# Gene Profiling, Clustering, and Networking

**Alfred O. Hero III**

*University of Michigan, Ann Arbor, MI*

*http://www.eecs.umich.edu/~hero*

**Mar. 2005**

1. Genomics, transcriptomics and gene microarrays
2. Preprocessing of gene microarray data
3. Screening differentially expressed genes
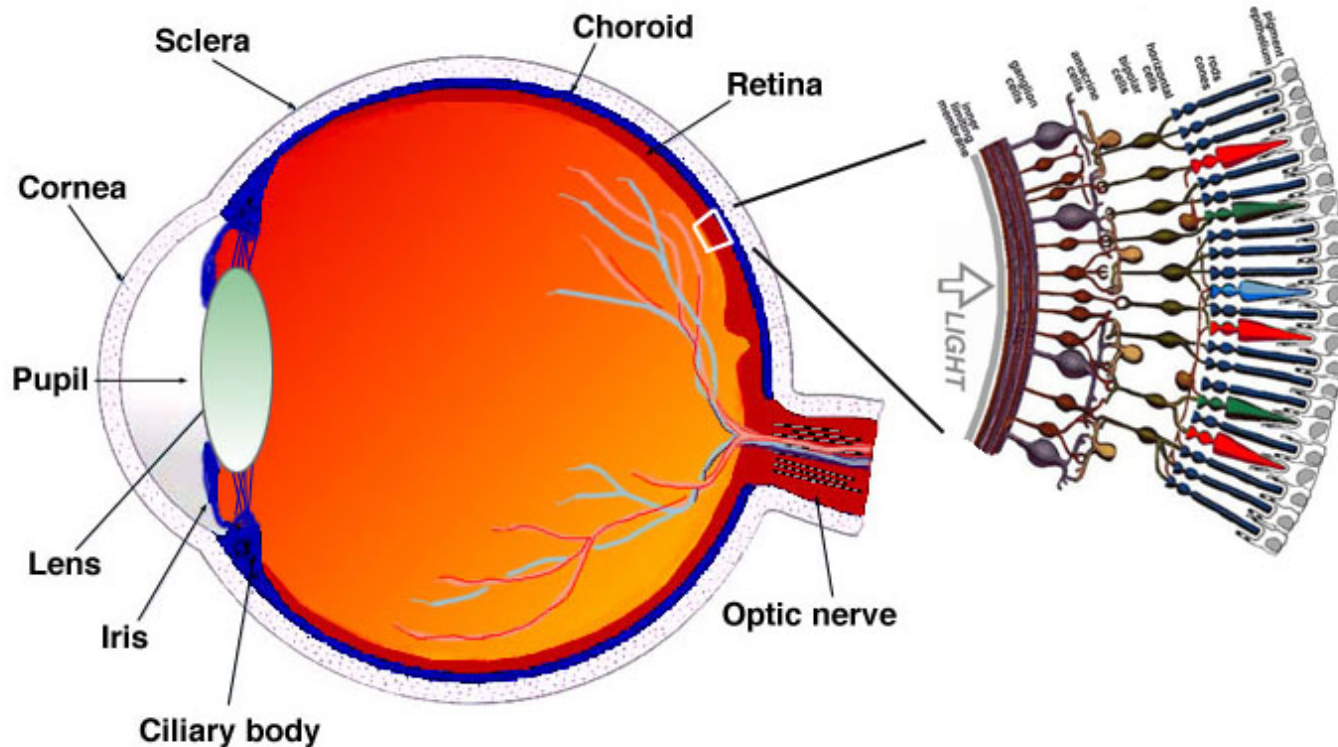4. Clustering gene co-regulation patterns
5. Conclusions

# Acknowledgements

- Anand Swaroop, Ophthalmology, UM
- David States, Bioinformatics, UM
- Alan Mears, Univ Ottawa (CA)
- Gilles Fleury, ESE, France
- Debashis Ghosh, Biostatistics, UM
- Terry Speed, Statistics, UCB
- Jindan Yu, BME, UM
- Dongxiao Zhu, Bioinformatics, UM
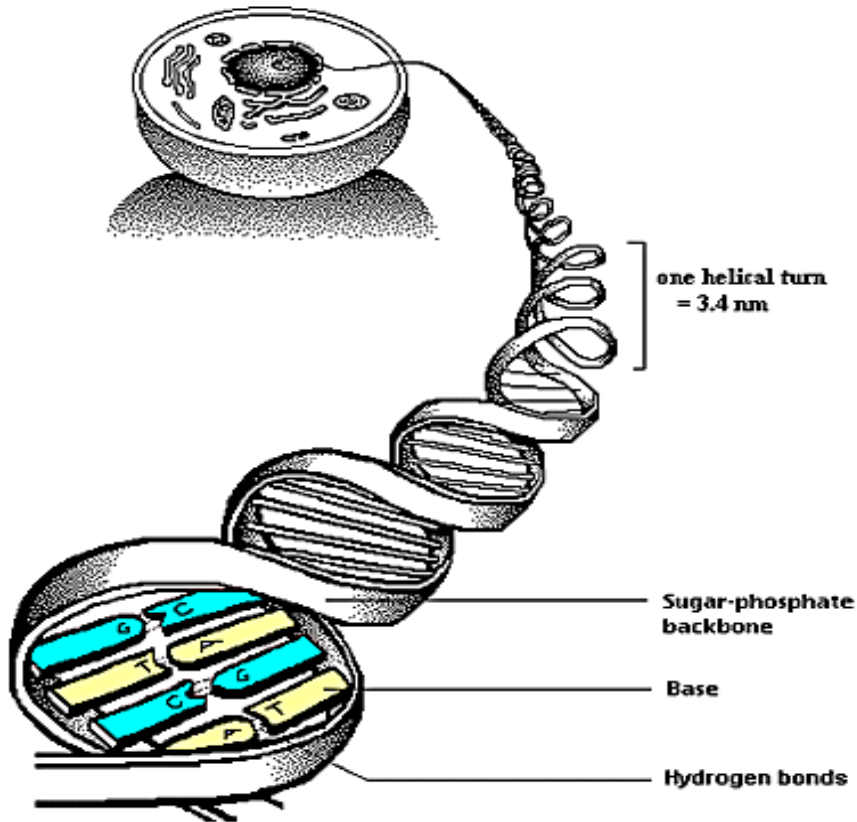- Yuezhou Jing, Dept. Statistics, UM

# Some Biological Questions

- What is genetic basis for photoreceptor development, aging, and degeneration?

- What are patterns of gene expression in the retina over time?
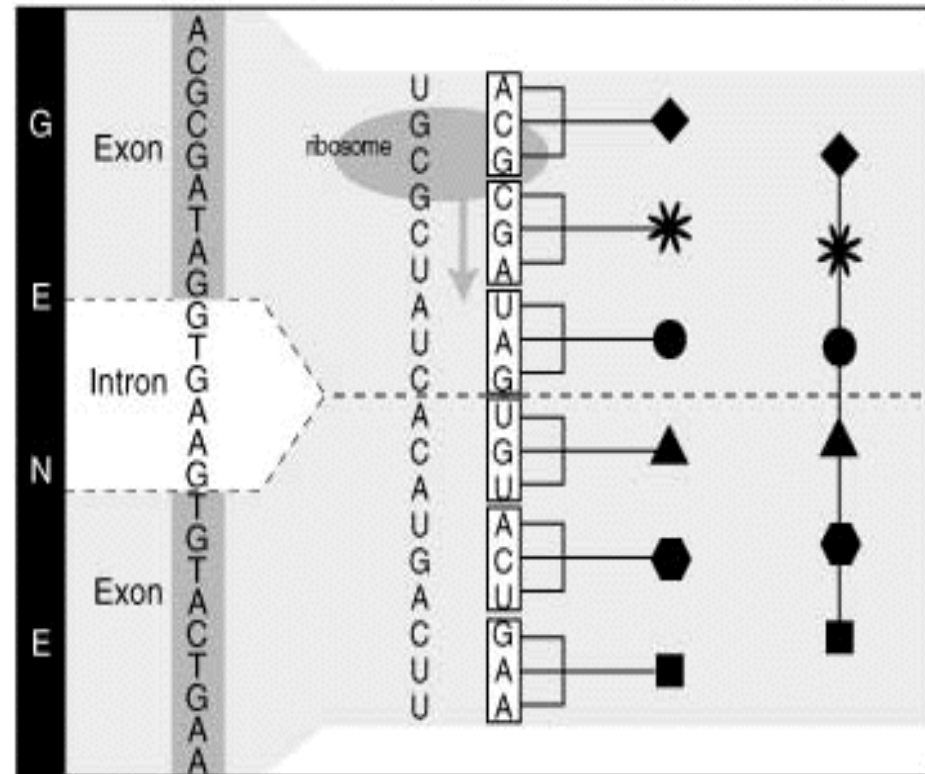
- What genes mediate development of rods and cones?

# 1. Genomics, Transcriptomics and Gene Microarrays



THE STRUCTURE OF DNA

one helical turn = 3.4 nm

Sugar-phosphate backbone

Base

Hydrogen bonds



Transcription → Translation
DNA → mRNA → tRNA → Amino Acid → Polypeptide chain

http://www-stat.stanford.edu/~susan/courses/s166/node2.html
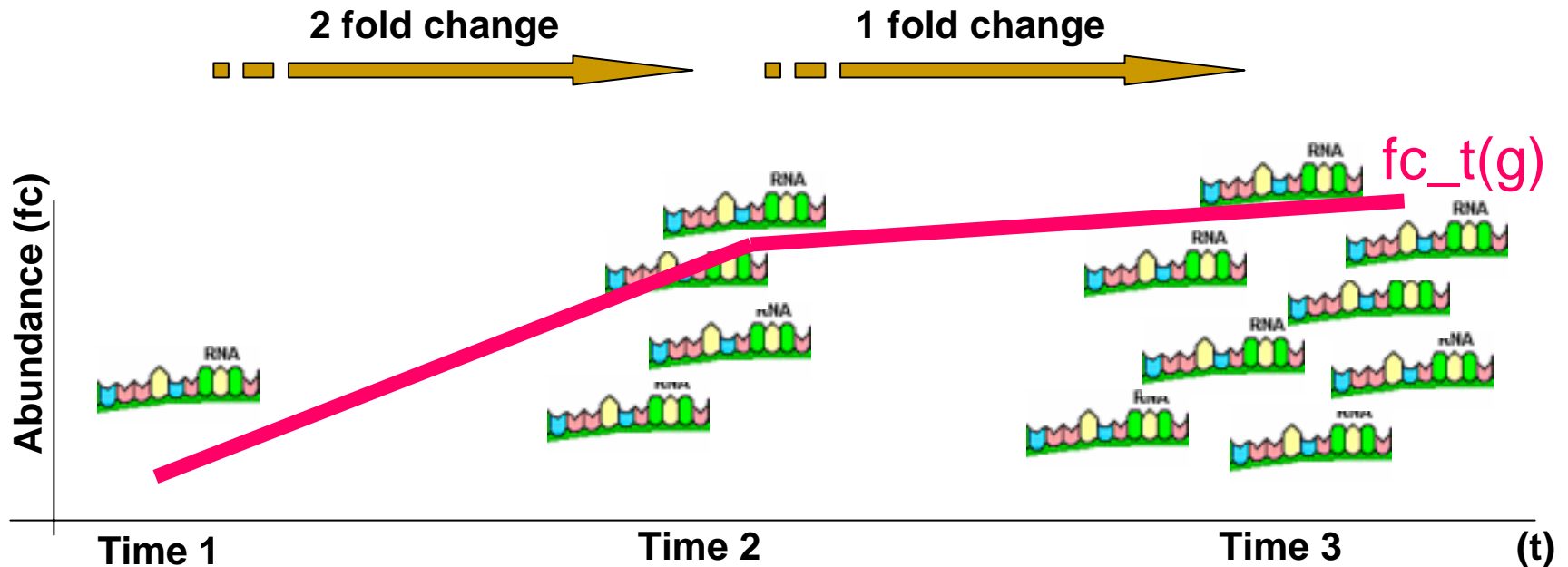
http://www.genome.gov/

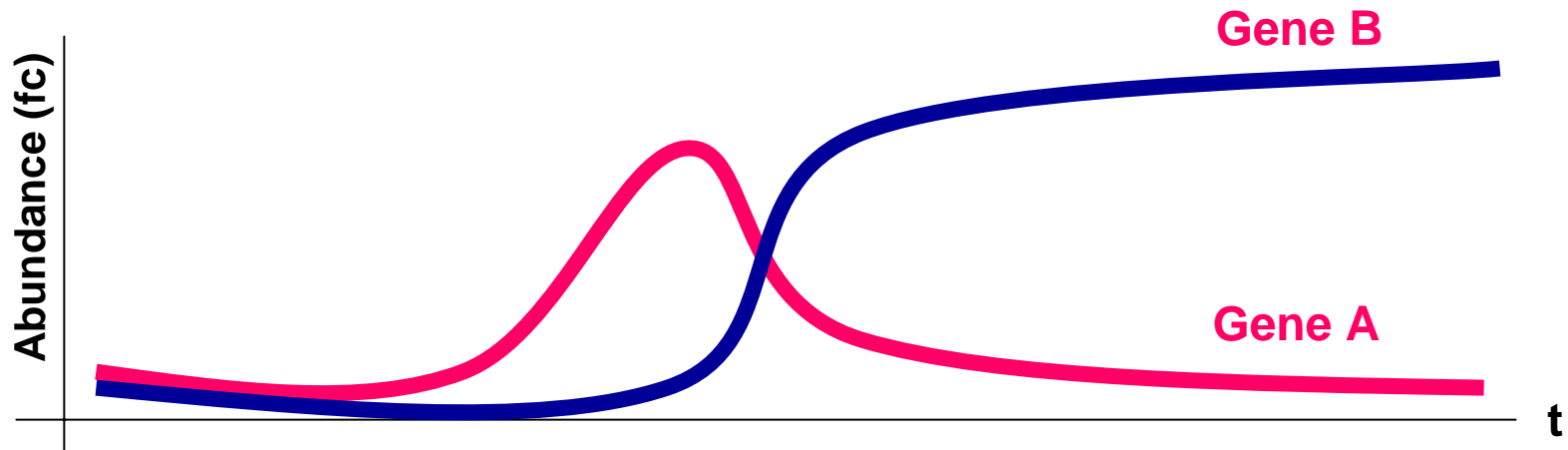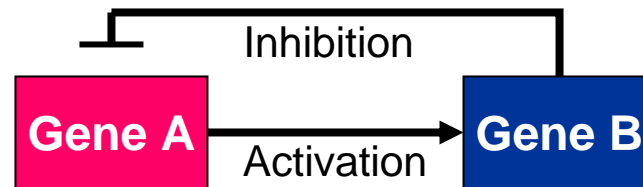# Transcriptomics: Gene expression profiling

*What is pattern of gene activation/inactivation over time, tissue, therapy, etc?*

**2 fold change**

**1 fold change**

Abundance (fc)

fc_t(g)

Time 1                    Time 2                    Time 3        (t)
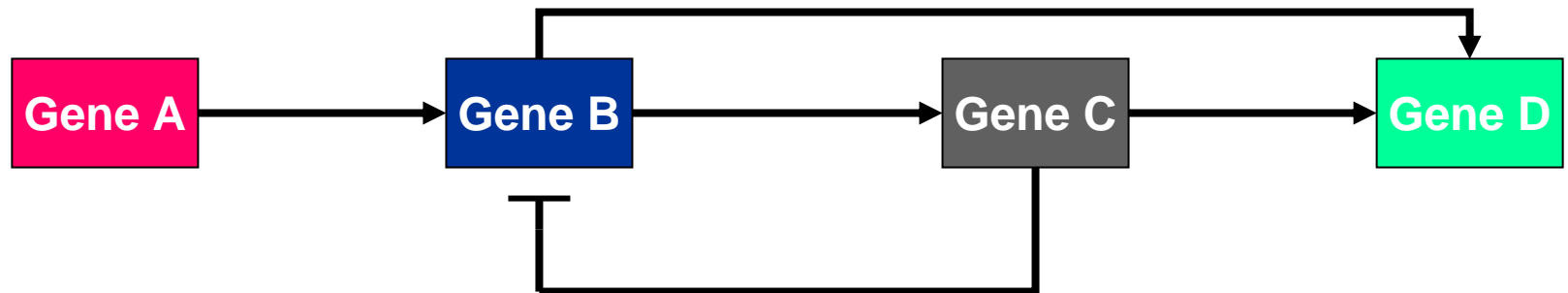
# Discovery of Genetic Circuits

*How do genes regulate (activate/inhibit)*

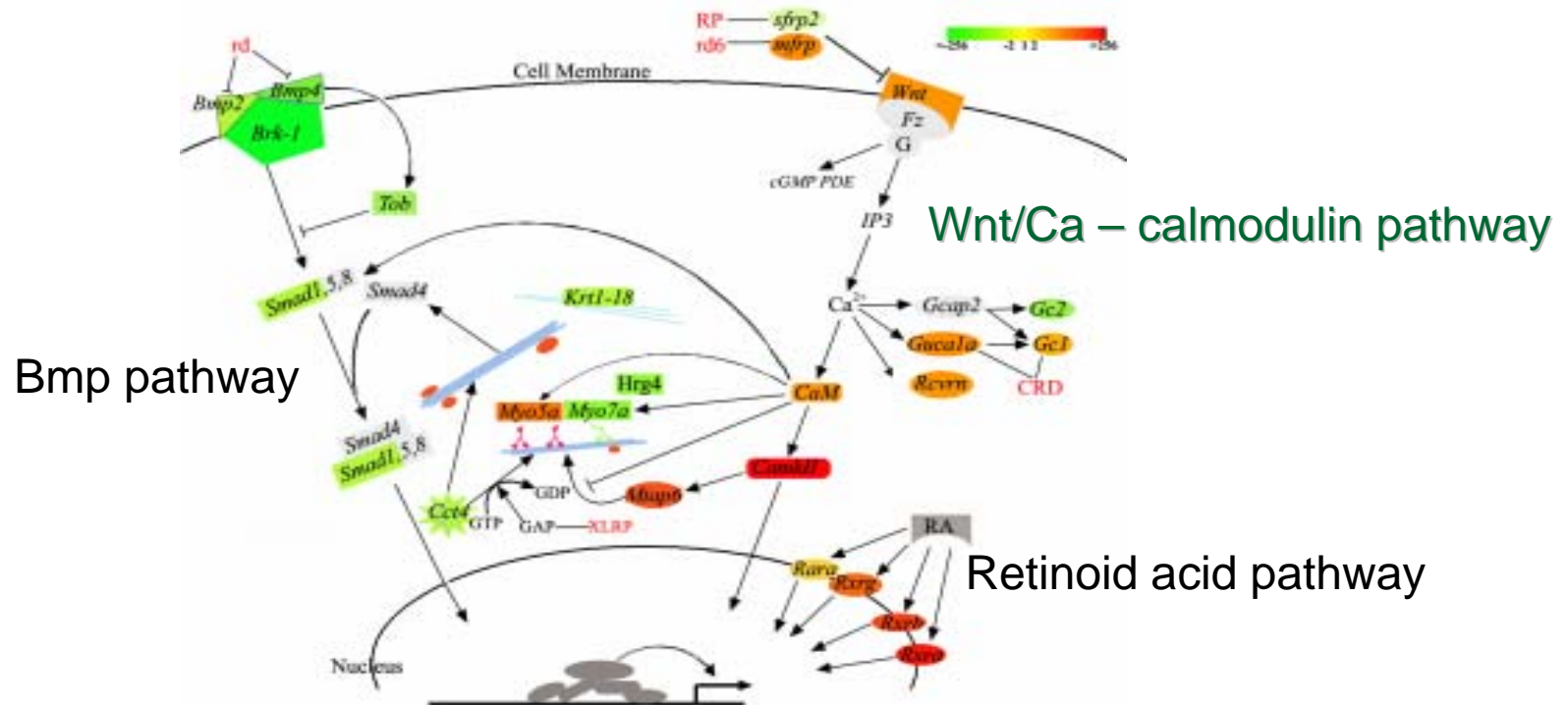*each other's expression levels over time?*

# Discovery of Genetic Pathways

*What sequence of gene interactions lead to a specific metabolic/structural (dys)function*

# Discovery of Gene Regulation Networks

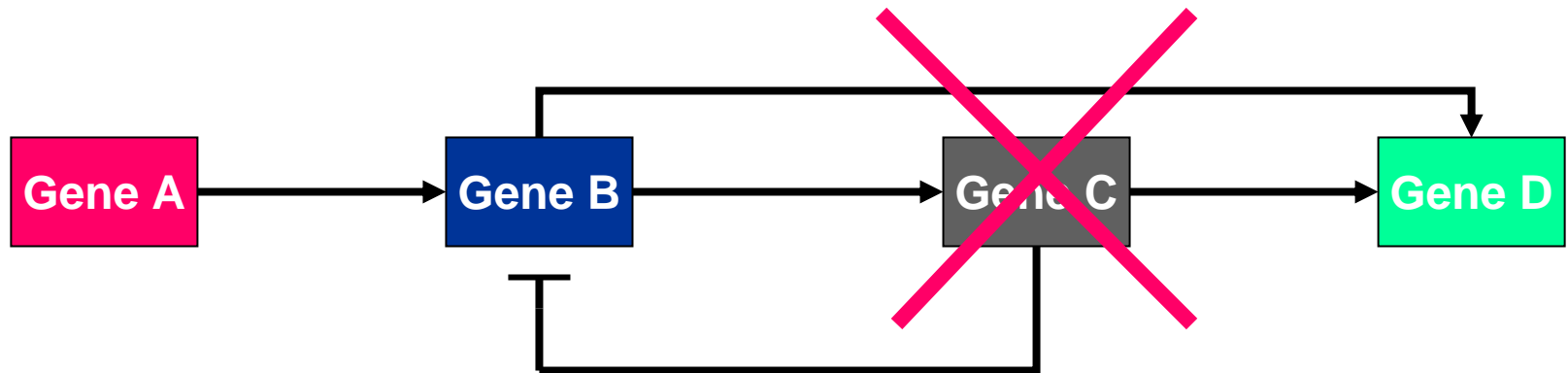*What are the networks of gene pathways that co-regulate gene expression of an organism?*



Wnt/Ca – calmodulin pathway

Bmp pathway

Retinoid acid pathway

Draft Pathways for Photoreceptor Function

Source: J. Yu, UM BioMedEng Thesis Proposal (2002)
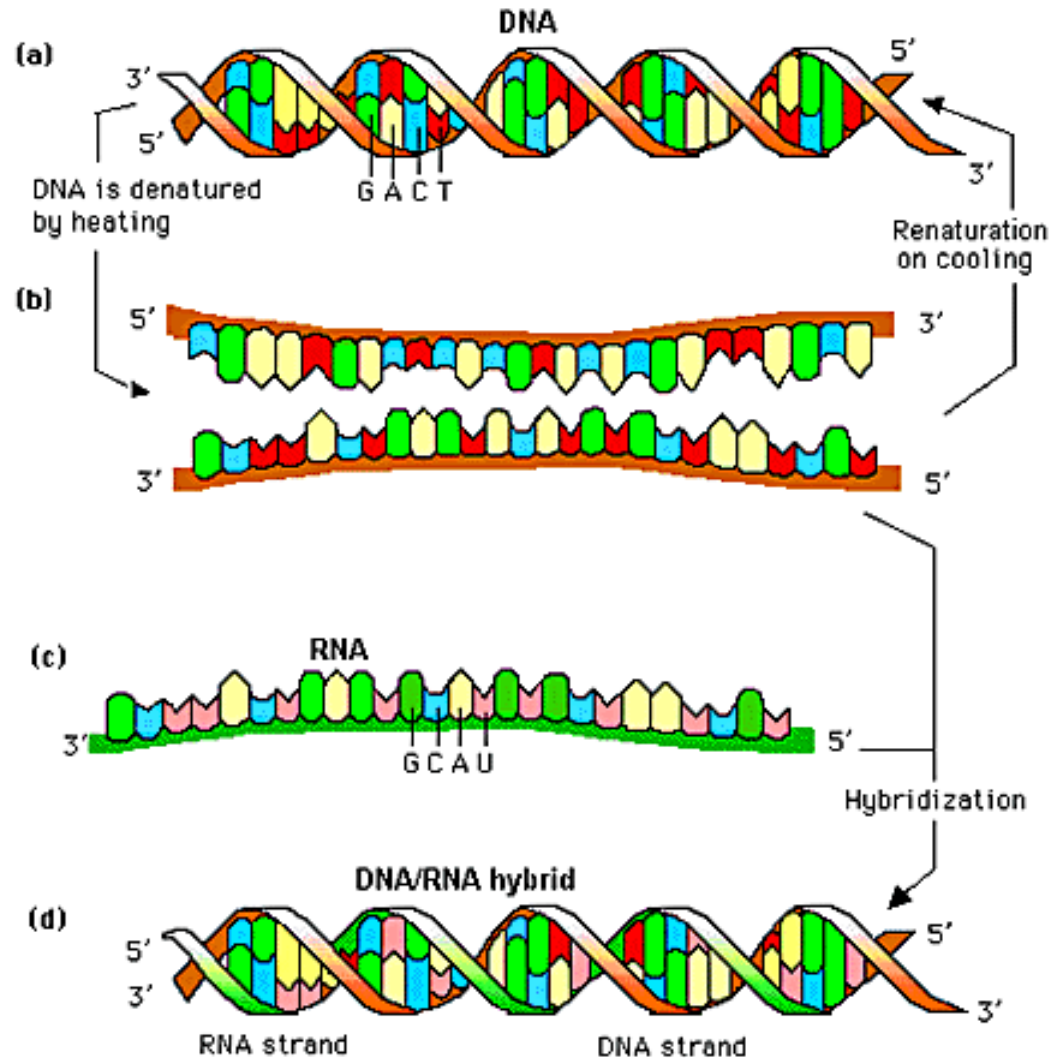
# Experimental Design for Structure Discovery

- Treatment level experiments: aging, starvation, drugs
- Gene knockout experiments: create a mutant organism



- Issues:
  - For a network of G genes require 2^G knockouts per time point to explore full co-regulation network.
  - Experimental replication is necessary ("large p small n")
  - There are other factors affecting gene expression: co-expression level,environment, protein-protein interactions…
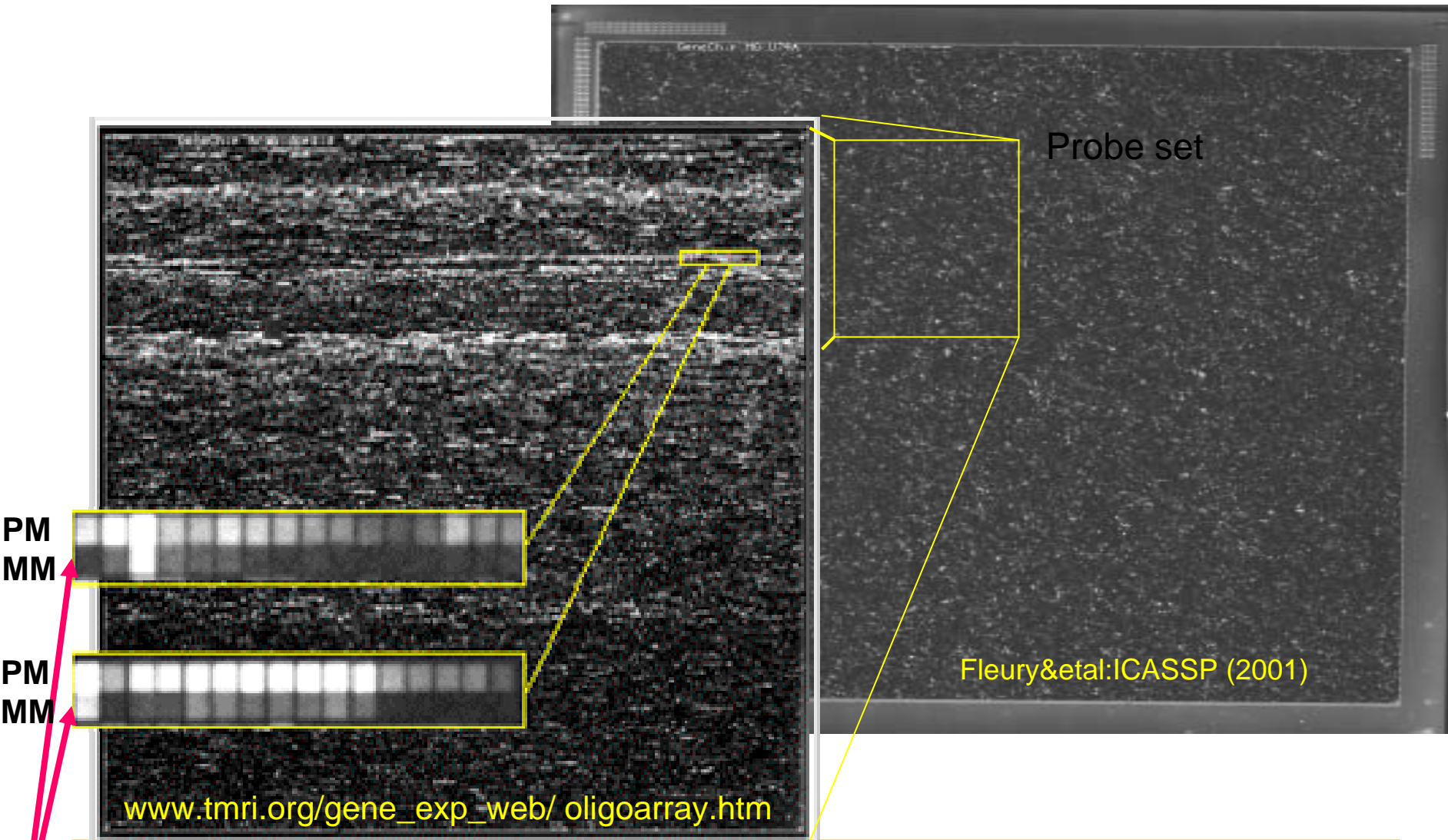
# Fundamental probing tool: hybridization



Nucleic Acid Hybridization

Plenary 2005

# Gene Microarrays

- Two principal gene microarray technologies:
  - Oligonucleotide arrays: (Affymetrix GeneChips)
    - Matched and mismatched oligonucleotide probe sequences photoetched on a chip
    - Dye-labeled RNA from sample is hybridized to chip
    - Abundance of RNA bound to each probe is laser-scanned
  - cDNA spotted arrays: (Brown/Botstein)
    - Specific complementary DNA sequences arrayed on slide
    - Dye-labeled sample mRNA is hybridized to slide
    - Presence of bound mRNA-cDNA pairs is read out by laser scanner

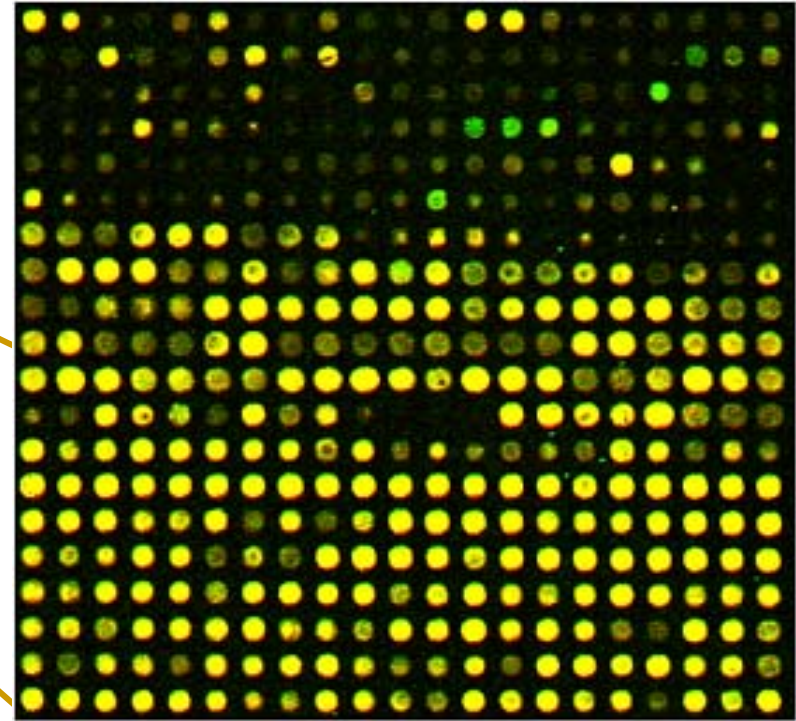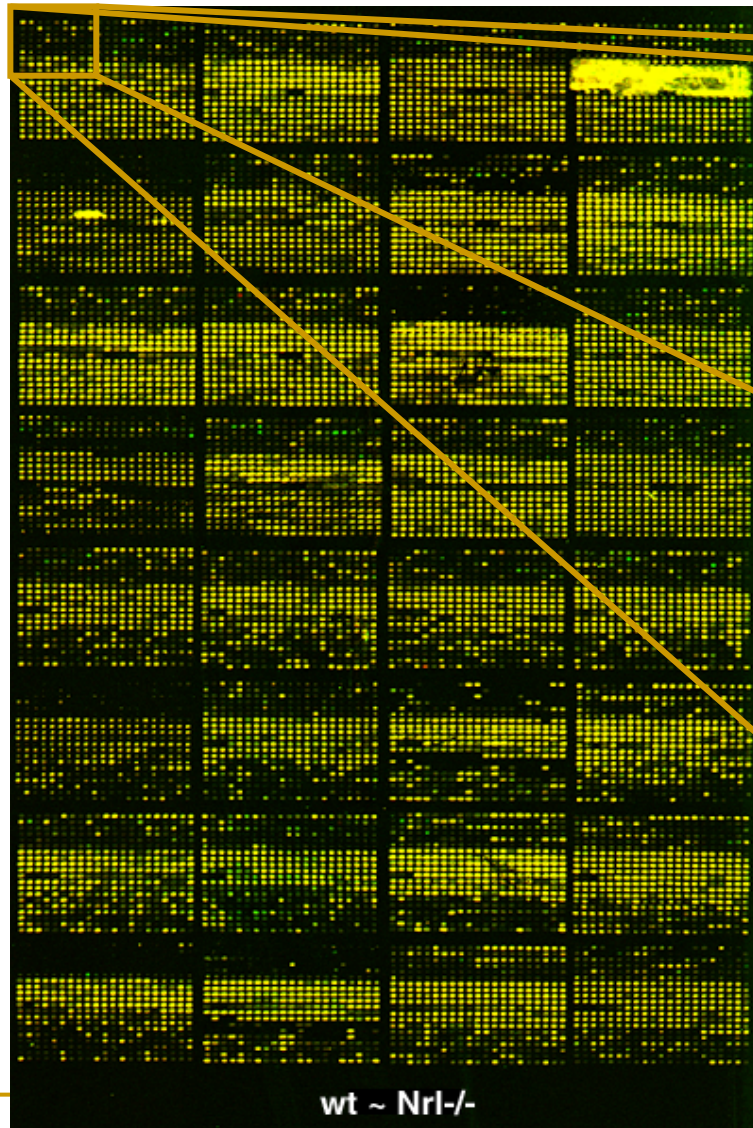- **10,000-50,000 genes can be probed simultaneously**
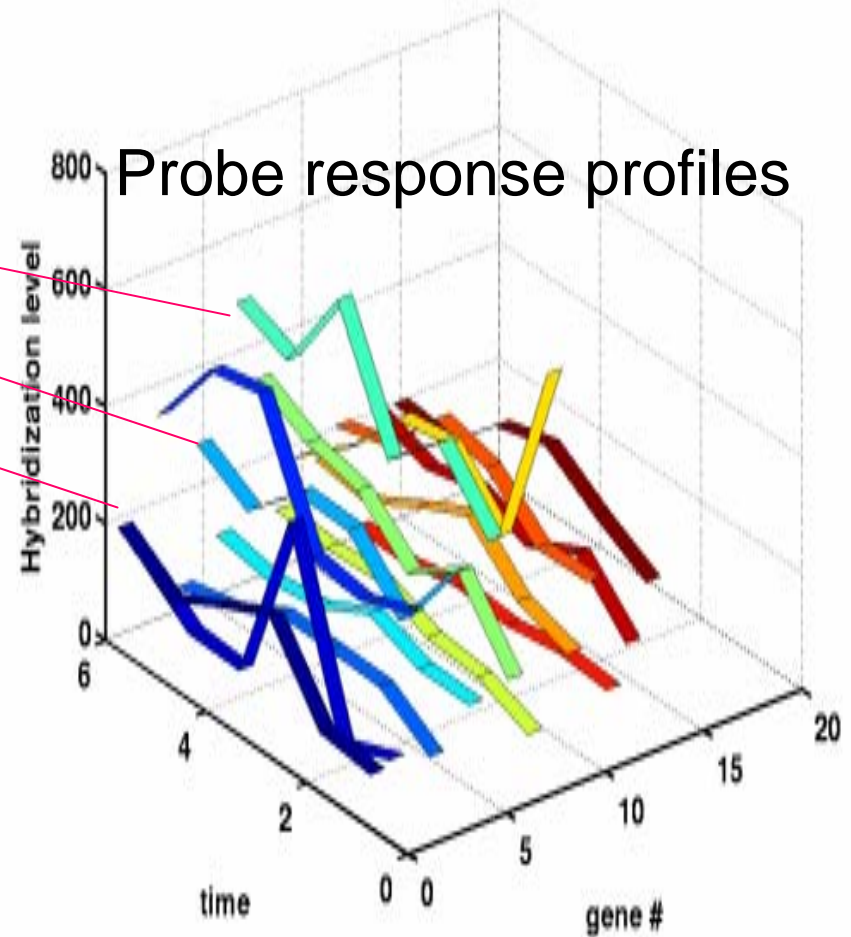
# Oligonucleotide GeneChip (Affymetrix)



Probe set

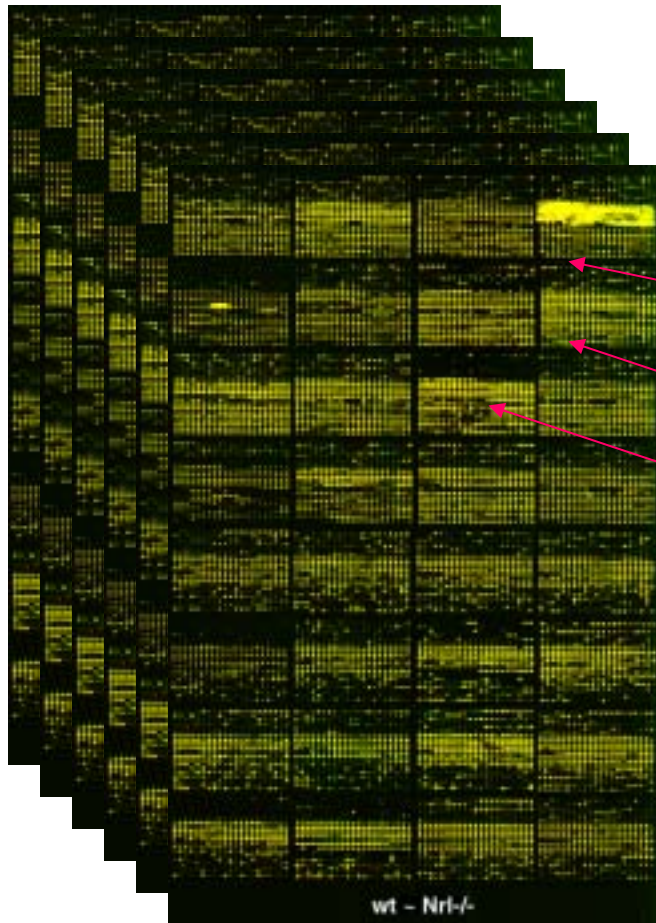GeneChip HG U74A

Fleury&etal:ICASSP (2001)

**PM**
**MM**

**PM**
**MM**

www.tmri.org/gene_exp_web/ oligoarray.htm

Two PM/MM Probe sets

# cDNA spotted array



wt ~ Nrl-/-

- **Treated sample (ko) labeled red (Cy5)**
- **Control (wt) labeled green (Cy3)**

# Add Treatment Dimension: Expression Profiles
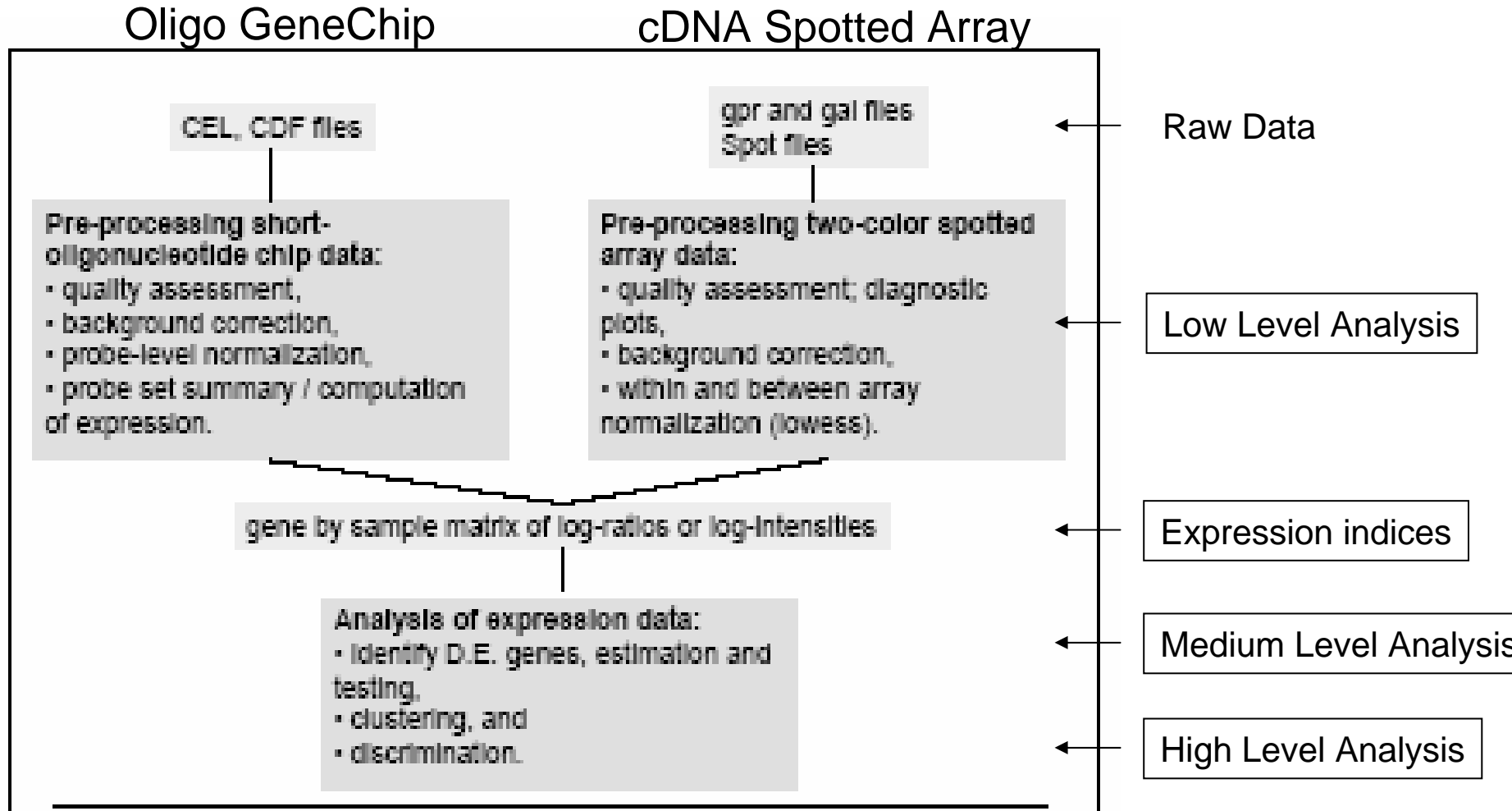


Probe response profiles

# Sources of Experimental Variability

- **Population** – wide genetic diversity
- **Cell lines** - poor sample preparation
- **Slide Manufacture** – slide surface quality, dust deposition
- **Hybridization** – sample concentration, wash conditions
- **Cross hybridization** – similar but different genes bind to same probe
- **Image Formation** – scanner saturation, lens aberrations, gain settings
- **Imaging and Extraction** – misaligned spot grid, segmentation

Microarray data is intrinsically statistical and replication is necessary
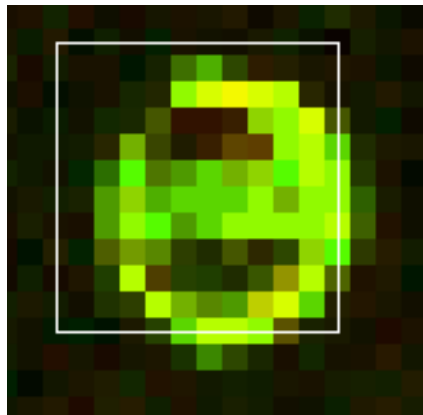
15

Churchill (2002)

# 2. Preprocessing of Gene Microarray Data

Oligo GeneChip                    cDNA Spotted Array

CEL, CDF files

gpr and gal files
Spot files                                    Raw Data

Pre-processing short-
oligonucleotide chip data:
· quality assessment,
· background correction,
· probe-level normalization,
· probe set summary / computation
of expression.

Pre-processing two-color spotted
array data:
· quality assessment; diagnostic
plots,
· background correction,
· within and between array
normalization (lowess).                        Low Level Analysis

gene by sample matrix of log-ratios or log-intensities          Expression indices

Analysis of expression data:
· Identify D.E. genes, estimation and
testing,
· clustering, and
· discrimination.

Medium Level Analysis

High Level Analysis

Source: Jean Yee Hwa Yang Statistical issues in design and analysis microarray experiment. (2003)

# Image Processing: cDNA Spot Extraction

- **Addressing** – Locate "center of description" for each spot
- **Spot Segmentation** – Classification of pixels either as signal or background.
- **Spot Quantification** – Estimation of hybridization level/ratio of spot



Grid misalignment



Laser Misalignment

Source: C. Ball, Stanford Microarray Database

**Refs**: Spotfire, ScanAnalyze, GenePix, Quantarray, Spot

# Spot Segmentation

- Threshold based
- Boundary based
  - Fixed circle
  - Adaptive circle (*used in QuantArray*)
  - Fixed Spot Mask (*used in ScanAlyze*)
- Region based
  - Seeded Region Growing (*used in Spot*)
- Active contours: level set algorithms
- Morphological operators: watershed segmentation

# Segmentation via Morphological Operators


Original Image


Alternate-Sequential Filtered


Watershed Transformed


Final Segmented Image

# A vs B Microarray Normalization

Exp A

Inverse

Unif Tran

Mean

Normalized A

Housekeeping Gene Selector

Normalized B

Unif Tran

Mean

Exp B

Inverse

# Pooled Microarray Normalization



Graphs are generated using R plot function hist() and boxplot()

Data: Lemon WJ et al. 2002

# Post-Normalization Histogram



Graphs are generated using <u>R</u> plot function hist() and boxplot()

Data: Lemon WL et al. 2002

# Extracting Expression Indices

- Each probe response level in microarray can be modeled via general mixed model

$$y_{gtr} = f_{gt}(\beta) + \rho_{gt}(\beta)Z_r + \sigma_{gt}(\beta)\epsilon_{gtr}$$

- g=gene probe index, t=timepoint, r=replicate
- $f_{gt}(\beta)$ is fixed effect
- $\sigma_{gt}(\beta)Z_r$ is random effect that may correlate t,g
- $\sigma_{gt}(\beta)\epsilon_{gtr}$ is noise component
- Special cases: MAS5, DChip, RMA. SMA, GEE
- Model similar to those used in array signal processing, statistical imaging, and other SP applications.

# 3. Screening Differentially Expressed Genes

12 knockout/wildtype mice in 3 groups of 4 subjects (24 GeneChips)

Knockout                                    Wildtype

$$\max_t\{\overline{K_t}(g) - \overline{W_t}(g)\} > \text{fcmin}$$

# Biological vs Statistical Significance

- **Biological significance** refers to foldchange being sufficiently large to be biologically meaningful or testable, e.g. testable by RT-PCR

$$|\text{fc}(g)| > \text{fcmin}$$

- **Statistical significance** refers to foldchange being different from zero

$$\text{fc}(g) \neq 0$$

Hero,Fleury,Mears,Swaroop:JASP2003

# Single Comparison Test

- Let $fc_t(g)$ = foldchange of gene 'g' at time point 't'.
- We wish to test the hypotheses:

$$H_0(g,t) \quad : \quad |\mathsf{fc}_t(g)| \leq |d|$$
$$H_1(g,t) \quad : \quad |\mathsf{fc}_t(g)| > |d|$$

- d = minimum acceptable difference (MAD)
- Method: confidence interval test

# Confidence Interval Test: Single Comparison

- Biologically&statistically **significant** differential response at 10% level of significance

$$f_{T_t(g)}(x|H_0)$$

$$T_t(g) = \frac{\overline{W}_t(g) - \overline{K}_t(g)}{\widehat{\sigma}_t(g)}$$

$-\mathcal{T}_{0.95}^{-1}$  **0**  $\mathcal{T}_{0.95}^{-1}$  **d**  $x$

**Conf. Interval on** $\text{fc}_t(g)$ **of level 1-alpha**

# Confidence Interval Test: Single Comparison

- Statistically significant but biologically **insignificant** fc



$$f_{T_t(g)}(x|H_0)$$

$-\mathcal{T}_{0.95}^{-1}$   **0**   $\mathcal{T}_{0.95}^{-1}$   **d**

**Conf. Interval on** $\mathrm{fc}_t(g)$ **of level 1-alpha**

# Multiple Comparisons: FWER, FDR

- **Pvalue, CI** apply to single comparison:



False positive

$$\mathcal{T}_{0.95}^{-1}$$

$$-\mathcal{T}_{0.95}^{-1}$$

g

- **FWER, FDR** and **FDRCI** depend on $\{T(g), g=1, \dots G\}$.

  - FWER: familywise error rate
    - Avg number of experiments yielding at least one false positive
  - FDR: false discovery rate (Benjamini&Hochburg:1996)
    - Avg proportion of false positives in experiments
  - FDRCI: $(1-\alpha)$ CI on discovered fc (Benjamini&Yekutieli:2002)
    - Avg. proportion of CIs that cover true fc in a given experiment

# Sorted FDRCI pvalues for ko/wt study



Sorted FDRCI p-values for various min fold changes

Legend:
- 0.32
- 0.58
- 0.85
- 1.00

$$T(g) = \max_t \frac{|\overline{W}_t(g) - \overline{K}_t(g)|}{\hat{\sigma}_t(g)}$$

$\alpha=0.2$

FDRCI p-value

Filtered genes at level (FDR=0.2, fc=0.32)

Ref:

# Screening Gene Expression Profiles

- Max foldchange is only one possible criterion of interest
- Objective: find the 250-300 genes having the most significant foldchanges wrt multiple criteria
- Example: Retinal aging study

# Multi-objective Optimization Approach

- Rarely does a linear order exist with respect to more than one ranking criterion, as in

$$|\mathsf{fc}_1(g_1)| > |\mathsf{fc}_1(g_2)| > \ldots > |\mathsf{fc}_1(g_p)|$$

- However, a partial order is usually possible

$$\{\mathsf{fc}_1(g), \ldots, \mathsf{fc}_6(g)\}_{g \in \mathcal{G}_1} > \ldots > \{\mathsf{fc}_1(g), \ldots, \mathsf{fc}_6(g)\}_{g \in \mathcal{G}_q}$$

# Illustration: two extreme cases

$$\xi_1(g) = fc_6(g) - fc_1(g) \text{ - end-to-end criterion}$$
$$\xi_2(g) = \min_t\{fc_t(g) - fc_{t-1}(g)\} \text{ -increasing criterion}$$

- A linear ordering exists
- No linear ordering exists



Optimum

# Pareto Front Analysis (PFA)

- Rank genes by peeling of successive Pareto Fronts

Hero&Fleury:VLSI04

# 4. Clustering  Gene Expression Patterns



- Gene expression levels over multiple conditions are required for pathway studies
- Requires symmetric similarity metric, e.g pairwise profile correlation

Source: Wing Wong Lab, Stanford (left)                    Swaroop Lab, Michigan (right)

# Drawback of Traditional Clustering

- Clustering using pairwise correlation fails to account for transitive co-expression (Zhou etal 2002)

# Extraction of Co-Regulation Circuits

**Gene A**

**Gene C**

**Gene D**

**Gene B**

$$p(\mathcal{X}) = \prod_{gt} p(x_{gt} | \mathcal{X}_{gt}^-)$$

| Gene A | Gene B | Gene C | Gene D |

# Modeling co-Regulation Networks

- ❏ Relevance networks
  - Edge = strong correlation
- ❏ Dependency networks
  - Directed edge = strong partial correlation
- ❏ Dynamical dependency networks
  - Directed edge = strong partial correlation
- ❏ Bayesian networks
  - Profiles are quantized to small number of bits
  - One bit quantization = boolean networks

# Network Constrained Clustering



- If topology were known could use to improve clustering
- Otherwise suffer from combinatorial explosion:

$$p = 2^{\binom{G}{2}}$$

- Soln: FDRCI edge screening

# FDRCI Edge Screen Procedure

- Fix FDR level and MAS level on discovered edges

- Construct FDRCI's of desired FDR level on edge strengths

- Accept edge if FDRCI exceeds MAS

**Pearson correlation coefficient**



Dongxiao Zhu, A. Hero, S. Qi, JCB, 2004.

# Yeast Galactose Metabolism Experiment

- 10 different yeast strains (9 gene knock-outs and 1 wild type) incubated in either GAL-inducing or non-inducing media (Ideker et al. 2001).

- 9 gene knock-outs are GAL1, GAL2, GAL3, GAL4, GAL5, GAL6, GAL7, GAL10, GAL80.

- 5935 gene 2-channel cDNA array. Reference channel is dilution "wild-type + galactose"



Dongxiao Zhu, A. Hero, S. Qi, JCB, 2004.

# Relevance Network Visualization
## (FDR $<= 0.05$, MAS $= 0.9$)



Dongxiao Zhu, A. Hero, S. Qi, JCB, 2004.

# Network Constrained Clustering



Clustering with prior distance matrix

Clustering with posterior (shortest-path) distance matrix

# Horizons: Transcriptomics/Proteomics Technology

- Higher throughput cDNA/GeneChip microarrays

- Suspension microarrays

- Microscale "Lab on a chip"

- Protein-protein arrays

- Nuclear magnetic resonance spectroscopy

- In vivos molecular imaging: reporter genes

# Where does SP fit?



Experiments

Biology

SP

Math Modeling

Statistics

CS

Processing

# Signal Processing Opportunites

- **There is room for new SP approaches**
    - Non-modularized analysis: task-driven and top-down?
    - Active waveform design: sequential design of experiments?
    - Internet Tomography: gene network topology discovery?
    - MIMO: spatio-temporal wideband array processing?
    - Channel optimization: optimal gene layout on microarray?
- **New technology is appearing that offers opportunities for SP'ers to develop models/algorithms**
- **There is still some low lying fruit!**
- **Collaboration with a biological scientist is essential in order to have impact**

# Where to learn more?



Genetics, the painless way, 1991



Historical overview by one of the pioneers, 2003



Basic undergraduate texbook, 1992

# Where to learn more?

Overview of microarray technology and analysis, 2003

Edited volume on principal statistical techniques of microarray analysis, 2003

Textbook aimed at biostatisticians, 2004

# Where to learn more?

Edited monograph on
random graphs
In nature, 2005

To appear soon!

To appear soon!

# Future venues



**GENSIPS 2005 – Newport**:

- Ray Liu, Jaako Astola
- Workshop dates: May 22 - 24, 2005
- Early registration ends March 30

**IEEE Transactions on Signal Processing Special Issue on Genomic Signal Processing**

- Submission deadline: May 1, 2005
- Publication date: Sept. 2006

# 5. Conclusions

- Gene filtering: accounting for biological and statistical significance

- Gene ranking: can involve optimization over multiple criteria

- Gene co-regulation networks: discover co-dependent gene profiles that can aid in clustering

- Statistical signal and image processing approaches can have impact

- References to UM work and software presented here: http://www.eecs.umich.edu/~hero/bioinfo.html