

# Lecture 12

## Segmentation

# Administrative

A2 is done

- Due Oct 28

A3 is out

- Due Nov 12 because Nov 11 is Veteran's Day

# Administrative

Recitation

- Ontologies

# Today's agenda

- Introduction to segmentation and clustering
- Gestalt theory for perceptual grouping
- Graph-based oversegmentation
- Agglomerative clustering

Reading:

Szeliski, 2<sup>nd</sup> edition, Chapter 7.5

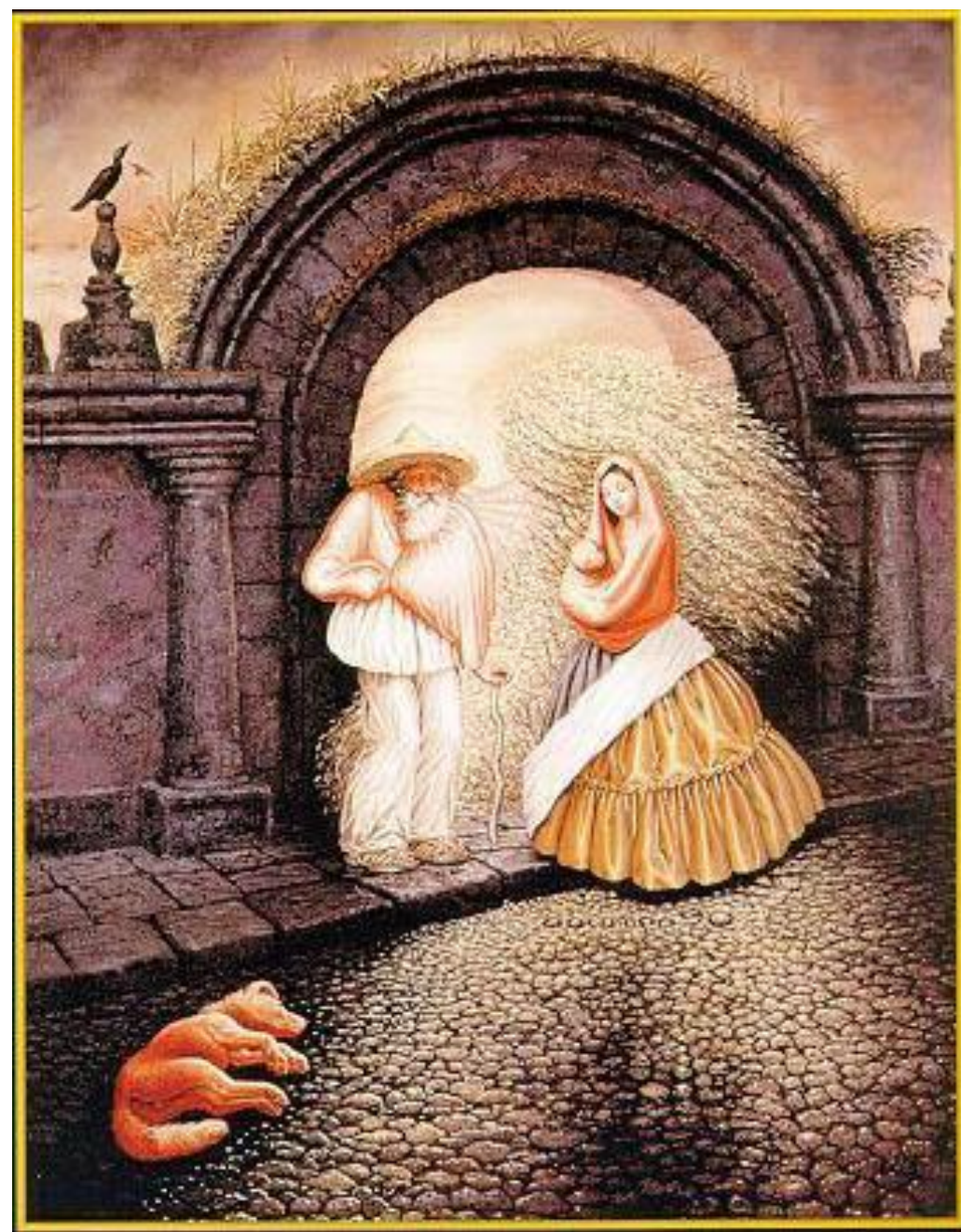
# Today's agenda

- Introduction to segmentation and clustering
- Gestalt theory for perceptual grouping
- Graph-based oversegmentation
- Agglomerative clustering

Reading:

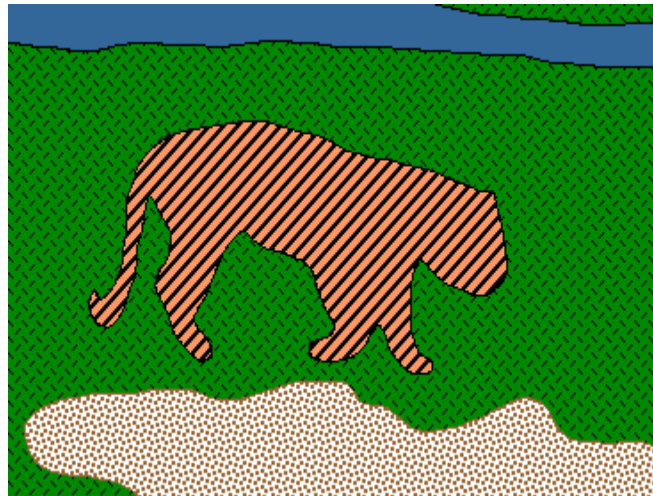
Szeliski, 2<sup>nd</sup> edition, Chapter 7.5

Q. What do you see?



# Image Segmentation

- Goal: identify groups of pixels that go together



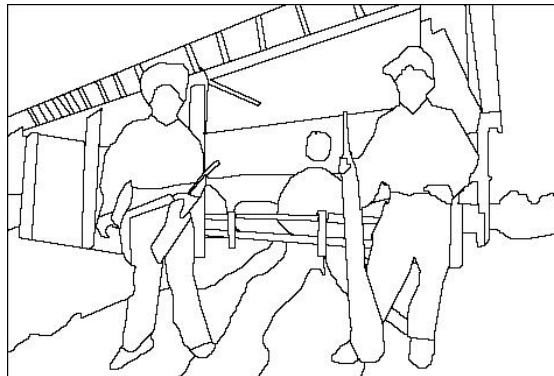
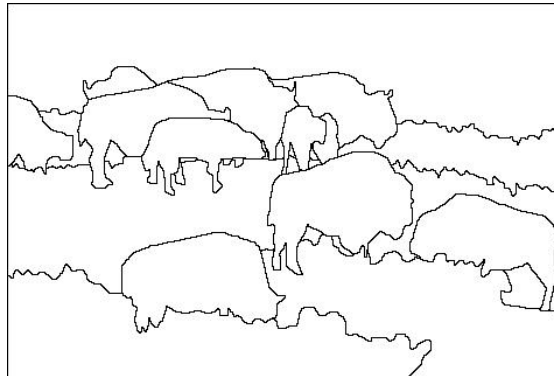
# The Goals of Segmentation

- Separate image into coherent “objects”

Image



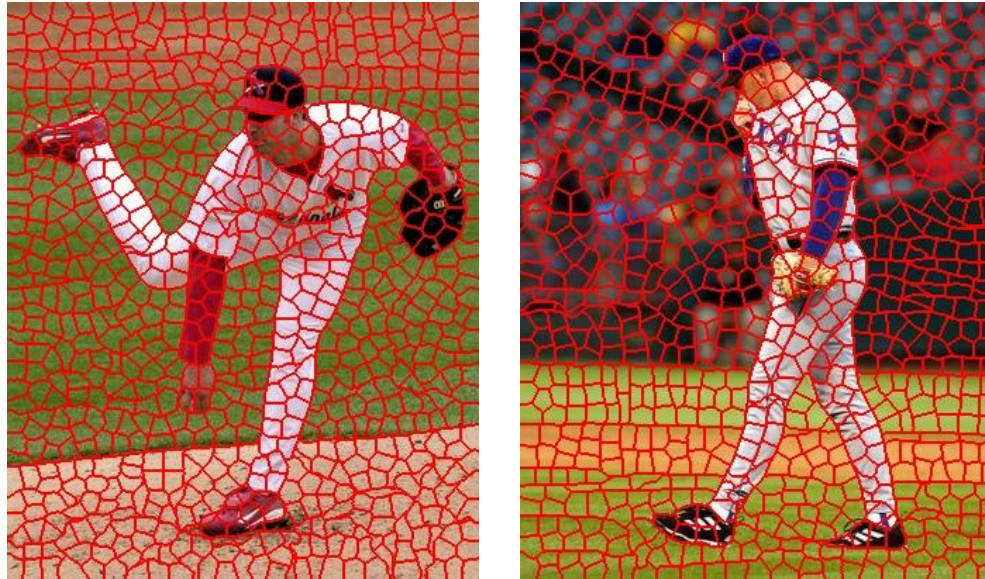
Human segmentation



# The Goals of Segmentation

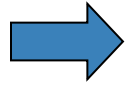
- Separate image into coherent “objects”
- Group together similar-looking pixels for efficiency of further processing

“superpixels”

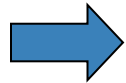


X. Ren and J. Malik. Learning a classification model for segmentation. ICCV 2003.

# Segmentation for efficiency



[Felzenszwalb and Huttenlocher 2004]



[Hoiem et al. 2005, Mori 2005]

[Shi and Malik 2001]

# Segmentation is used in Adobe photoshop to remove background



Rother et al. 2004

# Segment Anything [2023]

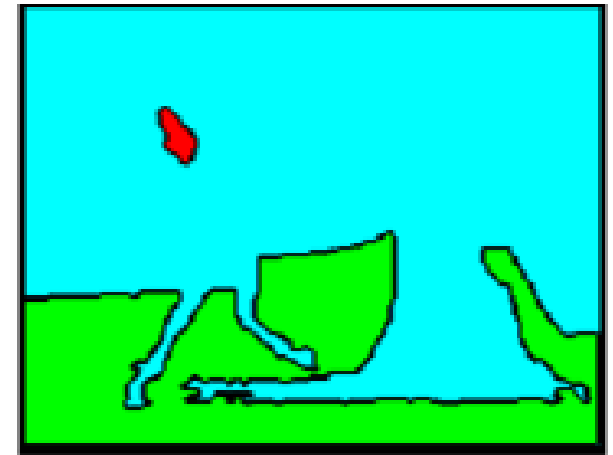
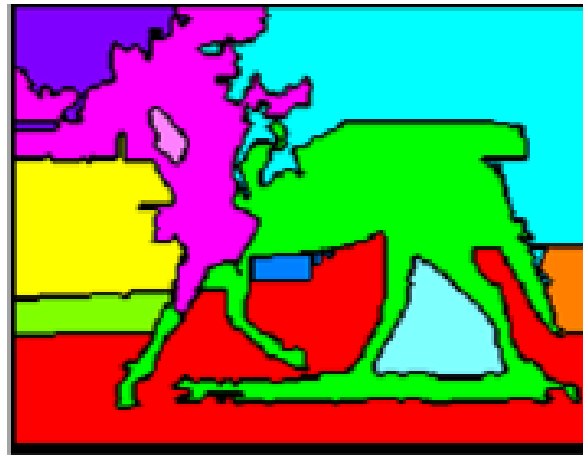
From Meta



# Levels of segmentations



Over-segmentation



Under-segmentation

# One way to think about “segmentation” is clustering

**Clustering:** group together similar data points and represent them with a single token

Key Challenges:

- 1) What makes two points/images/patches similar?
- 2) How do we compute an overall grouping from pairwise similarities?

# Why do we cluster?

- **Summarizing data**

- Look at large amounts of data
- Find clusters of pixels
- Represent each cluster of pixels with a HoG feature

- **Counting**

- Histograms of texture, color, SIFT vectors

- **Foreground-background separation**

- Separate the image into different regions

- **Prediction**

- Images in the same cluster may have the same labels

# How do we cluster?

- **Agglomerative clustering**

- Start with each point as its own cluster and iteratively merge the closest clusters

- **K-means**

- Iteratively re-assign points to the nearest cluster center

- **Mean-shift clustering**

- Estimate modes of pdf (probability density function)

# General ideas

- Tokens
  - Things that can be grouped together
  - (e.g. pixels, points, surface elements, etc., etc.)
- Bottom up clustering
  - tokens belong together because they are locally coherent
- Top down clustering
  - tokens belong together because they lie on the same visual entity (object, scene...)
- > These two are not mutually exclusive

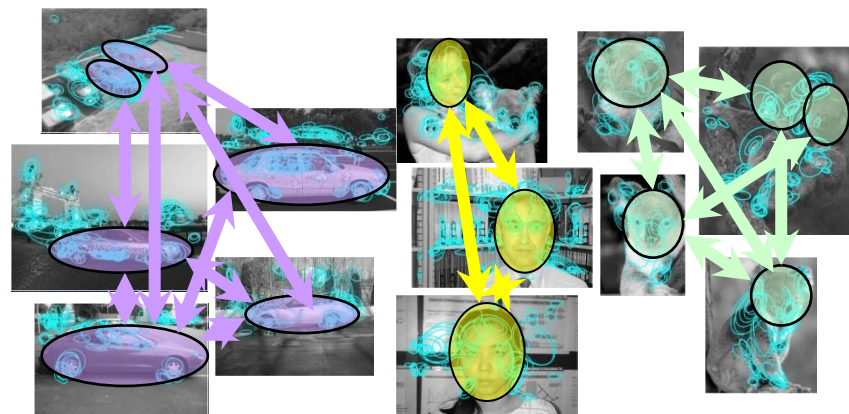
# Examples of Grouping in Vision



Determining image regions



Grouping video frames into shots



Object-level grouping



Figure-ground

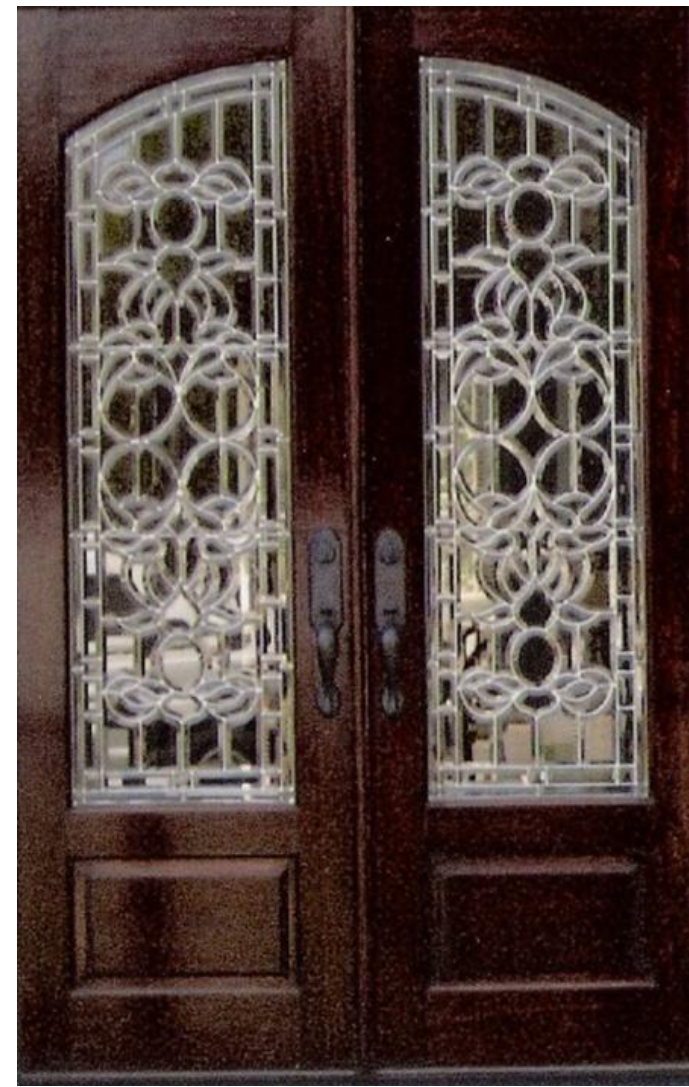
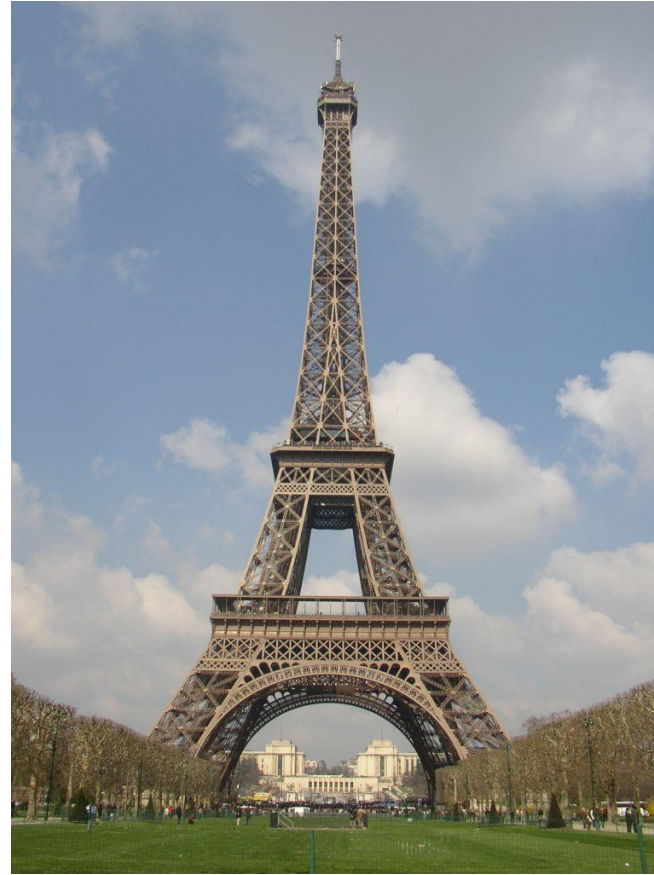
# Similarity



What things should  
be grouped?

What cues  
indicate groups?

# Symmetry



# Common Fate



Image credit: Arthus-Bertrand (via F. Durand)



(c) 2005 Heiko Burkhardt, iliano.com

# Proximity

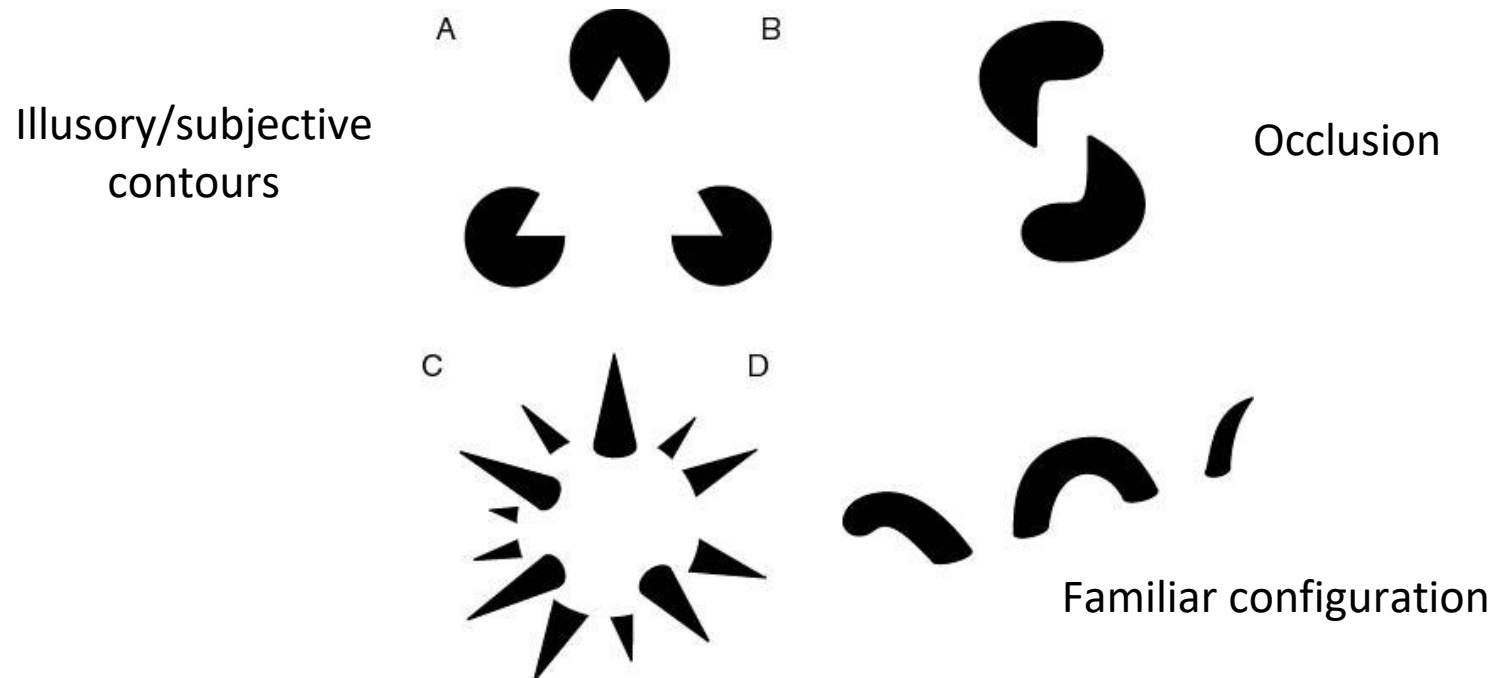


# What will we learn today?

- Introduction to segmentation and clustering
- **Gestalt theory for perceptual grouping**
- Graph-based oversegmentation
- Agglomerative clustering

# The Gestalt School

- Grouping is key to visual perception
- Elements in a collection can have properties that result from different **relationships (space, affordance, etc.)**
  - “The whole is greater than the sum of its parts”

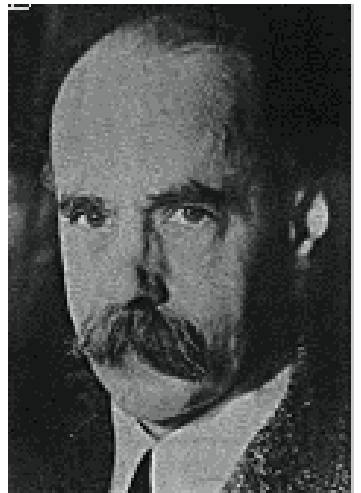


# Gestalt Theory

- Gestalt: whole or group
  - Whole is greater than sum of its parts
  - Relationships among parts can yield new properties/features
- Psychologists identified series of factors that predispose set of elements to be grouped (by human visual system)

***“I stand at the window and see a house, trees, sky. Theoretically I might say there were 327 brightnesses and nuances of colour. Do I have “327”? No. I have sky, house, and trees.”***

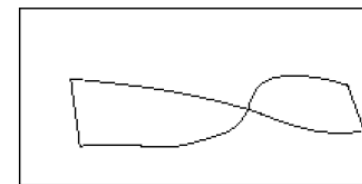
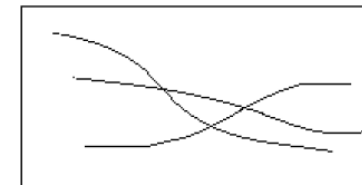
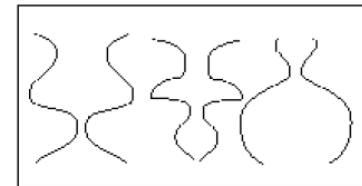
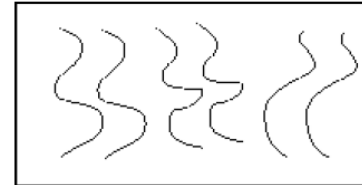
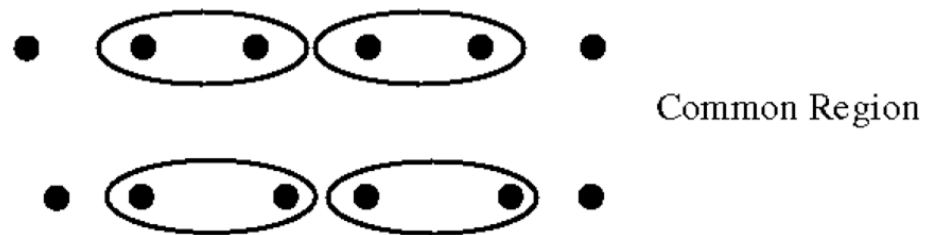
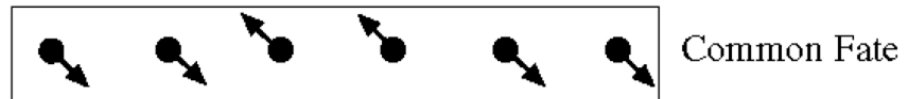
**Max Wertheimer  
(1880-1943)**



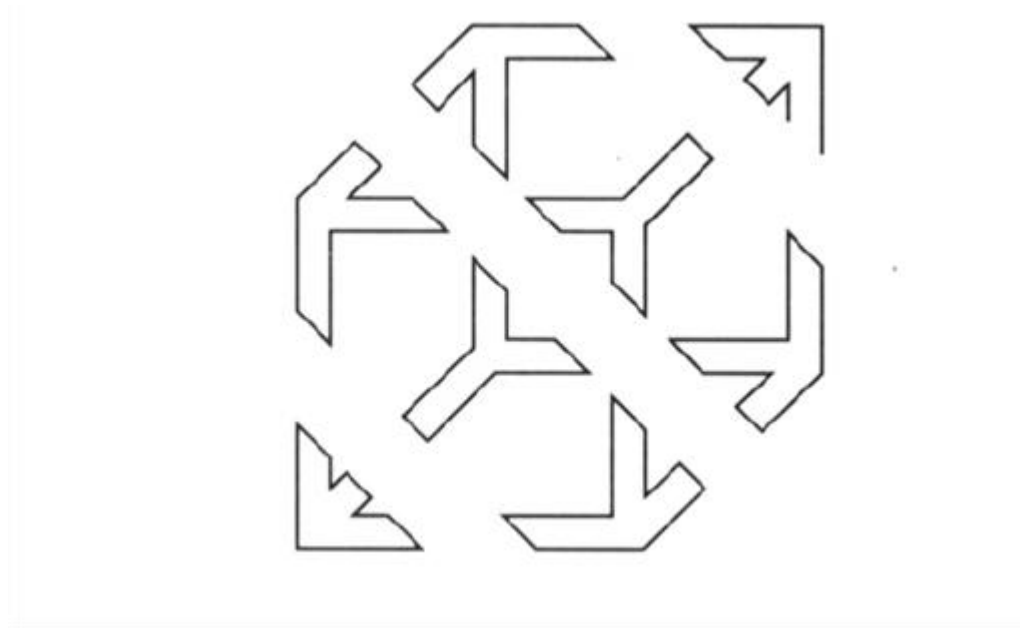
Untersuchungen zur Lehre von der Gestalt,  
*Psychologische Forschung*, Vol. 4, pp. 301-350, 1923  
<http://psy.ed.asu.edu/~classics/Wertheimer/Forms/forms.htm>

# Gestalt Factors

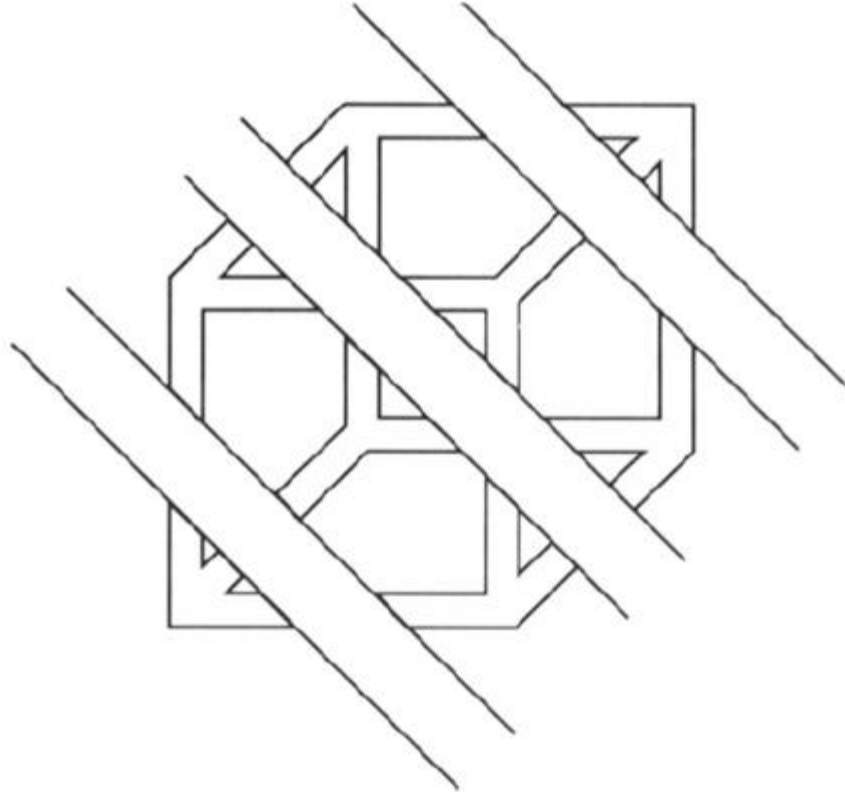
These factors make intuitive sense, but are very difficult to translate into algorithms.



# Continuity through Occlusion Cues

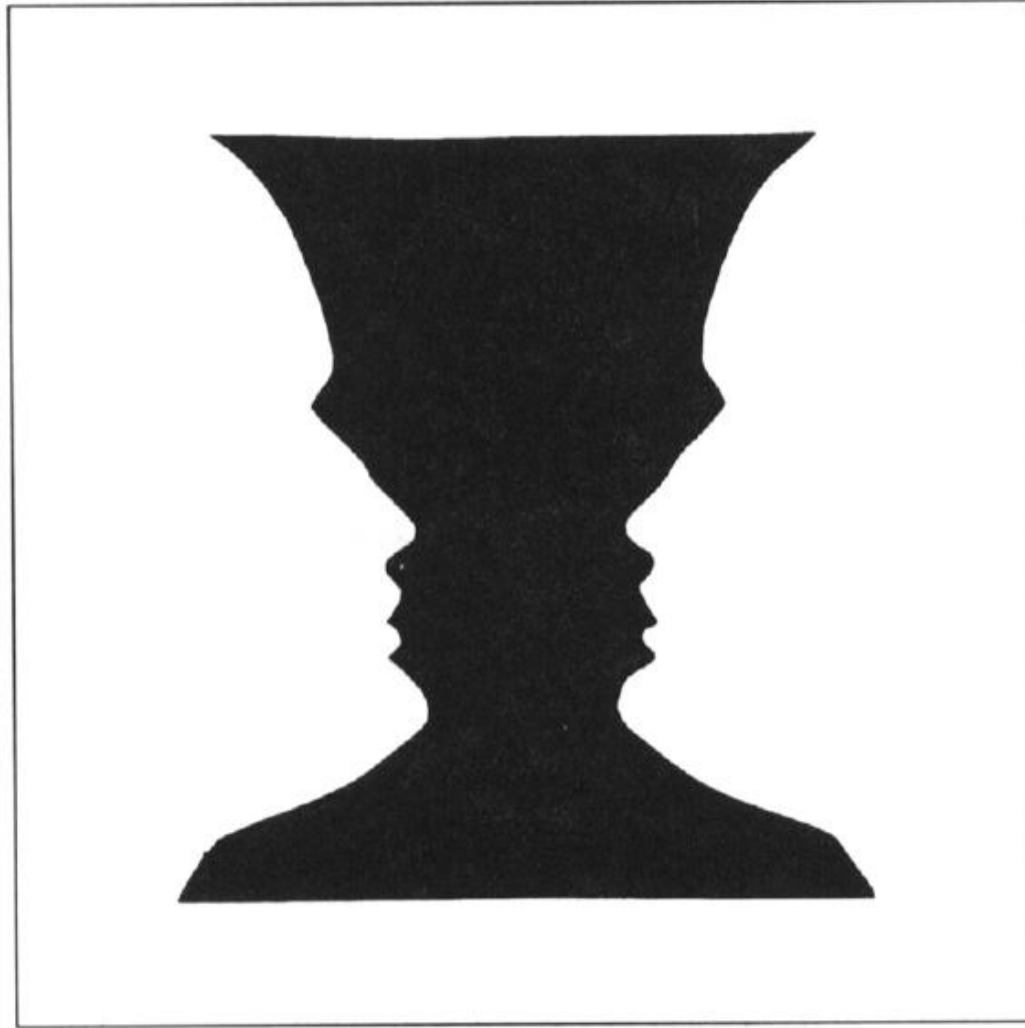


# Continuity through Occlusion Cues



Continuity, explanation by occlusion

# Figure-Ground Discrimination



# The Ultimate Gestalt?



# What will we learn today?

- Introduction to segmentation and clustering
- Gestalt theory for perceptual grouping
- **Graph-based oversegmentation**
- Agglomerative clustering

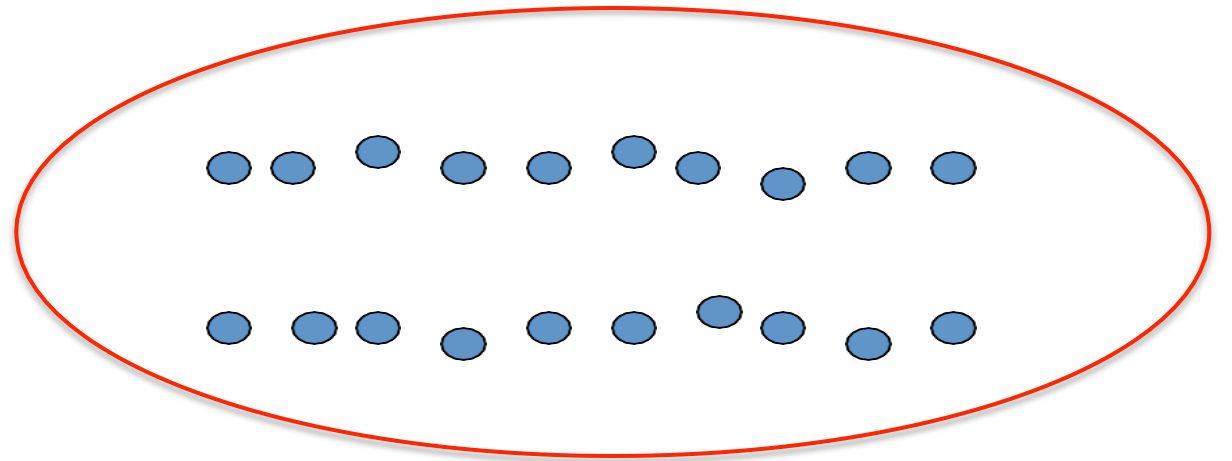
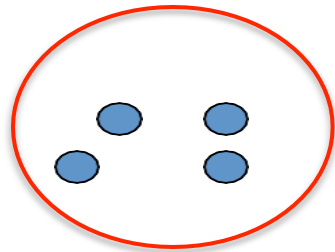
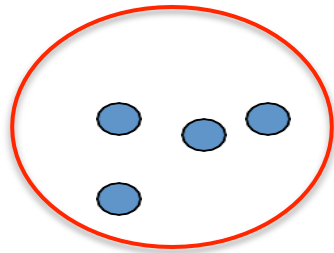


# Over-segmenting images

- Graph-based clustering for Image Segmentation
  - Introduced by *Felzenszwalb and Huttenlocher* in the paper titled *Efficient Graph-Based Image Segmentation*.



Imagine you have a set of pixels,  
how should you clustering them?

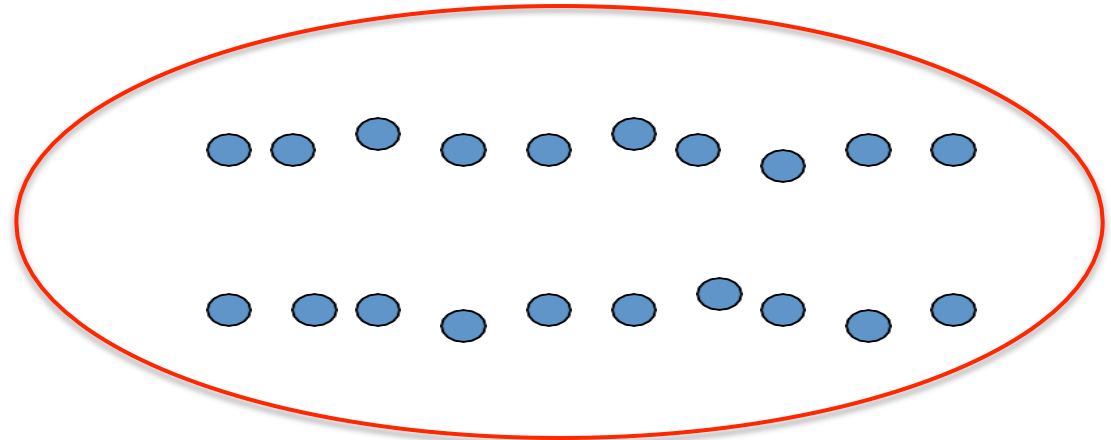
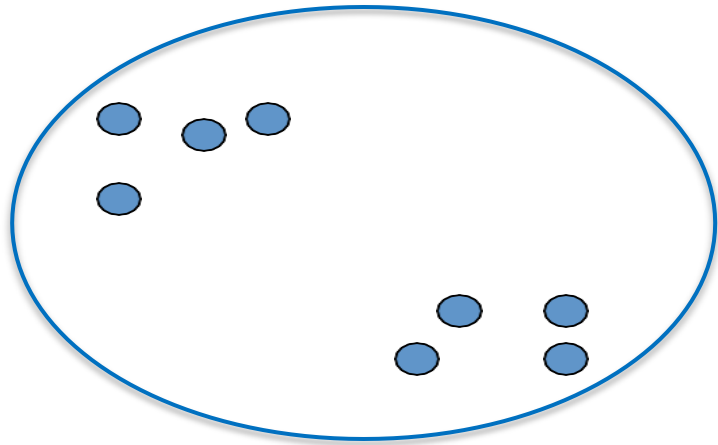


**Basic idea:** group together similar instances

Q. how do you measure similarity?

Q. do you need to measure similarity between every two pixels?

Imagine you have a set of pixels,  
how should you clustering them?

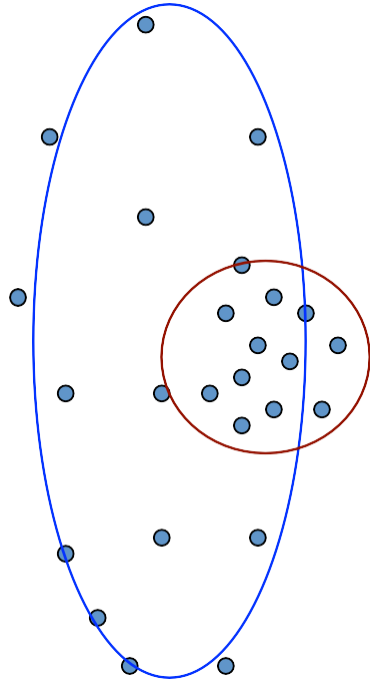


**Basic idea:** group together similar instances

Q. how do you measure similarity?

Q. do you need to measure similarity between every two pixels?

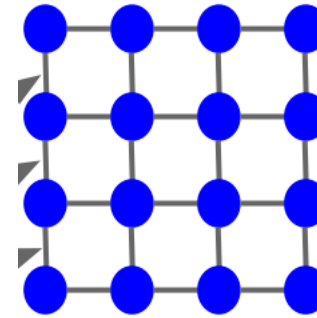
# Distances calculated using only (x,y) location of each pixel can be a bad idea



- Clusters may overlap
- Some clusters may be “wider” than others
- Distances can be deceiving!

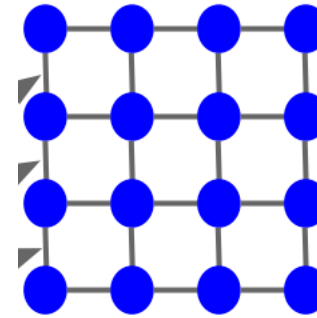
# Image as a Graph - Features and weights

- Every pixel is connected to its **8 neighboring pixels**
  - The edges between neighbors have **weights** that are determined by the distance between them.
  - Edge weights between pixels are determined using  **$\text{dist}(x, x')$**  distance in feature space.
    - where  **$x$**  and  **$x'$**  are two neighboring pixels
- 
- **Q. What is a good feature space?**

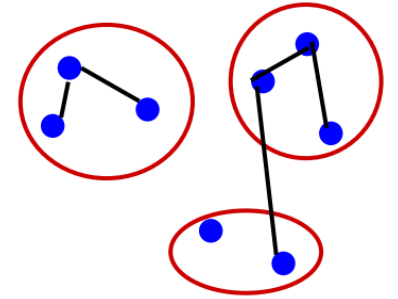


# What are good pixel features?

- Use **RGB values**?
  - $v = [r, g, b]$
  - It is 3-dimensional
- Use **location**?
  - $v = [x, y]$
  - 2-dim
- Use RGB + location?
  - $v = [x, y, r, g, b]$
  - 5-dim
- Use **gradient magnitude**?
  - $v = [df/dx, df/dy]$
  - 2-d



# Problem Formulation



- Graph  $G = (V, E)$
  - $V$  is set of nodes (i.e. pixels)
  - $E$  is a set of undirected edges between pairs of pixels
  - $\text{dist}(v_i, v_j)$  is the weight/distance of the edge between nodes  $v_i$  and  $v_j$ .
- 
- $S$  is a segmentation of a graph  $G$  such that  $G' = (V, E')$  where  $E' \subset E$ .
    - That is, **we keep all vertices**, but **select a subset  $E'$**  from all initial edges  $E$ .
  - $S$  divides  $G$  into  $G'$  such that it contains distinct clusters  $C$ .

# Weights of edges: distance measure

Clustering is an unsupervised learning method. Given items  $v_1, v_2, \dots, v_n \in \mathcal{R}^D$ , the goal is to group them into clusters.

We need a pairwise **distance/similarity function** between items, and sometimes the desired **number of clusters**.

When data (e.g. images, objects, documents) are represented by feature vectors, commonly used measures are:

- *Euclidean distance*.
- *Cosine similarity*.

# Defining Distance Measures

Let  $x$  and  $x'$  be two objects from the universe of possible objects.  
The distance (or similarity) between  $x$  and  $x'$  is a real number:

- The Euclidean distance is defined as  $dist(v_1, v_2) = \sqrt{\sum_i (v_{1i} - v_{2i})^2}$

- In contrast, the cosine similarity measure would be

$$\begin{aligned} dist(v_1, v_2) &= 1 - \cos(v_1, v_2) \\ &= 1 - \frac{v_1^T v_2}{||v_1|| \cdot ||v_2||} \end{aligned}$$

# The algorithm

The input is a graph  $G = (V, E)$ , with  $n$  vertices and  $m$  edges. The output is a segmentation of  $V$  into components  $S = (C_1, \dots, C_r)$ .

0. Sort  $E$  into  $\pi = (o_1, \dots, o_m)$ , by non-decreasing edge weight.
1. Start with a segmentation  $S^0$ , where each vertex  $v_i$  is in its own component.
2. Repeat step 3 for  $q = 1, \dots, m$ .
3. Construct  $S^q$  given  $S^{q-1}$  as follows. Let  $v_i$  and  $v_j$  denote the vertices connected by the  $q$ -th edge in the ordering, i.e.,  $o_q = (v_i, v_j)$ . If  $v_i$  and  $v_j$  are in disjoint components of  $S^{q-1}$  and  $w(o_q)$  is small compared to the internal difference of both those components, then merge the two components otherwise do nothing. More formally, let  $C_i^{q-1}$  be the component of  $S^{q-1}$  containing  $v_i$  and  $C_j^{q-1}$  the component containing  $v_j$ . If  $C_i^{q-1} \neq C_j^{q-1}$  and  $w(o_q) \leq MInt(C_i^{q-1}, C_j^{q-1})$  then  $S^q$  is obtained from  $S^{q-1}$  by merging  $C_i^{q-1}$  and  $C_j^{q-1}$ . Otherwise  $S^q = S^{q-1}$ .
4. Return  $S = S^m$ .

# Some results

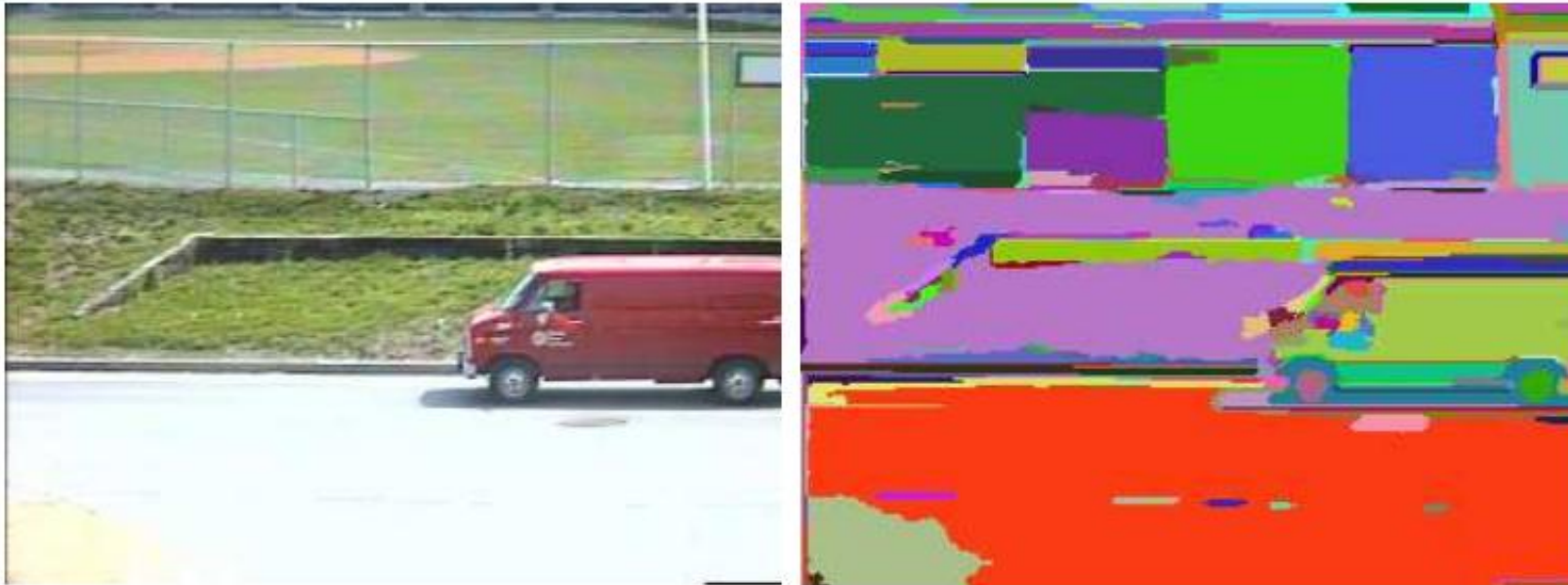


Figure 2: A street scene ( $320 \times 240$  color image), and the segmentation results produced by our algorithm ( $\sigma = 0.8$ ,  $k = 300$ ).

# More

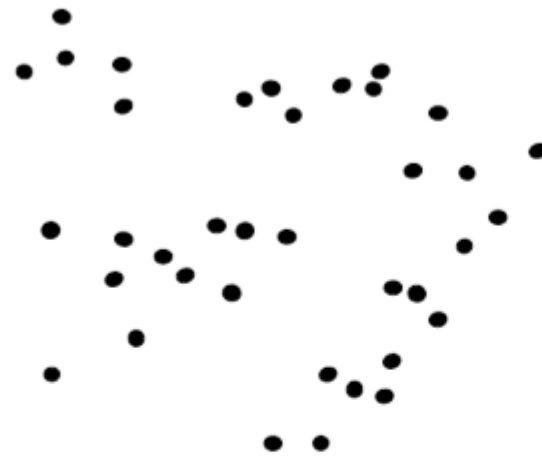


Figure 4: An indoor scene (image  $320 \times 240$ , color), and the segmentation results produced by our algorithm ( $\sigma = 0.8$ ,  $k = 300$ ).

# What will we learn today?

- Introduction to segmentation and clustering
- Gestalt theory for perceptual grouping
- Graph-based oversegmentation
- Agglomerative clustering

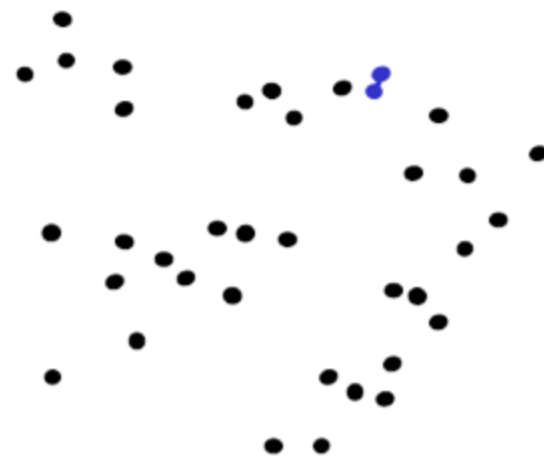
# Agglomerative clustering



1. Say "Every point is its own cluster"

Slide credit: Andrew Moore

# Agglomerative clustering

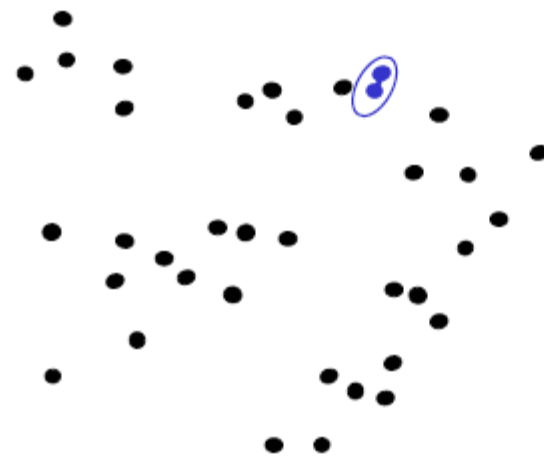


1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters



Slide credit: Andrew Moore

# Agglomerative clustering

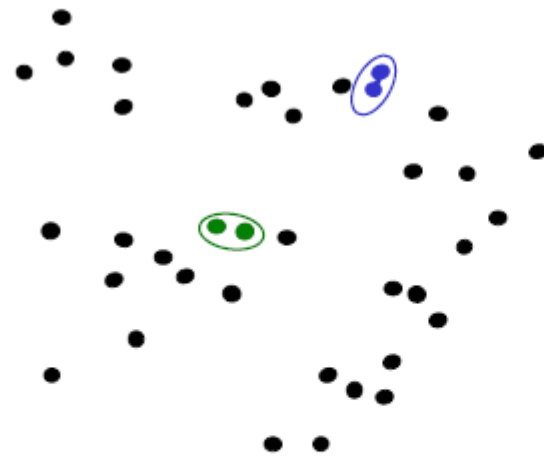


1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters
3. Merge it into a parent cluster



Slide credit: Andrew Moore

# Agglomerative clustering

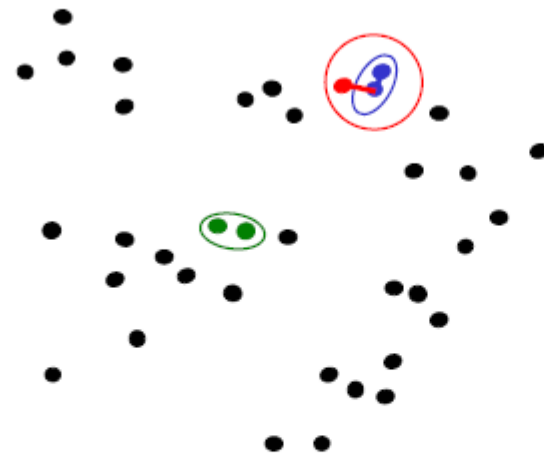


1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters
3. Merge it into a parent cluster
4. Repeat



Slide credit: Andrew Moore

# Agglomerative clustering



1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters
3. Merge it into a parent cluster
4. Repeat

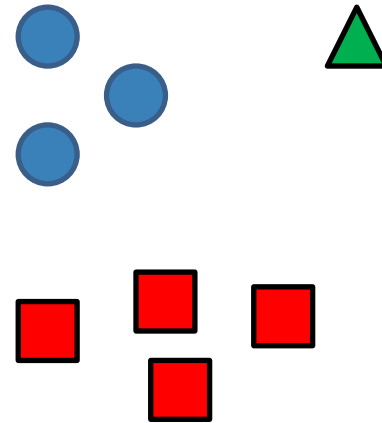


Slide credit: Andrew Moore

# Agglomerative clustering

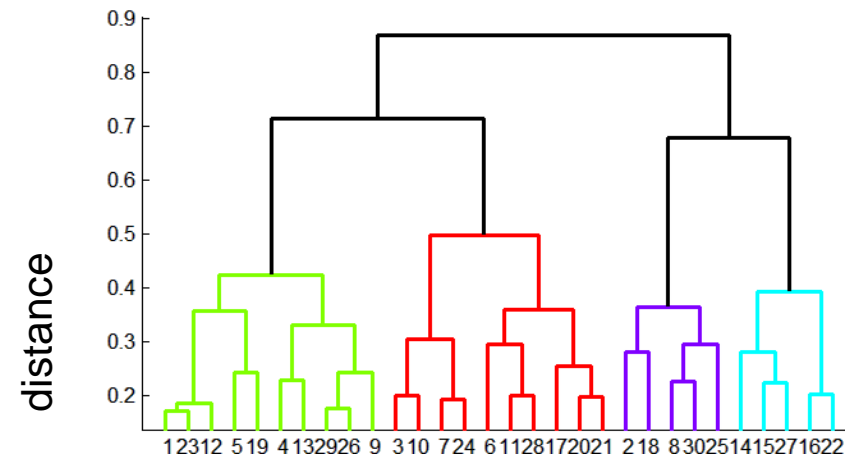
## How to define cluster similarity?

- Average distance between all pixels between the two cluster?
- Maximum distance?
- Minimum distance?
- Distance between means?



## How many clusters?

- Clustering creates a dendrogram (a tree)
- Threshold based on max number of clusters or based on distance between merges

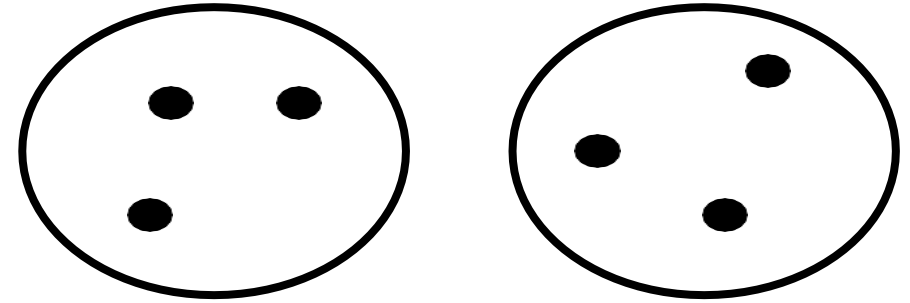


# Agglomerative Hierarchical Clustering - Algorithm

## Inputs:

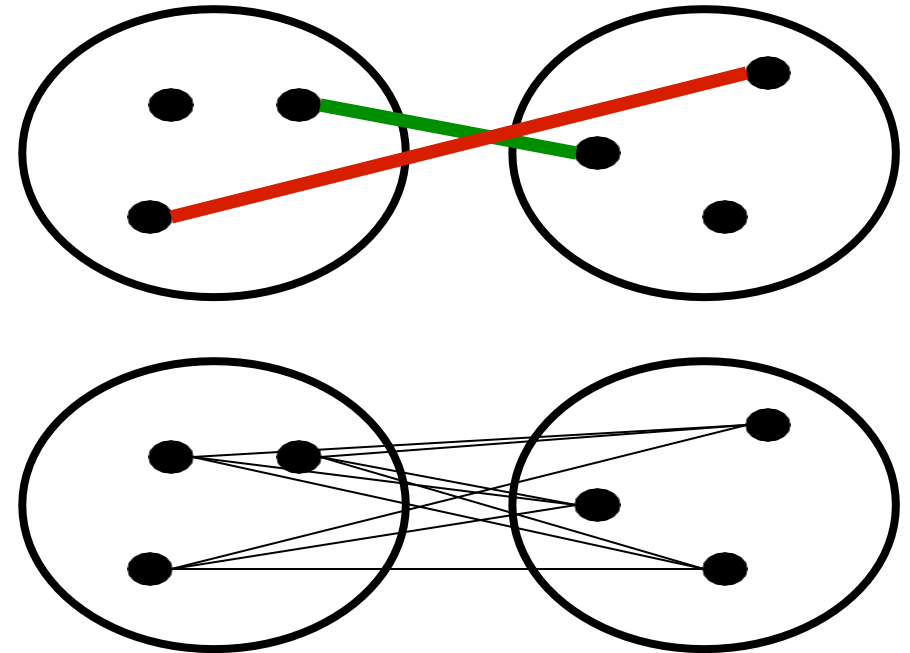
- An input image
  - Feature representation for each pixel
  - Distance metric  $\text{dist}(-,-)$
- 
- Initially, each pixel  $v_1, \dots, v_n$  is its own cluster  $C_1, \dots, C_n$
  - While True:
    - Find two nearest clusters according to  $\text{dist}(C_i, C_j)$
    - Merge  $C = (C_i, C_j)$
    - If only 1 cluster is left:
      - break

How should we define “closest” for clusters with multiple pixels already in it?



# How should we define “closest” for clusters with multiple pixels already in it?

- Closest pair  
(single-link clustering)
- Farthest pair  
(complete-link clustering)
- Average of all pairs

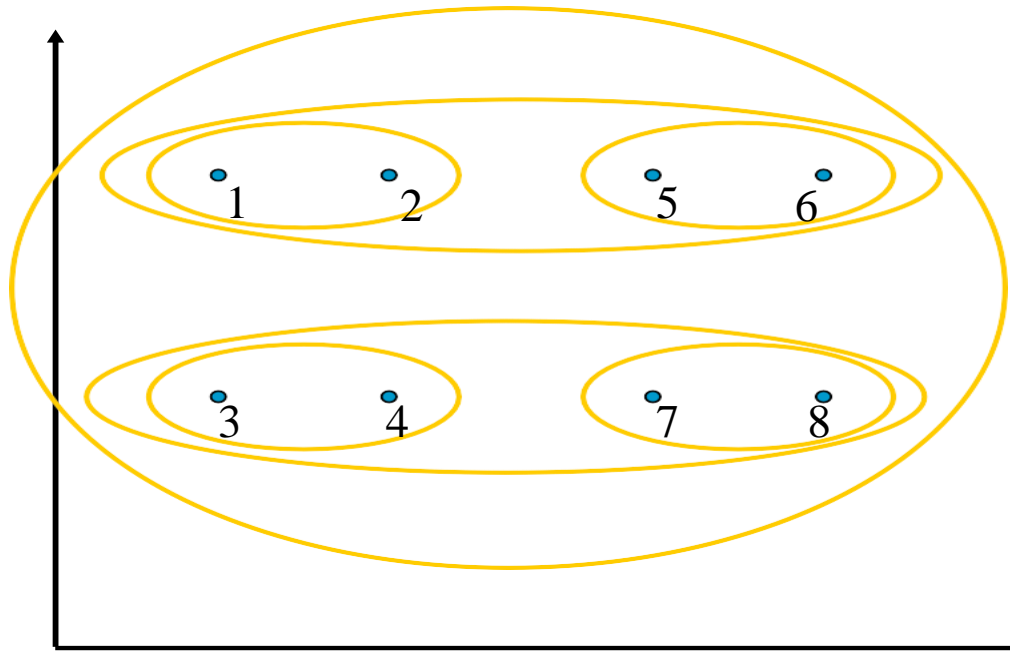


Different choices create different clustering behaviors

# How should we define “closest” for clusters with multiple pixels already in it?

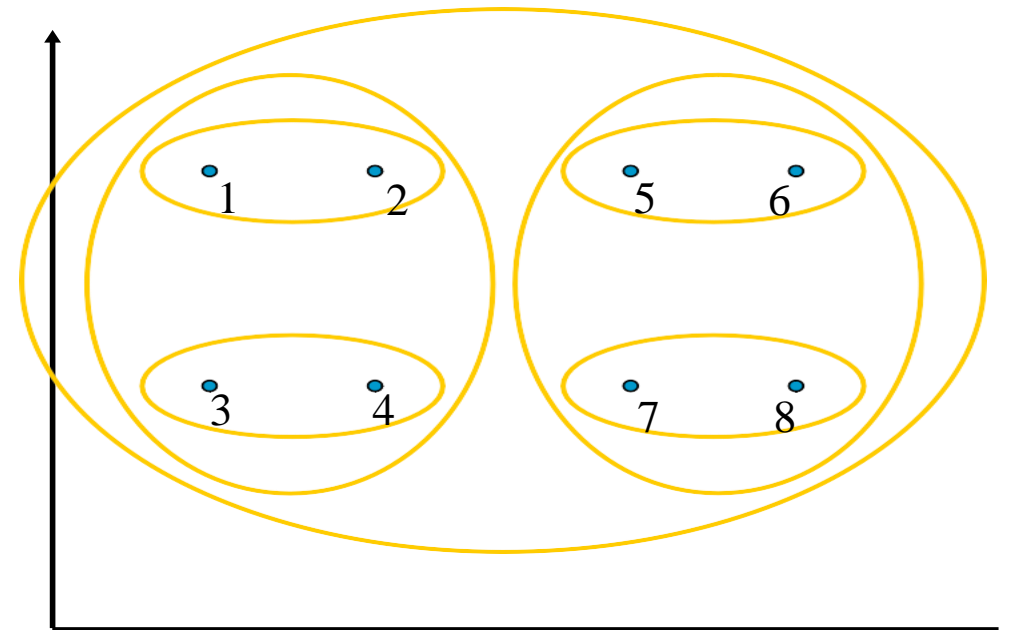
Closest pair

(single-link clustering)



Farthest pair

(complete-link clustering)

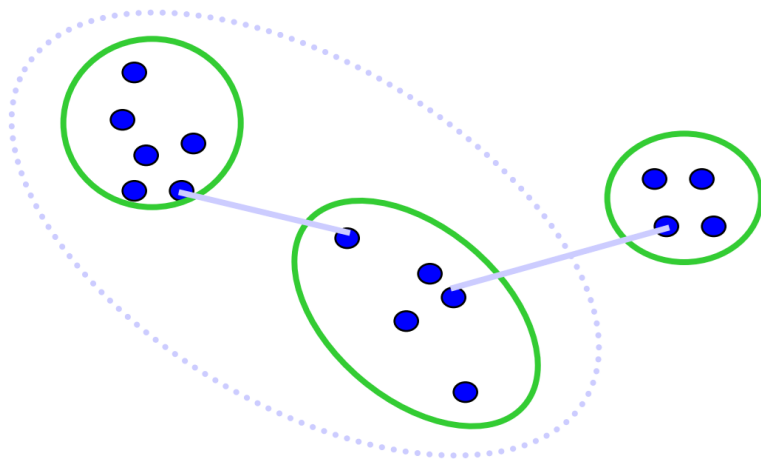


[Pictures from Thorsten Joachims]

## Single Linkage distance measure

$$\text{dist}(C_i, C_j) = \min_{v_i \in C_i, v_j \in C_j, (C_i, C_j) \in E} \text{dist}(v_i, v_j)$$

Connects the clusters based on the distance of their closest pixels  
It produces “long” clusters.

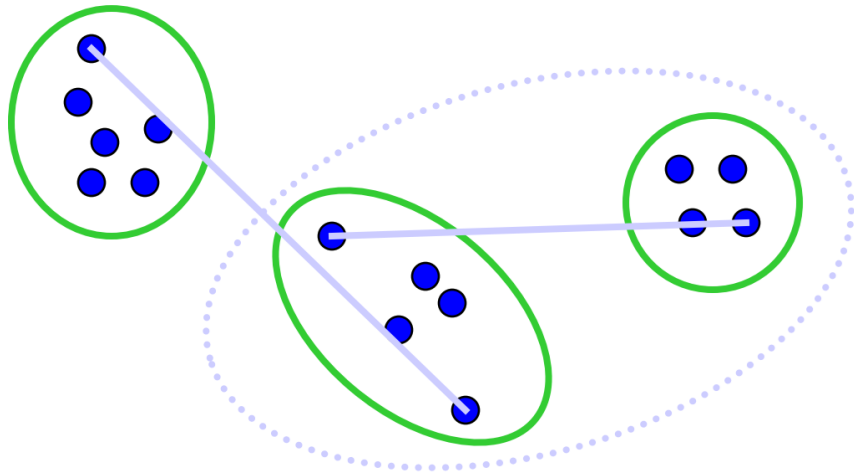


Long, skinny clusters

## Complete Link distance measure

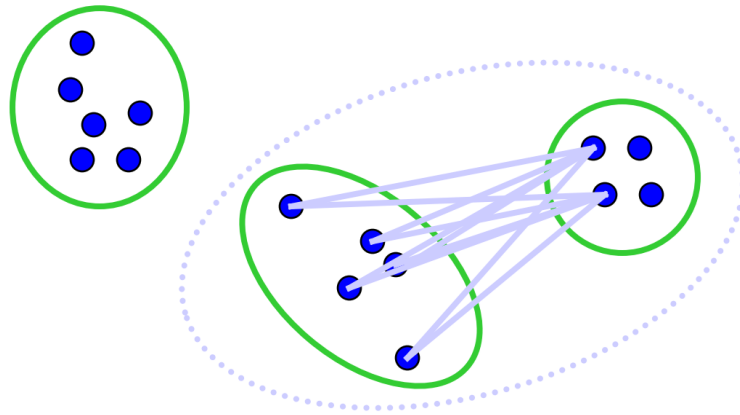
$$\text{dist}(C_i, C_j) = \max_{v_i \in C_i, v_j \in C_j, (C_i, C_j) \in E} \text{dist}(v_i, v_j)$$

Produces compact clusters that are similar in diameter



## Average Link distance measures

$$\text{dist}(C_i, C_j) = \frac{\sum_{v_i \in C_i, v_j \in C_j, (C_i, C_j) \in E} \text{dist}(v_i, v_j)}{|C_i||C_j|}$$



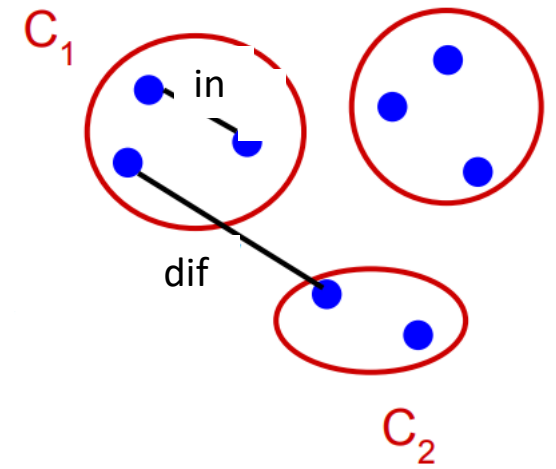
Robust against noise.

# Inlier-outlier linkage distance measure

$$\text{Merge}(C_1, C_2) = \begin{cases} \text{True} & \text{if } \text{dif}(C_1, C_2) < \text{in}(C_1, C_2) \\ \text{False} & \text{otherwise} \end{cases}$$

Where

- $\text{dif}(C_1, C_2)$  is the difference between two clusters.
- $\text{in}(C_1, C_2)$  is the internal difference in the clusters  $C_1$  and  $C_2$



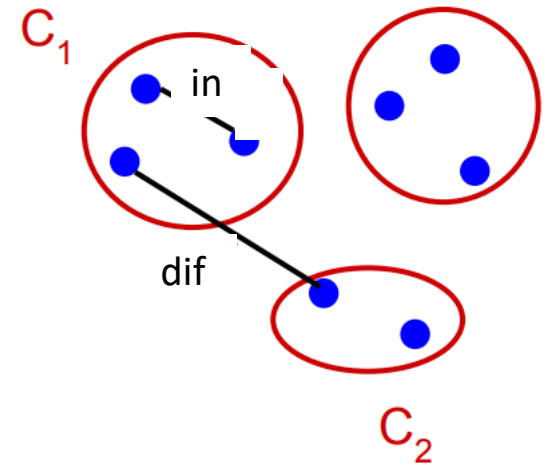
# Inlier-outlier linkage distance measure

$$\text{Merge}(C_1, C_2) = \begin{cases} \text{True} & \text{if } \text{dif}(C_1, C_2) < \text{in}(C_1, C_2) \\ \text{False} & \text{otherwise} \end{cases}$$

$$\text{dif}(C_i, C_j) = \min_{v_i \in C_i, v_j \in C_j, (C_i, C_j) \in E} \text{dist}(v_i, v_j)$$

Where

- $\text{dif}(C_1, C_2)$  is the difference between two clusters.
- $\text{in}(C_1, C_2)$  is the internal difference in the clusters  $C_1$  and  $C_2$



# Inlier-outlier linkage distance measure

$$\text{Merge}(C_1, C_2) = \begin{cases} \text{True} & \text{if } \text{dif}(C_1, C_2) < \text{in}(C_1, C_2) \\ \text{False} & \text{otherwise} \end{cases}$$

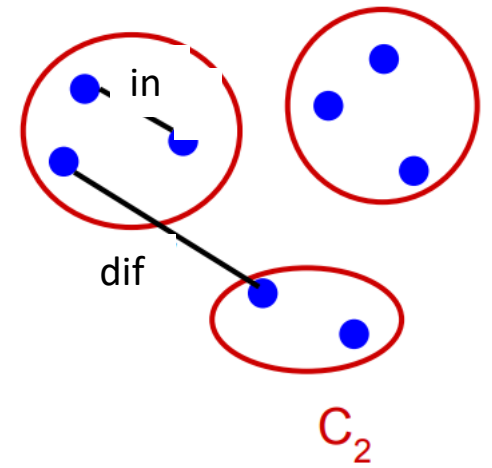
$$\text{dif}(C_i, C_j) = \min_{v_i \in C_i, v_j \in C_j, (C_i, C_j) \in E} \text{dist}(v_i, v_j)$$

Max weight edge  
in one cluster

$$\text{in}(C_i, C_j) = \min_{C \in \{C_i, C_j\}} \left[ \max_{v_i, v_j \in C} \left[ \text{dist}(v_i, v_j) + \frac{k}{|C|} \right] \right]$$

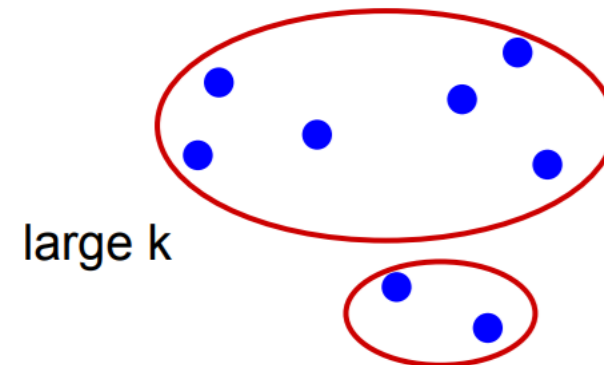
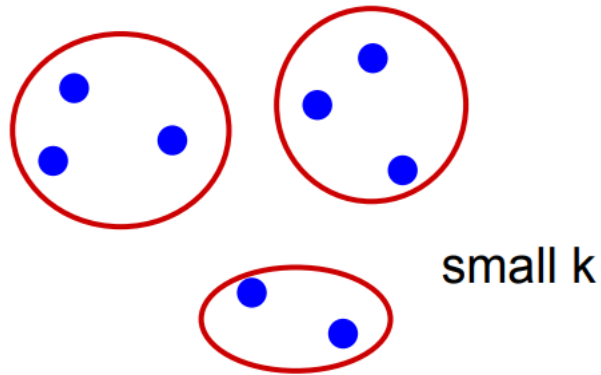
Where

- $\text{dif}(C_1, C_2)$  is the difference between two clusters.
- $\text{in}(C_1, C_2)$  is the internal difference in the clusters  $C_1$  and  $C_2$

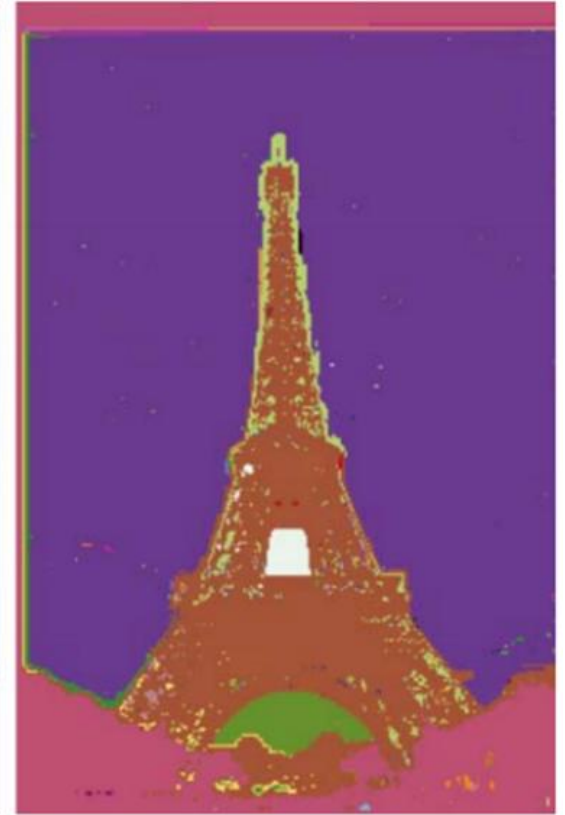


# inlier-outlier linkage for Segmentation

- $k/|C|$  sets the threshold by which the clusters need to be different from the internal pixels in a cluster.
- Effect of  $k$ :
  - **If  $k$  is large, it causes a preference for larger objects.**

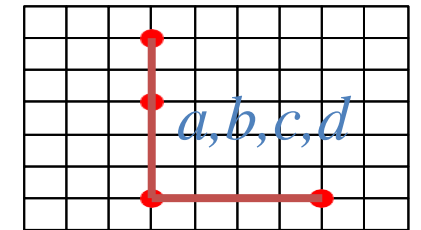
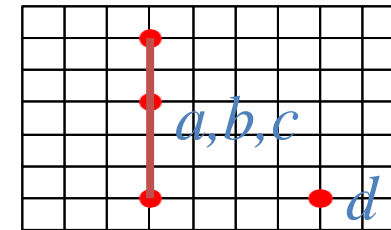
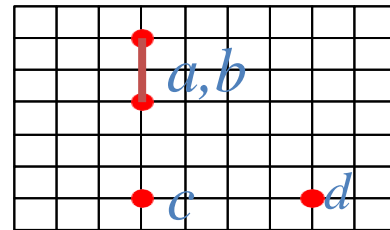
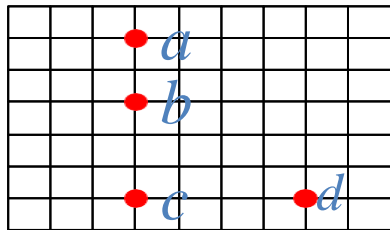


# Results



# How to implement single-linkage efficiently

Euclidean Distance



(1)  
*d*

(2)

(3)

	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>		5	6
2			
<i>b</i>		3	5
<i>c</i>			4

	<i>c</i>	<i>d</i>
<i>a, b</i>	3	5
<i>c</i>		4

	<i>d</i>
<i>a, b, c</i>	4

Distance Matrix

# Conclusions: Agglomerative Clustering

## Pros:

- Simple to implement, widespread application.
- Clusters have adaptive shapes.
- Provides a hierarchy of clusters.
- No need to specify number of clusters **in advance**.

## Cons:

- May have imbalanced clusters.
- Still have to choose number of clusters eventually for an application
- Does not scale well. Runtime of  $O(n^3)$ .
- Can get stuck at a local optima.

# Today's agenda

- Introduction to segmentation and clustering
- Gestalt theory for perceptual grouping
- Graph-based oversegmentation
- Agglomerative clustering

# Next time

K-means and mean shift