

Exploiting Local Indexing and Deep Feature Confidence Scores for Fast Image-to-Video Search

Savas Ozkan, Gozde Bozdagi Akar

Department of Electrical/Electronics Engineering, Middle East Technical University
06800, Ankara, Turkey

Abstract—The cost-effective visual representation and fast query-by-example search are two challenging goals that should be maintained for web-scale visual retrieval tasks on moderate hardware. This paper introduces a fast and robust method that ensures both of these goals by obtaining state-of-the-art performance for an image-to-video search scenario. Hence, we present critical enhancements to well-known indexing and visual representation techniques by promoting faster, better and moderate retrieval performance. We also boost the superiority of our method for some visual challenges by exploiting individual decisions of local and global descriptors at query time. For instance, local content descriptors represent copied/duplicated scenes with large geometric deformations such as scale, orientation and affine transformation. In contrast, the use of global content descriptors is more practical for near-duplicate and semantic searches. Experiments are conducted on a large-scale Stanford I2V dataset. The experimental results show that our method is useful in terms of complexity and query processing time for large-scale visual retrieval scenarios, even if local and global representations are used together. The proposed method is superior and achieves state-of-the-art performance based on the mean average precision (MAP) score of this dataset. Lastly, we report additional MAP scores after updating the ground annotations unveiled by retrieval results of the proposed method, and it shows that the actual performance.

Index Terms—Feature Indexing, Deep Features, Visual Content Retrieval, Image-to-Video Search

I. INTRODUCTION

In the last decades, we have witnessed an unprecedented proliferation of web-based multimedia data. The high-data population aggravates to retrieve a query from an extensive multimedia collection with a moderate hardware configuration.

Existing works [1]–[9] primarily focus on two goals: high performance and fast query. Visual representations computed in an offline step are not marked as one of the goals since it is done before the query is made. However, a fundamental question that needs to be addressed is that ‘*Is it enough to obtain the highest or fastest accuracy to deploy a complete retrieval system for users?*’. A plausible answer should be that, especially for larger multimedia databases (i.e., approaching real-world scenarios), solutions must consider hardware limitations before the querying to mitigate offline step complexity.

In this work, we present a visual multimedia retrieval method that aims to obtain high retrieval accuracy while keeping content presentation/indexing compact. Our main contributions are as follows:

- In this work, we use local and global visual features together. This step plays a crucial role by solving the

weaknesses of each feature set with the superiority of other methods. More specifically, local features tend to perform better with severe scale, rotation, and translation changes [1], [10]. Similarly, global features can provide better results for semantic tasks because part-based visual representations are utilized [5], [6], [11]. For this purpose, we combine the confidence scores of local and global features detected for the same scenes with a novel fusion technique. The direct fusion of confidence scores is not appropriate. Hence, our method aims to normalize the confidence scores of both local and global features first. It then merges them into a final ranking list by adaptively selecting a settling point for each query.

- In particular, the utilization of local and global descriptors for retrieval scenarios may conflict with low-cost computation limitations. Hence, we introduce improvements to well-known indexing and global pooling techniques; namely, Product Quantization (PQ) [12] and Compressed Fisher Vector (CFV) [13], to balance the workload and to promote modest visual description. In short, we propose a non-parametric weighting function to compute probabilistic similarity scores between local features for query and reference data in an asymmetric PQ space. Furthermore, we replace hand-crafted features [10] and sparse keypoints [1], [10] with densely sampled mid-level convolutional features. Notice that this step still has low computation workload, since deep features are densely estimated. In addition, semantic content can be represented precisely with deep features [5], [7], which differs from the local content (relevant to the goal of our work). Last but not least, we apply an approximate binary nearest neighboring search (NN) to make querying operation up to 6x faster for CFV with a minor decrease for accuracy.
- To this end, the proposed method enables fast query and low computational workloads for large-scale datasets while it outperforms the baselines by a large margin. Moreover, ground truth annotations for Stanford I2V dataset are updated that allow us to have a more reliable performance evaluation for future works.

II. RELATED WORK

Here, we review fundamentals and most recent studies related to visual content search tasks.

Local and Global Descriptors. With the advent of sparse local features [1], [10], this idea instantly becomes popular for



Fig. 1. Non-annotated scene samples are unveiled by our retrieval results on Stanford I2V dataset. Our method can retrieve queries with severe viewpoint and conditional changes at the top of the ranked list.

the visual representation domain. However, the usage of these features directly is not possible due to their dimensionality and large body. Capturing perfect descriptor relations by partitioning feature space into multiple clusters is a pioneering technique to ease this limitation [2]. However, to achieve discriminative hash codes, cluster size should be quite large even if an approximate search is utilized [3].

Product quantization [12], [14] reduces the cost of feature space partitioning significantly and mapping, since complexity is decreased by splitting a feature vector into multiple sub-vectors. Recently, the studies begin to focus on the adaptation of end-to-end learning techniques. Yi et al [4] propose a deep network pipeline for keypoint detection, orientation estimation and visual description in a unified manner. However, the main problem is that relatively high hardware configurations are needed to complete computations in an acceptable time. Recently, the studies concentrate on both computation and performance constraints [15].

On the other hand, orderless feature pooling (i.e., global descriptors) is an essential technique for retrieval tasks to produce a description from an independent set of features [11]. Similarly, large dimensionality aggravates the direct use of these descriptors. Although PCA-like methods can be used to transform descriptors to several principal axes, binary codes that are calculated with a simple threshold [13] provide several advantages as explained in [16]. Recently, neural networks are also used to compress feature vectors for effective search [17].

Similarly, Babenko et al. [7] compute deep features from different fragments of images and aggregate them with the VLAD descriptor. NetVlad [8] is a trainable generalization of VLAD to recognize the geo-location of images. Gordo et al. [9] improve performance on the visual landmark retrieval task by finetuning a deep model with a Siamese triplet loss. Lastly, [18] shows that local and global descriptors should be jointly used to eliminate possible outliers for accurate retrieval. However, note that these methods need labeled data to train/finetune high-level features for deep models, that is impractical for all datasets.

Low Offline Workload. As stated, a limited number of studies primarily consider low computational workload objectives in the literature [13], [19]. These solutions usually rely on representing a scene with a global visual descriptor. Also, the content around sparsely sampled points is used to achieve a reasonable computational workload. However, these descriptors tend to capture the content heavily from the background and repetitive parts [20]. Moreover, the discriminative power of representations (i.e., for binary codes) can decrease exponentially when the database size is increased [21]. Therefore, the single use of global descriptors remains weak for large-scale data, and additional representations should be exploited without sacrificing computational workload too much.

III. PROPOSED METHOD

Our main goal is to enhance the performance of visual retrieval tasks by exploiting the confidence scores of both local and global descriptors for the same scenes. Thus, the computational load should be low so that a large amount of data can be processed within a reasonable time.

Since our model is formulated as an image-based framework, keyframes are sampled at the start from a sequence of video frames $V = \{v_1, v_2, \dots, v_n\}$ with a simple heuristic rule. This rule enforces a uniform sampling strategy -1fps- and a constraint that each frame should not contain any sub-region with large motion variations. Otherwise, these frames are discarded.

A. Local Visual Content Representation

The local visual content of an input frame is computed from sparsely sampled key points. Then, these samples are converted into compact hash codes. In short, a two-stage quantization-based approach is utilized. Moreover, geometric consistency between local features is added with fast geometric filtering at the end.

Local Sparse Features. For local representation, we use Root SIFT [22] and Hessian Laplacian [1]. Since they are robust to scale and orientation changes, we expect to successfully

retrieve any query that presents strong geometric deformations (duplicate/copy) from the reference set.

In addition, to find geometric adjacency of local matches at the voting step, we store coordinate (x, y) , orientation θ and scale s coefficients of each local point in quantized forms.

Feature Indexing. As mentioned, the direct use of local features is impractical and should be converted into small representations. The idea of Bag-of-Word-like [2] methods is to quantize (Eq. 1) each feature vector $f_h \in \mathbb{R}^{128}$ (note that dimension of SIFT is 128 and k-means is used) to the closest center c_i from a pre-clustered feature space $C_{bow} \in \mathbb{R}^{D_{bow} \times 128}$.

$$q_b(f_h) = \min_i \|f_h - c_i\|_2, c_i \in C_{bow}, \quad (1)$$

However, the number of cluster centers should be adequately large (e.g. $D_{bow} > 100K$) in order to achieve discriminative space partitions [3]. This requirement leads to an increase in offline calculations and makes it difficult to represent visual content with local features.

In our work, a rational computation effort is achieved by encoding residual vector ($r = f_h - c_i$) with an additional quantizer $q_b(\cdot)$. Hence, relatively smaller cluster sizes (i.e., $D_{bow} \approx 10K$) can be selected. PQ space [12] $C_{pq} = \{C_{pq}^1, C_{pq}^2, \dots, C_{pq}^m\}$, $C_{pq}^k \in \mathbb{R}^{D_{pq} \times 128/m}$ is selected to maximize information bit per component by splitting residual vector into m non-overlapping sub-vectors $r = \{r_1, r_2, \dots, r_m\}$ as follows:

$$q_{pq}^k(r) = \min_i \|r_k - c_i\|_2, c_i \in C_{pq}^k, \forall k. \quad (2)$$

At the end, each feature vector is converted into two interrelated hash codes where $h_b = q_b$ and $h_{pq} = \{q_{pq}^1, q_{pq}^2, \dots, q_{pq}^m\}$, and they are stored in an inverted file structure based on their h_b values. Throughout this work, m and D_{pq} are empirically set to 8 and 256 as in [12]. In a nutshell, computation cost is exponentially reduced with low cost two-stage indexing scheme and local representation becomes applicable for large-scale databases.

Local Voting Scheme. We use a two-fold approach for the voting scheme: First, we estimate the best locally matching candidates based on similarities between hash codes. Later, outliers are eliminated based on the dominant geometric model between the query and reference frames.

Formally, in order to say that query and reference local points are similar, coarse hash codes should be the same ($h_b^q = h_b^r$), while residual similarities (h_{pq}^r and h_{pq}^q) must be within an error tolerance. Otherwise, the similarity score is set to zero and discarded from initial query candidates.

This work proposes a novel non-parametric score function for PQ Euclidean space to be used in the residual similarity calculation. Ultimately, this score function normalizes asymmetric Euclidean distance of two residual hash codes with a maximum asymmetric distance between all cluster centers.

Later, the final residual similarity score is equal to the average of all subvector scores:

$$w_{pq}(h_{pq}^r, h_{pq}^q) = \frac{1}{m} \sum_{k=1}^m \left(1 - \frac{1}{d_k} \|q_{pq}^{r,k} - q_{pq}^{q,k}\|_2 \right). \quad (3)$$

where d_k indicates the maximum asymmetric Euclidean distance for k^{th} subvector, i.e., $d_k = \max \|c_i - c_l\|_2, \forall i, l \in D_{pq}, c_i, c_l \in C_{pq}^k$. This non-parametric function enables us to assess similarities within $[0,1]$ ($w_{pq}(\cdot, \cdot) \in [0,1]$) rather than varying Euclidean distances. To this end, a threshold-based coefficient is used to select best matches.

After initial similarity scores are obtained by using Eq. (3), hard-similarity scores (these scores can be equal to either 0 or w_{pq}) are determined by applying a coarse threshold τ_{pq} . Practically, this step improves our method by selecting the best matches that yield high confidence scores, and it removes possible outliers immediately. Moreover, we prune 5% of the codewords according to their term frequencies to reduce the drawback of stop words [2] and to speed up the querying. Note that hard-similarity scores ($w_{pq}(\cdot, \cdot) > \tau_{pq}$ and $h_b^q = h_b^r$) are also weighted by this frequency term.

Later, we filter out outliers (i.e., the ones that not obey to dominant geometric model between the query and reference frames) by enforcing a 4-dof geometric constraint (affine transform might yield better results, yet it increases query time) (4) on initial query candidates:

$$\begin{pmatrix} x^q \\ y^q \\ 1 \end{pmatrix} = \begin{bmatrix} \tilde{s} \cos \tilde{\theta} & -\tilde{s} \sin \tilde{\theta} & t_x \\ \tilde{s} \sin \tilde{\theta} & \tilde{s} \cos \tilde{\theta} & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{pmatrix} x^r \\ y^r \\ 1 \end{pmatrix}. \quad (4)$$

where $\tilde{s} = s^q - s^r$ and $\tilde{\theta} = \theta^q - \theta^r$ are the differences of scale and orientation parameters for query and reference points. (t_x, t_y) is also spatial translation between query and reference local points, and their values are estimated from Eq. 4. This geometric model simply calculates a histogram using common parameter distributions of scale ($\log(\tilde{s})$), orientation ($\tilde{\theta}$) and translation ($(t_x + t_y)/\tilde{s}$) between query candidates of each frame. Later, the highest scored bin (sum of all hard-similarity scores that fit the dominant parameter distribution) yields the dominant geometric model between local points. To this end, the dominant value is set as a final local confidence score for a frame.

B. Global Visual Content Representation

Pretrained deep convolutional features are densely sampled and reduced by PCA to lower dimensions. Then, these features are aggregated with Fisher Kernel [11] and transformed into binary hash codes. These binary codes are compared with an approximate NN search setting in Hamming space to ease the querying step.

Dense Deep Convolution Features. We use densely sampled pre-trained deep convolutional features $f_d \in \mathbb{R}^{384}$ obtained at Alexnet-conv3 layer [5] by discarding zero paddings in convolution layers. As proved [5]–[7], deep features depict the semantic content of a scene more precisely than hand-crafted

features. By replacing hand-crafted features, we contain this semantic model with local structures to realize a complete retrieval method. Furthermore, we select conv3 features to keep the computational complexity low.

Later, these densely sampled features are mapped to 64-dimensional space by PCA. There are two main reasons: First, [7] shows that PCA-compressed deep features are robust since degrading the sparsity of features on a different visual set can improve their generalization capacity for unseen examples. Second, it provides time advantages in computations with feature pooling and voting stages.

Feature Pooling and Fast Voting Scheme. Deep features are aggregated with first-order Fisher Kernel [11] to estimate one compact representation $v \in \mathbb{R}^{64 \times D_{fk}}$ for each frame (D_{fk} is the number of Gaussian mixture components). Since dimensionality does not allow us to search and store them in large-scale databases, they are converted into binary codes b by applying a zero-bias threshold rather than quantization-based approaches. The main reason to select a threshold-based approach is that [16] shows that Euclidean-based quantization can be misleading for high dimensional representations. We also prove this assumption in our experiments.

In addition, even if converting a high-dimensional descriptor into a compact binary code speeds up the query time, there is still room to further eliminate the redundant calculations for our method. We replace the standard brute-force binary search with an approximate nearest neighboring (NN) scheme in this work.

In the proposed method, initial matched candidates for global binary codes are obtained based on KNN results in Hamming space (i.e., using inverted index structure the same as in local descriptors). Later, full distances (5) are computed by only comparing these candidates to improve confidence scores of query and reference frames.

$$w_b(b^r, b^q) = \begin{cases} g_h(b^r, b^q), & \text{if } b^r \text{ is in KNN of } b^q \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

where $g_h(\cdot, \cdot)$ is the score function, which returns normalized Hamming distance for two binary codes (normalized by the total number of bits, and similarly probabilistic scores are generated). Moreover, binary space is partitioned into 32 cluster centers throughout experiments.

C. Late Fusion

Until now, we compute a set of visual descriptors and construct two individual databases for local and global representations by depicting the same visual content. The most similar video scenes are retrieved separately by using these databases. Hence, we obtain two ranked lists, as illustrated in Fig. 2 for each query. From these lists, our objective is to obtain one final ranked list by fusing these confidence scores.

However, the fusion of these decisions is not straightforward. Even if confidence scores are in a similar range $[0, 1]$, there is no common score characteristic that can be directly exploited for all queries. Hence, confidence scores should be

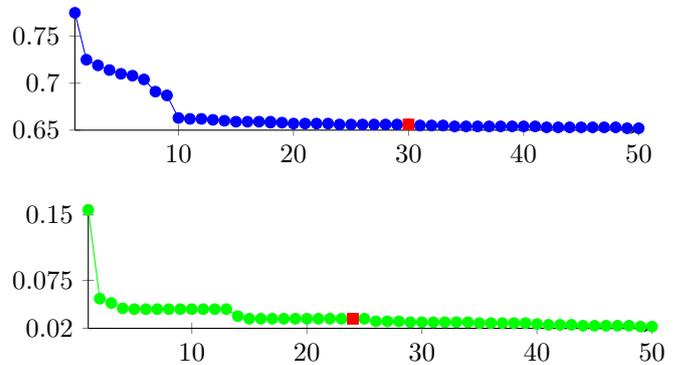


Fig. 2. Top confidence scores for two ranked lists obtained by global (blue) and local (green) descriptors. Red points indicate the adaptive settling points determined from each list.

normalized separately for each query before mapping these values to a final ranked list.

From these plots in Fig. 2, saturation of confidence scores for ranked lists shares similar characteristic (i.e., $0.75 - 0.65 \simeq 0.1$ and $0.15 - 0.02 \simeq 0.13$). Hence, this information can be exploited to fuse the scores. Therefore, a settling point must be initially determined from each list in order to normalize the scores. It can be done in two ways: 1) the last element of a list can be selected as a settling point, and scores in each list can be normalized by normalizing according to this value. 2) an adaptive point is determined from each list by using inner relations of confidence scores. Indeed, the second assumption yields better results, since it does not depend on the number of elements in ranked lists. Similarly, as in the assumption of query expansion [23], scores reflect an error characteristic after some points, and no content correlation is expected between the query and reference frames. Hence, our technique is inspired by the assumption of query expansion.

For this purpose, we iteratively calculate first-order score derivatives between all two consecutive confidence scores (i.e., subtracting one point from another). We then obtain an adaptively selectable point where gradient converges to a minimal number ϵ (e.g., $\epsilon=0.01$) after a period (after 10 elements). This point is accepted as a settling point, and all scores are normalized by subtracting this value.

Later, normalized local and global scores are merged by regarding their highest scores. Hence, the final ranked list is able to preserve both local structure and semantic similarities for each query in a unified form.

IV. EXPERIMENTS

Our experiments are conducted on Stanford I2V [19]. This dataset is particularly suitable for our method since it contains a large volume of videos collected from diverse news video archives to illustrate its actual capacity.

As stated in [19], Stanford I2V dataset is split into two versions, such that a lighter version contains a subset of query images and reference videos of a full version. Full and light versions consist of 3801 and 1035 hours of videos,

TABLE I
APPROXIMATE TIME (SEC) SPENT ON THE REPRESENTATION STAGE PER
FRAME ON A SINGLE CPU CORE.

Local Descriptor			
Keypoint	Descriptor	Indexing	Total
0.223	0.410	1.331	1.996
Global Descriptor			
Descriptor	PCA	Fisher Kernel	Total
1.163	0.005	0.193	1.187

respectively (Please do not confuse SI2V-4M or SI2V-600K in [21]). Moreover, the amount of query images is decreased by factor 3, from 229 to 78.

The provided script measures the performance in the evaluation step, and the mean Average Precision (MAP) scores are reported.

Computation Load. The main objective is to accomplish a method that computation load and query processing time are moderate while obtaining state-of-the-art retrieval accuracy for large-scale data. Table I illustrates average time requirements per frame. Observe that it is close to real-time (assume that 1fps keyframe is processed). This feature allows us to analyze the visual content of this dataset within several days on cheap CPU servers. In similar, memory requirement is negligible since all descriptors are converted to compact hash codes.

Impact of the Parameters. Our framework is composed of local and global ranking stages. Hence, first, we need to obtain the best parameters empirically for each stage.

For local visual representations, τ_{pq} defines error tolerance for PQ signature matches. For small values, noisy versions of valid signatures can be estimated correctly. However, this reduces the discriminative power of signatures and introduces lots of outliers. Therefore, we set D_{bow} as 5K and 10K to obtain trade-offs in terms of retrieval accuracy and query processing time for τ_{pq} . Fig. 4 illustrates mAP scores for

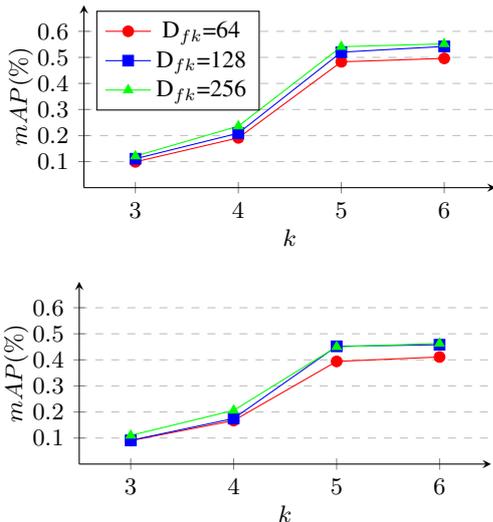


Fig. 3. Impact of k values for binary NN voting. Results for light (upper) and full (lower) sets are reported.

top 100 ranked scenes based on various τ_{pq} values. From the results, settling τ_{pq} around 0.72 yields the best performance for all configurations. Another important observation is that 5K scores drop down drastically for the full version of this dataset. This result can be demonstrated by the fact that the discriminative power of this representation is not adequate for smaller values of D_{bow} . Moreover, a larger cluster size D_{bow} ultimately provides an advantage in the querying stage by reducing operations due to an inverted index structure.

For global visual representations, D_{fk} and k are two parameters that users need to define. We select 64, 128 and 256 for the number of Gaussian mixture components (D_{fk}). Since k value determines redundancy in approximate binary code search, accuracy and query processing time are influenced inversely from this value. The accuracy saturates at $k = 5$ for all configurations, as shown in Fig. 3. This configuration speeds up search time approximately 6x faster. As expected, increasing the number of mixture components (D_{fk}) restores accuracy for both versions. However, the total number of comparisons, as well as storage requirements, becomes more extensive.

Impact of Late Fusion. We calculate mAP scores after fusing confidence scores of local and global descriptors in various combinations. Table II shows that late fusion boosts retrieval accuracy around 5% compared to their individual baselines estimated by local and global representations (i.e., Fig. 3 and Fig. 4). Also, the combination of 5K-256 obtains the best mAP accuracy for the light set. However, scores are decreased for the full version due to the failure of PQ. On the other hand, 10K-256 combination yields compatible scores for both full and light versions.

In the querying stage, the combination of 10K-64 yields the fastest response time. Since global representations are also stored in an inverted index structure, the use of high-dimensional representation increases the sparsity of code-words. The number of computations is reduced compared to 5K representation.

Baselines. We compare our performance with the reported baseline results in the literature (Table II). Initial database performance [19] on light and full versions are approximately 46% and 43% for mAP@100. Later, even if worse perfor-

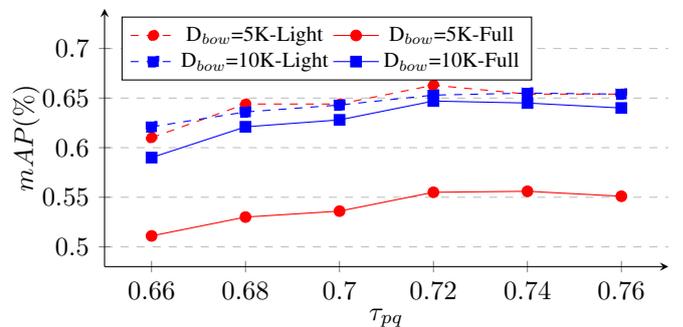


Fig. 4. Impact of τ_{pq} for top 100 retrieved scenes. Results are reported for local descriptors.

TABLE II
LATE FUSION RESULTS ON SI2V DATASET. LATENCY (SEC) IS MEASURED FOR 1000H REFERENCE VIDEOS PER QUERY. BEST RESULTS IN EACH PART ARE IN BOLD.

[D _{bow} -D _{fk}]	Light Dataset		Full Dataset		Latency Per 1000h
	mAP	mAP@1	mAP	mAP@1	
EH [24]	-	-	0.15	0.37	-
PHOG [24]	-	-	0.22	0.45	-
SCFV [19]	0.46	0.73	0.43	0.64	12.75 sec
BF-PI [21]	≈0.68	-	≈0.65	-	≈ 4.3 sec
RMAC [25]	-	-	≈0.66	-	-
ours[5K - 64]	0.667	0.769	0.582	0.716	17.11 sec
ours[5K - 128]	0.695	0.794	0.601	0.755	18.237 sec
ours[5K - 256]	0.707	0.782	0.622	0.755	19.253 sec
ours[10K - 64]	0.668	0.769	0.644	0.764	8.675 sec
ours[10K - 128]	0.679	0.782	0.663	0.786	9.802 sec
ours[10K - 256]	0.700	0.782	0.670	0.777	10.809 sec

mance is achieved, authors present simple yet somehow effective global representations [24]. More recently, [21] achieves an additional 21% mAP improvement compared to baseline scores by using a shot-based feature aggregation technique. RMAC based deep feature pooling is also adopted [25]. Lastly, the latency of [21] might be better than our method due to promoting frame-based assumption to shot-based assumption. However, remark that shot-based features introduce additional workloads for offline computations. Notice that our method obtains state-of-the-art performance on SI2V dataset.

Updating the Ground Truth Annotations. We update ground truth annotations for both full and light versions of SI2V dataset. These annotations are unveiled with our retrieval results (<https://github.com/savasozkan/i2v>).

Annotation pipeline of SI2V dataset [19] relies on an automated annotation process, as explained by the authors. Precisely, reference videos candidates are initially pruned with a time constraint (based on time tag of queries), and a feature-based matching technique is utilized before any human visual intervention. As a result, some of the scenes might be discarded unintentionally from annotation lists. Therefore, we manually examine our top retrieval results (up to 20 scenes per query) and find out that some of the retrieved results are non-annotated in the ground truth, even if they have strong semantic analogies and visual copies with query images. We illustrate some of the scene samples in Fig. 1.

Table III shows mAP scores recalculated with our updated ground truth annotations. The results introduce additional 5% improvements compared to Table II for light version of dataset. This result is profoundly critical since the actual performance of our method is even beyond the reported performance in the literature. Moreover, although there is a noticeable performance increase for the light set, this increase is not as much for the full version. The reason is that the selection capacity of global binary representation saturates for the larger set, as explained. This notion validates the importance of the joint use of local and global representations for retrieval tasks. To make fair comparisons, we also implement and test the baseline method proposed in [19], [24] on updated ground truth annotations. From the results, it provides only 2% and

TABLE III
LATE FUSION MAP AFTER UPDATING THE GROUND TRUTH ANNOTATIONS. BEST RESULTS IN EACH PART ARE IN BOLD.

[D _{bow} -D _{fk}]	Light Dataset		Full Dataset	
	mAP	mAP@1	mAP	mAP@1
[5K - 64]	0.697	0.794	0.577	0.720
[5K - 128]	0.735	0.833	0.607	0.755
[5K - 256]	0.755	0.846	0.624	0.764
[10K - 64]	0.708	0.807	0.648	0.768
[10K - 128]	0.729	0.820	0.667	0.786
[10K - 256]	0.755	0.833	0.681	0.790
EH [24]	-	-	0.19	0.42
SCFV [19]	0.48	0.76	0.44	0.68

1% additional improvement for [19] while 4% improvement is obtained for [24].

V. CONCLUSION

In this work, we introduce a visual search method for large-scale visual retrieval task. It exploits local and global descriptors together to represent visual data. The primary objective of the proposed method is to obtain moderate computation load and query time for large-scale datasets. Furthermore, performance is improved compared to baselines. We present critical contributions to the techniques for visual representation and feature hashing throughout the paper. In addition, we propose a novel technique to fuse local feature-based scores and deep global scores as a late fusion step. To show the superiority of our method, experiments are conducted on Stanford I2V dataset. As explained, it achieves the state-of-the-art mAP performance in the literature. Moreover, we update ground truth annotations for Stanford I2V based on the retrieval results of the proposed method. The final results show that the actual performance of our method is much better after updated ground truth annotations are used.

ACKNOWLEDGMENT

The authors gratefully acknowledge the support of NVIDIA Corporation with the donation of GPUs used for this research.

REFERENCES

- [1] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE transactions on pattern analysis and machine intelligence*, pp. 1615–1630, 2005.
- [2] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," *Proceedings of the IEEE international conference on computer vision*, pp. 1470–1477, 2003.
- [3] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8, 2007.
- [4] K.M. Yi, E. Trulls, V. Lepetit, and P. Fua, "Lift: Learned invariant feature transform," *European conference on computer vision*, pp. 467–483, 2016.
- [5] A. Krizhevsky, I. Sutskever, and G.E. Hinton, "Advances in neural information processing systems," *NIPS*, pp. 1097–1105, 2012.
- [6] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [7] A. Babenko and V. Lempitsky, "Aggregating local deep features for image retrieval," *Proceedings of the IEEE international conference on computer vision*, pp. 1269–1277, 2015.
- [8] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5297–5307, 2016.
- [9] A. Gordo, J. Almazán, J. Revaud, and D. Larlus, "Deep image retrieval: Learning global representations for image search," *European conference on computer vision*, pp. 241–257, 2016.
- [10] D.G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, pp. 91–110, 2004.
- [11] F. Perronnin, J. Sanchez, and T. Mensink, "Improving the fisher kernel for large-scale image classification," *European conference on computer vision*, pp. 143–156, 2010.
- [12] H. Jegou, M. Douze, and C. Schmid, "Product quantization for nearest neighbor search," *IEEE transactions on pattern analysis and machine intelligence*, pp. 117–128, 2011.
- [13] L.Y. Duan, J. Lin, J. Chen, T. Huang, and W. Gao, "Compact descriptors for visual search," *IEEE Multimedia*, pp. 30–40, 2014.
- [14] T. Ge, K. He, Q. Ke, and J. Sun, "Optimized product quantization," *IEEE transactions on pattern analysis and machine intelligence*, pp. 744–755, 2014.
- [15] Julieta Martinez, Shobhit Zakhmi, Holger H Hoos, and James J Little, "Lsq++: Lower running time and higher recall in multi-codebook quantization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 491–506.
- [16] H. Cevikalp, M. Elmas, and S. Ozkan, "Large-scale image retrieval using transductive support vector machines," *Computer Vision and Image Understanding*, 2017.
- [17] Stanislav Morozov and Artem Babenko, "Unsupervised neural quantization for compressed-domain similarity search," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 3036–3045.
- [18] Hyeonwoo Noh, Andre Araujo, Jack Sim, Tobias Weyand, and Bohyung Han, "Large-scale image retrieval with attentive deep local features," *Proceedings of the IEEE international conference on computer vision*, 2017.
- [19] A. Araujo, J. Chaves, D. Chen, R. Angst, and B. Girod, "Stanford I2V: a news video dataset for query-by-image experiments," *Proceedings of the 6th ACM Multimedia Systems Conference*, pp. 237–242, 2014.
- [20] A. Torii, J. Sivic, T. Pajdla, and M. Okutomi, "Visual place recognition with repetitive structures," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 883–890, 2013.
- [21] A. Araujo and B. Girod, "Large-scale video retrieval using image queries," *IEEE Transactions on Circuits and Systems for Video Technology*, 2017.
- [22] R. Arandjelovic and A. Zisserman, "Three things everyone should know to improve object retrieval," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2911–2918, 2012.
- [23] O. Chum, A. Mikulik, M. Perdoch, and J. Matas, "Total recall II: Query expansion revisited," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 889–896, 2011.
- [24] Gabriel de Oliveira Barra, Mathias Lux, and Xavier Giro-i Nieto, "Large scale content-based video retrieval with livre," in *2016 14th International Workshop on Content-Based Multimedia Indexing (CBMI)*. IEEE, 2016, pp. 1–4.
- [25] Noa Garcia and George Vogiatzis, "Asymmetric spatio-temporal embeddings for large-scale image-to-video retrieval," *BMVC*, 2018.