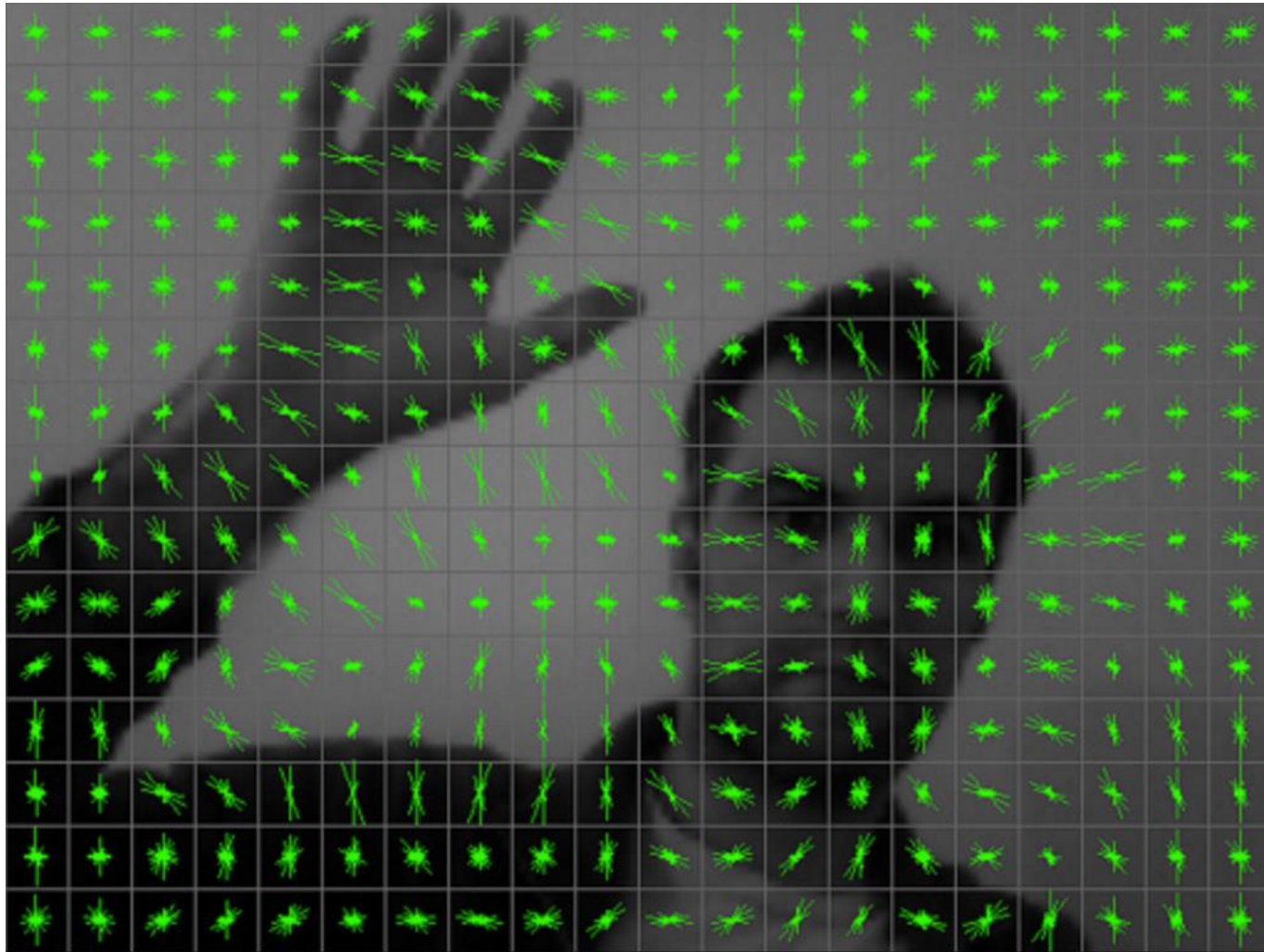
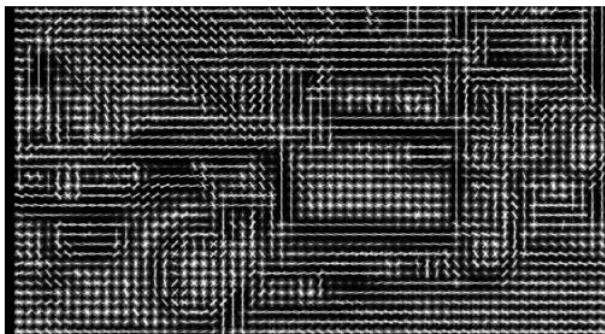


Hog



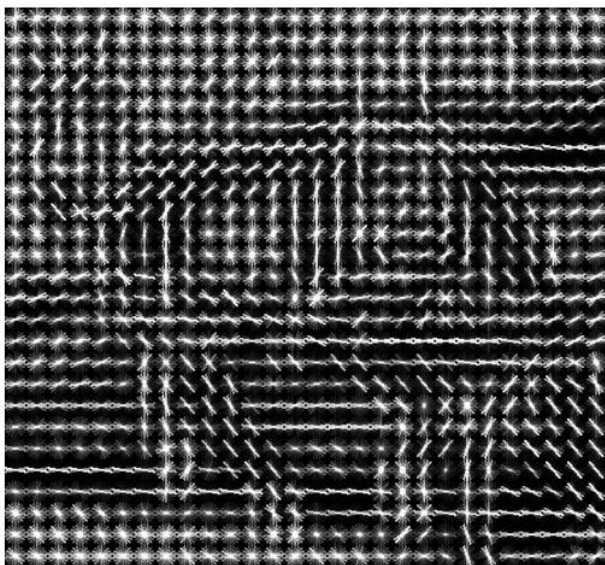
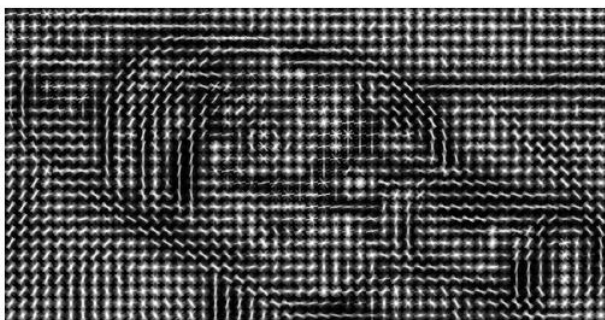
HOG [1]

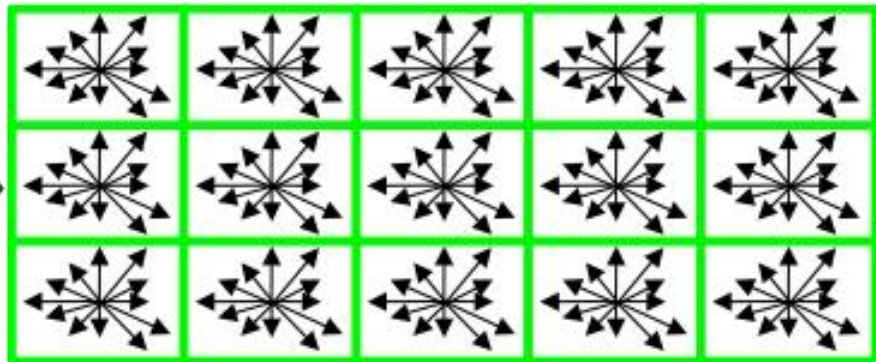
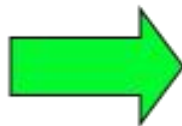
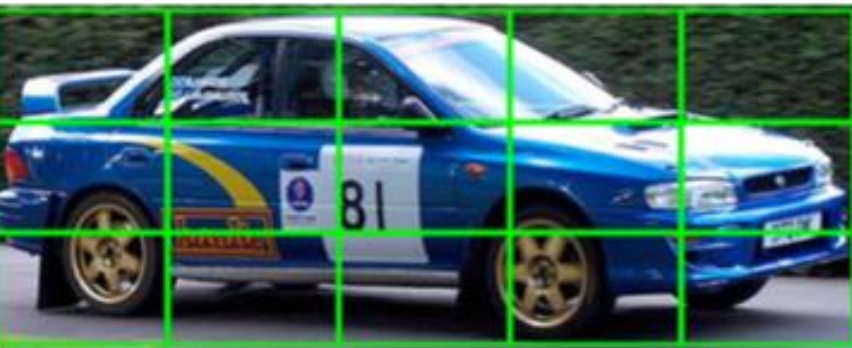


Inverse (Us)

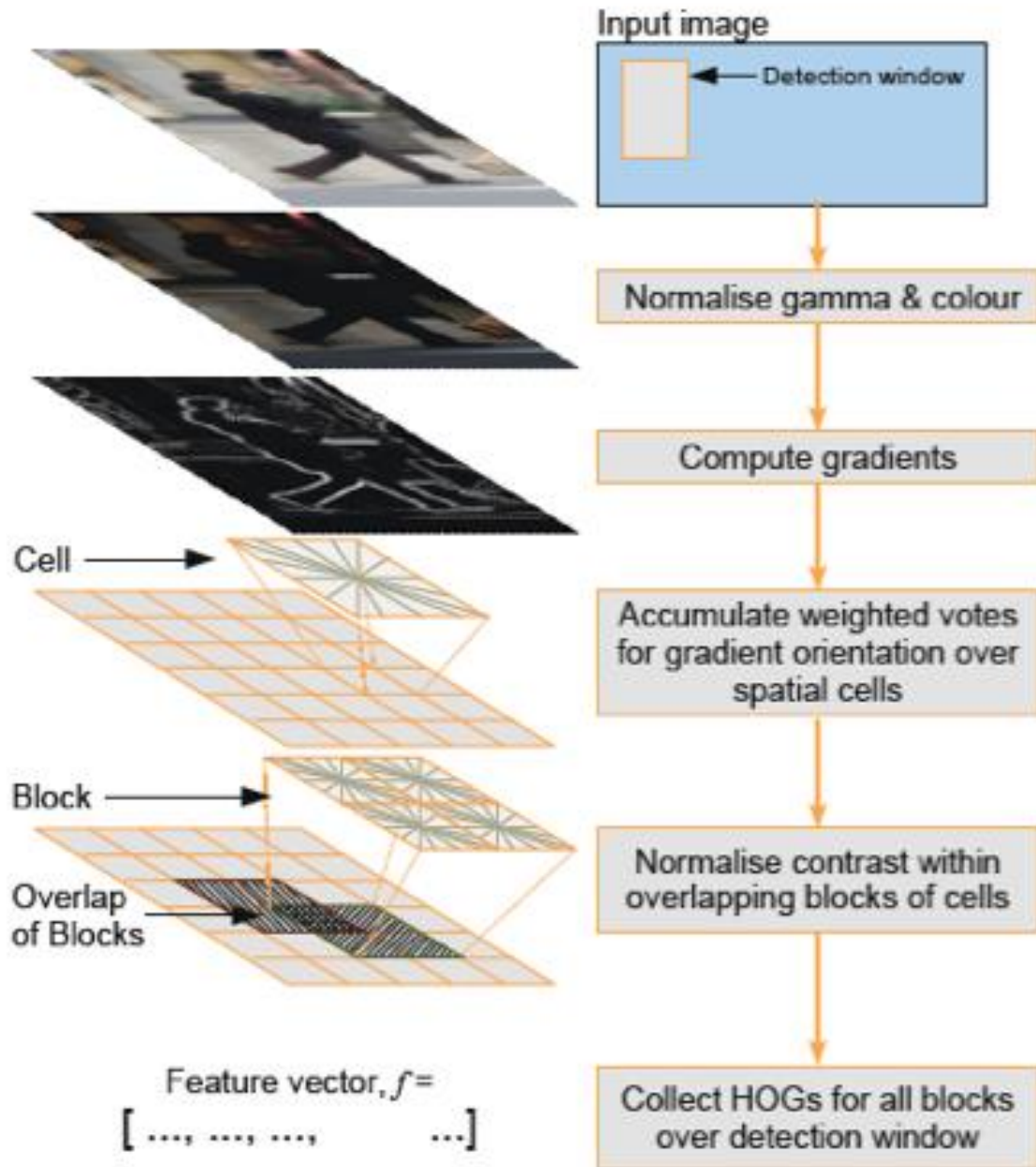


Original





HOG Features



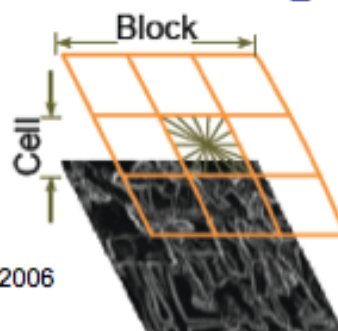
Create cell histograms

- Each pixel in cell casts weighted vote based on gradient magnitude centered there
 - Weighted by applying a Gaussian spatial window to each pixel before accumulating orientation votes into cells $\rightarrow (\sigma = .5 * \text{block width})$
- Votes are accumulated in 9 Histogram channels (orientation bins) spread evenly over 0-180 degrees (Or 0-360 degrees if signed values desired)

"Human Detection PHD Thesis" Navneet Dalal 2006

■ Descriptor Blocks

- To account for illumination/contrast changes the cells must be grouped into "blocks" and normalized
- HOG descriptor is a vector of components of normalized cell histograms from all the block regions
- Author's optimum R-HOG (10% miss rate)
 - 3 parameters



(a) R-HOG/SIFT

- 3x3 cell blocks
- 6x6 pixel cells
- 9 histogram channels (orientation bins)

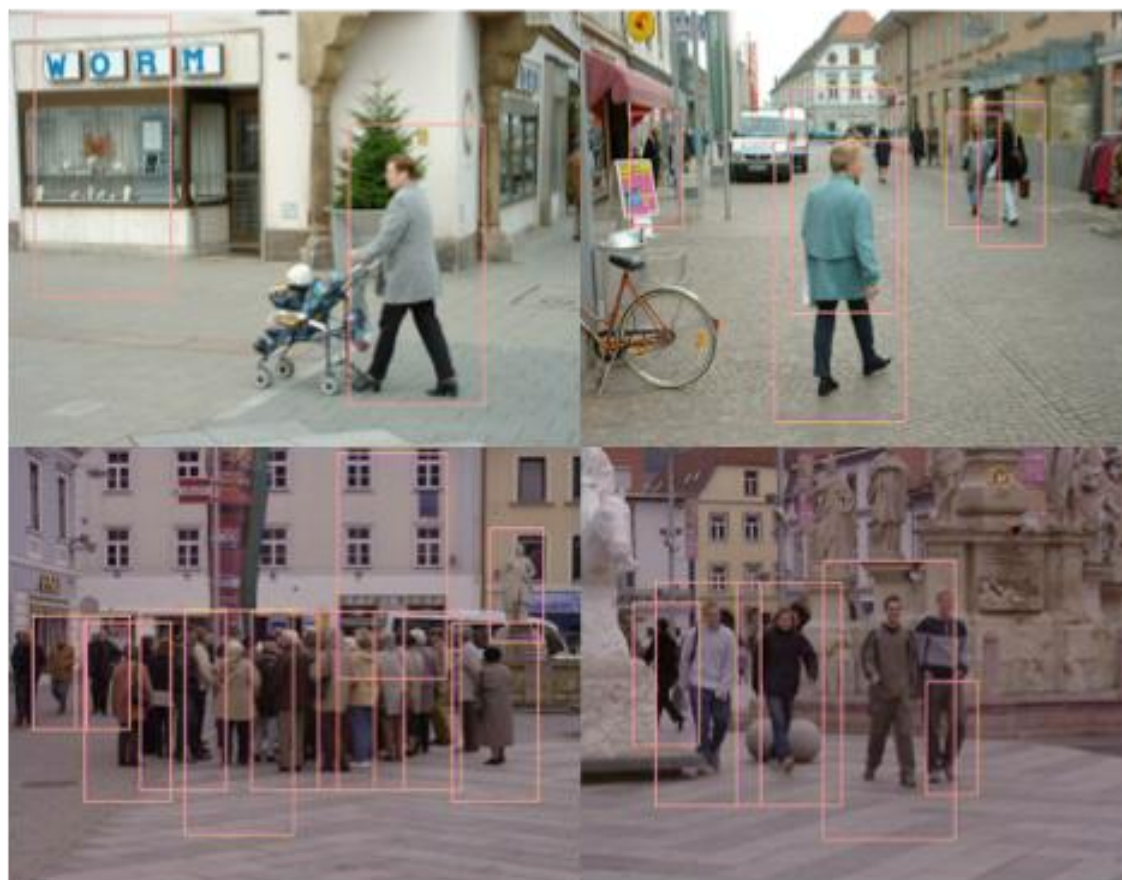
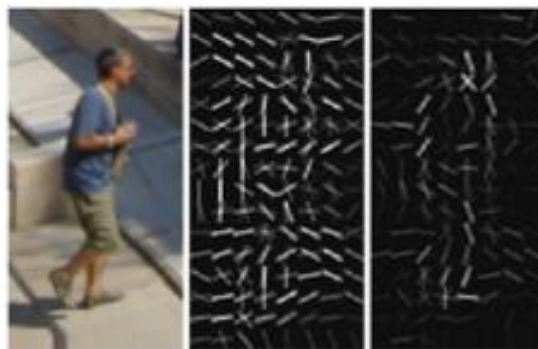
Normalize the Blocks

- V is vector containing non-normalized histogram data and e is a small constant (Not very important over the larger ranges – $1e^{-3}$ to $5e^{-2}$)

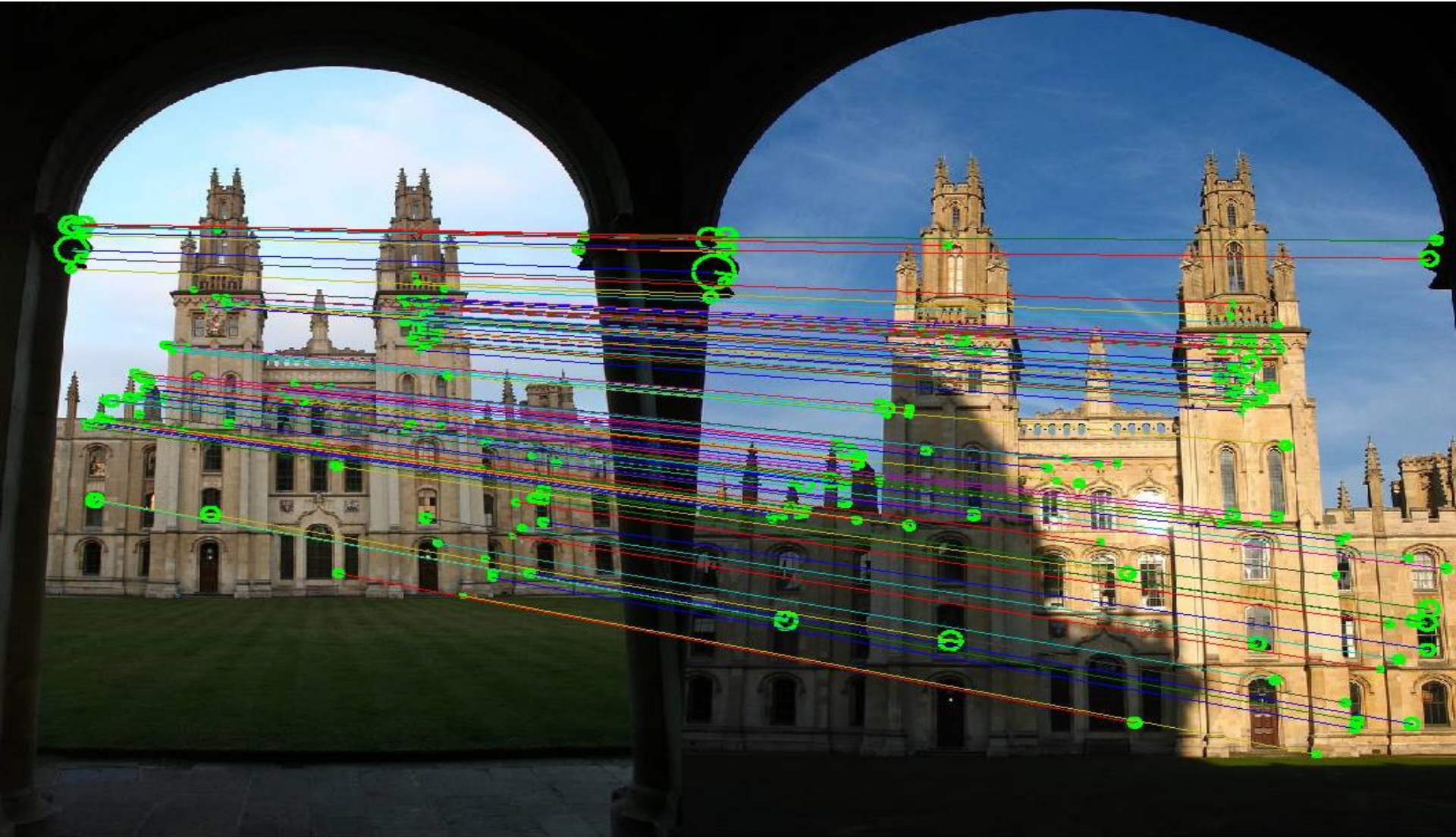
■ Typical Detector Window

- Authors used 64×128 detection window
- 16 pixels of margin around person on all four sides
- Decreasing window size or person size in image decreases performance

Importance weighted responses



SIFT



SIFT vector formation

- Thresholded image gradients are sampled over 16x16 array of locations in scale space
- Create array of orientation histograms
- 8 orientations x 4x4 histogram array = 128 dimensions

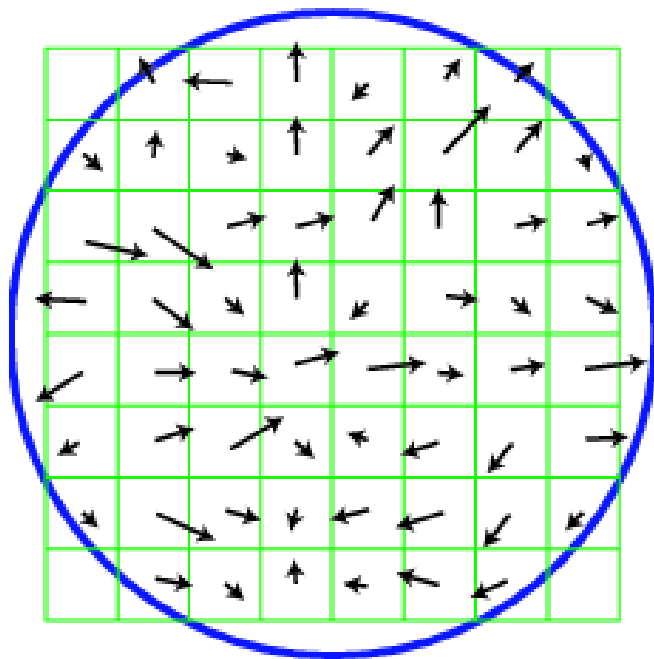
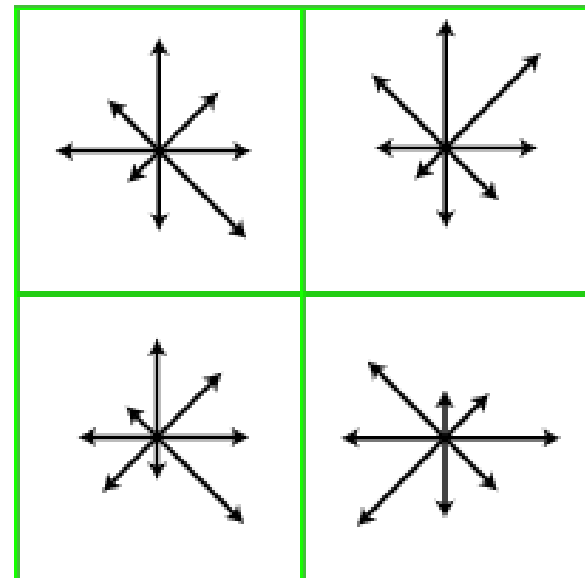


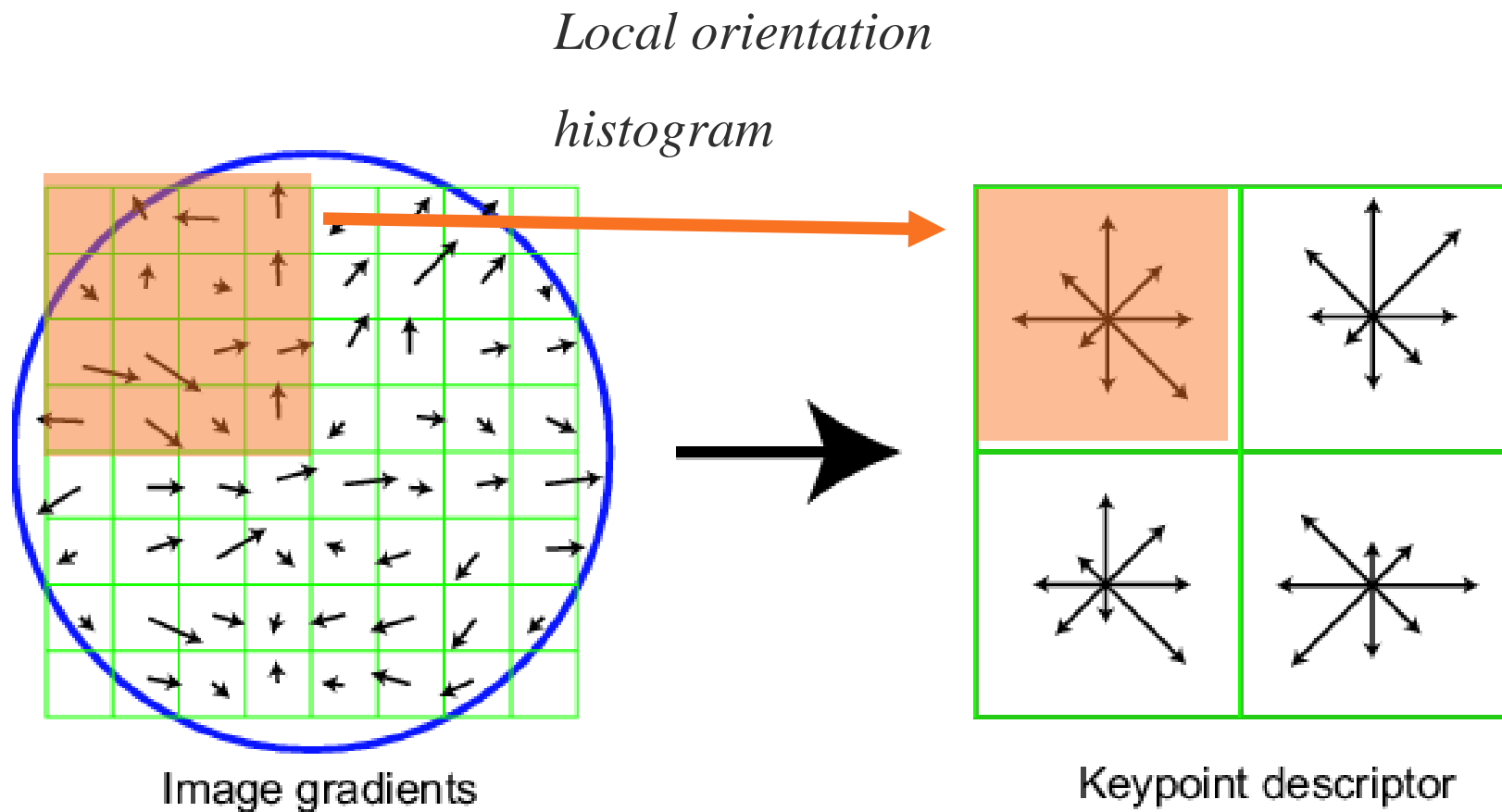
Image gradients



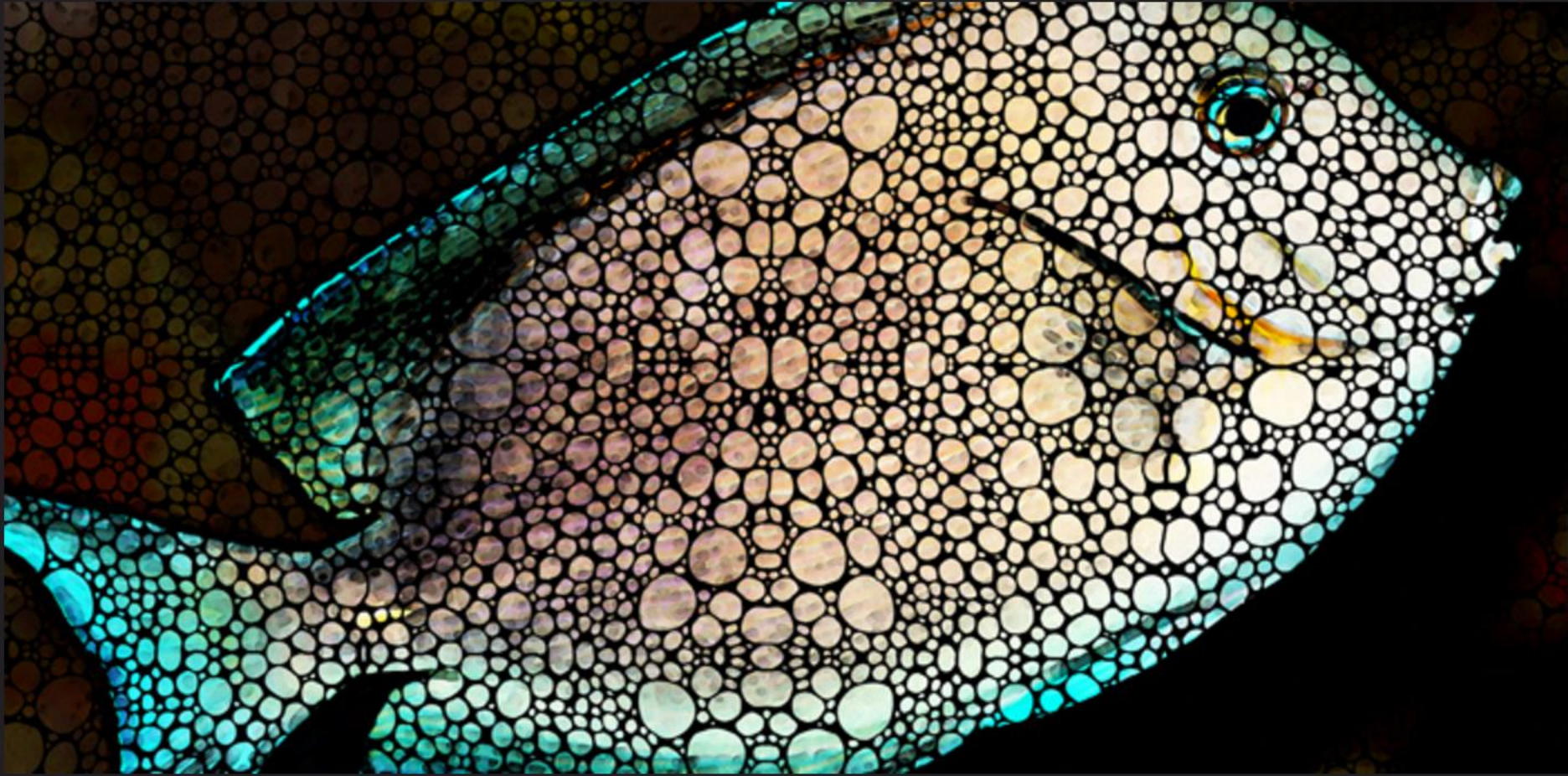
Keypoint descriptor

SIFT vector formation

- Orientation is defined relative to the orientation of the detected Sift feature

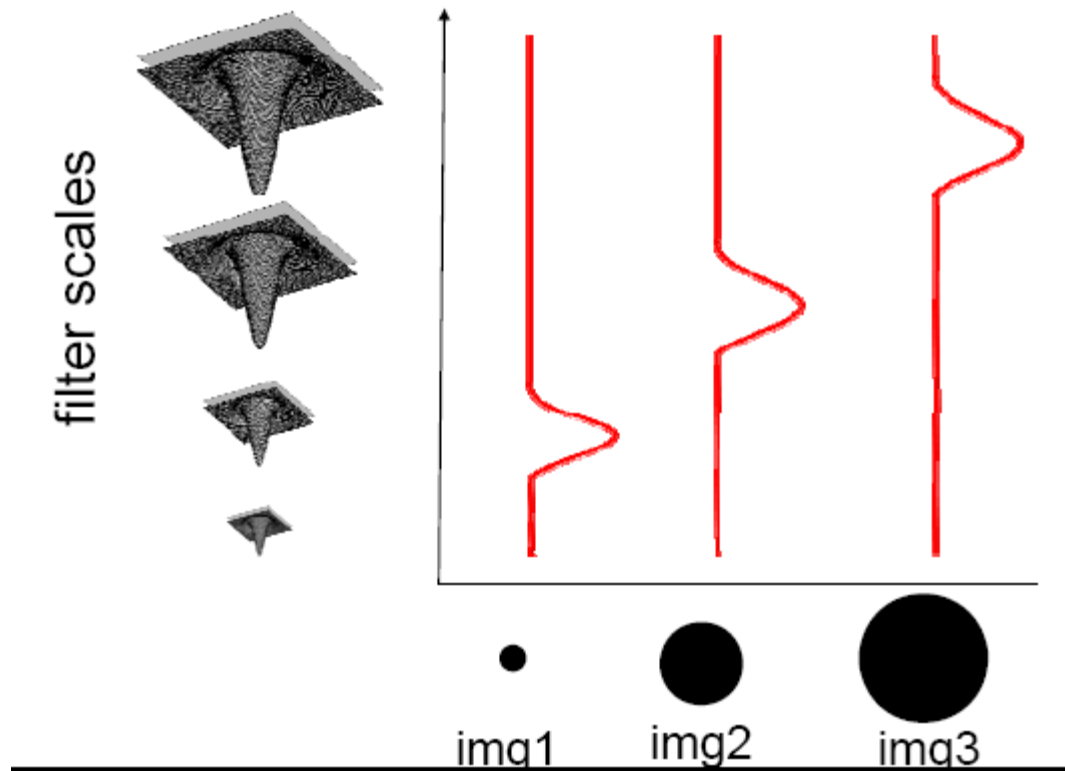


SIFT Fish



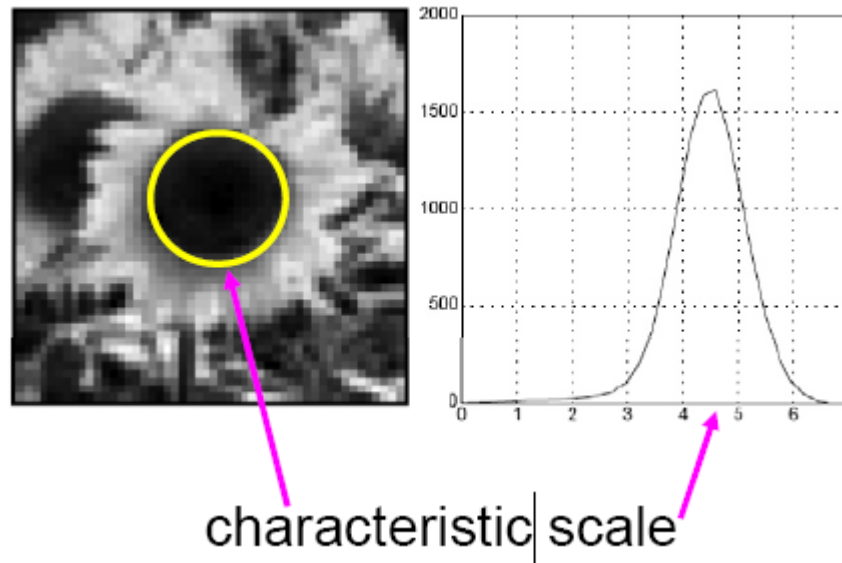
Sharon Cummings @ flickr

Laplacian-of-Gaussian = “*blob*” detector



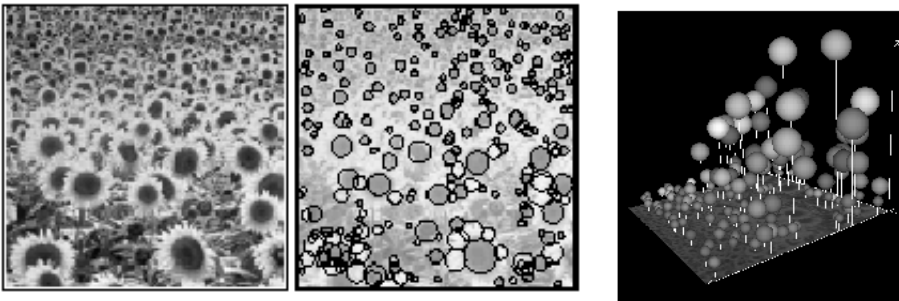
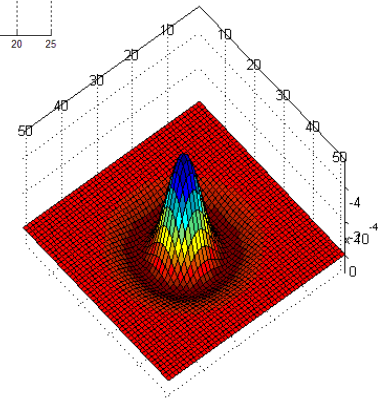
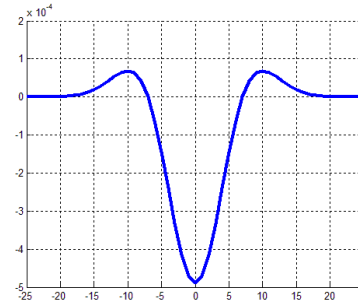
At a given point in the image:

- We define the *characteristic scale* as the scale that produces peak of Laplacian response



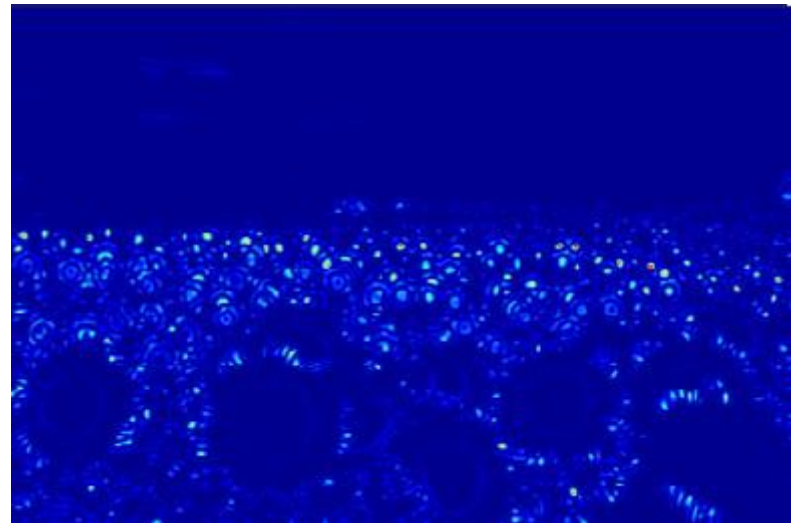
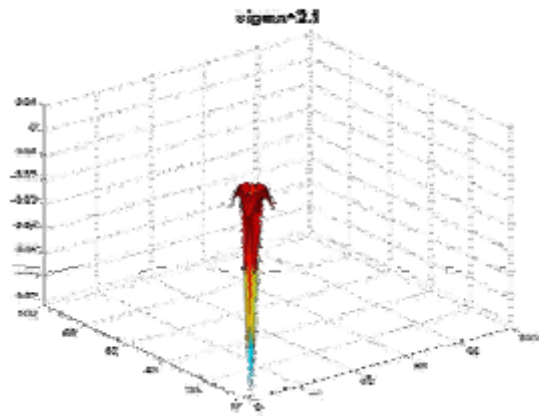
Lowe's Scale-space Interest Points

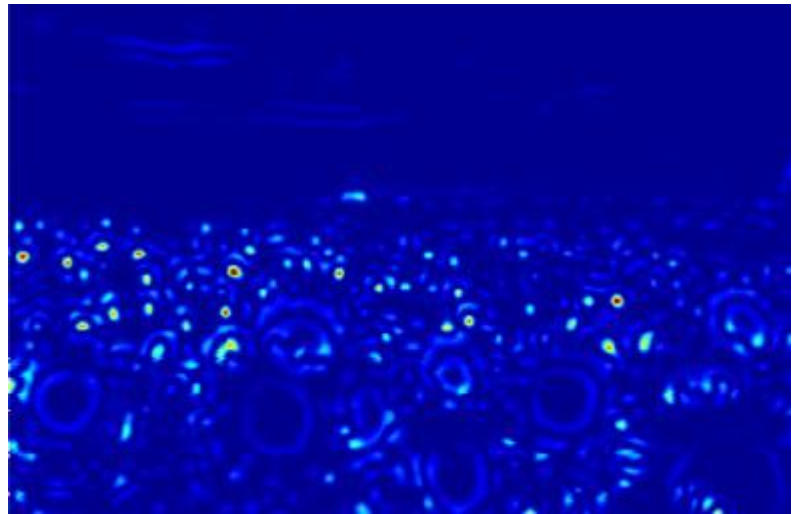
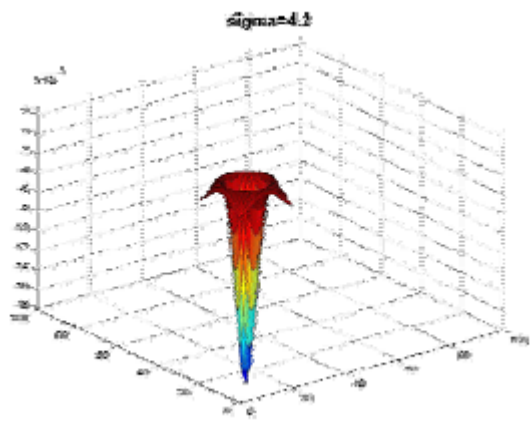
- **Laplacian of Gaussian kernel**
- **Scale-space detection**
 - Find local maxima across scale space
 - A good “blob” detector

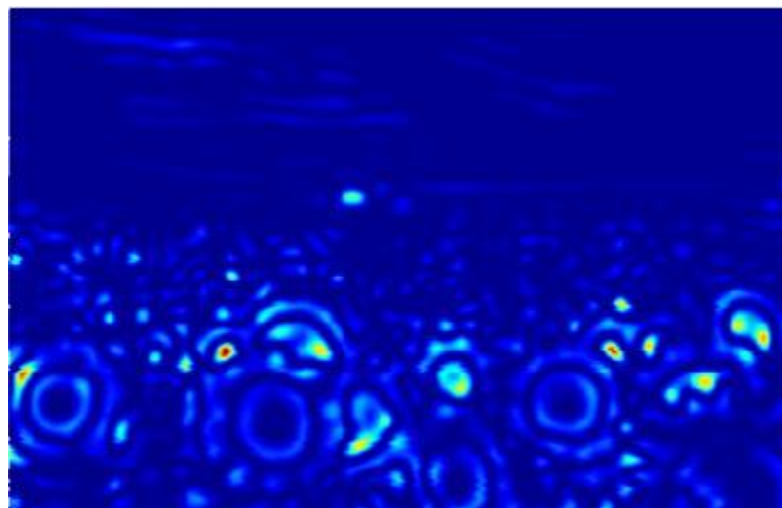
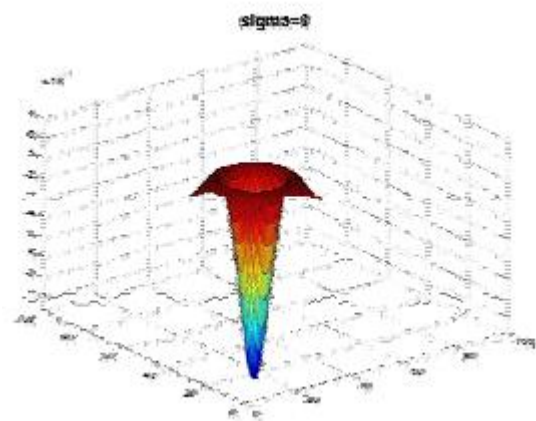


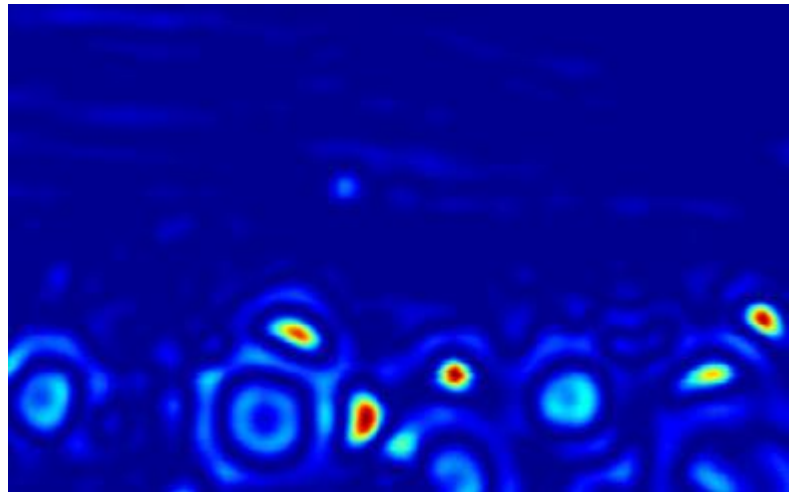
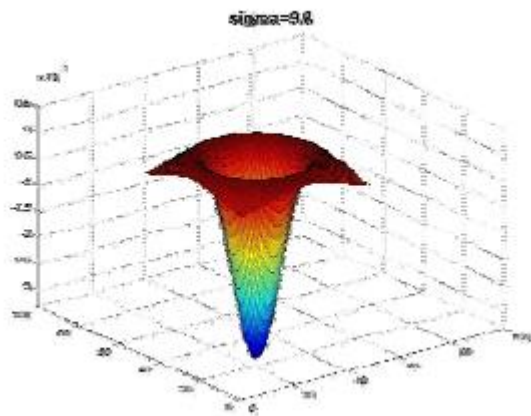
$$G(x, y, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{x^2+y^2}{\sigma^2}}$$

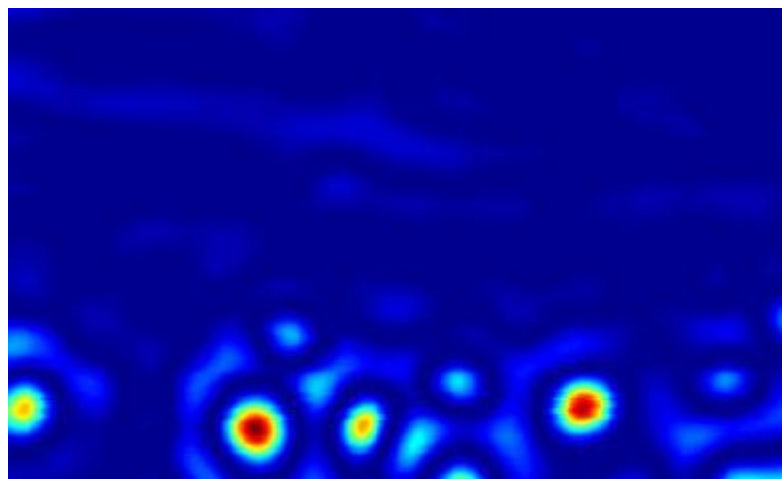
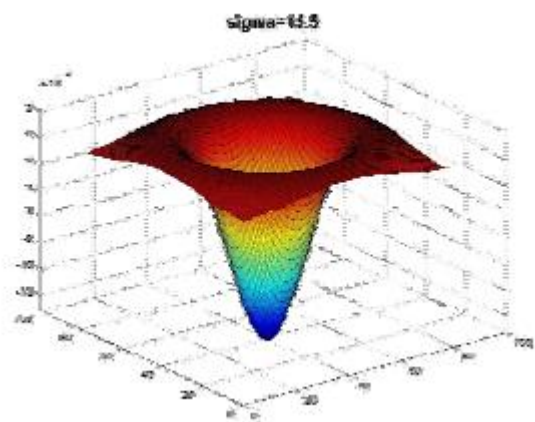
$$\nabla^2 G(x, y, \sigma) = \frac{\partial^2 G}{\partial x^2} + \frac{\partial^2 G}{\partial y^2}$$

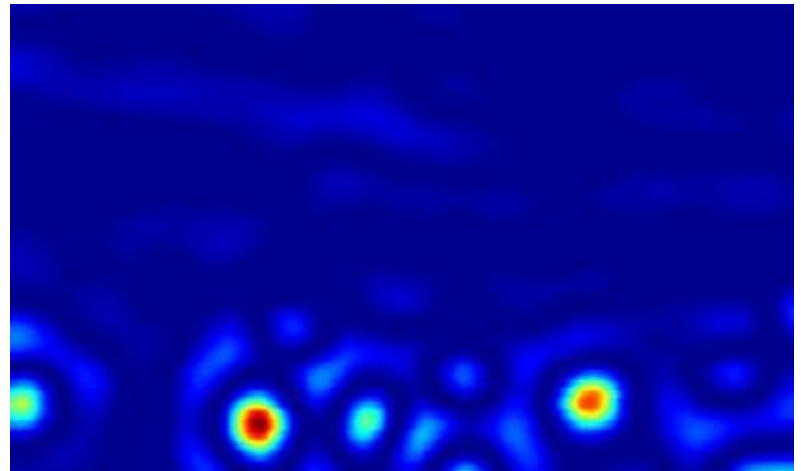
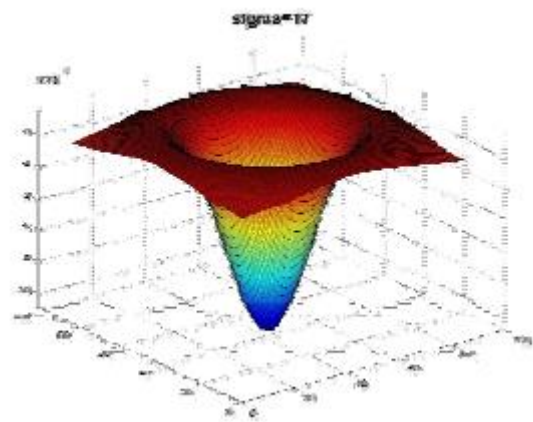




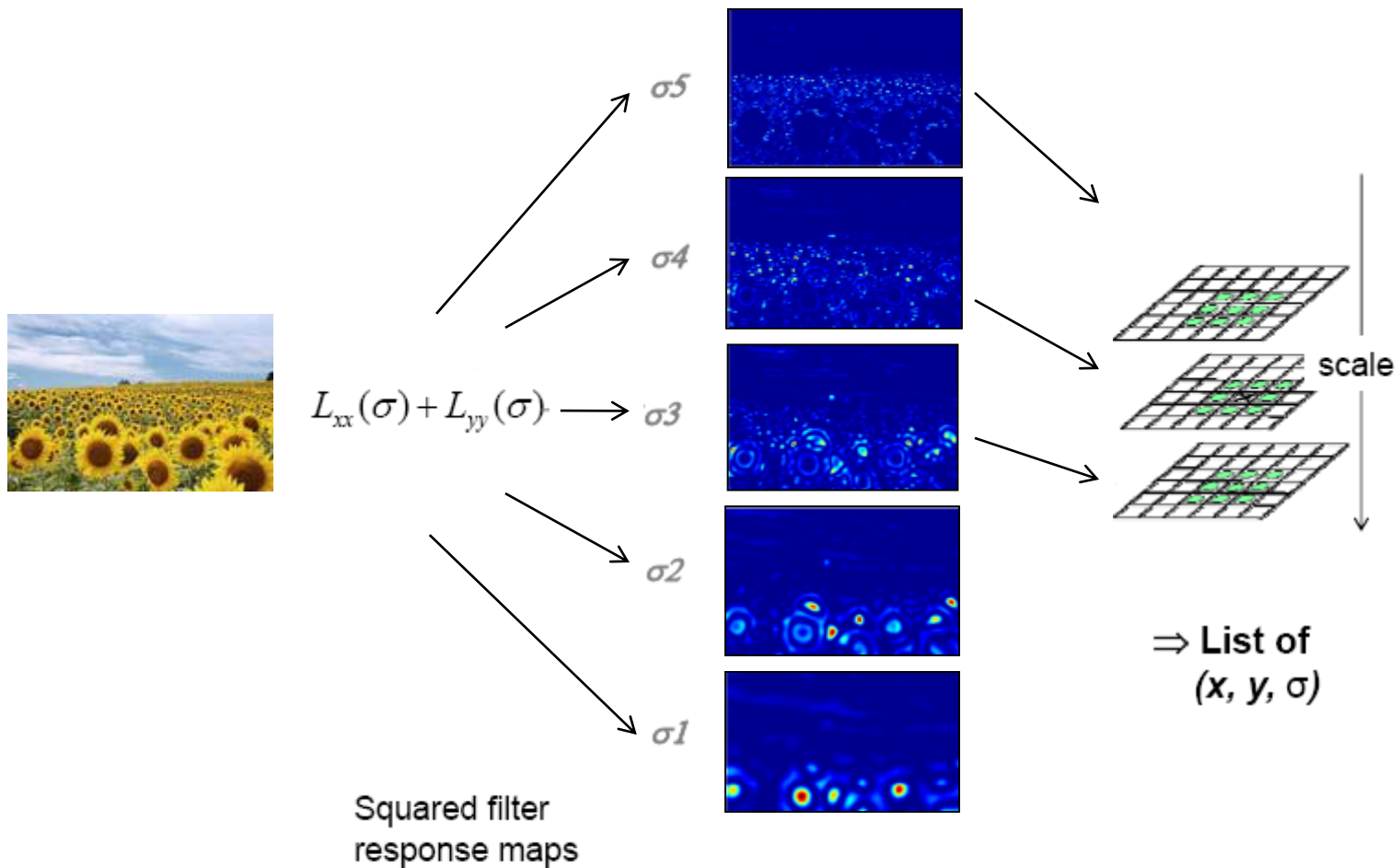




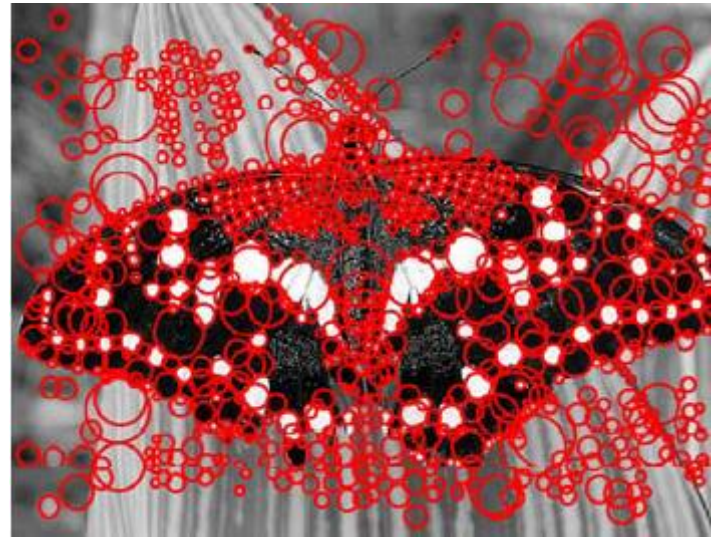




Scale-space blob detection



Scale-space blob detector: Example

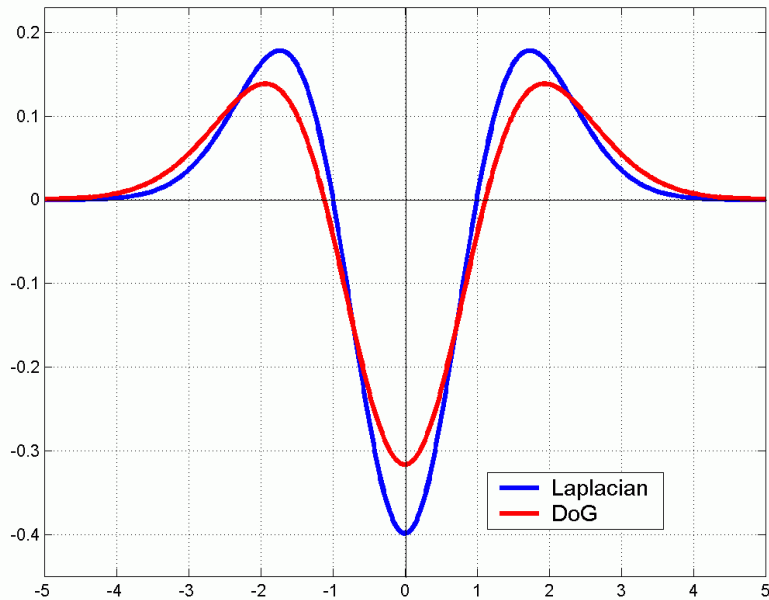


LoG V.S. DoG

$$\nabla^2 G_\sigma(x, y) = \left(\frac{x^2 + y^2}{\sigma^4} - \frac{2}{\sigma^2} \right) G_\sigma(x, y)$$

$$\underbrace{G(x, y, k\sigma) - G(x, y, \sigma)}_{\text{DoG}} \approx (k - 1)\sigma \frac{\partial G}{\partial \sigma} = (k - 1)\sigma^2 \underbrace{\nabla^2 G}_{\text{LoG}}$$

Lowe's Scale-space Interest Points: Difference of Gaussians



- Gaussian is an ad hoc solution of heat diffusion equation

$$\frac{\partial G}{\partial \sigma} = \sigma \nabla^2 G.$$

- Hence

$$G(x, y, k\sigma) - G(x, y, \sigma) \approx (k - 1)\sigma^2 \nabla^2 G.$$

- k is not necessarily very small in practice

Technical detail

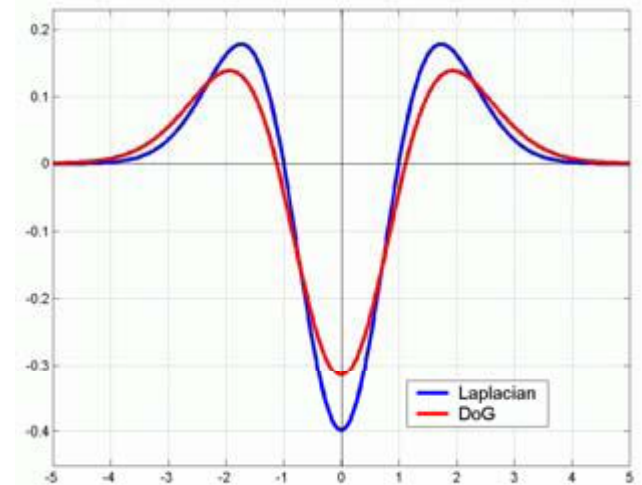
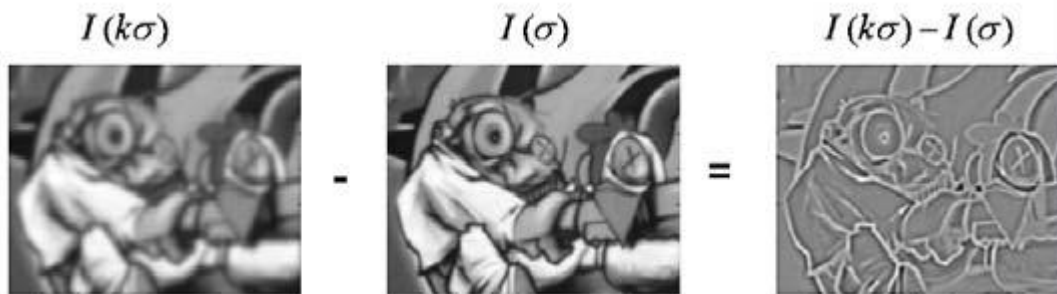
- We can approximate the Laplacian with a difference of Gaussians; more efficient to implement.

$$L = \sigma^2 (G_{xx}(x, y, \sigma) + G_{yy}(x, y, \sigma))$$

(Laplacian)

$$DoG = G(x, y, k\sigma) - G(x, y, \sigma)$$

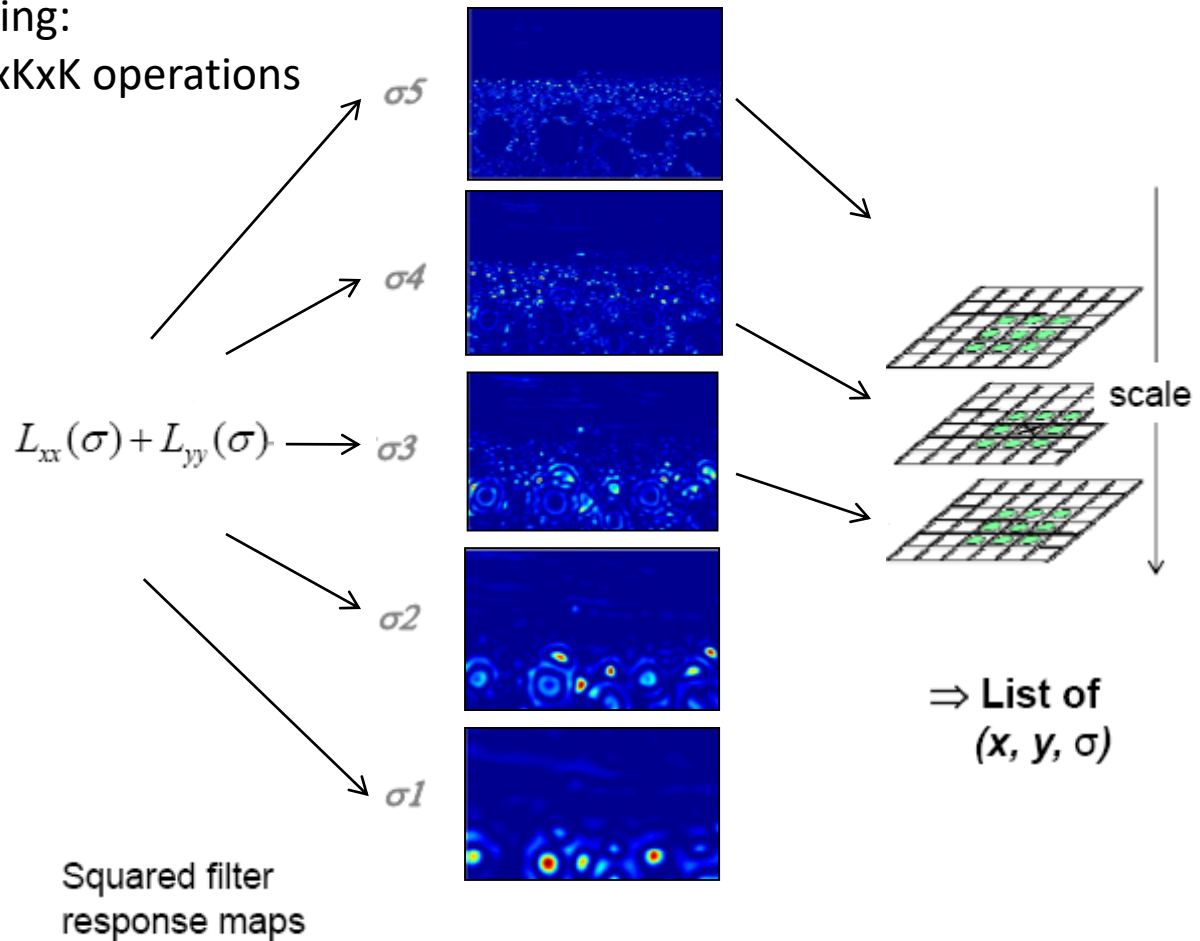
(Difference of Gaussians)



How many scales?

convolution computing:

At each scale: $M \times N \times K \times K$ operations



DoG Image Pyramid

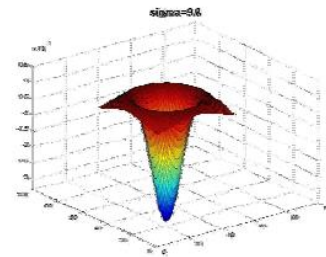
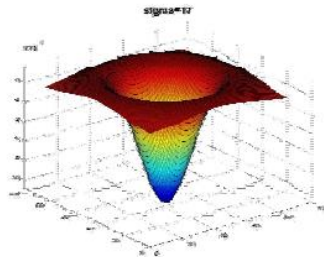
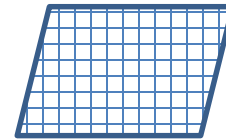
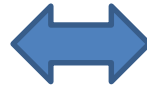
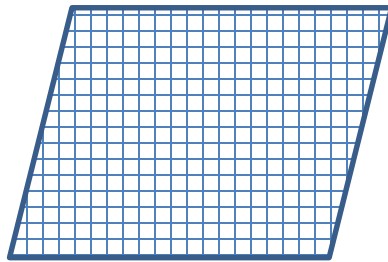
$$\sigma_0, k\sigma_0, k^2\sigma_0, k^3\sigma_0, k^4\sigma_0, k^5\sigma_0, k^6\sigma_0, \dots$$

$$\sigma_0 \rightarrow 2\sigma_0$$

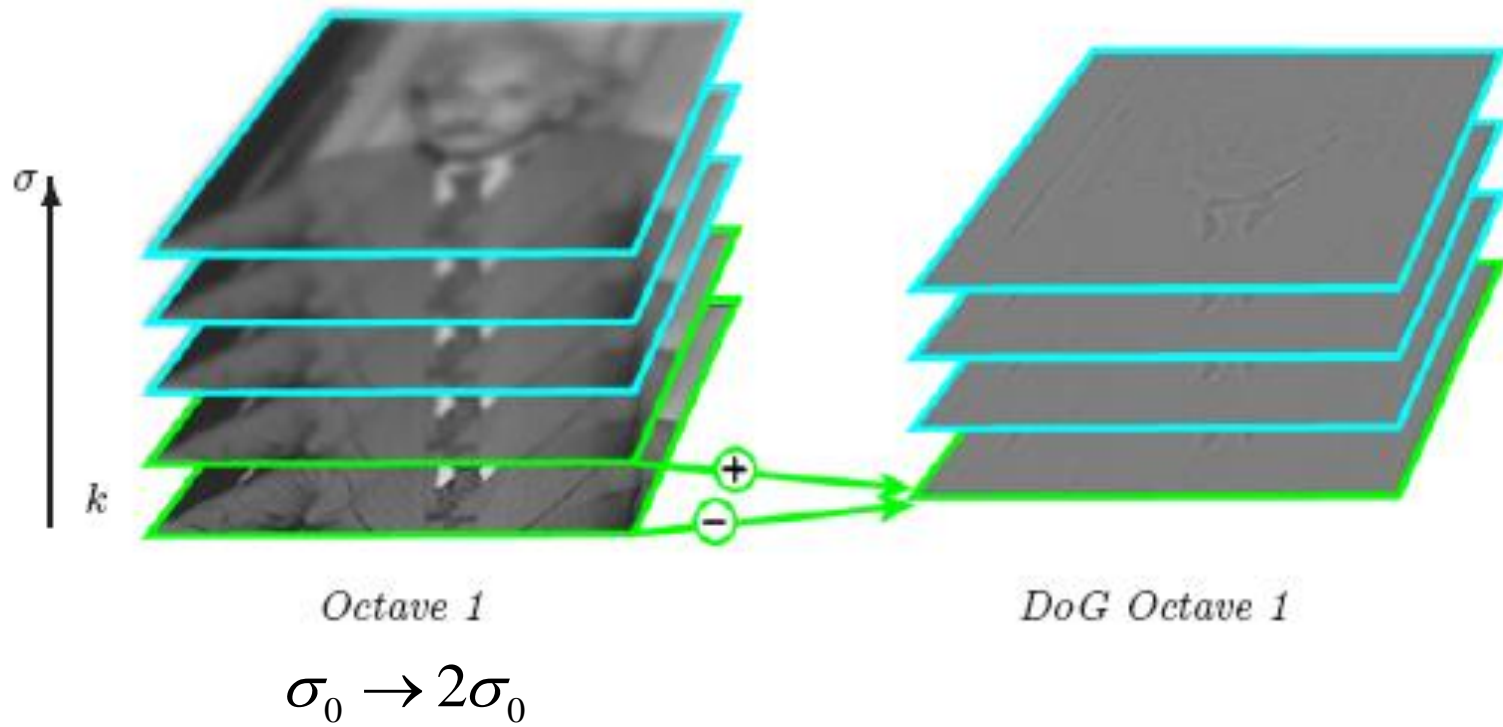
image $M \times N$, filter $2K \times 2K$

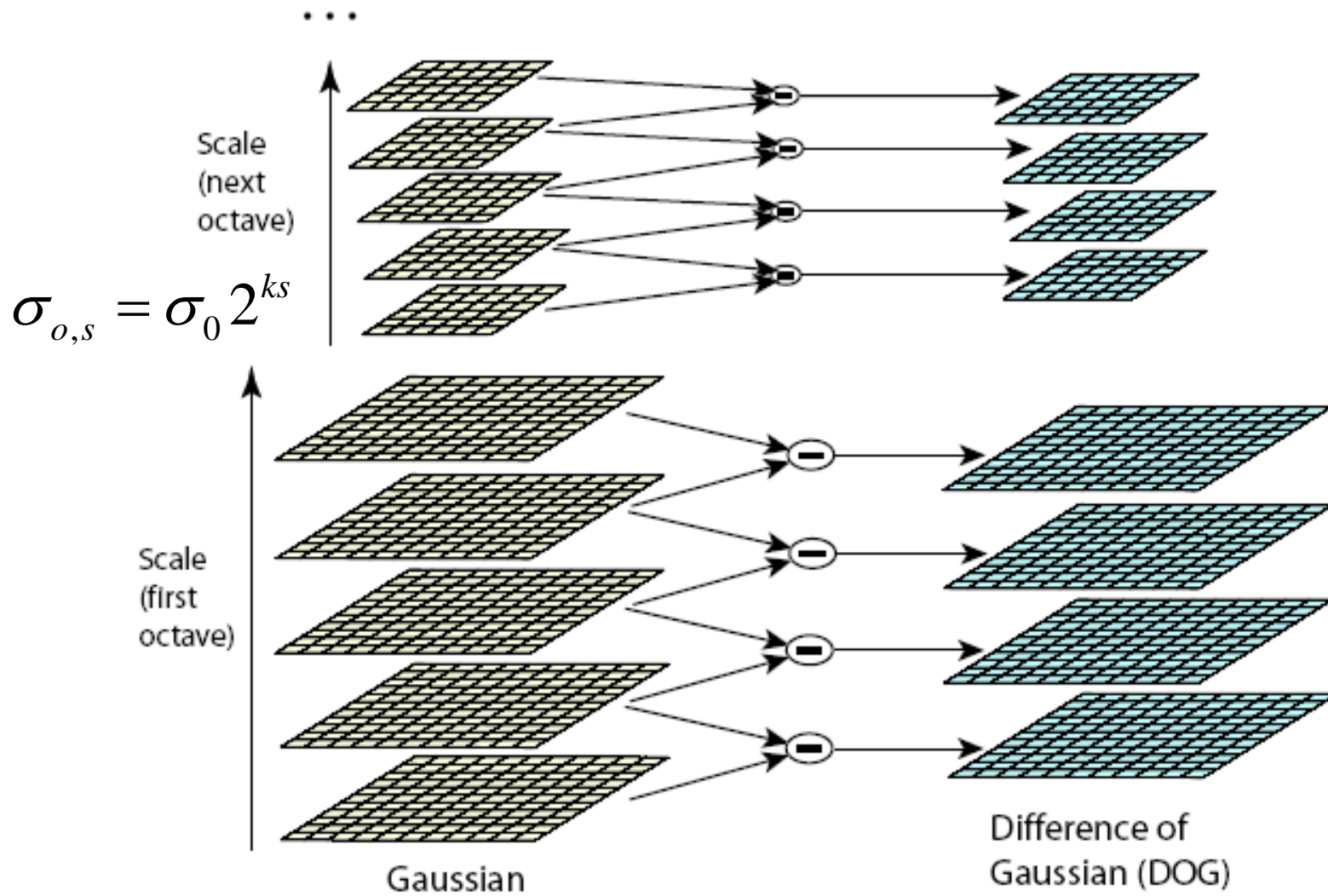


image $M/2 \times N/2$, filter, $K \times K$



DoG Image Pyramid



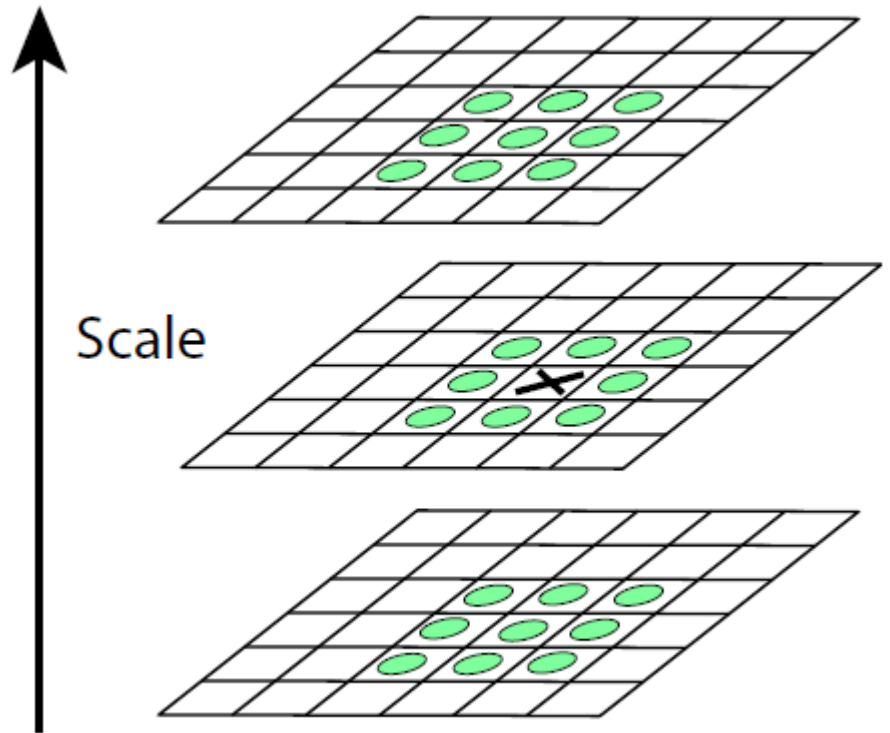


$$\sigma_{o,s} = \sigma_0 2^{ks}$$

$$\sigma_{o,s} = \sigma_0 2^{o+ks}$$

Local Extrema Detection

- Maxima and minima
- Compare x with its 26 neighbors at 3 scales



Frequency of sampling in scale

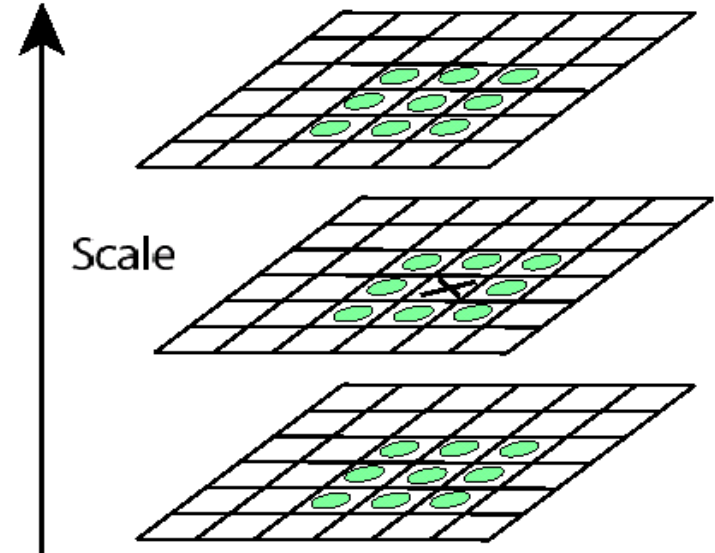
- s : intervals in each octave of scale space ($\sigma_0 \rightarrow 2\sigma_0$)
 - $k=2^{\{1/s\}}$
$$\sigma_{o,s} = \sigma_0 2^o k^s$$
- In order to cover a complete octave for extrema detection
 - $S = s+3$ Gaussian images are produced for each octave
 - $s: \{-1, S+1\}$
 - $s+2$ DoG images
 - s scales for extrema detection

SIFT Key point localization

- Detect maxima and minima of difference-of-Gaussian in scale space
- Fit a quadratic to surrounding values for sub-pixel and sub-scale interpolation (Brown & Lowe, 2002)
- Taylor expansion around point:
- Offset of extremum (use finite differences for derivatives):

$$D(\mathbf{x}) = D + \frac{\partial D^T}{\partial \mathbf{x}} \mathbf{x} + \frac{1}{2} \mathbf{x}^T \frac{\partial^2 D}{\partial \mathbf{x}^2} \mathbf{x}$$

$$\hat{\mathbf{x}} = -\frac{\partial^2 D^{-1}}{\partial \mathbf{x}^2} \frac{\partial D}{\partial \mathbf{x}}$$

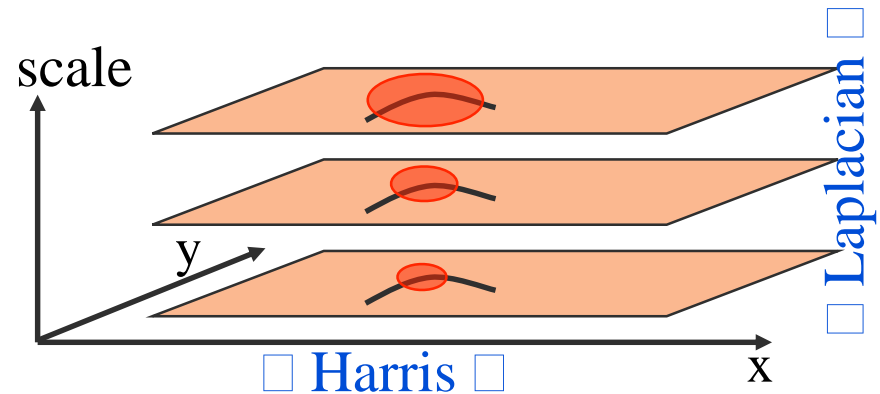


Scale Invariant Detectors

- **Harris-Laplacian**¹

Find local maximum of:

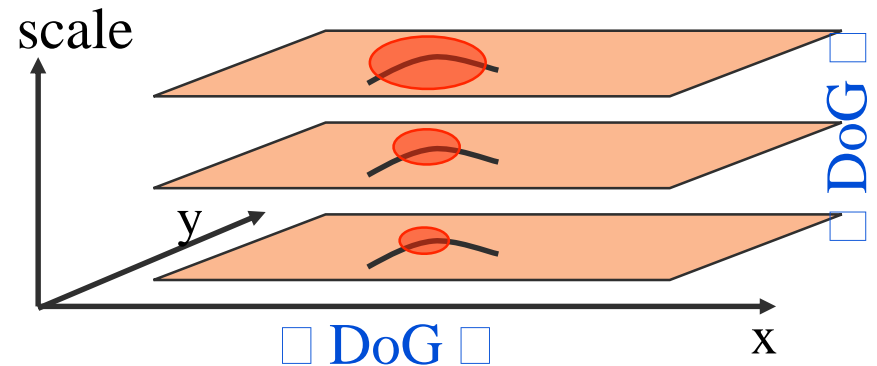
- Harris corner detector in space (image coordinates)
- Laplacian in scale



- **SIFT (Lowe)**²

Find local maximum of:

- Difference of Gaussians in space and scale

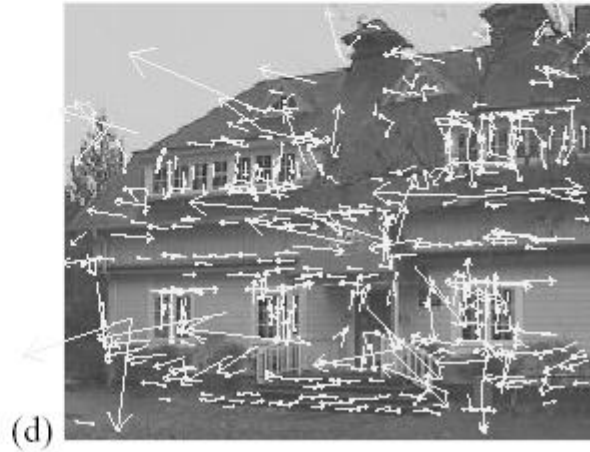


¹ K.Mikolajczyk, C.Schmid. “Indexing Based on Scale Invariant Interest Points”. ICCV 2001

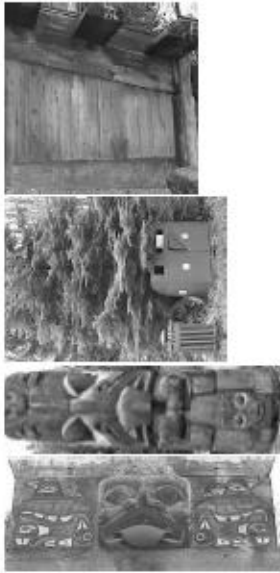
² D.Lowe. “Distinctive Image Features from Scale-Invariant Keypoints”. Accepted to IJCV 2004

Example of keypoint detection

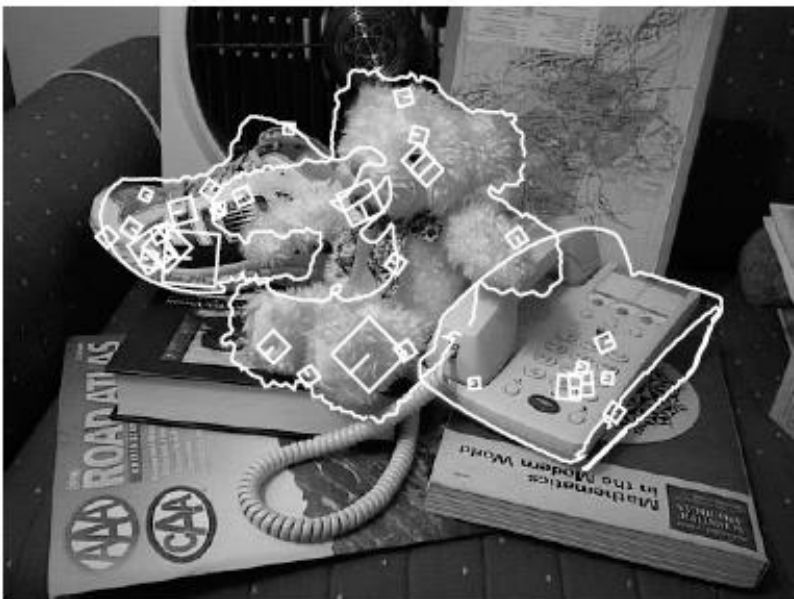
Threshold on value at DOG peak and on ratio of principle curvatures (Harris approach)



- (a) 233x189 image
- (b) 832 DOG extrema
- (c) 729 left after peak value threshold
- (d) 536 left after testing ratio of principle curvatures







SIFT vector formation

- Thresholded image gradients are sampled over 16x16 array of locations in scale space
- Create array of orientation histograms
- 8 orientations x 4x4 histogram array = 128 dimensions

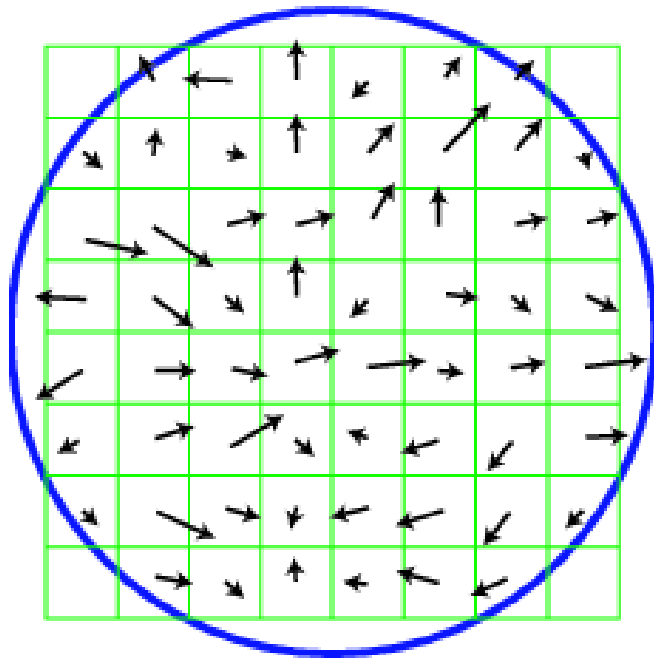
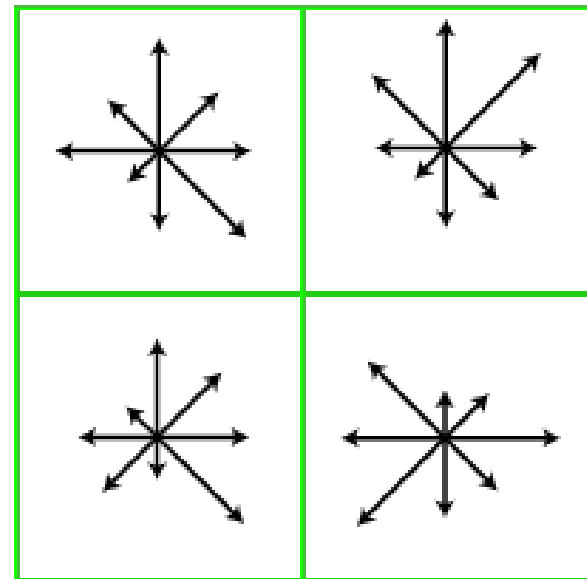


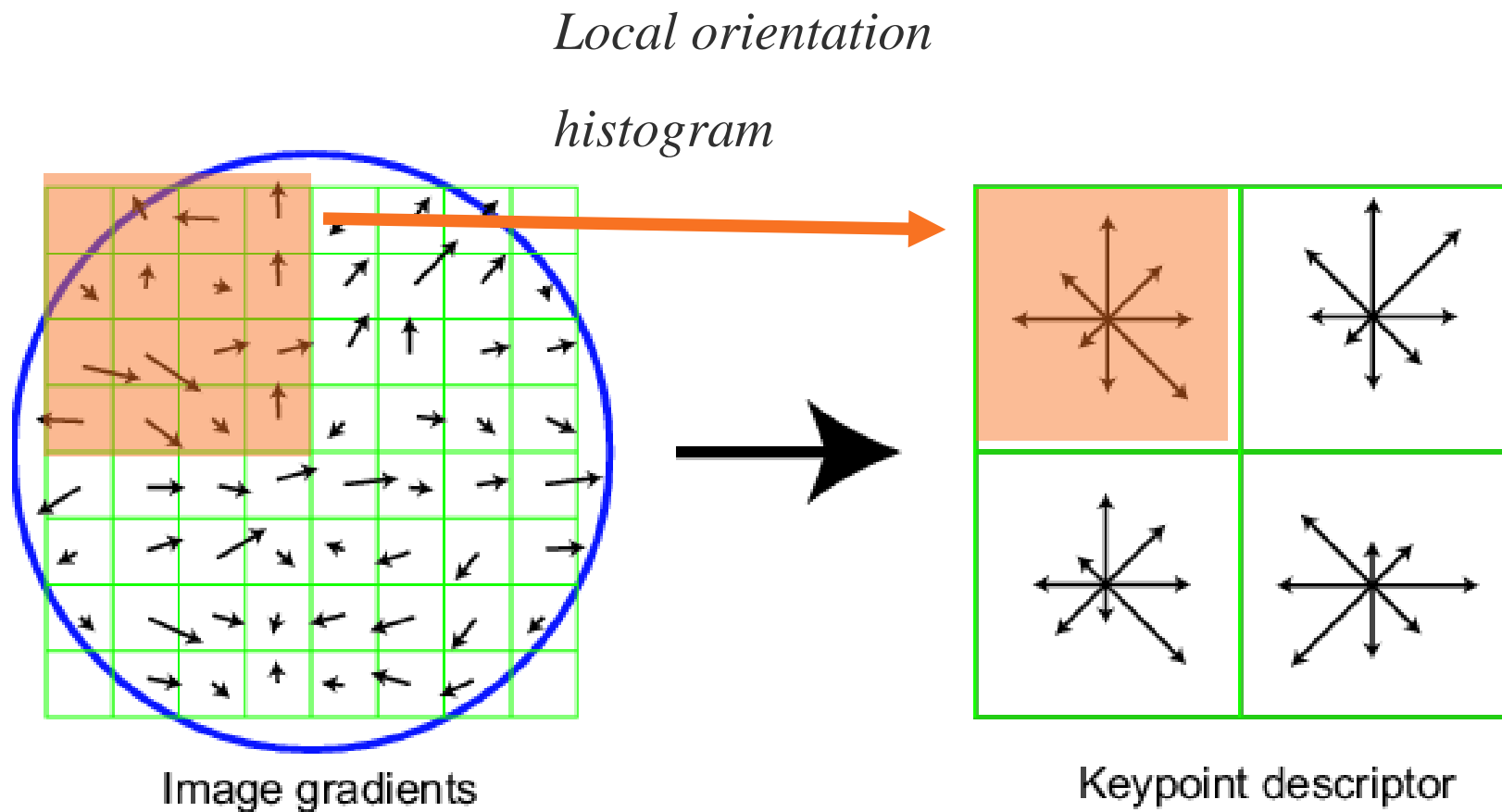
Image gradients



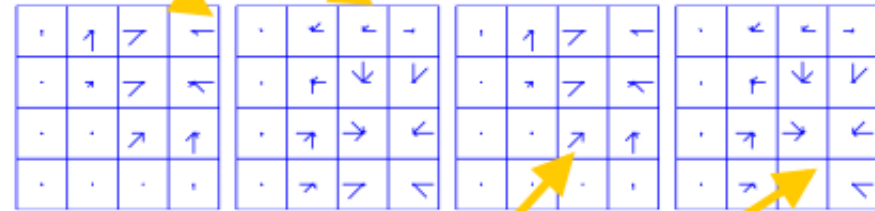
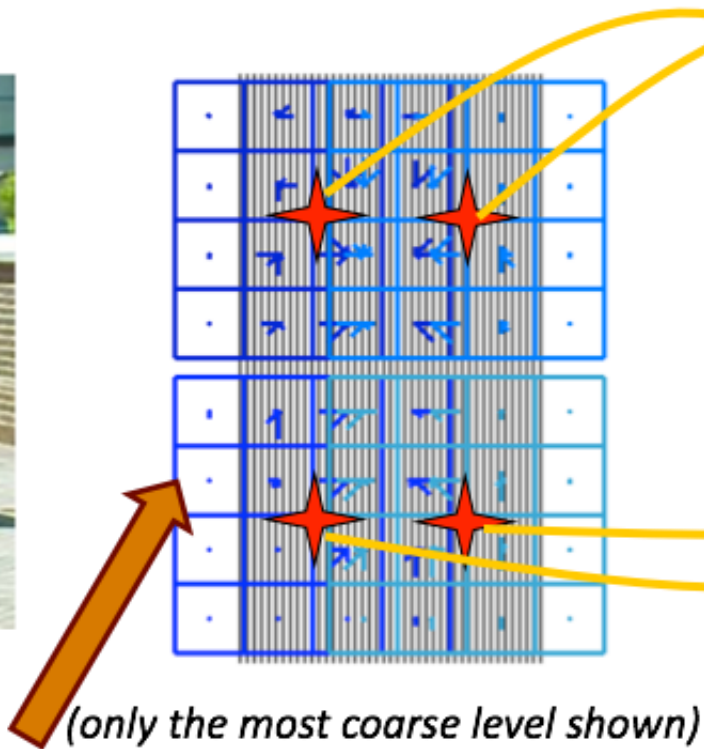
Keypoint descriptor

SIFT vector formation

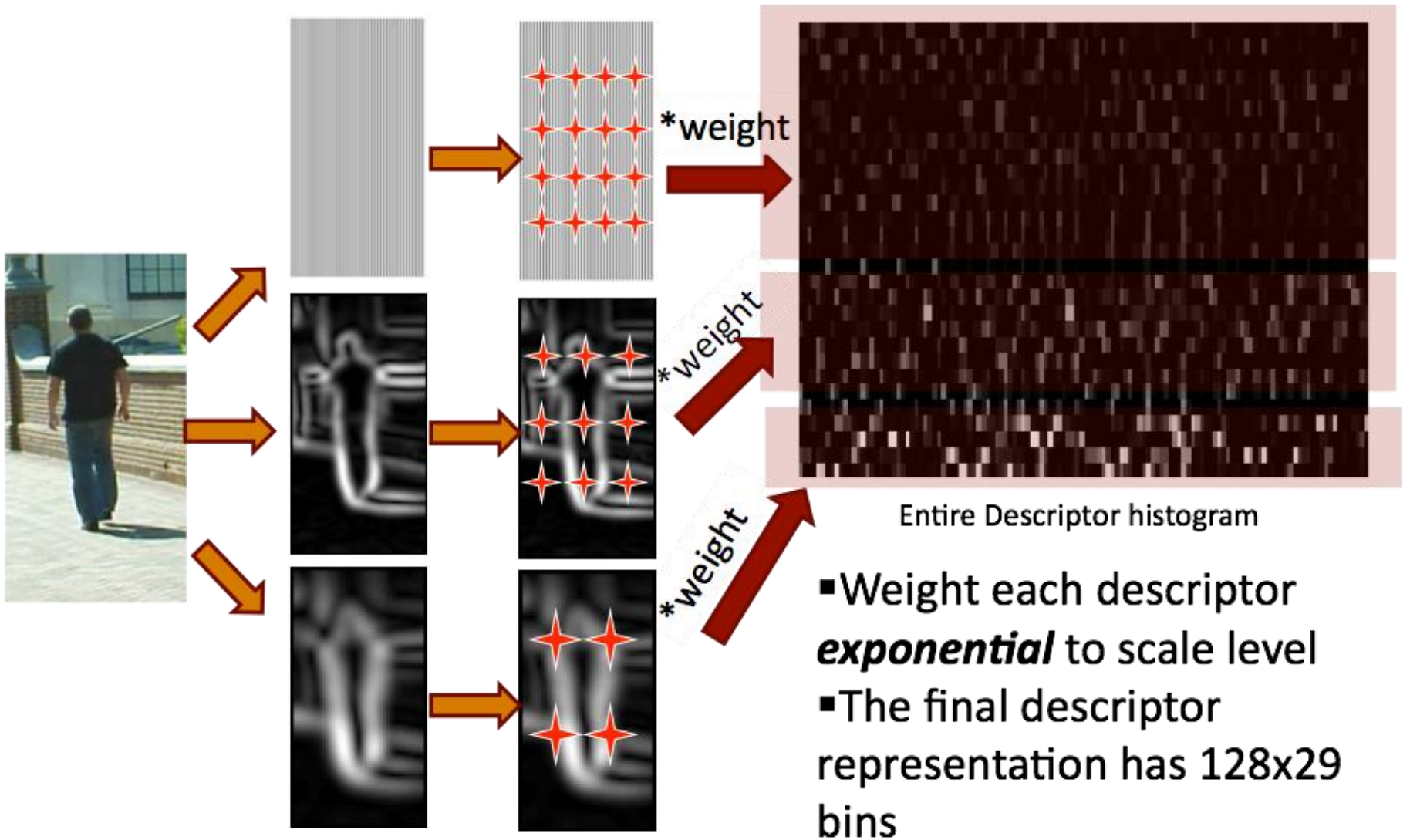
- Orientation is defined relative to the orientation of the detected Sift feature



Feature Extraction for Image Classifier



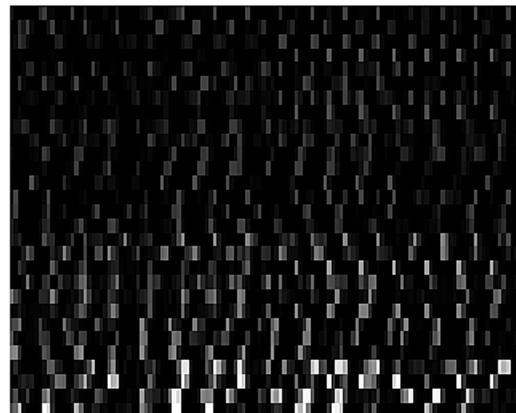
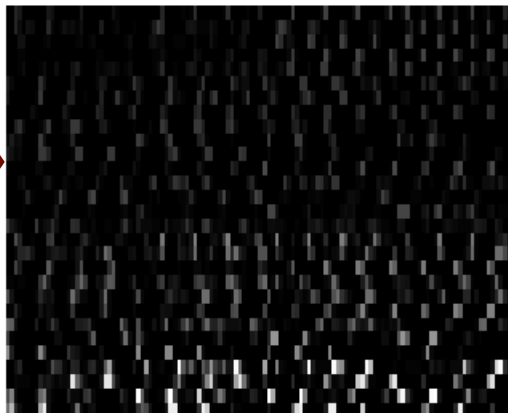
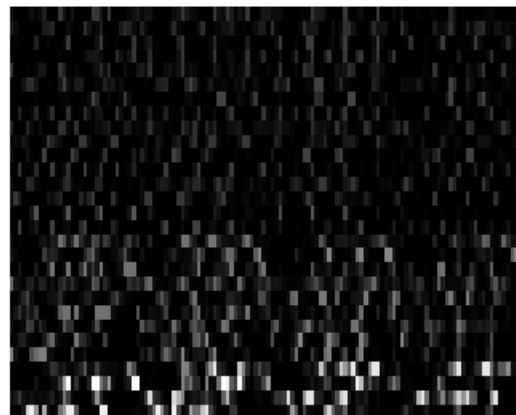
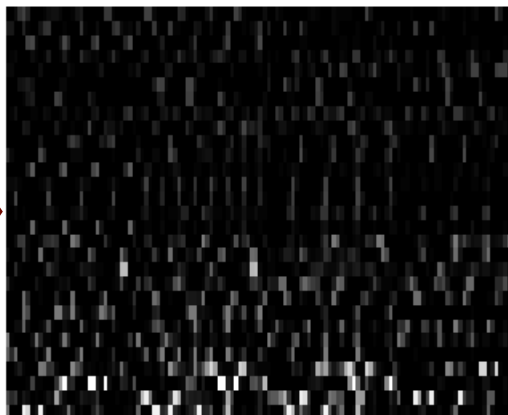
Descriptor histogram of Scale level 3



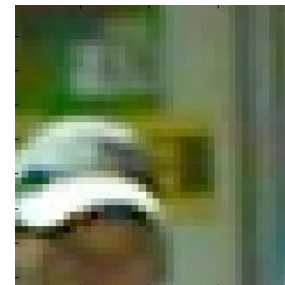
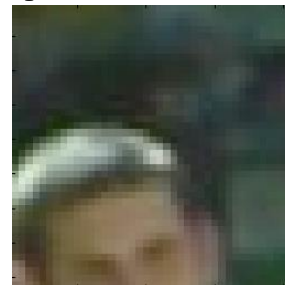
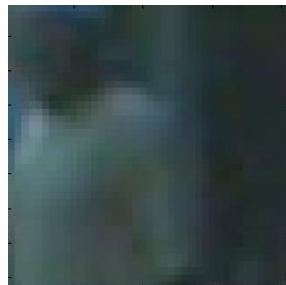
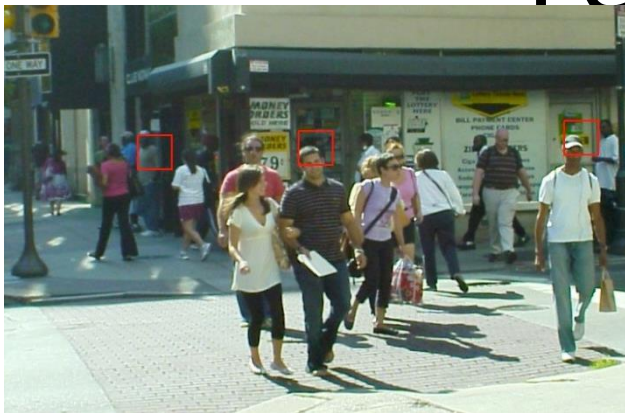
Feature Extraction for Image Classifier

Examples from training set

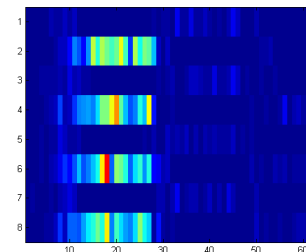
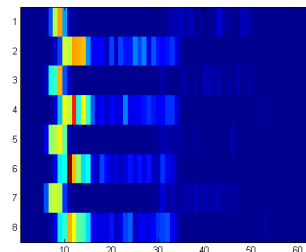
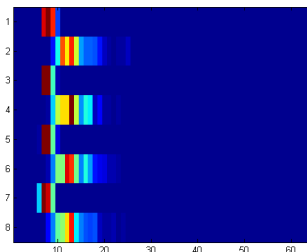
Examples from test set



Features Sample



1. Color Histogram



2. HOG feature

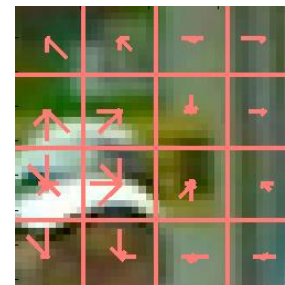
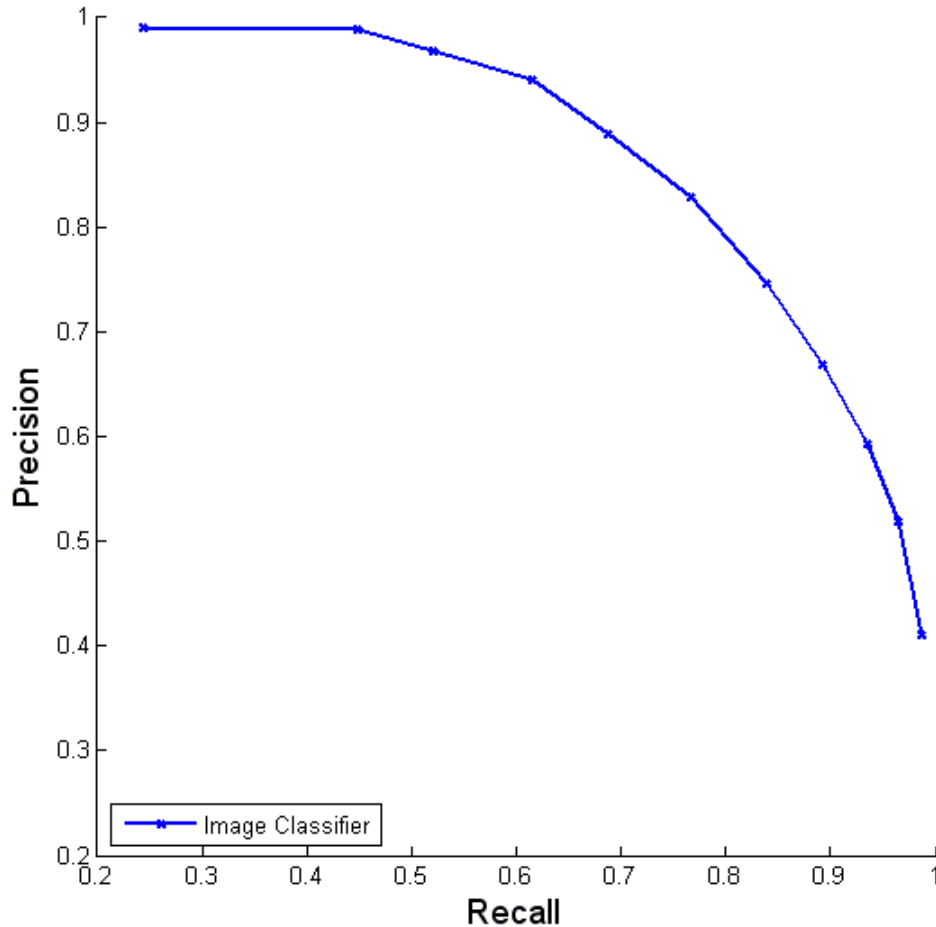


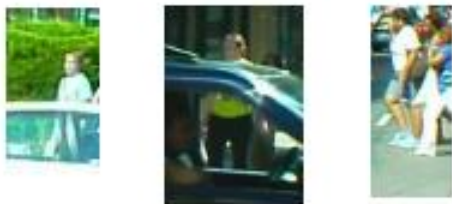
Image Classifier Result



- Precision-Recall curve of our Image classifier
- Trained/Tested on image set generated from stereo detection
- RBF kernel SVM used

Typical Missed Detections

•Occlusions



•Incomplete stereo detection



•Lack of training data



Typical False Positives

•Human-like Shapes & Clutters



Deep...

