



# Robotics Programming Laboratory

Bertrand Meyer  
Jiwon Shin

Lecture 10: Robot Perception

<http://pascallin.ecs.soton.ac.uk/challenges/VOC/databases.html#Caltech>



Given visual input, understand the information the input contains

- Object location: *object detection*
- Type of object: *object classification*
- Exact object name: *object recognition*
- Overall scene: *scene understanding*



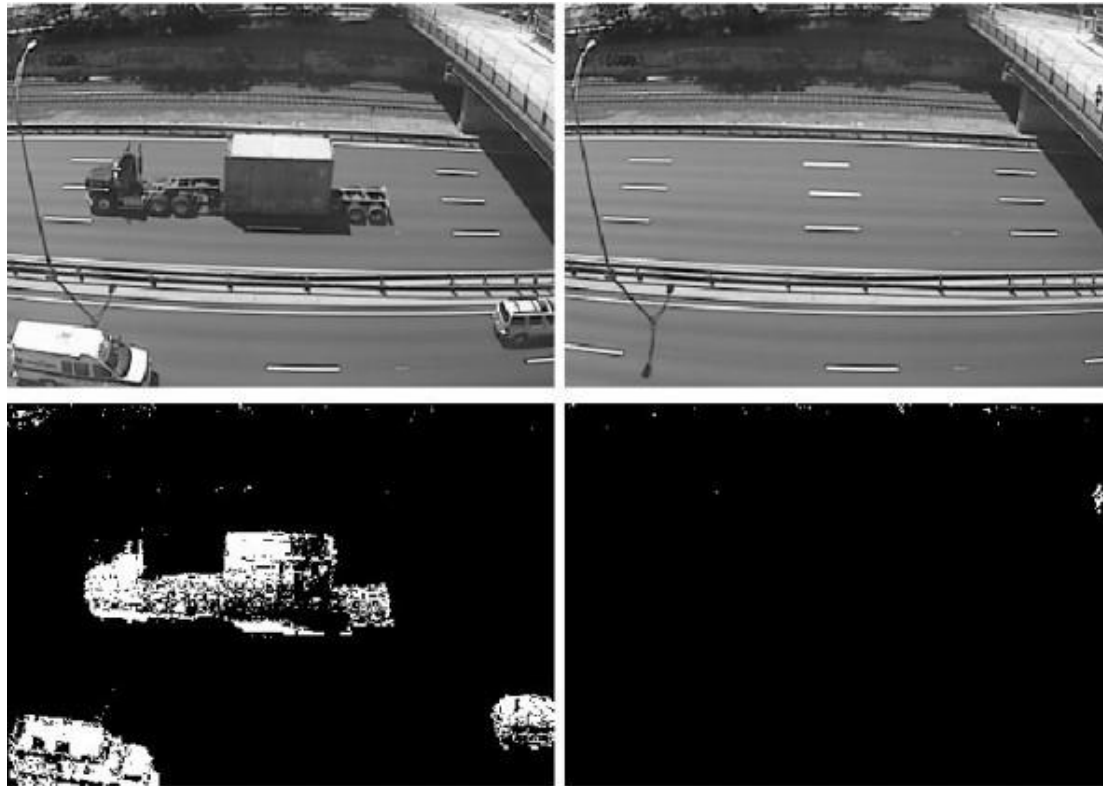
Segmentation: decomposition of an image into consistent regions

- Data that belong to the same region have similar properties
  - Similar color, texture, surface normal, etc.
- Data that belong to different regions have different properties
  - Different color, texture, surface normal, etc.
  
- Segmentation as clustering
  - Partitioning: divide an image into coherent regions
  - Grouping: group together elements of similar properties



- Divide an image into sensible regions using pixel intensity, color, texture, etc.
- Background subtraction
- Clustering
- Graph-based

# Background subtraction



[http://vip.bu.edu/files/2010/02/FDR\\_FPR\\_control\\_comparison1-594x636.jpg](http://vip.bu.edu/files/2010/02/FDR_FPR_control_comparison1-594x636.jpg)

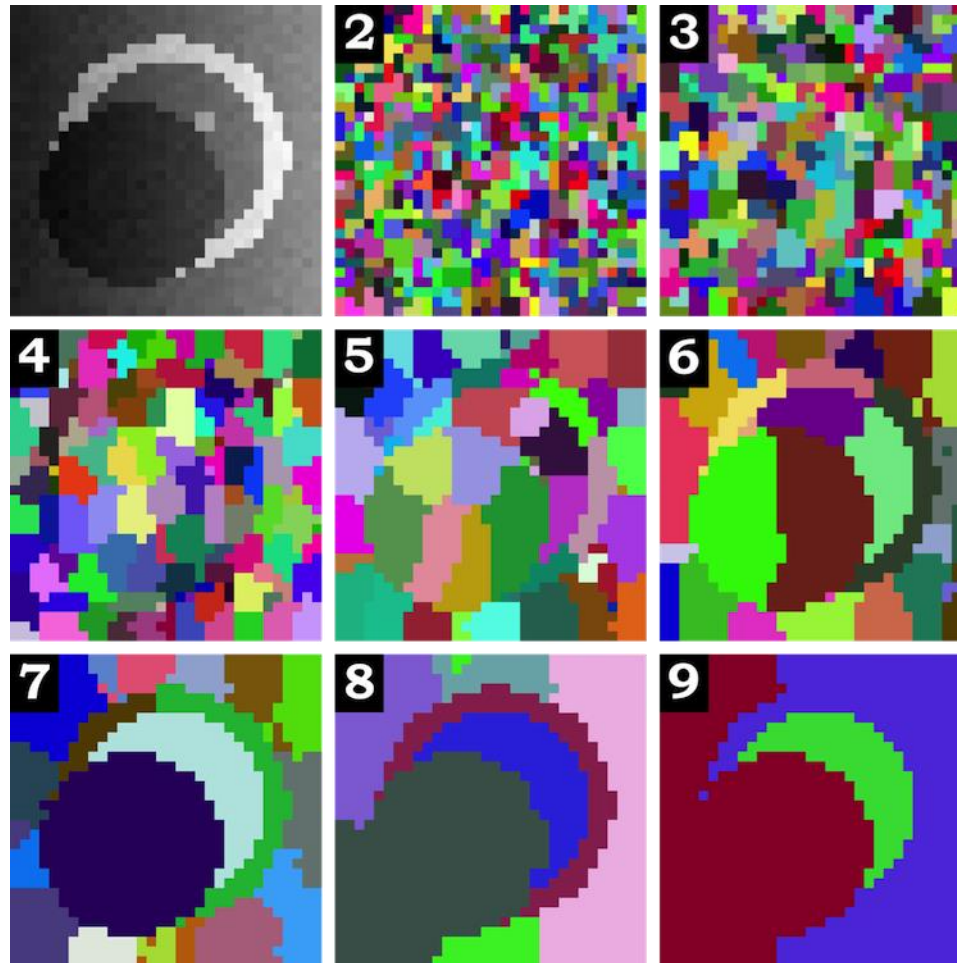


- Subtract an estimate of the appearance of the background from the image
- Consider areas of large absolute difference to be foreground

## Issues

- Obtaining a good estimate of the background is non-trivial
  - Changes in environment, lighting, weather, etc.
  - Use a moving average
- Threshold

# Agglomerative clustering



# Agglomerative clustering

---



- Consider each data point as a cluster
- Recursively merge the clusters with the smallest inter-cluster distance until the result is satisfactory

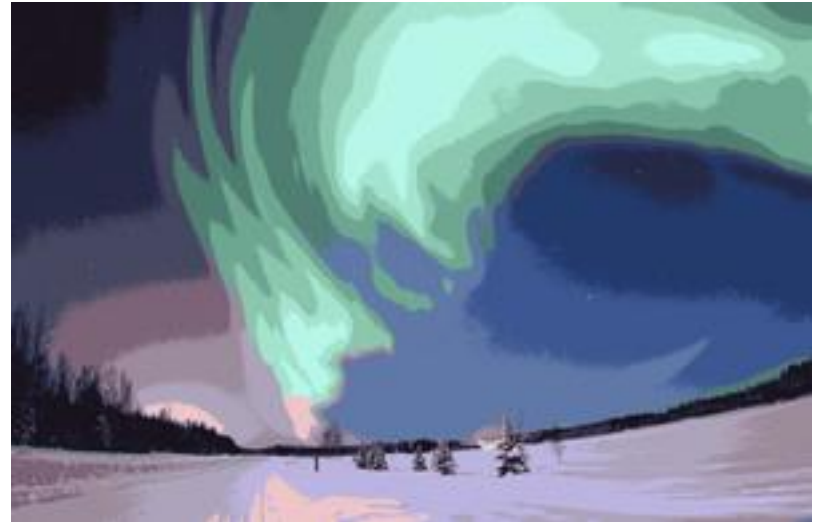
## Issues

- Inter-cluster distance
  - Distance between closest elements
  - Distance between farthest elements
  - Average distance between elements
- Number of clusters



# K-means clustering

---



[http://en.wikipedia.org/wiki/Segmentation\\_\(image\\_processing\)](http://en.wikipedia.org/wiki/Segmentation_(image_processing))



- Choose  $k$  data points as seed points
- Recursively assign each data point to the cluster whose center is the closest and recalculate the cluster mean until the center does not change

## Issues

- Segments are not connected in image
  - Using pixel coordinates would break up large regions
- Determining  $k$  is non-trivial

# Efficient graph-based image segmentation



Felzenszwalb, P. and Huttenlocher, D. 2004. "Efficient Graph-Based Image Segmentation"  
International Journal of Computer Vision, Volume 59, Number 2.

# Efficient graph-based image segmentation

---



- Represent image as a graph, each pixel being a node of a graph
- Edges are formed between neighboring pixels
- Merge the nodes such that nodes belonging to the same segment more similar to one another than nodes at the boundary of two segments



- Internal difference of a cluster  $c$ :

- $Int(C) = \max_{e \in MST(C,E)} w(e)$

- Difference between clusters  $c_1, c_2$ :

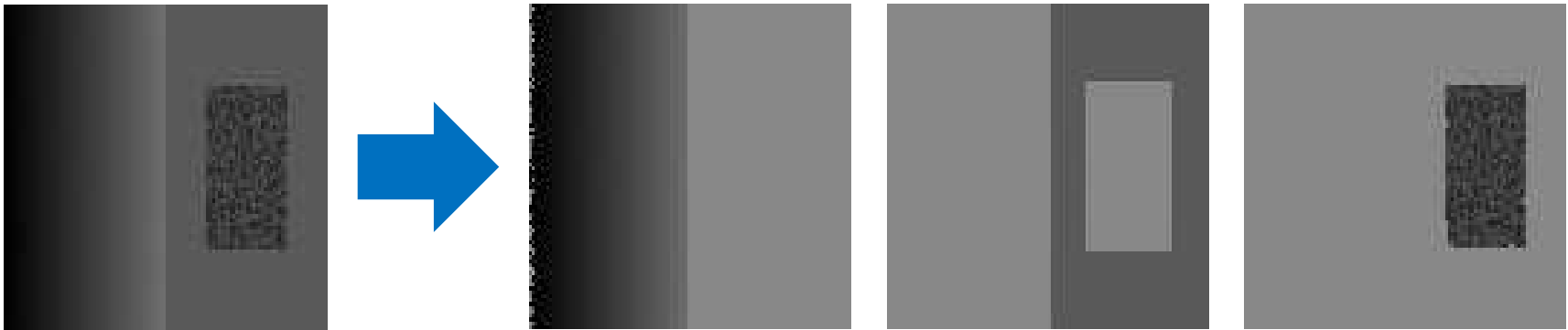
- $Dif(C_1, C_2) = \min_{v_i \in C_1, v_j \in C_2, (v_i, v_j) \in E} w((v_i, v_j))$

- Minimum internal difference:

- $MInt(C_1, C_2) = \min(Int(C_1) + \tau(C_1), Int(C_2) + \tau(C_2))$

- $\tau(C) = \frac{k}{|C|}$

- A boundary exists between  $c_1$  and  $c_2$  if  $Dif(C_1, C_2) > MInt(C_1, C_2)$



- Regions of consistent properties are grouped together

## Issues

- Number and quality of segments depend on the parameter  $k$ , smoothing factor, and minimum number of nodes



- Generally, we can use image segmentation algorithms by replacing intensity, color, or texture by surface normal
- Group together areas of consistent surface normal

Surface normal computation

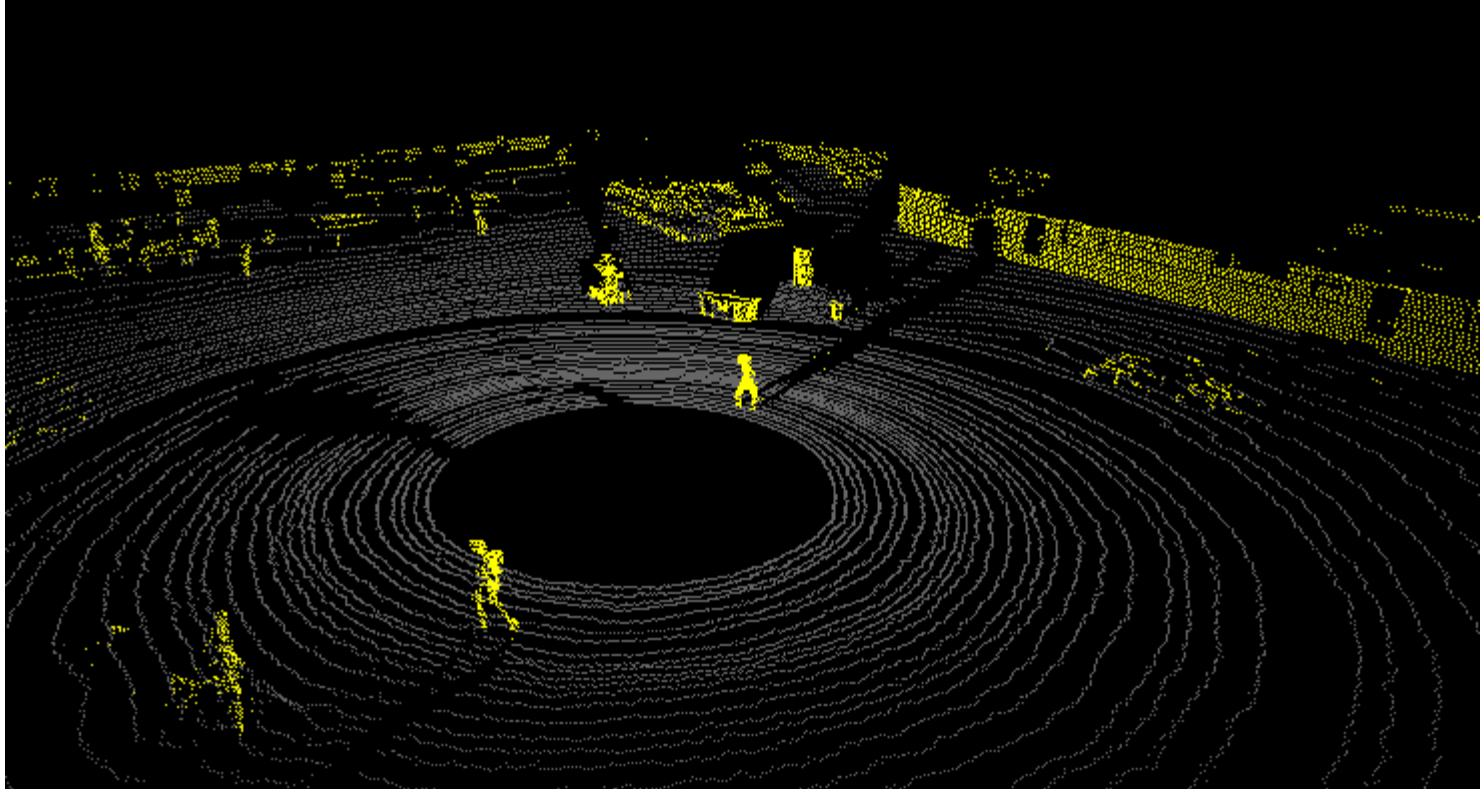
$$\mathbf{x}_u \equiv \partial \mathbf{x} / \partial u$$

$$\mathbf{x}_v \equiv \partial \mathbf{x} / \partial v$$

$$N = \frac{1}{|\mathbf{x}_u \times \mathbf{x}_v|} (\mathbf{x}_u \times \mathbf{x}_v)$$

# Ground segmentation

---



[http://www-personal.acfr.usyd.edu.au/p.morton/media/img/data\\_ground.png](http://www-personal.acfr.usyd.edu.au/p.morton/media/img/data_ground.png)



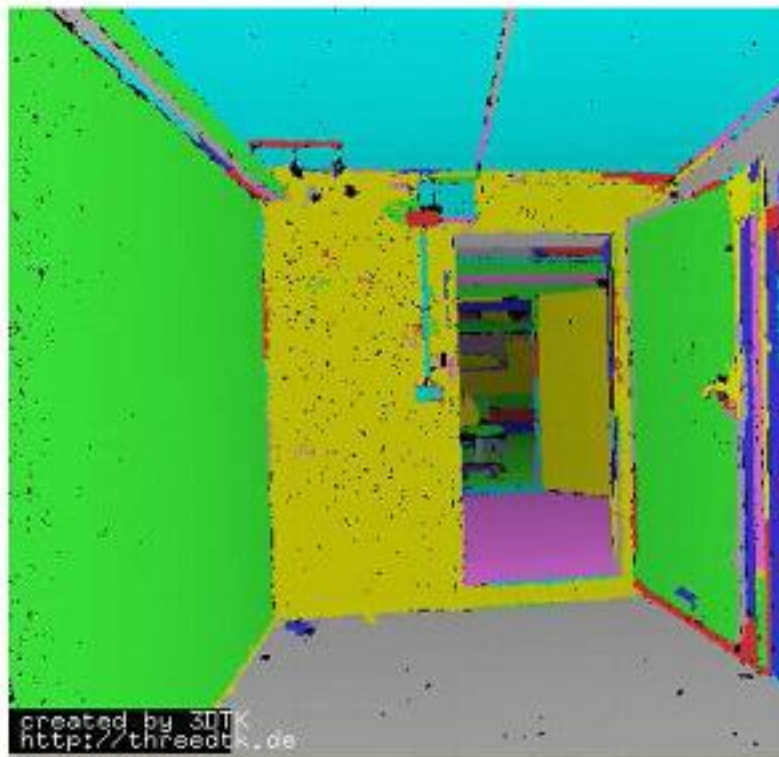


- Extract all points below a certain height

## Issues

- Data are noisy
  - Objects will also lose information
- Wall cannot be segmented out
- Ground is not always planar

# Plane segmentation





- Find a plane that minimize the average distance between a set of points and the surface
- Recursively merge the surface patches

## Issues

- Not every object is planar
  - Curved objects will be segmented into several segments



- Take a set of labeled examples
- Determine a rule that assign a label to any new example using the labeled examples
- Training dataset  $(\mathbf{x}_i, y_i)$ 
  - $\mathbf{x}_i$ : measurements of the properties of objects
  - $y_i$ : label
- Goal: given a new, plausible  $\mathbf{x}$ , assign it a label  $y$ .



$$p(k | \mathbf{x}) = \frac{p(\mathbf{x} | k) p(k)}{p(\mathbf{x})} \propto p(\mathbf{x} | k) p(k)$$

Given  $\mathbf{x}$

➤ Assign label  $k$  to  $\mathbf{x}$  if

➤  $p(k | \mathbf{x}) > p(i | \mathbf{x})$  for all  $i \neq k$  and  $p(k | \mathbf{x}) > \text{threshold}$

➤ Assign a random  $k$  label between  $k_1, \dots, k_j$  if

➤  $p(k_1 | \mathbf{x}) = \dots = p(k_j | \mathbf{x}) > p(i | \mathbf{x})$  for all  $i \neq k$

➤ Do not assign a label if

➤  $p(k | \mathbf{x}) > p(i | \mathbf{x})$  for all  $i \neq k$  and  $p(k | \mathbf{x}) \leq \text{threshold}$



Given  $x$

- Determine  $M$  training example that are nearest:  $x_1, \dots, x_M$
- Determine class  $k$  that has the largest representation  $n$  in the set
- Assign label  $k$  to  $x$  if  $n > \text{threshold}$
- Assign no label otherwise



Feature: a piece of information relevant for solving a computational task, e.g., locating an object in an image

- Raw data
- Histogram
- Pyramid of histograms
- Shape



- Compute a histogram of intensity or color
- Compute the correlation between example and test

## Issues

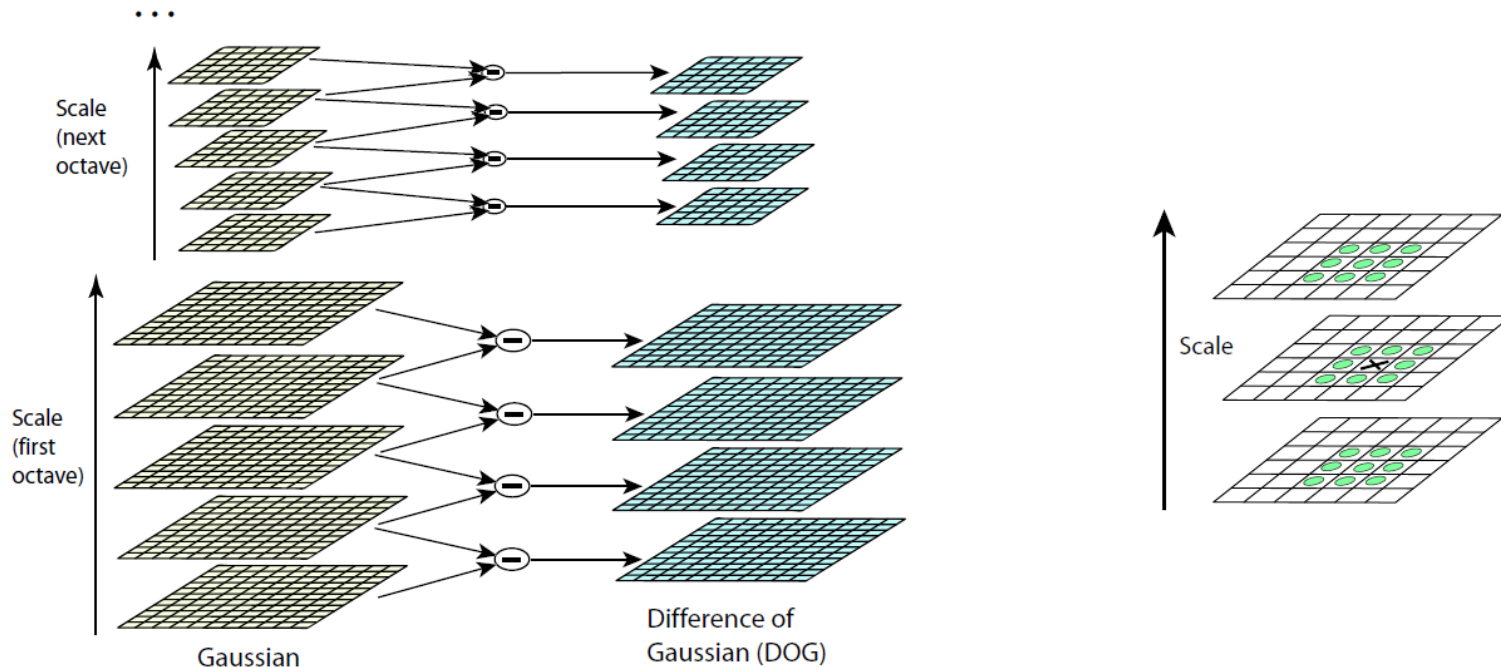
- Loss of the structural information
- Dimensionality



# Scale Invariant Feature Transform

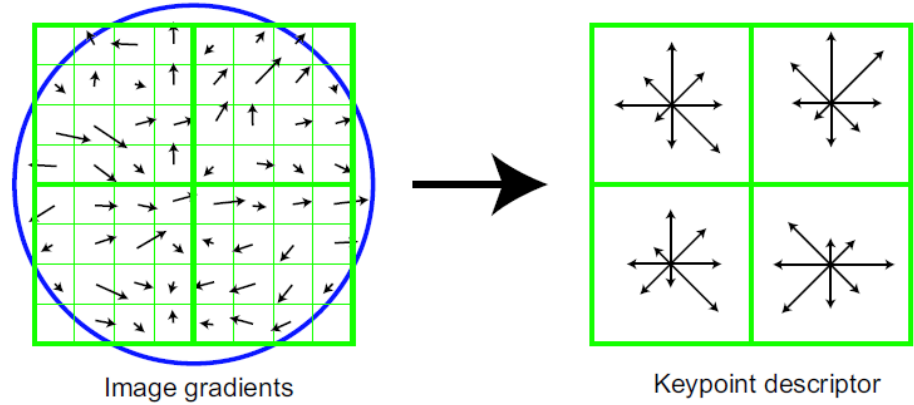


# Scale Invariant Feature Transform (SIFT)



- Identify locations and scales that are identifiable from different views of the same object
  - $L(x, y, \sigma) = G(x, y, \sigma) * I(x, y)$
  - $D(x, y, \sigma) = L(x, y, k\sigma) - L(x, y, \sigma)$
- Detect extrema (local minimum or maximum)

# Scale Invariant Feature Transform



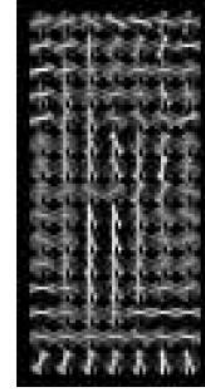
- Remove points of low contrast or poorly localized on an edge
- Orientation assignment

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

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

- Create a keypoint descriptor: 16 histograms (4x4 grid), each with 8 orientation bins, containing a total of 128 elements.

# Histogram of Oriented Gradient



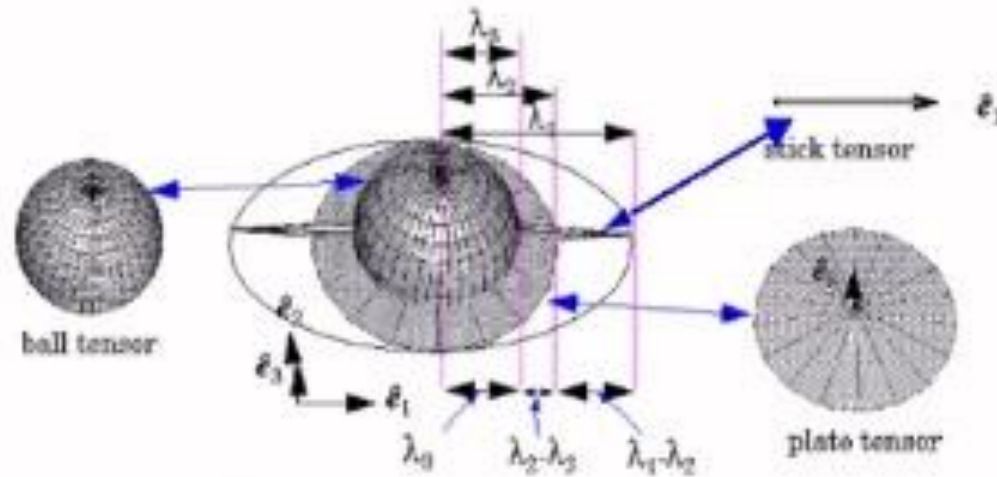
- Divide the image into small rectangular or radial cells
- Each cell accumulates a weighted local 1-D histogram of gradient directions over the pixels of the cell
- Normalize each cell by the energy over larger regions



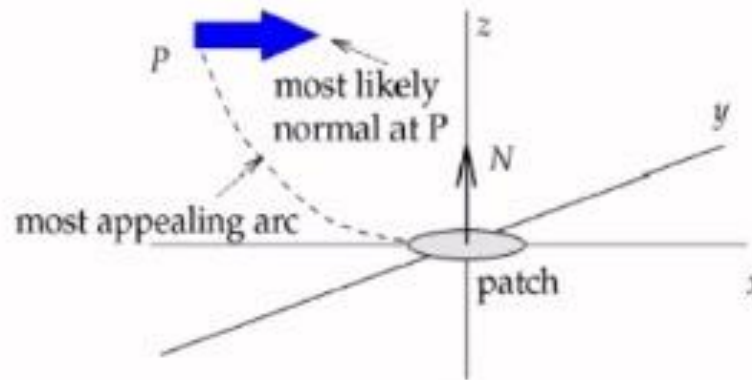
- Compute eigenvectors:  $\lambda_1, \lambda_2, \lambda_3$ 
  - Point/Spherical:  $\lambda_1 \approx \lambda_2 \approx \lambda_3$
  - Planar:  $\lambda_1 \approx \lambda_2 \gg \lambda_3$
  - Elongated:  $\lambda_1 \gg \lambda_2 \approx \lambda_3$

## Issues

- Many different objects have similar shape factor
- Shape factor of an object can depend on the point of view



- 2x2 or 3x3 matrix that captures both the orientation information and its confidence/saliency
  - Shape defines the type of information (point, surface, etc.)
  - Size represents the saliency
- Each token is first decomposed into the basis tensors, and then broadcasts its information to its neighbors.

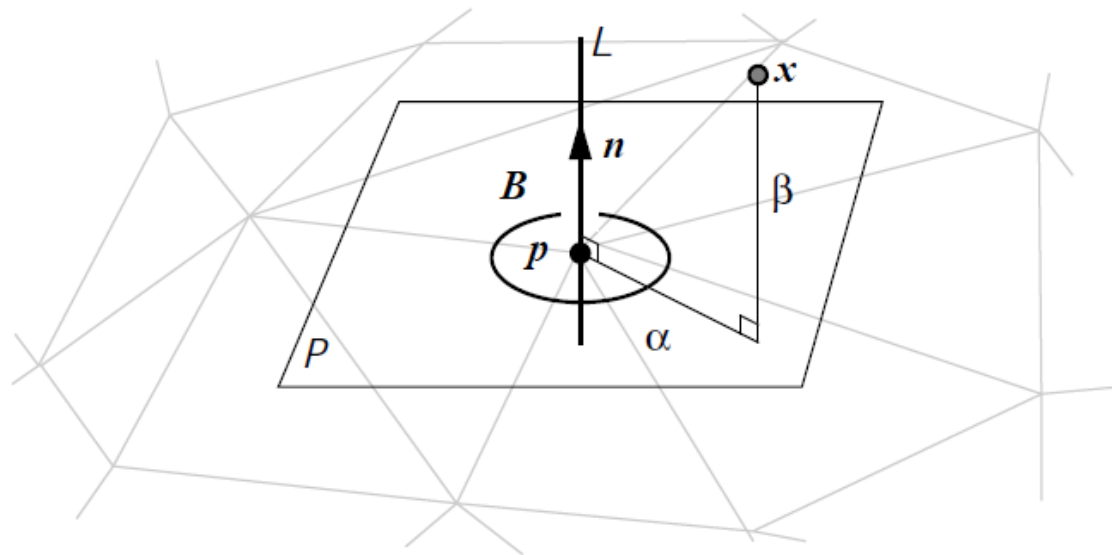


- The magnitude of the vote decays with distance and curvature:

$$V(d, \rho) = e^{-\frac{d^2 + c\rho^2}{\sigma^2}}$$

- $d$  is the distance along the smooth path
- $\rho$  is the curvature of the path
- $c$  controls the degree of decay
- $\sigma$  controls the size of the voting neighborhood
- Accumulate the votes by adding the matrices
- Analyze the tensor by eigen decomposition

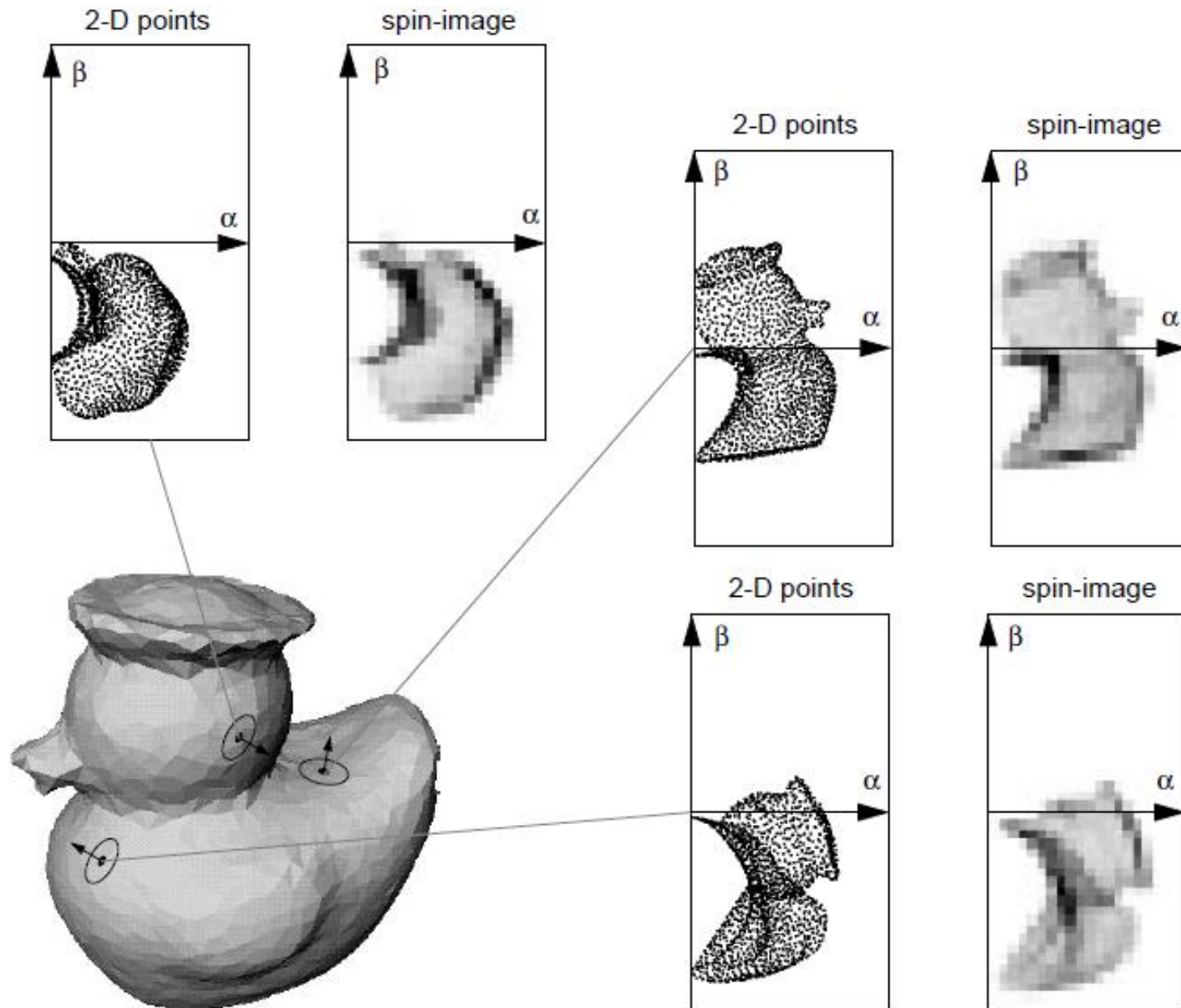
# Spin image



Johnson, A., Herbert, M., 1999. "Using Spin Images for Efficient Object Recognition in Cluttered 3D Scenes" IEEE Transactions on Pattern Analysis and Machine Intelligence, 21, (5).



# Spin image





- Collect a histogram of points
  - The resolution of the histogram
  - The size of the histogram
- To compare two spin images P and Q
  - Compute the correlation between two images

$$R(P, Q) = \frac{N \sum p_i q_i - \sum p_i \sum q_i}{\sqrt{(N \sum p_i^2 - (\sum p_i)^2)(N \sum q_i^2 - (\sum q_i)^2)}}$$

- Can also apply PCA, remove the mean spin image and compute the Euclidean norm