

GeoDesc: Learning Local Descriptors by Integrating Geometry Constraints

Zixin Luo^{1†}[0000-0001-6946-2826], Tianwei Shen^{1†}[0000-0002-3290-2258],
Lei Zhou^{1†}[0000-0003-4988-5084], Siyu Zhu^{1‡}[0000-0003-0293-0044],
Runze Zhang^{1†||}[0000-0001-9698-0178], Yao Yao^{1†}[0000-0001-9866-4291],
Tian Fang²[0000-0002-5871-3455], Long Quan¹[0000-0001-8148-1771]

¹ Hong Kong University of Science and Technology,
{zluoag, tshenaa, lzhouai, szhu, rzhangaj, yyaoag, quan}@cse.ust.hk
² Shenzhen Zhuke Innovation Technology (Altizure),
fangtian@altizure.com

Abstract. Learned local descriptors based on Convolutional Neural Networks (CNNs) have achieved significant improvements on patch-based benchmarks, whereas not having demonstrated strong generalization ability on recent benchmarks of image-based 3D reconstruction. In this paper, we mitigate this limitation by proposing a novel local descriptor learning approach that integrates geometry constraints from multi-view reconstructions, which benefits the learning process in terms of data generation, data sampling and loss computation. We refer to the proposed descriptor as GeoDesc, and demonstrate its superior performance on various large-scale benchmarks, and in particular show its great success on challenging reconstruction tasks. Moreover, we provide guidelines towards practical integration of learned descriptors in Structure-from-Motion (SfM) pipelines, showing the good trade-off that GeoDesc delivers to 3D reconstruction tasks between accuracy and efficiency.

Keywords: Local Features, Feature Descriptors, Deep Learning

1 Introduction

Computing local descriptors on interest regions serves as the subroutine of various computer vision applications such as panorama stitching [12], wide baseline matching [18], image retrieval [22], and Structure-from-Motion (SfM) [42, 26, 38, 41]. A powerful descriptor is expected to be invariant to both photometric and geometric changes, such as illumination, blur, rotation, scale and perspective changes. Due to the reliability, efficiency and portability, hand-crafted descriptors such as SIFT [14] have been influentially dominating this field for

[†]Zixin Luo, Lei Zhou and Yao Yao were summer interns, Tianwei Shen and Runze Zhang were interns at Everest Innovation Technology (Altizure).

[‡]Siyu Zhu is with Alibaba A.I. Labs since Oct. 2017.

^{||}Runze Zhang is the corresponding author.

more than a decade. Until recently, great efforts have been made on developing learned descriptors based on Convolutional Neural Networks (CNNs), which have achieved surprising results on patch-based benchmarks such as HPatches dataset [3]. However, on image-based datasets such as ETH local features benchmark [25], learned descriptors are found to underperform advanced variants of hand-crafted descriptors. The contradictory findings raise the concern of integrating those purportedly better descriptors in real applications, and show significant room of improvement for developing more powerful descriptors that generalize to a wider range of scenarios.

One possible cause of above contradictions, as demonstrated in [25], is the lack of generalization ability as a consequence of data insufficiency. Although previous research [4, 27, 28] discusses several effective sampling methods that produce seemingly large amount of training data, the generalization ability is still bounded to limited data sources, e.g., the widely-used Brown dataset [6] with only 3 image sets. Hence, it is not surprising that resulting descriptors tend to overfit to particular scenarios. To overcome it, research such as [29, 39] applies extra regularization for compact feature learning. Meanwhile, LIFT [34] and [19] seek to enhance data diversity and generate training data from reconstructions of Internet tourism data. However, the existing limitation has not yet been fully mitigated, while intermediate geometric information is overlooked in the learning process despite the robust geometric property that local patch preserves, e.g., the well approximation of local deformations [20].

Besides, we lack guidelines for integrating learned descriptors in practical pipelines such as SfM. In particular, the *ratio criterion*, as suggested in [14] and justified in [10], has received almost no individual attention or was considered inapplicable for learned descriptors [25], whereas it delivers excellent matching efficiency and accuracy improvements, and serves as the necessity for pipelines such as SfM to reject false matches and seed feasible initialization. A general method to apply ratio criterion for learned descriptors is in need in practice.

In this paper, we tackle above issues by presenting a novel learning framework that takes advantage of geometry constraints from multi-view reconstructed data. In particular, we address the importance of data sampling for descriptor learning and summarize our contributions threefold. i) We propose a novel batch construction method that simulates the pair-wise matching and effectively samples useful data for learning process. ii) Collaboratively, we propose a loss formulation to reduce overfitting and improve the performance with geometry constraints. iii) We provide guidelines about ratio criterion, compactness and scalability towards practical portability of learned descriptors.

We evaluate the proposed descriptor, referred to as GeoDesc, on traditional [9] and recent two large-scale datasets [3, 25]. Superior performance is shown over the state-of-the-art hand-crafted and learned descriptors. We mitigate previous limitations by showing consistent improvements on both patch-based and image-based datasets, and further demonstrate its success on challenging 3D reconstructions.

2 Related Works

Networks design. Due to weak semantics and efficiency requirements, existing descriptor learning often relies on shallow and thin networks, e.g., three-layer networks in DDesc [27] with 128-dimensional output features. Moreover, although widely-used in high-level computer vision tasks, max pooling is found to be unsuitable for descriptor learning, which is then replaced by L2 pooling in DDesc [27] or even removed in L2-Net [29]. To further incorporate scale information, DeepCompare [36] and L2-Net [29] use a two-stream central-surround structure which delivers consistent improvements at extra computational cost. To improve the rotational invariance, an orientation estimator is proposed in [35]. Besides of feature learning, previous efforts are also made on joint metric learning as in [36, 8, 7], whereas comparison in Euclidean space is more preferable by recent works [4, 27, 29, 34, 5] in order to guarantee its efficiency.

Loss formulation Various of loss formulations have been explored for effective descriptor learning. Initially, networks with a learned metric use softmax loss [36, 8] and cast the descriptor learning to a binary classification problem (similar/dissimilar). With weakly annotated data, [15] formulates the loss on keypoint bags. More generally, pair-wise loss [27, 34] and triplet loss [4, 7, 5] are used by networks without a learned metric. Both loss formulations encourage matching patches to be close whereas non-matching patches to be far-away in some measure space. In particular, triplet loss delivers better results [4, 7] as it suffers less overfitting [13]. For effective training, recent L2-Net [29] and Hard-Net [17] use the structured loss for data sampling which drastically improves the performance. To further boost the performance, extra regularizations are introduced for learning compact representation in [29, 39].

Evaluation protocol Previous works often evaluate on datasets such as [16, 31, 9]. However, those datasets either are small, or lack diversity to generalize well to various applications of descriptors. As a result, the evaluation results are commonly inconsistent or even contradictory to each other as pointed out in [3], which limits the application of learned descriptors. Two novel benchmarks, HPatches [3] and ETH local descriptor benchmark [25] have been recently introduced with clearly defined protocols and better generalization properties. However, inconsistency still exists in the two benchmarks, where HPatches [3] benchmark demonstrates the significant outperformance from learned descriptors over the handcrafted, whereas the ETH local descriptor benchmark [25] finds that the advanced variants of the traditional descriptor are at least on par with the learning-based. The inconclusive results indicate that there is still significant room for improvement to learn more powerful feature descriptors.

3 Method

3.1 Network architecture

We borrow the network in L2-Net [29], where the feature tower is constructed by eschewing pooling layers and using strided convolutional layers for in-network

downsampling. Each convolutional layer except the last one is followed by a batch normalization (BN) layer whose weighting and bias parameters are fixed to 1 and 0. The L2-normalization layer after the last convolution produces the final 128-dimensional feature vector.

3.2 Training data generation

Acquiring high quality training data is important in learning tasks. In this section, we discuss a practical pipeline that automatically produces well-annotated data suitable for descriptor learning.

2D correspondence generation. Similar to LIFT [34], we rely on successful 3D reconstructions to generate ground truth 2D correspondences in an automatic manner. First, sparse reconstructions are obtained from standard SfM pipeline [24]. Then, 2D correspondences are generated by projecting 3D point clouds. In general, SfM is used to filter out most mismatches among images.

Although verified by SfM, the generated correspondences are still outlier-contaminated from image noise and wrongly registered cameras. It happens particularly often on Internet tourism datasets such as [30, 23] (illustrated in Fig. 1(a)), and usually not likely to be filtered by simply limiting reprojection error. To improve data quality, we take one step further than LIFT by computing the visibility check based on 3D Delaunay triangulation [11] which is widely-used for outlier filtering in dense stereo. Empirically, 30% of 3D points will be discarded after the filtering while only points with high precision are kept for ground truth generation. Fig. 1(b) gives an example to illustrate its effect.

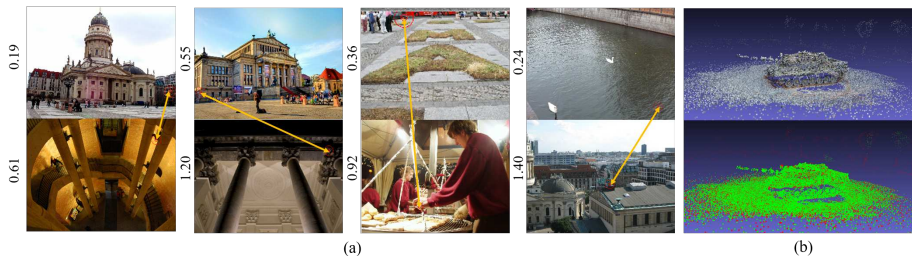


Fig. 1: (a) Outlier matches after SfM verification (by COLMAP [24]) on *Gendarmenmarkt* dataset [30]. The reprojection error (next to the image) cannot be used to identify false matches. (b) Reconstructed sparse point cloud (top), where points in red (bottom) indicate being filtered by Delaunay triangulation and only reliable points in green are kept. The number of points decreases from 75k to 53k after the filtering.

Matching patch generation. Next, the interest region of a 2D projection is cropped similar to LIFT, which is formulated by an similarity transformation

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \frac{k\sigma}{2} \cos(\theta) & \frac{k\sigma}{2} \sin(\theta) \\ -\frac{k\sigma}{2} \sin(\theta) & \frac{k\sigma}{2} \cos(\theta) \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \begin{pmatrix} x_i^t \\ y_i^t \end{pmatrix}, \quad (1)$$

where $(x_i^s, y_i^s), (x_i^t, y_i^t)$ are input and output regular sampling grids, and (x, y, σ, θ) are keypoint parameters (x, y coordinates, scale and orientation) from SIFT detector. The constant k is set to 12 as in LIFT, resulting in $12\sigma \times 12\sigma$ patches.

Due to the robust estimation of scale (σ) and orientation (θ) parameters of SIFT even in extreme cases [40], the resulting patches are mostly free of scale and rotation differences, thus suitable for training. In later experiments of image matching or SfM, we rely on the same cropping method to achieve scale and rotation invariance for learned descriptors.

3.3 Geometric similarity estimation

Geometries at a 3D point are robust and provide rich information. Inspired by the MVS (Multi-View Stereo) accuracy measurement in [37], we define two types of geometric similarity: patch similarity and image similarity, which will facilitate later data sampling and loss formulation.

Patch similarity. We define patch similarity S_{patch} to measure the difficulty to have a patch pair matched with respect to perspective changes. Formally, given a patch pair, we relate it to its corresponding 3D track P which is seen by cameras centering at C_i and C_j . Next, we compute the vertex normal P_n at P from the surface model. The geometric relationship is illustrated in Fig. 2(a). Finally, we formulate S_{patch} as

$$S_{patch} = s_1 s_2 = g(\angle C_i P C_j, \sigma_1) g(\angle C_i P P_n - \angle C_j P P_n, \sigma_2), \quad (2)$$

where s_1 measures the intersection angle between two viewing rays from the 3D track ($\angle C_i P C_j$), while s_2 measures the difference of incidence angles between a viewing ray and the vertex normal from the 3D track ($\angle C_i P P_n, \angle C_j P P_n$). The angle metric is defined as $g(\alpha, \sigma) = \exp(-\frac{\alpha^2}{2\sigma^2})$. As an interpretation, s_1 and s_2 measure the perspective change regarding a *3D point* and *local 3D surface*, respectively. The effect of S_{patch} is illustrated in Fig. 2(b).

The accuracy of s_1 and s_2 depends on sparse and mesh reconstructions, respectively, and is generally sufficient for its use as shown in [37]. The similarity does not consider scale and rotation changes as already resolved from Equation 1. Empirically, we choose $\sigma_1 = 15$ and $\sigma_2 = 20$ (in degree).

Image similarity. Based on the patch similarity, we define the image similarity S_{image} as the average patch similarity of the correspondences between an image pair. The image similarity measures the difficulty to match an image pair and can be interpreted as a measurement of perspective change. Examples are given in Fig. 2(c). The image similarity will be beneficial for data sampling in Sec. 3.4.

3.4 Batch construction

For descriptor learning, most existing frameworks take patch pairs (matching/non-matching) or patch triplets (reference, matching and non-matching) as input. As in previous studies, the convergence rate is highly dependent on being able to see

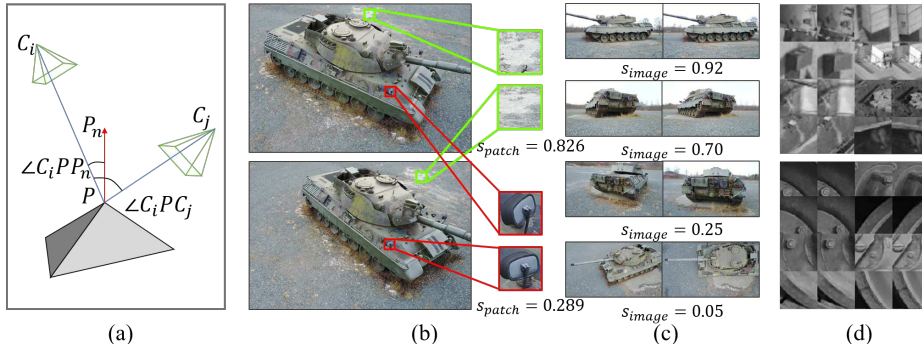


Fig. 2: (a) The patch similarity relies on the geometric relationship between cameras, tracks and surface normal. (b) The effect of patch similarity, which measures the difficulty to have a patch pair matched with respect to the perspective change. (c) The effect of image similarity, which measures the perspective change between an image pairs. (d) Batch data constructed by L2-Net [29] and HardNet [17] (top), whose in-batch patch pairs are often distinctive to each other and thus contribute nothing to the loss in the late learning (e.g., the margin-based loss). However, the batch data from the proposed batch construction method (bottom) consists of similar patch pairs due to the spatially close keypoints or repetitive patterns, which are considered harder to distinguish and thus raise greater challenges for learning

useful data [21]. Here, “useful” data often refers to patch pairs/triplets that produce meaningful loss for learning. However, the effective sampling of such data is generally challenging due to the intractably large number of patch pair/triplet combination in the database. Hence, on one hand, sampling strategies such as hard negative mining [27] and anchor swap [4] are proposed, while on the other hand, effective batch construction is used in [29, 7, 17] to compare the reference patch against all the in-batch samples in the loss computation.

Inspired by previous works, we propose a novel batch construction method that effectively samples “useful” data by relying on geometry constraints from SfM, including the image matching results and image similarity S_{image} , to simulate the pair-wise image matching and sample data. Formally, given one image pair, we extract a *match set* $X = \{(x_1, x_1^+), (x_2, x_2^+), \dots, (x_{N_1}, x_{N_1}^+)\}$, where N_1 is the set size and (x_i, x_i^+) is a matching patch pair surviving the SfM verification. A training batch is then constructed on N_2 match sets. Hence, the learning objective becomes to improve the matching quality for each match set. In Sec. 3.5, we will discuss the loss computation on each match set and batch data.

Compared with L2-Net [29] and HardNet [17] whose training batches are random sampled from the whole database, the proposed method produces harder samples and thus raises greater challenges for learning. As an example shown in Fig. 2(d), the training batch constructed by the proposed method consists of many similar patterns, due to the spatially close keypoints or repetitive textures. In general, such training batch has two major advantages for descriptor learning:

- It reflects the in-practice complexity. In real applications, image patches are often extracted between image pairs for matching. The proposed method simulates this scenario so that training and testing become more consistent.
- It generates hard samples for training. As observed in [4, 27, 17, 21], hard samples are critical to fast convergence and performance improvement for descriptor learning. The proposed method effectively generates batch data that is sufficiently hard, while not being overfitting as constructed on real matching results instead of model inference results [27].

To further boost the training efficiency, we exclude image pairs that are too similar to contribute to the learning. Those pairs are effectively identified by the image similarity defined in Sec. 3.3. In practice, we discard image pairs whose S_{image} are larger than 0.85 (e.g., the toppest pair in Fig. 2(c)), which results in a 30% shrink of training samples.

3.5 Loss formulation

We formulate the loss with two terms: structured loss and geometric loss.

Structured loss. The structured loss used in L2-Net [29] and HardNet [17] is essentially suitable to consume the batch samples constructed in Sec. 3.4. In particular, the formulation in HardNet based on the “hardest-in-batch” strategy and a distance margin shows to be more effective than the log-likelihood formulation in L2-Net. However, we observe successive overfitting when applying the HardNet loss to our batch data, which we ascribe to the too strong constraint of “hardest-in-batch” strategy. In this strategy, the loss is computed on the data sample that produces the largest loss, and a margin with a large value (1.0 in HardNet) is set to push the non-matching pairs away from matching pairs. In our batch data, we already effectively sample the “hard” data which is often visually similar, thus forcing a large margin is inapplicable and stalls the learning. One simple solution is to decrease the margin value, whereas the performance drops significantly in our experiments.

To avoid above limitation and better take advantage of our batch data, we propose the loss formulation as follows. First, we compute the structured loss for one match set. Given normalized features $\mathbf{F}_1, \mathbf{F}_2 \in \mathbb{R}^{N_1 \times 128}$ computed on match set X for all (x_i, x_i^+) , the cosine similarity matrix is derived by $\mathbf{S} = \mathbf{F}_1 \mathbf{F}_2^T$. Next, we compute $\mathbf{L} = \mathbf{S} - \alpha \mathbf{diag}(\mathbf{S})$ and formulate the loss as

$$E_1 = \frac{1}{N_1(N_1 - 1)} \sum_{i,j} (\max(0, l_{i,j} - l_{i,i}) + \max(0, l_{i,j} - l_{j,j})), \quad (3)$$

where $l_{i,j}$ is the element in \mathbf{L} , and $\alpha \in (0, 1)$ is the distance ratio mimicking the behavior of ratio test [14] and pushing away non-matching pairs from matching pairs. Finally, we take the average of the loss on each match set to derive the final loss for one training batch.

The proposed formulation is distinctive from HardNet in three aspects. First, we compute the cosine similarity instead of Euclidean distance for computational

efficiency. Second, we apply a *distance ratio margin* instead of a *fixed distance margin* as an adaptive margin to reduce overfitting. Finally, we compute the *mean* value of each loss element instead of the maximum (“hardest-in-batch”) in order to cooperate the proposed batch construction.

Geometric loss. Although E_1 ensures matching patch pairs to be distant from the non-matching, it does not explicitly encourage matching pairs to be close in its measure space. One simple solution is to apply a typical pair-wise loss in [27], whereas taking a risk of positive collapse and overfitting as observed in [13]. To overcome it, we adaptively set up the margin regarding the patch similarity defined in Sec. 3.3, serving as a soft constraint for maximizing the positive similarity. We refer to this term as *geometric loss* and formulate it as

$$E_2 = \sum_i \max(0, \beta - s_{i,i}), \beta = \begin{cases} 0.7 & s_{patch} \geq 0.5 \\ 0.5 & 0.2 \leq s_{patch} < 0.5 \\ 0.2 & \text{otherwise} \end{cases} \quad (4)$$

where β is the adaptive margin, $s_{i,i}$ is the element in S , namely, the cosine similarity of patch pair (x_i, x_i^+) , while s_{patch} is the patch similarity for (x_i, x_i^+) . We use $E_1 + \lambda E_2$ as the final loss, and empirically set α and λ to 0.4 and 0.2.

3.6 Training

We use image sets [30] as in LIFT [34], the SfM data in [23], and further collect several image sets to form the training database. Based on COLMAP [24], we run 3D reconstructions to establish necessary geometry constraints. Image sets that are overlapping with the benchmark data are manually excluded. We train the networks from scratch using Adam with a base learning rate of 0.001 and weight decay of 0.0001. The learning rate decays by 0.9 every 10,000 steps. Data augmentation includes randomly flipping, 90 degrees rotation and brightness and contrast adjustment. The match set size N_1 and batch size N_2 are 64 and 12, respectively. Input patches are standardized to have zero mean and unit norm.

4 Experiments

We evaluate the proposed descriptor on three datasets: the patch-based HPatches [3] benchmark, the image-based Heiny benchmark [9] and ETH local features benchmark [25]. We further demonstrate on challenging SfM examples.

4.1 HPatches benchmark

HPatches benchmark [3] defines three complementary tasks: patch verification, patch matching, and patch retrieval. Different levels of geometrical perturbations are imposed to form EASY, HARD and TOUGH patch groups. In the task of verification, two subtasks are defined based on whether negative pairs are sampled from images within the same (SAMESEQ) or different sequences (DIFFSEQ).

In the task of matching, two subtasks are defined based on whether the principle variance is viewpoint (VIEW) or illumination (ILLUM). Following [3], we use mean average precision (mAP) to measure the performance for all three tasks on HPatches split ‘full’.

We select five descriptors to compare: SIFT as the baseline, RSIFT [2] and DDesc [27] as the best-performing hand-crafted and learned descriptors concluded in [3]. Moreover, we experiment with recent learned descriptors L2-Net [29] and HardNet [17]. The proposed descriptor is referred to as GeoDesc.

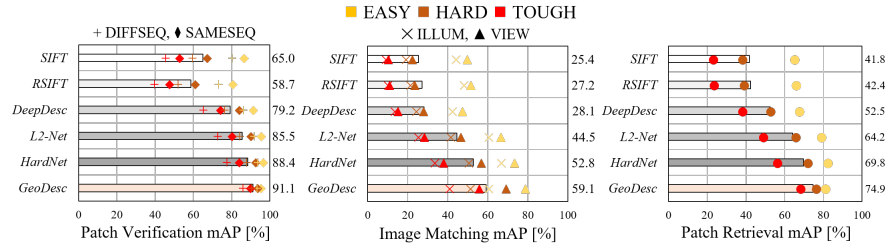


Fig. 3: Left to right: Verification, matching and retrieval results on HPatches dataset, split ‘full’. Results on different patch groups are colorized, while DIFFSEQ/SAMESEQ and ILLUM/VIEW denote the subtasks of verification and matching, respectively

As shown in Fig. 3, GeoDesc surpasses all the other descriptors on all three tasks by a large margin. In particular, the performance on TOUGH patch group is significantly improved, which indicates the superior invariance to large image changes of GeoDesc. Interestingly, comparing GeoDesc with HardNet, we observe some performance drop on EASY groups especially for illumination changes, which can be ascribed to the data bias for SfM data. Though applying randomness such as illumination during the training, we cannot fully mitigate this limitation which asks more diverse real data in descriptor learning.

In addition, we evaluate different configurations of GeoDesc on HPatches as shown in Tab. 1 to demonstrate the effect of each part of our method.

- *Config. 1*: the HardNet framework as the baseline.
- *Config. 2*: trained with the SfM data in Sec. 3.2. Compared with *Config. 1*, it is shown that crowd-sourced training data essentially improves the generalization ability. Meanwhile, on the other hand, *Config. 2* can be regarded as an extension of LIFT [34] with more advanced loss formulation.
- *Config. 3*: equipped with the proposed batch construction in Sec. 3.4. As discussed in Sec. 3.5, the “hardest-in-batch” strategy in HardNet is inapplicable to hard batch data and thus leads to performance drop compared with *Config. 2*. In practice, we need to adjust the margin value from 1.0 in HardNet to 0.6, otherwise the training will not even converge. Though trainable, the smaller margin value harms the final performance.
- *Config. 4*: equipped with the modified structured loss in Sec. 3.5. Notable performance improvements are achieved over *Config. 2* due to the collabora-

tive use of proposed methods, showing the effectiveness of simulating pair-wise matching and sampling hard data. Besides, replacing the *distance margin* with *distance ratio* can improve the training efficiency, as shown in Fig. 4.

- *Config. 5*: equipped with the geometric loss in Sec. 3.5. Further improvements are obtained over *Config. 4* as E_2 constrains the solution space and enhances the training efficiency.

To sum up, the “hardest-in-batch” strategy is beneficial when no other sampling is applied and most in-batch samples do not contribute to the loss. However, with harder batch data effectively constructed, it is advantageous to replace the “hardest-in-batch” and further boost the descriptor performance.

Table 1: Evaluation of different configurations of GeoDesc on HPatches. Modules are enabled if marked with “Y” otherwise with “-”. *SfM Data* denotes the training with our SfM data, *Batch Construct.* denotes the equipment of proposed batch construction, while E_1 and E_2 denote the use of proposed structured loss and geometric loss, respectively. The last configuration (*Config. 5*) is our best model with GeoDesc

No.	GeoDesc Configuration				HPatches Benchmark Tasks		
	<i>SfM Data</i>	<i>Batch Construct.</i>	E_1	E_2	Verification	Matching	Retrieval
1	-	-	-	-	88.4	52.8	69.8
2	Y	-	-	-	90.1	57.0	73.2
3	Y	Y	-	-	89.9	50.2	70.4
4	Y	Y	Y	-	90.9	58.5	74.5
5	Y	Y	Y	Y	91.1	59.1	74.9

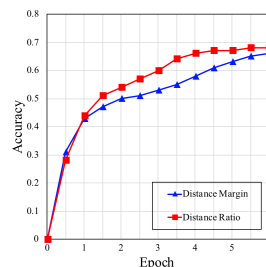


Fig. 4: Effect of taking *distance ratio* in loss computation. The metric is the validation accuracy of patch triplets with a margin of 0.5 by cosine similarity.

4.2 Heinly benchmark

Different from HPatches which experiments on image patches, the benchmark by Heinly *et al.* [9] evaluates pair-wise image matching regarding different types of photometric or geometric changes, targeting to provide practical insights for strengths and weaknesses of descriptors. We use two standard metrics as in [9] to quantify the matching quality. First, the *Matching Score* = $\#Inlier\ Matches / \#Features$. Second, the *Recall* = $\#Inlier\ Matches / \#True\ Matches$. Four descriptor are selected to compare: SIFT, the baseline hand-crafted descriptor; DSP-SIFT, the best hand-crafted descriptor even superior to the previous learning-based as evaluated in [25]; L2-Net and HardNet, the recent advanced learned descriptors. For fairness comparison, no ratio test and only cross check (mutual test) is applied for all descriptors.

Evaluation results are shown in Tab. 2. Compared with DSP-SIFT, GeoDesc performs comparably regarding image quality changes (compression, blur), while notably better for illumination and geometric changes (rotation, scale, viewpoint). On the other hand, GeoDesc delivers significant improvements on L2-Net

Table 2: Evaluation results on pair-wise image matching on benchmark by Heinly *et al.* [9] with respect to different types of image changes

	Matching Score in %					Recall in %				
	SIFT	DSP-SIFT	L2-Net	HardNet	GeoDesc	SIFT	DSP-SIFT	L2-Net	HardNet	GeoDesc
<i>JPEG</i>	31.9	35.1	25.7	27.0	34.7	60.7	66.9	49.0	51.5	66.1
<i>Blur</i>	12.4	14.3	9.1	11.3	14.4	41.0	47.3	30.1	37.4	47.7
<i>Exposure</i>	32.9	34.8	33.9	34.9	36.3	78.2	82.6	80.4	82.8	86.4
<i>Day-Night</i>	5.6	5.7	6.8	7.4	7.5	29.2	29.7	35.6	38.9	39.6
<i>Scale</i>	35.8	34.7	32.6	34.8	37.8	81.2	78.8	73.6	79.0	85.8
<i>Rotation</i>	56.3	49.1	55.9	57.4	59.8	82.4	71.8	81.9	84.0	87.6
<i>Scale-rotation</i>	12.6	12.0	10.7	12.1	14.3	29.6	28.1	25.0	28.5	33.7
<i>Planar</i>	23.8	24.8	25.6	27.4	29.1	48.2	50.4	51.9	55.6	59.1

and HardNet and particularly narrows the gap in terms of photometric changes, which makes GeoDesc applicable to different scenarios in real applications.

4.3 ETH local features benchmark

The ETH local features benchmark [25] focuses on image-based 3D reconstruction tasks. We compare GeoDesc with SIFT, DSP-SIFT and L2-Net, and follow the same protocols in [25] to quantify the SfM quality, including the number of registered images (*# Registered*), reconstructed sparse points (*# Sparse Points*), image observations (*# Observations*), mean track length (*Track Length*) and mean reprojection error (*Reproj. Error*). For fairness comparison, we apply no distance ratio test for descriptor matching and extract features at the same keypoints as in [25].

As observed in Tab. 3, first, GeoDesc performs best on *# Registered*, which is generally considered as the most important SfM metric that directly quantifies the reconstruction completeness. Second, GeoDesc achieves best results on *# Sparse Points* and *# Observations*, which indicates the superior matching quality in the early step of SfM. However, GeoDesc fails to get best statistics about *Track Length* and *Reproj. Error* as GeoDesc computes the two metrics on significantly larger *# Sparse Points*. In terms of datasets whose scale is small and have similar track number (*Fountain*, *Herzjesu*), GeoDesc gives the longest *Track Length*.

To sum up, GeoDesc surpasses both the previous best-performing DSP-SIFT and recent advanced L2-Net by a notable margin. In addition, it is noted that L2-Net also shows consistent improvements over DSP-SIFT, which demonstrates the power of taking structured loss for learned descriptors.

4.4 Challenging 3D reconstructions

To further demonstrate the effect of the proposed descriptor in a context of 3D reconstruction, we showcase selective image sets whose reconstructions fail or are in low quality with a typical SIFT-based 3D reconstruction pipeline but get significantly improved by integrating GeoDesc.

Table 3: Evaluation results on ETH local features benchmark [25] for SfM tasks

		# Images	# Registered	# Sparse Points	# Observations	Track Length	Reproj. Error
Fountain	<i>SIFT</i>	11	11	10,004	44K	4.49	0.30px
	<i>DSP-SIFT</i>		11	14,785	71K	4.80	0.41px
	<i>L2-Net</i>		11	16,119	78K	4.86	0.43px
	<i>GeoDesc</i>		11	16,687	83K	4.99	0.46px
Herzjesu	<i>SIFT</i>	8	8	4,916	19K	4.00	0.32px
	<i>DSP-SIFT</i>		8	7,760	32K	4.19	0.45px
	<i>L2-Net</i>		8	8,473	36K	4.27	0.47px
	<i>GeoDesc</i>		8	8,720	38K	4.34	0.55px
South Building	<i>SIFT</i>	128	128	62,780	353K	5.64	0.42px
	<i>DSP-SIFT</i>		128	110,394	664K	6.02	0.57px
	<i>L2-Net</i>		128	155,780	798K	5.13	0.58px
	<i>GeoDesc</i>		128	170,306	887K	5.21	0.64px
Madrid Metropolis	<i>SIFT</i>	1,344	440	62,729	416K	6.64	0.53px
	<i>DSP-SIFT</i>		476	107,028	681K	6.36	0.64px
	<i>L2-Net</i>		692	254,142	1,067K	4.20	0.69px
	<i>GeoDesc</i>		809	306,976	1,200K	3.91	0.66px
Gendarmenmarkt	<i>SIFT</i>	1,463	950	169,900	1,010K	5.95	0.64px
	<i>DSP-SIFT</i>		975	321,846	1,732K	5.38	0.74px
	<i>L2-Net</i>		1,168	667,392	2,611K	3.91	0.73px
	<i>GeoDesc</i>		1,208	779,814	2,903K	3.72	0.74px
Tower of London	<i>SIFT</i>	1,576	702	142,746	963K	6.75	0.53px
	<i>DSP-SIFT</i>		755	236,598	1,761K	7.44	0.64px
	<i>L2-Net</i>		1,049	558,673	2,617K	4.68	0.67px
	<i>GeoDesc</i>		1,081	622,076	2,852K	4.58	0.69px
Alamo	<i>SIFT</i>	2,915	743	120,713	1,384K	11.47	0.54px
	<i>DSP-SIFT</i>		754	144,341	1,815K	12.58	0.66px
	<i>L2-Net</i>		882	318,787	2,932K	9.17	0.76px
	<i>GeoDesc</i>		893	353,329	3,159K	8.94	0.84px
Roman Forum	<i>SIFT</i>	2,364	1,407	242,192	1,805K	7.45	0.61px
	<i>DSP-SIFT</i>		1,583	372,573	2,879K	7.73	0.71px
	<i>L2-Net</i>		1,537	708,794	4,530K	6.39	0.69px
	<i>GeoDesc</i>		1,566	770,363	5,051K	6.56	0.73px
Cornell	<i>SIFT</i>	6,514	4,999	1,010,544	6,317K	6.25	0.53px
	<i>DSP-SIFT</i>		4,946	1,177,916	7,233K	6.14	0.67px
	<i>L2-Net</i>		5,557	2,706,215	15,710K	5.81	0.72px
	<i>GeoDesc</i>		5,823	3,076,476	17,550K	5.70	0.96px

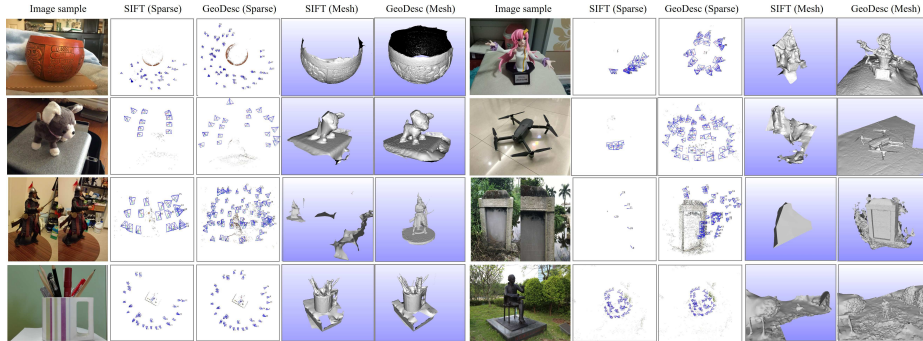


Fig. 5: Testing cases of challenging image sets, where a traditional SIFT-based reconstruction pipeline fails to apply but GeoDesc delivers significant improvement.

From examples shown in Fig. 5, it is clear to see the benefit of deploying GeoDesc in a reconstruction pipeline. First, by robust matching resistant to photometric and geometric changes, a complete sparse reconstruction registered

with more cameras can be obtained. Second, due to more accurate camera pose estimation, the final fined mesh reconstruction is then derived.

5 Practical Guidelines

In this section, we discuss several practical guidelines to complement the performance evaluation and provide insights towards real applications. Following experiments are conducted with 231 extra high-resolution image pairs, whose keypoints are downsampled to $\sim 10k$ per image. We use a single NVIDIA GTX 1080 GPU with TensorFlow [1], and forward each batch with 256 patches.

5.1 Ratio criterion

The ratio criterion [14] compares the distance between the first and the second nearest neighbor, and establishes a match if the former is smaller than the latter to some ratio. For SfM tasks, the ratio criterion improves overall matching quality, RANSAC efficiency, and seeds robust initialization. Despite those benefits, the ratio criterion has received little attention, or even been considered inapplicable to learned descriptors in previous studies [25]. Here, we propose a general method to determine the ratio that well cooperates with existing SfM pipelines.

The general idea is simple: the new ratio should function similarly as SIFT’s, as most SfM pipelines are parameterized for SIFT. To quantify the effect of the ratio criterion, we use the metric $Precision = \#Inlier\ Matches / \#Putative\ matches$, and determine the ratio that achieves similar $Precision$ as SIFT’s. As an example in Fig. 6, we compute the $Precision$ of SIFT and GeoDesc on our experimental dataset, and find the best ratio for GeoDesc is 0.89 at which it gives similar $Precision$ (0.70) as SIFT (0.69). This ratio is applied to experiments in Sec. 4.4 and shows robust results and compatibility in the practical SfM pipeline.

5.2 Compactness study

A compact feature representation generally indicates better performance with respective to discriminativeness and scalability. To quantify the compactness, we rely on the intermediate result in Principal Component Analysis (PCA). First, we compute the explained variance v_i which is stored in increasing order for each feature dimension indexed by i . Then we estimate the compact dimensionality (denoted as Compact-Dim) by finding the minimal k that satisfies $\sum_i^k v_i / \sum_i^D v_i \geq t$, where t is a given threshold and D is the original feature dimensionality. In this experiment, we set $t = 0.9$, so that the Compact-Dim can be interpreted as the minimal dimensionality required to convey more than 90% information of the original feature. Obviously, larger Compact-Dim indicates less redundancy, namely greater compactness.

As a result, the Compact-Dim estimated on 4 millions feature vectors for SIFT, DSP-SIFT, L2-Net and GeoDesc is 56, 63, 75 and 100, respectively. The ranking of Compact-Dim effectively responds to previous performance evaluations, where descriptors with larger Compact-Dim yield better results.

5.3 Scalability study

Computational cost. As evaluated in [3, 25], the efficiency of learned descriptors is on par with traditional descriptors such as CPU-based SIFT. Here, we further compare with GPU-based SIFT [32] to provide insights about practicability. We evaluate the running time in three steps. First, keypoint detection and canonical orientation estimation by SIFT-GPU. Next, patches cropping by Equ. 1. Finally, feature inference of image patches. The computational cost and memory demand are shown in Tab. 4, indicating that with GPU support, not surprisingly, SIFT (0.20s) is still faster than the learned descriptor (0.31s), with a narrow gap due to the parallel implementation. For applications heavily relying on matching quality (e.g., 3D reconstruction), the proposed descriptor achieves a good trade-off to replace SIFT.

Quantization. To conserve disk space, I/O and memory, we linearly map feature vectors of GeoDesc from $[-1, 1]$ to $[0, 255]$ and round each element to unsigned-char value. The quantization does not affect the performance as evaluated on Hatches benchmark.

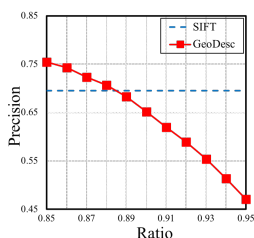


Fig. 6: Determine the ratio criterion of GeoDesc so that it has the same *Precision* as SIFT (at 0.89)

Table 4: Computational cost and memory demand of feature extraction of GeoDesc in three steps: SIFT-GPU extraction, patch cropping and feature inference. The total time cost is evaluated with three steps implemented in a parallel fashion

	SIFT Crop. Infer. Total			
<i>Device</i>	GPU	CPU	GPU	-
<i>Memory (GB)</i>	3.3	2.7	0.3	-
<i>Time (s)</i>	0.20	0.28	0.31	0.31

6 Conclusions

In contrast to prior work, we have addressed the advantages of integrating geometry constraints for descriptor learning, which benefits the learning process in terms of ground truth data generation, data sampling and loss computation. Also, we have discussed several guidelines, in particular, the *ratio criterion*, towards practical portability. Finally, we have demonstrated the superior performance and generalization ability of the proposed descriptor, GeoDesc, on three benchmark datasets in different scenarios. We have further shown the significant improvement of GeoDesc on challenging reconstructions, and the good trade-off between efficiency and accuracy to deploy GeoDesc in real applications.

Acknowledgment. This work is supported by T22-603/15N, Hong Kong ITC PSKL12EG02 and the Special Project of International Scientific and Technological Cooperation in Guangzhou Development District (No. 2017GH24).

References

1. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., et al.: Tensorflow: Large-scale machine learning on heterogeneous distributed systems. arXiv (2016)
2. Arandjelovic, R., Zisserman, A.: Three Things Everyone Should Know to Improve Object Retrieval. CVPR (2012)
3. Balntas, V., Lenc, K., Vedaldi, A., Mikolajczyk, K.: HPatches: A Benchmark and Evaluation of Handcrafted and Learned Local Descriptors. CVPR (2017)
4. Balntas, V., Riba, E., Ponsa, D., Mikolajczyk, K.: Learning Local Feature Descriptors with Triplets and Shallow Convolutional Neural Networks. BMVC (2016)
5. Balntas, V., Johns, E., Tang, L., Mikolajczyk, K.: PN-Net: Conjoined Triple Deep Network for Learning Local Image Descriptors. arXiv (2016)
6. Brown, M.A., Hua, G., Winder, S.A.J.: Discriminative Learning of Local Image Descriptors. PAMI (2011)
7. G, V.K.B., Carneiro, G., Reid, I.: Learning Local Image Descriptors with Deep Siamese and Triplet Convolutional Networks by Minimizing Global Loss Functions. CVPR (2016)
8. Han, X., Leung, T., Jia, Y., Sukthankar, R., Berg, A.C.: MatchNet - Unifying Feature and Metric Learning for Patch-based Matching. CVPR (2015)
9. Heinly, Jared, Dunn, Enrique, Frahm, Jan-Michael: Comparative Evaluation of Binary Features. ECCV (2012)
10. Kaplan, A., Avraham, T., Lindenbaum, M.: Interpreting the ratio criterion for matching sift descriptors. ECCV (2016)
11. Labatut, P., Pons, J.P., Keriven, R.: Efficient multi-view reconstruction of large-scale scenes using interest points, delaunay triangulation and graph cuts. ICCV (2007)
12. Li, S., Yuan, L., Sun, J., Quan, L.: Dual-feature warping-based motion model estimation. ICCV (2015)
13. Lin, J., Morere, O., Chandrasekhar, V., Veillard, A., Goh, H.: Deephash: Getting regularization, depth and fine-tuning right. arXiv (2015)
14. Lowe, D.G.: Distinctive Image Features from Scale-Invariant Keypoints. IJCV (2004)
15. Markuš, N., Pandžić, I.S., Ahlberg, J.: Learning local descriptors by optimizing the keypoint-correspondence criterion. ICPR (2016)
16. Mikolajczyk, K., Schmid, C.: A Performance Evaluation of Local Descriptors. PAMI (2005)
17. Mishchuk, A., Mishkin, D., Radenovic, F.: Working Hard to Know Your Neighbor's Margins: Local Descriptor Learning Loss. NIPS (2017)
18. Mishkin, D., Matas, J., Perdoch, M., Lenc, K.: Wxbs: Wide baseline stereo generalizations. BMVC (2015)
19. Mitra, R., Doiphode, N., Gautam, U., Narayan, S., Ahmed, S., Chandran, S., Jain, A.: A large dataset for improving patch matching. arXiv, year=2018
20. Morel, J.M., Yu, G.: Asift: A new framework for fully affine invariant image comparison. SIAM journal on imaging sciences
21. Movshovitz-Attias, Y., Toshev, A., Leung, T.K., Ioffe, S., Singh, S.: No Fuss Distance Metric Learning using Proxies. ICCV (2017)
22. Philbin, J., Chum, O., Isard, M., Sivic, J., Zisserman, A.: Object retrieval with large vocabularies and fast spatial matching. CVPR (2007)

23. Radenovic, F., Tolias, G., Chum, O.: CNN Image Retrieval Learns from BoW - Unsupervised Fine-Tuning with Hard Examples. ECCV (2016)
24. Schnberger, J.L., Frahm, J.M.: Structure-from-motion revisited. CVPR (2016)
25. Schönberger, J.L., Hardmeier, H., Sattler, T., Pollefeys, M.: Comparative evaluation of hand-crafted and learned local features. CVPR (2017)
26. Shen, T., Zhu, S., Fang, T., Zhang, R., Quan, L.: Graph-based consistent matching for structure-from-motion. European Conference on Computer Vision (2016)
27. Simo-Serra, E., Trulls, E., Ferraz, L., Kokkinos, I., Fua, P., Moreno-Noguer, F.: Discriminative Learning of Deep Convolutional Feature Point Descriptors. CVPR (2015)
28. Sohn, K.: Improved Deep Metric Learning with Multi-class N-pair Loss Objective. NIPS (2016)
29. Tian, BFY, Wu, F: L2-Net: Deep Learning of Discriminative Patch Descriptor in Euclidean Space. CVPR (2017)
30. Wilson, K., Snavely, N.: Robust Global Translations with 1DSfM. ECCV (2014)
31. Winder, S., Hua, G., Brown, M.: Picking the best daisy. CVPR (2009)
32. Wu, C.: Siftgpu: A gpu implementation of sift. <http://cs.unc.edu/~ccwu/siftgpu> (2007)
33. Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. arXiv (2018)
34. Yi, K.M., Trulls, E., Lepetit, V., Fua, P.: LIFT - Learned Invariant Feature Transform. ECCV (2016)
35. Yi, K.M., Verdie, Y., Fua, P., Lepetit, V.: Learning to Assign Orientations to Feature Points. CVPR (2015)
36. Zagoruyko, S., Komodakis, N.: Learning to compare image patches via convolutional neural networks. CVPR (2015)
37. Zhang, R., Li, S., Fang, T., Zhu, S., Quan, L.: Joint Camera Clustering and Surface Segmentation for Large-Scale Multi-view Stereo. ICCV (2015)
38. Zhang, R., Zhu, S., Fang, T., Quan, L.: Distributed very large scale bundle adjustment by global camera consensus. ICCV (2017)
39. Zhang, X., Yu, F.X., Kumar, S., Chang, S.F.: Learning spread-out local feature descriptors. CVPR (2017)
40. Zhou, L., Zhu, S., Shen, T., Wang, J., Fang, T., Quan, L.: Progressive large scale-invariant image matching in scale space. ICCV (2017)
41. Zhu, S., Fang, T., Xiao, J., Quan, L.: Local readjustment for high-resolution 3d reconstruction. CVPR (2014)
42. Zhu, S., Zhang, R., Zhou, L., Shen, T., Fang, T., Tan, P., Quan, L.: Very large-scale global sfm by distributed motion averaging. CVPR (2018)