

Talk given at University of Warsaw, January 2010

Topology of PPI Networks: Applications and Questions

Limsoon Wong
(Works with Hon Nian Chua & Guimei Liu)



2

Plan



- PPI network cleansing based on PPI topology
- PPI-based protein complex prediction
- PPI-based protein function prediction

PPI Network Cleansing Based on PPI Topology



4



Why Protein Interactions?

- Complete genomes are now available
- Knowing the **genes** is not enough to understand how biology **functions**
- **Proteins**, not genes, are responsible for many cellular activities
- Proteins function by **interacting** w/ other proteins and biomolecules

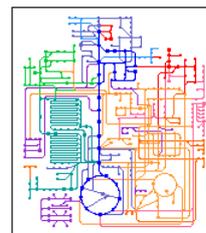
GENOME



PROTEOME



"INTERACTOME"



Slide credit: See-Kiong Ng

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

High-Tech Expt PPI Detection Methods

- Yeast two-hybrid assays
- Mass spec of purified complexes (e.g., TAP)
- Correlated mRNA expression
- Genetic interactions (e.g., synthetic lethality)
- ...

FACT: Generating *large amounts of* experimental data about protein-protein interactions can be done with ease.

Slide credit: See-Kiong Ng

Key Bottleneck

- Many high-throughput expt detection methods for protein-protein interactions have been devised
- But ...

High-throughput approach sacrifice quality for **quantity**:
 (a) limited or biased coverage: **false negatives**, &
 (b) high error rates : **false positives**

Slide credit: See-Kiong Ng

Noise in PPI Networks

Experimental method category ^a	Number of interacting pairs	Co-localization ^b (%)	Co-cellular-role ^b (%)
All: All methods	9347	64	49
A: Small scale Y2H	1861	73	62
A0: GY2H Uetz <i>et al.</i> (published results)	956	66	45
A1: GY2H Uetz <i>et al.</i> (unpublished results)	516	53	33
A2: GY2H Ito <i>et al.</i> (core)	798	64	40
A3: GY2H Ito <i>et al.</i> (all)	3655	41	15
B: Physical methods	71	98	95
C: Genetic methods	1052	77	75
D1: Biochemical, <i>in vitro</i>	614	87	79
D2: Biochemical, chromatography	648	93	88
E1: Immunological, direct	1025	90	90
E2: Immunological, indirect	34	100	93
2M: Two different methods	2360	87	85
3M: Three different methods	1212	92	94
4M: Four different methods	570	95	93

Sprinzak *et al.*, *JMB*, 327:919-923, 2003

Large disagreement betw methods

- High level of noise
- ⇒ Need to clean up before making inference on PPI networks

Measures that correlate with function homogeneity and localization coherence

- Two proteins participating in same biological process are more likely to interact
- Two proteins in the same cellular compartments are more likely to interact



- CD-distance
- FS-Weight

CD-distance & FS-Weight: Based on concept that two proteins with many interaction partners in common are likely to be in same biological process & localize to the same compartment

Czekanowski-Dice Distance (Brun et al, 2003)

- **Given a pair of proteins (u, v) in a PPI network**
 - N_u = the set of neighbors of u
 - N_v = the set of neighbors of v
- $CD(u,v) = \frac{2 | N_u \cap N_v |}{| N_u | + | N_v |}$
- **Consider relative intersection size of the two neighbor sets, not absolute intersection size**
 - Case 1: $|N_u| = 1, |N_v| = 1, |N_u \cap N_v| = 1, CD(u,v) = 1$
 - Case 2: $|N_u| = 10, |N_v| = 10, |N_u \cap N_v| = 10, CD(u,v) = 1$

Iterated CD-Distance (Liu et al, 2008)

- **Variant of CD-distance that penalizes proteins with few neighbors**

$$wL(u,v) = \frac{2 | N_u \cap N_v |}{| N_u | + \lambda_u + | N_v | + \lambda_v}$$

$$\lambda_u = \max\left\{0, \frac{\sum_{x \in G} |N_x|}{|V|} - |N_u|\right\}, \lambda_v = \max\left\{0, \frac{\sum_{x \in G} |N_x|}{|V|} - |N_v|\right\}$$

- **Suppose average degree is 4, then**
 - Case 1: $|N_u| = 1, |N_v| = 1, |N_u \cap N_v| = 1, wL(u,v) = 0.25$
 - Case 2: $|N_u| = 10, |N_v| = 10, |N_u \cap N_v| = 10, wL(u,v) = 1$

A thought...

$$wL(u,v) = \frac{2 |N_u \cap N_v|}{|N_u| + \lambda_u + |N_v| + \lambda_v}$$

- **Weight of interaction reflects its reliability**

⇒ **Can we get better results if we use this weight to recalculate the score of other interactions?**

Iterated CD-Distance (Liu et al, 2006)

- $wL^0(u,v) = 1$ if $(u,v) \in G$, otherwise $wL^0(u,v) = 0$

- $wL^1(u,v) = \frac{|N_u \cap N_v| + |N_u \cap N_v|}{|N_u| + \lambda_u + |N_v| + \lambda_v}$

- $wL^k(u,v) = \frac{\sum_{x \in N_u \cap N_v} wL^{k-1}(u,x) + \sum_{x \in N_u \cap N_v} wL^{k-1}(v,x)}{\sum_{x \in N_u} wL^{k-1}(u,x) + \lambda_u + \sum_{x \in N_v} wL^{k-1}(v,x) + \lambda_v}$

- $\lambda_u^k = \max\{0, \frac{\sum_{x \in V} \sum_{y \in N_x} wL^{k-1}(x,y)}{|V|} - \sum_{x \in N_u} wL^{k-1}(u,x)\}$

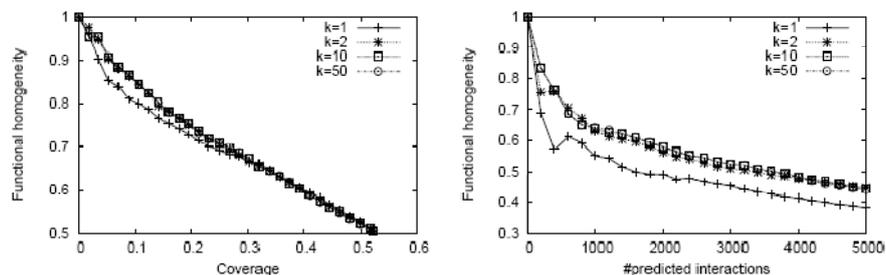
- $\lambda_v^k = \max\{0, \frac{\sum_{x \in V} \sum_{y \in N_x} wL^{k-1}(x,y)}{|V|} - \sum_{x \in N_v} wL^{k-1}(v,x)\}$

Validation

- **DIP yeast dataset**
 - Functional homogeneity is 32.6% for PPIs where both proteins have functional annotations and 3.4% over all possible PPIs
 - Localization coherence is 54.7% for PPIs where both proteins have localization annotations and 4.9% over all possible PPIs
- **Let's see how much better iterated CD-distance is over the baseline above, as well as over the original CD-distance/FS-weight**

How many iteration is enough?

Cf. ave functional homogeneity of protein pairs in DIP < 4%
ave functional homogeneity of PPI in DIP < 33%



- **Iterated CD-distance achieves best performance wrt functional homogeneity at k=2**
- **Ditto wrt localization coherence (not shown)**

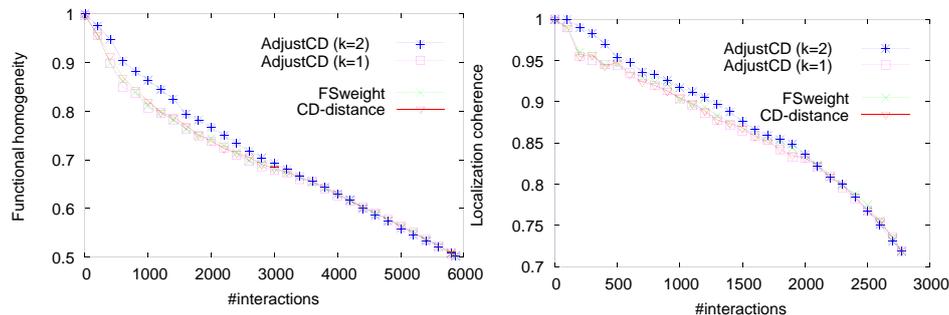
How many iteration is enough?

noise level	k	#common PPIs	avg_rank_diff	avg_score_diff
100%	1	5669	540.21	0.10
	2	5870	144.86	0.02
	20	5849	67.00	0.01
300%	1	5322	881.77	0.18
	2	5664	367.45	0.06
	20	5007	249.85	0.02
500%	1	5081	1013.14	0.23
	2	5502	625.46	0.12
	20	5008	317.33	0.05
1000%	k=1	4472	1187.10	0.28
	k=2	5101	1021.69	0.27
	k=20	5264	614.66	0.13

- Iterative CD-distance at diff k values on noisy network
⇒ # of iterations depends on amt of noise

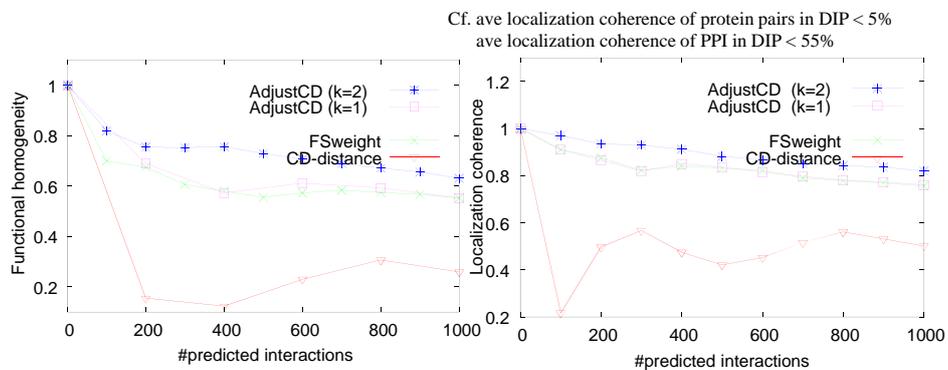
Identifying False Positive PPIs

Cf. ave localization coherence of protein pairs in DIP < 5%
ave localization coherence of PPI in DIP < 55%



- Iterated CD-distance is an improvement over previous measures for assessing PPI reliability

Identifying False Negative PPIs



- Iterated CD-distance is an improvement over previous measures for predicting new PPIs

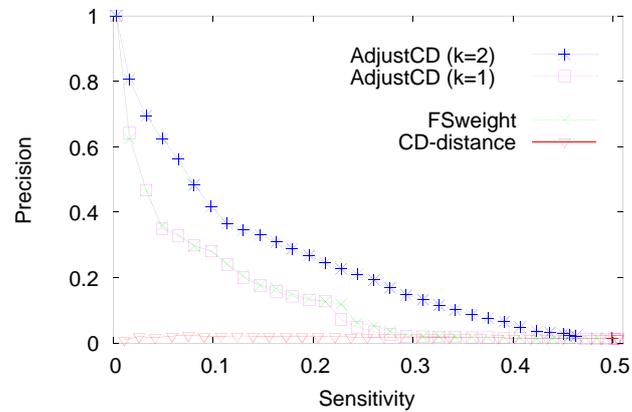
Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

5-Fold Cross-Validation

- **DIP core dataset**
 - Ave # of proteins in 5 groups: 986
 - Ave # of interactions in 5 training datasets: 16723
 - Ave # of interactions in 5 testing datasets: 486591
 - Ave # of correct answer interactions: 307
- **Measures:**
 - sensitivity = $TP / (TP + FN)$
 - specificity = $TN / (TN + FP)$
 - #negatives \gg #positives, specificity is always high
 - >97.8% for all scoring methods
 - precision = $TP / (TP + FP)$

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

5-Fold X-Validation



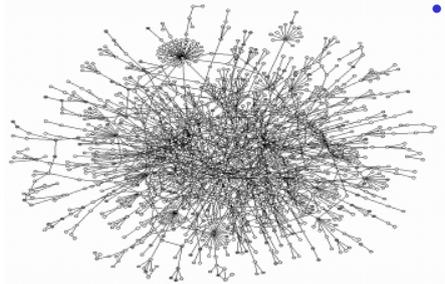
- Iterated CD-distance is an improvement over previous measures for identifying false positive & false negative PPIs

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

PPI-Based Protein Complex Prediction

Motivation

- **Nature of high-throughput PPI expts**
 - Proteins are taken out of their natural context!
- **Can a protein interact with so many proteins simultaneously?**
- **A big “hub” and its “spokes” should probably be decomposed into subclusters**
 - Each subcluster is a set proteins that interact in the same space and time
 - Viz., a protein complex



Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

PPI-Based Complex Prediction Algo

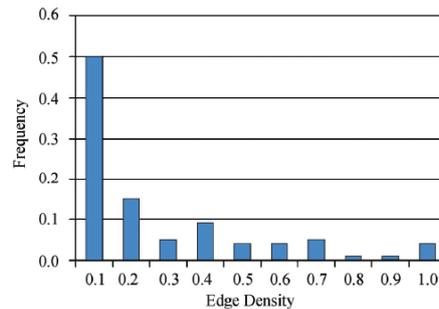
	RNSC	MCODE	MCL
Type	Clustering, local search cost based	Local neighborhood density search	Flow simulation
Multiple assignment of protein	No	Yes	No
Weighted edge	No	No	Yes

- **Issue: recall vs precision has to be improved**
 - ⇒ Does a “cleaner” PPI network help?
 - ⇒ How to capture “low edge density” complexes?

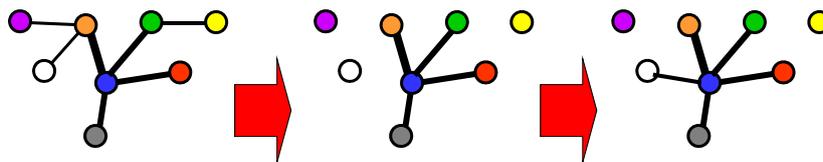
Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

PPI-Based Complex Prediction

- Recall & precision of protein complex prediction also have lots to be improved
- Does a “cleaner” PPI network help?
- How to capture “low edge density” complexes?
 - ⇒ Clique merging?
 - ⇒ Relative density?
 - ⇒ Core-n-attachment?



Cleaning PPI Network



- **Modify existing PPI network as follow**
 - Remove interactions with low weight
 - Add interactions with high weight
- **Then run RNSC, MCODE, MCL, ..., as well as our own method CMC**

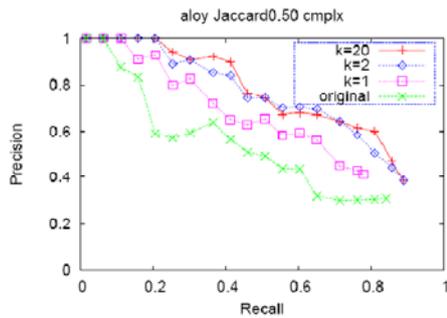
CMC: Clustering of Maximal Cliques

- **Remove noise edges in input PPI network by discarding edges having low iterated CD-distance**
- **Augment input PPI network by addition of missing edges having high iterated CD-distance**
- **Predict protein complex by finding overlapping maximal cliques, and merging/removing them**
- **Score predicted complexes using cluster density weighted by iterated CD-distance**

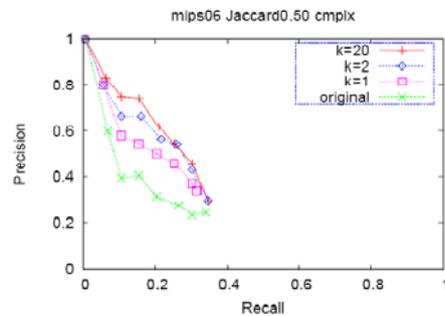
Validation Experiments

- **Matching a predicted complex S with a true complex C**
 - Vs: set of proteins in S
 - Vc: set of proteins in C
 - $\text{Overlap}(S, C) = |V_s \cap V_c| / |V_s \cup V_c|$
 - $\text{Overlap}(S, C) \geq 0.5$
- **Evaluation**
 - Precision = matched predictions / total predictions
 - Recall = matched complexes / total complexes
- **Datasets: combined info from 6 yeast PPI expts**
 - #interactions: 20461 PPI from 4671 proteins
 - #interactions with >0 common neighbor: 11487

Effecting of Cleaning on CMC



(a) Aloy, *match_thres*=0.50

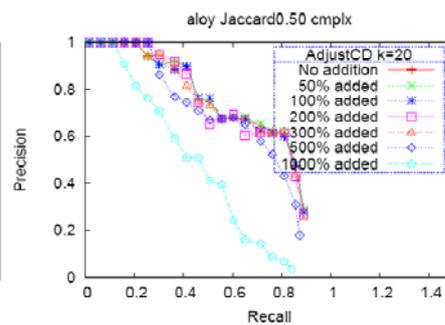
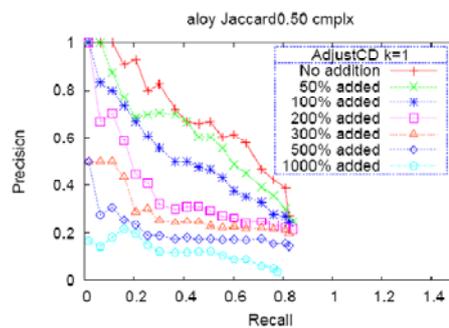


(b) MIPS, *match_thres*=0.50

- Cleaning by Iterated CD-distance improves recall & precision of CMC

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

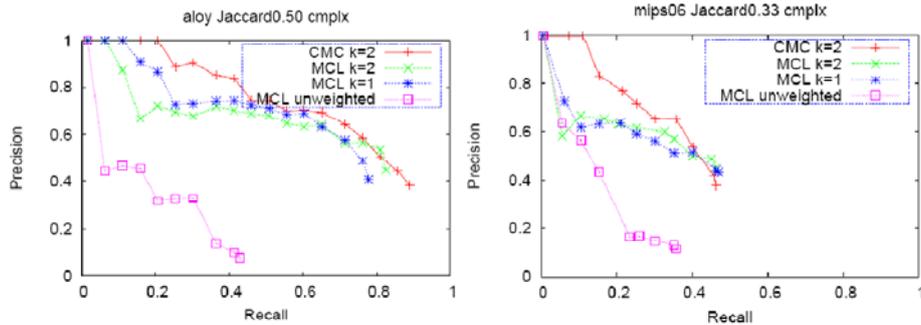
Noise Tolerance of CMC



- If cleaning is done by iterating CD-distance 20 times, CMC can tolerate up to 500% noise in the PPI network!

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Effect of Cleansing on MCL



- MCL benefits significantly from cleaning too
- Ditto for other protein complex prediction methods

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

CMC vs Others

scoring method: AdjustCD					match_thres=0.50							
clustering methods	k	#clusters	avg size	loc. score	Aloy (#complexes: 63)				MIPS (#complexes: 162)			
					#matched clusters	precision	#matched complexes	recall	#matched clusters	prec	#matched complexes	recall
CMC	0	172	9.83	0.823	53	0.308	53	0.841	42	0.244	55	0.340
	1	121	9.42	0.897	50	0.413	49	0.778	41	0.339	51	0.315
	2	148	8.50	0.899	57	0.385	56*	0.889	44	0.297	56*	0.346
	20	146	8.78	0.891	56	0.384	56*	0.889	43	0.295	56*	0.346
CFinder	0	103	13.84	0.528	39	0.379	38	0.603	34	0.330	40	0.247
	1	76	12.86	0.724	38	0.500	38	0.603	30	0.395	34	0.210
	2	95	11.66	0.713	44	0.463	43	0.683	36	0.379	46	0.284
	20	95	11.77	0.718	44	0.463	43	0.683	37	0.389	49	0.302
MCL	0	372	9.40	0.638	27	0.073	27	0.429	30	0.081	37	0.228
	1	120	10.18	0.848	49	0.408	49	0.778	40	0.333	51	0.315
	2	116	10.31	0.856	52	0.448	52	0.825	41	0.353	51	0.315
	20	110	10.75	0.849	49	0.445	49	0.778	37	0.336	47	0.290
MCode	0	61	7.31	0.849	20	0.328	20	0.317	18	0.295	22	0.136
	1	103	7.42	0.913	35	0.340	35	0.556	30	0.291	39	0.241
	2	88	8.67	0.897	34	0.386	34	0.540	29	0.330	39	0.241
	20	82	10.28	0.838	29	0.354	29	0.460	23	0.280	32	0.198

Table 3. The impact of the iterative scoring method on the performance of four clustering methods. For CMC, MCL and CFinder, we retain only the top-6000 interactions, and no new interactions are added. For MCode, we retain all the interactions with non-zero score and add top-3000 new interactions with the highest score. The 2nd column is the number of iterations k of the iterative scoring method, and $k=0$ means the PPI network is unweighted. The 3rd column is the number of clusters generated, the 4th and 5th column is the average size and co-localization score of generated clusters.

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Characteristics of Unmatched Clusters

- At $k = 2$...
- 85 clusters predicted by CMC do not match complexes in Aloy and MIPS
- Localization coherence score $\sim 90\%$
- 65/85 have the same informative GO term annotated to $> 50\%$ of proteins in the cluster

⇒ Likely to be real complexes

PPI-Based Protein Function Prediction

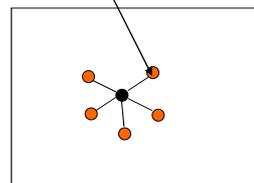
Protein Interaction Based Approaches

- **Neighbour counting** (Schwikowski et al, 2000)
 - Rank function based on freq in interaction partners
 - **Chi-square** (Hishigaki et al, 2001)
 - Chi square statistics using expected freq of functions in interaction partners
 - **Markov Random Fields** (Deng et al, 2003; Letovsky et al, 2003)
 - Belief propagation exploit unannotated proteins for prediction
 - **Simulated Annealing** (Vazquez et al, 2003)
 - Global optimization by simulated annealing
 - Exploit unannotated proteins for prediction
 - **Clustering** (Brun et al, 2003; Samanta et al, 2003)
 - Functional distance derived from shared interaction partners
 - Clusters based on functional distance represent proteins with similar functions
 - **Functional Flow** (Nabieva et al, 2004)
 - Assign reliability to various expt sources
 - Function “flows” to neighbour based on reliability of interaction and “potential”
- **Indirect Functional Assoc** (Chua et al, 2006)
 - Identification of reliable common interaction partners

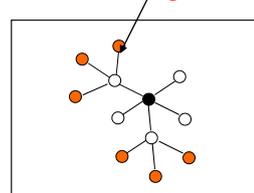
Functional Association Thru Interactions

- **Direct functional association:**
 - Interaction partners of a protein are likely to share functions w/ it
 - Proteins from the same pathways are likely to interact
- **Indirect functional association**
 - Proteins that share interaction partners with a protein may also likely to share functions w/ it
 - Proteins that have common biochemical, physical properties and/or subcellular localization are likely to bind to the same proteins

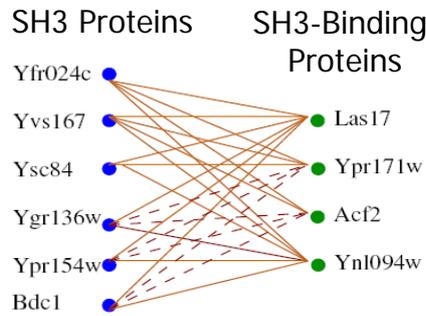
Level-1 neighbour



Level-2 neighbour

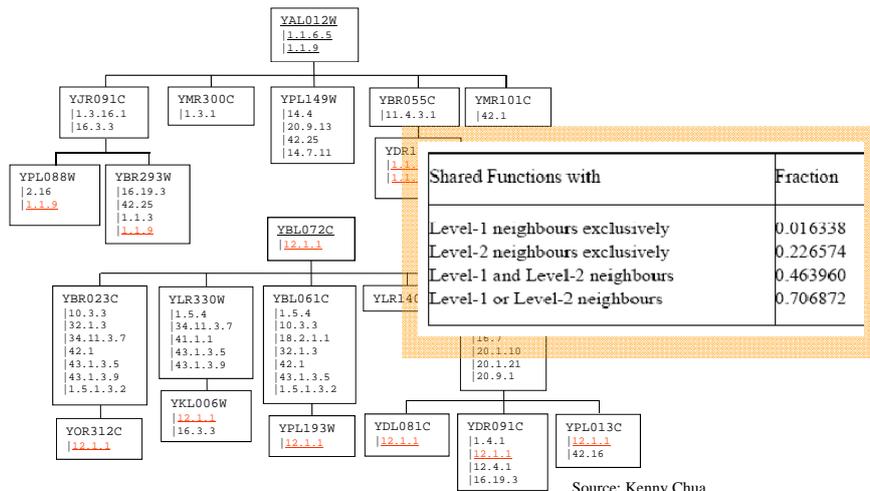


An Illustrative Case of Indirect Functional Association?



- Is indirect functional association plausible?
- Is it found often in real interaction data?
- Can it be used to improve protein function prediction from protein interaction data?

Freq of Indirect Functional Association



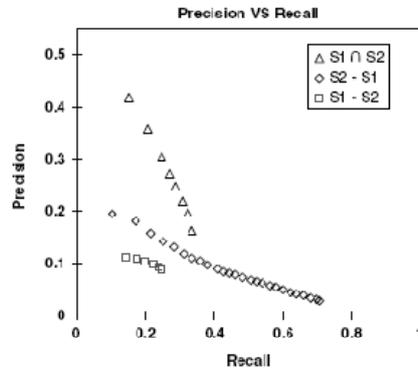
Source: Kenny Chua

Prediction Power By Majority Voting

- Remove overlaps in level-1 and level-2 neighbours to study predictive power of “level-1 only” and “level-2 only” neighbours
- Sensitivity vs Precision analysis

$$PR = \frac{\sum_i^K k_i}{\sum_i^K m_i} \quad SN = \frac{\sum_i^K k_i}{\sum_i^K n_i}$$

- n_i is no. of fn of protein i
- m_i is no. of fn predicted for protein i
- k_i is no. of fn predicted correctly for protein i



- ⇒ “level-2 only” neighbours performs better
- ⇒ L1 ∩ L2 neighbours has greatest prediction power

Functional Similarity Estimate: Czekanowski-Dice Distance

- Functional distance between two proteins (Brun et al, 2003)

$$D(u, v) = \frac{|N_u \Delta N_v|}{|N_u \cup N_v| + |N_u \cap N_v|}$$

- N_k is the set of interacting partners of k
- $X \Delta Y$ is symmetric diff betw two sets X and Y
- Greater weight given to similarity

⇒ Similarity can be defined as

$$S(u, v) = 1 - D(u, v) = \frac{2X}{2X + (Y + Z)}$$

Is this a good measure if u and v have very diff number of neighbours?

Functional Similarity Estimate: FS-Weighted Measure

- FS-weighted measure

$$S(u, v) = \frac{2|N_u \cap N_v|}{|N_u - N_v| + 2|N_u \cap N_v|} \times \frac{2|N_u \cap N_v|}{|N_v - N_u| + 2|N_u \cap N_v|}$$

- N_k is the set of interacting partners of k
- Greater weight given to similarity

⇒ Rewriting this as

$$S(u, v) = \frac{2X}{2X + Y} \times \frac{2X}{2X + Z}$$

Correlation w/ Functional Similarity

- Correlation betw functional similarity & estimates

Neighbours	CD-Distance	FS-Weight
S_1	0.471810	0.498745
S_2	0.224705	0.298843
$S_1 \cup S_2$	0.224581	0.29629

- Equiv measure slightly better in correlation w/ similarity for L1 & L2 neighbours

Reliability of Expt Sources

- **Diff Expt Sources have diff reliabilities**

- Assign reliability to an interaction based on its expt sources (Nabieva et al, 2004)

- **Reliability betw u and v computed by:**

$$r_{u,v} = 1 - \prod_{i \in E_{u,v}} (1 - r_i)$$

- r_i is reliability of expt source i ,
- $E_{u,v}$ is the set of expt sources in which interaction betw u and v is observed

Source	Reliability
Affinity Chromatography	0.823077
Affinity Precipitation	0.455904
Biochemical Assay	0.666667
Dosage Lethality	0.5
Purified Complex	0.891473
Reconstituted Complex	0.5
Synthetic Lethality	0.37386
Synthetic Rescue	1
Two Hybrid	0.265407

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Functional Similarity Estimate: FS-Weighted Measure with Reliability

- **Take reliability into consideration when computing FS-weighted measure:**

$$S_R(u, v) = \frac{2 \sum_{w \in (N_u \cap N_v)} r_{u,w} r_{v,w}}{\left(\sum_{w \in N_u - N_v} r_{u,w} + \sum_{w \in (N_u \cap N_v)} r_{u,w} (1 - r_{v,w}) \right) + 2 \sum_{w \in (N_u \cap N_v)} r_{u,w} r_{v,w}} \times \frac{2 \sum_{w \in (N_u \cap N_v)} r_{u,w} r_{v,w}}{\left(\sum_{w \in N_v - N_u} r_{v,w} + \sum_{w \in (N_u \cap N_v)} r_{v,w} (1 - r_{u,w}) \right) + 2 \sum_{w \in (N_u \cap N_v)} r_{u,w} r_{v,w}}$$

- N_k is the set of interacting partners of k
- $r_{u,w}$ is reliability weight of interaction betw u and v

⇒ **Rewriting**

$$S(u, v) = \frac{2X}{2X + Y} \times \frac{2X}{2X + Z}$$

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

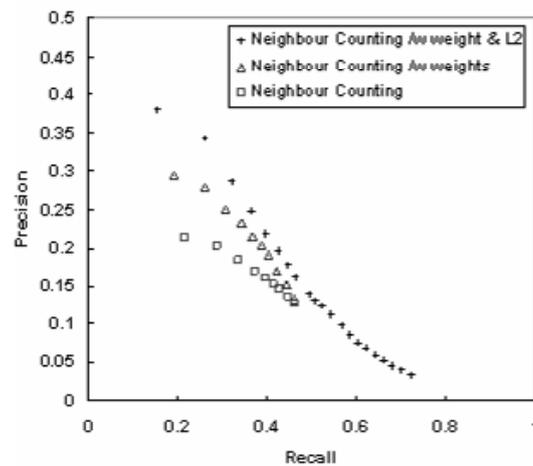
Integrating Reliability

- Equiv measure shows improved correlation w/ functional similarity when reliability of interactions is considered:

Neighbours	CD-Distance	FS-Weight	FS-Weight R
S_1	0.471810	0.498745	0.532596
S_2	0.224705	0.298843	0.375317
$S_1 \cup S_2$	0.224581	0.29629	0.363025

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

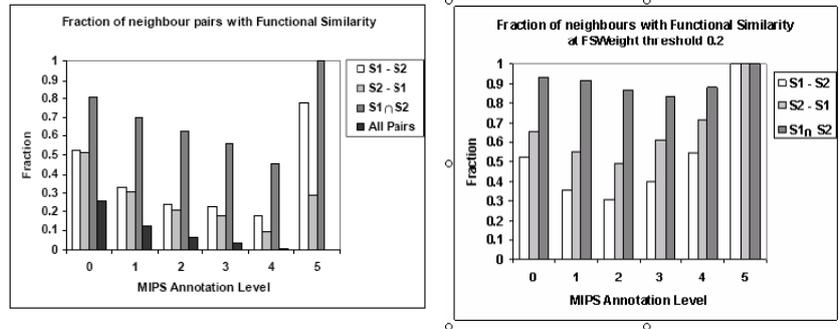
Improvement to Prediction Power by Majority Voting



Considering only neighbours w/ FS weight > 0.2

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Improvement to Over-Rep of Functions in Neighbours



Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Use L1 & L2 Neighbours for Prediction



• FS-weighted Average

$$f_x(u) = \frac{1}{Z} \left[\lambda r_{int} \pi_x + \sum_{v \in N_u} \left(S_{TR}(u, v) \delta(v, x) + \sum_{w \in N_v} S_{TR}(u, w) \delta(w, x) \right) \right]$$

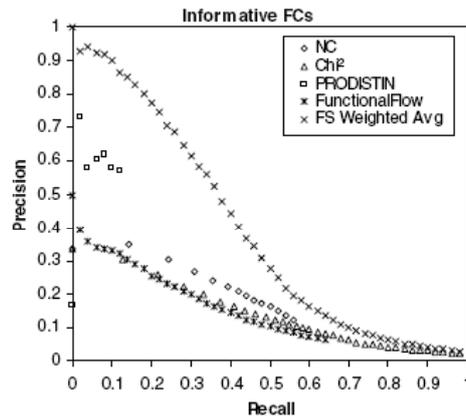
- r_{int} is fraction of all interaction pairs sharing function
- λ is weight of contribution of background freq
- $\delta(k, x) = 1$ if k has function x , 0 otherwise
- N_k is the set of interacting partners of k
- π_x is freq of function x in the dataset
- Z is sum of all weights

$$Z = 1 + \sum_{v \in N_u} \left(S_{TR}(u, v) + \sum_{w \in N_v} S_{TR}(u, w) \right)$$

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Performance of FS-Weighted Averaging

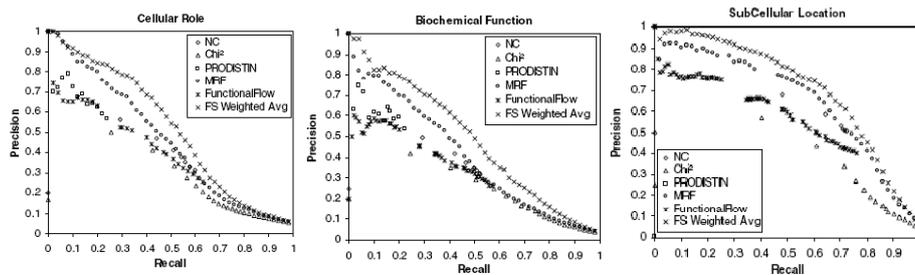
- LOOCV comparison with Neighbour Counting, Chi-Square, PRODISTIN



Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Performance of FS-Weighted Averaging

- Dataset from Deng et al, 2003
 - Gene Ontology (GO) Annotations
 - MIPS interaction dataset
- Comparison w/ Neighbour Counting, Chi-Square, PRODISTIN, Markov Random Field, FunctionalFlow

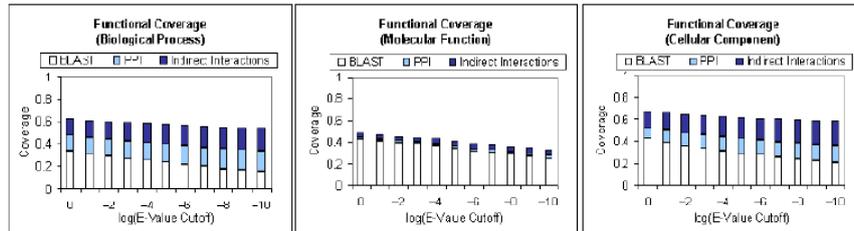


Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong



Freq of Indirect Functional Association in Other Genomes

D. melanogaster

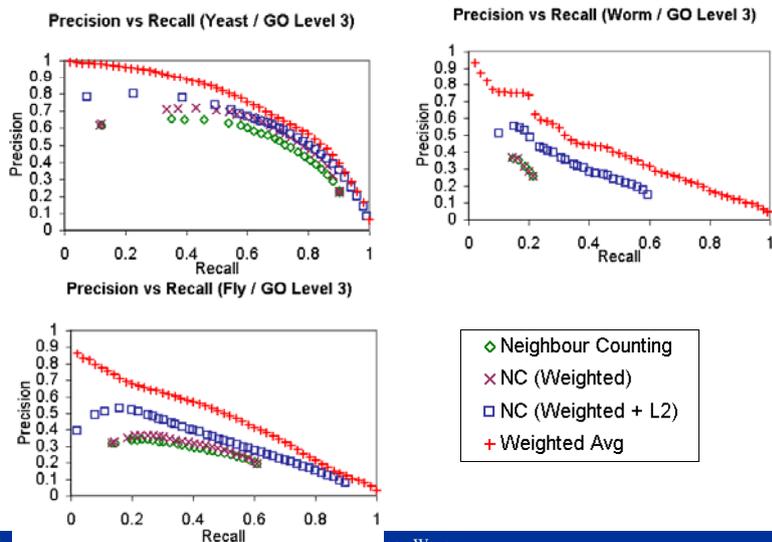


Genome	Annotation	S_1-S_2	S_2-S_1	$S_1 \cap S_2$	$S_1 \cup S_2$
<i>S. cerevisiae</i>	MIPS	0.007193	0.226574	0.463960	0.706872
<i>D. melanogaster</i>	GO	0.008801	0.168622	0.138138	0.315561
<i>C. elegans</i>	GO	0.007193	0.051237	0.061080	0.119510

Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong



Effectiveness of FS Weighted Averaging in Other Genomes



Talk at Warsaw University, January 2010. Copyright © 2010 by Limsoon Wong

Last Remarks



52

What have we learned?



- **Guilt by association of common interaction partners is useful for**
 - Cleansing high-throughput PPI network data
 - Predicting protein complexes
 - Inferring protein functional information
- **Acknowledgement**
 - Kenny Chua, Guimei Liu

Readings

- H. N. Chua, et al. "Exploiting indirect neighbours and topological weight to predict protein function from protein-protein interactions", *Bioinformatics*, 22:1623-1630, 2006
- H.N. Chua, et al. "Using Indirect Protein-Protein Interactions for Protein Complex Prediction", *JBCB*, 6(3):435--466, 2008
- H. N. Chua, L. Wong. "Increasing the Reliability of Protein Interactomes", *Drug Discovery Today*, 13(15/16):652--658, 2008
- G. Liu, et al. "Assessing and predicting protein interactions using both local and global network topological metrics", *Proc GIW2008*
- G. Liu, et al. "Complex Discovery from Weighted PPI Networks", *Bioinformatics*, 25(15):1891--1897, 2009.

Any Question?