

Voting-Based Grouping and Interpretation of Visual Motion

Mircea Nicolescu
Department of Computer Science
University of Nevada, Reno
Reno, NV 89557
mircea@cs.unr.edu

Gérard Medioni
Integrated Media Systems Center
University of Southern California
Los Angeles, CA 90089-0273
medioni@iris.usc.edu

Abstract

A main difficulty for estimating camera and scene geometry from a set of point correspondences is caused by the presence of false matches and independently moving objects. Given two images, after obtaining the matching points they are usually filtered by an outlier rejection step before being used to solve for epipolar geometry and 3-D structure estimation. In the presence of moving objects, image registration becomes a more challenging problem, as the matching and registration phases become interdependent. We propose a novel approach that decouples the above operations, allowing for explicit and separate handling of matching, outlier rejection, grouping, and recovery of camera and scene structure. The method is based on a voting-based computational framework for motion analysis, which determines an accurate representation in terms of dense velocities, segmented motion regions and boundaries by using only the smoothness of image motion, followed by the extraction of scene and camera 3-D geometry.

1. Introduction

The problem of recovering the 3-D scene structure and camera motion from two images has been intensively studied and it is considered well understood. However, most methods perform reasonably well only when: (i) the set of matches contains no outlier noise, and (ii) the scene is rigid – i.e., without objects having independent motions

However, the first assumption almost never holds, since image measurements are bound to be imperfect, and matching techniques will never produce accurate correspondences, mainly due to occlusion or lack of texture. If the second assumption is also violated by the presence of multiple independent motions, even robust methods [1][2] may become unstable, as the scene is no longer a static one. Depending on the size and number of the moving objects, these techniques may return a totally incorrect fundamental matrix. Furthermore, even if the

dominant epipolar geometry is recovered (for example, the one corresponding to the static background), motion correspondences are discarded as outliers.

The core inadequacy of most existing methods is that they attempt to enforce a global constraint – such as the epipolar one – on a data set which may include, in addition to noise, independent subsets that are subject to separate constraints. In this context, it is indeed very difficult to recover structure from motion and segment the scene into independently moving objects, if these two tasks are performed simultaneously.

In order to address these difficulties, we propose a novel approach that decouples the above operations, allowing for explicit and separate handling of matching, outlier rejection, grouping, and recovery of camera and scene structure. In the first step, we determine an accurate representation in terms of dense velocities (equivalent to point correspondences), segmented motion regions and boundaries, by using only the *smoothness of image motion* [3]. In the second step we proceed with the extraction of scene and camera 3-D geometry, separately on each rigid component of the scene.

The main advantage of our approach is that at the 3-D interpretation stage, noisy matches have been already rejected, and correct matches have been grouped according to the distinct motions in the scene. Therefore, standard methods can be reliably applied on each subset of matches in order to determine the 3-D camera and scene structure.

1.1 Related work

Linear methods, such as the Eight Point Algorithm [4][5] can be used for accurate estimation of the fundamental matrix, in the absence of noisy matches or moving objects.

In order to handle outlier noise, more complex, non-linear iterative optimization methods are proposed [1]. These techniques use objective functions, such as distance between points and corresponding epipolar lines, or gradient-weighted epipolar errors, to guide the

optimization process. Despite their increased robustness, iterative optimization methods in general require somewhat careful initialization for early convergence to the correct optimum.

RANSAC [2] consists of random sampling of a minimum subset with seven pairs of matching points for parameter estimation. The candidate subset that maximizes the number of inliers and minimizes the residual is the solution. Although considered one of the most robust methods, it is worth noting that RANSAC still requires a majority of the data to be correct, or else some statistical assumption is needed. If false matches and independent motions exist, many matching points on the moving objects are discarded as outliers.

In [6], Pritchett and Zisserman propose the use of local planar homographies, generated by Gaussian pyramid techniques. However, the homography assumption does not generally apply to the entire image.

1.2. Overview of our method

The first step of the proposed method formulates the motion analysis problem as an inference of motion layers from a noisy and possibly sparse point set in a 4-D space. In order to compute a dense set of matches (equivalent to a velocity field) and to segment the image into motion regions, we use an approach based on a *layered 4-D representation* of data, and a *voting scheme* for communication. First we establish candidate matches through a multi-scale, normalized cross-correlation procedure. Following a perceptual grouping perspective, each potential match is seen as a token characterized by four attributes – the image coordinates (x, y) in the first image, and the velocity with the components (v_x, v_y) .

Tokens are encapsulated as (x, y, v_x, v_y) points in the 4-D space, this being a natural way of expressing the spatial separation of tokens according to *both* velocities and image coordinates. In general, for each pixel (x, y) there can be several candidate velocities, so each 4-D point (x, y, v_x, v_y) represents a potential match.

Within this representation, smoothness of motion is embedded in the concept of surface saliency exhibited by the data. By letting the tokens communicate their mutual affinity through voting, noisy matches are eliminated as they receive little support, and distinct moving regions are extracted as smooth, *salient surface layers* in 4-D.

The second step interprets the image motion by estimating the 3-D scene structure and camera geometry. A rigidity test is performed on the matches within each region, to identify rigid objects, and also between objects, to merge those that move rigidly together but have separate image motions due to depth discontinuities. Finally, the epipolar geometry is estimated separately for each rigid component by using standard methods for parameter

estimation (such as RANSAC), and the scene structure and camera motion are recovered by using the dense velocity field.

2. The Tensor Voting framework

2.1. Voting in 2-D

The use of a voting process for feature inference from sparse and noisy data was formalized into a unified tensor framework by Medioni, Lee and Tang [7]. The input data is encoded as tensors, then support information (including proximity and smoothness of continuity) is propagated by voting. The only free parameter is the scale of analysis, which is indeed an inherent property of visual perception.

In the 2-D case, the salient features to be extracted are points and curves. Each token is encoded as a second order symmetric 2-D tensor, geometrically equivalent to an ellipse. It is described by a 2×2 eigensystem, where eigenvectors e_1 and e_2 give the ellipse orientation and eigenvalues λ_1 and λ_2 are the ellipse size. The tensor is represented as a matrix $S = \lambda_1 \cdot e_1 e_1^T + \lambda_2 \cdot e_2 e_2^T$.

An input token that represents a curve element is encoded as a *stick tensor*, where e_2 represents the curve tangent and e_1 the curve normal, while $\lambda_1=1$ and $\lambda_2=0$. An input point element is encoded as a *ball tensor*, with no preferred orientation, while $\lambda_1=1$ and $\lambda_2=1$.

The communication between tokens is performed through a voting process, where each token casts a vote at each site in its neighborhood. The size and shape of this neighborhood, and the vote strength and orientation are encapsulated in predefined voting fields (kernels), one for each feature type – there is a stick voting field and a ball voting field in the 2-D case. The fields are generated based only on the scale factor σ . Vote orientation corresponds to the smoothest local curve continuation from voter to recipient, while vote strength $VS(\vec{d})$ decays with distance $|\vec{d}|$ between them, and with curvature ρ :

$$VS(\vec{d}) = e^{-\left(\frac{|\vec{d}|^2 + \rho^2}{\sigma^2}\right)} \quad (1)$$

Figure 1(a) shows how votes are generated to build the 2-D stick field. A tensor P where curve information is

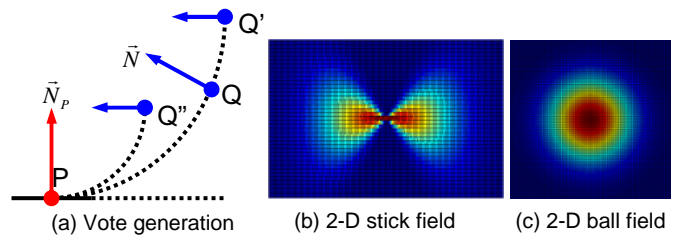


Figure 1. Voting in 2-D

locally known (illustrated by curve normal \vec{N}_p) casts a vote at its neighbor Q. The vote orientation is chosen so that it ensures a smooth curve continuation through a circular arc from voter P to recipient Q. To propagate the curve normal \vec{N} thus obtained, the vote $V_{stick}(\vec{d})$ sent from P to Q is encoded as a tensor according to:

$$V_{stick}(\vec{d}) = VS(\vec{d}) \cdot \vec{N}\vec{N}^T \quad (2)$$

Figure 1(b) shows the 2-D stick field, with its color-coded strength. When the voter is a ball tensor, with no information known locally, the vote is generated by rotating a stick vote in the 2-D plane and integrating all contributions. The 2-D ball field is shown in Figure 1(c).

At each receiving site, the collected votes are combined through simple tensor addition, producing generic 2-D tensors. During voting, tokens that lie on a smooth curve reinforce each other, and the tensors deform according to the prevailing orientation. Each tensor encodes the local orientation of geometric features (given by the tensor orientation), and their saliency (given by the tensor shape and size). For a generic 2-D tensor, its curve saliency is given by $(\lambda_1 - \lambda_2)$, the curve normal orientation by e_1 , while its point saliency is given by λ_2 . Therefore, the voting process infers curves and junctions simultaneously, while also identifying outliers (tokens that receive little support).

2.2. Extension to 4-D

Table 1 shows all the geometric features that appear in a 4-D space and their representation as *elementary* 4-D tensors, where n and t represent normal and tangent

Table 1. Elementary tensors in 4-D

Feature	λ_1	λ_2	λ_3	λ_4	e_1	e_2	e_3	e_4	Tensor
point	1	1	1	1	Any orth. basis				Ball
curve	1	1	1	0	n_1	n_2	n_3	t	C-Plate
surface	1	1	0	0	n_1	n_2	t_1	t_2	S-Plate
volume	1	0	0	0	n	t_1	t_2	t_3	Stick

Table 2. A generic tensor in 4-D

Feature	Saliency	Normals	Tangents
point	λ_4	none	none
curve	$\lambda_3 - \lambda_4$	e_1 e_2 e_3	e_4
surface	$\lambda_2 - \lambda_3$	e_1 e_2	e_3 e_4
volume	$\lambda_1 - \lambda_2$	e_1	e_2 e_3 e_4

vectors, respectively. Note that a surface in the 4-D space can be characterized by two normal vectors, or by two tangent vectors. From a *generic* 4-D tensor that results after voting, the geometric features are extracted as shown in Table 2.

3. Grouping into motion layers

We take as input two image frames that involve general motion – that is, both the camera and the objects in the scene may be moving. For illustration purposes, we give a description of our approach by using a specific example – the two images in Figure 2(a) are taken with a handheld moving camera, while the stack of books has also been moved between taking the two pictures.

Matching. For every pixel in the first image, the goal at this stage is to produce candidate matches in the second image. We use a normalized cross-correlation procedure, where all peaks of correlation are retained as candidates. Each candidate match is represented as a (x, y, v_x, v_y) point in the 4-D space of image coordinates and pixel velocities, with respect to the first image.

In order to increase the likelihood of including the correct match among the candidates, we repeat this process at multiple scales, by using different correlation window sizes. Small windows have the advantage of capturing fine detail, but produce considerable noise in areas lacking texture or having small repetitive patterns. Larger windows generate smoother matches, but their performance degrades in large areas along motion boundaries. We have experimented with a large range of window sizes, and found that best results are obtained by using only two or three different sizes, that should include at least a very small one. In practice we used three correlation windows, with 3x3, 5x5 and 7x7 sizes.

The resulting candidates appear as a cloud of (x, y, v_x, v_y) points in the 4-D space. Figure 2(b) shows the candidate matches. In order to display 4-D data, the last component of each 4-D point has been dropped – the 3 dimensions shown are x and y (in the horizontal plane), and v_x (the height). The motion layers can be already perceived as their tokens appear grouped in two layers surrounded by noisy matches.

Selection. Since no information is initially known, each potential match is encoded into a 4-D *ball tensor*. Then each token casts votes by using the 4-D *ball voting field*. During voting there is strong support between tokens that lie on a smooth surface (layer) – therefore, for each pixel (x, y) we retain the candidate match with the highest surface saliency $(\lambda_2 - \lambda_3)$, and we reject the others as wrong matches. By voting we also estimate the normals to layers at each token as e_1 and e_2 . A 3-D view of the dense layers is shown in Figure 2(c).

Segmentation. The next step is to group tokens into *regions*, by using again the smoothness constraint. We start from an arbitrary point in the image, assign a region label to it, and try to recursively propagate this label to all its image neighbors. In order to decide whether the label must be propagated, we use the smoothness of both velocity and layer orientation as a grouping criterion. Figure 2(d) illustrates the recovered v_x velocities within layers (dark corresponds to low velocity).

Boundary inference. The extracted layers may still be over or under-extended along the true object boundaries, typically due to occlusion. The boundaries of the extracted layers give us a good estimate for the position and overall orientation of the true boundaries. We combine this knowledge with monocular cues (intensity edges) from the original images in order to build a boundary saliency map within the uncertainty zone along the layers margins. At each location in this area, a 2-D stick tensor is generated, having an orientation normal to the image gradient, and a saliency proportional to the gradient magnitude.

The smoothness and continuity of the boundary is then enforced through a 2-D voting process, and the true boundary is extracted as the most salient curve within the saliency map. Finally, pixels from the uncertainty zone are reassigned to regions according to the new boundaries, and their velocities are recomputed. Figure 2(e) shows the refined motion boundaries, that indeed correspond to the actual object.

4. Three-dimensional interpretation

So far we have not made any assumption regarding the 3-D motion, and the only constraint used has been the *smoothness of image motion*. The observed image motion could have been produced by the 3-D motion of objects in the scene, or the camera motion, or both. Furthermore, some of the objects may suffer non-rigid motion.

For classification we used an algorithm introduced by McReynolds and Lowe [8], that verifies the potential

rigidity of a set of minimum six point correspondences from two views under perspective projection. The rigidity test is performed on a subset of matches within each object, to identify potentially rigid objects, and also across objects, to merge those that move rigidly together but have distinct image motions due to depth discontinuities.

The remaining task at this stage is to determine the 3-D object (or camera) motion, and the scene structure. Since wrong matches have been eliminated, and correct matches are already grouped according to the rigidly moving objects in the scene, standard methods for reconstruction can be reliably applied. For increased robustness, we chose to use RANSAC [2] to recover the epipolar geometry for each rigid object, followed by an estimation of camera motion and projective scene structure.

Multiple rigid motions. This case is illustrated by the BOOKS example in Figure 2, where two sets of matches have been detected, corresponding to the two distinct objects – the stack of books and the background. The rigidity test shows that, while each object moves rigidly, they cannot be merged into a single rigid structure. The two sets of recovered epipolar lines are illustrated in Figure 2(f), while the 3-D scene structure and motion are shown in Figure 2(g).

Single rigid motion. This is the stereo case, illustrated by the CANDY BOX example in Figure 3, where the scene is static and the camera is moving. Due to the depth disparity between the box and the background, they exhibit different image motions, and thus they have been segmented as two separate objects. However, the rigidity test shows that the two objects form a rigid configuration, and therefore the epipolar geometry estimation and scene reconstruction are performed on the entire set of matches. Along with the 3-D structure, Figure 3(g) also shows the two recovered camera positions.

Non-rigid motion. The FLAG example, shown in Figure 4, is a synthetic sequence where sparse random dots from the surface of a waving flag are displayed in two frames. The configuration is recognized as non-rigid, and

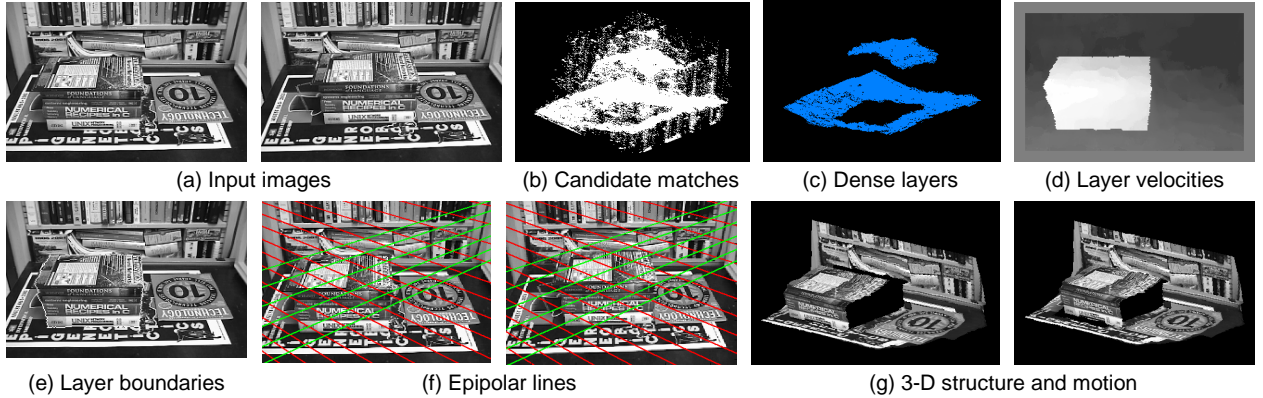


Figure 2. BOOKS sequence

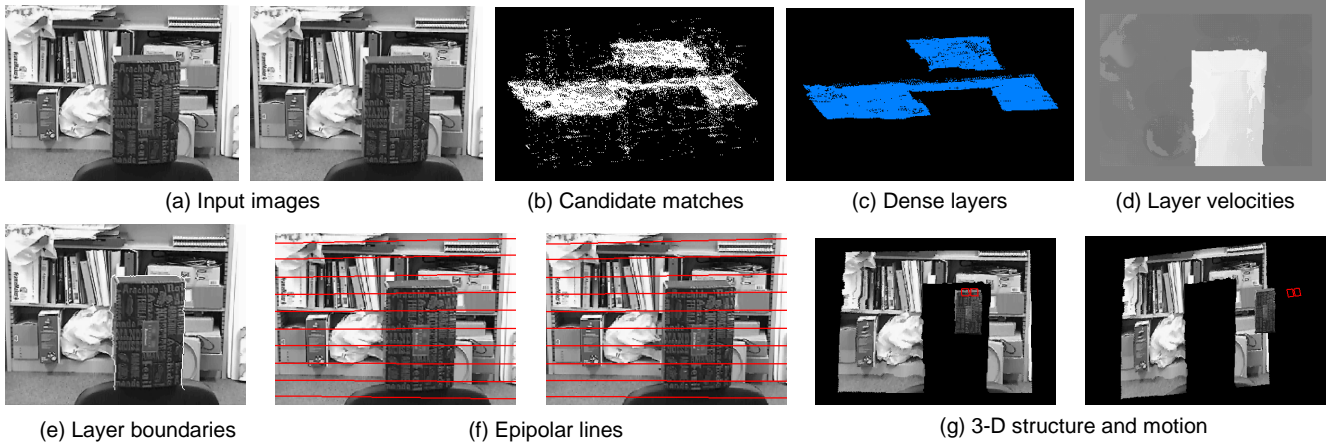


Figure 3. CANDY BOX sequence

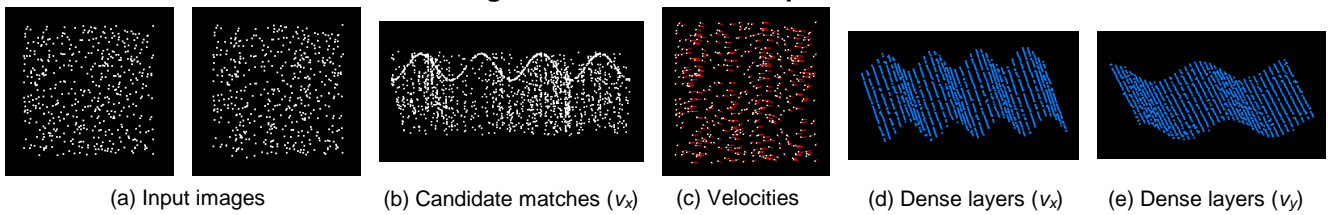


Figure 4. FLAG sequence

therefore no reconstruction is attempted. However, since the *image motion* is smooth, our framework is still able to determine correct correspondences, extract motion layers, segment non-rigid objects, and label them as such.

5. Conclusions

We have presented a novel approach that decouples grouping and interpretation of visual motion, allowing for explicit and separate handling of matching, outlier rejection, grouping, and recovery of camera and scene structure. The proposed framework is able to handle data sets containing large amounts of outlier noise, as well as multiple independently moving objects.

Our methodology for extracting motion layers allows for structure inference without using any prior knowledge of the motion model, based on the smoothness of image motion only, while consistently handling both smooth moving regions and motion discontinuities. The method is also computationally robust, being non-iterative, and does not depend on critical thresholds, the only free parameter being the scale of analysis.

We plan to extend our approach by incorporating information from multiple frames, and to study the possibility of using an adaptive scale of analysis in the voting process.

Acknowledgements

This research has been funded in part by the Integrated Media Systems Center, a National Science Foundation Engineering Research Center, Cooperative Agreement No. EEC-9529152, and by National Science Foundation Grant 9811883.

References

- [1] Z. Zhang, "Determining the Epipolar Geometry and Its Uncertainty: A Review", *IJCV*, 27(2), pp. 161-195, 1998.
- [2] P.H.S. Torr, D.W. Murray, "A Review of Robust Methods to Estimate the Fundamental Matrix", *IJCV*, 1997.
- [3] M. Nicolescu, G. Medioni, "4-D Voting for Matching, Densification and Segmentation into Motion Layers", *ICPR*, 2002.
- [4] H. C. Longuet-Higgins, "A Computer Algorithm for Reconstructing a Scene from Two Projections", *Nature*, 293:133-135, 1981.
- [5] R. I. Hartley, "In Defense of the 8-Point Algorithm", *Trans. PAMI*, 19(6), pp. 580-593, 1997.
- [6] P. Pritchett, A. Zisserman, "Wide Baseline Stereo Matching", *ICCV*, pp. 754-760, 1998.
- [7] G. Medioni, Mi-Suen Lee, Chi-Keung Tang, "A Computational Framework for Segmentation and Grouping", Elsevier Science, 2000.
- [8] D. McReynolds, D. Lowe, "Rigidity Checking of 3D Point Correspondences Under Perspective Projection", *Trans. PAMI*, 18(12), pp. 1174-1185, 1996.