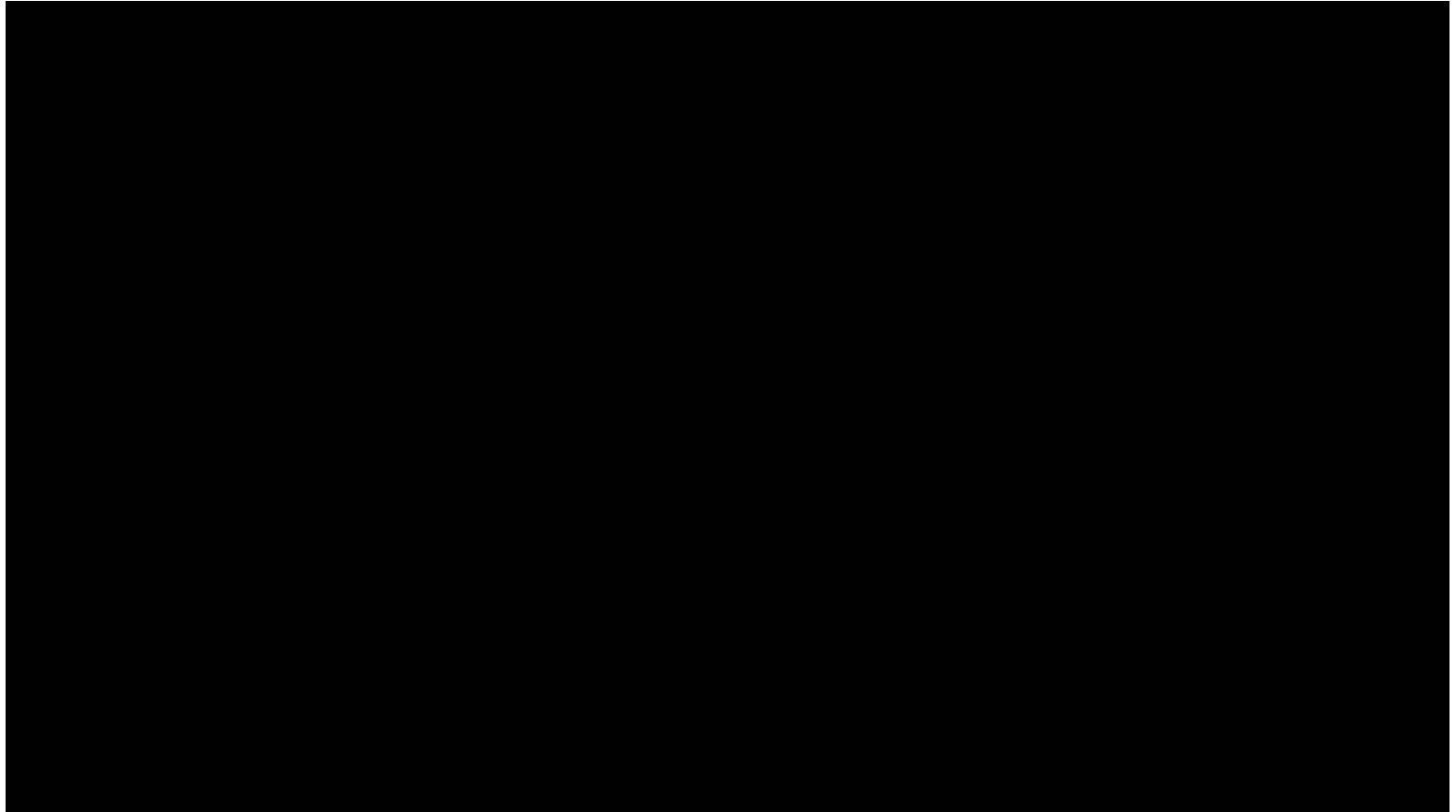
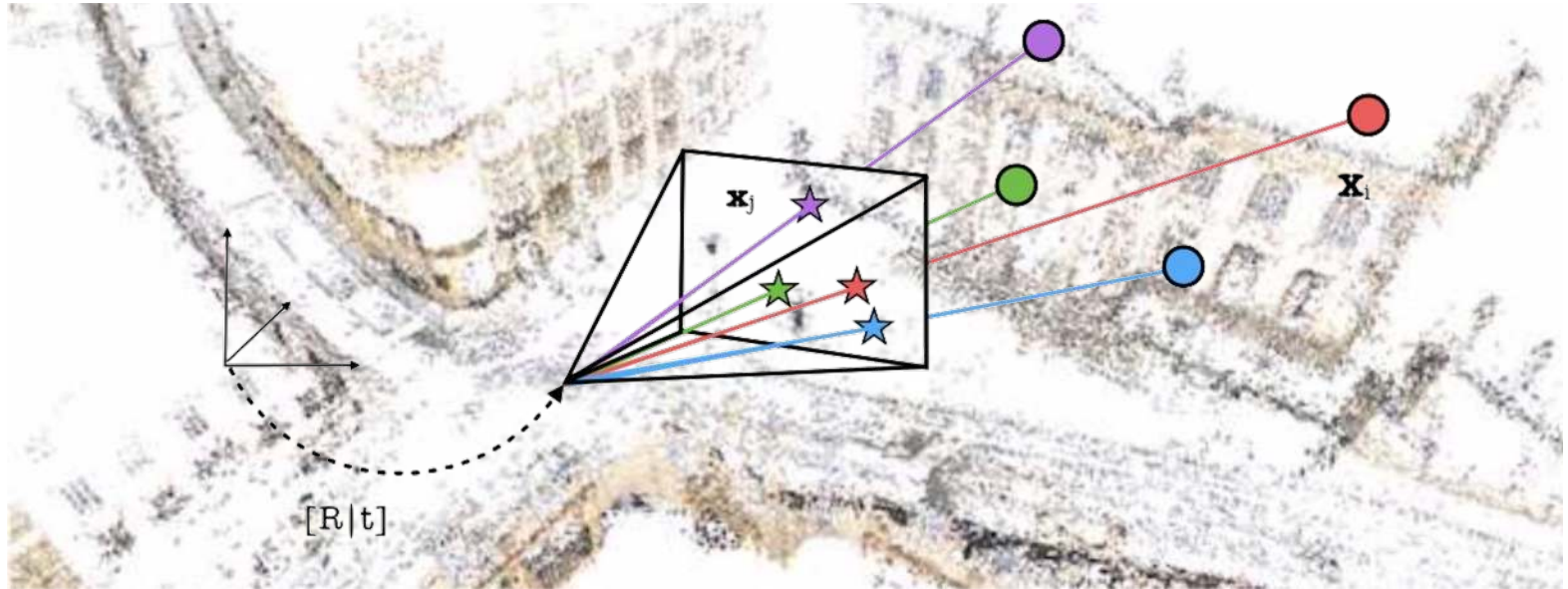


Deep Learning for Augmented Reality

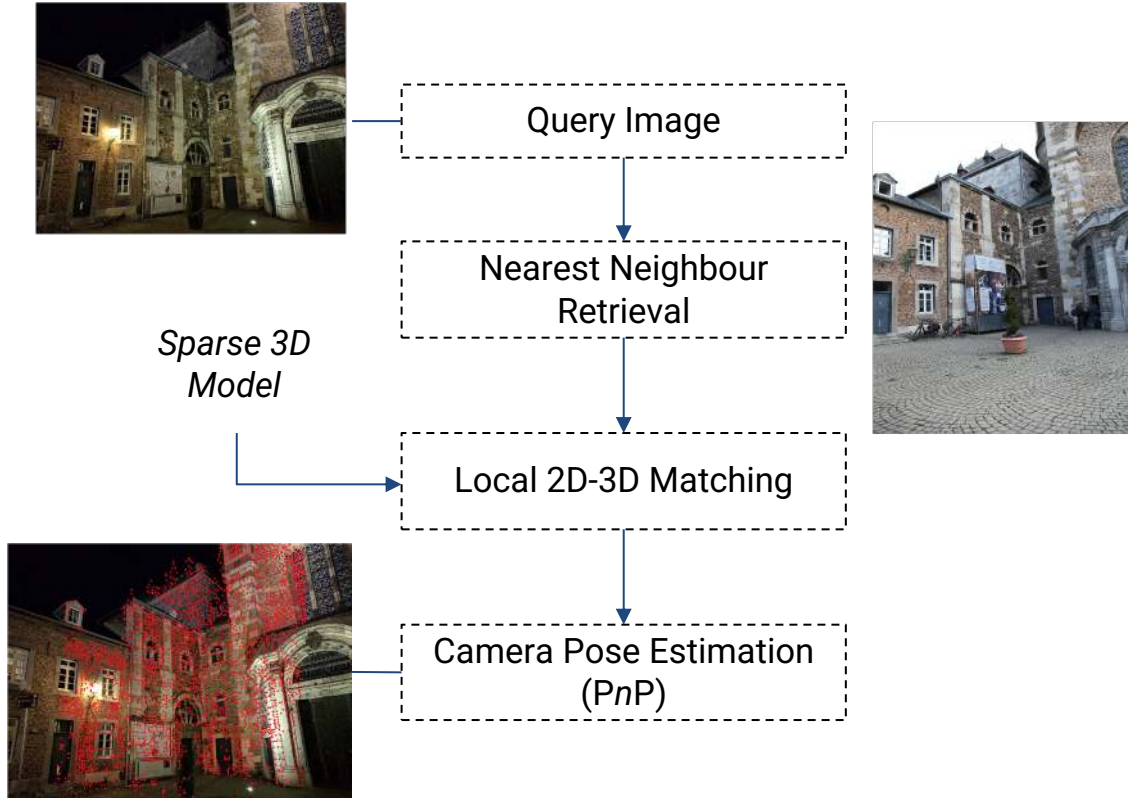
Vincent Lepetit

Image Retrieval





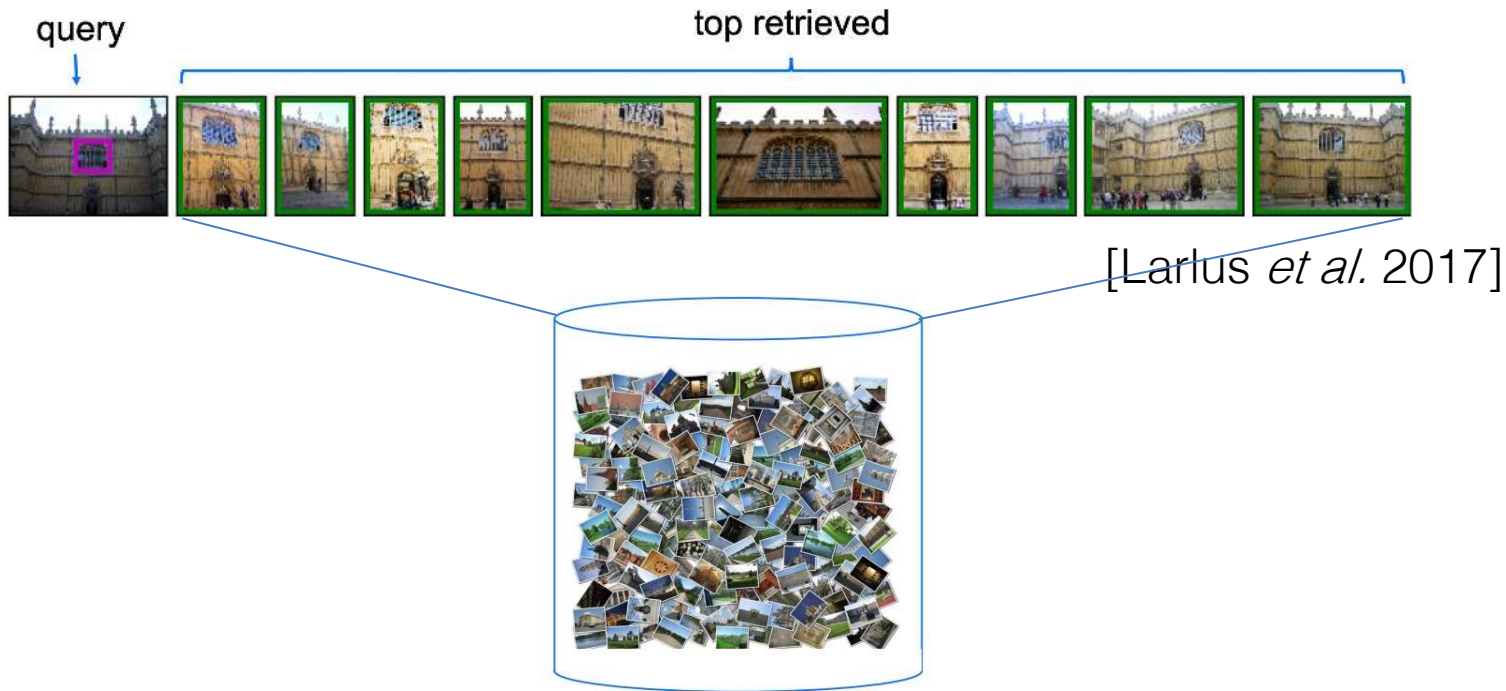
Hierarchical Localization



NetVLAD [Arandjelovic *et al.*, CVPR 2016]:
Application to Camera Localization

Localization by Image Retrieval

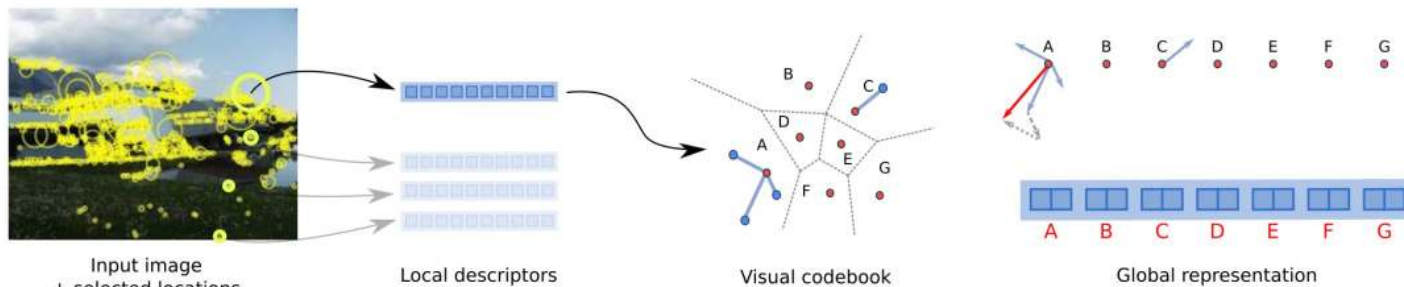
Matching a query image against a database



Aggregating Local Representations

A popular approach: VLAD (Vector of Locally Aggregated Descriptors)

- Assign local descriptors to *visual words*;
- Concatenate vectors for individual words by computing *residuals*;
- Store a 2D vector per cluster as part of final descriptor.

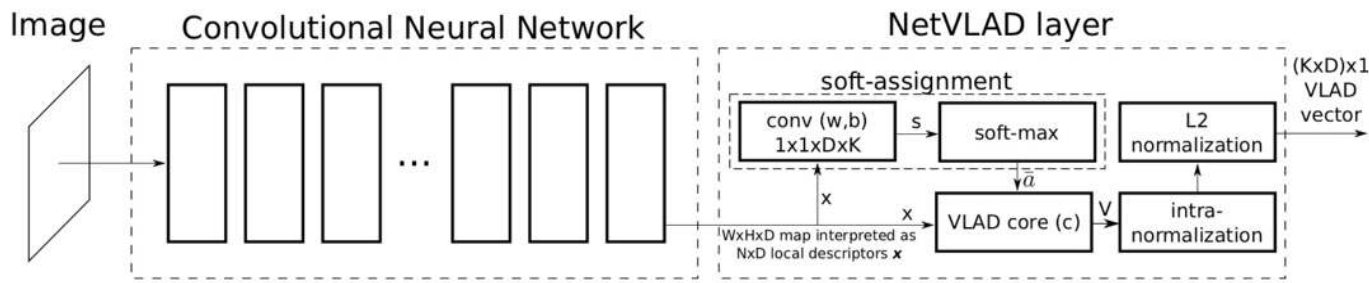


[Jégou et al. CVPR 2010]

Improving VLAD: *Learning* to extract local features and to aggregate them

NetVLAD: Apply VLAD on features learned end-to-end

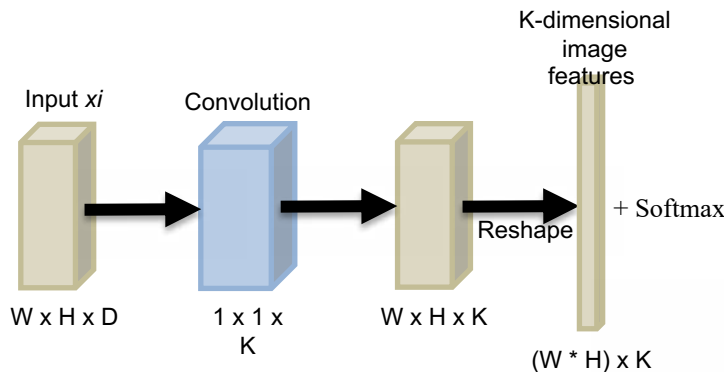
Define a differentiable VLAD layer, append it to a Siamese Network



[Arandjelovic et al. CVPR16]

NetVLAD in practice

$\{\mathbf{w}_k\}$, $\{b_k\}$ and $\{c_k\}$ are
sets of trainable parameters



"Soft" assignment of features to 1 of K cluster centers

$$\bar{a}_k(\mathbf{x}_i) = \frac{e^{\mathbf{w}_k^T \mathbf{x}_i + b_k}}{\sum_{k'} e^{\mathbf{w}_{k'}^T \mathbf{x}_i + b_{k'}}$$

Compute the residuals between the features and the centroid of their clusters (as in VLAD)

$(K \times D) \times 1$
descriptors

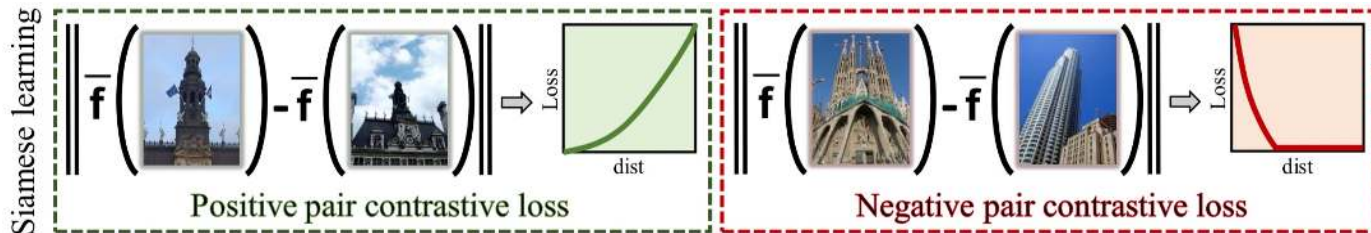
Normalization

$$V(j, k) = \sum_{i=1}^N \frac{e^{\mathbf{w}_k^T \mathbf{x}_i + b_k}}{\sum_{k'} e^{\mathbf{w}_{k'}^T \mathbf{x}_i + b_{k'}} (x_i(j) - c_k(j))$$

Training NetVLAD

The network is trained on pairs of images, either *positive* or *negative*;

- For positive pairs, minimize the distance between the output descriptors.
- For negative pairs, maximize it.



$$\mathcal{L}(i, j) = \begin{cases} \frac{1}{2} \|\bar{\mathbf{f}}(i) - \bar{\mathbf{f}}(j)\|^2, & \text{if } Y(i, j) = 1 \text{ [Radenovic et al. TPAMI2018]} \\ \frac{1}{2} (\max\{0, \tau - \|\bar{\mathbf{f}}(i) - \bar{\mathbf{f}}(j)\|\})^2, & \text{if } Y(i, j) = 0 \end{cases}$$

NetVLAD: Results

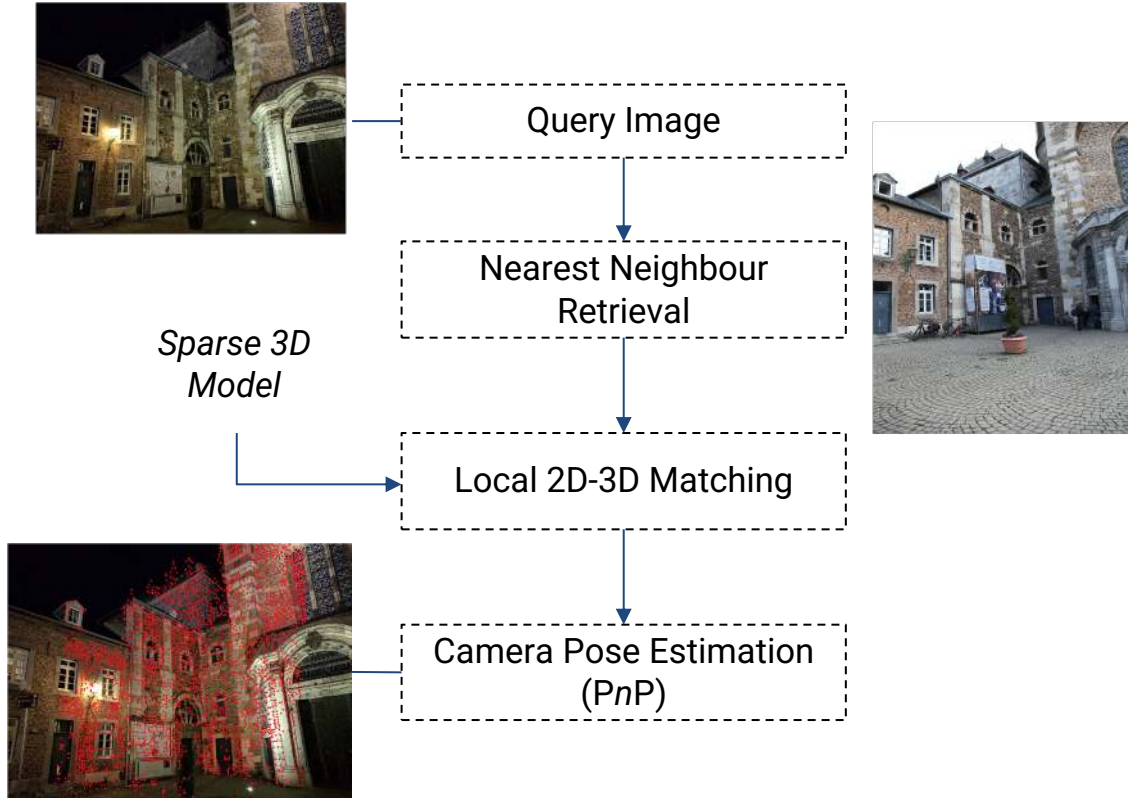
Query



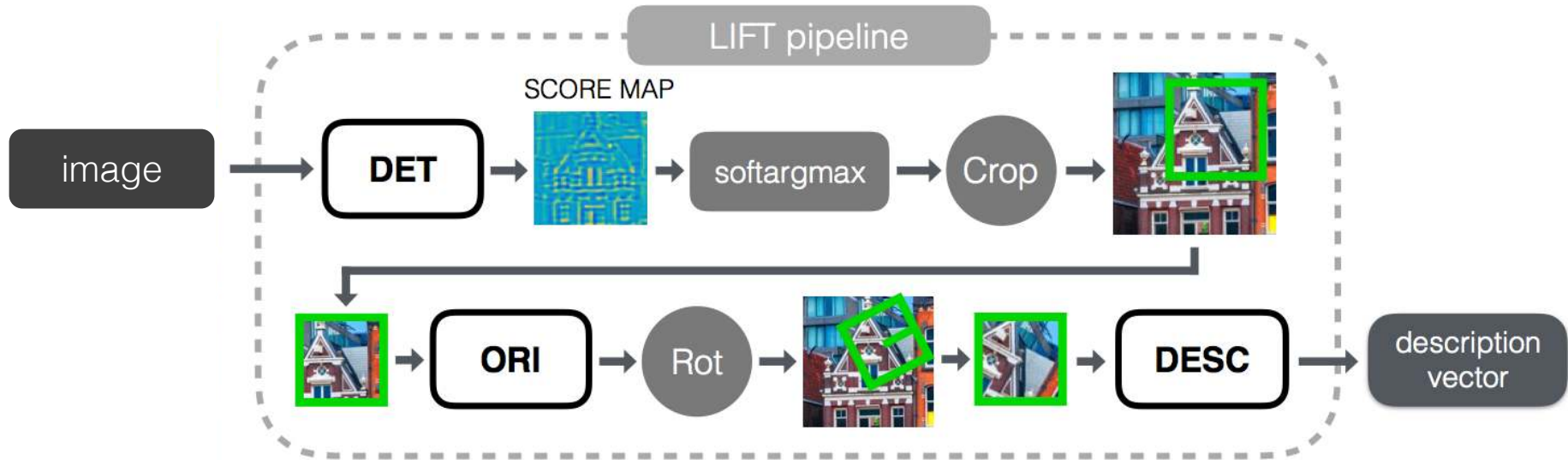
Ours

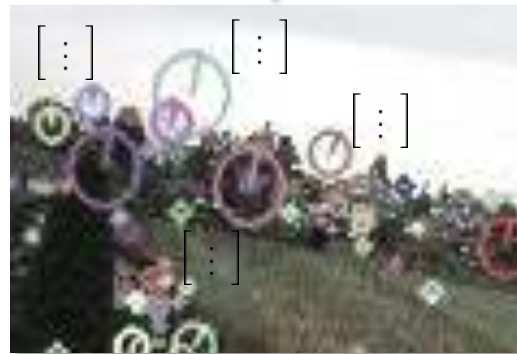
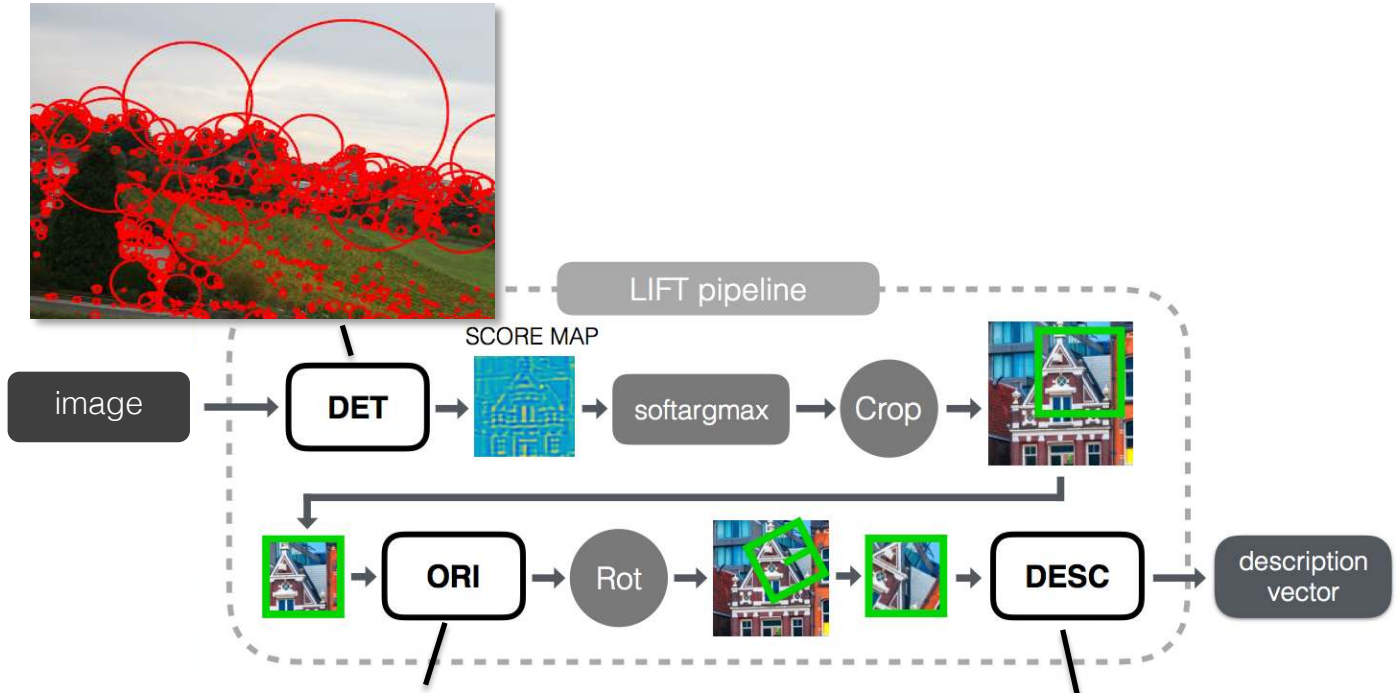


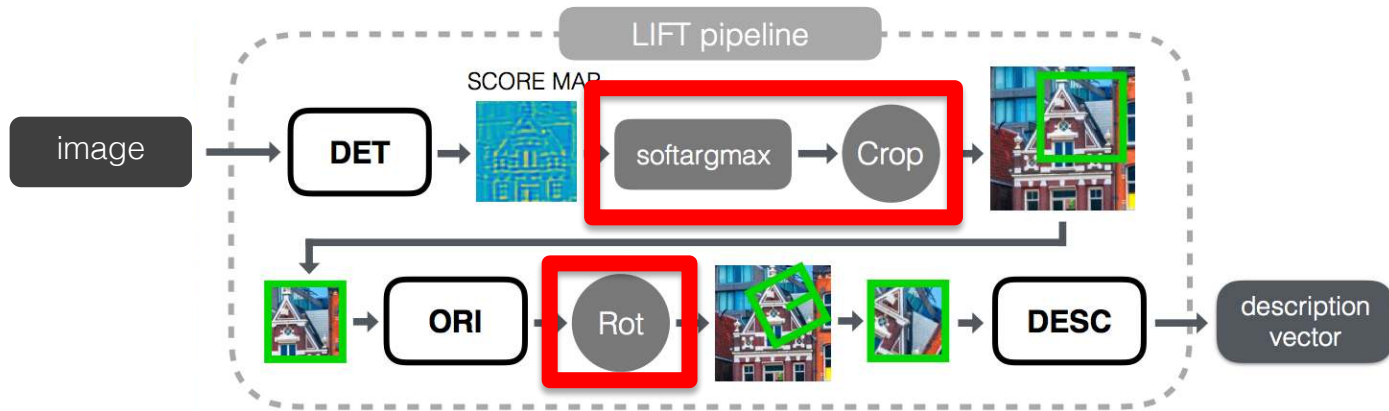
Hierarchical Localization



LIFT







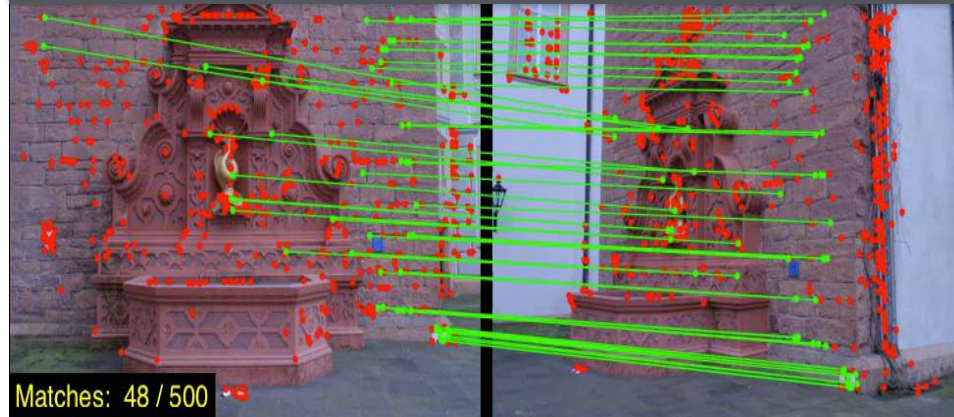
“Differentiable glue”

SIFT. Average: 60.2 matches

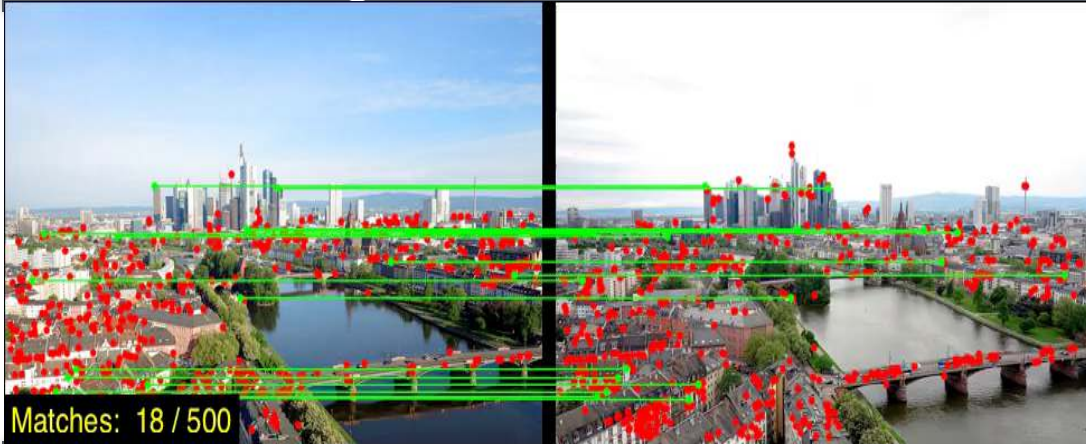


+64%

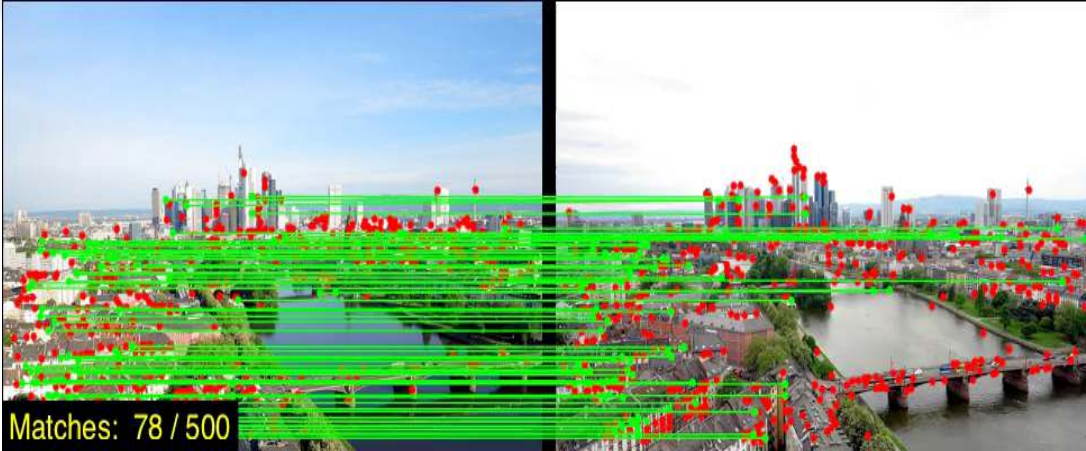
LIFT (Ours). Average: 98.6 matches



SIFT. Average: 23.1 matches

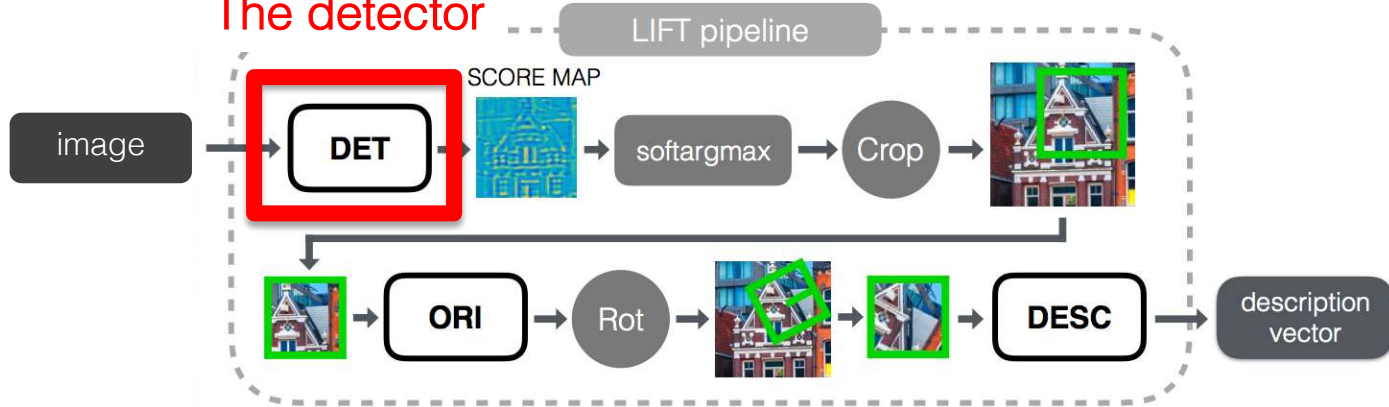


LIFT (Ours). Average: 60.6 matches



+162%

The detector





Detection under Severe Illumination Changes





Detection under Severe Illumination Changes



It's not going to work well even using very good descriptors !



Learning to Detect under Severe Illumination Changes

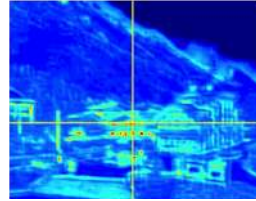




How the Detector is Used at Run-Time



input image



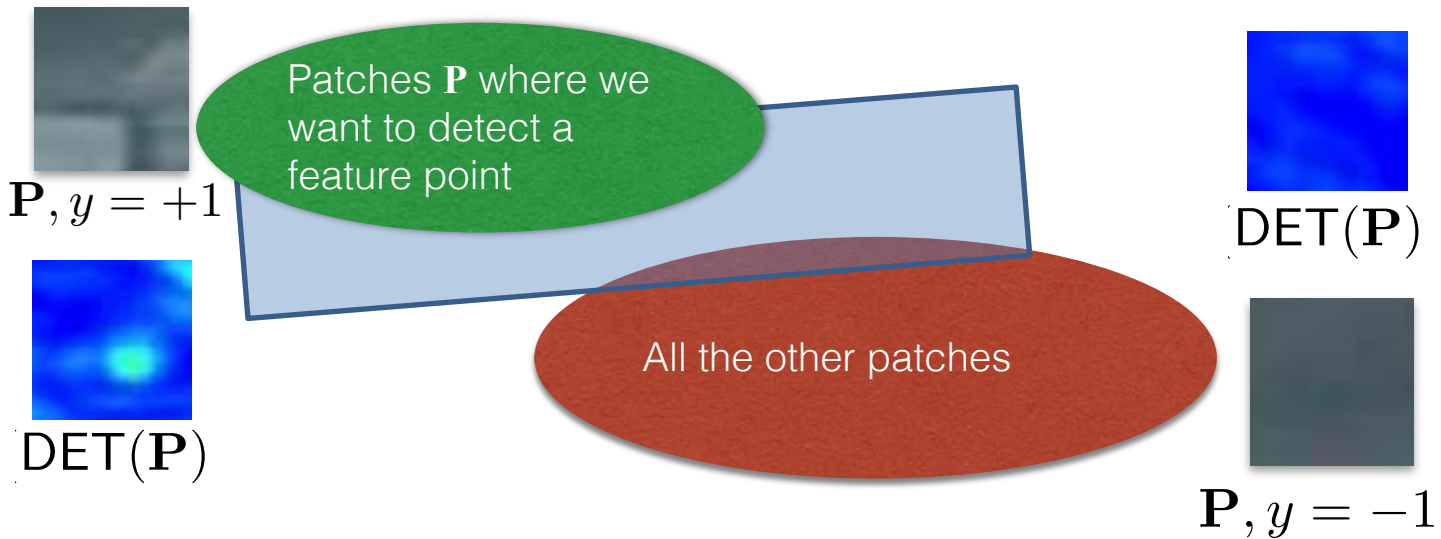
'score map'



feature points



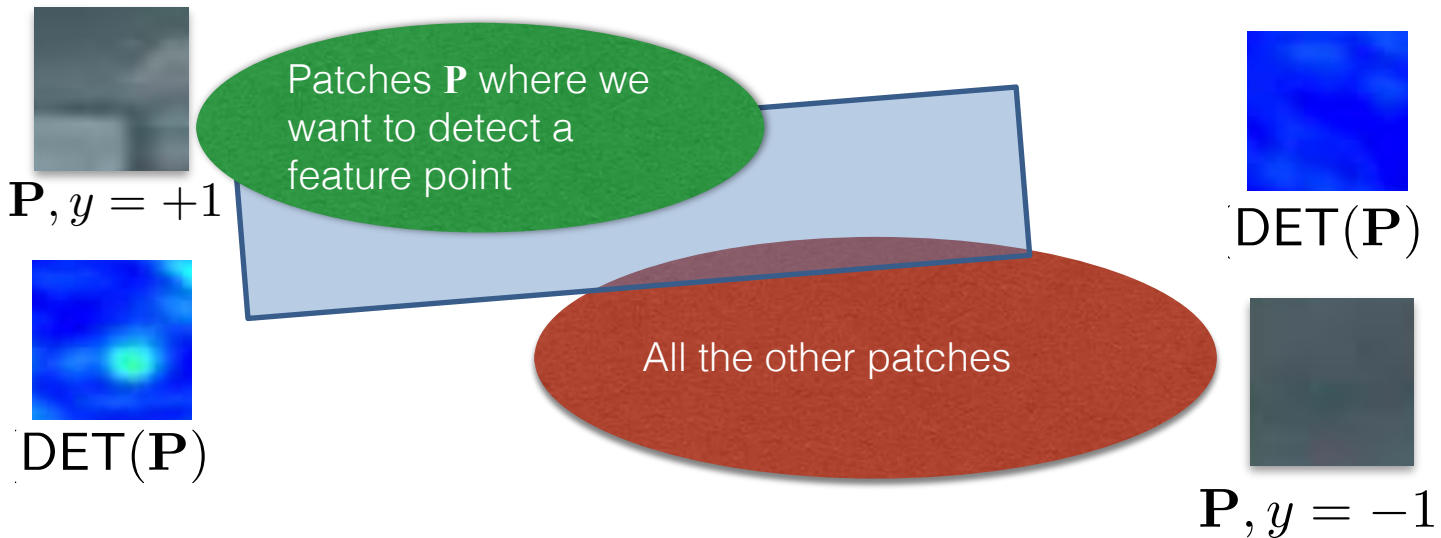
Cost Function (1)



$$\mathcal{L}_{\text{class}}(\mathbf{P}) = \max(0, 1 - y \max(\text{DET}(\mathbf{P})))^2, y \in \{-1, +1\}$$



Cost Function (1)

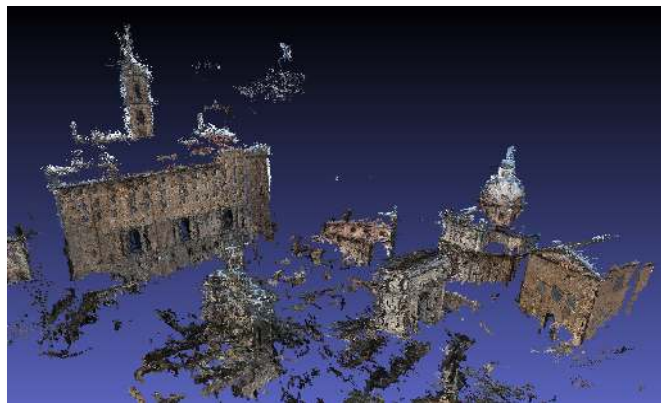


$$\mathcal{L}_{\text{class}}(\mathbf{P}) = \max(0, 1 - y \text{softmax}(\text{DET}(\mathbf{P})))^2, y \in \{-1, +1\}$$

Training with SfM Keypoints



Piccadilly (pic)

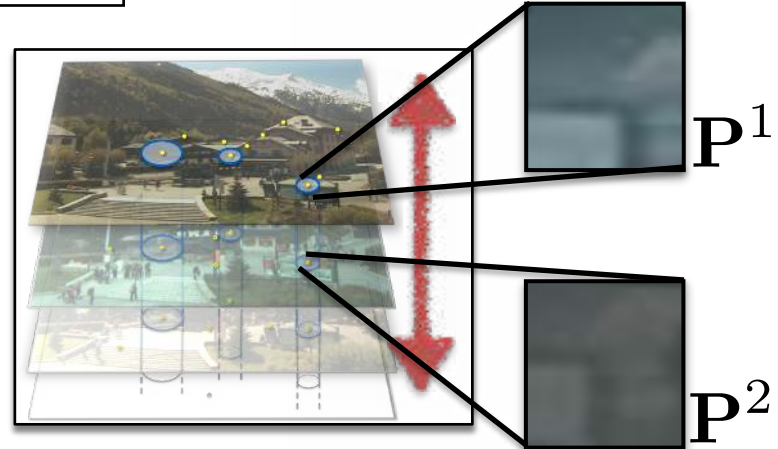


Roman Forum (rf)

- We need variability (illumination, perspective, etc). We build SfM reconstructions from **photo-tourism sets**.
- We keep only **points with SfM correspondences** as positive examples, that is, we **learn to find repeatable points**.

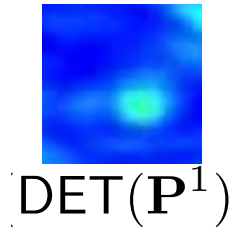
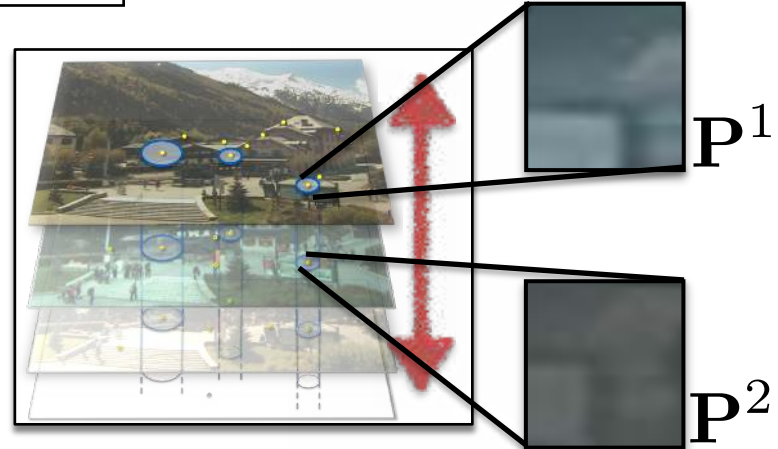


Cost Function (2)



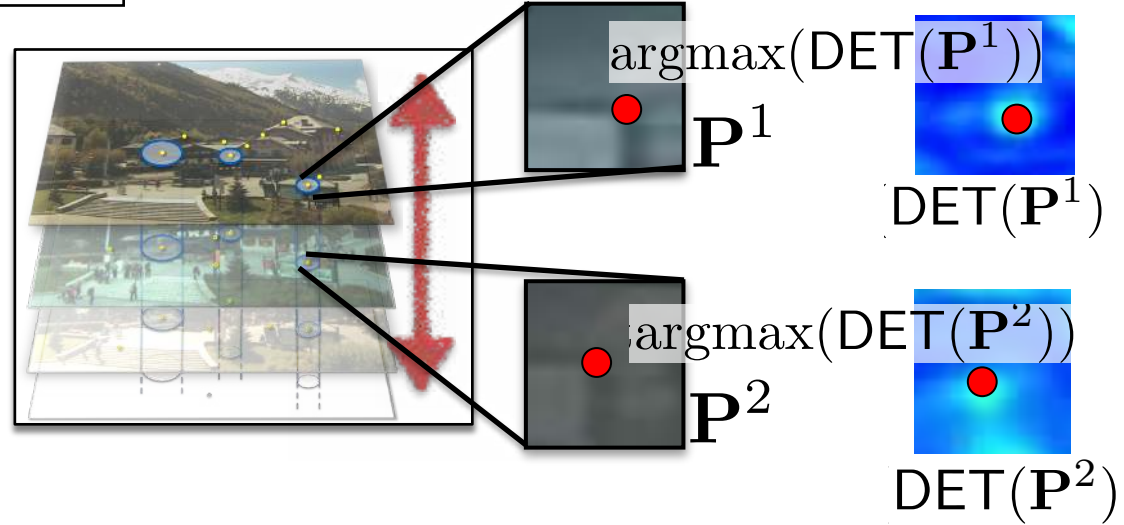


Cost Function (2)



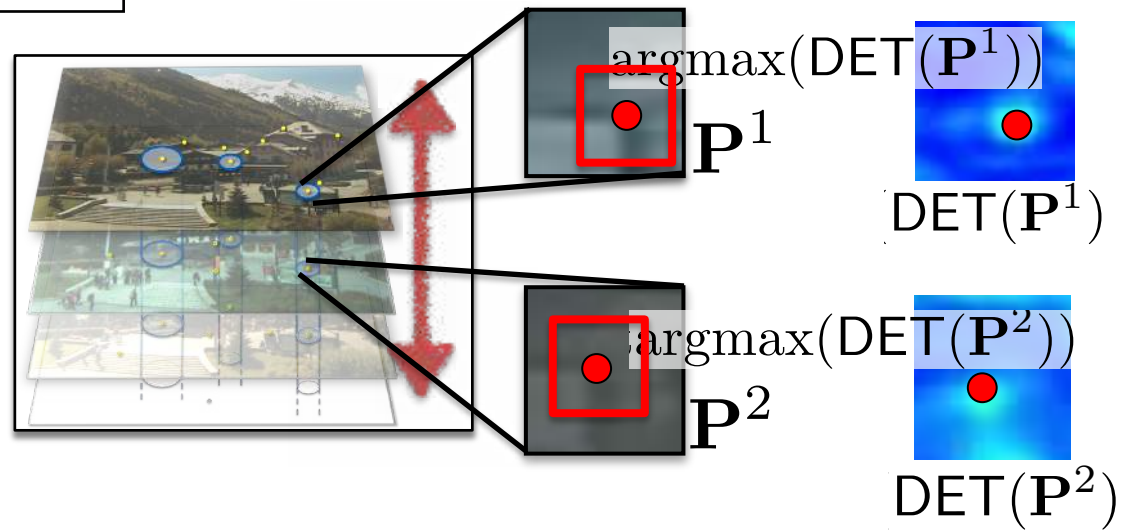


Cost Function (2)





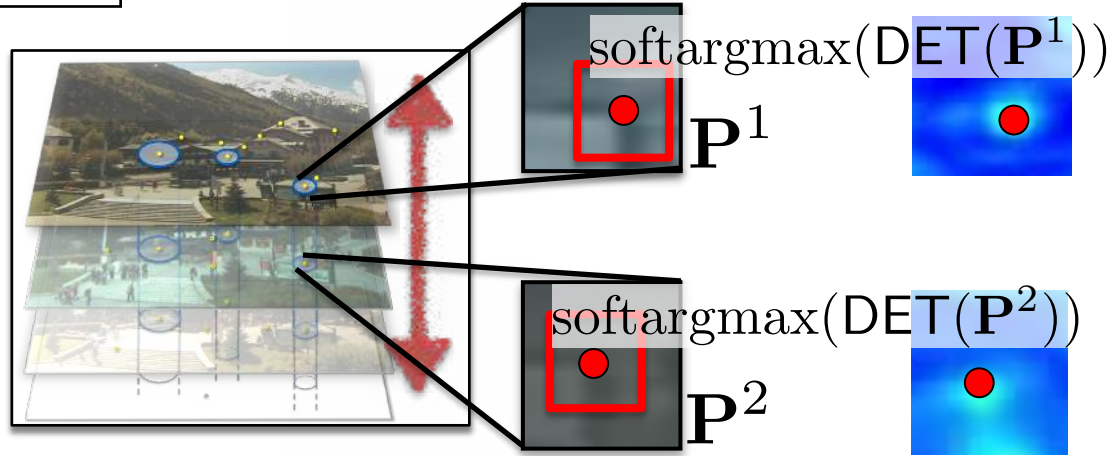
Cost Function (2)



$$\mathcal{L}_{\text{pair}}(\mathbf{P}^1, \mathbf{P}^2) = \left\| \text{DESC}(\text{Crop}(\mathbf{P}^1, \text{argmax}(\text{DET}(\mathbf{P}^1)))) - \text{DESC}(\text{Crop}(\mathbf{P}^2, \text{argmax}(\text{DET}(\mathbf{P}^2)))) \right\|^2$$



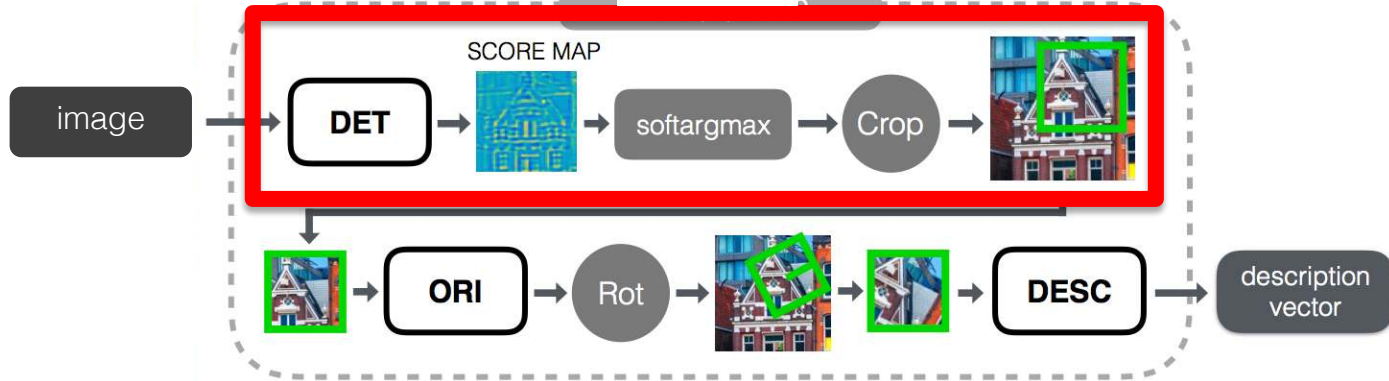
Cost Function (2)

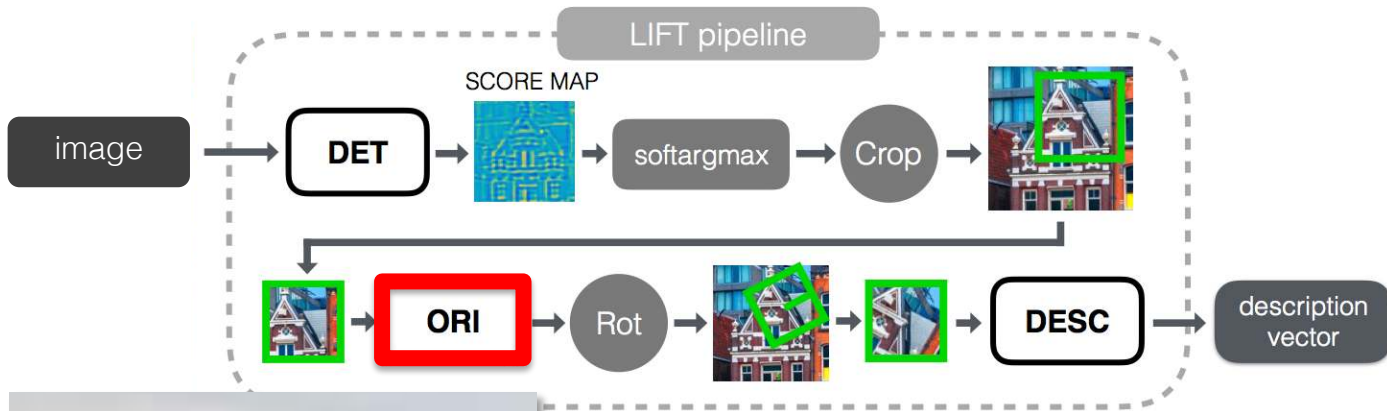


$$\mathcal{L}_{\text{pair}}(\mathbf{P}^1, \mathbf{P}^2) = \left\| \text{DESC}(\text{Crop}(\mathbf{P}^1, \underline{\text{softargmax}}(\text{DET}(\mathbf{P}^1)))) - \text{DESC}(\text{Crop}(\mathbf{P}^2, \underline{\text{softargmax}}(\text{DET}(\mathbf{P}^2)))) \right\|^2$$

$$\text{softargmax}(\mathbf{S}) = \frac{\sum_{\mathbf{x}} \exp(\beta \mathbf{S}(\mathbf{x})) \mathbf{x}}{\sum_{\mathbf{x}} \exp(\beta \mathbf{S}(\mathbf{x}))}$$

so far

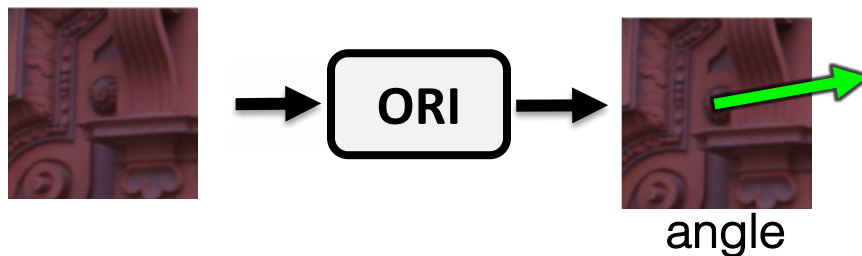






Learning Orientations Implicitly

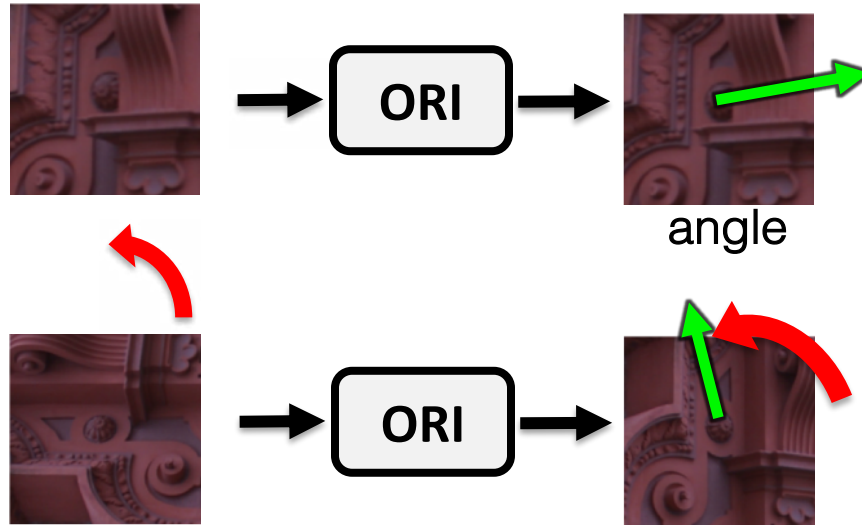
We want the orientation estimator to provide **consistent** results, regardless of imaging changes





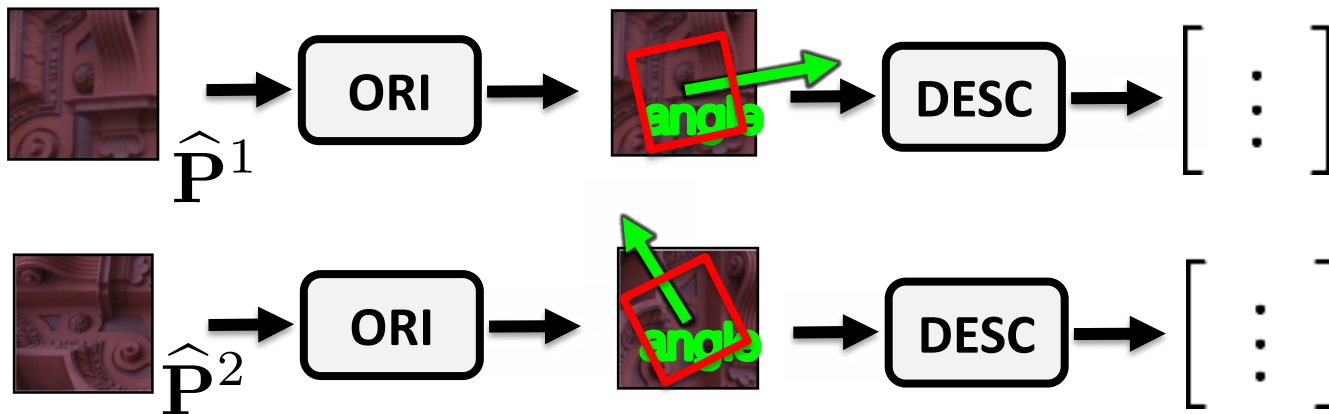
Learning Orientations Implicitly

We want the orientation estimator to provide **consistent** results, regardless of imaging changes





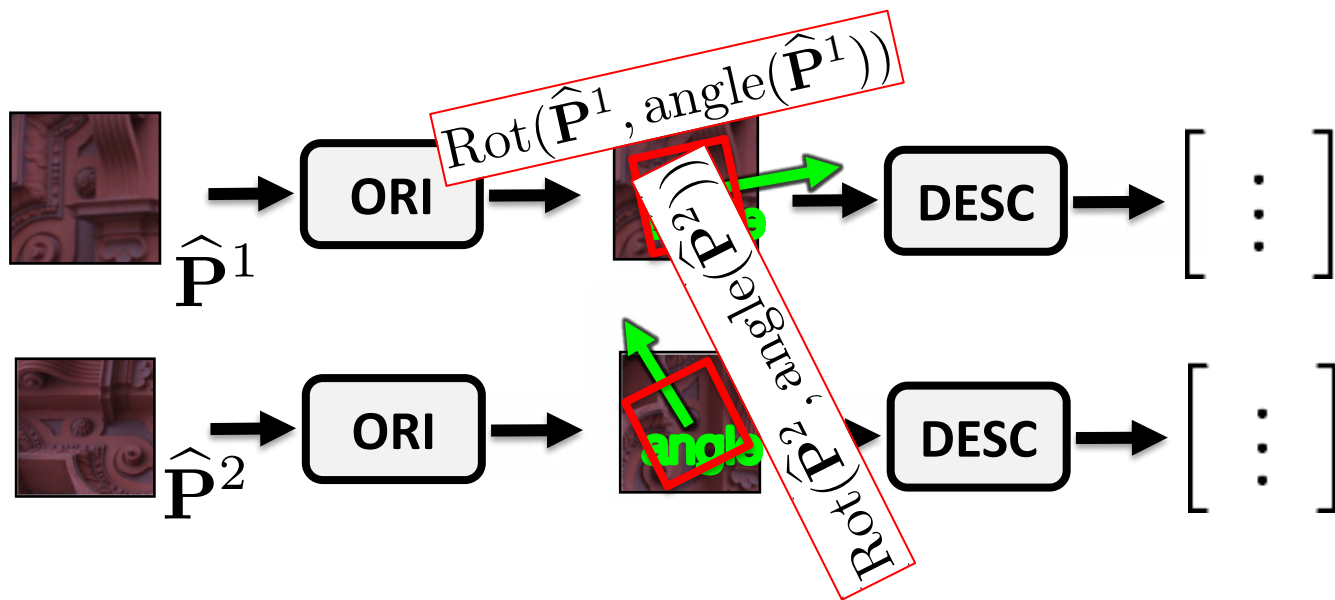
Learning Orientations Implicitly: A Siamese Network with a Twist





Learning Orientations Implicitly: A Siamese Network with a Twist

$$\mathcal{L}_{\text{pair}}(\hat{\mathbf{P}}^1, \hat{\mathbf{P}}^2) = \left\| \text{DESC}(\text{Rot}(\hat{\mathbf{P}}^1, \text{angle}(\hat{\mathbf{P}}^1))) - \text{DESC}(\text{Rot}(\hat{\mathbf{P}}^2, \text{angle}(\hat{\mathbf{P}}^2))) \right\|^2$$

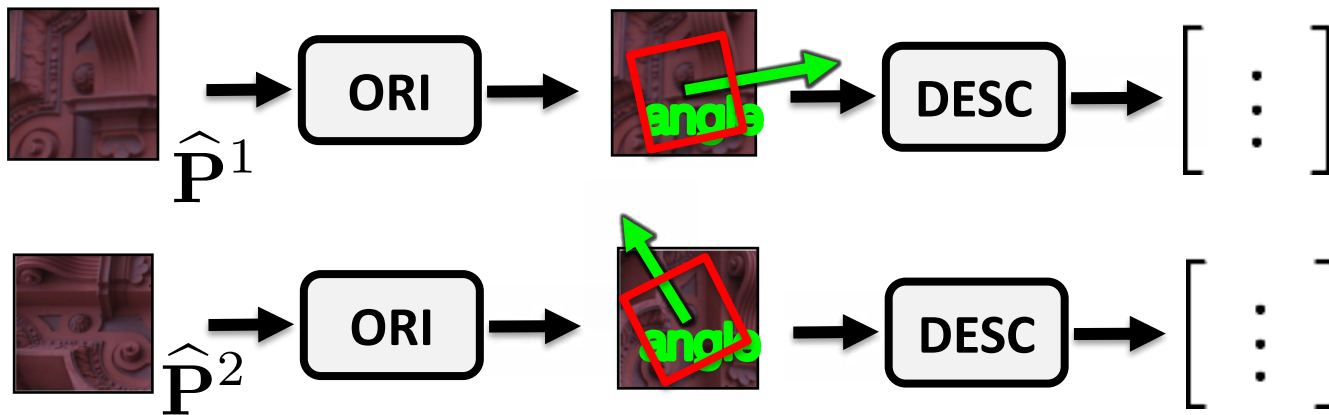




Learning Orientations Implicitly: A Siamese Network with a Twist

$$\mathcal{L}_{\text{pair}}(\hat{\mathbf{P}}^1, \hat{\mathbf{P}}^2) = \left\| \begin{array}{c} \text{DESC}(\text{Rot}(\hat{\mathbf{P}}^1, \text{angle}(\hat{\mathbf{P}}^1))) \\ - \\ \text{DESC}(\text{Rot}(\hat{\mathbf{P}}^2, \text{angle}(\hat{\mathbf{P}}^2))) \end{array} \right\|^2$$

with $\text{angle}(\hat{\mathbf{P}}) = \text{arctan2}(\text{ORI}(\hat{\mathbf{P}})^{[1]}, \text{ORI}(\hat{\mathbf{P}})^{[2]})$





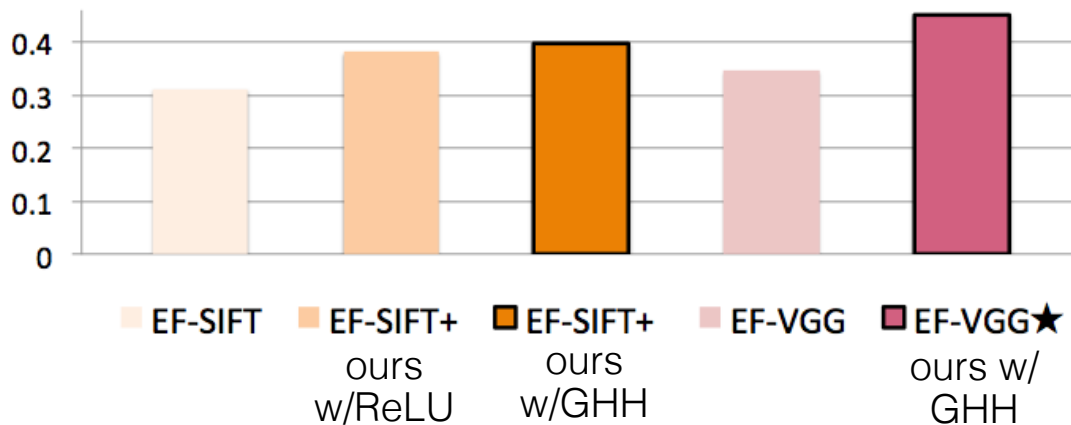
Performance Gain with Learned Orientations

mA

P

(higher is better)

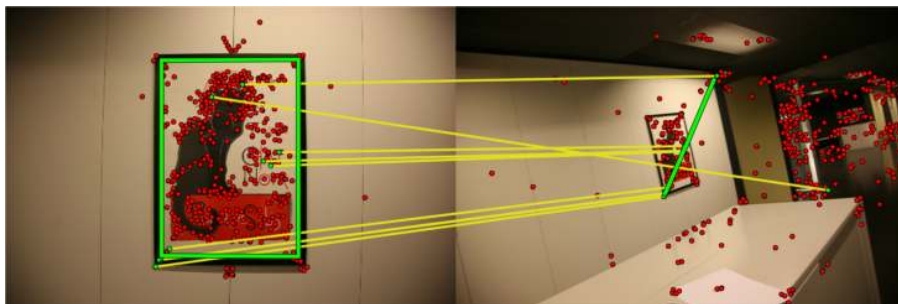
Average performance



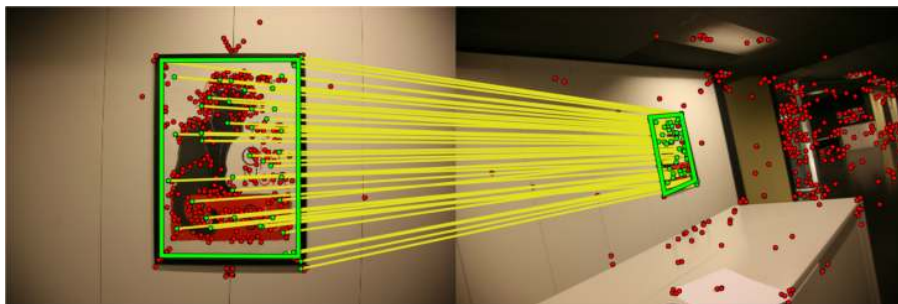
Descriptor matching performances (mAP) with nearest neighbor matching (Mikolajczyk & Schmid, IJCV'04).



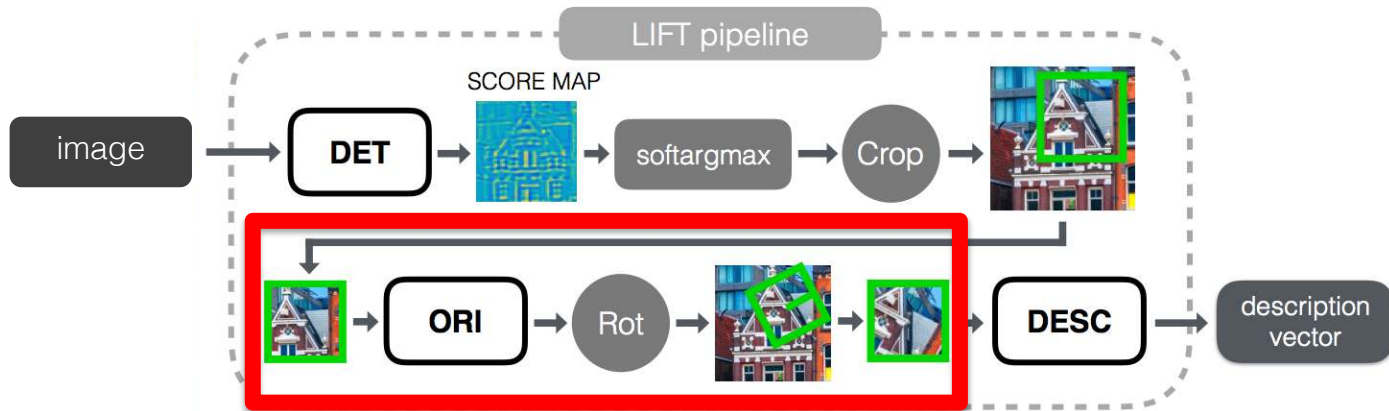
Learned Orientations

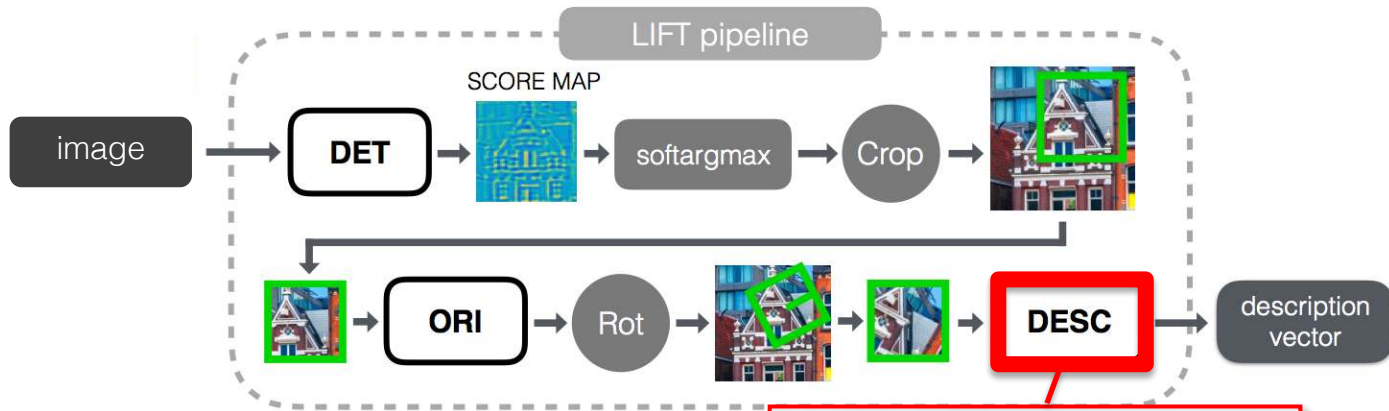


Dominant Gradient Orientations



Our Learned Orientations







Learning the Descriptor

- Positive pairs:



$\widehat{\mathbf{P}}^1$



$\widehat{\mathbf{P}}^2$

$$\mathcal{L}_{\text{pos}}(\widehat{\mathbf{P}}^1, \widehat{\mathbf{P}}^2) = \|\text{DESC}(\widehat{\mathbf{P}}^1) - \text{DESC}(\widehat{\mathbf{P}}^2)\|^2$$



Learning the Descriptor

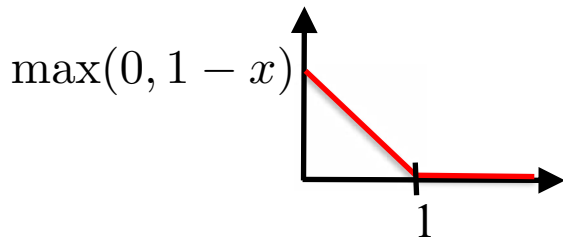
- Positive pairs:



$$\mathcal{L}_{\text{pos}}(\widehat{\mathbf{P}}^1, \widehat{\mathbf{P}}^2) = \|\text{DESC}(\widehat{\mathbf{P}}^1) - \text{DESC}(\widehat{\mathbf{P}}^2)\|^2$$

- Negative pairs:

$$\mathcal{L}_{\text{neg}}(\widehat{\mathbf{P}}^1, \widehat{\mathbf{P}}^3) = \max(0, 1 - \|\text{DESC}(\widehat{\mathbf{P}}^1) - \text{DESC}(\widehat{\mathbf{P}}^3)\|^2)$$



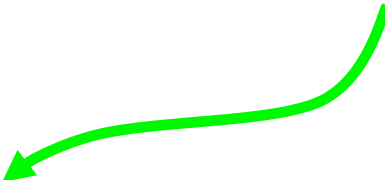
Hard example mining is very important for training

A Single, Global Cost Function

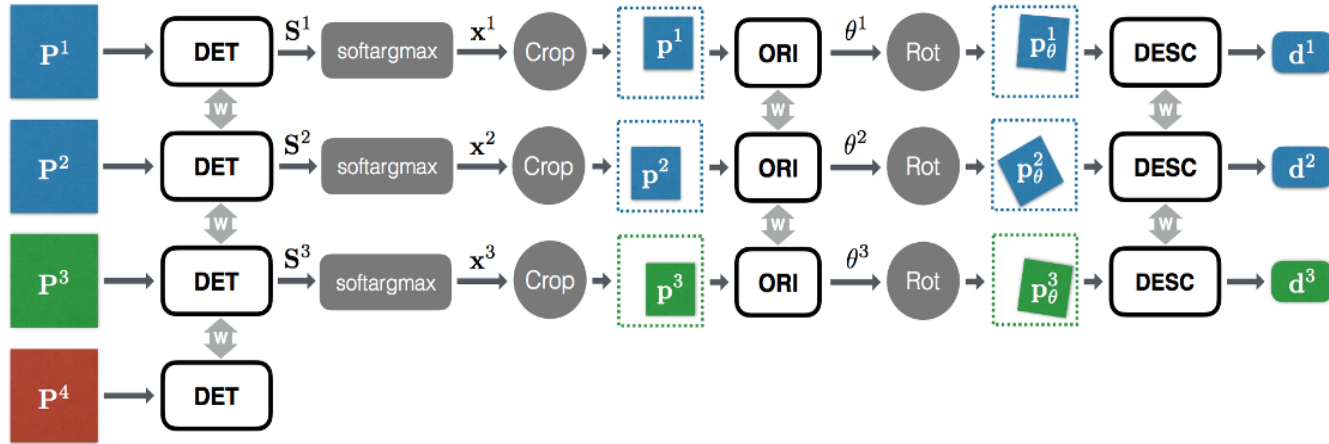
$$\min_{\{\text{DET}, \text{ORI}, \text{DESC}\}} \sum_{\{(\mathbf{P}, y)\}} \max(0, 1 - y \text{softmax}(\text{DET}(\mathbf{P})))^2 +$$

$$\sum_{(\mathbf{P}_1, \mathbf{P}_2)} \|\text{DESC}(G(\mathbf{P}^1, \text{softargmax}(\text{DET}(\mathbf{P}^1)))) - \text{DESC}(G(\mathbf{P}^2, \text{softargmax}(\text{DET}(\mathbf{P}^2))))\|^2 +$$

$$\sum_{(\mathbf{P}_1, \mathbf{P}_3)} \max(0, 1 - \|\text{DESC}(G(\mathbf{P}^1, \text{softargmax}(\text{DET}(\mathbf{P}^1)))) - \text{DESC}(G(\mathbf{P}^3, \text{softargmax}(\text{DET}(\mathbf{P}^3))))\|^2)$$


$$G(\mathbf{P}, \mathbf{x}) = \text{Rot}(\mathbf{P}, \mathbf{x}, \text{angle}_{\text{ORI}}(\text{Crop}(\mathbf{P}, \mathbf{x})))$$

Problem-Specific Training

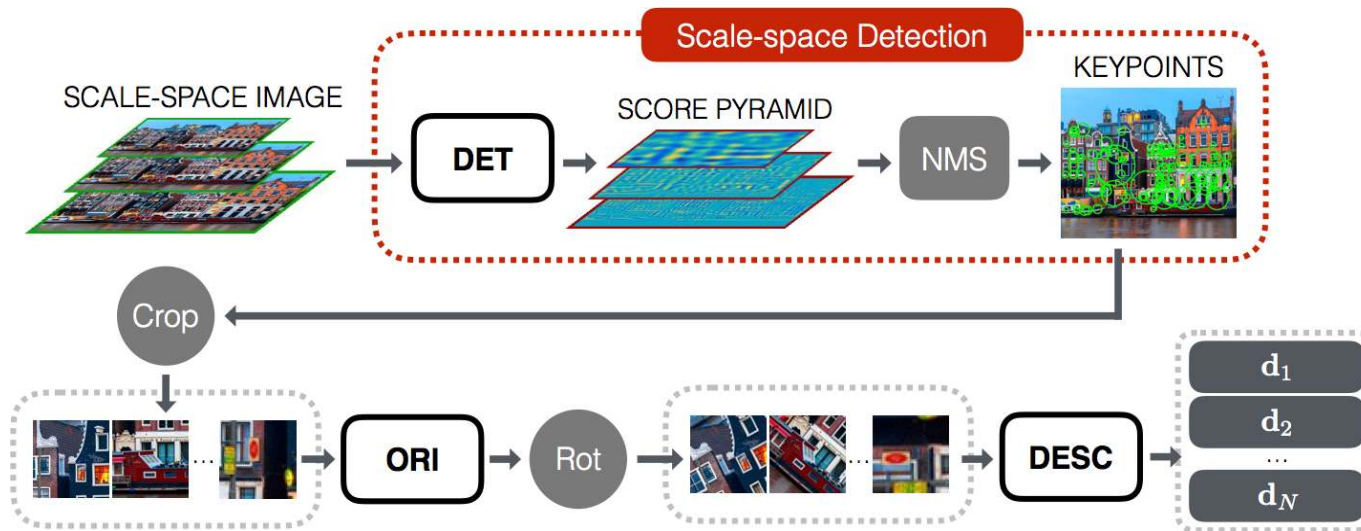


← LEARNING

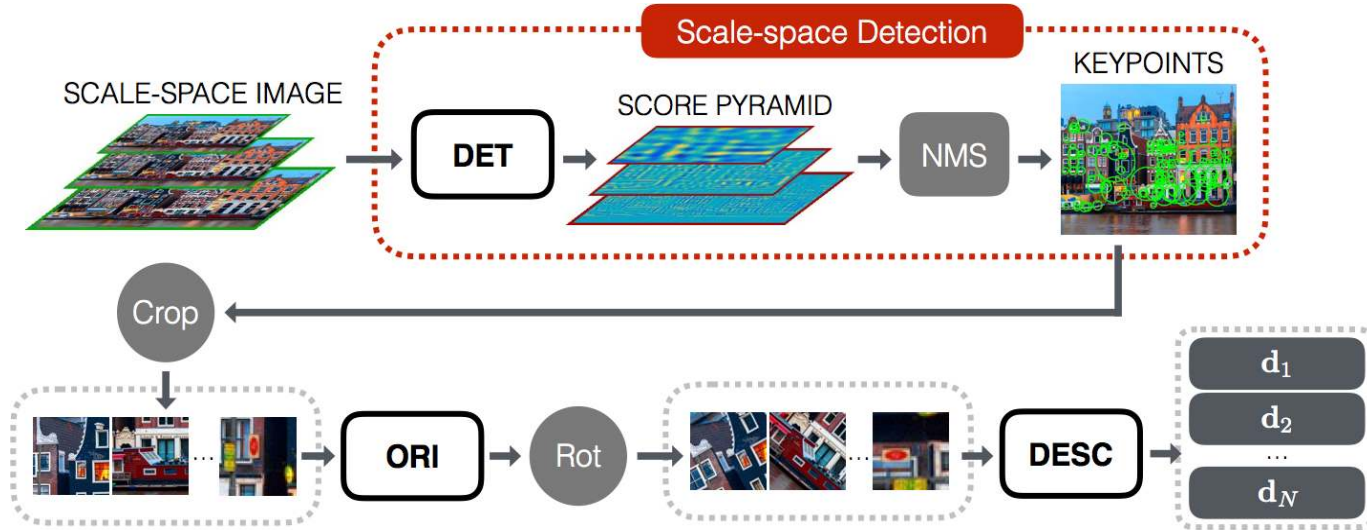
Run-Time Pipeline

Detector

- is purely convolutional, efficiently applied to the whole image;
- works in scale-space.



Run-Time Pipeline



The orientation estimator and the descriptor are applied only to keypoints.