

# **Digging into Self-Supervised Learning of Local Descriptors**

Iaroslav Melekhov\*

## **Problem statement**

To address the challenges, we propose a system shown in Fig. 2. All components are Fully-supervised CNN-based approaches for learning local image self-supervised and do not require human or costly machine guided labelling. per-pixel ground-truth descriptors require keypoint correspondence data which is difficult to acquire at scale. In this work, we focus on understanding the limitations of existing selfsupervised approaches and propose a set of improvements that combined lead to powerful feature descriptors.

## Challenges

Existing self-supervised approaches employ in-pair negative mining. However:

1. Supervised local descriptor learning methods show superior performance with in-batch negative mining. In self-supervised setting, it does not guarantee selection of true negatives. In randomized batch training does it have any impact?

2. Global mining has been accessible only to global representation learning methods. Local mining cost grows exponentially with dataset size. Since global descriptors are functions of local descriptors, can local mining across the dataset be approximated with global mining?



c) Stylization

Figure 1: Stylization creates realistic illumination variation.



Zakaria Laskar\*

Shuzhe Wang Juho Kannala Xiaotian Li

\*equal contribution

# Method



Figure 2: Proposed self-supervised local descriptor learning system

1. HN Mining block. This block selects the negatives for a given descriptor set. We use both in-pair and in-batch mining. The number of positives and negatives is balanced by selecting only the top-k negatives in terms of cosine similarity.

2. Image Retrieval. To enable global mining of local hard negative descriptors we use advances in global image retrieval. For each image in a training batch, in a coarse-to-fine manner, we first find hard-negative images in terms of global similarity. From the pool of global hard-negative images, the top-k hardnegative local descriptors are selected.

3. Stylization. Following Li et al. [37] we use a closed-form solution to image stylization (Fig. 1) based on a deep learning model trained in a self-supervised fashion by minimizing the sum of the reconstruction loss and perceptual loss. The method photorealistically transfers style of a reference photo to a content photo preserving local scene geometry.

# imelekhov.com/hndesc





Figure 3: Mean Matching accuracy on HPatches dataset. Proposed self-supervised methods suffixed with –U perform comparable or better than supervised counterparts suffixed with -S



Table 1: (Up) Results show superior performance of in-batch over in-pair negative mining. (Bottom) Selfsupervised methods perform at par with supervised counterparts on Aachen vislocalization benchmark

	Method	Supervision	Training data	Aachen v1.1 % localized queries					
		Duportision		Day (824 images)			Night (191 images)		
				$0.25m, 2^{\circ}$	$0.5m, 5^{\circ}$	5 <i>m</i> , 10°	$0.25m, 2^{\circ}$	0.5 <i>m</i> , 5°	5 <i>m</i> , 10°
Super	R2D2 [57]	OF	A+R	88.6	95.4	98.9	72.8	89.0	97.4
	R2D2*	OF	Α	87.7	94.7	98.7	69.6	86.4	95.3
	<b>CAPS</b> [76]	SL+RP	Μ	85.3	93.8	97.9	75.9	88.5	97.9
	R2D2-())	-	Α	87.4	94.9	98.3	63.9	80.1	92.1
Self-supervised	R2D2-(∰, □)	-	Α	88.0	94.8	98.2	70.2	86.4	95.8
	R2D2-(⊞, □)	-	Μ	87.4	94.7	98.3	72.3	88.5	97.4
	R2D2-(⊞, <b>□</b> )-gd	-	Μ	87.5	94.9	98.3	71.7	86.4	96.9
	R2D2-(Ⅲ, □)-selfgd	-	Μ	88.1	94.8	98.1	71.2	88.0	95.8
	R2D2-(⊞, <b>□</b> )	-	M+P	88.2	95.1	98.5	73.3	90.1	97.4
	CAPS-(I)	-	М	85.8	93.8	98.2	67.0	82.2	96.9
	CAPS-(Ⅲ, □)	-	Μ	85.1	93.2	97.8	71.7	87.4	<b>97.9</b>
	CAPS-(⊞, □)-gd	-	Μ	87.0	<b>93.8</b>	<b>98.3</b>	73.8	<b>89.0</b>	97.4
	CAPS-(Ⅲ, □)-selfgd	-	Μ	86.9	<b>93.8</b>	98.1	71.7	<b>89.0</b>	97.4
	CAPS-(⊞, <b>□</b> )	-	M+P	85.4	93.2	97.9	72.3	88.5	97.9



### Results



9	Negative sampling type						
	in-pair	in-batch					
М	55.38	58.69					
Н	29.67	33.17					
М	[94.29, 86.57, 78.43]	[95.71, 89.71, 83.29]					
Н	[82.86, 54.57, 42.29]	[85.71, 60.29, 45.71]					
 1px	0.239 / 0.425 / 0.332	0.254 / 0.439 / 0.346					
3px	0.585 / 0.677 / 0.631	0.630 / 0.707 / 0.669					
5px	0.648 / 0.742 / 0.695	0.706 / 0.784 / 0.745					
day	87.9 / 94.2 / 97.9	88.2 / 95.5 / 98.7					
night	66.5 / 79.1 / 91.6	68.1 / 83.8 / 94.8					