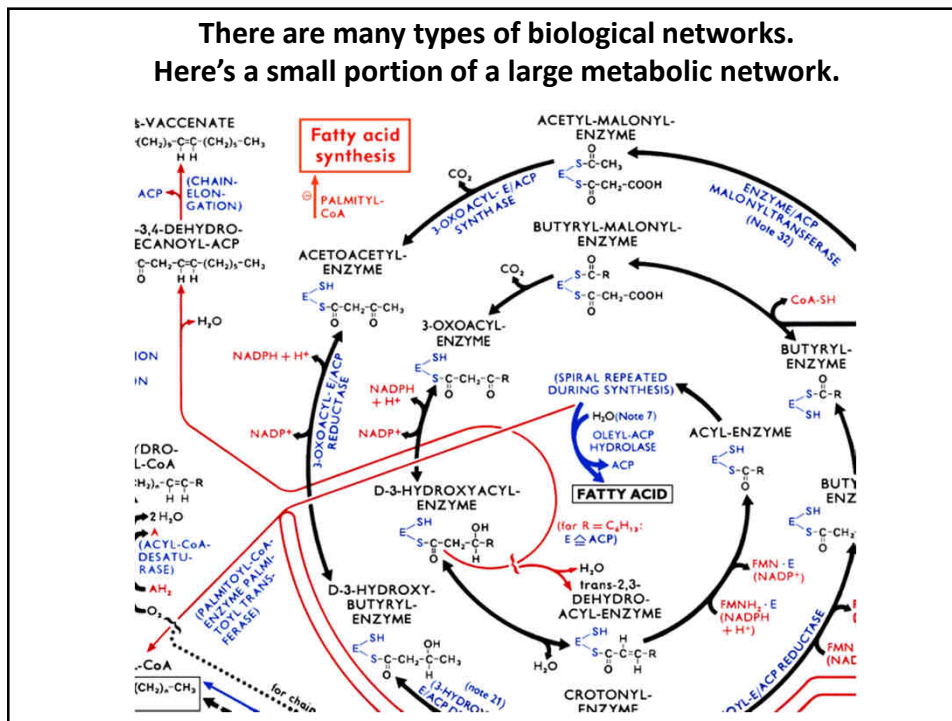
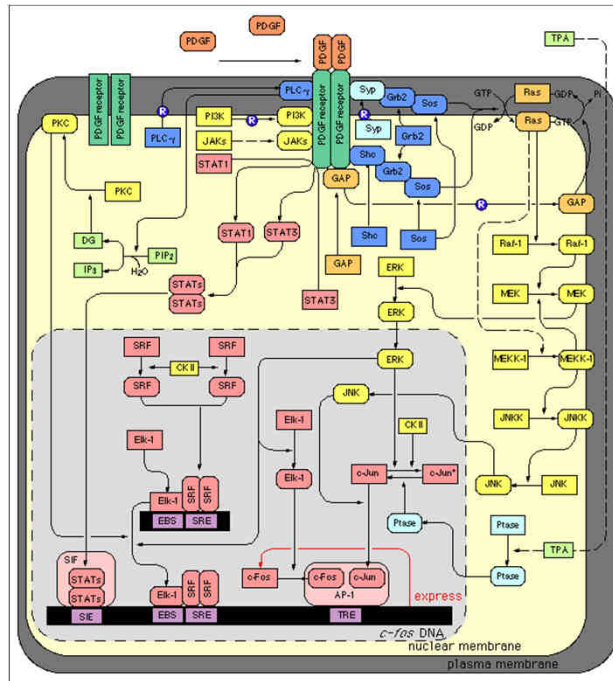


Network biology (& predicting gene function)

BCH339N Systems Biology / Bioinformatics
Edward Marcotte, Univ of Texas at Austin

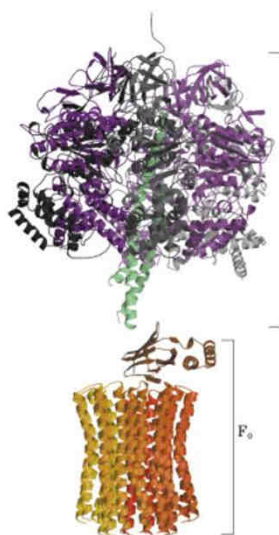


A typical
genetic
network

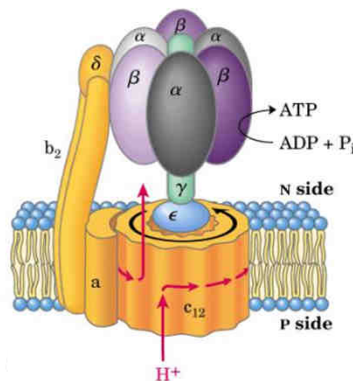


Contacts between proteins define protein interaction networks

X-ray structure
of ATP synthase

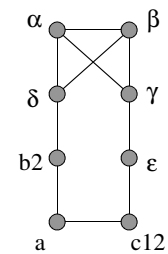


Schematic
version



Total set = protein complex
Sum of **direct** + **indirect**
interactions

Network
representation

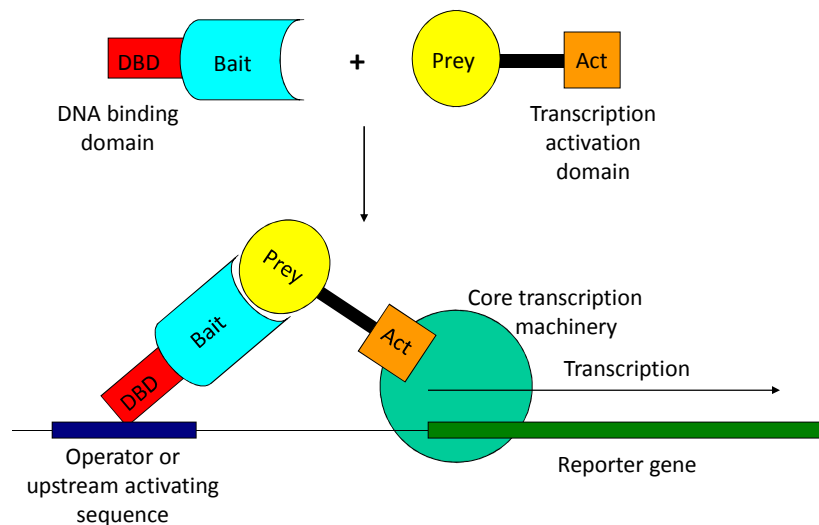


Let's look at some of the types of interaction data in more detail.

Some of these capture physical interactions, some genetic, some informational or logical.

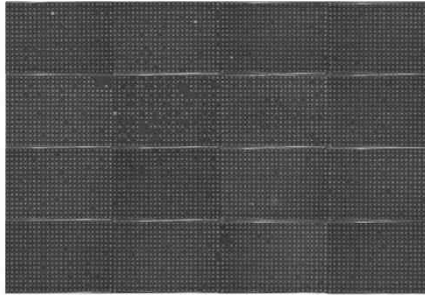
Pairwise protein interactions

In general, purifying proteins one at a time, mixing them, and assaying for interactions is far too slow & laborious. We need something faster! Hence, high-throughput screens, e.g. yeast two-hybrid assays

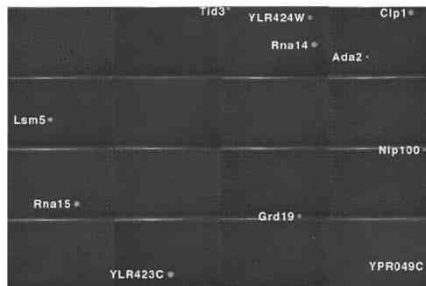


High-throughput yeast two-hybrid assays

Haploid yeast cells expressing activation domain-prey fusion proteins



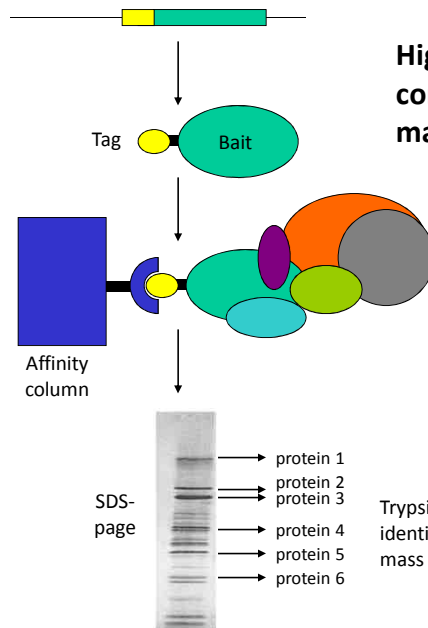
Diploid yeast probed with DNA-binding domain-Pcf11 bait fusion protein

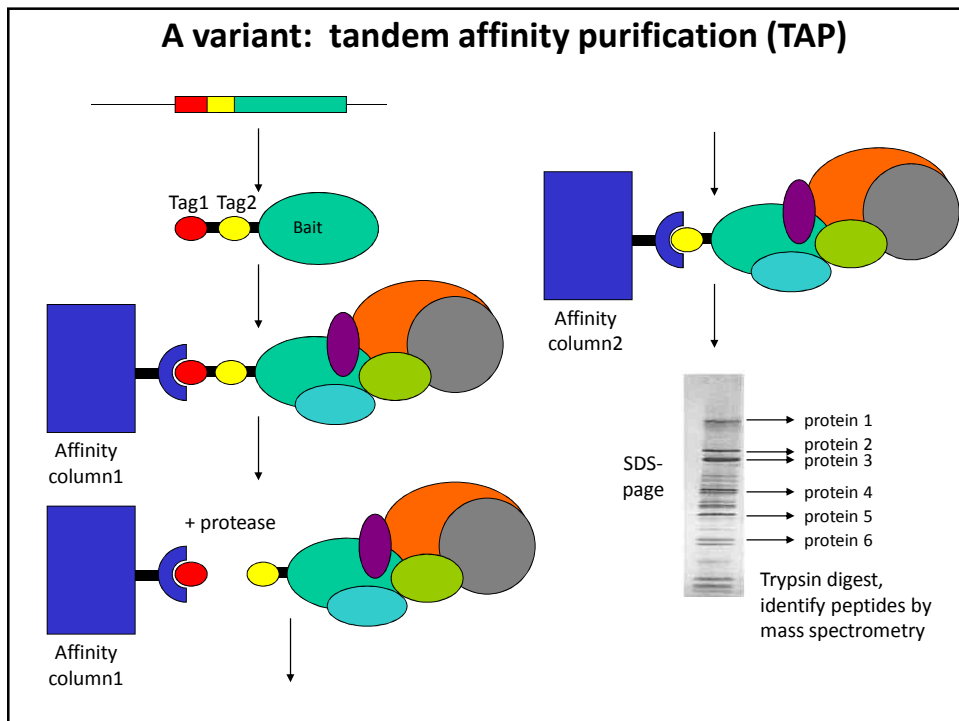
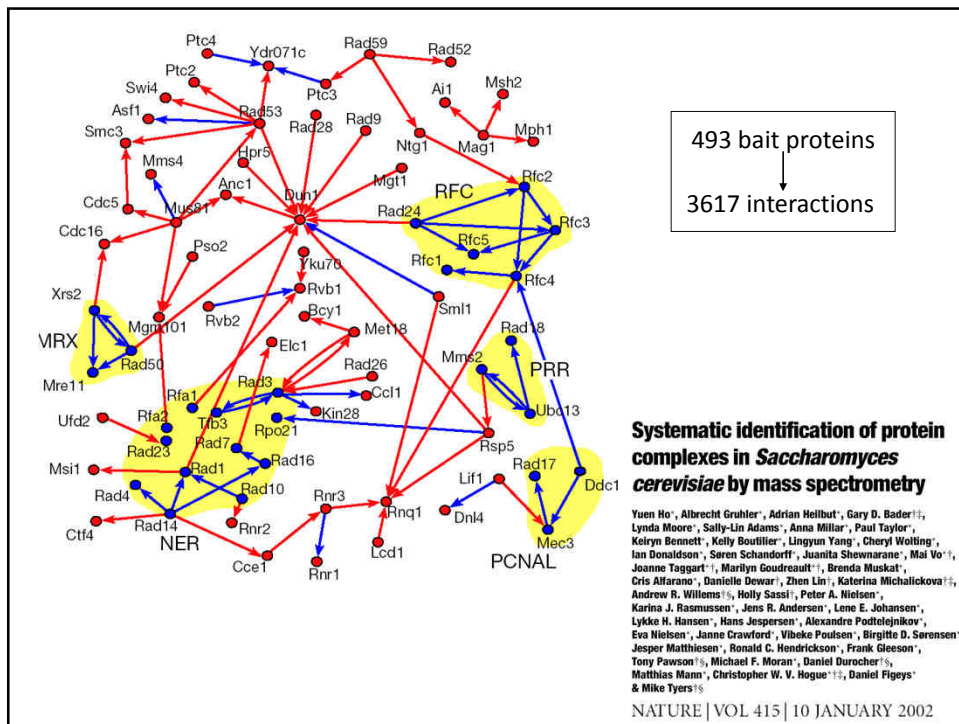


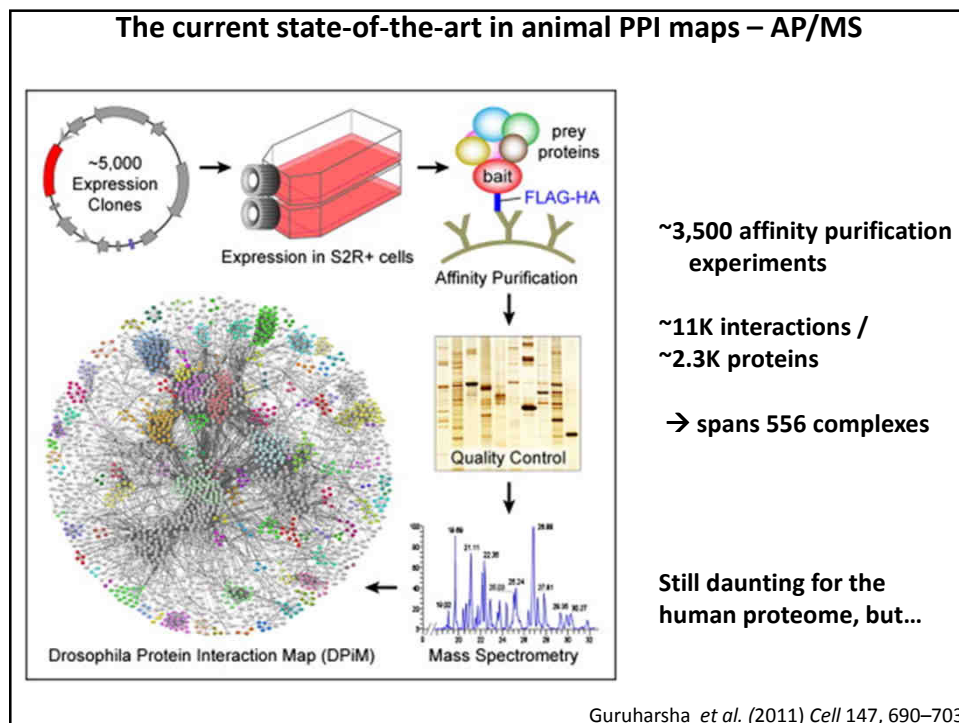
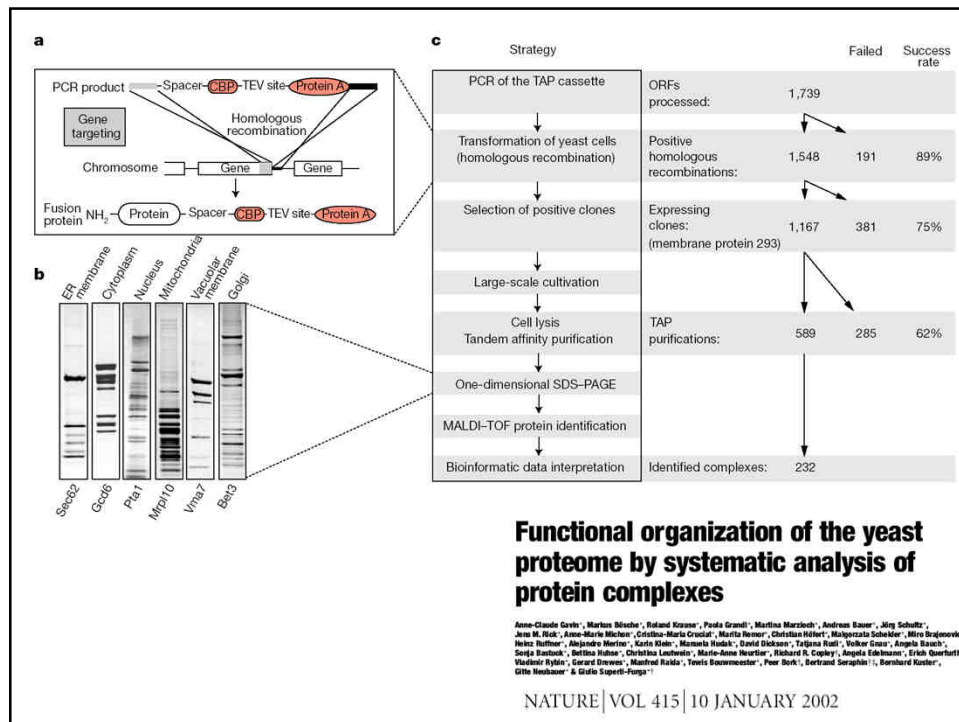
Uetz, Giot, *et al. Nature* (2000)

Protein complexes

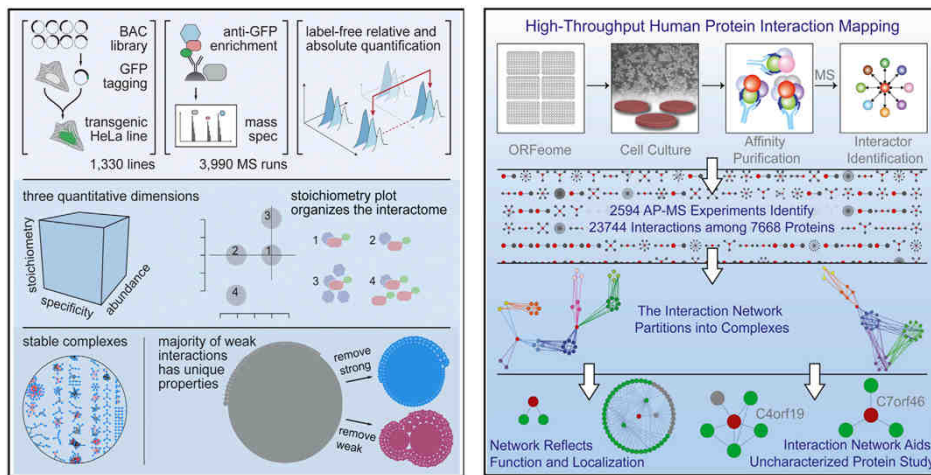
High-throughput complex mapping by mass spectrometry







The current state-of-the-art in human PPI maps – large scale AP/MS



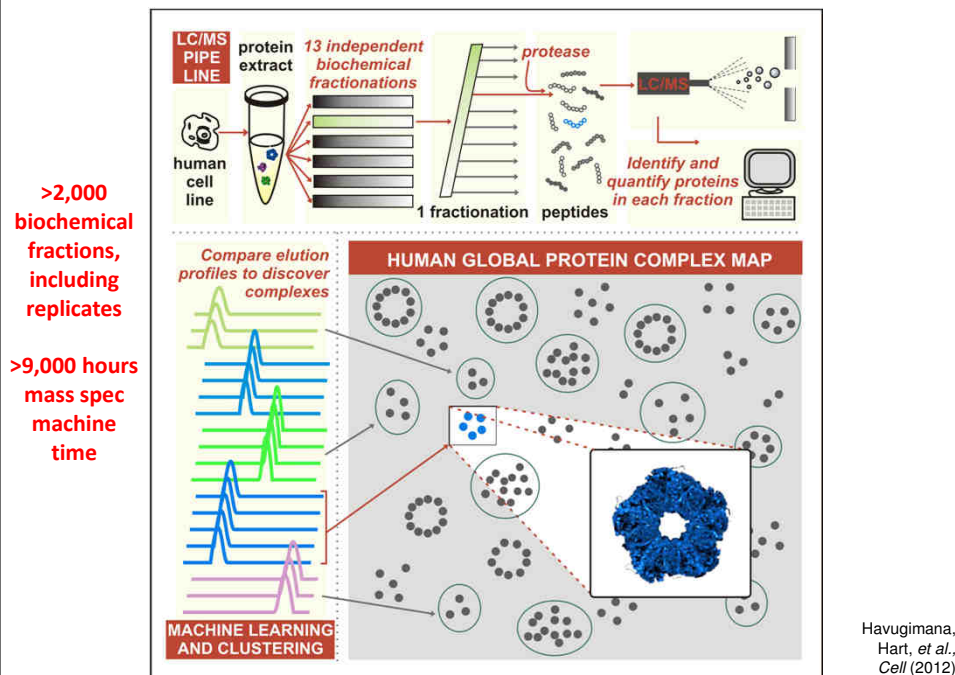
Hein *et al.*, *Cell* (2015) 163:712-23.

Huttlin *et al.*, *Cell* (2015) 162:425-440

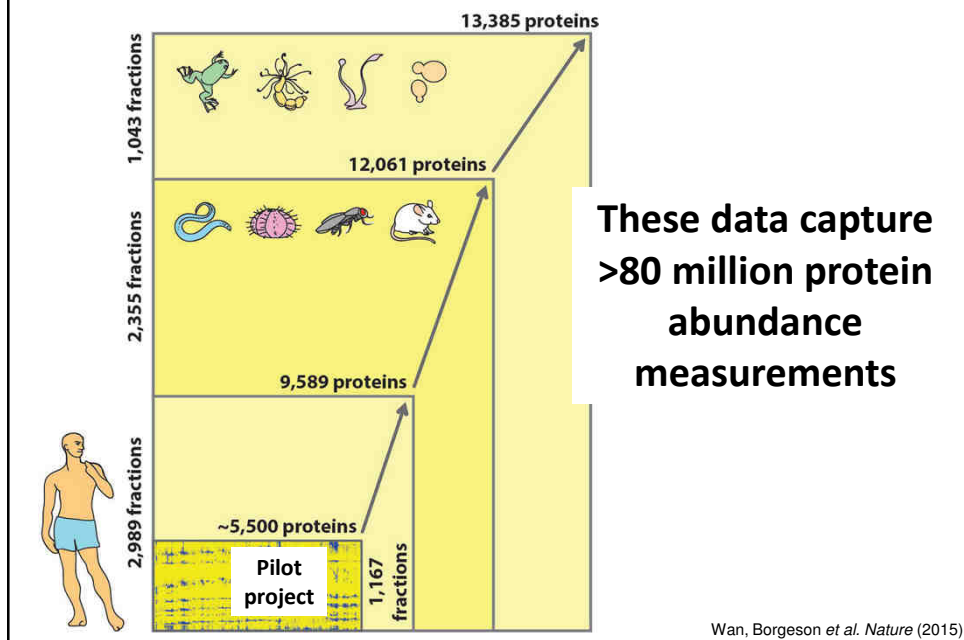
Huttlin *et al.*, *Nature* (2017) 545:505-509

Just in the past 3 years, nearly 6K affinity purification experiments on tagged human proteins expressed in cell lines

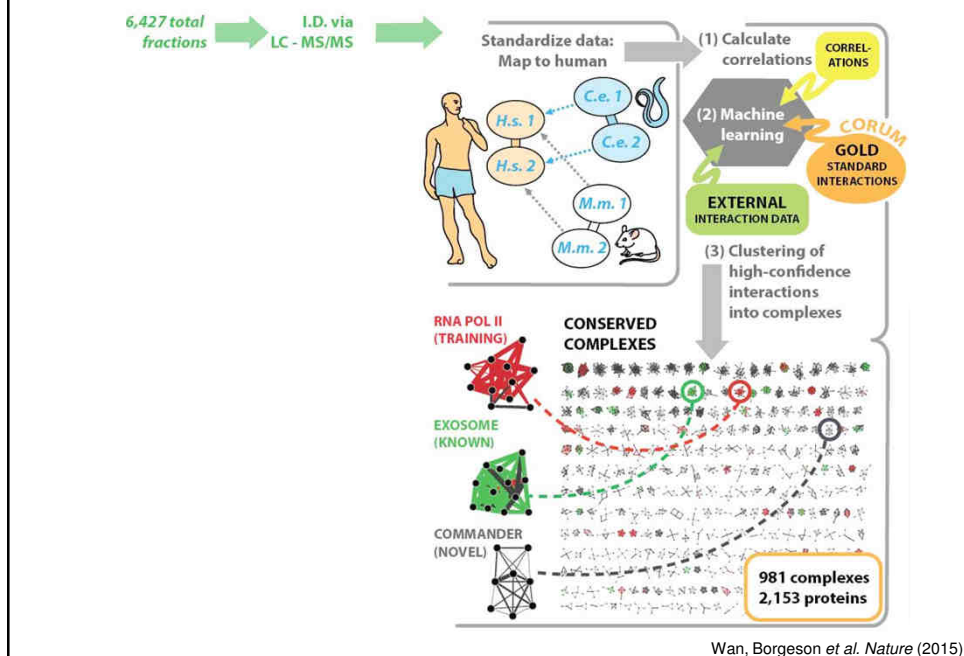
The current state-of-the-art in animal PPI maps – co-fractionation/MS



Now >6,400 CF/MS experiments across animals

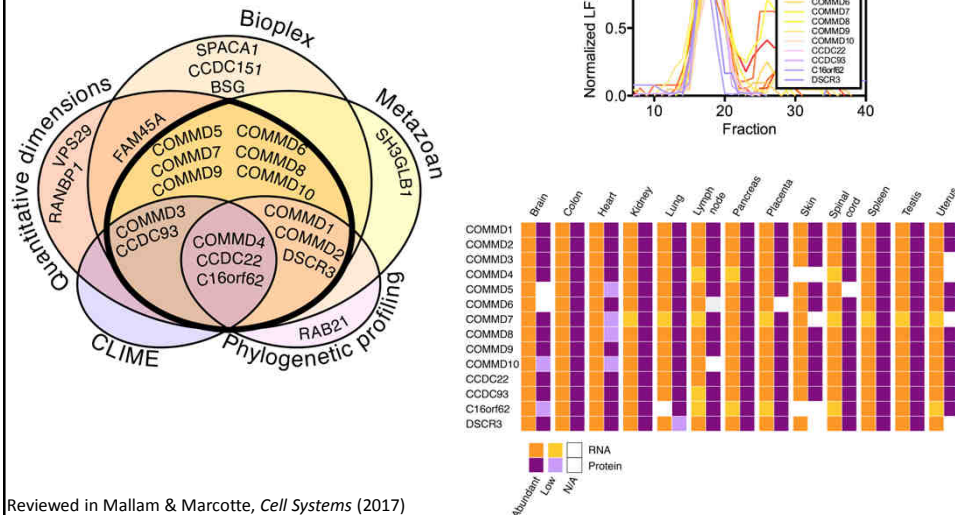


Extending the map across animals...



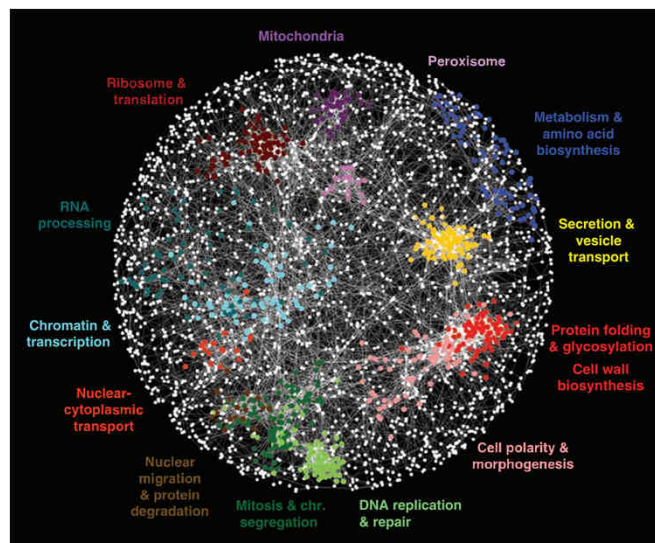
There are still lots of cellular machines left to find

e.g. the “Commander” complex, found in all 3 recent human PPI maps, a 600 kDa protein complex expressed in nearly every human cell type and tissue



Genetic interactions

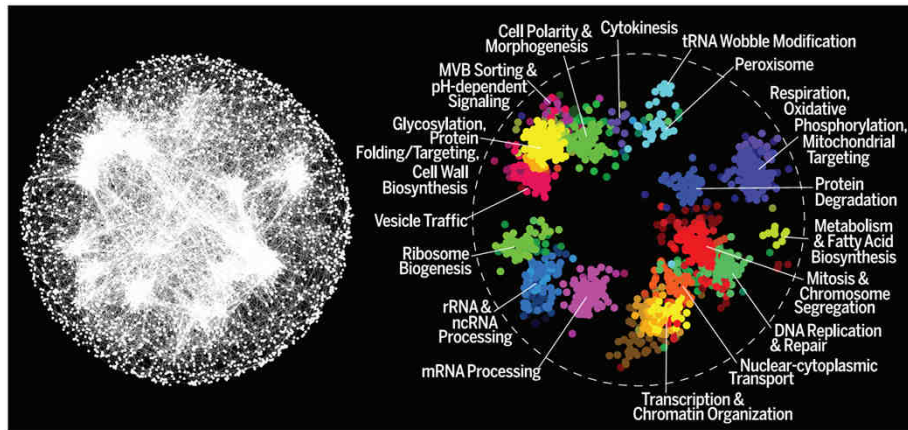
5.4 million gene-gene pairs assayed for synthetic genetic interactions in yeast



Costanzo *et al.*, *Science* 327: 425 (2010)

Genetic interactions, the 2016 version

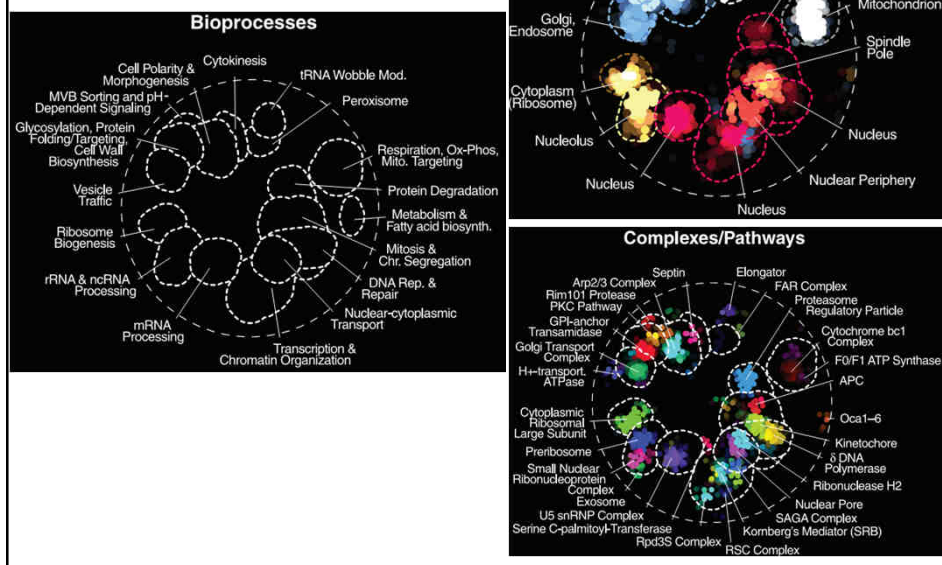
23 million gene-gene pairs assayed for synthetic genetic interactions in yeast, identifying ~550,000 negative and ~350,000 positive genetic interactions



A global network of genetic interaction profile similarities. (Left) Genes with similar genetic interaction profiles are connected in a global network, such that genes exhibiting more similar profiles are located closer to each other, whereas genes with less similar profiles are positioned farther apart. (Right) Spatial

Costanzo *et al.*, *Science* 353: 1381 (2016)

The global genetic interaction profile similarity network reveals a hierarchy of cellular function.



Comparative genomics

Functional relationships between genes impose subtle constraints upon genome sequences. Thus, genomes carry intrinsic information about the cellular systems and pathways they encode.

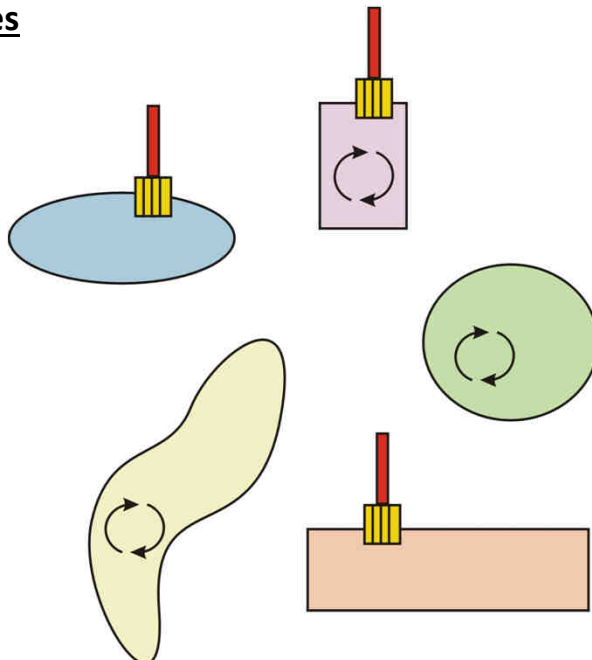
Linkages can be found from aspects of gene context, including:

- Distances between sequence elements
- Order of sequences
- Variation in order between organisms
- Regulatory sequences near genes
- Gene content of an organism
- Variation in gene content between organisms
- Fusions between genes from different organisms

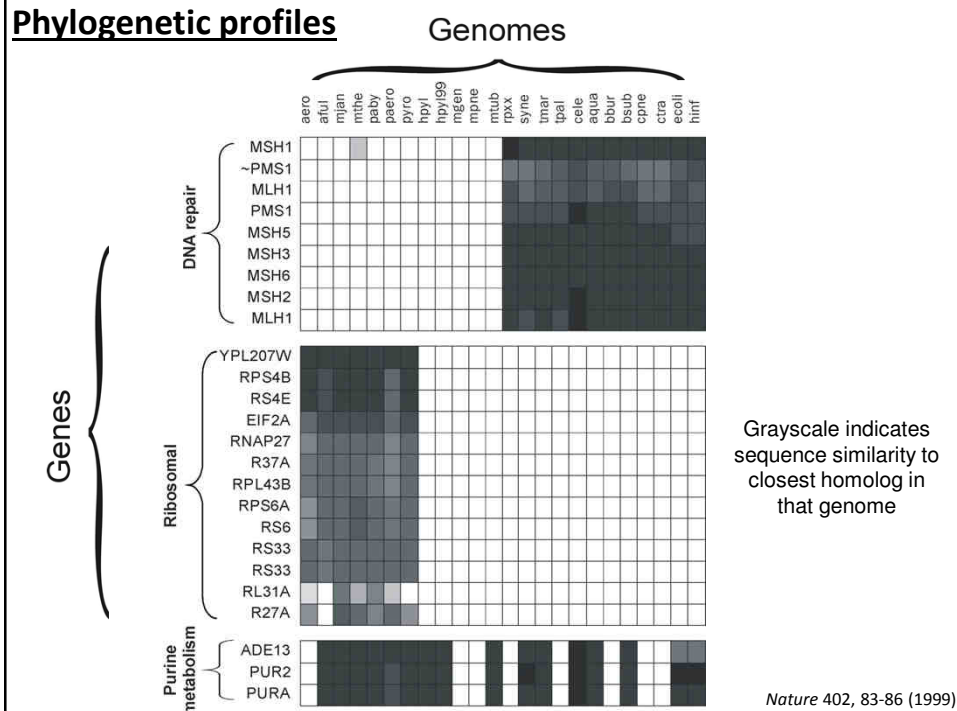
Phylogenetic profiles

Organisms with e.g. a flagellum have the necessary genes; those without tend to lack them.

Specific trends of gene presence/absence thus inform about biological processes.

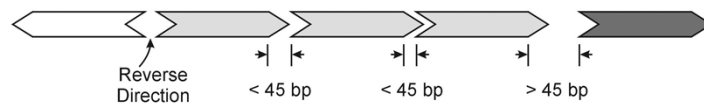


PNAS 96, 4285-4288 (1999)

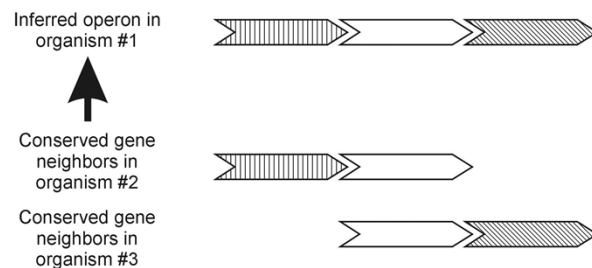


Operons and evolutionary conservation of gene order

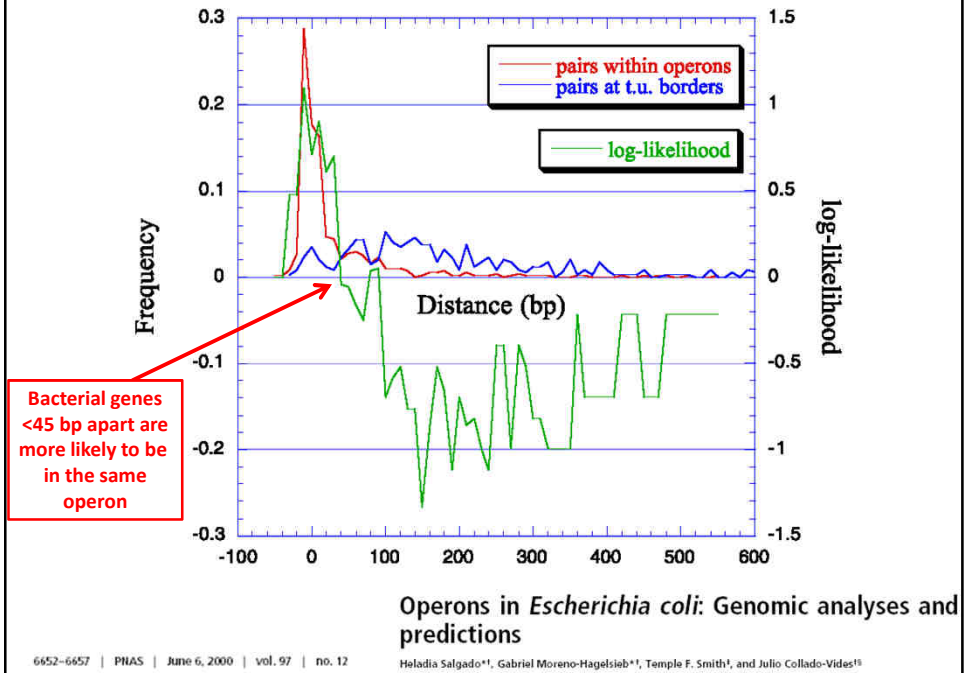
Prokaryotic operons tend to favor certain intergenic distances



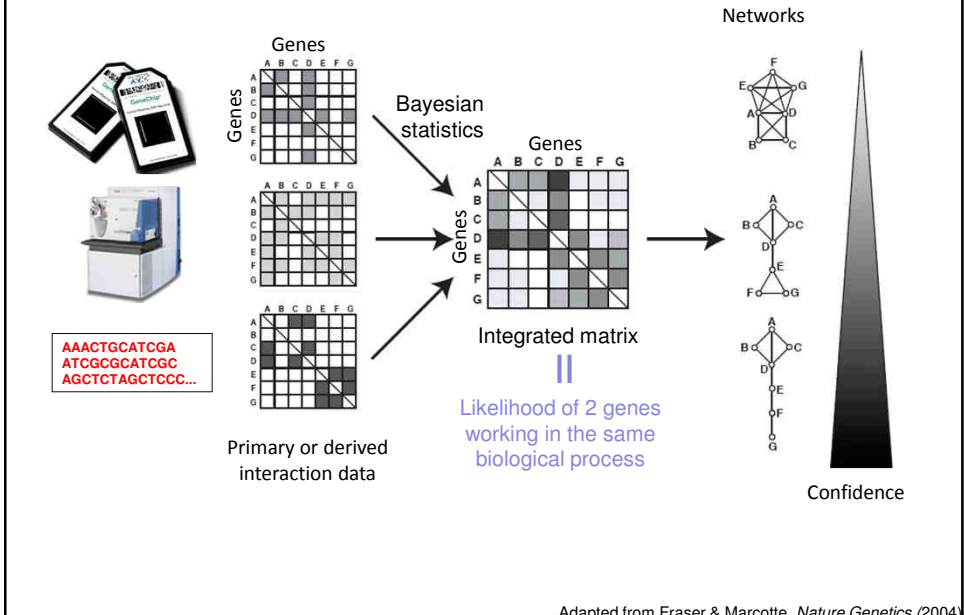
Conserved gene neighbors also reveal functional relationships

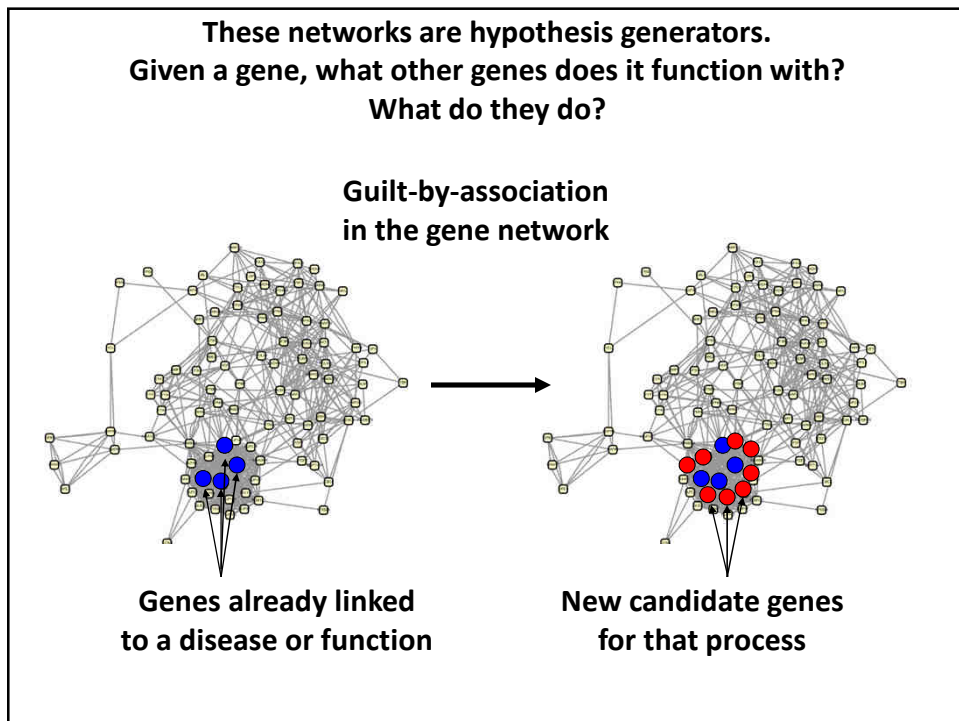
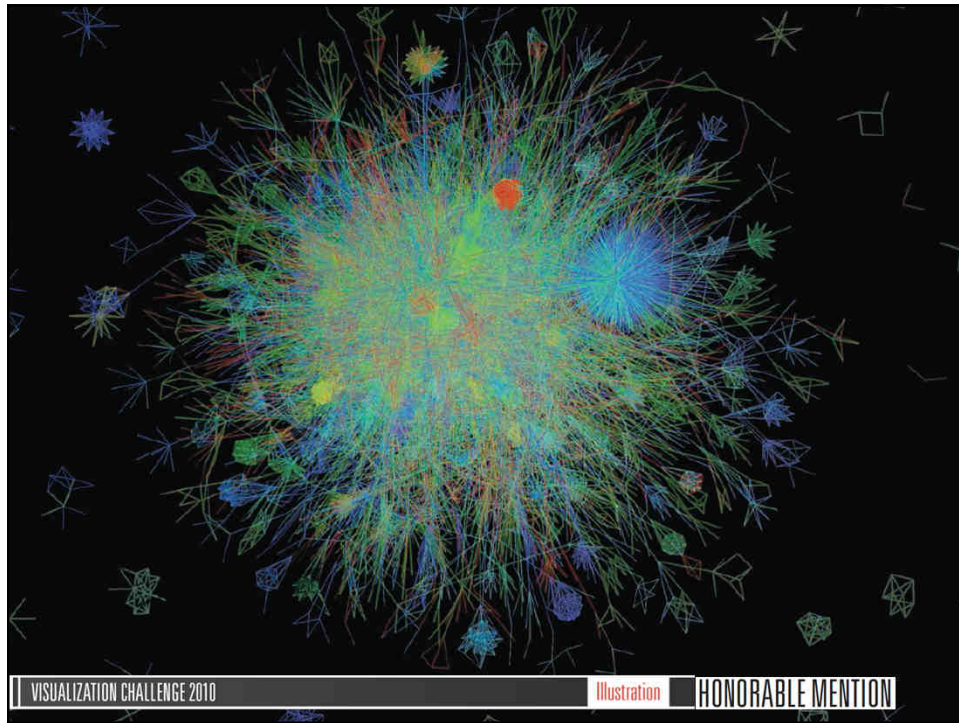


Again, such observations can be turned into pairwise scores:

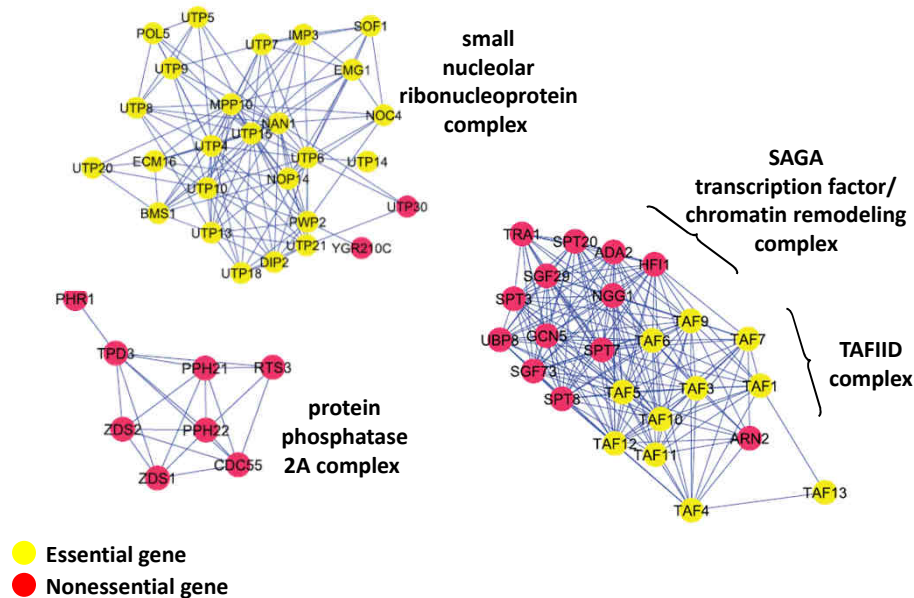


These sorts of data can be combined into functional gene networks



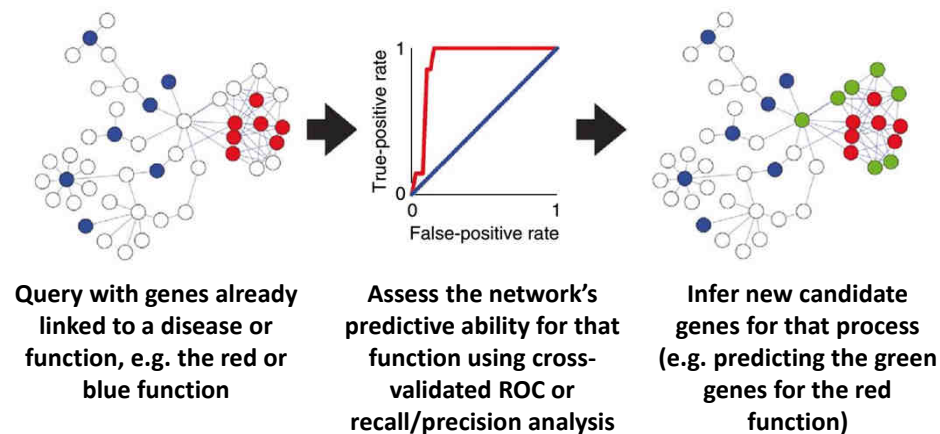


Gene networks frequently reflect functions, pathways, & phenotypes, e.g., lethality in yeast is linked to the molecular machine, not the gene



We can propagate annotations across the graph to infer new annotations for genes (network “guilt-by-association”, or GBA).

Testing how well this works on hidden, but known, cases let’s us measure how predictive it will be for new cases.

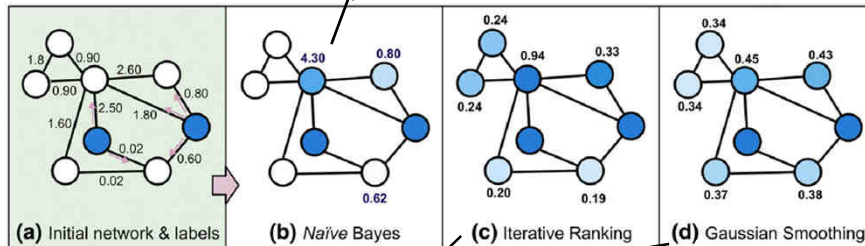


Lee, Ambaru *et al. Nature Biotechnology* 28:149-156 (2010)

Numerous algorithms exist for network GBA

Naïve Bayes assigns scores to neighboring nodes based on edges

Similar to Google's personalized PageRank



Network diffusion algorithms start with initial annotations and the graph topology, then propagate initial scores across the network, e.g. Gaussian smoothing tries to find scores:

$$f^{final} = \underset{f}{\operatorname{argmin}} \alpha \sum_i (f_i - f_i^0)^2 + (1 - \alpha) \sum_i \sum_j w_{ij} (f_i - f_j)^2$$

minimizing the difference between final and initial scores of a protein

& between a protein's score and that of each of its neighbors

Reviewed in Wang & Marcotte, *J Proteomics* (2010)

Calculating ROC curves

		Actual	
		P	N
Prediction	P'	True Positive	False Positive
	N'	False Negative	True Negative

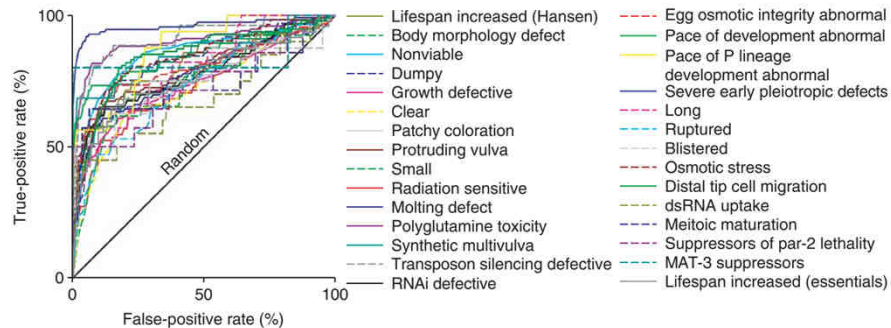
Basic idea: sort predictions from best to worst, plot TPR vs. FPR as you traverse the ranked list

$$\begin{aligned} \text{TPR} &= \text{TP} / P = \text{TP} / (\text{TP} + \text{FN}) \\ &= \text{True Positive Rate} \\ &= \text{Sensitivity, Recall} \end{aligned}$$

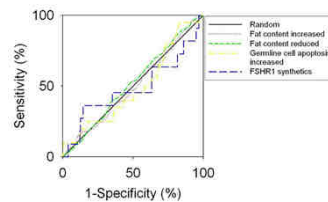
$$\begin{aligned} \text{FPR} &= \text{FP} / N = \text{FP} / (\text{FP} + \text{TN}) \\ &= \text{False Positive Rate} \\ &= 1 - \text{Specificity} \end{aligned}$$

Also useful to plot Precision [= TP / (TP + FP)] vs. Recall (= TPR)

For example, predicting genes linked with worm phenotypes in genome-wide RNAi screens



Some very poorly predicted pathways:



ROC analysis indicates the likely predictive power of the network for a system of interest.

A poor ROC → no better than random guessing.

Lee, Lehner *et al.*, *Nat Genet*, 40(2):181-8 (2008)

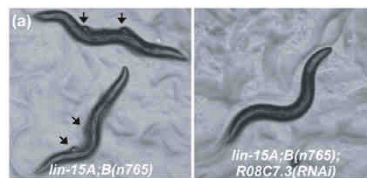
Remarkably, this strategy works quite well

Some examples of network-guided predictions:

In worms:

Genes that can reverse 'tumors' in a nematode model of tumorigenesis

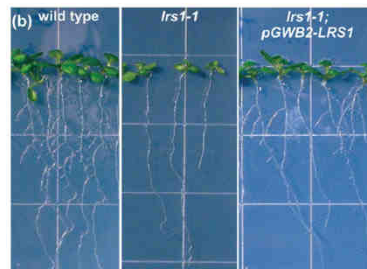
Lee, Lehner *et al.*
Nature Genetics (2008)



In Arabidopsis:

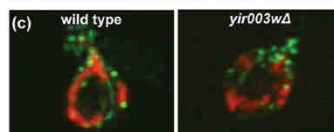
New genes regulating root formation

Lee, Ambaru *et al.*
Nature Biotech (2010)



In yeast: New mitochondrial biogenesis genes

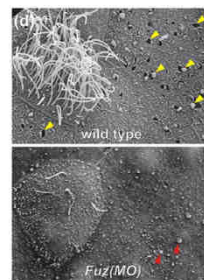
Hess *et al.*, *PLoS Genetics* (2009)



In mice/frogs:

Functions for a birth defect gene

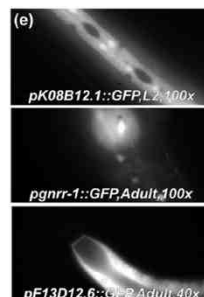
Gray *et al.*, *Nature Cell Biology* (2009)



In worms:

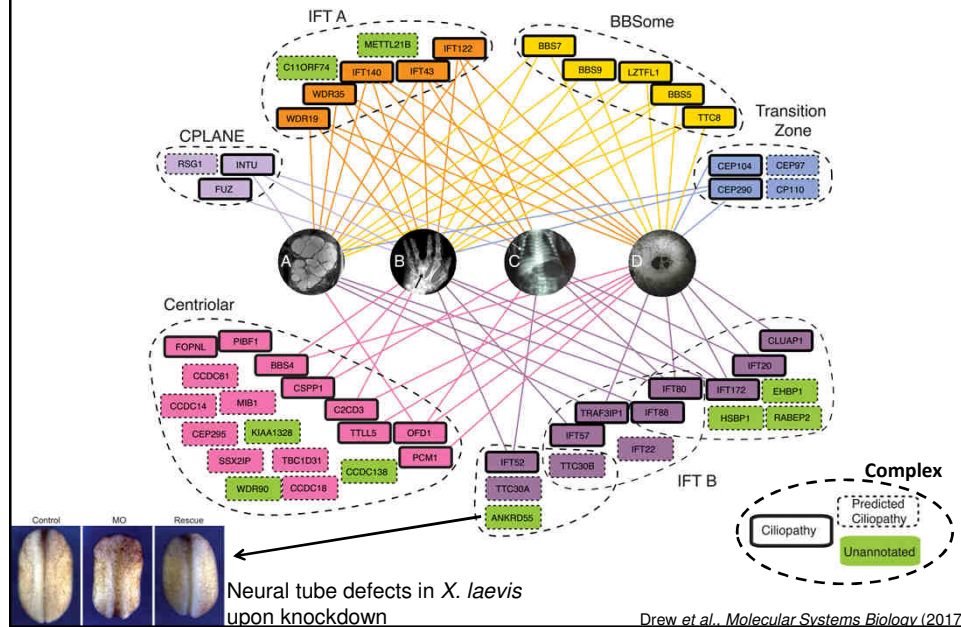
Predicting tissue specific gene expression

Chikina *et al.*, *PLoS Comp Biology* (2009)



Reviewed in Wang & Marcotte, *J Proteomics* (2010)

We use this approach routinely, e.g. a recent example predicting new ciliopathy genes from protein complexes



**Live demo of
STRING, BioGRID,
GeneMania,
functional networks
and Cytoscape**