# Computational Problems in Cancer Genomics

## Eli Upfal

## With Fabio Vandin and Ben Raphael

BROWN

# June 26, 2000 - Milestone for Humanity

**Announcing a "Milestone for Humanity--Decoding the Book of Life" at the White House Ceremony for the Completion of the Human Genome Project**

# A Milestone for Humanity?

**The New York Times**

*June 12, 2010*

**"A Decade Later, Genetic Map Yields Few New Cures"**

"Ten years after President Bill Clinton announced that the first draft of the human genome was complete, medicine has yet to see any large part of the promised result."

WHY?

# Functional Driven Sequencing - The Cancer Genome Atlas (TCGA)

Compare DNA of cancer and healthy tissue from the same patient - somatic mutation



Mutations and other genomic measurements

- Hundreds of cancer samples
- Dozens of cancer types

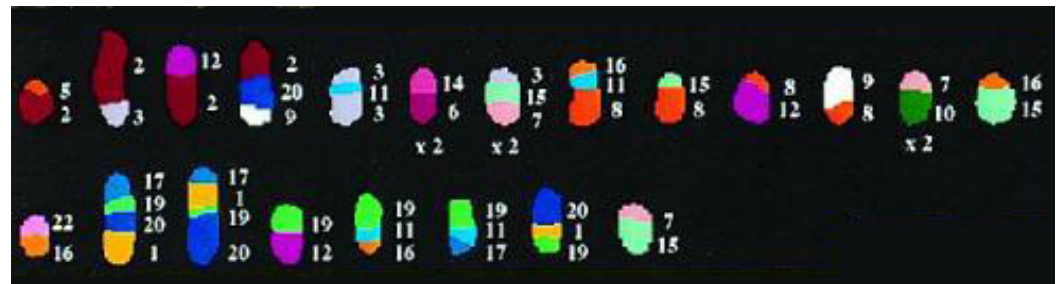**Statistical approach**: Find statistically significant *recurrent mutations*

# Cancer Genomes - Cancer is a disease of genome alterations

- Many mutations of various types
- Extensive diversity of mutations in tumors
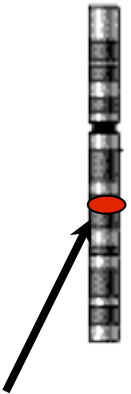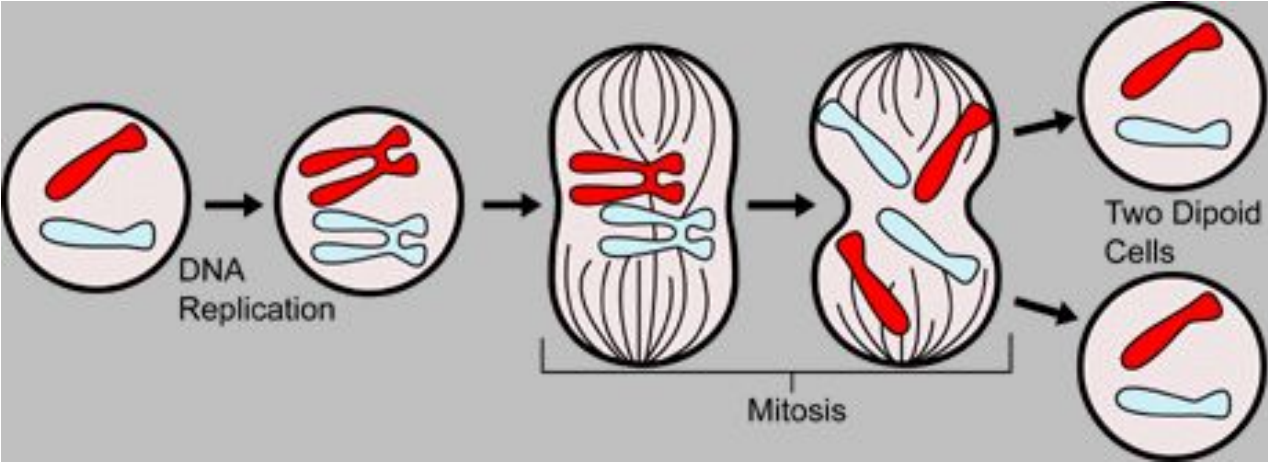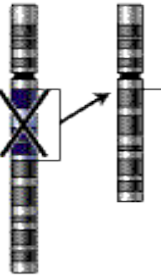  - Two tumors rarely (never?) have precisely the same set of somatic mutations



Leukemia



Breast

# DNA Replication and Mutation



Single Nucleotide variant

Copy number variants
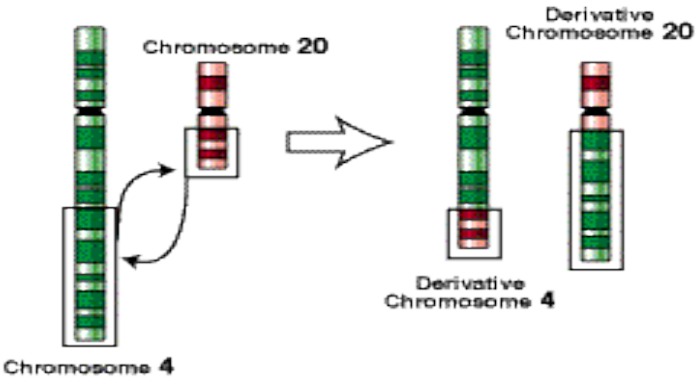
Structural variants

# Challenges in Cancer Genomics

Somatic

Human genome: ~3 billion letters

*Reads* of 30-1000 letters

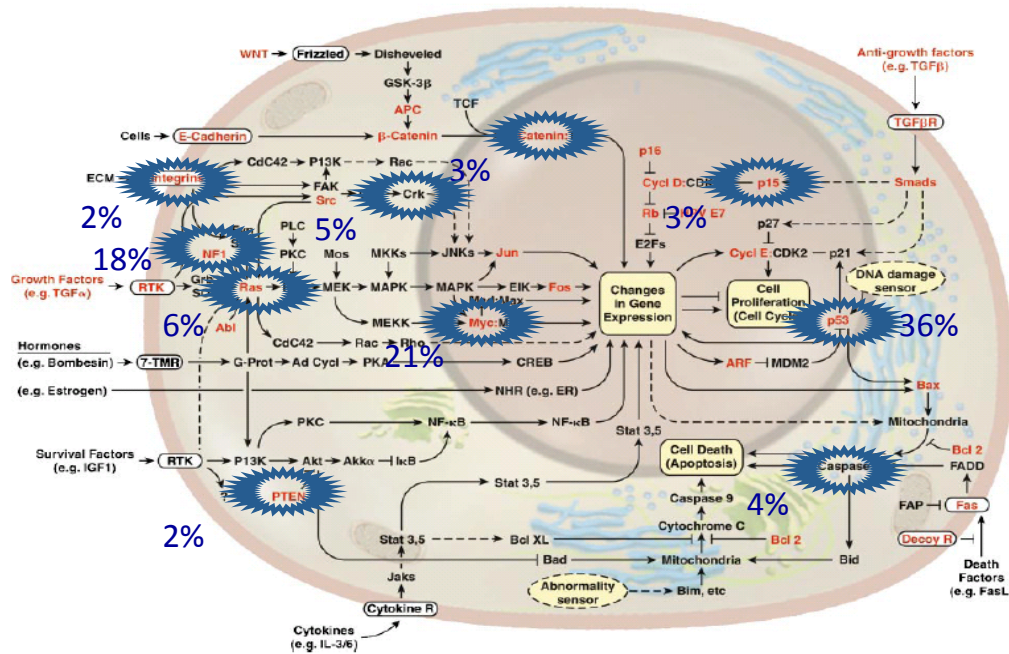**1. Measurement** of all somatic mutations

**2. Identify** functionally significant mutations

# Types of Mutations

- Driver mutations - functionally significant mutations (cause of the cancer)
- Passenger mutations – by product of the cancer process (faulty repair mechanism)

- Goal: identify the the driver mutations
- Problem: There is no small set of mutations that covers all patients

# Cancer is a disease of "pathways"
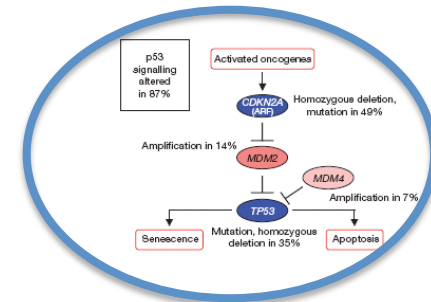


[Hanahan and Weinberg, Cell 2000]

## What pathways are altered/mutated?
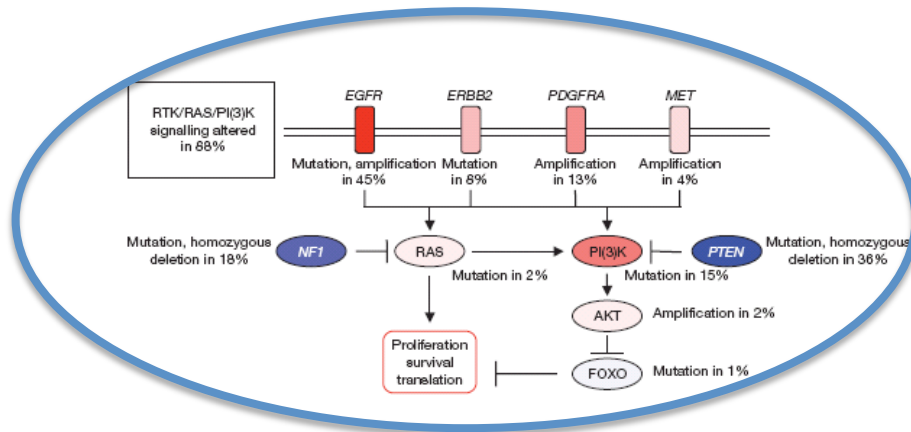
# Mutations data

- The *driver* mutations are found in pathways - sets of genes responsible for functions associated with cancer.

- *Passenger mutations are* random mutations that were not repaired because the repair mechanism in cancer cell is broken

# Finding Mutated Pathways
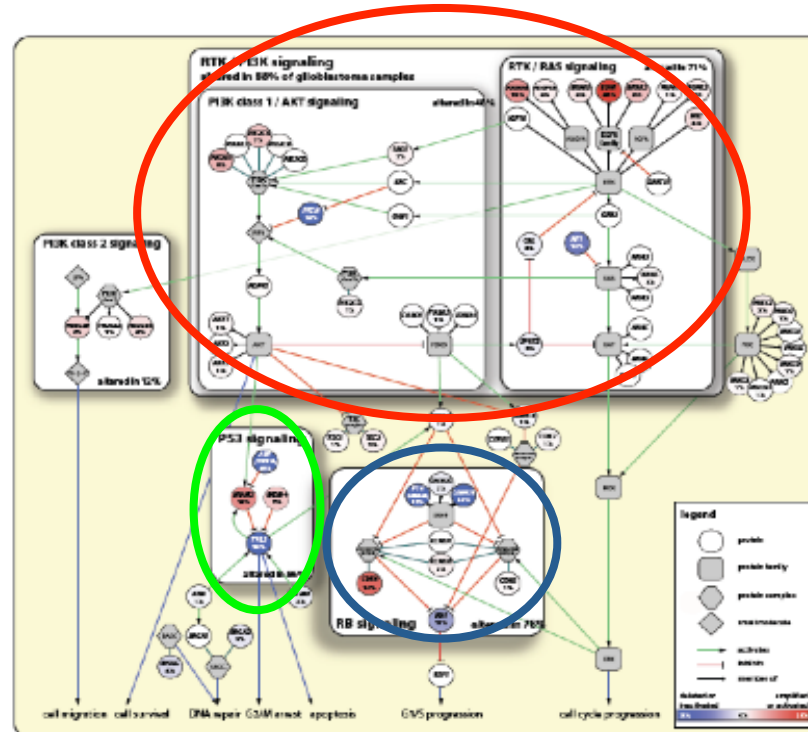
Standard practice: assess enrichment of mutations on known pathways



Only known pathways are tested!

# Finding Mutated Pathways

*Manually* constructed small network of interactions



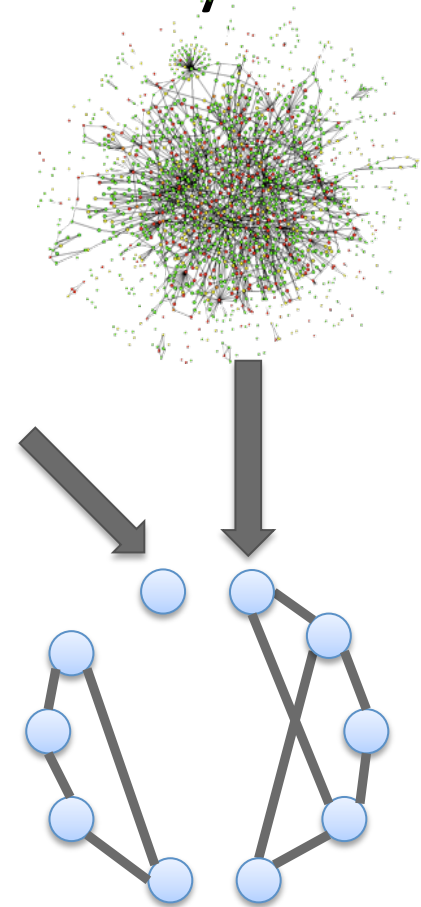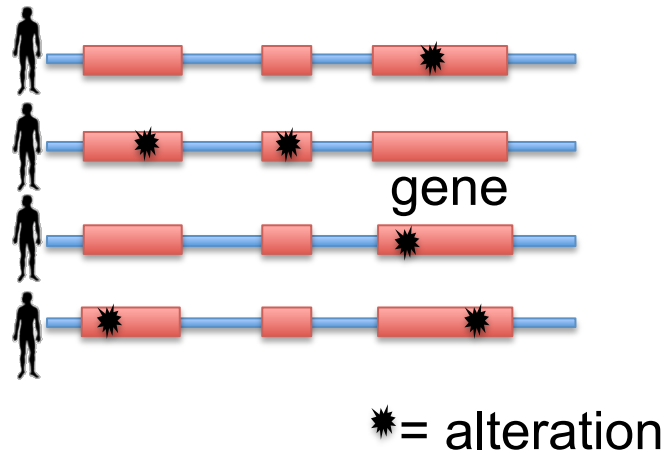[TCGA, Nature 2008]

Many genes not included!

# Network Methods

Use large interaction network to identify mutated subnewtorks

gene

$*$ = alteration

Networks are noisy!
Can we get reliable information?

# Problem

**Given**:

1. Large-scale interaction network
2. Mutation data from multiple cancer samples

gene

✳ = alteration

**Find**: Subnetworks mutated in a significant number of samples

# Problem Definition



**Given:**

1. Interaction network $G = (V, E)$

   $V$ = genes.  $E$ = interactions b/w genes

2. Binary alteration matrix

Samples

Genes



**Find**: Subnetworks mutated in a significant number of samples

- subnetwork = connected subgraph
- subnetwork *mutated* in sample if ≥ 1 gene mutated in sample

# Computational Formulation

For subnetwork $S$:

$N_S$ = number of samples in which $S$ is mutated with random alterations

$m$ = number of *observed* samples in which S mutated



**Goal**: Find $S$ such that $\Pr[N_S \geq m] < \varepsilon$ under suitable *null distribution*

# Mutated subnetworks: Naïve Method

**Find:** $S$ such that $\Pr[N_S \geq m] < \varepsilon$
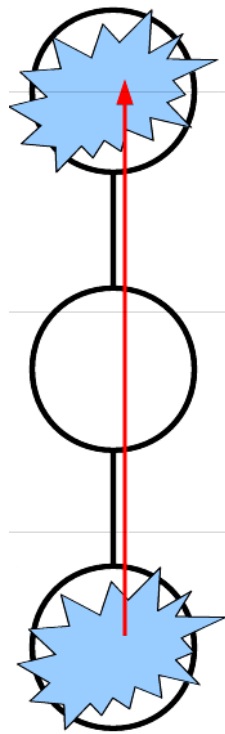under suitable null distribution

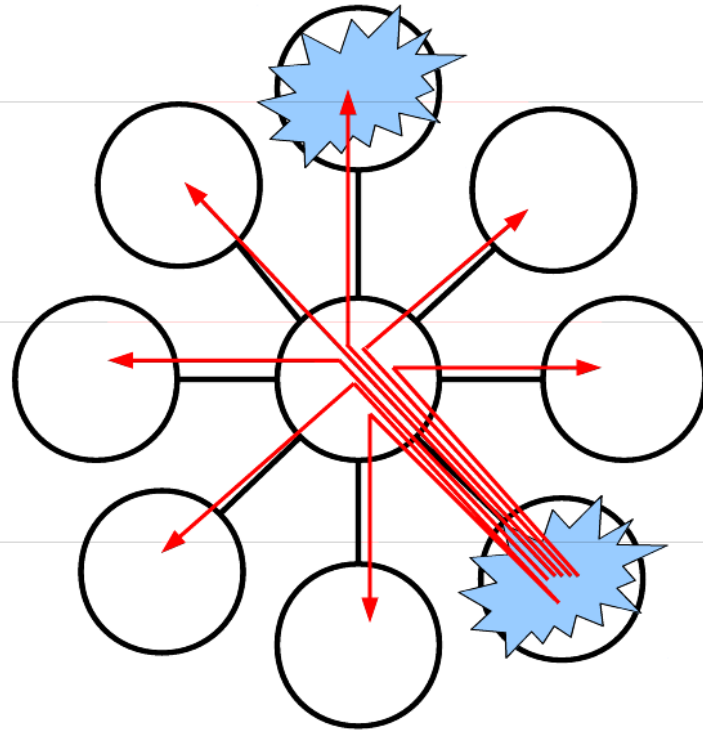**Naïve Method**: Test each $S$

**Problems**

1. Multiple hypothesis testing:
   > $10^{20}$ candidate subnetworks with < 6 genes
2. Network topology :
   TP53 has 238 neighbors in HPRD network

# (Local) Topology Matters



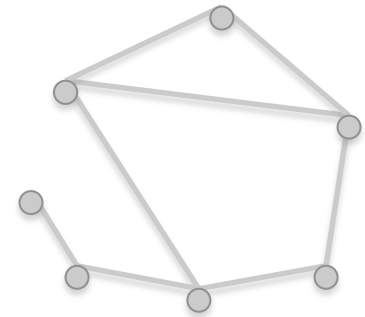Single path between mutated genes

Path between mutated genes is one of many through node.

# Our Contribution

1. *Methods* for *de novo* discovery of mutated subnetworks

   I.   Combinatorial model
   II.  Enhanced influence model


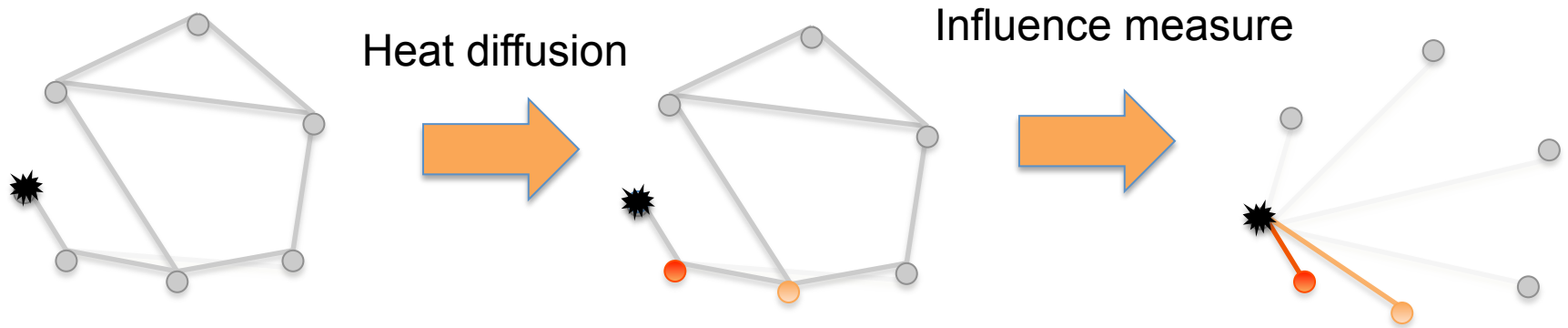2. Definition of *Influence Graph*:
   Identify subnetworks using both frequency of alteration *and* network topology


3. *Statistical tests* to assess the significance
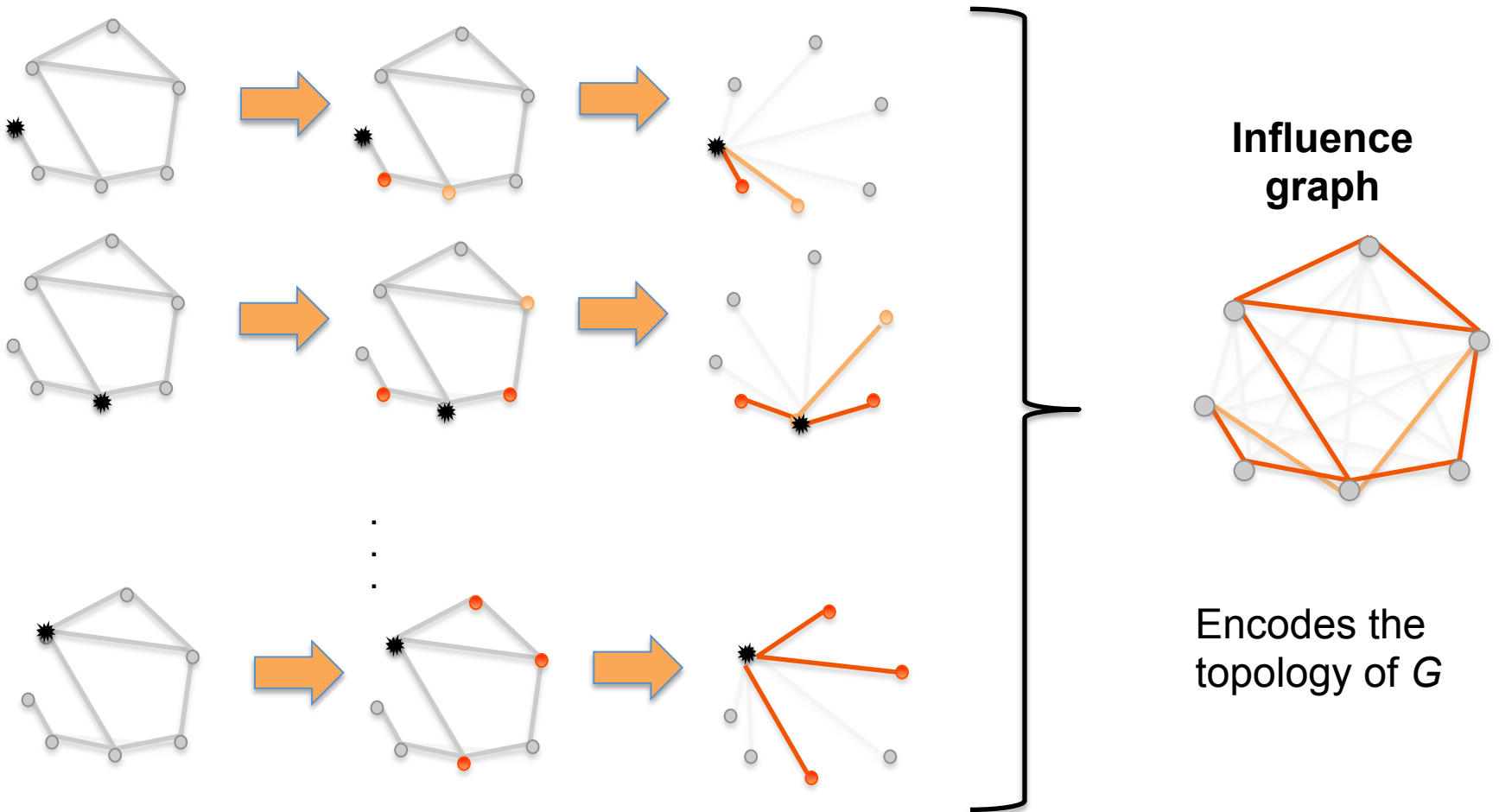
# Influence Graph

✹ alteration = unit source of heat



Heat diffusion

Influence measure

Easily derived from *Laplacian matrix* of *G*

# Influence Graph

✳ alteration = unit source of heat

**Influence graph**

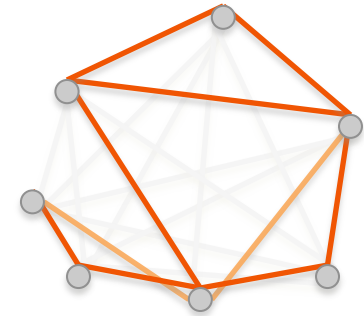Encodes the topology of *G*

# Heat equation

$f(t) = (f_1(t), \ldots, f_n(t))^\mathsf{T}$

 heat on vertices at time $t$.

$$\frac{df_i}{dt} = \sum_j a_{ij}(f_j(t) - f_i(t))$$

$df/dt = (A - D)\, f(t)$        $A = [a_{ij}] =$ **adjacency matrix** of $G$.

$f(t) = e^{-L\, t} f(0)$        $L = D - A =$ **Laplacian matrix** of $G$.

$e^{-L\, t}$ is **heat kernel** of $G$
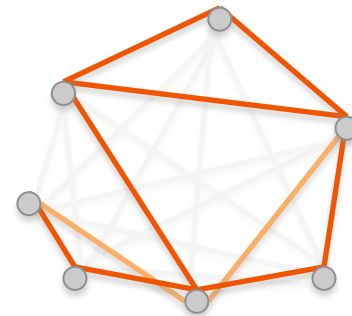
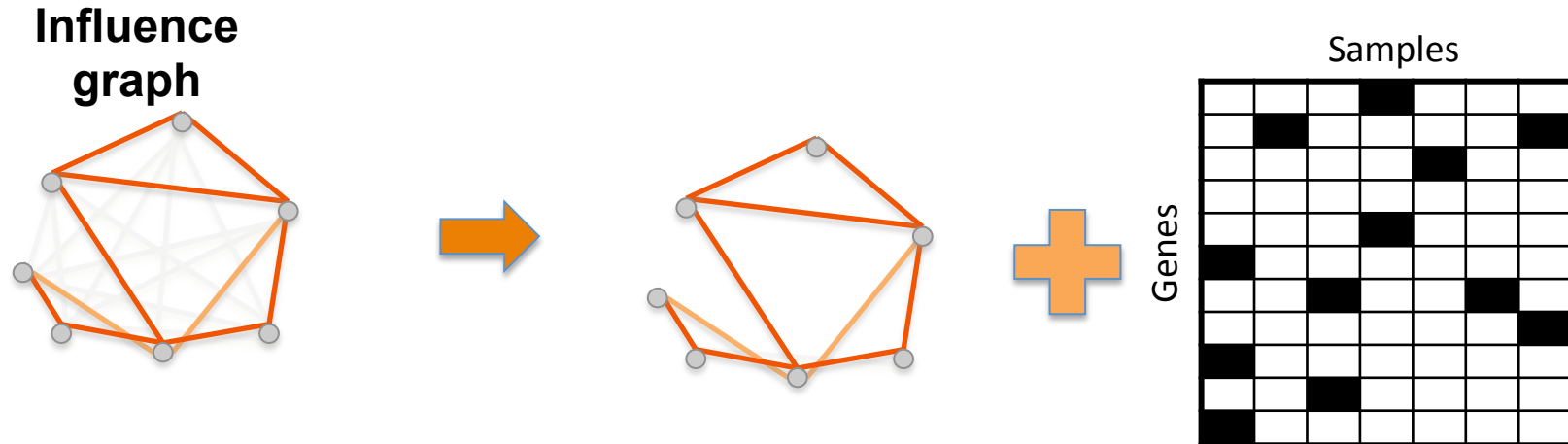# Discovering Significant Subnetworks

Two approaches:

1. Combinatorial Model

2. Enhanced Influence Model

Based on Influence Graph

Statistical tests to assess significance

# Combinatorial Model



Fix *K*: find the subnetwork with *K* genes mutated in the maximum number of samples

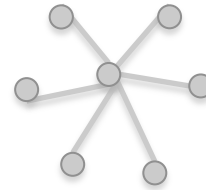**Connected maximum coverage** problem

("graph version" of maximum coverage problem – NP-Hard)

# Connected maximum coverage problem

1. **Thm.** NP-Hard for general graphs.

2. **Thm.** NP-Hard for star graphs.

3. **Thm.** 1 – 1/e approx. alg. for spider graphs

4. **Thm.** 1/(cr) approx. alg. for general graphs
   - c=(2e-1)/(e-1)
   - r= radius of the optimal solution in G

# Combinatorial Model: Statistical Test

Fix $K$: find the subnetwork with $K$ genes mutated in the maximum number of samples

testing the number of altered samples

only 1 hypothesis ➡ no multiple correction!

Limitation: inadequate representation of *heterogeneity* of cancer alterations

# Enhance Influence Model (EIM)



**Alteration Matrix**

Samples

Genes

**Influence Graph**

✺ = alterations
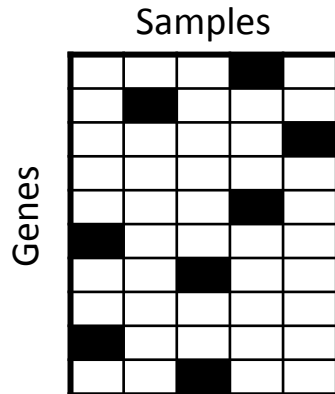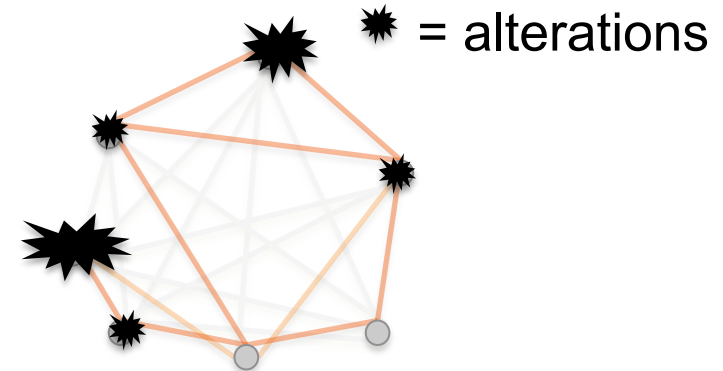
(1)

**Enhanced Influence**

Hot

Cold

**Extract "significantly hot" subnetworks**

(2)

Two-stage
multi-hypothesis
test

# EIM: Statistical test

$X_s$ = **number** of subnetworks with ≥ $s$ genes

   using "random" alteration matrix.

$H_0^s : X_s \geq \eta_s$, $s = 1, ..., N$ = # genes.

**2 subnetworks** with
2 or more genes

## Two-stage multi-hypothesis test

1.  Let $s^*$ = smallest $s$ where $H_0^s$ is rejected.

   $\Pr [ X_s \geq \eta_s ] < \alpha / N$  (Bonferroni correction)

   # hypotheses = #$s$ ≤ # measured genes.

# EIM: Statistical test



$X_s$ = **number** of subnetworks with ≥ $s$ genes using "random" alteration matrix.

$H_0^s : X_s \geq \eta_s$, $s = 1, ..., N = \#$ genes.

**2 subnetworks** with 2 or more genes

**Two-stage multi-hypothesis test**

2. Bound false discovery rate (FDR) for **list of identified** subnetworks.
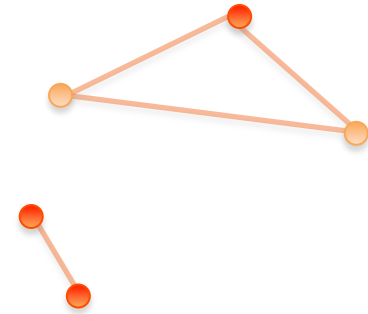
**Thm.** Fix $\beta_1, ..., \beta_N$ such that $\Sigma_i \beta_i \leq \beta$. Let $s*$ be smallest $s$ such that $\eta_s \geq E[X_s] / \beta_s$. If return all subnetworks of size ≥ $s*$ as significant, then FDR ≤ $\beta$.

# Two Stage Statistical Test

- Instead of testing the significance of the support of individual itemsets we test the significance of the <span style="color:red">number</span> of itemsets with a given support

- The null hypothesis distribution is specified by the Poisson approximation result

- Reduces the number of simultaneous tests

- More powerful test – less false negatives

[JACM 2012 - Kirsch, Mitzenmacher, Pietracaprina, Pucci, U, Vandin]

# Test I

- Define $\alpha_1, \alpha_2, \alpha_3, \ldots$ such that $\sum \alpha_i \leq$ <span style="color:red">$\alpha$</span>
- For $\mathbf{i} = \mathbf{0}, \ldots, \mathbf{log\ (s_{max} - s_{min}\ ) + 1}$
  - $\mathbf{s_i} = \mathbf{s_{min}} + \mathbf{2^i}$
  - $Q(\mathbf{k, s_i})$ = observed number of itemsets of size $\mathbf{k}$ and support $\geq \mathbf{s_i}$
  - $\mathbf{H_0(k, s_i)} = $ "$Q(\mathbf{k, s_i})$ conforms with $\mathbf{Poisson(\lambda_i)}$"

  - Reject $\mathbf{H_0(k, s_i)}$ if $\mathbf{p\text{-}value} < \alpha_i$

# Test I

- Let **s\*** be the smallest **s** such that
  $H_0(k,s)$ rejected by Test I
- With confidence level $\alpha$ the number of itemsets with support ≥ **s\*** is significant

- Some itemsets with support ≥ **s\*** could still be false positive

# Test II

- Define $\beta_1, \beta_2, \beta_3, \ldots$ such that $\Sigma \, \beta_i \le \beta$
- Reject $H_0(k, s_i)$ if:

  **p-value** $< \alpha_i$ and $Q(k, s_i) \ge \lambda_i / \beta_i$

- Let **s\*** be the minimum **s** such that $H_0(k, s)$ was rejected
- If we flag all itemsets with support $\ge$ **s\*** as significant, **FDR** $\le \beta$

# Proof

- $V_i$ = false discoveries if $H_0(k,s_i)$ first rejected
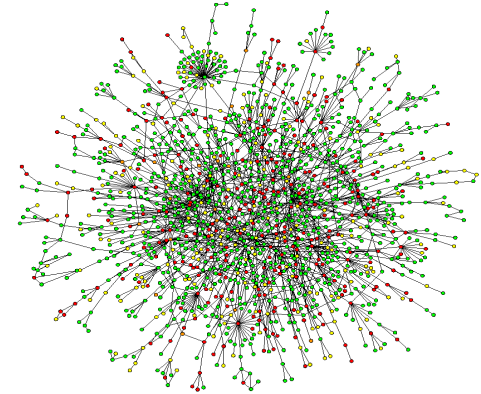- $E_i$ = "$H_0(k,s_i)$ rejected"

$$
\begin{aligned}
FDR &= \sum_{i=0}^{h-1} E\left[\frac{V_i}{Q_{k,s_i}}\right] \mathbf{Pr}(E_i, \bar{E}_{i-1}, \dots, \bar{E}_0) \\
&\leq \sum_{i=0}^{h-1} \frac{E[X_i \mid E_i\bar{E}_{i-1}, \dots, \bar{E}_0]}{\lambda_i/\beta_i} \mathbf{Pr}(E_i, \bar{E}_{i-1}, \dots, \bar{E}_0) \\
&= \sum_{i=0}^{h-1} \frac{\sum_j j\,\mathbf{Pr}(X_i = j, E_i, \bar{E}_{i-1}, \dots, \bar{E}_0)}{\lambda_i/\beta_i} \\
&\leq \sum_{i=0}^{h-1} \frac{\beta_i \lambda_i}{\lambda_i} \leq \sum_{i=0}^{h-1} \beta_i \leq \beta.
\end{aligned}
$$

# Experimental Results



## Interaction network

HPRD: 18796 nodes, 37107 edges

## Datasets

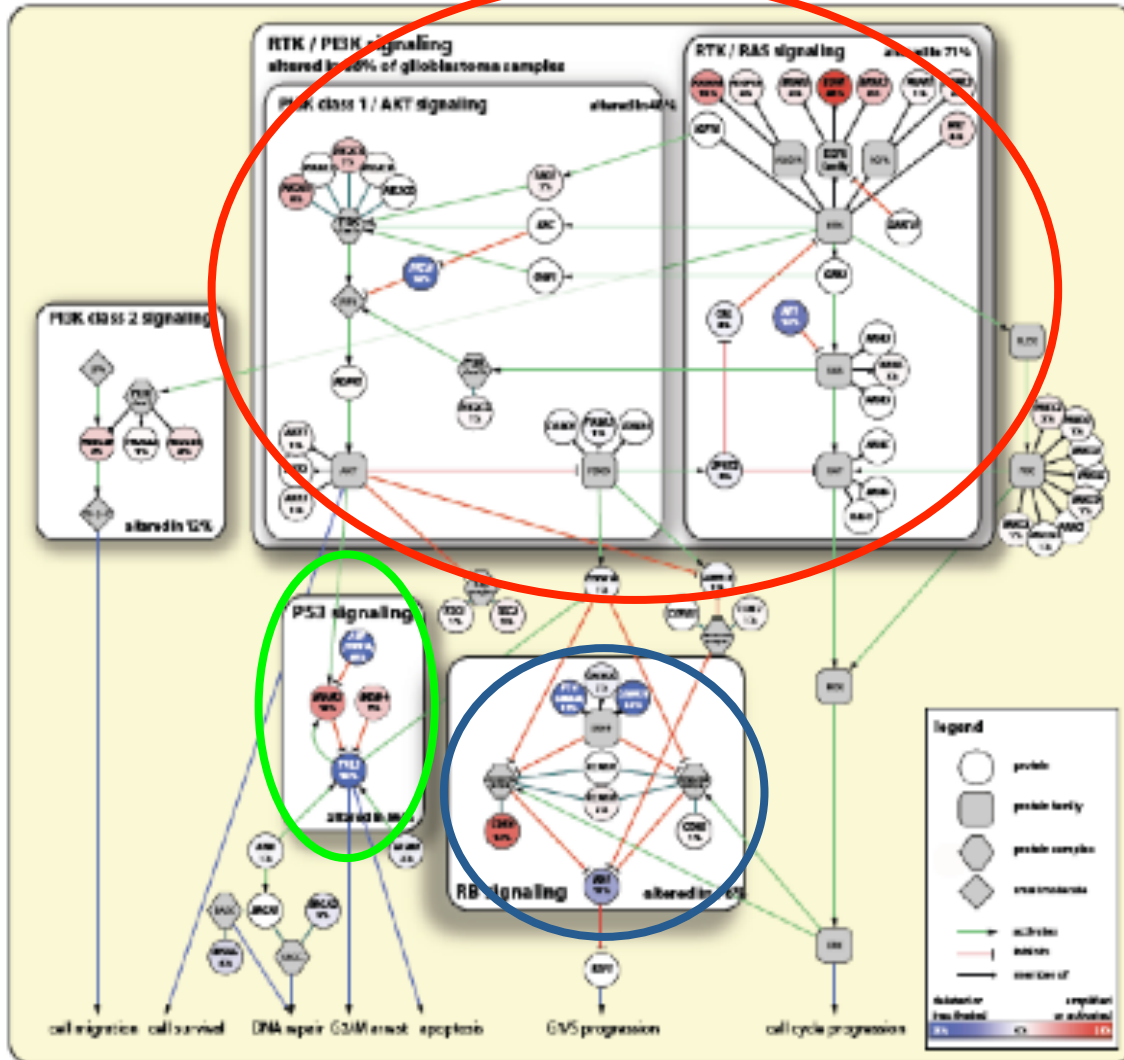1. **Glioblastoma Multiforme (GBM)** [TCGA, *Nature*, 2008]

   601 sequenced genes in 91 samples

   Array copy number data on *all* genes

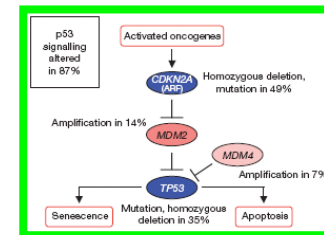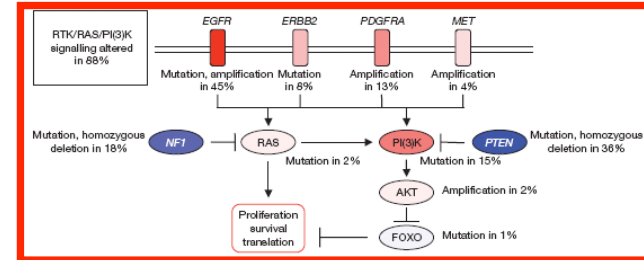2. **Lung Adenocarcinoma** [Ding et al., *Nature*, 2008]

   623 sequenced genes in 188 samples
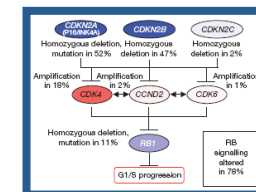
# GBM [TCGA, *Nature* 2008]



RTK/RAS/PI(3)K

p53

RB1

Manually created

Significant?

# GBM: Mutations + Copy number

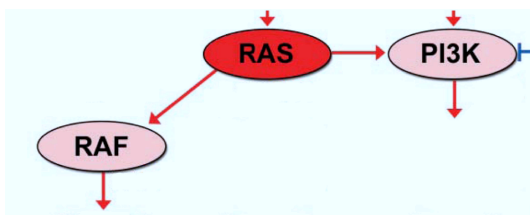| $s$ | #net $\geq s$ | $p$-val | Enrichment $p$-val | | |
| --- | --- | --- | --- | --- | --- |
| | | | RTK/RAS/PI(3)K | P53 | RB1 |
| 20 | 2 | $<10^{-2}$ | 0.69 | $2\times10^{-6}$ | $4\times10^{-8}$ |
| 26 | 1 | $5\times10^{-2}$ | $10^{-8}$ | - | - |

FDR <0.1

total enrichment for $s \geq 20$ : $p < 10^{-2}$

# Lung Adenocarcinoma
## [Ding et al., *Nature* 2008]

# Results: Lung Adenocarcinoma

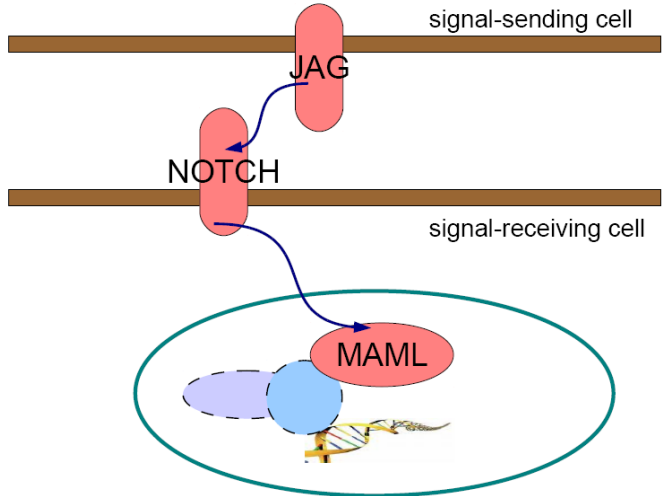|  |  |  | enrichment |
| :-: | :-: | :-: | :-: |
| $s$ | #net $\geq s$ | $p$-val | KEGG pathway/$p$-val |
| 6 | 3 | $<10^{-2}$ | Notch signaling/$2\times10^{-9}$ |
| 8 | 2 | $<10^{-2}$ | MAPK signaling/$3\times10^{-2}$ |
| 48 | 1 | $<10^{-2}$ | p53 signaling/$7\times10^{-4}$ |



FDR <0.07

total enrichment for $s \geq 6$ :
$p < 7\times10^{-9}$



Notch signaling

# Lung Adenocarcinoma: Notch



Implicated in a variety of cancers including lung
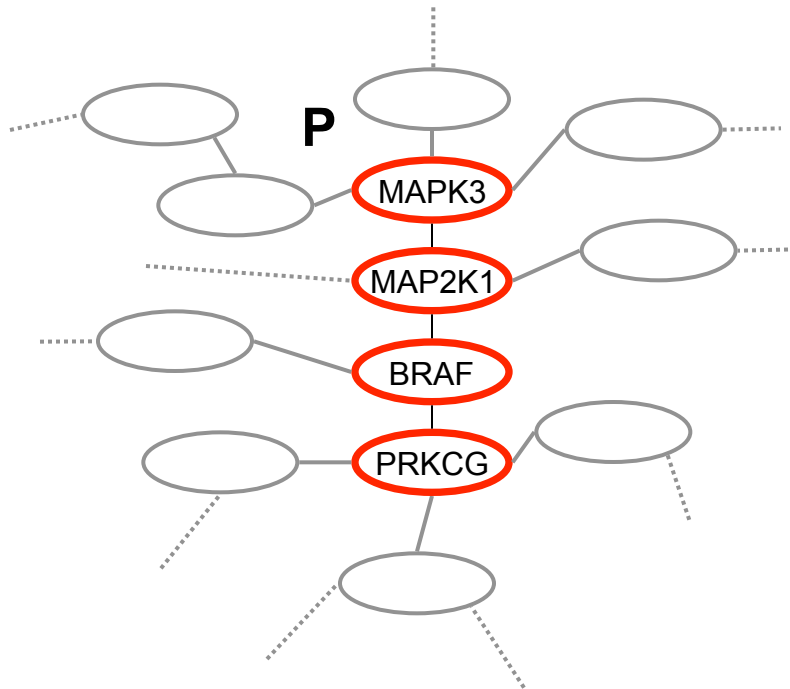
[Axelson, *Sem. Cancer Biol*. 2004, Collins et al., *Sem. Cancer Biol.* 2004]

Not reported in Ding et al. [*Nature* 2008]

| Gene | # samples |
|------|-----------|
| JAG2 | 3 |
| NOTCH2 | 1 |
| NOTCH3 | 2 |
| NOTCH4 | 3 |
| MAML1 | 3 |
| MAML2 | 1 |

# Simulated data



- **Graph**: KEGG pathway + random interactions
  - 258 genes
  - 1762 "real" edges
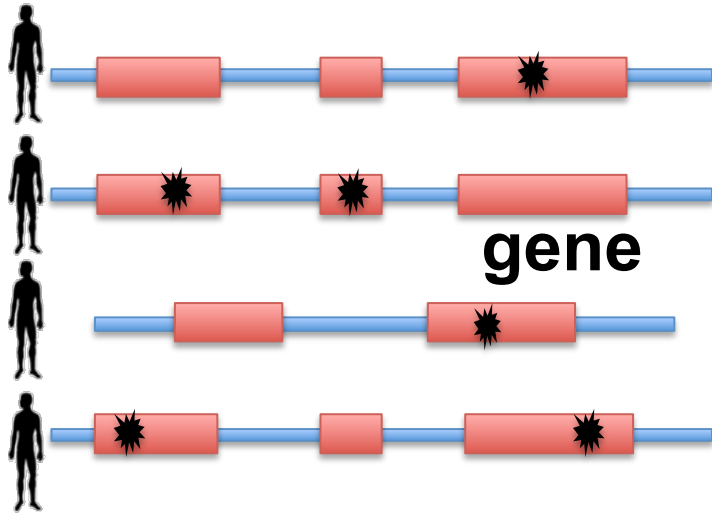  - 440 random edges
- **Alteration Matrix**
  - 30 tested genes including **P**
  - Random mutations (parameters from real data)
  - Mutations in **P** (17% of samples)

| s | #c.c.≥s | FDR | $p$-val |
|---|---------|-----|---------|
| 4 | 1 | $<10^{-2}$ | $<10^{-2}$ |

- removing mutations in **P**: nothing significant
- making BRAF hub: nothing significant

# Dendrix: Removing the Network

**Interaction network**

**Patients**

**Genes**

**Mutation matrix**

*HotNet*

Too many groups of genes to test exhaustively.
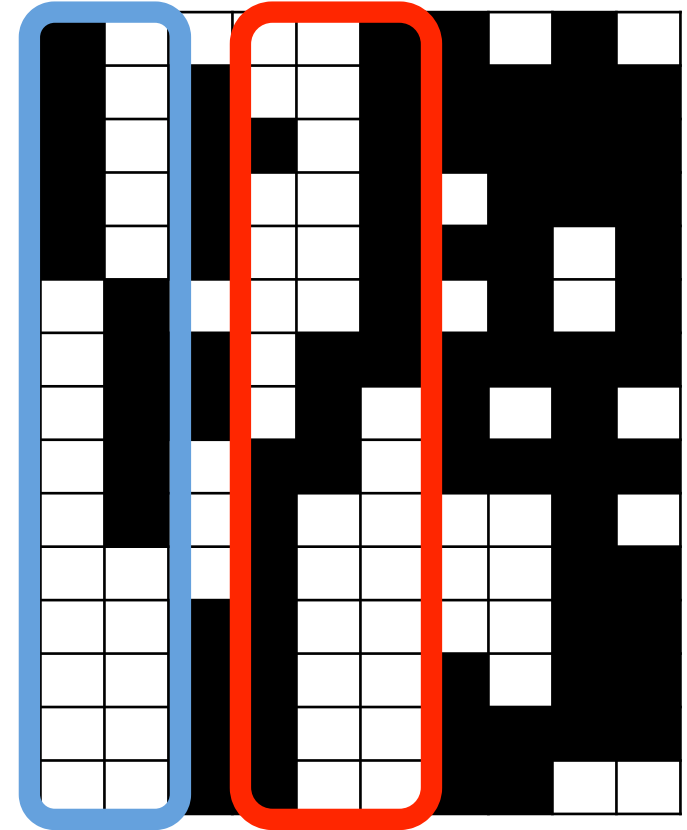
Networks are noisy.
Do we need them?

# Genomes



**: somatic mutation

# Mutation Matrix

genes

patients
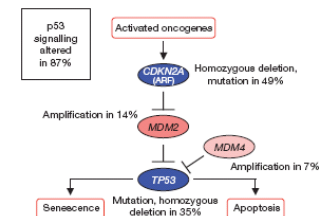
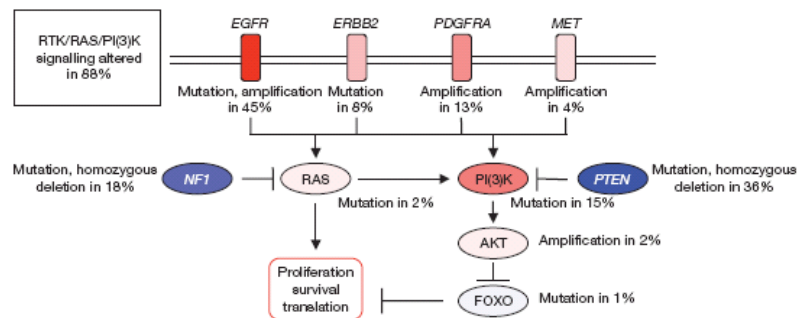Naïve: Test groups of genes
Too many hypotheses
Network reduced hypotheses. Other information?

# Pathways and Mutational Signatures

Driver mutations are rare.

→Cancer pathway has **exactly one** driver mutation (gene) per patient [REFs] [**Exclusivity**]



Most patients have mutation in pathway [**Coverage**]

# Properties of *driver* mutations

- **M** = pathway (set of genes)
- **n** = number of tested genes

- From current understanding of mutational process of cancer:

  - **Coverage**: Most samples have at least one mutation in **M**

  - **Exclusivity:** Most samples have no more than one mutation in **M**
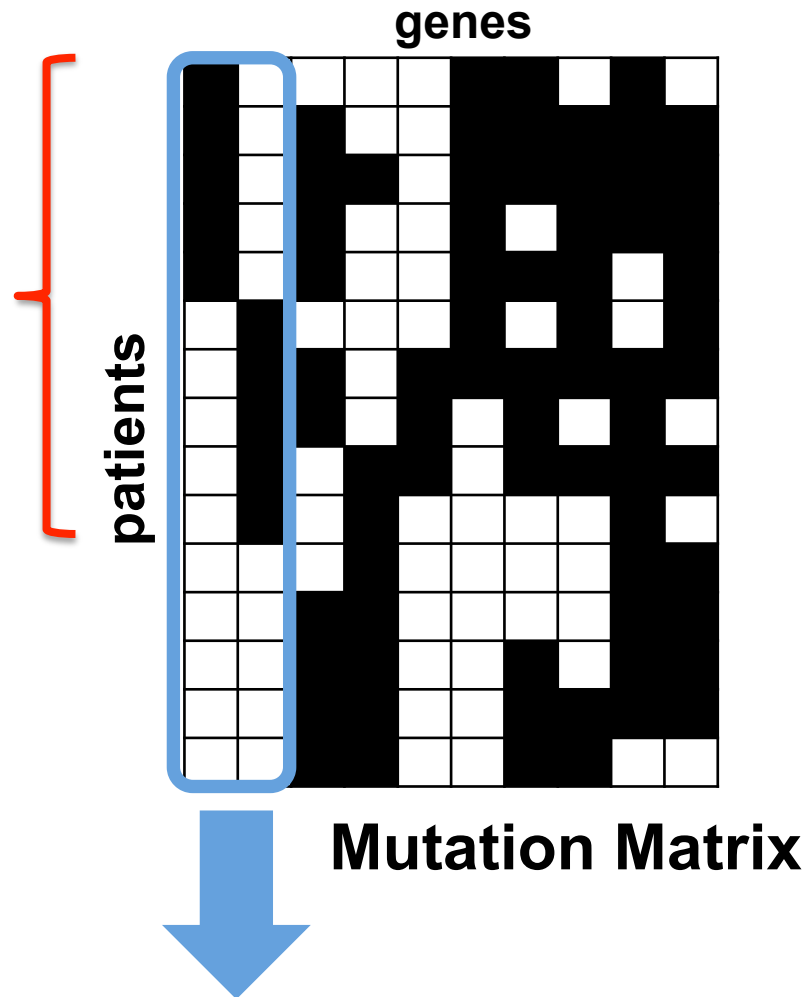
# Mutual Exclusivity and Coverage

**Coverage:**

$\Gamma(g)$ = {patients in which gene $g$ mutated}

$\Gamma(M) = U_i \, \Gamma(g_i) =$

{patients in which ≥ 1 of $\{g_1, g_2, ..., g_k\}$ is mutated}

genes

patients

**Mutation Matrix**

**Exclusive (Column) Submatrix**

# Mutual exclusivity and coverage

*Coverage:*

$\Gamma(g)$ = {patients in which gene $g$ mutated}

$\Gamma(M) = U_i\, \Gamma(g_i)$ = {patients in which ≥ 1 of $\{g_1, g_2, ..., g_k\}$ is mutated}

**Maximum Coverage Exclusive Submatrix Problem**: Given $k>0$, find the exclusive set $M$ of $k$ genes that maximizes $|\Gamma(M)|$

*Theorem* **Maximum Covering Exclusive Submatrix Problem** is NP-Hard.
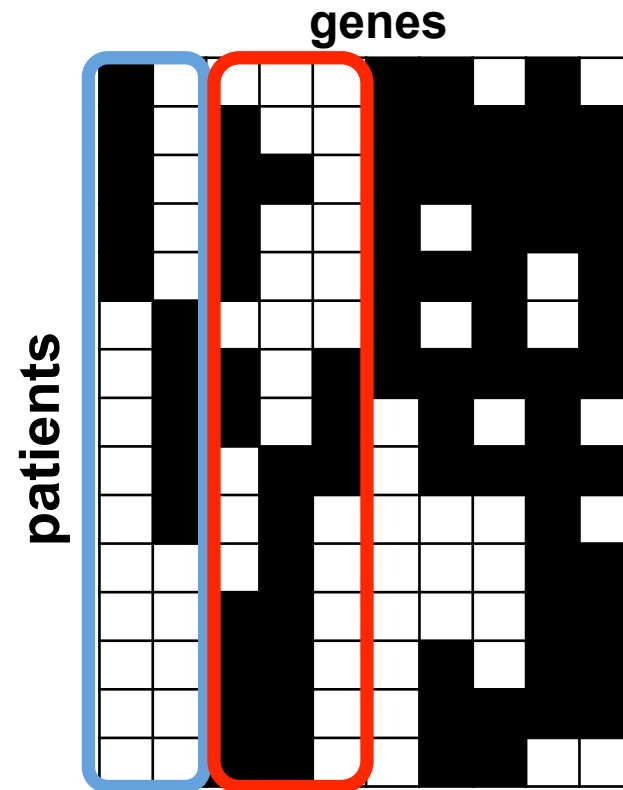
# Relaxing Constraints

For set **M** of genes:

**Coverage overlap:**

$\gamma(M) = \Sigma_i |\Gamma(g_i)| - |\Gamma(M)|$

$\gamma(M) = 0$ if and only if **M** is exclusive.

**Goal**: $|\Gamma(M)|$ large and $\gamma(M)$ small.



**"Approximately exclusive"**, high coverage submatrix

# Approximate Exclusivity

**Goal**: $\Gamma(M)$ large and $\gamma(M)$ small.

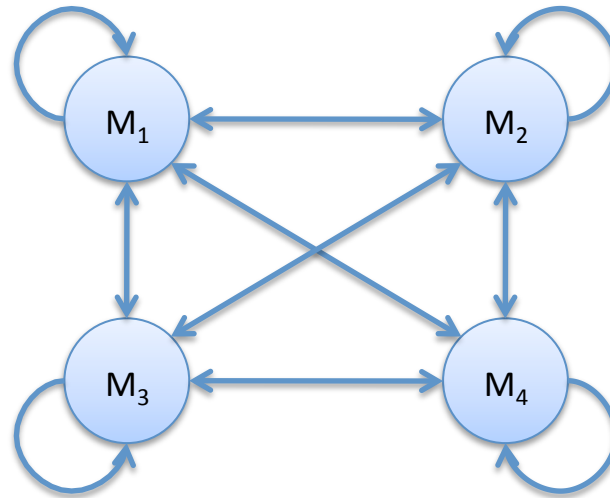Weight: $W(M) = |\Gamma(M)| - \gamma(M) = 2\,|\Gamma(M)| - \Sigma_i |\Gamma(g_i)|$

**Maximum Weight Submatrix Problem**:  Given $k>0$, find the set $M$ of $k$ genes that maximizes $W(M)$

*Thm.* **Maximum Weight Submatrix Problem** is NP-Hard.

# Markov Chain Monte Carlo

Sample gene sets $|M| = k$ according to $W(M)$

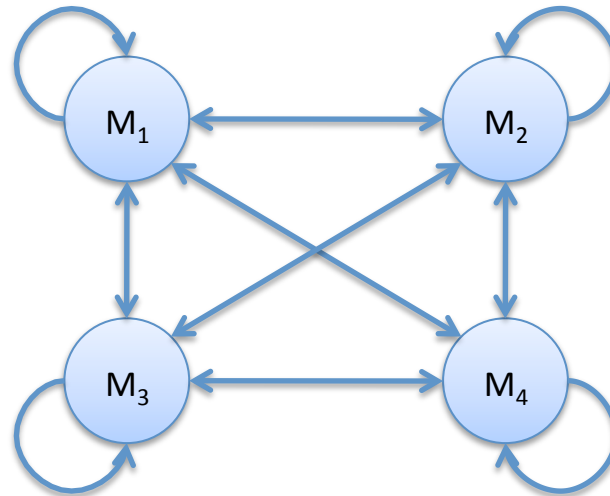Markov chain:
States = sets
M



Generate sequence of states:  $M^{(1)}, M^{(2)}, M^{(3)}, \ldots$

Markov Chain Convergence Thm: $M^{(i)} \rightarrow \pi$

# Metropolis-Hastings

Define transition probabilities of Markov Chain so
$\pi$ = desired distribution.

Markov chain:
States = sets
M



Distribution on gene sets: $\Pr[\boldsymbol{M}] \sim \boldsymbol{e^{c\,W(M)}}$

In general: no guarantees on rate of convergence

# MCMC approach

*Thm.* Markov Chain is rapidly mixing.

Returns a distribution on sets, not just optimal [max W($\mathbf{M}$)] set

No assumptions on distribution of mutations
- i.e. independence not necessary
- can handle various mutation types

# Experimental Results

- Simulated data

- Cancer data

    1. **Brain cancer (GBM)** [TCGA, *Nature* (2008)]
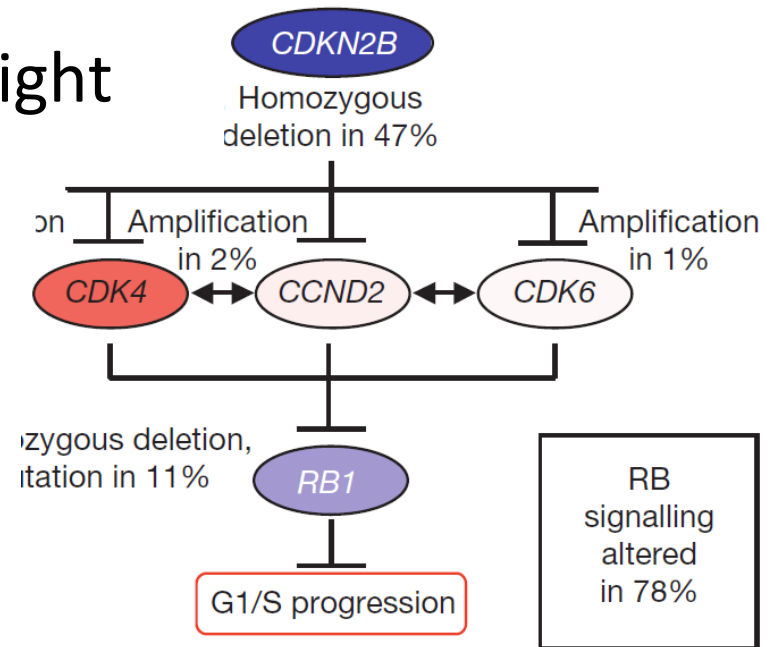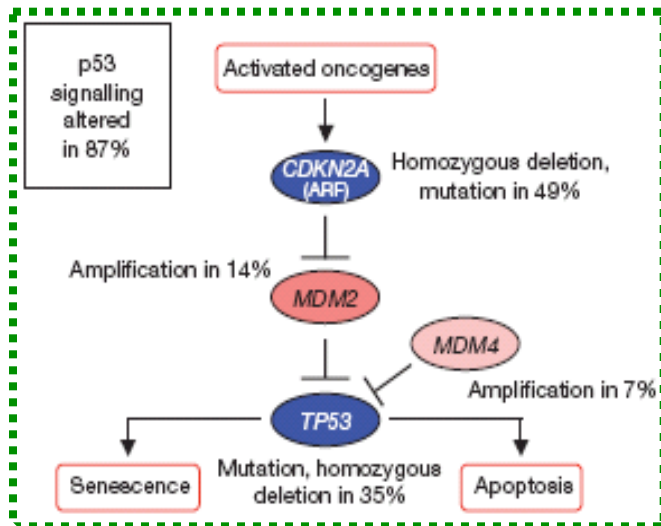
        601 sequenced genes in 84 samples
        Array copy number data on *all* genes

    2. **Lung Adenocarcinoma** [Ding et al., *Nature* (2008)]

        623 sequenced genes in 188 samples

# Brain Cancer (GBM)

- **M** = {CDKN2B, RB1, CDK4}
  - not the set with highest weight

- **M** = {TP53, CDKN2A}
  - p53 signaling pathway



From [TCGA, *Nature*, 2008]