

# Gene Expression

- Process of transcription and/or translation of a gene is called **gene expression**.
- Every cell of an organism has the same genetic material, but different genes are **expressed** at different times.
- Patterns of gene expression in a cell is indicative of its state.

# Hybridization

- If two complementary strands of DNA or mRNA are brought together, under appropriate experimental conditions they will **hybridize**.
- **A hybridizes** to **B**  $\Rightarrow$ 
  - **A** is reverse complementary to **B**, or
  - **A** is reverse complementary to a subsequence of **B**.
- It is possible to experimentally verify whether **A** hybridizes to **B**, by **labeling A** or **B** with a radioactive or fluorescent tag, followed by excitation by laser.

# Measuring gene expression

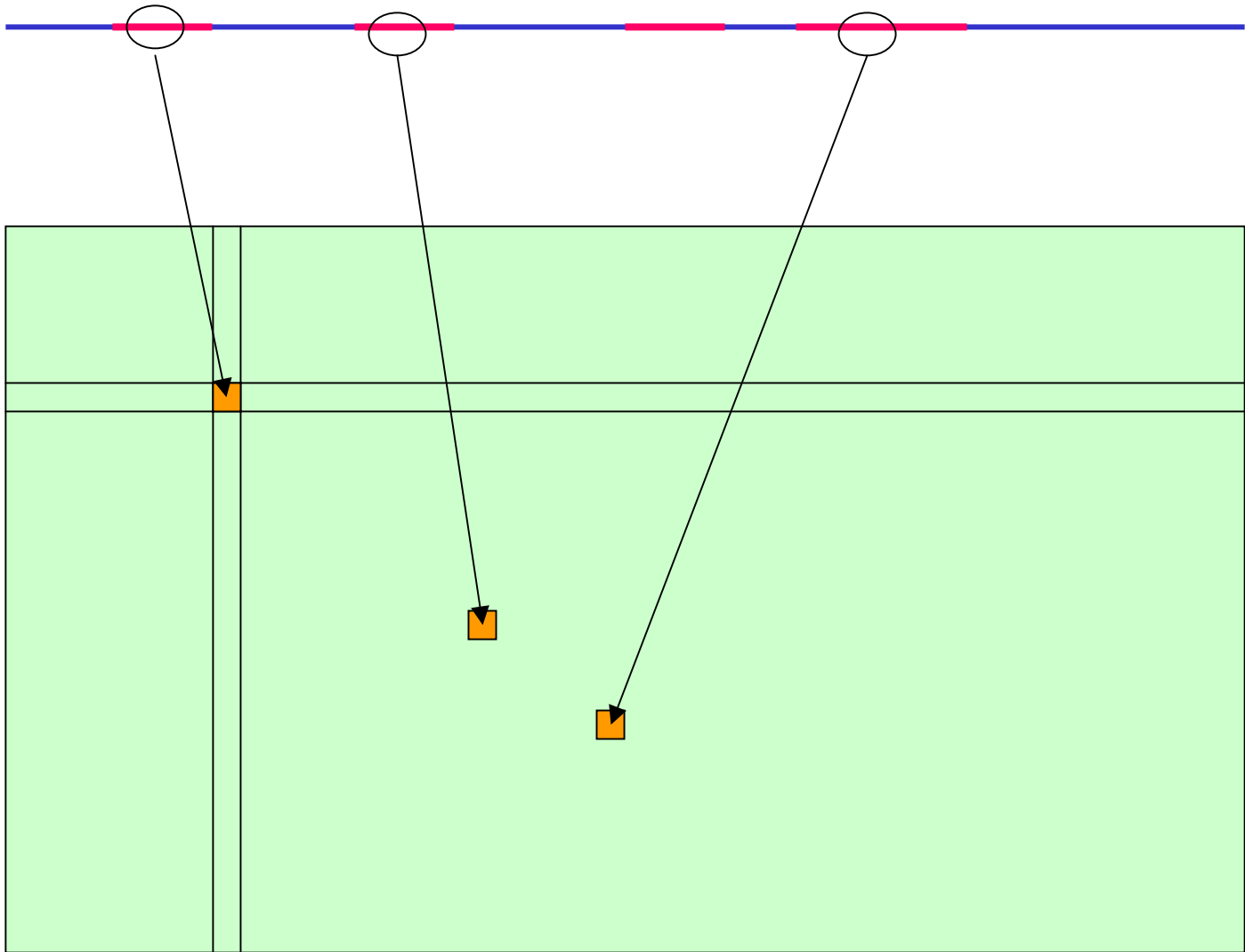
- Gene expression for a single gene can be measured by extracting mRNA from the cell and doing a simple **hybridization** experiment.
- Given a sample of cells, gene expression for every gene can be measured using a single **microarray** experiment.

# Microarray/DNA chip technology

- High-throughput method to study gene expression of thousands of genes simultaneously.
- Many applications:
  - Genetic disorders & Mutation/polymorphism detection
  - Study of disease subtypes
  - Drug discovery & toxicology studies
  - Pathogen analysis
  - Differing expressions over time, between tissues, between drugs, across disease states

# Microarray Data

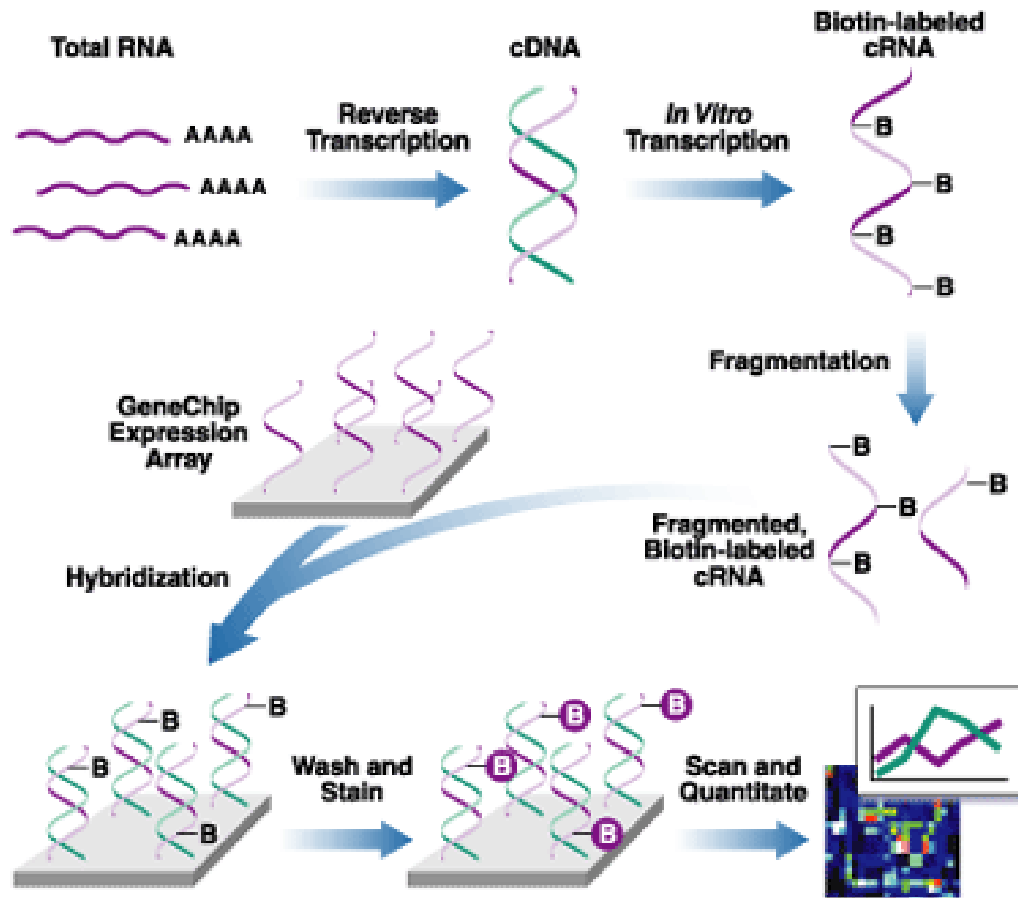
<b>Gene</b>	<b>Expression Level</b>
Gene1	
Gene2	
Gene3	
...	



# Microarray/DNA chips (Simplified)

- Construct **probes** corresponding to reverse complements of genes of interest.
- Microscopic quantities of probes placed on solid surfaces at defined spots on the chip.
- Extract mRNA from sample cells and **label** them.
- Apply labeled sample (mRNA extracted from cells) to every spot, and allow hybridization.
- Wash off unhybridized material.
- Use optical detector to measure amount of fluorescence from each spot.

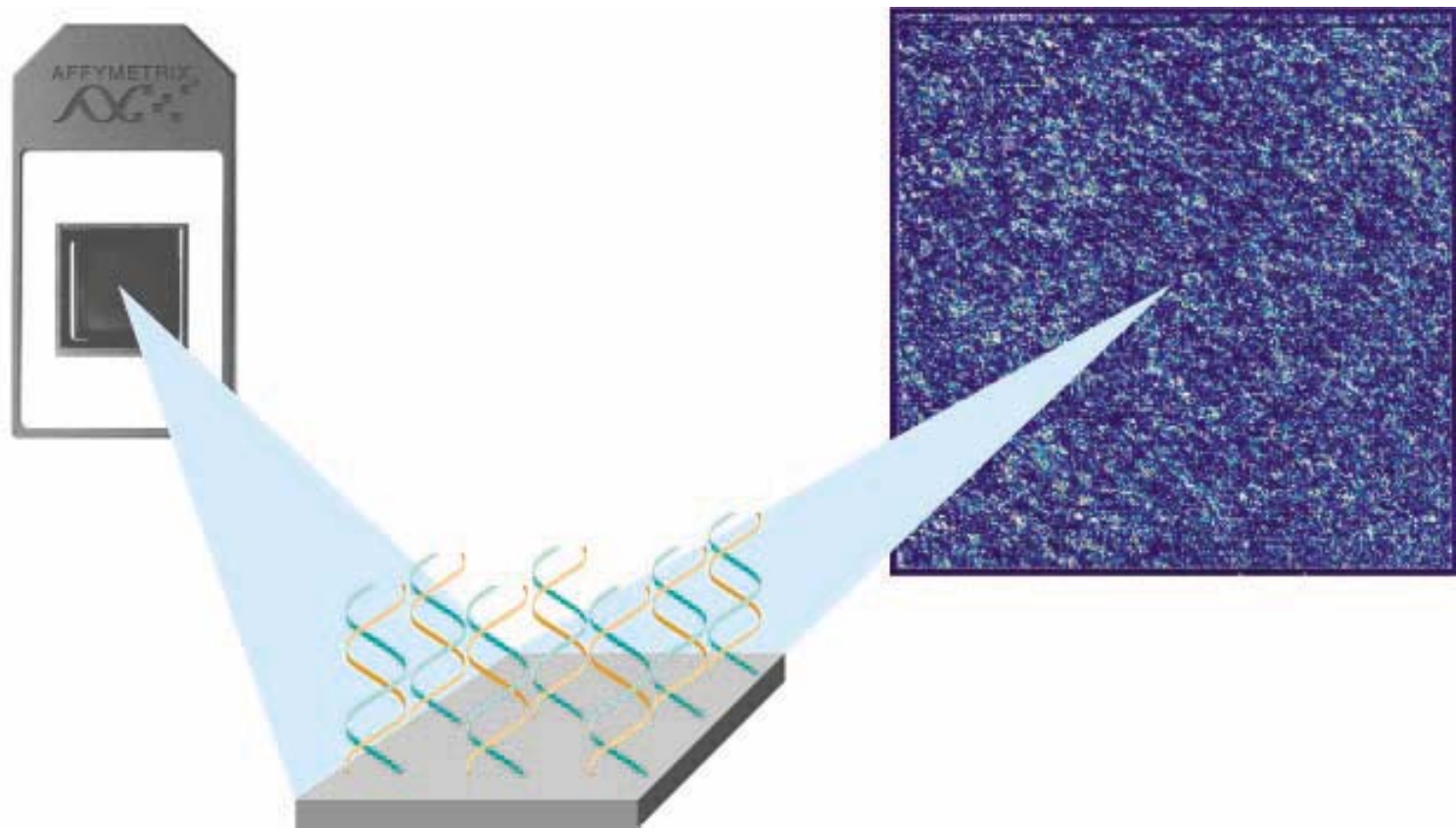
# Affymetrix DNA chip schematic

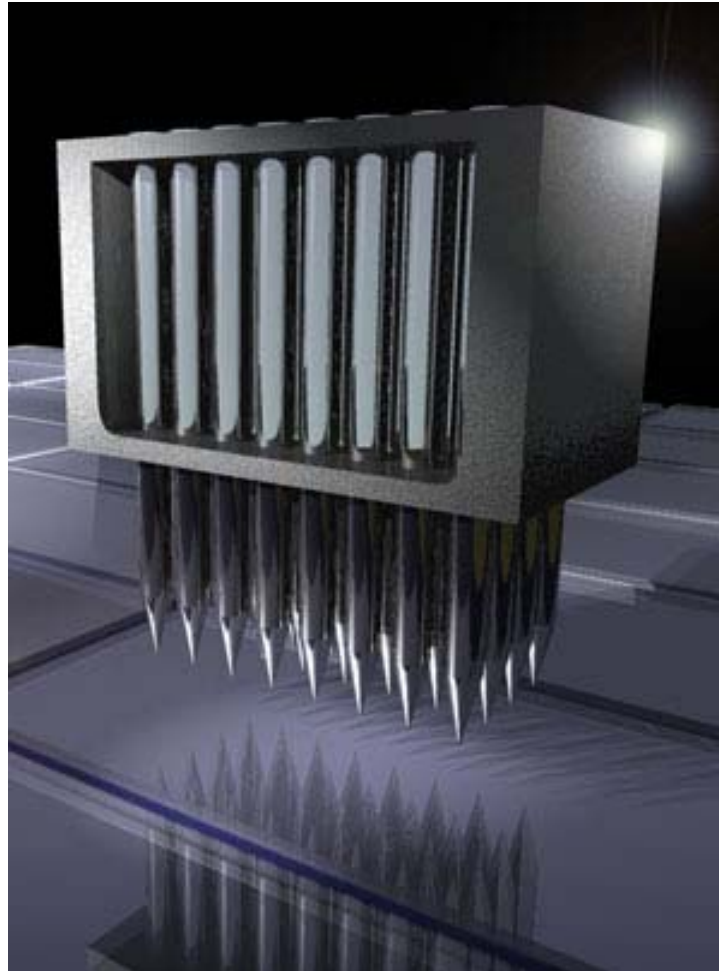


[www.affymetrix.com](http://www.affymetrix.com)



# DNA Chips & Images





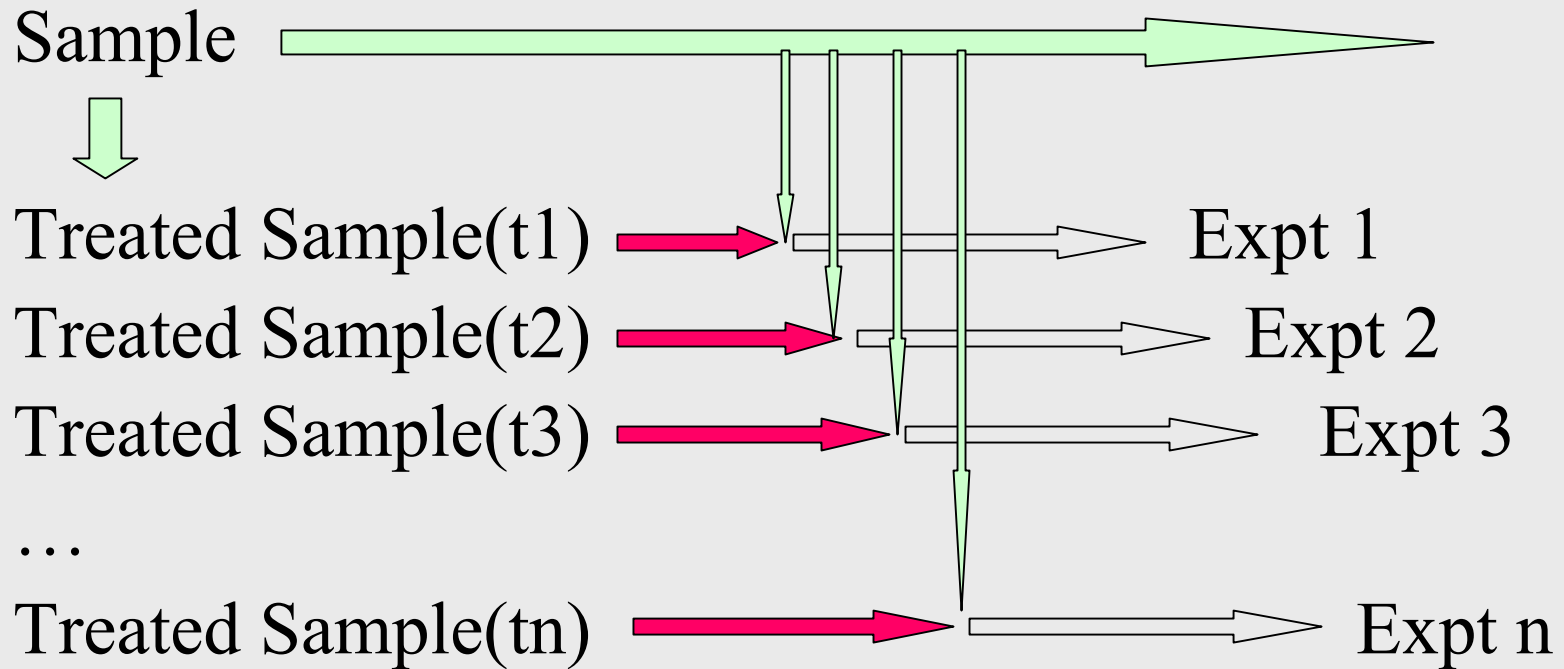
# Microarrays: competing technologies

- Affymetrix & Synteni/Stanford
- Differ in:
  - method to place DNA: Spotting vs. photolithography
  - Length of probe
  - Complete sequence vs. series of fragments

# How to compare 2 cell samples?

- mRNA from sample 1 is extracted and labeled with a **red fluorescent** dye.
- mRNA from sample 2 is extracted and labeled with a **green fluorescent** dye.
- Mix the samples and apply it to every spot on the microarray. Hybridize sample mixture to probes.
- Use optical detector to measure the amount of **green** and **red** fluorescence at each spot.

# Studying effect of a treatment over time



# Sources of Variations & Errors

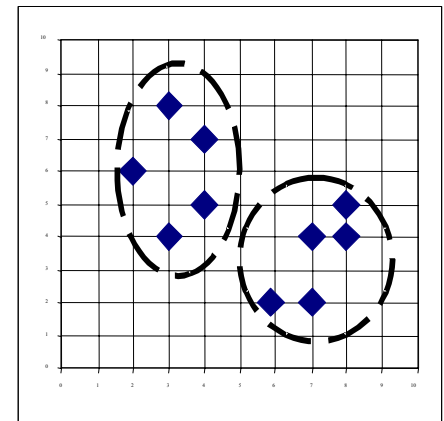
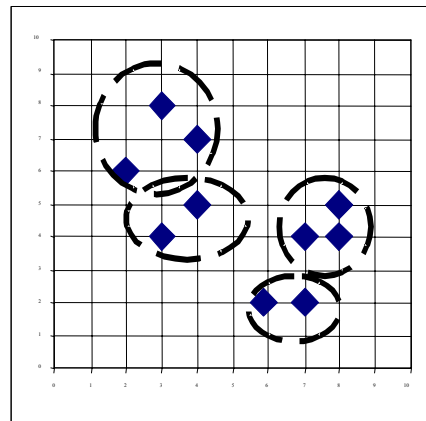
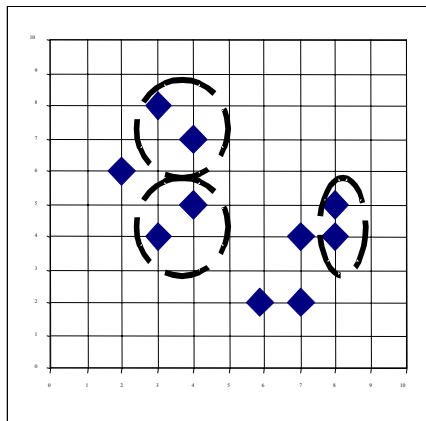
- Variations in cells/individuals.
- Variations in mRNA extraction, isolation, introduction of dye, variation in dye incorporation, dye interference.
- Variations in probe concentration, probe amounts, substrate surface characteristics
- Variations in hybridization conditions and kinetics
- Variations in optical measurements, spot misalignments, discretization effects, noise due to scanner lens and laser irregularities
- Cross-hybridization of sequences with high sequence identity.
- Limit of factor 2 in precision of results.

Need to Normalize data

# Clustering

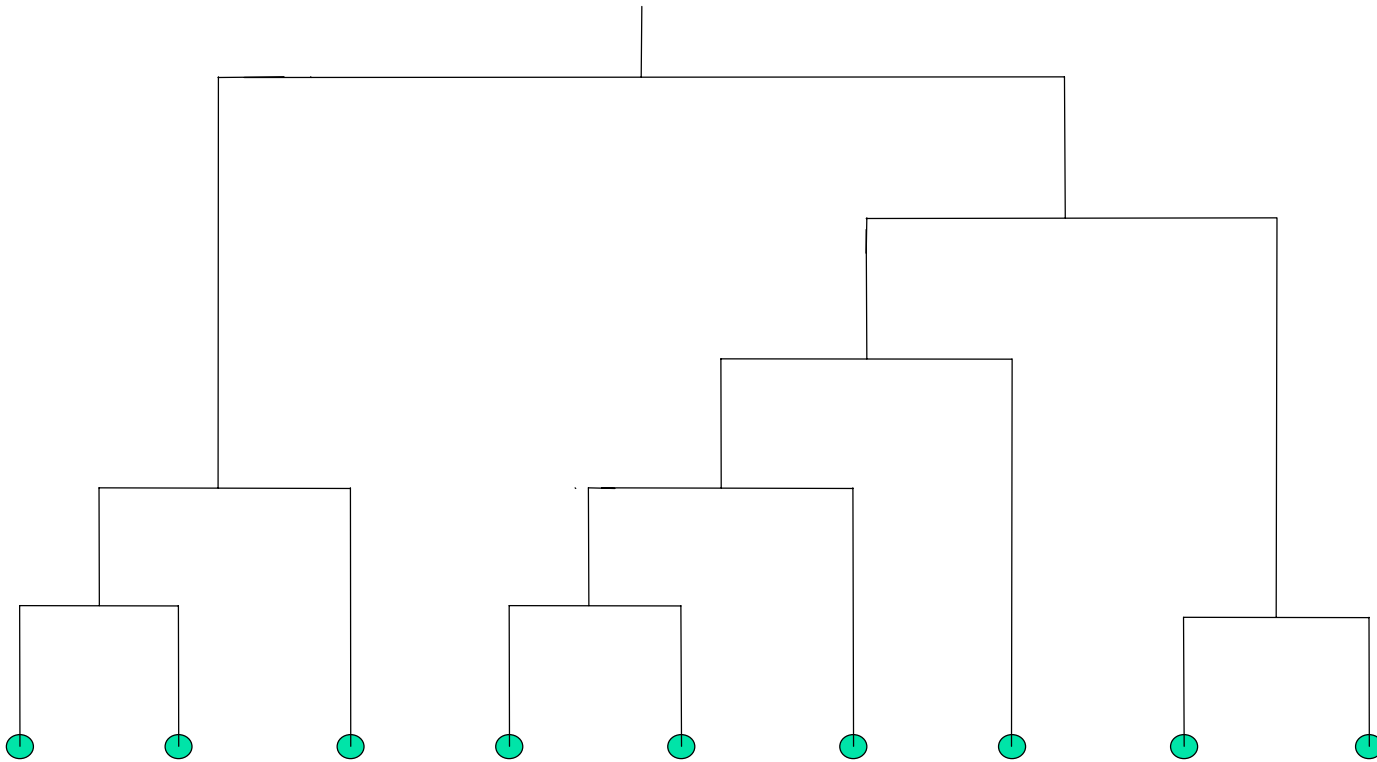
- Clustering is a general method to study patterns in gene expressions.
- Several known methods:
  - **Hierarchical Clustering** (Bottom-Up Approach)
  - **K-means Clustering** (Top-Down Approach)
  - **Self-Organizing Maps** (SOM)

# Hierarchical Clustering: Example





# A Dendrogram



# Hierarchical Clustering [Johnson, SC, 1967]

- Given  $n$  points in  $\mathbb{R}^d$ , compute the distance between every pair of points
- While (not done)
  - Pick closest pair of points  $s_i$  and  $s_j$  and make them part of the same cluster.
  - Replace the pair by an average of the two  $s_{ij}$

Try the applet at:

<http://www.cs.mcgill.ca/~papou/#applet>

# Distance Metrics

- For clustering, define a distance function:
  - Euclidean distance metrics

$$D_k(X, Y) = \left[ \sum_{i=1}^d (X_i - Y_i)^k \right]^{1/k}$$

k=2: Euclidean Distance

- Pearson correlation coefficient

$$\rho_{xy} = \frac{1}{d} \sum_{i=1}^d \left( \frac{X_i - \bar{X}}{\sigma_x} \right) \left( \frac{Y_i - \bar{Y}}{\sigma_y} \right)$$

$-1 \leq \rho_{xy} \leq 1$

**EXHIBIT 3.4** Joint Probability Model for the Ratings of Two People

(a)  $\rho_{XY} = 0$

x	y			Total
	1	2	3	
3	1/9	1/9	1/9	1/3
2	1/9	1/9	1/9	1/3
1	1/9	1/9	1/9	1/3
Total	1/3	1/3	1/3	1

(b)  $\rho_{XY} = \frac{1}{2}$

x	y			Total
	1	2	3	
3	1/18	1/18	4/18	1/3
2	1/18	4/18	1/18	1/3
1	4/18	1/18	1/18	1/3
Total	1/3	1/3	1/3	1

(c)  $\rho_{XY} = -\frac{1}{2}$

x	y			Total
	1	2	3	
3	4/18	1/18	1/18	1/3
2	1/18	4/18	1/18	1/3
1	1/18	1/18	4/18	1/3
Total	1/3	1/3	1/3	1

(d)  $\rho_{XY} = \frac{1}{3}$

x	y			Total
	1	2	3	
3	1/27	2/27	6/27	1/3
2	2/27	5/27	2/27	1/3
1	6/27	2/27	1/27	1/3
Total	1/3	1/3	1/3	1

(e)  $\rho_{XY} = -\frac{2}{3}$

x	y			Total
	1	2	3	
3	6/27	2/27	1/27	1/3
2	2/27	5/27	2/27	1/3
1	1/27	2/27	6/27	1/3
Total	1/3	1/3	1/3	1

(f)  $\rho_{XY} = \frac{2}{3}$

x	y			Total
	1	2	3	
3	1/36	2/36	9/36	1/3
2	2/36	8/36	2/36	1/3
1	9/36	2/36	1/36	1/3
Total	1/3	1/3	1/3	1

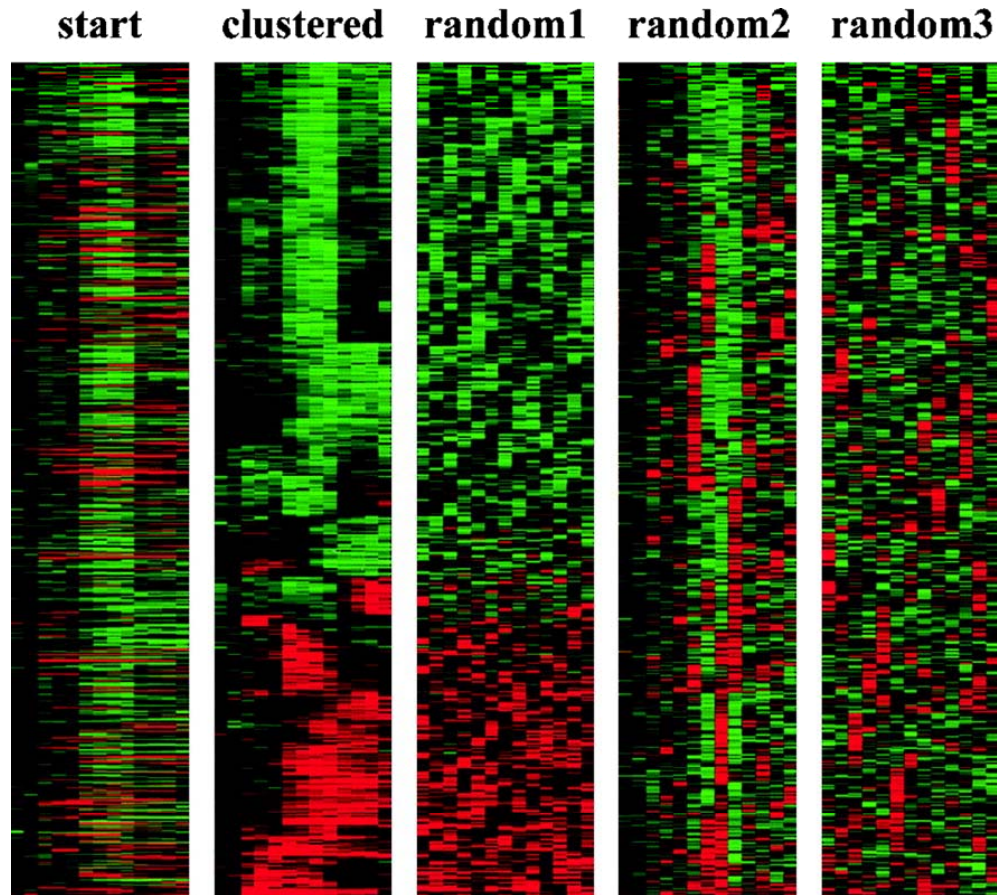
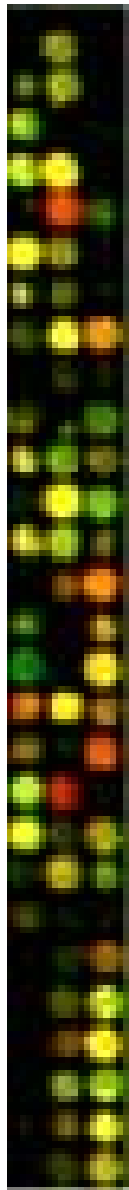
(g)  $\rho_{XY} = -\frac{1}{3}$

x	y			Total
	1	2	3	
3	9/36	2/36	1/36	1/3
2	2/36	8/18	2/18	1/3
1	1/36	2/36	9/36	1/3
Total	1/3	1/3	1/3	1

# Clustering of gene expressions

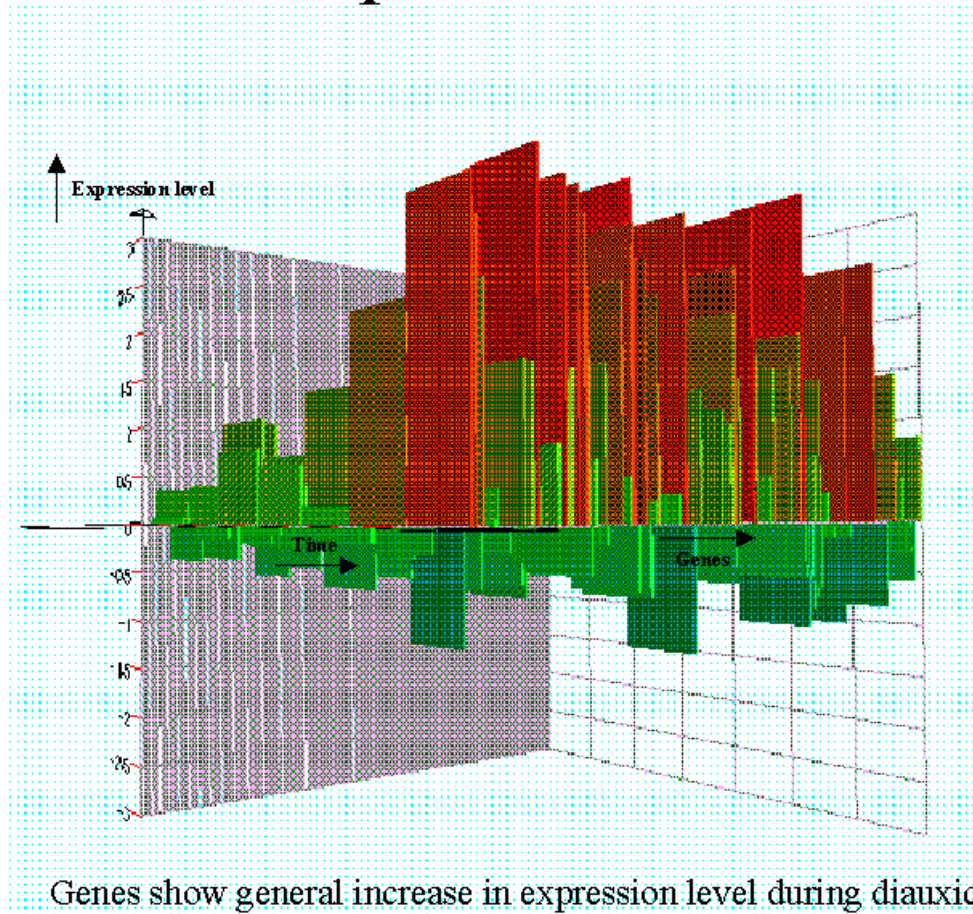
- Represent each gene as a vector or a point in **d**-space where **d** is the number of arrays or experiments being analyzed.

# Clustering Random vs. Biological Data



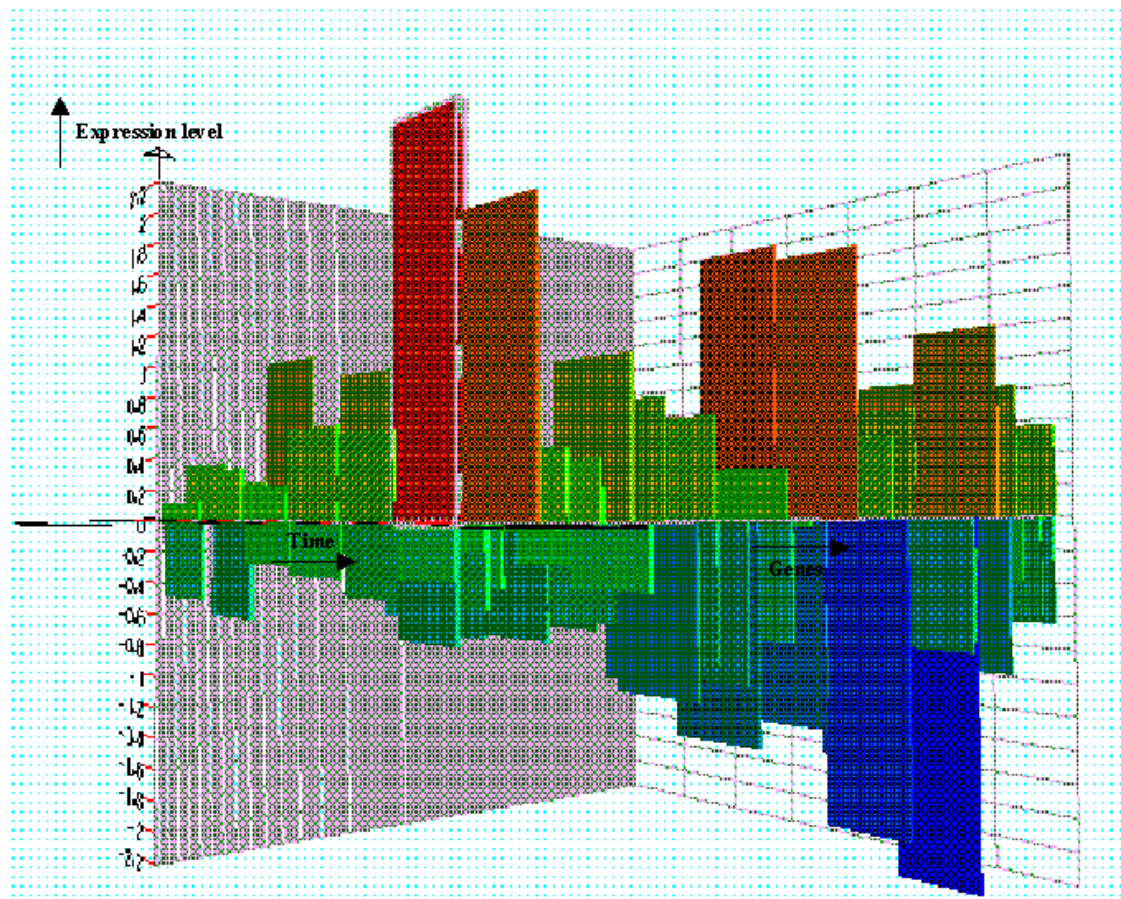
From Eisen MB, et al, *PNAS* 1998 95(25):14863-8

# Expression Profiles for Respiration Genes



Genes show general increase in expression level during diauxic shift

# Expression Profiles for Fermentation Genes



Bar two exceptions, genes show general decrease in expression level during diauxic shift



# Observations

- ◆ As glucose was depleted - Marked change in the global pattern of gene expression
- ◆ ~50% of differentially expressed genes have unknown function
- ◆ Genes with similar expression profiles had common promoters
- ◆ Expression patterns observed match those observed in other types of experiments

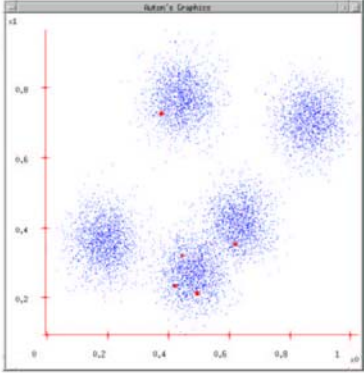
# K-Means Clustering: Example

Example from Andrew Moore's tutorial on Clustering.

Start

### K-means

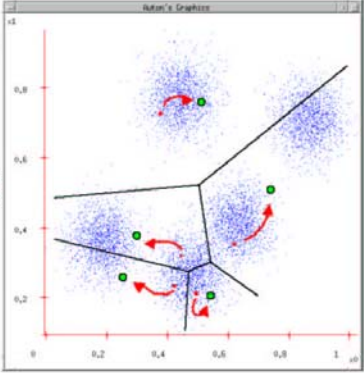
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 7

### K-means

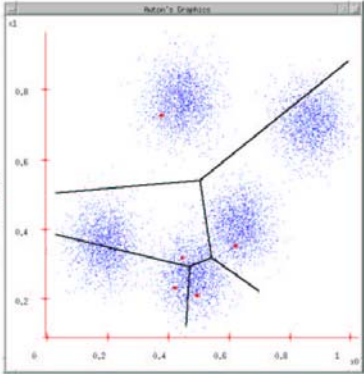
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 8

### K-means

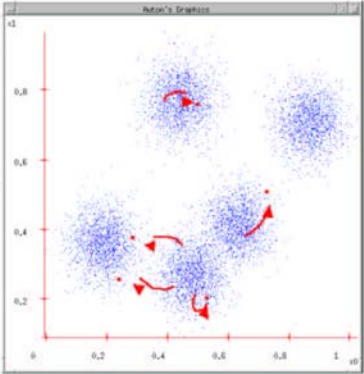
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 9

### K-means

1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



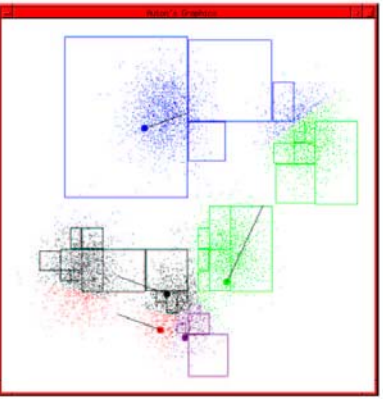
Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 10

## K-means Start

Advance apologies: in Black and White this example will deteriorate

Example generated by Dan Pelleg's super-duper fast K-means system:

*Dan Pelleg and Andrew Moore. Accelerating Exact k-means Algorithms with Geometric Reasoning. Proc. Conference on Knowledge Discovery in Databases 1999, (KDD99) (available on [www.autonlab.org/pap.html](http://www.autonlab.org/pap.html))*



Copyright © 2001, Andrew W. Moore

K-means and Hierarchical Clustering: Slide 11

## K-means continues

...



Copyright © 2001, Andrew W. Moore

K-means and Hierarchical Clustering: Slide 13

## K-means continues

...



Copyright © 2001, Andrew W. Moore

K-means and Hierarchical Clustering: Slide 12

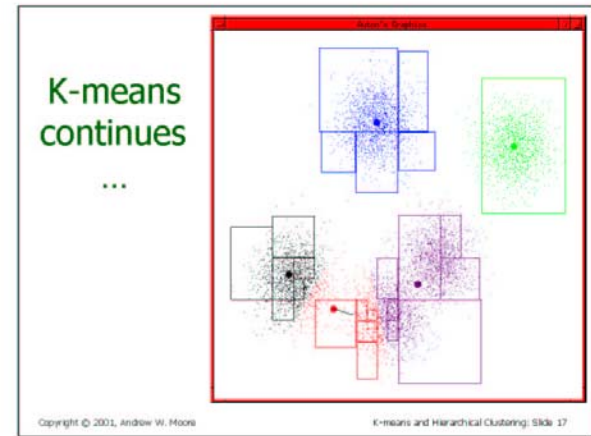
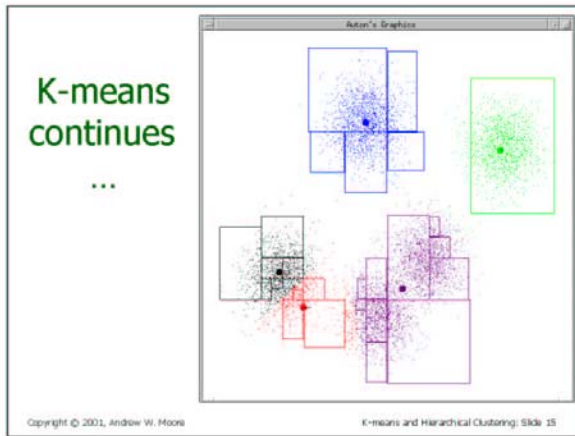
## K-means continues

...



Copyright © 2001, Andrew W. Moore

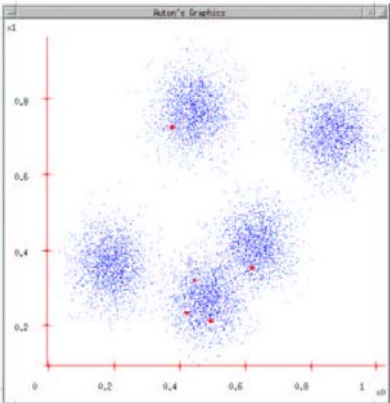
K-means and Hierarchical Clustering: Slide 14



Start

### K-means

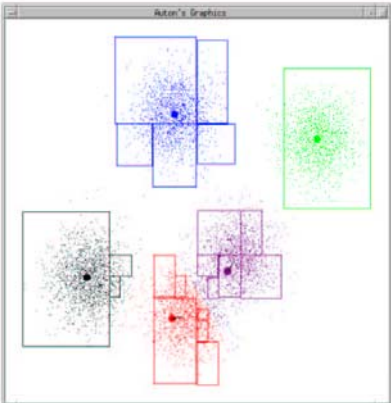
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 7

### K-means continues

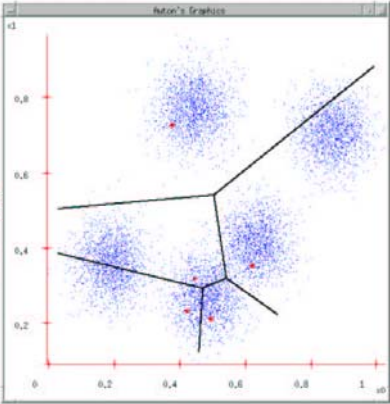
...



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 19

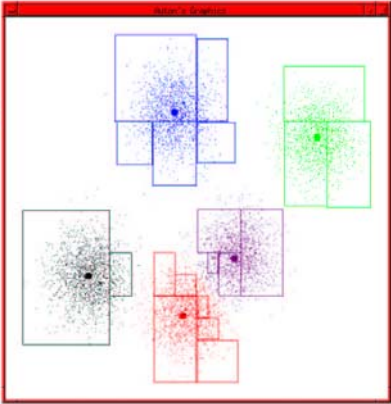
### K-means

1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 8

### K-means terminates



Copyright © 2001, Andrew W. Moore  
K-means and Hierarchical Clustering: Slide 20

End

# K-Means Clustering [McQueen '67]

## Repeat

- Start with randomly chosen cluster centers
- Assign points to give greatest increase in score
- Recompute cluster centers
- Reassign points

until (no changes)

Try the applet at: <http://www.cs.mcgill.ca/~bonnef/project.html>

# Comparisons

- Hierarchical clustering
  - Number of clusters not preset.
  - Complete hierarchy of clusters
  - Not very robust, not very efficient.
- K-Means
  - Need definition of a **mean**. Categorical data?
  - More efficient and often finds optimum clustering.