

SuperPoint

Friday, December 27, 2019 9:28 PM

- ↗ Fully convolutional n/w that computes pixel-level interest point locations and descriptors.
- ↗ state of the art homography estimation results on HPatches compared to LIFT, SIFT and ORB.

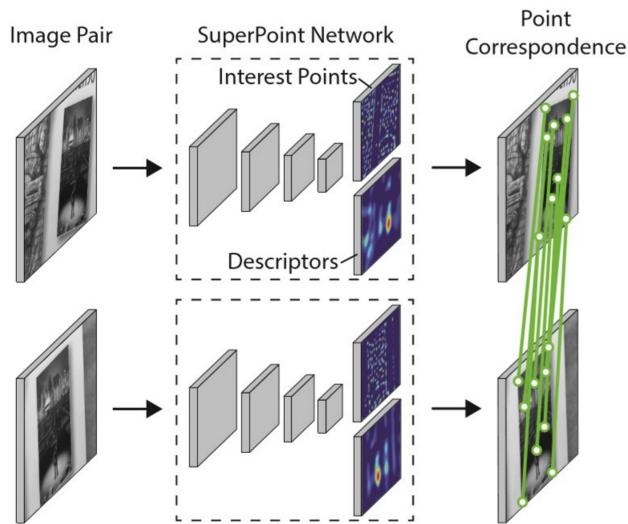


Figure 1. SuperPoint for Geometric Correspondences. We present a fully-convolutional neural network that computes SIFT-like 2D interest point locations and descriptors in a single forward pass and runs at 70 FPS on 480×640 images with a Titan X GPU.

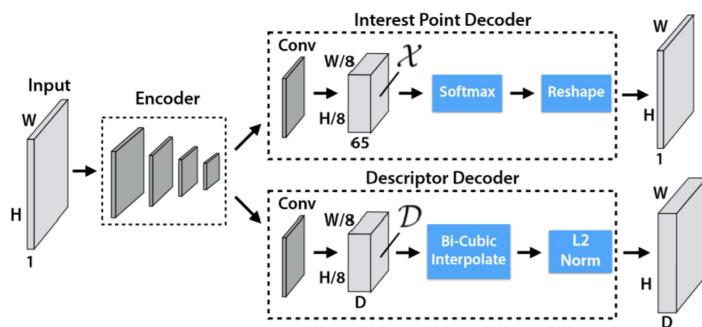


Figure 3. SuperPoint Decoders. Both decoders operate on a shared and spatially reduced representation of the input. To keep the model fast and easy to train, both decoders use non-learned upsampling to bring the representation back to $\mathbb{R}^{H \times W}$.

- 1 One encoder and 2 decoders with shared weights.
- 1 Encoder \rightarrow VGG-style with Max-pooling layers.
- 1 Interest point decoder \rightarrow uses sub-pixel convolution or depth to space (TF) or Pixel-shuffle (PyTorch).
 - \rightarrow this has much lower computation as compared to traditional deconv. layers.
- 1 Descriptor \rightarrow semi-dense, normalized to unit length, Bicubic interpolated to get a dense map.

2

$$\mathcal{L}(\mathcal{X}, \mathcal{X}', \mathcal{D}, \mathcal{D}'; Y, Y', S) = \mathcal{L}_p(\mathcal{X}, Y) + \mathcal{L}_p(\mathcal{X}', Y') + \lambda \mathcal{L}_d(\mathcal{D}, \mathcal{D}', S). \quad (1)$$

\mathcal{L}_p \Rightarrow Cross entropy loss on Ω_1 in contact point

- - Try using
pseudo-GT interest point
locations
- $L_D \rightarrow$ Descriptor loss b/w all
pairs of descriptors b/w 2
images
- 2) Synthetic pre-training
Synthetic shapes dataset with
quadrilaterals, triangles, lines and
ellipses. Use L, T and Y junctions
as interest points.
- Magic Point \rightarrow Only the detector
pathway.
This worked extremely well
on synthetic images and
generalized reasonably to
images with strong corners.
However didn't work well

However didn't work well
on natural images.

Hence a self-supervised approach
was used called Homography
Adaptation.

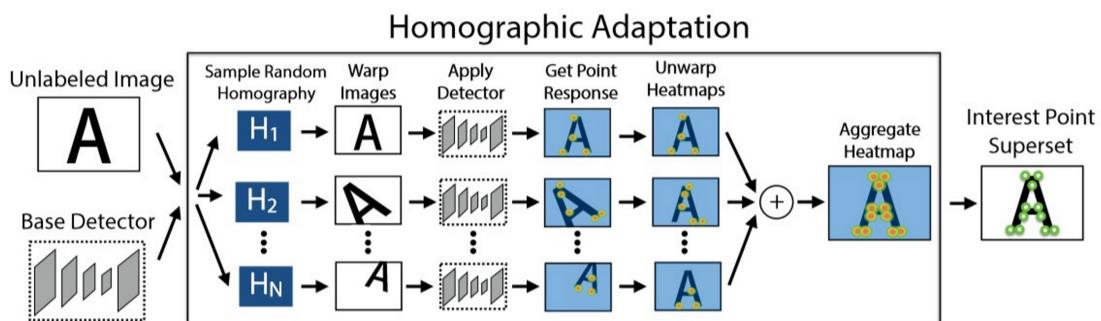


Figure 5. **Homographic Adaptation.** Homographic Adaptation is a form of self-supervision for boosting the geometric consistency of an interest point detector trained with convolutional neural networks. The entire procedure is mathematically defined in Equation 10.

$$\mathbf{x} = \mathcal{H}^{-1} f_{\theta}(\mathcal{H}(I)).$$

- The interest point should be covariant w.r.t transform
- An aggregation over a large number of H_i 's gives rise to SuperPoint

$$\hat{F}(I; f_{\theta}) = \frac{1}{N_h} \sum_{i=1}^{N_h} \mathcal{H}_i^{-1} f_{\theta}(\mathcal{H}_i(I)). \quad (10)$$

→ To have plausible camera transformations

H is decomposed into translation, scale, rotation (yaw) and symmetric distortion.

100 homographies for image are used.

→ Architecture →

VGG like 8 conv layers.

Encoder →

$$\begin{matrix} (3 \times 3) \\ \text{all} \end{matrix} C_{64} \rightarrow C_{64} \rightarrow M_{\downarrow 2} \rightarrow C_{64} \rightarrow C_{64} \rightarrow M_{\downarrow 2}$$

↓

$$C_{128} \leftarrow M_{\uparrow 2} \leftarrow C_{128} \leftarrow C_{128}$$

$$\begin{matrix} \downarrow \\ C_{128} \end{matrix} \longrightarrow M_{\uparrow 2}$$

Decoder →

$$\begin{matrix} C_{256} \\ 3 \times 3 \end{matrix} \rightarrow \begin{matrix} C \\ 65 \text{ or } 256 \\ 1 \times 1 \end{matrix}$$

$$C \rightarrow \text{Conv} \rightarrow \text{ReLU} \rightarrow \text{BN}$$

Trained on synthetic shapes.

$H = \text{Matrix} + (\text{Zero or } 1)$ is

Trained on synthetic shapes.
Then Magicpoint (prev. op) is
then trained with descriptor
head on MS-COCO with
Homography adaptation (100 H per
image).

Homography adaptation is done
again.

Data is augmented using Gaussian
noise, motion blur and brightness level
changes.

- For a 640×480 image to
sample 1000 feature points,
SuperPoint takes ~ 13 ms on
a TitanXp. (70 FPS).