

# Integrated person tracking using stereo, color, and pattern detection.

T. Darrell, G. Gordon, M. Harville, J. Woodfill  
Interval Research Corp.  
1801C Page Mill Road  
Palo Alto CA 94304

trevor,gaile,harville,woodfill@interval.com

<http://www.interval.com/papers/1998-021>

## Abstract

*We present an approach to real-time person tracking in crowded and/or unknown environments using multi-modal integration. We combine stereo, color, and face detection modules into a single robust system, and show an initial application in an interactive, face-responsive display. Dense, real-time stereo processing is used to isolate users from other objects and people in the background. Skin-hue classification identifies and tracks likely body parts within the silhouette of a user. Face pattern detection discriminates and localizes the face within the identified body parts. Faces and bodies of users are tracked over several temporal scales: short-term (user stays within the field of view), medium-term (user exits/reenters within minutes), and long term (user returns after hours or days). Short-term tracking is performed using simple region position and size correspondences, while medium and long-term tracking are based on statistics of user appearance. We discuss the failure modes of each individual module, describe our integration method, and report results with the complete system in trials with thousands of users.*

## 1 Introduction

The creation of displays or environments which passively observe and react to people is an exciting challenge for computer vision [4, 6]. Faces and bodies are central to human communication and yet machines have been largely blind to their presence in real-time, unconstrained environments.

Often, computer vision systems for person tracking exploit a single visual processing technique to locate and track user features. These systems can be non-robust to real-world conditions with multiple people and/or moving backgrounds. Additionally, tracking is usually performed only

over a single, short time scale: a person model is typically based only on an unbroken sequence of user observations, and is reset when the user is occluded or leaves the scene temporarily.

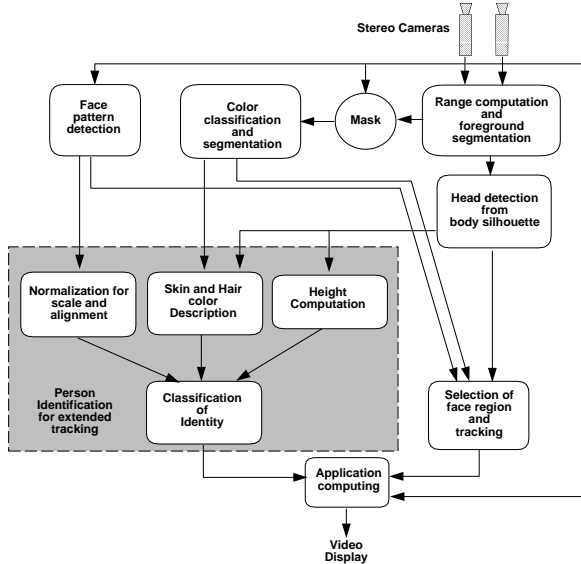
We have created a visual person tracking system which achieves robust performance through the integration of multiple visual processing modalities and by tracking over multiple temporal scales. With each modality alone it is possible to track a user under optimal conditions, but each also has, in our experience, substantial failure modes in unconstrained environments. Fortunately these failure modes are often independent, and by combining modules in simple ways we can build a system with overall robust performance.

In the following sections we describe our tracking framework and the three vision processing modalities used. We then describe an initial application of our system: a face-responsive, interactive video display. Finally we show the results of our system when deployed with naive users, and analyze both the qualitative success of the application and the quantitative performance of our tracking algorithms.

## 2 Tracking framework

A person tracking system for interactive environments has several desired criteria: it should operate in real-time, be robust to multiple users and changing background, provide a relatively rich visual description of the users, and be able to track people when they are occluded or momentarily leave the scene. We achieve these goals through the use of multi-modal integration and multi-scale temporal tracking.

We base our system on three primary visual processing modules: depth estimation, color segmentation, and intensity pattern classification (see Figure 1). As described in more detail below, depth information is estimated using a dense real-time stereo technique and allows easy segmentation of the user from other people and background objects. An intensity-invariant color classifier detects regions



**Figure 1. System overview showing the relationship of each modality with detection and short-term tracking, and with long-term tracking/identification.**

of flesh tone on the user and is used to identify likely body part regions such as face and hands. A face detection module is used to discriminate head regions from hands and other tracked body parts.

Figure 2 shows the output of the three vision processing modules. As a person tracker, each is individually fragile: notebooks are indistinguishable from faces in range silhouette, flesh color signs or clothes fool color-only trackers, and face pattern detectors typically are slower and only work with relatively canonical poses and expressions. However, when integrated together these modules can yield robust, fast tracking performance.

Tracking is performed in our system on three different time-scales: short-range (frame to frame while the person is visible), medium-range (when the person is momentarily occluded or leaves the field of view for a few minutes), and long range (when the person is absent for hours, days or more.) Long-term tracking can be thought of as a person identification task, where the database is formed from the set of previous users. For short-term tracking we simply compute region correspondences specific to each processing modality based on region position and size. Multi-modal integration is performed using the history of short-term tracked regions from each modality, yielding a representation of the user’s body shape and face location.

For medium and long-range tracking, we rely on a statistical model of multi-modal appearance to resolve correspondences between tracked users. In addition to body

shape and face location, and color of hair, skin, and clothes is recorded at each time step. We record the average value and covariance of represented features, and use them for matching. For medium-term tracking, lighting constancy and stable clothing color are assumed; for long-term tracking we adjust for changing lighting and do not include clothing in the match criteria.

In the next section, we discuss module specific processing, including classification, segmentation/grouping, and short-term tracking. Following that, we present our integration scheme, and correspondence method for medium and long-term tracking.

### 3 Mode-specific processing

Pixel-wise classification, grouping and short-term tracking are performed independently in each modality. Stereo processing outputs a user’s silhouette defined by range regions, color processing yields a set of skin color regions within range silhouette boundaries, and face processing returns a list of detected frontal face patterns; we describe each module in turn. Each mode also provides an independent estimate of head location and performs short-term tracking.

#### 3.1 User silhouette from dense stereo

To compute a set of user silhouettes, we rely on a dense real-time stereo system. Video from a pair of cameras is used to estimate dense range using a technique based on the census transform [8]; we have implemented the census algorithm on a single PCI card, multi-FPGA reconfigurable computing engine [9]. This stereo system is capable of computing 24 stereo disparities on 320 by 240 images at 42 frames per second, or approximately 77 million pixel-disparities per second. These processing speeds compare favorably with other real-time stereo implementations such as [3].

Our segmentation and grouping technique proceeds in several stages of processing, as illustrated in Figure 3. We first smooth the raw range signal to reduce the effect of low confidence stereo disparities using a morphological closing operator. We then compute the response of a gradient operator on the smoothed range data and threshold at a critical value based on the observed noise level in our disparity data. Connected components analysis is applied to these regions of smoothly varying range. We return all connected components whose area exceeds a minimum threshold.

The range processing module provides these user silhouettes, as well as estimates of head location. A candidate head is placed below the maxima of the range profile. Head position is refined in the integration stage, as described below.

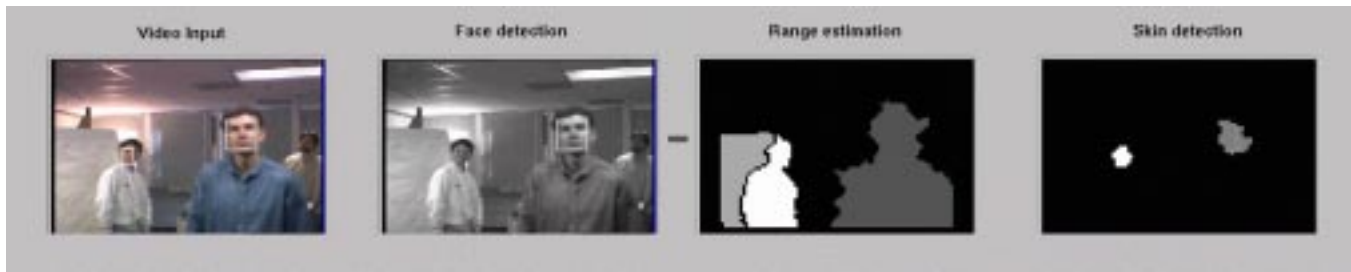


Figure 2. Output of vision processing modules: input image, face pattern detection output, connected components recovered from stereo range data, and flesh hue regions from skin hue classification. Boxes have been drawn on the faces of the two tracked users in the input image; the rightmost person in the image is beyond the workspace of the system.

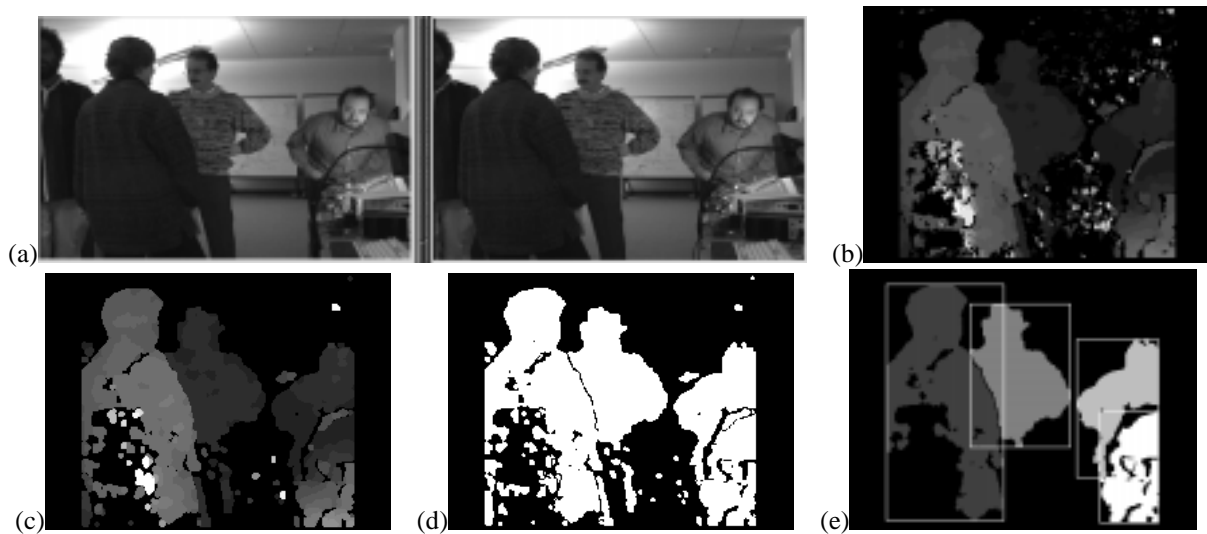


Figure 3. Stereo range processing to extract user silhouettes. (a) left/right image pair. (b) raw disparity computed using Census algorithm. (c) disparity after morphological smoothing. (d) regions of slowly varying disparity. (e) silhouettes recovered after connected components grouping.

Disparity estimation, segmentation, and grouping are repeated independently at each time step; range silhouettes are tracked from frame to frame based on position and size constancy. The centroid and size of each new range silhouette is compared to silhouettes from the previous time step. “Short-term” correspondences are indicated using a greedy algorithm starting with the closest unmatched region; for each new region the closest old region within a minimum threshold is marked as the correspondence matches.

### 3.2 Skin color localization

Skin color is a useful cue for tracking people’s faces and other body parts. We detect skin using a classification strategy which matches skin hue but is largely invariant to intensity or saturation, as this is robust to shading due to illumination and/or the absolute amount of skin pigment in a particular person.

We apply color segmentation processing to images obtained from one camera. Each image is initially represented with pixels corresponding to the red, green, and blue channels of the image, and is converted into a “log color-opponent” space. This space can directly represent the approximate hue of skin color, as well as its log intensity value. We convert  $(R, G, B)$  tuples into tuples of the form  $(\log(G), \log(R) - \log(G), \log(B) - (\log(R) + \log(G))/2)$ . Skin color is detected using a classifier with an empirically estimated Gaussian probability model of “skin” and “not-skin” in the log color-opponent color space. When a new pixel  $p$  is presented for classification, the likelihood ratio  $P(p = skin) / P(p = non-skin)$  is computed as a classification score. Our color representation is similar to that used in [2], but we estimate our classification criteria from examples rather than apply hand-tuned parameters. For computational efficiency at run-time, we precompute a lookup table over all possible color values.

After the lookup table has been applied, segmentation and grouping analysis are performed on the classification score image. Similar to the range case, we use morphological smoothing, threshold above a critical value, and apply connected component computation. However, there is one difference: before smoothing we apply the low-gradient mask from the *range* modality. This restricts color regions to be grown only within the boundary of range regions; if spurious background skin hue is present in the background it will not adversely affect the shape of foreground skin color regions.

As with range processing, classification, segmentation, and grouping are repeated at each time step. Short-term tracking is performed on recovered color regions based on similar centroid position and region size. When a color region changes size dramatically, we check to see if two regions merged, or if one region split into two. If so we

record the identity of the split or merged regions, to be used in the integration stage as described below.

Skin color regions that are above the midline of their associated range component, and which are appropriately sized at the given depth to be heads, are labeled as candidate heads and passed to the integration phase.

### 3.3 Face pattern detection

To distinguish head from hands and other body parts, and to localize the face within a region containing the head, we use pattern recognition methods which directly model the statistical appearance of faces based on intensity.

We based our implementation of this module on the CMU face detector [7] library. This library implements a neural network which models the appearance of frontal faces in a scene, and is similar to the pattern recognition approach described in [5]. Both methods are trained on a structured set of examples of faces and non-faces.

Face detection is initially applied over the entire image; when one or more detections are recorded, they are passed directly as candidate head locations to the integration phase. Short term tracking is implemented by focusing search in a new frame within windows around the detected locations in the previous frame. If a new detection is found within such a window it is considered to be in short-term correspondence with the previous detection; if no new detection is found and the detection in the previous frame overlapped a color or range region, then the face detection is updated to move with that region (as long as it persists).

## 4 Integrated Tracking

Our integration method is designed to take advantage of each module’s strengths: range is typically fast but coarse, color is fast and prone to false positives, and face pattern detection is slow and requires canonical pose and expression. We place priority on face detection hits, when available, and use color or range to update position from frame to frame.

For each range silhouette, we collect the range, color, and face detection candidate head features. As described above, when a candidate pattern detection head overlaps with a range or color candidate head, it persists and follows the range or color region. We record the relative offset of the face detection head with respect to the range or color head, and maintain that relationship in subsequent frames. This has the desired effect of allowing face detection to discriminate between head and hand regions in subsequent frames even when there may not be another face detection for several frames.

For each frame, we compute the location of a user’s head on the range silhouette as follows: if a face detection candidate head is present, we return it; otherwise we return any

location with overlapping range and color candidates, the location of the range candidate, or the location of a color candidate, in order of preference.

There is one special case in propagating face detection candidate heads. If the two color regions split or merge as described above, we take steps to allow the virtual face detection candidate head to follow the appropriate color region. We assume that the face is stationary between frames when deciding what color region to follow. If two regions have merged, the virtual detection follows the merged region, with offset such that the face’s absolute position on the screen is the same as the previous frame. If two regions have split, the face follows the region closest to its position in the previous frame. These heuristics are simple, but work in many cases where users are intermittently touching their face with their hands.

When the head location has been found, we update the estimate of head size. We have found that color is a relatively unreliable estimator of size; instead, we recompute size based on the results of the face detector and the range modules. When a face detection result has been found, we use it to determine the real size of the face. If no face detection hit has been found, we use an average model of real face size.

Our system can be configured in two modes: single- or multiple-person tracking. Single-person mode is most appropriate for interactive games or kiosks which are restricted to a single user; multiple-person is more appropriate for general surveillance and monitoring applications. In single person mode, we return only a single range silhouette; we initially choose the closest range region, and then follow that region until it is no longer tracked in the short-term.

## 5 Long-term tracking

When users are momentarily occluded or exit the scene, short-term tracking will fail since position and size correspondences in the individual modules are unavailable. To track users over medium and long-term time scales, we rely on statistical appearance models. Each visual processing module computes an estimate of certain user attributes, which are expected to be stable over longer time periods. These attributes are averaged as long as the underlying range silhouette continues to be tracked in the short-term, and used in a classification stage to establish medium and long-term correspondences.

Like multi-modal person detection and tracking, multi-modal person appearance classification is more robust than classification systems based on a single data modality. Height, color, and face pattern each offer independent classification data and are accompanied by similarly independent failure modes. Although face patterns are perhaps

the most common data source for current passive person classification methods, it is unusual to incorporate height or color information in identification systems because they do not provide sufficient discrimination to justify their use alone. However, combined with each other and with face patterns, height and color can provide important cues to disambiguate otherwise similar people, or help classify people when only degraded data is available in other modes.

### 5.1 Observed attributes

In the range module, we estimate the height of the user and use this as an attribute of identity. Height is obtained by computing the median value of the highest point of the a user silhouette in 3-D. In the color module, we compute the average color of the skin and hair regions; we plan to also add a histogram of clothing color. We define the hair region to be those pixels above the face but on the range silhouette; clothing can be defined as all other silhouette pixels not labeled as skin or hair.

In the face detector, we record an image of the actual face pattern wherever the detector records a hit. When a region is identified as a face based on the face pattern detection algorithm, the face pattern (greyscale subimage) in the target region is normalized and then passed to the classification stage. For optimal classification, we want the scale, alignment, and view of detected faces to be comparable. We resize the pattern to normalize for size, and discard images which are not in canonical pose or expression, which is determined by normalized correlation with an average canonical face.

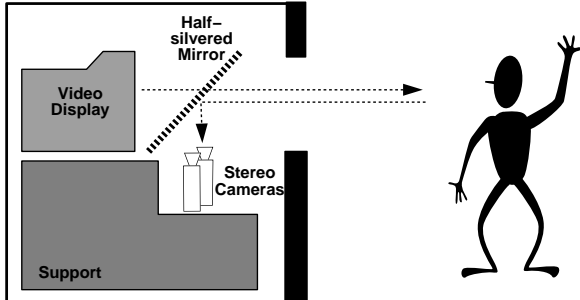
For “medium-term” tracking, e.g., over seconds or minutes of occlusion or absence, we rely on all of the above attributes. For “long-term” tracking, over hours or longer, we cannot rely on attributes which are not invariant with time of day or from day to day: we correct all color values with a mean color shift to account for changing illumination, and would exclude clothing color from the match computation.

### 5.2 Classification

In general, we compute statistics of these attributes while users are being tracked over the short-term, and compare against stored statistics of all previous tracked users.

When we observe a new person, we see if there is a previously tracked individual which could have generated the current observations. We find the previous individual most likely to have generated the new observations; if this probability is above a minimum threshold, we label the currently tracked region as being in correspondence with the previous individual. We integrate likelihood over time and modality: at time  $t$ , we find the identity estimate

$$u^* = \arg \max_j P(U_j | \omega) \tag{1}$$



**Figure 4. Display and viewing geometry: cameras and video-display share optical axis through a half-silvered mirror.**

where

$$P(U_j|\omega) = P(U_j|F_0, \dots, F_t, H_0, \dots, H_t, C_0, \dots, C_t) \quad (2)$$

where  $F_i, H_i$ , and  $C_i$  are the face pattern, height, and color observations at time  $i$ , and  $U_j$  are the saved statistics for person  $j$ . We restart time at  $t = 0$  when a new range silhouette is tracked. For the purposes of this discussion we assume  $P(U_j)$  is uniform across users. With Bayes rule and the assumption of modality independence, we have:

$$u^* = \arg \max_j ( P(F_0, \dots, F_t|U_j) P(H_0, \dots, H_t|U_j) P(C_0, \dots, C_t|U_j) ) \quad (3)$$

Since our observations are independent of the observed noise in our sensor and segmentation routines, the posterior probabilities at different times may be considered independent. With this we can incrementally compute probability in each modality:

$$P(F_0, \dots, F_t|U_j) = P(F_0, \dots, F_{t-1}|U_j) P(F_t|U_j) \quad (4)$$

and similarly for range and color data.

We collect mean and covariance data for the observed user color data, and mean and variance of user height; the likelihoods  $P(F_i|U_j)$  and  $P(C_i|U_j)$  are computed assuming a Gaussian density model. For face pattern data, we store the size- and position-normalized mean pattern for each user, and approximate  $P(F_t|C_p)$  with an empirically determined density which is a function of the normalized correlation of  $F_t$  with the the mean pattern for person  $j$ .

## 6 A Real-time Virtual Mirror Display

Our initial application of our integrated, multi-modal visual person tracking framework is to create a face-responsive visual display. We construct a video display

where cameras observe the user from the same optical axis as used by the display, and send estimates of the 3-D head position of observers of the screen to the application program. One application we have explored using this display is an interactive graphics experience in which users' faces are distorted in real-time. The effect is a virtual fun-house mirror, but in which only the face regions are distorted.

We create a virtual mirror by placing cameras so that they share the same optical axis as a video display, using a half-silvered mirror to merge the two optical paths. The cameras view the user through a 45-degree half mirror, so that the user can view a video monitor while also looking straight into (but not seeing) the cameras. Video from one camera is displayed on the monitor after the application of various computer graphics distortion effects, so as to create a virtual mirror effect. Figure 4 shows the display and viewing geometry of our apparatus. Using video texture mapping and the OpenGL graphics system, we have implemented graphics methods to distort faces on the screen using one of the following special effects: spherical expansion, spherical shrinking, swirl, lateral expansion, and a vertical melting effect. This creates a novel and entertaining interactive visual experience where users get immediate visual feedback from their tracked faces.

Our system is currently implemented using three computer systems (one PC, two SGI O2), a large NTSC video monitor, stereo video cameras, a dedicated stereo computation PC board, and the half-mirror imaging apparatus. The full tracking system, including all vision and graphics processing, runs at approximately 12Hz.

## 7 Results

We first demonstrated our system at the SIGGRAPH Conference from August 3-8, 1997 [1]. An estimated 5000 people over 6 days used our system (approximately two new users per minute, over 42 hours of operation). The goal of the system in this application was to identify the 3-D position and size of a single user's head in the scene, and apply a distortion effect in real-time only over the region of the image containing the user's face. The distorted image was then displayed on the virtual mirror screen. The system tracked the user while he or she was in the frame, and then switched to a new user.

Qualitatively, the system was a complete success. Our tracking results were able to localize video distortion effects on the user's face, and overall the system was interesting and fun for people to use. Figure 6 shows a typical final image displayed on the virtual mirror. The system performed well with both single users and crowded conditions; the background environment was quite visually noisy, with many spurious lighting effects being randomly projected throughout the conference hall, including onto the people



Figure 5. Color/Range stills of virtual mirror users collected during the SIGGRAPH '97 demonstration.



**Figure 6.** Example distortion output from virtual mirror application.

Modules Enabled			SIGGRAPH data	Lab data	Overall
<i>Color</i>	<i>Range</i>	<i>Pattern</i>			
✓	✓	✓	97%	96%	97%
	✓	✓	97%	95%	96%
✓	✓		97%	93%	95%
	✓		97%	90%	94%
✓		✓	92%	93%	92%
✓			90%	89%	90%
		✓	22% †	80%	44%

**Table 1.** Face detection and localization results on SIGGRAPH and Lab datasets using different combinations of input modules, ordered by increasing error rate. (†) The faces in the SIGGRAPH dataset were smaller than the size range the pattern module was trained to detect.

being tracked by our system.

## 7.1 Evaluation

We quantitatively evaluated the performance of our system using three off-line datasets: a set of stills captured at SIGGRAPH to evaluate detection performance, a set of stills of users in our laboratory, and a set of appearance statistics gathered from users in our laboratory who interacted with the system over several days. (Unfortunately we were not able to obtain observations of the same users across multiple days at the SIGGRAPH demonstration.)

We collected stills of users interacting with our system

every 15 seconds over a period of 3 hours at the SIGGRAPH demonstration. At each sample point we captured both a color image of the scene and a greyscale image of the output of the range module after disparity smoothing. We discarded images with no users present, yielding approximately 300 registered color/range pairs. Figure 5 shows examples of the collected stills. We also collected a similar set of approximately 200 registered range/color stills of users of the system while on display in our laboratory, similar to the images in Figures 2 and 3(a). Table 1 summarizes the single-person detection results we obtained on these test images. A correct match was defined when the corners of the estimated face region were sufficiently close to manually entered ground truth (within  $\frac{1}{4}$  of the face size). Overall, when all modules were functioning, we achieved a success rate of 97%; when the color and/or face detection module was removed, performance was still above 93%, indicating the power of the range cue for detecting likely head locations.

To evaluate our longer term tracking performance we used statistics gathered from 25 people in our laboratory who visited our display several times on different days. People’s hairstyle, clothing, and the exterior illumination conditions varied between the times data were collected. We tested whether our system was able to correctly identify users when they returned to the display. In general, our results were better for medium term tracking (intra-day) than for long term (inter-day) tracking, as would be expected. Table 2 shows the extended tracking results: the correct classification percentage is shown for each modality and for the combined observations from all modes. This table reflects the recognition rate using all of the data from each short-term tracking session: on average, users were tracked for 15 seconds before short-term tracking failed or they exited the workspace.

By integrating modes we were able to correctly establish correspondences between tracked users in all of the medium-term cases, which typically involved temporal gaps between 10 and 100 seconds. In the long-term cases, which typically reflected gaps of one day, integrated performance was 87%. A more complete description of medium- and long-term performance is shown in Figure 7 and Figure 10, respectively. These figures show the recognition rate vs rank threshold, i.e., the percentage of time the correct person was above a given rank in the ordered likelihood list of predicted users. We also measured our performance over time: Figures 8 and 11 compare the performance versus rank threshold at 4 different times during each testing session. Here we show only the multi-modal results; as expected, identification becomes more reliable over time as more data is collected. Figures 9 and 12 show the rank of the correct person over time, averaged across all test sessions; correct identification (average rank equals one) is almost al-



Performance	Medium-term (intra-day)	Long-term (inter-day)
Height	44%	20%
Color	84%	63%
Face pattern	84%	67%
Multi-modal	100%	87%

**Table 2. Extended tracking performance: correct identification rate at end of session.**

ways achieved within one second in the medium-term case, and within three seconds in the long-term case.

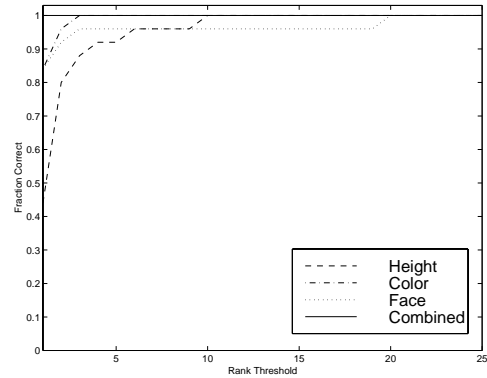
## 7.2 Discussion

We draw two main conclusions from the detection results; first, that range data is a powerful cue to detecting heads in complex scenes. Second, integration is useful: in almost every case, the addition of modules improved system performance. Performance was generally high, but individual module results varied considerably across datasets. In particular the face pattern module fared relatively poorly on the SIGGRAPH dataset. We believe that this is largely due to the small size and poor illumination of many of the faces in these images. Additionally, in both datasets our application encouraged people to make exaggerated expressions, which was beyond the scope of the training for this module.

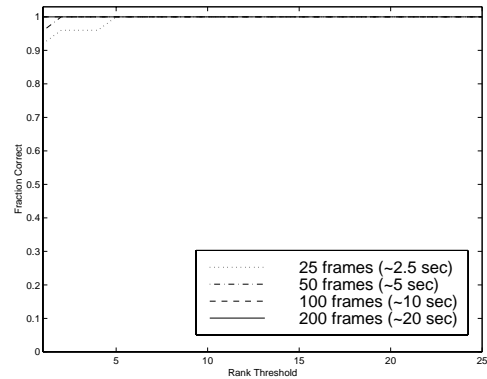
In contrast, for extended tracking it is clear from these results that the face pattern is the most valuable of the three modes when we consider all the data available during the session. Face pattern data is most discriminating at the *end* of the test session; however, the other modalities are dominant early in the session. The face detection module operates more slowly than the other modes, so the face pattern data is not available immediately and accumulates at a slower rate. Therefore, in the first few seconds the overall performance of the extended tracking system is due primarily to color and height data, and far exceeds the performance based on face pattern alone.

## 8 Conclusion

We have demonstrated a system which can respond to a user's face in real-time using completely passive and non-invasive techniques. Robust performance is achieved through the integration of three key modules: depth estimation to eliminate background effects, color classification for fast tracking, and pattern detection to discriminate the face from other body parts. We use descriptions of the user computed from the same modalities to track over longer time



**Figure 7. Medium-term tracking: performance vs rank threshold, results for each modality separately and then in combination.**

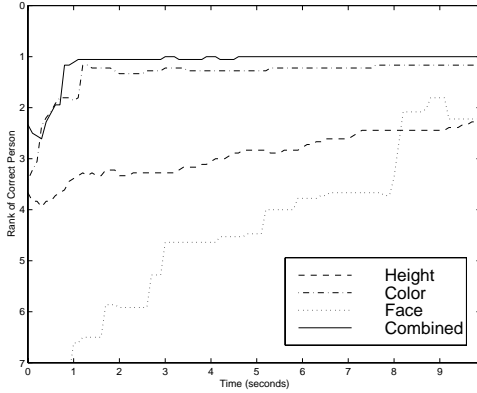


**Figure 8. Medium-term tracking: multi-modal performance vs rank threshold at 4 different time samples during a session.**

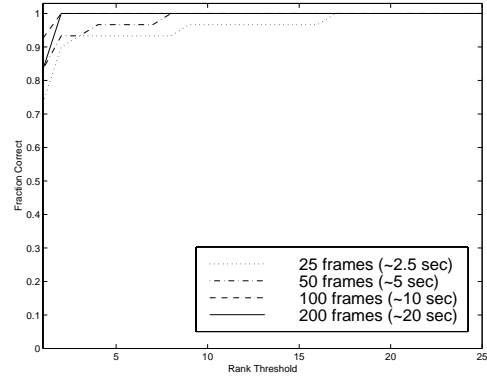
scales when the user is occluded or leaves the scene. Our system has application in interactive entertainment, telepresence/virtual environments, and intelligent kiosks which respond selectively according to the presence, pose, and identity of a user. We hope these and related techniques can eventually balance the I/O bandwidth between typical users and computer systems, so that they can control complicated virtual graphics objects and agents directly with their own expression.

## References

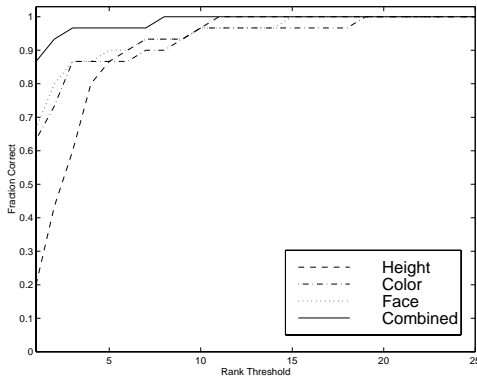
- [1] Darrell, T., Gordon, G., Woodfill, W., Baker, H., A Magic Morphin Mirror, SIGGRAPH '97 Visual Proceedings, ACM Press. 1997.
- [2] Margaret Fleck, David Forsyth, and Chris Bregler (1996) "Finding Naked People," European Confer-



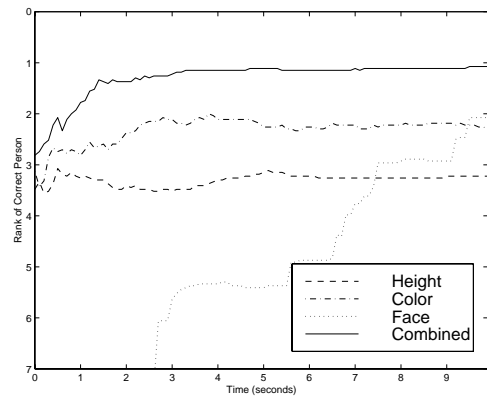
**Figure 9. Medium-term tracking: Average rank of correct person over time.**



**Figure 11. Long-term tracking: multi-modal performance vs rank threshold at 4 different time samples during a session.**



**Figure 10. Long-term tracking: performance vs rank threshold, results for each modality separately and then in combination.**



**Figure 12. Long-term tracking: Average rank of correct person over time.**

ence on Computer Vision , Volume II, pp. 592-602. 1996.

- [3] Kanade, T., Yoshida, A., Oda, K., Kano, H., and Tanaka, M., "A Video-Rate Stereo Machine and Its New Applications", Computer Vision and Pattern Recognition Conference, San Francisco, CA, 1996.
- [4] Maes, P., Darrell, T., Blumberg, B., and Pentland, A.P., "The ALIVE System: Wireless, Full-Body, Interaction with Autonomous Agents". ACM Multimedia Systems: Special Issue on on Multimedia and Multisensory Virtual Worlds, Sprint 1996.
- [5] Poggio, T., Sung, K.K., Example-based learning for view-based human face detection. Proceedings of the ARPA IU Workshop '94, II:843-850. 1994.
- [6] Rehg, J., Loughlin, M., and Waters, K., "Vision for a Smart Kiosk", Proc. IEEE Conf. Computer Vision and

Pattern Recognition, CVPR-97, pp. 690-696. IEEE Computer Society Press. 1997.

- [7] Rowley, H., Baluja, S., and Kanade, T., Neural Network-Based Face Detection, Proc. IEEE Conf. Computer Vision and Pattern Recognition, CVPR-96, pp. 203-207., IEEE Computer Society Press. 1996.
- [8] Zabih, R., and Woodfill, J., Non-parametric Local Transforms for Computing Visual Correspondence, Proceedings of the third European Conference on Computer Vision, Stockholm, pp. 151 - 158. May 1994.
- [9] Woodfill, J., and Von Herzen, B., Real-Time Stereo Vision on the PARTS Reconfigurable Computer, Proceedings IEEE Symposium on Field-Programmable Custom Computing Machines, Napa, pp. 242-250, April 1997.