

An Overview of Passive Vision Techniques

Steven M. Seitz
The Robotics Institute
Carnegie Mellon University
seitz@cs.cmu.edu, <http://www.cs.cmu.edu/~seitz>

This section surveys *passive* optical techniques for recovering scene shape and reflectance characteristics from images. The objective is to acquire images of a scene observed from different viewpoints and possibly different illuminations, and from these images to compute scene shape and reflectance at every surface point. Given estimates of shape we can fabricate duplicate 3D models using techniques like stereo lithography. The ability to capture reflectance data, in addition to shape, allows us to create graphical representations that can be composited into new environments, reilluminated, and rendered from new viewpoints (see SIGGRAPH 99 Course #39 on Image-Based Modeling and Rendering for more on this topic).

We use the term *passive* sensing to refer to the measurement of visible radiation that is already present in the scene, in contrast to *active* techniques that project light into the scene. Although active sensing can facilitate the computation of scene structure, as discussed in sections 6-9 of these course notes, active approaches are not always feasible, especially for modeling distant or fast-moving objects. In addition, current active techniques like laser range scanning tend to be more expensive, slower, and more intrusive than their passive counterparts. However, the best active methods generally produce more accurate reconstructions than is possible using passive techniques. Note that the term *active vision* has also been used to refer to methods in which the camera is controlled purposively to facilitate reconstruction and other tasks—this is not the meaning intended here.

More information on passive vision techniques and related topics can be found in Section 5 (Paul Debevec's notes).

Visual Cues

Over the years, researchers in human perception and computer vision have identified numerous cues for sensing shape and reflectance properties and a wide range of computation strategies have emerged for exploiting these cues for 3D inference. Practical vision algorithms have been developed to measure shape from:

- Texture
- Shading
- Focus
- Parallax
- Long-range motion
- Reflection
- Shadows
- Symmetry
- Inter-reflection
- Polarization

While the above list is by no means complete, it demonstrates a wealth of information from which 3D shape and reflectance can be computed. Nevertheless, scene reconstruction is an inverse problem and generally does not admit a unique solution, i.e., it is *ill-posed* [1]. Consequently, additional assumptions and heuristics are generally needed to make the problem tractable. Some prominent examples include:

- Ideal reflectance – surfaces in the scene are often assumed to satisfy ideal reflectance models. For instance, stereo and shape-from-shading techniques generally assume a perfect Lambertian (isotropic) reflection model with no transparency. Consequently these techniques perform poorly in the presence of specularities and other deviations from the model.
- Smoothness – imposing smoothness, or *regularization* functionals is a very common method for making ill-posed inverse problems well-posed [1]. Choosing the reconstruction that is *smoothest* yields a better-conditioned problem but has its own pitfalls, for example the tendency to smooth over sharp edges or miss thin structures in the scene.
- Ideal projection – simplified projection models like orthographic projection and ideal pinhole projection are used to make the reconstruction equations more tractable. Consequently, techniques that use these approximations pay a penalty in terms of accuracy and are not well-suited for applications that demand high-accuracy surface measurements.

Passive Techniques

This course will consider four prominent approaches to passive 3D Photography: stereo, structure from motion, shape from shading, and photometric stereo.

Stereo

When a point is imaged from two different viewpoints, its image projection undergoes a displacement from its position in the first image to that in the second image. The amount of displacement, alternatively called *parallax* or *optical flow*, is inversely proportional to distance and may therefore be used to compute 3D geometry. Given a correspondence between imaged points from two known viewpoints, it is possible to compute depth by triangulation – intersect the rays from each optical center through the point's projection on the image plane. The problem of establishing correspondence is a fundamental difficulty and is the subject of a large body of literature on stereo vision. One prominent approach is to correlate pixels of similar intensities in two images, using an assumption that each scene point reflects the same intensity of light in the two views. The included paper by Okutomi and Kanade [2] extends this correlation approach to two three or more images and demonstrates that using several cameras at different camera separations, or *baselines* yields a significant improvement in reconstruction accuracy.

Camera Calibration

The accuracy of stereo techniques depends critically on having precise knowledge of camera position, orientation, and internal parameters, i.e., focal length, aspect ratio, principle point, and distortion contributed by non-ideal lens optics. Therefore, calibrating these parameters is a key part of the success of any stereo system. The first paper included in this section, by Tsai [3], describes one of the most widely used algorithms for calibrating cameras. Tsai's approach images a calibration object with known geometry and solves for camera parameters and a radial distortion coefficient given the image projections of several points on the calibration object.

Structure from Motion

Rather than image a scene with two or more cameras, an alternative approach is to acquire images from a single moving camera and to reconstruct a 3D model from the resulting video sequence. An additional challenge in this approach is that the camera path must be estimated since precise calibration of freely moving cameras is extremely difficult (although see [4]). This problem of recovering scene geometry from the motion of points in the image plane of a moving camera is called *structure from motion*. As one of the classical problems in computer vision, there are numerous structure from motion algorithms.

Included in this volume is one such technique developed by Tomasi and Kanade [5]. Their elegant *factorization* method assembles point measurements in a measurement matrix that is factored into the product of a *motion* and *shape* matrix using singular value decomposition. This approach is based on an orthographic projection model.

Shape from Shading

Shape from shading is one of the simplest problems to state and one of the most complicated to solve: given a single intensity image of a smooth curved object, how can the shape of the object be recovered [6]? This problem and its solution was pioneered in the vision community by Berthold Horn in his 1970 doctoral dissertation [7, 8]. Given the intensity of a point in the image and a known directional light source, Lambert's law ($I = kNL$, where I is the intensity, k the reflectance, N the unit normal, and L the light source direction) yields a one-parameter family of solutions for the surface normal. Additional constraints are needed to make the problem well-posed—it is generally solved by assuming similarity of surface reflectance and orientation at nearby points.

Photometric Stereo

The difficulties with shape from shading may be mitigated by acquiring two or more images of the object under different illuminations. This is precisely the approach of photometric stereo, the subject of the third paper included in these course notes [9]. Each image provides one constraint on the normal, and therefore two images are sufficient to recover the normal up to a small number of possible solutions, and three images yield a unique solution for each image pixel. Photometric stereo enables relaxing the strong smoothness conditions imposed by classical shape from shading approaches, and therefore yields more reliable shape estimates.

References

- [1] T. Poggio, V. Torre, and C. Koch, "Computational vision and regularization theory," *Nature*, vol. 317, pp. 314–319, 1985.
- [2] M. Okutomi and T. Kanade, "A multiple-baseline stereo," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 15, no. 4, pp. 353–363, 1993.
- [3] R. Y. Tsai, "A versatile camera calibration technique for high-accuracy 3D machine vision metrology using off-the-shelf cameras and lenses," *IEEE Trans. Robotics and Automation*, vol. 3, no. 4, pp. 323–344, 1987.
- [4] M. Pollefeys, R. Koch, and L. V. Gool, "Self-calibration and metric reconstruction in spite of varying unknown internal camera parameters," in *Proc. Sixth Int. Conf. on Computer Vision*, pp. 90–91, 1998.
- [5] C. Tomasi and T. Kanade, "Shape and motion from image streams under orthography: A factorization method," *Int. J. of Computer Vision*, vol. 9, no. 2, pp. 137–154, 1992.
- [6] L. B. Wolff, S. A. Shafer, and G. E. Healey, eds., *Shape Recovery*. Physics-Based Vision: Principles and Practice, Boston, MA: Jones and Bartlett Publishers, 1992.
- [7] B. K. P. Horn, *Shape from Shading: A Method for Obtaining the Shape of a Smooth Opaque Object from One View*. PhD thesis, Massachusetts Institute Of Technology, Cambridge, MA, 1970.
- [8] B. K. P. Horn and M. Brooks, *Shape from Shading*. Cambridge, MA: MIT Press, 1989.
- [9] R. J. Woodham, "Photometric method for determining surface orientation from multiple images," *Journal of Optical Engineering*, vol. 19, no. 1, pp. 138–144, 1980.