

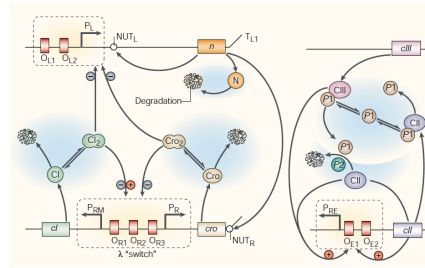
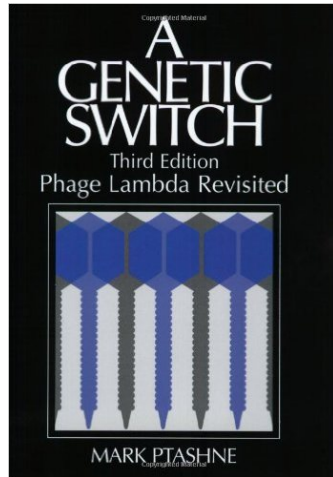
# Challenges in Integrative Modeling of Gene Regulation, Protein Signaling, and Metabolism



## Dynamic Modeling

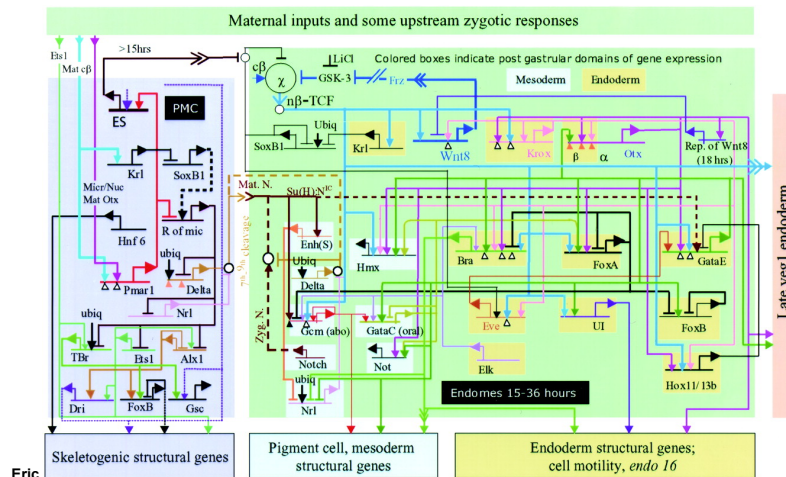
**STEP 0: Define the system**

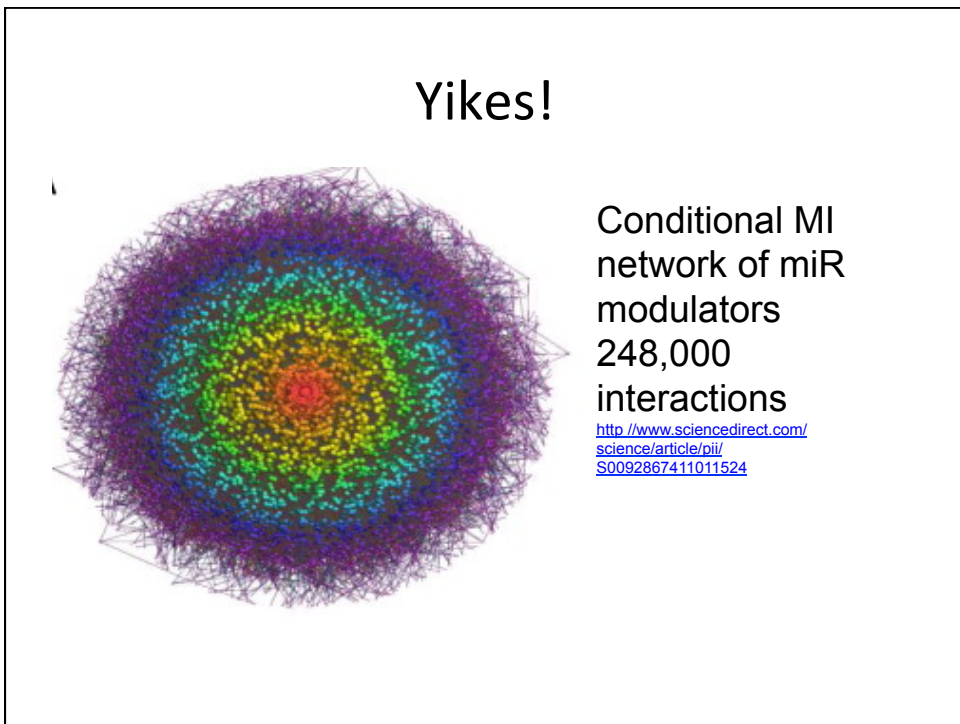
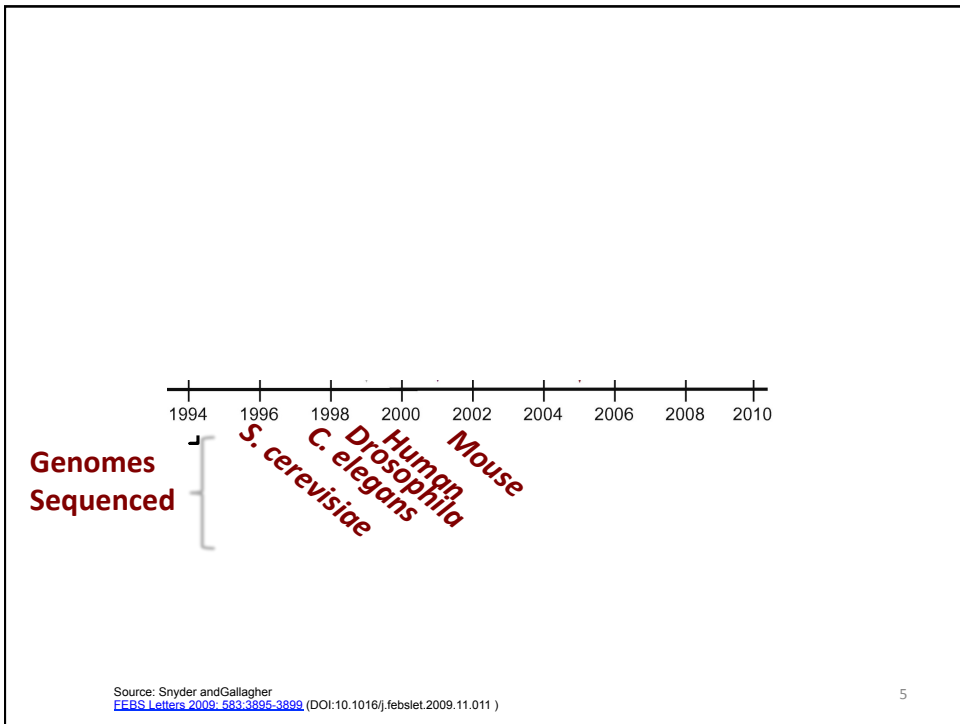
# The Genetic Switch (Then)

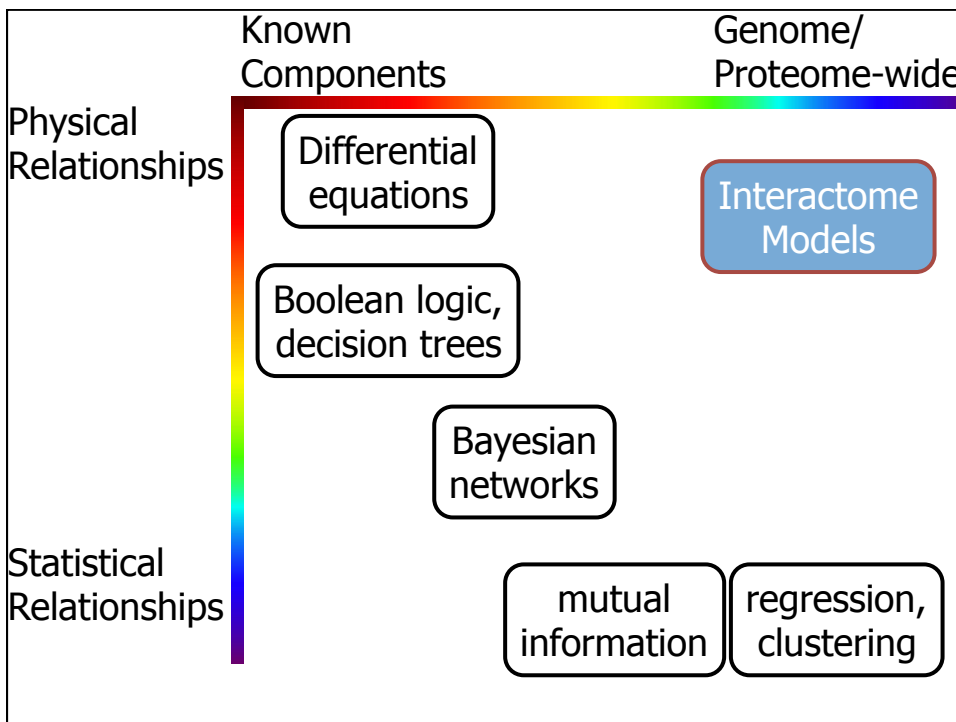
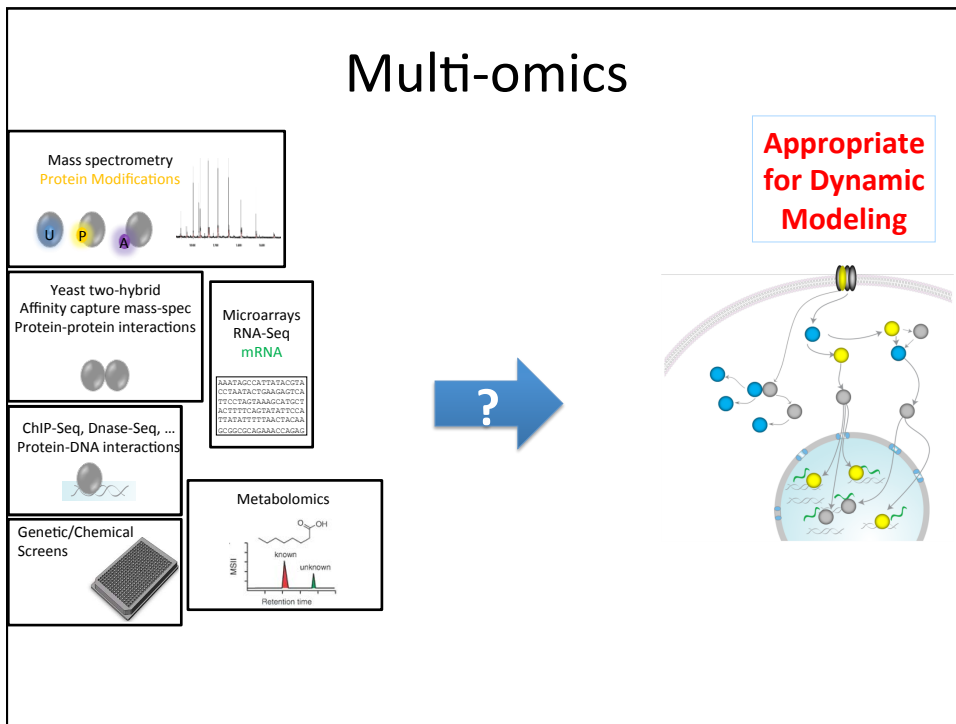


Nature Reviews Genetics 4, 471-477 (June 2003) | doi:10.1038/nrg1089

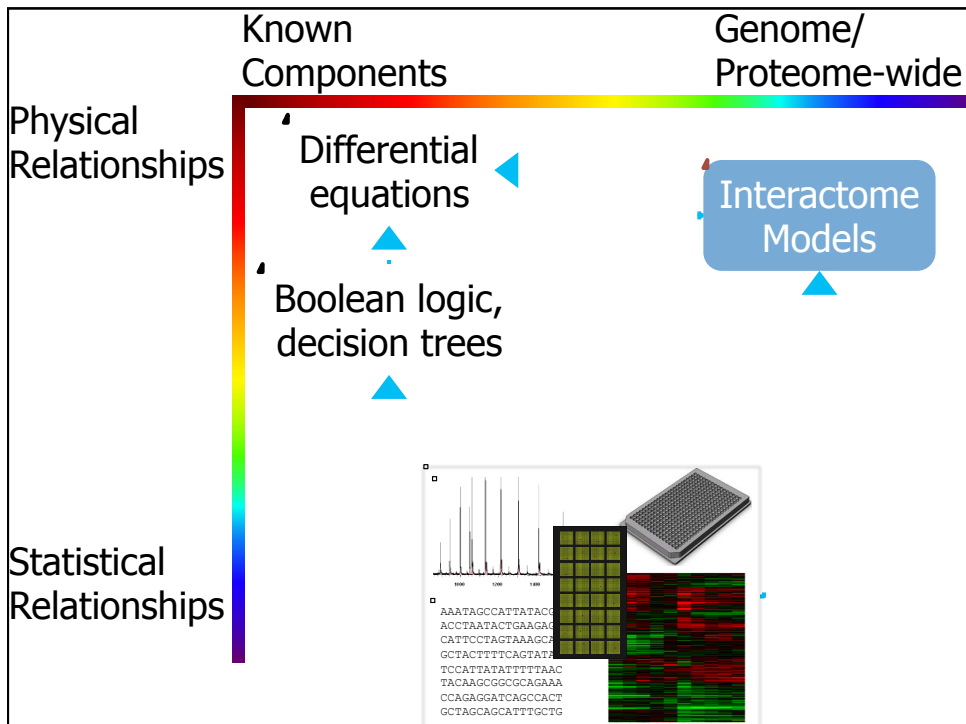
# Even Hand-built Models are Getting Complex!





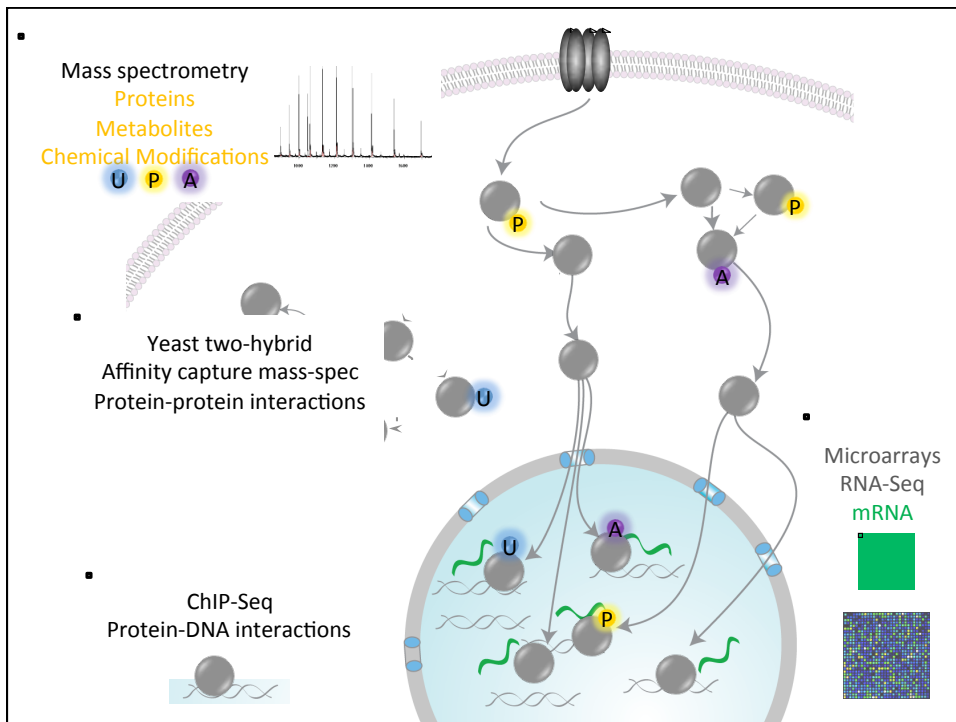
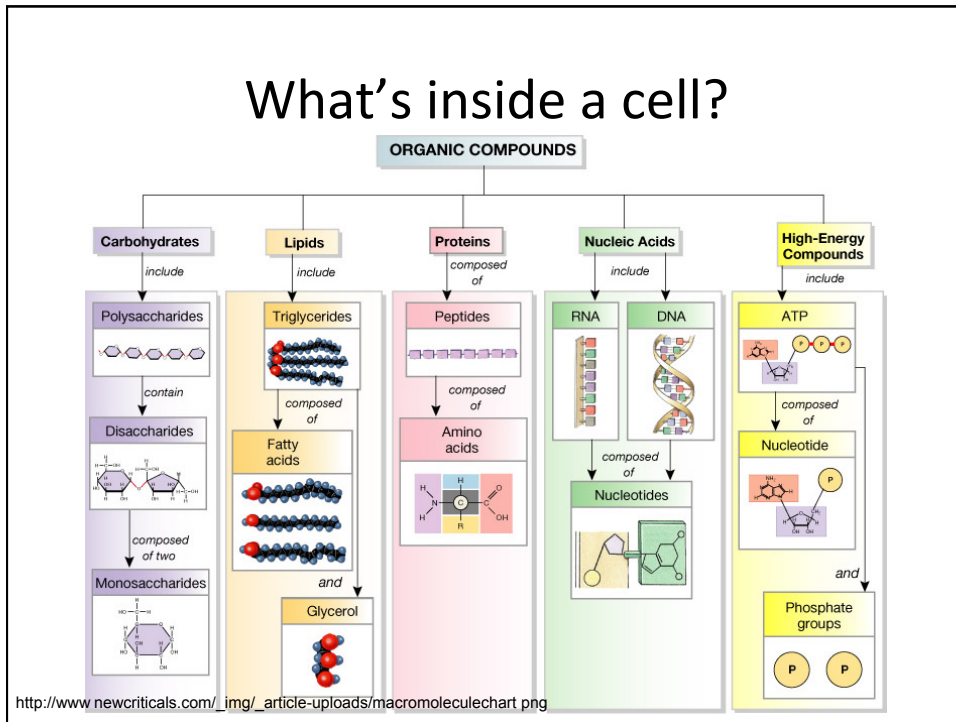




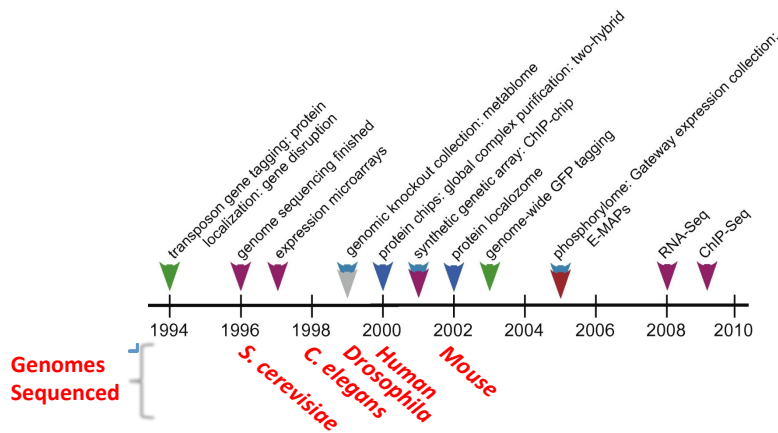


## Outline

- What molecules we can measure?
- How do we know which interact?
- How do we learn anything from these data?
  - Standard Approaches
  - Challenges
  - Network Methods
  - Toward Dynamic Models



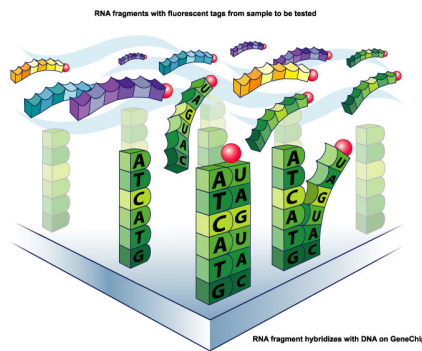
# Timeline of Omic Methods



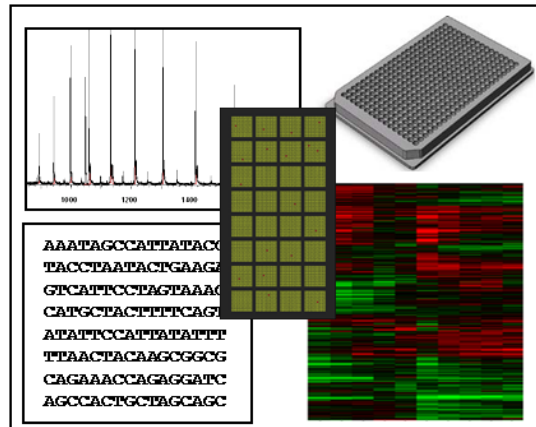
Source Snyder andGallagher  
[FEBS Letters 2009, 583, 3895-3899](https://doi.org/10.1016/j.febslet.2009.11.011) (DOI 10.1016/j.febslet.2009.11.011)

13

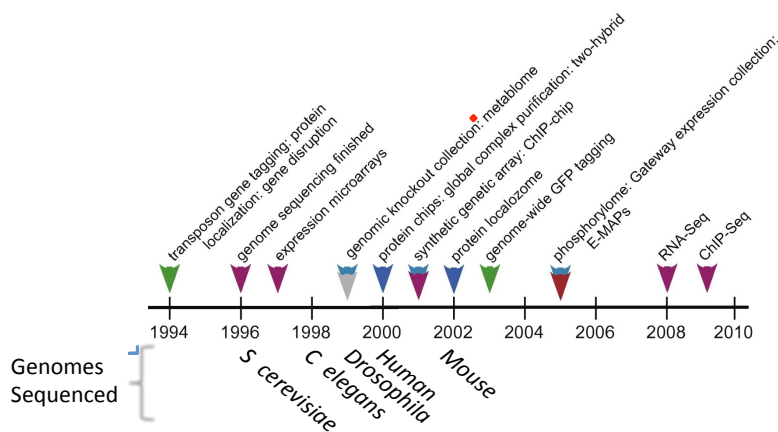
# mRNA levels



# What else can we measure?



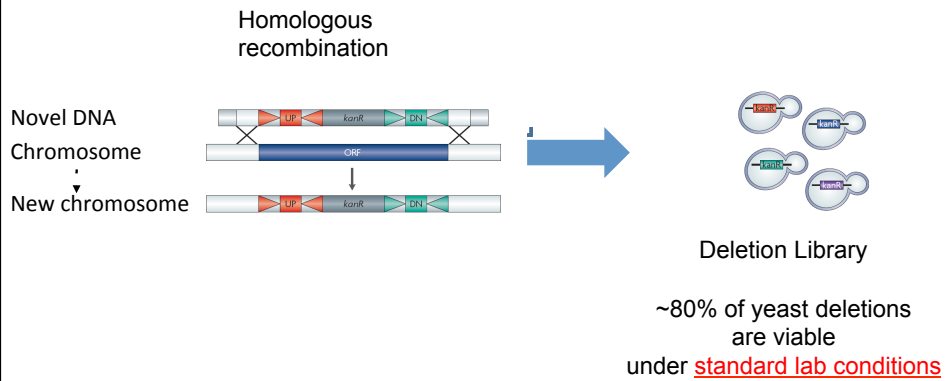
15



Source: Snyder and Gallagher  
[FEBS Letters 2009; 583:3895-3899](https://doi.org/10.1016/j.febslet.2009.11.011) (DOI:10.1016/j.febslet.2009.11.011)

16

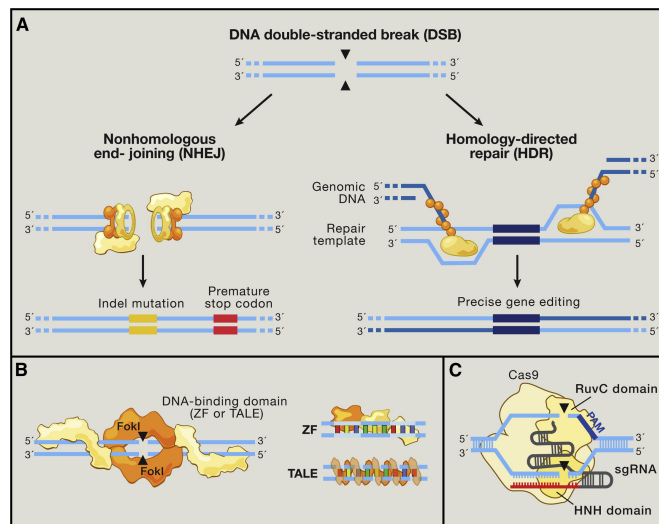
# Homologous Recombination



Boone et al. (2007) Nature Reviews Genetics

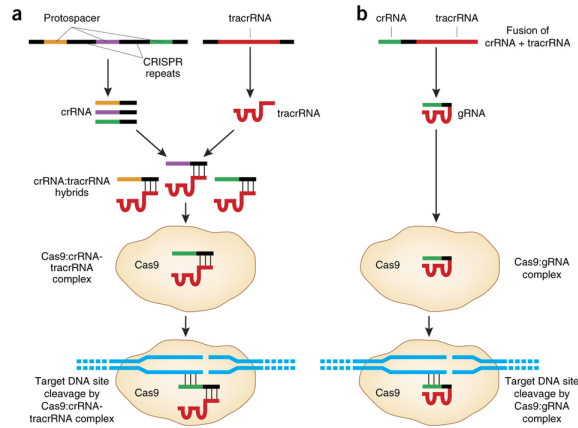
17

# Strategies for Genome Editing



Cell 2014 157, 1262-1278DOI: (10.1016/j.cell.2014.05.010)  
Copyright © 2014 Elsevier Inc. [Terms and Conditions](#)

# CRISPR/Cas9 is an RNA-guided DNA nuclease

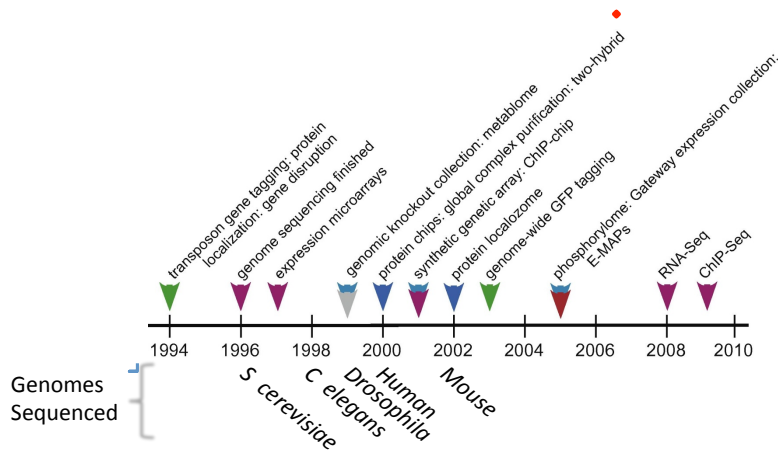


**CRISPR-Cas systems for editing, regulating and targeting genomes**

Jeffrey D Sander & J Keith Joung

*Nature Biotechnology* 32, 347–355 (2014) doi:10.1038/nbt.2842

19



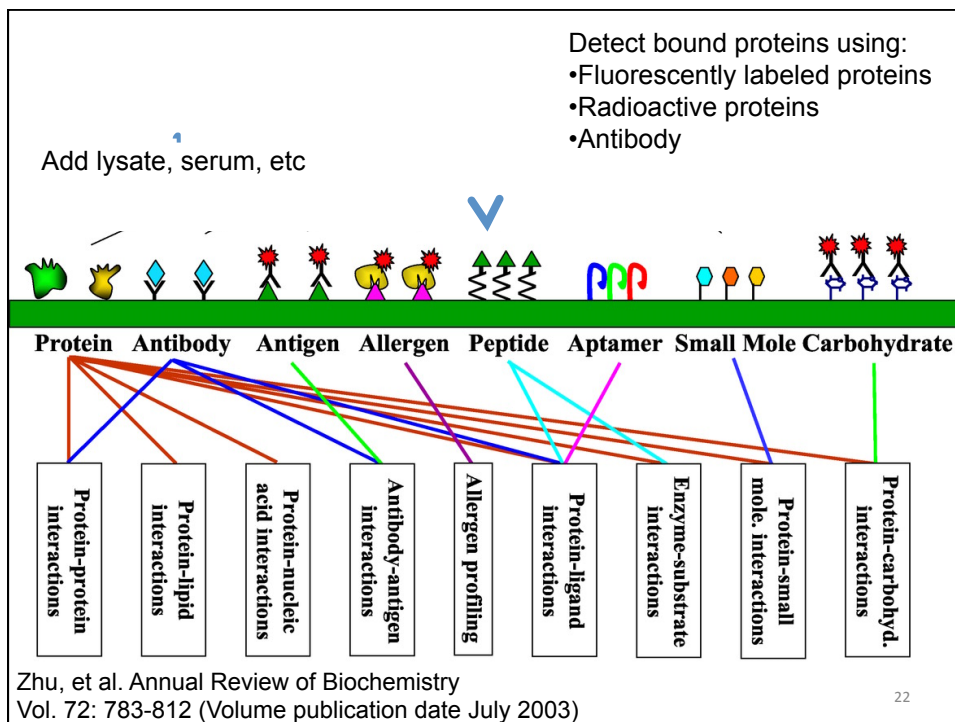
Source: Snyder and Gallagher  
 FEBS Letters 2009; 583:3895-3899 (DOI:10.1016/j.febslet.2009.11.011)

20

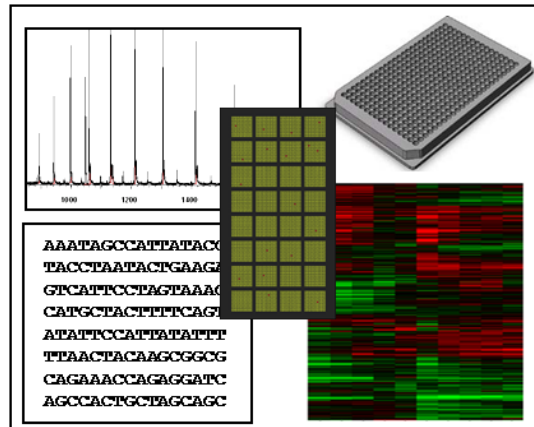
# Proteomic Arrays

- Can we extend the DNA array platform for proteomics?
  - What would you put on the array?
  - What would you use to probe the array?
  - How would you detect it?
  - What are the technical challenges?

21



## What else can we measure?



23

## Post-translational modifications

<a href="#">Acetylation</a>	<a href="#">Diphthamide</a>	<a href="#">S-palmitoleyl cysteine</a>
<a href="#">ADP-ribosylation</a>	<a href="#">FAD</a>	<a href="#">Phosphatidylethanolamine amidated glycine</a>
<a href="#">Allysine</a>	<a href="#">S-farnesyl cysteine</a>	<a href="#">Phosphorylation</a>
<a href="#">Amidation</a>	<a href="#">S-12-hydroxyfarnesyl cysteine</a>	<a href="#">Pyridoxal phosphate</a>
<a href="#">S-archaeol</a>	<a href="#">3-phenyllactic acid</a>	<a href="#">N6-poly(methylaminopropyl)lysine</a>
<a href="#">Beta-methylthiolation</a>	<a href="#">FMN conjugation (Cys)</a>	<a href="#">Phosphopantetheine</a>
<a href="#">Biotin</a>	<a href="#">FMN conjugation (Ser/Thr)</a>	<a href="#">Pyrrolidone carboxylic acid (Glu)</a>
<a href="#">Bromination</a>	<a href="#">FMN conjugation (His)</a>	<a href="#">Pyrrolysine</a>
<a href="#">N6-1-carboxyethyl lysine</a>	<a href="#">Formylation</a>	<a href="#">Pyrrolidone carboxylic acid</a>
<a href="#">Cholesterol</a>	<a href="#">Geranyl-geranylation</a>	<a href="#">Pyruvic acid (Cys)</a>
<a href="#">Cis-14-hydroxy-10,13-dioxo-7-heptadecenoic acid aspartate ester</a>	<a href="#">Gamma-carboxyglutamic acid</a>	<a href="#">Pyruvic acid (Ser)</a>
<a href="#">Citullination</a>	<a href="#">O-GlcNAc</a>	<a href="#">Sulfation</a>
<a href="#">C-Mannosylation</a>	<a href="#">Glucosylation (Glycation)</a>	<a href="#">1-thioglycine</a>
<a href="#">Cysteine sulfenic acid (-SOH)</a>	<a href="#">Glutathionylation</a>	<a href="#">Thyroxine</a>
<a href="#">Cysteine sulfinic acid (-SO<sub>2</sub>H)</a>	<a href="#">Hydroxylation</a>	<a href="#">2',4',5'-topoquinone</a>
<a href="#">Cysteine persulfide</a>	<a href="#">Hypusine</a>	<a href="#">Triiodothyronine</a>
<a href="#">Deamidation</a>	<a href="#">Lipoyl</a>	<a href="#">Trimethylation</a>
<a href="#">Deamidation followed by a methylation</a>	<a href="#">Methylation</a>	<a href="#">N6,N6,N6-trimethyl-5-hydroxylysine</a>
<a href="#">n-Decanoate</a>	<a href="#">Methionine sulfone</a>	
<a href="#">2,3-didehydroalanine (Ser)</a>	<a href="#">Myristoylation</a>	
<a href="#">2,3-didehydrobutyrine</a>	<a href="#">S-Nitrosylation</a>	
<a href="#">(Z)-2,3-didehydrotyrosine</a>	<a href="#">n-Octanoate</a>	
<a href="#">S-diacylglycerol cysteine</a>	<a href="#">Omega-hydroxyceramide glutamate ester</a>	
<a href="#">Dihydroxylation</a>	<a href="#">3-oxoalanine (Cys)</a>	
<a href="#">Dimethylation</a>	<a href="#">3-oxoalanine (Ser)</a>	
<a href="#">Dimethylation of proline</a>	<a href="#">2-oxobutanoic acid</a>	
	<a href="#">Palmitoylation</a>	

[http://ca.expasy.org/tools/findmod/findmod\\_masses.html](http://ca.expasy.org/tools/findmod/findmod_masses.html)



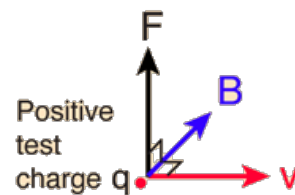


## Mass Spectrometry

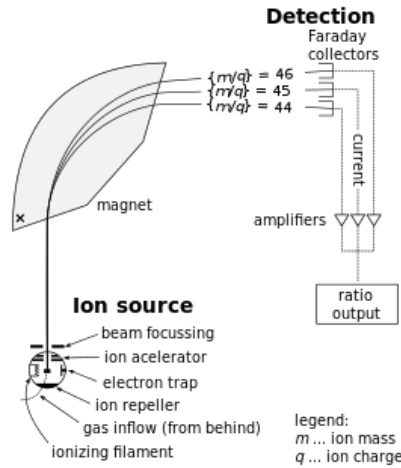
- Can identify proteins and other molecules by extremely accurate measurement of mass

## Force on a charge in a magnetic field

$$\vec{F} = q\vec{v} \times \vec{B}$$

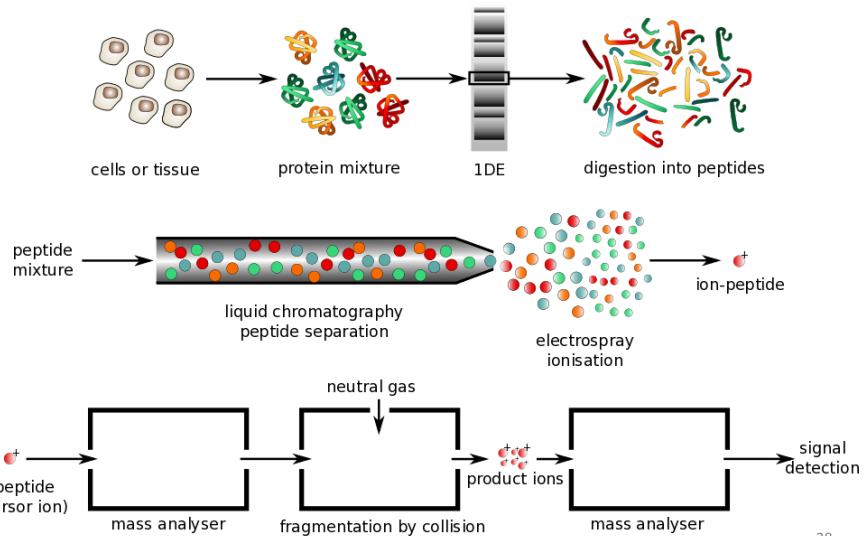


# MS separates by m/q

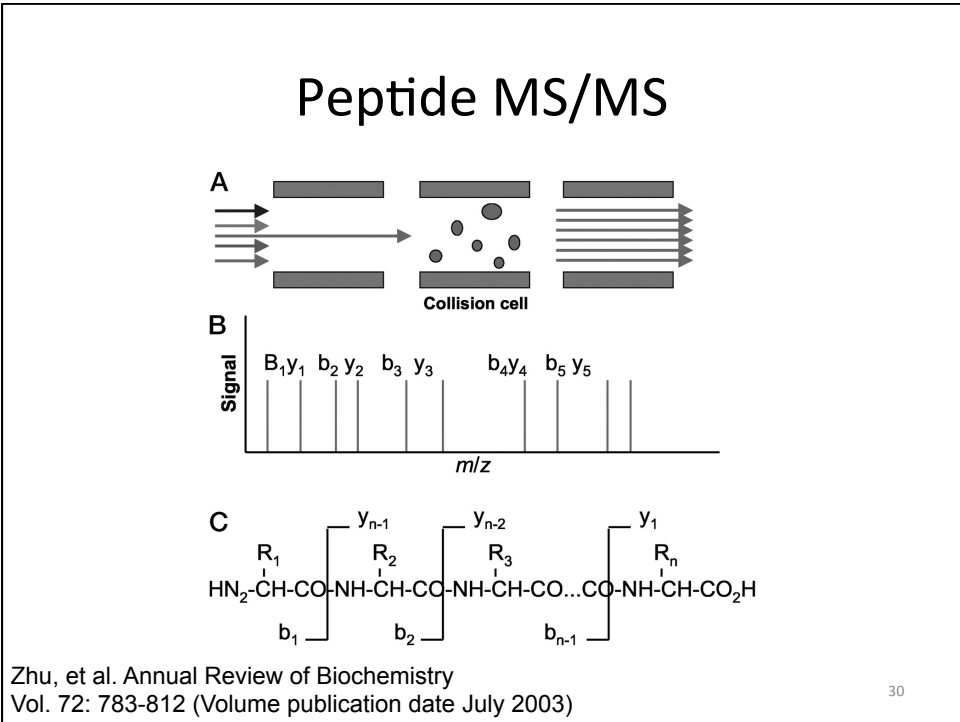
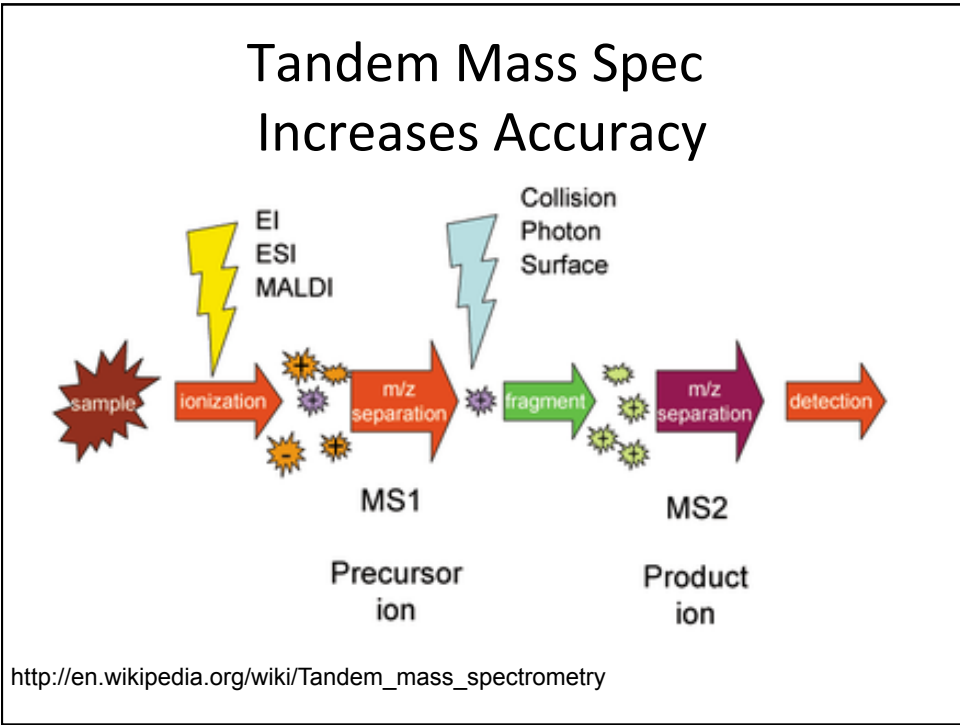


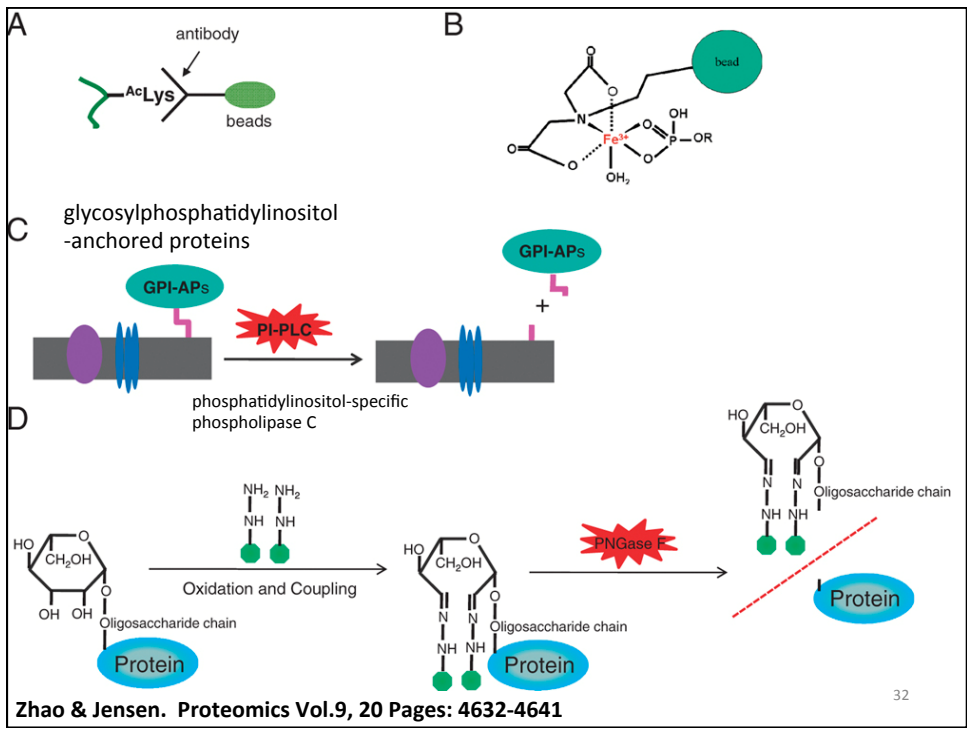
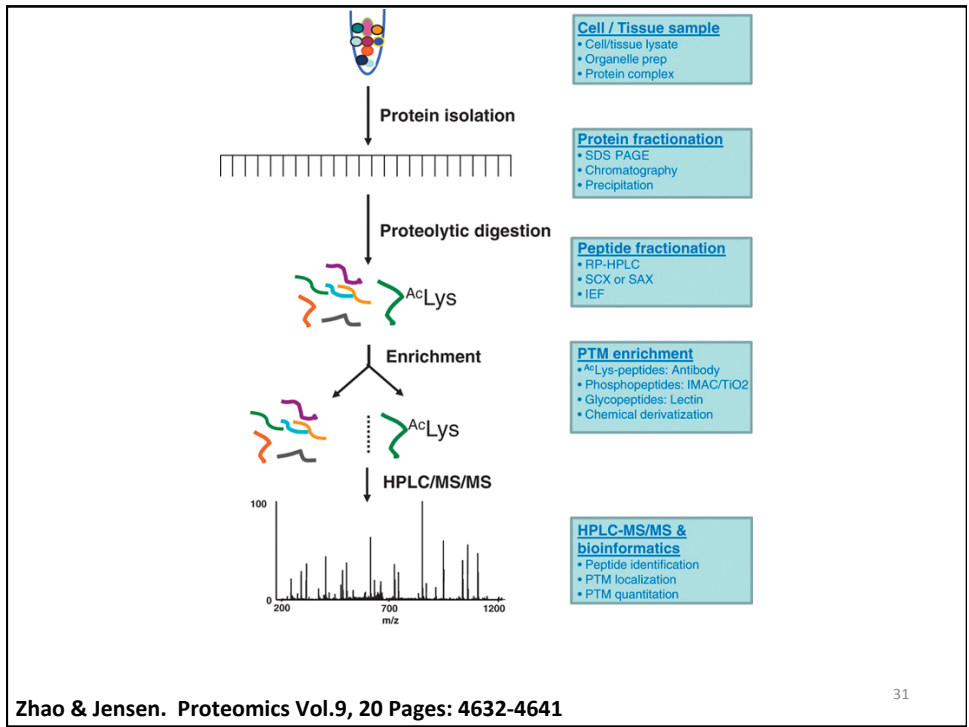
[http://en.wikipedia.org/wiki/Mass\\_spectrometry](http://en.wikipedia.org/wiki/Mass_spectrometry)

# MS cannot analyze very complex samples



[http://upload.wikimedia.org/wikipedia/commons/1/1f/Mass\\_spectrometry\\_protocol.png](http://upload.wikimedia.org/wikipedia/commons/1/1f/Mass_spectrometry_protocol.png)







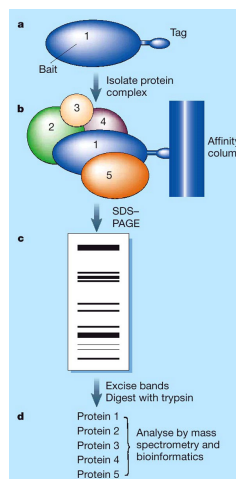
## Outline

- What molecules we can measure
- How do we know which interact?
- How do we learn anything from these data?

## Detecting protein-protein interactions

What are the likely false positives?

What are the likely false negatives?



Gavin, A.-C. *et al. Nature* **415**, 141-147 (2002).

Ho, Y. *et al. Nature* **415**, 180-183 (2002).

[Proteomics: Protein complexes take the bait](#)

Anuj Kumar and Michael Snyder  
*Nature* **415**, 123-124 (10 January 2002)  
 doi:10.1038/415123a

36

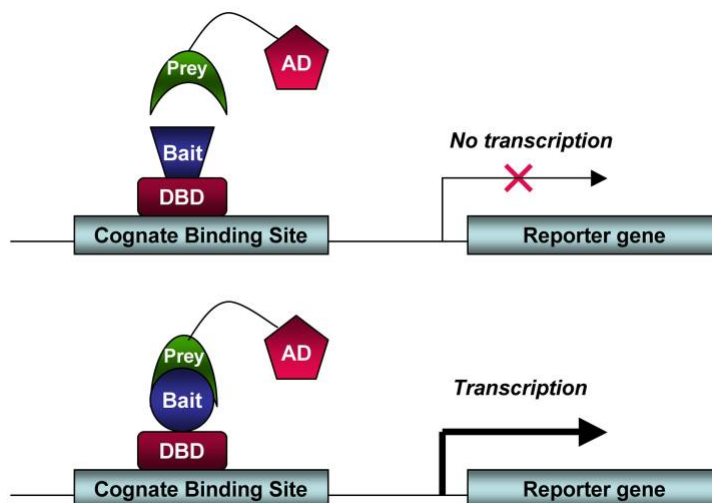
## Mass-spec for protein-protein interactions

- Extremely efficient method for detecting interactions
- Proteins are in their correct subcellular location.

Limitations?

- overexpression/tagging can influence results
- only long-lived complexes will be detected

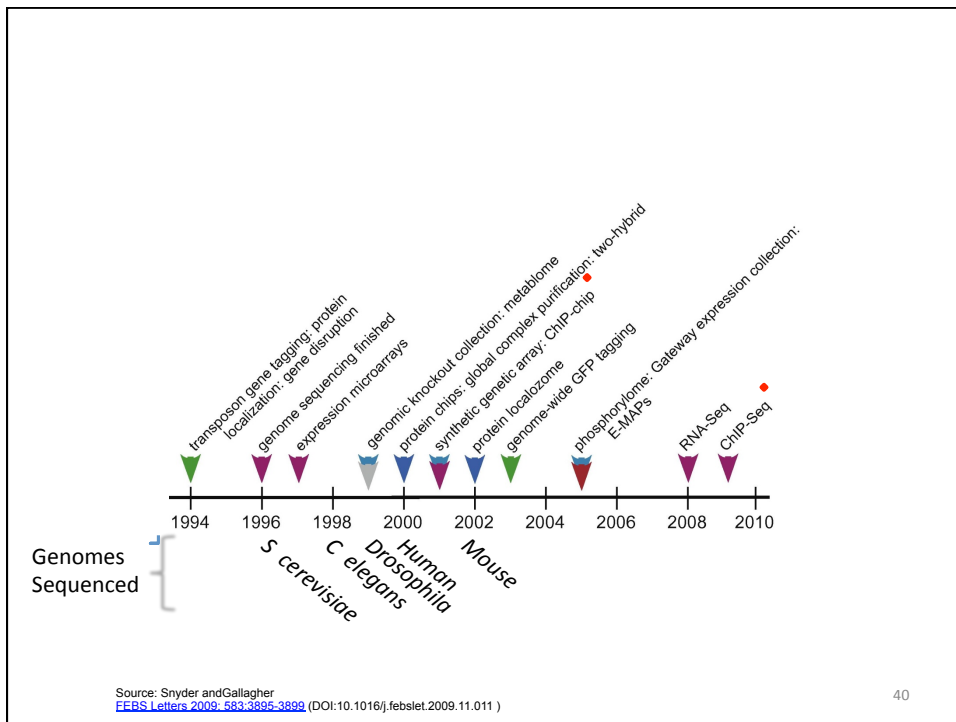
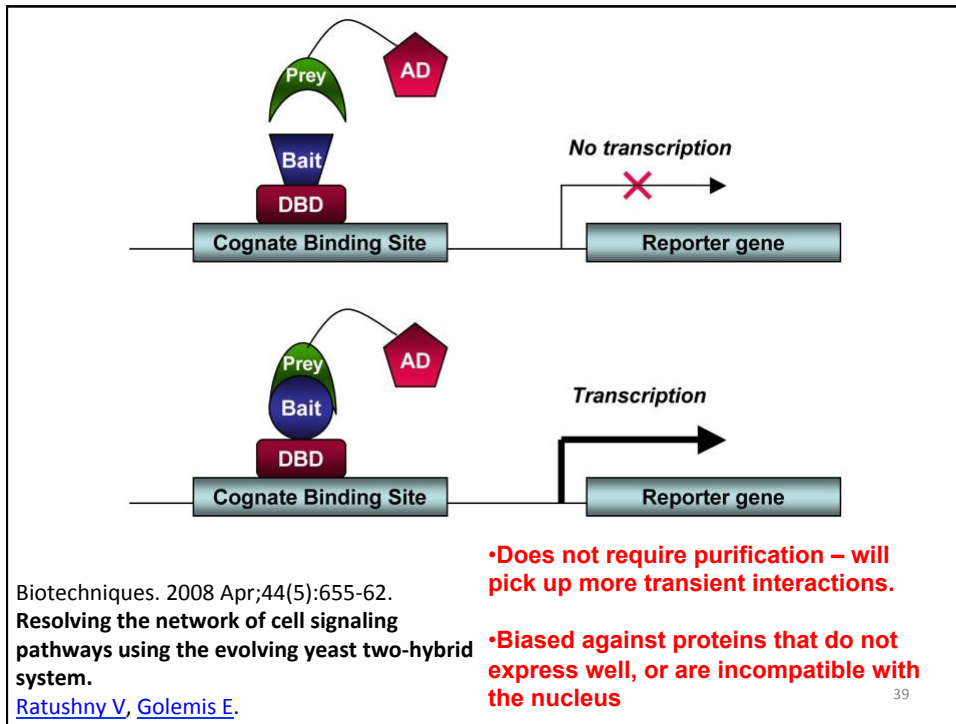
37



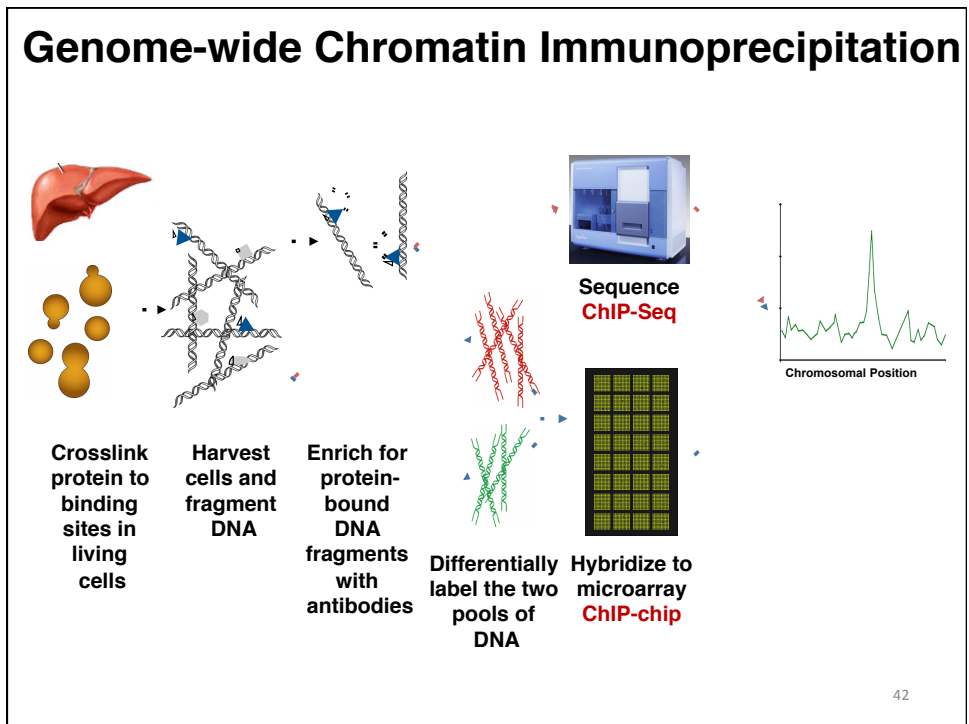
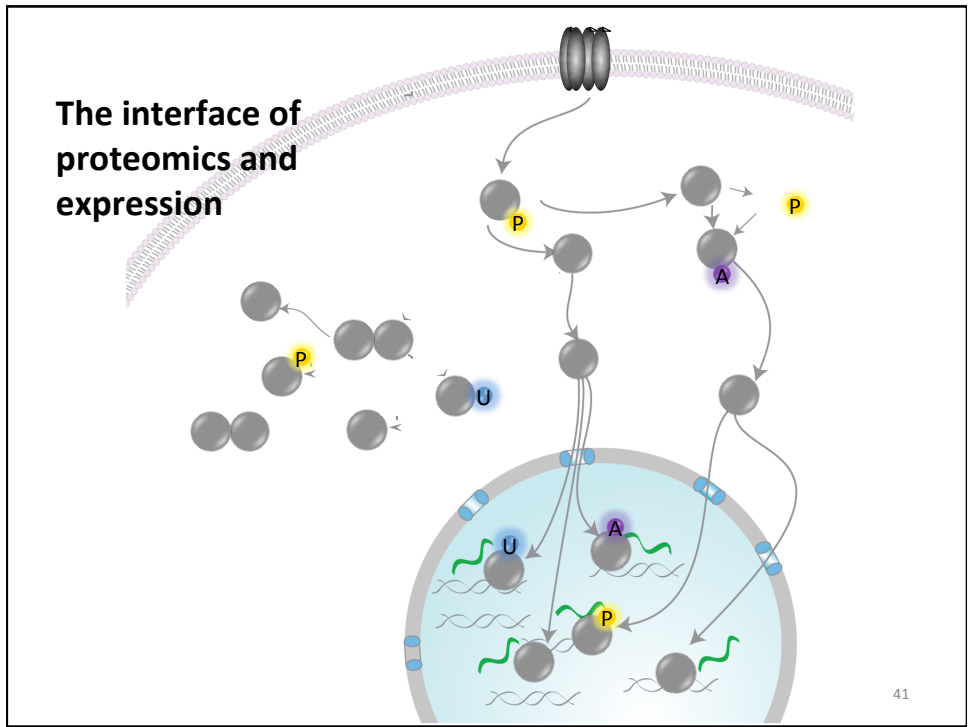
Biotechniques. 2008 Apr;44(5):655-62.  
 Resolving the network of cell signaling pathways using the evolving yeast two-hybrid system.  
[Ratushny V](#), [Golemis E](#).

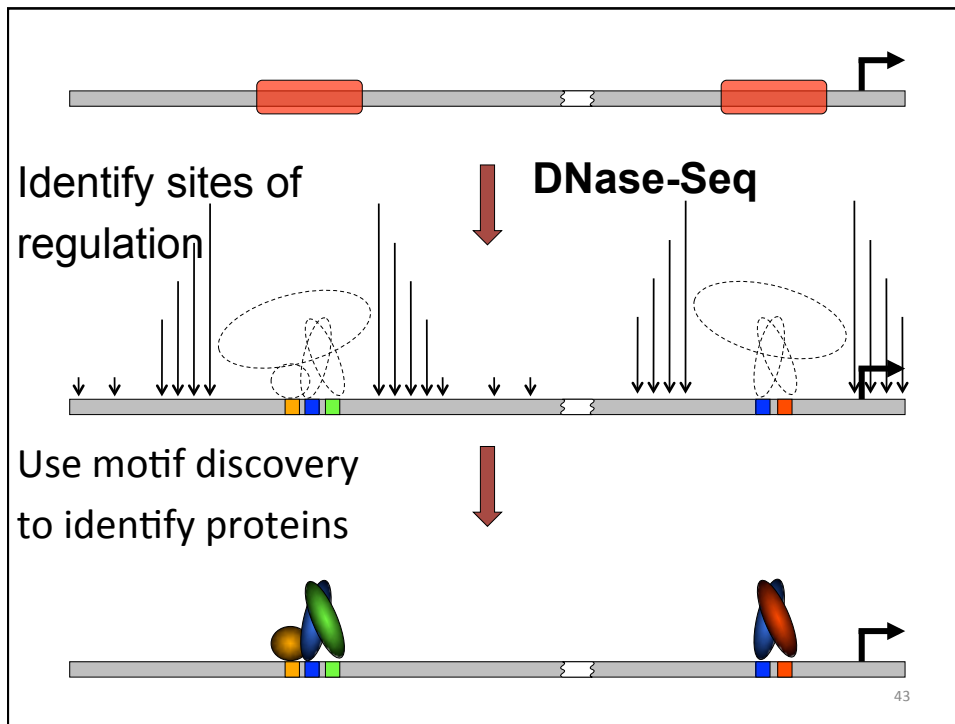
**How does this compare to mass-spec based approaches**

38



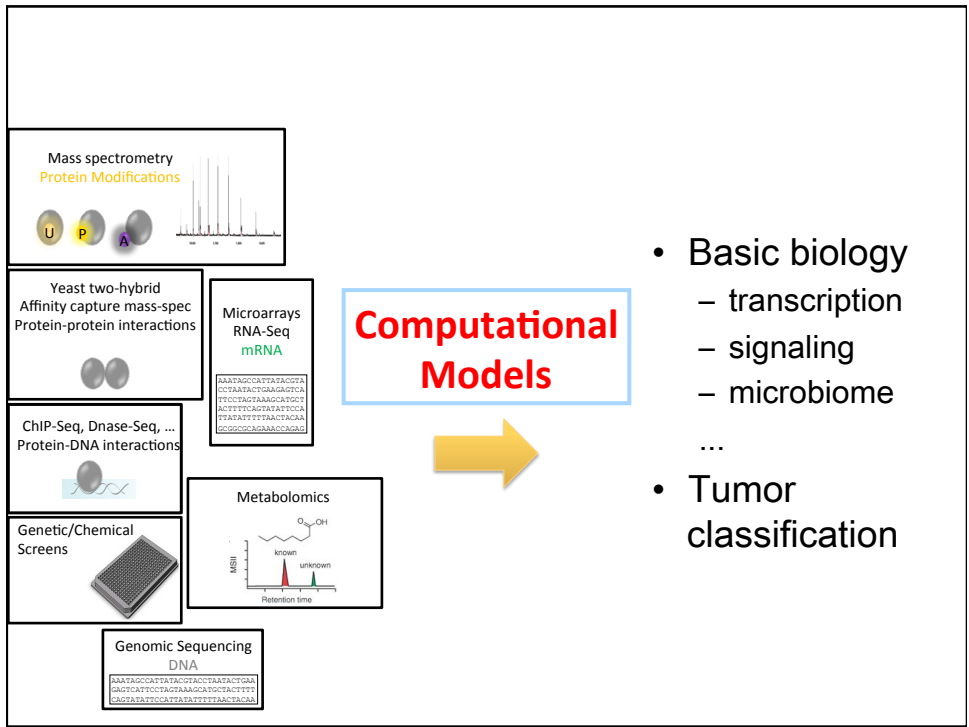






## Outline

- What molecules we can measure?
- How do we know which interact?
- How do we learn anything from these data?
  - Standard Approaches
  - Challenges
  - Network Methods
  - Toward Dynamic Models



The most common pediatric malignant brain tumor



- 70% survive but only 10% live independently as adults due to neurologic disability from the tumor and treatment

Polkinghorn and Tarbell. Medulloblastoma: tumorigenesis, current clinical paradigm, and efforts to improve risk stratification. Nat Clin Pract Oncol (2007) vol. 4 (5) pp. 295-304

## Mutations identified in genomics era of Medulloblastoma

**LETTER**

doi:10.1038/nature11329

## Medulloblastoma exome sequencing uncovers subtype-specific somatic mutations

Trevor J. Pugh<sup>1,2,3</sup>, Shyamal Dilhan Weeraratne<sup>3,4</sup>, Tenley C. Archer<sup>3,4</sup>, Daniel A. Pomeranz Krummel<sup>5</sup>, Daniel Auclair<sup>1</sup>, James Bochicchio<sup>1</sup>, Mauricio O. Carneiro<sup>1</sup>, Scott L. Carter<sup>1</sup>, Kristian Cibulskis<sup>1</sup>, Rachel L. Erlich<sup>1</sup>, Heidi Greulich<sup>1,2,3</sup>, Michael S. Lawrence<sup>1</sup>, Niall J. Lennon<sup>1</sup>, Aaron McKenna<sup>1</sup>, James Mehdritm<sup>1</sup>, Alex H. Ramos<sup>1,2,3</sup>, Michael G. Ross<sup>1</sup>, Carsten Russ<sup>1</sup>, Erica Shefter<sup>1</sup>, Andrey Sivachenko<sup>1</sup>, Brian Sogoleff<sup>1</sup>, Petar Stojanov<sup>1</sup>, Pablo Tamayo<sup>1</sup>, Jill P. Mesirov<sup>1</sup>, Vladimir Amanat<sup>1,2</sup>, Natalia Telder<sup>1,4</sup>, Soma Sengupta<sup>3,5</sup>, Jessica Pierre Francois<sup>3,4</sup>, Paul A. Northcott<sup>6</sup>, Michael D. Taylor<sup>6</sup>, Furong Yu<sup>1</sup>, Gerald R. Crabtree<sup>7,8</sup>, Amanda G. Kautzman<sup>1</sup>, Stacey B. Gabriel<sup>1</sup>, Gad Getz<sup>1</sup>, Natalie Hager<sup>1</sup>, David T. W. Jones<sup>9</sup>, Peter Lichter<sup>9</sup>, Stefan M. Pfister<sup>9</sup>, Thomas M. Roberts<sup>1,3</sup>, Matthew Meyerson<sup>1,3,10</sup>, Scott L. Pomeroy<sup>1,3,7</sup> & Yoon-Jae Cho<sup>1,3,7</sup>

doi:10.1038/nature11284

---

## Dissecting the genomic complexity underlying medulloblastoma

A list of authors and their affiliations appears at the end of the paper

**ARTICLE**

doi:10.1038/nature11213

## Novel mutations target distinct subgroups of medulloblastoma

Giles Robinson<sup>1,2,3,4</sup>, Matthew Parker<sup>1,4\*</sup>, Tanya A. Kranenburg<sup>1,2\*</sup>, Charles Lu<sup>1,5</sup>, Xiang Chen<sup>1,4</sup>, Li Ding<sup>1,5,6</sup>, Timothy N. Phoenix<sup>1,2</sup>, Erin Hedlund<sup>1,4</sup>, Lel Wei<sup>1,4,7</sup>, Xiaoyan Zhu<sup>1,2</sup>, Nader Chalhou<sup>1,2</sup>, Suzanne J. Baker<sup>1,2</sup>, Robert Huether<sup>1,4,8</sup>, tadhika Thiruvankatam<sup>1,2</sup>, Jianmin Wang<sup>1,9</sup>, Gang Wu<sup>1,4</sup>, Michael Rusch<sup>1,4</sup>, Xin Hong<sup>1,3</sup>, Ma<sup>1,4</sup>, John Easton<sup>1,4</sup>, Bhavin Vadodaria<sup>1,4</sup>, Arzu Omar-Thomas<sup>1,4</sup>, Tong Lin<sup>1,10</sup>, Paugh<sup>1,11</sup>, David Zhao<sup>1,4</sup>, Daisuke Kawachi<sup>1,12</sup>, Marlene F. Rousselet<sup>1,4</sup>, Ching C. Lau<sup>1,13</sup>, Eric Bouffet<sup>1,14</sup>, Tim Hassall<sup>1,15</sup>, Sridharan Gururangan<sup>1,16</sup>, ucinda L. Fulton<sup>1,3,6</sup>, David J. Dooling<sup>1,3,6</sup>, Kerri Ochoa<sup>1,3,6</sup>, Amar Gajjar<sup>1,3</sup>, n<sup>1,3,6,17</sup>, James R. Downing<sup>1,7</sup>, Jinghui Zhang<sup>1,4</sup> & Richard J. Gilbertson<sup>1,3,3</sup>

doi:10.1038/nature11327

---

## Subgroup-specific structural variation across 1,000 medulloblastoma genomes

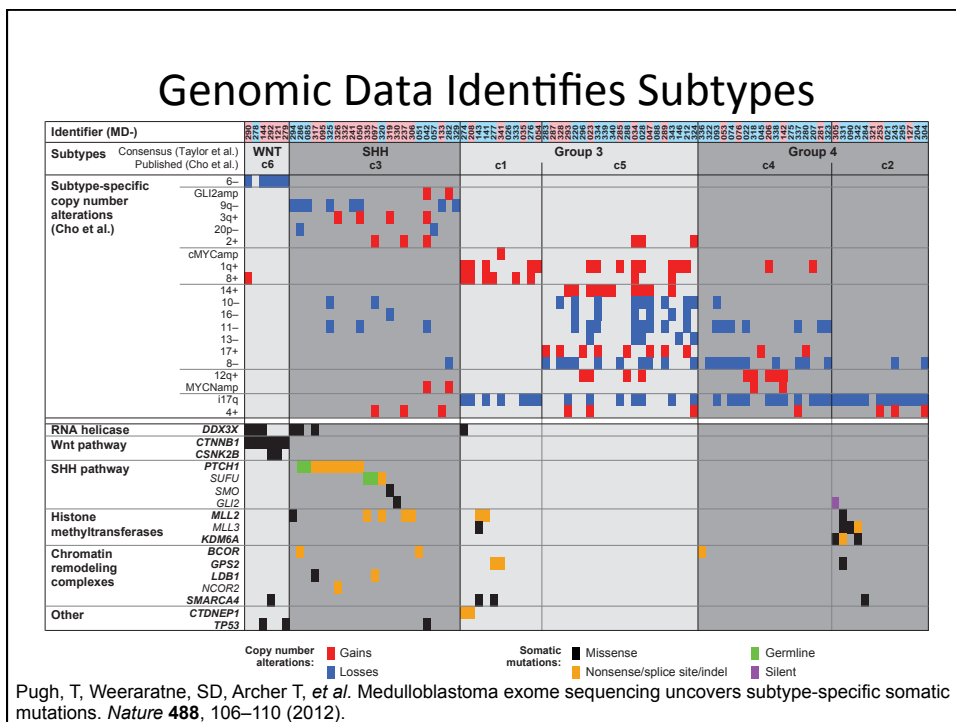
A list of authors and their affiliations appears at the end of the paper

**ARTICLE**

## Novel mutations target distinct subgroups of medulloblastoma

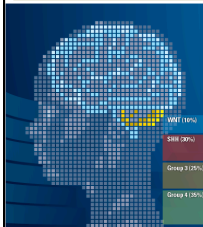
Giles Robinson<sup>1,2,3,4</sup>, Matthew Parker<sup>1,4\*</sup>, Tanya A. Kranenburg<sup>1,2\*</sup>, Charles Lu<sup>1,5</sup>, Xiang Chen<sup>1,4</sup>, Li Ding<sup>1,5,6</sup>, Timothy N. Phoenix<sup>1,2</sup>, Erin Hedlund<sup>1,4</sup>, Lel Wei<sup>1,4,7</sup>, Xiaoyan Zhu<sup>1,2</sup>, Nader Chalhou<sup>1,2</sup>, Suzanne J. Baker<sup>1,2</sup>, Robert Huether<sup>1,4,8</sup>, tadhika Thiruvankatam<sup>1,2</sup>, Jianmin Wang<sup>1,9</sup>, Gang Wu<sup>1,4</sup>, Michael Rusch<sup>1,4</sup>, Xin Hong<sup>1,3</sup>, Ma<sup>1,4</sup>, John Easton<sup>1,4</sup>, Bhavin Vadodaria<sup>1,4</sup>, Arzu Omar-Thomas<sup>1,4</sup>, Tong Lin<sup>1,10</sup>, Paugh<sup>1,11</sup>, David Zhao<sup>1,4</sup>, Daisuke Kawachi<sup>1,12</sup>, Marlene F. Rousselet<sup>1,4</sup>, Ching C. Lau<sup>1,13</sup>, Eric Bouffet<sup>1,14</sup>, Tim Hassall<sup>1,15</sup>, Sridharan Gururangan<sup>1,16</sup>, ucinda L. Fulton<sup>1,3,6</sup>, David J. Dooling<sup>1,3,6</sup>, Kerri Ochoa<sup>1,3,6</sup>, Amar Gajjar<sup>1,3</sup>, n<sup>1,3,6,17</sup>, James R. Downing<sup>1,7</sup>, Jinghui Zhang<sup>1,4</sup> & Richard J. Gilbertson<sup>1,3,3</sup>

doi:10.1038/nature11327

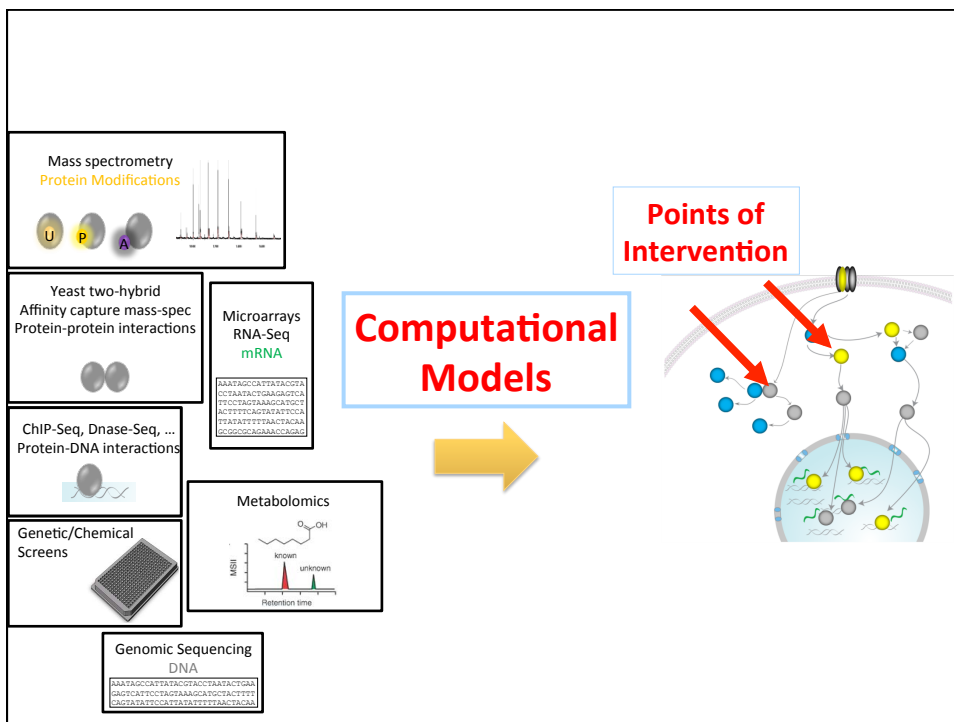


# Four Types of Medulloblastoma

	WNT (~10%)	SHH (~30%)	Group 3 (~25%)	Group 4 (~35%)
<b>Clinical features</b>				
Gender ratio (M/F)	~1/1	~1.5/1	~2/1	~3/1
Age distribution				
Histology	Classic; very rare LCA	Classic > desmoplastic/nodular > LCA > MBEN	Classic > LCA	Classic; rarely LCA
Metastasis at diagnosis	~5-10%	~15-20%	~40-45%	~35-40%
Overall survival (5 years)	~95%	~75%	~50%	~75%
Proposed cell of origin	Lower rhombic lip progenitor cells	CGNPs of the EGL and cochlear nucleus; neural stem cells of the SVZ	Prominin 1 <sup>+</sup> , lineage <sup>-</sup> neural stem cells; CGNPs of the EGL	Unknown



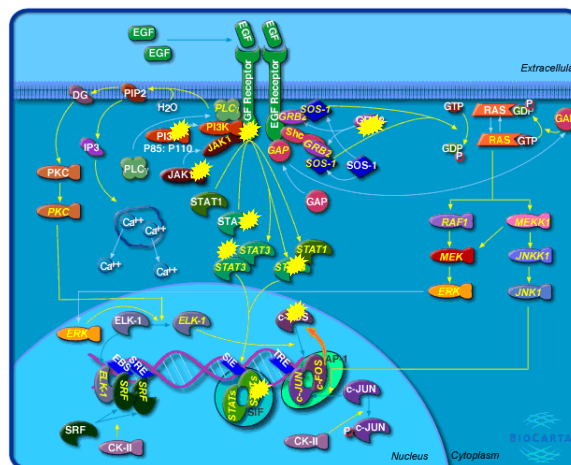
Northcott et al. Nat Rev Cancer. 2012 Dec;12(12):818-34. doi: 10.1038/nrc3410.



## Outline

- What molecules we can measure?
- How do we know which interact?
- How do we learn anything from these data?
  - Standard Approaches
  - Challenges
  - Network Methods
  - Toward Dynamic Models


## Map to Known Pathways



Name	Availability	Reference
<b>ORA tools</b>		
OrnV-Express	Web ( <a href="http://vortex.cs.wayne.edu">http://vortex.cs.wayne.edu</a> )	[45]
GeneMAPP	Standalone ( <a href="http://www.genemap.org">http://www.genemap.org</a> )	[11,71]
GoMiner	Standalone, Web ( <a href="http://discover.nci.nih.gov/gominer">http://discover.nci.nih.gov/gominer</a> )	[72,73]
FitGO	Web ( <a href="http://babelomics.bioinformatics">http://babelomics.bioinformatics</a> )	[74]
GOstat	Web ( <a href="http://gostat.wehi.edu.au">http://gostat.wehi.edu.au</a> )	[7]
FuncAssociate	Web ( <a href="http://fama.mshri.on.ca/funcassociate/">http://fama.mshri.on.ca/funcassociate/</a> )	[6]
GOToolBox	Web ( <a href="http://genome.crg.es/GOToolBox/">http://genome.crg.es/GOToolBox/</a> )	[10]
GeneMerge	Standalone, Web ( <a href="http://genemerge.cbc.umt.edu/">http://genemerge.cbc.umt.edu/</a> )	[9]
GOEAST	Web ( <a href="http://omicslab.genetics.ac.cn/GOEAST/">http://omicslab.genetics.ac.cn/GOEAST/</a> )	[75]
ClueGO	Standalone ( <a href="http://www.kclupmc.fr/cluego/">http://www.kclupmc.fr/cluego/</a> )	[76]
FunSpec	Web ( <a href="http://funspec.med.utoronto.ca/">http://funspec.med.utoronto.ca/</a> )	[77]
GABIAN	Web	[78]
GO-TermFinder	Standalone ( <a href="http://search.cpan.org/dist/GO-TermFinder/">http://search.cpan.org/dist/GO-TermFinder/</a> )	[8]
WebGestalt	Web ( <a href="http://bioinfo.vanderbilt.edu/webgestalt/">http://bioinfo.vanderbilt.edu/webgestalt/</a> )	[79]
agriGO	Web ( <a href="http://bioinfo.cau.edu.cn/agriGO/">http://bioinfo.cau.edu.cn/agriGO/</a> )	[80]
GOFA	Standalone, Web ( <a href="http://edkb.fda.gov/webstart/arraytrack/">http://edkb.fda.gov/webstart/arraytrack/</a> )	[81]
WEGO	Web ( <a href="http://wego.genomics.org.cn/cgi-bin/wego/index.pl">http://wego.genomics.org.cn/cgi-bin/wego/index.pl</a> )	[82]
<b>FCS tools</b>		
GSEA	Standalone ( <a href="http://www.broadinstitute.org/gsea/">http://www.broadinstitute.org/gsea/</a> )	[21,29]
igPathway	Standalone (BioConductor)	[22]
Category	Standalone (BioConductor)	[24]
SAFE	Standalone (BioConductor)	[30]
GlobalTest	Standalone (BioConductor)	[15]
PC072	Standalone (BioConductor)	[77]
SAM-GS	Standalone ( <a href="http://www.ualberta.ca/~yaru/software.html">http://www.ualberta.ca/~yaru/software.html</a> )	[83]
Catmap	Standalone ( <a href="http://bioinfo.thep.lu.se/catmap.html">http://bioinfo.thep.lu.se/catmap.html</a> )	[84]
T-profiler	Web ( <a href="http://www.t-profiler.org/">http://www.t-profiler.org/</a> )	[85]
FunCluster	Standalone ( <a href="http://cornellihenegar.info/FunCluster.htm">http://cornellihenegar.info/FunCluster.htm</a> )	[86]
GeneTrail	Web ( <a href="http://genetrail.bioinf.uni-st.de/">http://genetrail.bioinf.uni-st.de/</a> )	[87]
Gazer	Web	[88]
<b>PT-based tools</b>		
ScorePAGE	No implementation available	[37]
Pathway-Express	Web ( <a href="http://vortex.cs.wayne.edu">http://vortex.cs.wayne.edu</a> )	[38,39]
SPIA	Standalone (BioConductor)	[40]
NetGSA	No implementation available	[43]

doi:10.1371/journal.pcbi.1002375.t001

Khatri P, Sirota M, Butte AJ (2012) Ten Years of Pathway Analysis: Current Approaches and Outstanding Challenges. *PLoS Comput Biol* 8(2): e1002375. doi:10.1371/journal.pcbi.1002375  
<http://127.0.0.1:8081/ploscompbiol/article?id=info:doi/10.1371/journal.pcbi.1002375>



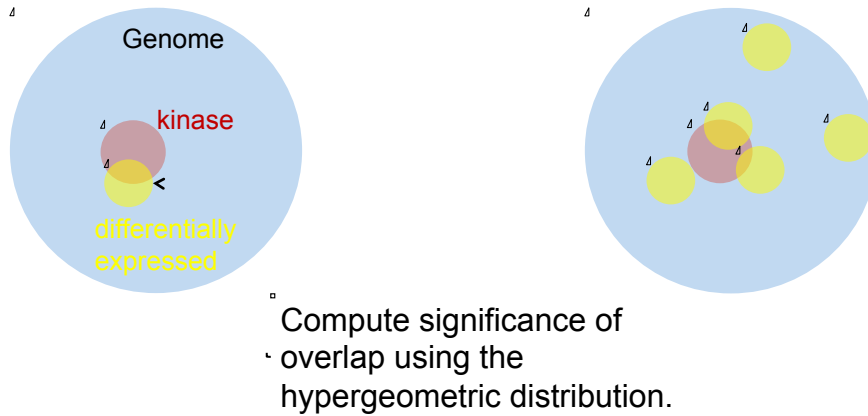
## Known Pathways



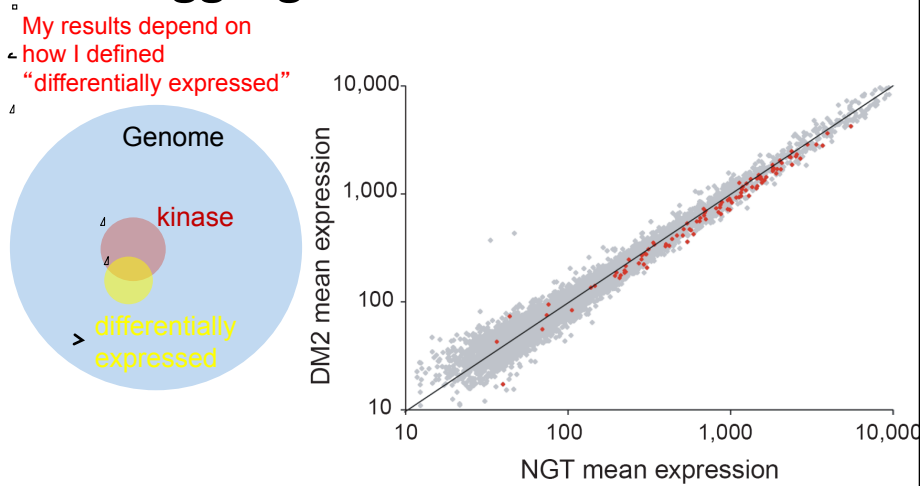
Controlled vocabulary to describe genes:

- Biological process
- Cellular component
- Molecular function

## Statistical Significance



## Aggregate score statistics



Mootha et al. (2003). *Nature Genetics* **34**, 267 – 273. doi:10.1038/ng1180



# Aggregate score statistics

<http://www.broadinstitute.org/gsea/>

**Overview**

**Gene Set Enrichment Analysis (GSEA)** is a computational method that determines whether a priori defined set of genes shows statistically significant, concordant differences between two biological states (e.g. phenotypes).

**What's New**

02/19/10: We have a new release of GSEA 2.0.5 that fixes the FTP problems that have been experienced recently. Please discontinue use of older versions and use the new version instead.

12/10/09: Leading Edge Analysis now works correctly in Release GSEA 2.0.5. There are no changes to the algorithm or functionality.

12/07/2009: Release GSEA 2.0.5 of the GSEA Java application is now available. The new release has been updated to work on 64-bit platforms. There are no changes to the algorithm or functionality. This update requires Java 6 (or all platforms).

**Getting Started**

A quick tutorial to get you up and running.

**Tools and Information**

**Downloads:** Implementations of GSEA plus additional resources to analyze, annotate and interpret enrichment results.

**Molecular Signatures Database:** A collection of gene sets for use with GSEA software and tools for exploring them.

**Documentation:** Information on the GSEA software, the GSEA algorithm.

**Registration**

Please register to download the GSEA software and view the MSigDB gene sets. After registering, you can log in at any time using your email address. Registration is free. Its only purpose is to help us track usage for reports to our funding agencies.

**Contributors**

GSEA is maintained by the GSEA team. Our thanks to our many contributors. Funded by: National Cancer Institute, National Institutes of Health, National Institute of General Medical Sciences.

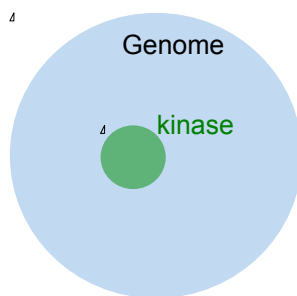
**Citing GSEA**

To cite your use of the GSEA software, please reference Subramanian, Tamayo, et al. (2005, PNAS 102, 15545-15550) and Hoadly, Lindgren, et al. (2010, Nat Genet 34, 207-215).

Broad Home | Cancer Genomics | Broad Site Map

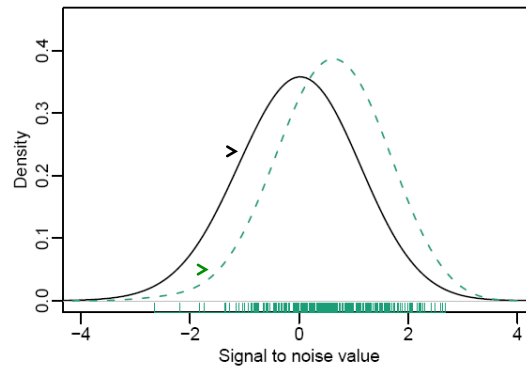
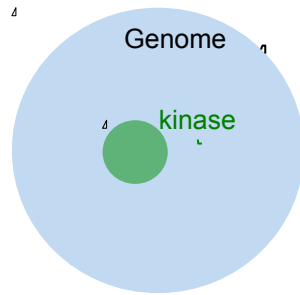
MSigDB Database v2.5 updated April 1, 2015  
©2015 Broad Institute. All rights reserved. December 12, 2015

# Aggregate score statistics



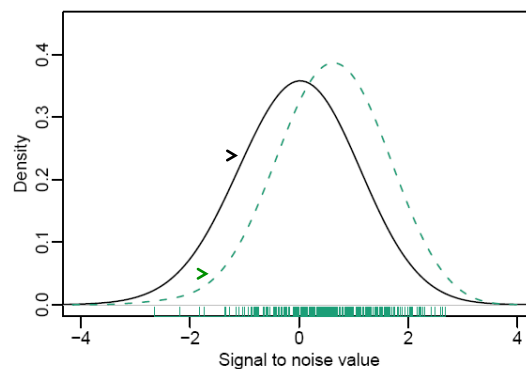
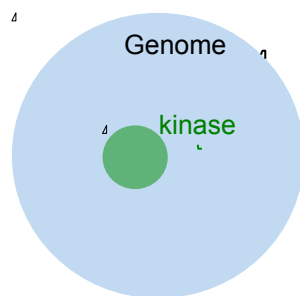
- For each list of genes in a database (such as GO)
  - compute distribution of t-statistics
  - compare this distribution to the overall distribution

## Aggregate score statistics



GSEA uses a Kolmogorov-Smirnov statistic to compare the distributions of t-statistics

## Aggregate score statistics



Irizarry, et al. argue for  $X^2$  and z-test

Gene set enrichment analysis made simple. (2009) Stat Methods Med Res

<http://www.bepress.com/jhubiostat/paper185/>

## Outline

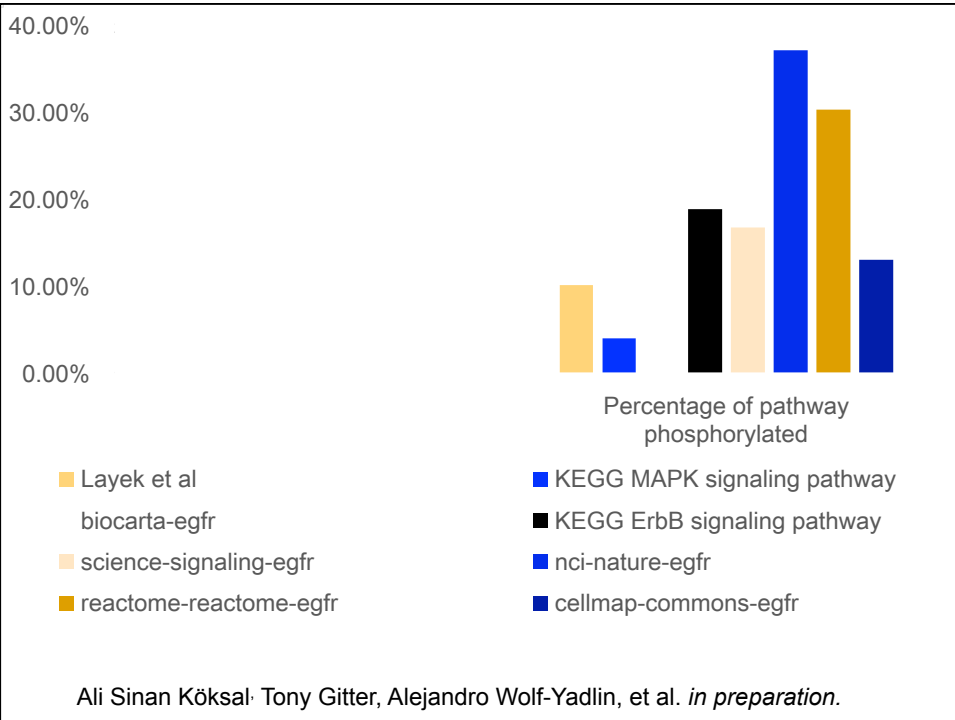
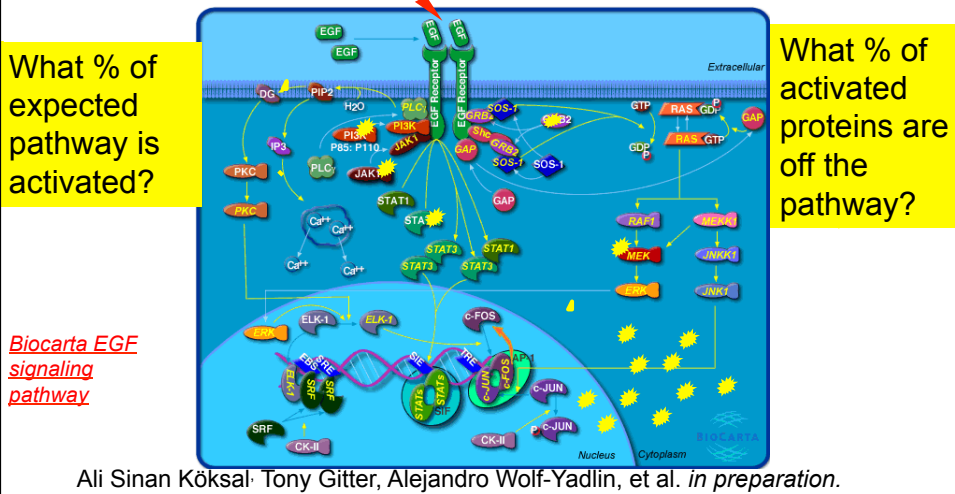
- What molecules we can measure?
- How do we know which interact?
- How do we learn anything from these data?
  - Standard Approaches
  - Challenges
  - Network Methods
  - Toward Dynamic Models

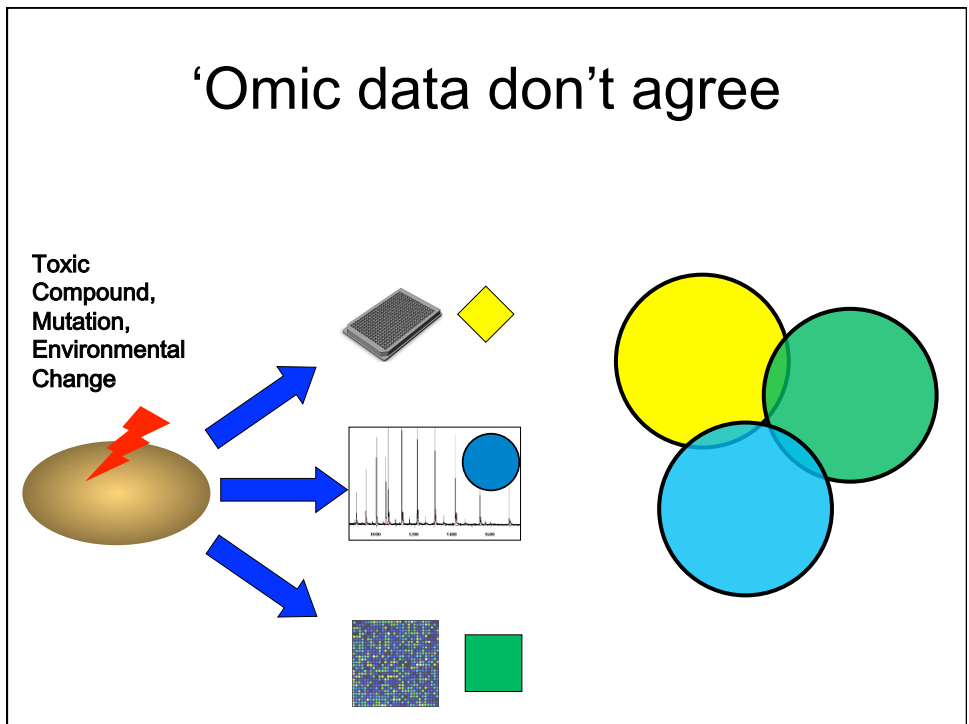
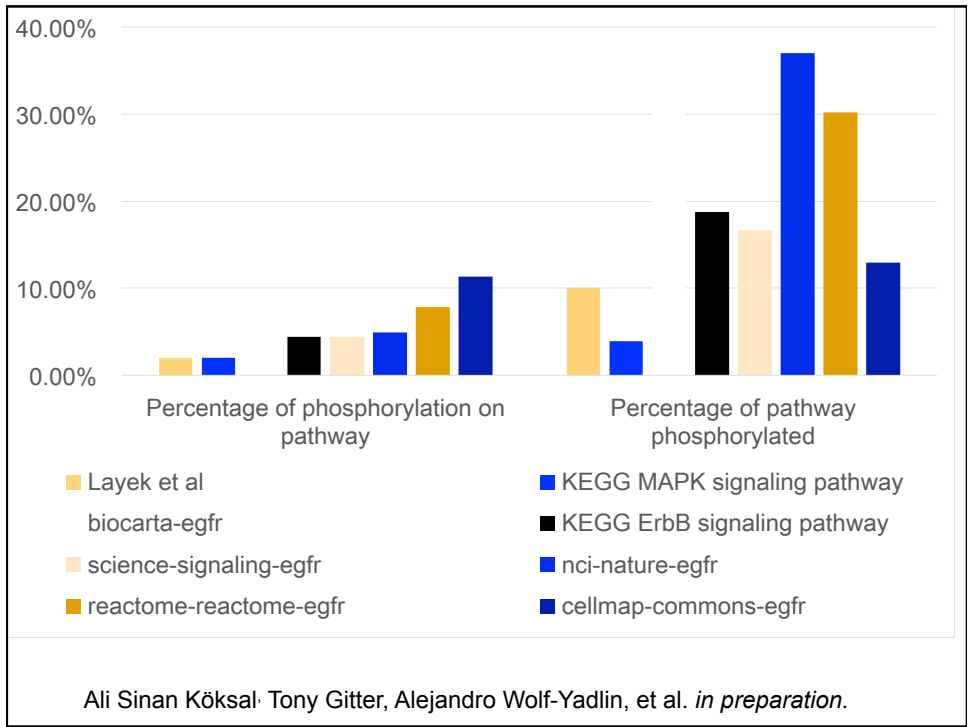
## End of Omics Review



“Getting an Education from MIT is like taking a drink from a Fire Hose.”  
Former MIT President Jerome Weisner

# Most 'Omic Hits Don't Lie in Known Pathways







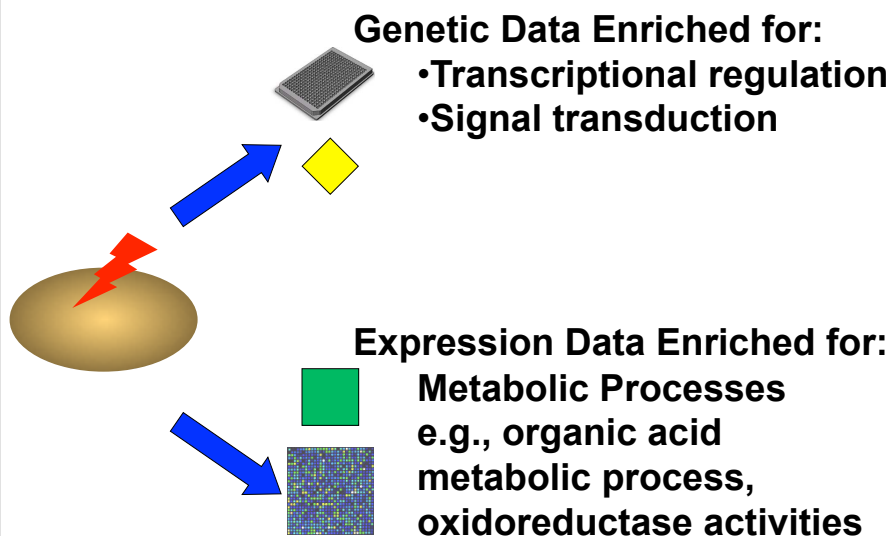
**Esti  
Yeger-Lotem**

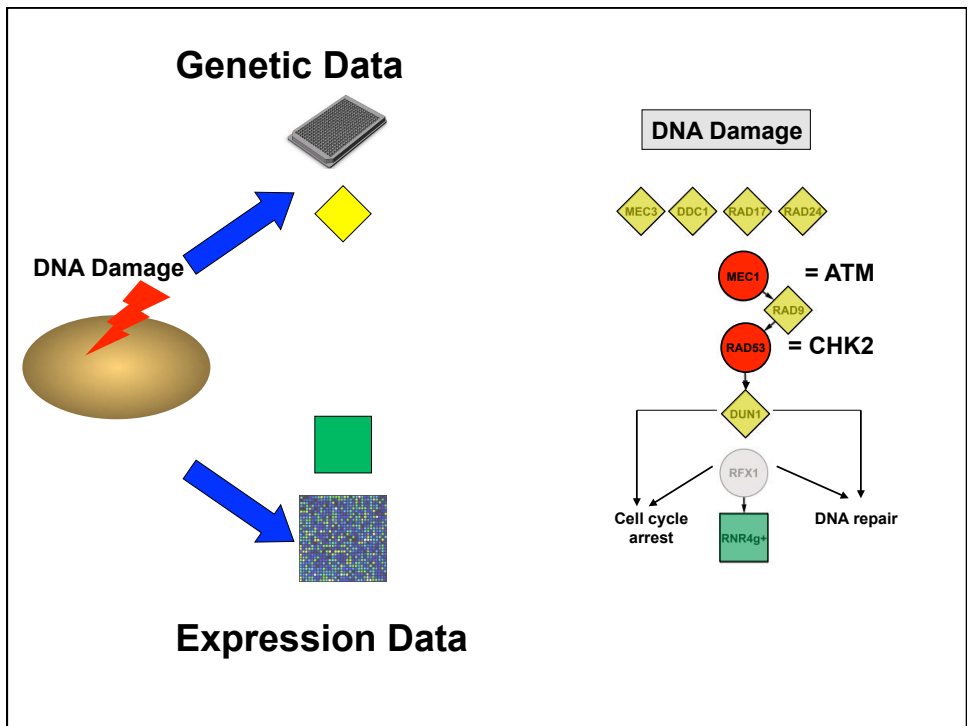
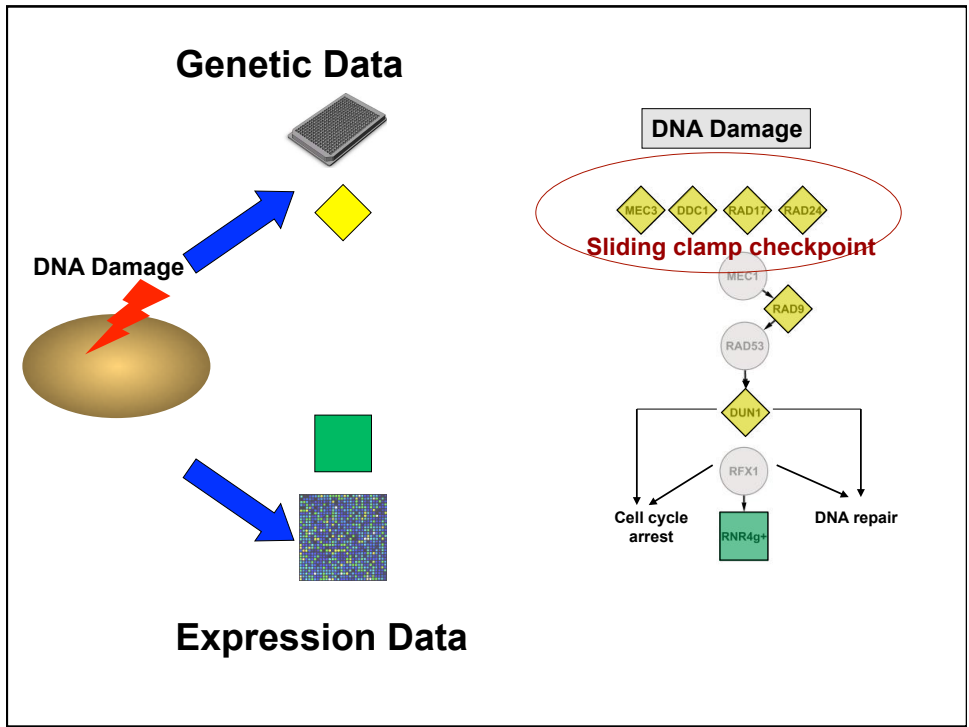
Senior Lecturer  
Ben-Gurion  
University  
National Institute for  
Biotechnology in the  
Negev  
Israel

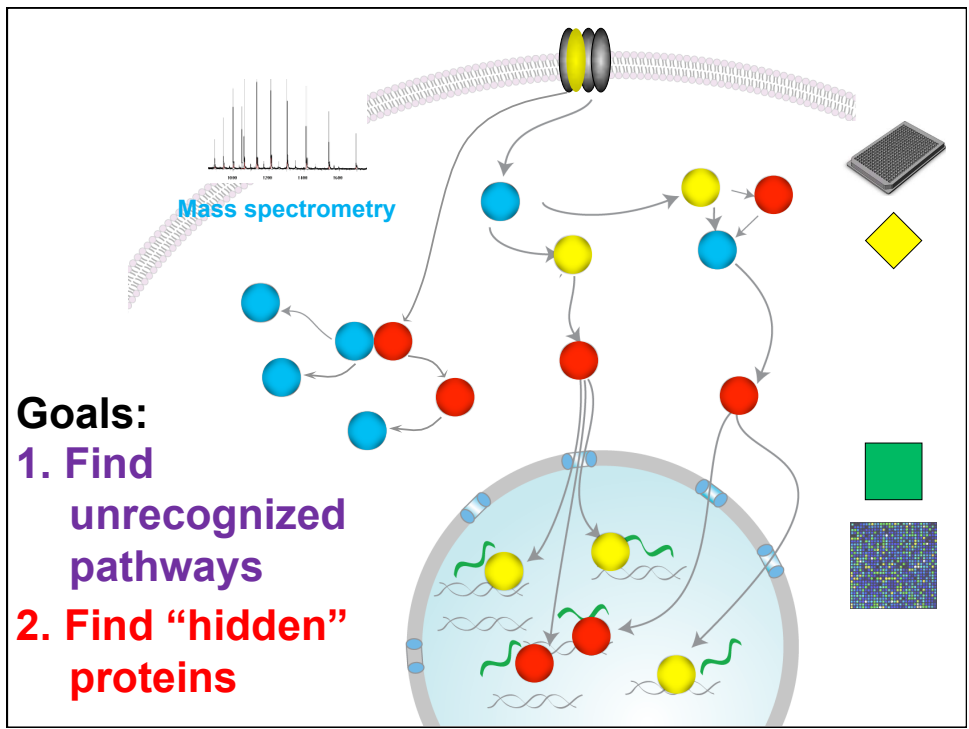
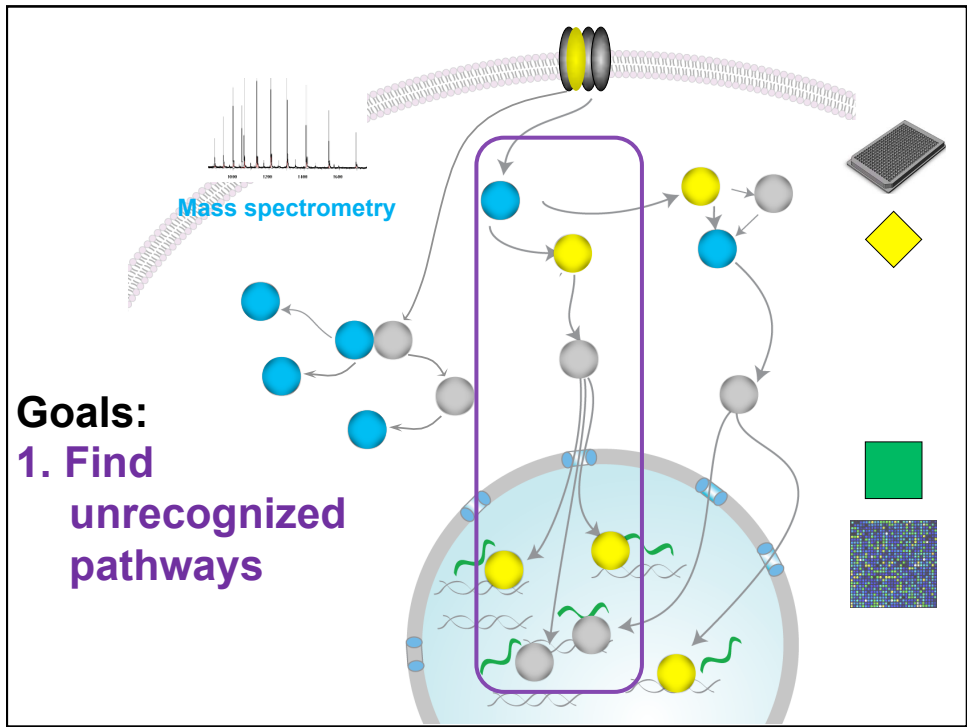
**Laura  
Riva**

Team Leader  
Center for Genomic  
Science  
Istituto Italiano di  
Tecnologia  
Italy

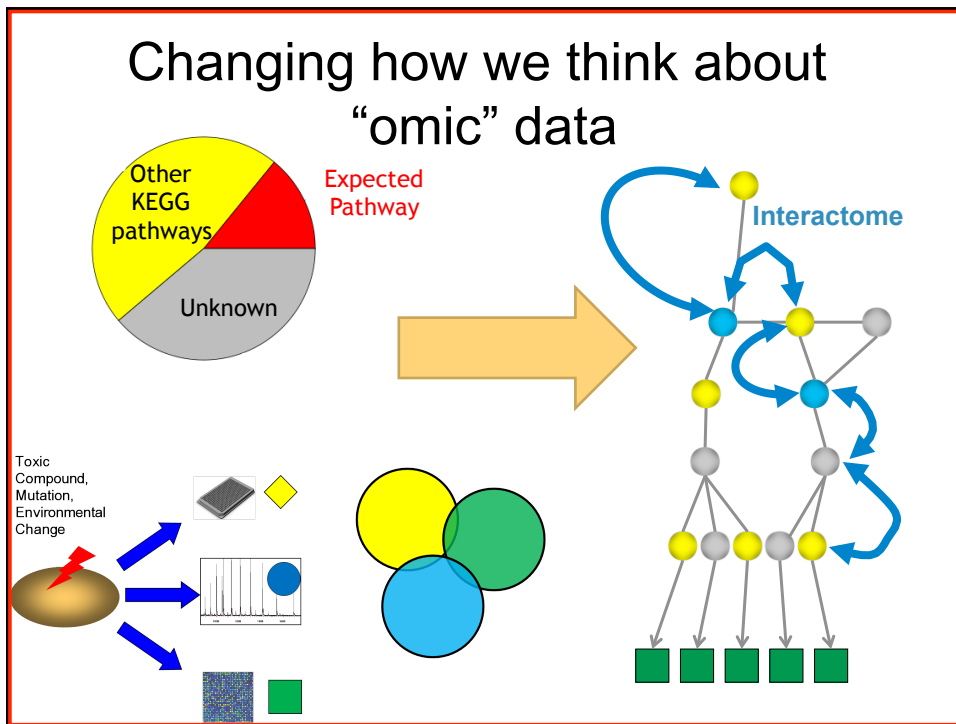
**For 156 perturbations:**









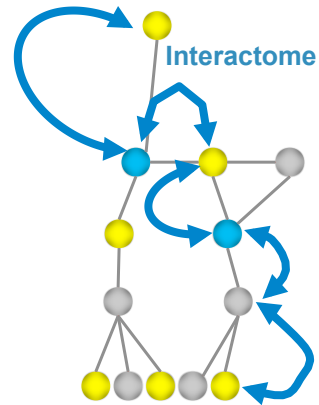


## Outline

- Network Modeling
  - Approach
  - Predicting in vitro targets using RNA-Seq, DNase-Seq, and phospho-proteomics

## Approach

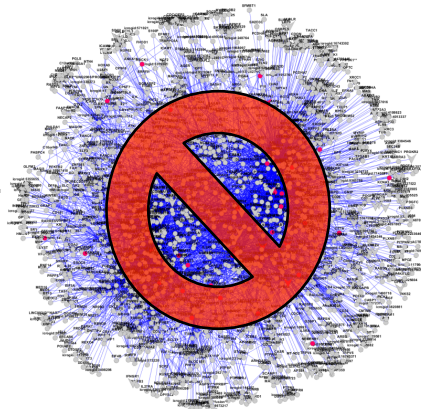
Map data onto a network of known interactions.

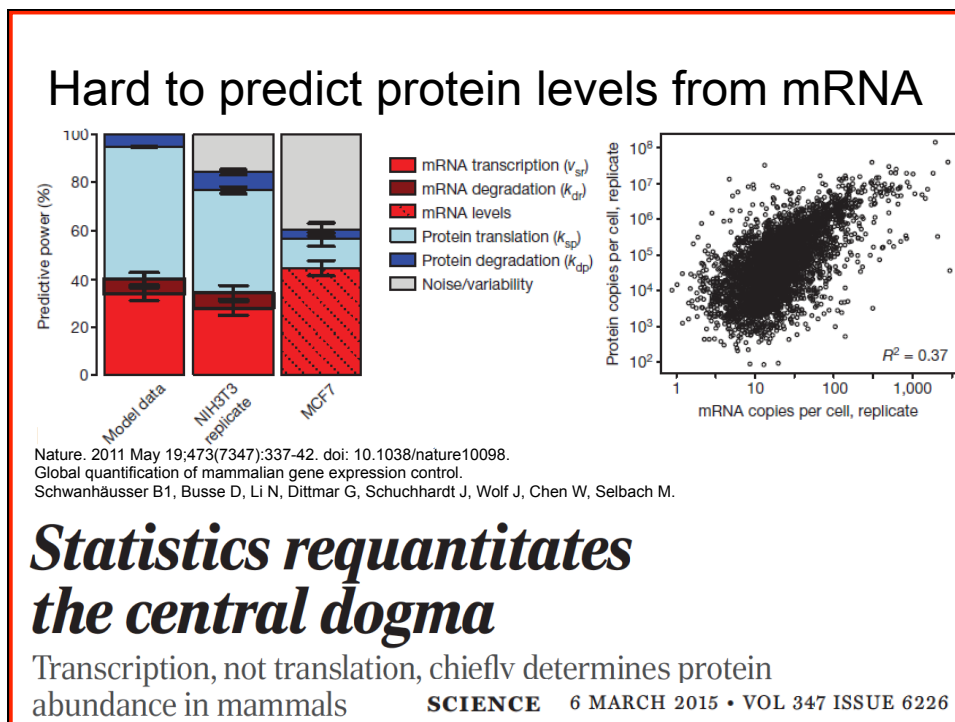
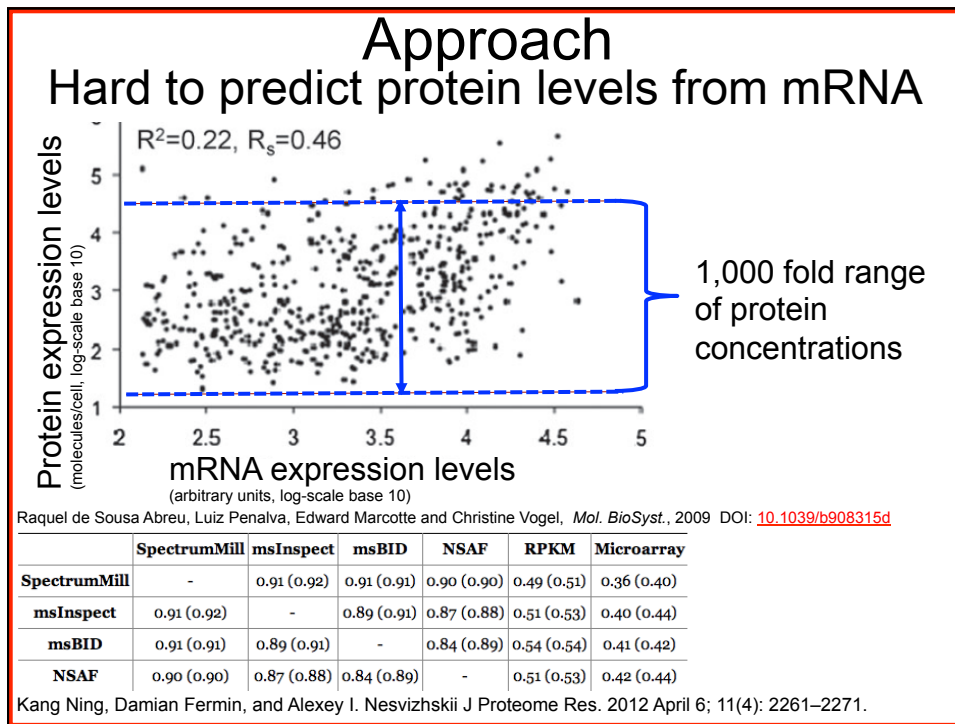


## Approach

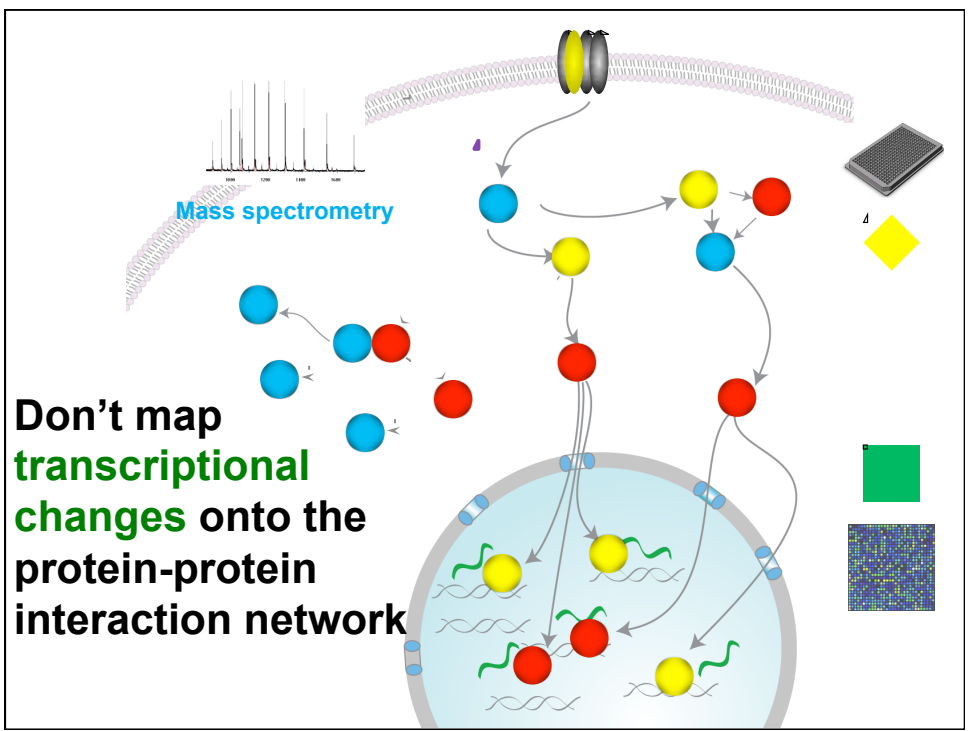
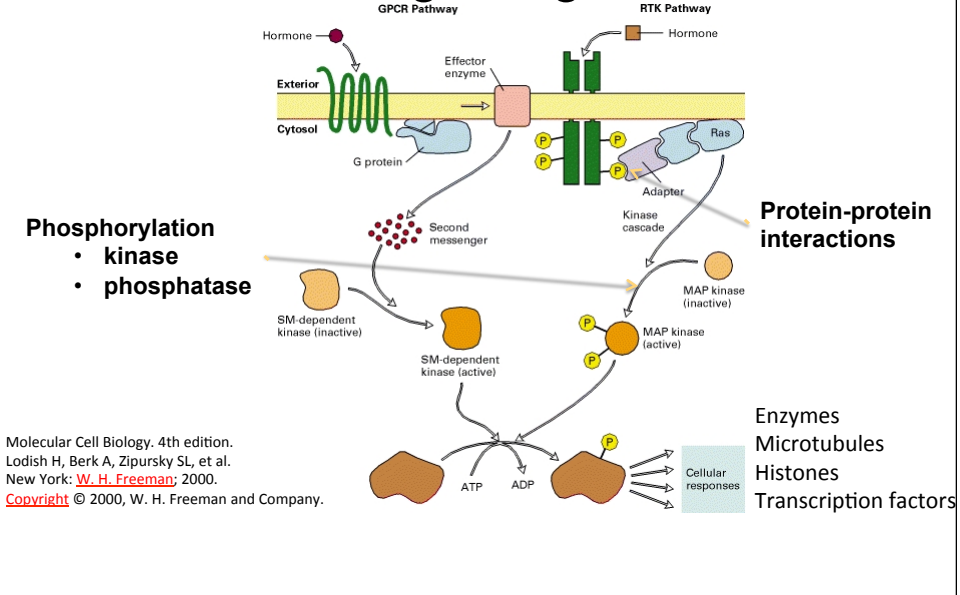
Lesson 1:  
Network models make sense of diverse data.

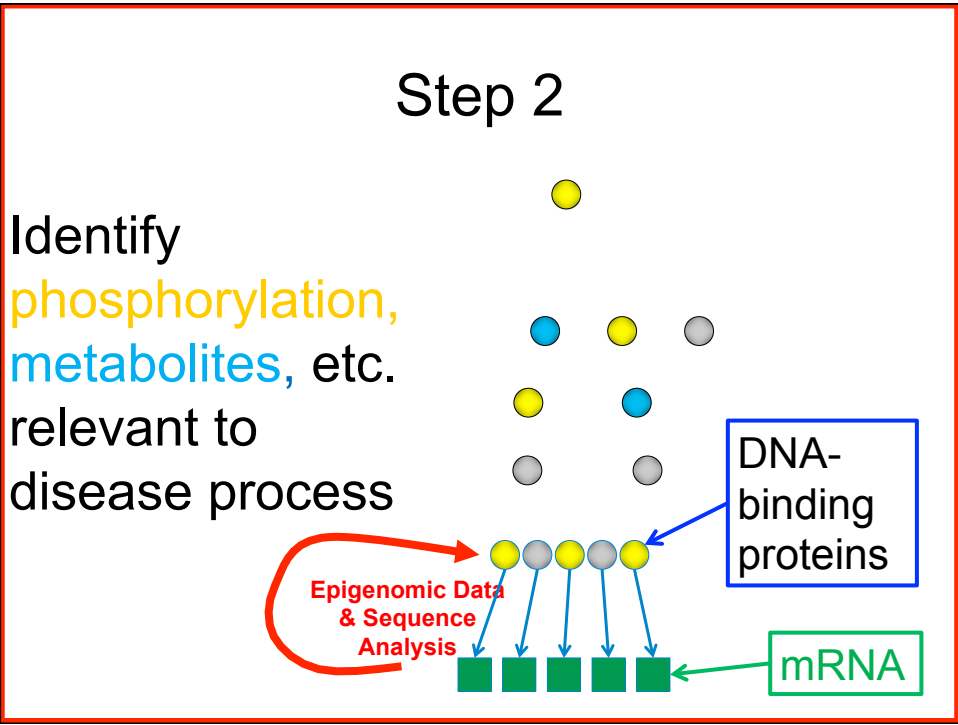
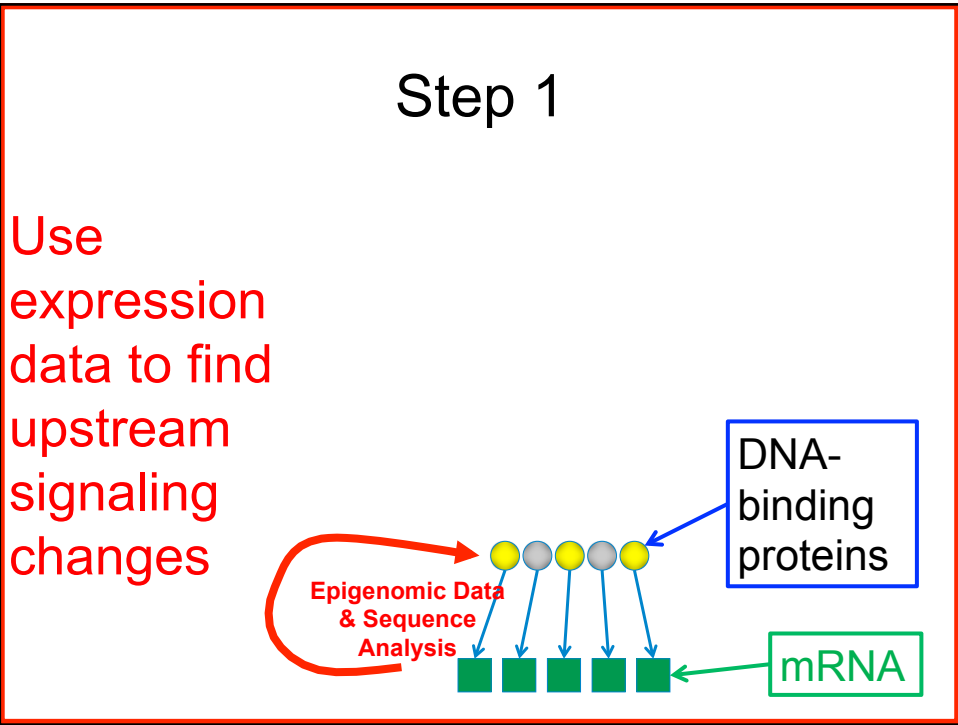
Lesson 2:  
Hairballs do not!  
Advanced network algorithms needed.

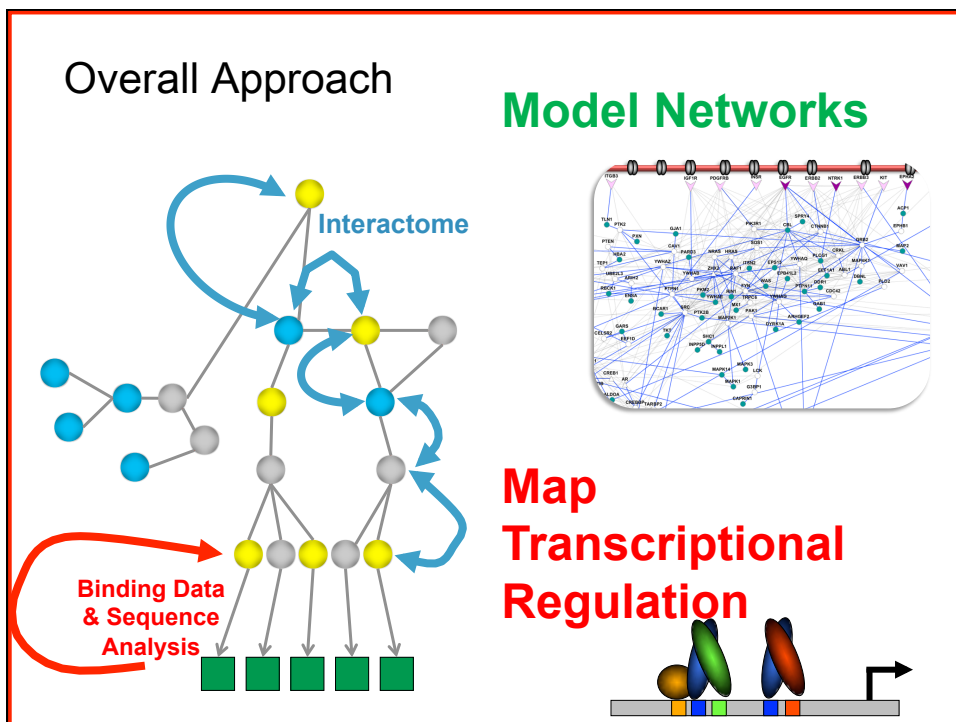
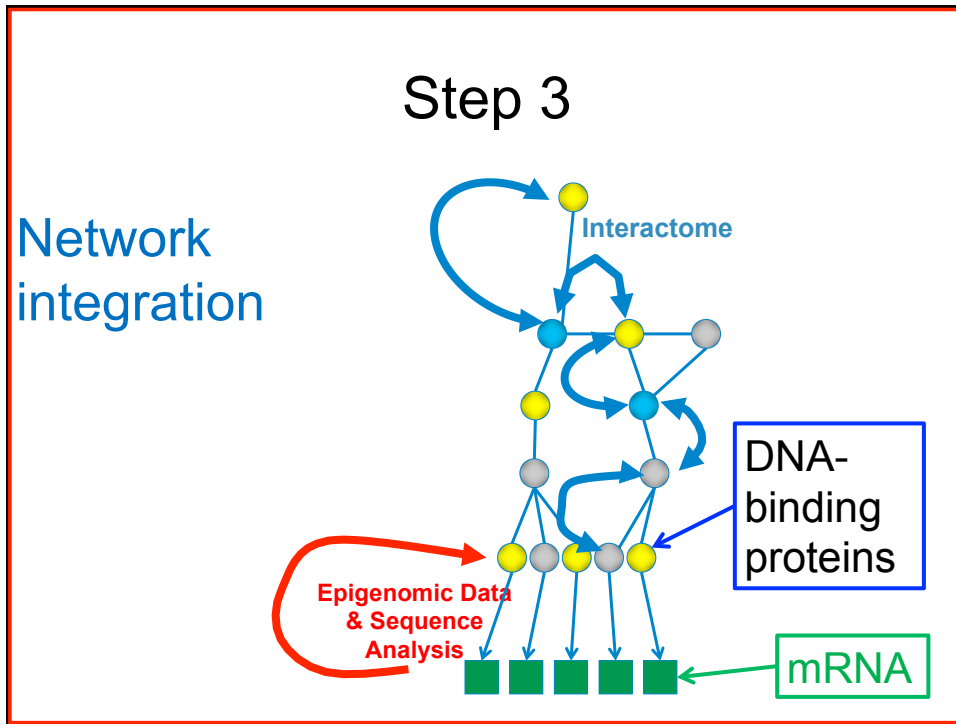




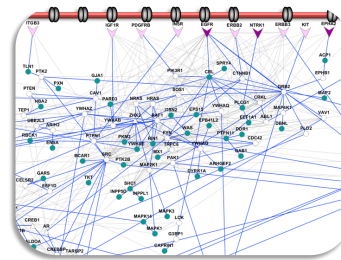
# .. and mRNA levels do not reveal signaling







# Outline



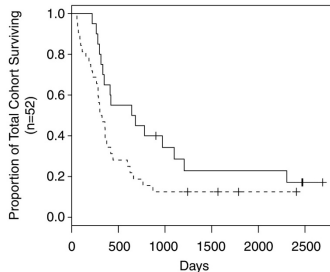
- Network Modeling
  - Approach
  - Predicting in vitro targets using RNA-Seq, DNase-Seq, and phospho-proteomics

Carol Huang

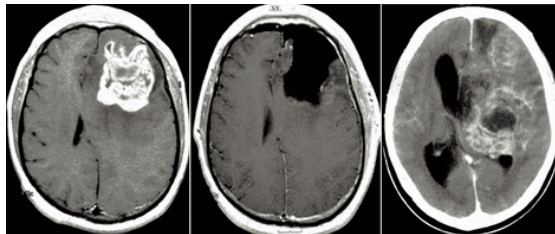


**Linking Proteomic and Transcriptional Data through the Interactome and Epigenome Reveals a Map of Oncogene-induced Signaling**  
 PLOS Computational Biology, 2013

Post-doctoral Associate  
 Salk Institute



Pope W B et al. Radiology 2008;249:268-277

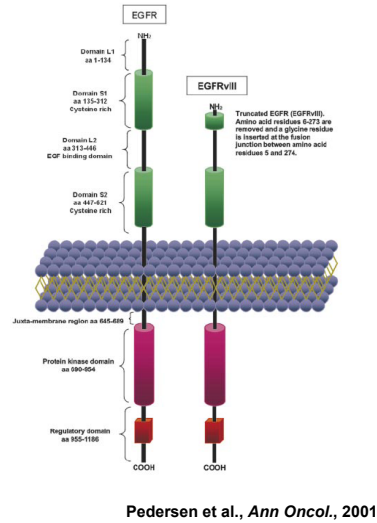


Presentation Post-op Recurrence

Weil RJ (2006) PLoS Med 3(1): e31.

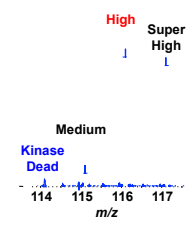
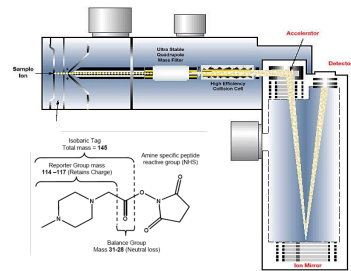
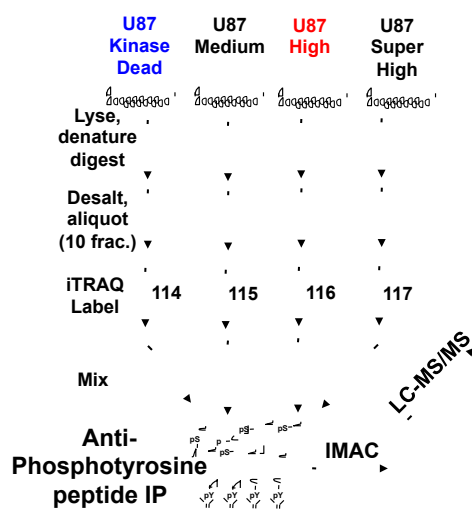
# EGFR variant III mutation

- Most common EGFR mutation in human cancer
- 60% of GBMs and 20% of anaplastic astrocytomas, also in lung, breast and prostate tumors
- EGFRvIII expression correlates with shorter life expectancies



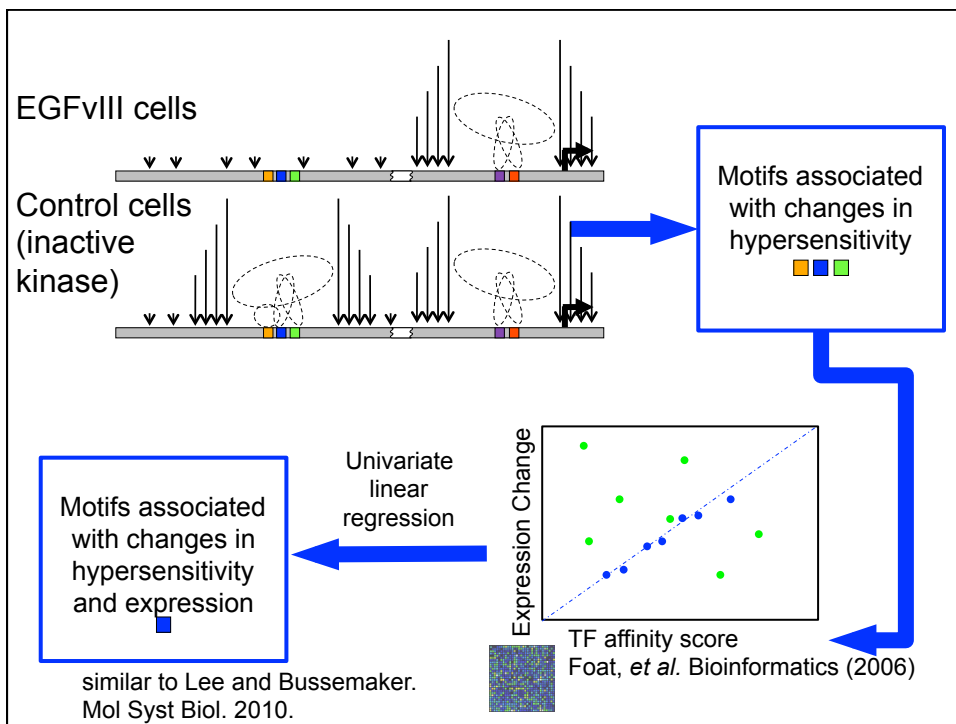
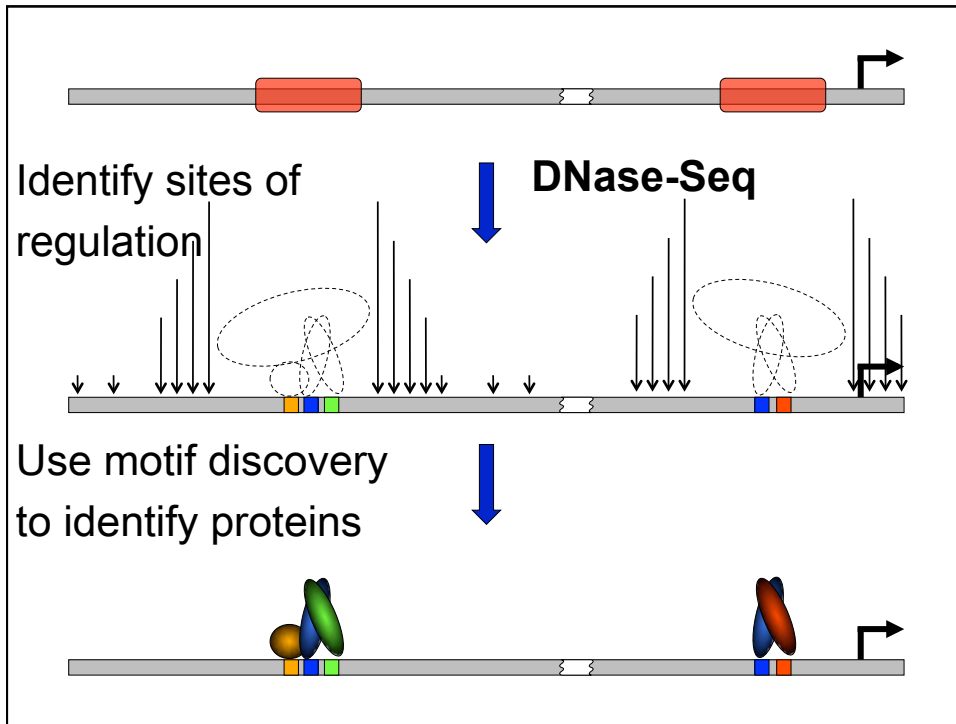
# Compare Tyrosine Phosphorylation by Mass Spec

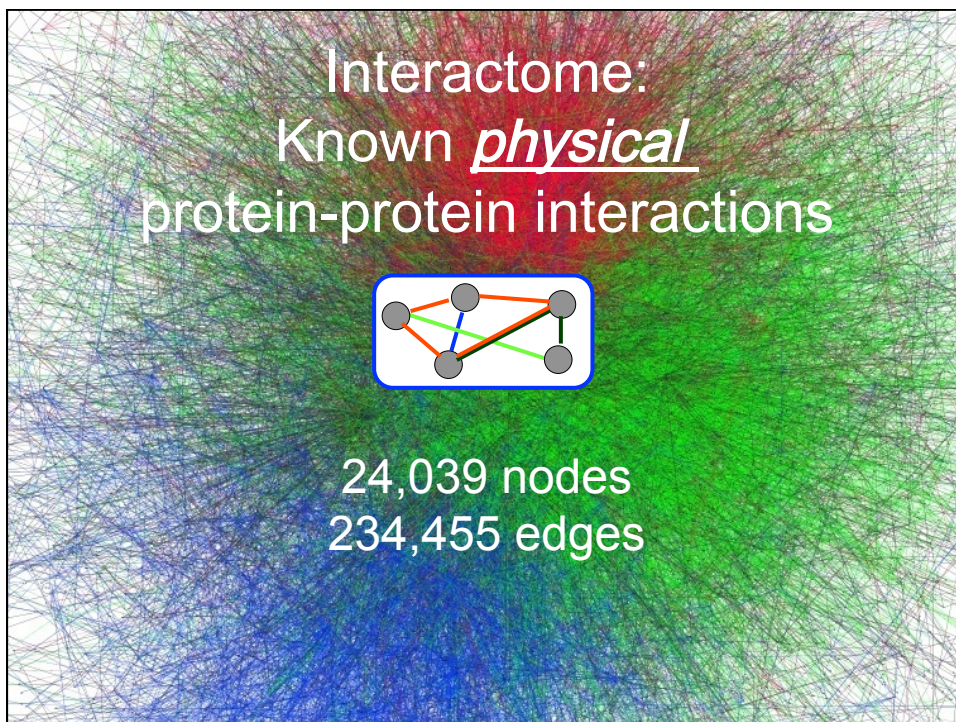
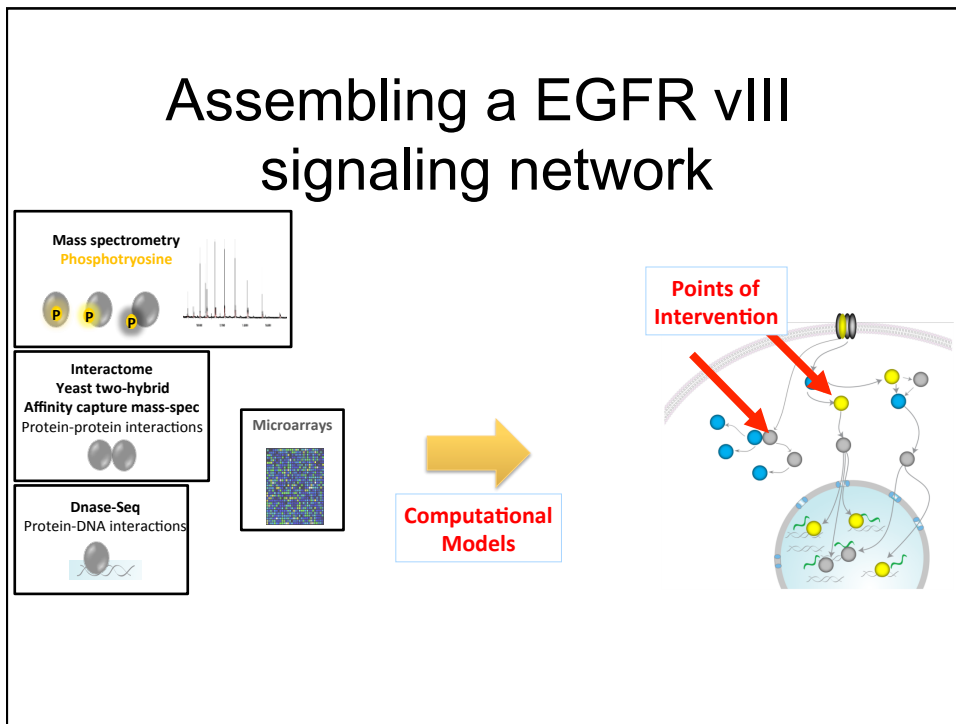
White Lab, MIT

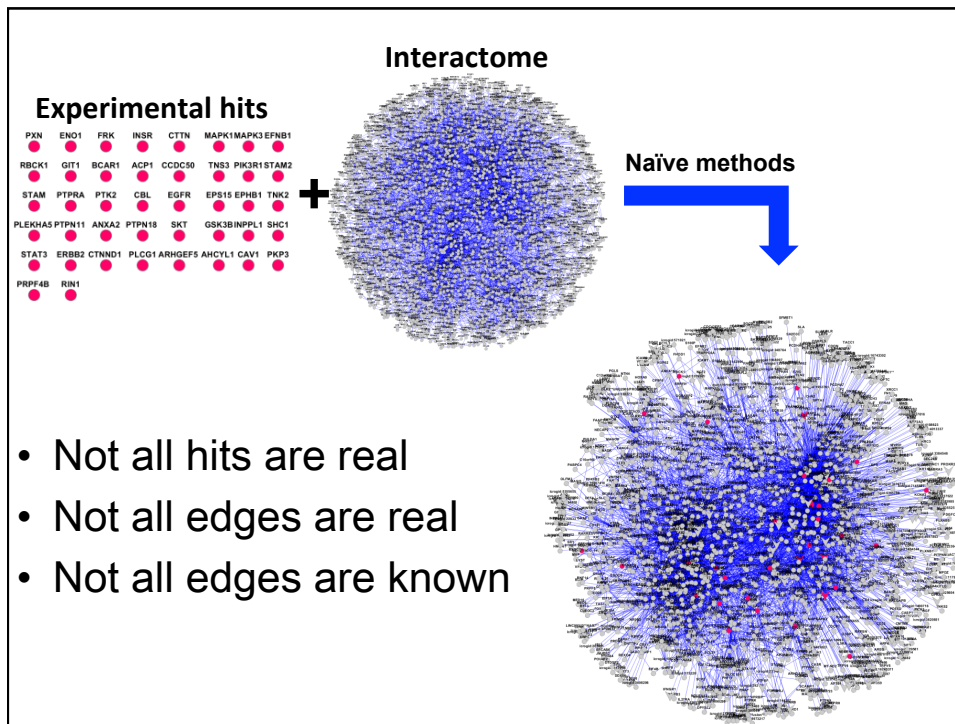


Huang et al., *PNAS* 2007







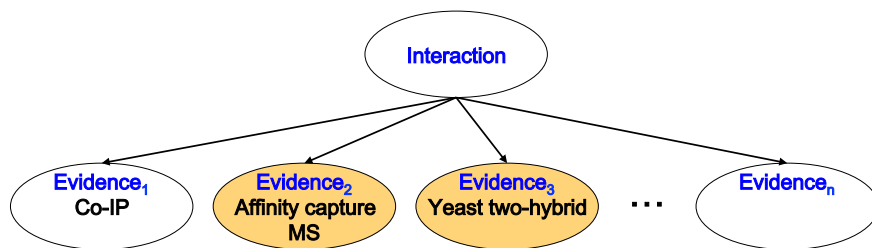
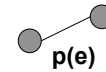


## A Better Approach

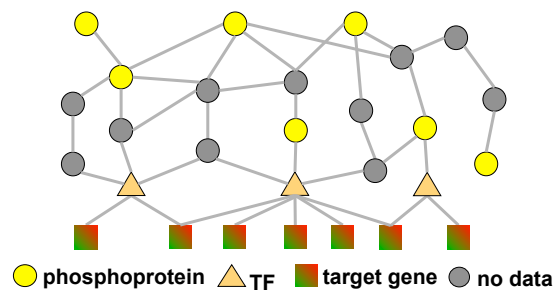
- Find most probable connections.
- **Only connect some** of the data to each other:
  - False positives in data
  - Missing interactions in interactome

## Probabilities Reflect Underlying Data

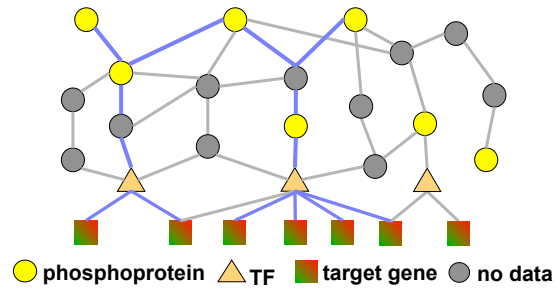
- Naïve Bayes approach integrates multiple types of evidence  
(Jansen *et al.* Science 2003; Myers et al. Genome Biology 2005.)
- Reliability determined by:
  - higher confidence experimental techniques
  - multiple experimental techniques



## Prize-collecting Steiner Tree

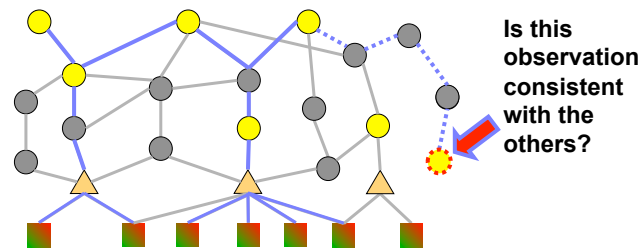


## Prize-collecting Steiner Tree



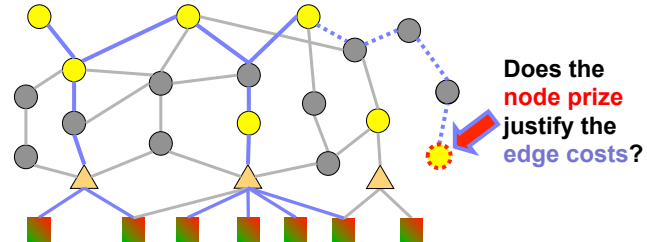
1. Don't Include All Data
2. Avoid Unlikely Interactions

## Prize-collecting Steiner Tree



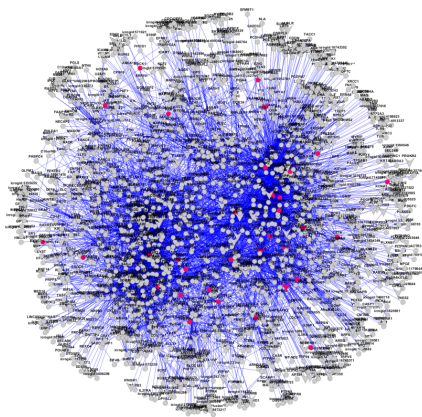
1. Don't Include All Data
2. Avoid Unlikely Interactions

## Prize-collecting Steiner Tree

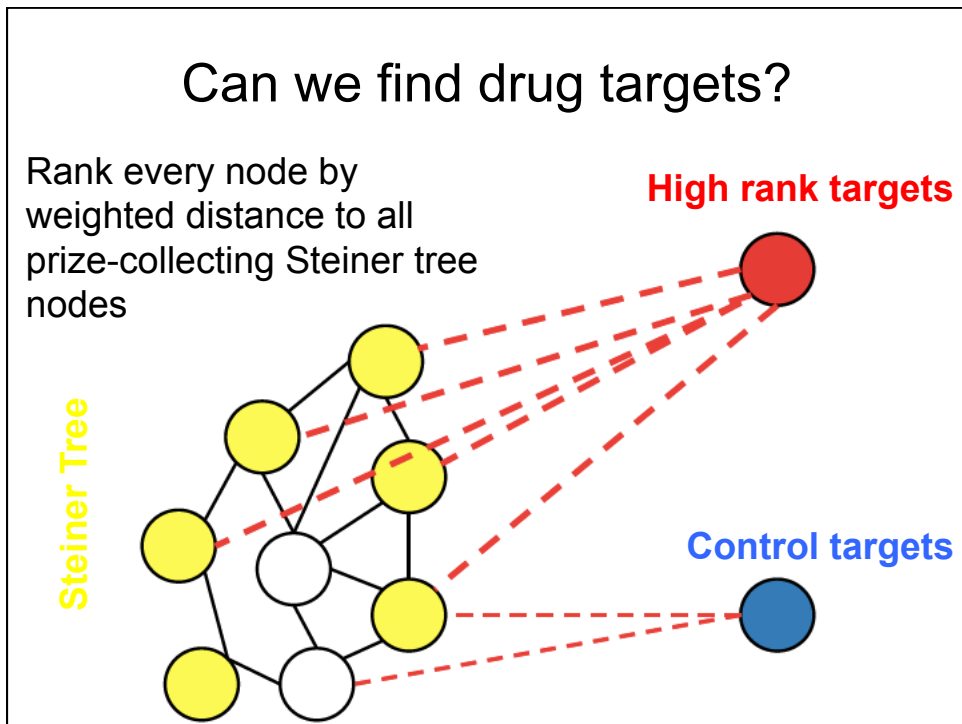
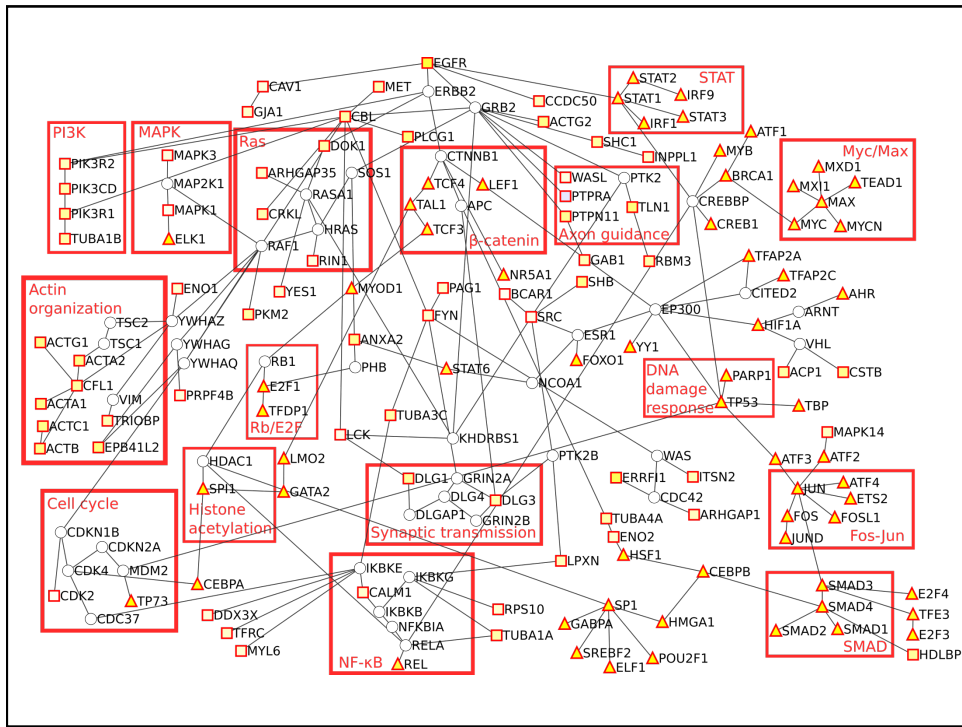


$$\sum_{v \text{ not in } T} \beta \text{prize}(v) + \sum_{e \text{ in } T} \text{cost}(e)$$

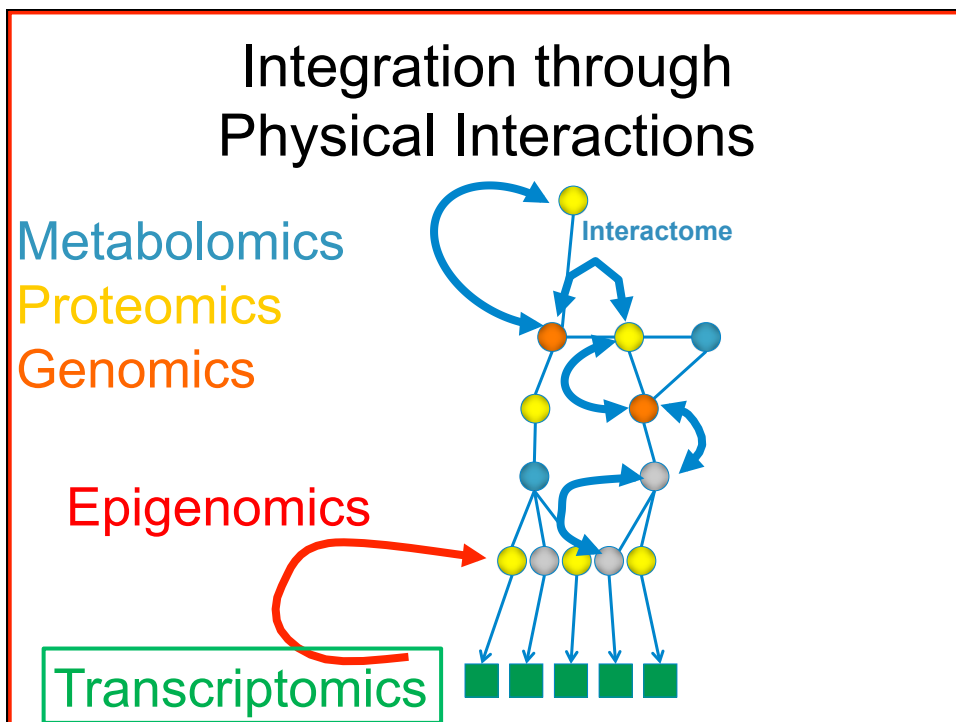
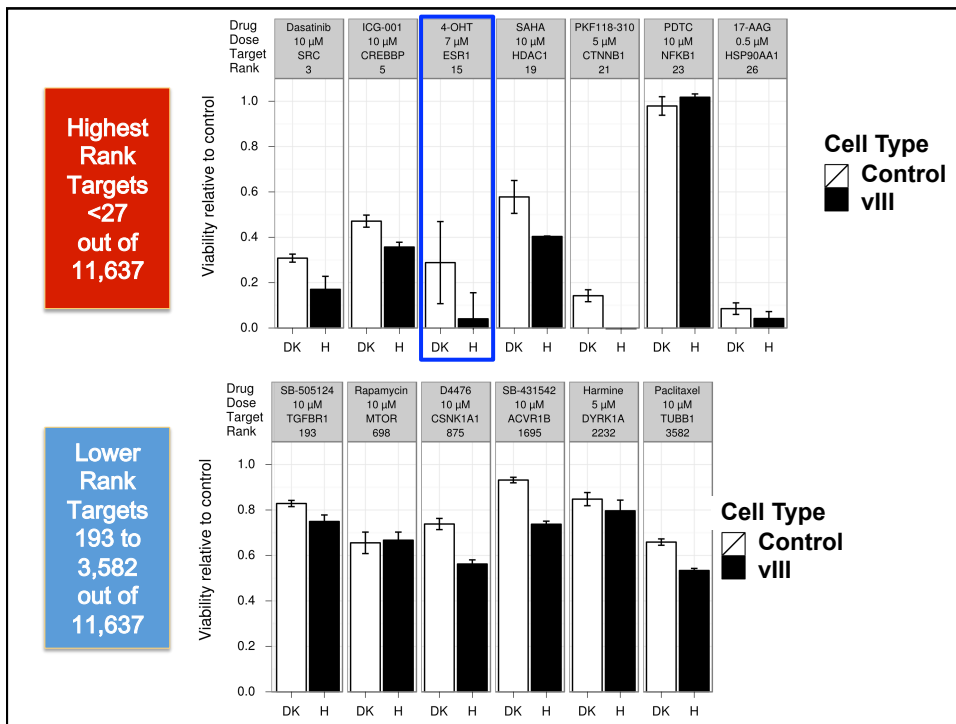
## Naïve Methods



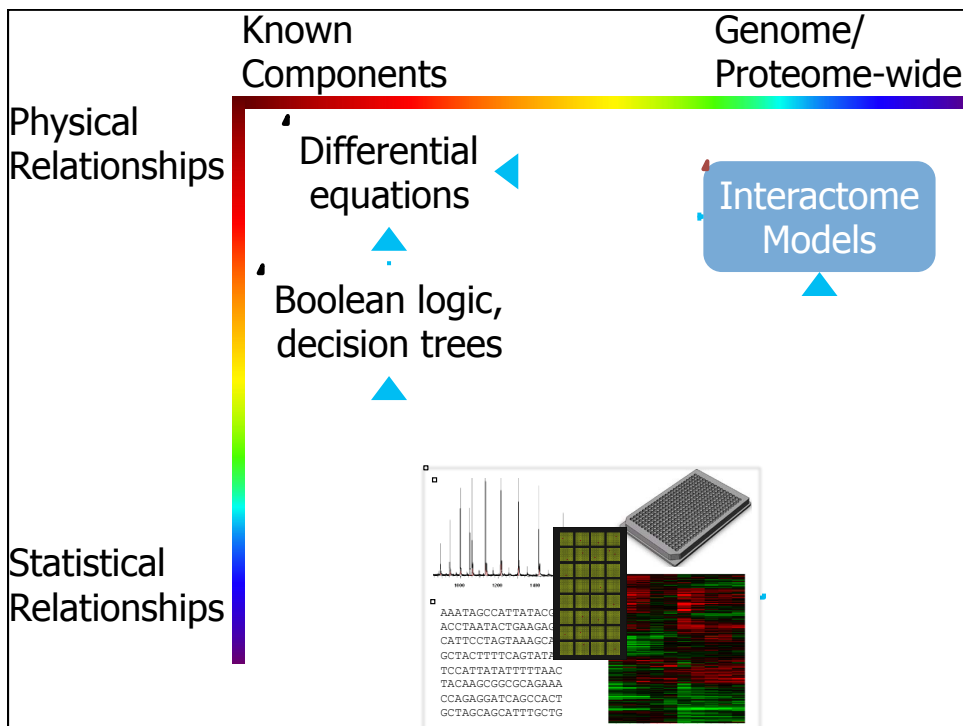
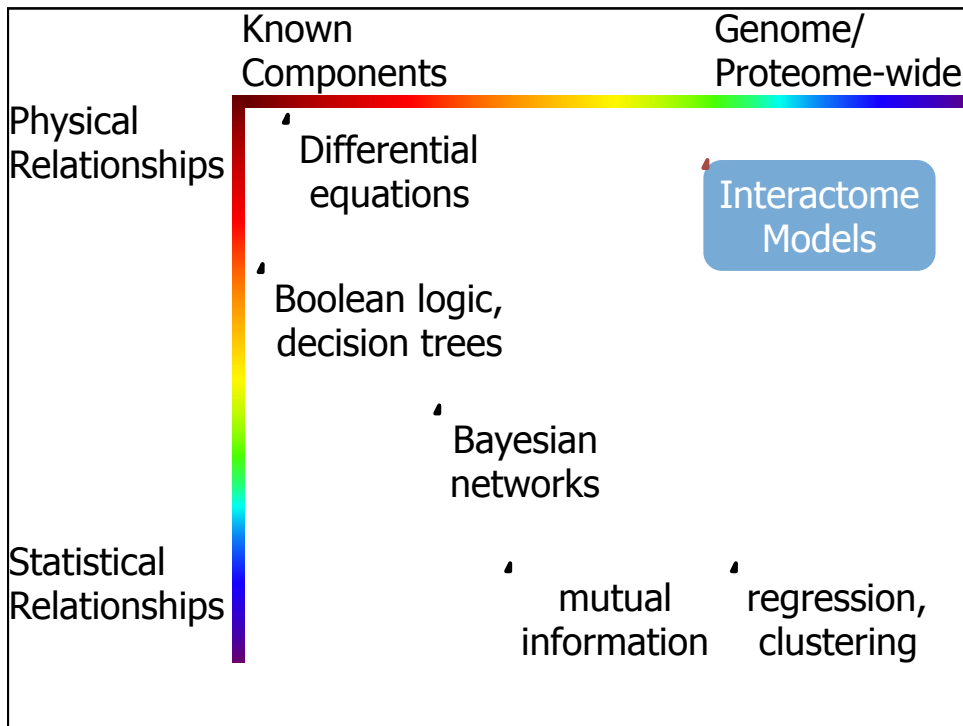
- >2,500 nearest neighbors of phosphoproteins
- >4,500 nearest neighbors of phosphoproteins +transcription factors

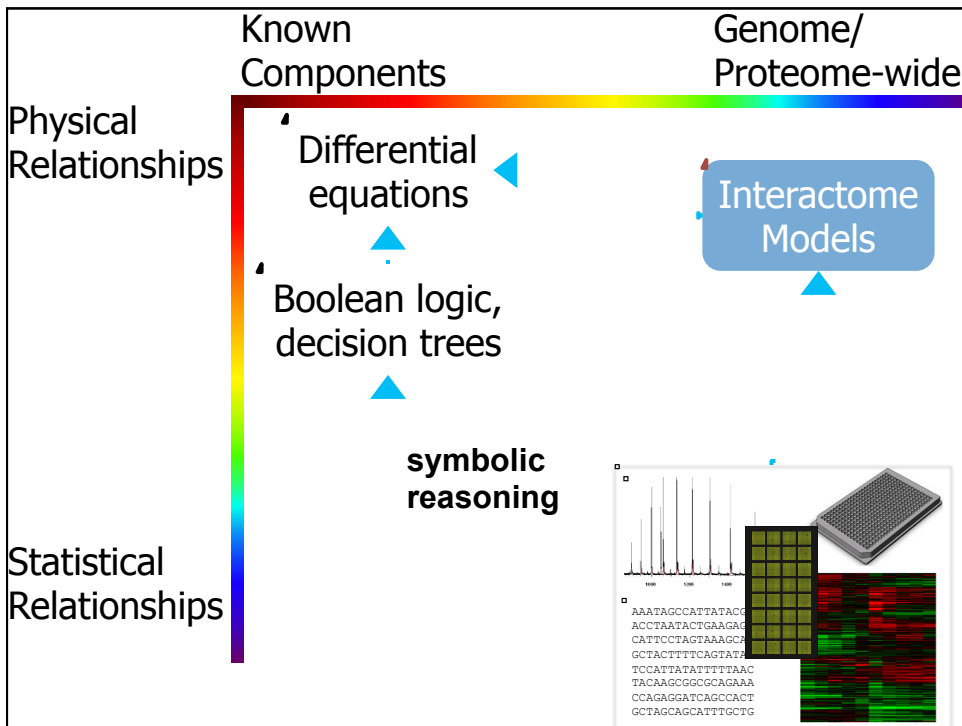




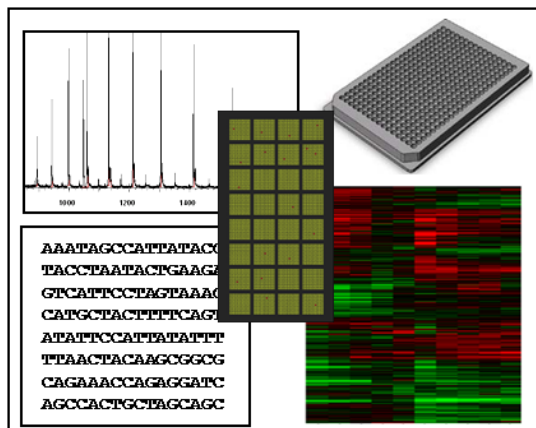








## Embracing the change



**THE FRAENKEL LAB**  
Biological Engineering

**Welcome to SteinerNet**

Revealing the Hidden Components in Regulatory and Signaling Networks by Integrating Proteomics, Transcriptome and Interactome Data

Sample files: [Interactome File](#), [Terminal file](#), and [Transcription factor to DNA file](#)

**Protein-Protein Interactions**

Upload an interactome here:

OR select a database interactome:

OR paste a list:

**Terminal Nodes with Penalties**

Upload terminals file here\*:

OR paste a list of proteins or genes with penalty values\*:

\*Please check the input format for terminals  
\*\*\*Terminal names must be consistent with naming in the interactome. Some external pages can be used to convert names if you need; such as DAVID, or HUGO

**(OPTIONAL) Transcription Factor to DNA Interactions**

Upload a file here:

OR select a database for TF to DNA interactions:

OR Paste the TF to DNA interactions:

Beta for Protein Terminals:

Beta for DNA Terminals:

**THE FRAENKEL LAB**  
Biological Engineering

**Omics Integrator**

Omics Integrator is package comprised of command-line tools designed to facilitate the integration of high-throughput datasets such as gene expression, phospho-proteomic data and the results from genetic screens. As shown below, garnet is used to identify transcription factors that give rise to gene expression changes using epigenetic data while forest integrates these data or other data by finding the minimum number of edges in a protein interaction network.

**Garnet**

Map epigenetic regions to expressed transcripts

Hypersensitive region (DNase-Seq) OR Epigenetic marks (ChIP-Seq) TSS

Hypersensitive DNA fragments OR Marked DNA fragments Expressed transcripts

Predict transcription factor binding using motifs

Regress transcription factor affinities to transcript changes

Epigenetic Regions

TF-affinity Score

TF Affinity Score

Select transcription factors as terminal nodes and assign prizes

**Forest**

Define terminal node set from experimental data and determine node prizes

Collect weighted interactome from literature or construct your own

Weight nodes in graph by experimental data

**Result**



Leslie  
Thompson



Forest  
White



David  
Housman



Jennifer  
Chayes



Christian  
Borgs



Riccardo  
Zecchina

