# Covering the Space of Tilts.
# Application to Affine Invariant Image Comparison[*]

### Mariano Rodríguez[†], Julie Delon[‡], and Jean-Michel Morel[†]

**Abstract.** We propose a mathematical method to analyze the numerous algorithms performing image matching by affine simulation (IMAS). To become affine invariant they apply a discrete set of affine transforms to the images, prior to the comparison of all images by a scale invariant image matching (SIIM), like SIFT (scale invariant feature transform). Obviously this multiplication of images to be compared increases the image matching complexity. Three questions arise: (a) what is the best set of affine transforms to apply to each image to gain full practical affine invariance? (b) what is the lowest attainable complexity for the resulting method? (c) how is the underlying SIIM method chosen? We provide an explicit answer and a mathematical proof of quasi-optimality of the solution to the first question. As an answer to (b) we find that the near-optimal complexity ratio between full affine matching and scale invariant matching is more than halved, compared to the current IMAS methods. This means that the number of key points necessary for affine matching can be halved, and that the matching complexity is divided by four for exactly the same performance. This also means that an affine invariant set of descriptors can be associated with any image. The price to pay for full affine invariance is that the cardinality of this set is around 6.4 times larger than for a SIIM.

**Key words.** image matching, space of tilts, affine invariance, scale invariance, local descriptors, affine normalization, SIFT, SURF, ASIFT, MODS

**AMS subject classifications.** 68T45, 65D19

**DOI.** 10.1137/17M1140509

**1. Introduction.** Image matching, which consists in detecting shapes common to two images, is a crucial issue for a large number of computer vision applications, such as scene recognition [60, 10, 51] and detection [15, 48], object tracking [65], robot localization [52, 59, 45], image stitching [2, 9], image registration [63, 32] and retrieval [18, 17], three dimensional modeling and reconstruction [14, 16, 61, 1], motion estimation [62], photo management [54], symmetry detection [34], or even image forgeries detection [13]. The problem has implementation variants depending on the setup. If, for example, the user knows that both compared images are related, the focus is on detecting the most reliable common set of shape descriptors. In the detection setup, an image is compared to a database of images and the question is to retrieve related images in the database. This is, for example, crucial for performing a video search [55]. Local shape descriptors must be extracted for this purpose, and this description

[†]CMLA, ENS Paris-Saclay, 94235 Cachan cedex, France (mariano.rodriguez@cmla.ens-cachan.fr, jean-michel.morel@cmla.ens-cachan.fr).
[‡]MAP5, University Paris Descartes, Paris 75006, France (julie.delon@parisdescartes.fr).

should be as invariant as possible to viewpoint changes and, of course, as sparse as possible. In our discussion most of the time we will refer to the simpler set up where two images are being compared. But the reduction of the number of descriptors is, of course, still more important for comparing an image to an image database as initially proposed in [53]. In this last reference, large sets of descriptors are sparsified by clustering techniques. This only indicates how important it is to reduce as much as possible the set of affine descriptors of each image.

*Detectors, descriptors, and affine invariance.* Given a query image of some physical object and a set of target images, the first goal of image matching is to decide if these target images contain a view of the same object. If the answer is positive, image matching aims at localizing this object in these target images. Deciding if the object is present is difficult and becomes especially tricky for large image databases, for which the control of false matches is crucial. Another difficulty of the matching problem comes from the change of camera viewpoints between images. In order to cope with these viewpoint changes, the whole matching process should be as invariant as possible to the resulting image deformations. As we shall develop, this requires affine invariance for the recognition process.

The classical approach to image matching consists in three steps: detection, description, and matching. First, keypoints are detected in the compared images. Second, regions around these points are described and encoded in local invariant descriptors. Finally, all these descriptors are compared and possibly matched. Using local descriptors yields robustness to context changes. Both the detection and description steps are usually designed to ensure some invariance to various geometrical or radiometric changes.

Local image point detectors are always translation invariant. While the venerable Harris point detector [19] is only invariant to translations and rotations, the Harris–Laplace [36], Hessian–Laplace [38], or DoG (difference-of-Gaussian) region detectors [33] are invariant to similarity transformations, i.e., translations, rotations, and scale changes. To ensure invariance to affine transforms, some authors have proposed moment-based region detectors [31, 6] including the Harris-affine and Hessian-affine region detectors [37, 38]. Locally affine invariant region detectors can also be based on edges [58, 57], intensity [56, 57], or entropy [21]. Finally, the detectors MSER ("maximally stable extremal region") [35] and LLD ("level line descriptor") [46, 47, 12] both rely on level lines. Yet the affine invariance of these detectors is limited by the fact that optical blur and affine transforms do not commute, as shown in [44]. Level-line-based detectors like MSER, therefore, are not fit to handle scale changes. Indeed, they do not take into account the effect of blur on the level line geometry [12].

In the last 15 years, numerous invariant image descriptors have been proposed in the literature, but the most well known and the most widely used remains the scale-invariant feature transform (SIFT), introduced by Lowe in his landmark paper [33]. SIFT makes use of a DoG region detector. It is fully invariant to similarities (see [43] for a mathematical proof of this fact). Each *SIFT descriptor* is composed of histograms of gradient orientation around a key point, invariant to local radiometric changes and to geometrical image similarities. As a result, the SIFT method can be considered as partially invariant to illumination and fully invariant to geometrical similarities. But its success is certainly also due to its robustness to reasonable viewpoint changes.

The superiority of SIFT based descriptors has been demonstrated in several comparative studies [39, 42]. As a consequence, many variants of the SIFT descriptor have emerged, among

which we can mention PCA-SIFT [23], GLOH (gradient location-orientation histogram) [39], SURF (speeded up robust features) [7], or RootSIFT [5]. The main claims of these variants are a lower complexity or a greater robustness to viewpoint changes. In the same vein, binary descriptors have also received much attention. Focusing on speed and efficiency, the BRIEF [11], BRISK [25], or LATCH [26] descriptors are compact and represented by sequences of bits, and can be compared more quickly than floating point descriptors like those used in SIFT. Descriptors based on nonlinear scale spaces, such as KAZE [3] or its accelerated version AKAZE [4], have also been proposed to locally adapt blur to the image data.

None of the previously mentioned state-of-the-art methods is fully affine invariant. The SIFT method does not cover the whole affine space and its performance drops under substantial viewpoint changes. SIFT and the other aforementioned descriptors cannot cope with viewpoint differences larger than 60° for planar objects [44, 40], and are still usable but much less efficient for angles larger than 45° [22]. We shall give and use here concrete measurements of their resilience to view angle changes.

To overcome this limitation, several simulation-based solutions have been recently proposed. The core idea of these algorithms, that we choose to call by the generic term *IMAS* (image matching by affine simulation), is to simulate a set of views from the initial images, by varying the camera orientation parameters. These simulations allow us to capture far stronger viewpoint angles than standard matching approaches, up to 88°. Among those IMAS algorithms, we can mention ASIFT [64], FAIR-SURF [49] and MODS [40].

A first suggestion to simulate affine distortions before applying a *SIIM* (scale invariant image matching) appeared in [50] where the authors proposed to simulate two tilts and two shear deformations followed by SIFT in a cloth motion capture application. As argued in [64, 40, 49], if a physical object has a smooth or piecewise smooth boundary, its views obtained by cameras in different positions undergo smooth apparent deformations. These regular deformations are locally well approximated by affine transforms of the image plane. By focusing on local image descriptors, the changes of aspect of objects can therefore be modeled by affine image deformations.

The problem of constructing affine invariant image descriptors by using an affine Gaussian scale space, that is equivalent to simulating affine distortions followed by the heat equation, has a long story starting with [20, 8, 27, 31]. The idea of affine shape adaptation underlying one of the methodologies for achieving affine invariance, was then in turn used as a base for the work on affine invariant interest points and affine invariant matching in [31, 6, 37, 38, 58, 57, 56]. The notion of an affine invariant reference frame was further developed in [29, 30]. Nevertheless, to the best of our knowledge, the direct constructions of affine invariant descriptors as fixed points for an iterative affine normalization process have never found a mathematical justification.

The first IMAS method provided with a mathematical proof of affine invariance is ASIFT [44, 64]. The authors of this paper proposed it as an affine invariant extension of SIFT and proved it to be fully affine invariant in a continuum model. The structure of ASIFT is generic in the sense that it can be implemented with any local descriptor, provided this descriptor has a robustness to viewpoint changes similar to SIFT descriptors. Unlike MSER, LLD, Harris-Affine, and Hessian-Affine, which attempt at normalizing all of the six affine parameters, ASIFT simulates three parameters and normalizes the rest. More specifically, ASIFT simulates the two camera axis parameters, and then applies SIFT which simulates

the scale and normalizes the rotation and the translation. Of the six parameters required for affine invariance, three are therefore simulated and three normalized.

Two recent successful methods follow the same affine simulation path. FAIR-SURF [49] combines the affine invariance of ASIFT and the efficiency of SURF. The MODS image comparison algorithm introduced in [40] also relies on this principle and affine simulations are generated on-demand if needed in the process of comparing two images. MODS employs a combination of different detectors when comparing images. It outperforms state-of-the-art image comparison approaches both in affine robustness and speed.

Other IMAS approaches without local descriptors have also been put up for template matching. FAsT-Match [24] delivers affine invariance by assuming that the template (a patch in the query image) can be recovered inside the target image by a *unique* affine map. Meaning there is no subjacent projective map to identify. Contrary to IMAS with local descriptors, the six required parameters to attain affine invariance are simulated instead of three of the present paper.

In this paper, we are interested in generic IMAS algorithms based on local descriptors and in their geometric optimization. In order to measure the degree of viewpoint change between different views of the same scene, we draw on the concept of *absolute* and *relative transition tilts*, previously introduced in [44, 64], and we illustrate why simulating large tilts on both compared images is necessary to obtain a fully affine invariant recognition. Indeed, transition tilts can in practice be much larger than absolute tilts, since they may behave like the square of absolute tilts.

The key question of IMAS methods is how to choose the list of affine transforms applied to the images before comparison. This list should be as short as possible to limit the computing time. But it should also sample the widest possible range of affine transforms. As we shall see, this question is closely related to the question of finding optimal coverings of the space of affine tilts. This question is formalized and solved in section 2, where we find nearly optimal coverings. Section 3 applies this result to IMAS algorithms. It first presents a complete mathematical theory of IMAS algorithms, proving that they are fully affine invariant under the assumption that the underlying SIIM has a (quantifiable) limited affine invariance. Section 4 gives an experimental validation. It starts by measuring the exact extent of affine invariance for several SIIMs and deduces the corresponding complexity required to attain full affine invariance from each. Section 5 is a conclusion.

**2. The space of affine tilts.** In this section, we introduce the space of tilts for planar affine transforms, and we look for optimal coverings of this space. Optimal coverings will be used in the next section to define an optimal discrete set of affine transformations as the basis for IMAS algorithms. The rest of this section can be read as a sequence of purely geometric results. However, the reader might prefer to keep in mind that the affine transforms considered here can be interpreted as different viewpoints of a camera or, more generally, as the transition from an image taken from a viewpoint to an image taken from another viewpoint. Indeed, given a frontal snapshot of a planar object $u(\mathbf{x}) = u(x, y)$, we can transition from any affine view $Bu$ of the same object to any other affine view $Au$ through the affine transformation $AB^{-1}$. This requires some notation. For any linear invertible map $A \in GL^+(2)$, we denote the affine transform $A$ of a continuous image $u(\mathbf{x})$ by $Au(\mathbf{x}) = u(A\mathbf{x})$. We recall classic notation

for three subsets of the general linear group $GL(2)$ of invertible linear maps of the plane:

$$GL^+(2) = \{A \in GL(2) \mid \det(A) > 0\},$$
$$GO^+(2) = \{A \in GL^+(2) \mid A \text{ is a similarity}\},$$
$$GL^+_*(2) = GL^+(2) \setminus GO^+(2),$$

where we call similarity any combination of a rotation and a zoom, and the symbol $\setminus$ denotes the set difference operator. Our central notion in the discussion is the *tilt* of an affine transform, which we now define.

### 2.1. Absolute tilts.

**Proposition 2.1** (see [44]). *Every $A \in GL^+_*(2)$ is uniquely decomposed as*

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

*where $R_1$, $R_2$ are rotations and $T_t = \left[\begin{smallmatrix} t & 0 \\ 0 & 1 \end{smallmatrix}\right]$ with $t > 1$, $\lambda > 0$, $\phi \in [0, \pi[$, and $\psi \in [0, 2\pi[$.*

*Remark* 2.2. A similar decomposition to (1) was also presented in [28] for small deformations around the identity.

*Remark* 2.3. It follows from this proposition that any affine map $A \in GL^+(2)$ is either uniquely decomposed as in (1) or is directly expressed as a similarity $\lambda R_1$.

Figure 1 shows a camera viewpoint interpretation of this affine decomposition where the longitude $\phi$ and latitude $\theta = arccos\frac{1}{t}$ characterize the camera's viewpoint angles, $\psi$ parameterizes the camera spin and $\lambda$ corresponds to the zoom. In the ideal affine model, the camera is supposed to stand at an infinite distance from a flat image $u$, so that the deformation of $u$ induced by the camera indeed is an affine map. But the above approximation is still valid provided the image's size is small with respect to the camera distance. In other terms the affine model is locally valid for each small and approximately flat patch of a physical surface photographed by a camera at some distance. Yet, the affine deformation of the object's aspect will be different for each of its patches. This explains why affine invariant recognition methods deal with local descriptors. The parameter $t$ defined above measures the so-called *absolute tilt* between the frontal view and a slanted view. The uniqueness of the decomposition in (1) justifies the next definition.
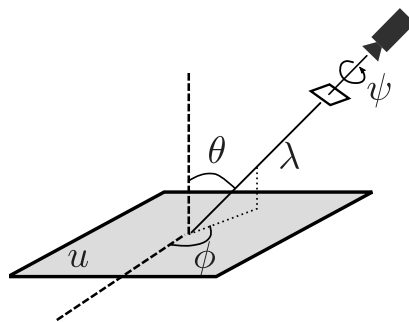


**Figure 1.** *Geometric interpretation of* (1).

**Definition 2.4.** *We call the absolute tilt of A the real number $\tau(A)$ defined by*

$$\begin{cases} GL^+(2) & \to & [1,\infty[\,, \\ A & \mapsto & \begin{cases} 1 & \text{if } A \in GO^+(2)\,, \\ t & \text{if } A \in GL_*^+(2)\,, \end{cases} \end{cases}$$

*where t is the parameter found when applying Proposition* 2.1 *to A.*

**Proposition 2.5.** *Let $A \in GL^+(2)$. Then*

$$\tau(A) = \sqrt{\frac{\lambda_1}{\lambda_2}} = \|\!|A|\!\|_2 \, \|\!|A^{-1}|\!\|_2\,,$$

*where $\lambda_1 \geq \lambda_2$ are the singular values of A and $\|\!|\cdot|\!\|_2$ is the usual Euclidean matrix norm.*

*Proof.* Since the case of a similarity is straightforward, we do not consider it here. Therefore, suppose that $A \in GL_*^+(2)$. Then, using (1) we can rewrite

$$A = R_1 \begin{pmatrix} \gamma_1 & 0 \\ 0 & \gamma_2 \end{pmatrix} R_2,$$

where $R_1, R_2$ are two rotations and $\gamma_1 \geq \gamma_2 > 0$. So

$$A^\star A = R_2^t \begin{pmatrix} \gamma_1^2 & 0 \\ 0 & \gamma_2^2 \end{pmatrix} R_2$$

whose eigenvalues are

$$\lambda_1 = \gamma_1^2 \text{ and } \lambda_2 = \gamma_2^2,$$

but $\gamma_1, \gamma_2 > 0$ imply

$$A = \sqrt{\lambda_2} R_1 \begin{pmatrix} \sqrt{\frac{\lambda_1}{\lambda_2}} & 0 \\ 0 & 1 \end{pmatrix} R_2$$

and, finally, $\tau(A) = \sqrt{\frac{\lambda_1}{\lambda_2}}$. In addition, it is well known that

$$\|\!|A|\!\|_2 = \sqrt{\rho(A^\star A)} = \sqrt{\lambda_1},$$

$$\|\!|A^{-1}|\!\|_2 = \sqrt{\rho\left((AA^\star)^{-1}\right)} = \frac{1}{\sqrt{\lambda_2}},$$

where $\rho(A^\star A)$ is the largest eigenvalue of $A^\star A$, i.e, the largest singular value of $A$. ■

**2.2. Transition tilts.** Image descriptors like those proposed in the SIFT method are invariant to translations, rotations, and Gaussian zooms, which in terms of the camera position interpretation (see Figure 1) correspond to a fronto-parallel motion of the camera, a spin of the camera, and to an optical zoom. We shall focus on the last part $T_t R_2$ of the decomposition (1) because it is the one that is imperfectly dealt with by SIIMs. SIIMs are instead able to detect objects *up to a similarity*. This leads us to the next definition.
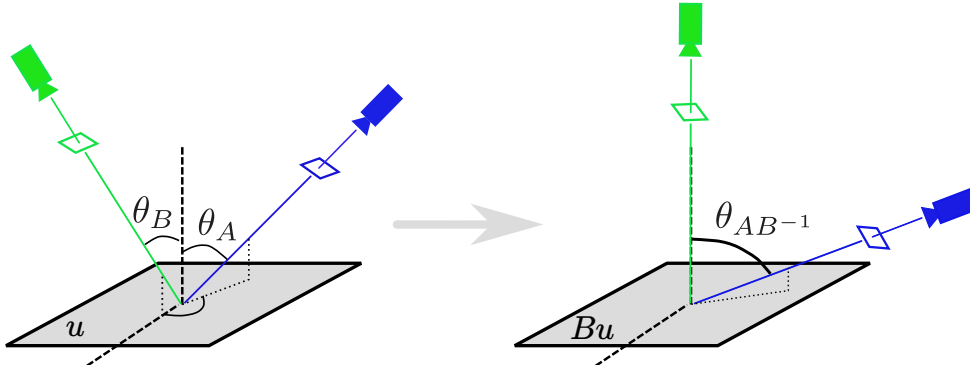
**Figure 2.** *Passage from transition tilts (left side) to absolute tilts (right side).*

**Definition 2.6.** *Let $A, B \in GL^{+}(2)$. Then we define the right equivalence relation $\sim$ as*

$$A \sim B \quad \Leftrightarrow \quad AB^{-1} \in GO^{+}(2).$$

*Remark* 2.7. It is important to notice here that the right and left equivalence relations do differ. For example, take

$$A = T_2 R_{\frac{\pi}{4}} \text{ and } B^{-1} = R_{\frac{\pi}{4}} T_2,$$

then

$$AB^{-1} = 2 R_{\frac{\pi}{2}} \in GO^{+}$$

whereas

$$B^{-1} A = R_{\frac{\pi}{4}} T_4 R_{\frac{\pi}{4}} \notin GO^{+}.$$

**Definition 2.8.** *Let $A, B \in GL^{+}(2)$. We call the* transition tilt *between $A$ and $B$ the absolute tilt of $AB^{-1}$, i.e.,*

$$\tau\left(AB^{-1}\right).$$

The transition tilt has an agreeable visual interpretation appearing in Figure 2. By formula (1) applied to $AB^{-1}$, passing from an image $Bu$ to an image $Au$ comprises a single non-Euclidean transformation, namely, the central tilt matrix $T_{\tau(AB^{-1})}$ which squeezes the image in the direction of $x$ after having rotated it. Thus the transition tilt measures the amount of image distortion caused by a change of view angle. We now state and give a brief proof of the formal properties of the transition tilt stated in [44].

**Proposition 2.9.** *For $A, B \in GL^{+}(2)$ we have*

1. $\tau\left(AB^{-1}\right) = 1 \Leftrightarrow A \sim B$;
2. $\tau(A) = \tau\left(A^{-1}\right)$;
3. $\tau\left(AB^{-1}\right) = \tau\left(BA^{-1}\right)$;
4. $\tau\left(AB^{-1}\right) \leq \tau(A)\,\tau(B)$;
5. $max\left\{\frac{\tau(A)}{\tau(B)}, \frac{\tau(B)}{\tau(A)}\right\} \leq \tau\left(AB^{-1}\right)$.

*Proof.*

(1)
$$\tau\left(AB^{-1}\right) = 1 \Leftrightarrow AB^{-1} = \lambda R \Leftrightarrow A = \lambda RB.$$

(2) By proposition 2.5,
$$\begin{aligned}
\tau\left(A\right) &= \left\|\left\|A\right\|\right\|_2 \left\|\left\|A^{-1}\right\|\right\|_2 \\
&= \tau\left(A^{-1}\right).
\end{aligned}$$

(3) From proof of 2 we have
$$\begin{aligned}
\tau\left(AB^{-1}\right) &= \tau\left(\left(AB^{-1}\right)^{-1}\right) \\
&= \tau\left(BA^{-1}\right).
\end{aligned}$$

(4) By Proposition 2.5
$$\begin{aligned}
\tau\left(AB^{-1}\right) &= \left\|\left\|AB^{-1}\right\|\right\|_2 \left\|\left\|\left(AB^{-1}\right)^{-1}\right\|\right\|_2 \\
&\leq \left\|\left\|A\right\|\right\|_2 \left\|\left\|B^{-1}\right\|\right\|_2 \left\|\left\|B\right\|\right\|_2 \left\|\left\|A^{-1}\right\|\right\|_2 \\
&= \tau\left(A\right)\tau\left(B\right).
\end{aligned}$$

(5) From proof of 4 we have
$$\begin{aligned}
\tau\left(A\right) &= \tau\left(AB^{-1}B\right) \\
&\leq \tau\left(AB^{-1}\right)\tau\left(B\right)
\end{aligned}$$

and the same relation for $B$. ∎

**Definition 2.10.** *We call the* space of tilts, *denoted by $\Omega$, the quotient $GL^+\left(2\right)/\sim$ where the equivalence relation $\sim$ has been given in Definition 2.6.*

This proposition completes Definition 2.6 and clarifies the geometrical interpretation of the space of tilts: an element in the space of tilts represents the set of all the camera spins and zooms associated with a certain tilt in a certain direction.

**Notation 2.11.** *Let $A \in GL^+\left(2\right)$. We denote by $[A]$ the equivalence class in the space of tilts associated with $A$, i.e.,*
$$[A] = \left\{B \in GL^+\left(2\right) \mid A \sim B\right\}.$$

**Definition 2.12.** *We denote by $i$ the canonical injection from the space of tilts to $GL^+\left(2\right)$ defined by*
$$i : \begin{cases} \Omega & \to & GL^+\left(2\right), \\ [A] & \mapsto & T_{\tau(A)}R_{\phi(A)}. \end{cases}$$

This injection filters out the canonical representative from each class which is a mere tilt in the $x$ direction.

*Remark* 2.13. Clearly, the function $i$ satisfies
$$[A] = [i\left([A]\right)]$$

and the space of tilts can be parameterized by picking these representative elements in each class as

$$\Omega = [Id] \bigcup \left\{ \bigcup_{(t,\phi) \in ]1,\infty[ \times [0,\pi[} [T_t R_\phi] \right\}.$$

The next proposition brings an additional justification to Definition 2.10. It means that the transition tilt does not depend on the choice of the class representative in the space of tilts.

**Proposition 2.14.** *Let $A$, $B$, $C$, $D \in GL^+(2)$ satisfying $C \in [A]$ and $D \in [B]$. Then*

$$\tau\left(AB^{-1}\right) = \tau\left(CD^{-1}\right).$$

*Proof.* Let $C \in [A]$, $D \in [B]$. We first remark that if either $A \in GO^+(2)$ or $B \in GO^+(2)$ then the transition tilt operation is, respectively, the absolute tilt of $D$ or $C$, which does not depend on the class representative.

So without loss of generality suppose $A, B \in GL_*^+(2)$. Then, by Proposition 2.1, they are rewritten in a unique way as

$$A = \lambda_A Q_A T_s R_A,$$
$$B = \lambda_B Q_B T_t R_B,$$

and the same result can be applied to the following two matrices,

$$(2) \qquad\qquad AB^{-1} = \lambda_{AB^{-1}} Q_{AB^{-1}} T_{\tau(AB^{-1})} R_{AB^{-1}},$$
$$T_s R_A R_B^{-1} T_t^{-1} = \alpha Q_3 T_{t_3} R_3.$$

Moreover

$$AB^{-1} = \lambda_A Q_A T_s R_A \left(\lambda_B Q_B T_t R_B\right)^{-1}$$
$$= \frac{\alpha \lambda_A}{\lambda_B} \underbrace{(Q_A Q_3)}_{\text{rotation}} T_{t_3} \underbrace{(R_3 Q_B^{-1})}_{\text{rotation}}.$$

Then, by the uniqueness of decomposition in (2) we have $T_{\tau(AB^{-1})} = T_{t_3}$, implying

$$\tau\left(AB^{-1}\right) = \tau\left(T_s R_A R_B^{-1} T_t^{-1}\right).$$

Again, the same methodology applied to

$$C = \lambda_C Q_C A$$
$$= \lambda_C \lambda_A Q_C Q_A T_s R_A$$

and

$$D = \lambda_D Q_D B$$
$$= \lambda_D \lambda_B Q_D Q_B T_t R_B$$

shows that

$$\tau\left(CD^{-1}\right) = \tau\left(T_s R_A R_B^{-1} T_t^{-1}\right) = \tau\left(AB^{-1}\right). \qquad\blacksquare$$

The next proposition follows directly from Proposition 2.9.

**Proposition 2.15.** *The function d*

$$d : \begin{cases} \Omega \times \Omega & \to & \mathbb{R}_+, \\ ([A],[B]) & \mapsto & \log \tau \left( AB^{-1} \right), \end{cases}$$

*is a metric acting on the space of tilts.*

*Proof.* First, $d$ is well defined thanks to Proposition 2.14 which ensures the independence from class representatives. Let us now prove the four metric axioms:

(1) By definition of the absolute tilt $\forall A, B \in GL^+(2)$ one has that $\tau \left( AB^{-1} \right) \geq 1$. This implies

$$d([A],[B]) \geq 0.$$

(2) By Proposition 2.9.1 $\forall A, B \in GL^+(2)$

$$\begin{aligned} d([A],[B]) = 0 &\Leftrightarrow \tau \left( AB^{-1} \right) = 1 \\ &\Leftrightarrow A \sim B \\ &\Leftrightarrow [A] = [B]. \end{aligned}$$

(3) $\forall A, B \in GL^+(2)$, Proposition 2.9.3 states that

$$\tau \left( AB^{-1} \right) = \tau \left( BA^{-1} \right)$$

which implies

$$d([A],[B]) = d([B],[A]).$$

(4) $\forall A, B, C \in GL^+(2)$, Proposition 2.9.4 assures that the following inequality holds:

$$\tau \left( BC^{-1} \left( AC^{-1} \right)^{-1} \right) \leq \tau \left( BC^{-1} \right) \tau \left( AC^{-1} \right).$$

As the logarithm is monotone in $[1, \infty[$, by simply applying it to both sides one obtains the triangular inequality for $d$. ∎

**2.3. Neighborhoods in the space of tilts.** Now that we have introduced the space of tilts and the adequate metric on this space to measure image distortion, we wish to explore optimal coverings for this space. We start by establishing closed formulas for disks in this two dimensional (2D) space.

**Theorem 2.16.** *Given an element of the space of tilts in canonical form $[T_t R(\phi_1)]$, the disk $\mathcal{B}\left([T_t R(\phi_1)], r\right)$ in the space of tilts centered at this element and with radius $r$ corresponds to the following set*

$$\left\{ [T_s R(\phi_2)] \mid G(t, s, \phi_1, \phi_2) \leq \frac{e^{2r} + 1}{2e^r} \right\},$$

*where*

$$G(t, s, \phi_1, \phi_2) = \left( \frac{\frac{t}{s} + \frac{s}{t}}{2} \right) \cos^2(\phi_1 - \phi_2) + \left( \frac{\frac{1}{st} + st}{2} \right) \sin^2(\phi_1 - \phi_2).$$
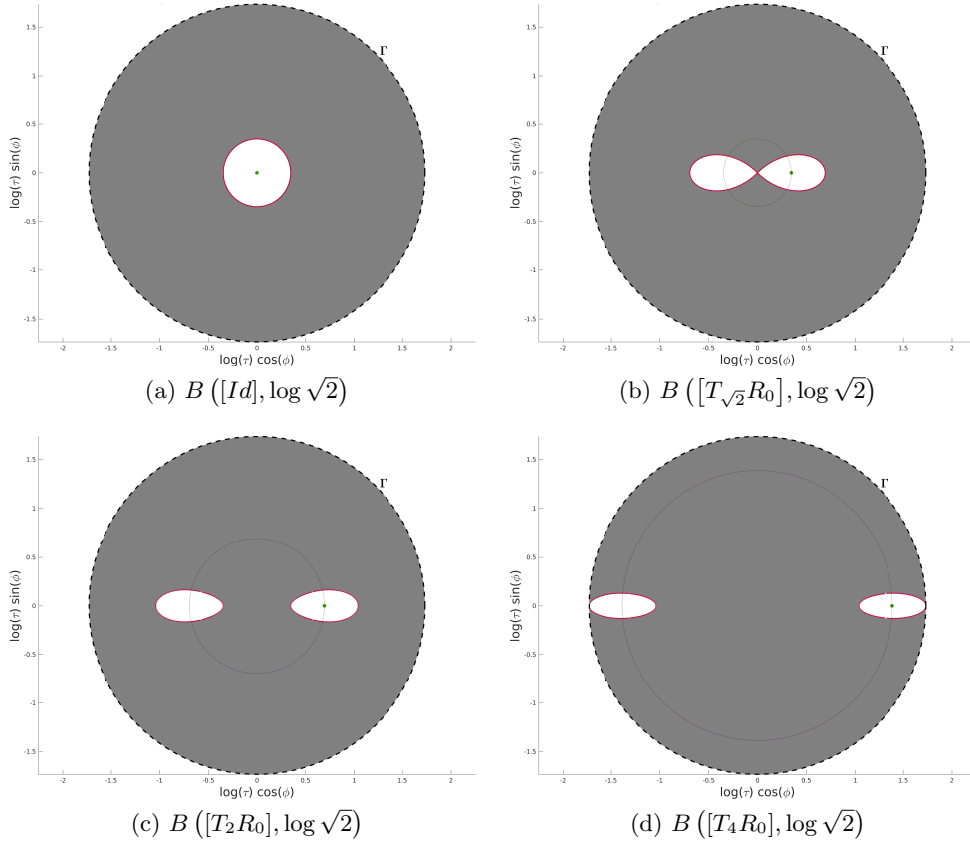
**Figure 3.**   *Polar coordinates.   Green point—affine transformation in question; dashed line—*
$\partial B([Id], \log 4\sqrt{2})$; *dotted line—equal tilts; red line—disk's boundary.*

The proof of this theorem is given in the appendix. Figure 3 displays such disks in polar coordinates $(\log \tau \cos(\phi), \log \tau \sin(\phi))$. This representation will be convenient to visualize region coverings defined by disks in the space of tilts. Figure 4 is illustrating an observation hemisphere, which displays in a geometric environment the space of tilts, the class of affine transformations in question (green dots), and their neighborhoods (black surfaces). Notice that green dots represent camera viewpoints as depicted in Figure 1. In both representations, the pairs $(\tau, \phi)$ and $(\tau, \phi + \pi)$ are denoting the same element of the space of tilts. This is easily interpreted: Two identical images of a planar scene are indeed obtained by an affine camera positioned with a $\pi$ longitude difference.

**Proposition 2.17.** *Let $A, B, C \in GL^{+}(2)$. Then*

$$[A]\, C = [AC]\,,$$

*i.e, classes in $\Omega$ are stable by right multiplication. Moreover,*

$$d\left([AC], [BC]\right) = d\left([A], [B]\right).$$

(a) $B\left([Id], \log \sqrt{2}\right)$



(b) $B\left(\left[T_{\sqrt{2}}R_0\right], \log \sqrt{2}\right)$



(c) $B\left(\left[T_2 R_0\right], \log \sqrt{2}\right)$



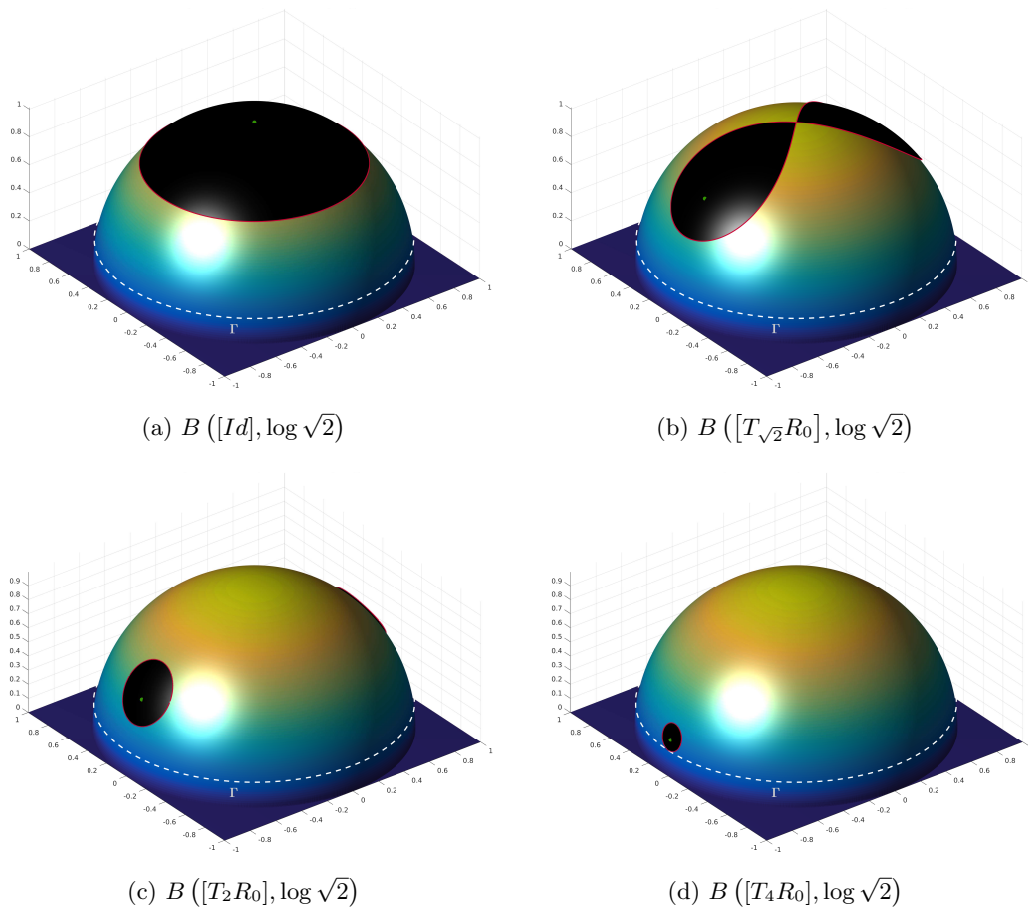(d) $B\left(\left[T_4 R_0\right], \log \sqrt{2}\right)$

**Figure 4.** *Perspective views. Green point—affine transformation in question; dashed line—* $\partial B([Id], \log 4\sqrt{2})$*; black surface—disk in question.*

*Proof.* (1) Proof of $[A]\, C \subset [AC]$.

$$B \in [A] \implies B = \lambda RA$$
$$\implies BC = \lambda RAC$$
$$\implies BC \in [AC].$$

(2) Proof of $[AC] \subset [A]\, C$.

$$D \in [AC] \implies D = \lambda RAC$$
$$\implies D \in [A]\, C.$$

(3)

$$d\left(\left[AC\right],\left[BC\right]\right) = \log\tau\left(AC\left(BC\right)^{-1}\right)$$
$$= \log\tau\left(AB^{-1}\right)$$
$$= d\left(A,B\right). \qquad\blacksquare$$

*Remark* 2.18. Proposition 2.17 guarantees that transition tilts remain unchanged by right compositions. Furthermore, as argued in the proof of Proposition 3.6, the right composition with an element $C \in GL^+\left(2\right)$ could be seen as a modification from a hypothetic frontal image $u$ to another hypothetic frontal image $C^{-1}u$. All this gives both motivation and meaning to the forthcoming Theorem 2.20.

*Remark* 2.19. One might also be interested in the way disks are transformed by left multiplication of elements belonging to $GL^+\left(2\right)$. Unfortunately, in general,

$$C\left[A\right] \neq \left[CA\right].$$

Take, for example, $C = A = T_t$ so

$$R_{\frac{\pi}{2}} = T_t\left(\frac{1}{t}R_{\frac{\pi}{2}}T_t\right) \notin \left[T_{t^2}\right].$$

Furthermore, for $C \in GL^+\left(2\right)$, one has

$$\tau\left(CAB^{-1}C^{-1}\right) = c_2\left(CAB^{-1}C^{-1}\right)$$
$$= \left|\left|\left|CAB^{-1}C^{-1}\right|\right|\right|_2 \left|\left|\left|C\left(AB^{-1}\right)^{-1}C^{-1}\right|\right|\right|_2$$
$$\leq \left|\left|\left|C\right|\right|\right|_2^2 \left|\left|\left|C^{-1}\right|\right|\right|_2^2 \left|\left|\left|AB^{-1}\right|\right|\right|_2 \left|\left|\left|\left(AB^{-1}\right)^{-1}\right|\right|\right|_2$$
$$= \tau\left(C\right)^2 \tau\left(AB^{-1}\right)$$

so, in general,

$$d\left(\left[CA\right],\left[CB\right]\right) \leq 2d\left(\left[C\right],\left[Id\right]\right) + d\left(\left[A\right],\left[B\right]\right).$$

The following theorem will be crucial in the next section to explain why IMAS algorithms are truly affine invariant.

Theorem 2.20. *Let*

$$\Gamma_1 = \mathcal{B}\left(\left[Id\right],\log\Lambda_1\right),$$
$$\Gamma_2 = \mathcal{B}\left(\left[Id\right],\log\Lambda_2\right),$$
$$\Gamma' = \mathcal{B}\left(\left[Id\right],\log\Lambda_2 r\right),$$

*be three neighborhoods of* $\left[Id\right]$ *in* $\Omega$, *where* $\Lambda_1, \Lambda_2, r \in \left[1,\infty\right[$, *and assume that* $\mathbb{S}_1, \mathbb{S}_2 \subset \Omega$ *are two* $\log r$-*coverings of* $\Gamma_1$ *and* $\Gamma'$, *i.e.,*

$$\Gamma_1 \subset \bigcup_{S \in \mathbb{S}_1} \mathcal{B}\left(S,\log r\right),$$
$$\Gamma' \subset \bigcup_{S \in \mathbb{S}_2} \mathcal{B}\left(S,\log r\right).$$

*Then, for every* $[A] \in \Gamma_1$, $[B] \in \Gamma_2$, *there exist* $C \in GL^+(2)$ *with* $\tau(C) \leq r$, $S_A \in \mathbb{S}_1$, *and* $S_B \in \mathbb{S}_2$ *such that*

$$d\left(S_A, \left[(AC)^{-1}\right]\right) = 0,$$

$$d\left(S_B, \left[(BC)^{-1}\right]\right) \leq \log r.$$

A sketch of Theorem 2.20 appears in Figure 5.

*Proof.* Let us set $C = A^{-1}i(S_A)^{-1}$, where $i$ appears in Definition 2.12.

(1) Proof of $d\left(S_A, \left[(AC)^{-1}\right]\right) = 0$. Proposition 2.9.2 directly implies

$$d([Id], [A]) = d\left([Id], \left[A^{-1}\right]\right).$$

Then, as $\mathbb{S}_1$ is a $\log r$-covering of $\Gamma_1$, there exists $S_A \in \mathbb{S}_1$ such that

$$\left[A^{-1}\right] \in \mathcal{B}(S_A, \log r),$$

meaning that the following inequality holds:

$$\begin{aligned}
d\left([Id], \left[A^{-1}i(S_A)^{-1}\right]\right) &= \log \tau\left(A^{-1}i(S_A)^{-1}\right) \\
&= d\left(\left[A^{-1}\right], S_A\right) \\
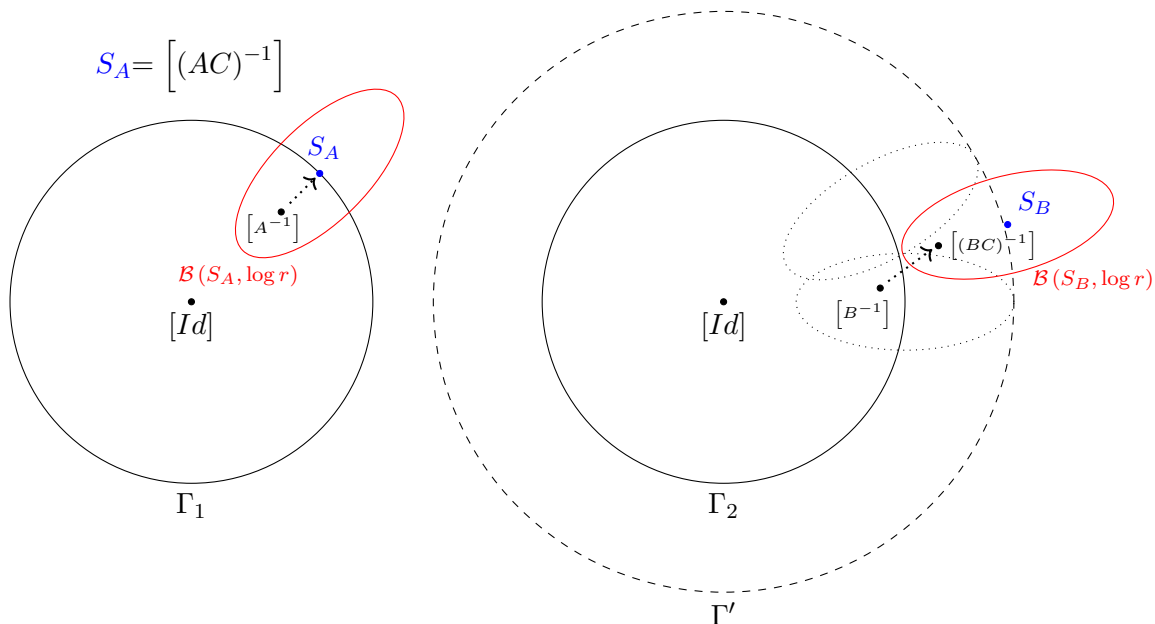&\leq \log r.
\end{aligned}$$



**Figure 5.** *Sketch of Theorem* 2.20.

Finally, as $d$ is a metric (by Proposition 2.15) we know

$$d\left(S_A, \left[(AC)^{-1}\right]\right) = d\left(S_A, [i\left(S_A\right)]\right) = 0.$$

(2) Proof of $d\left(S_B, \left[(BC)^{-1}\right]\right) \leq \log r$. By first using Proposition 2.9 followed by Proposition 2.15, we have

$$\tau\left(BC\right) \leq \tau\left(B\right)\tau\left(C^{-1}\right) = \Lambda_2 r$$
$$\Downarrow$$
$$d\left([Id], \left[(BC)^{-1}\right]\right) = \log\tau\left(BC\right) \leq \log\Lambda_2 r$$
$$\Downarrow$$
$$\left[(BC)^{-1}\right] \in \Gamma'.$$

Once more, as $\mathbb{S}_2$ is a $\log r$-covering of $\Gamma'$, there exists $S_B \in \mathbb{S}_2$ such that

$$\left[(BC)^{-1}\right] \in \mathcal{B}\left(S_B, \log r\right). \qquad \blacksquare$$

**3. Application: Optimal affine invariant image matching algorithms.** The theory and results presented above provide a well-suited geometrical framework for IMAS. This section gives the mathematical formalism and a mathematical proof that IMAS based algorithms are fully affine invariant, up to sampling errors. While the former sections only dealt with affine geometry, we now must introduce in the formalism the camera blur, as we shall deal with digital image recognition. Our goal is to define rigorously affine invariant recognition for digital images.

Consider a continuous and bounded image $u\left(\mathbf{x}\right)$ defined for every $\mathbf{x} = (x, y) \in \mathbb{R}^2$. All continuous image operators including the sampling will be written in capital letters $A$, $B$ and their composition as a mere juxtaposition $AB$.

**Definition 3.1.** *For any $A \in GL^+\left(2\right)$, we define the affine transform $A$ of a continuous image $u$ by*

$$Au(\mathbf{x}) := u(A\mathbf{x}).$$

*Homotheties and rotations acting on continuous images are similarly written as*

$$H_\lambda u\left(\mathbf{x}\right) = u\left(\lambda\mathbf{x}\right),$$
$$R_\phi u\left(\mathbf{x}\right) = u\left(R_\phi\mathbf{x}\right).$$

We now introduce a compact notation for the various convolutions with Gaussians. We shall denote by $\star_x$ the one dimensional (1D) convolution operator in the $x$-direction, i.e.,

$$G \star_x u\left(x, y\right) = \int_{\mathbb{R}} G\left(z\right)u\left(x - z, y\right)dz.$$

Similarly, we denote by $\star_y$ the 1D convolution operator in the $y$-direction. We denote by $\mathbb{G}_\sigma$, $\mathbb{G}_\sigma^x$, and $\mathbb{G}_\sigma^y$, respectively, the 2D and 1D convolution operators in the $x$ and $y$ directions with

$$G_{\mathbf{c}\sigma}(x,y) := \frac{1}{2\pi(\mathbf{c}\sigma)^2} e^{-\frac{x^2+y^2}{2(\mathbf{c}\sigma)^2}},$$

$$G_{\mathbf{c}\sigma}^x(x) := \frac{1}{\sqrt{2\pi}\mathbf{c}\sigma} e^{-\frac{x^2}{2(\mathbf{c}\sigma)^2}},$$

$$G_{\mathbf{c}\sigma}^y(y) := \frac{1}{\sqrt{2\pi}\mathbf{c}\sigma} e^{-\frac{y^2}{2(\mathbf{c}\sigma)^2}},$$

namely,

$$\mathbb{G}_\sigma u := G_{c\sigma} \star u,$$
$$\mathbb{G}_\sigma^x u := G_{c\sigma}^x \star_x u,$$
$$\mathbb{G}_\sigma^y u := G_{c\sigma}^y \star_y u.$$

Here the constant $c \geq 0.7$ is large enough to ensure that all convolved images, initially sampled at distance one, can be subsampled at Nyquist distance $\sigma$ without causing significant aliasing.

*Remark* 3.2. $\mathbb{G}_\sigma$ satisfies the semigroup property

(3) $$\mathbb{G}_\sigma \mathbb{G}_\beta = \mathbb{G}_{\sqrt{\sigma^2+\beta^2}}.$$

By a mere change of variables in the integral defining the convolution, the next formula holds and will be useful:

(4) $$\mathbb{G}_\sigma H_\gamma u = H_\gamma \mathbb{G}_{\sigma\gamma} u.$$

In the classic Shannon–Nyquist framework, we shall denote the image sampling operator (on a unary grid) by $\mathbf{S}_1$. Thus $\mathbf{S}_1 u$ is defined on the grid $\mathbb{Z}^2$. The Shannon–Whittaker interpolator of a digital image on $\mathbb{Z}^2$ will be denoted by $I$.

As developed in [64], the whole image comparison process, based on local features, can proceed as though images were (locally) obtained by using digital cameras that stand far away, at infinity. The geometric deformations induced by the motion of such cameras are affine maps. A model is also needed for the two main camera parameters not deducible from its position, namely, sampling and blur. The digital image is defined on the camera charge coupled device (CCD) plane. The pixel width can be taken as unit length, and the origin and axes chosen so that the camera pixels are indexed by $\mathbb{Z}^2$. The digital initial image is always assumed well sampled and obtained by a Gaussian blur with standard deviation around 0.8. In all that follows, $u_0$ denotes the (theoretical) infinite resolution image that would be obtained by a frontal snapshot of a plane object with infinitely many pixels. The digital image obtained by any camera at infinity is therefore formalized as $\mathbf{u} = \mathbf{S}_1 \mathbb{G}_1 A \mathcal{T} u_0$, where $A$ is *any* linear map with positive singular values and $\mathcal{T}$ any plane translation. Thus we can summarize the general image formation model with cameras at infinity as follows.
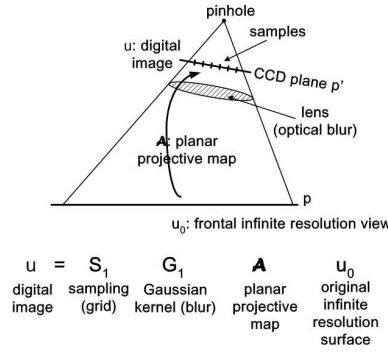
**Figure 6.** *The projective camera model $u = \mathbf{S}_1 \mathbb{G}_1 \mathcal{A} u_0$. $\mathcal{A}$ is a planar projective transform (a homography). $\mathbb{G}_1$ is an antialiasing Gaussian filtering. $\mathbf{S}_1$ is the CCD sampling.*

**Definition 3.3 (image formation model).** *Digital images of a planar object whose frontal infinite resolution image is $u_0$, obtained by a digital camera far away from the object, satisfy*

$$(5) \qquad\qquad \mathbf{u} =: \mathbf{S}_1 \mathbb{G}_1 A \mathcal{T} u_0,$$

*where $A$ is any linear map and $\mathcal{T}$ any plane translation. $\mathbb{G}_1$ denotes a Gaussian kernel broad enough to ensure no aliasing by 1-sampling, namely, $I\mathbf{S}_1 \mathbb{G}_1 A \mathcal{T} u_0 = \mathbb{G}_1 A \mathcal{T} u_0$.*

The image formation model in Definition 3.3 is illustrated in Figure 6.

**3.1. Inverting tilts.** We now formalize the notion of tilt. There are actually three different notions of tilt, that we must carefully distinguish.

**Definition 3.4.** *Given $t > 1$, the tilt factor, define then as follows:*
- *Geometric tilts:*

$$T_t^x u_0(x, y) =: u_0(tx, y),$$
$$T_t^y u_0(x, y) =: u_0(x, ty).$$

- *Simulated tilts (taking into account camera blur):*

$$\mathbb{T}_t^x v =: T_t^x \mathbb{G}_{\sqrt{t^2-1}}^x \star_x v,$$
$$\mathbb{T}_t^y v =: T_t^y \mathbb{G}_{\sqrt{t^2-1}}^y \star_y v.$$

- *Digital tilts (transforming a digital image $u$ into a digital image):*

$$\mathbf{u} \to \mathbf{S}_1 \mathbb{T}_t^x I\mathbf{u},$$
$$\mathbf{u} \to \mathbf{S}_1 \mathbb{T}_t^y I\mathbf{u}.$$

Digital tilts are the ones used in practice. It all adds up because the simulated tilt yields a blur permitting $\mathbf{S}_1$-sampling.

If $u_0$ is an infinite resolution image observed with a camera tilt of $t$ in the $x$ direction, the observed image is $\mathbb{G}_1 T_t^x u_0$. Our main problem is to reverse such tilts. This operation is, in

principle, impossible, because geometric tilts do not commute with blur. However, the first formula of Theorem 3.5 shows that $\mathbb{T}_t^y$ is, up to a zoom out, a pseudoinverse to $T_t^x$.

The meaning of this result is that a tilted image $\mathbb{G}_1 T_t^x u$ can be tilted back by tilting in the orthogonal direction. The price to pay is a $t$ zoom-out. The second relation in the theorem means that the application of the simulated tilt to an image that can be well sampled by $\mathbf{S}_1$ yields an image that keeps that well sampling property.

**Theorem 3.5.** *Let $t \geq 1$. Then*

$$\tag{6} \mathbb{T}_t^y \mathbb{G}_1 T_t^x = \mathbb{G}_1 H_t,$$

$$\tag{7} \mathbb{T}_t^y \mathbb{G}_1 = \mathbb{G}_1 T_t^y.$$

*Proof.* Since $H_t = T_t^y T_t^x$, (6) follows from (7) by composing both sides on the right by $T_t^x$. Let us now prove (7). We shall use the following obvious facts,

$$\tag{8} \mathbb{G}_1 = \mathbb{G}_1^x \mathbb{G}_1^y = \mathbb{G}_1^y \mathbb{G}_1^x,$$

which follows from the separability of the Gaussian and Fubini's theorem and the commutation

$$\tag{9} \mathbb{G}_1^x T_t^y = T_t^y \mathbb{G}_1^x$$

which is true because $\mathbb{G}_1^x$ and $T_t^y$ act separably on the variables $x$ and $y$. Using first (4) in the $y$ dimension, where $T_t^y$ is a mere homothety, and then successively (9), (8), the semigroup property for the Gaussians, and Definition 3.4 we get

$$T_t^y \mathbb{G}_t^y = \mathbb{G}_1^y T_t^y \Rightarrow$$

$$\mathbb{G}_1^x T_t^y \mathbb{G}_t^y = \mathbb{G}_1^x \mathbb{G}_1^y T_t^y \Rightarrow$$

$$T_t^y \mathbb{G}_t^y \mathbb{G}_1^x = \mathbb{G}_1 T_t^y \Rightarrow$$

$$T_t^y \mathbb{G}_{\sqrt{t^2-1}}^y \mathbb{G}_1^y \mathbb{G}_1^x = \mathbb{G}_1 T_t^y \Rightarrow$$

$$\mathbb{T}_t^y \mathbb{G}_1 = \mathbb{G}_1 T_t^y,$$

which proves (7). ∎

The meaning of Theorem 3.5 is that we can design an exact algorithm that simulates all inverse tilts for comparing two digital images. This algorithm handles two images $u = \mathbb{G}_1 A \mathcal{T}_1 w_0$ and $v = \mathbb{G}_1 B \mathcal{T}_2 w_0$ that are two snapshots from different viewpoints of a flat object whose front infinite resolution image is denoted by $w_0$.

**3.2. Proof that IMAS works.** In this section, the formal IMAS algorithm is duly presented (Algorithm 3.1). Our goal is to prove that it works. This proof is a direct application of the results introduced in the previous section. The algorithm and its proof rely on the formal assumption that there exists an image comparison algorithm able to compare image pairs with tilts lower than $r$. The core idea of IMAS algorithms is illustrated in Figure 7.
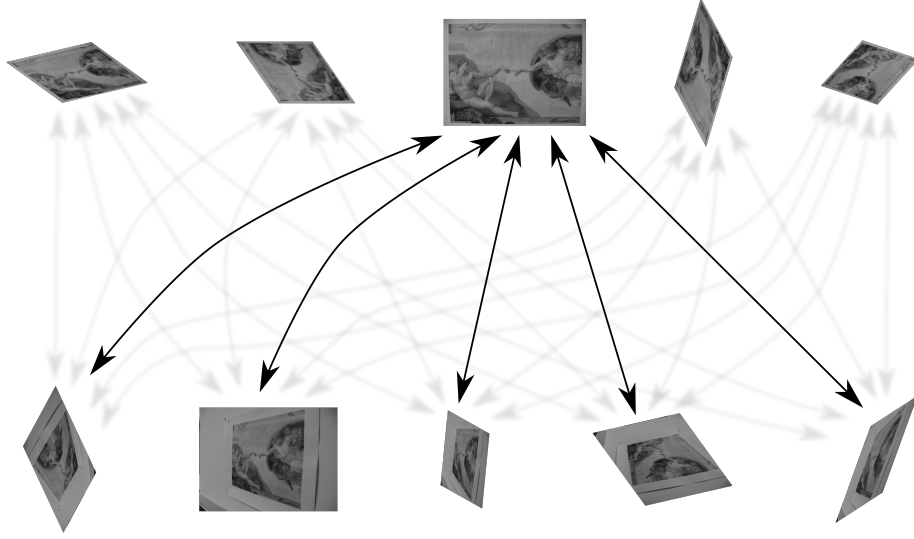
**Figure 7.** *IMAS algorithms start by applying a finite set of optical affine simulations to u and v, followed by pairwise comparisons.*

---

**Algorithm 3.1** Formal IMAS.

---

**Environment:**

Parameters and assumptions from Theorem 2.20 with

$$\mathbb{S}_i = \left\{ \left[ T^x_{t^i_k} R_{\phi^i_k} \right] \right\}_{k=1,\ldots,n_i}.$$

**Input:**

Query and target images: $u$ and $v$.

**Start:**

1: $\forall k = 1, \ldots, n_1$ do

$$u_k = \mathbb{T}^x_{t^1_k} R_{\phi^1_k} u.$$

2: $\forall k = 1, \ldots, n_2$ do

$$v_k = \mathbb{T}^x_{t^2_k} R_{\phi^2_k} v.$$

3: $\forall (k_1, k_2) \in \{1, \ldots, n_1\} \times \{1, \ldots, n_2\}$

$$M_{k_1,k_2} = \text{SIIM-Matches}(u_{k_1}, v_{k_2}).$$

**Output:**

$$M = \bigcup_{(k_1,k_2) \in \{1,\ldots,n_1\} \times \{1,\ldots,n_2\}} M_{k_1,k_2}.$$

---

Proposition 3.6. *Let $u$ and $v$ be, respectively, query and target images which are related by a transition tilt under $\Lambda_1 \Lambda_2$, i.e., there exist a continuous image $w_0$ and $A, B \in GL^+(2)$ with*

$$\tau(A) \le \Lambda_1 \text{ and } \tau(B) \le \Lambda_2$$

*such that*

(10)
$$u = \mathbb{G}_1 A \mathcal{T}_1 w_0 \text{ and } v = \mathbb{G}_1 B \mathcal{T}_2 w_0,$$

*where $\mathcal{T}_1, \mathcal{T}_2$ are planar translations. Then, under the assumptions of Theorem 2.20, the formal IMAS of Algorithm 3.1 generates two affine versions of the images $u$ and $v$ with a transition tilt lower than $r$.*

*Proof.* By Theorem 2.20 there exist $S_A \in \mathbb{S}_1$, $S_B \in \mathbb{S}_2$, and $C \in GL^+(2)$ with $\tau(C) \le r$ such that

$$d\left(S_A, \left[(AC)^{-1}\right]\right) = 0,$$
$$d\left(S_B, \left[(BC)^{-1}\right]\right) \le \log r.$$

Consider the slanted view of the frontal continuous image $w_0$ defined by $w_1 := C^{-1} w_0$. Then we can rewrite query and target images as

$$u = \mathbb{G}_1 AC \mathcal{T}_1 w_1 \text{ and } v = \mathbb{G}_1 BC \mathcal{T}_2 w_1.$$

By Proposition 2.17, the above modification keeps transitions tilts stable, i.e.,

$$d([AC], [BC]) = d([A], [B]),$$

so we can reason as if $w_1$ were the frontal image, instead of $w_0$.

Now, the formal IMAS Algorithm 3.1 will apply $i(S_A) = T^x_{t_A} R_{\phi_A}$ and $i(S_B) = T^x_{t_B} R_{\phi_B}$, respectively, on the query and target images. This is

1. $\mathbb{T}^x_{t_A} R_{\phi_A}$ to $u$, which yields

$$\tilde{u} = \mathbb{G}_1 i(S_A) AC \mathcal{T}_1 w_1$$
$$= \mathbb{G}_1 \lambda R \mathcal{T}_1 w_1;$$

2. $\mathbb{T}^x_{t_B} R_{\phi_B}$ to $v$, which yields

$$\tilde{v} = \mathbb{G}_1 i(S_B) BC \mathcal{T}_2 w_1.$$

But

$$d([Id], [i(S_B) BC]) = \log \tau(i(S_B) BC)$$
$$= d\left(S_B, \left[(BC)^{-1}\right]\right)$$
$$\le \log r$$

which proves that the affine relation between $\tilde{u}$ and $\tilde{v}$ involves a transition tilt under $r$. ∎

| Covered absolute tilts $(\tau(A) \leq \sqrt{r}\Lambda \text{ and } \tau(B) \leq \sqrt{r}\Lambda)$ | Attainable transition tilts $(\tau(AB^{-1}) \leq \Lambda^2)$ | Viewpoint angle $(arccos \frac{1}{\Lambda^2})$ |
|---|---|---|
| $\Lambda = 8$ | 64 | $89.1°$ |
| $\Lambda = 4\sqrt{2}$ | 32 | $88.2°$ |
| $\Lambda = 4$ | 16 | $86.4°$ |
| $\Lambda = 2\sqrt{2}$ | 8 | $82.8°$ |
| $\Lambda = 2$ | 4 | $75.5°$ |
| $\Lambda = \sqrt{2}$ | 2 | $60°$ |

*Remark* 3.7. Two $\log r$-coverings of the same region

$$\Gamma = \mathcal{B}([Id], \log \Lambda)$$

would then ensure that the formal IMAS Algorithm 3.1 manages to reduce transition tilts under $\frac{\Lambda^2}{r}$ between two images into transition tilts under $r$. A relation between covered absolute tilts, attainable transition tilts, and maximal viewpoint angle can be found in Table 1.

**3.3. Optimal discrete coverings in the space of tilts.** We now consider the problem of providing two optimal sets $\mathbb{S}_1, \mathbb{S}_2 \subset \Omega$ permitting the application of Theorem 2.20. These sets should ensure a minimal complexity for the IMAS algorithm. We thus need to define an optimality criterion. We observe that an IMAS algorithm simulates affine transformations on a digital image and then compares descriptors coming from those simulated versions. One would like to minimize the overall number of descriptor comparisons while maintaining the detection efficiency. This minimization *is not* equivalent to a minimization of the number of simulated versions being used. We shall base our efficiency criterion on two straightforward remarks. The first one is that if a digital image suffers a tilt $t$ in any direction, its area gets modified by a factor $\frac{1}{t}$. The second one is that the expected number of keypoints in a digital image is proportional to its area. Both remarks imply that the complexity of an IMAS algorithm will be given by the overall area of the simulated images being ultimately compared. This justifies the next definition.

Definition 3.8. *We call the* area ratio *of* $\mathbb{S}$ *(a finite set of elements in $\Omega$) the real number*

$$\sum_{S \in \mathbb{S}} \frac{1}{\tau(S)}.$$

The area ratio fixes the factor (larger than 1) by which the image area is being multiplied when summing the areas of all of its tilted versions. Then, as the ultimate goal is to reduce the number of key points comparisons, it is natural to look for a set $\mathbb{S}$ whose area ratio is close to the infimum among all $\log r$-coverings of $\Gamma$. Unfortunately, even in $\mathbb{R}^2$, the mathematical problem of finding a covering of a certain set with a minimum amount of disks is well known to be NP-hard. It is therefore difficult to find an optimal solution for our problem, and unlikely that it will be proved to be optimal even if it is. Fortunately, our search space in the set of $\log r$-coverings can be drastically reduced by imposing practical and theoretical constraints to $\mathbb{S}$. Those constraints follow from simple requirements for an image matching method.

**Definition 3.9.** *We shall say that a set $\mathbb{S} \in \Omega$ is feasible if and only if*

1. *$[Id] \in \mathbb{S}$;*
2. *there exist $n \in \mathbb{N}^+$ and*

$$(t_1, \ldots, t_n, \phi_1, \ldots, \phi_n) \in [1, \infty[^n \times \,]0, \pi]^n$$

   *such that*

$$\mathbb{S} \setminus \{[Id]\} = \bigcup_{i=1}^{n} \left\{ [T_{t_i} R_{k\phi_i}] \in \Omega \,|\, k = 0, \ldots, \left\lfloor \frac{\pi}{\phi_i} \right\rfloor \right\},$$

*where $\lfloor a \rfloor$ denotes the nearest integer less than or equal to a real number $a$.*

*Remark* 3.10. Definition 3.9.1 avoids an image resolution loss before comparison, an obvious requirement. Imposing groups of concentric equidistant tilts as in Definition 3.9.2 is a sound isotropy requirement.

**Definition 3.11.** *Set $\Gamma = \mathcal{B}([Id], \log \Lambda)$. A feasible set $\mathbb{S} \in \Omega$ with parameters*

$$(n, (t_1, \ldots, t_n, \phi_1, \ldots, \phi_n)) \in \mathbb{N}^+ \times [1, \infty[^n \times \,]0, \pi]^n$$

*is said to be optimal among feasible sets if and only if it realizes the minimal area ratio. In other words, optimal feasible sets are solutions of*

(11)
$$\underset{(n,(t_1,\ldots t_n,\phi_1,\ldots \phi_n)) \in \mathbb{N}^+ \times [1,\infty[^n \times ]0,\pi]^n}{\arg\min} 1 + \sum_{i=1}^{n} \frac{|J_{t_i,\phi_i}|}{t_i}$$

$$\text{subject to: } \Gamma \subset \mathcal{B}_{[Id]}^{\log r} \cup \left\{ \bigcup_{1 \le i \le n} \bigcup_{S \in J_{t_i,\phi_i}} \mathcal{B}_{[S]}^{\log r} \right\},$$

*where $J_{t_i,\phi_i}$ is the set of transformations of the form*

$$T_{t_i} R_{\phi_i}, T_{t_i} R_{2\phi_i}, \ldots, T_{t_i} R_{\lfloor \frac{\pi}{\phi_i} \rfloor \phi_i},$$

*$|J_{t_i,\phi_i}|$ is the cardinal of $J_{t_i,\phi_i}$, and $\mathcal{B}_{[S]}^{\log r}$ denotes $\mathcal{B}([S], \log r)$.*

Fortunately for our problem with the realistic values $\Lambda = 6$ and $r = 1.8$, $n = 2$ can be fixed, as easy heuristics indicate that any covering with $n > 2$ has a far too large area ratio. Thus our optimization in a realistic setting ends up being performed in dimension 4 for sets $(t_1, t_2, \phi_1, \phi_2)$. With $n$ thus fixed the optimization problem in (11) can be exhaustively optimized. In this minimization we deal with 4 dimensions and more specifically with $100^4$ feasible sets by sampling each parameter. This yields an almost exact discrete exhaustive optimization by sampling densely the explored set $(t_1, t_2, \phi_1, \phi_2)$ with 100 different values for each parameter. The next proposition describes the result of this optimization and verifies that it is indeed feasible.

**Proposition 3.12.** *There exists a feasible $\log 1.8$-covering, depicted in Figure 9(c), with area ratio equal to 6.34. It is an approximated solution of the optimization problem in (11) for $\Gamma = \{[T_t R_\phi] \,|\, t \le 6\}$, $n = 2$. Therefore, the infimum area ratio among all $\log 1.8$-coverings of $\{[T_t R_\phi] \,|\, t \le 6\}$ is lower than 6.34.*

*Proof.* We are dealing with 4 dimensions to minimize and more specifically with $100^4$ feasible sets. Computing area ratios for each feasible set is straightforward but validating the covering condition is a more involved computational issue. For the sake of clearness, the intersection of disk boundaries, which are composed at most of two elements for nonidentical disks, shall be denoted by

$$\Sigma_1 = \partial \mathcal{B}_{[T_{t_1}]}^{\log 1.8} \cap \partial \mathcal{B}_{[T_{t_1 R_{\phi_1}}]}^{\log 1.8}, \qquad\qquad \Sigma_2 = \partial \mathcal{B}_{[T_{t_2}]}^{\log 1.8} \cap \partial \mathcal{B}_{[T_{t_2 R_{\phi_2}}]}^{\log 1.8},$$

and their respective closest and farthest elements will be denoted by

$$\min \Sigma_1 := \arg \min_{S \in \Sigma_1} d\left(S, [Id]\right), \qquad\qquad \max \Sigma_1 := \arg \max_{S \in \Sigma_1} d\left(S, [Id]\right),$$
$$\min \Sigma_2 := \arg \min_{S \in \Sigma_2} d\left(S, [Id]\right), \qquad\qquad \max \Sigma_2 := \arg \max_{S \in \Sigma_2} d\left(S, [Id]\right).$$

In order to check if a feasible set does cover the specified region we propose to verify the following four conditions depicted in Figure 8:

1. $\Sigma_1 \neq \emptyset$ and $\Sigma_2 \neq \emptyset$.
2. $\min \Sigma_1$ must lie inside the ball $\mathcal{B}_{[Id]}^{\log 1.8}$, which ensures a covering of $\mathcal{B}_{[Id]}^{\log \tau(\max \Sigma_1)}$.
3. $\max \Sigma_2$ must lie outside the region $\Gamma$, which ensures a covering of the annulus defined by $\Gamma \setminus \mathcal{B}_{[Id]}^{\log \tau(\min \Sigma_2)}$.
4. For $\varepsilon$ small, all elements $S \in \mathbb{F}_\varepsilon$ must lie inside some disks of radius $\log\left(1.8 - \varepsilon\right)$, i.e.,

$$S \in \bigcup_{1 \leq i \leq 2} \bigcup_{S' \in J_{t_i, \phi_i}} \mathcal{B}_{[S']}^{\log(1.8 - \varepsilon)},$$

where $\mathbb{F}_\varepsilon$ is a finite $\varepsilon$-dense set of the annulus defined by

$$\mathcal{B}_{[Id]}^{\log \tau(\min \Sigma_2)} \setminus \mathcal{B}_{[Id]}^{\log \tau(\max \Sigma_1)}.$$

Notice that the fourth condition only ensures a $\log\left(1.8 - \varepsilon\right)$-covering up to an error

$$\varepsilon = \max_{S' \in \Gamma} \min_{S \in \mathbb{F}_\varepsilon} d\left(S, S'\right)$$

and so, by dilating back disks radii to 1.8, one ensures $\log 1.8$-coverings.

By using the procedure described above, an approximated solution to the optimization problem in (11) has been obtained. Its parameters can be found in Table 2. Its corresponding representation in the space of tilts appears in Figure 9(c). ∎

The procedure in the proof of Proposition 3.12 has also been applied to find more near optimal coverings appearing in Figure 9.

**4. Experimental validation.** We are now able to propose and evaluate for each SIIM method its IMAS, namely, its affine invariant extension. This affine invariant version relies on two facts. First, each SIIM identifies viewpoint changes, under a certain transition tilt threshold (that we shall estimate in this section). Second, any smooth map is locally approximable by an affine map. Hence, under the assumption that the surface of photographed
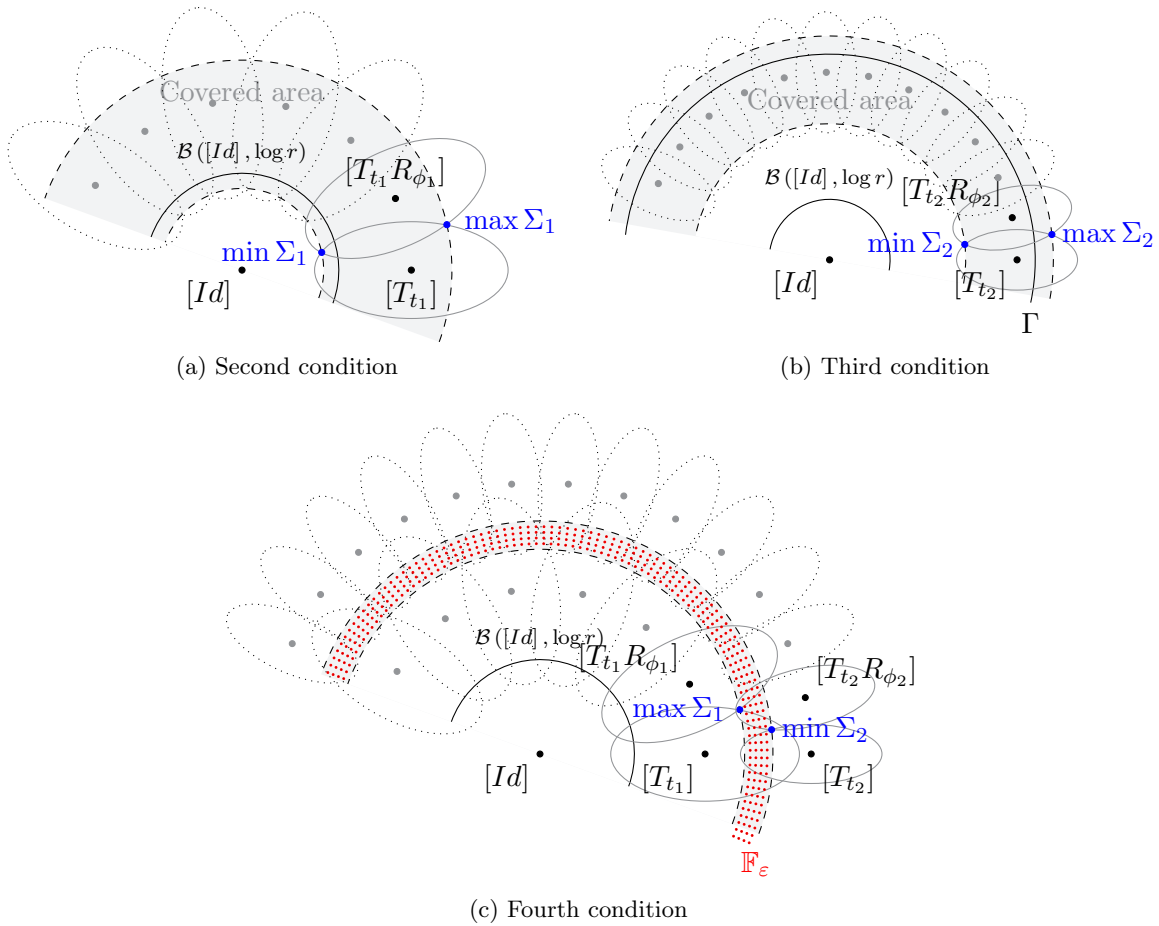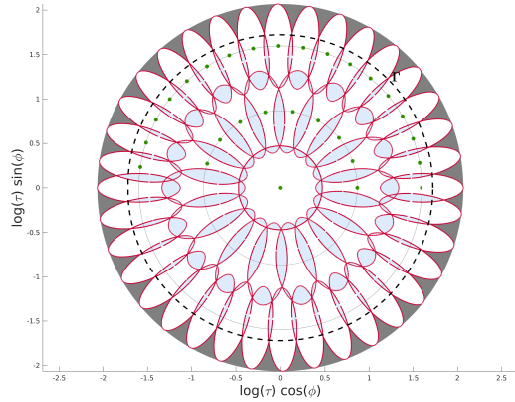
(a) Second condition

(b) Third condition

(c) Fourth condition

**Figure 8.** *Verifying covering conditions for feasible sets in Proposition* 3.12.

**Table 2**
*Approximated solution to the optimization problem in* (11).

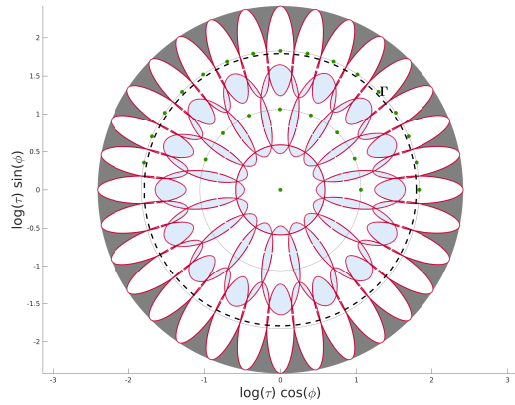| Parameter | Value |
|-----------|-------|
| $t_1^{opt}$ | 2.88447 |
| $\phi_1^{opt}$ | 0.394085 |
| $t_2^{opt}$ | 6.2197 |
| $\phi_2^{opt}$ | 0.196389 |

objects is locally smooth, all viewpoint changes can be understood as local transition tilt changes (see Figure 1). Third, once provided with a $\log r$-covering of $\Gamma = \Gamma'$, where $r$ is less than the transition tilt threshold of the SIIM, Proposition 3.6 states that Algorithm 3.1 offers an affine invariant version of the considered SIIM. Indeed, there is at least one pair of simulated images whose transition tilt is less than $r$, and on these two images the SIIM can succeed. The affine invariance property is ensured for transition tilt changes up to $\Lambda_1 \Lambda_2$, i.e., for viewpoint angle changes of about $\arccos(\frac{1}{\Lambda_1 \Lambda_2})$. We shall denote by $t_{\max}^{s_1 \times s_2}$ the associated maximum tilt tolerance with respect to a matching method for images with size larger than $s_1 \times s_2$.
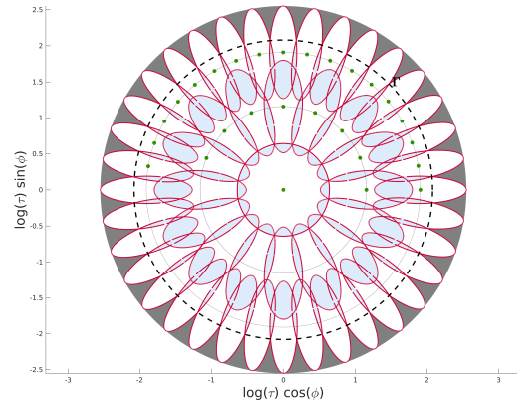
(a) Optimal $\log 1.6$-covering of $\{[T_t R_\phi] \mid t \leq 5.6\}$ with 28 affine simulations representing an area ratio of 8.42.

(b) Optimal $\log 1.7$-covering of $\{[T_t R_\phi] \mid t \leq 5.8\}$ with 25 affine simulations representing an area ratio of 7.06.

(c) Optimal $\log 1.8$-covering of $\{[T_t R_\phi] \mid t \leq 6\}$ with 25 affine simulations representing an area ratio of 6.34.

(d) Optimal $\log 1.9$-covering of $\{[T_t R_\phi] \mid t \leq 8\}$ with 27 affine simulations representing an area ratio of 6.18.

(e) Optimal $\log 2$-covering of $\{[T_t R_\phi] \mid t \leq 10\}$ with 32 affine simulations representing an area ratio of 6.02.

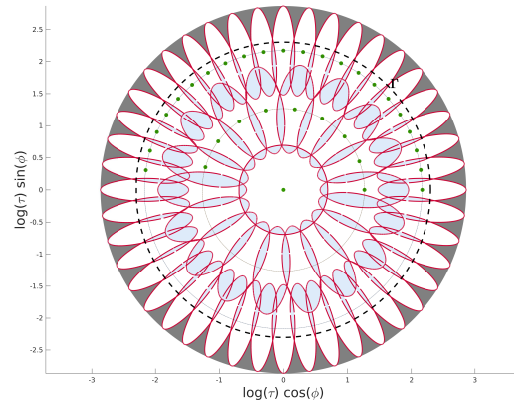**Figure 9.** *Near-optimal coverings in the space of tilts. Gray areas—uncovered; blue areas—covered by at least two disks; white areas—covered by only one disk.*

In our experiments, all SIIM methods were immersed in the same affine extension setup. The simulation of optical tilts, matching, and filtering were handled in the very same way. This setup received as a parameter the name of the base detector+extractor method to perform, then a brute force matcher was performed with the second-closest neighbor acceptance criterion proposed by Lowe in [33]. Finally, as presented in [44, 64], three main filters were applied: first, only unique matches were taken into account; second, groups of multiple-to-one and one-to-multiple matches were removed; finally, only matches coming from the most significant geometric model (if it existed!) were kept. In our case, as all tests were based on planar transformations, the ORSA homography detector [41] (a parameterless variant of RANSAC) was applied to filter out matches not compatible with the dominant homography.

All detectors, all extractors, and the matcher were taken from the Open Source Computer Vision (OPENCV) Library, version 3.2.0.

**4.1. Maximal tilt tolerance computation for each SIIM.** From the complexity viewpoint, the main quantitative parameter for extending a SIIM into an IMAS is its tilt tolerance. We do not question the invariance of descriptors with respect to zoom and rotations but rather how they perform against transition tilt changes incurred when matching, for example, $\mathbb{G}_1 Id\, u$ to $\mathbb{G}_1 T_t R_\phi u$, where $t \in [1, \infty[$ and $\phi \in [0, \pi[$.

We used the *tolerance image dataset* displayed in Figure 10 to evaluate the maximal tilt tolerance of each SIIM with respect to images of similar size. Images in this dataset have a fixed size and were selected to obtain a diversity of challenging scenarios. In order



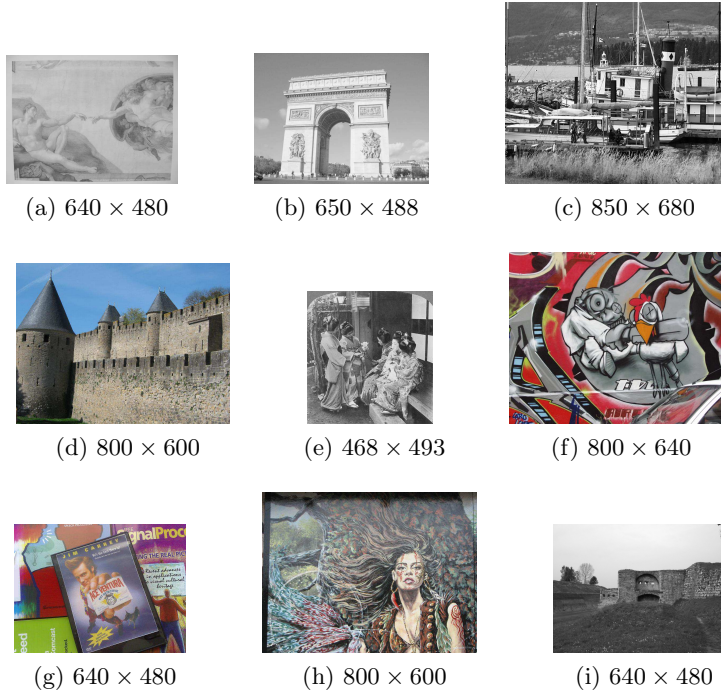(a) $640 \times 480$     (b) $650 \times 488$     (c) $850 \times 680$

(d) $800 \times 600$     (e) $468 \times 493$     (f) $800 \times 640$

(g) $640 \times 480$     (h) $800 \times 600$     (i) $640 \times 480$

**Figure 10.** *Tolerance image dataset.*

| (a) RootSIFT $\left(U_{\max}^{700\times550} = 2\right)$ | (b) SIFT $\left(U_{\max}^{700\times550} = 1.8\right)$ | (c) FREAK $\left(U_{\max}^{700\times550} = 1.8\right)$ |

| (d) AKAZE $\left(U_{\max}^{700\times550} = 1.7\right)$ | (e) BRISK $\left(U_{\max}^{700\times550} = 1.7\right)$ | (f) ORB $\left(U_{\max}^{700\times550} = 1.5\right)$ |

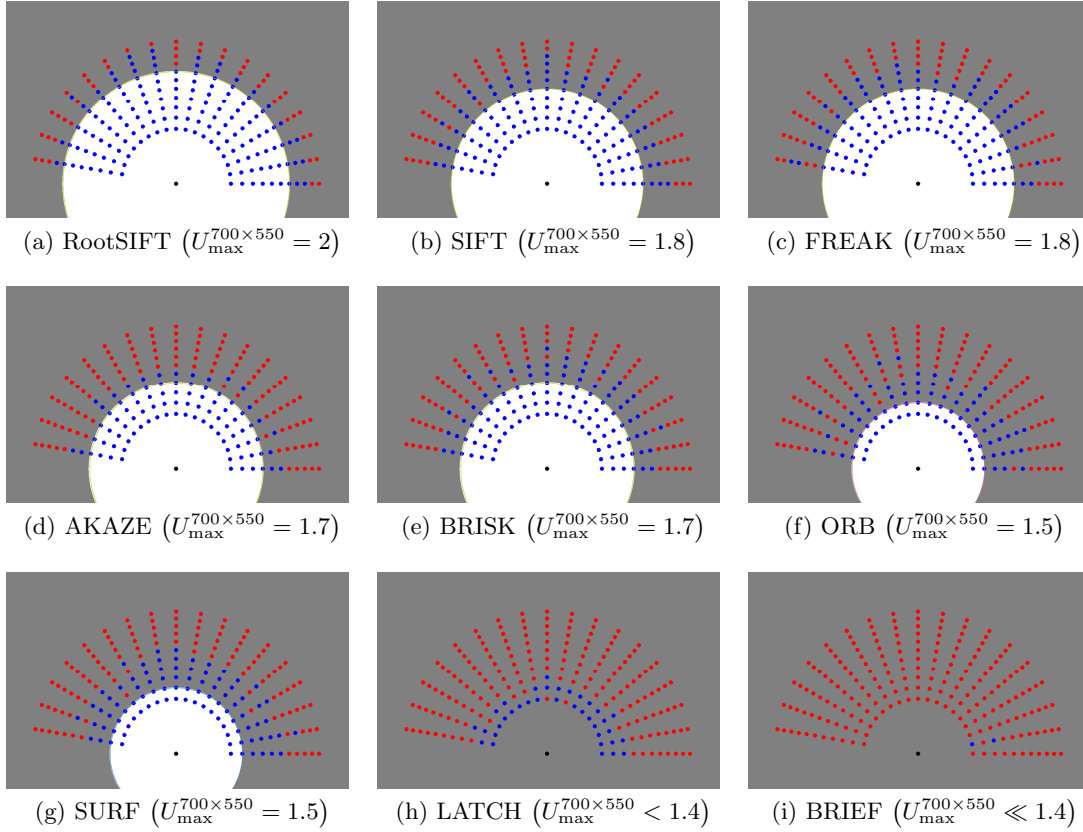| (g) SURF $\left(U_{\max}^{700\times550} = 1.5\right)$ | (h) LATCH $\left(U_{\max}^{700\times550} < 1.4\right)$ | (i) BRIEF $\left(U_{\max}^{700\times550} \ll 1.4\right)$ |

**Figure 11.** *Represented in the space of tilts, the associated upper bounds $(U_{\max}^{700\times550})$ for maximum tilt tolerances. Black dot—[Id]; colored dots stand for tested tilts $[T_t R_\phi]$, where $t \in \{1.4, 1.5, \ldots, 2.4\}$ and $\phi \in \{0, 10, \ldots, 170\}$; blue dots—attainable tilts for all images in the dataset; red dots—unattainable tilts for at least one image in the dataset; gray areas—$\left\{[T_t R_\phi] \,|\, t \geq U_{\max}^{700\times550}\right\}$; white areas—$\left\{[T_t R_\phi] \,|\, t \leq U_{\max}^{700\times550}\right\}$.*

to approximate $t_{\max}^{700\times550}$, we simulated optical tilts on the tolerance image dataset and then tested whether this affine simulation was identified by ORSA homography with a precision of 3 pixels. This test determined upper bounds $U_{\max}^{700\times550}$ depicted in Figure 11 for nine of the best state-of-the-art SIIMs.

This test yielded upper bounds for $t_{\max}^{700\times550}$, based on its application to nine images whose sizes are close to $700 \times 550$. Supposing a maximal angle error computation of $\frac{\pi}{10}$, we assumed that for each SIIM

$$t_{\max}^{700\times550} = \frac{U_{\max}^{700\times550}}{\frac{1}{\left|\cos\left(\frac{\pi}{10}\right)\right|}} \approx \frac{U_{\max}^{700\times550}}{1.05},$$

and constructed its affine invariant version with $\log t_{\max}^{700\times550}$-coverings.

**4.2. Affine invariant methods.** The matching process is as symmetric as possible. No significant changes should come along by interchanging the roles of the query and target images. In the case of IMAS algorithms this symmetry implies a unique set of optical tilts to
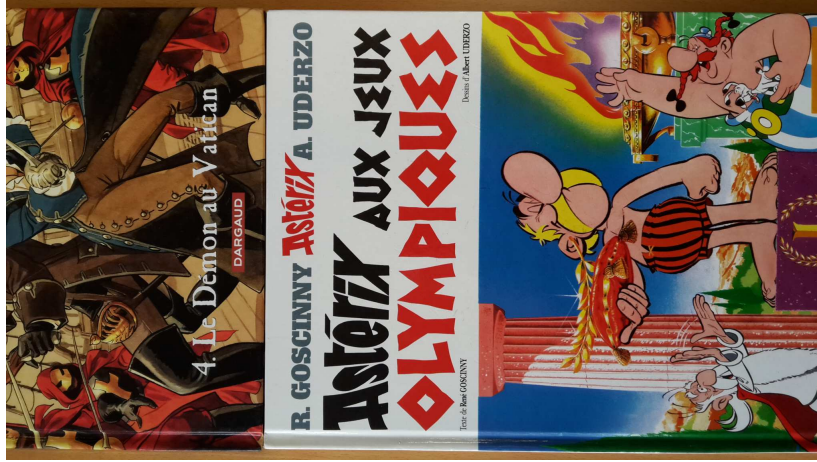
**Figure 12.** *Image $w_0$ ($3264 \times 1836$) for the IMAS efficiency test.*

simulate on both query and target images. Thus, if this unique set of optical tilts represents a $\log r$-covering of

$$\Gamma_1 = \Gamma' = \{[T_t R_\phi] \mid t \leq \Lambda\}$$

then Proposition 3.6 ensures that any IMAS based on a SIIM whose maximum tilt tolerance is greater than $r$ is able to identify all tilts under $\frac{\Lambda^2}{r}$ by simulating all affine maps in the $\log r$-covering.

Several coverings in the space of tilts have been proposed in [44, 64, 49, 40] for SIFT and SURF. Figure 16 displays these coverings. They are clearly not optimal. Indeed, most of these coverings do not really cover the region they were meant to, except for ASIFT [44, 64] (which instead is visually redundant) and for the affine DoG-SIFT version in [40].

In order to compare the efficiency of those coverings, query and target images were generated in a way so as to test Algorithm 3.1 to the limit, i.e., forcing the worst case scenario in which $\left[(BC)^{-1}\right]$ lies in $\Gamma' \setminus \Gamma_2$. We simulated the optical tilts on query and target images coming from one single image. This image, denoted by $w_0$ and appearing in Figure 12, was then used to compute the inputs of Algorithm 3.1 as follows:

- Query image (nonfixed tilt) $\mathbb{G}_1 A_{t,\phi} w_0$, where $A_{t,\phi} = R_\phi T_t R_{\frac{\pi}{2}}$.
- Target image (fixed tilt) $\mathbb{G}_1 B_\phi w_0$, where $B_\phi = R_{\phi+\frac{\pi}{2}} T_\Lambda$.

The veritable interest of these affine maps being the inverse maps they determine, namely,

$$\left[A_{t,\phi}^{-1}\right] = \left[T_t R_{\frac{\pi}{2}-\phi}\right],$$
$$\left[B_\phi^{-1}\right] = [T_\Lambda R_\phi],$$

which according to Proposition 2.9.4, attain maximal transition tilts for fixed tilts such as $t$ and $\Lambda$, i.e.,

$$\tau\left(A_{t,\phi}^{-1} B_\phi\right) = t\Lambda.$$

(a) Optimal Affine-SIFT ($r = 1.7$)
$\Gamma_1 = \{[T_t R_\phi] \,|\, t \leq 3.41\}$
$\Gamma' = \{[T_t R_\phi] \,|\, t \leq 5.8\}$

(b) ASIFT ($r = 1.8$)
$\Gamma_1 = \{[T_t R_\phi] \,|\, t \leq 3.05\}$
$\Gamma' = \{[T_t R_\phi] \,|\, t \leq 5.5\}$

(c) MEDIUM configuration for DoG-SIFT
($r = 1.8$)
$\Gamma_1 = \{[T_t R_\phi] \,|\, t \leq 5\}$
$\Gamma' = \{[T_t R_\phi] \,|\, t \leq 9\}$

(d) Optimal Affine-SURF ($r = 1.4$)
$\Gamma_1 = \{[T_t R_\phi] \,|\, t \leq 3.57\}$
$\Gamma' = \{[T_t R_\phi] \,|\, t \leq 5\}$

(e) FAIR-SURF - simulated tilts ($r = 1.5$)
$\Gamma_1 = \{[T_t R_\phi] \,|\, t \leq 3.77\}$
$\Gamma' = \{[T_t R_\phi] \,|\, t \leq 4\sqrt{2}\}$

(f) FAIR-SURF - fixed tilts ($r = 1.5$)
$\Gamma_1 = \{[T_t R_\phi] \,|\, t \leq 3.77\}$
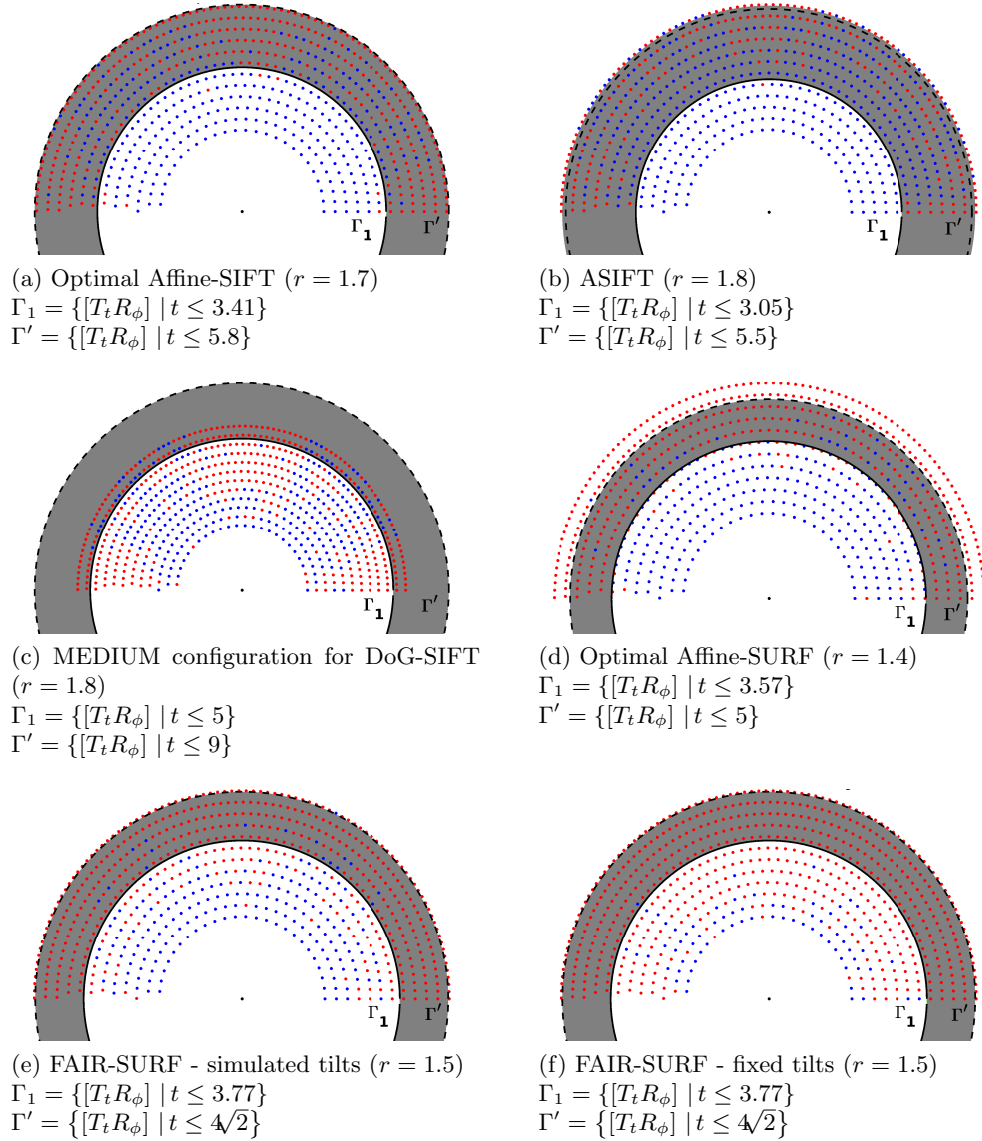$\Gamma' = \{[T_t R_\phi] \,|\, t \leq 4\sqrt{2}\}$

**Figure 13.** *Extreme test results. Black dot—$[Id]$; colored dots stand for $[A_{t,\phi}^{-1}]$ and belong to a fixed $\log 1.1$ uniform discretization of the annulus $\{[T_t R_\phi] \,|\, 2 \leq t \leq 4\sqrt{2}\}$. The angle $\phi$ implicitly fixes $\left[B_\phi^{-1}\right] = [T_\Lambda R_\phi]$, where $\Lambda = \arg\max_t [T_t R_\phi] \in \Gamma'$. Blue/Red dots—Success/Failure of ORSA homography in identifying the underlying affine map.*

When ORSA homography was able to identify the affine map that relates query and target images, we counted the event as a success. Clearly, if $\Gamma'$ and $\Gamma_2$ are truly $\log r$-covered then Proposition 3.6 implies that all tests for which $[A_{t,\phi}^{-1}] \in \Gamma_1$ should be counted as a success. Results in Figure 13 were as expected and highlight the importance of using the right coverings for extreme cases. Both ASIFT and optimal affine-SIFT were able to capture most of all transition tilts that Proposition 3.6 predicted, namely, those under $\frac{\Lambda^2}{r}$.

We must keep in mind that these log $r$-coverings depend on tilt tolerances found over images in Figure 10. Maximal tilt tolerances are linked to the size of images being compared and as a consequence the disks radius might grow or shrink proportionally to the minimum size of all simulated images. Moreover, Proposition 3.6 does not take into account discretization errors and relies on two main hypotheses:

1. The considered SIIM is truly rotation and zoom invariant.
2. For images similar to the input image, the SIIM under consideration has a maximal tilt tolerance not smaller than $r$.

As anticipated, the area ratio associated with a covering reliably evaluates the difference of performance between affine versions of the same matching method. Being proportionally linked to the total amount of keypoints, the area ratio of Definition 3.8 predicts the order of growth in computation time. For example, the SIFT keypoint computation part induced by the optimal covering in Figure 9(b) is twice as fast as the one induced by the ASIFT covering. The same goes for the matching part, only this time the optimal version is four times faster. Since both coverings cover about the same region, our optimal affine-SIFT supplants ASIFT with no qualitative matching loss.

Two examples of performance over query and target images from Figures 14 and 15 are, respectively, found in Tables 3 and 4. In Table 3, affine-ORB and affine-BRIEF both fail because of too many false matches. The best scores found by ORSA to identify meaningful homographies were, respectively, 16 out of 905 and 6 out of 1409. Code optimization, smart tweaks, and parallelism performance may vary from SIIM to SIIM and from IMAS to IMAS, which ultimately may lead to discrepant area ratio predictions on computation time. This is the case of SURF (and optimal affine-SURF) whose implementation uses several fine and clever optimizations. Nonetheless, the optimal affine-SIFT yields more matches for a lower computation time.

In Table 4 the reader will notice that affine-ORB has fewer matches than ORB itself, which might seem contradictory. This happens when postprocessing the matches, more specifically, when applying the second filter. The *multiple-to-one/one-to-multiple* filter, initially proposed in [44, 64], is meant to filter out undesired aberrant matches but, unfortunately, many good ones also get eliminated. In spite of this handicap, affine-ORB is able to catch more matches with higher transition tilts.

**5. Conclusion.** IMAS is acknowledged as the best methodology to match images of the same scene regardless of the viewpoint change. Its time complexity is one of the main drawbacks that has been widely criticized in the literature. The mathematical derivations in this paper imply that IMAS based methods really are affine invariant provided the base SIIM satisfies scale+rotation invariance, sufficient distinctiveness, and an acceptable viewpoint tolerance measured as its transition tilt. We have proved that, as summarized in Figure 16, all former IMAS methods are oversimulating optical tilts. We therefore have developed a method, finding for each SIIM an optimal IMAS method which only depends on the tilt tolerance of the SIIM. This led us to measure the tilt tolerance of a number of classic SIIMs. We found, for example, that the optimal IMAS extension of SIFT needs half as many descriptors and therefore is four times faster than ASIFT. This improvement applies to all state-of-the-art IMAS, that can be accelerated by a factor of four. Another consequence is that the set of affine descriptors associated with an image can be halved.

**Table 3**

*Matching methods performance over query and target images from Figure* 14. *The proposed matching methods in this paper appear in bold. Computations were performed on an Intel(R) Core(TM) i5-4210U CPU 1.70 GHz with 2 cores. M—Matches; ar—area ratio.*

|  | M | $ar$ | $ar^2$ | Keypoints (seconds) | Matching (seconds) | Filters (seconds) |
|---|---|---|---|---|---|---|
| SIFT | 0 | 1 | 1 | 0.69 | 0.70 | 0.18 |
| ASIFT | 1013 | 13.7 | 189.6 | 12.46 | 138.59 | 3.05 |
| **(Optimal) Affine-SIFT** | **795** | **7.06** | **49.8** | **6.04** | **29.61** | **1.39** |
| RootSIFT | 0 | 1 | 1 | 0.72 | 0.71 | 0.18 |
| **Affine-RootSIFT** | **658** | **6.9** | **47.6** | **5.05** | **20.70** | **1.44** |
| SURF | 0 | 1 | 1 | 1.01 | 0.79 | 0.19 |
| **(Optimal) Affine-SURF** | **471** | **14.82** | **219,6** | **12.53** | **35.24** | **1.40** |
| BRISK | 0 | 1 | 1 | 1.75 | 0.27 | 0.18 |
| **Affine-BRISK** | **421** | **8.42** | **70,89** | **18.95** | **8.68** | **2.06** |
| BRIEF | 0 | 1 | 1 | 0.05 | 0.01 | 0.19 |
| **Affine-BRIEF** | **0** | **14.82** | **219,6** | **4.20** | **2.18** | **6.08** |
| ORB | 0 | 1 | 1 | 0.05 | 0.02 | 0.17 |
| **Affine-ORB** | **0** | **14.82** | **219,6** | **4.34** | **5.13** | **3.25** |
| AKAZE | 0 | 1 | 1 | 0.42 | 0.13 | 0.21 |
| **Affine-AKAZE** | **194** | **8.42** | **70,89** | **5.00** | **6.23** | **3.74** |
| LATCH | 0 | 1 | 1 | 0.11 | 0.02 | 0.00 |
| **Affine-LATCH** | **37** | **14.82** | **219,6** | **4.52** | **2.16** | **0.17** |
| FREAK | 0 | 1 | 1 | 0.34 | 0.15 | 0.18 |
| **Affine-FREAK** | **145** | **7.06** | **49.8** | **4.37** | **2.38** | **1.94** |

**Table 4**

*Matching methods performance over query and target images from Figure* 15. *The proposed IMAS methods proposed here appear in bold. Computations were performed on an Intel(R) Core(TM) i5-4210U CPU 1.70 GHz with 2 cores. M—matches; ar—area ratio.*

|  | M | $ar$ | $ar^2$ | Keypoints (seconds) | Matching (seconds) | Filters (seconds) |
|---|---|---|---|---|---|---|
| SIFT | 102 | 1 | 1 | 0.23 | 0.01 | 0.09 |
| ASIFT | 317 | 13.7 | 189.6 | 5.43 | 1.68 | 0.47 |
| **(Optimal) Affine-SIFT** | **292** | **7.06** | **49.8** | **2.71** | **0.38** | **0.30** |
| RootSIFT | 110 | 1 | 1 | 0.25 | 0.01 | 0.09 |
| **Affine-RootSIFT** | **219** | **6.9** | **47.6** | **2.23** | **0.28** | **0.24** |
| SURF | 110 | 1 | 1 | 0.24 | 0.03 | 0.14 |
| **(Optimal) Affine-SURF** | **663** | **14.82** | **219,6** | **3.68** | **1.19** | **0.73** |
| BRISK | 29 | 1 | 1 | 1.57 | 0.00 | 0.04 |
| **Affine-BRISK** | **49** | **8.42** | **70,89** | **17.57** | **0.06** | **0.08** |
| BRIEF | 0 | 1 | 1 | 0.03 | 0.00 | 0.00 |
| **Affine-BRIEF** | **7** | **14.82** | **219,6** | **2.06** | **0.09** | **0.03** |
| ORB | 102 | 1 | 1 | 0.02 | 0.01 | 0.8 |
| **Affine-ORB** | **90** | **14.82** | **219,6** | **2.12** | **0.31** | **0.40** |
| AKAZE | 20 | 1 | 1 | 0.16 | 0.00 | 0.03 |
| **Affine-AKAZE** | **51** | **8.42** | **70,89** | **2.31** | **0.06** | **0.09** |
| LATCH | 54 | 1 | 1 | 0.07 | 0.01 | 0.04 |
| **Affine-LATCH** | **101** | **14.82** | **219,6** | **1.72** | **0.12** | **0.10** |
| FREAK | 124 | 1 | 1 | 0.14 | 0.01 | 0.10 |
| **Affine-FREAK** | **182** | **7.06** | **49.8** | **2.54** | **0.11** | **0.31** |

(a) $800 \times 640$            (b) $800 \times 640$

**Figure 14.** *Graffiti. Both images generate a large number of keypoints for most methods.*



(a) $600 \times 450$            (b) $600 \times 450$

**Figure 15.** *Adam. Both images generate a small number of keypoints for most methods.*

## 6. Appendix.

### 6.1. Proof of Theorem 2.16. By proposition 2.14 we know that
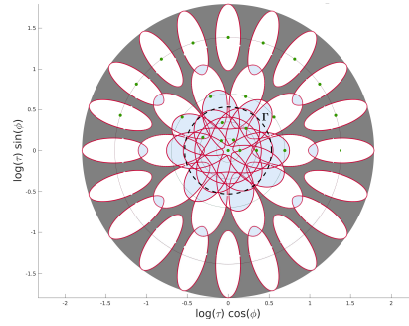
$$\tau\left(BA^{-1}\right) = \tau\left(i\left([B]\right) i\left([A]\right)^{-1}\right),$$

where $i$ is the injection in Definition 2.12. Thus, without loss of generality, we focus on computing the absolute tilt of

$$
\begin{aligned}
C &= T_t R_2 Q_2^{-1} T_s^{-1} \\
&= T_t R\left(\phi\right) T_s^{-1},
\end{aligned}
$$

where $R\left(\phi\right) = R_2 Q_2^{-1}$. Proposition 2.5 states that the ratio between the singular values of $C$ can be used to compute its absolute tilt.
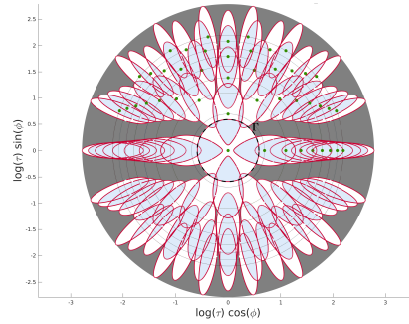
(a) Proposed covering for ASIFT in [44, 64]. This is a $\log 1.8$-covering of $\{[T_t R_\phi] \mid t \leq 5.5\}$ with 41 affine simulations representing an area ratio of 13.77.
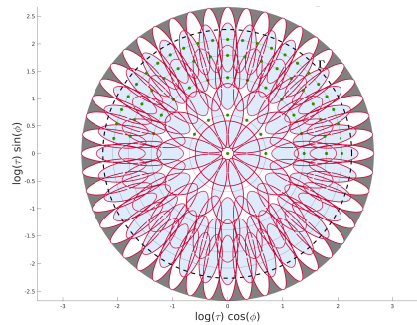
(b) Proposed covering for FAIR-SURF in [49], called fixed tilts. This is a $\log 1.5$-covering of $\{[T_t R_\phi] \mid t \leq 1.7\}$ with 23 affine simulations representing an area ratio of 11.42.
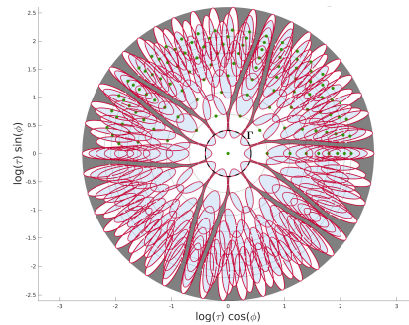
(c) Proposed covering for FAIR-SURF in [49], called simulated tilts. This is a $\log 1.5$-covering of $\{[T_t R_\phi] \mid t \leq 1.65\}$ with 41 affine simulations representing an area ratio of 13.77.

(d) Proposed covering in [40], called mediium configuration for DoG-SIFT. This is a $\log 1.8$-covering of $\{[T_t R_\phi] \mid t \leq 1.8\}$ with 45 affine simulations representing an area ratio of 9.

(e) Proposed covering in [40], called hard configuration for DoG-SIFT. This is a $\log 1.8$-covering of $\{[T_t R_\phi] \mid t \leq 9.6\}$ with 61 affine simulations representing an area ratio of 13.

(f) Proposed covering in [40], called hard configuration for SURF-SURF. This is a $\log 1.5$-covering of $\{[T_t R_\phi] \mid t \leq 1.5\}$ with 112 affine simulations representing an area ratio of 21.28.

**Figure 16.** *Examples of coverings found in the literature for maximum tilt tolerances as in Figure 11. Gray areas—uncovered; blue areas—covered by at least two disks; white areas—covered by only one disk.*

**6.1.1. Trace and determinant.** First, we start by computing the trace and determinant of

$$C^\star C = T_s^{-1} R\left(\phi\right)^{-1} T_t T_t R\left(\phi\right) T_s^{-1},$$

which are clearly

$$\det\left(C^\star C\right) = \frac{t^2}{s^2}$$

and

$$Tr\left(C^\star C\right) = \left(\frac{t^2}{s^2} + 1\right)\cos^2\phi + \left(\frac{1}{s^2} + t^2\right)\sin^2\phi.$$

**6.1.2. The eigenvalues of $C^\star C$.** Let $H = \left(\begin{smallmatrix} a & c \\ c & b \end{smallmatrix}\right) = C^\star C$ and $\lambda_+, \lambda_-$ being the biggest and smallest eigenvalues of $C^\star C$, respectively. It is well known that

$$Tr\left(H\right) = \lambda_+ + \lambda_-,$$
$$\det\left(H\right) = \lambda_+\lambda_-,$$

and even more that both $Tr$ and det also appear in the characteristic polynomial

$$|C^\star C - \lambda Id| = \lambda^2 - \lambda\left(a + b\right) + \left(ab - c^2\right),$$
$$= \lambda^2 - \lambda Tr\,H + \det H.$$

On the other hand, the eigenvalues of a symmetric positive definite matrix are in $\mathbb{R}$, which implies that $\sqrt{(Tr\,H)^2 - 4\det H} \geq 0$, and so one can write

$$\lambda_- = \frac{Tr\left(H\right) - \sqrt{\left(Tr\,H\right)^2 - 4\det H}}{2},$$
$$\lambda_+ = \frac{Tr\left(H\right) + \sqrt{\left(Tr\,H\right)^2 - 4\det H}}{2}.$$

Now, after some computations, the ratio between the biggest and smallest eigenvalues is

$$\frac{\lambda_+}{\lambda_-} = \frac{\left(\frac{Tr\,H}{2} + \frac{\sqrt{\left(Tr\,H\right)^2 - 4\det H}}{2}\right)^2}{\det H},$$

(12)
$$= \frac{s^2}{t^2}\left(\frac{g}{2} + \frac{\sqrt{g^2 - 4\frac{t^2}{s^2}}}{2}\right)^2,$$

where $g$ denotes the function

$$g\left(t, s, \phi\right) := Tr\left(C^{\star}C\right)$$
$$= \left(\frac{t^2}{s^2} + 1\right)\cos^2\phi + \left(\frac{1}{s^2} + t^2\right)\sin^2\phi.$$

**6.1.3. Computing $\tau\left(C\right)$.** Proposition 2.5 tells us that the absolute tilt of $C$ is

$$\tau\left(C\right) = \sqrt{\frac{\lambda_+}{\lambda_-}}$$
$$= \frac{s}{t}\left(\frac{g}{2} + \frac{\sqrt{g^2 - 4\frac{t^2}{s^2}}}{2}\right)$$
$$= \frac{s}{t}\frac{g}{2} + \sqrt{\left(\frac{s}{t}\frac{g}{2}\right)^2 - 1}$$
$$= G\left(s, t, \phi\right) + \sqrt{\left(G\left(s, t, \phi\right)\right)^2 - 1},$$

where

$$G\left(s, t, \phi\right) = \frac{s}{t}\frac{g\left(s, t, \phi\right)}{2}.$$

**6.1.4. Disks in the space of tilts.** Let $\boldsymbol{A} := [T_t R_2] \in \Omega$ be fixed and let us find conditions on $\boldsymbol{B} := [T_s Q_2] \in \Omega$ to satisfy

$$\boldsymbol{B} \in \mathcal{B}\left(\boldsymbol{A}, \log r\right)$$

which are clearly

$$d\left(\boldsymbol{A}, \boldsymbol{B}\right) = \log\tau\left(i\left(\boldsymbol{A}\right)i\left(\boldsymbol{B}\right)^{-1}\right) \leq \log r,$$
$$\Updownarrow$$
$$\tau\left(i\left(\boldsymbol{A}\right)i\left(\boldsymbol{B}\right)^{-1}\right) \leq r,$$

where $i$ is the injection in Definition 2.12. Thus, just by applying the above to $C := i\left(\boldsymbol{A}\right)i\left(\boldsymbol{B}\right)^{-1}$ we obtained

$$G\left(s, t, \phi\right) + \sqrt{\left(G\left(s, t, \phi\right)\right)^2 - 1} = \tau\left(AB^{-1}\right)$$
$$\leq r,$$

where $R\left(\phi\right) = R_2 Q_2^{-1}$. So

$$\sqrt{G^2 - 1} \leq r - G$$
$$\Updownarrow$$
$$G^2 - 1 \leq r^2 - 2rG + G^2$$
$$\Updownarrow$$
$$G \leq \frac{r^2 + 1}{2r}.$$

## REFERENCES

[1] S. AGARWAL, Y. FURUKAWA, N. SNAVELY, I. SIMON, B. CURLESS, S. M. SEITZ, AND R. SZELISKI, *Building Rome in a day*, Commun. 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*, Graph. ACM Trans., 25 (2006), pp. 853–861.

[3] P. F. ALCANTARILLA, A. BARTOLI, AND A. J. DAVISON, *KAZE features*, in European Conference on Computer Vision, Lecture Notes in Comput. Sci. 7577, 2012, pp. 214–227, https://doi.org/10.1007/978-3-642-33783-3_16.

[4] P. F. ALCANTARILLA, J. NUEVO, AND A. BARTOLI, *Fast explicit diffusion for accelerated features in nonlinear scale spaces*, British Machine Vision Conference, British Machine Vision Association, Durham, England, 2013.

[5] R. ARANDJELOVIC AND A. ZISSERMAN, *Three things everyone should know to improve object retrieval*, in Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, IEEE, Piscataway, NJ, 2012, pp. 2911–2918, https://doi.org/10.1109/CVPR.2012.6248018.

[6] A. BAUMBERG, *Reliable feature matching across widely separated views*, Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE Computer Society, Los Alamitos, CA, 2000, pp. 774–781.

[7] H. BAY, T. TUYTELAARS, AND L. VAN GOOL, *SURF: Speeded up robust features*, European Conference on Computer Vision, Springer, Berlin, 2006, pp. 404–417.

[8] J. BLOM, *Topological and Geometrical Aspects of Image Structure.*, Ph.D. Thesis, Utrecht, University of Utrecht, The Netherlands, 1992.

[9] M. BROWN AND D. LOWE, *Recognising panoramas*, in Proceedings the 9th International Conference on Computer Vision, IEEE Compter Society, Los Alamitos, CA, 2003, pp. 1218–1225.

[10] M. BROWN AND S. SÜSSTRUNK, *Multi-spectral SIFT for scene category recognition*, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2011, IEEE, Piscataway, NJ, 2011, pp. 177–184.

[11] M. CALONDER, V. LEPETIT, C. STRECHA, AND P. FUA, *BRIEF: Binary robust independent elementary features*, in European Conference on Computer Vision, Lecture Notes in Comput. Sci. 6314, Springer, Berlin, 2010, pp. 778–792, https://doi.org/10.1007/978-3-642-15561-1_56.

[12] F. CAO, J.-L. LISANI, J.-M. MOREL, P. MUSÉ, AND F. SUR, *A Theory of Shape Identification*, Springer, Berlin, 2008.

[13] D. COZZOLINO, G. POGGI, AND L. VERDOLIVA, *Efficient dense-field copy–move forgery detection*, IEEE Trans. Infor. Forensic. Secur., 10 (2015), pp. 2284–2297.

[14] O. FAUGERAS, *Three-Dimensional Computer Vision: A Geometric Viewpoint*, MIT Press, Cambridge, MA, 1993.

[15] G. FRITZ, C. SEIFERT, M. KUMAR, AND L. PALETTA, *Building detection from mobile imagery using informative SIFT descriptors*, in Scandinavian Conference on Image Analysis, Lecture Notes in Comput. Sci. 3540, 2005, pp. 629–638.

[16] A. GEIGER, J. ZIEGLER, AND C. STILLER, *Stereoscan: Dense 3d reconstruction in real-time*, in Intelligent Vehicles Symposium (IV), 2011, IEEE, Piscataway, NJ, 2011, pp. 963–968.

[17] Y. GONG, S. LAZEBNIK, A. GORDO, AND F. PERRONNIN, *Iterative quantization: A Procrustean approach to learning binary codes for large-scale image retrieval*, IEEE Trans. Pattern Anal. Mach. Intell., 35 (2013), pp. 2916–2929.

[18] J. S. HARE AND P. H. LEWIS, *Salient regions for query by image content*, Image and Video Retrieval: Third International Conference, CIVR, Springer, Berlin, 2004, pp. 317–325.

[19] C. HARRIS AND M. STEPHENS, *A combined corner and edge detector*, Alvey Vision Conference, Springer, Berlin, University of Manchester, Manchester, England, 1988, p. 50.

[20] T. IIJIMA, *Basic equation of figure and and observational transformation*, Syst. Comput. Controls, 2 (1971), pp. 70–77.

[21] T. KADIR, A. ZISSERMAN, AND M. BRADY, *An affine invariant salient region detector*, in European Conference on Computer Vision, Springer. Berlin, 2004, pp. 228–241.

[22] M. KARPUSHIN, *Local Features for RGBD Image Matching under Viewpoint Changes*, Ph.D. thesis, Telecom Paris Tech, Paris, 2016.

[23] Y. KE AND R. SUKTHANKAR, *PCA-SIFT: A more distinctive representation for local image descriptors*, Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Vol, 2, IEEE Computer Society, Los Alamitos, CA, 2004, pp. 506–513.

[24] S. KORMAN, D. REICHMAN, G. TSUR, AND S. AVIDAN, *Fast-Match: Fast Affine Template Matching*, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2013, IEEE, IEEE, Piscataway, NJ, 2013, pp. 1940–1947.

[25] S. LEUTENEGGER, M. CHLI, AND R. Y. SIEGWART, *BRISK: Binary robust invariant scalable keypoints*, in Proceedings of the IEEE International Conference on Computer Vision, IEEE, Piscataway, NJ, 2011, pp. 2548–2555, https://doi.org/10.1109/ICCV.2011.6126542.

[26] G. LEVI AND T. HASSNER, *LATCH: Learned arrangements of three patch codes*, in 2016 IEEE Winter Conference on Applications of Computer Vision, WACV 2016, IEEE, Piscataway, NJ, 2016, https://doi.org/10.1109/WACV.2016.7477723.

[27] T. LINDEBERG, *Scale-Space Theory in Computer Vision.*, Royal Institute of Technology, Stockholm, Sweden, 1993.

[28] T. LINDEBERG, *Direct estimation of affine image deformations using visual front-end operations with automatic scale selection*, in Proceedings of the Fifth International Conference on Computer Vision, 1995, IEEE Computer Society, Los Alamitos, CA, 1995, pp. 134–141.

[29] T. LINDEBERG, *Generalized Gaussian scale-space axiomatics comprising linear scale-space, affine scale-space and spatio-temporal scale-space*, J. Math. Imaging Vision, 40 (2011), pp. 36–81.

[30] T. LINDEBERG, *Invariance of visual operations at the level of receptive fields*, BMC Neurosci., 14 (Suppl. 1) (2013), P242.

[31] T. LINDEBERG AND J. GARDING, *Shape-adapted smoothing in estimation of 3-D depth cues from affine distortions of local 2-D brightness structure*, Proceedings of the ECCV, Pattern Anal. Mach. Intell., 1994, pp. 389–400.

[32] C. LIU, J. YUEN, AND A. TORRALBA, *SIFT flow: Dense correspondence across scenes and its applications*, IEEE Trans. Pattern Anal. Mach. Intell., 33 (2011), pp. 978–994.

[33] D. G. LOWE, *Distinctive image features from scale-invariant key points*, Int. J. Comput. Vis., 60 (2004), pp. 91–110.

[34] G. LOY AND J. O. EKLUNDH, *Detecting symmetry and symmetric constellations of features*, Proceedings of the ECCV, Vol. 2, Springer, Berlin, 2006, pp. 508–521.

[35] J. MATAS, O. CHUM, M. URBAN, AND T. PAJDLA, *Robust wide-baseline stereo from maximally stable extremal regions*, Image Vision Comput., 22 (2004), pp. 761–767.

[36] K. MIKOLAJCZYK AND C. SCHMID, *Indexing based on scale invariant interest points*, Proceedings of the ICCV, Vol. 1 , IEEE Computer Society, Los Alamitos, CA, 2001, pp. 525–531.

[37] K. MIKOLAJCZYK AND C. SCHMID, *An affine invariant interest point detector*, Proceedings of the ECCV, Vol. 1, IEEE Computer Society, Los Alamitos, CA, 2002, pp. 128–142.

[38] K. MIKOLAJCZYK AND C. SCHMID, *Scale and Affine Invariant Interest Point Detectors*, Int. J. Comput. Vis., 60 (2004), pp. 63–86.

[39] K. MIKOLAJCZYK AND C. SCHMID, *A performance evaluation of local descriptors*, IEEE Trans. Pattern Anal. Mach. Intell., 27 (2005), pp. 1615–1630.

[40] D. MISHKIN, J. MATAS, AND M. PERDOCH, *MODS: Fast and robust method for two-view matching.*, Comput. Vis. Image Underst., 141 (2015), pp. 81–93, http://doi.org/10.1016/j.cviu.2015.08.005.

[41] L. MOISAN, P. MOULON, AND P. MONASSE, *Automatic Homographic Registration of a Pair of Images, with A Contrario Elimination of Outliers*, IPOL J. Image Process. Online, 2 (2012), pp. 56–73, https://doi.org/10.5201/ipol.2012.mmm-oh.

[42] P. MOREELS AND P. PERONA, *Evaluation of features detectors and descriptors based on 3D objects*, Int. J. Comput. Vis., 73 (2007), pp. 263–284.

[43] J. M. MOREL AND G. YU, *On the consistency of the SIFT method*, Inverse Prob. Imaging, to appear.

[44] J.-M. MOREL AND G. YU, *ASIFT: A new framework for fully affine invariant image comparison*, SIAM J. Imaging Sci., 2 (2009), pp. 438–469.

[45] A. MURARKA, J. MODAYIL, AND B. KUIPERS, *Building local safety maps for a wheelchair robot using vision and lasers*, in Proceedings of the 3rd Canadian Conference on Computer and Robot Vision, IEEE, Piscataway, NJ, 2006.

[46]  P. Musé, F. Sur, F. Cao, and Y. Gousseau, *Unsupervised thresholds for shape matching*, in Proceedings of the International Conference on Image Processing, Vol. 2, IEEE, Piscataway, NJ, 2003, pp. 647–650.

[47]  P. Musé, F. Sur, F. Cao, Y. Gousseau, and J. M. Morel, *An a contrario decision method for shape element recognition*, Int. J. Comput. Vis., 69 (2006), pp. 295–315.

[48]  A. Negre, H. Tran, N. Gourier, D. Hall, A. Lux, and J. L. Crowley, *Comparative study of people detection in surveillance scenes*, Proceedings Structural, Syntactic and Statistical Pattern Recognition, Lecture Notes in Comput. Sci. 4109, Springer, Berlin, 2006, pp. 100–108.

[49]  Y. Pang, W. Li, Y. Yuan, and J. Pan, *Fully affine invariant SURF for image matching*, Neurocomputing, 85 (2012), pp. 6–10, https://doi.org/10.1016/j.neucom.2011.12.006.

[50]  D. Pritchard and W. Heidrich, *Cloth Motion Capture*, Comput. Graph. Forum, 22 (2003), pp. 263–271.

[51]  P. Scovanner, S. Ali, and M. Shah, *A 3-dimensional SIFT descriptor and its application to action recognition*, in Proceedings of the 15th International Conference on Multimedia, MULTIMEDIA '07, New York, 2007, ACM, New York, pp. 357–360, https://doi.org/10.1145/1291233.1291311.

[52]  S. Se, D. Lowe, and J. Little, *Vision-based mobile robot localization and mapping using scale-invariant features*, in Proceedings of the IEEE International Conference on Robotics and Automation, 2001, Vol. 2, IEEE, Piscataway, NJ, 2001, pp. 2051–2058.

[53]  J. Sivic, A. Zisserman, et al., *Video google: A text retrieval approach to object matching in videos.*, in ICCV, Vol. 2, IEEE Computer Society, Los Alamitos, CA, 2003, pp. 1470–1477.

[54]  N. Snavely, S. M. Seitz, and R. Szeliski, *Photo tourism: Exploring photo collections in 3D*, ACM Trans. Graph., 25 (2006), pp. 835–846.

[55]  C. Snoek, K. Sande, O. Rooij, B. Huurnink, J. Uijlings, M. v. Liempt, M. Bugalhoy, I. Trancosoy, F. Yan, M. Tahir, et al., *The MediaMill TRECVID 2009 semantic video search engine*, in TRECVID workshop, NIST, Gaitherburg, MD, 2009.

[56]  T. Tuytelaars and L. Van Gool, *Wide baseline stereo matching based on local, affinely invariant regions*, British Machine Vision Conference, Bristol, England, 2000, pp. 412–425.

[57]  T. Tuytelaars and L. Van Gool, *Matching widely separated views based on affine invariant regions*, Int. J. Comput. Vis. 59 (2004), pp. 61–85.

[58]  T. Tuytelaars and L. Van Gool, *Content-based image retrieval based on local affinely invariant regions*, International Conference on Visual Information and Information Systems, Springer, New York, 1999, pp. 493–500.

[59]  C. Valgren and A. J. Lilienthal, *SIFT, SURF and seasons: Appearance-based long-term localization in outdoor environments*, Robot. Auton. Syst., 58 (2010), pp. 149–156.

[60]  K. Van De Sande, T. Gevers, and C. Snoek, *Evaluating color descriptors for object and scene recognition*, IEEE Trans. Pattern Anal. Mach. Intell.,, 32 (2010), pp. 1582–1596.

[61]  M. Vergauwen and L. Van Gool, *Web-based 3D Reconstruction Service*, Mach. Vis. Appl., 17 (2005), pp. 411–426.

[62]  P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid, *Deepflow: Large displacement optical flow with deep matching*, in Proceedings of the IEEE International Conference on Computer Vision, IEEE, Piscataway, NJ, 2013, pp. 1385–1392.

[63]  G. Yang, C. V. Stewart, M. Sofka, and C. L. Tsai, *Registration of challenging image pairs: Refinement and region growing starting from a single keypoint correspondence*, IEEE Trans. Pattern Anal. Mach. Intell., 29 (2007), pp. 1973–1989.

[64]  G. Yu and J.-M. Morel, *ASIFT: An Algorithm for Fully Affine Invariant Comparison*, IIPOL J. Image Process. Online, 1 (2011), pp. 11–38, https://doi.org/10.5201/ipol.2011.my-asift.

[65]  H. Zhou, Y. Yuan, and C. Shi, *Object tracking using SIFT features and mean shift*, Comput. Vis. Image Underst., 113 (2009), pp. 345–352.