# Gene network inference



**Alberto de la Fuente**
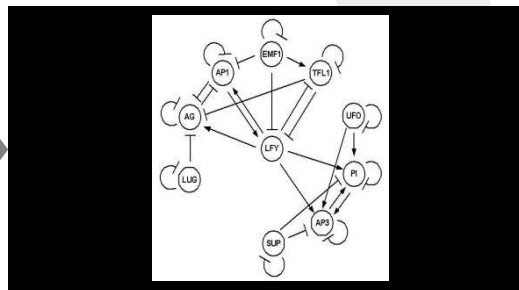
alf@crs4.it

**CRS4 Bioinformatica**

**Andrea Pinna**

**Nicola Soranzo**

**Vincenzo de Leo**
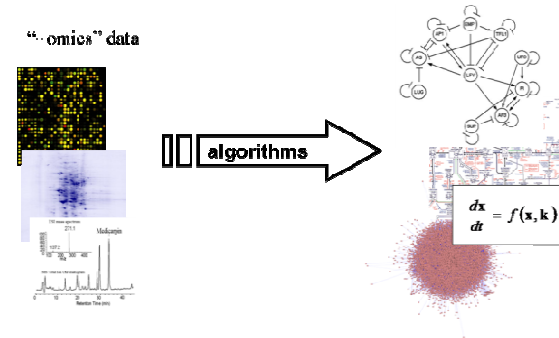
**GENOTYPE
+
ENVIROMENT** → **PHENOTYPE**
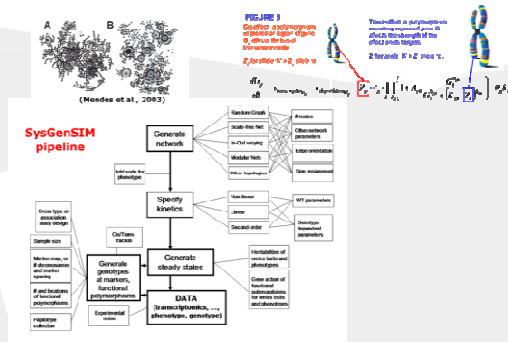
- Introduction to Gene networks

- Gene network inference

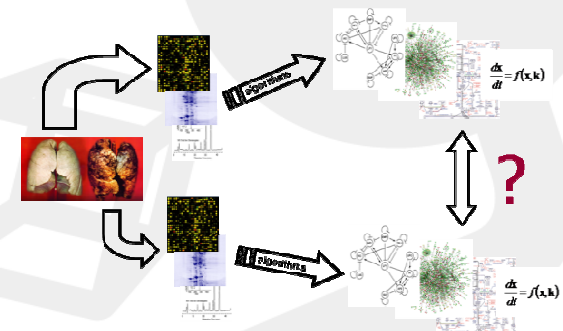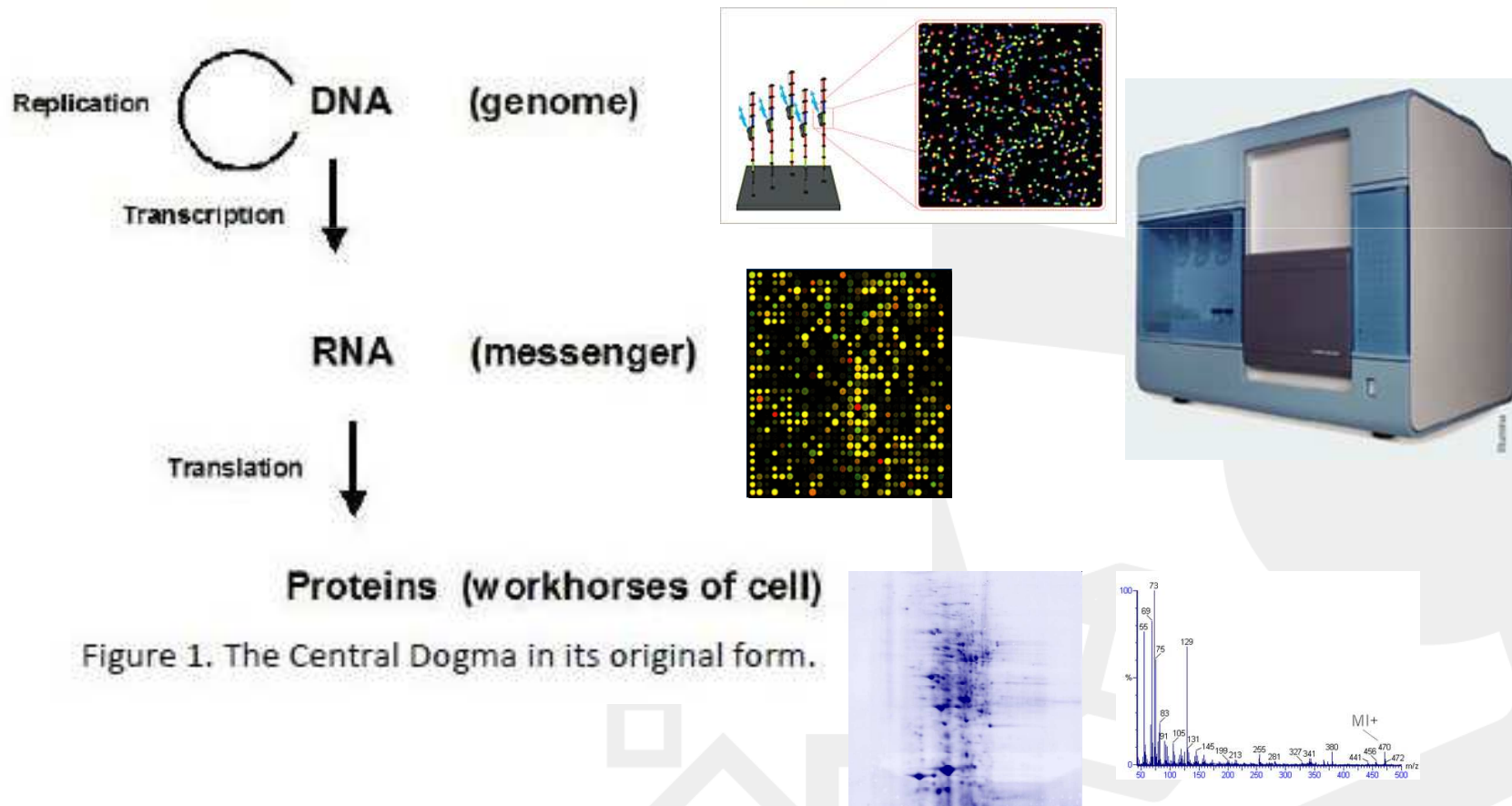- Evaluation of gene network inference algorithms

- Differential networking in disease

- **Introduction to Gene networks**

- Gene network inference

- Evaluation of gene network inference algorithms
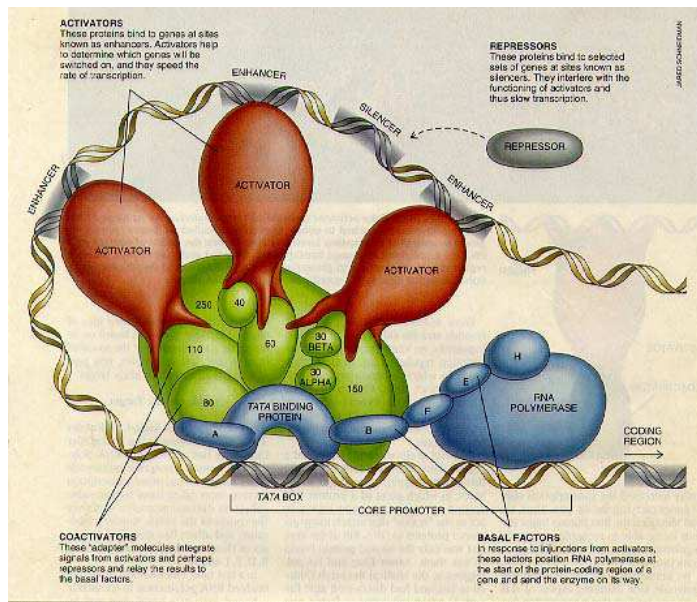
- Differential networking in disease

Figure 1. The Central Dogma in its original form.

```
GCCACATGTAGATAATTGAAACTGGATCCTCATCCCTCGCCTTGTACAAAAATCAACTCCAGATGGATCTAAG
ATTTAAATCTAACACCTGAAACCATAAAAATTCTAGGAGATAACACTGGCAAAGCTATTCTAGACATTGGCTT
AGGCAAAGAGTTCGTGACCAAGAACCCAAAAGCAAATGCAACAAAAACAAAAATAAATAGGTGGGACCTGATT
AAACTGAAAAGCCTCTGCACAGCAAAAGAAATAATCAGCAGAGTAAACAGACAACCCACAGAATGAGAGAAAA
TATTTGCAAACCATGCATCTGATGACAAAGGACTAATATCCAGAATCTACAAGGAACTCAAACAAATCAGCAA
GAAAAAAATAACCCCATCAAAAAGTGGGCAAAGGAATGAATAGACAATTCTCAAATATACAAATGGCCAATA
AACATACGAAAAACTGTTCAACATCACTAATTATCAGGGAAATGCAAATTAAAACCACAATGAGATGCCACCT
TACTCCTGCAAGAATGGCCATAATAAAAAAAATCAAAAAGAATAAATGTTGGTGTGAATGTGGTGAAAAGA
GAACACTTTGACACTGCTGGTGGGAATGGAAACTAGTACAACCACTGTGGAAAACAGTACCGAGATTTCTTAA
AGAACTACAAGTAGAACTACCATTTGATCCAGCAATCCCACTACTGGGTATCTACCCAGAGGAAAAGAAGTCA
TTATTTGAAAAGACACTTGTACATACATGTTTATAGCAGCACAATTTGCAATTGCAAAGATATGGAACCAGT
CTAAATGCCCATCAACCAACAAATGGATAAAGAAAATATGGTATATATACACCATGGAACACTACTCAGCCAT
AAAAAGGAACAAAATAATGGCAACTCACAGATGGAGTTGGAGACCACTATTCTAAGTGAAATAACTCAGGAAT
GGAAACCAAATATTGTATGTTCTCACTTATAAGTGGGAGCTAAGCTATGAGGACAAAAGGCATAAGAATTAT
ACTATGGACTTTGGGGACTCGGGGGAAAGGGTGGGAGGGGGATGAGGGACAAAAGACTACACATTGGGTGCAG
TGTACACTGCTGAGGTGATGGGTGCACCAAAATCTCAGAAATTACCACTAAAGAACTTATCCATGTAACTAAA
AACCACCTCTACCCAAATAATTTTGAAATAAAAAATAAAATATTTTAAAAGAACTCTTTAAAATAAATAAT
GAAAAGCACCAACAGACTTATGAACAGGCAATAGAAAAAATGAGAAATAGAAAGGAATACAAATAAAAGTACA
GAAAAAAAATATGGCAAGTTATTCAACCAAACTGGTAATTTGAAATCCAGATTGAAATAATGCAAAAAAAAGG
CAATTTCTGGCACCATGGCAGACCAGGTACCTGGATGATCTGTTGCTGAAACAACTGAAAATGCTGGTTAAA
ATATATTAACACATTCTTGAATACAGTCATGGCCAAAGGAAGTCACATGACTAAGCCCACAGTCAAGGAGTGA
GAAAGTATTCTCTACCTACCATGAGGCCAGGGCAAGGGTGTGCACTTTTTTTTTTCTTCTGTTCATTGAATAC
AGTCACTGTGTATTTTACATACTTTCATTTAGTCTTATGACAATCCTATGAACAAGTACTTTTAAAAAAATT
GAGATAACAGTTGCATACCGTGAAATTCATCCATTTAAAGTGAGCAATTCACAGGTGCAGCTAGCTCAGTCAG
CAGAGCATAAGACTCTTAAAGTGAACAATTCAGTGCTTTTTAGTATATTCACAGAGTTGTGCAACCATCACCA
CTATCTAATTGGTCTTAGTCTGTTTGGGCTGCCATAACAAAATACCACAAACTGGATAGCTCATAAACAACAG
GCATTTATTGCTCACAGTTCTAGAGGCTGGAAGTGCAAGATTAAGATGCCAGCAGATTCTGTGTCTGCTGAGG
GCCTGTTCCTCATAGAAGGTGCCCTCTTGCTGAATTCTCACATGGTGGAAGGGGGAAAACAAGCTTGCATTGC
```
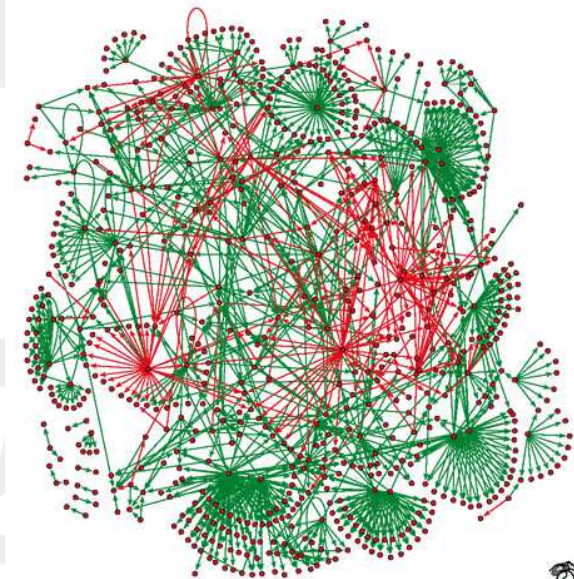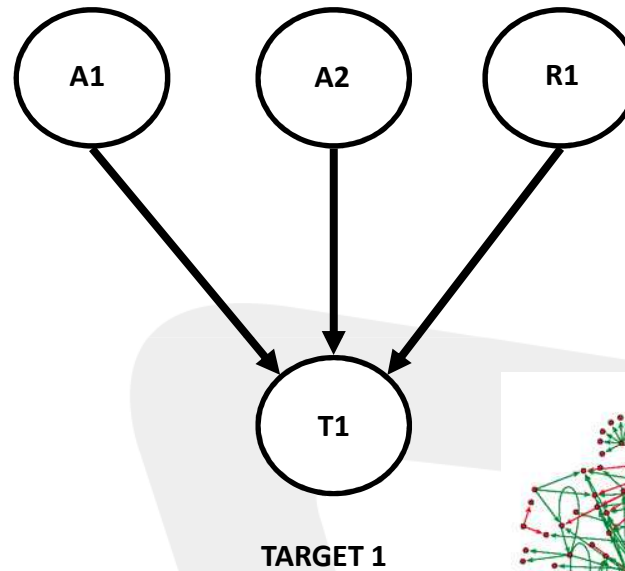
ACTIVATOR 1    ACTIVATOR 2    REPRESSOR 1

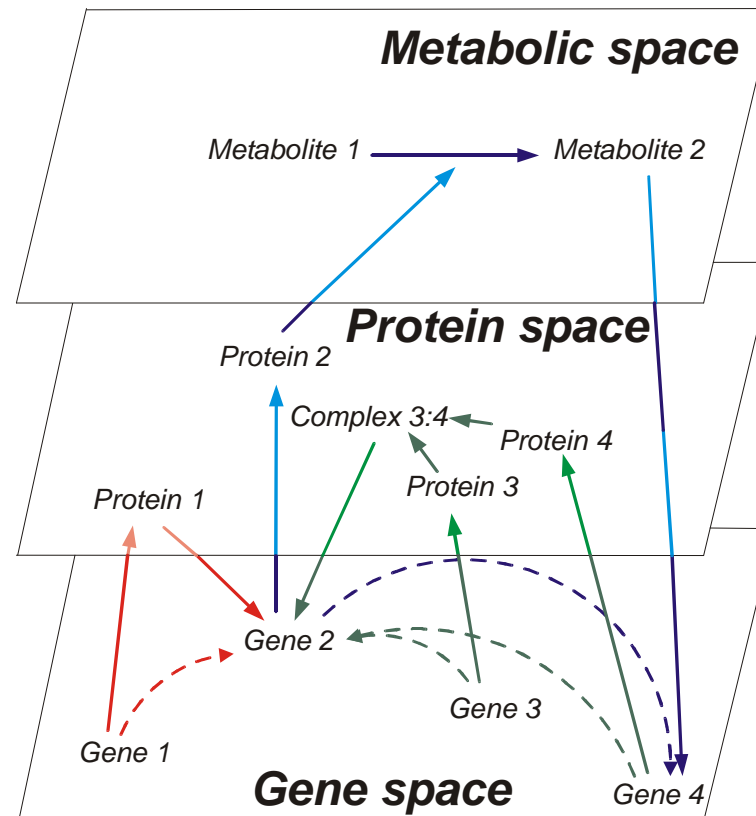A1    A2    R1

T1

TARGET 1

Transcription regulatory network in baker's yeast Saccharomyces Serevisiae
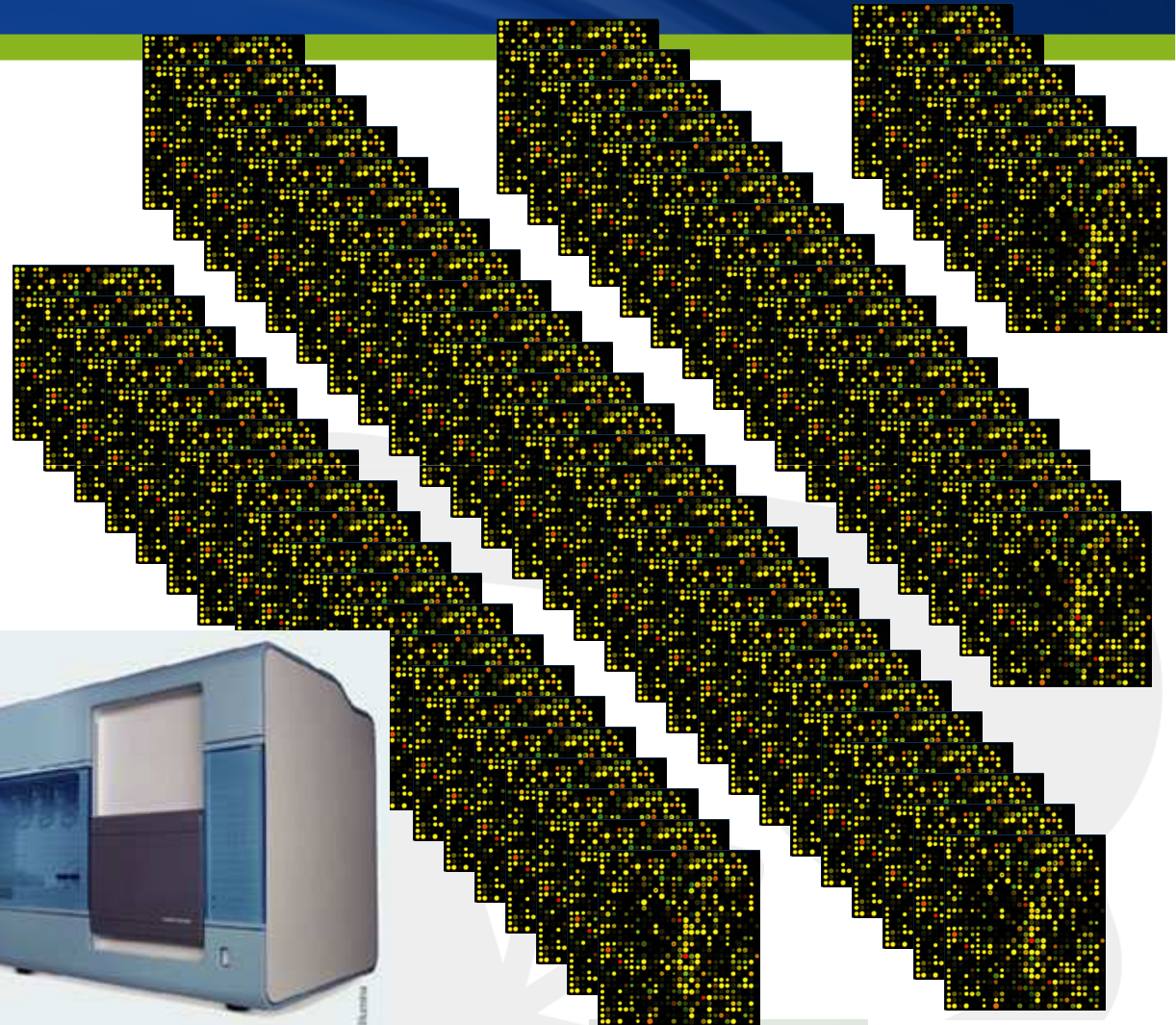
$$A = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{pmatrix}$$

$$A_W = \begin{pmatrix} a_{11} & 0 & 0 & 0 & a_{15} \\ a_{21} & a_{22} & 0 & 0 & 0 \\ 0 & a_{32} & a_{33} & 0 & 0 \\ 0 & 0 & a_{43} & a_{44} & 0 \\ 0 & a_{52} & 0 & 0 & a_{55} \end{pmatrix}$$

# Gene expression data



Matrix representation
of data:

$$\mathbf{X}_{p \times n}$$

(p = #genes, n = #observations)

- Introduction to Gene networks

- **Gene network inference**

- Evaluation of gene network inference algorithms

- Differential networking in disease

"~omics" data



**algorithms**

*Correlation, partial correlation, regression, linear Ordinary Differential Equations, graphical Gaussian models, perturbation analysis…*

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{k})$$
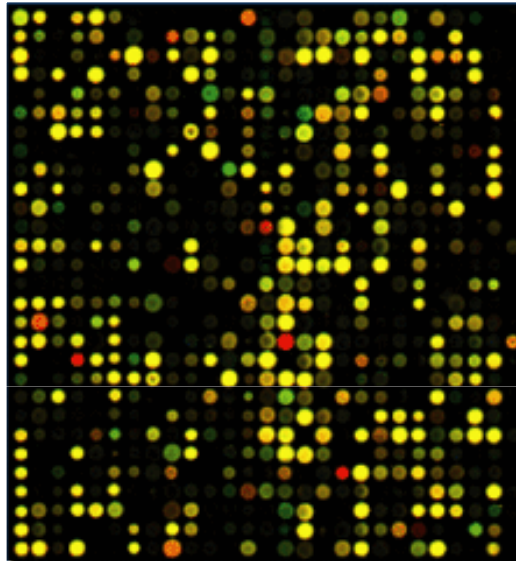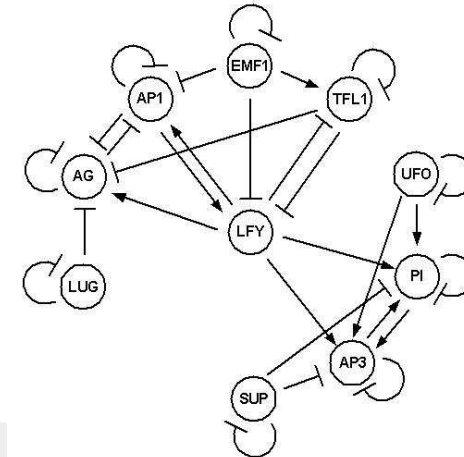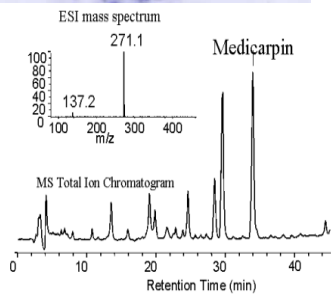
$$A = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{pmatrix}$$

$$A_W = \begin{pmatrix} a_{11} & 0 & 0 & 0 & a_{15} \\ a_{21} & a_{22} & 0 & 0 & 0 \\ 0 & a_{32} & a_{33} & 0 & 0 \\ 0 & 0 & a_{43} & a_{44} & 0 \\ 0 & a_{52} & 0 & 0 & a_{55} \end{pmatrix}$$

'Observational data'

Repeated measurements of a given tissue/cell type without experimental intervention

ALLOWS ONLY FOR INFERRING **UNDIRECTED** NETWORKS

'Perturbation data'

Creating targeted perturbations and measuring systems dynamic responses (steady states or time-series)

ALLOWS FOR INFERRING **DIRECTED** NETWORKS

Causal graph

$T_1$
$T_2$
$T_3$    $T_4$
$T_5$

Correlation graph

$T_1$
$T_2$
$T_3$    $T_4$
$T_5$

$$r_{xy.z} = \frac{r_{xy} - r_{xz}r_{yz}}{\sqrt{\left(1 - r_{xz}^2\right)\left(1 - r_{yz}^2\right)}}$$

$$r_{xy.zq} = \frac{r_{xy.z} - r_{xq.z}r_{yq.z}}{\sqrt{\left(1 - r_{xq.z}^2\right)\left(1 - r_{yq.z}^2\right)}}$$

Remove edges with zero (non significant) partial correlations

$$r_{T_1 T_5 . T_2} \approx 0$$

$$r_{T_3 T_4 . T_2} \approx 0$$

$$r_{T_2 T_5 . T_3 T_4} \approx 0$$

*The correlation between T3 and T4 disappears when conditioned on T2, because T2 is a causal parent of both T3 and T4*

**CRS4**
IDEAS BECOME LIFE

## Discovery of meaningful associations in genomic data using partial correlation coefficients

Alberto de la Fuente*, Nan Bing†, Ina Hoeschele and Pedro Mendes

Virginia Polytechnic Institute and State University, Virginia Bioinformatics Institute, 1880 Pratt Drive, Blacksburg, Virginia, 24061 USA

**ABSTRACT**
**Motivation:** A major challenge of systems biology is to infer biochemical interactions from large-scale observations, such as transcriptomics, proteomics and metabolomics. We propose to use a partial correlation analysis to construct approximate Undirected Dependency Graphs from such large-scale biochemical data. This approach enables a distinction between direct and indirect interactions of biochemical compounds, thereby inferring the underlying network topology

about the underlying network top
assumed that biochemical netwo
ected acyclic graphs (Friedman *e*
However, cyclic network structur
are ubiquitous in biology and are a
specific properties of living syste
should be independent of such as
We propose a method to cons
ted dependency graphs (UDGs) fr

Online journal ISSN - 1676-5680
**GMR** Genetics and Molecular Research
EVOLUTION AND TECHNOLOGY

## Genome-wide partial correlation analysis of *Escherichia coli* microarray data

D.F.T. Veiga[1]*, F.F.R. Vicente[1]*, M. Grivet[2], A. de la Fuente[3] and A.T.R. Vasconcelos[1]

de la Fuente A, Bing N, Hoeschele I and Mendes P. Discovery of meaningful associations in genomic data using partial correlation coefficients Bioinformatics, 2004, 20(18):3565-3574

Veiga, D.F., da Rocha Vicente, F.F., Grivet, M, de la Fuente, A., Ribeiro de Vasconcelos, A.T. (2007) Genome-wide Partial Correlation Analysis of Escherichia coli Microarray Data. Genetics and Molecular Research 6(4): 730-742

- Steady state perturbation data

Wild type

Over-expression Gene 1

Over-expression Gene 2

Over-expression Gene 3

Over-expression Gene *n*

- Time series data

Stress

**Data**:

- Steady state mRNA concentration/gene expression levels
  - » Wild-type
  - » Systematic single gene **knockdowns** or **over-expression**
    - » Heterozygous knockout
    - » Expression from plasmid

Measure gene-expression in unperturbed (WT) state

Perturb each gene and measure gene-expression responses

## Distinguish direct from indirect edges:

Algebraic relation between the deviation matrix **X** (perturbed levels – wild type levels) and the network matrix (encoding the network **A** of direct interactions)

$$
\begin{pmatrix}
a_{11} & 0 & 0 & 0 & a_{15} \\
a_{21} & a_{22} & 0 & 0 & 0 \\
0 & a_{32} & a_{33} & 0 & 0 \\
0 & 0 & a_{43} & a_{44} & a_{54} \\
0 & 0 & 0 & 0 & a_{55}
\end{pmatrix}
=
\begin{pmatrix}
\Delta x_{11} & 0 & 0 & 0 & \Delta x_{15} \\
\Delta x_{21} & \Delta x_{22} & 0 & 0 & \Delta x_{25} \\
\Delta x_{31} & \Delta x_{32} & \Delta x_{33} & 0 & \Delta x_{35} \\
\Delta x_{41} & \Delta x_{42} & \Delta x_{43} & \Delta x_{44} & \Delta x_{45} \\
0 & 0 & 0 & 0 & \Delta x_{55}
\end{pmatrix}^{-1}
$$

# Linear modeling approach

$$\frac{dx_i}{dt} = g_i(x_j; p_k) = g_i(x_j^0 + \Delta x_j; p_k^0 + \Delta p_k)$$

$$\approx g_i(x_j^0; p_k^0) + \sum_{j=1}^{n} \left.\frac{\partial g_i}{\partial x_j}\right|_{x^0,p^0} \Delta x_j$$

$$+ \sum_{k=1}^{p} \left.\frac{\partial g_i}{\partial p_k}\right|_{x^0,p^0} \Delta p_k$$

$$\Rightarrow \frac{d\Delta x_i}{dt} \approx \sum_{j=1}^{n} \left.\frac{\partial g_i}{\partial x_j}\right|_{x^0,p^0} \Delta x_j$$

$$+ \sum_{k=1}^{p} \left.\frac{\partial g_i}{\partial p_k}\right|_{x^0,p^0} \Delta p_k$$

$$\equiv \sum_{j=1}^{n} a_{ij}\Delta x_j + \sum_{k=1}^{p} r_{ik}\Delta p_k \quad (1)$$

$$\frac{d\Delta x_i}{dt} = \sum_{j}^{n} a_{ij}\Delta x_j + \Delta u_i$$

$$0 = \sum_{j}^{n} a_{ij}\Delta x_j + \Delta u_i \qquad \sum_{j}^{n} a_{ij}\Delta x_j = -\Delta u_i$$

$$\mathbf{JX} = -\mathbf{U}$$

$\mathbf{J} = \{a_{ij}\}$ Effect of gene $j$ on rate of change of gene $i$

$\mathbf{U} = \{u_{kk}\}$ Diagonal perturbation matrix

$\mathbf{X} = \{x_{ik}\}$ Change in gene $i$ expression after perturbation $k$

$$\mathbf{J} = -\mathbf{UX}^{-1}$$

$$\mathbf{R} = \mathbf{U}^{-1}\mathbf{J} = -\mathbf{X}^{-1}$$

THE CHALLENGES OF SYSTEMS BIOLOGY

## Inferring Gene Networks: Dream or Nightmare?

### Part 2: Challenges 4 and 5

Alan Scheinine, Wieslawa I. Mentzen, Giorgio Fotia,
Enrico Pieroni, Fabio Maggio, Gianmaria Mancosu,
and Alberto de la Fuente

*CRS4 Bioinformatica, Pula, Italy*

We describe several algorithms with winning performance in the Dialogue for Reverse Engineering Assessments and Methods (DREAM2) Reverse Engineering Competition 2007. After the gold standards for the challenges were released, the performance of the algorithms could be thoroughly evaluated under different parameters or alternative ways of solving systems of equations. For the analysis of Challenge 4, the "*In-silico*"

**Opinion**          *TRENDS in Genetics*  Vol.18 No.8 August 2002

## Linking the genes: inferring quantitative gene networks from microarray data

Alberto de la Fuente, Paul Brazhnik and Pedro Mendes

Trends Genet. 2002 Aug;18(8):395-8

Scheinine, A., Mentzen, W., Pieroni E., Fotia, G., Maggio, F., Mancosu, G. and de la Fuente, A. (2009) Inferring Gene Networks: Dream or nightmare? Part 2: Challenges 4 and 5. Annals of the New York Academy of Sciences 1158: 287301
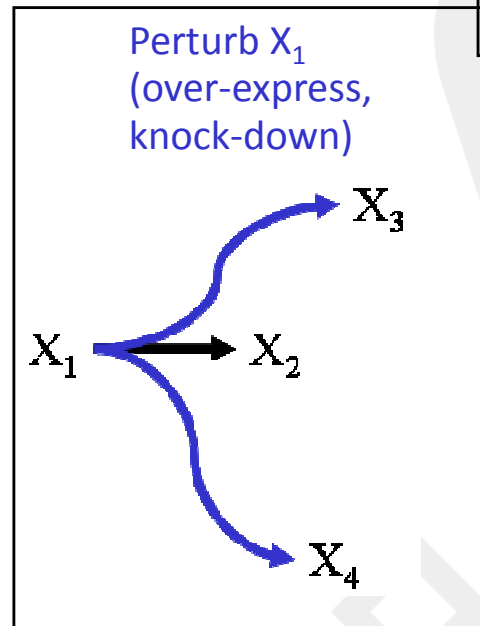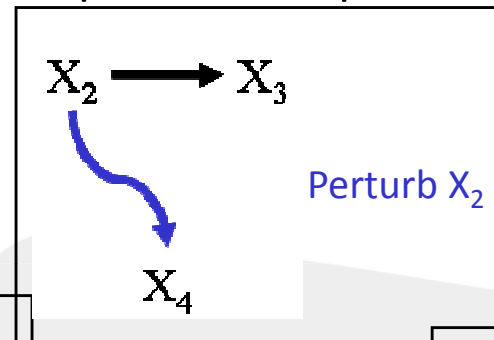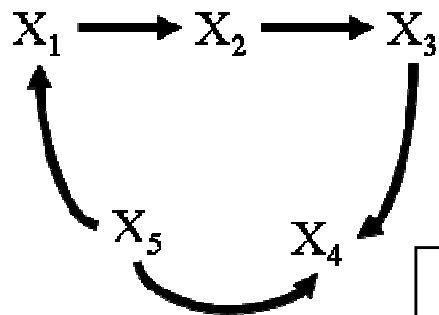
**Data**:

- Steady state mRNA concentration/gene expression levels
  - » Wild-type
  - » Systematic single gene **knock-outs**
    - » Complete removal of genes

- Weight estimation for edge $i \rightarrow j$: **change** in the mRNA level $x_{i,j}$ of gene $j$ after **knockout** of gene $i$
- **Z-score**:

$$W_{i,j} = \frac{x_{i,j} - \overline{x}_{\cdot,j}}{s_{\cdot,j}}$$

- The edge weight measures the total causal effect of a gene on another gene: *direct* or *mediated*?



- The initial network can have many feed-forward loops

  - Not essential for reachability

  - We want to rank them lower than "essential" edges

- Algorithm:

  1) Fix a **threshold** for weights and determine a network

  2) **Delete** feed-forward edges between strongly connected components of the network

  3) **Increase** the weight of remaining edges in W

W ➡  ➡ W*

**Result**:
Essential edges (solid) are ranked higher than feed-forward edges (dashed)

OPEN ACCESS Freely available online

PLoS one

# From Knockouts to Networks: Establishing Direct Cause-Effect Relationships through Graph Analysis

Andrea Pinna, Nicola Soranzo, Alberto de la Fuente*

Center for Advanced Studies, Research and Development (CRS4) Bioinformatica, Pula, Italy

## Abstract

*Background:* Reverse-engineering gene networks from expression profiles is a difficult problem for which a multitude of techniques have been developed over the last decade. The yearly organized DREAM challenges allow for a fair evaluation and unbiased comparison of these methods.

*Results:* We propose an inference algorithm that combines confidence matrices, computed as the standard scores from single-gene knockout data, with the down-ranking of feed-forward edges. Substantial improvements on the predictions can be obtained after the execution of this second step.

*Conclusions:* Our algorithm was awarded the best overall performance at the DREAM4 In Silico 100-gene network sub-challenge, proving to be effective in inferring medium-size gene regulatory networks. This success demonstrates once again the decisive importance of gene expression data obtained after systematic gene perturbations and highlights the usefulness of graph analysis to increase the reliability of inference.

*Citation:* Pinna A, Soranzo N, de la Fuente A (2010) From Knockouts to Networks: Establishing Direct Cause-Effect Relationships through Graph Analysis, PLoS

Pinna, A., Soranzo, N. and de la Fuente, A. (2010) From Knockouts to Networks: Establishing Direct Cause-Effect Relationships through Graph Analysis, PLoS ONE 5(10), e12912 (DREAM4 Special Collection)

Figure 7 from: GeneNetWeaver: In silico benchmark generation and performance profiling of network inference methods. Schaffter T, Marbach D, Floreano D. Bioinformatics (2011) 27 (16): 2263-2270.

# Natural genetic perturbations



**FIGURE 1**

(a) Parents

(c) Microarray per offspring

(b) Segregating population

(d) Markers per offspring

cDNA1

cDNA2

Marker B/b

Marker A/a

*TRENDS in Genetics*

Jansen, R.C., and Nap, J.P. (2001) Trends Genet. 17, 388-391

**Gene Network inference requires many perturbations**

**Experimental perturbations are difficult and costly**

**Use of naturally occurring genetic variations (perturbations)**

$$y_{in} = b_0 + b_1 x_{jn} + \varepsilon_{in}$$

$$y_{in} = b_0 + b_1 y_{jn} + b_2 x_{jn} + \varepsilon_{in}$$

x = genotype data (e.g. SNPs)
y = gene expression 'phenotypes'

# Gene Network Inference via Structural Equation Modeling in Genetical Genomics Experiments

Bing Liu,[*,†,1,2] Alberto de la Fuente[†,‡,1] and Ina Hoeschele[*,†,3]

*Department of Statistics, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24061, †Virginia Bioinformatics Institute, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24061-0477 and ‡CRS4 Bioinformatica, Parco Scientifico e Tecnologico, POLARIS, 09010 Pula (CA), Italy

## ABSTRACT
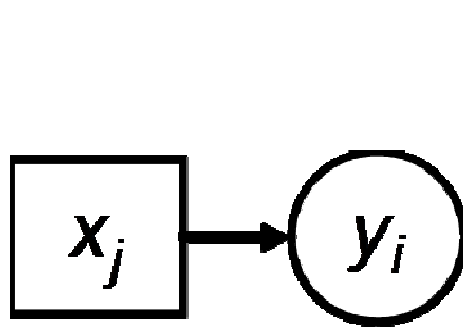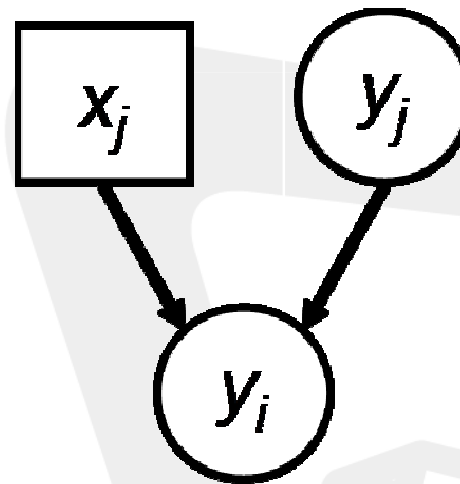
Our goal is gene network inference in genetical genomics or systems genetics experiments. For species where sequence information is available, we first perform expression quantitative trait locus (eQTL) mapping by jointly utilizing cis-, cis–trans-, and trans-regulation. After using local structural models to identify regulator–target pairs for each eQTL, we construct an encompassing directed network (EDN) by assembling all retained regulator–target relationships. The EDN has nodes corresponding to expressed genes and eQTL and directed edges from eQTL to cis-regulated target genes, from cis-regulated genes to cis–trans-regulated target genes, from trans-regulator genes to target genes, and from trans-eQTL to target

- Introduction to Gene networks

- Gene network inference

- **Evaluation of gene network inference algorithms**

- Differential networking in disease

A yeast synthetic network for in vivo assessment of reverse-engineering and modeling approaches. Cantone I, Marucci L, Iorio F, Ricci MA, Belcastro V, Bansal M, Santini S, di Bernardo M, di Bernardo D, Cosma MP. Cell. 2009 Apr 3;137(1):172-81. Epub 2009 Mar 26.

# In silico algorithm evaluation

## Precision-Recall example: AUC=0.82

Precision

Recall

Precision-Recall curve

| | actual class (expectation) | |
|---|---|---|
| predicted class (observation) | tp (true positive) Correct result | fp (false positive) Unexpected result |
| | fn (false negative) Missing result | tn (true negative) Correct absence of result |

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

## Dialogue for Reverse Engineering Assessments and Methods
http://wiki.c2b2.columbia.edu/dream/index.php/The_DREAM_Project

- DREAM2, best performer in:
  - Synthetic Five-Gene Network Inference
  - DREAM2 In Silico Network Challenge

- DREAM4, best performer in:
  - DREAM4 In Silico Network Challenge
    - Size 100 subchallenge

- DREAM5, honorary mention in:
  - Network inference challenge

* DREAM6: top 3 RNA-seq challenge

# SysGenSIM: Simulating Gene Network dynamics



A  B

(Mendes et al., 2003)

**FIGURE 3**

Cis-effect: a polymorphism in promotor region of gene $G_g$ affects the basal transcription rate.

$Z_g$ for allele 'A' > $Z_g$ allele 'a'.

Trans-effect: a polymorphism in coding region of gene $G_k$ affects the strength of the effect on its targets.

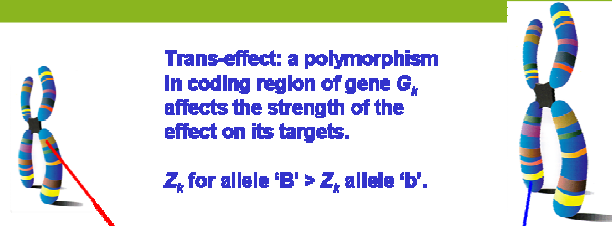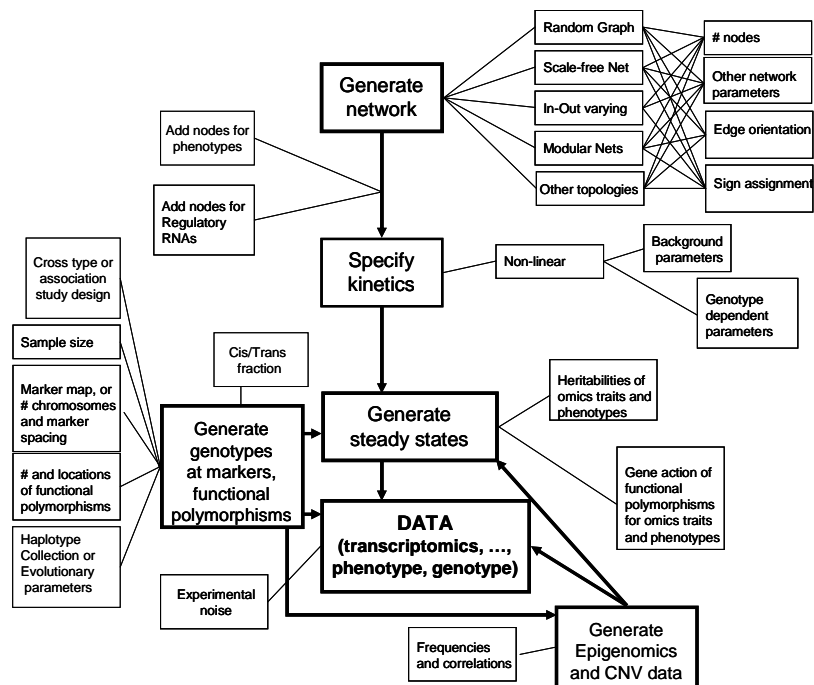$Z_k$ for allele 'B' > $Z_k$ allele 'b'.

$$\frac{dG_g}{dt} = v_{transcription_{G_g}} - v_{degradation_{G_g}} = \boxed{Z_g} \cdot V_g \cdot \prod_{k \in R_g}\left(1 + A_{gk}\frac{G_k^{h_{gk}}}{G_k^{h_{gk}} + \left(K_{gk}/\boxed{Z_k}\right)^{h_{gk}}}\right) - \theta_g k_g G_g$$

**Reason**: Many algorithms have been (and even more will be) proposed for Gene Network Inference: **need for unbiased evaluation**

SysGenSIM has been used to generate a challenge in **DREAM5, STAT-SEQ COST, Springer book**

Currently in MATLAB, but we want to reprogram in Python

Part of NIH project

## Flow diagram

- Generate network
  - Random Graph
  - Scale-free Net
  - In-Out varying
  - Modular Nets
  - Other topologies
    - # nodes
    - Other network parameters
    - Edge orientation
    - Sign assignment
- Add nodes for phenotypes
- Add nodes for Regulatory RNAs
- Specify kinetics
  - Non-linear
    - Background parameters
    - Genotype dependent parameters
- Cross type or association study design
- Sample size
- Marker map, or # chromosomes and marker spacing
- # and locations of functional polymorphisms
- Haplotype Collection or Evolutionary parameters
- Cis/Trans fraction
- Generate genotypes at markers, functional polymorphisms
- Generate steady states
  - Heritabilities of omics traits and phenotypes
  - Gene action of functional polymorphisms for omics traits and phenotypes
- **DATA (transcriptomics, ..., phenotype, genotype)**
- Experimental noise
- Frequencies and correlations
- Generate Epigenomics and CNV data

# SysGenSIM

## Simulating systems genetics data with SysGenSIM

Andrea Pinna[1], Nicola Soranzo[1], Ina Hoeschele[2,3] and Alberto de la Fuente[1,*]

[1]CRS4 Bioinformatica, 09010 Pula (CA), Italy, [2]Department of Statistics, Virginia Polytechnic Institute and State University, Blacksburg, VA 24061 and [3]Virginia Bioinformatics Institute, Virginia Polytechnic Institute and State University, Blacksburg, VA 24061-0477, USA

Associate Editor: Martin Bishop

**ABSTRACT**

**Summary:** SysGenSIM is a software package to simulate Systems Genetics (SG) experiments in model organisms, for the purpose of evaluating and comparing statistical and computational methods and their implementations for analyses of SG data [e.g. methods for expression quantitative trait loci (eQTL) mapping and network inference]. SysGenSIM allows the user to select a variety of network topologies, genetic and kinetic parameters to simulate SG data

known that the etraits of groups of genes share common regulators (DNA variants), which are more easily identified when associated with a group of etraits rather than with individual etraits. Several approaches to associating DNA variants with groups of etraits have recently been proposed (e.g. Chun and Keles, 2009; Lee *et al.*, 2009, 2006; Parkhomenko *et al.*, 2007; Waaijenborg *et al.*, 2008; Zhang *et al.*, 2010).

A major goal of SG studies is to reconstruct a causal network

COMMENTARY

# Verification of systems biology research in the age of collaborative competition

Pablo Meyer[1], Leonidas G Alexopoulos[2], Thomas Bonk[3], Andrea Califano[4], Carolyn R Cho[5], Alberto de la Fuente[6], David de Graaf[7], Alexander J Hartemink[8], Julia Hoeng[3], Nikolai V Ivanov[3], Heinz Koeppl[9], Rune Linding[10], Daniel Marbach[11], Raquel Norel[1], Manuel C Peitsch[3], J Jeremy Rice[1], Ajay Royyuru[1], Frank Schacherer[12], Joerg Sprengel[13], Katrin Stolle[3], Dennis Vitkup[4] & Gustavo Stolovitzky[1]

Collaborative competitions in which communities of researchers compete to solve challenges may facilitate more rigorous scrutiny of scientific results.
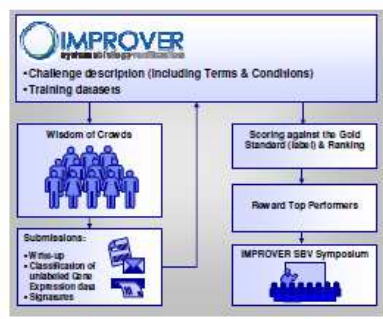
# IMPROVER

## Diagnostic Signature Challenge



- Challenge description (including Terms & Conditions)
- Training datasets
- Wisdom of Crowds
- Scoring against the Gold Standard (label) & Ranking
- Reward Top Performers
- Submissions:
  - Write-up
  - Classification of unlabeled Gene Expression data
  - Signatures
- IMPROVER SBV Symposium

The goal of the diagnostic signature challenge is to assess and verify computational approaches that classify clinical samples based on transcriptomics data.
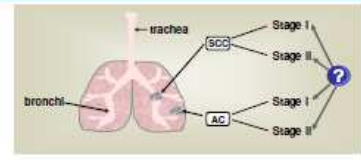
## Psoriasis Sub-Challenge



The challenge is to develop a classifier that differentiates healthy skin from that with psoriatic lesions.
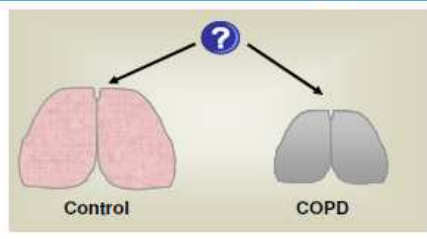
The classifier will be built by using publicly available gene expression data with their psoriasis-related clinical information (e.g. label). The classifier will be tested on an unpublished independent high quality dataset.

## Lung Cancer Sub-Challenge



The challenge is to classify lung cancer subtypes [Adenocarcinoma (AC) and Squamous Cell Carcinoma (SCC)] and their respective stages (I & II) based on transcriptomics data from tumor samples.

The classifier will be built by using publicly available gene expression data with the respective histo-pathological information. The classifier will be tested on an independent high quality dataset.
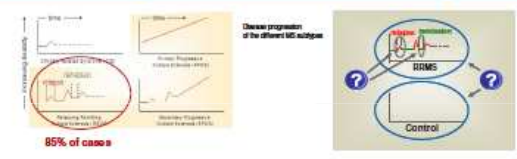
## Chronic Obstructive Pulmonary Disease Sub-Challenge



Control    COPD

The challenge is to develop a classifier that differentiates COPD vs control based on the airway transcriptome from clinical samples.

The classifier will be built by using publicly available gene expression data with clinical information. The classifier will be tested on an independent unpublished high quality dataset.

## Multiple Sclerosis Sub-Challenge



85% of cases

The challenge is to develop a classifier that differentiates clinical samples in two ways:
- control vs. multiple sclerosis
- relapsing vs remitting multiple sclerosis

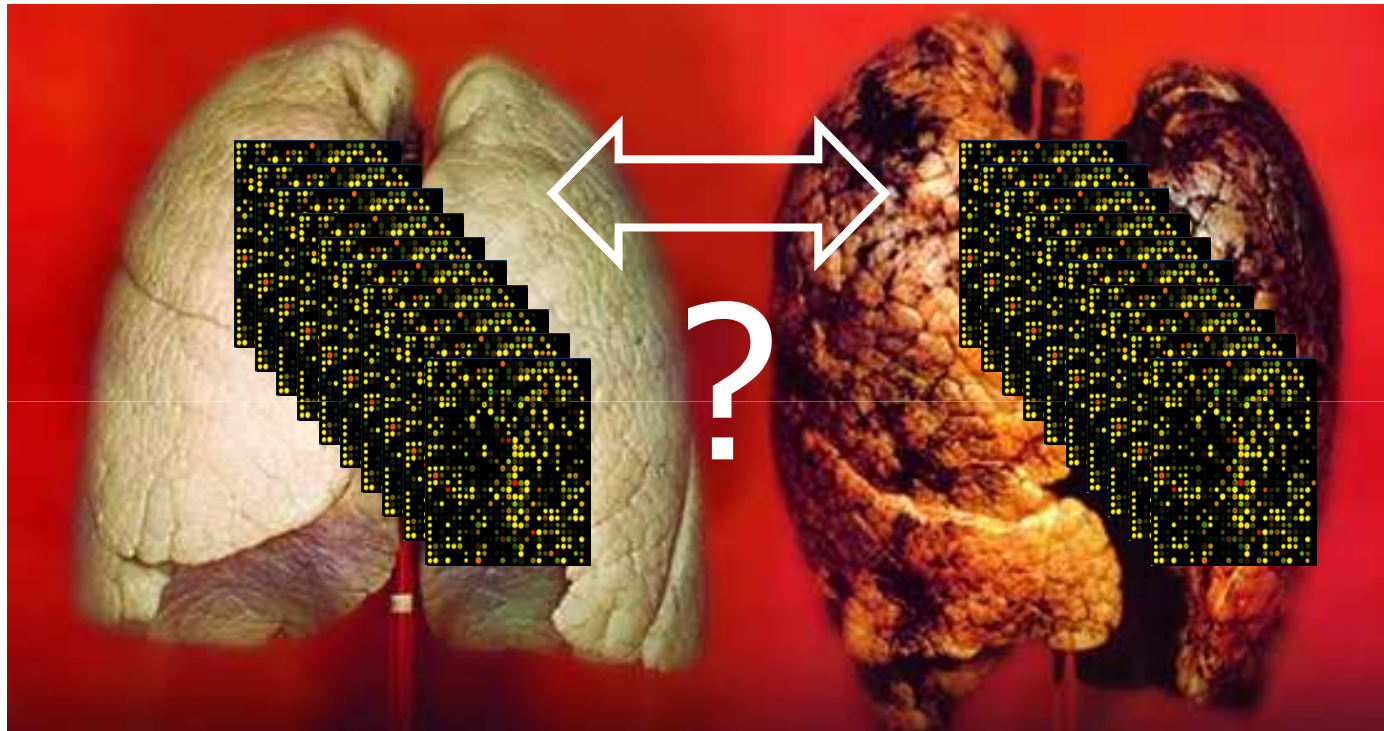based on transcriptome measured in Periph Mononuclear Cells (PBMC).

The classifier will be built by using publicly av expression data with clinical information. The be tested on two independent unpublished datasets.

## References

1. Meyer P. et al, Nature Biotechnology 29(9):811-815 (2011) systems biology research in the age of collaborative com
2. Marbach D. et al., Proc Natl Acad Sci U S A 107(14):6286- Revealing strengths and weaknesses of methods for gen inference.
3. Norel R. et al., Mol. Sys. Bio 7:537 (2011) The self-asses all be better than average?
4. Prill R.J. et al., PLoS ONE, 5(2):e9202 (2010) Towards a Assessment of Systems Biology Models: The DREAM3 C

- Website Launch
- Challenge Open
- Registration Open

- Scoring Complete
- Ranking Available
- Announce Best Performing Team

5th March 2012 → 30th May 2012 → 28th June 2012 → 28th/29th August 2012

- Submission Deadline for Predictions and Write-ups
- Challenge Closed

- SBV Symposium
- Award Top Performing Teams
- Share Results & Experiences

- Introduction to Gene networks

- Gene network inference

- Evaluation of gene network inference algorithms

- **Differential networking in disease**
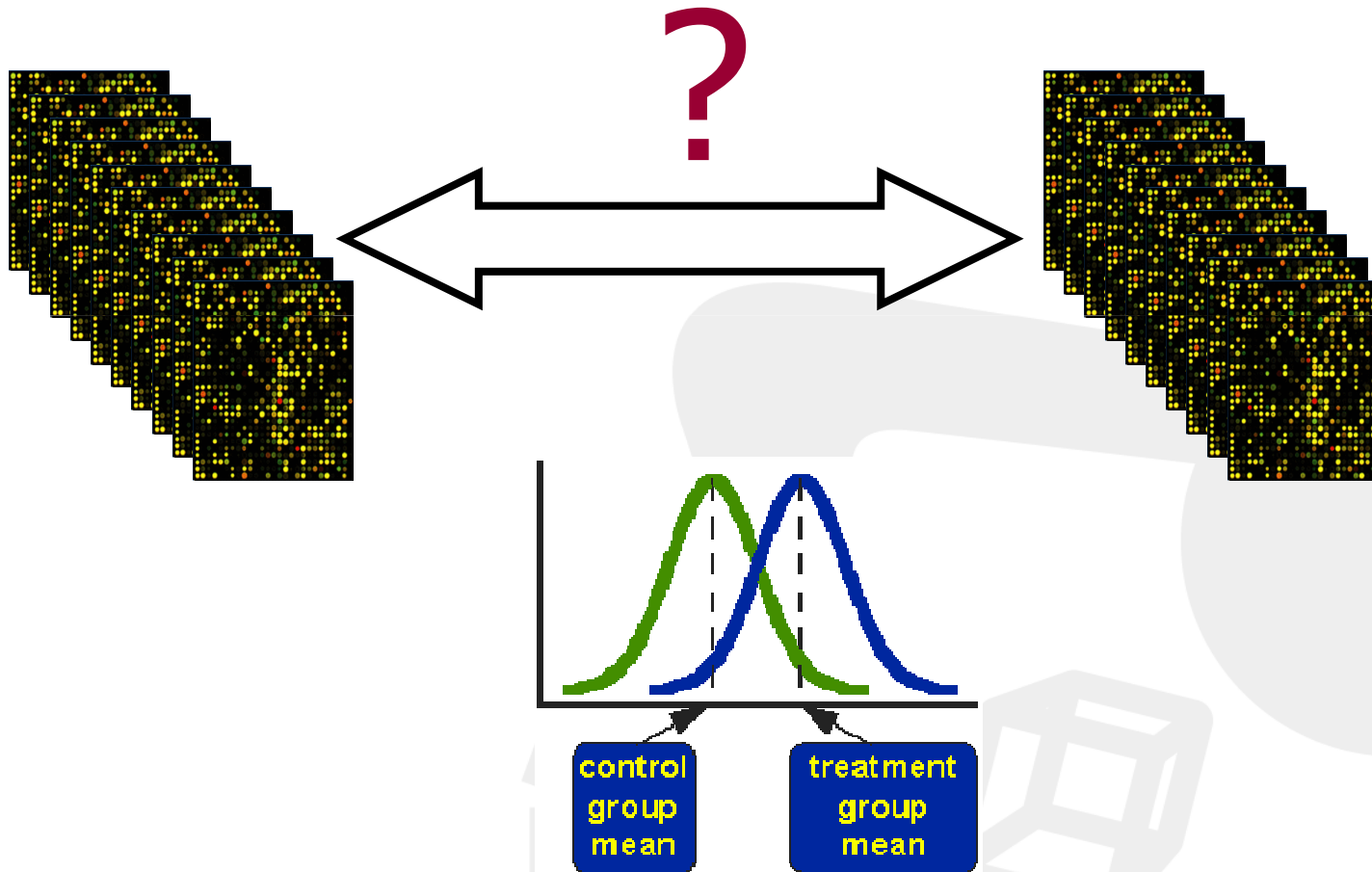
Group 1 (healthy tissue, treated with medicine, tumor stage X, etc.)
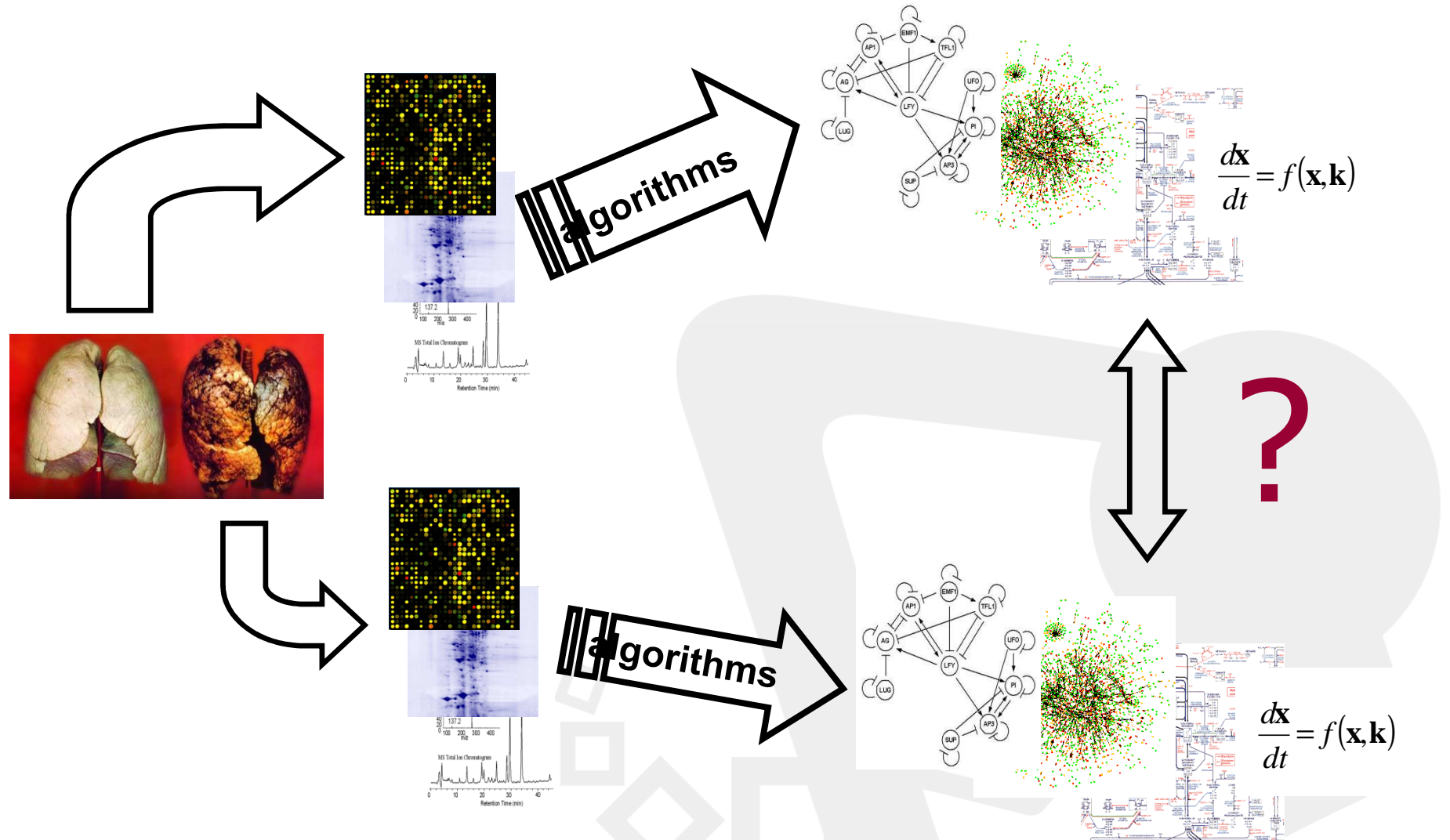
Group 2 (tumor tissue, not treated with medicine, tumor stage Y, etc.)

'Differential networking'

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{k})$$

algorithms

?

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{k})$$

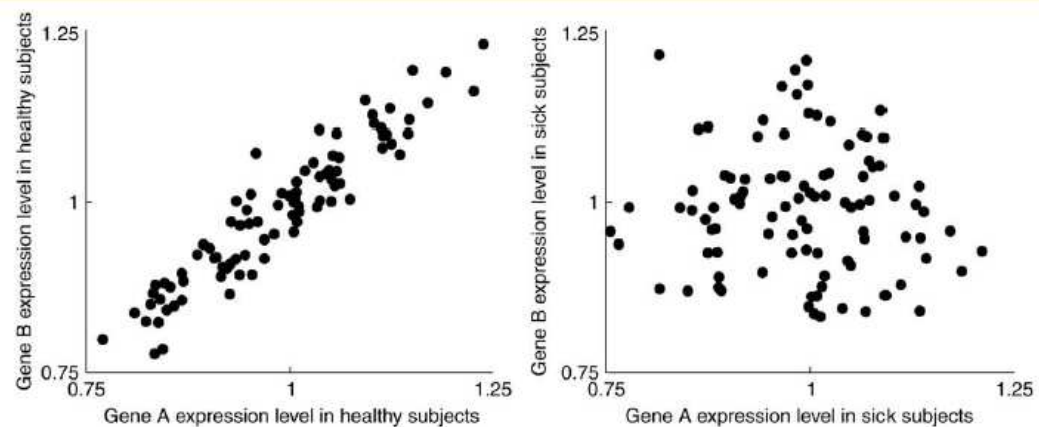algorithms

**Review**

**Cell** PRESS

# From 'differential expression' to 'differential networking' – identification of dysfunctional regulatory networks in diseases

**Alberto de la Fuente**

CRS4 Bioinformatica, Polaris Edificio 3, Località Piscina Manna, 09010 Pula (CA), Italy

Understanding diseases requires identifying the di text essentially take this approach to differential coexpression. Testing
ences between healthy and affected tissues. C
expression data have revolutionized the study of
eases by making it possible to simultaneously cons
thousands of genes. The identification of disease-as
ated genes requires studying the genes in the conte
the regulatory systems they are involved in. A major



*TRENDS in Genetics*

$$D = \sqrt{\frac{1}{p(p-1)/2} \times \sum_{i=1}^{p} \sum_{j=i+1}^{p} \left(r_{ij}^{healty} - r_{ij}^{sick}\right)^2}$$

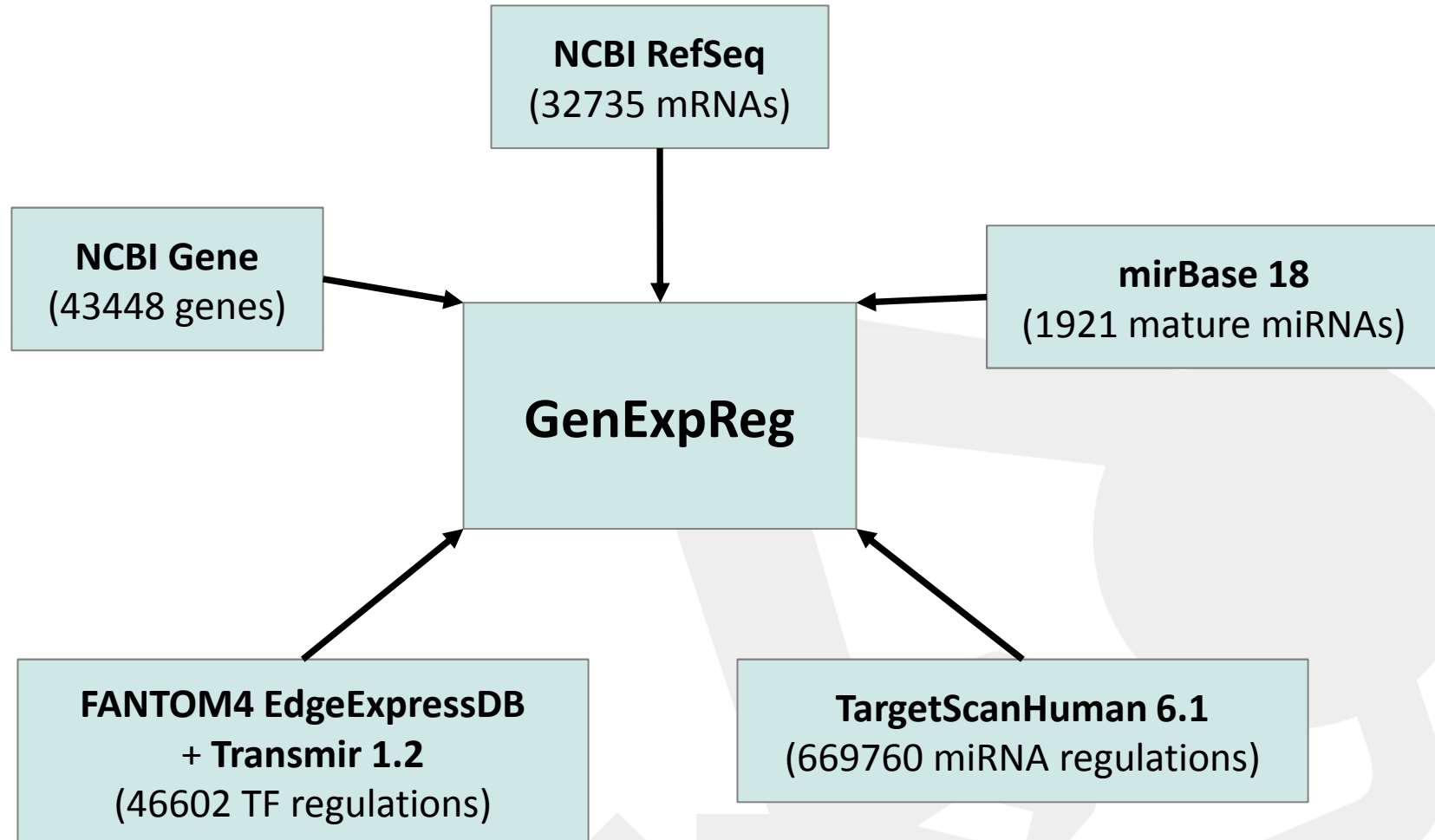AUPRC for knockout of 10 microRNAs



$$\frac{dG_j}{dt} = V_j \cdot \theta_j^{syn} - \lambda_j \cdot \theta_j^{deg} \cdot G_j \prod_{i=1}^{n} \left( 1 + A_{i,j}^{deg} \frac{G_i^{h_{i,j}^{deg}}}{G_i^{h_{i,j}^{deg}} + K_{i,j}^{deg}{}^{h_{i,j}^{deg}}} \right)$$

$$D_{L2}(S) = \sqrt{\frac{2}{|S|(|S|-1)} \sum_{i,j \in S, i<j} \left( \rho_1(i,j) - \rho_2(i,j) \right)^2}$$

$$D_{L1}(S) = \frac{2}{|S|(|S|-1)} \sum_{i,j \in S, i<j} \left| \rho_1(i,j) - \rho_2(i,j) \right|$$

**Bhattacharjee,A. *et al*. (2001) Classification of human lung carcinomas by mRNA expression profiling reveals distinct adenocarcinoma subclasses. *Proc. Natl Acad. Sci.*, 98, 13790-13795.**

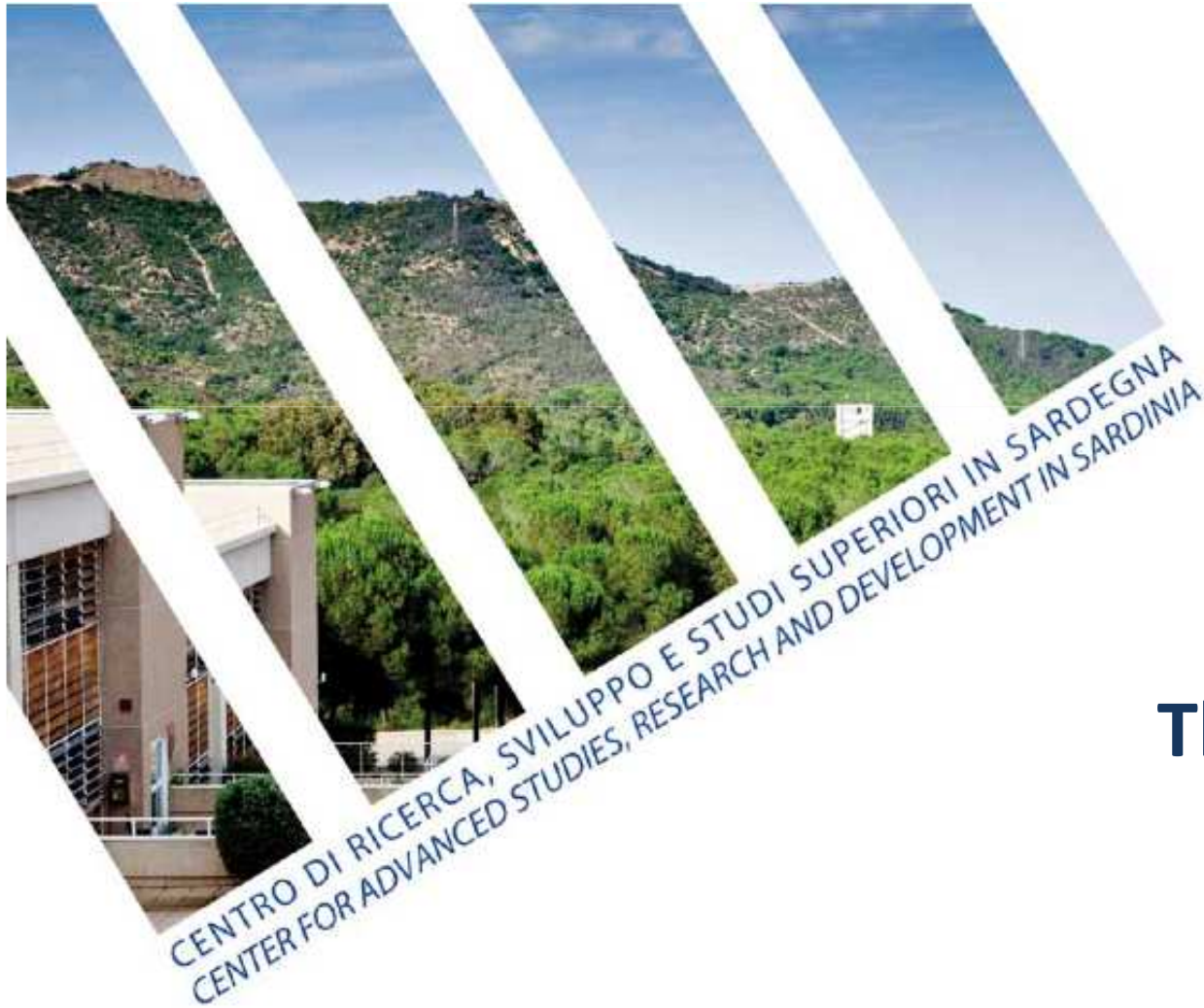| Family name | Seed | N. of target | N. of target | P-value for | Notes |
|---|---|---|---|---|---|
| miR-1293 | GGGUGGU | 73 | 23 | 0.0022 | |
| miR-28/28-3p | ACUAGAU | 77 | 19 | 0.0024 | upregulated in serum copy number of lung cancer patients w.r.t. healthy [1] |
| miR-1244 | AGUAGUU | 147 | 53 | 0.0027 | |
| miR-1269 | UGGACUG | 77 | 21 | 0.0048 | |
| miR-1224/1224-5p | UGAGGAC | 88 | 34 | 0.0050 | |
| miR-578 | UUCUUGU | 229 | 65 | 0.0052 | |
| miR-1305 | UUUCAAC | 414 | 106 | 0.0060 | |
| miR-433 | UCAUGAU | 207 | 63 | 0.0061 | |
| miR-205 | CCUUCAU | 288 | 92 | 0.0063 | highly specific marker for squamous cell lung carcinoma [2] and non-small cell lung cancer [3]; located in a region amplified in lung cancer; upregulated in lung cancer tissues w.r.t. noncancerous lung tissues [4] |
| miR-1237 | CCUUCUG | 177 | 42 | 0.0082 | |
| miR-520a-5p/525-5p | UCCAGAG | 296 | 79 | 0.0085 | |
| miR-582-3p | AACUGGU | 97 | 46 | 0.0086 | |
| miR-568 | UGUAUAA | 308 | 85 | 0.0087 | |
| miR-432 | CUUGGAG | 133 | 37 | 0.0090 | member of miR-127 cluster, which is downregulated in tumors [5] |
| miR-524-3p/525-3p | AAGGCGC | 38 | 10 | 0.0091 | |
| miR-513c | UCUCAAG | 223 | 64 | 0.0094 | |
| miR-370 | CCUGCUG | 239 | 52 | 0.0096 | downregulated after lung development [6] |

[1] Chen, X., et al. - Cell Res. 18(10) pp. 997–1006 – 2008
[2] Lebanony, D., et al. - J. Clinical Oncology 27(12) – pp. 2030-2037 – 2009
[3] Markou, A., et al. – Clin. Chem. 54(10) – pp. 1696-1704 – 2008
[4] Yanaihara, N., et al. - Cancer Cell 9(3) – pp. 189-198 – 2006
[5] Saito, Y., et al. - Cancer Cell 9(6) – pp. 435-443 – 2006
[6] Williams, A. E., et al. - Dev. Dyn. 236(2) – pp. 572-580 – 2007

**Thank you for your attention**