpret what he sees: shown a glove, the man calls it a "receptacle with five protuberances." The moral is that people see not only with the eyes, but with the brain as well. In other words, perception involves a whole and purposive cognitive process, and this is what this paper advocates in terms of technology and tools for IHCI.

REFERENCES

- [1] J. K. Aggarwal, Q. Cai, and B. Sabata, "Nonrigid motion analysis: Articulated and elastic motion," *Comput. Vis. Image Understanding*, vol. 70, pp. 142–156, 1998.
- [2] K. Akita, "Image sequence analysis of real world human motion," *Pattern Recognit.*, vol. 17, pp. 73–83, 1984.
- [3] M. W. Alibali and S. Goldin-Meadow, "Gesture-speech mismatch and mechanisms of learning: What the hands reveal about a child's state of mind," *Cogn. Psychol.*, vol. 25, pp. 468–523, 1993.
- [4] J. R. Anderson, C. F. Boyle, A. T. Corbett, and M. W. Lewis, "Cognitive modeling and intelligent tutoring," *Artif. Intell.*, vol. 42, pp. 7–49, 1990.
- [5] J. R. Anderson and C. Lebière, Eds., *Atomic Components of Thought*. Hillsdale, NJ: Erlbaum, 1998.
- [6] J. R. Anderson, M. Matessa, and C. Lebière, "ACT-R: A theory of higher-level cognition and its relation to visual attention," *Hum.-Comput. Interaction*, vol. 12, pp. 439–462, 1997.
- [7] N. I. Badler, C. B. Phillips, and B. L. Webber, *Simulating Humans*. New York: Oxford Univ. Press, 1993.
- [8] N. I. Badler and S. W. Smoliar, "Digital representations of human movement," *Comput. Surveys*, vol. 11, pp. 19–38, 1979.
- [9] R. Bajcsy and J. Kosecka, "The problem of signal and symbol integration: A study of cooperative mobile autonomous agent behavior," presented at the DAGM Symp., Bielefeld, Germany, 1995.
- [10] R. Barzel, *Physically-Based Modeling for Computer Graphics*. Boston, MA: Academic, 1992.
- [11] A. F. Bobick, "Movement, activity, and action: The role of knowledge in the perception of motion," *Philosophical Trans. Roy. Soc. London B*, vol. 352, pp. 1257–1265, 1997.
- [12] C. Bregler, "Learning and recognizing human dynamics in video sequences," in *Proc. Computer Vision and Pattern Recognition*, 1997, pp. 568–574.
- [13] B. Brehmer and D. Dörner, "Experiments with computer-simulated microworlds: Escaping both the narrow straits of the laboratory and the deep blue sea of the field study," *Comput. Hum. Behavior*, vol. 9, pp. 171–184, 1993.
- [14] C. Bregler and J. Malik, "Tracking people with twists and exponential maps," in *Proc. Computer Vision and Pattern Recognition*, 1998, pp. 8–15.
- [15] H. Buxton and R. Howarth, "Watching behavior: The role of context and learning," in *Proc. Int. Conf. Image Processing*, vol. 2, 1996, pp. 797–800.
- [16] M. D. Byrne and J. R. Anderson, "Perception and action," in *Atomic Components of Thought*, J. R. Anderson and C. Lebière, Eds. Hillsdale, NJ: Erlbaum, 1998, pp. 167–200.
- [17] C. Cedras and M. Shah, "Motion-based recognition: A survey," *Image Vis. Comput.*, vol. 13, pp. 129–155, 1995.
- [18] T.-J. Cham and J. M. Rehg, "A multiple hypothesis approach to figure tracking," in *Proc. Computer Vision and Pattern Recognition*, vol. 2, 1999, pp. 239–245.
- [19] E. Cohen, L. Namir, and I. M. Schlesinger, *A New Dictionary of Sign Language*. The Hague, The Netherlands: Mouton, 1977.
- [20] J. W. Davis and A. F. Bobick, "The representation and recognition of human movement using temporal templates," in *Proc. Computer Vision and Pattern Recognition*, 1997, pp. 928–934.
- [21] L. S. Davis, D. Harwood, and I. Haritaoglu, "Ghost: A human body part labeling system using silhouettes," in *Proc. ARPA Image Understanding Workshop*, 1998, pp. 229–235.
- [22] J. Deutscher, A. Blake, and I. Reid, "Articulated body motion capture by annealed particle filtering," in *Proc. Computer Vision and Pattern Recognition*, vol. 2, 2000, pp. 126–133.
- [23] Z. Duric, F. Li, and H. Wechsler, "Recognition of arm movements," in *Proc. Int. Conf. Automatic Face and Gesture Recognition*, Washington, DC, 2002, pp. 348–353.
- [24] N. Eshkol and A. Wachmann, *Movement Notation*, London, U.K.: Weidenfeld and Nicholson, 1958.
- [25] I. A. Essa and A. P. Pentland, "Facial expression recognition using a dynamic model and motion energy," in *Proc. Int. Conf. Computer Vision*, 1995, pp. 360–367.
- [26] —, "Coding, analysis, interpretation, and recognition of facial expressions," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 19, pp. 757–763, 1997.
- [27] P. F. Felzenszwalb and D. P. Huttenlocher, "Efficient matching of pictorial structures," in *Proc. Computer Vision and Pattern Recognition*, vol. 2, 2000, pp. 66–73.
- [28] D. M. Gavrila, "The visual analysis of human movement: A survey," *Comput. Vis. Image Understanding*, vol. 73, pp. 82–98, 1999.
- [29] D. M. Gavrila and L. S. Davis, "3D model-based tracking of humans in action: A multi-view approach," in *Proc. Computer Vision and Pattern Recognition*, 1996, pp. 73–80.
- [30] D. M. Gavrila and V. Philomin, "Real-time object detection for smart vehicles," in *Proc. Int. Conf. Computer Vision*, 1999, pp. 87–93.
- [31] S. Goldin-Meadow, M. W. Alibali, and R. B. Church, "Transitions in concept acquisition: Using the hand to read the mind," *Psychol. Rev.*, vol. 100, pp. 279–297, 1993.
- [32] K. Gould and M. Shah, "The trajectory primal sketch: A multi-scale scheme for representing motion characteristics," in *Proc. Computer Vision and Pattern Recognition*, 1989, pp. 79–85.
- [33] W. D. Gray, "Simulated task environments: The role of high-fidelity simulations, scaled worlds, synthetic environments, and microworlds in basic and applied cognitive research," *Cogn. Sci.*, vol. 2, no. 2, pp. 205–227, 2002.
- [34] W. D. Gray, P. Palanque, and F. Paternò, "Introduction to the special issue on: Interface issues and designs for safety-critical interactive systems," *ACM Trans. Computer-Human Interaction*, vol. 6, no. 4, pp. 309–310, 1999.
- [35] A. H. Guest, *Choreo-Graphics: A Comparison of Dance Notation Systems from the Fifteenth Century to the Present*. New York: Gordon and Breach, 1989.
- [36] S. Gutta, I. F. Imam, and H. Wechsler, "Hand gesture recognition using ensembles of radial basis function (RBF) networks and decision trees," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 11, pp. 845–872, 1997.
- [37] I. Haritaoglu, D. Harwood, and L. S. Davis, "W4S: A real-time system for detecting and tracking people," in *Proc. Computer Vision and Pattern Recognition*, 1998, pp. 962–968.
- [38] J. K. Hodgins and N. S. Pollard, "Adapting simulated behaviors for new characters," in *Proc. SIGGRAPH*, 1997, pp. 153–162.
- [39] D. Hogg, "Model based vision: A program to see a walking person," *Image Vis. Comput.*, vol. 1, pp. 5–20, 1983.
- [40] J. Huang, S. Gutta, and H. Wechsler, "Detection of human faces using decision trees," in *Proc. Int. Conf. Automatic Face and Gesture Recognition*, Killington, VT, 1996, pp. 248–252.
- [41] S. Ioffe and D. A. Forsyth, "Finding people by sampling," in *Proc. Int. Conf. Computer Vision*, 1999, pp. 1092–1097.
- [42] S. Jabri, Z. Duric, A. Rosenfeld, and H. Wechsler, "Detection and location of people in video images using adaptive fusion of color and edge information," in *Proc. Int. Conf. Pattern Recognition*, 2000.
- [43] B. Jähne, *Digital Image Processing*. Berlin, Germany: Springer-Verlag, 1997.
- [44] G. Johansson, "Visual perception of biological motion and a model for its analysis," *Perception and Psychophysics*, vol. 14, pp. 210–211, 1973.
- [45] —, "Spatio-temporal differentiation and integration in visual motion perception," *Psychol. Res.*, vol. 38, pp. 379–393, 1976.
- [46] S. X. Ju, M. J. Black, and Y. Yacoub, "Cardboard people: A parameterized model of articulated image motion," in *Proc. Int. Conf. Automatic Face and Gesture Recognition*, Nara, Japan, 1995, pp. 38–44.
- [47] S. K. Jung, "Motion Analysis of Articulated Object for Optical Motion Capture," Ph.D. dissertation, KAIST, Taejon, Korea, 1996. Tech. Memo 97–8.
- [48] I. Kakadiaris and D. Metaxas, "3D human body model acquisition from multiple views," in *Proc. Int. Conf. Computer Vision*, 1995, pp. 618–623.
- [49] D. E. Kieras and D. E. Meyer, "An overview of the EPIC architecture for cognition and performance with application to human-computer interaction," *Hum.-Comput. Interaction*, vol. 12, pp. 391–438, 1997.
- [50] M. L. Knapp and J. A. Hall, *Nonverbal Communication in Human Interaction*. New York: Harcourt Brace Jovanovich, 1997.
- [51] H.-J. Lee and Z. Chen, "Determination of 3D human body posture from a single view," *Comput. Vis. Graph. Image Process.*, vol. 30, pp. 148–168, 1985.
- [52] M. K. Leung and Y.-H. Yang, "Human body segmentation in a complex scene," *Pattern Recognit.*, vol. 20, pp. 55–64, 1987.
- [53] —, "First sight: A human body outline labeling system," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 17, pp. 359–377, 1995.
- [54] R. C. Lewontin, "The science of metamorphosis," *The New York Reviews of Books XXXVI (7)*, 1989.

- [55] G. L. Lohse and E. J. Johnson, "A comparison of two process tracing methods for choice tasks," *Org. Behavior Human Decision Processes*, vol. 68, pp. 28–43, 1996.
- [56] M. C. Lovett and C. D. Schunn, "Task representations, strategy variability and base-rate neglect," *J. Experimental Psychol.: General*, vol. 128, pp. 107–130, 1999.
- [57] S. J. McKenna, S. Jabri, Z. Duric, A. Rosenfeld, and H. Wechsler, "Tracking groups of people," *Comput. Vis. Image Understanding*, vol. 80, pp. 42–56, 2000.
- [58] S. J. McKenna, S. Jabri, Z. Duric, and H. Wechsler, "Tracking interacting people," in *Proc. Int. Conf. Automatic Face and Gesture Recognition*, Grenoble, France, 2000, pp. 348–353.
- [59] T. B. Moeslund and E. Granum, "A survey of computer vision-based human motion capture," *Comput. Vis. Image Understanding*, vol. 81, pp. 231–268, 2001.
- [60] S. A. Niyogi and E. H. Adelson, "Analyzing and recognizing walking figures in XYT," in *Proc. Computer Vision and Pattern Recognition*, 1994, pp. 469–474.
- [61] N. M. Oliver, B. Rosario, and A. P. Pentland, "A Bayesian computer vision system for modeling human interactions," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 22, pp. 831–843, 2000.
- [62] M. Oren, C. Papageorgiou, P. Sinha, E. Osuna, and T. Poggio, "Pedestrian detection using wavelet templates," in *Proc. Computer Vision and Pattern Recognition*, 1997, pp. 193–199.
- [63] V. I. Pavlovic, R. Sharma, and T. S. Huang, "Visual interpretation of hand gestures for human-computer interaction: A review," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 19, pp. 677–695, 1997.
- [64] J. W. Payne, J. R. Bettman, and E. J. Johnson, *The Adaptive Decision Maker*. New York: Cambridge Univ. Press, 1993.
- [65] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco, CA: Morgan Kaufmann, 1988.
- [66] V. Philomin, R. Duraiswami, and L. S. Davis, "Quasi-random sampling for condensation," in *Proc. Eur. Conf. Computer Vision*, vol. 2, 2000, pp. 134–149.
- [67] R. Picard, *Affective Computing*. Cambridge, MA: MIT Press, 1998.
- [68] R. Polana and R. Nelson, "Detection and recognition of periodic, nonrigid motion," *Int. J. Comput. Vis.*, vol. 23, pp. 261–282, 1997.
- [69] Annotated Computer Vision Bibliography, K. E. Price. [Online]. Available: <http://iris.usc.edu/Vision-Notes/bibliography/linebreakcontents.html>
- [70] J. R. Quinlan, *C4.5: Programs for Machine Learning*. San Francisco, CA: Morgan Kaufmann, 1993.
- [71] R. F. Rashid, "Toward a system for the interpretation of moving light displays," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-2, pp. 574–581, 1980.
- [72] J. M. Rehg and T. Kanade, "Model-based tracking of self occluding articulated objects," in *Proc. Int. Conf. Computer Vision*, 1995, pp. 612–617.
- [73] J. Rittscher and A. Blake, "Classification of human body motion," in *Proc. Int. Conf. Computer Vision*, 1999, pp. 634–639.
- [74] K. Rohr, "Toward model-based recognition of human movements in image sequences," *CVGIP: Image Understanding*, vol. 59, pp. 94–115, 1994.
- [75] D. D. Salvucci and J. R. Anderson, "Automated eye-movement protocol analysis," *Human-Computer Interaction*, vol. 16, pp. 39–86, 2001.
- [76] C. D. Schunn and L. M. Reder, "Another source of individual differences: Strategy adaptivity to changing rates of success," *J. Experimental Psychol.: General*, vol. 130, no. 1, pp. 59–76, 2001.
- [77] C. D. Schunn, S. Trickett, and J. G. Trafton, *What Gestures Reveal About the Scientist's Mind: Data Analyzes of Data Analysis*, ser. Krasnow Inst. Brown Bag Series. George Mason Univ., 1999.
- [78] B. Shneiderman and P. Maes, "Direct manipulation vs. interface agents," *Interactions*, vol. 4, pp. 643–661, 1997.
- [79] H. Sidenbladh, M. J. Black, and D. J. Fleet, "Stochastic tracking of 3D human figures using 2D image motion," in *Proc. Eur. Conf. Computer Vision*, 2000.
- [80] S. Sirohey, A. Rosenfeld, and Z. Duric, "A method of detecting and tracking irises and eyelids in video," *Pattern Recognit.*, vol. 35, pp. 1389–1401, 2002.
- [81] Y. Song, L. Goncalves, E. di Bernardo, and P. Perona, "Monocular perception of biological motion: Detection and labeling," in *Proc. Int. Conf. Computer Vision*, 1999, pp. 805–813.
- [82] T. Starner and A. Pentland, "Visual recognition of American Sign Language using hidden Markov models," in *Proc. Int. Workshop Automatic Face and Gesture Recognition*, Zurich, Switzerland, 1995, pp. 189–194.
- [83] S. Wachter and H. H. Nagel, "Tracking persons in monocular image sequences," *Comput. Vis. Image Understanding*, vol. 74, pp. 174–192, 1999.
- [84] H. Wechsler, *Computational Vision*. Orlando, FL: Academic, 1990.
- [85] C. R. Wren, A. Azarbayejani, T. Darrell, and A. P. Pentland, "Pfinder: Real-time tracking of the human body," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 19, pp. 780–785, 1997.
- [86] Y. Yacoob and L. S. Davis, "Recognizing human facial expressions from long image sequences using optical flow," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 18, pp. 636–642, 1996.
- [87] M.-H. Yang and N. Ahuja, "Recognizing hand gestures using motion trajectories," in *Proc. Computer Vision and Pattern Recognition*, vol. 1, 1999, pp. 466–472.
- [88] V. M. Zatsiorsky, *Kinematics of Human Motion*. Champaign, IL: Human Kinetics, 1997.
- [89] *Proc. Int. Workshop Automatic Face and Gesture Recognition*, Zurich, Switzerland, 1995.
- [90] *Proc. Int. Conf. Automatic Face and Gesture Recognition* Killington, VT, 1996.
- [91] *Proc. Int. Conf. Automatic Face and Gesture Recognition*, Nara, Japan, 1998.
- [92] *Proc. Int. Conf. Automatic Face and Gesture Recognition*, Grenoble, France, 2000.



Zoran Duric received the M.S. degree in electrical engineering from the University of Sarajevo, Bosnia and Herzegovina, in 1986 and the Ph.D. degree in computer science from the University of Maryland at College Park, in 1995.

He joined the faculty of George Mason University, Fairfax, VA, in the fall of 1997 as an Assistant Professor of Computer Science. He has published technical papers on various aspects of motion understanding in computer vision journals and conferences. Major focus of

his research has been the understanding of physical constraints on motions of vehicles and humans and using those constraints in building robust and efficient vision systems.



Wayne D. Gray received the Ph.D. degree from the University of California at Berkeley in 1979.

He is a prominent researcher in the fields of human-computer interaction (HCI), cognitive task analysis, computational models of embodied cognition, and human error. His first position was with the U.S. Army Research Institute where he worked on tactical team training (at the Monterey Field Unit) and later on the application of artificial intelligence (AI) technology to training for air-defense systems (HAWK) (at ARI-HQ

Alexandria, VA). He spent a post-doctoral year with Prof. J. R. Anderson's lab at Carnegie Mellon University, Pittsburgh, PA, before joining the AI Laboratory of NYNEX Science & Technology Division. At NYNEX, he applied cognitive task analysis and cognitive modeling to the design and evaluation of interfaces for large, commercial telecommunications systems. Since joining academe, he has received grants from government, industry, and private foundations. He is on the review board for the *Cognitive Science* journal and an Associate Editor for the journal *Human Factors* as well as for *ACM Transactions on Computer-Human Interaction*. He chaired the Fourth International Conference on Cognitive Modeling (ICCM-2001) and is co-chair of the 24 Annual Conference of the Cognitive Science Society. At Rensselaer Polytechnic Institute, Troy, NY, he is a Professor of Cognitive Science and Program Director of the Cognitive Science Ph.D. Program.



Ric Heishman (Student Member, IEEE) received the B.S. degree in computer engineering from the University of Cincinnati, Cincinnati, OH, in 1991 and the M.S. degree in information systems from American University, Washington, DC, in 1997. He is currently working toward the Ph.D. degree in information technology at George Mason University, Fairfax, VA.

He spent 20 years in various positions in the defense industry—with the U.S. Navy, IBM, Loral, and Lockheed Martin. He is presently an Associate Professor of Information Technology at Northern Virginia Community College, Manassas, and serves as Program Head for Computer Science and Microelectronics at the Manassas campus. His research is in the areas of human computer interaction and computer vision, focusing on monitoring affective and cognitive states in humans using eye-region biometrics.

Mr. Heishman is a member of the ACM and the IEEE Computer Society.



Fayin Li received the B.S. degree in electrical engineering from Huazhong University of Science and Technology, China, in 1996 and the M.S. degree in computer science from the Institute of Automation, Chinese Academy of Sciences, China, in 1999. He is currently working toward the Ph.D. degree at George Mason University, Fairfax, VA.

His research interests include automatic hand gesture recognition, video tracking and surveillance, image processing, pattern recognition, and model selection.



Azriel Rosenfeld received the Ph.D. degree in mathematics from Columbia University, New York, NY, in 1957.

He is a Distinguished University Professor emeritus and the founding Director of the Center for Automation Research at the University of Maryland at College Park. He also held Affiliate Professorships in the Departments of Computer Science, Electrical Engineering, and Psychology. He is a widely known researcher in the field of computer image analysis. He has published over

30 books and over 600 book chapters and journal articles, is an Associate Editor of over 25 journals, and has directed nearly 60 Ph.D. dissertations.

Dr. Rosenfeld is a Fellow of several professional societies and has won numerous professional society awards, and he has received several honorary doctoral degrees.



Michael J. Schoelles received the Ph.D. degree from George Mason University, Fairfax, VA, in 2002.

He is currently a Research Scientist in the Applied Research in Cognition and Human Factors Laboratory, George Mason University. His research focuses on computational modeling of human computer interactions. The modeling methodology combines an embodied unified cognitive theory with task knowledge and computer interface design constraints to produce

interactive behavior.



Christian Schunn received the Ph.D. degree in psychology from Carnegie–Mellon University, Pittsburgh, PA, in 1995.

He is currently a Research Scientist at the Learning Research and Development Center and an Assistant Professor in the Psychology Department and Intelligent Systems Program, University of Pittsburgh, Pittsburgh, PA. His research uses a combination of psychological experiments, computational cognitive modeling, and field observations to study the cognition

underlying adaptive behavior and scientific thinking, both as basic psychology research questions and also how they apply to improving education and improving human-computer interfaces. This research includes studying collaboration at a distance, expertise in science, spatial thinking in complex domains like the submarine sonar domain, and diagnosing and repairing problems in low-level leadership strategies through intelligent tutoring. He has coauthored over 40 scientific papers and chapters. He co-edited the book *Designing for Science: Implications from Professional, Instructional, and Everyday Science* (Hillsdale, NJ: Erlbaum, 2001). He has also co-chaired a number of international conferences, including Designing for Science in 1998, the 6th Annual ACT-R Workshop in 1999, the 4th International Conference on Cognitive Modeling in 2001, and the 24th Annual Meeting of the Cognitive Science Society in 2002.



Harry Wechsler (Fellow, IEEE) received the Ph.D. degree in computer science from the University of California at Irvine in 1975.

He is presently Professor of Computer Science and Director for the Center for Distributed and Intelligent Computation, George Mason University, Fairfax, VA. His research, in the field of intelligent systems, has been in the areas of perception (computer vision, automatic target recognition, signal and image processing), machine intelligence (pattern recognition, neural networks,

and data mining), evolutionary computation (genetic algorithms and animats), and human–computer intelligent interaction (face and hand gesture recognition, biometrics, video tracking and surveillance, and interpretation of human activity). He was Director for the NATO Advanced Study Institutes (ASI) on “Active Perception and RobotVision” (Maratea, Italy, 1989), “From Statistics to Neural Networks” (Les Arcs, France, 1993), and “Face Recognition: From Theory to Applications” (Stirling, UK, 1997), and he has served as co-chair for the International Conference on Pattern Recognition held in Vienna, Austria, in 1996. He has authored over 200 scientific papers and one book *Computational Vision* (New York: Academic, 1990), and he was the editor for *Neural Networks for Perception* (New York: Academic, 1991, vols. 1 and 2).

Dr. Wechsler is a Fellow of the International Association of Pattern Recognition.