# Simultaneous Identification of Multiple Driver Pathways in Cancer

## Max Leiserson
## ISMB 2014

Joint work with Benjamin Raphael, Roded Sharan and Dima Blokh.

# Clonal theory of cancer

Clonal Theory (Nowell 1976)

Passenger mutations

Driver mutation

Founder cell

Time (cell divisions)

Cell population

sequence genome

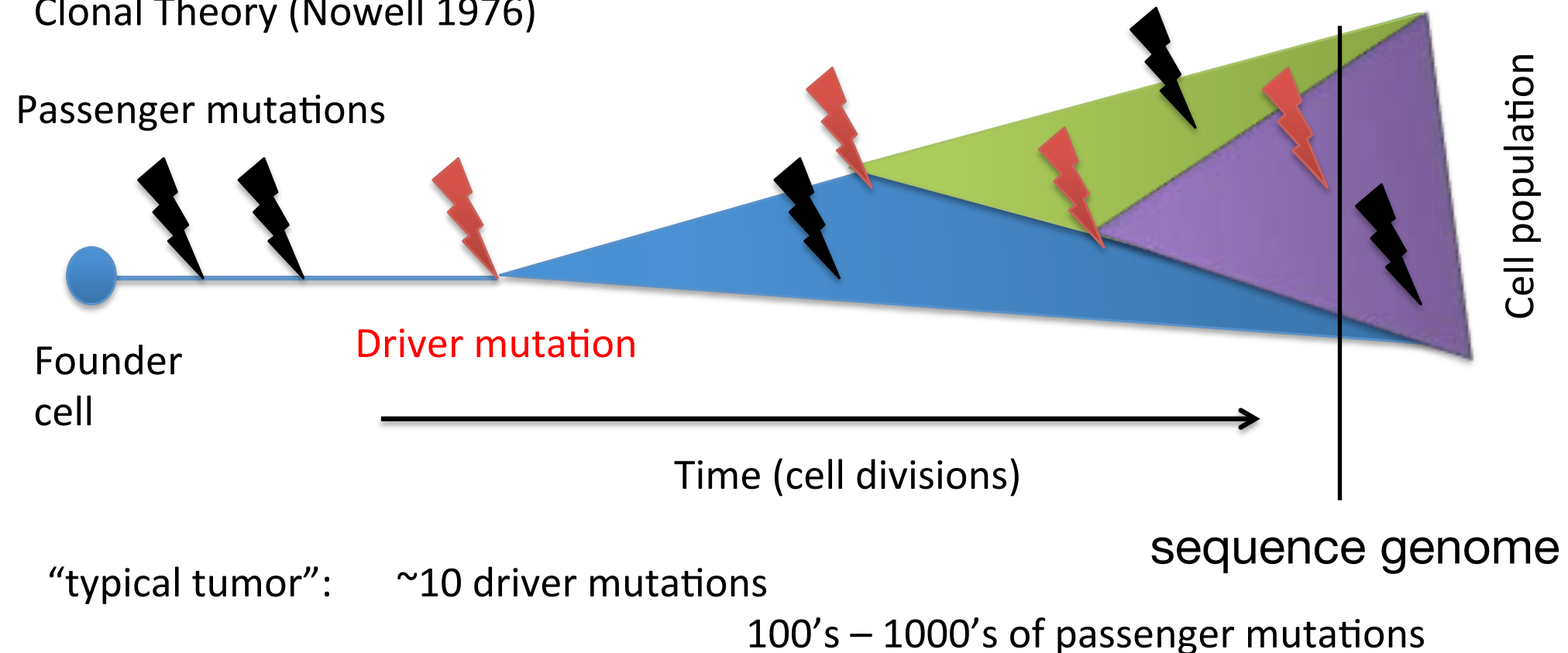"typical tumor":     ~10 driver mutations

100's – 1000's of passenger mutations

# Clonal theory of cancer

Clonal Theory (Nowell 1976)

Passenger mutations

Driver mutation

Founder cell

Time (cell divisions)

Cell population

sequence genome

"typical tumor":     ~10 driver mutations

100's – 1000's of passenger mutations

**Types of Variation in Tumor Genomes**

*Single Nucleotide Variants*

*Copy Number Variants*

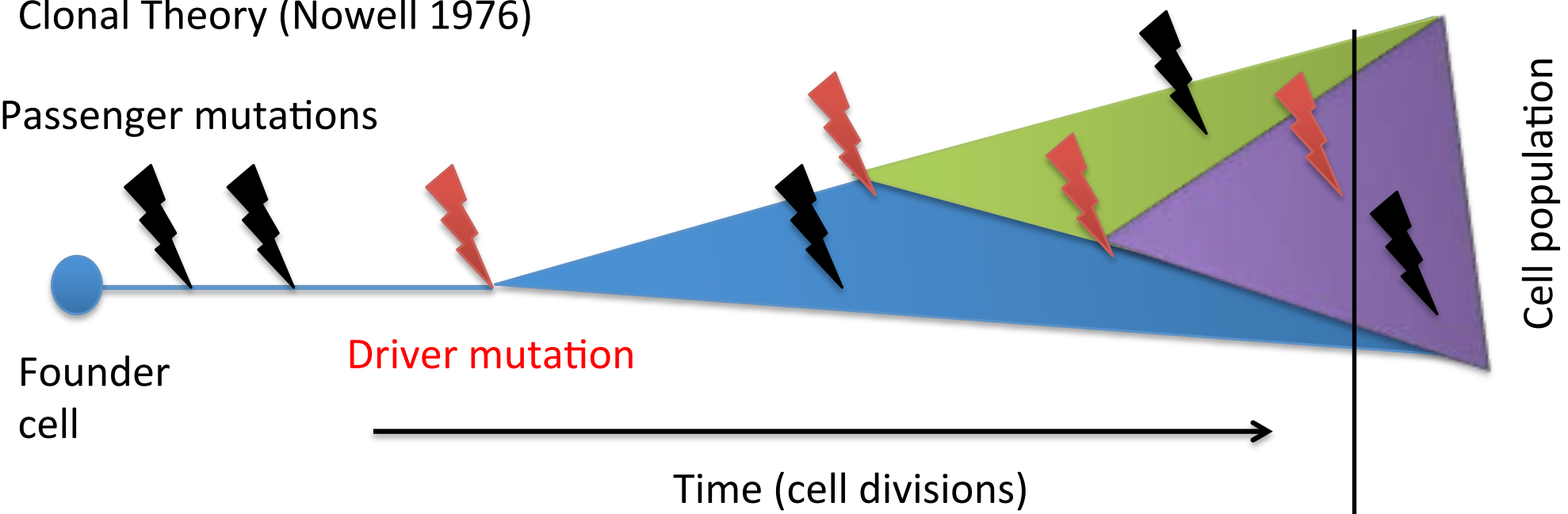**Healthy** ..ACGTCA**T**CATGA..

**Tumor** ..ACGTCA**G**CATGA..

*deletion*

*amplification*

# Clonal theory of cancer

Clonal Theory (Nowell 1976)

Passenger mutations

Driver mutation

Founder cell

Cell population

Time (cell divisions)

sequence genome

"typical tumor":  ~10 driver mutations

100's – 1000's of passenger mutations

## Types of Variation in Tumor Genomes

*Single Nucleotide Variants*

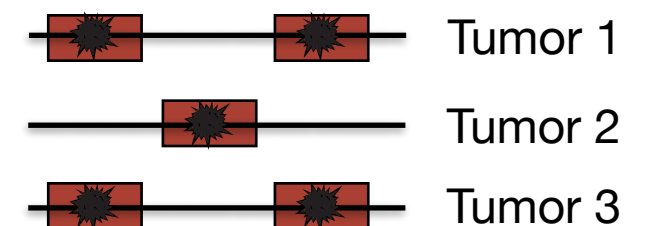*Copy Number Variants*

**Healthy** . . ACGTCA**T**CATGA . .

**Tumor** . . ACGTCA**G**CATGA . .

*deletion*

*amplification*

*Compare variation across tumors*
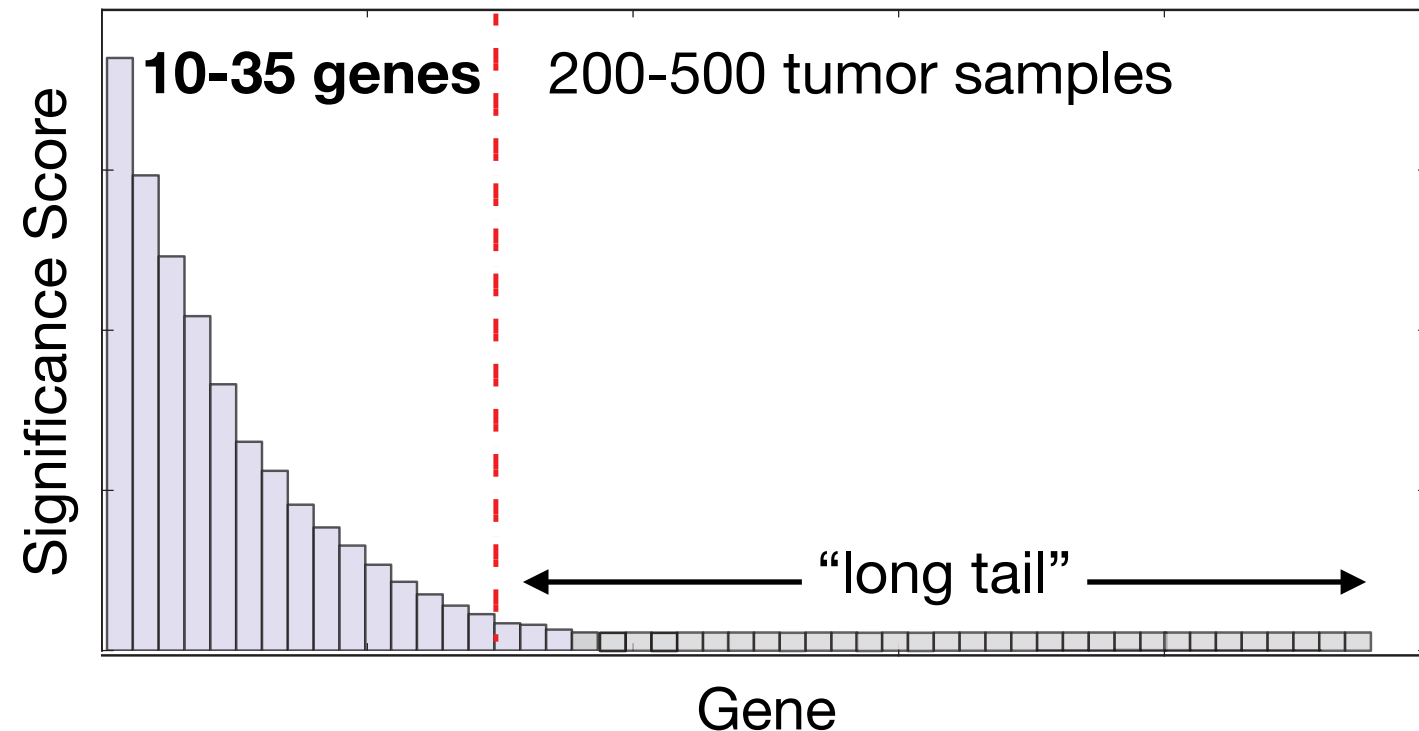
Tumor 1

Tumor 2

Tumor 3

= gene    = SNV / CNA

2

# Significantly mutated genes in cancer

**Significance Score**
Mutations weighted by:
- Gene length
- Mutation context
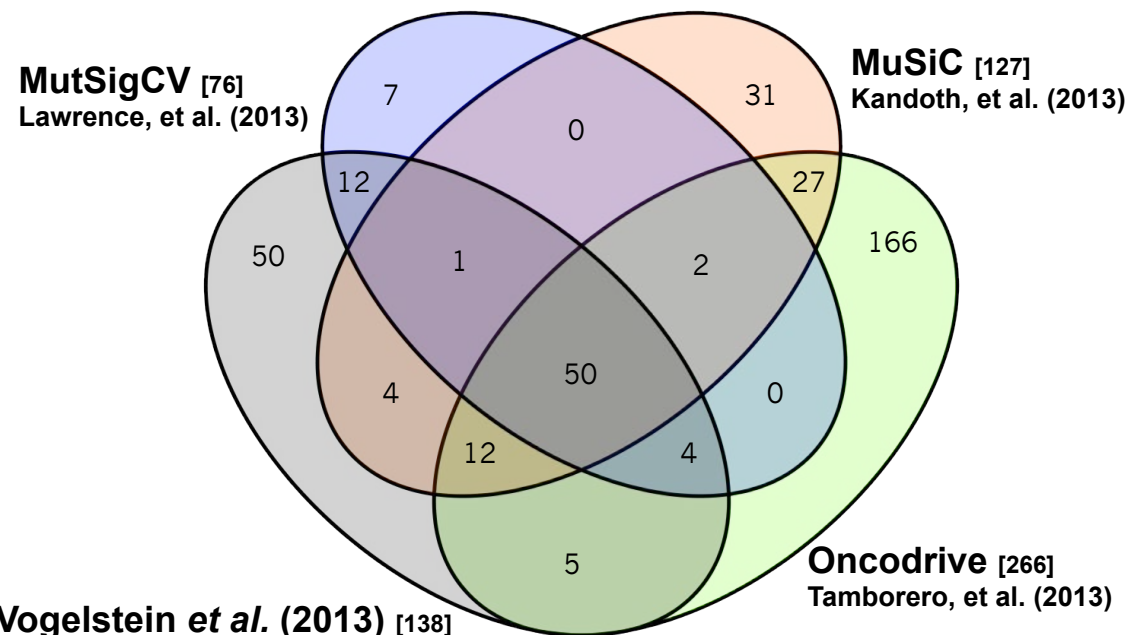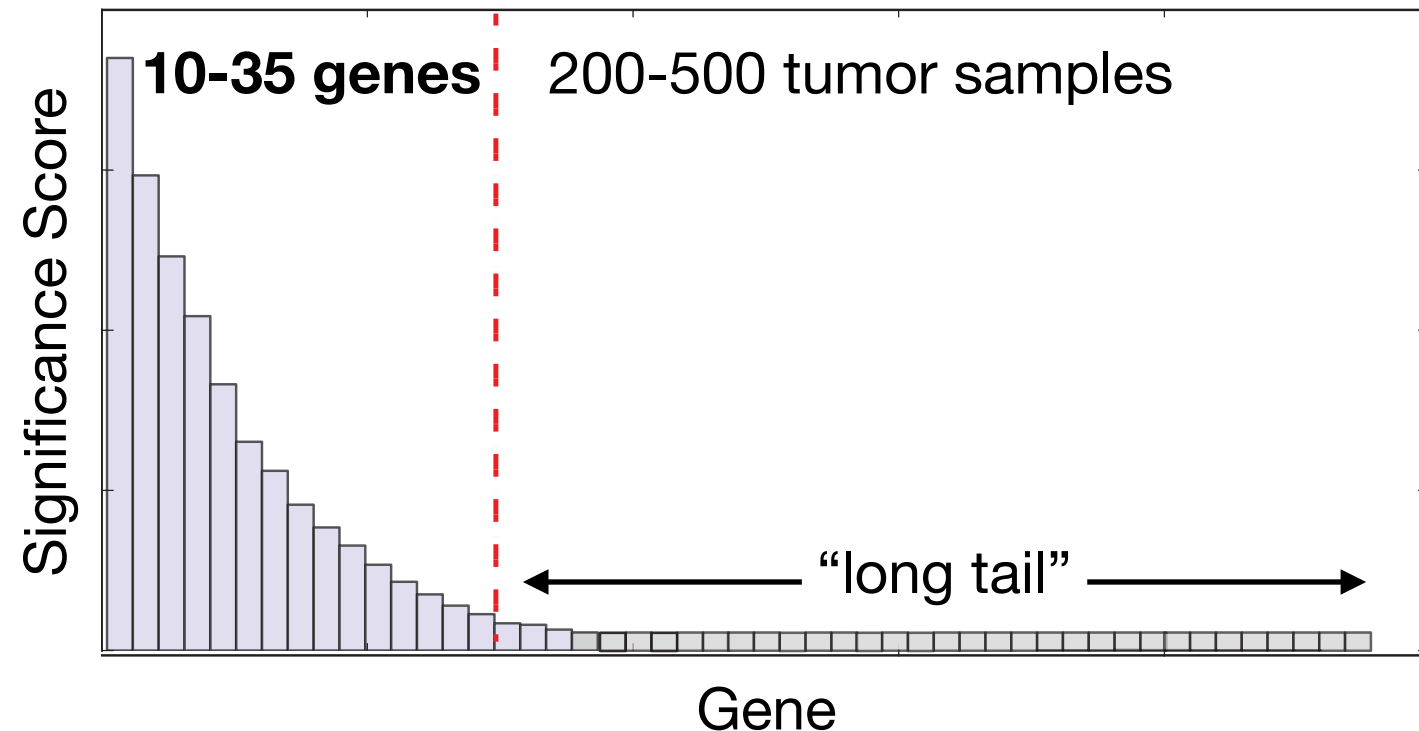- Expression level
- Replication timing
- …

# "Long tail" of mutated genes complicates finding driver mutations

**Significance Score**
Mutations weighted by:
- Gene length
- Mutation context
- Expression level
- Replication timing
- …





**Comparison of significantly mutated genes**
*TCGA Pan-Cancer Dataset* [TCGA Research Network, 2013]
>3000 tumor samples of twelve cancer types

# "Long tail" of mutated genes complicates finding driver mutations

**Significance Score**
Mutations weighted by:
- Gene length
- Mutation context
- Expression level
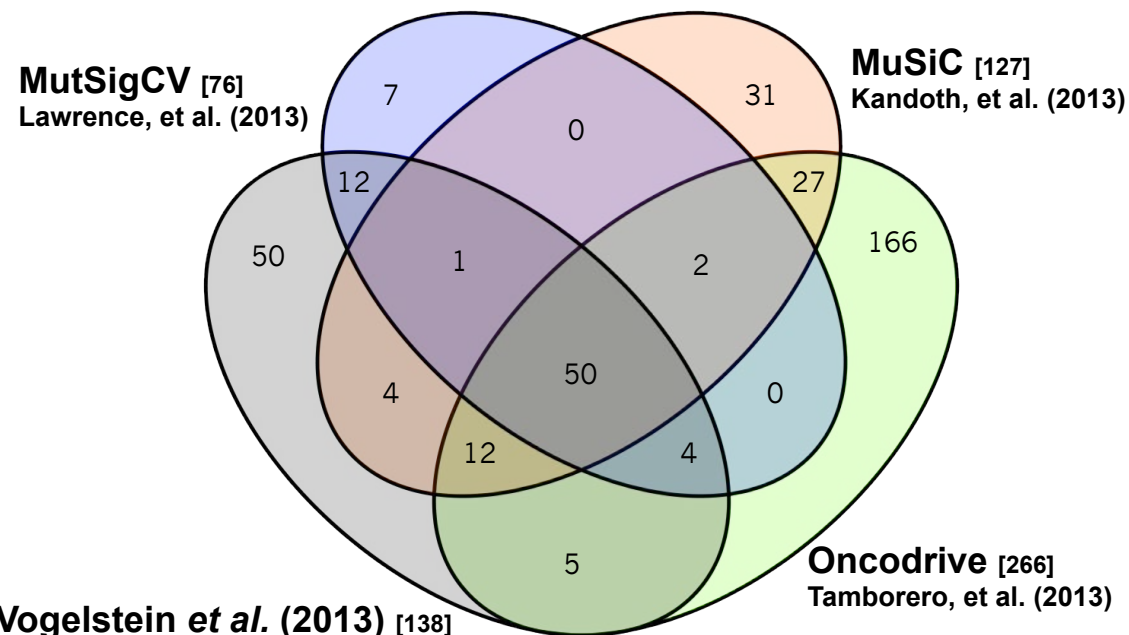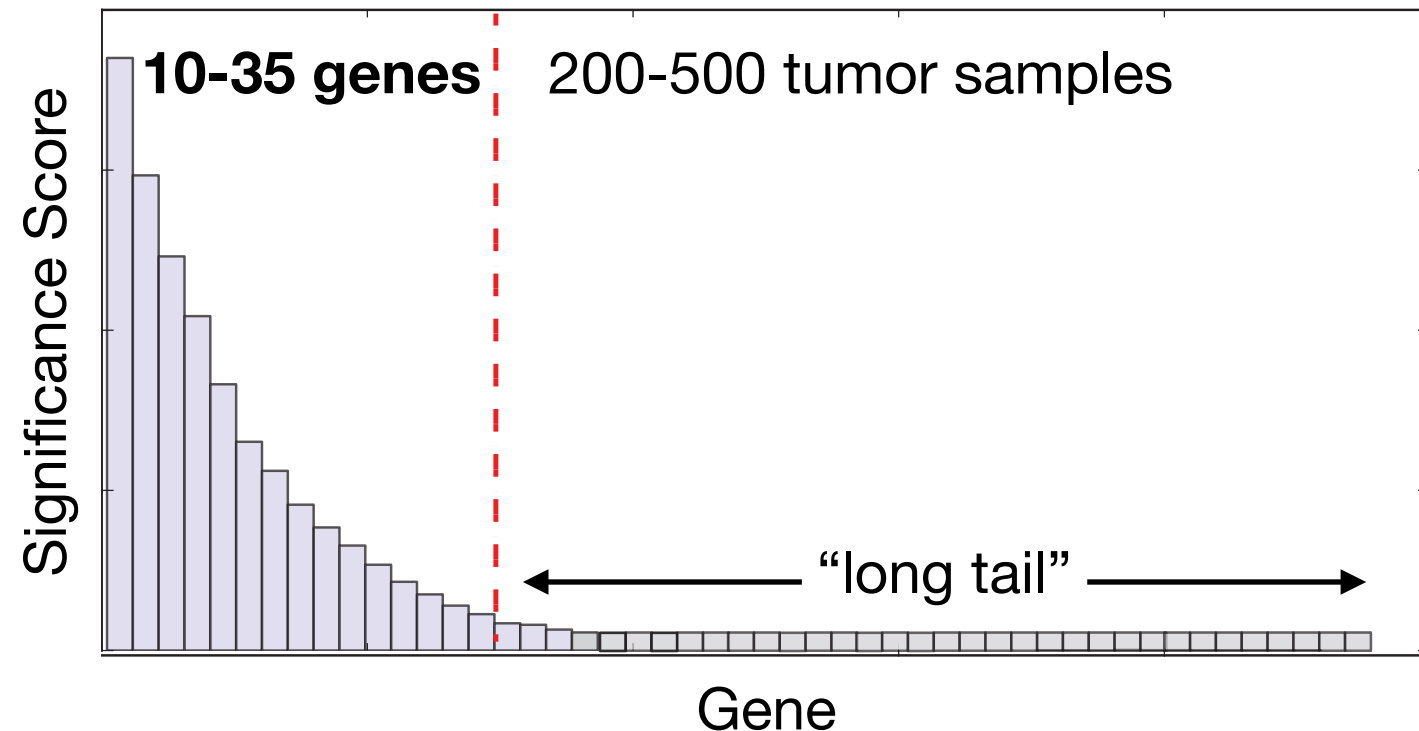- Replication timing
- …



**Comparison of significantly mutated genes**
*TCGA Pan-Cancer Dataset* [TCGA Research Network, 2013]
>3000 tumor samples of twelve cancer types

**50 genes agreed upon by all methods
Hundreds found by only one method**

# Mutations target pathways

**Significance Score**
Mutations weighted by:
- Gene length
- Mutation context
- Expression level
- Replication timing
- …



**Comparison of significantly mutated genes**
*TCGA Pan-Cancer Dataset* [TCGA Research Network, 2013]
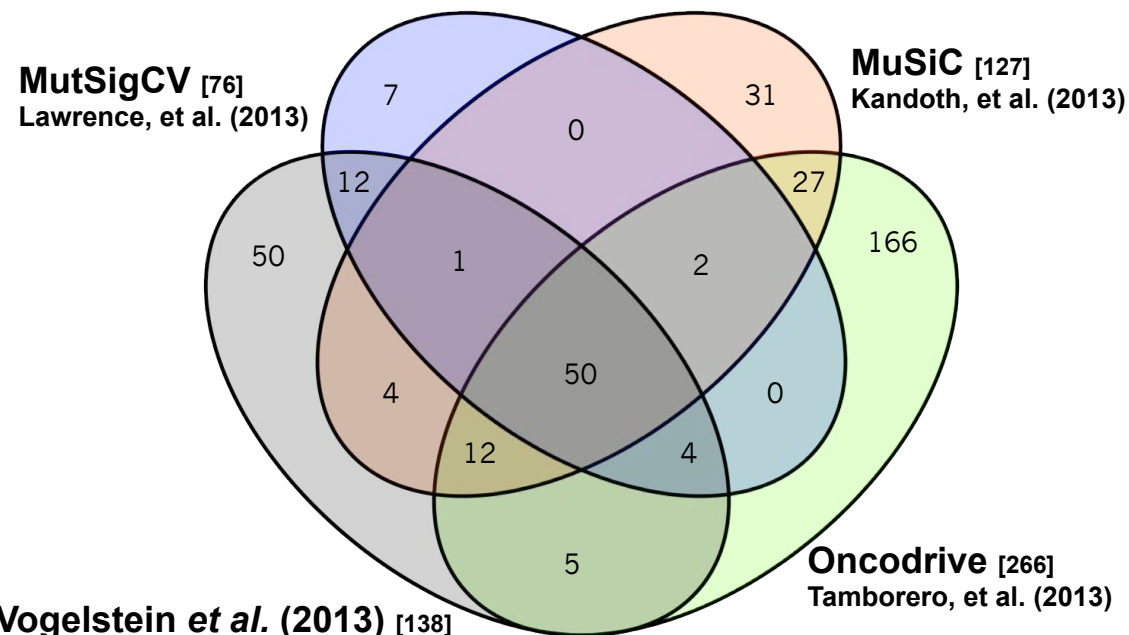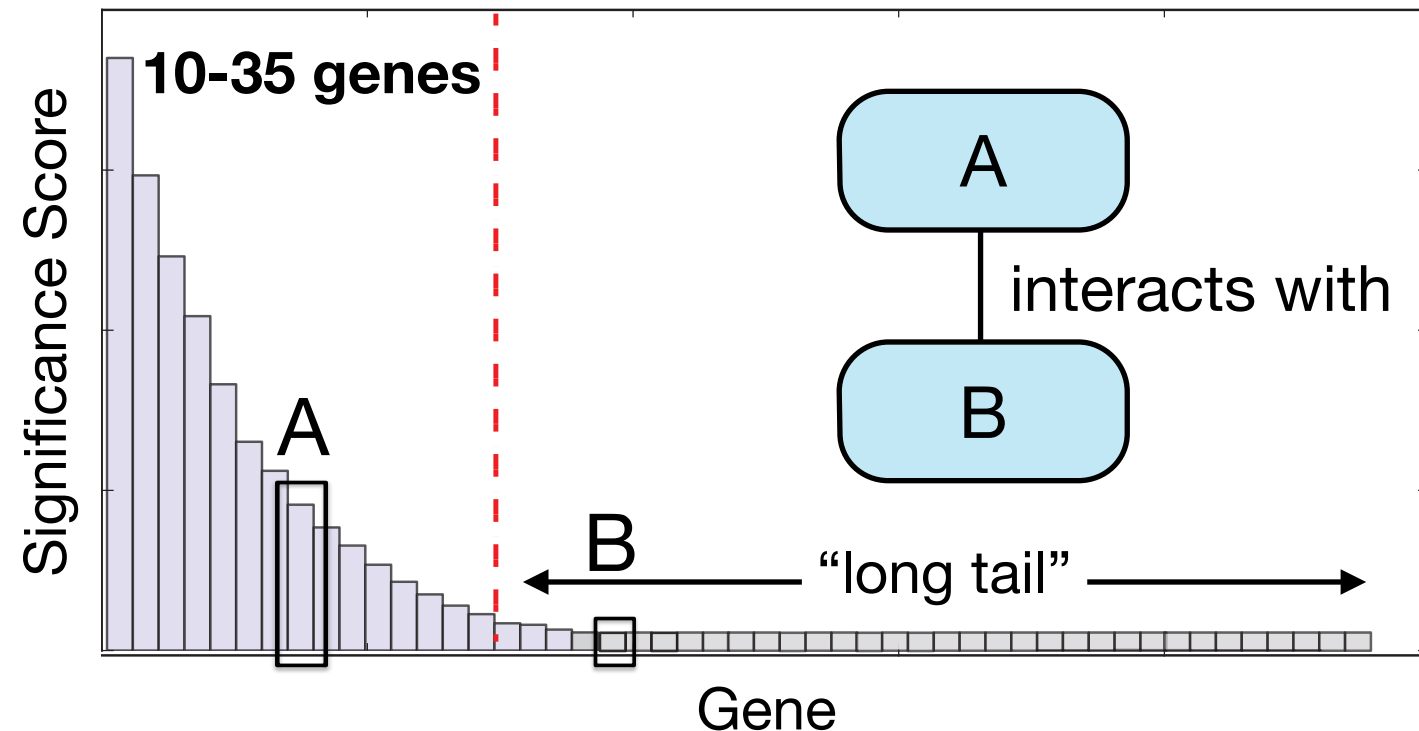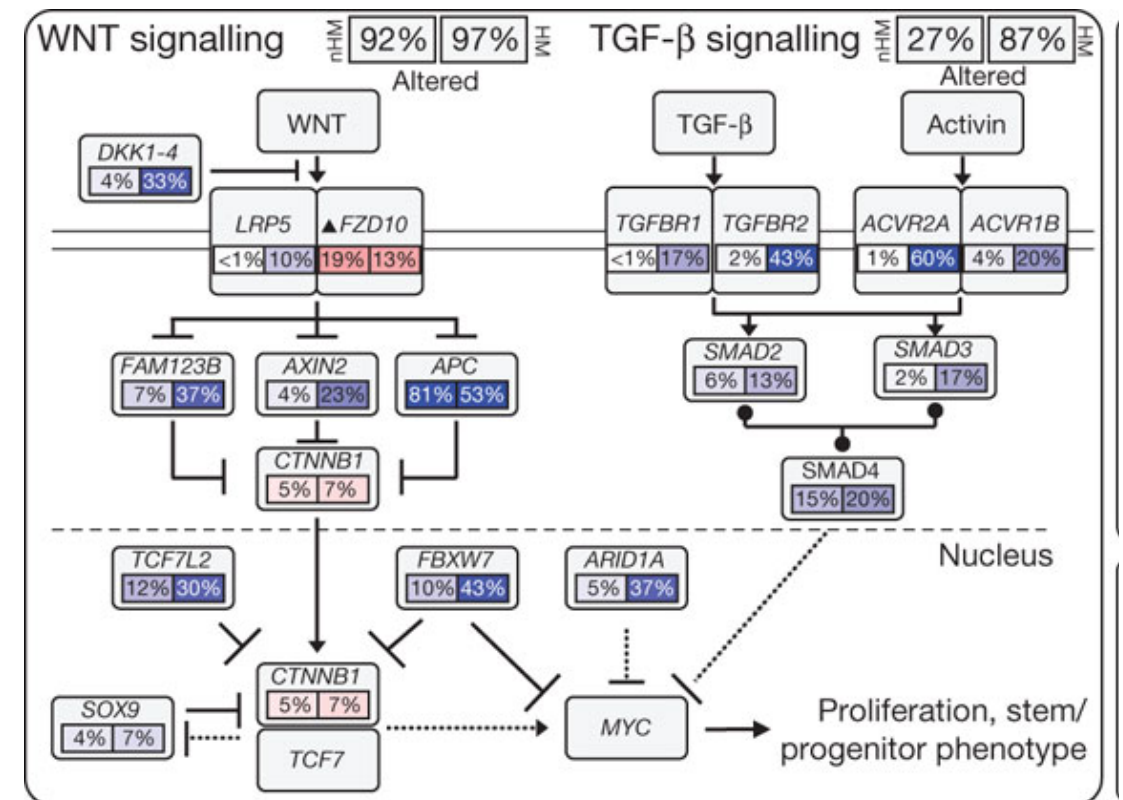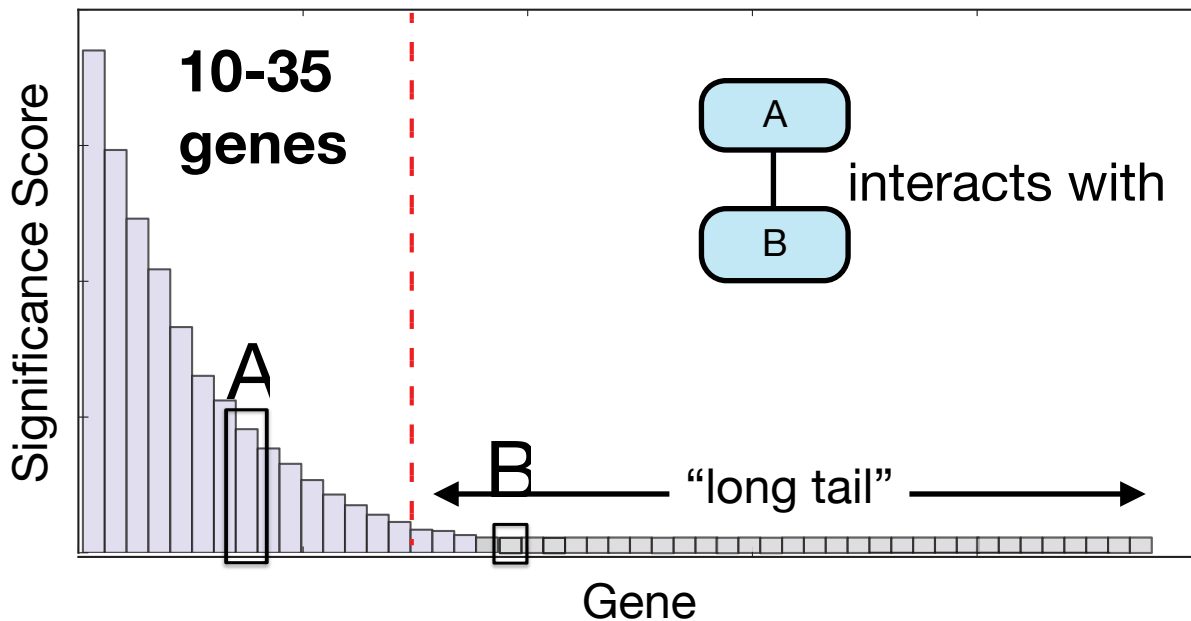>3000 tumor samples of twelve cancer types

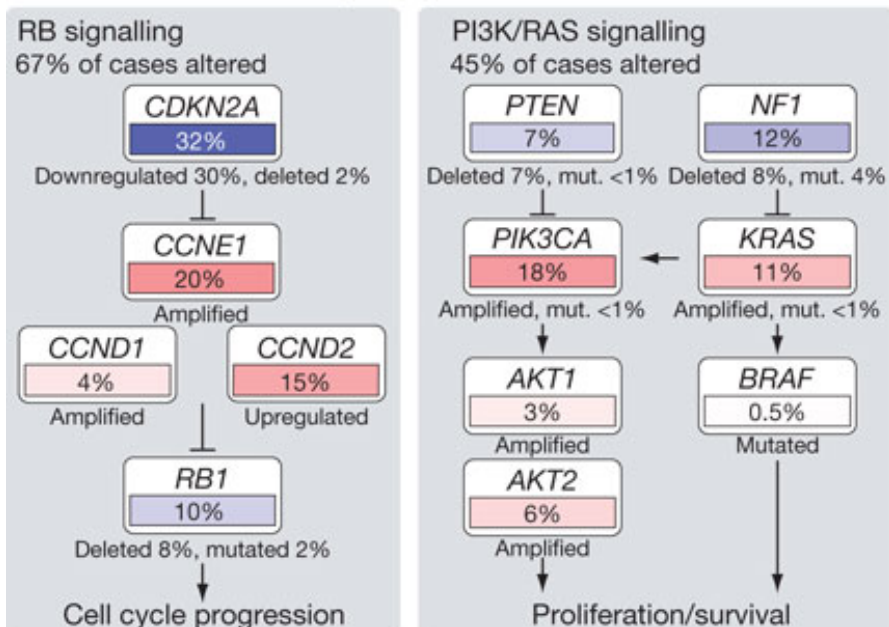**50 genes agreed upon by all methods**
**Hundreds found by only one method**

# Identifying significantly mutated pathways



10-35 genes

Significance Score

interacts with

"long tail"

Gene

TCGA Colorectal (*Nature* 2012)

TCGA Ovarian (*Nature* 2011)

# Identifying significantly mutated pathways



**10-35 genes**

Significance Score

A

B

"long tail"

Gene

A interacts with B



TCGA Colorectal (*Nature* 2012)



TCGA Ovarian (*Nature* 2011)

- Novel pathways
- *Crosstalk*
- Origin specific

# Cancer pathways often harbor *mutually exclusive* mutations

*Few* driver mutations distributed across *multiple* pathways

→ Approximately one driver mutation per pathway per patient

cell membrane

EGFR

RAS

RAF

MEK

MAPK

Transcription factors

# Cancer pathways often harbor *mutually exclusive* mutations

*Few* driver mutations distributed across *multiple* pathways

→ Approximately one driver mutation per pathway per patient

**1. Exclusivity**

[Thomas, et al. (2007)]

# Cancer pathways often harbor *mutually exclusive* mutations

*Few* driver mutations distributed across *multiple* pathways

→ Approximately one driver mutation per pathway per patient

**1. Exclusivity**

→ Many patients have a mutation in important cancer pathway
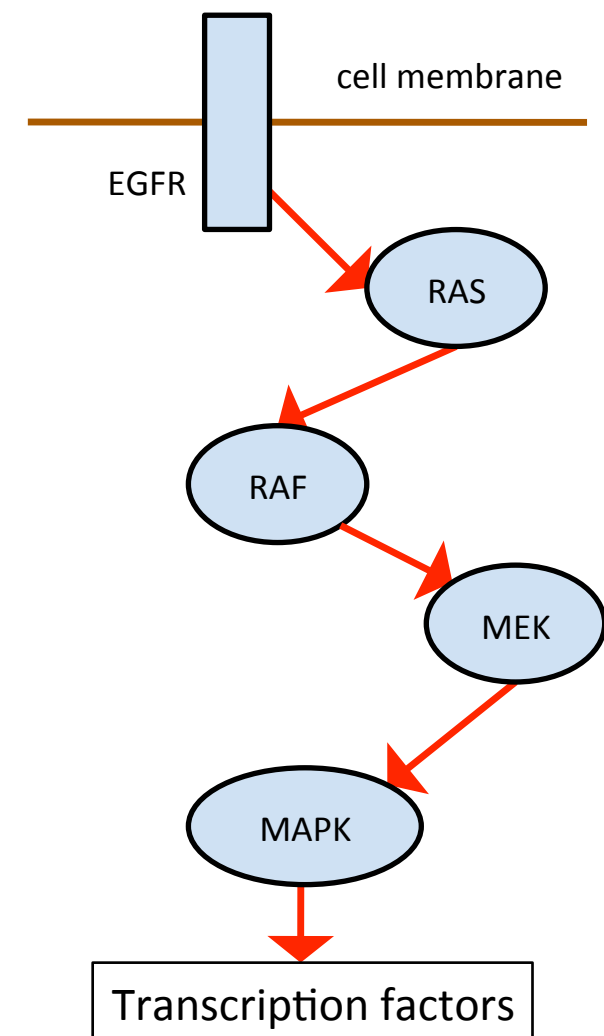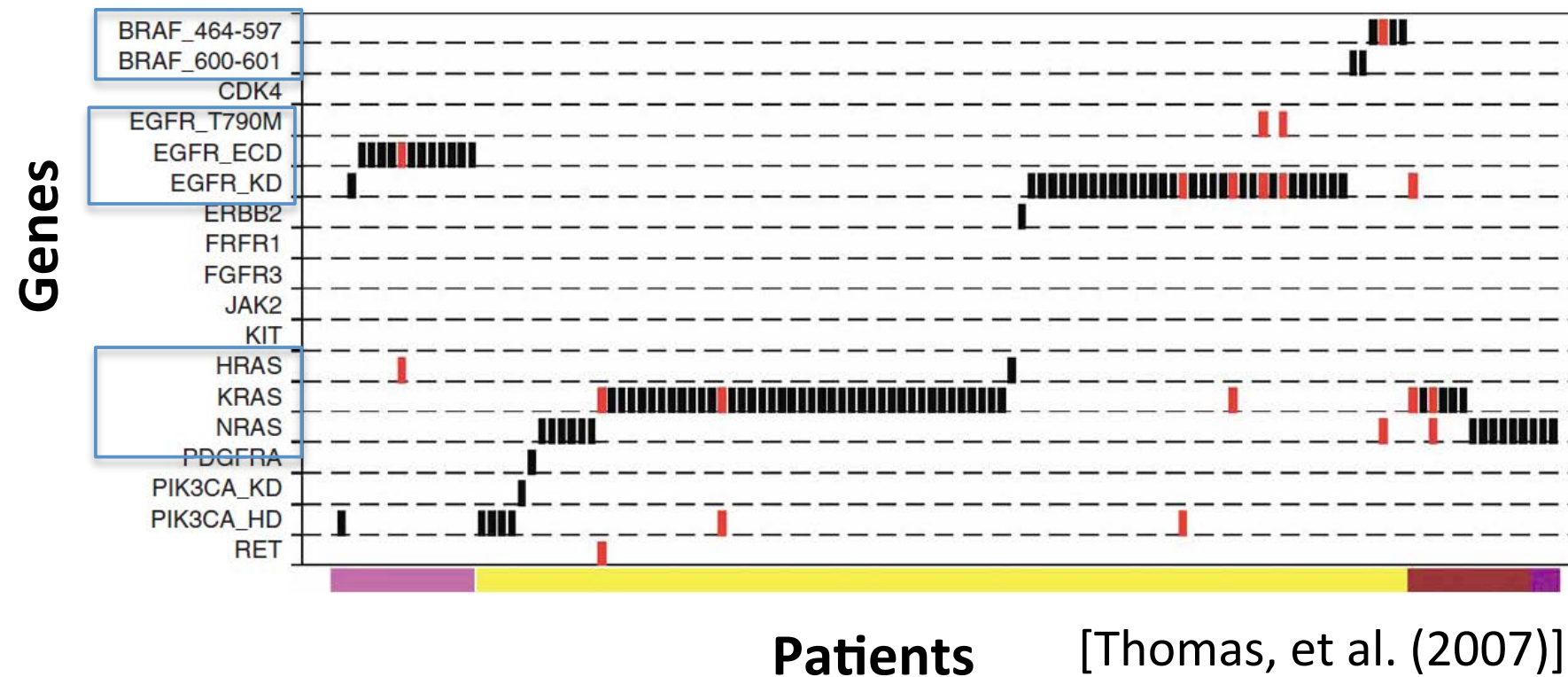
**2. High Coverage**

# De novo driver exclusivity (Dendrix)

*de novo:* **without prior biological information (pathways, interactions, etc.)**

**Goal**: Find sets **M** of genes with:
- *High coverage*: many patients with ≥ 1 mutation in **M**
- *Approximate exclusivity:* most patients have ≤ 1 mutation in **M**

Weight W(M)

- Finding optimal set is NP-Hard.
- MCMC algorithm samples sets in proportion to weight.



**Overlap(M)**   *patients*

*genes*

**Coverage(M)**

■ exclusive   ■ overlapping   ▪ none

Vandin *et al*. (2012) De novo discovery of mutated driver pathways in cancer. *Genome Research*, 22:375-85.

# Iter-Dendrix: multiple pathways

Cancer requires mutations to more than one pathway

**Find multiple pathways greedily.**



= gene *i* is mutated in patient *j*,
= otherwise.

# Iter-Dendrix: multiple pathways

Cancer requires mutations to more than one pathway

**Find multiple pathways greedily.**



*n* patients

*m* genes

■ = gene *i* is mutated in patient *j*,
■ = otherwise.

Vandin *et al*. (2012) De novo discovery of mutated driver pathways in cancer. *Genome Research*, 22:375-85.

# Iter-Dendrix: multiple pathways

Cancer requires mutations to more than one pathway

**Find multiple pathways greedily.**



= gene *i* is mutated in patient *j*,

= otherwise.

Vandin *et al*. (2012) De novo discovery of mutated driver pathways in cancer. *Genome Research*, 22:375-85.

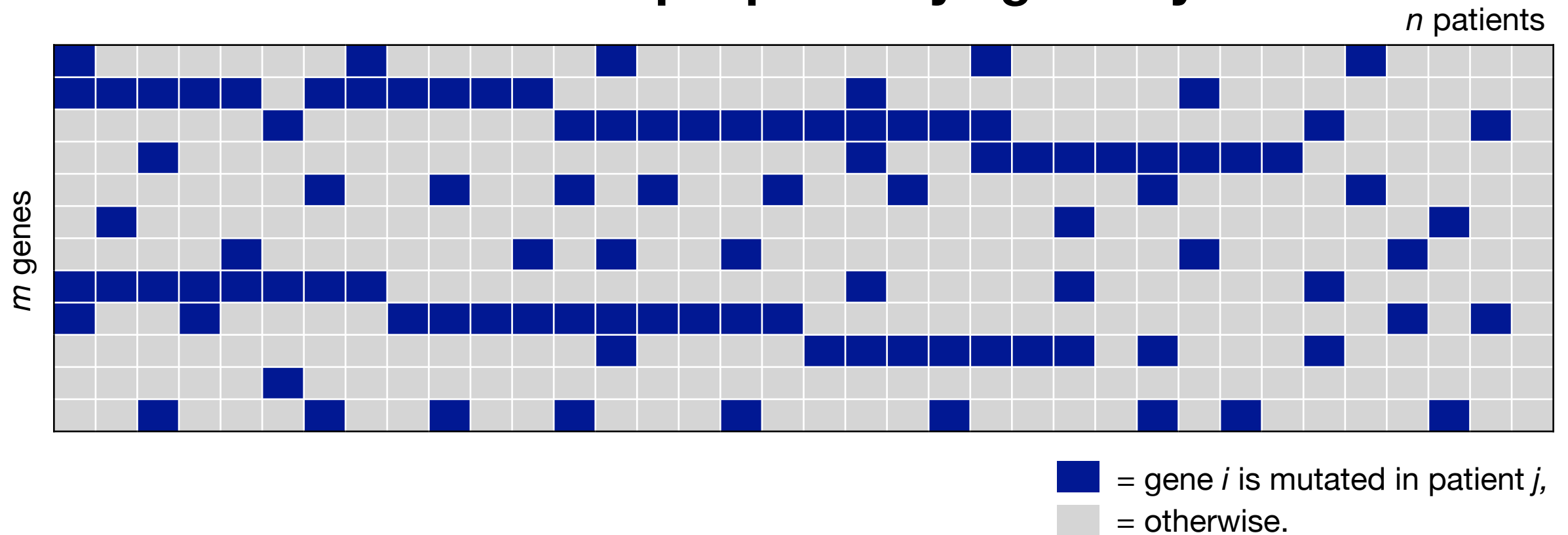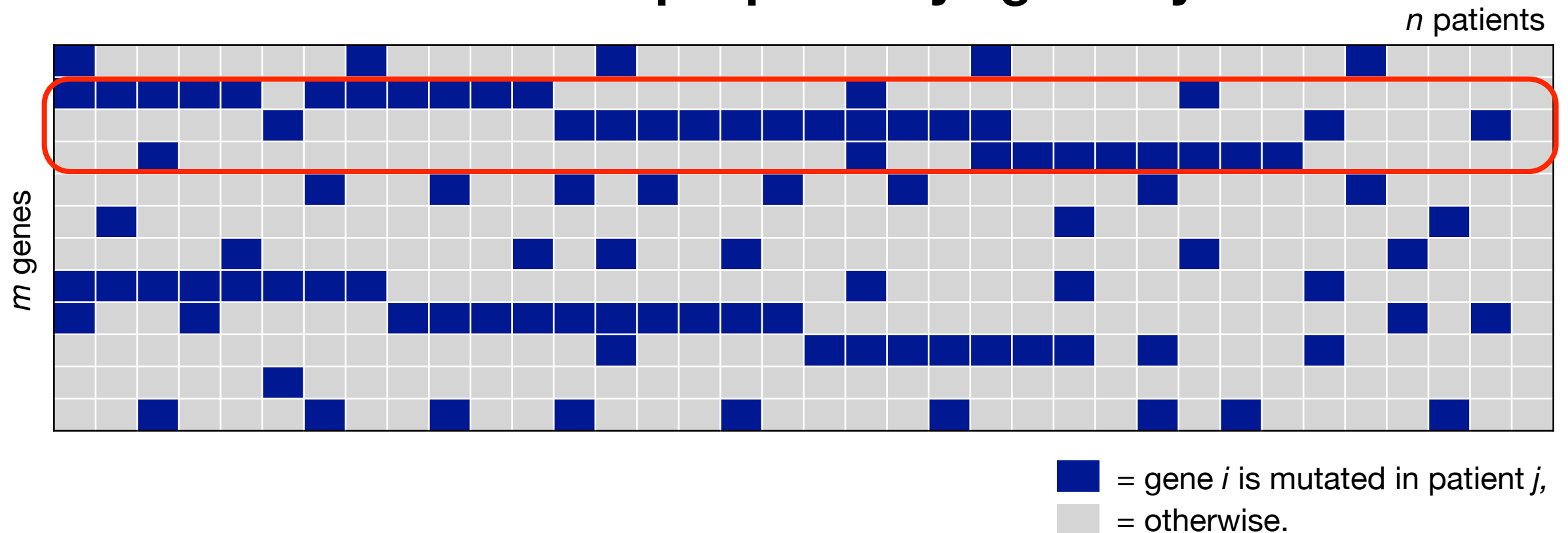# Iter-Dendrix: multiple pathways

Cancer requires mutations to more than one pathway
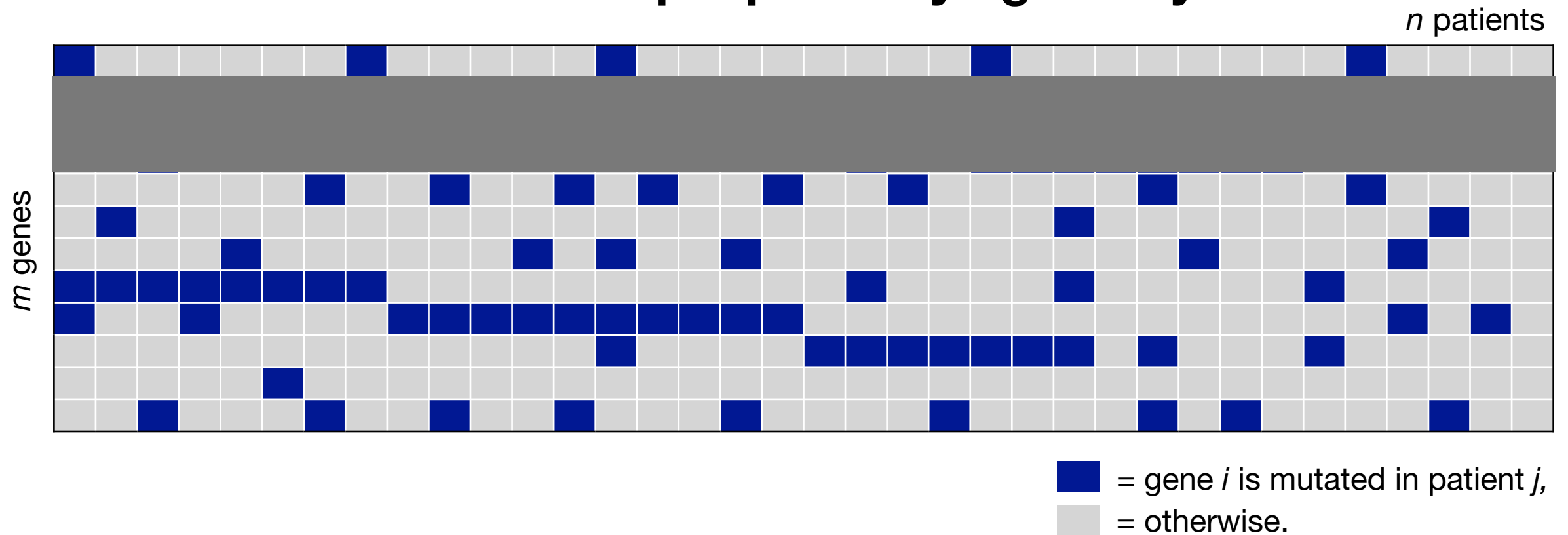
**Find multiple pathways greedily.**



= gene $i$ is mutated in patient $j$,

= otherwise.

Vandin *et al*. (2012) De novo discovery of mutated driver pathways in cancer. *Genome Research*, 22:375-85.

# Greedy can fail

Cancer requires mutations to more than one pathway

Find multiple pathways greedily.

**Groups most frequently mutated genes even without exclusivity.**



= gene *i* is mutated in patient *j*,
= otherwise.

# Greedy can fail

Cancer requires mutations to more than one pathway

Find multiple pathways greedily.

**Groups most frequently mutated genes even without exclusivity.**

# Greedy can fail

Cancer requires mutations to more than one pathway

Find multiple pathways greedily.

**Groups most frequently mutated genes even without exclusivity.**
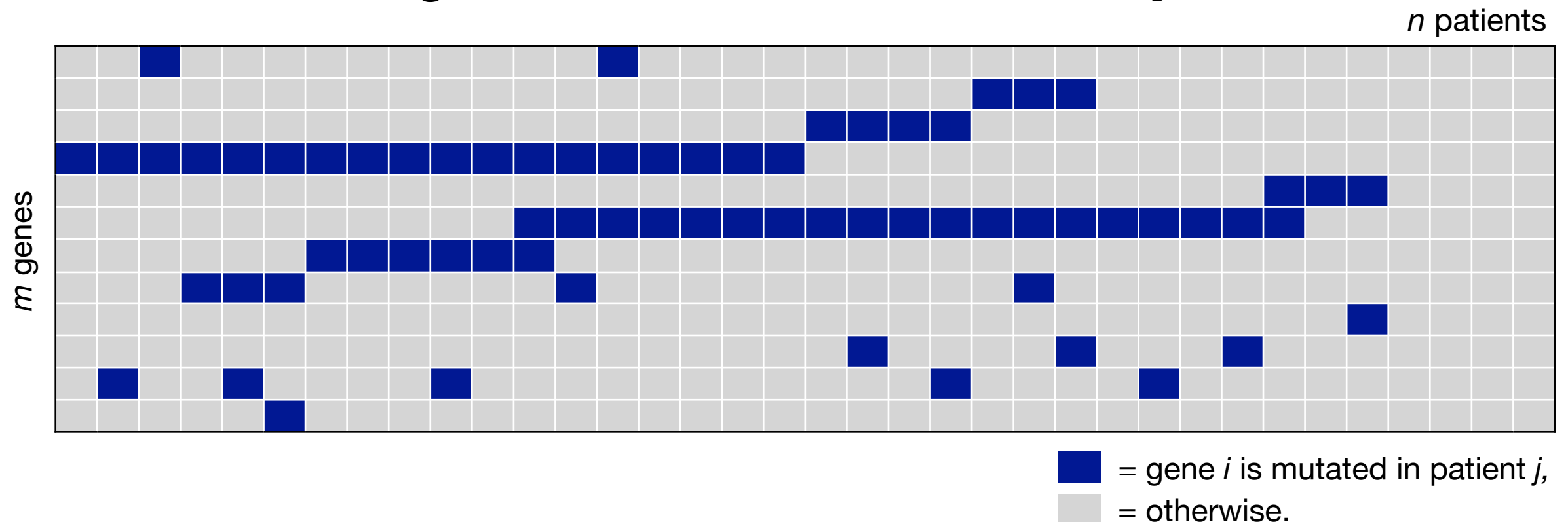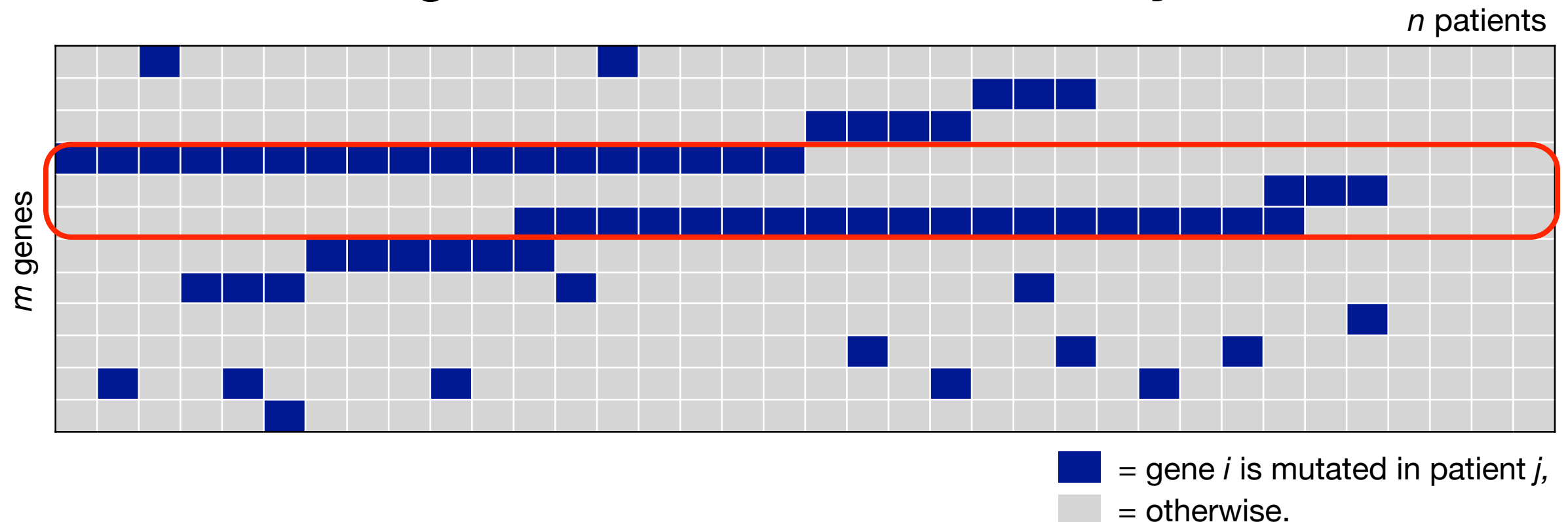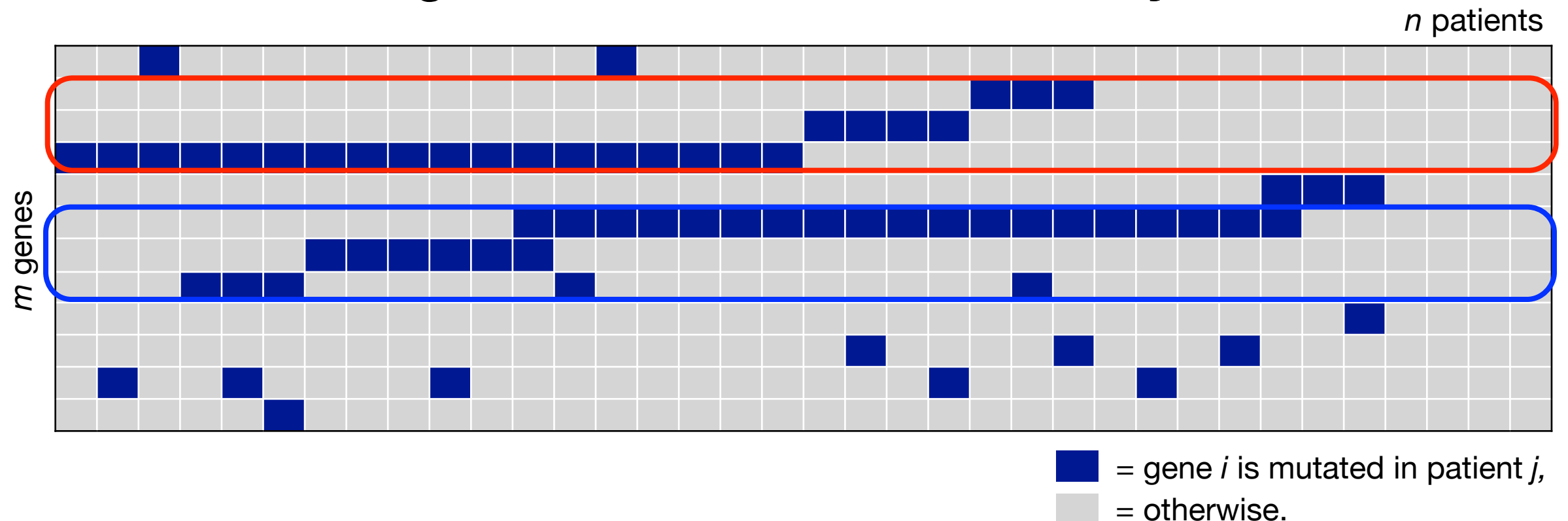


*n* patients

*m* genes

■ = gene *i* is mutated in patient *j*,
▢ = otherwise.

# Multi-Dendrix

Cancer requires mutations to more than one pathway
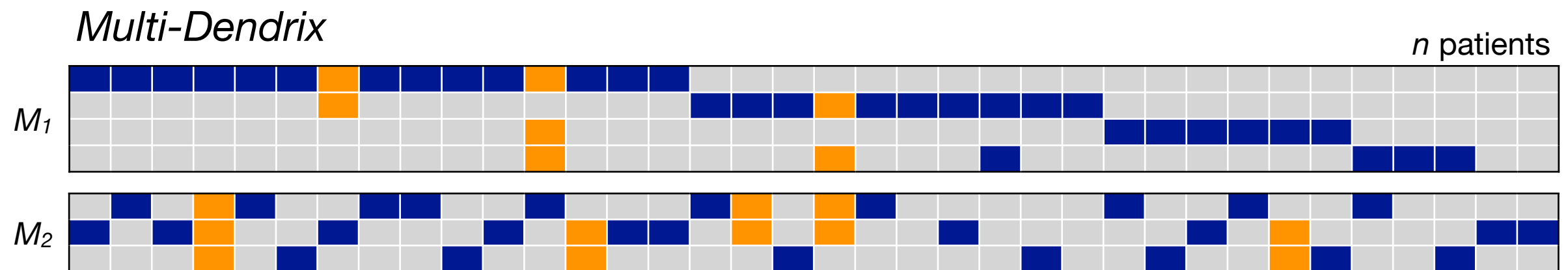
**Find <u>multiple</u> pathways simultaneously**

- ILP rapidly finds optimal solution
- Searches wide range of parameters to find stable collections of gene sets



*n* patients

*m* genes

■ = gene *i* is mutated in patient *j*,

□ = otherwise.

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Multi-Dendrix

Cancer requires mutations to more than one pathway



*Multi-Dendrix*

$n$ patients

$M_1$

$M_2$

Most samples have approximately one mutation in *each* of $t$ pathways.
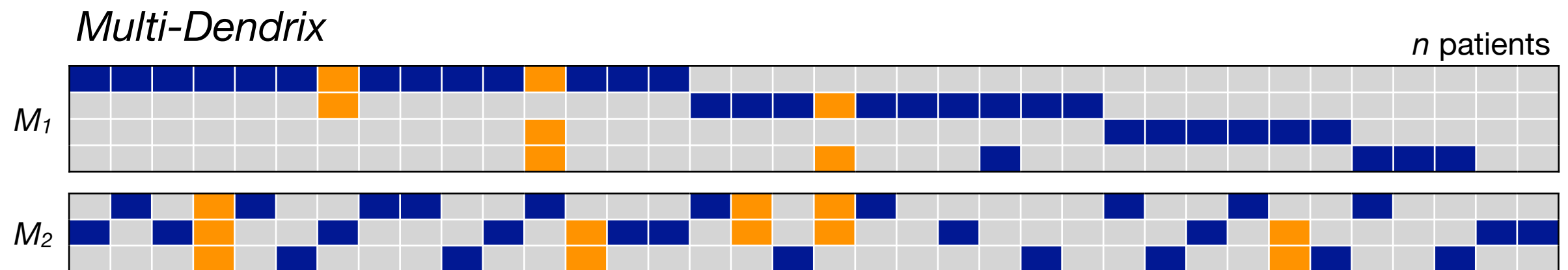
$$\textbf{maximize } W'_\alpha(\mathbf{M}) = \sum_{M \in \mathbf{M}} |\Gamma(M)| - \alpha\omega(M)$$

coverage     overlap

Integer Linear Program

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Multi-Dendrix

Cancer requires mutations to more than one pathway



*Multi-Dendrix*

*n* patients

$M_1$

$M_2$

Most samples have approximately one mutation in *each* of *t* pathways.

**maximize** $W'_\alpha(\mathbf{M}) = \sum_{M \in \mathbf{M}} |\Gamma(M)| - \alpha \omega(M)$

parameter to control "weight" of exclusivity

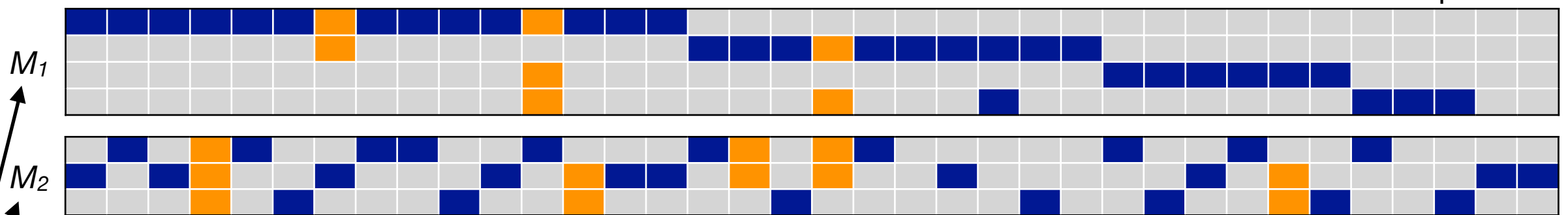coverage     overlap

Integer Linear Program

# Multi-Dendrix

Cancer requires mutations to more than one pathway



*Multi-Dendrix*

*n* patients

$M_1$

$M_2$

Most samples have approximately one mutation in *each* of *t* pathways.

$|M_i| \in [k_{\min}, k_{\max}]$

**maximize** $W'_\alpha(\mathbf{M}) = \sum_{M \in \mathbf{M}} |\Gamma(M)| - \alpha\omega(M)$

parameter to control "weight" of exclusivity

coverage     overlap

Integer Linear Program

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Contributions

A new algorithm, **Multi-Dendrix**, for identifying driver pathways *de novo:*

1. **Outperforms previous methods on simulated data in speed and accuracy**

2. Identifies gene sets that overlap known pathways in TCGA datasets

3. Ongoing work extending Multi-Dendrix to large datasets and overlapping pathways
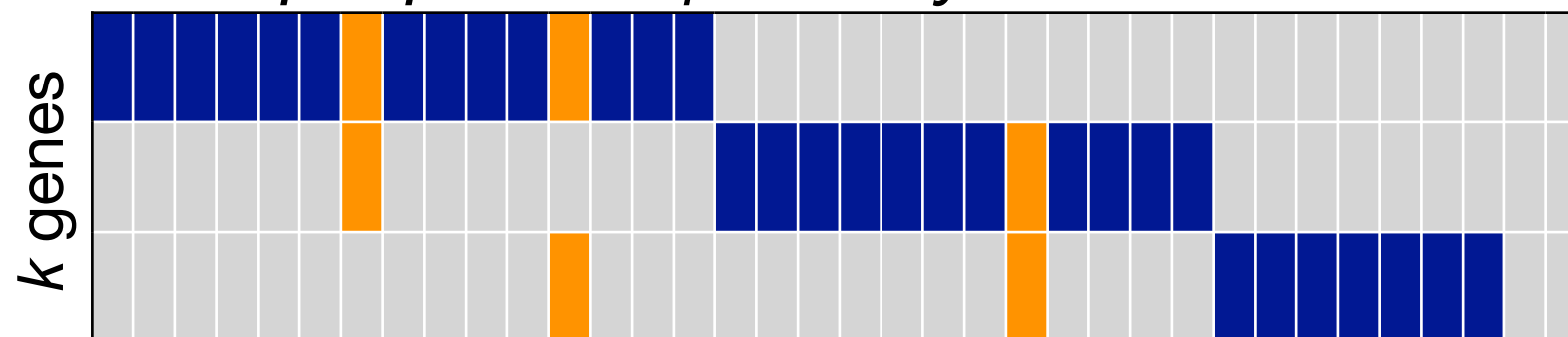
# Multi-Dendrix is significantly better on simulated data

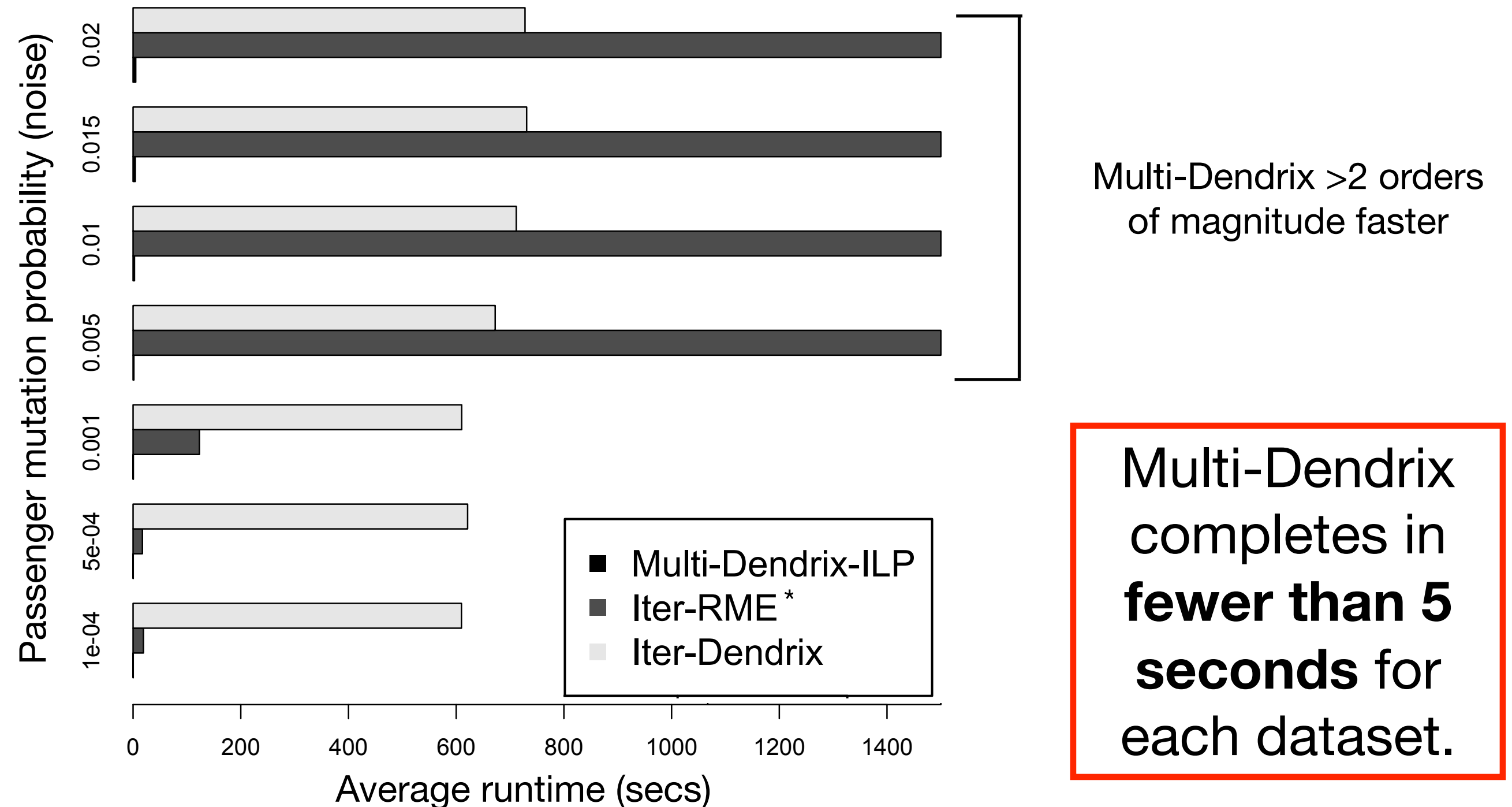|     | Avg. distance $d$ from planted pathways | | |
| --- | --- | --- | --- |
| $q$ | Multi-Dendrix | Iter-RME* | Iter-Dendrix |
| 0.0 | $0.02 \pm 0.19$ | $\mathbf{0.01 \pm 0.12}$ | $0.30 \pm 0.86$ |
| *0.0001* | *$0.02 \pm 0.18$* | ***$0.01 \pm 0.16$*** | *$0.30 \pm 0.86$* |
| 0.0005 | $\mathbf{0.04 \pm 0.23}$ | $0.10 \pm 0.40$ | $0.35 \pm 0.89$ |
| *0.001* | ***$0.10 \pm 0.35$*** | *$0.32 \pm 0.60$* | *$0.44 \pm 1.01$* |
| 0.005 | $\mathbf{0.44 \pm 0.71}$ | – | $0.75 \pm 1.07$ |
| 0.01 | $\mathbf{1.03 \pm 1.00}$ | – | $1.20 \pm 1.15$ |
| 0.015 | $\mathbf{1.68 \pm 1.16}$ | – | $1.78 \pm 1.26$ |
| 0.02 | $\mathbf{2.17 \pm 1.24}$ | – | $2.21 \pm 1.29$ |

$P < 0.01$

360 mutated genes in 160 patients

*Example planted pathway*



*k* genes

- High coverage and mutually exclusive
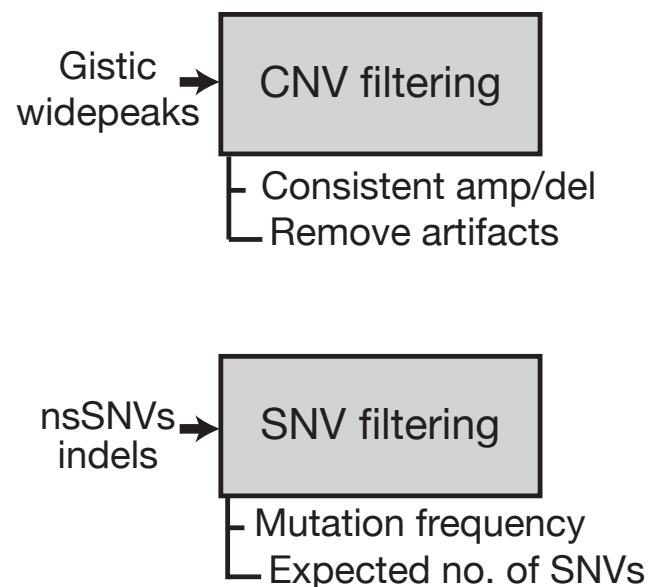- Add passenger mutations (noise) to all genes

*Miller *et al*. (2011). *BMC medical genomics*, **4**:34.

# Multi-Dendrix is significantly **faster** on simulated data



Multi-Dendrix >2 orders of magnitude faster

Multi-Dendrix completes in **fewer than 5 seconds** for each dataset.

*Miller *et al*. (2011). *BMC medical genomics*, **4**:34.

# Multi-Dendrix pipeline for identifying mutated cancer pathways

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Multi-Dendrix pipeline for identifying mutated cancer pathways

*0. Data preprocessing*

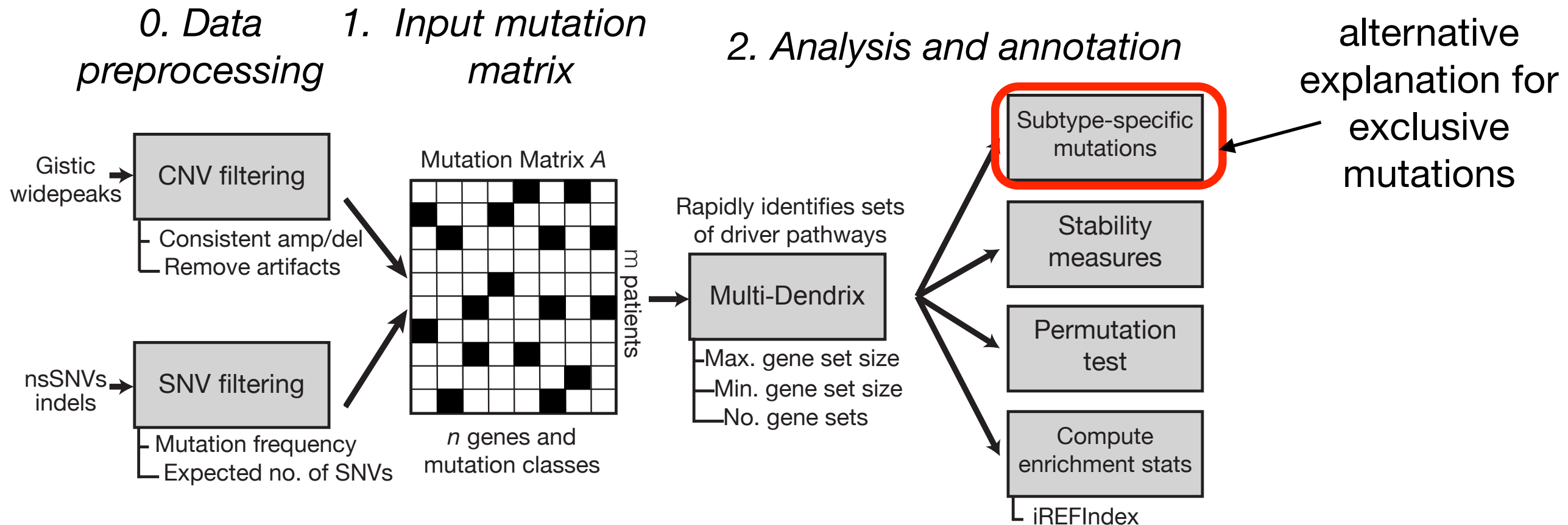# Multi-Dendrix pipeline for identifying mutated cancer pathways



*0. Data preprocessing*   *1. Input mutation matrix*

# Multi-Dendrix pipeline for identifying mutated cancer pathways



**0. Data preprocessing**

**1. Input mutation matrix**

**2. Analysis and annotation**

Gistic widepeaks → CNV filtering
- Consistent amp/del
- Remove artifacts

nsSNVs indels → SNV filtering
- Mutation frequency
- Expected no. of SNVs

Mutation Matrix $A$

$m$ patients

$n$ genes and mutation classes

Rapidly identifies sets of driver pathways

Multi-Dendrix
- Max. gene set size
- Min. gene set size
- No. gene sets

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Multi-Dendrix pipeline for identifying mutated cancer pathways

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Multi-Dendrix pipeline for identifying mutated cancer pathways

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

# Multi-Dendrix pipeline for identifying mutated cancer pathways



Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.

16

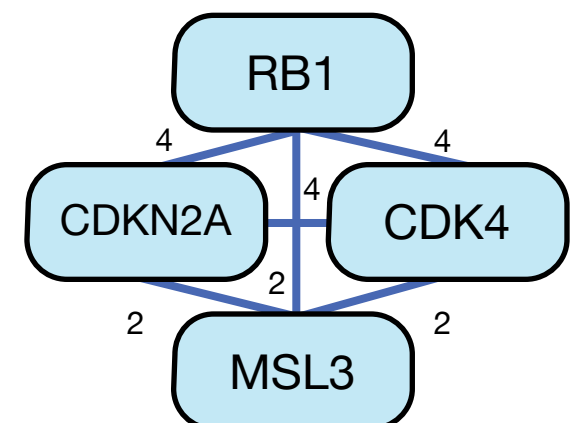# Multi-Dendrix pipeline for identifying mutated cancer pathways



**Direct Interactions Test**

- Measures enrichment of PPI interactions within individual gene sets or in a collection of gene sets

**Consensus modules**

- Run Multi-Dendrix across a range of parameters.
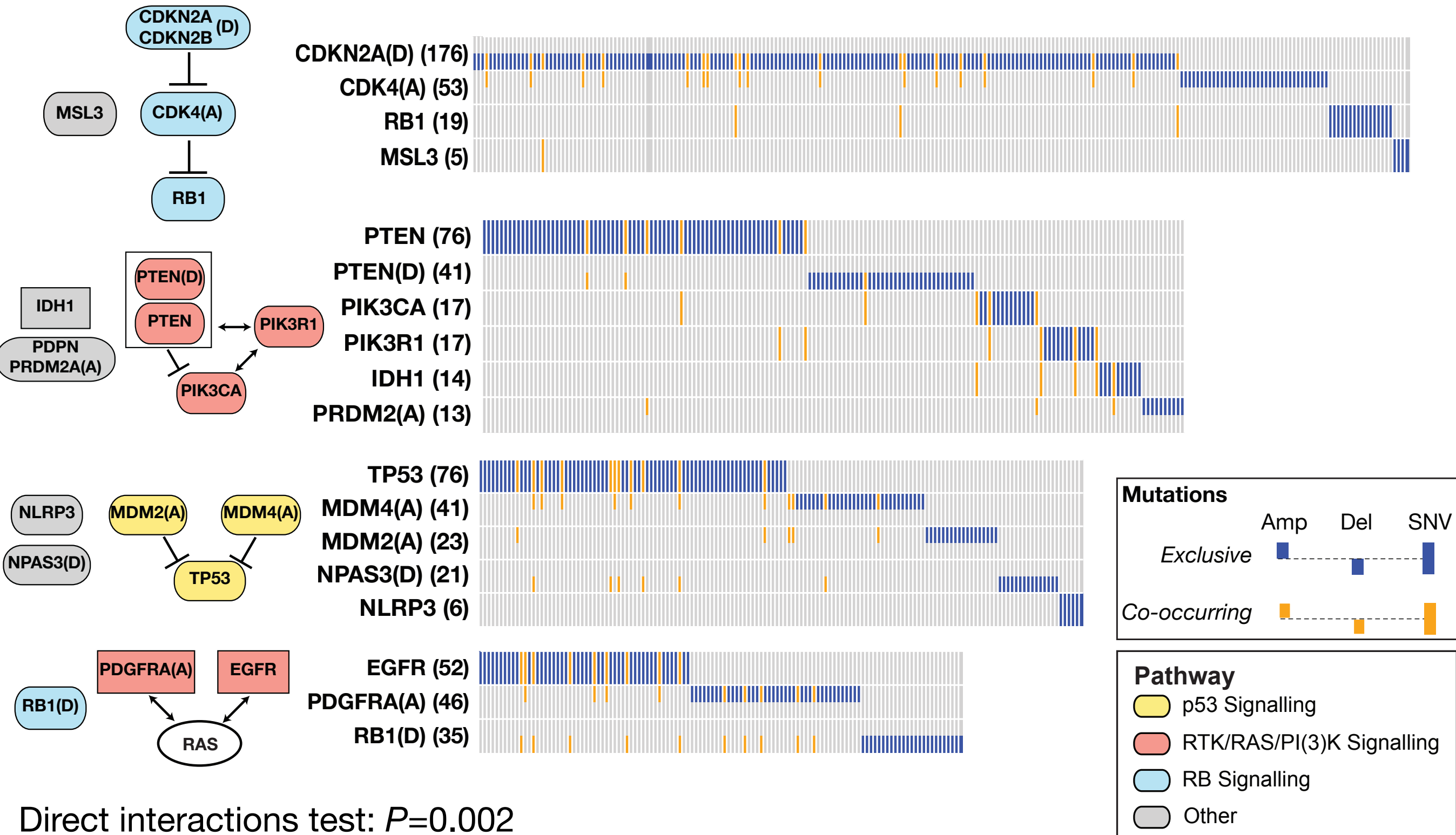- Identify the stable "modules" of genes that appear together multiple times.

Leiserson *et al*. (2013) Simultaneous Identification of Multiple Driver Pathways in Cancer. *PLoS Comp Bio* 9(**5**): e1003054.    16

# Contributions

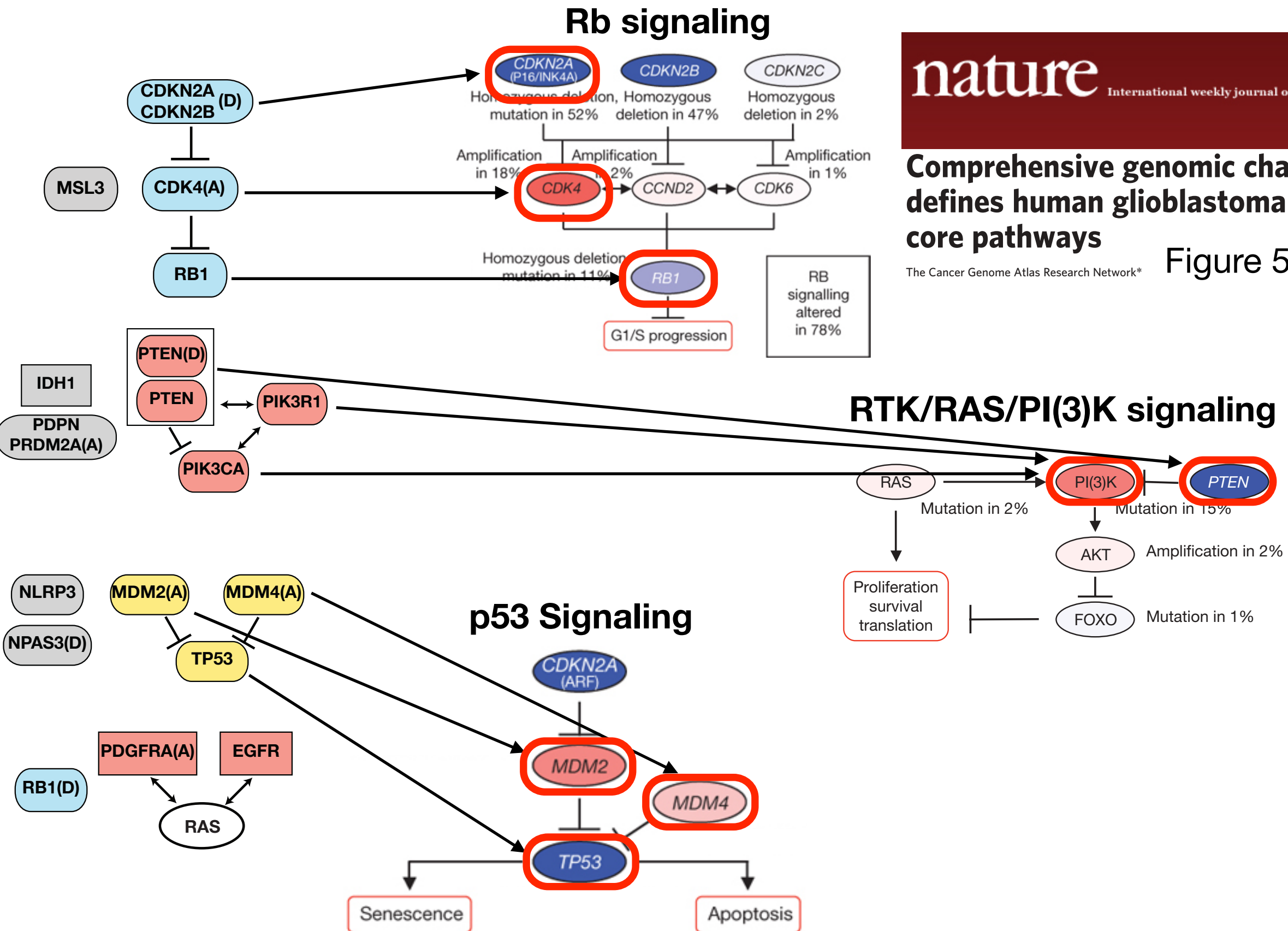A new algorithm, **Multi-Dendrix**, for identifying driver pathways *de novo:*

1. Outperforms previous methods on simulated data in speed and accuracy

2. **Identifies gene sets that overlap known pathways in TCGA datasets**

3. Ongoing work extending Multi-Dendrix to large datasets and overlapping pathways

# Results: Glioblastoma



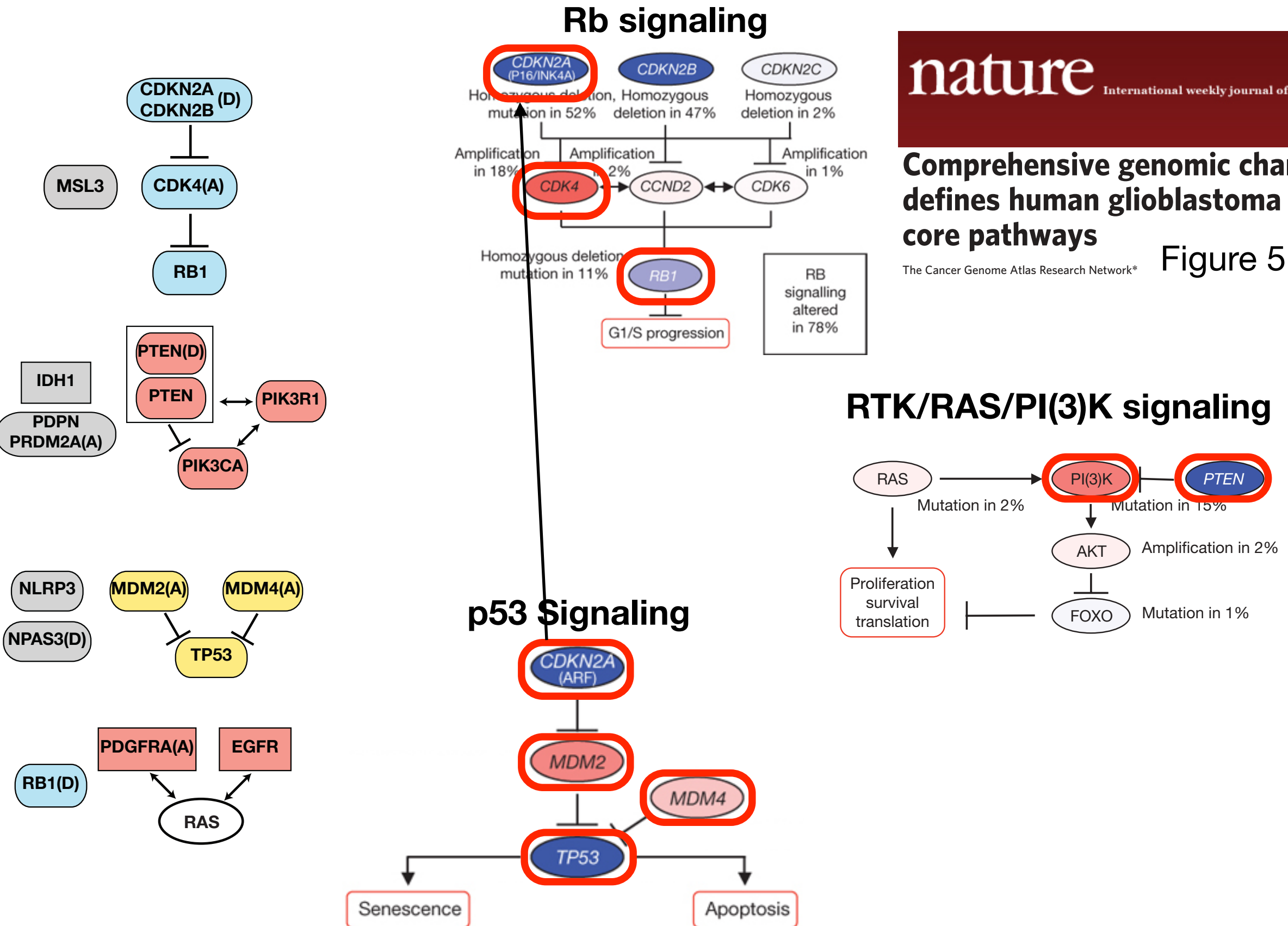*Mutation data:* 398 genes (events) in 261 patients

Direct interactions test: *P*=0.002

*TCGA (2008) Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature* **455**:1061-1088.

# Results: Glioblastoma



**Rb signaling**

**RTK/RAS/PI(3)K signaling**

**p53 Signaling**

Comprehensive genomic characterization defines human glioblastoma genes and core pathways
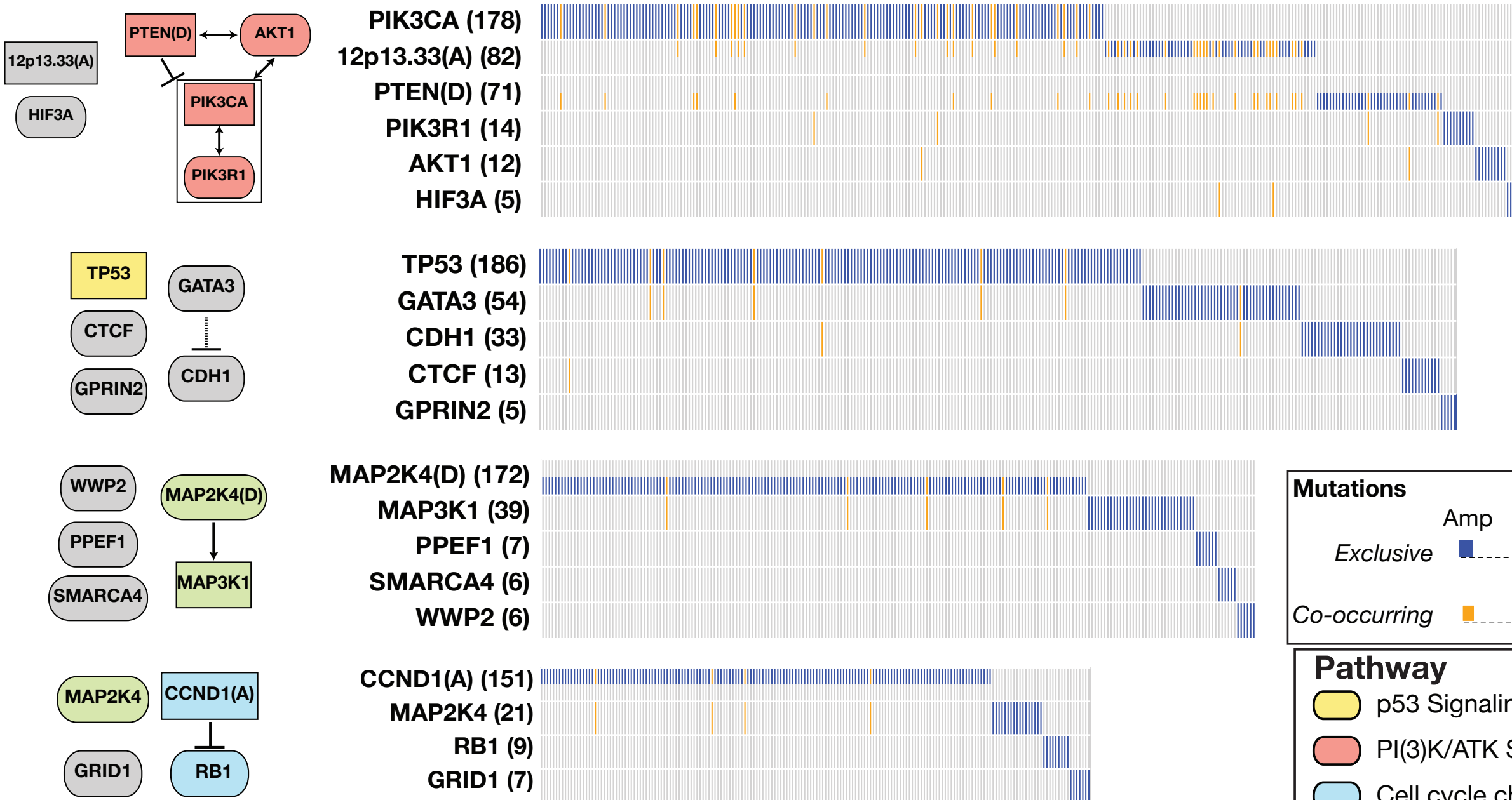
The Cancer Genome Atlas Research Network*

Figure 5

*TCGA (2008) Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature* **455**:1061-1088.

19

# Results: Glioblastoma



**Rb signaling**

**p53 Signaling**

**RTK/RAS/PI(3)K signaling**

**Comprehensive genomic characterization defines human glioblastoma genes and core pathways**

The Cancer Genome Atlas Research Network*

Figure 5

# Results: TCGA breast cancer

$$W'_\alpha(\mathbf{M}) = \sum_{M \in \mathbf{M}'} \Gamma(M) - \alpha\omega(M)$$

higher $\alpha \Rightarrow$ more exclusivity

**Multi-Dendrix (α=2.5)**

*Mutation data:* 375 genes (events) in 507 patients

Direct interactions test: $P < 0.001$

*TCGA. (2012) Comprehensive molecular portraits of human breast tumors. *Nature*, **490**:61-70.

20

# Contributions

A new algorithm, **Multi-Dendrix**, for identifying driver pathways *de novo:*

1. Outperforms previous methods on simulated data in speed and accuracy

2. Identifies gene sets that overlap known pathways in TCGA datasets

3. Ongoing work extending Multi-Dendrix to large datasets and overlapping pathways

# Multi-Dendrix MCMC

**Sample collections of gene sets in proportion to their weight**
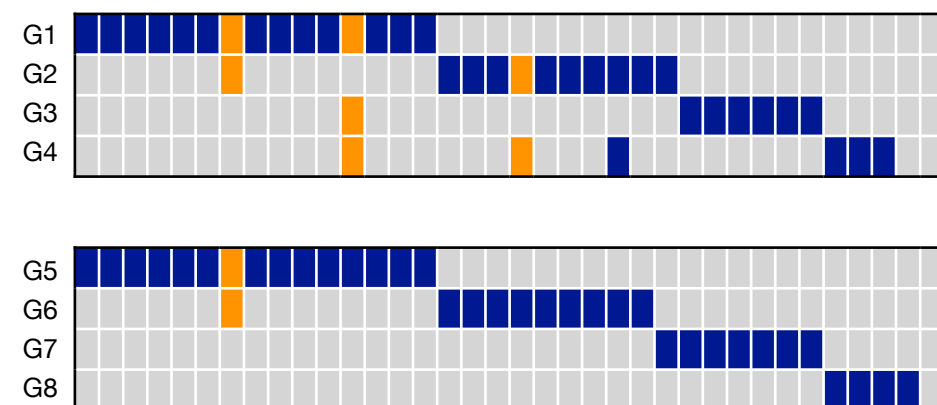


*Distribution of collections of gene sets*

Frequency — Weight (ascending)

High — Sampling frequency — Low

- Finds a distribution — optimal and suboptimal — of solutions

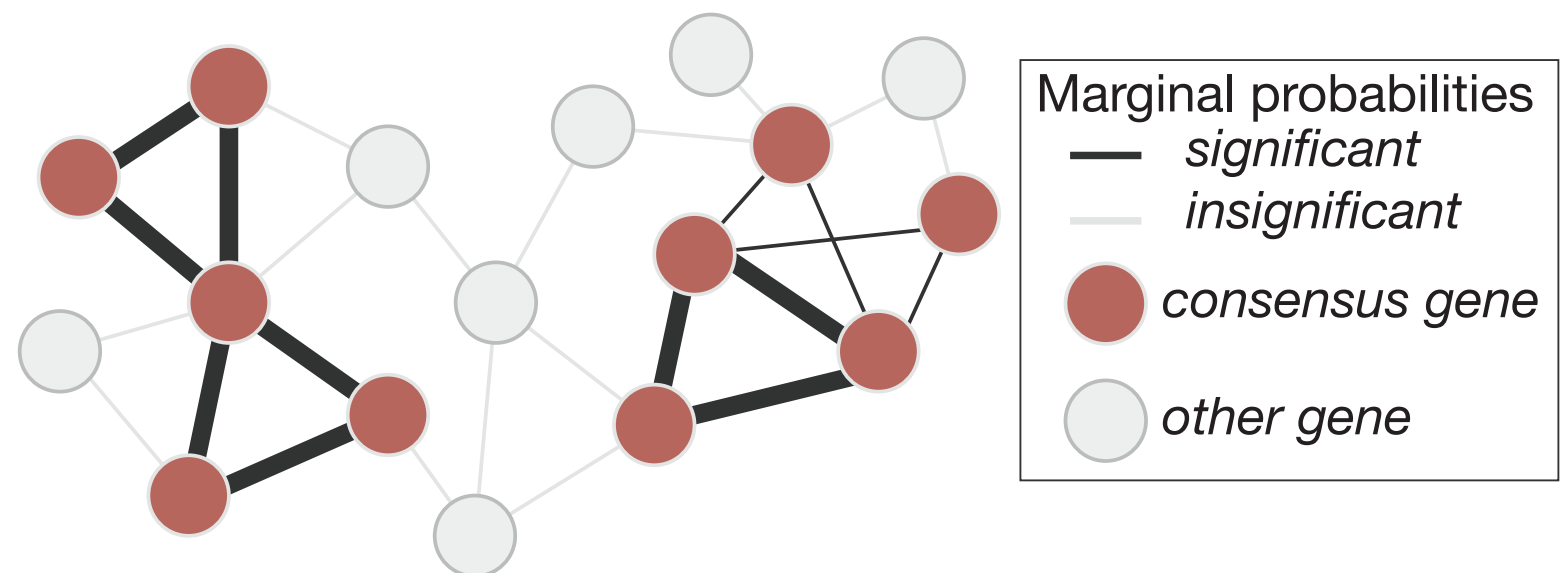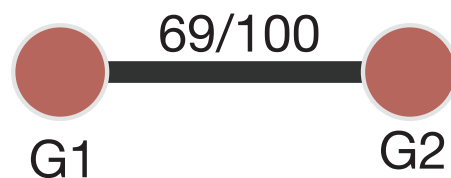- Newest version on GitHub, requirements are all open-source

*Sampling table*

| Geneset | Weight | Sampling Frequency |
|---|---|---|
| G1, G2, G3, G4 G5, G6, G7, G8 | 200 | 60 |
| G9, G10, G11, G12 G13, G14, G15, G16 | 170 | 30 |
| G1, G2, G3, G4 G5, G7, G15, G16 | 130 | 9 |
| G17, G18, G19, G20 G13, G14, G15, G16 | 80 | 1 |

*Top collection*



*Weight: 200*

# Marginal probability graph

- Marginal probability graph defines consensus subnetworks

- Edges $(u, v)$ are weighted by how often gene $u$ is sampled in the same gene set as gene $v$

| Geneset | Weight | Sampling Frequency |
|---|---|---|
| **G1**, **G2**, G3, G4 G5, G6, G7, G8 | 200 | 60 |
| **G1**, **G2**, G3, G4 G5, G7, G15, G16 | 130 | 9 |



*Complete, weighted marginal probability graph*

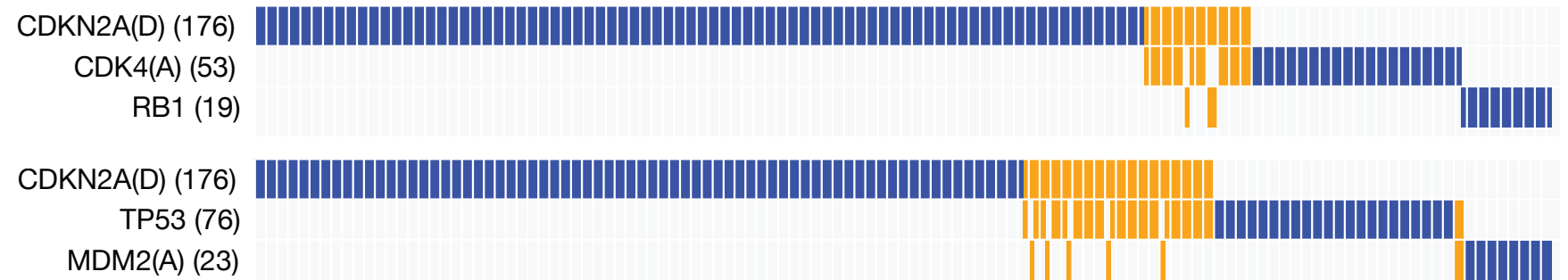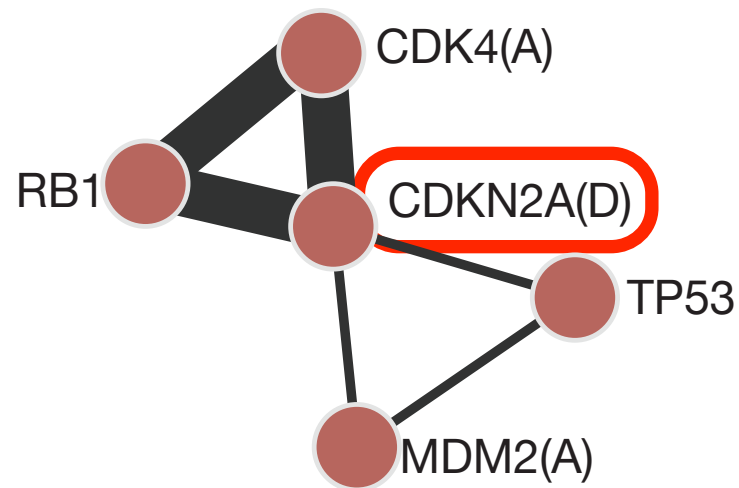→ unconstrained size and number of gene sets
→ gene sets can overlap

# Marginal probability graph

- Marginal probability graph defines consensus subnetworks

- Edges $(u, v)$ are weighted by how often gene $u$ is sampled in the same gene set as gene $v$

| Geneset | Weight | Sampling Frequency |
|---------|--------|--------------------|
| **G1**, **G2**, G3, G4 G5, G6, G7, G8 | 200 | 60 |
| **G1**, **G2**, G3, G4 G5, G7, G15, G16 | 130 | 9 |



Marginal probabilities
— *significant*
— *insignificant*
⬤ *consensus gene*
◯ *other gene*

69/100

G1 — G2

Hsin-Ta Wu

*Complete, weighted marginal probability graph*

→ unconstrained size and number of gene sets
→ gene sets can overlap

Leiserson, Wu *et al*. Dendrix++. (In preparation).

# Dendrix++ results: GBM

**Comprehensive genomic characterization defines human glioblastoma genes and core pathways**
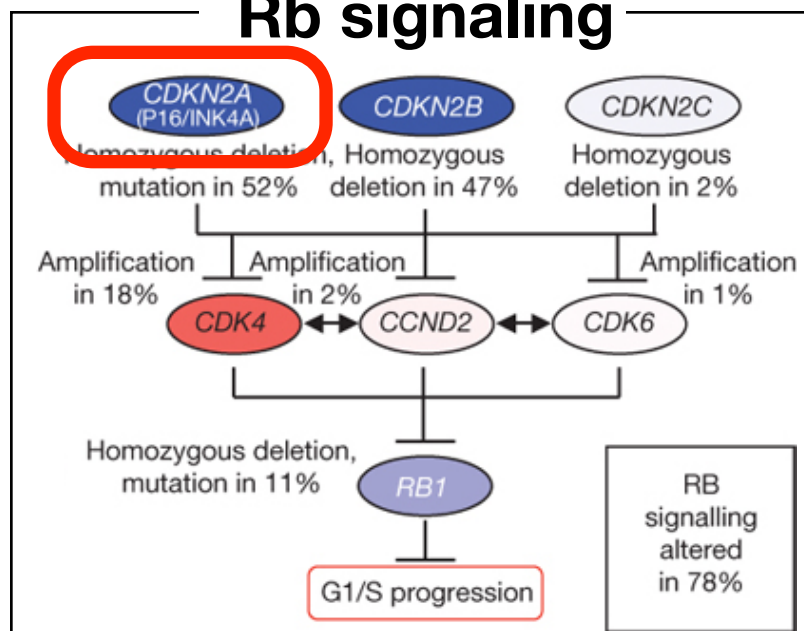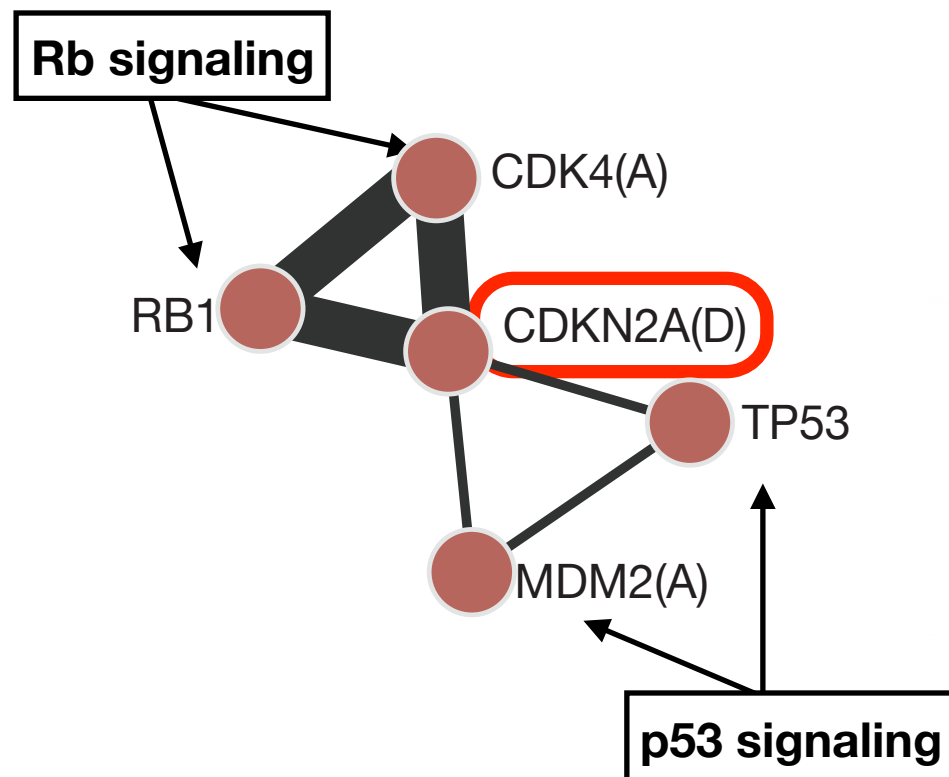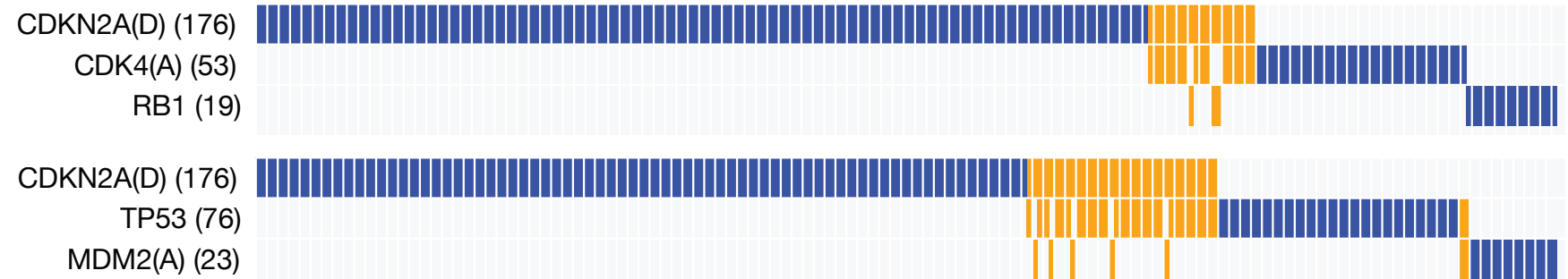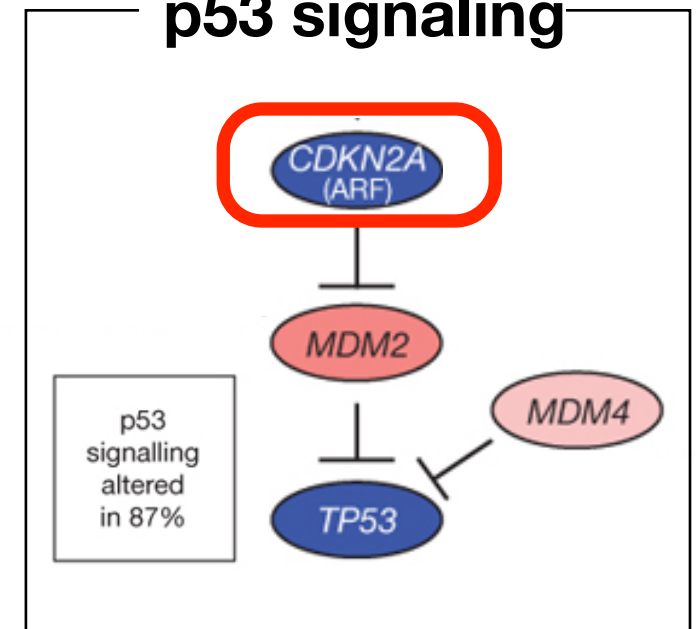
The Cancer Genome Atlas Research Network*



Leiserson, Wu *et al*. Dendrix++. (In preparation).
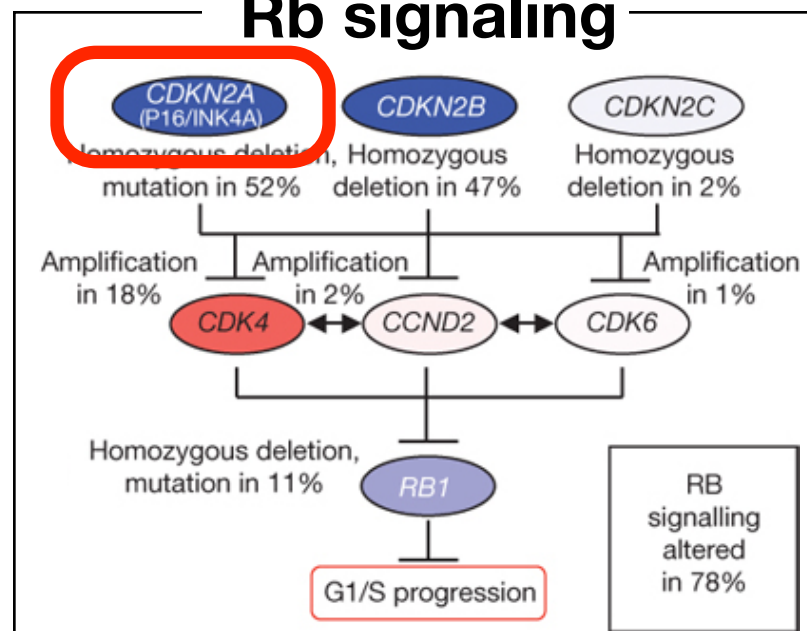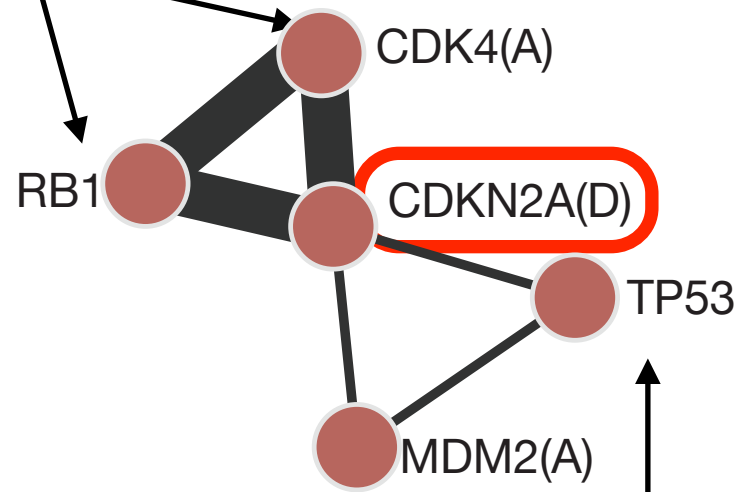
# Dendrix++ results: GBM



**Comprehensive genomic characterization defines human glioblastoma genes and core pathways**

The Cancer Genome Atlas Research Network*

Leiserson, Wu *et al*. Dendrix++. (In preparation).

# Dendrix++ results: GBM



**Comprehensive genomic characterization defines human glioblastoma genes and core pathways**

The Cancer Genome Atlas Research Network*

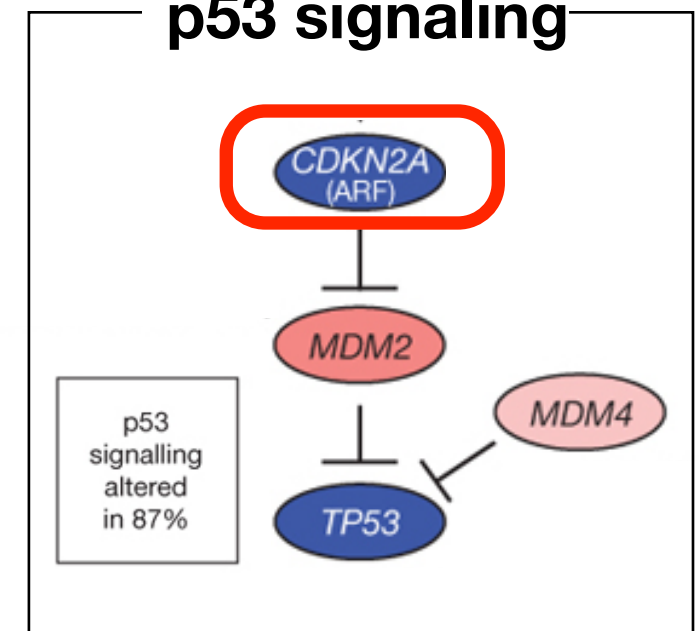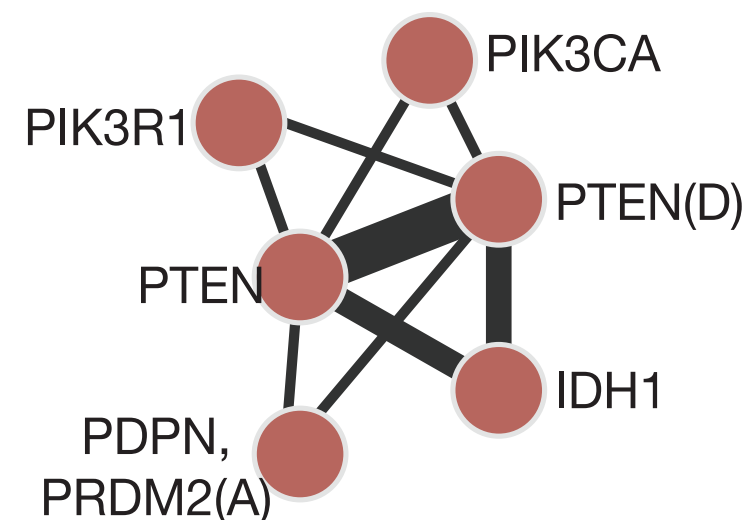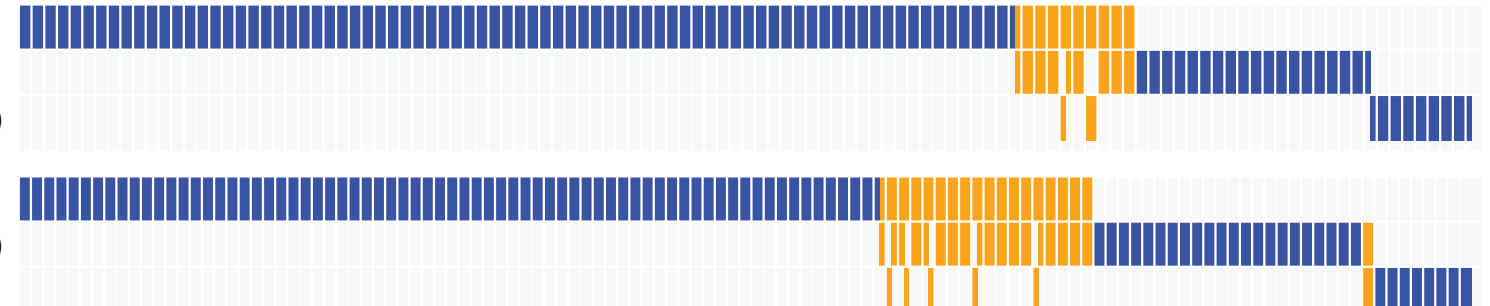Leiserson, Wu *et al*. Dendrix++. (In preparation).

# Summary

- Multi-Dendrix: Fast, exact ILP for identifying collections of gene sets with exclusive mutations

- Identifies modules that overlap multiple cancer pathways in glioblastoma and breast cancer

- Dendrix++: New algorithm that can identify more complex pathways (in preparation)

# Acknowledgements

## Research group

**Ben Raphael**
Jason Dobson
Iman Hajirasouliha
*Fabio Vandin*
Layla Oesper
*Hsin-Ta Wu*

BROWN

## Collaborators

**Roded Sharan**
**Dima Blokh**

אוניברסיטת תל-אביב
TEL AVIV UNIVERSITY

## Funding + Data

National Science Foundation
WHERE DISCOVERIES BEGIN

The Cancer Genome Atlas

NATIONAL INSTITUTES OF HEALTH