For written notes on this lecture, please read chapter 19 of The Practical Bioinformatician

KI1972: Applied Bioinformatics & Computational Biology Sequence Homology Interpretation

Limsoon Wong June 2006





Plan



- Recap of sequence alignment
- Guilt by association
- Active site/domain discovery
- Key mutation site discovery
- Guilt by other types of association
 - Genome phylogenetic profiling
 - Protfun
 - SVM-Pairwise

Very Brief Recap of Sequence Comparison/Alignment





Motivations for Seq Comparison

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

Karolinska Institutet

Sequence Alignment: Good Example

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

D >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
 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	FHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIV/HCSAGVGRTGTFVVIDAML	D
gi 2499753	FHFTGWPDHGVPYHATGLLSFIRRVKLSNPPSAGPIVVHCSAGAGRTGCYIVIDIML	D
gi 462550	YHYTQWPDMGVPEYALPVLTFVRRSSAARMPETGPVIVHCSAGVGRTGTYIVIDSML(2
gi 2499751	FHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPII.VHCSAGVGRTGTFIAIDRLI	Ŷ
gi 1709906	FQFTAWPDHGVPEHPTPFLAFLRRVKTCNPPDAGPMVVHCSAGVGRTGCFIVIDAMLI	Ξ
gi 126471	LHFTSWPDFGVPFTPIGMLKFLKKVKTLNPVHAGPIVVHCSAGVGRTGTFIVIDAMM	A
gi 548626	FHFTGWPDHGVPYHATGLLSFIRRVKLSNPPSAGPIVVHCSAGAGRTGCYIVIDIML	D
gi 131570	FHFTGWPDHGVPYHATGLLGFVRQVKSKSPPNAGPLVVHCSAGAGRTGCFIVIDIML	D
gi 2144715	FHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGTFIAIDRLI	Ŷ
	* * * * * * * * * * * * * * * * * * * *	

Conserved sites

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

Application of Sequence Comparison: Guilt-by-Association





Function Assignment to Protein Seq

SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE VT

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

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

Guilt-by-Association



- Compare the target sequence T with sequences
 S₁, ..., S_n of known function in a database
- Determine which ones amongst S₁, ..., 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





KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006



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



🞒 NCBI BLAST - M	licrosoft Internet Explorer		
File Edit View F	avorites Tools Help - 💌 🙆 🏠 🔎 Search ☆ Favorites 🧭	🔗 · 🍃 · 📴 🍇 🚳	NUS National University of Singapore
Address 🙆 http://ww	w.ncbi.nlm.nih.gov/BLAST/		Karolinska
Google -	💌 <u>G</u> Search 🝷 🥥 🌍 👹 Pagef	Bank 💁 6 blocked 🛛 🎸 Check 🔹 💐 AutoLink 🔹 🗐 Au	<u>ک</u> کی کھی کے کھی
S NCBI → BLAST		Latest news: 6 December 2005 : BLAST 2.2.13 released	
About • Getting started • News	The Basic Local Alignment Search Tool (BLAS sequences. The program compares nucleotide or pro the statistical significance of matches. BLAST can be between sequences as well as help identify members	T) finds regions of local similarity between tein sequences to sequence databases and calculates used to infer functional and evolutionary relationships s of gene families.	
• FAQs	Nucleotide	Protein	
More info NAR 2004 NCBI Handbook The Statistics of Sequence Similarity Scores Software	 Quickly search for highly similar sequences (megablast) Quickly search for divergent sequences (discontiguous megablast) Nucleotide-nucleotide BLAST (blastn) Search for short, nearly exact matches Search trace archives with megablast or discontiguous megablast 	 Protein-protein BLAST (blastp) Position-specific iterated and pattern-hit initiated BLAST (PSI- and PHI-BLAST) Search for short, nearly exact matches Search the conserved domain database (rpsblast) Protein homology by domain architecture (cdart) 	
 Downloads Developer info 	Translated	Genomes	
Dther resources References NCBI Contributors Mailing list Contact us 	 Translated query vs. protein database (blastx) Protein query vs. translated database (tblastn) Translated query vs. translated database (tblastx) 	 Human, mouse, rat, chimp, cow, pig, dog, sheep, cat Chicken, puffer fish, zebrafish Fly, honey bee, other insects Microbes, environmental samples Plants, nematodes Fungi, protozoa, other eukaryotes 	
ē)			

G Back 🔹 🕞 🛸	🞽 😰 🎧 🔎 Search 🎇 Favorites 🚱 🔯 🖓 😒 🔛 🔭 🛄 🛄 🦄	
Address 🕘 http://www.r	ncbi.nlm.nih.gov/BLAST/Blast.cgi?CMD=Web&LAYOUT=TwoWindows&AUTO_FORMAT=Semiauto&ALIGNMENTS	National Univer of Singapore
Google -	📉 🖸 Search 🔹 🏈 📚 🛷 PageHank 🕸 6 blocked 👋 Check 🔹 🛝 AutoLir	Karolins
$\boldsymbol{\varsigma}$		
S NCBI	protein–protein BLASI	
Nucleotide	Protein Translations Retrieve results for an RID	
	NRYVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWR	
Search	MIWEQNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFC	
	SAGVGRTGTFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLE	
	HYLYGDTELE 💌	
Cut and the second		
<u>Set subsequence</u>	From: To:	
Choose database	pr 💌	
Do CD-Search		
Now:	BLAST! or Reset query Reset all	
Options	for advanced blasting	
-		
Limit by entrez	or select from: All organisms	
query		

🌍 Back 🔹 🌍 🔹 📕	🔁 🎧 🎾 Search 🌱	🕇 Favorites 🛛 🔀	• 🍓 🖻 🛯 📙	1 🚳		尺 NU
Address 🙆 http://www.ncbi.nlm	.nih.gov/BLAST/Blast.cgi	1				
Google -	👻 <u>G</u> Search 🝷 🥝) 🥥 ø PageRank	🛂 6 blocked 🛛 🍄 Check 🗠	🔸 🌂 AutoLink	- 📲 AutoFill	🍋 Options 🏼 🥖
S NCBI	Protoin	formattin	g BLAST			
Nucleotide	Protein	Translations	Retrieve results for an F	RID		
Query = (302 letters) Putative conserved dom	ains have been detecte	d, click on the imag	e below for detailed re	esults.	250	202
	Ť.	100	DTDo	ļ	1	302
The request ID is 1146357 Format! or Reset all The results are estimated to be	911-10937-56421952334.E e ready in 9 seconds but may	LASTQ4 be done sooner.]			
Please press "FORMAT!" wh request results of a different s	en you wish to check your re earch by entering any other v	sults. You may change t valid request ID to see o	he formatting options for y ther recent jobs.	our result via the	e form below an	d press "FORMAT!" a

16



Homologs obtained by BLAST

Score E Sequences producing significant alignments: (bits) Value 62: e-177 gi|14193729|gb|AAK56109.1|AF332081 1 protein tyrosin phosph... 621 e-177 gi|126467|sp|P18433|PTRA HUMAN Protein-tyrosine phosphatase... 621 e-176 qi 4506303 [ref NP 002827.1] protein tyrosine phosphatase, r... qi|227294|prf||1701300A protein Tyr phosphatase 620 e-176 621 L e-176 gi|18450369|ref|NP 543030.1| protein tyrosine phosphatase, ... 61 e-176 gi[32067[emb]CAA37447.1] tyrosine phosphatase precursor [Ho... gi|285113|pir||JC1285 protein-tyrosine-phosphatase (EC 3.1.... 619 e-176 61: L e-176 gi 6981446 ref NP 036895.1 protein tyrosine phosphatase, r... 61 S e-174 gi 2098414 pdb 1YFO A Chain A, Receptor Protein Tyrosine Ph... 61 L e-174 qi|32313|emb|CAA38662.1| protein-tyrosine phosphatase [Homo... qi|450583|qb|AAB04150.1| protein tyrosine phosphatase >qi|4... 605 e-172 60 L e-172 gi|6679557|ref|NP_033006.1| protein tyrosine phosphatase, r... gi|483922|gb|AAA17990.1| protein tyrosine phosphatase alpha 599 e-170

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



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 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE 300 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE
- Sbjct: 442 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQVVFIYQALLEHYLYGDTELE 501



HSPs, E-Value, Bits, & P-Value

• HSPs

- A local alignment without gaps consists simply of a pair of equal length segments, one from each of the two sequences being compared.
- A segment pair whose score cannot be improved by extension or trimming is called high-scoring segment pairs or HSPs

• E-Value

- For large seq lengths *m* and *n*, the stats of HSP scores are characterized by two params, *K* and λ
- Expected number of HSPs with score > S is given by E = Kmne^{$-\lambda$ S}

Source: NCBI



HSPs, E-Value, Bit Score, & P-Value

Bit Score

- "Citing a raw score alone is like citing a distance without specifying feet, meters, or light years"
- Normalize raw score to S' = $(\lambda S \ln K)$ / ln 2 to get "bit score", which has a standard set of units
- E-value corresponds to bit score as $E = mn2^{-S'}$
- P-Value
 - Number of random HSPs with score \geq S is described by a Poisson distribution
 - \Rightarrow Chance of finding no HSPs with score \ge S is e^{-E}
 - \Rightarrow Prob of finding \geq 1 such HSP is P = 1 e^{-E}

Source: NCBI



Guilt-by-Association: Caveats

- Ensure that the effects of database size and composition have 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 10⁻⁶
- If database has 10⁷ seqs, then you expect 10⁷ * 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



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

Source: NCBI



Effect of Sequence Length



KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006



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



18 entries were found

ID	Organism	PIR	Swiss-Prot/TrEMBL	RefSeq/GenPept
<u>NF00181857</u>	Methanococcus jannaschii	<u>E64381</u> conserved hypothetical protein MJ0653	<u>Y653_METJA</u> Hypothetical protein MJ0653	<u>g1592300</u> inosine-5'-monophosphate dehydrogenase (guaE) <u>NP_247637</u> inosine-5'-monophosphate dehydrogenase (guaE)
<u>NF00187788</u>	Archaeoglobus fulgidus	G69355 MJ0653 homolog AF0847 ALT_NAMES: inosine-monophosphate dehydrogenase (guaB-1) homolog [misnomer]	<u>029411</u> INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-1)	<u>g2649754</u> 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>	Archae A partia	l list of IMPde	ydrogenase misn	omers ophosphate ive inophosphate ive
<u>NF00197776</u>	Thermo in CO	mplete genome public d	s remaining in so atabases	nophosphate d protein nonophosphate d protein
NF00414709	Methanothermobacter thermautotrophicus	ALT_NAMES: inosine-monophosphate dehydrogenase related protein V [misnomer]	O27294 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN V	onophosphate 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
NF00414837	Methanothermobacter thermautotrophicus	<u>H69232</u> MJ1225-related protein MTH992 <i>ALT_NAMES</i> : inosine-5'-monophosphate dehydrogenase related protein IX [misnomer]	O27073 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN IX	<u>g2622093</u> inosine-5'-monophosphate dehydrogenase related protein IX <u>NP_276127</u> inosine-5'-monophosphate dehydrogenase related protein IX
<u>NF00414969</u>	Methanothermobacter thermautotrophicus	<u>B69077</u> yhcV homolog 2 <i>ALT_NAMES:</i> inosine-monophosphate dehydrogenase related protein X [misnomer]	<u>027616</u> INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN X	<u>g2622697</u> inosine-5'-monophosphate dehydrogenase related protein X <u>NP_276687</u> inosine-5'-monophosphate dehydrogenase related brokein Xathy Wu

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

	**************************************	mstitu
	PCM00487: PD0C00391,IMP dehydrogenase / GMP reductase signature PE00478: IMP dehydrogenase / GMP reductase C terminus	
	PF00571: CBS domain	
004	PF01381: Helixturn-helix	
- Testar	PF01574: IMP dehydrogenase / GMP reductase N terminus	
ghaiga	A PF02195: ParB-like nuclease domain	
A31997	514	
5r000130)		
E70240		
SF000131)		
56000131)		
E64381 SF0004696)	404 404 194 IMPDH Misnomer in <i>Methanococcus jannaschii</i>	<u> </u>
E64381 SF004696)	404 404 194 IMPDH Misnomer in Methanococcus jannaschii	_

- 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

Source: Cathy Wu



Invalid Transitive Assignment



Root of invalid transitive assignment _

B⊨⇒	□ <u>H70468</u>	SF001258	051440	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) / phosphoribosyl-ATP pyro	<u>(EC</u> phosphatase	Aquifex aeolicus	Prok/other	594.3	4.8e-26	205	39.086	197	
	□ <u>\$76963</u>	<u>SF001258</u>	<u>039935</u>	(EC 3.6.1.31) [similarity] phosphoribosyl-AMP cyclohydrolase 3.5.4.19) / phosphoribosyl-ATP pyro (EC 3.6.1.31) [similarity]	(EC phosphatase	Synechocystis sp.	Prok/gram-	557.0	5.7e-24	230	39.175	194	
	T35073	SF029243	005738	probable phosphoribosyl-AMP cyclol	hydrolase	Streptomyces coelicolor	Prok/gram+	399.3	3.5e-15	128	42.157	102	
	□ <u>\$53349</u>	<u>SF001257</u>	<u>001188</u>	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) / phosphoribosyl-ATP pyroj (EC 3.6.1.31) / histidinol dehydrogen: 1.1.1.23)	<u>(EC</u> phosphatase ase (EC	Saccharomyces cerevisiae	Euk/fungi	384.1	2.5e-14	799	31.863	204	
A	□ <u>E69493</u>	SF029243	005738	phosphoribosyl-AMP cyclohydrolase 3.5.4.19) [similarity]	(EC	Archaeoglobus fulgidus	Archae	396.8	4.8e-15	108	47.778	90	
C	□ <u>G64337</u>	SF006833	030827	phosphoribosyl-ATP pyrophosphatas 3.6. [431] [similarity]	e (EC	Methanococcus jannaschii	Archae	246.9	1.1e-06	95	36.842	95	_
	D81178	SF006833	<u>101491</u>	phosphoribosyl-ATP pyrophosphatas 3.0.1.31) NMB0603 [similarity]	e (EC	Neisseria meninoitidis	Prok/oram-	239.9	2 fe-Nf	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			-
	□ <u>\$51513</u>	<u>SF001257</u>	001188	phosphoribosyl-AMP cyclohydrola 3.5.4.19) / phosphoribosyl-ATP py (EC 3.6.1.31) / histidinol dehydrog 1.1.1.23)		1	B (SFC	01258)	1			-
Ν	lis-as	ssign	me	nt	A	(SF029243)	X >		C	(SF	00683	3)	
0	f funo	ction	l		No I	MPDH doi	main				S	Sourc	e: Cathy Wu

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

Emerging Pattern Karolinska Institutet **Typical IMPDH** Functional IMPDH w/o CBS 60 PCM00487: PDOC00391.IMP dehvdrogenase / GMP reductase signature Exercise: How do you A DESCRIPTION OF THE OWNER OF THE PF00478: IMP dehvdrogenase / GMP reductase C terminus *** PF00571: CBS domain Recognize this kind of 00-00 PF01381: Helix-turn-helix 17-17-17-17 PF01574: IMP dehvdrogenase / GMP reductase N terminus problems? - And and and a start of the st PF02195: ParB-like nuclease domain A31997 514 (SF000130) detectoriete والمراد والمراد والمراد والمراد والمراد والمراد والمراد والمراد والمراد E70218 404 (SF000131) E64381 IMPDH Misnomer in Methanococcus jannaschii 194 (SF004696) viciojojojojojoj 669355 189 (SF004696) iajajajajajaj ala (ala) ala (ala) ala F69514 **IMPDH** Misnomers in Archaeoglobus fulgidus 183 (SF004694) 0000000 volotototo B69407 259 (SF004699) ոնորոներությունորորոն

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

Source: Cathy Wu

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

Application of Sequence Comparison: Active Site/Domain Discovery





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
 - \Rightarrow Emerging patterns relative to background
 - \Rightarrow Candidate active sites and/or domains
- Easier if sequences of distance homologs are used

Exercise: Why?



KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006



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

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 is has much more activity than PTP domain D2
- Why? And how do you figure that out?



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



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 ⇒ it is an EP ⇒ possible cause of D2's loss of function



Key Mutation Site: PTP D1 vs D2

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

2 2 22 2 QFHFHGWPEVGIPSDGKGMISIIAAVQKQQQQ-SGNHPITVHCSAGAGRTGTFCALSTVL OFHFTSWPDFGVPFTPIGMLKFLKKVKACNP--OYAGAIVVHCSAGVGRTGTFVVIDAML OFHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIML OYHYTOWPDMGVPEYALPVLTFVRRSSAARM--PETGPVLVHCSAGVGRTGTYIVIDSML OF HF TSWPDHGVPDTTDLL INFRYLVRDYMKOSPPESPILVHCSAGVGRTGTFIAIDRLI QFQFTAWPDHGVPEHPTPFLAFLRRVKTCNP--PDAGPMVVHCSAGVGRTGCFIVIDAML D1-OLHFTSWPDFGVPFTPIGMLKFLKKVKTLNP--VHAGPIVVHCSAGVGRTGTFIVIDAMM OFHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIML OFHFTGWPDHGVPYHATGLLGFVROVKSKSP--PNAGPLVVHCSAGAGRTGCFIVIDIML QFHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGTFIAIDRLI ***** **** Τ. **. *.*

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



Key Mutation Site: PTP D1 vs D2



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



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?

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

genome phylogenetic profiles protfun's feature profiles









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 Profiling: How it Works



Copyright 2006 © Limsoon Wong

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

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
 $W = \binom{N-z}{x-z} * \binom{N-z}{y-z}$
No. of ways to distribute
the remaining $x - z$ and $y - z$
occurrences over the remaining
 $N - z$ lineage's
 $W_z = \binom{N}{x} * \binom{N}{y}$
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 [†]	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: Why do proteins having high hamming distance also have this behaviour?

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006



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, phospharylation, 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?

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006

ProtFun: Evidence



Institutet



 Combinations of "features" seem to
 characterize some functional categories

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006



Karolinska

ProtFun: How it Works

	F	D = = = = i = ti = ti		
Abbrimation	Encoding	Description		
ec	single value	Extinction coefficient predicted by ExPASy ProtParam		
gravy	single value	Hydrophobicity predicted by ExPASy ProtParam		
nneg	single value	Number of negatively charged residues counted by ExPASy P	rotParam	
npos	single value	Number of positively charged residues counted by ExPASy Pr	otParam	
nglyc	potential in 5 bins	N-glycosylation sites predicted by NetNGlyc		
oglyc	potential-threshold in 10 bins	GaINAc O-glycosylations predicted by NetOGlyc		
pest	fraction in 10 bins	PEST rich regions identified by PESTfind		
phosST	potential in 10 bins	Serine and threonine phosporylations predicted by NetPhos		
phosY	potential in 10 bins	Tyrosine phosporylations predicted by NetPhos	xtract	feature
psipred	helix, sheet, coil in 5 bins	Predicted secondary structure from PSI-Pred p	rofile o	of protein
psort	20 probabilities	Subcellular location predtions by PSORT	sing va	rious
seg	fraction in 10 bins	Low-complexity regions identified by SEG	rodicti	on methods
signalp	meanS, maxY, log(cleavage pos)	Signal peptide predictions made by SignalP	reureu	on memous
tmhmm	inside, outside, membrane in 5 bins	Transmembrane helix predictions made by TMHMM		

Category	Hidden units	Input features
Amino acid biosynthesis	30	ec psipred psort tmhmm
	30	ec psipred tmhmm
A your of the output of	, 30	ec netoglyc psipred psort
Average the output of	30	gravy psipred psort
the 5 component ANN	S 30	oglyc psipred psort

KI 1972: Applied Bioinformatics & Computational Biology, Stockholm, June 2006



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, tranport & binding

This is in agreement w/ known functionality of these proteins



Karolinska Institutet

ProtFun: Performance





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.



median RFP

Copyright 2000 © Linisoon wong

Any Questions?







- 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
- RA Abagyan, S Batalov. Do aligned sequences share the same fold? *JMB*, 273(1):355-68, 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
- C. Wu, W. Barker. "A Family Classification Approach to Functional Annotation of Proteins", *The Practical Bioinformatician*, Chapter 19, pages 401—416, WSPC, 2004
- H.N. Chua, W.-K. Sung. <u>A better gap penalty for pairwise SVM</u>. 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