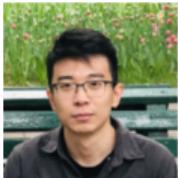


# SOSNET

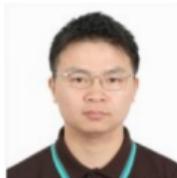
Second Order Similarity Regularization for Local Descriptor Learning



Yurun  
Tian



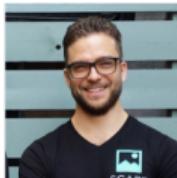
Xin  
Yu



Bin  
Fan



Fuchao  
Wu



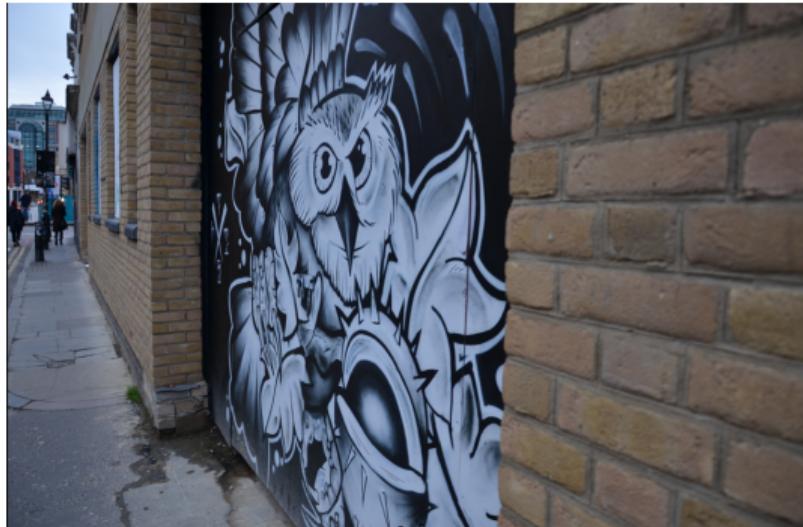
Huub  
Heijnen



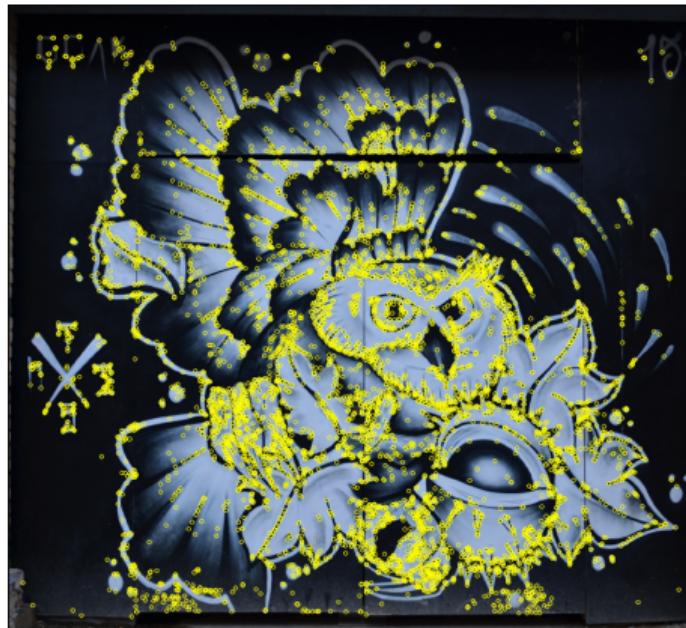
Vassileios  
Balntas

<http://research.scape.io/sosnet>  
20 June 2019

Problem: Matching two images from different viewpoints, under changing conditions.

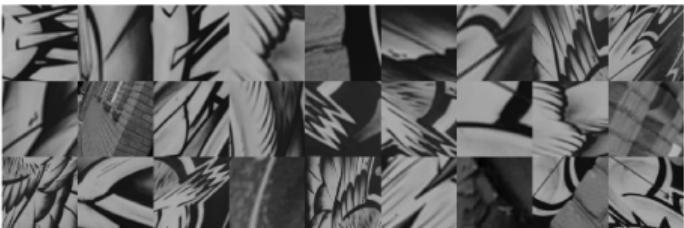


## Step 1: Extract Keypoints



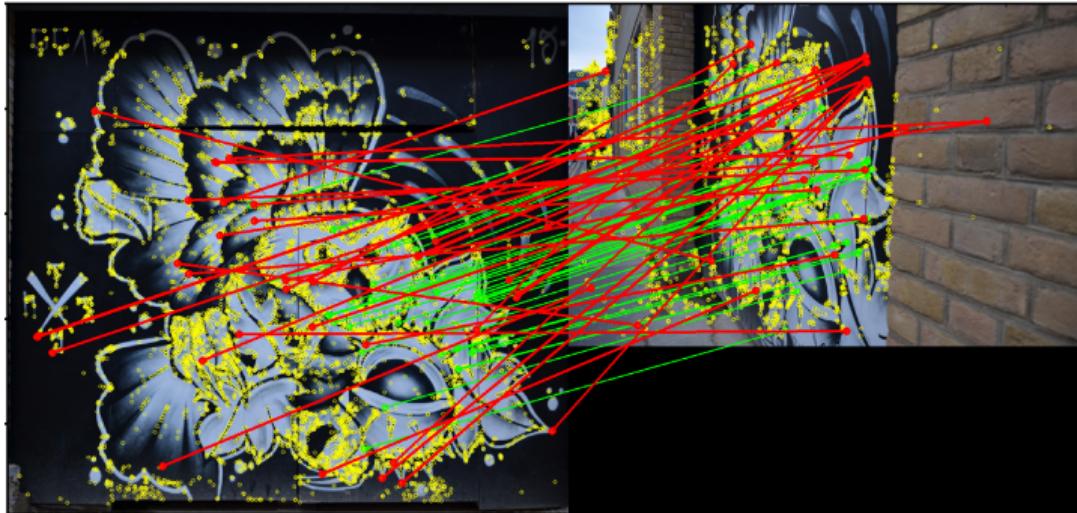
## Step 2: From Keypoints to Descriptors

Keypoints → Patches → Descriptors



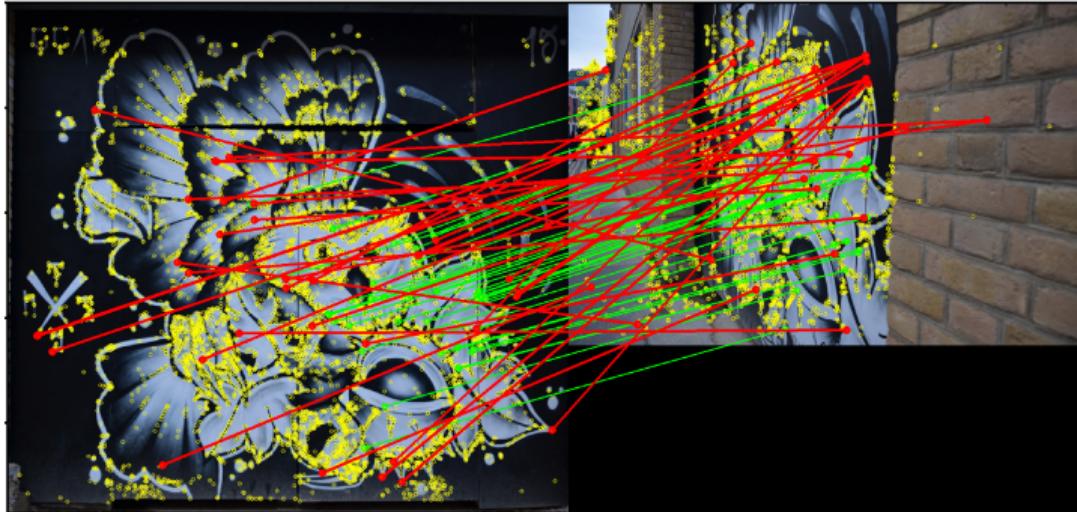
# Challenges

SIFT # Total Matches: 76 # Correct Matches: 40



## Challenges

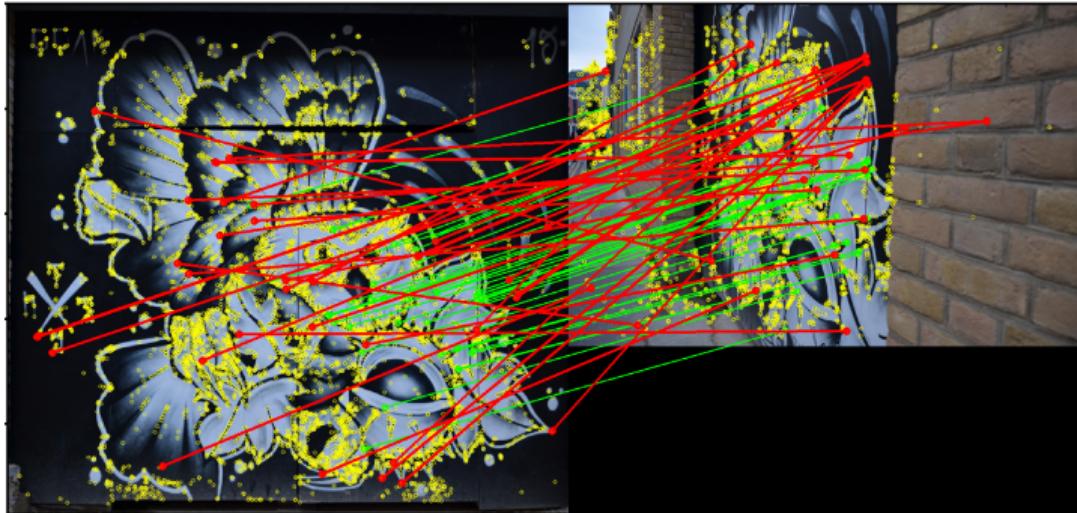
SIFT # Total Matches: 76 # Correct Matches: 40



- ▶ Keypoint repeatability

# Challenges

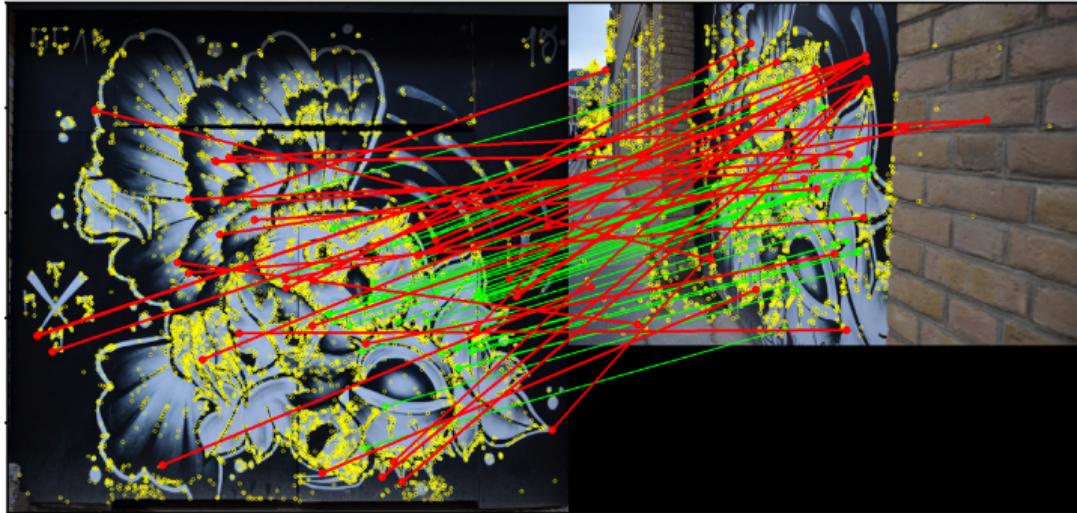
SIFT # Total Matches: 76 # Correct Matches: 40



- ▶ Keypoint repeatability
  - ▶ Design more robust keypoint detection methods

# Challenges

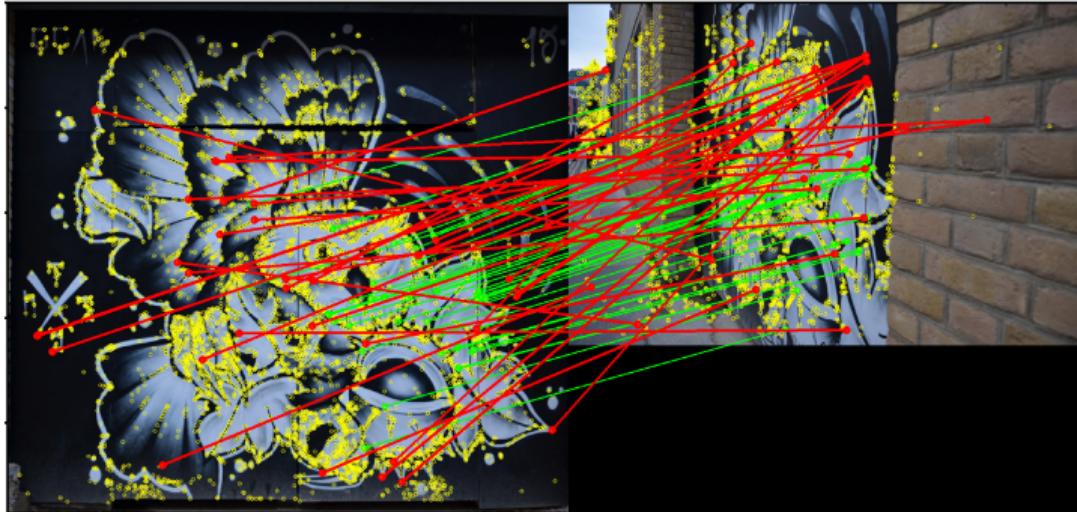
SIFT # Total Matches: 76 # Correct Matches: 40



- ▶ Keypoint repeatability
  - ▶ Design more robust keypoint detection methods
- ▶ Descriptor robustness

# Challenges

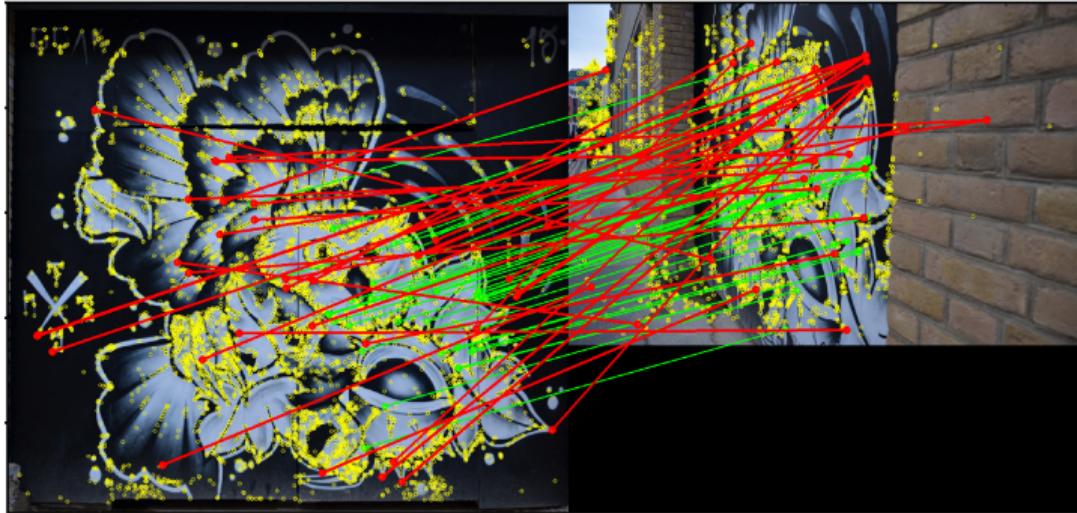
SIFT # Total Matches: 76 # Correct Matches: 40



- ▶ Keypoint repeatability
  - ▶ Design more robust keypoint detection methods
- ▶ Descriptor robustness
  - ▶ Design more robust patch description methods

# Challenges

SIFT # Total Matches: 76 # Correct Matches: 40



- ▶ Keypoint repeatability
  - ▶ Design more robust keypoint detection methods
- ▶ Descriptor robustness
  - ▶ Design more robust patch description methods [[SOSNet](#)]

## SOSNet: Methodology

We explore the potential of Second Order Similarity (SOS) in descriptor learning

- ▶ Previous work: optimise First Order Similarity (FOS)
  - ▶ Distances between descriptors
- ▶ This Work: Besides FOS, we can utilise Second Order Similarity (SOS) information to improve matching.
  - ▶ Distances between FOS distances
- ▶ FOS and SOS can be jointly optimized.

$$\mathcal{L} = \mathcal{L}_{\text{FOS}} + \mathcal{R}_{\text{sos}}$$

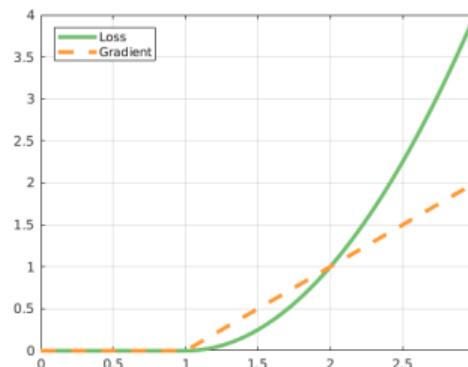
# Optimization of First Order Similarity (FOS)

Quadratic Hinge Triplet (QHT) Loss + hard sample mining as in HardNet<sup>1</sup>

$$\mathcal{L}_{\text{FOS}} = \frac{1}{N} \sum_{i=1}^N \max \left[ (0, t + d_i^{\text{pos}} - d_i^{\text{neg}})^2 \right], \quad (\text{QHT})$$

$$d_i^{\text{pos}} = d(\mathbf{x}_i, \mathbf{x}_i^+),$$

$$d_i^{\text{neg}} = \min_{\forall j, j \neq i} (d(\mathbf{x}_i, \mathbf{x}_j), d(\mathbf{x}_i, \mathbf{x}_j^+), d(\mathbf{x}_i^+, \mathbf{x}_j), d(\mathbf{x}_i^+, \mathbf{x}_j^+))$$



<sup>1</sup>Mishchuk et al., "Working hard to know your neighbor's margins: Local descriptor learning loss".

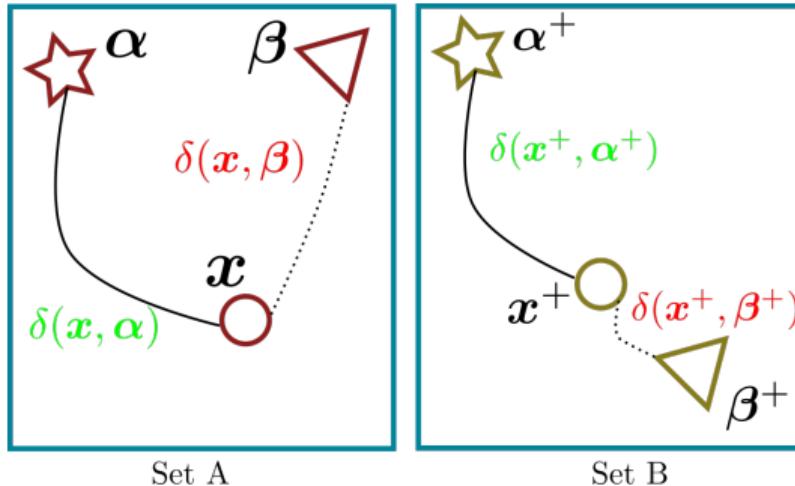
## Optimization of Second Order Similarity(SOS)

- ▶ Matching descriptors should have high SOS.

$$\mathcal{R}_{\text{SOS}} = \frac{1}{N} \sum_{i=1}^N d^{(2)}(\mathbf{x}_i, \mathbf{x}_i^+).$$

- ▶ We optimise second-order similarity by optimising the distance between distances.

## Optimization of Second Order Similarity(SOS)



Set A

Set B

$$\mathbf{y} = [\delta(\mathbf{x}, \boldsymbol{\alpha}), \delta(\mathbf{x}, \boldsymbol{\beta})]$$

$$\mathbf{y}^+ = [\delta(\mathbf{x}^+, \boldsymbol{\alpha}^+), \delta(\mathbf{x}^+, \boldsymbol{\beta}^+)]$$

$$d^{(2)}(\mathbf{x}, \mathbf{x}^+) = \|\mathbf{y} - \mathbf{y}^+\|$$

↓

$$d^{(2)}(\mathbf{x}, \mathbf{x}^+) = \sqrt{\{\delta(\mathbf{x}, \boldsymbol{\alpha}) - \delta(\mathbf{x}^+, \boldsymbol{\alpha}^+)\}^2 + \{\delta(\mathbf{x}, \boldsymbol{\beta}) - \delta(\mathbf{x}^+, \boldsymbol{\beta}^+)\}^2}$$

## kNN search strategy for the selection of reference descriptors

$$\mathcal{R}_{\text{sos}} = \frac{1}{N} \sum_{i=1}^N d^{(2)}(\mathbf{x}_i, \mathbf{x}_i^+).$$

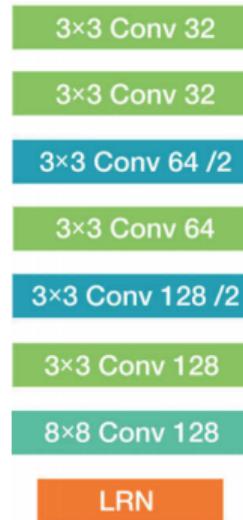
$$d^{(2)}(\mathbf{x}_i, \mathbf{x}_i^+) = \sqrt{\sum_{j \neq i}^N (d(\mathbf{x}_i, \mathbf{x}_j) - d(\mathbf{x}_i^+, \mathbf{x}_j^+))^2}$$

$$\mathbf{c}_i = \{z_j : \mathbf{x}_i \in K\text{NN}(\mathbf{x}_j) \vee \mathbf{x}_i^+ \in K\text{NN}(\mathbf{x}_j^+)\},$$
$$\forall j \in 1 \dots N, j \neq i$$



## Implementation

We use the network architecture of L2-Net<sup>2</sup>. The network is trained with Adam optimizer.



---

<sup>2</sup>Tian, Fan, and Wu, "L2-Net: Deep Learning of Discriminative Patch Descriptor in Euclidean Space".

# Performance

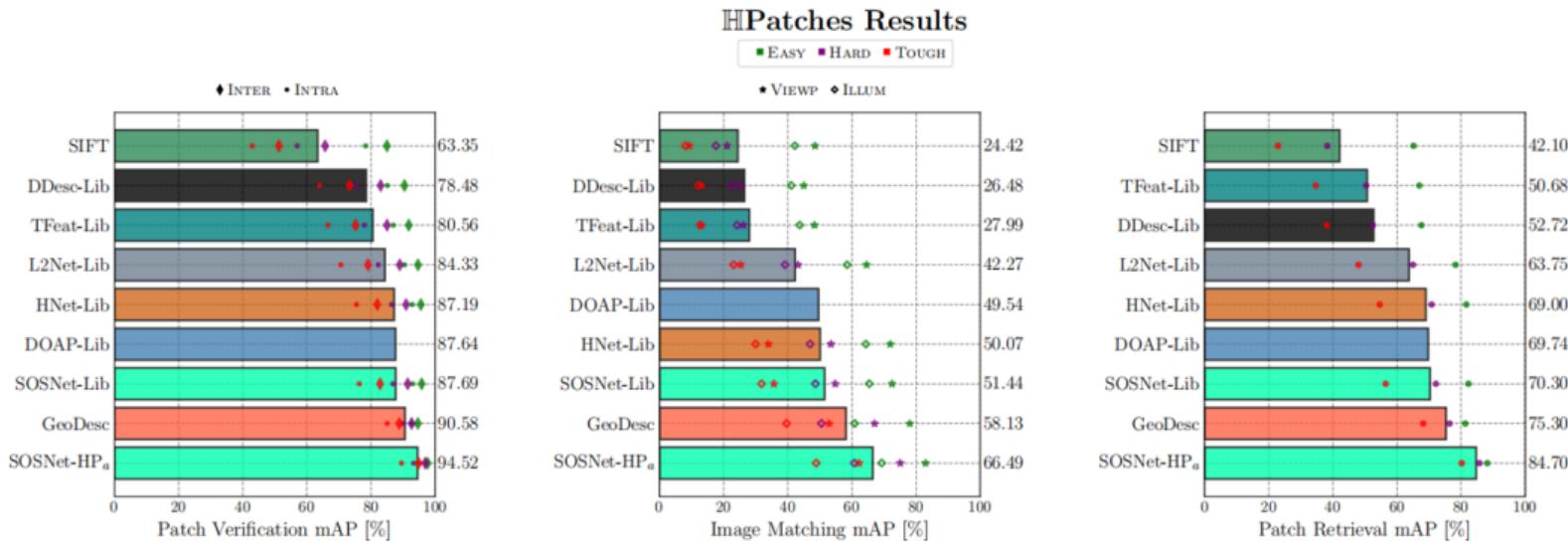
Performance of SOSNet on UBC dataset<sup>3</sup>.

Train	Notredame	Yosemite	Liberty	Yosemite	Liberty	Notredame	Mean
Test	Liberty		Notredame		Yosemite		
SIFT [20]	29.84		22.53		27.29		26.55
DeepDesc [32]	10.9		4.40		5.69		6.99
MatchNet [2]	7.04	11.47	3.82	5.65	11.6	8.70	8.05
L2Net [34]	3.64	5.29	1.15	1.62	4.43	3.30	3.24
CS L2Net [34]	2.55	4.24	0.87	1.39	3.81	2.84	2.61
HardNet [25]	1.47	<b>2.67</b>	0.62	0.88	2.14	1.65	1.57
HardNet-GOR [25, 39]	1.72	2.89	0.63	0.91	2.10	1.59	1.64
Michel <i>et al.</i> [16]	1.79	2.96	0.68	1.02	2.51	1.64	1.77
<b>SOSNet</b>	<b>1.25</b>	2.84	<b>0.58</b>	<b>0.87</b>	<b>1.95</b>	<b>1.25</b>	<b>1.46</b>
TFeat+ [2]	7.39	10.13	3.06	3.80	8.06	7.24	6.64
L2Net+ [34]	2.36	4.70	0.72	1.29	2.57	1.71	2.23
CS L2Net+ [34]	1.71	3.87	0.56	1.09	2.07	1.3	1.76
HardNet+ [25]	1.49	2.51	0.53	0.78	1.96	1.84	1.51
HardNet-GOR+ [25, 39]	1.48	2.43	0.51	0.78	1.76	1.53	1.41
DOAP+ [14]	1.54	2.62	0.43	0.87	2.00	1.21	1.45
DOAP-ST+ [14] [15]	1.47	2.29	0.39	0.78	1.98	1.35	1.38
<b>SOSNet+</b>	<b>1.08</b>	<b>2.12</b>	<b>0.35</b>	<b>0.67</b>	<b>1.03</b>	<b>0.95</b>	<b>1.03</b>
GeoDesc+ [21]	5.47		1.94		4.72		4.05
SOSNet-HP+	2.10		0.79		1.39		1.42

<sup>3</sup>Winder, Hua, and Brown, "Picking the best daisy".

# Performance

Performance of SOSNet on HPatches dataset<sup>4</sup>.



<sup>4</sup>Balntas et al., "HPatches: A Benchmark and Evaluation of Handcrafted and Learned Local Descriptors".

# Performance

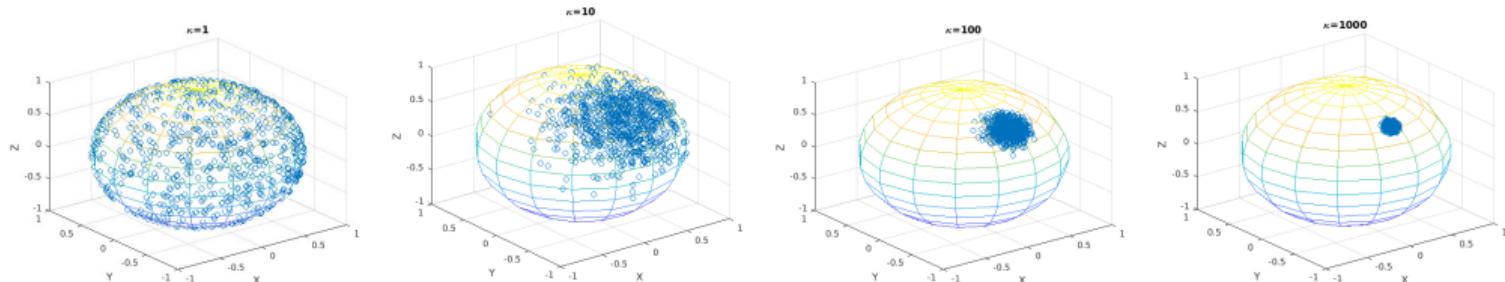
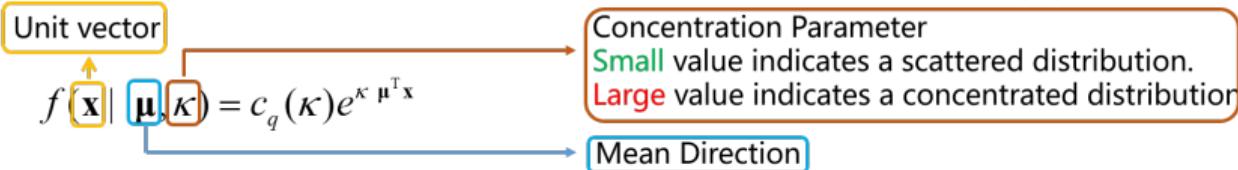
Performance of SOSNet on ETH dataset<sup>5</sup>.

		# Image	# Registered	# Sparse Points	# Observations	Track Length	Reproj. Error
<b>Fountain</b>	SIFT	11	11	14K	70K	4.79	0.39px
	DSP-SIFT	11		14K	71K	4.78	<b>0.37px</b>
	L2Net	11		17K	83K	4.88	0.47px
	GeoDesc	11		16K	83K	<b>5.00</b>	0.47px
	<b>SOSNet</b>	11		<b>17K</b>	<b>85K</b>	4.92	0.43px
<b>Herzjesu</b>	SIFT	8	8	7.5K	31K	4.22	<b>0.43px</b>
	DSP-SIFT	8		7.7K	32K	4.22	0.45px
	L2Net	8		9.5K	40K	4.24	0.51px
	Geodesc	8		9.2K	40K	<b>4.35</b>	0.51px
	<b>SOSNet</b>	8		<b>9.7K</b>	<b>41K</b>	4.26	0.53px
<b>South Building</b>	SIFT	128	128	108K	653K	<b>6.04</b>	<b>0.54px</b>
	DSP-SIFT	128		112K	666K	5.91	0.58px
	L2Net	128		170K	863K	5.07	0.63px
	GeoDesc	128		170K	887K	5.21	0.64px
	<b>SOSNet</b>	128		<b>178K</b>	<b>913K</b>	5.11	0.67px
<b>Madrid Metropolis</b>	SIFT	1344	500	116K	733K	6.32	<b>0.60px</b>
	DSP-SIFT		467	99K	649K	<b>6.52</b>	0.66px
	L2Net		692	254k	1067K	4.20	0.69px
	GeoDesc		809	306K	1200K	3.91	0.66px
	<b>SOSNet</b>		<b>844</b>	<b>335K</b>	<b>1411K</b>	4.21	0.70px
<b>Gendarmenmarkt</b>	SIFT	1463	1035	338K	1872K	<b>5.523</b>	<b>0.69px</b>
	DSP-SIFT		979	293K	1577K	5.381	0.74px
	L2Net		1168	667k	2611K	3.91	0.73px
	GeoDesc		<b>1208</b>	<b>770K</b>	<b>2903K</b>	3.72	0.74px
	<b>SOSNet</b>		1201	<b>816K</b>	<b>3255K</b>	3.984	0.77px

<sup>5</sup>Schönberger et al., “Comparative Evaluation of Hand-Crafted and Learned Local Features”.

# Analysis

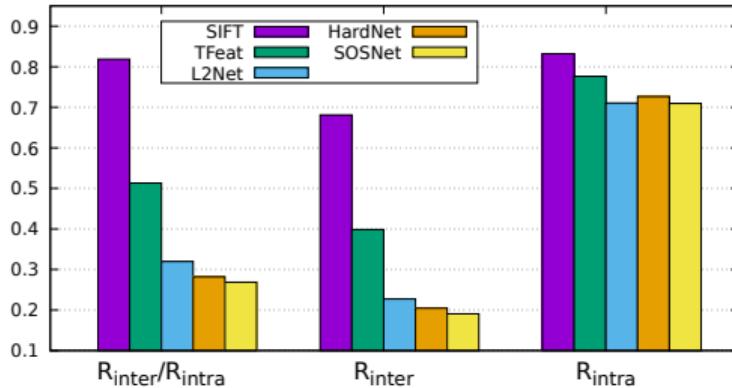
We define a von Mises-Fisher distribution performance measure ( $\rho$ )



Distribution of class centers	Distribution of elements inside a class
$k_{\text{inter}}$	$k_{\text{intra}}$
↓	↑

$$\rho = \frac{k_{\text{inter}}}{k_{\text{intra}}} \propto \frac{R_{\text{inter}}}{R_{\text{intra}}}$$

# Performance



- ▶ Observations
  - ▶ As performance increases, ratio  $\rho$  and  $R_{\text{inter}}$  drop monotonically;
  - ▶ SOSR helps to exploit more space of the hypersphere, Intra-class distribution may also become more scattered;
  - ▶ SIFT descriptors(normalized) tend to gather together in a small region of the hypersphere.
- ▶ Our SOS regularisation term, might potentially be useful in other metric learning problems

# The End

- ▶ Please visit our poster at 3.2 150 (15:20-18:00)
- ▶ Pre-trained network & code available

<http://research.scape.io/sosnet>

THANK YOU