# Machine Learning for Personalized Genomics

## Jean-Philippe Vert



*C3BI Kick-off meeting, Institut Pasteur,*
*Paris, March 16, 2015*

# Institut Curie / Inserm U900 / MINES ParisTech partnership



- A joint lab about ``Cancer computational genomics, bioinformatics, biostatistics and epidemiology''
- Located in Institut Curie, a major hospital and cancer research centre in Europe, and MINES ParisTech

# 4 teams + 1 platform

**Systems Biology (Barillot):**
- Modelling, simulating biological systems
- Building an *in silico* atlas of cancer pathways

**Clinical Biostatistics (Asselain / Paoletti):**
- Clinical trials for targeted therapies
- Predictive biomarkers

**Cancer Genetic Epidemiology (Andrieu):**
- Genetic and environmental factors in breast cancer

**Machine learning (Vert):**
- Learning from « big omics data » for personalized medicine
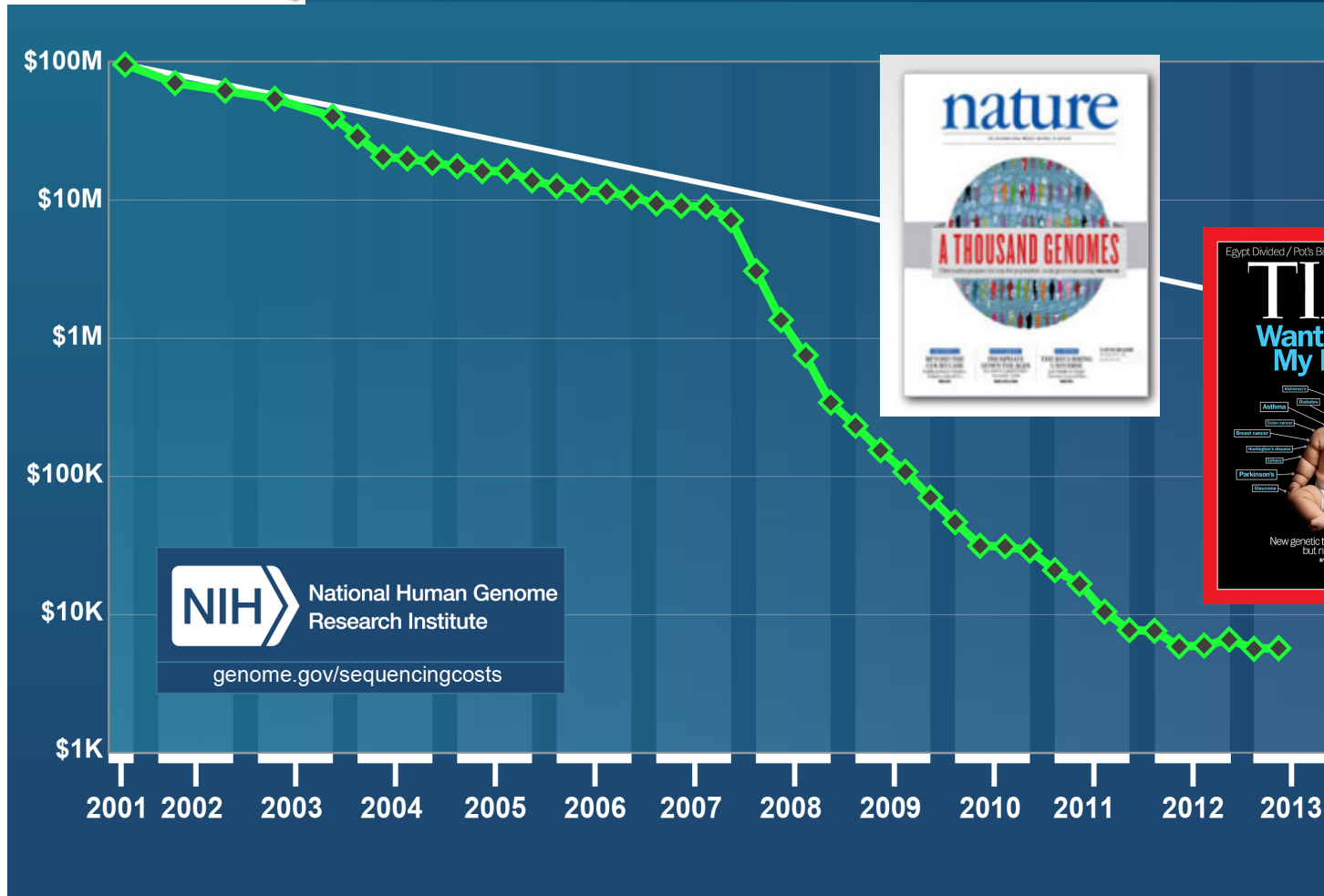
# Human genome project (1990-2003)

- Goal: sequence the 3,000,000,000 base pairs of the human genome
- Consortium of 20 laboratories, 6 countries
- 13 years, $3,000,000,000

# A flood of *omics* data
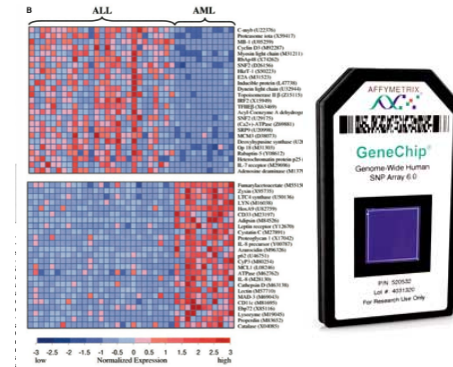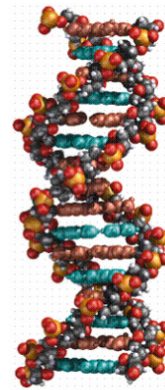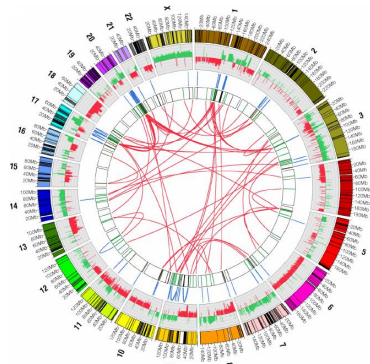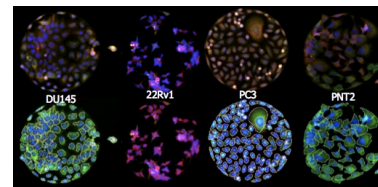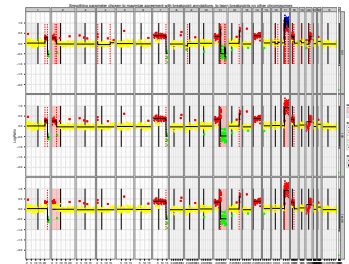


Publications

Interactome

Transcriptome
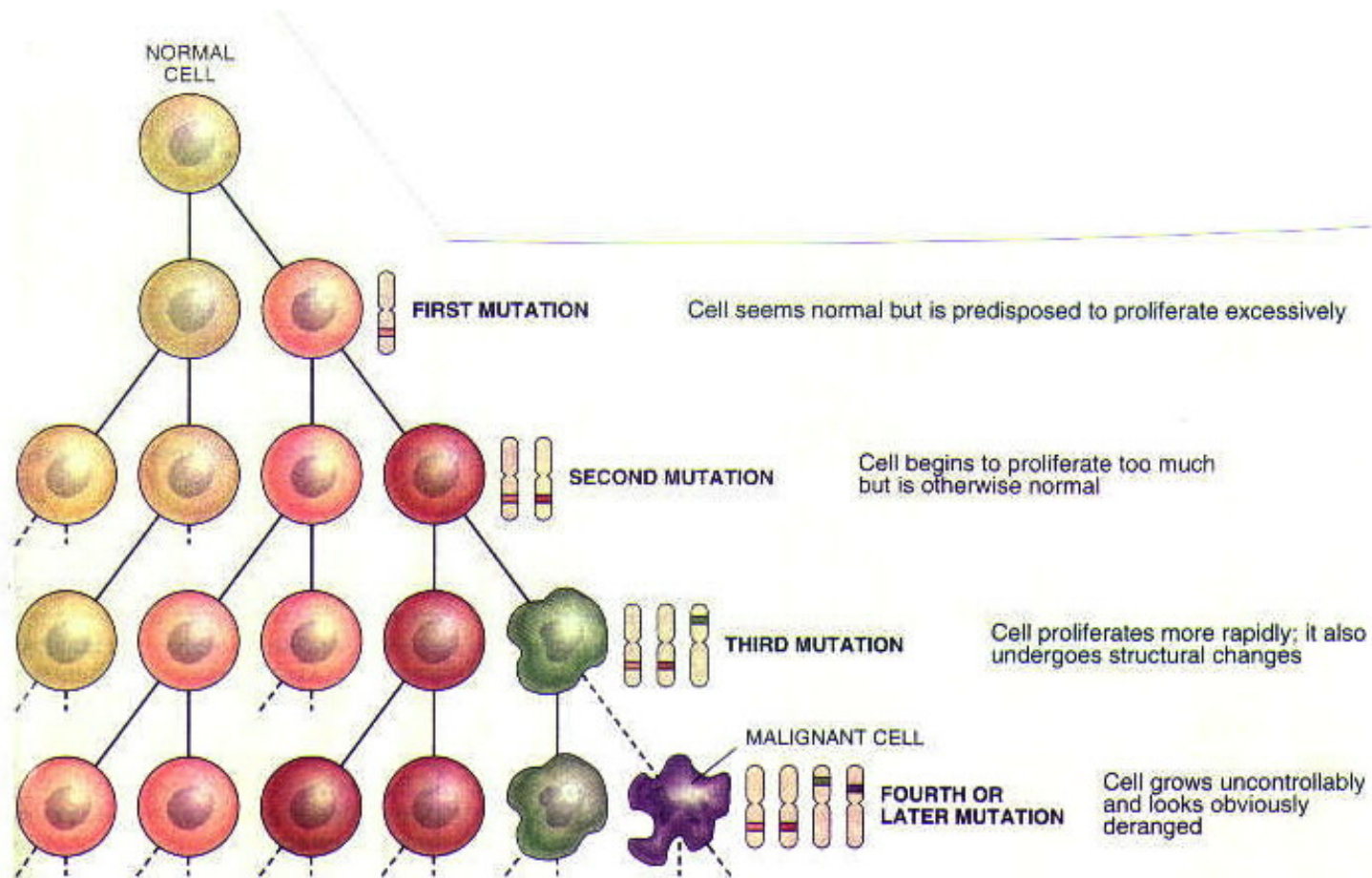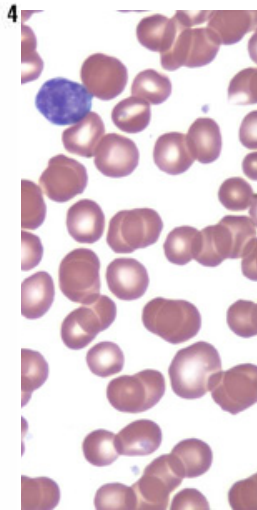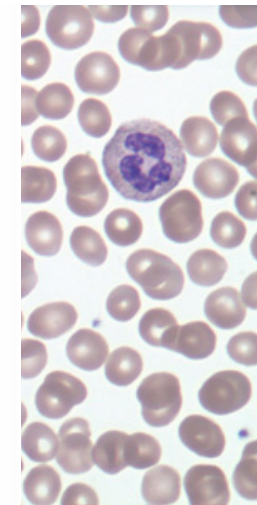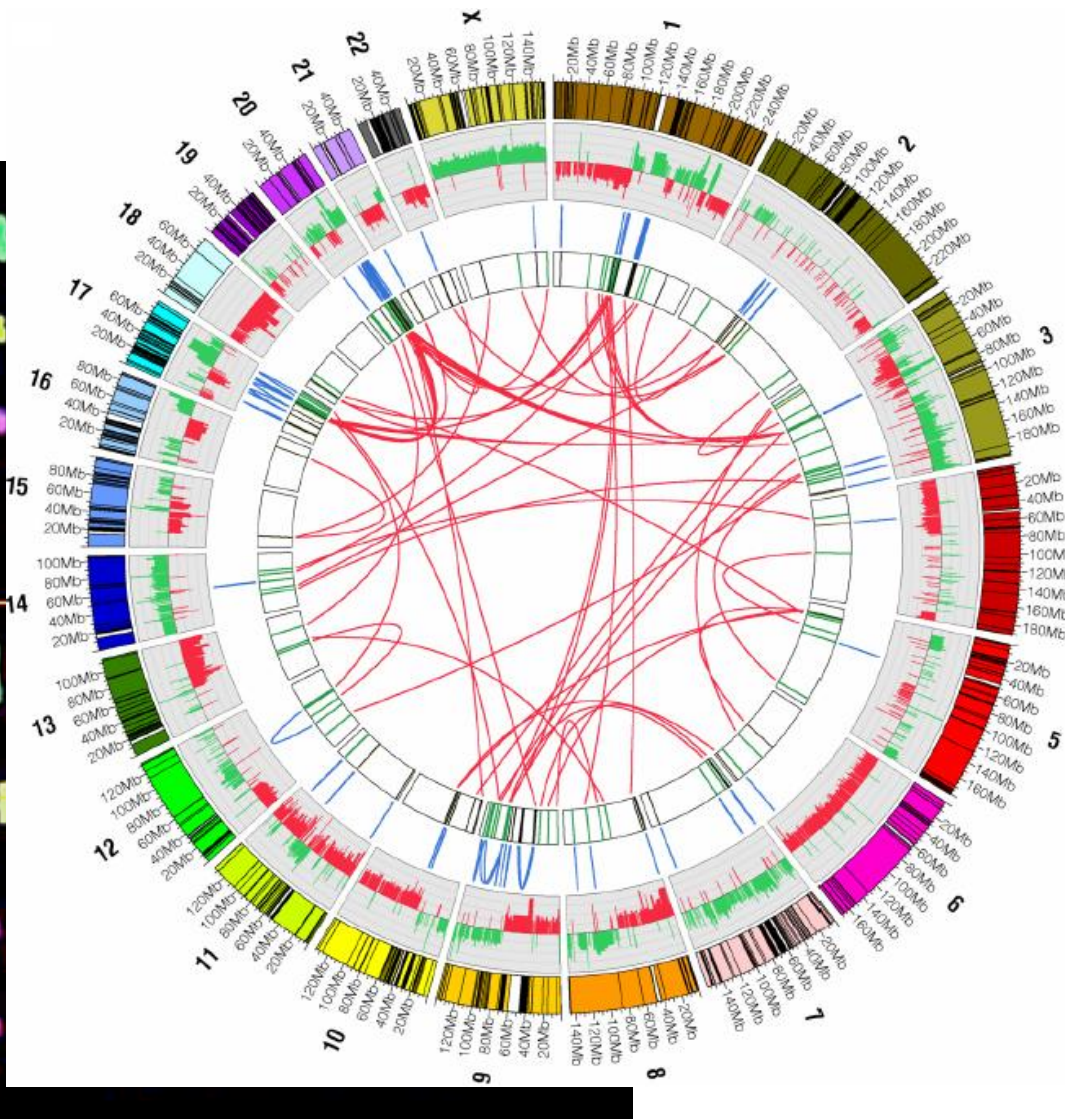
Genome

Epigenome

Mutations
Structural variations

Phenome

# All cancers are different



NORMAL CELL

**FIRST MUTATION** — Cell seems normal but is predisposed to proliferate excessively

**SECOND MUTATION** — Cell begins to proliferate too much but is otherwise normal

**THIRD MUTATION** — Cell proliferates more rapidly; it also undergoes structural changes

MALIGNANT CELL

**FOURTH OR LATER MUTATION** — Cell grows uncontrollably and looks obviously deranged

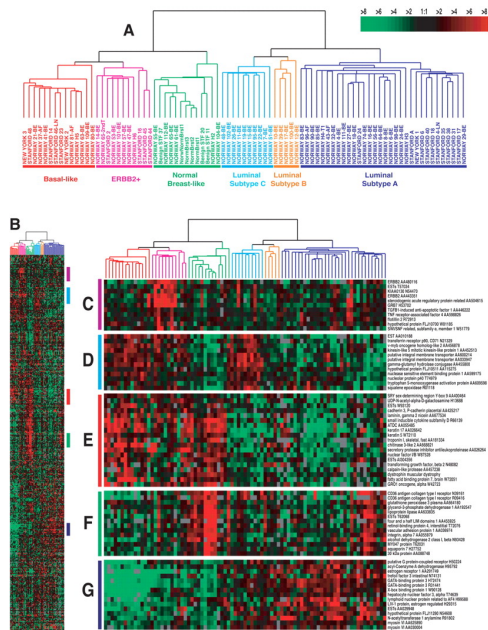# Cancer: different views

# Big data!

- http://aws.amazon.com/1000genomes/
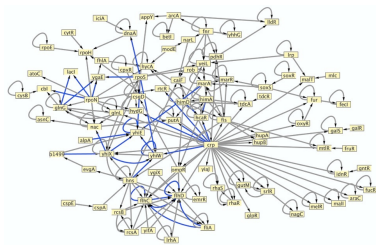
# P4 Medicine

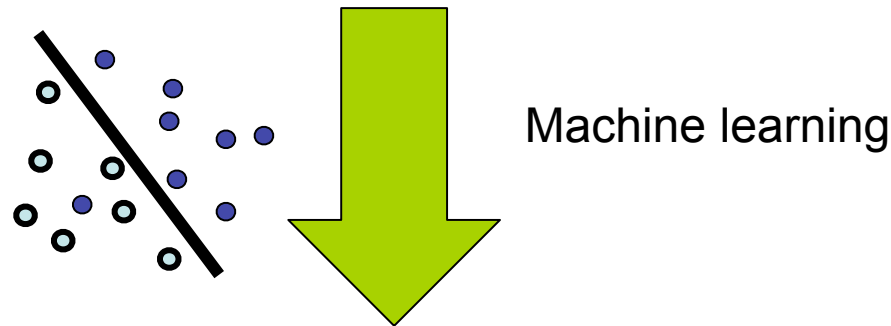● PREDICT ● PREVENT ● PERSONALIZE ● PARTICIPATE

# 23andMe

Predictive

## Opportunities

Prognosis

Diagnosis

Response to drugs

# Rationale of my team



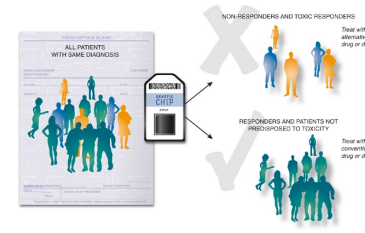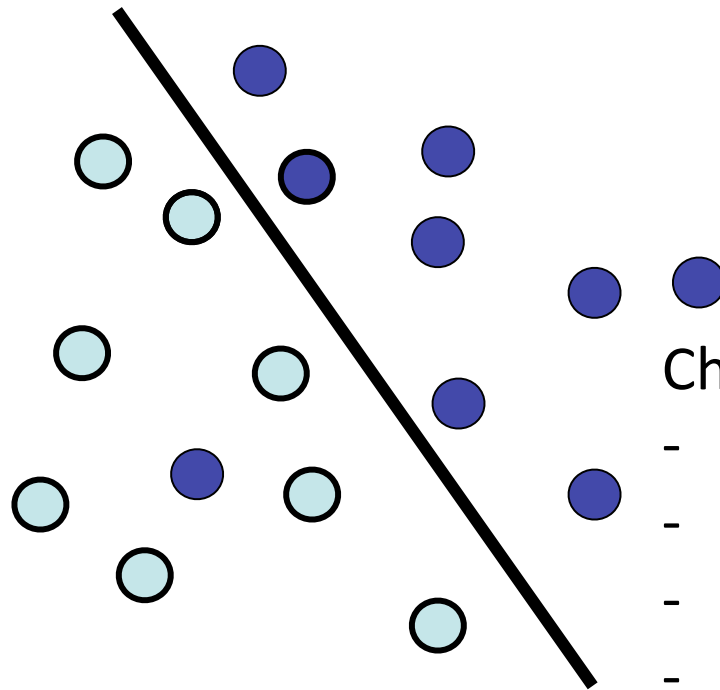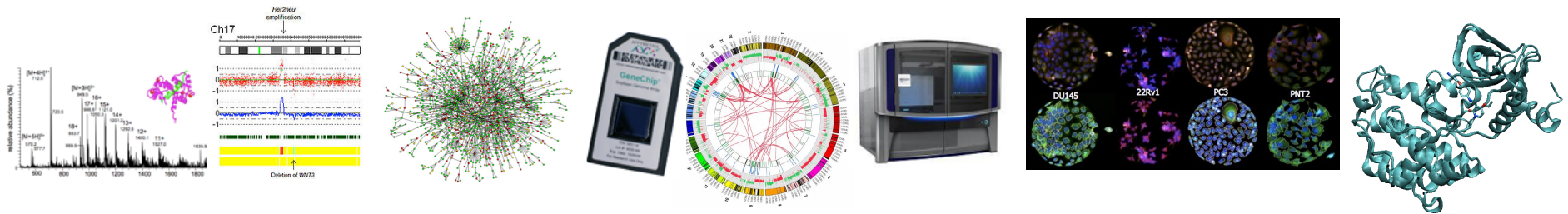Machine learning

*Mecanisms, drug targets*

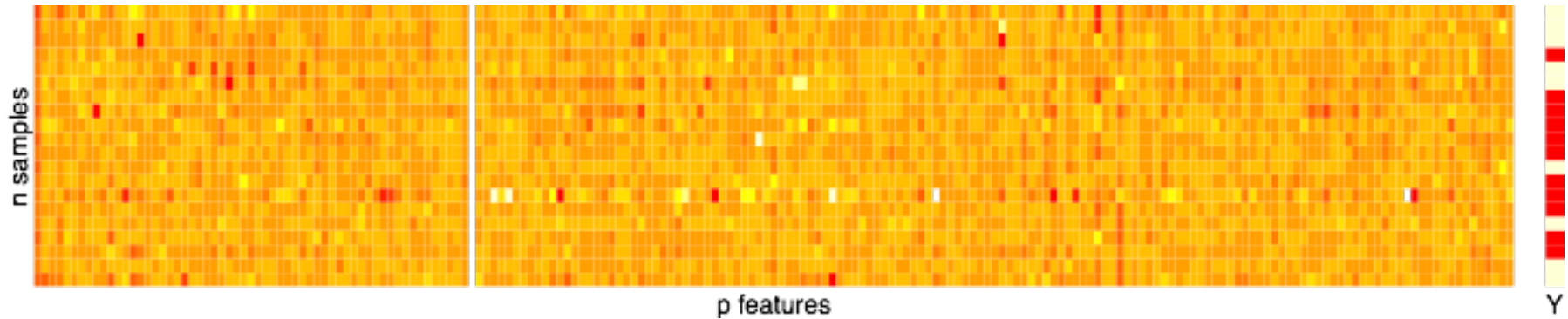*Drug design*

*Personalized medicine*

# Machine Learning?

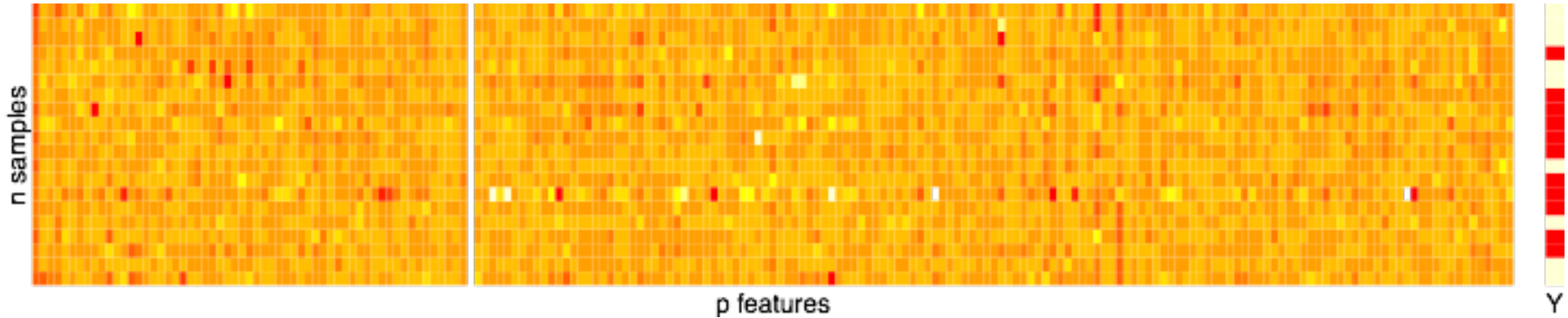

Challenges
- **High dimension**
- **Few examples**
- Structured data
- Efficient algorithms
- Interpretability
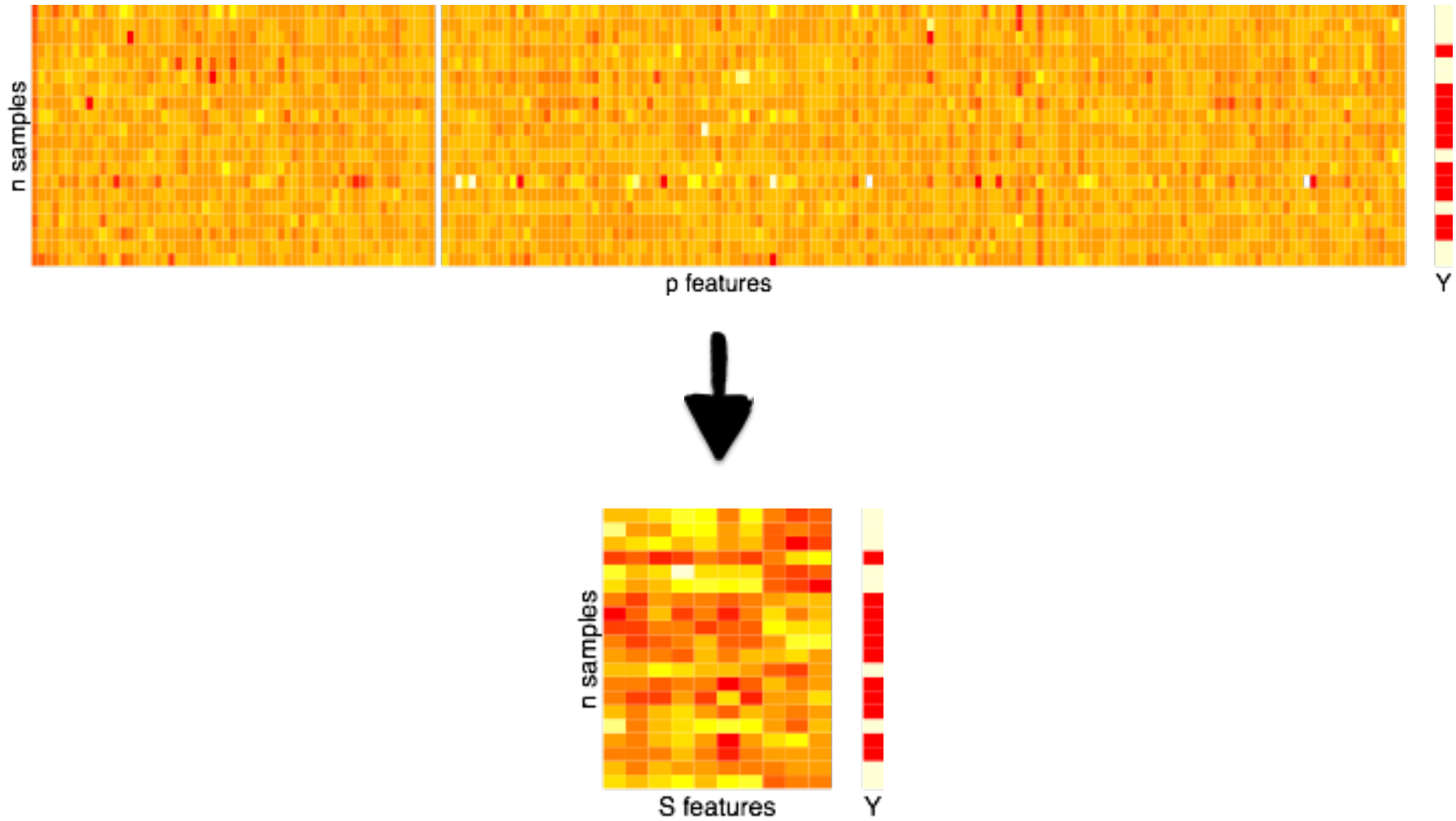
# Example: Patient stratification



n samples

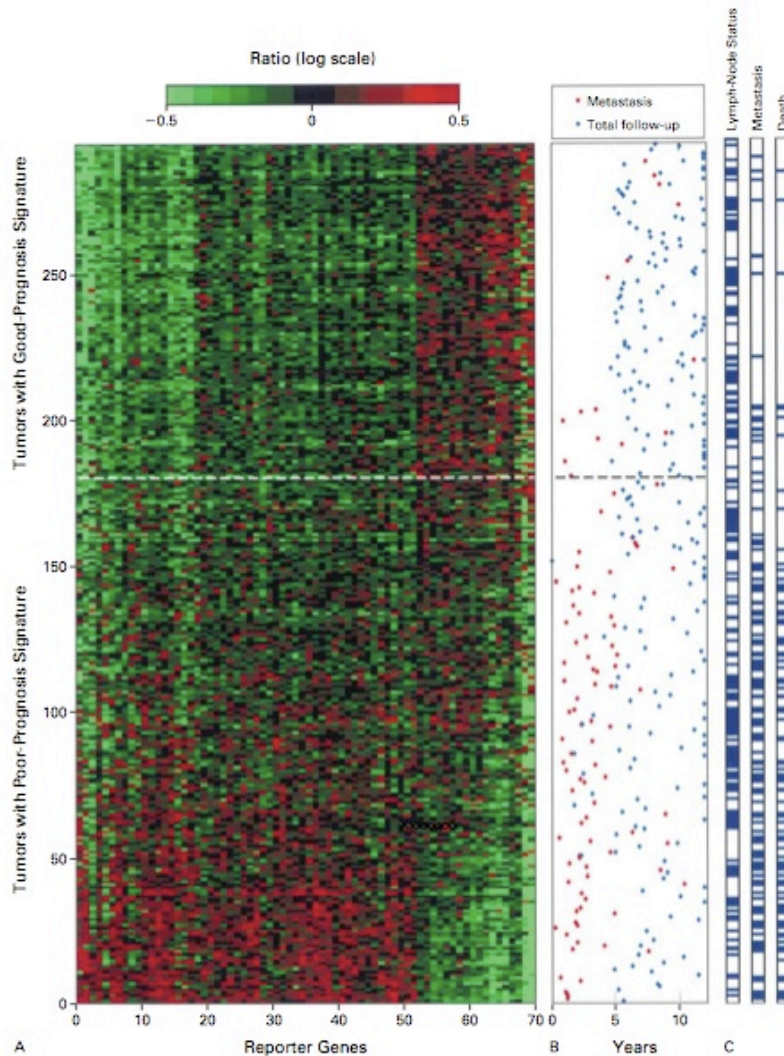p features

Y

Patients with same condition

DNA Profiling

Good responders

No Responders

Bad side effects

# Problem : n << p



**n = 1E2 ~ 1E4**
(patients)

**p = 1E4 ~ 1E7**
(genes, mutations,
copy numbers, …)

# Feature Selection
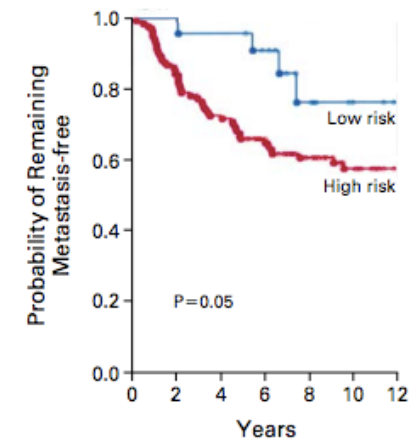
# Example:
# Breast cancer prognostic signature



*(Van de Vijver et al 2002)*

# But...

# Prior knowledge: gene network



*Can we « force » the signature to be « coherent » with a known network?*

# Example: the graph lasso

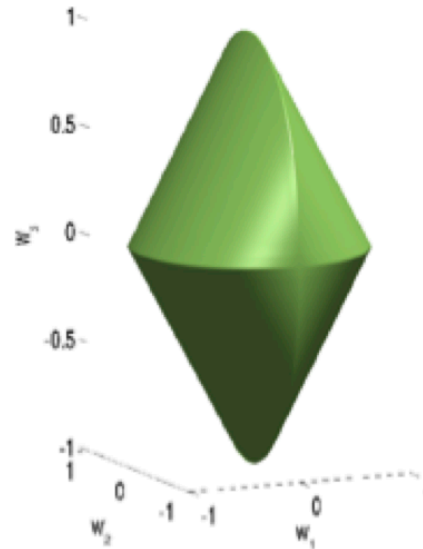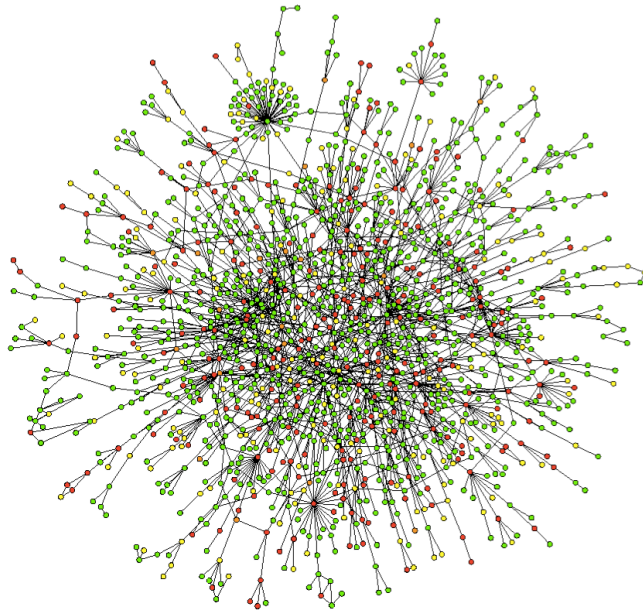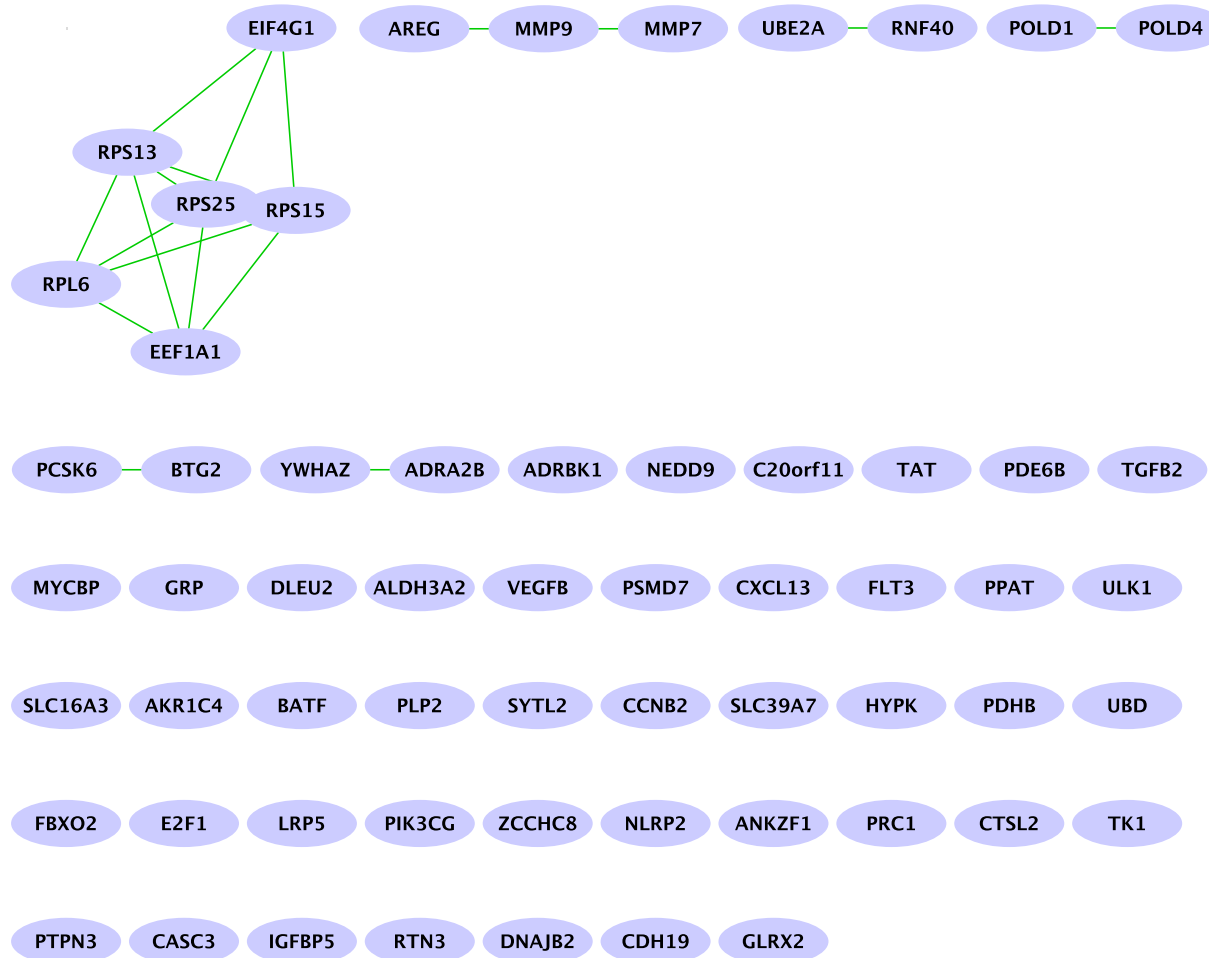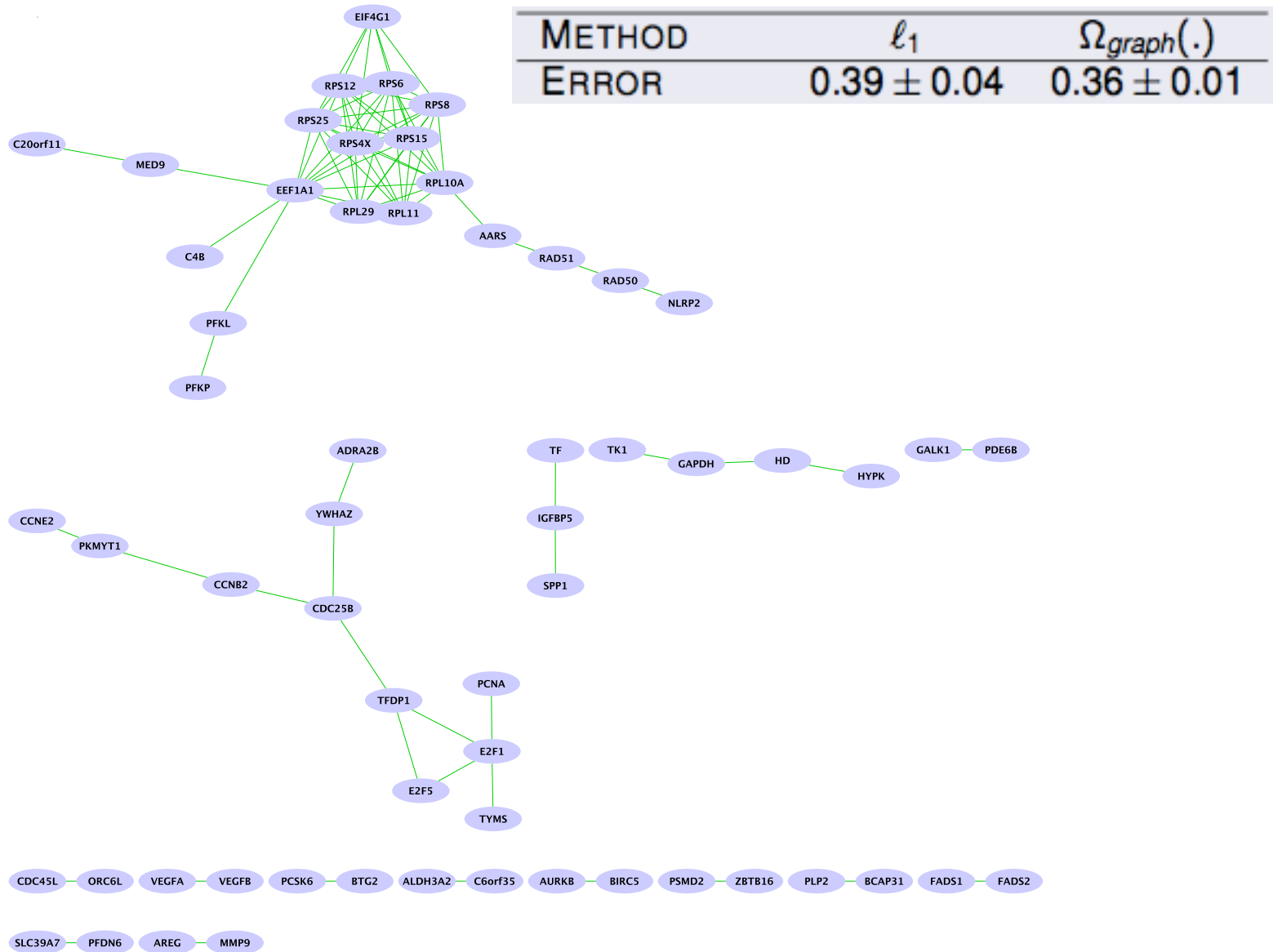- Step 1: Using the network, define a subset of « candidate » signatures



- Step 2: Among the candidates, find the best signature to explain the data
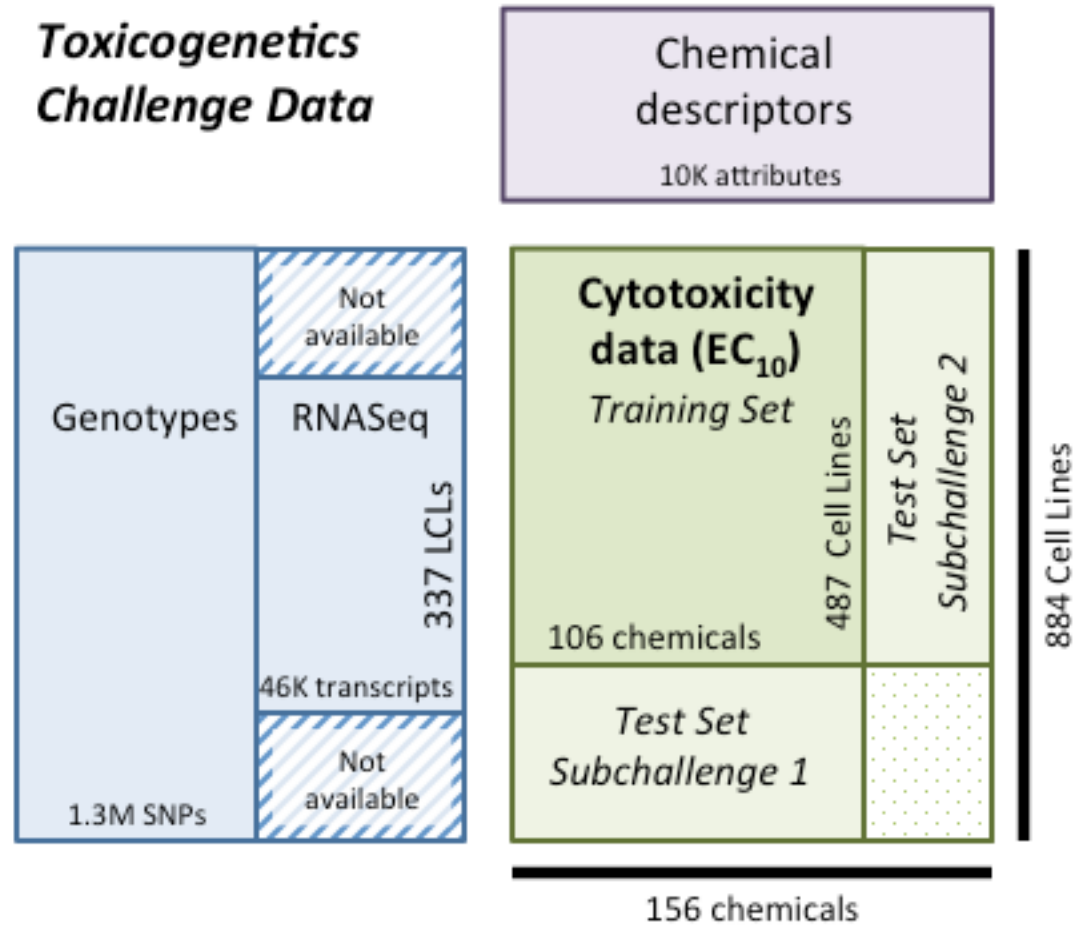
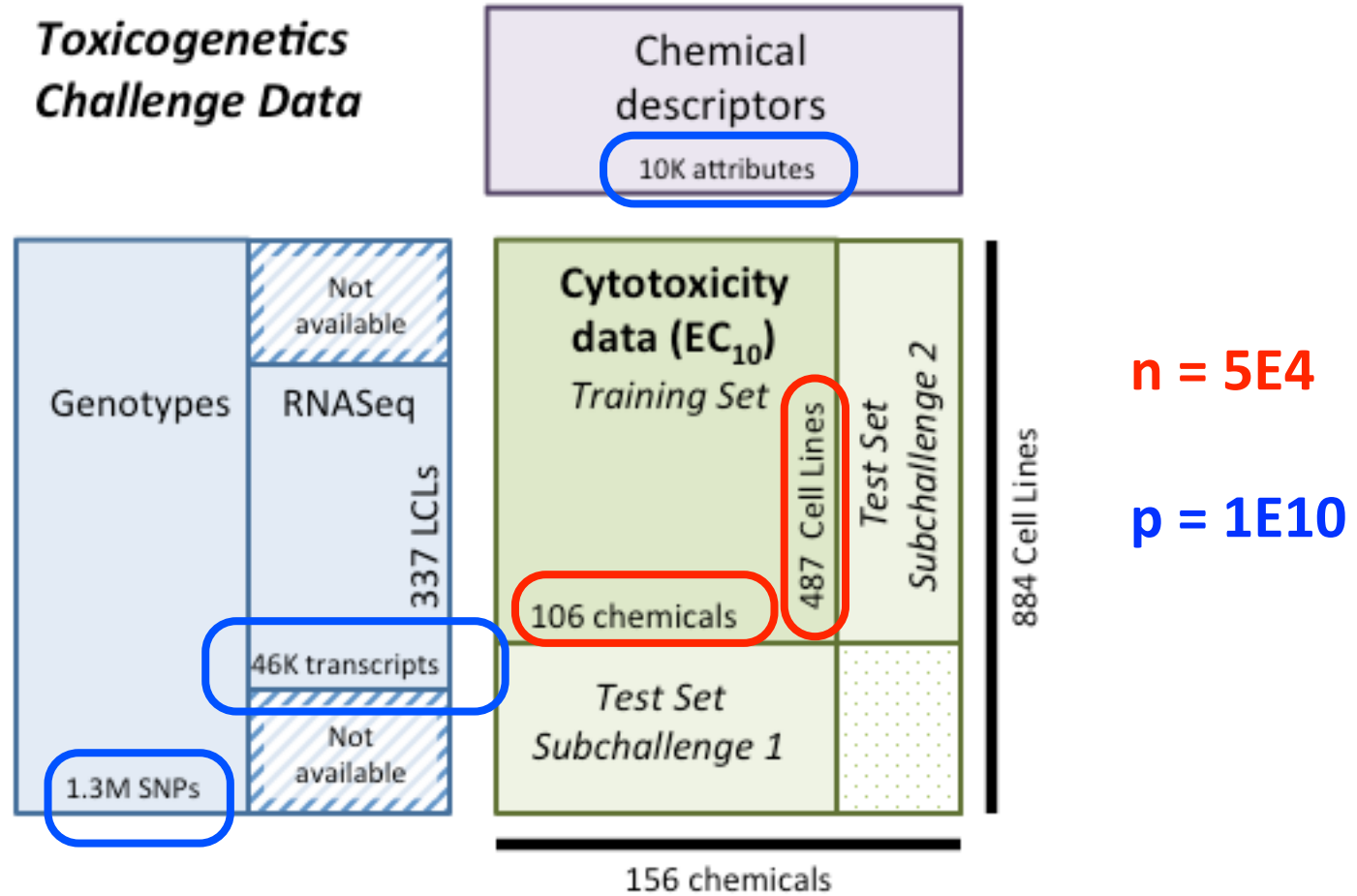*(Jacob et al 2009)*

# Classical signature

# The graph lasso signature



| METHOD | $\ell_1$ | $\Omega_{graph}(.)$ |
|---|---|---|
| ERROR | $0.39 \pm 0.04$ | $0.36 \pm 0.01$ |

# Example: Toxicogenetics / Pharmacogenomics

# Problem: n << p



**Toxicogenetics Challenge Data**

Chemical descriptors

10K attributes

Genotypes

RNASeq

Not available

337 LCLs

46K transcripts

Not available

1.3M SNPs

Cytotoxicity data ($EC_{10}$)

*Training Set*

487 Cell Lines

*Test Set Subchallenge 2*

884 Cell Lines

106 chemicals

*Test Set Subchallenge 1*

156 chemicals

n = 5E4

p = 1E10

# Crowd-sourcing initiatives

# Our approach



cell line descriptors

drug descriptors

kernelized →

Kcell

Kdrug

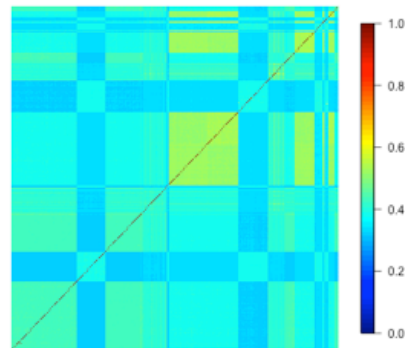kernel bilinear regression → $\hat{f}$

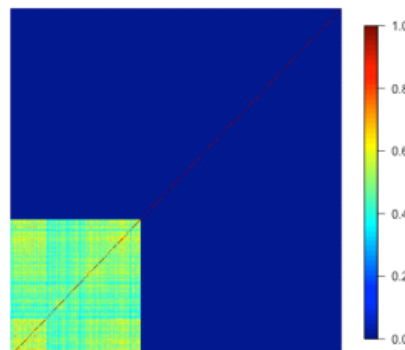# Cell line descriptors (30 kernels)
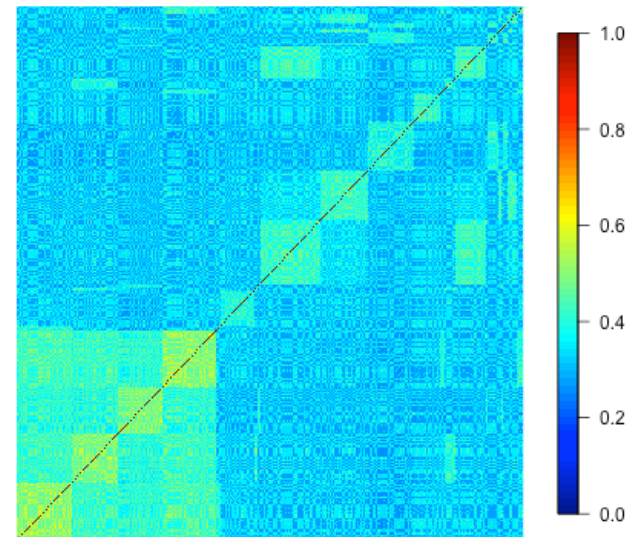


**Covariates**
. linear kernel

**SNPs**
. 10 gaussian kernels

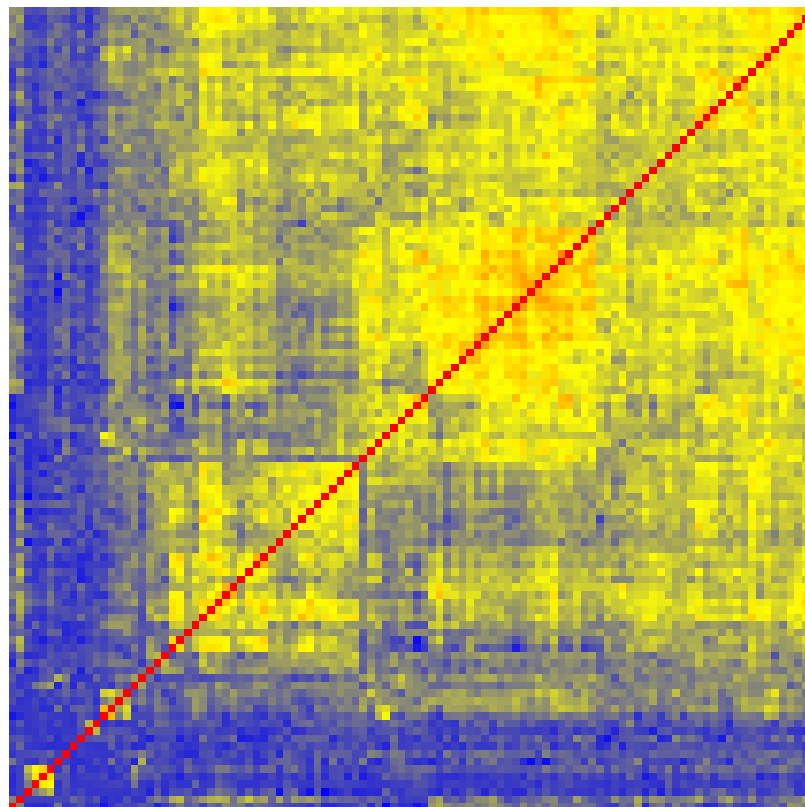**RNA-seq**
. 10 gaussian kernels
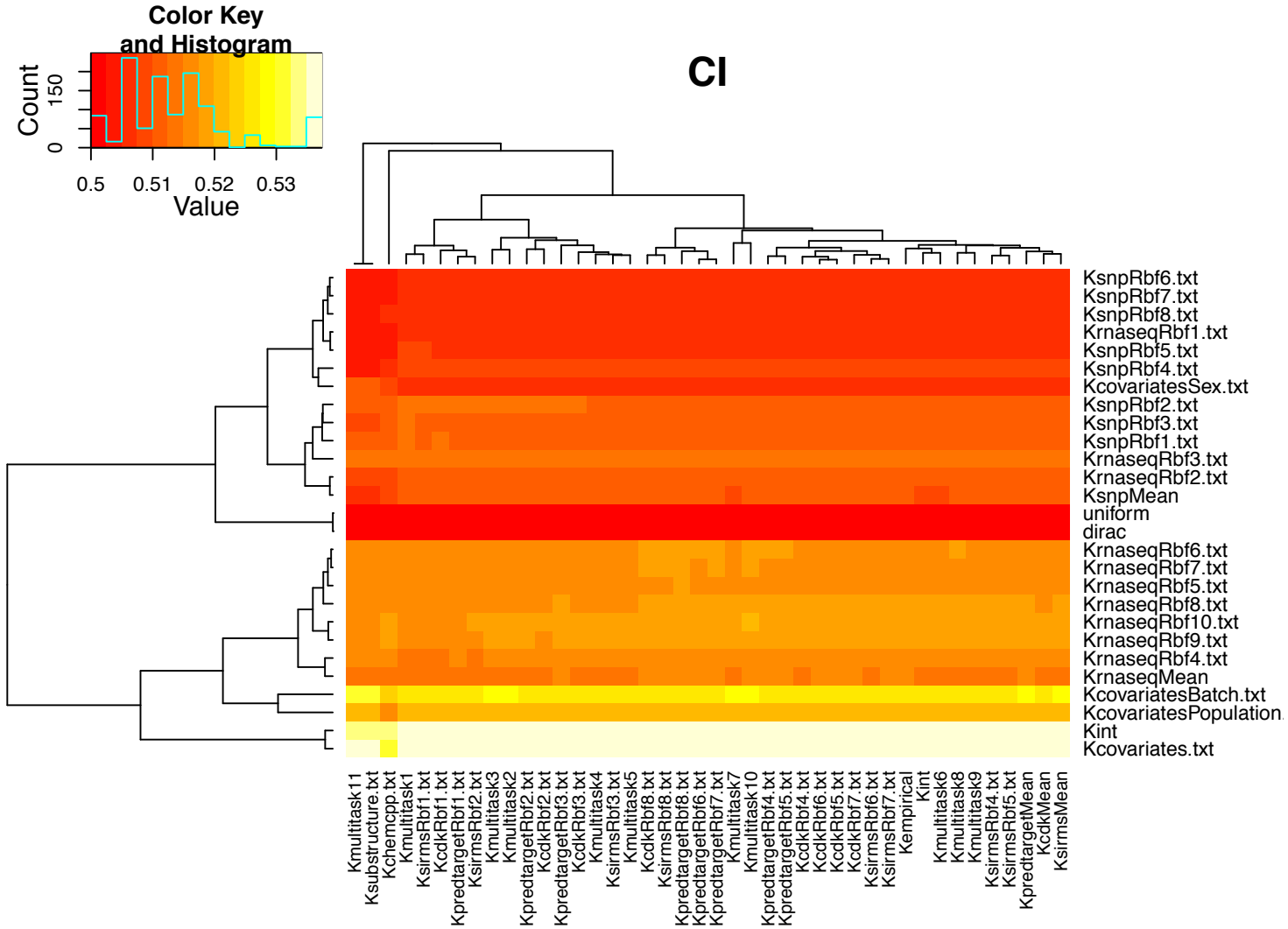
**Integrated kernel**

# Chemical descriptors (49 kernels)

- Descriptors of chemical structures

- Multitask kernels
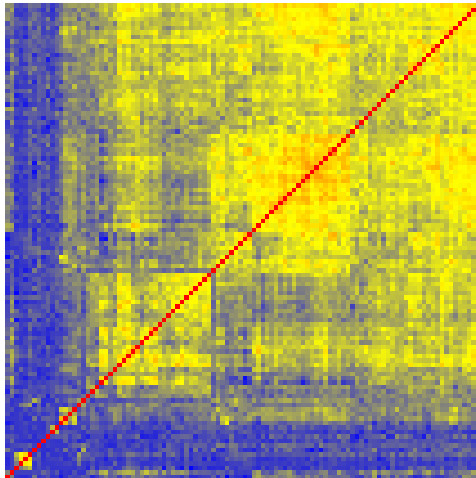
- Empirical correlation

- Integrated kernel

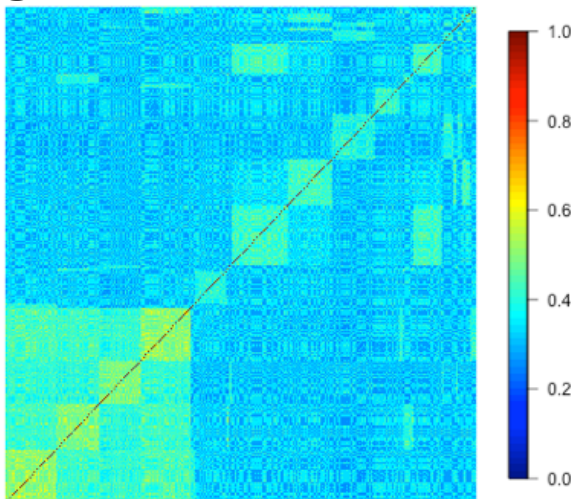# Learning occurs...

# Final submission (ranked 2ⁿᵈ)

**Empirical kernel on drugs**



**Integrated kernel on cell lines**



0.54



RECOMB/ISCB Conference on Regulatory and Systems Genomics, with DREAM Challenges 2013

TORONTO, ONTARIO
NOV 8 - 12, 2013

# Conclusion

- Lots of data due to technological progress

- **Opportunities**: precision medicine, quantitative biology

- **Challenges**: « small N », weak signal, complex systems

# Thanks!