# Our view on cDNA chip analysis from engineering informatics standpoint

**Chonghun Han**, **Sungwoo Kwon**

Intelligent Process System Lab

Department of Chemical Engineering

Pohang University of Science and Technology
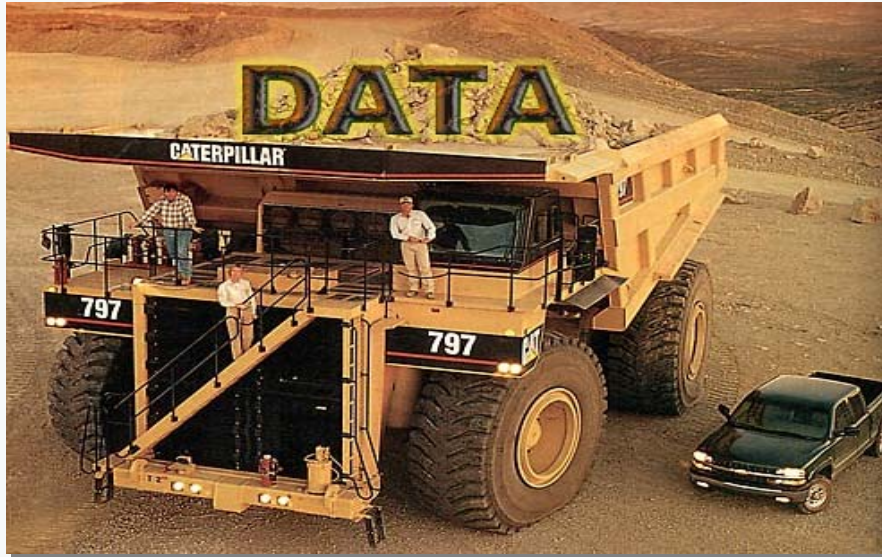
**POSTECH IPSE Lab.**

# Outline

- Introduction to Bioinformatics
- Introduction to cDNA chip
- Classification of Tumor Classes
- Identification of Marker Genes
- Conclusions and Future Works

# The Information Revolution



I need tools to extract important information from mountains of data

# Data Mining

■ **Data mining = 'exploration and analysis by automatic and semi-automatic means, of large quantities of data in order to discover meaningful patterns and rules'**

**Applications of data mining**

- ✓ **Search database**
- ✓ **Structural pattern recognition**
- ✓ **Medline abstract analysis**
- ✓ **DNA chip data analysis**

**Who has information and uses it wins**

**(Watterson, K)**

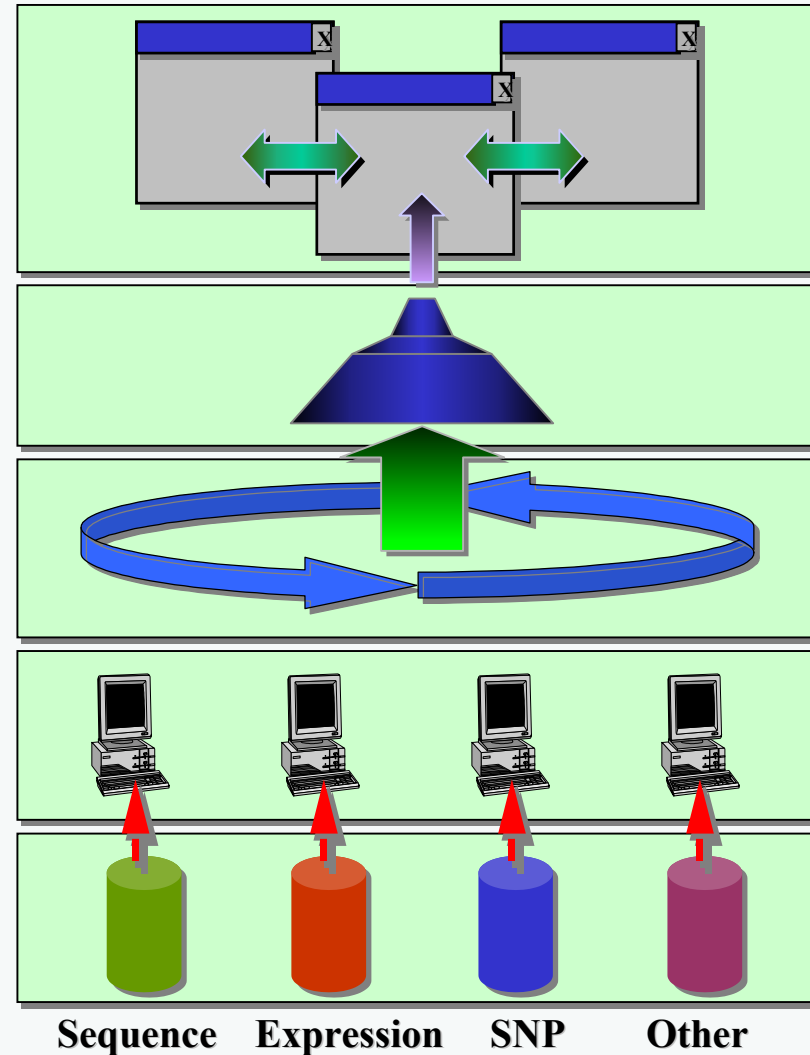# Data Mining + Domain Knowledge = Engineering Informatics

- **Informatics (정보 과학):** *the study of the structure, behavior, and interactions of natural and artificial computational systems*

- **Informatics = Information + Mathematics**

- **Application Areas**
  - ✓ **Bioinformatics (or Biomolecular Informatics)**
  - ✓ **Cheminformatics**
  - ✓ **Environmental Informatics**
  - ✓ **Medical Informatics**
  - ✓ **Neuro Informatics**
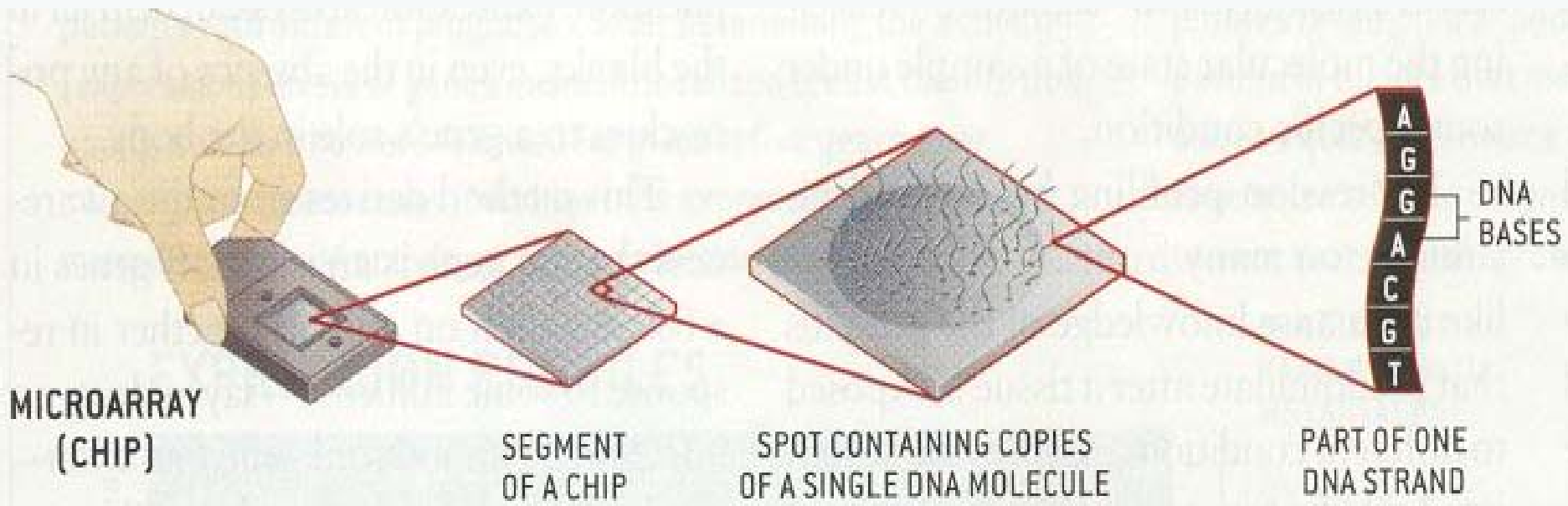  - ✓ **Process Informatics**
  - ✓ **Many More …**

# Bioinformatics

## Bioinformatics =
##     Biology + Informatics

- Artificial Intelligence
- Combinatorial Optimization
- Data Mining
- Digital Signal Processing
- Machine Learning
- Mathematical Modeling
- Multivariate Statistics
- Pattern Recognition
- System identification
- …

Sequence   Expression   SNP   Other

**POSTECH IPSE Lab.**

# DNA chip



MICROARRAY (CHIP) — SEGMENT OF A CHIP — SPOT CONTAINING COPIES OF A SINGLE DNA MOLECULE — PART OF ONE DNA STRAND — DNA BASES: A G G A C G T
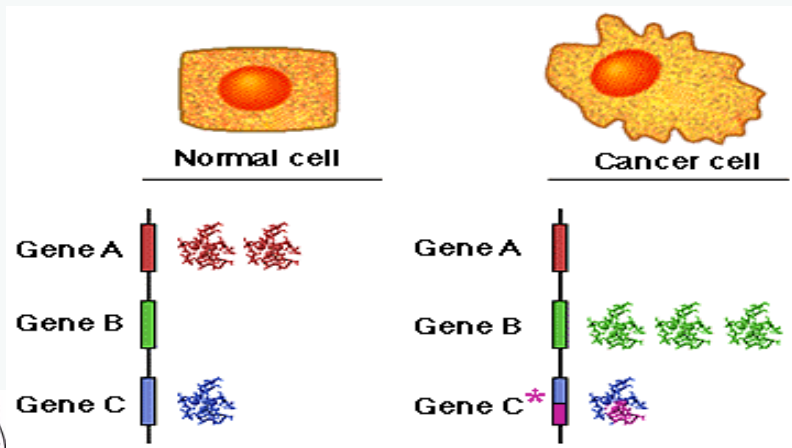
**POSTECH IPSE Lab.**

# The biological meaning of DNA chip

- Genome map is completed
  - ✓ Need to study functional genomic
- Know who, when, where, why, how much gene expressed
  - ✓ To classify different types of diseases (ex. Cancer types)
  - ✓ To understand the behavior of a biological system
  - ✓ To understand cell dynamics
- Can systematically disturb cell
- DNA chip experiment and data analysis are different matter
  - ✓ Methods of data analysis variant result of DNA chip experiment
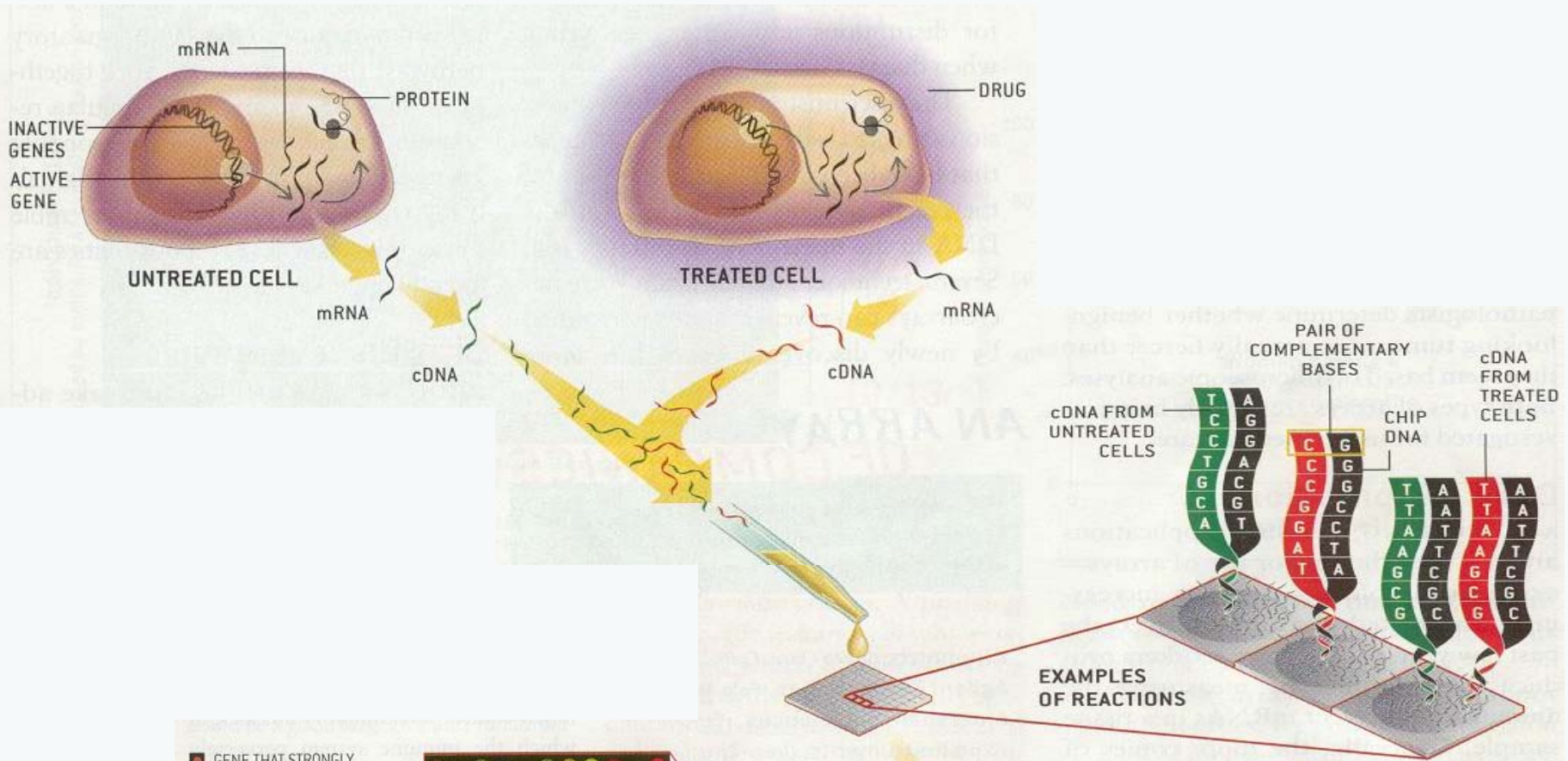  - ✓ Require suitable method of data analysis for DNA chip experiment objective

# Application of DNA chip

- **Analysis of gene expression and regulation**
  - ✓ Genetic network, pathway analysis, metabolic engineering
- **Disease diagnosis**
  - ✓ Molecular cancer classification, the discovery of disease subtype, The marker gene discovery
- **And many more…**

- **Cancer diagnosis**



Normal cell          Cancer cell
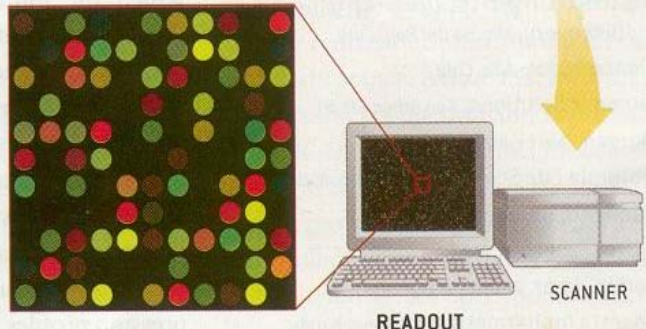
Gene A    Gene A
Gene B    Gene B
Gene C    Gene C*

✓ Because significantly different groups of genes are expressed by many type of cell, we can fingerprint characteristic cell

# cDNA chip: Lab Experiment



POSTECH IPSE Lab.

# cDNA chip: Data Analysis

Sample A    Sample B

Gene A

Gene B

Gene A

Gene B

**Statistical analysis (data mining)**

**Hierarchal clustering**

**K-means clustering**

**Self-organizing map**

**Neural network**

**Bayesian decision theory**

**Principal component analysis**

**And many more…**

Gene A

Gene B

Expression Level

Sample B    Sample A

Gene A

Gene B

**Molecular cancer classification**

Expression Level

Gene B

Sample B    Sample A

**Identification of the potential marker gene**

**POSTECH IPSE Lab.**

# cDNA chip: Procedure

**DNA chip experiment** → **For Functional study using data mining techniques**

**Biological validation**

**Informatics validation**

1. Expression profile data warehousing
2. Other database integration

**Data mining techniques**

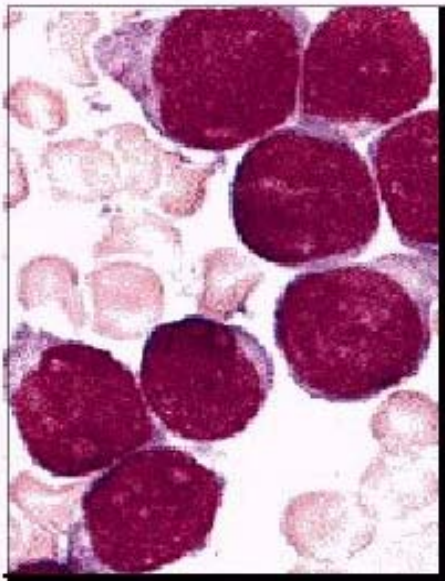**Inference new hypothesis for biological experiment**
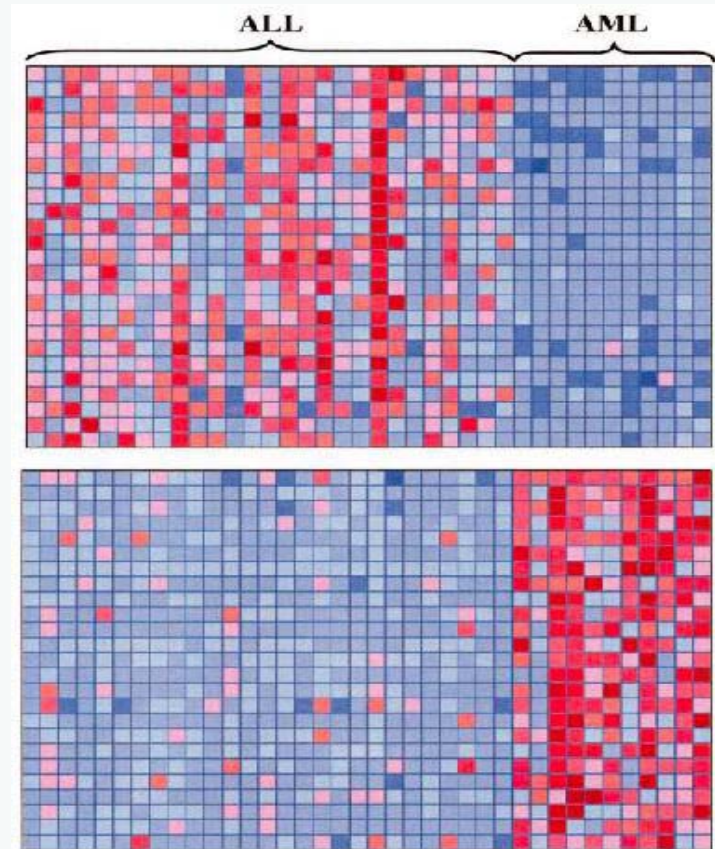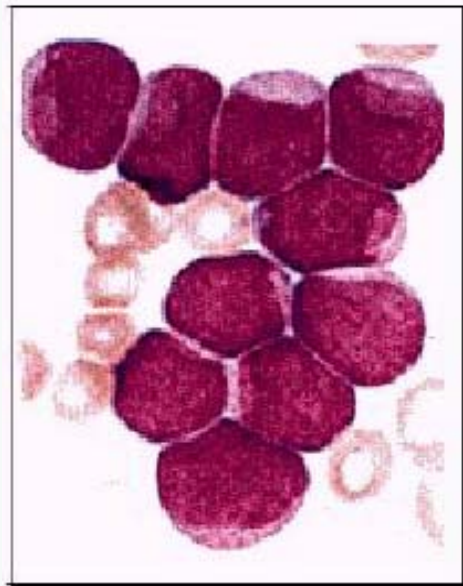
**POSTECH IPSE Lab.**

# Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

Moving from morphological to molecular classification

Acute lymphoblastic leukemia (ALL)
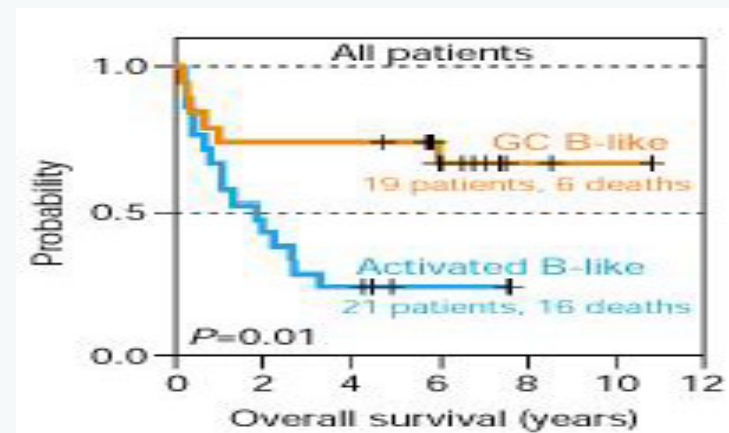
Acute myelogenous leukemia (AML)



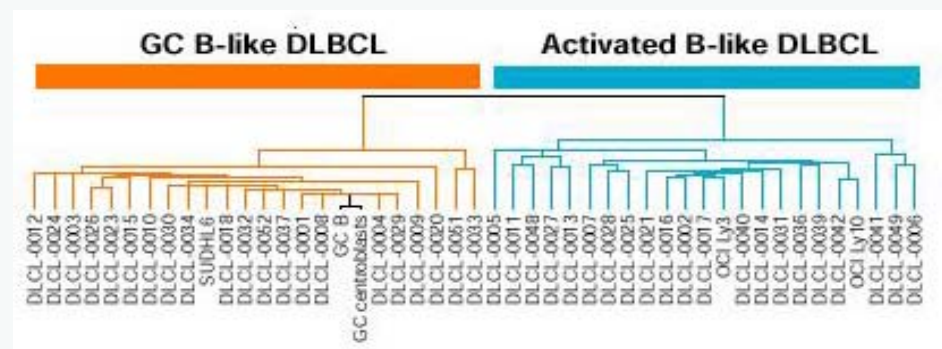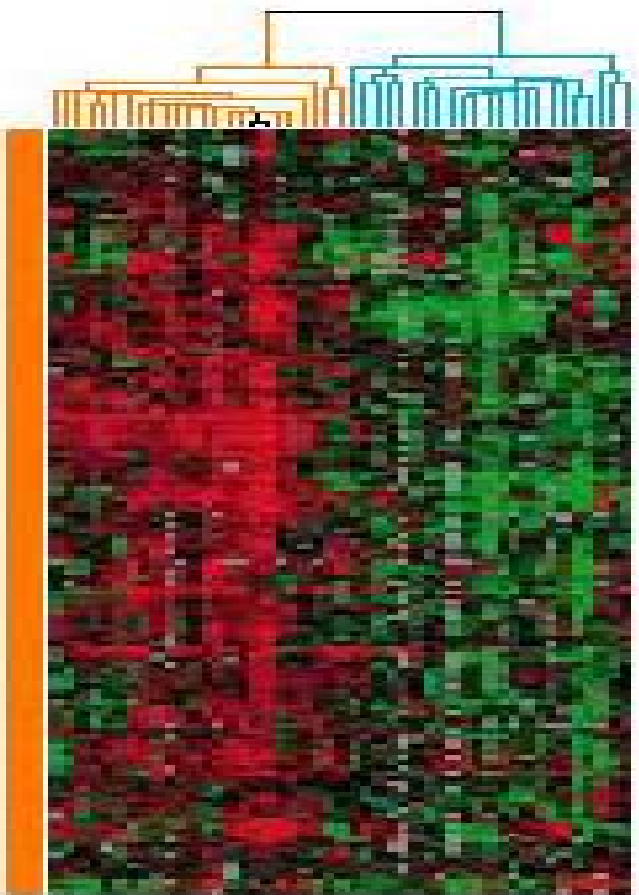**Golub.T.R.,et al. 1999. *Science,* 286, 531-537.**

# Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling

- There are new cancer class discovery, two molecularly distinct forms of B-cell lymphoma (DLBCL) that are composed of GC B-like and Activated B-like DLBCL



**Ash A. Alizadeh.,et al. 2000. *Nature*, 403, 503-511**
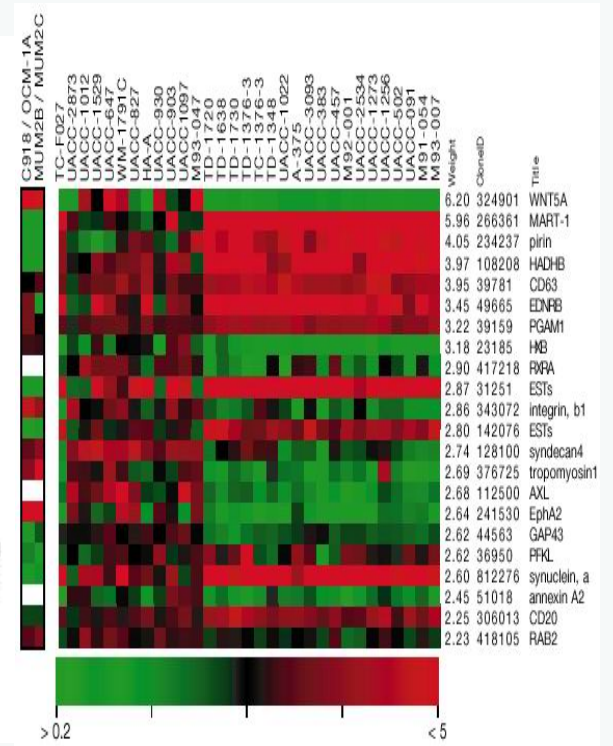
# Molecular Classification of Cutaneous Malignant Melanoma by Gene Expression Profile

## Limitation,

✓ Because weighting method based on univariative or bivariative statistical analysis, we can not capture correlated structure in the data

✓ When multi-class cancer classify, it is hard to know whether highly express or not



**M. Bittner.,et al. 2000. *Nature*, 406, 536-540.**

# Classification and Diagnostic Prediction of Cancers using Gene Expression Profiling and Artificial Neural Network

## Limitation,

✓ Because relevant gene extraction method based on univariative statistical analysis, we can not capture correlated structure in the data

# Gene-expression profiles in hereditary breast cancer

- Limitation,
  - ✓ Because relevant gene extraction method based on univariative statistical analysis, we can not capture correlated structure in the data
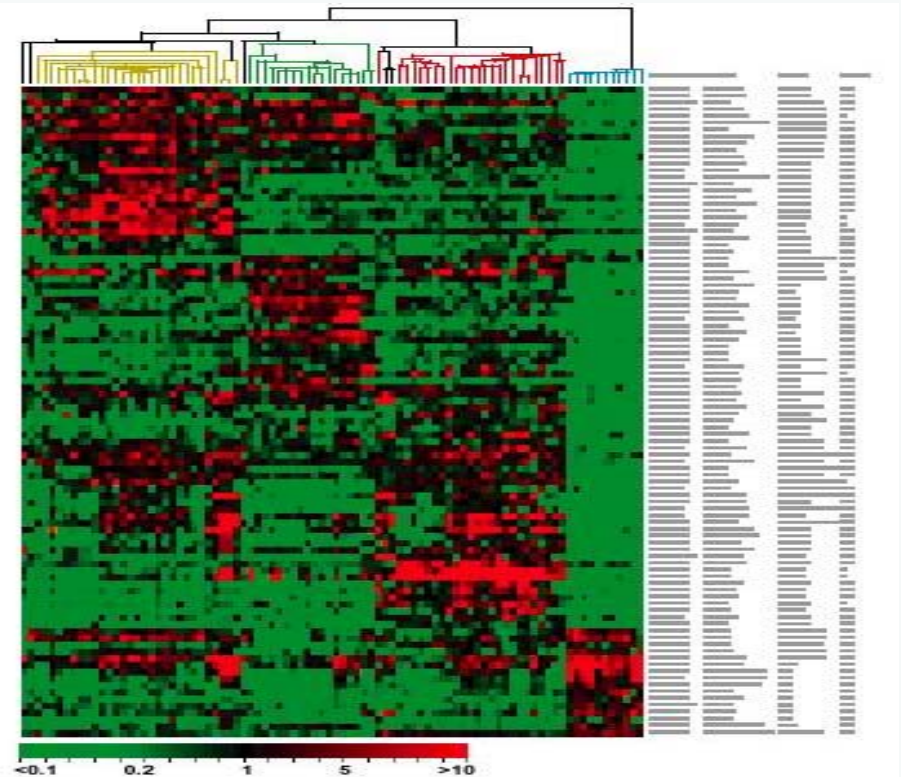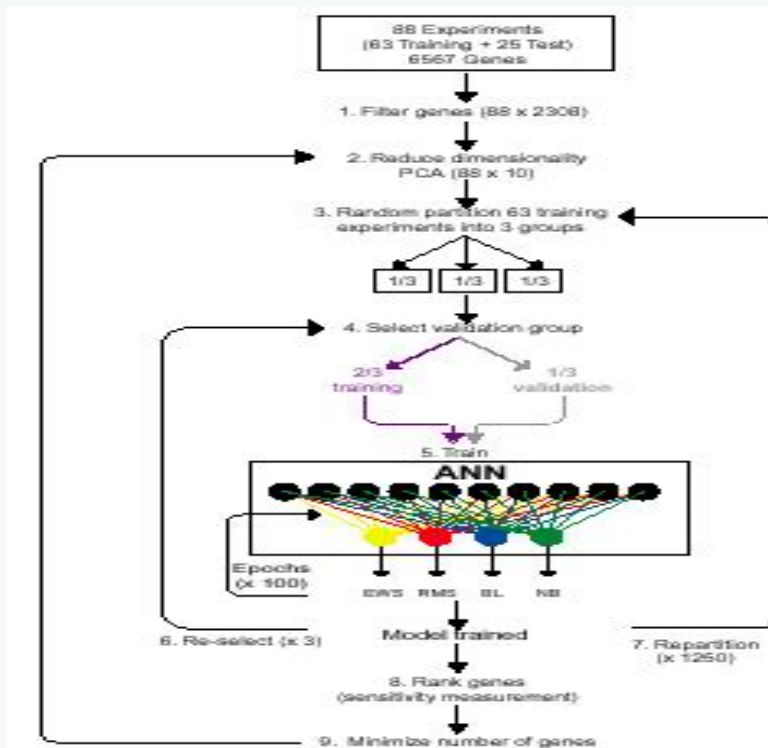  - ✓ When multi-class cancer classify, it is hard to know whether highly express or not



**Ingrid Hedenfalk ., et al. 2001. *N Eng j Med,* 344, 539-548.**

# Gene Selection for Sample Classification based on Gene Expression Data: Study of Sensitivity to Choice of Parameters of the GA/KNN Method

■ Limitation,

✓ It is difficult to determine parameter value

✓ When multi-class cancer classify, it is hard to know whether highly express or not

✓ Computing time take a long time



Leping. L .,et al. 2001. *bioinformatics*, 17, 1131-1142.

# Limitations and Improvements

## Limitations of Previous Approaches

- ✓ Small number of samples vs many variables
- ✓ Strong variable interaction
- ✓ Lack of interpretation based on biological meanings
- ✓ Limitation in the identification of marker genes due to the black box model
- ✓ Limitations due to univariate approaches
- ✓ Procedure of analysis are very complex and take a long time

## Improvements

- ✓ Overcome interaction of many variables
- ✓ Develop to a method to select potential marker genes
- ✓ Develop multivariate approach
- ✓ Develop simple and ease procedure of data analysis

# Proposed Procedure

**DNA Chip Data**

⬇

**PCA Analysis**

⬇

**Stepwise Discriminant Analysis**

⬇

**Bayesian Decision Theory (Classifier)**

⬇

**Contribution Analysis**

- **High dimensional data**
- **Highly correlated variables**

- **Data preprocessing**
  1. Dimension reduction
  2. Modeling of correlation structure
- **Feature selection**
  (Select highly discriminant PC)

- **Classification of tumor classes**

- **Select the potential marker genes**

# Major Steps

- **Stepwise Discriminant Analysis**
  - ✓ **Select subset where Wilks' lambda value is minimum**
  - ✓ **Maximize the discriminant power**

$$\Lambda = \frac{SS_w}{SS_t}$$

$SS_t$ : Class heterogeneity

$SS_w$ : Class homogeneity

- **Contribution Analysis**
  - ✓ **Discover potential marker genes to discriminate cancer classes**

$$C_j = \sum_{n=1}^{k} w_n \times p_{n,j} \times \left( t_{i,n} - t_{r,n} \right)$$

$C_j$ : the contribution of gene j

$p_{n,i}$ : the loading of the j-th gene on the n-th PC

$t_{i,n}$ : the average score of cancer class i

$t_{r,n}$ : the average score of reference cancer class

$w_n$ : the weight factor ( eigenvalue of n-th PC )

**POSTECH IPSE Lab.**

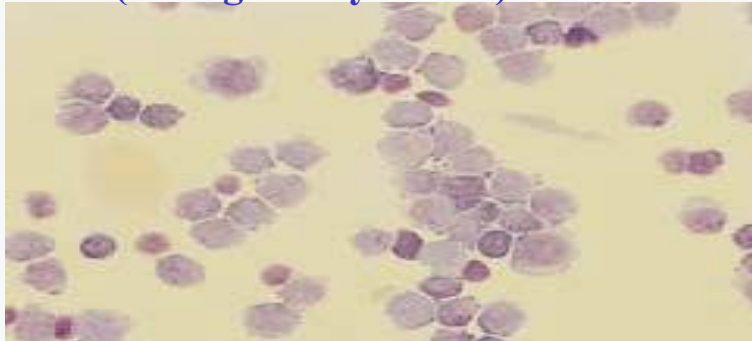# Case Study: Classification of Small Round Blue Cell Tumor

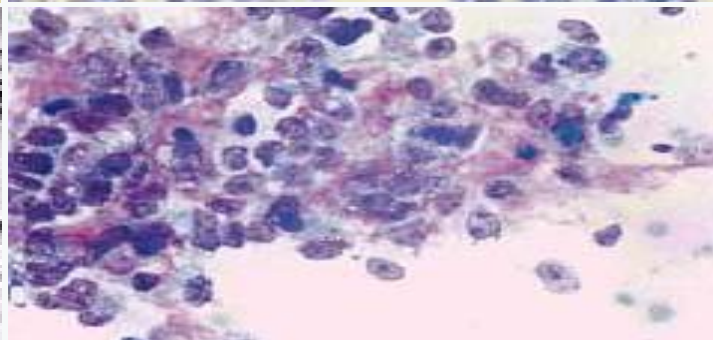- **Cancer DNA Chip data**
  - ✓ **Total samples : 88-by-2308 (samples-by-variables)**
  - ✓ **Training samples (63), Testing samples (20), Noise samples (5)**

- **Small round blue cell tumor**

**EWS (Ewing family tumor)**          **BL (non-Hodgkin lymphoma)**



**RMS (rhabdomyosarcoma)**          **NB (neuroblastoma)**

POSTECH IPSE Lab.

# Classification Results with SDA

**Classification power 100%**

| Bayesian decision theory | | Parametric method | | Nonparametric method | | |
|---|---|---|---|---|---|---|
| | | Linear discriminant function | Quadratic discriminant function | K- nearest neighbor | Kernel density | |
| | | | | | normal | biweight |
| Without SDA | Cross-validation of training set | 0.8359 | 0.75 | 0 | 0 | 0.45 |
| | Classification of the test set | 0.7833 | 0.75 | 0.7833 | 0.7833 | 0.3475 |
| Using SDA | Cross-validation of training set | 0.0238 | 0.5575 | 0 | 0 | 0.1051 |
| | Classification of the test set | 0 | 0.5417 | 0 | 0 | 0.1667 |

**Training sample: 63, test sample: 25**

# The number of identified potential marker genes

| Class | Number of genes identified using the proposed method | Number of genes identified (Khan et al., 2001) | Number of matched genes | Number of mismatched genes | Image ID number |
|---|---|---|---|---|---|
| EWS | 54 | 16 | 16 | 0 | |
| BL | 45 | 10 | 9 | 1 | 200814 |
| NB | 95 | 15 | 13 | 2 | 82225,  813266 |
| RMS | 68 | 20 | 14 | 6 | 788107,809901,122159 245330,246377,1409509 |
| Not BL | 61 | 12 | 8 | 4 | 45291, 204545 233721, 563673 |
| Not EWS | 12 | 1 | 1 | 0 | |
| Overlap | 24 | 0 | | | |
| Total | 311 | 74 | 61 | 13 | |

Khan et al. misjudgment 5 : image ID 82225, 813266, 233721, 245330, 122159
Redefine 2: image ID 45291, 563673
Overall trend agree 6 : image ID  204545, 788107, 1409509, 809901, 246377, 200814

**POSTECH IPSE Lab.**

# The expression profile of potential marker gene (1)

**Results are consistent with that of Khan et al., (2001)**

# The expression profile of potential marker gene (2)

**Not matched results**

**NB samples concurrently expressed in the BL,EWS, RMS**



**POSTECH IPSE Lab.**

## New discovered potential marker gene
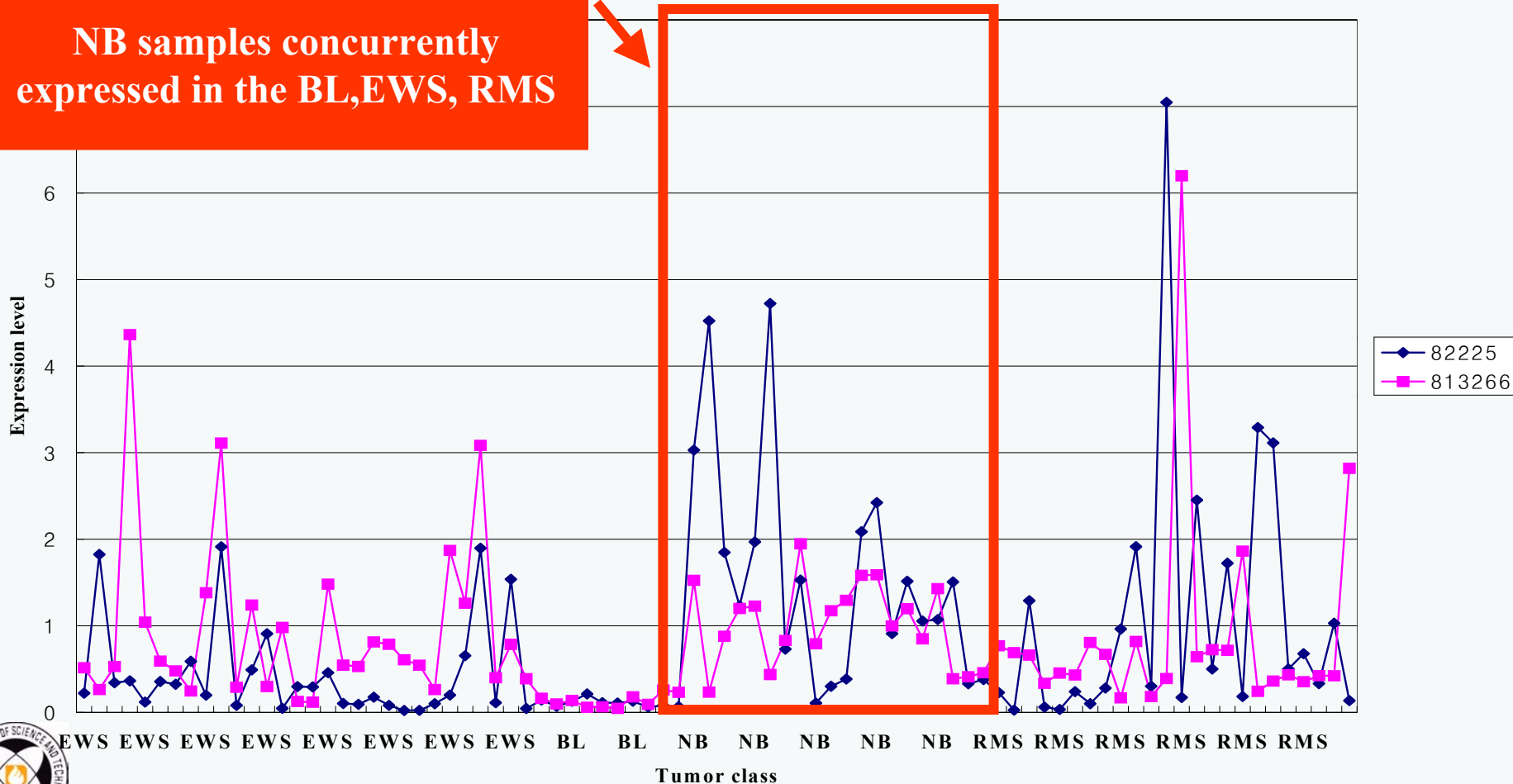


Genes expressed in RMS class

**POSTECH IPSE Lab.**

# Interpretation of the analysis results

■ **Hierarchal clustering based on 311 selected potential marker genes**

■ **Correct classification for each class**



Derived from same cell line (GICAN)

Derived from same cell line (ST486)

Noise sample

Noise sample of same cell (sk-muscle)

Samples

Genes

Highly expressed gene group of each tumor class

**POSTECH IPSE Lab.**

# Interpretation of the analysis results: Biological validation

## Marker genes for cancer classes

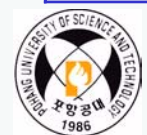| Gene Image ID | Cancer class | Biological gene function |
|---|---|---|
| 1435862 | EWS | antigen identified by monoclonal antibodies 12E7, F21 and O13 |
| 291756 | EWS | tubulin, beta, 5 |
| 43733 | EWS | glycogenin 2 |
| 52076 | EWS | olfactomedinrelated ER localized protein |
| 377731 | EWS | glutathione S-transferase M5 |
| 784224 | RMS | fibroblast growth factor receptor 4 |
| 470128 | RMS | Myosin IC |
| 296448 | RMS | insulin-like growth factor 2 (somatomedin A) |
| 207274 | RMS | Human DNA for insulin-like growth factor II (IGF-2); exon 7 and additional ORF |
| 461425 | RMS | Myogenesis |
| 377671 | RMS | integrin, alpha 7 |
| 823886 | RMS | Smooth muscle myosin heavy chain isoform SMemb [human, umbilical cord, fetal |

## Marker genes not matched for cancer classes

| Gene Image ID | Chip data Cancer class | Normal Cancer class | Gene Image ID | Chip data Cancer class | Normal Cancer class |
|---|---|---|---|---|---|
| 823886 | Not BL | RMS | 782488 | All class | Not NB |
| 897667 | EWS | RMS | 814773 | EWS | NB |
| 162208 | BL | RMS | 29054 | NB | RMS |
| 626502 | BL | RMS | 308231 | NB | RMS |
| 785793 | BL | RMS | 823886 | NB | RMS |
| 868304 | BL | RMS | 377048 | NB | RMS |
| 781018 | BL | RMS | | | |

# Contributions

- **Accurate multivariate classification method based on Bayesian method**

- **Potential marker gene selection method**

- **Simple and easy procedure for data analysis**

- **250 new candidate marker genes discovered**

- **new hypothesis testing based on the candidate marker genes for drug discovery or cancer research**

# Biotechnology meets data mining

- Time to dance!!!
- **Contacts** between the established 'data mining community' and 'bio/medical scientists' seem to be **rare**
- There will be more dances, and new biotechnology will be forthcoming as we learn the steps

**Coming dance !!!!**

POSTECH IPSE Lab.

# Questions ?



**Contacts and full paper request: sw74@postech.ac.kr**