# Efficient Sequential Correspondence Selection by Cosegmentation

Jan Čech, *Member, IEEE,* Jiří Matas, *Member, IEEE,* Michal Perďoch, *Member, IEEE*

**Abstract**—In many retrieval, object recognition and wide baseline stereo methods, correspondences of interest points (distinguished regions) are commonly established by matching compact descriptors such as SIFTs. We show that a subsequent cosegmentation process coupled with a quasi-optimal sequential decision process leads to a correspondence verification procedure that (i) has high precision (is highly discriminative) (ii) has good recall and (iii) is fast. The sequential decision on the correctness of a correspondence is based on simple statistics of a modified dense stereo matching algorithm. The statistics are projected on a prominent discriminative direction by SVM. Wald's sequential probability ratio test is performed on the SVM projection computed on progressively larger cosegmented regions. We show experimentally that the proposed Sequential Correspondence Verification (SCV) algorithm significantly outperforms the standard correspondence selection method based on SIFT distance ratios on challenging matching problems.

**Index Terms**—Correspondence, matching, verification, sequential decision, growing, cosegmentation, stereo, image retrieval, learning.

---

## 1 INTRODUCTION

MANY successful image retrieval, object recognition and wide baseline stereo methods exploit correspondences of distinguished regions[1]. Most real-world visual recognition problems are large scale where correspondences between regions from a query (test) image and many database (training) images of objects or scenes are sought. To achieve acceptable response times, large problems require the time complexity of the region matching process be sublinear in the size of the database; memory footprint of the database representation becomes a concern too. The standard solution is to describe regions with a compact descriptor such as SIFT [1] or some discretization of it (e.g. "visual words" [2]) and to store database image representations in a search tree (k-d [1], metric [3], k-means [4], [5], [6]).

The matching process typically proceeds as follows [7], [8], [1]. Distinguished regions are detected in the image and local affine or similarity covariant coordinate frames are constructed for each region. *Measurement regions*, i.e. image patches of typically rectangular or elliptical shape, are specified in terms of the local coordinate frames. For each region, a descriptor (such as SIFT) is computed from the signal in the measurement region, after both photometric and geometric normalization. Additionally, the descriptor may be compressed by quantization.

This process, schematically visualized in Fig. 2, has the

• *The authors are with Center for Machine Perception of Department of Cybernetics at Faculty of Electrical Engineering, Czech Technical University, Prague, Czech Republic.*
*E-mail: {cechj,matas,perdom1}@cmp.felk.cvut.cz*

1. Terms "transformation-covariant regions", "viewpoint invariant features", "interest points", "salient points" and "patches" also appear in the literature. We adopt the term "distinguished region" as a concise shortcut for the self-explanatory but unwieldy "repeatably-detectable transformation-covariant regions".
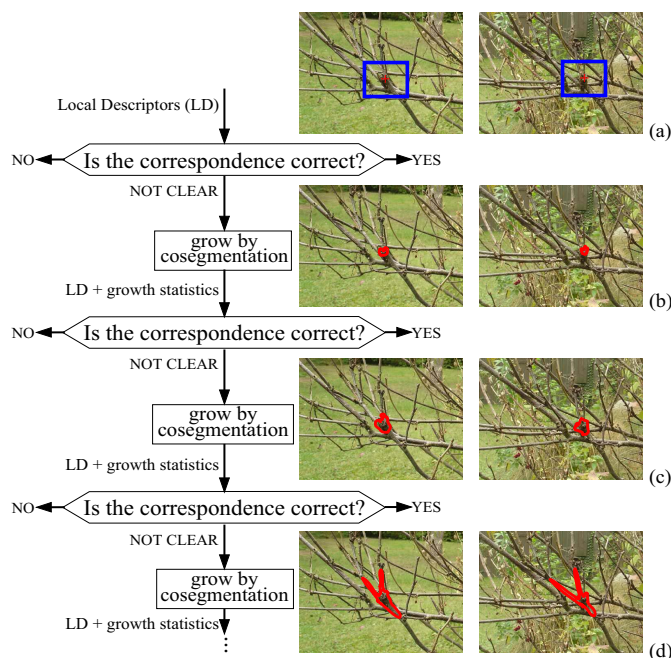


Fig. 1. The basic idea of the Sequential Correspondence Verification (SCV) algorithm. First, a decision is attempted on the basis of local descriptors (LD) computed from fixed shape measurement region (in blue, a). If the correctness of the correspondence cannot be decided reliably, the correspondence is successively grown by cosegmentation (in red, b-d), collecting additional evidence for the decision.

following main characteristics: (i) all steps are performed in individual images independently, (ii) the shape and size of the measurement region is a fixed function of the shape and size of the distinguished region and (iii) the descriptor has the same form for all regions, e.g. it is a vector in $R^d$. These properties facilitate fast sublinear region matching, e.g. via search tree or hashing.
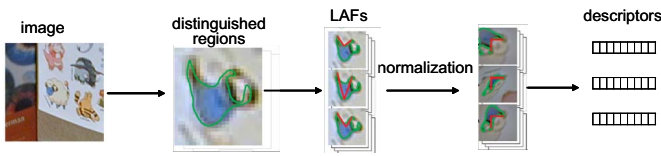
This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

2

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE



Fig. 2. Describing a distinguished region by a compact descriptor invariant to local geometric and photometric changes. "LAFs" stands for "Local Affine Frames".



(a) measurement region too large and/or with unsuitable shape

(b) measurement region too small

Fig. 3. Problems with a fixed shape and size of the measurement region.
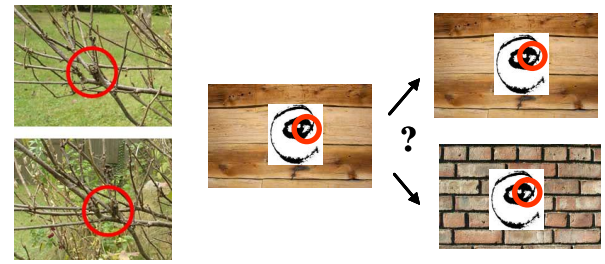
However, the fixed size and shape of the measurement region necessarily involves a compromise. In general, the larger the measurement region, the more discriminative the information inside is. On the other hand, large measurement regions may violate the local planarity assumption of wide-baseline matching methods, and are more likely to straddle object boundary or to be affected by occlusion. Moreover, they are more sensitive to localization errors of local frames. For "non-compact" objects, such as elongated and wiry ones, the compact fixed shape is problematic.

Consider, for instance, the two images depicted in Fig. 3(a). Only a very small circular or rectangular regions around the distinguished region on the branch will not include signal from the background, which is different for the two views. On the other hand, consider the images shown in Fig. 3(b). The measurement region inside the circle is too small since any descriptor computed from the region will be close to identical for both images on the right - the correct match cannot be reliably established.

A better estimate of correspondence quality (a prediction of it being correct) can be obtained by looking at both test and training image simultaneously, e.g. by attempting to expand the correspondence domains, which is illustrated in Fig. 1. The value of correspondence growing methods has been demonstrated in [9], [10], sometimes with impressive results, e.g. those achieved by the dual bootstrap method [11], [12]. Most approaches to simultaneous cosegmentation and registration focus on the problem of finding the largest corresponding domain and co-domain [11], [10], [13], [14].

Our objective is almost the opposite: given acceptable false positive and false negative rates, design the fastest possible test for correctness of a correspondence, based on cosegmentation of regions of progressively growing size. We formulate the problem as sequential decision making which is solved by performing Wald's sequential probability ratio test. The test is based on simple statistics of a modified dense stereo matching algorithm which are projected on a single prominent discriminative direction by a linear SVM.

Of course, we do not want to lose the excellent large-scale matching properties of descriptors based on measurement regions of fixed size and scale. The cosegmentation process is therefore only applied to *tentative correspondences* obtained by a sub-linear process, such as kD-tree search. In fact, if followed by correspondence

verification, any such process for generating tentative correspondences can be set to be much more permissive, outputting higher number of correspondences with lower inlier ratios but containing larger number of inliers. After filtering by simultaneous cosegmentation, inlier ratios are (more than) recovered and the larger number of inliers leads to higher recognition rates. We show on challenging problems that the selection of correspondences based on sequential cosegmentation is very efficient, runs near to real-time and significantly outperforms the standard correspondence process based on SIFT distance ratios, producing a higher number as well as higher percentage of correct correspondences.

Consequently, combinatorial procedures for estimation of a geometrically consistent subset of correspondences with time complexity sensitive to inlier ratios (polynomial dependence), e.g. RANSAC, should always adopt sequentially terminated cosegmentation as a pre-processing step.

The method scales well: the number of potential correspondences for a query image region can be controlled. If it is constant, the total time complexity of the region expansion process is independent of the size of the database and linear in the size of the input (number of regions in the query image). On a large scale retrieval experiment [4], we observed that the time needed to carry out the sequential procedure is not significant in comparison with the time needed for the initial indexing process for establishing tentative correspondences.

The rest of the paper is organized as follows. The method is described in Sec. 2 and the training data and learning procedure for the sequential classifier in Sec. 3. Experimental validation is presented in Sec. 4. Conclusions are summarized in Sec. 5.

This paper is a significantly extended and modified version of [15].

## 2 THE SEQUENTIAL CORRESPONDENCE VERIFICATION ALGORITHM

The motivation of the approach is to distinguish, as fast as possible, correct and incorrect correspondences via
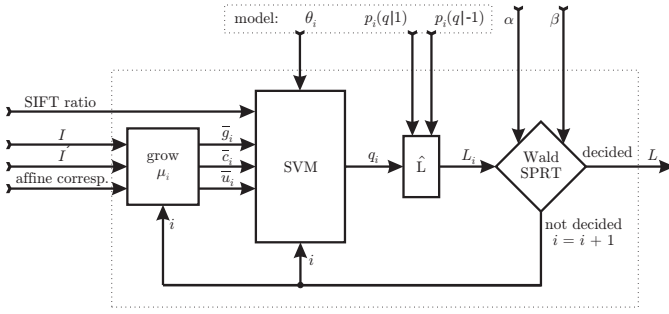
This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

ČECH *et al.*: EFFICIENT SEQUENTIAL CORRESPONDENCE SELECTION BY COSEGMENTATION 3



Fig. 4. The Sequential Correspondence Verification algorithm.

---

**Algorithm 1** Sequential Correspondence Verification (scv)

**Require:** images $\mathbf{I}, \mathbf{I}'$,
  correspondence with affine frame $(x, y, \mathbf{A})$,
  SIFT ratio $s_r$,
  false positive and false negatives rates $(\alpha, \beta)$,
  model: learned SVM parameters $\theta_i$,
      likelihoods $\mathbf{p}_i(q|+1)$, $\mathbf{p}_i(q|-1)$.

1.1: init_grow$(x, y, \mathbf{A})$.
1.2: **for** $i := 1$ **to** max. number of decision stages **do**
1.3: $\quad \mu_i = \begin{cases} 0, & i=1, \\ \gamma^{i-2}, & i>1. \end{cases}$
1.4: $\quad (\bar{g}_i, \bar{c}_i, \bar{u}_i) := \text{grow}(\mathbf{I}, \mathbf{I}', \mu_i)$.
1.5: $\quad q_i := \text{SVM}(s_r, \bar{g}_i, \bar{c}_i, \bar{u}_i, \theta_i)$.
1.6: $\quad L_i := \frac{\mathbf{p}_i(q_i|+1)}{\mathbf{p}_i(q_i|-1)}$.
1.7: $\quad$ **if** Wald SPRT$(L, \alpha, \beta)$ is conclusive **then break**.
1.8: **end for**
1.9: **return** likelihood ratio $L_i$ (of the last iteration).

---

dense matching, i.e. by a pixel-to-pixel correspondence growing algorithm. The requirements of high speed and quality of the decision process are contradictory. We therefore propose a quasi-optimal sequential decision algorithm that minimizes time to decision, given user-specified probabilities of false positive and false negative rates.

The two error rates control in an intuitive way the trade-off between the number of correspondences (high when false negative rate is low), inlier ratio (high when false positive rate is low) and the speed of the method (low for low error rates). Importance of these three factors differs in applications.

The proposed Sequential Correspondence Verification algorithm Alg. 1 (SCV) is overviewed in Fig. 4. The basic idea (see Fig. 1) is to perform Wald's Sequential Probability Ratio Test, having learned the distributions of elementary statistics generated by the growing process.

The algorithm proceeds in decision stages indexed by $i$. In the first stage, a fast dense stereo matching algorithm, Sec. 2.1, is initialized by a tentative correspondence of a pair of local affine frames. The verification proceeds by attempting to match discriminative, i.e. high variance, neighboring pixels. After a certain number of growing steps $\mu_i$, the cosegmentation process returns three simple statistics $(\bar{g}_i, \bar{c}_i, \bar{u}_i)$ characterizing the quality of the correspondence: the growth rate $\bar{g}_i$ – the size of the grown region divided by the maximum number of attempted growing steps $\mu_i$, the average correlation $\bar{c}_i$ of the region, and the average number of pixels violating the uniqueness $\bar{u}_i$, i.e. non-bijectivity matching.

The vector of statistics is projected by a linear SVM to a scalar quantity $q_i$ which avoids estimation of high-dimensional class conditional probabilities (likelihoods). Instead only the likelihoods of the projections $\mathbf{p}_i(q| + 1)$, $\mathbf{p}_i(q| - 1)$ of correct and incorrect correspondence classes are computed.

The region statistics are augmented with the first to the second nearest SIFT descriptor distance ratio $s_r$, a standard measure for selection of tentative correspondences [1]. We call $s_r$ the *SIFT ratio*. The Wald's Sequential Probability Ratio Test (SPRT) is performed on the likelihood ratio $L_i = \mathbf{p}_i(q_i| + 1)/\mathbf{p}_i(q_i| - 1)$. If the SPRT test is conclusive, the algorithm terminates and the correspondence is assigned the likelihood ratio $L_i$ of the de-

cision. Otherwise, another decision stage $i$ is performed, i.e. the cosegmentation is resumed with a new limit of attempted growing steps $\mu_i$, Alg. 1, Step 1.3, potentially producing more discriminative statistics, since it is based on more measurements. Note that $\mu_1 = 0$, which means the decision in the first stage is based solely on the SIFT ratio without growing. The process continues until the maximum number of decision stages $i$ is reached.

In our experiments, we set the maximum number of decision stages to 100 and the largest pixel growth is $\mu_{100} = 1000$ steps. We observed the error does not decrease after a larger growth, see Fig. 10(b). In principle, it is possible to perform Wald's SPRT test after each growing step, which would be the fastest strategy if the test execution time was zero. However, the statistics are average values, therefore in order to have a constant influence of new measurements, we propose growing in steps of a geometric sequence instead, see Step 1.3. The number of decision stages was set empirically to minimize the decision time based on our implementation. These considerations determine $\gamma$ in Step 1.3.

The choice of the three statistics of the growing process was driven mainly by computational requirements. The statistics were selected from a larger pool of easily computable characteristics. We experimented e.g. with geometric deviation from the transformation implied by the local affine frame correspondence, mean intensity difference of the regions, difference of Harris-like cornerness values. None of them discriminated well.

Beside benefiting computational speed, we attribute the good generalization of the sequential classifier to the simplicity of the characterization of the growing process. For instance, the sequential classifier performed almost equally well on correspondences established on distinguished regions other than those it was trained on, despite the fact that it is unlikely that image statistics in general are the same for different detection processes. Moreover, even after a very modest number of training examples, the classifier performed well on a very large

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

4

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE

---

**Algorithm 2** Initialize growth state (`init_grow`)

---

**Require:** affine correspondence $(x,y,\mathbf{A})$.
2.1: Initialize
    seeds: $\mathcal{S} := (x, y, \mathbf{A})$,
    matching maps: $\mathbf{T}(:,:) := 0$, $\mathbf{T}'(:,:) := 0$,
    counters: $K := G := C := U := 0$.
2.2: Compute the image correlation for all seeds $\mathbf{s} \in \mathcal{S}$.

---

test set. With satisfactory performance, we did not further investigate the feature selection problem (including feature number).

The linear SVM was our first choice of a projection method as it possessed the desirable properties of fast learning, fast execution and trivial implementation. Since it lead to a good classification rates, following the structural risk minimization principle, we did not test more complex classifiers.

### 2.1 Growing algorithm

The following algorithm explores regions around the input tentative correspondence. The growing mechanism is inspired by [16], [17], [13], [18].

Each correspondence defines a local affine mapping from the reference image $\mathbf{I}$ to the target image $\mathbf{I}'$. The mapping generates several pixel-to-pixel correspondences $\mathbf{s} = (x, y, \mathbf{A})$, where $(x, y)$ is a point in $\mathbf{I}$ with associated affine transformation $\mathbf{A}$ which maps the local neighborhood to the other image $\mathbf{I}'$:

$$x' = a_1 x + a_2 y + a_3,$$
$$y' = a_4 x + a_5 y + a_6, \tag{1}$$

or simply $(x', y') = \mathbf{A}(x, y)$.

The procedure is presented in pseudo-code as Alg. 3. The inputs are the images $\mathbf{I}, \mathbf{I}'$, the set of initial correspondences, the seeds $\mathcal{S}$ and the maximum number of growing steps $\mu$. The outputs consists of three statistics $\bar{g}, \bar{c}, \bar{u}$ which characterize the (in)correctness of the input correspondence. The growing algorithm keeps its state: the set of correspondence seeds $\mathcal{S}$, matching maps $\mathbf{T}, \mathbf{T}'$ and the counters $K, G, C, U$ that is initialized in Alg. 1, Step 1.1. In the initialization (Alg. 2), the image correlation $\text{corr}(\mathbf{s})$ of all initial seeds[2] $\mathbf{s} \in \mathcal{S}$ is computed, Step 2.2, as Moravec's normalized cross-correlation [22] of a $5 \times 5$ pixel window $\mathbf{w}$ centered at pixel $(x, y)$ in the reference image and window $\mathbf{w}'$ centered at $\mathbf{A}(x, y)$ in the target image, deformed according to the affinity[3] $\mathbf{A}$:

$$\text{corr}(\mathbf{s}) = \frac{2\text{cov}(\mathbf{w}, \mathbf{w}')}{\text{var}(\mathbf{w}) + \text{var}(\mathbf{w}')}, \tag{2}$$

---

2. In our experiments, this is realized by local affine frames (LAF) constructed on Maximally Stable Extremal Region [19], [20] (MSER) and Hessian Affine points [21]. We take the three point-to-point correspondences of a pair of LAFs as the initial seeds of the growing process.

3. The simplest and fastest interpolation, the nearest-neighbor, was used to compute $\mathbf{w}'$. More advanced interpolation did not bring a significant improvement.

---

**Algorithm 3** The Growing Algorithm (`grow`)

---

**Require:** images $\mathbf{I}$, $\mathbf{I}'$, maximum number of growing steps $\mu$, growth state.
3.1: **while** $K \le \mu$ and $\mathcal{S}$ not empty **do**
3.2:     $K := K + 1$.
3.3:     Draw the seed $\mathbf{s} \in \mathcal{S}$ of the best similarity $\text{corr}(\mathbf{s})$.
3.4:     **for** each of the best neighbors $\mathbf{t}_k^*$ in $\mathcal{N}_k(\mathbf{s})$:
        $\mathbf{t}_k^* = (x, y, \mathbf{A}) = \underset{\mathbf{t} \in \mathcal{N}_k(\mathbf{s})}{\arg\max} \text{corr}(\mathbf{t}), \ k \in \{1, 2, 3, 4\}$
    **do**
3.5:         $c := \text{corr}(\mathbf{t}_k^*)$,
3.6:         **if** $c \ge \tau$ and $\mathbf{T}(x, y) = 0$ **then**
3.7:             $G := G + 1$, $C := C + c$.
3.8:             **if** $\mathbf{T}'(\mathbf{A}(x, y)) = 1$ **then**
3.9:                 $U := U + 1$.
3.10:             **end if**
3.11:             Update the matching maps
            $\mathbf{T}(x, y) := \mathbf{T}'(\mathbf{A}(x, y)) := 1$ and
3.12:             the seed queue $\mathcal{S} := \mathcal{S} \cup \{\mathbf{t}_k^*\}$.
3.13:         **end if**
3.14:     **end for**
3.15: **end while**
3.16: **return** statistics: growth rate $\bar{g} := \frac{G}{\mu}$, average correlation $\bar{c} := \frac{C}{G}$, average uniqueness violation $\bar{u} := \frac{U}{G}$.

---

where $\text{cov}(\mathbf{w}, \mathbf{w}')$ is a covariance, $\text{var}(\mathbf{w})$ a variance.

Set $\mathcal{S}$ is organized as a priority queue according to the correlation. A seed is removed from the top of the queue, and for all its 4-neighbors (left, right, up, down) in the reference image, the best correlating candidate in $\mathcal{N}_k(s)$ is found (out of 9 possible positions in the target image), Alg. 3, Step 3.4, such that

$$\mathcal{N}_1(s) = \big\{(x - 1, y, \mathbf{A}_{c-1,r}) \mid c, r \in \{-1, 0, 1\}\big\},$$
$$\mathcal{N}_2(s) = \big\{(x + 1, y, \mathbf{A}_{c+1,r}) \mid c, r \in \{-1, 0, 1\}\big\},$$
$$\mathcal{N}_3(s) = \big\{(x, y - 1, \mathbf{A}_{c,r-1}) \mid c, r \in \{-1, 0, 1\}\big\},$$
$$\mathcal{N}_4(s) = \big\{(x, y + 1, \mathbf{A}_{c,r+1}) \mid c, r \in \{-1, 0, 1\}\big\}, \tag{3}$$

where

$$\mathbf{A}_{c,r} = \begin{bmatrix} a_1 & a_2 & a_3 + a_1 c + a_2 r \\ a_4 & a_5 & a_6 + a_4 c + a_5 r \end{bmatrix}. \tag{4}$$

If the highest correlation exceeds threshold $\tau$ and the point is unmatched so far in the reference image, then a new match is found, Step 3.6. Next, the counter for the region size $G$ is incremented and correlation value $c$ is added to sum $C$. If the pixel in the target image $\mathbf{I}'$ is already matched, the counter for uniqueness violation $U$ is incremented, Step 3.9. The binary matching maps $\mathbf{T}$ and $\mathbf{T}'$ are updated and the found match becomes a new seed. Up to four seeds are created in each growing step.

The process continues until there are no seeds in the queue or the algorithm is stopped when reaching the maximum number of growing steps $\mu$, Step 3.1.

We set correlation threshold $\tau = 0.5$, which was found experimentally, Fig. 9(b). On our training set, we observed the classification error rate for several growing

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

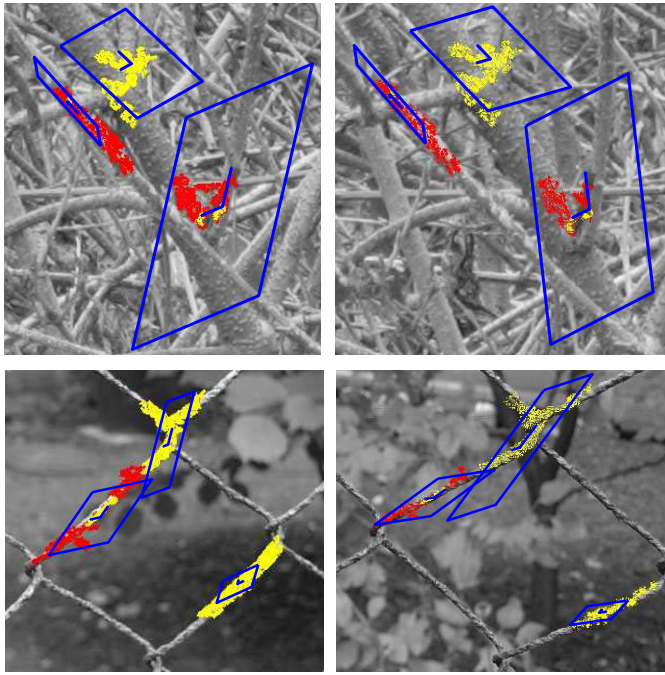ČECH *et al.*: EFFICIENT SEQUENTIAL CORRESPONDENCE SELECTION BY COSEGMENTATION 5

Fig. 5. Examples of region growth in the cosegmentation process. Selected pixels at the time of decision on correctness (in yellow) and after 1000 growing steps (in red). The red pixels are shown for visualization purposes only, their correspondence is not evaluated by the SCV algorithm. Corresponding local affine frames and measurement regions of SIFT descriptors are in blue.

steps $\mu$ as a function of $\tau$. The error is constant up to $\tau = 0.5$, then it rises. Lower thresholds make the decision process slower, due to worthless growth of incorrect correspondences.

Examples of region growth in the cosegmentation process are depicted in Fig. 5. Local affine frames are shown in blue as a pair of line segments. The three endpoints of the segments are the seed of the growing process. Blue parallelograms delineate measurement regions, i.e. parts of the image where SIFT descriptors are computed. Yellow marks pixels inside the region at the time of the decision on the correctness of the correspondence made by the SCV algorithm. Red marks pixels that would be chosen if the process was left to grow the maximum number of $\mu = 1000$ steps. Note that (i) the SCV decision is often reached after growing over a very small number of pixels and (ii) the shape of the region is data-dependent, preferring areas with edges and high variance of the signal where correlation response is high. Unsurprisingly, pixel correspondences follow correctly the 3D surfaces (branches of the shrub, parts of the fence); in fact, a small local disparity map is computed.

The measurement regions include large parts of the background that is different in the two images. It might be surprising that the regions in Fig. 5 are correctly matched, given that the first test in the sequential classifier is based on the SIFT ratio. We believe there are (at least) three reasons for the favorable outcome. First,

our test on the SIFT ratio is very permissive. Second, the centers of the regions are on corresponding 3D structures and SIFT applies a Gaussian weighting function that reduces influence of the outer parts of the parallelogram, which are not corresponding. Finally, SIFT is an array of histograms of gradients. The strong gradients are in correspondence in both pairs of images, areas without strong edges are irrelevant for the SIFT representation.

**Discussion.** Unlike Vedaldi and Soatto's region growing algorithm [9], Algorithm 3 includes no explicit regularization either of the mapping or of the shape of the cosegmented regions. The reason is that the algorithm grows only in informative areas with distinguishing signal (texture), so regularization is not needed. Areas without texture are ambiguous and do not help to distinguish correct and incorrect correspondences. Growth is restricted to unambiguous areas by requiring correlation statistic[4] to stay above threshold $\tau$, Step 3.6. Implicit surface smoothness is enforced. The disparity gradient change is constrained by (3), similar constraint is applied in [13].

In wide baseline dense stereo [17], [18], [23], local affine parameters $(a_1, a_2, a_4, a_5)$ representing a window deformation due to surface slant are optimized after each growing step in order to facilitate maximum growth on curved or projectively distorted surfaces. However, our goal is different: for correspondence verification the surface need not be grown too far. Therefore, in our algorithm, the parameters inherited from the initial seed are kept constant, which is significantly faster than the iterative optimization. Experiments show that a small imprecision of the local affine parameters is not critical, possibly due to the fact that effects of transformation errors are subsumed in disparity (gradients).

### 2.2 Statistical correspondence quality

Ideally, correspondence quality would be a function of the probability that a pair of grown patches is a projection of the same 3D surface, as calculated e.g. via MRF on the image grid by global methods in dense stereo [24]. However, finding the MAP solution is computationally intensive even for simple fields. Therefore, we use the efficient growing algorithm as a suboptimal solution and model the correspondence quality on the basis of elementary statistics that were empirically shown to discriminate correct and incorrect correspondences.

The class conditional probability densities of the adopted statistics are shown in Fig. 6. We observed that the growth rate $\bar{g}$ is typically larger for correct correspondences than for incorrect as reported by [9], exceptions include e.g. situations when a correct correspondence lies on a narrow surface or in cases of

---

4. Note, the Moravec's statistic is a zero mean normalized correlation, see (2). For areas without texture, after subtracting the mean values of signals in windows, the rest is an uncorrelated noise which results in a low value of the statistic.
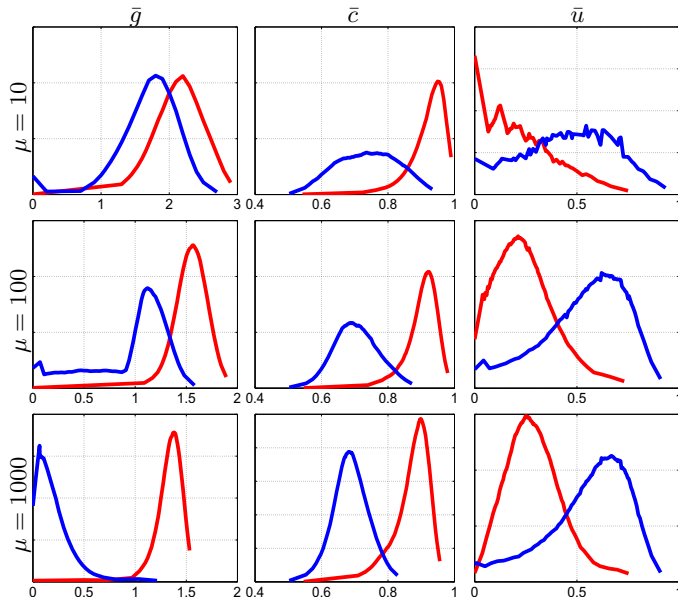
This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

6

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE

Fig. 6. Estimated class conditional probability density functions of cosegmentation statistics. From left to right: growth rate $\bar{g}$, average correlation $\bar{c}$, average uniqueness violation $\bar{u}$; from top to bottom: for 10, 100, 1000 growth cycles $\mu$. Correct correspondences (in red), incorrect correspondences (in blue). Note the gradual reduction in the overlap of the densities, especially in the two leftmost columns.
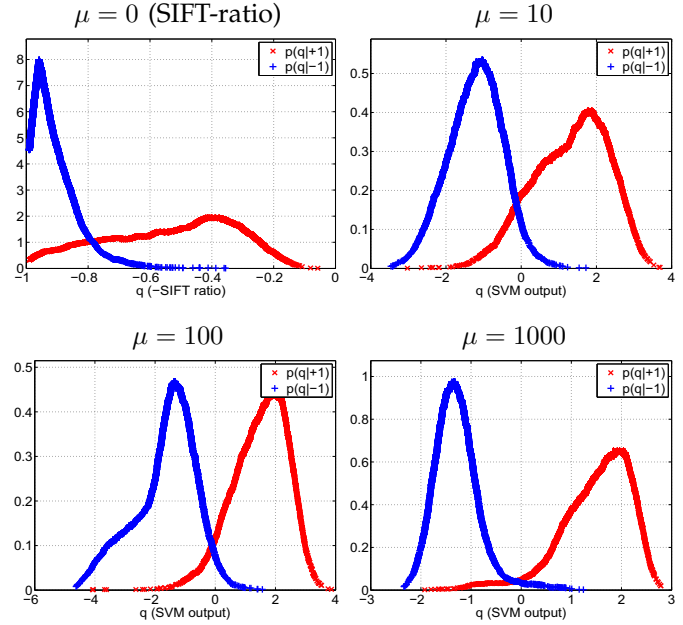


Fig. 7. Estimated class conditional probability densities for the oriented distance to the SVM hyperplane, i.e. for projections on a normal to the maximum margin hyperplane, for correct (red) and incorrect (blue) correspondences, for four stages of the sequential decision process. Note the decrease in the overlap of the distributions with increasing $\mu$.

partial occlusion. The average correlation in the region $\bar{c}$ is also typically higher for correct correspondences, but incorrect correspondences may accidentally have high correlation on repetitive or locally similar structures, especially on small regions. The average uniqueness violation $\bar{u}$ (deviation from bijective matching) when growing the region is also quite discriminative, see Fig. 6, right column. The statistics returned by the growing algorithm are combined with the ratio of the first to second closest distance of SIFT descriptors $s_r$ [1], a standard method.

The problem of estimating a high dimensional likelihood ratio is avoided by projecting the four dimensional feature vector onto a 1D scalar quantity $q_i = \mathrm{f}(s_r, \bar{g}_i, \bar{c}_i, \bar{u}_i)$ which expresses a confidence on correctness of the correspondence. This is done using a Support Vector Machine (SVM) trained on a set of positive and negative correspondences, see Sec. 3. The SVM finds a discriminative direction maximizing a margin in combination with a hinge loss (the training data are not separable). Projecting on this direction is an effective feature extraction procedure, suggested already by Vapnik [25] and popularized by e.g. Platt [26].

In consecutive decision stages $i$, the statistics are more discriminative, with the increase in the maximum number of growing steps $\mu_i$, Step 1.3. Thus a different SVM $\theta_i$ is trained for each decision stage $i$. The classification error due to the overlap of probability distributions is progressively decreasing, see plot in Fig. 10(b).

The likelihoods $\mathbf{p}_i(q|+1)$ and $\mathbf{p}_i(q|-1)$ of positive and negative class respectively were estimated by Parzen window method with a moving average kernel. The likelihoods estimated from our training set are shown for four decision stages $i$ in Fig. 7. In the first stage, there is no growth and the statistic is solely the SIFT ratio. Interestingly, the SIFT ratio threshold of $0.8$ suggested for accepting a correspondence by Lowe [1] is confirmed, being close to the equal-error operating point. Note that in the sequential process a significantly stricter test is applied in the first stage: only correspondences having SIFT ratio smaller than about $0.4$ are immediately accepted as correct, while the others are grown and decided in a later stage of a cascade, see Fig. 7 (top-left).

The likelihood ratio $L_i$, given the SVM output $q_i$, is computed using linearly interpolated estimates of class conditional probability.

## 2.3 Wald's sequential decision

Let $x$ be an object belonging to one of two classes $\{-1, +1\}$. In our case, the classified objects are correspondences and the classes are "correct" (1) and "incorrect" (-1). Next, let an ordering on the set of measurements $\{x_1, \ldots, x_n\}$ on $x$ be given. Here measurements $x_i(= q_i)$ are scalar values, oriented distances from SVM decision boundaries after growing step $i$.

A sequential decision strategy is a set of decision functions $S = \{S_1, \ldots, S_n\}$, where $S_i : \{x_1, \ldots, x_i\} \rightarrow \{-1, +1, \sharp\}$. The strategy $S$ makes one measurement at a time. The '$\sharp$' sign stands for a "continue" (do not decide yet). If a decision is '$\sharp$', $x_{i+1}$ is obtained and $S_{i+1}$ is evaluated. Otherwise, the output of $S$ is the class returned by $S_i$.
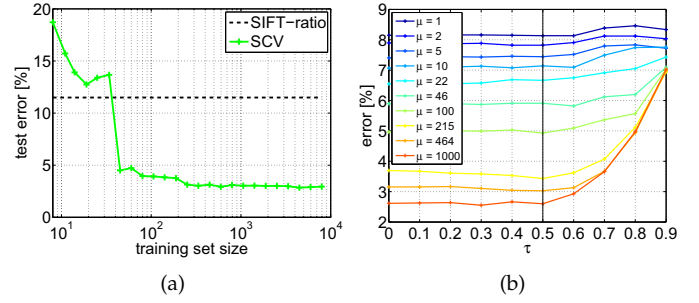
Fig. 8. The set of training images.



Fig. 9. (a) Test error as a function of a training set size. Note that only a very small training set with about 300 examples is required. (b) Cross-validated classification error as a function of the correlation threshold $\tau$ of the growing process plotted after several growing steps $\mu$. The higher the value of $\tau$, the faster the decision process. The selected value of $\tau = 0.5$ (highlighted) is the highest value with error close to minimal for all stages.

In two-class classification problems, errors of two kinds can be made by strategy $S$. Let us denote $\alpha_S$ the probability of rejecting a correct correspondence ($x$ belongs to $+1$ but is classified as $-1$) and $\beta_S$ the probability of accepting an incorrect correspondence ($x$ belongs to $-1$ but is classified as $+1$). A sequential strategy $S$ is characterized by its error rates $\alpha_S$ and $\beta_S$ and its average evaluation time $\bar{T}_S = E(T_S(x))$ where the expectation is over $p(x)$, and $\bar{T}_S$ is the expected evaluation time (or time-to-decision) for strategy. An optimal strategy for the sequential decision making problem is then defined as

$$S^* = \arg\min_S \bar{T}_S \qquad (5)$$
$$\text{s.t.} \quad \beta_S \leq \beta,$$
$$\alpha_S \leq \alpha$$

for specified $\alpha$ and $\beta$.

Wald [27] proved that the solution of the optimization problem (5) is the *sequential probability ratio test*.

**Sequential Probability Ratio Test.** Let $x$ be an object characterized by its hidden state (class) $y \in \{-1, +1\}$. The decision about the hidden state is based on successive measurements $x_1, x_2, \ldots$. Let the joint conditional density $p(x_1, \ldots, x_m | y = c)$ of the measurements $x_1, \ldots, x_m$ be known for $c \in \{-1, +1\}$.

SPRT is a sequential strategy $S^*$, which is defined as

$$S_m^* = \begin{cases} +1, & L_m \geq A \\ -1, & L_m \leq B \\ \sharp, & B < L_m < A \end{cases} \qquad (6)$$

where $L_m$ is the likelihood ratio

$$L_m = \frac{p(x_1, \ldots, x_m | y = +1)}{p(x_1, \ldots, x_m | y = -1)}. \qquad (7)$$

The constants $A$ and $B$ are set according to the required error of the first kind $\alpha$ and error of the second kind $\beta$. Optimal $A$ and $B$ are difficult to compute in practice, but tight bounds are easily derived. It can be shown that setting the thresholds $A$ and $B$ to

$$A = \frac{1 - \beta}{\alpha}, \qquad B = \frac{\beta}{1 - \alpha} \qquad (8)$$

is close to optimal [27].

In the SCV algorithm, we assume that all information about a correspondence is contained in the statistics from the last growth step: $p(q_i | y) = p(q_1, \ldots, q_i | y)$. Therefore only 1D PDFs are needed to carry out the SPRT test. Estimation of scalar PDFs poses no technical problems as discussed in previous section.

## 3 THE TRAINING PROCEDURE

The set of 24 image pairs used in a training set of correspondences is shown in Fig. 8. For all image pairs, MSERs were detected, LAFs constructed [20], [19] and SIFT descriptors computed on normalized patches. Standard wide-baseline matching was performed using SIFTs and a set of tentative correspondences was obtained. Finally, RANSAC was run on each pair of the set to estimate the epipolar geometry. We have manually re-labeled correspondences which were accidentally consistent with the epipolar geometry but were in fact incorrect[5]. The remaining inlier correspondences formed the positive subset of the training set, all other correspondences were inserted as negative examples.

Approximately 6200 positive and 9800 negative correspondences were obtained, which means that tentative correspondences on the training set had on average approximately 40% of inliers.

The ground truth set was split randomly into two equal parts, half for training, half for testing. The learning stage included linear SVM training and probability

---

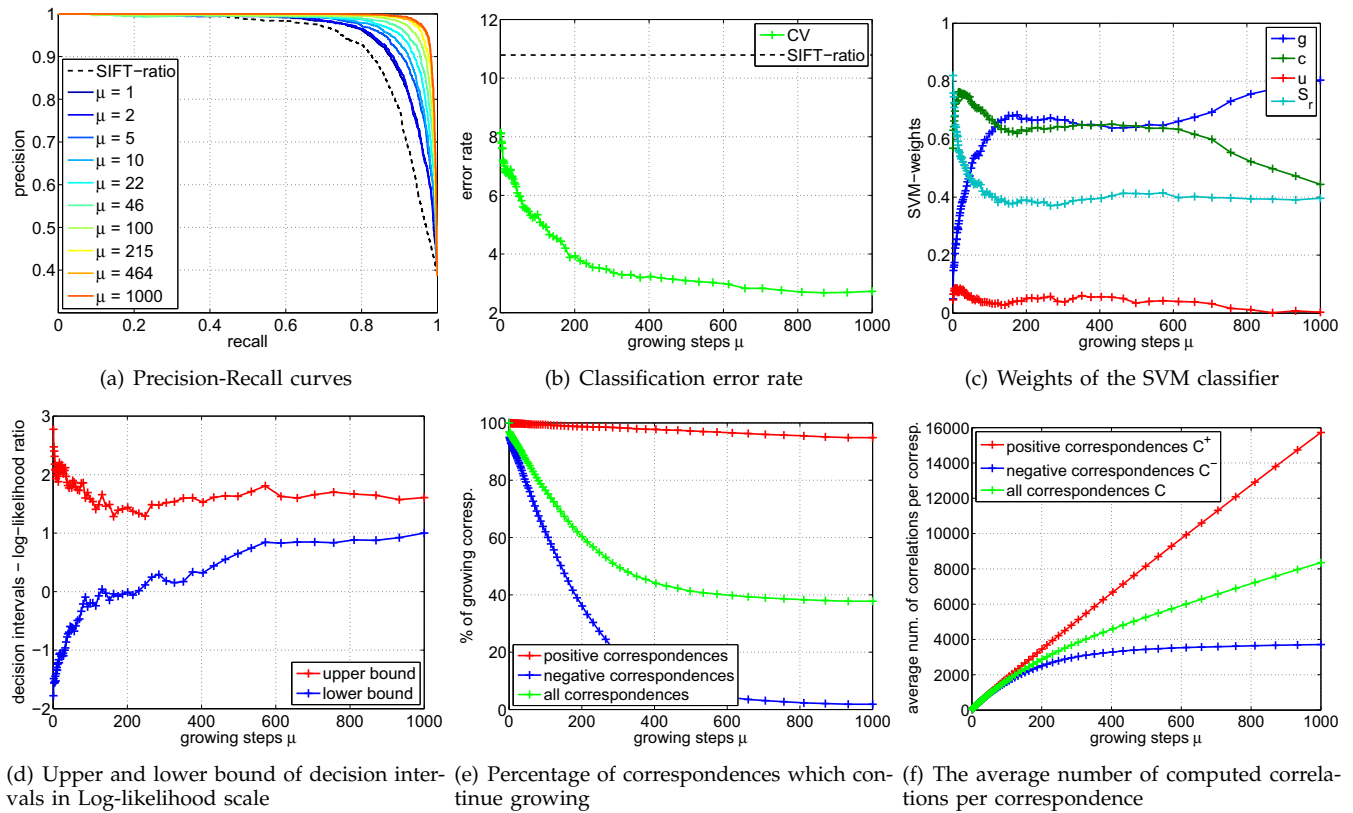5. A mismatch in tentative correspondence lying on a corresponding pair of epipolar lines is not detected by RANSAC.

(a) Precision-Recall curves

(b) Classification error rate

(c) Weights of the SVM classifier

(d) Upper and lower bound of decision inter- (e) Percentage of correspondences which con- (f) The average number of computed correla-
vals in Log-likelihood scale tinue growing tions per correspondence

Fig. 10. Properties of the SCV algorithm as a function of the number of growing steps $\mu$.

density estimation via Parzen windowing. For SVM learning, we used a recently published efficient algorithm [28]. Regularization constant $C$ of SVM with the hinge-loss criterion was estimated by crossvalidation (which set $C = 1$).

A priori, we had no idea about the necessary size of a training set for the correspondence classification problem. We therefore carried out the following test. A progressively larger portion of the training set was used to estimate the SVMs and likelihoods. The resulting sequential classifier (SCV) was applied to the test set. The observed classification error is plotted as a function of the training set size in Fig. 9(a). For reference, the classification error of the SIFT ratio is plotted too. The error does not improve significantly after about 300 samples, a surprisingly small number. We concluded that the size of our training set is sufficient.

Note that the Wald SPRT is a non-Bayesian technique based on conditional probabilities and its performance guarantees in terms of false positive and false negative rates hold for arbitrary prior probabilities. The insensitivity to the prior probability of (in)correct correspondence is an important property of the method, since wide range of inlier ratios is encountered in practical matching problems. In fact, the SCV procedure is extremely useful for matching problems where tentative correspondences have a very low inlier ratio and direct RANSAC application would require an astronomical number of samples. Such problems differ significantly

in inlier percentage of tentative correspondences from our training set, but in a non-Bayesian setting it does not matter. The training set must only be representative of the conditional probabilities of observations. Note also that all the parameters and learned models (SVM weights and likelihoods) were kept fixed throughout all the experiments.

## 4 EXPERIMENTS

### 4.1 Basic properties of the Sequential Correspondence Verification algorithm

We start the performance evaluation of the SCV algorithm by several experiments demonstrating some elementary properties of the algorithm. First, we measured discriminability of the SCV algorithm, i.e. its ability to distinguish correct and incorrect correspondences. The discriminability is characterized by a precision-recall curve which is computed as follows. The SCV algorithm assigns likelihood ratio $L$ to all $N$ correspondences in test set. The correspondences are sorted according to their likelihood ratio, $L_{(1)} \geq L_{(2)} \geq ... \geq L_{(N)}$. *Precision* is defined as $Q_n^+/n$, where $Q_n^+$ is the number of correct correspondences among $L_{(1)}, ..., L_{(n)}$. *Recall* is defined as $Q_n^+/Q_N^+$. The SCV algorithm is more discriminative than a standard ratio of SIFT descriptors, and the difference becomes more prominent with the number of growing steps $\mu$, Fig. 10(a).

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

ČECH *et al.*: EFFICIENT SEQUENTIAL CORRESPONDENCE SELECTION BY COSEGMENTATION 9
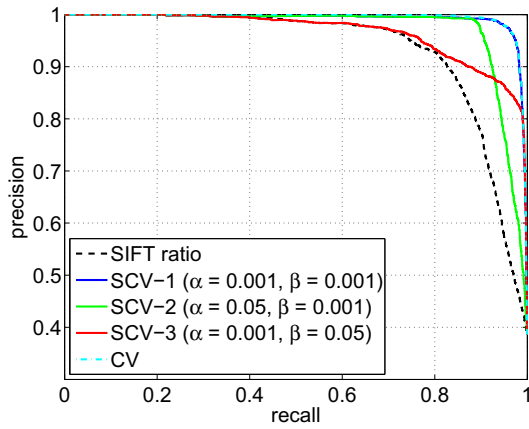


Fig. 11. Discriminability of the SCV algorithm. The precision-recall curves for SCV with various setting of false positive and false negative rates and for the SIFT ratio alone.

When a hard decision on the correctness of a correspondence is required, the likelihood ratio $L$ is thresholded. The classification error rate for the threshold $L = 1$ is plotted in Fig. 10(b). In the case of the SIFT-ratio the threshold is $0.8$ as discussed before.

Fig. 10(c) shows the evolution of the weights of the SVM classifier as a function of growing steps. The weights are the unit normal vector of the discriminative hyperplane, trained on the zero mean, unit variance normalized data. We show the absolute value of the weights which can be interpreted as an importance of a particular statistic. A large value means the statistic is important, while value near zero refers about a small influence. At the beginning, the SIFT-ratio $s_r$ has a large impact, together with the average correlation $\bar{c}$. With more growing steps, the weight of the growth rate $\bar{g}$ quickly rises and becomes the most import statistic. The weight of the average uniqueness $\bar{u}$ violation is quite low and decreases with more growing steps. The weight of the average correlation $\bar{c}$ also decreases which is probably due to an imprecise estimate of the LAFs which manifests itself by a lower correlation of correct correspondences far away from their seeds.

The next figure, Fig. 10(d), focuses on decision thresholds for Wald's SPRT; the upper and lower bounds of the indecision intervals for Wald's SPRT are plotted for $\alpha = 0.05$ and $\beta = 0.001$ in log-likelihood ratio scale. The undecided interval is shrinking with increasing growing steps $\mu$ due to lower error.

Fig. 10(e) shows the percentage of correspondences which are still growing after maximum number of growing steps $\mu$ is performed. Notice that almost all correct correspondences grow; 95% of correct correspondences can grow above the largest executed growth $\mu = 1000$. Incorrect correspondences stop growing much earlier since there are no high correlating neighbors and typically the algorithm finishes by exhausting the seed queue $S$ before the maximum number of growing steps is reached, see Alg. 3, Step 3.1.

Finally, Fig. 10(f) shows the average number of computed window correlations per correspondence. This quantity is closely related to the computational complexity of the algorithm. For correct correspondences, the number of correlations $C^+$ grows almost linearly with growing steps $\mu$, while for incorrect correspondences, the number of correlations $C^-$ saturates at $4000$. It means that for the largest growth $\mu = 1000$, negative correspondences are about four times faster to decide. This behavior is expected, since the algorithm stops growing the incorrect correspondences earlier, see Fig. 10(e).

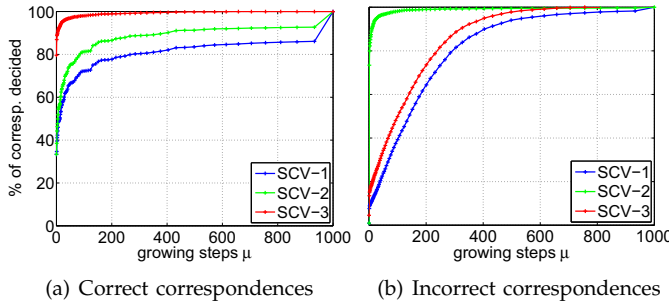## 4.2 The SCV efficiently increases discriminability

We show that the SCV algorithm is more discriminative than the SIFT ratio and that the sequential decision-making process speeds the algorithm significantly at the expense of a very small discriminability loss. The comparison of the SCV algorithm was carried out with various settings of Wald's SPRT parameters $(\alpha, \beta)$, Fig. 11. The SCV algorithm outperforms the SIFT ratio for all three settings. The SCV-1 ($\alpha = 0.001, \beta = 0.001$) is the most strict setting which has the highest discriminability. The SCV-2 ($\alpha = 0.05, \beta = 0.001$) allows more false negatives, while the SCV-3 ($\alpha = 0.001, \beta = 0.05$) more false positives, but they both are more efficient in terms of number of window correlations they had to compute.

In Fig. 12, three $(\alpha, \beta)$ settings of SCV algorithm are compared with the non-sequential version (CV), which does not decide until the last stage performing maximally $\mu = 1000$ growing steps. We measured the average number of window correlations per correspondence $C$ which had to be computed, and the percentage of correspondences decided (or stopped growing) in $i$-th stage of the algorithm after $\mu$ steps.

These values differ for correct and incorrect correspondences, so besides the mean values $C$ (which depends on the percentage of correct correspondences in the test set), we show the values for correct correspondences $C^+$ and wrong correspondences $C^-$ which differ as discussed in the previous subsection. The decision plots are shown for correct and incorrect correspondences as well.

The sequential decision speeds up the process by factor of more than two (SCV-1), more than five (SCV-3), or more than 9 in comparison to the non-sequential algorithm without losing much discriminability. The recall-precision curve in Fig. 11 of the non-sequential algorithm (CV) is almost identical to the SCV-1. Fig. 12 also shows that the SCV-2 with higher allowed false negative rate tends to decide negative correspondences in lower stages of the sequence speeding up the decision process by factor of more than 18, while the SCV-3 is speeding up the decision process for positive correspondences by factor of more than 100.

**Computational complexity.** The dominant operation in the SCV algorithm is correlation computation, other steps (SVM classification, Wald's SPRT) are negligible.

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

10

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE



(a) Correct correspondences      (b) Incorrect correspondences

| | $C/10^3$ | $C^+/10^3$ | $C^-/10^3$ |
|---|---|---|---|
| CV | 8.4 | 15.7 | 3.7 |
| SCV-1 | 3.5 | 3.3 | 3.6 |
| SCV-2 | 0.9 | 2.0 | 0.2 |
| SCV-3 | 1.7 | 0.1 | 2.7 |

Fig. 12. Efficiency of the SCV algorithm. Plots show the percentage of correspondences decided after growing $\mu$ steps. The table shows an average number of window correlations per correspondence.

The running time depends on the number of correspondences. Considering an example of 1000 tentative correspondences, each requiring on average $C = 900$ correlations, Fig. 12, we end up with approximately $10^6$ correlations per image pair; which is computed on recent CPU in about 0.25 seconds and about 20–100 times faster on a modern GPU. It usually takes about 0.5 second on a standard C2 2.4 GHz with our implementation, depending also on the ratio of correct correspondences, and on the setting of Wald's SPRT parameters.

## 4.3 SCV performance on Hessian affine points

Until now, all experiments have been carried out on correspondences of local affine frames on MSERs. We now show that SCV-algorithm performs equally well for verification of correspondences obtained from Hessian affine points [21].

For all training image pairs in Fig. 8, a set of tentative correspondences was generated from Hessian affine points and classified according to the ground-truth epipolar geometry, and split into a training and test set. This is the same procedure as described before for MSERs.

Fig. 13 shows that for Hessian affine points, the SCV algorithm (SCV-1) improves the recall-precision curve obtained by SIFT ratio matching. Moreover, the performance is virtually equal for the two cases when the algorithm is trained specifically on Hessian-affine points or when the SCV algorithm trained on MSERs is used. This was not expected a priori, as the image patches around the respective correspondences are quite different. But the (simple) statistics of the growth process initialized from pixel correspondences are preserved.

Interestingly, the ratio of the SIFT descriptors has slightly better discriminability on Hessian affine points
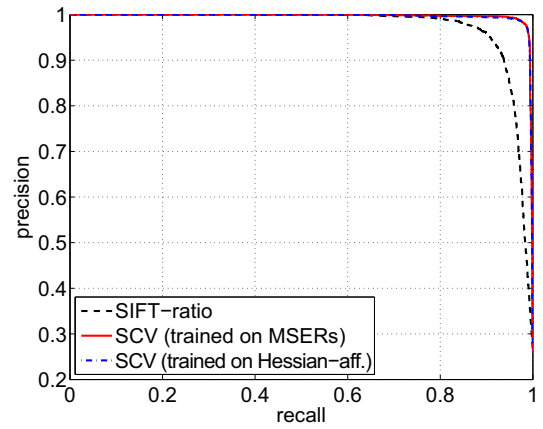


Fig. 13. Discriminability of the SCV algorithm for Hessian affine correspondences. The SCV algorithm performs equally well when trained on MSER correspondences or specifically on correspondences of Hessian affine points.

than on MSERs. The fact that MSERs are often detected on occlusion boundaries might play a role.

## 4.4 Challenging wide baseline stereo scenes

Results of correspondence selection on difficult wide baseline stereo scenes are shown in Fig. 14. These scenes are challenging due to small overlap, high degree of noise in the images (Raglan), complex 3D structure with many occlusions (Forsythia, Fence). In the Orange pair, matching is difficult since the background is locally similar (same grain of wood) but not the same (different location on the same table). To find the epipolar geometry at all, the matching process generating the tentative correspondences had to be very permissive, so that a sufficient number of correct correspondences was present among tentative correspondences. We allowed more than one-to-one mapping in tentative correspondences which lead to a high number of outliers (about 90 percent).

Plots in the last column of Fig. 14 show the precision among the best $n$ retrieved correspondences. This is important for progressive RANSAC procedure [29] which samples tentative correspondences according to preferences defined by the matching processes (approximately speaking in the order as sorted by the matcher of tentative correspondences). For correspondences sorted by the SCV algorithm, in all four scenes, the PROSAC procedure would terminate successfully after a single iteration, since a sufficient number of top correspondences is correct. This is neither the case when the ordering of tentative correspondences is given by the negative ratio of SIFT distances, nor the SIFT distances alone.

On the same images, we compared the sequential algorithm (SCV-2) and its non-sequential version (CV). For all the scenes, the results of SCV-2 are slightly worser than of non-sequential CV, but it is much faster. The two algorithms evaluated the following numbers of window correlations (SCV-2 vs. non-sequential CV): $0.5 \times 10^3$ vs. $2.5 \times 10^3$ (Raglan), $0.6 \times 10^3$ vs. $5.7 \times 10^3$ (Forsythia),
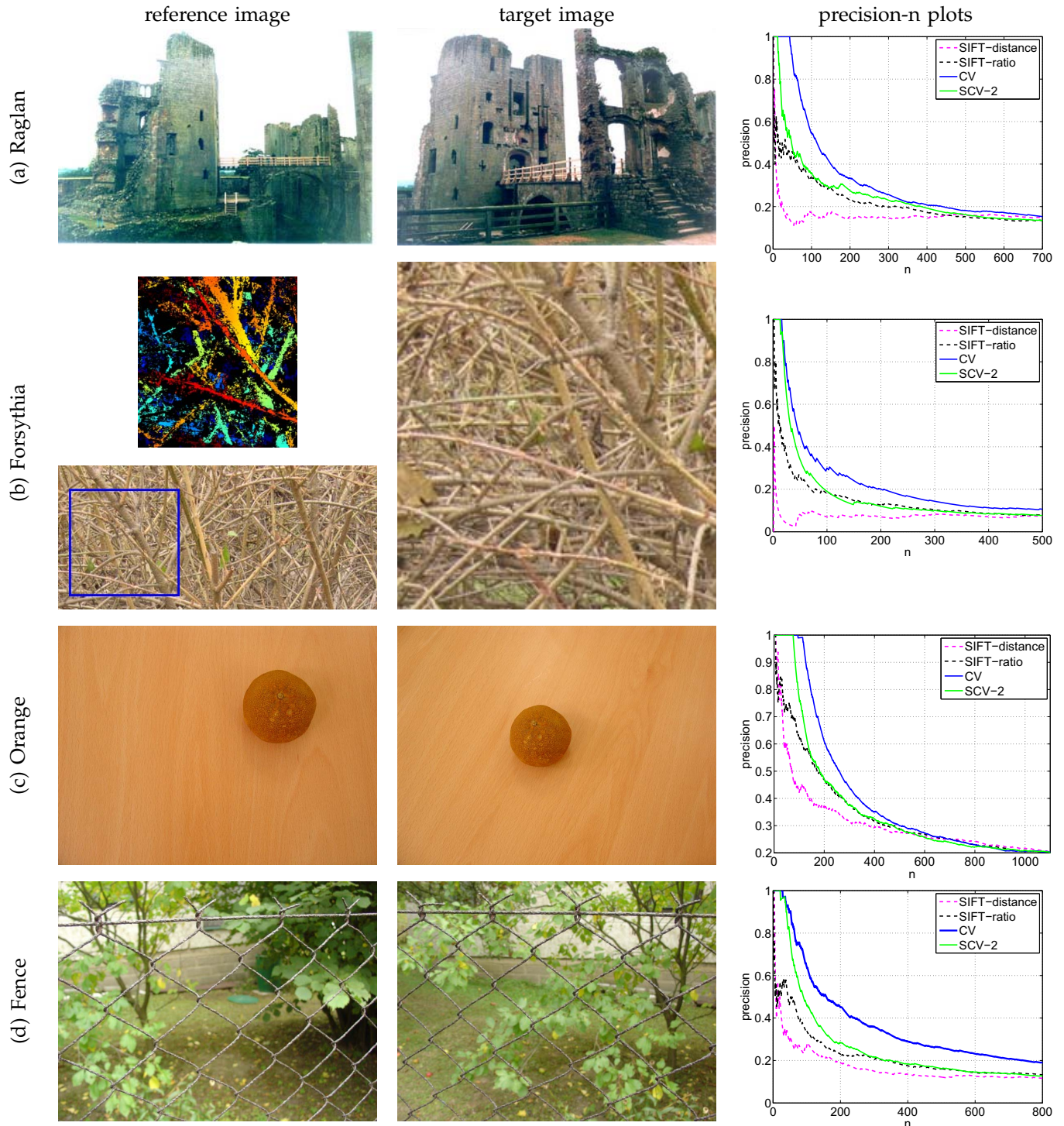
Fig. 14. Results on challenging wide-baseline scenes. For Forsythia, we show the color coded depth map of a common part (inside the blue frame) to demonstrate the 3D structure of the image pair. Notice that the Orange (c) at a different place on the table and no correct correspondence exist on the background. Although not obvious, almost the same part of the fence appears in both images of the Fence scene.
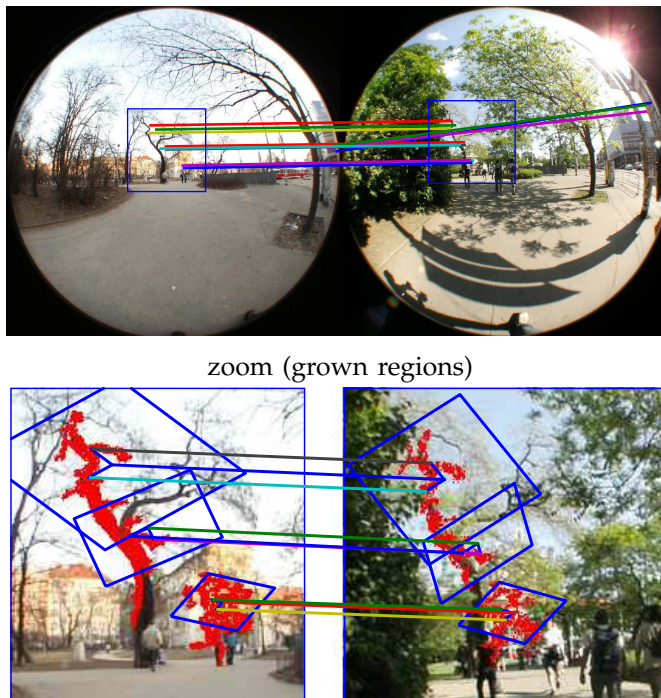
zoom (grown regions)



Fig. 15. Correspondences found by the SCV algorithm in challenging omnidirectional image pair with a 'wide temporal baseline'.

$1.2 \times 10^3$ vs. $6.3 \times 10^3$ (Orange), and $0.4 \times 10^3$ vs. $4.3 \times 10^3$ (Fence). The reason why the decision is even faster here than on the test set in the previous experiment is the high number of wrong correspondences which are faster to decide.

**Omnidirectional images.** The method was successfully tested on challenging image pairs like Fig. 15, obtained by a fish-eye camera. Besides a significant wide spatial baseline setup, the pair has a wide temporal baseline: the first image was captured in the winter when there were no leaves on trees, while the other in the summer at different day-time and in very different lighting conditions and shadows. Despite the difficult conditions, the SCV algorithm was able to find several correct correspondences, as can be checked visually; the ground-truth does not exist in this case. The algorithm selected four LAF correspondences (three correct, one incorrect) out of more than 2000 tentative correspondences. Note that the correspondences shown are selected solely by the SCV algorithm, i.e. before robust model fitting.

The method works because the initial affine transformation obtained from LAF correspondence locally approximates the non-linear neighborhood deformation in omnidirectional images. This is true for the central part of the images, while the problems occur at the boundary of spheres, where the distortion is not negligible. This is probably the reason for the incorrect correspondence which occurred close to boundary of the sphere.

## 4.5 Test on the Oxford dataset

We used a subset of the Oxford dataset[6] which has been used for performance evaluation of affine region detectors [30] and local descriptors [31]. The dataset consists of images which are distorted by various degradation: projective distortion (due to change of the camera position), image blur (from defocusing), JPEG compression artifacts and illumination changes. The ground-truth correspondences are known, since the database contains a homography mapping between the reference and distorted target images.

The input is a set of several hundreds tentative correspondences per each pair. The results for correspondence selection based on standard SIFT-ratio and on the SCV algorithm are shown in Fig. 16 as precision-recall curves. We can see the SCV algorithm is better in all cases. The most difficult distortion seems to be the blur, but it is destroying for SIFT as well. The projective distortion is well captured by local affine transformation, the illumination change does not also make serious problems, since the Moravec's correlation used in the growing algorithm is insensitive (however not fully invariant like NCC statistic) to affine illumination changes. Surprisingly, the SCV does not deteriorates that much with the JPEG compression. Although the images looks seriously (innaturally) corrupted, the frequencies preserved by the JPEG compression resulted in enough correlation.

## 4.6 Image retrieval

The benefits of SCV algorithm are demonstrated on a large scale image retrieval setup, using the data set from Nister and Stewenius benchmark [4]. It consists of 10200 images in groups of four that show the same object. In the benchmark experiment, each image becomes a query. For each query, the top $N$ images are returned, and a score is a computed that counts how many of the correct answers are in top $K$ images. In the benchmark $K$ is set to 4, giving the highest score 4, if the algorithm manages to retrieve as top four images the four instances of the object in the data set. Since the query image is also present in the data set, the worst score of algorithm returning only the query in top $K$ is 1. The overall performance of the algorithm is computed as the average score of all 10200 queries from the data set.

We reimplemented a part of the Nister's approach. MSERs [20] and LAFs [19] were computed on each of the images. Each of approximately 7.4 millions LAFs was described using SIFT [1] computed on an affine normalized patch. Then, similarly to visual words approach proposed by Sivic and Zisserman [2], we built visual a word vocabulary consisting of 1 million k-means in the SIFT descriptor space and assigned all the descriptors in the images to the nearest visual word. Each visual word in a given document is weighted using TFIDF (Term Frequency – Inverse Document Frequency) measure from

6. http://www.robots.ox.ac.uk/~vgg/research/affine/

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

ČECH *et al.*: EFFICIENT SEQUENTIAL CORRESPONDENCE SELECTION BY COSEGMENTATION
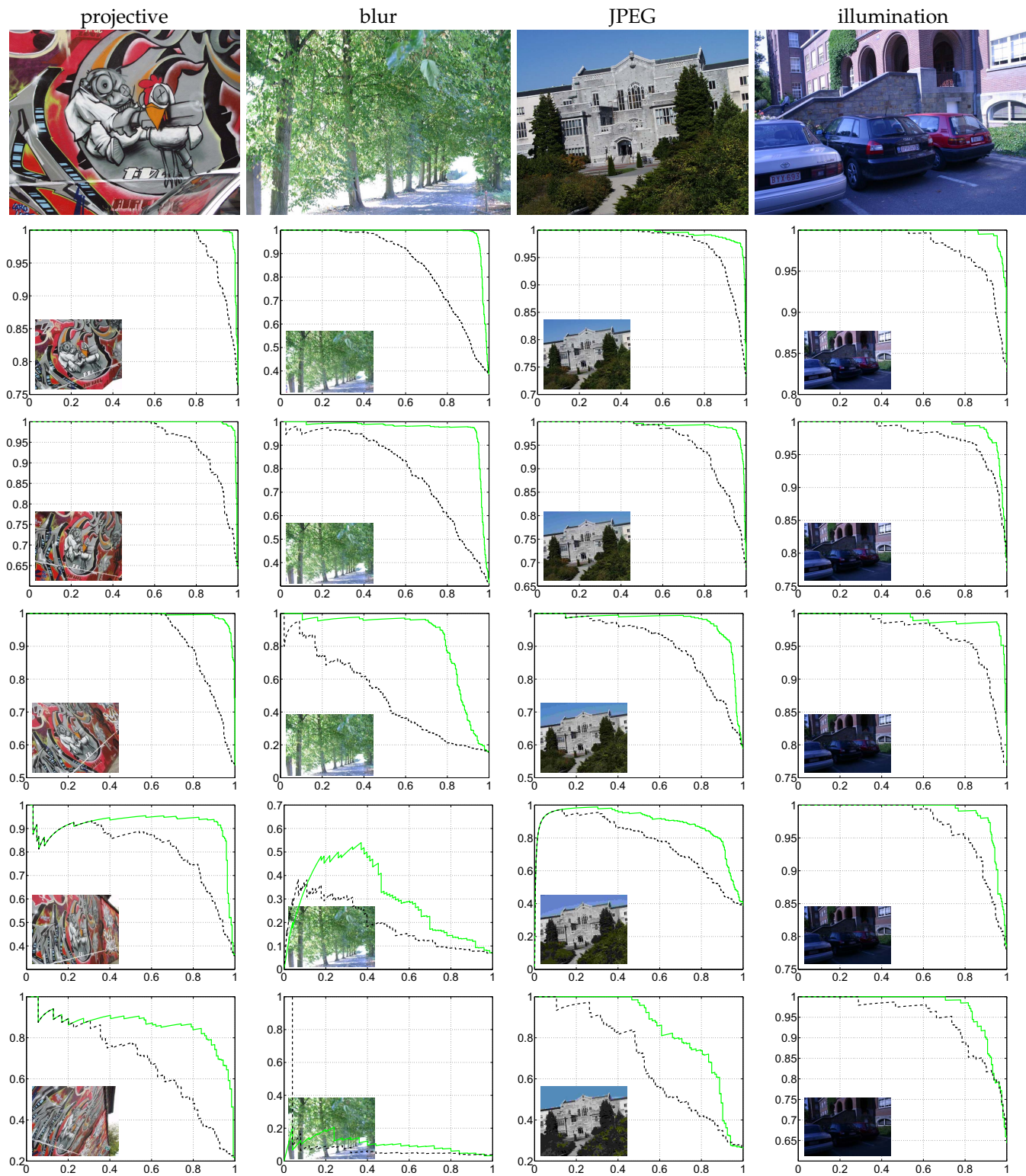13



Fig. 16. Results on the Oxford dataset. Precision-recall curves of the SCV algorithm (green) and of the SIFT ratio (black dashed), for images with an increasing degree of distortion. Reference images in the first row, target images are shown in graphs. See the electronic version of the article to better view of the distorted images.
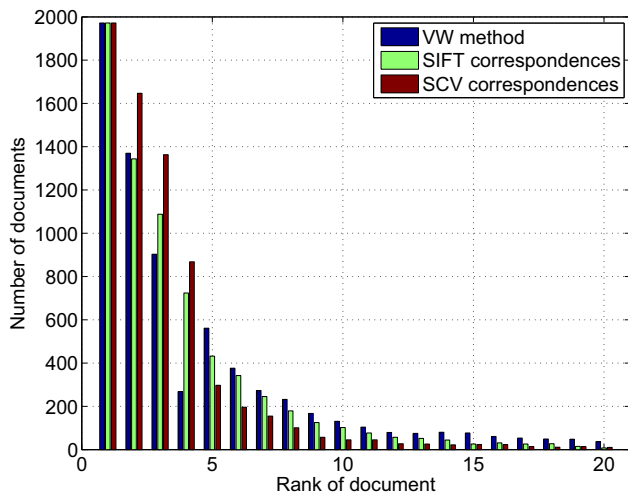
This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication.

14

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE



Fig. 17. Ranking of documents based on the visual word method, the number of SIFT correspondences with distance ratio $< 0.8$ and the number of SCV correspondences.

text retrieval. The similarity of two documents is then the $L_1$ distance between their vectors of the visual words weights. The top $K$ most similar documents are retrieved and an average score of 3.41 images retrieved per query is achieved.

To evaluate the performance of the SCV algorithm, we re-ranked retrieved images according to the number of SCV validated correspondences. Since for efficiency reasons the re-ranking could include only a small number of images (e.g. 20), we only considered queries that had at least one image of the retrieved object with rank 5 to 20. There are 1972 such query images. The overall score is 2.29 for these queries (note that these are difficult ones, the average on the full dataset is 3.41). The top 20 score, i.e. the average number of correct images among top 20, is 3.51 on this subset. This is the upper bound of the performance for a retrieval algorithm that resorts the top 20 retrieved images. The correspondences were verified by the SCV ($\alpha = 0.01, \beta = 0.001$) algorithm. Finally, new ranking was established according to the number of SCV correspondences found.

The performance of the SCV algorithm is compared in a histogram of ranks of the four correct images in answer to each query (see Fig. 17). Clearly, SCV significantly improves the ranking of the correct images bringing most of them to top 4. Its overall top 4 score on the selected query images is 5865 resulting in average 2.97, the average top 5 score is 3.12.

We also compared our method to the ranking based on SIFT correspondences (rank is based on the number of correspondences with SIFT distance ratio $< 0.8$). The overall top 4 score for SIFT correspondences is 5004 resulting in average 2.60 and the average top 5 score is 2.82.

Finally, we compared the achieved top 4 scores of both methods to the visual words method in Tab. 1. It shows the ranking is improved or unchanged with

TABLE 1

Comparison of the scores of the VW method and rankings based on SIFT and SCV correspondences.

| score | higher | same | lower |
|-------|--------|------|-------|
| SIFT  | 786    | 934  | 252   |
| SCV   | 1258   | 631  | 83    |

SCV in 95% of cases. Wrong ranking occurs typically for images of different objects with little texture (usually slightly blurred) on the same structured background. In this case, the most of, in fact correct, correspondences are found in the background which does not help retrieving a correct image.

## 5 CONCLUSIONS

We have presented the Sequential Correspondence Verification (SCV) algorithm which is able to efficiently distinguish correct and incorrect correspondences, via collecting statistics while cosegmenting gradually larger regions. We have shown this significantly benefits the matching process in challenging wide baseline scenes and improves results in a large scale image retrieval. The process is computationally efficient and very fast in practice, e.g. the method was successfully applied in a paper by Chum et al. [32] on large scale image retrieval algorithm with impressive results.

The SCV software is available for download at http://cmp.felk.cvut.cz/software/SCV.

## REFERENCES

[1]  D. Lowe, "Distinctive image features from scale-invariant key-points," *IJCV*, vol. 60, no. 2, pp. 91–110, 2004.
[2]  J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos." in *ICCV*, 2003, pp. 1470–1477.
[3]  B. Leibe, K. Mikolajczyk, and B. Schiele, "Efficient clustering and matching for object class recognition," in *BMVC*, 2006.
[4]  D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in *CVPR*, 2006, pp. 2161–2168.
[5]  J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *CVPR*, 2007.
[6]  O. Chum, J. Philbin, J. Sivic, M. Isard, and A. Zisserman, "Total recall: Automatic query expansion with a generative feature model for object retrieval," in *ICCV*, 2007.
[7]  T. Tuytelaars and L. V. Gool, "Wide baseline stereo matching based on local, affinely invariant regions," in *In BMVC*, 2000, pp. 412–425.
[8]  J. Matas, O. Chum, M. Urban, and T. Pajdla, "Robust wide-baseline stereo from maximally stable extremal regions," *Image and Vision Computing*, vol. 22, no. 10, pp. 761–767, 2004.
[9]  A. Vedaldi and S. Soatto, "Local features, all grown up," in *CVPR*, 2006, pp. 1753–1760.

[10] V. Ferrari, T. Tuytelaars, and L. van Gool, "Simultaneous object recognition and segmentation from single or multiple image views," *IJCV*, vol. 67, no. 2, pp. 159–188, 2006.

[11] G. Yang, C. V. Stewart, M. Sofka, and C.-L. Tsai, "Registration of challenging image pairs: Initialization, estimation, and decision," *IEEE Trans. on PAMI*, vol. 29, no. 11, pp. 1973–1989, 2007.

[12] C. V. Stewart, C.-L. Tsai, and B. Roysam, "The dual-bootstrap iterative closest point algorithm with application to retinal image registration," *Medical Imaging*, vol. 22, no. 11, pp. 1379–1394, 2003.

[13] M. Lhuillier and L. Quan, "Match propagation for image-based modeling and rendering," *PAMI*, vol. 24, no. 8, pp. 1140–1146, 2002.

[14] C. Rother, V. Kolmogorov, T. Minka, and A. Blake, "Cosegmentation of image pairs by histogram matching – incorporating a global constraint into MRFs," in *CVPR*, 2006.

[15] J. Cech, J. Matas, and M. Perdoch, "Efficient sequential correspondence selection by cosegmentation," in *CVPR*, 2008.

[16] J. Cech and R. Sara, "Efficient sampling of disparity space for fast and accurate matching," in *BenCOS Workshop, CVPR*, 2007.

[17] G. P. Otto and T. K. W. Chau, "'Region-growing' algorithm for matching of terrain images," *IVC*, vol. 7, no. 2, pp. 83–94, 1989.

[18] Z. Megyesi, G. Kos, and D. Chetverikov, "Dense 3D reconstruction from images by normal aided matching," *Machine Graphics and Vision*, vol. 15, pp. 3–28, 2006.

[19] S. Obdrzalek and J. Matas, "Sub-linear indexing for large scale object recognition," in *BMVC*, 2005.

[20] J. Matas, O. Chum, M. Urban, and T. Pajdla, "Robust wide baseline stereo from maximally stable extremal regions," in *BMVC*, 2002, pp. 384–393.

[21] F. Schaffalitzky and A. Zisserman, "Multi-view matching for unordered image sets, or how do i organize my holiday snaps?," in *ECCV*, 2002, pp. 414–431.

[22] H. P. Moravec, "Towards automatic visual obstacle avoidance," in *IJCAI*, 1977, p. 584.

[23] J. Kannala, E. Rahtu, S. S. Brandt, and J. Heikkila, "Object recognition and segmentation by non-rigid quasi-dense matching," in *CVPR*, 2007.

[24] V. Kolmogorov and R. Zabih, "Computing visual correspondence with occlusions using graph cuts," in *ICCV*, 2001, pp. 508–515.

[25] V. N. Vapnik, *The Nature of Statistical Learning Theory*. Springer Verlag, 1995.

[26] J. C. Platt, "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods," in *Advances in Large Margin Classifiers*. MIT Press, 1999, pp. 61–74.

[27] A. Wald, *Sequential analysis*. New York: Dover, 1947.

[28] V. Franc and S. Sonnenburg, "OCAS optimized cutting plane algorithm for support vector machines," in *ICML*, 2008.

[29] O. Chum and J. Matas, "Matching with PROSAC - progressive sample consensus," in *CVPR*, 2005, pp. 220–226.

[30] K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. Van Gool, "A comparison of affine region detectors," *IJCV*, vol. 65, no. 1/2, pp. 43–72, 2005.

[31] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE Trans. on PAMI*, vol. 27, no. 10, pp. 1615–1630, 2005.

[32] O. Chum, M. Perdoch, and J. Matas, "Geometric min-hashing: Finding a (thick) needle in a haystack," in *CVPR*, 2009.

**Jan Čech** received the MSc degree (with honors) in cybernetics from the Faculty of Electrical Engineering, Czech Technical University, Prague, Czech Republic, in 2002. He received the PhD degree in artificial intelligence and biocybernetics from the same institution, in 2009. His research interests include stereoscopic vision, 3D reconstruction, image understanding and interpretation. He serves as a reviewer for several conferences in the field. He is a member of the IEEE.

**Jiří Matas** received the MSc degree (with honors) in cybernetics from the Czech Technical University, Prague, Czech Republic, in 1987 and the PhD degree from the University of Surrey, Guildford, United Kingdom, in 1995. From 1991 to 1997, he was a research fellow in the Center for Vision, Speech, and Signal Processing at the University of Surrey. In 1997, he joined the Center for Machine Perception at the Czech Technical University. Since 1997, he has held various positions at these two institutions. He has published more than 100 papers in refereed journals and conference proceedings. His publications have more than 800 citations in the Science Citation Index. He received the best paper prize at the British Machine Vision Conferences in 2002 and 2005 and at the Asian Conference on Computer Vision in 2007. He has served in various roles at major international conferences (e.g., IEEE International Conference on Computer Vision (ICCV), IEEE Conference on Computer Vision and Pattern Recognition (CVPR), International Conference on Pattern Recognition (ICPR), and Neural Information Processing Symposium (NIPS)), cochairing the European Conference on Computer Vision (ECCV) 2004 and CVPR 2007. He is an associate editor-in-chief of the *IEEE Transactions on Pattern Analysis and Machine Intelligence* and a member of the editorial board of the *International Journal of Computer Vision*. He is a member of the IEEE.

**Michal Perďoch** received the BSc degree in software engineering from the Faculty of Electrical Engineering and Information Technology, Slovak University of Technology, Bratislava, Slovak Republic, in 2001 and the MSc degree in computer science from the Faculty of Electrical Engineering, Czech Technical University, Prague, Czech Republic, in 2004. He is currently a PhD student at the Center for Machine Perception. His current research interests include low level image processing, feature detection and description, object recognition and large scale object retrieval. He is a member of the IEEE.