

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

# Knowledge Discovery Techniques for Bioinformatics, Part VI: Sequence Homology Interpretation

Limsoon Wong



Lecture at Yang Ming National University, Taipei, June 2006

Copyright © 2006 by Limsoon Wong

2

## Plan



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

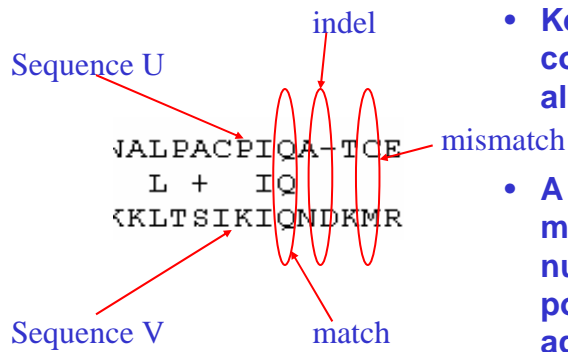
Lecture at Yang Ming National University, Taipei, June 2006

Copyright 2006 © Limsoon Wong

# Very Brief Recap of Sequence Comparison/Alignment



## Sequence Alignment



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



# Sequence Alignment: Poor Example

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

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

```

                60      70      80      90      100
Amicyanin      MPHNVHVFVAGVVGAAALKGFMMKKEQAYSLETFTEAGTYDYHCTFHPFMRGKVVVE
                :...: . :...: ::
Ascorbate Oxidase ILQRQTFWADGTASISQCAINPGETFFYNFTVDNPGTFFYHGHLMQRSAGLYGSLI
                70      80      90      100      110      120

```

No obvious match between Amicyanin and Ascorbate Oxidase



# Sequence Alignment: Good Example

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

```

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

```

```

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

```

```

Query: 1  MKPGRLASIALAIIFLPMVPAHAATIEITMENLVISPTVEVSAKVGDTIRWVNKDVFAHT 60
          MK G L  ++      MA PA AATIE+T++ LV SP  V AKVGDTI WVN DV AHT
Sbjct: 1  MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNDVVAHT 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|      FHFTS W P D F G V P F T P I G M L K F L K K V K A C N P -- Q Y A G A I V V H C S A G V G R T G T F V V I D A M L D
gi|2499753|     FHFTG W P D H G V P Y H A T G L L S F I R R V K L S N P -- P S A G P I V V H C S A G A G R T G C Y I V I D I M L D
gi|462550|      Y H Y T Q W P D H G V P E Y A L P V L T F V R R S S A A R M -- P E T G P V I V H C S A G V G R T G T Y I V I D S M L Q
gi|2499751|     FHFTS W P D H G V P D T T D L L I N F R Y L V R D Y M K Q S P P E S P I L V H C S A G V G R T G T F I A I D R L I Y
gi|1709906|     F Q F T A W P D H G V P E H P T P F L A F L R R V K T C N P -- P D A G P M V V H C S A G V G R T G C F I V I D A M L E
gi|126471|      L H F T S W P D F G V P F T P I G M L K F L K K V K T L N P -- V H A G P I V V H C S A G V G R T G T F I V I D A M M A
gi|548626|      F H F T G W P D H G V P Y H A T G L L S F I R R V K L S N P -- P S A G P I V V H C S A G A G R T G C Y I V I D I M L D
gi|131570|      F H F T G W P D H G V P Y H A T G L L G F V R Q V K S K S P -- P N A G P L V V H C S A G A G R T G C F I V I D I M L D
gi|2144715|     F H F T S W P D H G V P D T T D L L I N F R Y L V R D Y M K Q S P P E S P I L V H C S A G V G R T G T F I A I D R L I Y
..* *** **      . *                               ..***** **.. ** ..
```

Conserved sites

## Application of Sequence Comparison: Guilt-by-Association

## Function Assignment to Protein Sequence

SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR  
 YVNILPYDHSRVHLTPVEGVPSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMWE  
 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD  
 VTNRKPQRLITQFHFTSWPDFGVPTPIGMLKFLKVKACNPQYAGAIVVHCSAGVGRGTG  
 TFVVIDAMLDMMHSEKVDVYGFVSRIRAQRQMVQTDMQYVFIYQALLEHYLYGDTELE  
 VT

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

## Guilt-by-Association

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

## Guilt-by-Association

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

**Poor Sequence Alignment**

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

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

```

Amicyanin      60  70  80  90 100
KPRVYVAVDYLGAALGRHKKGGAYLITFFLAGDIDICTYRFFRGGVVV
Ascorbate Oxidase 1LQKQTFWAGDTASTSQCAIHPDCTFFYFVDPOTFFYFKLQKDSAGLVD
                    70  80  90 100 110
    
```

No obvious match between Amicyanin and Ascorbate Oxidase

Discard this function as a candidate

**Good Sequence Alignment**

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

```

g114613921.g114613921.11 unknown protein [Mesochorus loti]
g114613921.g114613921.11 unknown protein [Mesochorus loti]
Length = 105

Score = 105 bits (242), Expect = 1e-22
Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)

Query: 1  MDPQLALALAIPLPMAYAAATLITDNEFLIDTEYGAWDTLRYDQVFAIT 60
           M G L  **  M A P A AATLIT** L F D F  Y AKWDTI  W N D F AIT
Sbjct: 1  MWAHLLELALALALPMAYAAATITDTLIDTEYGAWDTLRYDQVFAIT 60
    
```

good match between Amicyanin and unknown M. loti protein

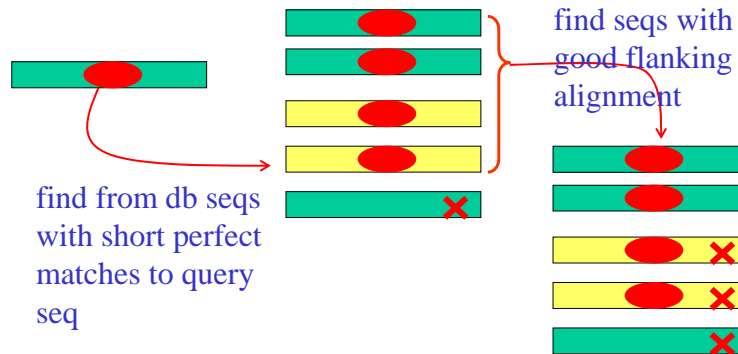
Assign to  $T$  same function as homologs

Confirm with suitable wet experiments

## BLAST: How It Works

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

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





## Homologs obtained by BLAST

Sequences producing significant alignments:	Score (bits)	E Value
<a href="#">gi 14193729 cb AAK56109.1 AF332081.1</a> protein tyrosin phosph...	62.1	e-177
<a href="#">gi 126467 sp P18433 PTRA_HUMAN</a> Protein-tyrosine phosphatase...	62.1	e-177
<a href="#">gi 4506303 ref NP_002827.1</a> protein tyrosine phosphatase, r...	62.1	e-176
<a href="#">gi 227294 prf I1701300A</a> protein Tyr phosphatase	62.0	e-176
<a href="#">gi 18450369 ref NP_543030.1</a> protein tyrosine phosphatase, ...	62.1	e-176
<a href="#">gi 32067 emb CAA37447.1</a> tyrosine phosphatase precursor [Ho...	61.1	e-176
<a href="#">gi 285113 pir JJC1285</a> protein-tyrosine-phosphatase (EC 3.1....	61.9	e-176
<a href="#">gi 6981446 ref NP_036895.1</a> protein tyrosine phosphatase, r...	61.1	e-176
<a href="#">gi 2098414 pdb 1YFO A</a> Chain A, Receptor Protein Tyrosine Ph...	61.5	e-174
<a href="#">gi 32313 emb CAA38662.1</a> protein-tyrosine phosphatase [Homo...	61.1	e-174
<a href="#">gi 450583 gb AA04150.1</a> protein tyrosine phosphatase >gi 4...	60.5	e-172
<a href="#">gi 6679557 ref NP_033006.1</a> protein tyrosine phosphatase, r...	60.1	e-172
<a href="#">gi 483922 gb AAA17990.1</a> protein tyrosine phosphatase alpha	59.9	e-170

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



## Example Alignment with PTP $\alpha$

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

```

Query: 1  SFSTNRKYPPLPVDKLEEE INRRMADDNKLFREEFNALPACPIQATCEAAASXXXXXXXXR 60
          SFSTNRKYPPLPVDKLEEE INRRMADDNKLFREEFNALPACPIQATCEAAS      R
Sbjct: 202 SFSTNRKYPPLPVDKLEEE INRRMADDNKLFREEFNALPACPIQATCEAAASKEENKEKNR 261

Query: 61  YVNILPYDHSRVHLTPVEGVPSDYINASFINGYQEKKNFIAAQGPKEETVNDFWRMIWE 120
          YVNILPYDHSRVHLTPVEGVPSDYINASFINGYQEKKNFIAAQGPKEETVNDFWRMIWE
Sbjct: 262 YVNILPYDHSRVHLTPVEGVPSDYINASFINGYQEKKNFIAAQGPKEETVNDFWRMIWE 321

Query: 121 QNTATIVMVTNLKERKECKCAQYWPDQGCWYGNVRSVEDVTVLVDYTVRKFCIQQVGD 180
          QNTATIVMVTNLKERKECKCAQYWPDQGCWYGNVRSVEDVTVLVDYTVRKFCIQQVGD
Sbjct: 322 QNTATIVMVTNLKERKECKCAQYWPDQGCWYGNVRSVEDVTVLVDYTVRKFCIQQVGD 381

Query: 181 VTNRKPQRLITQFHFTSWPDFGVFFPTIGMLKFLKVKACNPQYAGAI VVHCSAGVGRGTG 240
          VTNRKPQRLITQFHFTSWPDFGVFFPTIGMLKFLKVKACNPQYAGAI VVHCSAGVGRGTG
Sbjct: 382 VTNRKPQRLITQFHFTSWPDFGVFFPTIGMLKFLKVKACNPQYAGAI VVHCSAGVGRGTG 441

Query: 241 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCCMQVQTDMQYVF IYQALLEHYLYGDTELE 300
          TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCCMQVQTDMQYVF IYQALLEHYLYGDTELE
Sbjct: 442 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCCMQVQTDMQYVF IYQALLEHYLYGDTELE 501
  
```

## Guilt-by-Association: Caveats

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

## Interpretation of P-value

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

**Question: In a room with 366 people. What is the probability that one of them has the same birthday as you? And what is the probability that two of them have the same birthday?**



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



18 entries were found

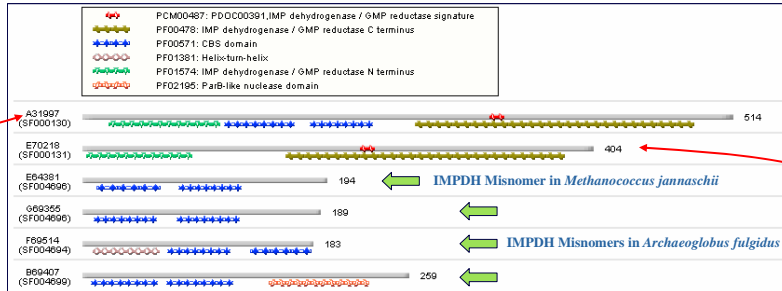
ID	Organism	PIR	Swiss-Prot/TrEMBL	RefSeq/GenPept
NF00181857	Methanococcus jannaschii	E64351 conserved hypothetical protein MJ0653	Y853_MET1A Hypothetical protein MJ0653	g1322300 inosine-5'-monophosphate dehydrogenase (guaB) NP_247671 inosine-5'-monophosphate dehydrogenase (guaB)
NF00187788	Archaeoglobus fulgidus	G69355 MJ0653 homolog AF0847 ALT_NAME5: inosine-monophosphate dehydrogenase (guaB-1) homolog [misnomer]	G28411 INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-1)	g2649754 inosine monophosphate dehydrogenase (guaB-1) NP_069631 inosine monophosphate dehydrogenase (guaB-1)
NF00188267	Archaeoglobus fulgidus	E69511 yhcV homolog 2 ALT_NAME5: inosine-monophosphate dehydrogenase (guaB-2) homolog [misnomer]	G28149 INOSINE MONOPHOSPHATE DEHYDROGENASE (GUAB-2)	g2652310 inosine monophosphate dehydrogenase (guaB-2) NP_070743 inosine monophosphate dehydrogenase (guaB-2)
NF00188697	Archaeo			phosphate rve inophosphate rve
NF00197776	Thermo			inophosphate d protein monophosphate d protein
NF00414702	Methanothermobacter thermautotrophicus	D67093 methanothermobacter thermautotrophicus ALT_NAME5: inosine-monophosphate dehydrogenase related protein V [misnomer]	G27204 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN V	dehydrogenase related protein V NP_276334 inosine-5'-monophosphate dehydrogenase related protein V
NF00414811	Methanothermobacter thermautotrophicus	D69035 MJI232 protein homolog MTH126 ALT_NAME5: inosine-5'-monophosphate dehydrogenase related protein VII [misnomer]	G26229 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN VII	g2621166 inosine-5'-monophosphate dehydrogenase related protein VII NP_273269 inosine-5'-monophosphate dehydrogenase related protein VII
NF00414837	Methanothermobacter thermautotrophicus	E69221 MJI221-related protein MTH993 ALT_NAME5: inosine-5'-monophosphate dehydrogenase related protein IX [misnomer]	G27073 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN IX	g2622023 inosine-5'-monophosphate dehydrogenase related protein IX NP_276127 inosine-5'-monophosphate dehydrogenase related protein IX
NF00414969	Methanothermobacter thermautotrophicus	E69077 yhcV homolog 2 ALT_NAME5: inosine-monophosphate dehydrogenase related protein X [misnomer]	G27616 INOSINE-5'-MONOPHOSPHATE DEHYDROGENASE RELATED PROTEIN X	g2622027 inosine-5'-monophosphate dehydrogenase related protein X NP_276607 inosine-5'-monophosphate dehydrogenase related protein X

**A partial list of IMPdehydrogenase misnomers in complete genomes remaining in some public databases**

Lecture at Yang Ming National University, Taipei, June 2006

Copyright 2006 © Limsoon Wong

## IMPDH Domain Structure



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

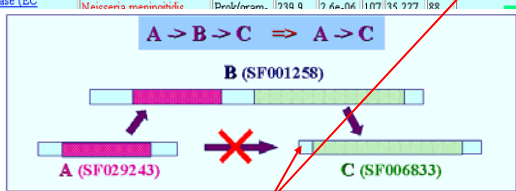
Lecture at Yang Ming National University, Taipei, June 2006

Copyright 2006 © Limsoon Wong

# Invalid Transitive Assignment

Root of invalid transitive assignment

B →	H70468	SF001258	051440	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	
	S76963	SF001258	039935	phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) [similarity]	Synechocystis sp.	Prok/gram-	557.0	5.7e-24	230	39.175	194	
	T35073	SF029243	005738	probable phosphoribosyl-AMP cyclohydrolase	Streptomyces coelicolor	Prok/gram+	399.3	3.5e-15	128	42.157	102	
	S53349	SF001257	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 →	E69493	SF029243	005738	phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) [similarity]	Archaeoglobus fulgidus	Archae	396.8	4.8e-15	108	47.778	90	
C →	G64337	SF006833	030827	phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.1) [similarity]	Methanococcus jannaschii	Archae	246.9	1.1e-06	95	36.842	95	
	D81178	SF006833	101491	phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) NMB0603 [similarity]	Nitrososphaera meningioides	Prok/gram-	730.0	7.4e-04	1107	34.277	88	
	G81925	SF006833	101491	phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) NMA0807 [similarity]								
	S51513	SF001257	001188	phosphoribosyl-AMP cyclohydrolase (EC 3.5.4.19) / phosphoribosyl-ATP pyrophosphatase (EC 3.6.1.31) / histidinol dehydrogenase (EC 1.1.1.23)								



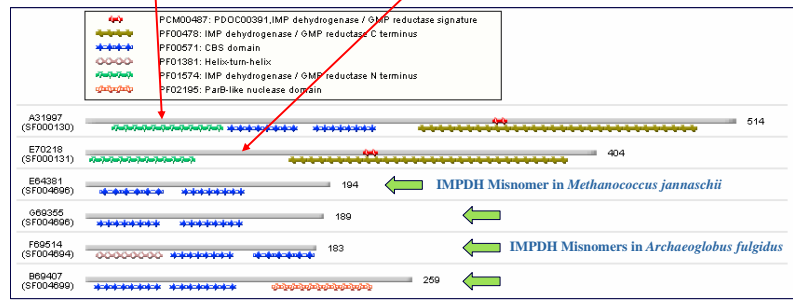
Mis-assignment of function

No IMPDH domain

# Emerging Pattern

Typical IMPDH

Functional IMPDH w/o CBS



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

# Application of Sequence Comparison: Active Site/Domain Discovery



Lecture at Yang Ming National University, Taipei, June 2006

Copyright © 2006 by Limsoon Wong

2.2

## 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**

Lecture at Yang Ming National University, Taipei, June 2006

Copyright 2006 © Limsoon Wong



## Multiple Alignment of PTPs

```

gi|126467|      FHFTSVPDFGVPFTP I GMLKFLKKVKACNP--QYAGAIVVHCSAGVGRTGTFVVIDAML D
gi|2499753|     FHFTGWPDHGVPYHATGLLSF IRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIML D
gi|462550|      YHYTQWPDMGVPEYALPVLTFVRRSSAARM--PETGPVLVHCSAGVGRTGTYIVIDSMLQ
gi|2499751|     FHFTSVPDHGVPD TDDLINFRYLVRD YMKQSPPEP SILVHCSAGVGRTGTF IAIDRLIY
gi|1709906|     FQFTA WPDHGVP EHP T PFLAFLRRVKT CNP--PDAGPMVVHCSAGVGRTGCF IVIDAMLE
gi|126471|      LHFTSVPDFGVPFTP I GMLKFLKKVKTLNP--VHAGPIVVHCSAGVGRTGTF IVIDAMMA
gi|548626|      FHFTGWPDHGVPYHATGLLSF IRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIML D
gi|131570|      FHFTGWPDHGVPYHATGLLGFVRQVKS KSP--PNAGPLVVHCSAGAGRTGCF IVIDIML D
gi|2144715|     FHFTSVPDHGVPD TDDLINFRYLVRD YMKQSPPEP SILVHCSAGVGRTGTF IAIDRLIY
..* *** **          . *                ..***** ***. . . ** ..

```

- 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

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

genome phylogenetic profiles  
profun's feature profiles  
SVM Pairwise

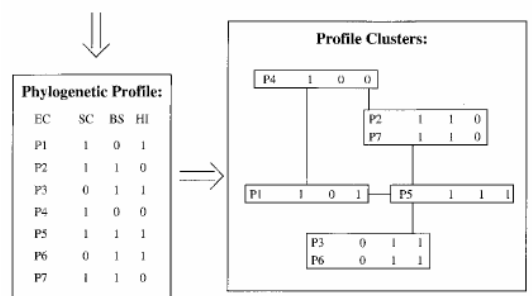
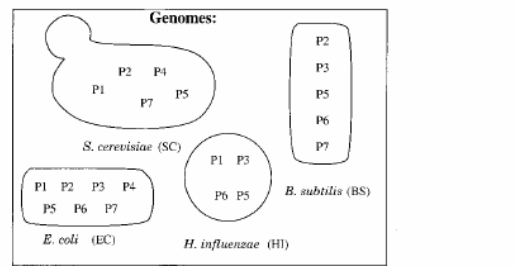




# Phylogenetic Profiling

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

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

## 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

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

where

$$\begin{aligned}
 w_z &= \binom{N}{z} \\
 \bar{w}_z &= \binom{N-z}{x-z} * \binom{N-z}{y-z} \\
 W &= \binom{N}{x} * \binom{N}{y}
 \end{aligned}$$

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

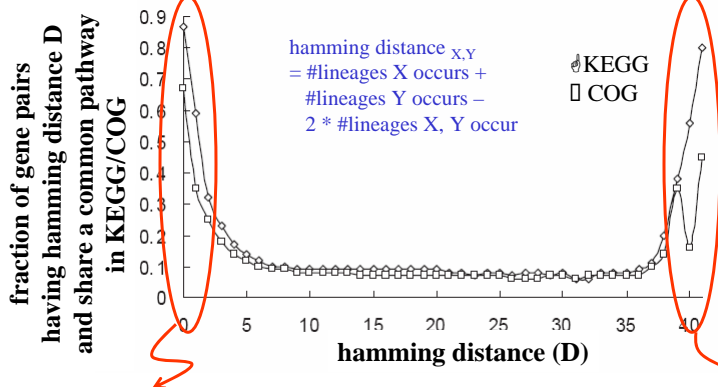
Keyword	No. of non-homologous proteins in group	No. neighbors in keyword group	No. neighbors in random group
Ribosome	60	197	27
Transcription	36	17	10
rRNA synthase and ligase	26	11	5
Membrane proteins <sup>†</sup>	25	89	5
Flagellar	21	89	3
Iron, ferric, and ferritin	19	31	2
Galactose metabolism	18	31	2
Molybdoferin and Molybdenum, and molybdoferin	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



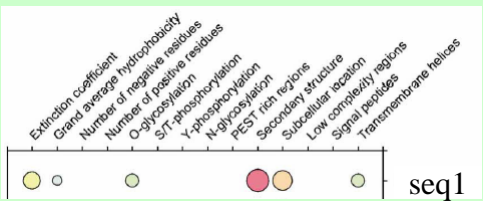
- Proteins having low hamming distance (thus highly similar phylogenetic profiles) tend to share common pathways
- Exercise: Why do proteins having high hamming distance also have this behaviour?

## The ProtFun Approach



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

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



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

# ProtFun: How it Works

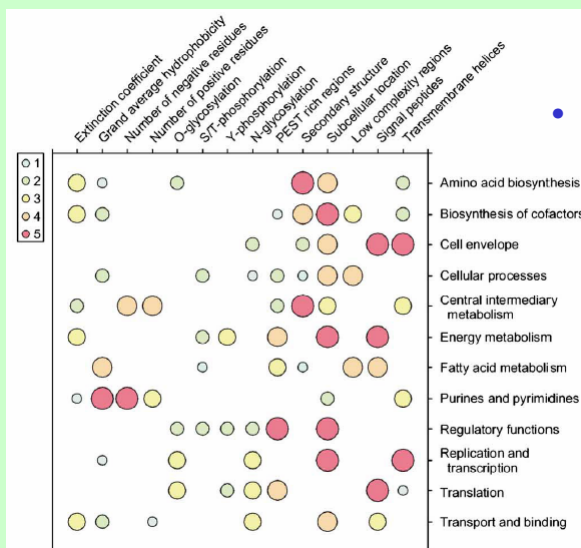
Abbreviation	Encoding	Description
ec	single value	Extinction coefficient predicted by <a href="#">ExpASy ProtParam</a>
gravy	single value	Hydrophobicity predicted by <a href="#">ExpASy ProtParam</a>
nneg	single value	Number of negatively charged residues counted by <a href="#">ExpASy ProtParam</a>
npos	single value	Number of positively charged residues counted by <a href="#">ExpASy ProtParam</a>
nglyc	potential in 5 bins	N-glycosylation sites predicted by <a href="#">NetNGlyc</a>
oglyc	potential-threshold in 10 bins	GalNAc O-glycosylations predicted by <a href="#">NetOGlyc</a>
pest	fraction in 10 bins	PEST rich regions identified by <a href="#">PESTfind</a>
phosST	potential in 10 bins	Serine and threonine phosphorylations predicted by <a href="#">NetPhos</a>
phosY	potential in 10 bins	Tyrosine phosphorylations predicted by <a href="#">NetPhos</a>
psipred	helix, sheet, coil in 5 bins	Predicted secondary structure from <a href="#">PSI-Pred</a>
psort	20 probabilities	Subcellular location predictions by <a href="#">PSORT</a>
seg	fraction in 10 bins	Low-complexity regions identified by SEG
signalp	meanS, maxY, log(cleavage pos)	Signal peptide predictions made by <a href="#">SignalP</a>
tmhmm	inside, outside, membrane in 5 bins	Transmembrane helix predictions made by <a href="#">TMHMM</a>

Extract feature profile of protein using various prediction methods

Category	Hidden units	Input features
Amino acid biosynthesis	30	ec psipred psort tmhmm
	30	ec psipred tmhmm
	30	ec netoglyc psipred psort
	30	gravy psipred psort
	30	oglyc psipred psort

Average the output of the 5 component ANNs

# ProtFun: Evidence



• Combinations of "features" seem to characterize some functional categories



## ProtFun: Example Output

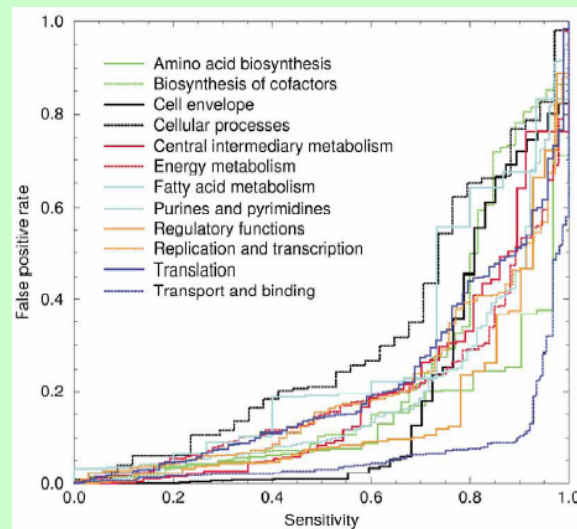
	Prion	A4	TTHY
Amino acid biosynthesis	0.011	0.011	0.011
Biosynthesis of cofactors	0.041	0.161	0.034
Cell envelope	0.146	0.804	0.698
Cellular processes	0.027	0.027	0.051
Central intermediary metabolism	0.047	0.139	0.059
Energy metabolism	0.029	0.023	0.046
Fatty acid metabolism	0.017	0.017	0.023
Purines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

- At the seq level, Prion, A4, & TTHY are dissimilar

- ProtFun predicts them to be cell envelope-related, tranport & binding

- This is in agreement w/ known functionality of these proteins

## ProtFun: Performance





# SVM-Pairwise Framework

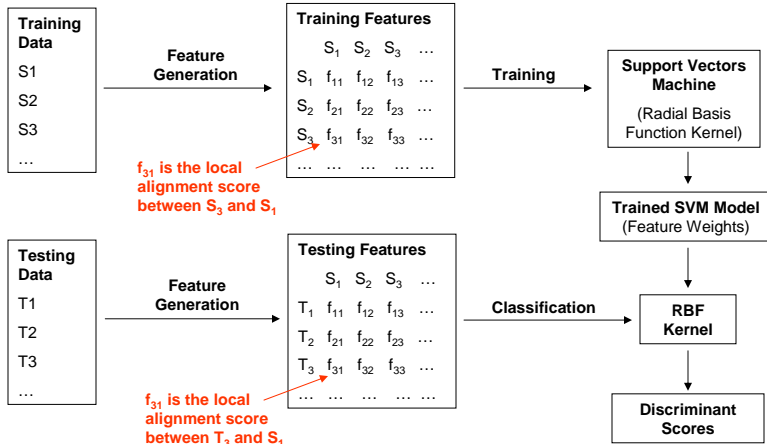
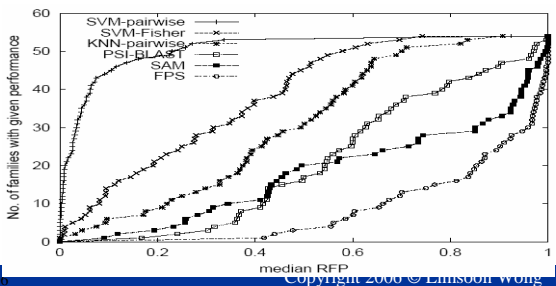
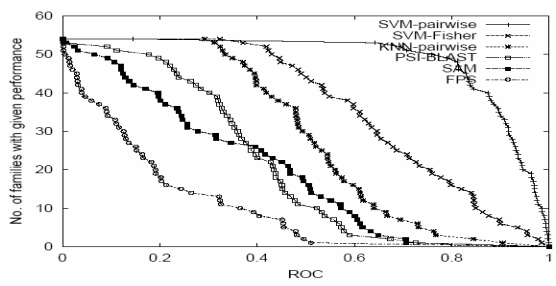


Image credit: Kenny Chua

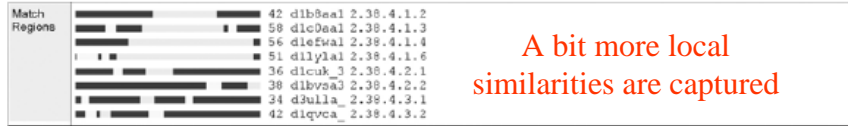
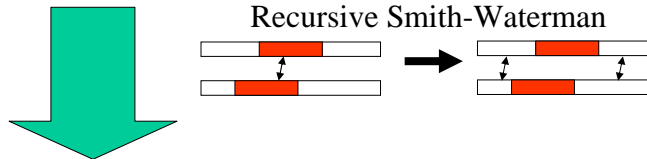
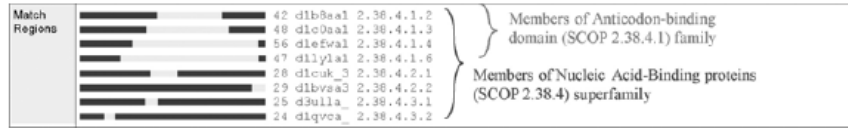


# Performance of SVM-Pairwise

- **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.

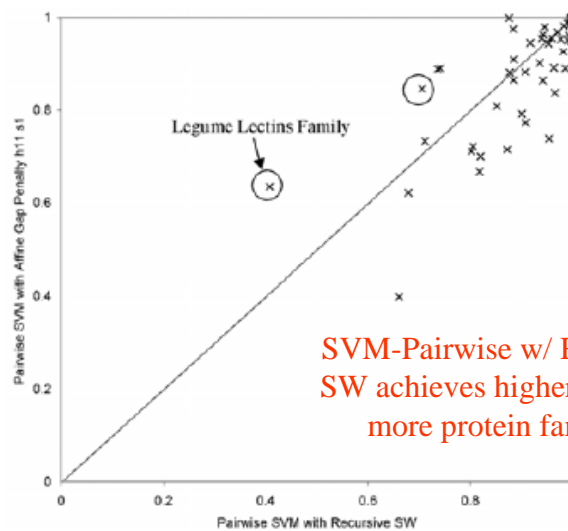


# A Refinement to Capture Multiple Local Similarities



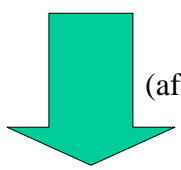
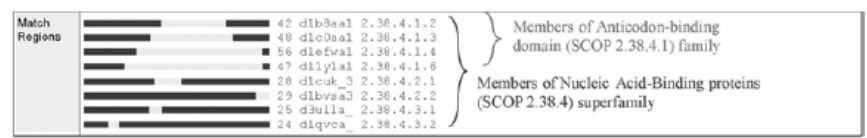
A bit more local similarities are captured

# ROC 2-D Plot of SVM-Pairwise w/ vs w/o Recursive SW

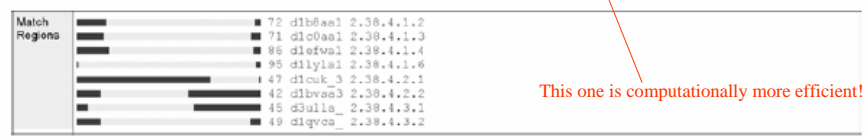


SVM-Pairwise w/ Recursive SW achieves higher ROC on more protein families

## Simpler Refinement to Capture Multiple Local Similarities

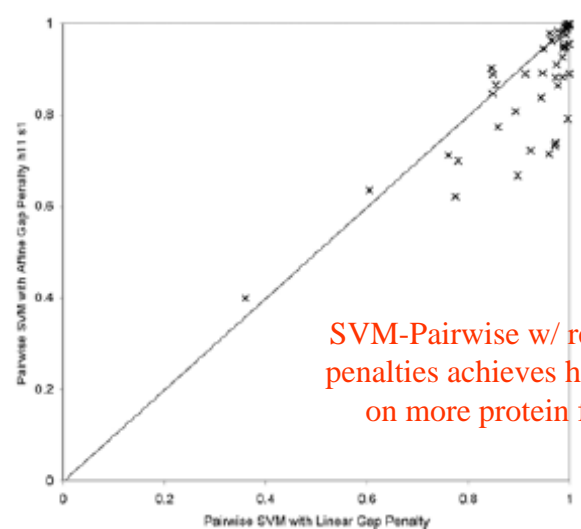


relax gap penalties  
 (affine gap penalty of open = -4, extend = -1)  
 (or linear gap penalty of -4)



This one is computationally more efficient!

## ROC 2-D Plot of SVM Pairwise w/ vs w/o Relaxed Penalties



SVM-Pairwise w/ relaxed gap penalties achieves higher ROC on more protein families

# 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

```
>g1|0000|PTPA-D2  
EEEFKILTSIKIONDKIKTGMLPANIKIKNVLQIIPYEFHWIIPVIGAGEIDTDYNASF  
IDGYRQKDSYIASQQLLETIEDFURNIIEWESCSIVELTELEERQQRCAQYTPSDOLV  
SYGDITVELKEEKECESTTVRDLVYNTREKESRQIQEFYKQWPEYQIPSDGKQKLSII  
AAVQRQQQSSQWEPITVBCSAGAQRTOTYFCALSTVLERVKAEQILDYVQTVICLRLQRPK  
EYQTLKQYEFCTYKVVQETIDAFSDYANFK
```

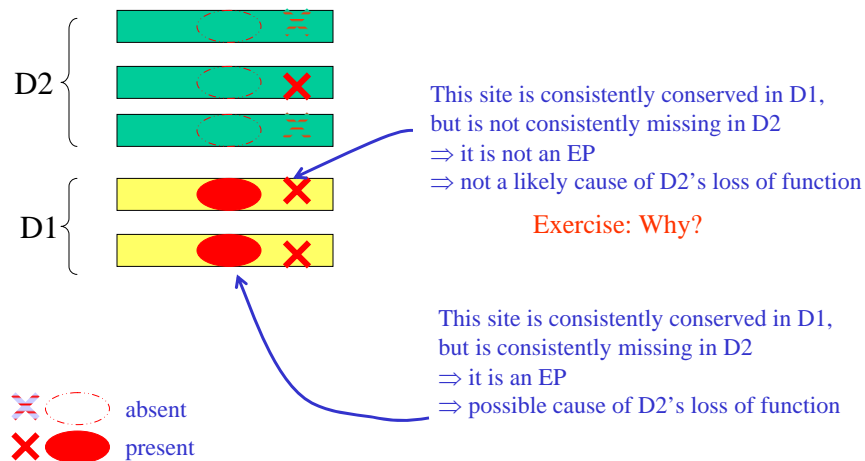
- Some PTPs have 2 PTP domains
- PTP domain D1 is has much more activity than PTP domain D2
- Why? And how do you figure that out?

## 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



## Key Mutation Site: PTP D1 vs D2

```

      ? ! ?           ?           ?           ? ??
gi|00000|P D2 QFHFHGWPVEVGIPSDGKMISIIAAVQKQQQ--SGNHPITVHCSAGAGRTGTFPCALSTVL
gi|126467| QFHFTSWPDFGVVFTPIGMLKFLKKVKACNP--QYAGAIVVHCSAGVGRGTFFVVIDAML
gi|2499753 QFHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCVIVIDIML
gi|462550| QYHYTQWPDMGVPEYALPVLTFVRRSSAARM--PETGPVLVHCSAGVGRGTGYIVIDSML
gi|2499751 QFHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPEPILVHCSAGVGRGTGTFIAIDRLI
gi|1709906 D1 QFQFTAWPDHGVPYHATGLLSFIRRVKLSNP--PDAGPMVVHCSAGVGRGTGCFIVIDAML
gi|126471| QLHFTSWPDFGVVFTPIGMLKFLKKVKTLP--VHAGPIVVHCSAGVGRGTGTFIVIDAMM
gi|548626| QFHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCVIVIDIML
gi|131570| QFHFTGWPDHGVPYHATGLLGFVRQVKSASP--PNAGPLVVHCSAGAGRTGCFIVIDIML
gi|2144715 QFHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPEPILVHCSAGVGRGTGTFIAIDRLI
      * ..  ** .*. *

```

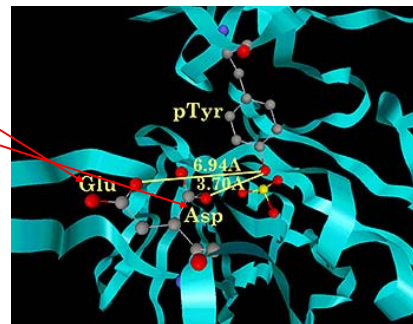
- Positions marked by “!” and “?” are likely places responsible for reduced PTP activity
  - All PTP D1 agree on them
  - All PTP D2 disagree on them

## Key Mutation Site: PTP D1 vs D2

```

      ? ! ?
gi|00000|P D2 QFHFHGWPVEVGIPSDGK
gi|126467| QFHFTSWPDFGVVFTPI
gi|2499753 QFHFTGWPDHGVPYHAT
gi|462550| QYHYTQWPDMGVPEYAL
gi|2499751 QFHFTSWPDHGVPDTTD
gi|1709906 D1 QFQFTAWPDHGVPYHAT
gi|126471| QLHFTSWPDFGVVFTPI
gi|548626| QFHFTGWPDHGVPYHAT
gi|131570| QFHFTGWPDHGVPYHAT
gi|2144715 QFHFTSWPDHGVPDTTD
      * ..  ** .*. *

```



- 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

Any Questions?





## 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



## 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. [A better gap penalty for pairwise SVM](#). Proc. APBC05, pages 11-20
- T. Jaakkola, M. Diekhans, and D. Haussler. A discriminative framework for detecting remote homologies. *JCB*, 7(1-2):95—11, 2000