

SIFT Matching by Context Exposed

Fabio Bellavia

Abstract—This paper investigates how to step up local image descriptor matching by exploiting matching context information. Two main contexts are identified, originated respectively from the descriptor space and from the keypoint space. The former is generally used to design the actual matching strategy while the latter to filter matches according to the local spatial consistency. On this basis, a new matching strategy and a novel local spatial filter, named respectively blob matching and Delaunay Triangulation Matching (DTM) are devised. Blob matching provides a general matching framework by merging together several strategies, including rank-based pre-filtering as well as many-to-many and symmetric matching, enabling to achieve a global improvement upon each individual strategy. DTM alternates between Delaunay triangulation contractions and expansions to figure out and adjust keypoint neighborhood consistency. Experimental evaluation shows that DTM is comparable or better than the state-of-the-art in terms of matching accuracy and robustness. Evaluation is carried out according to a new benchmark devised for analyzing the matching pipeline in terms of correct correspondences on both planar and non-planar scenes, including several state-of-the-art methods as well as the common SIFT matching approach for reference. This evaluation can be of assistance for future research in this field.

Index Terms—Keypoint matching, SIFT, local image descriptors, local spatial filters, Delaunay triangulation, RANSAC, image context.



1 INTRODUCTION

KEYPOINT correspondences play a crucial role in many computer vision algorithms dealing with spatial localization. These include [1]: Structure from Motion (SfM), image stitching, large-scale image retrieval and Simultaneous Localization And Mapping (SLAM), whose practical applications are more and more affecting everyday life in providing assistance or for mere entertainment. The emerging autonomous driving systems, together with the various applications of the augmented reality from medicine to gaming, represent some relevant examples in that sense.

This state of things has granted an active interest on this research topic over the decades, continuously evolving side by side with the novel advancements and challenges arising in the field. In this scenario, despite its age, the Scale Invariant Feature Transform (SIFT) [2], both as keypoint detector and as local image descriptor, is in good health. As a matter of fact, SIFT is still popular [3], [4] and revisited [3], [5], [6]. Moreover, SIFT has obtained satisfactory results in recent benchmarks [3], [4], [7] and many applications still rely on it [8]. At the current state of the research, the keypoint extraction process and the computation of the associated local image descriptors, that must be synergically used to establishing correspondences, seem to have reached somewhat their limits when the matching process is considered untied from the context provided by the source images [9]. This observation is reflected in the progress done by deep learning descriptors [9]–[13] in conjunction with the availability of ever more big datasets [4], [14]–[16]. In this paper two different matching contexts are discussed.

The first one is provided by the *descriptor space*. Mutual Nearest Neighbor (NN) and the Nearest Neighbor Ratio (NNR) [2], two of the most commonly employed matching strategies are examples of this context exploitation. NN requires a match to be the best on both the input pair images, i.e. NN considers a *inter-relation* in the descriptor space.

NNR selects matches according to the ratio between the first and second best distances inside the reference image, i.e. NNR considers a *intra-relation* in the descriptor space.

The second matching context is the one provided by the *keypoint space* inside the images. This kind of scene knowledge includes patch relative orientations [6] and keypoint spatial relations [17], [18], which have been exploited to successfully disambiguate matches, boosting the final matching accuracy. Model constraints such those imposed by planar scenes and epipolar geometry [19], mainly built upon the Random SAmple COnsensus (RANSAC) [20]–[25], also operate on the keypoint space and generally represent the final post-filtering step of the entire process. Although very effective, this last kind of model-based matching is somewhat less general since it is not valid in all situations, such as in the case of moving or deformable objects that violate the constraint of a single rigid structure [26].

The recent evolution of deep networks jointly training a keypoint detector and descriptor [27], [28] has brought fresh attention to the matching step, inherently correlated to scene context, as the next step to be included in an all-in-one deep matching network [9], [29], [30] or as a replacement for RANSAC [31]–[34]. In this paper existent and novel general matching strategies exploiting image context in terms of both descriptor and keypoint spaces are presented and discussed. As baseline, their applications to SIFT are considered, but results on other keypoint detectors and descriptors are also shown. The contributions of this study can be mainly divided into three parts:

- Descriptor space context: Starting from a known greedy matching strategy using one-to-one NN together with NNR, several potential enhancements that include rank-based filtering as well as many-to-many and symmetric matching are investigated, extensively combined, and evaluated. The overall matching strategy obtained by merging altogether the best approaches, named blob matching, is proved

F. Bellavia is with the Department of Mathematics and Computer Science, Università degli Studi di Palermo, Via Archirafi, 34, 90123 Palermo, Italy, e-mail: fabio.bellavia@unipa.it .

to give more correct correspondences with respect to each individual matching strategy it is based upon.

- **Keypoint space context:** A new robust and effective local spatial filter named Delaunay Triangulation Matching (DTM) is designed. DTM associates matches greedily, according to the consistency of the keypoint spatial local neighborhoods. Neighborhood relations are obtained on the basis of Delaunay triangulations on the source images. DTM can guarantee comparable or better matching results with respect to the state-of-the-art.
- **Benchmarking image matching by context:** A new evaluation aimed at comparing state-of-the-art spatial matching strategies is devised, evaluating their behaviors on different keypoint detectors and recent descriptors on both planar and non-planar scenarios. In the non-planar case, when no 3D data are available, ground-truth correct matches are checked according to [3]. With respect to other evaluations merely relying on the epipolar distance [35], this setup avoid incorrect associations due to epipolar ambiguity and, differently from benchmarks based on camera pose estimation error [4], provides direct evaluation scores of the matches.

The rest of the paper is organized as follows: Related work is presented in Sec. 2, blob matching is defined in Sec. 3, DTM in Sec. 4, and proposed benchmark setup and results are discussed in Sec. 5. Finally, conclusions and future work are outlined in Sec. 6.

2 RELATED WORK

Assigning correspondences requires the close cooperation between three main subjects: The keypoint detector, the local image descriptor and the matching strategy.

A good keypoint detector for this task must be able to extract distinctive yet repeatable characteristic points on the input images. Moreover, the number of detected keypoints should be sufficient to provide a good coverage of relevant structures of the scene but, at the same times, the number of keypoints should be limited so as to make the computational process feasible and to reduce the chance of false matches. Keypoint distinctiveness decreases as the number of detected keypoints increases, especially in the presence of repeated structures on the scene.

Classical keypoint detectors usually extract corners or blob-like structures, the reader can refer to [1] for a general overview. Recently, keypoint detectors based on deep learning have started to emerge, especially in frameworks where detectors and descriptors sharing a common network are jointly estimated [27], [28], [30], or by combining handcrafted and deep learned filters [36]. Experimental evaluation in Sec. 5 will be carried out considering the state-of-the-art SIFT and the HarrisZ detector [37]. SIFT is a well known blob-like DoG multi-scale detector, while HarrisZ is an affine multi-scale corner detector. Relying on Harris corners computed on scale enhanced gradient derivatives, HarrisZ uses an adaptive filter response and a rough edge mask to select keypoints. In recent image matching challenges [38], a matching pipeline built on HarrisZ corners provided better

results than other keypoint detectors based on DoG (i.e. SIFT) keypoints.

Local image descriptors are used to obtain a meaningful numerical vector encoding the distinctive attributes of a keypoint patch, i.e. the local keypoint neighborhood. Ideally, descriptors for the same keypoint undergoing both geometrical or color distortions must be close in the descriptor vector space, and the opposite must hold for distinct keypoints. A trade-off between the descriptor tolerance to image deformation and its discriminability is often required, since high descriptor invariance decreases descriptor discriminability.

Local image descriptors can be divided into handcrafted and data-driven, the reader can refer to [1] for a general overview. The popular handcrafted SIFT descriptor is based on the gradient orientation histogram. Among data-driven descriptors, deep descriptors nowadays outperform any others [3], [4], [7] thanks to the modern GPU hardware capabilities and the availability of large training datasets [4], [14]–[16]. Triplet loss [10], hard negative mining [11], second-order similarity [13], geometric constraint integration [12] and jointly detector-descriptor optimization [9], [27], [28], [30] are some of the techniques employed to get state-of-the-art results.

The matching strategy evaluation carried out in Sec. 5 will employ the handcrafted descriptors RootSIFT [5] and the double square-rooting shifting Gradient Local Orientation Histogram (RootsGLOH2) [39]. Furthermore, the deep Second Order Similarity Network (SOSNet) [13] and HardNet2 [40] will be used. RootSIFT improves upon SIFT by replacing the Euclidean distance with the Hellinger’s distance and, in addition to RootSIFT, RootsGLOH2 is able to better handle patch orientation estimation. Deep-based SOSNet and HardNet2 are the current state-of-the-art, while RootsGLOH2 has been shown to be among the current best handcrafted descriptor [39].

Patch normalization is the interface between keypoint extraction and local descriptor computation, and it is often addressed together with this last step. The most common patch normalization approach is the one employed by SIFT, yet other approaches exists [41]–[43]. In particular, it has been reported in [3], [44] that the orientation estimation is one the most critical aspect that needs to be handled. The deep patch orientation assignment designed in [42] was proved to improve the matching accuracy noticeably [3] and will be employed for the evaluation presented in Sec. 5.

Matches are assigned by the pairwise inspection of the similarity in the descriptor space and optionally by considering the keypoint displacement in the images. The most common pipeline uses mutual NN or NNR matching followed by RANSAC. Using a symmetric variant of NNR [3] or considering the first geometrical inconsistent match in NNR [45] have been shown to generally improve the matching process. Further improvements have been observed when considering many-to-many putative matches instead of constraining matches to be one-to-one [46]. Many-to-many matches can be also related to the employment of multiple synthesized views to enrich the candidate matches [47].

Correspondence filtering on the basis of spatial constraints is a wide research topic, the reader can refer to [17], [18] for a more comprehensive presentation. The scene model can provide effective constraints, as in the case of pla-

nar and epipolar geometries. Although quite powerful, robust model regression strategies relying on RANSAC [20]–[24], can be limiting when violating the assumption of rigid scenes such in the case of moving and deformable objects [26], unless handling multiple models [48]. More non-global and relaxed spatial constraints have been designed to overpass this limitation, also able to boost RANSAC in terms of both efficiency and accuracy when employed to pre-filter candidate matches. When feasible, executing RANSAC after any other kind of match selection is always the best practice.

An early example of spatial filtering is the topological filter [49], which checks the relative positions of the matched keypoints across the images. Better filtering approaches consider the local consistency around the keypoints of the pair defining the match. Local neighborhood can be defined by using a fixed circular radius and measuring consistency in terms of the number of other matches having both keypoints on the respective images falling into the neighborhoods associated to the considered match. In addition, the relative local image transformation between corresponding keypoints inducted by patch normalization can provide a further consistency check to be exploited to refine local neighborhoods. This can be done by comparing the transformation parameters or by considering the reprojection errors according to the transformations. Circular neighborhoods refined according to patch-based consistency on local similarity transformations are checked in [22] to pre-filter matches before RANSAC. Grids can also be employed to define neighborhoods efficiently, as for the Grid-based Motion Statistics (GMS) [50]. By defining neighborhoods as 3×3 grid blocks while measuring the consistency block-wise, GMS is able to take into account the relative positions of the matches inside the neighborhoods. Moreover, since the best neighborhood size is generally not known a priori, GMS makes use of different grid sizes to get a multi-scale approach. Another possible choice, employed in the Locality Preserving Matching (LPM) [26] is to define circular neighborhoods by considering only the closest k matches in terms of keypoint proximity. Unlike GMS, LPM limits the neighbors to those matches having motion flow vectors similar to that of considered match, and considers multiple values of k to be more robust. In the case of the Guided LPM (GLPM) [51], NNR match pre-filtering precedes LPM. LPM neighborhood definition is also employed by the Learning for Mismatch Removal (LMR) [52] to extrapolate local correspondence relations to learn how to classify correspondences. Neighborhoods can be also defined in terms of the edges of the Delaunay triangulation, as for the Progressive Feature Matching (PFM) [53], employing affine patch-based neighborhood consistency to cast the problem as a Markov random field function optimization. The Progressive Graph Matching (PGM) [54] uses instead another graph matching formulation with neighborhoods defined by the k closest matches and affine patch-based consistency, while in [24] graph-based optimization of fixed circular neighborhoods is used to improve the best selected RANSAC model.

Spatial filtering can be also formulated as the estimation and outlier rejection of the motion vector field inducted by the scene on the images. In that sense, both PFM and PGM estimate local affine motion fields from sparse corre-

spondences. According to this idea, grid neighborhood is also employed in [55] to collect motion hypotheses and discard the outlier matches violating them. More complex approaches such as the Locally Linear transforming (LTT) [56], the Vector Field Consensus (VFC) [57] and the Bilateral Model (BM) [58] simultaneously estimate a smooth motion field from the putative matches while removing outliers. To a certain extend, also the Sequential Correspondence Verification (SCV) [59], which considers the region growing progression around the local affine patch neighborhoods to make decisions about the correctness of a match, can be framed in this kind of design.

Other solutions relying on the spatial relations between correspondences are designed as the clustering of the motion vector field. In [23] images are segmented according to their spatial and motion information prioritizing RANSAC model sampling from large, more consistent groups. In a similar spirit, the recent Adaptive Locally-Affine Matching (AdaLAM) [25] executes multiple local affine RANSACs inside the circular neighborhoods of seed matches, i.e. those matches with high confidence according to their descriptor similarity. The Robust Feature Matching Spatial Clustering of Applications with Noise (RFM-SCAN) [60] uses instead clustering to directly identify outlier correspondences not belonging to any cluster. A game theoretic formulation is instead proposed in [61], which iteratively outputs the best clusters only on the basis of the patch-based local transformation consistency in terms of reprojection error, disregarding keypoint proximity¹.

More recently, deep networks have reached the state-of-the-art in the case of planar or epipolar geometry scene constraint matching [31]–[34]. In [31] the context normalization is introduced to exploit contextual information while preserving permutation equivariance. The network architecture of [32] defines instead grouping operations to supply each correspondence with data from its nearest correspondences only in terms of patch-based local transformation consistency as in [61]. In addition to context normalization, the Order-Aware Network (OANet) [33] adds network layers to learn how to cluster unordered sets of correspondences so as to incorporate the data context and the spatial correlation, while the Attentive Context Network (ACNe) [34] employs local and global attention to exclude outliers from context normalization. Attentional graph neural networks is instead proposed in [29] to infer global and local spatial correlations. More recently, the network designed in [62] was able to get keypoint-free correspondences for a coarse-to-fine matching strategy, with initial rough dense correspondences provided by exploiting self and cross attention inferred by transformers.

3 BLOB MATCHING

Blob matching mainly works by filtering matches through descriptor space context heuristics. Before defining the blob matching, the base matching strategies it relies upon are reviewed and discussed. The notation adopted hereafter assumes that $D \in \mathbb{R}^{n \times m}$ is a matrix such that the entry D_{ij}

1. However, given two matches, it usually holds that the closer are the respective keypoints on the same image, the more consistent is their patch-based consistency according to the reprojection error.

is the distance between the descriptors associated to the i -th keypoint on the first image I_1 and the j -th keypoint on the second image I_2 . D_{ij^w} denotes the w -th lowest value found on row i , i.e. the w -th best descriptor association. Likewise, $D_{i^w j}$ is the w -th lowest value on column j , and the subscript \Downarrow denotes the extraction of the index pair (i, j) from D_{ij} , done element-wise in the case of a set of matches. NN is the basic way to associate correspondences

$$S_{NN} = \{D_{ij^1}\}_{\Downarrow\Downarrow} \quad (1)$$

Here, the index spans $1 \leq i \leq n$ and $1 \leq j \leq m$ are omitted for convenience. Mutual NN constrains even more S_{NN} by requiring that selected matches must be the best in both images, i.e. simultaneously on both the rows and columns of D

$$S_{mNN} = \{D_{ij^1} = D_{i^1 j}\}_{\Downarrow\Downarrow} \quad (2)$$

It is easy to see that $S_{mNN} \subseteq S_{NN}$. NNR [2] considers instead a matrix D' such that

$$D'_{ij} = \frac{D_{ij}}{D_{ij^2}} \quad (3)$$

and the threshold t_r , usually set to 0.8, so that

$$S_{NNR} = \{D'_{ij^1} \leq t_r\}_{\Downarrow\Downarrow} \quad (4)$$

NNR can be related to the triplet matching learning adopted by deep descriptors starting from [10]. Notice that $D'_{ij} \in [0, 1]$ by definition, and $D'_{ij^1\Downarrow} = D_{ij^1\Downarrow}$ for any row since scalar multiplication does not affect the ordering. NNR is usually more accurate than mutual NN since NNR relative values can better express the image context than NN absolute values. Nevertheless, S_{mNN} is symmetric since it takes both the images as reference, and provides a one-to-one matching relation. By contrast, S_{NNR} considers only the first image I_1 as reference, since the denominator in Eqs. 2 is computed on rows, providing a one-to-many matching relation. In order to relax the strict requirements of mutual NN, the following greedy strategy, to the best of the author's knowledge first mentioned in [63], can be employed. As notation, D_k is the k -th best entry of the matrix D considered as linear vector and G is the set of matches, initially empty. At each iteration $1 \leq k \leq n \times m$, D_k is included into G if both $D_{k\Downarrow} \notin G_{\Downarrow}$ and $D_{k\Downarrow} \notin G_{\Downarrow}$, where the subscripts \Downarrow and \Downarrow denote the operators that extract the row and column index of the entry, respectively. The final matching set is

$$S_{gNN} = G_{\Downarrow\Downarrow} \quad (5)$$

Under the assumption that D entries have unique values, it holds that $S_{mNN} \subseteq S_{gNN}$ and $|S_{gNN}| = \min(n, m)$. If a mutual match (i, j) of S_{mNN} has not been included in S_{gNN} , there would have been a match (\tilde{i}, \tilde{j}) with $i = \tilde{i}$ or $j = \tilde{j}$ in a previous iteration such that $D_{\tilde{i}\tilde{j}} < D_{ij}$. This is a contradiction since $D_{ij} = D_{i^1 j} = D_{ij^1}$ by definition. Figure 1a shows the differences between S_{NN} , S_{mNN} and S_{gNN} .

Mutual or greedy NN can be combined with NNR, in the sense that matches are extracted by NN but sorted and eventually filtered according to the NNR ranking. In this case, when the greedy NN is employed with NNR, one have to replace Eq. 2 with an alternative definition

$$D'_{ij} = \frac{D_{ij}}{D_{ij^2}} \quad (6)$$

D_{ij^w} and $D_{i^w j}$ are the lowest w -th values greater or equal to D_{ij} on the j -th column and on the i -th row, respectively. This further constraint is necessary since D_{ij} may not be equal to D_{ij^1} so that $D'_{ij} > 1$. In order to improve the matching process, a symmetric NNR is proposed in [3] as the harmonic mean between the two entries obtained by swapping the reference image, corresponding to operate on the matrix transpose D^T of D

$$D''_{ij} = \frac{2D'_{ij}(D^T)'_{ji}}{D'_{ij} + (D^T)'_{ji}} \quad (7)$$

In the particular case of Eq. 6 the harmonic mean becomes

$$D''_{ij} = \frac{2D_{ij}}{D_{ij^2} + D_{i^2 j}} \quad (8)$$

The First Geometrically Inconsistent NN (FGINN) [45] is another possible improvement to NNR

$$D'_{ij} = \frac{D_{ij}}{D_{ij^2_{\circledast}}} \quad (9)$$

Here, the second lowest value in the denominator of Eq. 2 is intended among those keypoints far at least $t_o = 10$ pixels from the keypoint j in the corresponding image, denoted as $D_{ij^2_{\circledast}}$. The choice of the second lowest value according to FGINN is shown in Fig. 1b. Although the keypoint position implies to work on the keypoint space, this approach mainly deals with the descriptor space and it is discussed here. Finally, in [46], many-to-many match relations have been shown to improve the recall of the matching process, leading to better samples for the RANSAC hypothesis generation.

Aimed at incorporating all the matching strategies discussed so far, the blob matching is now formulated according to the following steps.

- 1) The similarity matrix D is pre-filtered so that only matches appearing among the f best matches for both or any of the input pair images will be considered [4], giving rise respectively to one of the sets

$$F_{\cap} = \{D_{ij} \leq D_{ij^f}\} \cap \{D_{ij} \leq D_{i^f j}\} \quad (10)$$

$$F_{\cup} = \{D_{ij} \leq D_{ij^f}\} \cup \{D_{ij} \leq D_{i^f j}\} \quad (11)$$

As additional notation, the subscript \cap or \cup will be appended to the value of f to shortly refers to F_{\cap} or F_{\cup} , respectively, while F will indicate one of the two sets indistinctly.

- 2) Surviving matches are then filtered according to the greedy approach, modified to take into account the first f' best matches instead of only the first one. In detail, after being sorted by increasing values, $D_{ij} \in F$ is added to the multiset G if both $(G \cup D_{ij})_{\Downarrow}$ and $(G \cup D_{ij})_{\Downarrow}$ do not contain elements counted more than f' times (see again Fig. 1a).
- 3) NNR-like similarity values for the matches in G are obtained by taking into account a many-to-many scheme but also FGINN. Equation 6 is modified into

$$D'_{ij} = D \quad (12)$$

D					setup	G_{\downarrow}
1.6	2.5	1	4	2.3	NN by row	$S_{NN} = \{(1, 3), (2, 2), (3, 4), (4, 2), (5, 3), (6, 1), (7, 1)\}$
4.2	0.5	1.7	3	1.1	NN by column	$S_{NN} = \{(1, 3), (2, 2), (2, 5), (3, 4), (7, 1)\}$
5.1	3.5	3.1	1.2	2	$f = 1_{\cap}, f' > 0$	$S_{mNN} = \{(2, 2), (1, 3), (3, 4), (7, 1)\}$
2.8	0.6	2.1	4.1	5	$f = 1_{\cup}, f' = 1$	$S_{mNN} = \{(2, 2), (1, 3), (3, 4), (7, 1)\}$
4.4	3.4	2.4	4.3	4.5	$f \geq \Omega, f' = 1$	$S_{gNN} = \{(2, 2), (1, 3), (3, 4), (7, 1), (6, 5)\}$
3.2	5.5	5.8	6.1	3.6	$f > 0_{\cup}, f' = 1$	$S_{gNN} = \{(2, 2), (1, 3), (3, 4), (7, 1), (6, 5)\}$
1.3	6	3.7	2.7	1.4	$f = 3_{\cap}, f' = 1$	$\{(2, 2), (1, 3), (3, 4), (7, 1)\}$
					$f = 1_{\cup}, f' = 2$	$\{(2, 2), (4, 2), (1, 3), (2, 5), (3, 4), (7, 1), (5, 3), (6, 1)\}$
					$f = 3_{\cap}, f' = 2$	$\{(2, 2), (4, 2), (1, 3), (2, 5), (3, 4), (7, 1), (7, 5), (1, 1), (4, 3)\}$
					$f \geq \Omega, f' = 2$	$\{(2, 2), (4, 2), (1, 3), (2, 5), (3, 4), (7, 1), (7, 5), (1, 1), (4, 3), (5, 4)\}$
					$f = 3_{\cup}, f' = 2$	$\{(2, 2), (4, 2), (1, 3), (2, 5), (3, 4), (7, 1), (7, 5), (1, 1), (4, 3), (5, 4)\}$

$\Omega = \max(n, m)$ with $D \in \mathbb{R}^{n \times m}$, in this case $F_{\cap} = F_{\cup}$

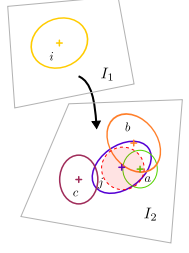


Fig. 1. (a) The candidate match sets G_{\downarrow} for different blob matching configurations on a toy example input D . Some setups may always lead to the same G_{\downarrow} set, regardless of D , as for the third and fourth setup rows. This is not the case of the seventh setup row, where the same G_{\downarrow} output is due to the choice of D . (b) Visual representations of NNR and FGINN. If it holds for D that $D_{ij} \leq D_{ia} \leq D_{ib} \leq D_{ic}$ and the spatial configuration of the associated patch ellipses is the one reported in the figure, for the match (i, j) then D_{ij}/D_{ia} , D_{ij}/D_{ib} and D_{ij}/D_{ic} are the values of NNR, FGINN using keypoint distance and FGINN using the overlap error, respectively (see the text for details, best viewed in color).

where \mathcal{D} can be further specified as

$$\mathcal{D}_{\geq} = \frac{D_{ij}}{D_{ij_{\geq}^{\odot}}} \quad (13)$$

The symbols \geq and \odot are the same as for Eqs. 6 and 13, respectively. The threshold t_o can also express a relative threshold and not only an absolute pixel distance. In the first case, a filtering based on the overlap error between the elliptical patches is considered instead of the original FGINN criterion based on the keypoint center distance (see again Fig. 1b). This allows to rely on non-absolute values. A further specialization for \mathcal{D} is also considered

$$\mathcal{D}_{*}^{+} = \frac{D_{ij}}{D_{ij} + D_{ij_{*}^{\odot}}} \quad (14)$$

The $*$ symbol is a placeholder for \geq , and can be possible empty, since \geq is not strictly required to accommodate values into the range $[0, 1]$. This happens because the whole column or row spans can be considered instead of limiting the selection only to the values greater than the current one as it was required for $D_{ij_{\geq}^{\odot}}$ in Eq. 6.

- 4) Lastly, in order to provide the final similarity score \bar{D}_{ij} for matches in G to be used for sorting or thresholding the keypoint pairs, a function \mathcal{W} is applied so as to combine the two possible matching similarity values obtained when considering each image as reference. In detail, using the first image as reference corresponds to employ D'_{ij} , while using the other image corresponds to $(D^{\top})'_{ji}$ so that

$$\bar{D}_{ij} = \mathcal{W}(D'_{ij}, (D^{\top})'_{ji}) \quad (15)$$

$\mathcal{W}(a, b)$ can simply be the projection on one of the two arguments ($\mathcal{W}(a, b) = a$ or $\mathcal{W}(a, b) = b$), the harmonic mean of Eq. 7 ($\mathcal{W}(a, b) = (2ab)/(a + b)$), and the minimum ($\mathcal{W}(a, b) = \min(a, b)$) or maximum ($\mathcal{W}(a, b) = \max(a, b)$) values of the two arguments, likewise respectively the intersection and union of the matching sets employed in [4].

The first two steps of blob matching extract the set G_{\downarrow} of the candidate matches, ranked by the remaining two steps.

Clearly, it may happen that different initial configurations lead to the same G_{\downarrow} set (see again Fig. 1a). It comes out in Sec. 5.2 that the best blob matching configuration is $f = 1_{\cup}, f' = 5, \mathcal{D}^{+}$ with $t_o = 10$ px or $t_o = 75\%$, and $\mathcal{W}(a, b) = (2ab)/(a + b)$. Notice also that blob matching is quite general: By setting $f = 1$ only mutual NN matches are considered, moving up f' from 1 when $f > 1$ the one-to-one match relation becomes a many-to-many relation, FGINN is turned off when $t_o = \infty$, and only the first image is used as reference when $\mathcal{W}(a, b) = a$.

4 DELAUNAY TRIANGULATION MATCHING

As discussed in Sec. 2, spatial filtering relies on the concept of the image local neighborhood, often assumed to be isotropic, circular or squared, for an efficient and fast computation. In the general case, the optimal circular radius needed to define the neighborhood is not known a priori and can vary among different image regions due to the non-homogeneous distribution of keypoints on the image. In order to alleviate this issue, the neighborhood estimation can proceed in a coarse-to-fine manner using different radius, or it can consider the k closest keypoints constrained by the similarity of the motion flow or by the inducted local patch-based transformations. DTM employs an alternative neighborhood definition based on the Delaunay triangulation, which naturally fits into the keypoint distribution of the image and its structure, implicitly providing a sort of dynamic neighborhood without requiring to supply the neighborhood size. Unlikely [53], appeared after the original submission of this manuscript, DTM considers Delaunay triangulations from both the images of the input pair to get more consistent and refined neighborhoods as intersection. This approach was already proposed in [3] with the aim of benchmarking descriptors for growing up good matches from an initial set of ground-truth correspondences.

When computing the Delaunay triangulation, only boundaries need some attention. Boundary edges of any Delaunay triangulation correspond to the convex hull of the considered keypoints, and triangles for keypoints on the edges of convex hull are generally not well-shaped and do not lead to appropriate neighborhoods. Adding the image corners does not solve the problem as well as breaking the image canvas borders into multiple lines. A

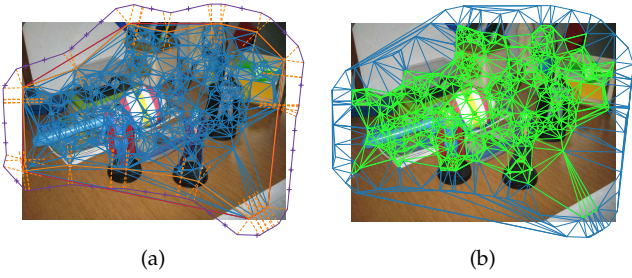


Fig. 2. (left) original Delaunay triangulation (blue), convex hull (red), alpha shape boundary edges (solid orange), border fattening (dashed orange), and final contour edges (purple). (right) final Delaunay triangulation with keypoint-only marked edges (green) (see text for details, best viewed in color and zoomed in).

more feasible solution is to expand the boundary edges and split them, where the boundary can be determined by the convex hull or, better, by alpha-shapes², which also relies on Delaunay triangulation. A visual explanation of the different boundary choices is reported in Appendix A, while the alpha shape border computation adopted for DTM is illustrated in Fig. 2. Being V a keypoint set, alpha-shape boundary edges (solid orange) are extracted and expanded. Specifically, from the two vertexes of each alpha-shape edge, the four points at a fixed distance s from the edge³ lying on the line perpendicular to it are included in the set B' (dashed orange). Alpha-shape boundary edges for $V \cup B'$ are then extracted and break down into segments of length s , whose vertexes provide the final boundary point set B for the triangulation (purple).

The use of Delaunay triangulation in DTM is twofold: In the first DTM₁ stage, Delaunay triangulation is combined with the greedy matching strategy employed to define $S_{g,NN}$ in Eq. 5 to iteratively prune matches. In the second DTM₂ stage, it is employed to grow up consistent matches from the previously surviving matches. Note that DTM input is only the set of keypoint correspondences and their similarity in the descriptor space context, without considering additional patch-based local transformation information for consistency. This make DTM more general and potentially robust when the relative patch-based local transformations are unavailable, unreliable or unable to provide a good approximation of the specific input motion field, e.g. when relying on patch-based similarity or affine local transformations in case of severe perspective distortions.

DTM₁ repeatedly and greedily removes inconsistent neighbor matches and restores the consistent ones, progressively adjusting neighborhoods. It is known that both the neighborhood graph and the minimum spanning tree are subsets of the Delaunay triangulation, which also maximizes the minimum angle of each triangle mesh. This late property gives rise to neighborhoods well spread among all directions. Moreover, as shown in Appendix B, in the ideal case Delaunay-based neighborhoods have intuitively better chances to contain correspondences consistent with the considered match and hence to restore accidentally removed good matches.

2. Using the Matlab `boundary` function with default parameters.

3. The value of s is empirically set to the minimum between the width and the height of the image divided by 10.

Given an initial set of matches M^0 , each iteration i of DTM₁ iteratively prunes M^i until $M^i = M^{i-1}$ (see Fig. 3) as described by the following steps:

- 1) Extract keypoint locations for current (surviving) matches. Set $K_1^i = M_{\downarrow 1}^{i-1}$ and $K_2^i = M_{\downarrow 2}^{i-1}$.
- 2) Construct the current Delaunay triangulation of each image. Round-off keypoint coordinates of K_1 to define the vertex sets $V_1 = \{([k_x], [k_y]) \in K_1\}$. From V_1 , compute the boundary set B_1^i using alpha-shapes as described before, and build the Delaunay triangulation \mathcal{T}_1 for I_1 from $V_1 \cup B_1^i$ (see Fig. 3, left column). Analogously, define V_2, B_2^i and \mathcal{T}_2 (see Fig. 3, right column). Vertex collapsing by rounding-off avoids many too small triangles that can slow-down the computation.
- 3) Define the local non-isotropic neighborhoods. For each vertex $v \in V_1$, define A_v^1 as the set of vertexes adjacent to v in the triangulation \mathcal{T}_1 , including v itself. Define also $M_{A_v^1}^i$ as the set of matches in M^{i-1} each having the keypoint lying on I_1 collapsed into a vertex in A_v^1 by rounding-off, as described before. A_v^2 and $M_{A_v^2}^i$ are defined analogously for I_2 .
- 4) Rank matches according to their coherence. Assign a rank $r(m)$ to matches $m \in M^{i-1}$ collapsed into vertex pair (v_l, v_r) , by sorting them first according to their increasing descriptor similarity and then by the decreasing cardinalities $|M_{A_{v_l}^1}^i \cap M_{A_{v_r}^2}^i|$.
- 5) Contract the Delaunay triangulations. Set $T = \emptyset$, $M' = M^{i-1}$, and add the match $m \in M'$ ranked first according to $r(m)$ to T . Then remove this match m from M' as well as matches in $(M_{A_{v_l}^1}^i \cup M_{A_{v_r}^2}^i) \setminus (M_{A_{v_l}^1}^i \cap M_{A_{v_r}^2}^i)$, where m collapsed into (v_l, v_r) . Repeat until $M' = \emptyset$ (see Fig. 3b,d).
- 6) Expand the Delaunay triangulations. Define M^i as the union of the sets $(M_{A_{v_l}^1}^i \cap M_{A_{v_r}^2}^i)$, each one obtained from the collapsing pair (v_l, v_r) a match $m \in T$ corresponds to (see again Fig. 3b,d).

The set $M^{\bar{i}}$ at the last iteration \bar{i} contains the matches survived to the Delaunay triangulation “pulses” (see Fig. 3e). The convergence is always guarantee since by construction the cardinality of M^i cannot increase with the iterations i . In the worst case, no sufficient matches for the Delaunay triangulations are found in the last iteration, and DTM outputs no correspondences. Notice that DTM₁ greedy formulation does not involve parameters, unlikely other approaches requiring parameters to define energy minimization cost functions and criteria to stop the execution either to select the initial or final matches.

DTM₂ employs the Delaunay triangulations of the correspondences survived to DTM₁ to approximate the motion field and hopefully restore consistent matches previously discarded. Delaunay triangulation allows a weak model assumption for motion field, i.e. correct matches should be inside corresponding triangular meshes, without any explicit motion field characterization as imposed by other formulations.

In order to pick up good the matches accidentally discarded by DTM₁ since surrounded only by wrong matches, DTM₂ proceeds in reverse order, starting from iteration

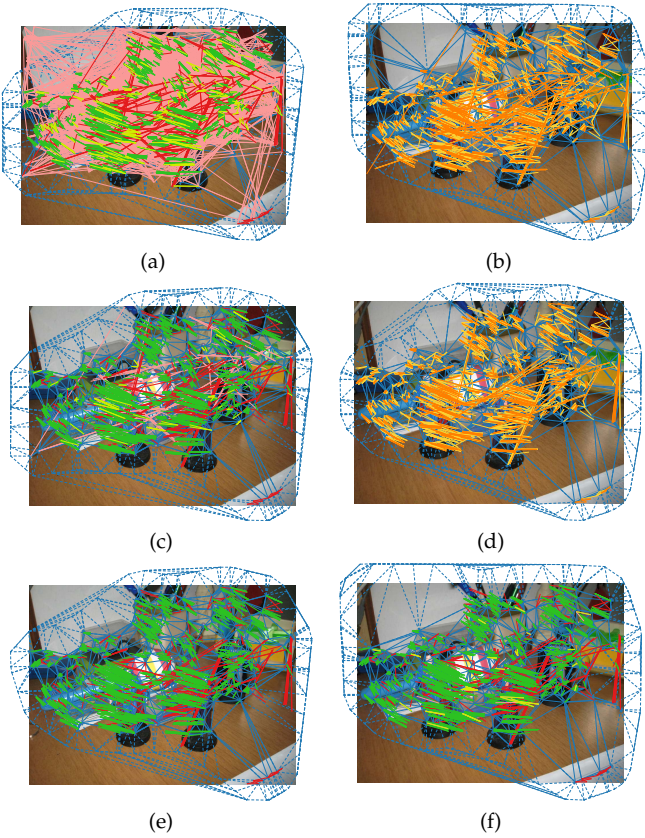


Fig. 3. DTM computation steps. The first and second images of the input pair are superimposed on the left and on the right, respectively, together with their associated Delaunay triangulations (blue). Initial matches are obtained from the best configuration given by HarrisZ+SOSNet with blob matching (see Sec. 5). The first and second iterations i of DTM_1 are reported as the first and second rows, respectively. For these rows, on the left, the clusters of vector flows for the retained (green and red) and pruned (yellow and light red) matches are shown. The clusters of correct (green and yellow) and wrong (red and light red) matches can be established as well, according to the evaluation protocol described in Sec. 5. Corresponding contraction (orange) and expansion (yellow) clusters are indicated on the right. Image (e) reports the final filtered matches at the fourth last iteration \bar{i} of DTM_1 , while image (f) shows the final matches after the last iteration $i = 0$ of DTM_2 . In this last image the colored clusters indicate the matches before DTM_2 (green and red) and those added (yellow and light red), correct (green and yellow) or wrong (red and light red) (see text for details, best viewed in color and zoomed in).

$i = \bar{i} - 1$ downto 0, with $E^{\bar{i}-1} = M^{\bar{i}} = M^{\bar{i}-1}$ by the definition of \bar{i} in DTM_1 :

- 1) Construct the Delaunay triangulations for the (estimated) good matches. At step i , compute the Delaunay triangulations \mathcal{T}'_1 and \mathcal{T}'_2 as before, but from the collapsed vertex set of E^i plus the boundary sets B^i_1 and B^i_2 computed in the DTM_1 stage.
- 2) Add coherent matches according to the Delaunay triangulations. Initially set $E^{i-1} = E^i$. For each match m in $M_i \setminus M^{i-1}$, find the triangle W_1 of \mathcal{T}'_1 where the corresponding keypoint of m on I_1 falls into. If the other keypoint of m on I_2 falls into any triangle on I_2 formed by the corresponding vertexes of W_1 , and this also holds when swapping the role of the images, add the match to E^{i-1} .

The final set $E = E^0$ contains the enhanced matches (see Fig. 3f). Notice that DTM_2 uses the boundary sets B_1 and

B_2 computed in DTM_1 in order to increase the chance to include previously discarded keypoints close to the boundary. Like the previous stage, no parameters get involved in the computation.

5 EVALUATION

5.1 Setup

The evaluation pipeline is composed by the following steps: keypoint extraction, local descriptor computation, descriptor matching, local spatial filtering and model fitting. For the keypoint extraction the SIFT and HarrisZ detectors are considered, while RootSIFT, RootsGLOH2, HardNet2 and SOSNet are employed as local descriptors. SIFT and RootSIFT are included as baselines. Patch orientation is estimated according to [42] for all descriptors, except for RootsGLOH2 that needs no orientation-adjusted patches. Descriptor matching, the next step of the pipeline, is achieved only by blob matching, since it can behave as the common descriptor matching strategies with a proper tuning of the parameters. The goal of this evaluation step is to check the advantages offered by match pre-filtering, many-to-many matches, and the alternative distance definitions and combinations (referring in order to the four different steps of blob matching in Sec. 3). The next step of the pipeline evaluates spatial matching filters, including the proposed DTM and fourteen state-of-the-art filters, learned or not: LMR, LPM, GLPM, GMS, VFC, LLT, RFM-SCAN, AdaLAM, OANet, ACNe, PFM, PGM, SCV and BM (see Sec. 2). For reference, the standard 0.8 NNR threshold is also included, although not properly a local spatial filter, and indicated as ‘th’. For the last step of the pipeline, being q the minimal number of matches required to estimate the model⁴, only a simple global model estimation using uniquely one sample made up of the $3 \times q$ top-ranked surviving matches at the previous step is employed. Although this approach, named 1SAC (one SAmples Consensus) is quite naive, it can give insights on more complex RANSAC approaches. Notice that except for th, GLPM and DTM (step 4 of stage DTM_1 in Sec. 4), other pruning methods do not take into account the descriptor context, i.e. the descriptor similarity.

For the evaluation both planar and non-planar scenes are considered, the latter being more complex due to the inclusion of spatial discontinuities caused by occlusion and parallax. In the case of planar scenes, the 15 sequences employed in [3] were used. Each sequence is made up of 6 images where the first one is fixed as reference, for a total of $19 \times (6 - 1) = 75$ image pairs. In the case of non-planar scenes, two different datasets with distinct evaluation protocols are considered. The first one, explicitly referred to as the non-planar dataset, includes the 72 image pairs from [3], and additionally 27 image pairs already known in the computer vision community, for a total of $72 + 27 = 99$ image pairs belonging to 61 different scenes. For each image pairs of this dataset, a sparse set of hand-taken ground-truth correspondences and occluded points is available. Appendix C shows the scenes contained in the planar and non-planar datasets. The second dataset,

4. The value of q is respectively set to 4, 8 for homography, fundamental matrix.

SUN3D [14], contains 415 indoor sequences, whose only 401 of them were supplied with the additional data required to extrapolate almost dense ground-truth correspondences. The training sequences of ACNe, corresponding to half of the SUN3D dataset are also included in the evaluation, since as reported in the additional material no relevant differences were observed when taken into account. For each sequence in SUN3D, a maximum of 30 image pairs, uniformly distributed among the sequence time interval, were chosen, for a total of 11231 pairs. Images making each SUN3D pair correspond to a time step of 80 frames, unless the maximal visual overlap between the images is less than 25%, in this case the frame step is lower.

The proposed evaluation relies on the computation of ground-truth matches. On the planar dataset, ground-truth correspondences can be easily obtained by estimating the homography that maps one-to-one points across the two images [15], [63]. On the non-planar dataset, unless 3D data are supplied, only a point-to-line mapping through the fundamental matrix is available according to the epipolar geometry. A common approach to evaluate RANSAC-like methods defines inlier matches (hence ground-truth matches) according to the distance between a point and the epipolar line imaging the corresponding point to be matched in the other image (here denoted as method A). However, this approach can lead to many false positive ground-truth matches due to the ambiguity of the map. Other approaches measure the error on the fundamental matrix obtained at the end of matching process with respect to a ground-truth one, which can be estimated by hand-taken correspondences [64] or on the basis of robust SfM approaches [4], [38]. The error distance between corresponding ground-truth and estimated epipolar lines can be considered to judge the goodness of the matching, statistically [65], or on the basis of a uniform image sampling [64]. The former approach has been employed in a recent evaluation [35]. Alternatively, the pose error of the camera [4] is obtained from the estimated fundamental matrix (to be understood in a broad sense). Nevertheless, all these evaluations give an indirect measure on the goodness of matches that does not guarantee a true evaluation of the matching process. This happens since there is no direct and clear formula relating fundamental matrix correctness and correct matches, so that from the perspective of evaluating the accuracy of the correspondences, this kind of solutions can be associated to method A. Notice that although SfM can estimate depth data to classify correct matches, these are in most cases incomplete or contain errors, so that a further possible solution can be addressed by limiting the evaluation to synthetic datasets [38]. Finally, another approach is to compute ground-truth matches by relying on additional constraints inducted for instance by employing more images [66] or by manually limiting the spatial localization of the matches [3]. This last solution is adopted in following evaluation according to the protocol described in [3] based on the approximated overlap error (denoted as method B), extended to cope with the issues discussed hereafter. True or approximated patch overlap can lead to the presence of large patches with low overlap error but where the keypoint center distance measured in pixels not acceptable by visual inspection. To cope with this, homography reprojection and the epipolar line distance errors

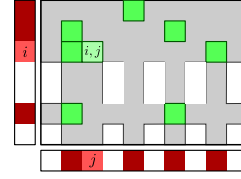


Fig. 4. Visual representation of the normalized recall on a given set of many-to-many matches (i, j) (green). $|Z_{\downarrow}| = 4$ and $|Z_{\leftarrow}| = 5$ by counting the numbers of elements (red) in the row and column outer sets, respectively. In the example $|Z| = \min(4, 5) = 4$ (see text for details, best viewed in color and zoomed in).

are employed (limited to 30 pixels) as additional constraints to the true or approximated overlap errors (limited to 50%), providing visually adequate results (method C). Moreover, taking into account that not all the local spatial filters handle the shape and size of the patch, an evaluation based only on the keypoint center distance disregarding the patch overlap can be conceived to define ground-truth matches (method D). Specifically, in the planar case the distance between a keypoint and the reprojection of the corresponding keypoint in the other image must not exceed 15 pixels. For the non-planar case the epipolar distance must not exceed 15 pixels but the matches must also pass the check defined on the basis of their spatial localization [3]. With respect to method D, method C obtains a lower number of correct matches due to the scale constraints according to the patch shapes. A comparison of the ground-truth estimation methods A, B, C and D is reported in Appendix D, upon which after visual inspection method D has been selected for the planar and non-planar datasets. For the SUN3D dataset, an analogous of method D is employed. More specifically, each of the corresponding keypoints of a match is reprojected from one image to the other by exploiting depth and the extrinsic camera matrix data. The match is considered correct if at least in one case the reprojection error is less than 15 pixels or, in the remote eventuality that no depth estimation is available, if the maximum epipolar error is less than 15 pixels.

Lastly, the definition of the recall is adjusted to handle many-to-many matches and avoiding apparent boosted results. In particular, using the same notation of Sec. 3, given the set Z of good matches according to the ground-truth, the associated number of correct matches necessary to compute the recall is defined as $\min(|Z_{\downarrow}|, |Z_{\leftarrow}|)$ instead of $|Z|$, so that multiple keypoint associations are only taken into account once (see Fig. 4). The precision is instead computed as usual. Notice that only in the case of mutual one-to-one matches $|Z| = \min(|Z_{\downarrow}|, |Z_{\leftarrow}|)$, so that the new recall definition extends the standard one. This sort of normalization is also employed to compose the top-ranked sample in 1SAC to handle many-to-many matches.

The evaluation code and data are freely available to support the reproducibility of the results ⁵.

5.2 Blob matching results

Figure 5a shows the heat map depicting the mean Average Precision (mAP) of blob matching for different setups. The

5. <https://sites.google.com/view/fbellavia>

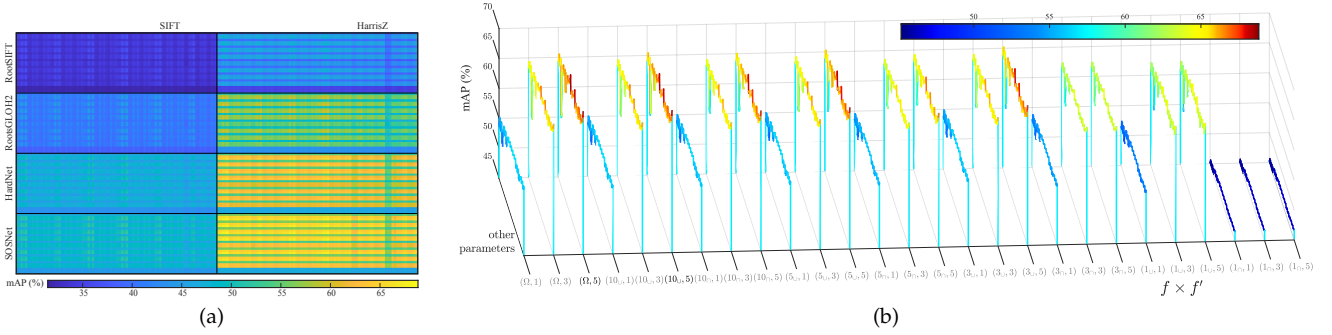


Fig. 5. (a) mAP for blob matching averaged on the image pairs of the planar and non-planar datasets for different setups and (b) for HarrisZ+SOSNet only, with ground-truth matches estimated by method D (see text for details, best viewed in color and zoomed in).

mAP values are averaged on all the image pairs considering the planar and non-planar datasets, with ground-truth estimated according to method D. For method C, based on the original patch overlap, a similar ranking has been obtained (not shown). Each rectangular area in Fig. 5a corresponds to a different detector+descriptor pair. Inside each rectangle, a row corresponds to a different $f \times f'$ pair with $f \in \{1_U, 1_\cap, 3_U, 3_\cap, 5_U, 5_\cap, 10_U, 10_\cap, \Omega\}$ and $f' \in \{1, 3, 5\}$ (see Fig. 1 for the definition of Ω). Likewise, each column represents the remaining parameter triplets $t_o \times \mathcal{D} \times \mathcal{W}$ with $t_o \in \{\infty, 50\%, 75\%, 99\%, 5 \text{ px}, 10 \text{ px}\}$, $\mathcal{D} = \{\mathcal{D}_\geq, \mathcal{D}_\leq, \mathcal{D}^+\}$ and $\mathcal{W}(a, b) \in \{a, b, \min(a, b), \max(a, b), (2ab)/(a+b)\}$. According to Fig. 5a, HarrisZ provides better mAP results than SIFT, probably due to the more strict keypoint selection criteria of HarrisZ with respect to SIFT. The average number of ground-truth matches is 847/701, 750/532 in the case of the planar/non-planar datasets for HarrisZ and SIFT, respectively. Moreover, confirming previous benchmarks SOSNet and HardNet provide the best accuracy results followed by RootsGLOH2 and RootSIFT.

Under the same blob matching setup, mAP correlation between different detector+descriptor pairs is high. Specifically, the correlation is more than 90% with the exception of SIFT+RootsGLOH2 with respect to any other detector+descriptor pairs, for which is yet higher than 60%. According to these observations, only the best pair HarrisZ+SOSNet is chosen for a more detailed analysis. Figure 5b plots the mAP values for the HarrisZ+SOSNet pair. The $f \times f'$ pairs $(\Omega, 5)$ and $(10_U, 5)$ are those providing the best mAP values (respectively 68.8% and 68.7%, highlighted in Fig. 5b), while the corresponding one-to-one matching setup $(\Omega, 1)$ and $(10_U, 1)$ obtain mAP values around 57%. This suggests that one-to-one matching can discard a lot of correct candidate matches. Moreover, the close results obtained by $f = \Omega$ and $f = 10_U$ indicate that it is sufficient to inspect only the first 10 top-ranked matches when designing a matching strategy. This observation can be exploited to improve the computational efficiency since many wrong matches can be discarded a priori.

Figure 6a plots the mAP values for the remaining blob matching parameters for the best configuration found so far. For reference, the corresponding plots in the case of SIFT+RootSIFT is reported in Fig. 6b. Among NNR-like similarities, \mathcal{D}^+ is generally the best choice. This is reasonably expected, since \mathcal{D}^+ is designed for many-to-many matches ($f' = 5$ in the evaluated configurations), while \mathcal{D}_\leq does not

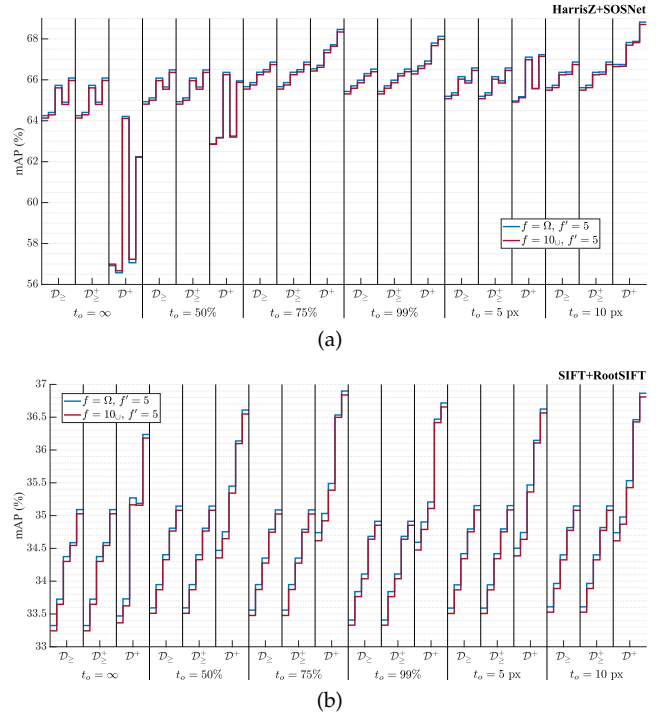


Fig. 6. Detailed mAP of blob matching with (a) HarrisZ+SOSNet and (b) SIFT+RootSIFT for different setups averaged on the whole dataset. $\mathcal{W}(a, b)$ values corresponding to a , b , $\min(a, b)$, $\max(a, b)$ and $(2ab)/(a+b)$ are reported in order inside each vertical sub-band (see text for details, best viewed in color and zoomed in).

provide any improvement with respect to \mathcal{D}_\geq . FGINN with $t_o = 75\%$ or $t_o = 10 \text{ px}$ provides substantial improvements with respect to the base setup with no FGINN ($t_o = \infty$), but only when the setup includes \mathcal{D}^+ . Inside each sub-band of the plots, mAP values are reported in order considering the different \mathcal{W} choices. While there is no evident difference using one image or another as reference, their combined distances improve the results. In particular, the minimum seems to achieve better results when FGINN has no or small ranges on the opposite of the maximum, and in any case the harmonic mean equals or surpasses the best among the previous \mathcal{W} choices. As observed in [4], the more the matches are discriminative, that happens when FGINN is employed with a sufficient range, the more their combination by union, which is equivalent to use the maximum, is better. The opposite holds for their intersection, which is equivalent to use the minimum. According to these evaluation the best

blob matching setup is $f = 10_{\cup}$, $f' = 5$, \mathcal{D}^+ with $t_o = 10$ px or $t_o = 75\%$, and $\mathcal{W}(a, b) = (2ab)/(a + b)$.

5.3 Delaunay Triangulation Matching results

The scatter plots of Fig. 7 refer to the average precision and recall values of the evaluated local spatial filters, without or with 1SAC as post-processing, on the planar, non-planar and SUN3D datasets. For the planar and non-planar datasets, results are reported by considering: The global behavior (Fig. 7a), i.e. averaging the results over all the considered detectors, descriptors and blob matching setups; the baseline configuration as reference (Fig. 7b), i.e. SIFT+RootSIFT with the one-to-one NNR greedy matching obtained by setting $f = \Omega$, $f' = 1$, $t_o = \infty$, \mathcal{D}_{\geq} , $\mathcal{W} = a$; the best configuration in terms of mAP (Fig. 7c), i.e. HarrisZ+SOSNet with blob matching setup $f = 10_{\cup}$, $f' = 5$, $t_o = 0.75$, \mathcal{D}^+ , $\mathcal{W} = 2ab/(a + b)$. For the SUN3D dataset (Fig. 7d), results are reported for the best configuration, also replacing HarrisZ with SIFT, and ACNe* indicates that ACNE is trained with the SUN3D indoor dataset instead of the YFCC100M outdoor dataset [33]. For the planar and non-planar datasets, detailed average statistics including the precision and recall, the number of correct and output matches, the number of times a method failed, and the running time are reported in Appendix E. The recall is computed by considering only ground-truth matches from the specific blob matching setup used in the pipeline. No precision/recall aggregated measures, such as mAP or F_{β} score are considered in the evaluation. On one hand, mAP requires that the number of output matches should be approximately the same for all methods since it is very sensitive to the recall, otherwise the highest scores would be assigned to the methods providing more output matches, including the initial blob matching. On the other hand, the choice of the β parameter in F_{β} can be questionable, as well as the choice of the recall normalization factor (see again Appendix E for further details), promoting one method or another without reflecting their effective performances. Nevertheless, the mAP and the F_1 and $F_{0.5}$ scores are reported for completeness in the additional material.

The average number of ground-truth matches per image pair, limited to blob matching only, for the planar/non-planar dataset are 551/429, 367/257, 681/575 for the global, reference and best configurations, and 963/691, 997/696, 928/686 when considering all the possible matches. According to the number of correct retrieved matches, the planar case is easier than the non-planar case. For SUN3D, the average number of ground-truth matches per image pair, limited to blob matching, is 481 for the best configuration and 229 when replacing HarrisZ with SIFT. Given the best configuration HarrisZ+SOSNet, the relative distribution of the spatial filters over the plots for the non-planar and SUN3D datasets is quite similar. By inspecting Fig. 7, all methods improve the precision with respect to the blob matching, that obviously achieve the highest recall.

Focusing on the results without 1SAC, the full-stage DTM (DTM₂) provides high levels of precision and recall, comparable with OANet, ACNe and PGM. OANet obtains somewhat better precision than the other match filters. Notice also that ACNe trained on the same kind of images to

be processed (ACNe* on SUN3D) achieves boosted performances, being or not the training set included (see Sec. 5.1 and the additional material). The precision of the first stage of DTM alone (DTM₁) is similar to that of the complete DTM, but the recall is lower, underlining the goodness of the DTM₂ stage. LPM and LMR obtain recall values equivalent to those of the previously mentioned spatial filters, but with lower precision in the baseline configuration, while they regain in terms of precision at expense of the recall in the best configuration. GPLM behaves likewise LPM but with better recall in the best configuration. AdaLAM achieves very high precision but lower recall with respect to the previous spatial filters, unless running on planar scenes. Simple thresholding (th) gets precision values similar to AdaLAM but with lower recall. BM recall and precision are in-between those of AdaLAM and simple thresholding, while PFM obtains a lower recall than the simple thresholding for a similar precision. GMS improves upon the simple thresholding in the case of the best configuration, which is probably closer to the original setup GMS was designed for in terms of the number of input matches and keypoint type. VFC gets very high recall and a reasonable precision on the planar dataset, while for the other datasets it achieves still very high recall but low precision. RFM-SCAN is generally equal or worse than VFC. Finally, LLT and SVC obtained in this evaluation the worst results.

1SAC, without being a full RANSAC, is able to improve the precision with an acceptable loss in terms of recall. This effect decreases when simultaneously the precision is high and the recall is low, such as for AdaLAM. It is reasonable that 1SAC post-process does not affect AdaLAM, which consists of multiple local RANSACs. Moreover, 1SAC generally reduces the gain in terms of precision of OANet with respect to DTM, PGM and ACNe. Notice that OANet and ACNe network design and training takes more or less explicitly into account the same global model constraints based on the epipolar geometry of 1SAC, not considered in the design of the other spatial filters. Blob matching after 1SAC post-filtering achieves almost the best results in terms of both precision and recall in the planar case. This holds because planar images are relatively easy, homographies provide one-to-one point maps between images, and 1SAC re-filters the output of the local spatial filters. These observations indicate that homographies are correctly estimated in any case, so that the highest recall is obtained as more input matches are provided. With 1SAC, going from the baseline to the best configuration, a general expansion towards the top-right area of scatter plots can be observed. This agrees with the fact that the absolute number of correct matches gets roughly doubled while the total number of output matches becomes five times (f' changes from 1 to 5) as moving from the baseline to the best configuration, so that top-ranked matches employed by 1SAC becomes less contaminated by outliers⁶.

Blob matching precision of each plot gives also an indication of the average inlier ratio, so that SIFT+SOSNet on SUN3D (about 8% of average inlier ratio according to

6. Note that in general the average precision is not equal to the ratio between the average number of correct matches and the average number of output matches.

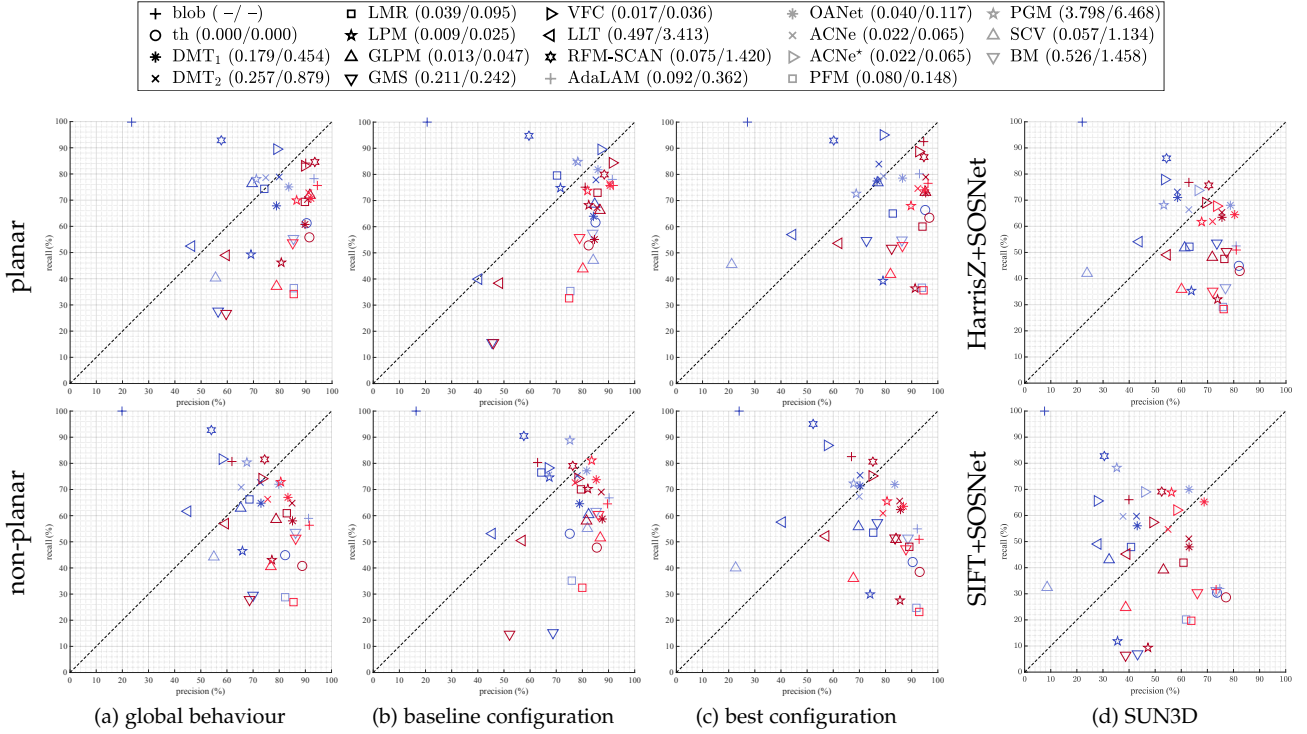


Fig. 7. Average precision and recall values of the local spatial filters on the planar, non-planar and SUN3D datasets. The recall is computed with respect to ground truth matches obtained by blob matching, results without/with 1SAC are in blue/red. Average running times (s) for the baseline/best configuration are reported alongside in the legend (see text for details, best viewed in color and zoomed in).

blob matching) is more contaminated by outliers than the other presented plots. On this setup, OANet followed by ACNe* suffer less the outlier contamination. Nevertheless, their differences with the other spatial filters after 1SAC are quite reduced since they implicitly contains epipolar constraints, unlike other methods. Under this consideration DTM neighborhoods based only on triangulation without considering any form of consistency based on patch relative local transformations are quite robust. More detailed histograms supporting these observations according to the inlier ratio can be found in Appendix E.

The number of times the local spatial filters failed to get at least one correct match as output, reported on Appendix E, can provide further clues about the robustness of each match filter. For both planar or non-planar scenes, excluding obviously the blob matching, DTM, OANet, ACNe and PGM are among those methods which fail less. Notice also that when 1SAC is applied more failures arise, due to the reduced number of the candidate matches.

A visual qualitative analysis on the best configuration according to the examples reported in Appendix E agrees with the quantitative results discussed above. It can be noted that for DTM, with respect to other spatial filters such as ACNe and PGM, the wrong matches often concern the image regions near the triangulation boundary, which lacks a neighborhood covering in all directions. These wrong matches are generally removed by 1SAC which, unlike the case of the other spatial filters, seems to remove less inliers in the case of DTM, maybe due to the kind of outliers.

Average running times are reported in Fig. 7 alongside the legend for the baseline/best configurations, implying respectively one-to-one or many-to-many matching relations. Results have been obtained with Ubuntu 20.04 running on

an Intel Core I9 10900K with 64 GB of RAM equipped with a NVIDIA GeForce RTX 2080 Ti GPU. The original code was used for each implementation, only modified to works with blob matches as input. The code is implemented in Python or Matlab with different level of optimization, ranging to Matlab mex C functions to GPU parallel optimization. From the baseline to the best configuration the running time increases proportionally to the number of processed matches by at least a linear factor, theoretically expected to be around $\sum_{i=1}^{f'} \frac{i}{f'} = 3$ since $f' = 5'$ for the many-to-many match in the best configuration. The code of PGM, BM and LLT is the slowest, since implemented in Matlab with almost no optimization. DTM and GMS, respectively written in Matlab and Python without any type of optimization, follow in the list, and the faster code of the remaining spatial filters contains several optimizations. Better running times are expected for DTM after code optimization since the various steps of DTM can be highly parallelized and, in the same manner of LPM, portions of DTM code would benefit of mex C code rewriting⁷.

6 CONCLUSIONS AND FUTURE WORK

This paper analyzes the problem of the local image descriptor matching according to two possible contexts that can be used to characterize the images, and to improve the number of correct correspondences.

The first context is provided by the descriptor space. The novel general blob matching strategy is designed to incorporate different approaches for a clear and detailed analysis of the aspects that characterize the basic matching

⁷ According to the author's experience Matlab is not the best environment for graph manipulations.

strategies. According to the evaluation, pre-filtering, many-to-many matches, two-way comparisons using symmetric distances and a good choice of the second best match in NNR can improve the matching process.

The second context is provided by the keypoint space, i.e. the actual image space. A new local spatial filter named DTM is proposed. DTM extracts spatial neighborhood relations between keypoints by Delaunay triangulation, alternating triangulation contractions and expansions to remove inconsistent matches and to include consistent matches. DTM is robust and obtains comparable or better results than the state-of-the-art. Furthermore, DTM neighborhoods do not rely on parameters to be defined, but are implicitly derived by the keypoint distribution onto the images. Moreover, DTM does not require patch relative local transformations for validate the neighborhood consistency.

Although blob matching and DTM mainly operate respectively on the descriptor and space contexts, they both betray contaminations from the opposite contexts, underlining the need of fully integrating different contexts to go beyond the current state-of-the-art. To be noted that blob matching and DTM have been employed recently as part of a very competitive matching pipeline which achieved among the best results in the recent Image Matching Challenge 2021 (IMC2021) and SimLocMatch contests [38].

Finally, a comprehensive evaluation of the main phases of the matching pipeline is carried out, based on a new benchmark, focusing on the estimation of correct matches and not on their effects on the scene. The analysis considers both blob and corner like keypoints, among the current best local image descriptors, several image matching strategies, state-of-the-art local spatial filters, and also the simple model-based filter. It clearly emerges that the combination of the different methods can offer a clear advantage with respect to the baseline SIFT matching strategy.

As future work, further possibilities for merging matching strategies will be investigated, as well as mesh-based applications of the triangulation in order to grow up matches and obtain semi-dense correspondences. Additionally, it would be interesting to analyze how triangulation-based neighborhoods can be used for clustering and spatially characterizing the objects in the scene. Further research directions will be also aimed at improving the benchmark, by extending the datasets with more scenes and by designing better error metrics to compare the different approaches.

ACKNOWLEDGMENT

The Titan Xp used for this research was generously donated by the NVIDIA Corporation.

This research is funded by the Italian Ministry of Education and Research (MIUR) under the program PON Ricerca e Innovazione 2014-2020, cofunded by the European Social Fund (ESF), CUP B74I18000220006, id. proposta AIM 1875400, linea di attività 2, Area Cultural Heritage.

REFERENCES

- [1] R. Szeliski, *Computer Vision: Algorithms and Applications*, 2nd edition. Springer, 2021.
- [2] D. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [3] F. Bellavia and C. Colombo, "Is there anything new to say about SIFT matching?" *International Journal of Computer Vision*, vol. 128, pp. 1847–1866, 2020.
- [4] Y. Jin, D. Mishkin, A. Mishchuk, J. Matas, P. Fua, K. M. Yi, and E. Trulls, "Image matching across wide baselines: From paper to practice," *International Journal of Computer Vision*, vol. 129, pp. 517–547, 2021.
- [5] R. Arandjelović and A. Zisserman, "Three things everyone should know to improve object retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 2911–2918.
- [6] F. Bellavia and C. Colombo, "Rethinking the sGLOH descriptor," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 931–944, 2018.
- [7] —, "Which is Which? Evaluation of local descriptors for image matching in real-world scenarios," in *International Conference on Computer Analysis of Images and Patterns (CAIP)*, 2019. [Online]. Available: <http://cvg.dsi.unifi.it/wisw.caip2019>
- [8] J. L. Schönberger and J. M. Frahm, "Structure-from-Motion revisited," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [9] Z. Luo, T. Shen, L. Zhou, J. Zhang, Y. Yao, S. Li, T. Fang, and L. Quan, "ContextDesc: Local descriptor augmentation with cross-modality context," *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [10] V. Balntas, E. Riba, D. Ponsa, and K. Mikolajczyk, "Learning local feature descriptors with triplets and shallow convolutional neural networks," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2016, pp. 119.1–119.11.
- [11] A. Mishchuk, D. Mishkin, F. Radenovic, and J. Matas, "Working hard to know your neighbor's margins: Local descriptor learning loss," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems (NIPS)*, 2017, pp. 4829–4840.
- [12] Z. Luo, T. Shen, L. Zhou, S. Zhu, R. Zhang, Y. Yao, T. Fang, and L. Quan, "Geodesc: Learning local descriptors by integrating geometry constraints," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [13] Y. Tian, X. Yu, B. Fan, F. Wu, H. Heijnen, and V. Balntas, "SOSNet: Second order similarity regularization for local descriptor learning," in *CVPR*, 2019.
- [14] J. Xiao, A. Owens, and A. Torralba, "SUN3D: A database of big spaces reconstructed using sfm and object labels," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2013, pp. 1625–1632.
- [15] V. Balntas, K. Lenc, A. Vedaldi, and K. Mikolajczyk, "HPatches: A benchmark and evaluation of handcrafted and learned local descriptors," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3852–3861.
- [16] Z. Li and N. Snavely, "MegaDepth: Learning single-view depth prediction from internet photos," in *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [17] C. Zhao, Z. Cao, J. Yang, K. Xian, and X. Li, "Image feature correspondence selection: A comparative study and a new contribution," *IEEE Transactions on Image Processing*, vol. 29, pp. 3506–3519, 2020.
- [18] J. Ma, J. Jiang, A. Fan, J. Jiang, and J. Yan, "Image matching from handcrafted to deep features: A survey," *International Journal of Computer Vision*, vol. 129, pp. 23–79, 2021.
- [19] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2000.
- [20] M. Fischler and R. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [21] O. Chum and J. Matas, "Matching with PROSAC - progressive sample consensus," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005, pp. 220–226.
- [22] T. Sattler, B. Leibe, and L. Kobbelt, "SCRAMSAC: Improving RANSAC's efficiency with a spatial consistency filter," in *Proceedings of the IEEE International Conference on Computer Vision (ICPR)*, 2009, pp. 2090–2097.
- [23] N. Ni, J. Hailin, and F. Dellaert, "GroupSAC: Efficient consensus in the presence of groupings," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2009, pp. 2193–2200.

- [24] D. Barath and J. Matas, "Graph-Cut RANSAC," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6733–6741.
- [25] L. Cavalli, V. Larsson, M. R. Oswald, T. Sattler, and M. Pollefeys, "Handcrafted outlier detection revisited," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 770–787.
- [26] J. Ma, J. Zhao, J. Jiang, H. Zhou, and X. Guo, "Locality preserving matching," *International Journal of Computer Vision*, vol. 127, no. 5, pp. 512–531, May 2019.
- [27] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018.
- [28] M. J. Tyszkiewicz, P. Fua, and E. Trulls, "DISK: Learning local features with policy gradient," in *Proceedings of the 32nd Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [29] P. E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "Super-Glue: Learning feature matching with graph neural networks," in *arXiv*, 2019.
- [30] Y. Tian, V. Balntas, T. Ng, A. B. Laguna, Y. Demiris, and K. Mikolajczyk, "D2D: keypoint extraction with describe to detect approach," in *Proceedings of the 15th Asian Conference on Computer Vision (ACCV)*, 2020.
- [31] K. M. Yi, E. Trulls, Y. Ono, V. Lepetit, M. Salzmann, and P. Fua, "Learning to find good correspondences," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [32] C. Zhao, Z. Cao, C. Li, X. Li, and J. Yang, "NM-Net: Mining reliable neighbors for robust feature correspondences," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 215–224.
- [33] J. Zhang, D. Sun, Z. Luo, A. Yao, L. Zhou, T. Shen, Y. Chen, H. Liao, and L. Quan, "Learning two-view correspondences and geometry using order-aware network," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 5844–5853.
- [34] W. Sun, W. Jiang, E. Trulls, A. Tagliasacchi, and K. M. Yi, "ACNe: Attentive context normalization for robust permutation-equivariant learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [35] J. W. Bian, Y. H. Wu, J. Zhao, Y. Liu, L. Zhang, M. M. Cheng, and I. Reid, "An evaluation of feature matchers for fundamental matrix estimation," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2019.
- [36] A. Barroso-Laguna, E. Riba, D. Ponsa, and K. Mikolajczyk, "KeyNet: Keypoint detection by handcrafted and learned CNN filters," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
- [37] F. Bellavia, D. Tegolo, and C. Valenti, "Improving Harris corner selection strategy," *IET Computer Vision*, vol. 5, no. 2, pp. 86–96, 2011.
- [38] "Image Matching Workshop (IMW) challenge at (CVPR2021)," 2021. [Online]. Available: <https://image-matching-workshop.github.io>
- [39] F. Bellavia and C. Colombo, "RootsGLOH2: embedding RootSIFT "square rooting" in sGLOH2," *IET Computer Vision*, 2020.
- [40] M. Pultar, D. Mishkin, and J. Matas, "Leveraging outdoor webcams for local descriptor learning," in *Proceedings of Computer Vision Winter Workshop (CVWW) 2019*, 2019.
- [41] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: an efficient alternative to SIFT or SURF," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2011, pp. 2564–2571.
- [42] K. Yi, Y. Verdie, P. Fua, and V. Lepetit, "Learning to assign orientations to feature points," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 1–8.
- [43] 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.
- [44] B. Fan, F. Wu, and Z. Hu, "Rotationally invariant descriptors using intensity order pooling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 10, pp. 2031–2045, 2012.
- [45] K. Lenc, J. Matas, and D. Mishkin, "A few things one should know about feature extraction, description and matching," in *Proceedings of the Computer Vision Winter Workshop (CVWW)*, 2014, pp. 67–74.
- [46] W. Zhang and J. Kosecka, "Generalized RANSAC framework for relaxed correspondence problems," in *Proceeding of the International Symposium on 3D Data Processing, Visualization, and Transmission (3DPVT)*, 2006, pp. 854–860.
- [47] D. Mishkin, J. Matas, and M. Perdoch, "MODS: Fast and robust method for two-view matching," *Computer Vision and Image Understanding*, 2015.
- [48] D. Barath and J. Matas, "Progressive-X: Efficient, anytime, multi-model fitting algorithm," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [49] V. Ferrari, T. Tuytelaars, and Luc Val Gool, "Wide-baseline multiple-view correspondences," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2003.
- [50] J. W. Bian, W. Y. Lin, Y. Liu, L. Zhang, S. K. Yeung, M. M. Cheng, and I. Reid, "GMS: Grid-based motion statistics for fast, ultra-robust feature correspondence," *International Journal of Computer Vision*, vol. 128, pp. 1580–1593, 2020.
- [51] J. Ma, J. Jiang, H. Zhou, J. Zhao, and X. Guo, "Guided locality preserving feature matching for remote sensing image registration," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 8, pp. 4435–4447, 2018.
- [52] J. Ma, X. Jiang, J. Jiang, J. Zhao, and X. Guo, "LMR: Learning a two-class classifier for mismatch removal," *IEEE Transactions on Image Processing*, vol. 28, no. 8, pp. 4045–4059, 2019.
- [53] S. Lee, J. Lim, and I. H. Suh, "Progressive feature matching: Incremental graph construction and optimization," *IEEE Transactions on Image Processing*, vol. 29, pp. 6992–7005, 2020.
- [54] M. Cho and K. M. Lee, "Progressive graph matching: Making a move of graphs via probabilistic voting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 398–405.
- [55] A. Geiger, J. Ziegler, and C. Stiller, "Stereoscan: Dense 3D reconstruction in real-time," in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2011.
- [56] J. Ma, H. Zhou, J. Zhao, Y. Gao, J. Jiang, and J. Tian, "Robust feature matching for remote sensing image registration via locally linear transforming," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 12, pp. 6469–6481, 2015.
- [57] J. Ma, J. Zhao, J. Tian, A. L. Yuille, and Z. Tu, "Robust point matching via vector field consensus," *IEEE Transactions on Image Processing*, vol. 23, no. 4, pp. 1706–1721, 2014.
- [58] W. Y. D. Lin, M. M. Cheng, J. Lu, H. Yang, M. N. Do, and P. Torr, "Bilateral functions for global motion modeling," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 341–356.
- [59] J. Cech, J. Matas, and M. Perdoch, "Efficient sequential correspondence selection by cosegmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1568–1581, 2010.
- [60] X. Jiang, J. Ma, J. Jiang, and X. Guo, "Robust feature matching using spatial clustering with heavy outliers," *IEEE Transactions on Image Processing*, vol. 29, pp. 736–746, 2020.
- [61] A. Albarelli, E. Rodolà, and A. Torsello, "Imposing semi-local geometric constraints for accurate correspondences selection in structure from motion: A game-theoretic perspective," *International Journal of Computer Vision*, vol. 92, pp. 36–53, 2012.
- [62] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, "LoFTR: Detector-free local feature matching with transformers," 2021.
- [63] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615–1630, 2005.
- [64] D. Tegolo and F. Bellavia, "noRANSAC for fundamental matrix estimation," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2011.
- [65] Z. Zhang, "Determining the epipolar geometry and its uncertainty: A review," *International Journal of Computer Vision*, vol. 27, no. 2, pp. 161–195, 1998.
- [66] P. Moreels and P. Perona, "Evaluation of features detectors and descriptors based on 3d objects," *International Journal of Computer Vision*, vol. 73, no. 3, pp. 263–284, 2007.

APPENDIX A DTM BORDER COMPUTATION

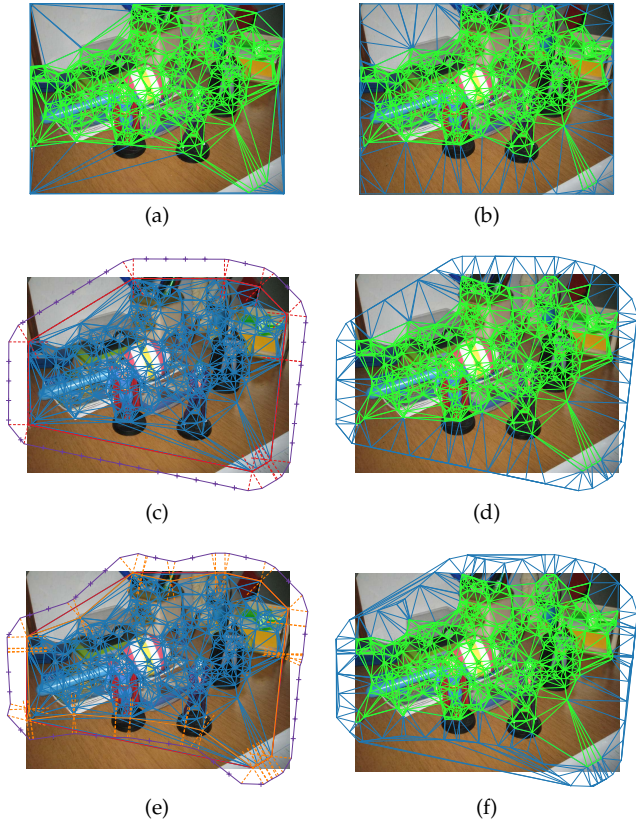


Fig. 8. Delaunay triangulations by including image corners (a) and by breaking the canvas borders into multiple lines (b). (c) Convex hull (solid red) fattening (dashed red) in order to obtain expanded contour edges (purple) with the resulting Delaunay triangulation (d). (e) Analogous results using alpha-shape boundary edges (orange) and (f) the final Delaunay triangulation employed in DTM. Green Delaunay edges are used to distinguish the edges connecting keypoints from the others (see Sec. 4 for details, best view in color and zoomed in).

APPENDIX B NEIGHBORHOOD FORMULATION DIFFERENCES

Figure 10 shows the intuition behind the different neighborhood formulations, discussed in the main text, on the toy example scene of Fig. 9. Assuming that dots are detected as keypoints, the aim is to check whether for the foreground object the keypoint neighborhood contains at least another foreground keypoints, i.e. belonging to the same motion field cluster. As neighborhood parameters, the circular radius is set equal to $1/4$ of the foreground square side in the frontal view, while $k = 3$. These parameters are set empirically. It can be noted that the Delaunay-based neighborhood allows to have more neighborhoods containing keypoints of the same motion field cluster than the other neighborhood formulations in all the views, avoiding to setup parameters in the neighborhood design. Clearly, background keypoints get included, but their number decreases when considering the intersection of the corresponding neighborhoods among the images, as shown in Fig. 11. Notice that apart from the Delaunay-based neighborhood, intersection causes a great

drop of neighborhoods within the same motion field cluster. Moreover, patch-based local similarity or affine transformations for this scene is very unlikely to be useful in checking consistency.

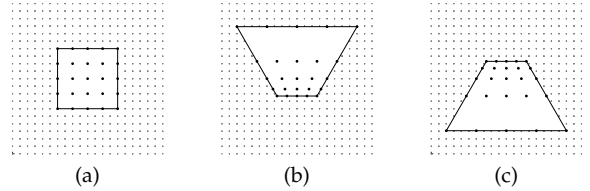


Fig. 9. (a) Frontal view of a toy example scene with a static background (gray dots) and a planar surface as foreground (black edges and dots). Other views of the scene with the foreground plane slanted (b) upward and (c) downward (best view in color and zoomed in).

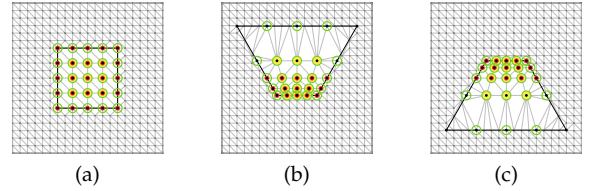


Fig. 10. Neighborhoods for the foreground keypoints of each scene view of Fig. 9 containing at least another foreground keypoint, i.e. sharing the motion field cluster. These neighborhoods are indicated respectively in red, yellow and green for the circular radius neighborhood, the closest k nearest neighborhood and the Delaunay-based neighborhood. Delaunay triangulation edges are in gray (best view in color and zoomed in).



Fig. 11. The intersections of the corresponding neighborhoods between (a) Fig. 10a and Fig. 10b, and (b) Fig. 10b and Fig. 10c on the first view according to Fig. 10 notation (best view in color and zoomed in).

APPENDIX C EVALUATION DATASET

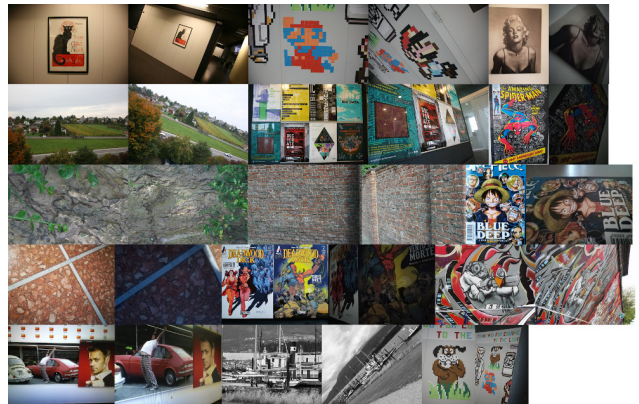


Fig. 12. Exemplar image thumbnails for each scene of the planar dataset. Each scene is made up of six images, only two of these are shown (see Sec. 5.1 for details, best viewed in color and zoomed in).



Fig. 13. Exemplar image thumbnails for each scene of the non-planar dataset. Each scene is made up of two or three images, only two of these are shown (see Sec. 5.1 for details, best viewed in color and zoomed in).

APPENDIX D GROUND TRUTH ESTIMATION METHODS

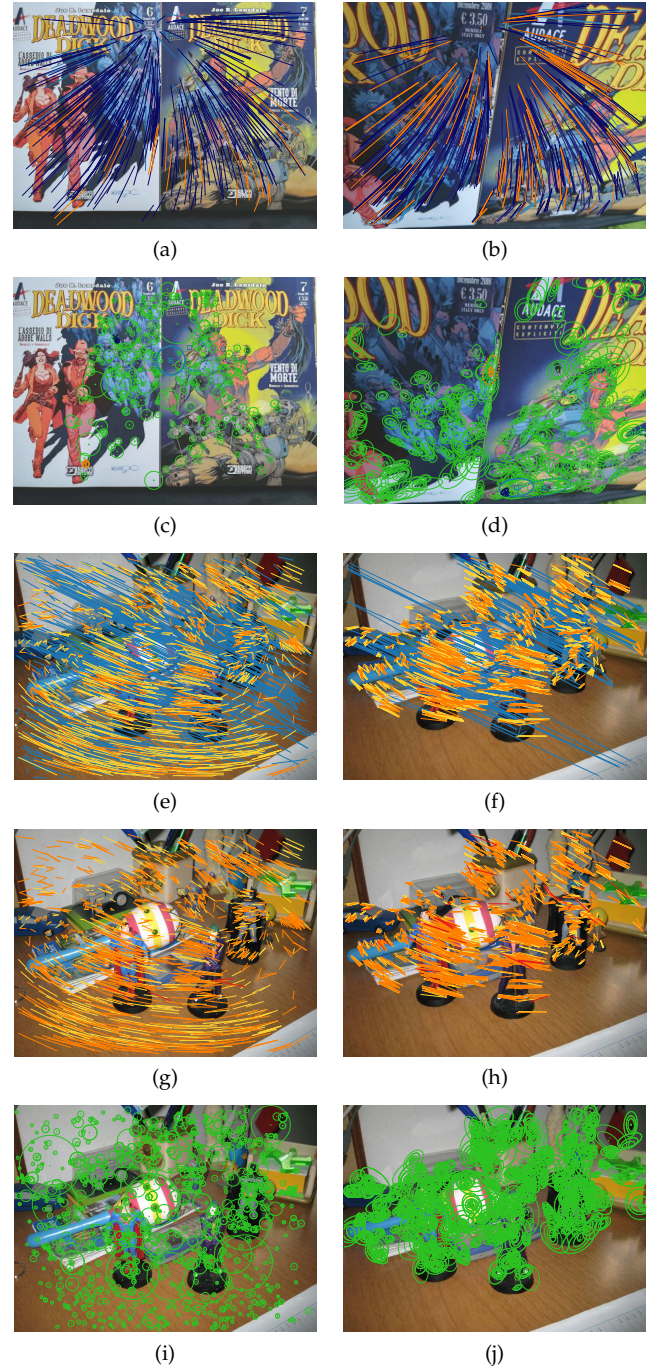


Fig. 14. Clusters of motion flows representing matches and keypoint support regions superimposed on planar and non-planar image pairs according to different ground-truth estimation methods when using SIFT (left column) and HarrisZ (right column). (first row) Method B (blue and light blue), method C (blue) and method D (blue and orange). (second row) Descriptor support regions for each keypoint match found in the previous row. The worst matches according to the reprojection error for each colored cluster are evidenced, the dashed ellipses are reprojected from the other image and the overlap error is computed by doubling the shape size. (third row) Method A (light blue and yellow) by setting the epipolar line error distance threshold to 7 pixels and method D (yellow and orange). (fourth row) Method B (yellow, red and light red), C (yellow and red) and D (yellow and orange). (fifth row) Descriptor support regions for each correspondence found by methods C and D (see Sec. 5.1, best viewed in color and zoomed in).

APPENDIX E

DTM EVALUATION

Tables 1-3 report the results of the local spatial filters evaluation for the planar and non-planar datasets, without or with 1SAC applied as post-processing. The following statistics, averaged on the image pairs, are included: The precision, the recall, the number of correct and output matches, the number failures (i.e. the number of times the spatial filter did not provide at least one ground-truth match), and the running time. The recall is computed with respect to the ground-truth matches of the specific pipelined blob matching. The tables also include recall*, that considers overall the ground-truth matches among all the tested configurations. For the comparison, the former recall definition is more reasonable, since the local spatial filters cannot include matches outside those provided to them as input, i.e. the blob matching output. Nevertheless, recall* can give clues of the whole pipeline performances. More detailed data can be found in the additional material.

Figures 15-22 provide a visual qualitative comparison of the local spatial filters evaluated in Sec. 5.3, which agrees with the quantitative analysis presented in the main text. Clearly, there is no method that gives always, and with respect to all the evaluated criteria, the best results. Anyways, DTM appears to be quite robust and generally provides good results, comparable with the best spatial filters both in terms of precision and recall. Notice also that, as reported in Sec. 4, DTM does not require to set parameters to define the neighborhoods or the criteria for the selection of the initial and output matches. From the visual comparison, with respect to other spatial filters such as ACNe and PGM, DTM wrong matches (red and light red) often concern image areas near the triangulation boundary, which lacks a neighborhood covering in all directions. These wrong matches are generally removed by 1SAC (light red), that unlike the case of the other spatial filters, for DTM seems to remove less inliers (yellow), maybe due to the kind of outliers. In order to be robust, inliers should be able to survive to 1SAC (green) and to be well distributed on the image area. These observations can better appreciated for instance in the LeuvenB, Fountain01, Valbonne, Harvard Restroom and MIT Office image pairs.

Figure 23 shows the precision and recall histograms according to different inlier ranges for the corresponding scatter plots of Fig. 7. In addition, further histograms indicating the proportion of image pairs having inliers within a specific range for each setup are reported. From the analysis of these histograms, the same observations discussed in Sec. 5.3 can be drawn out. In particular, it appears that OANet and ACNe suffer less than DTM of outlier contamination. Nevertheless, these differences after 1SAC are reduced. OANet and ACNe implicitly contains epipolar constraints which are forced into DTM by 1SAC. Under this consideration DTM appears still quite robust.

TABLE 1
Global statistics for the local spatial filters. The average planar/non-planar number of ground-truth matches per image pair is 457/382 (blob, for recall), 792/602 (all matches, for recall*)

		precision (%)	recall (%)	recall* (%)	correct matches	output matches	failures	time (s)	
planar	no model	blob	23.48	99.79	56.23	457	3806	0	-
		th	90.29	61.21	36.49	310	370	0	0.000
		DTM ₁	78.78	67.84	40.05	356	593	1	0.283
		DTM ₂	79.82	78.85	45.84	399	651	0	0.519
		LMR	74.14	74.32	43.11	384	588	2	0.065
		LPM	69.05	49.26	27.99	259	362	4	0.017
		GLPM	69.48	76.34	45.46	398	952	3	0.029
		GMS	56.51	27.62	18.37	200	365	28	0.253
		VFC	78.91	89.48	51.52	440	781	4	0.031
		LLT	46.39	52.47	31.63	298	591	28	1.665
		RFM-SCAN	57.71	92.84	52.65	438	1237	0	0.774
		AdaLAM	93.11	78.23	46.40	395	542	3	0.226
		OANet	83.47	75.06	42.22	362	552	1	0.082
		ACNe	74.69	78.61	45.33	385	638	0	0.047
	PFM	85.40	36.39	21.79	180	195	8	0.114	
	PGM	71.10	78.06	43.78	341	486	6	5.051	
	SCV	55.52	40.30	24.35	240	916	4	0.658	
	BM	85.27	55.35	33.06	328	445	9	0.892	
	1SAC	blob	89.94	84.15	49.54	418	620	4	-
		th	91.35	55.81	33.76	294	340	5	-
		DTM ₁	89.68	60.72	36.57	333	482	5	-
		DTM ₂	90.61	70.26	41.72	372	529	5	-
		LMR	89.64	69.35	40.69	366	475	6	-
		LPM	80.62	46.18	26.51	248	311	12	-
		GLPM	91.54	71.87	42.98	378	534	4	-
		GMS	59.58	26.73	17.79	193	316	29	-
		VFC	89.56	83.11	48.39	417	614	5	-
		LLT	59.60	48.94	29.77	281	408	28	-
RFM-SCAN		93.43	84.57	49.13	414	612	2	-	
AdaLAM		94.41	75.59	44.86	380	513	3	-	
OANet		91.95	70.67	40.21	346	472	4	-	
ACNe		91.10	72.79	42.67	366	493	5	-	
PFM	85.38	34.15	20.63	174	186	10	-		
PGM	86.48	69.92	40.02	318	365	7	-		
SCV	78.80	37.10	22.55	228	277	14	-		
BM	85.00	53.69	32.11	316	419	10	-		
non-planar	no model	blob	19.86	100.00	60.99	382	3834	0	-
		th	82.26	44.89	28.71	185	235	1	0.000
		DTM ₁	73.01	64.77	41.34	269	501	0	0.271
		DTM ₂	72.76	72.71	45.90	298	546	0	0.512
		LMR	68.60	66.21	40.85	269	471	1	0.067
		LPM	65.92	46.48	26.73	177	264	3	0.017
		GLPM	65.26	62.86	40.62	271	644	4	0.018
		GMS	69.98	29.54	20.40	161	304	21	0.251
		VFC	58.20	81.63	51.75	345	897	6	0.037
		LLT	45.01	61.68	40.20	279	737	23	1.602
		RFM-SCAN	54.03	92.71	57.27	365	1145	0	0.735
		AdaLAM	91.21	58.95	37.50	254	339	3	0.165
		OANet	79.83	72.00	44.49	282	468	1	0.083
		ACNe	65.47	70.86	43.98	281	564	0	0.047
	PFM	82.27	28.81	18.46	112	125	8	0.100	
	PGM	67.62	80.37	50.02	312	539	1	7.119	
	SCV	54.99	44.20	28.35	186	812	2	0.729	
	BM	86.24	53.59	35.00	244	351	5	0.729	
	1SAC	blob	61.91	80.71	51.25	334	740	1	-
		th	88.80	40.77	26.46	172	204	2	-
		DTM ₁	85.08	58.01	37.57	248	394	2	-
		DTM ₂	84.92	64.78	41.50	273	427	1	-
		LMR	82.89	60.92	38.00	253	351	3	-
		LPM	77.19	43.01	24.93	167	218	6	-
		GLPM	78.79	58.55	37.98	255	410	5	-
		GMS	68.63	27.88	19.28	153	267	25	-
		VFC	73.39	74.14	47.31	318	573	7	-
		LLT	59.41	56.94	37.28	260	463	24	-
RFM-SCAN		74.42	81.51	51.33	333	631	1	-	
AdaLAM		91.50	56.32	35.85	243	322	3	-	
OANet		83.33	66.98	41.65	267	409	1	-	
ACNe		75.58	66.27	41.41	267	440	1	-	
PFM	85.46	26.93	17.41	106	116	8	-		
PGM	80.59	72.88	45.76	290	402	2	-		
SCV	76.81	40.53	26.25	173	265	4	-		
BM	86.28	51.35	33.56	234	325	7	-		

TABLE 2

Local spatial filter statistics for the baseline configuration. The average planar/non-planar number of ground-truth matches per image pair is 293/214 (blob, for recall), 740/519 (all matches, for recall*)

		precision (%)	recall (%)	recall* (%)	correct matches	output matches	failures	time (s)			
planar	no model	blob	20.63	100.00	37.60	293	1329	0	-		
		th	85.02	61.62	25.99	212	227	1	0.000		
		DTM ₁	84.23	63.99	25.69	228	244	2	0.168		
		DTM ₂	85.19	77.91	31.02	265	284	2	0.240		
		LMR	70.29	79.58	31.90	274	320	6	0.036		
		LPM	71.61	74.84	29.54	266	311	3	0.008		
		GLPM	84.71	68.62	29.37	258	275	7	0.014		
		GMS	45.48	15.69	7.31	92	94	40	0.260		
		VFC	87.12	89.57	34.56	282	298	4	0.019		
		LLT	40.34	39.94	16.89	190	227	37	0.485		
	ISAC	RFM-SCAN	59.56	94.78	35.96	283	482	0	0.075		
		AdaLAM	91.45	78.08	32.41	266	273	5	0.090		
		OANet	85.96	81.80	30.78	248	267	2	0.043		
		ACNe	77.48	84.92	33.13	266	313	0	0.023		
		PFM	75.38	35.38	14.92	122	127	15	0.076		
		PGM	78.24	84.82	32.65	259	283	7	2.759		
		SCV	84.15	47.13	19.92	198	209	6	0.445		
		BM	83.91	57.48	24.36	238	245	11	0.532		
		non-planar	no model	blob	16.41	100.00	38.91	214	1370	0	-
				th	75.17	53.09	22.95	130	155	2	0.000
DTM ₁	78.93			64.57	27.04	155	179	2	0.166		
DTM ₂	78.08			75.28	31.30	178	209	2	0.241		
LMR	64.29			76.50	31.95	182	257	4	0.040		
LPM	67.36			74.63	31.33	180	240	2	0.008		
GLPM	82.37			60.41	27.19	161	180	4	0.010		
GMS	68.75			15.22	7.54	55	56	30	0.266		
VFC	66.94			78.26	32.25	189	243	7	0.018		
LLT	45.50			53.14	23.45	154	227	28	0.531		
ISAC	RFM-SCAN		57.61	90.48	35.60	200	370	0	0.070		
	AdaLAM		90.24	66.81	28.78	168	176	4	0.064		
	OANet		81.66	77.16	30.30	167	192	1	0.043		
	ACNe		67.16	75.56	30.78	175	242	1	0.023		
	PFM		75.95	35.13	15.53	83	91	12	0.075		
	PGM		75.23	88.81	36.34	204	256	2	3.732		
	SCV		81.93	55.10	23.72	136	156	2	4.555		
	BM		85.25	61.61	26.68	163	177	6	0.424		
	non-planar		no model	blob	62.87	80.34	33.83	194	269	2	-
				th	85.52	47.78	21.33	124	132	4	-
DTM ₁		87.61		58.72	25.18	147	157	2	-		
DTM ₂		87.23		69.04	29.26	169	182	2	-		
LMR		79.46		70.05	29.97	174	191	5	-		
LPM		82.10		70.23	29.89	174	192	3	-		
GLPM		81.59		57.97	26.31	157	167	11	-		
GMS		52.14		14.63	7.31	54	55	47	-		
VFC		78.29		74.22	30.87	181	201	8	-		
LLT		56.95		50.47	22.39	148	169	29	-		
ISAC		RFM-SCAN	76.32	79.06	32.86	189	225	1	-		
		AdaLAM	89.59	64.51	27.85	164	171	5	-		
		OANet	85.30	73.76	29.32	163	178	1	-		
		ACNe	77.22	72.62	29.86	170	202	1	-		
		PFM	79.97	32.42	14.64	79	84	13	-		
		PGM	83.54	81.13	33.78	191	214	3	-		
		SCV	86.84	51.50	22.49	130	139	3	-		
		BM	86.15	60.41	26.25	160	170	7	-		

TABLE 3

Local spatial filter statistics for the baseline configuration. The average planar/non-planar number of ground-truth matches per image pair are 632/572 (blob, for recall), 844/685 (all matches, for recall*)

		precision (%)	recall (%)	recall* (%)	correct matches	output matches	failures	time (s)			
planar	no model	blob	27.18	100.00	74.59	632	6081	0	-		
		th	95.17	66.38	50.93	449	758	0	0.000		
		DTM ₁	76.75	77.48	59.05	526	1510	0	0.477		
		DTM ₂	77.45	83.94	63.57	561	1592	0	0.902		
		LMR	82.77	65.04	49.39	448	929	0	0.096		
		LPM	79.02	39.27	30.69	288	684	0	0.026		
		GLPM	77.05	76.81	58.77	523	1518	0	0.062		
		GMS	72.71	54.69	43.07	433	1329	9	0.246		
		VFC	79.13	95.09	71.88	626	1877	1	0.035		
		LLT	44.49	56.98	43.67	328	1354	26	3.583		
	ISAC	RFM-SCAN	60.18	92.96	69.89	610	2525	0	1.534		
		AdaLAM	92.98	80.22	61.24	536	1175	1	0.429		
		OANet	86.50	78.58	58.77	525	1173	0	0.119		
		ACNe	79.21	79.33	60.47	524	1308	0	0.067		
		PFM	93.90	36.65	28.39	236	293	3	0.165		
		PGM	68.78	72.63	54.10	420	728	5	4.974		
		SCV	21.14	45.54	35.79	353	2136	0	1.203		
		BM	86.26	54.88	42.69	439	954	8	1.654		
		non-planar	no model	blob	94.59	92.54	69.53	593	1527	1	-
				th	96.75	63.42	48.66	431	715	1	-
DTM ₁	95.35			72.95	55.80	500	1211	1	-		
DTM ₂	95.40			78.96	60.03	532	1278	1	-		
LMR	94.13			60.04	45.79	417	793	2	-		
LPM	91.31			36.40	28.60	270	577	4	-		
GLPM	95.24			73.16	56.10	498	1218	1	-		
GMS	82.35			51.73	40.84	410	1077	10	-		
VFC	92.69			88.75	67.37	588	1510	2	-		
LLT	62.13			53.63	41.14	307	825	26	-		
ISAC	RFM-SCAN		94.69	86.62	65.54	574	1485	1	-		
	AdaLAM		96.25	76.58	58.44	512	1085	1	-		
	OANet		94.94	74.08	55.64	500	1037	2	-		
	ACNe		92.35	74.60	57.12	499	1073	3	-		
	PFM		94.54	35.65	27.56	230	284	3	-		
	PGM		89.74	67.95	50.76	393	545	5	-		
	SCV		81.88	41.76	32.89	335	548	12	-		
	BM		86.46	52.75	41.04	420	881	9	-		
	non-planar		no model	blob	24.01	100.00	81.61	572	5933	0	-
				th	90.37	42.26	36.15	262	424	0	0.000
DTM ₁		70.25		71.38	59.59	432	1291	0	0.430		
DTM ₂		70.25		75.44	62.88	455	1348	0	0.856		
LMR		75.29		53.49	44.65	328	698	0	0.093		
LPM		74.09		29.90	25.29	180	451	1	0.024		
GLPM		69.64		55.79	47.42	348	1002	3	0.032		
GMS		76.84		57.30	48.49	379	1131	5	0.237		
VFC		57.76		86.87	72.55	535	2055	4	0.037		
LLT		40.50		57.53	49.15	376	1646	27	3.243		
ISAC		RFM-SCAN	52.30	95.06	77.99	557	2508	0	1.305		
		AdaLAM	92.15	54.95	46.67	353	691	2	0.295		
		OANet	83.54	71.94	59.18	423	949	0	0.114		
		ACNe	69.94	67.33	55.65	396	1047	0	0.063		
		PFM	91.80	24.73	21.15	145	177	2	0.131		
		PGM	67.76	72.30	59.92	417	829	1	7.961		
		SCV	22.66	40.03	34.16	248	1530	1	1.024		
		BM	88.58	51.39	44.14	349	737	4	1.261		
		non-planar	no model	blob	66.97	82.58	68.72	497	1624	0	-
				th	93.05	38.48	32.99	240	379	1	-
DTM ₁	85.71			62.32	52.51	387	976	1	-		
DTM ₂	85.40			65.59	55.22	407	1018	0	-		
LMR	89.08			48.10	40.48	301	560	0	-		
LPM	85.54			27.57	23.47	169	375	3	-		
GLPM	84.01			50.91	43.42	318	782	3	-		
GMS	83.66			51.43	43.67	340	917	5	-		
VFC	74.84			75.19	63.26	472	1351	4	-		
LLT	57.15			52.21	44.77	344	1002	27	-		
ISAC	RFM-SCAN		75.16	80.64	67.08	489	1466	0	-		
	AdaLAM		92.85	50.94	43.32	324	628	2	-		
	OANet		86.86	63.55	52.70	381	811	1	-		
	ACNe		79.01	60.91	50.67	365	850	0	-		
	PFM		92.89	23.17	19.84	137	165	2	-		
	PGM		80.57	65.48	54.63	388	646	1	-		
	SCV		67.69	35.99	30.89	225	473	2	-		
	BM		87.80	47.29	40.65	319	659	6	-		

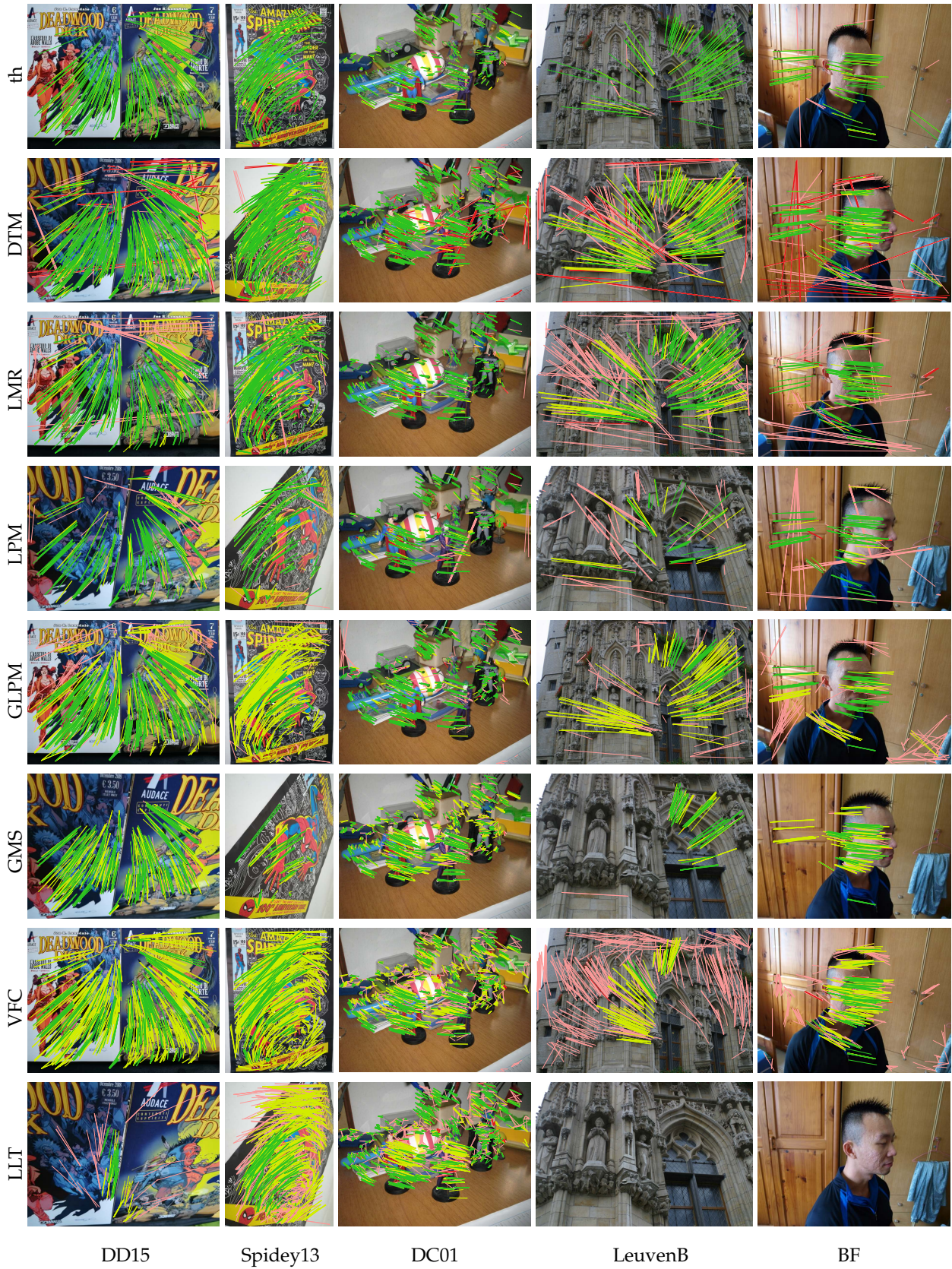


Fig. 15. Planar and non-planar local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. Image indexes are reported as suffix when the sequence contains more than two images. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).

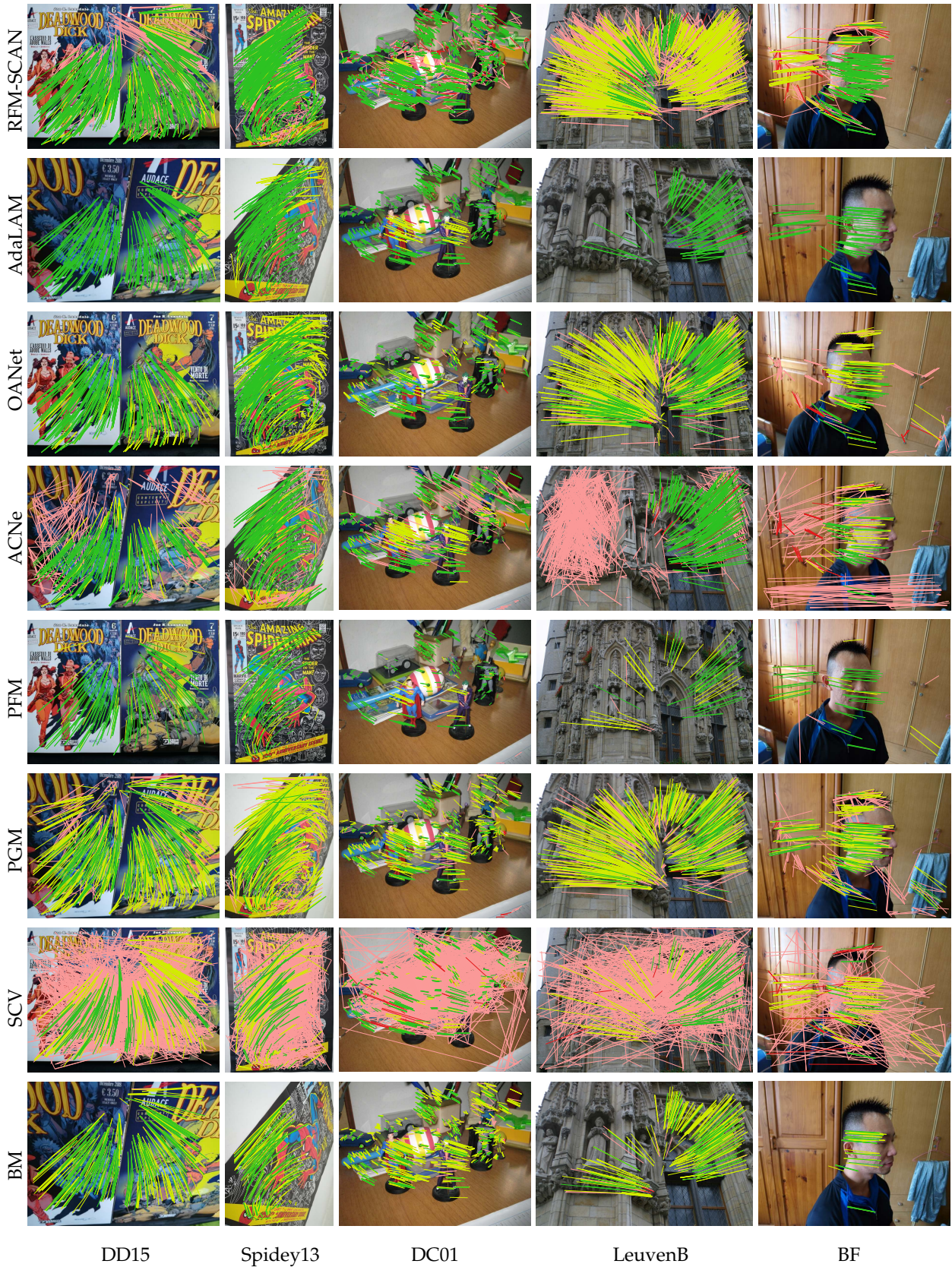


Fig. 16. Planar and non-planar local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. Image indexes are reported as suffix when the sequence contains more than two images. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).

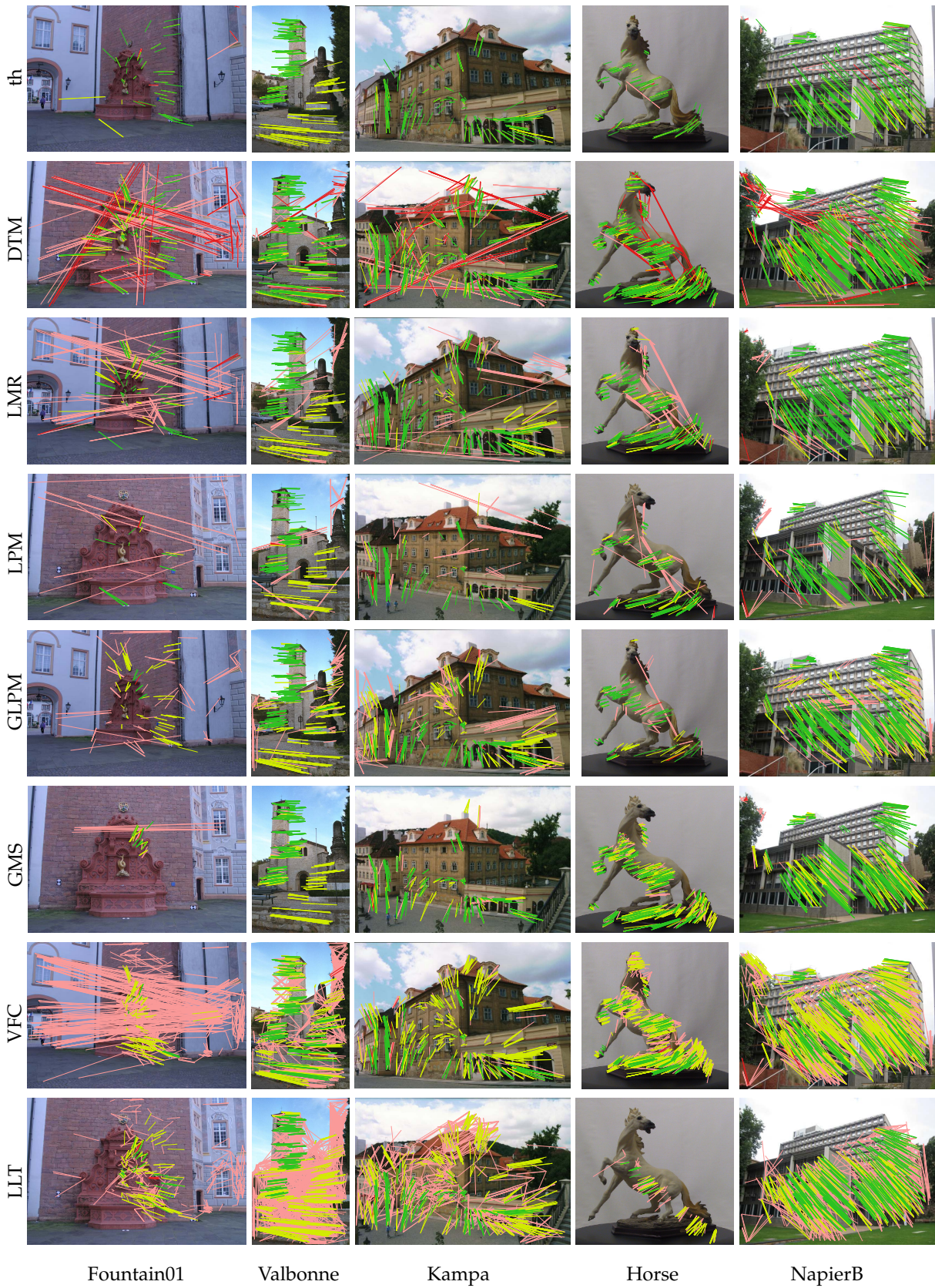


Fig. 17. Planar and non-planar local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. Image indexes are reported as suffix when the sequence contains more than two images. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).

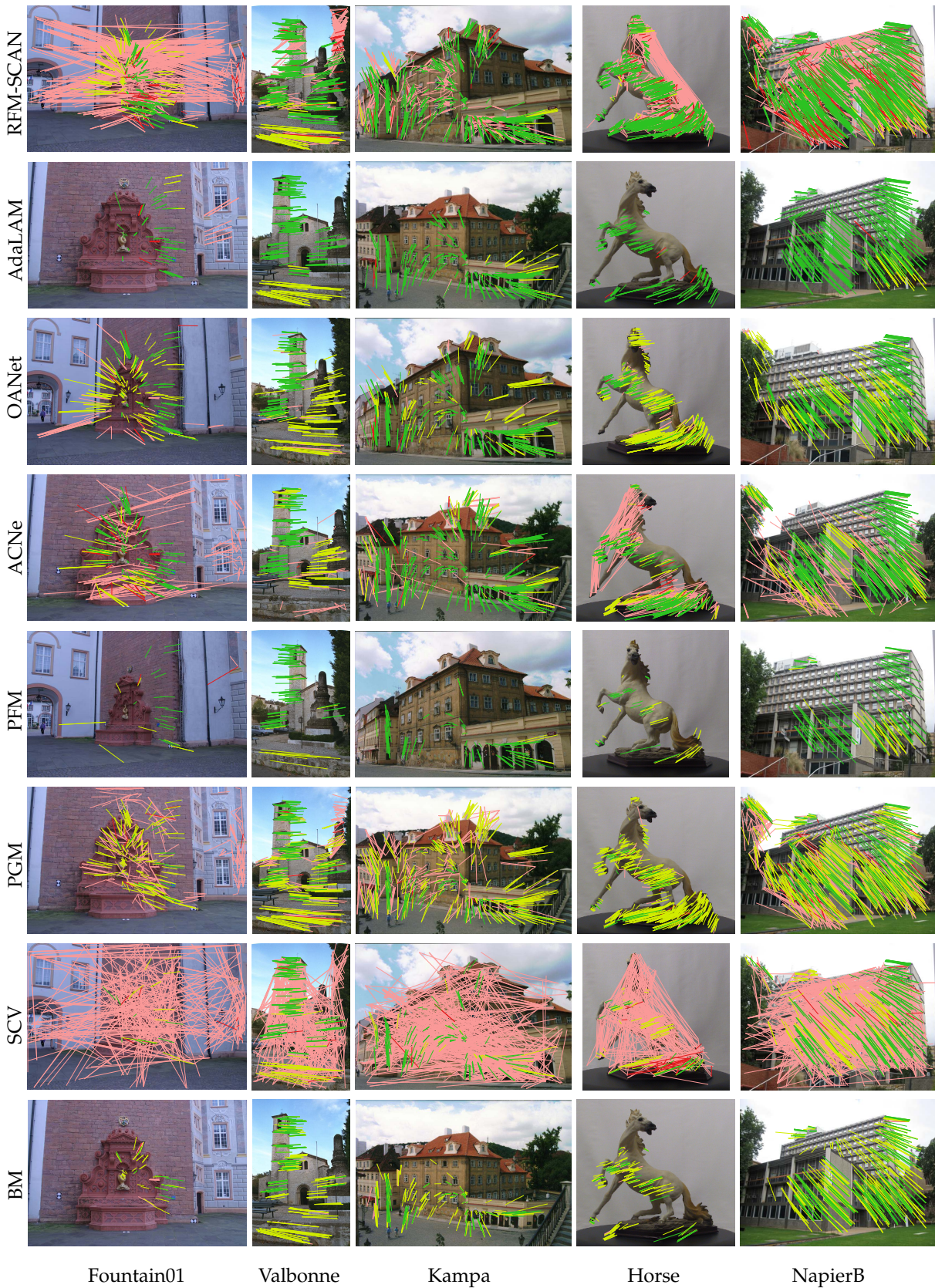


Fig. 18. Planar and non-planar local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. Image indexes are reported as suffix when the sequence contains more than two images. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).



Fig. 19. SUN3D local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).



Fig. 20. SUN3D local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).

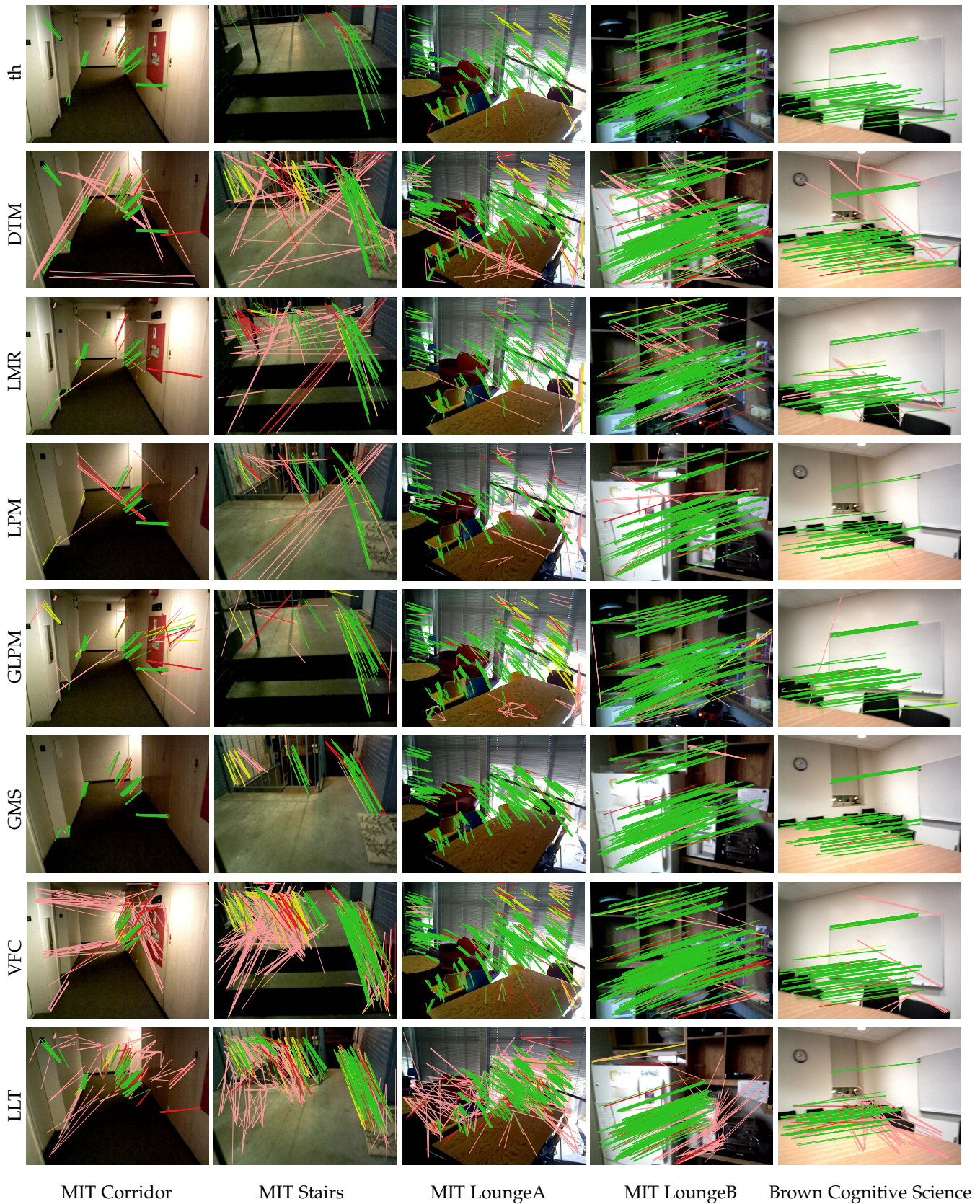


Fig. 21. SUN3D local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).

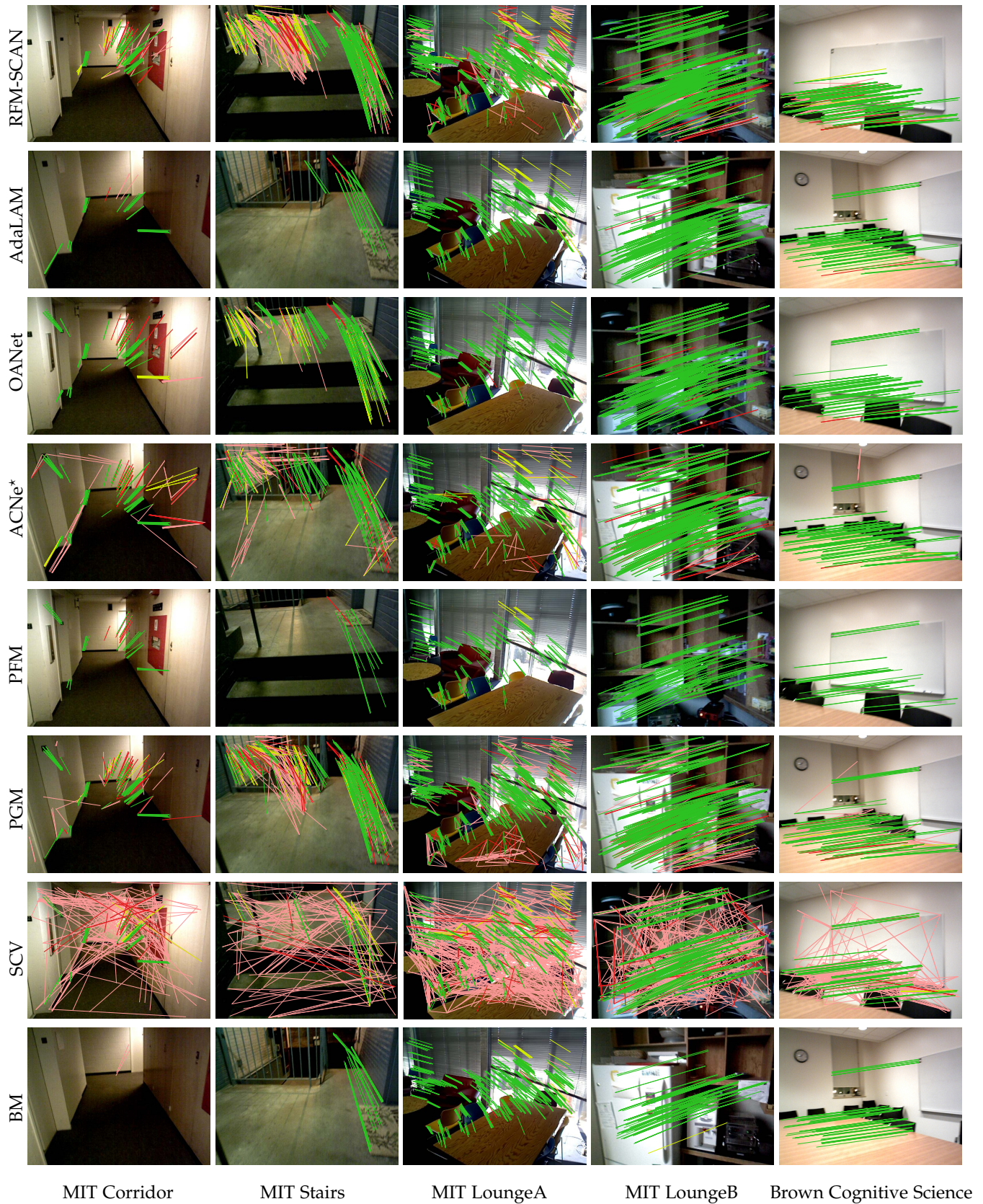


Fig. 22. SUN3D local spatial filter matches according to the best configuration setup, the images of the input pair alternate among the rows. For each method inlier (yellow, green) and outlier (red and light red) clusters are shown, as well as the 1SAC filtered matches (green, red) (see Sec. 5.3, best viewed in color and zoomed in).

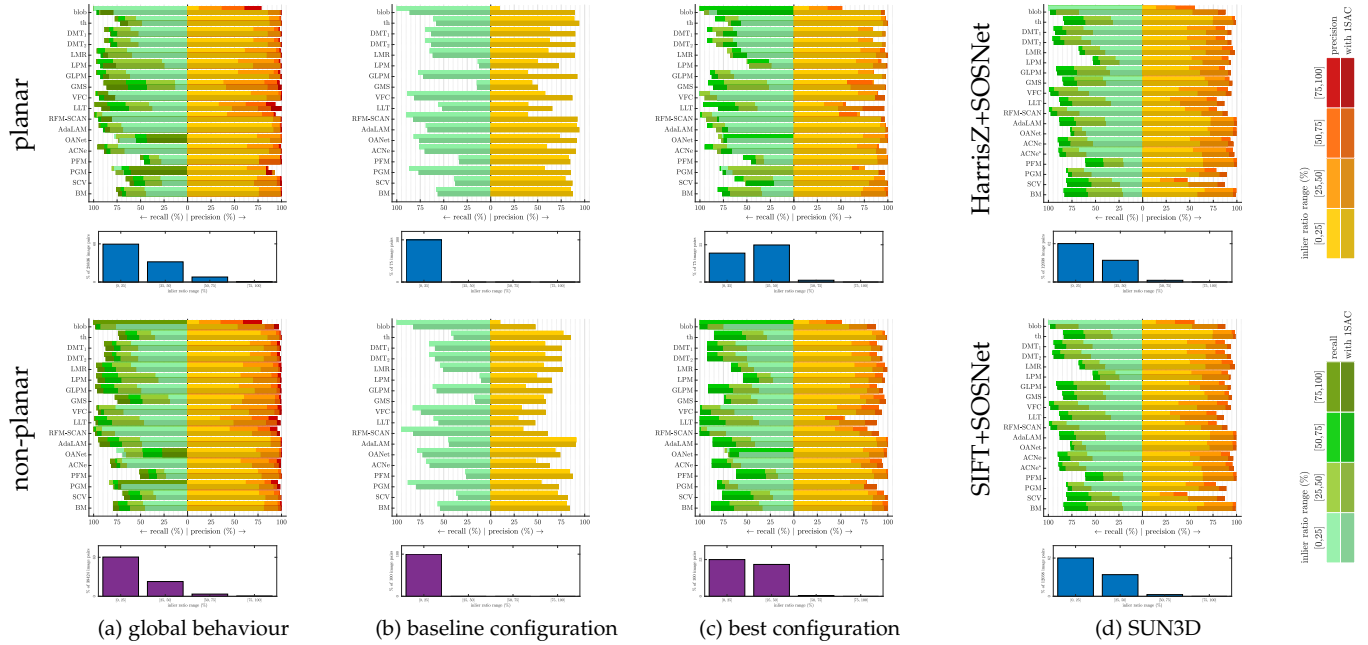


Fig. 23. Average precision/recall values of the local spatial filters according to the inlier ratio. The recall is computed with respect to the ground truth matches found by blob matching (see text for details, best viewed in color and zoomed in).