



# Machine learning in high-throughput genomics and proteomics

(Part I)

Xuegong Zhang, Ph.D.  
Professor of Pattern Recognition and Bioinformatics  
Tsinghua University  
zhangxg@tsinghua.edu.cn

## Outline

- Background
  - Microarray and mass-spectrometry technologies
  - Typical bioinformatics problems
  - Basic concepts of machine learning
- Microarray data mining with learning machines
  - Classification, clustering and gene selection
  - Problems, solutions, and possible pitfalls
- Further topics
  - Important considerations and open problems
  - Systems biology study after microarray data mining
  - New types of arrays

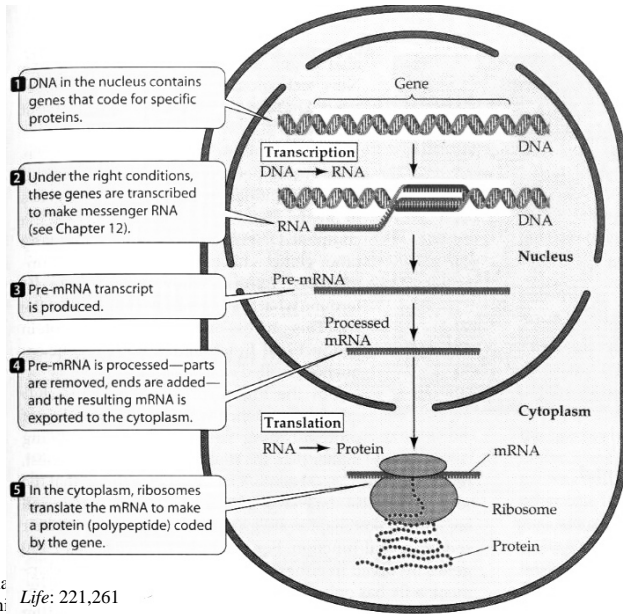
## Target Audience:

- Researchers/students in information science, math or statistics who are interested in biological applications and bioinformatics
- Biologists who want to learn more about the underlying concepts behind the widely-used bioinformatics tools

## Part I Biological Backgrounds

- Microarray and mass spectrometry: technologies and data
  - Typical bioinformatics topics
- Basic concepts of Machine Learning

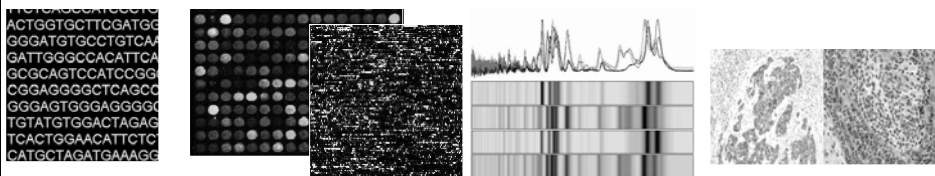
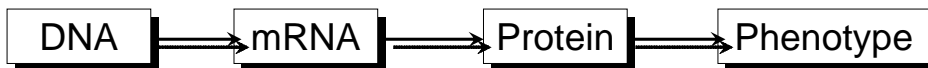
# The Central Dogma



Xuegong Zhang  
Tsinghua University *Life*: 221,261

5

# The Central Dogma



DNA Microarrays

Xuegong Zhang  
Tsinghua University

6

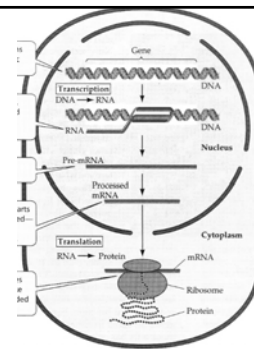
# Gene Expression and Microarrays

Xuegong Zhang  
Tsinghua University

7

## Gene Expression

- Gene → Primary Transcript  
→ nuc. mRNA → cytosolic mRNA  
→ protein → protein activity
- Basic cellular processes are realized by tightly regulated gene expression programs
- Different cell types in a multicellular organism express different sets of genes at different time and with different quantities



Xuegong Zhang  
Tsinghua University

8

## Measuring mRNA abundance in cytosol (*gene expression levels*)

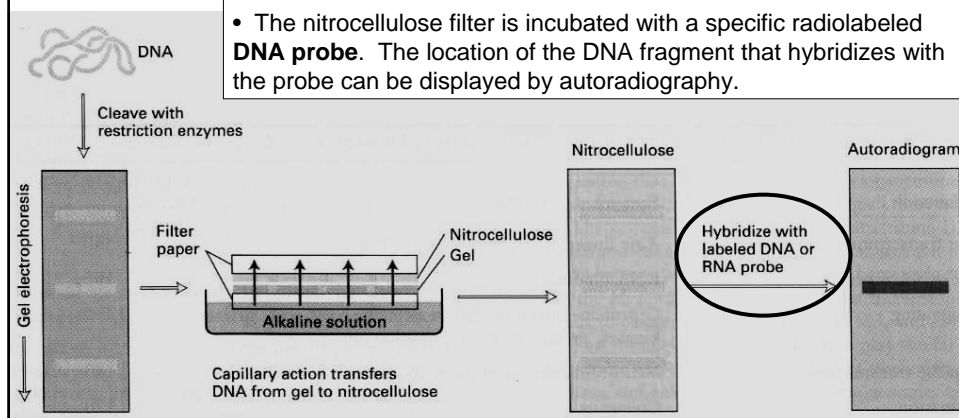
- Northern blotting
  - The traditional method
  - **Considered as the gold standard**
  - Gene by gene
  - Very time consuming
- Microarray

Xuegong Zhang  
Tsinghua University

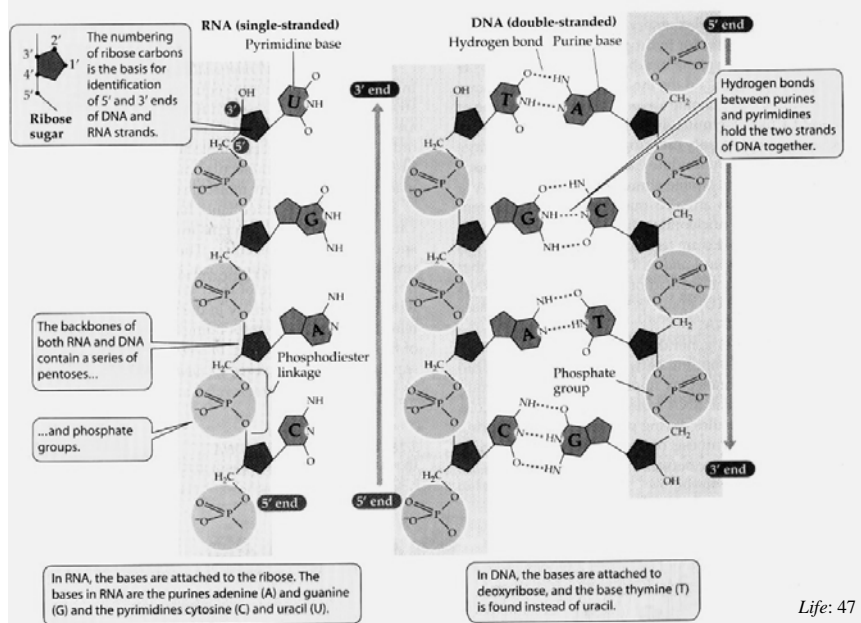
9

**Southern Blotting** [E.M. Southern, 1975, JMB, 98:508] for detecting the presence of specific DNA sequences following gel electrophoresis of a complex mixture of restriction fragments.

- The DNA to be analyzed is digested with restriction enzymes and then separated by agarose gel electrophoresis.
- The DNA fragments in the gel are denatured with alkaline solution and transferred onto a nitrocellulose filter or nylon membrane by blotting, preserving the distribution of the DNA fragments in the gel.
- The nitrocellulose filter is incubated with a specific radiolabeled **DNA probe**. The location of the DNA fragment that hybridizes with the probe can be displayed by autoradiography.



## The Key Principle: Complementary Base-Pairing



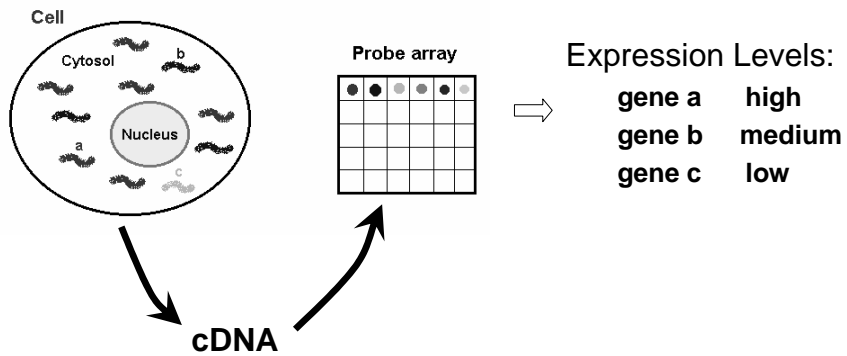
- Northern Blotting

- Detecting RNA fragments, instead of DNA fragments
- RNA fragments are treated with formaldehyde to ensure linear conformation
- The amount of a specific RNA in a sample can be estimated from a Northern blot

- Western Blotting

- Detecting a particular protein in a mixture
- The probe used is antibodies
- Also called *immunoblotting*

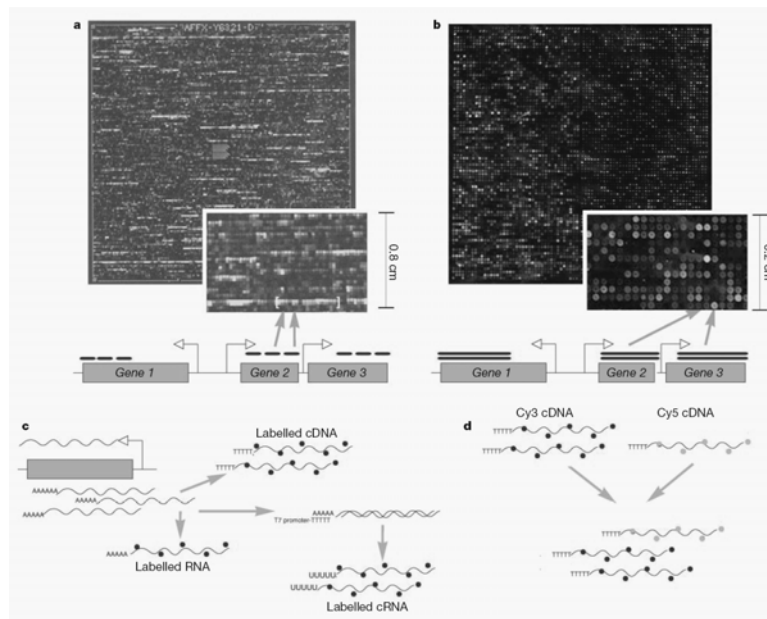
## Going high-throughput: the basic idea of microarrays



Xuegong Zhang  
Tsinghua University

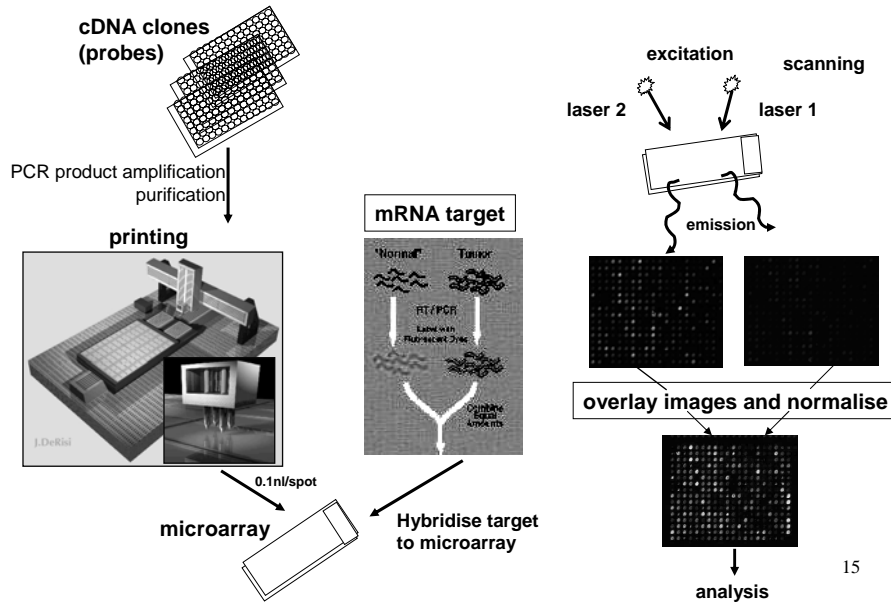
13

## Two ways of making microarrays

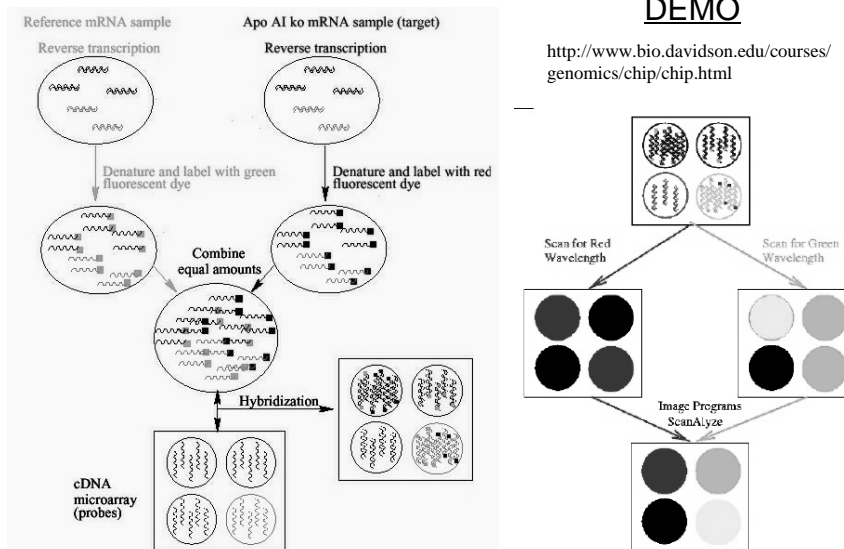


D.J. Lockhart & E.A. Winzler, Genomics, gene expression and DNA arrays, *Nature*, 405(15): 827-836, 2000

# Printed cDNA microarrays



# Printed cDNA microarrays

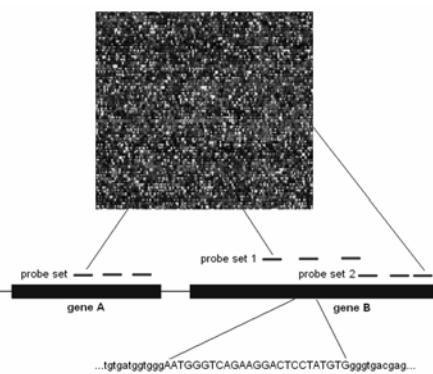




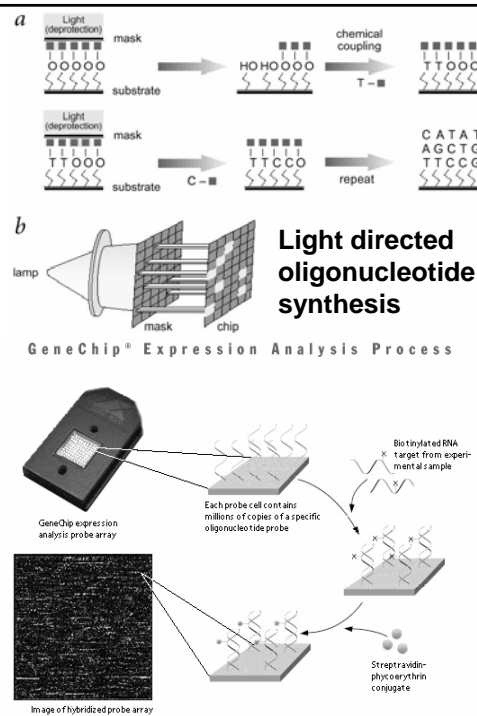
# Oligonucleotide microarrays



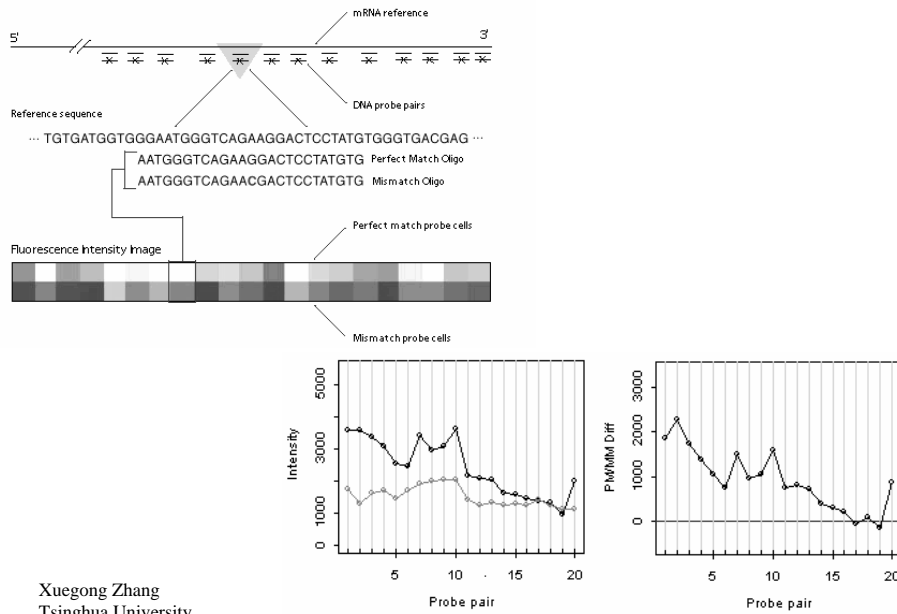
## The GeneChip<sup>®</sup> microarray



Xuegong Zhang  
Tsinghua University

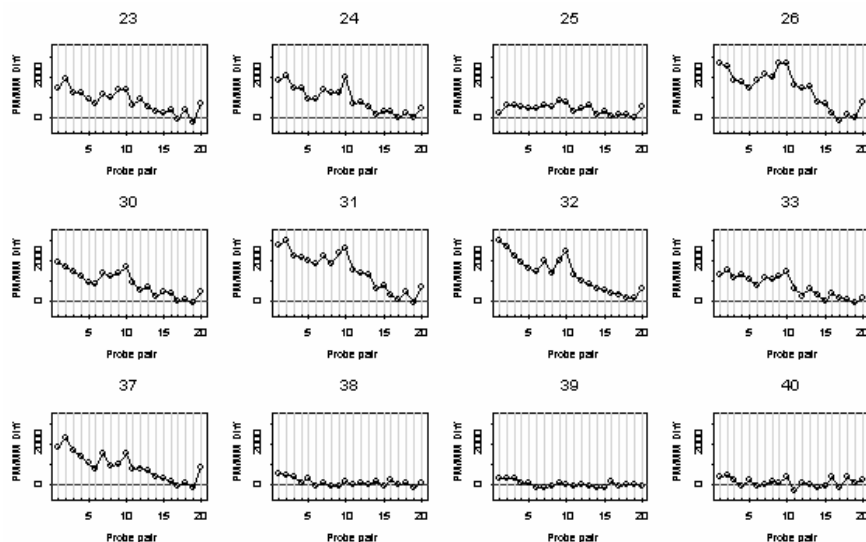


## PM/MM Probes, Probe-pairs, Probe-sets



Xuegong Zhang  
 Tsinghua University

## Example: Data for one gene in many arrays



20

## Getting the expression from the probes: the Li-Wong Model

An multiplicative model for one gene :

$$Y_{ij} = PM_{ij} - MM_{ij} = \theta_i \phi_j + \varepsilon_{ij}$$

$\theta_i$  : expression level index for array  $i$  ,  
comparable to PM/MM difference level,  
assumed to be proportional to true expression level

$\phi_j$  : probe sensitivity index for probe pair  $j$  ,  
subject to constraint  $\sum_j \phi_j^2 = J$  for identifiability

$\varepsilon_{ij}$  : random normal errors  
Tsinghua University

Cheng Li & Wing H. Wong, Model-based analysis of oligonucleotide arrays: Expression index computation and outlier detection, *PNAS*, vol.98, no.1, pp.31-36, 2001

## dChip (DNA-Chip Analyzer)

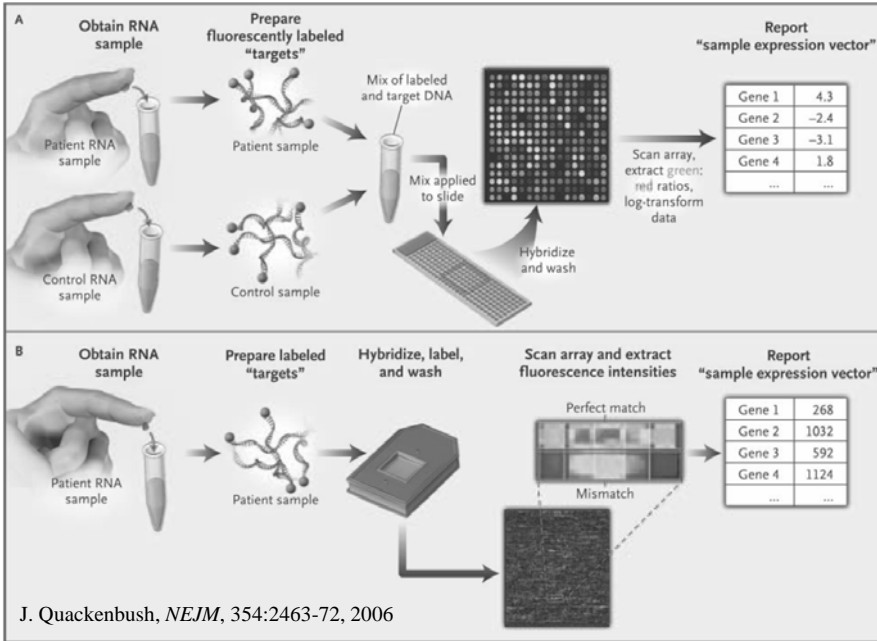
[www.dchip.org](http://www.dchip.org)

- Data Processing
  - ...
- Probe-level Analysis
  - Normalize arrays
  - View data
  - Model-based expression values
  - Outlier array
  - Exon and tiling array analysis
  - ...
- High-level Analysis
  - Clustering
  - Time course analysis
  - Classification
  - Significant genes
  - Gene mapping
  - Pathway analysis
  - ...
- SNP Array Analysis
  - LOH and copy number
  - Linkage analysis
  - ...

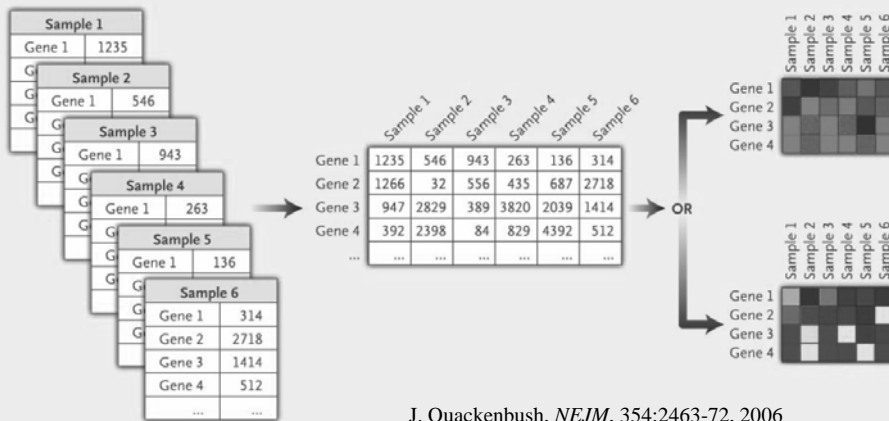
Xuegong Zhang  
Tsinghua University

22

# Overview of microarray analysis



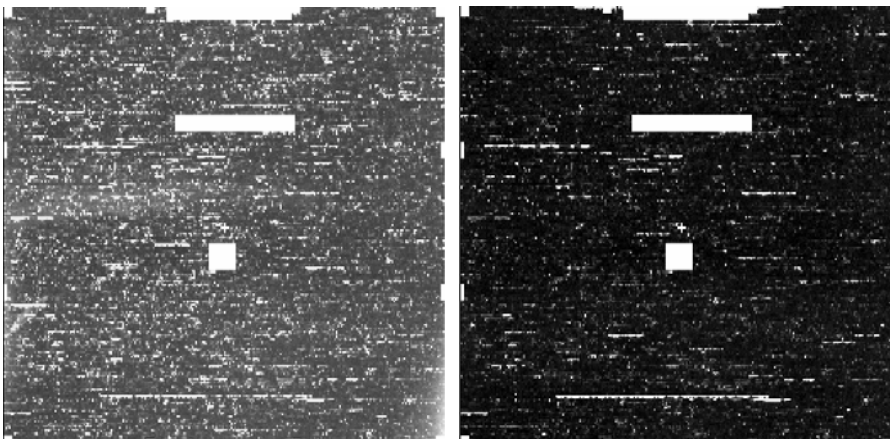
# The Expression Matrix



## Example Data Sheets

The image displays two screenshots of Microsoft Excel spreadsheets. The top screenshot shows a table with columns labeled A through J. Row 1 is a header row with labels: A: probe set, B: gene, C: N1 call, D: N3 call, E: N4 call, F: N5 call, G: N1 call, H: N3 call, I: N4 call, J: N5 call. Rows 2-5 contain data for probes AFFX-BioEJ04423, AFFX-BioEJ04423, AFFX-BioEJ04423, and AFFX-BioEJ04423. The bottom screenshot shows a similar table with columns labeled A through J. Row 1 is a header row with labels: A: probe set, B: gene, C: A, D: B, E: C, F: D, G: E, H: F, I: G, J: H. Rows 2-22 contain data for various probes including HG4243-HIzinc fing, HG4245-HIForkhead, HG4258-HIcyclin-de, HG4263-HIkiller ce, HG4272-HImet protc, HG429-HIcB-cell gr, HG4297-HIactivated, HG4310-HIretinol-b, HG4316-HITransketc, HG4319-HIRibosomal, HG4321-HIAHNAK-rel, HG4332-HIzinc fing, HG4333-HIzinc fing, HG4336-HIBacterici, HG4390-HIRibosomal, HG4411-HIMucin, Ge, HG4433-HICyclin D1, HG4458-HIImmunogl, and HG4460-HIImmunogl.

## The Problem of Normalization



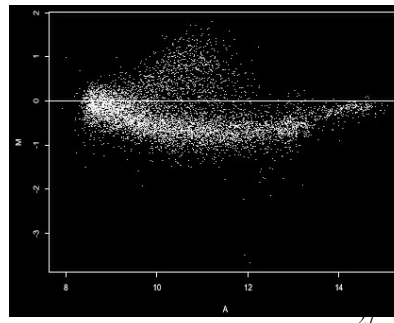
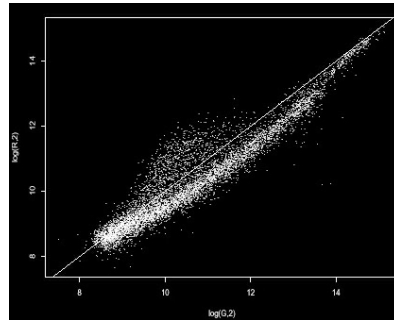
## The M-A plot

- LogR vs. LogG plot
- M vs. A plot (M-A plot)

$$A = (\log(\text{Cy5}) + \log(\text{Cy3})) / 2$$

$$M = \log(\text{Cy5} / \text{Cy3})$$

- An MA-plot amounts to a  $45^\circ$  counterclockwise rotation of a logR vs. logG plot followed by scaling



Xuegong Zhang  
Tsinghua University

## References for Normalization

- Housing keeping genes / spot-in genes
- Average or Median
- Rank-invariant genes

Xuegong Zhang  
Tsinghua University

28

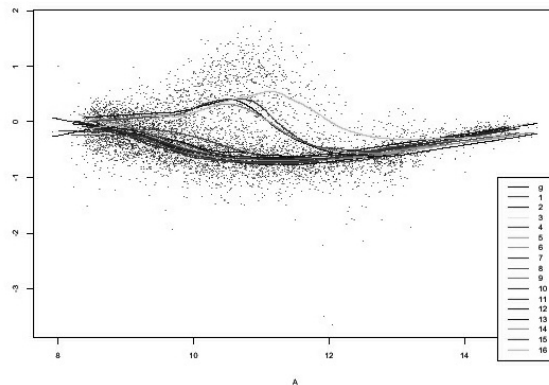
### Example: Pin-wise Normalization (Dudoit et al. 2000)

- Pin-wise normalization using all the genes
- Requires the assumption that up- and down-regulated genes with similar average intensities are roughly cancelled out or otherwise most genes remain unchanged
- Lowess (locally weighted least squares) smooth

$$\log R/G \rightarrow$$

$$\log R/G - c_j(A)$$

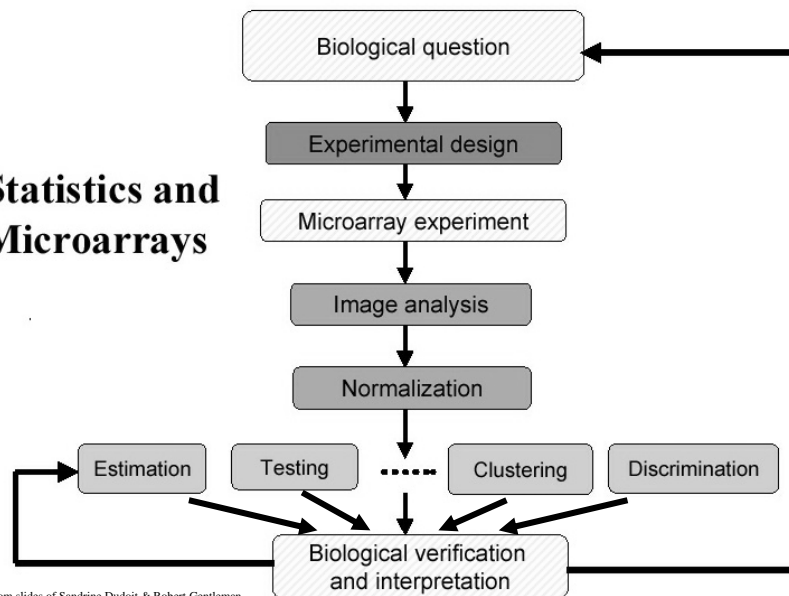
$$= \log(k_j(A) R/G)$$



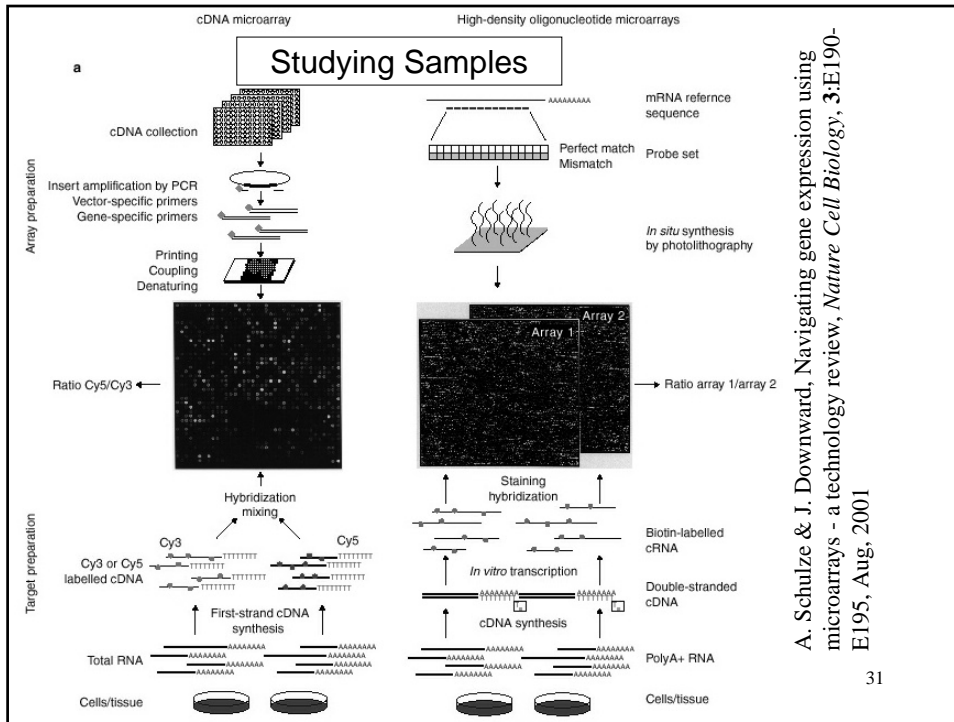
Xuegong Zhang  
Tsinghua University

### Bioinformatics analysis of microarrays

#### Statistics and Microarrays

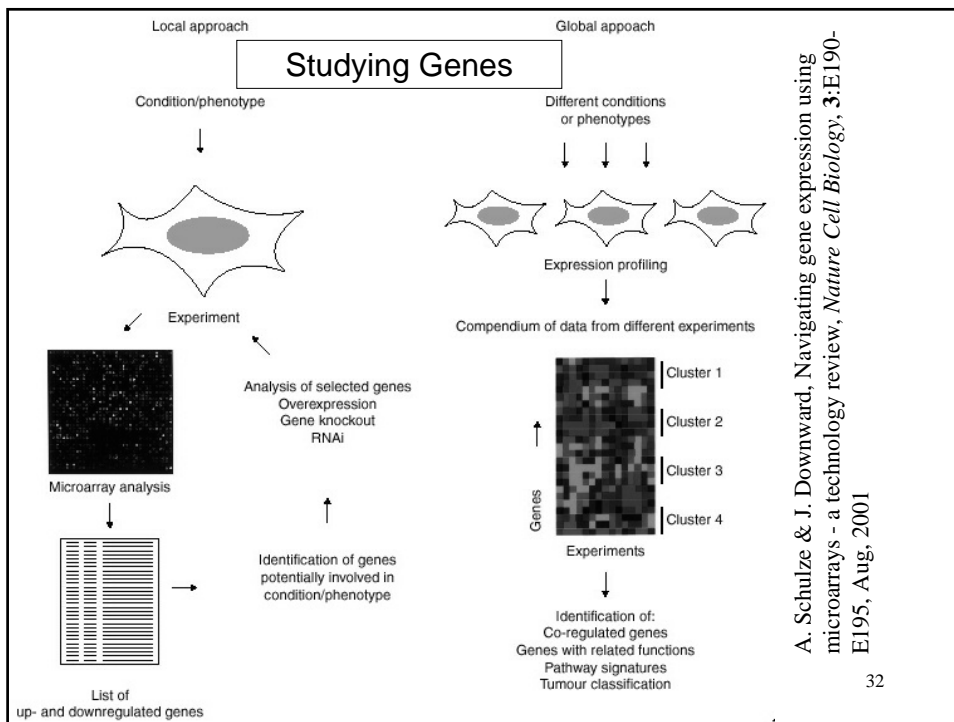


From slides of Sandrine Dudoit & Robert Gentleman



A. Schulze & J. Downward, Navigating gene expression using microarrays - a technology review, *Nature Cell Biology*, 3:E190-E195, Aug. 2001

31

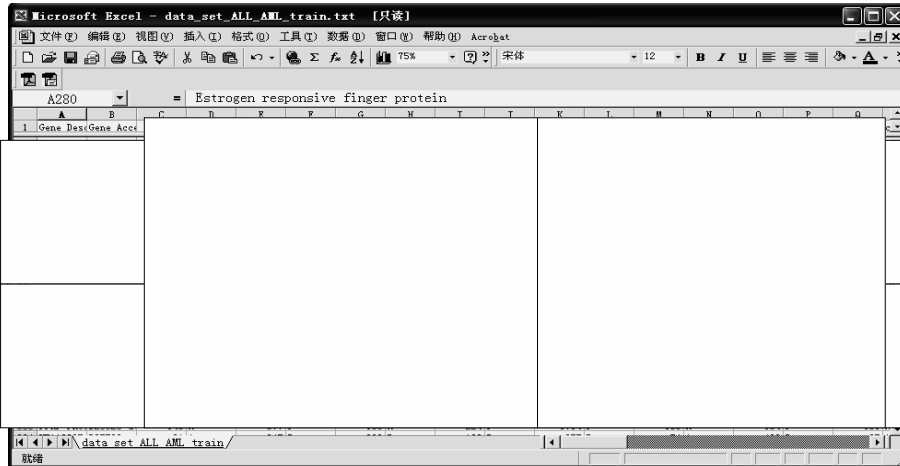


A. Schulze & J. Downward, Navigating gene expression using microarrays - a technology review, *Nature Cell Biology*, 3:E190-E195, Aug. 2001

32



## Classification study with multiple arrays



Xuegong Zhang  
Tsinghua University

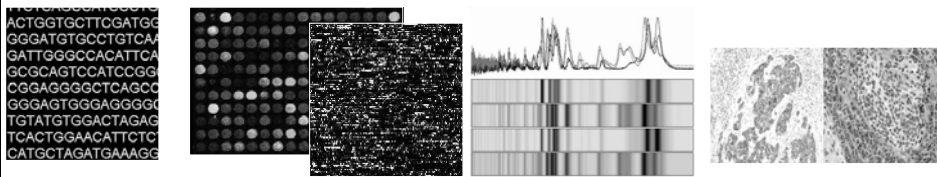
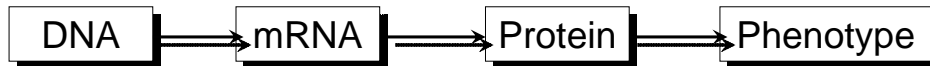
33

## Proteomics and Mass Spectrometry

Xuegong Zhang  
Tsinghua University

34

## Revisit the Central Dogma



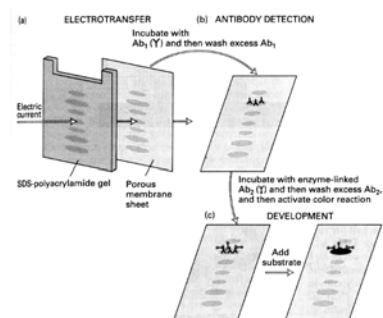
Mass Spectrometry

Xuegong Zhang  
Tsinghua University

35

## Detecting Proteins

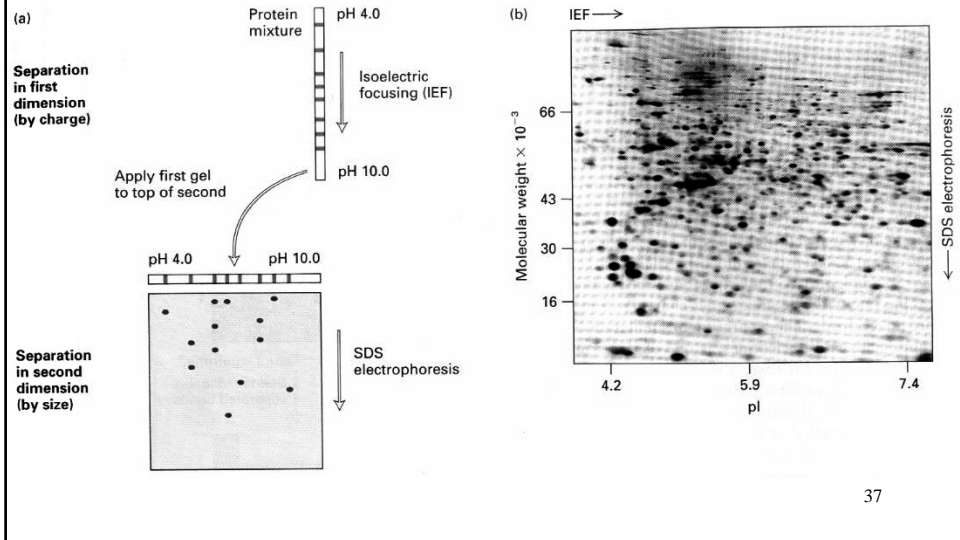
- Southern Blotting
  - Detecting DNA sequences
- Northern Blotting
  - Detecting RNA fragments
- Western Blotting
  - Detecting a particular protein in a mixture
  - The probes used are antibodies
  - Also called *immunoblotting*



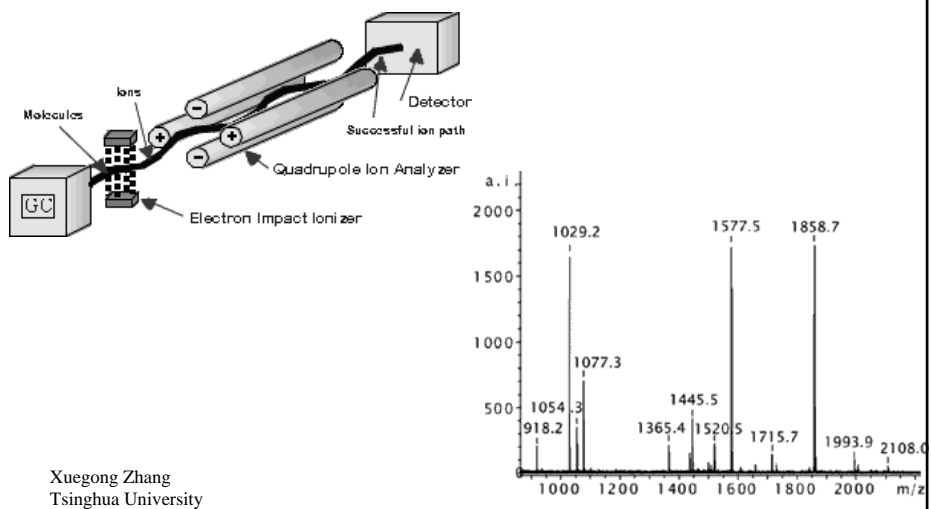
Xuegong Zhang  
Tsinghua University

36

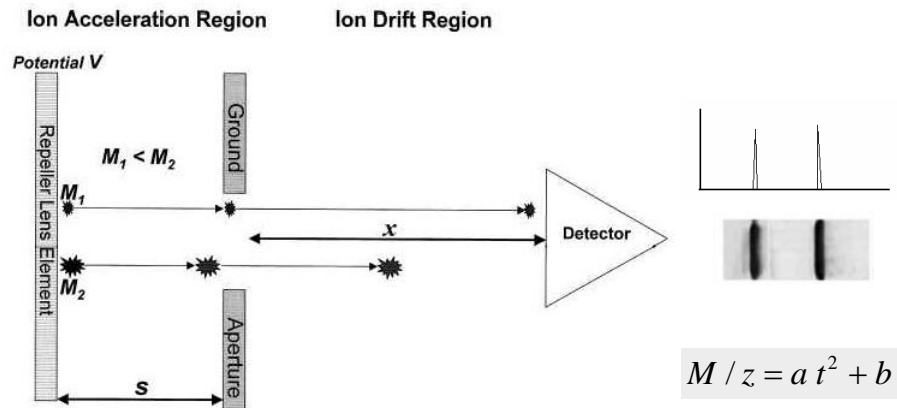
# 2D gel electrophoresis



# Mass spectrometry



# TOF-MS (Time-Of-Flight Mass Spectrometry)

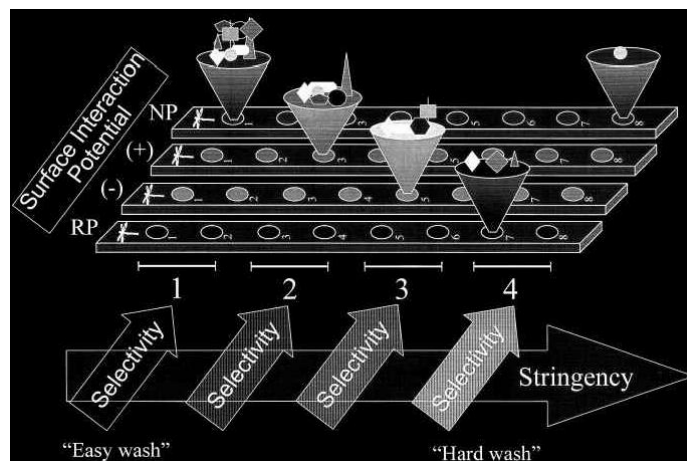


A mass spectrometry system uses 3 components – an ionization source, an analyser and a detector

39

# SELDI-TOF-MS: ProteinChip® from Ciphergen

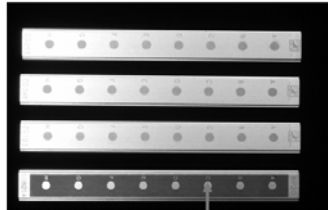
Surface Enhanced Laser Desorption/Ionization - Time Of Flight - Mass Spectrometry



40

# SELDI-TOF-MS

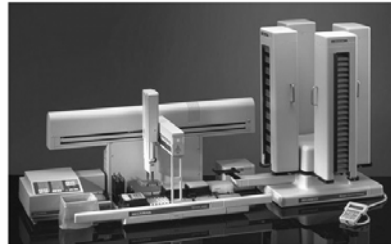
## Serum Sample Loading



One Microliter of Serum

### Robotic handling

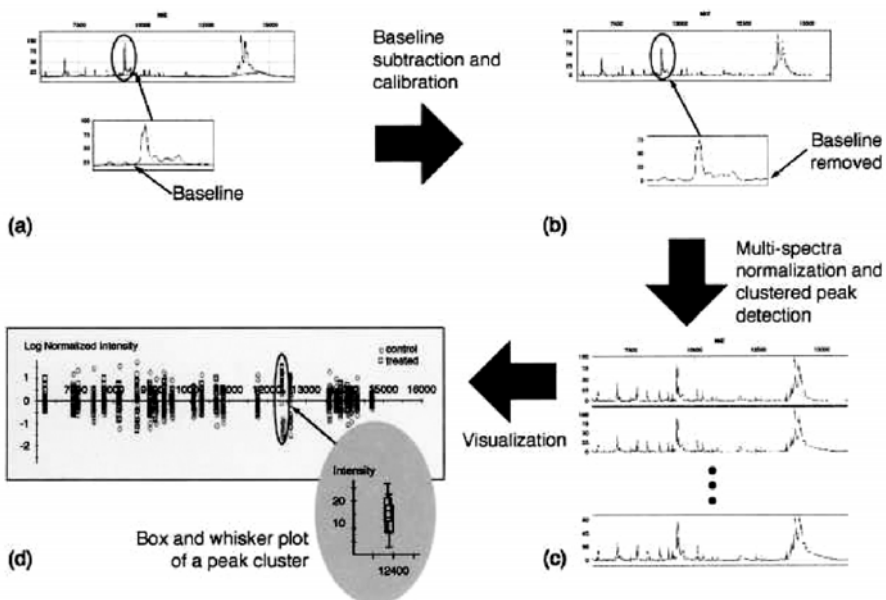
- Improved Reproducibility
- Better Sample Handling
- Increased Throughput
- Reduce Cross Contamination



X Vincent Fusaro and Sally Ross  
T

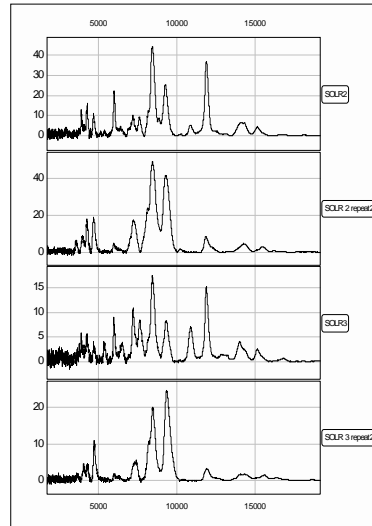
1

## Typical Data Processing Steps



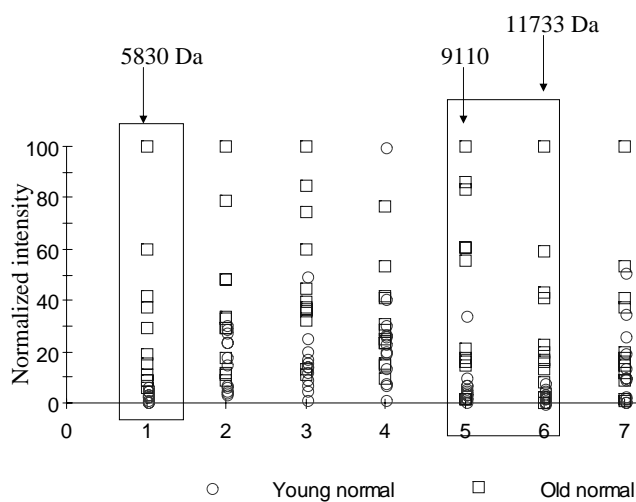
# Data Processing

- Baseline subtraction
- Filtering
- Calibration
- Normalization
- Peak Detection
- Peak Alignment and Biomarker Detection
- Assessment of the Expression Indexes



Xuegong Zhang  
Tsinghua University

# Biomarker discovery



Tsinghua University

# Mass-spectrometry data sheets

- Mass Spectrums
  - Time serieses
  - Could be discretized as expression profiles
- Biomarkers
  - Protein profiles at certain "meaningful" molecular weights
  - Still looks like gene expression data in format

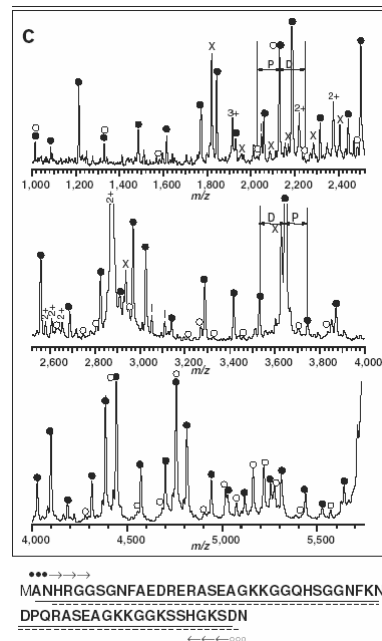
marker #	Case 1	Case 2	Case 3	...
1	123.4	345.6	23.4	
2	23.1	1234.6	3245.2	
3	123.2	987.4	654.3	
...				

m/z	Case 1	Case 2	...
2345.6	123.4	345.6	
2347.2	23.1	1234.6	
2349.1	123.2	987.4	
...			

Xuegong Zhang  
Tsinghua University

45

# Peptide Identification and AA Sequencing with MS/MS



Xuegong Zhang  
Tsinghua University

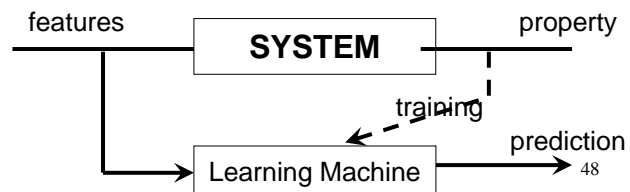
# Basic Concepts of Machine Learning

Xuegong Zhang  
Tsinghua University

47

## Machine Learning (ML): basic concepts

- A system (unknown but existent)
- A set of observable features (input)
- An interested property (output)
- The Learning Machine
- Predictions
- Training Samples
- The Learning Algorithm



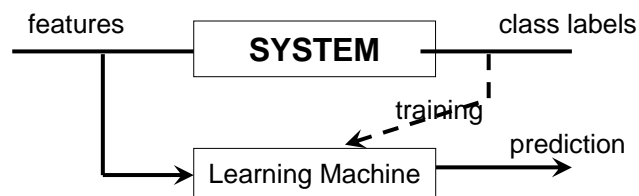
Xuegong Zhang  
Tsinghua University

48



## Pattern Recognition (PR)

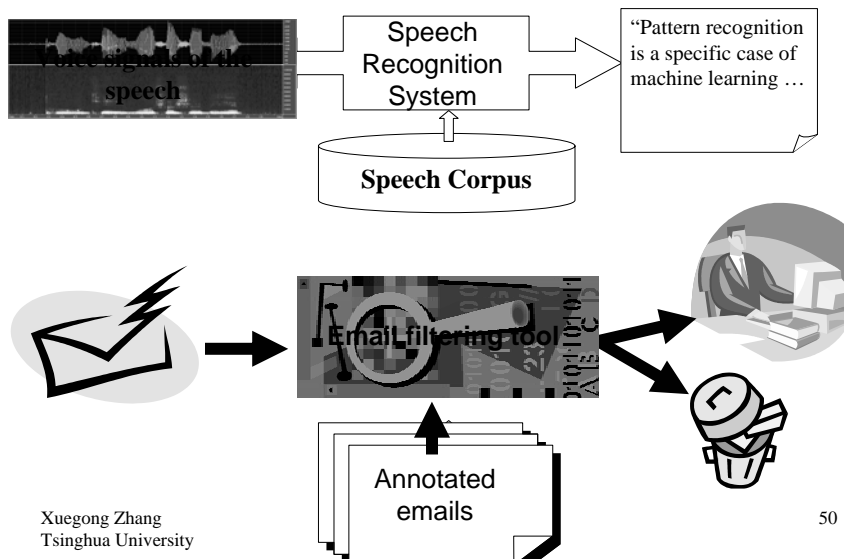
- ML for classification properties,  
e.g., normal vs. cancer, good vs. bad prognosis,  
subtypes of a disease, .....
- Pattern Classification



Xuegong Zhang  
Tsinghua University

49

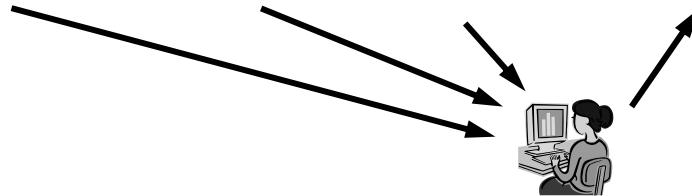
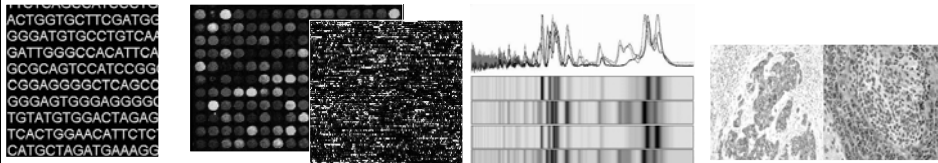
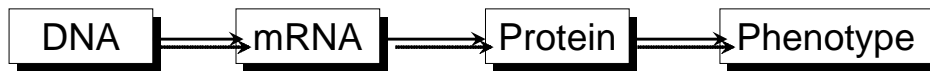
## Examples of PR systems



Xuegong Zhang  
Tsinghua University

50

## Let Machines to Learn the Central Dogma



Xuegong Zhang  
Tsinghua University

51