

Extension 3: Vision

Text by Golan Levin

A well-known anecdote relates how, sometime in 1966, the legendary artificial intelligence pioneer Marvin Minsky directed an undergraduate student to solve “the problem of computer vision” as a summer project.¹ This anecdote is often resuscitated to illustrate how egregiously the difficulty of computational vision has been underestimated. Indeed, nearly forty years later the discipline continues to confront numerous unsolved (and perhaps unsolvable) challenges, particularly with respect to high-level “image understanding” issues such as pattern recognition and feature recognition. Nevertheless, the intervening decades of research have yielded a great wealth of well-understood, low-level techniques that are able, under controlled circumstances, to extract meaningful information from a camera scene. These techniques are indeed elementary enough to be implemented by novice programmers at the undergraduate or even high-school level.

Computer vision in interactive art

The first interactive artwork to incorporate computer vision was, interestingly enough, also one of the first interactive artworks. Myron Krueger’s legendary *Videoplace*, developed between 1969 and 1975, was motivated by his deeply felt belief that the entire human body ought to have a role in our interactions with computers. In the *Videoplace* installation, a participant stands in front of a backlit wall and faces a video projection screen. The participant’s silhouette is then digitized and its posture, shape, and gestural movements analyzed. In response, *Videoplace* synthesizes graphics such as small “critters” which climb up the participant’s projected silhouette, or colored loops drawn between the participant’s fingers. Krueger also allowed participants to paint lines with their fingers, and, indeed, entire shapes with their bodies; eventually, *Videoplace* offered more than fifty compositions and interactions. *Videoplace* is notable for many “firsts” in the history of human-computer interaction. Some of its interaction modules allowed two participants in mutually remote locations to participate in the same shared video space, connected across the network—an implementation of the first multiperson virtual reality, or, as Krueger termed it, an “artificial reality.” *Videoplace*, it should be noted, was developed before the mouse became the ubiquitous desktop device it is today, and was (in part) created to demonstrate interface alternatives to the keyboard terminals that dominated computing so completely in the early 1970s.

Messa di Voce (p. 511), created by this text’s author in collaboration with Zachary Lieberman, uses whole-body vision-based interactions similar to Krueger’s, but combines them with speech analysis and situates them within a kind of projection-based augmented reality. In this audiovisual performance, the speech, shouts, and

songs produced by two abstract vocalists are visualized and augmented in real time by synthetic graphics. To accomplish this, a computer uses a set of vision algorithms to track the locations of the performers' heads; this computer also analyzes the audio signals coming from the performers' microphones. In response, the system displays various kinds of visualizations on a projection screen located just behind the performers; these visualizations are synthesized in ways that are tightly coupled to the sounds being spoken and sung. With the help of the head-tracking system, moreover, these visualizations are projected such that they appear to emerge directly from the performers' mouths.

Rafael Lozano-Hemmer's installation *Standards and Double Standards* (2004) incorporates full-body input in a less direct, more metaphorical context. This work consists of fifty leather belts, suspended at waist height from robotic servomotors mounted on the ceiling of the exhibition room. Controlled by a computer vision-based tracking system, the belts rotate automatically to follow the public, turning their buckles slowly to face passers-by. Lozano-Hemmer's piece "turns a condition of pure surveillance into an 'absent crowd' using a fetish of paternal authority: the belt."²

The theme of surveillance plays a foreground role in David Rokeby's *Sorting Daemon* (2003). Motivated by the artist's concerns about the increasing use of automated systems for profiling people as part of the "war on terrorism," this site-specific installation works toward the automatic construction of a diagnostic portrait of its social (and racial) environment. Rokeby writes: "The system looks out onto the street, panning, tilting and zooming, looking for moving things that might be people. When it finds what it thinks might be a person, it removes the person's image from the background. The extracted person is then divided up according to areas of similar colour. The resulting swatches of colour are then organized [by hue, saturation and size] within the arbitrary context of the composite image" projected onsite at the installation's host location.³

Another project themed around issues of surveillance is *Suicide Box*, by the Bureau of Inverse Technology (Natalie Jeremijenko and Kate Rich). Presented as a device for measuring the hypothetical "despondency index" of a given locale, the *Suicide Box* nevertheless records very real data regarding suicide jumpers from the Golden Gate Bridge. According to the artists, "The *Suicide Box* is a motion-detection video system, positioned in range of the Golden Gate Bridge, San Francisco, in 1996. It watched the bridge constantly and when it recognized vertical motion, captured it to a video record. The resulting footage displays as a continuous stream the trickle of people who jump off the bridge. The Golden Gate Bridge is the premiere suicide destination in the United States; a 100-day initial deployment period of the *Suicide Box* recorded 17 suicides. During the same time period the Port Authority counted only 13."⁴ Elsewhere, Jeremijenko has explained that "the idea was to track a tragic social phenomenon which was not being counted—that is, doesn't count."⁵ The *Suicide Box* has met with considerable controversy, ranging from ethical questions about recording the suicides, to disbelief that the recordings could be real. Jeremijenko, whose aim is to address the hidden politics of technology, has pointed out that such attitudes express a recurrent theme—"the inherent suspicion of artists working with material evidence"—evidence obtained, in this case, with the help of machine vision-based surveillance.

Considerably less macabre is Christian Möller's clever Cheese installation (2003), which the artist developed in collaboration with the California Institute of Technology and the Machine Perception Laboratories of the University of California, San Diego. Motivated, perhaps, by the culture shock of his relocation to Hollywood, the German-born Möller directed "six actresses to hold a smile for as long as they could, up to one and a half hours. Each ongoing smile is scrutinized by an emotion recognition system, and whenever the display of happiness fell below a certain threshold, an alarm alerted them to show more sincerity."⁶ The installation replays recordings of the analyzed video on six flat-panel monitors, with the addition of a fluctuating graphic level-meter to indicate the strength of each actress' smile. The technical implementation of this artwork's vision-based emotion recognition system is quite sophisticated.

As can be seen from these examples, artworks employing computer vision range from the highly formal and abstract to the humorous and sociopolitical. They concern themselves with the activities of willing participants, paid volunteers, or unaware strangers. They track people of interest at a wide variety of spatial scales, from extremely intimate studies of their facial expressions, to the gestures of their limbs, to the movements of entire bodies. The examples above represent just a small selection of notable works in the field and of the ways in which people (and objects) have been tracked and dissected by video analysis. Other noteworthy artworks that use machine vision include Marie Sester's *Access*; Joachim Sauter and Dirk Lüsebrink's *Zerseher* and *Bodymover*; Scott Snibbe's *Boundary Functions* and *Screen Series*; Camille Utterback and Romy Achituv's *TextRain*; Jim Campbell's *Solstice*; Christa Sommerer and Laurent Mignonneau's *A-Volve*; Danny Rozin's *Wooden Mirror*; Chico MacMurtrie's *Skeletal Reflection*, and various works by Simon Penny, Toshio Iwai, and numerous others. No doubt many more vision-based artworks remain to be created, especially as these techniques gradually become incorporated into developing fields like physical computing and robotics.

Elementary computer vision techniques

To understand how novel forms of interactive media can take advantage of computer vision techniques, it is helpful to begin with an understanding of the kinds of problems that vision algorithms have been developed to address, and of their basic mechanisms of operation. The fundamental challenge presented by digital video is that it is computationally "opaque." Unlike text, digital video data in its basic form—stored solely as a stream of rectangular pixel buffers—contains no intrinsic semantic or symbolic information. There is no widely agreed upon standard for representing the content of video, in a manner analogous to HTML, XML, or even ASCII for text (though some new initiatives, notably the MPEG-7 description language, may evolve into such a standard in the future). As a result, a computer, without additional programming, is unable to answer even the most elementary questions about whether a video stream contains a person or object, or whether an outdoor video scene shows daytime or nighttime, et cetera. The discipline of computer vision has developed to address this need.

Many low-level computer vision algorithms are geared to the task of distinguishing which pixels, if any, belong to people or other objects of interest in the scene. Three elementary techniques for accomplishing this are frame differencing, which attempts to locate features by detecting their movements; background subtraction, which locates visitor pixels according to their difference from a known background scene; and brightness thresholding, which uses hoped-for differences in luminosity between foreground people and their background environment. These algorithms, described in the following examples, are extremely simple to implement and help constitute a base of detection schemes from which sophisticated interactive systems may be built.

Example 1: Detecting motion (p. 556)

The movements of people (or other objects) within the video frame can be detected and quantified using a straightforward method called frame differencing. In this technique, each pixel in a video frame F_1 is compared with its corresponding pixel in the subsequent frame F_2 . The difference in color and/or brightness between these two pixels is a measure of the amount of movement in that particular location. These differences can be summed across all of the pixels' locations to provide a single measurement of the aggregate movement within the video frame. In some motion detection implementations, the video frame is spatially subdivided into a grid of cells, and the values derived from frame differencing are reported for each of the individual cells. For accuracy, the frame differencing algorithm depends on relatively stable environmental lighting, and on having a stationary camera (unless it is the motion of the camera that is being measured).

Example 2: Detecting presence (p. 557)

A technique called background subtraction makes it possible to detect the presence of people or other objects in a scene, and to distinguish the pixels that belong to them from those that do not. The technique operates by comparing each frame of video with a stored image of the scene's background, captured at a point in time when the scene was known to be empty. For every pixel in the frame, the absolute difference is computed between its color and that of its corresponding pixel in the stored background image; areas that are very different from the background are likely to represent objects of interest. Background subtraction works well in heterogeneous environments, but it is very sensitive to changes in lighting conditions and depends on objects of interest having sufficient contrast against the background scene.

Example 3: Detection through brightness thresholding (p. 559)

With the aid of controlled illumination (such as backlighting) and/or surface treatments (such as high-contrast paints), it is possible to ensure that objects are considerably darker or lighter than their surroundings. In such cases objects of interest can be distinguished based on their brightness alone. To do this, each video pixel's brightness is compared to a threshold value and tagged accordingly as foreground or background.



Example 1. Detects motion by comparing each video frame to the previous frame. The change is visualized and is calculated as a number.



Example 2. Detects the presence of someone or something in front of the camera by comparing each video frame with a previously saved frame. The change is visualized and is calculated as a number.



Example 3. Distinguishes the silhouette of people or objects in each video frame by comparing each pixel to a threshold value. The circle is filled with white when it is within the silhouette.



Example 4. Tracks the brightest object in each video frame by calculating the brightest pixel. The light from the flashlight is the brightest element in the frame, therefore the circle follows it.

Example 4: Brightness tracking (p. 560)

A rudimentary scheme for object tracking, ideal for tracking the location of a single illuminated point (such as a flashlight), finds the location of the single brightest pixel in every fresh frame of video. In this algorithm, the brightness of each pixel in the incoming video frame is compared with the brightest value yet encountered in that frame; if a pixel is brighter than the brightest value yet encountered, then the location and brightness of that pixel are stored. After all of the pixels have been examined, then the brightest location in the video frame is known. This technique relies on an operational assumption that there is only one such object of interest. With trivial modifications, it can equivalently locate and track the darkest pixel in the scene, or track multiple, differently colored objects.

Of course, many more software techniques exist, at every level of sophistication, for detecting, recognizing, and interacting with people and other objects of interest. Each of the tracking algorithms described above, for example, can be found in elaborated versions that amend its various limitations. Other easy-to-implement algorithms can compute specific features of a tracked object, such as its area, center of mass, angular orientation, compactness, edge pixels, and contour features such as corners and cavities. On the other hand, some of the most difficult to implement algorithms, representing the cutting edge of computer vision research today, are able (within limits) to recognize unique people, track the orientation of a person's gaze, or correctly identify facial expressions. Pseudocodes, source codes, or ready-to-use implementations of all of these techniques can be found on the Internet in excellent resources like Daniel Huber's Computer Vision Homepage, Robert Fisher's HIPR (Hypermedia Image Processing Reference), or in the software toolkits discussed on pages 554-555.

Computer vision in the physical world

Unlike the human eye and brain, no computer vision algorithm is completely general, which is to say, able to perform its intended function given any possible video input. Instead, each software tracking or detection algorithm is critically dependent on certain unique assumptions about the real-world video scene it is expected to analyze. If any of these expectations is not met, then the algorithm can produce poor or ambiguous results or even fail altogether. For this reason, it is essential to design physical conditions in tandem with the development of computer vision code, and to select the software techniques that are most compatible with the available physical conditions.

Background subtraction and brightness thresholding, for example, can fail if the people in the scene are too close in color or brightness to their surroundings. For these algorithms to work well, it is greatly beneficial to prepare physical circumstances that naturally emphasize the contrast between people and their environments. This can be achieved with lighting situations that silhouette the people, for example, or through the use of specially colored costumes. The frame-differencing technique, likewise, fails to detect people if they are stationary. It will therefore have very different degrees of

success detecting people in videos of office waiting rooms compared with, for instance, videos of the Tour de France bicycle race.

A wealth of other methods exist for optimizing physical conditions in order to enhance the robustness, accuracy, and effectiveness of computer vision software. Most are geared toward ensuring a high-contrast, low-noise input image. Under low-light conditions, for example, one of the most helpful such techniques is the use of infrared (IR) illumination. Infrared, which is invisible to the human eye, can supplement the light detected by conventional black-and-white security cameras. Using IR significantly improves the signal-to-noise ratio of video captured in low-light circumstances, and can even permit vision systems to operate in (apparently) complete darkness. Another physical optimization technique is the use of retroreflective marking materials, such as those manufactured by 3M Corporation for safety uniforms. These materials are remarkably efficient at reflecting light back toward their source of illumination and are ideal aids for ensuring high-contrast video of tracked objects. If a small light is placed coincident with the camera's axis, objects with retroreflective markers will be detected with tremendous reliability.

Finally, some of the most powerful physical optimizations for machine vision can be made without intervening in the observed environment at all, through well-informed selections of the imaging system's camera, lens, and frame-grabber components. To take one example, the use of a "telecentric" lens can significantly improve the performance of certain kinds of shape-based or size-based object recognition algorithms. For this type of lens, which has an effectively infinite focal length, magnification is nearly independent of object distance. As one manufacturer describes it, "an object moved from far away to near the lens goes into and out of sharp focus, but its image size is constant. This property is very important for gauging three-dimensional objects, or objects whose distance from the lens is not known precisely."⁷ Likewise, polarizing filters offer a simple, nonintrusive solution to another common problem in video systems, namely glare from reflective surfaces. And a wide range of video cameras are available, optimized for conditions like high-resolution capture, high-frame-rate capture, short exposure times, dim light, ultraviolet light, and thermal imaging. It pays to research imaging components carefully.

As we have seen, computer vision algorithms can be selected to best negotiate the physical conditions presented by the world, and likewise, physical conditions can be modified to be more easily legible to vision algorithms. But even the most sophisticated algorithms and the highest-quality hardware cannot help us find meaning where there is none, or track an object that cannot be described in code. It is therefore worth emphasizing that some visual features contain more information about the world, and are also more easily detected by the computer, than others. In designing systems to "see for us," we must not only become freshly awakened to the many things about the world that make it visually intelligible to us, but also develop a keen intuition about their ease of computability. The sun is the brightest point in the sky, and by its height also indicates the time of day. The mouth cavity is easily segmentable as a dark region, and the circularity of its shape is also closely linked to vowel sound. The pupils of the eye emit an easy-to-track infrared retroreflection, and they also indicate a person's direction

of gaze. Simple frame differencing makes it easy to track motion in a video. The *Suicide Box* (p. 548) uses this technique to dramatic effect.

Tools for computer vision

It can be a rewarding experience to implement machine vision techniques from scratch using code such as the examples provided in this section. To make this possible, the only requirement of one's software development environment is that it should provide direct read-access to the array of video pixels obtained by the computer's frame-grabber. Hopefully, the example algorithms discussed earlier illustrate that creating low-level vision algorithms from first principles isn't so hard. Of course, a vast range of functionality can also be obtained immediately from readily available solutions. Some of the most popular machine vision toolkits take the form of plug-ins or extension libraries for commercial authoring environments geared toward the creation of interactive media. Such plug-ins simplify the developer's problem of connecting the results of the vision-based analysis to the audio, visual, and textual affordances generally provided by such authoring systems.

Many vision plug-ins have been developed for Max/MSP/Jitter, a visual programming environment that is widely used by electronic musicians and VJs. Originally developed at the Parisian IRCAM research center in the mid-1980s and now marketed commercially by the California-based Cycling'74 company, this extensible environment offers powerful control of (and connectivity between) MIDI devices, real-time sound synthesis and analysis, OpenGL-based 3D graphics, video filtering, network communications, and serial control of hardware devices. The various computer vision plug-ins for Max/MSP/Jitter, such as David Rokeby's SoftVNS, Eric Singer's Cyclops, and Jean-Marc Pelletier's CV.Jit, can be used to trigger any Max processes or control any system parameters. Pelletier's toolkit, which is the most feature-rich of the three, is also the only one that is freeware. CV.Jit provides abstractions to assist users in tasks such as image segmentation, shape and gesture recognition, motion tracking, etc., as well as educational tools that outline the basics of computer vision techniques.

Some computer vision toolkits take the form of stand-alone applications and are designed to communicate the results of their analyses to other environments (such as Processing, Director, or Max) through protocols like MIDI, serial RS-232, UDP, or TCP/IP networks. BigEye, developed by the STEIM (Studio for Electro-Instrumental Music) group in Holland, is a simple and inexpensive example. BigEye can track up to 16 objects of interest simultaneously, according to their brightness, color, and size. The software allows for a simple mode of operation, in which the user can quickly link MIDI messages to many object parameters, such as position, speed, and size. Another example is the powerful EyesWeb open platform, a free system developed at the University of Genoa. Designed with a special focus on the analysis and processing of expressive gesture, EyesWeb includes a collection of modules for real-time motion tracking and extraction of movement cues from human full-body movement; a collection of modules for analysis of occupation of 2D space; and a collection of modules for extraction of features

from trajectories in 2D space. EyesWeb's extensive vision affordances make it highly recommended for students.

The most sophisticated toolkits for computer vision generally demand greater familiarity with digital signal processing, and require developers to program in compiled languages like C++ rather than languages like Java, Lingo, or Max. The Intel Integrated Performance Primitives (IPP) library for example, is among the most general commercial solutions available for computers with Intel-based CPUs. The OpenCV library, by contrast, is a free, open source toolkit with nearly similar capabilities and a tighter focus on commonplace computer vision tasks. The capabilities of these tools, as well as all of those mentioned above, are continually evolving.

Processing includes a basic video library that handles getting pixel information from a camera or movie file, as demonstrated in the examples included with this text. The computer vision capabilities of Processing are extended by libraries like Myron, which handles video input and has basic image processing capabilities. Other libraries connect Processing to EyesWeb and OpenCV. They can be found on the libraries page of the Processing website: www.processing.org/reference/libraries.

Conclusion

Computer vision algorithms are increasingly used in interactive and other computer-based artworks to track people's activities. Techniques exist that can create real-time reports about people's identities, locations, gestural movements, facial expressions, gait characteristics, gaze directions, and other characteristics. Although the implementation of some vision algorithms requires advanced understanding of image processing and statistics, a number of widely used and highly effective techniques can be implemented by novice programmers in as little as an afternoon. For artists and designers who are familiar with popular multimedia authoring systems like Macromedia Director and Max/MSP/Jitter, a wide range of free and commercial toolkits are also available that provide ready access to more advanced vision functionalities.

Since the reliability of computer vision algorithms is limited according to the *quality* of the incoming video scene and the definition of a scene's quality is determined by the specific algorithms that are used to analyze it, students approaching computer vision for the first time are encouraged to apply as much effort to optimizing their physical scenario as they do to their software code. In many cases, a cleverly designed physical environment can permit the tracking of phenomena that might otherwise require much more sophisticated software. As computers and video hardware become more available, and software-authoring tools continue to improve, we can expect to see the use of computer vision techniques increasingly incorporated into media-art education and into the creation of games, artworks, and many other applications.

Notes

1. <http://mechanism.ucsd.edu/~bill/research/mercier/2ndlecture.pdf>.
2. <http://www.fundacion.telefonica.com/at/rlh/eproyecto.html>.

3. <http://homepage.mac.com/davidrokeby/sorting.html>.
4. <http://www.bureauit.org/sbox>.
5. <http://www.wired.com/news/culture/0,1284,64720,00.html>.
6. <http://www.christian-moeller.com>.
7. <http://www.mellesgriot.com/pdf/pg11-19.pdf>.

Code

Video can be captured into Processing from USB cameras, IEEE 1394 cameras, or video cards with composite or S-video input devices. The examples that follow assume you already have a camera working with Processing. Before trying these examples, first get the examples included with the Processing software to work. Sometimes you can plug a camera in to your computer and it will work immediately. Other times it's a difficult process involving trial-and-error changes. It depends on the operating system, the camera, and how the computer is configured. For the most up-to-date information, refer to the Video reference on the Processing website: www.processing.org/reference/libraries.

Example 1: Detecting motion

```
// Quantify the amount of movement in the video frame using frame-differencing

import processing.video.*;

int numPixels;
int[] previousFrame;
Capture video;

void setup(){
  size(640, 480); // Change size to 320 x 240 if too slow at 640 x 480
  video = new Capture(this, width, height, 24);
  numPixels = video.width * video.height;
  // Create an array to store the previously captured frame
  previousFrame = new int[numPixels];
}

void draw() {
  if (video.available()) {
    // When using video to manipulate the screen, use video.available() and
    // video.read() inside the draw() method so that it's safe to draw to the screen
    video.read(); // Read the new frame from the camera
    video.loadPixels(); // Make its pixels[] array available

    int movementSum = 0; // Amount of movement in the frame
    loadPixels();

    for (int i = 0; i < numPixels; i++) { // For each pixel in the video frame...
```