

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015

# Local Feature Augmentation with Cross-Modality Context for Geometric Matching

Anonymous CVPR submission

Paper ID 325

## Abstract

Most existing studies on learning local features focus on appropriate descriptions of individual image patches, whereas neglecting spatial location relationship of descriptors in the image. In this paper, we bridge isolated patches by their keypoint coordinates, and go beyond the local detail by introducing context awareness to augment raw local feature. Specifically, we propose a unified learning framework that leverages cross-modality contextual information, consisting of (i) visual context from high-level image understandings, and (ii) geometric context from 2D keypoint distribution. Moreover, we propose an effective technique to alleviate the scale affects of N-pair loss. The proposed augmentation scheme costs only 6% extra time compared with raw local feature description, but improves remarkably on several large-scale benchmarks with diversified scenes, which demonstrates both strong generalization and practicality in geometric matching applications.

## 1. Introduction

Designing powerful local feature descriptor is a fundamental problem in applications such as wide-baseline matching [24], image retrieval [27], and structure-from-motion (SfM) [33]. Despite of notable achievement by recent advance, the performance of state-of-the-art learned descriptors is observed to be somewhat saturated on standard benchmarks. As shown in Fig. 1a, due to visual repetitiveness, the matching process often finds nearest neighbors that are hardly distinguishable from true matches unless validated by geometry, e.g., homography. Essentially, such visual ambiguity may not be easily resolved with only local information. In this spirit, we propose to enhance the feature description with prior knowledge, which we refer to as introducing *context awareness* to augment local feature.

As a common practice, a multi-scale-like architecture helps to capture *visual context* of different levels, termed as aggregating domains of multiple sizes by DSP-SIFT [8] and

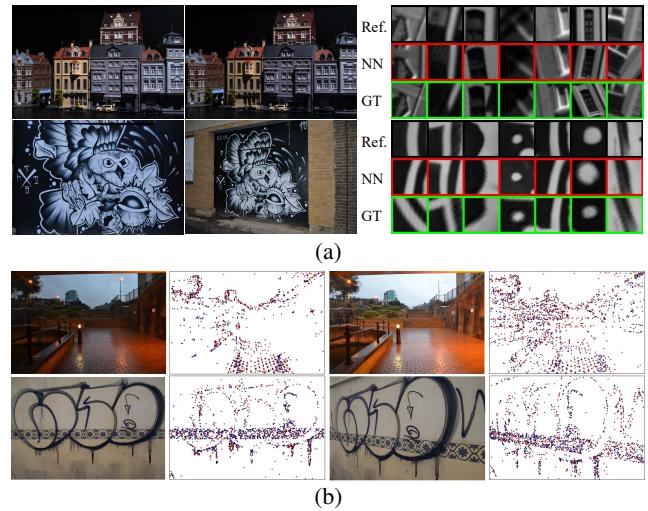


Figure 1: (a) Saturated results on standard benchmark [2] by the recent advance [23]. The search of nearest neighbors (NN) returns false matches though visually similar to groundtruth (GT), indicating the limitation of relying on only local visual information. (b) 2D keypoints distribute structurally, on which we human beings are capable of establishing coarse matches even without color information.

adopted by recent learned advances [45, 19, 38], of which the key challenge lies with the proper accuracy-efficiency trade-off and selection of domain size in order not to deviate the teachings from scale-space theory [22] or sample out-of-boundary pixels. On the other hand, visual context has essentially incorporated high-level image understandings, e.g., image retrieval [30]. It is an open problem whether such off-the-shelf context representation can be effectively leveraged for enhancing local feature description.

In addition to visual information, it would be interesting to exploit context in other modality. In particular, as shown in Fig. 1b, since keypoint is principally designed to be repeatable in the same underlying scene, its distribution thus has revealed comprehensive structure that allows we human beings to establish coarse matches even without color information. In essence, keypoints can be bridged by their coordinates, from which we can explore *geometric context*

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

108 to help to alleviate the ambiguity of local visual feature.  
109  
110 Thus far, we have discussed two context candidates, re-  
111 ferred to as *visual context* and *geometric context* that incor-  
112 porate high-level visual understandings and geometric cues  
113 of 2D keypoint distribution, respectively. Instead of learn-  
114 ing a completely new descriptor, in the present work, we  
115 target to flexibly leverage the above context awareness to  
116 augment off-the-shelf raw local features, in which process  
117 we consider the key challenges threefold:

- 118 • Proper integration of local geometric feature and high-  
119 level visual understandings. As keypoint description re-  
120 quires sub-pixel accuracy, the integration is not supposed  
121 to mix up raw representation of local details.
- 122 • Instability of 2D keypoints. Due to image appearance  
123 changes, keypoint distribution often suffers from substan-  
124 tial variations of sparsity, non-uniformity or perspective,  
125 which raises a challenge on acquiring strong invariance  
126 property of the context encoder.
- 127 • Effective learning scheme. Input signals and features in  
128 different modalities are supposed to be efficiently pro-  
129 cessed and aggregated in a unified framework.

131 Finally, regarding practicability, the augmentation is not  
132 supposed to be too complex to introduce excessive com-  
133 putational cost, as the local feature description is often re-  
134 garded as part of preprocessing in practical pipelines.

135 Although contextual information has been widely ap-  
136 plied in computer vision tasks, the challenges faced by lo-  
137 cal feature learning are substantially different, posing many  
138 non-trivial technical and systematic issues to overcome.  
139 In this paper, we propose a unified augmentation scheme  
140 that leverages and aggregates cross-modality information,  
141 of which the contributions are summarized threefold: 1) a  
142 novel *visual context encoder* that integrates high-level vi-  
143 sual understandings from regional image representation, 2)  
144 a novel *geometric context encoder* that exploits geometric  
145 cues from raw 2D keypoint distribution. 3) A novel learn-  
146 ing scheme and a new loss to boost the learning effective-  
147 ness and mitigate the scale affects in N-pair loss. To our  
148 best knowledge, it is the first work that emphasizes the im-  
149 portance of context awareness in 2D local feature learning.

150 The proposed augmentation scheme is extensively eval-  
151 uated and achieves state-of-the-art results on several  
152 large-scale benchmarks, including patch-level homography  
153 dataset, image-level wild outdoor/indoor scenes and 3D re-  
154 construction image sets, with only 6% extra time cost com-  
155 pared with raw local description, demonstrating both strong  
156 generalization ability and practicability.

## 157 2. Related Work

158 **Learned local features.** Initially, local descriptors are  
159 jointly learned with a new metric [9, 45], which is later

162 simplified as direct comparison in Euclidean space [35,  
163 43, 3, 19, 1]. More recently, attention is drawn on effi-  
164 cient training data sampling [38, 25, 11], effective regulari-  
165 zations [38, 46], and geometric shape estimation of input  
166 patches [26, 7]. However, most of above methods take *indi-  
167 vidual* image patches as input, whereas in the present work,  
168 we aim to make use of contextual cues beyond the local de-  
169 tail and take advantage of features in multiple modalities.

170 **Context awareness.** Although widely introduced in many  
171 tasks, context awareness has received little attention in  
172 learning 2D local descriptors. In terms of visual context,  
173 the central-surround (CS) structure [45, 19, 38] leverages  
174 multi-scale information to boost the performance, whereas  
175 sacrificing computational efficiency due to doubled extrac-  
176 tion time and doubled feature dimensionality. Regarding se-  
177 mantics, one practice [18] designs a new comparison metric  
178 and describes features directly from histogram of semantic  
179 labels. On the other hand, a family of studies has focused  
180 on finding semantic correspondences [40, 31] across *diff-  
181 ferent* objects of the same category, of which the purpose is  
182 substantially different from our case. Beside of visual in-  
183 formation, a recent advance [44] explores to encode motion  
184 context for identifying outlier from image correspondences,  
185 i.e., 4-d coordinate pairs, whereas we aim to exploit geo-  
186 metric context from *single* image without any reference.  
187 Overall, encoding proper context is non-trivial and still un-  
188 clear in 2D local feature learning.

189 **Point feature learning.** In the present work, one of our  
190 goals is to explore geometric features from keypoint dis-  
191 tribution, we thus resort to PointNet [28] and its vari-  
192 ants [29, 5, 44] to consume unordered points. Although  
193 great success has been shown in learning tasks on 3D points,  
194 there are only few studies exploiting the potential outcome  
195 of 2D keypoints. In essence, keypoint structure is not intui-  
196 tively meaningful and robust, as being highly dependent on  
197 the performance of interest point detectors and strongly af-  
198 fected by image variations. However, in descriptor learning,  
199 we consider the keypoint location as an important cue that  
200 bridge each individual local feature, constructing a unified  
201 instance that reveals high-level contextual information.

202 **Loss formulation.** Recent local descriptors are often  
203 evolved with advanced variants of N-pair losses. Initially,  
204 L2-Net [38] adopts a log-likelihood formulation, which is  
205 later extended by HardNet [25] with hard negative triplet  
206 margin loss. Furthermore, GeoDesc [23] applies an adap-  
207 tive margin value to improve the convergence in terms of  
208 different data sampling strategies, where AffNet [7] ap-  
209 proaches the same issue by fixing the distance of hardest  
210 negative sample in the training. Meanwhile, on the other  
211 hand, DOAP [11] extends the N-pair loss to a list-wise rank-  
212 ing loss, while [17] points out and studies the scale affects in  
213 N-pair losses. Principally, the loss is supposed to encourage  
214 similar patches to be close while dissimilar ones to be dis-

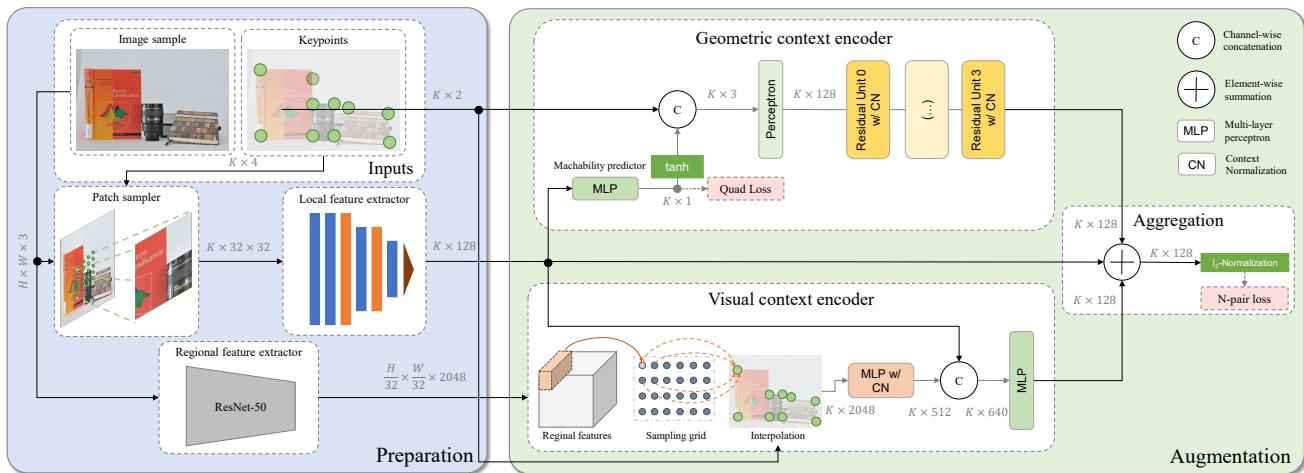


Figure 2: The network architecture of local feature augmentation.

tant in the descriptor space. In this spirit, we aim to further resolve scale affects in an self-adaptive manner, without the need of complex heuristics or manual tuning.

### 3. Local Feature Augmentation

**Overview.** As illustrated in Fig. 2, the proposed framework consists of two main modules: *preparation* (left) and *augmentation* (right). The *preparation* module provides input signals in different modalities (raw local feature, high-level visual feature and keypoint location), which is then fed to the *augmentation* module and aggregated into compact feature descriptions. At test time, the augmentation needs to be performed once per image, resulting in  $K$  feature vectors for  $K$  respective keypoints.

#### 3.1. Preparation

**Patch sampler.** This module takes images and their keypoints as input, producing  $32 \times 32$  gray-scale patches. Akin to [43, 23], the patch is sampled as applying similarity transformation parameterized by keypoints (coordinates, orientation and scale) from the SIFT detector, implemented by a spatial transformer [16]. The patch has the same size with the supporting region of SIFT descriptor.

**Local feature extractor.** This module takes image patches as input, producing 128-d feature descriptions as output. We borrow the lightweight 7-layer convolutional networks, as used in several recent works [38, 25, 23].

**Regional feature extractor.** In contrast to aggregating features of different domain sizes [45, 19, 38], in the present work, we fix the sampling scale of patches, and exploit contextual cues by inspiration of well-studied regional representation [39, 30]. Without the loss of generality, we reuse features from an off-the-shelf deep image retrieval model of ResNet-50 [12], which is pretrained as [30]. As in [39],

feature maps are extracted from the last bottleneck block, across which each response is regarded as a regional feature vector effectively corresponding to a particular region in the image. As a result, we derive regional features of  $\frac{H}{32} \times \frac{W}{32} \times 2048$ , where  $H$  and  $W$  denote original image height and width. The aggregation of regional and local features will be later discussed in Sec. 3.3.

#### 3.2. Geometric context encoder

This module takes  $K$  unordered points as input, and outputs  $K$  corresponding feature vectors. Each input point is represented as 2D coordinate of the keypoint, and can be associated with other attributes.

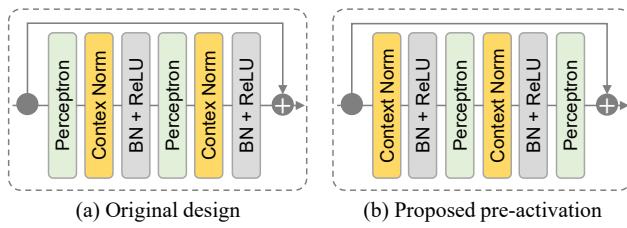
**2D point processing.** At first glance, 2D keypoints are inappropriate to serve as robust contextual cues, as the presence is heavily dependent on image appearance and thus affected by various image variations. As a result, keypoints depicting the same scene may suffer from significant density or structure changes, as examples shown in Fig. 1b. Hence, acquiring strong invariance property is the key challenge when designing the context encoder.

Initially, we attempt to approach the goal by PointNet [28] and its variants [29, 5]. Although having shown great success on processing 3D point clouds, the family of PointNet methods fails to achieve consistent improvement in terms of processing 2D points. Instead, we resort to [44], which originally focuses on outlier rejection in image matching and consumes *putative matches* (4-d coordinate pairs) of image pairs as their network input. In particular, we aim to further explore the potential of context normalization (CN) proposed in [44], and extend its use of processing 2D points in *single* image.

Formally, context normalization is a non-parametric operation that simply normalizes feature maps according to their distribution, expressed as  $\hat{o}_i^l = \frac{(o_i^l - \mu^l)}{\sigma}$ , where  $o_i^l$  is

324 the output of  $i$ -th point in layer  $l$ , and  $\mu^l, \sigma^l$  are mean and  
 325 standard deviation of the output in layer  $l$ . To equip the operation,  
 326 we follow the residual architecture in [44], where  
 327 each residual unit is built with two perceptrons followed by  
 328 context and batch normalization, as illustrated in Fig. 3a.  
 329

330 However, the above construction leads to a *non-negative*  
 331 output from the residual function that impacts the repre-  
 332 sentational ability as investigated in [13] and witnessed in  
 333 our experiments. Following the teachings of [13], we pro-  
 334 pose to re-arrange the operations in residual unit by adopt-  
 335 ing the *pre-activation* construction and showing its appli-  
 336 cability with context normalization as in Fig. 3b. We then  
 337 construct four such units as the final geometric context  
 338 encoder, as shown in Fig. 2. We will show that this simple  
 339 revision plays an important role to ease the optimization.



340 Figure 3: Different designs of residual unit with context  
 341 normalization, where the proposed construction improves  
 342 by a considerable margin than its original counterpart.  
 343

344 **Matchability predictor.** In 3D point cloud processing,  
 345 low-level color and normal [28] information, or more com-  
 346 plexly, geometric attributes such as point pair features [5]  
 347 are often adopted to enhance the representation. Simi-  
 348 larly, associating 2D coordinate input with other mean-  
 349 ingful attributes would be promising to boost the performance.  
 350 However, due to the substantial variations, e.g. perspective  
 351 change, it is non-trivial to define appropriate intermediate  
 352 attributes on 2D points.

353 Although this issue has been merely discussed, we draw  
 354 inspiration from [10], which poses a problem named *match-  
 355 ability prediction* and targets to decide *whether a keypoint  
 356 descriptor is matchable before the matching stage*. In prac-  
 357 tice, the learned matchability can be used to guide the key-  
 358 point sampling and accelerate the matching without sacri-  
 359 ficing accuracy. In contrast to criteria such as cornerness  
 360 or edgeness of a keypoint, matchability implies high-level  
 361 prior knowledge derived from data. Dependent on individ-  
 362 ual keypoint, matchability can be naturally associated as a  
 363 representative attribute in addition to 2D coordinate.

364 Previously in [10], the matchability prediction is casted  
 365 as a binary classification problem, taking individual de-  
 366 scriptor as input and inferred by a random forest. In the  
 367 present work, we approach this problem with deep learn-  
 368 ing techniques, and apply a more strict constraint requir-  
 369 ing consistent prediction between images. Similar to [32],  
 370 we resort to an unsupervised learning scheme that aims

371 to rank points. Specifically, given  $K$  corresponding key-  
 372 point pairs  $(p_1^n, p_2^n)$ ,  $n \in [1, K]$  from an image pair, we  
 373 first extract local features  $(f_1^n, f_2^n)$  of each keypoint, then  
 374 construct *feature quadruples* as  $(f_1^i, f_1^j, f_2^i, f_2^j)$ , satisfying  
 375  $i, j \in [1, K], i \neq j$  and holding that:

$$\begin{cases} H(f_1^i) > H(f_1^j) & \& H(f_2^i) > H(f_2^j) \\ & \text{or} & \\ H(f_1^i) < H(f_1^j) & \& H(f_2^i) < H(f_2^j) \end{cases}, \quad (1)$$

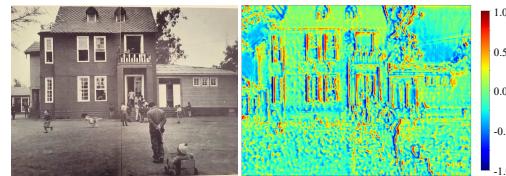
376 where  $H(\cdot)$  absorbs the raw local feature into a single real-  
 377 valued matchability, implemented as standard multi-layer  
 378 perceptrons (MLPs). Here, Cond. 1 aims to preserve a rank-  
 379 ing of each keypoint, hence improves the repeatability of  
 380 prediction. The condition can be re-written as:

$$R(f_1^i, f_1^j, f_2^i, f_2^j) = (H(f_1^i) - H(f_1^j))(H(f_2^i) - H(f_2^j)) > 0, \quad (2)$$

381 the final objective can be obtained with hinge loss:

$$\mathcal{L}_{quad} = \frac{1}{K(K-1)} \sum_{i,j,i \neq j} \max(0, 1 - R(f_1^i, f_1^j, f_2^i, f_2^j)). \quad (3)$$

382 In the proposed framework, the matchability is learned  
 383 as an auxiliary task, which is then activated by a tanh and  
 384 associated with keypoint coordinate before fed into the en-  
 385 coder, as in Fig. 2. Beside of Eq. 3, the gradient from final  
 386 augmented feature will flow through the matchability pre-  
 387 dictor, allowing a joint optimization of the entire encoder.



388 Figure 4: Visualization of matchability responding to the  
 389 entire image (best viewed in color).

### 3.3. Visual context encoder

390 This module takes regional features of  $\frac{H}{32} \times \frac{W}{32} \times 2048$   
 391 in Sec. 3.1,  $K$  raw local features and their location as input,  
 392 producing  $K$  corresponding feature vectors. As widely-  
 393 adopted in other tasks, the main purpose of this module is to  
 394 integrate visual information in different levels, e.g., seman-  
 395 tics. In our context, the key issue is to handle the regional  
 396 features and keypoints of different numbers. One option as  
 397 in [5] is to concatenate global representation of entire image  
 398 on raw local features, where the global feature in our frame-  
 399 work can be derived by applying Maximum Activations of  
 400 Convolutions (MAC) aggregation [30] which simply max-  
 401 pools over all dimensions per feature map. However, such  
 402 compact representation does not suffice to strengthen the  
 403 local description, due to the lack of distinctiveness.

To better preserve the regional distinction, we associate regional features to a regular sampling grid on the image, then interpolate  $\frac{H}{32} \times \frac{W}{32}$  grid points at coordinates of the  $K$  keypoints. For interpolation, we use the inverse distance weighted average based on  $k$  nearest neighbors (in default we use  $k = 4$ ), formulated as:

$$\mathbf{f}(\hat{p}_i) = \frac{\sum_{j=1}^k w(p_j) \mathbf{f}(p_j)}{\sum_{j=1}^k w(p_j)}, \text{ and } w(p_j) = \frac{1}{d(\hat{p}_i, p_j)}, \quad (4)$$

where  $\mathbf{f}(\cdot)$  is the regional feature located at a certain grid point.  $\hat{p}_i, i \in [1, N]$  indicates interpolated point, while  $p_j, j \in [1, \frac{H}{32} \times \frac{W}{32}]$  indicates original grid point. Next, the dimensionality is reduced by applying point-wise MLPs, where we also insert CN after each perceptron in order to capture global context. Finally, raw local features are concatenated and further mapped by MLPs, forming the final 128-d features. The above process is illustrated in Fig. 2.

### 3.4. Feature aggregation

To this end, we have obtained two types of contextual feature, and aim to properly aggregate with raw local features. Similar to the CS structure, one option is to concatenate them together and forms, in our case, 384-d ( $128 \times 3$ ) feature for each keypoint. However, the increased dimensionality will introduce excessive computational cost in the matching stage of  $\mathcal{O}(n^2)$  complexity. Hence, as shown in Tab. 2, we propose to combine different feature streams into a single vector by summing and L2-normalizing them in the end, i.e., without the change of feature dimensionality. Beside of the simplicity, such strategy allows flexible use of the augmentation. For example, in situations where regional features are not available, one may aggregate with only geometric context without the need of retraining the model.

## 4. Learning Scheme

### 4.1. N-pair loss with softmax temperature

N-pair losses have been primarily used by recent works. Empirically, the subtractive hinge loss [25, 23, 7] has reported better performance, of which the main idea is to push similar samples away from dissimilar ones to a certain *margin* in the descriptor space. However, setting the appropriate margin is tricky, which does not always assure convergence as observed in [23, 7]. More generally, the aspects of making a good loss is studied in [17], from which guidelines are provided about tuning loss coefficients on particular dataset. In this spirit, we aim to ease the pain of parameter searching in [17], and obtain scale-aware loss that allows fast convergence despite of different data sampling strategies.

We use the log-likelihood version of N-pair loss [38] as a base, which originally does not involve any tunable parameter. Formally, given L2-normalized feature descriptors

$\mathbf{F}_1 = [\mathbf{f}_1^1 \mathbf{f}_1^2 \dots \mathbf{f}_1^N]^T, \mathbf{F}_2 = [\mathbf{f}_2^1 \mathbf{f}_2^2 \dots \mathbf{f}_2^N]^T \in \mathbb{R}^{N \times 128}$ , the distance matrix  $\mathbf{D} = [d_{ij}]_{N \times N}$  can be obtained by  $\mathbf{D} = \sqrt{2(1 - \mathbf{F}_1 \mathbf{F}_2^T)}$ . By applying, e.g., row-wise softmax, we derive the final loss as:

$$\mathcal{L}_{N\text{-pair}} = - \sum_i \log s_{ii}, \quad (5)$$

where  $[s_{ij}]_{N \times N} = \text{softmax}(2 - \mathbf{D})$ .

Noted that since input features are L2-normalized, the resulting  $d_{ij}$  is bounded between 0 and 2, which causes convergence issues due to the scale sensitivity of softmax function [15]. Similarly, we introduce a single scalar parameter  $\alpha$ , referred to as *softmax temperature*, to amend the inability of re-scaling the input. The loss now becomes:

$$[s_{ij}]_{N \times N} = \text{softmax}(\alpha(2 - \mathbf{D})), \quad (6)$$

where  $\alpha$  is initialized to 1 and regularized with the same weight decay in the network, which does not require any manual tuning or complex heuristics. In the proposed framework, the loss is computed on augmented features.

### 4.2. Learning with noisy keypoints

The training of proposed framework, apparently, needs to be conducted between image pairs instead of isolated patches, which is referred to as simulating image matching in [23]. However, the simulation in [23] is not complete, as it considers only matchable keypoints that are paired with correspondences, whereas in real situation, only a subset of keypoints is repeatable in other images. In practice, as illustrated in Fig. 5, we divide groundtruth keypoints generated from SfM as in [43, 23] into three categories: i) *Matchable*: repeatable and verified by SfM; ii) *Undiscovered*: repeatable but did not survive the SfM. iii) *Unrepeatable*: unable to be re-detected in other images.

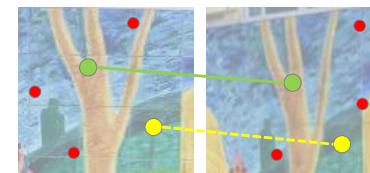


Figure 5: We divide keypoints after SfM into three categories: matchable (green), undiscovered (yellow) and unrepeatable (red), and aim to perform a complete simulation in training that incorporates all three types of keypoints.

In the present work, a complete simulation with undiscovered and unrepeatable keypoints (referred to as *noisy keypoints*) is necessary, otherwise the generalization ability cannot be acquired as the training considers only ideal setting which is inconsistent with real applications. Formally, given index sets  $C_m = \{i_1, \dots, i_{K_m}\}$  and  $C_n =$

540 { $i_1, \dots, i_{K_n}$ }, where  $K_m$  and  $K_n$  are numbers of matchable  
 541 and noisy keypoints for an image pair, the losses of Eq. 3  
 542 and Eq. 5 are now rewritten as:  
 543

$$\begin{aligned} 544 \quad \mathcal{L}'_{quad} &= \\ 545 \quad \frac{1}{K_m(K_m - 1)} \sum_{i,j \in C_m, i \neq j} &\max(0, 1 - R(\mathbf{f}_1^i, \mathbf{f}_1^j, \mathbf{f}_2^i, \mathbf{f}_2^j)), \\ 546 \quad & \end{aligned} \quad (7)$$

548 and

$$\begin{aligned} 549 \quad \mathcal{L}'_{N-pair} &= - \sum_{i \in C_m} \log s_{ii}. \\ 550 \quad & \end{aligned} \quad (8)$$

552 Subsequently, adding noisy keypoints will influence  
 553  $\mathcal{L}_{N-pair}$  as all keypoints have been cross-paired, while it  
 554 also affects the encoding of geometric context. Finally, the  
 555 total loss is obtained by:

$$\begin{aligned} 556 \quad \mathcal{L}_{total} &= \mathcal{L}'_{N-pair} + \lambda \mathcal{L}'_{quad}, \\ 557 \quad & \end{aligned} \quad (9)$$

559 where we choose  $\lambda = 1$  in the experiment.

## 560 5. Experiments

### 561 5.1. Implementation

564 **Training details.** Although the framework is end-to-end  
 565 trainable, we fix the local and regional feature extractors in  
 566 Sec. 3.1 during the training, in order to make a clear demon-  
 567 stration of the proposed augmentation scheme. We train  
 568 the networks using SGD with a base learning rate of 0.05,  
 569 weight decay of 0.0001 and momentum at 0.9. The learn-  
 570 ing rate exponentially decays by 0.1 for every 100k steps.  
 571 The batch size is set to 2, and each time 1024 keypoints are  
 572 randomly sampled including random numbers of matchable  
 573 and noisy keypoints. Input patches are standardized to have  
 574 zero mean and unit norm, while input keypoint coordinates  
 575 are normalized to  $[-1, 1]$  regarding the image size.

576 **Training dataset.** Although UBC Phototour [4] is used as a  
 577 common practice, this dataset consists of only three scenes  
 578 with limited diversity of keypoint distribution. In order to  
 579 achieve better generalization ability, we resort to large-scale  
 580 photo-tourism dataset [41, 30] as in [43, 23], and generate  
 581 groundtruth matches from SfM. We manually exclude the  
 582 data that is used in the evaluation.

583 **Data augmentation.** For image patches, we apply random  
 584 affine transformations including rotation, anisotropic scal-  
 585 ing and translation w.r.t. the detection scale. For keypoint  
 586 augmentation, we perturb the coordinate with random ho-  
 587 mography transformation as in [6].

### 589 5.2. Evaluation datasets

591 **Homography dataset.** HPatches [2] is a large-scale patch  
 592 dataset for evaluating local features regarding illumination  
 593 and viewpoint changes. As groundtruth homographies and

594 raw images are provided, HPatches can also be used to eval-  
 595 uate image matching performance, which we accordingly  
 596 refer to as HPSequences as in [20], consisting of 116 se-  
 597 quences and 580 image pairs.

598 **Wild dataset.** Similar to settings in [44], we use out-  
 599 door YFCC100M [37] (1000 pairs) and indoor SUN3D [42]  
 600 (539 pairs). The two datasets additionally introduce varia-  
 601 tions such as self-occlusions, and in particular, repetitive or  
 602 feature-poor patterns in indoor scenes, which is generally  
 603 considered challenging for sparse keypoint methods.

604 **SfM dataset.** Following [34], we evaluate on SfM dataset  
 605 such as well-known *Fountain* and *Herzjesu* [36], or land-  
 606 mark collections [41]. We integrate the proposed frame-  
 607 work into SfM pipeline, i.e., COLMAP [33], and use the  
 608 keypoints provided in [34] to compute the local features.

## 609 5.3. Evaluation protocols

610 **Patch level.** Following the same protocols of HPatches [2],  
 611 we use mean average precision (mAP) for its three subtasks  
 612 , including patch verification, matching, and retrieval.

613 **Image level.** For HPSequences, we use  $Recall = \# Correct$   
 614  $Matches / \# Correspondences$  defined in [14], to quantify  
 615 the image matching performance, where  $\# Correct matches$   
 616 are matches found by nearest neighbor searching and ver-  
 617 ified by groundtruth geometry, e.g., homography, while  $\#$   
 618  $Correspondences$  are matches that should have been iden-  
 619 tified by the given keypoint locations. Following [14], a  
 620 match point is determined to be correct if it is within 2.5  
 621 pixels from the wrapped keypoint in the reference image.  
 622 We use a standard SIFT detector to localize the keypoints,  
 623 of which the number is randomly sampled to 2048. For  
 624 YFCC100M [37] and SUN3D [42], we follow the same set-  
 625 ting in [44] and report the median number of inlier matches  
 626 after RANSAC for each dataset.

627 **Reconstruction level.** For clarity, we report metrics in [34]  
 628 that quantify the completeness of SfM, including the num-  
 629 ber of registered images ( $\# Registered$ ), sparse points ( $\#$   
 630  $Sparse Points$ ) and image observations ( $\# Observations$ ).

## 631 5.4. Ablation study

### 632 5.4.1 Design of context encoder

633 In this section, we evaluate two splits of HPSequences [2]:  
 634 *illumination* ( $i$ ) and *viewpoint* ( $v$ ), regarding different image  
 635 transformations. We report *Recall* as defined in Sec. 5.3.  
 636 If not specified, we use GeoDesc [23] as a baseline model  
 637 (*baseline (GeoDesc)*) to extract raw local features, whose  
 638 parameters are fixed during the training of augmentation.

639 **Visual context.** We compare four designs, including i) *CS*  
 640 (*256-d*): the central-surround (CS) structure [45, 19, 38] as  
 641 described in Sec. 2, which leverages visual information of  
 642 different domain sizes. ii) *w/ global feature*: the integra-  
 643 tion with global features [5], which is originally designed

648	Vsual context encoder			Geometric context encoder			Comparison with other methods		702
649	Strategy	Recall i/v		Network architecture	Recall i/v		Method	Recall i/v	703
650	baseline (GeoDesc)	59.27	71.44	baseline (GeoDesc)	59.27	71.44	SIFT [22]	47.41	53.19
651	CS (256-d) [45, 19, 38]	59.64	71.47	PointNet [28]	59.61	71.16	L2-Net [38]	47.55	54.10
652	w/ global feature [5]	58.92	71.22	w/ CN (pre) + xy	61.05	72.47	HardNet [25]	57.61	63.45
653	w/ regional feature	63.45	73.57	w/ CN (pre) + xy + raw local feature	60.61	72.69	GeoDesc [23]	59.27	71.44
654	<b>w/ regional feature + CN</b>	<b>63.79</b>	<b>73.63</b>	w/ CN (orig.) + xy + matchability	59.94	71.25	<b>multi-context</b>	<b>65.35</b>	<b>74.70</b>
655				w/ CN (pre) + xy + matchability	<b>61.52</b>	<b>72.83</b>	<b>multi-context+</b>	<b>65.76</b>	<b>75.64</b>

Table 1: Comparisons on HPSequences [2] of different designs of visual and geometric context encoder, and the performance of entire augmentation scheme. ‘i/v’ denotes two evaluations on *illumination* and *viewpoint* sequences, respectively.

for improving 3D local descriptors. iii) *w/ regional feature*: the proposed integration with interpolated regional features, and its variant iv) *w/ regional feature + CN*: with context normalization to exploit global visual information.

As shown in Tab. 1 (left columns), the CS structure [45, 19, 38] delivers only marginal improvements despite of the doubled dimensionality. On the other hand, though being effective in 3D descriptor learning, the integration with global features [5] instead harms the performance, which we attribute to the weak relevance of local geometric and global semantic features. Finally, the proposed integration with interpolated regional features clearly shows advantages as it better preserves distinctions from a smaller visual scale. To amend the loss of global awareness, we show that the performance can be further boosted by equipping context normalization to associate regional features.

**Geometric context.** We study five options: i) PointNet-like architecture, i.e., segmentation networks in [28] without the final classifier. ii) Pre-activated context normalization (CN) networks in Sec. 3.2 with 2D xy input, and its variants iii) with additional raw local feature input or iv) with matchability. We also compare the use of pre-activation of the residual unit in context normalization networks.

As presented in Tab. 1 (middle columns), though widely used in processing 3D points, PointNet [28] does not work well in our context, where the similar phenomenon is also observed in [44] when processing 2D correspondences. Besides, it is noticed that input with additional raw local feature does not help to boost the performance, which we attribute to the weak relevance between local features as extracted from different levels of scale space pyramid. Instead, the cooperation with matchability is beneficial, as matchability is more interpretable as a high-level abstraction of local feature. Finally, the pre-activation is clearly a preferable alternative than the original design in this task.

**Integration with multiple context.** Finally, we evaluate the full augmentation with both visual and geometric context (*multi-context*). As shown in Tab. 1 (right columns), the simple summation aggregation in Sec. 3.4 effectively takes advantage of both context, delivering remarkable improvements over the state of the art.

#### 5.4.2 Efficacy of softmax temperature

To make a clear demonstration on HPatches [2], we train *only* the local feature extractor with the proposed loss and adopt image matching simulation as in Sec. 4.2. We compare different losses including: i) the proposed loss and ii) its original form [38] without scale temperature, also iii) the loss in [17] with its original parameters.

	SIFT	L2-Net	HardNet	GeoDesc	w/ loss [38]	w/ loss [17]	proposed	
Verification, mAP [%]								
Easy	80.0	91.4	93.6	<b>94.0</b>	83.3	87.5	93.6	
Hard	59.2	83.7	87.6	91.8	78.9	82.1	<b>91.8</b>	
Tough	44.9	72.1	77.3	87.6	72.6	73.9	<b>88.4</b>	
Mean	61.4	82.4	86.2	91.1	78.3	81.2	<b>91.3</b>	
Matching, mAP [%]								
Easy	46.7	63.2	69.7	69.5	35.2	52.8	<b>70.1</b>	
Hard	20.4	43.7	53.3	60.0	22.8	40.8	<b>62.3</b>	
Tough	0.09	26.4	35.4	48.0	13.9	28.1	<b>51.5</b>	
Mean	25.5	44.5	52.8	59.1	23.9	40.5	<b>61.3</b>	
Retrieval, mAP [%]								
Easy	64.7	78.6	81.9	80.7	56.2	72.2	<b>81.8</b>	
Hard	37.9	65.44	71.7	75.9	47.0	65.2	<b>78.4</b>	
Tough	22.7	48.6	55.8	67.9	37.3	54.8	<b>72.1</b>	
Mean	41.7	64.2	69.8	74.9	46.8	64.0	<b>77.4</b>	

Table 2: Evaluation results on HPatches [2] of three complementary tasks: patch verification, matching and retrieval.

As shown in Tab. 2, the proposed loss delivers notable improvements over the previous best-performing GeoDesc [23] under similar training settings except for the loss design. Besides, the proposed loss clearly shows better convergence compared with [38] and [17]. Although we suspect that the loss in [17] may perform better with careful parameter searching, the proposed loss is advantageous due to its self-adaptivity without the need of complex heuristics or manual tuning. In addition, once equipped with the resulting model as a base, the augmentation results can be further improved by a healthy margin, denoted as *multi-context+* in Tab. 1 (right columns). We will use this model to complete the following experiments.

#### 5.5 Generalization

**Wild dataset.** The evaluation results on two challenging datasets (*outdoor* YFCC100M [37] and *indoor* SUN3D [42]) are presented in Tab. 3. The proposed multi-context augmentation delivers  $\sim 35\%$  and  $\sim 125\%$  improvements over the previous state of the art, which effectively

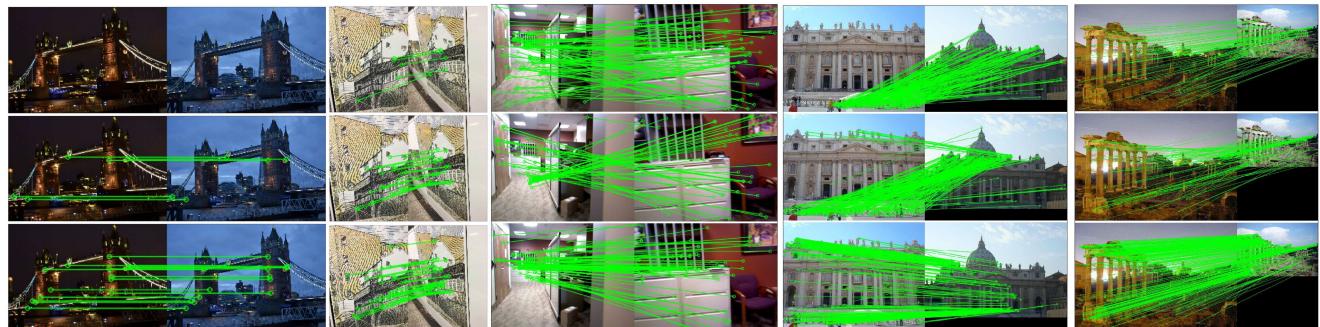


Figure 6: Matching results after RANSAC in different challenging scenarios. The augmented feature not only helps to find more inlier matches, but allows a more accurate recovery of camera geometry.

demonstrates the strong generalization ability of the learned context features in practical scenes.

	SIFT	L2-Net	HardNet	GeoDesc	Ours
median number of inlier matches					
indoor	138	153	239	271	<b>365</b>
outdoor	168	173	219	214	<b>482</b>

Table 3: Evaluation results on wild datasets: *indoor* SUN3D [42] and *outdoor* YFCC100M [37].

**SfM dataset.** We further demonstrate the improvement in complex SfM pipeline. As shown in Tab. 4, the integration of augmented feature generalizes well among different scenes even in large-scale SfM tasks, meanwhile consistently boosts the completeness of sparse reconstruction. Some matching results are presented in Fig. 6, and more visualizations can be found in the appendix.

	# Images	# Registered	# Sparse Points	# Observations
Fountain	SIFT	11	11	10,004
	GeoDesc		11	16,687
	Ours	11	<b>16,965</b>	<b>84K</b>
Herzjesu	SIFT	8	8	4,916
	GeoDesc		8	8,720
	Ours	8	<b>9,429</b>	<b>40K</b>
South Building	SIFT	128	128	62,780
	GeoDesc		128	170,306
	Ours	128	<b>174,359</b>	<b>893K</b>
Roman Forum	SIFT	2,364	1,407	242,192
	GeoDesc		1,566	770,363
	Ours		<b>1,571</b>	<b>848,319</b>
Alamo	SIFT	2,915	743	120,713
	GeoDesc		893	353,329
	Ours		<b>921</b>	<b>424,348</b>

Table 4: Evaluation results on SfM dataset [34].

## 5.6. Discussions

**Invariance property.** We again use *Recall* and evaluate on Heinly benchmark [14] to quantify the invariance property. As shown in Tab 5, the proposed method improves remarkably over the previous best-performing descriptor, except for some minor underperformance regarding *Rotation* change where images are rotated up to 180°, which may be caused by the essential inability of being fully rotation-invariant especially for the regional feature extractor.

**Computational cost.** Towards practicability, we only use basic and shallow MLPs or non-parametric context normalization in our framework design, which thus introduces little

	SIFT	GeoDesc	Ours
JPEG	60.7	66.1	<b>77.8</b>
Blur	41.0	47.7	<b>56.9</b>
Exposure	78.2	86.4	<b>87.8</b>
Day-Night	29.2	39.6	<b>44.6</b>
Scale	81.2	85.8	<b>87.9</b>
Rotation	82.4	<b>87.6</b>	87.3
Scale-Rotation	29.6	33.7	<b>38.0</b>
Planar	48.2	59.1	<b>61.3</b>

Table 5: Evaluation results regarding different transformations on Heinly benchmark [14].

computational overhead, i.e., ~6% time cost on a NVIDIA GTX 1080 compared with raw local feature, as reported in Tab. 6. Besides, we consider that regional features are often off-the-shelf in practical pipelines, e.g., from a retrieval model deployed in SfM pipeline for accelerating image matching, which can be thus reused without introducing extra cost, achieving system-level efficiency and integrity.

	Preparation		Augmentation		
	local feat.	regional feat.	geo. context	vis. context	multi-context
Time (ms)	351	49	8	14	21
FLOPs (B)	802.9	123.4	3.0	13.9	16.9
Params (M)	2.4	24.5	0.2	3.1	3.3

Table 6: The computational cost of proposed framework. The evaluation tests 10k keypoints and 896 × 896 images.

## 6. Conclusion

In contrast to current trends, we have addressed the importance of introducing *context awareness* in learning local descriptors. The augmentation framework takes keypoint location, raw local and high-level regional feature as input, from which two types of context are encoded, including *geometric* and *visual* context. The training process adopts a novel learning scheme and a new loss that is self-adaptive to the task difficulty. We have extensively evaluated the proposed framework on several large-scale datasets that cover diverse practical scenes, outperforming the state of the art by a significant margin and showing strong generalization and the practicability of the proposed method.

864

## References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] V. Balntas, E. Johns, L. Tang, and K. Mikolajczyk. Pn-net: Conjoined triple deep network for learning local image descriptors. In *arXiv*, 2016. 2
- [2] V. Balntas, K. Lenc, A. Vedaldi, and K. Mikolajczyk. Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. In *CVPR*, 2017. 1, 6, 7
- [3] V. Balntas, E. Riba, D. Ponsa, and K. Mikolajczyk. Learning local feature descriptors with triplets and shallow convolutional neural networks. In *BMVC*, 2016. 2
- [4] M. Brown and D. G. Lowe. Automatic panoramic image stitching using invariant features. In *IJCV*, 2007. 6
- [5] H. Deng, T. Birdal, and S. Ilic. Ppfnet: Global context aware local features for robust 3d point matching. In *CVPR*, 2018. 2, 3, 4, 6, 7
- [6] D. DeTone, T. Malisiewicz, and A. Rabinovich. Deep image homography estimation. In *arXiv*, 2016. 6, 11
- [7] J. M. Dmytro Mishkin, Filip Radenovic. Repeatability is not enough: learning discriminative affine regions via discriminability. In *ECCV*, 2018. 2, 5
- [8] J. Dong and S. Soatto. Domain-size pooling in local descriptors: Dsp-sift. In *CVPR*, 2015. 1
- [9] X. Han, T. Leung, Y. Jia, R. Sukthankar, and A. C. Berg. Matchnet: Unifying feature and metric learning for patch-based matching. In *CVPR*, 2015. 2
- [10] W. Hartmann, M. Havlena, and K. Schindler. Predicting matchability. In *CVPR*, 2014. 4
- [11] K. He, Y. Lu, and S. Sclaroff. Local descriptors optimized for average precision. In *CVPR*, 2018. 2
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 4
- [14] J. Heinly, E. Dunn, and J.-M. Frahm. Comparative evaluation of binary features. In *ECCV*, 2012. 6, 8
- [15] E. Hoffer, I. Hubara, and D. Soudry. Fix your classifier: the marginal value of training the last weight layer. In *ICLR*, 2018. 5
- [16] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *NIPS*, 2015. 3
- [17] M. Keller, Z. Chen, F. Maffra, P. Schmuck, and M. Chli. Learning deep descriptors with scale-aware triplet networks. In *CVPR*, 2018. 2, 5, 7
- [18] N. Kobyshev, H. Riemenschneider, and L. Van Gool. Matching features correctly through semantic understanding. In *3DV*, 2014. 2
- [19] B. Kumar, G. Carneiro, I. Reid, et al. Learning local image descriptors with deep siamese and triplet convolutional networks by minimising global loss functions. In *CVPR*, 2016. 1, 2, 3, 6, 7
- [20] K. Lenc and A. Vedaldi. Large scale evaluation of local image feature detectors on homography datasets. In *BMVC*, 2018. 6
- [21] S. Li, L. Yuan, J. Sun, and L. Quan. Dual-feature warping-based motion model estimation. In *CVPR*, 2015.
- [22] D. G. Lowe. Distinctive image features from scale-invariant keypoints. 2004. 1, 7, 11
- [23] Z. Luo, T. Shen, L. Zhou, S. Zhu, R. Zhang, Y. Yao, T. Fang, and L. Quan. Geodesc: Learning local descriptors by integrating geometry constraints. In *ECCV*, 2018. 1, 2, 3, 5, 6, 7, 11
- [24] J. Matas, O. Chum, M. Urban, and T. Pajdla. Robust wide-baseline stereo from maximally stable extremal regions. In *Image and vision computing*, 2004. 1
- [25] A. Mishchuk, D. Mishkin, F. Radenovic, and J. Matas. Working hard to know your neighbor's margins: Local descriptor learning loss. In *NIPS*, 2017. 2, 3, 5, 7
- [26] K. Moo Yi, Y. Verdie, P. Fua, and V. Lepetit. Learning to assign orientations to feature points. In *CVPR*, 2016. 2
- [27] D. Nister and H. Stewenius. Scalable recognition with a vocabulary tree. In *CVPR*, 2006. 1
- [28] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017. 2, 3, 4, 7
- [29] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NIPS*, 2017. 2, 3
- [30] F. Radenović, G. Tolias, and O. Chum. Cnn image retrieval learns from bow: Unsupervised fine-tuning with hard examples. In *ECCV*, 2016. 1, 3, 4, 6, 11
- [31] I. Rocco, R. Arandjelovic, and J. Sivic. End-to-end weakly-supervised semantic alignment. In *CVPR*, 2018. 2
- [32] N. Savinov, A. Seki, L. Ladicky, T. Sattler, and M. Pollefeys. Quad-networks: unsupervised learning to rank for interest point detection. In *CVPR*, 2017. 4
- [33] J. L. Schönberger and J.-M. Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 1, 6
- [34] J. L. Schönberger, H. Hardmeier, T. Sattler, and M. Pollefeys. Comparative evaluation of hand-crafted and learned local features. In *CVPR*, 2017. 6, 8
- [35] E. Simo-Serra, E. Trulls, L. Ferraz, I. Kokkinos, P. Fua, and F. Moreno-Noguer. Discriminative learning of deep convolutional feature point descriptors. In *CVPR*, 2015. 2
- [36] C. Strecha, W. Von Hansen, L. Van Gool, P. Fua, and U. Thoennessen. On benchmarking camera calibration and multi-view stereo for high resolution imagery. In *CVPR*, 2008. 6
- [37] B. Thomee, D. A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, and L.-J. Li. Yfcc100m: The new data in multimedia research. In *CACM*, 2016. 6, 7, 8, 11
- [38] Y. Tian, B. Fan, F. Wu, et al. L2-net: Deep learning of discriminative patch descriptor in euclidean space. In *CVPR*, 2017. 1, 2, 3, 5, 6, 7
- [39] G. Tolias, R. Sicre, and H. Jégou. Particular object retrieval with integral max-pooling of cnn activations. In *ICLR*, 2016. 3
- [40] N. Ufer and B. Ommer. Deep semantic feature matching. In *CVPR*, 2017. 2
- [41] K. Wilson and N. Snavely. Robust global translations with 1dsfm. In *ECCV*, 2014. 6
- [42] J. Xiao, A. Owens, and A. Torralba. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *ICCV*, 2013. 6, 7, 8, 11

- 972 [43] K. M. Yi, E. Trulls, V. Lepetit, and P. Fua. Lift: Learned 1026  
973 invariant feature transform. In *ECCV*, 2016. 2, 3, 5, 6 1027  
974 [44] K. M. Yi, E. Trulls, Y. Ono, V. Lepetit, M. Salzmann, and 1028  
975 P. Fua. Learning to find good correspondences. In *CVPR*, 1029  
976 2018. 2, 3, 4, 6, 7, 11 1030  
977 [45] S. Zagoruyko and N. Komodakis. Learning to compare im- 1031  
978 age patches via convolutional neural networks. In *CVPR*, 1032  
979 2015. 1, 2, 3, 6, 7 1033  
980 [46] X. Zhang, X. Y. Felix, S. Kumar, and S.-F. Chang. Learning 1034  
981 spread-out local feature descriptors. In *ICCV*, 2017. 2 1035  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

## A. Supplementary appendix

### A.1 Implementation details

**Network architecture details.** In terms of the matchability predictor, we construct simple 4-layer MLPs whose output node numbers are 128, 32, 32, 1, respectively. The visual context encoder is composed of two MLPs, located before/after the concatenation with raw local features. We insert context normalization only into the former MLPs in the way shown in Fig. 7, while insertion in the latter one is observed to harm the performance.

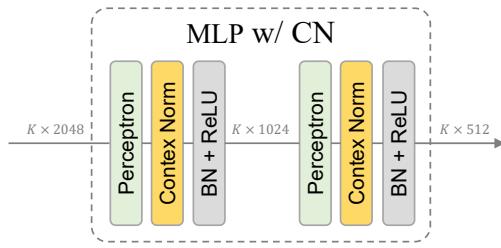


Figure 7: The construction of MLP module with context normalization in visual context encoder.

**Performance of the retrieval model.** The retrieval model is trained on *Large-Scale Landmark Recognition Challenge*<sup>1</sup>, which consists of more than 1M landmark images. Instead of adopting the training scheme in [30], we find that the model pretrained on landmark classification task (containing 15K classes) suffices to produce satisfactory results in practice. We have validated the retrieval model with MAC aggregation on standard Oxford dataset, resulting in mAP of 0.83 on par with [30] whose mAP is 0.80. The performance of proposed augmentation is expected to be further boosted with the evolution of its high-level feature extractor, which we leave as a future work to further achieve the system-level efficacy.

**Keypoint coordinate augmentation.** Similar to [6], we choose to use the 4-point parameterization, which represents a homography as follows:

$$H_{4point} = \begin{Bmatrix} u_1 + \Delta u_1 & v_1 + \Delta v_1 \\ u_2 + \Delta u_2 & v_2 + \Delta v_2 \\ u_3 + \Delta u_3 & v_3 + \Delta v_3 \\ u_4 + \Delta u_4 & v_4 + \Delta v_4 \end{Bmatrix},$$

where  $(u_1, v_1), (u_2, v_2), (u_3, v_3), (u_4, v_4)$  are four corner points at  $(-1, 1), (1, 1), (-1, -1), (1, -1)$ , and  $\Delta u_i, \Delta v_i$  are random variables between  $(-s, s)$ . One can easily convert  $H_{4point}$  to a standard  $3 \times 3$  homography by, e.g., normalized Direct Linear Transform (DLT) algorithm. We choose  $s = 0.5$ , which means that the keypoint set can be

<sup>1</sup><https://landmarkscvprw18.github.io/>

perturbed by a maximum of one quarter of the total image size. We then apply the random homography on the keypoint coordinate before fed into geometric context encoder.

### A.2 Training with softmax temperature

We plot the growth of softmax temperature and its respective loss decrease in Fig. 8. As can be seen, the softmax temperature fast grows at the beginning and gradually converges to a constant value. As mentioned in Sec. 4.1, the softmax temperature is regularized with the same network weight decay, whereas we have observed that the eschewing of regularization does not harm the performance but results in a larger temperature value, e.g.,  $\sim 42$ .

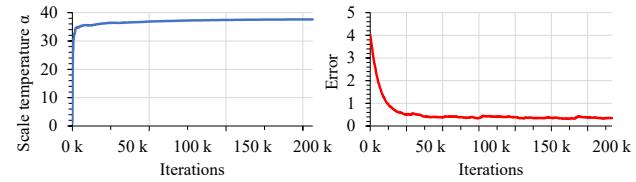


Figure 8: Left: the growth of scale temperature. Right: the respective decrease of loss.

### A.3 Ratio test

In previous experiments of image matching, we did not apply any outlier rejection (e.g., cross check, ratio test [22]) for all methods for fair comparison, whereas the early rejection is generally crucial and necessary to later geometry computation, e.g., recovering camera pose. In particular, ratio test [22] has demonstrated remarkable success, we thus follow the practice in [23] to determine the ratio criteria of proposed augmented feature. Specifically, given *# Correct Matches* defined in Sec. 5.3, we test on HPSequence and aim to find a proper ratio that achieves *Precision* = *# Putative Matches* / *# Correct Matches* similar to SIFT. As a result, we choose 0.89 for the proposed descriptor.

	SIFT	GeoDesc	Ours
<i>mAP of pose (error threshold 20°)</i>			
indoor	37.4	41.8	<b>42.9</b>
outdoor	17.9	20.5	<b>22.5</b>

Table 7: Pose evaluation on wild datasets with ratio test applied: *indoor* SUN3D [42] and *outdoor* YFCC100M [37].

To demonstrate the efficacy of obtained ratio, we evaluate on the wild indoor/outdoor data [42, 37] with an error metric of relative camera pose. Following the protocols defined in [44], we use mean average precision (mAP) of a certain threshold (e.g.,  $20^\circ$ ) to quantify the error of rotation and translation. For comparison, we use ratio criteria of 0.80 for SIFT [22] and 0.89 for GeoDesc [23], and present

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

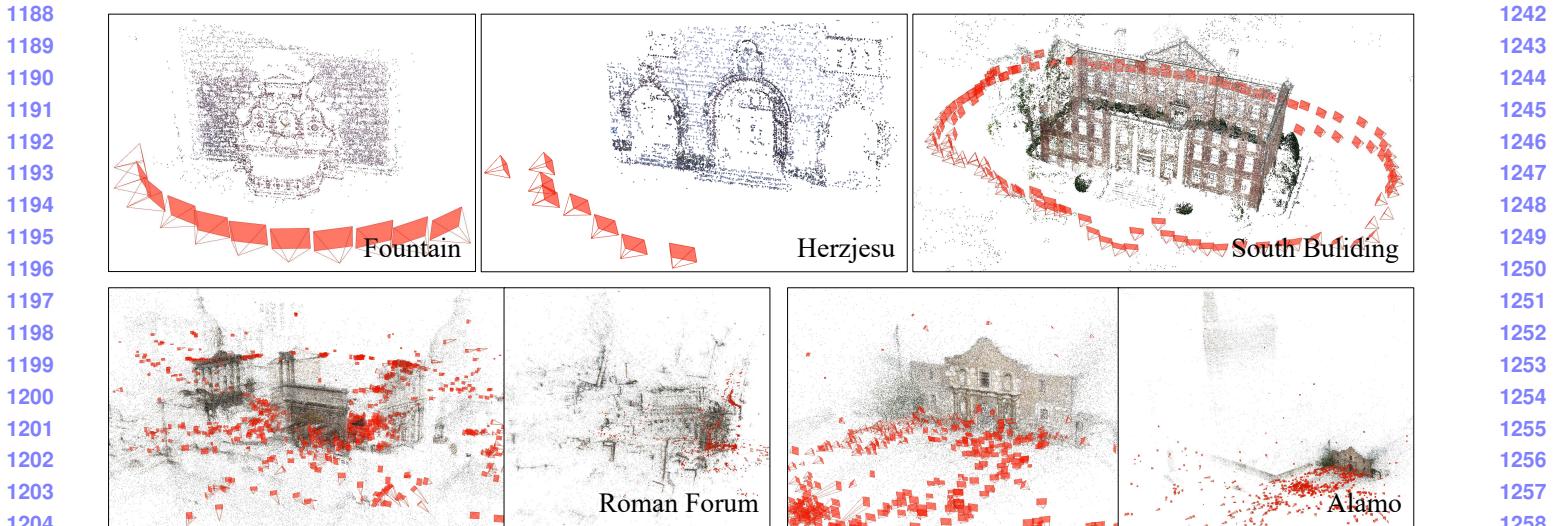


Figure 9: Visualizations of SfM results of Sec. 5.5.

evaluation results in Tab. 7, showing consistent improvements on pose estimation with proper outlier rejection.

#### A.4 Different domain sizes

Somewhat counter-intuitively, the CS structure improves marginally on image matching as reported in Tab. 1. To further study this phenomenon, we compare the patch sampling from different domain sizes, including the original SIFT’s ( $1\times$ ) as used in previous experiments, half ( $0.5\times$ ) or double ( $2\times$ ) sizes. We also compare the aggregation of multiple sizes, i.e., the original and halved ( $1+0.5\times$ ) or the original and doubled ( $1+2\times$ ). Instead of concatenating features as used by CS structure, we apply the simple summing-and-normalizing aggregation in Sec. 3.4 to avoid increasing the dimensionality. We experiments with our *context+* model, and as shown in Tab. 8, when only single size is adopted, the original ‘ $1\times$ ’ performs best as consistent with the training setting. In addition, when combining a larger size, we can further boost the proposed method by a considerable margin, yet leading to excessive computational cost and doubling the inference time. In practice, it is compatible with the proposed framework and can be applicable where high accuracy is in demand.

domain size	Recall i/v	
0.5×	59.27	68.19
2×	62.17	71.67
1×	<b>65.76</b>	<b>75.64</b>
(1 + 0.5)×	65.32	75.05
(1 + 2)×	<b>66.72</b>	<b>77.30</b>

Table 8: The efficacy of extracting local features from different domain sizes.

#### A.5 Invariance of density change

We further demonstrate the robustness regarding density change on HPSequences, of which images are feature-rich and have keypoints up to 15k. Beside of sampling keypoints of different numbers, we consider a more challenging case where *no sampling* and *all detected keypoints* are used. As presented in Fig. 10, the proposed method delivers consistent improvements in terms of all cases, which demonstrates the reliable invariance property acquired by context encoders.

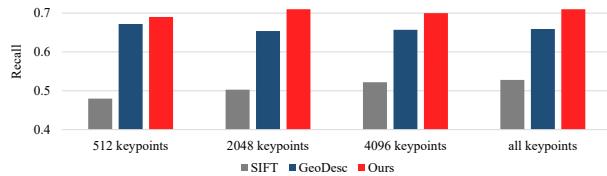
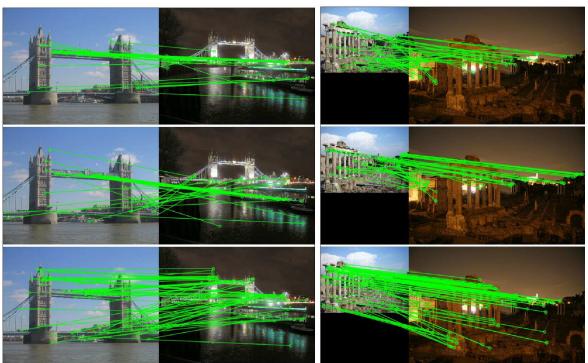


Figure 10: The performance of proposed augmentation scheme regarding density change of keypoints.

#### A.6 More visualizations

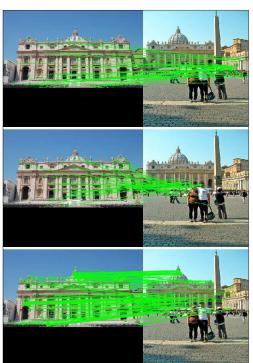
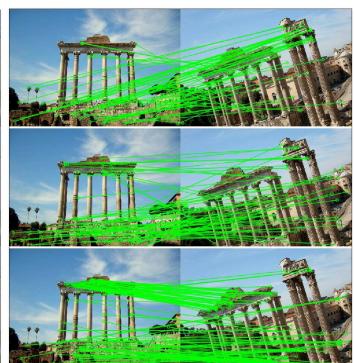
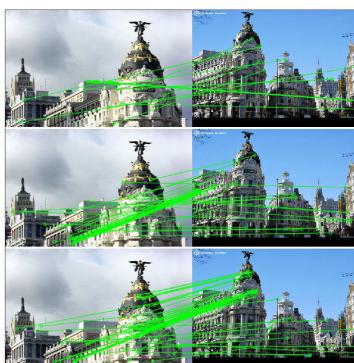
We have provided more visualizations of previous experiments in Fig. 9 (SfM results in Sec. 5.5) and Fig. 11 (image matching results w.r.t different image transformations).

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307



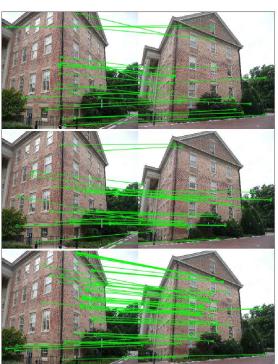
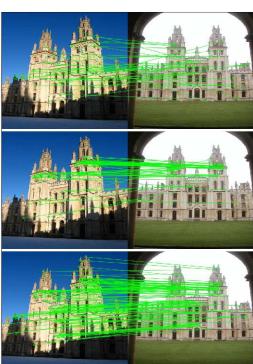
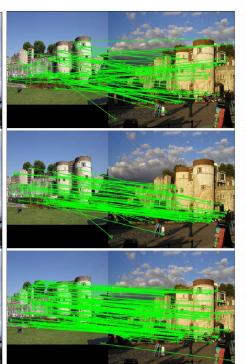
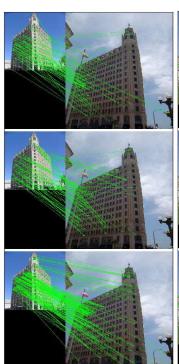
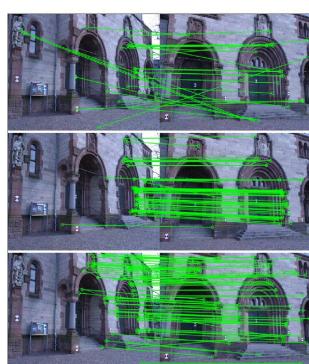
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361

1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320



1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374

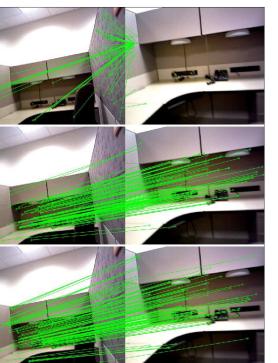
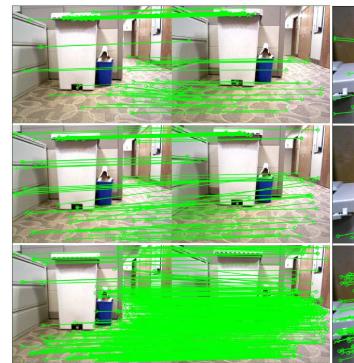
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334



1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388

1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347

### Perspective change



1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401

1348  
1349

### Indoor scene (repetitive or texture-less pattern)

Figure 11: Image matching results after RANSAC. From top to bottom: SIFT, GeoDesc and proposed augmented feature.