

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
ARTIFICIAL INTELLIGENCE LABORATORY
and
CENTER FOR BIOLOGICAL INFORMATION PROCESSING
WHITAKER COLLEGE

A.I. Memo No. 1297
C.B.I.P. Memo No. 64

June 1991

Recovering Heading for Visually-Guided Navigation

Ellen C. Hildreth

ABSTRACT: We present a model for recovering the direction of heading of an observer who is moving relative to a scene that may contain self-moving objects. The model builds upon an algorithm proposed by Rieger and Lawton (1985), which is based on earlier work by Longuet-Higgins and Prazdny (1981). The algorithm uses velocity differences computed in regions of high depth variation to estimate the location of the *focus of expansion*, which indicates the observer's heading direction. We relate the behavior of the proposed model to psychophysical observations regarding the ability of human observers to judge their heading direction, and show how the model can cope with self-moving objects in the environment. We also discuss this model in the broader context of a navigational system that performs tasks requiring rapid sensing and response through the interaction of simple task-specific routines.

©Massachusetts Institute of Technology (1991)

ACKNOWLEDGEMENTS: This paper describes research done at the Artificial Intelligence Laboratory and the Center for Biological Information Processing at the Massachusetts Institute of Technology. Support for the A. I. Laboratory's research is provided in part by the Advanced Research Projects Agency of the Department of Defense under Office of Naval Research contract N00014-85-K-0124. The Center's support is provided in part by the Office of Naval Research, Cognitive and Neural Sciences Division, the National Science Foundation (IRI-8719394 and IRI-8657824) and the McDonnell Foundation.

INTRODUCTION

Relative movement in the changing visual image provides a primary cue to the three-dimensional (3-D) structure and motion of object surfaces, and the movement of the observer relative to the scene, allowing biological systems to navigate quickly and efficiently through the environment. The range of tasks that use relative motion information imposes different demands on the speed, precision and completeness with which image motion must be measured and analyzed. Some tasks require precise 3-D models of the structure and motion of objects in the environment, and careful planning of observer motions. Examples include high-performance navigation tasks, such as negotiating through narrow channels or walking a tightrope; and fine manipulation tasks, such as assembling a model or threading a needle. Other tasks, however, require rapid sensing of rough environmental layout, followed by quick reflexive responses of the observer. Examples of this latter type include high-speed navigation tasks, such as racing through a cluttered environment toward a desired target; and tasks posed by fast-paced sports, such as dodging objects or other players, making quick moves to intercept a ball, or rapidly adjusting posture to maintain balance.

For tasks that require rapid sensing and response, there is often no time to construct elaborate 3-D models of the world. This is especially true for biological systems, which must rely on neural hardware that is very slow compared to today's high-speed computer hardware. Consideration of the demands of such tasks suggests that the human visual system may use specialized routines for performing reflexive actions in response to rapid changes in the environment. These routines may use only partial or qualitative information about image motion that is most critical to the task performed. Such critical information must be extracted from the changing visual image both reliably and with minimal computation. These low-level reflexive routines might form a primitive base upon which more elaborate strategies for high-performance navigation or fine object manipulation are built.

This approach has been developed previously in the domain of mobile robot navigation (Brooks, 1986; Brooks, Flynn & Marill, 1987; Aloimonos, Weiss & Bandopadhyay, 1988; Aloimonos, 1990). Brooks' mobile robots use a control architecture that decomposes the overall navigation behavior into an independent set of specific task-achieving modules. These modules incorporate specialized routines to avoid obstacles, wander, explore, monitor changes, build maps, and so on. Each module uses only simple sensory information that is most critical for achieving its desired goal, and the set of modules together provide the robot with flexible, intelligent behavior. Described as a "subsumption architecture," the system also embodies the idea that more sophisticated control behavior can be achieved by building upon more primitive mechanisms that remain intact.

The Medusa system developed by Aloimonos and his colleagues (Aloimonos et al., 1988; Aloimonos, 1990) follows a similar approach, in which a loosely coupled set of

task-oriented processes are combined to yield a range of navigational behaviors. Sensory information is obtained from an active camera system, inertial sensors and a robot arm, and separate modules compute partial image motion information, detect independently moving objects, isolate image regions that indicate approaching objects, track targets, intercept objects, and so on. Individual processes use simple algorithms specialized to their particular task.

In this paper, we first briefly consider the computation of three critical properties: (1) the 3-D direction of heading of an observer relative to object surfaces, (2) the time-to-collision between an observer and an approaching surface, and (3) the locations of object boundaries defined by discontinuities in image motion. These three properties are essential for tasks such as high-speed navigation that require rapid sensing and response, and ultimately these properties must be considered together in a system capable of performing such tasks effectively. We then focus on the computation of the 3-D direction of translation of an observer relative to object surfaces. After presenting some theoretical preliminaries, we review existing perceptual literature regarding the ability of human observers to judge heading direction. We then consider existing algorithms for performing this computation in light of these perceptual observations. This analysis leads us to focus on a particular model proposed by Rieger and Lawton (1985) that exhibits some of the behavior observed in human judgements of heading and also fits well into the overall approach described above. We present some modifications and extensions to Rieger and Lawton's model that are aimed primarily at improving its performance in the presence of image noise and allowing it to cope with self-moving objects in the scene. The results of computer simulations with this model address its behavior when applied to visual patterns similar to those used in perceptual studies and synthetic images of scenes containing self-moving objects. Finally, we discuss a number of questions that arise from this work that could form the basis for further perceptual experiments in this area.

CRITICAL INFORMATION FOR NAVIGATIONAL TASKS

In this section we briefly consider the interaction between three critical processes that must underlie tasks such as navigation: (1) computation of the 3-D direction of translation of the observer relative to object surfaces, (2) assessment of the time-to-collision of the observer with approaching surfaces, and (3) the segmentation of the scene into distinct objects on the basis of motion discontinuities. With regard to segmentation, it is critical to distinguish between objects that are stationary with respect to the background and those that undergo their own self-movement. We argue informally that these three computations taken together are essential, even for the most basic navigational behavior, and that effective navigation can be performed with only these three properties. Many of the observations here are straightforward and have been considered previously in the design of navigation systems. Before focusing on the computation of observer heading,

however, it is important to consider the broader context in which this information is used, because this context places additional demands on the heading computation.

Consider an observer moving rapidly through a cluttered scene toward a moving or stationary target, while avoiding obstacles in its path. Clearly the observer must continually assess its 3-D direction of translation relative to its target, in order to make constant, correct adjustments of its heading direction to maintain a trajectory toward the target. In principle, either the absolute or relative directions of translation of the observer and target could be computed, but for the purpose of tracking, a minimal system must at least be able to judge reliably whether the observer is heading to the left or right of the target, and the precision with which the observer makes this qualitative judgement should increase as the observer's heading becomes closer to the direction toward the target.

In addition to monitoring heading direction, the observer performing a tracking task must continually adjust his forward speed, in order to insure that he is gaining on the target. In principle, one could try to assess the absolute 3-D distance between observer and target, and the absolute speed of both, but for a minimal system, it may be sufficient to monitor the expected time-to-collision between the observer and target, which essentially depends on the ratio between 3-D distance and speed. If time-to-collision is non-decreasing, then the observer must increase his forward speed to intercept the target. Note that from motion information alone, it is only possible to recover the ratio of speed and distance, rather than absolute parameters. Perceptual studies suggest that human observers can, in fact, judge time-to-collision in contexts where there is no information regarding absolute distance or speed (Todd, 1981; Schiff & Detwiler, 1979; McLeod & Ross, 1983; Simpson, 1988). There is also behavioral data that suggests that time-to-collision estimates are used in the control of complex motor behavior (Lee, 1974, 1976, 1980; Lee & Reddish, 1981; Lee, Lishman & Thomson, 1982; Lee, Young, Reddish, Lough & Clayton, 1983). These perceptual and behavioral studies have proposed mechanisms to extract time-to-collision information from simple image motion measurements, without a complete solution to the structure-from-motion problem.

The observer must also monitor his heading relative to other object surfaces in order to detect potential collisions with stationary or moving objects in the scene. The judgement of time-to-collision is again essential, as the observer should only initiate an avoidance behavior if an object surface is moving directly toward the observer and its expected time-to-collision is small. Without an assessment of time-to-collision, the observer is likely to initiate non-essential avoidance behavior. The magnitude of the time-to-collision estimate is also useful for determining how rapidly the observer must react to an impending collision.

Finally, judgements of relative heading and time-to-collision alone are not sufficient to support effective navigation. It is also essential to determine the locations of object boundaries, from discontinuities in motion or other visual properties. Such boundaries are used in a variety of ways. First, the rapid detection of motion discontinuities quickly

draws the observer's attention to regions of the image containing objects that could potentially collide with the observer, and allows the observer to segment a target from a moving background. Second, the detection and localization of object boundaries allows an assessment of the overall size and shape of relevant objects in the scene. If an object is moving directly toward the observer, this information is essential for determining an appropriate avoidance movement that successfully steers the observer clear of the approaching object. If the object is a target being tracked, knowledge of its size and shape allows an assessment of its rough center of mass, which can serve as the focus of the observer's approach.

Finally, segmentation is essential for computing relative heading and time-to-collision reliably and accurately, as it allows the observer to integrate only those motion measurements contained within single objects to compute their properties of motion. Without this segmentation, the computation of 3-D motion parameters can be degraded by the inclusion of motion measurements from adjacent object surfaces undergoing different motions relative to the observer. This is especially problematic for the case of small objects that may be moving directly toward the observer. Also, patterns of movement created by multiple objects undergoing self-motion can mimic velocity patterns that would normally be created by critical situations such as a directly approaching object. For example, a set of objects arrayed around a circle and moving away from the center of the circle mimic the pure expansion that is characteristic of an approaching object. The detection of object boundaries from motion discontinuities allows the distinction of these situations. For obstacle avoidance, it is further useful to distinguish whether an approaching surface is stationary relative to the background, or undergoing its own motion, because self-moving objects may undergo accelerative components of motion.

Many algorithms have been proposed for the detection and localization of object boundaries from motion discontinuities, which detect these boundaries either before, during, or after the computation of 2-D image velocities (for example, Reichardt & Poggio, 1980; Hildreth, 1984; Adiv, 1985; Mutch & Thompson, 1985; Thompson, Mutch & Berzins, 1985; Schunck, 1986; Spoerri & Ullman, 1987; Hutchinson, Koch, Luo & Mead, 1988; Waxman & Wohn, 1988; Wohn & Waxman, 1990; for review, see Hildreth & Koch, 1987). Heuristics have been suggested for distinguishing stationary and self-moving objects (Jain, 1984; Heeger & Hager, 1988; Zhang, Faugeras & Ayache, 1988; Burt, Bergen, Hingorani, Kolczinski, Lee, Leung, Lubin & Schvaytser, 1989; Enkelmann, 1990; Frazier & Nevatia, 1990; Nelson, 1990; Thompson & Pong, 1990), some of which are based on the behavior of motion measurements around motion discontinuities (Thompson & Pong, 1990). Human observers are very sensitive to relative movement (for review, see Nakayama (1985)), although it appears that a large difference in direction and speed of motion may be required to localize a boundary accurately (Hildreth, 1984). Human observers can also detect very small objects relative to a moving background (Hildreth, 1984).

**DERIVING 3-D DIRECTION OF TRANSLATION
— THEORETICAL PRELIMINARIES**

In this section we present the basic equations relating image motion measurements to the parameters of translation and rotation of the observer relative to the scene. The formulation here assumes that the observer is moving relative to a stationary scene, but the same geometric relationships hold locally for the case where an object is moving rigidly relative to the observer.

When an observer moves relative to the environment, he induces a pattern of movement on the surface of the eye due to his translation and rotation relative to objects in the scene. Assume for now that the observer is moving and the environment is static, and that a coordinate system is fixed with respect to the observer, with the Z -axis directed along the optical axis. The translation of the observer can be expressed in terms of translation along three orthogonal directions, which we will denote by the vector $\mathbf{t} = (t_x, t_y, t_z)^T$. Similarly, the rotation of the observer can be expressed in terms of rotation around three orthogonal axes, which we will denote by the vector $\mathbf{w} = (w_x, w_y, w_z)^T$. Let the position of a point P in space be given by the coordinate vector $\mathbf{r} = (X, Y, Z)^T$. Then the 3-D velocity of P in the observer's coordinate frame is given by:

$$\mathbf{V} = (\dot{X}, \dot{Y}, \dot{Z})^T = -\mathbf{t} - \mathbf{w} \times \mathbf{r}$$

where

$$\dot{X} = -t_x - w_y Z + w_z Y$$

$$\dot{Y} = -t_y - w_z X + w_x Z$$

$$\dot{Z} = -t_z - w_x Y + w_y X.$$

If we assume perspective projection of 3-D velocity \mathbf{V} onto the image plane, with a focal length for the projection of 1, then the projection of point P onto the image (x, y) is given by:

$$x = \frac{X}{Z} \quad y = \frac{Y}{Z}.$$

The projected velocities in the image plane (\dot{x}, \dot{y}) are then given by:

$$\dot{x} = \frac{-t_x + x t_z}{Z} + w_x x y - w_y (x^2 + 1) + w_z y$$

$$\dot{y} = \frac{-t_y + y t_z}{Z} + w_x (y^2 + 1) - w_y x y - w_z x.$$

The first term represents the component of image velocity due to the translation of the observer and depends on the depth Z to each point in the scene. The remaining terms represent the component of velocity due to the observer's rotation, and depend only on the rotation parameters and image location.

The translational component alone yields a radial pattern of velocity, which in the case of forward translation, emanates from a single location in the image referred to as the *focus of expansion* (FOE), corresponding to the observer's direction of heading. Note that this translational component depends on the ratios of the three translation parameters to depth Z . Thus it is not possible from motion information alone to recover the absolute translation and depth parameters.

THE PERCEPTION OF OBSERVER TRANSLATION

Although the image velocity field contains components of motion due to both the observer's rotation and translation, psychophysical studies have concentrated on our ability to measure direction of translation. Navigation tasks impose severe demands on our ability to perform this computation. Cutting (1986) has shown that under reasonable assumptions, we require an accuracy of about 1° of visual arc in our judgement of heading in order to avoid obstacles successfully while running and driving, as well as performing more challenging tasks such as downhill skiing and aircraft landing. This section reviews perceptual studies of the ability of human observers to judge their direction of translation, which suggest that the human visual system can, in fact, achieve this degree of accuracy under the best conditions. We summarize these studies in some detail, as they will form the basis for the computer simulations described later.

A series of experiments by Warren and his colleagues (Warren & Hannon, 1988, 1990; Warren, Morris & Kalish, 1988) measured the accuracy with which observers can judge their heading direction in computer displays that simulate movement toward a planar surface or 3-D cloud of random dots. The first experiments simulated movement along a ground plane extending to a visible horizon. A target vertical line segment was located somewhere along the horizon, and the subjects' task was to judge their direction of heading relative to the vertical target. That is, the simulated heading would differ from the direction of the target by varying angular differences, and the subject had to judge whether the heading on a particular trial was to the left or right of the vertical line segment. In the initial experiments, the vertical target was visible throughout the motion of the points, but in subsequent experiments, the target only appeared after the points stopped moving. A number of factors were varied in these experiments, including the orientation of the plane relative to the viewer, the observer's speed and direction of heading, dot density and the temporal extent of the motion. In later studies, movement was simulated relative to a random cloud of dots distributed in a 3-D volume of space. The field of view in the experimental displays was usually 40° horizontal by 32° vertical,

and the size of individual elements did not change with their position or movement in depth.

The general conclusion of the studies by Warren and his colleagues is that human observers can judge their heading direction with an accuracy of $1^\circ - 2^\circ$ of visual angle, for a variety of surface types and under a range of experimental conditions. Performance is the same, regardless of whether the vertical target line is visible during the movement of the points. Observers perform better when higher speeds of translation are simulated, consistent with earlier observations by Johnston, White and Cumming (1973), and Carel (1961).

Warren and Hannon (1988, 1990) compared performance under three conditions: (1) the observer fixated a stationary marker on the display, and the displays only simulated pure translation of the observer, (2) the observer tracked a moving point in the display, thus introducing a rotational component of motion, and (3) the display itself contained both translational and rotational components of motion, and the observer was required to maintain stationary fixation. For conditions (2) and (3), the same flow pattern appears on the surface of the eye, but in condition (2), rotational information could, in principle, be derived from extraretinal eye movement signals, while in condition (3), such information must be derived from visual input alone. Subjectively, observers could not distinguish between displays corresponding to conditions (2) and (3), and in the latter case, there was a strong illusion of the eye actually moving. It was found that for the case of movement toward a ground plane, or movement toward a random cloud of dots, there was essentially no difference in performance between these three conditions. Heading was computed with an accuracy of $1 - 2^\circ$ in all three conditions. When simulating translation perpendicular to a plane, however, performance still reached this level of accuracy in the first two conditions, but was at chance for the third condition, in which the rotational component was added to the movements of the dots in the display. Subjectively, observers perceived themselves as moving toward the point of fixation, which would correspond to a center of outflowing motion in this case. Similar observations regarding movement toward a frontoparallel plane were made in other studies (Llewellyn, 1971; Johnston et al., 1973; Regan & Beverley, 1982; Rieger & Toet, 1985; Cutting, 1986). This observation suggests first, that extraretinal information regarding eye rotation is used in the analysis of heading direction, and second, that the passive decoupling of the rotational and translational components of motion from visual input alone requires differential motion produced by elements at different depths.

Warren and Hannon (1990) also examined the influence of dot density for simulated movement toward a 3-D cloud of dots. The number of dots visible at the beginning of each trial was either 6, 12, 25 or 50. The overall field of view was kept constant at $40^\circ \times 32^\circ$, so the density of dots changed with their number. An average "neighborhood size" was calculated that assumed the values 6° , 4° , 2° and 1° , for 6, 12, 25 and 50 dots, respectively. A neighborhood size of 2° , for example, was defined to mean that

there was an average of three or more pairs of dots with angular separations less than or equal to 2° , but not three pairs separated by less than or equal to 1° , in the first frame of the display. When the added rotational flow was generated by the subject tracking a dot on the display, there was no change in performance with dot density (confirming earlier observations by Warren et al. (1988)), but when rotational flow was added to the movements of the points, there was some degradation of performance with lower densities, which became significant at the lowest density (only 6 dots in the display, with a neighborhood size of 6°). Thus, observers could accurately judge heading direction when presented with a relatively sparse, discontinuous flow field.

The experiments by Warren and Hannon (1988, 1990) and Warren et al. (1988) used a total viewing time of about 3 seconds, with image sequences of about 50 frames. It was later found that for the case of pure translation of the observer, there is no deterioration in performance if the number of frames is reduced, until only 2 – 3 frames are presented (Warren, Griesar, Blackwell, Kalish & Hatsopoulos, 1990). There is still about 3° of accuracy for only two frames, with significant improvement when a third frame is added. Experiments were also conducted in which the total duration of the motion was about 3 seconds, but the lifetimes of individual dots were varied. With a total duration of about 3 seconds and lifetime of only two frames, subjects could achieve about 1° of accuracy in judging heading direction, again for the case of pure translation. For the case in which a rotational component is added to the motions of the points, a more extended time period may be needed to recover observer heading accurately (W. Warren, personal communication).

The visual system can also tolerate significant noise, with performance degrading smoothly with increased amounts of noise. Warren et al. (1990) found, for example, that in the case of pure translation of the observer, subjects could still judge heading direction with an average error of 2.6° when the directions of motion of individual points were randomly perturbed within an envelope of 90° . This result suggests that the heading computation may involve significant spatial pooling of image motion measurements.

Cutting (1986) examined observers' ability to determine their direction of translation toward a field of vertical lines that were placed on three frontoparallel planes that were separated in depth. Subjects were asked to judge whether the display simulated a view that was to the left, right or in the direction of heading. When the lines were placed at the same depth, subjects performed at chance, and their heading accuracy improved with an increased separation of the lines in depth. The best accuracy achieved corresponded to an angle between gaze and heading directions of about 1.25° .

Rieger and Toet (1985) measured subjects' ability to judge their heading direction relative to two frontoparallel planes of random dots placed at different depths. Translational and rotational components of motion were combined in the movements of the points on the display. The following parameters were varied in these experiments: (1) the direction of translation (possible directions were separated by 2.5° or 5°), (2) the

separation in depth between the two planes (either 0 or 9 meters), (3) the magnitude of the rotational component of observer motion, and (4) the size of the field of view (either 10° or 20°). Dot density was high, with an average of 700 dots in a single static frame. For the case of a single plane, performance degraded rapidly as the magnitude of simulated rotation was increased, similar to previous studies. When the points were placed at different depths, however, subjects could reliably judge heading direction at both resolutions (although performance was worse at the finer resolution), over the range of angular rotations tested, and with little degradation with the size of the field of view.

To summarize the perceptual experiments, we make the following observations regarding the human recovery of direction of translation:

- Human observers can achieve an accuracy of about $1^\circ - 2^\circ$ of visual angle at judging heading direction, with or without the presence of a target in the environment.
- Performance improves with higher speeds of translation.
- Performance improves when surfaces span a greater range of depth.
- Extraretinal information regarding eye rotation is used in the recovery of heading direction.
- Heading direction can be judged reliably in the presence of significant amounts of noise in the image motion measurements.
- For the case of pure translation, heading direction can be recovered accurately in a relatively short time of 2 or 3 frames, with accuracy increasing with time.
- Heading direction can also be recovered in a context where the rotational and translational flows must be passively decoupled from visual input alone. This decomposition
 - (1) requires differential motion produced by elements at different depths,
 - (2) can be performed successfully with sparse, discontinuous flow fields, and
 - (3) requires only a relatively small field of view, at least as small as 10° .

The next section examines computational models for the recovery of observer motion in light of the above observations.

THE COMPUTATION OF DIRECTION OF TRANSLATION

Computational methods for recovering the direction of translation of an observer relative to a scene can be broadly divided into two classes, depending on whether they rely on discrete or continuous image motion measurements. In the discrete approach, a set of isolated image features are tracked over time, and their sequence of positions forms the input to a system of equations whose solution depends on the parameters of 3-D structure and motion. In the continuous approach, it is usually assumed that an

instantaneous 2-D velocity field is available at one or more instants of time, and the image velocities, together with their spatial or temporal derivatives, are used to solve for 3-D structure and motion parameters.

Many examples of the discrete approach present theoretical results regarding the minimal number of motion measurements required to compute 3-D structure and motion parameters uniquely (for example, Ullman, 1979; Prazdny, 1980; Longuet-Higgins, 1981, 1984; Tsai & Huang, 1984a,b; Faugeras, Lustman & Toscani, 1987; Aloimonos & Brown, 1989; Weng, Huang & Ahuja, 1989). The direct application of the mathematical results suggests possible algorithms for recovering these parameters, but computer experiments with these algorithms indicate that they may be vulnerable to error in the image motion measurements. The ability of the human visual system to judge heading direction accurately for a few, sparse features in motion suggests that the underlying computation must be able to derive movement parameters from discrete motion measurements, but unlike existing algorithms, the human system can tolerate large amounts of noise in these sparse measurements. Recent algorithms that use discrete motion measurements over a more extended time period exhibit better performance (Ullman, 1984; Broida & Challappa, 1986; Shariat, 1986; Faugeras et al., 1987). Extended time appears to be necessary for the human system to decouple rotational and translational components of motion on the basis of visual input alone (W. Warren, personal communication).

Continuous approaches that use spatial derivatives of velocity require a locally continuous velocity field, or one that is sufficiently dense that interpolation can be used to approximate the continuous field (for example, Longuet-Higgins & Prazdny, 1981; Koenderink & Van Doorn, 1976; Waxman & Ullman, 1985; Subbarao, 1988; Waxman & Wohn, 1988). Other recent models have used the theory of planar dynamical systems as a basis for recovering information about 3-D motion and structure (Verri, Giroso & Torre, 1989), where the time evolution of the structure of the flow field in the vicinity of singularities (such as the FOE) is used to recover motion parameters. These continuous methods may have difficulty with the sparse and discontinuous velocity fields used in perceptual studies. Some of these techniques also require accurate velocity measurements, which make them vulnerable to noise. Methods that rely directly on spatial and temporal derivatives of image intensity (Negahdaripour & Horn, 1987, 1989; Horn & Weldon, 1988; Heel, 1990a,b) may have difficulty coping with the impoverished displays of isolated dots used in perceptual studies.

There are other velocity based approaches that do not require a continuous velocity field (for example, Bruss & Horn, 1983; Ballard & Kimball, 1983; Jain, 1983; Lawton, 1983; Adiv, 1985; Burger & Bhanu, 1990; Heeger & Jepson, 1990). Some of these methods use an optimization approach, in which 3-D motion parameters are computed that yield a velocity field that best fits the observed image velocities in the least-squares sense, and integrate a large number of image motion measurements, yielding less sensitivity to error. The human system, however, does not require extensive spatial integration to compute

heading direction accurately; in contrast, it can cope with both a small number of motion measurements and a relatively small field of view.

Finally, some methods make direct use of information about motion parallax, that is, the relative motion of features at different depths, to derive 3-D motion and structure (Longuet-Higgins & Prazdny, 1981; Rieger & Lawton, 1985; Cutting, 1986). The difference in velocity between two points that are nearby in the image, but separated in depth, depends largely on the translational parameters of observer motion and can be used directly to infer the direction of translation. The explicit reliance of these methods on depth variation in the scene makes them appealing from the perspective of the human system, which fails for the case of the perpendicular approach to a plane.

To summarize, it appears that most existing models do not exhibit the basic properties of the human recovery of direction of translation. None of these models have been shown to yield the accuracy of $1^\circ - 2^\circ$ of visual angle seen in human judgements of heading, over a range of viewing conditions. Some current models could be modified to cope with some of the range of conditions considered in perceptual studies, but the need to cope with sparse, noisy and discontinuous motion fields, and the failure of the human system with the frontoparallel plane, seems to rule out many models on more fundamental grounds.

THE RIEGER AND LAWTON MODEL

This section describes the algorithm proposed by Rieger and Lawton (1985), which is based on earlier work by Longuet-Higgins and Prazdny (1981). This class of models begins with the observation that at the location of a discontinuity in depth, there will be a discontinuity in the translational component of the image velocity field because of the dependence of this component on depth, while the rotational component will be roughly constant across the boundary. Furthermore, if we construct a field of vectors that represent the differences in velocity across these boundaries, these vectors will be oriented approximately along the lines connecting their image location with the focus of expansion (the so-called *translational field lines*), and therefore should all point to the FOE.

Longuet-Higgins and Prazdny suggested an algorithm based on the above observations that uses instantaneous spatial derivatives of velocity to recover the FOE. This original algorithm proved to be quite sensitive to error in the image velocity measurements. A robust algorithm that uses this observation to extract the FOE must take into account the fact that accurate velocity measurements may not be available immediately to either side of a depth discontinuity. Rieger and Lawton (1985) presented an algorithm that addresses this problem. The basic steps of the algorithm are as follows. First, the differences between each local image velocity and every other velocity measured within a restricted neighborhood are computed. From the resulting distribution of velocity differ-

ence vectors, the dominant orientation of the vectors is computed, and this information is preserved only at locations where the distribution of velocity differences is strongly anisotropic. Such points will typically arise at locations where there is a strong depth variation in some direction. The result of this first stage is a set of directions at a number of points in the image, which are all roughly aligned with the translational field lines. The FOE is then calculated as the “best-fit” intersection point for all the resulting vector directions. Once the FOE is determined, then the direction of the translational component of motion is known at every location in the image, so that any motion in the original flow field that is perpendicular to this direction must be due to the rotation of the observer. From these perpendicular motions, the best rotational parameters are inferred (a similar strategy is used by Burger & Bhanu (1990)). Once the rotational parameters are estimated, the full rotational flow field can be computed and subtracted from the original flow field to obtain the full translational component of the flow field. Finally, the relative depth at every point can be computed from knowledge of the FOE and magnitude of the translational component of motion at each location.

The algorithm proposed by Rieger and Lawton is appealing for a number of reasons. First, it provides a rough initial estimate of the direction of translation, independent of the rotation parameters and 3-D shape of the surface. As we discussed earlier, heading direction is a critical property of observer motion for navigation that must be computed with high accuracy and speed. We also noted that it is important to detect object boundaries from motion discontinuities as soon as possible, and these are precisely the locations that provide the best information for this algorithm. Another appealing aspect of this algorithm is its simplicity and reliance on primitive image motion information, such as velocity differences, that require little computation. The fact that it does not rely critically on the solution of optimization problems is also an advantage. Optimization is used to some extent at each step of the algorithm, but the information being computed at each of these steps could be obtained to a close approximation with non-iterative techniques.

One question that arises regarding the Rieger and Lawton algorithm as it stands is whether it can achieve the degree of accuracy of human performance measured across the range of conditions that have been considered in perceptual studies. Simulation results presented by Rieger and Lawton (1985) suggest that the resulting heading accuracy may be at least within a factor of two or three of the needed accuracy. An especially challenging aspect of human performance, however, is its ability to cope with sparse displays. The average angular separation between points in Warren and Hannon’s (1990) study is large compared to the neighborhood sizes used in Rieger and Lawton’s simulations. Over larger distances, the basic observations that the model relies upon become less valid. The computer simulations presented later suggest that this algorithm can yield the desired accuracy for the particular conditions of the perceptual experiments, with reasonable assumptions about the available precision of image motion measurements.

BUILDING UPON THE RIEGER AND LAWTON MODEL

From a computational standpoint, the most severe limitation of Rieger and Lawton's model is that it does not cope with self-moving objects in the environment. The difference in velocity across the boundary between a self-moving object and stationary background, or between two self-moving objects, in general does not yield vectors that are oriented along the translational field lines that emanate from the true FOE. Combining these differences with those obtained along the boundaries between stationary surfaces can yield significant error in the computed FOE location, especially if self-moving objects cover a large part of the visual field. Thus, it becomes necessary either to detect self-moving objects explicitly or to remove their influence on the FOE computation by some implicit means.

This section first considers a different method for performing the FOE computation in Rieger and Lawton's model that allows self-moving objects to be present in the scene and helps to isolate the boundaries of such objects. In this context, we also summarize previous methods for coping with self-moving objects. We then present some additional modifications to other stages of the algorithm that improve its performance in the presence of error in the image motion measurements. The results of computer simulations with the algorithm described here are presented in the next section.

Coping with Self-Moving Objects

We first consider existing methods for detecting and coping with self-moving objects in the scene. One approach assumes that the camera is stationary, so that significant image motion indicates self-moving objects (for example, Jain, Militzer & Nagel, 1977; Jain, Martin & Aggarwal, 1979; Anderson, Burt & Van der Wal, 1985; Dinstein, 1988; Bouthemy & Lelande, 1990). A variation on this approach considered by Burt et al. (1989) implicitly recovers global camera motion parameters by attempting to stabilize regions of the image, analogous to eye tracking in the human system. Once the image motion due to the actual camera motion is largely removed, any significant motions that remain are likely to be due to self-moving objects. A second approach assumes that the camera undergoes pure translation, so that any self-moving objects violate the expected pure expansion of the image (for example, Jain, 1984). If 3-D depth data is available, then inconsistency between image velocities, estimated observer motion and depth data can also signal self-moving objects (for example, Thompson & Pong, 1990). Nelson (1990) shows that it is possible to detect such inconsistencies from partial information about image motion and observer motion. Nelson also notes that the motion of objects due to the observer's motion tends to change slowly over time, while self-moving objects can sometimes generate rapidly changing patterns of motion that can be used to detect their presence.

A more general strategy is to compute an initial set of observer motion parameters, either by combining all available data or by performing separate computations within limited image regions, and then to find areas of the scene that move relative to the observer in a way that is inconsistent with the global motion parameters (for example, Heeger & Hager, 1988; Zhang et al., 1988). With these latter approaches, if all motion information is used initially, the recovery of observer motion parameters can be degraded by the inconsistent motions of self-moving objects, especially if self-moving objects cover a significant portion of the visual field. On the other hand, the use of spatially local information can yield inaccuracy due to the limited field of view. Thompson, Lechleider & Stuck (1991) present a variation on this approach that uses a technique from robust statistics (Huber, 1981) to compute global motion parameters in the presence of so-called "outliers," which are data that deviate significantly from consistency with the true parameters. Image motions resulting from self-moving objects are treated as outliers, and the least median squares algorithm (Rousseeuw & Leroy, 1987; Meer, Mintz, Kim & Rosenfeld, 1991) is used to compute motion parameters in a way that detects these potential outliers. Thompson et al. (1991) note that self-moving objects whose projected image motion is close to the motion that is expected from the observer's global translation and rotation are difficult to detect with this technique.

We present here a different strategy for detecting and coping with self-moving objects that builds upon the Rieger and Lawton algorithm. We first summarize the basic strategy in general terms and then elaborate on the motivation and details. The scheme first computes local velocity differences and determines the dominant orientation of the distribution of velocity differences within a small neighborhood of each point, as in the Rieger and Lawton model. The orientations, θ_i , are preserved for the next stage of the computation only at points where the distribution of velocity differences is strongly anisotropic. As noted earlier, most of the θ_i measurements preserved at this stage are derived from points on or near depth discontinuities, or along surfaces such as the ground plane, whose angle of slant relative to the image plane is large.

Some portion of the θ_i measurements will point roughly toward the true FOE location, while θ_i measurements obtained in the vicinity of self-moving objects or those with high error will be oriented in arbitrary directions. Assuming that self-moving objects do not cover a large part of the visual field, we can obtain a good initial guess of the rough location of the FOE by looking for limited image regions for which a large percentage of the θ_i measurements point toward locations within the region. In particular, we consider how much evidence exists to support the FOE being located within a large set of possible image regions, then choose the region (or regions) with maximum support and use the θ_i measurements that provide this maximum support to derive an FOE estimate. This strategy is analogous to the Hough transform technique used extensively in computer vision (Ballard & Brown, 1982).

In more detail, after the θ_i measurements have been derived, the visual image is

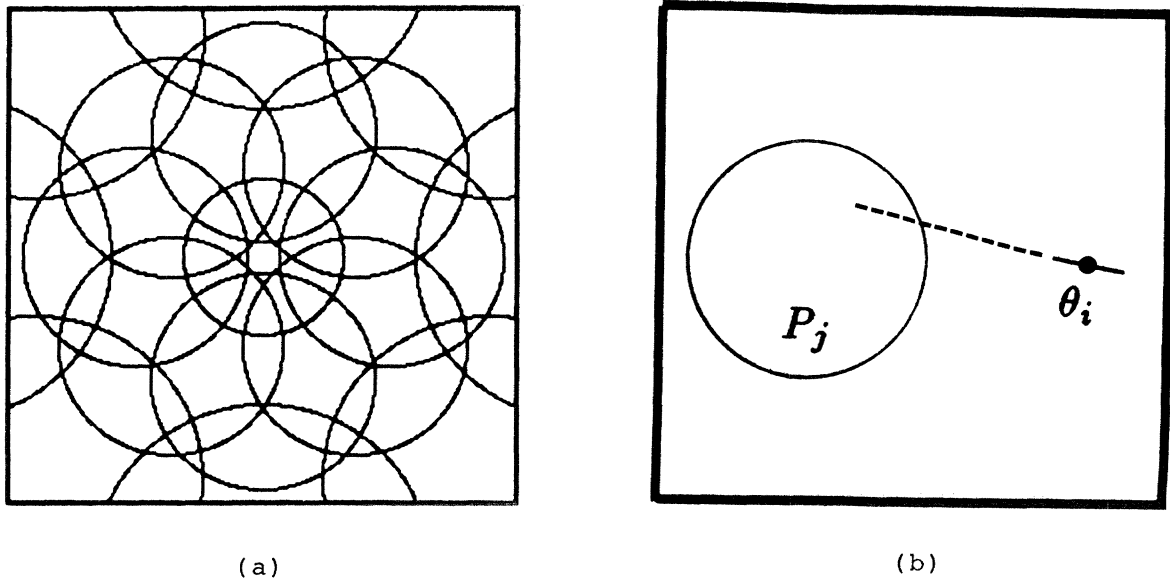


Figure 1. (a) A set of overlapping circular patches that represent regions of the image that could contain the FOE. (b) Positive evidence for the FOE being located within P_j is given by a measurement θ_i if a line from the point that contains the vector defined by θ_i intersects P_j .

divided into a set of overlapping, circular patches that represent possible regions within which the FOE may be located, as shown in Fig. 1a. For each patch P_j , we collect all of the positive evidence for the FOE being located within P_j . Positive evidence comes from points whose orientation θ_i lies along a line that intersects P_j , as shown in Fig. 1b. If the true FOE is located within the patch P_j , then velocity differences computed within the surfaces of stationary objects or along boundaries between two stationary objects should yield positive evidence. In this case, points at which an orientation θ_i is obtained that does not yield positive evidence for the FOE being located within P_j either lie within or near the boundaries of self-moving objects, or they are projected from stationary regions of the scene, but result in significant error in the computation of θ_i . If the true FOE is not located within P_j , there will still be a number of points that yield an orientation θ_i that incorrectly provides positive evidence for an FOE in P_j , but the percentage of points yielding such false positive evidence should be substantially reduced.

For each patch P_j , if a sufficiently large percentage of the available θ_i yield positive evidence for the FOE being located within P_j , then the set of θ_i estimates yielding this positive evidence is used to generate a hypothesized FOE location. If this hypothesized FOE is located well within the patch, it is preserved for later consideration. If multiple FOE hypotheses remain after this stage, they are reconciled to obtain a single FOE location by considering the extent of the positive evidence in their support, their

goodness-of-fit to the computed θ_i , and the proximity of the multiple hypotheses.

The reasoning behind this strategy is that by combining only those θ_i measurements that yield positive evidence for an FOE being located within restricted patches, we significantly reduce the degradation in the FOE computation that can result from the presence of self-moving objects and from large errors in the θ_i estimates. When patches that contain the true FOE are considered, self-moving objects and large errors in the θ_i computation are likely to result in θ_i estimates that do not yield positive evidence and hence do not enter into the FOE computation. Patches that do not contain the true FOE are likely to yield significantly less positive evidence and therefore do not lead to an FOE hypothesis.

As shown in Fig. 1a, the circular patches may increase in size with distance from the center of the image. This serves both to minimize the total number of patches needed to cover the image and to allow the FOE to be computed more accurately when it is located toward the center of the image. Reducing the total number of patches reduces the amount of computation required to test the set of patches for possible FOE locations. The desire to compute the FOE more accurately toward the center of the image is motivated in part by properties of human visual processing. Human observers judge their heading direction most accurately when their eyes are pointed in the direction of heading, and the spatial resolution of processing in general increases toward the center of the eye. Thus heading direction is derived most accurately when the FOE lies near the center of the visual image.

The determination of whether a particular measurement θ_i is consistent with the FOE being located within a patch P_j requires a simple computation. We began with the Rieger and Lawton model in part because of the simplicity of the criterion for determining whether the image motion around a point is consistent with a restricted window of FOE locations. We can either determine whether the orientation θ_i falls within a limited cone of directions defined by the two lines running through the underlying point and tangent to the circular boundary of P_j , or whether the perpendicular distance from the center of P_j to the line containing the vector in the direction θ_i is less than the radius of P_j . The measurements of θ_i obtained from points within P_j are not included in the positive evidence for P_j , because the size of the translational component of velocity is usually very small in the vicinity of the FOE, yielding velocity differences that are not reliable indicators of the location of the FOE. We also limit the overall extent of the region from which θ_i measurements are considered for P_j , because the range of consistent orientations θ_i becomes too small for points very distant from P_j , requiring too much accuracy in their estimate.

After the set of θ_i that yield positive evidence for a given patch are computed, we determine whether there is sufficient evidence to combine these θ_i measurements to derive an FOE hypothesis. In particular, we calculate the percentage of all θ_i measurements that yield positive evidence and compare this percentage to a threshold. This threshold

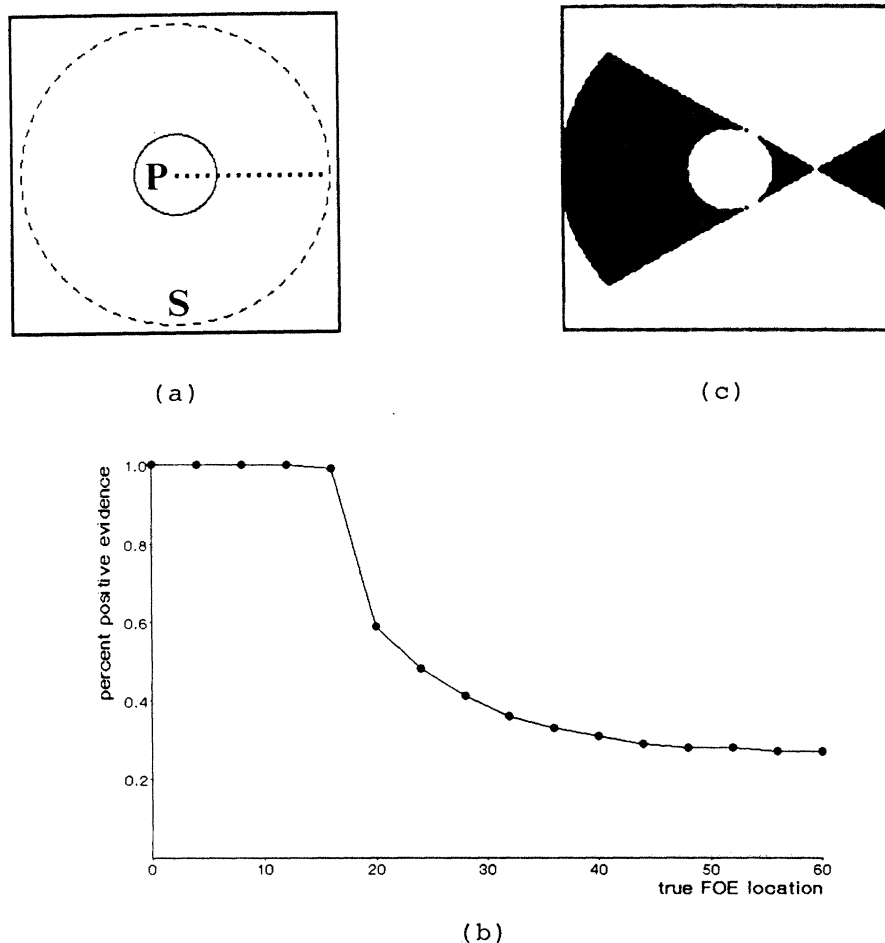


Figure 2. (a) We consider the positive evidence for the FOE being located within the central patch P from θ_i measurements that could be obtained within the larger annular region S , for the set of true FOE locations indicated by the solid dots. The radius of P is 16 pixels, and radius of S is 64 pixels. (b) Graph of the percentage of θ_i measurements that would provide positive evidence for the FOE being located within P as a function of the true location of the FOE. (c) Given the patch P with radius 16, and a true FOE located 32 pixels to the right of the center of P , the points that could yield positive evidence for the FOE being located within P are shown in black.

must be large enough to minimize the number of false hypotheses generated from patches that do not contain the true FOE, while at the same time allowing a significant portion of the visual field to contain self-moving objects. Thus the choice of threshold here is governed in part by what percentage of points yielding θ_i measurements are expected to be within or near the boundaries of self-moving objects, and in part by what percentage of points from stationary regions of the scene are expected to yield false positive evidence for inappropriate FOE locations. With regard to the first factor, we note that if too much of the visual field contains self-moving surfaces, human observers do not judge their heading correctly.

Fig. 2 addresses the second of the factors mentioned above. We consider a patch P at the center of the image, as shown in Fig. 2a, and determine the positive evidence that could be obtained for FOE locations within P , for different true locations of the FOE. Evidence is considered from all points lying within a large circular region surrounding P , and we assume for this example that every point in the image yields a measurement of θ_i that is directed along translational field lines emanating from the true FOE. The graph in Fig. 2b shows the percentage of θ_i measurements that represent positive evidence for FOE locations within P for the set of true FOE locations indicated with solid circles in Fig. 2a. When the true FOE is located within P , 100% of all θ_i measurements yield positive evidence, but as the true FOE moves outside P , the percentage of points that could, in theory, yield positive evidence for FOE locations inside P drops rapidly. (For the simulations presented in this paper, we required that 40–50% of the θ_i measurements yield positive evidence for a particular patch P_j , in order to generate an FOE hypothesis from P_j .) Fig. 2c shows a map of the points that could yield positive evidence for the FOE being located within P when the true FOE is located outside P , as described in the figure legend. If, in a particular scene, all of the available measurements of θ_i happen to fall within the regions shown in black in Fig. 2c, then it could appear that there is significant positive evidence for an FOE within P , and the set of θ_i measurements would be combined to generate an FOE hypothesis. If the true FOE is located outside P , the estimate obtained here may not have as good a fit to the θ_i measurements as the FOE hypothesis generated from a correct patch. In general, however, it is possible for a skewed spatial distribution of the available θ_i measurements to yield an inappropriate FOE estimate.

Self-moving objects can also yield false positive evidence for an FOE being located within a given patch P_j , especially if an object undergoes a significant translation toward or away from P_j . If the true FOE is not located within P_j , then the added θ_i measurements from self-moving objects are likely to yield an FOE hypothesis that does not yield a good fit to the θ_i measurements. Even for the patch that contains the true FOE, self-moving objects with significant translation near but not along the true translational field lines can distort the computation of the FOE location. We assume that this situation is rare, and note that when it does occur, it is unlikely to persist for an extended period of time, or over an extended region of the image.

Due in part to the overlap of adjacent patches (see Fig. 1a), valid FOE hypotheses may emerge from multiple patches. If there is a single FOE location that both accounts for a significantly larger percentage of the θ_i measurements and yields a significantly better goodness-of-fit to these measurements, then this FOE location is considered to be the best current guess. Multiple FOE locations that are close to one another can be averaged together to yield a current estimate. If, however, there are multiple FOE hypotheses that have strong support and are distant from one another, it may be possible to resolve the global FOE through an analysis of possible self-moving objects in the scene,

which we consider next.

If there is significant positive evidence for the FOE being located within a patch P_j , then those points that do not yield positive evidence can be used to detect possible self-moving objects. In particular, extended, connected groups of such points can signal a self-moving object. Isolated points or small groups of points yielding negative evidence are more likely to be the consequence of error in the θ_i computation. Some points within or near the boundaries of self-moving objects will yield false positive evidence for an FOE within P_j . If, however, such points are connected to an extended region of points yielding negative evidence, we assume that they represent a continuation of a self-moving object and generate a new FOE hypothesis with these points removed, as long as their removal does not then lead to an insufficient percentage of the θ_i measurements yielding positive evidence for an FOE in P_j .

Finally, we note that a coarse-to-fine strategy can be used, in which larger patch sizes are used first to obtain a rough estimate of the region (or regions) likely to contain the global FOE, and the size of the patches is then successively reduced to refine the estimated FOE location. At each scale, a current estimate (or estimates) could be obtained, and smaller patches could then be centered on the current estimate. Such a coarse-to-fine strategy provides a rapid assessment of the rough FOE location and reduces the total amount of computation required to obtain a more precise estimate.

Recent work in the area of robust statistics provides a number of techniques for deriving global parameters in the presence of significant outliers in the data (for example, Rousseeuw & Leroy, 1987; Meer et al., 1991). Similar to the scheme proposed by Thompson et al. (1991), the θ_i measurements derived from self-moving objects could be considered outliers and general techniques such as the least median squares algorithm could be applied to the full set of θ_i measurements to compute an FOE estimate and detect the “outlying” self-moving objects. The approach presented here, however, takes better advantage of the geometrical relationship between θ_i measurements obtained from stationary and self-moving objects and requires far less computation.

Other Modifications of the Rieger and Lawton Model

This section considers some additional modifications aimed primarily at improving the performance of the Rieger and Lawton algorithm in the presence of error in the image motion measurements. These modifications include limited temporal smoothing of the image velocities, a different strategy for computing the dominant orientations, θ_i , that effectively filters the local distributions of velocity differences, and a method for refining the θ_i measurements at a later stage.

If the errors in the instantaneous 2-D velocities of moving features are uncorrelated from one moment to the next, then smoothing or averaging of the velocity measurements over time can improve their quality. This temporal smoothing should cover a limited

time window, however, as the observer's heading can change over a long time interval. In the simulations presented in the next section, velocity measurements with added noise that were obtained at two different times were averaged together. This limited smoothing took place prior to the computation of velocity differences, and significantly improved the quality of these difference estimates.

The local distribution of velocity differences can be computed in one of two ways. First, we could take the difference between the velocity of a point p_i and that of each neighboring point p_j within some distance of p_i , to obtain a set of velocity differences associated with the location of p_i . If p_i has n such neighbors, then the distribution will contain at most n differences. A second option is to consider fixed neighborhoods distributed over the image, and to compute the difference in velocity between every pair of points that falls within each neighborhood. In this case, if there are n points within a given neighborhood, then there will be at most $\frac{n(n-1)}{2}$ velocity differences computed within the neighborhood. Both strategies were used in the simulations described in the next section. For the simulations with the sparse dot patterns used in perceptual experiments, all pairs of points within fixed neighborhoods were used to obtain the local distributions of velocity differences, while the simulations with images on dense grids used only the differences between single locations and their neighbors.

To obtain estimates of the dominant orientations, θ_i , we first note that the distribution of velocity differences computed at a point or within a neighborhood that lies in the vicinity of a depth discontinuity or on a surface with a substantial slant in depth will typically cover a range of directions, as shown in Fig. 3a. Differences between the velocity of two points that lie at significantly different depths will be larger and have an orientation that is roughly along the translational field line that is directed toward the FOE. There will be some deviation from the orientation of the true translational field line, due to error in the velocity measurements or to the spatial separation between the two points, which yields added differences in velocity due to the rotation of the observer. (If the magnitude of the observer's rotation is not too large, the latter differences will be small.) Differences obtained from pairs of points at a similar depth will typically be smaller and have directions that are randomly distributed around the full 360° range. These latter difference measurements can degrade the computation of the dominant orientation if all of the difference measurements are considered together. To reduce this degradation, we only combine velocity differences within two opposite ranges of 90° , as shown in Fig. 3b, and choose the particular ranges that yield the largest ratio between the overall weight of the differences obtained within and outside of these ranges. Estimates of θ_i are preserved only at locations at which this ratio is above a specified threshold, indicated a strong anisotropy in the directions of the velocity differences. The θ_i themselves are computed by finding a line that best fits the set of difference vectors in the least-squares sense; that is, the sum of the squared distances of the endpoints of the vectors from this line is minimized.

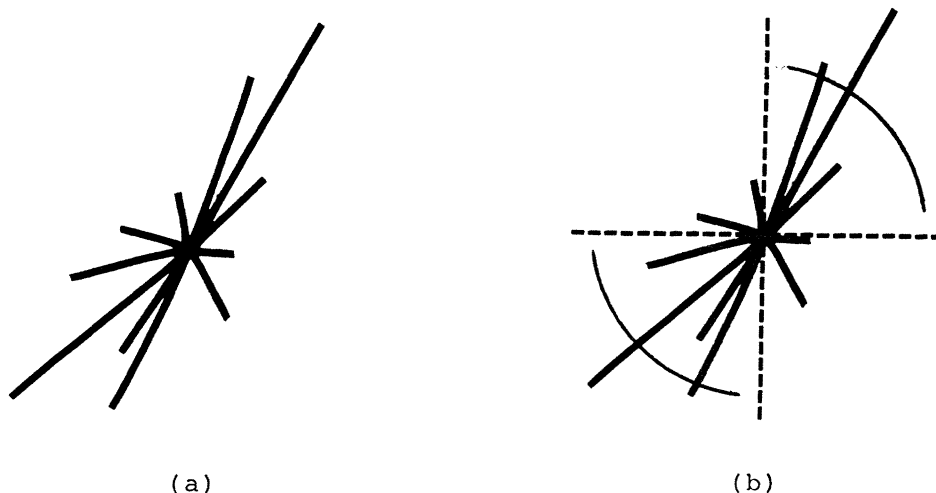


Figure 3. (a) A typical distribution of velocity differences obtained at a point that is near a depth discontinuity or located on a highly slanted surface. The larger vectors represent the difference in velocity between this point and other points lying at significantly different depths, and are directed roughly along the translational field line. Other vectors represent the difference between the velocity at this point and that of other points located at similar depths. The aim is to compute the dominant direction of these differences. (b) We find two opposite 90° ranges of orientations that separate the differences in a way that maximizes the ratio between the sum of the lengths of the velocity differences lying within and outside of these ranges.

Finally, we note that the θ_i estimates can be improved after an initial FOE estimate is obtained. An initial estimate of the location of the FOE gives rise to a set of predicted translational field lines, along which local velocity differences should lie. We can then “filter” the local distributions of velocity differences to emphasize differences whose direction is closer to the orientation of the translational field lines. A new FOE location can be computed based on the computation of new dominant orientations of the filtered local velocity differences. In principle, the same strategy can be applied over time. The location of the FOE can now change over time, so it becomes necessary to estimate the rotational component of motion as well, in order to predict the displacement of the FOE in the image due to the observer’s rotation. This can be done, for example, in the way that Rieger and Lawton (1985) propose. At each new moment in time, the current estimate of the location of the FOE can be used to weigh local velocity differences in the computation of a new FOE. A better estimate of the FOE should then result in a better estimate of the rotational component of motion, yielding progressive improvement over an extended sequence of images.

COMPUTER SIMULATIONS

This section presents the results of computer simulations with the algorithm proposed in the previous section. We consider aspects of the human recovery of heading direction and the use of the algorithm for computer vision systems.

Simulations with the Model Applied to Perceptual Displays

This section summarizes the results of computer simulations with our extension of the Rieger and Lawton (1985) model, applied to visual patterns similar to those used in the perceptual studies described earlier. We used synthetic image data corresponding to displays of discrete points whose image motion is determined by the translation and rotation of an observer relative to a random-dot surface in space. The motions of the dots on the image plane were computed analytically, and these movements, with or without added noise, formed the input to the model for heading recovery described in the previous section. Perspective projection was used throughout this analysis.

We first summarize the conditions of the perceptual experiments by Warren and his colleagues that we approximately simulated here:

- *Observer's translation:* The observer translates in the horizontal plane, with a heading direction spanning a range within 6° to the left and right of straight ahead. For most experiments, translational speed was 1.9 m/sec.
- *Observer's rotation:* The typical range of simulated angular velocity of the eye was $0.3 - 0.7^\circ/\text{sec}$, covering the full range of 2-D directions. (Note that this amount of rotation is small.)
- *Field of view:* 40° horizontal \times 32° vertical.
- *Temporal extent:* Most experiments used a total viewing time of about 3 seconds, with a frame rate of 15 frames/sec. The simulations presented here, however, used average displacements computed from only the first three image frames.
- *Ground plane:* The observer's simulated eye height was 1.6m and points covered a plane extending 37.3m in front of the observer. The spatial distribution of the points was uniform on the plane, creating a non-uniform distribution in the image, due to perspective projection.
- *3-D cloud:* Points were placed randomly within a depth range of 6.9m - 37.3m.
- *Frontoparallel plane:* A plane was placed at a distance of 9.3m in front of the observer.
- *Number of dots:* In most experiments, there was an average of 63 dots at the beginning of the movement.

In these experiments, observers were asked to judge only the horizontal component of motion. Additional error in the perception of the vertical component of heading would indicate a larger overall heading error. The accuracy of $1^\circ - 2^\circ$ measured in perceptual experiments refers to the horizontal component alone.

The computer simulations also considered the following conditions: (1) points placed on two frontoparallel planes, whose absolute and relative depths were varied, (2) variation in the absolute and relative range of depth for the 3-D cloud, (3) wider heading angles, with directions ranging up to 30° to the left and right of straight ahead, (4) larger rotational components, corresponding to an angular velocity of the eye up to $10^\circ/\text{sec}$, and (5) a smaller field of view of 20° . Some of these latter issues were motivated by the studies of Rieger and Toet (1985) and Cutting (1986). Note that with a very large rotational component, the relative difference between the velocities at nearby locations due to the translational component becomes very small, reducing the signal available for recovering the direction of heading.

In the computer simulations, we placed thresholds on both the absolute image velocity and on the velocity differences that were considered detectable. The threshold used for absolute velocity was $1^\circ/\text{sec}$ and the threshold on velocity differences was 10% (see Nakayama (1985) for a review of data on human thresholds). Values falling below these thresholds did not enter into the computation of heading direction. There will be noise in the velocity estimates, but it is not clear what is a reasonable level of noise to expect for the visual system. In the simulations, we initially explored the question of what level of noise in the velocity measurements would yield a heading accuracy of about $2^\circ - 3^\circ$, for the case of translation relative to the ground plane and the overall conditions of the perceptual experiments summarized above. (It is expected that the greater heading accuracy of $1^\circ - 2^\circ$ measured for the human visual system could be obtained by extending the heading computation further in time.) We found that this accuracy could be achieved with an average error in speed of about 25% and average error in the direction of velocity of about 25° . Error was introduced as Gaussian distributed perturbations of the direction and speed of velocity, and the average error was recorded for the actual configurations used in the simulations. An average error in speed of 25% and in velocity direction of 25° was then used throughout the remaining simulations. A small amount of temporal smoothing was performed to reduce the overall sensitivity of the algorithm to this error in the initial velocities. In particular, for each configuration of points, three image frames were generated using a particular set of conditions, the two pairs of adjacent frames were each used to compute the image velocities of the points, noise was added to the two resulting velocity fields, and the two velocity fields were then averaged to yield a final set of image velocities from which the velocity differences were computed.

Although the scene consisted of a single rigid surface in these simulations, we used the strategy described in the previous section for computing the FOE location in the presence of self-moving objects, in order to reduce the sensitivity of the FOE computation to error

in the θ_i estimates. Three circular and overlapping patches representing possible locations of the FOE were centered on heading directions located at 6° , 0° and -6° from straight ahead, and each covered an area of radius 6° . Thus heading angles computed by the algorithm could cover a range from -12° to 12° in the horizontal direction. Preliminary simulations suggested that image patches outside of the regions covered by these three patches yield significantly less positive evidence, and therefore need not be included in this analysis. If more than one patch yielded a predicted FOE location, we first checked whether one estimate was significantly better than the others, in that it had significantly more positive evidence and better fit to the θ_i measurements. If this was not the case, then the multiple predictions were averaged together to yield a final estimate.

The results of a set of simulations with the ground plane are summarized in Table 1. Each data point represents an average of the results obtained from 100 examples (that is, 100 different random configurations of points). The full set of parameters used for the first example (top entry in Table 1) is given in the legend; other entries indicate only the value of the parameter that was different from the first example. Assuming a ground speed for the observer of 1.9m/sec and presentation rate of 15 frames/sec, this corresponds to 0.127 m/frame of observer translation. Similarly, an angular velocity range of $0.3 - 0.7^\circ$ /sec for the simulated eye rotation corresponds to a range of $0.02^\circ - 0.05^\circ$ per frame. This initial range of angular velocities that was used in psychophysical studies is very small. We also conducted simulations with rotations drawn from the range of $5^\circ - 10^\circ$ /sec. The field of view is defined to be the total width of the field in the horizontal direction. For each configuration of points, a simulated heading direction was chosen randomly from the range of 6° to the left and right of straight ahead. For the simulations shown in Table 1, velocity differences were computed for any pair of velocity measurements falling within a neighborhood of 6° of one another.

From this initial set of simulations, it can be seen that direction judgements improve with higher speed of observer translation and higher density of points, and degrade with higher error in the velocity differences and a higher angular velocity of the eye. If the density of points is kept relatively constant, the field of view has little effect on heading accuracy. These factors interact with one another. For example, with the limited field of view, higher angular rotations yield significant degradation in the direction computation, but if the field of view and number of points were increased, a more accurate heading direction could be obtained for higher rotation speeds. Most simulation results reported in the literature use fairly large rotational components, which often yields significant error; such rotations may also yield larger error in human judgements of heading. Overall, the heading accuracy remains high for the range of conditions explored here.

In general, as the velocity difference errors increase, there can be substantial error in the local computations of the dominant orientation of the distribution of velocity differences within image neighborhoods. If these measurements are distributed over a

parameters	horizontal
initial parameters	2.5
7.6 m/sec	2.2
20° field of view, 60 points	2.6
40° field of view, 30 points	4.0
20° field of view, 30 points	2.7
40% average speed error 40° average direction error	3.9
5°–10°/sec rotation range	4.4

Table 1. The results of simulations with the Rieger and Lawton model, applied to images generated by an observer moving along a ground plane. Average errors, in degrees, are given for the horizontal component of heading. The top entry gives results for the following parameters: observer speed of 1.9 m/sec, 40° field of view, 60 points, 6° heading range, 0.3 – 0.7°/sec rotation range, 25% average error in image speed, and 25° average error in the direction of image velocity.

large field, however, the overall computation of the FOE can still be accurate. There is a characteristic asymmetry in the pattern of errors obtained over the visual field. In particular, the directions of the dominant orientation of local velocity differences usually point to the right of the FOE in the right half of the visual field and to the left of the FOE in the left half of the visual field. With a roughly uniform distribution of points in the horizontal direction, these errors effectively cancel one another out in the overall computation of the FOE. The same observation holds true in the vertical direction. An implication of this observation is that if the distribution of θ_i measurements is strongly skewed within the visual field, a characteristic error in the heading computation can result.

The results of some additional simulations with the 3-D cloud and two planes of dots are shown in Table 2. For all of these simulations, the field of view was a square of size 40°, which is somewhat larger than the 40° × 32° field of view used in the perceptual experiments by Warren and his colleagues. The results of simulations with the ground plane suggest that the density of points is a critical factor in determining the accuracy of recovered heading. Because of the somewhat larger field of view used in the simulations here, we used displays of 80 points, rather than 60, in order to keep the density of points similar to that used in the perceptual experiments. The other parameters used in these simulations are listed in the legend for Table 2. Overall, it can be seen that similar heading accuracy can be obtained for the 3-D cloud and two planes. In general, accuracy degrades as absolute depth is increased, but improves as the overall range of depth is increased.

Errors increase slightly when more oblique heading directions are simulated. In general, heading direction is underestimated, in that it is closer to straight ahead relative to the true direction of heading. Again, an increased field of view can reduce the errors for more oblique headings. Errors increase significantly for very sparse patterns containing only 10 points, largely because the image neighborhoods over which the velocity differences are computed contain very few pairs of points from which to compute the θ_i measurements. For the case of the frontoparallel plane, the errors were very large. For headings chosen within a 6° range of directions around straight ahead, the average heading error was 5.0° in the horizontal direction.

parameters	horizontal
3-D cloud, depth range 7-40m	2.3
3-D cloud, depth range 15-32m	4.0
3-D cloud, depth range 7-40m 10 points	5.0
two planes, 5m and 25m	1.5
two planes, 10m and 20m	2.6
two planes, 20m and 40m	3.7
two planes, 5m and 25m $6^\circ - 12^\circ$ heading range	1.8

Table 2. The results of simulations with the Rieger and Lawton model, applied to images generated by an observer moving toward a 3-D cloud of points or two frontoparallel planes separated in depth. Unless specified above, parameters were as follows: observer speed 1.9 m/sec, 40° field of view, 80 points, 6° heading range, $0.3 - 0.7^\circ$ /sec rotation range, 25% average error in image speed, and 25° average error in the direction of image velocity.

Simulations with Self-Moving Objects

This section presents the results of simulations with the algorithm applied to synthetic image sequences containing multiple objects, some of which undergo their own self-motion. For each example, a known velocity field was first generated from a known depth map and movement parameters for the observer and objects. Noise was added to the image velocities, in the form of Gaussian distributed perturbations of their speed and direction. The algorithm was then applied to the noisy velocity field to recover the location of the FOE and to detect self-moving objects.

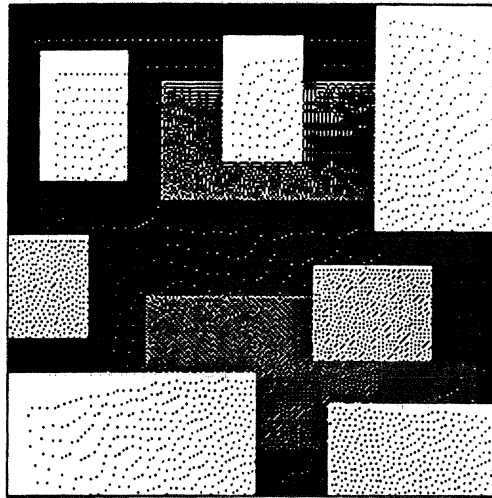


Figure 4. A synthetically generated depth map, with brightness encoding depth (a dithered image is shown, so that the density of black and white dots conveys different brightness levels). Depths range from 75 to 250 units.

A depth map for the scene that formed the basis of these experiments is shown in Fig. 4. Brightness encodes depth, with darker objects located further from the observer. (A dithered image is shown, so that the density of black and white dots conveys different brightness levels.) The scene consists of a set of planar surface patches of different 3-D orientations positioned over a distance of 75–250 units from the observer. From this known depth map and a set of known parameters for the observer's rotation and translation, an image velocity field was computed. An example of an original velocity field is shown in Fig. 5a. The velocities are sampled from an array of size 128×128 . Noise was then added to yield velocity fields such as that shown in Fig. 5b. Before computing the velocity differences, the velocities were then averaged spatially over a neighborhood of size 3×3 pixels, in order to reduce the sensitivity to noise of the subsequent velocity differences.

The distribution of velocity differences was then computed for each image location. The distribution at a given location consisted of the differences in velocity between this location and every other location within a neighborhood of radius 4 pixels. The dominant orientation, θ_i , of this distribution was computed using the scheme described in the previous section, and these θ_i measurements were preserved at locations where the distribution of local velocity differences was strongly anisotropic. For one set of observer and object motion parameters, a map of all the locations at which the θ_i were initially preserved is shown in Fig. 6a. Isolated θ_i measurements that do not belong to a connected patch of at least 10 pixels were then removed, as it was assumed that the most appropriate θ_i estimates to use for the FOE computation would occur in the vicinity of

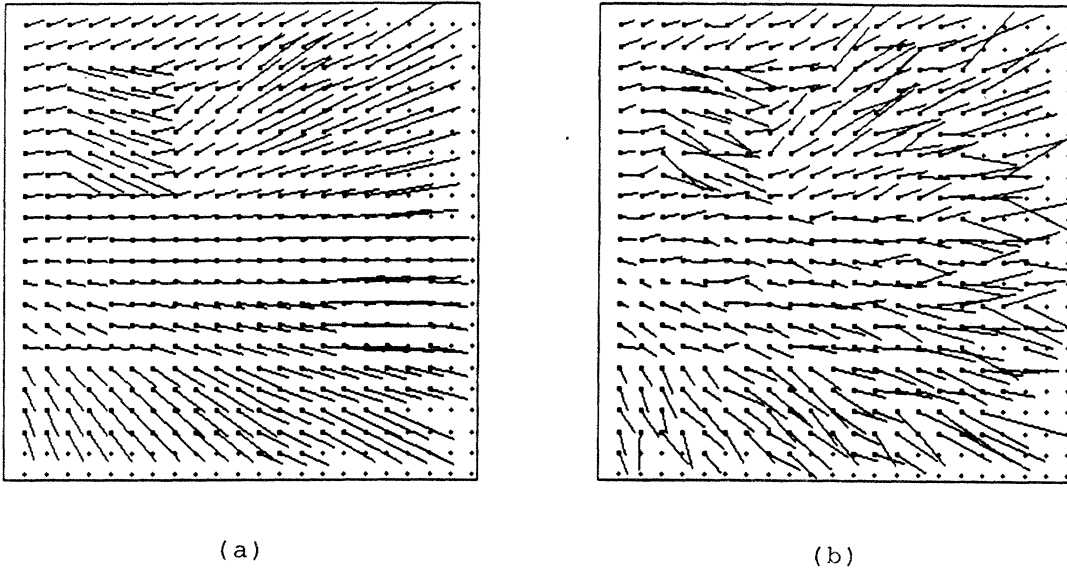


Figure 5. (a) An ideal velocity field obtained from the known depth map shown in Fig. 4 and known observer motion parameters. (b) The velocity field with added noise.

extended boundaries. The locations of the θ_i that remain after this filtering step are shown in Fig. 6b. These remaining measurements are concentrated around the locations of boundaries and over the surface of the object in the upper right corner of the image, which has a large slant. Fig. 6c shows the dominant orientations that are computed at a sample of the image locations. The true FOE is located near the upper right corner of the image, and the two objects highlighted in Fig. 6d are self-moving. It can be seen that there is significant error in the θ_i measurements, as those vectors in Fig. 6c that are not located in the vicinity of the two self-moving objects should, in theory, all point toward the FOE.

To compute the location of the FOE, the image was carved up into overlapping circular patches, as suggested in the previous section. In these simulations, the patches had a radius of 24 pixels and were centered at locations spaced by 24 pixels. For each patch P_j , the set of θ_i measurements yielding positive evidence for the FOE being located within P_j was then determined. If at least 50% of the θ_i measurements yielded positive evidence, a hypothesized FOE was computed from these measurements. If multiple FOE hypotheses emerged, they were reconciled to obtain a single FOE location by considering the extent of the positive evidence in their support, their goodness-of-fit to the computed θ_i , and the proximity of the multiple hypotheses. Fig. 7 shows the true (solid circles) and computed (open circles) FOE locations for 6 different choices of the observer translation parameters, and for rotation parameters, $= (w_x, w_y, w_z) = (0.0, 2.0, 0.0)$ (these rotation

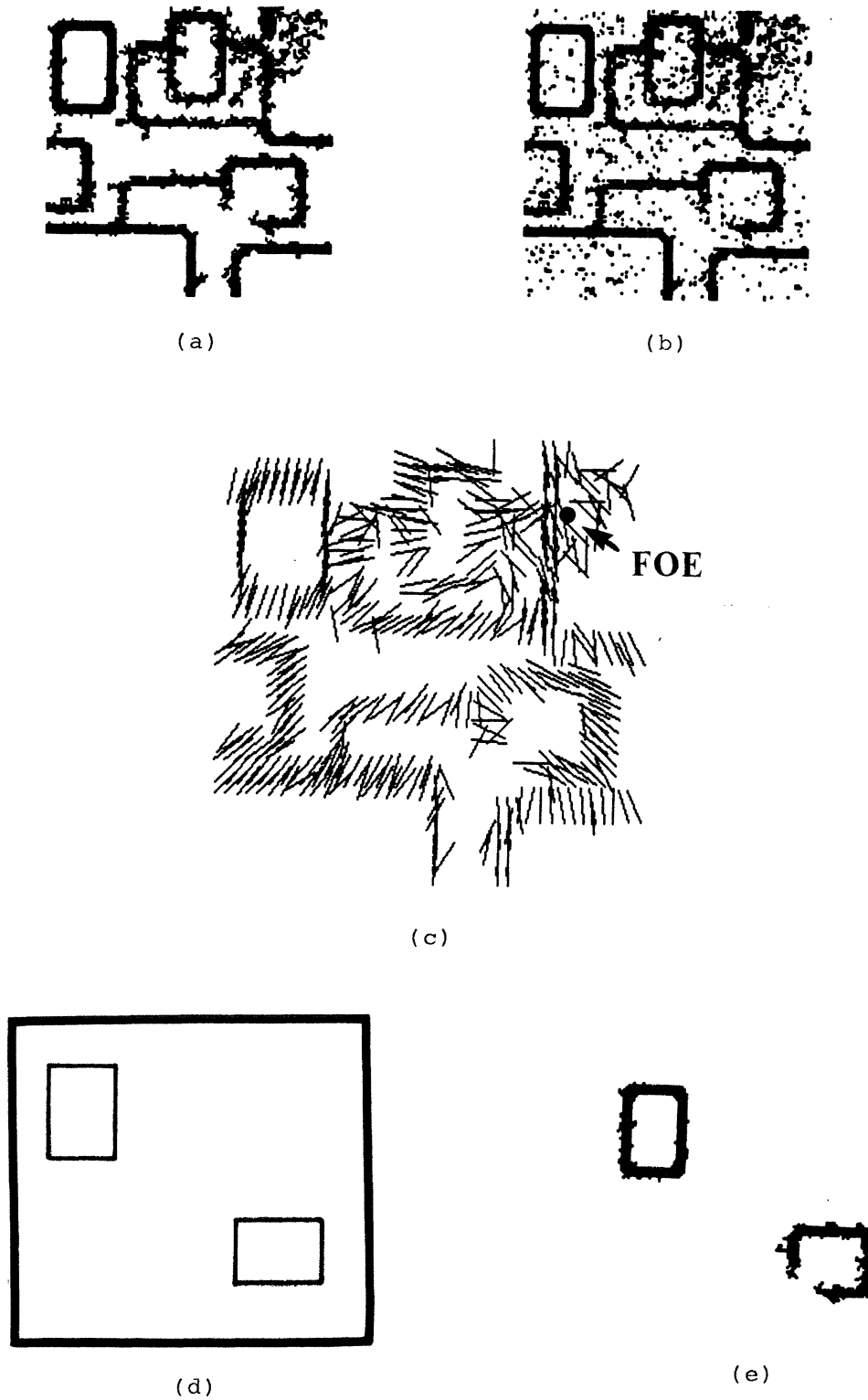


Figure 6. (a) A map of all the locations where θ_i were derived from local velocity difference distributions with strong anisotropy. (b) Isolated θ_i measurements are removed. (c) A sampling of the dominant orientations, θ_i . The true FOE is located in the upper right corner. (d) The locations of two objects in the scene that are self-moving. (e) Locations where θ_i measurements were obtained that indicate self-moving objects.

parameters were used to generate the velocity fields shown in Fig. 5). The error in the final FOE estimates is small, given the large error in the input velocity fields and the θ_i estimates.

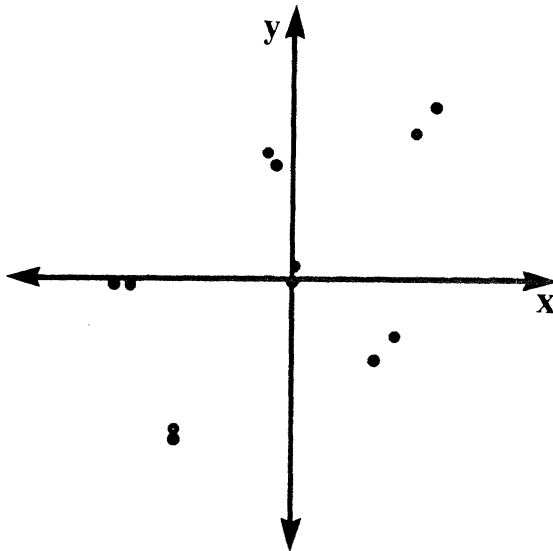


Figure 7. True FOE locations (solid circles) are compared to the FOE locations derived from the algorithm (open circles) for six choices of the observer translation parameters. The full extent of the horizontal and vertical axes correspond to an image distance of 128 pixels.

Once an initial estimate for the FOE location was obtained, extended regions yielding negative evidence were isolated as possibly indicating self-moving objects. For the example shown in Fig. 6, the patch that yielded the most positive evidence is located in the upper right corner of the image. The θ_i measurements that were not directed toward this patch were isolated, and extended, connected groups of such measurements were hypothesized to correspond to self-moving objects. Fig. 6e shows the final locations found to arise from self-moving objects, which correspond correctly to the two self-moving objects in the scene.

SUMMARY AND CONCLUSIONS

This paper first briefly considered the computation of three critical properties for low-level navigation tasks: (1) the 3-D direction of heading of an observer relative to object surfaces, (2) the time-to-collision between an observer and an approaching surface, and (3) the locations of object boundaries defined by discontinuities in image motion. We argued that these three properties are essential for tasks that require rapid sensing and response, and ultimately should be considered together in a system capable of performing such tasks effectively. We then focused on the computation of the 3-D direction of trans-

lation of an observer relative to object surfaces. Consideration of perceptual observations regarding the human recovery of heading direction and existing computational models led us to examine the model proposed by Rieger and Lawton (1985) in more detail. We explored some extensions to the Rieger and Lawton model that yield improvement of its performance in the presence of error in the image motion measurements and allow it to cope with scenes containing multiple moving surfaces. The results of computer simulations with this modified model applied to visual patterns similar to those used in perceptual studies suggest that it exhibits much of the basic behavior of the human system.

The style of model developed here was chosen also because it fits into the overall framework that we are pursuing for the visual processing mechanisms that underlie low-level navigation. We argued that because of the demands of navigational tasks requiring rapid sensing and response, the human visual system may use specialized routines that use only partial or qualitative information regarding motion in the image or in the scene that can be computed reliably with minimal computation, and which is critical to performing a specific task. In the model presented here, simple measurements of velocity differences within local image neighborhoods are used to compute only the direction of observer heading, independent of the observer's rotation or scene layout. Velocity differences in regions of significant depth variation provide a direct cue to the observer's heading that can be exploited with relatively little computation. This partial information about heading direction can then be used directly by routines that detect potential collisions or track objects in the scene. Furthermore, because velocity differences will be significant along discontinuities in depth that occur along the boundaries of stationary and self-moving objects, they can also be exploited to detect these boundaries. We have shown that the heading computation itself can embody a strategy for detecting the boundaries of self-moving objects. This boundary information can also be used by routines that detect potential collisions, to determine the overall size and shape of relevant objects. Once an object has been isolated in the scene, the rate of change of the size of its bounding contour can be used to assess its time-to-collision (for example, Lee, 1980; Todd, 1981).

A number of additional questions regarding the human perception of heading direction arise from the analysis of the model presented here, which can be explored through further perceptual experiments. Among these are the following:

- Does accuracy in judging heading direction decrease with more oblique headings, and is there a general tendency to underestimate oblique headings? Is the size of the field of view more critical for the accurate judgement of oblique headings?
- Is there degradation of heading judgements when larger angular rotations are simulated, and is the size of the field of view critical in this case?
- Does an asymmetric spatial distribution of points yield characteristic errors in heading judgements, as suggested by the simulations?

- Is there a systematic degradation in heading accuracy with a smaller depth range and larger absolute depth?

It would also be useful to examine the accuracy of our judgement of the vertical component of heading direction, in order to assess our overall precision at performing the heading computation. Other experimental questions arise regarding the recovery of observer heading in the presence of motion discontinuities and self-moving objects. In particular:

- What is the effect of self-moving objects in the field of view on the accuracy of heading judgements?
- Is there any difference in performance, depending on whether the boundaries of a self-moving object yield immediately perceptible motion discontinuities?
- How much of the image must contain significant depth variation? Suppose, for example, that the image contains a single object (a small frontoparallel plane) in front of a larger frontoparallel plane in the background. How large must the closer object be, and how much does it need to be separated in depth from its background, in order to yield accurate heading judgements?
- How much deviation in direction of image motion must a self-moving object undergo, relative to the motion direction expected from the observer's motion alone, in order to detect its presence?

Further experimental work that addresses these questions is critical to assessing the appropriateness of a model of the type explored here as a description of the recovery of heading direction by the human system.

Acknowledgement: I thank Shimon Ullman and Eric Grimson for valuable comments on a draft of this paper.

References

- Adiv, G. (1985). Determining three-dimensional motion and structure from optical flow generated by several moving objects. *IEEE Trans. Patt. Anal. Machine Intell., PAMI-7*, 384-401.
- Aloimonos, J. (1990). Purposive and qualitative active vision. *Proc. AAAI Workshop on Qualitative Vision*, Boston, 1-5.
- Aloimonos, J. & Brown, C. M. (1989). On the kinetic depth effect. *Biol. Cybern.*, *60*, 445-455.
- Aloimonos, J., Weiss, I. & Bandopadhyay, A. (1988). Active vision. *Int. J. Comp. Vis.*, *1*, 333-356.
- Anderson, C. H., Burt, P. J. & van der Wal, G. S. (1985). Change detection and tracking using pyramid transform techniques. *Proc. SPIE Conf. on Intelligent Robots and Computer Vision*, Boston, MA, 300-305.
- Ballard, D. H. & Brown, C. M. (1982). *Computer Vision*. Englewood Cliffs, NJ: Prentice-Hall.
- Ballard, D. H. & Kimball, O. A. (1983). Rigid body motion from depth and optical flow. *Comp. Vis. Graph. Image Proc.*, *22*, 95-115.
- Bouthemy, P. & Lalande, P. (1990). Detection and tracking of moving objects based on a statistical regularization method in space and time. *Proc. First European Conf. Comp. Vision*, O. Faugeras (ed.), Antibes, France, Berlin: Springer-Verlag, 307-311.
- Broida, T. J. & Challappa, R. (1986). Estimation of object motion parameters from noisy images. *IEEE Trans. Patt. Anal. Machine Intell., PAMI-8*, 90-99.
- Brooks, R. A. (1986). A robust layered control system for a mobile robot. *IEEE J. Robotics and Automation, RA-2*, April, 14-23.
- Brooks, R. A., Flynn, A. M. & Marill, T. (1987). Self calibration of motion and stereo vision for mobile robot navigation. *MIT Artif. Intell. Lab. Memo 984*.
- Bruss, A. R. & Horn, B. K. P. (1983). Passive navigation. *Comp. Vis. Graph. Image Proc.*, *21*, 3-20.
- Burger, W. & Bhanu, B. (1990). Estimating 3-D egomotion from perspective image sequences. *IEEE Trans. Patt. Anal. Machine Intell., PAMI-12*, 1040-1058.
- Burt, P. J., Bergen, J. R., Hingorani, R., Kolczinski, R., Lee, W. A., Leung, A., Lubin, J. & Shvaytser, H. (1989). Object tracking with a moving camera, an application of dynamic motion analysis. *Proc. IEEE Workshop on Visual Motion*, Irvine, CA, March, 2-12.
- Carel, W. L. (1961). Visual factors in the contact analog. Ithaca, NY: General Electric Advanced Electronics Center Pub. R61 ELC60, 1-65.
- Cutting, J. E. (1986). *Perception with an Eye Towards Motion*. Cambridge: MIT Press.

- Dinstein, I. (1988). A new technique for visual motion alarm. *Patt. Recog. Letters*, 8, 347.
- Enkelmann, W. (1990). Obstacle detection by evaluation of optical flow fields from image sequences. *Proc. First European Conf. Comp. Vis.*, O. Faugeras (ed.), Antibes, France, Berlin: Springer-Verlag, 134-138.
- Faugeras, O. D., Lustman, F. & Toscani, G. (1987). Motion and structure from motion from point and line matches. *Proc. First Int. Conf. on Comp. Vision*, London, June, 25-34.
- Frazier, J. & Nevatia, R. (1990). Detecting moving objects from a moving platform. *Proc. DARPA Image Understanding Workshop*, Pittsburgh, PA, San Mateo: Morgan Kaufmann, 348-355.
- Heeger, D. J. & Hager, G. (1988). Egomotion and the stabilized world. *Proc. 2nd Int. Conf. Comp. Vision*, Tampa, FL, 435-440.
- Heeger, D. J. & Jepson, A. (1990). Visual perception of three-dimensional motion. *MIT Media Lab. Tech. Rep.*, 124.
- Heel, J. (1990a). Dynamical motion vision. *Robotics and Autonomous Systems*, 6(1).
- Heel, J. (1990b). Direct estimation of structure and motion from multiple frames. *MIT Artif. Intell. Lab. Memo*, 1190.
- Hildreth, E. C. (1984). *The Measurement of Visual Motion*, Cambridge: MIT Press.
- Hildreth, E. C. & Koch, C. (1987). The analysis of visual motion: From computational theory to neuronal mechanisms. *Ann. Rev. Neurosci.*, 10, 477-533.
- Horn, B. K. P. & Weldon, E. J. (1988). Direct methods for recovering motion. *Int. J. Comp. Vis.*, 2, 51-76.
- Huber, P. J. (1981). *Robust Statistics*. New York: John Wiley & Sons.
- Hutchinson, J., Koch, C., Luo, J. & Mead, C. (1988). Computing motion using analog and binary resistive networks. *IEEE Computer*, 21, 52-63.
- Jain, R. C. (1983). Direct computation of the focus of expansion. *IEEE Trans. Patt. Anal. Machine Intell.*, PAMI-5, 58-63.
- Jain, R. C. (1984). Segmentation of frame sequences obtained by a moving observer. *IEEE Trans. Patt. Anal. Machine Intell.*, PAMI-6, 624-629.
- Jain, R., Martin, W. N. & Aggarwal, J. K. (1979). Extraction of moving object images through change detection. *Proc. Sixth Int. Joint Conf. Artif. Intell.*, Tokyo, 425-428.
- Jain, R., Militzer, D. & Nagel, H. H. (1977). Separating non-stationary from stationary scene components in a sequence of real world TV images. *Proc. Fifth Int. Joint Conf. Artif. Intell.*, Cambridge, MA, 425-428.
- Johnston, I. R., White, G. R. & Cumming, R. W. (1973). The role of optical expansion patterns in locomotor control. *J. Exp. Psych.*, 86, 311-324.
- Koenderink, J. J. & Van Doorn, A. J. (1976). Local structure of movement parallax of the plane. *J. Opt. Soc. Am.*, 66, 717-723.

- Lawton, D. T. (1983). Processing translational motion sequences. *Comp. Graph. Image Proc.*, *22*, 116–144.
- Lee, D. N. (1974). Visual information during locomotion. *Perception: Essays in Honor of James Gibson*, R. B. Macleod, H. Pick (eds.), Ithaca, NY: Cornell Univ. Press.
- Lee, D. N. (1976). A theory of visual control of braking based on information about time-to-collision. *Perception*, *5*, 437–459.
- Lee, D. N. (1980). The optic flow field: The foundation of vision. *Phil. Trans. Roy. Soc. Lond. B*, *290*, 169–179.
- Lee, D. N., Lishman, J. R. & Thomson, J. A. (1982). Regulation of gait in long jumping. *J. Exp. Psych: Human Perc. Perf.*, *8*, 448–459.
- Lee, D. N. & Reddish, P. E. (1981). Plummeting gannets: a paradigm of ecological optics. *Nature*, *293*, 293–294.
- Lee, D. N., Young, D. S., Reddish, P. E., Lough, S. & Clayton, T. M. H. (1983). Visual timing in hitting an accelerating ball. *Quart. J. Exp. Psych.*, *35A*, 333–346.
- Llewellyn, K. R. (1971). Visual guidance of locomotion. *J. Exp. Psych.*, *91*, 245–261.
- Longuet-Higgins, H. C. (1981). A computer algorithm for reconstructing a scene from two projections. *Nature*, *293*, 133–135.
- Longuet-Higgins, H. C. (1984). Visual ambiguity of a moving plane. *J. Opt. Soc. Am. A*, *1*, 1215.
- Longuet-Higgins, H. C. & Prazdny, K. (1981). The interpretation of moving retinal images. *Proc. R. Soc. London Ser. B*, *208*, 385–397.
- McLeod, R. W. & Ross, H. E. (1983). Optic-flow and cognitive factors in time-to-collision estimates. *Perception*, *12*, 417–423.
- Meer, P., Mintz, D., Kim, D. Y. & Rosenfeld, A. (1991). Robust regression methods for computer vision: A review. *Int. J. Comp. Vision*, *6*, 59–70.
- Mutch, K. M. & Thompson, W. B. (1985). Analysis of accretion and deletion at boundaries in dynamic scenes. *IEEE Trans. Patt. Anal. Mach. Intell.*, *PAMI-7*, 133–138.
- Nakayama, K. (1985). Biological motion processing: A review. *Vision Res.*, *25*, 625–660.
- Negahdaripour, S. & Horn, B. K. P. (1987). Direct passive navigation. *IEEE Trans. Patt. Anal. Machine Intell.*, *PAMI-9*, 168–176.
- Negahdaripour, S. & Horn, B. K. P. (1989). A direct method for locating the focus of expansion. *Comp. Vis. Graph. Image Proc.*, *46*, 303–326.
- Nelson, R. C. (1990). Qualitative detection of motion by a moving observer. *Univ. Rochester Comp. Science Tech. Rep. 341*, April.
- Prazdny, K. (1980). Egomotion and relative depth map from optical flow. *Biol. Cyber.*, *36*, 87–102.
- Regan, D. M. & Beverley, K. I. (1982). How do we avoid confounding the direction we are looking and the direction we are moving? *Science*, *215*, 194–196.

- Reichardt, W. & Poggio, T. (1980). Figure-ground discrimination by relative movement in the visual system of the fly. Part I: Experimental results. *Biol. Cybern.*, *35*, 81-100.
- Rieger, J. H. & Lawton, D. T. (1985). Processing differential image motion. *J. Opt. Soc. Am. A*, *2*, 354-360.
- Rieger, J. H. & Toet, L. (1985). Human visual navigation in the presence of 3D rotations. *Biol. Cybern.*, *52*, 377-381.
- Rousseeuw, P. & Leroy, A. (1987). *Robust Regression and Outlier Detection*. New York: John Wiley & Sons.
- Schiff, W. & Detwiler, M. L. (1979). Information used in judging impending collision. *Perception*, *8*, 647-658.
- Schunck, B. G. (1986). The motion constraint equation for optical flow. *Proc. Int. J. Conf. Patt. Recog.*, 20-22.
- Shariat, H. (1986). The motion problem: How to use more than two frames. PhD. thesis, Dept. Elec. Eng., Univ. Southern Calif.
- Simpson, W. A. (1988). Depth discrimination from optic flow. *Perception*, *17*, 497-512.
- Spoerri, A. N. & Ullman, S. (1987). The early detection of motion boundaries. *MIT Artif. Intell. Lab. Memo*, 935.
- Subbarao, M. (1988). Interpretation of visual motion: A computational study. *Research Notes in Artificial Intelligence*, San Mateo: Morgan Kaufmann.
- Thompson, W. B., Lechleider, P. & Stuck, E. R. (1991). Detecting moving objects using the rigidity constraint. *IEEE Trans. Patt. Anal. Machine Intell.*, in press.
- Thompson, W. B., Mutch, K. M. & Berzins, V. (1985). Dynamic occlusion analysis in optical flow fields. *IEEE Trans. Patt. Anal. Machine Intell.*, *PAMI-7*, 374-383.
- Thompson, W. B. & Pong, T. C. (1990). Detecting moving objects. *Int. J. Comp. Vis.*, *4*, 39-57.
- Todd, J. T. (1981). Visual information about moving objects. *J. Exp. Psych.: Human Perc. Perf.*, *7*, 795-810.
- Tsai, R. Y. & Huang, T. S. (1984a). Estimating three-dimensional motion parameters of a rigid planar patch: III. Finite point correspondences and the three-view problem. *IEEE Trans. Acoust. Speech Signal Proc.*, *ASSP-32*, 213-220.
- Tsai, R. Y. & Huang, T. S. (1984b). Uniqueness and estimation of three-dimensional motion parameters of rigid objects with curved surfaces. *IEEE Trans. Patt. Anal. Machine Intell.*, *PAMI-6*, 13-27.
- Ullman, S. (1979). *The Interpretation of Visual Motion*. Cambridge: MIT Press.
- Ullman, S. (1984). Maximizing rigidity: the incremental recovery of 3-D structure from rigid and rubbery motion. *Perception*, *13*, 255-274.
- Verri, A., Giroso, F. & Torre, V. (1989). Mathematical properties of the two-dimensional motion field: from singular points to motion parameters. *J. Opt. Soc. Am. A*, *6*, 698-712.

- Warren, W. H., Griesar, W., Blackwell, A. W., Kalish, M. & Hatsopoulos, N. G. (1990). On the sufficiency of the velocity field for perception of heading. *Manuscript in preparation*.
- Warren, W. H. & Hannon, D. J. (1988). Direction of self-motion is perceived from optical flow. *Nature*, 336, 162-163.
- Warren, W. H. & Hannon, D. J. (1990). Eye movements and optical flow. *J. Opt. Soc. Am. A*, 7, 160-169.
- Warren, W. H., Morris, M. W. & Kalish, M. (1988). Perception of translational heading from optical flow. *J. Exp. Psych.: Human Perc. Perf.*, 14, 646-660.
- Waxman, A. M. & Ullman, S. (1985). Surface structure and 3D motion from image flow: a kinematic analysis. *Int. J. Robotics Res.*, 4, 72-94.
- Waxman, A. M. & Wohn, K. (1988). Image flow theory: A framework for 3-D inference from time-varying imagery. In *Advances in Computer Vision*, C. Brown, (ed.), New Jersey: Erlbaum.
- Weng, J., Huang, T. S. & Ahuja, N. (1989). Motion and structure from two perspective views: Algorithms, error analysis and error estimation. *IEEE Trans. Patt. Anal. Machine Intell.*, PAMI-11, 451-476.
- Wohn, K. & Waxman, A. M. (1990). The analytic structure of image flows: Deformation and segmentation. *Comp. Vision Graphics image Proc.*, 49, 127-151.
- Zhang, Z., Faugeras, O. D. & Ayache, N. (1988). Analysis of a sequence of stereo scenes containing multiple moving objects using rigidity constraints. *Proc. 2nd Int. Conf. Comp. Vision*, Tampa, FL, 177-186.