



HAL
open science

Robust homography estimation from local affine maps

Mariano Rodríguez, Gabriele Facciolo, Jean-Michel Morel

► **To cite this version:**

Mariano Rodríguez, Gabriele Facciolo, Jean-Michel Morel. Robust homography estimation from local affine maps. 2022. hal-03579839

HAL Id: hal-03579839

<https://hal.science/hal-03579839>

Preprint submitted on 18 Feb 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Robust homography estimation from local affine maps

Mariano Rodríguez¹, Gabriele Facciolo¹, and Jean-Michel Morel¹

¹ Université Paris-Saclay, ENS Paris-Saclay, CNRS, Centre Borelli, F-94235, Cachan, France
({mariano.rodriguez, gabriele.facciolo, jean-michel.morel}@ens-paris-saclay.fr)

Abstract

The corresponding point coordinates determined by classic image matching approaches define local zero-order approximations of the global mapping between two images. But the patches around keypoints typically contain more information, which may be exploited to obtain a first-order approximation of the mapping, incorporating local affine maps between corresponding keypoints. Several methods have been proposed in the literature to compute this first-order approximation. In this paper we present several modifications of the RANSAC (RANDOM SAmple Consensus) algorithm [18], that uses affine approximations and *a-contrario* procedures to improve the homography estimation between a pair of images. The *a-contrario* methodology provides a definition of the soundness of an estimation and allows for adaptive thresholds of inlier/outlier discrimination. These approaches outperform the state-of-the-art for different choices of image descriptors and image datasets, and permit to increase the probability of success in identifying image pairs in challenging matching databases.

Keywords: homography, image comparison, image matching, robust estimation, RANSAC, affine invariance, scale invariance, local descriptors, affine normalization, SIFT, convolutional neural networks.

1 Introduction

Image matching consists in establishing correspondences between different images. This problem is recognized as difficult, especially under severe viewpoint changes between images. This is a fundamental step in many computer vision and image processing applications such as scene recognition [7, 13, 22, 45, 46, 65, 66, 73, 74, 79] and detection [19, 52], object tracking [81], robot localization [5, 49, 53, 67, 72], image stitching [2, 6], image registration [31, 78] and retrieval [21, 23], 3D modeling and reconstruction [1, 16, 20, 58, 75], motion estimation [76], photo management [9, 28, 68, 77], symmetry detection [33] or even image forgeries detection [10].

State-of-the-art image matching algorithms usually consist of three parts: *detector*, *descriptor* and *matching step*. They first detect points of interest in the compared images and select a region around each point of interest, and then associate an invariant descriptor or feature to each region. Correspondences may thus be established by matching the descriptors. Detectors and descriptors should be as invariant as possible.

Local image detectors can be classified by their incremental invariance properties. All of them are translation invariant. The Harris point detector [24] is also rotation invariant. The Harris-Laplace, Hessian-Laplace and the DoG (Difference-of-Gaussian) region detectors [17, 32, 36, 38] are invariant to rotations and changes of scale. Based on the AGAST [34] corner score, BRISK [29] performs a 3D nonmaxima suppression and a series of quadratic interpolations to extract the BRISK keypoints;

both detections aim at quickly providing rotation and scale invariances. Some moment-based region detectors [4,30] including the Harris-Affine and Hessian-Affine region detectors [37,38], an edge-based region detector [70,71], an intensity-based region detector [69,70], an entropy-based region detector [27], and two level line-based region detectors MSER (“maximally stable extremal region”) [35] and LLD (“level line descriptor”) [8,50,51] are designed to be invariant to affine transforms. MSER, in particular, has been demonstrated to have often better performance than other affine invariant detectors, followed by Hessian-Affine and Harris-Affine [39]. In his pivotal paper [32], Lowe proposed a scale-invariant feature transform (SIFT) that is invariant to image scaling and rotation and partially invariant to illumination and viewpoint changes. The SIFT method combines the DoG region detector that is rotation, translation and scale invariant (a mathematical proof of its scale invariance is given in [47]) with a descriptor based on the gradient orientation distribution in the region, which is partially illumination and viewpoint invariant [32]. These two stages of the SIFT method will be called respectively *SIFT detector* and *SIFT descriptor*. The SIFT detector is *a priori* less invariant to affine transforms than the Hessian-Affine and the Harris-Affine detectors [36,38].

The apparent deformations of objects caused by changes of the camera position can be locally approximated by affine maps, which explains why robust affine invariant methods allow to capture strong homography deformations. A possible way to obtain affine invariance is through the recently proposed Affnet [42]. Affnet is a Convolutional Neural Network (CNN) that was first conceived for improving the normalized representations of Hessian-Affine [38]. As proposed in [42], these normalized representations are then described and match by HardNet [40]. More recently the AID descriptor [62], a CNN-based patch descriptor trained to capture affine invariance, is able to cope directly with strong viewpoint deformations. Still, it seems that more classical Image Matching by Affine Simulation (IMAS) methods [41,48,54,60,63] provide the best affine invariances of them all [64]. Therefore, IMAS methods might be more suited for very strong viewpoint differences, although the price to pay is a heavier computational load.

First-order approximations of the local geometry, or simply affine approximations (see Figure 1), can be easily obtained from affine detectors like MSER [35], Harris-Affine [37] or Hessian-Affine [38]. Similarly, the SIFT detector can also be armed with local affine approximations [61]. When estimating homographies from sets of correspondences with the RANSAC algorithm [18], the use of first-order approximations allows to increase the performance in homography estimation. This has already been proposed in [57] by composing normalized affine maps provided by the Hessian Laplace detector. This information can be replaced with the one provided by Affnet [42] or LOCATE [61] since they have been shown to produce more accurate affine maps. In addition, a modification in the RANSAC consensus step has been proposed in [61], encouraging geometry consistency. Instead of defining inliers only through location agreement, the authors also consider the agreement in tilt, rotations and scale of the local affine maps.

A well established way of automatically estimating the 2D homography relating two images, see [25] p.123, is:

1. **Detection and description.** Compute and describe interest points in each image.
2. **Putative correspondences.** Compute a set of interest point matches based on similarity between their descriptors.
3. **RANSAC robust estimation.** Choose the homography η with the largest number of inliers.
4. **Optimal estimation.** Re-estimate η from all correspondences classified as inliers, by minimizing the Maximum likelihood cost function, see [25] p.95.
5. **Guided matching.** Further interest point correspondences can be determined using the estimated η or the affine approximation around each match.

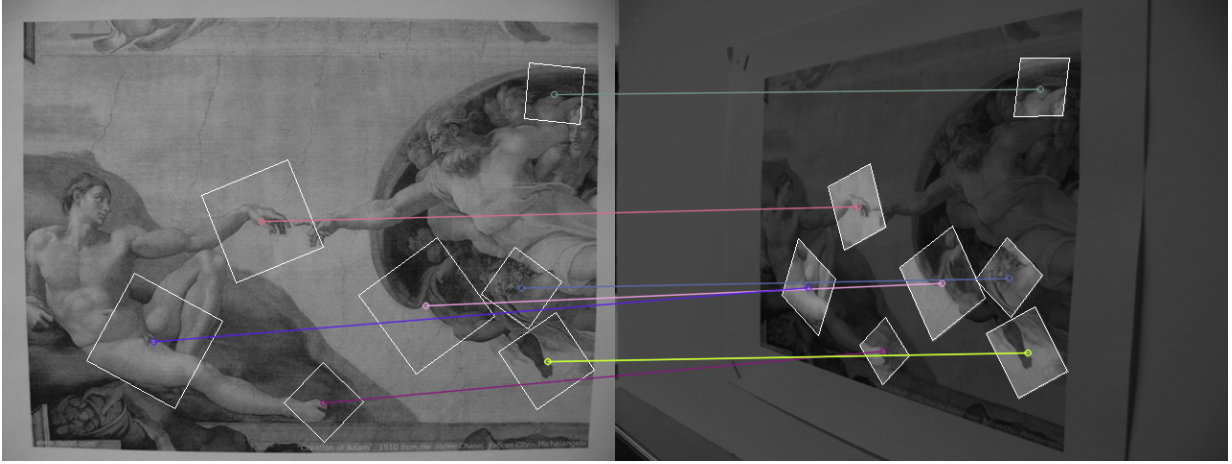


Figure 1: Some correspondences together with local affine approximations of local geometry. Patches on the target are warped versions of their corresponding query patch.

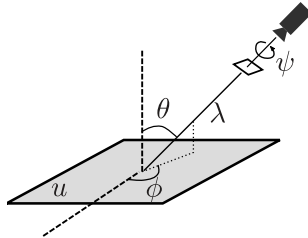


Figure 2: Geometric interpretation of Equation (1).

In this paper, we focus on some available improvements for the step 3. In particular, we highlight the advantages of three local affine approximations in improving the robustness of RANSAC homography estimations. They also enable RANSAC to reduce the amount of iterations by generating homography candidates with only two matching pairs instead of four. Furthermore, incorporating a *a-contrario* methodology [11] will, as in ORSA [43], result in: first, a threshold for inlier/outlier discrimination that is adaptive; second, a measure of the soundness of an estimation. In order to validate and measure the impact of these improvements, all modifications to RANSAC will be tested under several choices of detectors, descriptors and matchers which determine the first two steps).

The last two steps (optimal estimation and guided matching) could be seen as optional, and can be iterated until the number of correspondences is stable. Note that the local affine information could also be exploited in these steps. In the case of guided matching, the affine information have already been used to predict location and camera parameters from affine invariant detectors [14, 15] and also from local geometry estimators [61].

The rest of this paper is organized as follows. Section 2 summarizes a formal methodology for approximating locally the viewpoint changes induced by motion of real cameras. Three methods for computing those local approximations are introduced in Section 3. Section 4 presents the modified RANSAC model. The proposed methods are illustrated with experiments in Section 5. Concluding remarks are presented in Section 6.

2 Affine maps and homographies

As stated in [48, 59], a digital image \mathbf{u} obtained by any camera at infinity is modeled as $\mathbf{u} = \mathbf{S}_1 \mathbb{G}_1 A u$, where \mathbf{S}_1 is the image sampling operator (on a unitary grid), A is an affine map, u is a continuous image and \mathbb{G}_δ denotes the convolution by a Gaussian kernel broad enough to ensure no aliasing is induced by the δ -sampling. This model takes into account the blur incurred when tilting or zooming a view. Note that \mathbb{G}_1 and A generally do not commute.

Let \mathcal{A} denote the set of affine maps and define $Au(x) = u(Ax)$ for $A \in \mathcal{A}$, where x is a 2D vector and Ax denotes function evaluation, $A(x)$. We define the set of invertible orientation preserving affinities

$$\mathcal{A}^+ = \{L + v \in \mathcal{A} \mid \det(L) > 0\}$$

where L is a linear map and v a translation vector. We call \mathcal{S} the set of similarity transformations, which are any combination of translations, rotations and zooms. Lastly, we define the set

$$\mathcal{A}_*^+ = \mathcal{A}^+ \setminus \mathcal{S},$$

where we exclude pure similarities. As it was pointed out in [48], every $A \in \mathcal{A}_*^+$ is uniquely decomposed as

$$A = \lambda R_1(\psi) T_t R_2(\phi), \quad (1)$$

where R_1, R_2 are rotations, $T_t = \begin{bmatrix} t & 0 \\ 0 & 1 \end{bmatrix}$ with $t > 1$, $\lambda > 0$, $\phi \in [0, \pi)$ and $\psi \in [0, 2\pi)$. Furthermore, the above decomposition comes with a geometric interpretation (see Figure 2) where the longitude ϕ and latitude $\theta = \arccos \frac{1}{t}$ characterize the camera's viewpoint angles (or tilt), ψ parameterizes the camera roll and λ corresponds to the camera zoom. The so-called optical affine maps involving a tilt t in the z -direction and zoom λ are formally simulated by:

$$\mathbf{u} \mapsto \mathbf{S}_1 A \mathbb{G}_{\sqrt{t^2-1}}^z \mathbb{G}_{\sqrt{\lambda^2-1}} I \mathbf{u},$$

where I is the Shannon-Whittaker interpolator and the superscript z indicates that the operator takes place only in the z -direction. We denote the corresponding operator by

$$\mathbb{A} := \mathbf{S}_1 A \mathbb{G}_{\sqrt{t^2-1}}^z \mathbb{G}_{\sqrt{\lambda^2-1}} I.$$

The operator \mathbb{A} is not always invertible and therefore its application might incur into a loss of information. We refer to [62] for an example where no optical transformation \mathbb{A} is found between two views.

2.1 Local affine approximation of homographies

Let $H = (h_{ij})_{i,j=1,\dots,3}$ be the 3×3 matrix associated to the homography $\eta(\cdot)$. Let \mathbf{x} be the homogeneous coordinates vector associated to the image point $x = (x_1, x_2)$ around which we want to determine the local affine map. We denote by

$$y = (y_1, y_2) = \left(\frac{(H\mathbf{x})_1}{(H\mathbf{x})_3}, \frac{(H\mathbf{x})_2}{(H\mathbf{x})_3} \right) = \eta(x)$$

the image of x by the homography η .

The first order Taylor approximation of η at x leads to

$$\eta(x + z) = v + L(x + z) + o(\|z\|). \quad (2)$$

More specifically, if $x = (0, 0)$, we know that

$$y_i(z_1, z_2) = (h_{i1}z_1 + h_{i2}z_2 + h_{i3}) \left(\frac{1}{h_{33}} - \frac{h_{31}}{h_{33}^2}z_1 + -\frac{h_{32}}{h_{33}^2}z_2 + o(\|z\|) \right), \quad i = 1, 2.$$

Then, by polynomial identification in the Taylor formula

$$v + L(z) = \frac{1}{h_{33}} \begin{pmatrix} h_{13} \\ h_{23} \end{pmatrix} + \left[\frac{1}{h_{33}} \begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix} - \frac{1}{h_{33}^2} \begin{pmatrix} h_{13} \\ h_{23} \end{pmatrix} \begin{pmatrix} h_{31} & h_{32} \end{pmatrix} \right] \begin{pmatrix} z_1 \\ z_2 \end{pmatrix},$$

where

$$\frac{1}{h_{33}} \begin{pmatrix} h_{13} \\ h_{23} \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}.$$

If $x \neq (0, 0)$, a simple change of variables $z \rightarrow z + x$ would lead us back to the case $x = (0, 0)$. Notice that the resulting homography,

$$\tilde{\eta}(z) = \eta(z + x),$$

has an associated matrix determined by columns,

$$H_{\tilde{\eta}} = \left[H \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad H \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad H \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} \right].$$

This brief computation shows that the vector v and the matrix L are given by

$$L = \begin{bmatrix} \frac{h_{11}-y_1h_{31}}{h_{31}x_1+h_{32}x_2+h_{33}} & \frac{h_{12}-y_1h_{32}}{h_{31}x_1+h_{32}x_2+h_{33}} \\ \frac{h_{21}-y_2h_{31}}{h_{31}x_1+h_{32}x_2+h_{33}} & \frac{h_{22}-y_2h_{32}}{h_{31}x_1+h_{32}x_2+h_{33}} \end{bmatrix}, \quad (3)$$

$$v = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - Lx. \quad (4)$$

This derivation allows us to compute the exact local affine approximation for a given homography. This will be useful for Section 4.2-4.3 and to assess the accuracy of our method when using annotated datasets.

3 Computing local affine approximations

A SIFT-like patch is simply the square crop at the origin of some similarity transformation (translation, rotation and zoom) of the original image. This also stands true for affine invariant image matching methods, for which a patch can be considered as the square crop at the origin of a tilted version of the original image followed by some similarity transformation. If two patches are a match (i.e. their descriptors are similar), then, through the assumption of locality and Taylor's formula, we can approximate the geometry transformation by an affine map.

Consider two square patches, P_q and P_t , coming, for example, from the Gaussian pyramid of the query and target images, respectively. Let c_q and c_t be their centers expressed in image coordinates. Let also A_q be the affine map that converts from the query image domain to patch coordinates; likewise A_t converts from target to patch coordinates. Note that, in the case of SIFT-like patches, the affinities A_q and A_t are pure similarities, combining just the translation, rotation and zoom corresponding to the location, orientation and scale associated to SIFT-like keypoints. Lastly, in order to locally approximate the transformation between query and target images (centered at c_q and c_t), we only need the affine map relating P_q and P_t , denoted by A . Figure 3 illustrates the affine

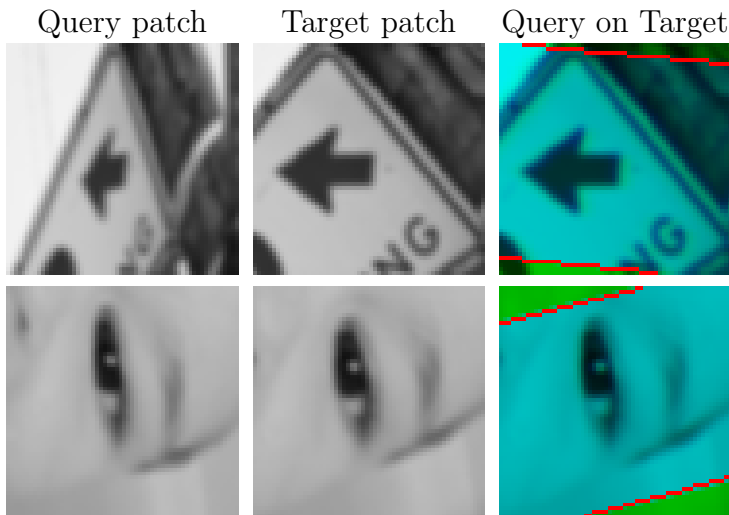


Figure 3: Two pairs of patches used as query and target input patches (columns 1-2). Each pair differs only by an affine map. Blue and green channels in the 3rd column correspond to the target patch and a warped version of the query patch (the red line delimits its borders); the red channel is filled with zeros.

differences between two patches and how the query is transformed into the target. All in all, around c_q , the local affine map transforming the query into the target (in image coordinates) is written as

$$A_{q \rightarrow t} = A_t^{-1} A A_q. \quad (5)$$

The same procedure described above can be derived for Harris-Affine and Hessian-Affine region detectors [37, 38]. Only that A_q, A_t are not restrained to similarities but to affine maps in general.

3.1 Affine connections between patches

We now list different choices of methods for computing A , the affine map locally connecting the query patch to the target patch.

3.1.1 The naive method

If for some reason no strong tilt deformations are expected between query and target patches, then, a fair assumption would be for A to be the identity in all cases. Unfortunately, this is rarely the case for real life images. Nevertheless, it is a simplification worth trying because it entails no additional computational complexity. Formally, we set

$$A := Id. \quad (6)$$

3.1.2 The Affnet method

The Affnet [42] method was conceived to predict normalizing ellipse shapes for single patches based on a 3-variable parametrization. Figure 4 depicts the passage from Affnet affine maps to the affine map transforming query into target (i.e. A). The connection provided by two Affnet-normalizing affine maps for the query and target patches is richer than each normalizing transformation. Indeed, for different choices of $A_1 = T_1 R_1$ and $A_2 = T_2 R_2$ one would need the four parameters (zoom, camera

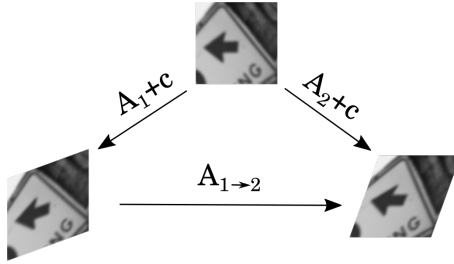


Figure 4: Passage from Affnet affine maps (A_1, A_2) to the connecting mapping $A_{1 \rightarrow 2}$. The center of the normalized patch (on top) corresponds to the origin in normalized coordinates.

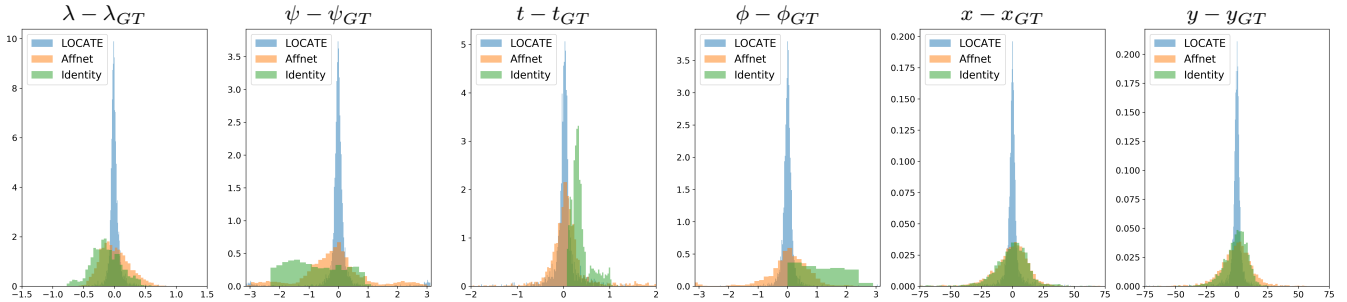


Figure 5: Affine error prediction in terms of the affine decomposition of Equation 1 (namely zoom λ , camera rotation ψ , tilt t , tilt direction ϕ , and translation x, y), for the LOCATE method, the Affnet method [42] and the identity map method. The [62] dataset is used; it contains 3352 patch pairs with corresponding ground truth. The sub-index GT means ground truth, conversely, no sub-index stands for estimated parameters.

rotation, tilt and tilt direction) in Equation 1 in order to express $A_2 A_1^{-1}$. However, Affnet does not estimate translations. Formally, we set

$$A := A_2 (A_1^{-1} \mathbf{x} - A_1^{-1} \mathbf{c}) + \mathbf{c}, \quad (7)$$

where \mathbf{c} denotes the center of patch domain and A_i are the estimated affine maps by Affnet.

3.1.3 The LOCATE method

The LOCAL Affine Transform Estimator (LOCATE) network presented in [61] directly estimates the affine transform A_{LOCATE} that maps the query patch into the target patch. LOCATE simultaneously tracks the direct and inverse maps which significantly improves the network performance in predicting local approximating affine maps. This network was trained exclusively with simulated patches from an affine camera model. We expect it to generalize the affine world to all sorts of geometry as long as the Taylor approximation holds. LOCATE is able to estimate all six parameters composing the affine map. Formally, we set

$$A := A_{\text{LOCATE}}. \quad (8)$$

3.2 Precision

In order to measure the precision in a realistic environment of these three methods (Naive, Affnet and LOCATE) we used the viewpoint dataset presented in [62], consisting of five pairs of images with their ground truth homographies and 3352 true matches. Notice that Equations 3-4 allow us to compute ground truth local affine maps around each match. Figure 5 shows the accuracy of Naive

(Identity), Affnet [42] and LOCATE, represented by error density functions with respect to the affine decomposition appearing in Equation 1. Ideally, we expect a Dirac delta function centered at 0 for a perfect method. This is approximately true for the LOCATE method. Note in Figure 5 that translations from the Affnet [42] method do not quite match those from the Identity method; this difference can be explained by the connecting mapping itself (see Equation 7) which is different from $A_2A_1^{-1}\mathbf{x}$. As expected, LOCATE is more precise than Affnet [42]. Indeed, Affnet analyzes one patch at a time, whereas LOCATE has access to both patches simultaneously. However, in practice, using Affnet involves fewer computations. This trade-off must be resolved depending on the application.

4 Robust homography estimation

The standard RANSAC algorithm computes the parameters fitting a mathematical model from observed data in the presence of outliers. Numerous improvements have been proposed in the literature for RANSAC, see [43, 44, 55, 56], but the core idea remains the same.

In the case of homography estimation, the classic RANSAC algorithm returns the homography η_j computed at iteration j having the largest consensus of inliers among all iterations. The j -iteration can be described in two steps:

1. (Fitting) Randomly select s matches $(x^i \leftrightarrow y^i)_{i=1,\dots,s}$ from the set of all matches (M_T) and compute the homography η_j that yields the best fit.
2. (Consensus) Determine, with respect to an error function ξ_{η_j} , the matches from M_T in consensus with η_j .

In this paper, two error functions ξ_η are used. First, the symmetric transfer error,

$$\xi_\eta^4(x \leftrightarrow y) := \left\| \begin{pmatrix} \eta_j(x) - y \\ x - \eta_j^{-1}(y) \end{pmatrix}_{4 \times 1} \right\|_{l_2}. \quad (9)$$

Note that the 4 dimensional vector inside the norm, corresponding to the concatenation of two vectors of dimension two. The second proposal is a function measuring the classic symmetric transfer error as well as the affine coherence,

$$\xi_\eta^8(x \leftrightarrow y) := \left\| \begin{pmatrix} \eta_j(x) - y \\ x - \eta_j^{-1}(y) \\ \alpha \left(A_E^{(x \leftrightarrow y)}, A_H^{(x \leftrightarrow y)} \right) - \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix} \end{pmatrix}_{8 \times 1} \right\|_{l_2}, \quad (10)$$

where the α -vector was introduced in [61] and will be revisited here in the following sections. The vector inside the norm has eight dimensions, corresponding to the concatenation of two vectors of dimension two and a vector of dimension four.

A classic way of determining inliers in step 2 of the RANSAC algorithm is by thresholding the function ξ_η . Another way is through the Number of False Alarms (NFA), based on the *a-contrario* procedure [11]. Let $\varepsilon_i = \xi_\eta(m_{i_i})$ be the ordered errors of matches in M_T with respect to the current testing homography η , i.e.,

$$\varepsilon_1 \leq \varepsilon_2 \leq \dots \leq \varepsilon_{|M_T|}.$$

As explained in [43], the NFA of the testing homography η and its k smallest error matches from M_T is defined as,

$$NFA(\eta) = (|M_T| - s) \binom{|M_T|}{k} \binom{k}{s} \mathbb{P}(\xi_\eta(\mathbf{m}) \leq \varepsilon_k)^{(k-s)} \quad (11)$$

where $\mathbb{P}(\xi_\eta(\mathbf{m}) \leq \varepsilon_k)^{(k-s)}$ corresponds to the probability of k matches to fall within an error threshold ε_k with respect to ξ_η . This probability expression is explained through independence among matches and the power $k - s$ accounts for the fact that s matches (the ones used for estimating η) are automatically inliers with respect to η . If there are k inliers, potentially all s out of them yield the correct configuration; explaining the term $\binom{k}{s}$. Also, there are $\binom{|M_T|}{k}$ possible subsets of M_T with k elements. The number k of inliers is usually not known in advance, so all values of k are tested (from $s + 1$ to $|M_T|$), which explains the factor $(|M_T| - s)$ in Equation 11. In practice, the errors (measured by ξ_η) of all data terms are collected and sorted: ε_1 to $\varepsilon_{|M_T|}$. For each possible k , we compute the NFA as in Equation 11, and keep only the minimum of them all, provided it is below some threshold, usually set to 1 in the a-contrario methodology.

4.1 The benchmark RANSAC

Usually the steps 1-2 only take into account point coordinates. If no further improvements are applied, this defines a base RANSAC, that we denote by $RANSAC_{base}$. Since the homography matrix have eight degrees of freedom and each match defines two equations, then the number of matches must be at least $s = 4$.

The a-contrario $RANSAC_{base}$ is equivalent to the ORSA Homography estimation method presented in [43]. The probability term in Equation 11 does not need to be exact for the NFA to work, and can be approximated by assuming a uniform distribution in \mathbb{R}^4 ,

$$\mathbb{P}(\xi_\eta^4(\mathbf{m}) \leq \varepsilon) \approx \frac{\varepsilon^4 \frac{\pi^2}{2}}{w_q h_q w_t h_t}, \quad (12)$$

where we find in the numerator the volume of a sphere of radius ε in \mathbb{R}^4 and in the denominator the volume of the set with all possible coordinates for the query (of size $[w_q, h_q]$) and target (of size $[w_t, h_t]$) images.

Algorithm 1 details the operations taking place in $RANSAC_{base}$.

4.2 Homography fitting from local affine maps in RANSAC

From Section 2.1 we know how to locally approximate a homography by an affine map. Conversely, the problem of determining a homography from a set of affine maps at different locations was addressed in [3, 57]. Let $x \leftrightarrow y$ be a match and $L = (l_{ij})_{i,j=1,2}$ the linear map in Equation 2. Then the unknown homography η must satisfy

$$E_{6 \times 9} \cdot \vec{h} = \vec{0}, \quad (13)$$

where $E_{6 \times 9}$ is the matrix

$$\begin{bmatrix} 1 & & & -y_1 - l_{11}x_1 & -l_{11}x_2 & -l_{11} & & & \\ & 1 & & -l_{12}x_1 & -y_1 - l_{12}x_2 & -l_{12} & & & \\ & & 1 & -y_2 - l_{21}x_1 & -l_{21}x_2 & -l_{21} & & & \\ & & & 1 & -l_{22}x_1 & -y_2 - l_{22}x_2 & -l_{22} & & \\ x_1 & x_2 & 1 & & -y_1x_1 & -y_1x_2 & -y_1 & & \\ & & & x_1 & x_2 & 1 & & -y_2x_1 & -y_2x_2 & -y_2 \end{bmatrix}, \quad (14)$$

and $\vec{h} = [h_{11}, h_{12}, h_{13}, h_{21}, h_{22}, h_{23}, h_{31}, h_{32}, h_{33}]^T$ is a vectorized version of the matrix H associated to η . The first four rows of $E_{6 \times 9}$ are determined by Equation (3) and the last two are the classic equations derived from rewriting $\eta(x) = y$ in terms of $H\mathbf{x} = \mathbf{y}$.

Algorithm 1 RANSAC_{base}

input:
 M_T - set of all matches.

parameters:
 N_{iters} - Number of iterations.

 κ - Spatial inliers threshold.

useNFA - A bool stating if NFA measure should be used.

start:
foreach $j \in \{1, \dots, N_{iters}\}$ **do**

 Randomly select $m_1 = (x^1 \leftrightarrow y^1), \dots, m_4 = (x^4 \leftrightarrow y^4)$ from M_T

 // Homography fitting by the Direct Linear Transformation (DLT) algorithm \triangleright see [25] p.88

$$E_i = \begin{bmatrix} x_1^i & x_2^i & 1 & & & & -y_1^i x_1^i & -y_1^i x_2^i & -y_1^i \\ & & & x_1^i & x_2^i & 1 & -y_2^i x_1^i & -y_2^i x_2^i & -y_2^i \end{bmatrix}_{2 \times 9}, \quad \triangleright \text{where } z^i = (z_1^i, z_2^i)$$

$$E = \begin{bmatrix} E_1 \\ \vdots \\ E_4 \end{bmatrix}_{8 \times 9}$$

 \vec{h}_j is the unit singular vector corresponding to the smallest singular value of E . \triangleright where \vec{h}_j is a vectorized version of the matrix H_j associated to the homography η_j .

// Selection of inliers with respect to the symmetric transfer error

$$I_j = \{(x \leftrightarrow y) \in M_T \mid \xi_\eta^4((x \leftrightarrow y)) < \kappa\}$$

if useNFA **then**

$$\varepsilon_{j,\cdot} = \text{SORT} \left(\left\{ \xi_\eta^4((x \leftrightarrow y)) \right\}_{(x \leftrightarrow y) \in I_j} \right)$$

foreach $k \in \{3, \dots, |I_j|\}$ **do**

$$\left[\left[NFA_k^j = (|M_T| - 4) \binom{|M_T|}{k} \binom{k}{4} \left(\frac{\varepsilon_{j,k}^4 \frac{\pi^2}{2}}{w_q h_q w_t h_t} \right)^{k-4} \right] \right]$$

if useNFA **then**

$$j^*, k^* = \arg \min_{j,k} NFA_k^j$$

if $NFA_{k^*}^{j^*} < 1$ **then**
 $\lfloor M_I$ is the subset of I^{j^*} achieving the first k^* smaller values in ε_{j^*} .

else
 $\lfloor M_I = \emptyset$
else

$$j^* = \arg \max_j |I_j|$$

 $\lfloor M_I = I_{j^*}$ **if** $|I_{j^*}| > 2$ **else** \emptyset
return M_{j^*}, η_{j^*}

Clearly, two matches with their corresponding local affine maps already over-determine the homography matrix. Indeed, putting those equations together provides us with 12 equations

$$\begin{bmatrix} E_1 \\ E_2 \end{bmatrix}_{12 \times 9} \cdot \vec{h} = \vec{0},$$

where E_i denotes the matrix E appearing in Equation 13 for each match. To avoid the solution $\vec{h} = \vec{0}$ we look for a unitary vector \vec{h} minimizing

$$\left\| \begin{bmatrix} E_1 \\ E_2 \end{bmatrix} \cdot \vec{h} \right\|,$$

see [25] (algorithm 4.1, p.91) for computational details.

We call $RANSAC_{2pts}$ a $RANSAC_{base}$ version in which the classic homography fitting of step 1 is replaced by the homography fitting of this section. Note that $RANSAC_{2pts}$ only needs two samples at each iteration ($s = 2$). These samples carry additional information consisting in affine approximations of local geometry deformations around them. The affine approximations can be provided by one of the methods introduced in Section 3.

The a-contrario $RANSAC_{2pts}$ differs little from the a-contrario $RANSAC_{base}$. Indeed, the only difference being $s = 2$ in Equation 11. Algorithm 2 details all operations taking place in $RANSAC_{2pts}$.

4.3 Affine consensus for RANSAC homography

When matching two image patches, the transformation that relates them may not be consistent with the global transformation of the scene. This can be due to the presence of symmetric objects or even to failures in the matching process. For instance, suppose that two patches centered at the same scene location but with incoherent rotations are identified by a matching method. The symmetry issue is easy to address as usually we should have encountered as many keypoints as degrees of symmetry around the center; so at least two rotations will correspond. However, aberrant matches are not treated by the matching method nor by RANSAC. This problem can be circumvented by imposing consistency between the local approximating affine maps and the testing homography proposed by RANSAC.

To impose local geometry consistency, most existing works [41,80] propose to measure the incurred error in mapping keypoints of a match $x \leftrightarrow y$, e.g. $\|y - A(x)\| + \|x - A^{-1}(y)\|$. Unlike them, [61] proposes to enforce geometry consistency directly on the transformations parameters given by Equation 1. Finally, we redefine the consensus set of the RANSAC model by imposing geometry consistency as in [61].

Inliers are now defined as follows. Let A_E and A_H be, respectively, the estimated affine map by one of the methods introduced in Section 3 and the testing affine map computed from the testing homography (using Equation (3)). Let also $[\lambda_E, \psi_E, t_E, \phi_E]$ and $[\lambda_H, \psi_H, t_H, \phi_H]$ be, respectively, the affine parameters of A_E and A_H . We define the α -vector between A_E and A_H as:

$$\alpha(A_E, A_H) = \left[\max\left(\frac{\lambda_E}{\lambda_H}, \frac{\lambda_H}{\lambda_E}\right), \angle(\psi_E, \psi_H), \max\left(\frac{t_E}{t_H}, \frac{t_H}{t_E}\right), \angle(\phi_E, \phi_H) \right], \quad (15)$$

where $\angle(\cdot, \cdot)$ denotes the angular difference. To test consistency between A_E and A_H we add to the classic threshold on the Euclidean distance, four more thresholds on the α -vector. A perfect match would result in an α -vector equal to $[1, 0, 1, 0]$. If we assume independence on each dimension, the resulting probability of a match passing all thresholds is the multiplication of individual probabilities. With this in mind, we claim that rough thresholds are enough to obtain good performances and that

Algorithm 2 RANSAC_{2pts}

input:
 M_T - set of all matches.

 A_T - Local affine approximations associated to each match in M_T .

parameters:
 N_{iters} - Number of iterations.

 κ - Spatial inliers threshold.

useNFA - A bool stating if NFA measure should be used.

start:
foreach $j \in \{1, \dots, N_{\text{iters}}\}$ **do**

 Randomly select $m_1 = (x^1 \leftrightarrow y^1)$ and $m_2 = (x^2 \leftrightarrow y^2)$ from M_T . Let also $L_1 + v_1, L_2 + v_2 \in A_T$ be their associated affine maps.

 // Homography fitting by the Direct Linear Transformation (DLT) algorithm ▷ see [25] p.88

$$E_i = \begin{bmatrix} 1 & & & -y_1^i - l_{11}^i x_1^i & -l_{11}^i x_2^i & -l_{11}^i \\ & 1 & & -l_{12}^i x_1^i & -y_1^i - l_{12}^i x_2^i & -l_{12}^i \\ & & 1 & -y_2^i - l_{21}^i x_1^i & -l_{21}^i x_2^i & -l_{21}^i \\ & & & 1 & -l_{22}^i x_1^i & -l_{22}^i \\ x_1^i & x_2^i & 1 & -y_1^i x_1^i & -y_1^i x_2^i & -y_1^i \\ & & x_1^i & x_2^i & 1 & -y_2^i x_1^i & -y_2^i x_2^i & -y_2^i \end{bmatrix}_{6 \times 9}, \quad \triangleright \quad \begin{matrix} \text{where} \\ z^i = (z_1^i, z_2^i), \\ L_i = (l_{lk}^i)_{l,k=1,2}. \end{matrix}$$

$$E = \begin{bmatrix} E_1 \\ E_2 \end{bmatrix}_{12 \times 9}$$

 \vec{h}_j is the unit singular vector corresponding to the smallest singular value of E . ▷ where \vec{h}_j is a vectorized version of the matrix H_j associated to the homography η_j .

// Selection of inliers with respect to the symmetric transfer error

$$I_j = \{(x \leftrightarrow y) \in M_T \mid \xi_\eta^4((x \leftrightarrow y)) < \kappa\}$$

if useNFA **then**

$$\begin{aligned} & \varepsilon_{j,\cdot} = \text{SORT} \left(\left\{ \xi_\eta^4((x \leftrightarrow y)) \right\}_{(x \leftrightarrow y) \in I_j} \right) \\ & \text{foreach } k \in \{3, \dots, |I_j|\} \text{ do} \\ & \quad \left[\text{NFA}_k^j = (|M_T| - 2) \binom{|M_T|}{k} \binom{k}{2} \left(\frac{\varepsilon_{j,k}^4 \frac{\pi^2}{2}}{w_q h_q w_t h_t} \right)^{k-2} \right] \end{aligned}$$

if useNFA **then**

$$j^*, k^* = \arg \min_{j,k} \text{NFA}_k^j$$

if $\text{NFA}_{k^*}^{j^*} < 1$ **then**
 $\lfloor M_I$ is the subset of I^{j^*} achieving the first k^* smaller values in $\varepsilon_{j^*,\cdot}$.

else
 $\lfloor M_I = \emptyset$
else

$$j^* = \arg \max_j |I_j|$$

 $\lfloor M_I = I_{j^*}$ **if** $|I_{j^*}| > 2$ **else** \emptyset
return M_{j^*}, η_{j^*}

there is no need to optimize them. Thus, we propose to further refine inliers by accepting only those matches also satisfying

$$\alpha(A_E, A_H) < \left[2, \frac{\pi}{4}, 2, \frac{\pi}{8}\right], \quad (16)$$

where the above logical operation is true if and only if it holds true for each dimension.

We call $RANSAC_{affine}$ the version of $RANSAC_{2pts}$ that includes the affine consensus presented in this section. The NFA (Equation 11) for the a-contrario $RANSAC_{affine}$ is determined by: two samples to fit the testing homography, $s = 2$; and a probability term that is approximated by a uniform random variable in a hyperrectangle of \mathbb{R}^8 ,

$$\mathbb{P}(\xi_\eta^8(\mathbf{m}) \leq \varepsilon) \approx \frac{\varepsilon^8 \frac{\pi^4}{4!}}{w_q h_q w_t h_t 12^2 \pi^2}, \quad (17)$$

where on top we find the volume of a sphere of radius ε in \mathbb{R}^8 and at the bottom the volume of the hyperrectangle $[0, w_q - 1] \times [0, h_q - 1] \times [0, w_t - 1] \times [0, h_t - 1] \times [0, 11]^2 \times [0, \pi]^2$.

Algorithm 3 details all operations taking place in $RANSAC_{affine}$.

5 Experiments

In order to quantify the benefits of the local affine approximations for the homography estimation problem, we will compare $RANSAC_{2pts}$ and $RANSAC_{affine}$ to $RANSAC_{base}$. The three variants of RANSAC presented in this work are based on the original RANSAC algorithm and do not include recent modifications proposed in the literature (e.g. RANSAC USAC [55], etc). Even if they are not entirely comparable to ORSA [43] nor RANSAC USAC [55], a line on each experiment will be added for ORSA or USAC as benchmark to compare against the state-of-the-art. The reader should keep in mind that most improvements proposed in RANSAC USAC [55] can also be applied to $RANSAC_{2pts}$ and $RANSAC_{affine}$.

All experiments in this section were conducted on four well known datasets for homography estimation. Those datasets are: EF [82], EVD [41], OxAff [39] and SymB [26]. The EF dataset presents challenging non-linear lighting variations and occlusions. Both EVD and OxAff datasets present several pairs incurring in strong viewpoint differences. In particular, viewpoint differences between pairs in the EVD dataset are extreme and most matching methods will struggle to find correct matches. Lastly, the SymB dataset consists of painting-to-photo pairs, which will be challenging for detectors, descriptors and local geometry estimators that were not intended to be used under this circumstances. All datasets include ground truth homographies that were used to verify the accuracy.

Four combinations of state-of-the-art detectors and descriptors were used as starting point for all RANSACs. These choices are: SIFT [32] + AID [62]; SIFT [32] + HardNet [40]; HessianAffine [38] + AID [62]; HessianAffine [38] + HardNet [40]. Once local features were detected and matched, then each homography estimation method was applied and we declared a success if at least 80% of inliers (in consensus with the estimated homography) were in consensus with the ground truth homography. Four metrics are reported: the number of successes; the number of correctly matched image pairs; the average number of correct inliers; and the average pixel error. Notice that all these metric indicators are computed by first thresholding the symmetric transfer error ($\xi_{\eta_{gt}}^4(\cdot) \leq \kappa$) with respect to the ground truth homography η_{gt} of all matches in consensus with the estimated homography. These last two elements are nothing more than the output: M_I, η_{j^*} . The two steps of RANSAC (fitting and consensus) are iterated a 1000 times for each of the RANSAC variants, except for ORSA that is iterated a 10000 times and USAC that adapts the number of iterations at each execution.

Algorithm 3 RANSAC_{affine}

input:
 M_T - set of all matches.

 A_T - Local affine approximations associated to each match in M_T .

parameters:
 N_{iters} - Number of iterations.

 κ - Spatial inliers threshold.

 $\vec{\alpha}_{\text{max}}$ - Affine inliers thresholds (a 4-sized vector).

useNFA - A bool stating if NFA measure should be used.

start:
foreach $j \in \{1, \dots, N_{\text{iters}}\}$ **do**

 Randomly select $m_1 = (x^1 \leftrightarrow y^1)$ and $m_2 = (x^2 \leftrightarrow y^2)$ from M_T . Let also $L_1 + v_1, L_2 + v_2 \in A_T$ be their associated affine maps.

 // Homography fitting by the Direct Linear Transformation (DLT) algorithm ▷ see [25] p.88

$$E_i = \begin{bmatrix} 1 & & & -y_1^i - l_{11}^i x_1^i & -l_{11}^i x_2^i & -l_{11}^i \\ & 1 & & -l_{12}^i x_1^i & -y_1^i - l_{12}^i x_2^i & -l_{12}^i \\ & & 1 & -y_2^i - l_{21}^i x_1^i & -l_{21}^i x_2^i & -l_{21}^i \\ & & & -l_{22}^i x_1^i & -y_2^i - l_{22}^i x_2^i & -l_{22}^i \\ x_1^i & x_2^i & 1 & -y_1^i x_1^i & -y_1^i x_2^i & -y_1^i \\ & & & x_1^i & x_2^i & 1 \\ & & & -y_2^i x_1^i & -y_2^i x_2^i & -y_2^i \end{bmatrix}_{6 \times 9}, \quad \triangleright \quad \begin{array}{l} \text{where} \\ z^i = (z_1^i, z_2^i), \\ L_i = (l_{lk}^i)_{l,k=1,2}. \end{array}$$

$$E = \begin{bmatrix} E_1 \\ E_2 \end{bmatrix}_{12 \times 9}$$

 \vec{h}_j is the unit singular vector corresponding to the smallest singular value of E . ▷ where \vec{h}_j is a vectorized version of the matrix H_j associated to the homography η_j .

// Selection of inliers with respect to the symmetric transfer error and the affine information

$$S_j = \{(x \leftrightarrow y) \in M_T \mid \xi_\eta^4((x \leftrightarrow y)) < \kappa\}$$

$$T_j = \{(x \leftrightarrow y) \in M_T \mid \alpha(A_E^{(x \leftrightarrow y)}, A_H^{(x \leftrightarrow y)}) < \vec{\alpha}_{\text{max}}\}$$

▷ $A_E^{(x \leftrightarrow y)}$ is the associated affine map to $(x \leftrightarrow y)$ in A_T ; and $A_H^{(x \leftrightarrow y)}$ is the best approximating affine map at $(x \leftrightarrow y)$ computed from the testing homography η_j using Equation 3.
if useNFA **then**

$$\varepsilon_{j,\cdot} = \text{SORT} \left(\left\{ \xi_\eta^8((x \leftrightarrow y)) \right\}_{(x \leftrightarrow y) \in S_j} \right)$$

foreach $k \in \{3, \dots, |S_j|\}$ **do**

$$\left[\text{NFA}_k^j = (|M_T| - 2) \binom{|M_T|}{k} \binom{k}{2} \left(\frac{\varepsilon_{j,k}^8 \pi^4}{w_q h_q w_t h_t 12^2 \pi^2} \right)^{k-2} \right]$$

if useNFA **then**

$$j^*, k^* = \arg \min_{j,k} \text{NFA}_k^j$$

if $\text{NFA}_{k^*}^{j^*} < 1$ **then**
 M_I is the subset of S^{j^*} achieving the first k^* smaller values in $\varepsilon_{j^*,\cdot}$.

else
 $M_I = \emptyset$
else

$$j^* = \arg \max_j |S_j \cap T_j|$$

 $M_I = S_{j^*} \cap T_{j^*}$ **if** $|S_{j^*} \cap T_{j^*}| > 2$ **else** \emptyset
return M_I, η_{j^*}

Table 1 Homography estimation performances for RANSAC USAC [55], RANSAC_{base}, RANSAC_{2pts} and RANSAC_{affine}. Four combinations of detectors and descriptors are used: SIFT [32] + AID [62]; SIFT [32] + HardNet [40]; HessianAffine [38] + AID [62]; HessianAffine [38] + HardNet [40]. Affine approximations are provided by each method presented in Section 3.1; ‘None’ states that no affine information is provided. Each RANSAC runs for 1000 internal iterations, except for USAC that adapts the number of iterations at each execution. To measure probability of success, all RANSACs were run 20 times on resulting matches from each pair of images. Legend: S - the number of successes (bounded by $20 \times \boxed{\text{number}}$); the number of correctly matched image pairs; inl. - the average number of correct inliers; AvE - the average pixel error. All these metric indicators are computed by first thresholding the symmetric transfer error with respect to the ground truth homography of all matches in consensus with the estimated homography. The $\boxed{\text{numbers}}$ of image pairs in a dataset are boxed.

Detector + Descriptor	Affine maps	Homography Estimator	EF dataset [82]				EVD dataset [41]				OxAff dataset [39]				SymB dataset [26]			
			S	$\boxed{33}$	inl.	AvE	S	$\boxed{15}$	inl.	AvE	S	$\boxed{40}$	inl.	AvE	S	$\boxed{46}$	inl.	AvE
SIFT + AID	None	RANSAC _{base}	182	19	85	6.7	17	1	53	7.1	714	39	1651	4.8	220	22	484	6.6
		USAC	176	10	144	5.3	0	0	0	-	659	35	1887	4.6	200	11	767	6.6
	LOCATE	RANSAC _{2pts}	367	24	91	6.7	21	2	71	6.9	783	40	1603	5.0	382	31	348	7.0
		RANSAC _{affine}	400	28	44	6.3	21	2	58	6.4	794	40	939	4.6	417	35	168	7.0
	Affnet	RANSAC _{2pts}	341	26	71	6.5	19	1	42	8.1	737	39	1217	4.7	365	32	252	6.9
		RANSAC _{affine}	350	28	31	6.1	19	1	28	6.9	765	40	539	4.3	393	35	120	6.8
	Naive	RANSAC _{2pts}	355	26	80	6.5	14	1	31	8.6	756	40	1299	4.8	402	33	297	7.0
		RANSAC _{affine}	381	27	49	6.5	19	3	11	9.1	721	38	816	4.7	450	37	151	7.0
SIFT + HardNet	None	RANSAC _{base}	500	25	48	3.3	0	0	0	-	760	38	784	2.2	580	29	114	2.9
		USAC	500	25	48	3.3	0	0	0	-	780	39	763	2.2	580	29	113	2.9
	LOCATE	RANSAC _{2pts}	560	28	43	3.6	0	0	0	-	780	39	764	2.3	620	31	106	3.0
		RANSAC _{affine}	560	28	31	3.7	0	0	0	-	780	39	532	2.2	620	31	71	3.3
	Affnet	RANSAC _{2pts}	560	28	43	3.5	0	0	0	-	780	39	712	2.1	620	31	94	2.9
		RANSAC _{affine}	560	28	26	3.4	0	0	0	-	780	39	370	2.0	620	31	59	3.1
	Naive	RANSAC _{2pts}	560	28	43	3.6	0	0	0	-	780	39	717	2.2	620	31	105	3.0
		RANSAC _{affine}	580	29	38	3.4	0	0	0	-	729	38	553	2.2	620	31	80	2.8
HessAff + AID	None	RANSAC _{base}	74	8	33	4.8	6	2	13	5.9	582	34	161	2.4	92	12	87	3.6
		USAC	63	4	68	4.9	0	0	0	-	580	29	170	2.3	67	4	135	3.3
	LOCATE	RANSAC _{2pts}	214	20	30	5.3	25	2	12	7.2	693	38	142	2.5	227	22	49	4.2
		RANSAC _{affine}	215	21	14	4.3	24	2	7	5.4	689	38	84	2.1	222	24	33	3.7
	Affnet	RANSAC _{2pts}	202	17	25	5.0	7	2	9	7.9	659	35	123	2.5	225	22	42	4.2
		RANSAC _{affine}	216	19	10	4.4	20	2	5	6.9	668	37	57	2.2	224	24	22	3.7
	Naive	RANSAC _{2pts}	185	17	27	5.0	3	1	8	7.5	660	36	126	2.5	219	22	39	4.3
		RANSAC _{affine}	178	17	13	4.2	13	2	4	6.9	657	37	58	2.2	198	20	22	3.7
HessAff + HardNet	None	RANSAC _{base}	418	21	27	3.8	0	0	0	-	721	37	134	1.6	384	22	36	3.6
		USAC	431	23	26	3.7	0	0	0	-	720	37	134	1.5	402	22	35	3.6
	LOCATE	RANSAC _{2pts}	535	27	22	3.6	0	0	0	-	760	38	127	1.6	517	26	27	3.5
		RANSAC _{affine}	535	27	13	3.2	0	0	0	-	760	38	78	1.5	555	28	18	2.9
	Affnet	RANSAC _{2pts}	516	26	21	3.5	0	0	0	-	760	38	121	1.5	460	24	28	3.3
		RANSAC _{affine}	543	29	11	3.4	20	1	2	1.9	760	38	61	1.4	563	30	13	3.5
	Naive	RANSAC _{2pts}	517	26	21	3.2	0	0	0	-	760	38	123	1.5	499	25	25	3.7
		RANSAC _{affine}	554	29	11	3.4	20	1	2	1.9	760	38	61	1.4	560	29	13	3.8

Table 2 A-contrario homography estimation performances for ORSA [43], RANSAC_{base}, RANSAC_{2pts} and RANSAC_{affine}. Four combinations of detectors and descriptors are used: SIFT [32] + AID [62]; SIFT [32] + HardNet [40]; HessianAffine [38] + AID [62]; HessianAffine [38] + HardNet [40]. Affine approximations are provided by the LOCATE method presented in Section 3.1; ‘None’ states that no affine information is provided. Each RANSAC runs for **1000** internal iterations, except for ORSA that runs for **10000** iterations. To measure probability of success, all RANSACs were run 20 times on resulting matches from each pair of images. Legend: S - the number of successes (bounded by $20 \times \boxed{\text{number}}$); the number of correctly matched image pairs; inl. - the average number of correct inliers; AvE - the average pixel error. All these metric indicators are computed by first thresholding the symmetric transfer error with respect to the ground truth homography of all matches in consensus with the estimated homography. The $\boxed{\text{numbers}}$ of image pairs in a dataset are boxed.

Detector + Descriptor	Affine maps	A-contrario Homography Estimator	EF dataset [82]				EVD dataset [41]				OxAff dataset [39]				SymB dataset [26]			
			S	$\boxed{33}$	inl.	AvE	S	$\boxed{15}$	inl.	AvE	S	$\boxed{40}$	inl.	AvE	S	$\boxed{46}$	inl.	AvE
SIFT + AID	None	RANSAC _{base}	201	19	80	6.4	13	1	53	7.2	712	39	1539	4.4	221	21	489	6.5
		ORSA	279	21	119	6.2	19	1	172	6.8	783	40	1519	4.7	408	30	469	7.0
	LOCATE	RANSAC _{2pts}	363	26	89	6.5	21	2	80	6.9	780	40	1520	4.8	388	33	337	7.0
		RANSAC _{affine}	377	27	66	5.9	22	2	63	7.5	780	40	1208	3.9	380	31	272	6.5
SIFT + HardNet	None	RANSAC _{base}	460	23	48	3.2	0	0	0	-	740	37	702	1.6	500	25	114	2.6
		ORSA	506	26	44	3.0	0	0	0	-	760	38	686	1.6	579	29	99	2.6
	LOCATE	RANSAC _{2pts}	500	25	46	3.2	0	0	0	-	760	38	726	1.9	580	29	104	2.8
		RANSAC _{affine}	500	25	43	3.2	0	0	0	-	760	38	705	1.9	580	29	96	2.7
HessAff + AID	None	RANSAC _{base}	85	10	32	4.5	4	1	15	7.4	583	35	139	1.9	95	12	78	3.5
		ORSA	138	12	39	4.3	15	1	41	6.2	651	35	130	1.5	174	16	58	2.7
	LOCATE	RANSAC _{2pts}	209	19	30	4.9	19	2	12	6.3	687	38	133	2.1	221	21	47	4.0
		RANSAC _{affine}	187	18	26	4.1	10	1	10	4.6	678	37	118	1.8	215	20	41	3.4
HessAff + HardNet	None	RANSAC _{base}	405	21	25	3.2	0	0	0	-	716	36	123	1.3	363	20	35	3.0
		ORSA	432	22	25	3.2	0	0	0	-	716	36	123	1.3	402	21	32	2.9
	LOCATE	RANSAC _{2pts}	453	23	24	3.6	0	0	0	-	740	37	124	1.3	415	21	32	3.5
		RANSAC _{affine}	430	22	23	3.2	0	0	0	-	740	37	117	1.3	380	19	31	3.2

In Section 4 we have introduced, for the proposed RANSACs, a procedure allowing to order and validate estimations with respect to a measure of statistical significance, the NFA. However, the typical thresholding for determining inliers is less costly and might be a preferred option depending on the application. With this in mind, experiments are separated in two (Subsection 5.1 and Subsection 5.2) and are presented in form of tables. Each table should be analyzed by blocks. Each block consists of a fixed matching method and dataset. The reader should pay special attention to the first two columns: total number of successes and total number of image pairs identified at least once. The best method should have these two indicators as high as possible, for most choices of matching method and dataset (i.e. blocks). In all cases, the gap between RANSAC_{base} and RANSAC_{2pts}/RANSAC_{affine} will give a measure of the improvement provided by the local affine approximations.

5.1 Fixed thresholds for inlier discrimination

Table 1 shows a comparison of homography estimation methods using fixed thresholds for inlier/outlier discrimination. The performance of each RANSAC (combined with the affine approxi-

mations of Section 3.1) is tested on the four aforementioned datasets for the four choices of detectors and descriptors presented above. The RANSAC USAC [55] method is added to Table 1 in order to compare the proposed RANSACs against a well established state-of-the-art method.

The USAC [55] method shows better performances than $\text{RANSAC}_{\text{base}}$ when equipped with the HardNet [40] descriptor. This situation is inverted for the AID [62] descriptor. A plausible explanation is the higher rate of false positive matches for the AID [62] descriptor with respect to the HardNet [40] descriptor, which might be harming the performance of some steps (like the $T_{d,d}$ and Bail-Out tests, among others) present in USAC [55].

The reader will also note the performance drop of the LOCATE method in estimating affine approximations between image pairs in the case of the SymB [26] dataset. Indeed, LOCATE misreads the information when analyzing the painting-to-photo patch pairs, an invariance for which LOCATE was not trained. The Affnet [42] method is less affected by these painting-to-photo image pairs as it analyzes separately each patch, so it is the structure that provides the invariance. However, the Affnet and Naive methods show similar performances under the SymB [26] dataset.

Incorporating the affine information to the homography fitting step allowed to boost the total number of successes in retrieving ground truth homographies for almost any configuration of detector and descriptor. Indeed, having decreased the sample size (from 4 to 2) has increased the probability that at least one of the 1000 random samples that were drawn while iterating is free from outliers. This implies that the homography fitting step is more likely to capture the true homography in fewer iterations. Therefore, the processing time spent in computing local approximating affine maps could be compensated later on by decreasing the number of internal iterations. Furthermore, we have observed that, in general, even if $\text{RANSAC}_{\text{affine}}$ produces less apparent inliers, the quality of those matches yields a higher probability of success for the same number of internal iterations. Moreover, the affine approximations provided by all methods presented in Section 3.1 often resulted in an added value.

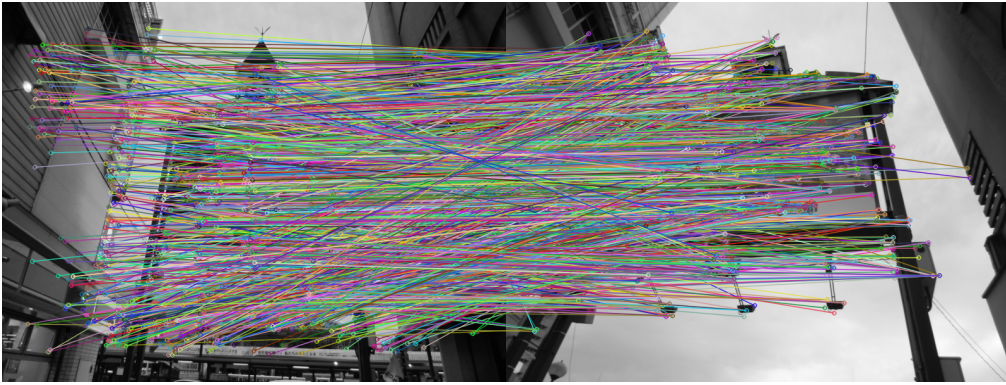
5.2 Adaptive thresholds for inlier discrimination

A comparison of the proposed a-contrario RANSACs is provided in Table 2. ORSA [43] is used for benchmarking. Note that ORSA runs for 10000 iterations whereas the proposed methods runs for 1000 iterations. Nevertheless, $\text{RANSAC}_{2\text{pts}}$ and $\text{RANSAC}_{\text{affine}}$ attain comparable (sometimes better) results with respect to ORSA [43]. This points out that the affine information is making up for the 9000 never-done iterations in the proposed a-contrario RANSACs.

The reader will notice in Table 1 that HessianAffine [38] + HardNet [40] combined with two of the proposed methods have systematically exhibit 2 correct matches out of 3 inliers of the estimated homography; whereas in Table 2 the a-contrario methods did not validate those 3 matches to be of statistical significance.

6 Conclusions

In this paper we reviewed three methods for estimating local affine maps between images. They provide first-order approximations of local geometry. This information is proved to be beneficial for homography estimation, for which we presented several RANSAC versions that systematically improved results in four well known datasets [26, 39, 41, 82]. The proposed RANSACs regularly improved the number of successes in retrieving the ground truth homographies with respect to a baseline RANSAC, and, in a minor degree, with respect to well-established homography estimation methods like USAC [55] and ORSA [43]. The Number of False Alarms (NFA) from the *a-contrario* procedure [11] helps us measure the soundness of estimated homographies and allows for adaptive



(a) Resulting matches from the SIFT-AID [62] method.



(b) Transforming query into target with the best scoring homography found by the a-contrario $\text{RANSAC}_{\text{affine}}$ equipped with the LOCATE method.



(c) Matches in consensus with the above homography.



(d) 10 random matches from (c) with their estimated affine approximations.

Figure 6: Visual results associated with this demo.

thresholds of inlier/outlier discrimination. The computations needed for estimating local affine maps around each match can be compensated later on by reducing the number of internal iterations of these RANSAC algorithms.

References

- [1] SAMEER AGARWAL, YASUTAKA FURUKAWA, NOAH SNAVELY, IAN SIMON, BRIAN CURLESS, STEVEN M SEITZ, AND RICHARD SZELISKI, *Building rome in a day*, Communications of the ACM, 54 (2011), pp. 105–112.
- [2] A AGARWALA, M AGRAWALA, M COHEN, D SALESIN, AND R SZELISKI, *Photographing long scenes with multi-viewpoint panoramas*, International Conference on Computer Graphics and Interactive Techniques, (2006), pp. 853–861.
- [3] DANIEL BARATH AND LEVENTE HAJDER, *Novel ways to estimate homography from local affine transformations*, (2016).
- [4] A. BAUMBERG, *Reliable feature matching across widely separated views*, CVPR, 1 (2000), pp. 774–781.
- [5] M. BENNEWITZ, C. STACHNISS, W. BURGARD, AND S. BEHNKE, *Metric Localization with Scale-Invariant Visual Features Using a Single Perspective Camera*, European Robotics Symposium, (2006).
- [6] M BROWN AND D LOWE, *Recognising panoramas*, in Proc. the 9th Int. Conf. Computer Vision, October, 2003, pp. 1218–1225.
- [7] MATTHEW BROWN AND SABINE SÜSTRUNK, *Multi-spectral SIFT for scene category recognition*, in Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on, IEEE, 2011, pp. 177–184.
- [8] F. CAO, J.-L. LISANI, J.-M. MOREL, P. MUSÉ, AND F. SUR, *A Theory of Shape Identification*, Springer Verlag, 2008.
- [9] E. Y. CHANG, *EXTENT: fusing context, content, and semantic ontology for photo annotation*, Proceedings of the 2nd international workshop on Computer vision meets databases, (2005), pp. 5–11.
- [10] DAVIDE COZZOLINO, GIOVANNI POGGI, AND LUISA VERDOLIVA, *Efficient dense-field copy-move forgery detection*, IEEE Transactions on Information Forensics and Security, 10 (2015), pp. 2284–2297.
- [11] A. DESOLNEUX, L. MOISAN, AND J.-M. MOREL, *From Gestalt Theory to Image Analysis*, Springer, 2008.
- [12] D. DETONE, T. MALISIEWICZ, AND A. RABINOVICH, *Deep image homography estimation*, arXiv preprint arXiv:1606.03798, (2016).
- [13] Q FAN, K BARNARD, A AMIR, A EFRAT, AND M LIN, *Matching slides to presentation videos using SIFT and scene background matching*, Proceedings of the 8th ACM international workshop on Multimedia information retrieval, (2006), pp. 239–248.

- [14] EREZ FARHAN AND RAMI HAGEGE, *Geometric expansion for local feature analysis and matching*, SIAM Journal on Imaging Sciences, 8 (2015), pp. 2771–2813.
- [15] EREZ FARHAN, ELAD MEIR, AND RAMI HAGEGE, *Local Region Expansion: a Method for Analyzing and Refining Image Matches*, Image Processing On Line, 7 (2017), pp. 386–398.
- [16] O FAUGERAS, *Three-Dimensional Computer Vision: A Geometric Viewpoint*, MIT Press, 1993.
- [17] L FÉVRIER, *A wide-baseline matching library for Zeno*, Internship report, www.di.ens.fr/~fevrier/papers/2007-InternsipReportILM.pdf, (2007).
- [18] M. A. FISCHLER AND R. C. BOLLES, *Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography*, Communications of the ACM, 24 (1981), pp. 381–395.
- [19] G FRITZ, C SEIFERT, M KUMAR, AND L PALETTA, *Building detection from mobile imagery using informative SIFT descriptors*, Lecture Notes in Computer Science, pp. 629–638.
- [20] ANDREAS GEIGER, JULIUS ZIEGLER, AND CHRISTOPH STILLER, *Stereoscan: Dense 3d reconstruction in real-time*, in Intelligent Vehicles Symposium (IV), 2011 IEEE, Ieee, 2011, pp. 963–968.
- [21] YUNCHAO GONG, SVETLANA LAZEBNIK, ALBERT GORDO, AND FLORENT PERRONNIN, *Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 35 (2013), pp. 2916–2929.
- [22] I GORDON AND D G LOWE, *What and Where: 3D Object Recognition with Accurate Pose*, Lecture Notes in Computer Science, 4170 (2006), p. 67.
- [23] J S HARE AND P H LEWIS, *Salient regions for query by image content*, Image and Video Retrieval: Third International Conference, CIVR, (2004), pp. 317–325.
- [24] C HARRIS AND M STEPHENS, *A combined corner and edge detector*, Alvey Vision Conference, 15 (1988), p. 50.
- [25] R. HARTLEY AND A. ZISSERMAN, *Multiple view geometry in computer vision*, Cambridge university press, 2003.
- [26] D. C. HAUAGGE AND N. SNAVELY, *Image matching using local symmetry features*, in 2012 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2012, pp. 206–213.
- [27] T KADIR, A ZISSERMAN, AND M BRADY, *An Affine Invariant Salient Region Detector*, in ECCV, 2004, pp. 228–241.
- [28] B N LEE, W Y CHEN, AND E Y CHANG, *Fotofiti: web service for photo management*, Proceedings of the 14th annual ACM international conference on Multimedia, (2006), pp. 485–486.
- [29] STEFAN LEUTENEGGER, MARGARITA CHLI, AND ROLAND Y. SIEGWART, *BRISK: Binary Robust invariant scalable keypoints*, in Proceedings of the IEEE International Conference on Computer Vision, 2011, pp. 2548–2555.

- [30] T. LINDBERG AND J. GARDING, *Shape-adapted smoothing in estimation of 3-D depth cues from affine distortions of local 2-D brightness structure*, ECCV, (1994), pp. 389–400.
- [31] CE LIU, JENNY YUEN, AND ANTONIO TORRALBA, *Sift flow: Dense correspondence across scenes and its applications*, IEEE transactions on pattern analysis and machine intelligence, 33 (2011), pp. 978–994.
- [32] D. LOWE, *Distinctive image features from scale-invariant keypoints*, IJCV, 60 (2004), pp. 91–110.
- [33] G LOY AND J O EKLUNDH, *Detecting symmetry and symmetric constellations of features*, Proceedings of ECCV, 2 (2006), pp. 508–521.
- [34] ELMAR MAIR, GREGORY D. HAGER, DARIUS BURSCHKA, MICHAEL SUPPA, AND GERHARD HIRZINGER, *Adaptive and generic corner detection based on the accelerated segment test*, in Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), vol. 6312 LNCS, 2010, pp. 183–196.
- [35] J. MATAS, O. CHUM, M. URBAN, AND T. PAJDLA, *Robust wide-baseline stereo from maximally stable extremal regions*, IVC, 22 (2004), pp. 761–767.
- [36] K. MIKOLAJCZYK AND C. SCHMID, *Indexing based on scale invariant interest points*, ICCV, 1 (2001), pp. 525–531.
- [37] —, *An affine invariant interest point detector*, ECCV, 1 (2002), pp. 128–142.
- [38] —, *Scale and Affine Invariant Interest Point Detectors*, IJCV, 60 (2004), pp. 63–86.
- [39] 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, 65 (2005), pp. 43–72.
- [40] ANASTASHIA MISHCHUK, DMYTRO MISHKIN, FILIP RADENOVIC, AND JIRI MATAS, *Working hard to know your neighbor’s margins: Local descriptor learning loss*, in Advances in Neural Information Processing Systems, 2017, pp. 4826–4837.
- [41] D. MISHKIN, J. MATAS, AND M. PERDOCH, *MODS: Fast and robust method for two-view matching.*, CVIU, 141 (2015), pp. 81–93.
- [42] D. MISHKIN, F. RADENOVIC, AND J. MATAS, *Repeatability is not enough: Learning affine regions via discriminability*, in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 284–300.
- [43] L. MOISAN, P. MOULON, AND P. MONASSE, *Automatic Homographic Registration of a Pair of Images, with A Contrario Elimination of Outliers*, IPOL, 2 (2012), pp. 56–73.
- [44] —, *Fundamental Matrix of a Stereo Pair, with A Contrario Elimination of Outliers*, IPOL, 6 (2016), pp. 89–113.
- [45] P MOREELS AND P PERONA, *Common-frame model for object recognition*, Neural Information Processing Systems, (2004).
- [46] —, *Evaluation of Features Detectors and Descriptors based on 3D Objects*, IJCV, 73 (2007), pp. 263–284.

- [47] J M MOREL AND G YU, *On the consistency of the SIFT Method*, Tech. Report Prepublication, to appear in Inverse Problems and Imaging (IPI), CMLA, ENS Cachan, 2008.
- [48] J.-M. MOREL AND G. YU, *ASIFT: A new framework for fully affine invariant image comparison*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 438–469.
- [49] A MURARKA, J MODAYIL, AND B KUIPERS, *Building Local Safety Maps for a Wheelchair Robot using Vision and Lasers*, in Proceedings of the The 3rd Canadian Conference on Computer and Robot Vision, IEEE Computer Society Washington, DC, USA, 2006.
- [50] P. MUSÉ, F. SUR, F. CAO, AND Y. GOUSSEAU, *Unsupervised thresholds for shape matching*, ICIP, (2003).
- [51] P. MUSÉ, F. SUR, F. CAO, Y. GOUSSEAU, AND J.-M. MOREL, *An A Contrario Decision Method for Shape Element Recognition*, IJCV, 69 (2006), pp. 295–315.
- [52] A NEGRE, H TRAN, N GOURIER, D HALL, A LUX, AND J L CROWLEY, *Comparative study of People Detection in Surveillance Scenes*, Structural, Syntactic and Statistical Pattern Recognition, Proceedings Lecture Notes in Computer Science, 4109 (2006), pp. 100–108.
- [53] D NISTER AND H STEWENIUS, *Scalable recognition with a vocabulary tree*, CPVR, (2006), pp. 2161–2168.
- [54] Y. PANG, W. LI, Y. YUAN, AND J. PAN, *Fully affine invariant SURF for image matching*, Neurocomputing, 85 (2012), pp. 6–10.
- [55] R. RAGURAM, O. CHUM, M. POLLEFEYS, J. MATAS, AND J.-M. FRAHM, *USAC: a universal framework for random sample consensus*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 35 (2013), pp. 2022–2038.
- [56] M. RAIS, G. FACCILOLO, E. MEINHARDT-LLOPIS, MOREL J.-M., BUADES A., AND COLL B., *Accurate motion estimation through random sample aggregated consensus*, CoRR, abs/1701.05268 (2017).
- [57] CAROLINA RAPOSO AND JOAO P BARRETO, *Theory and practice of structure-from-motion using affine correspondences*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 5470–5478.
- [58] F RIGGI, M TOEWS, AND T ARBEL, *Fundamental Matrix Estimation via TIP-Transfer of Invariant Parameters*, Proceedings of the 18th International Conference on Pattern Recognition (ICPR'06)-Volume 02, (2006), pp. 21–24.
- [59] M. RODRIGUEZ, J. DELON, AND J.-M. MOREL, *Covering the space of tilts. application to affine invariant image comparison*, SIIMS, 11 (2018), pp. 1230–1267.
- [60] —, *Fast affine invariant image matching*, IPOL, 8 (2018), pp. 251–281.
- [61] MARIANO RODRIGUEZ, GABRIELLE FACCILOLO, RAFAEL GROMPONE VON GIOI, PABLO MUSÉ, AND JULIE DELON, *Robust estimation of local affine maps and its applications to image matching*, in WACV, 2020.
- [62] M. RODRIGUEZ, G. FACCILOLO, R. GROMPONE VON GIOI, P. MUSÉ, J.-M. MOREL, AND J. DELON, *Sift-aid: boosting sift with an affine invariant descriptor based on convolutional neural networks*, in ICIP, Sep 2019.

- [63] M. RODRIGUEZ AND R. GROMPONE VON GIOI, *Affine invariant image comparison under repetitive structures*, in ICIP, Oct 2018, pp. 1203–1207.
- [64] MARIANO RODRÍGUEZ, GABRIELLE FACCILOLO, RAFAEL GROMPONE VON GIOI, PABLO MUSÉ, JULIE DELON, AND JEAN-MICHEL MOREL, *Cnn-assisted coverings in the space of tilts: best affine invariant performances with the speed of cnns*, in ICIP, Oct 2020.
- [65] J RUIZ-DEL SOLAR, P LONCOMILLA, AND C DEVIA, *A New Approach for Fingerprint Verification Based on Wide Baseline Matching Using Local Interest Points and Descriptors*, Lecture Notes in Computer Science, 4872 (2007), p. 586.
- [66] PAUL SCOVANNER, SAAD ALI, AND MUBARAK SHAH, *A 3-dimensional SIFT descriptor and its application to action recognition*, in MULTIMEDIA '07: Proceedings of the 15th international conference on Multimedia, New York, NY, USA, 2007, ACM, pp. 357–360.
- [67] S SE, D LOWE, AND J LITTLE, *Vision-based mobile robot localization and mapping using scale-invariant features*, Robotics and Automation, 2001. Proceedings 2001 ICRA. IEEE International Conference on, 2 (2001).
- [68] N SNAVELY, S M SEITZ, AND R SZELISKI, *Photo tourism: exploring photo collections in 3D*, ACM Transactions on Graphics (TOG), 25 (2006), pp. 835–846.
- [69] T. TUYTELAARS AND L. VAN GOOL, *Wide baseline stereo matching based on local, affinely invariant regions*, BMVC, (2000), pp. 412–425.
- [70] ———, *Matching Widely Separated Views Based on Affine Invariant Regions*, IJCV, 59 (2004), pp. 61–85.
- [71] T. TUYTELAARS, L. VAN GOOL, AND OTHERS, *Content-based image retrieval based on local affinely invariant regions*, Int. Conf. on Visual Information Systems, (1999), pp. 493–500.
- [72] CHRISTOFFER VALGREN AND ACHIM J LILIENTHAL, *SIFT, SURF & seasons: Appearance-based long-term localization in outdoor environments*, Robotics and Autonomous Systems, 58 (2010), pp. 149–156.
- [73] KOEN VAN DE SANDE, THEO GEVERS, AND CEES SNOEK, *Evaluating color descriptors for object and scene recognition*, IEEE transactions on pattern analysis and machine intelligence, 32 (2010), pp. 1582–1596.
- [74] M VELOSO, F VON HUNDELSHAUSEN, AND P E RYBSKI, *Learning visual object definitions by observing human activities*, in Proc. of the IEEE-RAS Int. Conf. on Humanoid Robots,, 2005, pp. 148–153.
- [75] M. VERGAUWEN AND L. VAN GOOL, *Web-based 3D Reconstruction Service*, Machine Vision and Applications, 17 (2005), pp. 411–426.
- [76] PHILIPPE WEINZAEPFEL, JEROME REVAUD, ZAID HARCHAOU, AND CORDELIA SCHMID, *Deepflow: Large displacement optical flow with deep matching*, in Proceedings of the IEEE International Conference on Computer Vision, 2013, pp. 1385–1392.
- [77] K YANAI, *Image collector III: a web image-gathering system with bag-of-keypoints*, Proc. of the 16th Int. Conf. on World Wide Web, (2007), pp. 1295–1296.

- [78] G YANG, C V STEWART, M SOFKA, AND C L TSAI, *Alignment of challenging image pairs: Refinement and region growing starting from a single keypoint correspondence*, IEEE Trans. Pattern Anal. Machine Intell., (2007).
- [79] J YAO AND W K CHAM, *Robust multi-view feature matching from multiple unordered views*, Pattern Recognition, 40 (2007), pp. 3081–3099.
- [80] YEFENG ZHENG AND DAVID DOERMANN, *Robust point matching for nonrigid shapes by preserving local neighborhood structures*, IEEE transactions on pattern analysis and machine intelligence, 28 (2006), pp. 643–649.
- [81] HUIYU ZHOU, YUAN YUAN, AND CHUNMEI SHI, *Object tracking using SIFT features and mean shift*, Computer vision and image understanding, 113 (2009), pp. 345–352.
- [82] C. L. ZITNICK AND K. RAMNATH, *Edge foci interest points*, in 2011 International Conference on Computer Vision, IEEE, 2011, pp. 359–366.