

Overview

- ▶ We consider kernelized patch encodings \rightarrow vectorized patch representation
- ▶ Revisiting the SIFT descriptor with match kernels
- ▶ Jointly encode pixel gradient and position continuously (no quantization)
- ▶ Fast dominant orientation alignment without descriptor re-computation

Match kernels

Local descriptor representation: set of pixel attributes

$$\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_i, \dots\}, \mathbf{x}_i \in \mathbb{R}^d$$

e.g. $d = 1$ for grayscale patches; $d = 4$ in this work

Match kernel to compare two patches

$$K(\mathcal{X}, \mathcal{Y}) = \beta(\mathcal{X})\beta(\mathcal{Y}) \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{y} \in \mathcal{Y}} k(\mathbf{x}, \mathbf{y})$$

Vector mapping $\psi: \mathbb{R}^d \rightarrow \mathbb{R}^D$, such that $k(\mathbf{x}, \mathbf{y}) = \langle \psi(\mathbf{x}) | \psi(\mathbf{y}) \rangle$

Patch representation:

$$\mathbf{X} = \beta(\mathcal{X}) \sum_{\mathbf{x} \in \mathcal{X}} \psi(\mathbf{x}), \quad (\text{such that } \|\mathbf{X}\| = 1)$$

Match kernel computation via inner product

$$K(\mathcal{X}, \mathcal{Y}) = \beta(\mathcal{X})\beta(\mathcal{Y}) \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{y} \in \mathcal{Y}} \langle \psi(\mathbf{x}) | \psi(\mathbf{y}) \rangle = \langle \mathbf{X} | \mathbf{Y} \rangle$$

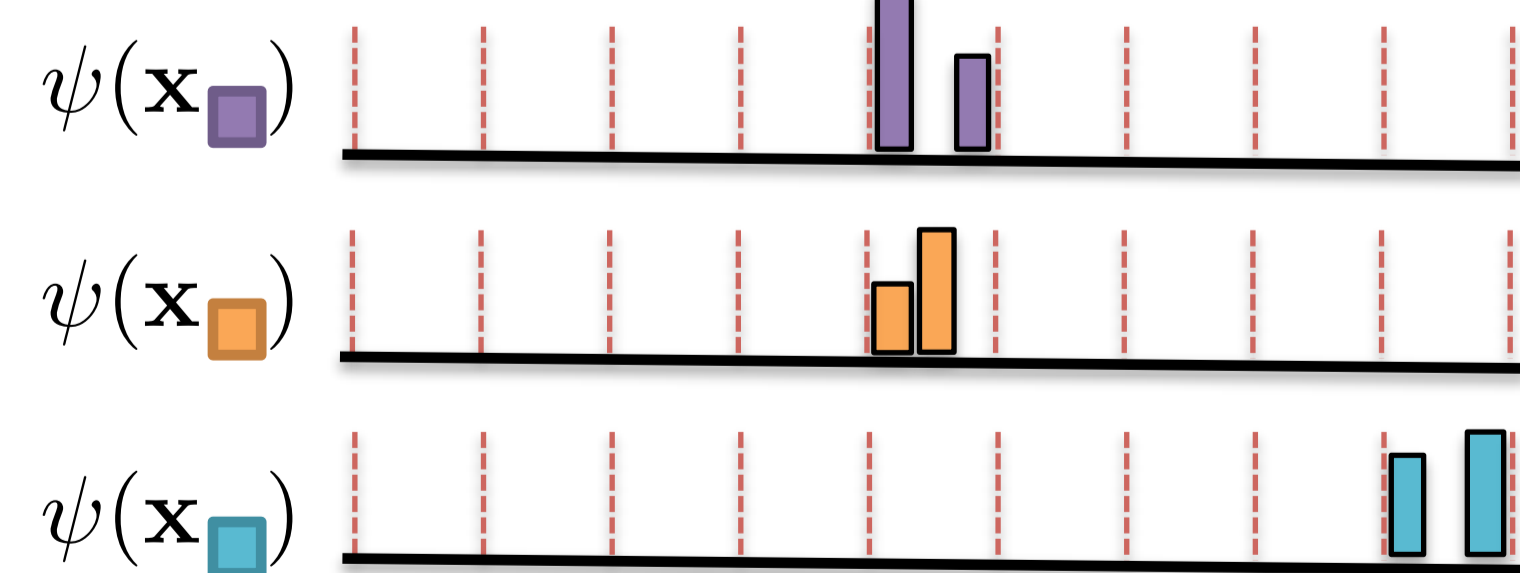
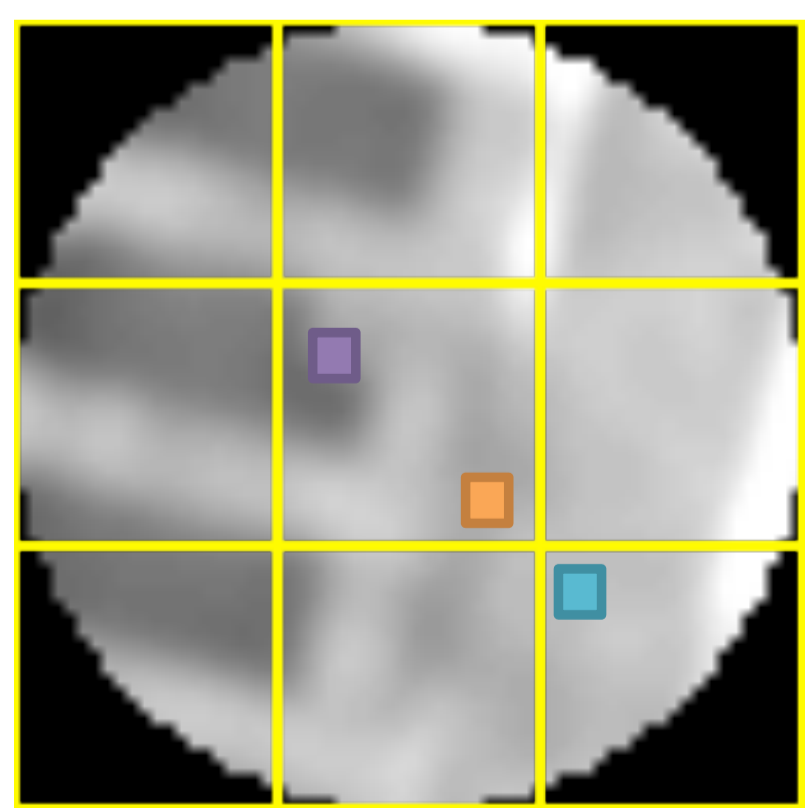
We seek to

Identify a kernel function for patch elements that reflects their resemblance in terms of gradients and their proximity in terms of spatial position. Create vectorized local descriptor representation aggregating multiple pixel attributes, such that inner product of two such patch representations approximates a match kernel of the form:

$$K(\mathbf{X}^*, \mathbf{Y}^*) = \beta(\mathcal{X}^*)\beta(\mathcal{Y}^*) \sum_{\mathbf{x} \in \mathcal{X}^*} \sum_{\mathbf{y} \in \mathcal{Y}^*} \tilde{m}_x \tilde{m}_y k_\theta(\theta_x, \theta_y) k_\varphi(\varphi_x, \varphi_y) k_\rho(\tilde{\rho}_x, \tilde{\rho}_y) \approx \langle \mathbf{X}^* | \mathbf{Y}^* \rangle \quad (1)$$

- ▶ (m_x, θ_x) : magnitude and gradient at pixel \mathbf{x}
- ▶ (φ_x, ρ_x) : polar coordinates of pixel \mathbf{x} with respect to the patch center
- ▶ $k_{\theta, \varphi, \rho}$: functions measuring consistency in gradients orientation, in positions on coordinate φ and on ρ , respectively

SIFT: kernel view



gradient sampling from grid

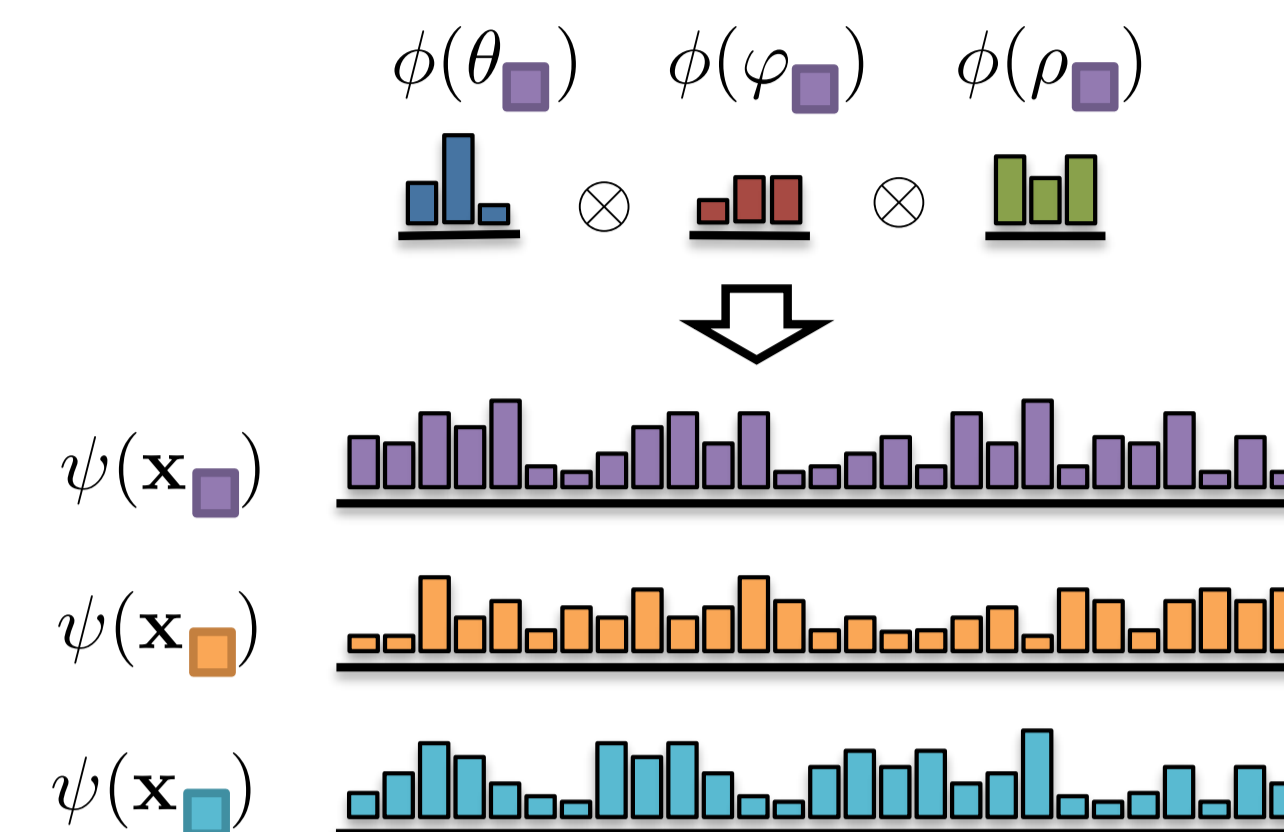
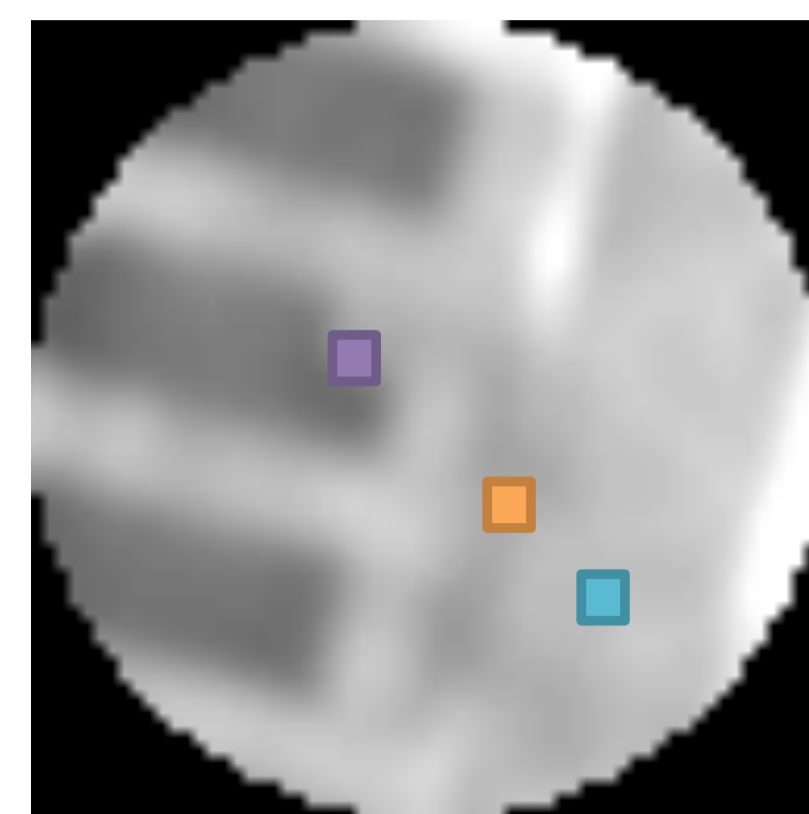
pixel gradient embedding

- ▶ Matching pixels with relatively similar positions and gradients orientations
- ▶ The hard assignment in the quantization process inserts some artifacts
- ▶ Losses in the selectivity of the similarity function

Method

Our approach

- ▶ Individual pixel sampling
- ▶ Continuous embedding of pixel attributes
- ▶ Aggregation into dense pixel description

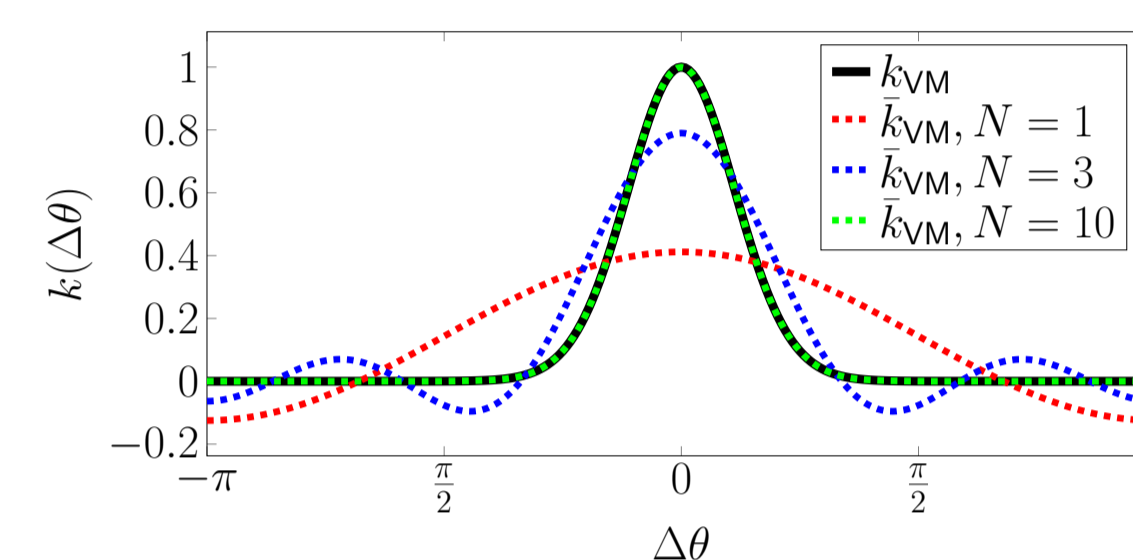


Weighting function for relative angles (normalized Von Mises)

$$k_{VM}(\Delta\theta) = \frac{\exp(\kappa \cos(\Delta\theta)) - \exp(-\kappa)}{2 \sinh(\kappa)}$$

Fourier series approximation:

$$\tilde{k}_{VM}^N(\Delta\theta) = \sum_{n=0}^N \gamma_n \cos(n\Delta\theta)$$



Feature map for angle

Define a mapping from angle θ to a vector $\phi(\theta): [-\pi, \pi] \rightarrow \mathbb{R}^{2N+1}$

N : number of truncated components/frequencies from the Fourier series approximation

$$\phi(\theta) = (\sqrt{\gamma_0}, \sqrt{\gamma_1} \cos(\theta), \sqrt{\gamma_1} \sin(\theta), \dots, \sqrt{\gamma_N} \cos(N\theta), \sqrt{\gamma_N} \sin(N\theta))^T \quad (2)$$

Compare vectors of angle via inner product \rightarrow approximate angle similarity

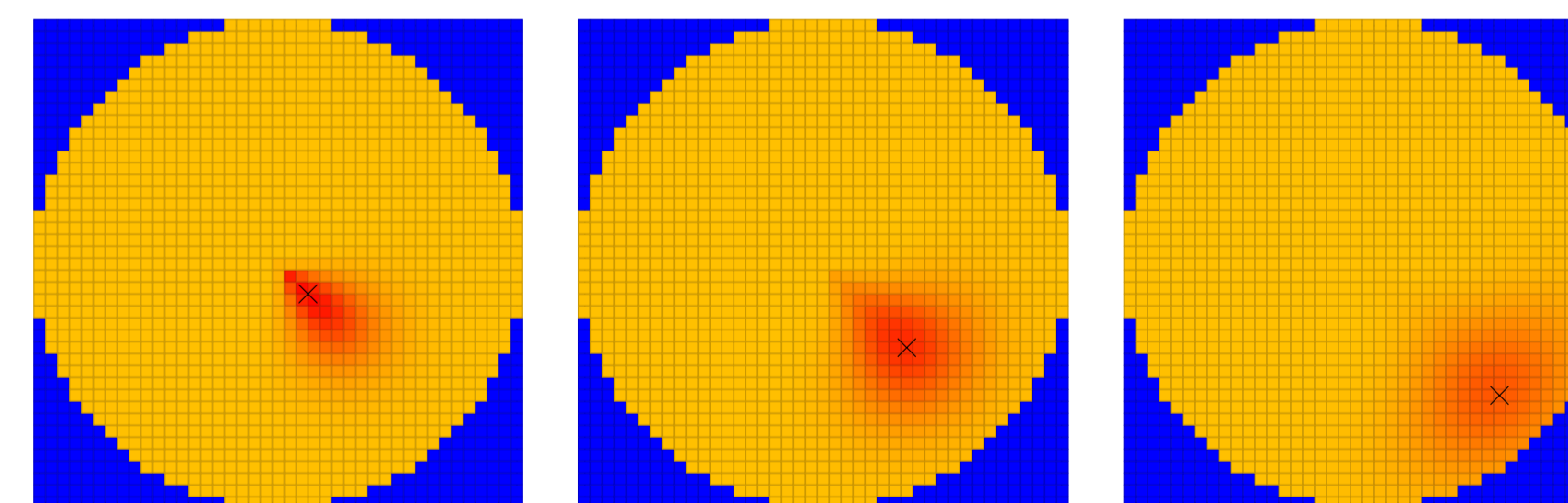
$$\begin{aligned} \phi(\theta_1)^T \phi(\theta_2) &= \gamma_0 + \sum_{n=1}^N \gamma_n (\cos(n\theta_1) \cos(n\theta_2) + \sin(n\theta_1) \sin(n\theta_2)) \\ &= \sum_{n=0}^N \gamma_n \cos(n(\theta_1 - \theta_2)) = \tilde{k}_{VM}^N(\theta_1 - \theta_2) \approx k_{VM}(\theta_1 - \theta_2) \end{aligned}$$

Embedding and aggregating pixel attributes

θ and φ are angles and are trivially mapped with (2).

The radius ρ of each pixel \mathbf{x} is mapped to an angle by $\tilde{\rho}_x = \rho_x \pi$, with $\rho_x \in [0, 1]$ and then embedded with (2).

2D weighting function for 3 sample pixels. Colors reflect spatial similarity to other pixels. Red corresponds to maximum similarity



Final local descriptor representation:

$$\mathbf{X}^* = \beta(\mathcal{X}^*) \sum_{\mathbf{x} \in \mathcal{X}^*} \psi(\mathbf{x}) = \beta(\mathcal{X}^*) \sum_{\mathbf{x} \in \mathcal{X}^*} \tilde{m}_x \phi(\theta_x) \otimes \phi(\varphi_x) \otimes \phi(\tilde{\rho}_x)$$

$$\langle \mathbf{X}^* | \mathbf{Y}^* \rangle \propto \sum_{\mathbf{x} \in \mathcal{X}^*} \psi(\mathbf{x})^T \sum_{\mathbf{y} \in \mathcal{Y}^*} \psi(\mathbf{y}) = \sum_{\mathbf{x} \in \mathcal{X}^*} \sum_{\mathbf{y} \in \mathcal{Y}^*} \psi(\mathbf{x})^T \psi(\mathbf{y}) \approx K(\mathbf{X}^*, \mathbf{Y}^*)$$

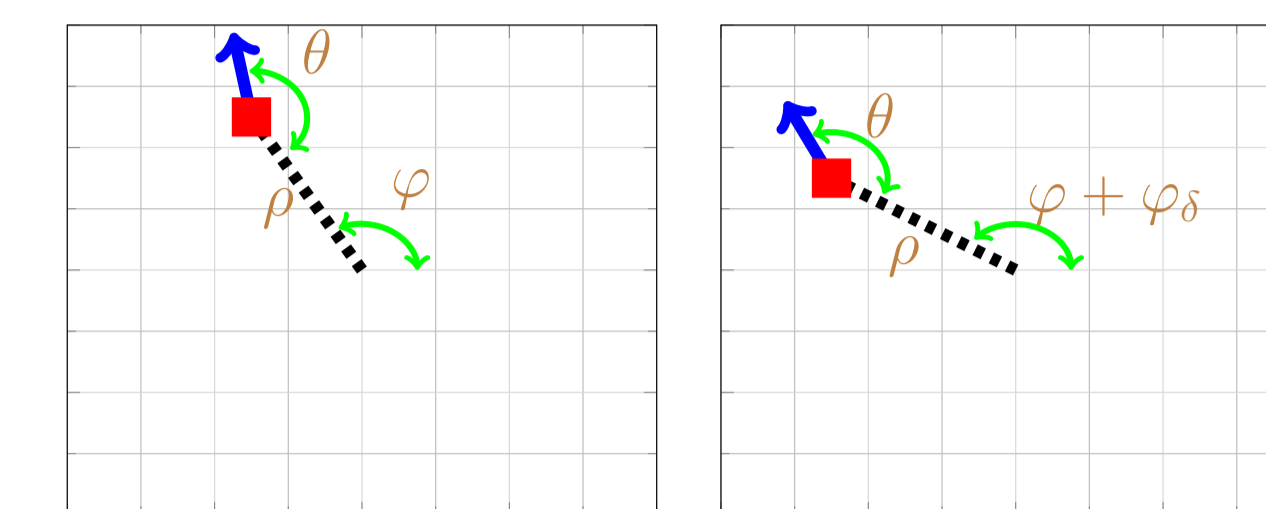
References

- ▶ L. Bo, X. Ren, and D. Fox. Kernel descriptors for visual recognition. In *NIPS*, 2010.
- ▶ G. Tolias, T. Furon, and H. Jégou. Orientation covariant aggregation of local descriptors with embeddings. In *ECCV*, 2014.
- ▶ H. Jégou, M. Douze, and C. Schmid. Hamming embedding and weak geometric consistency for large scale image search. In *ECCV*, 2008.
- ▶ A. Vedaldi and A. Zisserman. Efficient additive kernels via explicit feature maps. *Trans. PAMI*, 2012.

Rotation matching

The method assumes upright objects

Evolution of pixel attributes when rotating patch by angle φ_δ



Angle φ is the only pixel attribute that changes

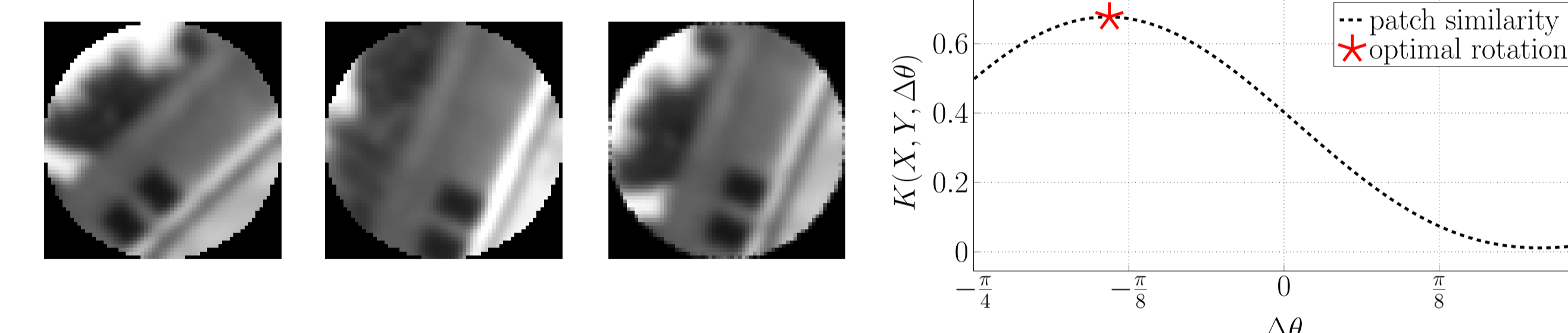
$$\mathbf{X}^* = [\mathbf{X}_0^{*\top}, \mathbf{X}_{1,c}^{*\top}, \mathbf{X}_{1,s}^{*\top}, \dots, \mathbf{X}_{N,c}^{*\top}, \mathbf{X}_{N,s}^{*\top}]^\top$$

- ▶ \mathbf{X}_0^* : associated to constant terms
- ▶ $\mathbf{X}_{N,c}^*, \mathbf{X}_{N,s}^*$: corresponding to *cosine* and *sine* terms for frequency N

$$\begin{aligned} \mathbf{X}_{\delta,n,c}^* &= \mathbf{X}_{n,c}^* \cos n\varphi_\delta + \mathbf{X}_{n,s}^* \sin n\varphi_\delta \\ \mathbf{X}_{\delta,n,s}^* &= -\mathbf{X}_{n,c}^* \sin n\varphi_\delta + \mathbf{X}_{n,s}^* \cos n\varphi_\delta \end{aligned}$$

$$\begin{aligned} \langle \mathbf{X}_\delta^* | \mathbf{Y}^* \rangle &= \langle \mathbf{X}_0^* | \mathbf{Y}_0^* \rangle + \sum_{n=1}^{N_\varphi} \cos(n\varphi_\delta) (\langle \mathbf{X}_{n,c}^* | \mathbf{Y}_{n,c}^* \rangle + \langle \mathbf{X}_{n,s}^* | \mathbf{Y}_{n,s}^* \rangle) \\ &\quad + \sum_{n=1}^{N_\varphi} \sin(n\varphi_\delta) (-\langle \mathbf{X}_{n,c}^* | \mathbf{Y}_{n,s}^* \rangle + \langle \mathbf{X}_{n,s}^* | \mathbf{Y}_{n,c}^* \rangle) \end{aligned}$$

- ▶ \mathbf{X}_δ^* : aggregated representation of rotated query
- ▶ $\langle \mathbf{X}_\delta^* | \mathbf{Y}^* \rangle$: trigonometric polynomial with coefficients independent from φ_δ



Patch A Patch B Aligned Patch A Similarities for rotations of Patch A

Experiments

Hypothesis test on Brown patch dataset

False positive rate (%) at 95% recall. Learning type: N -none, US -unsupervised, S -supervised.

Train/Test	RootSIFT	RootSIFT-PCA	Simonyan	Brown	$KD_{2,3,1}$	$KD_{3,3,1}$	$KD_{3,2,2}$ -PCA
ND Lib		19.34	12.42	16.85	21.58	20.06	13.17
Yos Lib	29.64	19.95	14.58	18.27			14.53
ND Yos		18.37	10.08	13.55	11.07	9.66	8.31
Lib Yos	26.69	19.52	11.18	N/A			9.65
Lib ND		13.90	7.22	N/A	7.99	7.00	6.36
Yos ND	22.06	13.98	6.82	11.98			6.50
Mean	26.14	17.51	10.38	15.16	13.55	12.24	9.75
Dimensions	128	80	73-77	29-36	105	147	80
Learning	N	US	S	S	N	N	US

Nearest neighbors

Recall computed at R top ranked patches for 1000 randomly selected patch queries on Notredame dataset (80d)

R	1	5	10	100	1000	10000
RootSIFT	8.5	24.4	33.0	62.9	79.7	90.6
RootSIFT-PCA	8.8	23.9	32.7	61.4	78.4	90.4
$KD_{3,3,1}$ (No rotations)	9.1	24.7	34.6	64.9	80.8	91.3
$KD_{3,3,1}$ (16 rotations)	8.8	26.2	37.3	68.3	84.4	93.1
$KD_{3,3,1}$ - PCA	9.4	24.9	35.2	66.4	82.9	92.4