

Supervised Maximum Likelihood Weighting of Composite Protein Networks for Protein Complex Prediction

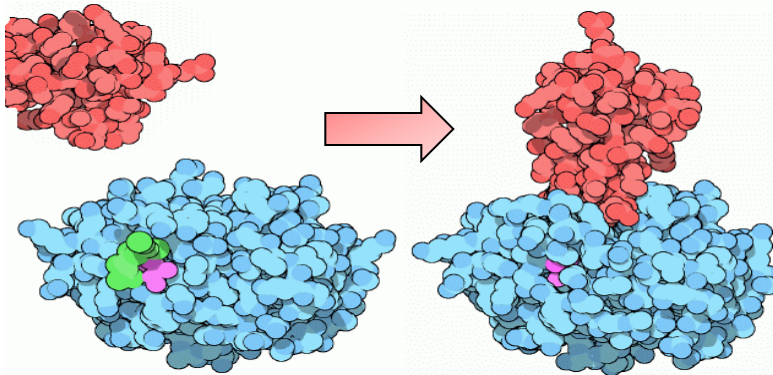
Limsoon Wong



Outline

- **Background**
- **Supervised Weighting of Composite Networks**
 - Data integration
 - Supervised edge weighting
 - Clustering
- **Results**
 - Prediction accuracy
 - Semantic coherence
 - Examples
- **Conclusion**

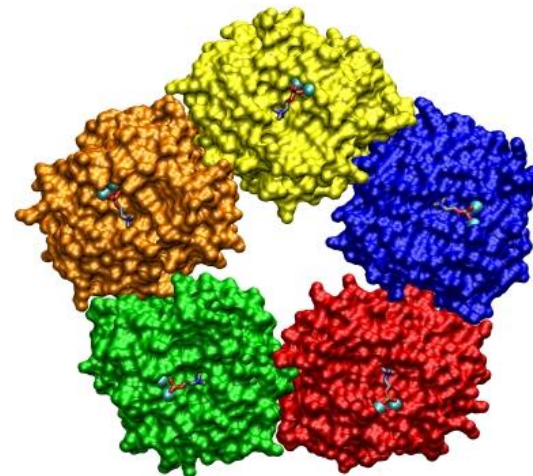
“Assemblies” of Interacting Proteins



Individual proteins come together and interact

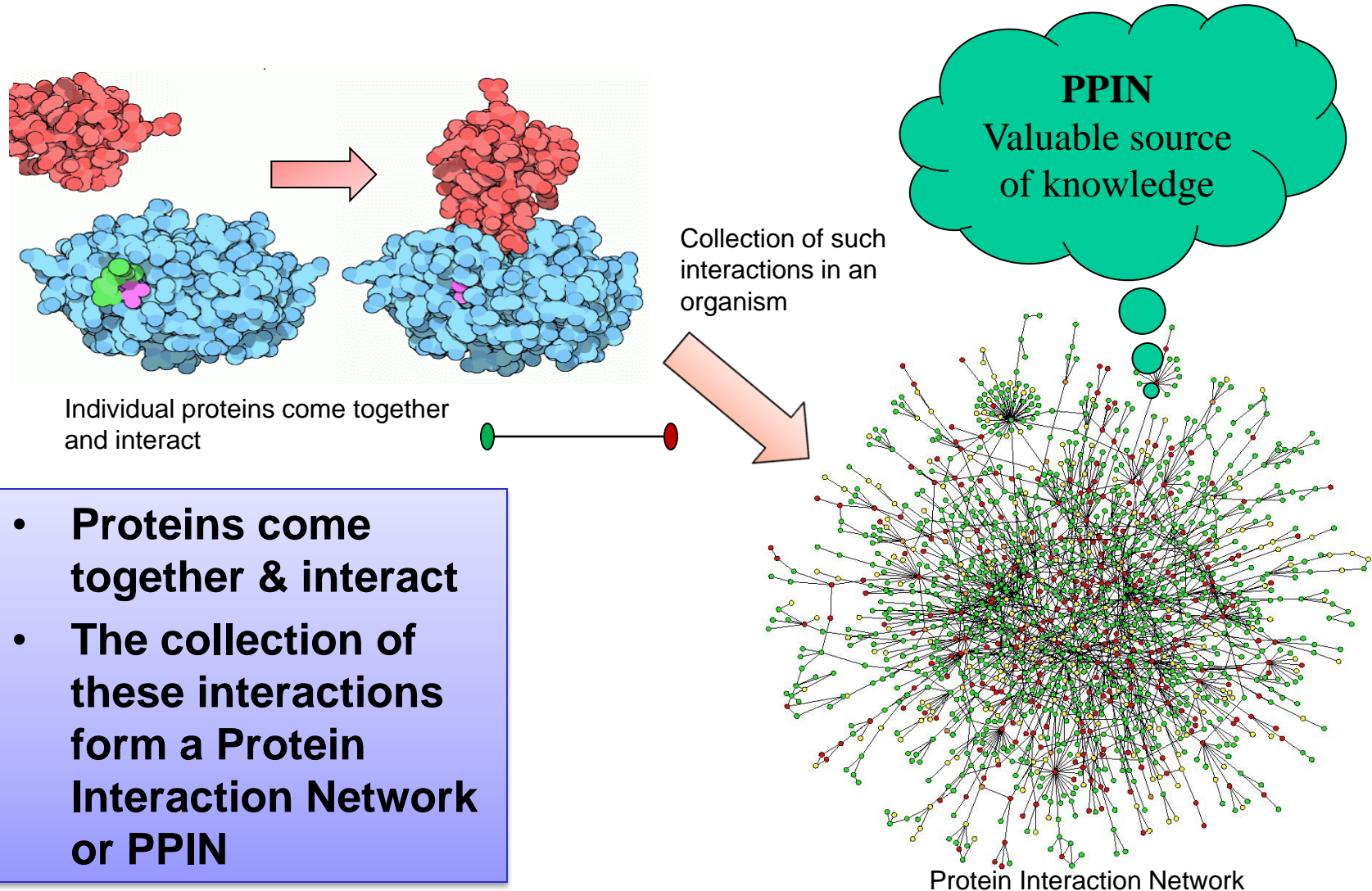
- **Proteins interact to form “protein assemblies”**
- **These assemblies are like “protein machines”**
 - Highly coordinated parts
 - Highly efficient

- **Protein assemblies**
 - Complexes
 - Functional modules
 - Intricate, ubiquitous, control many biological processes

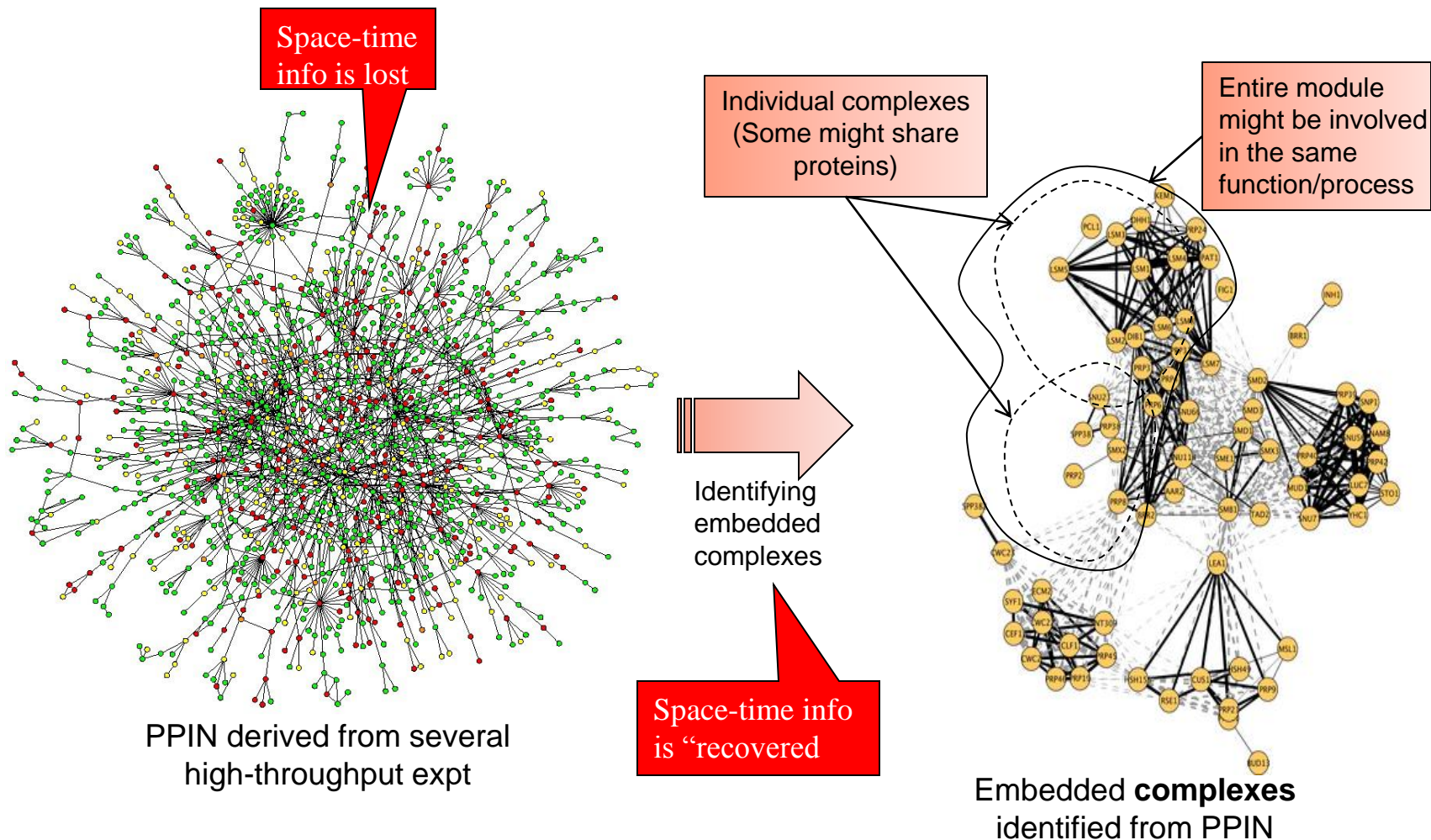


Protein assembly of multiple proteins

Protein Interaction Networks



Detection & Analysis of Protein Complexes in PPIN



Difficulties

- **Typical complex discovery method: Predict dense subgraphs in PPIN as complexes**

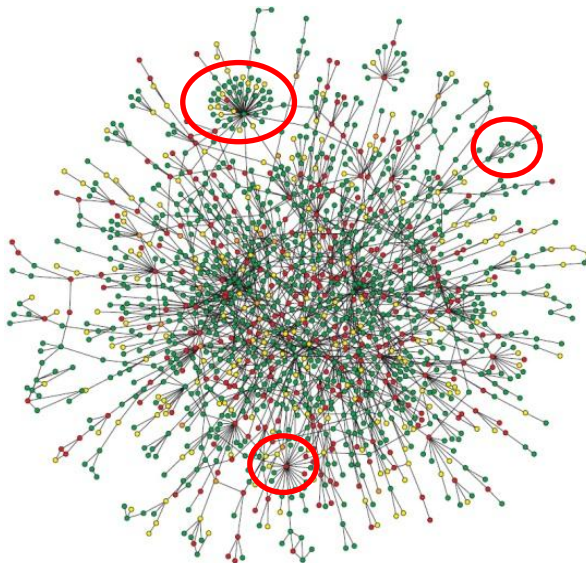
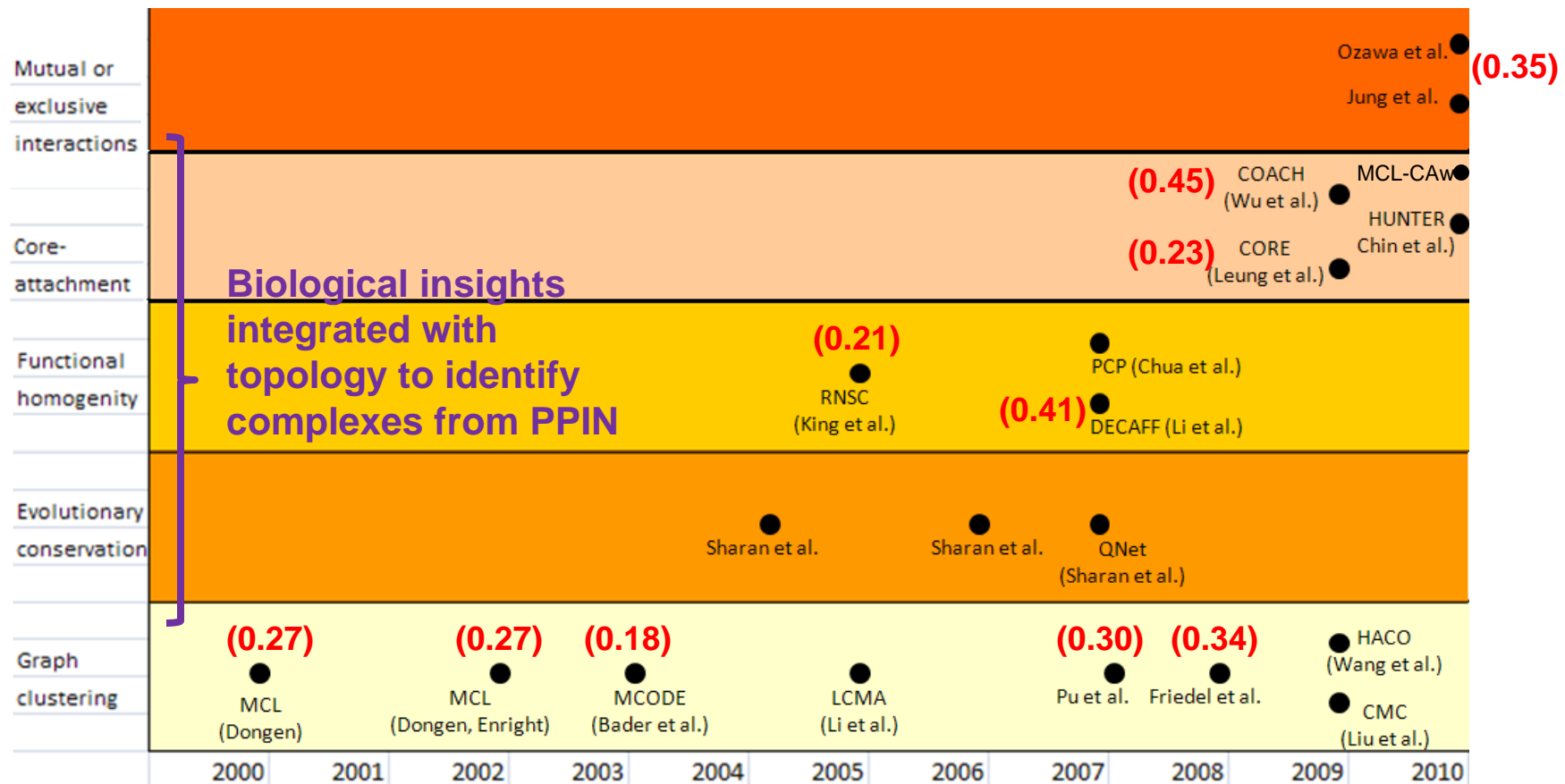


Image Source: Barabási and Oltvai, 2004

- **Noise in PPI data**
 - Spuriously-detected interactions (false positives), and missing interactions (false negatives)
- **Transient interactions**
 - Many proteins that actually interact are not from the same complex, they bind temporarily to perform a function
- **Not all proteins in the same complex may actually interact with each other**

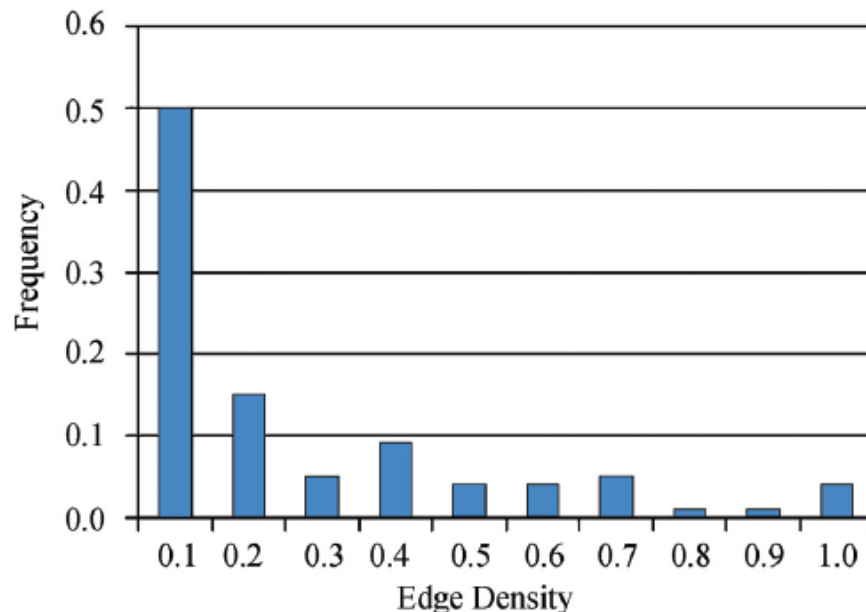
Protein Complex Prediction Methods



- Adding biological info improves F1

Challenges

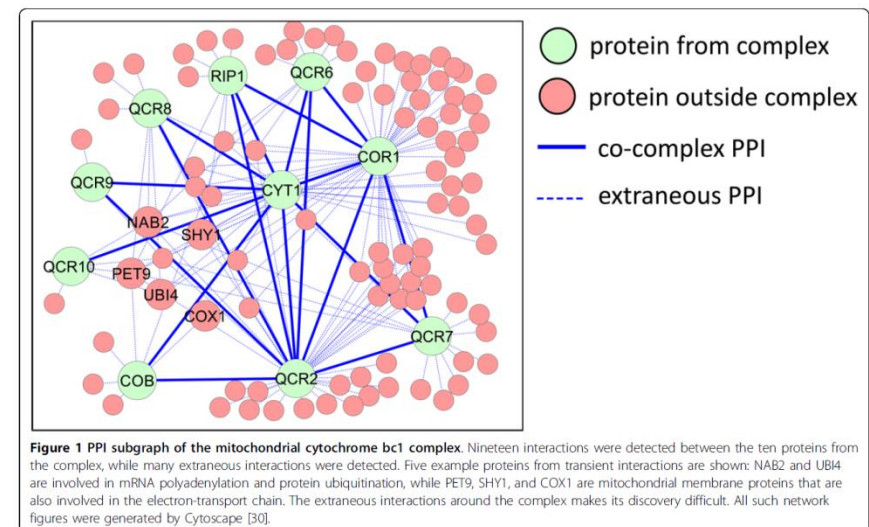
- Recall & precision of protein complex prediction algo's have lots to be improved



- Does a “cleaner” PPI network help?
- How to capture “high edge density” complexes that overlap each other?
- How to capture “low edge density” complexes?
- How to capture small complexes?

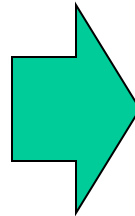
Cytochrome BC1 Complex

- Involved in electron-transport chain in mitochondrial inner membrane
- Discovery of this complex from PPI data is difficult
 - Sparseness of the complex's PPI subnetwork
 - Only 19 out of 45 possible interactions were detected between the complex's proteins
 - Many extraneous interactions detected with other proteins outside the complex
 - E.g., UBI4 is involved in protein ubiquitination, and binds to many proteins to perform its function.



- **Key idea to deal with sparseness**

Augment physical PPI network with other forms of linkage that suggest two proteins are likely to integrate



Supervised Weighting of Composite Networks (SWC)

- **Data integration**
- **Supervised edge weighting**
- **Clustering**

Overview of SWC

1. Integrate diff data sources to form composite network
 2. Weight each edge based on probability that its two proteins are co-complex, using a naïve Bayes model w/ supervised learning
 3. Perform clustering on the weighted network
- **Advantages**
 - Data integration increases density of complexes
 - **co-complex proteins are likely to be related in other ways even if they do not interact**
 - Supervised learning
 - **Allows discrimination betw co-complex and transient interactions**
 - Naïve Bayes' transparency
 - **Model parameters can be analyzed, e.g., to visualize the contribution of diff evidences in a predicted complex**

1. Integrate multiple data sources

- **Composite network: Vertices represent proteins, edges represent relationships between proteins**
- **There is an edge betw proteins u , v , if and only if u and v are related according to any of the data sources**

Data source	Database	Scoring method
PPI	BioGRID, IntACT, MINT	Iterative AdjustCD.
L2-PPI (indirect PPI)	BioGRID, IntACT, MINT	Iterative AdjustCD
Functional association	STRING	STRING
Literature co-occurrence	PubMed	Jaccard coefficient

	Yeast			Human		
	# Pairs	% co-complex	coverage	# Pairs	% co-complex	coverage
PPI	106328	5.8%	55%	48098	10%	14%
L2-PPI	181175	1.1%	18%	131705	5.5%	20%
STRING	175712	5.7%	89%	311435	3.1%	27%
PubMed	161213	4.9%	70%	91751	4.3%	11%
All	531800	2.1%	98%	522668	3.4%	49%

2. Supervised edge-weighting

- Treat each edge as an instance, where features are data sources and feature values are data source scores, and class label is “co-complex” or “non-co-complex”

PPI	L2 PPI	STRING	Pubmed	Class
0	0.56	451	0	“co-complex”
0.1	0	25	0	“non-co-complex”

- Supervised learning:

- Discretize each feature (Minimum Description Length discretization⁷)
- Learn maximum-likelihood parameters for the two classes:

$$P(F = f|co - comp) = \frac{n_{c,F=f}}{n_c} \quad P(F = f|non - co - comp) = \frac{n_{\neg c,F=f}}{n_{\neg c}}$$

for each discretized feature value f of each feature F

- Weight each edge e with its posterior probability of being co-complex:

$$\begin{aligned}
 &weight(e) \\
 &= P(co - comp|F_1 = f_1, F_2 = f_2, \dots) \\
 &= \frac{P(F_1 = f_1, F_2 = f_2, \dots | co - comp)P(co - comp)}{Z} \\
 &= \frac{\prod_i P(F_i = f_i | co - comp)P(co - comp)}{Z} \\
 &= \frac{\prod_i P(F_i = f_i | co - comp)P(co - comp)}{\prod_i P(F_i = f_i | co - comp)P(co - comp) + \prod_i P(F_i = f_i | non - co - comp)P(non - co - comp)}
 \end{aligned}$$

3. Complex Discovery

- **Weighted composite network used as input to clustering algorithms**
 - CMC, ClusterONE, IPCA, MCL, RNSC, HACO
 - **Predicted complexes scored by weighted density**
-
- **The clustering algo's generate clusters with low overlap**
 - Only 15% of clusters are generated by two or more algo's
- ⇒ **Voting-based aggregative strategy, COMBINED:**
- Take union of clusters generated by the diff algo's
 - Similar clusters from multiple algo's are given higher scores
 - **If two or more clusters are similar (Jaccard ≥ 0.75), then use the highest scoring one and multiply its score by the # of algo's that generated it**

Experiments

- **Weighting approaches:**
 - SWC vs BOOST, TOPO, STR, NOWEI
- **Evaluate performance on the 6 clustering algos and the COMBINED clustering strategy**
- **Real complexes for training and testing: CYC200814 for yeast, CORUM15 for human**
- **Evaluation**
 - How well co-complex edges are predicted
 - How well predicted complexes match real complexes

Evaluation wrt Co-Complex Prediction

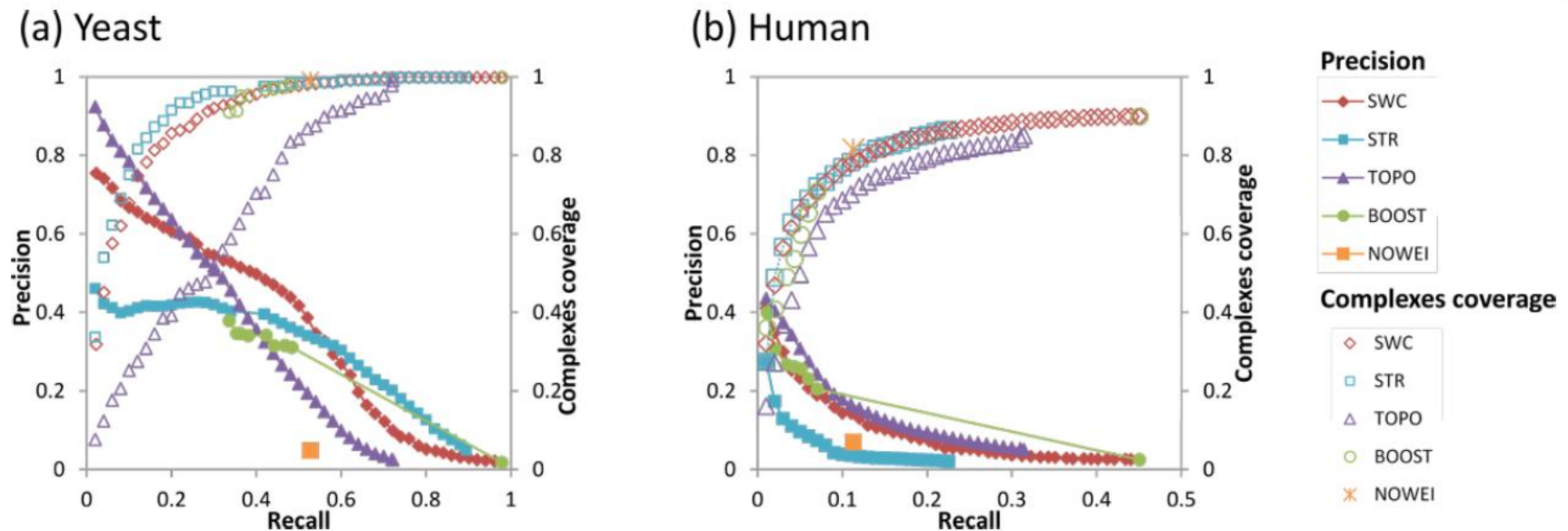
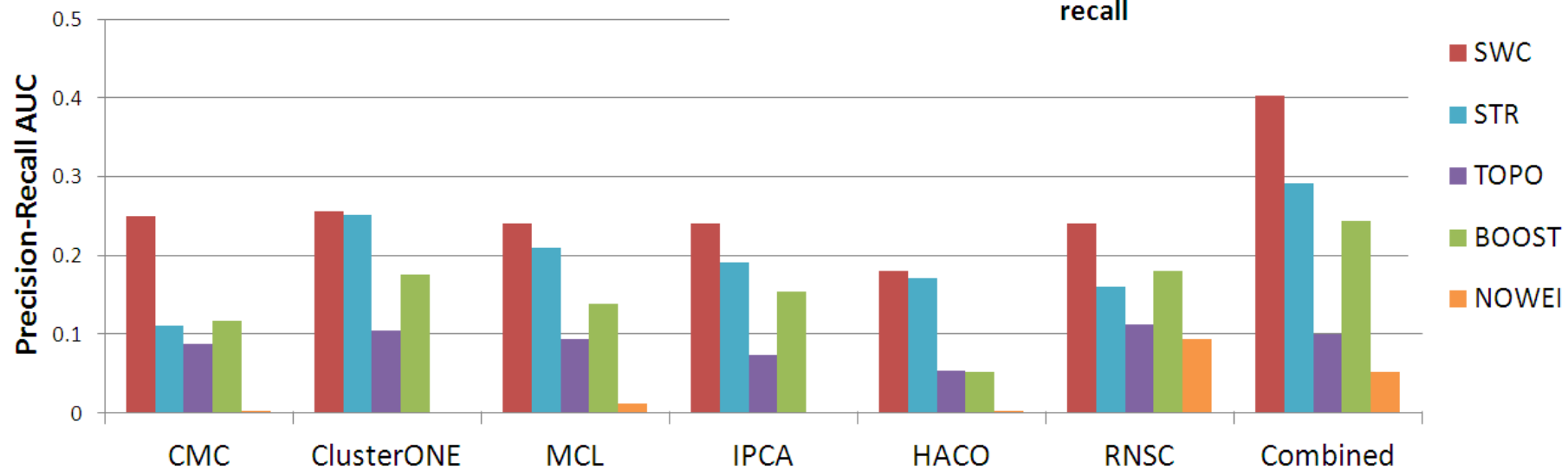
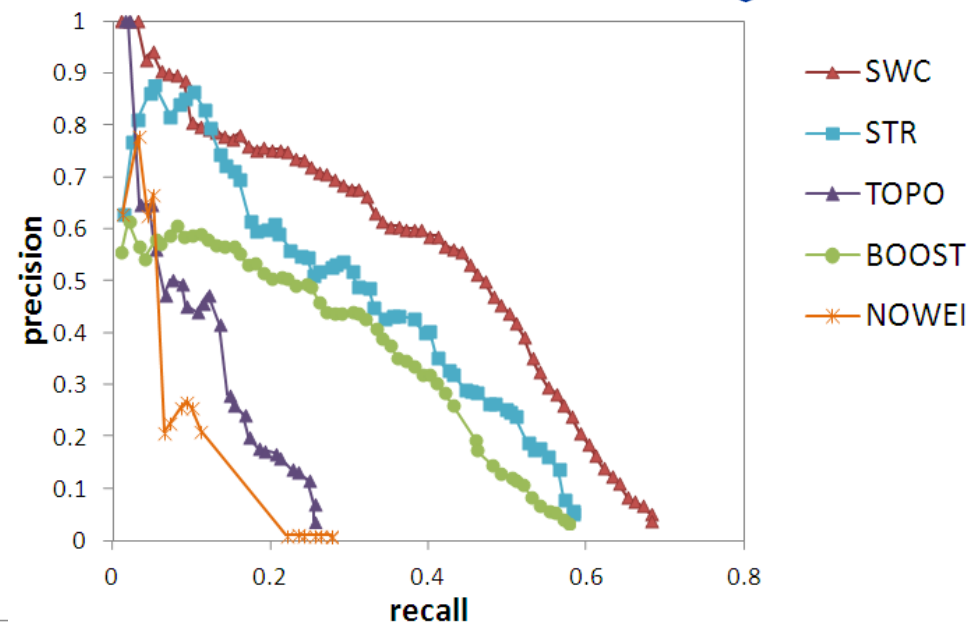
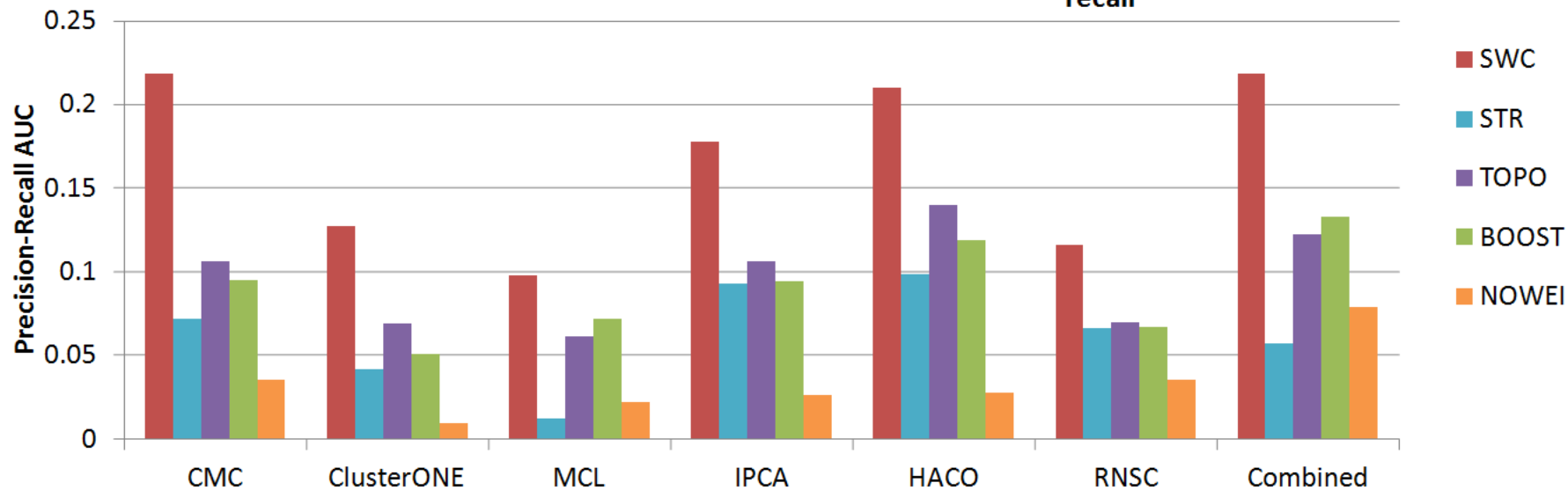
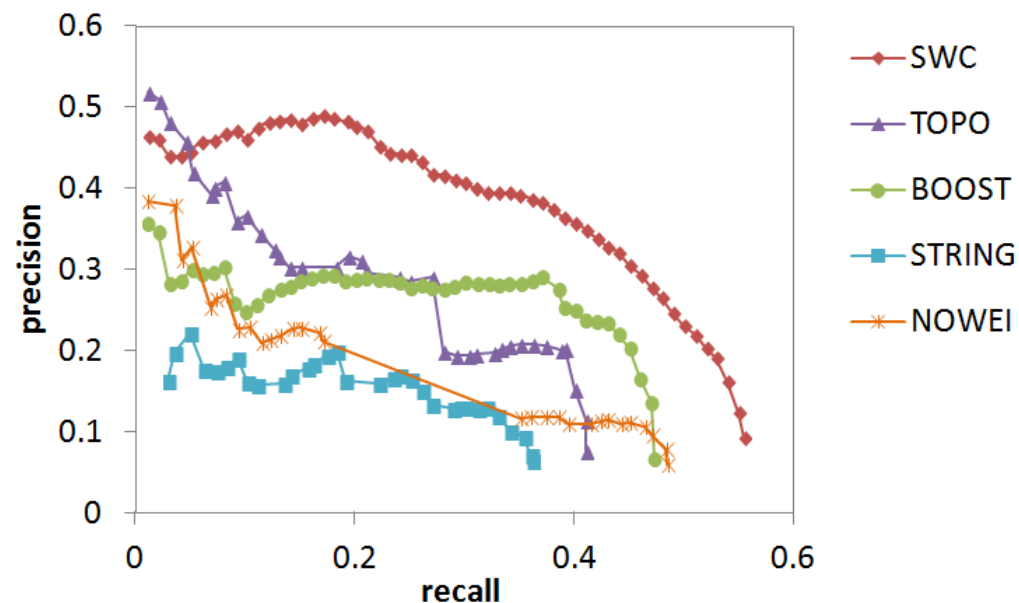


Figure 2 Precision-recall graph for classification of co-complex edges using the five weighting schemes. (a) Classification of yeast co-complex edges. SWC and BOOST achieve the highest recall through data integration. TOPO has high precision for its top-scoring edges, but these are clustered in a few complexes. SWC achieves higher precision than STR, except when too many edges are considered. BOOST classifies edges categorically, giving high scores to one set of edges with about 50% recall and 35% precision, and low scores to the remainder. (b) Classification of human co-complex edges. Recall and precision for human is much lower than for yeast. TOPO has higher precision than SWC, but its predicted edges are clustered in fewer complexes. BOOST classifies edges categorically, and its high-scoring edges achieve 7% recall, with comparable precision with SWC. NOWEI has slightly higher precision than STR, which has the lowest precision.

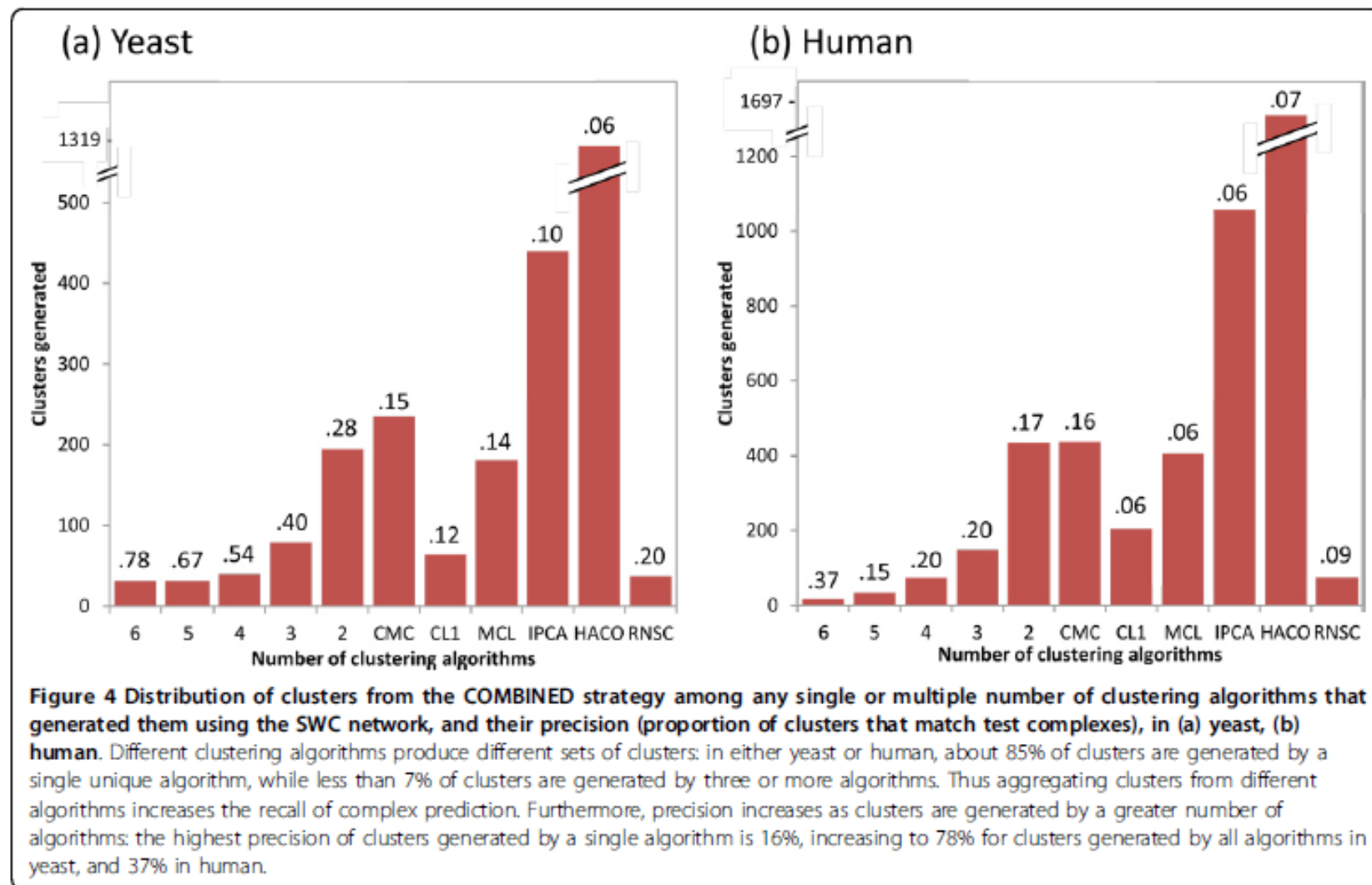
Evaluation wrt Yeast Complex Prediction



Evaluation wrt Human Complex Prediction



Why the “COMBINED” Strategy?



Power of the “COMBINED” Strategy

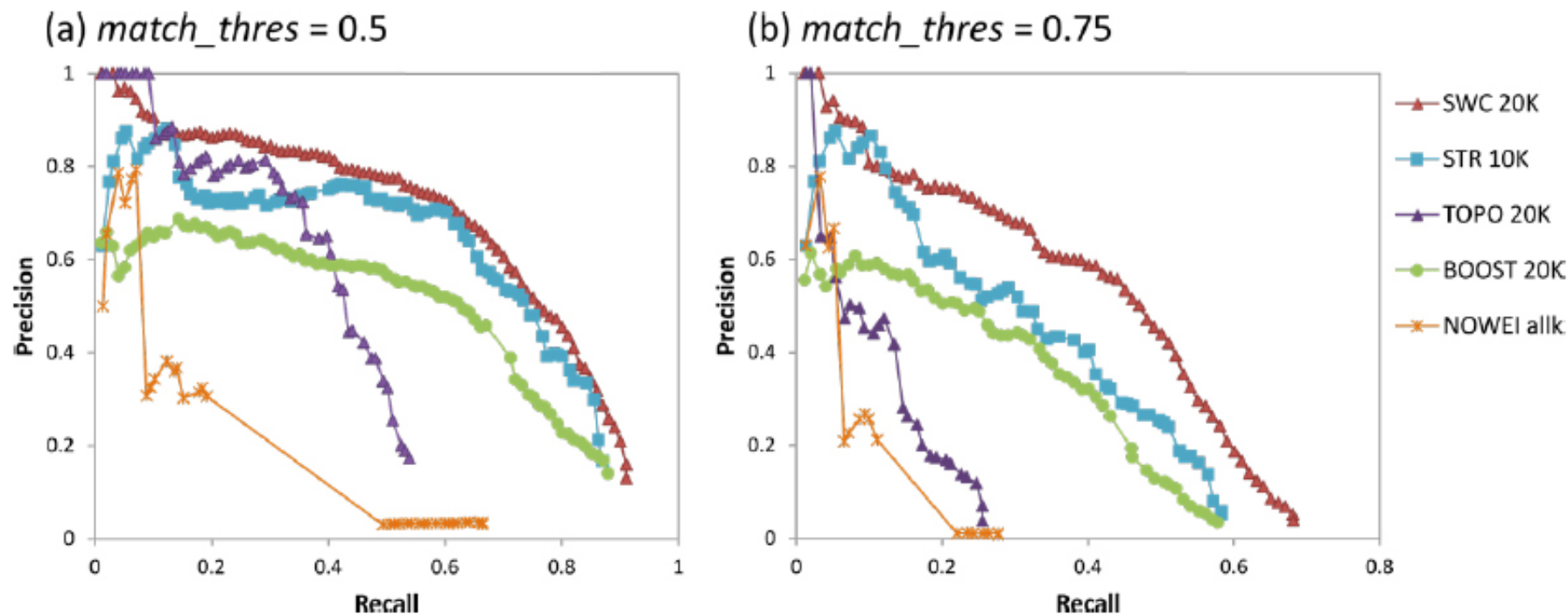


Figure 5 Precision-recall graphs for yeast complex prediction using the five weighting approaches with the COMBINED clustering strategy, using $k = 20000$ for SWC, TOPO, and BOOST, $k = 10000$ for STR, and $k = \text{all edges}$ for NOWEI. (a) *match_thres* = 0.5, (b) *match_thres* = 0.75. SWC achieves the highest recall, with the highest precision at almost all recall levels, especially with the stricter *match_thres* = 0.75. Thus it outperforms all other weighting approaches, especially at predicting complexes with fine granularity.

Power of the “COMBINED” Strategy

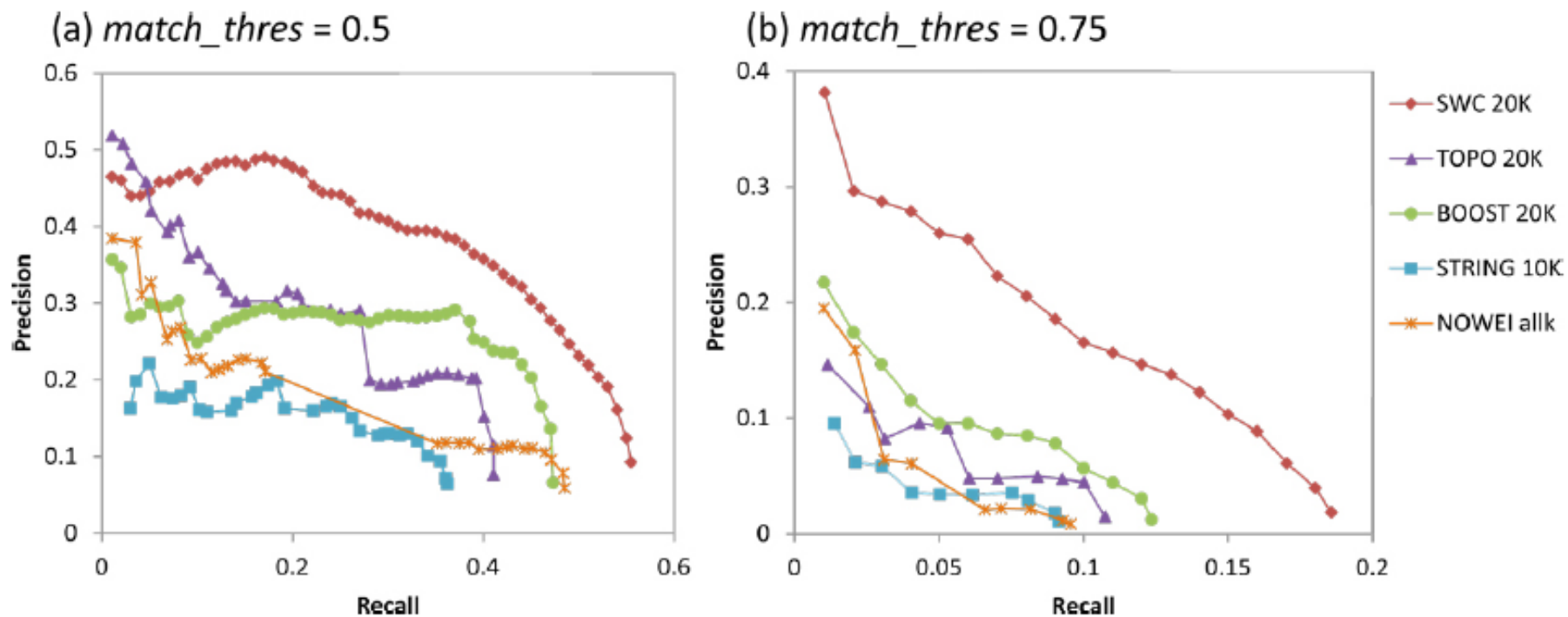
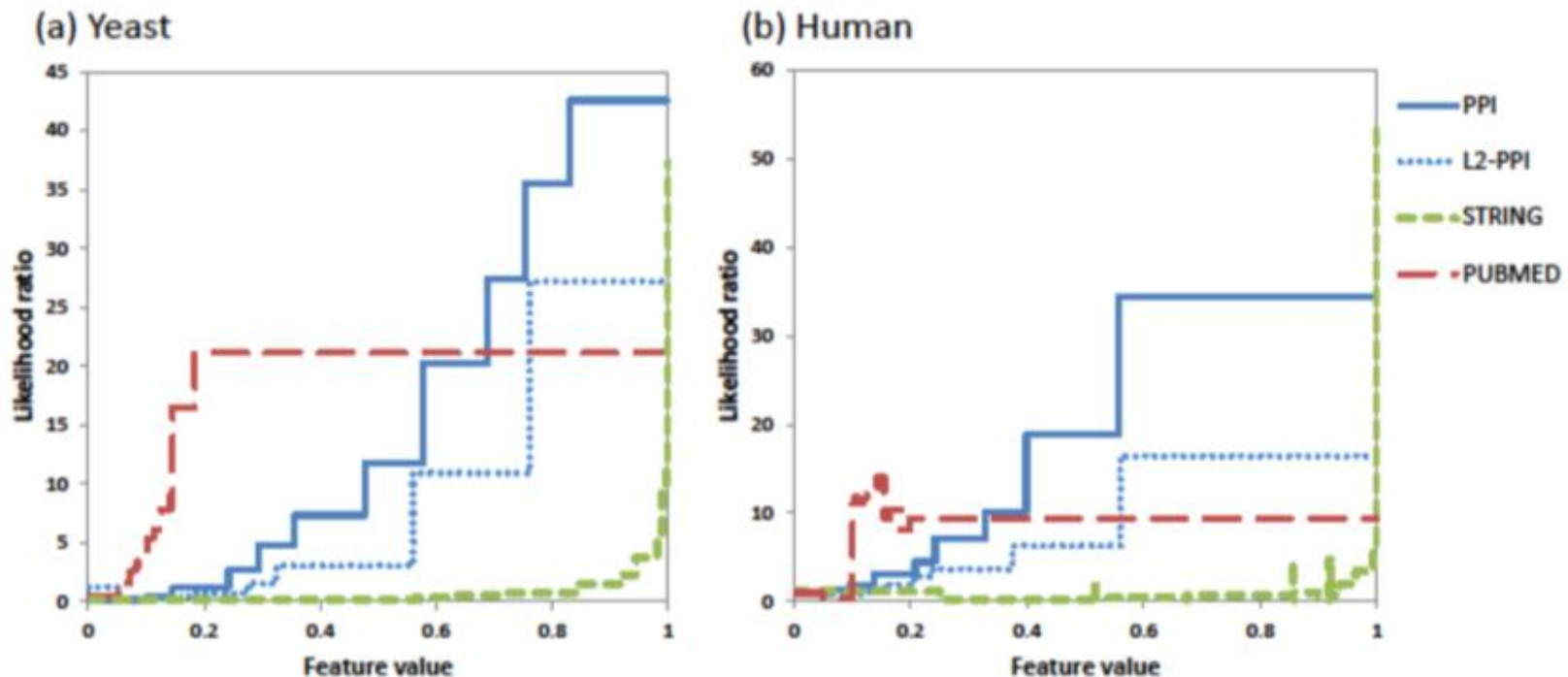


Figure 7 Precision-recall graphs for human complex prediction using the five weighting approaches for the COMBINED clustering strategy. SWC achieves the highest recall with the highest precision at almost all recall levels, especially with the stricter *match_thres* = 0.75, where SWC recalls at least 50% more test complexes compared to the other approaches and maintains almost twice the precision throughout its recall range. Thus it outperforms all other weighting approaches, especially at predicting complexes with ne granularity.

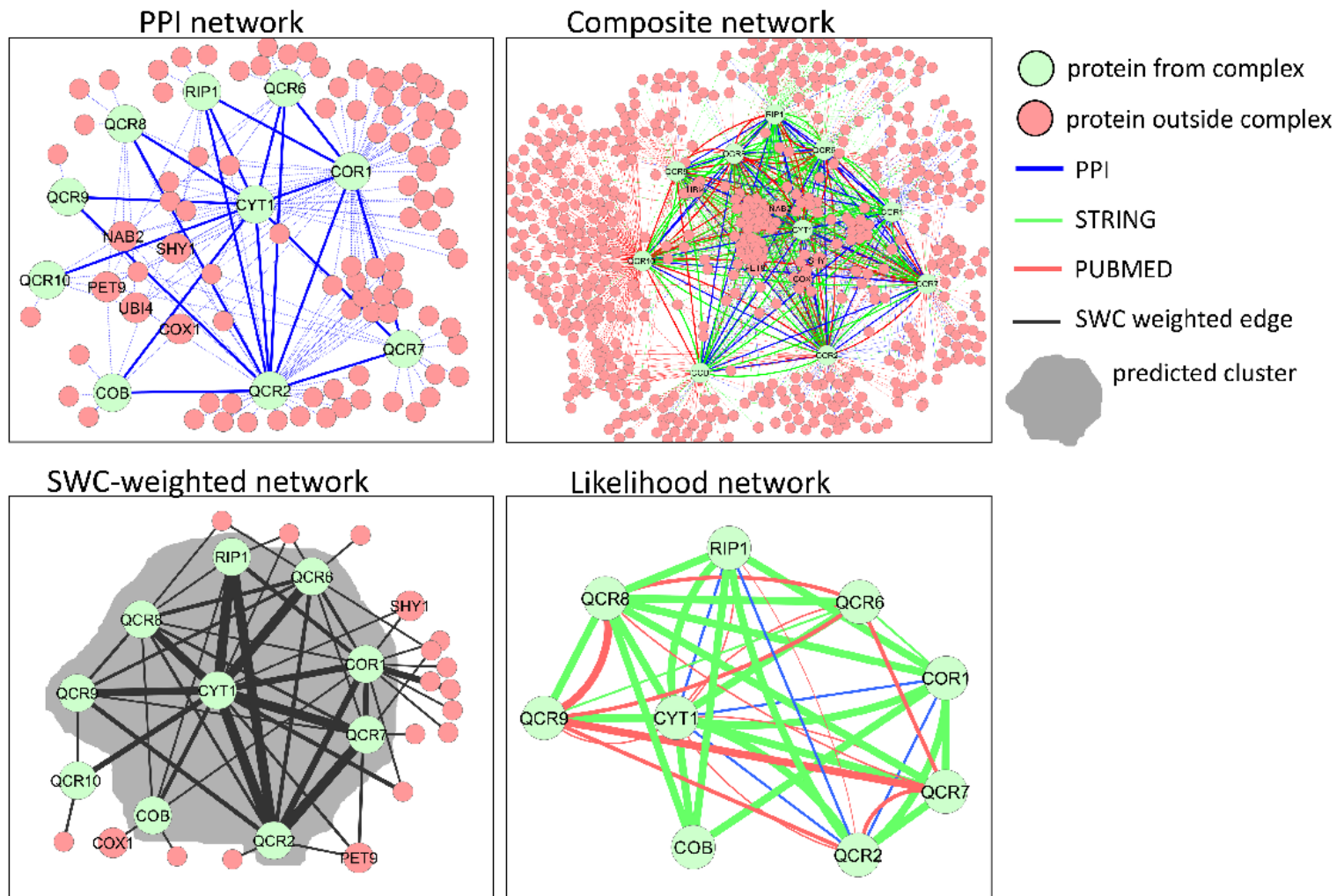
Evidence: Co-complexness Likelihoods

- “Co-complexness strength” of a feature F with score f can be expressed as:

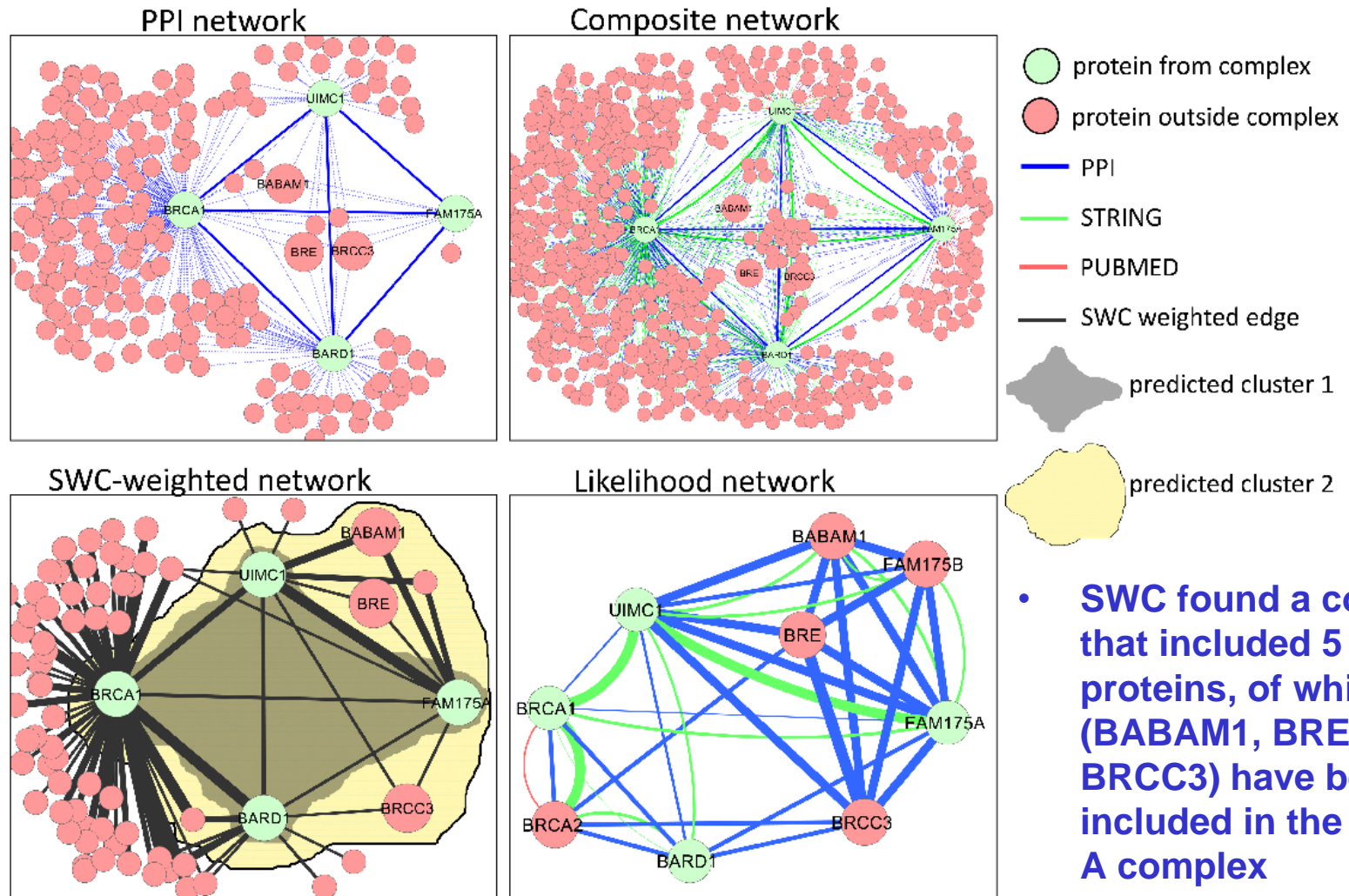
$$\text{likelihood ratio} = \frac{P(F = f | \text{co-complex})}{P(F = f | \text{non-co-complex})}$$



Example: Yeast BC1 Complex

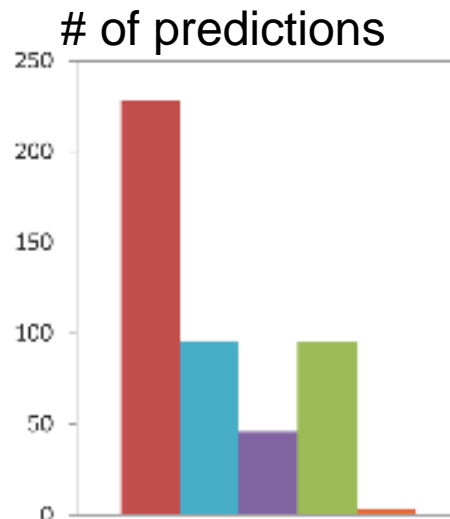


Example: Human BRCA1-A complex

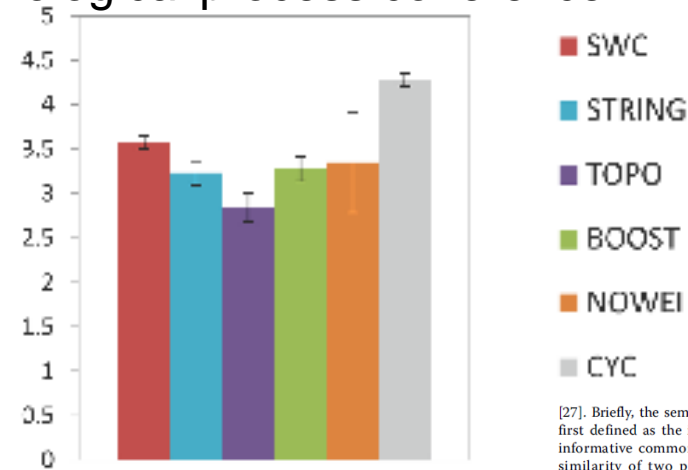


High-Confidence Predicted Complexes

Yeast

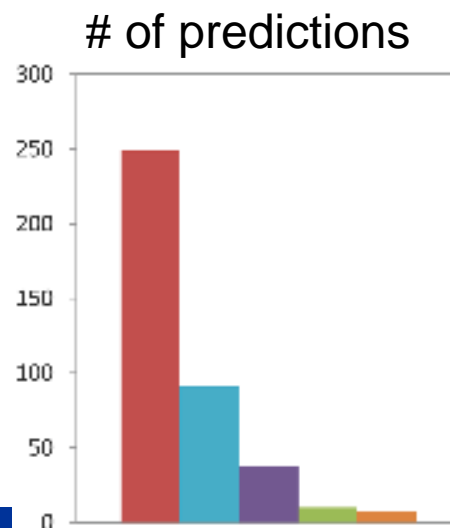


Biological process coherence

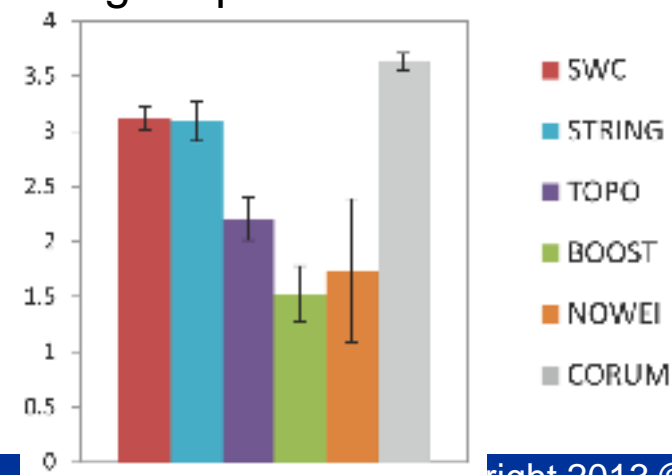


[27]. Briefly, the semantic similarity of two GO terms is first defined as the information content of their most informative common ancestor. Next, the BP semantic similarity of two proteins is defined as the highest semantic similarity between their two sets of annotated BP terms. Then, we define the BP semantic coherence of a predicted complex as the average BP semantic similarity between every pair of proteins in that complex (likewise for CC and MF).

Human

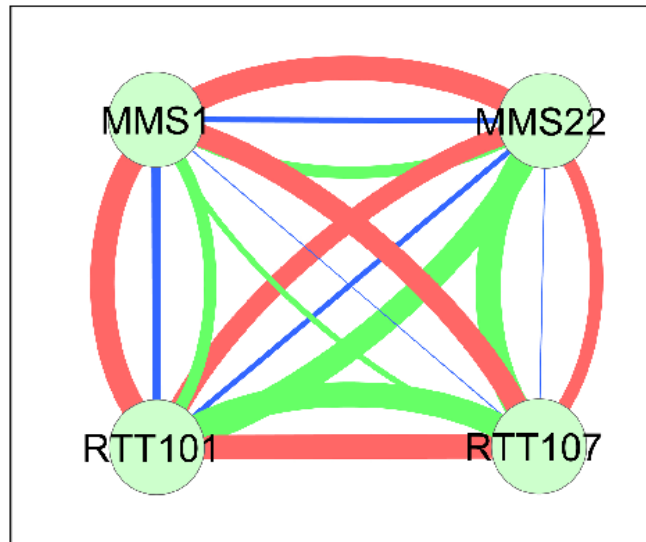


Biological process coherence

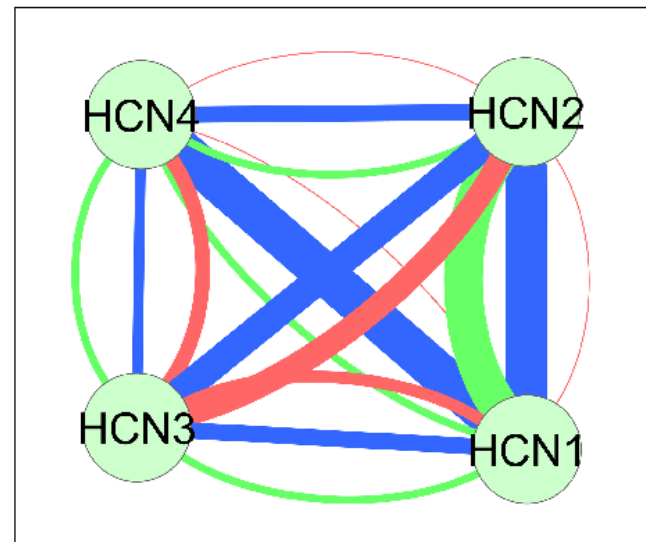


Two Novel Predicted Complexes

(a) Yeast



(b) Human



— PPI
— STRING
— PUBMED

- **Novel yeast complex:** Annotated w/ DNA metabolic process and response to stress, forms a complex called Cul8-RING which is absent in our ref set
- **Novel human complex:** Annotated w/ transport process, Uniprot suggests it may be a subunit of a potassium channel complex

Novel complexes predicted

Yeast

Biological process	# complexes
Protein metabolic process	49
RNA metabolic process	36
DNA metabolic process	15
Small molecule metabolic process	23
Regulation of metabolic process	11
Regulation of gene expression	8
Organelle organization	40
Transport	43
Response to stress	20
Response to chemical stimulus	7
Cell cycle process	11

Human

Biological process	# complexes
Protein metabolic process	32
RNA metabolic process	29
DNA metabolic process	4
Small molecule metabolic process	19
Regulation of metabolic process	74
Regulation of gene expression	34
Organelle organization	19
Transport	38
Response to stress	28
Response to chemical stimulus	32
Cell cycle process	14

Conclusions

- **Naïve-Bayes data-integration to predict co-complexed proteins**
 - Use of multiple data sources increases density of complexes
 - Supervised learning allows discrimination betw co-complex and transient interactions
- **Tested approach using 6 clustering algo's**
 - Clusters produced by diff algo's have low overlap, combining them gives greater recall
 - Clusters produced by more algo's are more reliable

Acknowledgement

- Yong et al. Supervised maximum-likelihood weighting of composite protein networks for complex prediction. *BMC Systems Biology*, 6(Suppl 2):S13, 2012



Yong Chern Han