

## Programming for Biology Protein Evolution / Similarity Searching

### What BLAST Does / Why BLAST works

Bill Pearson  
[wrp@virginia.edu](mailto:wrp@virginia.edu)

1

### *Sequence Similarity - Conclusions*

- *Homologous* sequences share a common ancestor, but most sequences are *non-homologous*
- Always compare Protein Sequences
- Sequence Homology can be reliably inferred from statistically significant similarity (non-homology cannot from non-similarity)
- Homologous proteins share common structures, but not necessarily common functions
- Sequence statistical significance estimates are accurate (verify this yourself)  $10^{-6} < E() < 10^{-3}$  is statistically significant
- Scoring matrices set evolutionary look back horizons - not every discovery is distant
- PSI-BLAST can be more sensitive, but with lower statistical accuracy

2

*Establishing homology from  
statistically significant similarity*

**Why BLAST works**

- For most proteins, homologs are easily found over long evolutionary distances (500 My – 2 By) using standard approaches (BLAST, FASTA)
- Difficult for distant relationships or very short domains
- Most default search parameters are optimized for distant relationships and work well

3

**This talk is not about:**

- *Alignment*
  - Alignment quality may be more sensitive to parameter choice
  - Multiple sequences for biologically accurate alignments
- *Inferring Protein Function*
  - Homology (common ancestry) implies common structure (guaranteed), not necessarily common function
  - Homologs have different functions
  - Non-homologs have similar (or identical) functions
- *The best sequences for building evolutionary trees*
  - Protein sequences are clearly best for establishing homology, but DNA sequences may be better for resolving recent divergence

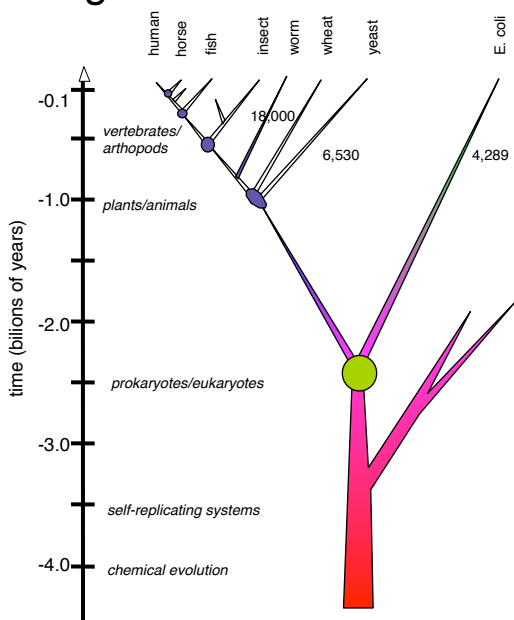
4

## Protein Evolution and Sequence Similarity

- What is Homology and how do we recognize it?
- How do we measure sequence similarity – alignments and scoring matrices?
- DNA vs protein comparison
- Alignment Algorithms/Local sequence alignments
- Similarity scoring matrices
- When are we certain that an alignment is significant - similarity score statistics?
- When to trust similarity statistics?
- Improving sensitivity with PSI-BLAST

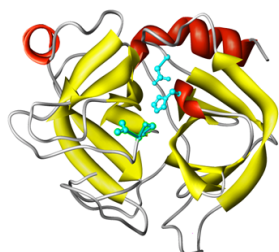
5

## Homologues share a common ancestor



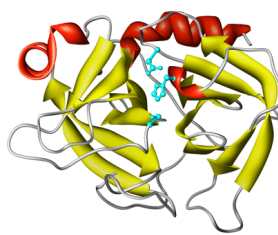
6

## When do we infer homology?

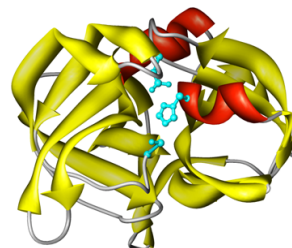


Bovine trypsin (5ptp)  
 Structure:  $E() < 10^{-23}$ ;  
 RMSD 0.0 Å  
 Sequence:  $E() < 10^{-84}$   
 100% 223/223

Homology  $\Leftrightarrow$  structural similarity  
 ? sequence similarity



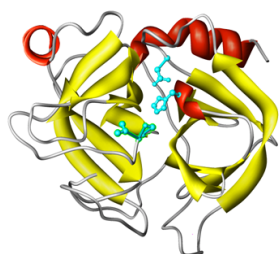
S. griseus trypsin (1sgt)  
 $E() < 10^{-14}$  RMSD 1.6 Å  
 $E() < 10^{-19}$  36%; 226/223



S. griseus protease A (2sga)  
 $E() < 10^{-4}$ ; RMSD 2.6 Å  
 $E() < 2.6$  25%; 199/181

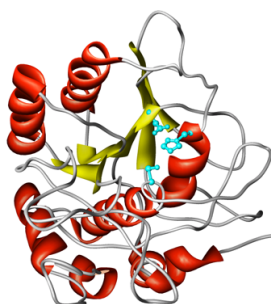
7

## When can we infer non-homology?



Bovine trypsin (5ptp)  
 Structure:  $E() < 10^{-23}$   
 RMSD 0.0 Å  
 Sequence:  $E() < 10^{-84}$   
 100% 223/223

Non-homologous proteins have  
 different structures



Subtilisin (1sbt)  
 $E() > 100$   
 $E() < 280$ ; 25% 159/275



Cytochrome c4 (1etp)  
 $E() > 100$   
 $E() < 5.5$ ; 23% 171/190

8

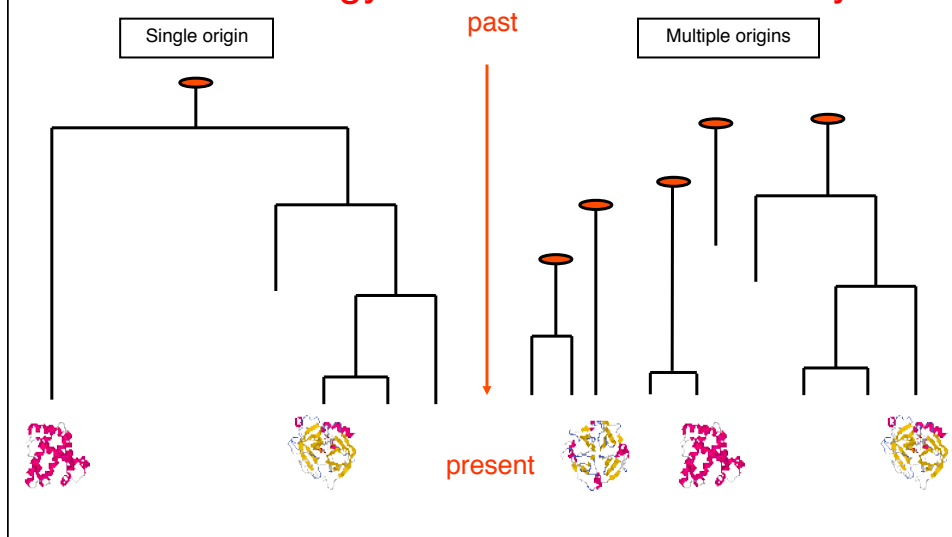
## Homology is confusing I: Homology defined Three(?) Ways

- Proteins/genes/DNA that share a common ancestor
- Specific positions/columns in a multiple sequence alignment that have a 1:1 relationship over evolutionary history
  - sequences are *50% homologous* ???
- Specific (morphological/functional) characters that share a recent divergence (clade)
  - bird/bat/butterfly wings are/are not homologous

9

## Homology is confusing II: Are All Sequences Homologous?

**No Homology without excess similarity**



## Homology from (sequence/structure) similarity

- Sequences are inferred to share a common ancestor based on statistically significant *excess* similarity. Any evidence of *excess* similarity can be used to infer homology
- Lack of evidence *cannot* be used to infer non-homology.
  - Proteins with different structures are non-homologous
- There are always two alternative hypotheses: homology (common ancestry), or independence – one must weigh the evidence for each hypothesis (independence is the *null* hypothesis).

11

## What BLAST does:

Similarity  $\overset{?}{\rightleftharpoons}$  Homology

## Why BLAST works:

Statistical  $\overset{?}{\rightleftharpoons}$  Biological  
Significance  $\rightleftharpoons$  Significance

Divergence  $\overset{?}{\rightleftharpoons}$  Convergence

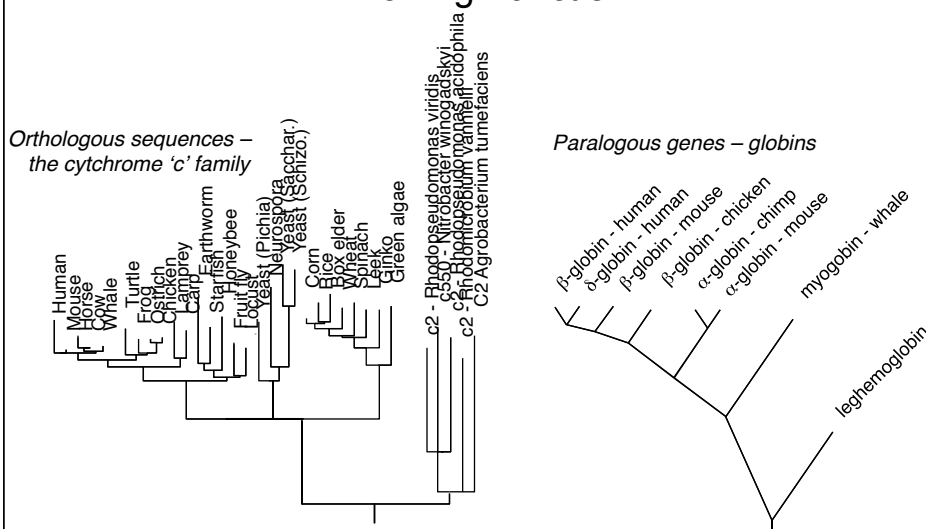
12

### E. coli proteins vs Human – Ancient Protein Domains

expect	%_id	alen	E coli descr	Human descr	sp_name
2.7e-206	53.8	944	glycine decarboxylase, P	Glycine dehydrogenase [de	GCSP_HUMAN
1.2e-176	59.5	706	methylmalonyl-CoA mutase	Methylmalonyl-CoA mutase,	MUTA_HUMAN
3.8e-176	50.6	803	glycogen phosphorylase [E	Glycogen phosphorylase, 1	PHS1_HUMAN
9.9e-173	55.6	1222	B12-dependent homocystein	5-methyltetrahydrofolate-	METH_HUMAN
1.8e-165	41.8	1031	carbamoyl-phosphate synth	Carbamoyl-phosphate synth	CPSM_HUMAN
5.6e-159	65.7	542	glucosephosphate isomeras	Glucose-6-phosphate isome	G6PI_HUMAN
8.1e-143	53.7	855	aconitate hydratase 1 [Esch	Iron-responsive element b	IRE1_HUMAN
2.5e-134	73.0	459	membrane-bound ATP syntha	ATP synthase beta chain,	ATPB_HUMAN
3.3e-121	55.8	550	succinate dehydrogenase,	Succinate dehydrogenase [	DHSA_HUMAN
1.5e-113	60.6	401	putative aminotransferase	Cysteine desulfurase, mit	NFS1_HUMAN
4.4e-111	60.9	460	fumarase C= fumarate hydr	Fumarate hydratase, mitoc	FUMH_HUMAN
1.5e-109	56.1	474	succinate-semialdehyde de	Succinate semialdehyde de	SSDH_HUMAN
3.6e-106	44.7	789	maltodextrin phosphorylas	Glycogen phosphorylase, m	PHS2_HUMAN
1.4e-102	53.1	484	NAD+-dependent betaine al	Aldehyde dehydrogenase, E	DHAG_HUMAN
3.8e-98	53.0	449	pyridine nucleotide trans	NAD(P) transhydrogenase,	NNTM_HUMAN
5.8e-96	49.9	489	glycerol kinase [Escheric	Glycerol kinase, testis s	GKP2_HUMAN
2.1e-95	66.8	328	glyceraldehyde-3-phospat	Glyceraldehyde 3-phospat	G3P2_HUMAN
5.0e-91	62.5	368	alcohol dehydrogenase cla	Alcohol dehydrogenase cla	ADHX_HUMAN
6.7e-91	56.5	393	protein chain elongation	Elongation factor Tu, mit	EFTU_HUMAN
9.5e-91	56.6	392	protein chain elongation	Elongation factor Tu, mit	EFTU_HUMAN
2.2e-89	59.1	369	methionine adenosyltransf	S-adenosylmethionine synt	METK_HUMAN
6.5e-88	53.3	422	enolase [Escherichia coli	Alpha enolase (2-phospho-	ENOA_HUMAN
9.2e-88	43.3	536	NAD-linked malate dehydro	NADP-dependent malic enzy	MAOX_HUMAN
7.3e-86	55.5	389	2-amino-3-ketobutyrate Co	2-amino-3-ketobutyrate co	KBL_HUMAN
5.2e-83	44.4	543	degrades sigma32, integra	AFG3-like protein 2 (Para	AF32_HUMAN

13

### Orthologs and Paralogs – Inferring Function

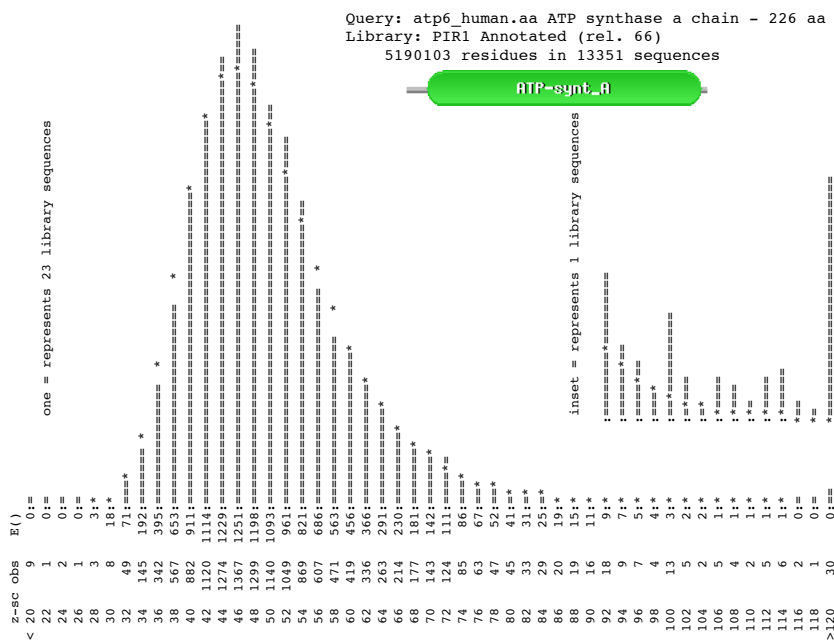


14

## Protein Evolution and Sequence Similarity

- What is Homology and how do we recognize it?
- How do we measure sequence similarity – alignments and scoring matrices?
- DNA vs protein comparison
- Alignment Algorithms/Local sequence alignments
- Similarity scoring matrices
- When are we certain that an alignment is significant - similarity score statistics?
- When to trust similarity statistics?
- Improving sensitivity with PSI-BLAST

15



16



## Inferring Homology from Statistical Significance

- Real **UNRELATED** sequences have similarity scores that are indistinguishable from **RANDOM** sequences
- If a similarity is NOT **RANDOM**, then it must be NOT **UNRELATED**
- Therefore, NOT **RANDOM** (statistically significant) similarity must reflect **RELATED** sequences

17

```

Query: atp6_human.aa ATP synthase a chain - 226 aa
Library: 5190103 residues in 13351 sequences
The best scores are:
      ( len)  s-w bits E(13351)  %_id  %_sim  alen
sp|P00846|ATP6_HUMAN ATP synthase a chain (AT ( 226) 1400 325.8 5.8e-90 1.000 1.000 226
sp|P00847|ATP6_BOVIN ATP synthase a chain (AT ( 226) 1157 270.5 2.5e-73 0.779 0.951 226
sp|P00848|ATP6_MOUSE ATP synthase a chain (AT ( 226) 1118 261.7 1.2e-70 0.757 0.916 226
sp|P00849|ATP6_XENLA ATP synthase a chain (AT ( 226) 745 176.8 4.0e-45 0.533 0.847 229
sp|P00851|ATP6_DROYA ATP synthase a chain (AT ( 224) 473 115.0 1.7e-26 0.378 0.721 222
sp|P00854|ATP6_YEAST ATP synthase a chain pre ( 259) 428 104.7 2.3e-23 0.353 0.694 232
sp|P00852|ATP6_EMENI ATP synthase a chain pre ( 256) 365 90.4 4.8e-19 0.304 0.691 230
sp|P14862|ATP6_COCHE ATP synthase a chain (AT ( 257) 353 87.7 3.2e-18 0.313 0.650 214
sp|P68526|ATP6_TRITI ATP synthase a chain (AT ( 386) 309 77.6 5.1e-15 0.289 0.651 235
sp|P05499|ATP6_TOBAC ATP synthase a chain (AT ( 395) 309 77.6 5.2e-15 0.283 0.635 233
sp|P07925|ATP6_MAIZE ATP synthase a chain (AT ( 291) 283 71.7 2.3e-13 0.311 0.667 180
sp|P0AB98|ATP6_ECOLI ATP synthase a chain (AT ( 271) 178 47.9 3.2e-06 0.233 0.585 236
sp|P0C2Y5|ATPI_ORYSA Chloroplast ATP synth (A ( 247) 144 40.1 0.00062 0.242 0.580 231
sp|P06452|ATPI_PEA Chloroplast ATP synthase a ( 247) 143 39.9 0.00072 0.250 0.586 232
sp|P27178|ATP6_SYNY3 ATP synthase a chain (AT ( 276) 142 39.7 0.00095 0.265 0.571 170
sp|P06451|ATPI_SPIOL Chloroplast ATP synthase ( 247) 138 38.8 0.0016 0.242 0.580 231
sp|P08444|ATP6_SYN6 ATP synthase a chain (AT ( 261) 127 36.3 0.0095 0.263 0.557 167
sp|P69371|ATPI_ATRBE Chloroplast ATP synthase ( 247) 126 36.0 0.01 0.221 0.571 231
sp|P06289|ATPI_MARPO Chloroplast ATP synthase ( 248) 126 36.0 0.011 0.240 0.575 167
sp|P30391|ATPI_EUGGR Chloroplast ATP synthase ( 251) 123 35.4 0.017 0.257 0.579 214

sp|P19568|TLCA_RICPR ADP,ATP carrier protein ( 498) 122 35.0 0.043 0.243 0.579 152

sp|P24966|CYB_TAYTA Cytochrome b ( 379) 113 33.0 0.13 0.234 0.532 158
sp|P03892|NU2M_BOVIN NADH-ubiquinone oxidored ( 347) 107 31.7 0.31 0.261 0.479 211
sp|P68092|CYB_STEAT Cytochrome b ( 379) 104 31.0 0.54 0.277 0.547 137
sp|P03891|NU2M_HUMAN NADH-ubiquinone oxidored ( 347) 103 30.8 0.58 0.201 0.537 149
sp|P00156|CYB_HUMAN Cytochrome b ( 380) 102 30.5 0.74 0.268 0.585 205
sp|P15993|AROP_ECOLI Aromatic amino acid tr ( 457) 103 30.7 0.78 0.234 0.622 111
sp|P24965|CYB_TRANA Cytochrome b ( 379) 101 30.3 0.87 0.234 0.563 158
sp|P29631|CYB_POMTE Cytochrome b ( 308) 99 29.9 0.95 0.274 0.584 113
sp|P24953|CYB_CAPHI Cytochrome b ( 379) 99 29.8 1.2 0.236 0.564 140

```

18

**ATP-synt\_0**

```

>>sp|P0AB98|ATP6_ECOLI ATP synthase a chain (ATPase protein 6) g (271 aa)
s-w opt: 178 Z-score: 218.2 bits: 47.9 E(): 3.2e-06
Smith-Waterman score: 178; 23.3% identity (58.5% similar) in 236 aa overlap (8-222:45-264)

          10          20          30          40
human      MNENLFASFIAPTILGLPAVLIILFPPLLIPTSKYLINNRITITQQ
E coli     NMTPQDYIGHHLNNLQLDLRTRFSLVDQNPFPATFWTINIDSMFFSVVLGL---LFLVLFRSVAKKATSG--VPGKFQTAIE
             10      20      30      40      50      60      70      80

          50          60          70          80          90          100          110
human      WLIKLTSKQMMTMHNTKGRWLSMLVSLIIFIATTNLLGLLP-----HSF-----TPTTQLSMNLAMAIPWAG
E coli     LVIGFVNGSVKDMYHGKSKLIAPLALTI FVWVFLMNLMDLLPIDLLPYIAEHVGLPALRVVPSADVNVVTLSMALGVF--
             90      100     110     120     130     140     150

          120          130          140          150          160          170          180
human      TVIMGFRSKIKNALAHFLPQGTPTPL----IPMLVIIETISLLIQPMALAVRLTANITAGHLLMHLIGSATLAMSTINL
E coli     -ILILFYSIKMKGIGGGFTKELTQPFNHWFIPVNLILEGVLSLLSKPVSLLGRLFGNMYAGELIFILIAGLLPWWWSQWIL
             160     170     180     190     200     210     220     230

          190          200          210          220
human      PSTLIIIFTILILLTILEIAVALIQAYVFTLLVLSLYLHDNT
E coli     NVPWAI FHLIIIT-----LQAFIFMVLTIIVLSMASEEH
             240     250               260     270
    
```

19

### The PAM250 matrix

Cys	12																								
Ser	0	2																							
Thr	-2	1	3																						
Pro	-1	1	0	6																					
Ala	-2	1	1	1	2																				
Gly	-3	1	0	-1	1	5																			
Asn	-4	1	0	-1	0	0	2																		
Asp	-5	0	0	-1	0	1	2	4																	
Glu	-5	0	0	-1	0	0	1	3	4																
Gln	-5	-1	-1	0	0	-1	1	2	2	4															
His	-3	-1	-1	0	-1	-2	2	1	1	3	6														
Arg	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6													
Lys	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5												
Met	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6											
Ile	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5										
Leu	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6									
Val	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	5								
Phe	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9							
Tyr	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10						
Trp	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	1					
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W					

20

```

>>sp|P30391|ATPI_EUGGR Chloroplast ATP synthase a chain precursor (251 aa)
s-w opt: 123 Z-score: 151.3 bits: 35.4 E(): 0.017
Smith-Waterman score: 123; 25.7% identity (57.9% similar) in 214 aa overlap (21-222:50-243)

          10          20          30          40          50          60
human      MNENLFASFIAPTILGLPAAVLIIILFPPLLIPTSKYLINNRLLITQQWLKIKLTSKQMMTM
          .:. . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Euglena VNMFISGIFQIANVEVGQHFYWSILGFQIHGQVLINSWIVILLIIGF--LSIYTKNL--TLVPANKQIFIELVTEFITDI
          10          20          30          40          50          60          70          80

          70          80          90          100         110         120
human      HNTK-GRT---NSLMLVSLIIFIATTNLLG-LLPHSFT--PTTQL---SMNLAMAIPLWAGTVIMGFRSKI-KNALAHF
          .:. . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Euglena  SKTQIGKEYSKWVPYIGTMFLFIFVSNWSGALIPWKIIELPNGLGAPTNDINTAGLAILTSLAYFYAGLNKKGGLTYF
          90          100         110         120         130         140         150         160

          130         140         150         160         170         180         190         200
Human     LPQGTPTPLIPMLVVIETISLLIQPMALAVRLTANITAGHLLMHLIGSATLAMSTINLPSTLIIFTILILLTILEIAVAL
          .:. . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Euglena  KKYVQPTPILLPINILEDFT---KPLSLSFRLFGNILADELVVAVLVSL-----VP--LIVPVPLIFLGLF---TSG
          170         180         190         200         210         220

          210         220
human     IQAYVFTLLVSLYLDHNT
          .:. . . . . . . . . .
Euglena  IQALIFATLSGSGYIGEAMEGHH
          230         240         250
    
```

```

Query: atp6 human.aa ATP synthase a chain - 226 aa
Library: 5190103 residues in 13351 sequences

The best scores are:
(sp|P00846|ATP6_HUMAN ATP synthase a chain (AT (226) 1400 325.8 5.8e-90 1.000 1.000 226)
(sp|P00847|ATP6_BOVIN ATP synthase a chain (AT (226) 1157 270.5 2.5e-73 0.779 0.951 226)
(sp|P00848|ATP6_MOUSE ATP synthase a chain (AT (226) 1118 261.7 1.2e-70 0.757 0.916 226)
(sp|P00849|ATP6_XENLA ATP synthase a chain (AT (226) 745 176.8 4.0e-45 0.533 0.847 229)
(sp|P00851|ATP6_DROYA ATP synthase a chain (AT (224) 473 115.0 1.7e-26 0.378 0.721 222)
(sp|P00854|ATP6_YEAST ATP synthase a chain pre (259) 428 104.7 2.3e-23 0.353 0.694 232)
(sp|P00852|ATP6_EMENI ATP synthase a chain pre (256) 365 90.4 4.8e-19 0.304 0.691 230)
(sp|P14862|ATP6_COCHE ATP synthase a chain (AT (257) 353 87.7 3.2e-18 0.313 0.650 214)
(sp|P68526|ATP6_TRITI ATP synthase a chain (AT (386) 309 77.6 5.1e-15 0.289 0.651 235)
(sp|P05499|ATP6_TOBAC ATP synthase a chain (AT (395) 309 77.6 5.2e-15 0.283 0.635 233)
(sp|P07925|ATP6_MAIZE ATP synthase a chain (AT (291) 283 71.7 2.3e-13 0.311 0.667 180)
(sp|P0AB98|ATP6_ECOLI ATP synthase a chain (AT (271) 178 47.9 3.2e-06 0.233 0.585 236)
(sp|P0C2Y5|ATPI_ORYSA Chloroplast ATP synth (A (247) 144 40.1 0.00062 0.242 0.580 231)
(sp|P06452|ATPI_PEA Chloroplast ATP synthase a (247) 143 39.9 0.00072 0.250 0.586 232)
(sp|P27178|ATP6_SYNY3 ATP synthase a chain (AT (276) 142 39.7 0.00095 0.265 0.571 170)
(sp|P06451|ATPI_SPIOL Chloroplast ATP synthase (247) 138 38.8 0.0016 0.242 0.580 231)
(sp|P08444|ATP6_SYN6 ATP synthase a chain (AT (261) 127 36.3 0.0095 0.263 0.557 167)
(sp|P69371|ATPI_ATRBE Chloroplast ATP synthase (247) 126 36.0 0.01 0.221 0.571 231)
(sp|P06289|ATPI_MARPO Chloroplast ATP synthase (248) 126 36.0 0.011 0.240 0.575 167)
(sp|P30391|ATPI_EUGGR Chloroplast ATP synthase (251) 123 35.4 0.017 0.257 0.579 214)

sp|P19568|TLCA_RICPR ADP,ATP carrier protein (498) 122 35.0 0.043 0.243 0.579 152

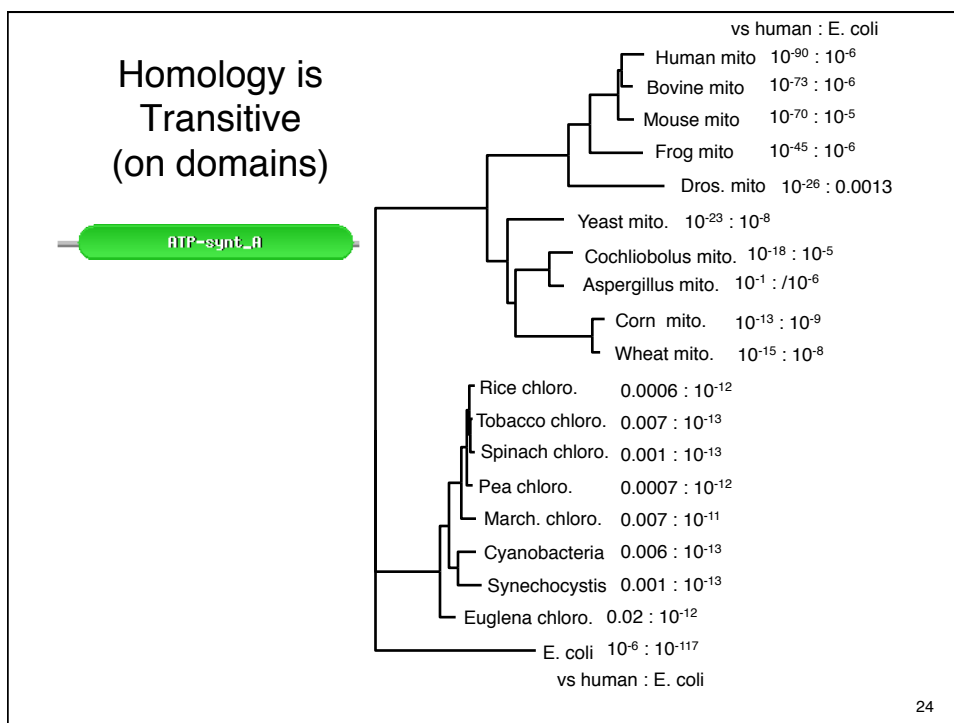
sp|P24966|CYB_TAYTA Cytochrome b (379) 113 33.0 0.13 0.234 0.532 158
sp|P03892|NU2M_BOVIN NADH-ubiquinone oxidored (347) 107 31.7 0.31 0.261 0.479 211
sp|P68092|CYB_STEAT Cytochrome b (379) 104 31.0 0.54 0.277 0.547 137
sp|P03891|NU2M_HUMAN NADH-ubiquinone oxidored (347) 103 30.8 0.58 0.201 0.537 149
sp|P00156|CYB_HUMAN Cytochrome b (380) 102 30.5 0.74 0.268 0.585 205
sp|P15993|AROP_ECOLI Aromatic amino acid tr (457) 103 30.7 0.78 0.234 0.622 111
sp|P24965|CYB_TRANA Cytochrome b (379) 101 30.3 0.87 0.234 0.563 158
sp|P29631|CYB_POMTE Cytochrome b (308) 99 29.9 0.95 0.274 0.584 113
sp|P24953|CYB_CAPHI Cytochrome b (379) 99 29.8 1.2 0.236 0.564 140
    
```

Query: atp6\_ecoli.aa ATP synthase a - 271 aa  
Library: 5190103 residues in 13351 sequences

The best scores are:

	( len)	s-w bits	E(13351)	%_id	%_sim	alen	
sp P0AB98 ATP6_ECOLI	ATP synthase a chain (AT ( 271)	1774	416.8	3.e-117	1.000	1.000	271
sp P06451 ATPI_SPIOL	Chloroplast ATP synthase ( 247)	274	70.4	5.8e-13	0.270	0.616	211
sp P69371 ATPI_ATRBE	Chloroplast ATP synthase ( 247)	271	69.7	9.3e-13	0.270	0.607	211
sp P08444 ATP6_SYNP6	ATP synthase a chain (AT ( 261)	271	69.7	9.9e-13	0.267	0.600	240
sp P06452 ATPI_PEA	Chloroplast ATP synthase a ( 247)	266	68.5	2.1e-12	0.274	0.614	223
<b>sp P30391 ATPI_EUGGR</b>	<b>Chloroplast ATP synthase ( 251)</b>	<b>265</b>	<b>68.3</b>	<b>2.5e-12</b>	<b>0.298</b>	<b>0.596</b>	<b>225</b>
sp P0C2Y5 ATPI_ORYSA	Chloroplast ATP synthase ( 247)	260	67.2	5.4e-12	0.259	0.603	239
sp P27178 ATP6_SYNY3	ATP synthase a chain (AT ( 276)	260	67.1	6.1e-12	0.264	0.578	258
sp P06289 ATPI_MARPO	Chloroplast ATP synthase ( 248)	250	64.8	2.7e-11	0.261	0.621	211
sp P07925 ATP6_MAIZE	ATP synthase a chain (AT ( 291)	215	56.7	8.7e-09	0.259	0.578	232
sp P68526 ATP6_TRITI	ATP synthase a chain (AT ( 386)	209	55.3	3.1e-08	0.259	0.603	239
sp P00854 ATP6_YEAST	ATP synthase a chain pre ( 259)	204	54.2	4.5e-08	0.235	0.578	277
sp P05499 ATP6_TOBAC	ATP synthase a chain (AT ( 395)	189	50.7	7.8e-07	0.220	0.582	268
<b>sp P00846 ATP6_HUMAN</b>	<b>ATP synthase a chain (AT ( 226)</b>	<b>178</b>	<b>48.2</b>	<b>2.5e-06</b>	<b>0.237</b>	<b>0.589</b>	<b>236</b>
sp P00852 ATP6_EMENI	ATP synthase a chain pre ( 256)	178	48.2	2.8e-06	0.209	0.590	244
sp P00849 ATP6_XENLA	ATP synthase a chain (AT ( 226)	173	47.1	5.5e-06	0.261	0.630	165
sp P00847 ATP6_BOVIN	ATP synthase a chain (AT ( 226)	172	46.8	6.5e-06	0.233	0.581	236
sp P14862 ATP6_COCHE	ATP synthase a chain (AT ( 257)	171	46.6	8.7e-06	0.204	0.608	265
sp P00848 ATP6_MOUSE	ATP synthase a chain (AT ( 226)	166	45.5	1.7e-05	0.259	0.617	193
sp P00851 ATP6_DROYA	ATP synthase a chain (AT ( 224)	139	39.2	0.0013	0.225	0.549	253
sp P24962 CYB_STELO	Cytochrome b ( 379)	125	35.9	0.021	0.223	0.575	193
sp P09716 US17_HCMVA	Hypothetical protein HVL ( 293)	109	32.3	0.21	0.260	0.565	131
sp P68092 CYB_STEAT	Cytochrome b ( 379)	109	32.2	0.27	0.211	0.562	194
sp P24960 CYB_ODOHE	Cytochrome b ( 379)	104	31.1	0.61	0.210	0.555	200
sp P03887 NULM_BOVIN	NADH-ubiquinone oxidored ( 318)	98	29.7	1.3	0.287	0.545	167
sp P24992 CYB_ANTAM	Cytochrome b ( 379)	99	29.9	1.4	0.192	0.565	193

23



## Homology and Domains – Histone deacetylase PCAF

The best scores are:

		s-w bits	E(362341)	%_id	%_sim	alen
PCAF_HUMAN	Histone acetyltransferase PCAF;	( 832)	4876 1092	0	1.000	832
PCAF_MOUSE	Histone acetyltransferase PCAF;	( 813)	4507 1010	0	0.929	0.974 817
GCNL2_HUMAN	General control of amino acid synthesis protein 5-l	( 837)	3535 793.	0	0.716	0.864 821
GCN5_YEAST	Histone acetyltransferase GCN5	( 439)	1049 240.	5.2e-62	0.469	0.743 354
GCN5_ARATH	Histone acetyltransferase GCN5; AtGCN5	( 568)	956 219.	1.2e-55	0.435	0.733 375
BPTF_HUMAN	Nucleosome-remodeling factor subunit BPTF	(3046)	369 88.3	2.4e-15	0.495	0.773 97
NU301_DROME	Nucleosome-remodeling factor subunit NURF301	(2669)	359 86.2	9.3e-15	0.511	0.787 94
CECR2_HUMAN	Cat eye syndrome critical region protein 2	(1484)	306 74.6	1.6e-11	0.371	0.771 105
BRD4_HUMAN	Bromodomain-containing protein 4; HUNK1 protein	(1362)	288 70.6	2.3e-10	0.379	0.681 116
BRDT_MACFA	Bromodomain testis-specific protein	( 947)	270 66.7	2.3e-09	0.353	0.690 116
FSH_DROME	Homeotic protein female sterile; Fragile-chorion memb	(2038)	276 67.8	2.4e-09	0.341	0.651 129
BRDT_HUMAN	Bromodomain testis-specific protein; RING3-like prot	( 947)	266 65.9	4.3e-09	0.345	0.690 116
Y0777_DICDI	Bromodomain-containing protein DDB_G0280777	(1823)	260 64.3	2.5e-08	0.385	0.725 109
BRDT_MOUSE	Bromodomain testis-specific protein; RING3-like prot	( 956)	247 61.6	8.1e-08	0.328	0.647 116
BAZ2B_HUMAN	Bromodomain adjacent to zinc finger domain protein	(1972)	247 61.3	2e-07	0.343	0.695 105
TAF1_DROME	Transcription initiation factor TFIID subunit 1; Tra	(2129)	230 57.5	3.1e-06	0.349	0.689 106
B2_SCHPO	Bromodomain-containing protein C631.02	( 727)	217 55.0	5.9e-06	0.320	0.587 172
BRD9_XENLA	Bromodomain-containing protein 9	( 527)	214 54.5	6.2e-06	0.292	0.579 171
GTE6_ARATH	Transcription factor GTE6; Protein GENERAL TRANSCRIP	( 369)	201 51.7	2.9e-05	0.290	0.601 183
BAZ1B_MOUSE	Bromodomain adjacent to zinc finger domain protein	(1479)	212 53.7	3.1e-05	0.302	0.583 139
K2_SCHPO	Bromodomain-containing protein C1450.02	( 578)	204 52.2	3.3e-05	0.310	0.628 113
TAF1_HUMAN	Transcription initiation factor TFIID subunit 1; Tra	(1872)	212 53.6	4.2e-05	0.339	0.678 115
BAZ1B_HUMAN	Bromodomain adjacent to zinc finger domain protein	(1483)	209 53.0	5e-05	0.397	0.705 78
TIF1A_HUMAN	Transcription intermediary factor 1-alpha; TIF1-al	(1050)	206 52.5	5.1e-05	0.384	0.698 86
BDF2_YEAST	Bromodomain-containing factor 2	( 638)	200 51.3	6.9e-05	0.304	0.607 168

25

## Homology and Domains – Histone deacetylase PCAF

The best scores are:

		E(362341)	alen
PCAF_HUMAN	Histone acetyl ( 832)	0	832
GCN5_YEAST	Histone acetyl ( 439)	5.2e-62	354
BPTF_HUMAN	Nucleosome-rem (3046)	2.4e-15	97
CECR2_HUMAN	Cat eye syndr (1484)	1.6e-11	105
GTE6_ARATH	Transcription ( 369)	2.9e-05	183

26

## Protein Evolution and Sequence Similarity

- What is Homology and how do we recognize it?
- How do we measure sequence similarity – alignments and scoring matrices?
- DNA vs protein comparison
- Alignment Algorithms/Local sequence alignments
- Similarity scoring matrices
- When are we certain that an alignment is significant - similarity score statistics?
- When to trust similarity statistics?
- Improving sensitivity with PSI-BLAST

27

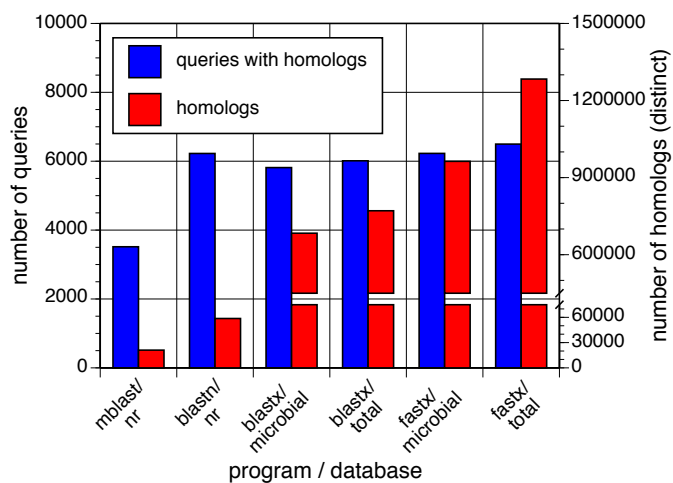
### DNA vs protein sequence comparison

The best scores are:

		DNA E(188,018)	tfastx3 E(187,524)	prot. E(331,956)
DMGST	D.melanogaster GST1-1	1.3e-164	4.1e-109	1.0e-109
MDGST1	M.domestica GST-1 gene	2e-77	3.0e-95	1.9e-76
LUCGLTR	Lucilia cuprina GST	1.5e-72	5.2e-91	3.3e-73
MDGST2A	M.domesticus GST-2 mRNA	9.3e-53	1.4e-77	1.6e-62
MDNF1	M.domestica nf1 gene. 10	4.6e-51	2.8e-77	2.2e-62
MDNF6	M.domestica nf6 gene. 10	2.8e-51	4.2e-77	3.1e-62
MDNF7	M.domestica nf7 gene. 10	6.1e-47	9.2e-77	6.7e-62
AGGST15	A.gambiae GST mRNA	3.1e-58	4.2e-76	4.3e-61
CVU87958	Culicoides GST	1.8e-41	4.0e-73	3.6e-58
AGG3GST11	A.gambiae GST1-1 mRNA	1.5e-46	2.8e-55	1.1e-43
BMO6502	Bombyx mori GST mRNA	1.1e-23	8.8e-50	5.7e-40
AGSUGST12	A.gambiae GST1-1 gene	2.3e-16	4.5e-46	5.1e-37
MOTGLUSTRA	Manduca sexta GST	5.7e-07	2.5e-30	8.0e-25
RLGSTARGN	R.leguminosarum gsta	0.0029	3.2e-13	1.4e-10
HUMGSTT2A	H. sapiens GSTT2	0.32	3.3e-10	2.0e-09
HSGSTT1	H.sapiens GSTT1 mRNA	7.2	8.4e-13	3.6e-10
ECAE000319	E. coli hypothet. prot.	—	4.7e-10	1.1e-09
MYMDCMA	Methyl. dichlorometh. DH	—	1.1e-09	6.9e-07
BCU19883	Burkholderia maleylacetate red.	—	1.2e-09	1.1e-08
NFU43126	Naegleria fowleri GST	—	3.2e-07	0.0056
SP505GST	Sphingomonas paucim	—	1.8e-06	0.0002
EN1838	H. sapiens maleylaceto. iso.	—	2.1e-06	5.9e-06
HSU86529	Human GSTZ1	—	3.0e-06	8.0e-06
SYCCPNC	Synechocystis GST	—	1.2e-05	9.5e-06
HSEF1GMR	H.sapiens EF1g mRNA	—	9.0e-05	0.00065

28

## Improving search strategies (windshield splatter metagenomics)



- always use protein/translated DNA comparisons
- smaller databases are more sensitive

## Similarity Searching II

1. What question to ask?
2. What program to use?
3. What database to search?
4. How to avoid mistakes (what to look out for)
5. When to do something different (changing scoring matrices)

## 1. What question to ask?

- Is there an homologous protein (a protein with a similar structure)?
- Does that homologous protein have a similar function?
- Does XXX genome have YYY (kinase, GPCR, ...)?

### Questions not to ask:

- Does this DNA sequence have a similar regulatory element (too short – never significant)?
- Does (non-significant) protein have a similar function/modification/antigenic site?

31

## 2. What program to run?

- What is your query sequence?
  - protein – BLAST (NCBI), SSEARCH (EBI)
  - protein coding DNA (EST) – BLASTX (NCBI), FASTX (EBI)
  - DNA (structural RNA, repeat family) – BLASTN (NCBI), FASTA (EBI)
- Does XXX genome have YYY (protein)?
  - TBLASTN YYY vs XXX genome
  - TFASTX YYY vs XXX genome
- Does my protein contain repeated domains?
  - LALIGN (UVa <http://fasta.bioch.virginia.edu>)

32



NCBI  
BLAST  
Server

blast.ncbi.nlm.nih.gov

BLAST finds regions of similarity between biological sequences. [more...](#)

**New** Aligning Multiple Protein Sequences? Try the **COBALT** Multiple Alignment Tool.

**BLAST Assembled Genomes**

Choose a species genome to search, or [list all genomic BLAST databases](#).

- [Human](#)
- [Mouse](#)
- [Rat](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Bos taurus](#)
- [Danio rerio](#)
- [Drosophila melanogaster](#)
- [Gallus gallus](#)
- [Pan troglodytes](#)
- [Microbes](#)
- [Apis mellifera](#)

**Basic BLAST**

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a <b>nucleotide</b> database using a <b>nucleotide</b> query <i>Algorithms: blastn, megablast, discontinuous megablast</i>
<a href="#">protein blast</a>	Search <b>protein</b> database using a <b>protein</b> query <i>Algorithms: blastp, psi-blast, phi-blast</i>
<a href="#">blastx</a>	Search <b>protein</b> database using a <b>translated nucleotide</b> query
<a href="#">tblastn</a>	Search <b>translated nucleotide</b> database using a <b>protein</b> query
<a href="#">tblastx</a>	Search <b>translated nucleotide</b> database using a <b>translated nucleotide</b> query

**Specialized BLAST**

Choose a type of specialized search (or database name in parentheses.)

- Make specific primers with [Primer-BLAST](#)
- Search [trace archives](#)
- Find [conserved domains](#) in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)

## NCBI BLAST Server

blast.ncbi.nlm.nih.gov

### Basic BLAST

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a <b>nucleotide</b> database using a <b>nucleotide</b> query <i>Algorithms: blastn, megablast, discontinuous megablast</i>
<a href="#">protein blast</a>	Search <b>protein</b> database using a <b>protein</b> query <i>Algorithms: blastp, psi-blast, phi-blast</i>
<a href="#">blastx</a>	Search <b>protein</b> database using a <b>translated nucleotide</b> query
<a href="#">tblastn</a>	Search <b>translated nucleotide</b> database using a <b>protein</b> query
<a href="#">tblastx</a>	Search <b>translated nucleotide</b> database using a <b>translated nucleotide</b> query

What is wrong with this picture?

Always compare protein sequences

NCBI  
BLAST  
Server

The screenshot shows the NCBI BLAST web interface. At the top, there are navigation tabs: Home, Recent Results, Saved Strategies, and Help. Below this is the 'NCBI/BLAST/blastp suite' header with sub-tabs for blastn, blastp, blastx, tblastn, and tblastx. The main section is titled 'Enter Query Sequence' and contains a large text input field for 'Enter accession number, gi, or FASTA sequence'. To the right of this field are 'Clear' and 'Query subrange' options with 'From' and 'To' input boxes. Below the text field is a file upload section with a 'Choose File' button and 'no file selected' text. There is also a 'Job Title' input field and a checkbox for 'Align two or more sequences'. The 'Choose Search Set' section includes a 'Database' dropdown menu set to 'Non-redundant protein sequences (nr)', an 'Organism' input field with an 'Exclude' checkbox, and an 'Entrez Query' input field. The 'Program Selection' section has radio buttons for 'blastp (protein-protein BLAST)', 'PSI-BLAST (Position-Specific Iterated BLAST)', and 'PHI-BLAST (Pattern Hit Initiated BLAST)'. At the bottom, there is a 'BLAST' button and a checkbox for 'Show results in a new window'.

Searching at the EBI  
www.ebi.ac.uk/Tools/sss/

The screenshot shows the EBI Sequence Similarity Searching (SSS) tool page. The page title is 'Sequence Similarity Searching' and the URL is 'www.ebi.ac.uk/Tools/sss/'. Below the title, there is a 'BLAST' section with a list of tools: 'NCBI BLAST', 'WU-BLAST', 'PSI-BLAST', 'FASTA', 'SSEARCH', 'PSI-Search', and 'GGSEARCH'. Each tool has a brief description and a list of available databases. For example, 'NCBI BLAST' is available for Protein, Nucleotide, and Vectors. 'FASTA' is available for Protein, Nucleotide, Proteomes, Genomes, and Whole Genome Shotgun. 'SSEARCH' is available for Protein, Nucleotide, Proteomes, Genomes, and Whole Genome Shotgun. 'PSI-Search' is available for ASD Protein, ASD Nucleotide, LGIC Protein, and LGIC Nucleotide. 'GGSEARCH' is available for Protein and Nucleotide. The page number '36' is visible in the bottom right corner.

## Searching at the EBI – ssearch

EBI > Tools > Similarity & Homology

**FASTA/SSEARCH/GGSEARCH/GLSEARCH - Protein Similarity Search**

Provides sequence similarity searching against protein databases using the FASTA and SSEARCH programs. **SSEARCH** does a rigorous Smith-Waterman search for similarity between a query sequence and a database. **GGSEARCH** compares a protein or DNA sequence to a sequence database producing global-global alignment (Needleman-Wunsch). **GLSEARCH** compares a protein or DNA sequence to a sequence database. **FASTA** can be very specific when identifying long regions of low similarity especially for highly diverged sequences. You can also conduct sequence similarity searching against [nucleotide databases](#) or complete [proteome/genome](#) databases using the [FASTA programs](#).

[Download Software](#)

PROGRAM	DATABASES	RESULTS	SEARCH TITLE	YOUR EMAIL
SSEARCH	Protein	Interactive	Sequence	
	UniProt Knowledgebase UniProtKB/Swiss-Prot UniProt Clusters 100% UniProt Clusters 100% (SEG filter)			
MATRIX	GAP OPEN	GAP EXTEND	EXPECTATION UPPER VALUE	EXPECTATION LOWER VALUE
BLOSUM62	-10	-2	10.0	default
SCORES	ALIGNMENTS	SEQUENCE RANGE	DATABASE RANGE	FILTER
50	50	START-END	START-END	none
				Regress

Enter or Paste a  Sequence in any format:

Upload a file:  no file selected

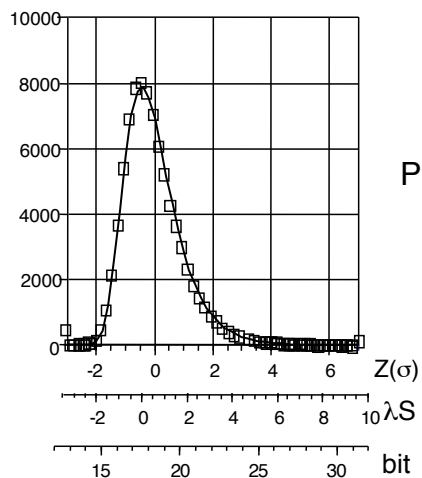
37

### 3. What database to search?

- Search the smallest comprehensive database likely to contain your protein
  - vertebrates – human proteins (40,000)
  - fungi – *S. cerevisiae* (6,000)
  - bacteria – *E. coli*, gram positive, etc. (<100,000)
- Search a richly annotated protein set (SwissProt, 450,000)
- Always search NR (> 12 million) *LAST*
- Never Search “GenBank” (DNA)

38

### Why smaller databases are better – statistics



$$S' = \lambda S_{\text{raw}} - \ln K m n$$

$$S_{\text{bit}} = (\lambda S_{\text{raw}} - \ln K) / \ln(2)$$

$$P(S' > x) = 1 - \exp(-e^{-x})$$

$$P(S_{\text{bit}} > x) = 1 - \exp(-mn2^{-x})$$

$$E(S' > x \text{ ID}) = P D$$

$$P(B \text{ bits}) = m n 2^{-B}$$

$$P(40 \text{ bits}) = 1.5 \times 10^{-7}$$

$$E(40 \mid D=4000) = 6 \times 10^{-4}$$

$$E(40 \mid D=12E6) = 1.8$$

39

### Statistical Significance and Database Size

atp6\_human vs E. coli  
 >>ref|NP\_290377.1| F0F1 ATP synthase subunit [E. coli] (271 aa)  
 s-w opt: 178 Z-score: 188.8 bits: 42.4 E(): 4.4e-05  
 Smith-Waterman score: 178; 23.3% identity (58.5% similar) in 236 aa overlap (8-222:45-264)

Database	Entries	Length	E()	hits	time (s)
E. coli	4,237	1.3 E 06	1.5 E-06*	1	< 0.5
S. cerevisiae	5,866	2.9 E 06	2.1 E-06	1	< 0.5
Human	38,114	18.4 E 06	1.2 E-05	1	1.1
Swiss Prot	4.3 E 05	1.5 E 08	2.4 E-05*	393	7.1
Refseq NP only	7.1 E 05	2.6 E 08	0.00017*	504	10.8
Refseq	7.3 E 06	2.5 E 09	0.0017*	2767	124
NR	9.9 E 06	3.4 E 09	0.0032*	7773	151

40

## NCBI – selecting sequences with Entrez

NCBI/ BLAST/ blastp suite

blastn blastp blastx tblastn tblastx

BLASTP programs search protein databases using a protein query. [more...](#)

**Enter Query Sequence**

Enter accession number, gi, or FASTA sequence [Clear](#) **Query subrange**

From

To

Or, upload file  no file selected

**Job Title**

Enter a descriptive title for your BLAST search

**Align two or more sequences**

**Choose Search Set**

**Database** Reference proteins (refseq\_protein)

**Organism** Optional human (taxid:9606)  Exclude

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

**Entrez Query** Optional

Enter an Entrez query to limit search

41

## Similarity Searching II

1. What question to ask?
2. What program to use?
3. What database to search?
4. How to avoid mistakes (what to look out for)
5. When to do something different

42

## Inferring Homology from Statistical Significance

- Real **UNRELATED** sequences have similarity scores that are indistinguishable from **RANDOM** sequences
- If a similarity is NOT **RANDOM**, then it must be NOT **UNRELATED**
- Therefore, NOT **RANDOM** (statistically significant) similarity must reflect **RELATED** sequences

43

## Smith-Waterman (ssearch)

The best scores are:

			s-w	bits	E(115640)	%_id	alen
GTM1_MOUSE	Glutathione S-trans	( 218)	1497	363.5	2e-100	1.000	218
GTM2_CHICK	Glutathione S-trans	( 220)	958	234.9	1.1e-61	0.619	218
GTP_HUMAN	Glutathione S-trans	( 210)	356	91.2	1.8e-18	0.308	211
PGD2_MOUSE	Glutathione-req.	( 199)	262	68.8	9.7e-12	0.319	204
GTA1_MOUSE	Glutathione S-trans	( 223)	229	60.9	2.6e-09	0.284	225
SC1_OCTDO	S-crystallin 1 OL1	( 215)	228	60.7	3.0e-09	0.269	219
GTS_MUSDO	Glutathione S-trans	( 241)	228	60.6	3.4e-09	0.264	201
GTS1_CAEEL	Prob. Glut. S-trans	( 210)	220	58.8	1.1e-08	0.284	225
GTS_OMMSL	Glutathione S-trans	( 203)	196	53.0	5.5e-07	0.258	209
GTH3_ARATH	Glutathione S-trans	( 215)	142	40.1	0.0045	0.310	126
GTT2_HUMAN	Glutathione S-trans	( 244)	132	37.7	0.027	0.257	167
GT24_DROME	Glutathione S-trans	( 216)	131	37.5	0.028	0.255	153
YFCG_ECOLI	Hypothetical GST	( 215)	112	33.0	0.64	0.235	187
YJY1_YEAST	hypothetical 30.5	( 261)	110	32.4	*1.1*	0.248	149
DCMA_METS1	dichloromethane DM	( 267)	103	30.8	3.7	0.214	210
YA42_HAEIN	Hypothetical prot.	( 617)	108	31.7	*4.6*	0.283	120
GTO1_RAT	Glutathione trans	( 241)	100	30.1	5.4	0.234	158
DP41_BACHD	DNA polymerase I	( 413)	104	30.8	*5.4*	0.234	184
GTH1_WHEAT	Glutathione S-trans	( 229)	98	29.6	7.0	0.246	171
LGUL_SOYBN	Lactoylglutathione	( 219)	97	29.4	7.8	0.200	190

Highest scoring unrelated sequence E() ~ 1.0

44

## Looking for mistakes; BLASTP at NCBI

NCBI/BLAST/blastp suite **BLAST Human Sequences**

blastn **blastp** blastx tblastn tblastx

BLASTP programs search protein databases using a protein query. [more...](#)

**Enter Query Sequence**

Enter accession number(s), gi(s), or FASTA sequence(s)  [Clear](#) [Query subrange](#)

From   
To

Or, upload file  no file selected

Job Title   
Enter a descriptive title for your BLAST search

**Choose Search Set**

Database  29315 sequences

Exclude  Models (XM/XP)  Uncultured/environmental sample sequences

Entrez Query   
Enter an Entrez query to limit search

**Program Selection**

Algorithm  blastp (protein-protein BLAST)  
 PSI-BLAST (Position-Specific Iterated BLAST)  
 PHI-BLAST (Pattern Hit Initiated BLAST)  
[Choose a BLAST algorithm](#)

**BLAST** Search database RefSeq protein using Blastp (protein-protein BLAST)  
 Show results in a new window

45

## BLAST results 1 - gstt1\_drome

BLAST Results - 7PAPKRC016

1 record to records matching [entrez query: human\[orgn\]](#).

[Search Strategies](#) [Formatting options](#) [Download](#)

1) [\[594\] \[sp\] P20432.1 \[GSTT1\\_DROME\]](#) Database Name [GP/9606.9558/RefSeq\\_protein](#)  
Accession: P20432.1; Description: Full=Glutathione S-transferase 1-1; AltName: Full=DDT-rochlorinase; AltName: Full=GST class-theta acid; Description: Homo sapiens RefSeq protein  
Program: BLASTP 2.2.25+ [Citation](#)

[Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Multiple alignment](#)

Putative conserved domains have been detected, click on the image below for detailed results.

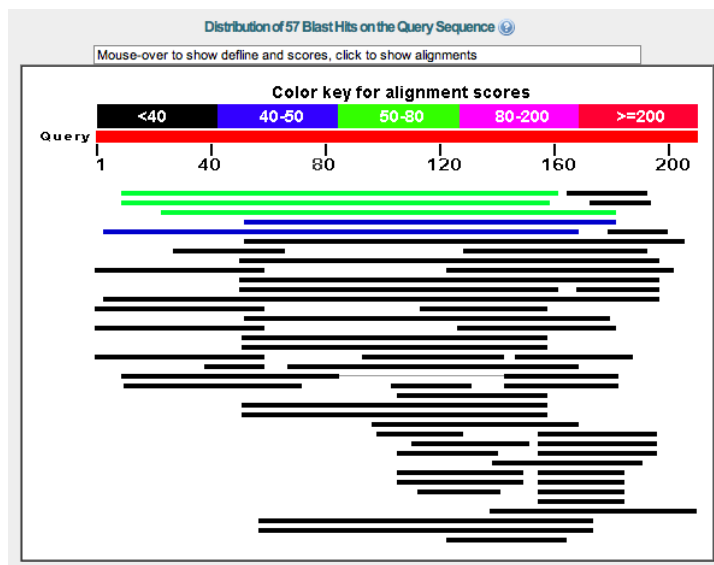
seq.

binding site (G-site) dimer interface substrate binding pocket (theta) N-terminal domain interface

Thioredoxin-like superfamily **GST\_N\_Delta\_Epsilon** **GST\_C\_Delta\_Epsilon** GST\_C\_family superfamily

46

### BLAST results 2 - gstt1\_drome



47

### BLAST results 3 - gstt1\_drome

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Links
<a href="#">NP_000844.2</a>	glutathione S-transferase theta-1 [Homo sapiens]	<a href="#">76.3</a>	76.3	72%	2e-19	<a href="#">UGM</a>
<a href="#">NP_000845.1</a>	glutathione S-transferase theta-2 [Homo sapiens] >ref[NP_665877.1]	<a href="#">71.6</a>	71.6	71%	1e-17	<a href="#">GM</a>
<a href="#">NP_665877.1</a>	maleylacetoacetate isomerase isoform 1 [Homo sapiens]	<a href="#">56.6</a>	56.6	75%	2e-12	<a href="#">UGM</a>
<a href="#">NP_001504.2</a>	maleylacetoacetate isomerase isoform 3 [Homo sapiens]	<a href="#">47.4</a>	47.4	61%	2e-09	<a href="#">UGM</a>
<a href="#">NP_001503.1</a>	glutathione S-transferase A4 [Homo sapiens]	<a href="#">40.0</a>	40.0	78%	1e-06	<a href="#">UGM</a>
<a href="#">NP_004271.1</a>	eukaryotic translation elongation factor 1 epsilon-1 isoform	<a href="#">38.5</a>	38.5	73%	3e-06	<a href="#">UGM</a>
<a href="#">NP_006294.2</a>	aminoacyl tRNA synthase complex-interacting multifunctio	<a href="#">37.7</a>	37.7	30%	1e-05	<a href="#">UGM</a>
<a href="#">NP_665683.1</a>	glutathione S-transferase A1 [Homo sapiens]	<a href="#">37.4</a>	37.4	69%	1e-05	<a href="#">UGM</a>
<a href="#">NP_036270.1</a>	protein AATF [Homo sapiens]	<a href="#">35.0</a>	35.0	37%	8e-05	<a href="#">UGM</a>
<a href="#">NP_000837.3</a>	glutathione S-transferase A2 [Homo sapiens]	<a href="#">34.7</a>	34.7	69%	8e-05	<a href="#">UGM</a>
<a href="#">NP_000838.3</a>	glutathione S-transferase A3 [Homo sapiens]	<a href="#">33.9</a>	33.9	53%	2e-04	<a href="#">UGM</a>
<a href="#">NP_714543.1</a>	glutathione S-transferase A5 [Homo sapiens]	<a href="#">32.7</a>	32.7	92%	3e-04	<a href="#">UGM</a>
<a href="#">NP_671488.1</a>	blood vessel epicardial substance [Homo sapiens] >ref[NP_001129122.1]	<a href="#">33.1</a>	33.1	21%	4e-04	<a href="#">UGM</a>
<a href="#">NP_001129122.1</a>	eukaryotic translation elongation factor 1 epsilon-1 isoform	<a href="#">31.6</a>	31.6	60%	6e-04	<a href="#">GM</a>
<a href="#">NP_665878.2</a>	maleylacetoacetate isomerase isoform 2 [Homo sapiens]	<a href="#">30.4</a>	30.4	26%	0.002	<a href="#">UGM</a>
<a href="#">NP_671489.1</a>	glutathione S-transferase Mu 4 isoform 2 [Homo sapiens]	<a href="#">28.1</a>	28.1	50%	0.012	<a href="#">UGM</a>
<a href="#">NP_000841.1</a>	glutathione S-transferase Mu 4 isoform 1 [Homo sapiens]	<a href="#">28.1</a>	28.1	50%	0.014	<a href="#">UGM</a>
<a href="#">NP_000245.2</a>	methionine synthase [Homo sapiens]	<a href="#">27.3</a>	27.3	23%	0.031	<a href="#">UEGM</a>
<a href="#">NP_001182566.1</a>	hypothetical protein LOC100500938 [Homo sapiens]	<a href="#">25.4</a>	25.4	18%	0.036	<a href="#">UG</a>
<a href="#">XP_003403539.1</a>	PREDICTED: glutathione S-transferase theta-4-like [Homo	<a href="#">26.6</a>	26.6	48%	0.038	<a href="#">G</a>
<a href="#">NP_061845.2</a>	ganglioside-induced differentiation-associated protein 1 iso	<a href="#">26.6</a>	52.7	55%	0.053	<a href="#">UGM</a>
<a href="#">NP_057460.3</a>	ankyrin repeat and FYVE domain-containing protein 1 isofo	<a href="#">26.6</a>	26.6	13%	0.055	<a href="#">UGM</a>
<a href="#">NP_001035808.1</a>	ganglioside-induced differentiation-associated protein 1 iso	<a href="#">26.2</a>	26.2	19%	0.067	<a href="#">UGM</a>
<a href="#">NP_150648.2</a>	vacuolar protein sorting-associated protein 13A isoform A	<a href="#">26.2</a>	26.2	28%	0.068	<a href="#">UGM</a>
<a href="#">NP_001018047.1</a>	vacuolar protein sorting-associated protein 13A isoform C	<a href="#">26.2</a>	26.2	28%	0.071	<a href="#">UGM</a>
<a href="#">NP_001280.3</a>	chloride intracellular channel protein 2 [Homo sapiens]	<a href="#">25.8</a>	25.8	19%	0.071	<a href="#">UGM</a>

48



## BLAST results 4 - gstt1\_drome

Homolog? [>ref|NP\\_036270.1| UGM](#) protein AATF [Homo sapiens]  
 Length=560

[GENE ID: 26574 AATF](#) | apoptosis antagonizing transcription factor  
 [Homo sapiens] (Over 10 PubMed links)

Score = 35.0 bits (79), Expect = 8e-05, Method: Compositional matrix adjust.  
 Identities = 24/79 (30%), Positives = 34/79 (43%), Gaps = 7/79 (9%)

```

Query 123 ADPEAFKKIEAAFEFLNTFLEGGDYAAGDSLTVADIALVATVSTFEVAKFEISKYANVNR 182
          + A ++ F EG+D GD L V I +A+ S + K K +
Sbjct 22 ADPEADPEEATAARVIDRFDEGED-GEGLVVGSIKRLASASLLDTRKRYCGKTSRKA 80

Query 183 WYENAKKVTGWENWAGC 201
          W E+ WE+ G
Sbjct 81 WNEDE-----WEQTLPGS 93
    
```

Homolog? [>ref|NP\\_000837.3| UGM](#) glutathione S-transferase A2 [Homo sapiens]  
 Length=222

[GENE ID: 2939 GSTA2](#) | glutathione S-transferase alpha 2 [Homo sapiens]  
 (Over 10 PubMed links)

Score = 34.7 bits (78), Expect = 8e-05, Method: Compositional matrix adjust.  
 Identities = 42/181 (23%), Positives = 71/181 (39%), Gaps = 58/181 (32%)

```

Query 51 HTIPTLVNDGFALWESRAIQVYLVEKYKTDLSLYPKCPKRAVINQRLL--YFDMGTLY-- 106
          +P + +G L ++RAI Y+ KY +LY K K++A+I+ + D+G +
Sbjct 53 QQVPMVEIDGMKLVQTRAILNYIASKY----NLYGKDIKEKALIDMYIEGIADLGEMILL 108

Query 107 -----QSFANYYPQVFAKAPADPEAFKKIEAAFEFLNTFLEGGDYAA 149
          + N Y+P AF+K+ + GQDY
Sbjct 109 LPFSQPEEQDAKLALIQEKTKNRYFP-----AFEKVLKS-----HGQDYLV 149

Query 150 GDSLTVADIALV-----ATVSTF---EVAKFEISKYANVNRWYENAKKVTGPWE 195
          G+ L+ ADI LV + +S+F + K IS V ++ + P +
Sbjct 150 GNKLSRADIHVVELLYVEELDSSLISFLLKALKTRISNLPVTKKFLQGPSRKPMPD 209

Query 196 E 196
          E
Sbjct 210 E 210
    
```

49

## BLAST results – validating statistics

Re-search vs SwissProt at NCBI site  
 Query: AATF – initial  $E() < 10^{-4}$

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Links
<a href="#">Q9N61.1</a>	RecName: Full=Protein AATF; AltName: Full=Apoptosis-antagonizing transcription factor	1137	1137	100%	0.0	<a href="#">G</a>
<a href="#">Q9JXC4.1</a>	RecName: Full=Protein AATF; AltName: Full=Apoptosis-antagonizing transcription factor	794	794	99%	0.0	<a href="#">G</a>
<a href="#">Q9QYW0.1</a>	RecName: Full=Protein AATF; AltName: Full=Apoptosis-antagonizing transcription factor	752	752	99%	0.0	<a href="#">G</a>
<a href="#">Q5ZIM6.1</a>	RecName: Full=Protein AATF; AltName: Full=Apoptosis-antagonizing transcription factor	643	643	97%	0.0	<a href="#">G</a>
<a href="#">Q9VM95.1</a>	RecName: Full=Protein AATF-like	140	140	58%	1e-34	<a href="#">M</a>
<a href="#">Q55E65.2</a>	RecName: Full=Putative uncharacterized protein DDB_G02	129	129	59%	5e-31	
<a href="#">Q6BXX1.2</a>	RecName: Full=Protein BFR2	111	111	60%	6e-25	
<a href="#">Q7S6P8.1</a>	RecName: Full=Protein bfr-2	112	112	61%	9e-25	
<a href="#">Q6C9G2.1</a>	RecName: Full=Protein BFR2	109	109	59%	2e-24	
<a href="#">Q9US05.1</a>	RecName: Full=Protein bfr2	100	100	59%	1e-21	<a href="#">G</a>
<a href="#">Q6CTS8.1</a>	RecName: Full=Protein BFR2	94.0	94.0	67%	4e-19	
<a href="#">Q5ACL9.2</a>	RecName: Full=Protein BFR2	92.8	92.8	58%	1e-18	
<a href="#">Q6FSD4.1</a>	RecName: Full=Protein BFR2	91.7	91.7	58%	2e-18	
<a href="#">Q4WMI1.1</a>	RecName: Full=Protein bfr2	88.6	88.6	70%	2e-17	<a href="#">G</a>
<a href="#">Q06631.1</a>	RecName: Full=Protein BFR2; AltName: Full=Brefeldin A receptor	87.4	87.4	60%	5e-17	<a href="#">G</a>
<a href="#">Q3AW04.1</a>	RecName: Full=Protein bfr2	86.7	86.7	64%	1e-16	
<a href="#">Q4P9V5.1</a>	RecName: Full=Protein BFR2	85.5	85.5	61%	3e-16	
<a href="#">Q4I3Z7.1</a>	RecName: Full=Protein BFR2	82.8	82.8	60%	2e-15	
<a href="#">Q7SEZ2.1</a>	RecName: Full=Protein BFR2	81.6	81.6	60%	3e-15	<a href="#">G</a>
<a href="#">POCL91.1</a>	RecName: Full=Protein BFR2 >sp POCL90.1 BFR2_CRYNJ R	43.9	82.4	21%	0.003	

50

## BLAST results – validating statistics

Re-search vs SwissProt at NCBI site

Query: GSTA1 – initial  $E() < 10^{-4}$

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Links
P09210.4	RecName: Full=Glutathione S-transferase A2; AltName: Fu	459	449	100%	4e-162	<a href="#">G</a> <a href="#">M</a>
P08263.3	RecName: Full=Glutathione S-transferase A1; AltName: Fu	427	427	100%	2e-153	<a href="#">G</a> <a href="#">M</a>
Q78TV2.1	RecName: Full=Glutathione S-transferase A5; AltName: Fu	399	399	100%	1e-142	<a href="#">G</a>
Q16772.3	RecName: Full=Glutathione S-transferase A3; AltName: Fu	399	399	100%	2e-142	<a href="#">G</a>
Q28035.3	RecName: Full=Glutathione S-transferase A1; AltName: Fu	374	374	100%	1e-132	<a href="#">G</a>
Q18929.4	RecName: Full=Glutathione S-transferase A2; AltName: Fu	373	373	100%	4e-132	<a href="#">G</a> <a href="#">M</a>
P51781.2	RecName: Full=Glutathione S-transferase alpha M14; AltN	367	367	99%	7e-130	<a href="#">G</a> <a href="#">M</a>
Q08892.2	RecName: Full=Glutathione S-transferase Yc; AltName: Fu	350	350	99%	4e-123	<a href="#">G</a> <a href="#">M</a>
Q64200.1	RecName: Full=Glutathione S-transferase A6; AltName: Fu	349	349	100%	8e-123	<a href="#">G</a> <a href="#">M</a>
Q08863.1	RecName: Full=Glutathione S-transferase alpha I; AltNam	340	340	99%	3e-119	<a href="#">G</a> <a href="#">M</a>
P30115.2	RecName: Full=Glutathione S-transferase A3; AltName: Fu	339	339	99%	1e-118	<a href="#">G</a>
P04904.3	RecName: Full=Glutathione S-transferase alpha-3; AltNam	338	338	99%	3e-118	<a href="#">G</a> <a href="#">M</a> <a href="#">S</a>
P13745.2	RecName: Full=Glutathione S-transferase A1; AltName: Fu	338	338	99%	4e-118	<a href="#">G</a>
P46418.2	RecName: Full=Glutathione S-transferase alpha-5; AltNam	331	331	99%	1e-115	<a href="#">G</a> <a href="#">M</a>
P10648.2	RecName: Full=Glutathione S-transferase A2; AltName: Fu	326	326	100%	2e-113	<a href="#">G</a>
P04903.2	RecName: Full=Glutathione S-transferase alpha-2; AltNam	323	323	100%	2e-112	<a href="#">G</a> <a href="#">M</a>
P09592.3	RecName: Full=Glutathione S-transferase alpha-1; AltNam	323	323	100%	2e-112	<a href="#">G</a> <a href="#">M</a>
P81706.1	RecName: Full=Glutathione S-transferase A; Short=GST A	313	313	97%	1e-108	<a href="#">G</a>
Q08393.2	RecName: Full=Glutathione S-transferase; AltName: Full=	305	305	100%	4e-105	<a href="#">G</a>
Q08392.1	RecName: Full=Glutathione S-transferase; AltName: Full=	295	295	99%	4e-101	<a href="#">G</a>
P28697.2	RecName: Full=Glutathione S-transferase 3; AltName: Full	287	287	99%	4e-98	<a href="#">G</a> <a href="#">M</a>
P24472.3	RecName: Full=Glutathione S-transferase A4; AltName: Fu	283	283	100%	1e-96	<a href="#">G</a>
P14942.2	RecName: Full=Glutathione S-transferase alpha-4; AltNam	278	278	100%	1e-94	<a href="#">G</a> <a href="#">M</a>
P80894.1	RecName: Full=Glutathione S-transferase; AltName: Full=	267	267	100%	3e-90	<a href="#">G</a>
Q5E900.1	RecName: Full=Glutathione S-transferase A4; AltName: Fu	264	264	99%	7e-89	<a href="#">G</a> <a href="#">M</a>
Q15217.1	RecName: Full=Glutathione S-transferase A4; AltName: Fu	261	261	99%	4e-88	<a href="#">G</a>
P80931.2	RecName: Full=Glutathione S-transferase P; AltName: Full	96.3	96.3	86%	2e-23	<a href="#">G</a>
Q29408.3	RecName: Full=Glutathione S-transferase P 10; AltName:	95.9	95.9	87%	3e-23	<a href="#">G</a>
P12354.2	RecName: Full=Glutathione S-transferase P; AltName: Full	95.1	95.1	85%	4e-23	<a href="#">G</a>
P81842.1	RecName: Full=Glutathione S-transferase P 1; AltName: F	94.0	94.0	90%	1e-22	<a href="#">G</a>
P46426.1	RecName: Full=Glutathione S-transferase; AltName: Full=	93.2	93.2	91%	3e-22	<a href="#">G</a>
Q60550.3	RecName: Full=Glutathione S-transferase P; AltName: Full	92.8	92.8	85%	4e-22	<a href="#">G</a>
P46424.2	RecName: Full=Glutathione S-transferase P; AltName: Full	92.4	92.4	85%	6e-22	<a href="#">G</a>
P46427.1	RecName: Full=Glutathione S-transferase 2; AltName: Full	89.7	89.7	91%	6e-21	<a href="#">G</a>
P28801.2	RecName: Full=Glutathione S-transferase P; AltName: Full	89.7	89.7	85%	6e-21	<a href="#">G</a> <a href="#">M</a>
Q18598.3	RecName: Full=Glutathione S-transferase; AltName: Full=	89.0	89.0	94%	1e-20	<a href="#">G</a>

51

## Unrelated ≠ Random (low complexity)

Search with complete grou\_drome:

The best scores are:

				opt	bits	E(14548)
RGHUB1	GTP-binding regulatory protein beta-1	chai	( 341)	237	46.6	3.5e-05
RGBOB1	GTP-binding regulatory protein beta-1	chai	( 341)	237	46.6	3.5e-05
RGHUB3	GTP-binding regulatory protein beta-3	chai	( 341)	233	46.0	5.2e-05
RGMSB4	GTP-binding regulatory protein beta-4	chai	( 341)	232	45.8	5.7e-05
PIHUPF	salivary proline-rich glycoprotein precurs		( 252)	224	44.5	*0.00010*
RGFFB	GTP-binding regulatory protein beta chain		( 347)	223	44.5	0.00014
PIRT3	acidic proline-rich protein precursor - rat		( 207)	199	40.8	*0.0011*
PIHUB6	salivary proline-rich protein precursor PR		( 393)	203	41.6	*0.0012*
CGBO2S	collagen alpha 2(I) chain - bovine (fragme		( 403)	195	40.5	*0.0027*
WMBEW6	capsid protein - human herpesvirus 1 (stra		( 636)	192	40.2	*0.0051*

Search with seg-ed grou\_drome: (low complexity regions removed)

The best scores are:

				opt	bits	E(14548)	
RGHUB3	GTP-binding regulatory protein beta-3	chai	( 341)	233	56.5	3.6e-08	
RGMSB4	GTP-binding regulatory protein beta-4	chai	( 341)	232	56.3	4.1e-08	
RGHUB2	GTP-binding regulatory protein beta-2	chai	( 341)	228	55.5	7.2e-08	
RGBOB1	GTP-binding regulatory protein beta-1	chai	( 341)	225	54.9	1.1e-07	
RGFFB	GTP-binding regulatory protein beta chain		( 347)	223	54.5	1.5e-07	
BVBVMS	MSI1 protein - yeast (Saccharomyces cerevi		( 423)	135	37.0	*0.033*	
ERHUAH	coatomer complex alpha chain homolog - hum		(1225)	134	37.1	*0.088*	
A28468	chromogranin A precursor - human		( 458)	122	34.4	*0.21*	
RGOOBE	GTP-binding regulatory protein beta chain		( 342)	120	33.9	0.22	52

### pseg removes low-complexity regions

>gi|17380405|sp|P16371|GROU\_DROME Groucho protein (Enhancer of split M9/10)

```

1-8      MYSPVVRH
paagppppgpp 9-19
20-131   IKFTIADTLERIKEEFNFLQAQYHSIKLEC
          EKLSNEKTEMQRHYVVMYEMSYGLNVMHK
          QTEIAKRLNTLINQLLPFLQADHQQVQLQA
          VERAKQVTMQEELNLIIGQQIHA
qqvppgppppmg 132-143
144-281  ALNPPGALGATMGLPHGPQGLLNKPPPEHHR
          PDIKPTGLEGPAAAEERLRNSVSPADREKY
          RTRSPLDIENDSKRRKDEKLEDEGEKSDQ
          DLVVDVANEMESHSPRPNGEHVSMVDRRE
          SLNGERLEKPPSSSGIKQE
rppsrsgsssrstps 282-297
298-310  LKTKDMEKPGTPG
akartptpnaaapagvnpk 311-330
qmmpqgppppagypgapyqrpa 331-351
352-719  DPYQRPPSDPAYGRPPMPYDPHAHVRTNG
          IPHPSALTGGKPAYSFHMNGESLQPVPPFP
          PDALVGVGIPRHARQINTLSHGEVCAVTI
          SNPTKYVYTGKGCVKVWDISQPGNKNPVS
          QLDCLQRDNYIRSVKLLPDGRTLIVGGEAS
          NLSIWDLASPTPRIKAELTSAAPACYALAI
          SPDSKVCFCSCSDGNI AVWDLHNEILVRQF
          QGHTDGASCIDISPDGSRLLWTGGLDNTVRS
          WDLREGRLQQHDFSSQIFSLGYCPTGDWL
          AVGMENSHVEVLHASKPKDYQLHLHESCVL
          SLRFAACGKWFVSTGKDNLLNARWTPYGAS
          IFQSKETS SVLSCDITDDKYIVTGS GDKK
          ATVYEVII
    
```

53

### BLAST remove low complexity - gstt1\_drome

**Scoring Parameters**

Matrix: BLOSUM62

Gap Costs: Existence: 11 Extension: 1

Compositional adjustments: Conditional compositional score matrix adjustment

---

**Filters and Masking**

Filter:  Low complexity regions

Mask:  Mask for lookup table only  Mask lower case letters

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Links
NP_000844.2	glutathione S-transferase theta-1 [Homo sapiens]	76.3	76.3	72%	8e-17	<a href="#">UGM</a>
NP_000845.1	glutathione S-transferase theta-2 [Homo sapiens] >ref[NP_000845.1]	71.6	71.6	71%	5e-15	<a href="#">UGM</a>
NP_665877.1	maleylacetoacetate isomerase isoform 1 [Homo sapiens]	56.6	56.6	75%	7e-10	<a href="#">UGM</a>
NP_001504.2	maleylacetoacetate isomerase isoform 3 [Homo sapiens]	47.4	47.4	61%	8e-07	<a href="#">UGM</a>
NP_001903.1	glutathione S-transferase A4 [Homo sapiens]	49.0	40.0	78%	5e-04	<a href="#">UGM</a>
NP_004271.1	eukaryotic translation elongation factor 1 epsilon-1 isoform	38.5	38.5	73%	0.001	<a href="#">UGM</a>
NP_006294.2	aminoacyl tRNA synthase complex-interacting multifunctional	37.7	37.7	30%	0.004	<a href="#">UGM</a>
NP_665683.1	glutathione S-transferase A1 [Homo sapiens]	37.4	37.4	69%	0.004	<a href="#">UGM</a>
NP_036270.1	protein AATF [Homo sapiens]	35.0	35.0	37%	0.030	<a href="#">UGM</a>
NP_000837.3	glutathione S-transferase A2 [Homo sapiens]	34.7	34.7	69%	0.031	<a href="#">UGM</a>
NP_000838.3	glutathione S-transferase A3 [Homo sapiens]	33.9	33.9	53%	0.056	<a href="#">UGM</a>
NP_714543.1	glutathione S-transferase A5 [Homo sapiens]	32.7	32.7	92%	0.13	<a href="#">UGM</a>
NP_671488.1	blood vessel epicardial substance [Homo sapiens] >ref[NP_000845.1]	33.1	33.1	21%	0.13	<a href="#">UGM</a>
NP_001129122.1	eukaryotic translation elongation factor 1 epsilon-1 isoform	31.6	31.6	60%	0.24	<a href="#">UGM</a>
NP_665878.2	maleylacetoacetate isomerase isoform 2 [Homo sapiens]	30.4	30.4	26%	0.73	<a href="#">UGM</a>
NP_671489.1	glutathione S-transferase Mu 4 isoform 2 [Homo sapiens]	28.1	28.1	50%	4.6	<a href="#">UGM</a>
NP_000841.1	glutathione S-transferase Mu 4 isoform 1 [Homo sapiens]	28.1	28.1	50%	5.2	<a href="#">UGM</a>

## Validating homologs/statistics

- In general, BLASTP statistical estimates are accurate
- The most common errors occur because of low-complexity regions, or biased amino-acid composition
- To confirm statistical accuracy, find the highest scoring non homolog
  - No need to test every hit, test hits that are surprising
  - Confirm homology/non-homology by searching against a different comprehensive database, e.g. SwissProt, or refseq.
  - Non-homologs will find many significant members of other families, but not the family you are testing for
- Statistical estimates can be confirmed with shuffles (see ISMB2000 tutorial, [fasta.bioch.virginia.edu/fasta\\_www2](http://fasta.bioch.virginia.edu/fasta_www2) shuffle link)

55

## Scoring matrices

- Scoring matrices can set the evolutionary look-back time for a search
  - Lower PAM (PAM10/MDM10 ... PAM60) for closer (10% ... 50% identity)
  - Higher BLOSUM for higher conservation (BLOSUM50 distant, BLOSUM80 conserved)
- Shallow scoring matrices for short domains/short queries (metagenomics)
  - Matrices have “bits/position” (score/position), 40 aa at 0.7 bits/position (BLOSUM62) means 28 bit max score (50 bits significant)
- Deep scoring matrices allow alignments to continue, possibly outside the homologous region

56

### Finding Domains – Local alignments: calmodulin

```

46.1% identity in 76 aa overlap (1-76:77-149); score: 222 E(10000): 2.7e-10
      10      20      30      40      50      60
mchu MADQLTEEQIAEFKEAFSLFDKDGDTITTKELGTVMRSLGQNPTEAELQDMINEVDADG
      : : .::: .::: .::: .::: .::: .::: .::: .::: .::: .::: .::: .::: .:::
mchu MKDTSSEEEI---REAFRVFDKDGNGYISAAELRHVMTNLGKLTDEEVDEMIREADIDG
      80      90      100     110     120     130

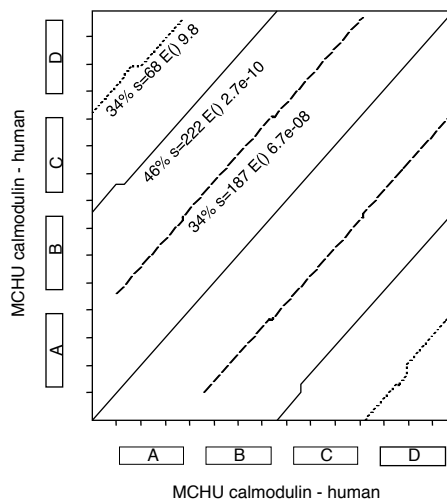
      70
mchu NGTIDFPEFLTMMARK
      .: .: .: .: .: .:
mchu DGQVNYEEFVQMMTAK
      140

34.3% identity in 105 aa overlap (11-111:47-147); score: 187 E(10000): 6.7e-08
      20      30      40      50      60
mchu AEFKEAFSLFDKDGDTITTKELGTVM-RSLGQNPTEAELQDMINEVDADGNGTIDFPEF
      : : .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .:
mchu AELQDMINEVDADGNGTIDFPEFLTMMARKMKDTSSEEEIREAFRVFDKDGNGYISAAEL
      50      60      70      80      90      100
      70      80      90      100     110
mchu ---LTMMARKMKDTSSEEEIREAFRVFDKDGNGYISAAELRHVMT
      .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .:
mchu RHVMTNLGKLTDEEVDEMIREA---DIDGGQVNYEEFVQMMT
      110     120     130     140

34.2% identity in 38 aa overlap (1-37:113-146); score: 68 E(10000): 9.8
      10      20      30
mchu MADQLTEEQIAEF-KEAFSLFDKDGDTITTKELGTVM
      .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .: .:
mchu LGEKLTDEEVDEMIREA---DIDGGQVNYEEFVQMM
      120     130     140
    
```

57

### Repeated domains with local alignments



58

## More about scoring matrices ...

### PAM series:

- Evolutionary model - extrapolated from PAM1
- PAM20: 20% change (mammals)
- PAM250: 250% change (<20% identity)
- Gap penalties should vary
- shallow matrices (PAM10-40) for short sequences and short distances

### BLOSUM series

- Empirically determined, no extrapolation (no model)
- BLOSUM45-50 - distant (1/3 bits)
- BLOSUM80 -very highly conserved (not small change), high info/position
- BLOSUM62 - 1/2 bits

59

## Where do scoring matrices come from?

### Pam40

```

A R N D E I L
A 8
R -9 12
N -4 -7 11
D -4 -13 3 11
E -3 -11 -2 4 11
I -6 -7 -7 -10 -7 12
L -8 -11 -9 -16 -12 -1 10

```

### Pam250

```

A R N D E I L
A 2
R -2 6
N 0 0 2
D 0 -1 2 4
E 0 -1 1 3 4
I -1 -2 -2 -2 -2 5
L -2 -3 -3 -4 -3 2 6

```

$q_{ij}$  : replacement frequency at PAM40, 250

$$q_{R:N(40)} = 0.000435$$

$$p_R = 0.051$$

$$q_{R:N(250)} = 0.002193$$

$$p_N = 0.043$$

$$I_2 S_{ij} = \lg_2 (q_{ij}/p_i p_j) \quad I_e S_{ij} = \ln(q_{ij}/p_i p_j) \quad p_R p_N = 0.002193$$

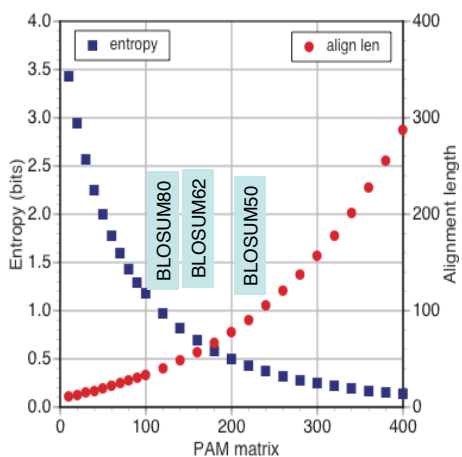
$$I_2 S_{R:N(40)} = \lg_2 (0.000435/0.00219) = -2.333$$

$$I_2 = 1/3; S_{R:N(40)} = -2.333/I_2 = -7$$

$$I S_{R:N(250)} = \lg_2 (0.002193/0.002193) = 0$$

60

## PAM matrices and alignment length



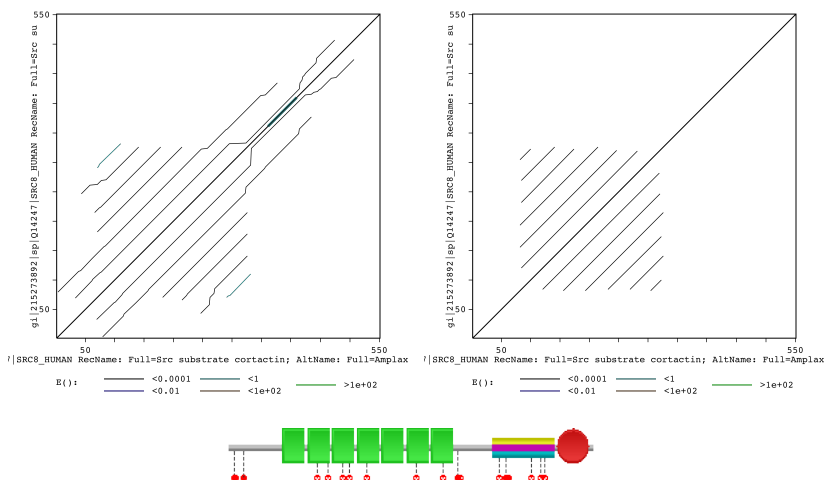
Short domains require “shallow” scoring matrices

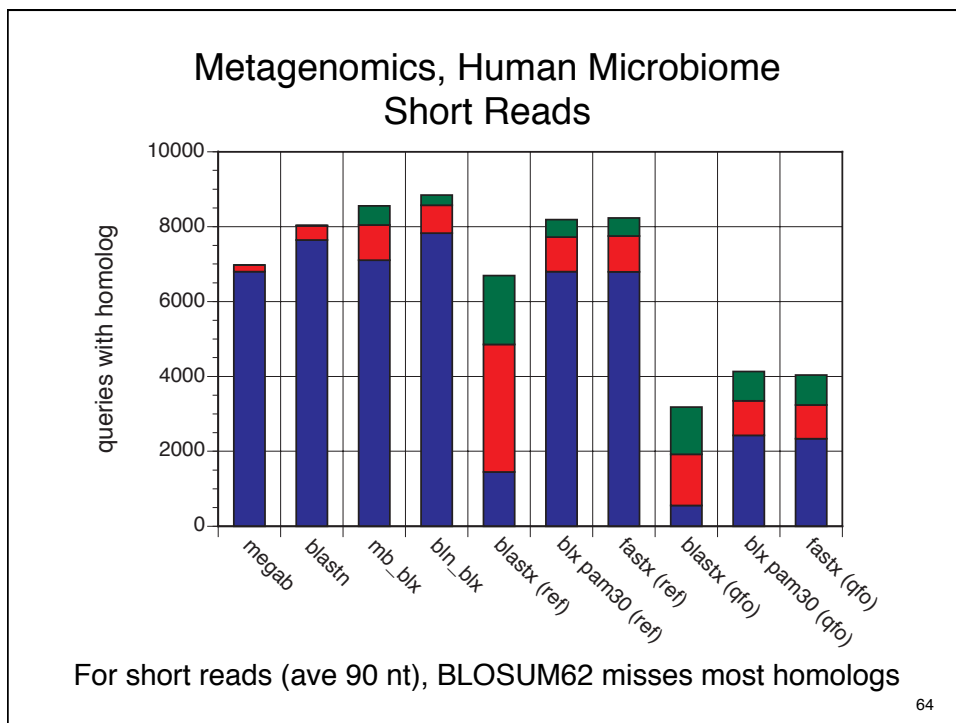
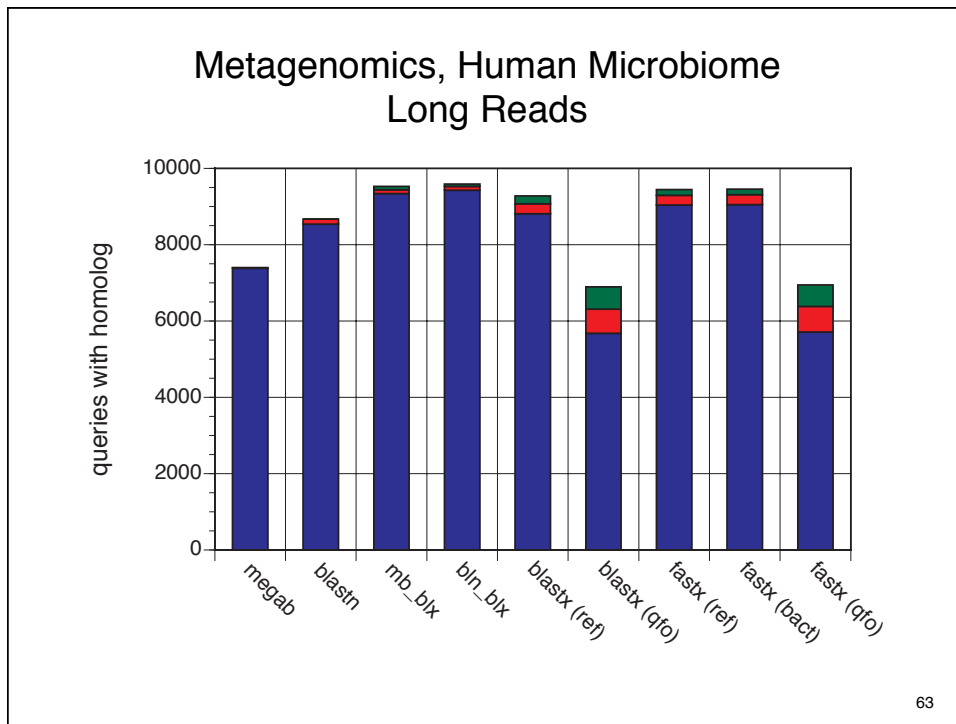
61

## Scoring matrices affect alignment lengths

BLOSUM62 -11/-1

MD20 -26/-4







## *Scoring Matrices - Summary*

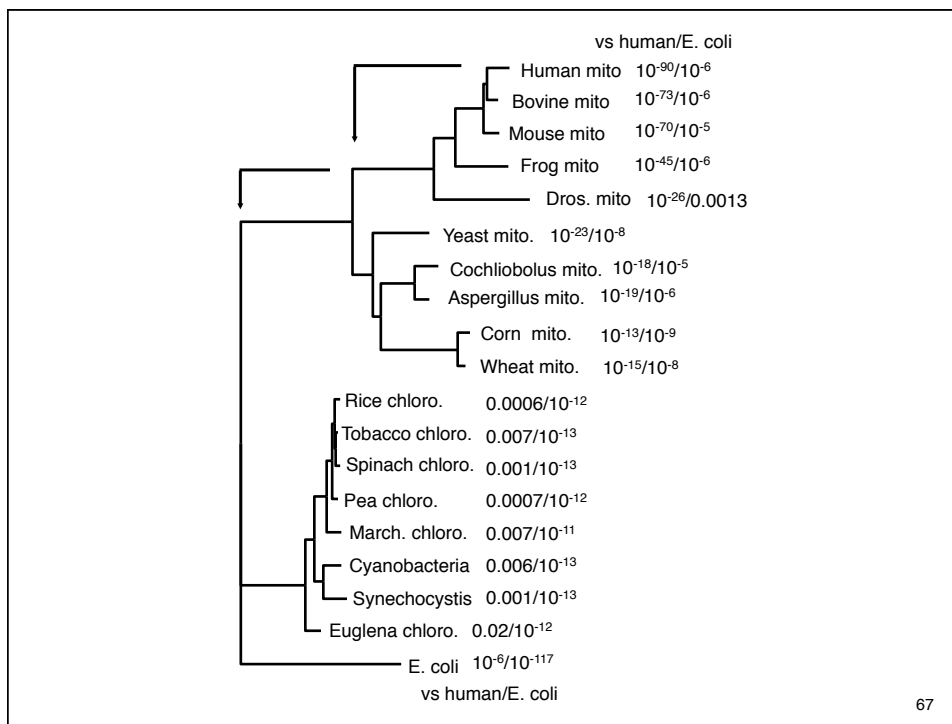
- PAM and BLOSUM matrices greatly improve the sensitivity of protein sequence comparison – low identity with significant similarity
- PAM matrices have an evolutionary model - lower number, less divergence – lower=closer; higher=more distant
- BLOSUM matrices are sampled from conserved regions at different average identity – higher=more conservation
- Short alignments require shallow matrices
- Shallow matrices set maximum look-back time

65

## Similarity Searching II

1. What question to ask?
2. What program to use?
3. What database to search?
4. How to avoid mistakes (what to look out for)
5. When to do something different
6. PSI-BLAST – the most sensitive method

66



### ATP synthase - matrices, gaps, algorithms

```

Matrix:
Gap open/extend      BLOSUM50      BLOSUM62      BLASTP
                    -10/-2       -11/-1
The best scores are:  bits E(13351) bits E(13351) bits E(
ATP6_HUMAN ATP synthase a chai 297.7 1.7e-81 373.6 2.4e-104 296 3e-81
ATP6_BOVIN ATP synthase a chai 252.4 7.2e-68 310.7 2.0e-85 253 2e-68
ATP6_MOUSE ATP synthase a chai 246.4 4.5e-66 302.9 4.4e-83 245 5e-66
ATP6_XENLA ATP synthase a chai 111.9 1.4e-25 125.9 8.7e-30 142 9e-35
ATP6_YEAST ATP synthase a ch  78.7 1.6e-15  90.1 5.7e-19  93 5e-20
ATP6_EMENI ATP synthase a chai 66.3 8.4e-12  76.6 6.8e-15  75 2e-14
ATP6_DROYA ATP synthase a chai 65.6 1.2e-11  75.4 1.4e-14 101 2e-22
ATP6_COCHE ATP synthase a cha 53.6 5.5e-08  60.6 4.6e-10  75 1e-14
ATP6_ECOLI ATP synthase a ch  45.1 2.2e-05  49.1 1.4e-06  42 1e-04
ATP6_TRITI ATP synthase a ch  45.0 3.3e-05  50.7 6.5e-07  83 5e-17
ATP6_TOBAC ATP synthase a chai 40.4 0.00084  47.0 8.6e-06  80 3e-16
ATP6_MAIZE ATP synthase a chai 39.6 0.001  44.9 2.6e-05
ATPI_PEA Chloroplast ATP syn 35.8 0.013  38.0 0.0028
ATPI_SPIOL Chloroplast ATP syn 35.5 0.015  38.0 0.0028
ATPI_ATRBE Chloroplast ATP s  34.0 0.044  36.3 0.0086
ATPI_MARPO Chloroplast ATP syn 33.2 0.075  34.3 0.036
*HBA_ODOVI Hemoglobin subunit a 31.9 0.11*
*AROP_ECOLI Aromatic amino ac 32.1 0.31 31.4 0.5 *
ATPI_EUGGR Chloroplast ATP syn 31.1 0.32 32.2 0.15
ATP6_SYN6 ATP synthase a chai 31.1 0.34 31.8 0.21
TLCA_RICPR ADP,ATP carrier pro 31.5 0.49 29.7 1.7
ATP6_SYNY3 ATP synthase a chai 30.6 0.51 31.8 0.22 28 1.9
ATPI_ORYSA Chloroplast ATP 30.1 0.65 32.2 0.15
*GLUC_MYOSC Glucagon precursor 28.7 0.65 34.4 0.013*
*VP6_BPPH6 Protein P6 29.1 0.85 28.6 1.3*
*GLUC_LEFSP Glucagon precursor 27.7 1. 32.7 0.033*
*ADH1_MOUSE Alcohol dehydrogena 29.8 1.2 34.4 0.013*

```

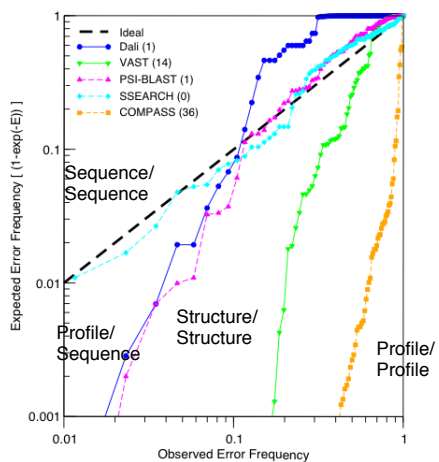
68



### Position-Specific Scores ATP Synthase, 4 iterations

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	bits/pos
BL62 Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2	0.70
46 Q	-2	-1	-2	-2	-4	6	0	1	0	-4	-3	-1	-2	-1	-3	-1	-2	6	4	-3	0.74
%	0	0	0	0	0	54	0	12	0	0	0	0	0	0	0	0	0	13	20	0	
47 Q	-1	-1	3	3	-3	3	3	-2	3	-4	-4	-1	-3	-4	-2	2	-1	-4	-2	-3	0.51
%	0	0	13	20	0	16	19	0	8	0	0	0	0	0	0	24	0	0	0	0	
56 Q	-2	-1	-2	-2	-3	5	2	-4	-1	4	-1	-1	-1	-2	-3	-2	-2	-3	-2	0	0.51
%	0	0	0	0	0	46	13	0	0	41	0	0	0	0	0	0	0	0	0	0	
97 Q	-2	-1	0	-2	-4	4	0	-3	8	-4	-4	-1	-2	-3	-3	-1	-2	-3	0	-4	1.11
%	0	0	0	0	0	35	0	0	65	0	0	0	0	0	0	0	0	0	0	0	
131 Q	3	-1	-1	-1	-2	5	2	-2	-1	-3	-3	0	-2	-4	-2	1	-1	-3	-3	-2	0.52
%	44	0	0	0	0	36	11	0	0	0	0	0	0	0	0	9	0	0	0	0	
152 Q	-2	6	-1	-2	-4	4	0	-3	-1	-4	-3	1	-2	-4	-3	-1	-2	-4	-3	-3	1.00
%	0	77	0	0	0	23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
210 Q	-2	0	-1	-1	-4	7	1	-3	0	-4	-3	1	-1	-4	-2	-1	-2	-3	-2	-3	1.13
%	0	0	0	0	0	100	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

### Accuracy of statistical estimates



- SSEARCH (Smith-Waterman) provides very accurate statistical estimates
- PSI-BLAST can provide estimates that off by 10–100-fold

## Why does PSI-BLAST fail?

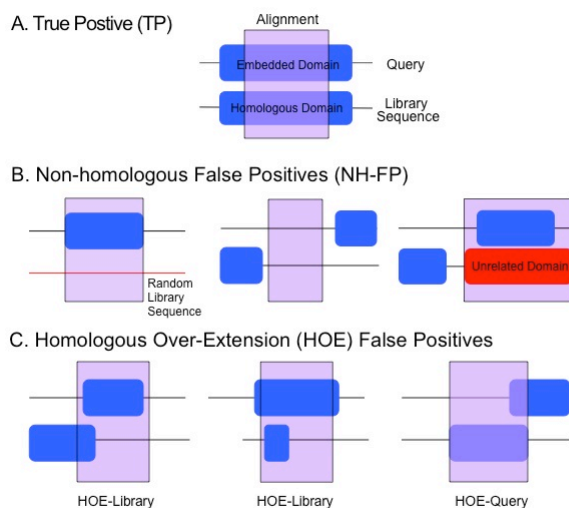


Figure 1

## Sensitive searches with PSI-BLAST

- PSI-BLAST improves sensitivity by building a Position Specific Scoring Matrix (PSSM)
  - models ancestral sequence (consensus distribution)
  - similar to PFAM HMM (but less sophisticated weights, gaps)
- Sensitivity improves with additional iterations
  - model moves to base of tree
- Statistical estimates are difficult
  - once a sequence is in, it is “significant” - validation must be done before a sequence is included
- Very diverse families may not produce a well defined PSSM
  - similar problems with HMMs have led to “clans”

## Sequence Similarity - Conclusions

- Homologous sequences share a common ancestor, but most sequences are non-homologous
- Always compare Protein Sequences
- Sequence Homology can be reliably inferred from statistically significant similarity (non-homology cannot from non-similarity)
- Homologous proteins share common structures, but not necessarily common functions
- Sequence statistical significance estimates are accurate (verify this yourself)  $10^{-6} < E() < 10^{-3}$  is statistically significant
- Scoring matrices set evolutionary look back horizons - not every discovery is distant
- PSI-BLAST can be more sensitive, but with lower statistical accuracy

75

## Discussion questions

1. What is the difference between similarity and homology? When does high identity not imply homology? What conclusions can be drawn from homology?
2. What is the difference between homology and common ancestry?
3. When the *M. janaschii* genome was first sequenced, Venter and his colleagues stated that almost 60% of the open reading frames (proteins or genes) were novel to this organism. (For eubacterial like *E. coli* or *H. influenzae*, a similar number would be 20 - 40%.) On what would they base such a statement? Is it likely to be correct?
4. Name two reasons why protein sequence comparison is more effective (longer evolutionary look-back time) than DNA sequences?
5. What is the range of an expectation value (E()-value)? If you compare a sequence to 50,000 random (unrelated) sequences, what should the expectation value for the highest of the 50,000 similarity scores be (on average)?
6. In a sequence similarity database search, you identify a statistically significant similarity ( $E() < 0.005$ ), but the alignment is relatively short (50 aa). How might you determine whether the alignment reflects a genuine homology, or a random sequence match?
7. How can a sequence be homologous if you search a small database (e.g. human, 40,000 sequences), but not share significant similarity if you search a complete database (>4 million sequences)?
8. What scoring matrix should be used to identify protein orthologs that have diverged over the past 100 My (e.g. human/mouse)?
9. What scoring matrix should be used when comparing Illumina 90 nt reads against a protein database?

76