# CAP 5510: Introduction to Bioinformatics
# CGS 5166: Bioinformatics Tools

## Giri Narasimhan

ECS 254; Phone: x3748

giri@cis.fiu.edu

www.cis.fiu.edu/~giri/teach/BioinfS11.html

# In Memoriam



Isabel Melo, Accountant
Feb 16, 2011

# 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 the right 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.
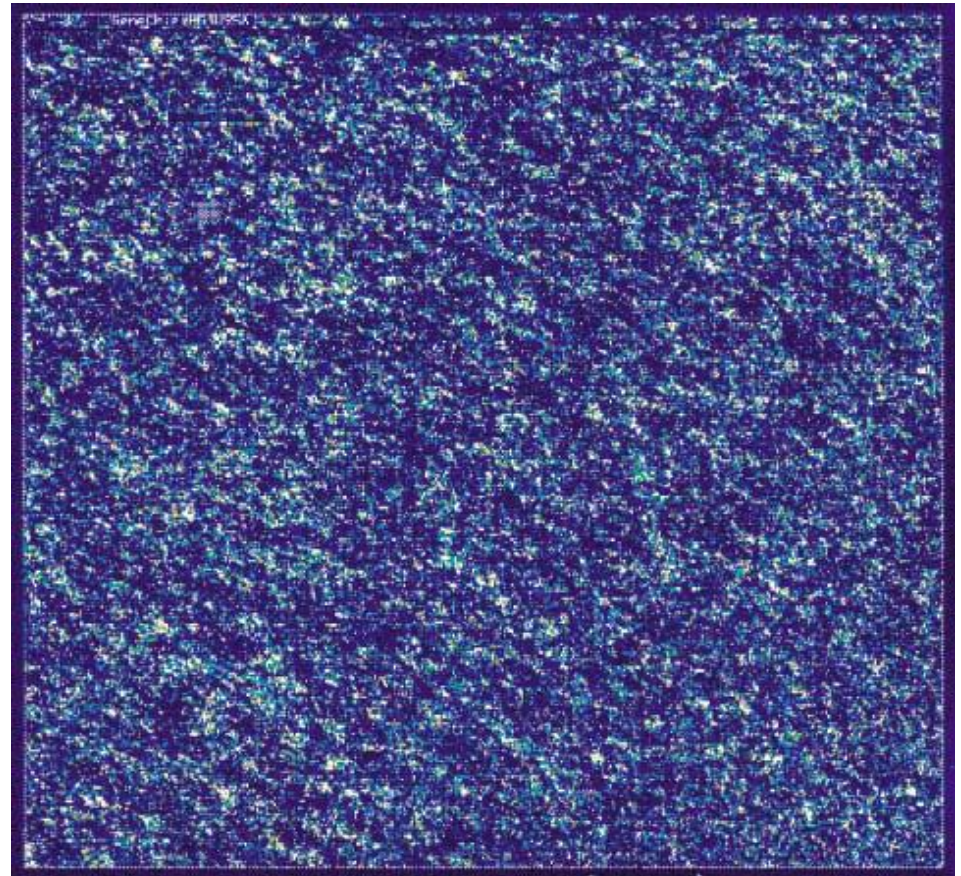
# Microarray/DNA chip technology

❑ 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
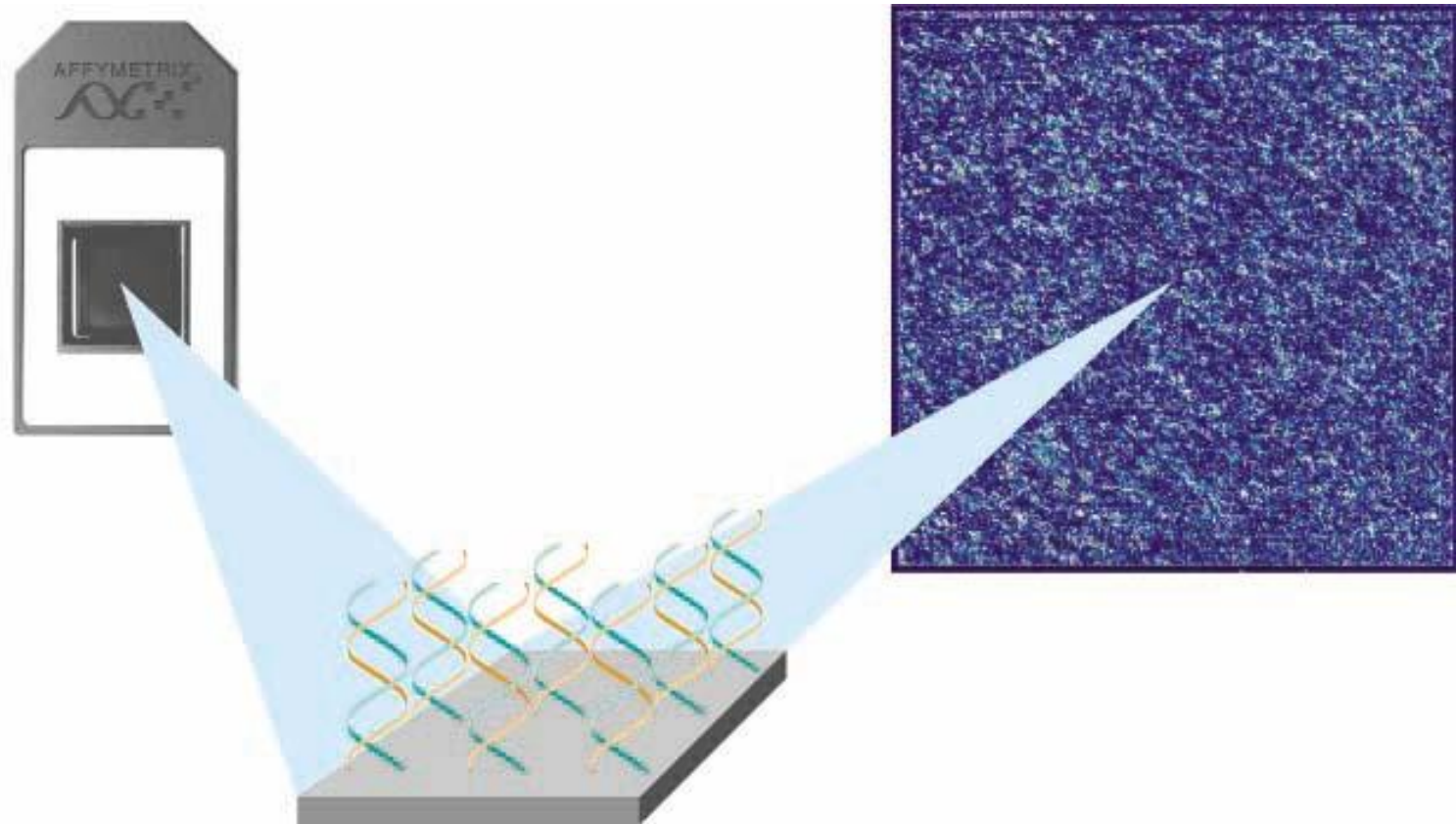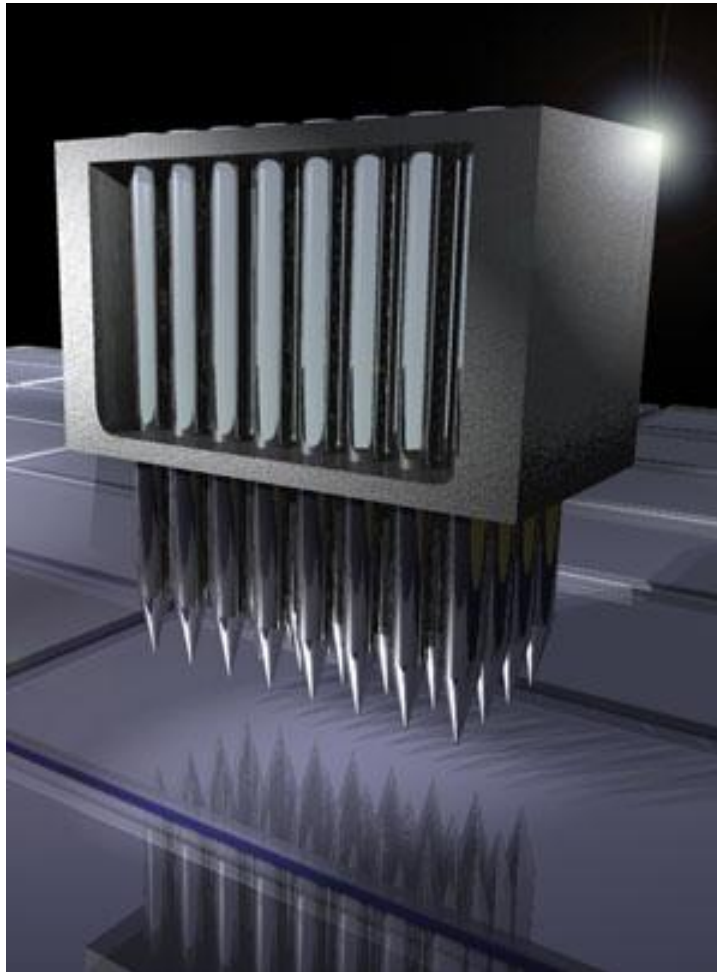
# Microarray Data

| Gene | Expression Level |
|:---:|:---:|
| Gene1 | |
| Gene2 | |
| Gene3 | |
| … | |

# Gene Chips

# DNA Chips & Images

# Gene g



Probe 1    Probe 2    ...    Probe N
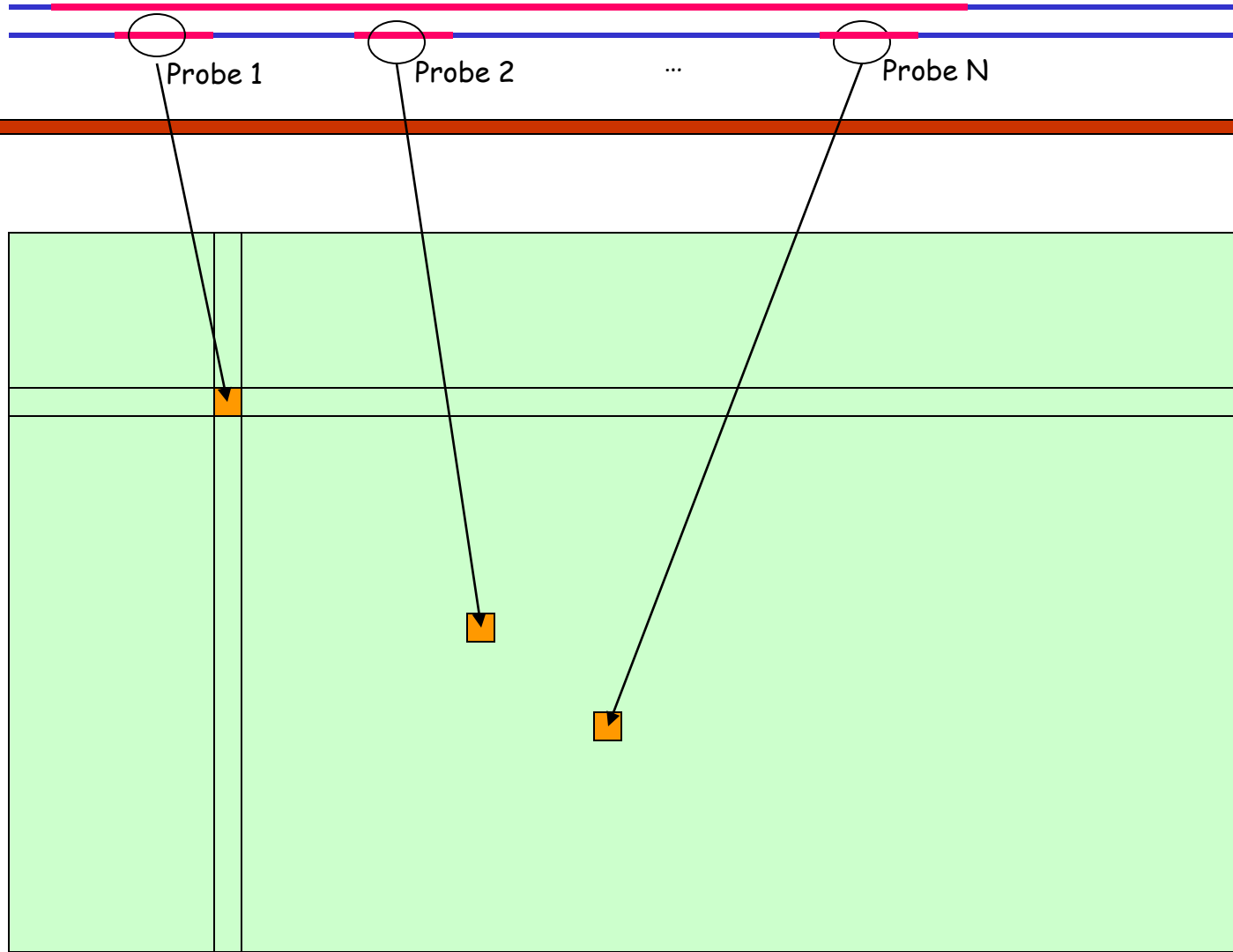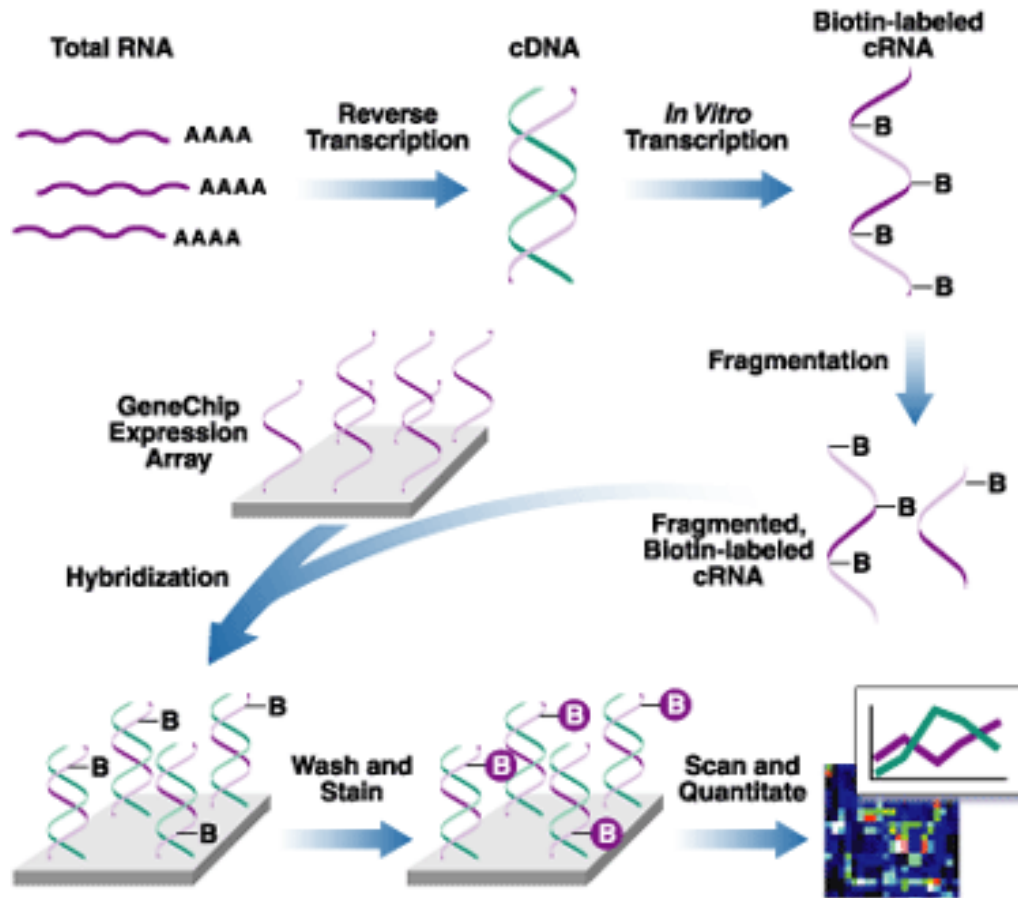
CAP5510    11

# 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.
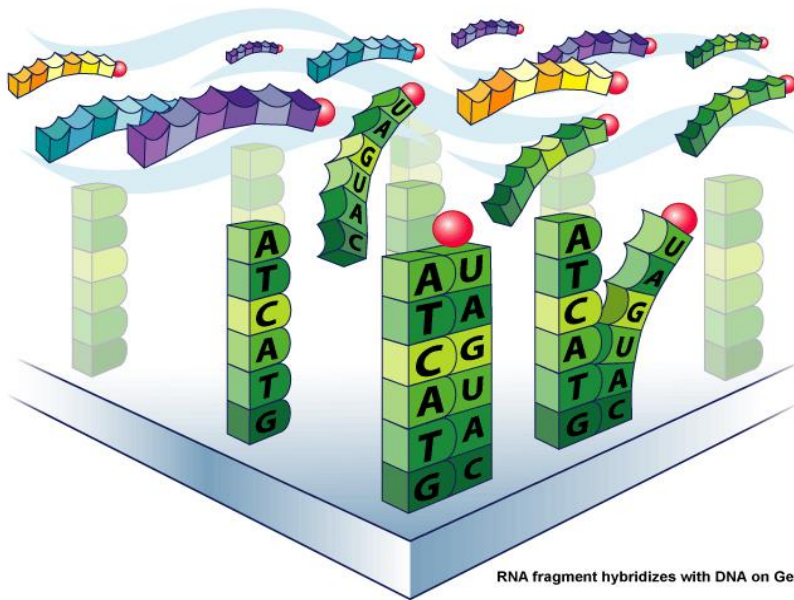- ❑ Use optical detector to measure amount of fluorescence from each spot.

# Affymetrix DNA chip schematic

www.affymetrix.com

# What's on the slide?



RNA fragments with fluorescent tags from sample to be tested

RNA fragment hybridizes with DNA on GeneChip® array

Shining a laser light at GeneChip® array causes tagged DNA fragments that hybridized to glow

Non-hybridized DNA

Hybridized DNA
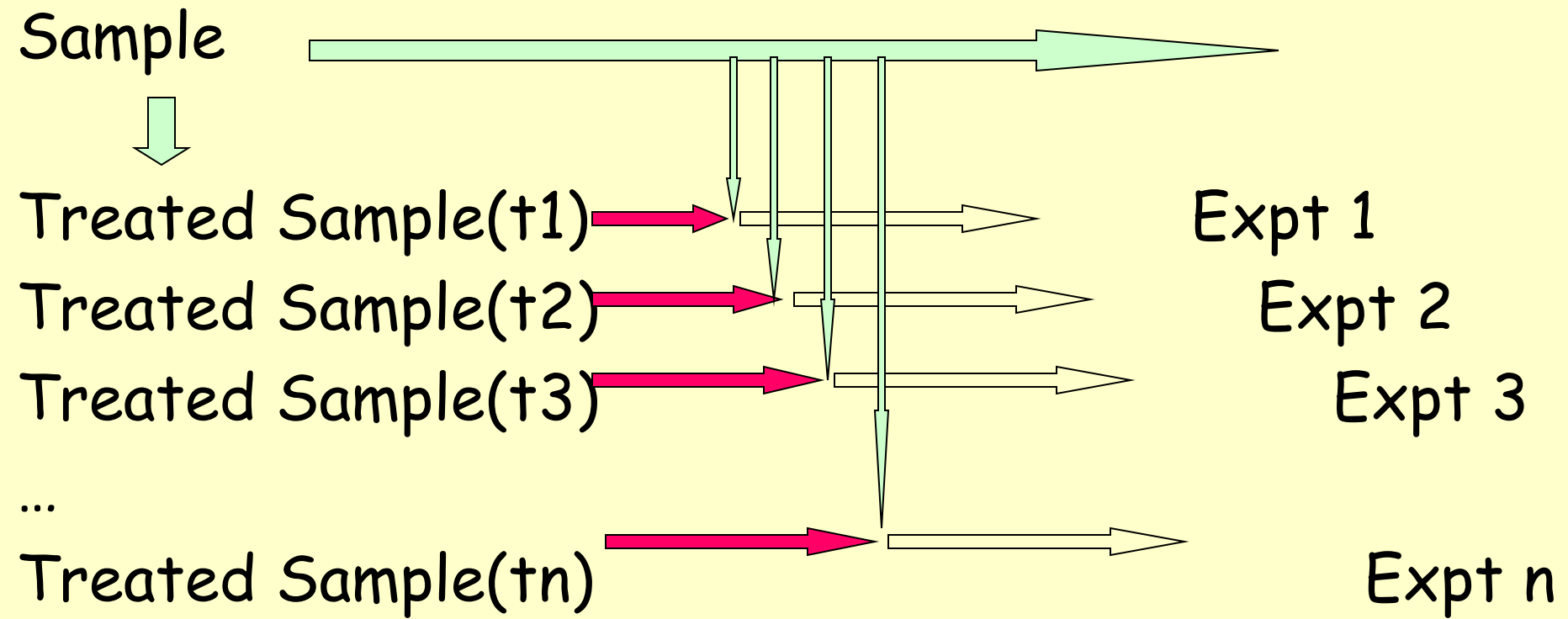
# Microarrays: competing technologies

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

# Study effect of treatment over time

Sample

Treated Sample(t1)         Expt 1

Treated Sample(t2)         Expt 2

Treated Sample(t3)         Expt 3

…

Treated Sample(tn)         Expt n

# 2-color DNA microarray

AFGC

Treated          Control

Normalization

Data extraction

Scanning

mRNA          mRNA

Simultaneous
hybridization

Cy5 Probe Cy3 Probe

# How to compare 2 cell samples with Two-Color Microarrays?

- ❑ 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.

# Sources of Variations & Experimental 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
- ❑ Variation changes with intensity: larger variation at low or high expression levels

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

CAP5510

# Hierarchical Clustering [Johnson, SC, 1967]

❑ Given $n$ points in $R^d$, 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://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletH.html

# 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

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

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

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

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

| x | y 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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) $\rho_{XY} = \frac{4}{9}$

| x | y 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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) $\rho_{XY} = -\frac{5}{9}$

| x | y 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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 1 | y 2 | y 3 | Total |
|---|---|---|---|---|
| 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



start     clustered     random1     random2     random3

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

# K-Means Clustering: Example

Example from Andrew Moore's tutorial on Clustering.