

Stereoscopic Tracking of Bodies in Motion

Roberto CIPOLLA * Masanobu YAMAMOTO

Robotics Research Group, University of Oxford Electrotechnical Laboratory, Japan

In this paper we present a fast, highly efficient and robust visual tracking process for multiple moving objects, using stereo image sequences taken from a stationary camera.

*The algorithm assumes that object motion is restricted to a horizontal plane (for example motion of cars on roads or humans walking). Dense stereo image sequences and the **Visualised Locus** method [1] [2] (in which each image sequence is first sampled to produce a 2D spatio-temporal cross-section image) are used to ensure temporal correspondence without search. Edge segments in the left and right spatio-temporal images are then matched. Additional stereo matching constraints are derived by using motion and temporal continuity to reduce the number of ambiguous matches. Speed is achieved by only processing a single spatio-temporal cross-section image from each image sequence.*

The algorithm succeeds in tracking objects in space and time moving against arbitrarily complex backgrounds and in the presence of occlusion, disappearance and reappearance of object features. The output of the algorithm is the 3D position of moving objects as a function of time.

In real dynamic scene analysis the considerable changes in image structure that can occur as a result of object motion (such as the occlusion, disappearance and reappearance of object features) are major obstacles to the successful application of existing "Shape from Stereo" [3] and "Shape from Monocular Motion" [4] techniques to the problem of 3D object tracking from visual data. Monocular image sequence analysis, for example, has the following inherent difficulties: the **Temporal Correspondence** problem in obtaining optical flow locally [4] or matching tokens in discrete views; a **Speed-Scale ambiguity** which makes it impossible to determine 3D structure and motion in absolute terms for a monocular observer viewing unfamiliar object; and a restriction to rigid body motions which usually require the segmentation of images into parts corresponding to objects with the same rigid body motion. These methods perform poorly with respect to **accuracy**, **sensitivity to noise**, and **robustness** in the face of errors. This is because it is difficult to estimate optical flow accurately [5], or to extract the position of feature points such as corners in the image. In particular motion of objects towards the camera (motion in depth) is sensitive to noise since the image velocities are small. Also features can not lie in a plane (and some other special 2nd de-

gree surfaces) since coplanar points lead to a degenerate system of equations [6]. In the image sequence analysis of human motion, for example, the features of interest such as the limbs are nearly coplanar. Human motion is not strictly rigid body motion. An additional practical problem is that most applications of visual tracking require the real - time processing of large volumes of data. For most existing algorithms this requires special purpose hardware.

Stereo vision can be used independently of motion analysis in a dynamic environment to determine the trajectory of an object in space by taking successive stereo snapshots, determining object locations at each instant and combining these locations into a trajectory [7], [8]. The most difficult part of the processing concerns the Correspondence problem: what to match and how to match it [3].

The correspondence problem is particularly difficult in the presence of occluding boundaries and semi-transparent surfaces such as fences or windows. With dynamic scenes the disappearance and reappearance of object and image features may make matching and the interpretation of 3D structure instantaneously ambiguous or impossible. An additional drawback in repeating the stereo vision processing at each time instant is that large amounts of data processing (both time and volume) are required. A considerable saving in data processing can be achieved by exploiting knowledge of motion. Ideally only the edges of objects in motion need be matched and the depth map can be updated. Of greater interest however is whether *motion* can interact with stereopsis at the level of matching to help disambiguate false matches. Poggio and Poggio [3] note that image motion may be able to aid stereo in the matching process.

Motion and Stereo Fusion

Both stereo vision and structure from motion have typically been treated as separate parallel processes. However each has inherent difficulties and so it seems logical to attempt to combine them into a single system. Jenkin and Tsotos [9] and Waxman and Duncan [10] have attempted to unify stereo and motion analysis in a manner which helps to overcome the others shortcomings.

Jenkin and Tsotos [9] describe a stereo vision system which will track special extracted object feature points in 3D space over time. It uses the 3D interpretation of the feature point velocities to help in the stereo matching process. Even with a sparse feature set (extracted with

*Roberto Cipolla acknowledges the support of The IBM UK Science Centre.

the Moravec interest operator [7]) the algorithm requires an extensive search in both time and space and this makes it unsuitable for real-time processing.

Waxman and Duncan [10] have proposed an integrated stereo-motion analysis beginning with the determination of image flow and using correlation between relative image flow (*binocular difference flow*) and stereo disparity locally in establishing correspondence. However this method suffers from the problems of estimating optical flow and it requires finely textured smooth surfaces which can be approximated locally as planar. Its application to long real image sequences has not been tested.

As mentioned above temporal and stereo correspondence are severe problems and limit the usefulness of existing algorithms to applications of visual tracking. An additional practical problem is the need for real time processing. If, however, we assume a restricted class of motions and process a dense sequence of images so that temporal continuity is guaranteed (as in Bolles' Epipolar-Plane Image Analysis [11]) it is possible to greatly simplify (and in a special case avoid) the temporal correspondence problem. If in addition we use motion and temporal continuity as additional matching constraints, the stereo correspondence problem can also be considerably simplified.

In this paper we present a fast, highly efficient and robust visual tracking process for multiple moving objects, from stereo image sequences taken from a stationary camera. Objects moving against arbitrarily complex backgrounds and in the presence of occlusion, disappearance and reappearance of object features are tracked in space and time.

The algorithm presented assumes that object motion is restricted so that its height above a horizontal plane is constant (for example motion of cars on roads, humans walking etc). We show that dense stereo image sequences and the **Visualised Locus** method [1] [2] (in which each image sequence is first sampled to produce 2D spatio-temporal cross-section images) can be used to ensure temporal correspondences automatically without search. Additional stereo matching constraints are derived by using motion and temporal continuity to reduce the number of ambiguous matches. Computational speed is achieved by only processing a single spatio-temporal cross-section image from each image sequence.

Temporal Correspondence: The Visualised Locus

Theory

Consider the perspective projection of an object point onto a planar screen normal to the z -axis. We model the camera as a pin hole with centre at $C(0, H, f)$ and with focal length f . A point $P(X, Y, Z)$ on an object will be projected onto a point $p(x, y)$ on the image plane such

that (figure 1a):

$$\frac{X}{x} = \frac{Y - H}{y} = \frac{Z}{f} \quad (1)$$

If the object moves the projected point's position is a function of time, $p(x, y, t)$:

$$x(t) = \frac{fX(t)}{Z(t)} \quad (2)$$

and in the image y direction:

$$y(t) = \frac{f(Y(t) - H)}{Z(t)} \quad (3)$$

If a sequence of images is taken in rapid succession and piled up sequentially with time, we can construct a 3-dimensional spatio-temporal image (figure 1b). If temporal continuity from image to image is ensured the image of object point P , $p(x, y, t)$, in *general* forms a 3D locus in this 3D image. This locus is referred to as the **Visualised Locus** since it is the locus of the projection as a function of time and if the object point P is occluded or goes out of the view of the camera the image point p disappears [1] [2].

In the case in which the motion of P is constrained to move on a plane $Y(t) = H$ (ie. at the same height as the optical centre), the visualised locus lies on a plane in the 3D image. For this *special* case, the projection of a point in space is constrained to the same single scan line (raster) of the image at all times. The 2D image defined by this plane is a spatio-temporal image. It is a $x-t$ cross-section of the 3D image at $y = 0$. It can be synthesized by storing the scan line containing the optical centre ($y = 0$) of each image of the sequence and arranging them sequentially in order of time (figure 3).

Other researchers have constructed similar images: [12], [13], [11]. As Bolles et al noted, even though the spatial images which were used to construct it contain complex shapes and intensity changes (figure 2), the spatio-temporal image as a consequence of smooth motion is composed of simpler image structures, regions and edge segments (figure 4).

Higher level properties can also be inferred from the synthesized image. The relationship between loci of points on different objects can be used to determine their relative motions. The disappearance and reappearance of a locus indicates occlusion and by looking at neighbouring loci, the occluding object can be determined.

Unlike Bolles' Epipolar-Image analysis, however, the spatio-temporal cross-section images used in this method are generated from a stationary camera with multiple moving objects and non-linear motion. Accurate knowledge of the motion is not required. The original contribution of this paper is the stereo matching of edge segments between spatio-temporal cross-section images generated from the left and right camera image sequences.

Extraction of Visualised Loci

The loci of points on the boundaries of regions with high contrast in an image appear as edges in the synthesized spatio-temporal cross-section image. The extraction of these loci thus involve procedures for edge detection and enhancement; the fitting of line and curve segments to these edges; and the merging of segments belonging to the same locus based on temporal continuity. These algorithms are described in [2]. The latter step is required because due to occlusion, or camouflage (object has same intensity as background) and noise the extracted loci may be fragmented. Where possible the fragments are linked by linear interpolation across "missing" segments. This procedure is based on proximity of the fragment end points and attempting to ensure smooth continuous loci. This latter step is not always possible. It is not essential but as mentioned below the disambiguating power of the proposed matching algorithm increases with the length of the extracted loci.

The visualised loci of points on stationary objects appear as straight edges with no time gradient in the synthesized cross-section image . It is therefore very easy to distinguish between the loci of stationary objects and moving objects, regardless of the complexity of the background.

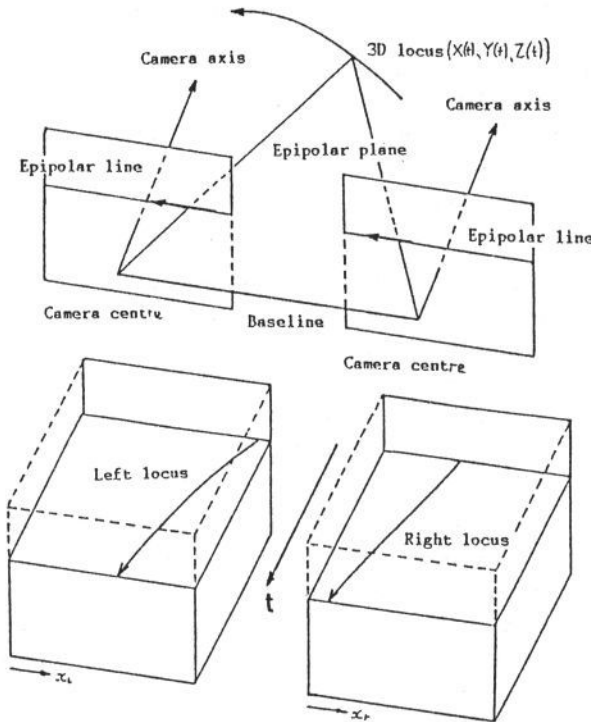


Figure 1. Stereo visualised locus method:

a) camera geometry b) 3D spatio-temporal images with cross-sections showing visualised loci

Stereo Correspondence Stereo Visualised Locus Method

If we observe a scene with a stereo pair of cameras and synthesise a spatio-temporal 3D image for both the left and right image sequences there exist 2 visualised loci corresponding to the same point $P : p_l[x_l(t), y_l(t)]$ and $p_r[x_r(t), y_r(t)]$. If these left and right loci can be correctly matched (stereo correspondence found) the depth as a function of time, $Z(t)$, and hence the position ($X(t)$ and $Y(t)$) can be determined for a calibrated camera set by triangulation.

For the special case in which the motion of the body is constrained to a horizontal plane and the cameras are at the same height, the 2 visualised loci of an object point at the same height as the optical centres are constrained to a single plane in the 3D block of data (figure 1). The visualised loci ($x_l(t)$ and $x_r(t)$) can then be automatically extracted from the 2D spatio-temporal cross-section image for each image sequence.

The dynamic stereo correspondence problem is then to match the visualised loci in the left spatio-temporal cross-section image with those in the right image. The stereo matching algorithm is an extension of edge-based techniques to edges with time as an additional dimension.

Parallel camera geometry is not necessary. Any camera geometry can be used that ensures the camera scan lines are horizontal (i.e. camera tilt and convergence are allowed). This is because the epipolar plane for an object point at the height of the optical centres will always be horizontal.

Search for correspondence

For cameras with parallel optical axes, the relative position of corresponding points in the left and right image are geometrically constrained by the **epipolar constraint** and by:

$$x_{pl} \geq x_{pr} \quad (4)$$

which is the **ordering constraint** relative to a point at infinity. This constraint expresses the portion of the epipolar line in the other image on which a potential match may be found. An additional constraint is that of **uniqueness**: each matching primitive should match at most one primitive from the other image.

In the dynamic stereo problem the epipolar constraint is used to generate the left and right spatio-temporal cross-section images – corresponding rasters are epipolar lines for a given time (figure 3). The matching primitives are portions of the visualised loci and we extend the above conventional constraints to encompass **temporal continuity**. Namely instead of matching at a single time (i.e. matching edge pixels along corresponding rasters of the cross-section images), we apply the matching constraint to the visualised locus for all time. The search for corresponding points in left and right images is also con-

siderably reduced by only considering the loci of moving objects.

The search initially finds portions of loci in the left and right images which co-exist (overlap) in time (epipolar constraint); which satisfy the constraint of equation (4) for *all* common times and which have similar x-direction profiles. Multiple correspondences are reduced by ensuring that each locus in the left image has at least one corresponding candidate for matching in the right image (uniqueness). If ambiguous matches still exist pairs are chosen which have maximum time overlap and whose disparity changes smoothly with time.

These matching constraints were sufficient for our experiments. Unambiguous correspondence is not guaranteed.

Although it is possible for the matching strategy to be used on static edges there is no advantage over existing stereo algorithms in using the proposed matching algorithm for static edges. In fact the algorithm's disambiguating power will be poor with static edges since in this case it is strictly matching along epipolar lines only and it will be inferior to algorithms which use figural continuity [14] or local support [15].

Results

Eight frame samples of an office scene containing 2 people moving along the floor are shown in figure 2. The scene was observed for 8.5s (512 images (fields) at video rate) with a stereo pair of calibrated TV cameras [16] with a long baseline of 0.340 m and optical centres at 1.2m. One person is occluded for part of the observation period. 512×512 spatio-temporal cross-section image are produced from 1 scan line (average of a swath of 7 lines for robustness) from each image for both left and right camera image sequences (figure 4). This operation can be performed in real time. For both left and right cross-section images: Edge detection (cpu time 2-3s); Segment labelling and curve fitting (cpu time 10s); and the extraction and description of the visualised loci (cpu time 3s) are performed on a Sun-3/260 workstation (figure 5). The extraction includes merging fragments of loci across "missing" edges to give long continuous loci (shown as dashed lines). Correspondence and determination of depth was then carried out. The algorithm's output is a plan view of the 3D locus (X-position and depth, Z, Y-position assumed fixed) as a function of time. This is shown in figure 6. Solid lines correspond to visible features. The dashed lines correspond to features which disappear and reappear as a consequence of occlusion or camouflage or noise.

Additional examples of stereo image sequences of office scenes and the tracking algorithm output are presented in figure 7 and figure 8 (output) for a scene in which one person overtakes the other and figure 9 and 10 for a person moving towards the camera in front of a parallel bar fence.

In the last example the object features are continually

disappearing and reappearing and the loci are modulated by an approximately periodic function as a consequence of arm motion modulating the body width as seen by the cameras. At any time instant existing static stereo matching algorithms would fail to find the correct correspondence or produce ambiguous results. The algorithm presented successfully tracks the moving object behind the fence by using temporal continuity to predict where hidden points are based on the object's motion history.

Conclusions

A method using stereo image sequences to automatically detect and track the 3D motion of objects has been presented. It overcomes the problems of large data storage and processing by processing the single spatio-temporal cross-section image of the image sequence which contains the visualised loci (the locus of image points in time) of object points at the same height as the optical centres of the cameras. This is possible if the motion is restricted to a horizontal plane and makes unnecessary the need for searching for temporal correspondence. Stereo correspondence between left and right camera images is also greatly simplified by using motion as a cue to suppress the background and by modifying the stereo matching constraints to encompass temporal continuity.

The algorithm succeeds in tracking objects in space and time moving against arbitrarily complex backgrounds and in the presence of occlusion, disappearance and reappearance of object features. The correspondence techniques used establish matches using only partial information and make predictions where invisible (hidden) points are given their past motion histories and the motion of their visible neighbours. Stereo matching in the presence of partly transparent objects in the foreground (for example fences with parallel bars) has been demonstrated.

The algorithm is robust to calibration errors, image noise and deviations from perfect rigidity (eg. presence of extremal boundaries of curved surfaces or the motion of human body). It also works well with large stereo baselines and for long image sequences. It was tested on a variety of scenes to track the 3D motions of humans moving along the ground in various directions relative to the camera.

The motion tracking scheme presented will only correctly track points whose true motion is in the horizontal plane containing the camera optical centres. The resulting algorithms are simple and fast and can handle object occlusion, disappearance and reappearance. The output is sufficient for determining the number of moving objects in a scene and their general motions.

The methods presented can be extended to general 3D motions. This however involves a more difficult temporal correspondence problem in order to extract the 3D visualised loci from the spatio-temporal block of data. The stereo correspondence problem will however be simplified.



Figure 2. 8 samples from left camera image sequence (video rate)

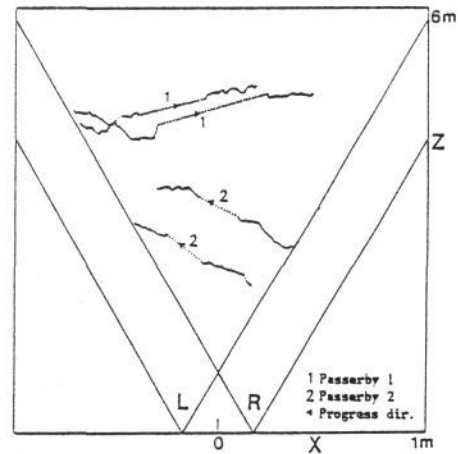


Figure 6. Output of tracking algorithm showing 3D loci of moving objects (plan view) in scene of figure 2. Two loci are shown for each body. One person is moving from left to right while the other is moving from right to left and away from the camera. Solid lines correspond to visible features. Dashed lines correspond to features which are occluded or have the same intensity as the background.

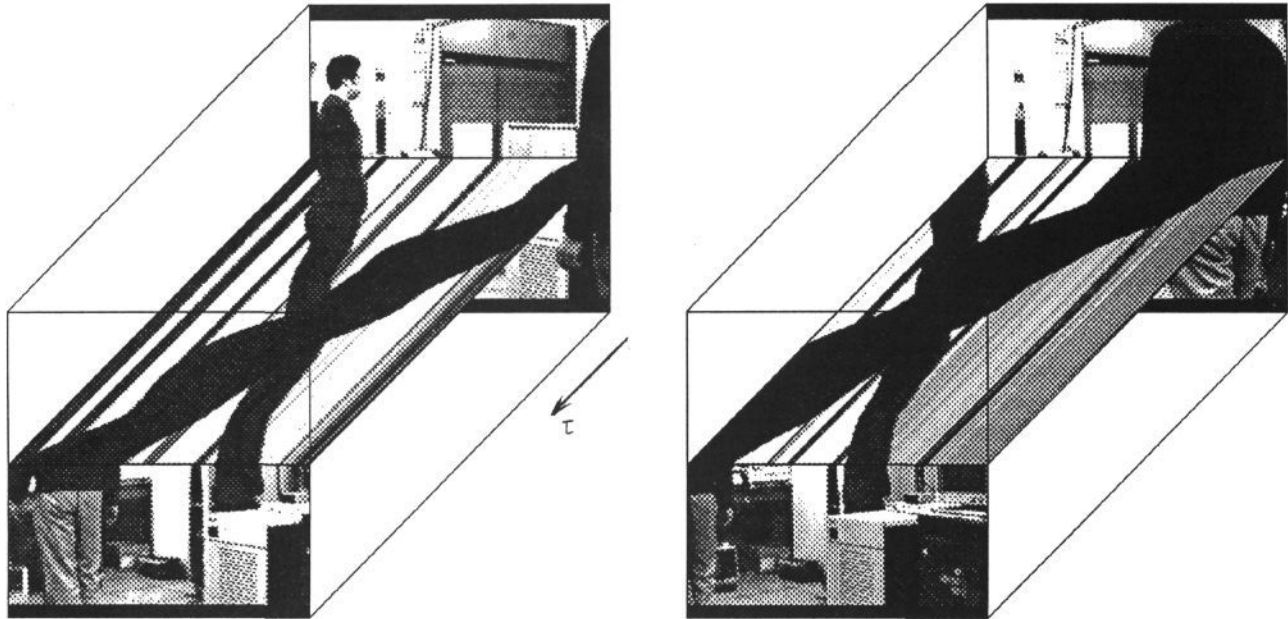


Figure 3. Left and Right 3D spatio-temporal images

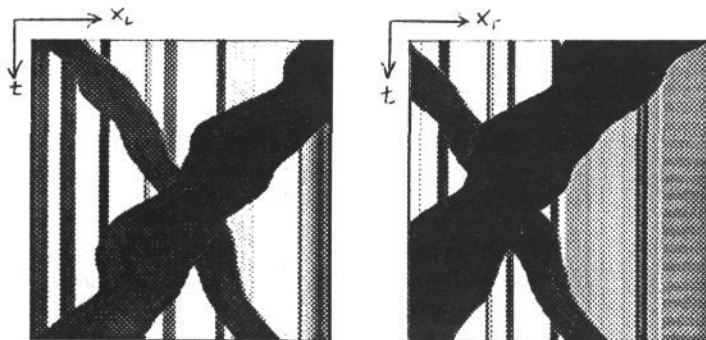


Figure 4. Left and Right 2D cross-section spatio-temporal images showing visualised loci of object features at the same height as the camera centres

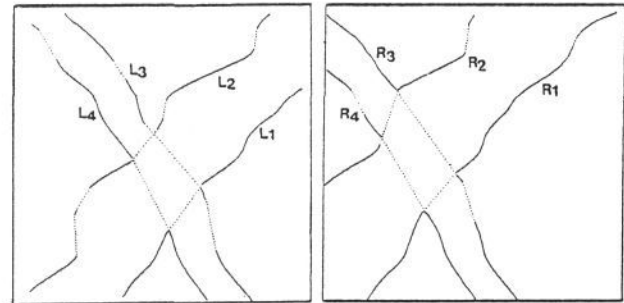


Figure 5: Extracted Loci on left and right spatio-temporal cross-section images. Full lines correspond to fragments of the visualised loci and are linked by dotted lines where the visualised locus disappears due to occlusion, camouflage or noise.



Figure 7. Left 3D spatio-temporal image for scene in which one person overtakes the other

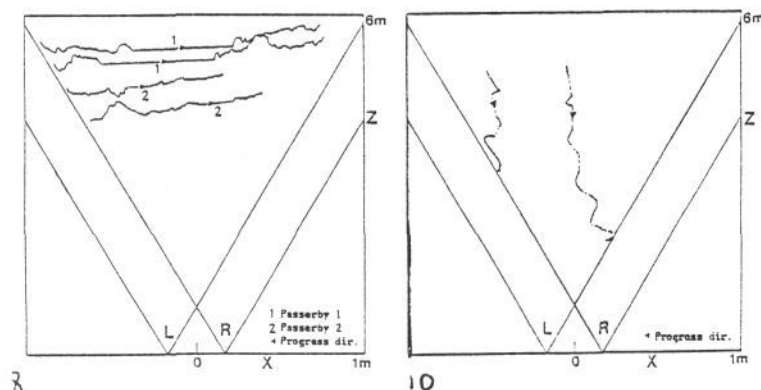


Figure 8. 3D loci of moving bodies (plan view): One person (labelled 1) overtakes the other (2).

Figure 10. 3D loci of a person approaching camera behind a fence. The 2 loci shown correspond to the extrema of the left and right arms. They are modulated by an approximately periodic function due to arm motion.



Figure 9. First and last image of left image sequence of a person approaching the cameras behind a fence.

References

- [1] M. Yamamoto. "Motion analysis using the visualized locus method." *Trans. of Information Processing Society of Japan*, vol.22,no.5:442-449, 1981. (in Japanese).
- [2] R. Cipolla and M. Yamamoto. *Image Sequence Analysis of Human Motion*. Technical Report TR-88-11, Electrotechnical Laboratory, 1988.
- [3] G.F. Poggio and T. Poggio. "The analysis of stereopsis." *Annual Review of Neuroscience*, 7:379-412, 1984.
- [4] S. Ullman. *The interpretation of visual motion*. MIT Press, Cambridge, USA, 1979.
- [5] Barron. J. *A Survey of approaches for determining optic flow, environmental layout and egomotion*. Technical Report RBCV-TR-84-5, University of Toronto, 1984.
- [6] H.C. Longuet-Higgins. "The visual ambiguity of a moving plane." *Proc.R.Soc.Lond.*, B223:165-175, 1984.
- [7] H.P. Moravec. "Visual mapping by a robot rover." In *Proc. of the 6th International Joint Conference on Artificial Intelligence*, pages 598-600, 1979.
- [8] N. Ayache and O.D. Faugeras. "Building, registration and fusing noisy visual maps." In *Proc. 1st Int. Conf. on Computer Vision*, London, 1987.
- [9] M. Jenkin and J.K. Tsotsos. "Applying temporal constraints to the dynamic stereo problem." *Computer, Vision Graphics and Image Processing*, vol.33:16-32, 1986.
- [10] A.M. Waxman and J.H. Duncan. "Binocular image flows : steps toward stereo-motion fusion." *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol.PAMI-8,no.6:715-729, 1986.
- [11] R.C. Bolles, H.H. Baker, and D.H. Marimont. "Epipolar-plane image analysis: an approach to determining structure." *International Journal of Computer Vision*, vol.1:7-55, 1987.
- [12] N.J. Bridwell and T.S. Huang. "A discrete spatial representation for lateral motion stereo." *Computer Vision, Graphics and Image Processing*, 21:33-57, 1983.
- [13] E.H. Adelson and J.R. Bergen. "Spatio-temporal energy models for the perception of motion." *J. Optical Soc. of America*, vol A2,no.2:284-299, 1985.
- [14] J.E.W. Mayhew and J.P. Frisby. "Psychological and computational studies leading towards a theory of human stereopsis". *Artif. Intell.*, 17:349-387, 1981.
- [15] S.B. Pollard, J.E.W. Mayhew, and J.P. Frisby. "PMF: A Stereo Correspondence Algorithm Using A Disparity Gradient". *Perception*, 14:449-470, 1985.
- [16] S. Ganapathy. "Decomposition of transformation matrices for robot vision." In *Proc. of IEEE Conference on Robotics*, pages 130-139, 1984.