

# A Guest Lecture for CS6280: Guilt by Association

**Limsoon Wong**  
**1 March 2006**



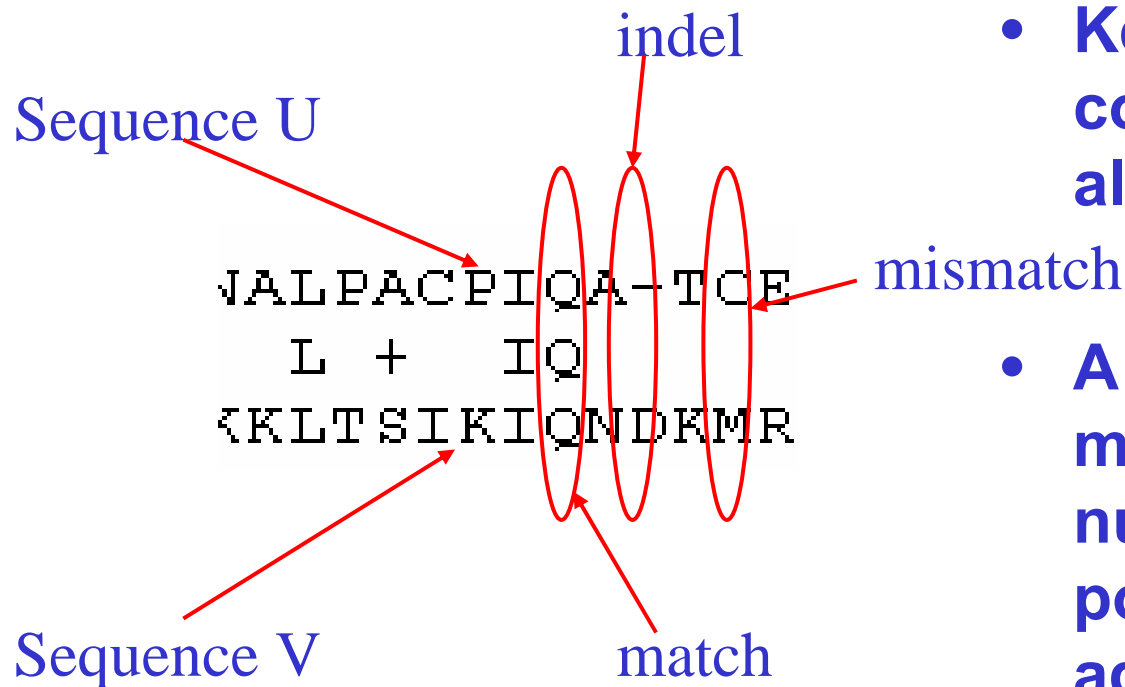
# Plan

- **Recap of sequence alignment**
- **Guilt by association**
- **What if no homology of known function is found?**
  - **Guilt by different types of association!**
    - **Genome phylogenetic profiling**
    - **Protfun**
    - **SVM-Pairwise**
    - **Protein-protein interactions**

# Very Brief Recap of Sequence Comparison/Alignment



# Sequence Alignment



- **Key aspect of seq comparison is seq alignment**
- **A seq alignment maximizes the number of positions that are in agreement in two sequences**

# Sequence Alignment: Poor Example

- **Poor seq alignment shows few matched positions**  
 ⇒ **The two proteins are not likely to be homologous**

**Alignment by FASTA of the sequences of amicyanin and domain 1 of ascorbate oxidase**

	60	70	80	90	100
Amicyanin	MPHNVH FVAGVLGEAALKGPM MKKEQAYSLTFTEAGTYDYHCTPHPFMRGKV VVE				
			...	.	...
Ascorbate Oxidase	ILQRGTPWADGTASISQCAINPGETFFYNFTVDNPGTFFYHGH LGMQRSAGLYGSLI				
	70	80	90	100	110 120

No obvious match between  
 Amicyanin and Ascorbate Oxidase

# Sequence Alignment: Good Example

- Good alignment usually has clusters of extensive matched positions
- ⇒ The two proteins are likely to be homologous

```
□ >gil13476732|ref|NP\_108301.1| unknown protein [Mesorhizobium loti]  
gil14027493|dbj|BAB53762.1| unknown protein [Mesorhizobium loti]  
Length = 105
```

```
Score = 105 bits (262), Expect = 1e-22  
Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)
```

```
Query: 1   MKPGRLASIALAIIFLPMVPAHAATIEITMENLVISPTEVSAKVGDTIRWVNKDVF AHT 60  
          MK G L  ++      MA PA AATIE+T++ LV SP  V AKVGDTI WVN DV AHT  
Sbjct: 1   MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNNDVVAHT 60
```

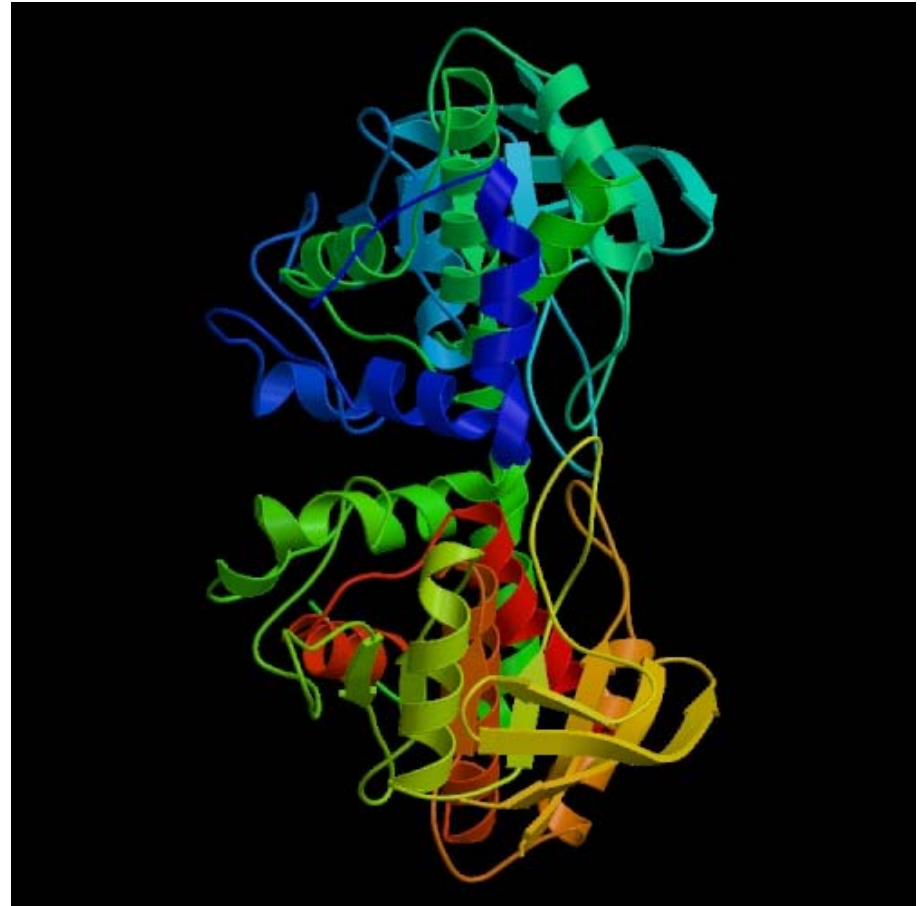
good match between  
Amicyanin and unknown M. loti protein

# Guilt-by-Association



## A protein is a ...

- A protein is a large complex molecule made up of one or more chains of amino acids
- Protein performs a wide variety of activities in the cell





# Function Assignment to Protein Sequence

SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR  
YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE  
QNTATIVMTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD  
VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG  
TFVVIDAMLDMMSERKVDVYGFVSRIRAQRCQMVQTD MQYVFIYQALLEHYLYGDTELE  
VT

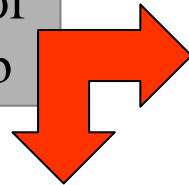
- How do we attempt to assign a function to a new protein sequence?

# Guilt-by-Association

- Compare the target sequence  $T$  with sequences  $S_1, \dots, S_n$  of known function in a database
- Determine which ones amongst  $S_1, \dots, S_n$  are the mostly likely homologs of  $T$
- Then assign to  $T$  the same function as these homologs
- Finally, confirm with suitable wet experiments

# Guilt-by-Association

Compare  $T$  with seqs of known function in a db



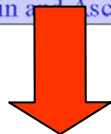
### Poor Sequence Alignment

- Poor seq alignment shows few matched positions  
⇒ The two proteins are not likely to be homologous

**Alignment by FASTA of the sequences of amicyanin and domain 1 of ascorbate oxidase**

	60	70	80	90	100
Amicyanin	MPHNHFVAGVLGEAALKGPMKKEQAYSLTFTEAGTYDYHCTPHPFMRGKVVV				
		...	...	...	...
Ascorbate Oxidase	ILQRGTPWADGTASISQCAINPGETFFYNEFTVDNPGTFFYHGHLMQRSAGLYG				
	70	80	90	100	110

No obvious match between Amicyanin and Ascorbate Oxidase



Discard this function as a candidate

### Good Sequence Alignment

- Good alignment usually has clusters of extensive matched positions  
⇒ The two proteins are likely to be homologous

```

>gi113476732|ref|NP_108301.1| unknown protein [Mesorhizobium loti]
gi114027493|dbj|BAB53762.1| unknown protein [Mesorhizobium loti]
Length = 105

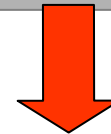
Score = 105 bits (262), Expect = 1e-22
Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)

Query: 1 MKPQRLASIALAIIFLPMVPAHAATIEITMENLVISPTESAKVGDITRWVNDVFAHT 60
      MK G L ++ MA PA AATIE+T++ LV SP V AKVGDIT WVN DV AHT
Sbjct: 1 MKAGALIRLSVLALALMAFAAAATIEVTIDKLVSFATVEAKVGDITIEWVNDVFAHT 60
  
```

good match between  
Amicyanin and unknown M. loti protein



Assign to  $T$  same function as homologs

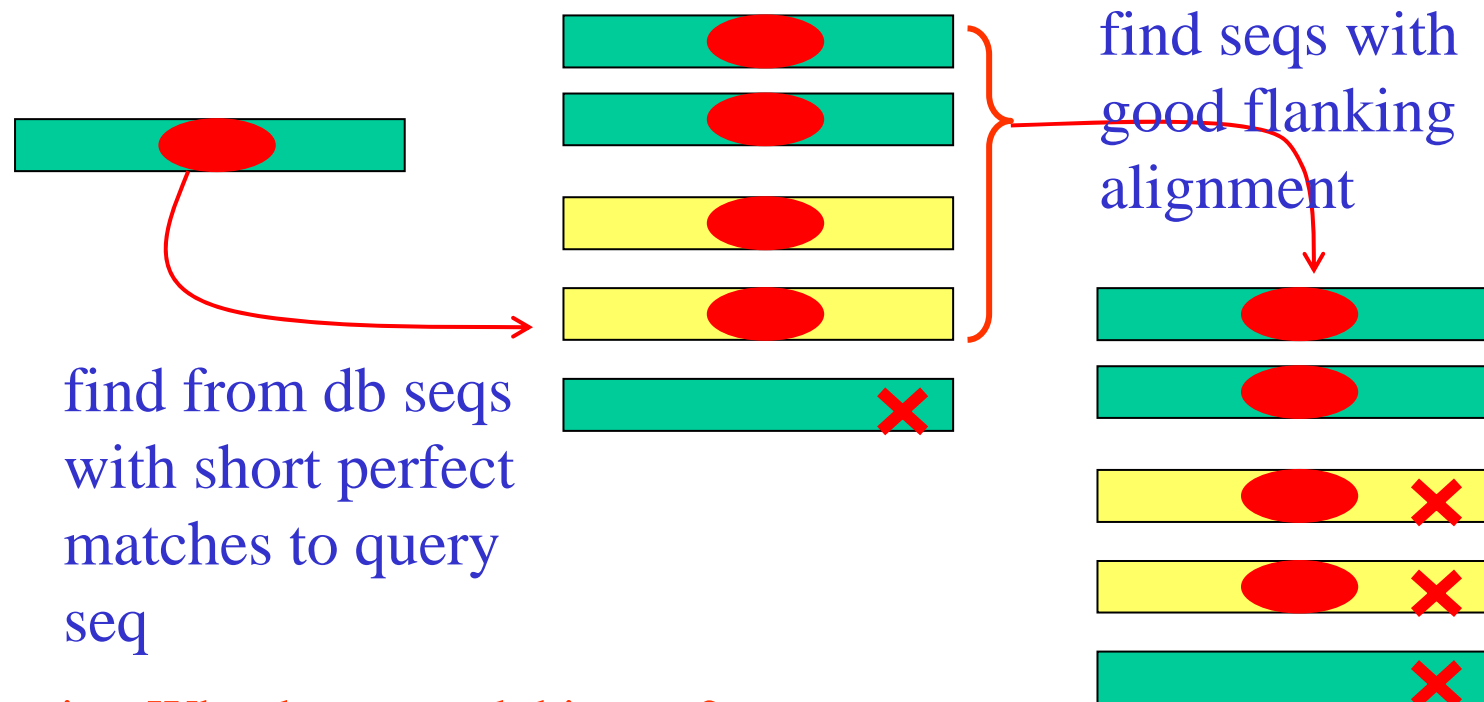


Confirm with suitable wet experiments

# BLAST: How It Works

Altschul et al., *JMB*, 215:403--410, 1990

- **BLAST is one of the most popular tool for doing “guilt-by-association” sequence homology search**



Exercise: Why do we need this step?

# Homologs obtained by BLAST

Sequences producing significant alignments:		Score (bits)	E Value
<a href="#">gi 14193729 gb AAK56109.1 AF332081_1</a>	protein tyrosin phosph...	621 <b>L</b>	e-177
<a href="#">gi 126467 sp P18433 PTRA_HUMAN</a>	Protein-tyrosine phosphatase...	621 <b>L</b>	e-177
<a href="#">gi 4506303 ref NP_002827.1 </a>	protein tyrosine phosphatase, r...	621 <b>L</b>	e-176
<a href="#">gi 227294 prf  1701300A</a>	protein Tyr phosphatase	620	e-176
<a href="#">gi 18450369 ref NP_543030.1 </a>	protein tyrosine phosphatase, ...	621 <b>L</b>	e-176
<a href="#">gi 32067 emb CAA37447.1 </a>	tyrosine phosphatase precursor [Ho...	611 <b>L</b>	e-176
<a href="#">gi 285113 pir  JC1285</a>	protein-tyrosine-phosphatase (EC 3.1....	619	e-176
<a href="#">gi 6981446 ref NP_036895.1 </a>	protein tyrosine phosphatase, r...	611 <b>L</b>	e-176
<a href="#">gi 2098414 pdb 1YFO A</a>	Chain A, Receptor Protein Tyrosine Ph...	61 <b>S</b>	e-174
<a href="#">gi 32313 emb CAA38662.1 </a>	protein-tyrosine phosphatase [Homo...	61 <b>L</b>	e-174
<a href="#">gi 450583 gb AAB04150.1 </a>	protein tyrosine phosphatase >gi 4...	605	e-172
<a href="#">gi 6679557 ref NP_033006.1 </a>	protein tyrosine phosphatase, r...	60 <b>L</b>	e-172
<a href="#">gi 483922 gb AAA17990.1 </a>	protein tyrosine phosphatase alpha	599	e-170

- Thus our example sequence could be a protein tyrosine phosphatase  $\alpha$  (PTP $\alpha$ )

# Example Alignment with PTP $\alpha$

Score = 632 bits (1629), Expect = e-180  
 Identities = 294/302 (97%), Positives = 294/302 (97%)

```

Query: 1   SPSTNRKYPPLPVDKLEEE INRRMADDNKLFREEFNALPACPIQATCEAASXXXXXXXXXR 60
          SPSTNRKYPPLPVDKLEEE INRRMADDNKLFREEFNALPACPIQATCEAAS      R
Sbjct: 202 SPSTNRKYPPLPVDKLEEE INRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR 261

Query: 61  YVNILPYDHSRVHLTPVEGVPSDYINASF INGYQEKKNF IAAQGPKEETVNDFWRMIWE 120
          YVNILPYDHSRVHLTPVEGVPSDYINASF INGYQEKKNF IAAQGPKEETVNDFWRMIWE
Sbjct: 262 YVNILPYDHSRVHLTPVEGVPSDYINASF INGYQEKKNF IAAQGPKEETVNDFWRMIWE 321

Query: 121 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD 180
          QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD
Sbjct: 322 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD 381

Query: 181 VTNRKPQRLITQFHFTSWPDFGVPFITP IGMLKFLKKVKACNPQYAGAI VVHCSAGVGRTG 240
          VTNRKPQRLITQFHFTSWPDFGVPFITP IGMLKFLKKVKACNPQYAGAI VVHCSAGVGRTG
Sbjct: 382 VTNRKPQRLITQFHFTSWPDFGVPFITP IGMLKFLKKVKACNPQYAGAI VVHCSAGVGRTG 441

Query: 241 TFWVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTD MQYVF IYQALLEHYLYGDTELE 300
          TFWVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTD MQYVF IYQALLEHYLYGDTELE
Sbjct: 442 TFWVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTD MQYVF IYQALLEHYLYGDTELE 501
  
```

# Guilt-by-Association: Caveats

- **Ensure that the effect of database size has been accounted for**
- **Ensure that the function of the homology is not derived via invalid “transitive assignment”**
- **Ensure that the target sequence has all the key features associated with the function, e.g., active site and/or domain**

# Interpretation of P-value

- Seq. comparison progs, e.g. BLAST, often associate a P-value to each hit
  - P-value is interpreted as prob that a random seq has an equally good alignment
  - Suppose the P-value of an alignment is  $10^{-6}$
  - If database has  $10^7$  seqs, then you expect  $10^7 * 10^{-6} = 10$  seqs in it that give an equally good alignment
- ⇒ Need to correct for database size if your seq comparison prog does not do that!

Exercise: Name a commonly used method for correcting p-value for a situation like this

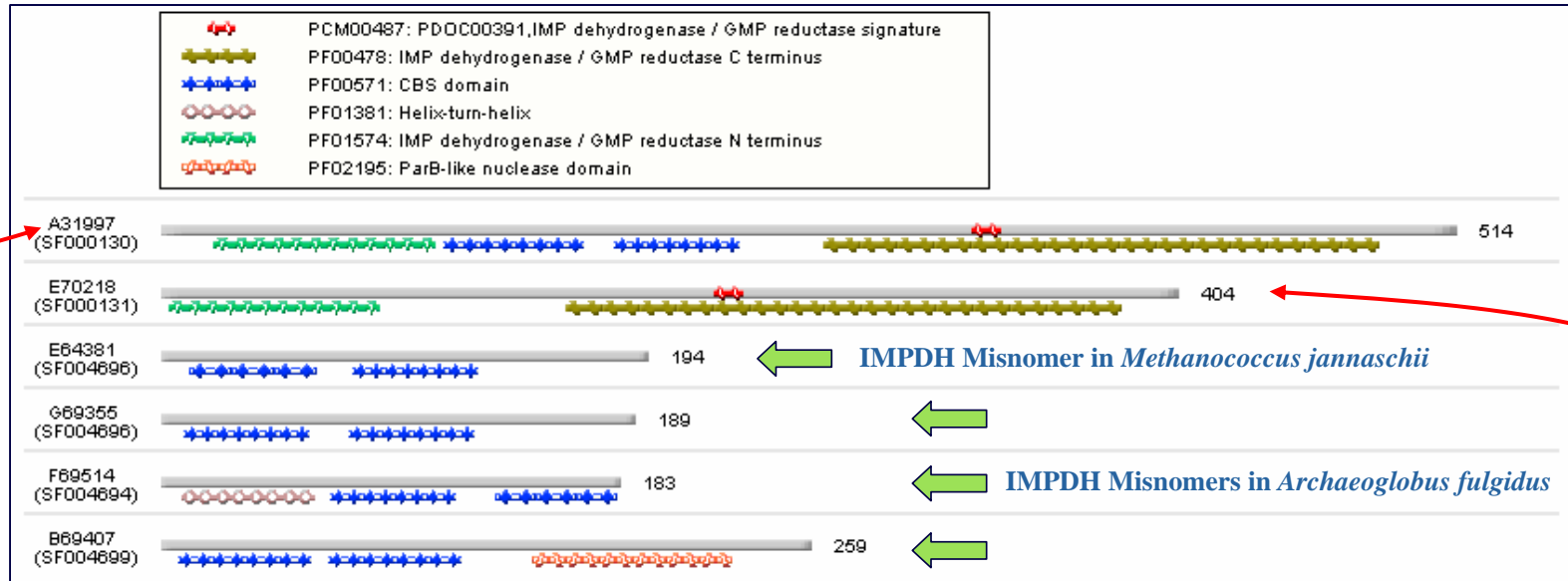


# Examples of Invalid Function Assignment: The IMP Dehydrogenases (IMPDH)

18 entries were found

ID	Organism	PIR	Swiss-Prot/TrEMBL	RefSeq/GenPept
<a href="#">NF00181857</a>	Methanococcus jannaschii	<a href="#">E64381</a> conserved hypothetical protein MJ0653	<a href="#">Y653_METJA</a> Hypothetical protein MJ0653	<a href="#">g1592300</a> inosine-5'-monophosphate dehydrogenase (guaB) <a href="#">NP_247637</a> inosine-5'-monophosphate dehydrogenase (guaB)
<a href="#">NF00187788</a>	Archaeoglobus fulgidus	<a href="#">G69355</a> MJ0653 homolog AF0847 <i>ALT_NAMES</i> : inosine-monophosphate dehydrogenase (guaB-1) homolog [misnomer]	<a href="#">O29411</a> INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-1)	<a href="#">g2649754</a> inosine monophosphate dehydrogenase (guaB-1) <a href="#">NP_069681</a> inosine monophosphate dehydrogenase (guaB-1)
<a href="#">NF00188267</a>	Archaeoglobus fulgidus	<a href="#">F69514</a> yhcV homolog 2 <i>ALT_NAMES</i> : inosine-monophosphate dehydrogenase (guaB-2) homolog [misnomer]	<a href="#">O28162</a> INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-2)	<a href="#">g2648410</a> inosine monophosphate dehydrogenase (guaB-2) <a href="#">NP_070943</a> inosine monophosphate dehydrogenase (guaB-2)
<a href="#">NF00188697</a>	Archaeoglobus fulgidus	<b>A partial list of IMP dehydrogenase misnomers in complete genomes remaining in some public databases</b>		inosine-5'-monophosphate dehydrogenase related protein V
<a href="#">NF00197776</a>	Thermoplasma volcanum			inosine-5'-monophosphate dehydrogenase related protein VII
<a href="#">NF00414709</a>	Methanothermobacter thermautotrophicus			inosine-5'-monophosphate dehydrogenase related protein IX
<a href="#">NF00414811</a>	Methanothermobacter thermautotrophicus			inosine-5'-monophosphate dehydrogenase related protein X
<a href="#">NF00414837</a>	Methanothermobacter thermautotrophicus	<a href="#">H69232</a> MJ1225-related protein MTH992 <i>ALT_NAMES</i> : inosine-5'-monophosphate dehydrogenase related protein IX [misnomer]	<a href="#">O27073</a> INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN IX	<a href="#">g2622093</a> inosine-5'-monophosphate dehydrogenase related protein IX <a href="#">NP_276127</a> inosine-5'-monophosphate dehydrogenase related protein IX
<a href="#">NF00414969</a>	Methanothermobacter thermautotrophicus	<a href="#">B69077</a> yhcV homolog 2 <i>ALT_NAMES</i> : inosine-monophosphate dehydrogenase related protein X [misnomer]	<a href="#">O27616</a> INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN X	<a href="#">g2622697</a> inosine-5'-monophosphate dehydrogenase related protein X <a href="#">NP_276687</a> inosine-5'-monophosphate dehydrogenase related protein X










# IMPDH Domain Structure



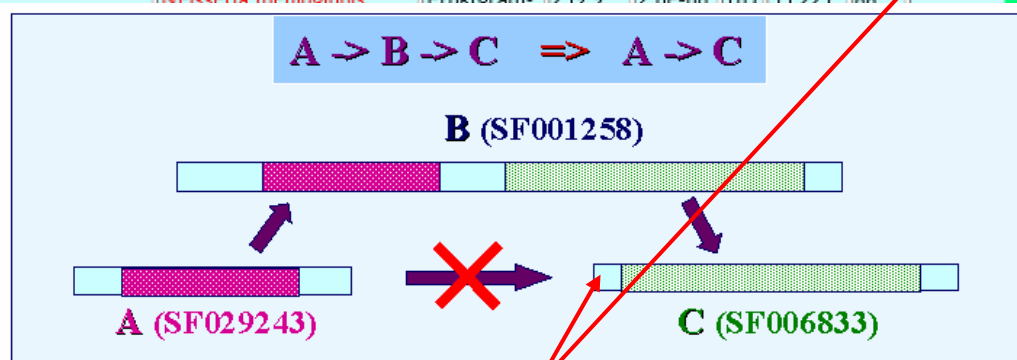
- Typical IMPDHs have 2 IMPDH domains that form the catalytic core and 2 CBS domains.
- A less common but functional IMPDH (E70218) lacks the CBS domains.
- Misnomers show similarity to the CBS domains

# Invalid Transitive Assignment

Root of invalid transitive assignment

<b>B</b> →	<input type="checkbox"/> <a href="#">H70468</a>	<a href="#">SF001258</a>	<a href="#">051440</a>	<a href="#">phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) [similarity]</a>	<a href="#">Aquifex aeolicus</a>	Prok/other	594.3	4.8e-26	205	39.086	197	
	<input type="checkbox"/> <a href="#">S76963</a>	<a href="#">SF001258</a>	<a href="#">039935</a>	<a href="#">phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) [similarity]</a>	<a href="#">Synechocystis sp.</a>	Prok/gram-	557.0	5.7e-24	230	39.175	194	
	<input type="checkbox"/> <a href="#">T35073</a>	<a href="#">SF029243</a>	<a href="#">005738</a>	<a href="#">probable phosphoribosyl-AMP cyclohydrolase</a>	<a href="#">Streptomyces coelicolor</a>	Prok/gram+	399.3	3.5e-15	128	42.157	102	
	<input type="checkbox"/> <a href="#">S53349</a>	<a href="#">SF001257</a>	<a href="#">001188</a>	<a href="#">phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) / histidinol dehydrogenase (EC 1.1.1.23)</a>	<a href="#">Saccharomyces cerevisiae</a>	Euk/fungi	384.1	2.5e-14	799	31.863	204	
<b>A</b> →	<input type="checkbox"/> <a href="#">E69493</a>	<a href="#">SF029243</a>	<a href="#">005738</a>	<a href="#">phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) [similarity]</a>	<a href="#">Archaeoglobus fulgidus</a>	Archae	396.8	4.8e-15	108	47.778	90	
<b>C</b> →	<input type="checkbox"/> <a href="#">G64337</a>	<a href="#">SF006833</a>	<a href="#">030827</a>	<a href="#">phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) [similarity]</a>	<a href="#">Methanococcus jannaschii</a>	Archae	246.9	1.1e-06	95	36.842	95	
	<input type="checkbox"/> <a href="#">D81178</a>	<a href="#">SF006833</a>	<a href="#">101491</a>	<a href="#">phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) NMB0603 [similarity]</a>	<a href="#">Neisseria meningitidis</a>	Prok/gram-	239.9	2.6e-06	107	35.227	88	
	<input type="checkbox"/> <a href="#">G81925</a>	<a href="#">SF006833</a>	<a href="#">101491</a>	<a href="#">phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) NMA0807 [similarity]</a>								
	<input type="checkbox"/> <a href="#">S51513</a>	<a href="#">SF001257</a>	<a href="#">001188</a>	<a href="#">phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) / histidinol dehydrogenase (EC 1.1.1.23)</a>								

Mis-assignment  
of function

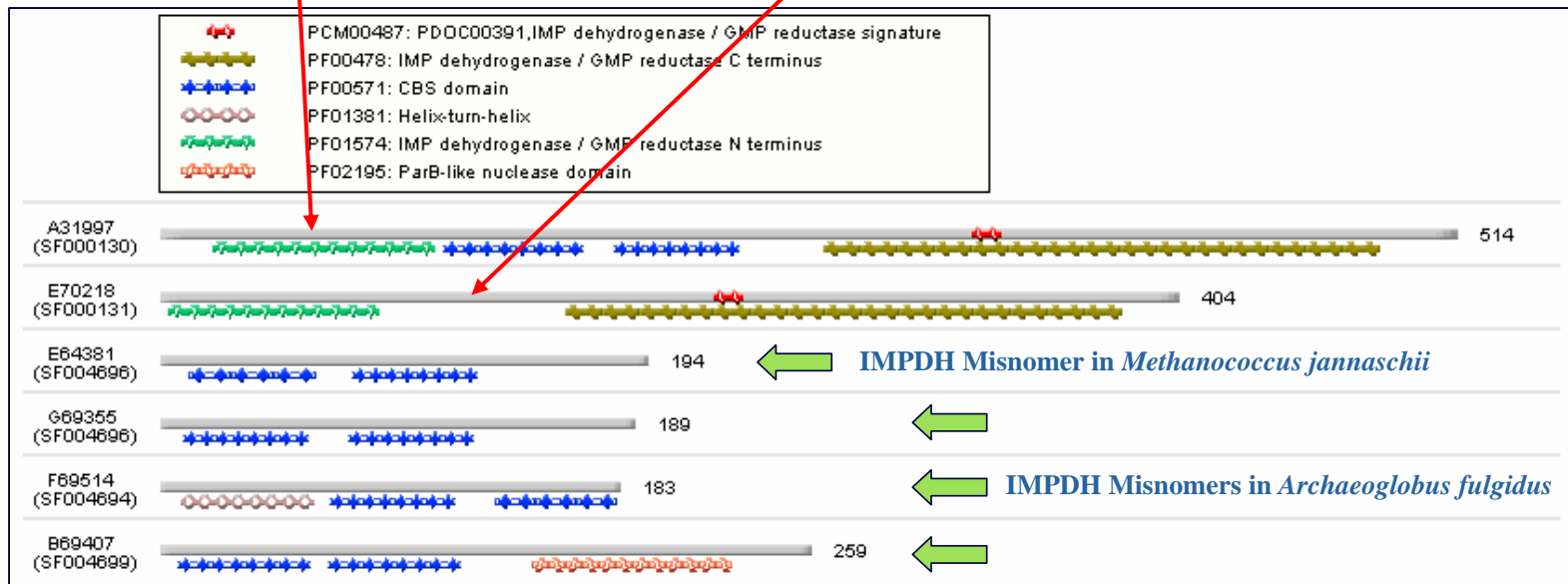


No IMPDH domain

# Emerging Pattern

Typical IMPDH

Functional IMPDH w/o CBS



- Most IMPDHs have 2 IMPDH and 2 CBS domains
  - Some IMPDH (E70218) lacks CBS domains
- ⇒ IMPDH domain is the emerging pattern

# Guilt-by-Association: What if no homolog of known function is found?

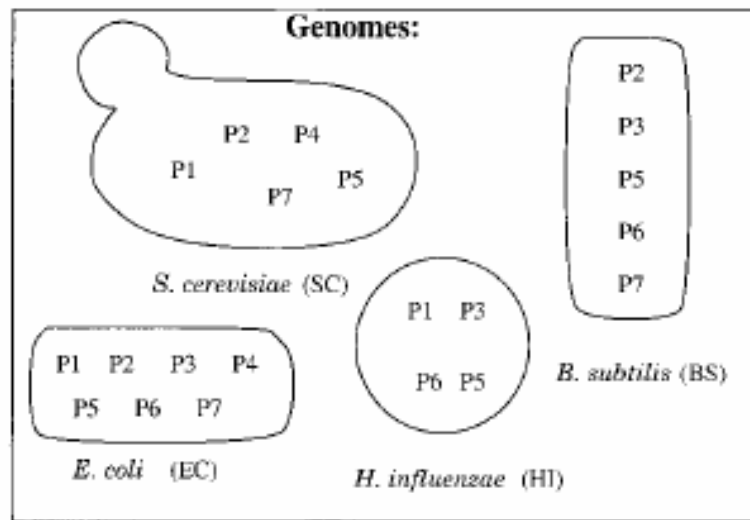
**genome phylogenetic profiles**  
**protfun's feature profiles**  
**SVM Pairwise**  
**Level-2 Neighbours**



# Phylogenetic Profiling

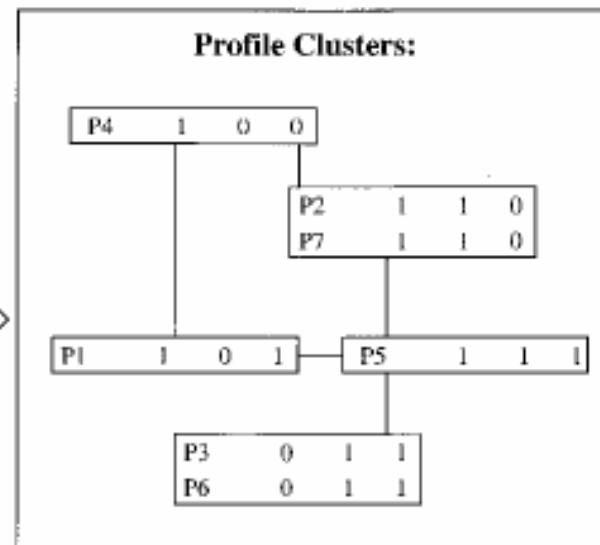
Pellegrini et al., *PNAS*, 96:4285--4288, 1999

- **Gene (and hence proteins) with identical patterns of occurrence across phyla tend to function together**
- ⇒ **Even if no homolog with known function is available, it is still possible to infer function of a protein**



**Phylogenetic Profile:**

	EC	SC	BS	HI
P1	1	0	1	
P2	1	1	0	
P3	0	1	1	
P4	1	0	0	
P5	1	1	1	
P6	0	1	1	
P7	1	1	0	



**Conclusion:** P2 and P7 are functionally linked,  
 P3 and P6 are functionally linked

# Phylogenetic Profiling: How it Works

# Phylogenetic Profiling: P-value

The probability of observing by chance  $z$  occurrences of genes  $X$  and  $Y$  in a set of  $N$  lineages, given that  $X$  occurs in  $x$  lineages and  $Y$  in  $y$  lineages is

$$P(z|N, x, y) = \frac{w_z * \overline{w}_z}{W}$$

where

**No. of ways to distribute  $z$  co-occurrences over  $N$  lineage's**

**No. of ways to distribute the remaining  $x - z$  and  $y - z$  occurrences over the remaining  $N - z$  lineage's**

$$w_z = \binom{N}{z}$$

$$\overline{w}_z = \binom{N - z}{x - z} * \binom{N - z}{y - z}$$

$$W = \binom{N}{x} * \binom{N}{y}$$

**No. of ways of distributing  $X$  and  $Y$  over  $N$  lineage's without restriction**



# Phylogenetic Profiles: Evidence

Pellegrini et al., *PNAS*, 96:4285--4288, 1999

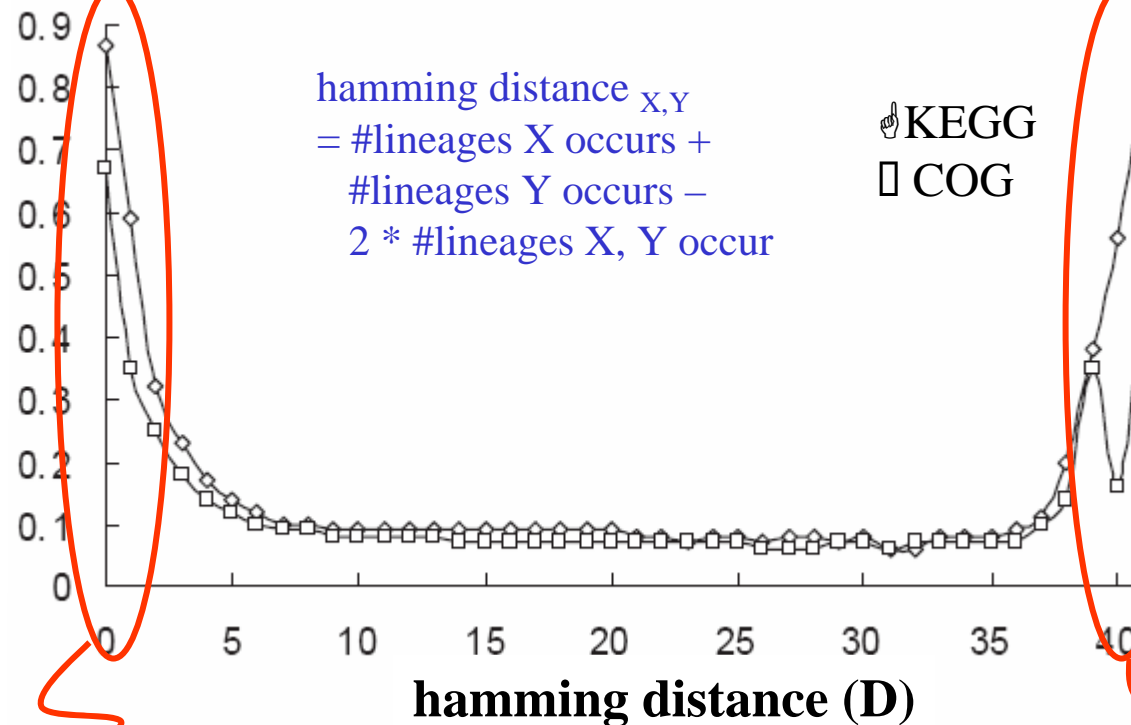
Keyword	No. of non-homologous proteins in group	No. neighbors in keyword group	No. neighbors in random group
Ribosome	60	197	27
Transcription	36	17	10
tRNA synthase and ligase	26	11	5
Membrane proteins*	25	89	5
Flagellar	21	89	3
Iron, ferric, and ferritin	19	31	2
Galactose metabolism	18	31	2
Molybdoterin and Molybdenum, and molybdoterin	12	6	1
Hypothetical <sup>†</sup>	1,084	108,226	8,440

- **E. coli proteins grouped based on similar keywords in SWISS-PROT have similar phylogenetic profiles**

# Phylogenetic Profiling: Evidence

Wu et al., *Bioinformatics*, 19:1524--1530, 2003

fraction of gene pairs  
having hamming distance D  
and share a common pathway  
in KEGG/COG



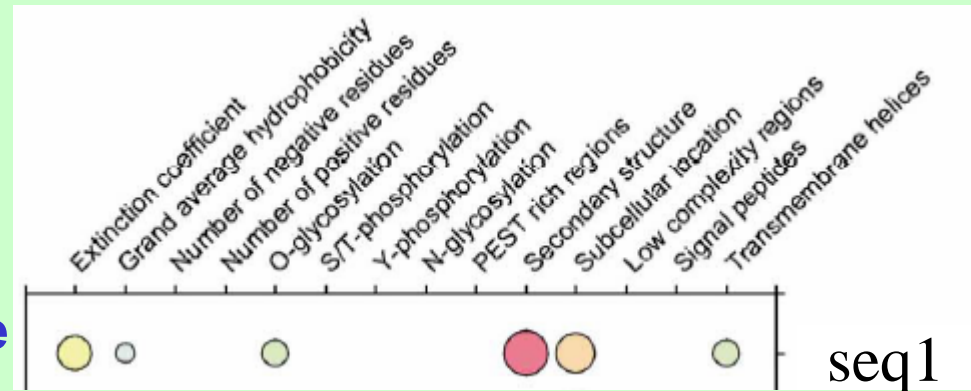
- Proteins having low hamming distance (thus highly similar phylogenetic profiles) tend to share common pathways

Exercise: Why do proteins having high hamming distance also have this behaviour?

# The ProtFun Approach

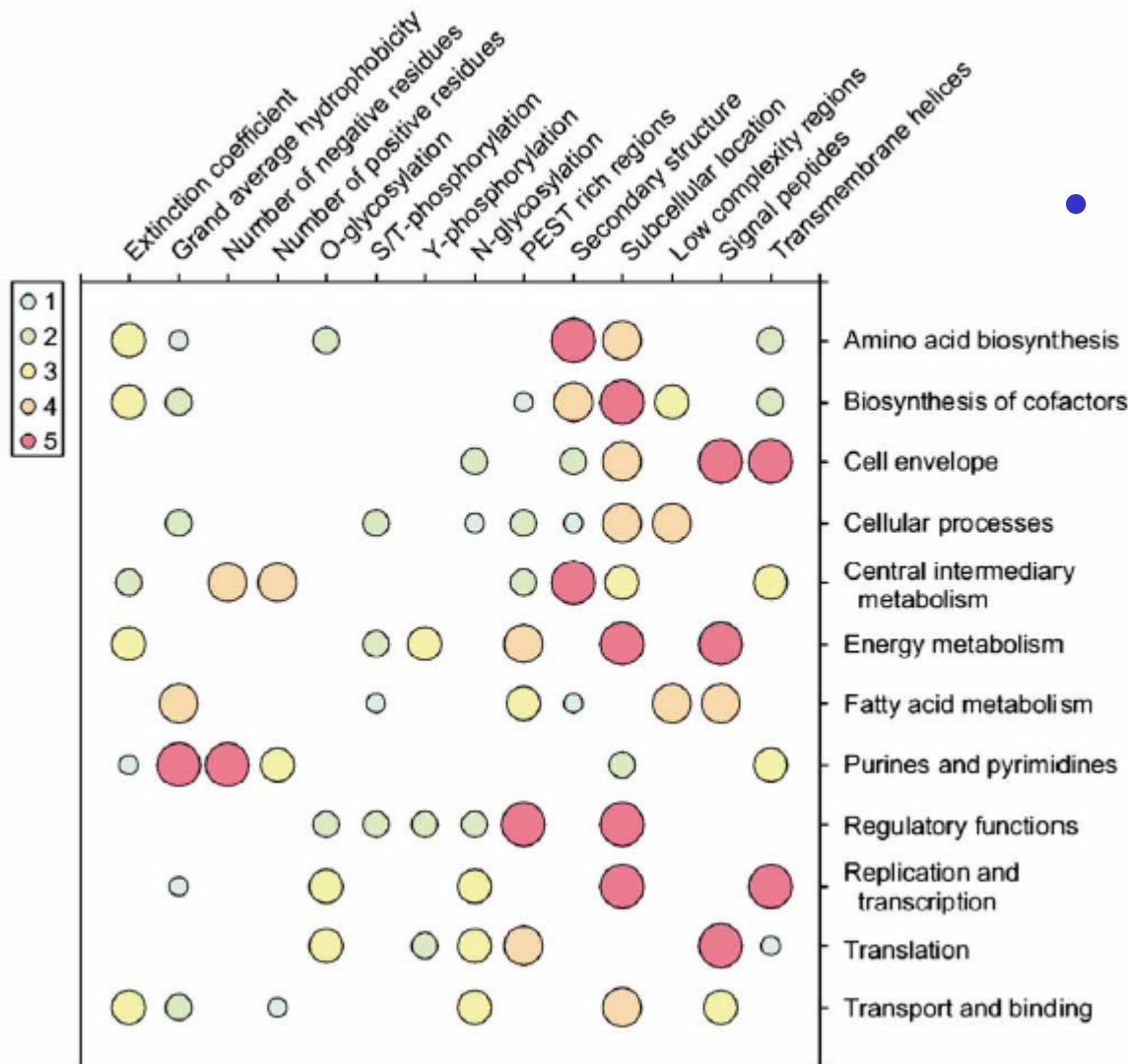
Jensen, *JMB*, 319:1257--1265, 2002

- A protein is not alone when performing its biological function
- It operates using the same cellular machinery for modification and sorting as all other proteins do, such as glycosylation, phosphorylation, signal peptide cleavage, ...
- These have associated consensus motifs, patterns, etc.



- Proteins performing similar functions should share some such “features”
- ⇒ Perhaps we can predict protein function by comparing its “feature” profile with other proteins?

# ProtFun: Evidence



- Combinations of “features” seem to characterize some functional categories

# ProtFun: Example Output

	Prion	A4	TTHY
Amino acid biosynthesis	0.011	0.011	0.011
Biosynthesis of cofactors	0.041	0.161	0.034
Cell envelope	0.146	0.804	0.698
Cellular processes	0.027	0.027	0.051
Central intermediary metabolism	0.047	0.139	0.059
Energy metabolism	0.029	0.023	0.046
Fatty acid metabolism	0.017	0.017	0.023
Purines and pyrimidines	0.528	0.417	0.153
Regulatory functions	0.013	0.014	0.014
Replication and transcription	0.020	0.029	0.040
Translation	0.035	0.027	0.032
Transport and binding	0.831	0.827	0.812
Enzyme	0.233	0.367	0.227
Non-enzyme	0.767	0.633	0.773
Oxidoreductase (EC 1.-.-.-)	0.070	0.024	0.055
Transferase (EC 2.-.-.-)	0.031	0.208	0.037
Hydrolase (EC 3.-.-.-)	0.101	0.090	0.208
Isomerase (EC 4.-.-.-)	0.020	0.020	0.020
Ligase (EC 5.-.-.-)	0.010	0.010	0.010
Lyase (EC 6.-.-.-)	0.017	0.078	0.017

- At the seq level, Prion, A4, & TTHY are dissimilar
- ProtFun predicts them to be cell envelope-related, transport & binding
- This is in agreement w/ known functionality of these proteins

# SVM-Pairwise Framework

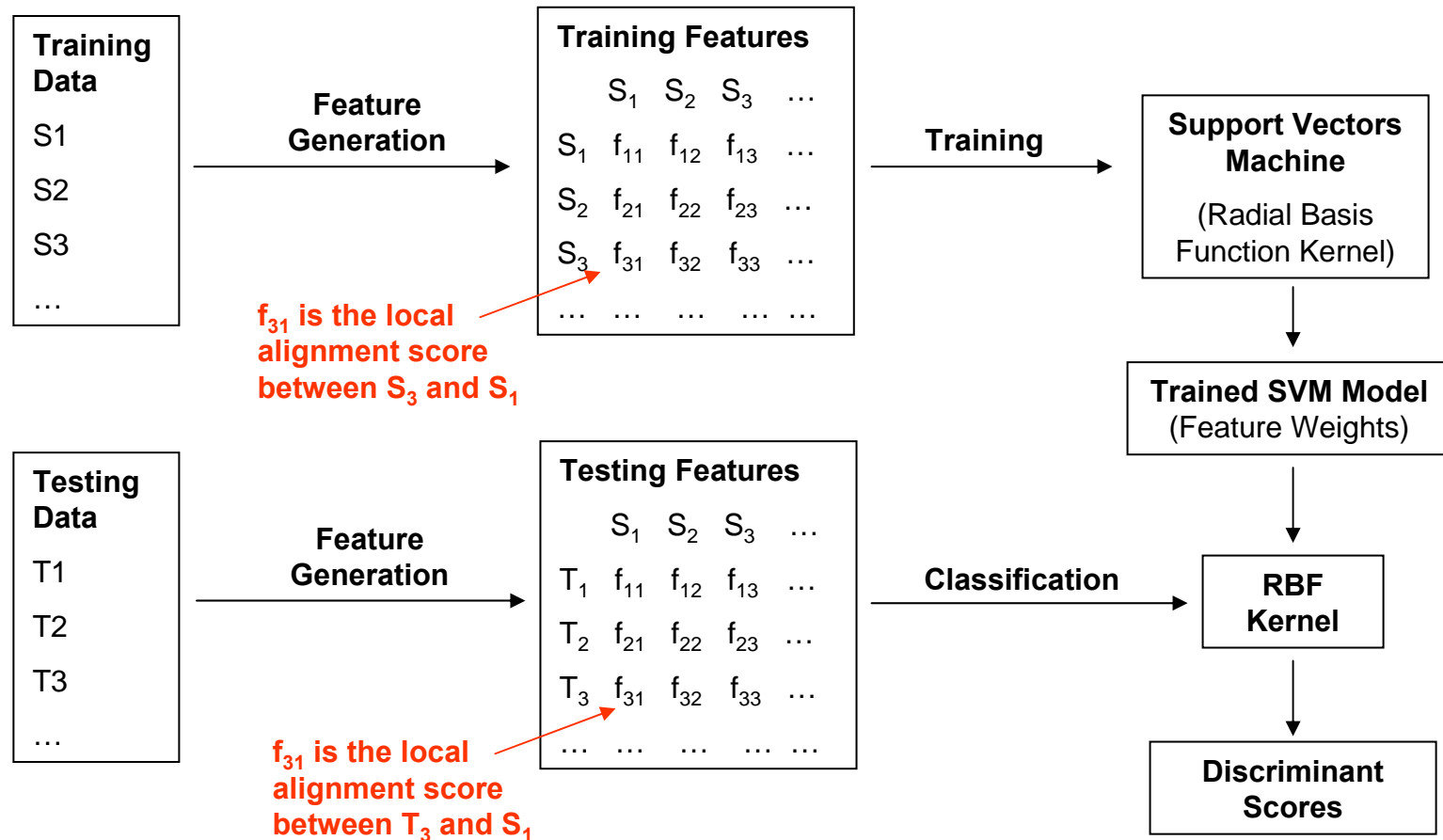
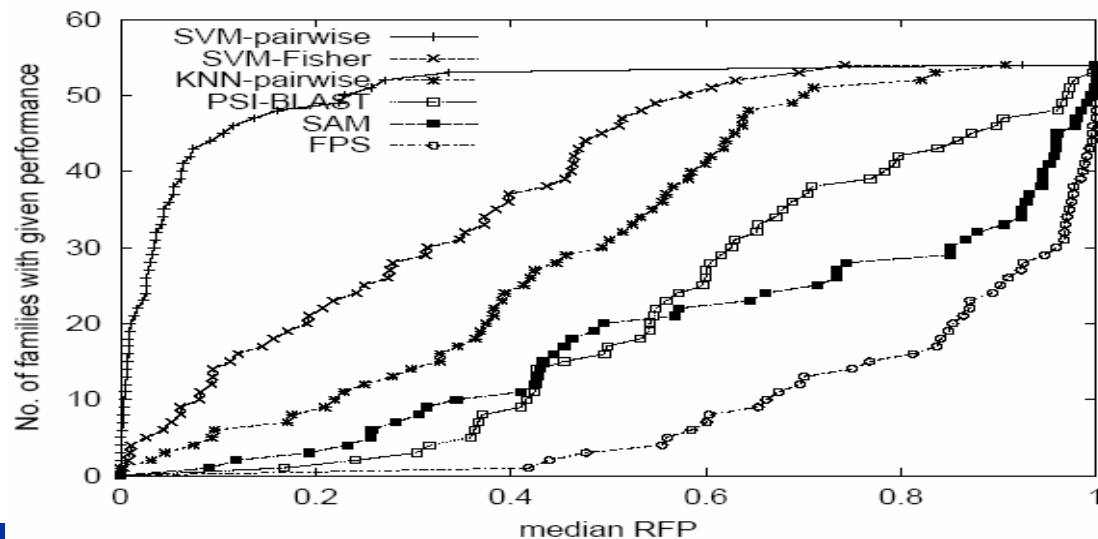
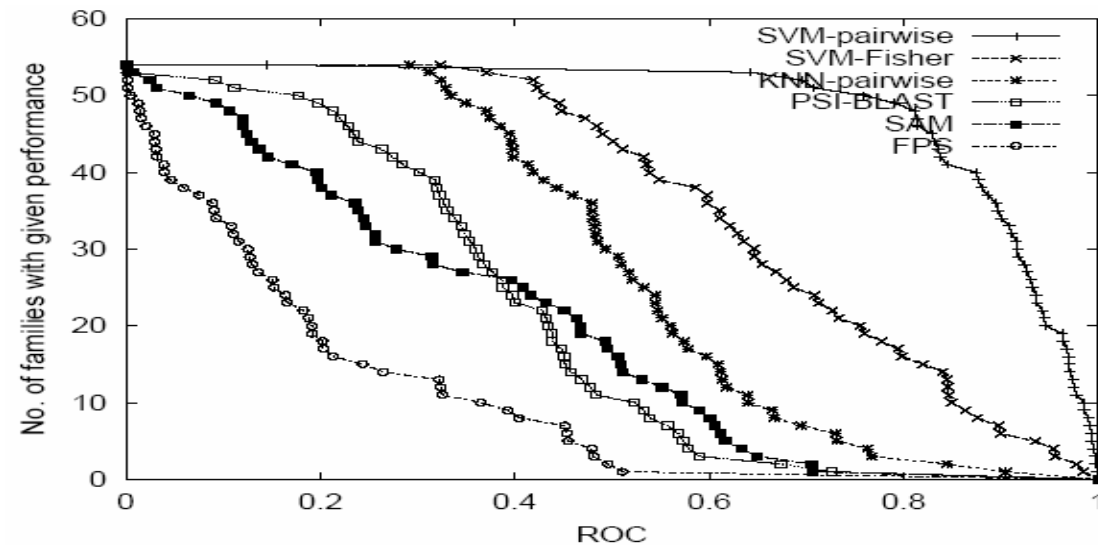


Image credit: Kenny Chua

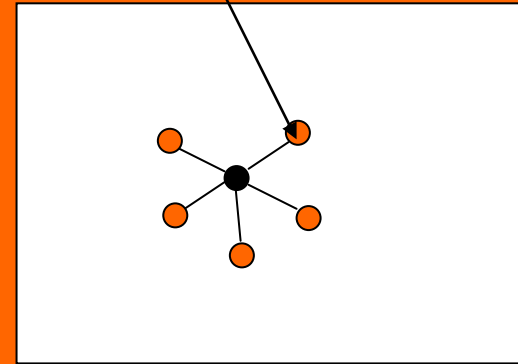
# Performance of SVM-Pairwise

- **Receiver Operating Characteristic (ROC)**
  - The area under the curve derived from plotting true positives as a function of false positives for various thresholds.
- **Rate of median False Positives (RFP)**
  - The fraction of negative test examples with a score better or equals to the median of the scores of positive test examples.

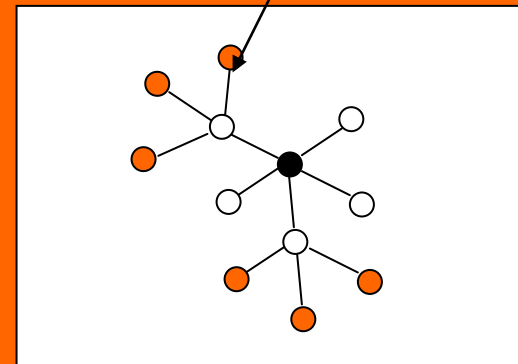


# Protein Function Prediction from Protein Interactions

Level-1 neighbour

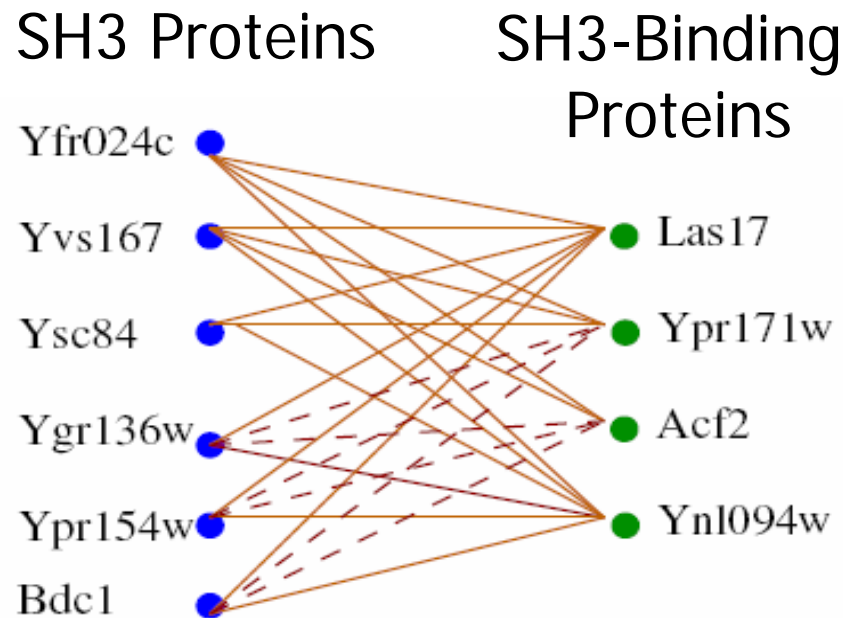


Level-2 neighbour



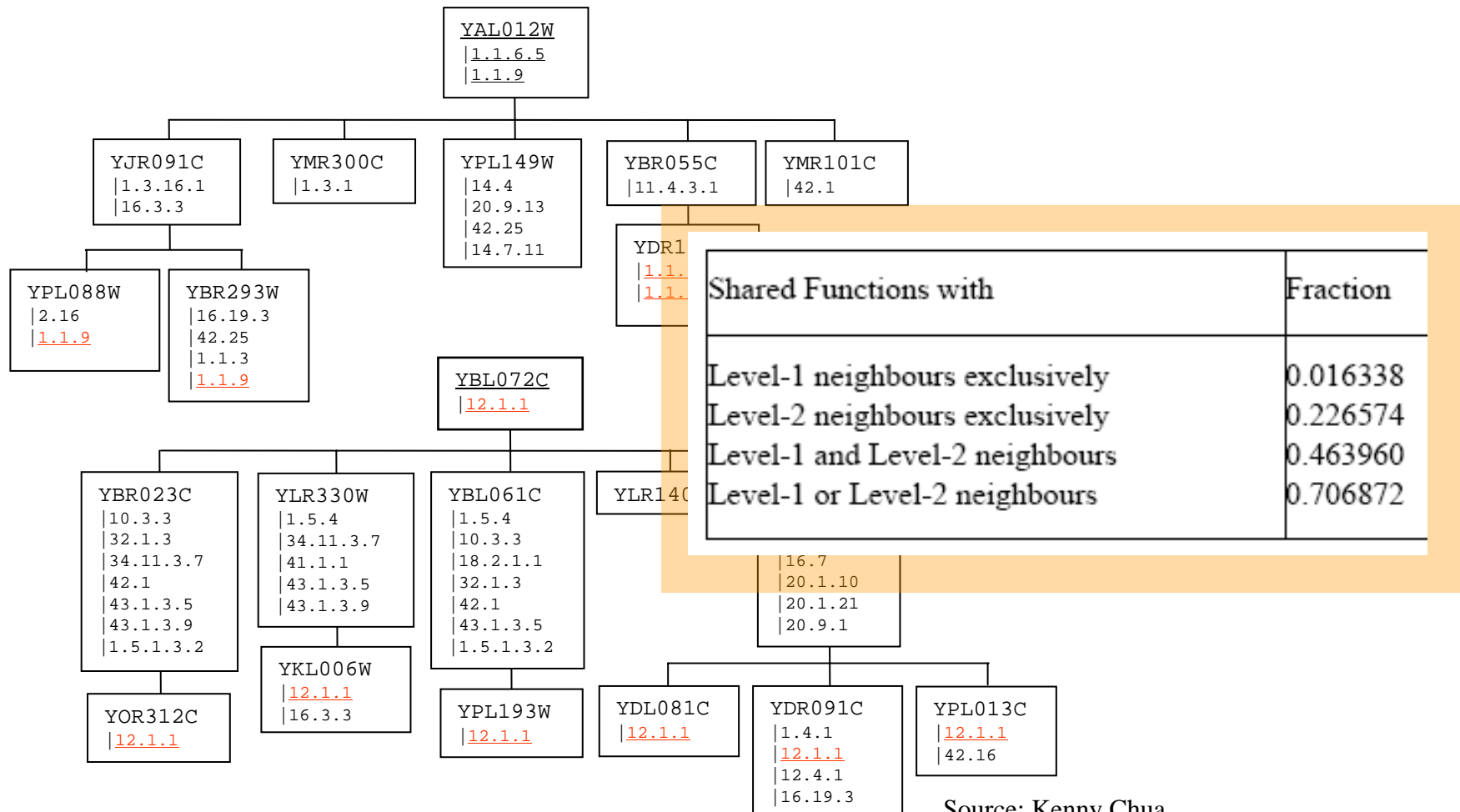


# An illustrative Case of Indirect Functional Association?



- 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

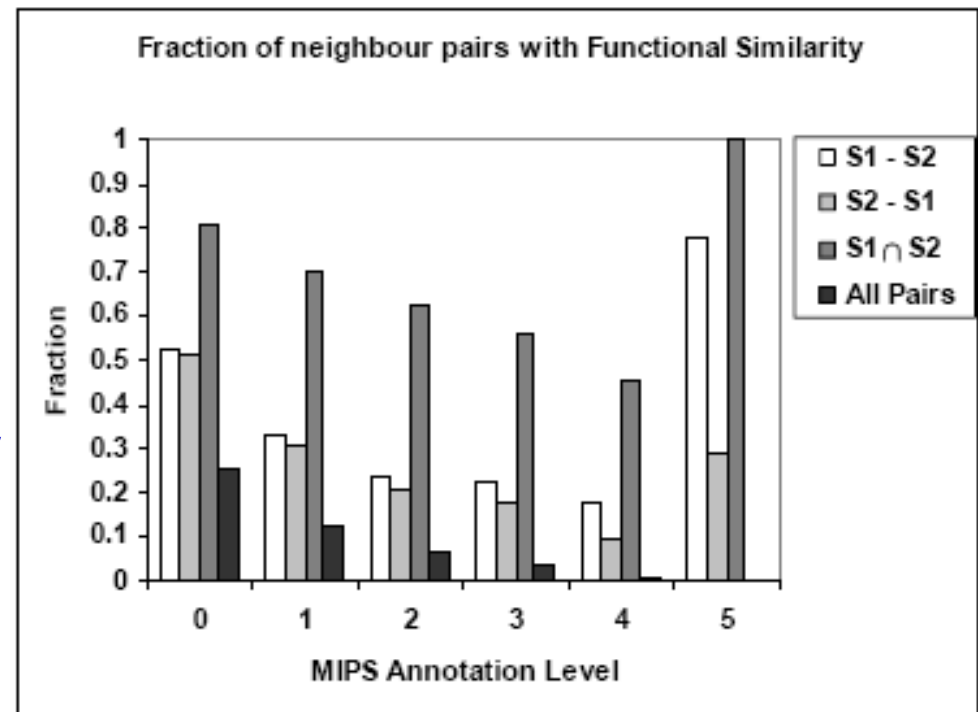
# Over-Rep of Functions in Neighbours

- Functional Similarity:**

$$S(i, j) = \frac{|F_i \cap F_j|}{|F_i \cup F_j|}$$

- where  $F_k$  is the set of functions of protein  $k$

- L1  $\cap$  L2 neighbours show greatest over-rep**
- L3 neighbours show little observable over-rep**



Source: Kenny Chua

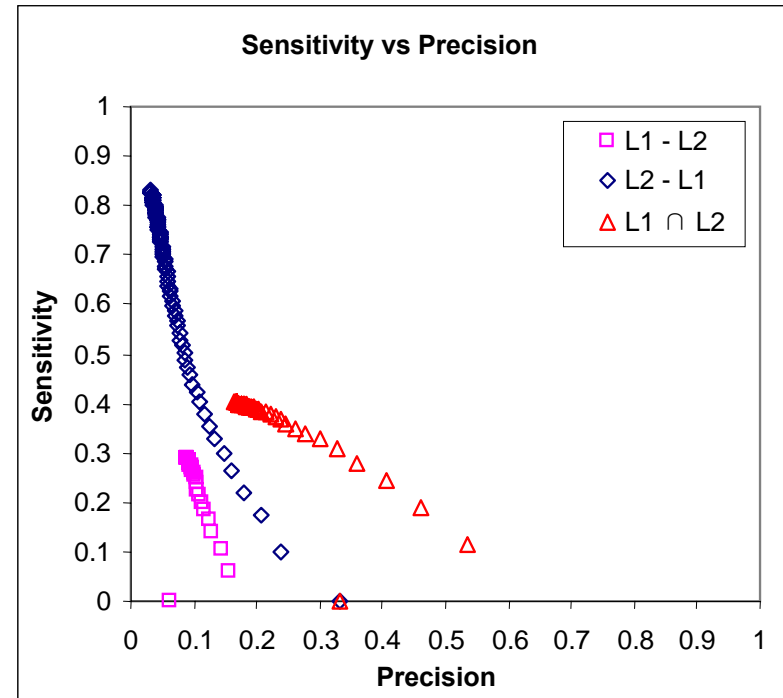
# Prediction Power By Majority Voting

Source: Kenny Chua

- 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

# Use L1 & L2 Neighbours for Prediction

- **Weighted Average**

- Over-rep of functions in L1 and L2 neighbours
- Each observation of L1 or L2 neighbour is summed

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

- $S_{TR}(u, v)$  is an “index” for function xfer betw u and v,
- $\delta(k, x) = 1$  if k has function x, 0 otherwise
- $N_k$  is the set of interacting partners of k
- $\pi_x$  is freq of function x in the dataset
- $\lambda$  is contribution of background freq to the score
- $r_{\text{int}}$  is fraction of all interaction pairs that share some functions

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

Source: Kenny Chua

# 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**

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

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

Source: Kenny Chua

# Functional Similarity Estimate: Modified Equiv Measure

- **Modified Equivalence 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**

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

Exercise: What else should we consider in this formula?

Source: Kenny Chua

# 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

Source: Kenny Chua



# An “Index” for Function Transfer Based on Reliability of Interactions

- Take reliability into consideration when computing Equiv Measure:

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

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

# Functional Similarity Estimate: Transitive Functional Association



- If protein  $u$  is similar to protein  $w$ , and protein  $w$  is similar to protein  $v$ , proteins  $u$  and  $v$  may show some degree of similarity
- So we estimate functional similarity betw  $u$  and  $v$  by product of functional similarity betw  $u$  and  $w$ , and that between  $w$  and  $v$ :

$$S_{TR}(u, v) = \max \left( S_R(u, v), \max_{w \in N_u} S_R(u, w) S_R(w, v) \right)$$

# Correlation with Functional Similarity

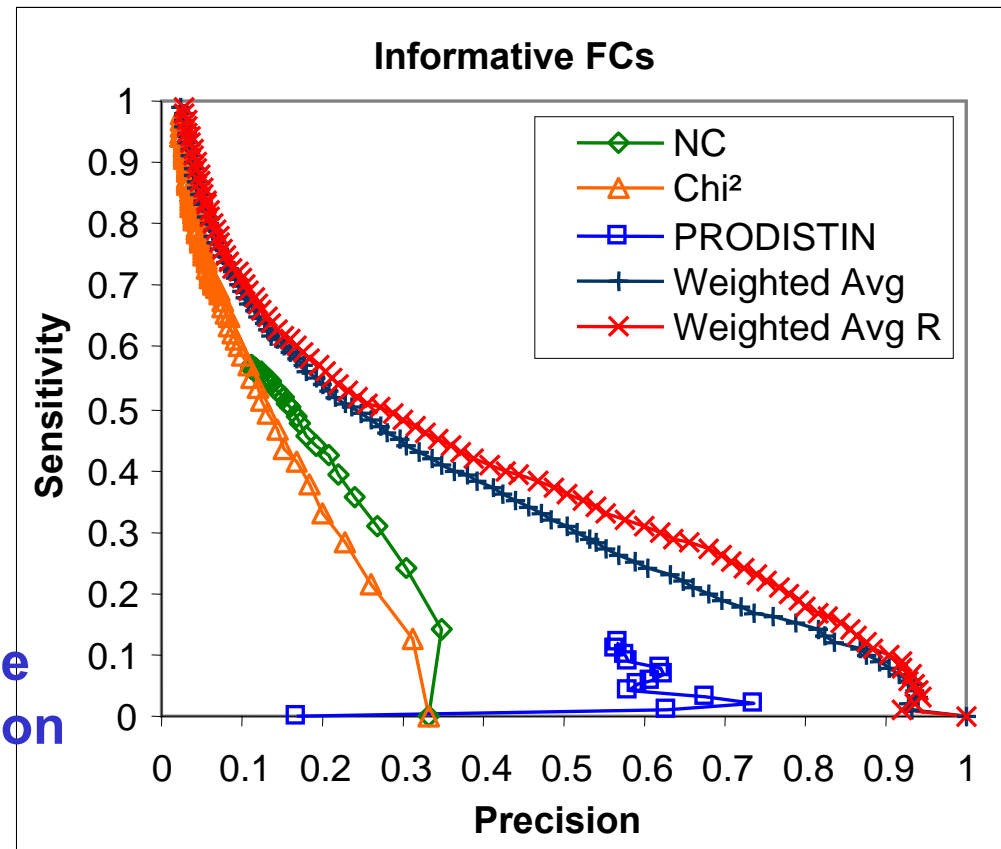
- **Equiv measure shows improved correlation w/ functional similarity when reliability of interactions & transitive association is considered:**

Neighbours	CD-Distance	FS-Weight	FS-Weight R	Transitive FS-Weight R
$S_1$	0.471810	0.498745	0.532596	<b>0.532626</b>
$S_2$	0.224705	0.298843	0.375317	<b>0.381966</b>
$S_1 \cup S_2$	0.224581	0.29629	0.363025	<b>0.369378</b>

Source: Kenny Chua

# Performance Evaluation

- Prediction performance improves after incorporation of interaction reliability
- ⇒ Indirect functional association is plausible
- ⇒ It is found often in real interaction data
- ⇒ It can be used to improve protein function prediction from protein interaction data



Source: Kenny Chua

Any Questions?



# Acknowledgements

- **Some of the slides are based on slides given to me by Kenny Chua**

# References

- T.F.Smith & X.Zhang. “The challenges of genome sequence annotation or ‘The devil is in the details’”, *Nature Biotech*, 15:1222--1223, 1997
- D. Devos & A.Valencia. “Intrinsic errors in genome annotation”, *TIG*, 17:429--431, 2001
- S.F.Altshcul et al. “Basic local alignment search tool”, *JMB*, 215:403--410, 1990
- S.F.Altschul et al. “Gapped BLAST and PSI-BLAST: A new generation of protein database search programs”, *NAR*, 25(17):3389--3402, 1997
- H.N. Chua, W.-K. Sung. [A better gap penalty for pairwise SVM.](#) Proc. APBC05, pages 11-20
- T. Jaakkola, M. Diekhans, and D. Haussler. A discriminative framework for detecting remote homologies. *JCB*, 7(1-2):95—11, 2000

# References

- S.E.Brenner. “Errors in genome annotation”, *TIG*, 15:132--133, 1999
- M. Pellegrini et al. “Assigning protein functions by comparative genome analysis: Protein phylogenetic profiles”, *PNAS*, 96:4285--4288, 1999
- J. Wu et al. “Identification of functional links between genes using phylogenetic profiles”, *Bioinformatics*, 19:1524--1530, 2003
- L.J.Jensen et al. “Prediction of human protein function from post-translational modifications and localization features”, *JMB*, 319:1257--1265, 2002
- C. Wu, W. Barker. “A Family Classification Approach to Functional Annotation of Proteins”, *The Practical Bioinformatician*, Chapter 19, pages 401—416, WSPC, 2004