

# Augmented Realities Integrating User and Physical Models

Thad Starner, Bernt Schiele, Bradley J. Rhodes, Tony Jebara,  
Nuria Oliver, Joshua Weaver and Alex Pentland

**Abstract.** *The obvious advantage of wearable computing is mobility; it also offers the user a certain intimacy with augmented realities. A model of the user is as important as a model of the physical world for creating a seamless, unobtrusive interface while avoiding “information overload.” This paper summarizes some of the current projects at the MIT Media Laboratory that explore user and physical environment modeling.*

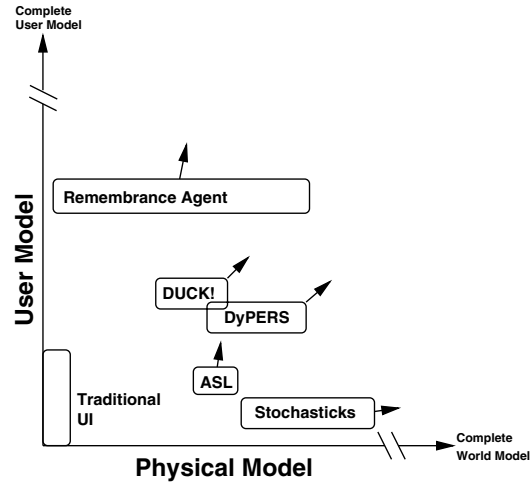
## 1 Introduction

Wearable computers have the potential to *see* as the user sees, *hear* as the user hears, and experience the life of the user from his or her perspective. They are also physically and mentally more intimate with the user than are other kinds of computers and they may be used for extended periods of time, allowing a unique opportunity to model users' usage patterns and habits. This increase in user and environmental information can lead to more intelligent and fluid interfaces that incorporate the physical world.

This paper summarizes several of the current wearable computing and augmented reality research projects at the MIT Media Laboratory that explore the dimensions of user and physical modeling. (See Figure 1.) For more complete information on a particular project, the reader is encouraged to refer to the original papers.

## 2 Remembrance Agent

The Remembrance Agent (RA) is a program that continuously watches over the shoulder of the user of a wearable computer; it displays one-line summaries of notes files, old e-mail, papers, and other text information that might be relevant to the user's current context [Rhodes 97]. These



**Figure 1.** Wearable computing projects positioned along conceptual axes of the degree of involvement of a user model versus the extent to which the project models the physical environment of the user. The arrows indicate the direction of current research.

summaries are listed on a heads-up display, so that the wearer can view the information with a quick glance. The full text can be retrieved using a one-handed keyboard.

The wearable version of the RA uses physical sensors to model the user's environment. The RA continuously monitors these sensors, as well as the notes being entered by the user. It uses this information to suggest documents, from a set of pre-indexed text, that are most relevant to the user's current situation. For example, a user's context might be described by a combination of the current time of day and day of the week (provided by the wearable system's clock), location (provided by an infrared beacon in the room), the person to whom the user is speaking (provided by an active badge), and the subject of the conversation (as indicated by the notes being taken). The suggestions provided on the RA are based on a combination of all of these elements, using text-retrieval techniques similar to those used in Web search engines.

## 2.1 Augmented Reality Remembrance Agent

The wearable RA uses an overlay display, but does not register its annotations with specific objects or locations in the real world—as one might expect from a full AR system. In many cases, such a real-world fixed dis-

play wouldn't even make sense, since suggestions are often conceptually relevant to the current situation without being relevant to a specific object or location. In order to examine augmented reality interfaces, a different version of the RA was implemented using a desktop computer, HUD, and head-mounted camera. The wearer of the system viewed the world through the camera, and the camera output also went to an SGI reality engine. Around the room were colored tags, which a vision system could detect. The size and shape of the tags were used to determine their distance and orientation (see Figure 2). Color-coding was used to identify each tag.



**Figure 2.** Multiple graphical overlays aligned through visual tag tracking. The color code of each tag provides a unique ID for the object. In addition, the vision process tracks each tag in 2.5D as the head-mounted camera moves.

Whenever a tag was detected, the code number was looked up in a table, and a pre-computed message was overlaid on top of the object being viewed. The system could detect the distance of the tag, allowing more information to be displayed as the user came closer to a tagged object. Thus, the act of approaching a tag was equivalent to clicking on a hypertext link.

On top of this system, the Remembrance Agent showed the top suggestions for an individual tag. This created a two-level information system, with some information being provided by the infrastructure (tied to the tag via a lookup table), and some information provided from the user's own files via the Remembrance Agent.

## 2.2 Dynamic Personal Enhanced Reality Agent

A recent extension of the system described above uses a generic object recognizer to identify objects instead of tags. The system, called *Dynamic Personal Enhanced Reality System* (DyPERS [Schiele 99]), retrieves audio and video clips based on associations with real objects. The generic object



**Figure 3.** A DyPERS user listening to a guide during a test art gallery tour

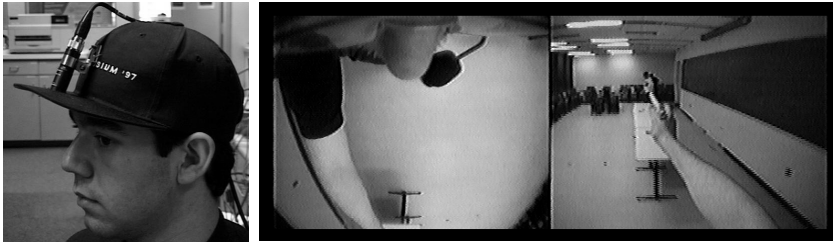
recognizer is based on a probabilistic recognition system [Schiele 96], which is capable of discriminating more than 100 objects in the presence of major occlusions, scalings, and rotations. While 100 objects is not enough to be practical in an unconstrained environment, the number of possible objects can be significantly reduced using the location of the user, time of day, and other available information. (See Figure 3.)

### 3 User-observing Wearable Cameras

In the previous section, head-mounted camera systems face forward and observe the same region as the user's own eyes. By changing the angle of the camera so that it points downward, we can use it to track the user's own body. This allows the user's hands, feet, torso, and lips to be observed without the gloves or body suits associated with virtual reality gear.

Two projects currently use this camera orientation. The first attempts to translate American Sign Language into English by tracking the user's hands via the downward-looking camera. The wearable ASL recognition system outperforms equivalent desk-based camera systems and most dataglove-based systems in recognition accuracy. For five-word sentences composed from a 40-word lexicon, the system achieves 96.8% word accuracy with an unrestricted grammar (any word is possible, any number of times, in any order) [Starner 98b].

The second system, DUCK!, begins to demonstrate how such methods may be useful for augmented realities. Using only two video views, DUCK! tracks the wearer's current location and task. DUCK!'s domain is limited to the real-space game Patrol, which is played by MIT students every weekend in a campus building. Participants are divided into teams and aggressively hunt each other with rubber suction-cup dart guns through 14 rooms or areas. DUCK! monitors the average color and luminance values of the floor,



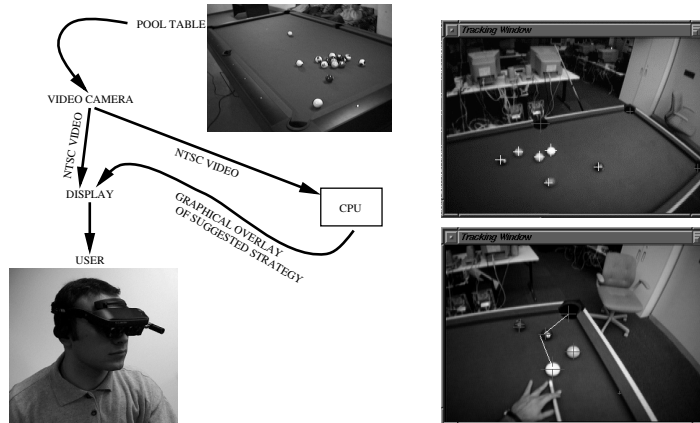
**Figure 4.** Left: The Patrol cap with two cameras. The larger, visible camera is mounted facing downward. The second camera faces forward and is hidden by the brim. Right: Images from the Patrol cap. The left and right images are from the downward-looking and forward-looking cameras, respectively.

the scene in front of the player, and the player's nose as a lighting calibration image. These images are pulled from the two video streams at six frames per second in order to determine location. In 24.5 minutes of training video and 19.3 minutes of test video, the system performed with 82% accuracy [Starner 98a]. (See Figure 4.) DUCK! also attempts to discriminate between the player's tasks by identifying hand gestures representing aiming/shooting, reloading, and everything else, through a combination of a generic object recognition system [Schiele 96] and the HMMs used in the ASL task above. Preliminary results show 86% accuracy in distinguishing these classes. Other user actions, such as standing, walking, running, and scanning the environment, can be considered as tasks that can run concurrently with other actions. Future work will address these tasks.

While still in preliminary stages, the systems described above suggest how context perception may be used in augmented reality interfaces. Through heads-up displays, the players can keep track of the team's positions, even to the extent of "seeing through walls" to what may be occurring several rooms away. If aim and reload gestures are recognized by a particular player's system, his position can be highlighted in the rest of the team's displays, indicating that he needs aid. Furthermore, when the computer recognizes a player to be in battle, the computer should inhibit his interface in order to avoid interruptions.

## 4 Stochasticks

Stochasticks is a practical application of wearable computing and augmented reality which enhances the game of billiards [Jebara 97]. (See Figure 5.) Probabilistic color models and symmetry operations are used to localize the table, pockets, and balls through a video camera near the user's



**Figure 5.** Left: The system components. Right Top: Finding the balls. Right Bottom: Suggested shot.

eyes. The system classifies the objects of interest and ranks each possible shot in order to determine its relative usefulness. The system allows the user to proceed through a regular pool game while it automatically determines strategic shots. The resulting trajectories are rendered as graphical overlays on a head-mounted live video display. The wearable video output and the computer vision system provide an integration of real and virtual environments which enhances the experience of playing and learning the game of billiards, without encumbering the player.

#### 4.1 The System

A wearable computer is the hardware platform for the system and it includes a head-mounted display, a head-mounted video camera, and a central processing unit. The head-mounted camera is a miniature ELMO 2.2-mm video camera mounted on the heads-up display and aligned with the orientation of the eyes. Thus, the user's head direction will automatically direct the camera to areas of interest in the scene. The head-mounted display consists of a Virtual I/O 3D display (or a Sony Glasstron), which allows the CPU to project semi-transparent imagery into each eye via two separate CRTs at about 10 Hz.

Once the ball position is known, the easiest possible shot for a given player is computed, and the shot trajectory is projected onto the user's eye. At this point, we are undertaking a performance analysis of the overall system. The reliability of the algorithm is being investigated as well as its accuracy for both 2D and 3D overlays.

## 5 Conclusion

We have demonstrated how computer vision can be incorporated into augmented realities. We have discussed self-observing wearable camera systems that identify the user's gestures and location in a variety of conditions. Finally, through the projects presented, we have shown how modeling of both the user and the physical world play an important role in augmented realities.

## References

- [Jebara 97] T. Jebara, C. Eyster, J. Weaver, T. Starner, and A. Pentland. "Stochasticks: Augmenting the Billards Experience with Probabilistic Vision and Wearable Computers," *International Symposium on Wearable Computers*, IEEE Press, 1997.
- [Rhodes 97] B. Rhodes. "The Wearable Remembrance Agent: A System for Augmenting Memory," *Personal Technologies*, Vol. 1, No. 1, Springer, 1997.
- [Schiele 96] B. Schiele and J Crowley. "Probabilistic Object Recognition Using Multidimensional Receptive Field Histograms," *International Conference on Pattern Recognition*, Vol. B, IEEE Press, pp. 50-54, August 1996.
- [Schiele 99] B. Schiele, N. Oliver, T. Jebara, and A. Pentland. "An Interactive Computer Vision System, DyPERS: Dynamic and Personal Enhanced Reality System," *International Conference on Computer Vision Systems*, Springer, 1999.
- [Starner 98a] T. Starner, B. Schiele, and A. Pentland. "Visual Conextual Awareness in Wearable Computing," *Second International Symposium on Wearable Computers*, IEEE Press, 1998.
- [Starner 98b] T. Starner, J. Weaver, and A. Pentland. "Real-time American Sign Language Recognition Using Desk and Wearable Computer-based Video," *IEEE Trans. Patt. Analy. and Mach. Intell.*, To appear 1998.