## Introduction to Protein Function Prediction

Wong Limsoon



National University of Singapore



- Sequence alignment
- Guilt by association
- Key mutation site discovery
- What if no homology of known function is found?

## Sequence comparison / alignment



National University of Singapore

### Motivations for seq comparison

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

Thus, sequence comparison is a foundation for inferring function, active site, and key mutations

# Sequence alignment

Key aspect of seq comparison is seq alignment

A seq alignment maximizes the number of positions that agree in two sequences



### Sequence alignment: Poor example

Poor seq alignment shows few matched positions  $\Rightarrow$  The 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 ILORGTPWADGTASISQCAINPGETFFYNFTVDNPGTFFYHGHLGMQRSAGLYGSLI 70 80 90 100 110 120 No obvious match between Amicyanin and Ascorbate Oxidase

### Sequence Alignment: Good example

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

D >gil13476732|refINP\_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
 MKPGRLASIALATIFLPMAVPAHAATIEITMENLVISPTEVSAKVGDTIRWVNKDVFAHT
 60

 MK G L
 ++
 MA PA AATIE+T++
 LV SP V AKVGDTI WVN DV AHT

 Sbjct:
 1
 MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNNDVVAHT
 60

Good match betw Amicyanin & an 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	FHFTSWPDFGVPFTPIG	MLKFLKKVKACNP-	QYAGAIV	HCS	GVGRTGTFV	/VIDAMLD
gi 2499753	FHFTGWPDHGVPYHATGI	LLSFIRRVKLSNP-	PSAGPIV <mark>V</mark>	HCSA	GAGRTGCYI	IVIDIMLD.
gi 462550	YHYTQWPDMGVPEYALPV	/LTFVRRSSAARM-	PETGPVI <mark>V</mark>	HCSA	GVGRTGTYI	<b>VIDSMLQ</b>
gi 2499751	FHFTSWPDHGVPDTTDLI	LINFRYLVRDYMK(	2SPPESPII <mark>.</mark> V	HCSA	GVGRTGTFI	LAIDRLIY
gi 1709906	FQFTAWPDHGVPEHPTPH	FLAFLRRVKTCNP-	PDAGPM <mark>V</mark> V	HCSA	GVGRTGCFI	IVIDAMLE
gi 126471	LHFTSWPDFGVPFTPIG	MLKFLKKVKTLNP-	VHAGP I <mark>v</mark> v	HCSA	GVGRTGTFI	<b>IVIDAMMA</b>
gi 548626	FHFTGWPDHGVPYHATGI	LLSFIRRVKLSNP-	PSAGPIVV	HCSA	GAGRTGCYI	<b>VIDIMLD</b>
gi 131570	FHFTGWPDHGVPYHATGI	LLGFVRQVKSKSP-	PNAGPLVV	HCSA	GAGRTGCFI	<b>VIDIMLD</b>
gi 2144715	FHFTSWPDHGVPDTTDLI	LINFRYLVRDYMK(	2SPPESPIL <mark>V</mark>	HCSA	GVGRTGTFI	LAIDRLIY
	+ +++ +++	+			+ ++++	+ +

Conserved sites

# Guilt by association



# National University of Singapore

### **Proteins**

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



### Function assignment to protein seq

SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE VT

How to assign a function to a new protein sequence?

### In the course of evolution...



### Evolution takes time ...

Let a = AFPHQHRVP Let b = PQVYNIMKE

Suppose each generation differs from the previous by 1 residue

What is the max difference between the 2<sup>nd</sup> generation of a

What is the min difference between the 2<sup>nd</sup> generation of a and b?

### The triumph of logic

In the course of evolution...



Evolution takes time ...

Let a = AFPHQHRVP Let b = PQVYNIMKE

Suppose each generation differs from the previous by 1 residue

What is the max difference between the  $2^{nd}$  generation of a

What is the min difference between the  $2^{nd}$  generation of a and b?

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

### **Discussion**

### How to guess the function of a protein?

### **Guilt by association: Caveats**

Ensure that the effect of database size has been accounted for

Ensure that the function of the homolog 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 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  $10^7$  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!

### 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) 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

BLAST employs the SEG algo to filter low complexity regions from proteins before executing a search

Source: NCBI

### **Effect of sequence length**



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

18 entries were found										
D	Organism	PIR	Swiss-Prot/TrEMBL	RefSeq/GenPept						
NF00181857	Methanococcus jannaschii	<mark>E64381</mark> conserved hypothetical protein MJ0653	Y653_METJA Hypothetical protein MJ0653	g <u>1592300</u> inosine-5'-monophosphate dehydrogenase (guaB) <u>NP_247637</u> inosine-5'-monophosphate dehydrogenase (guaB)						
<u>NF00187788</u>	Archaeoglobus fulgidus	G69355 MJ0653 homolog AF0847 ALT_NAMES: inosine-monophosphate dehydrogenase (guaB-1) homolog [misnomer]	<u>O29411</u> INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-1)	g2649754 inosine monophosphate dehydrogenase (guaB-1) <u>NP_069681</u> inosine monophosphate dehydrogenase (guaB-1)						
<u>NF00188267</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)						
<u>NF00188697</u>	Archaeoglobus fulgidus	B69407 MJ0188 homolog ALT_NAMES: inosine monophosphate dehvdrogenase homolog [misnomer]	O29009 Hypothetical protein AF1259	g2649320 inosine monophosphate dehydrogenase, putative <u>NP_070087</u> inosine monophosphate putative						
	A partial list of IMPdehydrogenase misnomers in -5-monophosphate									
<u>NF00197776</u>	<sup>Thermo</sup> complete	genomes remaini	ng in some public d	atabases <sup>monophosphate</sup>						
<u>NF00414709</u>	Methanothermobacter thermautotrophicus	<u>G69030</u> MJ0653 homolog MTH1226 <i>ALT_NAMES</i> : inosine-monophosphate dehydrogenase related protein V [misnomer]	O27294 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN V	g2622337 inosine-5'-monophosphate dehydrogenase related protein V <u>NP_276354</u> inosine-5'-monophosphate dehydrogenase related protein V						
<u>NF00414811</u>	Methanothermobacter thermautotrophicus	D69035 MJ1232 protein homolog MTH126 ALT_NAMES: inosine-5'-monophosphate dehydrogenase related protein VII [misnomer]	O26229 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN VII	<u>g2621166</u> inosine-5'-monophosphate dehydrogenase related protein VII <u>NP_275269</u> inosine-5'-monophosphate dehydrogenase related protein VII						
<u>NF00414837</u>	Methanothermobacter thermautotrophicus	H69232 MJ1225-related protein MTH992 <i>ALT_NAMES</i> : inosine-5'-monophosphate dehydrogenase related protein IX [misnomer]	O27073 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN IX	g2622093 inosine-5'-monophosphate dehydrogenase related protein IX <u>NP_276127</u> inosine-5'-monophosphate dehydrogenase related protein IX						
<u>NF00414969</u>	Methanothermobacter thermautotrophicus	B69077 yhcV homolog 2 ALT_NAMES: inosine-monophosphate dehydrogenase related protein X [misnomer]	<u>027616</u> INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN X	g2622697 inosine-5'-monophosphate dehydrogenase related protein X <u>NP_276687</u> inosine-5'-monophosphate dehydrogenase related protein X						

### **IMPDH domain structure**

	中       PCM00487: PD0C00391,IMP dehydrogenase / GMP reductase signature         中日       PF00478: IMP dehydrogenase / GMP reductase C terminus         中日       PF00571: CBS domain         ●       PF01381: Helix-turn-helix         PF01574: IMP dehydrogenase / GMP reductase N terminus         PF01574: IMP dehydrogenase / GMP reductase N terminus         ゆごについい         ゆごについい
A31997 (SF000130)	
E70218 (SF000131)	<i>ֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈֈ</i>
E64381 (SF004696)	194 IMPDH Misnomer in Methanococcus jannaschii
G69355 (SF004696)	
F69514 (SF004694)	OCOCOCOC ******************************
869407 (SF004699)	

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>	<u>051440</u>	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) / phosphoribosyl-ATP pyro (EC 3.6.1.31) [similarity]	e (EC ophosphatase	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]	<u>e (EC</u> ophosphatase	Synechocystis sp.	Prok/gram-	557.0	5.7e-24	230	39.175	194	
	T35073	SF029243	005738	probable phosphoribosyl-AMP cyclo	hydrolase	Streptomyces coelicolor	Prok/gram+	399.3	3.5e-15	128	42.157	102	
	□ <u>\$53349</u>	<u>SF001257</u>	001188	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) / phosphoribosyl-ATP pyro (EC 3.6.1.31) / histidinol dehydroger 1.1.1.23)	<u>e (EC</u> pphosphatase nase (EC	Saccharomyces cerevisiae	Euk/fungi	384.1	2.5e-14	799	31.863	204	
A⊨>	E69493	SF029243	<u>005738</u>	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) [similarity]	<u>e (EC</u>	Archaeoglobus fulgidus	Archae	396.8	4.8e-15	108	47.778	90	
C⇒	□ <u>G64337</u>	SF006833	<u>030827</u>	phosphoribosyl-ATP pyrophosphata: 3.6.1_31) [similarity]	se (EC	Methanococcus jannaschii	Archae	246.9	1.1e-0ó	95	36.842	95	
	D81178	<u>SF006833</u>	<u>101491</u>	phosphoribosyl-ATP pyrophosphata: 3.0.1.31) NMB0603 [similarity]	se (EC	Neisseria meninoitidis	Prok/oram-	239 9	2 fie-0fi	107	35 227	88	
	□ <u>G81925</u>	SF006833	<u>101491</u>	hosphoribosyl-ATP pyrophosphat 3.6.1.31) NMA0807 [similarity]		$A \rightarrow B$	-> C =	=> .	A -> (	С			-
phosphoribosyl-AMP cyclohydrola			<b>B</b> (SF001258)										
	□ <u>\$51513</u>	<u>SF001257</u>	001188	(EC 3.6.1.31) / histidinol dehydrog 1.1.1.23)		7			X				
Mis-assignment of function					╳	-							
			A (SF029243) C (SF006833)										
			No IMPDH domain										
						II DII UUIIIa							

### **Emerging pattern**



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

# Key mutation site discovery



# National University of Singapore

### Sequence from a typical PTP domain D2

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

### Discussion

Some PTPs have 2 PTP domains

PTP domain D1 has more activity than PTP domain D2

Why? How do you figure that out?

### Key mutation site: PTP D1 vs D2

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



Positions marked by "!" and "?" are likely places responsible for reduced PTP activity

PTP D1 agree on them; PTP D2 disagree on them

### Key mutation site: PTP D1 vs D2

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



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

### **Confirmation by mutagenesis**

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

Why do you need this 2-way expt?

# What if no homolog of known function is found?



National University of Singapore

### 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 physicochemico properties(e.g., PROTFUN) Similarity of gene expression profiles Similarity of protein-protein interaction partners

### Hidden Markov Model for biological seqs

https://webpages.math.luc.edu/~tobrien/courses/bioinf/krogh.pdf

# A multiple alignment of sequences of the same family/function

A C A - - - A T G T C A A C T A T C A C A C - - A G C A G A - - - A T C A C C G - - A T C For each protein family, create multiple alignment of its members

(I use DNA only because of the smaller number of letters than for amino acids).

### HMM derived from the multiple alignment



### Scoring by HMM

https://webpages.math.luc.edu/~tobrien/courses/bioinf/krogh.pdf



The prob depends on seq length; not very convenient score for interpretation

Log odds obtained by comparing to a null model of the same length

log-odds for sequence 
$$S = \log \frac{P(S)}{0.25^L} = \log P(S) - L \log 0.25$$
.

For a 4-letter alphabet, each letter has 0.25 prob to occur

### **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

 $\Rightarrow$  Even if no homolog with known function is available, it is still possible to infer function of a protein!

### 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

$$w_{z} = \binom{N}{z}$$
$$\overline{w_{z}} = \binom{N-z}{x-z} * \binom{N-x}{y-z}$$
$$W = \binom{N}{x} * \binom{N}{y}$$

### **Phylogenetic profiles: Evidence**

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

Kanword	No. of non- homologous proteins in	No. neighbors in keyword	No. neighbors in random
Keywold	group	group	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



Explain the two red ovals. Any surprise there?

## Concluding remarks



# National University of Singapore

### **Discussion**

What have we learned?

### **Suggested readings**

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

### **Suggested readings**

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

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