For written notes on this lecture, please read chapter 19 of *The Practical Bioinformatician* and *Hawkins & Kihara, JBCB 5(1):1-30, 2007* 

#### CS2220: Introduction to Computational Biology Unit 5: Sequence Homology Interpretation

Wong Limsoon







- Recap of sequence alignment
- Guilt by association
- Active site/domain discovery
- What if no homology of known function is found?
  - Genome phylogenetic profiling
  - SVM-Pairwise
  - Protein-protein interactions
- Key mutation site discovery



# Brief recap of sequence comparison / alignment

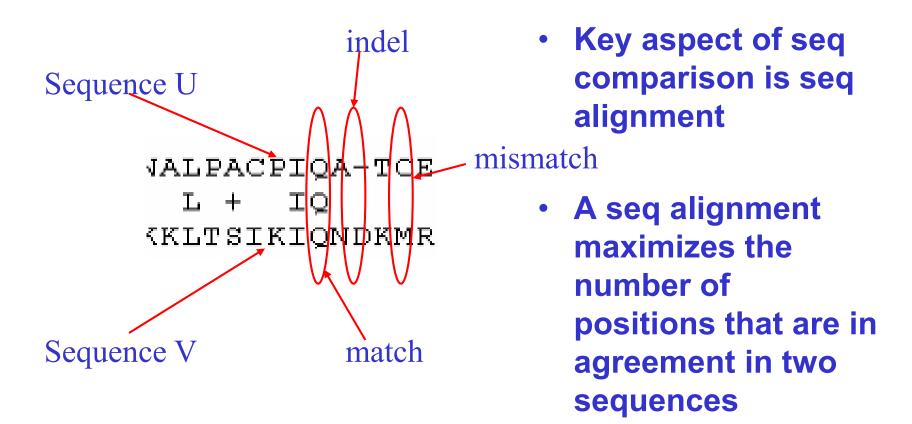


#### Motivations for seq comparison 55

- DNA is blue print for living organisms
- $\Rightarrow$  Evolution is related to changes in DNA
- ⇒ By comparing DNA sequences we can infer evolutionary relationships between the sequences w/o knowledge of the evolutionary events themselves
- Foundation for inferring function, active site, and key mutations

#### Sequence alignment





# Sequence alignment: Poor examp

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 MPHNVHFVAGVLGEAALKGPMMKKEOAYSLTFTEAGTYDYHCTPHPFMRGKVVVE :: Ascorbate Oxidase ILORGTPWADGTASISOCAINPGETFFYNFTVDNPGTFFYHGHLGMORSAGLYGSLI 70 80 90 100 110 120 No obvious match between Amicyanin and Ascorbate Oxidase



6

# Sequence Alignment: Good exampter Singapore

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

D >gi|13476732|ref|NP\_108301.1| unknown protein [Mesorhizobium loti]
gi|14027493|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 MKPGRLASIALAIIFLPMAVPAHAATIEITMENLVISPTEVSAKVGDTIRWVNKDVFAHT 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



### Multiple alignment: An example

- Multiple seq alignment maximizes number of positions in agreement across several seqs
- seqs belonging to same "family" usually have more conserved positions in a multiple seq alignment
- gi|126467| FHFTSWPDFGVPFTPIGMLKFLKKVKACNP--OYAGAIV/HCSAGVGRTGTFVVIDAMLD gi|2499753 FHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIMLD YHYTQWPDMGVPEYALPVLTFVRRSSAARM--PETGPVIVHCSAGVGRTGTYIVIDSMLQ gi|462550| FHFTSWPDHGVPDTTDLLINFRYLVRDYMKOSPPESPILVHCSAGVGRTGTFIAIDRLIY gi|2499751 qi|1709906 FOF TAWPDHGVPEHPTPFLAFLRRVKTCNP--PDAGPMVVHCSAGVGRTGCF IVIDAMLE LHFTSWPDFGVPFTPIGMLKFLKKVKTLNP--VHAGPIVVHCSAGVGRTGTFIVIDAMMA gi|126471| gi|548626| FHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIMLD FHFTGWPDHGVPYHATGLLGFVRQVKSKSP--PNAGPLVVHCSAGAGRTGCFIVIDIMLD gi|131570| gi|2144715 FHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGTFIAIDRLIY \* \*\*\* \*\*\* \*\*\*\*\* \*

Conserved sites

8

Application of sequence comparison: Guilt-by-association



#### A protein is a ...



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





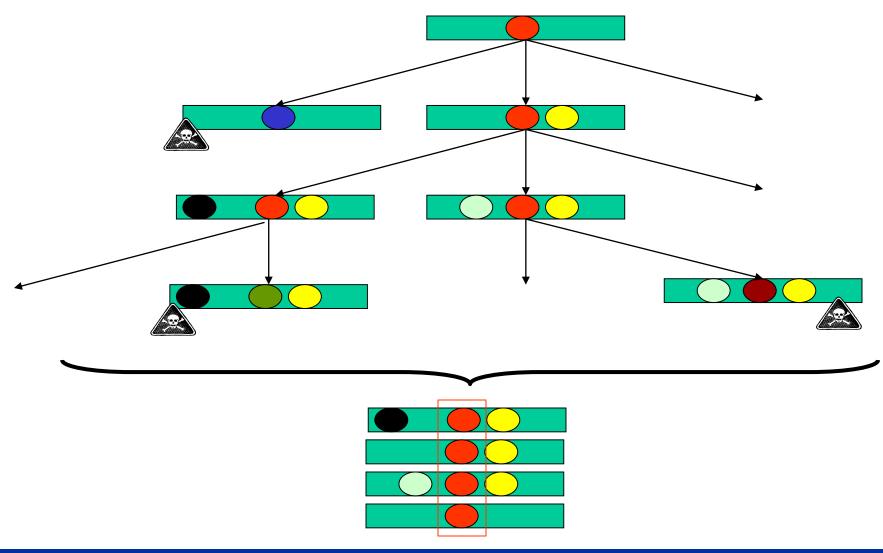
SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE VT

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

#### In the course of evolution...



2



#### CS2220, AY17/18

#### Copyright 2017 © Wong Limsoon



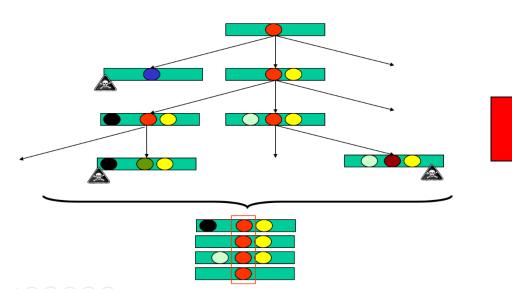


Let a = AFPHQHRVP Let **b** = PQVYNIMKE Suppose each generation differs from the previous by 1 residue What is the average difference between the 2<sup>nd</sup> generation of a What is the average difference between the 2<sup>nd</sup> generation of a and b?

#### The triumph of logic



In the course of evolution...



Two proteins inheriting their function from a common ancestor have very similar amino acid sequences





15

# How can we guess the function of a protein?





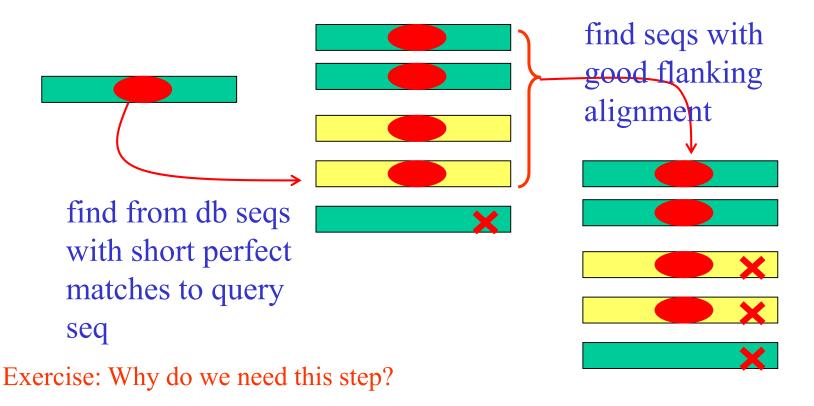
Copyright 2017 © Wong Limsoon

#### BLAST: How it works Altschul et al., *JMB*, 215:403--410, 1990



9

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



CS2220, AY17/18



# Homologs obtained by BLAST

Sequences producing significant alignments:	Score (bits)	E Value
<u>qi 14193729 qb AAK56109.1 AF332081_1</u> protein tyrosin phosph	<u>62:</u> L	e-177
<u>gi 126467 sp P18433 PTRA_HUMAN</u> Protein-tyrosine phosphatase	<u>621 L</u>	e-177
<u>gi 4506303 ref NP_002827.1 </u> protein tyrosine phosphatase, r <u>gi 227294 prf  1701300A</u> protein Tyr phosphatase	621 L	e-176 e-176
<u>gi 18450369 ref NP_543030.1 </u> protein tyrosine phosphatase,	<u>621 L</u>	e-176
<u>qi 32067 emb CAA37447.1 </u> tyrosine phosphatase precursor [Ho <u>qi 285113 pir  JC1285</u> protein-tyrosine-phosphatase (EC 3.1	61: 619	e-176 e-176
<u>gi 6981446 ref NP_036895.1 </u> protein tyrosine phosphatase, r	<u>61;</u>	e-176
gi 2098414 pdb 1YFO A Chain A, Receptor Protein Tyrosine Ph	<u>61</u> S	e-174
<u>qi 32313 emb CAA38662.1 </u> protein-tyrosine phosphatase [Homo <u>qi 450583 qb AAB04150.1 </u> protein tyrosine phosphatase >gi 4	61 L 605	e-174 e-172
<u>qi 6679557 ref NP_033006.1 </u> protein tyrosine phosphatase, r <u>qi 483922 qb AAA17990.1 </u> protein tyrosine phosphatase alpha	<u>60.</u> 599	e-172 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%)

- Sbjct: 202 SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR 261
- Query: 61 YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE 120 YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE
- Sbjct: 262 YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE 321
- Query: 121 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD 180 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD
- Sbjct: 322 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD 381
- Query: 181 VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG 240 VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG
- Sbjct: 382 VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG 441
- Query: 241 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQVVFIYQALLEHVLYGDTELE 300 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQVVFIYQALLEHVLYGDTELE
- Sbjct: 442 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE 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

# Law of large numbers



- Suppose you are in a room with 365 other people
- Q: What is the prob that a specific person in the room has the same birthday as you?
- A: 1/365 = 0.3%

- Q: What is the prob that there is a person in the room having the same birthday as you?
- A:  $1 (364/365)^{365} = 63\%$
- Q: What is the prob that there are two persons in the room having the same birthday?
- A: 100%

### 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<sup>-6</sup>
- If database has 10<sup>7</sup> seqs, then you expect 10<sup>7</sup> \* 10<sup>-6</sup> = 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

Note:  $P = 1 - e^{-E}$ 

## Lightning does strike twice!



- Roy Sullivan, a former park ranger from Virgina, was struck by lightning 7 times
  - 1942 (lost big-toe nail)
  - 1969 (lost eyebrows)
  - 1970 (left shoulder seared)
  - 1972 (hair set on fire)
  - 1973 (hair set on fire & legs seared)
  - 1976 (ankle injured)

CS2220, AY17/18

- 1977 (chest & stomach burned)



• September 1983, he committed suicide

Cartoon: Ron Hipschman Data: David Hand

# Effect of seq compositional bias

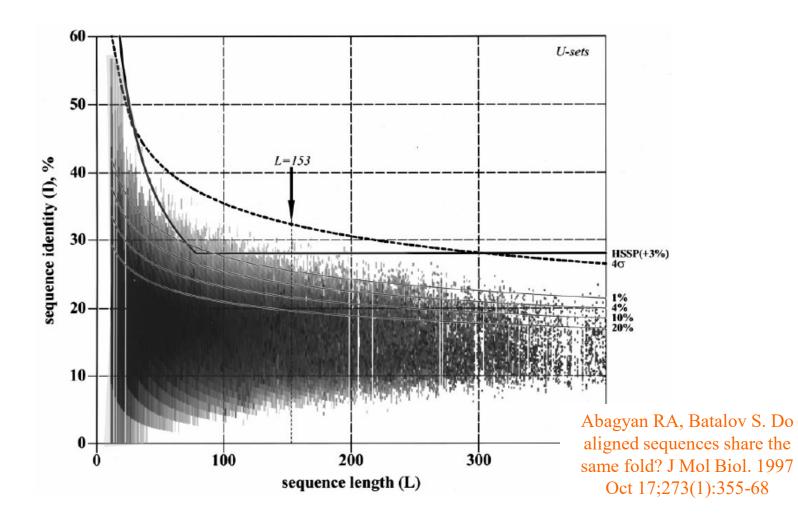
- One fourth of all residues in protein seqs occur in regions with biased amino acid composition
- Alignment of two such regions achieves high score purely due to segment composition
- ⇒ While it is worth noting that two proteins contain similar low complexity regions, they are best excluded when constructing alignments
- E.g., by default, BLAST employs the SEG algo to filter low complexity regions from proteins before executing a search

Source: NCBI



28

#### Effect of sequence length



#### Examples of invalid function assignment: IMP dehydrogenases (IMPDH)



29

		18 ent	ries were found				
ID	Organism	PIR	Swiss-Prot/TrEMBL	RefSeq/GenPept			
1F00181857	Methanococcus jannaschii	<u>E64381</u> conserved hypothetical protein MJ0653	<u>Y653_METJA</u> Hypothetical protein MJ0653	g <u>1592300</u> inosine-5'-monophosphate dehydrogenase (guaB) <u>NP_247637</u> inosine-5'-monophosphate dehydrogenase (guaB)			
F00187788	Archaeoglobus fulgidus	G69355 MJ0653 homolog AF0847 ALT_NAMES: inosine-monophosphate dehydrogenase (guaB-1) homolog [misnomer]	O29411 INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-1)	<u>g2649754</u> inosine monophosphate dehydrogenase (guaB-1) <u>NP_069681</u> inosine monophosphate dehydrogenase (guaB-1)			
<u>F00188267</u>	Archaeoglobus fulgidus	F69514 yhcV homolog 2 ALT_NAMES: inosine-monophosphate dehydrogenase (guaB-2) homolog [misnomer]	<u>028162</u> INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-2)	<u>g2648410</u> inosine monophosphate dehydrogenase (guaB-2) <u>NP_070943</u> inosine monophosphate dehydrogenase (guaB-2)			
IF00188697			nydrogenase mis	tve			
1200197776	Thermo	• •	s remaining in so latabases	nonophosphate 1 protein			
	Thermo Methanothermobacter thermautotrophicus	• •	S remaining in So atabases 027294 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN V	nonophosphate			
1F00414709	Methanothermobacter	ALT_NAMES: inosine-monophosphate	O27294 INOSINE-5-MONOPHOSPHATE	nonophosphate d protein onophosphate dehydrogenase related protein V <u>NP_276354</u> inosine-5'-monophosphate			
VF00414709 VF00414811 VF00414837	Methanothermobacter thermautotrophicus Methanothermobacter	Deputie of the second s	<b>atabases</b> O27294       INOSINE-5'-MONOPHOSPHATE         DEHYDROGENASE RELATED PROTEIN V         O26229       INOSINE-5'-MONOPHOSPHATE	nonophosphate d protein onophosphate dehydrogenase related protein V <u>NP_276354</u> inosine-5'-monophosphate dehydrogenase related protein V <u>g2621166</u> inosine-5'-monophosphate dehydrogenase related protein VII <u>NP_275269</u> inosine-5'-monophosphate			

CS2220, AY17/18

#### Copyright 2017 © Wong Limsoon



#### **IMPDH domain structure**

	949	PCM00487: PDOC003	91,IMP dehydrogen <i>a</i> se / GMP re	eductase signature				
	and the second sec	PF00478: IMP dehydrogenase / GMP reductase C terminus						
	*****	PF00571: CBS domain						
	00-00	PF01381: Helix-turn-h	elix					
	e fan fan fan fan fan	PF01574: IMP dehydr	ogenase / GMP reductase N term	inus				
	գինիսիներ	PF02195: ParB-like nu	clease domain					
A31997 (SF000130)	0.000		tototototo statetatatatat	514				
. ,								
E70218 (SF000131)	-			404				
(01000101)								
E64381			194	IMPDH Misnomer in <i>Methanococcus jannaschii</i>				
(SF004696)	0\$=\$0\$=\$0\$	tato statotatato stat	<b>N</b>					
			189					
G69355								
G69355 (SF004696)		tak statetatetak						
	ojetojetek	tak stalatalatak	402	IMPDH Mispomors in Arabaaaalobus fulaidus				
(SF004696)			183	IMPDH Misnomers in Archaeoglobus fulgidus				
(SF004696) F69514			<b>183</b>	IMPDH Misnomers in Archaeoglobus fulgidus				

- 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

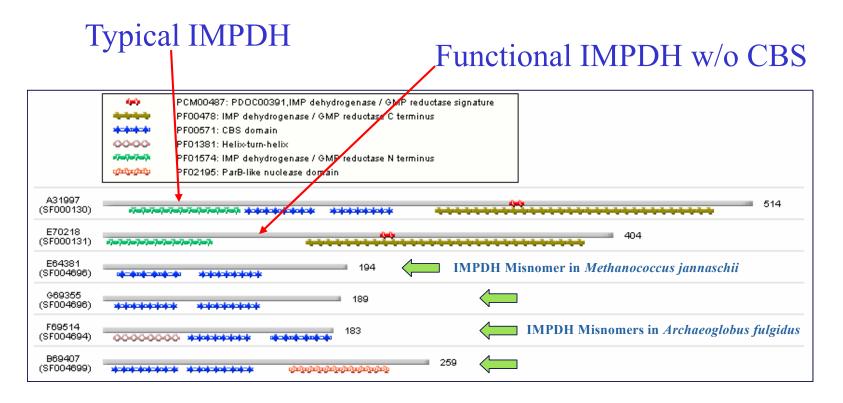


#### Root of invalid transitive assignment \_\_\_\_

B⇒	□ <u>H70468</u>	<u>SF001258</u>		phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) [similarity]		Aquifex aeolicus	Prok/other	594.3	4.8e-26	205	39.086	197		
	□ <u>\$76963</u>	<u>SF001258</u>	<u>039935</u>	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) / phosphoribosyl-ATP pyro (EC 3.6.1.31) [similarity]		Synechocystis sp.	Prok/gram-	557.0	5.7e-24	230	39.175	194		
	T35073	SF029243	<u>005738</u>	probable phosphoribosyl-AMP cyclohydrolase		Streptomyces coelicolor	Prok/gram+	399.3	3.5e-15	128	42.157	102		
	□ <u>\$53349</u>	<u>SF001257</u>	001188	phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) / histidinol dehydrogenase (EC 1.1.1.23)		Saccharomyces cerevisiae	Euk/fungi	384.1	2.5e-14	799	31.863	204		
A⊫>	□ <u>E69493</u>	SF029243	<u>005738</u>	phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) [similarity]		Archaeoglobus fulgidus	Archae	396.8	4.8e-15	108	47. <i>7</i> 78	90		
C	□ <u>G64337</u>	SF006833	<u>030827</u>	phosphoribosyl-ATP pyrophosphatas 3.6.1_31) [similarity]	se (EC	Methanococcus jannaschii	Archae	246.9	1.1e-06	95	36.842	95	, .	
	D81178	<u>SF006833</u>	<u>101491</u>	phosphoribosyl-ATP pyrophosphatas 3.6.1.31) NMB0603 [similarity]	se (EC	Neisseria meninoitidis	Prok/oram-		2 ńe-Nń		35 227	88		
	G81925 SF006833 101491 hosphoribosyl-ATP pyrophosphat 3.6.1.31) NMA0807 [similarity]		$A \rightarrow B \rightarrow C \implies A \rightarrow C$								-			
	□ <u>\$51513</u>	<u>SF001257</u>		phosphoribosyl-AMP cyclohydrola 3.5.4.19) / phosphoribosyl-ATP py (EC 3.6.1.31) / histidinol dehydrog 1.1.1.23)		B (SF001258)								-
Mis-assignment			A (SF029243) C (SF006833)											
of function			No IMPDH domain											
CS2220, AY17/18							Co	pyrigł	nt 2	017 ©	) Wo	ong Lin	nsoon	

# Emerging pattern





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

# Application of sequence comparison: Active site / domain discovery



### Discover active site and/or domai



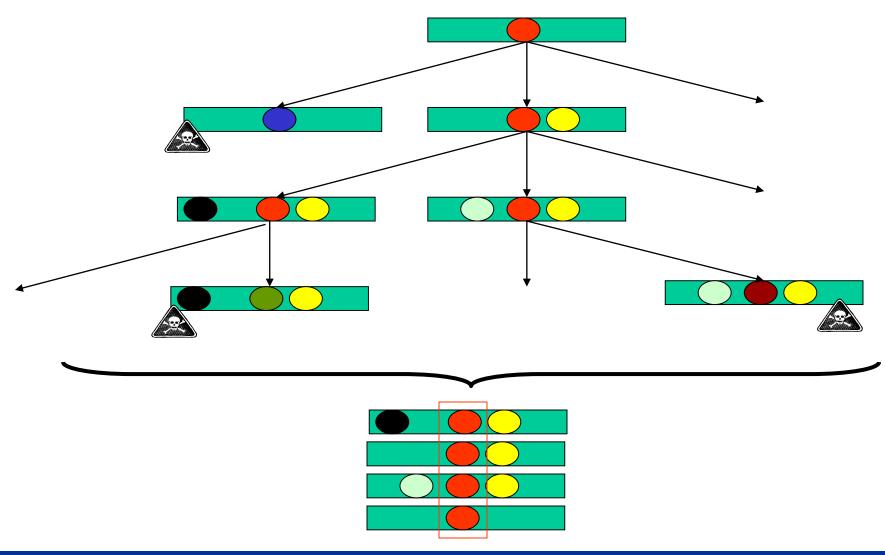
- How to discover the active site and/or domain of a function in the first place?
  - Multiple alignment of homologous seqs
  - Determine conserved positions
  - $\Rightarrow$  Emerging patterns relative to background
  - $\Rightarrow$  Candidate active sites and/or domains
- Easier if sequences of distance homologs are used

Exercise #2: Why?

#### In the course of evolution...



35



#### CS2220, AY17/18

Copyright 2017 © Wong Limsoon



# Multiple alignment of PTPs

gi 126467	FHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTGTFVVIDAMLD
gi 2499753	FHFTGWPDHGVPYHATGLLSFIRRVKLSNPPSAGPIVVHCSAGAGRTGCYIVIDIMLD
gi 462550	YHYTQWPDMGVPEYALPVLTFVRRSSAARMPETGPVLVHCSAGVGRTGTYIVIDSMLQ
gi 2499751	FHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGTFIAIDRLIY
gi 1709906	FQFTAWPDHGVPEHPTPFLAFLRRVKTCNPPDAGPMVVHCSAGVGRTGCFIVIDAMLE
gi 126471	LHFTSWPDFGVPFTPIGMLKFLKKVKTLNPVHAGPIVVHCSAGVGRTGTFIVIDAMMA
gi 548626	FHFTGWPDHGVPYHATGLLSFIRRVKLSNPPSAGPIVVHCSAGAGRTGCYIVIDIMLD
gi 131570	FHFTGWPDHGVPYHATGLLGFVRQVKSKSPPNAGPLVVHCSAGAGRTGCFIVIDIMLD
gi 2144715	FHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGTFIAIDRLIY
	··* *** *** · * ··* ··* ··* ··* ··

- Notice the PTPs agree with each other on some positions more than other positions
- These positions are more impt wrt PTPs
- Else they wouldn't be conserved by evolution
- $\Rightarrow$  They are candidate active sites

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



### What if there is no useful seq homologicational University of Singapore

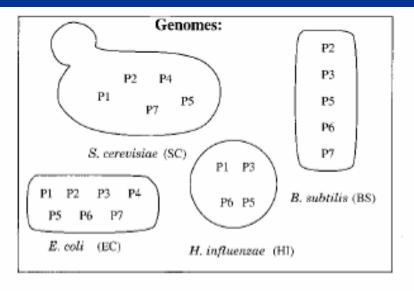
- Guilt by other types of association!
  - Domain modeling (e.g., HMMPFAM)
  - ✓ Similarity of phylogenetic profiles
  - ✓ Similarity of dissimilarities (e.g., SVM-PAIRWISE)
  - Similarity of subcellular co-localization & other physico-chemico properties(e.g., PROTFUN)
  - Similarity of gene expression profiles
  - ✓ Similarity of protein-protein interaction partners
  - Fusion of multiple types of info

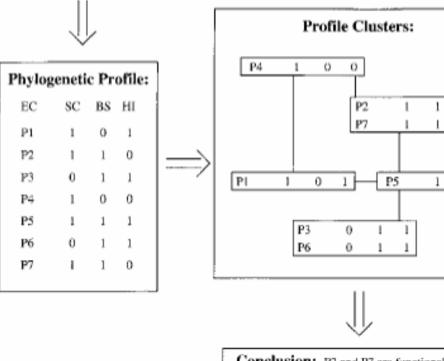
. . .

#### Phylogenetic profiling Pellegrini et al., PNAS, 96:4285--4288, 1999



- Genes (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







40

# Phylogenetic profiling: How it works

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

0

0

1 L



Copyright 2017 © Wong Limsoon

## 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

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

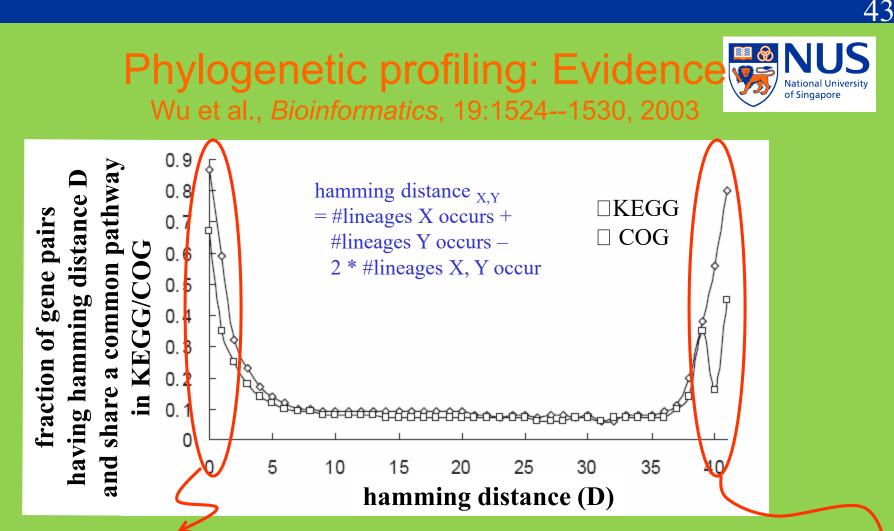
where

No. of ways to distribute 
$$z$$
  
co-occurrences over  $N$   
lineage's  
No. of ways to distribute  
 $W = \binom{N-z}{x-z} * \binom{N-x}{y-z}$   
No. of ways to distribute  
the remaining  $x - z$  and  $y - z$   
occurrences over the remaining  
 $N - z$  lineage's  
No. of ways to distribute  
 $W = \binom{N}{x} * \binom{N}{y}$   
No. of ways of  
distributing X and Y  
over N lineage's  
without restriction

# 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



 Proteins having low hamming distance (thus highly similar phylogenetic profiles) tend to share common pathways Exercise #3: Why do proteins having high hamming distance also have this behaviour?

Copyright 2017 © Wong Limsoon



### Guilt by association of dissimilarities



Differences of "unknown" to other fruits are same as "apple" to other fruits



	Drange <sub>1</sub>	ana <sub>1</sub>	
Apple <sub>1</sub>	Color = red vs orange	Color = red vs yellow	
	Skin = smooth vs rough	Skin = smooth vs smooth	
	Size = small vs small	Size = small vs small	
	Shape = round vs round	Shape = round vs oblong	
Orange <sub>2</sub>	Color = orange vs orange	Color = orange vs yellow	
	Skin = rough vs rough	Skin = rough vs smooth	
	Size = small vs small	Size = small vs small	
	Shape = round vs round	Shape = round vs oblong	
Unknown <sub>1</sub>	Color = red vs orange	Color = red vs yellow	
5	Skin = smooth vs rough	Skin = smooth vs smooth	
	Size = small vs small	Size = small vs small	
<b>– –</b>	Shape = round vs round	Shape = round vs oblong	



45

#### **SVM-Pairwise framework**

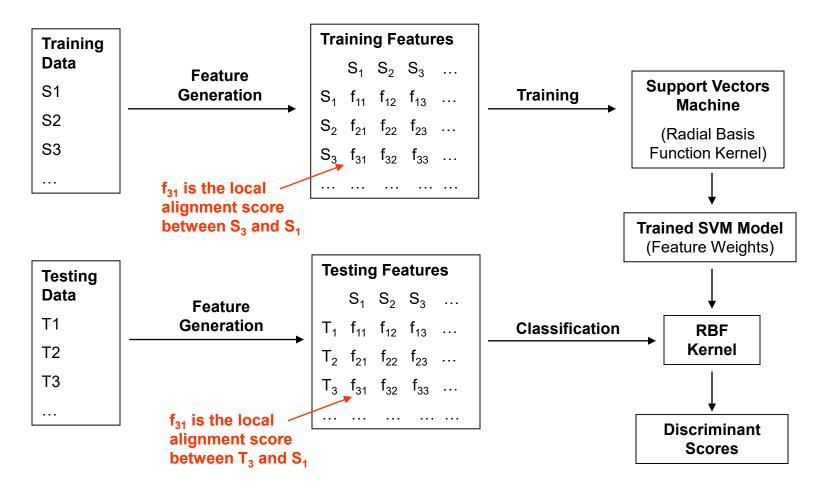


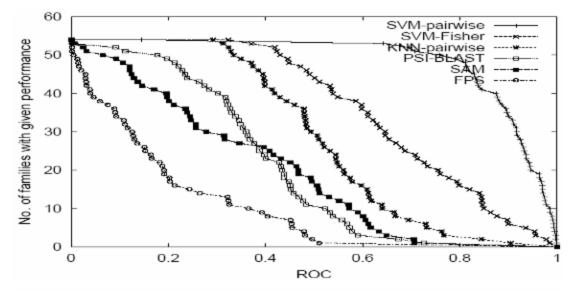
Image credit: Kenny Chua

### **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)**

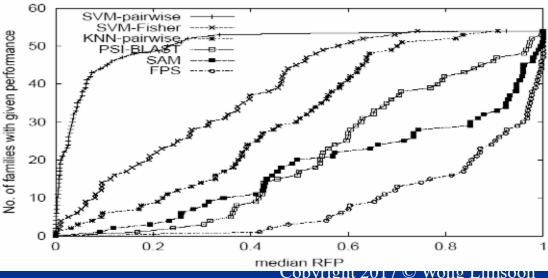
CS2220, AY17/18

 The fraction of negative test examples with a score better or equals to the median of the scores of positive test examples.

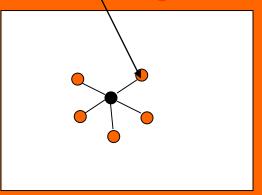


46

of Singapore

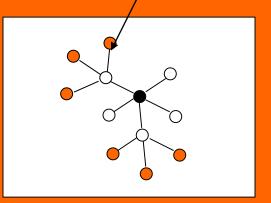


#### Level-1 neighbour



### Protein function prediction from protein interactions







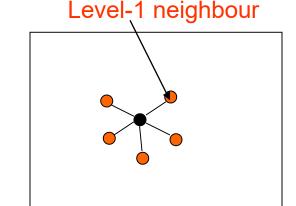
## Functional association thru interaction singapore

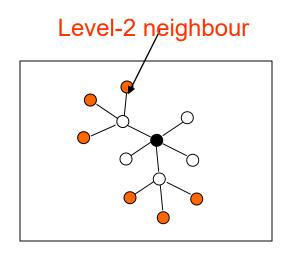
#### • Direct functional association:

- 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

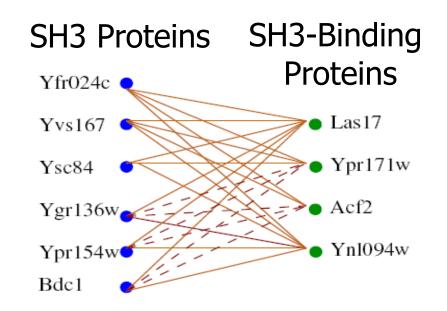




# An illustrative case of indirect functional association?

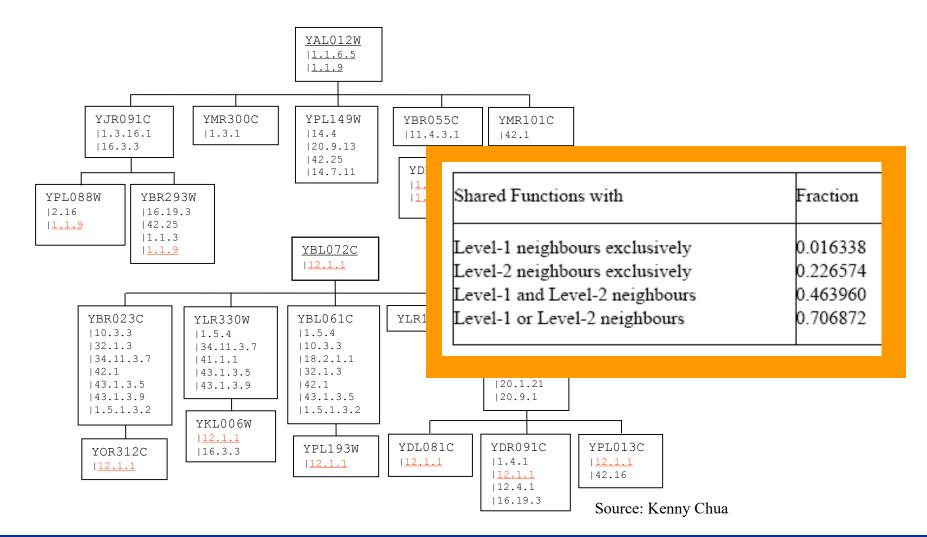


49



- 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



#### CS2220, AY17/18

#### Copyright 2017 © Wong Limsoon

50

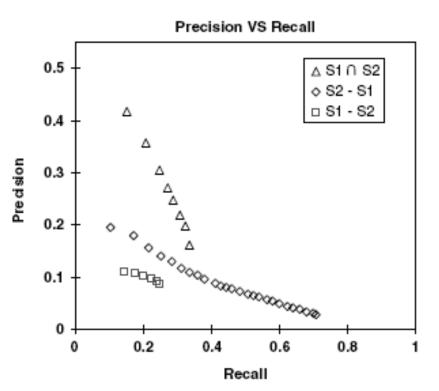


- Remove overlaps in level-1 and level-2 neighbours to study predictive power of "level-1 only" and "level-2 only" neighbours
- Sensitivity vs Precision analysis

$$PR = \frac{\sum_{i}^{K} k_{i}}{\sum_{i}^{K} m_{i}} \quad SN = \frac{\sum_{i}^{K} k_{i}}{\sum_{i}^{K} n_{i}}$$

- n<sub>i</sub> is no. of fn of protein i
- m<sub>i</sub> is no. of fn predicted for protein i
- k<sub>i</sub> is no. of fn predicted correctly for protein i

CS2220, AY17/18



- ⇒ "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<sub>k</sub> 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

 $\Rightarrow$  Similarity can be defined as

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

CS2220, AY17/18

S(u,v) = 1 - D(u,v) =

2X

2X + (Y + Z)



53

### 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<sub>k</sub> is the set of interacting partners of k
- Greater weight given to similarity
- $\Rightarrow$  Rewriting this as

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

# Correlation w/ functional similarit

Correlation betw functional similarity & estimates

Neighbours	CD-Distance	FS-Weight	]
$egin{array}{c} S_1 \ S_2 \ S_1 \cup S_2 \end{array}$	0.471810 0.224705 0.224581	0.498745 0.298843 0.29629	(

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



55

### 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<sub>i</sub> is reliability of expt source i,
- E<sub>u,v</sub> 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

#### CS2220, AY17/18

#### Exercise #4



56

Can you think of things a biologist can do to assess the overall reliability of a PPI screening assay / source?

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

Functional similarity estimate: FS-weighted measure with reliabilit



 Take reliability into consideration when computing FS-weighted measure:

$$S_{R}(u,v) = \frac{2\sum_{w \in (N_{u} \cap N_{v})} r_{u,w}r_{v,w}}{\left(\sum_{w \in N_{u} - N_{v}} r_{u,w} + \sum_{w \in (N_{u} \cap N_{v})} r_{u,w}(1 - r_{v,w})\right) + 2\sum_{w \in (N_{u} \cap N_{v})} r_{u,w}r_{v,w}} \times \frac{2\sum_{w \in (N_{u} \cap N_{v})} r_{u,w}r_{v,w}}{\left(\sum_{w \in N_{v} - N_{u}} r_{v,w} + \sum_{w \in (N_{u} \cap N_{v})} r_{v,w}(1 - r_{u,w})\right) + 2\sum_{w \in (N_{u} \cap N_{v})} r_{v,w}r_{v,w}}}$$

- N<sub>k</sub> is the set of interacting partners of k
- r<sub>u.w</sub> is reliability weight of interaction betw u and v

 $\Rightarrow$  **Rewriting** 

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

#### Integrating reliabilities

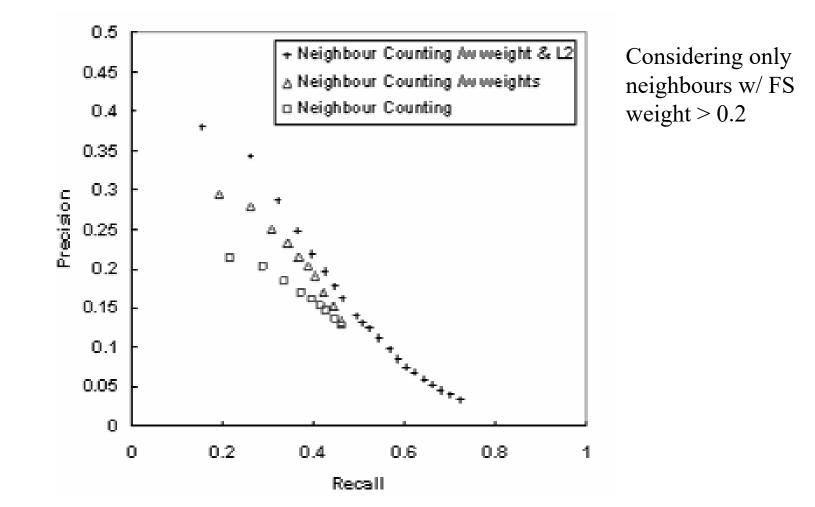


58

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

Neighbours	CD-Distance	FS-Weight	FS-Weight R
$egin{array}{c} S_1 \ S_2 \ S_1 \cup S_2 \end{array}$	0.224705	0.298843	0.532596 0.375317 0.363025

# Improvement to prediction power by majority voting

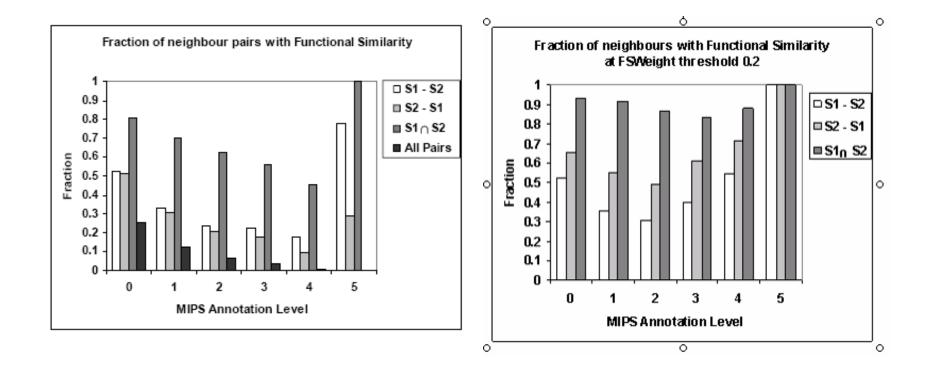


59



60

# Improvement to over-rep of functions in neighbours



## Use L1 & L2 neighbours for predict of Singapore

#### FS-weighted Average

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

- *r<sub>int</sub>* is fraction of all interaction pairs sharing function
- $\lambda$  is weight of contribution of background freq
- $\delta(\mathbf{k}, \mathbf{x}) = 1$  if k has function x, 0 otherwise
- N<sub>k</sub> is the set of interacting partners of k
- $\pi_x$  is freq of function x in the dataset
- Z is sum of all weights

CS2220, AY17/18

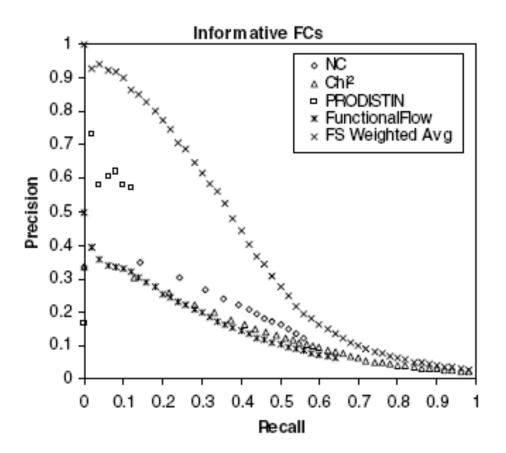
$$Z = 1 + \sum_{v \in N_u} \left( S_{TR}(u, v) + \sum_{w \in N_v} S_{TR}(u, w) \right)$$

1

6

## Performance of FS-weighted average Singapore

 LOOCV comparison with Neighbour Counting, Chi-Square, PRODISTIN

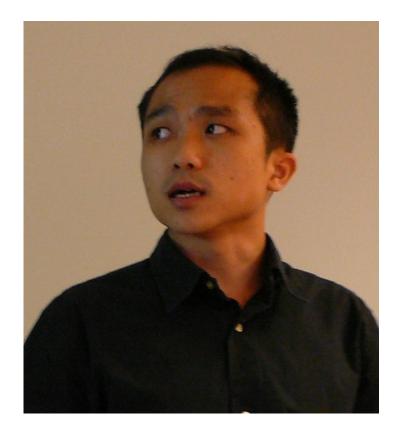


62

## About the inventor: Chua Hon Nia

#### Chua Hon Nian

- PhD, NUS, 2008
- Postdoc at Harvard
  & Univ of Toronto
- 49<sup>th</sup> hottest paper in Computer Science published in 2006
- Winner, DREAM2
   challenge PPI
   subnetwork, 2007
- Now Data Scientist at Data Robot



Application of sequence comparison: Key mutation site discovery



#### Identifying key mutation sites K.L.Lim et al., *JBC*, 273:28986--28993, 1998



#### Sequence from a typical PTP domain D2

>gi|00000|PTPA-D2 EEEFKKLTSIKIQNDKMRTGNLPANMKKNRVLQIIPYEFNRVIIPVKRGEENTDYVNASF IDGYRQKDSYIASQGPLLHTIEDFWRMIWEWKSCSIVMLTELEERGQEKCAQYWPSDGLV SYGDITVELKKEEECESYTVRDLLVTNTRENKSRQIRQFHFHGWPEVGIPSDGKGMISII AAVQKQQQQSGNHPITVHCSAGAGRTGTFCALSTVLERVKAEGILDVFQTVKSLRLQRPH MVQTLEQYEFCYKVVQEYIDAFSDYANFK

- Some PTPs have 2 PTP domains
- PTP domain D1 has much more activity than PTP domain D2
- Why? And how do you figure that out?

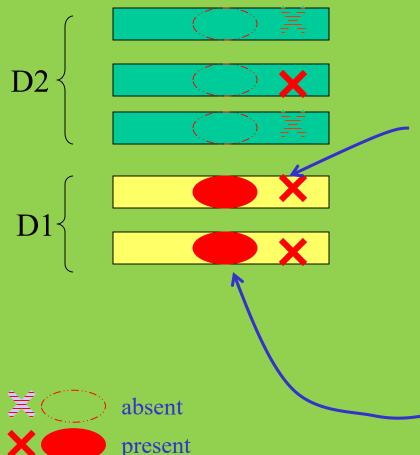
### Emerging patterns of PTP D1 vs De National University of Singapore

- Collect example PTP D1 sequences
- Collect example PTP D2 sequences
- Make multiple alignment A1 of PTP D1
- Make multiple alignment A2 of PTP D2
- Are there positions conserved in A1 that are violated in A2?
  - These are candidate mutations that cause PTP activity to weaken
- Confirm by wet experiments

66

### Emerging patterns of PTP D1 vs





This site is consistently conserved in D1, but is not consistently missing in D2 ⇒ it is not an EP ⇒ not a likely cause of D2's loss of function

#### Exercise #5: Why?

This site is consistently conserved in D1, but is consistently missing in D2 ⇒ it is an EP ⇒ possible cause of D2's loss of function





gi|00000|P gi|126467| gi|2499753 gi|462550| gi|2499751 gi|1709906 gi|126471| gi|548626| gi|131570| gi|2144715

2 2 2 22 OFHFHGWPEVGIPSDGKGMISIIAAVOKOOOO-SGNHPITVHCSAGAGRTGTFCALSTVL QFHFTSWPDFGVPFTPIGMLKFLKKVKACNP--QYAGAIVVHCSAGVGRTGTFVVIDAML OFHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIML OYHYTOWPDMGVPEYALPVLTFVRRSSAARM--PETGPVLVHCSAGVGRTGTYIVIDSML OF HFTSWPDHGVPDTTDLLINFRYLVRDYMKOSPPESPILVHCSAGVGRTGTFIAIDRLI QFQFTAWPDHGVPEHPTPFLAFLRRVKTCNP--PDAGPMVVHCSAGVGRTGCFIVIDAML QLHFTSWPDFGVPFTPIGMLKFLKKVKTLNP--VHAGPIVVHCSAGVGRTGTFIVIDAMM OFHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIML OFHFTGWPDHGVPYHATGLLGFVROVKSKSP--PNAGPLVVHCSAGAGRTGCFIVIDIML QFHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGTFIAIDRLI \*\*\*\*\* \*\*\*\* \* \*\*. \*.\*

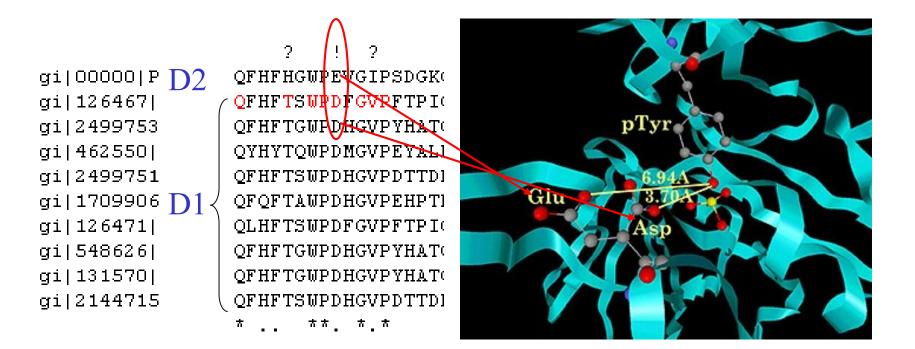
- Positions marked by "!" and "?" are likely places responsible for reduced PTP activity
  - All PTP D1 agree on them
  - All PTP D2 disagree on them

68



69

#### Key mutation site: PTP D1 vs D2



 Positions marked by "!" are even more likely as 3D modeling predicts they induce large distortion to structure





- What wet experiments are needed to confirm the prediction?
  - Mutate  $E \rightarrow D$  in D2 and see if there is gain in PTP activity
  - Mutate D  $\rightarrow$  E in D1 and see if there is loss in PTP activity

Exercise: Why do you need this 2-way expt?



# About the inventor: Prasanna Kolat

#### Prasanna Kolatkar

- Research Fellow,
   BIC, NUS, 1997 1999
- Currently Senior
   Scientist at Qatar
   Biomedical
   Research Institute



#### Concluding remarks



#### What have we learned?



- General methodologies & applications
  - Guilt by association for protein function inference
  - Invariants for active site discovery
  - Emerging patterns for mutation site discovery
- Important tactics
  - Genome phylogenetic profiling
  - SVM-Pairwise
  - Protein-protein interactions

#### Any question?



#### 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
- K.L.Lim et al. "Interconversion of kinetic identities of the tandem catalytic domains of receptor-like protein tyrosine phosphatase PTP-alpha by two point mutations is synergist and substrate dependent", *JBC*, 273:28986--28993, 1998
- 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

#### 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

#### References



- H.N. Chua, W.-K. Sung. <u>A better gap penalty for pairwise SVM</u>. Proc. APBC05, pages 11-20
- Hon Nian Chua, Wing Kin Sung, Limsoon Wong. Exploiting Indirect Neighbours and Topological Weight to Predict Protein Function from Protein-Protein Interactions. *Bioinformatics*, 22:1623-1630, 2006.
- T. Jaakkola, M. Diekhans, and D. Haussler. A discriminative framework for detecting remote homologies. *JCB*, 7(1-2):95-114, 2000
- T. Hawkins and D. Kihara. Function prediction of uncharacterized proteins. *JBCB*, 5(1):1-30, 2007