

Appeared at EHCI '95, Grand Targhee, August '95.

Vision for man machine interaction

*J. L. Crowley
IMAG-GRAVIR
46 Ave Félix Viallet, 38031 Grenoble, France
email: jlc@imag.fr*

*J. Coutaz
IMAG-CLIPS
B.P 53, 38041 Grenoble CEDEX 9, France
email: coutaz@imag.fr*

Abstract

Computer vision provides a powerful tool for the interaction between man and machine. The barrier between physical objects (paper, pencils, calculators) and their electronic counterparts limits both the integration of computing into human tasks, and the population willing to adapt to the required input devices. Computer vision, coupled with video projection using low cost devices, makes it possible for a human to use any convenient object, including fingers, as digital input devices. In such an "augmented reality", information is projected onto ordinary objects and acquired by watching the way objects are manipulated. In the first part of this paper we describe experiments with techniques for watching the hands and recognizing gestures.

Vision of the face is an important aspect of human-to-human communication. We have been experimenting with the use of computer vision to "watch the face". In the second part of this paper we describe techniques for detecting, tracking and recognizing faces. When combined with real time image processing and active control of camera parameters, these techniques can greatly reduce the communications bandwidth required for video-phone and video-conference communications.

Keywords

Computer Vision, Multi-modal Man Machine Interaction, Eigen-Faces, Tracking, Principal Components Analysis

1 VISION AND MAN-MACHINE INTERACTION

One of the effects of the continued exponential growth in available computing power has been an exponential decrease in the cost of hardware for real time computer vision. This trend has been accelerated by the recent integration of image acquisition and processing hardware for multi-media applications in personal computers. Lowered cost has meant more wide-spread experimentation in real time computer vision, creating a rapid evolution in robustness and reliability and the development of architectures for integrated vision systems (Crowley 1984).

Man-machine interaction provides a fertile applications domain for this technological evolution. The barrier between physical objects (paper, pencils, calculators) and their electronic counterparts limits both the integration of computing into human tasks, and the population willing to adapt to the required input devices. Computer vision, coupled with video projection using low cost devices, makes it possible for a human to use any convenient object, including fingers, as digital input devices. Computer vision can permit a machine to track, identify and watch the face of a user. This offers the possibility of reducing bandwidth for video-telephone applications, for following the attention of a user by tracking his fixation point, and for exploiting facial expression as an additional information channel between man and machine.

Traditional computer vision techniques have been oriented toward using contrast contours (edges) to describe polyhedral objects. This approach has proved fragile even for man-made objects in a laboratory environment, and inappropriate for watching deformable non-polyhedral objects such as hands and faces. Thus man-machine interaction requires computer vision scientists to "go back to basics" to design techniques adapted to the problem. The following sections describe experiments with techniques for watching hands and faces.

2. WATCHING HANDS: GESTURE AS AN INPUT DEVICE

Human gesture serves three functional roles (Cadoz 1994): semiotic, ergotic, and epistemic.

- The *semiotic* function of gesture is to communicate meaningful information. The structure of a semiotic gesture is conventional and commonly results from shared cultural experience. The good-bye gesture, the American sign language, the operational gestures used to guide airplanes on the ground, and even the vulgar "finger", each illustrates the semiotic function of gesture.
- The *ergotic* function of gesture is associated with the notion of work. It corresponds to the capacity of humans to manipulate the real world, to create artefacts, or to change the state of the environment by "direct manipulation". Shaping pottery from clay, wiping dust, etc. result from ergotic gestures.
- The *epistemic* function of gesture allows humans to learn from the environment through tactile experience. By moving your hand over an object, you appreciate its structure, you may discover the material it is made of, as well as other properties.

All three functions may be augmented using an *instrumentt*: Examples include a handkerchief for the semiotic good-bye gesture, a turn-table for the ergotic shape-up gesture of pottery, or a dedicated artefact to explore the world (for example, a retro-active system such as the pantograph [Ramstein 94] to sense the invisible).

In Human Computer Interaction, gesture has been primarily exploited for its ergotic function: typing on a keyboard, moving a mouse, and clicking buttons. The epistemic role of gesture has emerged effectively from pen computing and virtual reality: ergotic gestures applied to an

electronic pen, to a data-glove or to a body-suit are transformed into meaningful expressions for the computer system. Special purpose interaction languages have been defined, typically 2-D pen gestures as in the Apple Newton, or 3-D hand gestures to navigate in virtual spaces or to control objects remotely (Baudel 1993).

With the exception of the electronic pen and the keyboard which both have their non-computerized counterparts, mice, data-gloves, and body-suits are “artificial add-on’s” that wire the user down to the computer. They are not real end-user instruments (as a hammer would be), but convenient tricks for computer scientists to sense human gesture.

We claim that computer vision can transform ordinary artefacts and even body parts into effective input devices. Krueger’s seminal work on the videoplace (Krueger 1991), followed recently by Wellner’s concept of digital desk (Wellner 1993) show that the camera can be used as a non-intrusive sensor for human gesture. However, to be effective the processing behind the camera must be fast and robust. The techniques used by Krueger and Wellner are simple concept demonstrations. They are fast but fragile and work only within highly constrained environments.

We are exploring advanced computer vision techniques to non-intrusively observe human gesture in a fast and robust manner. In the next section, we present FingerPaint, an experiment in the use of cross-correlation as a means of tracking *natural pointing devices* for a digital desk. By “natural pointing device”, we mean a bare finger or any real world artefact such as a pen or an eraser.

2.1 Projecting the workspace

In the digital desk a computer screen is projected onto a physical desk using a video-projector, such as a liquid-crystal "data-show" working with standard overhead projector. A video-camera is set up to watch the workspace such that the surface of the projected image and the surface of the imaged area coincide. This coincidence can not match "pixel to pixel" unless the camera and projector occupy the same physical space and use the same optics. Since this is impossible, it is necessary to master the transformation between the real workspace, and the imaged area. This transformation is a mapping between two planes.

The projection of a plane to another plane is an affine transformation. Thus the video projector can be used to project a reference frame onto the physical desk in the form of a set of points. The camera image of these four points permits the calibration of 6 coefficients (A, B, C, D, E, F) which transform image coordinates (i, j) to workspace coordinates (x, y).

$$x = A i + B j + C \qquad y = D i + E j + F.$$

The visual processes required for the digital desk are relatively simple. The basic operation is tracking of a pointing device, such as a finger, a pencil or an eraser. Such tracking should be supported by methods to determine what device to track and to detect when tracking has failed. A method is also required to detect the equivalent of "mouse-down" and “mouse-up” events.

The tracking problem can be expressed as: "Given an observation of an object at time t , determine the most likely position of the same object at time $t+\Delta T$ ". If different objects can be used as a pointing device, then the system must include some form of "trigger" which includes presentation of the pointing device to the system. The observation of the pointing device gives a small neighbourhood, $w(n, m)$, of an image $p(i, j)$. This neighbourhood will serve as a "reference template". The tracking problem can then be expressed as, given the position of the pointing device in the k th image, determine the most likely position of the pointing device in the $k+1$ th image. The size of the tracked neighbourhood must be determined such that the

neighbourhood includes a sufficiently large portion of the object to be tracked with a minimum of the background.

We have experimented with a number of different approaches to tracking pointing devices: these include color, correlation tracking, principal components and active contours (snakes) (Berard 1994) (Martin 1995). The active contour model (Kass 1987) presented problems which we believe can be resolved, but which will require additional experiments. Our current demonstration uses cross-correlation and principal components analysis.

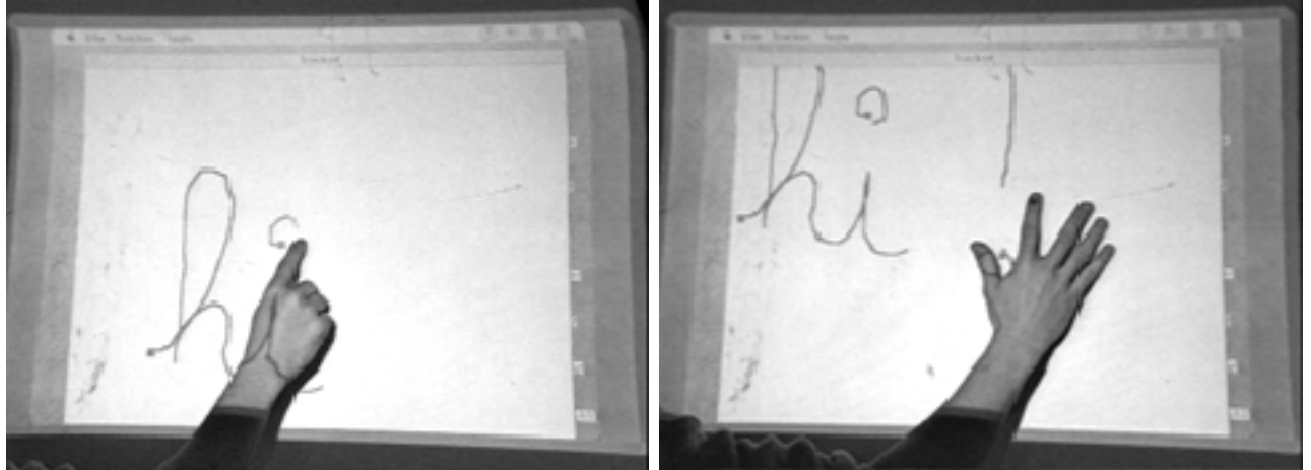


Figure 1 Drawing and placing with "Fingerpaint".

2.2 FingerPaint

As a simple demonstration, we constructed a program called "FingerPaint"*. FingerPaint runs on an Apple Quadra AV/840 and uses a work-space projected with an overhead projector using a liquid-crystal display "data-show". A CCD camera with an 18mm lens observes this workspace and provides visual input. "Finger down" and "finger up" events are simulated using the space bar of the keyboard but they could be sensed using a microphone attached to the surface of the desk. As illustrated in Figure 1, any "natural pointing device" such as a finger can be used to draw pictures and letters, or to move a drawing.

The image at time $(k+1)\Delta T$ to be searched will be noted as $p_{k+1}(i, j)$. The search process can generally be accelerated by restricting the search to a region of this image, denoted $s(i, j)$, and called a "Region of Interest". Our system uses a square search region of size M by M centered on the location where the reference template was detected in the previous image.

The robustness of the tracking system is reasonable but, as discussed below, its performance is inadequate with respect to Fitt's law (Card 1983). Preliminary experiments with local users indicate however that the current performance is acceptable for investigation purposes. In addition, the widespread availability of image acquisition and processing hardware adequate for real time correlation should alleviate our current performance problem. Most importantly, this demonstration has permitted us to explore the problems involved in watching gesture.

* The finger paint system has been implemented by Francois Berard.

2.3 Correlation as a tracking technique for the digital desk

The basic operation for a digital desk application is tracking some pointing device, such as a finger, a pencil or an eraser. The tracking problem can be expressed as: "Given an observation of an object at time t , determine the most likely location of the same object at time $t+\Delta T$ ". The pointing device can be modelled as a reference template. The reference template is a small neighbourhood, i.e., a window $w(m, n)$ of a picture $p(i, j)$ obtained at some prior time, t . The reference template is compared to an image neighborhood (i, j) , by computing a sum of squared differences (SSD) between the N by N template and the neighborhood of the image whose upper left corner is at (i, j) .

$$SSD(i, j) = \sum_{m=0}^{N-1} \sum_{n=0}^{N-1} (p_k(i+m, j+n) - w(m, n))^2$$

A SSD value is computed for each pixel, (i, j) , within the $M \times M$ search region. A perfect match between the template and the neighborhood gives a SSD value of 0. Such a perfect match is rarely obtained because of differences and appearance due to lighting and other effects. These effects can be minimized by normalising the energy in the template and the neighborhood.

Completing the squares of the SSD equation gives three terms, which can be written as:

$$SSD(i, j) = E_P^2(i, j) + E_W^2 - 2 \langle p_k(i+m, j+n), w(m, n) \rangle$$

The term $\langle p_k(i+m, j+n), w(m, n) \rangle$ is the inner product of the template $w(m, n)$ with the neighborhood $p_k(i, j)$. The term $E_P^2(i, j)$ represents the energy in the image neighborhood and E_W^2 is the energy in the reference window. The neighborhood and reference window may be normalized by dividing by $E_P^2(i, j)$ and E_W^2 , to give a normalized cross correlation, or NCC.

$$NCC(i, j) = \frac{\langle p_k(i+m, j+n), w(m, n) \rangle}{E_P^2(i, j) E_W^2}$$

Normalised cross correlation produces a peak with a value of 1.0 at a perfect match between window and neighborhood, and is relatively robust in the presence of noise, changes in scale and gray level, and image deformations (Martin and Crowley 1995). Hardware exists for computing a normalized cross correlation at video rates. However, in software, it is more efficient to use SSD.

Implementing cross-correlation by SSD requires solving practical problems such as determining the sizes of the reference template and of the search region, triggering and breaking tracking, and updating the reference template.

2.4 The size of the reference template

The size of the reference template must be determined such that it includes a sufficiently large portion of the device to be tracked and a minimum of the background. If the template window is too large, correlation may be corrupted by the background. On the other extreme, if the template is composed only of the interior of the pointing device then the reference template will be relatively uniform, and a high correlation peak will be obtained with any uniform region of the

image, including other parts of the pointing device. For a reasonable correlation peak, the reference template size should be just large enough to include the boundary of the pointing device, which contains the information used for detection and localisation.

In fingerprint, our workspace is of size 40 cm by 32 cm. This surface is mapped onto an image of 192 x 144 pixels, giving pixel sizes of 2 mm by 2.2 mm. At this resolution, a finger gives a correlation template of size 8 by 8 pixels or 16mm by 18mm, as shown in figure 2.

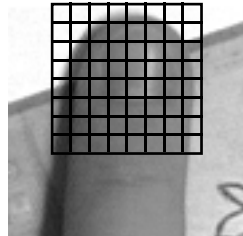


Figure 2 Reference template for a finger.

2.5 The size of the search region

The size M of the search region depends on the speed of the pointing device. Considerations based on Fitt's law (Card 1983) indicate a need for tracking speeds of up to 180 cm/sec. To verify this, we performed an experiment in which a finger was filmed making typical pointing movements in our workspace. The maximum speed observed in this experiment was $V_m = 139$ cm/sec. Expressed in pixels this gives $V_m = 695$ pixels/sec.

Given an image processing cycle time of ΔT seconds and a maximum pointer speed of V_m pixels/sec, it is possible to specify that the pointing device will be found within a radius of $M = \Delta T V_m$ pixels of its position in the previous frame. For images of 192 x 144 pixels, our built-in digitizer permits us to register images at a maximum frame rate of 24 frames per second, giving a cycle time of $\Delta T_{\max} = 41.7$ msec. This represents an upper limit on image acquisition speed which is attainable only if image tracking were to take no computation time.

The computational cost of cross-correlation is directly proportional to the number of pixels in the search region. Reducing the number of pixels will decrease the time needed for the inner loop of correlation by the same amount. This, in turn, increases the number of times that correlation can be operated within a unit time, further decreasing the region over which the search must be performed. Thus there is an inverse relation between the width of the search region, M , and the maximum tracking speed, V_m . The smaller the search region, the faster the finger movement that can be tracked, up to a limit set by the digitizing hardware.

The fastest tracking movement can be expected at a relatively small search region. This is confirmed by experiments. To verify the inverse relation between M and V_m , we systematically varied the size of the search region from $M = 10$ to 46 pixels and measured the cycle time that was obtained. The maximum speed of 126 pixels/sec is obtained with $M=26$. Although this is 5.5 times less than the maximum desirable speed (i.e., 695 pixels/sec), the system is quite usable to perform drawing and placements in a "natural" way.

2.6. Triggering and breaking tracking

When tracking is not active, the system monitors an N by N pixel "tracking trigger", $T_k(i,j)$, located in the lower right corner of the workspace. As each image is acquired at time k , the contents of this tracking trigger are subtracted from the contents at the previous image

$k-1$. This creates a difference image as shown in figure 3. The energy of the difference image is computed as

$$E_k = \sum_{m=0}^{N-1} \sum_{n=0}^{N-1} (T_k(m,n) - T_{k-1}(m,n))^2$$

When a pointing device enters the tracking trigger, the energy rises above a threshold. In order to assure that the tracking device is adequately positioned, the system waits until the difference energy drops back below the threshold before acquiring the reference template. At that point, the contents of the tracking trigger, $T_k(m, n)$ are saved as a reference image, and the tracking process is initiated.



Figure 3 Temporal difference of images in the reference square.

Tracking continues as long as the minimum value of SSD remains below a relatively high threshold. However, it can happen that the tracker locks on to a pattern on the digital desk (for example a photo of the pointing device!). To cover this eventuality, if the tracked location of the pointer stops moving for more than a few seconds (say 10), the system begins again to observe the difference energy in the tracking trigger. If the trigger energy rises above threshold, the tracker will break the previous track and re-initialise the reference pattern with the new contents of the tracking trigger.

2.7. Updating the reference mask

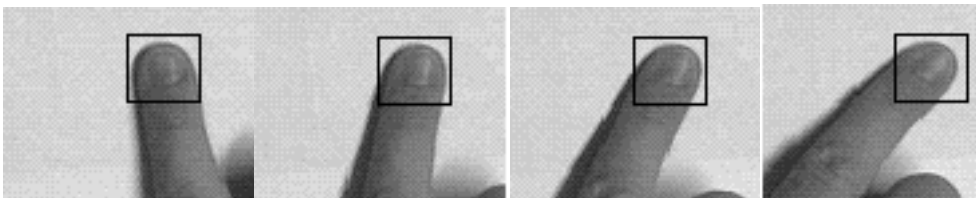


Figure 4 : Change in reference template as a function of finger orientation.

As the user moves the pointing device around the workspace, there is a natural tendency for the device to rotate, as shown in figure 4. This, in turn, may cause loss of tracking. In order to avoid loss of tracking, the smallest SSD value from each search is compared to a threshold. If the smallest SSD rises above this threshold, then the reference template is updated using the contents of the image at time $k-1$ at the detected position.

Tracking fingertips is an example of a simple fast visual process which can be used to change the nature of the interaction between man and machine. Vision can also be used to make the machine aware of the user by detecting, tracking and watching his face.

3 FACES: DETECTING, TRACKING, AND WATCHING THE USER

Face to face communication plays an important role in human to human communication. Thus it is natural to assume that an important quantity of non-verbal information can be obtained for man-machine interaction by watching faces. However, even more than hands, face interpretation poses difficult problems for established machine vision techniques. In this section we briefly report on experiments with simple techniques for detecting, tracking and interpreting images of faces. The key to robustness in such tracking and interpretation is the integration of complementary techniques.

3.1 Why watch a face?

Detection and interpretation of a face image can have a number of applications in machine vision. The most obvious use is to know whether a person is present in front of a computer screen. This makes a cute, but very expensive, screen saver. It is also possible to use face recognition as a substitute for a login password, presenting a person with his preferred workspace as soon as he appears in front of the computer system. Slightly more practical is the use of computer vision to watch the eyes and lips of a user. Eye tracking can be used to determine whether the user is looking at the computer screen and to which part his fixation is posed. This could conceivably be used to activate the currently active window in an interface. Observing the mouth to detect lip movements can be used to detect speech acts in order to trigger a speech recognition system.

None of the above uses would appear to be compelling enough to justify the cost of a camera and digitizer. However, there is an application for which people are ready to pay the costs of such hardware: video communications. Recognizing and tracking faces can have several important uses for the applications of video telephones and video conferencing. We are currently experimenting with combining face interpretation with a rudimentary sound processing system to determine the origin of spoken words and associate speech with faces. Each of the applications which we envisage require active computer control of the direction (pan and tilt) zoom (focal length), focus and aperture of a camera. Fortunately, such cameras are appearing on the market.

In the video-telephone application, we use an active camera to regulate zoom, pan, tilt, focus and aperture so as to keep the face of the user centered in the image at the proper size and focus, and with an appropriate light level. Such active camera control is not simply for esthetics. Keeping the face at the same place, same scale and same intensity level can dramatically reduce the information to be transmitted. One possible such coding is to define (on-line) a face space using principle components analysis (defined below) of the sample images from the last few minutes. Once the face basis vectors are transmitted, subsequent images can be transmitted as a short vector of face space coefficients. Effective use of this technique is only possible with active camera control. Other image codings can also be accelerated if the face image is normalized in position, size, gray level and held in focus.

In the video-conference scenario, a passive camera with a wide angle lens provides a global view of the persons seated around a conference table. An active camera provides a closeup of whichever person is speaking. When no one is speaking, and during transitions of the close-up camera, the wide-angle camera view can be transmitted. Face detection, operating on the wide-angle images, can be used to estimate the location of conference participants around the table. When a participant speaks, the high resolution camera can zoom onto the face of the speaker and hold the face in the center of the image.

What are the technologies required for the above applications? Both scenarios require an ability to detect and track faces, and an ability to servo control pan, tilt, zoom, focus and aperture so as to maintain a selected face in the desired position, scale and contrast level. From a hardware

standpoint, such an application requires a camera for which these axes can be controlled. Such camera heads are increasingly appearing on the market. For example, we have purchased a small RS232 controllable camera from a Japanese manufacturer which produces excellent color images for little more than the price of a normal color camera.

A second hardware requirement is the ability to digitize and process images at something close to video-rates. The classic bottle-neck here is communication of the image between the frame-grabber and the processor. Fortunately, the rush to multi-media applications has pushed a number of vendors to produce workstations in which a framegrabber is linked to the processor by such a high speed bus. Typical hardware available for a reasonable cost permits acquisition of up to 20 frames per second at full image size and full video rates for reduced resolution images. Adding simple image processing can reduce frame rates to 2 to 10 Hz (depending on image resolution). Such workstations are suitable for concept demonstrations and for experiments needed to define performance specifications. An additional factor of 2 (18 months) in band-width and processing power will bring us to full video-rates.

The questions we ask in the laboratory are: What software algorithms can be used for face detection, tracking and recognition, and what are the systems concepts needed to tie these processes together. Systems concepts have been the subject of our ESPRIT Basic Research Project "Vision as Process", described in the book (Crowley and Christensen 1994) or the paper (Crowley 94) for more details.

3.2 Detection: Finding a face with color

One of the simplest methods for detecting pixels which might be part of a face is to look for skin color. Human skin has a particular hue and saturation. The intensity, however, will vary as a function of the relative direction to the illumination. Of course, the perceived hue and saturation will be the product of the color of the ambient light and the skin color (Schiele and Weibul 1995).

We have found that candidate pixels for faces and hands can be detected very rapidly using a normalized color histogram. Histogram color can be normalised for changes in intensity by dividing the color vector by the luminance. This permits us to convert a color vector [R, G, B] having three dimensions into a vector [r, g] of normalised color having two dimensions. The normalized color histogram $H(r, g)$ provides a fast means of skin detection. The histogram is initialised by observing a patch of skin and counting the number of times each normalized color value occurs. The histogram contains the number of pixels which exhibit a particular color vector [r, g]. Dividing by the total number of pixels in the histogram, N , gives the probability of obtaining a particular vector given that the pixel observes skin.

$$P(\text{color} | \text{skin}) = \frac{1}{N} H(r, g)$$

Bayes rule can be used to determine the probability of skin given the color has values [r, g].

$$P(\text{skin} | \text{color}) = \frac{P(\text{skin})}{P(\text{color})} P(\text{color} | \text{skin})$$

$P(\text{skin})$ can be taken as constant. $P(\text{color})$ is the global statistic for the vector [r, g]. In practice this ratio is often approximated by a constant. The result at each pixel is the estimate of the probability of skin. An example (unfortunately printed here in black and white) is shown in figure 5a through 5d.



Figure 5 a. Black and white rendition of color image of Y.H. Berne



Figure 5 b. Probability of skin in color image of Y.H. Berne.

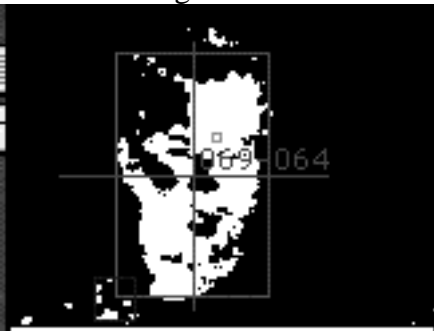


Figure 5 c. Thresholded skin probability with bounding box of connected components



Figure 5d Bounding box for the face.

Histogram matching can provide a very fast indicator that a face is present at a part of an image. However, reliability requires that this information be confirmed by another detection means. Such a means can be provided by a blink detector.

3.3 Finding a face by blink detection

A human must periodically blink to keep his eyes moist. Blinking is involuntary and fast. Most people do not notice when they blink. However, detecting a blinking pattern in an image sequence is an easy and reliable means to detect the presence of a face. Blinking provides a space-time signal which is easily detected and unique to faces. The fact that both eyes blink together provides a redundancy which permits blinking to be discriminated from other motions in the scene. The fact that the eyes are symmetrically positioned with a fixed separation provides a means to normalize the size and orientation of the head.

We have built a simple blink detector which works as follows: As each image is acquired, the previous image is subtracted. The resulting difference image generally contains a small boundary region around the outside of the head. If the eyes happened to be closed in one of the two images, there are also two small roundish regions over the eyes where the difference is significant.

The difference image is thresholded, and a connected components algorithm is run on the thresholded image. A bounding box is computed for each connected components. A candidate for an eye must have a bounding box within a particular horizontal and vertical size. Two such candidates must be detected with a horizontal separation of a certain range of sizes, and little vertical difference in the vertical separation. When this configuration of two small bounding boxes is detected, a pair of blinking eyes is hypothesized. The position in the image is determined from the center of the line between the bounding boxes. The distance to the face is measured from the separation. This permits to determine the size of a window which is used to

extract the face from the image. This simple technique has proven quite reliable for determining the position and size of faces.

3.4 Tracking the face by correlation (SSD)

Detection of a face by color is fast and reliable, but not always precise. Detection by blinking is precise, but requires capturing an image pair during a blink. Correlation can be used to complete these two techniques and to hold the face centered in the image as the head moves. The crucial question is how large a window to correlate (N) and over how large a search region to search.

We have found that a 7 by 7 search region extracted from the center of a face is sufficient to determine the precise location of the face in the next image. The size of the search region is determined by the speed with which a face can move between two frames. Calling on the finger tracking results above, we optimised the search region to maximize the frequency of image acquisition. With our current hardware this is provided by a search region to 26 by 26 pixels, but must be adjusted when several image processes are running in parallel.

Correlation can also be used in isolation to track a face. As a test of the robustness of correlation, we acquired a sequence of images of a head turning (shown in figure 6 below). We then correlated the center image from the sequence to produce the correlation graph shown in figure 7. The graph shows the results of zero-mean normalized correlation for the 40th face image with the other images in this set.

