# Comparative evaluation of 2D feature correspondence selection algorithms

Chen Zhao, Jiaqi Yang, Yang Xiao, and Zhiguo Cao

arXiv:1904.13383v1 [cs.CV] 30 Apr 2019

*Abstract*—Correspondence selection aiming at seeking correct feature correspondences from raw feature matches is pivotal for a number of feature-matching-based tasks. Various 2D (image) correspondence selection algorithms have been presented with decades of progress. Unfortunately, the lack of an in-depth evaluation makes it difficult for developers to choose a proper algorithm given a specific application. This paper fills this gap by evaluating eight 2D correspondence selection algorithms ranging from classical methods to the most recent ones on four standard datasets. The diversity of experimental datasets brings various nuisances including zoom, rotation, blur, viewpoint change, JPEG compression, light change, different rendering styles and multi-structures for comprehensive test. To further create different distributions of initial matches, a set of combinations of detector and descriptor is also taken into consideration. We measure the quality of a correspondence selection algorithm from four perspectives, i.e., precision, recall, F-measure and efficiency. According to evaluation results, the current advantages and limitations of all considered algorithms are aggregately summarized which could be treated as a "user guide" for the following developers.

*Index Terms*—2D feature correspondence, feature matching, correspondence selection, inliers

## I. INTRODUCTION


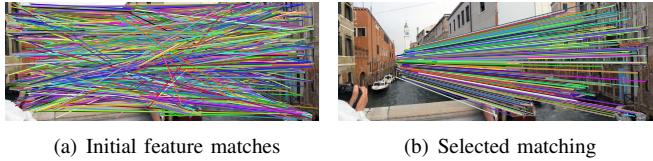
(a) Initial feature matches　　　(b) Selected matching

Fig. 1. An exemplar illustration of the 2D feature correspondence selection, where the colorized lines represent correspondences between two images. (a) Initial feature matches generated by brute-force descriptor matching. (b) Feature matches after correspondence selection.

Feature correspondence selection is a fundamental and critical task in computer vision and robotics. It is the basis for a wide range of applications, such as structure-from-motion [1], simultaneous localization and mapping [2], tracking [3], image stitching [4], and object recognition [5], to name just a few.

The main purpose of correspondence selection is retrieving as many as correct correspondences (also known as *inliers*) from the initial feature matches. Usually, this task is under the background of feature matching. The general process of feature matching starts by detecting representative points, namely keypoints, for two images to be matched. Then, local descriptors such as SIFT [6] and ORB [7] are employed to perform feature description for those keypoints. To build the connection between two images, keypoints with similar feature descriptors are matched, generating a set of raw feature matches. However, the initial feature matches often suffer from severe wrong matches (as shown in Fig. 1(a)) due to the limited distinctiveness of feature descriptors or/and external interferences such as noise and occlusion. This problem makes correspondence selection a necessity for accurate feature matching. Fig. 1(b) shows that those matches after correspondence selection are far more consistent than the initial feature matches. This consensus allows massive high-level vision tasks. For instance, homography, affine and essential matrices can be estimated from those consistent correspondences, thus allowing us to compute the transformation between two images and warp them into a unified coordinate system [4]. Other applications also involve camera parameter estimation [1] and object tracking [3]. Nonetheless, the correspondence selection problem is difficult in real applications due to several factors, e.g., zoom, rotation, blur, viewpoint change, JPEG compression, light change, different rendering styles, multi-structures, and etc. Different scenarios will also lead to different distributions of feature matches which are linearly non-separable.

To address these problems, many approaches that have been presented during the past two decades can be divided into two categories [8]: parametric and non-parametric methods. (i) For parametric methods, they seek consistent correspondences grounded on parametric geometric models. Typical methods include the random sample consensus (RANSAC) [9], the progressive sample consensus (PROSAC) [10], the universal framework for random sample consensus (USAC) [11], and etc. (ii) For non-parametric methods, they are independent from parametric model assumptions. Some of them search correspondence inliers via either feature similarity constraint or geometric constraint, such as the nearest neighbor similarity ratio (NNSR) [6], spectral technique (ST) [12], game-theoretic matching (GTM) [13], graph-based affine invariant matching (GAIM) [14] and locality preserving matching (LPM) [15]. There are also constraint-independent non-parametric methods such as identifying point correspondences by correspondence function (ICF) [16], vector field consensus (VFC) [8], grid-based motion statistics (GMS) [17] and coherence based decision boundaries (CODE) [18]. With the wealth of existing correspondence selection methods, however, it is on the one hand difficult for developers to choose the most proper method given a specific application and on the other hand confusing for researchers to compare these methods under different

C. Zhao, Jiaqi Yang, Y. Xiao, and Zhiguo Cao was with School of Artificial Intelligence and Automation, Huazhong University of Science and Technology, Wuhan, 430074, China e-mail: (hust_zhao@hust.edu.cn, jqyang@hust.edu.cn, Yang_Xiao@hust.edu.cn, zgcao@hust.edu.cn) (Corresponding author: Zhiguo Cao.)

conditions. This problem is mainly due to the fact that most methods were tested under a specific application scenario and compared with a limited number of baselines.

Some performance evaluations in the field of image feature matching also exist. For instance, Mikolajczyk et al. [19] and Heinly et al. [20] evaluated the performance of several 2D feature descriptors. Aans et al. [21] investigated the performance of 2D feature detectors. Moreels et al. [22] performed an aggregated evaluation of both 2D detectors and descriptors. In addition to feature detectors and descriptors, Raguram et al. [23] tested the performance of a set of random sample consensus methods including the popular RANSAC and its variants. However, all these evaluations are either not in line with 2D correspondence selection or not comprehensive enough for an in-depth comparison. First, the critical step in correspondence selection is finding correspondence consensus, while feature detection and description aim at building high-quality initial feature correspondences (such quality is difficult to be guaranteed without correspondence selection [6]). Second, the performance of non-parametric approaches and some recent algorithms remains unclear (only parametric methods were tested in [23]).

In these regards, we present the first comprehensive evaluation, to the best of our knowledge, for 2D correspondence selection from different perspectives in a uniform experimental framework. The considered methods in our evaluation range from classical algorithms to the most recent ones, typically covering both parametric and non-parametric approaches. To be specific, RANSAC [9] and USAC [11] are selected from the parametric family, as RANSAC is arguably the most popular parametric approach and USAC is a well-known modified version of RANSAC. As for non-parametric methods, we choose NNSR [6] as a representative of those methods based on descriptor similarity constraints. ST [12], GTM [13] and LPM [15] are selected as they all rely on the geometric consensus. VFC [8] and the recent GMS [17] are taken into consideration since they eliminate outliers from the perspective of statistical measures. In order to compare those methods from different perspectives, we choose four standard datasets, i.e., VGG [24], Heinly [20], Symbench [25], AdelaideRMF [26], as experimental platforms under the motivation to test those correspondence selection methods' overall performance when faced with a variety of nuisances rather than in their favoring circumstances. For instance, geometric constraints may turn to be vulnerable under rigid/non-rigid transformations such as zoom and rotation; feature similarity constraints are suspicious when the image undergoes blur and light changes; parametric models (the homography matrix) can hardly cope with scenes with parallax (all above conclusions have been verified in Sect. IV). The considered datasets well cover these concerns. To be specific, VGG is a hybrid dataset containing challenges including zoom, rotation, blur, JPEG compression, light and viewpoint change. Heinly contains pure zoom and rotation. Symbench involves scenes with light changes and varying rendering styles. AdelaideRMF possesses viewpoint change and multi-structures, resulting in parallax. The behavior of each method is quantitatively measured using precision, recall and F-measure [15], [17], [27]. In addition, the performance under

preselected correspondences (with higher inlier ratios) and different detector-descriptor combinations are also accessed to test their flexibility with respect to the inlier ratio and correspondence distribution changes. Finally, the efficiency with respect to different scales of initial feature matches are examined. According to the experimental outcomes, we make an aggregated summary of the current advantages and limitations of our evaluated methods as well as their suitable applications.

In a nutshell, the contributions of this paper are threefold:
- A review and the core computation steps of eight state-of-the-art 2D correspondence selection algorithms are presented.
- We comprehensively evaluate and compare the performance, the robustness to a variety of perturbations and the efficiency of each algorithm on four standard datasets consisting of hundreds of images with zoom, rotation, blur, viewpoint change, JPEG compression, light change, different rendering styles and multi-structures.
- Instructive summarizations including merits, demerits and suitable applications of the tested methods are given that can be served as a "user guide" for the developers.

The remainder of this paper is organized as follows. Sect. II gives a review of 2D correspondence selection algorithms and relevant evaluations. Sect. III presents the core computation steps of eight state-of-the-art approaches. Sect. IV describes the experimental setup including datasets, criteria and implementation details of the evaluated methods. Qualitative and quantitative experimental results are shown in Sect. V. Summary and discussion are presented in Sect. VI. Conclusions are finally drawn in Sect. VII.

## II. RELATED WORK

This section briefly reviews the prior works of 2D correspondence selection including both parametric and non-parametric categories. Relevant evaluations in the field of feature matching are also discussed.

### A. Correspondence selection methods

For parametric methods, the most well-known algorithm is arguably RANSAC presented by Fischler et al. [9]. RANSAC iteratively explores the space of model parameters by randomly sampling and estimates the most reliable model based on the maximum number of inliers. Then, outliers can be removed using the generated model. Several variants of RANSAC such as MLESAC [28], LO-RANSAC [29], PROSAC [10] and USAC [11] were proposed in the following decades. MLESAC employs the maximum likelihood estimation rather than the inlier count to check the solutions. LO-RANSAC inserts an optimization process where the generated model is refined by the subset of inliers. A weighted sampling step is adopted instead of random sampling in PROSAC. This method sorts the raw correspondences by matching quality and generates hypotheses from the most promising correspondences. USAC extends the standard hypothesize-and-verify structure in RANSAC and presents a universal framework that integrates advantages of previous parametric methods. In

addition, some other approaches relying on local parametric structures have also been developed, such as agglomerative correspondence clustering (ACC) [30], multi-structures robust fitting (Multi-GS) [31], Hough voting and inverted Hough voting (HVIV) [32]. ACC uses Hessian-affine detector [33], which is invariant to affine transformations, to estimate the local homography matrix as constraints. The initial correspondences are then clustered based on the constraints, and the clusters with inliers are supposed to be larger than the ones constituted by outliers. Multi-GS generates a series of tentative hypotheses by random sampling and considers that two correspondences from the same local structure are inliers if they share a common list of hypotheses. HVIV employs the BPLR detector [34] to cluster correspondences and estimates the homographic transformation for each correspondence as well. The most plausible correspondence in each cluster is then selected using normalized kernel density estimation.

For non-parametric methods, their theoretical foundations are not always the same. A widely-used strategy is exploiting the consistency information of local geometric structures or appearance (feature similarity). Specifically, Lowe et al. [6] proposed a nearest neighbor similarity ratio (NNSR) method that assigns a penalty equaling to the ratio of the closest to the second-closest feature distance to each correspondence and treats those correspondences with low ratios as inliers. Leordeanu et al. [12] presented spectral technique (ST), where an affinity matrix is built using pairwise geometric constraints to remove mismatches in conflict with the most credible correspondences. Albarelli et al. [13] casted the selection of correspondences in a game theoretic framework, known as game-theoretic matching (GTM), where a natural selection process allows corresponding points that satisfy a mutual distance constraint to thrive. Cho et al. [35] presented reweighted random walk algorithm (RRWM) for graph matching. An associated graph between two sets of candidate correspondences is drawn at first, and reliable nodes indicating the consistent correspondences in this graph are then selected by the reweighted random walk algorithm. Ma et al. [15] proposed locality preserving matching (LPM) to improve inlier selection by maintaining the local neighborhood structures of those potential true matches. Some non-parametric approaches that formulate the correspondence selection problem as a statistics problem have also been used, e.g., vector field consensus (VFC) [8] and grid-based motion statistics (GMS) [17]. VFC supposes that the noise around inliers and outliers falls in different distributions. This approach estimates the probability of inliers by the maximum likelihood estimation for parameters in the mixture probabilistic model. Additionally, GMS rejects false matches by counting the quantity of matches in small neighborhoods and achieves real-time performance with an efficient grid-based score estimator.

### B. Other evaluations

In the feature matching field, some evaluations of 2D/3D local descriptors and detectors have been performed. For instance, Mikolajczyk et al. [19] evaluated the performance of 2D feature descriptors under transformations of rotation, zoom, viewpoint change, blur, JPEG compression, light change and keypoint localization errors. Moreels et al. [22] conducted an evaluation of several groups of 2D feature detectors and descriptors on images captured from the same 3D object with different viewpoints and lighting conditions. Heinly et al. [20] performed an evaluation of several 2D binary descriptors, aiming at testing their descriptiveness under different feature detectors on several scenes with illumination change, viewpoint change, pure camera rotation and pure scale change. Aans et al. [21] investigated the performance of several 2D feature detectors on a particular dataset wherein each scene was depicted from 119 camera positions with a range of light directions. In 3D domain, Tombari et al. [36] compared two categories (i.e., fixed-scale and adaptive-scale) of 3D feature detectors in terms of distinctiveness, repeatability and efficiency under the nuisances of viewpoint changes, clutter, occlusions and noise. Guo et al. [37] tested the descriptiveness, robustness, compactness and efficiency of ten local geometric descriptors on eight datasets with radius variations, varying mesh resolution, Gaussian noise and etc. More relevant to our work is the evaluation performed by Raguram et al. [23], where RANSAC and a set of its variants were examined under different ratios of inliers. This paper, compared with [23], considers both parametric and non-parametric methods as well as a variety of nuisances for more comprehensive evaluation.

### III. Considered methods

Eight 2D correspondence selection algorithms including two parametric ones, i.e., RANSAC [9] and USAC [11], and six non-parametric ones, i.e., NNSR [6], ST [12], GTM [13], VFC [8], GMS [17], LPM [15], are considered in our evaluation. Before introducing their theories, we give some general notations for better readability.

Given two images $(I, I^{'})$ to be matched, keypoints and local feature descriptors are computed for them as $(\mathcal{K}, \mathcal{K}^{'})$ and $(\mathcal{F}, \mathcal{F}^{'})$, respectively. This procedure can be accomplished using off-the-shelf detectors and descriptors, e.g., SIFT [6]. To generate initial feature matches $\mathcal{C}$, keypoints are matched with each other based on feature similarity, i.e., a correspondence (match) in $\mathcal{C}$ is defined as $c = \{\mathbf{x}, \mathbf{x}^{'}, \arg\max_{\mathbf{f}^{'}} s_{\mathcal{F}(\mathbf{f}, \mathbf{f}')}\}$ with $\mathbf{x} \in \mathcal{K}$, $\mathbf{x}^{'} \in \mathcal{K}^{'}$, $\mathbf{f} \in \mathcal{F}$, $\mathbf{f}^{'} \in \mathcal{F}^{'}$ and $s_{\mathcal{F}}$ being the feature similarity score. The objective of correspondence selection is digging out the maximum consensus (inlier) subset $\mathcal{C}_{inlier} \subseteq \mathcal{C}$. Core principles and computation steps of evaluated algorithms are given as follows.

**Nearest Neighbor Similarity Ratio [6].** NNSR directly utilizes descriptor similarities to remove less distinctive matches. Specifically, the term equaling to the ratio of the closest to the second-closest feature distance to each correspondence is used as a penalty. Therefore, a correspondence is judged as inlier if

$$\frac{\| \mathbf{f} - \mathbf{f}_1^{'} \|_2}{\| \mathbf{f} - \mathbf{f}_2^{'} \|_2} \leq t_{nnsr}, \tag{1}$$

where $t_{nnsr} \in [0, 1]$, $\|\cdot\|_2$ hereinafter denotes the $L_2$ norm (this distance metric is suggested in [6]), $\mathbf{f}_1^{'}$ and $\mathbf{f}_2^{'}$ represent the most and the second most similar feature descriptors of $\mathbf{f}$, respectively. Values of threshold $t_{nnsr}$ and other mentioned thresholds in the following are presented in Table I.

**Random Sample Consensus [9].** RANSAC follows a hypothesize-and-verify framework by repeating procedures of random sampling and checking to maximize the object function. For 2D correspondence selection, the desired parametric model is usually a plane homography matrix or a fundamental matrix. Taking the homography matrix as an example, it first randomly samples several correspondences (at least 4) from $\mathcal{C}$ and generates the model hypothesis $\mathbf{H}_i$ for those samples at the $i$th iteration. Then, the hypothesis $\mathbf{H}_i$ is verified via the following object function

$$O_i = \sum_{c \in \mathcal{C}} h_i(c), \tag{2}$$

where $h(\cdot)$ is a binary function defined as

$$h_i(c) = \begin{cases} 1, & \text{if } \|\mathbf{x}' - \rho\left(\mathbf{H}_i \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}\right)\|_2 \leq t_{ransac} \\ 0, & \text{otherwise} \end{cases}, \tag{3}$$

with $\rho([a_1 \ a_2 \ a_3]^T) = [a_1/a_3 \ a_2/a_3]^T$ and $t_{ransac}$ being a threshold that determines the accuracy of a judged inlier. Above steps are repeated $n_{ransac}$ times and the model with the maximum object function is selected as the final model $\mathbf{H}^\star$. Correspondences agreeing with $\mathbf{H}^\star$ (producing 1 values using Eq. 3) are identified as inliers.

**Spectral Technique [12].** ST locates the most reliable element by matrix decomposition. It assumes that the connection among correct matches is much tighter than the one among mismatches. Based on this assumption, ST first builds an adjacency matrix $\mathbf{A}$ as

$$a_{ij} = \min\left(\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2}{\|\mathbf{x}'_i - \mathbf{x}'_j\|_2}, \frac{\|\mathbf{x}'_i - \mathbf{x}'_j\|_2}{\|\mathbf{x}_i - \mathbf{x}_j\|_2}\right), \tag{4}$$

where $a_{ij} \in \mathbf{A}$ is the affinity between $c_i$ and $c_j$. Second, the principle eigenvector $\mathbf{v}_{st}$ of $\mathbf{A}$ is computed using the singular value decomposition algorithm. Third, the maximum element in $\mathbf{v}_{st}$ is selected as $v_i$ indicating $c_i$ being the most reliable correspondence. Fourth, set $v_i$ to zero and remove other components of $\mathcal{C}$ that are in conflict with $c_i$, i.e.,

$$a_{ij} \leq t_{st}, \tag{5}$$

where $t_{st}$ is a predefined threshold. By repeating the third and fourth steps until $\mathcal{C}$ is empty or $v_i = 0$, the correspondences related to all elements selected from $\mathbf{v}_{st}$ are determined as inliers.

**Game Theory Matching [13].** GTM concentrates on extracting correspondences being consistent to the majority of $\mathcal{C}$. Specifically, this strategy interprets the filtering process as a game-theoretic framework where players attempt to obtain high payoffs. At the beginning of this game, every two players extracted from a large population choose a pair of correspondences (served as strategies in this context) from $\mathcal{C}$. Then they will receive a payoff linearly correlated to the coherence between these correspondences. The player who gets high payoffs will receive higher supports. In general, as the game going on, players will prefer to select more reliable correspondences to pursue higher pay-offs.

Given a pair of correspondences $(c_i, c_j)$, the payoff function is defined as

$$\Pi_{ij} = e^{-\lambda_{gtm} \max(|T_i(\mathbf{x}_i) - T_j(\mathbf{x}_i)|, |T_i(\mathbf{x}_j) - T_j(\mathbf{x}_j)|)}, \tag{6}$$

where $\lambda_{gtm}$ is a selectivity parameter, $|\cdot|$ represents the $L_1$ norm and $T_i(\mathbf{x})$ is the similarity transformation estimated by (similarly for $T_j(\mathbf{x})$)

$$T_i(\mathbf{x}) = \rho\left(\mathbf{H}_{c_i} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}\right), \tag{7}$$

where $\mathbf{H}_{c_i}$ is the homographic transformation of $c_i$. Note that this algorithm particularly requires the local affine transformation cue to compute the pay-off function. Next, the payoff matrix $\mathbf{P}_{gtm}$ with the element in the $i$th row and $j$th column that is defined as

$$p_{ij} = \begin{cases} \Pi_{ij}, & \text{if } i \neq j \\ 0, & \text{otherwise} \end{cases}, \tag{8}$$

can be generated. The population vector $\mathbf{q}$ is updated by the evolutionary stable states algorithm (ESS's) [38] as

$$q_i(k+1) = q_i(k)\frac{(\mathbf{P}_{gtm}\mathbf{q}(k))_i}{\mathbf{q}(k)^T \mathbf{P}_{gtm}\mathbf{q}(k)}, \tag{9}$$

where $q_i$ represents the element in the $i$th row of $\mathbf{q}$ and $k$ is the iteration number. After $n_{gtm}$ iterations, a correspondence $c_i$ is identified as inlier if its corresponding $q_i$ is higher than a threshold $t_{gtm}$.

**Universal RANSAC [11].** USAC integrates a universal framework for RANSAC, where each original step is optimized by referring to the advantages of previous parametric approaches such as PROSAC [10], SPRT test [39] and LO-RANSAC [29]. Further, this algorithm inserts degeneracy and local optimization processes after generating the minimal-sample model.

During the sampling step, USAC uses a weighted sampling algorithm named PROSAC [10], where the initial correspondences are reordered at first based on the descending sort order of brute-force matching scores and correspondences with higher scores are preserved. At the checking stage of the model (homography matrix or fundamental matrix), a correspondence is judged as inlier by Eq. 3 with the threshold $t_{\mathbf{H}}$ or by the equation

$$\frac{\left(\mathbf{y}'^T \mathbf{F}_i \mathbf{y}\right)^2}{(\mathbf{F}_i \mathbf{y})_1^2 + (\mathbf{F}_i \mathbf{y})_2^2 + \left(\mathbf{F}_i^T \mathbf{y}'\right)_1^2 + \left(\mathbf{F}_i^T \mathbf{y}'\right)_2^2} \leq t_{\mathbf{F}}, \tag{10}$$

where $\mathbf{F}_i$ is the $i$th hypothetic fundamental matrix, $t_{\mathbf{F}}$ is the threshold and $\mathbf{y} = [\ \mathbf{x} \ 1 \ ]^T$ (similarly for $\mathbf{y}'$). After generating the minimal-sample model, USAC verifies whether the model is interesting by the SPRT test [39]. The likelihood ratio can be computed after evaluating $n$ correspondences as

$$\xi_n = \prod_{i=1}^{n} \frac{p(r_i|\mathbf{H}_b)}{p(r_i|\mathbf{H}_g)}, \tag{11}$$

where $\mathbf{H}_g$ and $\mathbf{H}_b$ respectively represent a "good" model and a "bad" model, $r_i$ is equal to 1 if $c_i$ is consistent with the generated model and 0 otherwise, $p(1|\mathbf{H}_g)$ is approximated by

the inlier ratio and $p(r_i|\mathbf{H}_b)$ follows a Bernoulli distribution. If the $\xi_n$ is higher than an adaptive threshold, the model will be discarded. When fitting the fundamental matrix by epipolar geometry constraint, USAC utilizes DEGENSAC [40] for degeneracy. It assumes that the generated model is often incorrect in the context of images containing a dominant scene plane. Accordingly, DEGENSAC employs a homographic transformation to reject the generated fundamental model if there are five or more sampled correspondences lying on the same plane. Eventually, USAC adds a local optimization (LO-RANSAC [29]) to refine the minimal-sample model. It re-samples correspondences only from the set of selected inliers and refines the previous model by the sampling subset. This whole process is repeated until achieving confidence in solution or iterations reach the upper bound $n_{usac}$.

**Vector Field Consensus [8].** VFC interpolates a vector field where the posteriori probability of a correct correspondence is estimated by the Bayes rule.

For a correspondence $c_i$, the transformation to a motion field is expressed as $(\mathbf{x}_i, \mathbf{x}_i^{'}) \rightarrow (\mathbf{u}_i, \mathbf{v}_i)$, where $\mathbf{u}_i = \mathbf{x}_i$ and $\mathbf{v}_i = \mathbf{x}_i^{'} - \mathbf{x}_i$. In this motion field, VFC holds the assumption that the noise around inliers indicated by $z_i = 1$ follows the Gaussian distribution and the noise around outliers indicated by $z_i = 0$ follows the uniform distribution. Thus, the probability is a mixture model given by

$$p(\mathcal{U}|\mathcal{V}, \theta) = \prod_{i=1}^{N} \left( \frac{\gamma}{(2\pi\sigma^2)^{D/2}} e^{-\frac{\|\mathbf{v}_i - \mathbf{f}_{vfc}(\mathbf{u}_i)\|_2}{2\sigma^2}} + \frac{1-\gamma}{a} \right), \quad (12)$$

where $\theta = \left\{ \mathbf{f}_{vfc}, \sigma^2, \gamma \right\}$ is a set of unknown parameters, $\mathbf{f}_{vfc}$ is the vector field expected to be recovered, $\gamma$ is the mixing coefficient of the mixture probability model, i.e, $p(z_i = 1) = \gamma$, $\mathcal{U}$ and $\mathcal{V}$ respectively are sets of $\mathbf{u}$ and $\mathbf{v}$, $\sigma$ is the uniform standard deviation of Gaussian distribution, $\frac{1}{a}$ is the probability density of the uniform distribution and $D$ is the dimension of the output space. VFC employs the EM [41] algorithm to deal with the maximum likelihood estimation with latent variables. At E-step, the diagonal element of a diagonal matrix $\mathbf{P}$, i.e., $p_i = p(z_i = 1|\mathbf{u}_i, \mathbf{v}_i, \theta)$, can be computed by the Bayes rule

$$p_i = \frac{\gamma e^{-\frac{\|\mathbf{v}_i - \mathbf{f}_{vfc}(\mathbf{u}_i)\|_2}{2\sigma^2}}}{\gamma e^{-\frac{\|\mathbf{v}_i - \mathbf{f}_{vfc}(\mathbf{u}_i)\|_2}{2\sigma^2}} + (1-\gamma)\frac{(2\pi\sigma^2)^{D/2}}{a}}. \quad (13)$$

At M-step, a coefficient matrix $\mathbf{C}$ is created first by

$$(\mathbf{K}_{Gauss} + \lambda_{vfc}\sigma^2\mathbf{P}^{-1})\mathbf{C} = \mathcal{V}, \quad (14)$$

where $\mathbf{K}_{Gauss}$ is a matrix consisting of the Gaussian kernel $k(\mathbf{u}_i, \mathbf{u}_j) = e^{-\beta\|\mathbf{u}_i - \mathbf{u}_j\|_2}$ and $\lambda_{vfc}$ is a regularization constant. Second, the vector field $\mathbf{f}_{vfc}$ is estimated by

$$\mathbf{f}_{vfc}(\mathbf{u}) = \sum_{i=1}^{N} k(\mathbf{u}, \mathbf{u}_i)\mathbf{c}_i, \quad (15)$$

where $\mathbf{c}_i \in \mathbf{C}$. Third, values of $\sigma^2$ and $\gamma$ are updated by

$$\sigma^2 = \frac{(\mathcal{V} - \mathcal{F}_{vfc})^T \mathbf{P}(\mathcal{V} - \mathcal{F}_{vfc})}{D \cdot \text{tr}(\mathbf{P})}, \quad (16)$$

and

$$\gamma = \text{tr}(\mathbf{P})/N, \quad (17)$$

where $\mathcal{F}_{vfc} = \left( \mathbf{f}_{vfc}(\mathbf{u}_1)^T, ...\mathbf{f}_{vfc}(\mathbf{u}_N)^T \right)^T$. The E-step and M-step are repeated until parameters are converged. Finally, the inlier set is generated as

$$\mathcal{C}_{inlier} = \left\{ c_i : p_i > t_{vfc} \right\}, \quad (18)$$

where $t_{vfc}$ is a predefined threshold.

**Grid-based Motion Statistics [17].** GMS proves that besides feature descriptiveness, feature number also contributes to the quality of correspondences. It supposes that the quantity of correspondences in a small neighborhood around a true match is larger than that around a false match under the smooth motion. In over-large neighborhoods, regions are divided into multiple small region pairs where distributions of correspondence number are approximated by Binomial distributions. Given a correspondence $c_i$, the joint statistical distribution is modeled as

$$S_i \sim \begin{cases} B(Kn, p_t), & \text{if } c_i \text{ is inlier} \\ B(Kn, p_f), & \text{otherwise} \end{cases}, \quad (19)$$

where $S_i$ is the total number of correspondences in a region pair $(a, b)$ around $c_i$, $K$ is the quantity of small region pairs, $p_t$ is the probability that the nearest neighbor of each keypoint in $a$ is located in $b$ under the condition that $a$ and $b$ view the same location, and $p_f$ is the probability provided that $a$ and $b$ view the different locations. $p_t$ and $p_f$ can be estimated by

$$p_t = \delta + (1-\delta)\zeta m/M, \quad (20)$$

and

$$p_f = \zeta(1-\delta)(m/M), \quad (21)$$

where $\delta$ is the probability of a correspondence being correct, $m$ is the amount of keypoints in region $b$, $M$ is the size of $\mathcal{K}^{'}$ in $I^{'}$, and $\zeta$ is a factor added to balance deviations caused by repeated structures. A quantitative score is next designed to evaluate the distinction between two distributions as

$$P = \frac{m_t - m_f}{s_t + s_f}, \quad (22)$$

where $m$ is the mean value and $s$ is the standard deviation. This equation can be simplified as

$$P \propto \sqrt{Kn}, \quad (23)$$

where the distinction is positive correlated to the number of correspondences.

In addition, to incorporate this approach into a real-time system, a fast gird-based score estimator is developed as follows. First, $I$ and $I^{'}$ are divided into $20 \times 20$ non-overlapping cells. Second, for each cell in $I$, the cell containing the maximum amount of correspondences is grouped in $I^{'}$. Third, in cell-pair $(i, j)$ as well as its small neighborhoods (eight cell-pairs), $S_{ij}$ is estimated as

$$S_{ij} = \sum_{k=1}^{K=9} \left| \chi_{i^k j^k} \right|, \quad (24)$$

where $|\chi|$ is the amount of correspondences in the cell-pair $(i^k, j^k)$. All correspondences in $(i, j)$ are judged as inliers if

$S_{ij} > t_{gms}$, where $t_{gms}$ is a threshold approximated by $\alpha \sqrt{n_i}$ with $\alpha$ being a given parameter and $n_i$ being the average (of the nine cell-pairs) amount of correspondences.

**Locality Preserving Matching [15].** This algorithm removes mismatches by digging out the local geometric structure consensus. With the hypothesis that the local structure around a correspondence may not change freely, a cost function is defined as

$$
L(\mathcal{K}_{inlier}, \lambda_{lpm}) = \sum_{i \in \mathcal{K}_{inlier}} \left( \sum_{j|\mathbf{x}_j \in \mathcal{K}_{\mathbf{x}_i}} \left( d\left(\mathbf{x}_i, \mathbf{x}_j\right) - d\left(\mathbf{x}_i', \mathbf{x}_j'\right) \right)^2 \right.
$$
$$
\left. + \sum_{j|\mathbf{x}_j' \in \mathcal{K}_{\mathbf{x}_i'}'} \left( d\left(\mathbf{x}_i, \mathbf{x}_j\right) - d\left(\mathbf{x}_i', \mathbf{x}_j'\right) \right)^2 \right) + \lambda_{lpm} \left( N - |\mathcal{K}_{inlier}| \right),
$$
(25)

where $\lambda_{lpm}$ is a regularization parameter, $d$ is the Euclidean distance between two keypoints, $\mathcal{K}_{\mathbf{x}_i}$ and $\mathcal{K}_{\mathbf{x}_i'}'$ respectively are sets of the $k$ nearest neighbors of $\mathbf{x}_i$ and $\mathbf{x}_i'$, $N$ is the size of $\mathcal{K}$, and $\mathcal{K}_{inlier}$ is an inlier subset of $\mathcal{K}$. Under non-rigid transformations such as deformation, the absolute distance in Eq. 25 may not be preserved well. To address this issue, LPM converts the cost function to

$$
L(\mathcal{W}, \lambda_{lpm}) = \sum_{i=1}^{N} w_i \left( \sum_{j|\mathbf{x}_j \in \mathcal{K}_{\mathbf{x}_i}} d\left(\mathbf{x}_i', \mathbf{x}_j'\right) \right.
$$
$$
\left. + \sum_{j|\mathbf{x}_j' \in \mathcal{K}_{\mathbf{x}_i'}'} d\left(\mathbf{x}_i, \mathbf{x}_j\right) \right) + \lambda_{lpm} \left( N - \sum_{i=1}^{N} w_i \right),
$$
(26)

where $\mathcal{W}$ is a set of indicators where $w_i = 1$ indicates the inlier and $w_i = 0$ otherwise. This equation can be further reorganized by merging the related items of $w_i$ as

$$
L(\mathcal{W}, \lambda_{lpm}) = \sum_{i=1}^{N} w_i \left( l_i - \lambda_{lpm} \right) + \lambda_{lpm} N, \quad (27)
$$

where

$$
l_i = \sum_{j|\mathbf{x}_j \in \mathcal{K}_{\mathbf{x}_i}} d\left(\mathbf{x}_i', \mathbf{x}_j'\right) + \sum_{j|\mathbf{x}_j' \in \mathcal{K}_{\mathbf{x}_i'}'} d\left(\mathbf{x}_i, \mathbf{x}_j\right) \quad (28)
$$

is a constraint item measuring the local geometric structure changes. With the objective of minimizing the cost function, a correspondence with the cost, i.e., $l_i > \lambda_{lpm}$, is negative. For this purpose, the correct correspondence set is determined by

$$
w_i = \begin{cases} 1, & \text{if } l_i \leq \lambda_{lpm} \\ 0, & \text{otherwise} \end{cases}, i = 1, \ldots, N. \quad (29)
$$

## IV. EXPERIMENTAL SETUP

The experimental setup is introduced detailedly in this section. First, we list implementations and parameter settings of the evaluated methods. Second, characteristics of four datasets, the experimental criteria and the deployment are formulated.

TABLE I
PARAMETER SETTINGS AND IMPLEMENTATIONS OF EVALUATED ALGORITHMS, WHERE *pix* REPRESENTS THE PIXEL UNIT.

| No. | Algorithm | Implementation | Parameters | Setting |
|---|---|---|---|---|
| 1 | NNSR [6] | OPENCV | $t_{nnsr}$ | Adaptive [42] |
| 2 | RANSAC [9] | OPENCV | $t_{ransac}$ | 10*pix* |
| | | | $n_{ransac}$ | 2000 |
| 3 | ST [12] | MATLAB | $t_{st}$ | 0.3 |
| 4 | GTM [13] | OPENCV | $t_{gtm}$ | Adaptive [42] |
| | | | $n_{gtm}$ | 100 |
| | | | $\lambda gtm$ | 0.0001 |
| 5 | USAC [11] | OPENCV | $n_{usac}$ | 850000 |
| | | | $t_{\mathbf{H}}$ | 10*pix* |
| | | | $t_{\mathbf{F}}$ | 1.5*pix* |
| 6 | VFC [8] | OPENCV | $\beta$ | 0.1 |
| | | | $\lambda_{vfc}$ | 3 |
| | | | $t_{vfc}$ | 0.75 |
| | | | $\gamma$ | 0.9 |
| 7 | GMS [17] | OPENCV | $\alpha$ | 4 |
| 8 | LPM [15] | MATLAB | $\lambda_{lpm}$ | 6 |
| | | | $k$ | 4 |

TABLE II
PROPERTIES OF THE EXPERIMENTAL DATASETS.

| Dataset | Challenges | Matching pairs |
|---|---|---|
| VGG [24] | Zoom, rotation, blur, viewpoint change, light change and JPEG compression | 40 |
| Symbench [25] | Light change, different rendering styles | 46 |
| Heinly [20] | Zoom and rotation | 29 |
| AdelaideRMF [26] | Multi-structures, viewpoint change | 38 |

### A. Implementations

In our experiments, Hessian-affine detector [33] and SIFT descriptor [6] (a popular detector-descriptor combination [22]) are employed in default for image keypoint detection and description. Notice that another reason for using the Hessian-affine detector is that the evaluated GTM method requires local affine information, while we also consider different detector-descriptor combinations in Sect. V-C. The initial correspondence set $\mathcal{C}$ is generated by brute-force matching, i.e., greedy comparison of two feature sets. Parameters and implementations for each algorithm are listed in Table I. Notably, for NNSR and GTM we set $t_{nnsr}$ and $t_{gtm}$ adaptively using the OTSU [42] algorithm to reduce thresholding errors as proper thresholds may vary in different scenarios even in different images.

All those methods are implemented in OPENCV or MATLAB with a PC equipped with a 3.2GHz processor and 8GB memory.

### B. Datasets

We perform our experiments on four datasets, i.e., VGG [24], Symbench [25], Heinly [20], and AdelaideRMF [26]. Exemplar images from these datasets and a brief summarization of their inherited nuisances are shown in Fig. 2 and Table II, respectively.

**The VGG dataset [24].** VGG is a hybrid dataset involving eight scenes. Each scene consists of six images with the first image being the reference one with respect to the others. Challenges including blur, viewpoint change, zoom, rotation, light change, and JPEG compression exist in this dataset. The ground-truth is the homography matrix $\mathbf{H}$, indicating that the

(a) VGG



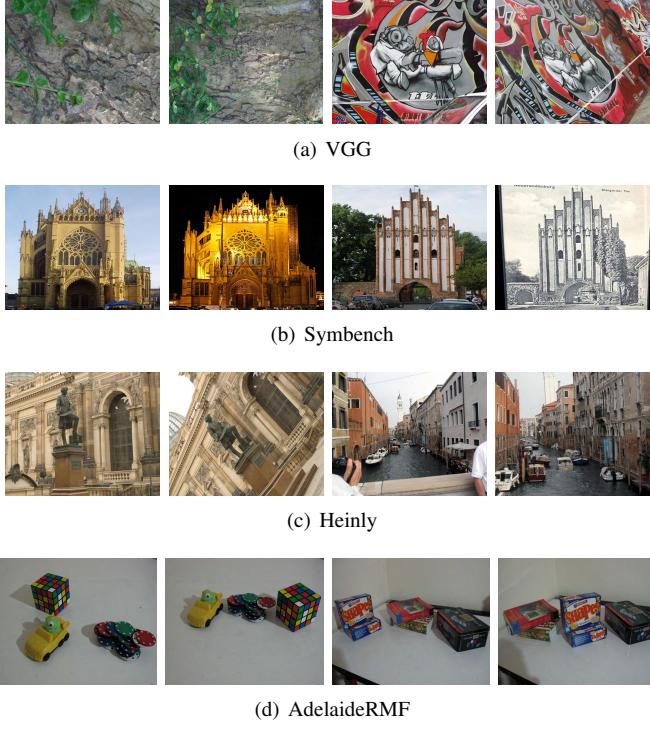(b) Symbench



(c) Heinly



(d) AdelaideRMF

Fig. 2. Sampled image pairs from four experimental datasets. (a) VGG [24], (b) Symbench [25], (c) Heinly [20], and (d) AdelaideRMF [26].

transformation between two images on each scene satisfies the plane homographic constraint.

**The Symbench dataset [25].** The Symbench dataset is composed of 46 image pairs. Each pair includes the same object with light change or different rendering styles. The homographic transformation **H** of each image pair is given as the ground-truth.

**The Heinly dataset [20].** The Heinly dataset comprises images with dense or sparse viewpoint change, illumination, pure large-scale zoom or rotation. Considering that nuisances of viewpoint change and illumination have been covered in the other three datasets, we choose a subset of Heinly containing 29 pairs of image shot on 4 scenes with the specific challenges, i.e., pure zoom or rotation, to perform a more targeted test. The ground-truth is provided as the homographic transformation.

**The AdelaideRMF dataset [26].** AdelaideRMF includes 38 pairs of image with viewpoint change and multi-structures. The keypoint coordinates of initial correspondences are provided and the ground-truth correspondences are manually labeled in this dataset.

Motivations of employing these datasets can be summarized as: (i) The eight scenes in the VGG dataset cover a peculiar wide range of interferences such as the rigid/non-rigid transformation and image quality variation. Both the generality to diverse conditions and the robustness to a specific nuisance can be assessed on this dataset. (ii) The focus of Symbench is the image quality variation caused by light change and different rendering styles that give rise to potential errors of feature detection and description. The performance in the context of image quality variation can be specifically evaluated. (iii)

The subset of Heinly is selected with the aim of testing the performance under the condition of a geometrical structure deformation (pure zoom or rotation). (iv) AdelaideRMF aims at evaluating the performance of those correspondence selection algorithms where plane homographic constraint fails and multiple consistent correspondence sets are involved due to multi-structures. All above peculiarities make the evaluation benchmarks complementary to each other and allow us to find prominent algorithms under a specific nuisance.

### C. Criteria

The performance of evaluated algorithms is measured via precision, recall and F-measure as in [15], [17], [27]. First, we denote the selected correspondence set, the ground-truth correspondence set and the correct subset in the selected correspondence set as $\mathcal{C}_{inlier}$, $\mathcal{C}_{inlier}^{GT}$ and $\mathcal{C}_{inlier}^{correct}$, respectively. Then, the precision, recall and F-measure are respectively defined as

$$\text{Precision} = \frac{\left|\mathcal{C}_{inlier}^{correct}\right|}{\left|\mathcal{C}_{inlier}\right|}, \tag{30}$$

$$\text{Recall} = \frac{\left|\mathcal{C}_{inlier}^{correct}\right|}{\left|\mathcal{C}_{inlier}^{GT}\right|}, \tag{31}$$

and

$$\text{F-measure} = \frac{2\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}, \tag{32}$$

where $|\cdot|$ denotes the cardinality of a set. A correspondence $c = \{\mathbf{x}, \mathbf{x}'\}$ belongs to $\mathcal{C}_{inlier}^{GT}$ if

$$\left\| \mathbf{x}_i' - \rho \left( \mathbf{H}_{gt} \left[ \begin{array}{c} \mathbf{x}_i \\ 1 \end{array} \right] \right) \right\|_2 \leq t_{gt}, \tag{33}$$

where $\mathbf{H}_{gt}$ is the ground-truth homography matrix and $t_{gt}$ is a threshold set to $10pix$ ($pix$ being the unit of pixel) that controls the upper bound of the accuracy of a true inlier in our experiments.

Similarly, a correct correspondence in $\mathcal{C}_{inlier}$ is defined as

$$\left\| \mathbf{x}_i' - \rho \left( \mathbf{H}_{gt} \left[ \begin{array}{c} \mathbf{x}_i \\ 1 \end{array} \right] \right) \right\|_2 \leq \tau \tag{34}$$

with $\tau$ being the matching tolerance. We vary $\tau$ from $1pix$ to $t_{gt}$ with an interval of $1pix$, thus generating a curve [17], [27].

### D. Experimental deployment

Our experiments are deployed as follows. In Sect. V-A, the overall performance of the evaluated algorithms in different scenarios, i.e., the four experimental datasets, is tested. In Sect. V-B, the performance with preselected correspondences by NNSR, i.e., commonly employed to improve the inlier ratio of initial matches [8], [23], [43], [44], is tested on the four datasets. In Sect. V-C, different detector-descriptor combinations are considered to examine the performance variation of correspondence selection algorithms. Notice that different combinations of detector and descriptor are desired in different application contexts [22], [33] and will result in different distributions and inlier ratios. In Sect. V-D, the robustness to different nuisances, i.e., blur, viewpoint change,
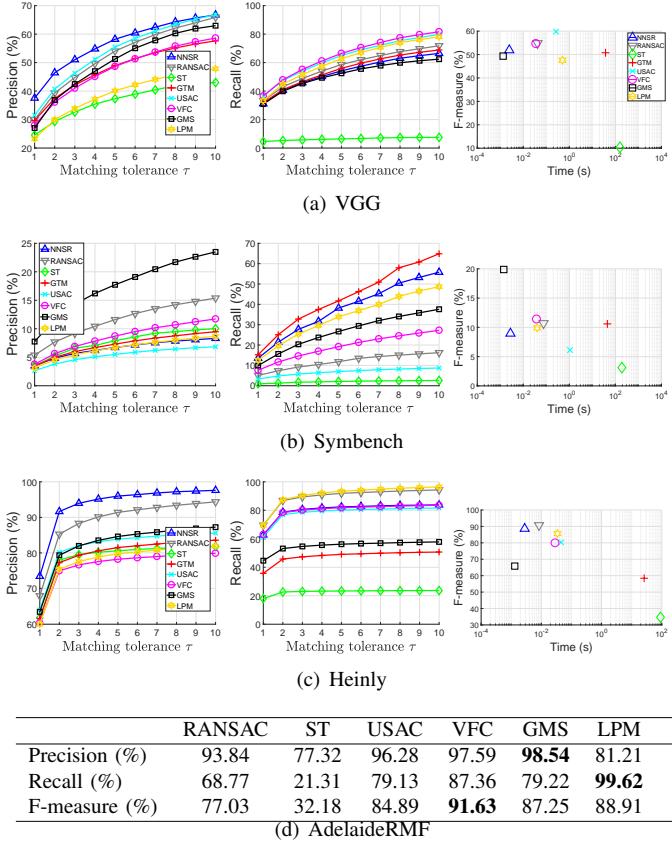
(a) VGG

(b) Symbench

(c) Heinly

|  | RANSAC | ST | USAC | VFC | GMS | LPM |
|---|---|---|---|---|---|---|
| Precision (%) | 93.84 | 77.32 | 96.28 | 97.59 | **98.54** | 81.21 |
| Recall (%) | 68.77 | 21.31 | 79.13 | 87.36 | 79.22 | **99.62** |
| F-measure (%) | 77.03 | 32.18 | 84.89 | **91.63** | 87.25 | 88.91 |

(d) AdelaideRMF

Fig. 3. Performance of the evaluated algorithms on four datasets, i.e., (a) VGG, (b) Symbench, (c) Heinly, and (d) AdelaideRMF, in terms of precision, recall and F-measure under different matching tolerance $\tau$. The maximum values of precision, recall and F-measure are shown in bold face on the AdelaideRMF dataset.

zoom, rotation, light change, and JPEG compression, is independently examined on the VGG dataset. In Sect. V-E, we address concerns about the efficiency in those algorithms by examining their overall time cost on different datasets paired with the speed comparison under different scales of initial matches. Finally, some representative visual results of the evaluated algorithms are shown in Sect. V-F.

## V. RESULTS

Following the experimental arrangement in Sect. IV-D, this section presents the corresponding results together with necessary discussions and explanations.

### A. Performance on the different datasets

In the following, we show the precision, recall and F-measure performance of our evaluated algorithms on different datasets, i.e., under different scenarios. In particular, the overall precision, recall and F-measure curves are shown in Fig. 3 for aggregately view and the F-measure scores for each image pair on the four datasets are shown in Fig. 4 to give a more detailed view. We mainly discuss the performance based on Fig. 3.

*1) Performance on the VGG dataset:* Fig. 3(a) shows outcomes on the VGG dataset. It is interesting to see that NNSR achieves the best precision performance, being marginally better than USAC, RANSAC and GMS. This result is due to the fact that the feature distinctiveness cue is rather selective with rich-textured images, e.g., images in the VGG dataset. On the down side, feature distinctiveness is sometimes ambiguous and not a robust constraint as we can see that the recall of NNSR is just mediocre. It indicates that many correct correspondences have been filtered by NNSR. For ST and LPM, they are generally inferior to the others on this dataset in terms of the F-measure. That is because ST may fail to locate the main cluster in the spectral domain if the ourlier ratio is large, resulting in quite poor recall performance. LPM achieves much better recall performance than ST, while its precision performance is surpassed by most compared ones. It arises from the loose constraint employed in LPM. Overall, USAC is the best method on this dataset. Explanation behind is that USAC is a parametric method and the parametric model of each image pair existed in this dataset can be properly fitted.

*2) Performance on the Symbench dataset:* Fig. 3(b) presents results on the Symbench dataset. All methods suffer a clear drop in performance on this dataset when compared with that on the VGG dataset, which is attributed to light change and various rendering styles. More specifically, we observed that the average inlier ratio of initial correspondences on this dataset is lower than $10\%$. As previously explained, the feature distinctiveness constraint strongly relies on the discriminative power of the local feature descriptor. However, the rendering style variation makes it fairly challenging to maintain descriptiveness in this case. As a result, NNSR delivers very poor precision performance. Another significant difference compared to that on the VGG dataset is USAC's performance. One can see that USAC returns the most and the second most inferior precision and recall performance, respectively. That is because USAC may find empty inlier sets in some cases when its average estimated scores decreases owing to the multiple constraints in this algorithm [11]. In general, GMS and VFC are the two most well-behaved methods after referring their F-measure rankings. A common trait of these two algorithms is that both of them are independent from the descriptor similarity.

*3) Performance on the Heinly dataset:* Fig. 3(c) presents results on the Heinly dataset. Image pairs on this dataset only contain pure zoom or rotation, and we can observe that all methods obtain relatively decent performance on this dataset. In terms of precision, NNSR and RANSAC neatly outperform the others. Regarding recall, LPM and RANSAC are the two best ones. Note that the reason for the high recall of LPM is that most inliers are selected with the loose constraint designed by this algorithm. For NNSR and RANSAC, the former one is attributed to the high distinctiveness of SIFT (we will see its performance variation with less distinctive descriptors in Sect. V-C), whereas the latter one is owing to the powerful homography fitting ability of RANSAC. GMS, due to its sensitivity to large degrees of rotation [17], shows worse results compared to its performance on the VGG and Symbench datasets.
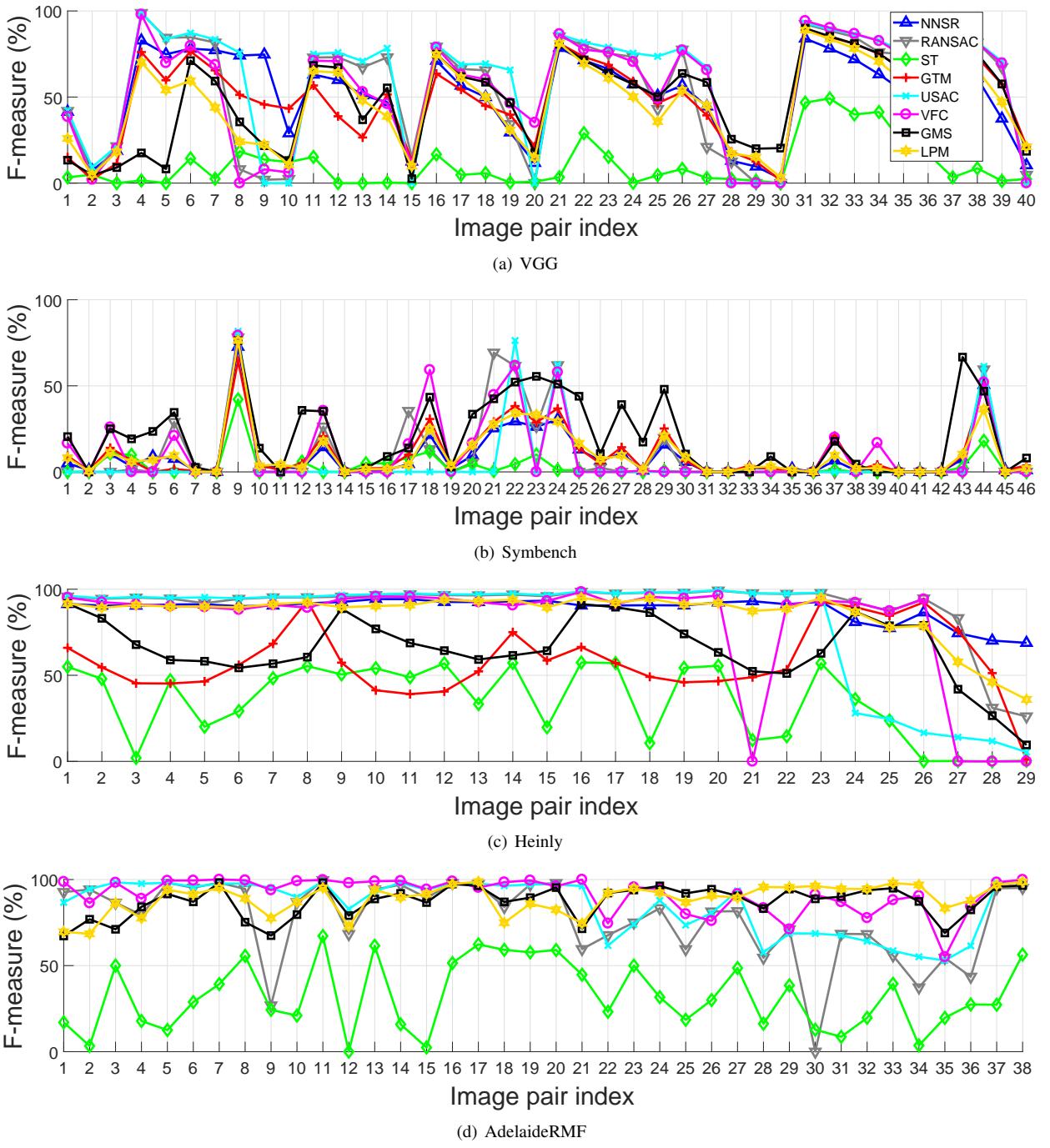
Fig. 4. F-measure performance tested on each image pair on the (a) VGG, (b) Symbench, (c) Heinly, and (d) AdelaideRMF datasets for all evaluated algorithms. Here, the matching tolerance $\tau$ is set as 5.

*4) Performance on the AdelaideRMF dataset:* Fig. 3(d) presents results on the AdelaideRMF dataset. Two explanations should be given on this dataset. First, as only manual labeled ground-truth correspondences are available, we present the exact scores rather than curves with respect to matching tolerance for each method. Second, the keypoints on this dataset are not located by image detectors. Rather, they were labeled manually. Thus, GTM requiring local affine information and NNSR based on auto-detected keypoints are not assessed on this dataset. Since each scene in this dataset contains multiple planes, the fundamental matrix based on the

epipolar geometry constraint is employed to approximate the parametric model for RANSAC and USAC. By observing the scores in Fig. 3(d), one can see that GMS, LPM and VFC achieve the best precision, recall and F-measure performance, respectively. All the three methods are non-parametric. This is reasonable since the AdelaideRMF contains multi-structures, and the parametric assumption for methods like RANSAC and USAC will fail in this case.

*5) Overall performance:* By weighing up the results presented in Fig. 3 and Fig. 4, we can draw the following conclusions. First, the performance of all correspondences

selection algorithms is affected by the initial inlier ratio. For instance, the performance of all algorithms deteriorates dramatically on the Symbench dataset with less than 10% inliers. Second, NNSR simply relying on feature's distinctiveness produces pleasurable results if images are well-textured and clean. Third, parametric approaches, i.e., RANSAC and USAC, prefer the context that the transformation between two images can be well fitted by a parametric model. While non-parametric algorithms perform better in situations without large degrees of rigid/non-rigid transformation. Overall, VFC and RANSAC are the two best algorithms under across-dataset experiments.

### B. Performance on selected matches



(a) Brute force matching
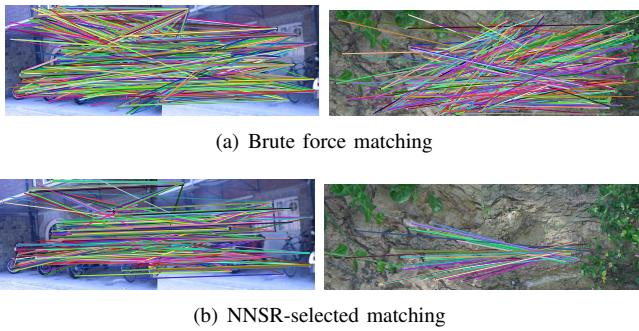


(b) NNSR-selected matching

Fig. 5. Examples of correspondence sets obtained via brute-force matching (a) and NNSR selection (b). Sample image pairs are taken from the VGG dataset.



(a)

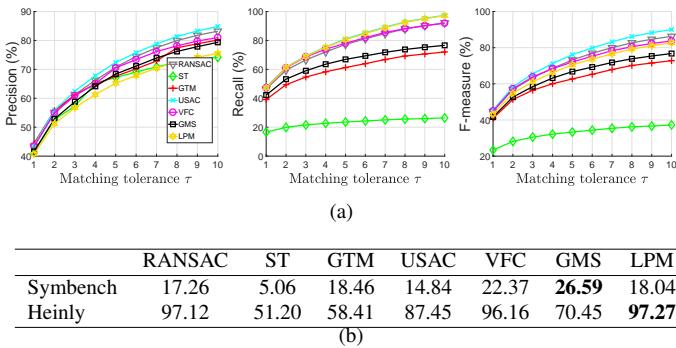| | RANSAC | ST | GTM | USAC | VFC | GMS | LPM |
|---|---|---|---|---|---|---|---|
| Symbench | 17.26 | 5.06 | 18.46 | 14.84 | 22.37 | **26.59** | 18.04 |
| Heinly | 97.12 | 51.20 | 58.41 | 87.45 | 96.16 | 70.45 | **97.27** |

(b)

Fig. 6. Performance of evaluated algorithms (except NNSR) on selected matches on the (a) VGG, (b) Symbench and Heinly datasets. For aggregated view, precision, recall and F-measure curves are shown for the VGG dataset, and F-measure (%) performance under $\tau = 5$ is shown for the Symbench and Heinly datasets. The AdelaideRMF dataset is not tested as NNSR fails to work on this dataset.

Many existing works [8], [23], [43], [44] first prune false correspondences via NNSR and then use parametric or non-parametric methods to for further selection. This experiment then checks this scenario. Remarkably, since NNSR fails to work on the AdelaideRMF dataset, this dataset is not considered in this test. Fig. 5 shows the difference between correspondences before and after applying NNSR, and results using NNSR-selected correspondences for selection are shown in Fig. 6.

On the VGG dataset shown in Fig. 6(a), one can see that the performance of all methods has been improved using

| | | NNSR | RANSAC | ST | USAC | VFC | GMS | LPM |
|---|---|---|---|---|---|---|---|---|
| SIFT + | Symbench | 9.34 | 3.02 | 1.22 | 3.27 | 11.56 | **11.64** | 9.36 |
| SIFT | Heinly | 94.13 | 95.43 | 33.82 | **98.75** | 83.08 | 40.29 | 89.84 |
| ORB + | Symbench | 4.70 | 5.27 | 2.13 | 3.33 | 3.00 | **11.62** | 6.31 |
| ORB | Heinly | 57.62 | 58.98 | 17.57 | 56.45 | 56.30 | 50.24 | **60.30** |
| ASIFT + | Symbench | 7.00 | 7.15 | 3.29 | 4.54 | 14.42 | **17.48** | 12.54 |
| ASIFT | Heinly | 69.31 | **92.31** | 27.47 | 78.72 | 78.75 | 44.21 | 88.21 |
| BLOB + | Symbench | 4.62 | 2.35 | 0.95 | 1.97 | **6.25** | 0.50 | 2.19 |
| FREAK | Heinly | 68.63 | **76.25** | 20.32 | 74.45 | 68.71 | 4.30 | 66.64 |

NNSR-selected matches compared to brute-force matches in Fig. 3(a). Particularly, USAC manages to be the best method regarding precision, recall and F-measure. Also, gaps between most curves excluding that of ST are relatively small. On the Symbench and Heinly datasets, GMS and LPM respectively achieve the best overall performance, where LPM even produces an extremely high F-measure score, i.e., 97.27%, on the Heinly dataset. We can infer that LPM adapts well to initial correspondence sets with high inlier ratio.

### C. Performance under different detectors and descriptors
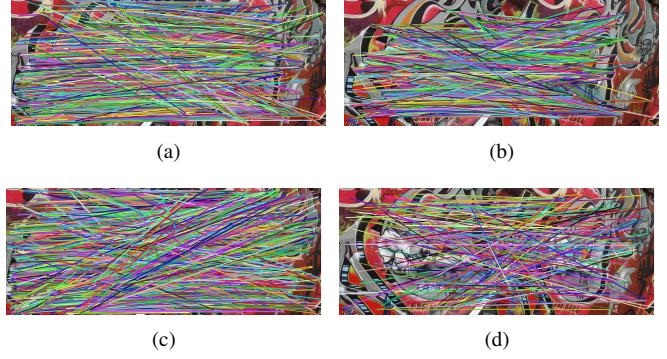


(a)



(b)



(c)



(d)

Fig. 7. Initial correspondences using (a) SIFT + SIFT, (b) ORB + ORB, (c) ASIFT + ASIFT and (d) BLOB + FREAK on an exemplar image pair taken from the VGG dataset.

In addition to Hessian-affine + SIFT, we also consider four other popular detector-descriptor combinations, i.e., SIFT + SIFT [6], ORB + ORB [7], ASIFT + ASIFT [45], and BLOB [46] + FREAK [47]. Fig. 7 shows the initial correspondences with these combinations on a sample image pair. Note that GTM is excluded in this test as it requires local affine information and these detectors do not provide this information. Also, the AdelaideRMF dataset is not considered due to human-labeled keypoints. The results are reported in Fig. 8 and Table III.

A common characteristic of these results is that the best correspondence selection algorithm generally varies with combinations of detector and descriptor. While we can still find some consistencies, e.g., the VFC method achieves pleasurable performance on the VGG dataset in spite of the descriptor-detector combinations. The performance of some methods fluctuates dramatically. For example, NNSR ranks the first with SIFT + SIFT while performs poorly using ASIFT +
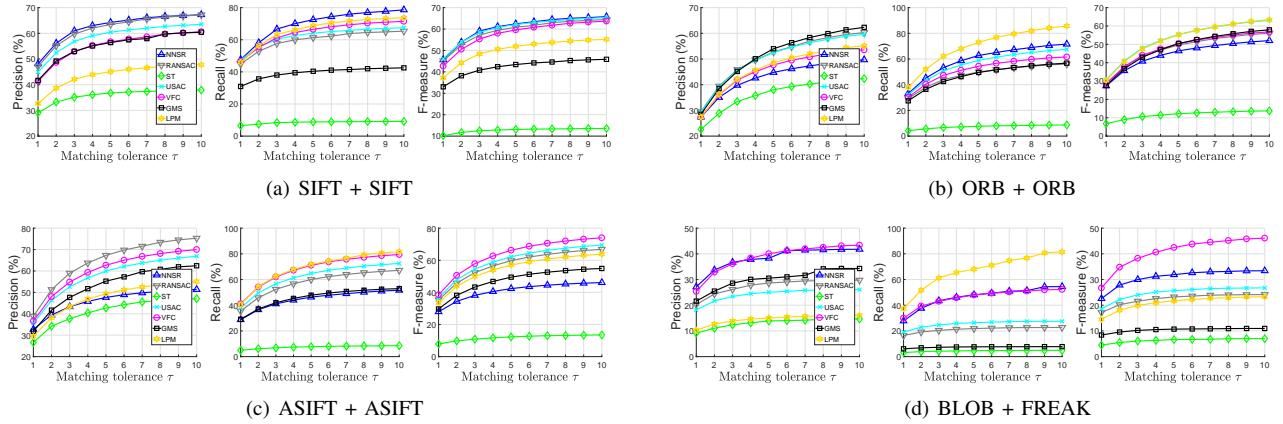
Fig. 8. Performance of evaluated algorithms on the VGG dataset using four different detector-descriptor combinations, i.e., (a) SIFT + SIFT, (b) ORB + ORB, (c) ASIFT + ASIFT, and (d) BLOB + FREAK.

ASIFT on the VGG dataset. On the Symbench and Heinly datasets, GMS and RANSAC are two prominent methods under different kinds of detector-descriptor combinations.

### D. Robustness



Fig. 9. Sample image pairs from the 8 sub-categories of the VGG dataset including (a) zoom and rotation, (b) blur, (c) zoom and rotation, (d) viewpoint change, (e) light change, (f) blur, (g) JPEG compression and (h) viewpoint change.

In this section, we independently evaluate the robustness of these algorithms to a specific nuisance, e.g., zoom, rotation, blur, viewpoint change, light change and JPEG compression on the VGG dataset. Some exemplar images with different nuisances are exhibited in Fig. 9. Results are shown in Table IV.

Under zoom and rotation (case1 and case3), USAC and RANSAC, i.e., two parametric methods, behave the best

(F-measure is referred) mainly attributed to that zoom and rotation are faint impact on homography fitting. Under blur (case2 and case6), GMS and NNSR outperform others. GMS is independent from feature similarity constraint, thus making it rational. For NNSR, it is still explicable as SIFT is very robust to blur. Regarding viewpoint change (case4 and case8), USAC and VFC are the best methods. Note that VFC generally delivers good performance under all kinds of nuisances, being benefited from the consensus search in the non-parametric field. USAC also achieves the best performance under light change (case5) and JPEG compression (case7), being the one that is robust to the broadest categories of nuisances.
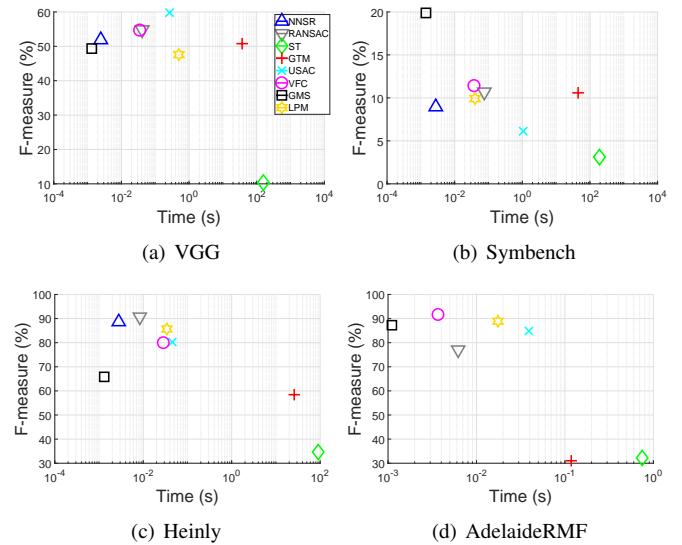
### E. Efficiency



Fig. 10. Efficiency *v.s.* F-measure plots on the (a) VGG, (b) Symbench, (c) Heinly and (d) AdelaideRMF datasets. The efficiency-axis is shown logarithmically for clarity.

To provide an overview of the evaluated methods by taking both selection performance and efficiency into consideration, we present the efficiency *v.s.* F-measure plots on the four

TABLE IV
ROBUSTNESS RESULTS OF EVALUATED ALGORITHMS AGAINST DIFFERENT NUISANCES WITH $\tau = 5$. THE BEST RESULT IS EXPRESSED IN BOLD FACE.

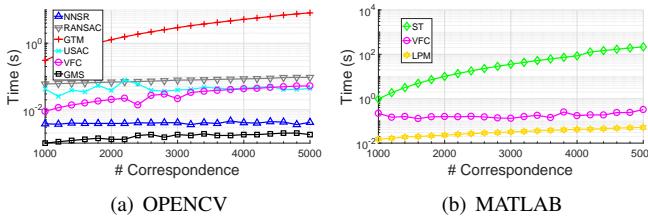| | | NNSR | RANSAC | ST | GTM | UASC | VFC | GMS | LPM |
|---|---|---|---|---|---|---|---|---|---|
| Case1 | Precision | **81.16** | 76.11 | 17.98 | 43.22 | 77.38 | 67.19 | 63.61 | 42.50 |
| (zoom and rotation) | Recall | 77.68 | 92.86 | 4.51 | 79.60 | **99.05** | 86.11 | 11.45 | 83.54 |
| | F-measure | 77.35 | 82.42 | 6.56 | 53.69 | **84.48** | 74.27 | 18.46 | 54.75 |
| Case2 | Precision | **74.57** | 36.87 | 44.23 | 67.00 | 49.66 | 29.44 | 41.71 | 27.73 |
| (blur) | Recall | **79.39** | 41.85 | 8.23 | 56.74 | 60.00 | 51.41 | 50.45 | 54.75 |
| | F-measure | **71.87** | 38.71 | 13.46 | 61.12 | 54.30 | 35.27 | 45.54 | 35.86 |
| Case3 | Precision | 61.53 | **70.54** | 15.97 | 44.92 | 67.41 | 49.38 | 58.57 | 44.59 |
| (zoom and rotation) | Recall | 57.91 | 83.28 | 1.97 | 52.16 | 79.95 | **99.22** | 57.21 | 76.43 |
| | F-measure | 53.74 | **74.81** | 3.50 | 44.83 | 73.12 | 61.91 | 56.35 | 55.10 |
| Case4 | Precision | 51.77 | 55.58 | 37.21 | 50.94 | **63.01** | 57.86 | 57.05 | 45.08 |
| (viewpoint change) | Recall | 61.63 | 66.38 | 3.52 | 68.55 | 79.73 | **97.08** | 75.52 | 83.97 |
| | F-measure | 51.75 | 58.56 | 6.41 | 55.69 | 70.23 | **71.23** | 64.55 | 56.56 |
| Case5 | Precision | 76.28 | 81.44 | 61.90 | 68.90 | **83.76** | 71.99 | 64.89 | 57.65 |
| (light change) | Recall | 63.75 | 86.97 | 6.76 | 80.35 | **100** | **100** | 87.95 | 84.46 |
| | F-measure | 68.00 | 82.34 | 11.61 | 73.94 | **91.11** | 82.49 | 74.37 | 67.90 |
| Case6 | Precision | 31.90 | 45.33 | 24.95 | 33.45 | 32.23 | 31.18 | **57.10** | 26.72 |
| (blur) | Recall | **69.13** | 27.06 | 2.57 | 39.10 | 40.00 | 40.00 | 47.00 | 66.81 |
| | F-measure | 31.49 | 28.86 | 4.34 | 34.29 | 35.68 | 35.02 | **50.80** | 35.82 |
| Case7 | Precision | 89.47 | 87.07 | 89.46 | 80.66 | **89.59** | 89.48 | 79.87 | 75.87 |
| (JPEG compression) | Recall | 61.17 | 97.41 | 28.59 | 94.42 | **100** | **100** | 96.70 | 93.38 |
| | F-measure | 72.42 | 91.81 | 43.07 | 86.88 | **94.43** | 94.26 | 87.25 | 83.43 |
| Case8 | Precision | 67.27 | 74.42 | 52.34 | 72.05 | 73.03 | 72.40 | **80.86** | 62.67 |
| (viewpoint change) | Recall | 61.08 | 79.03 | 4.02 | 80.12 | 80.00 | 79.51 | 73.10 | **81.76** |
| | F-measure | 58.39 | 76.23 | 7.36 | 73.42 | **76.33** | 75.74 | 76.08 | 69.64 |



(a) OPENCV    (b) MATLAB

Fig. 11. Speed comparison of evaluated algorithms with respect to different numbers of initial correspondences. (a) and (b) present the results of methods implemented in OpenCV and MATLAB, respectively. To give a better comparison, VFC is implemented in both platforms. The time-axis is shown logarithmically for clarity.

experimental datasets in Fig. 10. Owing to fast execution speed and overall decent performance, GMS strikes a good balance between selection performance and efficiency.

In order to further test an algorithms's efficiency regarding different numbers of initial correspondences, i.e., the number of initial correspondences may vary in different applications or with different feature detectors, we vary the amount of initial correspondences from 1000 to 5000 and record the average speed of the eight methods. This experiment has been repeated for 10 rounds and average statistics are retained. Because codes of these algorithms are implemented either in OpenCV (C++) or MATLAB, we assess methods within the same platform independently. In addition, the VFC method is evaluated on both platforms and can be a reference for comparing across-platform methods. Results are reported in Fig. 11.

For methods implemented in OpenCV, the efficiency of GMS is beyond all others. That is because GMS involves a grid framework for fast scoring. NNSR ranks the second, as only sort operation is needed to rank correspondences. RANSAC is slightly slower than USAC, and the core time consumption of both methods is dedicated to hypothesis generation-verification. GTM, with the computational complexity of $O(n^2)$ ($n$ being the number of input correspondences), is significantly slower than the other five methods. The margin is rather significant as the number of correspondences increases. For methods implemented in MATLAB, LPM is very efficient as it relies on a simple yet efficient strategy by preserving local neighborhood structure. ST is the most inefficient method, being slower than others by tens of magnitude with dense correspondences. It is due to the fact that the time consumption for computing eigenvalues increases exponentially with the size of the affinity matrix.

### F. Visual results

To obtain a qualitative sense of outputs of evaluated algorithms, we present several visual results of these algorithms on the four experimental datasets in Fig. 12.

Two main observations can be made from the figure. First, distributions of selected correspondences by different algorithms are generally different from each other. For instance, few correspondences are found by GTM on the *bread* in Fig. 12(d). However, NNSR and LPM get plenty of correspondences on it. Second, the quantity of selected correspondences also varies with different methods. In particular, LPM manages to return dense correspondences on most datasets, while ST seeks out much less than others.

### VI. SUMMARY AND DISCUSSION

To give a quick guidance for developers regarding proper algorithms in a specific case, we list the superior and inferior correspondence selection in Table V. Also, peculiarities inherited to each evaluated algorithm are presented as follows:

- **NNSR** is arguably the most straightforward strategy to select correspondences. Its key strength is that repeatable patterns can be removed reliably in certain circumstances, provided that its employed feature detectors can locate the keypoints accurately and descriptors possess strong discriminative power, e.g., SIFT. Also, the high execution

TABLE V
SUMMARIZATION OF SUPERIOR AND INFERIOR CORRESPONDENCE SELECTION ALGORITHMS IN DIFFERENT SCENARIOS BASED ON THE EVALUATION RESULTS. NOTE THAT THIS CONCLUSION IS DRAWN UPON THE F-MEASURE, I.E., THE AGGREGATE PERFORMANCE REGARDING BOTH PRECISION AND RECALL. KEYPOINT DETECTOR AND DESCRIPTOR ARE ABBREVIATED TO DET AND DES, RESPECTIVELY.

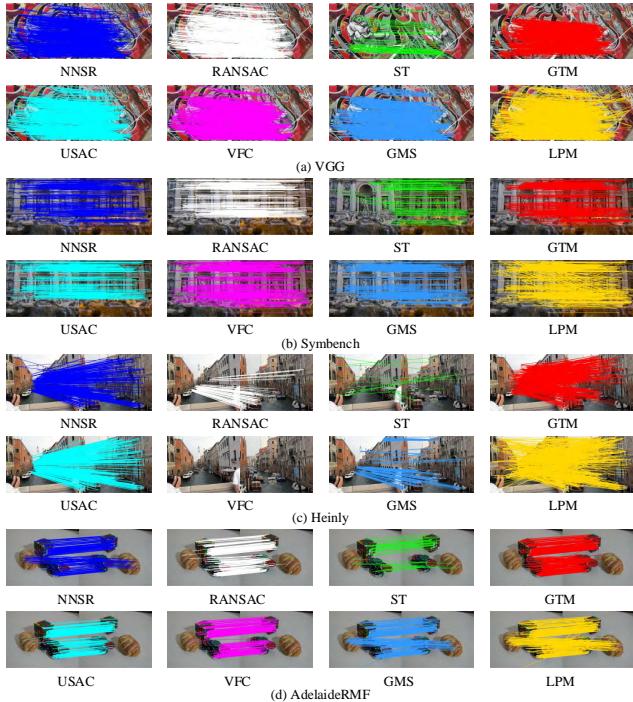| Scenarios | | Superior methods | Inferior methods |
|---|---|---|---|
| Datasets | VGG | USAC, RANSAC, VFC | ST |
| | Symbench | GMS | ST, USAC |
| | Heinly | RANSAC, NNSR, LPM | ST, GTM, GMS |
| | AdelaideRMF | VFC, LPM | ST, RANSAC |
| NNSR pre-selection | VGG | USAC, RANSAC, LPM | ST |
| | Symbench | GMS, VFC | ST |
| | Heinly | LPM, RANSAC, VFC | ST, GMS |
| Det/Des combinations | SIFT+SIFT | USAC, NNSR, GMS | ST, RANSAC |
| | ORB+ORB | LPM, GMS, USAC | ST, VFC |
| | ASIFT+ASIFT | VFC, RANSAC, GMS | ST, USAC, NNSR |
| | BLOB+FREAK | VFC, NNSR, RANSAC | ST, GMS, USAC |
| Robustness | Zoom and rotation | USAC, RANSAC | ST, GTM, GMS |
| | Blur | NNSR, GMS | ST, RANSAC |
| | Viewpoint change | USAC, VFC | ST, NNSR |
| | Light change | USAC, VFC | ST |
| | JPEG compression | USAC, VFC | ST, NNSR |
| Efficiency | | GMS, NNSR | ST, GTM |



Fig. 12. Visual results of evaluated algorithms on examplar image pairs respectively taken from the (a) VGG, (b) Symbench, (c) Heinly and (d) AdelaideRMF datasets. For the best view, lines with different colors represent results of different algorithms.

speed makes it suitable for real-time or near real-time systems. However, the limitation of NNSR is obvious because of the simple descriptor similarity constraint. It is vulnerable when image quality is low (e.g., facing with light change, blur, exposure, and style-transfer) and texture information is limited.

- **RANSAC** and **USAC**, i.e., two evaluated parametric approaches, can fit the parametric models including the homography and fundamental matrices between two images effectively, with the premise that the image pair has homography or epipolar geometry constraint. Thus, they are prior options in such circumstances. Nevertheless, such assumption also brings drawbacks, e.g., when non-rigid objects are captured in images with large scale of parallax or the pure rotation between two camera positions, resulting in the failure of RANSAC and USAC. Further, the reliable models may not be generated by limited iterations with high outlier ratios, which will give rise to expensive time cost. For RANSAC, the minimal-sample models sometimes fall into the local optimization. USAC optimizes over RANSAC, though, it does not guarantee convergence and may produce an empty inlier set due to strict constraints.

- **ST** and **GTM** are methods relying on the affinity matrix computed from initial matches. We can find that these two methods are relatively time-consuming, especially for the ST method. The performance of GTM is much better than ST, mainly because GTM employs local affine information to judge the compatibility of two correspondences. While ST is based on rigid constraint. ST, when inputted with high-quality correspondences, is able to achieve high precision performance (as verified in Sect. V-B). These two methods are optional for off-line applications desiring high precision and with high-quality input.

- **LPM** rejects outliers by the local structure consistency. The constraint item in LPM is relatively loose, resulting in high recall yet relatively low precision. LPM prefers scenarios where the geometric structure information is well preserved between the same local pattern in the image pairs, e.g., small degrees of rigid transformations. Similar to NNSR, it relies strongly on the discriminative power of the feature descriptor. In other words, retrieving the local consistency can be problematic if the local region contains too few inliers. We therefore suggest to choose LPM in the context that has well preserved geometric structures and requires dense correspondences.

- **VFC**, as revealed by our experiment, is the most robust method under all tested scenarios. This is attributed to the fact that VFC is independent from the feature similarity

and parametric models. Specifically, it performs inlier selection in a vector field. VFC generalizes well under different application contexts and can cope with various kinds of nuisances, especially for viewpoint change, light change and JPEG compression.

- **GMS**, similar to VFC, is also independent from the feature similarity and parametric models. However, it assumes that the motion between two images is smooth. Accordingly, it behaves unsatisfactory for image pairs undergoing large degrees of rotation. While if the motion smoothness assumption holds, its performance is superior even for correspondence set with very limited number of inlier, e.g., correspondences generated from the Symbench dataset. Another attractive merit of GMS is the ultra fast execution speed even under several thousands of initial correspondences, making it a prior selection for real-time applications.

## VII. CONCLUSIONS

This paper has comprehensively evaluated eight state-of-the-art image correspondence selection algorithms, covering both parametric and non-parametric families. The experiments addressed several critical issues regarding correspondence selection, e.g., different application scenarios (datasets), inputs from different combinations of feature detector and descriptor, robustness under various challenging conditions including zoom, rotation, blur, viewpoint change, JPEG compression, light change, different rendering styles and multi-structures, and efficiency. Advantages and limitations, in light of experimental outcomes, are summarized so as to guide developers to choose a proper algorithm given a specific scenario.

Remarkably, the performance of most existing algorithms changes dramatically in different scenarios and most methods fail to achieve satisfactory results when the inlier ratio of the initial correspondence set is low. We therefore believe the research should towards the development of correspondence selection algorithms with well generality and be robust to a low inlier rate.

## ACKNOWLEDGMENT

## REFERENCES

[1] N. Snavely, S. M. Seitz, and R. Szeliski, "Modeling the world from internet photo collections," *International Journal of Computer Vision*, vol. 80, no. 2, pp. 189–210, 2008.

[2] S. Benhimane and E. Malis, "Real-time image-based tracking of planes using efficient second-order minimization," in *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 1, 2004, pp. 943–948.

[3] S. Hare, A. Saffari, and P. H. Torr, "Efficient online structured output learning for keypoint-based object tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 1894–1901.

[4] M. Brown and D. G. Lowe, "Automatic panoramic image stitching using invariant features," *International Journal of Computer Vision*, vol. 74, no. 1, pp. 59–73, 2007.

[5] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2. IEEE, 1999, pp. 1150–1157.

[6] ——, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.

[7] 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*, 2011, pp. 2564–2571.

[8] 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.

[9] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[10] O. Chum and J. Matas, "Matching with prosac-progressive sample consensus," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 220–226.

[11] R. Raguram, O. Chum, M. Pollefeys, J. Matas, and J. M. Frahm, "Usac: A universal framework for random sample consensus," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 2022–2038, 2013.

[12] M. Leordeanu and M. Hebert, "A spectral technique for correspondence problems using pairwise constraints," in *Proceedings of the IEEE International Conference on Computer Vision*, 2005, pp. 1482–1489.

[13] 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. 97, no. 1, pp. 36–53, 2012.

[14] T. Collins, P. Mesejo, and A. Bartoli, *An analysis of errors in graph-based keypoint matching and proposed solutions*. Springer International Publishing, 2014.

[15] J. Ma, J. Zhao, H. Guo, J. Jiang, H. Zhou, and Y. Gao, "Locality preserving matching," in *Proceedings of the International Joint Conference on Artificial Intelligence*, 2017, pp. 4492–4498.

[16] X. Li and Z. Hu, "Rejecting mismatches by correspondence function," *International Journal of Computer Vision*, vol. 89, no. 1, pp. 1–17, 2010.

[17] J. Bian, W.-Y. Lin, Y. Matsushita, S.-K. Yeung, T. D. Nguyen, and M.-M. Cheng, "Gms: Grid-based motion statistics for fast, ultra-robust feature correspondence," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[18] W. Y. Lin, F. Wang, M. M. Cheng, S. K. Yeung, P. H. S. Torr, M. N. Do, and J. Lu, "Code: Coherence based decision boundaries for feature correspondence," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 1–1, 2017.

[19] 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.

[20] J. Heinly, E. Dunn, and J. M. Frahm, "Comparative evaluation of binary features," in *Proceedings of the European Conference on Computer Vision*, 2012, pp. 759–773.

[21] H. Aans, A. Dahl, and K. Steenstrup Pedersen, "Interesting interest points: A comparative study of interest point performance on a unique data set," *International Journal of Computer Vision*, vol. 97, no. 1, pp. 18–35, 2012.

[22] 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.

[23] R. Raguram, J. M. Frahm, and M. Pollefeys, "A comparative analysis of ransac techniques leading to adaptive real-time random sample consensus," in *Proceedings of the European Conference on Computer Vision*, 2008, pp. 500–513.

[24] K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. V. Gool, "A comparison of affine region detectors," *International Journal of Computer Vision*, vol. 65, no. 1-2, pp. 43–72, 2005.

[25] N. Snavely and D. C. Hauagge, "Image matching using local symmetry features," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 206–213.

[26] H. S. Wong, T. J. Chin, J. Yu, and D. Suter, "Dynamic and hierarchical multi-structure geometric model fitting," in *Proceedings of the IEEE International Conference on Computer Vision*, 2011, pp. 1044–1051.

[27] 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*. Springer, 2014, pp. 341–356.

[28] P. H. S. Torr and A. Zisserman, "Mlesac: A new robust estimator with application to estimating image geometry," *Computer Vision and Image Understanding*, vol. 78, no. 1, pp. 138–156, 2000.

[29] O. Chum, J. Matas, and J. Kittler, "Locally optimized ransac," *Pattern Recognition*, pp. 236–243, 2003.

[30] M. Cho, J. Lee, and K. M. Lee, "Feature correspondence and deformable object matching via agglomerative correspondence clustering," in *Proceedings of the IEEE International Conference on Computer Vision*. IEEE, 2009, pp. 1280–1287.

[31] T.-J. Chin, J. Yu, and D. Suter, "Accelerated hypothesis generation for multi-structure robust fitting," *Proceedings of the European Conference on Computer Vision*, pp. 533–546, 2010.

[32] H. Y. Chen, Y. Y. Lin, and B. Y. Chen, "Robust feature matching with alternate hough and inverted hough transforms," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 2762–2769.

[33] K. Mikolajczy and C. Schmid, "Scale & affine invariant interest point detectors," *International Journal of Computer Vision*, vol. 60, no. 1, pp. 63–86, 2004.

[34] J. Kim and K. Grauman, "Boundary preserving dense local regions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2011, pp. 1553–1560.

[35] M. Cho, J. Lee, and K. M. Lee, "Reweighted random walks for graph matching," in *Proceedings of the European Conference on Computer Vision*, 2010, pp. 492–505.

[36] F. Tombari, S. Salti, and L. Di Stefano, "Performance evaluation of 3d keypoint detectors," *International Journal of Computer Vision*, vol. 102, no. 1-3, pp. 198–220, 2013.

[37] Y. Guo, M. Bennamoun, F. Sohel, M. Lu, J. Wan, and N. M. Kwok, "A comprehensive performance evaluation of 3d local feature descriptors," *International Journal of Computer Vision*, vol. 116, no. 1, pp. 66–89, 2016.

[38] J. W. Weibull, *Evolutionary game theory*. MIT press, 1997.

[39] J. Matas and O. Chum, "Randomized ransac with sequential probability ratio test," in *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2. IEEE, 2005, pp. 1727–1732.

[40] O. Chum, T. Werner, and J. Matas, "Two-view geometry estimation unaffected by a dominant plane," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 772–779 vol. 1.

[41] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the Royal Statistical Society*, vol. 39, no. 1, pp. 1–38, 1977.

[42] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 9, no. 1, pp. 62–66, 1979.

[43] J. Yang, Z. Cao, and Q. Zhang, "A fast and robust local descriptor for 3d point cloud registration," *Information Sciences*, vol. 346, pp. 163–179, 2016.

[44] J. Yang, Q. Zhang, and Z. Cao, "Multi-attribute statistics histograms for accurate and robust pairwise registration of range images," *Neurocomputing*, vol. 251, pp. 54–67, 2017.

[45] J. M. Morel and G. Yu, "Asift: A new framework for fully affine invariant image comparison," *Siam Journal on Imaging Sciences*, vol. 2, no. 2, pp. 438–469, 2009.

[46] T. Lindeberg, "Feature detection with automatic scale selection," *International Journal of Computer Vision*, vol. 30, no. 2, pp. 79–116, 1998.

[47] A. Alahi, R. Ortiz, and P. Vandergheynst, "Freak: Fast retina keypoint," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 510–517.