

UR2KiD: Unifying Retrieval, Keypoint Detection, and Keypoint Description without Local Correspondence Supervision

Tsun-Yi Yang^{*,1,2}¹Scape Technologies

tsun-yi@scape.io

Duy-Kien Nguyen^{*,1,3}²National Taiwan University

kien@vision.is.tohoku.ac.jp

Huub Heijnen¹

huub@scape.io

Vassileios Balntas¹³Tohoku University

vassileios@scape.io

Abstract

In this paper, we explore how three related tasks, namely keypoint detection, description, and image retrieval can be jointly tackled using a single unified framework, which is trained without the need of training data with point to point correspondences. By leveraging diverse information from sequential layers of a standard ResNet-based architecture, we are able to extract keypoints and descriptors that encode local information using generic techniques such as local activation norms, channel grouping and dropping, and self-distillation. Subsequently, global information for image retrieval is encoded in an end-to-end pipeline, based on pooling of the aforementioned local responses. In contrast to previous methods in local matching, our method does not depend on pointwise/pixelwise correspondences, and requires no such supervision at all i.e. no depth-maps from an SfM model nor manually created synthetic affine transformations. We illustrate that this simple and direct paradigm, is able to achieve very competitive results against the state-of-the-art methods in various challenging benchmark conditions such as viewpoint changes, scale changes, and day-night shifting localization.

1. Introduction

Image matching is one of the most important research topics in computer vision with several applications such as 3D reconstruction [48, 49, 51, 52], visual tracking [16, 61, 22] and SLAM [36, 37]. Several hand-crafted feature descriptors have been proposed [27, 64, 44, 1, 24, 7], and have been widely utilized in state-of-the-art systems. Similarly to other areas of computer vision, deep learning has recently influenced this area, with a plethora of works focusing on learning deep patch descriptors in a supervised manner [55, 56, 4, 62, 6, 32, 17, 35] exhibiting superior performance compared to the “classical” hand-crafted meth-

^{*}Equal contribution. This work was mostly done in Scape Technologies research internship period.



(a) Pre-trained ResNet101



(b) D2-Net [15]



(c) UR2KID (ours)

Figure 1. Extremely challenging image matching scenario with severe scale change and significant scene difference between day and night. The proposed UR2KID method is able to utilize a common network structure to achieve state-of-the-art results.

ods.

Recently, an important direction of research with significant impact has been combining the whole matching pipeline into a single end-to-end process [15, 39, 65, 14]. This enables the network to take advantage of several nuance factors that are related and contribute to matching. The detect-and-describe concept [15] jointly encodes the keypoints and the descriptors in a single feature map. However, a significant drawback of such methods is that the training process is strongly supervised either by pixel-level correspondences from dense Structure-from-Motion (SfM) models which are extremely costly to produce [25], or the application of manual synthetic affine transformations which unfortunately do not exhibit all the complex deformations seen in the real world [14].

In very large-scale applications such as city-scale localization, exhaustive brute force matching of all possible image pairs with local descriptors is an extremely costly and non-scalable process. An intuitive way to limit the candidate pool of images, thus making the problem tractable, is to utilise global image retrieval to limit the search space, and subsequently perform re-ranking using more accurate image matching and geometric verification methods [46, 47].

Nevertheless, previous studies focused on separate concepts for training or optimizing local matching pipeline and global descriptor retrieval as several independent parts. For instance, local matching methods are normally evaluated using local correspondences, and are frequently split into evaluation of keypoint robustness [23] and descriptor matching performance [5]. However, such local methods neglect global context information that is by definition encoded into global descriptors for the task of image retrieval. Our key observation is that the tasks of image retrieval and image matching are highly interconnected, and a suitable optimization process can tackle both. For example, parts of the image that are suitable for global description such as buildings and static structures, would also normally be suitable for local matching. On the other hand, people, trees and cars, are normally unsuitable for both problems. Despite their seemingly obvious relation, these tasks have been tackled either completely independently [55, 56, 23, 41], or with minimal interaction [38, 46].

In this paper, we propose to unify feature encoding in terms of both local and global information, using a multi-task learning approach. Unlike traditional patch based methods, our local matching feature maps are suitable for both keypoint detection and description. Multi-task information is embedded in a single network which is also responsible for learning a global descriptor suitable for image retrieval. We focus on the following contributions:

- We present a multi-task method for global retrieval and local matching embedded within a single network with a training process that does not rely on local pixel level correspondence ground truth.
- We introduce a novel method to aggregate feature map descriptors, namely group-concept detect-and-describe (GC-DAD). Using our method, we show that a standard deep network trained on ImageNet for a classification task can match or outperform state-of-the-art end-to-end trained matching methods.
- We show that by using cross dimensional self-distillation along with the proposed training method with no pixel-level correspondences, our method is able to acquire low dimensional local descriptor which is more robust against scale changes, viewpoint changes, and day-night shifting localization as shown in Figure.1. Our combined global and local descrip-

tor, is able to outperforms state-of-the-art localization methods.

2. Related Work

In this section, we briefly discuss the classical patch based matching pipeline, the state-of-the-art end-to-end image matching methods, and the relation between local and global matching.

Classical matching pipeline. The classical matching pipeline can be deconstructed in the following components: keypoint detection, description, and matching. **a) Keypoint detection:** For sparse matching, finding robust keypoints is the first crucial part. Several properties should be concurrently satisfied, such as scale & affine invariance, and keypoint repeatability [33, 31, 27, 34, 44, 34, 23, 43]. **b) Keypoint description:** Given a keypoint location and a corresponding scale associated with it, a patch can be rectified around it, and subsequently described either by hand-crafted method such as SIFT [27], SURF [9, 8], LIOP [58], ORB [44, 36, 37] or by deep learned patch description methods such as [55, 56, 4, 62, 6, 32, 17, 35].

End-to-end matching pipelines. In the classical pipeline, the main individual components (i.e. keypoint detection and description) are often tackled separately which might lead to sub-optimal results. Recent approaches try to tackle this issue by introducing end-to-end matching pipelines. For preserving the differentiability over the whole matching pipeline, the authors of LIFT [65] replace the keypoint detection and affine estimation by a spatial transformer network [20] and the non-maxima suppression by a soft argmax. By using a dense SfM model during training, viewpoint and lighting conditions are resolved by training a siamese network. Similarly, Lf-Net [39] and Superpoint [14] use synthetic affine transformations for generating image pairs as training data, which allows the network learning to be based on point-wise pixel level ground truth local correspondences. D2-Net [15] is also trained based on pointwise ground-truth from a set of dense SfM models [25], and its main contribution, was the proposed detect-and-describe process which entails using same feature map for keypoint detection and description. ELF [10] exploits the possibility of pre-trained network through back-propagated saliency map to get robust detector and matching results.

Retrieval methods. Where feature matching uses pixel level descriptions and correspondences, image retrieval uses a single (global) description for the whole image. [54, 3, 40, 42, 41, 2, 59, 60]. The similarity between two images can then be computed very efficiently with a single distance computation, something that makes the process useful for retrieving a pool of possible matching candidates in large-scale dataset, with respect to a given query image. Recent state-of-the-art methods aggregate feature maps into

Method	Training dataset	Ground truth label			Local		Global
		How to get	Img pair	Loc. corres	Detector	Descriptor	Descriptor
Keypoint detection							
QuadNet	DTU robot image dataset	3D	-	✓	✓	-	-
Key.Net	ImageNet ILSVRC 2012	Self	✓	-	✓	-	-
Descriptor learning							
HardNet	UBC/Brown dataset	MVS	-	✓	-	✓	-
SOSNet	UBC/Brown dataset	MVS	-	✓	-	✓	-
Matching pipeline							
LIFT	Piccadilly Circus dataset	SfM	✓	✓	✓	✓	-
Superpoint	MS-COCO	Self	✓	✓	✓	✓	-
LF-Net	ScanNet, 25 photo-tourism	SfM	✓	✓	✓	✓	-
D2-Net	MegaDepth dataset	SfM	✓	✓	✓	✓	-
R2D2	Oxford, Paris, Aachen (scene specific)	SfM, Flow, Style	✓	✓	✓	✓	-
ELF	ImageNet pre-trained	-	-	-	▲	▲	-
Retrieval							
NetVLAD	Google street view	T.M.	✓	-	-	-	✓
GeM, DAME	SfM-120k (cluster+3D)	SfM	✓	-	-	-	✓
Multi-task method							
DELf	Landmark dataset	Class	-	-	▲	▲	✓
HF-Net	Google landmark, BDD	Teacher	✓	✓	✓	✓	✓
ContextDesc	Photo-tourism, aerial dataset	SfM	✓	✓	✓	✓	▲
UR2KID	ImageNet pre-trained, MegaDepth	SfM	✓	-	✓	✓	✓

Table 1. Overall comparison of related methods. ▲ indicates that this specific task is not directly optimized by this method.

global descriptors using an ℓ_3 -norm in Generalized-Mean (GeM) [42] and more recently a learnable dynamic ℓ_p -norm in DAME [63]. Detect-to-retrieve [54] adopts an extra detection module for sub-region feature extraction before the final aggregation.

Multi-task methods. Considering both local correspondence and global description aim to induce a similarity metric between two images, joint learning local matching and global retrieval might be beneficial for both sides. For the purpose of landmark classification using global descriptors, local masks and image pyramids are adopted in DELF [38]. Here, even though local matching is explored as a side-product, it is not directly optimized for the task. HF-Net [46] adopts a teacher-student distillation for both matching and retrieval in order to achieve fast and robust localization. ContextDesc [28] encodes both visual and geometric context along with the information from off-the-shelf retrieval network trained on Google-landmarks [38] into the description.

In Table 1 we present a systematic categorization of related techniques with respect to training methodologies and method outcomes. Despite the fact that some of these methods can be used for both matching and retrieval (e.g., DELF, HF-Net, Contextdesc) there is no simple unified framework for learning the both tasks concurrently. In the following section we will introduce the proposed UR2KiD method

that aims to explore several problems related to learning a simple unified multi-task method that is suitable for both image matching and retrieval.

3. Methodology

Given an image I , the extraction of the visual information into a suitable mathematical representations also referred to as visual descriptors, can be formulated in two ways, either globally or locally. The global descriptors are extracted using the entire image, and are normally utilised in retrieval problems. The local descriptors, use information from small regions of the image, and are normally used for image matching scenarios. For each image I , the output of these processes, is a set of keypoints k_i , a set of descriptors d_i corresponding to index i , and a single global descriptor d_g .

Unlike previous methods that treat local and global descriptors separately, our goal is to generate both global and local description within a single network pipeline P_{joint}

$$d_g, \{k_i, d_i\} = P_{joint}(I). \quad (1)$$

The given training supervision is the image level pair $(a, p, n_1, n_2, \dots, n_k)$ where (a, p) is the positive pair and (a, n_k) is the negative pairs without the pixelwise or pointwise matching correspondences. The image pairs come

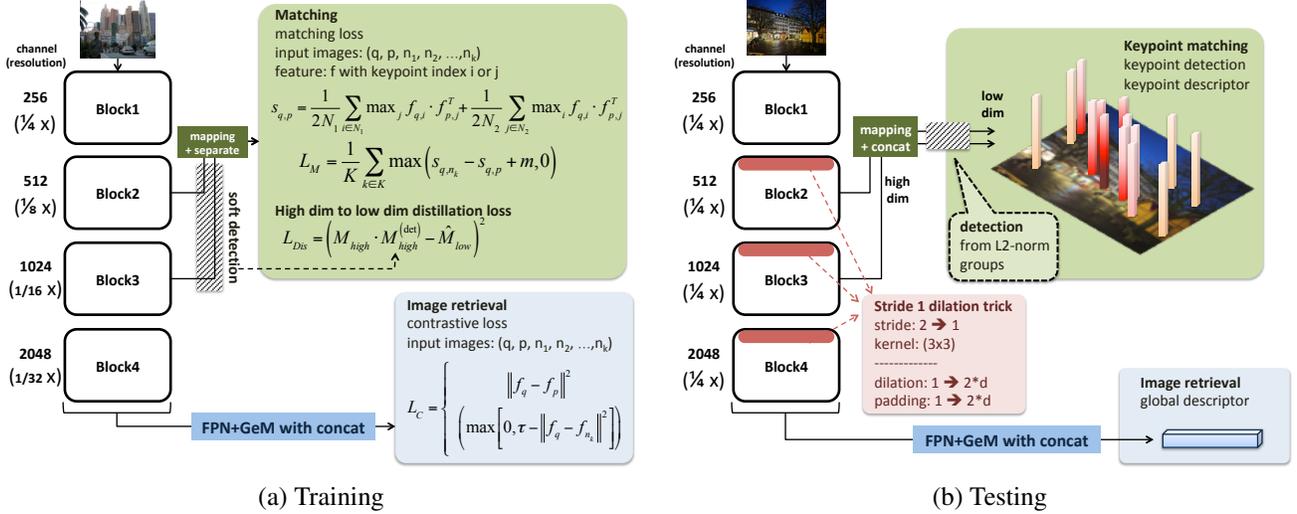


Figure 2. Overview of the proposed method.

from SfM model and the whole process does not depend on any human labeled ground truth.

We now briefly describe the overview of our method. In Figure. 2, the training and testing pipelines are demonstrated, and we present a general visual outline of our method using blocks to represent different parts of the system. As a simplification, we adopt the block design of the ResNet for the high-level network explanation (e.g. capital B as Block. B1 for Block1.). Our goal is to utilize the feature map hierarchy of a convolutional neural network, which has known to contain different levels of semantics, and produce multi-level representations of an image with discriminative features for image retrieval problem. We use a pretrained CNN (ResNet101 [18] trained on ImageNet [13, 45]) to extract multi-level visual features.

The architecture is composed of the following components: multi-level feature extraction for local descriptors from ResNet blocks and feature pyramid network (FPN) [26]. The final global descriptor is the concatenation of several different GeM [42] pooling results from different layers.

3.1. Local keypoint and description

For the local matching paradigm, we first start by formalising the the detect-and-describe (DAD) method introduced in D2-Net [15]. To determine the i -th keypoint coordinate (x_i, y_i) out of feature map F , the keypoint confidence of such location is computed by thresholding the maximum response of the feature map $F_{(x_i, y_i)}^{c^*}$ across channel dimension (i.e. $c^* = \arg \max_c F_{(x_i, y_i)}^c$) along with local non-maximum suppression and edge confidence thresholding of Harris corner detector, and the second-order spatial displacement first described in [27].

Our method is mainly inspired by D2-Net [15] and Sime-

oni et al [50] who propose a way to explore the activations of CNNs as keypoint detectors, and use the parameters of classical detectors such as Harris or MSER [30] on the activation to match two images channel by channel. However, our key difference with Simeoni et al [50], is that their method doesn't exploit the discriminative information for the descriptor which limits their method to only retrieval re-ranking. Based on the idea of independent matching for each channel, we refine the process by aggregating the concept from different channels as **Group-Concept Detect-and-Describe (GC-DAD)** for both keypoint detection and description.

Unlike D2-Net, which takes the maximum value across feature map channel, our method depends on the L2 response as the importance of the keypoint. For extracting different concepts, we would like to divide the feature map into different groups as independent concepts. The collection of the groups is represented by $\{F\}^g$ with g as the group index. For example, if a feature map F contains K channels and we want to uniformly divide it into G groups, then $F^{g=1}$ is corresponding to channel $1 \sim \lfloor \frac{K}{G} \rfloor$, and F^g is corresponding to channel $((g-1) * \lfloor \frac{K}{G} \rfloor + 1) \sim (\lfloor \frac{K}{G} \rfloor * g)$. We may compute the i -th keypoint location (x_i, y_i) out of the L2 response map of feature map F^g and combining the keypoint sets from every group by setting a L2 response threshold on the L2-norm of each group $F_{(x_i, y_i)}^g$ along with Harris edge threshold and the 2-nd order displacement similarly to [15].

To illustrate the power of the GC-DAD process, we show in Figure 3 results for the Aachen dataset for the day-night shifting camera localization [47]. Our method is strong enough even by using a pretrained ImageNet off-the-shelf ResNet101 without any extra training, while others requires self-learning [14] or fine-tuning [15].

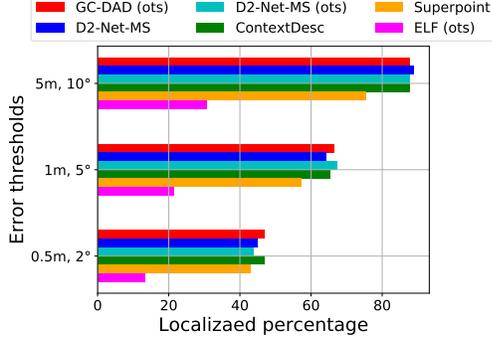


Figure 3. Aachen day-night localization benchmark [47] results. Comparison between 3 state of the art end-to-end matching methods, and our group-concept detect-and-describe paradigm (GC-DAD) with a simple pre-trained ImageNet ResNet101 network [13, 45, 18]. Surprisingly, by carefully utilising the feature maps with the proposed GC-DAD method, even a pre-trained ImageNet ResNet101 is capable of defeating every state-of-the-art local matching method without any fine-tuning or further training.

There are some drawbacks of performing GC-DAD on high dimensional descriptors such as descriptors aggregated from B2+B3 blocks of ResNet101 (more information on this will be discussed in the implementation details). The final descriptor is still memory costly since GC-DAD is only responsible for detecting the keypoints and the channel dimension out of the given feature map is preserved. In terms of practicality, a low dimensional descriptor is much more suitable for large-scale applications. Apart from that, manually selected groups is not the optimal choice since there is no clear clue about how to aggregate them. To that end, we explore a dimensionality reduction mapping with a channel dropping method in the next section.

Concept dropout dimension reduction A straight forward way to implement the grouping would be dividing the channel into several parts evenly with the same dimension. However, each channel contains different information and concepts. There is no clear evidence about which channels should be grouped together. Another trivial solution would be learning multiple mapping layers as aggregation learners for groups, but such setting is also redundant since several mapping layers have to be defined for the groups.

To avoid the aforementioned problems of non-trivial grouping and high dimensional descriptor, we adopt 2D dropout on the feature maps and perform a convolutional mapping for dimension reduction. By randomly dropping the feature channels in the training phase, it’s conceptually selecting different groups without manual selection, and brings the diverse concepts into the matching and avoid overfitting.

$$\hat{F} = \text{Conv}(\text{Drop}_{2D}(F)) \quad (2)$$

with feature map F as the high dimensional input for the dimension reduction, and \hat{F} as the low dimensional results.

Matching loss for affinity matrix For training the discriminative local descriptor, given different input images a, p , and the corresponding feature maps \hat{F}_a, \hat{F}_p along with the matching affinity matrix M can be computed by the inner product.

$$M = \hat{F}_a \cdot \hat{F}_p^T \quad (3)$$

By taking the maximum value along the column and row of the matching affinity matrix, the average score s can be computed

$$s_{a,p} = \frac{1}{2N_1} \sum_{i \in N_1} \max_j \hat{f}_{a,i} \cdot \hat{f}_{p,j}^T + \frac{1}{2N_2} \sum_{j \in N_2} \max_i \hat{f}_{a,i} \cdot \hat{f}_{p,j}^T \quad (4)$$

with i, j as the index of the column and row. A margin loss L_M is adopted as

$$L_M = \frac{1}{K} \sum_{k \in K} \max(s_{a,n_k} - s_{a,p} + m, 0). \quad (5)$$

which indicates the score of the positive pair is higher than the score of the negative pair by a margin m in terms of matching. Note that optimizing the matching loss does not require pixel level nor pointwise correspondence supervision. The local supervision we used is based on image pairs, similarly to the ones used for training image retrieval methods. This supervision is very weak comparing to the state-of-the-art methods in Table 1.

Cross dimension distillation During the feature mapping, the low dimensional descriptor \hat{F} cannot capture the whole information from the high dimensional descriptor F by only using a matching loss. Normally, distillation [19, 46] is a good way for transferring the information from a teacher model to a student model. Nevertheless, HF-Net adopts direct distillation which requires same descriptor dimensions for both teacher and student output, and in addition, the teacher models require training on extra dataset.

Here we propose a new way for distilling high dimension information into a low dimensional descriptor indirectly.

$$L_{Dis} = \left(M_{\text{high}} \cdot M_{\text{high}}^{(\text{det})} - \hat{M}_{\text{low}} \right)^2 \quad (6)$$

We refer the high dimensional one as teacher and the low dimensional one as student. The distillation happens on the matching affinity matrix M_{high} and \hat{M}_{low} with the corresponding soft detection score $M_{\text{high}}^{(\text{det})}$ for the high dimensional feature F as described in [15]. The affinity matrix is irrelevant to the feature dimension which is similar to [57]. In order to avoid the training of low dimension descriptor affect the high dimension one, we adopt backbone freezing or gradient cutdown for the training for the distillation. Comparing to HF-Net, our teacher descriptor is based on the same backbone with concatenation trick (B2,B3), while HF-Net [46] require extra teacher (DOAP [17], NetVLAD [2]) for the training.

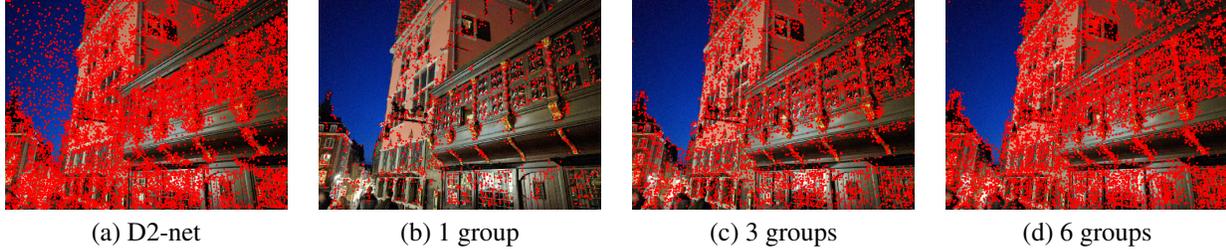


Figure 4. Grouping channel keypoint detection comparison. D2-net chooses the maximum value of each channel which results in high density of keypoints. On the other hand, our feature channel grouping technique combined with the feature channel L2-norm based thresholding can help us to control different levels of keypoint activation. Results here are based on the high dimensional teacher detection.

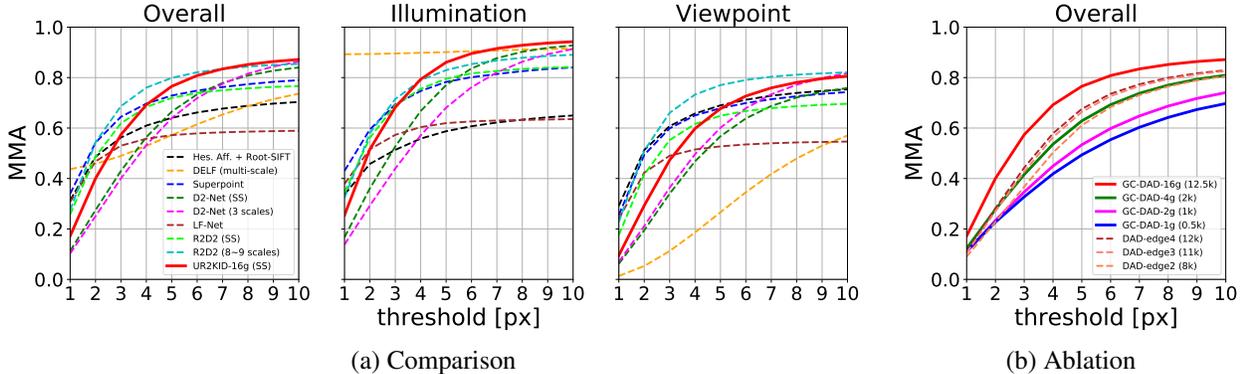


Figure 5. (a) Hpatches dataset mean matching accuracy (MMA) comparison. We can observe that the proposed method outperforms all other methods at the threshold of 4 pixels and above. (b) Ablation study for DAD and GC-DAD for the same network ResNet101. We can observe that more groups lead to increased performance while DAD has limited performance even with more keypoints (number in bracket).

3.2. Global description

Generalized-mean pooling Proposed in [42], the global representation is given by the generalized mean operation (GeM) as

$$d_g = \left(\frac{1}{hw} \sum_{s \in hw} x_{c,s}^p \right)^{\frac{1}{p}}. \quad (7)$$

with hw as the total spatial dimension of the input feature map, $x_{c,s}$ as the feature from channel c and spatial dimension s where $s \in [0, hw]$ and $p = 3$ as described in GeM [42]. Following the previous works [42, 15], we use a siamese architecture and train a two-branch network. Each branch shares the same network’s architecture and parameters.

The final global descriptor is the concatenation of all the FPN output along with Eq.7 pooling. To train the network on non-matching ($y = 0$) and matching pairs ($y = 1$), we employ the contrastive loss.

$$L_C = \begin{cases} \frac{1}{2} \|d_g(a) - d_g(p)\|^2 & y = 1 \\ \frac{1}{2} (\max\{0, \tau - \|d_g(a) - d_g(n_k)\|\})^2 & y = 0 \end{cases} \quad (8)$$

where τ is the margin.

3.3. Joint local and global training

By considering the global context and the local keypoint information, the final loss for jointly optimizing the local and global tasks is

$$L = L_M(F_{B2}) + L_M(F_{B3}) + L_M(\hat{F}) + L_C + \lambda L_{Dis}. \quad (9)$$

Considering that the distillation loss λL_{Dis} is not directly optimizing the metric, we use a relatively small hyperparameter $\lambda = 0.1$ parameter value for controlling it.

3.4. Implementation details

For the feature map representation that is used, we feed a single-scale image of an arbitrary size, and output four feature maps at multiple levels, in a fully convolution style. We do this to exploit all of the information from low to high level to the global pooling method. We conjecture that different concept levels will provide more useful features for a global image representation.

We first extract the features from the last four residual blocks in bottom-up path way (after the ReLU and before the pooling layers), which are in different sizes as described in Figure 2. Given the bottom-up feature map C_r with layer r , the final representation we used for each layer is F_r as

described in the following equation.

$$F_r = \text{ReLU}(\text{Conv}_{3 \times 3}(C_r) + \text{Upsampling}(F_{r+1})) \quad (10)$$

which follows the top-down merging in FPN [26]. The ReLU is applied to each feature map to make sure that it is non-negative.

Our local feature representation is made by the B2 and B3 residual blocks as described in Figure 2 (B refers to a Block) which are optimized jointly in training stage. During testing, in order to maintain multi-scale information, we take the feature maps from B2 & B3 and concatenate them. The spatial differences between B2 and B3 caused by pooling are resolved by the dilation trick in Figure 2 (b) which replaces stride 2 by stride 1 along with dilated convolution kernel. High resolution of the output feature map is essential for keypoint detection and matching.

4. Experiment

In this section, we present several experimental results that illustrate the power of the proposed GC-DAD method. In addition, we present ablation studies related to several parameters of our method.

4.1. Dataset and Setting

Training Megadepth [25] dataset was adopted for training which contains 1,070,468 images from 196 different scenes and reconstructed by COLMAP [48, 49] along with their depth maps and intrinsics / extrinsics matrices.

We use the Adam optimizer in training with the parameters $\alpha, \beta = (0.9, 0.99)$. During the training procedure, we make the learning rate α decay at i -th epoch with an exponential rate of $\exp^{-0.1i}$. We treat each training sample as a tuple of one query, one positive and five negative images.

For the training, we use exactly the same image pairs as D2-Net [15] from Megadepth dataset [25] for the sake of fair comparison. However, we do not utilize the point-wise correspondences, or the depth map ground truth for our training method. All models are trained up to 100 epochs. The batch size is set to 5 tuples, and the margin τ is set to 0.85. For each training epoch, around 6k tuples are selected.

Testing To evaluate the performance of the local matching pipeline, and of the global image retrieval, we examine several different datasets for testing a set of representative scenarios.

(a) **Hpatches** dataset [5] is the most well-known image matching dataset for identifying the matching robustness against different illumination and viewpoint changes. We compute the MMA (mean matching accuracy) as indicated in the D2-Net [15] along with other state-of-the-art methods for verifying the performance of our method. (b) **Aachen Day-Night** dataset contains 98 night-time query images in the testing dataset, along with 20 relevant images in day-time with known ground truth camera poses.

	day	night
Protocol 1 (Pre-defined query candidates)		
ELF	-	13.3 / 21.4 / 30.6
SuperPoint	-	42.8 / 57.1 / 75.5
DELf (new)	-	39.8 / 61.2 / 85.7
D2-Net (single)	-	44.9 / 66.3 / 88.8
D2-Net (multi)	-	44.9 / 64.3 / 88.8
R2D2 (web)	-	43.9 / 61.2 / 77.6
R2D2 (aachen day)	-	45.9 / 65.3 / 86.7
ContextDesc	-	46.9 / 65.3 / 87.8
UR2KID (single)	-	46.9 / 67.3 / 88.8
Protocol 2 (Global retrieval for candidate ranking)		
ESAC (50 experts)	42.6 / 59.6 / 75.5	3.1 / 9.2 / 11.2
AS	57.3 / 83.7 / 96.6	19.4 / 30.6 / 43.9
NV+Superpoint	79.7 / 88.0 / 93.7	40.8 / 56.1 / 74.5
HF-Net	75.7 / 84.3 / 90.9	40.8 / 55.1 / 72.4
NV+D2-Net (single)	79.7 / 89.3 / 94.8	41.8 / 63.3 / 81.6
UR2KID (single)	79.9 / 88.6 / 93.6	45.9 / 64.3 / 83.7

Table 2. Aachen day-night comparison with the localization threshold for day: (0.25m, 2°) / (0.5m, 5°) / (5m, 10°), and night: (0.5m, 2°) / (1m, 5°) / (5m, 10°). Our method is able to achieve top results in both scenarios.

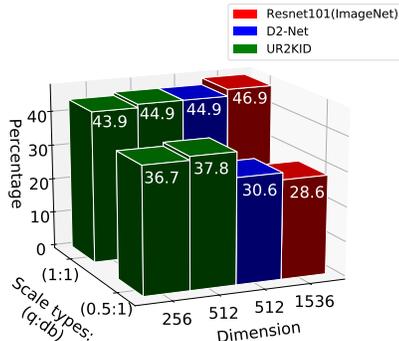


Figure 6. Scale comparison for Aachen dataset with localization threshold at night: (0.5m, 2°). The query versus database scales are shown as (q, db). We can observe that our method is very robust against scale changes.

After the keypoint and descriptor extraction based on 4479 database images, COLMAP pipeline is adopted for the 3D reconstruction. We follow the evaluation protocol from [47] and D2-Net [15], and the percentage of the queries localized within a given error bound on the estimated camera position and orientation is reported. Both keypoint detection and descriptor matching contribute to the camera localization task. (c) **Oxford5k, Paris6k** dataset which consist of 55 query images with bounding boxes, with the images exhibiting significant background noise such as people, trees etc. 5,062 building images are captured in Oxford, and 6,412 images of landmarks in Paris. Following the GeM setting, the global descriptor is tested in both datasets.

Dim			$r_q : r_{db} = 1 : 1$	$r_q : r_{db} = 0.5 : 1$			$r_q : r_{db} = 0.5 : 0.5$			$r_q, r_{db} \in [0.25, 1]$				
UR2KID (TS): teacher detect / student desc														
B2	B3	Fix	0.5m	1m	5m	0.5m	1m	5m	0.5m	1m	5m	0.5m	1m	5m
256	256	✓	46.9	67.3	88.8	31.6	49.0	66.3	35.7	56.1	78.6	35.7	54.1	74.5
256	256	-	44.9	66.3	88.8	28.6	49.0	67.3	36.7	63.3	81.6	37.8	62.2	83.7
128	128	✓	44.9	68.4	88.8	27.6	44.9	63.3	34.7	51.0	76.5	30.6	51.0	69.4
128	128	-	41.8	65.3	86.7	28.6	43.9	61.2	34.7	53.1	77.6	34.7	51.0	71.4
UR2KID (SS): student detect / student desc														
256	256	✓	44.9	69.4	88.8	37.8	52.0	73.5	41.8	64.3	86.7	39.8	62.2	80.6
256	256	-	43.9	67.3	87.8	34.7	46.9	73.5	38.8	60.2	84.7	35.7	62.2	79.6
128	128	✓	43.9	66.3	86.7	36.7	51.0	73.5	37.8	59.2	84.7	38.8	59.2	82.7
128	128	-	42.9	66.3	87.8	33.7	49.0	69.4	40.8	56.1	82.7	37.8	59.2	79.6
d2-net (512)			44.9	66.3	88.8	30.6	45.9	65.3	36.7	58.2	80.6	38.8	55.1	80.6

Table 3. Ablation study for mapped student feature output. The "Fix" option means the weights before ResNet block2 (B2) and block3 (B3) are frozen or not. The "Dim" means the final mapped output local descriptor dimension.

	Oxf5k	Par6k
SIFT	51.64	52.23
Geodesc	54.98	55.02
Contextdesc	65.03	64.53
LIFT	54.0	53.6
NetVLAD	71.6	79.7
UR2KID (megadepth)	82.03	91.94
GeM	88.17	92.6
DAME	88.24	93.0
UR2KID (sfm120k)	88.75	93.0
DELTA	90.0	95.7

Table 4. Image retrieval mean average precision (mAP) comparison for Oxford5k and Paris6k datasets. Our method not only achieves state-of-the-art on the previous matching benchmarks, but also has competitive performance in global landmark retrieval.

4.2. Hpatches dataset

Similar to the previous state-of-the-art methods in patch descriptor or matching pipeline, we evaluate the proposed UR2KID on the Hpatches dataset. In Figure 5(a), we demonstrate the comparison over multiple state-of-the-art methods in both illumination changes and viewpoint changes. Considering that our training does not depend on any special augmentation and pointwise correspondence ground-truth supervision, it is remarkable to see UR2KID outperform other methods such as Superpoint [14] and Lf-Net [39] by a margin in illumination changes. Notice that compared to D2-Net [15], which is the most relevant method comparing to ours, a clear advantage can be observed in the experimental results considering our method is only based on single scale while D2-Net [15] uses multiple scale inputs for boosting the performance in viewpoint changes. Among the state-of-the-art methods, DELF [38] is the most extreme case considering it's the most robust one against illumination change while it is also the most vulnerable one against viewpoint changes. Note that R2D2 [43]

achieves very high performance on the viewpoint changes with very large amount of scales (8 ~ 9 scales) which consume linear amount of time against the sampled scales. On the other hand, our method can outperform them on the single scale constraint.

Ablation studies In Figure 5(b), we show the ablation study among different group choices as described in Section 3.1. Generally speaking, increasing the keypoint number by tuning the threshold is capable of generating more correspondence candidates. However, as shown in Figure 4, there are a lot of portion of the keypoints detected by D2-Net which are not on the target buildings and cause false matches. For a fair comparison, we compare DAD and GC-DAD with different parameters on the same trained network, ResNet101, based on UR2KID pipeline and the keypoint numbers are also shown in Figure 5(b). We can observe that Our method with increased number of groups is able to capture the different concepts across the channel while maintaining good correspondence rate.

4.3. Aachen Day-Night dataset

Comparing to a pure local matching benchmark such as HPatches, a localization benchmark such as Aachen-day-night [47] involves more steps including local matching, geometric verification, triangulation, bundle adjustment, and solving the PnP problem. A structure-from-motion model is built based on the database images while the goal of the query is to recover the camera pose from the given query image. It's a more challenge standard for the sparse local matching based methods because local matching is the first step of the SfM pipeline and the error will be amplified after going through the aforementioned steps.

In Table 2, we demonstrate two different evaluation protocols. The first one is supported by the ground-truth query-database pair candidates. In this case, it is only required to match the local descriptor sets between the query image and the provided candidates from the database images. The pro-

posed UR2KID is the only method that optimizing the local matches without pointwise supervision and the best performance is obtained, with a single scale.

The second protocol is based on global descriptor retrieval and there is no given query-database pair. We follow the similar setting as described in [15, 47, 46], and the top-20 retrieved images are consider as the query-database pairs for the local matching. Using NetVLAD with D2-Net single scale for the localization, the performance degrades a little comparing to protocol 1 due to the fact that retrieved images may not form the optimal query-database pair comparing to the ground-truth. ESAC [12] is the family of unifying partial pipeline of the localization process. Nevertheless, it still suffers from the training data overfitting similar to PoseNet [21] and DSAC [11]. By using the global descriptor from our multi-task framework, UR2KID is able to provide both retrieval candidates and the local matching for the localization task. Our method is comparable in day case and outperforms the other methods at night.

Ablation studies In order to identify the strength of the proposed method, we examine ablation studies in Figure 6 and Table 3 based on the protocol 1 in Aachen dataset.

Despite the localization task is more challenging comparing to the pure matching benchmark, the dataset design is still non-realistic enough considering we cannot be sure about the target building size is large enough in the taken query image for localization. Therefore, we make different ratios between the query and the databased images into $(r_q : r_{db}) = (1 : 1)$ and $(r_q : r_{db}) = (0.5 : 1)$ in Figure 6. As we previously stated, the ResNet101 is indeed very robust in the localization task with our GC-DAD paradigm with high dimensional local descriptors. However, the performance degrades severely when it comes to $(r_q : r_{db}) = (0.5 : 1)$ with massive scale changes and similar degradation is observed in D2-Net. On the other hand, our method is not only robust in $(r_q : r_{db}) = (1 : 1)$, and the performance in $(r_q : r_{db}) = (0.5 : 1)$ case is also much higher comparing the others even with much lower dimension descriptor (e.g. 256, 512).

In Table 3, more details about the training process along with different scale changes variants are compared between our UR2KID and D2-Net considering they are highly related. Four different ratio combinations are discussed. $(r_q : r_{db}) = (1 : 1)$, $(0.5 : 1)$, $(0.5 : 0.5)$ and a random ratio combination between $[0.25, 1]$. We examine different mapping dimension (256,128 per block), detector choices (teacher or student), and the frozen weight choice during training (freeze block2 and block3 or not). Among those combinations, we can see that the trained student detector with student descriptor achieved the best results in the different variants, while the teacher detector achieves the best results when there is no scale change between query and database images. In every case, freezing the weight during

training is the best option which indicates that fine-tuning the mapping layer only is enough for learning the low dimensional descriptor. On the other hand, D2-Net performs poorly comparing to ours, giving 5% – 7% worse results for severe scale changes.

4.4. Oxford5k, Paris6k dataset

In Table 4, we compare different methods with our global descriptor in the retrieval framework. Based on the previous ablation study we know that the localization performance can be optimized when we freeze the weights and fine-tune the mapping layer only, and the retrieval task is also trained based on such setting. Our method is better than most visual word retrieval method such as ContextDesc [28], GeoDesc [29], and LIFT [65] when trained on the megadept dataset. However, the performance is only comparable or lower than the retrieval task specific methods such as GeM [42], DAME [63] and DELF [38]. As shown in Table 4, the performance is greatly improved by using SfM120k dataset as suggested by GeM [42]. In Oxford5k dataset, there are significant scale differences of the landmarks between the query image and the dataset reference image as discussed in [63] which is consistent with the SfM120k dataset, while Paris6k suffers less from that point of view. Megadept is suitable for matching as in D2-Net [15] by estimating the overlap ratio, but it is not suitable for global retrieval training.

These competitive results indicate that it is possible to embed local and global tasks into one single network.

5. Conclusion

We propose a multi-task framework for unifying global context along with local descriptors suitable for both retrieval and local matching tasks. The method exploits full image pairs instead of pointwise supervision during training while exhibiting state-of-the-art matching and localization results. In addition, we explore how significant scale changes affect the localization benchmark and identify that previous state-of-the-art methods are vulnerable against scale changes between query and database. Compared to other methods, our method is more robust against illumination and viewpoint changes, day-night shifting, and scale changes. As many new strong backbones (e.g. EfficientNet [53]) have been recently developed, we believe that our work is important for simplifying multi-task methods into a single network with very weak label information and can inspire similar future studies.

Acknowledgement This work was supported partially by the Computer Vision Lab of Tohoku University.

References

- [1] Alexandre Alahi, Raphael Ortiz, and Pierre Vanderghelynst. Freak: Fast retina keypoint. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 510–517. Ieee, 2012. 1
- [2] Relja Arandjelovic, Petr Gronát, Akihiko Torii, Tomás Pajdla, and Josef Sivic. Netvlad: CNN architecture for weakly supervised place recognition. In *CVPR*, 2016. 2, 3.1
- [3] Artem Babenko, Anton Slesarev, Alexander Chigorin, and Victor S. Lempitsky. Neural codes for image retrieval. In *ECCV*, 2014. 2
- [4] Vassileios Balntas, Edward Johns, Lilian Tang, and Krystian Mikolajczyk. Pn-net: Conjoined triple deep network for learning local image descriptors. 2016. 1, 2
- [5] Vassileios Balntas, Karel Lenc, Andrea Vedaldi, Tinne Tuytelaars, Jiri Matas, and Krystian Mikolajczyk. Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 1, 4.1
- [6] Vassileios Balntas, Edgar Riba, Daniel Ponsa, and Krystian Mikolajczyk. Learning local feature descriptors with triplets and shallow convolutional neural networks. 1, 2
- [7] Vassileios Balntas, Lilian Tang, and Krystian Mikolajczyk. Bold-binary online learned descriptor for efficient image matching. In *CVPR*, 2015. 1
- [8] Herbert Bay, Andreas Ess, Tinne Tuytelaars, and Luc J. Van Gool. Speeded-up robust features (SURF). *CVIU*, 2008. 2
- [9] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *European conference on computer vision*, pages 404–417. Springer, 2006. 2
- [10] Assia Benbihi, Matthieu Geist, and Cédric Pradalier. Elf: Embedded localisation of features in pre-trained cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7940–7949, 2019. 2
- [11] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. DSAC - differentiable RANSAC for camera localization. In *CVPR*, 2017. 4.3
- [12] Eric Brachmann and Carsten Rother. Expert sample consensus applied to camera re-localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7525–7534, 2019. 4.3
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3, 3
- [14] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *CVPRw*, 2018. 1, 2, 3.1, 4.2
- [15] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-net: A trainable cnn for joint detection and description of local features. *arXiv preprint arXiv:1905.03561*, 2019. 1, 1, 2, 3.1, 3.1, 3.2, 4.1, 4.1, 4.2, 4.3, 4.4
- [16] Steffen Gauglitz, Tobias Höllerer, and Matthew Turk. Evaluation of interest point detectors and feature descriptors for visual tracking. *International journal of computer vision*, 94(3):335, 2011. 1
- [17] Kun He, Yan Lu, and Stan Sclaroff. Local descriptors optimized for average precision. In *CVPR*, 2018. 1, 2, 3.1
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 3
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 3.1
- [20] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. 2
- [21] Alex Kendall, Matthew Grimes, and Roberto Cipolla. PoseNet: A convolutional network for real-time 6-dof camera relocalization. In *ICCV*, 2015. 4.3
- [22] Han-Ul Kim, Dae-Youn Lee, Jae-Young Sim, and Chang-Su Kim. Sowp: Spatially ordered and weighted patch descriptor for visual tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3011–3019, 2015. 1
- [23] Axel Barroso Laguna, Edgar Riba, Daniel Ponsa, and Krystian Mikolajczyk. Key. net: Keypoint detection by handcrafted and learned cnn filters. *arXiv preprint arXiv:1904.00889*, 2019. 1, 2
- [24] Stefan Leutenegger, Margarita Chli, and Roland Siegwart. Brisk: Binary robust invariant scalable keypoints. In *2011 IEEE international conference on computer vision (ICCV)*, pages 2548–2555. Ieee, 2011. 1
- [25] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2041–2050, 2018. 1, 2, 4.1
- [26] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 3, 3.4
- [27] David G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 2004. 1, 2, 3.1
- [28] Zixin Luo, Tianwei Shen, Lei Zhou, Jiahui Zhang, Yao Yao, Shiwei Li, Tian Fang, and Long Quan. Contextdesc: Local descriptor augmentation with cross-modality context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2527–2536, 2019. 2, 4.4
- [29] Zixin Luo, Tianwei Shen, Lei Zhou, Siyu Zhu, Runze Zhang, Yao Yao, Tian Fang, and Long Quan. Geodesc: Learning local descriptors by integrating geometry constraints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 168–183, 2018. 4.4
- [30] Jiri Matas, Ondrej Chum, Martin Urban, and Tomás Pajdla. Robust wide-baseline stereo from maximally stable extremal regions. *Image and vision computing*, 22(10):761–767, 2004. 3.1
- [31] Krystian Mikolajczyk and Cordelia Schmid. An affine invariant interest point detector. In *European conference on computer vision*, pages 128–142. Springer, 2002. 2

- [32] Anastasiia Mishchuk, Dmytro Mishkin, Filip Radenovic, and Jiri Matas. Working hard to know your neighbor's margins: Local descriptor learning loss. In *NIPS*, 2017. 1, 2
- [33] Dmytro Mishkin, Filip Radenovic, and Jiri Matas. Repeatability is not enough: Learning affine regions via discriminability. In *ECCV*, 2018. 2
- [34] Kwang Moo Yi, Yannick Verdie, Pascal Fua, and Vincent Lepetit. Learning to assign orientations to feature points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 107–116, 2016. 2
- [35] Arun Mukundan, Giorgos Tolias, and Ondrej Chum. Explicit spatial encoding for deep local descriptors. In *CVPR*, 2019. 1, 2
- [36] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015. 1, 2
- [37] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, 2017. 1, 2
- [38] Hyeonwoo Noh, Andre Araujo, Jack Sim, and Tobias Weyand. Large-Scale Image Retrieval with Attentive Deep Local Features. In *ICCV*, 2017. 1, 2, 4.2, 4.4
- [39] Yuki Ono, Eduard Trulls, Pascal Fua, and Kwang Moo Yi. Lf-net: learning local features from images. In *NIPS*, 2018. 1, 2, 4.2
- [40] Filip Radenović, Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Revisiting Oxford and Paris : Large-Scale Image Retrieval Benchmarking. In *CVPR*, 2018. 2
- [41] Filip Radenović, Giorgos Tolias, and Ondřej Chum. CNN Image Retrieval Learns from BoW: Unsupervised Fine-Tuning with Hard Examples. In *ECCV*, 2016. 1, 2
- [42] Filip Radenović, Giorgos Tolias, and Ondrej Chum. Fine-tuning CNN Image Retrieval with No Human Annotation. *PAMI*, 2018. 2, 3, 3.2, 3.2, 4.4
- [43] Jerome Revaud, Philippe Weinzaepfel, César De Souza, Nœ Pion, Gabriela Csurka, Johann Cabon, and Martin Humenberger. R2d2: Repeatable and reliable detector and descriptor. *arXiv preprint arXiv:1906.06195*, 2019. 2, 4.2
- [44] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary R Bradski. Orb: An efficient alternative to sift or surf. In *ICCV*. Citeseer, 2011. 1, 2
- [45] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 3, 3
- [46] P.-E. Sarlin, C. Cadena, R. Siegwart, and M. Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. 2019. 1, 2, 3.1, 3.1, 4.3
- [47] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, et al. Benchmarking 6dof outdoor visual localization in changing conditions. In *CVPR*, 2018. 1, 3.1, 3, 4.1, 4.3
- [48] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 4.1
- [49] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 1, 4.1
- [50] Oriane Siméoni, Yannis Avrithis, and Ondrej Chum. Local features and visual words emerge in activations. In *CVPR*, 2019. 3.1
- [51] Chris Sweeney. Theia multiview geometry library: Tutorial & reference. <http://theia-sfm.org>. 1
- [52] Chris Sweeney, Victor Fragoso, Tobias Höllerer, and Matthew Turk. Large scale sfm with the distributed camera model. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 230–238. IEEE, 2016. 1
- [53] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019. 5
- [54] Marvin Teichmann and Jack Sim. Detect-to-Retrieve: Efficient Regional Aggregation for Image Search. In *CVPR*, 2019. 2
- [55] Yurun Tian, Bin Fan, and Fuchao Wu. L2-net: Deep learning of discriminative patch descriptor in euclidean space. In *CVPR*, 2017. 1, 2
- [56] Yurun Tian, Xin Yu, Bin Fan, Fuchao Wu, Huub Heijnen, and Vassileios Balntas. Sosnet: Second order similarity regularization for local descriptor learning. In *CVPR*, 2019. 1, 2
- [57] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1365–1374, 2019. 3.1
- [58] Zhenhua Wang, Bin Fan, and Fuchao Wu. Local intensity order pattern for feature description. In *2011 International Conference on Computer Vision*, pages 603–610. IEEE, 2011. 2
- [59] Jian Xu, Cunzhao Shi, Chengzuo Qi, Chunheng Wang, and Baihua Xiao. Unsupervised Part-based Weighting Aggregation of Deep Convolutional Features for Image Retrieval. In *AAAI*, 2018. 2
- [60] Jian Xu, Chunheng Wang, Chengzuo Qi, Cunzhao Shi, and Baihua Xiao. Unsupervised Semantic-based Aggregation of Deep Convolutional Features. *TIP*, 2019. 2
- [61] Hanxuan Yang, Ling Shao, Feng Zheng, Liang Wang, and Zhan Song. Recent advances and trends in visual tracking: A review. *Neurocomputing*, 74(18):3823–3831, 2011. 1
- [62] Tsun-Yi Yang, Jo-Han Hsu, Yen-Yu Lin, and Yung-Yu Chuang. Deepcd: Learning deep complementary descriptors for patch representations. In *ICCV*, 2017. 1, 2
- [63] Tsun-Yi Yang, Duy Kien Nguyen, Huub Heijnen, and Vassileios Balntas. Dame web: Dynamic mean with whitening ensemble binarization for landmark retrieval without human annotation. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2, 4.4

- [64] Tsun-Yi Yang, Yen-Yu Lin, and Yung-Yu Chuang. Accumulated stability voting: A robust descriptor from descriptors of multiple scales. In *CVPR*, 2016. [1](#)
- [65] Kwang Moo Yi, Eduard Trulls, Vincent Lepetit, and Pascal Fua. Lift: Learned invariant feature transform. In *European Conference on Computer Vision*, pages 467–483. Springer, 2016. [1](#), [2](#), [4.4](#)