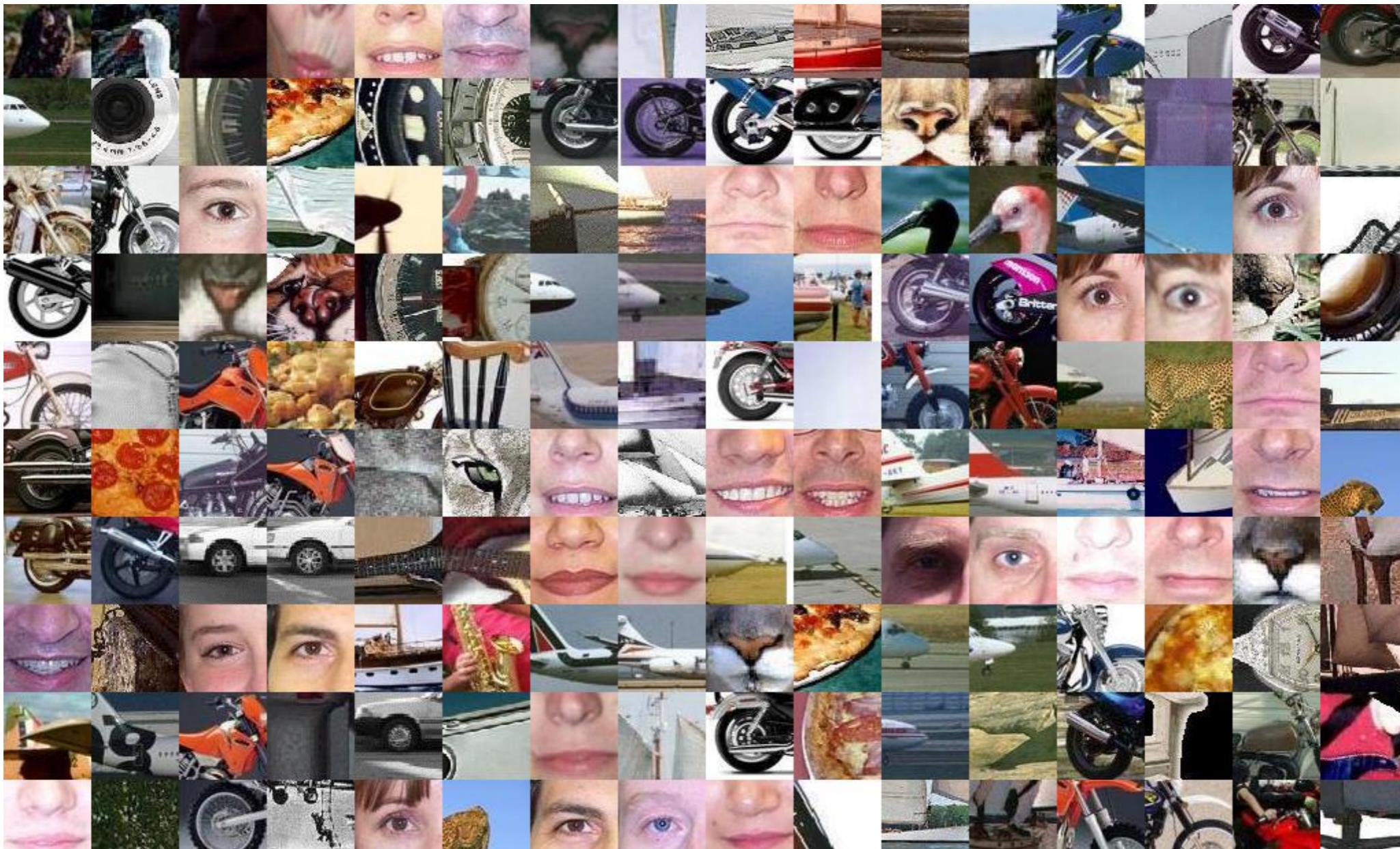


Designing descriptors



Course announcements

- Homework 0 and homework 1 are available on course website.
 - Homework 1 is due on February 7th.
 - Any questions about the homeworks?
 - How many of you have looked at/started/finished homework 0?
 - How many of you have looked at/started/finished homework 1?
- There was a small correction to homework 1.
 - Posted on Piazza.
 - Make sure to download the latest version.

Overview of today's lecture

Leftover from lecture 5:

- Multi-scale detection.

New in lecture 6:

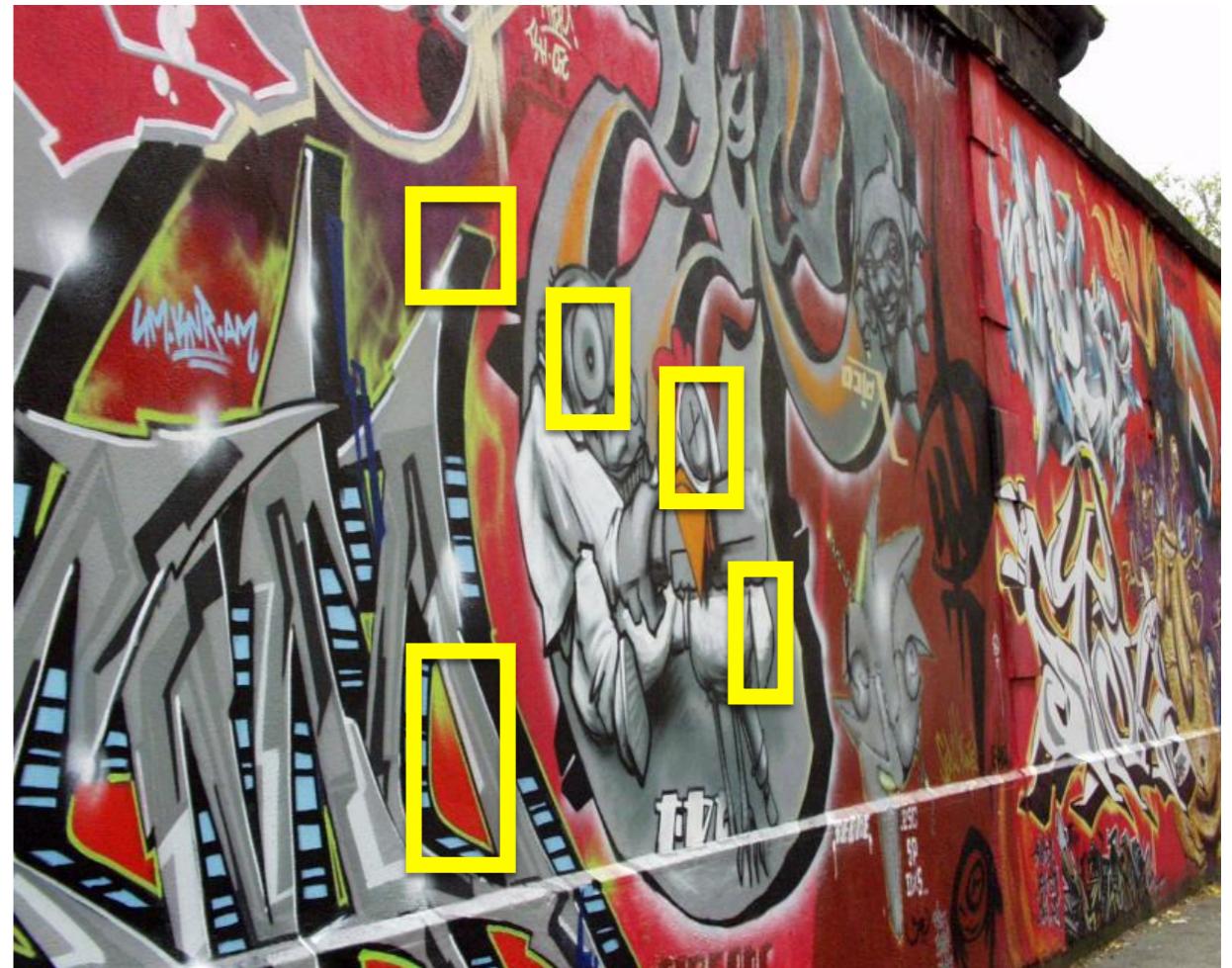
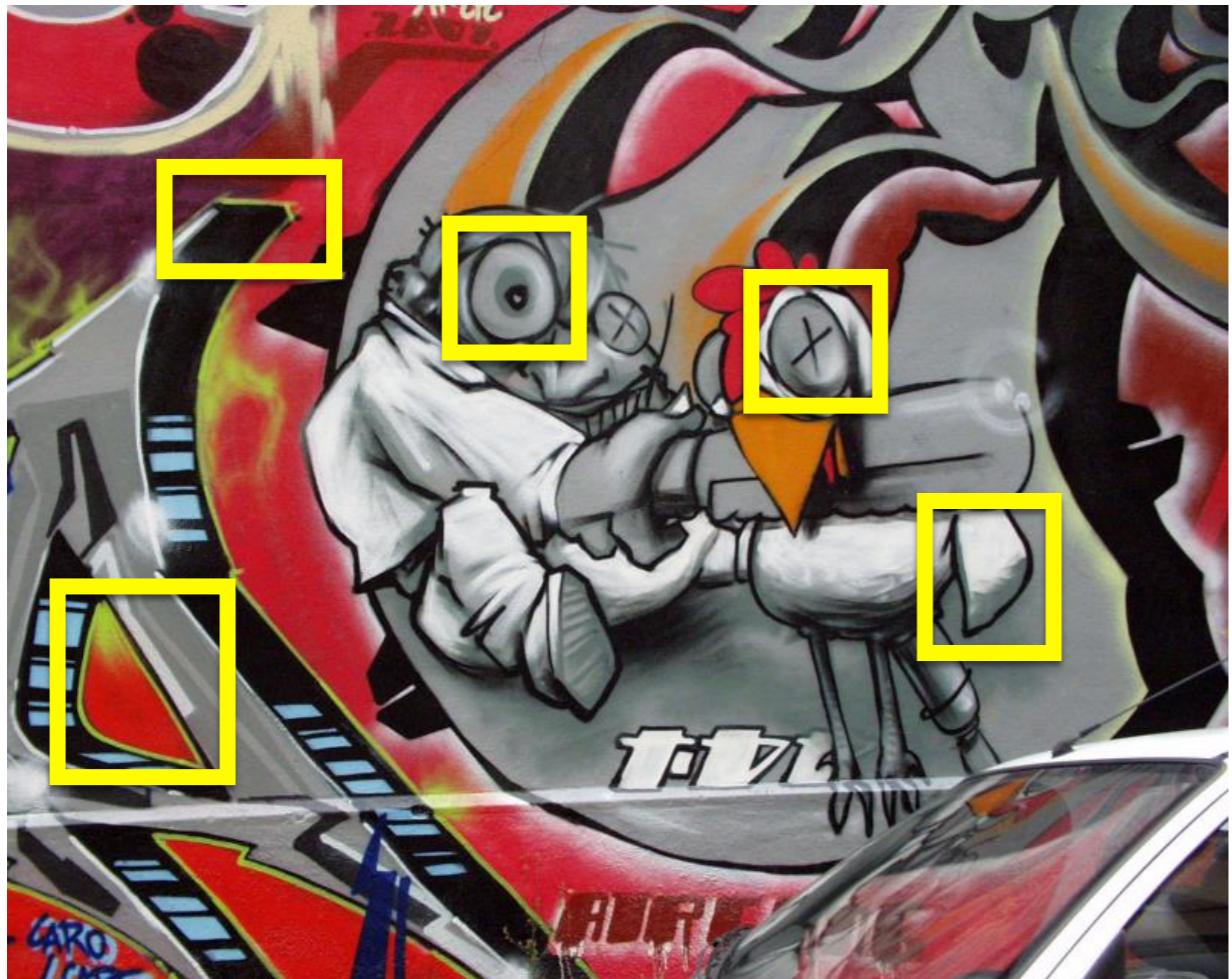
- Why do we need feature descriptors?
- Designing feature descriptors.
- MOPS descriptor.
- GIST descriptor.
- Histogram of Textons descriptor.
- HOG descriptor.
- SURF descriptor.
- SIFT.

Slide credits

Most of these slides were adapted from:

- Kris Kitani (16-385, Spring 2017).

Why do we need feature
descriptors?



*If we know where the good features are,
how do we match them?*

Object instance recognition



Schmid and Mohr 1997



Sivic and Zisserman, 2003



Rothganger et al. 2003



Lowe 2002

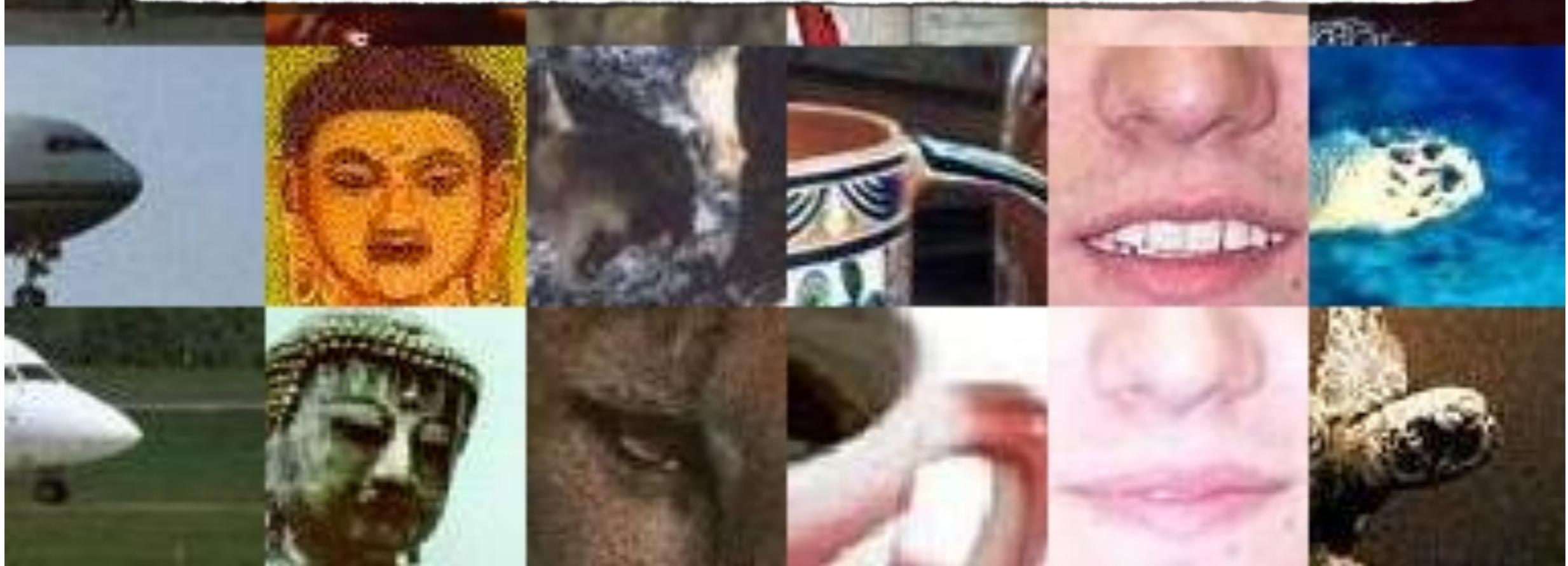
Image mosaicing





How do we describe an image patch?

Patches with similar content should have similar descriptors.

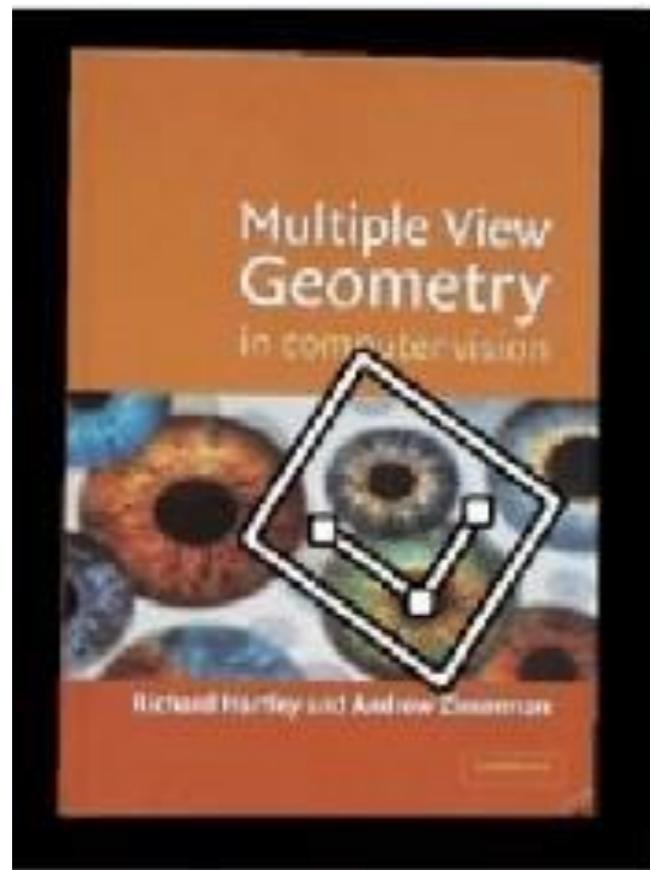


Designing feature descriptors

Photometric transformations



Geometric transformations



objects will appear at different scales,
translation and rotation



What is the best descriptor for an image feature?

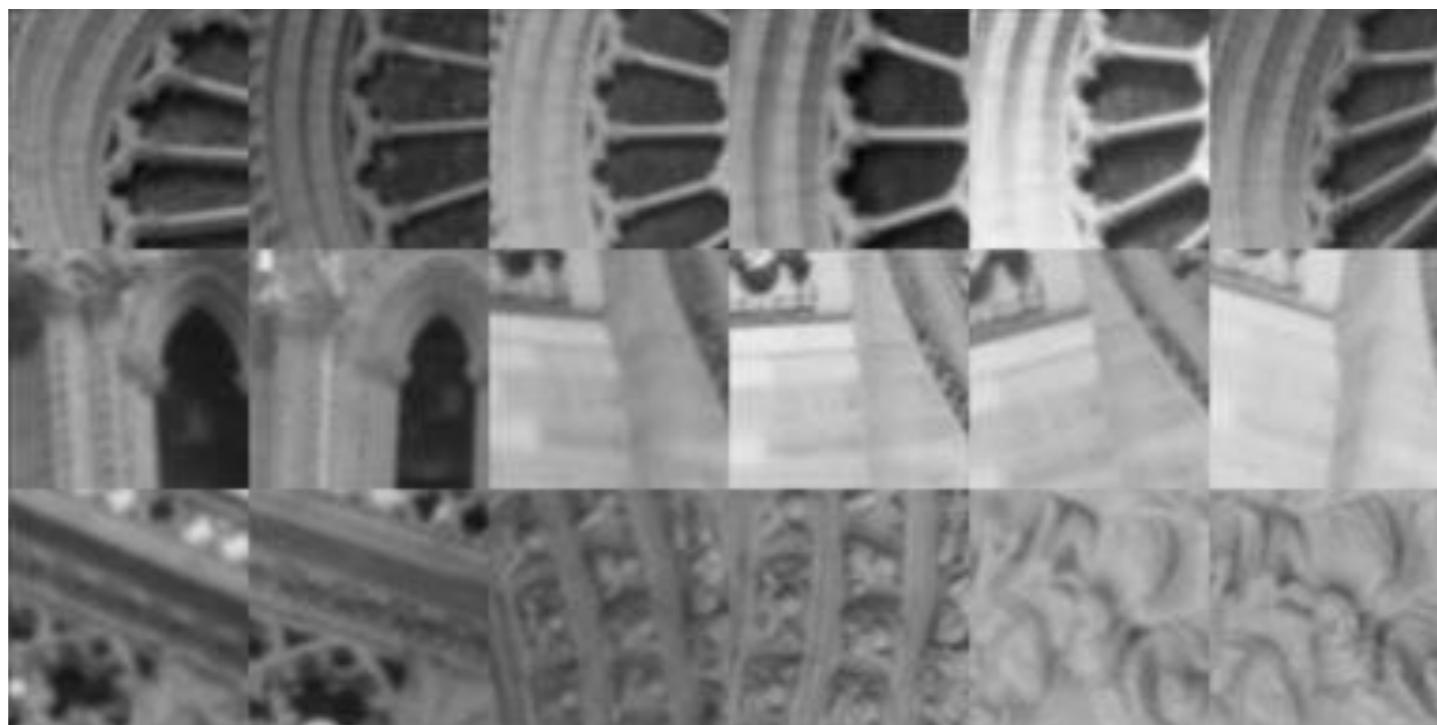


Image patch

Just use the pixel values of the patch



Perfectly fine if geometry and appearance is unchanged
(a.k.a. template matching)

Tiny Images



Just down-sample it!
Simple, fast, robust to small affine transforms.



Image patch

Just use the pixel values of the patch



Perfectly fine if geometry and appearance is unchanged
(a.k.a. template matching)

What are the problems?

Image patch

Just use the pixel values of the patch



Perfectly fine if geometry and appearance is unchanged
(a.k.a. template matching)

What are the problems?

How can you be less sensitive to absolute intensity values?

Image gradients

Use pixel differences

1	2	3
4	5	6
7	8	9



$$(\quad - \quad + \quad + \quad - \quad - \quad + \quad)$$

vector of x derivatives

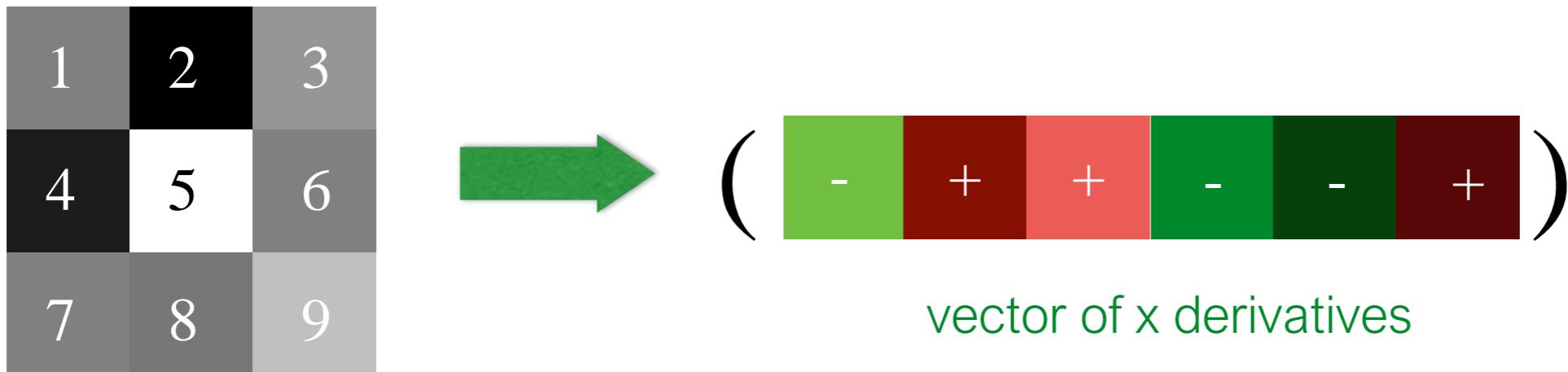
'binary descriptor'

Feature is invariant to absolute intensity values

What are the problems?

Image gradients

Use pixel differences



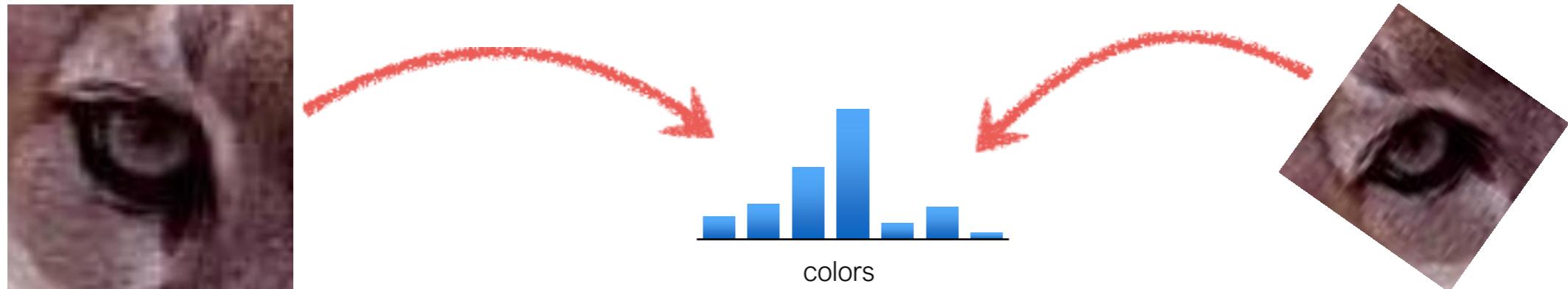
Feature is invariant to absolute intensity values

What are the problems?

How can you be less sensitive to deformations?

Color histogram

Count the colors in the image using a histogram

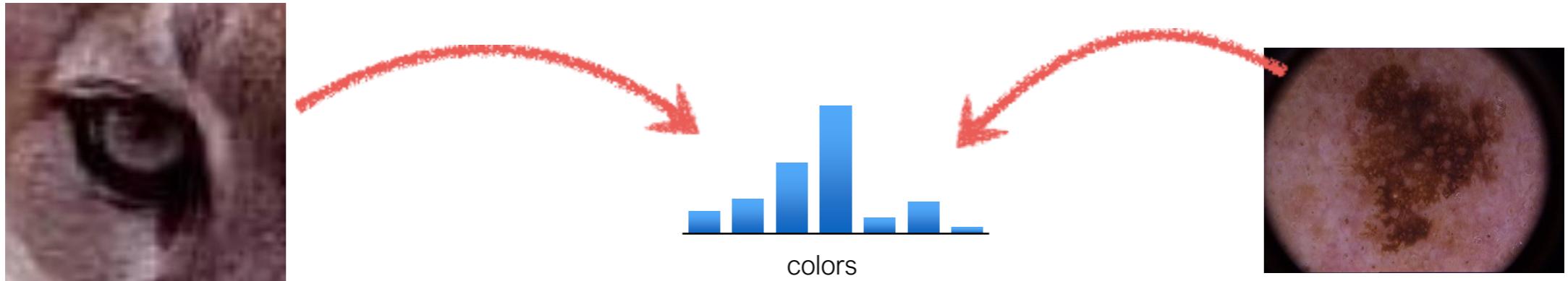


Invariant to changes in scale and rotation

What are the problems?

Color histogram

Count the colors in the image using a histogram

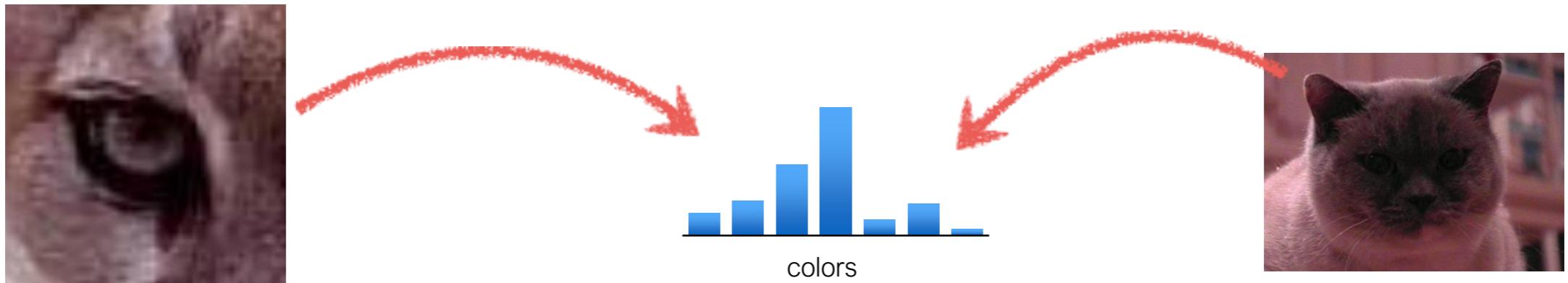


Invariant to changes in scale and rotation

What are the problems?

Color histogram

Count the colors in the image using a histogram



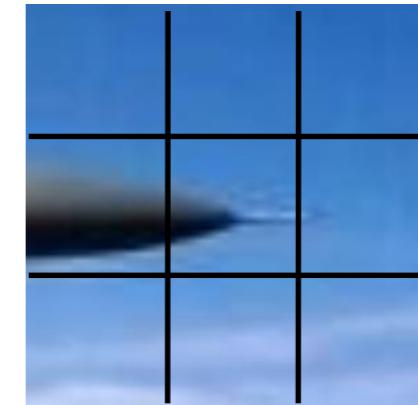
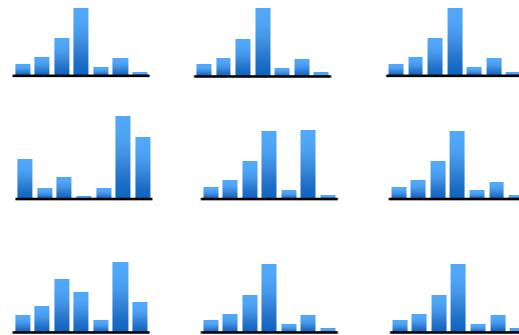
Invariant to changes in scale and rotation

What are the problems?

How can you be more sensitive to spatial layout?

Spatial histograms

Compute histograms over spatial ‘cells’

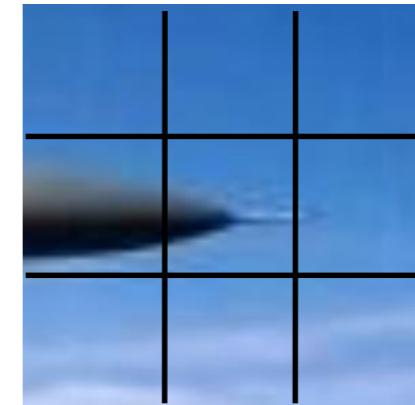
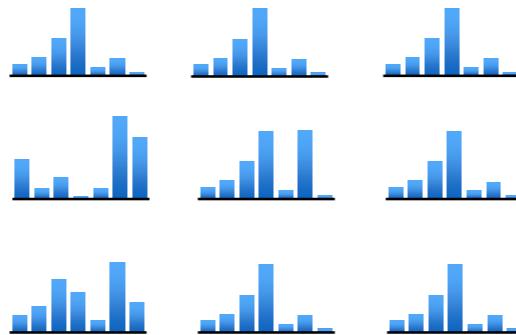


Retains rough spatial layout
Some invariance to deformations

What are the problems?

Spatial histograms

Compute histograms over spatial ‘cells’



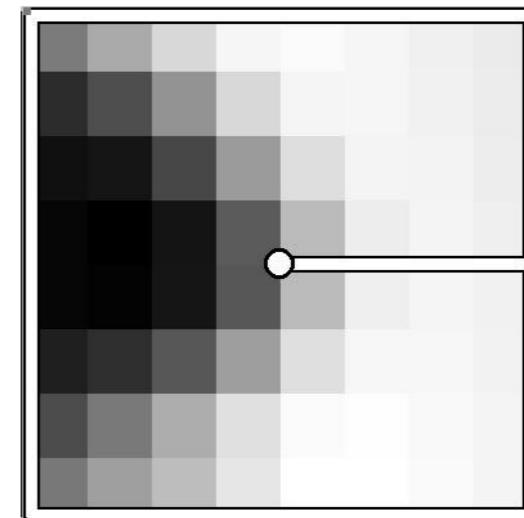
Retains rough spatial layout
Some invariance to deformations

What are the problems?

How can you be completely invariant to rotation?

Orientation normalization

Use the dominant image gradient direction to normalize the orientation of the patch



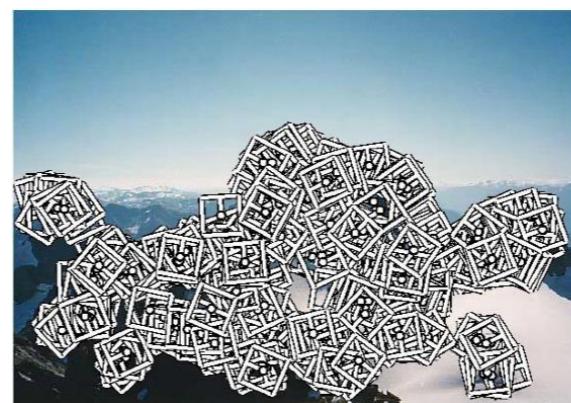
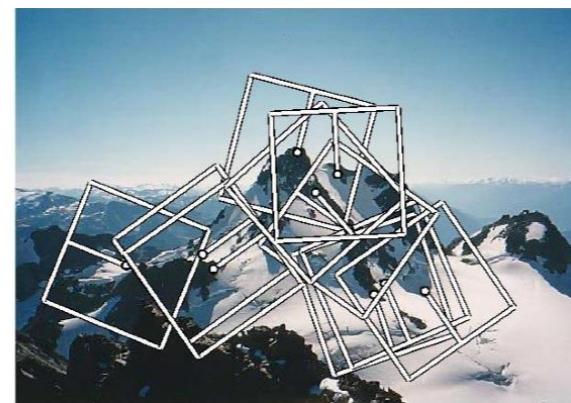
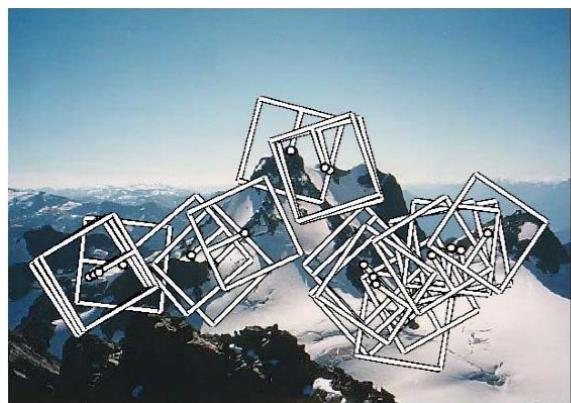
save the orientation angle θ along with (x, y, s)

What are the problems?

MOPS descriptor

Multi-Scale Oriented Patches (MOPS)

Multi-Image Matching using Multi-Scale Oriented Patches. M. Brown, R. Szeliski and S. Winder.
International Conference on Computer Vision and Pattern Recognition (CVPR2005). pages 510-517



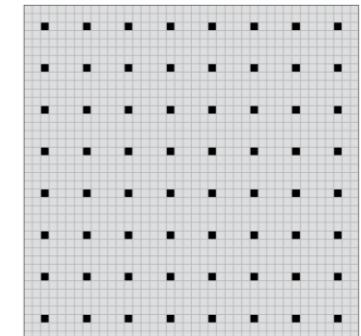
Multi-Scale Oriented Patches (MOPS)

Multi-Image Matching using Multi-Scale Oriented Patches. M. Brown, R. Szeliski and S. Winder.
International Conference on Computer Vision and Pattern Recognition (CVPR2005). pages 510-517

Given a feature (x, y, s, θ)

Get 40×40 image patch, subsample
every 5th pixel

(*what's the purpose of this step?*)



Subtract the mean, divide by standard
deviation

(*what's the purpose of this step?*)

Haar Wavelet Transform

(*what's the purpose of this step?*)

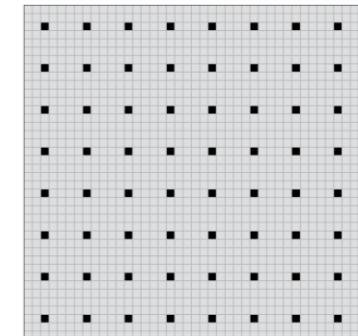
Multi-Scale Oriented Patches (MOPS)

Multi-Image Matching using Multi-Scale Oriented Patches. M. Brown, R. Szeliski and S. Winder.
International Conference on Computer Vision and Pattern Recognition (CVPR2005). pages 510-517

Given a feature (x, y, s, θ)

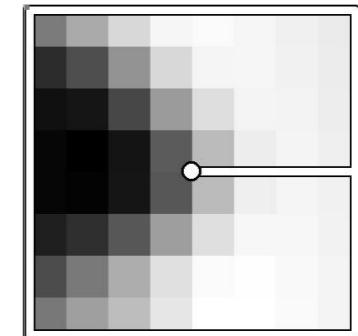
Get 40×40 image patch, subsample
every 5th pixel

(low frequency filtering, absorbs localization errors)



Subtract the mean, divide by standard
deviation

(*what's the purpose of this step?*)



Haar Wavelet Transform
(*what's the purpose of this step?*)

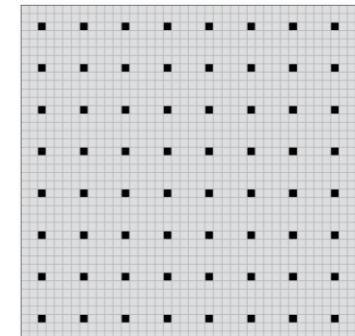
Multi-Scale Oriented Patches (MOPS)

Multi-Image Matching using Multi-Scale Oriented Patches. M. Brown, R. Szeliski and S. Winder.
International Conference on Computer Vision and Pattern Recognition (CVPR2005). pages 510-517

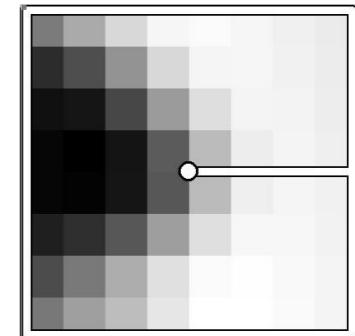
Given a feature (x, y, s, θ)

Get 40×40 image patch, subsample
every 5th pixel

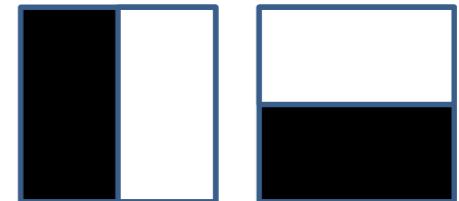
(low frequency filtering, absorbs localization errors)



Subtract the mean, divide by standard
deviation
(removes bias and gain)



Haar Wavelet Transform
(what's the purpose of this step?)



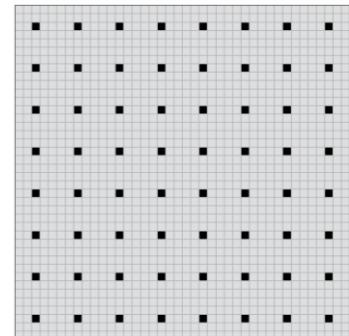
Multi-Scale Oriented Patches (MOPS)

Multi-Image Matching using Multi-Scale Oriented Patches. M. Brown, R. Szeliski and S. Winder.
International Conference on Computer Vision and Pattern Recognition (CVPR2005). pages 510-517

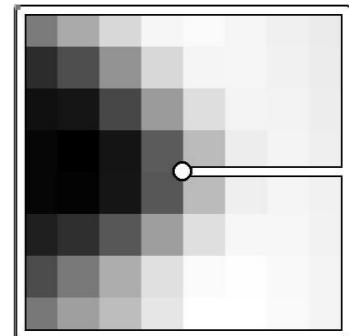
Given a feature (x, y, s, θ)

Get 40×40 image patch, subsample
every 5th pixel

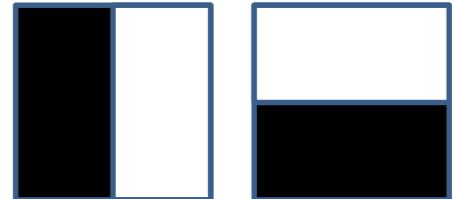
(low frequency filtering, absorbs localization errors)



Subtract the mean, divide by standard
deviation
(removes bias and gain)



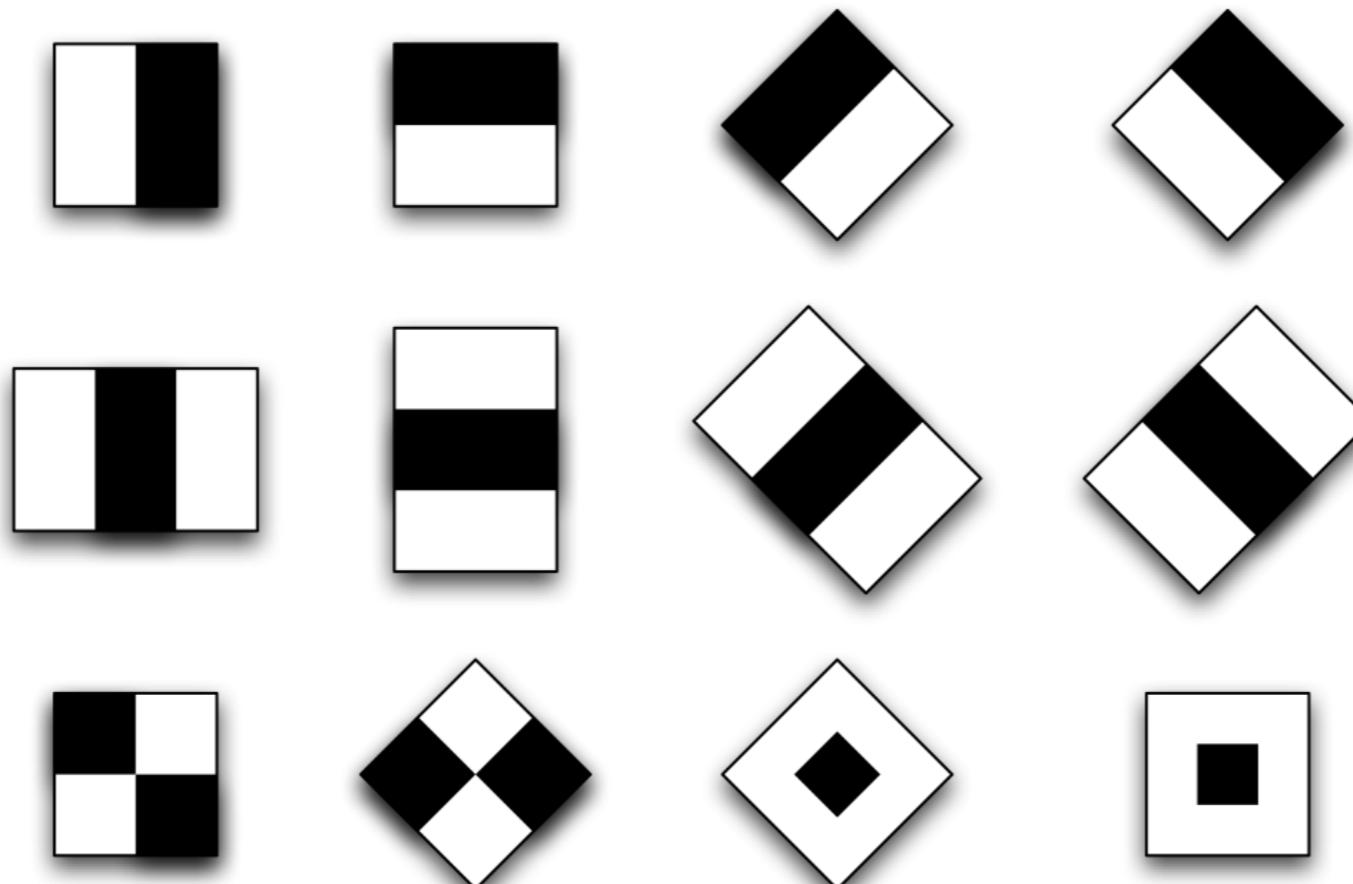
Haar Wavelet Transform
(low frequency projection)



Haar Wavelets

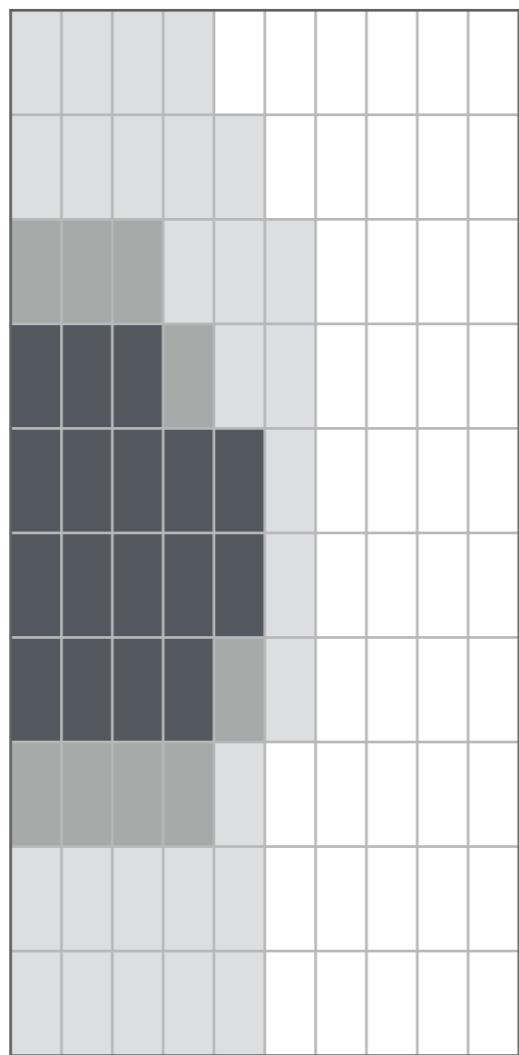
(actually, Haar-like features)

Use responses of a bank of filters as a descriptor

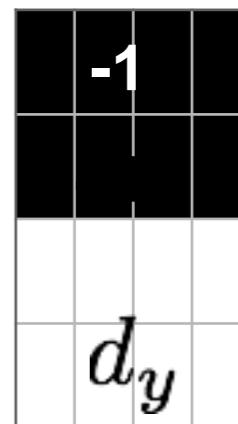
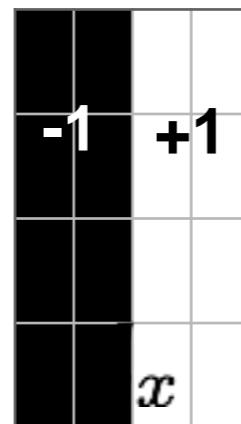


Haar wavelet responses can be computed with filtering

image patch



Haar wavelets filters



Haar wavelet responses can be
computed **efficiently** (in constant time)
with integral images

Integral Image

$I(x, y)$	$A(x, y)$																		
<table border="1"><tr><td>1</td><td>5</td><td>2</td></tr><tr><td>2</td><td>4</td><td>1</td></tr><tr><td>2</td><td>1</td><td>1</td></tr></table> <p>original image</p>	1	5	2	2	4	1	2	1	1	<table border="1"><tr><td>1</td><td>6</td><td>8</td></tr><tr><td>3</td><td>12</td><td>15</td></tr><tr><td>5</td><td>15</td><td>19</td></tr></table> <p>integral image</p>	1	6	8	3	12	15	5	15	19
1	5	2																	
2	4	1																	
2	1	1																	
1	6	8																	
3	12	15																	
5	15	19																	

$$A(x, y) = \sum_{x' \leq x, y' \leq y} I(x', y')$$

Integral Image

$I(x, y)$	$A(x, y)$																		
<p>original image</p> <table border="1"><tr><td>1</td><td>5</td><td>2</td></tr><tr><td>2</td><td>4</td><td>1</td></tr><tr><td>2</td><td>1</td><td>1</td></tr></table>	1	5	2	2	4	1	2	1	1	<table border="1"><tr><td>1</td><td>6</td><td>8</td></tr><tr><td>3</td><td>12</td><td>15</td></tr><tr><td>5</td><td>15</td><td>19</td></tr></table> <p>integral image</p>	1	6	8	3	12	15	5	15	19
1	5	2																	
2	4	1																	
2	1	1																	
1	6	8																	
3	12	15																	
5	15	19																	

$$A(x, y) = \sum_{x' \leq x, y' \leq y} I(x', y')$$

Can find the **sum** of any block using **3** operations

$$A(x_1, y_1, x_2, y_2) = A(x_2, y_2) - A(x_1, y_2) - A(x_2, y_1) + A(x_1, y_1)$$

What is the sum of the bottom right 2x2 square?

$$A(x_1, y_1, x_2, y_2) = A(x_2, y_2) - A(x_1, y_2) - A(x_2, y_1) + A(x_1, y_1)$$

$I(x, y)$

1	5	2
2	4	1
2	1	1

image

$A(x, y)$

1	6	8
3	12	15
5	15	19

integral image

$$\begin{aligned} A(1, 1, 3, 3) &= A(3, 3) - A(1, 3) - A(3, 1) + A(1, 1) \\ &= 19 - 8 - 5 + 1 \\ &= 7 \end{aligned}$$

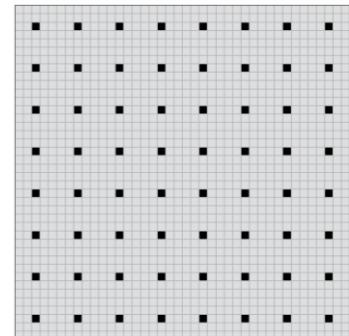
Multi-Scale Oriented Patches (MOPS)

Multi-Image Matching using Multi-Scale Oriented Patches. M. Brown, R. Szeliski and S. Winder.
International Conference on Computer Vision and Pattern Recognition (CVPR2005). pages 510-517

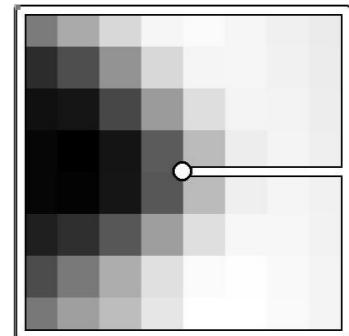
Given a feature (x, y, s, θ)

Get 40×40 image patch, subsample
every 5th pixel

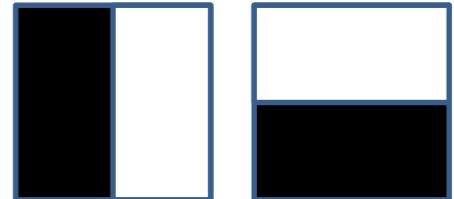
(low frequency filtering, absorbs localization errors)



Subtract the mean, divide by standard
deviation
(removes bias and gain)



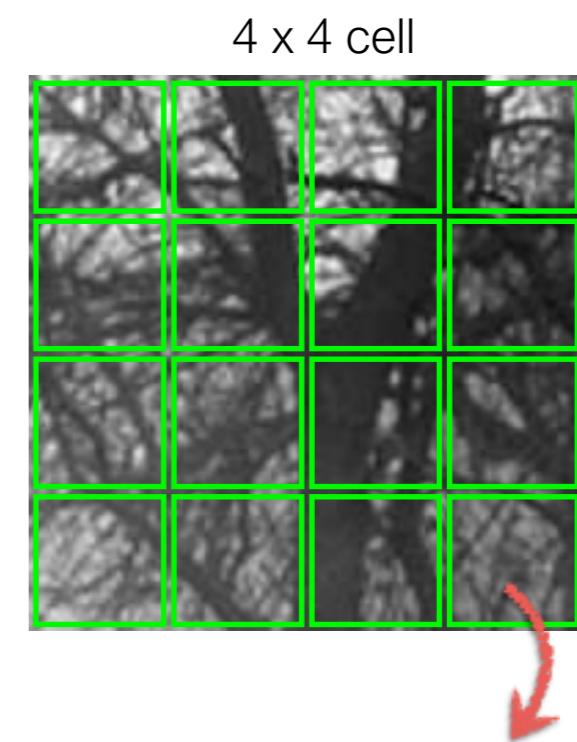
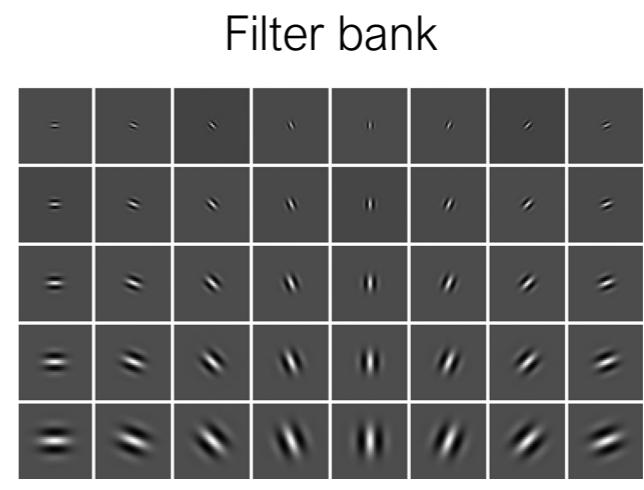
Haar Wavelet Transform
(low frequency projection)



GIST descriptor

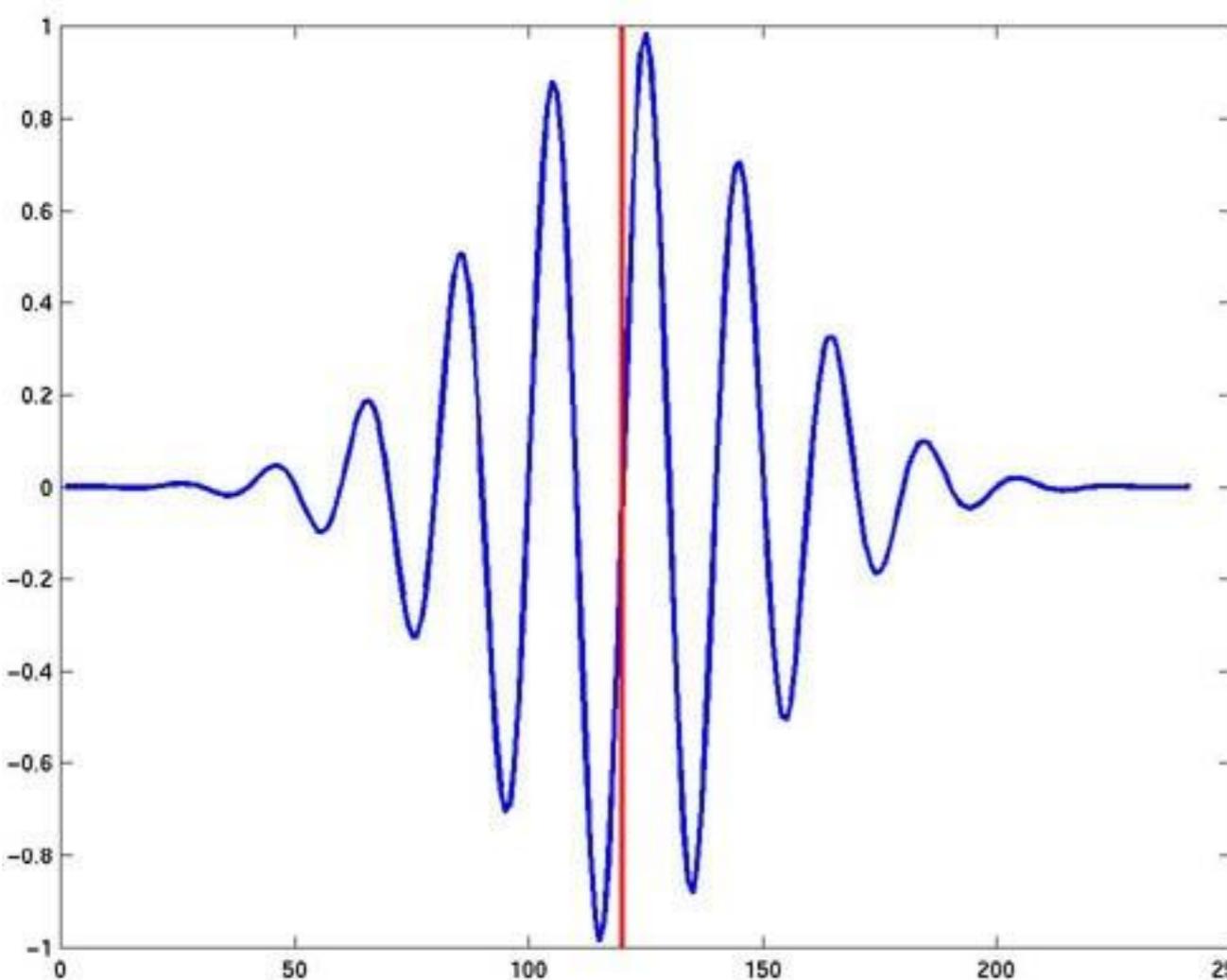
GIST

1. Compute filter responses (filter bank of Gabor filters)
2. Divide image patch into 4×4 cells
3. Compute filter response averages for each cell
4. Size of descriptor is $4 \times 4 \times N$, where N is the size of the filter bank

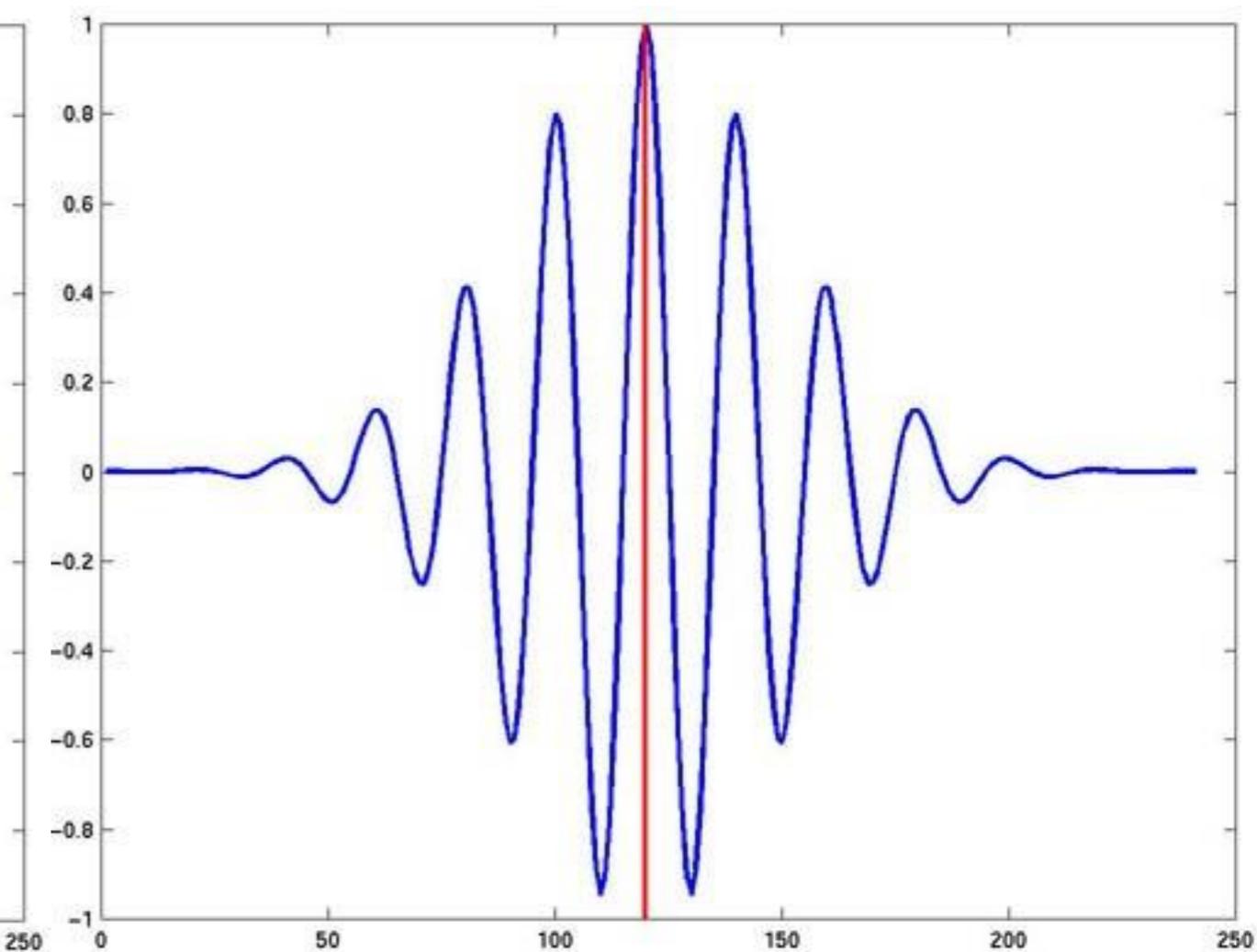


Gabor Filters

(1D examples)



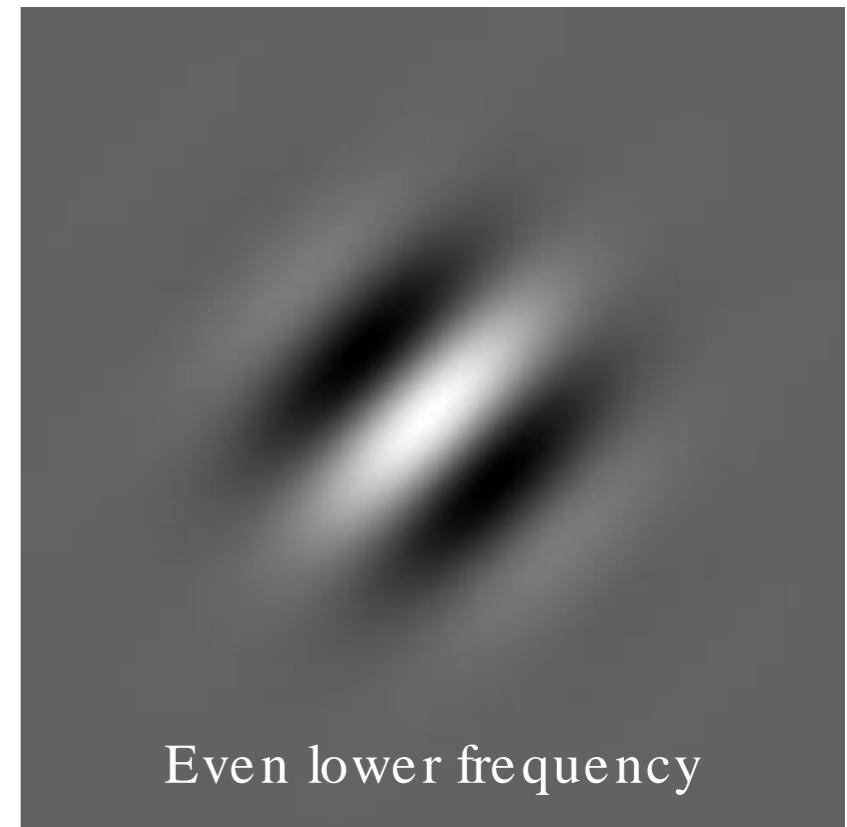
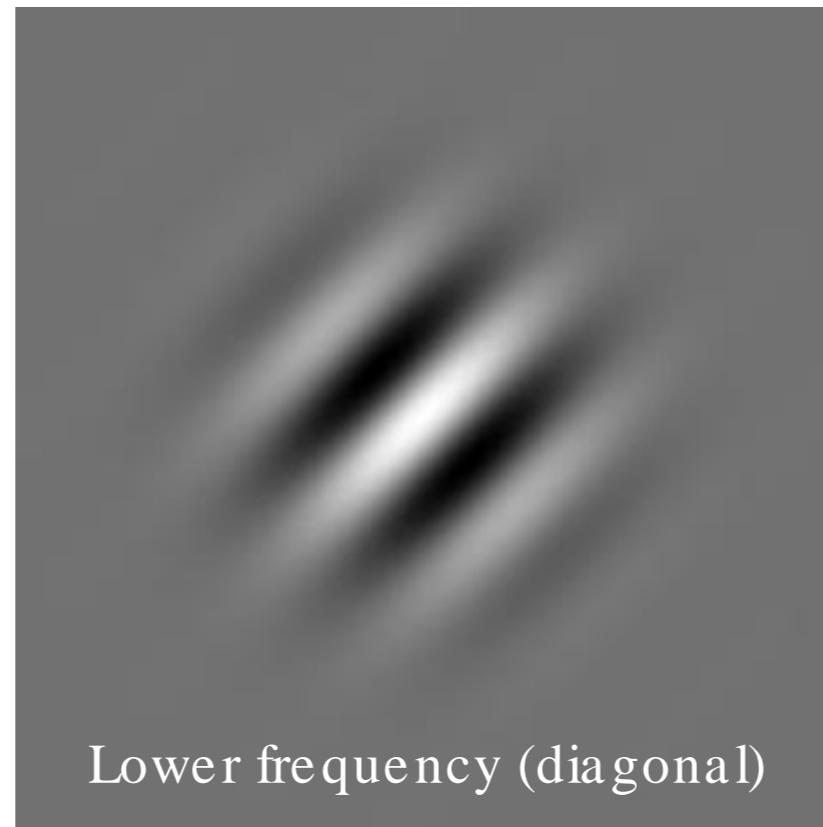
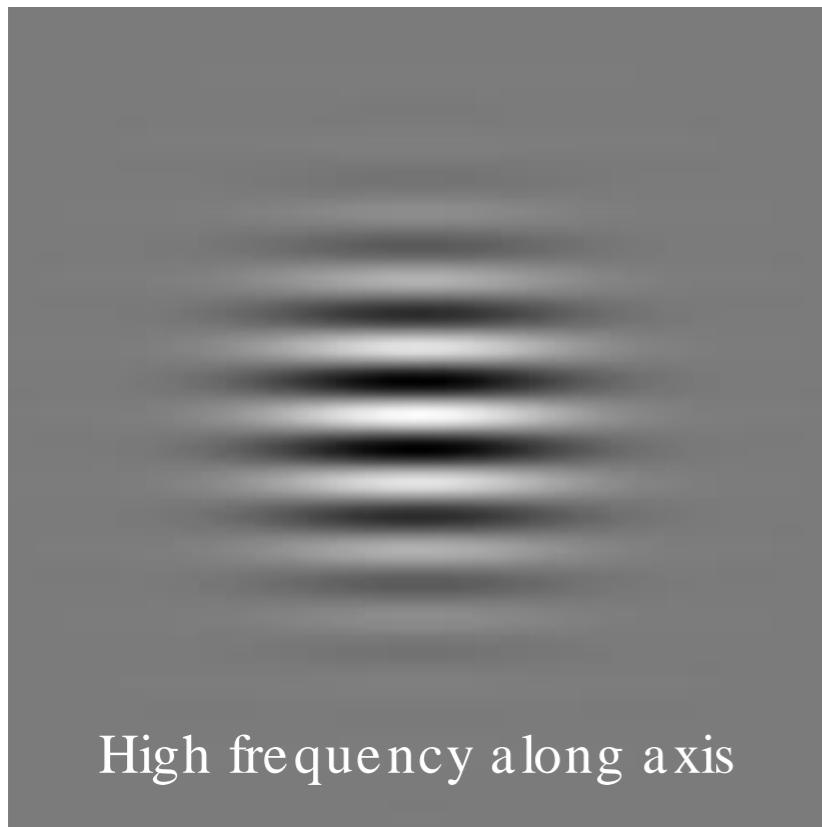
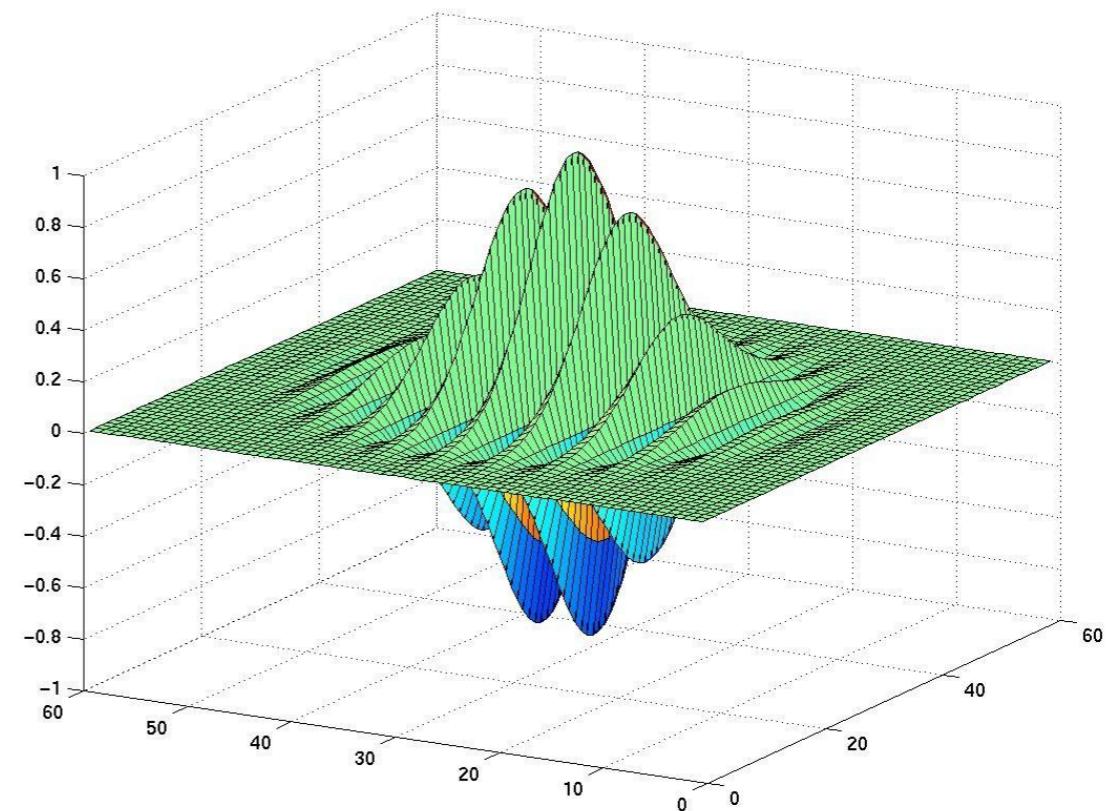
$$e^{-\frac{x^2}{2\sigma^2}} \sin(2\pi\omega x)$$

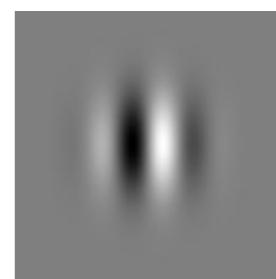
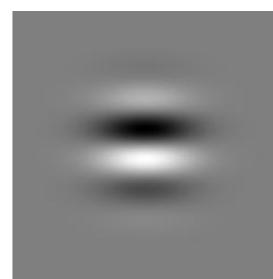
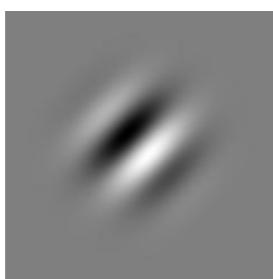
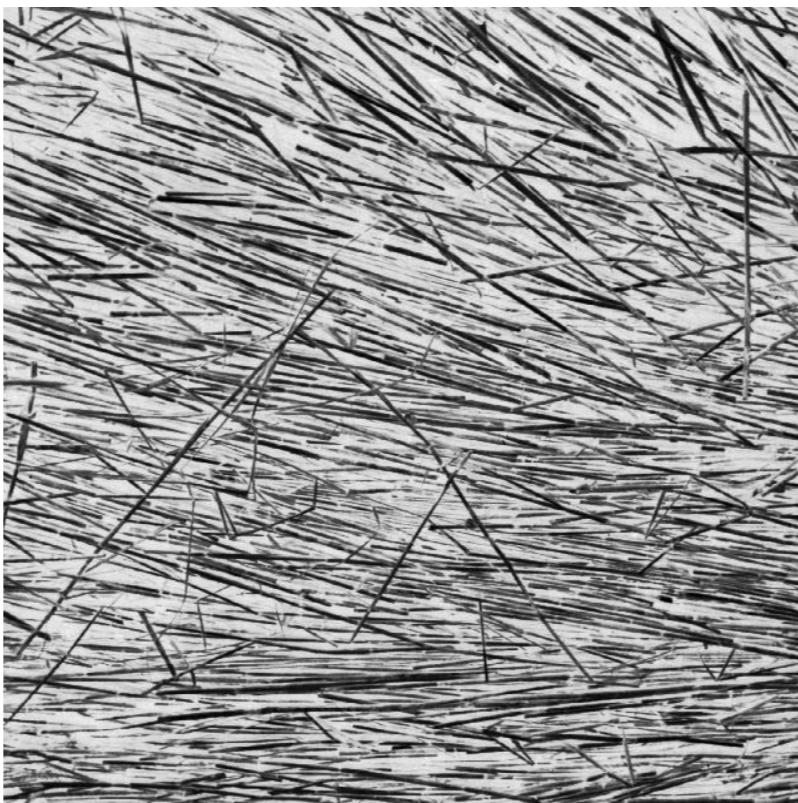


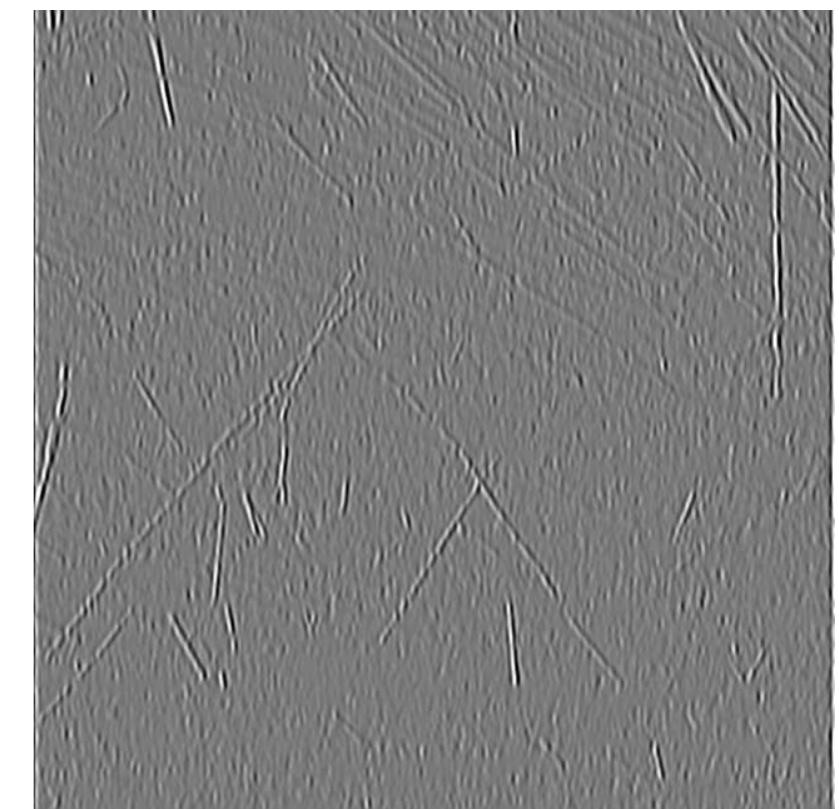
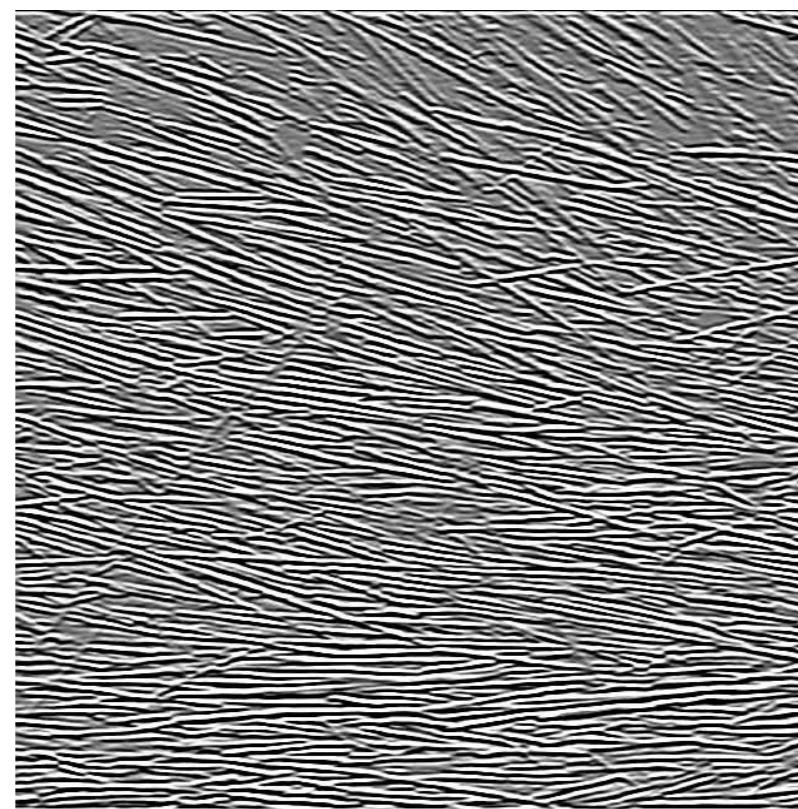
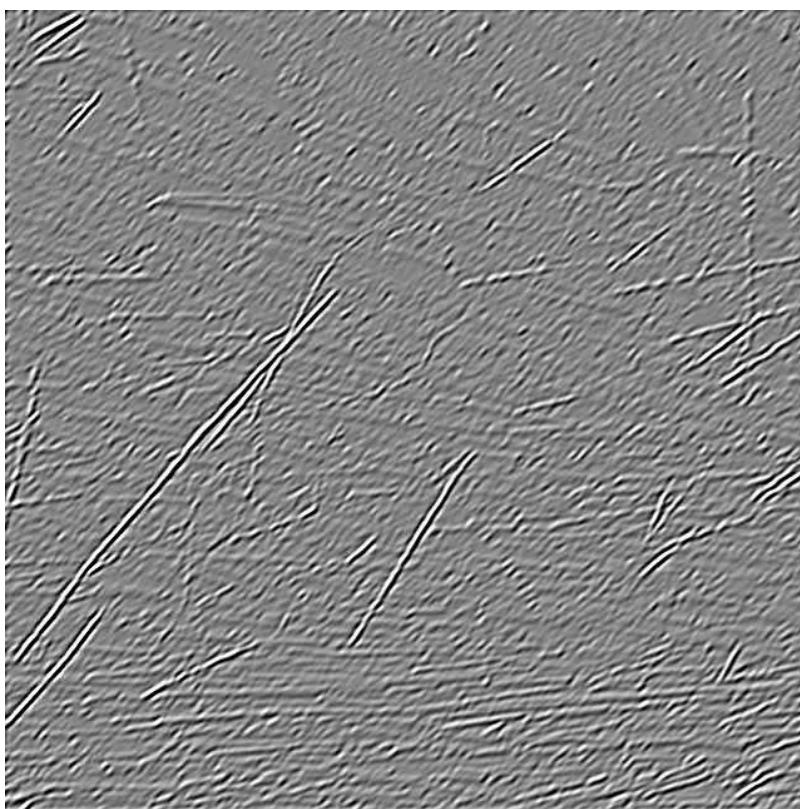
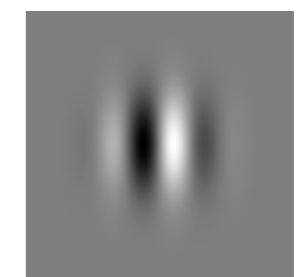
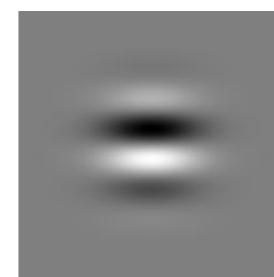
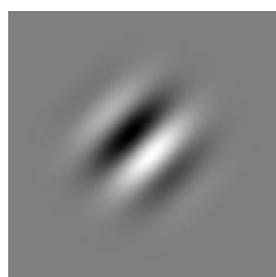
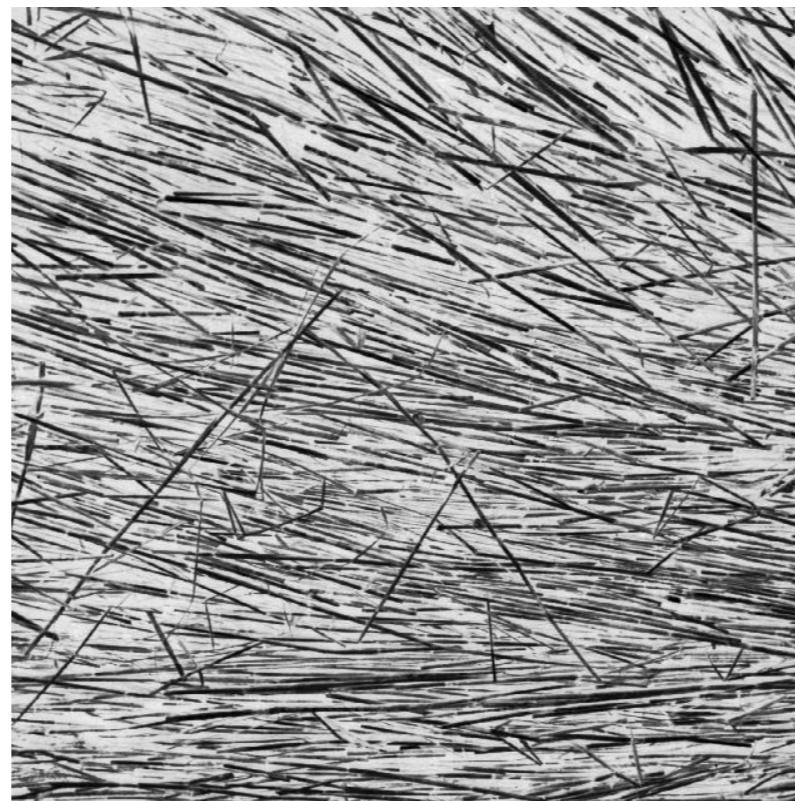
$$e^{-\frac{x^2}{2\sigma^2}} \cos(2\pi\omega x)$$

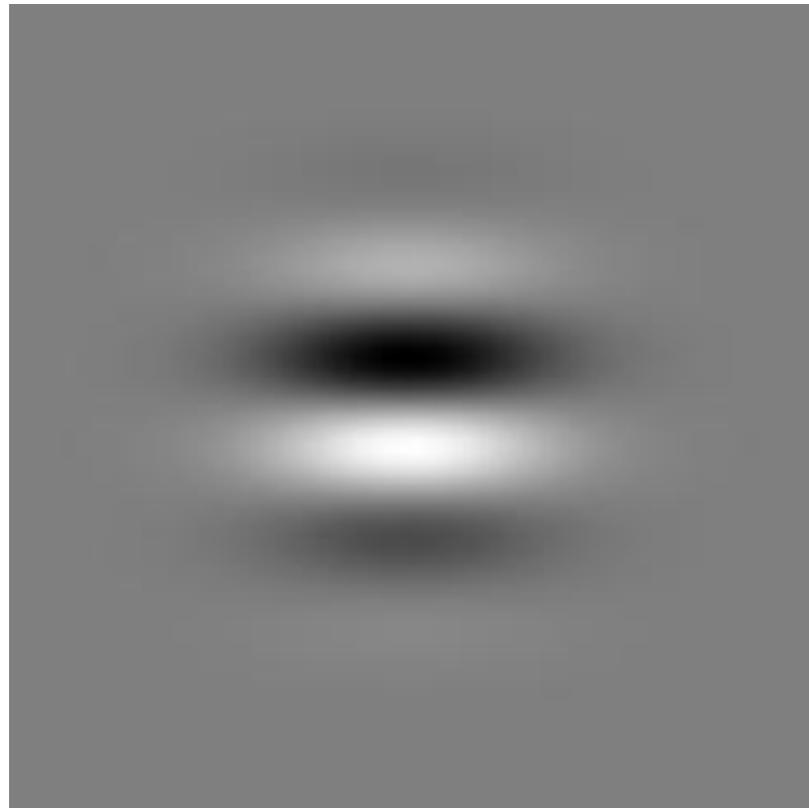
2D Gabor Filters

$$e^{-\frac{x^2+y^2}{2\sigma^2}} \cos(2\pi(k_x x + k_y y))$$

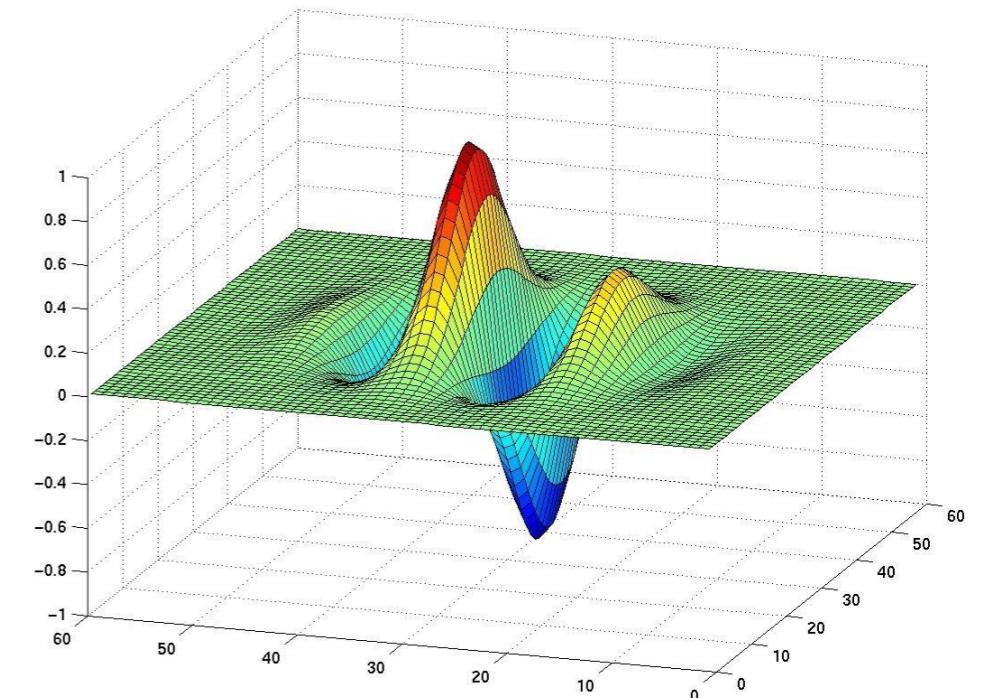




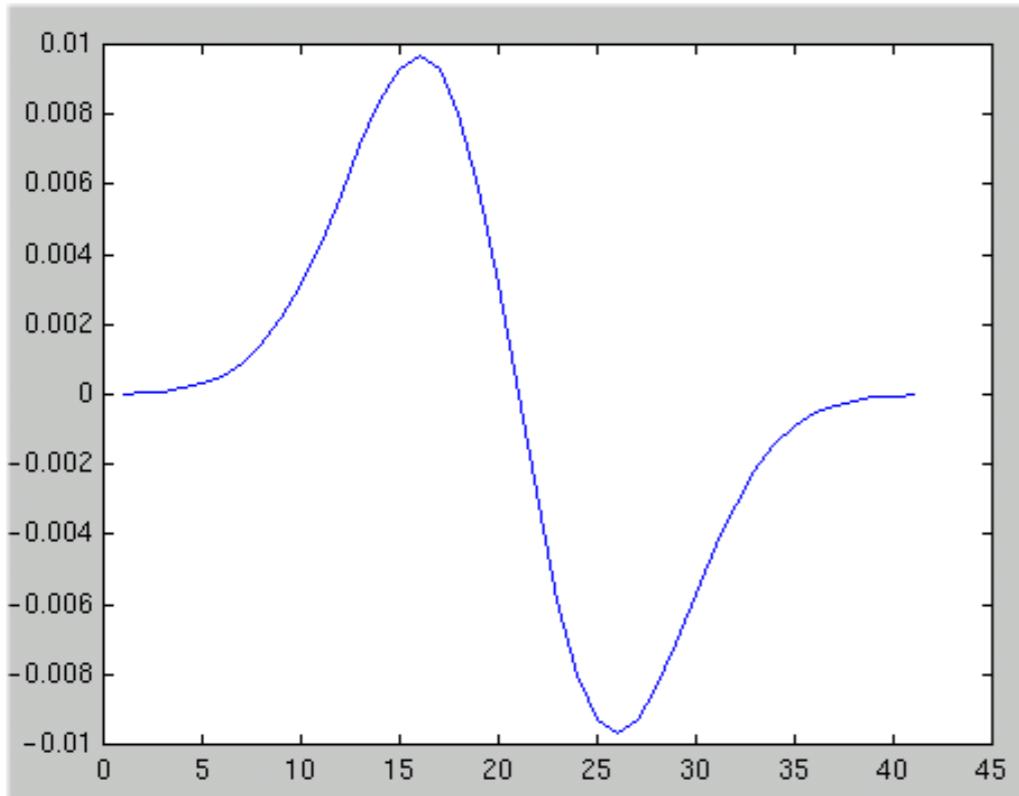




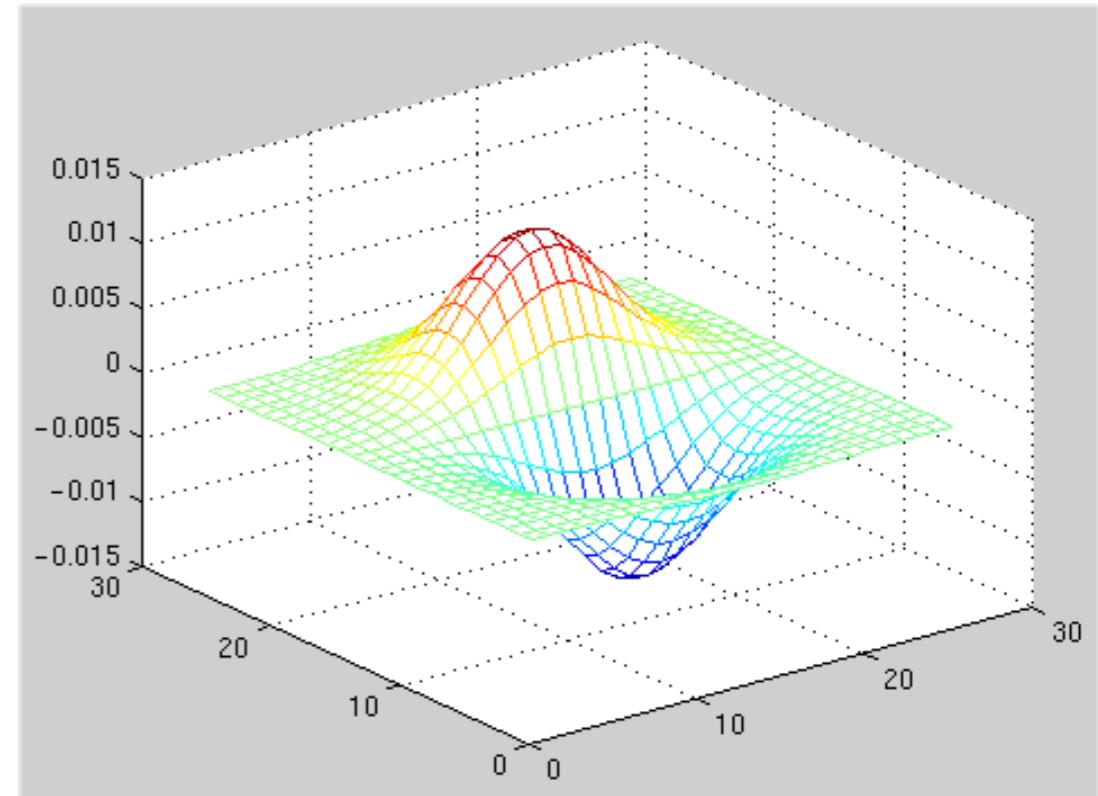
Odd
Gabor
filter

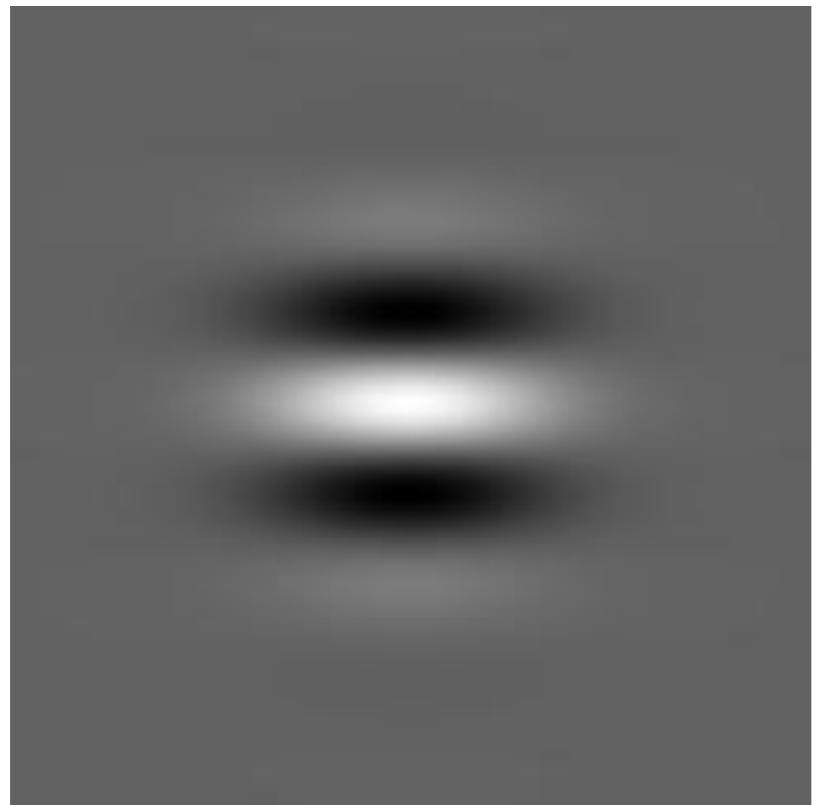


... looks a lot like...

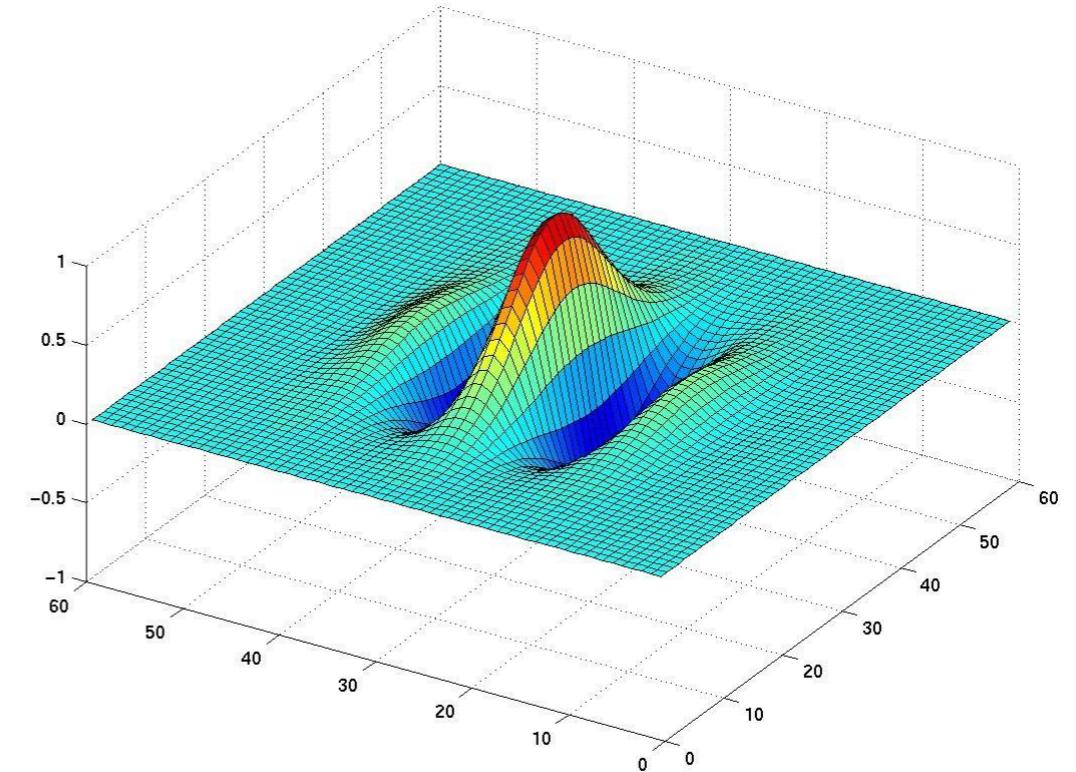


Gaussian
Derivative

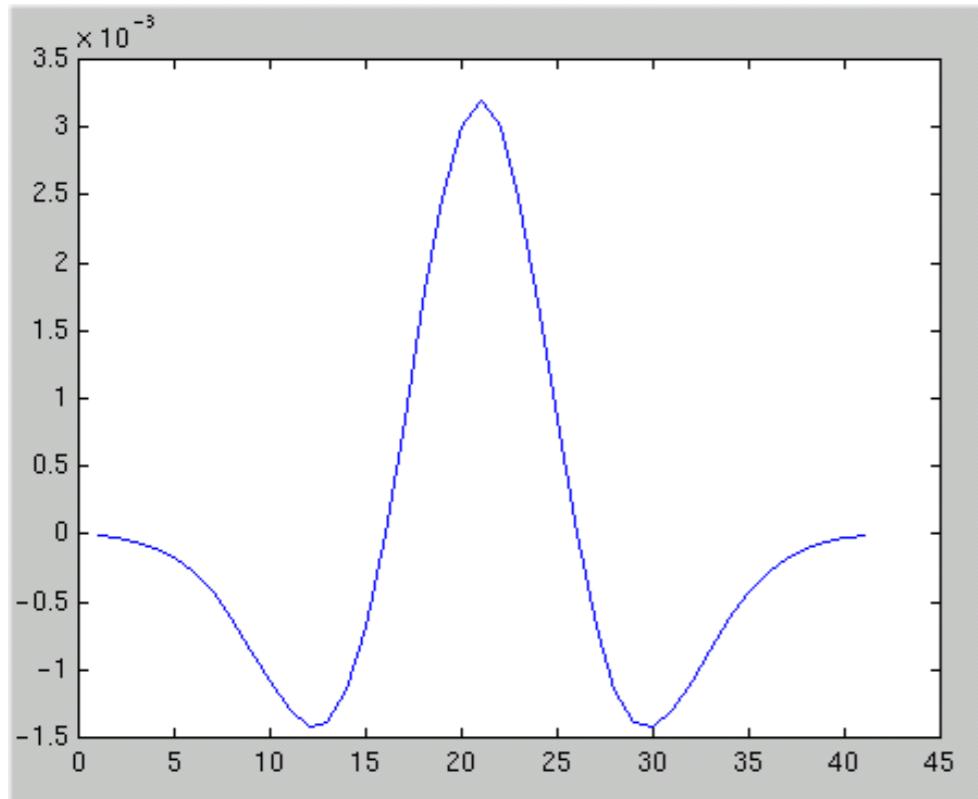




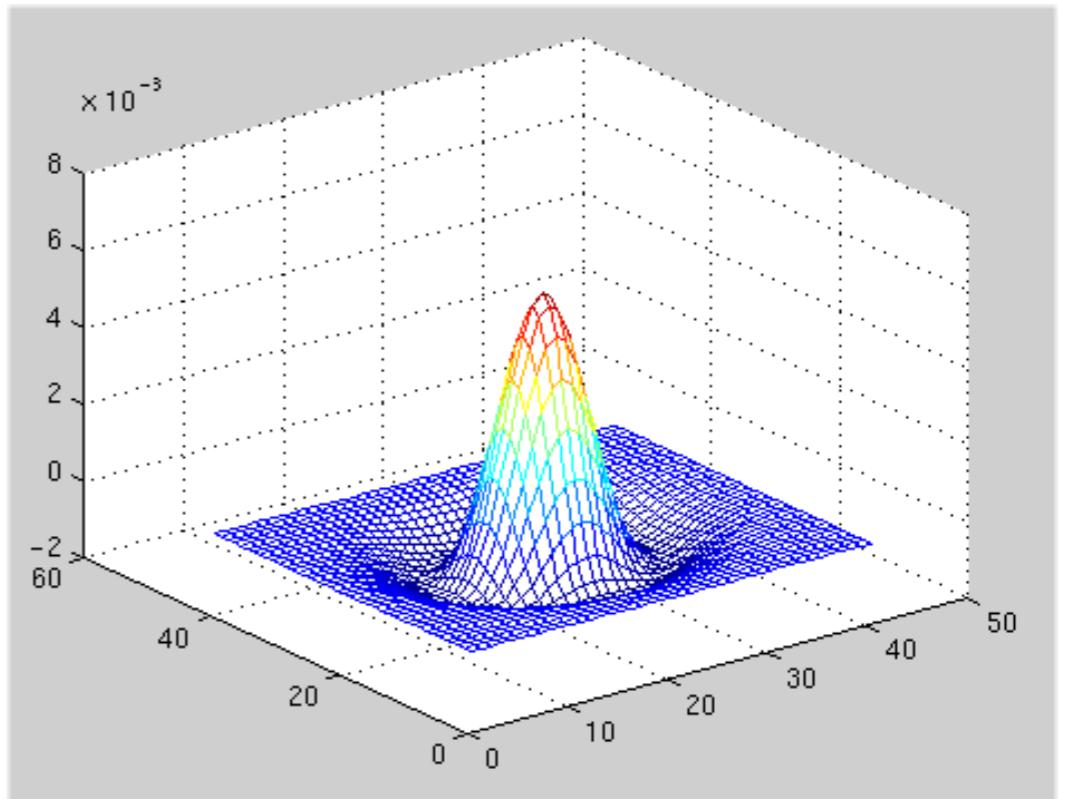
Even
Gabor
filter



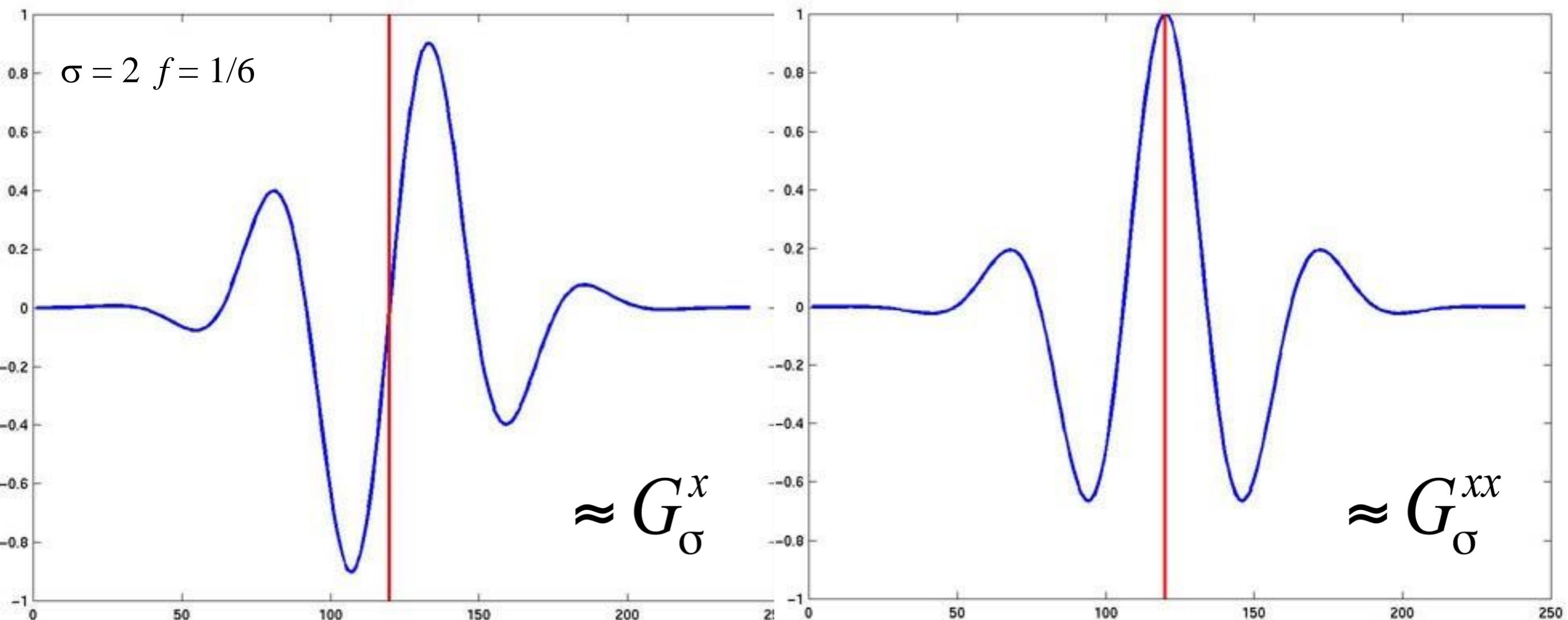
... looks a lot like...



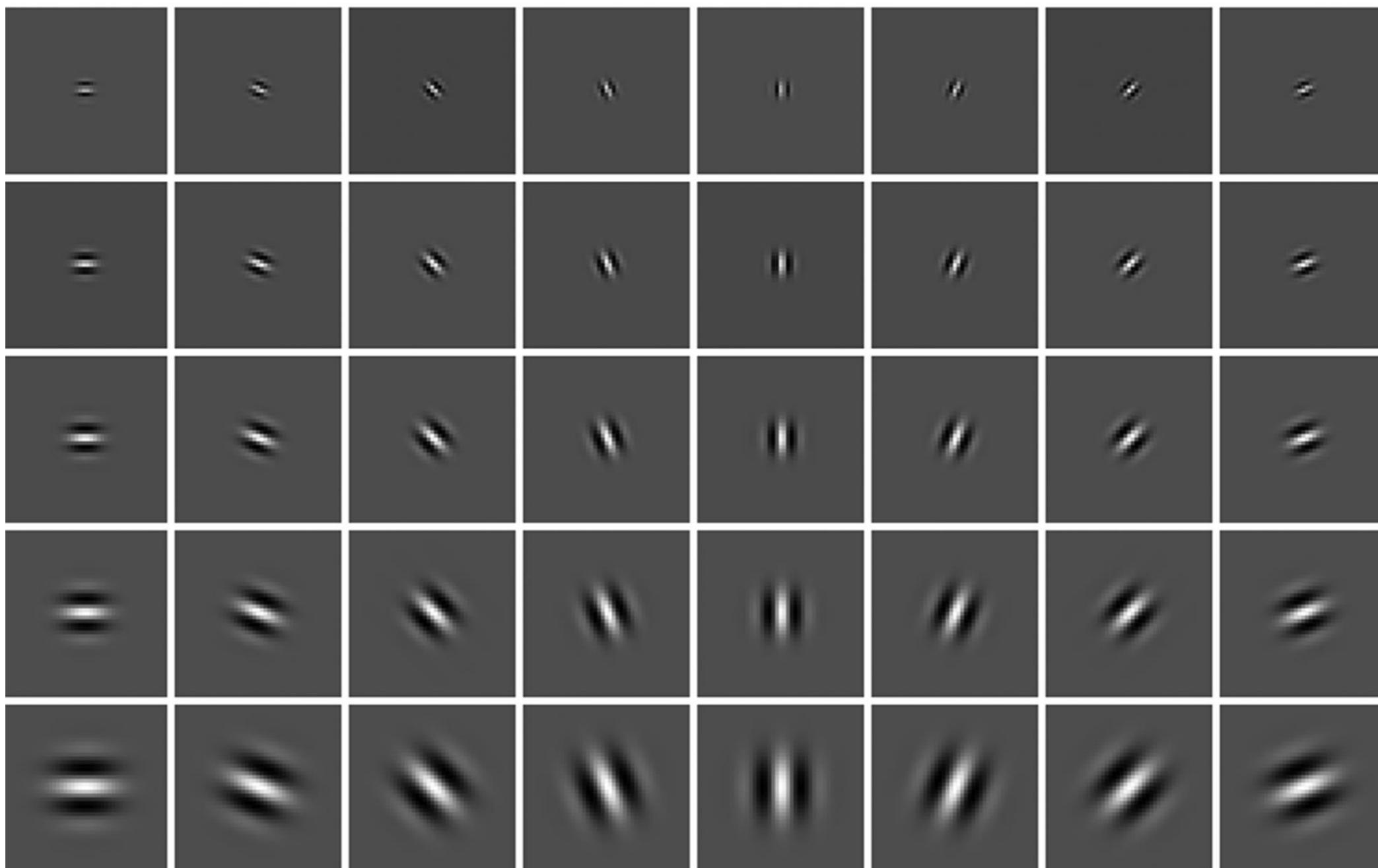
Laplacian



If scale small compared to inverse frequency, the Gabor filters become derivative operators

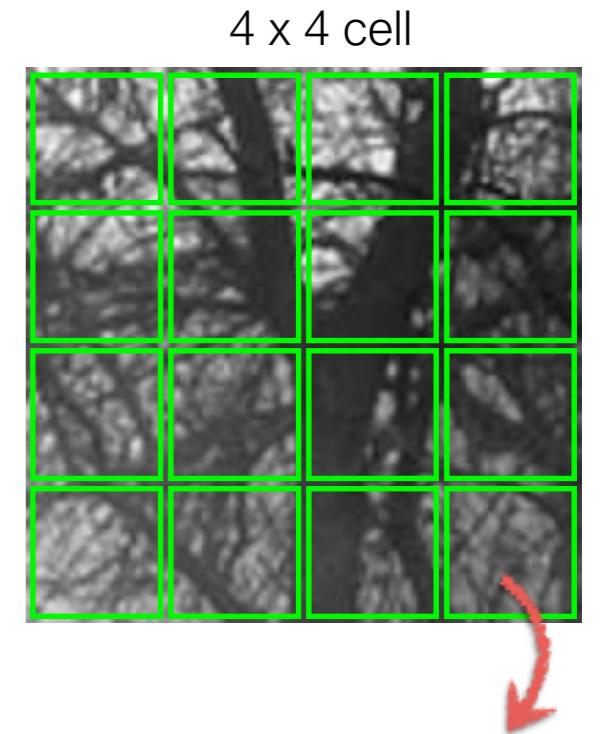
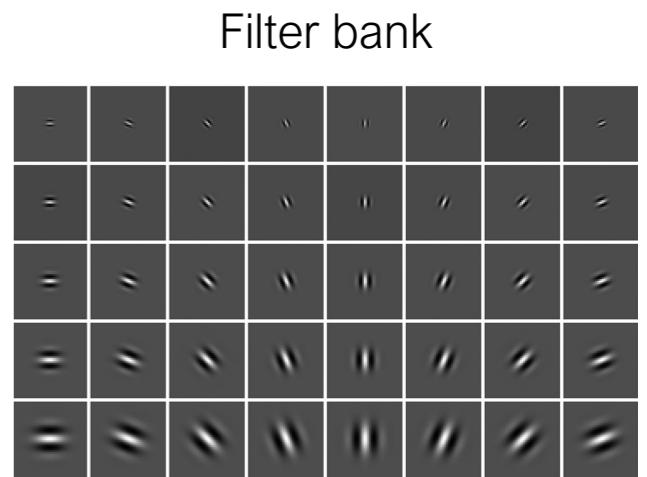


Directional edge detectors



GIST

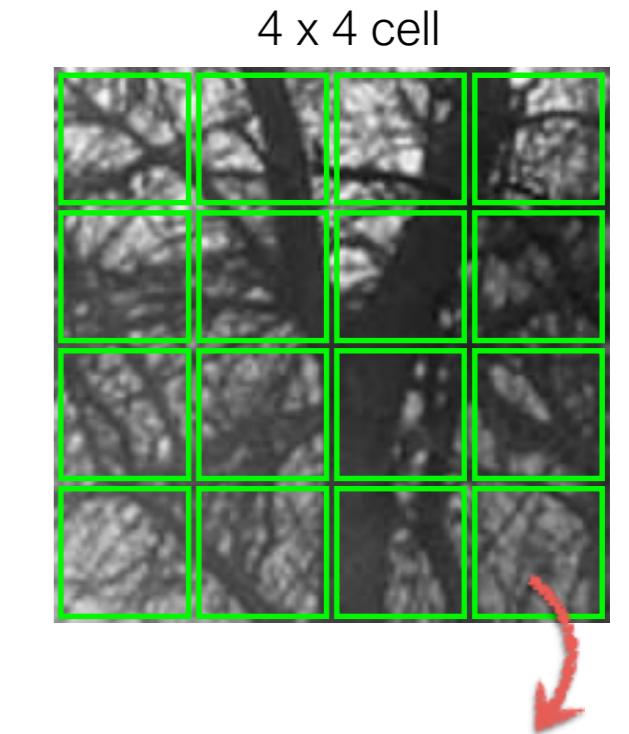
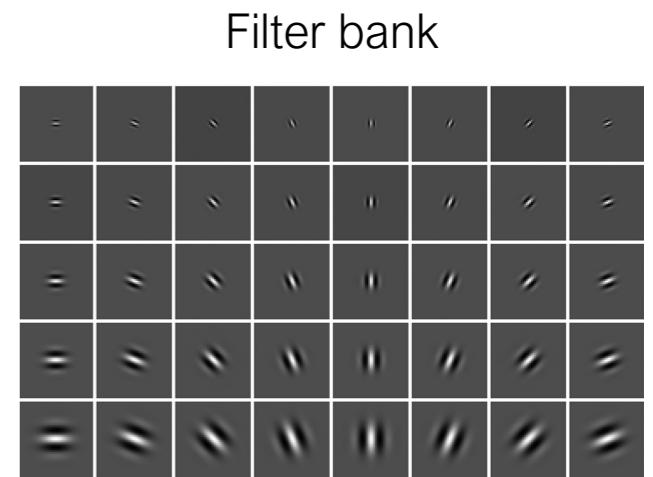
1. Compute filter responses (filter bank of Gabor filters)
2. Divide image patch into 4×4 cells
3. Compute filter response averages for each cell
4. Size of descriptor is $4 \times 4 \times N$, where N is the size of the filter bank



What is the GIST descriptor encoding?

GIST

1. Compute filter responses (filter bank of Gabor filters)
2. Divide image patch into 4×4 cells
3. Compute filter response averages for each cell
4. Size of descriptor is $4 \times 4 \times N$, where N is the size of the filter bank



What is the GIST descriptor encoding?

Rough spatial distribution of image gradients

Histogram of Textons descriptor

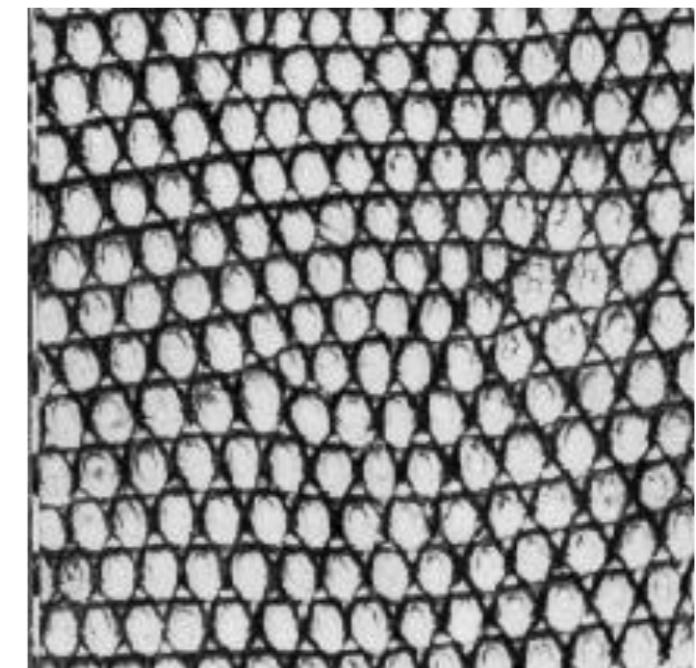
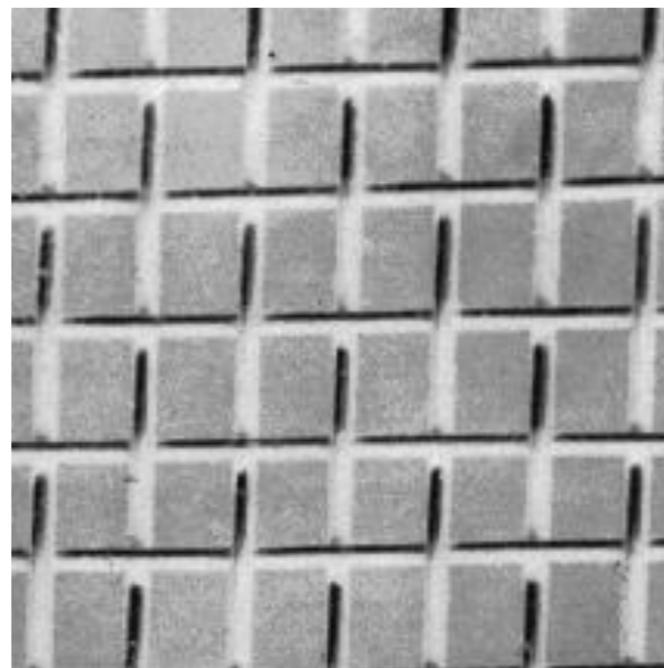
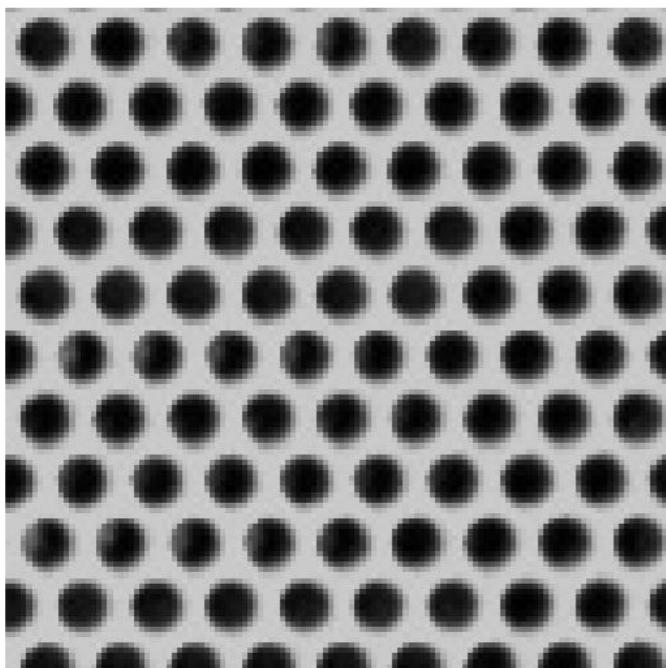
Textons

Julesz. Textons, the elements of texture perception, and their interactions. Nature 1981

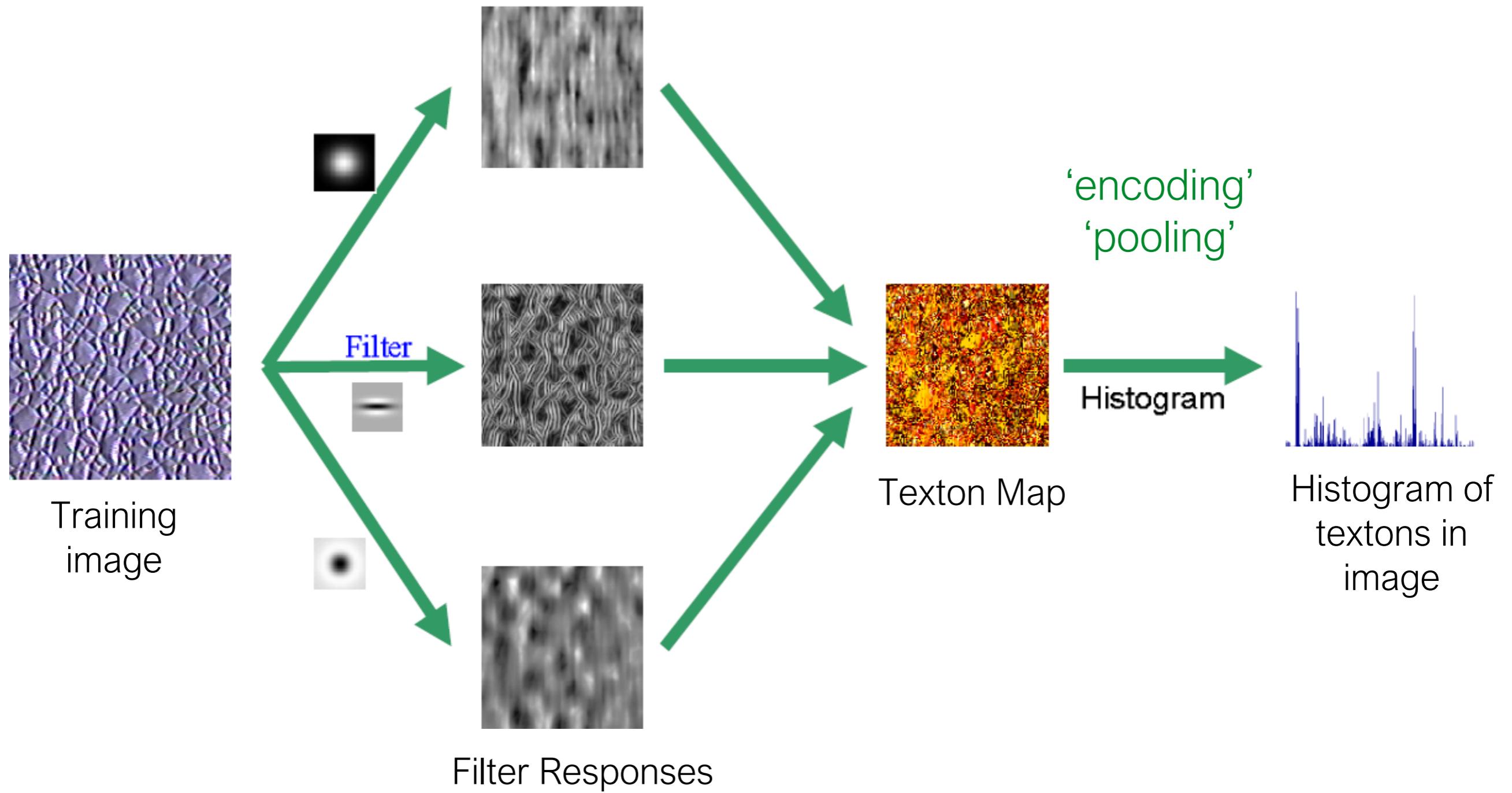
Texture is characterized by the repetition of basic elements or ***textons***



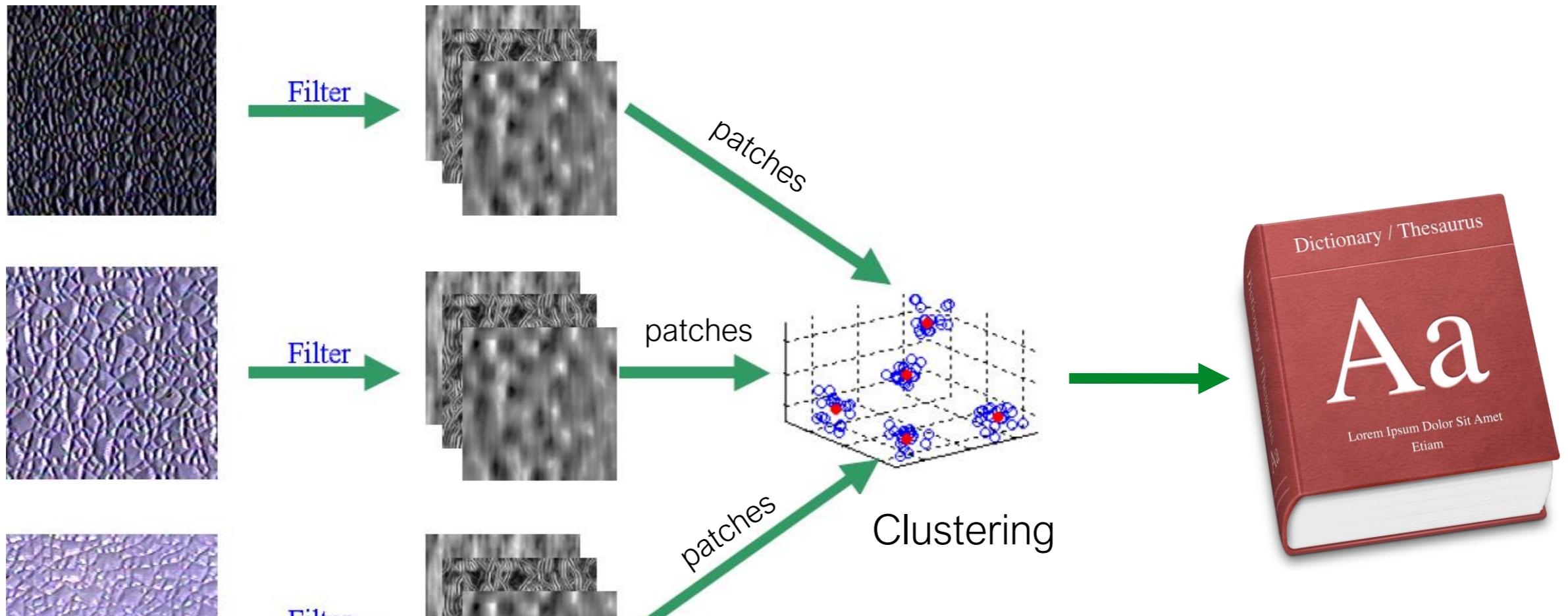
For stochastic textures, it is the identity of the ***textons***, not
their spatial arrangement, that matters



Histogram of Textons descriptor



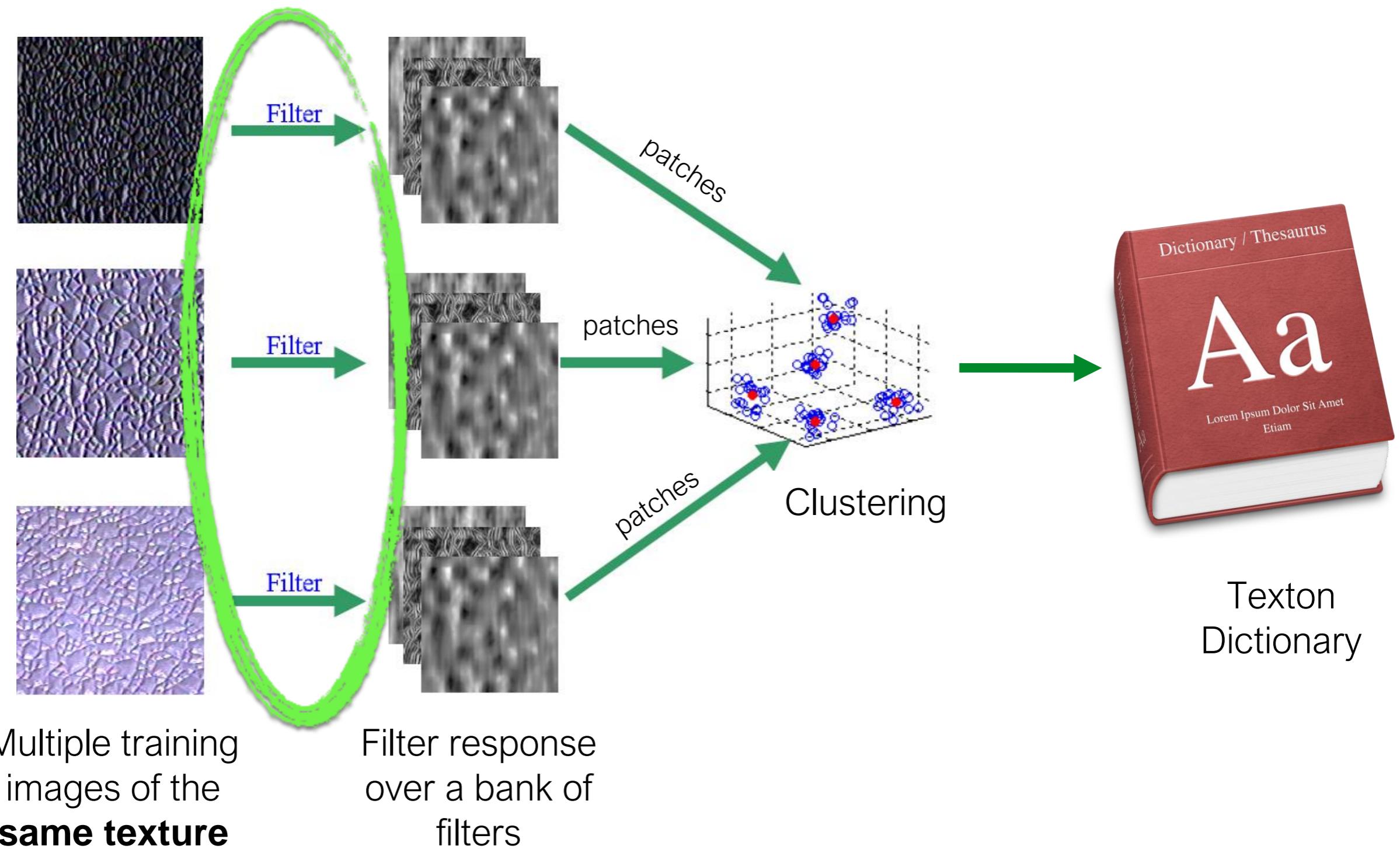
Learning Textons from data



Multiple training
images of the
same texture

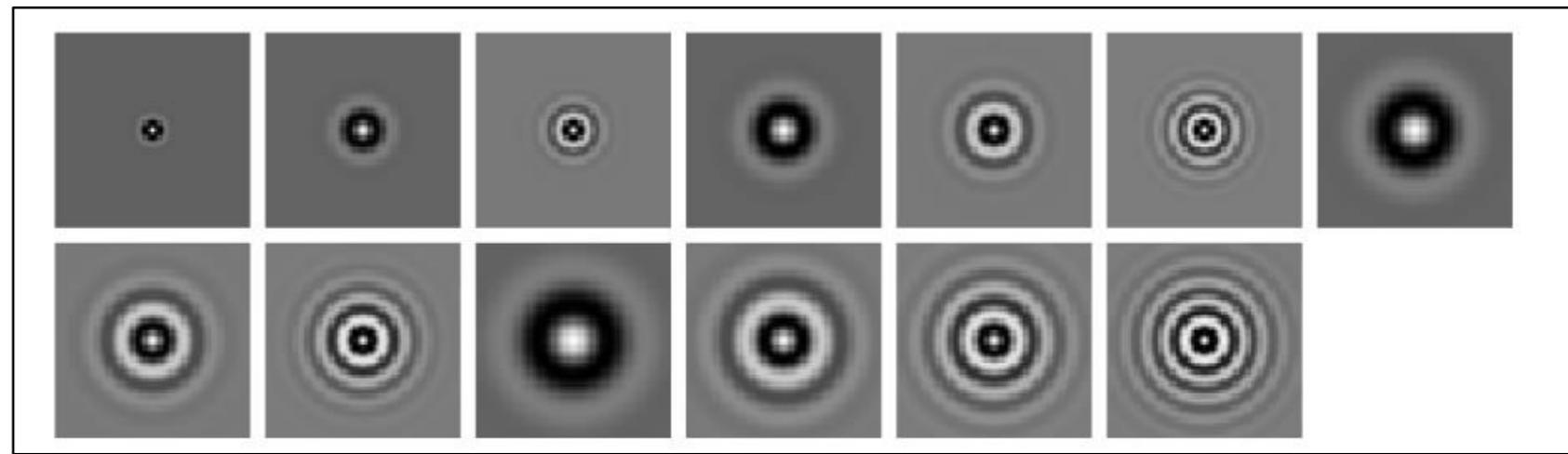
Filter response
over a bank of
filters

Learning Textons from data

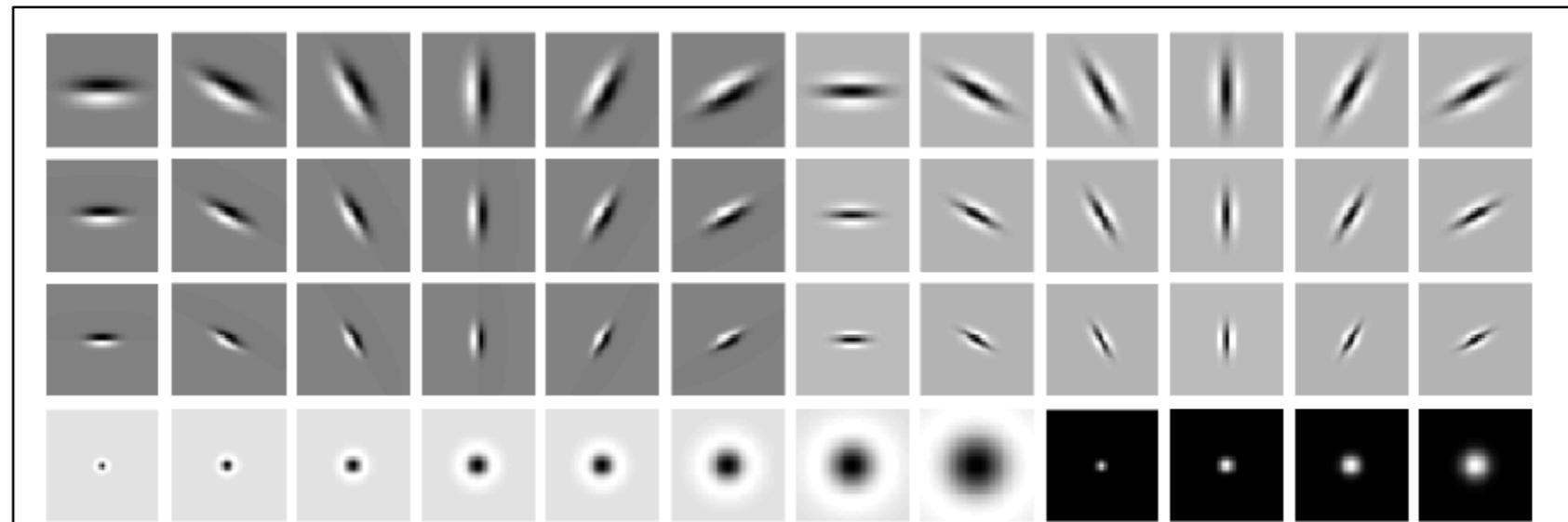


Example of Filter Banks

Isotropic Gabor



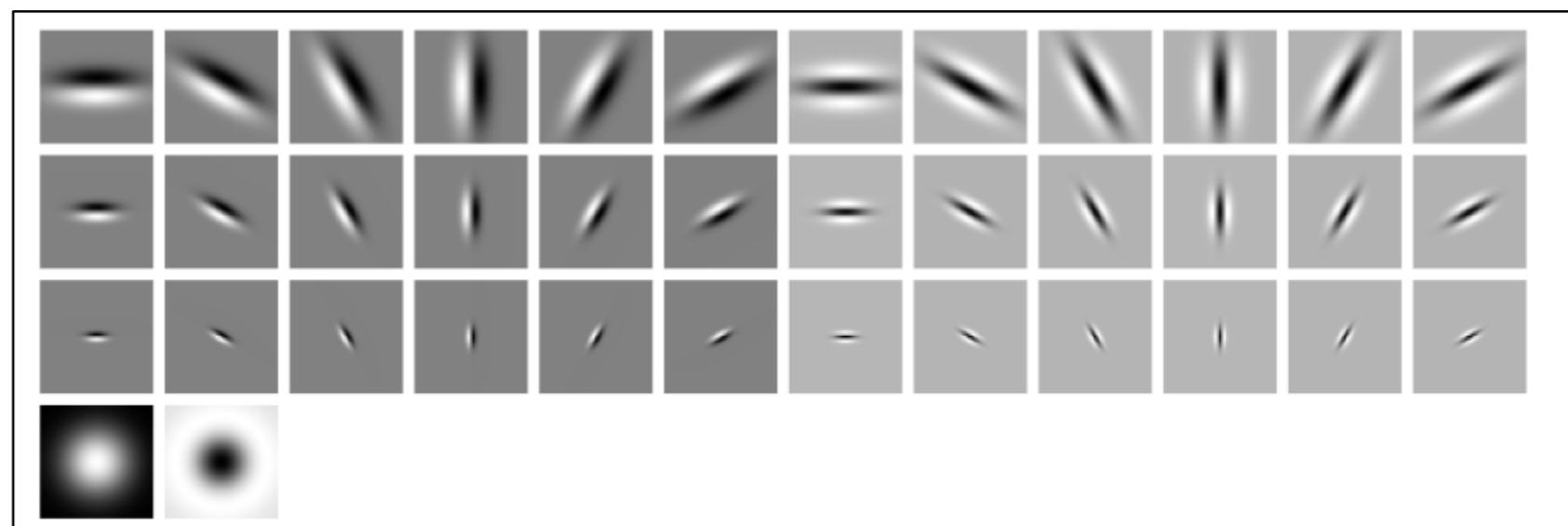
Gaussian derivatives at different scales and orientations



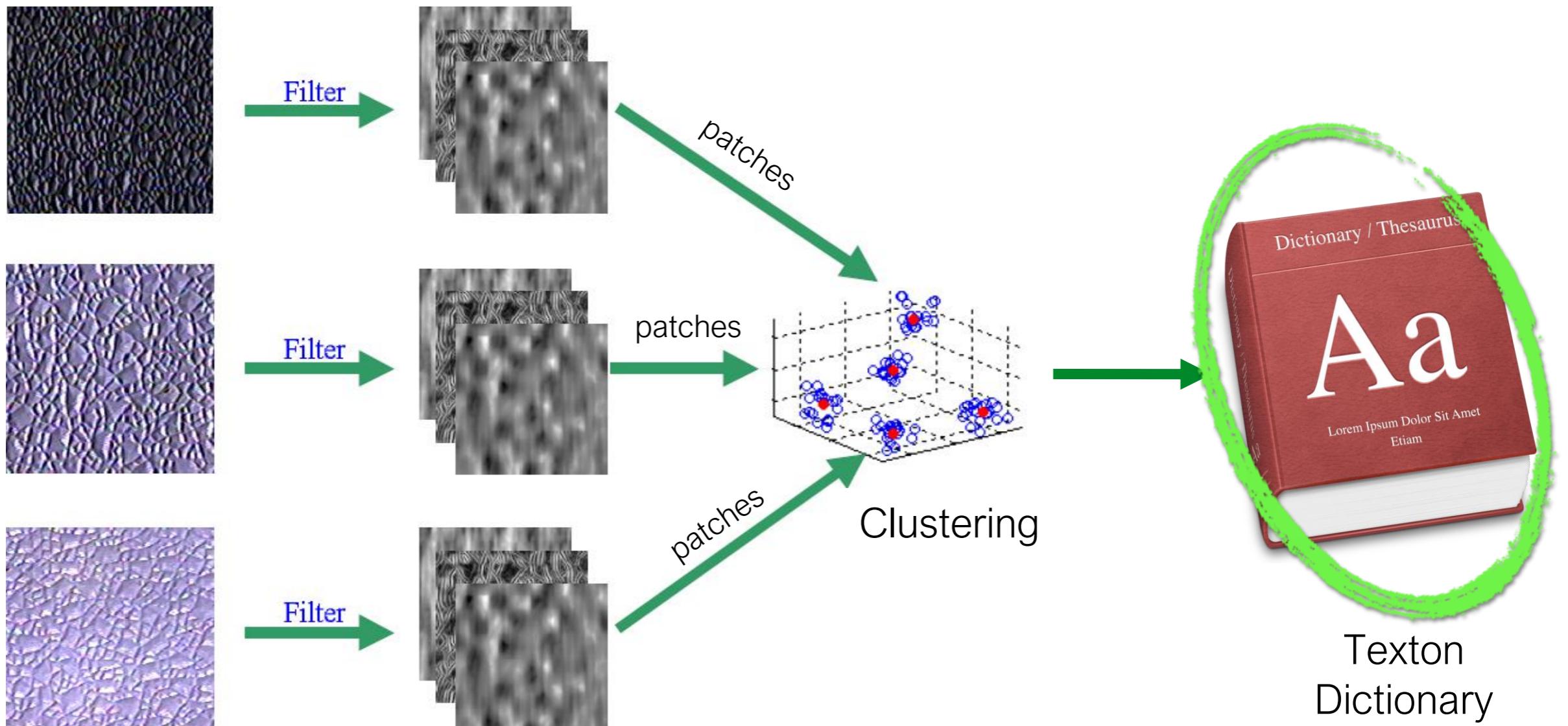
'S'

'LM'

'MR8'



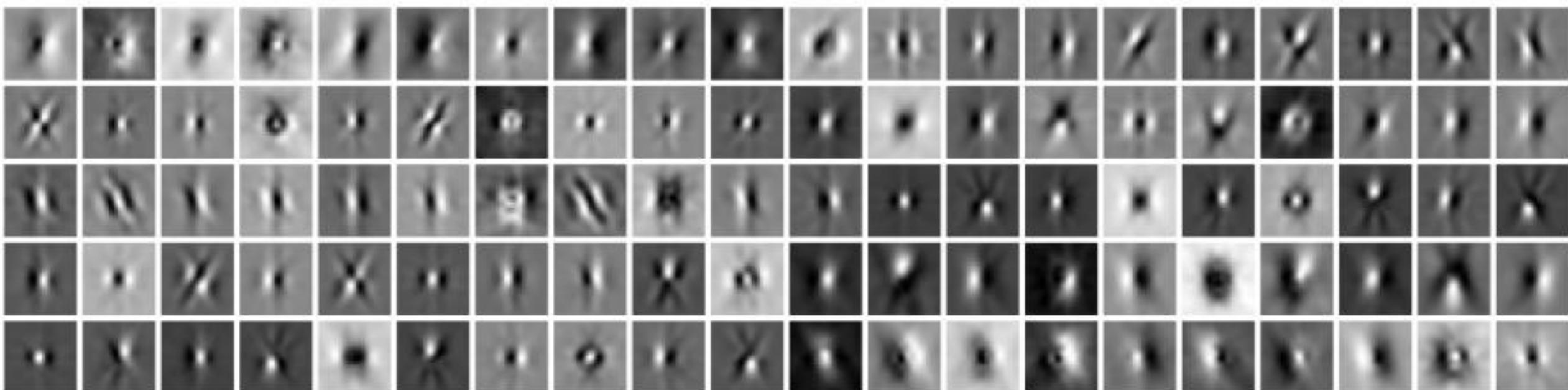
Learning Textons from data



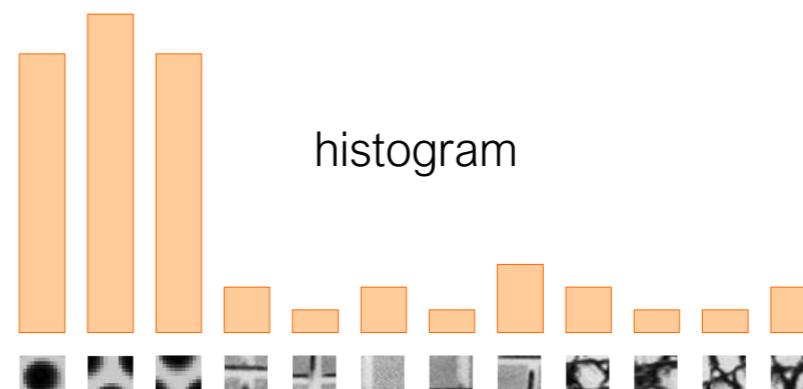
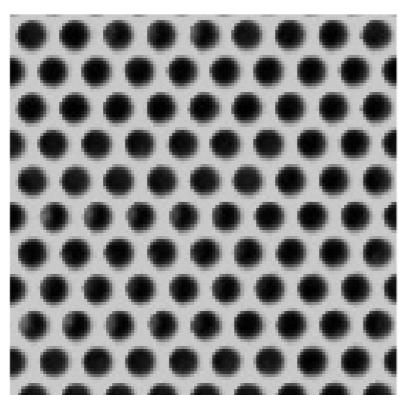
Multiple training
images of the same
texture

Filter response
over a bank of
filters

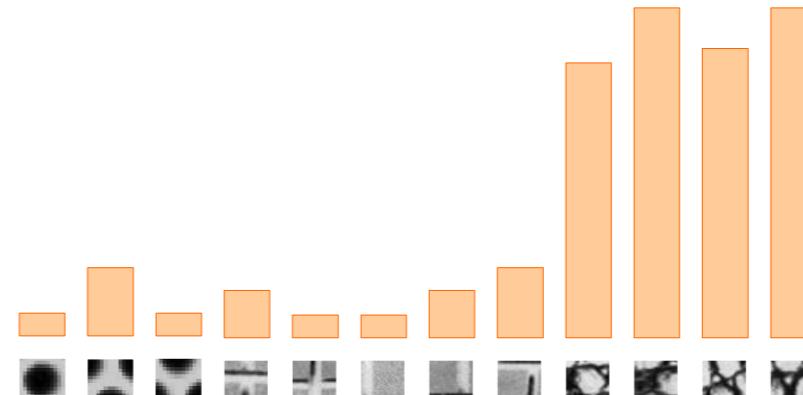
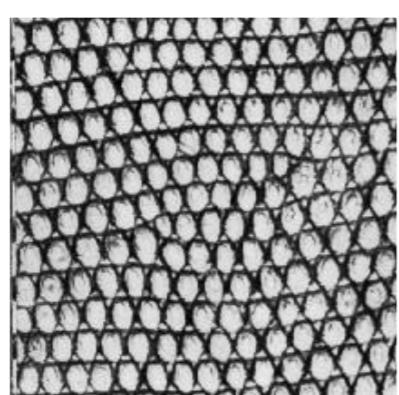
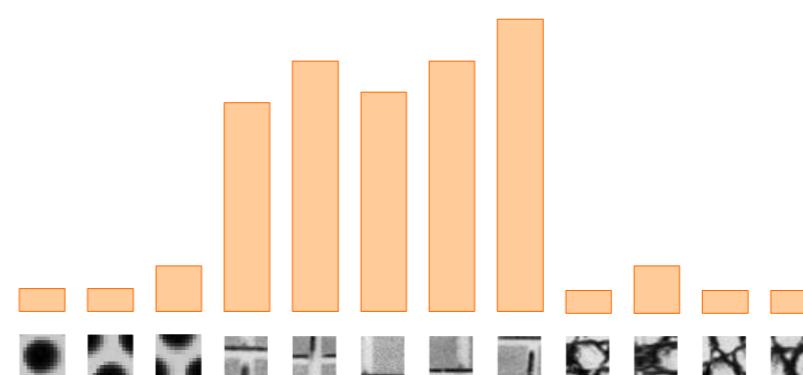
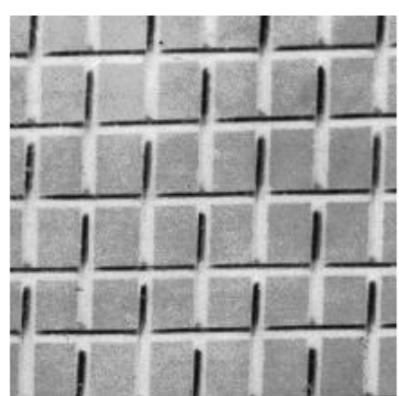
Texton Dictionary



Malik, Belongie, Shi, Leung. Textons, Contours and Regions: Cue Integration in Image Segmentation. ICCV 1999.



Universal texton dictionary



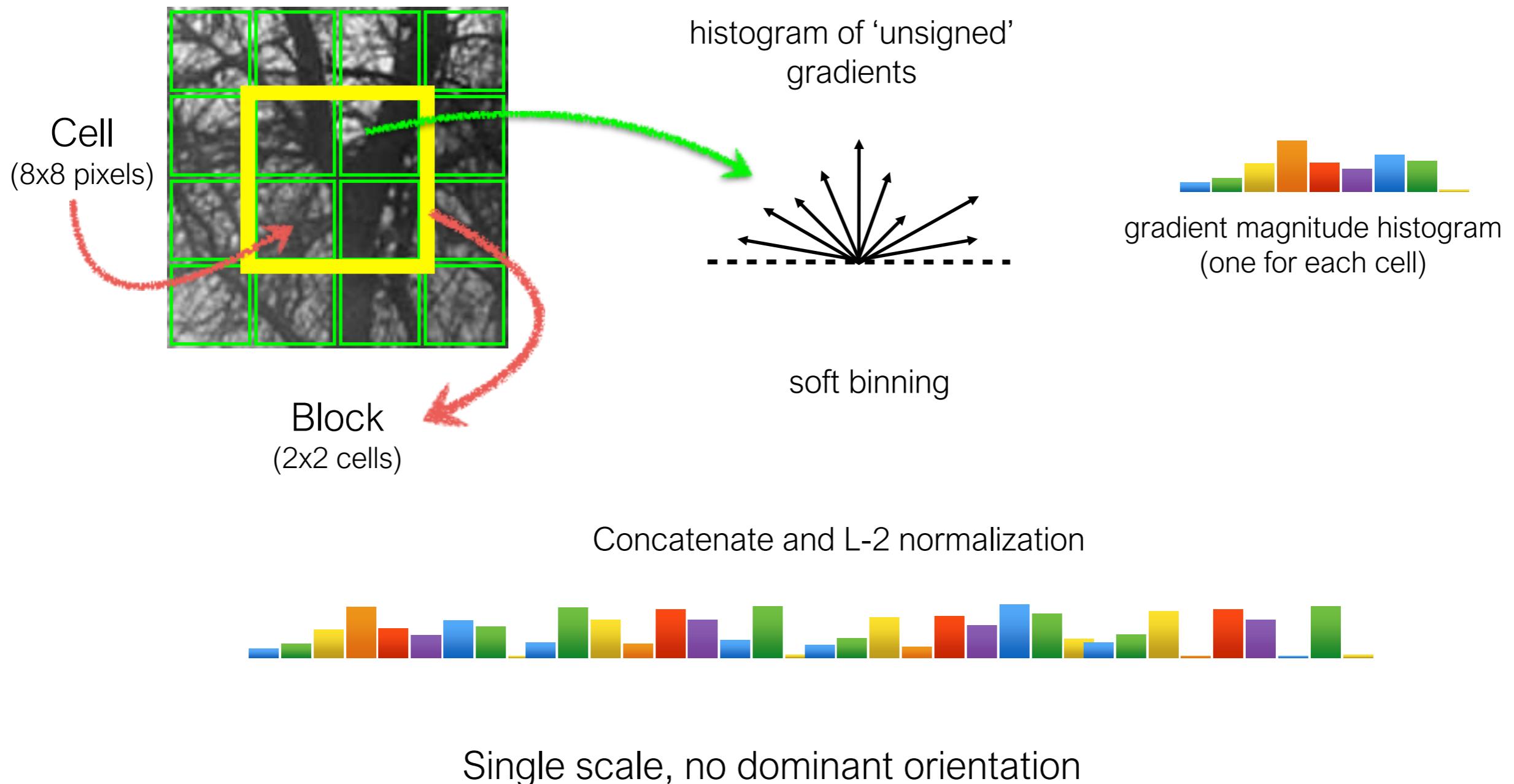
Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

HOG descriptor

HOG



Dalal, Triggs. **Histograms of Oriented Gradients** for Human Detection. CVPR, 2005



Pedestrian detection

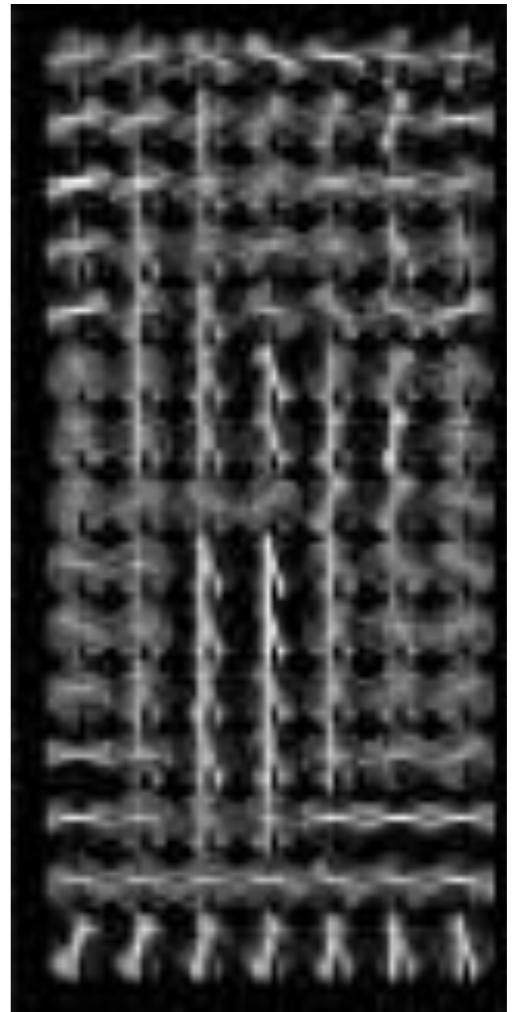
1 cell step size

128 pixels
16 cells
15 blocks



$$15 \times 7 \times 4 \times 36 = 3780$$

visualization



64 pixels
8 cells
7 blocks

Redundant representation due to overlapping blocks
How many times is each inner cell encoded?



SURF descriptor

SURF

(‘Speeded’ Up Robust Features)

Compute Haar wavelet response at each pixel in patch

center of detected feature

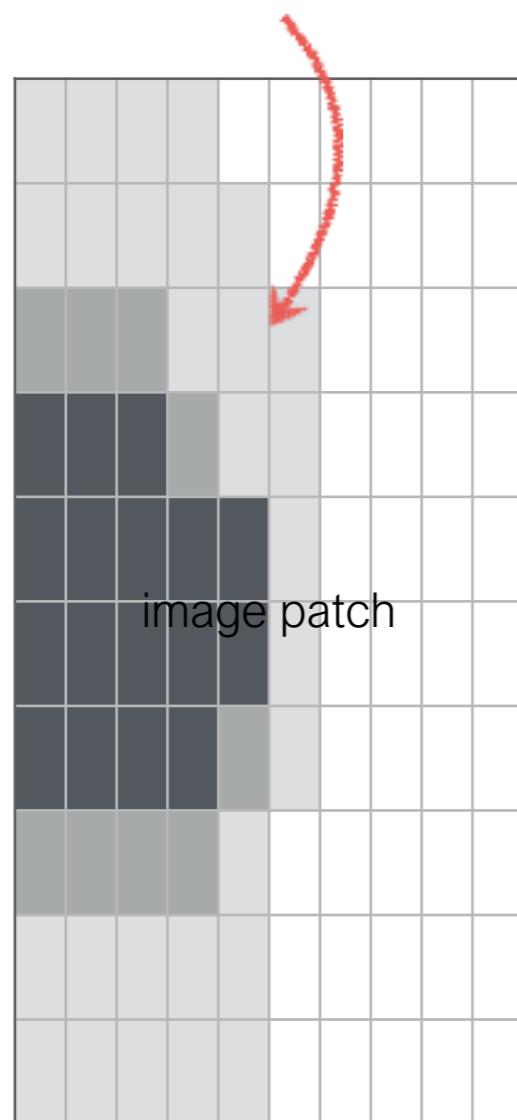
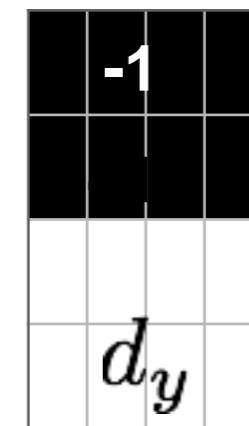
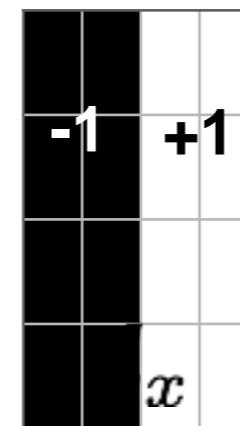


image patch

Haar wavelets filters



(Gaussian weighted from center)

How would do you compute the filter response?

SURF

(‘Speeded’ Up Robust Features)

Compute Haar wavelet response at each pixel in patch

center of detected feature

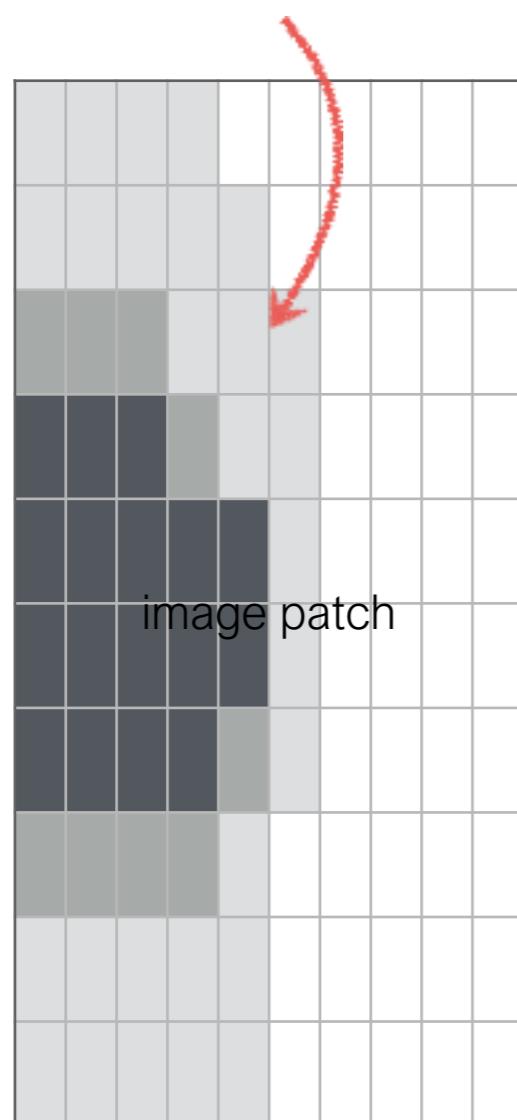
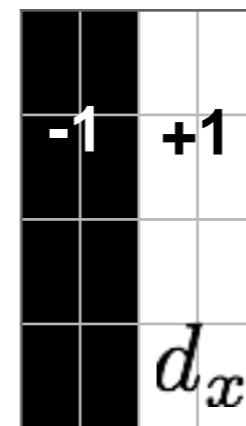
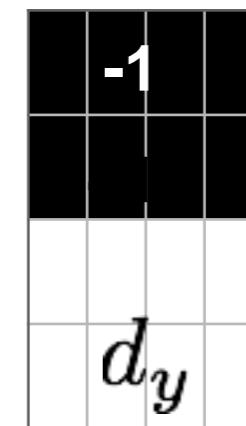


image patch

Haar wavelets filters



d_x



d_y

(Gaussian weighted from center)

How would you compute the filter response?

Filtering using a sliding window can be slow

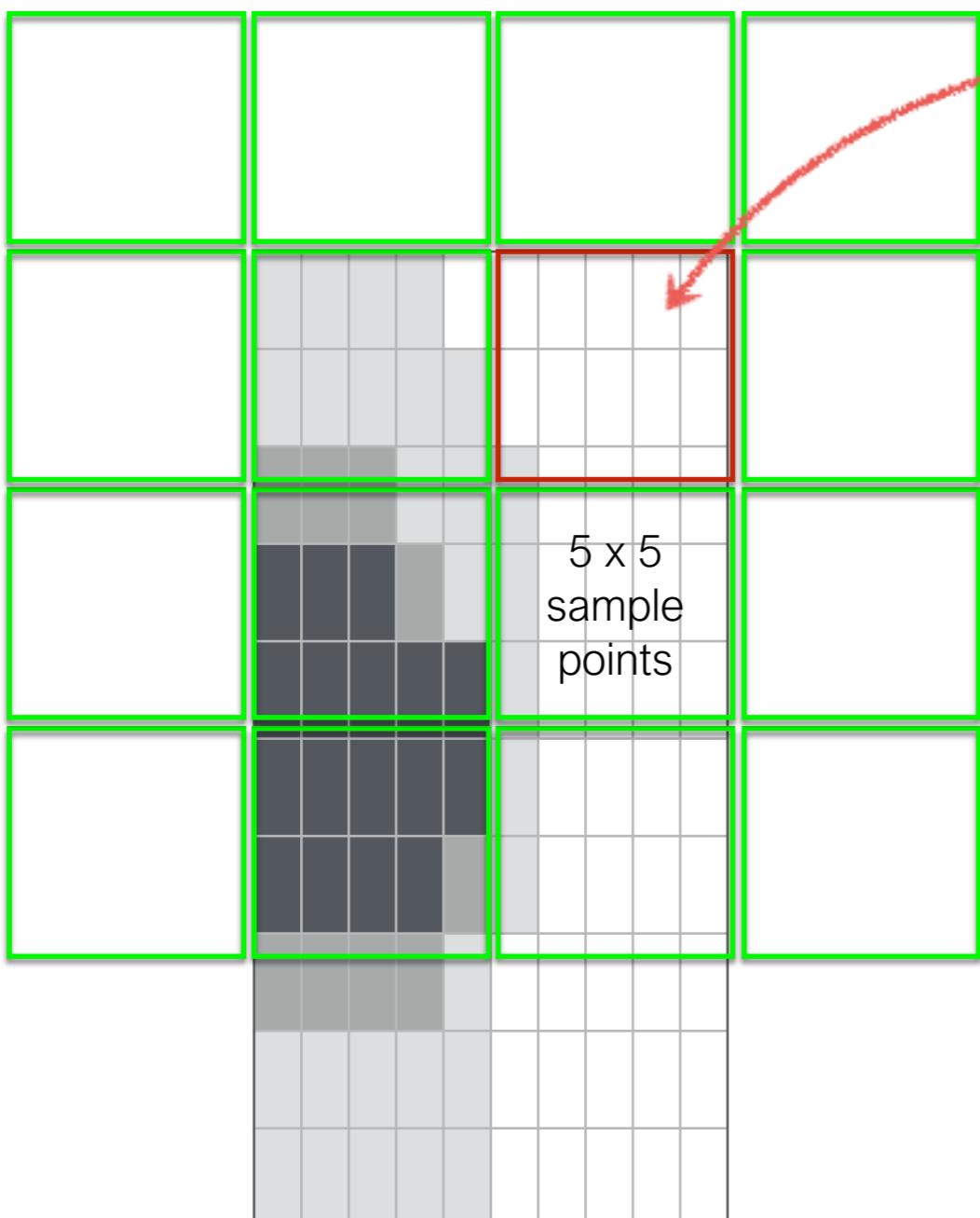
Haar wavelets are just sums over blocks

Use integral images for efficiency (6 operations)

SURF

(‘Speeded’ Up Robust Features)

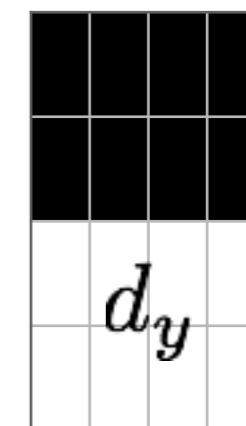
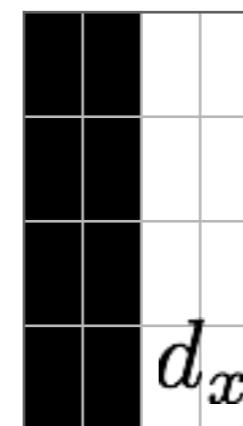
4 x 4 cell grid



Each cell is represented by 4 values:

$$\left[\sum d_x, \sum d_y, \sum |d_x|, \sum |d_y| \right]$$

Haar wavelets filters
(Gaussian weighted from center)

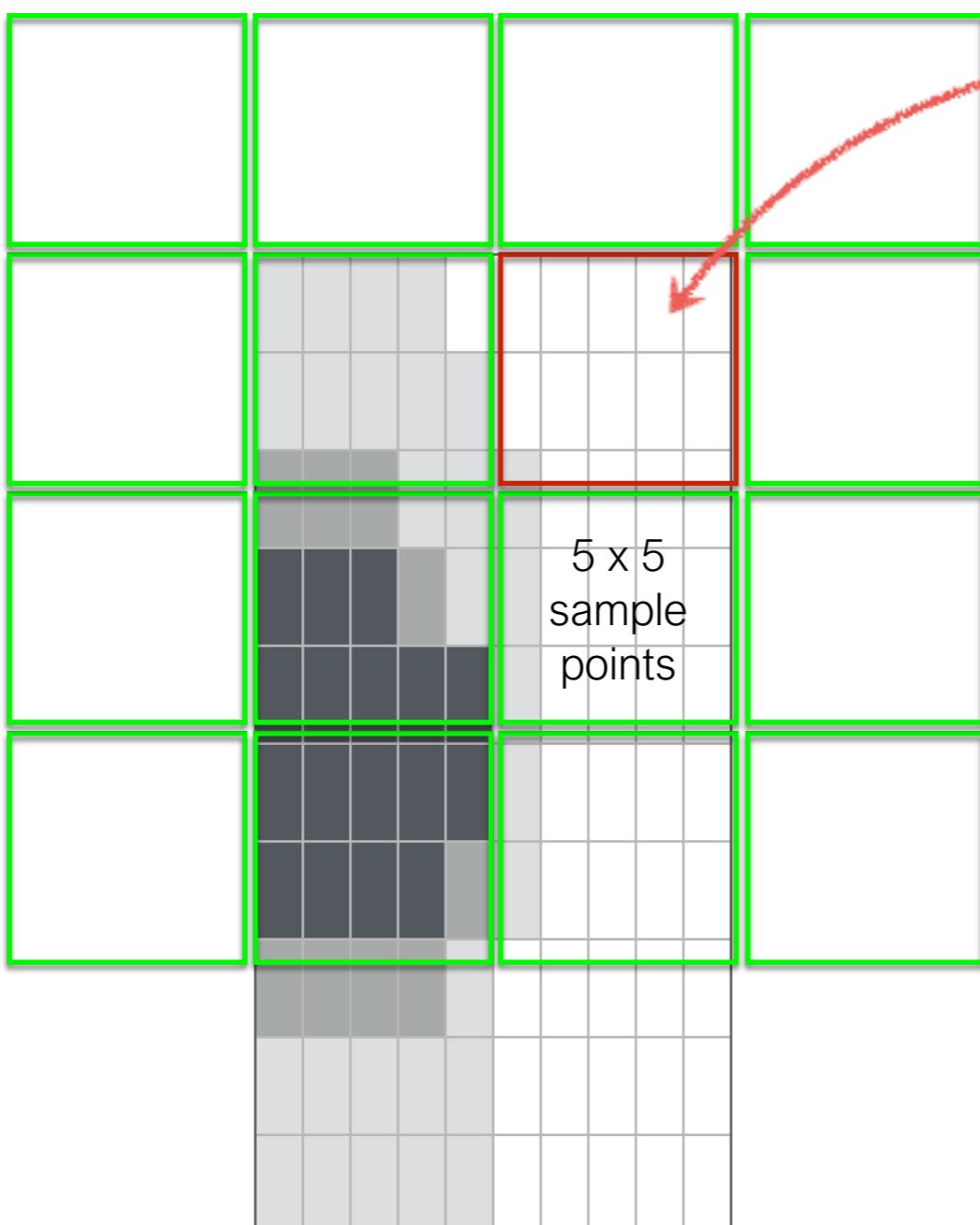


How big is the SURF descriptor?

SURF

(‘Speeded’ Up Robust Features)

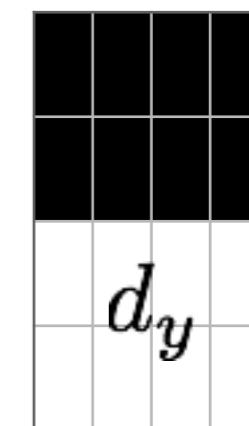
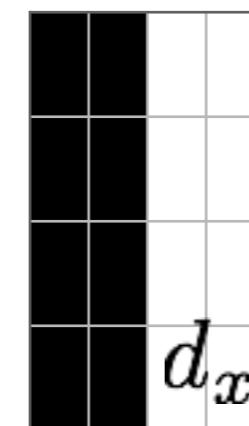
4 x 4 cell grid



Each cell is represented by 4 values:

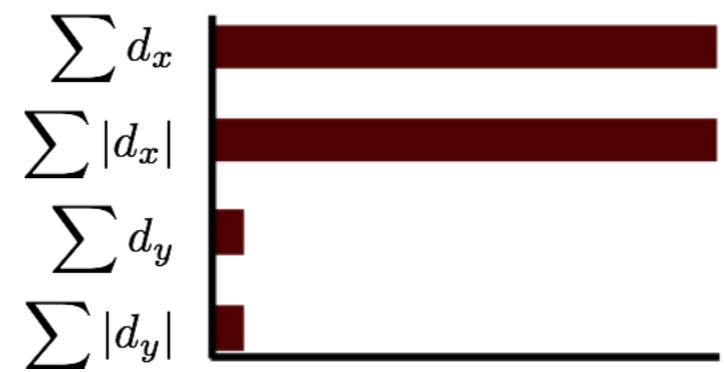
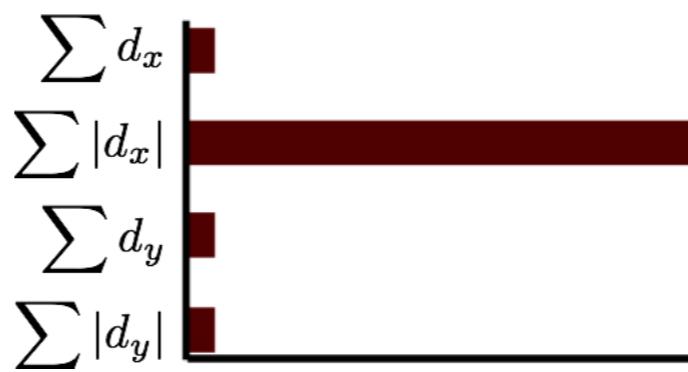
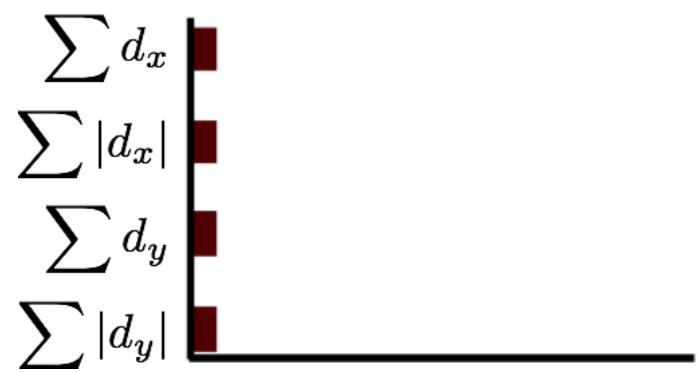
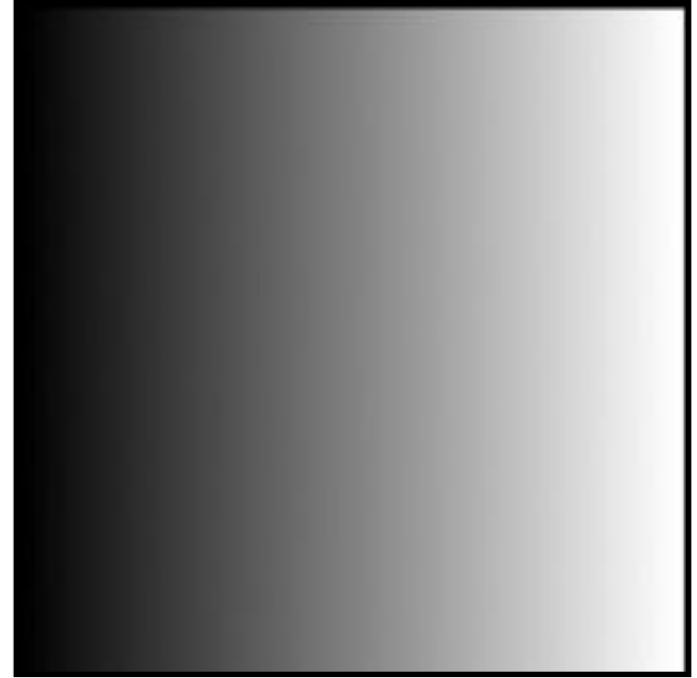
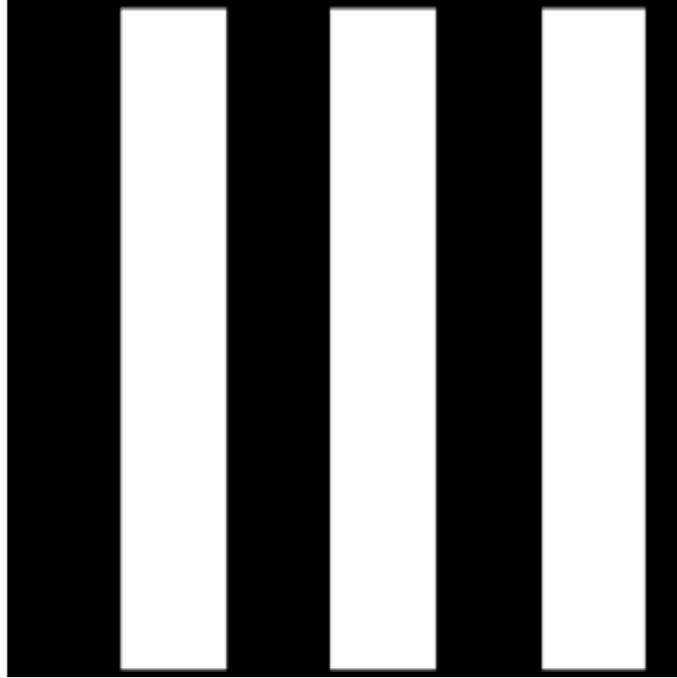
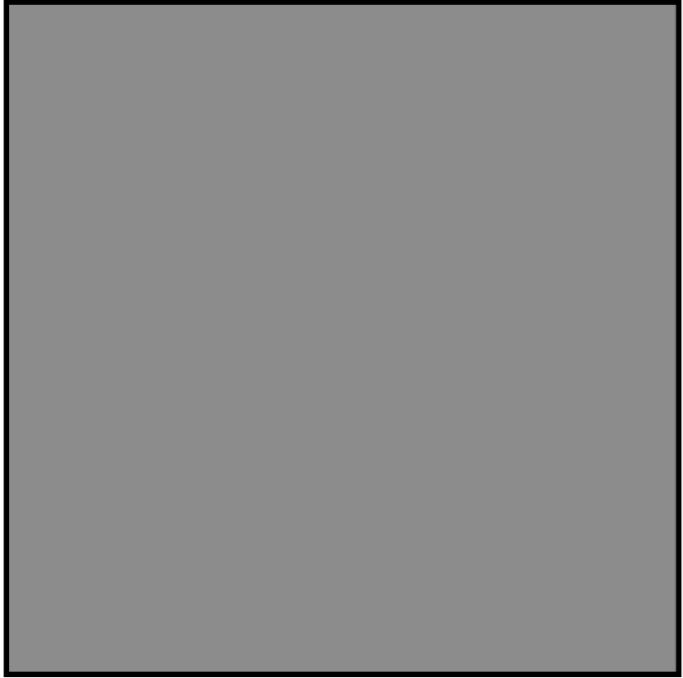
$$[\sum d_x, \sum d_y, \sum |d_x|, \sum |d_y|]$$

Haar wavelets filters
(Gaussian weighted from center)



How big is the SURF descriptor?

64 dimensions



SIFT



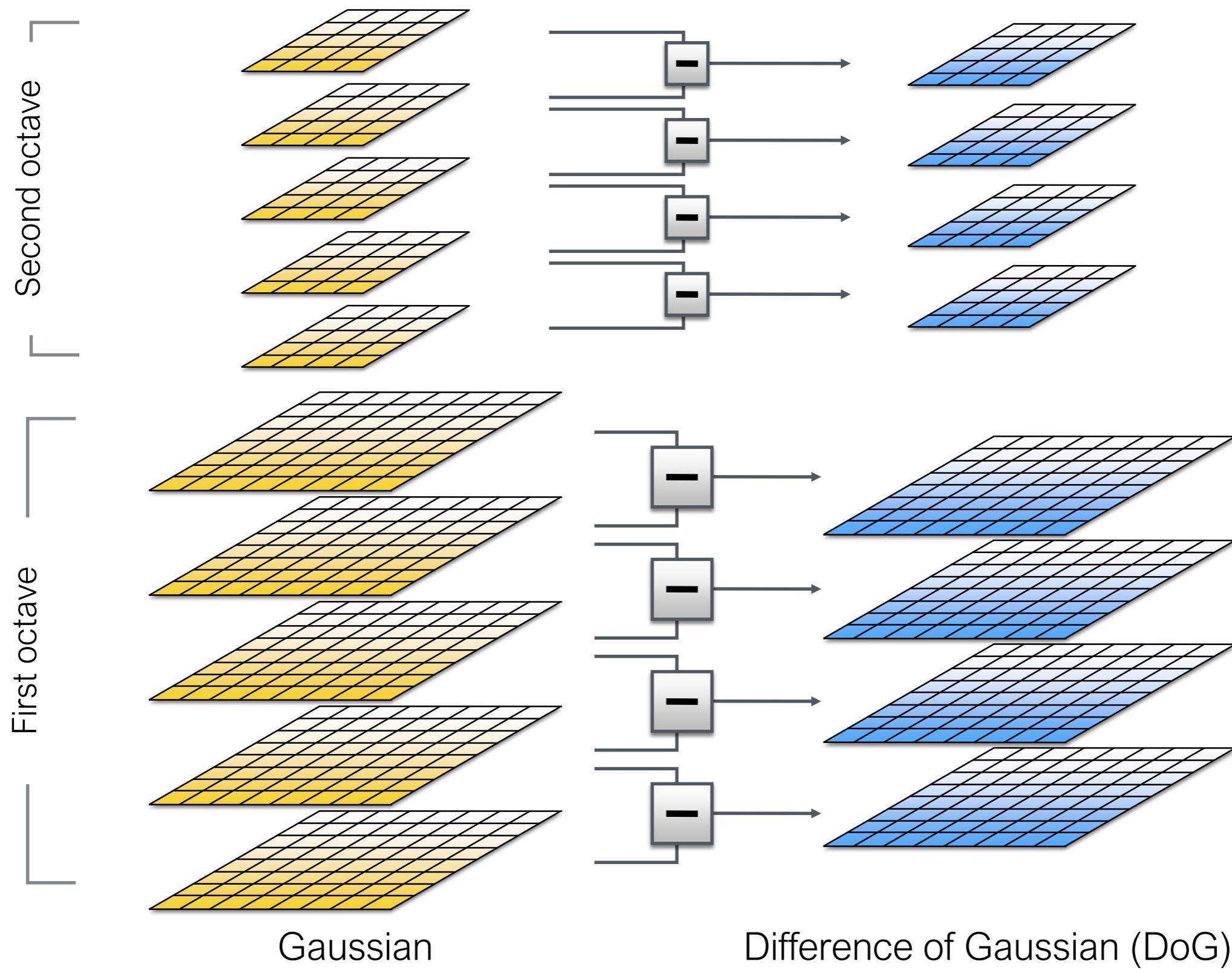
SIFT

(Scale Invariant Feature Transform)

SIFT describes both a **detector** and **descriptor**

1. Multi-scale extrema detection
2. Keypoint localization
3. Orientation assignment
4. Keypoint descriptor

1. Multi-scale extrema detection



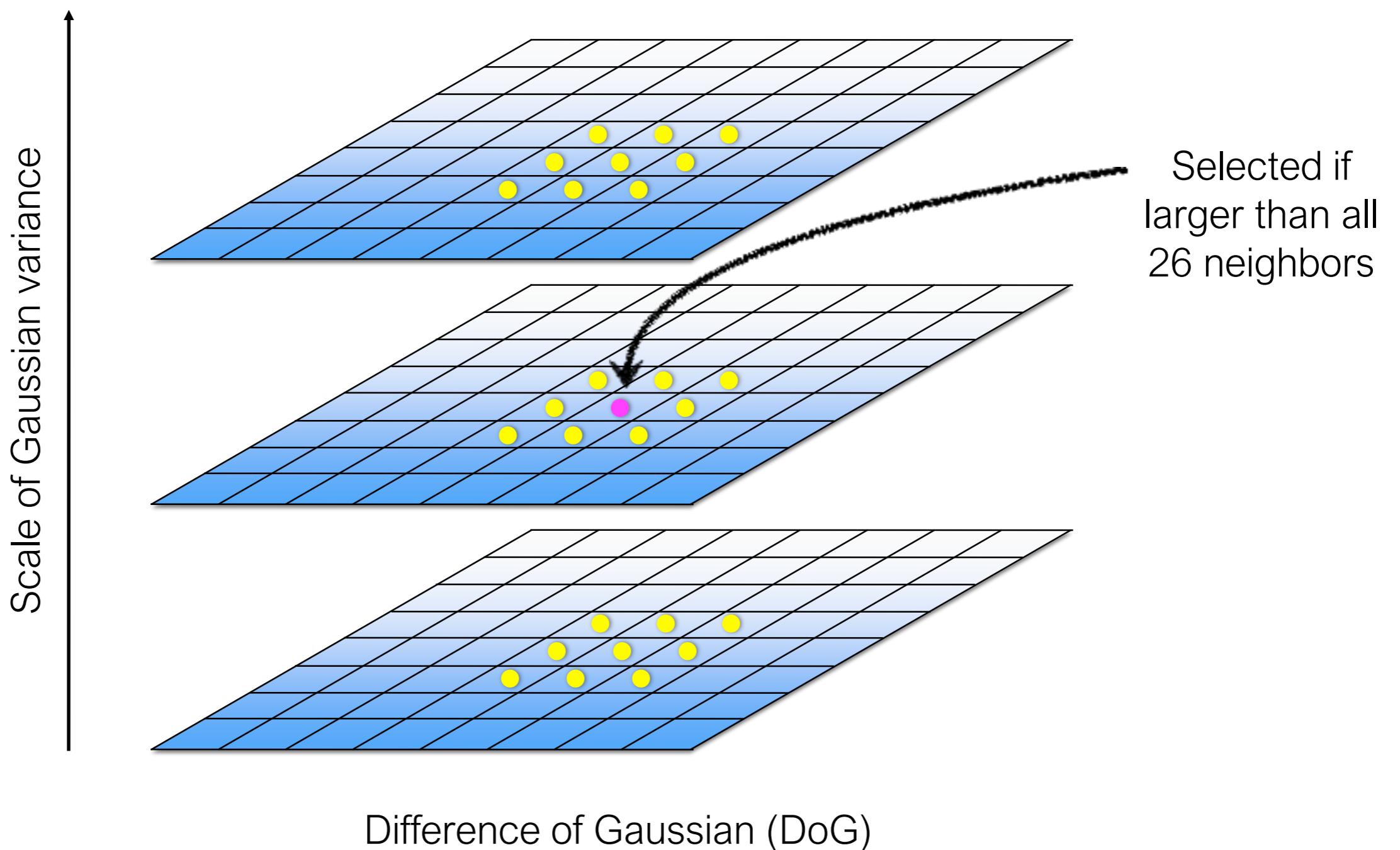


Gaussian



Laplacian

Scale-space extrema



2. Keypoint localization

2nd order Taylor series approximation of DoG scale-space

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

$$\mathbf{x} = \{x, y, \sigma\}$$

Take the derivative and solve for extrema

$$\mathbf{x}_m = - \frac{\partial^2 f}{\partial \mathbf{x}^2}^{-1} \frac{\partial f}{\partial \mathbf{x}}$$

Additional tests to retain only strong features

3. Orientation assignment

For a keypoint, \mathbf{L} is the **Gaussian-smoothed** image with the closest scale,

$$m(x, y) = \sqrt{(L(x + 1, y) - L(x - 1, y))^2 + (L(x, y + 1) - L(x, y - 1))^2}$$

x-derivative y-derivative

$$\theta(x, y) = \tan^{-1}((L(x, y + 1) - L(x, y - 1)) / (L(x + 1, y) - L(x - 1, y)))$$

Detection process returns

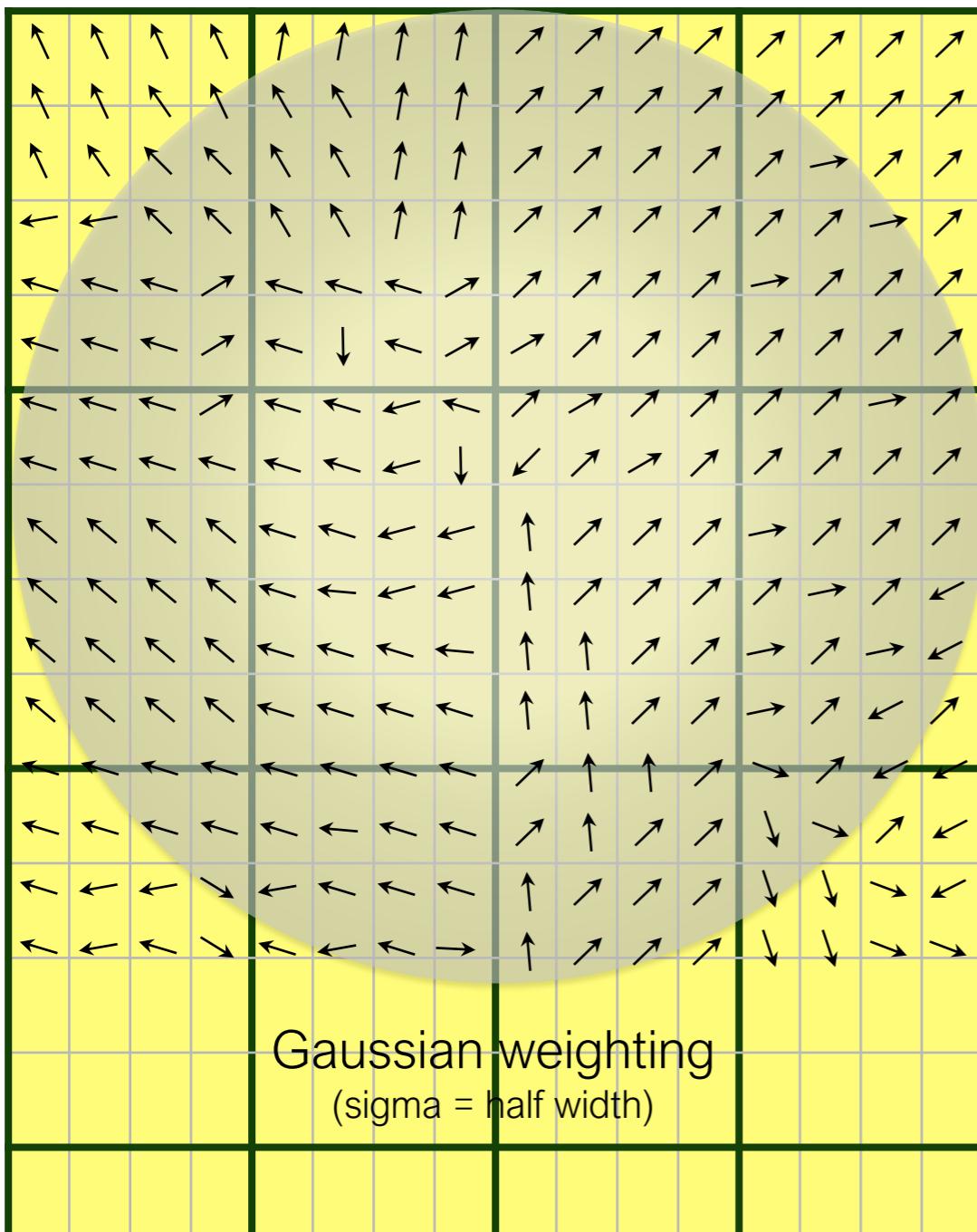
$$\{x, y, \sigma, \theta\}$$

location scale orientation

4. Keypoint descriptor

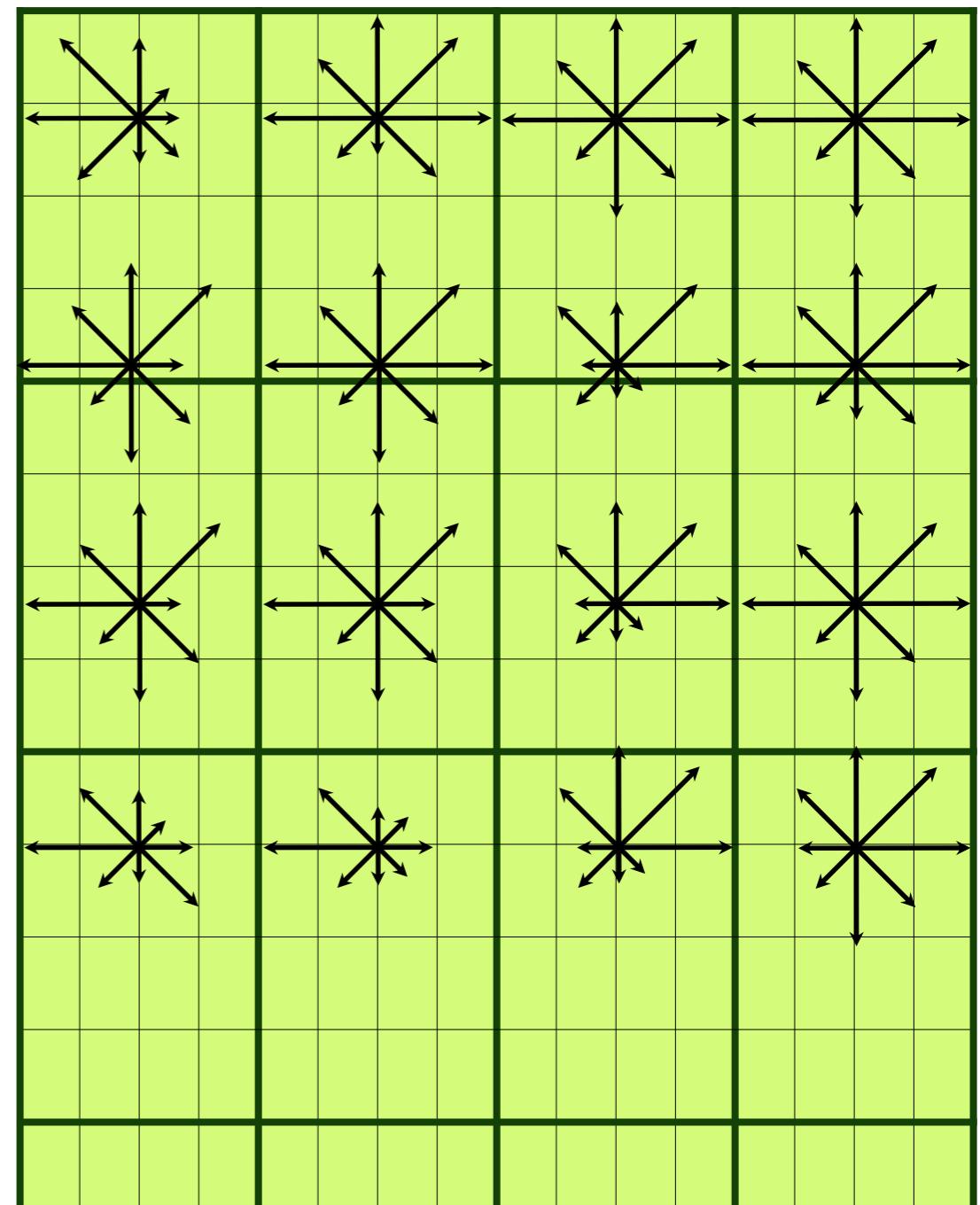
Image Gradients

(4 x 4 pixel per cell, 4 x 4 cells)

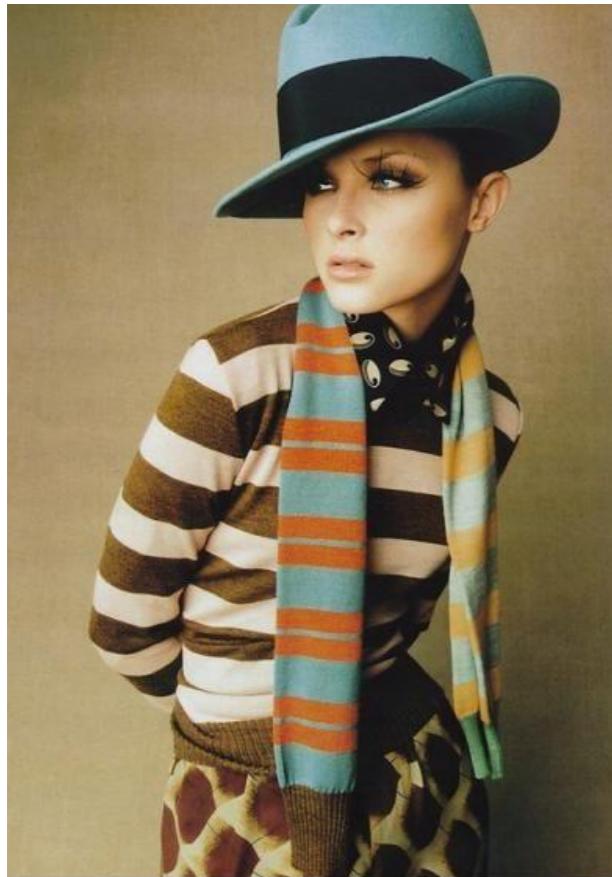


SIFT descriptor

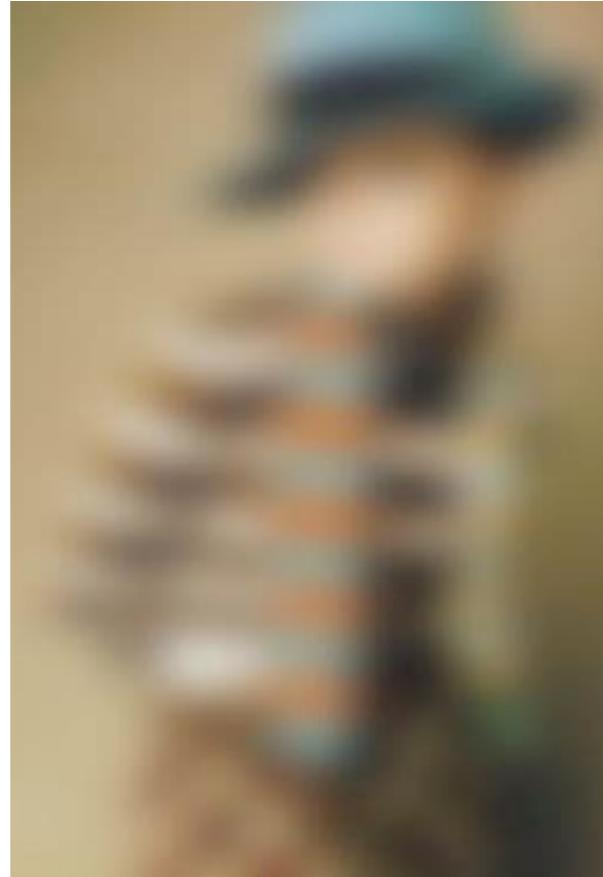
(16 cells x 8 directions = 128 dims)



Discriminative power



Raw pixels



Sampled



Locally orderless



Global histogram

Generalization power



References

Basic reading:

- Szeliski textbook, Sections 4.1.2, 14.1.2.