

Comparative Evaluation of Binary Features

Jared Heinly, Enrique Dunn, and Jan-Michael Frahm

The University of North Carolina at Chapel Hill
{jheinly,dunn,jmf}@cs.unc.edu

Abstract. Performance evaluation of salient features has a long-standing tradition in computer vision. In this paper, we fill the gap of evaluation for the recent wave of binary feature descriptors, which aim to provide robustness while achieving high computational efficiency. We use established metrics to embed our assessment into the body of existing evaluations, allowing us to provide a novel taxonomy unifying both traditional and novel binary features. Moreover, we analyze the performance of different detector and descriptor pairings, which are often used in practice but have been infrequently analyzed. Additionally, we complement existing datasets with novel data testing for illumination change, pure camera rotation, pure scale change, and the variety present in photo-collections. Our performance analysis clearly demonstrates the power of the new class of features. To benefit the community, we also provide a website for the automatic testing of new description methods using our provided metrics and datasets (www.cs.unc.edu/feature-evaluation).

Keywords: binary features, comparison, evaluation

1 Introduction

Large-scale image registration and recognition in computer vision has led to an explosion in the amount of data being processed in simultaneous localization and mapping [1], reconstruction from photo-collections [2, 3], object recognition [4], and panorama stitching [5] applications. With the increasing amount of data in these applications, the complexity of robust features becomes a hindrance. For instance, storing high dimensional descriptors in floating-point representation consumes significant amounts of memory and the time required to compare descriptors within large datasets becomes longer. Another factor is the proliferation of camera-enabled mobile devices (e.g. phones and tablets) that have limited computational power and storage space. This further necessitates features that compute quickly and are compact in their representation.

This new scale of processing has driven several recent works that propose binary feature detectors and descriptors, promising both increased performance as well as compact representation [6–8]. Therefore, a comparative analysis of these new, state-of-the-art techniques is required. At the same time, the analysis must embed itself into the large body of existing analyses to allow comparison.

In this paper, we provide such an analysis. We rely on established evaluation metrics and develop a new taxonomy of all features. To evaluate traditional and

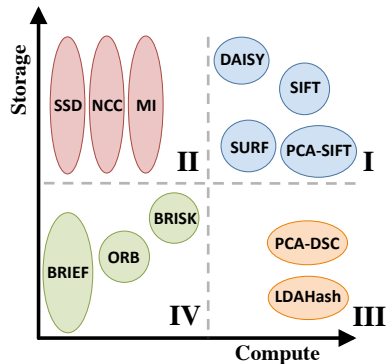


Fig. 1. A taxonomy of descriptors based on their computational and storage requirements: I: Real Value Parameterization [14, 19–21], II: Patch-Based [17], III: Binarized [22, 23], and IV: Binary [6–8].

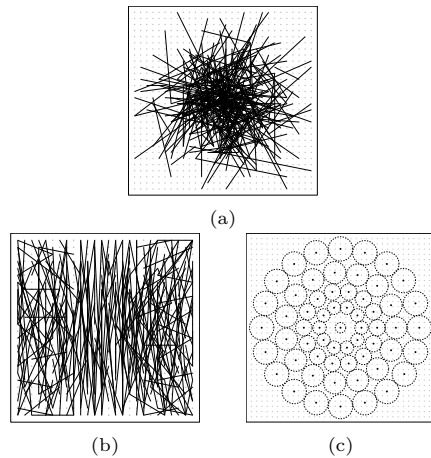


Fig. 2. Example patterns of the (a) BRIEF, (b) ORB, and (c) BRISK descriptors.

binary features, we propose a comprehensive set of metrics that also overcomes limitations in the previously performed evaluations. Additionally, when selecting datasets, we chose standard existing image sequences, but supplemented them with our own custom datasets to test for missing aspects. The new complete benchmarking scenes cover a wide variety of challenges including changes in scale, rotation, illumination, viewpoint, image quality, and occlusion due to the change in viewpoint of non-planar geometry. We also decouple the detection and description phases, and evaluate the pairwise combinations of different detectors and descriptors. As a result of our analysis, we provide practical guidelines for the strengths and weaknesses of binary descriptors. Finally, we provide an evaluation website for the automatic benchmarking of novel features (using the results of new detectors and descriptors on the datasets used in our analysis).¹

2 Related Work

Feature performance (detection, description, and matching) is important to many computer vision applications. In 2005, Mikolajczyk et al. [9] evaluated affine region detectors, and looked to define the repeatability and accuracy of several affine covariant region detectors. They also provided a set of benchmark image sequences (the Oxford dataset) to test the effects of blur, compression, exposure, scale/rotation, and perspective change, which we leverage for compatibility. Also in 2005, Mikolajczyk and Schmid performed an evaluation of local descriptors [10], comparing complex description techniques and several region detectors. It defined two important metrics, *recall* and *1 - precision*. Later, Moreels and Perona [11] evaluated several popular detectors and descriptors by

¹ www.cs.unc.edu/feature-evaluation

analyzing their performance on 3D objects. Strecha et al. [12] published a dense 3D dataset, which we use in our evaluation as it provides 3D LIDAR-based geometry and camera poses. Finally, Aanæs et al. [13] evaluated detectors using a large dataset of known camera positions, controlled illumination, and 3D models.

While many evaluations address the performance of a feature, the recent surge in large-scale feature-based applications draws attention to their runtime requirements. We isolate a feature’s computation expense (detection, description, or matching) and the amount of memory required (to store and use) to establish a taxonomy for feature classification (see Figure 1).

Small-scale applications can often afford large computational and memory requirements (corresponding to real value parameterization). Techniques in this category rely on a parameterization of an image region, where each dimension is a floating-point type (or a discretization of a float, excluding binary). These techniques, examined in [10], use image gradients, spatial frequencies, etc. to describe the local image patch and to test for similarity by using the L^2 norm, Mahalanobis distance, etc. These descriptors have proven to be effective, and tackle issues such as scale, rotation, viewpoint, or illumination variation. The popular SIFT [14] is in this class. However, increased complexity and robustness comes with an increase in computation and storage requirements. High performance, parallel hardware (e.g. graphics processors) can be used to mitigate higher computational expenses, as shown in [15, 16], but even then, descriptor computation can still be the most time-consuming aspect of a system [2].

Therefore, to reduce the computational bottleneck of a system, we address patch-based descriptors. These methods use an image patch surrounding the feature to directly represent it. Distance measures such as sum of squared differences (SSD), normalized cross correlation (NCC), or mutual information (MI) are used to compute pair similarity [17]. The pixels in the patch must be stored, with quadratically increasing requirements for larger patches. However, in large-scale databases (for recognition or reconstruction [18]) the biggest constraints are the matching speed, bandwidth, and the storage required by the descriptor.

The next region in our taxonomy, binarized descriptors, consists of techniques that have high computational but low storage requirements. Binarized descriptors rely on hashing techniques to reduce high-dimensional, real-value parameterizations into compact binary codes [22–25]. While reducing storage constraints, it also speeds up comparison times through the use of the Hamming distance measure. However, computational requirements are still high as the full real-value parameterization must be computed before the hashing can occur.

The final region in our taxonomy is the binary descriptors. These descriptors have a compact binary representation and limited computational requirements, computing the descriptor directly from pixel-level comparisons. This makes them an attractive solution to many modern applications, especially for mobile platforms where both compute and memory resources are limited.

Descriptor	Detector	Rotation Invariant	Scale Invariant
BRIEF	Any	No	No
ORB	FAST	Yes	No
BRISK	AGAST	Yes	Yes

Table 1. Overview of the basic properties of the binary descriptors.

3 Survey of Binary Descriptors

Recent works have focused on providing methods to directly compute binary descriptors from local image patches. We survey these proposed techniques below.

Common Principles: In our analysis, we found that all of the recent binary descriptors possess the same following properties:

- the descriptor is built from a set of pairwise intensity comparisons
- each bit in the descriptor is the result of exactly one comparison
- the sampling pattern is fixed (except for possible scale and rotation)
- Hamming distance is used as a similarity measure

While the binary descriptors use these basic principles, each adds its own unique properties to achieve its design goals. We detail the differences between the binary descriptors (Table 1 highlights their key properties).

BRIEF (Binary Robust Independent Elementary Features) are proposed by Calonder et al. [6], and are the simplest of the methods. It uses a sampling pattern consisting of 128, 256, or 512 comparisons (equating to 128, 256, or 512 bits), with sample points selected randomly from an isotropic Gaussian distribution centered at the feature location (see Figure 2(a). Calonder et al. [6] use the SURF detector negating their computational gain, but BRIEF can be coupled with any other detector. Calonder et al. suggest to use BRIEF with the efficient CenSurE detector [26]. Given its simple construction and compact storage, BRIEF has the lowest compute and storage requirements.

ORB (Oriented FAST and Rotated BRIEF) was proposed by Rublee et al. [7], and overcomes the lack of rotation invariance of BRIEF. ORB computes a local orientation through the use of an intensity centroid [27], which is a weighted averaging of pixel intensities in the local patch assumed not be coincident with the center of the feature. The orientation is the vector between the feature location and the centroid. While this may seem unstable, it is competitive with the single orientation assignment employed in SIFT [28].

The sampling pattern employed in ORB uses 256 pairwise intensity comparisons, but in contrast to BRIEF, is constructed via machine learning, maximizing the descriptor’s variance and minimizing the correlation under various orientation changes (see Figure 2(b) for an example).

When using this descriptor, ORB proposes to use the FAST corner detector [29], noting that FAST does not provide a good measure of cornerness and lacks robustness to multi-scale features. In order to combat this, the Harris [30] corner measure is applied at each keypoint location to provide non-maximal suppression within the image, and a limited image scale pyramid is used to detect keypoints of varying sizes. However, no non-maxima suppression is used between the scales, resulting in potential duplicate detections within different pyramid levels.

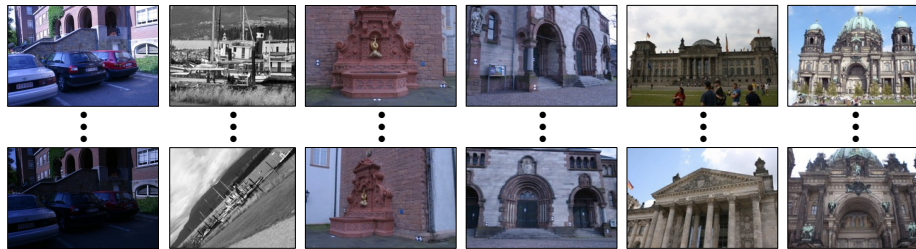


Fig. 3. Example images from the evaluation datasets. From left to right, the datasets are Leuven, Boat, fountain-P11, Herz-Jesu-P8, Reichstag, and Berliner-Dom.

BRISK (Binary Robust Invariant Scalable Keypoints) was proposed by Leutenegger et al. [8] and provides both scale and rotation invariance. In order to compute the feature locations, it uses the AGAST corner detector [31], which improves FAST by increasing speed while maintaining the same detection performance. For scale invariance, BRISK detects keypoints in a scale-space pyramid, performing non-maxima suppression and interpolation across all scales.

To describe the features, the authors turn away from the random or learned patterns of BRIEF and ORB, and instead use a symmetric pattern. Sample points are positioned in concentric circles surrounding the feature, with each sample point representing a Gaussian blurring of its surrounding pixels. The standard deviation of this blurring is increased with the distance from the center of the feature (see Figure 2(c) for an illustration). This may seem similar to the DAISY descriptor [20], but the authors point out that DAISY was designed specifically for dense matching, and captures more information than is needed for keypoint description.

To determine orientation, several long-distance sample point comparisons (e.g. on opposite sides of the descriptor pattern) are used. For each long-distance comparison, the vector displacement between the sample points is stored and weighted by the relative difference in intensity. Then, these weighted vectors are averaged to determine the dominant gradient direction of the patch. The sampling pattern is then scaled and rotated, and the descriptor is built up of 512 short-distance sample point comparisons (e.g. a sample point and its closest neighbors) representing the local gradients and shape within the patch.

Overall, BRISK requires significantly more computation and slightly more storage space than either BRIEF or ORB. This places it in the higher compute, higher storage region of the binary descriptor category of our taxonomy.

4 Evaluation

We analyze the performance characteristics of the three recent binary descriptors (BRIEF, ORB, and BRISK), while using state-of-the-art full parametrization descriptors (SIFT and SURF [19]) as a baseline. Besides traditional descriptor performance metrics, we also evaluate the correlation between the detector and descriptor with respect to matching performance. This has received some attention

in the past [32], but is important to investigate given its significant practical implications in descriptor post-processing methods (for instance, RANSAC-based estimations [33]). One of the contributions of our comparative analysis is that we test the performance of the studied descriptors with a diverse set of keypoint detectors (Harris, MSER, FAST, ORB, BRISK, SURF, and SIFT).

Datasets: To ensure our work is compatible with existing analyses, we use existing datasets to evaluate the performance of the binary descriptors under various transformations. Specifically, we relied on the Oxford dataset provided and described by Mikolajczyk et al. [9] (for evaluations of the effects of image blur, exposure, JPEG compression, combined scale and rotation, and perspective transformations of planar geometry), and the fountain-P11 and Herz-Jesu-P8 datasets from Strecha et al. [12] (for evaluating the effects of perspective transformations of non-planar geometry). Moreover, we complement the existing datasets with our own to test for pure rotation, pure scaling, illumination changes, and the challenges posed by photo collection datasets such as white balance, auto-exposure, image quality, etc. These various datasets also enable us to isolate the effects of each transformation, or in some cases pairs of transformations.

Performance Metrics: Mikolajczyk et al. [9, 10] propose to use the metrics of *recall*, *repeatability*, and $1 - \textit{precision}$. They describe useful characteristics of a feature’s performance, and are widely used as standard measures. However, there are some subtleties to a feature’s performance that are missed by only using these measures. For instance, they fail to capture information about the spatial distribution of the features, as well as the frequency of candidate matches.

We wanted to not only use a comprehensive set of metrics that allow us to embed our analysis into the existing body of work, but we also aimed at evaluating parameters relevant to algorithms relying on the features. As such, we propose a set a five different metrics: *putative match ratio*, *precision*, *matching score*, *recall*, and *entropy*.

The putative match ratio, $\textit{Putative Match Ratio} = \# \textit{Putative Matches} / \# \textit{Features}$, addresses the selectivity of the descriptor and describes what fraction of the detected features will be initially identified as a match (though potentially incorrect). We define a putative match to be a single pairing of keypoints, where a keypoint cannot be matched to more than one other. Keypoints that are outside of the bounds of the second image (once they are transformed based on the true camera positioning) are not counted.

The value of the putative match ratio is directly influenced by the matching criteria. A less restrictive matching criteria will generate a higher putative match ratio, whereas a criteria that is too restrictive will discard potentially valid matches and will decrease the putative match ratio. Another influencing factor is the distinctiveness of the descriptors under consideration. If many descriptors are highly similar (have small distance values between them), this creates confusion in the matching criteria and can drive down the putative match ratio.

The precision, $\textit{Precision} = \# \textit{Correct Matches} / \# \textit{Putative Matches}$ [10]), defines the number of correct matches out of the set of putative matches

(the inlier ratio). In this equation, the number of correct matches are those putative matches that are geometrically verified based on the known camera positions. The ratio has significant performance consequences for robust estimation modules that use feature matches, such as RANSAC [33], where execution times increase exponentially as the inlier ratio decreases. It is also influenced by many of the same factors that influenced the putative match ratio, but the consequences are different. For instance, while a less restrictive matching criteria will increase the putative match ratio, it will decrease the precision as a higher number of incorrect matches will be generated. Additionally, highly similar descriptors, which drove down the putative match ratio, will also decrease the precision, as confusion in the matching step will also generate a higher number of incorrect matches.

The matching score, $Matching\ Score = \#Correct\ Matches / \#Features$ [10] is equivalent to the multiplication of the putative match ratio and precision. It describes the number of initial features that will result in correct matches, and like the previous two metrics, the matching score can be influenced by indistinct descriptors and the matching criteria. Overall, the matching score describes how well the descriptor is performing and is influenced by the descriptor’s robustness to transformations present in the data.

Recall, $Recall = \#Correct\ Matches / \#Correspondences$ (defined in [10]), quantifies how many of the possible correct matches were actually found. The correspondences are the matches that should have been identified given the keypoint locations in both images. While this value is dependent on the detector’s ability to generate correspondences, recall shares the same influences as the matching score. For instance, a low recall could mean that the descriptors are indistinct, the matching criterion is too strict, or the data is too complex.

The final metric, entropy (used by Zitnick and Ramnath [34]) addresses the influence of the detector on a descriptor. The purpose of this metric is to compute the amount of spread or randomness in the spatial distribution of the keypoints in the image. This is important as too little spread increases the possibility of confusion in the descriptor matching phase due to keypoint clusters.

To compute the entropy, we create a 2D evenly-spaced binning of the feature points. Each point’s contribution to a given bin is weighted by a Gaussian relative to its distance to the bin’s center. A bin $b(p)$ at position $p = (x, y)$ equals $b(p) = \frac{1}{Z} \sum_{m \in M} G(\|p - m\|)$ where m is a keypoint in the full set M of detected keypoints, and G is a Gaussian. A constant of $1/Z$ is added so that the sum of all bins evaluates to 1. This binning allows us to compute the entropy: $Entropy = \sum_p -b(p) \log b(p)$. Even though entropy is dataset dependent, the relative value of the entropy is useful in identifying detector spatial distribution behaviors such as non-random clustering of keypoint detections (as seen in Figure 4).

We can now quantify how many features were reported as matching (putative match ratio), how many actually matched (precision and matching score), how many were matched out of those possible (recall), and the spread of the detector’s keypoints (entropy). Next, we will describe our evaluation framework.

Test Setup: For each dataset, we detect and describe features in each of its images, and use the features in the first image as a reference when matching to all further images in the set. In order to test the detector/descriptor pairings, we had to address several subtleties. For instance, there was a mismatch when a scale invariant descriptor was combined with a detector that was not scale invariant, and vice versa. Additionally, combining detectors and descriptors that were both scale invariant was not trivial, as they typically each use their own method of defining what a feature scale is. In both cases, we simply discarded the scale information, and computed the descriptor at the native image resolution.

In addition, we addressed mismatches in rotation invariance between various detectors and descriptors. We solved this by isolating the orientation computation to the description phase, so that the descriptor overrides any orientation provided by the detector. We next discuss the matching criteria that we use.

Match Criteria: To compute putative matches, we adopt a ratio style test that has proven to be effective [11], [14]. This test compares the ratio of distances between the two best matches for a given keypoint, and rejects the match if the ratio is above a threshold of 0.8 for all tests (the same used in [14]).

To determine the correctness of a match, we use ground truth data to warp keypoints from the first image of the dataset into all remaining images. The warping is achieved either by homographies (provided in the Oxford dataset [9] as well as our non-photo-collection supplemental datasets) or by using ground truth 3D geometry (provided in the Strecha dataset [12]) to project the points into the known cameras. In both cases, match points that are within 2.5 pixels of each other are assumed to be correct. This threshold was chosen empirically as it provided a good balance between pixel-level corner detection (such as Harris [30] and FAST [29]), and the centers of blobs in blob-style detectors (SIFT [14] and SURF [19]). For the photo-collection datasets, we use known camera positions to project points from the first image as epipolar lines in the other images, and once again apply a 2.5 pixel distance threshold test. While we could have used the three camera arrangement test described by Moreels and Perona [11], we opted for the single epipolar line approach as it allowed us to compare only a given pair of images, without the need for a third image to provide correspondences and closely mimics typical uses of features.

5 Analysis and Results

In our analysis we tested the individual effects of various geometric and photometric transformations, as well as several of their combinations to gain a better and more complete understanding of the detectors' and descriptors' performance. Figure 5 provides the results for all tested dataset categories. The individual values making up each measurement (# putative matches, features, etc.) were first summed across each pairwise comparison, and one final average was computed.

Detector Performance: As mentioned before, entropy can be used as a measure of the randomness of the keypoint locations, penalizing detectors that spa-

Detector/Descriptor	BRIEF	ORB	BRISK	SURF	SIFT	Harris	MSER	FAST
Avg # Features	n/a	13427	7771	3766	4788	2543	693	8166
Detector Entropy	n/a	12.10	12.33	12.26	12.34	11.84	10.74	12.52
Detector ms/image	n/a	17	43	377(19)	572(25)	78(4.7)	117	2.7
Descriptor μ s/feature	4.4(0.4)	4.8	12.9	143(6.6)	314(19)	n/a	n/a	n/a
Storage bytes/feature	16,32,64	32	64	64(256)	128(512)	n/a	n/a	n/a

Table 2. This table provides statistics for the detectors and descriptors used in our evaluation. For storage, we used a 32 byte BRIEF, and the values in parenthesis for SURF and SIFT are the number of required bytes if the descriptors are stored as floats. For the timings, values in parenthesis are GPU implementations ([15], [16], or our own).



Fig. 4. This figure shows the distribution of keypoints (for two datasets) from the Harris (left image in the pair) and FAST (right image in the pair) corner detectors. The FAST detector not only detects more keypoints, but has a higher entropy for its detections.

tially cluster their keypoints. In order to compute the entropy across all of the datasets, we perform a weighted average of the individual entropies to account for the different contributions of the datasets (results are in Table 2).

We see several notable attributes. First, FAST has the highest entropy, which can be attributed to its good spread of points and the sheer number of detections that occurred. A higher number of keypoints does not necessarily correspond to increased entropy, but it can help if more spatial bins in the entropy equation are populated. The second group of detectors, SIFT, BRISK, and SURF, have the next highest entropies, which is expected as Zitnick’s and Ramnath’s evaluation [34] noted that a blob-style detector (difference of Gaussian) have high entropies.

On the other end, Harris and MSER reported the lowest two entropy values. The primary reason for this is a lower average number of keypoint detections (especially in the case of MSER), which leads to less spatial bins being populated. While these detectors are still very viable, it is important to note their lower detection rates and potential differences in keypoint distribution (one such example is provided in Figure 4).

Descriptor Performance: One of the most compelling motivations for the use of binary descriptors is their efficiency and compactness. Table 2 shows the timing results for the various detectors and descriptors that we used in our system. The implementations for ORB and BRISK were obtained from the authors, while all others came from OpenCV 2.3 [35]. The code was run on a computer with an Intel Xeon 2.67GHz processor, 12GB of RAM, NVIDIA GTX 285, and Microsoft Windows 7. From the results, we can see that the binary descriptors (and many of their paired detectors) are an order of magnitude faster than the SURF or SIFT alternatives. For mobile devices, while the overall

timings would change, a speedup would still be realized as the binary descriptors are algorithmically more efficient. Table 2 also lists the storage requirements for the descriptors. In some cases, the binary descriptors can be constructed such that they are on par with efficient SURF representations, but overall, binary descriptors reduce storage to a half or to a quarter. We are assuming that the real value parameterization descriptors are stored in a quantized form (1 byte per dimension). If instead they are stored as floating-point values, the storage savings of binary features are even more significant.

Figure 5 shows the results for our evaluation of detector and descriptor pairings. Upon a close inspection, several interesting conclusions become apparent.

Non-Geometric Transforms: Non-geometric transforms consist of those that are image-capture dependent, and do not rely on the viewpoint (e.g. blur, JPEG compression, exposure, and illumination). Challenges in these datasets involve less-distinct image gradients and changes to the relative difference between pixels due to compression or lighting change.

Upon inspecting the results, BRIEF’s performance is very favorable. It outperforms ORB and BRISK in many of the categories (even when compared to SURF or SIFT), though the precision is noticeably lower than ORB or BRISK. The key insight to BRIEF’s good performance is that it has a fixed pattern (no scale or rotation invariance), and is designed to accurately represent the underlying gradients of an image patch (which are stable for monotonic color transformations). In regard to the lower precision, as we mentioned before, the precision is decreased when the matching criteria are either less restrictive, or the features are not distinct enough. We enforced the same matching criteria for each descriptor, hence performance differences are due to the lack of distinctiveness. The pattern used by ORB is specifically trained to maximize distinctiveness, and BRISK’s descriptor is much more complex (larger number of samples with each sample being a comparison of blurred pixels), which allows it to be much more certain of the matches it generates. Therefore, it is not surprising that BRIEF’s precision would suffer slightly compared to ORB or BRISK.

For the choice of detector, both SURF and Harris perform well under these non-geometric transformations. The key is that both are very repeatable in these images. SURF (being a blob detector) does well for the blur dataset, as the smoothing of the gradients negatively impacts corner detectors like Harris. SIFT (also being a blob detector) does not do as well as the accuracy of its keypoint detections decreases as the scale of the blob increases. This is an artifact of the downsampling in the image pyramid, where SURF overcomes this by applying larger filters at the native resolution. For the non-blurred datasets, Harris does well because high gradient change still exists at corners under color changes.

Affine Image Transforms: The affine image transforms that we used in our testing consist of image plane rotation and scaling. Overall, SIFT was the best, but ORB and BRISK still performed well. As is expected, BRIEF performs poorly under large rotation or scale change compared to the other descriptors (as it is not designed to handle those effects). For pure scale change, ORB is better than BRIEF although there is no rotation change in the image. This can

be attributed to the limited scale space used by the ORB detector, allowing it to detect FAST corners at several scales. However, of all the binary descriptors, BRISK performs the best, as it is scale invariant.

For pure rotations, the ORB detector/descriptor combination performs better than the BRISK detector/descriptor. However, the FAST detector paired with BRISK is a top performer (as well as the Harris detector paired with ORB), competing on the same level as SIFT. The insight to this result is that both FAST and Harris are not scale invariant, and excel because there is no scale change in the images. The better performances of BRISK and ORB when paired with other detectors are unexpected, but only help to highlight the importance of considering various combinations when deciding on a detector and descriptor.

Finally, when analyzing combined scale and rotation changes, BRISK takes the lead of the binary descriptors. However, SIFT and SURF have competitive matching scores, but significantly higher recalls. This difference in matching score and recall speaks to the performance of the detectors. They detected almost the same percentage of correct matches (the matching score), SURF and SIFT's number of correspondences must have been lower in order for their recall to be higher. This means that for basic scale and rotation combinations, BRISK's detector is more repeatable than SURF or SIFT's.

Turning our attention back to BRIEF, one way to overcome BRIEF's sensitivity to rotation when computing on mobile devices is via inertial sensor-aligned features [36] which aligns a descriptor with the current gravity direction of the mobile device (assuming that the device has an accelerometer), allowing for repeatable detections of the same feature as the device changes orientation.

Perspective Transforms are the result of changes in viewpoint. The biggest issues faced in these sets of images are occlusions due to depth edges, as well as perspective warping. For binary descriptors, BRIEF surprisingly has a slight lead in recall and matching score over ORB and BRISK, given its limited complexity. However, in most of the perspective evaluation datasets, there is no significant change in scale, and the orientation of the images are the same. This is not the case in the photo-collection datasets, which include a higher variety of feature scales. Even then, BRIEF still takes the lead over the other binary descriptors. So, even though there is considerable viewpoint change throughout all of the perspective datasets, the upright nature of all of the images allows the fixed orientation of BRIEF to excel.

Additionally, it is interesting to note that the Harris corner detector does well when coupled with BRIEF. Harris is not scale invariant, but the nature of corners is that there will be high gradient change even under perspective distortion allowing for repeatable detections.

Another observation is that the precision of BRISK is once again higher than its binary competitors. This points to the descriptiveness of BRISK, allowing it to achieve very high quality matches. Once again, however, SIFT provides the best performance, and proves to be the most robust to changes in perspective.

6 Conclusion

From the analysis of our results, we highlight several pertinent observations, each of which was driven by the metrics used in our evaluation. Moreover, by proposing and integrating a comprehensive set of measures and datasets, we are able to gain insights into the various factors that influence a feature’s behavior.

First, consider BRIEF’s performance under non-geometric and perspective transforms. In both cases, BRIEF excelled because the images contained similar scales and orientations. This highlights the importance of leveraging any additional knowledge about the data being processed (e.g. similar scales/orientations or orientation from inertial sensors) as it impacts the performance of the binary descriptors. The key idea is that a binary descriptor will suffer in performance when it takes into account a transform not present in the data.

As a result of our evaluation of detector/descriptor pairings, we observed that the best performance did not always correspond to the original authors’ recommendations. This leads us to our second observation that by considering the effect a detector has on a descriptor, we enable the educated assignment of detectors to descriptors depending on the expected properties of the data.

Finally, for all datasets except the non-geometric transforms, SIFT was the best. Although achieving a gain in matching rate, some applications using SIFT may unnecessarily forfeit the significant speed gains made possible by binary descriptors (Table 2). This performance tradeoff makes binary features a viable choice whenever the application provides robustness against additional outliers.

7 Acknowledgments

This work was supported by NSF IIS-0916829 and the Department of Energy under Award DE-FG52-08NA28778.

References

1. Chang, H. J., et al.: P-SLAM: Simultaneous Localization and Mapping With Environmental-Structure Prediction. *IEEE Trans. Robot.* 23(2), 281–293 (2007)
2. Frahm, J.M., et al.: Building Rome on a Cloudless Day. *ECCV*, 368–381 (2010)
3. Snavely, N., Seitz, S. M., Szeliski, R.: Photo Tourism: Exploring Photo Collections in 3D. *SIGGRAPH Conference Proceedings*, 835–846 (2006)
4. Nister, D., Stewenius, H.: Scalable Recognition with a Vocabulary Tree. *CVPR*, 2161–2168 (2006)
5. Brown, M., Lowe, D. G.: Recognising Panoramas. *ICCV*, 1218–1225 (2003)
6. Calonder, M., Lepetit, V., Fua, P.: BRIEF: Binary Robust Independent Elementary Features. *ECCV*, 778–792 (2010)
7. Rublee, E., Rabaud, V., Konolige, K., Bradski, G.: ORB: An Efficient Alternative to SIFT or SURF. *ICCV*, 2564–2571 (2011)
8. Leutenegger, S., Chli, M., Siegwart, R.: BRISK: Binary Robust Invariant Scalable Keypoints. *ICCV*, 2548–2555 (2011)

9. Mikolajczyk, K., et al.: A Comparison of Affine Region Detectors. *IJCV* 65(1-2), 43–72 (2005)
10. Mikolajczyk, K., Schmid, C.: A Performance Evaluation of Local Descriptors. *IEEE Trans. PAMI* 27(10), 1615–1630 (2005)
11. Moreels, P., Perona, P.: Evaluation of Features Detectors and Descriptors based on 3D Objects. *IJCV* 73(3), 263–284 (2007)
12. Strecha, C., et al.: On Benchmarking Camera Calibration and Multi-View Stereo for High Resolution Imagery. *CVPR*, 1–8 (2008)
13. Aanæs, H., et al.: Interesting Interest Points. *IJCV* 97(1), 18–35 (2012)
14. Lowe, D. G.: Distinctive image features from scale-invariant keypoints. *IJCV* 60(2), 91–110 (2004)
15. Wu, C.: SiftGPU. <http://cs.unc.edu/~ccwu/siftgpu> (2007)
16. Schulz, A., et al.: CUDA SURF - A real-time implementation for SURF. <http://www.d2.mpi-inf.mpg.de/surf> (2010)
17. Hirschmüller, H., Scharstein, D.: Evaluation of Cost Functions for Stereo Matching. *CVPR*, 1–8 (2007)
18. Agarwal, S., et al.: Building Rome in a Day. *ICCV*, 72–79 (2009)
19. Bay, H., et al.: Speeded-Up Robust Features (SURF). *Comp. Vis. and Image Understanding* 110(3), 346–359 (2008)
20. Tola, E., Lepetit, V., Fua, P.: DAISY: An Efficient Dense Descriptor Applied to Wide Baseline Stereo. *IEEE Trans. PAMI* 32(5), 815–830 (2010)
21. Ke, Y., Sukthankar, R.: PCA-SIFT: a more distinctive representation for local image descriptors. *CVPR*, 506–513 (2004)
22. Strecha, C., et al.: LDAHash: Improved Matching with Smaller Descriptors. *IEEE Trans. PAMI* 34(1), 66–78 (2012)
23. Yeo, C., Ahammad, P., Ramchandran, K.: Coding of Image Feature Descriptors for Distributed Rate-efficient Visual Correspondences. *IJCV* 94(3), 267–281 (2011)
24. Raginsky, M., Lazebnik, S.: Locality-Sensitive Binary Codes from Shift-Invariant Kernels. *Advances in Neural Info. Processing Systems*, 1509–1517 (2009)
25. Gong, Y., Lazebnik, S.: Iterative Quantization: A Procrustean Approach to Learning Binary Codes. *CVPR*, 817–824 (2011)
26. Agrawal, M., Konolige, K., Blas, M.: CenSurE: Center Surround Extremas for Realtime Feature Detection and Matching. *ECCV*, 102–115 (2008)
27. Rosin, P. L.: Measuring Corner Properties. *Comp. Vis. and Image Understanding*, 291–307 (1999)
28. Gauglitz, S., et al.: Improving Keypoint Orientation Assignment. *BMVC*, (2011)
29. Rosten, E., Porter, R., Drummond, T.: Faster and Better: A Machine Learning Approach to Corner Detection. *IEEE Trans. PAMI* 32(1), 105–119 (2010)
30. Harris, C., Stephens, M.: A Combined Corner and Edge Detector. *Proc. of The Fourth Alvey Vision Conference*, 147–151 (1988)
31. Mair, E., et al.: Adaptive and Generic Corner Detection Based on the Accelerated Segment Test. *ECCV*, 183–196 (2010)
32. Dahl, A. L., Aanæs, H., Pedersen, K. S.: Finding the Best Feature Detector-Descriptor Combination. *3DIMPVT*, 318–325 (2011)
33. Fischler, M. A., Bolles, R. C.: Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography. *Communications of the ACM* 24(6), 381–395 (1981)
34. Zitnick, L., Krishnan, R.: Edge Foci Interest Points. *ICCV*, 359–366 (2011)
35. Bradski, G.: The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, (2000)
36. Kurz, D., BenHimane, S.: Inertial sensor-aligned visual feature descriptors. *CVPR*, 161–166 (2011)

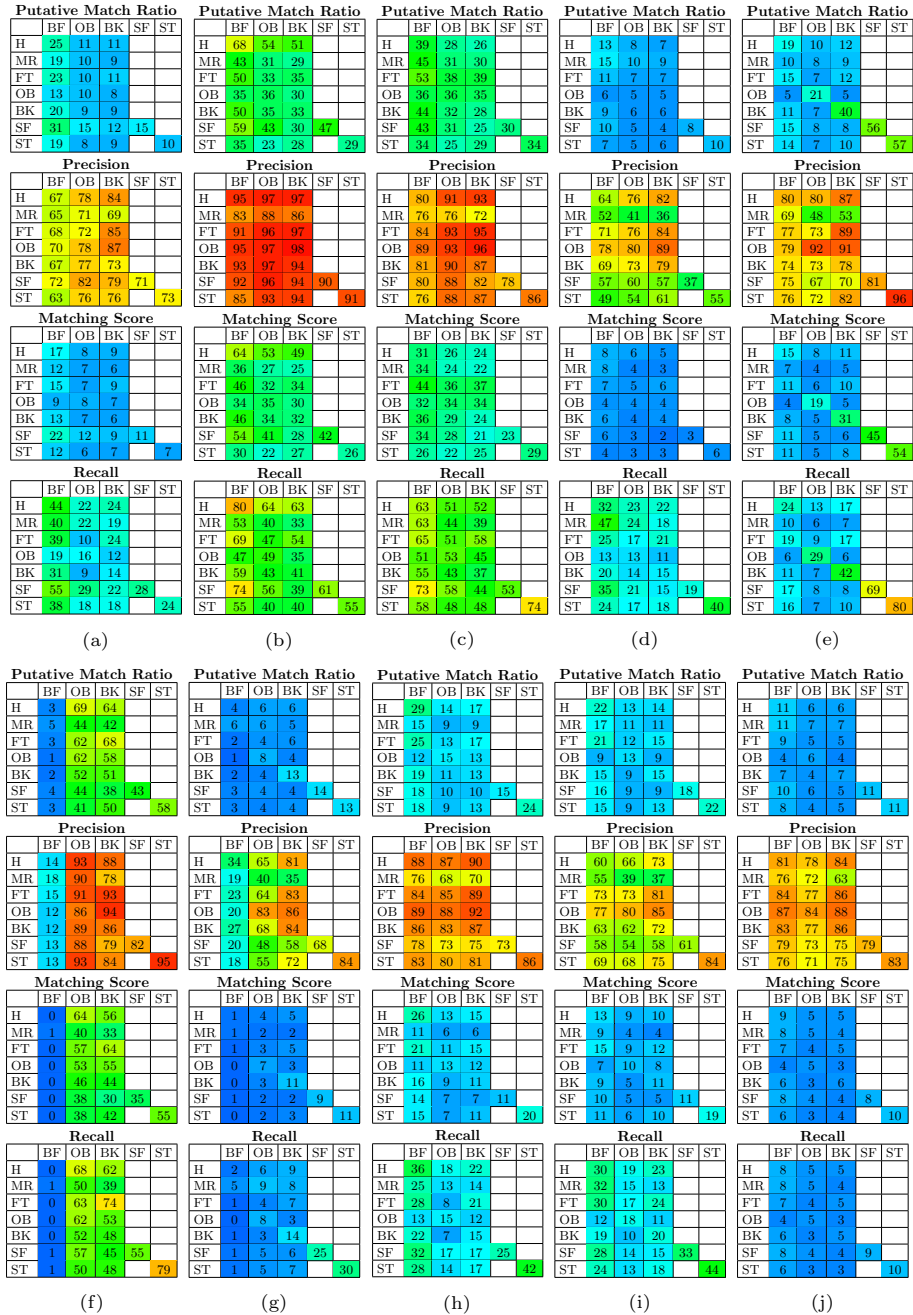


Fig. 5. Results for (a) blur, (b) JPEG compression, (c) exposure, (d) day-to-night illumination, (e) scale, (f) rotation, (g) scale and rotation, (h) perspective with planar scene, (i) perspective with non-planar scene, (j) photo-collection. Rows are detectors, columns are descriptors: H=Harris, MR=MSER, FT=FAST, BF=BRIEF, OB=ORB, BK=BRISK, SF= SURF, ST=SIFT. Values are in percent (blue 0%, red 100%).