# Theoretical limitations of massively parallel biology

## Genetic network analysis – gene and protein expression measurements

**Zoltan Szallasi**
**Children's Hospital**
**Informatics Program**
**Harvard Medical School**
**www.chip.org**

**Vipul Periwal**
Gene Network Sciences Inc.

**Mattias Wahde**
Chalmers University

**John Hertz** (Nordita)
**Greg Klus** (USUHS)

**How much information is needed to solve a given problem ?**

⟷

**How much information is (or will be) available ?**

**Conceptual limitations**

**Practical limitations**

- **Finding transcription factor binding sites based on primary sequence information**

- **SNP <> disease association**

# What are the problems we want to solve ?

So far the "DNA chip" revolution has been mainly technological:

The principles of measurements (e.g. complementary hybridization) have not changed.

It is not clear yet whether a conceptual revolution is approaching as well ?

potential breakthrough questions:

- can we perform efficient, non-obvious reverse engineering ?

- can we identify non-dominant cooperating factors ?

- can we predict truly new subclasses of tumors based on gene expression patterns ?

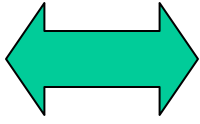- can we perform meaningful (non-obvious & predictive) forward modeling

1. **Reverse engineering time series measurements**

2. **Identification of novel classes or separators in gene expression matrices in a statistically significant manner**

3. **Potential use of artificial neural nets (machine learning) in the analysis of gene expression matrices.**

**Biological research has been based on the discovery of strong dominant factors.**
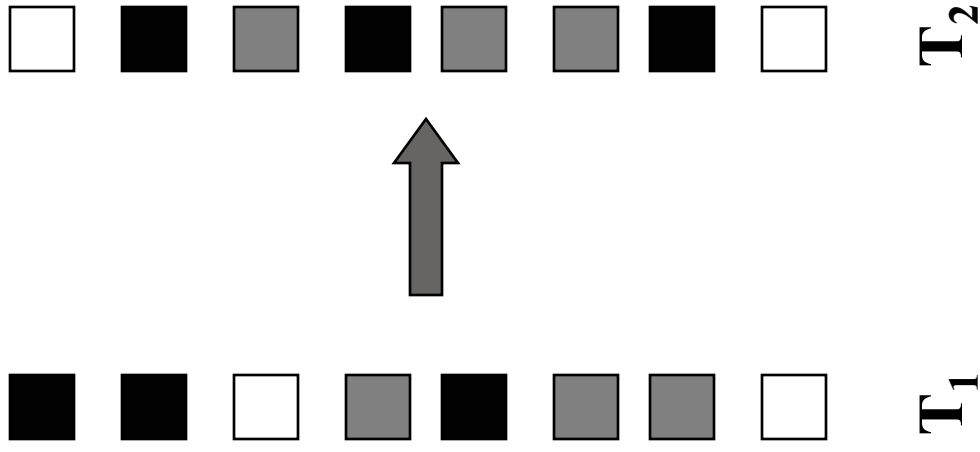
**More than methodological issue ?**

**Robust network based on stochastic processes**

$\longleftrightarrow$

**Strong dominant factors**

# The Principle of Reverse Engineering of Genetic Regulatory Networks from time series data:

**Determine a set of regulatory rules that can produce the gene expression pattern at $T_2$ given the gene expression pattern at the previous time point $T_1$**

$T_1$

$T_2$

## Continuous modeling:

$$x_i(t+1) = g\left(b_i + \sum_j w_{ij} x_j(t)\right)$$

(Mjolsness et al, 1991 - connectionist model;
Weaver et al., 1999, - weight matrix model;
D'Haeseleer et al., 1999, - linear model;
Wahde & Hertz, 1999 - coarse-grained reverse engineering)

at least as many time points as genes: T-1>N+2
(Independently regulated entities)

**For differential equations with r parameters 2r+1 experiments are enough for identification (E.D.Sontag, 2001)**

## How much information is needed for reverse engineering?

Boolean fully connected

$$2^N$$

Boolean, connectivity K

$$K\,2^K\,\log(N)$$

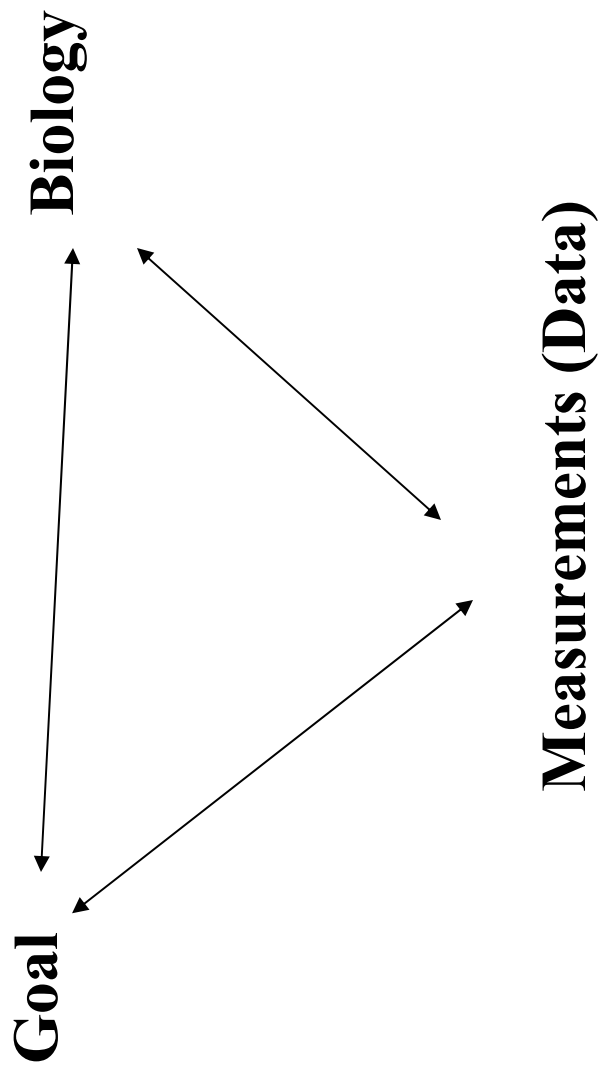Boolean, connectivity K, linearly separable rules

$$K\,\log(N/K)$$
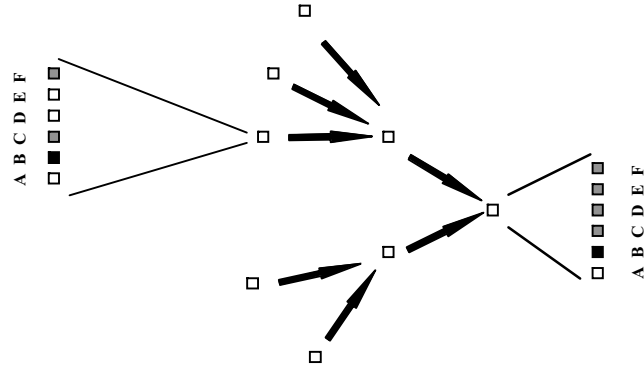
Pairwise correlation

$$\log(N)$$

N = number of genes
K = average regulatory input/gene

Biology

Measurements (Data)

Goal

# Biological factors that will influence our ability to perform successful reverse engineering.

(1) the stochastic nature of genetic networks ,

(2) the effective size of genetic networks ,

(3) the compartmentalization of genetic networks,

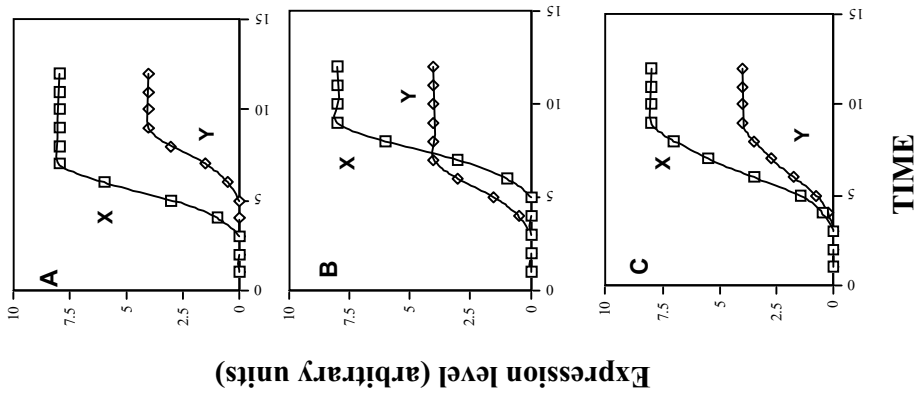# 1. The prevailing nature of the genetic network

## The effects of stochasticity:

1. It can conceal information (How much ?)

2. The lack of sharp switch on/off kinetics can reduce useful information of gene expression matrices.

(For practical purposes genetic networks might be considered as deterministic systems ?)

**Expression level (arbitrary units)**

**TIME**

## 2. The effective size of the genetic network:

### How large is our initial directed graph ?
### (It is probably not that large.)

We might have a relatively well defined deterministic cellular network with not more than 10 times the number of total genes.

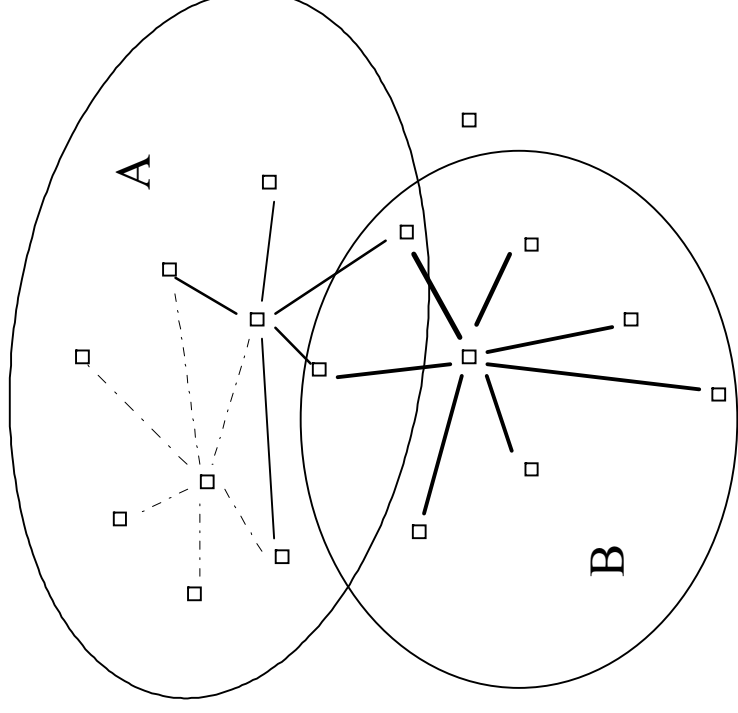$$N_{bic} < 10 \times N_{gene}$$

10,000-20,000 active genes per cell

Splice variants < > modules

# 3. The compartmentalization (modularity) of the genetic network:

The connectivity of the initial directed gene network graph

Low connectivity - better chance for computation.

**Genetic networks exhibit:**

**Scale–free properties (Barabasi et al.)**

**Modularity**

**Flatness**

## (Useful) Information content of measurements is influenced by the inherent nature of living systems

We can sample only a subspace of all gene expression patterns (gene expression space), because:

1. the system has to survive
(83% of the genes can be knocked out in S. cerevisiae)

2. Gene-expression matrices (i.e. experiments) are coupled
Cell cycle of yeast under different conditions

**Data:**

A reliable detection of 2-fold differences seems to be the practical limit of massively parallel quantitation. (estimate: optimistic and not cross-platform)

Population averaged measurements

GeneChip HG U95A

[Gene 5000:33700_at] [#196 C]

## The useful information content of time series measurements depend on:

1. Measurement error (conceptual and technical limitations, such as normalization)

2. Kinetics of gene expression level changes (lack of sharp switch on/off kinetics – stochasticity ?)

3. Number of genes changing their expression level.

4. The time frame of the experiment.

Level of gene expression

Measurements with error bars

Time window

Time

A rational experiment will sample gene-expression according to a time-series in which each consecutive time point is expected to produce at least as large expression level difference as the error of measurement: approximately 5 min intervals in yeast, 15-30 min intervals in mammalian cells.

**P = K log(N/K)  (John Hertz, Nordita)**

**P : gene expression states**
**N: size of network**
**K: average number of regulatory interactions**

Applying all this to cell cycle dependent gene expression measurements by cDNA microarray one can obtain 1-2 orders of magnitude less information than expected in an ideal situation. (Szallasi, 1998)

Can we identify non-dominant cooperating factors ?

Can we predict truly new subclasses of tumors based on gene expression patterns ?

How much data is needed ?

How much data will be available ?

Samples

M93-007
M91-054
UACC-091
UACC-502
UACC-1256
UACC-1273
UACC-2534
M92-001
UACC-457
UACC-383
UACC-3093
A-375
UACC-1022
TD-1384
TC-1376-3
TD-1376-3
TD-1730
TD-1638
TD-1720
M93-47
UACC-1097
UACC-903
UACC-930
H-A
UACC-827
WM1791-C
UACC-647
UACC-1529
UACC-1012
UACC-2873
TC-F027

141562 724112 768357 79629 490306 842906 826173 33051 344282 357970 140966 245489 141171 35236 809727 137531 203240 109265 365060 563873 840942 627541 814306 714426 120544 138021 297439 134829 151418 296754 281843 358531 809848 502891 293292 49591 242037 282310

## Analysis of massively parallel data sets

**Unsupervised** – avoiding artifacts in random data sets

avoiding artifacts in data sets retaining the internal data structure

**Supervised**

# INFORMATION REQUIREMENT

# Consistently mis-regulated genes in random matrices

"E" different samples

"N"-gene microarray

$M_i$ genes mis-regulated in the "i"-th sample,

K consistently mis-regulated across all E samples.

What is the probability that (at least) K genes were mis-regulated by chance ?

$$P(E, k \geq K) = 1 - \sum_{i=0}^{K-1} P(E,k)$$

**Where P(E,k) is the probability that exactly k genes are consistently mis-regulated**

$$P(E,k) = \sum_{j=k}^{M} \frac{\binom{j}{k} * \binom{N-j}{M-k} * P(E-1,j)}{\binom{N}{M}}$$

$$P(2,k) = \frac{\binom{M}{k} * \binom{N-M}{M-k}}{\binom{N}{M}}$$

**If N>>M, then**

$$P(E, k \geq K) \approx \frac{\binom{N}{K} * \binom{N-K}{M-K}^E}{\binom{N}{M}^E}$$

or

$$P(E, k \geq K) \approx \binom{N}{k} * (q^E)^k * (1 - q^E)^{N-k}$$

**For a K gene separator:**

$$n_K = \binom{N}{K} * \left(1 - q^K\right)^E$$

| N | M | E | K | $n_K$ simulated | $n_K$ calculated |
|---|---|---|---|---|---|
| 500 | 100 | 4 | 3 | 1172455±123637 | 1174430 |
| 500 | 100 | 8 | 3 | 69630 ± 17487 | 66605 |
| 300 | 50 | 15 | 3 | 760 ± 579 | 785 |
| 200 | 40 | 20 | 4 | 2032 ± 1639 | 1713 |

## how many cell lines do we need in order to avoid accidental separators ?

for N=10000    M=1000        for p<0.001
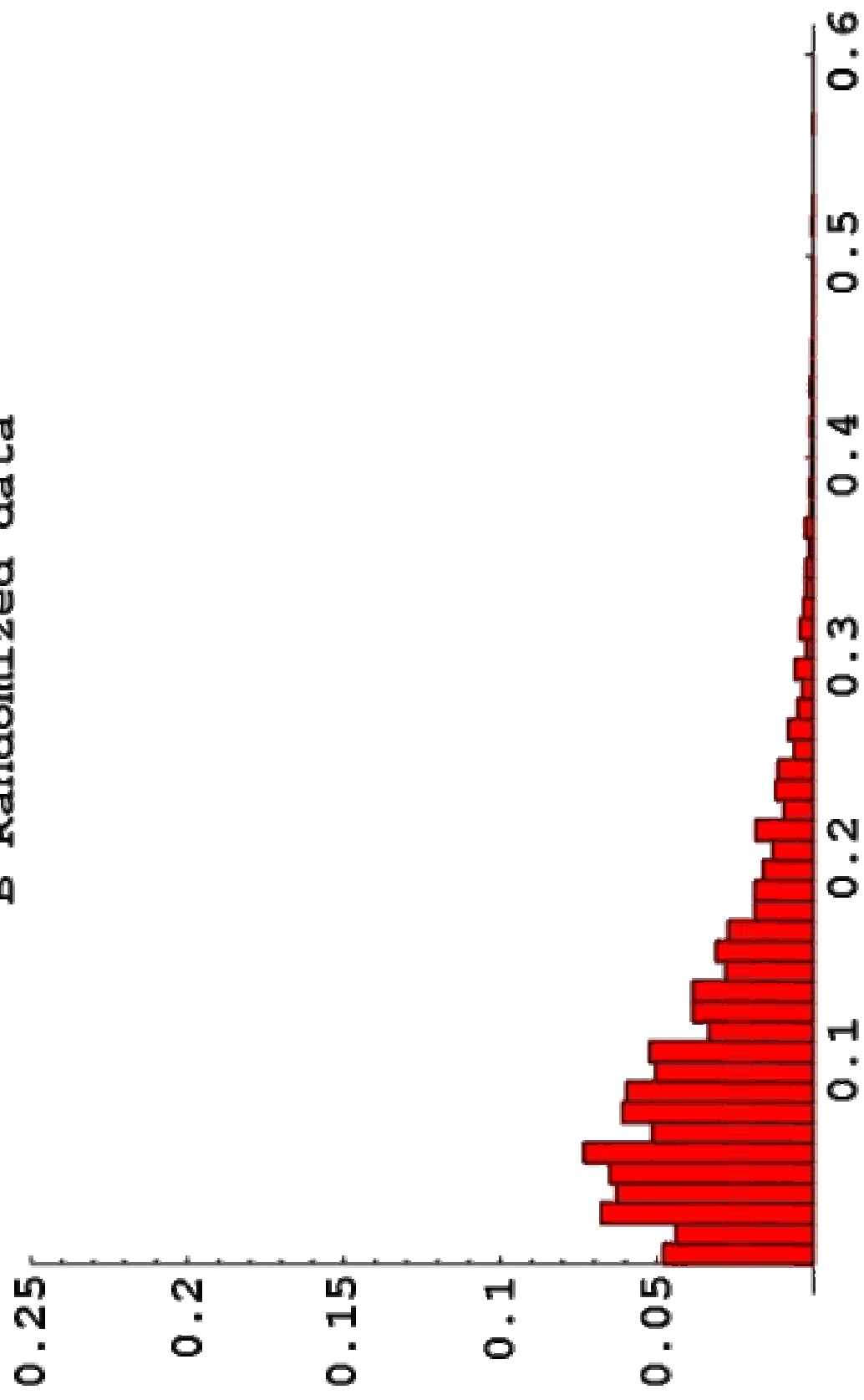
K=1    E=7

### Higher order separator

K=2    E=15

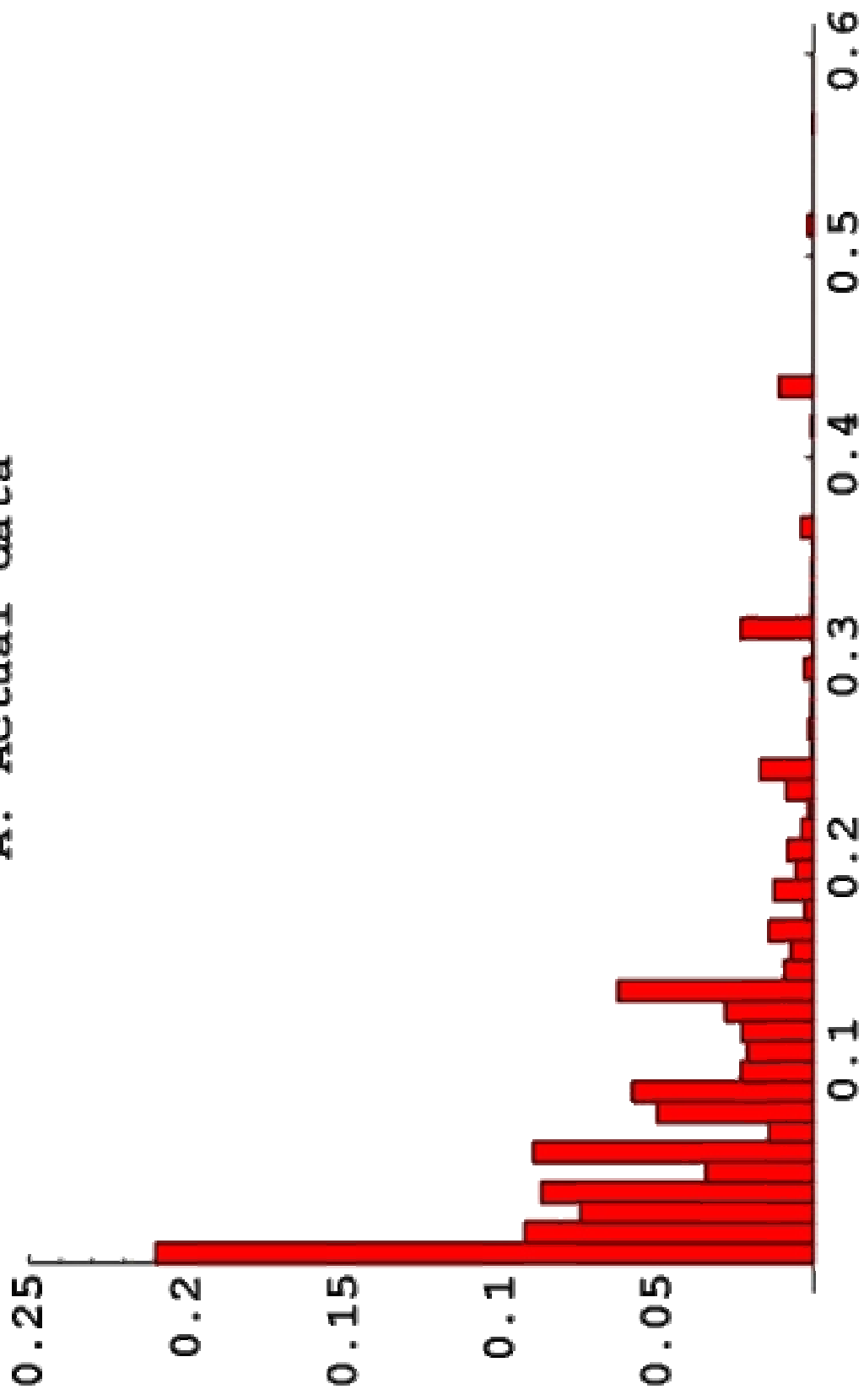K=3    E=25

K=4    E=38

K=5    E=54

K=6    E=73

B Randomized data

# Genes are not independently regulated

## A. Actual data

**Generative models (gene expression operator) will simulate realistic looking gene expression matrices ?**

- the number of genes that can be mis-regulated
- the independence of gene mis-regulation.

| | $N_1$ | $N_2$ | $N_3$ .... | $N_i$ | $T_1$ | $T_2$ | $T_3$ .......... | $T_i$ |
|---|---|---|---|---|---|---|---|---|
| gene$_1$ | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 |
| gene$_2$ | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| gene$_3$ | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| gene$_4$ | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 |

# Algorithm to extract Boolean separators from a gene expression matrix.

**U. Alon data set (colon tumors) : N=2000, $M_{average}$ =180 K=2**

| E | Alon data | calc. | Num. sim. |
|---|---|---|---|
| 10 | 708 | 131 | 130 |
| 11 | 120 | ~1 | 1 |
| 12 | 45 | $8.6 \times 10^{-3}$ | $8.6 \times 10^{-3}$ |
| 13 | 3 | $7.0 \times 10^{-5}$ | - |
| 14 | 3 | $5.6 \times 10^{-7}$ | - |
| 15 | 1 | $4.6 \times 10^{-9}$ | - |
| 16 | 1 | $3.7 \times 10^{-11}$ | - |

**Generative model: 4+/-2 separators**

Unclustered **A** · Clustered **B** · Random Unclustered **C** · Random Clustered **D**

features, while preserving step-like changes in intensity. The features were arranged in the order they appear in the EST sequence, the PM-MM intensities in a moving window of five features were sorted, and the filtered intensity was given by the mean of the middle three sorted intensities. The total intensity

coefficient of x a
The binary tre
to reorganize the
this end, we incl
in a determinist

Fi
norm
are t
the s
the h
so it
is 1.
scale
62 ti
the c
on t
(C)
data
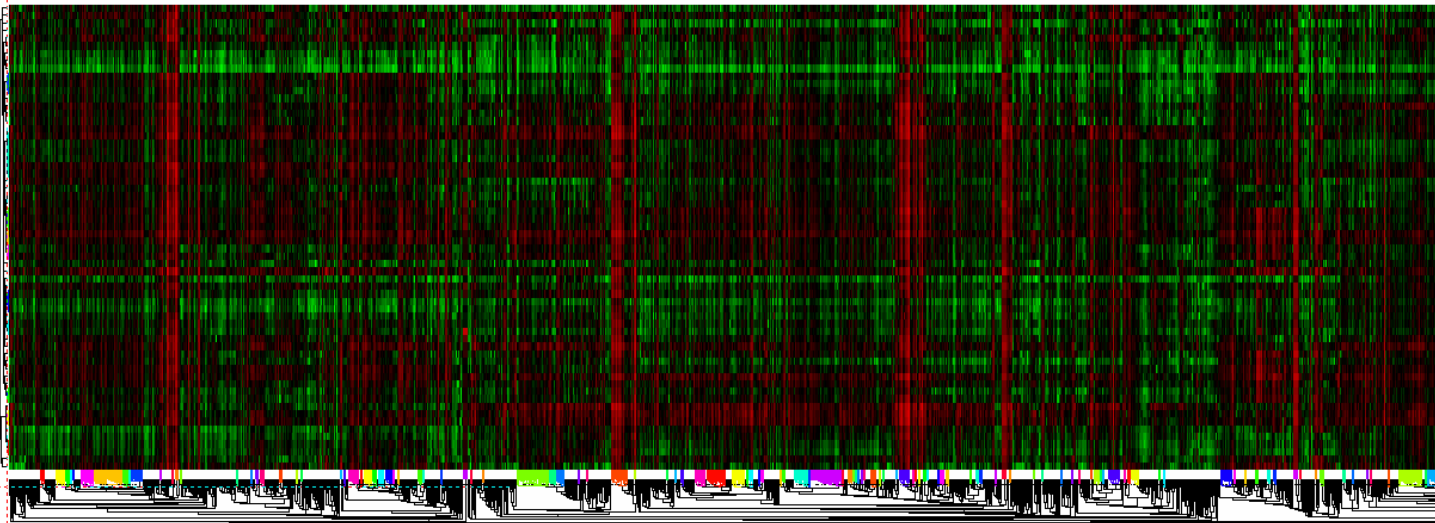in t
rand
algo
are
colo

# Pearson-disproportion of an array:

$$PE(\mathbf{y}) = \sum_{i=1}^{r} \sum_{j=1}^{c} \frac{(y_{ij} - \frac{m_i n_j}{N})^2}{\frac{m_i n_j}{N}}$$

$y_{ij}$ = gene expression level in the $i$th row and $j$th column

$$m_i = \sum_{j=1}^{c} y_{ij} \qquad n_j = \sum_{i=1}^{r} y_{ij}$$

$$N = \sum_{i} m_i = \sum_{j} n_j$$

Random matrices with the same intensity distribution and same (or larger) disproportion measure as the original matrix (Monte Carlo simulations)

**Generative models (random matrices retaining internal data structure) will help to determine the required sample number for statistically meaningful identification of classes and separators.**

# Machine learning – Artificial Neural Nets in the analysis Cancer associated gene expression matrices

**A**

**B**

ER Cluster

A

**58 Experiments**
**(47 Training + 11 Test)**
**6728 Genes**

1. Filter genes (58 x 3389)

2. Reduce dimensionality PCA (58 x 10)

3. Random partition 47 training experiments into 3 groups

1/3  1/3  1/3

4. Select validation group

2/3 training

1/3 validation

5. Train

**ANN**

Epochs (x 100)

0: ER-    1: ER+

Model trained

6. Re-select (x 3)

7. Repartition (x 200)

8. Rank genes (sensitivity measurement)
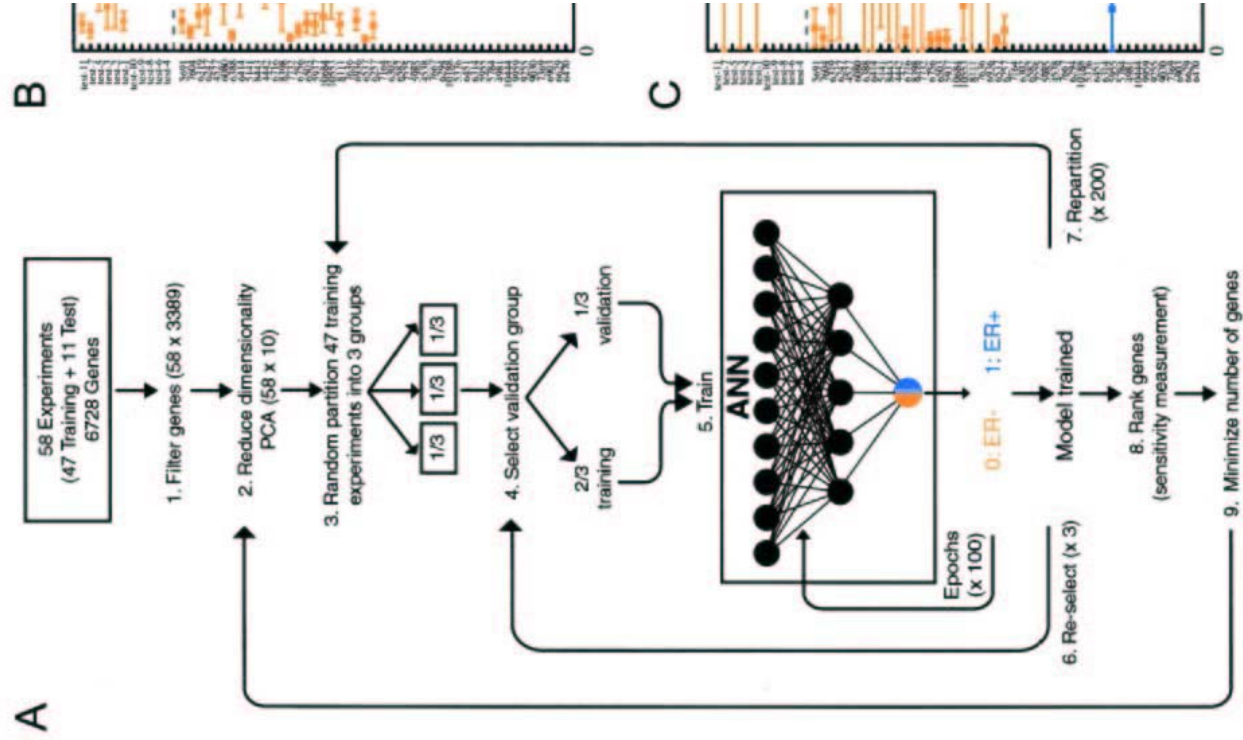
9. Minimize number of genes

B

C

Fig. 1. Classification of ER+ and ER− tumors using ANNs and gene expression patterns. A expression and spot area reduced the number of genes to 3389 (1). PCA further reduced the dime into three groups (3). Two of these groups were used for training and one for validation (4) using process was repeated so that all three groups were used for validation (6). The random partition
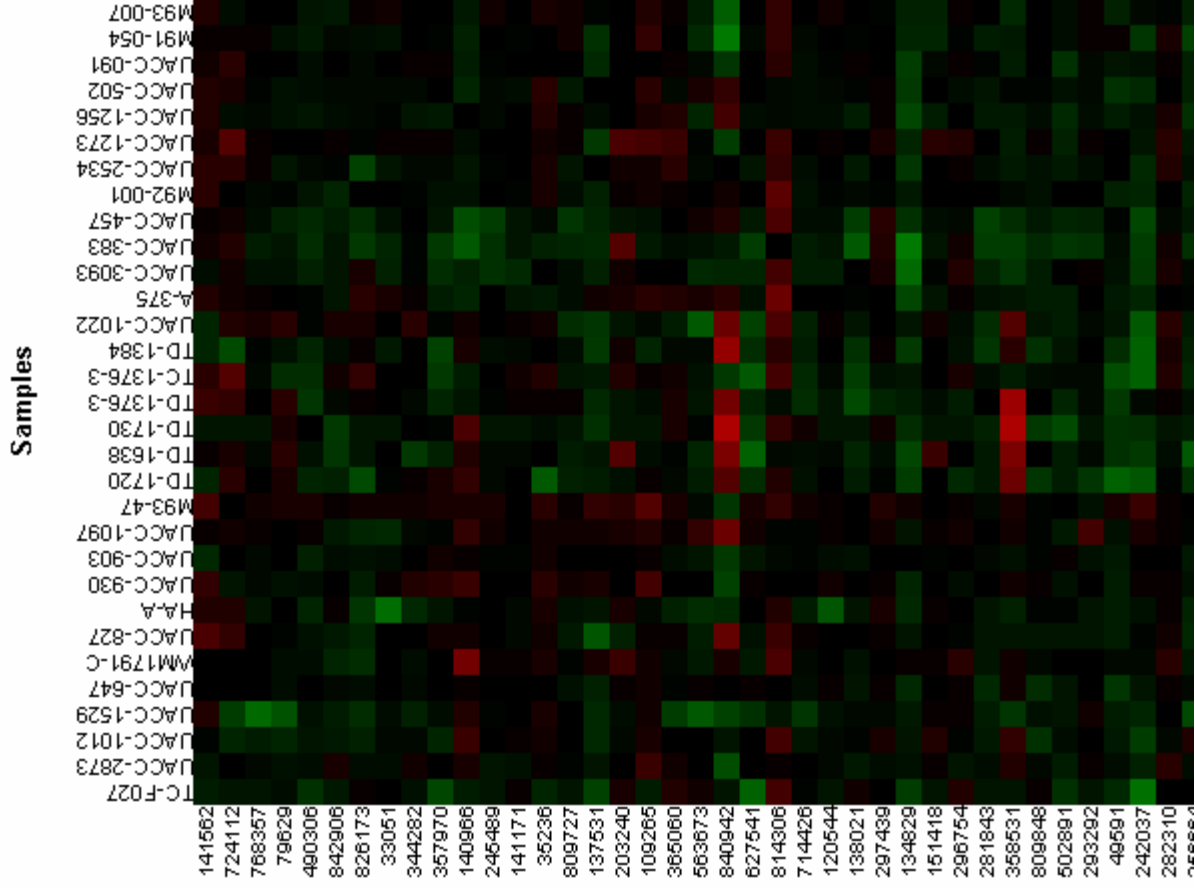
P. Meltzer,
J. Trent
M. Bittner

ANN (artificial neural nets) work well when a large number of samples is available relative to the number of variables

(e.g. for the pattern recognition of hand written digits one can create a huge number of sufficiently different samples).

In biology there might be two limitations:

1. the number of samples might be quite limited, at least relative to the complexity of the problems (The cell has to survive)

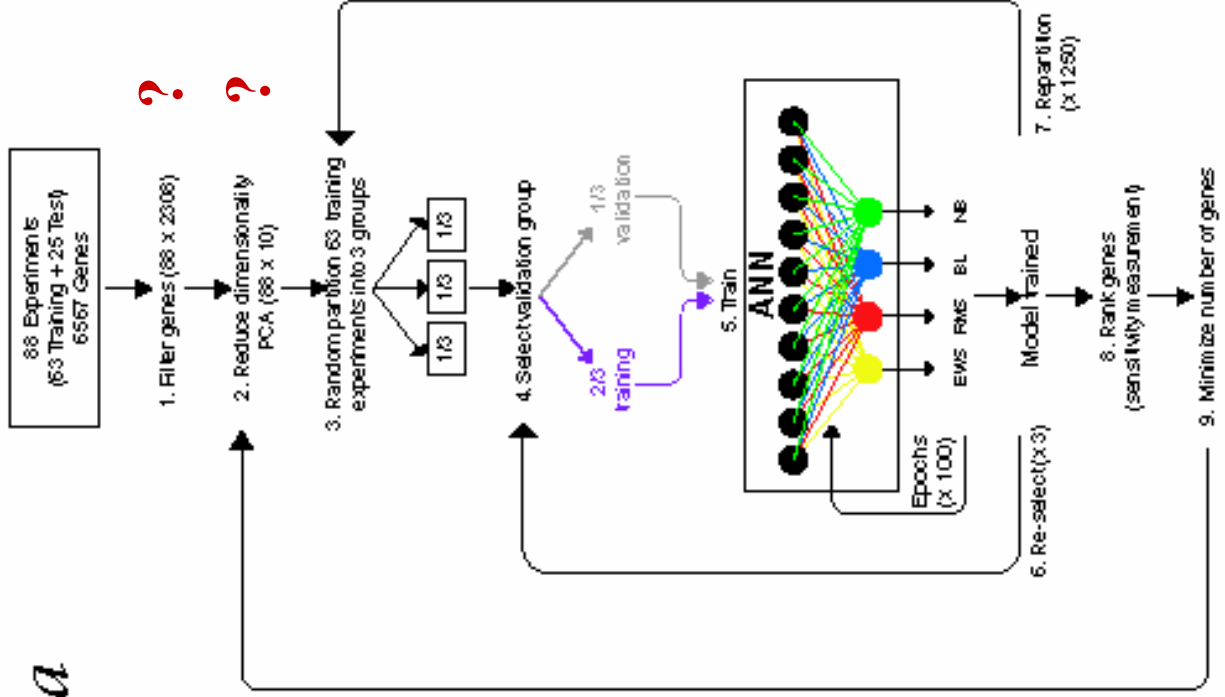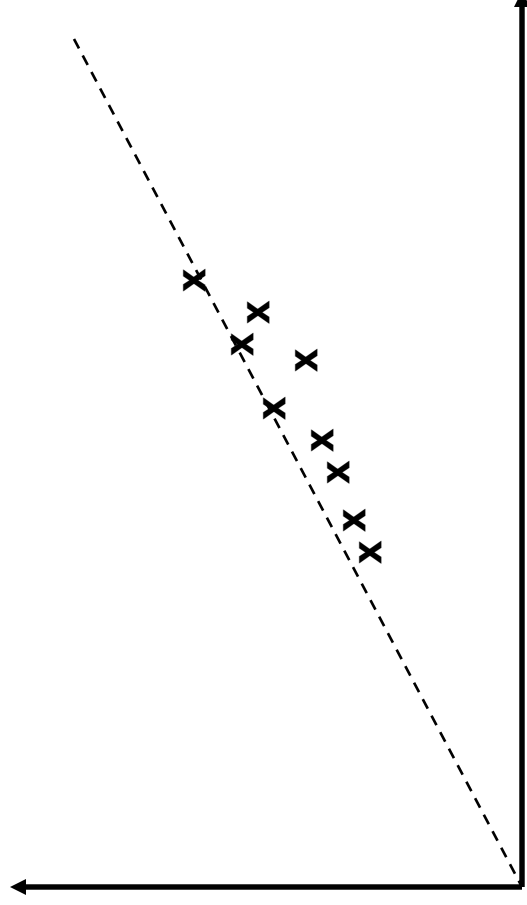2. There might be a practical limit to collecting certain types of samples

Summed square error

Number of misclassified samples

*a*

*b*

*c*

88 Experiments
(63 Training + 25 Test)
6567 Genes

1. Filter genes (88 x 2308)

2. Reduce dimensionality
PCA (88 x 10)

3. Random partition 63 training
experiments into 3 groups

1/3   1/3   1/3

4. Select validation group

1/3
validation

2/3
training

5. Train

ANN

Epochs
(x 100)

EWS   RMS   BL   NB

Model Trained

8. Rank genes
(sensitivity measurement)

9. Minimize number of genes

6. Re-select (x3)

7. Repartition
(x 1250)

# Reducing dimensionality
## Principal component analysis
### retain variance

# The risk of reducing dimensionality by PCA

A

B

ER Cluster

**(Rosetta)
83% accuracy
with 70 genes**

**Simple genetic
algorithm by us:
93% with 3 genes**

Figure 2 Supervised classification on prognosis signatures. **a**, Use of prognostic reporter genes ... prognostic classifier with optimal accuracy; dashed line, with optimized sensitivity. Above the dashed line, patients have a good prognosis signature, below the dashed line the ...