

# Advanced Database Searching

February 13, 2008

Jonathan Pevsner, Ph.D.  
Introduction to Bioinformatics  
Johns Hopkins University

## Copyright notice

Many of the images in this powerpoint presentation are from *Bioinformatics and Functional Genomics* by Jonathan Pevsner (ISBN 0-471-21004-8). Copyright © 2003 by John Wiley & Sons, Inc.

These images and materials may not be used without permission from the publisher. We welcome instructors to use these powerpoints for educational purposes, but please acknowledge the source.

The book has a homepage at <http://www.bioinfbook.org> including hyperlinks to the book chapters.

## Outline of tonight's lecture

### Specialized BLAST sites

Finding distantly related proteins: PSI-BLAST  
PHI-BLAST

Profile Searches: Hidden Markov models

BLAST-like tools for genomic DNA  
PatternHunter  
Megablast  
BLAT  
BLASTZ

BLAST for gene discovery: Find-a-gene

## Specialized BLAST servers

Species-specific BLAST sites

Molecule-specific BLAST sites

Specialized algorithms (WU-BLAST 2.0)

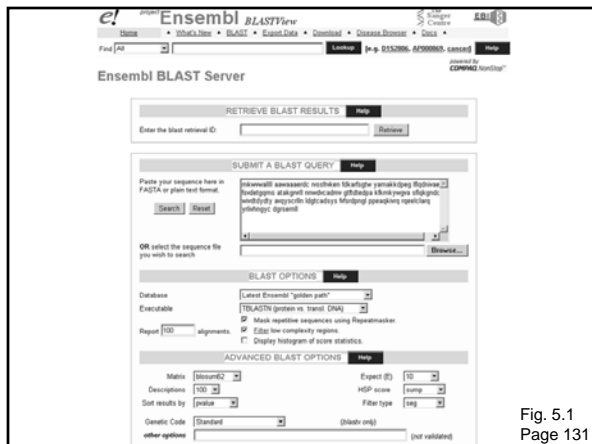
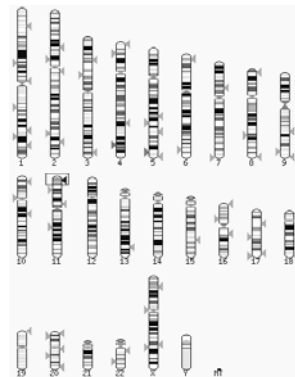


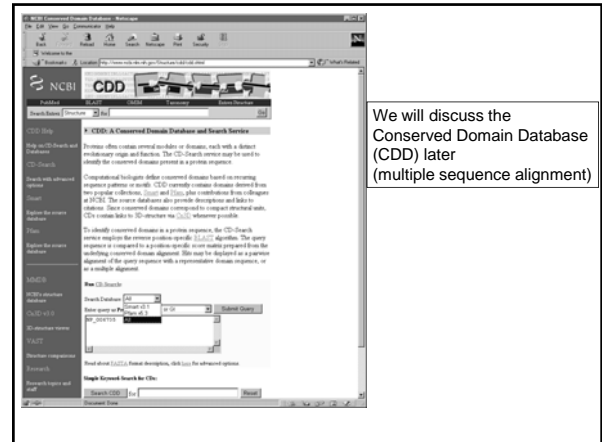
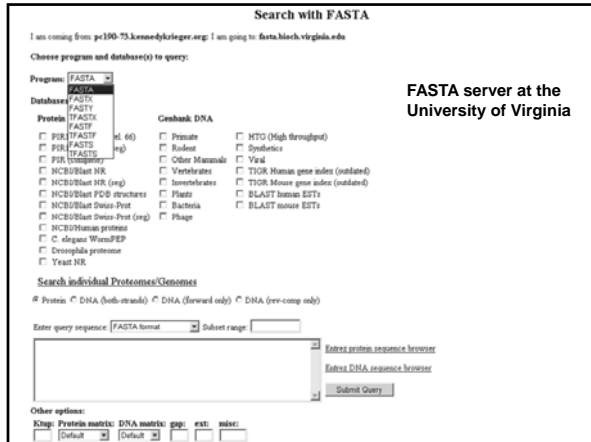
Fig. 5.1  
Page 131

## Ensembl BLAST output includes an ideogram



2e Fig. 5.1





### Outline of today's lecture

---

- Specialized BLAST sites
- Finding distantly related proteins: PSI-BLAST  
PHI-BLAST
- Profile Searches: Hidden Markov models
- BLAST-like tools for genomic DNA
  - PatternHunter
  - Megablast
  - BLAT
  - BLASTZ
- BLAST for gene discovery: Find-a-gene

### Position specific iterated BLAST: PSI-BLAST

---

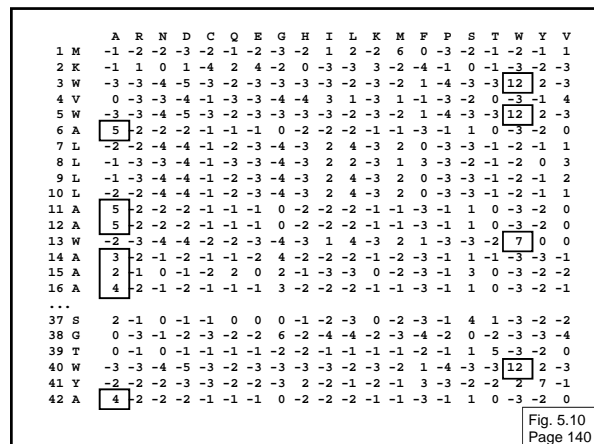
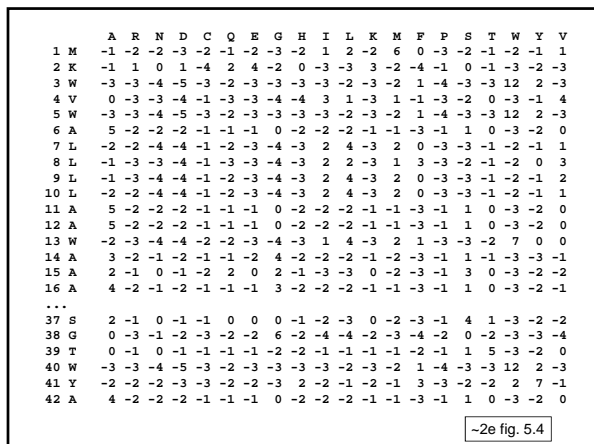
The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a scoring matrix that is customized to your query.

### PSI-BLAST is performed in five steps

---

- [1] Select a query and search it against a protein database
- [2] PSI-BLAST constructs a multiple sequence alignment then creates a "profile" or specialized position-specific scoring matrix (PSSM)

<p>730496 66 200679 63 206589 34 2136812 2 132408 65 267584 44 267585 44 8777408 63 6687453 60 10697027 81 13645517 1 13925316 38 131649 65</p>	<p>66 63 34 2 65 44 44 63 60 81 1 38 65</p>	<p>FTVDENGQHSATAKGRVRLFRNVDCADHIGSFDTDEPARKFKHYGVASFQKGNDDH 125              FSVDEKQHNSATAKGRVRLLSNVVVCADHVGTFDTDEPARKFKHYGVASFQKGNDDH 122              FSVDEKQHNSATAKGRVRLLSNVVVCADHVGTFDTDEPARKFKHYGVASFQKGNDDH 93              HSATAKGRVRLLSNVVVCADHVGTFDTDEPARKFKHYGVASFQKGNDDH 53              FRIEDNGKTTATAGRVRILLKLELCANVVGTFETNDPAKFKHYGVGALALERGLDDH 124              FSVDESGKVTATAGRVRILLNVVHCANVVGTFETDPAKFKHYGVGAAATLQGNDDH 103              FSVDSGKVTATAQGRVILLNVVHCANVVGTFETDPAKFKHYGVGAAATLQGNDDH 103              FTVEEDGHTATAGRVRILLNVVHCADHRAFTETDPAKFKHYGVGAAATLQGNDDH 122              FKVEEDGHTATAGRVRILLNVVHCANVVGTFETDPAKFKHYGVGAAATLQGNDDH 119              FKVEEDGHTATATGRVILLNVVHCANVVGTFETDPAKFKHYGVGAAATLQGNDDH 140              NVGTFDTDEPARKFKHYGVGAAATLQGNDDH 32              FSVDSGKHTATAQGRVILLNVVHCANVVGTFETDPAKFKHYGVGAAATLQGNDDH 97              YTVEEDGHTASSGRVKLFGVVICADHAAQYTPPTPAKHNTYQGLASTLSSGGINNY 126</p>
		<p>↑                    ↑                    ↑                    ↑                    ↑</p> <p><b>R,I,K      C      D,E,T      K,R,T      N,L,Y,G</b></p>



### PSI-BLAST is performed in five steps

- [1] Select a query and search it against a protein database
- [2] PSI-BLAST constructs a multiple sequence alignment then creates a "profile" or specialized position-specific scoring matrix (PSSM)
- [3] The PSSM is used as a query against the database
- [4] PSI-BLAST estimates statistical significance (E values)

Page 138

### Results of a PSI-BLAST search

Iteration	# hits	# hits > threshold
1	104	49
2	173	96
3	236	178
4	301	240
5	344	283
6	342	298
7	378	310
8	382	320

Table 5-4  
Page 141

### PSI-BLAST is performed in five steps

- [1] Select a query and search it against a protein database
- [2] PSI-BLAST constructs a multiple sequence alignment then creates a "profile" or specialized position-specific scoring matrix (PSSM)
- [3] The PSSM is used as a query against the database
- [4] PSI-BLAST estimates statistical significance (E values)
- [5] Repeat steps [3] and [4] iteratively, typically 5 times. At each new search, a new profile is used as the query.

Page 138

### Results of a PSI-BLAST search

Iteration	# hits	# hits > threshold
1	104	49
2	173	96
3	236	178
4	301	240
5	344	283
6	342	298
7	378	310
8	382	320

Table 5-4  
Page 141

### PSI-BLAST search: human RBP versus RefSeq, iteration 1

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

	Score	E	UC
	(Bits)	Value	
<input checked="" type="checkbox"/> ref NP_004735.2	388	1e-111	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001628.1	57.4	7e-09	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001018059.1	36.2	0.019	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001624.1	35.8	0.021	<input checked="" type="checkbox"/>

Run PSI-Blast Iteration 2

Sequences with E-value WORSE than threshold

	Score	E	UC
	(Bits)	Value	
<input type="checkbox"/> ref NP_000597.1	33.9	0.077	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_076222.1	28.5	3.8	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_066015.2	27.3	7.5	<input checked="" type="checkbox"/>

Run PSI-Blast Iteration 2

### PSI-BLAST search: human RBP versus RefSeq, iteration 2

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

	Score	E	UC
	(Bits)	Value	
<input checked="" type="checkbox"/> ref NP_004735.2	388	1e-102	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001628.1	149	2e-36	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001018059.1	124	5e-32	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001624.1	123	2e-29	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref XP_001129927.1	70.0	1e-12	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_044626.1	69.3	2e-12	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_000945.1	42.5	1e-04	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_076222.1	39.6	0.002	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_048564.1	32.2	0.002	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001018076.1	28.5	0.003	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_000597.1	26.9	0.010	<input checked="" type="checkbox"/>

Run PSI-Blast Iteration 3

Sequences with E-value WORSE than threshold

	Score	E	UC
	(Bits)	Value	
<input type="checkbox"/> ref NP_002288.1	31.5	0.48	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_004524.2	30.4	1.0	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_075903.2	30.0	1.2	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_055983.1	29.2	2.1	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_001033982.1	28.8	2.8	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_040146.2	28.4	3.4	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_055297.1	28.1	4.4	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_045184.1	27.3	9.3	<input checked="" type="checkbox"/>

Run PSI-Blast Iteration 3

### PSI-BLAST search: human RBP versus RefSeq, iteration 3

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

	Score	E	UC
	(Bits)	Value	
<input checked="" type="checkbox"/> ref NP_004735.2	318	2e-99	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_000597.1	180	6e-34	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001628.1	133	7e-32	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_076222.1	128	2e-30	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001018059.1	119	1e-27	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001624.1	112	2e-25	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001129927.1	69.8	7e-10	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_044626.1	69.4	8e-10	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_000945.1	59.1	4e-09	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_048564.1	42.7	2e-04	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_001018076.1	42.3	3e-04	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_045184.1	41.5	4e-04	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_055297.1	38.4	0.003	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_055396.1	36.5	0.016	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> ref NP_002288.1	34.2	0.039	<input checked="" type="checkbox"/>

Run PSI-Blast Iteration 4

Sequences with E-value WORSE than threshold

	Score	E	UC
	(Bits)	Value	
<input type="checkbox"/> ref NP_048631.2	31.1	0.66	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_001033982.1	30.7	0.82	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_026741.1	30.3	0.99	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_004565.1	30.3	0.99	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_048573.2	27.6	5.9	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_001074.1	27.6	5.9	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_026742.1	27.6	5.9	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_000977.1	27.2	6.5	<input checked="" type="checkbox"/>
<input type="checkbox"/> ref NP_004534.2	27.2	9.6	<input checked="" type="checkbox"/>

Run PSI-Blast Iteration 4

### RBP4 match to ApoD, PSI-BLAST iteration 1

```
>> ref|NP_001628.1|  apolipoprotein D precursor [Homo sapiens]
Length=189

Score = 57.4 bits (137), Expect = 3e-07, Method: Composition-based stats.
Identities = 47/151 (31%), Positives = 78/151 (51%), Gaps = 39/151 (25%)

Query 29 VYENFKARFSGTIVYAMAKDPEGLFDLQDNIYAEFVDTGHSATGAPVFLMMHVC 88
      +VEMFD ++ G VY + +K F I A ++E G +++++LM ++
Sbjct 33 YQENFDVWYKLGWYVEI-EKIPITFENGRCIQANTSLMNG-----KIKYLMQ-ELR 82

Query 69 ANMVGTFIDTE-----DPAKFSKCY-NGVASFQKQGDHIVDTYDTATVTC 130
      AD QT E ++AK ++K+ U + S ++W TDYF YK TDC
Sbjct 83 AD--GTVNOIEGATFVNLTEFAKLVKFSVDFPMS-----APVILATYVNTATVTC 134

Query 139 ---RLMLDGTADSTVFFVSDRNGLPPE 165
      +L ++D +++++ +R+R LPPE
Sbjct 135 TCIQLFMDV-----FAMILARPN-LPPE 150
```

2e Fig. 5.6

### RBP4 match to ApoD, PSI-BLAST iteration 2

```
>> ref|NP_001628.1|  apolipoprotein D precursor [Homo sapiens]
Length=189

Score = 175 bits (443), Expect = 1e-42, Method: Composition-based stats.
Identities = 45/163 (27%), Positives = 77/163 (47%), Gaps = 31/163 (19%)

Query 14 GSGRAERDCRVSSFRVZNFVDFKARFSGTIVYAMAKDPEGLFDLQDNIYAEFVDTGHS 73
      GSA + + VEMFD ++ G VY + +K F I A ++E G G++
Sbjct 18 AKGQAFHLGKCFKPFVQENFDVWYKLGWYVEI-EKIPITFENGRCIQANTSLMNGKIKV 76

Query 74 TAK-----GRVFLMMHVCADHMGVTFIDTEFAKFSKCY-NGVASFQKQGDHIVDT 127
      + G V + T + +PAK ++K+ U + S ++W TDYF Y
Sbjct 77 LMQELRADGTVMQIEG-----EATFVNLTEFAKLVKFSVDFPMS-----APVILATY 123

Query 120 DTVYTAVYTCR---LLMLDGTADSTVFFVSDRNGLPPEA 166
      DTVF YK TDC I ++D +++++ +R+R LPPE
Sbjct 124 DVENYALVYTCRITQLFMDV-----FAMILARPN-LPPE 159
```

2e Fig. 5.6

### RBP4 match to ApoD, PSI-BLAST iteration 3

```
>> ref|NP_001628.1|  apolipoprotein D precursor [Homo sapiens]
Length=189

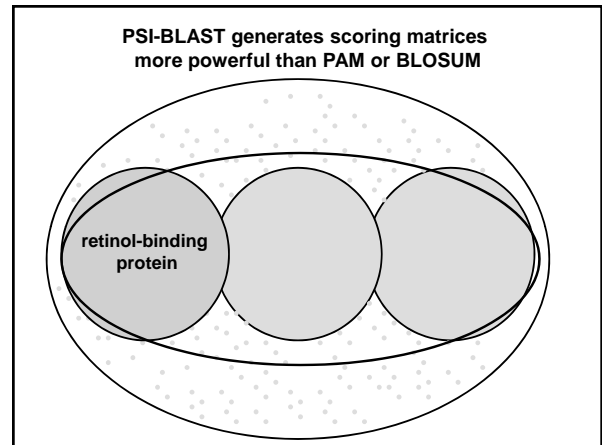
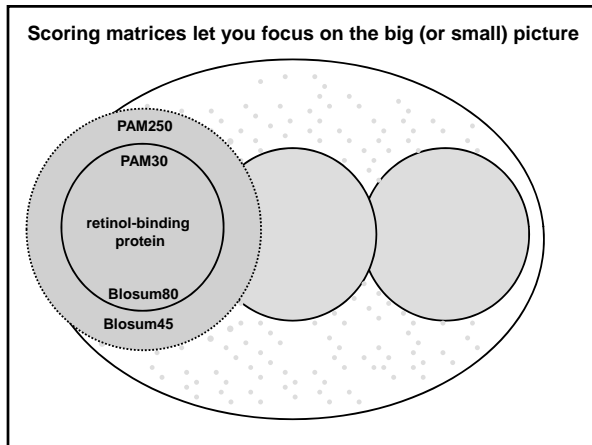
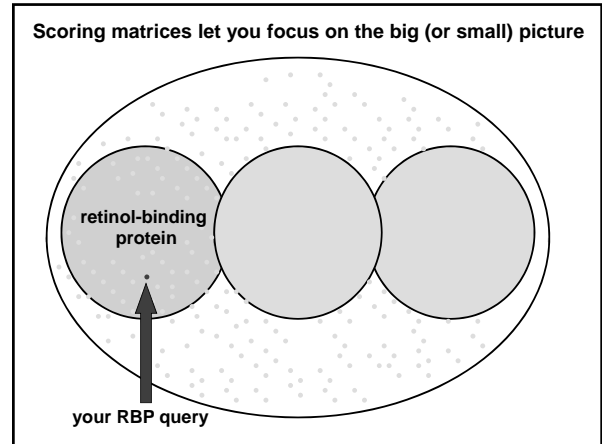
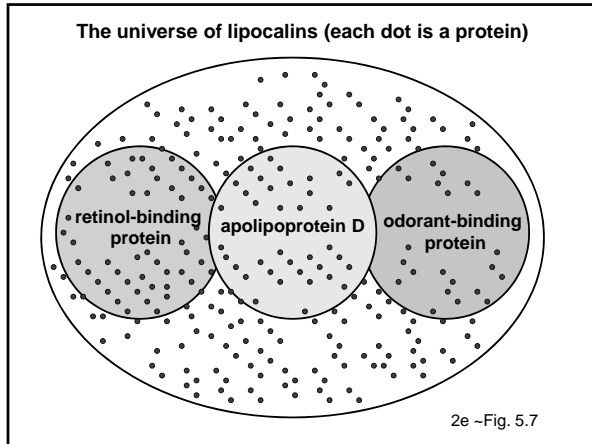
Score = 146 bits (368), Expect = 6e-34, Method: Composition-based stats.
Identities = 41/163 (25%), Positives = 76/163 (46%), Gaps = 20/163 (12%)

Query 14 GSGRAERDCRVSSFRVZNFVDFKARFSGTIVYAMAKDPEGLFDLQDNIYAEFVDTGHS 73
      GSA + + VEMFD ++ G VY + +K F I A ++E G G++
Sbjct 18 AKGQAFHLGKCFKPFVQENFDVWYKLGWYVEI-EKIPITFENGRCIQANTSLMNGKIKV 76

Query 74 TAKGRVFLMMHVCADHMGVTFIDTEFAKFSKCY-NGVASFQKQGDHIVDTYDTY 132
      + +R + + T + +PAK ++K+ U + S ++W TDYF Y
Sbjct 77 LMQ-ELRADGTVMQI-EKATFVNLTEFAKLVKFSVDFPMS-----APVILATYKMY 128

Query 133 AVYTCR---LLMLDGTADSTVFFVSDRNGLPPEAGKIVR 171
      +R TDC L ++D +++++ +R+R F +
Sbjct 129 ALVYTCRITQLFMDV-----FAMILARPN-LPPE 165
```

2e Fig. 5.6



**PSI-BLAST: performance assessment**

---

Evaluate PSI-BLAST results using a database in which protein structures have been solved and all proteins in a group share  $\leq 40\%$  amino acid identity.

Page 143

**PSI-BLAST: the problem of corruption**

---

PSI-BLAST is useful to detect weak but biologically meaningful relationships between proteins.

The main source of false positives is the spurious amplification of sequences not related to the query. For instance, a query with a coiled-coil motif may detect thousands of other proteins with this motif that are not homologous.

Once even a single spurious protein is included in a PSI-BLAST search above threshold, it will not go away.

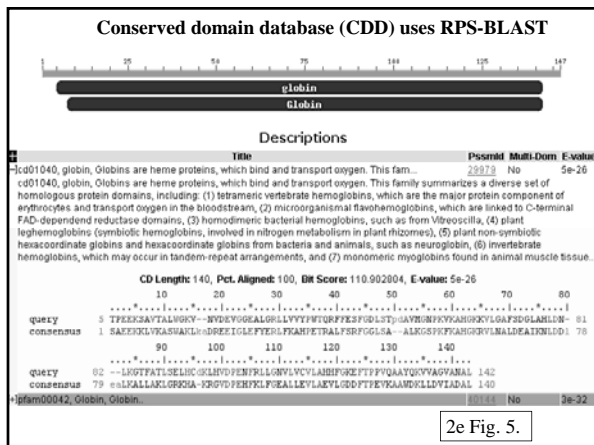
Page 144

### PSI-BLAST: the problem of corruption

Corruption is defined as the presence of at least one false positive alignment with an E value < 10<sup>-4</sup> after five iterations.

Three approaches to stopping corruption:

- [1] Apply filtering of biased composition regions
- [2] Adjust E value from 0.001 (default) to a lower value such as E = 0.0001.
- [3] Visually inspect the output from each iteration. Remove suspicious hits by unchecking the box.



### PHI-BLAST: Pattern hit initiated BLAST

Launches from the same page as PSI-BLAST

Combines matching of regular expressions with local alignments surrounding the match.

### PHI-BLAST: Pattern hit initiated BLAST

Launches from the same page as PSI-BLAST

Combines matching of regular expressions with local alignments surrounding the match.

Given a protein sequence S and a regular expression pattern P occurring in S, PHI-BLAST helps answer the question: What other protein sequences both contain an occurrence of P and are homologous to S in the vicinity of the pattern occurrences? PHI-BLAST may be preferable to just searching for pattern occurrences because it filters out those cases where the pattern occurrence is probably random and not indicative of homology.

### Align three lipocalins (RBP and two bacterial lipocalins)

```

1
ecb1c MRLPLVAAA TAAFLVVACS SPTPPRGVTV VNNFDKRYL GTWYEIARFD 50
vc MRAIFLLCS V...LLNGCL G..MPESVXP VSDFELNNYL GKWYEVARLD
hsrbp ---MKVWVAL LLLAAWAAE RDCRVSSFRV KNFDFKARFS GTWYAMAKKD
    
```





### Human RBP PSI-BLAST search against bacteria

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:	Score (Bits)	E Value
g111947045 ref YP_01613353.1  outer membrane lipoprotein (l...	45.7	0.005

Run PSI-Blast Iteration 2

### PHI-BLAST: RBP search against bacteria

Sequences with pattern at position 48 and E-value BETTER than threshold

Sequences producing significant alignments:	Score (Bits)	E Value
ref YP_01613353.1  outer membrane lipoprotein (lipocalin) [Al...	24.7	1e-05
ref YP_212549.1  hypothetical protein BF2935 [Bacteroides fragil...	22.5	5e-05
ref NP_813404.1  putative sugar nucleotide epimerase [Bactero...	22.1	7e-05
ref YP_100376.1  putative sugar nucleotide epimerase [Bactero...	21.4	1e-04
ref YP_677044.1  outer membrane lipoprotein (lipocalin) [Cyto...	20.3	3e-04
ref YP_01006814.1  lipoprotein B1c [Prochlorococcus marinus str...	19.2	5e-04
ref YP_001101630.1  outer membrane lipoprotein (lipocalin) [H...	18.9	7e-04
ref YP_01244727.1  conserved hypothetical protein [Flavobacte...	18.8	7e-04
ref YP_341541.1  outer membrane lipoprotein (lipocalin) [Faeu...	18.8	7e-04
ref YP_01065301.1  lipoprotein B1c [Vibrio sp. MED222]	18.1	0.001
ref YP_01578373.1  Lipocalin-like [Delftia acidovorans SPH-1]	17.3	0.002
ref YP_01474066.1  hypothetical protein YKx2w_02009361 [Vibrio s...	16.6	0.003
ref YP_00592335.1  Lipocalin-related protein and Bos/Can/Equ ...	16.2	0.004

Run PSI-Blast Iteration 2

2e Fig. 5.10

### Outline of today's lecture

---

Specialized BLAST sites

Finding distantly related proteins: PSI-BLAST  
PHI-BLAST

Profile Searches: Hidden Markov models

BLAST-like tools for genomic DNA  
PatternHunter  
Megablast  
BLAT  
BLASTZ


BLAST for gene discovery: Find-a-gene

### Multiple sequence alignment to profile HMMs

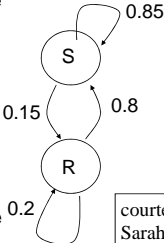
- in 90's people began to see that aligning sequences to profiles gave much more information than pairwise alignment alone.
- Hidden Markov models (HMMs) are "states" that describe the probability of having a particular amino acid residue at arranged in a column of a multiple sequence alignment
- HMMs are probabilistic models
- Like a hammer is more refined than a blast, an HMM gives more sensitive alignments than traditional techniques such as progressive alignments

Page 325

### Simple Markov Model



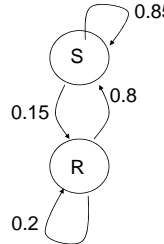
Rain = dog may not want to go outside  
Sun = dog will probably go outside



Markov condition = no dependency on anything but nearest previous state ("memoryless")

courtesy of Sarah Wheelan

### Simple Hidden Markov Model



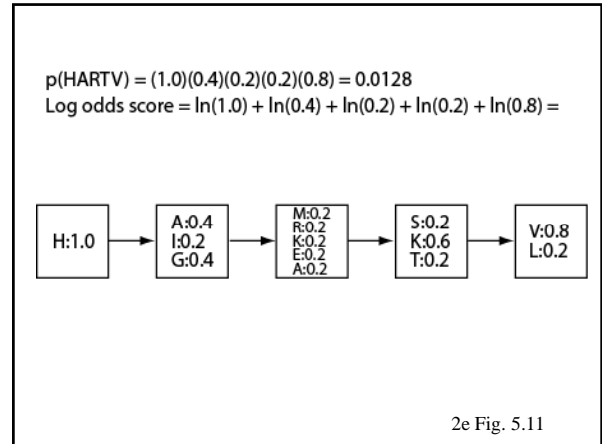
$P(\text{dog goes out in rain}) = 0.1$   
 $P(\text{dog goes out in sun}) = 0.85$

Observation: YNNYYNNNNY  
(Y=goes out, N=doesn't go out)

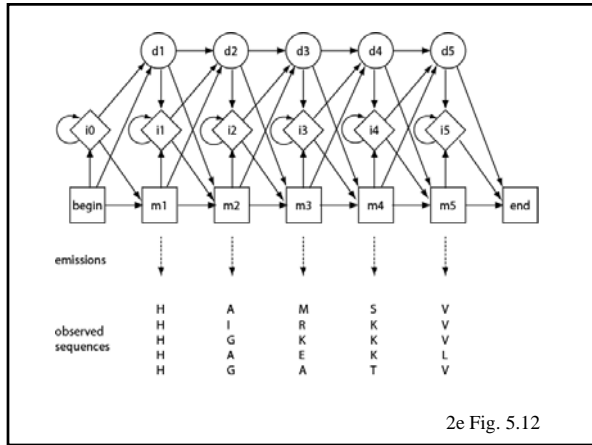
What is underlying reality (the hidden state chain)?

		Probability	position				
			1	2	3	4	5
1D8U	HAMSV	p(H)	1.0				
1OJ6A	HIRKV	p(A)		0.4			
2hhbB	HGKKV	p(I)		0.2			
1FSL	HAEKL	p(G)		0.4			
2MM1	HGATV	p(M)			0.2		
		p(R)			0.2		
		p(K)			0.2		
		p(E)			0.2		
		p(A)			0.2		
		p(S)				0.2	
		p(K)				0.6	
		p(T)				0.2	
		p(V)					0.8
		p(L)					0.2

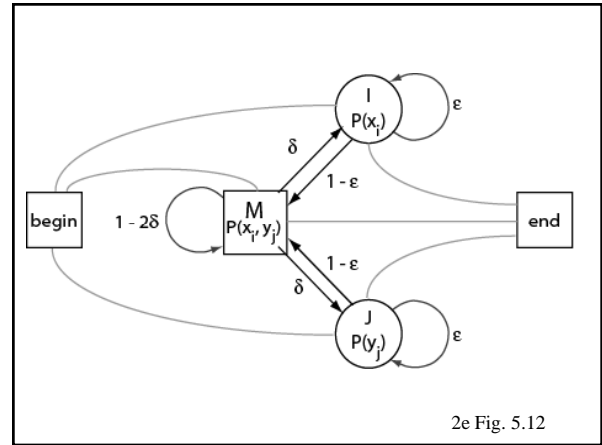
2e Fig. 5.11



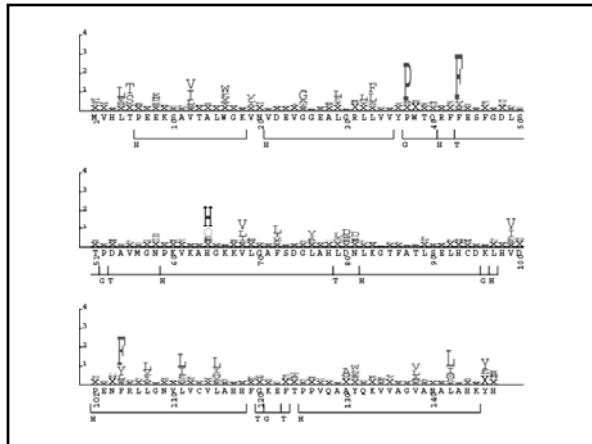
2e Fig. 5.11



2e Fig. 5.12



2e Fig. 5.12



**HMMER: build a hidden Markov model**

- Determining effective sequence number ... done. [4]
- Weighting sequences heuristically ... done.
- Constructing model architecture ... done.
- Converting counts to probabilities ... done.
- Setting model name, etc. ... done. [x]

Constructed a profile HMM (length 230)

- Average score: 411.45 bits
- Minimum score: 353.73 bits
- Maximum score: 460.63 bits
- Std. deviation: 52.58 bits

Fig. ~5.13

### HMMER: calibrate a hidden Markov model

HMM file: lipocalins.hmm  
 Length distribution mean: 325  
 Length distribution s.d.: 200  
 Number of samples: 5000  
 random seed: 1034351005  
 histogram(s) saved to: [not saved]  
 POSIX threads: 2

HMM : x  
 mu : -123.894508  
 lambda : 0.179608  
 max : -79.334000

### HMMER: search an HMM against GenBank

Scores for complete sequences (score includes all domains):

Sequence	Description	Score	E-value	N
gi 20888903 ref XP_129259.1	(XM_129259) ret	461.1	1.9e-133	1
gi 132407 sp P04916 RETB_RAT	Plasma retinol-	458.0	1.7e-132	1
gi 20548126 ref XP_005907.5	(XM_005907) sim	454.9	1.4e-131	1
gi 5803139 ref NP_006735.1	(NM_006744) ret	454.6	1.7e-131	1
gi 20141667 sp P02753 RETB_HUMAN	Plasma retinol-	451.1	1.9e-130	1
gi 16767588 ref NP_463203.1	(NC_003197) out	318.2	1.9e-90	1

gi|5803139|ref|NP\_006735.1|: domain 1 of 1, from 1 to 195: score 454.6, E = 1.7e-131  
 \*->mkwVnkLLLaLagvfgaAeRdAafsvgkCrvpsPFGfVkeNFDv  
 mkwV--LILLLa + +sAeRd Crvs+ fVkeNFD+  
 gi|5803139 1 MKWVWALLLLAA--W--AAAEER-----CVVSS----FRVKEKDFK 33  
 erylGtWYeIaKkDpsFBrGLllgkItAeySlEhGsMaataeGrivVL  
 ++GQWf++sKkDp E GL-lgd-I--keS+++E-G-MaataeGr++L  
 gi|5803139 34 ARFSGTWYMAKEDP--E-GLFLQDNIvAEFVSDETQMSATAKGRVLL 80  
 eNkelcADkvGVtqIGeaaevLItadPakIkIkyGvAsfIgpGfdy  
 N+++cAdV+QVt++E dPak-kKyGvAsfIgpGfd+  
 gi|5803139 81 NNNDVCDVMADMTFTDTE-----DPAKFKMKYGVAsFLQKGNDDH 120

### HMMER: search an HMM against GenBank match to a bacterial lipocalin

gi|16767588|ref|NP\_463203.1|: domain 1 of 1, from 1 to 177: score 318.2, E = 1.9e-90  
 \*->mkwVnkLLLaLagvfgaAeRdAafsvgkCrvpsPFGfVkeNFDv  
 M-LLa +A a ++Af++++C+sp+PF+G++V++NFD+  
 gi|1676758 1 ---NRLLPVVA-----AVTA-AFLVACSSPTPPKGVTVVNFDA 36  
 erylGtWYeIaKkDpsFBrGLllgkItAeySlEhGsMaataeGrivVL  
 +rylGtWYeIa+ D++FEGl + +LaySl++ +G+i+V+  
 gi|1676758 37 KRYLGTWYELARLDHRPERGL---EQVTATYSLRD-----DGGINVI 75  
 eNkelcADkvGVtqIGeaaevLItadPakIkIkyGvAsfIgpGfdy  
 Nk+++D+ +++ +EG+a ++t+ P +++lK+ Sf++P++++y  
 gi|1676758 76 -NKGYNDR-EMWQKTEGKA---YPTGSPNRAALKV----SFPQPPYGY 116

### HMMER: search an HMM against GenBank

Scores for complete sequences (score includes all domains):

Sequence	Description	Score	E-value	N
gi 3041715 sp P27485 RETB_PIG	Plasma retinol-	614.2	1.6e-179	1
gi 89271 pir [A39486	plasma retinol-	613.9	1.9e-179	1
gi 20888903 ref XP_129259.1	(XM_129259) ret	608.8	6.8e-178	1
gi 132407 sp P04916 RETB_RAT	Plasma retinol-	608.0	1.1e-177	1
gi 20548126 ref XP_005907.5	(XM_005907) sim	607.3	1.9e-177	1
gi 20141667 sp P02753 RETB_HUMAN	Plasma retinol-	605.3	7.2e-177	1
gi 5803139 ref NP_006735.1	(NM_006744) ret	600.2	2.6e-175	1

gi|5803139|ref|NP\_006735.1|: domain 1 of 1, from 1 to 199: score 600.2, E = 2.6e-175  
 \*->mkwVnkLLLaLagvfgaAeRdAafsvgkCrvpsPFGfVkeNFDv  
 m+WwL+LLea+ a+AEEDCrvsFRVKEKDFKARFSGWYAAK  
 gi|5803139 1 MKWVWALLLLAA--W--AAAEERCVSSFRVKEKDFKARFSGWYAAK 45  
 KDPEGLFLQDNIvAEFVSDEKQhmsATAKGRVRLlNnWdVcAdmVgtFlD  
 KDPEGLFLQDNIvAEFVSDE+G+msATAKGRVRLlNnWdVcAdmVgtFlD  
 gi|5803139 46 KDPEGLFLQDNIvAEFVSDETQMSATAKGRVRLlNnWdVcAdmVgtFlD 95  
 tEDPAKFKMKYGVAsFLqKgnDDHwIldtDvdtfAvqYsCRllnLDGtC  
 tEDPAKFKMKYGVAsFLqKgnDDHwIldtDvdtfAvqYsCRllnLDGtC  
 gi|5803139 96 TEDPAKFKMKYGVAsFLQKGNDDHwIVDvDYDTAVQYsCRllnLDGtC 145

### Outline of today's lecture

- Specialized BLAST sites
- Finding distantly related proteins: PSI-BLAST  
PHI-BLAST
- Profile Searches: Hidden Markov models

BLAST-like tools for genomic DNA  
 PatternHunter  
 Megablast  
 BLAT  
 BLASTZ

BLAST for gene discovery: Find-a-gene

### BLAST-related tools for genomic DNA

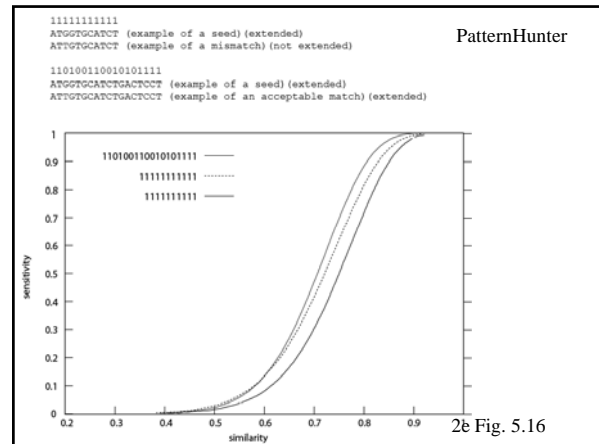
- The analysis of genomic DNA presents special challenges:
- There are exons (protein-coding sequence) and introns (intervening sequences).
  - There may be sequencing errors or polymorphisms
  - The comparison may be between related species (e.g. human and mouse)

## BLAST-related tools for genomic DNA

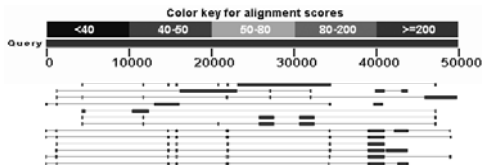
Recently developed tools include:

- MegaBLAST at NCBI.
- BLAT (BLAST-like alignment tool). BLAT parses an entire genomic DNA database into words (11mers), then searches them against a query. Thus it is a mirror image of the BLAST strategy. See <http://genome.ucsc.edu>
- SSAHA at Ensembl uses a similar strategy as BLAT. See <http://www.ensembl.org>

Page 136



## MegaBLAST at NCBI

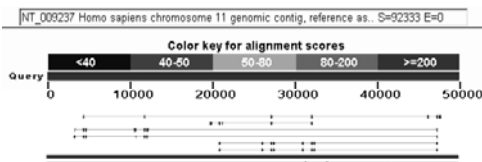


## MegaBLAST

Sequence producing significant alignments:  
(Click headers to sort columns)

Accession	Description	Max score	Total score	Query coverage	t value	Max ident
U01029.1	Panjo pygmaeus gamma-1 and gamma-2 globin genes, cDNA	1,055e+04	1,204e+04	21%	0.0	95%
U18029.1	Orangutan (P.pygmaeus) beta- and eta-globin pseudogenes	1,035e+04	1,156e+04	15%	0.0	94%
U02025.1	Orangutan epsilon-globin gene with Alu repeats in flanking region	6547	8190	10%	0.0	96%
U18236.1	Orangutan beta- and delta-globin gene intergenic region	3122	5899	7%	0.0	96%
U18233.1	Orangutan delta globin gene, complete cds	2818	4516	5%	0.0	97%
U18232.1	Orangutan gamma-2-fetal globin gene, complete cds	2550	4424	9%	0.0	94%
U18231.1	Orangutan gamma-1-fetal globin gene, complete cds	4320	6667	9%	0.0	94%

## MegaBLAST: 50 kilobases of the globin locus



## To access BLAT, visit <http://genome.ucsc.edu>

UCSC Genome Bioinformatics

Genomes Gene Sorter Blat PCR Tables FAQ Help

**About the UCSC Genome Bioinformatics Site**

This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also shows the CPGP (CpG Islands) regions in 12 species and provides a portal to the ENCODE project.

**News**

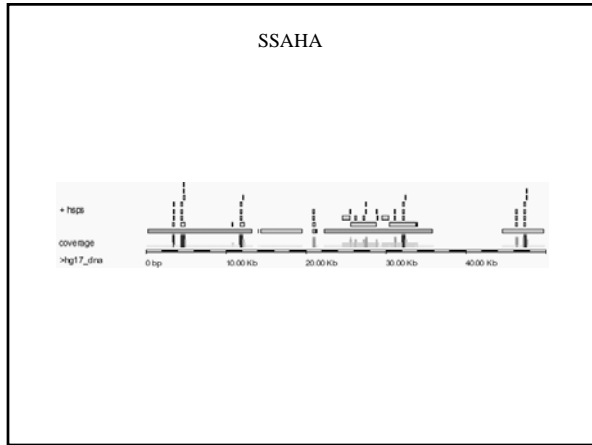
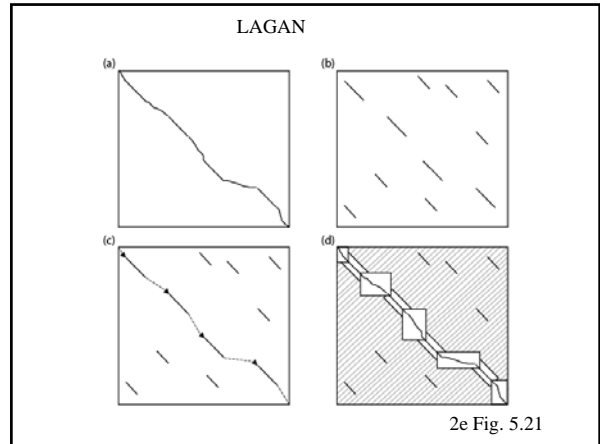
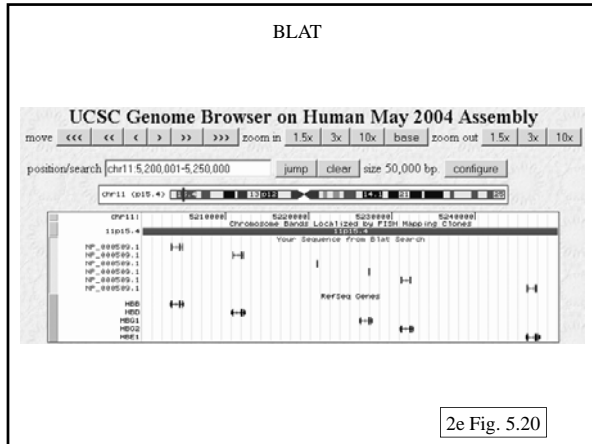
18 September 2004 - Tetraodon Genome Assembly in Genome Browser

The *Tetraodon* v7 Tetraodon euphratica genome assembly is now available in the UCSC Genome Browser and Blat server. The assembly, UCSC version webQ1 dated Feb. 2004, is the result of a collaboration between Genoscope and the Broad Institute of MIT and Harvard.

The v7 assembly was constructed using the whole genome shotgun (WGS) approach, resulting in a sequence coverage of about 7 X. The assembly contains 45,609 contigs and 23,773 scaffolds generated by the Arachne program and covers more than 90% of the genome.

"BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 33 bases, and sometimes find them down to 20 bases. BLAT on proteins finds sequences of 80% and greater similarity of length 20 amino acids or more. In practice DNA BLAT works well on primates, and protein blat on land vertebrates." --BLAT website





**Outline of today's lecture**

---

- Specialized BLAST sites
- Finding distantly related proteins: PSI-BLAST  
PHI-BLAST
- Profile Searches: Hidden Markov models
- BLAST-like tools for genomic DNA
  - PatternHunter
  - Megablast
  - BLAT
  - BLASTZ

BLAST for gene discovery: Find-a-gene

**BLAST for gene discovery**

---

You can use BLAST to find a "novel" gene

Page 147

**BLAST for gene discovery**

---

You can use BLAST to find a "novel" gene

You will need to do this for 40% of your grade.

In the first 8 years of this course,  
everyone has succeeded at this exercise.

Page 147



```

Sequences producing significant alignments:
Score E
(bits) Value
gml1Ranger_60118.typhi_Salmonella typhi CT18 Salmonella typ... 37 0.13
gml1W92C_310271.papua_P_97A.D.16410 Salmonella paratyphi A ... 36 0.17
gml1W92C_922871.stm12-A2A.Contig1366 Salmonella typhisurium L... 35 0.29
gml1W92C_922871.stm12-A2A.Contig1366 Salmonella typhisurium L... 35 0.29
gml1A0001741 Escherichia coli O157:H7EDL93, complete genome 35 0.50
gml1A0000071 Escherichia coli O157:H7, complete genome 35 0.50
gml1D00261.EC051 Escherichia coli K-12 M91655 complete genome 35 0.50
gml1D10C_081501.gshublin_6141_12_23 Salmonella dublin unfinis... 33 1.1
gml1TIG2_57591.shahavirale_ENTK0718 1600 2000 1800 12 751 Ent... 32 2.5
gml1TIG2_1581ident_10127 Treponema denticola unfinished fra... 32 3.2
gml1TIG2_C.tspidum_3499 Chlorobium tepidum unfinished fragm... 31 5.5
gml1TIG2_57591.shahavirale_ENTK0718 1600 2000 1800 20 863 Ent... 31 5.5
gml1TIG2_1581ident_10143 Treponema denticola unfinished fra... 31 7.2
gml1Ranger_14961.clostridicille.Contig1558 Clostridium difficile ... 31 7.2
gml1SRSTC_54761C.nibicane.Contig1-1956 size=11146 (bases 2.... 30 9.3
msb1A002161.EB028 Bacillus subtilis complete genome 30 9.3
gml1SRSTC_54761C.nibicane.Contig-2286 size=21641 (EP01) (N... 30 9.3

Alignments
>gml1Ranger_60118.typhi_Salmonella typhi CT18 Salmonella typhi unfinished fragment
Length = 4890036
Score = 36.6 bits (83), Expect = 0.13
Identities = 22/82 (26%), Positives = 38/82 (45%)
Frame = -2
Query: 27 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 86
W NFD R+ GTWY +A+ D + + A +S+ + G + + EG W
Sbjct: 4558772 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK- 4558596

Query: 87 ADMVGTFTDTEPAREKRYKRW 108
+ FT + + A E+ +G
Sbjct: 4558595 TEKGAVTTPSPRAALKVSYFFG 4558530
    
```

(Page 150)

```

Sequences producing significant alignments:
Score E
(bits) Value
gml1Ranger_60118.typhi_Salmonella typhi CT18 Salmonella typ... 37 0.13
gml1W92C_310271.papua_P_97A.D.16410 Salmonella paratyphi A ... 36 0.17
gml1W92C_922871.stm12-A2A.Contig1366 Salmonella typhisurium L... 35 0.29
gml1W92C_922871.stm12-A2A.Contig1366 Salmonella typhisurium L... 35 0.29
gml1A0001741 Escherichia coli O157:H7EDL93, complete genome 35 0.50
gml1A0000071 Escherichia coli O157:H7, complete genome 35 0.50
gml1D00261.EC051 Escherichia coli K-12 M91655 complete genome 35 0.50
gml1D10C_081501.gshublin_6141_12_23 Salmonella dublin unfinis... 33 1.1
gml1TIG2_57591.shahavirale_ENTK0718 1600 2000 1800 12 751 Ent... 32 2.5
gml1TIG2_1581ident_10127 Treponema denticola unfinished fra... 32 3.2
gml1TIG2_C.tspidum_3499 Chlorobium tepidum unfinished fragm... 31 5.5
gml1TIG2_57591.shahavirale_ENTK0718 1600 2000 1800 20 863 Ent... 31 5.5
gml1TIG2_1581ident_10143 Treponema denticola unfinished fra... 31 7.2
gml1Ranger_14961.clostridicille.Contig1558 Clostridium difficile ... 31 7.2
gml1SRSTC_54761C.nibicane.Contig1-1956 size=11146 (bases 2.... 30 9.3
msb1A002161.EB028 Bacillus subtilis complete genome 30 9.3
gml1SRSTC_54761C.nibicane.Contig-2286 size=21641 (EP01) (N... 30 9.3

Alignments
>gml1Ranger_60118.typhi_Salmonella typhi CT18 Salmonella typhi unfinished fragment
Length = 4890036
Score = 36.6 bits (83), Expect = 0.13
Identities = 22/82 (26%), Positives = 38/82 (45%)
Frame = -2
Query: 27 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 86
W NFD R+ GTWY +A+ D + + A +S+ + G + + EG W
Sbjct: 4558772 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK- 4558596

Query: 87 ADMVGTFTDTEPAREKRYKRW 108
+ FT + + A E+ +G
Sbjct: 4558595 TEKGAVTTPSPRAALKVSYFFG 4558530
    
```

this is a good candidate for a novel gene/protein

```

Alignments
>g1147754841emb1|CAE42626.1| (Y17716) hypothetical protein (Klebsiella oxytoca)
Length = 177
Score = 161 bits (406), Expect = 9e-40
Identities = 75/81 (92%), Positives = 70/81 (95%)
Query: 1 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 60
VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 60
Sbjct: 30 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 89

Query: 61 EKGAVTTPSPRAALKVSYFFG 81
EKGAVTTP P+RAALKVSYFFG
Sbjct: 90 EKGAVTTPSPRAALKVSYFFG 110

>g1124977021|ep1|Q460361|BLC_CITF OUTER MEMBRANE LIPOPROTEIN BLC PRECURSOR
BL13110101|E11114210| outer membrane lipoprotein - Citrobacter freundii
g11743346|ep1|AAC46355.1| (U12727) lipocalin precursor [Citrobacter freundii]
Length = 177
Score = 159 bits (401), Expect = 6e-39
Identities = 74/81 (91%), Positives = 77/81 (94%)
Query: 1 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 60
VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 60
Sbjct: 30 VVNNFDAKRYLGTWYIEIARLDRHFERFEGLEQVATTYSLRDDGGININRGTPNFAREWQK 89

Query: 61 EKGAVTTPSPRAALKVSYFFG 81
EKGAVTTP P+ AALKVSYFFG
Sbjct: 90 EKGAVTTPSPRAALKVSYFFG 110
    
```

A blastp nr search confirms that the *Salmonella* query is closely related to other lipocalins

(Page 150)

## BLAST for gene discovery

You can use BLAST to find a "novel" gene

Ideally, try to find a new gene this week or next. I will provide sample projects from last year.