

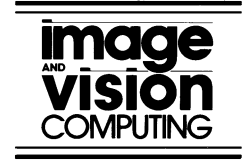


ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)



Image and Vision Computing 22 (2004) 761–767



[www.elsevier.com/locate/imavis](http://www.elsevier.com/locate/imavis)

# Robust wide-baseline stereo from maximally stable extremal regions

J. Matas<sup>a,b,\*</sup>, O. Chum<sup>a</sup>, M. Urban<sup>a</sup>, T. Pajdla<sup>a</sup>

<sup>a</sup>Department of Cybernetics, Center for Machine Perception, CTU Prague, Karlovo nám 13 CZ 121 35, Czech Republic

<sup>b</sup>CVSSP, University of Surrey, Guildford GU2 7XH, UK

Received 12 March 2003; received in revised form 11 February 2004; accepted 12 February 2004

## Abstract

The wide-baseline stereo problem, i.e. the problem of establishing correspondences between a pair of images taken from different viewpoints is studied.

A new set of image elements that are put into correspondence, the so called *extremal regions*, is introduced. Extremal regions possess highly desirable properties: the set is closed under (1) continuous (and thus projective) transformation of image coordinates and (2) monotonic transformation of image intensities. An efficient (near linear complexity) and practically fast detection algorithm (near frame rate) is presented for an affinely invariant stable subset of extremal regions, the maximally stable extremal regions (MSER).

A new robust similarity measure for establishing tentative correspondences is proposed. The robustness ensures that invariants from multiple measurement regions (regions obtained by invariant constructions from extremal regions), some that are significantly larger (and hence discriminative) than the MSERs, may be used to establish tentative correspondences.

The high utility of MSERs, multiple measurement regions and the robust metric is demonstrated in wide-baseline experiments on image pairs from both indoor and outdoor scenes. Significant change of scale ( $3.5 \times$ ), illumination conditions, out-of-plane rotation, occlusion, locally anisotropic scale change and 3D translation of the viewpoint are all present in the test problems. Good estimates of epipolar geometry (average distance from corresponding points to the epipolar line below 0.09 of the inter-pixel distance) are obtained.

© 2004 Elsevier B.V. All rights reserved.

**Keywords:** Wide-baseline stereo; Distinguished regions; Maximally stable extremal regions; MSER; Robust metric

## 1. Introduction

Finding reliable correspondences in two images of a scene taken from arbitrary viewpoints viewed with possibly different cameras and in different illumination conditions is a difficult and critical step towards fully automatic reconstruction of 3D scenes [5]. A crucial issue is *the choice of elements whose correspondence is sought*. In the wide-baseline set-up, local image deformations cannot be realistically approximated by translation or translation with rotation and a full affine model is required. Correspondence cannot be therefore established by comparing regions of a fixed (Euclidean) shape like rectangles or circles since their shape is not preserved under affine transformation.

In most images there are regions that can be detected with high repeatability since they possess some distinguishing, invariant and stable property. We argue that such regions of, in general, data-dependent shape, called distinguished regions (DRs) in the paper, may serve as the elements to be put into correspondence either in stereo matching or object recognition.

The first contribution is the introduction of a new set of DRs, the so called *extremal regions*. Extremal regions have two desirable properties. The set is closed under continuous one-to-one (and thus perspective) transformation of image coordinates and, secondly, it is closed under monotonic transformation of image intensities. An efficient (near linear complexity) and practically fast detection algorithm is presented for an affinely invariant stable subset of extremal regions, the maximally stable extremal regions (MSER). Robustness of a particular type of DR depends on the image data and must be tested experimentally. Successful wide-baseline experiments on indoor and outdoor datasets presented in Section 4 demonstrate the potential of MSERs.

\* Corresponding author. Address: Department of Cybernetics, Center for Machine Perception, CTU Prague, Karlovo nám 13 CZ 121 35, Czech Republic. Tel.: +420-2-24357212; fax: +420-2-24357385.

E-mail addresses: [matas@cmp.felk.cvut.cz](mailto:matas@cmp.felk.cvut.cz) (J. Matas); [chum@cmp.felk.cvut.cz](mailto:chum@cmp.felk.cvut.cz) (O. Chum).

Reliable extraction of a manageable number of potentially corresponding image elements is a necessary but certainly not a sufficient prerequisite for successful wide-baseline matching. With two sets of DRs, the matching problem can be posed as a search in the correspondence space [4]. Forming a complete bipartite graph on the two sets of DRs and searching for a globally consistent subset of correspondences is clearly out of question for computational reasons. Recently, a whole class of stereo matching and object recognition algorithms with common structure has emerged [1,3,7,9,10,13,15,18,20,21]. These methods exploit *local invariant descriptors* to limit the number of tentative correspondences. Important design decisions at this stage include: (1) the choice of measurement regions, i.e. the parts of the image on which invariants are computed, (2) the method of selecting tentative correspondences given the invariant description and (3) the choice of invariants.

Typically, DRs or their scaled version serve as measurement regions and tentative correspondences are established by comparing invariants using Mahalanobis distance [14,16,21]. As a second novelty of the presented approach, a robust similarity measure for establishing tentative correspondences is proposed to replace the Mahalanobis distance. The robustness of the proposed similarity measure allows us to use invariants from a collection of measurement regions, even some that are much larger than the associated DR. Measurements from large regions are either very discriminative (it is very unlikely that two large parts of the image are identical) or completely wrong (e.g. if orientation or depth discontinuity becomes part of the region). The former helps establishing reliable tentative (local) correspondences, the influence of the latter is limited due to the robustness of the approach.

Finding epipolar geometry (EG) consistent with the largest number of tentative (local) correspondences is the final step of all wide-baseline algorithms. RANSAC has been by far the most widely adopted method since [19]. The presented algorithm takes novel steps to increase the number of matched regions and the precision of the EG. The rough EG estimated from tentative correspondences is used to guide the search for further region matches. It restricts location to epipolar lines and provides an estimate of affine mapping between corresponding regions. This mapping allows the use of correlation to filter out mismatches. The process significantly increases precision of the EG estimate; the final average inlier distance-from-epipolar-line is below 0.1 pixel. For details see Section 3.

*Related work.* Since the influential paper by Schmid and Mohr [16] many image matching and wide-baseline stereo algorithms have been proposed, most commonly using Harris interest points as DRs. Tell and Carlsson [18] proposed a method where line segments connecting Harris interest points form measurement regions. The measurements are characterised by scale invariant Fourier coefficients. The Harris interest detector is stable over a range of scales, but defines no scale or affine invariant measurement

region. Baumberg [1] applied an iterative scheme originally proposed by Lindeberg and Gårding [6] to associate affine-invariant measurement regions with Harris interest points. In [10], Mikolajczyk and Schmid show that a scale-invariant MR can be found around Harris interest points. In [11], the approach was combined with Baumberg's iteration to obtain an affine-invariant detector. In [13], Pritchett and Zisserman form groups of line segments and estimate local homographies using parallelograms as measurement regions. Tuytelaars and Van Gool introduced two new classes of affine-invariant DRs, one based on local intensity extrema [21] the other using point and curve features [20]. In the latter approach, DRs are characterised by measurements from inside an ellipse, constructed in an affine invariant manner. Lowe [7] describes the 'Scale Invariant Feature Transform' approach which produces a scale and orientation-invariant characterisation of interest points.

The rest of the paper is structured as follows. MSER are defined and their detection algorithm is described in Section 2. In Section 3, details of a novel robust matching algorithm are given. Experimental results on outdoor and indoor images taken with an uncalibrated camera are presented in Section 4. Presented experiments are summarized and the contributions of the paper are reviewed in Section 5.

## 2. Maximally stable extremal regions

In this section, we introduce a new type of image elements useful in wide-baseline matching—the *Maximally Stable Extremal Regions*. The regions are defined solely by an extremal property of the intensity function in the region and on its outer boundary.

The concept can be explained informally as follows. Imagine all possible thresholdings of a gray-level image  $I$ . We will refer to the pixels below a threshold as 'black' and to those above or equal as 'white'. If we were shown a movie of thresholded images  $I_t$ , with frame  $t$  corresponding to threshold  $t$ , we would see first a white image. Subsequently black spots corresponding to local intensity minima will appear and grow. At some point regions corresponding to two local minima will merge. Finally, the last image will be black. The set of all connected components of all frames of the movie is the set of all maximal regions; minimal regions could be obtained by inverting the intensity of  $I$  and running the same process. The formal definition of the MSER concept and the necessary auxiliary definitions are given in Table 1.

In many images, local binarization is stable over a large range of thresholds in certain regions. Such regions are of interest since they possess the following properties:

- Invariance to affine transformation of image intensities.
- *Covariance to adjacency preserving* (continuous) transformation  $T : \mathcal{D} \rightarrow \mathcal{D}$  on the image domain.

Table 1  
Definitions used in Section 2

Image  $I$  is a mapping  $I : \mathcal{D} \subset \mathbb{Z}^2 \rightarrow \mathcal{S}$ . Extremal regions are well defined on images if:

1.  $\mathcal{S}$  is totally ordered, i.e. reflexive, antisymmetric and transitive binary relation  $\leq$  exists. In this paper only  $S = \{0, 1, \dots, 255\}$  is considered, but extremal regions can be defined on, e.g. real-valued images ( $\mathcal{S} = R$ )

2. An adjacency (neighbourhood) relation  $A \subset \mathcal{D} \times \mathcal{D}$  is defined. In this paper 4-neighbourhoods are used, i.e.  $p, q \in \mathcal{D}$  are adjacent ( $pAq$ ) iff  $\sum_{i=1}^4 |p_i - q_i| \leq 1$

Region  $\mathcal{Q}$  is a contiguous subset of  $\mathcal{D}$ , i.e. for each  $p, q \in \mathcal{Q}$  there is a sequence  $p, a_1, a_2, \dots, a_n, q$  and  $pAa_1, a_1Aa_2, \dots, a_nAq$

(Outer) Region Boundary  $\partial\mathcal{Q} = \{q \in \mathcal{D} \setminus \mathcal{Q} : \exists p \in \mathcal{Q} : qAp\}$ , i.e. the boundary  $\partial\mathcal{Q}$  of  $\mathcal{Q}$  is the set of pixels being adjacent to at least one pixel of  $\mathcal{Q}$  but not belonging to  $\mathcal{Q}$

Extremal Region  $\mathcal{Q} \subset D$  is a region such that for all  $p \in \mathcal{Q}, q \in \partial\mathcal{Q} : I(p) > I(q)$  (maximum intensity region) or  $I(p) < I(q)$  (minimum intensity region)

Maximally Stable Extremal Region (MSER). Let  $\mathcal{Q}_1, \dots, \mathcal{Q}_{i-1}, \mathcal{Q}_i, \dots$  be a sequence of nested extremal regions, i.e.  $\mathcal{Q}_i \subset \mathcal{Q}_{i+1}$ . Extremal region  $\mathcal{Q}_i$  is maximally stable iff  $q(i) = |\mathcal{Q}_{i+\Delta} \setminus \mathcal{Q}_{i-\Delta}| / |\mathcal{Q}_i|$  has a local minimum at  $i^*$  ( $|\cdot|$  denotes cardinality).  $\Delta \in \mathcal{S}$  is a parameter of the method

- *Stability*, since only extremal regions whose support is virtually unchanged over a range of thresholds is selected.
- *Multi-scale detection*. Since no smoothing is involved, both very fine and very large structure are detected.
- The set of all extremal regions can be *enumerated* in  $O(n \log \log n)$ , where  $n$  is the number of pixels in the image.

Enumeration of extremal regions proceeds as follows. First, pixels are sorted by intensity. The computational complexity of this step is  $\mathcal{O}(n)$  if the cardinality of the set  $S$  of image intensities is small, e.g. the typical  $\{0, \dots, 255\}$ , since the sort can be implemented as BINSORT [17]. After sorting, pixels are placed in the image (either in decreasing or increasing order) and the list of connected components and their areas is maintained using the efficient union-find algorithm [17]. The complexity of our union-find implementation is  $\mathcal{O}(n \log \log n)$ , i.e. almost linear<sup>1</sup>. Importantly, the algorithm is very fast in practice. The MSER detection takes only 0.14 s on a Linux PC with the Athlon XP 1600 + processor for an  $530 \times 350$  image ( $n = 185,500$ ).

The process produces a data structure storing the area of each connected component as a function of intensity. A merge of two components is viewed as termination of existence of the smaller component and an insertion of all pixels of the smaller component into the larger one. Finally, intensity levels that are local minima of the rate of change of the area function are selected as thresholds producing MSER. In the output, each MSER is represented by position

<sup>1</sup> Even faster (but more complex) connected component algorithms exist with  $O(n\alpha(n))$  complexity, where  $\alpha$  is the inverse Ackerman function;  $\alpha(n) \leq 4$  for all practical  $n$ .

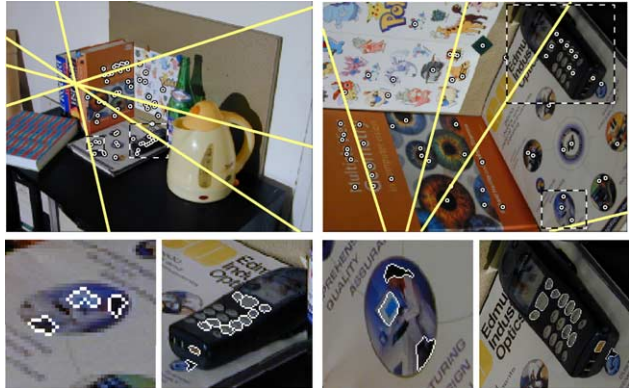


Fig. 1. Bookshelf. Estimated epipolar geometry on indoor scene with significant scale change. In the cutouts, the change in the resolution of detected DRs is clearly visible.

of a local intensity minimum (or maximum) and a threshold. Examples of MSERs are shown in Figs. 1, 2 and 5.

Notes. Although the set of extremal regions is covariant with any one-to-one continuous transformation of the image domain and thus covariant to projective transformation, the process of the selection of the maximally stable subset is affine-covariant. The MSERs are therefore only affine-covariant.

The structure of the above algorithm and of an efficient watershed algorithm [22] is essentially identical. However, the structure of the output of the two algorithms is different. The watershed is a partitioning of  $\mathcal{D}$ , i.e. a set of regions  $\mathcal{R}_i : \bigcup \mathcal{R}_i = D, \mathcal{R}_j \cap \mathcal{R}_k = \emptyset$ . In watershed computation, focus is on the thresholds where regions merge (and two watersheds touch). Such threshold are of little interest here, since they are highly unstable—after merge, the region area jumps. In MSER detection, we seek a range of thresholds that leaves the watershed basin effectively unchanged. Detection of MSER is also related to thresholding. Every extremal region is a connected component of a thresholded image. However, no global or ‘optimal’ threshold is sought, all thresholds are tested and the stability of the connected components evaluated. The output of the MSER detector is not a binarized image. For some parts of the image, multiple



Fig. 2. Valbonne. Estimated epipolar geometry and points associated to the matched regions are shown in the first row. Cutouts in the second row show matched bricks.

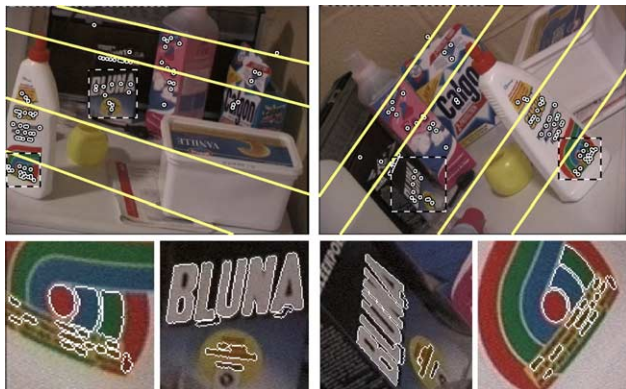


Fig. 3. Wash. Epipolar geometry and dense matched regions with fully affine distortion.

stable thresholds exist and a system of nested subsets is output in this case. Finally we remark that MSERs can be defined on any image (even high-dimensional) whose pixel values are from a totally ordered set.

### 3. The proposed robust wide-baseline algorithm

*Distinguished region detection.* As a first step, the DRs are detected—the MSERs computed on the intensity image (MSER+) and on the inverted image (MSER−).

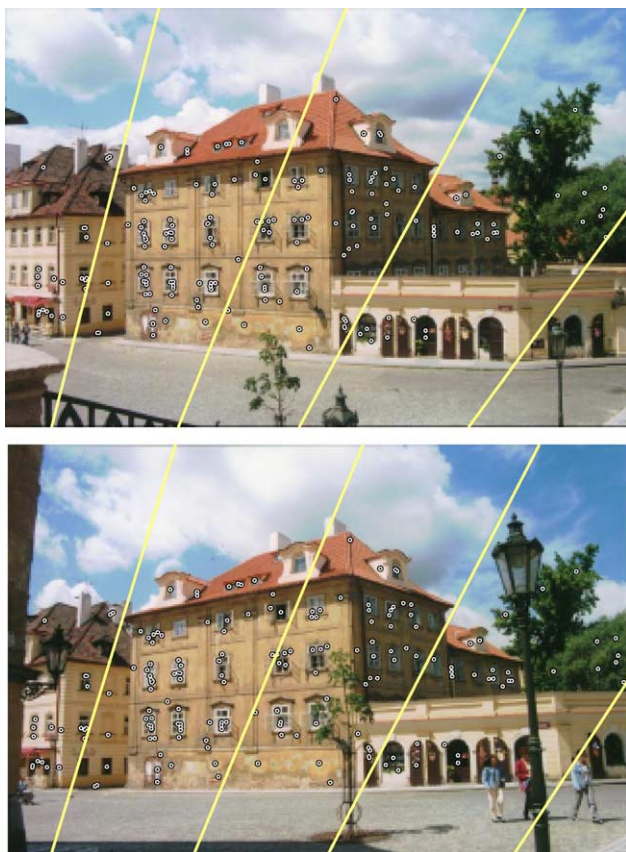


Fig. 4. Estimated EG on an outdoor scene.

*Measurement regions.* A measurement region of arbitrary size may be associated with each DR, if the construction is affine-covariant. Smaller measurement regions are both more likely to satisfy the planarity condition and not to cross a discontinuity in depth or orientation. On the other hand, small regions are less discriminative, i.e. they are much less likely to be unique. Increasing the size of a measurement region carries the risk of including parts of background that are completely different in the two images considered. Clearly, the optimal size of a MR depends on the scene content and it is different for each DR. In [21], Tuytelaars and Van Gool double the elliptical DR to increase discriminability, while keeping the probability of crossing object boundaries at an acceptable level.

In the proposed algorithm, measurement regions are selected at multiple scales: the DR itself, 1.5, 2 and 3 times scaled convex hull of the DR. Since matching is accomplished in a robust manner, we benefit from the increase of distinctiveness of large regions without being severely affected by clutter or non-planarity of the DR's pre-image. This is a novelty of our approach. Commonly, Mahalanobis distance has been used in MR matching. However, the non-robustness of this metric means that matching may fail because of a single corrupted measurement (this happened in the experiments reported below).

*Invariant description.* In all experiments, rotational invariants (based on complex moments) [8] were used after applying a transformation that diagonalises the regions covariance matrix of the DR. In combination, this is an affinely invariant procedure. Combination of rotational and affinely invariant generalised colour moments [12] gave a similar result. On their own, the affine invariants failed on problems with a large scale change.

*Robust matching.* A measurement taken from an almost planar patch of the scene with stable invariant description will be referred to as a 'good measurement'. Unstable measurements or those computed on non-planar surfaces or



Fig. 5. Cylindrical box. Epipolar geometry (top) and matched regions (bottom left). Fully affine distortion, a non-planar object, textured surface and strong specular reflections are present in the scene. SHOUT (bottom right), a scene with a change of illumination spectral power distribution.

at discontinuities in depth or orientation will be referred to as ‘corrupted measurements’.

The robust similarity is computed as follows. For each measurement  $M_A^i$  on region  $A$ ,  $k$  regions  $B_1, \dots, B_k$  from the other image with the corresponding  $i$ th measurement  $M_{B_1}^i, \dots, M_{B_k}^i$  nearest to  $M_A^i$  are found and a vote is cast suggesting correspondence of  $A$  and each of  $B_1, \dots, B_k$ . The votes are summed over all measurements.

The DRs with the largest number of votes are the candidates for tentative correspondences. Experimentally, we found that  $k$  set to 1% of the number of regions gives good results. The number of regions is typically in the  $10^2 - 10^3$  range and  $k$  is thus between 1 and 10. In the current implementation 216 invariants at each scale, i.e. a total of 864 measurements are used ( $i \in [1, 864]$ ,  $i$  runs through all scales and all invariants). The 216 rotational invariants are described in detail in [8]. The choice of four scales was made by trial and error and as a compromise between speed and performance.

Probabilistic analysis of the likelihood of the success of the procedure is not simple, since the distribution of invariants and their noise is image-dependent. We therefore only suppose that corrupted measurements spread their votes randomly, not conspiring to create a high score and that good measurements are more likely to vote for correct matches.

*Tentative correspondences using correlation.* Invariant description is used as a preliminary test. The final selection of tentative correspondences is based on correlation. First, transformations that diagonalise the covariance matrix of the DRs are applied. The resulting circular regions are correlated (for all relative rotations). This procedure is done efficiently in polar coordinates for different sizes of circles.

Rough EG is estimated by applying RANSAC to the centres of gravity of DRs. Subsequently, the precision of the EG estimate is significantly improved by the following process. First, an affine transformation between pairs of potentially corresponding DRs, i.e. the DRs consistent with the rough EG, is computed. Correspondence of covariance matrices defines an affine transformation up to a rotation. The rotation is determined from epipolar lines [2]. Next, DR correspondences are pruned and only those with correlation of their transformed images above a threshold are selected. In the next step, RANSAC is applied again, but this time with a very narrow threshold. The final improvement of the EG is achieved by adding to RANSAC inliers DR pairs whose convex hull centres are EG-consistent. Commonly, DRs differ in minute differences that render their centres of gravity inconsistent with the fine EG, but the centres of the convex hulls are precise enough. The precision of the final EG, estimated linearly by the eight point algorithm (without bundle adjustment or radial distortion correction) is surprisingly high. The average distance of inliers from epipolar line is below 0.1 pixel, see Table 3.

## 4. Experiments

The following experiments were conducted:

*Bookshelf*, (Fig. 1). The BOOKSHELF scene tests performance under a very large scale change. The corresponding DRs in the left view are confined only to a small part of the image since the rest of the scene is not visible in the second view. Different resolution of detected features is evident in the close-up.

*Valbonne*, (Fig. 2). This outdoor scene has been analysed in the literature [13,14]. Repetitive patterns such as bricks are present. The part of the scene visible in both views covers a small fraction of the image.

*Wash*, (Fig. 3). Results on this image set have been presented in [21]. The camera undergoes significant translation and rotation. The ordering constraint is notably violated, objects appear on different backgrounds.

*Kampa*, (Fig. 4), is an example of an urban outdoor scene. A relatively large fraction of the images is covered by changing sky. Repeating windows made matching difficult.

*Cylindrical box*, (Fig. 5, top and bottom left), shows a metal box on a textured floor. The regions matched on the box demonstrate performance on a non-planar surface. A significant change of illumination and a strong specular reflection is present in the second image that was taken with a flash (this strongly decreases the number of MSER + ).

*Shout*, (Fig. 5, bottom right). This scene has been used in [21]. Since the spectral power distribution of the illumination and the position of light sources is significantly different, we included the test to demonstrate performance in variable illumination conditions.

Results are summarized in Tables 2 and 3. Table 2 shows the number of detected DRs in the left  $\times$  right images for both types of the DRs (MSER – and MSER +). The number of tentative correspondences is given in the last column of Table 2. Table 3 shows the number of correspondences established in different stages of the algorithm. Column ‘TC’ repeats the number of tentative correspondences. Column ‘rough EG’ displays the number of tentative correspondences consistent with the rough estimate of the EG. The ratio of ‘TC’ and ‘rough EG’ determines the speed of the RANSAC algorithm.

Table 2  
Numbers of DRs detected in the left and right images in the ‘left DRs’  $\times$  ‘right DRs’ format

No. of	MSER –	MSER +	TC
Bookshelf	511 $\times$ 908	349 $\times$ 488	85
Valbonne	906 $\times$ 1012	761 $\times$ 950	49
Wash	1026 $\times$ 714	542 $\times$ 448	171
Kampa	1015 $\times$ 914	659 $\times$ 652	303
Cyl. box	1043 $\times$ 627	788 $\times$ 39	63
Shout	298 $\times$ 348	80 $\times$ 93	151

The number of tentative correspondences is given in the TC column.

Table 3  
Experimental results

	TC	Rough EG	Rough $d_{\perp}$	EG + corr	Fine EG	Fine $d_{\perp}$	Miss
Bookshelf	85	25	0.48	151	63	0.09	1
Valbonne	49	27	0.17	180	82	0.08	0
Wash	171	42	0.34	220	86	0.08	2
Kampa	303	78	0.34	422	185	0.08	2
Cyl. box	63	23	0.15	102	67	0.09	3
Shout	151	44	0.43	220	86	0.08	1

For details, see the text at the beginning of Section 4.

After establishing the ‘rough EG’ the so-called ‘guided matching’ step is applied [2,5]. In the process of finding tentative correspondences, at most a single corresponding region is associated with one DR. Often this association is erroneous, for instance if there is a repetitive pattern in the scene. Moreover some DRs are not matched at all since they fail the ‘mutually nearest’ requirement. Given the ‘rough EG’, even rather imprecise, the process of finding tentative matches can be revisited. The original tentative correspondences are discarded and all potential matches consistent with the ‘rough EG’ are selected. The matching now need not rely on rotational invariants, since epipolar lines passing through a pair of matching regions define their relative orientation [2]. The ‘guided’ tentative correspondences are therefore selected using correlation.

The column headed ‘EG + corr’ gives the number of correspondences consistent with rough EG that passed the correlation test. Notice that the numbers are much higher than those in the ‘rough EG’ column. The final number of correspondences is given in the penultimate column ‘fine EG’. Average distances from epipolar lines are presented in columns ‘rough  $d_{\perp}$ ’ and ‘fine  $d_{\perp}$ ’. We can see, that the precision of the estimated EG is very high, much higher than the precision of the rough EG. The last column shows the number of mismatches (found manually).

## 5. Conclusions

A new method for wide-baseline matching was proposed. The three main novelties are: the introduction of MSERs, robust matching of local features and the use of multiple scaled measurement regions.

The MSERs are sets of image elements, closed under the affine transformation of image coordinates and invariant to affine transformation of intensity. An efficient (near linear complexity) and practically fast detection algorithm was presented. The stability and high utility of MSERs was demonstrated experimentally. Another novelty of the approach is the use of a *robust similarity measure* for establishing tentative correspondences. Due to the robustness, we were able to consider invariants from *multiple*

*measurement regions*, even some that were significantly larger (and hence probably discriminative) than the associated MSER.

Good estimates of EG were obtained on challenging wide-baseline problems with the robustified matching algorithm operating on the output produced by the MSER detector. The average distance from corresponding points to the epipolar line was below 0.09 of the inter-pixel distance. Significant change of scale ( $3.5 \times$ ), illumination conditions, out-of-plane rotation, occlusion, locally anisotropic scale change and 3D translation of the viewpoint are all present in the test problems. Test images included both outdoor and indoor scenes, some already used in published work.

## Acknowledgements

The authors were supported by the European Union project IST-2001-32184, by the Grant Agency of the Czech Republic project GACR 102/02/1539 and by the Austrian Ministry of Education project CONEX GZ 45.535. The SHOUT and WASH images were kindly made available by Tinne Tuytelaars.

## References

- [1] A. Baumberg, Reliable feature matching across widely separated views, in: CVPR'00, 2000, pp. 1:774–781.
- [2] O. Chum, T. Werner, T. Pajdla, Joint orientation of epipoles, in: Proceedings of BMVC'03, vol. 1, BMVA, London, UK, September 2003, pp. 73–82.
- [3] Y. Dufournaud, C. Schmid, R. Horaud, Matching images with different resolutions, in: CVPR'00, 2000, pp. 1:612–618.
- [4] W.E.L. Grimson, Object Recognition, MIT Press, Cambridge, MA, 1990.
- [5] R. Hartley, A. Zisserman, Multiple View Geometry in Computer Vision, Cambridge University Press, Cambridge, UK, 2000.
- [6] T. Lindeberg and J. Gårding, “Shape-adapted smoothing in estimation of 3-D depth cues from affine distortions of local 2-D structure”, in Proc. 3rd European Conference on Computer Vision, vol. 800 of Lecture Notes in Computer Science, (Stockholm, Sweden), pp. 389–400, Springer Verlag, 1994.
- [7] D. Lowe, Object recognition from local scale-invariant features, in: ICCV'99, 1999, pp. 1150–1157.
- [8] J. Matas, P. Bílek, O. Chum, Rotational invariants for wide-baseline stereo, in: Proceedings of CVWW'02, February 2002, pp. 296–305.
- [9] J. Matas, Š. Obdržálek, O. Chum, Local affine frames for wide-baseline stereo, in: Proceedings of ICPR, vol. 4, IEEE CS, August 2002, pp. 363–366.
- [10] K. Mikolajczyk, C. Schmid, Indexing based on scale invariant interest points, in: Eighth International Conference on Computer Vision, Vancouver, Canada, 2001.
- [11] K. Mikolajczyk, C. Schmid, An affine invariant interest point detector, in: Proceedings of ECCV, vol. 1, 2002, pp. 128–142.
- [12] F. Mindru, T. Moons, L. van Gool, Recognizing color patterns irrespective of viewpoint and illumination, in: CVPR'99, 1999, pp. 1:368–373.

- [13] P. Pritchett, A. Zisserman, Wide baseline stereo matching, in: Proceedings of 6th International Conference on Computer Vision, Bombay, India, January 1998, pp. 754–760.
- [14] F. Schaffalitzky, A. Zisserman, Viewpoint invariant texture matching and wide baseline stereo, in: Eighth International Conference on Computer Vision, Vancouver, Canada, 2001.
- [15] F. Schaffalitzky, A. Zisserman, Multi-view matching for unordered image sets, or ‘How do I organize my holiday snaps?’, in: Proceedings of ECCV’02, vol. 1, Springer, Berlin, 2002, pp. 414–431.
- [16] C. Schmid, R. Mohr, Local grayvalue invariants for image retrieval, PAMI 19 (5) (1997) 530–535.
- [17] R. Sedgewick, Algorithms, second ed., Addison-Wesley, Reading, MA, 1988.
- [18] D. Tell, S. Carlsson, Wide baseline point matching using affine invariants computed from intensity profiles, in: ECCV’00, 2000.
- [19] P. Torr, A. Zisserman, Robust parameterization and computation of the trifocal tensor, in: BMVC’96, page Motion and Active Vision, 1996.
- [20] T. Tuytelaars, L.V. Gool, Content-based image retrieval based on local affinely invariant regions, in: Proceedings of Third International Conference on Visual Information Systems, 1999, pp. 493–500.
- [21] T. Tuytelaars, L.V. Gool, Wide baseline stereo based on local, affinely invariant regions, in: M. Mirmehdi, B. Thomas (Eds.), Proceedings of the British Machine Vision Conference BMVC’00, London, UK, 2000, pp. 412–422.
- [22] L. Vincent, P. Soille, Watersheds in digital spaces: an efficient algorithm based on immersion simulations, IEEE Transactions on Pattern Analysis and Machine Intelligence 13 (6) (1991) 583–598.