# The Diver Project: Interactive Digital Video Repurposing

**Roy Pea, Michael Mills, Joseph Rosen, and Kenneth Dauber**
*Stanford University*

**Wolfgang Effelsberg**
*University of Mannheim, Germany*

**Eric Hoffert**
*Versatility Software*

**The Digital Interactive Video Exploration and Reflection (Diver) system lets users create virtual pathways through existing video content using a virtual camera and an annotation window for commentary. Users can post their Dives to the WebDiver server system to generate active collaboration, further repurposing, and discussion.**

With the inexorable growth of low-cost consumer video electronics and broadband networking, communication-centric, symmetric video-point-of-view sharing will eventually supplant broadcast-centric and asymmetric approaches to video content and repurposing. In sheer volume, the use of existing video to share perspectives will likely far outstrip video repurposing for commercial purposes and for different device platforms.

Thus, automatic, real-time repurposing of content designed for one technical environment configuration (a device or network, for example) to fit other technical environments is but one, albeit important, facet of the content repurposing problem. Another fundamental issue in content repurposing is user-driven content reuse, prompted by users' desire to reinterpret content and to communicate and share their interpretations.

The Digital Interactive Video Exploration and Reflection (Diver) system provides facilities for dynamic time/space cropping and annotation of archival video to support digital video repurposing. Diver lets users create point-of-view video tours of original content in a way that supports sharing, collaboration, and knowledge building.

Diver works like a video camera: a user points, pans, tilts, and zooms with the virtual viewfinder to create an infinite number of perspectives on a video record. Although other digital video editing and effects tools offer similar repurposing features (Adobe Premiere's zoom transition effects and cropping and Adobe After Effects' side-by-side group clip compositing, for example), none offer a comparable intuitive and easy-to-use interface, and few can infinitely repurpose audio and video without forcing the user to wait through a rendering process.

## Diver overview

Much of our work on Diver involves recording human activities such as teachers and students in a classroom or research group meetings using a conventional or panoramic digital video camera. Such scenarios involve a broad range of complex human interactions of interest to researchers in the learning and social sciences, such as meeting behaviors or instructional discourse.

A video-analysis phase follows the video-recording phase.[1] Using a mouse, users can zoom, pan, and tilt a virtual camera on an overview of the source video. The virtual camera dynamically crops still image clips, or records multiframe video pathways through the panoramic video to create a Diver worksheet. A Dive consists of a set of reorderable panels (inspired by VideoNoter),[2] each containing a key frame or thumbnail representing a clip, and a text field that can contain an accompanying annotation. Source video can support an unlimited number of Dives.

After creating a Dive, a user can upload it to WebDiver, a Web site for interactive browsing, searching, and display of video clips and collaborative commentary on Dives. Diver automatically packages the Dive for WebDiver as an XML document with associated media files. The Dive's Web representation includes the key frames, annotations, and video clips. Colleagues can thus share Dives over the Internet.

Figure 1 gives a schematic representation of the recording, Diving, and Web-sharing phases.

## Architecture and implementation

Because 4:3-aspect-ratio digital videos are so

pervasive, they often serve as Diver's source material. However, it's difficult to use such video records in social science research. The videographer's choice of what to include (and what not to include) limits the usefulness of videos recorded with conventional digital cameras not only for the original researcher, who can't recover events relevant to the analysis occurring off camera, but also for other researchers who might be interested in different aspects of the filmed event.
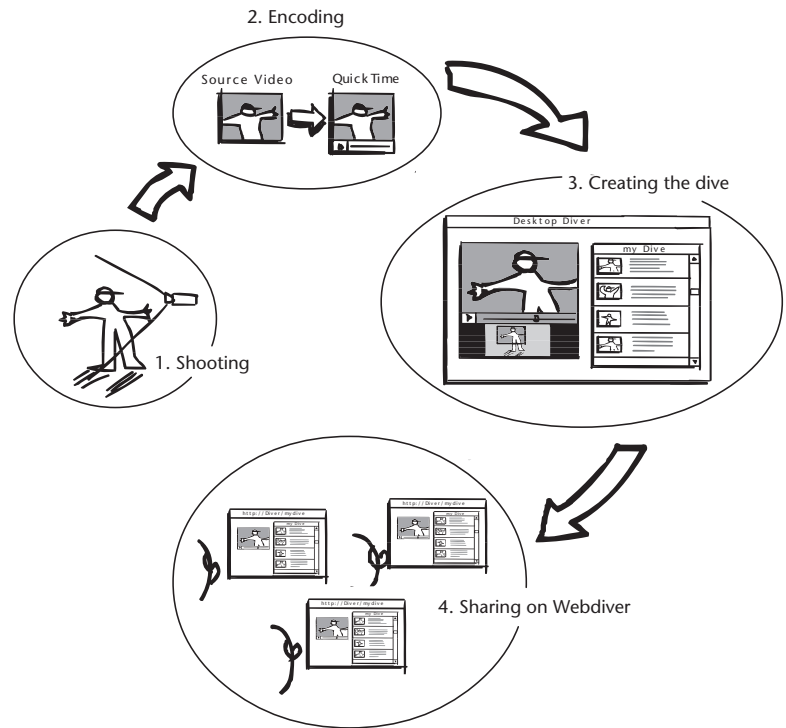
By supporting panoramic video, Diver avoids this problem. Panoramic video has been used in teleconferencing[3] and in surveillance, or as one component in an immersive simulation of an event, such as a football game or concert.[4] We argue that panoramic video is of particular interest in all settings where the video material's purpose and application isn't precisely defined or known at recording time.

## Video recording system

Image resolution is a major technical challenge in recording panoramic video. For a user to zoom into an area of the panoramic video and still see sufficient detail, the camera must capture the video at the highest possible resolution.

In principle, a single-lens system can record panoramic video: a spherical mirror, somewhat similar to a fisheye lens, collects the light from 360 degrees and redirects it onto a camera's charge-coupled device (CCD) chip. The system can then store, compress, and edit the camera's digital video stream like any other video stream. Such systems are relatively inexpensive and readily available. Examples include the BeHere camera (http://www.behere.com) and the camera used in Ricoh's Portable Meeting Recorder.[5] However, single-lens systems produce low-resolution images: the number of pixels on the CCD chip, designed for a small angle, must now suffice for surround vision—that is, a much wider angle. Consequently, a user zooming in on an area in the video will see a low-quality image.

Several recent multilens systems—FullView (http://fullview.com), Fuji's FlyCam,[6] and Microsoft's RingCam,[7] for example—address this problem. They capture scenes using multiple lenses and CCDs with slightly overlapping fields of vision and stitch the frames together to construct panoramic images. However, these systems introduce a new challenge: stitching and dewarping multiple frames in the computer, perhaps in real time. Multilens systems running at full resolution generate a very high bit rate, often too

high to be transferred over a single peripheral component interconnect (PCI) bus.

Diver uses a FullView camera system with five mirrors and five cameras, and this system feeds the camera streams to a PC where FullView software stitches and dewarps the video streams for live preview (at a reduced frame rate and/or reduced resolution, and always without sound). The FullView system uses three parallel PCI buses on the PC processing the streams to achieve such a high data transfer rate.

Because of these limitations, we've developed a tape-based recording system capable of recording at the full resolution and full frame rate of NTSC-DV (digital video), which is desirable for repurposing panoramic video (that is, $720 \times 480 \times 5$ cameras, at ~30 fps). A system for capturing multichannel sound with a microphone array generates corresponding audio.

To ensure compatibility with 4:3 aspect-ratio video and panoramic video content at a variety of aspect ratios, Diver supports a general approach to video analysis, collaboration, and interaction, allowing for a diversity of aspect ratios and resolutions for recorded video material. This provides a flexible approach for recording, analyzing, and sharing a broad range of video material independently of the specifics of the captured video material.



*Figure 1. Overview of the Diver video exploration and reflection system.*
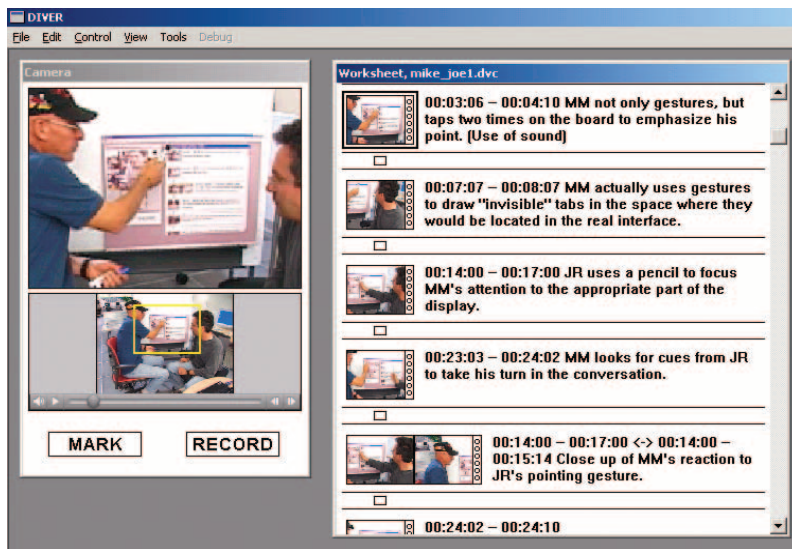
Figure 2. Diver user interface. The overview window (bottom left) shows the full video source. The magnified viewing window (upper left) shows a selected image from the scene. The annotation window, a.k.a. Dive worksheet (right) lets users comment on the frames or path movies they create.

## User interface

The Diver software consists of the Diver desktop application and the WebDiver server system. The desktop application is a native Windows application coded in C++ that builds on Apple's QuickTime tools and uses the QuickTime video format and QuickTime browser plug-in to display and control the video.

The design rationale for Diver's user interface is rooted in cognitive science and human–computer interaction research. Our goal for the tool is to augment the fundamental and social activity of human visual interpretation and communication about what's interesting or relevant to a given purpose—namely, looking at a complex scene, noticing events, and commenting on the focus of attention.

The Diver interface consists of three distinct regions, each corresponding to an element in the look-notice-comment cycle, as Figure 2 illustrates.

The *overview* window (look region) provides access to the original source material and provides standard video controllers.

The *magnified viewing* window (notice region) shows the part of the original scene that the user selects by dragging and resizing a rectangle overlay (the virtual camera) across the overview window, panning over and zooming into regions of interest.

This overview/magnified viewing interface aims to satisfy the well-known dual-awareness principle of human–computer user interface design: it lets the user zoom in on details of an information source while maintaining orientation and context.[8,9] Other researchers have explored related navigational interfaces in experimental[10–12] and commercial applications of panoramic video technology (FullView and BeHere, for example).

As a user drags the virtual camera rectangle across the video source overview, the viewing window gives a dynamically updated and magnified view of the circumscribed region. Two buttons make the viewing selection a flexible authoring tool. *Mark* takes a temporally tagged snapshot of the viewing window contents and automatically creates a new panel in the screen's third region: the Dive worksheet (comment region). In addition to a video thumbnail representing a marked frame or clip, a worksheet panel contains a text field in which users can describe the selection's noteworthy characteristics.

Pressing the *record* button creates a new panel. At the same time, moving the selection rectangle across the overview window creates a recording of the spatiotemporal selection shown in the viewing window. Thus the user can record a dynamic path through the original source video. Pressing the record button again ends the recording. The user can annotate the recorded pathway clip by entering text in the Dive panel.

In creating a unique point-of-view virtual tour of the original video source, the Diver user creates a Dive. The Dive can serve as a storytelling device for playing back both the marked locations and the dynamically cropped pathway recordings through the original video source and the records' annotations.

Panels contain the data elements constituting a Dive (see Figure 2):

- *static and dynamic time and space markers*—a static marker points to a single frame in the video (temporal coordinate) and the location of the virtual viewfinder in that frame (spatial coordinates); dynamic markers result from recording the virtual camera's time/space path as the user guides it through the source video.

- *thumbnail representations of time and space video markers*—thumbnail-sized images copied from the source video at the marked time and space.

- *annotation text*—the commentary associated with the marker (that is, the text the user typed into the panel).

- *time codes*—codes indicating the temporal point or range of video contained within a panel selection. Static markers take a single time code; dynamic markers take a range of codes.

- *time/space cropped video clips*—thumbnails in a panel that also symbolize portions of the source video. Double-clicking a thumbnail (or dragging and dropping it back into the viewing window) repositions the source video at the corresponding time/space coordinates.

**Full-resolution panoramic video recorder**

Diver's user interface supports panoramic video. As Figure 3 demonstrates, the user gets a 360-degree horizontal view of the unfolding scene while the video plays as a "peeled back" cylinder with a wide-aspect ratio.

To accommodate even higher resolution video capture than native FullView live preview and direct-to-disk video can achieve, we developed an alternative full-resolution video tape-based recorder. Our general-purpose Diver recorder system captures the FullView camera array output at full resolution ($720 \times 480 \times 5$) and full frame rate (~30 fps) to five DV tape recorders (VTRs). The Diver recorder system later transfers select video tape segments to a PC, where our software reassembles in postprocessing the multiple recordings into seamless full-resolution panoramic movies encoded in QuickTime and ready for Diving.

**Virtual path movies**

The Diver file format supports the efficient creation of Dives. As a user drags the virtual camera viewfinder across the source overview window, Diver records time (the video time) and space (a viewing region's position within the video's border) references, as Figure 4 shows. The virtual viewfinder in Figure 4 is rectangular, but it can also be circular or any selection shape. Marking parallel video regions of interest using multiple concurrent pathways with the appropriate user interface is also feasible.

What's recorded when the user presses the record button is not a new video; rather, it's the digital video collaboratory (.dvc) file information needed to replay the dynamically cropped pathway that the virtual camera traced through the original video frame sequence.

When a user plays back a .dvc file, Diver resets the current video time and space to the previ-
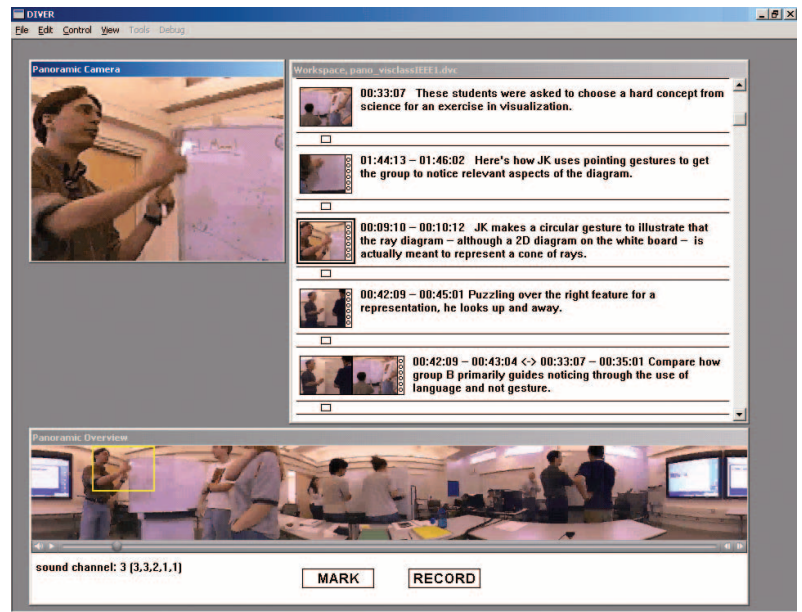
ously saved coordinates. Before displaying each frame of video in the virtual camera window, Diver scales the video's borders (in real time, on-the-fly) to match only the viewing area coordinates associated with the same frame (video time) in the recorded virtual camera path.

Using the .dvc list recording scheme instead of rendering new video clips has two significant advantages:

- Virtual video clips eliminate the generation of redundant video files, greatly reducing disk storage requirements.

- No rendering time means vastly improved performance. Users can instantly create and play back dynamic path videos without long video-rendering delays.

Dives are thus extremely lightweight glosses on original source content yet they can contain rich and compelling viewing experiences.
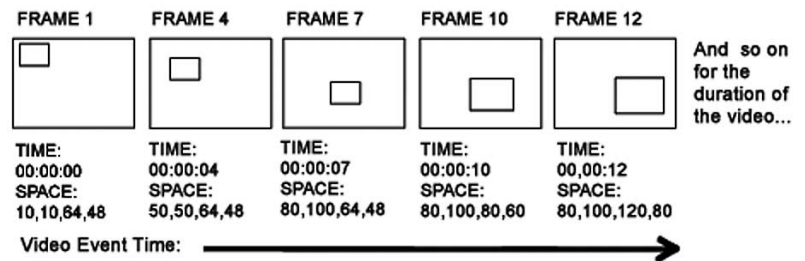


*Figure 4. Graphical representation of the .dvc file format.*

```
<dive id= owner=>
   <name></name>
   <annotation></annotation>
   <video type= id=>
      <video_file bytes=></video_file>
   </video>
   <panel id= order=>
      <annotation></annotation>
      <video type=>
        <time_in></time_in>
        <time_out></time_out>
        <source_video></source_video>
        <video_file bytes=></video_file>
        <thumbnail_file
           bytes=></thumbnail_file>
      </video>
   </panel>
</dive>
```

The structure and extensibility of the XML-based .dvc file format lets users reformat Dives into many display and interaction environments. For example, on a lightweight device supporting only text and image display, a user could extract and show just the text annotations and video image thumbnails.

**Infinite panoramic audio mixing**

The Diver panoramic recording system captures up to 10 independent audio channels, each mapped to an associated field of view captured by the FullView camera. (Near-term enhancements will incorporate 16 or more independent audio channels into the panoramic recording array.) Diver embeds these multichannel sound recordings as navigable QuickTime audio tracks in the postprocessed panoramic movies. When playing back the recordings, Diver can steer panoramic sound so only the tracks associated with the virtual viewfinder's current focus become audible. Thus, Diver users can easily repurpose audio recordings and dynamically create infinite virtual audio mixes.

**WebDiver server system**

Users can upload their Dives to a WebDiver server where others can share them, promoting dialogue and collaboration around video clips. Transforming content for the Web occurs in a series of client- and server-side automated tasks invisible to the user. The user chooses the Export to WebDiver menu option in the desktop Diver application, initiating text, image, and video-clip transformation processes.

The Diver system formats and transfers annotation text, time, and space markers and other metadata associated with each Dive (file names, file sizes, and so on) to the WebDiver server in an XML format. At the server, PHP scripts and a MySQL database store and retrieve Dive content for delivery to client browsers as dynamically generated and JavaScript-enhanced HTML pages. Diver transforms thumbnail images from video bitmaps to JPEG files.

Because QuickTime plug-ins don't support real-time playback of time/space cropped movies, accommodating dynamic playback of such movies in Web pages across multiple platforms and browsers would require a custom-coded plug-in for each of the most popular operating systems and browsers. As a possible cross-platform alternative, we experimented with a Diver Java applet architecture, but found that QuickTime for Java in its current implementation isn't robust enough to support time/space cropped movies. Rather than write our own custom plug-ins, we rendered all Dive video into standalone QuickTime video clips. Users can play the videos directly from the WebDiver Web site using standard QuickTime players.

During the video clip export process, Diver uses the .dvc file to trace the virtual camera's path through the time/space cropped video. It rebuilds the QuickTime clips frame by frame by copying time/space regions of the source video into the new standalone files. Diver also converts the original video source file to QuickTime in a compressed format of smaller geometry. These conversions constitute a technical video repurposing step in Diver.

Figure 5 shows the XML schema for a Dive.

When the desktop Diver finishes the transformation, the application automatically opens the user's Web browser and dynamically generates and displays a new Web page containing a file transfer and HTML form frameset. Submitting the form completes the Dive transmission (represented by the XML file containing Dive metadata, image thumbnails, and automatically generated video files) from the desktop to the Web server.

When a user later opens a Dive on the WebDiver site, shown in Figure 6, Diver reconstitutes the content as HTML and JavaScript, dynamically generating the Web pages from the media files and XML data. WebDiver pages in the browser look similar to desktop Dives and contain the same data and media elements.

In addition to the Dive's author, other WebDiver users can explore the full source video used to create it. Moreover, users can collaborate

on Dives by adding comments to the panels in the manner of a threaded discussion. WebDiver users can also perform keyword and metadata searches across all annotations.
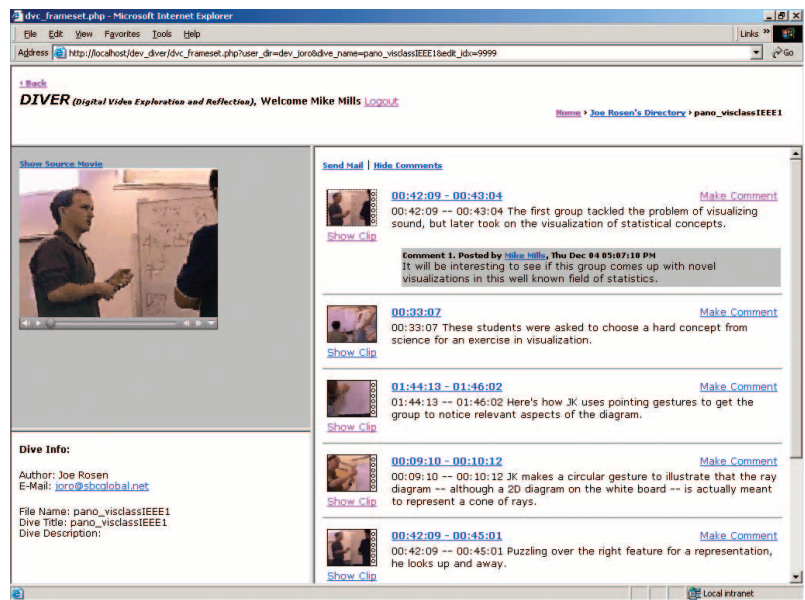
## Diver usage scenarios and implications

Multimedia records are an increasingly important type of data for researchers in many fields because they capture the complexity of "mind in context"—human interactions and behaviors in situations such as the classrooms and the home. Researchers today lack tools for sharing or discussing these data readily with fellow scholars and practitioners. Diver seeks to bridge this gap, enabling collaborative analysis of learning and teaching video records in a distributed community of researchers and practitioners.

For example, social scientists can select and annotate video clips of human behavior from the workplace, home, or learning settings, and share their analyses with distant colleagues. They can investigate areas in which video has proven useful, such as patterns of student participation in instructional discourse, uses of diagrams, small group collaboration, and teaching.[13-16] Used diagnostically, Diver can tease out differences in how novices and experts notice social, behavioral, and physical phenomena in video. Researchers can use Dive panels as structured templates for assessing trainees in various fields; they can direct users to particular Diver clips and ask them to answer specific questions. Local or remote experts can provide critiques as threaded comments in WebDiver. For training purposes, experts in a domain could construct Diver videos from Dives created to serve as think-aloud reflective commentaries, offering "guided noticing" help for novices.

More generally, Diver is relevant to any human activity domain with dynamic or static media in which community engagement in analysis and education using those media might bear fruit. For example, virtual classroom visits are increasingly used in teacher education for professional development (Teachscape.com, Lessonlab.com, and Teachfirst.com, for example).

To ensure that Diver reflects user needs, and to introduce Diver into relevant user communities, we recently held a design kitchen in the learning and social sciences departments at Stanford University. Faculty and graduate students used Diver to examine their own video data while the Diver team observed their use of the tool and recorded their reactions. The initial user reaction was positive, with several research



*Figure 6. A Dive from the WebDiver site, as it appears in a browser window.*

teams already incorporating Diver into their work with digital video data—for example, in studies of collaborative learning in mathematics, group discussions of text comprehension, uses of diagrams in physics education, and teacher preparation.

We're working on a pilot project with several California universities in which K–12 teachers use Diver to document and reflect on their teaching practices in terms of specific quality rubrics for planning, instruction, assessment, and reflection.

## Outlook

Although our current work focuses on video records in learning research and educational practices, Diver can aid collaborative analysis of a broad array of visual data records, including simulations, 2D and 3D animations, and static works of art, photography, and text. In addition to the social and behavioral sciences, substantive application areas include medical visualization, astronomic data or cosmological models, military satellite intelligence, and ethnology and animal behavior.

Diver-style user-centered video repurposing might also prove compelling for popular media with commercial application involving sports events, movies, television shows, and video gaming.

Future technical development includes possible enhancements to the interface to support simultaneous display of multiple Dives on the same source content, a more fluid two-way relation between desktop Diver and WebDiver, and

solutions to the current limitations on displaying and authoring time/space cropped videos in a browser context. These developments support the tool's fundamentally collaborative, communication-oriented nature. **MM**

## References

1. B. Jordan and A. Henderson, "Interaction Analysis: Foundations and Practice," *J. Learning Sciences*, vol. 4, 1995, pp. 39-103.
2. J. Roschelle, R.D. Pea, and R. Trigg, *VideoNoter: A Tool for Exploratory Video Analysis*, tech. report no. 17, Inst. for Research on Learning, Palo Alto, Calif., 1990.
3. D. Kimber, J. Foote, and S. Lertsithichai, "FlyAbout: Spatially Indexed Panoramic Video," *Proc. ACM Multimedia 2001*, ACM Press, 2001, pp. 339-347.
4. T. Pintaric, U. Neumann, and A. Rizzo, "Immersive Panoramic Video," *Proc. ACM Multimedia 2000*, ACM Press, 2000, pp. 493-494.
5. D. Lee et al., "Portable Meeting Recorder," *Proc. ACM Multimedia 2002*, ACM Press, 2002, pp. 493-502.
6. X. Sun et al., "Panoramic Video Capturing and Compressed Domain Virtual Camera Control," *Proc. ACM Multimedia 2001*, ACM Press, 2001, pp. 329-338.
7. R. Cutler et al., "Distributed Meetings: A Meeting Capture and Broadcasting System," *Proc. ACM Multimedia 2002*, ACM Press, 2002, pp. 503-512.
8. S.K. Card, J.D. Mackinlay, and B. Shneiderman, eds., *Readings in Information Visualization: Using Vision to Think,"* Morgan Kaufmann, 1999.
9. M. Mills et al., "A Magnifier Tool for Video Data," *Proc. Computer–Human Interactions*, ACM Press, 1992, pp. 93-98.
10. L. Teodosio and M.I. Mills, "Panoramic Overviews for Navigating Real-World Scenes," *Proc. ACM Multimedia 1993*, ACM Press, 1993, pp. 359-364.
11. B. Prihavec and F. Solina, "User Interface for Video Observation over the Internet," *J. Network and Computer Applications*, vol. 21, 1998, pp. 219-237.
12. Y. Rui, G. Gupta, and J.J. Cadiz, "Viewing Meetings Captured by an Omni-Directional Camera," tech. report MSR-TR-2000-97, Microsoft Corp., 2000.
13. B. Barron, "When Smart Groups Fail," *J. Learning Sciences*, vol. 12, 2003, pp. 307-359.
14. J. Frederiksen et al., "Video Portfolio Assessment: Creating a Framework for Viewing the Functions of Teaching," *Educational Assessment*, vol. 5, no. 4, 1998, pp. 225-297.
15. M. Lampert and D. Loewenberg-Ball, *Teaching, Multimedia and Mathematics: Investigations of Real Practice*, Teachers' College Press, 1998.
16. M. Ulewicz and A. Beatty, eds., *The Power of Video Technology in Int'l Comparative Research in Education,* Nat'l Research Council, Board on Int'l Comparative Studies in Education, Board on Testing and Assessment, Center for Education, Washington, D.C., 2001.

**Roy Pea** is a professor of education and learning sciences at Stanford University, director of the Stanford Center for Innovations in Learning (SCIL), and cofounder of Teachscape.com. His research interests include applying video analysis technology to studies in the learning sciences and technology design. Pea has a DPhil in developmental psychology from Oxford University as a Rhodes Scholar. He is president-elect of the International Society of the Learning Sciences, and a fellow of the National Academy of Education, the American Psychological Society, and the Center for Advanced Study in the Behavioral Sciences.

**Michael Mills** is SCIL's design director and a cognitive scientist. His research interests include interface, product and interface design, user studies, and teaching. Mills has a PhD in commu-

nication with a cognitive science specialization from McGill University. As principal scientist at Apple Computer, he was instrumental in the development of QuickTime and QuicktimeVR, holds 6 interface design patents in digital video, and has authored many articles on interface design.

**Joseph Rosen** is SCIL's senior software engineer. His research interests include software programming, interface design, digital interactive video, computer graphics, and animation. Rosen has a master's degree in interactive telecommunications from New York University. He also contributed to early QuickTime research at Apple Computer in the late 1980s.

**Kenneth Dauber** is SCIL's acting deputy research director. His research interests include sociology of culture and science. Dauber has a PhD in sociology from University of Arizona. He was a software architect and director of learning technologies at Unext.com, an Internet distance education company.

**Wolfgang Effelsberg** is a professor of computer science at the University of Mannheim, with research and teaching interests in multimedia technology and computer networks. Effelsberg has a Dr-Ing degree in computer science from the Technical University of Darmstadt, Germany. He is a member of and frequent contributor to the publications of the ACM, the IEEE and Gesellschaft für Informatik.

**Eric Hoffert** is the chief executive officer of Versatility Software. His research interests include digital media, collaboration, and communications. Hoffert has MS and BS degrees in computer science from New York University. He holds 11 patents and has published many papers for the IEEE and the ACM.

Readers may contact Pea at Stanford University, Stanford Center for Innovations in Learning, Wallenberg Hall (Room 232), 450 Serra Mall, Stanford CA 94305-2055, roypea@stanford.edu.