

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015054
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

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 practicability in geometric matching applications.

1. Introduction

Designing powerful local feature descriptor is a fundamental problem in applications such as wide-baseline matching [23], image retrieval [26], and structure-from-motion (SfM) [32]. 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]

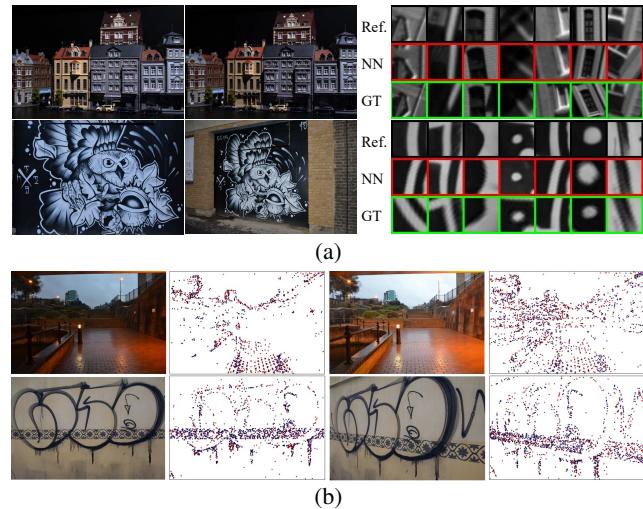


Figure 1: (a) Saturated results on standard benchmark [2] by the recent advance [22]. The search of nearest neighbors (NN) returns false matches though visually similar to groundtruth (GT), indicating the limitation of relying on only local visual information. (b) 2D keypoints distribute structurally, on which we human beings are capable of establishing coarse matches even without color information.

and adopted by recent learned advances [44, 19, 37]. Beside of the challenge on properly selecting domain sizes and the concern of deviating the teachings from scale-space theory [21], a naïve multi-scale implementation may cost excessive computation such as doubled inference time and doubled dimensionality as in [44, 19, 37]. Seeking for a more acceptable accuracy-efficiency trade-off, we instead resort to well-studied high-level image representation that has essentially incorporated rich contextual information, and aim to effectively enhance the local feature description with off-the-shelf visual understandings.

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

108 beings to establish coarse matches even without color information.
109 In essence, keypoints can be bridged by their
110 coordinates, from which we can explore *geometric context*
111 to help to alleviate the ambiguity of local visual feature.
112

113 Thus far, we have discussed two context candidates, re-
114 ferred to as *visual context* and *geometric context* that incor-
115 porate high-level visual understandings and geometric cues
116 of 2D keypoint distribution, respectively. Instead of learning
117 a completely new descriptor, in the present work, we
118 target to flexibly leverage the above context awareness to
119 augment off-the-shelf raw local features, in which process
120 we consider the key challenges threefold:
121

- Proper integration of local geometric feature and high-level visual understandings. As keypoint description requires sub-pixel accuracy, the integration is not supposed to mix up raw representation of local details.
- Instability of 2D keypoints. Due to image appearance changes, keypoint distribution often suffers from substantial variations of sparsity, non-uniformity or perspective, which raises a challenge on acquiring strong invariance property of the context encoder.
- Effective learning scheme. Input signals and features in different modalities are supposed to be efficiently processed and aggregated in a unified framework.

134 Finally, regarding practicability, the augmentation is not
135 supposed to be too complex to introduce excessive compu-
136 tational cost, as the local feature description is often re-
137 garded as part of preprocessing in practical pipelines.
138

139 Although contextual information has been widely ap-
140 plied in computer vision tasks, the challenges faced by lo-
141 cal feature learning are substantially different, posing many
142 non-trivial technical and systematic issues to overcome. In
143 this paper, we propose a unified augmentation scheme that
144 leverages and aggregates cross-modality context, of which
145 the contributions are summarized threefold: 1) a novel *vi-
146 sual context encoder* that integrates high-level visual un-
147 derstandings from *regional image representation*, a technique
148 often used by image retrieval [29], 2) a novel *geometric con-
149 text encoder* that exploits geometric cues from raw 2D key-
150 point distribution while being robust to complex variations.
151 3) A novel learning scheme with a new loss that requires no
152 manual parameter tuning and improves the convergence by
153 mitigating scale affects in N-pair loss. To our best knowl-
154 edge, it is the first work that emphasizes the importance of
155 context awareness in 2D local feature learning.

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

2. Related Work

162 **Learned local features.** Initially, local descriptors are
163 jointly learned with a new metric [9, 44], which is later
164 simplified as direct comparison in Euclidean space [34,
165 42, 3, 19, 1]. More recently, attention is drawn on effi-
166 cient training data sampling [37, 24, 11], effective regulari-
167 zations [37, 46], and geometric shape estimation of input
168 patches [25, 7]. However, most of above methods take *indi-
169 vidual* image patches as input, whereas in the present work,
170 we aim to make use of contextual cues beyond the local de-
171 tail and take advantage of features in multiple modalities.
172

173 **Context awareness.** Although widely introduced in many
174 tasks, context awareness has received little attention in
175 learning 2D local descriptors. In terms of visual context, the
176 central-surround (CS) structure [44, 19, 37] leverages multi-
177 scale information by additionally extracting features from
178 the central part of patch to boost the performance, whereas
179 sacrificing computational efficiency due to doubled extrac-
180 tion time and doubled feature dimensionality. Regarding
181 semantics, previous practice [18] designs a new comparison
182 metric and describes features directly from histogram of seman-
183 tic labels. On the other hand, a family of studies has fo-
184 cused on finding semantic correspondences [39, 30] across
185 *different* objects of the same category, of which the purpose
186 is substantially different from our geometric matching. Be-
187 side of visual information, a recent advance [43] explores
188 to encode motion context for identifying outlier from image
189 correspondences, i.e., 4-d coordinate pairs, whereas we aim
190 to exploit geometric context from *single* image without any
191 reference. Overall, encoding proper context is non-trivial
192 and still unclear in 2D local feature learning.
193

194 **Point feature learning.** In the present work, one of our
195 goals is to explore geometric features from keypoint dis-
196 tribution, we thus resort to PointNet [27] and its vari-
197 ants [28, 5, 43] to consume unordered points. Although
198 great success has been shown in learning tasks on 3D points,
199 there are only few studies exploiting the potential outcome
200 of 2D keypoints. In essence, keypoint structure is not intui-
201 tively meaningful and robust, as being highly dependent on
202 the performance of interest point detectors and strongly af-
203 fected by image variations. However, in descriptor learning,
204 we consider the keypoint location as an important cue that
205 bridge each individual local feature, constructing a unified
206 instance that reveals high-level contextual information.
207

208 **Loss formulation.** Recent local descriptors are often
209 evolved with advanced variants of N-pair losses. Initially,
210 L2-Net [37] adopts a log-likelihood formulation, which is
211 later extended by HardNet [24] with hard negative triplet
212 margin loss. Furthermore, GeoDesc [22] applies an adap-
213 tive margin value to improve the convergence in terms of
214 different data sampling strategies, where AffNet [7] ap-
215 proaches the same issue by fixing the distance of hardest

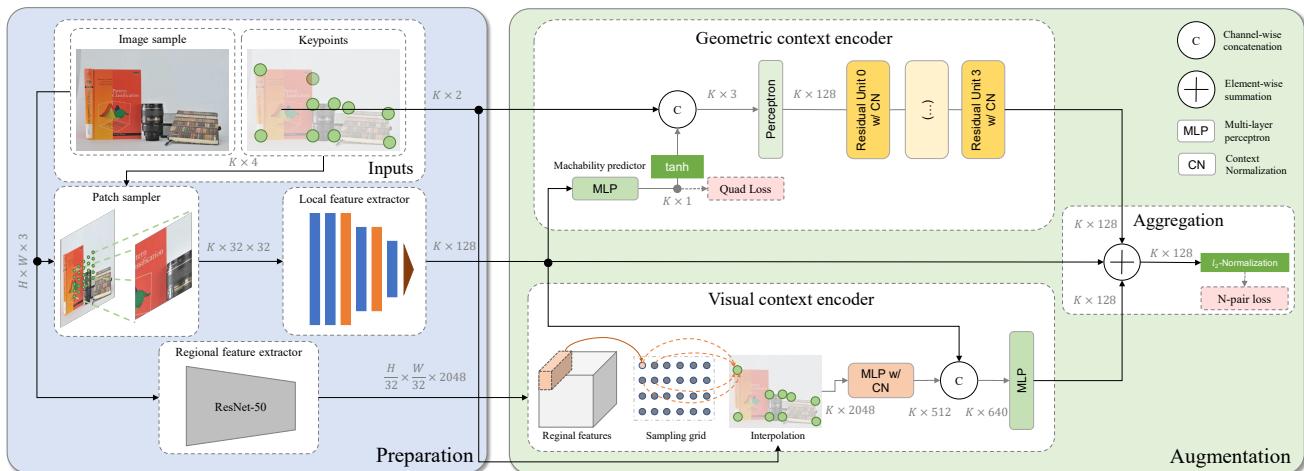


Figure 2: The network architecture of local feature augmentation.

negative sample in the training. Meanwhile, on the other hand, DOAP [11] extends the N-pair loss to a list-wise ranking loss, while [17] points out and studies the scale affects in N-pair losses while introducing additional manual tuning of hyper-parameters. Principally, the loss is supposed to encourage similar patches to be close while dissimilar ones to be distant 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 32×32 gray-scale patches. Akin to [42, 22], 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 [37, 24, 22].

Regional feature extractor. In contrast to aggregating features of different domain sizes [44, 19, 37], in the present

work, we fix the sampling scale of patches, and exploit contextual cues by inspiration of well-studied regional representation in image retrieval tasks [38, 29]. Without the loss of generality, we reuse features from an off-the-shelf deep image retrieval model of ResNet-50 [12]. As in [38], 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 keypoint is inappropriate to serve as robust contextual cues, as its presence is heavily dependent on image appearance and thus affected by various image variations. As a result, keypoint distribution depicting the same scene may suffer from significant density or structure variations, as examples shown in Fig. 1b. Hence, acquiring strong invariance property is the key challenge when designing the context encoder.

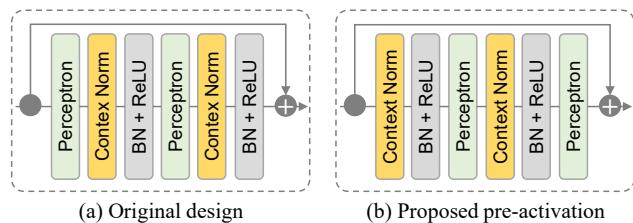
Initially, we attempt to approach the goal by PointNet [27] and its variants [28, 5]. Although having shown great success on processing 3D point clouds, the family of PointNet methods fails to achieve consistent improvement in terms of processing 2D points (Sec. 5.4.1). Instead, we resort to [43], which 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

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

324 normalization (CN) proposed in [43], and extend its usability
 325 to processing 2D points in *single* image.
 326

327 Formally, context normalization is a non-parametric operation
 328 that simply normalizes feature maps according to their distribution, expressed as $\hat{o}_i^l = \frac{(o_i^l - \mu^l)}{\sigma^l}$, where o_i^l is
 329 the output of i -th point in layer l , and μ^l, σ^l are mean and standard deviation of the output in layer l . To equip the
 330 operation, we borrow the residual architecture in [43], where each residual unit is built with two perceptrons followed by
 331 context and batch normalization, as illustrated in Fig. 3a.
 332

333 However, the above construction leads to a *non-negative*
 334 output from the residual branch that impacts the representational ability as investigated in [13] and witnessed in our
 335 experiments (Sec. 5.4.1). Following the teachings of [13], we propose to re-arrange the operations in residual unit by
 336 adopting the *pre-activation* construction, which is compatible with the context normalization as presented in Fig. 3b.
 337 We then 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



353 Figure 3: Different designs of residual unit with context
 354 normalization, where the proposed construction improves
 355 by a considerable margin than its original counterpart.
 356

357 Intuitively, the observation that context normalization
 358 (CN) performs better than PointNet in our task can be
 359 interpreted as the ‘inexactness’ nature of CN, which makes a
 360 global association of keypoints with simple statistics about
 361 the distribution and is thus less sensitive to the instability of
 362 2D keypoints. In principle, the ‘inexactness’ suffices to ex-
 363 tract the context prior as we essentially expect to establish
 364 only *coarse* matches on keypoint distribution.
 365

366 **Matchability predictor.** In 3D point cloud processing,
 367 low-level color and normal [27] information or more com-
 368 plex geometric attributes [5] are often adopted to enhance
 369 the representation. Similarly, associating 2D coordinate in-
 370 put with other meaningful attributes would be promising
 371 to boost the performance. However, due to the substantial
 372 variations, e.g. perspective change, it is non-trivial to define
 373 appropriate intermediate attributes on 2D points.
 374

375 Although this issue has been merely discussed, we draw
 376 inspiration from [10], which poses a problem named *match-*
 377 *ability prediction* and targets to decide whether a keypoint
 descriptor is matchable before the matching stage. In practice,
 the learned matchability can be used to guide the key-

378 point sampling and accelerate the matching without sacri-
 379 ficing accuracy. In contrast to criteria such as cornerness
 380 or edgeness of a keypoint, matchability implies high-level
 381 prior knowledge derived from data. Dependent on individ-
 382 ual keypoint, matchability can be naturally associated as a
 383 representative attribute in addition to 2D coordinate.
 384

385 Previously in [10], the matchability prediction is casted
 386 as a binary classification problem, taking individual de-
 387 scriptor as input and inferred by a random forest. In the
 388 present work, we approach this problem with deep learn-
 389 ing techniques, and apply a more strict constraint requiring
 390 consistent prediction between images. Similar to [31, 45],
 391 we resort to an unsupervised learning scheme that aims
 392 to rank points. Specifically, given K corresponding key-
 393 point pairs (p_1^n, p_2^n) , $n \in [1, K]$ from an image pair, we
 394 first extract local features (f_1^n, f_2^n) of each keypoint, then
 395 construct *feature quadruples* as $(f_1^i, f_1^j, f_2^i, f_2^j)$, satisfying
 396 $i, j \in [1, K], i \neq j$ and holding that:
 397

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

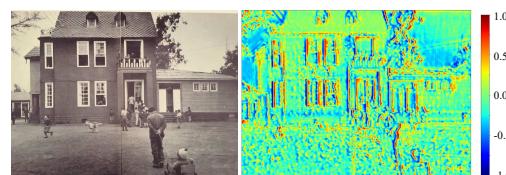
398 where $H(\cdot)$ absorbs the raw local feature into a single real-
 399 valued matchability, implemented as standard multi-layer
 400 perceptrons (MLPs). Here, Cond. 1 aims to preserve a rank-
 401 ing of each keypoint, hence improves the repeatability of
 402 prediction. The condition can be re-written as:
 403

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

404 the final objective can be obtained with hinge loss:
 405

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

406 In the proposed framework, the matchability is learned
 407 as an auxiliary task, which is then activated by a tanh and
 408 associated with keypoint coordinate before fed into the en-
 409 coder, as in Fig. 2. Beside of Eq. 3, the gradient from final
 410 augmented feature will flow through the matchability pre-
 411 dictor, allowing a joint optimization of the entire encoder.
 412 We illustrate the efficacy of matchability in Fig. 4.
 413



414 Figure 4: Visualization of matchability responding to the
 415 entire image (best viewed in color).
 416

432 3.3. Visual context encoder 486

433 This module takes regional features of $\frac{H}{32} \times \frac{W}{32} \times 2048$ 487
 434 in Sec. 3.1, K raw local features and their location as input, 488
 435 producing K corresponding feature vectors. As widely- 489
 436 adopted in other tasks, the main purpose of this module is to 490
 437 integrate visual information in different levels, e.g., semantics. 491
 438 In our context, the key issue is to handle the regional 492
 439 features and keypoints of different numbers. One option as 493
 440 in [5] is to concatenate global representation of entire image 494
 441 on raw local features, where the global feature in our frame- 495
 442 work can be derived by applying Maximum Activations of 496
 443 Convolutions (MAC) aggregation [29] which simply max- 497
 444 pools over all dimensions per feature map. However, such 498
 445 compact representation can mess up the raw local descrip- 499
 446 tion, due to the lack of distinctiveness (Sec. 5.4.1). 500
 447

448 To better preserve the regional distinction, we associate 501
 449 regional features to a regular sampling grid on the image, 502
 450 then interpolate $\frac{H}{32} \times \frac{W}{32}$ grid points at coordinates of the 503
 451 K keypoints. For interpolation, we use the inverse distance 504
 452 weighted average based on k nearest neighbors (in default 505
 453 we use $k = 4$), formulated as:

$$454 \mathbf{f}(\hat{p}_i) = \frac{\sum_{j=1}^k w(p_j) \mathbf{f}(p_j)}{\sum_{j=1}^k w(p_j)}, \text{ and } w(p_j) = \frac{1}{d(\hat{p}_i, p_j)}, \quad (4)$$

455 where $\mathbf{f}(\cdot)$ is the regional feature located at a certain grid 506
 456 point. \hat{p}_i , $i \in [1, N]$ indicates interpolated point, while p_j , 507
 457 $j \in [1, \frac{H}{32} \times \frac{W}{32}]$ indicates original grid point. Next, the 508
 458 dimensionality is reduced by applying point-wise MLPs, 509
 459 where we also insert CN after each perceptron in order to 510
 460 capture global context. Finally, raw local features are 511
 461 concatenated and further mapped by MLPs, forming the final 512
 462 128-d features. The above process is illustrated in Fig. 2. 513
 463

464 3.4. Feature aggregation with raw local feature 514

465 To aggregate the above two types of contextual features, 515
 466 similar to the CS structure, one option is to concatenate 516
 467 them together and forms features of, in our case, 384-d 517
 468 (128×3). However, the increased dimensionality will 518
 469 introduce excessive computational cost in the matching stage of 519
 470 $\mathcal{O}(n^2)$ complexity. Instead, as shown in Tab. 2, we propose 520
 471 to combine different feature streams into a single vector by 521
 472 summing and L2-normalizing them in the end, i.e., without 522
 473 the change of feature dimensionality. Beside of the simplicity, 523
 474 such strategy allows flexible use of the augmentation. 524
 475 For example, in situations where regional features are not 525
 476 available, one may aggregate with only geometric context 526
 477 without the need of retraining the model. 527

480 4. Learning Scheme 528

481 4.1. N-pair loss with softmax temperature 529

482 N-pair losses have been primarily used by recent works. 530
 483 Empirically, the subtractive hinge loss [24, 22, 7] has re- 531

484 ported better performance, of which the main idea is to push 532
 485 similar samples away from dissimilar ones to a certain margin 533
 486 in the descriptor space. However, setting the appropriate 534
 487 margin is tricky, which does not always assure convergence 535
 488 as observed in [22, 7]. More generally, the aspects of 536
 489 making a good loss are studied in [17], from which guidelines 537
 490 are provided about tuning loss coefficients on particular 538
 491 dataset. In this spirit, we aim to ease the pain of parameter 539
 492 searching in [17], and obtain an adaptive loss that allows 540
 493 fast convergence regardless of the learning difficulty. 541
 494

495 We use the log-likelihood version of N-pair loss [37] as 542
 496 a base, which originally does not involve any tunable 543
 497 parameter. Formally, given L2-normalized feature descriptors 544
 498 $\mathbf{F}_1 = [f_1^1 f_1^2 \dots f_1^N]^T, \mathbf{F}_2 = [f_2^1 f_2^2 \dots f_2^N]^T \in \mathbb{R}^{N \times 128}$, 545
 499 the distance matrix $\mathbf{D} = [d_{ij}]_{N \times N}$ can be obtained by 500
 500 $\mathbf{D} = \sqrt{2(1 - \mathbf{F}_1 \mathbf{F}_2^T)}$. By applying, e.g., row-wise softmax, 501
 501 we derive the final loss as:

$$502 \mathcal{L}_{N\text{-pair}} = - \sum_i \log s_{ii}, \quad (5)$$

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

504 Noted that since input features are L2-normalized, the 505 resulting d_{ij} is bounded between 0 and 2, which causes 506 convergence issues due to the scale sensitivity of softmax 507 function [15]. Similarly, we introduce a single scalar 508 parameter α , referred to as *softmax temperature*, to amend the 509 inability of re-scaling the input. The loss now becomes:

$$510 [s_{ij}]_{N \times N} = \text{softmax}(\alpha(2 - \mathbf{D})), \quad (6)$$

511 where α is initialized to 1 and regularized with the same 512 weight decay in the network, which does not require any 513 manual tuning or complex heuristics. In the experiments in 514 Sec. 5.4.2, we show this simple technique improves drastically 515 than the original formulation in [37], whose performance 516 we believe is hindered due to the scale sensitivity of 517 softmax. In the proposed framework, we compute the N- 518 pair loss on augmented features, and obtain the total loss:

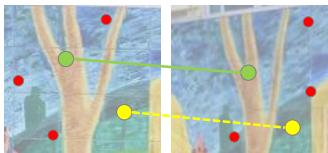
$$519 \mathcal{L}_{total} = \mathcal{L}_{N\text{-pair}} + \lambda \mathcal{L}_{quad}, \quad (7)$$

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

521 4.2. Learning with noisy keypoints 524

522 The training of proposed framework, apparently, needs 523 to be conducted between image pairs instead of isolated 524 patches, which is referred to as simulating image matching 525 in [22]. However, the simulation in [22] is not complete, as 526 it considers only matchable keypoints that are paired with 527 correspondences, whereas in real situation, only a subset of 528 keypoints is repeatable in other images. In practice, as il- 529 lustrated in Fig. 5, we divide keypoints obtained from SfM 530 as in [42, 22] into three categories: i) *Matchable*: repeatable 531 and verified by SfM; ii) *Undiscovered*: repeatable but 532

540 did not survive the SfM. iii) *Unrepeatable*: unable to be re-
 541 detected in other images. In the training, we randomly sam-
 542 ple a number of matchable keypoints as well as some undis-
 543 covered and unrepeatable (*noisy keypoints*) to have a com-
 544 plete simulation, which is necessary to acquire the gen-
 545 eralization ability, otherwise the training considers only ideal
 546 setting which is inconsistent with real applications. In this
 547 setting, the loss is still computed on matchable keypoints,
 548 of which the formal expression is provided in the appendix.
 549



550 Figure 5: We divide keypoints after SfM into three cate-
 551 gories: matchable (green), undiscovered (yellow) and un-
 552 repeatable (red), and aim to perform a complete simulation in
 553 training that incorporates all three types of keypoints.
 554

5. Experiments

5.1. Implementation

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

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 [40, 29] as in [42, 22], 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].
 588

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 [43], we use out-
 599 door YFCC100M [36] (1000 pairs) and indoor SUN3D [41]
 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 [33], we evaluate on SfM dataset
 605 such as well-known *Fountain* and *Herzjesu* [35], or land-
 606 mark collections [40]. We integrate the proposed frame-
 607 work into SfM pipeline, i.e., COLMAP [32], and use the
 608 keypoints provided in [33] to compute the local features.
 609

5.3. Evaluation protocols

Patch level. Following the same protocols of HPathces [2],
 613 we use mean average precision (mAP) for its three subtasks
 614 , including patch verification, matching, and retrieval.
 615

Image level. For HPSequences, we use $Recall = \# Correct$
 616 $Matches / \# Correspondences$ defined in [14], to quantify
 617 the image matching performance, where $\# Correct matches$
 618 are matches found by nearest neighbor searching and ver-
 619 ified by groundtruth geometry, e.g., homography, while $\#$
 620 $Correspondences$ are matches that should have been iden-
 621 tified by the given keypoint locations. Following [14], a
 622 match point is determined to be correct if it is within 2.5
 623 pixels from the wrapped keypoint in the reference image.
 624 We use a standard SIFT detector to localize the keypoints,
 625 of which the number is randomly sampled to 2048. For
 626 YFCC100M [36] and SUN3D [41], we follow the same set-
 627 ting in [43] and report the median number of inlier matches
 628 after RANSAC for each dataset.
 629

Reconstruction level. For clarity, we report metrics in [33]
 630 that quantify the completeness of SfM, including the num-
 631 ber of registered images ($\# Registered$), sparse points ($\#$
 632 $Sparse Points$) and image observations ($\# Observations$).
 633

5.4. Ablation study

5.4.1 Design of context encoder

In this section, we evaluate two splits of HPSequences [2]: *illumination* (*i*) and *viewpoint* (*v*), regarding different image transformations. We report *Recall* as defined in Sec. 5.3. If not specified, we use GeoDesc [22] as a baseline model (*baseline (GeoDesc)*) to extract raw local features, whose parameters are *fixed* during the training of augmentation.
 637

Visual context. We compare four designs, including i) *CS*
 638 (256-d): the central-surround (CS) structure [44, 19, 37] as
 639 described in Sec. 2, which leverages visual information of
 640

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 [21]	47.41	53.19	704
651	CS (256-d) [44, 19, 37]	59.64	71.47	PointNet [27]	59.61	71.16	L2-Net [37]	47.55	54.10	705
652	w/ global feature [5]	58.92	71.22	w/ CN (pre) + xy	61.05	72.47	HardNet [24]	57.61	63.45	706
653	w/ regional feature	63.45	73.57	w/ CN (pre) + xy + raw local feature	60.61	72.69	GeoDesc [22]	59.27	71.44	707
654	w/ regional feature + CN	63.79	73.63	w/ CN (orig.) + xy + matchability	59.94	71.25	multi-context	65.35	74.70	708
655				w/ CN (pre) + xy + matchability	61.52	72.83	multi-context+	65.76	75.64	709

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.

different domain sizes. ii) *w/ global feature*: the integration with global features [5], which is originally designed 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 [44, 19, 37] delivers only marginal improvements despite of the doubled dimensionality. By contrast, 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. Moreover, to strengthen global context 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 [27] 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 [27] does not work well in our context, where the similar phenomenon is also observed in [43] 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 improve-

ments over the state of the art.

5.4.2 Efficacy of softmax temperature

To make a clear demonstration on HPathces [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 [37] without scale temperature, also iii) the loss in [17] with its original parameters.

	SIFT [21]	L2-Ne [37]	HardNet [24]	GeoDesc [22]	w/ loss [37]	w/ loss [17]	proposed
	Verification, mAP [%]						
Easy	80.0	91.4	93.6	94.0	83.3	87.5	93.6
Hard	59.2	83.7	87.6	91.8	78.9	82.1	91.8
Tough	44.9	72.1	77.3	87.6	72.6	73.9	88.4
Mean	61.4	82.4	86.2	91.1	78.3	81.2	91.3
	Matching, mAP [%]						
Easy	46.7	63.2	69.7	69.5	35.2	52.8	70.1
Hard	20.4	43.7	53.3	60.0	22.8	40.8	62.3
Tough	0.09	26.4	35.4	48.0	13.9	28.1	51.5
Mean	25.5	44.5	52.8	59.1	23.9	40.5	61.3
	Retrieval, mAP [%]						
Easy	64.7	78.6	81.9	80.7	56.2	72.2	81.8
Hard	37.9	65.44	71.7	75.9	47.0	65.2	78.4
Tough	22.7	48.6	55.8	67.9	37.3	54.8	72.1
Mean	41.7	64.2	69.8	74.9	46.8	64.0	77.4

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 [22] under similar training settings except for the loss design. Besides, the proposed loss clearly shows better convergence compared with [37] 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 [36] and *indoor* SUN3D [41]) 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

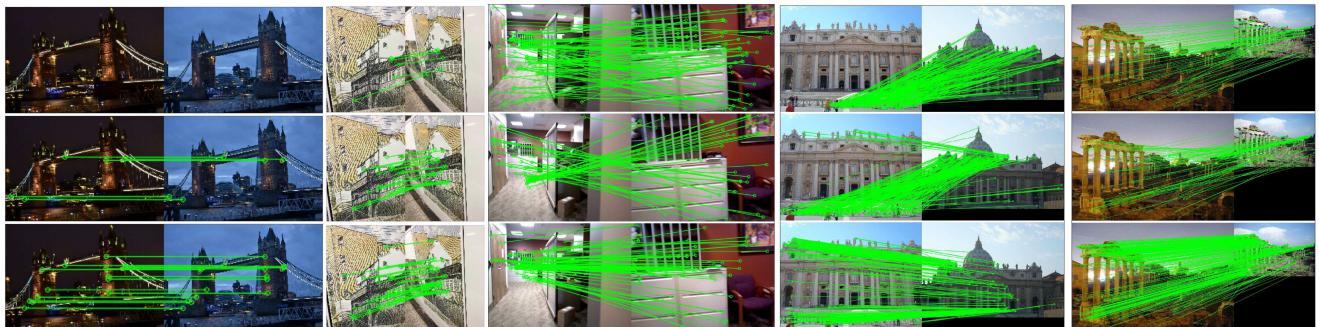


Figure 6: Matching results after RANSAC in different challenging scenarios. From top to bottom: SIFT, GeoDesc and ours. The augmented feature not only helps to find more inlier matches, but allows a more accurate recovery of camera geometry.

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	365
outdoor	168	173	219	214	482

Table 3: Evaluation results on wild datasets: *indoor* SUN3D [41] and *outdoor* YFCC100M [36].

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	10,004	44K
	GeoDesc	11	16,687	83K
	Ours	11	16,965	84K
Herzjesu	SIFT	8	4,916	19K
	GeoDesc	8	8,720	38K
	Ours	8	9,429	40K
South Building	SIFT	128	62,780	353K
	GeoDesc	128	170,306	887K
	Ours	128	174,359	893K
Roman Forum	SIFT	2,364	1,407	1,805K
	GeoDesc		1,566	5,051K
	Ours		1,571	5,484K
Alamo	SIFT	2,915	743	120,713
	GeoDesc		893	353,329
	Ours		921	424,348

Table 4: Evaluation results on SfM dataset [33].

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 normal-

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

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

ization in our framework design, which thus introduces little 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
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

References

- [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, 5, 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, 3
- [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. 3, 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] D. G. Lowe. Distinctive image features from scale-invariant keypoints. 2004. 1, 7, 11, 12
- [22] 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, 12
- [23] 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
- [24] 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
- [25] K. Moo Yi, Y. Verdie, P. Fua, and V. Lepetit. Learning to assign orientations to feature points. In *CVPR*, 2016. 2
- [26] D. Nister and H. Stewenius. Scalable recognition with a vocabulary tree. In *CVPR*, 2006. 1
- [27] 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
- [28] 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
- [29] F. Radenović, G. Tolias, and O. Chum. Cnn image retrieval learns from bow: Unsupervised fine-tuning with hard examples. In *ECCV*, 2016. 2, 3, 5, 6, 11
- [30] I. Rocco, R. Arandjelovic, and J. Sivic. End-to-end weakly-supervised semantic alignment. In *CVPR*, 2018. 2
- [31] 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
- [32] J. L. Schönberger and J.-M. Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 1, 6
- [33] J. L. Schönberger, H. Hardmeier, T. Sattler, and M. Pollefeys. Comparative evaluation of hand-crafted and learned local features. In *CVPR*, 2017. 6, 8
- [34] 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
- [35] 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
- [36] 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, 12
- [37] 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
- [38] G. Tolias, R. Sicre, and H. Jégou. Particular object retrieval with integral max-pooling of cnn activations. In *ICLR*, 2016. 3
- [39] N. Ufer and B. Ommer. Deep semantic feature matching. In *CVPR*, 2017. 2
- [40] K. Wilson and N. Snavely. Robust global translations with 1dsfm. In *ECCV*, 2014. 6
- [41] 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, 12

- 972 [42] 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 [43] 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 [44] 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 [45] L. Zhang and S. Rusinkiewicz. Learning to detect features in 1034
981 texture images. In *CVPR*, 2018. 4 1035
- 982 [46] X. Zhang, X. Y. Felix, S. Kumar, and S.-F. Chang. Learning 1036
983 spread-out local feature descriptors. In *ICCV*, 2017. 2 1037
- 984 1038
- 985 1039
- 986 1040
- 987 1041
- 988 1042
- 989 1043
- 990 1044
- 991 1045
- 992 1046
- 993 1047
- 994 1048
- 995 1049
- 996 1050
- 997 1051
- 998 1052
- 999 1053
- 1000 1054
- 1001 1055
- 1002 1056
- 1003 1057
- 1004 1058
- 1005 1059
- 1006 1060
- 1007 1061
- 1008 1062
- 1009 1063
- 1010 1064
- 1011 1065
- 1012 1066
- 1013 1067
- 1014 1068
- 1015 1069
- 1016 1070
- 1017 1071
- 1018 1072
- 1019 1073
- 1020 1074
- 1021 1075
- 1022 1076
- 1023 1077
- 1024 1078
- 1025 1079

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098

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.

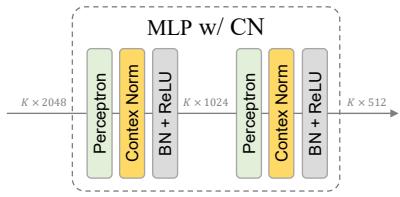


Figure 7: The construction of MLP module with context normalization in visual context encoder.

Performance of the retrieval model. The retrieval model is trained on *Large-Scale Landmark Recognition Challenge*¹, which consists of more than 1M landmark images. Instead of adopting the training scheme in [29], 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 [29] 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 3×3 homography by, e.g., normalized Direct Linear Transform (DLT) algorithm. We choose $s = 0.5$, which means that the keypoint set can be 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.

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

Loss computation with noisy keypoints. Formally, given index sets $C_m = \{i_1, \dots, i_{K_m}\}$ and $C_n = \{i_1, \dots, i_{K_n}\}$, where K_m and K_n are numbers of matchable and noisy keypoints for an image pair, the losses of Eq. 3 and Eq. 5 are now rewritten as:

$$\mathcal{L}'_{quad} = \frac{1}{K_m(K_m - 1)} \sum_{i,j \in C_m, i \neq j} \max(0, 1 - R(\mathbf{f}_1^i, \mathbf{f}_1^j, \mathbf{f}_2^i, \mathbf{f}_2^j)), \quad (8)$$

and

$$\mathcal{L}'_{N-pair} = - \sum_{i \in C_m} \log s_{ii}. \quad (9)$$

Subsequently, adding noisy keypoints will influence \mathcal{L}_{N-pair} as all keypoints have been cross-paired, while it also affects the encoding of geometric context.

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., ~ 42 .

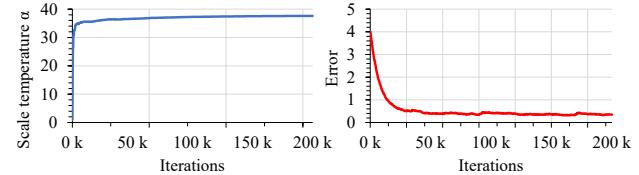


Figure 8: Left: the growth of scale temperature. Right: the respective decrease of loss.

A.3 Ratio test

In previous experiments of image matching, we did not apply any outlier rejection (e.g., cross check, ratio test [21]) 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 [21] has demonstrated remarkable success, we thus follow the practice in [22] 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.

To demonstrate the efficacy of obtained ratio, we evaluate on the wild indoor/outdoor data [41, 36] with an error metric of relative camera pose. Following the protocols defined in [43], we use mean average precision (mAP) of a certain threshold (e.g., 20°) to quantify the error of rotation

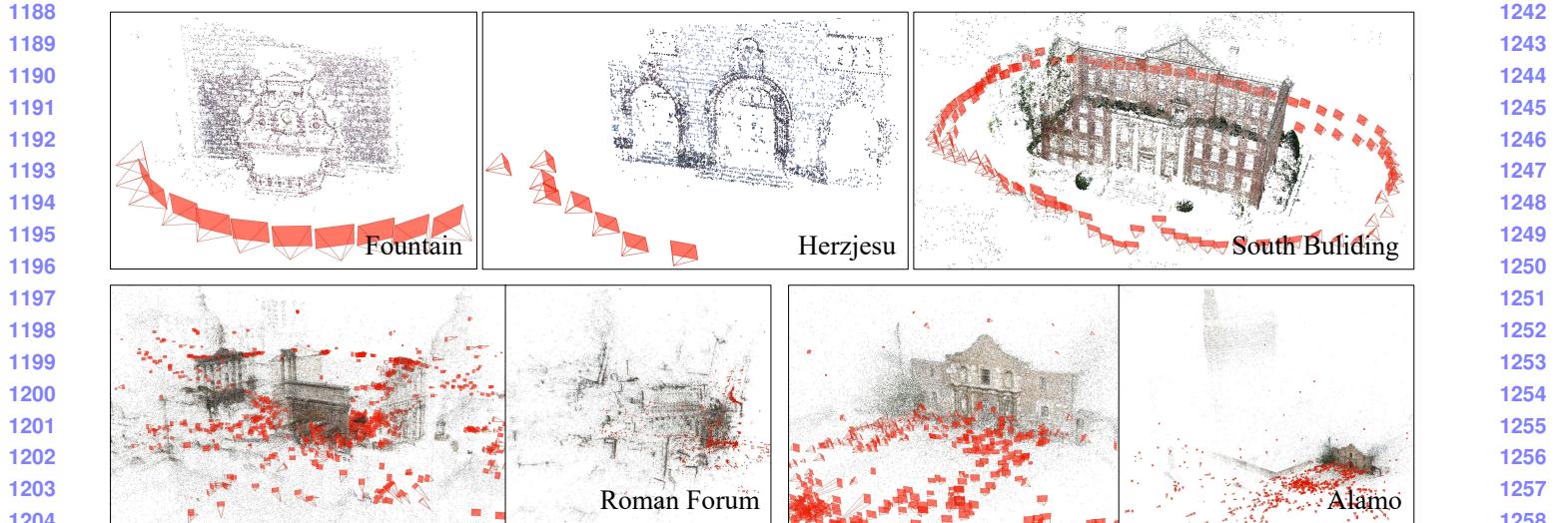


Figure 9: Visualizations of SfM results of Sec. 5.5.

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

Table 7: Pose evaluation on wild datasets with ratio test applied: *indoor* SUN3D [41] and *outdoor* YFCC100M [36].

and translation. For comparison, we use ratio criteria of 0.80 for SIFT [21] and 0.89 for GeoDesc [22], and present 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
$0.5\times$	59.27
$2\times$	62.17
$1\times$	65.76
$(1+0.5)\times$	65.32
$(1+2)\times$	66.72
	75.64
	75.05
	77.30

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.

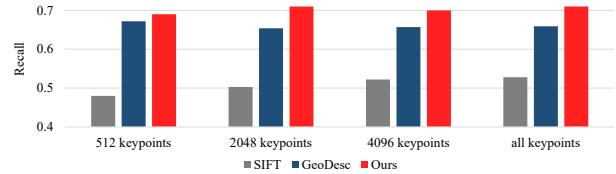
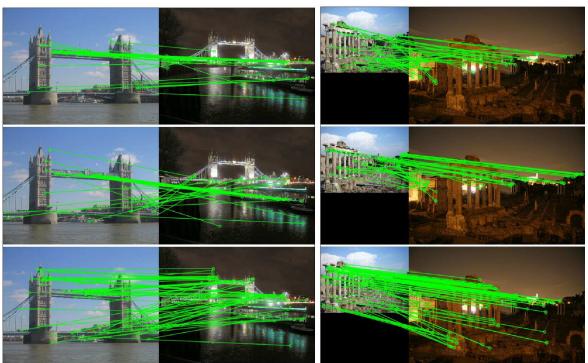


Figure 10: The performance of proposed augmentation scheme regarding density change of keypoints.

A.6 More visualizations

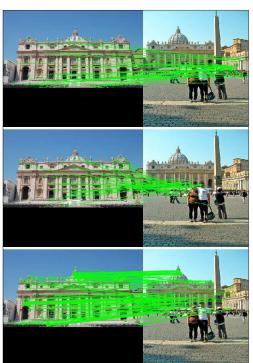
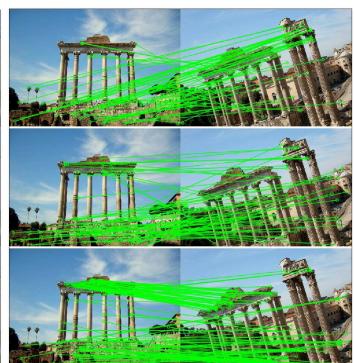
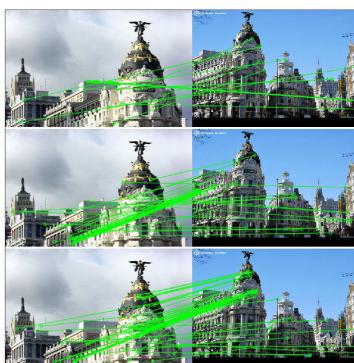
We have provided more visualizations of previous experiments in Fig. 9 (SfM results in Sec. 5.5) and Fig. 11 (image matching results w.r.t different image transformations).

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



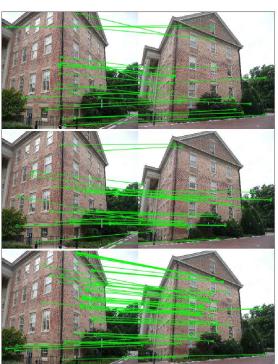
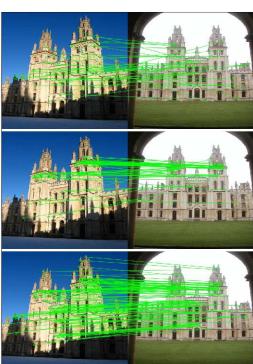
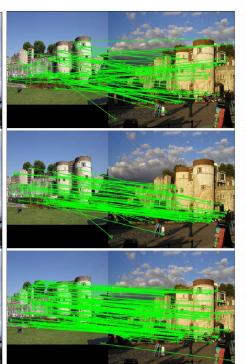
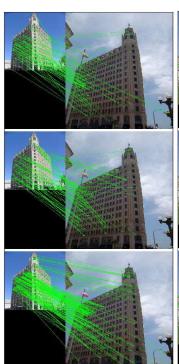
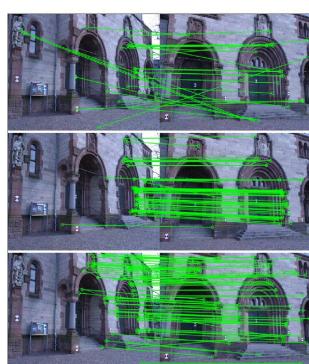
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361

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

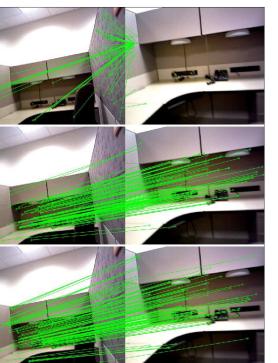
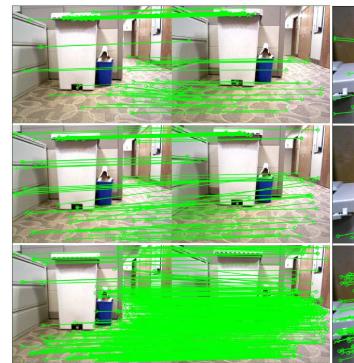
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334



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.