

# Local Gradient, Global Matching, Piecewise-Smooth Optical Flow

Ming Ye

Department of Electrical Engineering  
University of Washington  
Seattle, WA 98195

Robert M. Haralick

Department of Computer Science  
CUNY Graduate Center  
New York, NY 10016

## Abstract

*In this paper we discuss a hybrid technique for piecewise-smooth optical flow estimation. We first pose optical flow estimation as a gradient-based local regression problem and solve it under a high-breakdown robust criterion. Then taking the output from the first step as the initial guess, we recast the problem in a robust matching-based global optimization framework. We have developed novel fast-converging deterministic algorithms for both optimization problems and incorporated a hierarchical scheme to handle large motions. This technique inherits the good sub-pixel accuracy from the local gradient approach and the insensitivity to local perturbation and derivative quality from the global matching approach, and it overcomes the limitations of both. Significant advantages over competing techniques are demonstrated on various standard synthetic and real image sequences.*

## 1. Introduction

Optical flow is a measure of 2D image velocity. It is of traditional importance in computer vision for 3D motion and structure analysis [10] and also receives increasing interest in video coding and computer graphics [14]. The demand on accuracy and efficiency in many real-world applications keeps optical flow estimation an active field of research.

Usually optical flow estimation techniques are either matching-based or gradient-based, and use either local or global optimization schemes. We have developed a hybrid technique which incorporates *all* of the above. It is composed of two steps, the first being a gradient-based local regression method, and the second being a matching-based global optimization method with the output from the first step as its initial guess. Hierarchical processing is adopted to handle large motions [3].

The gradient-based local regression method uses a high-breakdown robust criterion, namely Least Median of Squares (LMS) or Least Trimmed Squares (LTS) [13], to solve the optical flow constraint (Eq. 3) [11]. We approximate the criterion by a novel deterministic iterative algo-

rithm whose complexity adapts to local outlier contaminations. It converges faster and achieves more stable accuracy than the random sampling algorithm previously used in optical flow estimation [1, 12, 15, 19].

The precision of gradient-based techniques saturates near motion boundaries, where the quality of derivatives becomes extremely poor. The smoothing effect of the hierarchical process even worsens the situation. To achieve high boundary fidelity, we use the result from the gradient-based step as the initial guess and minimize a robust matching-based global energy. The energy is designed so that each flow vector minimizes the forward *or* backward warping error and maintain smoothness with the *majority* of its neighbors. The high-quality initialization enables the use of a fast-converging deterministic optimization procedure which results in good sub-pixel accuracy.

The following section provides a review of optical flow estimation research and motivates our study. Section 3 briefly describes the gradient-based local regression method, and Section 4 explains the matching-based global optimization method. Experimental results on various synthetic and real image sequences and comparison with other techniques are given in Section 5. Finally Section 6 concludes our work and points out future work directions.

## 2. Optical Flow Estimation

### 2.1. Brightness Conservation Constraint

The fundamental assumption enabling optical flow estimation is *brightness conservation*, i.e.,

$$I(x, y, t) = I(x + u\Delta t, y + v\Delta t, t + \Delta t), \quad (1)$$

where  $V = (u, v)'$  is the optical flow vector. Depending on what variation of Eq.(1) is used, optical flow estimation methods are classified into two main categories, *matching-based* and *gradient-based*<sup>1</sup>. Matching-based methods make direct use of Eq.(1). They can handle large motions and

<sup>1</sup>Frequency/phase-based methods are close to frequency-domain equivalents of the above two methods [2].

avoid tricky derivative computation, but straightforward implementations usually suffer from poor sub-pixel accuracy [2]. Gradient based methods use the linear approximation of Eq. 1

$$I_x v_x + I_y v_y + I_t = 0, \quad (2)$$

a.k.a the *Optical Flow Constraint Equation (OFCE)* [9], where  $(I_x, I_y, I_t)'$  is the spatiotemporal image intensity gradient. They have been the most popular because of the relatively low complexity and good accuracy, but they meet difficulties with large motions and derivative evaluation.

## 2.2. Flow Smoothness Constraint

For each pixel the brightness constraint (Eq.1 or 2) forms one constraint on  $(u, v)$ . Additional constraints come from the *flow field smoothness assumption*, which means neighboring pixels experience consistent motion. Based on how smoothness is imposed, approaches are further divided into two types, *local parametric* and *global optimization*. Local parametric methods assume that within a certain neighborhood, which could be as large as the entire image, the flow field is described by a parametric model [3], with the simplest and most popular model being piecewise constant. Their accuracy and efficiency are among the best according to various comparative studies [2, 7], but they degrade or fail when local information becomes insufficient or unreliable. Global optimization methods cast optical flow estimation in a regularization framework—every vector satisfies its brightness constraint while maintaining coherence with its neighbors [9]. Such approaches are less sensitive to local constraint quality, but existing techniques of this type, even those using robust estimators [4], tend to oversmooth the flow field.

Most traditional techniques assume only one of these approaches. A detailed survey and comparison of these techniques is given by Barron et al. [2].

## 2.3. Robust Methods

In fact many existing techniques are able to produce reasonably good results when their assumptions approximately hold, and the real challenge in motion estimation is to achieve high robustness against assumption violations especially motion discontinuities.

One attempt at enhancing robustness is to use *model-based* techniques [14, 18], which make explicit assumptions about objects and motions in the scene, as opposed to *general* flow estimation techniques which only assume piecewise-smooth flow. Model-based formulations usually involve large-scale nonconvex optimization problems. In practice the solution often boils down to a procedure alternating between general motion estimation and interpretation, and the achievable accuracy largely depends on the

initial motion estimates from a general method [16]. Moreover, precisely modeling motion fields is difficult. Some motions (e.g. facial expression) do not have explicit models, and some have models constantly varying with time [17, 16]. Therefore, although model-based techniques have achieved some success, they are not a replacement for general estimation techniques.

Another approach, which is applicable to general optical flow estimation, is to use robust statistics [13]. Most traditional techniques solve the estimation problems in the Least-Squares (LS) sense. LS criteria have little tolerance to assumption violations, and they form a major source of gross errors. This problem has been widely recognized and led to extraordinary efforts of replacing LS estimators in traditional techniques by more robust ones. For instance, [4] uses an M-estimator in the Horn-Schunck method [9]; [1] uses an LMS estimator in the Lucas-Kanade (LK) method [11]. Our technique applies robust criteria to piecewise-smooth motion estimation.

## 3. Gradient-Based Local Regression

The first step in the proposed technique is a gradient-based local regression method. Following LK [11], we group  $n$  OFCEs around each pixel, forming a linear equation in  $V$

$$AV = b \quad (3)$$

$$A = \begin{pmatrix} I_{x1} & I_{y1} \\ \vdots & \vdots \\ I_{xn} & I_{yn} \end{pmatrix} \quad b = - \begin{pmatrix} I_{t1} \\ \vdots \\ I_{tn} \end{pmatrix}.$$

The LK method solves it in an LS sense. More robust criteria have also been experimented with, including M [4], LMS [1] and LTS [19]. Among them we find high-breakdown robust criteria [13], e.g. LMS and LTS, are most appropriate [19].

So far all optical flow techniques using high-breakdown criteria [1, 12, 15, 19] adopt the bootstrap-like algorithm given in [13]—the estimate with the best criterion value is picked from a random pool of trial estimates. They uniformly apply the algorithms to all pixels disregarding the actual amount of outliers.

By taking advantage of the piecewise smoothness property of optical flow fields and the selection capability of robust estimators, we have proposed a deterministic adaptive algorithm. Basically starting from LS estimates, we iteratively choose neighbors' values as trial solutions, select inliers using LMS, calculate updated solutions using LS on inliers and the associated criterion values, and finally keep the solution with the best criteria value [13]. This method in effect provides an estimator whose complexity depends on the actual outlier contamination. It is faster and has more stable accuracy than algorithms based on random sampling.

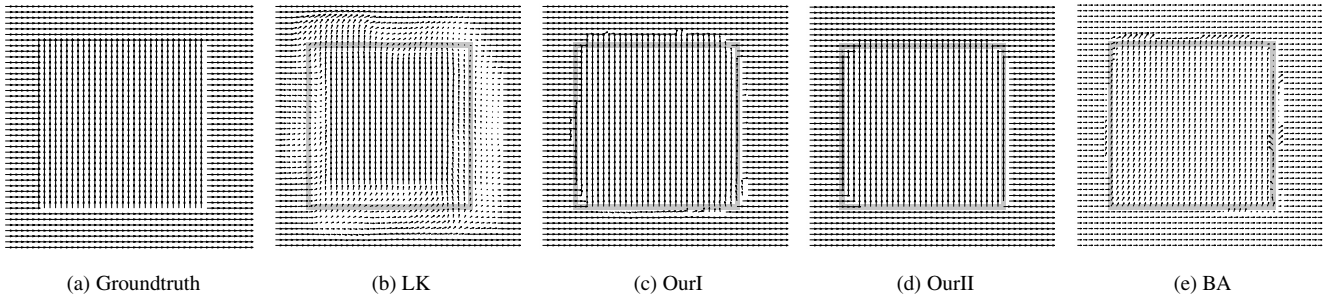


Figure 1: TS sequence results

The performance of the LS and LMS methods are compared on a synthetic sequence Translating Squares (TS). It contains two squares both translating at 1 pixel/frame. Vector plots around the motion boundary (shadowed) are given in Fig. 1. The result from our first method has a much more clear-cut motion boundary than that from LK. However the position of the boundary still has noticeable offsets from the truth. These errors are due to derivative estimation failure, which is inevitable to all gradient-based methods. In order to tackle this problem, we further enhance the result by a matching based approach.

## 4. Matching-Based Global Optimization

Global optimization techniques minimize a global energy  $E = \sum_{\text{all pixels } i} E_B(V_i) + \lambda E_S(V_i)$  where  $V_i$  is the flow vector at pixel  $i$ .  $E_S$  represents flow smoothness, normally defined on the difference between each vector and its neighbors;  $E_B$  represents brightness conservation, defined in matching-based approaches on the intensity error (Eq. 1) between two frames; and  $\lambda$  controls their relative importance. Various formulations of  $E_S$  have been presented to account for discontinuities; in particular, the notion of weak continuity [6] has been popular [4]. Pointing out that such formulations are special cases of robust estimation problems, Black [5] poses both  $E_B$  and  $E_S$  using M-estimators.

Beside those mentioned in Section 2, many serious problems exist with current global matching formulations. (i) They are *very* hard to solve, especially in a robust estimation framework. Simulated annealing type of methods [8] have to be used [5], which are extremely slow and unstable. (ii) The advantage of matching-based methods at motion discontinuities is usually lost with crude hierarchical schemes. (iii) Their difficulty with sub-pixel accuracy still has not been effectively dealt with. (iv) Pixels at occlusions and image borders might not be visible in both frames and thus the matching criteria result in gross errors [16]. (v) Various parameters including  $\lambda$ , and scales in M-estimator based methods, have to be manually tuned for meaningful outputs. As a consequence of all above problems, so far no such methods to our knowledge consistently outperform

gradient-based counterparts.

In the follows we present a novel approach which overcomes most of above limitations and achieves highly competitive performance.

### 4.1. Global Energy Design

**Matching Error.** We observe that without aliasing, *all pixels in a frame are visible in the previous or the next frame*. This means in order to find correspondence for all pixels, at least *three* successive frames instead of two should be examined. Therefore we propose the normalized matching error

$$E_B(V_i) = \begin{cases} \frac{2e_n(V_i)}{I_i + I_n(V_i)}, & e_p(V_i) > e_n(V_i) \\ \frac{2e_p(V_i)}{I_i + I_p(V_i)}, & \text{otherwise} \end{cases}$$

where  $I_p(V_i), I_n(V_i)$  are warped intensities in the previous and the next frames respectively, and  $e_p(V_i) = |I_i - I_p(V_i)|, e_n(V_i) = |I_i - I_n(V_i)|$ . This error measure not only allows matching against occlusions, but also provides a means to detect such situations.

**Smoothness Error.** This term requires a vector  $V_i$  to be smooth with the majority of its 8-connected neighbors  $V_j, j \in N_i$ . We calculate error vector energies  $e_{V_i}(j) = |V_i - V_j|^2$  for all  $V_j$ 's, select those consistent with  $V_i$ ,  $e_{V_i}(\text{inliers})$ , using the LMS-LS procedure (Section 3), and compute the normalized smoothness error as

$$E_S(V_i) = \overline{e_{V_i}(\text{inliers})} / (|V_i|^2 + 1).$$

This scheme rejects outliers adaptively according to local flow variation. In contrast, previous uses of M-estimators select outliers by a fixed global scale and have difficulties with scenes having a wide range of motions (Fig. 6(b)).

Finally our global energy is simply

$$E = \sum_{\text{all pixels } i} E_B(V_i) + E_S(V_i) \quad (4)$$

Notice that no tuning parameters such as  $\lambda$  exist in this measure. By use of normalization and adaptive scaling, different sources of errors are automatically balanced, and hence

this method has good local accuracy and consistent performance on different data.

## 4.2. Minimization

Each flow estimate  $V_i$  affects the global energy  $E$  through its own pixel energy  $E_B(V_i) + E_S(V_i)$  and smoothness energies of a few neighbors. The set of affected pixels is called a *clique* [8]. The minimum of  $E$  is reached when all cliques have the lowest energy. Starting from a field of initial estimates, we may in turn perturb each estimate, observe its clique energy change, and accept the candidate if it leads to a total energy decrease. Two critical problems in this procedure are how to generate candidates and when to accept a candidate. Greedy algorithms accept a candidate iff. the clique energy reduces. Stochastic algorithms occasionally accept a “bad” candidate to avoid local optima. Candidates may be generated according to a fixed global schedule, or randomly by sampling local distribution [6].

With a good initial guess seldom available, random schemes usually have to be used [5]. However, equipped with the high-quality initial estimates from the robust local gradient method (Section 3), we are able to minimize the global energy using a *simple fastest-descent* algorithm.

We first calculate the clique energy for all pixels. Then we repeatedly visit each pixel examining whether a trial estimate from a candidate set results in a lower pixel energy. The candidate set consists of the 8-connected neighbors and their average which were updated in the last visit. Once seeing a pixel energy decrease, we accept the candidate and update the clique energy. This process continues until no pixel is updated. It converged quickly in all our past experiments.

As it is clear in Fig. 1(c), estimates from the first step have excellent accuracy away from discontinuities. Drawing candidates from them, together with the averaging update, passes on the good sub-pixel accuracy to the estimates in the global matching step.

## 4.3. Overall Algorithm

We employ a hierarchical process [3] to cope with large motions and expedite convergence. We create a  $P$ -level image pyramid  $I^p, p = 0, \dots, P - 1$  and start estimation from the top (coarsest) level  $P - 1$  with a zero initial flow field. At each level  $p$ , we warp images  $I^p$  using the initial flow  $V_0^p$  obtaining images  $I_w^p$ . On  $I_w^p$  we calculate the residual flow  $\Delta V^p$  using the local gradient method and add it to  $V_0^p$  yielding  $V_1^p$ . Then we refine  $V_1^p$  by applying the global matching method to  $I^p$ , resulting in the final flow estimate on Level  $p, V_2^p$ , which is projected down to Level  $p - 1$  as its initial flow field  $V_0^{p-1}$ . Finally the flow estimate on the original resolution is  $V_2^0$ .

Hierarchical schemes have lots of limitations which are often overlooked. The projection and warping operations

oversmooth the flow field; they often even become the accuracy bottleneck especially at discontinuities. Errors in coarser levels are magnified and propagated to finer levels and are generally irreversible [5, 16]. These problems are much alleviated by our global refinement step—it works on the original pyramid images and corrects gross errors caused by derivative computation, projection and warping.

## 5. Experiments and Analysis

We calculate derivatives from a first-order spatiotemporal facet model on a support of size  $3 \times 3 \times 3$  [19]. Optical flow is estimated on the middle frame of every three frames. The number of pyramid levels is empirically determined. Bilinear interpolation is used for image warping. The constant flow window size is fixed at  $9 \times 9$ . No image pre-smoothing is done. We handle estimates at image borders such that they also have good accuracy and the resulting flow field is of the same size as the original image. For the experiments reported here, we use only trial values which are at least  $T_1$  and  $T_2$  pixels different from the current estimate in the propagation procedures described in Section 3 and 4.2 respectively.  $T_1, T_2$  are used to speed up the computation but they do not affect the results. We set  $T_1 = 0.01$  for synthetic data and  $T_1 = T_2 = 0.05$  in all other cases.

Results are given for four techniques: (i) LS local gradient (LS), (ii) LMS local gradient (OurI), (iii) global matching (OurII) and (iv) Black and Anandan’s robust regularization method (BA) [4]. BA’s code was provided by Michael Black with all parameters set as in [4]. It calculates flow on the second of two frames. (i) is a modified version of Lucas and Kanade’s [11]; and (ii) is an improved version of Bab-Hadiashar and Suter’s method [1]. All experiments are carried out on a PIII 500MHz PC running Solaris.

### 5.1. Synthetic Sequences

Four image sequences with flow groundtruth are used for quantitative comparison. Two error measures are reported. One is the angular error  $e_\angle$  used in [2]. The other is the error vector magnitude measure  $e_{| \cdot |} = |V - V_0|/|V_0|$ , where  $V_0$  is the correct flow vector. We also report the consumed CPU time in seconds to give a rough idea on speeds. All measures are summarized in Table 1. Generally, consistent observations of smaller errors indicate better quality. But due to the simplicity of the data and the crudeness of the error measures, the numbers should not be taken literally to claim quantitative merits.

The TS sequence was introduced in Section 3. Vector plots for OurII and BA are given in Fig. 1. OurII has almost perfect results except the rounded corners, where the background motion becomes dominant. BA produces poor accuracy on this one—it achieves global smoothness with the sacrifice of local fidelity.

Three synthetic sequences are obtained from John Barrow [2], Translating Tree (TT), Diverging Tree (DT) and Yosemite (YOS). TT and DT simulate translational camera motion with respect to a textured planar surface. TT's motion is horizontal and DT's is divergent. YOS's motion is mostly divergent. The cloud part is excluded from evaluation. See [2] for detailed descriptions of these data. We use 2 levels of pyramid for TT and DT, and 3 levels for YOS.



Figure 2: TT, DT middle frame



Figure 3: YOS middle frame

Data	Technique	$e_r$ ( $^\circ$ )	$e_{ v }$ (%)	time (sec)
TS	LS	6.14	15.12	0
	BA	8.12	21.08	1
	OurI	1.09	2.65	0
	OurII	0.32	0.79	1
TT	LS	2.36	5.48	1
	BA	2.61	6.62	11
	OurI	1.39	3.22	7
	OurII	0.49	1.53	8
DT	LS	6.12	18.33	1
	BA	6.57	20.70	11
	OurI	5.00	16.14	10
	OurII	4.43	17.25	14
YOS	LS	3.69	12.68	4
	BA	2.77	10.16	62
	OurI	3.42	11.10	40
	OurII	2.87	9.94	55

Table 1: Quantitative measures

## 5.2. Real Sequences

Results on two well-know real image sequences, Hamburg Taxi (TAXI), and Flower Gargen (FG), are given below.

TAXI [2] has four moving objects: three cars at image speeds 3.0, 1.0, 3.0 pixels/frame (from left to right) respectively, and a pedestrian in the upper-left walking rightwards at about 0.3 pixels/frame. Two levels of pyramid are used. Vector plots with enough details do not fit the page, so we display the horizontal flow component as intensity images in Fig. 6. Brighter pixels represent larger speeds to the right. In LS the flow fields of the vehicles have severely invaded into the background. BA's result is completely smoothed

out. Better performance might be obtained by tuning parameters. But smoothing seems to be inevitable in scenes of such diverse motions as long as rigid global parameters are used. OurI preserves motion boundaries better but still shows smoothing incurred by the hierarchical process. Our-II gives the best result.

Motion in FG is caused by camera translation. Three main layers are observed due to the scene depth corresponding to the tree in front (maximum speed as large as about 6 pixels/frame), the garden and the rest background. This version of FG was obtained from Michael Black. It has a dark strip at the right border. Three levels of pyramid are used. Horizontal flows are given in Fig. 6. In all results the motion of the tree twigs smears into the background. The reason is that, the background sky has little texture and thus any flow estimate yield good matching (aperture problem); and the errors are further enlarged by the hierarchical process. BA and OurII work much better than LS and OurI. But BA still has considerable oversmoothing between every pair of layers, and one twig in the upper-right is missing. OurII shows a very clear layered look. The results are superior to those produced by other general method [16], and are highly competitive to those from model or layer based techniques [14, 18].



Figure 4: TAXI middle frame



Figure 5: FG middle frame

## 6. Conclusions And Discussion

This paper has presented a novel hybrid optical flow estimation technique. A gradient-based local regression method is employed to produce an initial solution, and then the result is refined in a matching-based global optimization framework. This technique has inherited the merits of both approaches and overcome their limitations. In particular, high sub-pixel accuracy of the local gradient method is passed on to the global matching method, and the latter corrects most errors caused by failure in gradient evaluation and smoothing in the hierarchical process. Our results on various standard testing sequences show significant advantages over previous general techniques and high competitiveness with model/layer based ones.

As robust estimation has formed a trend in computer vision, its efficient and effective implementation remains a great challenge. Therefore the contribution of our work *not only* lies in the new robust formulations, *but* more impor-

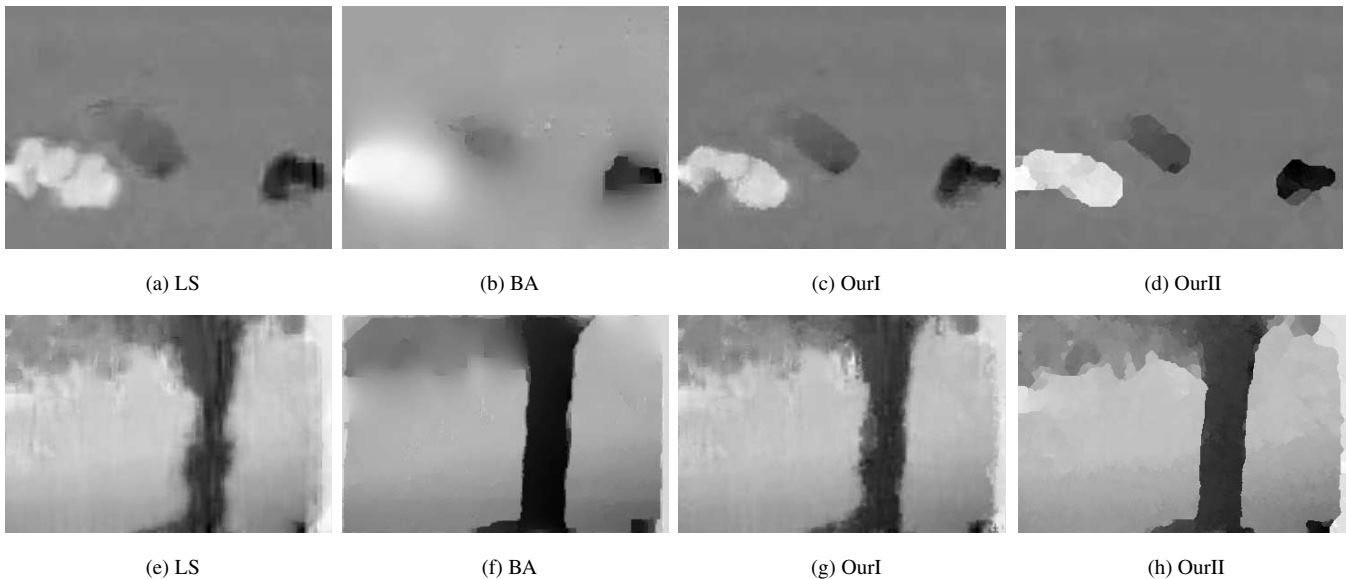


Figure 6: TAXI (1st row) and FG (2nd row): intensity images of x-component

tantly in the optimization schemes which achieve the designed robustness. We are probably the first to approximate a high-breakdown robust criterion by a deterministic algorithm with adaptive complexity. We define a matching criterion on *three* instead of two frames to avoid the visibility problem caused by occlusions [16]. Comparing to previous related work, our robust regularization method has the advantages of automatically balancing different energy terms, and solving the large-scale nonconvex problem with a simple greedy method. The above are all accomplished by carefully integrating standard robust methods with the characteristics of the optical flow application.

We are currently extending the research in a number of directions. First the reported algorithm exposes the viability of reliable general motion estimation from a hybrid approach; but many details still need careful consideration. The matching error is able to detect occlusions, and thus can guide image warping. We have already carried out some experiments in this direction and obtained preliminary success. As far as application is concerned, we are now exploring automatic motion interpretation, particularly model selection [17], motion segmentation [14], layered representation [18], from the flow estimates. We expect these tasks be much eased with the high-quality input. In addition, the proposed robust approaches might find applications in other computer vision problems such as image restoration.

## References

- [1] A. Bab-Hadiashar and D. Suter, "Robust Optical Flow Estimation", *IJCV*, Vol. 29, No. 1, pp. 59-77, 1998.
- [2] J.L. Barron, S. S. Beauchemin and D. J. Fleet, "Performance of Optical Flow Techniques", *IJCV*, Vol. 12, No. 1, pp. 43-77, 1994.
- [3] J.R. Bergen, P. Anandan, K.J. Hanna and R. Hingorani, "Hierarchical Model-Based Motion Estimation", *Proc. ECCV*, pp. 237-252, 1992.
- [4] M. J. Black and P. Anandan, "The Robust Estimation of Multiple Motions: Parametric and Piecewise-Smooth Flow Fields", *CVIU*, Vol. 63, No. 1, pp. 75-104, 1996.
- [5] M. J. Black, "Robust Incremental Optical Flow", *Ph.D. Thesis, Yale Univ.*, 1992; *Research Report YALEU/DCS/RR-923*.
- [6] A. Blake and A. Zisserman, *Visual Reconstruction*, MIT Press, 1987.
- [7] B. Galvin, B. McCane, K. Novins, D. Mason and S. Mills, "Recovering Motion Fields: An Evaluation of Eight Optical Flow Algorithms", *Proc. BMVC*, 1998.
- [8] S. Geman and D. Geman, "Stochastic Relaxation, Gibbs Distributions and Bayesian Restoration of Images", *PAMI*, No. 6, pp. 721-741, 1984.
- [9] B.K.P. Horn and B.G. Schunck, "Determining optical flow", *AI*, Vol. 17, pp. 185-203, 1981.
- [10] T. Jebara, A. Azarbayejani and A. Pentland, "3D Structure From 2D Motion", *IEEE SP Magazine*, pp. 66-84, May 1999.
- [11] B.D. Lucas and T. Kanade, "An Iterative Image-Registration Technique with an Application to Stereo Vision", *DARPA Proc. IUW*, pp. 121-130, 1981.
- [12] E. P. Ong and M. Spann, "Robust Optical Flow Computation Based on Least-median-of-Squares", *IJCV*, Vol. 31, No. 1, pp. 51-82, 1999.
- [13] P. J. Rousseeuw and A. M. Leroy, *Robust Regression and Outlier Detection*, John Wiley and Sons, 1987.
- [14] H. S. Sawhney and S. Ayer, "Compact representations of videos through dominant and multiple motion estimation," *PAMI*, Vol. 18, No. 8, pp. 814-830, 1996.
- [15] D-G. Sim and R-H Park, "Robust Reweighted MAP Motion Estimation", *PAMI*, Vol. 20, No. 4, pp. 353-365, 1998.
- [16] R. Szeliski, "A Multi-View Approach to Motion and Stereo", *Proc. CVPR*, Vol. I, pp. 157-163, 1999.
- [17] P.H.S. Torr, "Geometric Motion Segmentation and Model Selection", *Royal Soc., J. Lasenby et al. (eds.)*, pp. 1321-1340, 1998.
- [18] P.H.S. Torr, R. Szeliski and P. Anandan, "An Integrated Bayesian Approach to Layer Extraction from Image Sequences", *Proc. ICCV*, pp. 983-990, 1998.
- [19] M. Ye and R.M. Haralick, "Optical Flow From A Least-Trimmed Squares Based Adaptive Approach", *Proc. ICPR*, Vol. 3, pp. 1052-1055, 2000.