

Golan Levin

Computer vision for artists and designers: pedagogic tools and techniques for novice programmers

Received: 17 January 2005 / Accepted: 19 August 2005 / Published online: 17 June 2006
© Springer-Verlag London Limited 2006

Abstract This article attempts to demystify computer vision for novice programmers through a survey of new applications in the arts, system design considerations, and contemporary tools. It introduces the concept and gives a brief history of computer vision within interactive art from Myron Kruger to the present. Basic techniques of computer vision such as detecting motion and object tracking are discussed in addition to various software applications created for exploring the topic. As an example, the results of a 1-week machine vision workshop are presented to show how designers are able to apply their skills toward creating novel uses of these technologies. The article concludes with a listing of code for basic computer vision techniques.

Keywords Computer vision · Machine vision · Interactive art · Artistic applications · Authoring tools · Education

1 Introduction

“Computer vision” refers to a broad class of algorithms that allow computers to make intelligent assertions about digital images and video. Historically, the creation of computer vision systems has been regarded as the exclusive domain of expert researchers and engineers in the fields of signal processing and artificial intelligence. Likewise, the scope of application development for computer vision technologies, perhaps constrained by conventional structures for research funding, has generally been limited to military and law-enforcement purposes. Recently, however, improvements in software development tools for student programmers and interactive-media artists—in combination with the rapid growth of open-source code-sharing communities, predictable increases in PC processor speeds, and plummeting costs of digital video hardware—have made widespread artistic experimentation with computer vision techniques a reality. The result is a proliferation of new practitioners with an abundance of new application ideas, and the incorporation of

G. Levin
School of Art, Carnegie Mellon University,
CFA-300, 5000 Forbes Avenue, Pittsburgh, PA 15213-3890, USA
E-mail: golan@flong.com · Tel.: +1-917-5207456

computer vision techniques into the design vocabularies of novel artworks, games, home automation systems, and other areas. This article attempts to demystify computer vision for novice programmers, through a survey of new applications in the arts, system design considerations, and contemporary tools.

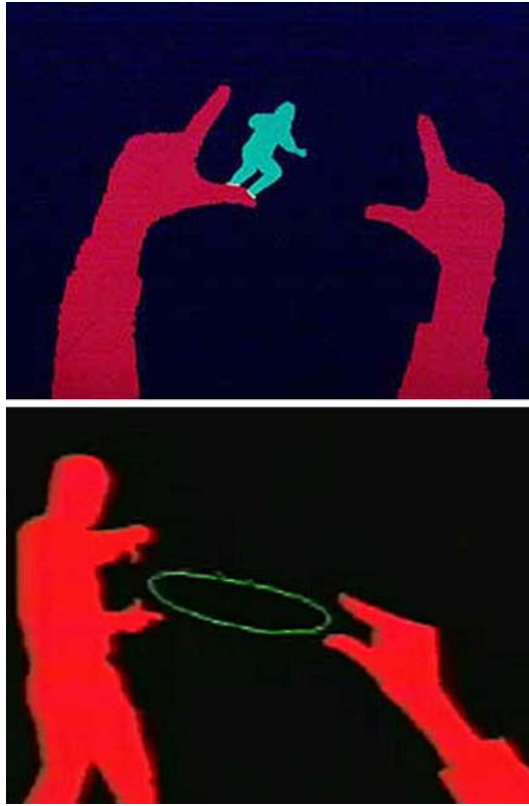
A well-known anecdote relates how, sometime in 1966, the legendary Artificial Intelligence pioneer Marvin Minsky directed an undergraduate student to solve “the problem of computer vision” as a summer project (Bechtel 2003). This anecdote is often resuscitated to illustrate how egregiously the difficulty of computational vision has been underestimated. Indeed, nearly 40 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.

This paper attempts to demystify computer vision for novice programmers, emphasizing the use of vision-based detection and tracking techniques in the interactive media arts. The next section of this article introduces some of the ways in which computer vision has found artistic applications outside of industrial and military research. Section 3 presents an overview of several basic but widely-used vision algorithms, with example code included in appendices at the end of the article. Although it is easy to suppose that sophisticated software is all one needs to create a computer vision system, Sect. 4 makes the case that a well-prepared physical environment can dramatically improve algorithmic performance and robustness. The remaining sections present a brief survey of several artist-friendly new computer vision toolkits, and an example of a student project, developed by novice programmers in a workshop structured around the considerations presented in this article.

2 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 over 50 different compositions and interactions (Figs. 1–2).

Videoplace was notable for many “firsts” in the history of human-computer interaction. Some of its interaction modules, for example the ones illustrated



Figs. 1–2 Interaction modules from Myron Krueger's *Videoplace*, 1969–1975

here, allowed two participants in mutually remote locations to participate in the same shared video space, connected across the network—an implementation of the first multi-person virtual reality, or, as Krueger termed it, an “artificial reality”. Videoplace, it should be noted, was developed before Douglas Engelbart’s mouse became the ubiquitous desktop device it is today, and was (in part) created to demonstrate interface alternatives to the keyboard terminals which dominated computing so completely in the early 1970s. Remarkably enough, the original Videoplace system is still operational as of this writing.

Messa di Voce (2003), created by this article’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 which are tightly coupled to the sounds being spoken and

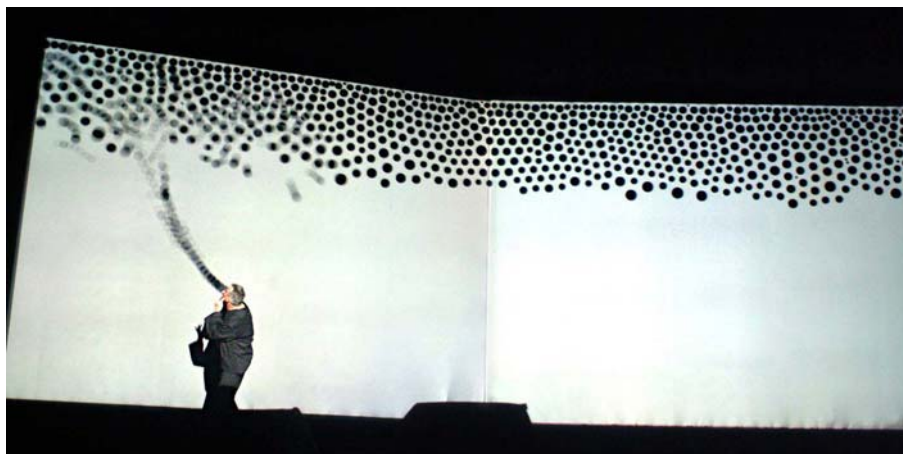


Fig. 3 Vocalist Jaap Blonk performing the *Messa di Voce* interactive software by Golan Levin and Zachary Lieberman (2003)

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 (Levin and Lieberman) (Fig. 3).

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 50 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



Fig. 4 Rafael Lozano-Hemmer's *Standards and Double Standards* (2004)

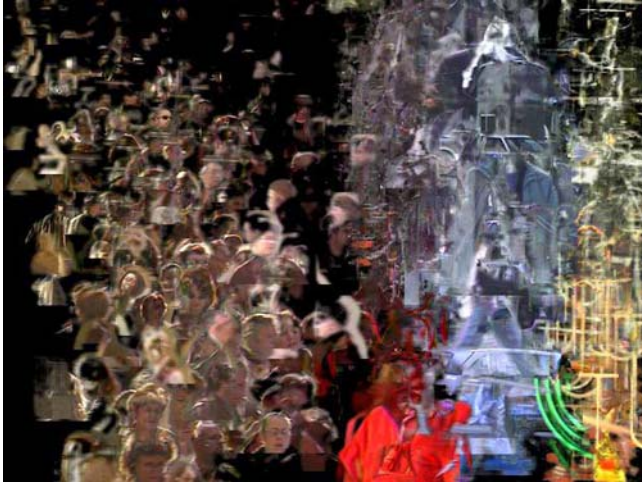


Fig. 5 A composite image assembled by David Rokeby's *Sorting Daemon* (2003)

“turns a condition of pure surveillance into an ‘absent crowd’ using a fetish of paternal authority: the belt” (Lozano-Hemmer) (Fig. 4).

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 (Rokeby) (Fig. 5).

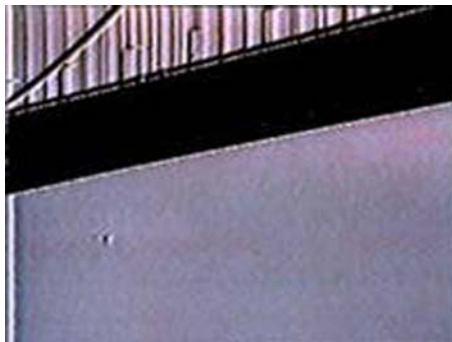


Fig. 6 *Suicide Box* by the Bureau of Inverse Technology (1996)

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 recognised 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.” (Bureau of Inverse Technology). Elsewhere, Jeremijenko has explained that “the idea was to track a tragic social phenomenon which was not being counted—that is, doesn’t count” [Shachtman]. The Suicide Box has met with considerable controversy, ranging from ethical questions about recording the suicides, to others disbelieving 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 (Fig. 6).



Figs. 7–8 Stills from *Cheese*, an installation by Christian Möller (2003)

Considerably less macabre is Christian Möller's clever *Cheese* installation (2003), which the artist developed in collaboration with 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 an hour and a half. 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" (Möller). 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 (Figs. 7–8).

As can be seen from the examples above, 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. And 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, and to movements of entire bodies. The examples above represent just a small selection of notable works in the field, and of ways in which people (and objects) have been tracked and dissected by video analysis. Other noteworthy artworks which 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.

3 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 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 this 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, etc. 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 below, are extremely simple to implement and help constitute a base of detection schemes from which sophisticated interactive systems may be built. (Complete implementations of these algorithms, written in the popular *Processing* flavor of Java, appear in code listings at the end of this article.)

3.1 Detecting motion (Code listing 1)

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 F1 is compared with its corresponding pixel in the subsequent frame F2. 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, in order 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 which is being measured).

3.2 Detecting presence (Code listing 2)

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 which belong to them from those which 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 which 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.

3.3 Detection through brightness thresholding (Code listing 3)

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 of interest are considerably darker than, 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 as foreground or background accordingly.

3.4 Simple object tracking (Code listing 4)

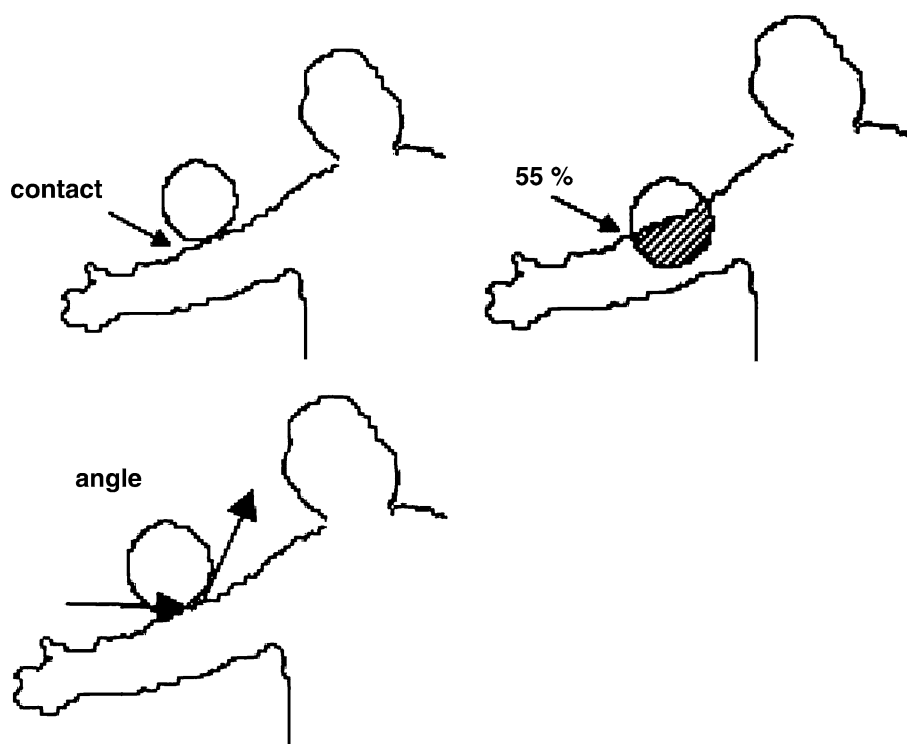
A rudimentary scheme for object tracking, ideal for tracking the location of a single illuminated point (such as a flashlight), finds the location of the single brightest pixel in every fresh frame of video. In this algorithm, the brightness of each pixel in the incoming video frame is compared with 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.

3.5 Basic interactions

Once a person's body pixels have been located (through the aid of techniques like background subtraction and/or brightness thresholding), this information can be used as the basis for graphical responses in interactive systems. In a 2003 Master's thesis, *Unencumbered Full Body Interaction in Video Games*, Jonah Warren presents an elegant vocabulary of various essential interaction techniques which can use this kind of body-pixel data. These schema are useful in "mirror-like" contexts, such as Myron Krueger's Videoplace, or a game like the PlayStation *Eye-Toy*, in which the participant can observe his own image or silhouette composited into a virtual scene (Figs. 9–11).

Three of the interactions Warren identifies and explains are the *Contact* interaction, which can trigger an event when a user's digital silhouette comes into contact with a graphic object; the *Overlap* interaction, which is a continuous metric based on the percentage of pixels shared between a user's silhouette and a graphic object; and the *Reflect* interaction, which computes the angle of reflection when a moving object strikes the user's silhouette (and deflects the object appropriately). Documentation of several charming games which make use of these interactions can be found in Warren's site. As Warren explains it, the implementation of these interactions requires little more than counting pixels (Warren).

Naturally, 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 which 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



Figs. 9–11 *Contact, Overlap, Reflect*: Examples of “Unencumbered Full-Body Interactions” identified by Jonah Warren (Warren)

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, and/or ready-to-use, executable implementations of all of these techniques can be found on the Internet in excellent resources like Daniel Huber’s *Computer Vision Homepage* (Huber), Robert Fisher’s *HIPR (Hypermedia Image Processing Reference)* (Fisher), or in the software toolkits discussed in Sect. 5, below.

4 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/or to select software techniques which are best 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 which 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, and 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 exists for optimizing physical conditions in order to enhance the robustness, accuracy and effectiveness of computer vision software. Most are geared towards 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. IR, 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 towards 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 gaging three-dimensional objects, or objects whose distance from the lens is not known precisely" (Melles Griot). Likewise, polarizing filters offer a simple, non-intrusive solution to another common problem in video systems, namely glare from reflective surfaces. And a wide range of video cameras is available, optimized for conditions like high-resolution capture, high-frame-rate capture, short exposure times, dim light, ultraviolet light, or thermal imaging. Clearly, 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 highest-quality hardware cannot help us find meaning where there is none, or track an object which 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 which 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 IR retroreflection, *and* they also indicate a person’s direction of gaze. Or in the dramatic case of Natalie Jeremijenko’s *Suicide Box*, discussed earlier: vertical motion in the video frame is easy to find through simple frame-differencing, *and* (in a specific context) it can be a stark indicator of a tragic event. In judging which features in the world are most profitably selected for analysis by computer vision, we will do well to select those graphical facts about the world which not only are easy to detect, but also simplify its semantic understanding.

5 Computer vision in multimedia authoring tools

The last decade has witnessed a profound transformation in the ease-of-use of software authoring tools for art and design students, and for novice programmers generally. While multimedia authoring environments are now commonly used to create interactive experiences for the World Wide Web, it is now equally common that these tools are used to create art installations, performances, commercial kiosks, and interactive industrial design prototypes. With the gradual incorporation of live video cameras into the palette of available computer inputs, the demand for straightforward computer vision capabilities has grown as well.

It can be an especially rewarding experience to implement machine vision techniques directly from first principles, using code such as the examples provided in this article. 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. *Processing* is one such environment, which, through an abundance of graphical capabilities, is extremely well suited to the electronic arts and visual design communities. Used worldwide by students, artists, designers, architects, and researchers for learning, prototyping, and production, Processing obtains live video through a QuickTime-based interface, and allows for fast manipulations of pixel buffers with a Java-based scripting language (Fry). The examples which appear in this article are written in Processing code.

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 immediately obtained from readymade, “off-the-shelf” solutions. Some of the most popular machine vision toolkits take the form of “plug-ins” or extension libraries for commercial authoring environments geared towards 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.

Aficionados of Macromedia's popular *Director* software, for example, can choose vision plug-ins ("Xtras") such as Danny Rozin's *TrackThemColors*, and Joshua Nimoy's *Myron* (named in honor of Myron Krueger). Rozin's inexpensive plug-in can track multiple objects in the video according to their chroma or brightness (Rozin). Nimoy's newer toolkit, which is freeware and open source, provides more detailed data about the tracked objects in the scene, such as their bounding quads and contour pixels (Nimoy). Through *Director*, the features detected by these Xtras can be linked to the control of sound playback, 2D and 3D graphics, text, and serial communications.

Many vision plug-ins have been developed for *Max/MSP/Jitter*, a visual programming environment which 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 (Cycling'74). 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 which 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 (Pelletier).

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 (STEIM). 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 (Camurri). *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 interpreted 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 [Intel]. The OpenCV library, by contrast, is a free, open-source toolkit with nearly similar capabilities, and a tighter focus on

commonplace computer vision tasks (Davies). The capabilities of these tools, as well as all of those mentioned above, are continually evolving.

5.1 An example: *LimboTime*

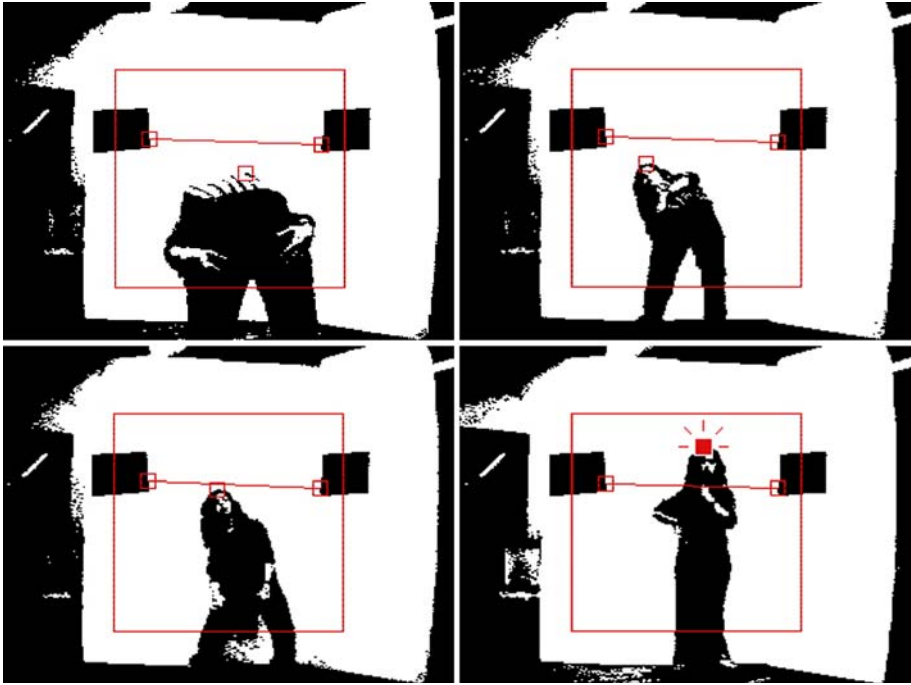
In October 2004, I conducted a workshop in machine vision for young artists and designers at the Benetton Fabbrica center in Treviso, Italy. The first day of the workshop covered the art-historical uses of computer vision, presented in Sect. 2, and the design considerations discussed in Sects. 3, 4; on the second day, the participants broke into small groups and were charged with the task of designing and programming a vision-based system “from scratch”. *Processing* was used as the development environment; the workshop participants were, for the most part, either novice programmers or intermediate-level programmers with little exposure to machine vision techniques.

LimboTime is an example of an interactive game which emerged from this workshop. In *LimboTime*, two players stand at opposite sides of the field of vision of a regular webcam, with their fingers or arms extended towards each other. The game system locates their extended fingers, and connects a horizontal line between them. This line continually connects these points, even if the players move their hands. A third player then tries to “limbo” underneath the imaginary line created by the first two players. The application tracks the vertical position of the “limboer” relative to the imaginary line; if the limboer goes above the limbo line, then the system sounds an alarm and the limboer is retired. If the limboer can pass completely under the line, however, then her companions lower their line-generating hands somewhat, and the process begins again (Figs. 12–15).

LimboTime is a simple game, conceived and implemented in a single afternoon. Its implementation grew organically from its creators’ discovery of a wall-size sheet of white Foamcore in the scrap closet. Realizing that they possessed an ideal environment for brightness-based thresholding, they used this technique in order to locate the games’ three players against their background. Detecting the players’ hands and heads was then accomplished with simple heuristics, e.g. the limboer’s head is the topmost point of the middle (non-side-touching) blob of black pixels. *LimboTime* was just one of the many interesting applications developed in the Fabbrica workshop; other examples included a system to detect and record the arrivals and departures of birds in a nearby tree, and a system which allowed its creators to “paint” an image by waving their glowing mobile phones around a dark room.

6 Conclusion

Computer vision algorithms are increasingly used in interactive and other computer-based artworks to track people’s activities. Techniques exist which 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 require advanced understandings of image processing and statistics, a number



Figs. 12–15 Stills captured from the *LimboTime* game system developed by workshop participants at the Fabbrica research center in Treviso, Italy. The participant shown attempts to pass below an imaginary line drawn between the two black rectangles. If she crosses above the line, the game rings an alarm. The black rectangles are temporary substitutes for the extended hands of two other participants (not shown)

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 additionally available which 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 which 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.

6.1 Code listing 1: frame differencing

```
// Processing code for detecting and quantifying
// the amount of movement in the video frame
// using a simple Frame-Differencing technique.

int video_width = 320;
int video_height = 240;
int num_pixels = (video_width * video_height);
int previous_frame[];

//-----
void setup(){
    size(320, 240);

    previous_frame = new int [num_pixels];    // Allocate memory for storing the previous
    for (int i=0; i<num_pixels; i++){          // frame of video, and initialize this buffer
        previous_frame[i] = 0;                 // with blank values.
    }
    beginVideo(video_width, video_height, 30);
}

//-----
void loop(){

    int curr_color, prev_color;                // Declare variables to store a pixel's color.
    float curr_r, curr_g, curr_b;             // Declare variables to hold current color values.
    float prev_r, prev_g, prev_b;             // Declare variables to hold previous color values.
    float diff_r, diff_g, diff_b;             // Declare variables to hold computed differences.
    float movement_sum = 0;                   // This stores the amount of movement in this frame.

    for (int i=0; i<num_pixels; i++){          // For each pixel in the video frame,
        curr_color = video.pixels[i];          // Fetch the current color in that location,
        prev_color = previous_frame[i];         // and also the previous color from that spot.

        curr_r = red    (curr_color);          // Extract the red, green, and blue components
        curr_g = green  (curr_color);          // of the current pixel's color.
        curr_b = blue   (curr_color);

        prev_r = red    (prev_color);          // Extract the red, green, and blue components
        prev_g = green  (prev_color);          // of the previous pixel's color.
        prev_b = blue   (prev_color);

        diff_r = abs (curr_r - prev_r);        // Compute the difference of the red values.
        diff_g = abs (curr_g - prev_g);        // Compute the difference of the green values.
        diff_b = abs (curr_b - prev_b);        // Compute the difference of the blue values.

        movement_sum += diff_r+diff_g+diff_b;  // Add these differences to the running tally.
        pixels[i] = color(diff_r,diff_g,diff_b); // Render the difference image to the screen.
        previous_frame[i] = curr_color;        // Swap the current information into the previous.
    }

    println(movement_sum);                    // Print out the total amount of movement
}
```

6.2 Code listing 2: background subtraction

```
// Processing code for detecting the presence of people and objects in the frame using a
// simple Background-Subtraction technique. Initialize the background by pressing a key.

int video_width = 320;
int video_height = 240;
int num_pixels = (video_width * video_height);
int background_img[];

//-----
void setup(){
  size(320, 240);
  beginVideo(video_width, video_height, 30);

  background_img = new int [num_pixels]; // Allocate memory for storing the background
  for (int i=0; i<num_pixels; i++){ // image, and initialize this buffer
    background_img[i] = 0; // with blank values.
  }
}

//-----
void keyPressed(){
  // When a key is pressed, capture the background
  for (int i=0; i<num_pixels; i++){ // image into the background_img buffer, by copying
    background_img[i] = video.pixels[i]; // each of the current frame's pixels into it.
  }
}

//-----
void loop(){
  int curr_color, bkgd_color; // Declare variables to store a pixel's color.
  float curr_r, curr_g, curr_b; // Declare variables to hold current color values.
  float bkgd_r, bkgd_g, bkgd_b; // Declare variables to hold previous color values.
  float diff_r, diff_g, diff_b; // Declare variables to hold computed differences.

  float presence_sum = 0; // This stores, for the current frame, the difference
  // between the current frame and the stored background

  for (int i=0; i<num_pixels; i++){ // For each pixel in the video frame,
    curr_color = video.pixels[i]; // Fetch the current color in that location,
    bkgd_color = background_img[i]; // and also the color of the background in that spot.

    curr_r = red (curr_color); // Extract the red, green, and blue components
    curr_g = green (curr_color); // of the current pixel's color.
    curr_b = blue (curr_color);

    bkgd_r = red (bkgd_color); // Extract the red, green, and blue components
    bkgd_g = green (bkgd_color); // of the background pixel's color.
    bkgd_b = blue (bkgd_color);

    diff_r = abs (curr_r - bkgd_r); // Compute the difference of the red values.
    diff_g = abs (curr_g - bkgd_g); // Compute the difference of the green values.
    diff_b = abs (curr_b - bkgd_b); // Compute the difference of the blue values.

    presence_sum += diff_r+diff_g+diff_b; // Add these differences to the running tally.
    pixels[i] = color(diff_r,diff_g,diff_b); // Render the difference image to the screen.
  }

  println (presence_sum); // Print out the total amount of movement
}
```

6.3 Code listing 3: brightness thresholding

```
// Processing code for determining whether a test location (such as the cursor)
// is contained within the silhouette of a dark object.

int video_width = 320;
int video_height = 240;
int num_pixels = (video_width * video_height);

//-----

void setup(){
    size(320, 240);
    beginVideo(video_width, video_height, 30);
}

//-----

void loop(){

    int black = color(0,0,0);                // Declare some constants for colors.
    int white = color(255,255,255);
    int threshold = 127;

    int pix_val;                             // Declare variables to store a pixel's color.
    float pix_bri;

    // Split the image into dark and light areas:
    for (int i=0; i<num_pixels; i++){        // For each pixel in the video frame,
        pix_val = video.pixels[i];           // fetch the current color in that location,
        pix_bri = brightness (pix_val);      // and compute the brightness of that pixel.

        if (pix_bri > threshold){            // Now "binarize" the video into black-and-white,
            pixels[i] = white;               // depending on whether each pixel is brighter
        } else {                             // or darker than a threshold value.
            pixels[i] = black;
        }
    }

    // Test a location to see where it is contained.
    int test_val = get (mouseX, mouseY);     // Fetch the pixel at the test location (the cursor),
    float test_bri = brightness (test_val);  // and compute its brightness.
    if (test_bri > threshold){               // If the test location is brighter than threshold,
        fill (255,0,0);                     // draw a RED ellipse at the test location.
        ellipse (mouseX-10, mouseY-10, 20,20);
    } else {                                // Otherwise,
        fill (0,255,0);                     // draw a GREEN ellipse at the test location.
        ellipse (mouseX-10, mouseY-10, 20,20);
    }
}
```

6.4 Code listing 4: simple object tracking

```
// Processing code for tracking the brightest pixel in a live video signal

int video_width  = 320;

int video_height = 240;

//-----

void setup(){

    size(320, 240);

    beginVideo(video_width, video_height, 30);

}

//-----

void loop(){

    image(video, 0, 0, width, height);           // Draw the webcam video onto the screen.

                                                    // Declare some numbers to be computed later:

    int brightest_x = 0;                          // the x-coordinate of the brightest video pixel
    int brightest_y = 0;                          // the y-coordinate of the brightest video pixel
    float brightest_val = 0;                      // the brightness of the brightest video pixel

                                                    // Now search for the brightest pixel:

    for (int y=0; y<video_height; y++){          // For each row of pixels in the video image,
        for (int x=0; x<video_width; x++){        // and for each pixel in the y'th row,
            int index = y*video_width + x;        // compute each pixel's index in the video,
            int pix_val = video.pixels[index];    // fetch the color stored in that pixel,
            float pix_bri = brightness(pix_val);  // and determine the brightness of that pixel.
            if (pix_bri > brightest_val){          // If that value is brighter than any previous,
                brightest_val = pix_bri;           // then store the brightness of that pixel,
                brightest_y = y;                   // as well as its (x,y) location.
                brightest_x = x;
            }
        }
    }

    fill(255,0,0);                                // Set the fill color to red, and then
    ellipse( brightest_x-10, brightest_y-10, 20,20); // draw a circle at the brightest pixel.

}
```

7 About the author

Golan Levin is an artist, performer and engineer interested in developing artifacts and events which explore supple new modes of reactive expression. His work focuses on the design of systems for the creation, manipulation and performance of simultaneous image and sound, as part of a more general inquiry into the formal language of interactivity, and of nonverbal communications

protocols in cybernetic systems. Through performances, digital artifacts, and virtual environments, often created with a variety of collaborators, Levin applies creative twists to digital technologies that highlight our relationship with machines, make visible our ways of interacting with each other, and explore the intersection of abstract communication and interactivity. Levin is assistant professor of Electronic Time-Based Art at Carnegie Mellon University, Pittsburgh.

Acknowledgments The author is grateful to Andy Cameron of Fabbrica Center, Treviso, for organizing the workshop which led to this article; Jonathan Harris, Francesca Granato and Juan Ospina for their *LimboTime* project and its documentation; David Rokeby and Rafael Lozano-Hemmer for their kind permissions; and Casey Reas, editor of this volume, who made this article possible in numerous different ways.

References

- Camurri, Antonio et al Eyesweb Vision-oriented software development environment. Laboratorio di Informatica Musicale, University of Genoa, Italy <http://www.eyesweb.org/>
- Cycling'74 Inc. Max/MSP/Jitter. Graphic software development environment. <http://www.cycling74.com/>
- Davies, Bob et al OpenCV Open-source computer vision library. <http://www.sourceforge.net/projects/opencvlibrary/>
- Fry, Ben and Reas, Casey Processing. Software development environment. <http://www.processing.org/>
- Nimoy, Joshua. Myron Xtra (plug-in) for macromedia director and processing. <http://www.webcamxtra.sourceforge.net/>
- Pelletier, Jean-Marc CV.Jit. Extension library for Max/MSP/Jitter. <http://www.iamas.ac.jp/~jovan02/cv/>
- Rokeby, David SoftVNS Extension library for Max/MSP/Jitter. <http://www.homepage.mac.com/davidrokeby/softVNS.html>
- Rozin, Danny TrackThemColors Xtra (plug-in) for macromedia director. <http://www.smoothware.com/track.html>
- Singer, Eric Cyclops Extension library for Max/MSP/Jitter. <http://www.cycling74.com/products/cyclops.html>
- STEIM (studio for electro-instrumental music). BigEye. Video analysis software. <http://www.steim.org/steim/bigeye.html>
- Bechtel, William (2003) The cardinal mercier lectures at the Catholic University of Louvain: An exemplar neural mechanism: the brain's visual processing system. Ch. 2 p. 1 <http://www.mechanism.ucsd.edu/~bill/research/mercier/2ndlecture.pdf>
- Fisher Robert, et al The hypermedia image processing reference (HIPR) <http://www.homepages.inf.ed.ac.uk/rbf/HIPR2/index.html>
- Fisher Robert, et al CVonline: The evolving, distributed, non-proprietary, on-line compendium of computer vision. <http://www.homepages.inf.ed.ac.uk/rbf/CVonline/>
- Huber, Daniel et al The computer vision homepage. <http://www-2.cs.cmu.edu/~cil/vision.html>
- Krueger Myron (1991) Artificial Reality II. Addison-Wesley Professional
- Levin Golan, Lieberman Zachary (2003) Messa di Voce. Interactive installation, <http://www.tmema.org/messa>
- Levin, Golan and Lieberman, Zachary (2004) "In-Situ speech visualization in real-time interactive installation and performance." In: Proceedings of the 3rd international symposium on non-photorealistic animation and rendering, June 7-9, Annecy, France http://www.flong.com/writings/pdf/messa_NPAR_2004_150dpi.pdf
- Lozano-Hemmer, Rafael Standards and double standards. Interactive installation. <http://www.fundacion.telefonica.com/at/rh/eproyecto.html>
- Melles Griot Corporation. Machine vision lens fundamentals. <http://www.mellesgriot.com/pdf/pg11-19.pdf>
- Moeller Christian (2003) Cheese. Installation artwork, <http://www.christian-moeller.com/>

- Rokeby David (2003) Sorting daemon. Computer-based installation, <http://www.homepage.mac.com/davidrokeby/sorting.html>
- Shachtman, Noah "Tech and Art Mix at RNC Protest". Wired News, 8/27/2004 <http://www.wired.com/news/culture/0,1284,64720,00.html>
- Sparacino, Flavia (2001) "(Some) computer vision based interfaces for interactive art and entertainment installations". INTER_FACE Body Boundaries, ed. Emanuele Quinz, Anomalie, n.2, Paris, France, Anomos
http://www.sensingplaces.com/papers/Flavia_isea2000.pdf
- Warren Jonah (2003) Unencumbered full body interaction in video games. Master's Thesis, Parsons School of Design (Unpublished)
http://www.a.parsons.edu/~jonah/jonah_thesis.pdf