# Large scale near-duplicate image retrieval using Triples of Adjacent Ranked Features (TARF) with embedded geometric information

Sergei Fedorov[*1] and Olga Kacher[†1]

[1]ABBYY, Moscow, Russia

October 25, 2018

**Abstract**

Most approaches to large-scale image retrieval are based on the construction of the inverted index of local image descriptors or visual words. A search in such an index usually results in a large number of candidates. This list of candidates is then re-ranked with the help of a geometric verification, using a RANSAC algorithm, for example. In this paper we propose a feature representation, which is built as a combination of three local descriptors. It allows one to significantly decrease the number of false matches and to shorten the list of candidates after the initial search in the inverted index. This combination of local descriptors is both reproducible and highly discriminative, and thus can be efficiently used for large-scale near-duplicate image retrieval.

## 1   Introduction

Most systems for large-scale image retrieval are designed for a scenario in which a query image is searched for within a large image set (hundreds of thousands or millions of images), and the goal is to find its partial duplicates. To solve this task most state-of-the-art algorithms make use of the Bag of Words (BoW) representation of an image, where each local descriptor extracted from an image is represented with a visual word, and implement an effective search in the inverted index (see for example [1]). One of the problems of all such algorithms is the large list of initial candidates, which needs to be further pruned with geometrical checks, usually with a RANSAC algorithm. For large image sets the candidate list can be very long, which

---

[*]sfedorov@abbyy.com
[†]okacher@abbyy.com

1

makes it infeasible to perform geometric verification for each candidate. In order to increase the discriminative power of the features extracted from an image, and to decrease the size of the initial candidates list, different ways of embedding spatial information into the features were suggested. Among those are bundling features [2], bundle min-hashing [3], nested SIFT [4], geometry-preserving visual phrases [5], and weak geometrical consistency [6]. These approaches allow one to significantly increase the effectiveness of the duplicates search in a large image set.

For the closely related task of logo recognition there have also been proposed methods for adding geometric information to the inverted index. In [7] triples of SIFT features were used for logo detection, and in [8] authors proposed multi-scale Delaunay triangulation for grouping local features into triples. In both papers [7, 8] a set of training images with the same logo is required for the construction of a logo model, a key point by which these approaches differ from ours, in which a single image is used for extracting triples of local features and indexing.

In this paper we consider the task of finding duplicates between two large image sets, both with hundreds of thousands or millions of images. One possible approach is to build a search index using the first image set, and then perform the search within it for each image from the second image set using one of the methods discussed above. This approach requires a large number of search queries, and is absolutely infeasible when the search index does not fit in the main memory. The problem here is the large list of initial candidates. For example, nested SIFT would first generate the list of all images with matched bounding features, and only then this list will be pruned by the requirement that there should be similar member features. Finding duplicates between two large image sets is in fact equivalent to considering all possible pairs of images, thus the initial list of candidates would be vast.

To address this problem we propose a composite local feature, which we have named TARF. A TARF descriptor is a combination of three local descriptors, and because of this it is highly discriminative, thus it allows us to significantly decrease the number of false matches and to shorten the list of candidates after the initial search in an inverted index. In order to increase the reproducibility of TARFs we propose a method of their extraction which takes into account the scores of constituent feature points, described in detail in the next section. Adding geometric information to this feature is straightforward and further increases its discriminative power, allowing us to achieve a very low false positives rate. This makes it possible to use TARF descriptors for the task of finding duplicates between two image sets even for very large image sets that don't fit in the main memory.

# 2 TARF features

In this paper we propose a new composite local feature, Triple of Adjacent Ranked Features (TARF). These triples of feature points are used to construct complex descriptors, which include three local descriptors, as well as the geometric layout of the three points. It is essential that these complex descriptors should be a) highly distinctive, and b) reproducible. To achieve these aims, we propose the following scheme for extracting TARFs from an image.

## 2.1 Extracting TARFs

Triples contain heterogeneous feature points. One feature point is detected with a blob detector, such as SIFT [9] or SURF [10]. Two other points are found with a corner detector, such as the BRISK detector [11]. We will denote these points as BLOB and CORNER. A group of three points gives a richer description of the local image area, as compared to a single feature point, and includes very distinctive geometric features. We first extract BLOB feature points from an image, and then each BLOB feature point is used to select nearby CORNER feature points. Enumerating all such triples would lead to a large number of combinations, most of which would not be very reproducible, because even if one of the three points is missing on a duplicate image, the whole triple will be lost. So, it is essential to select only those triples that have a good chance of being reproduced on a duplicate. To achieve this, we make use of the score of CORNER points (given by the CORNER points detector), modified to take into account position of the CORNER points relative to a BLOB feature point.

More precisely, the modified score is calculated as

$$
\begin{cases}
S^* = \exp\left[-0.5\left(\frac{d-d_0}{\sigma_d}\right)^2\right]\exp\left[-0.5\left(\frac{R_c-R_0}{\sigma_R}\right)^2\right]S \\
d_0 = 0.5R_b, \quad \sigma_d = 0.15R_b \\
R_0 = 0.33R_b, \quad \sigma_R = 0.15R_b
\end{cases}
, \qquad (1)
$$

where $R_c$ and $R_b$ are CORNER and BLOB feature points radii correspondingly, $d$ – distance between centers of CORNER and BLOB, $S$ and $S^*$ are original and modified CORNER scores. Parameters $d_0$, $\sigma_d$, $R_0$ and $\sigma_R$ were determined experimentally; the values in (1) correspond to TARFs based on SIFT/BRISK. These parameters play an important role for the quality of the TARF extractor. All CORNER points which lie in the vicinity of a BLOB feature point are sorted by modified score $S^*$, and $N_{top}$ CORNER feature points are taken. The value $N_{top} = 7$ was found to give best results. Next a global score threshold $S_0^*$ is determined in such a way that taking all CORNER points with modified score $S^* < S_0^*$ gives a total of $N_0$ triples:
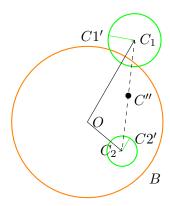
Figure 1: TARF descriptor. The descriptor is composed of BLOB point $B$ with the center at $O$ and radius $R_b$ and two CORNER points, $C_1$ and $C_2$. Segments $C_1C_1'$ and $C_2C_2'$ are related to the directions of CORNER feature points. $C''$ is the middle point of the line segment $C_1C_2$.

$$N_0 = \sum_{i \in \text{all BLOB points}} n_i(n_i - 1)/2, \qquad (2)$$

where $n_i$ is the number of CORNER points in the vicinity of $i$-th BLOB, with modified score below threshold, $S^* < S_0^*$. $N_0$ is a parameter of the detector, usually $N_0 = 2500..3000$

The exact algorithm for extracting TARF features:

---

**Algorithm 1** Extracting TARFs

---

1: Extract BLOB feature points with BLOB detector
2: Extract CORNER feature points with low threshold (for BRISK detector threshold 5 was used)
3: **for** each BLOB point **do**
4:     Choose all CORNER points from the neighbourhood of BLOB feature point
5:     Re-rank CORNER points according to (1)
6:     Choose top $n_*$ CORNER points ($n_* = 7$)
7: **end for**
8: Find global threshold for the scores of CORNER points. For each BLOB select $n_i$ top CORNER feature points in its neighbourhood with the score above global threshold, $n_i \le n_*$
9: Add all combinations (BLOB point + 2 different CORNER points from the list, determined in step 8) to the list of extracted TARFs
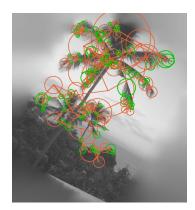
---

Figure 2: Matched TARF features on a partial duplicate images with 30 degrees rotation

## 2.2 Matching TARFs

Matching TARF features is straightforward: descriptors of all three constituent points should be matched, and the geometrical layout of three feature points should be similar. Geometrical features describing the TARF layout may be chosen in different ways. Our choice of geometrical features is given below ($R_b$ is the radius of the BLOB feature point):

1. Angle between vectors $OC_1$ and $OC_2$, $\alpha = \angle(C_1OC_2)$

2. Angle between vectors $OC_1$ and $C_1C_1'$, $\beta_1 = \angle(OC_1C_1')$

3. Angle between vectors $OC_2$ and $C_2C_2'$, $\beta_2 = \angle(OC_2C_2')$

4. Ratio $\epsilon_1 = OC_1/OC_2$

5. Ratio $\epsilon_2 = C_1C_2/R_b$

6. Ratio $\epsilon_3 = OC''/R_b$

Fig. 1 shows the geometric layout of feature points in TARF. An example of partial duplicates with matched TARF descriptors is given in Fig. 2.

## 2.3 Matching TARFs in the inverted index

Two TARFs are fully matched if all the three constituent descriptors are matched and their geometric layout is similar. Geometric layout is checked by verifying that all the geometric features listed above differ by less than a threshold:

$$
\begin{cases}
|\alpha^{(1)} - \alpha^{(2)}| < t_1, \\
|\beta_k^{(1)} - \beta_k^{(2)}| < t_1, & k = 1, 2 \\
|\epsilon_k^{(1)} - \epsilon_k^{(2)}| < t_\epsilon, & k = 1, 2, 3
\end{cases}
\tag{3}
$$

Thresholds were determined experimentally, $t_1 = 20°$, $t_\epsilon = 0.1$.

### 2.3.1 Matching descriptors with visual words

For large scale image retrieval it's essential to be able to match three constituent descriptors of TARFs simultaneously. The descriptors are represented in the inverted index with visual words. Three visual words, assigned to three descriptors, are packed into one 32-bit integer. Constituent descriptors of TARFs are matched if these 32-bit integers are exactly equal. Moreover, each visual word has an inverse document frequency (IDF) score associated with it. Several visual words with the lowest score (i.e. belonging to the largest clusters of descriptors) are placed into the stop list, and TARFs with such constituent descriptors are ignored.

One of the problems is that such a representation with visual words makes the probability of matching two similar descriptors significantly lower, as compared to matching with e.g. LSH (locality sensitive hashing), even for small vocabulary sizes. To illustrate this, Table 1 shows experimental results for probabilities of matching descriptors by various methods. For a set of 500 images we generated transformed duplicates; transformations included rotations and downscale. Feature points were extracted from original images and from duplicates. Those feature points that are at the same position and scale (up to image transformation) are considered to be reproduced. The reproducibility rate $R$ is the percentage of feature points that were extracted from the original image and were reproduced on the transformed image. Depending on the image transformations this rate is $R = 30..40\%$ for both SIFT and BRISK detectors. Next, pairs of descriptors were collected for the reproduced feature points, taking one descriptor from the original image, and another from the transformed image. These pairs of descriptors were matched using various methods, and the percentage of matched pairs is given in the TPR column. Thus, the overall probability of finding a match for a particular feature point from the original image on a transformed image is given by $p = R \cdot \text{TPR}$. The probability of false matches is given in the FPR column, i.e FPR is the match rate for pairs of descriptors calculated at random uncorrelated points of images.

As seen from Table 1, matching with visual words results in a TPR that is lower compared to matching by threshold, even for small vocabulary sizes. Though usually this TPR is sufficient, sometimes it may be desirable to increase the recall. This can be achieved by training two different vocabularies for each type of descriptor, and representing each descriptor with two visual words (different vocabularies are trained using clusterization of descriptors extracted from different sets of images). In this case descriptors are matched if any of the two visual words is matched. This makes TPR significantly higher and comparable with matching by a threshold, though FPR is of course much higher (see Table 1).

If this approach with several vocabularies of visual words is taken, then each TARF descriptor will be represented by $2 \times 2 \times 2 = 8$ 32-bit integers.

| matching method | FPR | TPR |
|---|---|---|
| threshold, L2 distance 0.4 | $3 \cdot 10^{-4}$ | 0.75 |
| visual words, 128 words | $8.5 \cdot 10^{-3}$ | 0.65 |
| visual words, 256 words | $4.3 \cdot 10^{-3}$ | 0.61 |
| visual words, 512 words | $2.2 \cdot 10^{-3}$ | 0.57 |
| visual words, 2 groups of 128 words | $1.3 \cdot 10^{-2}$ | 0.76 |
| visual words, 2 groups of 256 words | $6.6 \cdot 10^{-3}$ | 0.73 |
| visual words, 2 groups of 512 words | $3.4 \cdot 10^{-3}$ | 0.7 |

SIFT

| matching method | FPR | TPR |
|---|---|---|
| threshold, hamming distance 90 | $7 \cdot 10^{-4}$ | 0.57 |
| visual words, 128 words | $10^{-2}$ | 0.43 |
| visual words, 256 words | $5.2 \cdot 10^{-3}$ | 0.38 |
| visual words, 512 words | $2.8 \cdot 10^{-3}$ | 0.35 |
| visual words, 2 groups of 128 words | $1.9 \cdot 10^{-2}$ | 0.6 |
| visual words, 2 groups of 256 words | $9.6 \cdot 10^{-3}$ | 0.54 |
| visual words, 2 groups of 512 words | $5.1 \cdot 10^{-3}$ | 0.5 |

BRISK

Table 1: Descriptor match rates.

This requires 8 tables in the inverted index. Descriptors from two TARFs are considered to be matched if any of the eight 32-bit integers matches exactly. This can be viewed as a modification of LSH, with different hashes being generated by different vocabularies of visual words. Other TARF representations may also be considered: a) a BLOB descriptor is represented with 1 visual word, and a CORNER descriptor is represented with 2 visual words ($1 \times 2 \times 2$) or b) 2 visual words for BLOB, and 1 visual word for CORNER, $2 \times 1 \times 1$.

Equation (4) illustrates how a TARF descriptor can be represented with visual words when using two vocabularies for each type of constituent descriptor.

$$
\begin{cases}
\text{BLOB point descriptor} \rightarrow (w_1^b, w_2^b) \\
\text{1st CORNER point descriptor} \rightarrow (w_1^{c_1}, w_2^{c_1}) \\
\text{2nd CORNER point descriptor} \rightarrow (w_1^{c_2}, w_2^{c_2})
\end{cases}
\tag{4}
$$

$$
\begin{aligned}
&\text{TARF} \rightarrow (w_1^b, w_2^b) \times (w_1^{c_1}, w_2^{c_1}) \times (w_1^{c_2}, w_2^{c_2}) = \\
&(\{w_1^b, w_1^{c_1}, w_1^{c_2}\}, \{w_1^b, w_1^{c_1}, w_2^{c_2}\}, ..., \{w_2^b, w_2^{c_1}, w_2^{c_2}\})
\end{aligned}
$$

| matching method | FPR | TPR |
|---|---|---|
| threshold, geometric features discarded | $2.5 \cdot 10^{-5}$ | 0.7 |
| threshold | $4 \cdot 10^{-8}$ | 0.68 |
| visual words, 1x1x1, no stop list | $4 \cdot 10^{-9}$ | 0.28 |
| visual words, 1x1x1, stop lists of size 10 | $1.5 \cdot 10^{-9}$ | 0.25 |
| visual words, 2x1x1, stop lists of size 10 | $2.5 \cdot 10^{-9}$ | 0.27 |
| visual words, 1x2x2, stop lists of size 10 | $4 \cdot 10^{-9}$ | 0.41 |
| visual words, 2x2x2, stop lists of size 10 | $5 \cdot 10^{-9}$ | 0.44 |

Table 2: TARF match rates. First row corresponds to matching without checks of geometric layout. Other rows correspond to fully matched TARFs.

## 2.4 TARFs reproducibility and match rates

We have tested TARF features on the same set of 500 images with transformed duplicates as described in the previous section. TARFs that were extracted in the same places, i.e. in which all three constituent feature points were found at the same positions and scale (up to transformations), were considered to be reproduced. The reproducibility of TARFs is $R_T = 8\%$. Though this is a low number, it is still much higher than the reproducibility of randomly chosen triples of feature points, and this is due to the proposed scheme of TARFs extraction.

Table 2 shows the FPR and TPR for TARFs. The probability of false matches is given in the FPR column. TPR is defined as the probability that a pair of reproduced TARF descriptors is matched.

Normally TARFs are matched fully, including checks of geometric layout (3). In order to evaluate the importance of geometric features, we have tested TARFs matching without performing these checks, i.e. discarding geometric features and matching only constituent decriptors. Results are presented in the first row of the table. It can be seen that discarding geometric features leads to an FPR that is several orders of magnitude higher than the FPR for full TARF matching, yet it is still quite low at $2.5 \cdot 10^{-5}$. Including geometric features makes TARFs highly discriminative, with the FPR as low as $10^{-9}..10^{-8}$. Using stop lists for the visual words further reduces FPR, making it several times lower, and almost does not affect TPR. The TPR depends on the matching method; it has the largest value of 0.68 if the constituent descriptors are matched by threshold. In the case of matching constituent descriptors with visual words it is lower and varies from 0.25 to 0.44; the TPR is higher when several dictionaries of visual words are used.

## 3 Experimental results

The search of near-duplicate images is performed as follows. First, a search index is built. It is an inverted index of three visual words packed into 32-bit

| number of matched TARFs | probability | |
|:---:|:---:|:---:|
| | 1x1x1 scheme<br>no stop lists | 1x1x1 scheme<br>stop lists of size 10 |
| 0 | 0.985 | 0.995 |
| 1 | 0.013 | 0.004 |
| 2 | $10^{-3}$ | $10^{-4}$ |
| 3 | $3 \cdot 10^{-4}$ | $10^{-5}$ |

Table 3: Probability of a given number of matched TARFs on a pair of non-duplicate images

integers. Geometric data is stored in the index as auxiliary data. Three IDF scores of constituent visual words are summed up to give the IDF score of the TARF.

After finding all the images that have at least one TARF matching a query image, the score of the candidate images is calculated as the sum of IDF scores of all matched TARFs. A precision-recall curve can be obtained by varying the threshold for this score. Matched images can be filtered by applying a geometric model. We use a standard RANSAC algorithm. RANSAC converges very fast, because of the extremely low number of outliers, so we make only 10 iterations of RANSAC.

An important property of TARFs is a very low false positive rate. Table 3 shows the probability of finding a given number of matched TARF features for a pair of non-duplicate images. It can be seen that the probability decreases fast as the the number of false positive matches increases. It should be noted, that the use of stop lists for visual words is essential and reduces the number of false candidate images by the order of magnitude. If one builds a search index with 100K images using the 1x1x1 scheme with stop lists of size 10, and then searches for a query image in this index, the result of such a search will have on average as little as 400 false candidate images with one matched TARF descriptor, 10 false candidates with two matched TARFs and 1 false candidate with three matched TARFs. So, checking all the candidates with more sophisticated algorithms, e.g. with geometric verification by a RANSAC algorithm, is feasible even for very large search indices.

In order to evaluate image retrieval based on TARFs we have constructed an image set of 1853 images, and four sets of near-duplicates of these images, with the following modifications: a) downscale to small size 30k pixels (150x200), b) rotation 30 degrees, c) crop 70% (30 % of an image retained) and d) strong gaussian blur with $\sigma = 4$. Original images were added to the search index and diluted with the 100K distracting images, taken from the Flickr 100k Dataset [12]. For each of these modifications we have evaluated the precision-recall curve, and have measured AP (average precision), de-
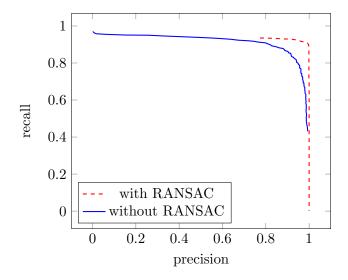
Figure 3: The effect of geometric verification with RANSAC on the precision-recall curve. Query images are duplicates downscaled to 30k pixels; original images are diluted with 100K distracting images.

fined as the area under the precision-recall curve. Figure 3 illustrates how geometric verification with RANSAC affects the precision-recall curve.

We have varied different parameters of the image retrieval system, and measured the resulting performance. The default parameters were: BLOB – SIFT descriptors, 1 visual word (dictionary with 256 words, stop list of size 10); CORNER – BRISK descriptors, 1 visual word (dictionary with 128 words, stop list of size 10); 3000 TARF descriptors per image; RANSAC verification, 10 iterations; 100K distracting images in the search index.

1. Number of distracting images

   Fig. 4 shows the dependence of AP on the number of distracting images in the search index. It is a key feature of TARF that the number of false positive candidates is very low, so even without geometric verification by RANSAC the precision remains high for a large number of distracting images.

2. Number of dictionaries of visual words

   Fig. 5 shows the dependence of AP on the number of dictionaries of visual words. 2x2x2 – two dictionaries for BLOB descriptors, and two dictionaries for CORNER descriptors, 8 tables in the inverted index; 1x2x2 – one dictionary for BLOB descriptors, and two dictionaries for CORNER descriptors, 4 tables in the inverted index; 2x1x1 – 2 tables in the inverted index; 1x1x1 – 1 table in the inverted index.
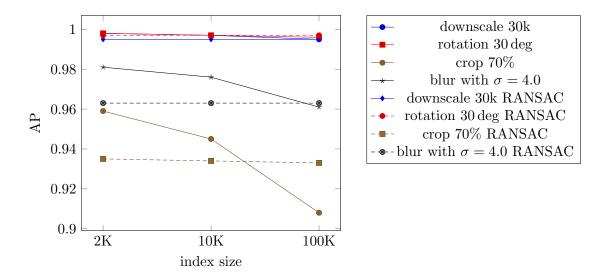
10

Figure 4: Average precision as a function of the size of the index (i.e. number of distracting images).

Best results are achieved for the 2x2x2 configuration, and the 1x2x2 configuration is nearly as good, thus one visual word is sufficient for BLOB point representation. However, using two visual words to represent each CORNER point is useful and increases AP.

3. Types of descriptors used

Fig. 6 illustrates the dependence of AP on the type of detectors/descriptors used to build TARF. For BLOB points we have tested SIFT, SURF, SURF detector + BRISK descriptor, and AKAZE [13]. For COR-NER points we have always used BRISK. With the only exception of AKAZE feature points, which appear to give low recall for crop 70%, all other configurations give similar results.

4. Number of TARF descriptors per image

Fig. 7 shows the dependence of AP on the number of TARF descriptors extracted from each image; the number of descriptors varies from 1000 to 4000. This graph justifies the choice of 3000 as a default number of TARF descriptors extracted from an image.

Table 4 shows experimental results for AP on popular **Holidays** and **Copydays** datasets [6], with 100K distracting images. The AP for **Holidays** and **Copydays strong** is low due of the low recall. Low recall on these datasets can be explained by the fact that TARF descriptors are matched only if the part of the image is reproduced almost exactly, up to scale/rotation. So, TARF descriptors are very sensitive to the viewpoint change, geometrical distortions of an image, etc. **Holidays** and **Copydays**
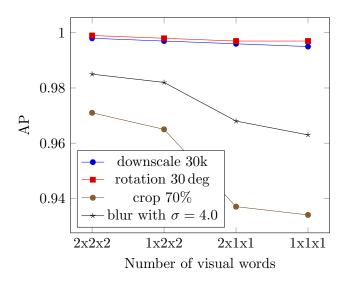
11

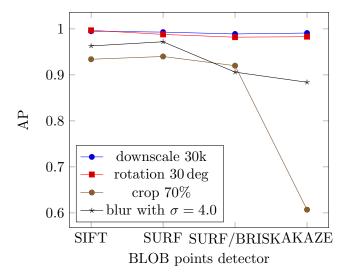Figure 5: Average precision as a function of the number of visual words.



Figure 6: Average precision vs types of detectors and descriptors of BLOB feature points.
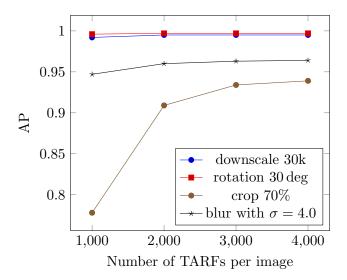
Figure 7: Average precision as a function of the number of TARFs extracted from each image.

| image set | average precision (AP) |
|---|---|
| Holidays | 0.278 |
| Copydays, strong | 0.466 |
| Copydays, crop 30 | 1.0 |
| Copydays, crop 70 | 0.916 |
| Copydays, crop 80 | 0.697 |
| Copydays, jpeg 50 | 1.0 |
| Copydays, jpeg 30 | 0.999 |
| Copydays, jpeg 10 | 0.98 |

Table 4: Average precision for commonly used datasets

**strong** image sets include many duplicates which don't have sufficiently large undistorted areas.

## 4  Conclusion

In this paper we have proposed a novel composite image feature – TARF, which is a combination of three local features. This feature is highly discriminative, because 1) TARF incorporates three local descriptors, thus the probability of a false positive match of a TARF is of the order of a third power of the probability of a false positive match for a single descriptor, which is a very low probability. In other words, it is very improbable that a random combination of three descriptors would match a given TARF. 2) TARF embeds the geometric relationships between three constituent local

features.

A method to effectively build an inverted index of TARF features is proposed. Each combination of three descriptors is represented by one 32-bit integer, and geometric information is added to the search index as auxiliary data.

It is shown that this approach makes it possible to build a system for large scale near-duplicates search, with high recall and a low false positives rate.

# References

[1] J. Sivic and A. Zisserman. Video google: a text retrieval approach to object matching in videos. In *Computer Vision, 2003. Proceedings. Ninth IEEE International Conference on*, pages 1470–1477 vol.2, Oct 2003.

[2] Zhong Wu, Qifa Ke, M. Isard, and Jian Sun. Bundling features for large scale partial-duplicate web image search. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 25–32, June 2009.

[3] Stefan Romberg and Rainer Lienhart. Bundle min-hashing for logo recognition. In *Proceedings of the 3rd ACM Conference on International Conference on Multimedia Retrieval*, ICMR '13, pages 113–120, New York, NY, USA, 2013. ACM.

[4] Pengfei Xu, Lei Zhang, Kuiyuan Yang, and Hongxun Yao. Nested-sift for efficient image matching and retrieval. *MultiMedia, IEEE*, 20(3):34–46, July 2013.

[5] Yimeng Zhang, Zhaoyin Jia, and Tsuhan Chen. Image retrieval with geometry-preserving visual phrases. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 809–816, June 2011.

[6] Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Hamming embedding and weak geometric consistency for large scale image search. In Andrew Zisserman David Forsyth, Philip Torr, editor, *European Conference on Computer Vision*, volume I of *LNCS*, pages 304–317. Springer, oct 2008.

[7] Stefan Romberg. *Aggregating Local Features into Bundles for High-Precision Object Retrieval*. PhD thesis, 2014.

[8] Yannis Kalantidis, Lluis Garcia Pueyo, Michele Trevisiol, Roelof van Zwol, and Yannis Avrithis. Scalable triangulation-based logo recognition. In *Proceedings of the 1st ACM International Conference on Multimedia Retrieval*, ICMR '11, pages 20:1–20:7, New York, NY, USA, 2011. ACM.

[9] David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60:91–110, 2004.

[10] Herbert Bay, Andreas Ess, Tinne Tuytelaars, and Luc Van Gool. Speeded-up robust features (surf). *Computer Vision and Image Understanding*, 110(3):346 – 359, 2008. Similarity Matching in Computer Vision and Multimedia.

[11] S. Leutenegger, M. Chli, and R.Y. Siegwart. Brisk: Binary robust invariant scalable keypoints. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2548–2555, Nov 2011.

[12] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman. Object retrieval with large vocabularies and fast spatial matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2007.

[13] P. F. Alcantarilla, J. Nuevo, and A. Bartoli. Fast explicit diffusion for accelerated features in nonlinear scale spaces. In *British Machine Vision Conf. (BMVC)*, 2013.