

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 Lecture 6: Sequence Homology Interpretation

Limsoon Wong



Plan

- 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

Copyright 2011 © Limsoon Wong

Very Brief Recap of Sequence Comparison/Alignment

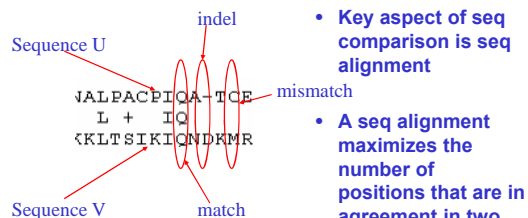


Motivations for Sequence Comparison

- DNA is blue print for living organisms
 - ⇒ 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

Copyright 2011 © Limsoon Wong

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

Copyright 2011 © Limsoon Wong

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

```

Amicyanin      40      60      80      100
                DNNNPPPPGSGGAAAGPNNHQAASFEYAGTWNKTFRRKSHVVE
Ascorbate Oxidase 120GTPFADGASISQAINCGTTTNTFWYTFYTFNAGLQNDKQDQSL
                70      90      110      130
  
```

No obvious match between
Amicyanin and Ascorbate Oxidase

Copyright 2011 © Limsoon Wong

Sequence Alignment: Good Example

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

☐ >g113476732|ref|NP_108301.1| unknown protein [Mesorhizobium loti]
 g114027493|db|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 MKPORLASIALAIIFLPMAVAHAATIEITMENLVISPEVSAKVODTIRVNDVFAHT 60
 MK G L ++ MA PA AATIE++ LV SP V AKVODTI VVN DV AHT
 Sbjct: 1 MKAGALIELSWLAALALMAFAAAATIEVTIDKLVSFATVEAKVODTIEVNDVFAHT 60

good match between
 Amicyanin and unknown M. loti protein

Copyright 2011 © Limsoon Wong

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

```

g1|126467| FHTSWPDPGVFPPTIGMLKFLKKVKACNP--QYAGAIIVHCSAGVGRGTGTFVVIDAMLD
g1|2499753| FHTGWPDHGVFPYHATGLLSFIRKVKLENP--PSAGPIVHCSAGAGRTGCTYVIDIMLD
g1|462550| YHYTQWPDHGVPEYALPVLTFRSSAARM--PETGPIVHCSAGVGRGTGTFVVIDSMLQ
g1|2499751| FHTSWPDPHGVPTDILLINFRYLVRDYKOSPPESPILVHCSAGVGRGTGTFVVIDRLIY
g1|1709906| FQFTAWPDHGVPEHPTFLAFLERVKTNP--PDAGPIVHCSAGVGRGTGTFVVIDAHLE
g1|126471| LHFTSWPDPGVFPPTIGMLKFLKKVKLENP--VHAGPIVHCSAGVGRGTGTFVVIDAHMA
g1|548626| FHTGWPDHGVFPYHATGLLSFIRKVKLENP--PSAGPIVHCSAGAGRTGCTYVIDIMLD
g1|1315701| FHTGWPDHGVFPYHATGLLSFIRKVKLENP--PSAGPIVHCSAGAGRTGCTYVIDIMLD
g1|2144715| FHTSWPDPHGVPTDILLINFRYLVRDYKOSPPESPILVHCSAGVGRGTGTFVVIDRLIY
..* ** * ..
  
```

Conserved sites

Copyright 2011 © Limsoon Wong

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
- Protein performs a wide variety of activities in the cell



Copyright 2011 © Limsoon Wong

Function Assignment to Protein Sequence

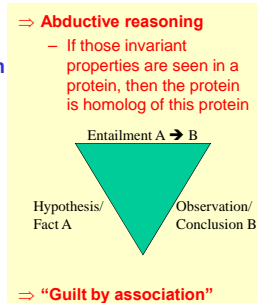
SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR
 YVNILPYDHSRVHLTPVEGVDPDSYINASFINGYQEKNFIAAQGPKEETVNDFFWRMIWE
 QNTATIVMVTNLKERKECKCAQYWPDPQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD
 VTNRKPKRLITQFHTSWPDPGVFPPTIGMLKFLKKVKACNPQYAGAIIVHCSAGVGRGT
 TFVVIDAMLDMHMERKVDVYGFVSRIARQCRCMVQTDMDQYVFYIYQALLEHYLYGDTLE
 VT

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

Copyright 2011 © Limsoon Wong

Invariant and Abductive Reasoning

- Function is determined by 3D struct of protein & environment protein is in
- Constraints imposed by 3D struct & environment give rise to "invariant" properties observed in proteins having the ancestor with that function



Copyright 2011 © Limsoon Wong

13

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

Copyright 2011 © Limsoon Wong

14

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 *amylase* and domain 1 of *acetylcholinesterase*

No obvious match between *Ameyase* and *acetylcholinesterase*

Good Sequence Alignment

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

Alignment by FASTA of the sequences of *amylase* and domain 1 of *acetylcholinesterase*

Ameyase and acetylcholinesterase

Assign to T same function as homologs

Confirm with suitable wet experiments

Discard this function as a candidate

Copyright 2011 © Limsoon Wong

15

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

find from db seqs with short perfect matches to query seq

find seqs with good flanking alignment

Exercise: Why do we need this step?

Copyright 2011 © Limsoon Wong

16

Homologs obtained by BLAST

Sequences producing significant alignments:	Score (bits)	E Value
gi14192728 gb AAK54109.1 AF331081.1 protein tyrosine phosph...	62.1	e-177
gi1264471 gb F184331 F184331.HUMAN Protein-tyrosine phosphatase...	62.1	e-177
gi14506303 ref NP_002827.1 protein tyrosine phosphatase, x...	62.0	e-176
gi12372241 ref U17813004 protein Tyr phosphatase	62.0	e-176
gi1184503 ref NP_043030.1 protein tyrosine phosphatase, ...	61.9	e-176
gi1320671 emb CA37447.1 tyrosine phosphatase precursor [Ho...	61.9	e-176
gi1285113 p P10C1285 protein-tyrosine-phosphatase [EC 3.1.1...	61.9	e-176
gi16981446 ref NP_036895.1 protein tyrosine phosphatase, x...	61.8	e-176
gi12029414 emb U17501A Chain A, Receptor Protein Tyrosine Ph...	61.8	e-174
gi132131 emb CA32666.1 protein-tyrosine phosphatase (Homo...	61.8	e-174
gi14506303 ref NP_002827.1 protein tyrosine phosphatase, x...	60.5	e-172
gi1679587 ref NP_033006.1 protein tyrosine phosphatase, x...	60.4	e-172
gi1483322 gb AAK17990.1 protein tyrosine phosphatase alpha	59.9	e-170

- Thus our example sequence could be a protein tyrosine phosphatase α (PTP α)

Copyright 2011 © Limsoon Wong

17

Example Alignment with PTP α

Score = 632 bits (1629), Expect = e-130
Identities = 294/302 (97%), Positives = 294/302 (97%)

```

Query: 1  SPSTNRKPPPLPVLKLEEE INSRMADONLFPREFNALPAC IQATCEAAS 60
Subject: 202 SPSTNRKPPPLPVLKLEEE INSRMADONLFPREFNALPAC IQATCEAAS 261

Query: 61  YVNLPTNRKPPPLPVLKLEEE INSRMADONLFPREFNALPAC IQATCEAAS 120
Subject: 262 YVNLPTNRKPPPLPVLKLEEE INSRMADONLFPREFNALPAC IQATCEAAS 321

Query: 121 QNTATITVMTNLKEREKCAQVDFDQCTVGVVRSYEDVTVLVDITVKKPCIQVSD 180
Subject: 322 QNTATITVMTNLKEREKCAQVDFDQCTVGVVRSYEDVTVLVDITVKKPCIQVSD 381

Query: 161 VTRNRKPPPLPVLKLEEE INSRMADONLFPREFNALPAC IQATCEAAS 240
Subject: 362 VTRNRKPPPLPVLKLEEE INSRMADONLFPREFNALPAC IQATCEAAS 441

Query: 241 TPVVDIAMDNDIERRVDVDFVSRIDARQDQVQVDMQVFIQALLSIVLPDTELE 300
Subject: 442 TPVVDIAMDNDIERRVDVDFVSRIDARQDQVQVDMQVFIQALLSIVLPDTELE 501
  
```

Copyright 2011 © Limsoon Wong

18

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

Copyright 2011 © Limsoon Wong

19

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%

Copyright 2011 © Limsoon Wong

20

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 $10^7 \times 10^{-6} = 10$ seqs in it that give an equally good alignment
- ⇒ Need to correct for database size if your seq comparison prog does not do that!

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

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

Copyright 2011 © Limsoon Wong

21

2.2 Review of BLAST statistics

In order to derive a way to combine the E-values reported by PIR, BLAST, etc. different counts are enough. It is useful to review how these E-values are actually calculated. If a query sequence with q residues is aligned against a random target sequence of n residues using the Smith-Waterman (Smith and Waterman, 1970) local alignment algorithm, the best alignment score Z can be considered a random variable. In the absence of gaps, the random variable has been proven (Karlin and Altschul, 1990, 1991; Karlin and Ewens, 1993) to follow an extreme value distribution

$$P\left(\sum_{i=1}^n X_i \leq Z\right) = \exp(-\lambda Z m^{-1/2}) \quad (1)$$

if the query and target sequences are sufficiently long. Here, λ and m are two parameters that depend on the query sequence, the target sequence and the random sequence used to derive the random target sequences.

From the literature of gaps, Karlin and Altschul (1990) and Altschul et al. (1990) have shown that the best alignment score Z can be considered a random variable. In the absence of gaps, the random variable has been proven (Karlin and Altschul, 1990, 1991; Karlin and Ewens, 1993) to follow an extreme value distribution

$$P\left(\sum_{i=1}^n X_i \leq Z\right) = \exp(-\lambda Z m^{-1/2}) \quad (2)$$

BLAST and FASTA (Pearson, 1990) use the expected number of random hits of score Z or higher, instead of P -values. These two quantities are different since the P -value is the probability of having at least one hit at the given significance level, while the E -value quantifies the expected number of hits. The general relationship between these two quantities for random statistics with exponential tails is the Pomeroy formula

$$P = 1 - \exp(-E) \quad (3)$$

By comparing with Equation (2) we find

$$P_{\text{BLAST}} = 1 - \exp(-P) = P \approx E \quad (4)$$

It is important to note that this is true for one single sequence comparison, where we select the highest score. The convenient property of E -value is that the logarithm is additive for repeating the alignment for every sequence in a database of N random sequences, and thus simply $N \times E$ sequences, just multiply N is a simplification of the expected number of hits for a single sequence by the number of sequences, i.e.

$$P_{\text{BLAST}} = 1 - \exp(-N \times E) \approx N \times E \quad (5)$$

Copyright 2011 © Limsoon Wong

22

Lightning Does Strike Twice!

- Roy Sullivan, a former park ranger from Virginia, 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)
 - 1977 (chest & stomach burned)
- September 1983, he committed suicide

Cartoon: Ron Hipschman
Data: David Hand

Copyright 2011 © Limsoon Wong

23

Effect of Seq Compositional Bias

- One fourth of all residues in protein seqs occur in regions with biased amino acid composition
- Alignments 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

Copyright 2011 © Limsoon Wong

24

Effect of Sequence Length

Abagyan RA, Batalov S. Do aligned sequences share the same fold? J Mol Biol. 1997 Oct 17;273(1):355-68

Copyright 2011 © Limsoon Wong

25

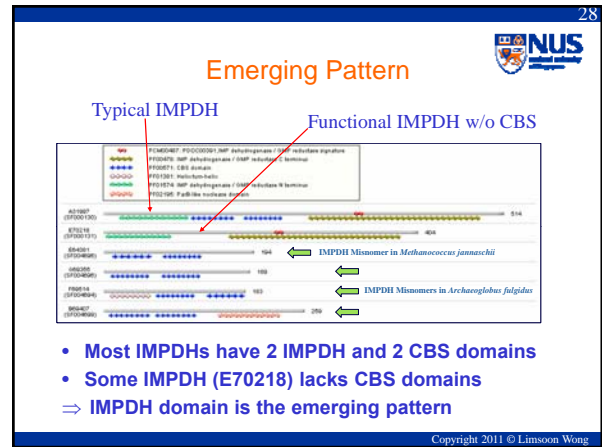
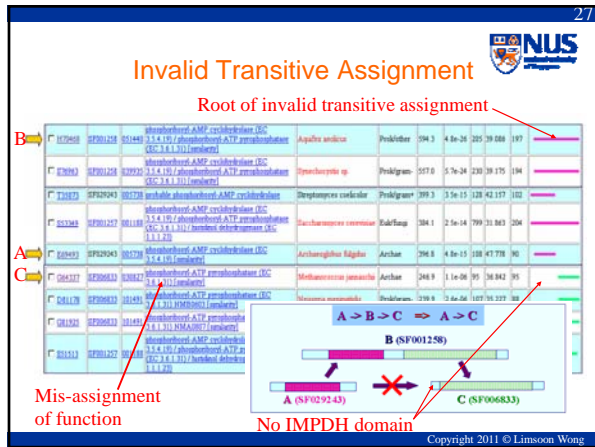
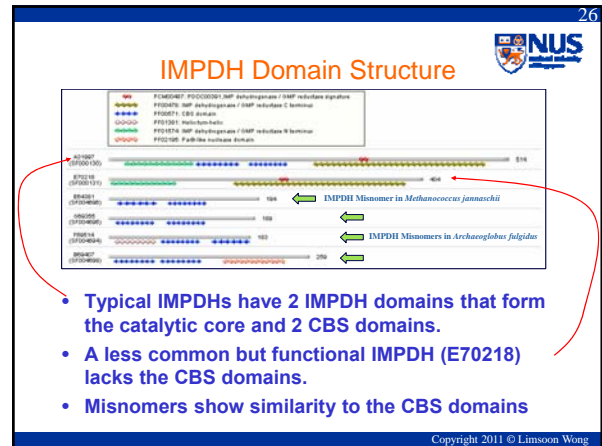
Examples of Invalid Function Assignment:
The IMP Dehydrogenases (IMPDH)

10 entries were found

ID	Organism	PIR	Swiss-Prot/TrEMBL	RefSeq/GenPept
NP0011852	<i>Methanococcus jannaschii</i>	04501 conserved hypothetical protein (M0603)	04501 conserved hypothetical protein (M0603)	04501 conserved hypothetical protein (M0603)
NP0011700	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04
NP0011852	<i>Archaeoglobus fulgidus</i>	04501 M0603 homolog APT04	04501 M0603 homolog APT04	04501 M0603 homolog APT04

A partial list of IMP dehydrogenase misnomers in complete genomes remaining in some public databases

Copyright 2011 © Limsoon Wong



Application of Sequence Comparison: Active Site/Domain Discovery

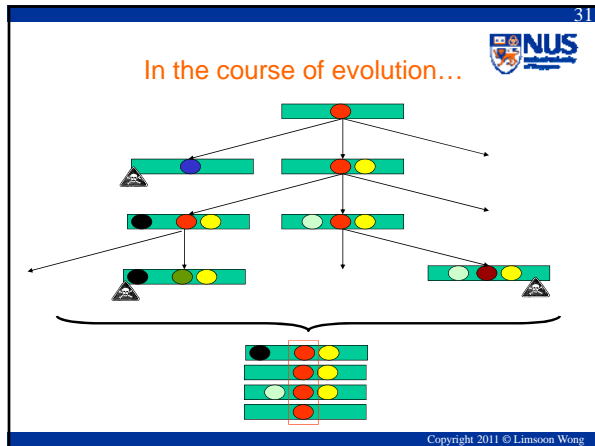
NUS

Discover Active Site and/or Domain

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

Exercise: Why?

Copyright 2011 © Limsoon Wong



32

Multiple Alignment of PTPs

g1|126467| FHFTSWPDFGVPTPIOMLKFLLKKVKACNP--QYAGAIUVHCSAGVGRGTOTFVVIDAMLD
 g1|2499753| FHFTGWDPDHVPYHATGLLSFIRRVKLSNP--PSAGPIUVHCSAGAGTGTCTIVIDIMLD
 g1|462550| YHTTQWPDHVPYHATGLLSFIRRVKLSNP--PSAGPIUVHCSAGAGTGTCTIVIDIMLD
 g1|2499751| FHFTSWDPDHVPDITDLLINFRVLRDYNKQSPPEPILVHCSAGVGRGTOTFVVIDAMLD
 g1|1709906| FQFTAMPDGHVPYHATGLLSFIRRVKLSNP--PDAGPIUVHCSAGVGRGTCTIVIDAMLD
 g1|126471| LHFTSWPDFGVPTPIOMLKFLLKKVKTLNP--VHAGPIUVHCSAGVGRGTCTIVIDAMLD
 g1|548626| FHFTGWDPDHVPYHATGLLSFIRRVKLSNP--PSAGPIUVHCSAGAGTGTCTIVIDIMLD
 g1|1315701| FHFTGWDPDHVPYHATGLLSFIRRVKLSNP--PSAGPIUVHCSAGAGTGTCTIVIDIMLD
 g1|2144715| FHFTSWDPDHVPDITDLLINFRVLRDYNKQSPPEPILVHCSAGVGRGTOTFVVIDAMLD

- Notice the PTPs agree with each other on some positions more than other positions
- These positions are more imp't wrt PTPs
- Else they wouldn't be conserved by evolution

⇒ They are candidate active sites

Copyright 2011 © Limsoon Wong

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

Copyright 2011 © Limsoon Wong

34

What if there is no useful seq homolog?

- **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

Copyright 2011 © Limsoon Wong

35

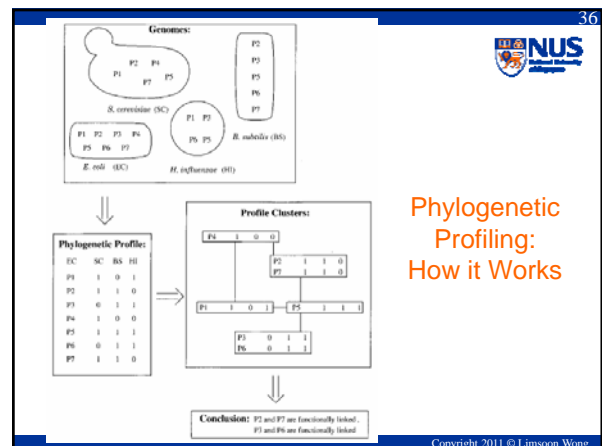
Phylogenetic Profiling

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

Copyright 2011 © Limsoon Wong



Protein Function Prediction from Protein Interactions

Level-1 neighbour

Level-2 neighbour

NUS National University of Singapore

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

NUS National University of Singapore

Copyright 2011 © Limsoon Wong

An illustrative Case of Indirect Functional Association?

SH3 Proteins SH3-Binding Proteins

Yfr024c Las17
Yvs167 Ypr171w
Ysc84 Act2
Ygr136w Ynl094w
Ypr154w Bdc1

- 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?

NUS National University of Singapore

Copyright 2011 © Limsoon Wong

Freq of Indirect Functional Association

Shared Functions with

	Frequency
Level-1 neighbours exclusively	0.016138
Level-2 neighbours exclusively	0.226474
Level-1 and Level-2 neighbours	0.706877

Source: Kenny Chua

NUS National University of Singapore

Copyright 2011 © Limsoon Wong

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

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

- 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

NUS National University of Singapore

Copyright 2011 © Limsoon Wong

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?

NUS National University of Singapore

Copyright 2011 © Limsoon Wong

49

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

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

Copyright 2011 © Limsoon Wong

50

Correlation w/ Functional Similarity

- Correlation between 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

Copyright 2011 © Limsoon Wong

51

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 between 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 between 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

Copyright 2011 © Limsoon Wong

52

Functional Similarity Estimate: FS-Weighted Measure with Reliability

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

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

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

⇒ Rewriting

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

Copyright 2011 © Limsoon Wong

53

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

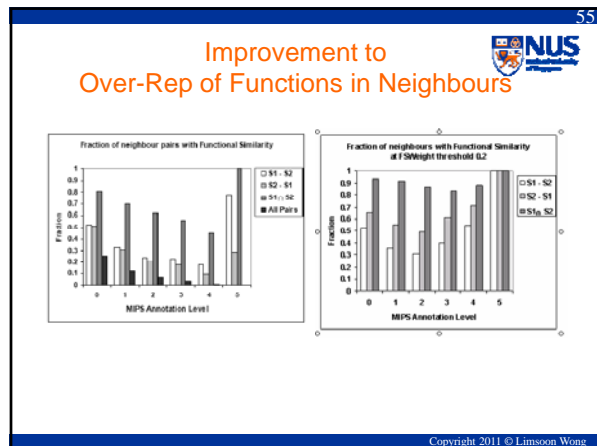
Copyright 2011 © Limsoon Wong

54

Improvement to Prediction Power by Majority Voting

Considering only neighbours w/ FS weight > 0.2

Copyright 2011 © Limsoon Wong



56

Use L1 & L2 Neighbours for Prediction

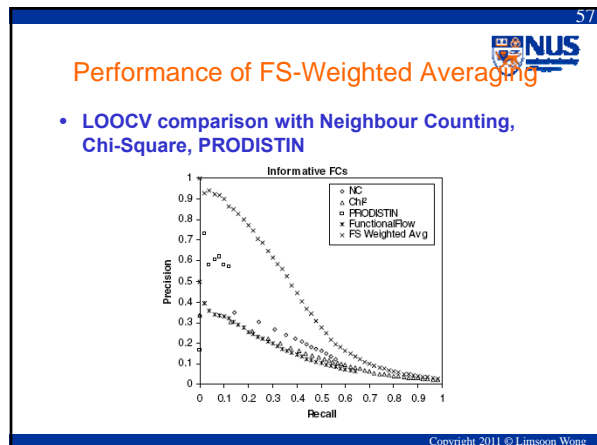
- **FS-weighted Average**

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

- 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

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

Copyright 2011 © Limsoon Wong



58

About the Inventor: Chua Hon Nian

- **Chua Hon Nian**
 - PhD, NUS, 2008
 - Postdoc at Harvard & Univ of Toronto
 - 49th hottest paper in Computer Science published in 2006
 - Winner, DREAM2 challenge PPI subnetwork, 2007

Copyright 2011 © Limsoon Wong

Application of Sequence Comparison: Key Mutation Site Discovery

60

Identifying Key Mutation Sites

K.L.Lim et al., JBC, 273:28986-28993, 1998

Sequence from a typical PTP domain D2

```
>g1 | 00000 | PTP1-D2
EEEEKLLTSIKIDKERTGHPANKEKKRVLQIPYEFHVIIPVERGEEDTFVMAISF
IDGTRQDSYILASQGPLLETIEDFUREIWEISCSIVELTELEKRGQKCLQFVPSDGLV
SYGDIIVLEKKEKCESTYVRDLLVINTREKESQRIQOFHFEQUPFVQIPSDGKGHSII
AAVQKQQQSGMHPITVECSAGLQRTOTFCALSTVLEKVKLEGILDFVQTVESLRQRP
EVQTLKQTEFCYKVVQYIDAFSDYIMFK
```

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

Copyright 2011 © Limsoon Wong

61

Emerging Patterns of PTP D1 vs D2

- 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

Copyright 2011 © Limsoon Wong

62

Emerging Patterns of PTP D1 vs D2

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

Exercise: Why?

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

X absent
 X present

Copyright 2011 © Limsoon Wong

63

Key Mutation Site: PTP D1 vs D2

g1|00000|P D2 QFHFGVPEVGIPSDGKMSIIAAVQKQQQ--SGNHPIIVHCSAGAGRTGTFALSTVL
 g1|126467| QFHFTSVPDFGVPTTIGMLKFLKKVKACNP--QYAGAIIVHCSAGVGTGTFFVVDAML
 g1|2499753| QFHFTGVPDGHGVPYHATGLLSFIRVKLSNP--PSAGPIVVHCSAGAGRTGTCYIVIDIML
 g1|462550| QYHVTQVPDGHGVPYALPVLTFFVRSAAAR--PETGPVLVHCSAGVGTGTGTIVIDSML
 g1|2499751| QFHFTSVPDGHGVPDITLLINFRVLVDYHQSPPESPILVHCSAGVGTGTFFIADRLI
 g1|1709906| QFQTAVPDGHGVPFPLFLRVRKTCNP--PDAGPIVVHCSAGVGTGTCTIVIDAML
 g1|126471| QLHFTSVPDFGVPTTIGMLKFLKKVKTLNP--VHAGPIVVHCSAGVGTGTCTIVIDAMH
 g1|548626| QFHFTGVPDGHGVPYHATGLLSFIRVKLSNP--PSAGPIVVHCSAGAGRTGTCYIVIDIML
 g1|131570| QFHFTGVPDGHGVPYHATGLLGFVRQVSKESP--PNAGPLVVHCSAGAGRTGTCYIVIDIML
 g1|2144715| QFHFTSVPDGHGVPDITLLINFRVLVDYHQSPPESPILVHCSAGVGTGTFFIADRLI
 * .. ** .*

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

Copyright 2011 © Limsoon Wong

64

Key Mutation Site: PTP D1 vs D2

g1|00000|P D2 QFHFGVPEVGIPSDGKMSIIAAVQKQQQ--SGNHPIIVHCSAGAGRTGTFALSTVL
 g1|126467| QFHFTSVPDFGVPTTIGMLKFLKKVKACNP--QYAGAIIVHCSAGVGTGTFFVVDAML
 g1|2499753| QFHFTGVPDGHGVPYHATGLLSFIRVKLSNP--PSAGPIVVHCSAGAGRTGTCYIVIDIML
 g1|462550| QYHVTQVPDGHGVPYALPVLTFFVRSAAAR--PETGPVLVHCSAGVGTGTGTIVIDSML
 g1|2499751| QFHFTSVPDGHGVPDITLLINFRVLVDYHQSPPESPILVHCSAGVGTGTFFIADRLI
 g1|1709906| QFQTAVPDGHGVPFPLFLRVRKTCNP--PDAGPIVVHCSAGVGTGTCTIVIDAML
 g1|126471| QLHFTSVPDFGVPTTIGMLKFLKKVKTLNP--VHAGPIVVHCSAGVGTGTCTIVIDAMH
 g1|548626| QFHFTGVPDGHGVPYHATGLLSFIRVKLSNP--PSAGPIVVHCSAGAGRTGTCYIVIDIML
 g1|131570| QFHFTGVPDGHGVPYHATGLLGFVRQVSKESP--PNAGPLVVHCSAGAGRTGTCYIVIDIML
 g1|2144715| QFHFTSVPDGHGVPDITLLINFRVLVDYHQSPPESPILVHCSAGVGTGTFFIADRLI
 * .. ** .*

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

Copyright 2011 © Limsoon Wong

65

Confirmation by Mutagenesis Expt

- 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?

Copyright 2011 © Limsoon Wong

66

About the Inventor: Prasanna Kolatkar

- Prasanna Kolatkar
 - Research Fellow, BIC, NUS, 1997-1999
 - Currently Group Leader at GIS

Copyright 2011 © Limsoon Wong

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

Copyright 2011 © Limsoon Wong

Any Question?



Acknowledgements



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

Copyright 2011 © Limsoon Wong

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.Altshul et al. "Basic local alignment search tool", *JMB*, 215:403--410, 1990
- S.F.Altshul et al. "Gapped BLAST and PSI-BLAST: A new generation of protein database search programs", *NAR*, 25(17):3389--3402, 1997

Copyright 2011 © Limsoon Wong

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

Copyright 2011 © Limsoon Wong



References

- H.N. Chua, W.-K. Sung. [A better gap penalty for pairwise SVM](#). 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—11, 2000
- T. Hawkins and D. Kihara. Function prediction of uncharacterized proteins. *JBCB*, 5(1):1-30, 2007