

Visual Recognition: Instance Level Recognition

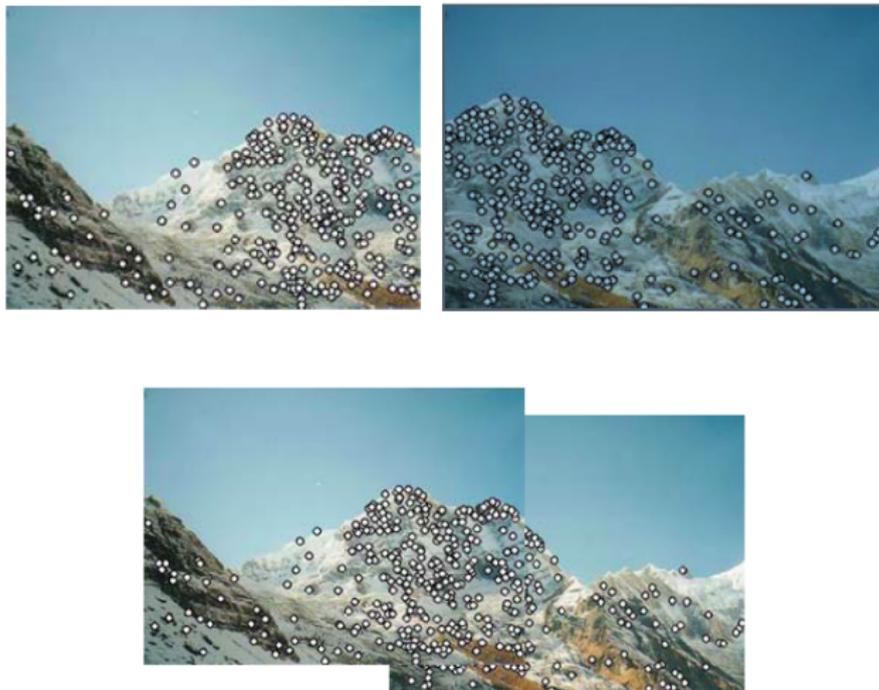
Raquel Urtasun

TTI Chicago

Jan 19, 2012

Local features for **instance-level** recognition

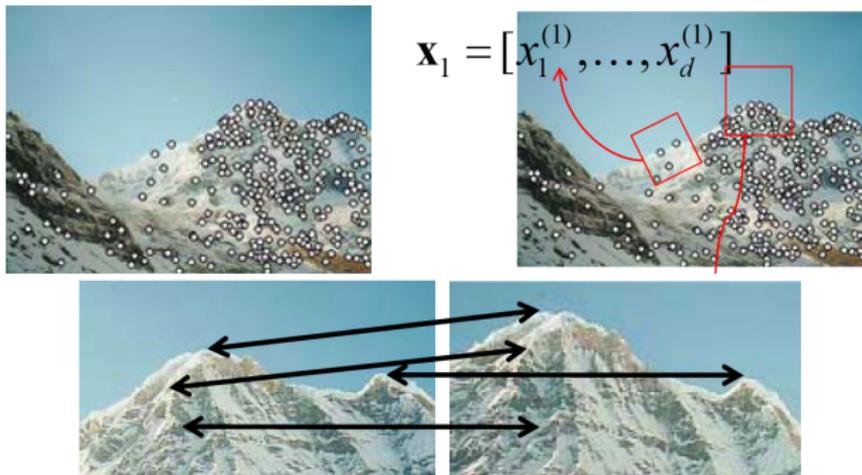
Application Example: Image stitching



[Source: K. Grauman]

Local features

- **Detection:** Identify the interest points.
- **Description:** Extract vector feature descriptor around each interest point.
- **Matching:** Determine correspondence between descriptors in two views.
- **Tracking:** alternative to matching that only searches a small neighborhood around each detected feature.



[Source: K. Grauman]

Goal: interest operator repeatability

- We want to detect (at least some of) the same points in both images.
- We have to be able to run the detection procedure independently per image.

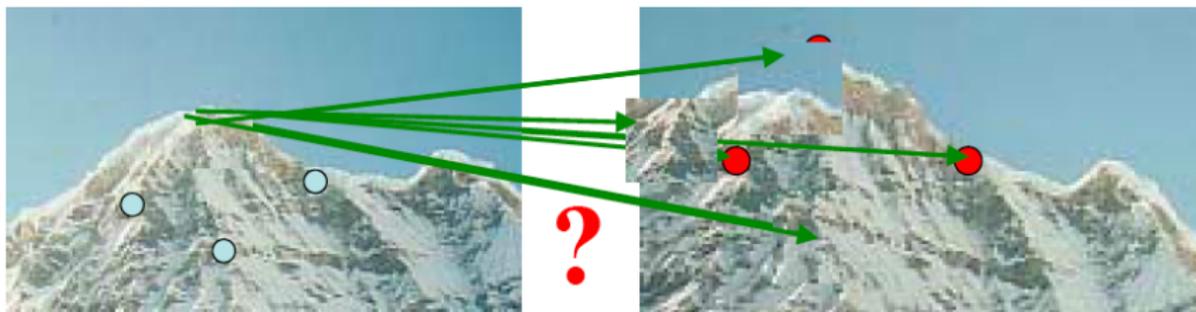


Figure: No chance to find the true matches

[Source: K. Grauman]

Goal: descriptor distinctiveness

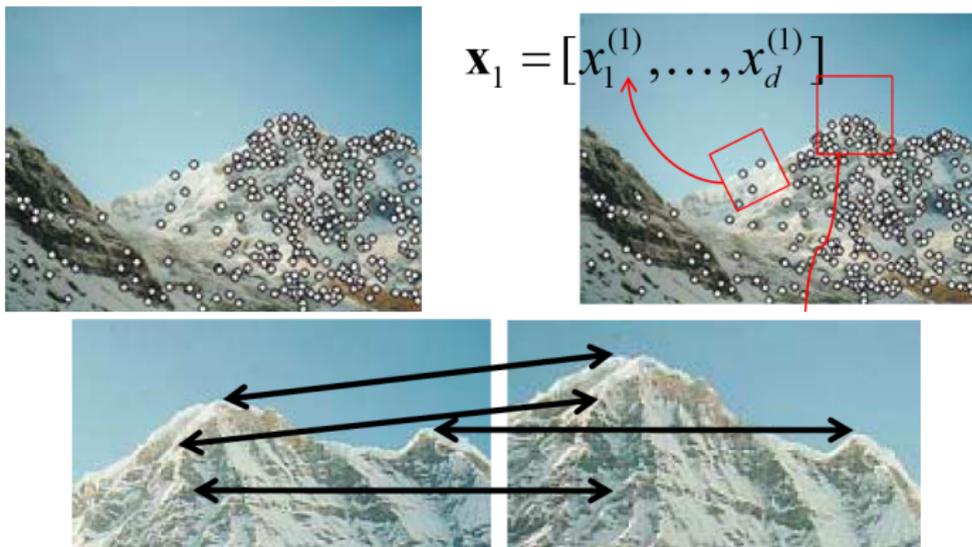
- We want to be able to reliably determine which point goes with which.
- Must provide some invariance to geometric and photometric differences between the two views.



[Source: K. Grauman]

Local features

- **Detection:** Identify the interest points.
- **Description:** Extract vector feature descriptor around each interest point.
- **Matching:** Determine correspondence between descriptors in two views.



[Source: K. Grauman]

What points to choose?



[Source: K. Grauman]

What points to choose?

- Textureless patches are nearly impossible to localize.
- Patches with large contrast changes (gradients) are easier to localize.

What points to choose?

- Textureless patches are nearly impossible to localize.
- Patches with large contrast changes (gradients) are easier to localize.
- But straight line segments at a single orientation suffer from the aperture problem, i.e., it is only possible to align the patches along the direction normal to the edge direction.

What points to choose?

- Textureless patches are nearly impossible to localize.
- Patches with large contrast changes (gradients) are easier to localize.
- But straight line segments at a single orientation suffer from the aperture problem, i.e., it is only possible to align the patches along the direction normal to the edge direction.
- Gradients in at least two (significantly) different orientations are the easiest, e.g., corners.

What points to choose?

- Textureless patches are nearly impossible to localize.
- Patches with large contrast changes (gradients) are easier to localize.
- But straight line segments at a single orientation suffer from the aperture problem, i.e., it is only possible to align the patches along the direction normal to the edge direction.
- Gradients in at least two (significantly) different orientations are the easiest, e.g., corners.

Corners as distinctive interest points

- We should easily recognize the point by looking through a small window.
- Shifting a window in any direction should give a large change in intensity.

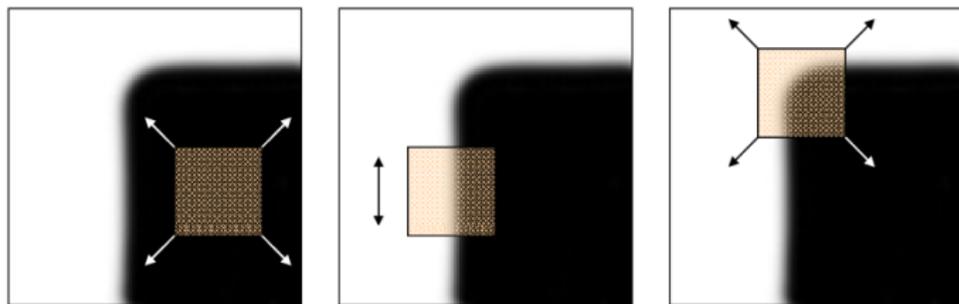


Figure: (left) flat region: no change in all directions, (center) edge: no change along the edge direction, (right) corner: significant change in all directions

[Source: Alyosha Efros, Darya Frolova, Denis Simakov]

A Simple Matching Criteria

- Compare two image patches using (weighted) summed square difference

$$E_{WSSD}(\mathbf{u}) = \sum_i w(\mathbf{p}_i)[I_1(\mathbf{p}_i + \mathbf{u}) - I_0(\mathbf{p}_i)]^2$$

with I_0 and I_1 two images being compared, $\mathbf{u}(u_x, u_y)$ a displacement vector, $w(\mathbf{p})$ a spatially varying weighting function, and the summation i is over all the pixels in the patch.

- We do not know which other image locations the feature will end up being matched against.

A Simple Matching Criteria

- Compare two image patches using (weighted) summed square difference

$$E_{WSSD}(\mathbf{u}) = \sum_i w(\mathbf{p}_i)[I_1(\mathbf{p}_i + \mathbf{u}) - I_0(\mathbf{p}_i)]^2$$

with I_0 and I_1 two images being compared, $\mathbf{u}(u_x, u_y)$ a displacement vector, $w(\mathbf{p})$ a spatially varying weighting function, and the summation i is over all the pixels in the patch.

- We do not know which other image locations the feature will end up being matched against.
- We can only compute how stable this metric is with respect to small variations in position u by comparing an image patch against itself.

A Simple Matching Criteria

- Compare two image patches using (weighted) summed square difference

$$E_{WSSD}(\mathbf{u}) = \sum_i w(\mathbf{p}_i)[I_1(\mathbf{p}_i + \mathbf{u}) - I_0(\mathbf{p}_i)]^2$$

with I_0 and I_1 two images being compared, $\mathbf{u}(u_x, u_y)$ a displacement vector, $w(\mathbf{p})$ a spatially varying weighting function, and the summation i is over all the pixels in the patch.

- We do not know which other image locations the feature will end up being matched against.
- We can only compute how stable this metric is with respect to small variations in position u by comparing an image patch against itself.
- This is the **auto-correlation function**

$$E_{AC}(\Delta\mathbf{u}) = \sum_i w(\mathbf{p}_i)[I_0(\mathbf{p}_i + \Delta\mathbf{u}) - I_0(\mathbf{p}_i)]^2$$

A Simple Matching Criteria

- Compare two image patches using (weighted) summed square difference

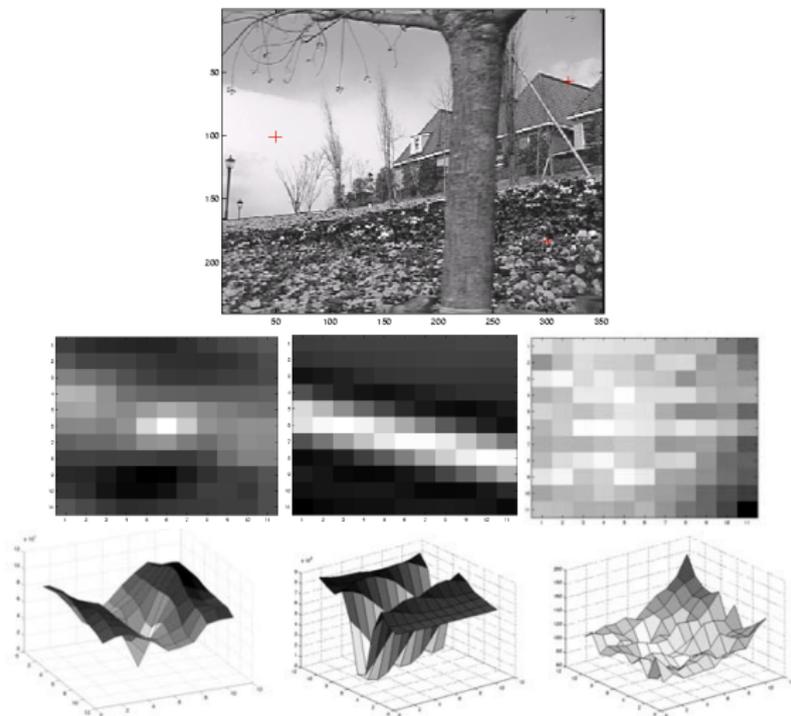
$$E_{WSSD}(\mathbf{u}) = \sum_i w(\mathbf{p}_i)[I_1(\mathbf{p}_i + \mathbf{u}) - I_0(\mathbf{p}_i)]^2$$

with I_0 and I_1 two images being compared, $\mathbf{u}(u_x, u_y)$ a displacement vector, $w(\mathbf{p})$ a spatially varying weighting function, and the summation i is over all the pixels in the patch.

- We do not know which other image locations the feature will end up being matched against.
- We can only compute how stable this metric is with respect to small variations in position u by comparing an image patch against itself.
- This is the **auto-correlation function**

$$E_{AC}(\Delta\mathbf{u}) = \sum_i w(\mathbf{p}_i)[I_0(\mathbf{p}_i + \Delta\mathbf{u}) - I_0(\mathbf{p}_i)]^2$$

Which one is better?



[Source: R. Szeliski]

How to select?

- Using a Taylor Series expansion $l_0(\mathbf{p}_i + \Delta\mathbf{u}) \approx l_0(\mathbf{p}_i) + \nabla l_0(\mathbf{p}_i)$ we can approximate the autocorrelation as

$$\begin{aligned} E_{AC}(\Delta\mathbf{u}) &= \sum_i w(\mathbf{p}_i) [l_0(\mathbf{p}_i + \Delta\mathbf{u}) - l_0(\mathbf{p}_i)]^2 \\ &\approx \sum_i w(\mathbf{p}_i) [l_0(\mathbf{p}_i) + \nabla l_0(\mathbf{p}_i)\Delta\mathbf{u} - l_0(\mathbf{p}_i)]^2 \\ &= \sum_i w(\mathbf{p}_i) [\nabla l_0(\mathbf{p}_i)\Delta\mathbf{u}]^2 \\ &= \Delta\mathbf{u}^T \mathbf{A} \Delta\mathbf{u} \end{aligned}$$

with

$$\nabla l_0(\mathbf{p}_i) = \left(\frac{\partial l_0}{\partial x}, \frac{\partial l_0}{\partial y} \right) (\mathbf{p}_i)$$

the image gradient.

- Gradient can be computed with the filtering techniques we saw, e.g., derivatives of Gaussians.

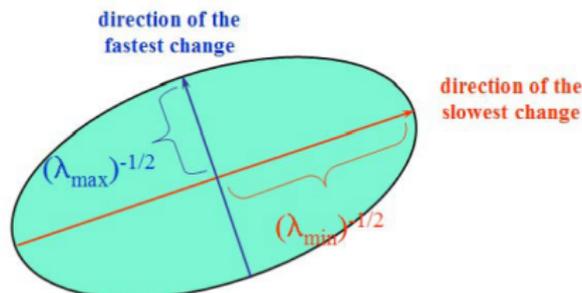
More on selection

- The autocorrelation is $E_{AC}(\Delta \mathbf{u}) = \Delta \mathbf{u}^T \mathbf{A} \Delta \mathbf{u}$, with

$$\mathbf{A} = \sum_u \sum_v w(u, v) \begin{bmatrix} l_x^2 & l_x l_y \\ l_y l_x & l_y^2 \end{bmatrix} = w * \begin{bmatrix} l_x^2 & l_x l_y \\ l_y l_x & l_y^2 \end{bmatrix}$$

where we have replaced the weighted summations with discrete convolutions with the weighting kernel w .

- \mathbf{A} can be interpreted as a tensor where the outer products of the gradients are convolved with a weighting function.
- Eigenvalues a notion of uncertainty



[Source: R. Szeliski]

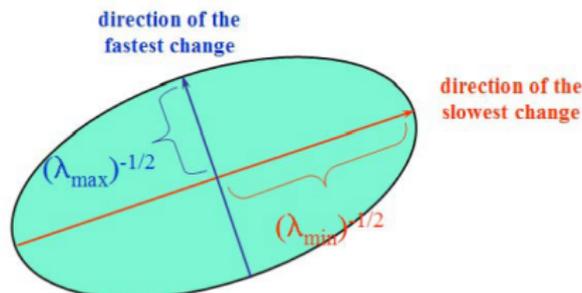
More on selection

- The autocorrelation is $E_{AC}(\Delta \mathbf{u}) = \Delta \mathbf{u}^T \mathbf{A} \Delta \mathbf{u}$, with

$$\mathbf{A} = \sum_u \sum_v w(u, v) \begin{bmatrix} l_x^2 & l_x l_y \\ l_y l_x & l_y^2 \end{bmatrix} = w * \begin{bmatrix} l_x^2 & l_x l_y \\ l_y l_x & l_y^2 \end{bmatrix}$$

where we have replaced the weighted summations with discrete convolutions with the weighting kernel w .

- \mathbf{A} can be interpreted as a tensor where the outer products of the gradients are convolved with a weighting function.
- Eigenvalues a notion of uncertainty



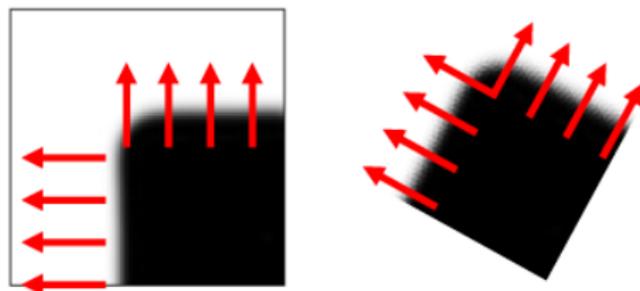
[Source: R. Szeliski]

Eigenvalues a notion of uncertainty

- \mathbf{A} is symmetric

$$\mathbf{A} = \mathbf{U} \begin{bmatrix} \lambda_0 & 0 \\ 0 & \lambda_1 \end{bmatrix} \mathbf{U}^T \quad \text{with} \quad \mathbf{A}\mathbf{u}_j = \lambda_j\mathbf{u}_j$$

- The eigenvalues of \mathbf{A} reveal the amount of intensity change in the two principal orthogonal gradient directions in the window.
- How is this matrix for



[Source: R. Szeliski]

Local Feature Selection Criteria

- Shi and Tomasi, 94 proposed the smallest eigenvalue of \mathbf{A} , i.e., $\lambda_0^{-1/2}$.
- Harris and Stephens, 88 is rotationally invariant and downweights edge-like features where $\lambda_1 \gg \lambda_0$

$$\det(\mathbf{A}) - \alpha \text{trace}(\mathbf{A})^2 = \lambda_0 \lambda_1 - \alpha (\lambda_0 + \lambda_1)^2$$

Local Feature Selection Criteria

- Shi and Tomasi, 94 proposed the smallest eigenvalue of \mathbf{A} , i.e., $\lambda_0^{-1/2}$.
- Harris and Stephens, 88 is rotationally invariant and downweights edge-like features where $\lambda_1 \gg \lambda_0$

$$\det(\mathbf{A}) - \alpha \text{trace}(\mathbf{A})^2 = \lambda_0 \lambda_1 - \alpha (\lambda_0 + \lambda_1)^2$$

- Triggs, 04 suggested

$$\lambda_0 - \alpha \lambda_1$$

also reduces the response at 1D edges, where aliasing errors sometimes inflate the smaller eigenvalue.

Local Feature Selection Criteria

- Shi and Tomasi, 94 proposed the smallest eigenvalue of \mathbf{A} , i.e., $\lambda_0^{-1/2}$.
- Harris and Stephens, 88 is rotationally invariant and downweights edge-like features where $\lambda_1 \gg \lambda_0$

$$\det(\mathbf{A}) - \alpha \text{trace}(\mathbf{A})^2 = \lambda_0 \lambda_1 - \alpha (\lambda_0 + \lambda_1)^2$$

- Triggs, 04 suggested

$$\lambda_0 - \alpha \lambda_1$$

also reduces the response at 1D edges, where aliasing errors sometimes inflate the smaller eigenvalue.

- Brown et al, 05 use the harmonic mean

$$\frac{\det(\mathbf{A})}{\text{trace}(\mathbf{A})} = \frac{\lambda_0 \lambda_1}{\lambda_0 + \lambda_1}$$

Local Feature Selection Criteria

- Shi and Tomasi, 94 proposed the smallest eigenvalue of \mathbf{A} , i.e., $\lambda_0^{-1/2}$.
- Harris and Stephens, 88 is rotationally invariant and downweights edge-like features where $\lambda_1 \gg \lambda_0$

$$\det(\mathbf{A}) - \alpha \text{trace}(\mathbf{A})^2 = \lambda_0 \lambda_1 - \alpha (\lambda_0 + \lambda_1)^2$$

- Triggs, 04 suggested

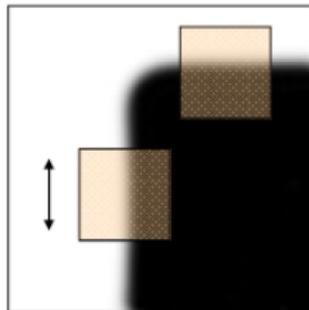
$$\lambda_0 - \alpha \lambda_1$$

also reduces the response at 1D edges, where aliasing errors sometimes inflate the smaller eigenvalue.

- Brown et al, 05 use the harmonic mean

$$\frac{\det(\mathbf{A})}{\text{trace}(\mathbf{A})} = \frac{\lambda_0 \lambda_1}{\lambda_0 + \lambda_1}$$

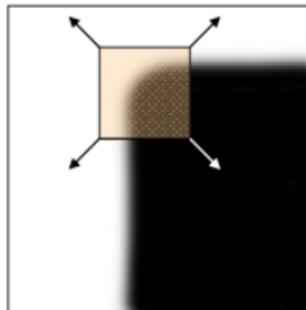
Type of responses



“edge”:

$$\lambda_1 \gg \lambda_2$$

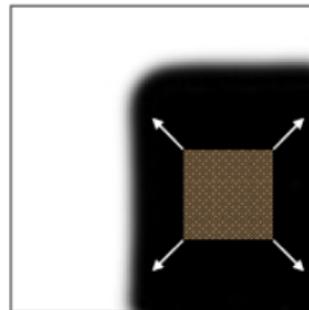
$$\lambda_2 \gg \lambda_1$$



“corner”:

λ_1 and λ_2 are large,

$$\lambda_1 \sim \lambda_2;$$



“flat” region

λ_1 and λ_2 are
small;

[Source: K. Grauman]

Harris Corner detector

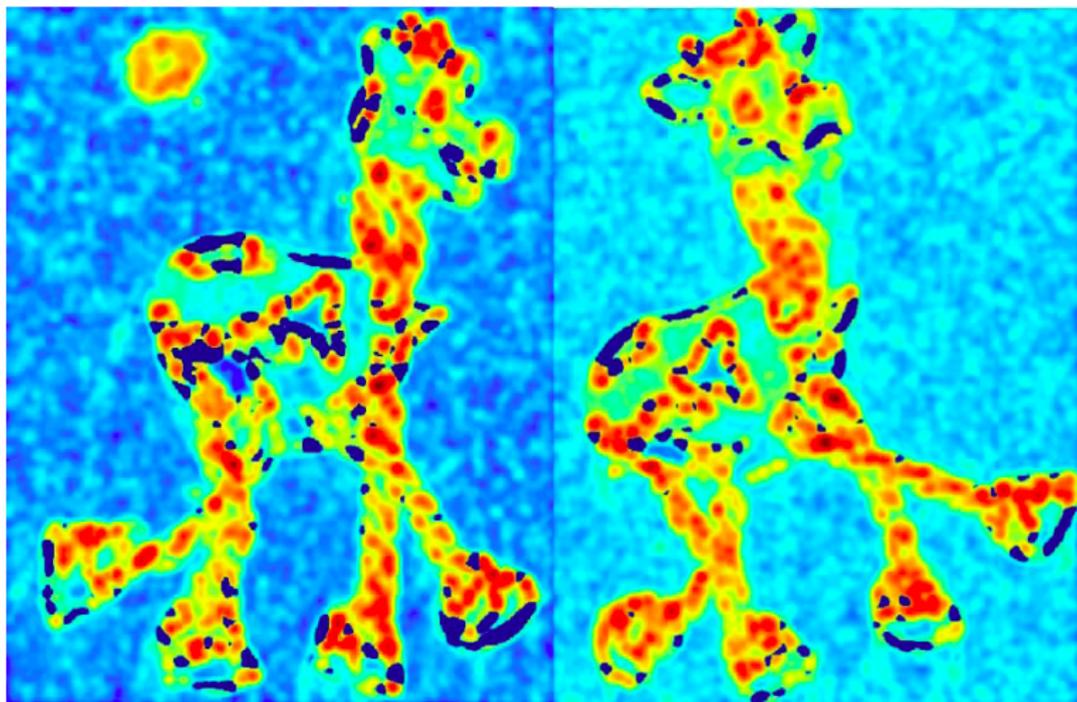
- 1 Compute \mathbf{A} for each image window to get their cornerness scores.
- 2 Find points whose surrounding window gave large corner response ($f >$ threshold).
- 3 Take the points of local maxima, i.e., perform non-maximum suppression.

Example



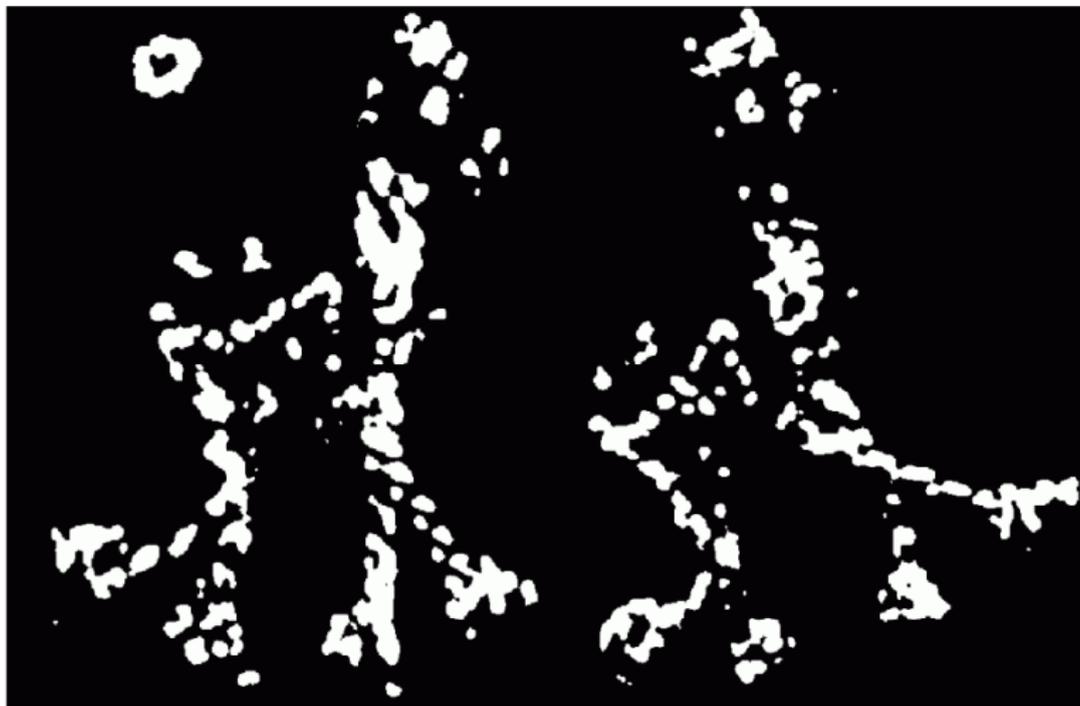
[Source: K. Grauman]

1) Compute Cornerness



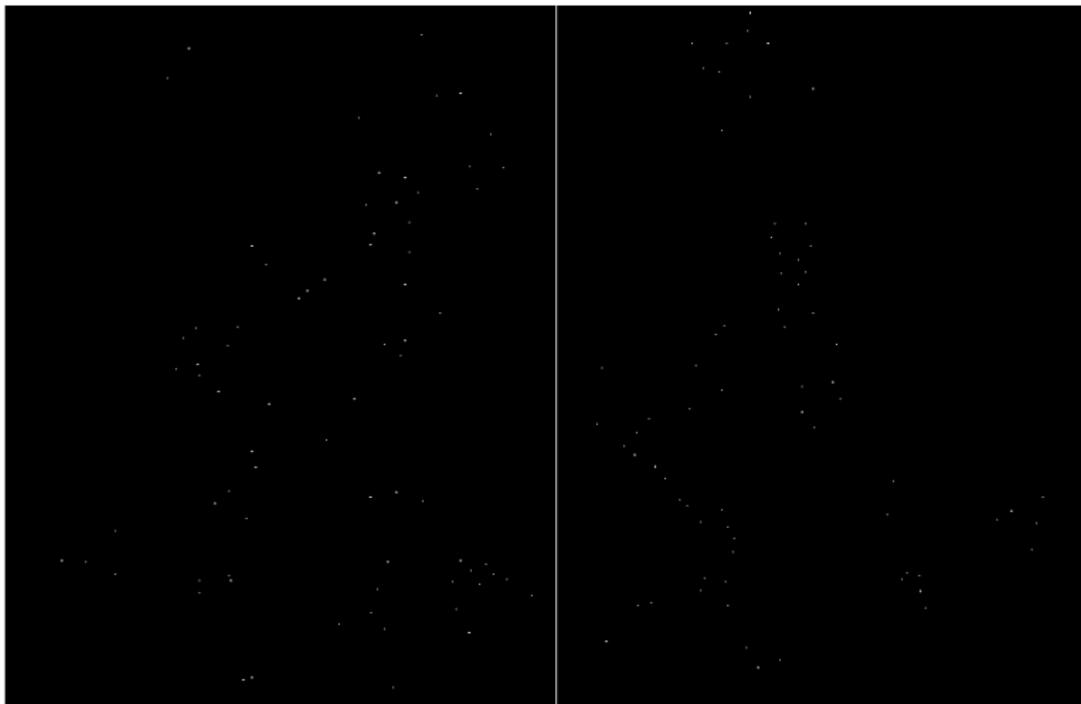
[Source: K. Grauman]

2) Find High Response



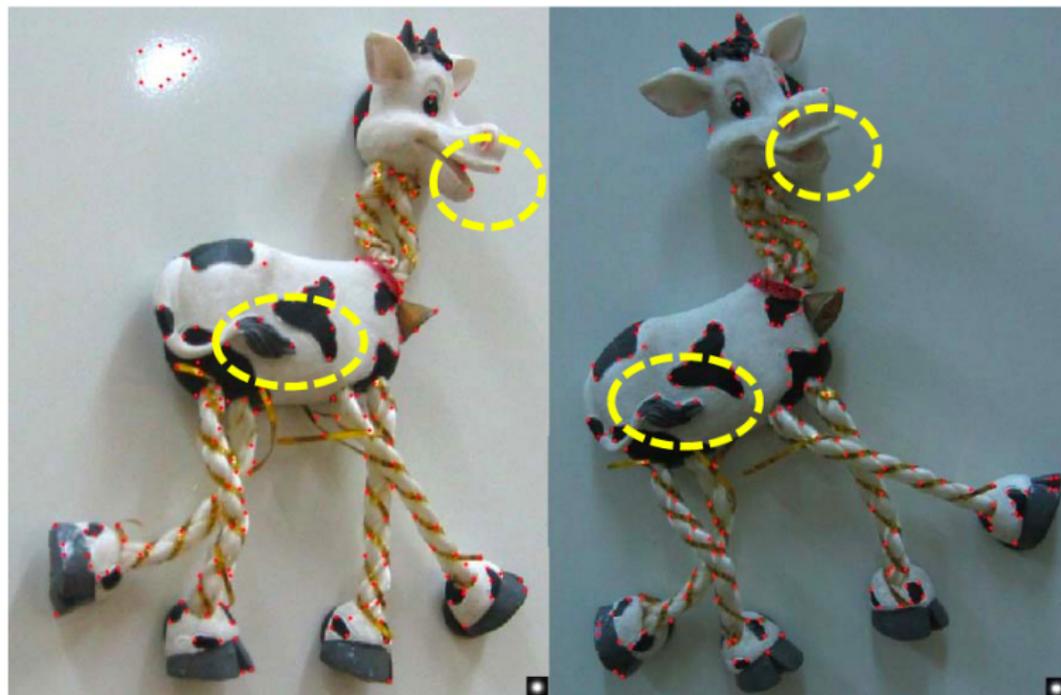
[Source: K. Grauman]

3) Non-maxima Suppression



[Source: K. Grauman]

Results

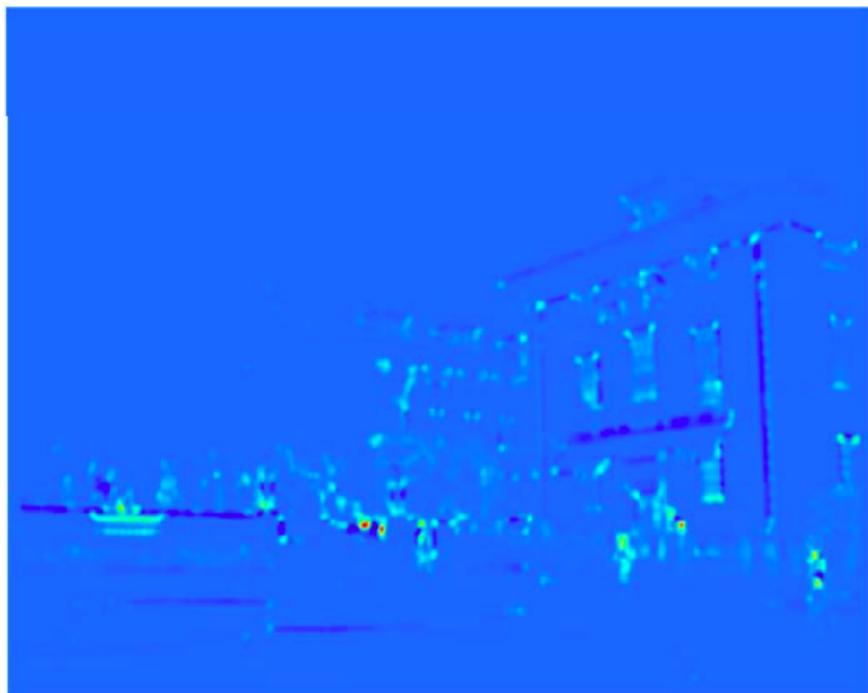


[Source: K. Grauman]

Another Example



[Source: K. Grauman]



[Source: K. Grauman]

Interest Points



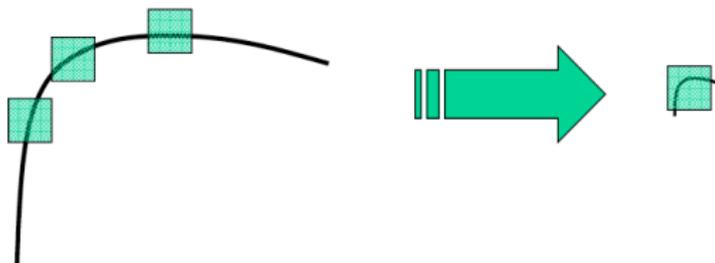
[Source: K. Grauman]

Properties of Harris Corner Detector

- Rotation invariant?

$$\mathbf{A} = w * \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix} = \mathbf{U} \begin{bmatrix} \lambda_0 & 0 \\ 0 & \lambda_1 \end{bmatrix} \mathbf{U}^T \quad \text{with} \quad \mathbf{A} \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

- Scale Invariant?



All points will be classified as **edges**

Corner !

[Source: K. Grauman]

Scale invariant interest points

How can we independently select interest points in each image, such that the detections are repeatable across different scales?

- Extract features at a variety of scales, e.g., by using multiple resolutions in a pyramid, and then matching features at the same level.

Scale invariant interest points

How can we independently select interest points in each image, such that the detections are repeatable across different scales?

- Extract features at a variety of scales, e.g., by using multiple resolutions in a pyramid, and then matching features at the same level.
- When does this work?

Scale invariant interest points

How can we independently select interest points in each image, such that the detections are repeatable across different scales?

- Extract features at a variety of scales, e.g., by using multiple resolutions in a pyramid, and then matching features at the same level.
- When does this work?
- More efficient to extract features that are stable in both location and scale.

Scale invariant interest points

How can we independently select interest points in each image, such that the detections are repeatable across different scales?

- Extract features at a variety of scales, e.g., by using multiple resolutions in a pyramid, and then matching features at the same level.
- When does this work?
- More efficient to extract features that are stable in both location and scale.
- Find scale that gives local maxima of a function f in both position and scale.



$$f(I_{i_1 \dots i_m}(x, \sigma)) = f(I_{i_1 \dots i_m}(x', \sigma'))$$

Scale invariant interest points

How can we independently select interest points in each image, such that the detections are repeatable across different scales?

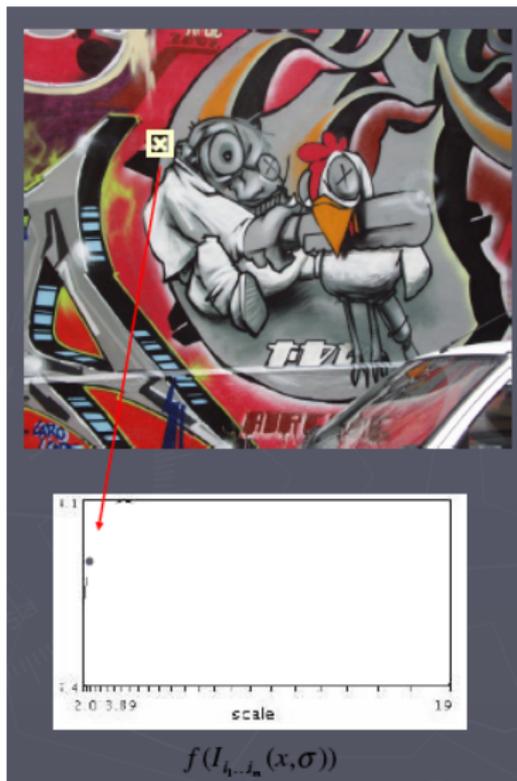
- Extract features at a variety of scales, e.g., by using multiple resolutions in a pyramid, and then matching features at the same level.
- When does this work?
- More efficient to extract features that are stable in both location and scale.
- Find scale that gives local maxima of a function f in both position and scale.



$$f(I_{i_1 \dots i_m}(x, \sigma)) = f(I_{i_1 \dots i_m}(x', \sigma'))$$

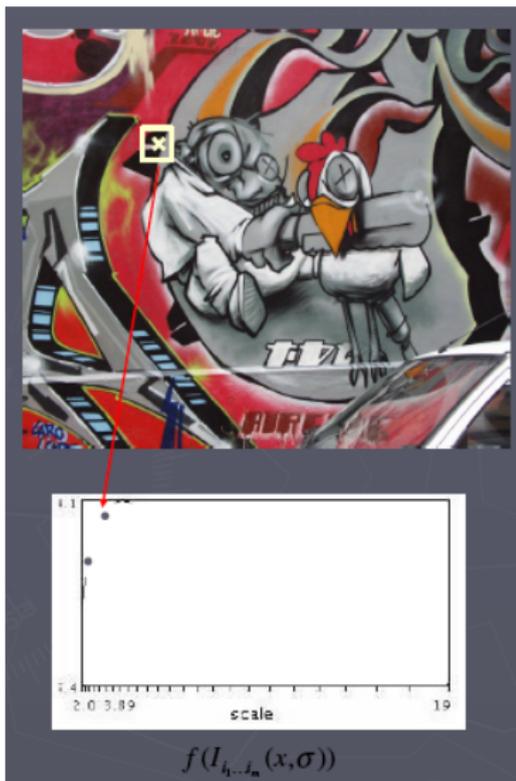
Automatic Scale Selection

Function responses for increasing scale (scale signature).



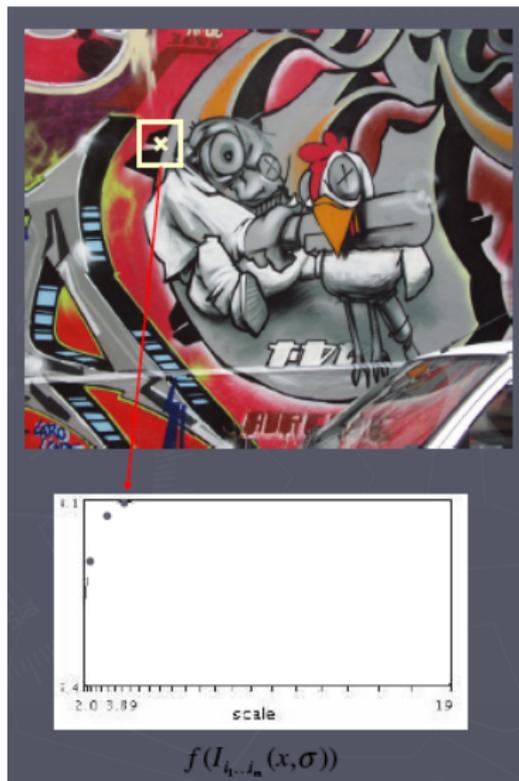
Automatic Scale Selection

Function responses for increasing scale (scale signature).



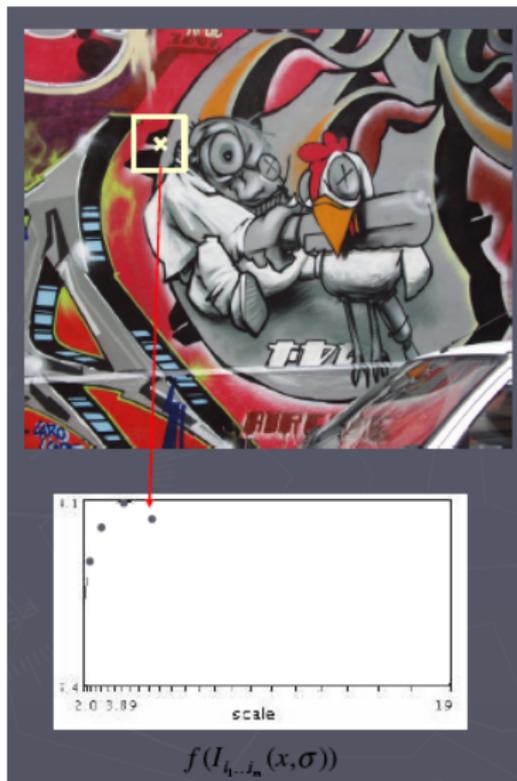
Automatic Scale Selection

Function responses for increasing scale (scale signature).



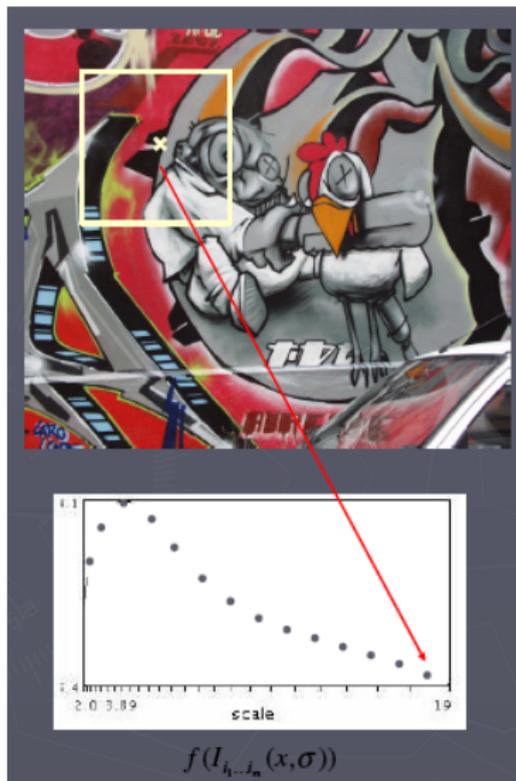
Automatic Scale Selection

Function responses for increasing scale (scale signature).



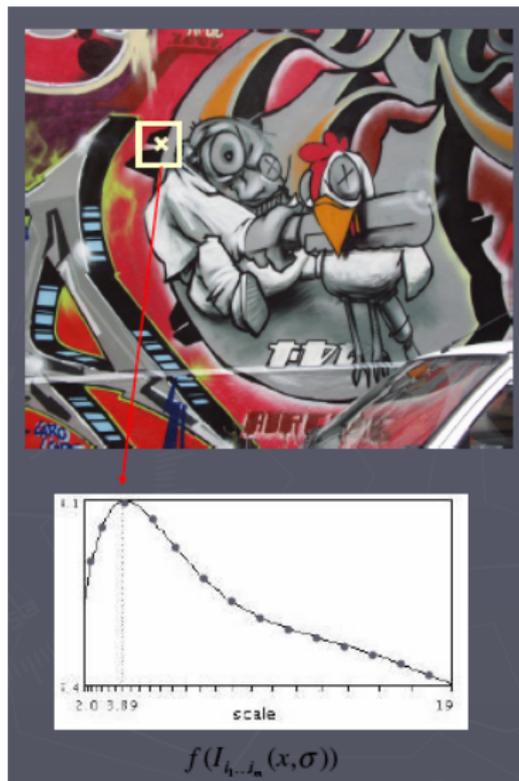
Automatic Scale Selection

Function responses for increasing scale (scale signature).



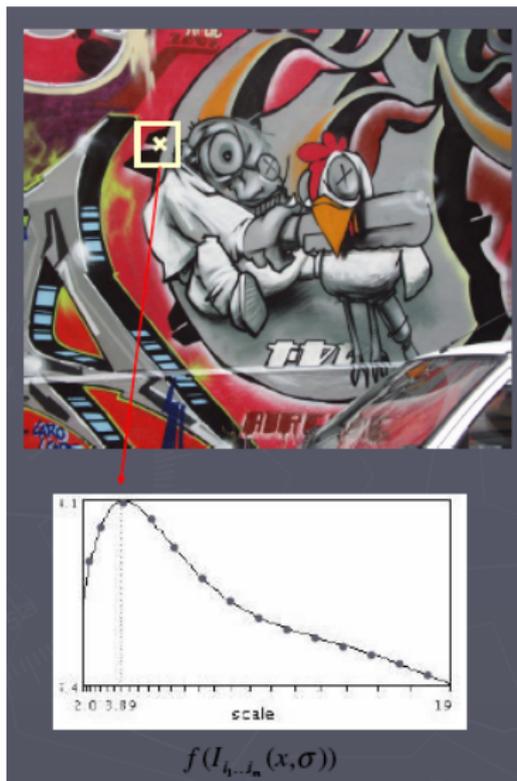
Automatic Scale Selection

Function responses for increasing scale (scale signature).



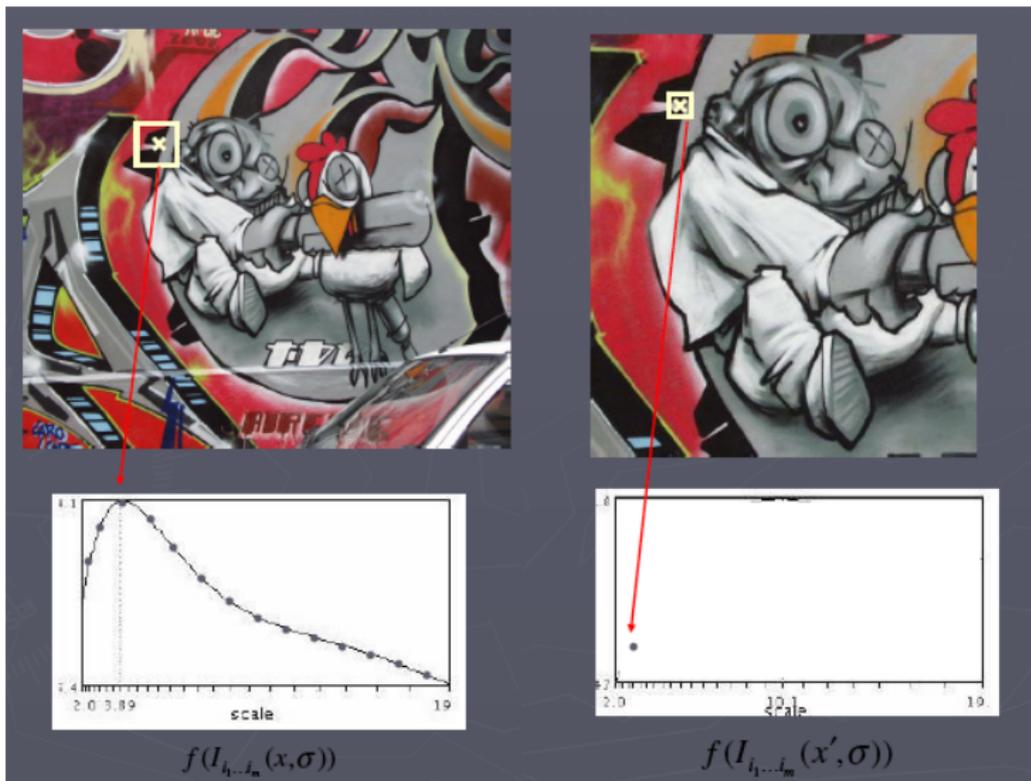
Automatic Scale Selection

Function responses for increasing scale (scale signature).



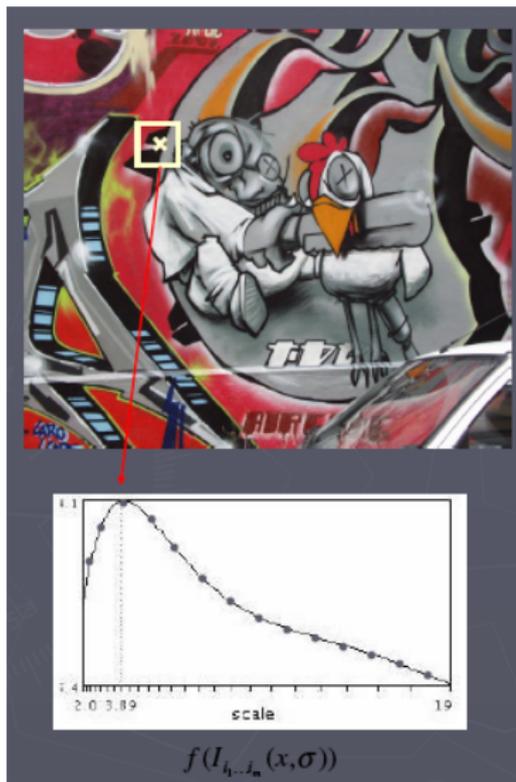
Automatic Scale Selection

Function responses for increasing scale (scale signature).



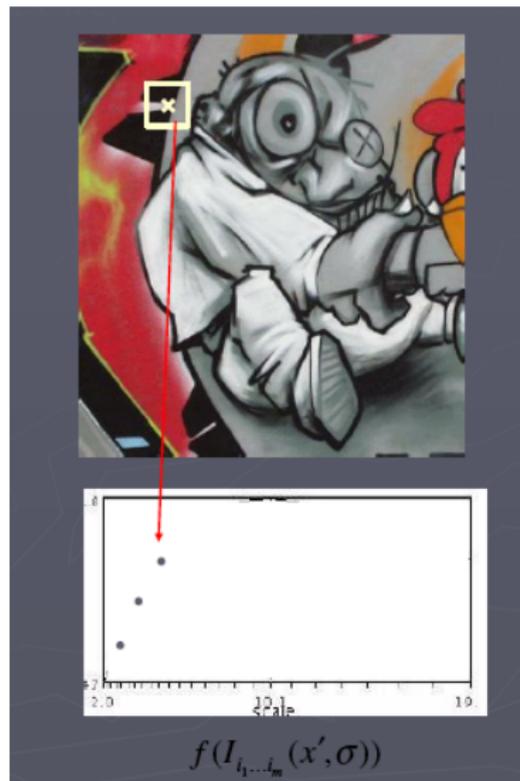
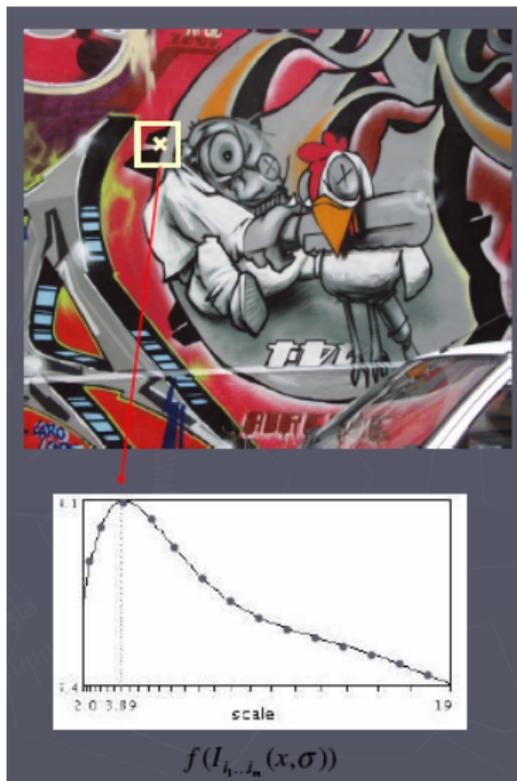
Automatic Scale Selection

Function responses for increasing scale (scale signature).



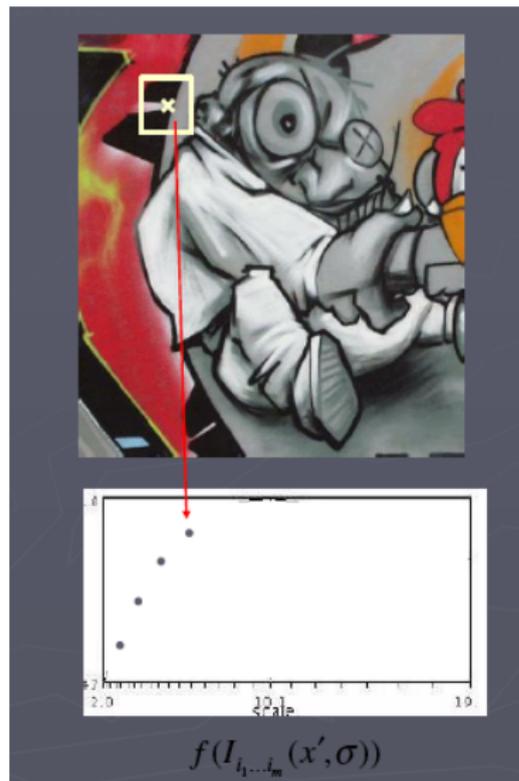
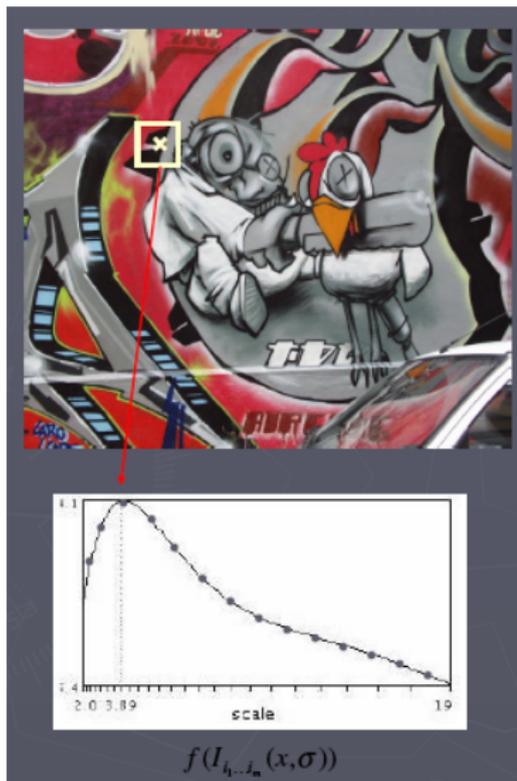
Automatic Scale Selection

Function responses for increasing scale (scale signature).



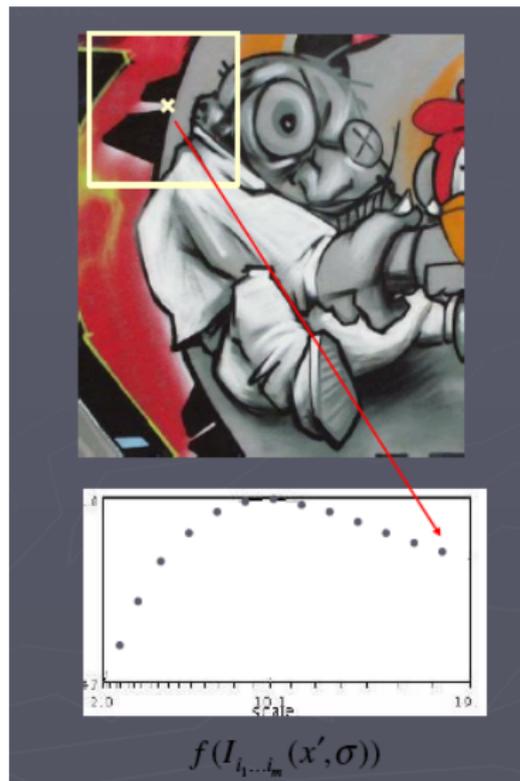
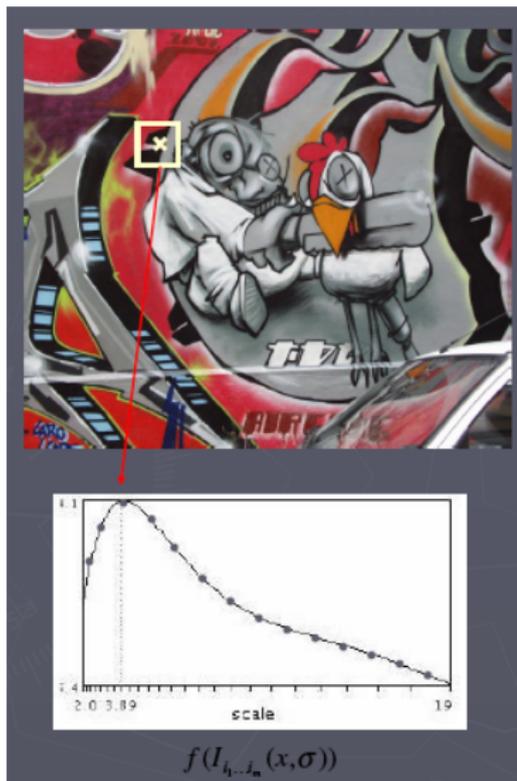
Automatic Scale Selection

Function responses for increasing scale (scale signature).



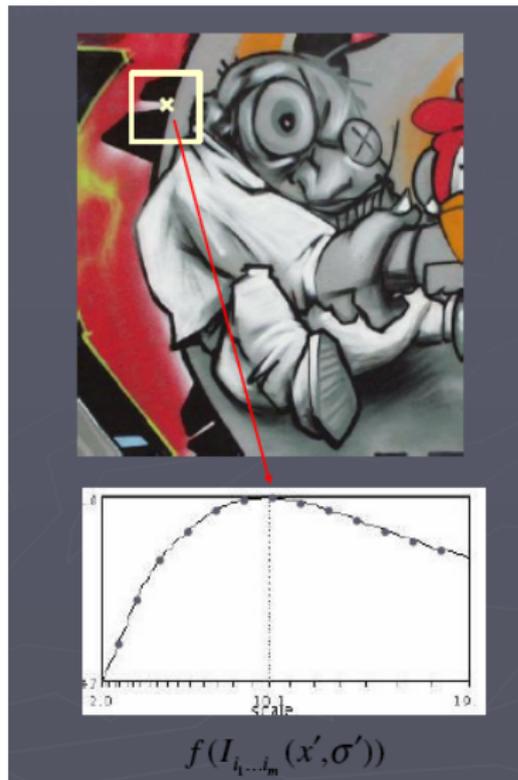
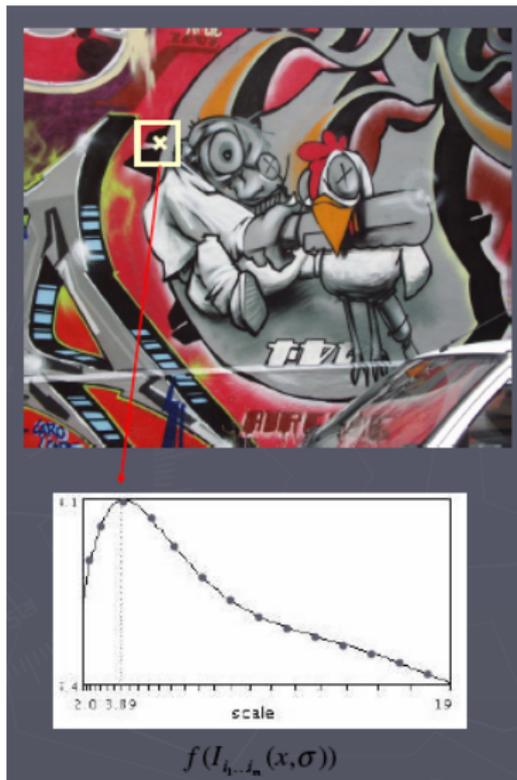
Automatic Scale Selection

Function responses for increasing scale (scale signature).



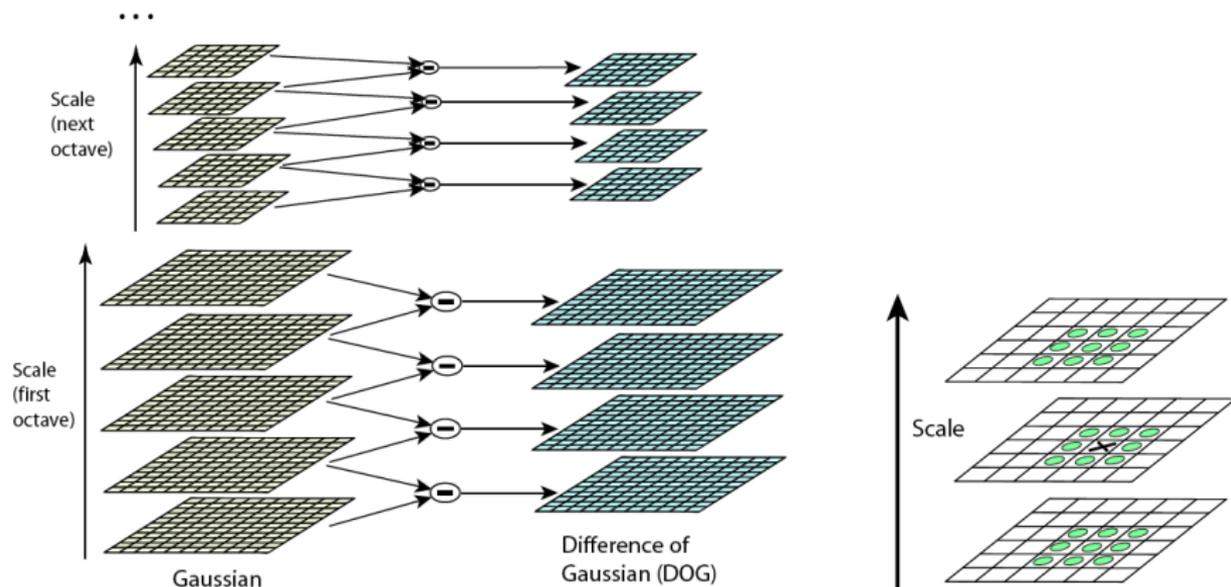
Automatic Scale Selection

Function responses for increasing scale (scale signature).



What can the signature function be?

- Lindeberg (1998): extrema in the Laplacian of Gaussian (LoG).
- Lowe (2004) proposed computing a set of sub-octave Difference of Gaussian filters looking for 3D (space+scale) maxima in the resulting structure.

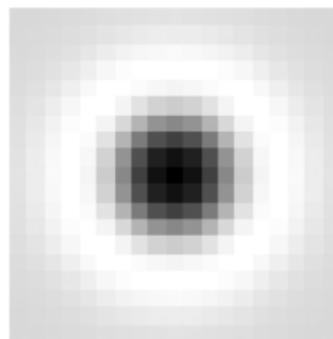
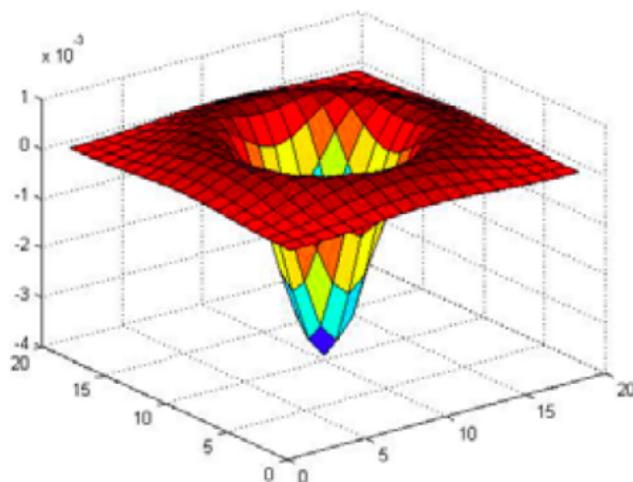


[Source: R. Szeliski]

Blob detection

- Laplacian of Gaussian: Circularly symmetric operator for blob detection in 2D

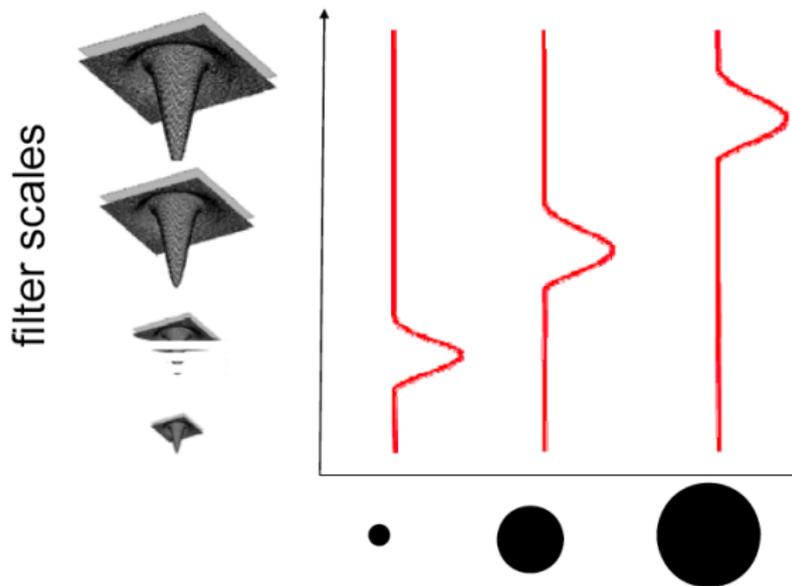
$$\nabla^2 g = \frac{\partial^2 g}{\partial x^2} + \frac{\partial^2 g}{\partial y^2}$$



[Source: K. Grauman]

Blob detection in 2D: scale selection

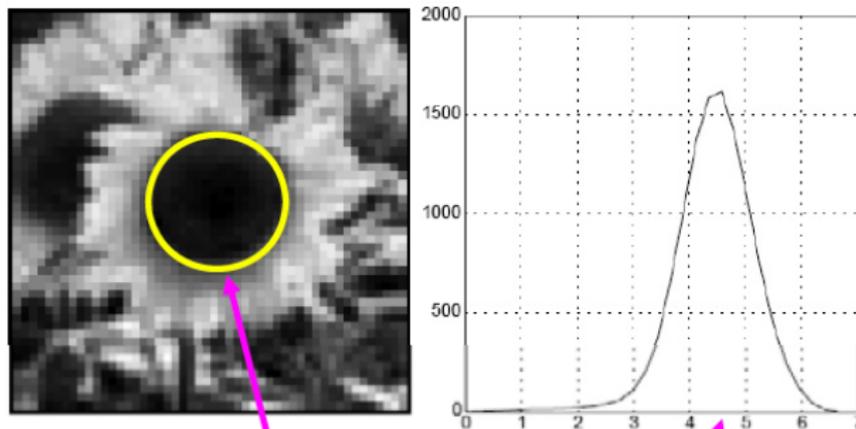
Laplacian-of-Gaussian = blob detector



[Source: B. Leibe]

Characteristic Scale

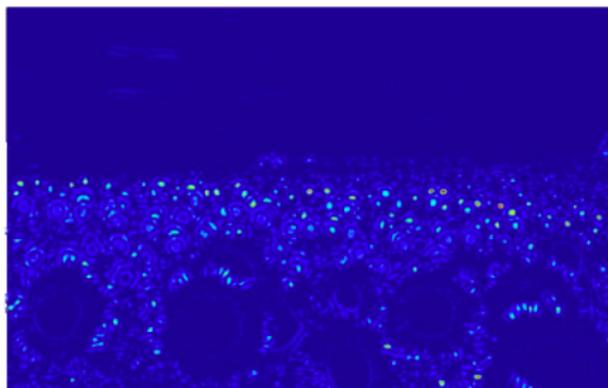
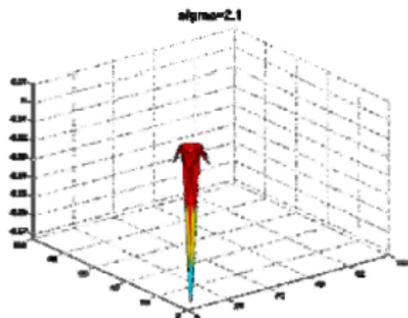
- We define the **characteristic scale** as the scale that produces peak of Laplacian response



characteristic scale

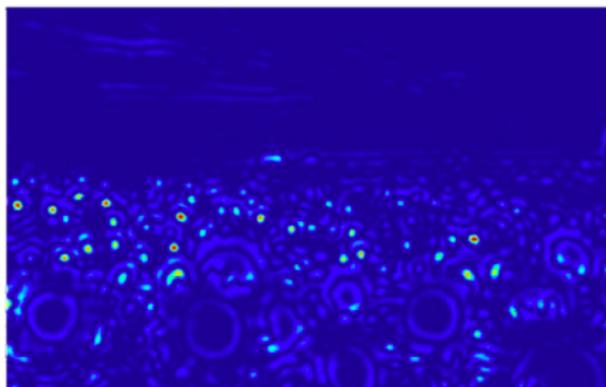
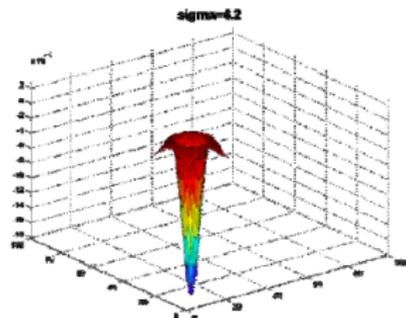
[Source: S. Lazebnik]

Example



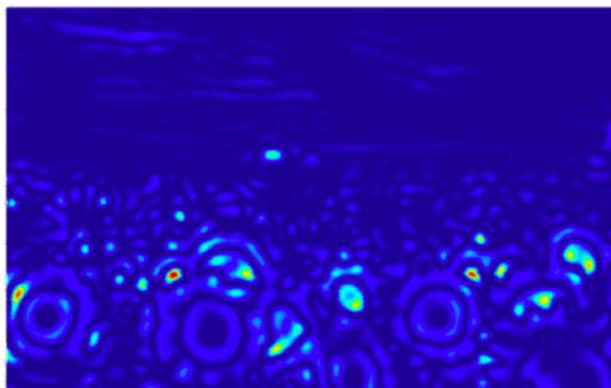
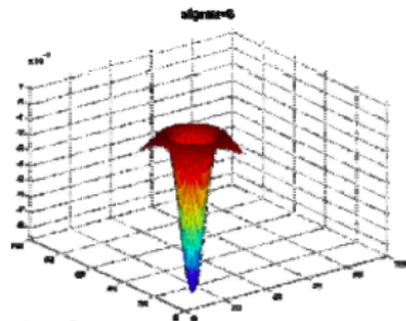
[Source: K. Grauman]

Example



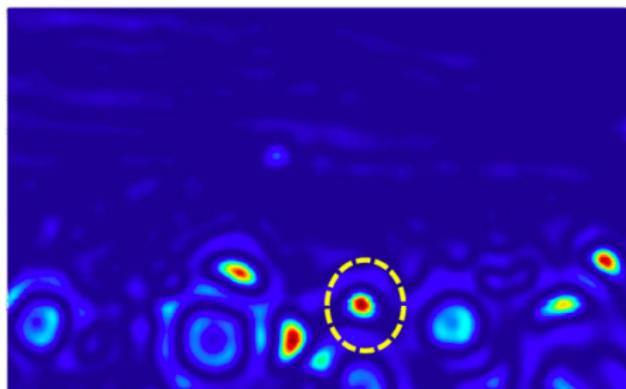
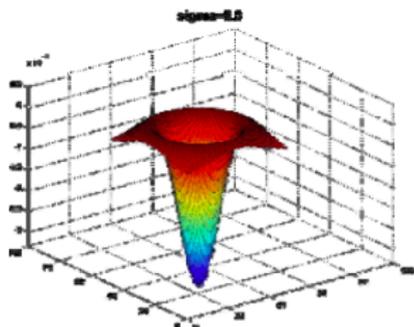
[Source: K. Grauman]

Example



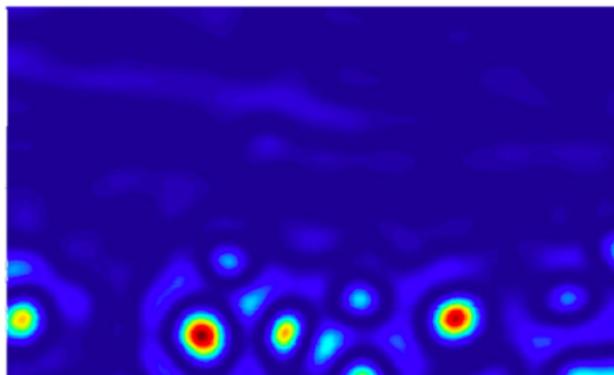
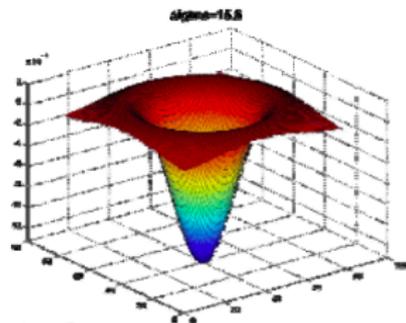
[Source: K. Grauman]

Example



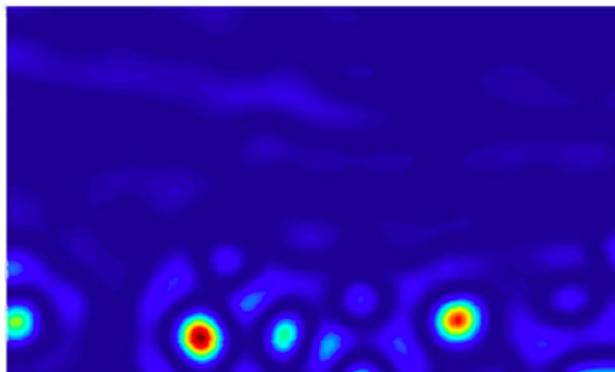
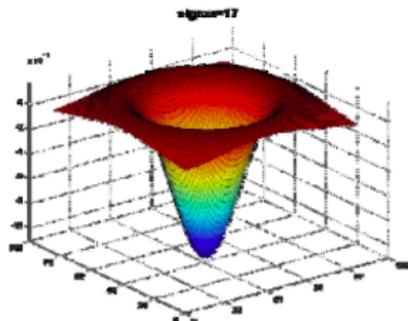
[Source: K. Grauman]

Example



[Source: K. Grauman]

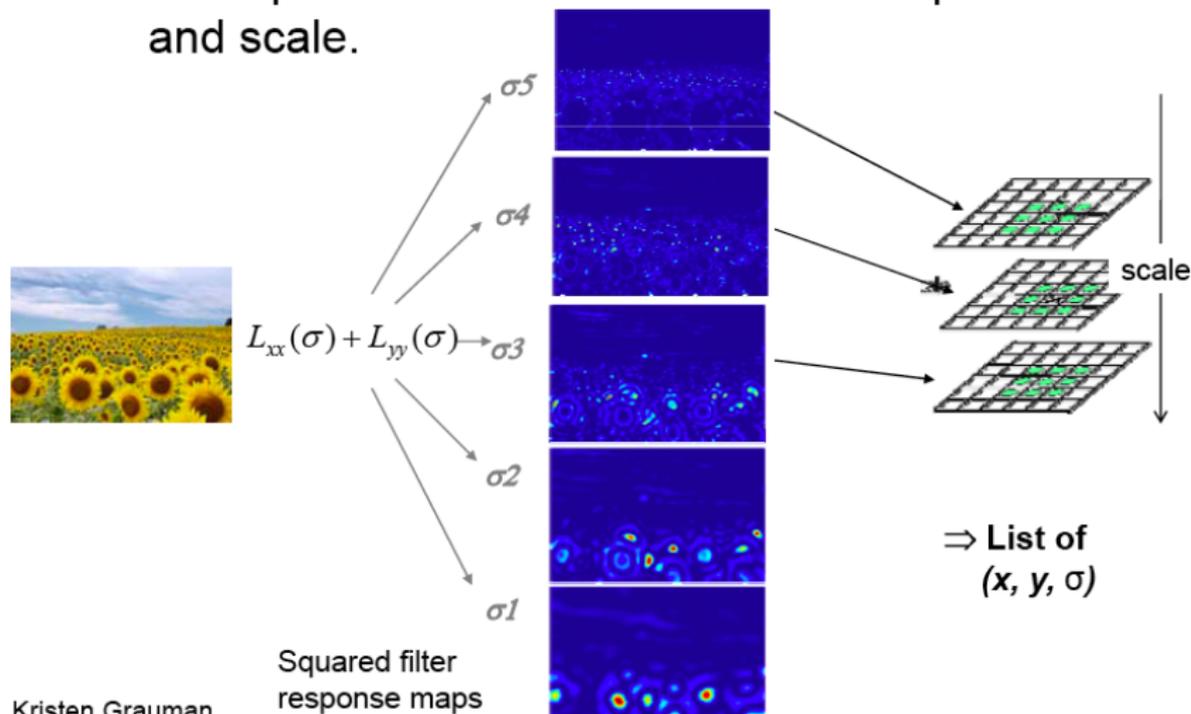
Example



[Source: K. Grauman]

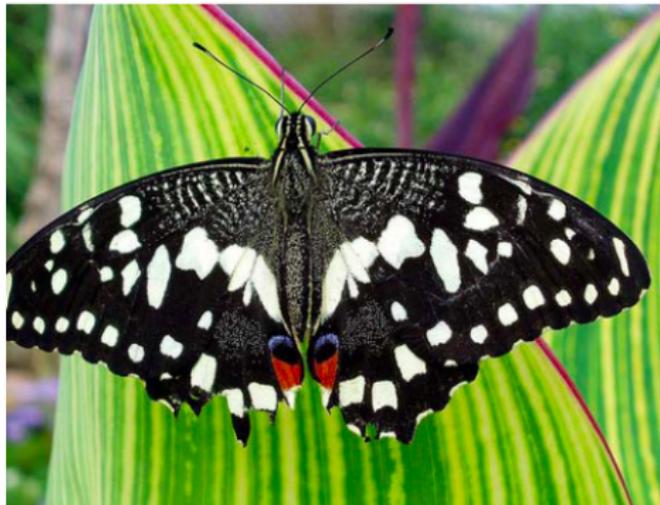
Scale invariant interest points

Interest points are local maxima in both position and scale.



Kristen Grauman

Example



[Source: S. Lazebnik]

Fast approximation

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

(Laplacian)

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

(Difference of Gaussians)

$I(k\sigma)$



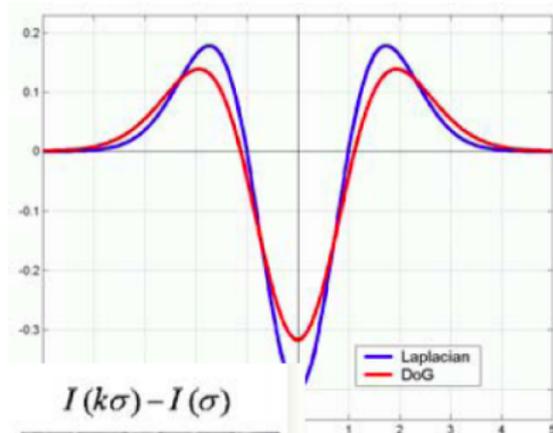
-

$I(\sigma)$



=

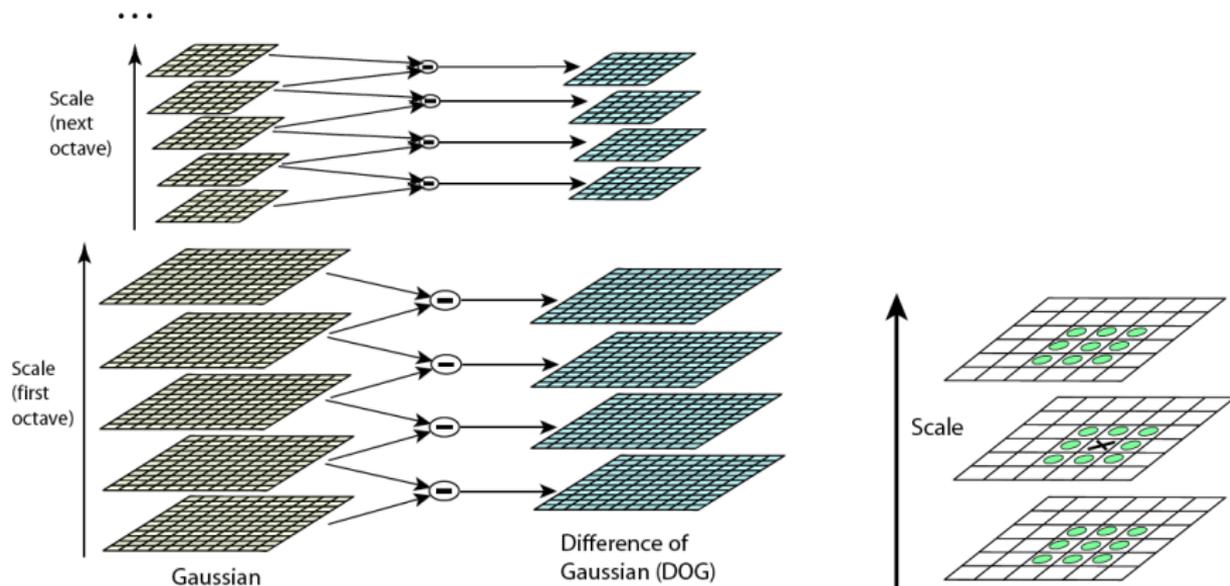
$I(k\sigma) - I(\sigma)$



[Source: K. Grauman]

Lowe's DoG

- Lowe (2004) proposed computing a set of sub-octave Difference of Gaussian filters looking for 3D (space+scale) maxima in the resulting structure



[Source: R. Szeliski]

Laplacian vs Hessian

- Laplacian of Gaussians is scale invariant.
- Simple and efficient.
- But fires more on edges than determinant of hessian



Properties of the ideal feature

- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.

Properties of the ideal feature

- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.
- **Robust:** noise, blur, discretization, compression, etc. do not have a big impact on the feature.

Properties of the ideal feature

- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.
- **Robust:** noise, blur, discretization, compression, etc. do not have a big impact on the feature.
- **Distinctive:** individual features can be matched to a large database of objects.

Properties of the ideal feature

- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.
- **Robust:** noise, blur, discretization, compression, etc. do not have a big impact on the feature.
- **Distinctive:** individual features can be matched to a large database of objects.
- **Quantity:** many features can be generated for even small objects.

Properties of the ideal feature

- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.
- **Robust:** noise, blur, discretization, compression, etc. do not have a big impact on the feature.
- **Distinctive:** individual features can be matched to a large database of objects.
- **Quantity:** many features can be generated for even small objects.
- **Accurate:** precise localization.

Properties of the ideal feature

- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.
- **Robust:** noise, blur, discretization, compression, etc. do not have a big impact on the feature.
- **Distinctive:** individual features can be matched to a large database of objects.
- **Quantity:** many features can be generated for even small objects.
- **Accurate:** precise localization.
- **Efficient:** close to real-time performance.

Properties of the ideal feature

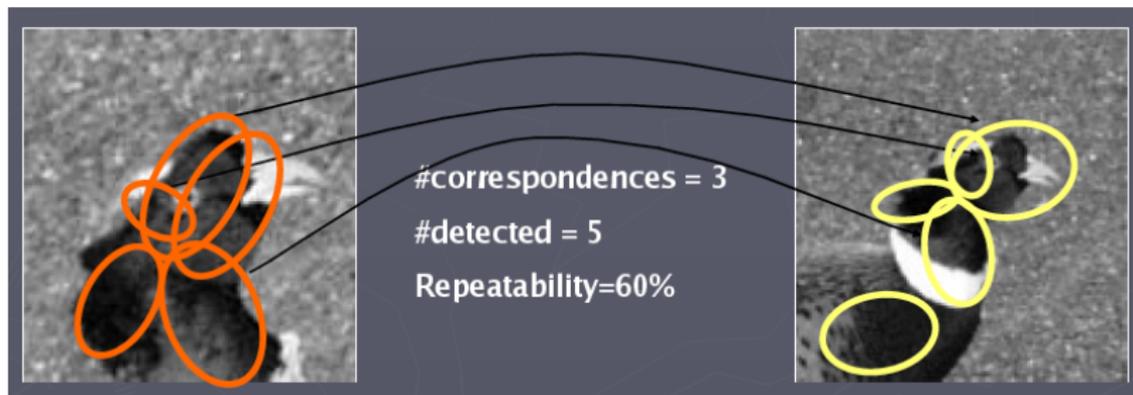
- **Local:** features are local, so robust to occlusion and clutter (no prior segmentation).
- **Invariant:** to certain transformations, e.g, scale, rotation.
- **Robust:** noise, blur, discretization, compression, etc. do not have a big impact on the feature.
- **Distinctive:** individual features can be matched to a large database of objects.
- **Quantity:** many features can be generated for even small objects.
- **Accurate:** precise localization.
- **Efficient:** close to real-time performance.

A lot of other interest point detectors

- Hessian
- Lowe: DoG
- Lindeberg: scale selection
- Miikolajczyk & Schmid: Hessian/Harris-Laplacian/Affine
- Tuytelaars & Van Gool: EBR and IBR
- Matas: MSER
- Kadir & Brady: Salient Regions
- Speeded-Up Robust Features (SURF) of Bay et al.

Evaluation criteria: repeatability

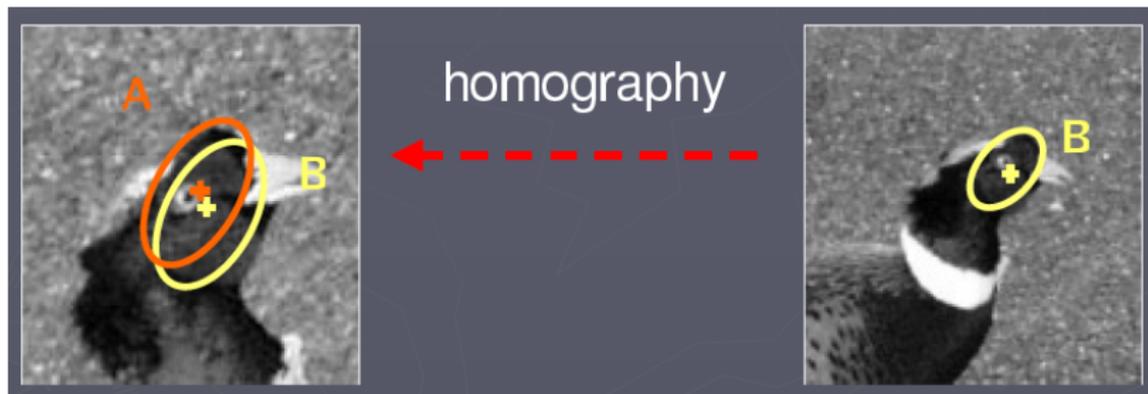
- Repeatability rate: percentage of detected that have correct corresponding points
- What's the problem of this?



[Source: T. Tuyttellaars]

Evaluation criteria: repeatability

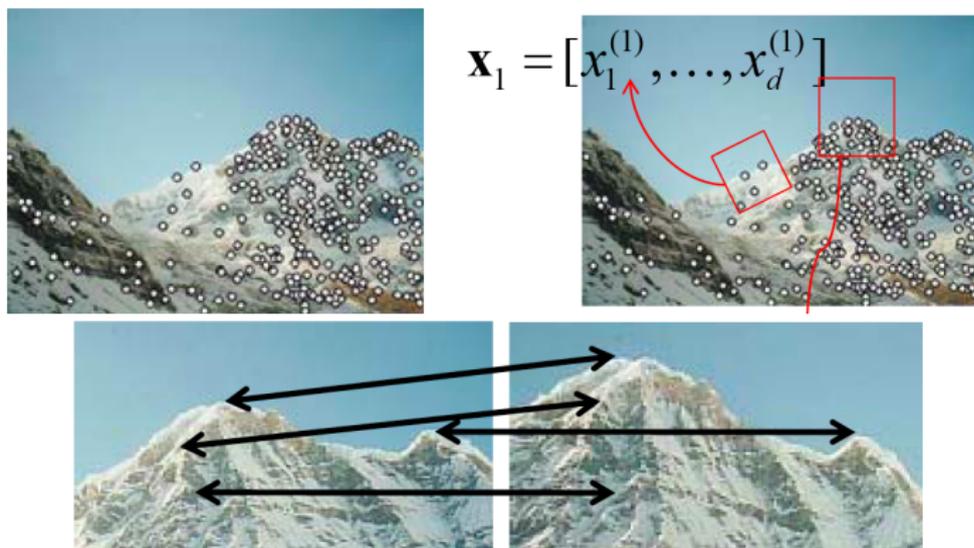
- Two points are in correspondence if the intersection over union is bigger than a certain threshold.



[Source: T. Tuytellaars]

Local features

- **Detection:** Identify the interest points.
- **Description:** Extract vector feature descriptor around each interest point.
- **Matching:** Determine correspondence between descriptors in two views.



[Source: K. Grauman]

The ideal feature descriptor

- Repeatable (invariant/robust)
- Distinctive
- Compact
- Efficient

Invariances



[Source: T. Tuytelaars]

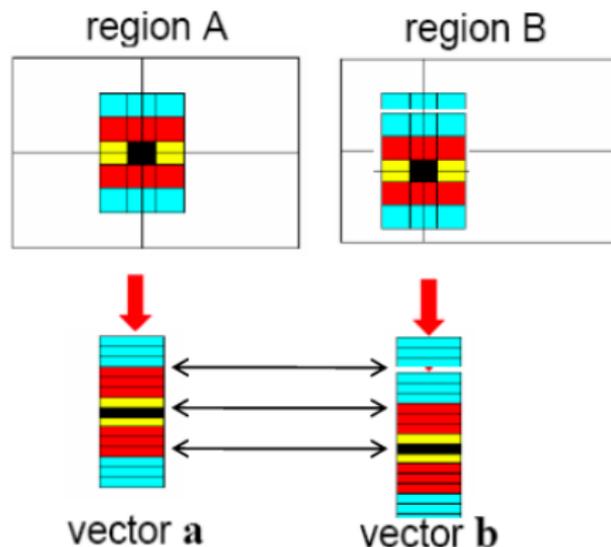
Invariances



[Source: T. Tuytelaars]

Raw Pixels as Local Descriptors

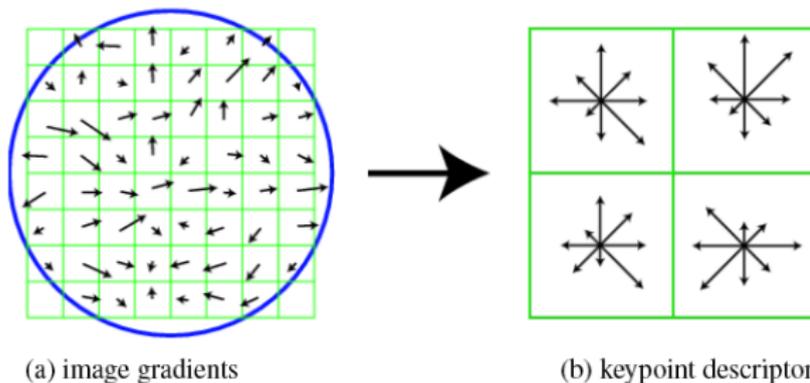
- The simplest way is to write down the list of intensities to form a feature vector, and normalize them (i.e., mean 0, variance 1).
- But this is very sensitive to even small shifts, rotations.



[Source: K. Grauman]

SIFT descriptor [Lowe 2004]

- Compute the gradient at each pixel in a 16×16 window around the detected keypoint, using the appropriate level of the Gaussian pyramid at which the keypoint was detected.
- Doweight gradients by a Gaussian fall-off function (blue circle) to reduce the influence of gradients far from the center.
- In each 4×4 quadrant, compute a gradient orientation histogram using 8 orientation histogram bins.



[Source: R. Szeliski]

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.
- To reduce the effects of contrast or gain (additive variations are already removed by the gradient), the 128-D vector is normalized to unit length.

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.
- To reduce the effects of contrast or gain (additive variations are already removed by the gradient), the 128-D vector is normalized to unit length.
- To further make the descriptor robust to other photometric variations, values are clipped to 0.2 and the resulting vector is once again renormalized to unit length.

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.
- To reduce the effects of contrast or gain (additive variations are already removed by the gradient), the 128-D vector is normalized to unit length.
- To further make the descriptor robust to other photometric variations, values are clipped to 0.2 and the resulting vector is once again renormalized to unit length.
- Great engineering effort!

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.
- To reduce the effects of contrast or gain (additive variations are already removed by the gradient), the 128-D vector is normalized to unit length.
- To further make the descriptor robust to other photometric variations, values are clipped to 0.2 and the resulting vector is once again renormalized to unit length.
- Great engineering effort!
- Why subpatches?

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.
- To reduce the effects of contrast or gain (additive variations are already removed by the gradient), the 128-D vector is normalized to unit length.
- To further make the descriptor robust to other photometric variations, values are clipped to 0.2 and the resulting vector is once again renormalized to unit length.
- Great engineering effort!
- Why subpatches?
- Why does SIFT have some illumination invariance?

SIFT descriptor [Lowe 2004]

- To reduce the effects of location and dominant orientation misestimation, each of the original 256 weighted gradient magnitudes is softly added to nearby bins.
- The resulting 128 non-negative values form a raw version of the SIFT descriptor vector.
- To reduce the effects of contrast or gain (additive variations are already removed by the gradient), the 128-D vector is normalized to unit length.
- To further make the descriptor robust to other photometric variations, values are clipped to 0.2 and the resulting vector is once again renormalized to unit length.
- Great engineering effort!
- Why subpatches?
- Why does SIFT have some illumination invariance?

Making descriptor rotation invariant

- Rotate patch according to its dominant gradient orientation
- This puts the patches into a canonical orientation

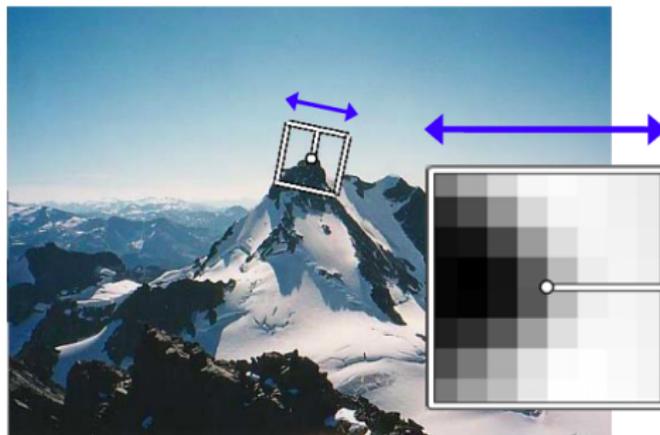


Figure: Figure from M. Brown

[Source: K. Grauman]

SIFT descriptor [Lowe 2004]

Extraordinarily robust matching technique

- Changes in viewpoint: up to about 60 degree out of plane rotation
- Changes in illumination: sometimes even day vs. night
- Fast and efficient: can run in real time
- Lots of code available



[Source: S. Seitz]

Example

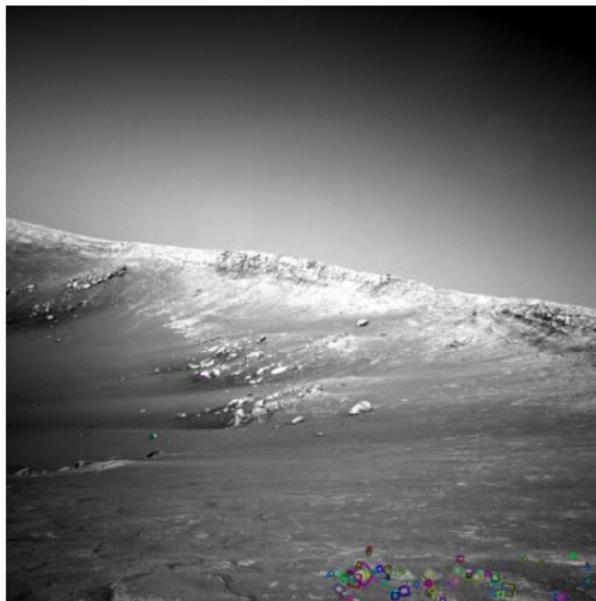


Figure: NASA Mars Rover images with SIFT feature matches

[Source: N. Snavely]

Invariant to

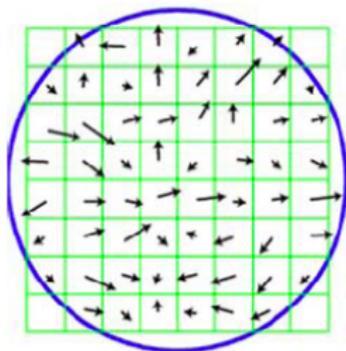
- Scale
- Rotation

Partially invariant to

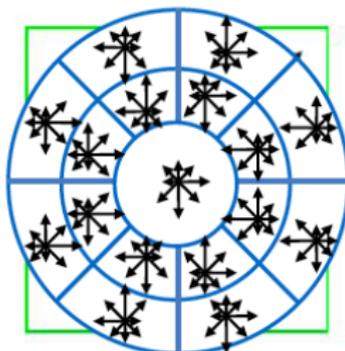
- Illumination changes
- Camera viewpoint
- Occlusion, clutter

Gradient location-orientation histogram (GLOH)

- Developed by Mikolajczyk and Schmid (2005): variant on SIFT that uses a log-polar binning structure instead of the four quadrants.
- The spatial bins are 11, and 15, with eight angular bins (except for the central region), for a total of 17 spatial bins and 16 orientation bins.
- The 272D histogram is then projected onto a 128D descriptor using PCA trained on a large database.



(a) image gradients



(b) keypoint descriptor

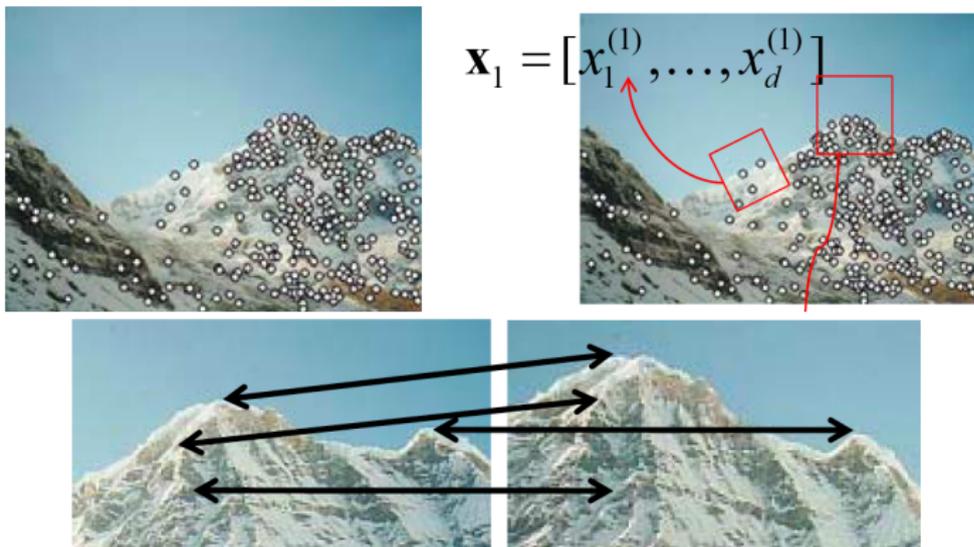
[Source: R. Szeliski]

Other Descriptors

- Steerable filters
- moment invariants
- complex filters
- shape contexts
- PCA-SIFT
- HOG
- SURF
- DAISY

Local features

- **Detection:** Identify the interest points.
- **Description:** Extract vector feature descriptor around each interest point.
- **Matching:** Determine correspondence between descriptors in two views.



[Source: K. Grauman]

Matching local features

Once we have extracted features and their descriptors, we need to match the features between these images.

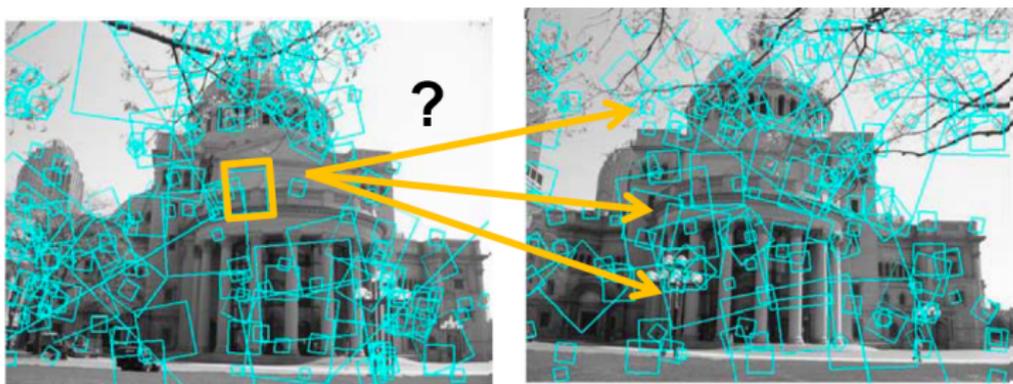
- Matching strategy: which correspondences are passed on to the next stage
- Devise efficient data structures and algorithms to perform this matching



Figure: Images from K. Grauman

Matching local features

- To generate candidate matches, find patches that have the most similar appearance (e.g., lowest SSD)
- Simplest approach: compare them all, take the closest (or closest k , or within a thresholded distance)



[Source: K. Grauman]

Ambiguous matches

- At what SSD value do we have a good match?
- To add robustness, consider ratio of distance to best match to distance to second best match
 - If low, first match looks good.
 - If high, could be ambiguous match.



[Source: K. Grauman]

Matching SIFT Descriptors

- Nearest neighbor (Euclidean distance)
- Threshold ratio of nearest to 2nd nearest descriptor

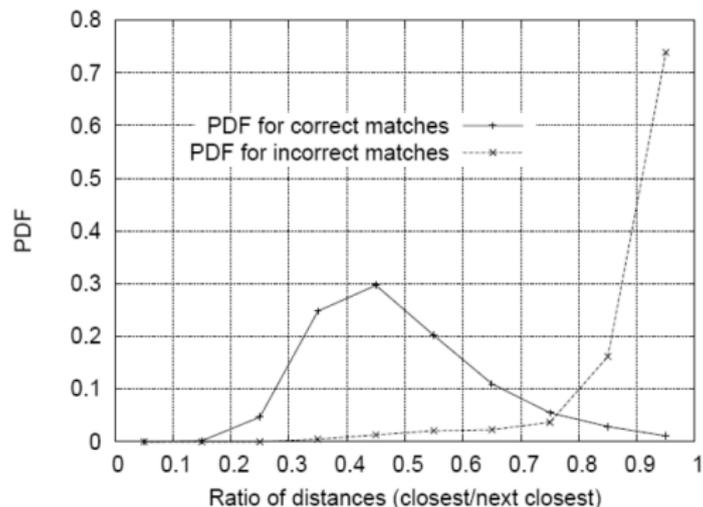


Figure: Images from D. Lowe

[Source: K. Grauman]

Which threshold to use?

- Setting the threshold too high results in too many false positives, i.e., incorrect matches being returned.
- Setting the threshold too low results in too many false negatives, i.e., too many correct matches being missed

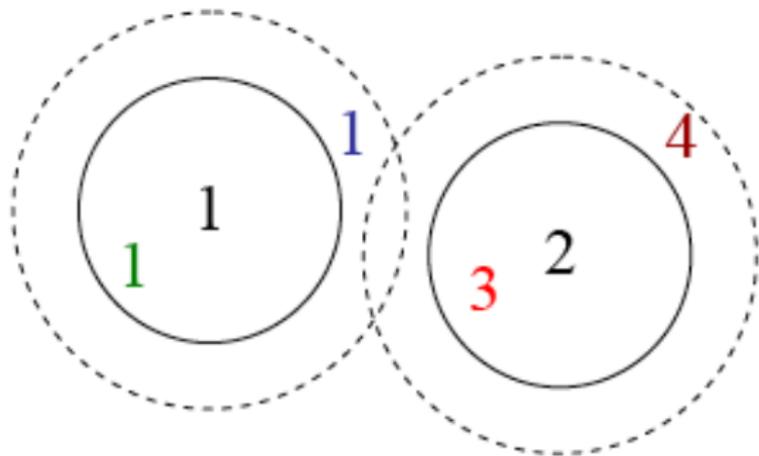


Figure: Images from R. Szeliski

How to quantize how good is our matching?

- TP: true positives, i.e., number of correct matches
- FN: false negatives, matches that were not correctly detected
- FP: false positives, proposed matches that are incorrect
- TN: true negatives, non-matches that were correctly rejected.

$$\text{True positive rate (recall)} \quad TPR = \frac{TP}{TP + FN} = \frac{TP}{P}$$

$$\text{True negative rate} \quad TNR = \frac{FP}{FP + TN} = \frac{FP}{N}$$

$$\text{positive predictive value (precision)} \quad PPV = \frac{TP}{TP + FP} = \frac{TP}{P'}$$

$$\text{accuracy} \quad ACC = \frac{TP + TN}{P + N}$$

Measuring performance

- Any particular matching strategy (at a particular threshold or parameter setting) can be rated by the TPR and FPR numbers
- We want $TPR=1$ and $FPR=0$.
- As we vary the matching threshold, we obtain a family of such points, i.e., receiver operating characteristic (ROC curve)
- The closer this curve lies to the upper left corner, the better its performance.

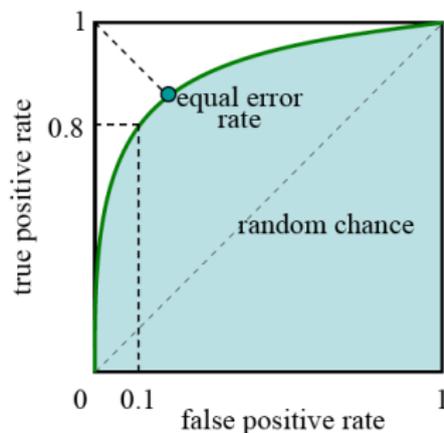


Figure: Images from R. Szeliski

Measuring performance

- Area under the curve (AUC) is a way to summarize ROC with 1 number.
- Mean average precision, which is the average precision (PPV) as you vary the threshold.
- The equal error rate is sometimes used as well.

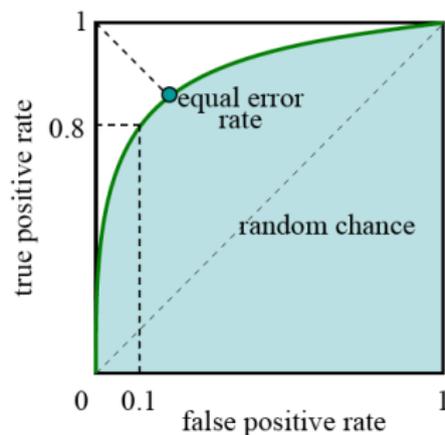


Figure: Images from R. Szeliski

Applications of local invariant features

- Wide baseline stereo
- Motion tracking
- Panoramas
- Mobile robot navigation
- 3D reconstruction
- Recognition

[Source: K. Grauman]

Wide Baseline Stereo



[Source: T. Tuytelaars]

Recognizing the Same Object



Schmid and Mohr 1997



Sivic and Zisserman, 2003



Rothganger et al. 2003



Lowe 2002

[Source: K. Grauman]

Motion Tracking



Figure: Images from J. Pilet

Interest point detection

- Harris corner detector
- Laplacian of Gaussian, automatic scale selection
- Difference of Gaussians

Invariant descriptors

- Rotation according to dominant gradient direction
- Histograms for robustness to small shifts and translations (SIFT descriptor)
- Polar coordinate descriptors GLOH.

Category-level recognition

Recognizing or retrieving specific objects

- Example: Visual search in feature films

Visually defined query

“Find this clock”



“Find this place”



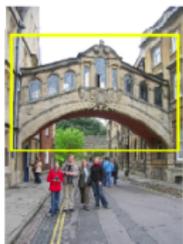
“Groundhog Day” [Rammis, 1993]



[Source: J. Sivic]

Recognizing or retrieving specific objects

- Example: Search photos on the web for particular places



Find these landmarks

...in these images and 1M more

[Source: J. Sivic]



Google Goggles

Use pictures to search the web. [▶ Watch a video](#)



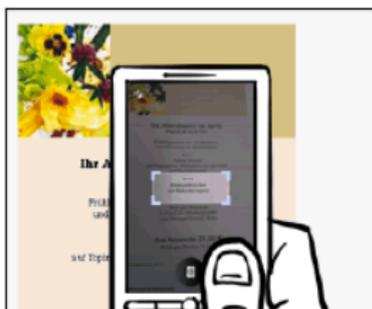
Get Google Goggles

Android (1.6+ required)
Download from Android Market.

[Send Goggles to Android phone](#)

New: iPhone (iOS 4.0 required)
Download [from the App Store](#).

[Send Goggles to iPhone](#)

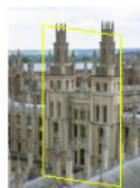
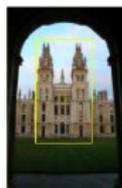


Why is it difficult?

- Want to find the object despite possibly large changes in scale, viewpoint, lighting and partial occlusion.
- We can't expect to match such varied instances with a single global template...



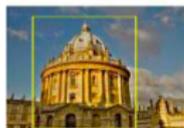
Scale



Viewpoint



Lighting

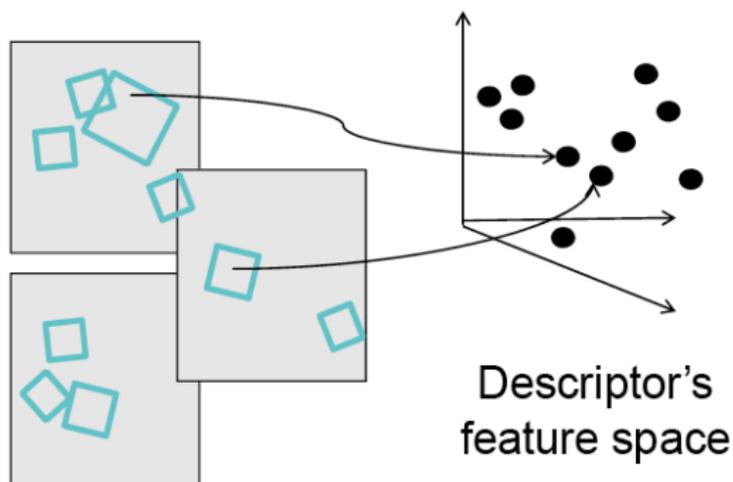


Occlusion

[Source: J. Sivic]

Indexing local features

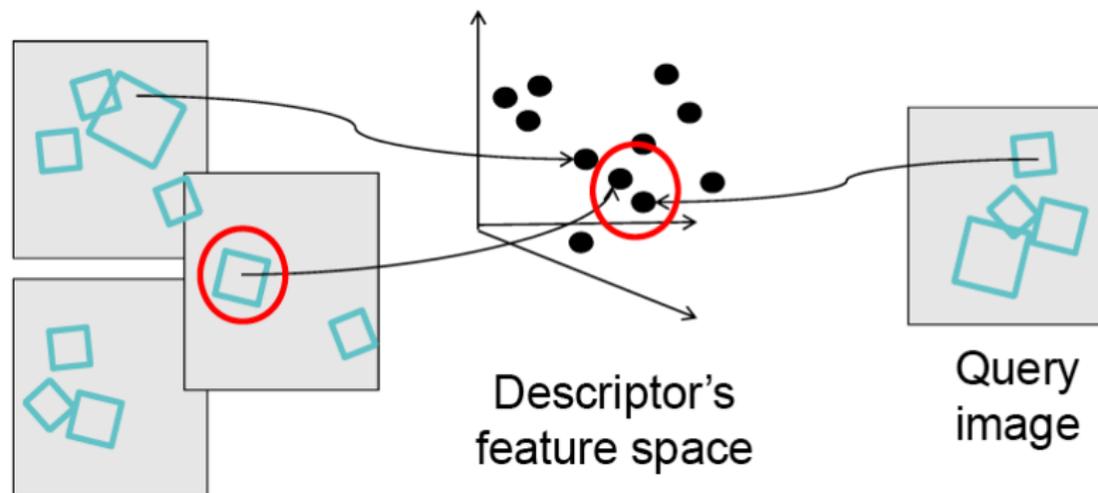
- Each patch / region has a descriptor, which is a point in some high-dimensional feature space (e.g., SIFT)



[Source: K. Grauman]

Indexing local features

- It can have millions of features to search.



[Source: K. Grauman]

Indexing local features: inverted file index

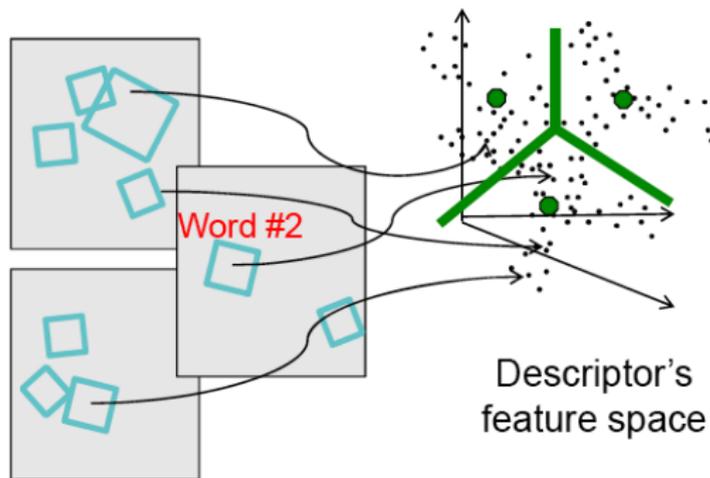
- For text documents, an efficient way to find all pages on which a word occurs is to use an index.
- We want to find all images in which a feature occurs.
- To use this idea, we'll need to map our features to visual words.
- Why?

Index		
*Along I-75, From Detroit to Florida, inside back cover	Buttery Center, McClure, 134	Diving Lenses, 85
*Shore I&L, From Boston to Florida, inside back cover	CAA (see AAA)	Duval County, 163
1909 Spanish Trail Roadway, 101-102, 104	CCC, The, 111, 112, 113, 126, 142	Eac, Gallo, 176
AAA (Basic info - 100 Access, 86 AAA (and CAA), 83	Cs vZur, 147	Edison, Thomas, 152
AAA National Office, 88	Catowahatchee River, 152	Edin AFJ, 116-118
Alabama, 124	Name, 150	Eggt House, 178
Alabama, 124	Catonwood Hotel Seashore, 173	Elewan, 144-145
County, 131	Canon Creek Airport, 130	Emerson Point Inn, 120
Andale River, 143	Canopy Road, 106, 169	Emergency Carbores, 63
Atapaha, Name, 126	Cape Catawba, 174	Eggsby, 142, 143, 157, 159
Atford B Miley Gardens, 106	Castle Sea Marine, 169	Escambia Bay, 119
Atgator Farm, St Augustine, 189	Cave Diving, 131	Edge (J-10), 119
Atgator Inn (attribution), 167	Cape Coral, Name, 150	County, 102
Atgator, Buddy, 155	Celibration, 89	Edwin, 159
Atgators, 103, 126, 136, 147, 156	Charlotte County, 149	Dewlight, 90, 95, 109, 140, 154-160
Aurassia Island, 170	Charlotte Harbor, 150	Drawing of, 156, 181
Aurassia, 108-109, 148	Chautauque, 116	Waldie MA, 160
Austchecha River, 112	Chapay, 114	Wander Gardens, 154
Austrian Mus of Art, 136	Name, 115	Falling Waters SP, 115
Azules, 152	Chocowatch, Name, 115	Factory of Flight 50
Aztec-Night, 94	Cinco Museum, Ringing, 147	Fayer Dyles SP, 171
Art Museum, Ringing, 147	Cities, 88, 97, 126, 136, 140, 189	Fawn Forest, 168
Artie Beach Club, 182	City/Place, W Palm Beach, 100	Fawn, Prescribed, 148
Aucilla River Project, 106	City Maps,	Fisherman's Village, 751
Babcock-Wash NMAA, 151	FL Lighthouse Expeds, 194-196	Fisher County, 171
Bahia Mir Marina, 184	Jacksonville, 163	Fisher, Henry, 97, 165, 167, 171
Baker County, 99	Kissimmee Expeds, 192-193	Florida Aquarium, 186
Banfield Malman, 162	Miami Expressways, 194-195	Florida,
Barge Canal, 127	Oslando Expressways, 190-193	12,000 years ago, 167
Basin Line Expy, 80	Pensacola, 98	Caves SP, 114
Basin Cultural Mnt, 89	Tallahassee, 191	Map of all Expressways, 2-3
Barnard Center, 136	Tampa St, Petersburg, 63	Map of Natural History, 184
Big T, 165	St. Augustine, 191	National Cemetery, 141
Big Cypress, 155, 158	Old War, 103, 105, 171, 186, 141	Part of Alaska, 177
	Oliverwater Marine Aquarium, 187	Platform, 187
	Older County, 154	Sheep's Boys Camp, 126
	Callie, Barren, 152	Spain's Hall of Fame, 120
	Cultural Spanish Quarters, 168	Sun 'n Fun Museum, 97
	Columbia County, 101-120	Supreme Court, 167
	Copains Building Material, 165	Florida's Turnpike (TP), 178, 189
	Conkrew Swamp, Name, 154	23 mile Strip Maps, 66
	Crabtree, 65	Administration, 18
	Croft Trap II, 144	Coast System, 190
	Croft Tractor, Florida, 88, 95, 132	Earl Services, 188
	Crosswater Expy, 11, 20, 86, 143	HEPT, 76, 181, 190
	Culter Break, 184	History, 189
	Dade Ditch, 143	Name, 189
	Dade, Maj Francis, 139-140, 161	Service Plaza, 190

[Source: K. Grauman]

Indexing local features: inverted file index

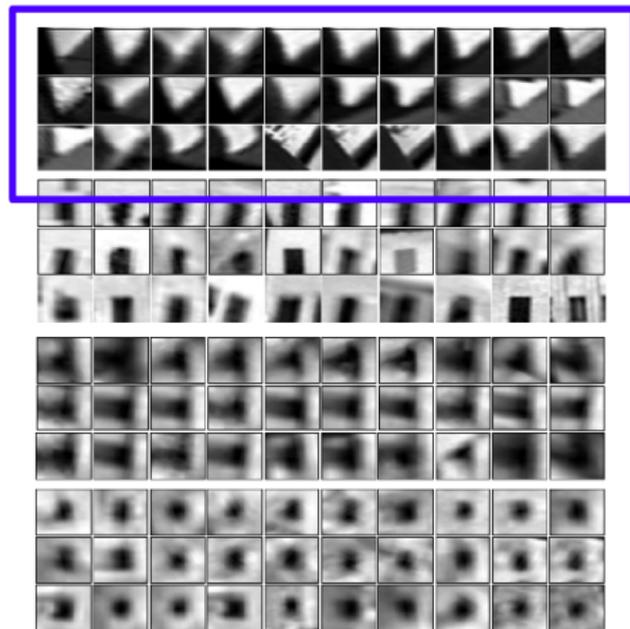
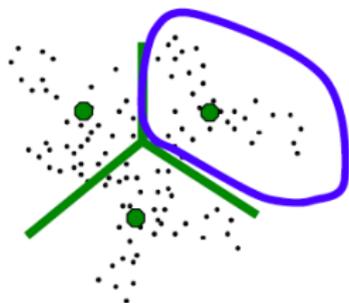
- Map high-dimensional descriptors to tokens/words by quantizing the feature space.
- Quantize via clustering, let cluster centers be the prototype words.
- Determine which word to assign to each new image region by finding the closest cluster.



[Source: K. Grauman]

Visual words

- Each group of patches belongs to the same visual word.



Visual vocabulary formation issues

- Vocabulary size, number of words
- Sampling strategy: where to extract features?
- Clustering / quantization algorithm
- Unsupervised vs. supervised
- What corpus provides features (universal vocabulary?)