

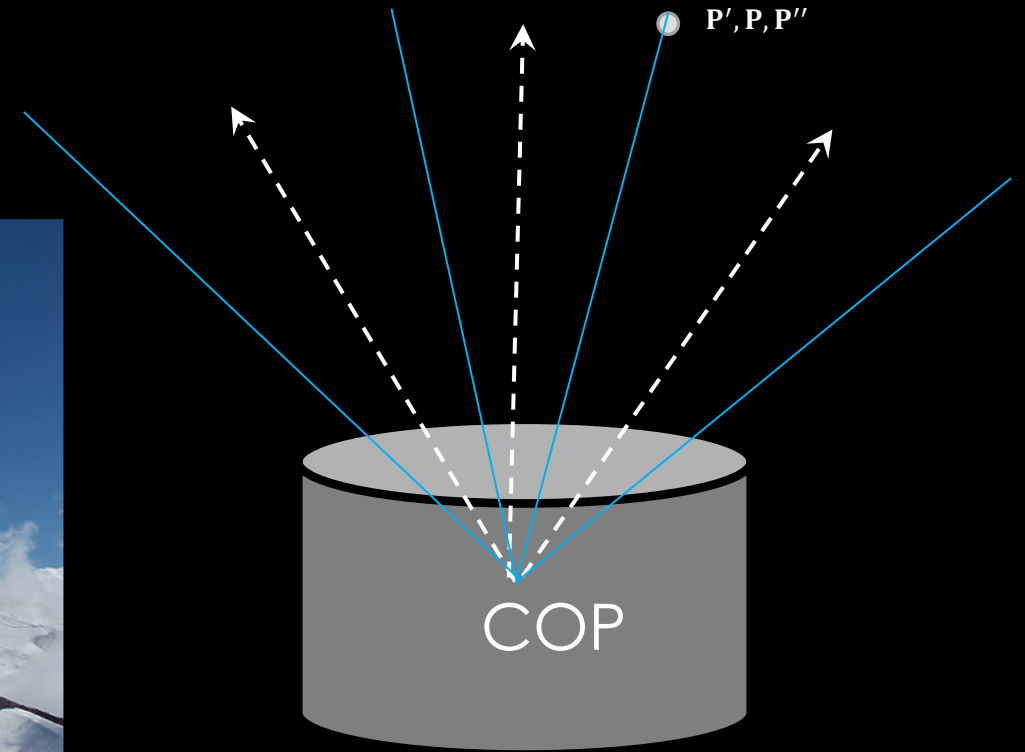
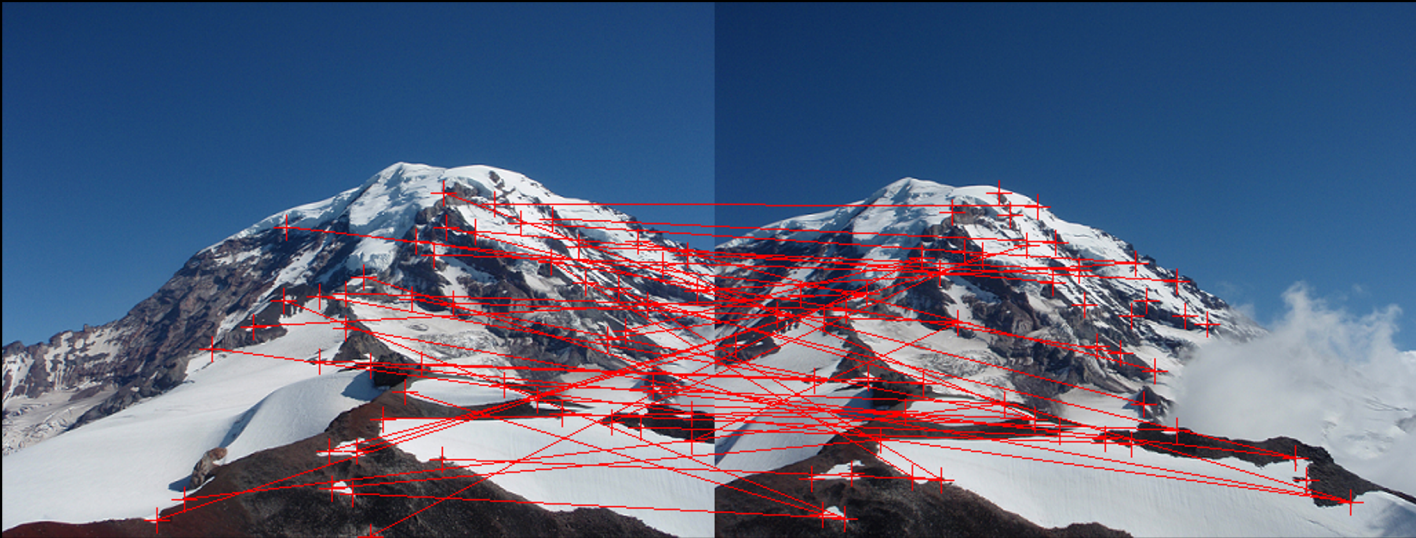
Start Recording

Recap: Long Range Image Matching and Its Applications

May 14th, 2020

Matching patches: Corners + Descriptors!

- Harris Corners + Patch Colors / Intensities as Descriptors
- Homography and Other Transformations for aligning with RANSAC
- Choosing a Projection Surface such as Cylinders
- Creating Panoramas

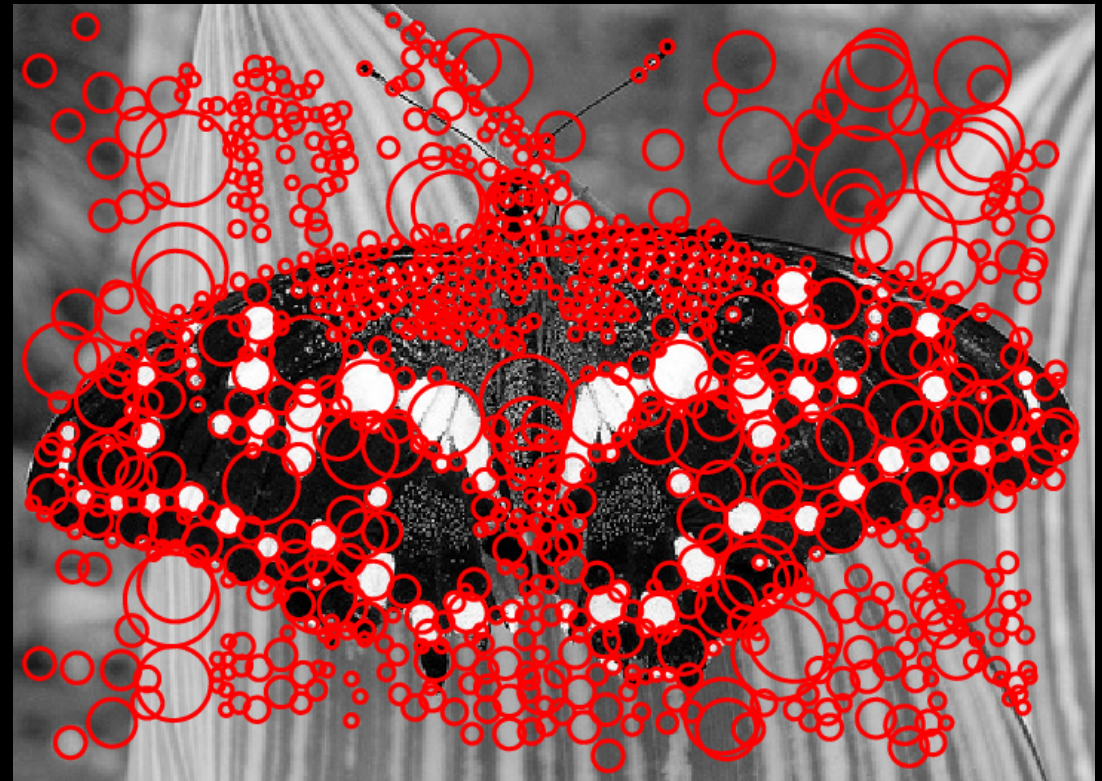
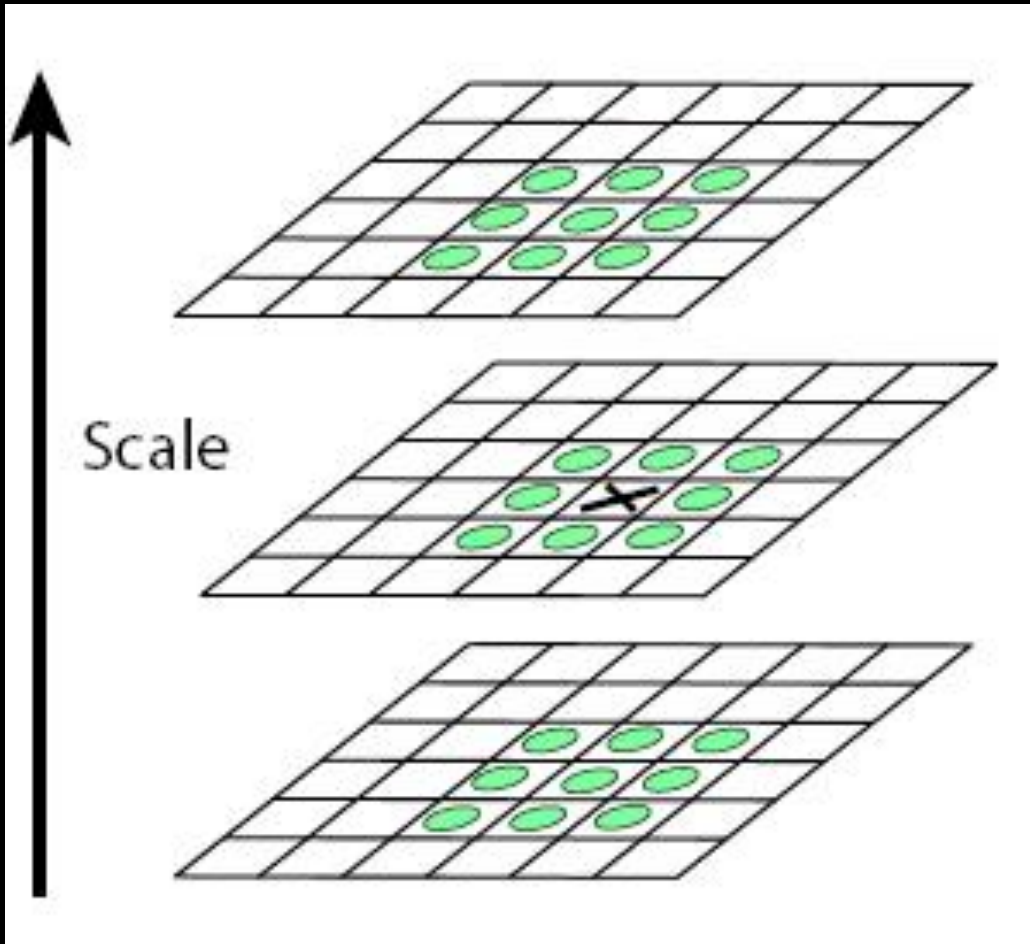


Wide Baseline Matching

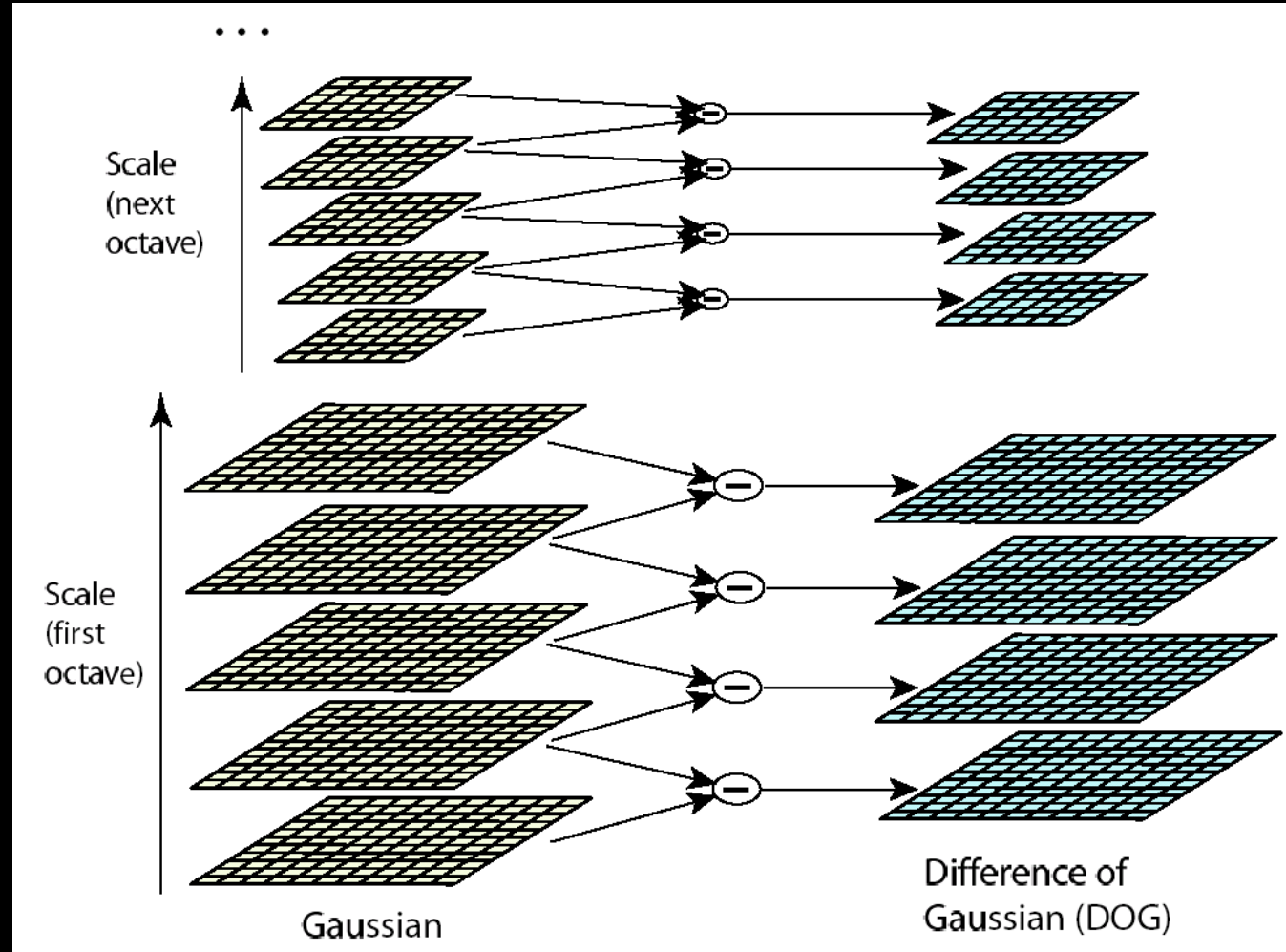


Scale-space blob detector

1. Convolve image with scale-normalized Laplacian at several scales
2. Find maxima of squared Laplacian response in scale-space

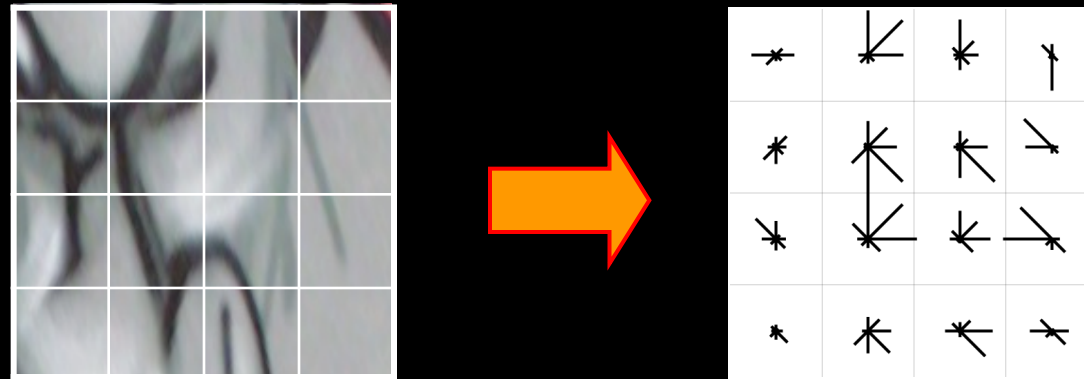


Efficient implementation



Feature descriptors: SIFT

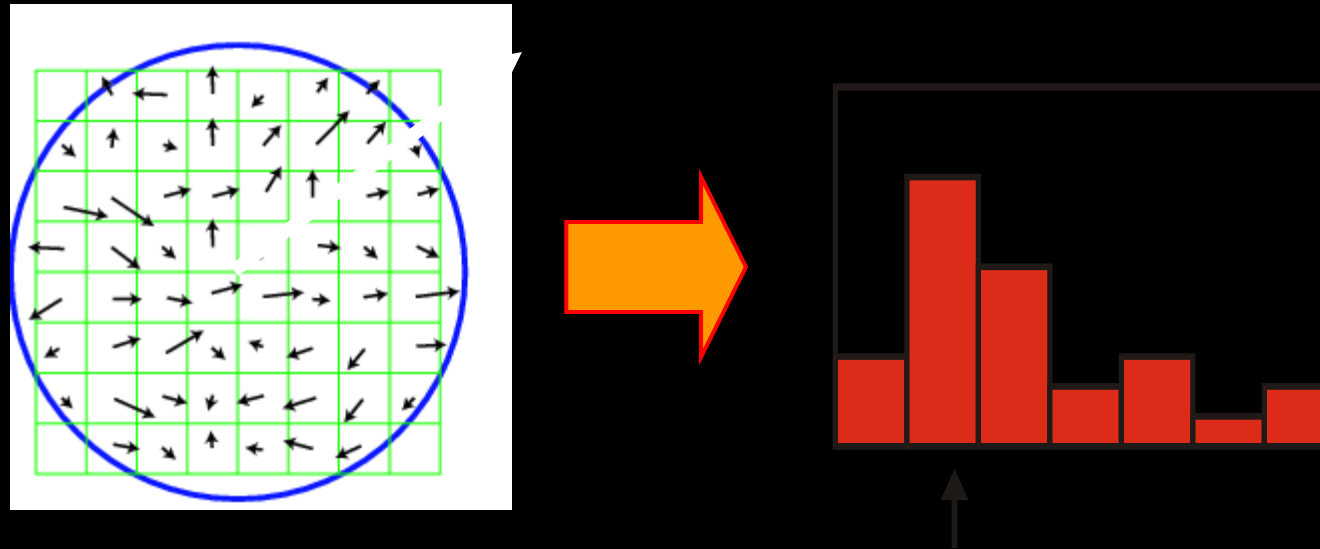
- Descriptor computation:
 - Divide patch into 4x4 sub-patches
 - Compute histogram of gradient orientations (8 reference angles) inside each sub-patch
 - Resulting descriptor: $4 \times 4 \times 8 = 128$ dimensions



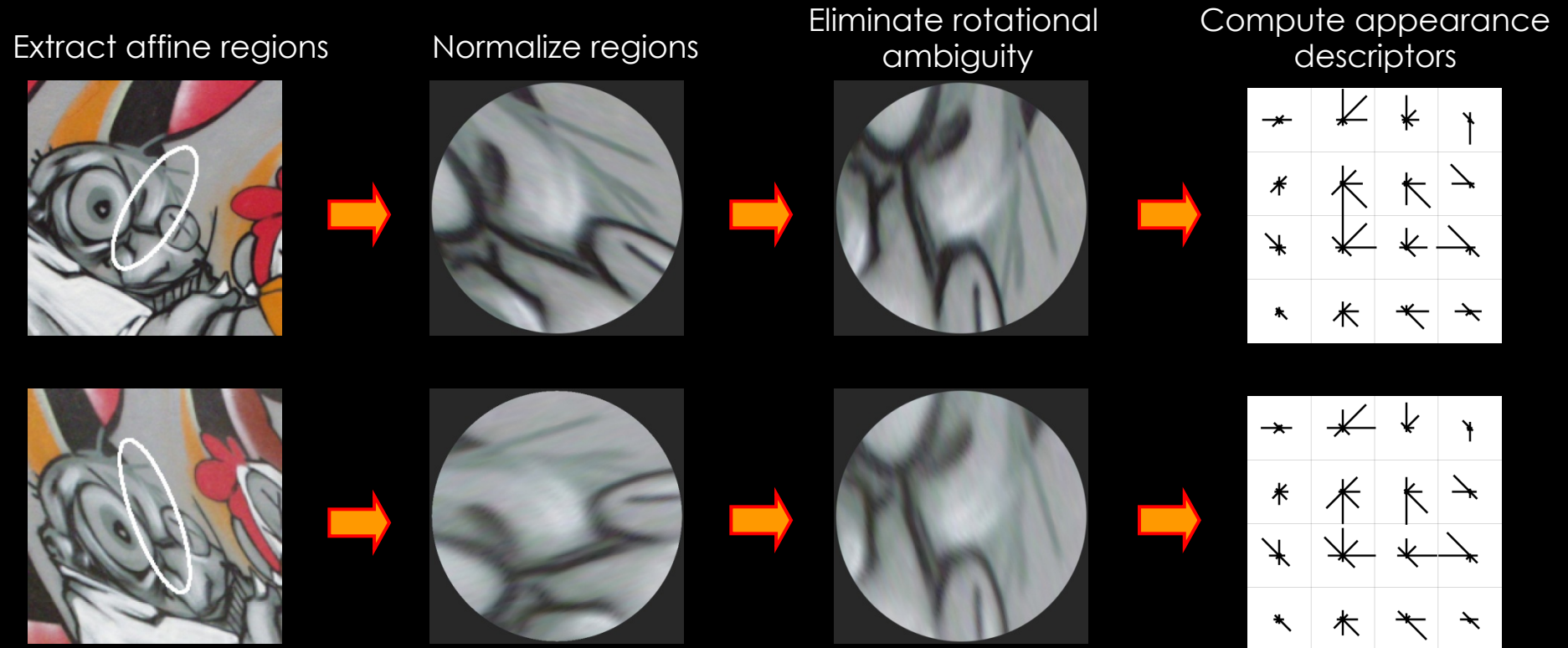
Rotational Normalization of SIFT Feature

To assign a unique orientation to circular image windows:

- Create histogram of local gradient directions in the patch
- Assign canonical orientation at peak of smoothed histogram

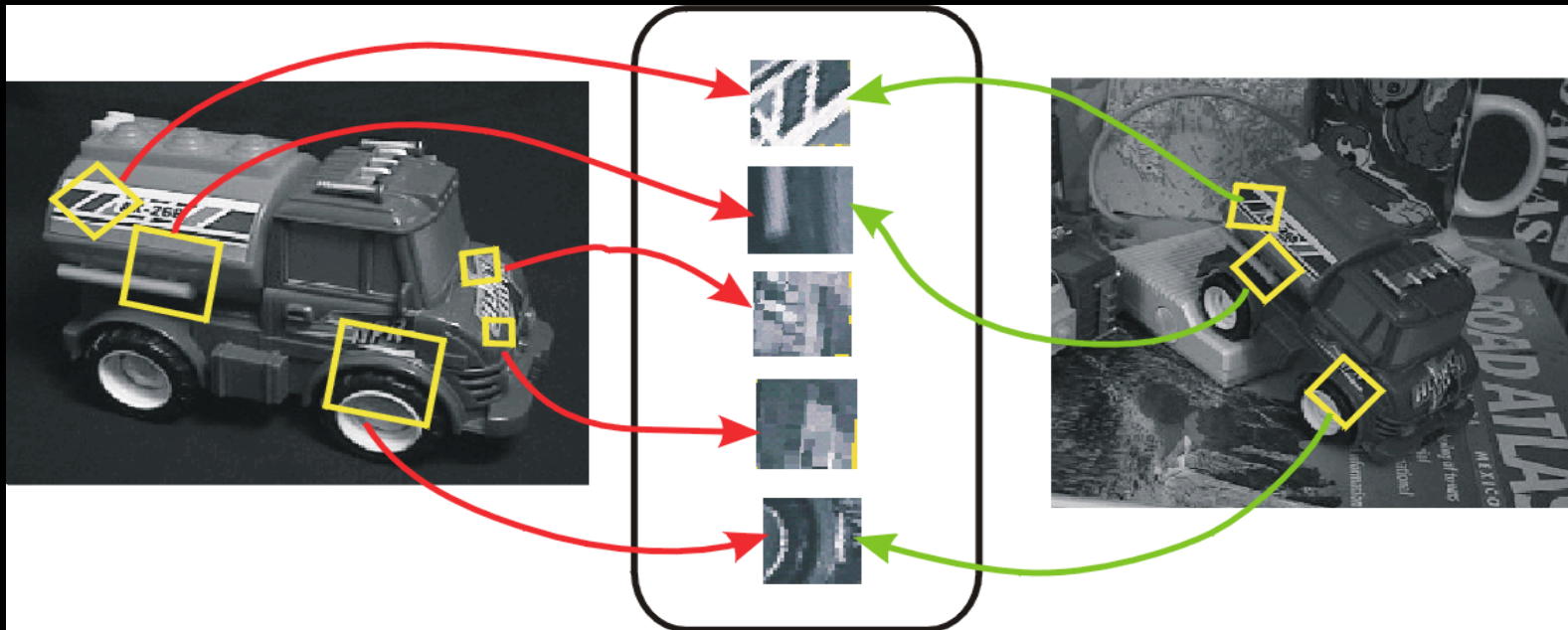


Normalization: From covariant regions to invariant features



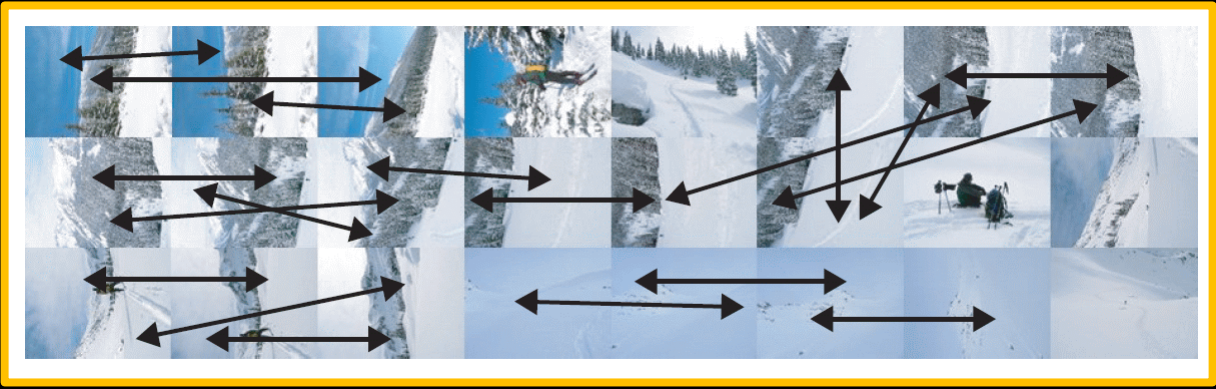
Invariance vs. covariance

- **Invariance:**
 - $\text{features}(\text{transform}(\text{image})) = \text{features}(\text{image})$
- **Covariance:**
 - $\text{features}(\text{transform}(\text{image})) = \text{transform}(\text{features}(\text{image}))$



Covariant detection => invariant description

Application: Panoramas from a Jumble of Pictures



- Extract SIFT features from all images
- Find Pairwise Homographies
- Find Connected Components over the pair connections
- Bundle adjust the connected component to find image to panorama transformation
- Render panorama with blending

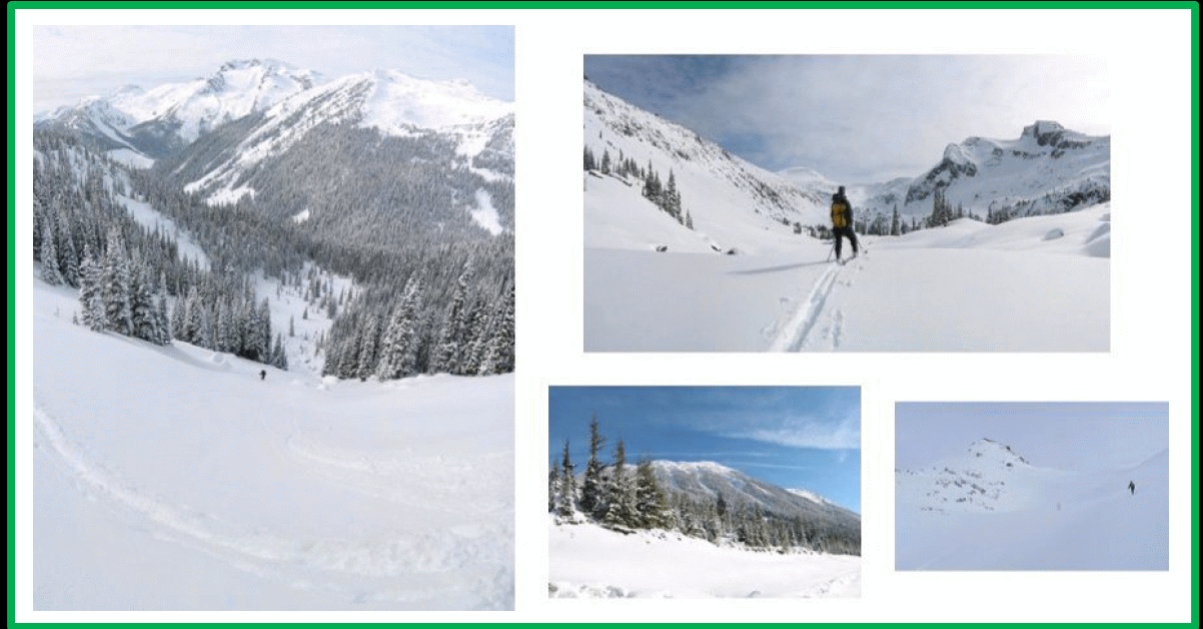


Photo Tourism / Photosynth



The slide features the University of Washington Computer Science & Engineering logo in the top left, the Microsoft logo in the top right, and the title "Photo Tourism" with the subtitle "Exploring photo collections in 3D". It contains three panels: (a) a grid of 25 photos of Notre-Dame de Paris; (b) a sparse 3D model of the cathedral's facade; and (c) a 3D viewer interface showing the cathedral and a photo gallery.

University of Washington
Computer Science & Engineering

Photo Tourism

Exploring photo collections in 3D

Microsoft

(a)

(b)

(c)

- Automatically computes each photo's viewpoint, and
- A sparse 3D model of the scene
- Explorer interface enables interactively moving in 3D space by seamlessly transitioning between photographs

Application: Scalable Images based Search

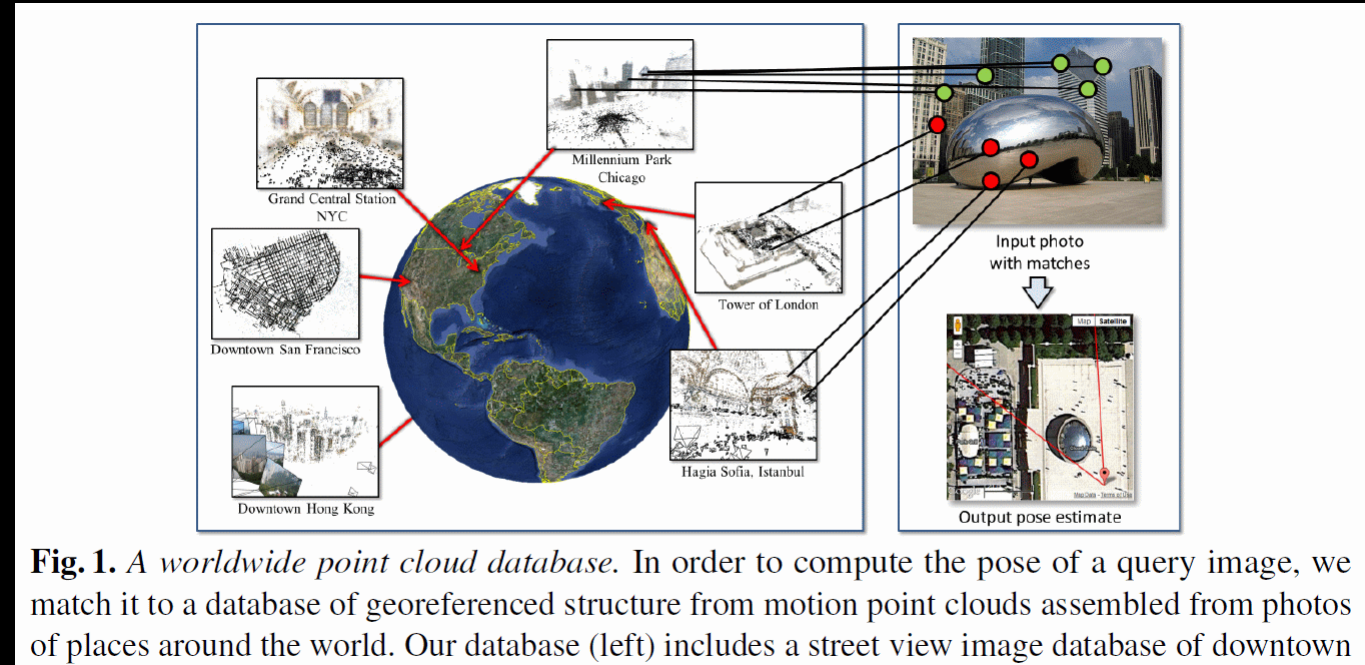


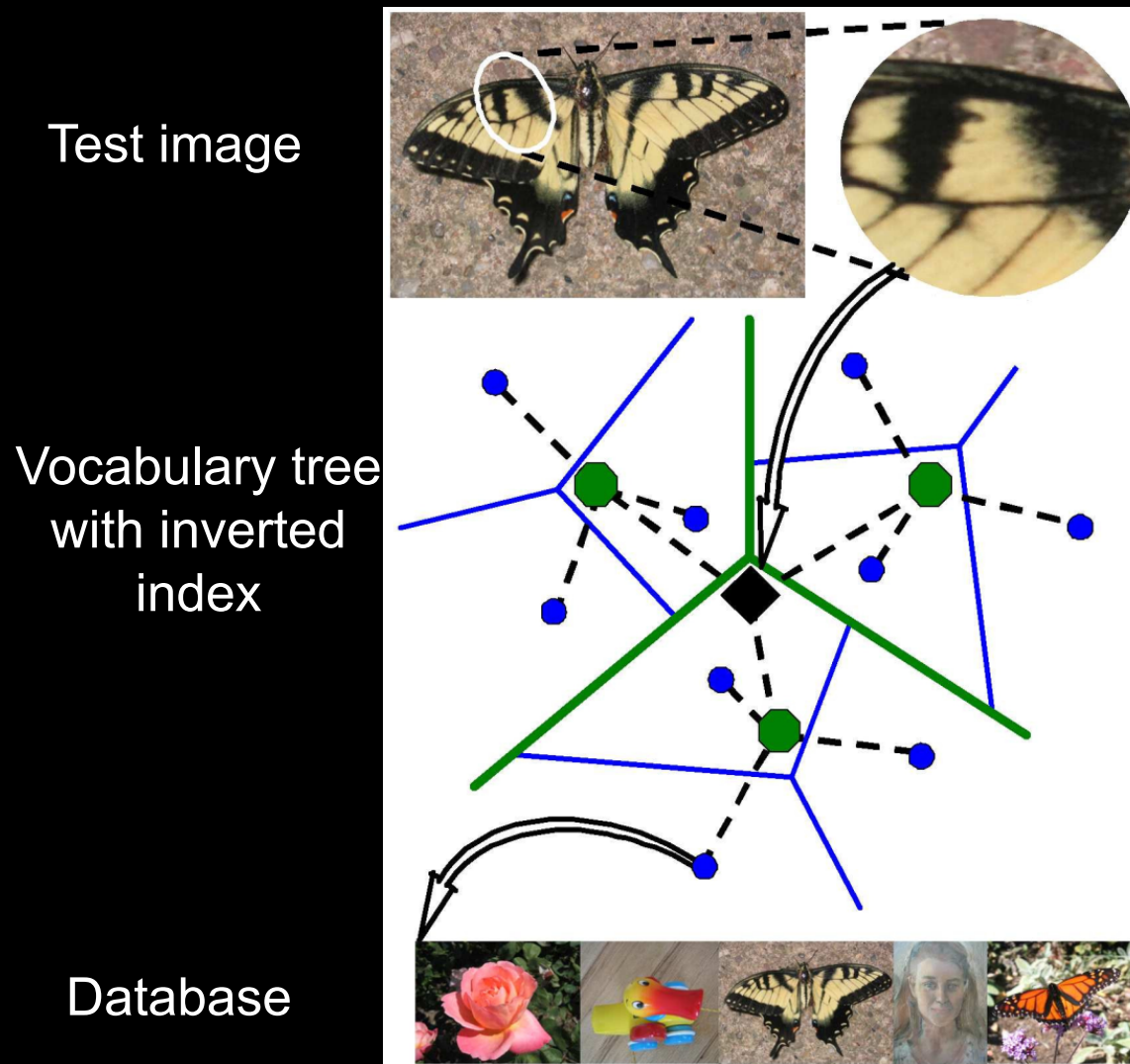
Fig. 1. A worldwide point cloud database. In order to compute the pose of a query image, we match it to a database of georeferenced structure from motion point clouds assembled from photos of places around the world. Our database (left) includes a street view image database of downtown

David G. Lowe. "Distinctive image features from scale-invariant keypoints." *IJCV* 60 (2), pp. 91-110, 2004.

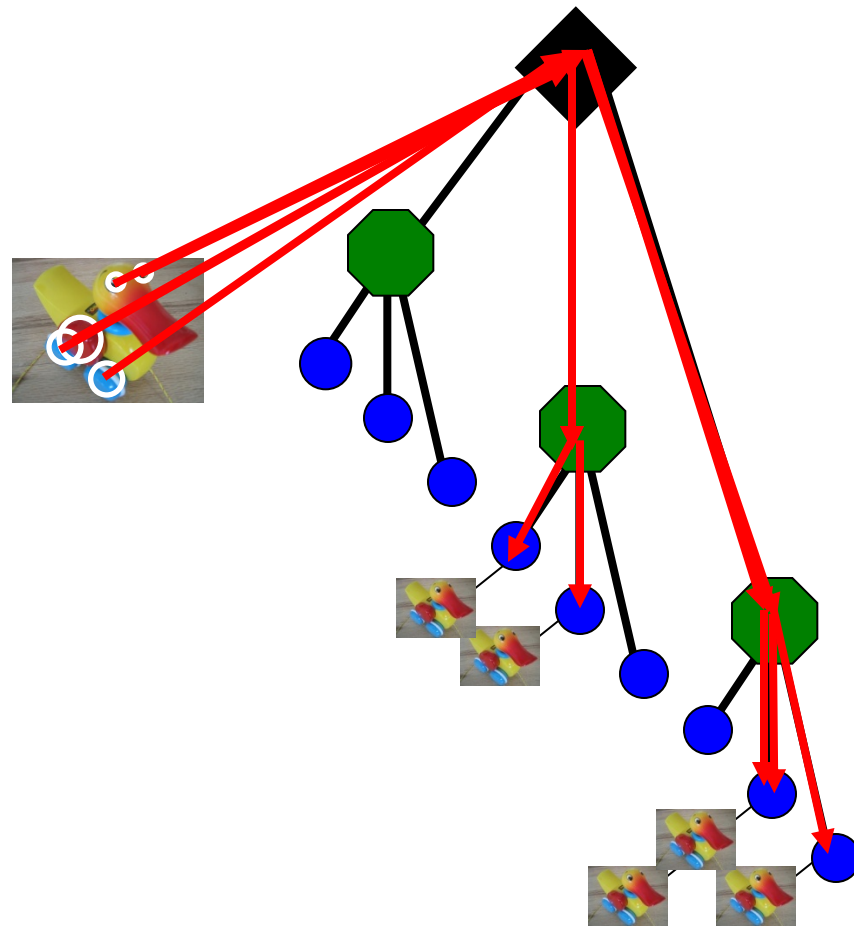
Worldwide Pose Estimation
Yunpeng Li, Noah Snavely, Dan Huttenlocher, Pascal Fua
ECCV 2012

- Find location of a Query image by matching against a large database of images indexed with their respective locations
- Find instances of objects / images in a database of images

Efficient indexing technique: Vocabulary trees

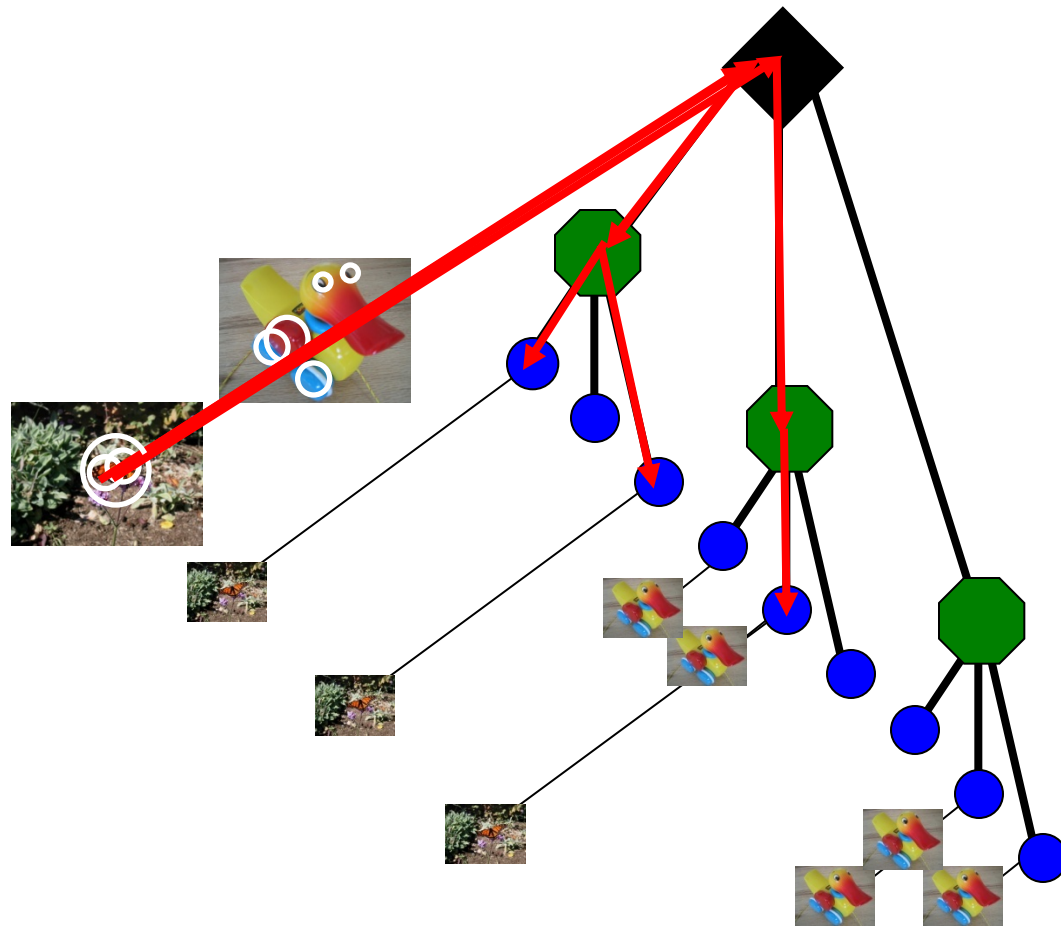


Model images



Populating the vocabulary tree/inverted index

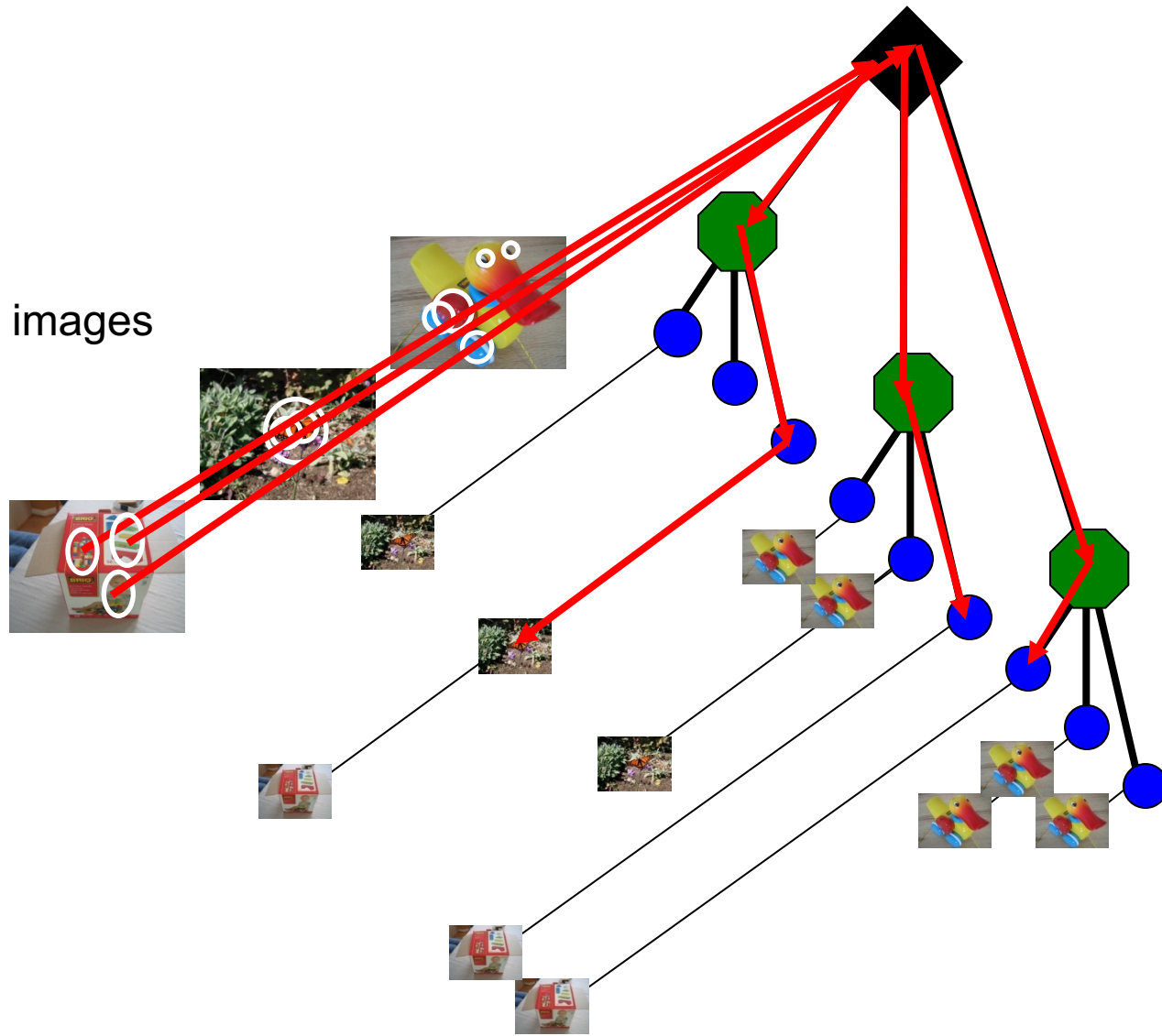
Model images



Populating the vocabulary tree/inverted index

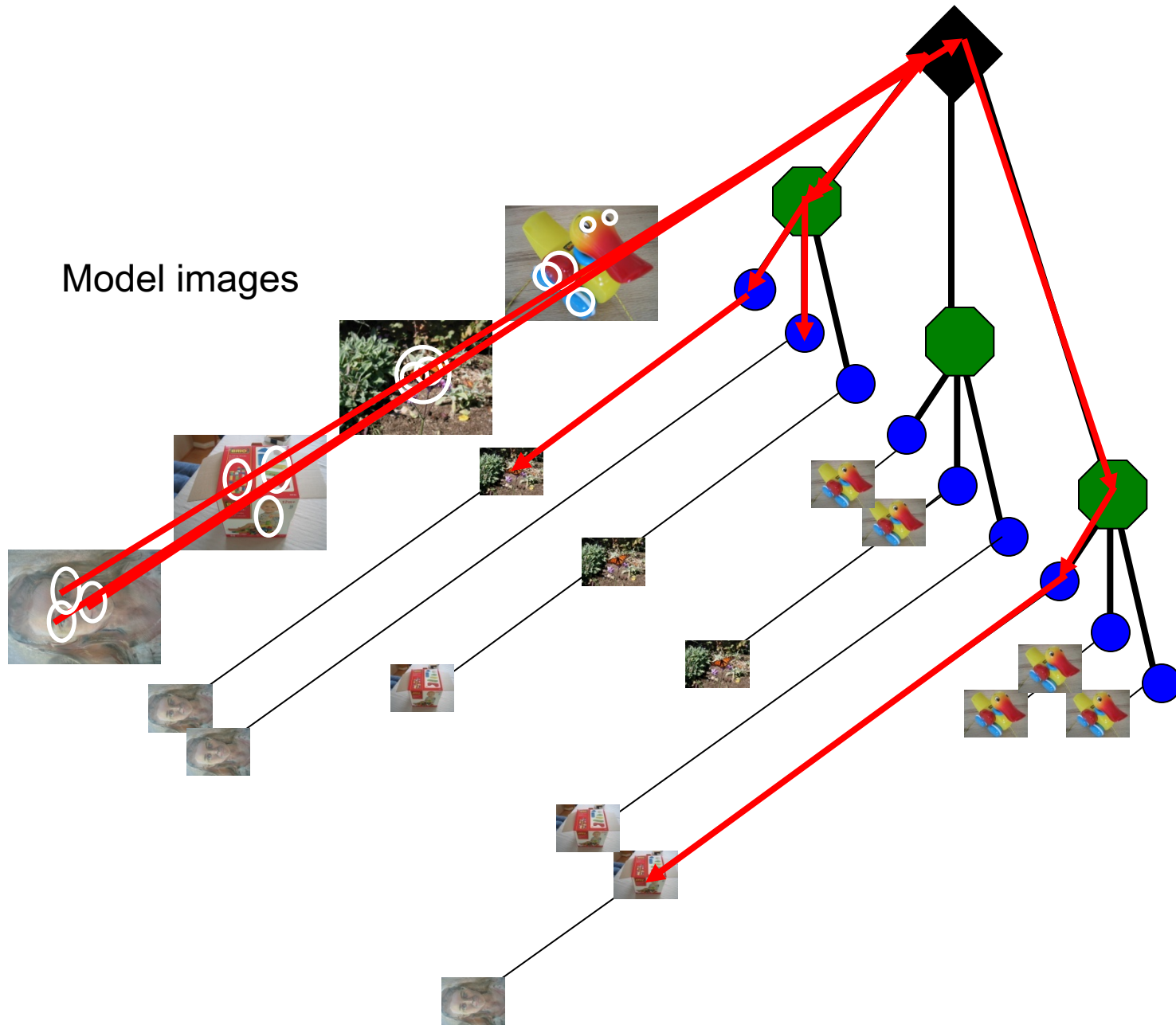
Slide credit: D. Nister

Model images



Populating the vocabulary tree/inverted index

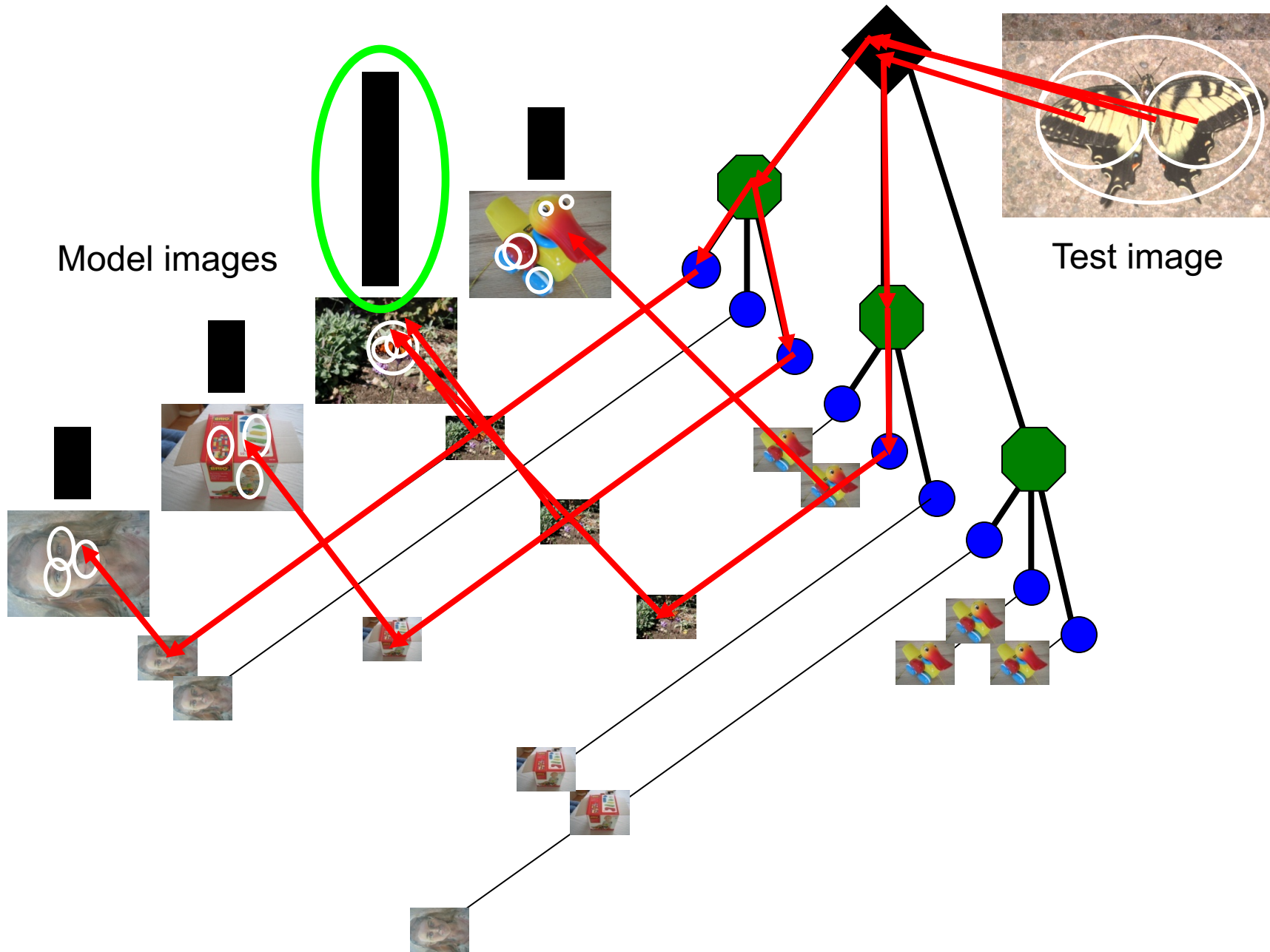
Slide credit: D. Nister



Model images

Populating the vocabulary tree/inverted index

Slide credit: D. Nister



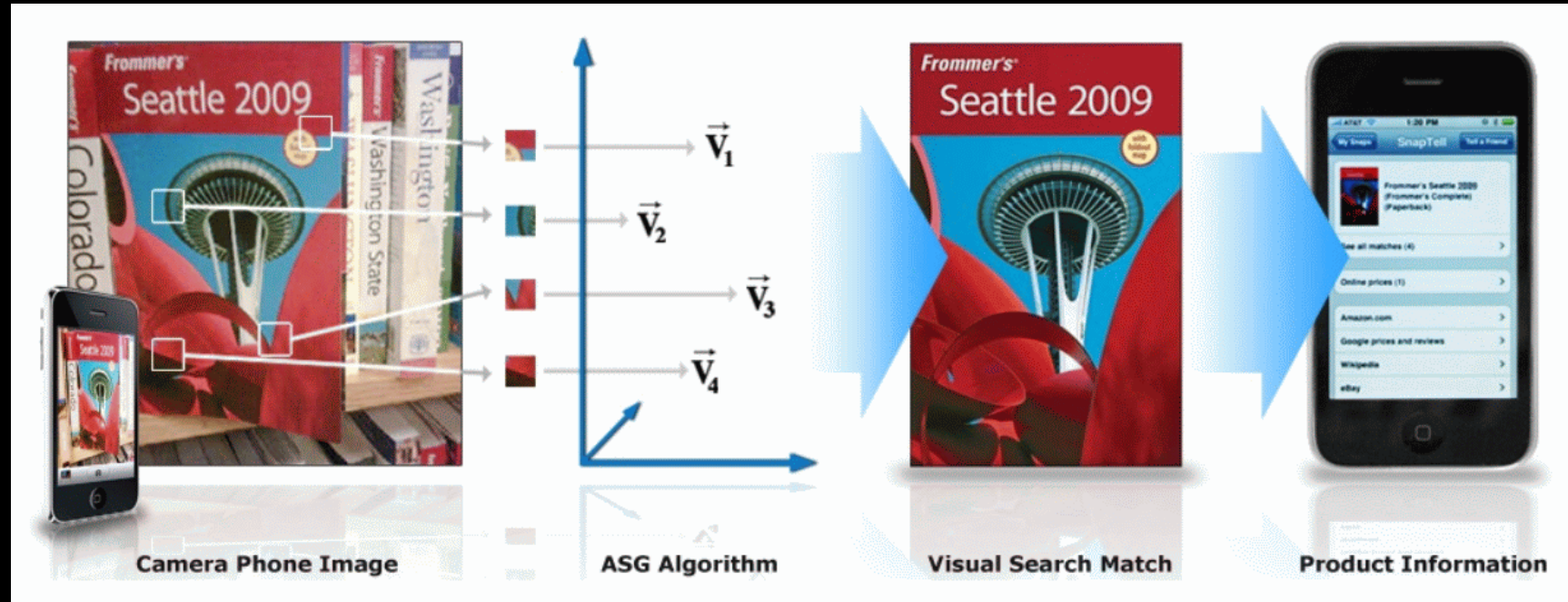
Model images

Test image

Looking up a test image

Slide credit: D. Nister

SnapTell : Visual Product Search (Acquired by A9 / Amazon)



- ASG algorithm (Accumulated Signed Gradient) is a SIFT like feature
- <https://computervisionblog.wordpress.com/tag/image-recognition/>
- Many Others too...

Google Cloud Anchors

- Google 'Cloud Anchors' will help synchronize group AR experiences across iOS and Android devices
- Employ Wide Baseline Matching
- <https://mediafocus.biz/google-cloud-anchors-will-help-synchronize-group-ar-experiences-across-ios-and-android-devices/>



Minecraft Earth via Azure Spatial Anchors



<https://www.theverge.com/2019/11/12/20961639/minecraft-earth-now-available-early-access-us-ios-android>

<https://youtu.be/AQEizp-VrVU>

Spatial Anchors are built and queried with Wide Baseline Matching and SfM / SLAM

How do Hand-Crafted Features Compare with Learned Features?

Comparative Evaluation of Hand-Crafted and Learned Local Features

Johannes L. Schönberger¹ Hans Hardmeier¹ Torsten Sattler¹ Marc Pollefeys^{1,2}

¹ Department of Computer Science, ETH Zürich ² Microsoft Corp.

{jsch,harhans,sattlert,pomarc}@inf.ethz.ch

CVPR 2017

- “Hand-crafted features still perform on par or better than recent learned features for image-based reconstruction.
- The current generation of learned descriptors shows a high variance across different datasets and applications.
- The next generation of learned descriptors needs more training data.”

Viewpoint and Illumination Variations Dataset

HPatches: A benchmark and evaluation of handcrafted and learned local descriptors

Vassileios Balntas*
Imperial College London
v.balntas@imperial.ac.uk

Karel Lenc*
University of Oxford
karel@robots.ox.ac.uk

Andrea Vedaldi
University of Oxford
vedaldi@robots.ox.ac.uk

Krystian Mikolajczyk
Imperial College London
k.mikolajczyk@imperial.ac.uk

<https://github.com/hpatches/hpatches-dataset>



Figure 1. Examples of image sequences; note the diversity of scenes and nuisance factors, including viewpoint, illumination, focus, reflections and other changes.

- **Reproducible, patch-based:** Descriptor evaluation should be done on patches to eliminate the detector related factors.
- **Diverse:** Representative of many different scenes and image capturing conditions.
- **Real:** Real data more challenging than a synthesized one due to nuisance factors that cannot be modelled in image transformations.
- **Large:** For accurate and stable evaluation; to provide substantial training sets for learning based descriptors.
- **Multitask:** Use cases, from matching image pairs to image retrieval.

How do Hand-Crafted Features Compare with Learned Features?

Image Matching across Wide Baselines: From Paper to Practice

Yuhe Jin¹ Dmytro Mishkin² Anastasiia Mishchuk³
Jiří Matas² Pascal Fua³ Kwang Moo Yi¹ Eduard Trulls⁴

¹University of Victoria ²Czech Technical University in Prague ³École Polytechnique Fédérale de Lausanne ⁴Google Research

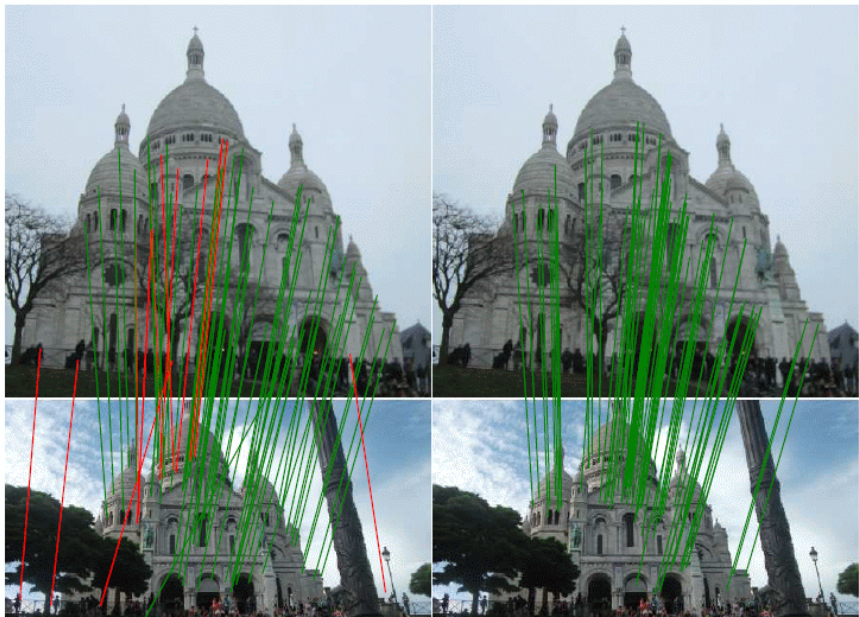


Figure 1. Every paper claims to outperform the state of the art. Is this possible, or an artifact of insufficient validation? On the left, we show stereo matches obtained with **D2-Net** (2019) [33], a state-of-the-art local feature, using OpenCV RANSAC with its default settings. On the right, we show **SIFT** (1999) [48] with a carefully tuned MAGSAC [29] – notice how the latter performs much better. We fill this gap with a new, modular benchmark for sparse image matching, with dozens of built-in methods.

Contributions

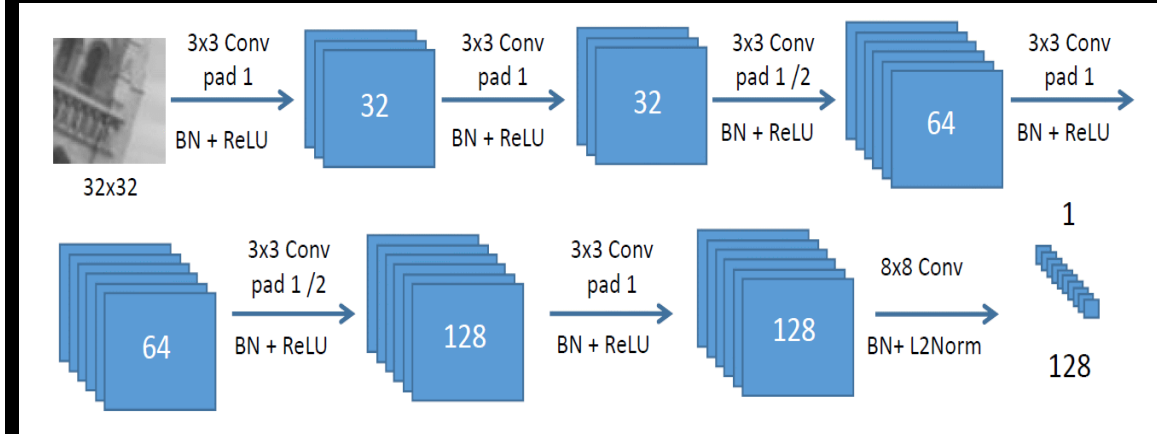
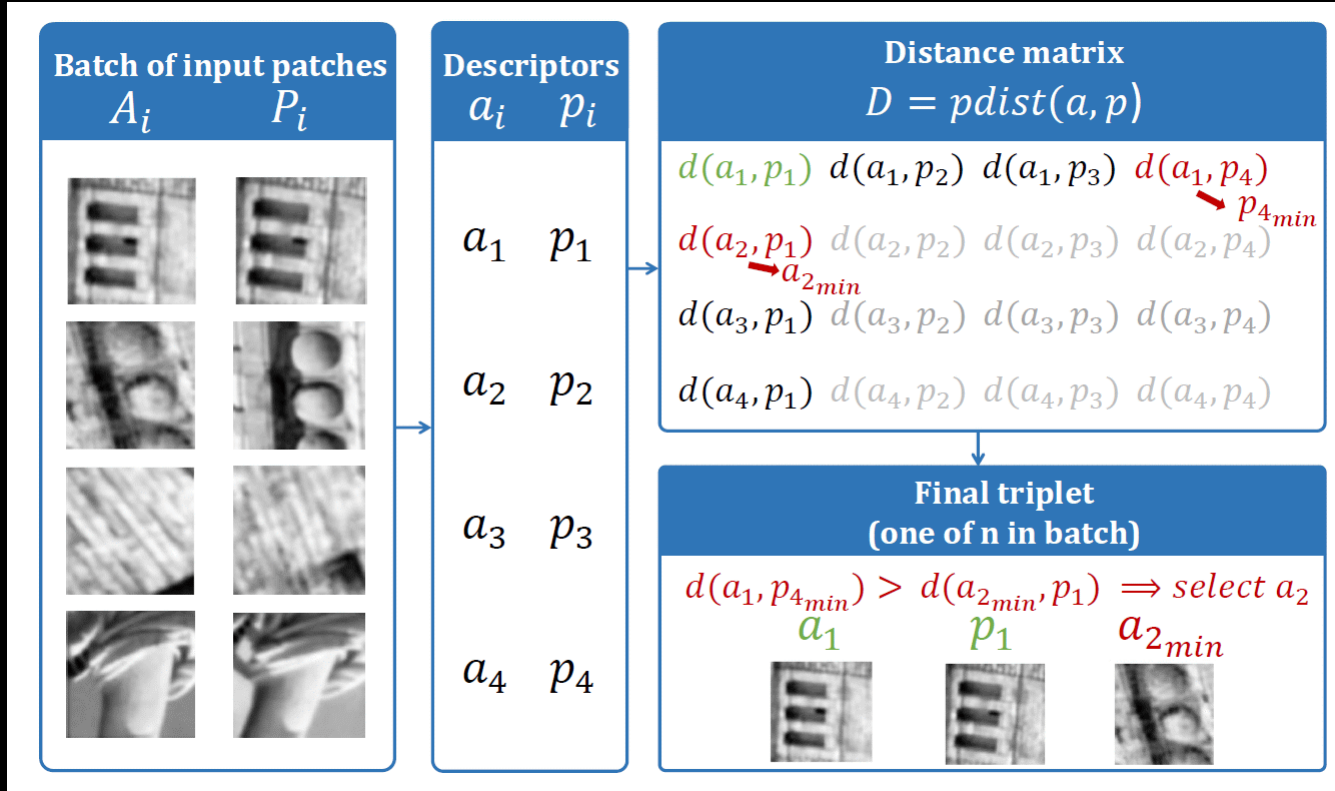
- Dataset with 30k images with depth maps and ground truth poses
- A modular pipeline incorporating dozens of methods for feature extraction and matching, and pose estimation
- Two downstream tasks – stereo and multi-view reconstruction – evaluated with downstream and intermediate metrics
- A thorough study of dozens of methods and techniques, hand-crafted and learned, and their combination, along with a procedure for hyper-parameter selection

Multiview Results

Method	NL [↑]	SR [↑]	RC [↑]	TL [↑]	mAP(5°) [↑]	mAP(10°) [↑]	ATE [↓]	Rank
CV-SIFT	2567.4	89.1	95.6	3.51	.3772	.4626	.7475	9
CV- $\sqrt{\text{SIFT}}$	2798.4	91.2	96.2	3.62	.4148	.5053	.6908	8
SURF	2421.7	84.5	93.9	3.11	.2591	.3353	.8259	13
AKAZE	3258.8	89.4	95.9	3.45	.3388	.4256	.7604	10
ORB	2341.1	84.5	92.1	3.09	.1905	.2529	.9165	16
DoG-HardNet	2834.1	91.5	96.3	3.79	.4644	.5599	.6669	1
L2Net	2413.7	88.0	95.2	3.70	.4383	.5308	.6656	6
Key.Net-HardNet	3755.3	96.6	98.0	3.89	.4438	.5456	.6688	4
Geodesc	2631.8	89.4	95.9	3.72	.4325	.5258	.6729	7
ContextDesc	2223.8	90.2	96.4	3.66	.4393	.5354	.6697	5
SOSNet	2681.6	90.1	96.4	3.81	.4650	.5592	.6583	2
LogPolarDesc	3029.7	90.1	95.6	3.79	.4622	.5565	.6657	3
SuperPoint (2k)	762.5	83.0	92.7	3.76	.2959	.3814	.7767	11
LF-Net (2k)	1014.6	80.0	89.8	3.63	.2936	.3723	.7517	12
D2-Net (SS)	3302.8	90.0	95.8	3.17	.2056	.2933	.8595	15
D2-Net (MS)	4022.9	93.7	97.0	3.05	.2143	.3149	.8335	14

- (NL) Number of 3D Landmarks
- (SR) Success Rate (%) in the 3D reconstruction across ‘bags’
- (RC) Ratio of Cameras (%) registered in a ‘bag’
- (TL) Track Length or number of observations per landmark;
- mAP at 5-10degs
- (ATE) Absolute Trajectory Error.

HardNet: Uses SIFT like Best-to-Next Distance for Training Descriptors



$$L = \frac{1}{n} \sum_{i=1, n} \max(0, 1 + d(a_i, p_i) - \min(d(a_i, p_{j_{min}}), d(a_{k_{min}}, p_i)))$$

- Choose the hardest negative, i.e. minimum distance to a negative.
- Make gap between the matching and hard non-matching to be a maximum

“Hot Topic”



Image Matching: Local Features & Beyond
CVPR 2020 Workshop

A Contemporary Example of Learned Features

SuperPoint: Self-Supervised Interest Point Detection and Description

CVPR 2018

- Self-supervised framework for training interest point detectors and descriptors
- Fully-convolutional model operates on full-sized images and jointly computes pixel-level interest point locations and associated descriptors in one forward pass.

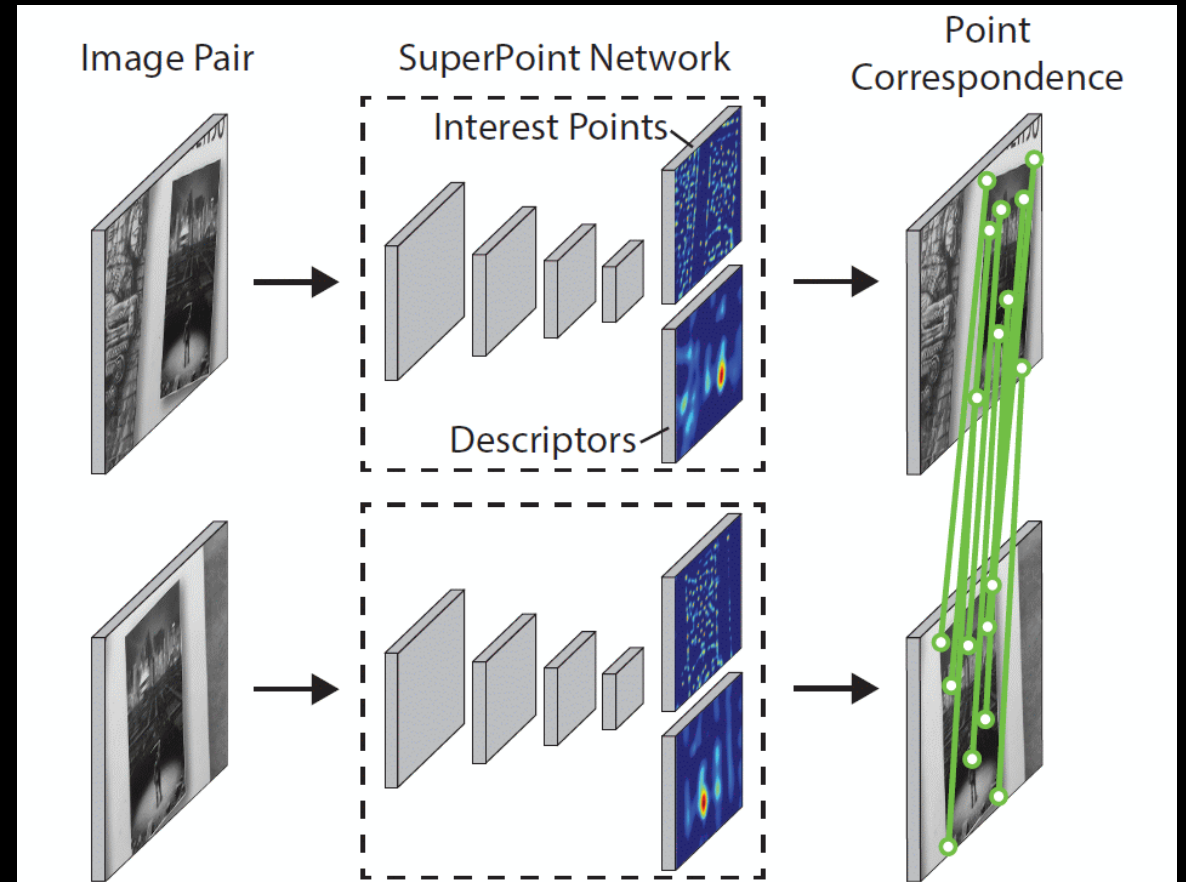


Figure 1. **SuperPoint for Geometric Correspondences.** We present a fully-convolutional neural network that computes SIFT-like 2D interest point locations and descriptors in a single forward pass and runs at 70 FPS on 480×640 images with a Titan X GPU.

Self-Supervised Training

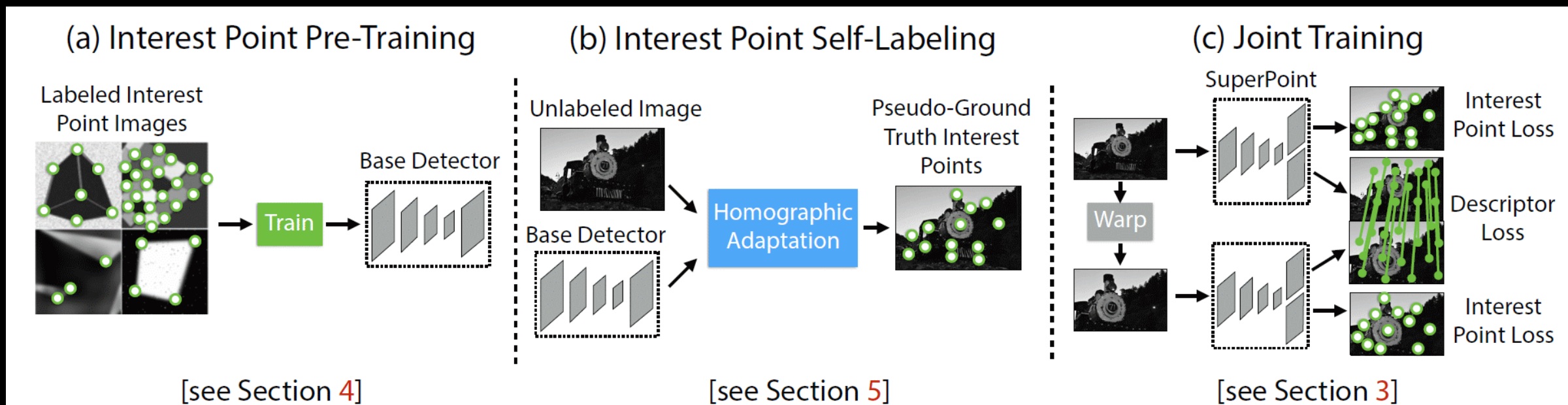


Figure 2. **Self-Supervised Training Overview.** In our self-supervised approach, we (a) pre-train an initial interest point detector on synthetic data and (b) apply a novel Homographic Adaptation procedure to automatically label images from a target, unlabeled domain. The generated labels are used to (c) train a fully-convolutional network that jointly extracts interest points and descriptors from an image.

Superpoint Architecture

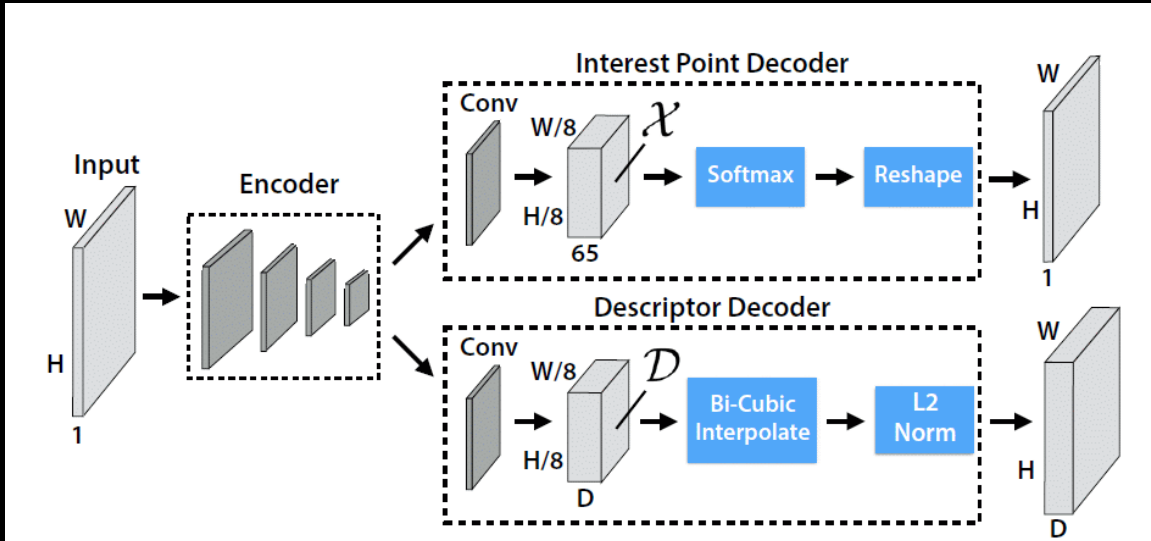


Figure 3. **SuperPoint Decoders.** Both decoders operate on a shared and spatially reduced representation of the input. To keep the model fast and easy to train, both decoders use non-learned upsampling to bring the representation back to $\mathbb{R}^{H \times W}$.

- The interest point detector head computes $H_c \times W_c \times 65$ and outputs a tensor sized $H \times W$.
- The 65 channels correspond to local, non-overlapping 8×8 grid regions of pixels plus an extra “no interest point” dustbin.
- After a channel-wise softmax, the dustbin dimension is removed and $H_c \times W_c \times 64 \rightarrow H \times W$ reshape is performed.
- The descriptor head computes $H_c \times W_c \times D$ and outputs a tensor sized $H \times W \times D$.
- To output a dense map of L2-normalized fixed length descriptors, first output a semi-dense grid of descriptors (e.g., one every 8 pixels).
- The decoder then performs bicubic interpolation of the descriptor and then L2-normalizes to unit length.

Joint Geometric and Classification Loss

$$\mathcal{L}(\mathcal{X}, \mathcal{X}', \mathcal{D}, \mathcal{D}'; Y, Y', S) = \mathcal{L}_p(\mathcal{X}, Y) + \mathcal{L}_p(\mathcal{X}', Y') + \lambda \mathcal{L}_d(\mathcal{D}, \mathcal{D}', S).$$

- The interest point detector loss function \mathcal{L}_p is a fully convolutional cross-entropy
- The descriptor loss is applied to all pairs of descriptor cells, (h, w) and (h', w')
- The homography-induced correspondence between the (h, w) cell and the (h', w') cell can be written as:
- The descriptor loss is given by:

$$\mathcal{L}_p(\mathcal{X}, Y) = \frac{1}{H_c W_c} \sum_{\substack{h=1 \\ w=1}}^{H_c, W_c} l_p(\mathbf{x}_{hw}; y_{hw}),$$

where

$$l_p(\mathbf{x}_{hw}; y) = -\log \left(\frac{\exp(\mathbf{x}_{hw} y)}{\sum_{k=1}^{65} \exp(\mathbf{x}_{hw} k)} \right).$$

$$s_{hwh'w'} = \begin{cases} 1, & \text{if } \|\widehat{\mathcal{H}}\mathbf{p}_{hw} - \mathbf{p}_{h'w'}\| \leq 8 \\ 0, & \text{otherwise} \end{cases}$$

$$\mathcal{L}_d(\mathcal{D}, \mathcal{D}', S) =$$

$$\frac{1}{(H_c W_c)^2} \sum_{\substack{h=1 \\ w=1}}^{H_c, W_c} \sum_{\substack{h'=1 \\ w'=1}}^{H_c, W_c} l_d(\mathbf{d}_{hw}, \mathbf{d}'_{h'w'}; s_{hwh'w'}),$$

where

$$l_d(\mathbf{d}, \mathbf{d}'; s) = \lambda_d * s * \max(0, m_p - \mathbf{d}^T \mathbf{d}') + (1 - s) * \max(0, \mathbf{d}^T \mathbf{d}' - m_n).$$

Comparative Results

	57 Illumination Scenes		59 Viewpoint Scenes	
	NMS=4	NMS=8	NMS=4	NMS=8
<i>SuperPoint</i>	.652	.631	.503	.484
<i>MagicPoint</i>	.575	.507	.322	.260
<i>FAST</i>	.575	.472	.503	.404
<i>Harris</i>	.620	.533	.556	.461
<i>Shi</i>	.606	.511	.552	.453
<i>Random</i>	.101	.103	.100	.104

Table 3. **HPatches Detector Repeatability.** SuperPoint is the most repeatable under illumination changes, competitive on viewpoint changes, and outperforms MagicPoint in all scenarios.

	Homography Estimation			Detector Metrics		Descriptor Metrics	
	$\epsilon = 1$	$\epsilon = 3$	$\epsilon = 5$	Rep.	MLE	NN mAP	M. Score
<i>SuperPoint</i>	.310	.684	.829	.581	1.158	.821	.470
<i>LIFT</i>	.284	.598	.717	.449	1.102	.664	.315
<i>SIFT</i>	.424	.676	.759	.495	0.833	.694	.313
<i>ORB</i>	.150	.395	.538	.641	1.157	.735	.266

Table 4. **HPatches Homography Estimation.** SuperPoint outperforms LIFT and ORB and performs comparably to SIFT using various ϵ thresholds of correctness. We also report related metrics which measure detector and descriptor performance individually.

- SIFT performs well for sub-pixel precision homographies and has the lowest mean localization error (MLE).
- SuperPoint scores strongly in descriptor-focused metrics such as nearest neighbor mAP and matching score (M. Score)