

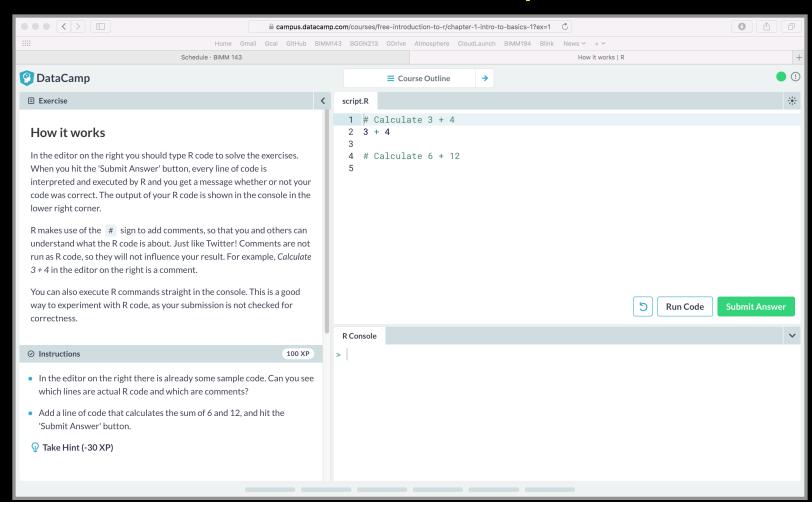
#### **Homework: DataCamp Signup!**

Install R and RStudio (see website)

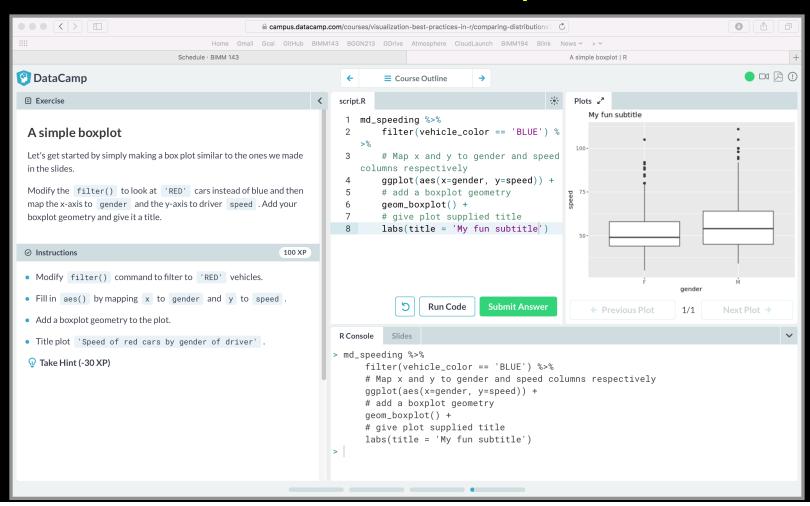
Complete the Introduction to R course on DataCamp (Check Piazza for your DataCamp invite and sign up with your UCSD email (i.e. first part of your email address) please.

Let me know **NOW** if you don't have access to DataCamp!

## DataCamp



## DataCamp



#### **Homework: DataCamp Signup!**

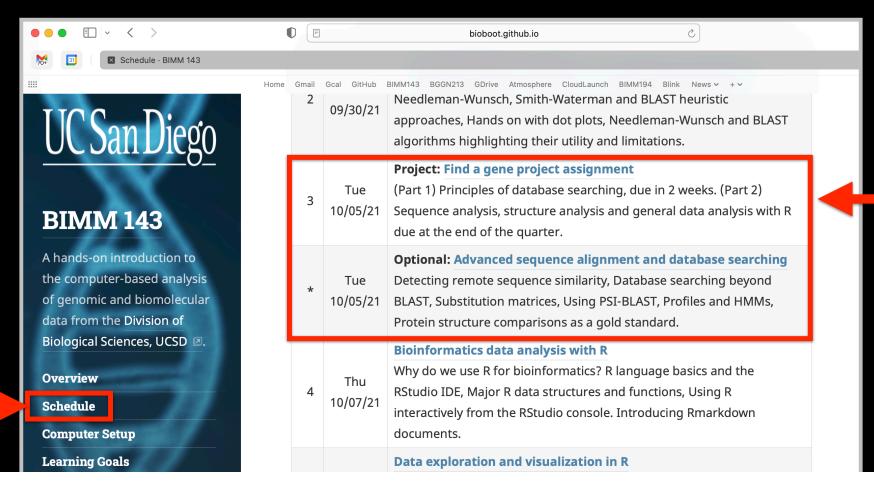
Install R and RStudio on YOUR laptop (see website)

Complete the Introduction to R course on DataCamp (Check Piazza for your DataCamp invite and sign up with your UCSD email (i.e. first part of your email address) please.

Let me know **NOW** if you don't have access to DataCamp!

#### Class 3: Hands-on section

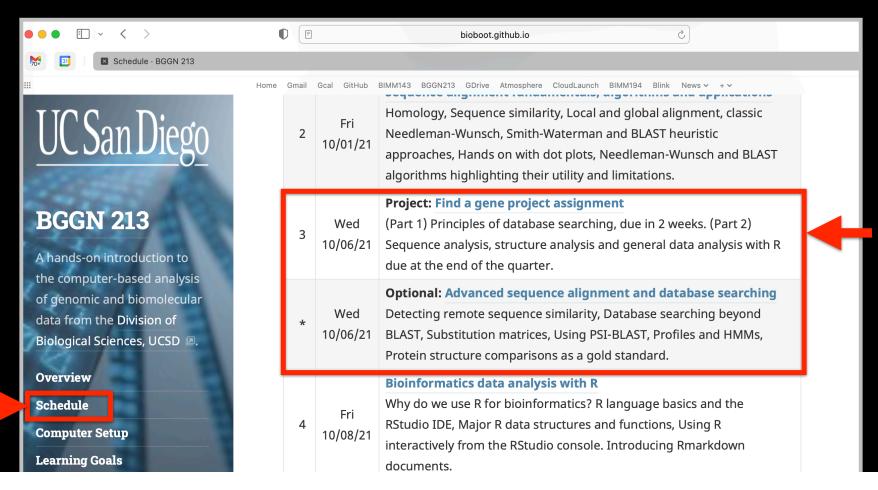
http://thegrantlab.org/bimm143/





#### Class 3: Hands-on section

http://thegrantlab.org/bggn213/



• A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project description) and the <u>example report</u> for format and content guidance.

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project description) and the <u>example report</u> for format and content guidance.
  - → Your responses to questions Q1-Q4 are due 12pm San Diego time on Monday of week 4 (Oct 20th, 10/20/25).
  - → The complete assignment, including responses to all questions, is due 12pm Monday of week 10 (Dec 8th, 12/08/25).

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project description) and the <a href="mailto:example report">example report</a> for format and content guidance.
  - Your responses to questions Q1-Q4 are due 12pm San Diego time on Monday of week 4 (Oct 13th, 10/13/25).
  - → The complete assignment, including responses to all questions, is due 12pm Monday of week 10 (Dec 8th, 12/08/25).

#### Questions:

[Q1] Tell me the name of a protein you are interested in. Include the species and the accession number. This can be a human protein or a protein from any other species as long as it's function is known

If you do not have a favorite protein, select human RBP4 or KIF11. Do not use beta globin as this is in the worked example report that I provide you with online.

[Q2] Perform a BLAST search against a DNA database, such as a database consisting of genomic DNA or ESTs. The BLAST server can be at NCBI or elsewhere. Include details of the BLAST method used, database searched and any limits applied (e.g. Organism).

Also include the output of that BLAST search in your document. If appropriate, change the font to Courier size 10 so that the results are displayed neatly. You can also screen capture a BLAST output (e.g. alt print screen on a PC or on a MAC press \( \frac{2}{3} \)-shift-4. The pointer becomes a bulls eye. Select the area you wish to capture and release. The image is saved as a file called Screen Shot [].png in your Desktop directory). It is **not** necessary to print out all of the blast results if there are many pages.

On the BLAST results, clearly indicate a match that represents a protein sequence, encoded from some DNA sequence, that is homologous to your query protein. I need to be able to inspect the pairwise alignment you have selected, including the E value and score. It should be labeled a "genomic clone" or "mRNA sequence", etc. - but include no functional annotation.

In general, [Q2] is the most difficult for students because it requires you to have a "feel" for how to interpret BLAST results. You need to distinguish between a perfect match to your query (i.e. a sequence that is not "novel"), a near match (something that might be "novel", depending on the results of [Q4]), and a non-homologous result.

If you are having trouble finding a novel gene try restricting your search to an organism that is poorly annotated.

[Q3] Gather information about this "novel" protein. At a minimum, show me the protein sequence of the "novel" protein as displayed in your BLAST results from [Q2] as FASTA format (you can copy and paste the aligned sequence subject lines from your BLAST result page if necessary) or translate your novel DNA sequence using a tool called EMBOSS Transeq at the EBI. Don't forget to translate all six reading frames; the ORF (open reading frame) is likely to be the longest sequence without a stop codon. It may not start with a methionine if you don't have the complete coding region. Make sure the sequence you provide includes a header/subject line and is in traditional FASTA format.

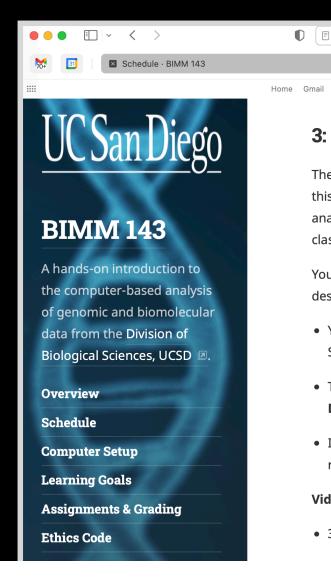
Here, tell me the name of the novel protein, and the species from which it derives. It is very unlikely (but still definitely possible) that you will find a novel gene from an organism such as *S. cerevisiae*, human or mouse, because those genomes have already been thoroughly annotated. It is more likely that you will discover a new gene in a genome that is currently being sequenced, such as bacteria or plants or protozoa.

[Q4] Prove that this gene, and its corresponding protein, are novel. For the purposes of this project, "novel" is defined as follows. Take the protein sequence (your answer to [Q3]), and use it as a query in a blastp search of the nr database at NCBI.

- If there is a match with 100% amino acid identity to a protein in the database, from the same species, then your protein is NOT novel (even if the match is to a protein with a name such as "unknown"). Someone has already found and annotated this sequence, and assigned it an accession number.
- If the top match reported has less than 100% identity, then it is likely that your protein is novel, and you have succeeded.
- If there is a match with 100% identity, but to a different species than the one you started with, then you have likely succeeded in finding a novel gene.
- If there are no database matches to the original query from [Q1], this indicates that you have partially succeeded: yes, you may have found a new gene, but no, it is not actually homologous to the original query. You should probably start over.

[Q5] Generate a multiple sequence alignment with your novel protein, your original query protein, and a group of other members of this family from different species. A typical number of proteins to use in a multiple sequence alignment for this assignment purpose is a minimum of 5 and a maximum of 20 - although the exact number is up to you. Include the multiple sequence alignment in your report. Use Courier font with a size appropriate to fit page width.

Side-note: Indicate your sequence in the alignment by choosing an appropriate name for each sequence in the input unaligned sequence file (i.e. edit the sequence file so that the species, or short common, names (rather than accession numbers) display in the output alignment and in the subsequent answers below). The goal in this step is to create an interesting an alignment for building a phylogenetic tree that illustrates species divergence.



#### 3: (Project) Find a Gene Assignment Part 1

bioboot.github.io

The **find-a-gene project** is a required assignment for BIMM-143. The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.

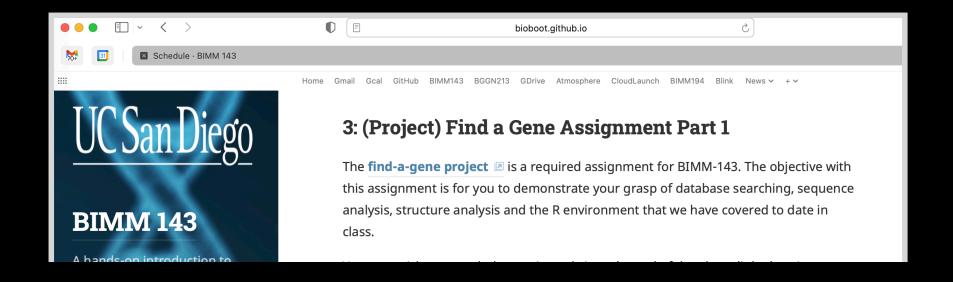
Ç

You may wish to consult the scoring rubric at the end of the above linked project description and the **example report** for format and content guidance.

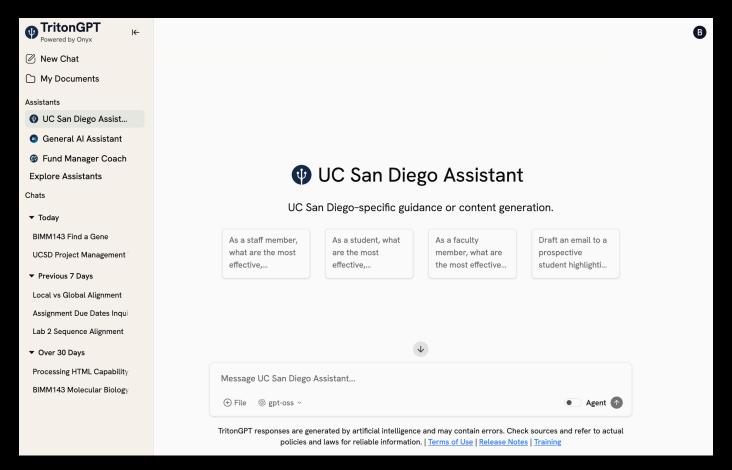
- Your responses to questions Q1-Q4 are due **Tuesday Oct 19th** (10/19/21) at 12pm San Diego time.
- The complete assignment, including responses to all questions, is due Thursday
   Dec 2nd (12/02/21) at 12pm San Diego time.
- In both instances your PDF format report should be submitted to GradeScope. Late responses will not be accepted under any circumstances.

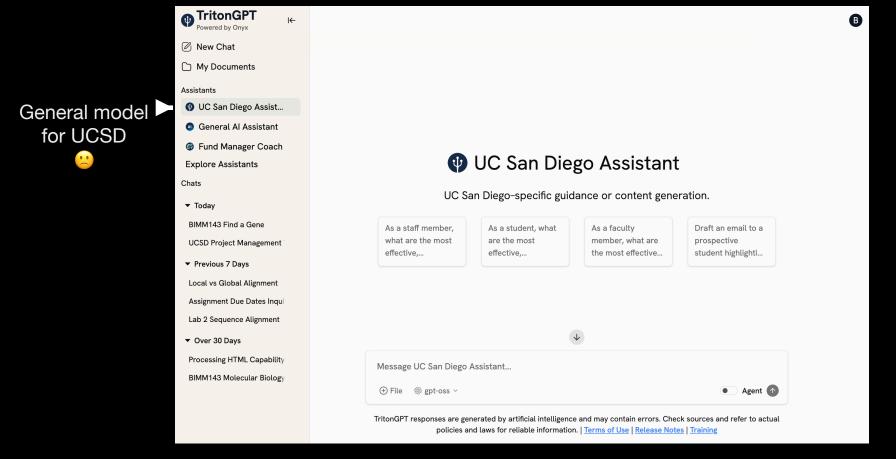
#### **Videos:**

• 3.1 - Project introduction 🗷 Please note: due dates may differ from those in video.



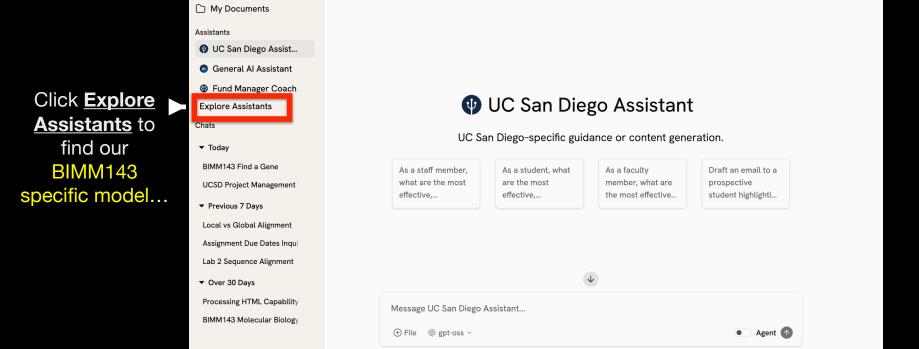
- Your responses to questions Q1-Q4 are due 12pm San Diego time on Monday of week 4 (Oct 20th, 10/20/25).
- The complete assignment, including responses to all questions, is due 12pm Monday of week 10 (Dec 8th, 12/08/25).





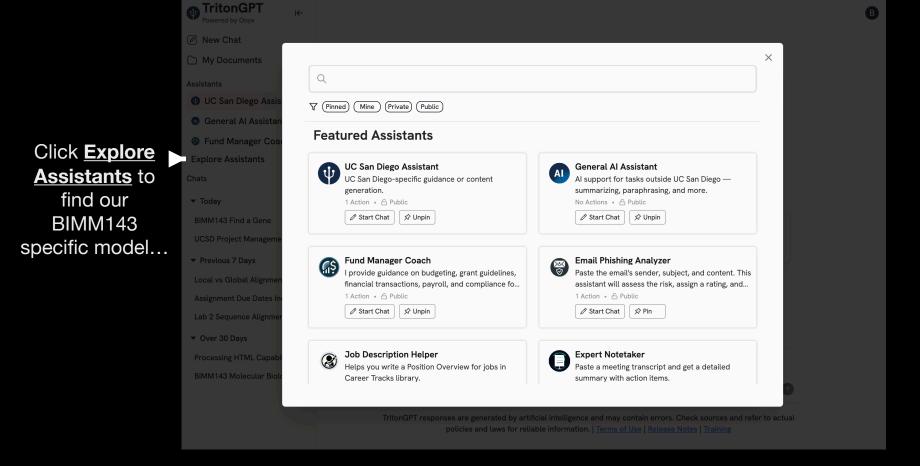
TritonGPT responses are generated by artificial intelligence and may contain errors. Check sources and refer to actual policies and laws for reliable information. | <u>Terms of Use</u> | <u>Release Notes</u> | <u>Training</u>

B

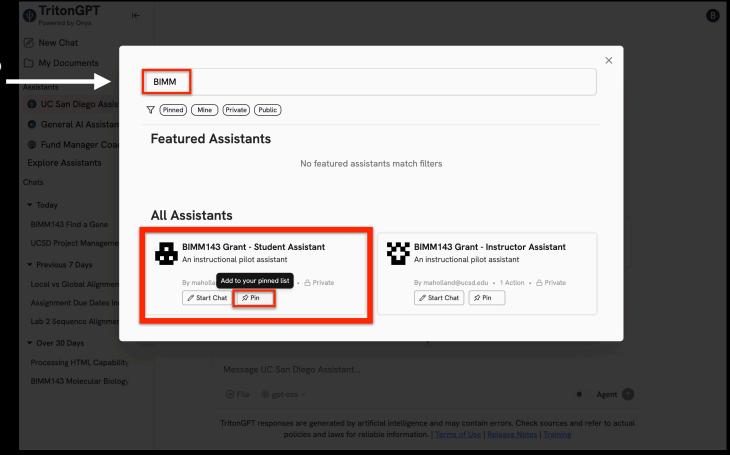


TritonGPT
Powered by Onyx

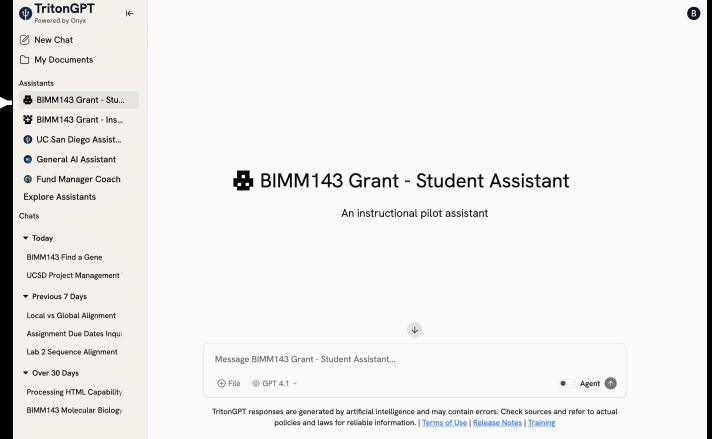
New Chat



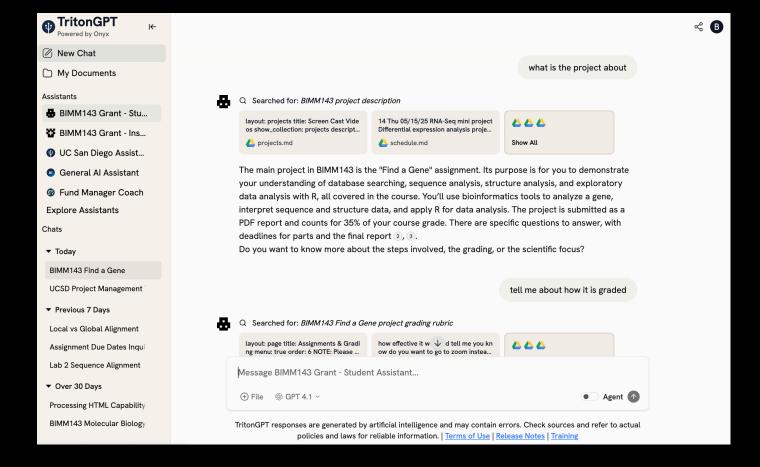
Filter or scroll to find "BIMM143 Grant - Student Assistant"



Pin and drag to ► top of list...

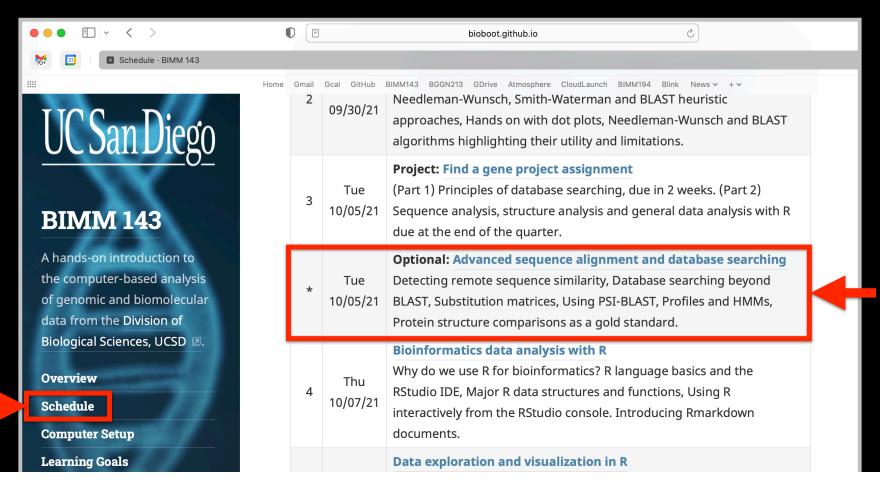


Ask away...



#### Class 3: Hands-on section

http://thegrantlab.org/bimm143/



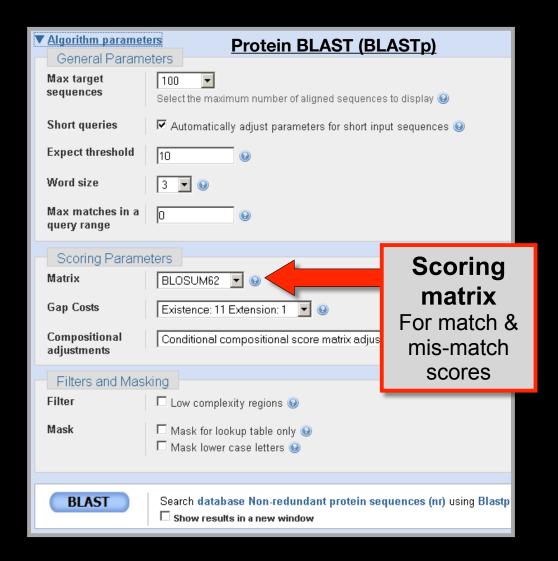
Sequence 1 Sequence 2		GATTAC									
		GTCGACGC				G T C G A C G C					
Match Score Mismatch Score  1							G T C G A C G C G A T T A C Score = -4				
		G	Т	C	G	A	C	G	С		
	0	-2	-4	-6	-8	-10	-12	-14	-16		
G	-2	1	<b>←</b> -1	<b>4</b> -3	-5	-6 + 1		o a match G) = -5		re from Upper cell -2 (The Gap score) =	
A	-4	<b>↑</b> -1	0	<b>*</b> <b>+</b> -2	<b>×</b> <b>←</b> -4		e from S 2 (The C	ide cell Sap score		ning (max) score is -5	
T	-6	-3	0	-1	<b>×</b> <b>↓</b> -3	+ -5	-5	<b>+</b> -7	<b>←</b> -9		
Т	-8	-5	<b>K ↑</b> -2	-1	-2	<b>×</b> <b>←</b> -4	<b>X</b> <b>4</b> -6	<b>K</b> -6	<b>K</b> <b>←</b> -8		
A	-10	<b>↑</b> -7	-4	<b>× →</b> -3	-2	× -1	<b>+</b> -3	<b>+</b> -5	<b>*</b>		
C	-12	<b>+</b> -9	<b>+</b> -6	<b>K</b> -3	<b>× ↑</b> -4	<b>× ↑</b> -3	0	<b>+</b> -2	<b>K</b>		

Barry J Grant.

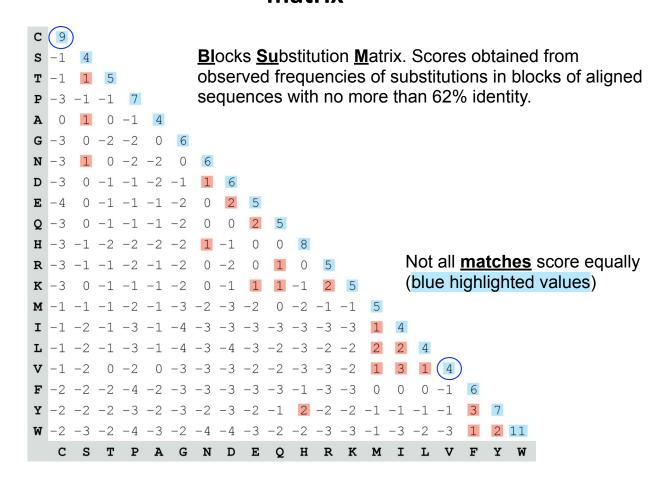
NW App Link

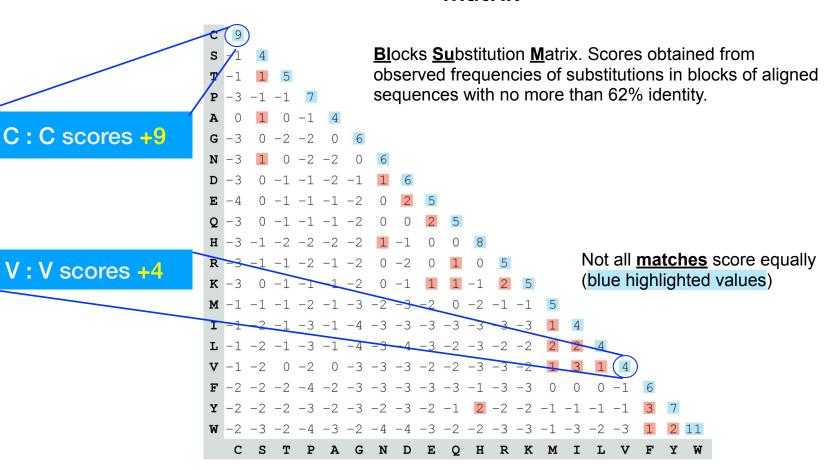
#### **Key Question:**

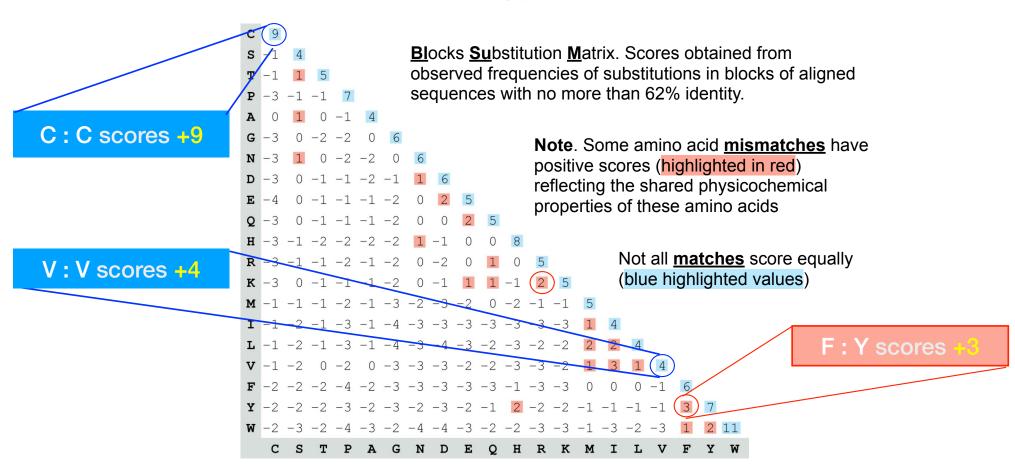
Q. Where do our alignment match and mis-match scores typically come from?



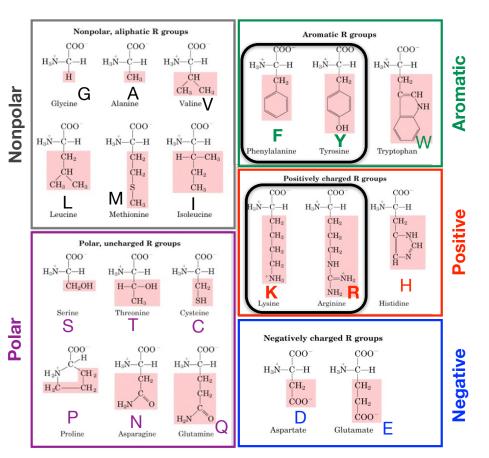
```
Blocks Substitution Matrix. Scores obtained from
s - 1 4
                   observed frequencies of substitutions in blocks of aligned
                   sequences with no more than 62% identity.
A 0 1 0 -1 4
G -3 0 -2 -2 0 6
Y -2 -2 -2 -3 -2 -3 -2 -1 2 -2 -2 -1 -1 -1 -1 3 7
```



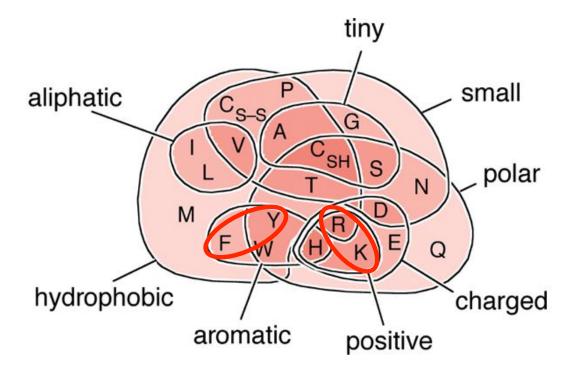




# Protein scoring matrices reflect the properties of amino acids



# Protein scoring matrices reflect the properties of amino acids



**Key Trend**: High scores for amino acids in the same "biochemical group" and low scores for amino acids from different groups.

N.B. BLOUSM62 does not take the local context of a particular position into account

(i.e. all like substitutions are scored the same regardless of their location in the molecules).

We will revisit this later...

#### YOUR TURN!

 There are four required and one optional hands-on sections including:

1.	Limits of using BLAST	[~10 mins]
2.	Using PSI-BLAST	[~30 mins]
3.	Examining conservation patterns  — BREAK [15 mins]—	[~20 mins]
1		[. 10 minol
4.	[Optional] Using HMMER	[~10 mins]
5.	Divergence of protein sequence and structure	[~25 mins]

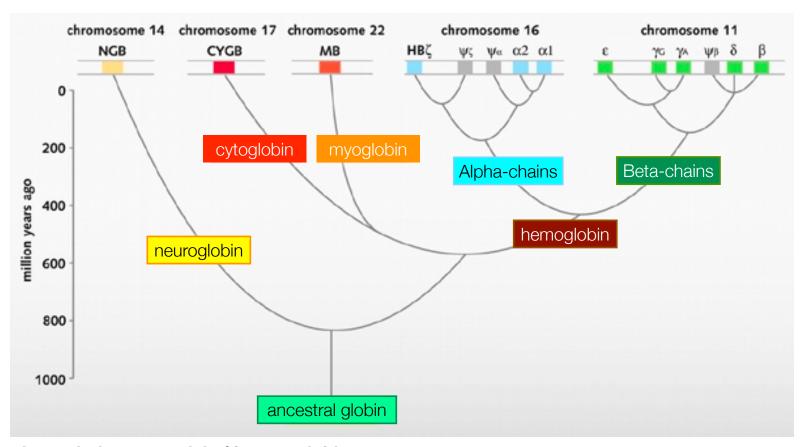
- ▶ Please do answer the last review question (Q20).
- ▶ We encourage <u>discussion</u> at your **Table** and on **Piazza**!

#### YOUR TURN!

 There are four required and one optional hands-on sections including:

1.	Limits of using BLAST	[~10 mins]
2.	Using PSI-BLAST	[~30 mins]
3.	Examining conservation patterns	[~20 mins]
	— BREAK [15 mins]—	
4.	[Optional] Using HMMER	[~10 mins]
5.	Divergence of protein sequence and structure	[~25 mins]

- ▶ Please do answer the last review question (Q20).
- ▶ We encourage <u>discussion</u> at your **Table** and on **Piazza**!



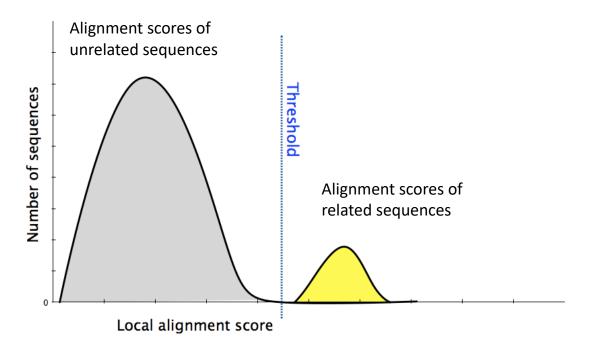
#### An evolutionary model of human globins.

The different locations of globin genes in human chromosomes are reported at the top of the figure, distinguishing between the functional genes (in color) and the pseudogenes (in grey).

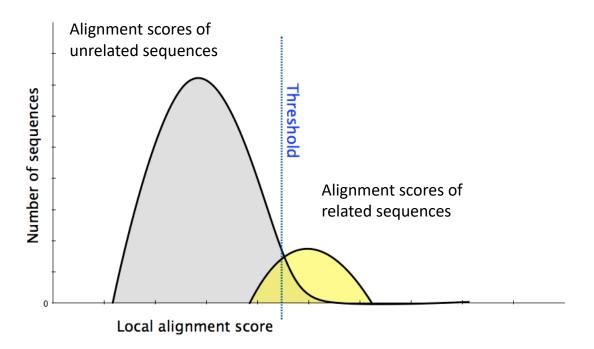
# **Question:**

Q. Can we find and align these homologous globins using SW approaches such as BLAST?

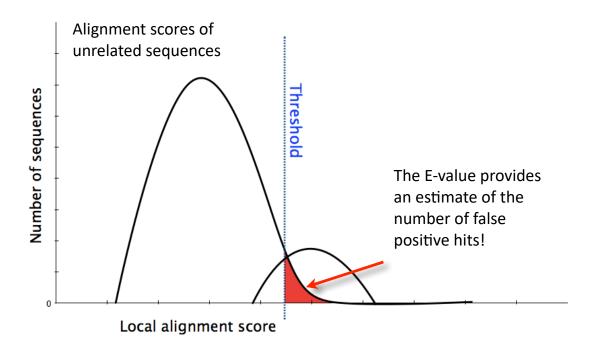
• Ideally, a threshold separates all query related sequences (yellow) from all unrelated sequences (gray)



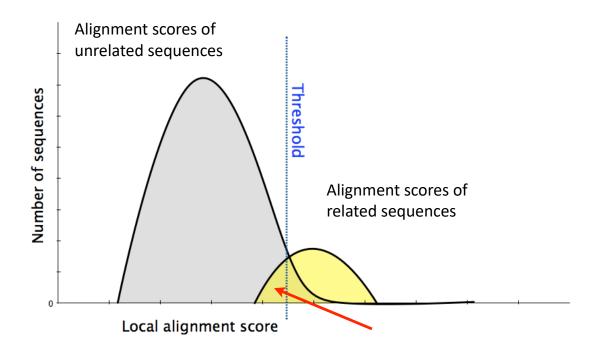
- Unfortunately, often both score distributions overlap
  - The E value describes the expected number of hits with a score above the threshold if the query and database are unrelated



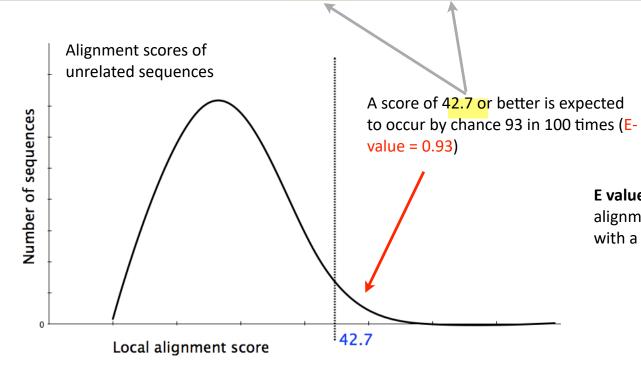
- Unfortunately, often both score distributions overlap
  - The E value describes the expected number of hits with a score above the threshold if the query and database are unrelated



- Maybe myoglobin, cytoglobin, neuroglobin etc. are found but not reported because of our E-value cutoff?
  - Lets change the cutoff and see...



Description	Max score	Query cover	E value	Max ident	Accession
hemoglobin subunit beta	284	100%	0	100%	NP_000510.1
hemoglobin subunit delta	240	100%	0	75.5%	NP_005321.1
hemoglobin subunit alpha	114	97%	0	43.45%	NP_000508.1
probable ATP-dependent RNA helicase	42.7	10%	0.93	32%	XP_011530405.1



**E value**: The number of alignments expected by chance with a particular score

# YOUR TURN!

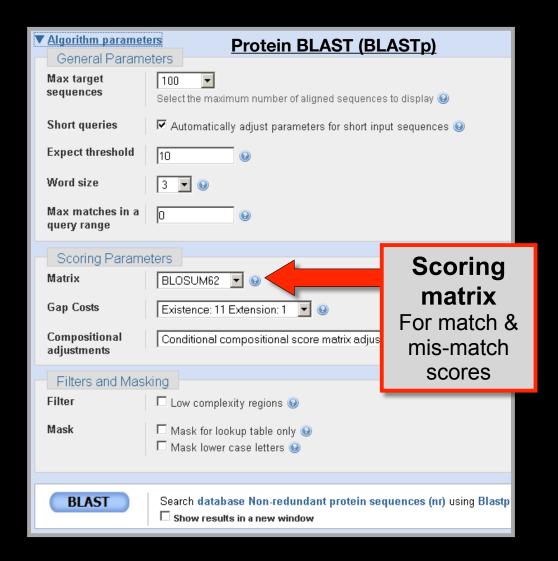
 There are four required and one optional hands-on sections including:

1.	Limits of using BLAST	[~10 mins]
2.	Using PSI-BLAST	[~30 mins]
3.	Examining conservation patterns  — BREAK [15 mins]—	[~20 mins]
4.	[Optional] Using HMMER	[~10 mins]
5.	Divergence of protein sequence and structure	[~25 mins]

- ▶ Please do answer the last review question (Q20).
- ▶ We encourage <u>discussion</u> at your **Table** and on **Piazza**!

Recall: BLOUSM62 does not take the local context of a particular position into account

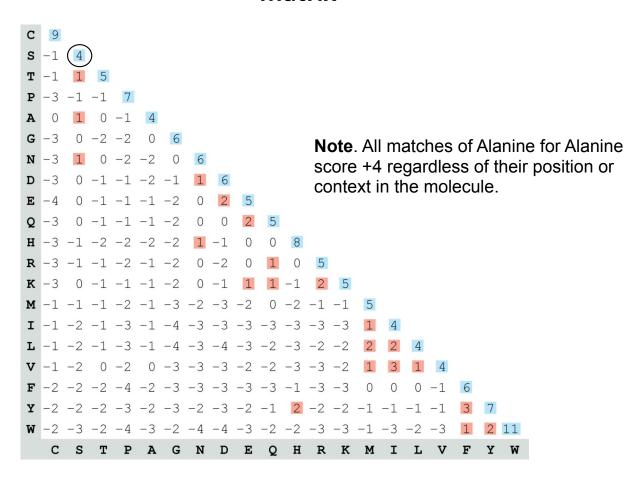
(i.e. all like substitutions are scored the same regardless of their location in the molecules).



# By default BLASTp match scores come from the BLOSUM62 matrix

```
Blocks Substitution Matrix. Scores obtained from
s - 1 4
                   observed frequencies of substitutions in blocks of aligned
                   sequences with no more than 62% identity.
A 0 1 0 -1 4
G -3 0 -2 -2 0 6
Y -2 -2 -2 -3 -2 -3 -2 -1 2 -2 -2 -1 -1 -1 -1 3 7
```

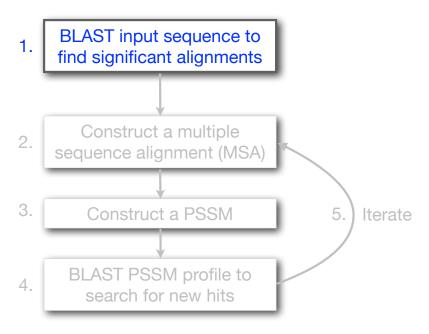
# By default BLASTp match scores come from the BLOSUM62 matrix



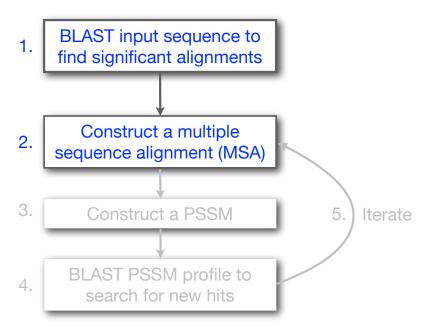
 The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a <u>scoring matrix that is</u> <u>customized to your query</u>

- The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a <u>scoring matrix that is</u> <u>customized to your query</u>
  - PSI-BLAST constructs a multiple sequence alignment from the results of a first round BLAST search and then creates a "<u>profile</u>" or specialized **position-specific** scoring matrix (<u>PSSM</u>) for subsequent search rounds

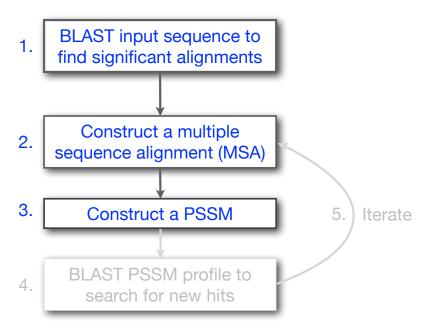
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST

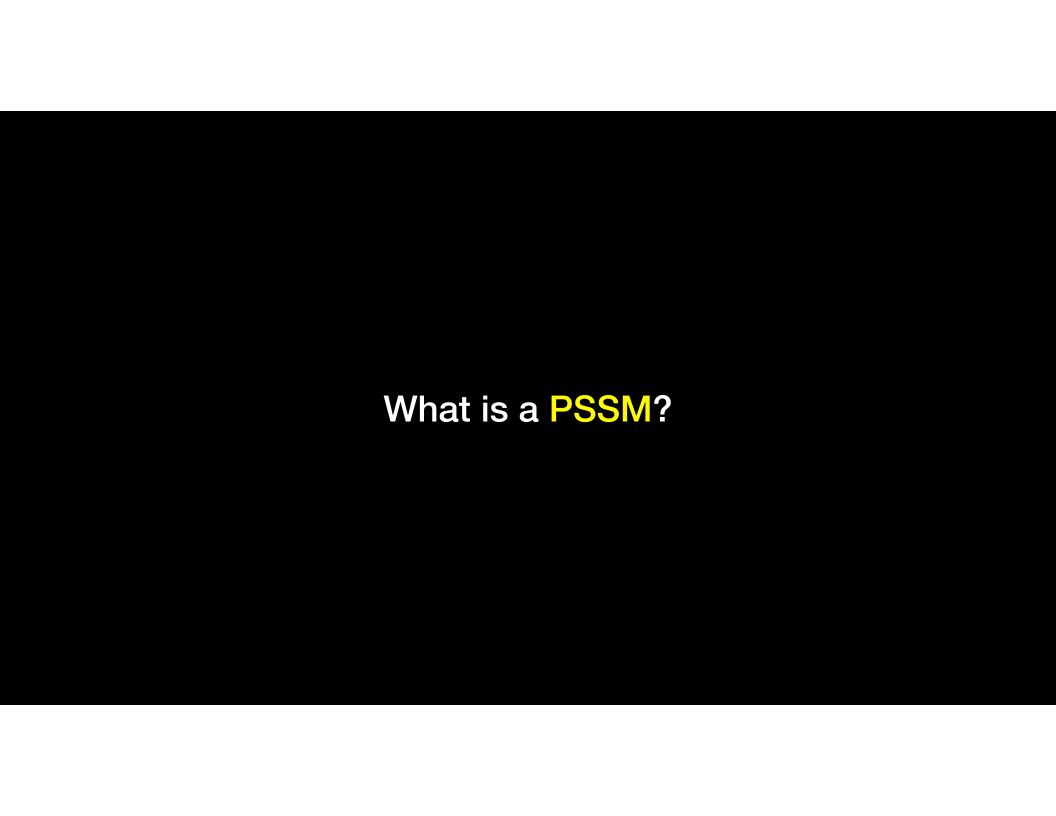


Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



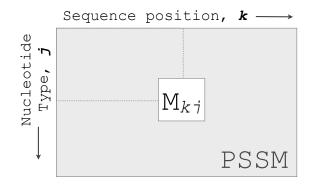


# What are PSSM sequence profiles?

A sequence profile is a **position-specific scoring matrix** (or **PSSM**, often pronounced 'possum') that gives a *quantitative* description of a set of aligned sequences.

PSSMs assign a score to a query sequence and are widely used for database searching.

A simple PSSM has as many columns as there are positions in the alignment, and either 4 rows (one for each DNA nucleotide) or 20 rows (one for each amino acid).

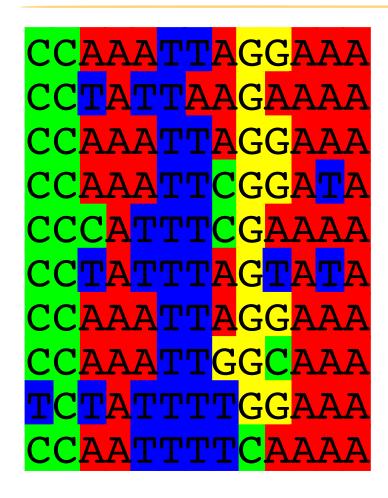


$$M_{kj} = \log\left(\frac{p_{kj}}{p_j}\right)$$

 $\mathbf{M}_{kj}$  score for the jth nucleotide at position k

 $\mathbf{p}_{kj}$  probability of nucleotide j at position k

 $\mathbf{p}_{j}$  "background" probability of nucleotide j



Here we have **10 aligned** transcription factor binding site nucleotide sequences

That span **13 positions** (i.e. columns of nucleotides).

We will build a 13 x 4 **PSSM** (k=13, j=4).

CCAAATTAGGAAA
CCTATTAAGAAAA
CCAAATTCGGATA
CCCATTTCGAAAA
CCTATTTAGTATA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTGGCAAA
TCTATTTTGGAAA
CCAA

# First we will build an alignment Counts matrix Position k = 1 2 3 4 5 6 7 8 9 10 11 12 13 A: C: G: T:

CC	AAA	TTA	GG	AAA
CC	TAT	'T <mark>A</mark> A	<mark>G</mark> A	AAA
CC	AAA	TT <mark>A</mark>	GG	A <mark>A</mark> A
CC	AAA	TTC	GG	A <mark>T</mark> A
CC	C <mark>A</mark> T	'TT <mark>C</mark>	G <mark>A</mark>	AAA
CC	TAT	TT.	G <mark>T</mark>	A <mark>T</mark> A
CC	AAA	TTA	GG	AAA
CC.	<u>A</u> AA	TTG	GC	AAA
TC	TAT	TTT	GG	AAA
CC	AA <mark>T</mark>	TTT	CA	AAA

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:													
C:													
G:													
T:													

CCAAATTAGGAAA
CCTATTAAGAAAA
CCAAATTAGGAAA
CCAAATTCGGATA
CCCATTTCGAAAA
CCTATTTAGTATA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTGGCAAA
TCTATTTTGGAAA

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0												
C:	9												
G:	0												
T:	1												

CC	AAA <mark>TT</mark> A <mark>GG</mark> A	AA
CC	<mark>T</mark> A <mark>TT</mark> AA <mark>G</mark> AA	AA
	AAA <mark>TT</mark> AGGA	
CC	AAA <mark>TT</mark> CGGA <mark>!</mark>	<b>C</b> A
CC	<mark>CATTT</mark> CGAA <i>I</i>	AA
CC	<mark>TATTTAGTA</mark> :	<b>C</b> A
	AAA <mark>TT</mark> AGGA	
	<mark>AAA<mark>TT</mark>GGC</mark> A <i>I</i>	
	<mark>TA</mark> TTTT <mark>GG</mark> A	
CC	AA <mark>TTTT</mark> CAA <i>I</i>	AΑ

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0												
C:	9												
G:	0												
T:	1												
Consensus	С												

CCAAATTAGGAAA
CCTATTAAGAAAA
CCAAATTAGGAAA
CCAAATTCGGATA
CCCATTTCGAAAA
CCTATTTAGTATA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTGGCAAA
CCAAATTGGCAAA
CCAAATTTGGAAA
CCAATTTTCAAAA

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0											
C:	9	10											
G:	0	0											
T:	1	0											
Consensus	С	С											

CCAAATTAGGAAA
CCTATTAAGAAAA
CCAAATTAGGAAA
CCAAATTCGGATA
CCCATTTCGAAAA
CCTATTTAGTATA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTGGCAAA
CCAAATTTAGCAAA
CCAAATTTAGAAA
CCAAATTTT

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6										
C:	9	10	1										
G:	0	0	0										
T:	1	0	3										
Consensus	С	С	Α										

CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTCGGATA
CCCATTTCGAAAA
CCTATTTAGTATA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTAGGAAA
CCAAATTGGCAAA
TCTATTTTGGAAA
CCAAATTT

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0
Consensus	С	С	Α	Α	[AT]	Т	Т	Α	G	G	Α	Α	Α

CCAAATTAGGAAA
CCTATTAAGAAAA
CCAAATTCGGATA
CCCATTTCGAAAA
CCTATTTAGTATA
CCAAATTGGAAA
CCTATTTAGGAAA
CCAAATTGGAAA
CCAAATTGGCAAA
CCAAATTGGCAAA
CCAAATTTGGAAA
CCAATTTTCAAAA

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0
Consensus	С	С	Α	Α	[AT]	Т	Т	Α	G	G	Α	Α	Α

Average Profile (Frequency) matrix:

Often we will not
communicate with
the count matrix
but rather the
derived average
profile (a.k.a.
frequency matrix).
, , , , , , , , , , , , , , , , , , ,

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
<b>A</b> :	0	0	0.6	1	0.5	0	0.1	0.5	0	0.3	1	0.8	1
C:	0.9	1	0.1	0	0	0	0	0.2	0.1	0.1	0	0	(
G:	0	0	0	0	0	0	0	0.1	0.9	0.5	0	0	(
T:	0.1	0	0.3	0	0.5	1	0.9	0.2	0	0.1	0	0.2	(
Consensus	С	С	Α	Α	[AT]	Т	Т	Α	G	G	Α	Α	A

CC	AAA <mark>TT</mark> A <mark>GG</mark> AAA
CC	<mark>T</mark> A <mark>TT</mark> AA <mark>GA</mark> AAA
CC	AAA <mark>TT</mark> A <mark>GG</mark> AAA
	<mark>AAA</mark> TT <mark>CGGA</mark> TA
	<mark>CA</mark> TTT <mark>CG</mark> AAAA
	<mark>TATTTAGTATA</mark>
	AAA <mark>TT</mark> AGGAAA
	<mark>AAA<mark>TT</mark>GGC</mark> AAA
	<mark>TA</mark> TTTT <mark>GG</mark> AAA
CC	AA <mark>TTTT</mark> CAAAA

#### **Alignment Counts matrix:**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0
Consensus	С	С	Α	Α	[AT]	Т	Т	Α	G	G	Α	Α	Α

#### Or the "score $(M_{kj})$ matrix" = PS**S**M

- $C_{kj}$  Number of jth type nucleotide at position k
- **Z** Total number of aligned sequences
- **p**<sub>j</sub> "background" probability of nucleotide *j*
- $\mathbf{p}_{kj}$  probability of nucleotide j at position k

$$M_{kj} = \log\left(\frac{p_{kj}}{p_j}\right) \quad p_{kj} = \frac{C_{kj} + p_j}{Z + 1}$$

$$M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right)$$

Alignment Matrix: Cki 
 1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12

 0
 0
 6
 10
 5
 0
 1
 5
 0
 3
 10
 8

 9
 10
 1
 0
 0
 0
 0
 2
 1
 1
 0
 0

 0
 0
 0
 0
 0
 0
 1
 9
 5
 0
 0

 1
 0
 3
 0
 5
 10
 9
 2
 0
 1
 0
 2
 Position k = 13 A: 10 C: G: **k=1**, **j=A**:  $M_{kj} = \log \left( \frac{C_{kj} + p_j / Z + 1}{p_j} \right) = \log \left( \frac{0 + 0.25 / 10 + 1}{0.25} \right) = -2.4$ **k=1**, **j=C**:  $M_{kj} = \log \left( \frac{C_{kj} + p_j / Z + 1}{p_j} \right) = \log \left( \frac{9 + 0.25 / 10 + 1}{0.25} \right) = 1.2$ **k=1**, **j=T**:  $M_{kj} = \log \left( \frac{C_{kj} + p_j / Z + 1}{p_j} \right) = \log \left( \frac{1 + 0.25 / 10 + 1}{0.25} \right) = -0.8$ PSSM: Mki 
 1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12

 -2.4
 -2.4
 -0.8
 0.6
 -2.4
 0.2
 1.3
 1.1
 Position k = 13 A: 1.3 1.2 1.3 -0.8 -2.4 -2.4 -2.4 -2.4 -0.2 -0.8 -0.8 -2.4 -2.4 C: -2.4 -2.4 -2.4 -2.4 -2.4 -2.4 -2.4 -0.8 1.2 0.6 -2.4 -2.4 G: -2.4 0.2 -2.4 0.6 1.3 1.2 -0.2 -2.4 -0.8 -2.4 -0.8 -2.4 T: -0.2 -2.4

# Scoring a test sequence

#### Query Sequence

#### **CCTATTTAGGATA**

#### PSSM:

## Scoring a test sequence

Query Sequence

#### **CCTATTTAGGATA**

#### PSSM:

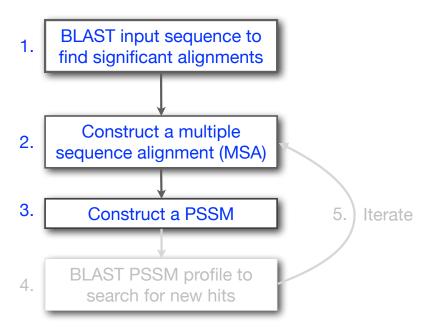
Q. Does the query sequence match the DNA sequence profile?

## Scoring a test sequence...

Query Sequence Best Possible Sequence PSSM: Position k = 1 2 3 5 6 7 10 11 12 13 1.3 0.6 -2.4 0.8 -2.4 -0.8 0.6 -2.4 0.2 1.3 1.1 1.3 A: -2.4 1.3 -0.8 -2.4 -2.4 -2.4 -2.4 C: 1.2 -0.2 -0.8 -0.8 -2.4 -2.4 -2.4 -2.4 G: -2.4 -2.4 -2.4 -2.4 -2.4 -2.4 -0.8 1.2 0.6 -2.4 -2.4 -2.4 T: -0.8 -2.4 0.2 -2.4 0.6 1.3 1.2 -0.2 -2.4 -2.4 -0.2 -2.4 -0.8 Max Score: C C Т Α G G Α Α

**A.** Following method in Harbison *et al.* (2004) Nature 431:99-104 Heuristic threshold for match = 60% x Max Score =  $(0.6 \times 13.8 = 8.28)$ ; 11.9 > 8.28; Therefore our query is a potential TFBS!

Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



# Inspect the blastp output to identify empirical "rules" regarding amino acids tolerated at each position

```
730496
              FTVDENGOMSATAKGRVRLFNNWDVCADMIGSFTDTEDPAKFKMKYWGVASFLQKGNDDH 125
200679
         63
              FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTFTDTEDPAKFKMKYWGVASFLORGNDDH 122
206589
         34
              FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTFTDTEDPAKFKMKYWGVASFLORGNDDH 93
2136812
         2
                      MSATAKGRVRLLNNWDVCADMVGTFTDTEDPAKFKMKYWGVASFLQKGNDDH 53
132408
         65
              FKIEDNGKTTATAKGRVRILDKLELCANMVGTFIETNDPAKYRMKYHGALAILERGLDDH 124
267584
         44
              FSVDESGKVTATAHGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAASYLQTGNDDH 103
267585
         44
              FSVDGSGKVTATAQGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAAAYLQSGNDDH 103
8777608
         63
              FTIHEDGAMTATAKGRVIILNNWEMCADMMATFETTPDPAKFRMRYWGAASYLQTGNDDH 122
6687453
         60
              FKVEEDGTMTATAIGRVIILNNWEMCANMFGTFEDTEDPAKFKMKYWGAAAYLOTGYDDH 119
10697027 81
              FKVQEDGTMTATATGRVIILMNWEMCANMFGTFEDTEEPARFKMKYWGAAAYLQTGYDDH 140
13645517 1
                                           MVGTFTDTEDPAKFKMKYWGVASFLQKGNDDH 32
13925316 38
              FSVDGSGKMTATAOGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAAAYLOSGNDDH 97
131649
         65
              YTVEEDGTMTASSKGRVKLFGFWVICADMAAQYTDPTTPAKMYMTYQGLASYLSSGGDNY 126
```

N,M,L,Y,G

#### GLB2\_LUMTE/31-141 . . EELDHLQVQHE . . . GRKIPD AQLEHLRQQHI...KLGITG GLB2 TYLHE/32-143 GLB3 LAMSP/30-141 TQLAHLASQHS... GLB\_TUBTU/29-139 AQLAHLKSQHA SHLGHLADQHI. . QRKGVTK GLB4 LUMTE/36-146 SLIDHLAEQHK. . ARAGFKT GLB4\_TYLHE/33-143 GLB3\_TYLHE/33-143 EELKHLARQHR. . ERSGVKA GLB1\_TYLHE/30-137 QALAHYAAFHK...QFGTIP MYG CYPCA/23-136 MYG ALLMI/27-143 EVLKPLAKSHA MYG GALGA/23-138 MYG HETPO/23-138 TNVKELADTHI HBAM\_LITCT/26-129 . . . KYQDLHT HBA3\_PLEWA/27-137 . QALSKLSDLHA HBAZ\_CAPHI/27-137 SALSKLSELHA HBB\_HETPO/25-136 PHFVELSKKHY HBB\_SQUAC/25-137 HBB2\_XENLA/26-142 SSLQQLSKIHA HBB1 CYGMA/26-142 DAYAELSTLHS HBB\_LEPPA/25-142 GHLANLSHLHS HBB ALLMI/25-141 HBB0 MOUSE/26-142 GAFASLSZLHC BKLHVBP HBBN\_AMMLE/20-136 HBB\_LITCT/19-135 AYYAKLSERHS . . GELHVDP HBB1\_XENBO/25-141 HBA\_LEPPA/26-138 HHLNKLAEKHG. . KGLLVDP HBA1 TORMA/26-136 HBA\_HETPO/33-143 THLHKLATFHG. . SELKVDP HBA SQUAC/26-136 GHLDPLAVLHG. . TTLCVDP HBAD ERYML/26-136 GTLSQLSDLHA..YNLRVDP [HPKSAD] Н

. . EELDHLQVQHE . . . GRKIPD GLB2 LUMTE/31-141 . . AQLEHLRQQHI. . . KLGITG GLB2 TYLHE/32-143 GLB3 LAMSP/30-141 . . T Q L A H L A S Q H S . . . S R G V S A GLB\_TUBTU/29-139 AQLAHLKSQHA...ERNIKA SHLGHLADQHI..QRKGVTK GLB4 LUMTE/36-146 SLIDHLAEQHK..ARAGFKT GLB4 TYLHE/33-143 EELKHLARQHR. . ERSGVKA GLB3 TYLHE/33-143 GLB1\_TYLHE/30-137 QALAHYAAFHK... MYG\_CYPCA/23-136 . . A I L K P L A T T H A . . N T H K I A L EVLKPLAKSHA..LEHKIPV MYG ALLMI/27-143 MYG\_GALGA/23-138 QPVKALAATHI..TTHKIPP MYG HETPO/23-138 TNVKELADTHI..NKHKIPP KYQDLHT..NKLKLSS HBAM\_LITCT/26-129 HBA3\_PLEWA/27-137 QALSKLSDLHA..YNLRVDP HBAZ CAPHI/27-137 SALSKLSELHA..YVLRVDP SQFTDLSKKHA..EELHVDV HBB HETPO/25-136 PHFVELSKKHY. . EELHVDP HBB SQUAC/25-137 SSLQQLSKIHA..TE HBB2 XENLA/26-142 DAYAELSTLHS HBB1 CYGMA/26-142 HBB LEPPA/25-142 GHLANLSHLHS. . EKLH HBB ALLMI/25-141 GHFANLSKLHC . . EKFHVDP HBB0\_MOUSE/26-142 ETFAHLSELHC . . DKLHADP GAFASLSZLHC HBBN AMMLE/20-136 AYYAKLSERHS. . GELHVDP HBB\_LITCT/19-135 HBB1\_XENBO/25-141 GYYAQLSKYHS..ETLHVDP SCLHTLSEKHA. . RELMVDP HBA LEPPA/26-138 HBA1\_TORMA/26-136 HHLNKLAEKHG. . KGLLVDP THLHKLATFHG. . SELKVDP HBA\_HETPO/33-143 GHLDPLAVLHG..TTLCVDP HBA\_SQUAC/26-136 HBAD ERYML/26-136 . . GTLSQLSDLHA . . YNLRVDP [HPKSAD] Н

## 20 amino acids

All the 7 L 8 L amino 9 L 10 L acids from 11 A position 1 12 A 13 W to the end 14 A of your 15 A 16 A query seq. 37 S 38 H 39 T 40 W

1 M 2 K

41 H 42 A

#### 20 amino acids

All the amino acids from position 1 to the end of your query seq.

1 M 2 K

7 L

8 L

9 L

10 L

11 A

12 A

13 W

14 A

15 A

16 A

37 S

40 W

42 A

```
H
```

#### **Key Point:**

PSSM "scores" differ for the same aminoacid type depending on where it is in your protein sequence

> - i.e. sores are POSITION DEPENDENT!

#### 20 amino acids

A R N D C Q E G H I L K M F P S T W Y V
-1 -2 -2 -3 -2 -1 -2 -3 -2 1 2 -2 6 0 -3 -2 -1 -2 -1 1
-1 1 0 1 -4 2 4 -2 0 -3 -3 3 -2 -4 -1 0 -1 -3 -2 -3
-3 -3 -4 -5 -3 -2 -3 -3 -3 -2 -3 -2 1 -4 -3 -3 12 2 -3
-3 -3 -4 -5 -3 -2 -3 -3 -3 -3 -2 -3 -2 1 -4 -3 -3 12 2 -3

The PSI-BLAST **PSSM** is essentially a query customized scoring matrix that is more sensitive than BLOSUM.

#### **Key Point:**

PSSM "scores" differ for the same aminoacid type depending on where it is in your protein sequence

- i.e. sores are POSITION DEPENDENT!

amino
acids from
position 1
to the end
of your
query seq.

All the

1 M 2 K

6 A

9 L

10 L

11 A

12 A

13 W

14 A

15 A

16 A

37 S

40 W

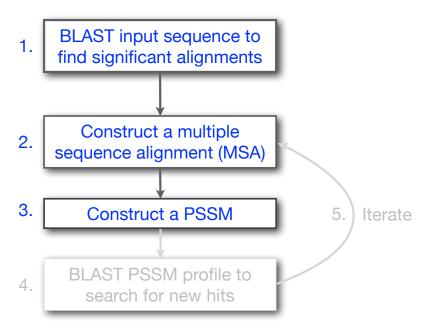
42 A

7 L

8 L

## PSI-BLAST: Position-Specific Iterated BLAST

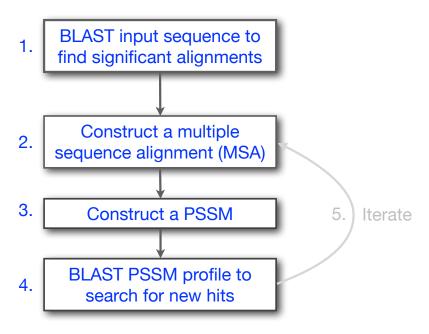
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



(see Altschul et al., Nuc. Acids Res. (1997) 25:3389-3402)

## PSI-BLAST: Position-Specific Iterated BLAST

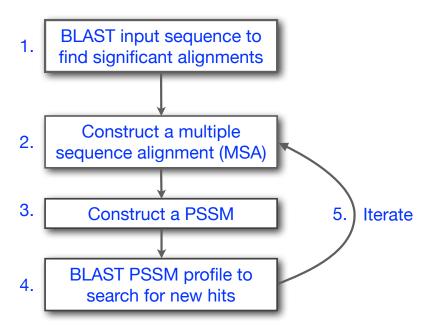
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



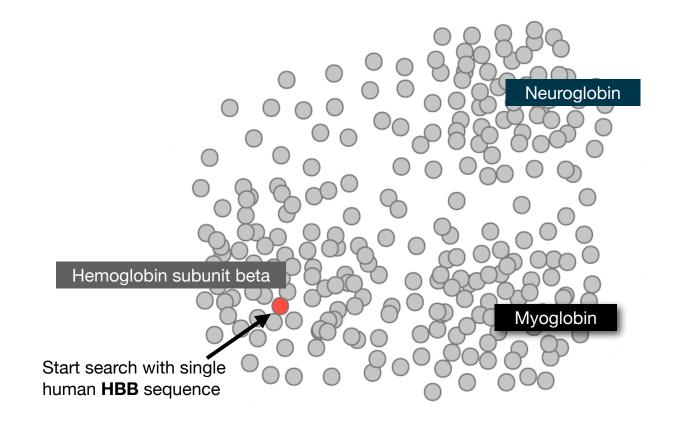
(see Altschul et al., Nuc. Acids Res. (1997) 25:3389-3402)

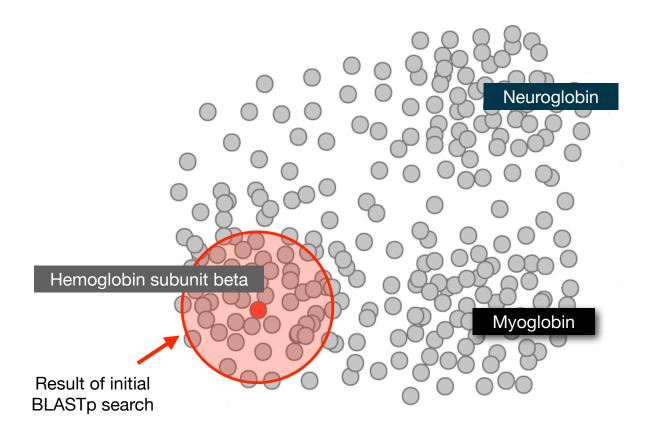
## PSI-BLAST: Position-Specific Iterated BLAST

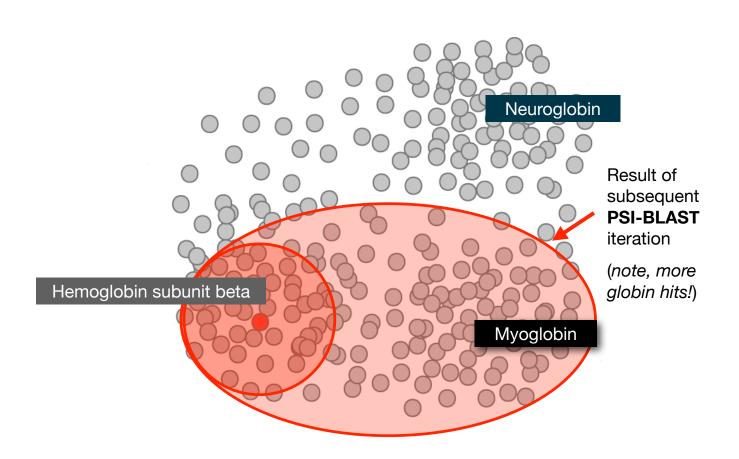
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST

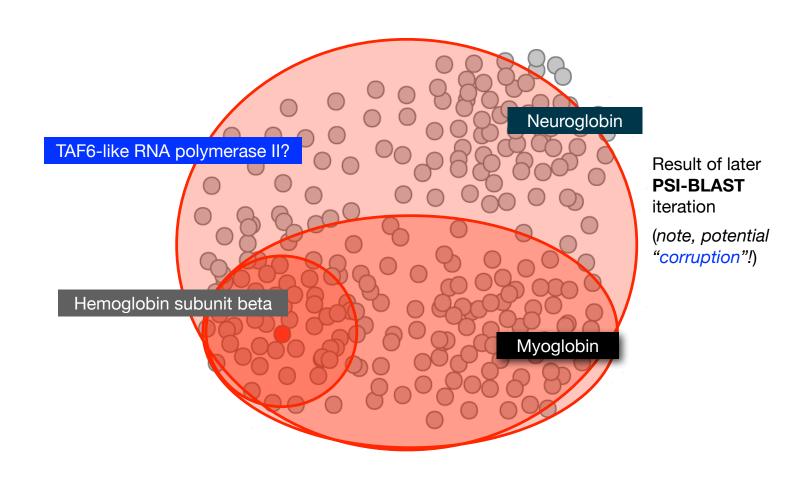


(see Altschul et al., Nuc. Acids Res. (1997) 25:3389-3402)









Description	Max score	Total score	Query	E value	Ident	Accession
hemoglobin subunit beta [Homo sapiens]	301	301	100%	2e-106	100%	NP_000509.1
hemoglobin subunit delta [Homo sapiens]	284	284	100%	7e-100	93%	NP_000510.1
hemoglobin subunit epsilon [Homo sapiens]	240	240	100%	2e-82	76%	NP_005321.1
hemoglobin subunit gamma-2 [Homo sapiens]	235	235	100%	2e-80	73%	NP_000175.1
hemoglobin subunit gamma-1 [Homo sapiens]	232	232	100%	3e-79	73%	NP_000550.2
hemoglobin subunit alpha [Homo sapiens]	114	114	97%	7e-33	43%	NP_000508.1
hemoglobin subunit zeta [Homo sapiens]	100	100	97%	3e-27	36%	NP_005323.1

Description	Max score	Total score	Query	E value	Ident	Accession
hemoglobin subunit beta [Homo sapiens]	301	301	100%	2e-106	100%	NP_000509.1
hemoglobin subunit delta [Homo sapiens]	284	284	100%	7e-100	93%	NP_000510.1
hemoglobin subunit epsilon [Homo sapiens]	240	240	100%	2e-82	76%	NP_005321.1
hemoglobin subunit gamma-2 [Homo sapiens]	235	235	100%	2e-80	73%	NP_000175.1
hemoglobin subunit gamma-1 [Homo sapiens]	232	232	100%	3e-79	73%	NP_000550.2
hemoglobin subunit alpha [Homo sapiens]	114	114	97%	7e-33	43%	NP_000508.1
hemoglobin subunit zeta [Homo sapiens]	100	100	97%	3e-27	36%	NP_005323.1
myoglobin [Homo sapiens]	80.5	80.5	5 97%	2e-19	26%	NP_005359.1
neuroglobin [Homo sapiens]	54.7	54.7	92%	2e-09	23%	NP_067080.1

New relevant globins found only by PSI-BLAST

Description	Max score	Total score	Query	E value	Ident	Accession	
hemoglobin subunit beta [Homo sapiens]	301	301	100%	2e-106	100%	NP_000509.1	
hemoglobin subunit delta [Homo sapiens]	284	284	100%	7e-100	93%	NP_000510.1	
hemoglobin subunit epsilon [Homo sapiens]	240	240	100%	2e-82	76%	NP_005321.1	
hemoglobin subunit gamma-2 [Homo sapiens]	235	235	100%	2e-80	73%	NP_000175.1	
hemoglobin subunit gamma-1 [Homo sapiens]	232	232	100%	3e-79	73%	NP_000550.2	
hemoglobin subunit alpha [Homo sapiens]	114	114	97%	7e-33	43%	NP_000508.1	
hemoglobin subunit zeta [Homo sapiens]	100	100	97%	3e-27	36%	NP_005323.1	
myoglobin [Homo sapiens]	80.5	80.5	97%	2e-19	26%	NP_005359.1	;
neuroglobin [Homo sapiens]	54.7	54.7	92%	2e-09	23%	NP_067080.1	
myoglobin [Homo sapiens]	159	159	97%	3e-50	26%	NP_005359.1	
hemoglobin subunit alpha [Homo sapiens]	151	151	97%	3e-47	42%	NP_000508.1	
hemoglobin subunit mu [Homo sapiens]	147	147	97%	6e-46	35%	NP_001003938.1	
hemoglobin subunit theta-1 [Homo sapiens]	147	147	97%	2e-45	37%	NP_005322.1	
neuroglobin [Homo sapiens]	134	134	92%	3e-40	23%	NP_067080.1	,
PREDICTED: cytoglobin isoform X2 [Homo sapiens]	115	115	66%	3e-33	25%	XP_016879605.1	
PREDICTED: microtubule cross-linking factor 1 isoform X1 [Homo sapie	46.3	46.3	27%	7e-06	39%	XP_011523942.1	
PREDICTED: microtubule cross-linking factor 1 isoform X4 [Homo sapie	46.3	46.3	27%	7e-06	39%	XP_005258156.1	

Inclusion of irrelevant hits can lead to PSSM corruption

# YOUR TURN!

 There are four required and one optional hands-on sections including:

1.	Limits of using BLAST	[~10 mins]
2.	Using PSI-BLAST	[~30 mins]
3.	Examining conservation patterns  — BREAK [15 mins]—	[~20 mins]
4.	[Optional] Using HMMER	[~10 mins]
5.	Divergence of protein sequence and structure	[~25 mins]

- ▶ Please do answer the last review question (Q20).
- ▶ We encourage <u>discussion</u> at your **Table** and on **Piazza**!

```
✓Query_73613
                          MVHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTORFFE-SFGDLSTPDAVM-GNPKVKAHGKKVLGAF
                 1
✓NP 000510.1
                 1
                          MVHLTPEEKTAVNALWGKV--NVDAVGGEALGRILVVYPWTORFFE-SFGDLSSPDAVM-GNPKVKAHGKKVLGAF
✓NP 000175.1
                 1
                          MGHFTEEDKATITSLWGKV--NVEDAGGETLGRLLVVYPWTORFFD-SFGNLSSASAIM-GNPKVKAHGKKVLTSL
✓NP 000509.1
                 1
                          MVHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTORFFE-SFGDLSTPDAVM-GNPKVKAHGKKVLGAF
✓NP 005321.1
                 1
                          MVHFTAEEKAAVTSLWSKM--NVEEAGGEALGRLLVVYPWTQRFFD-SFGNLSSPSAIL-GNPKVKAHGKKVLTSF

√NP 000550.2

                 1
                          MGHFTEEDKATITSLWGKV--NVEDAGGETLGRLLVVYPWTORFFD-SFGNLSSASAIM-GNPKVKAHGKKVLTSL
✓ NP 005323.1
                 1
                          -MSLTKTERTIIVSMWAKISTQADTIGTETLERLFLSHPQTKTYFP-HF-----DLHpGSAQLRAHGSKVVAAV

√NP 000508.1

                 1
                          -MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFP-HF-----DLShGSAOVKGHGKKVADAL

✓ XP 005257062.1 1

                      [15]SEELSEAERKAVQAMWARLYANCEDVGVAILVRFFVNFPSAKQYFS-QFKHMEDPLEME-RSPQLRKHACRVMGAL
✓NP 001003938.1 1
                          --MLSAQERAQIAQVWDLIAGHEAQFGAELLLRLFTVYPSTKVYFP-HL-----SACQ-DATQLLSHGQRMLAAV
✓NP 005322.1
                 1
                          -MALSAEDRALVRALWKKLGSNVGVYTTEALERTFLAFPATKTYFS-H-----LDLSpGSSQVRAHGOKVADAL

√NP 599030.1

                 1
                      [15]SEELSEAERKAVOAMWARLYANCEDVGVAILVRFFVNFPSAKOYFS-OFKHMEDPLEME-RSPOLRKHACRVMGAL

✓ XP 016879605.1 1

                          -----MEDPLEME-RSPOLRKHACRVMGAL
✓NP 001349775.1 1
                          -MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFD-KFKHLKSEDEMK-ASEDLKKHGATVLTAL

√NP 067080.1

                 1
                          ---MERPEPELIROSWRAVSRSPLEHGTVLFARLFALEPDLLPLFOYNCROFSSPEDCL-SSPEFLDHIRKVMLVI
✓NP 001369741.1 1
                                                                 ----MK-ASEDLKKHGATVLTAL
✓Query_73613
                 73
                      SDGLAHLDNLKGT---FATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH
                                                                                                         147
✓ NP 000510.1
                 73
                      SDGLAHLDNLKGT---FSQLSELHCDKLHVDPENFRLLGNVLVCVLARNFGKEFTPQMQAAYQKVVAGVANALAHKYH
                                                                                                          147
✓ NP 000175.1
                 73
                      GDAIKHLDDLKGT---FAQLSELHCDKLHVDPENFKLLGNVLVTVLAIHFGKEFTPEVQASWQKMVTGVASALSSRYH
                                                                                                          147

√NP 000509.1

                 73
                      SDGLAHLDNLKGT---FATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH
                                                                                                          147
✓NP 005321.1
                 73
                      GDAIKNMDNLKPA---FAKLSELHCDKLHVDPENFKLLGNVMVIILATHFGKEFTPEVQAAWQKLVSAVAIALAHKYH
                                                                                                          147
✓NP 000550.2
                 73
                      GDATKHLDDLKGT---FAOLSELHCDKLHVDPENFKLLGNVLVTVLAIHFGKEFTPEVOASWOKMVTAVASALSSRYH
                                                                                                          147
✓ NP 005323.1
                 68
                      GDAVKSIDDIGGA---LSKLSELHAYILRVDPVNFKLLSHCLLVTLAARFPADFTAEAHAAWDKFLSVVSSVLTEKYR
                                                                                                          142

√NP 000508.1

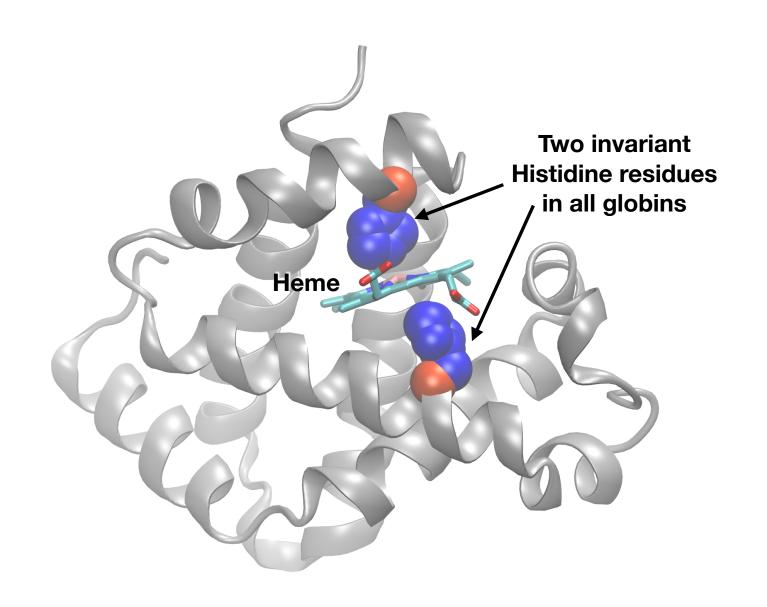
                 68
                      TNAVAHVDDMPNA---LSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR
                                                                                                          142

✓ XP 005257062.1

                 90
                      NTVVENLHDPDKVssvLALVGKAHALKHKVEPVYFKILSGVILEVVAEEFASDFPPETQRAWAKLRGLIYSHVTAAYK [35]
                                                                                                          202
✓NP 001003938.1
                 67
                      GAAVOHVDNLRAA---LSPLADLHALVLRVDPANFPLLIOCFHVVLASHLODEFTVOMOAAWDKFLTGVAVVLTEKYR
                                                                                                          141
✓NP 005322.1
                 68
                      SLAVERLDDLPHA---LSALSHLHACOLRVDPASFOLLGHCLLVTLARHYPGDFSPALQASLDKFLSHVISALVSEYR
                                                                                                          142
✓NP 599030.1
                      NTVVENLHDPDKVssvLALVGKAHALKHKVEPVYFKILSGVILEVVAEEFASDFPPETQRAWAKLRGLIYSHVTAAYK[23]
                 90
                                                                                                         190

✓ XP 016879605.1

                 25
                      NTVVENLHDPDKVssvLALVGKAHALKHKVEPVYFKILSGVILEVVAEEFASDFPPETORAWAKLRGLIYSHVTAAYK [35]
                                                                                                         137
✓NP 001349775.1
                 74
                      GGILKKKGHHEAE---IKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYK[ 6]
                                                                                                         154
✓NP 067080.1
                 73
                      DAAVTNVEDLSSLeeyLASLGRKHRA-VGVKLSSFSTVGESLLYMLEKCLGPAFTPATRAAWSOLYGAVVOAMSRGWD[ 2]
                                                                                                         151
✓NP 001369741.1 19
                      GGILKKKGHHEAE---IKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYK[ 6]
                                                                                                         99
```



# HW Q: Make your own PSSM

You can work on it now in-class!

# YOUR TURN!

 There are four required and one optional hands-on sections including:

1.	Limits of using BLAST	[~10 mins]
2.	Using PSI-BLAST	[~30 mins]
3.	Examining conservation patterns  — BREAK [15 mins]—	[~20 mins]
4.	[Optional] Using HMMER	[~10 mins]
5.	Divergence of protein sequence and structure	[~25 mins]

- ▶ Please do answer the last review question (Q20).
- ▶ We encourage <u>discussion</u> at your **Table** and on **Piazza**!

# Problems with PSSMs: Positional dependencies

Do not capture positional dependencies

WEIRD
WEIRD
WEIQH
WEIRD
WEIQH

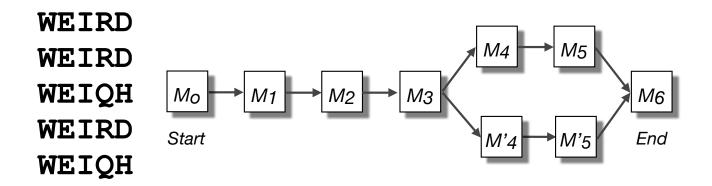
D					0.6
Е		-			
Ι					0.4
ı			I		
Q				0.4	
R				0.6	
W	I				

**Note**: We <u>never</u> see **QD** or **RH**, we only see **RD** and **QH**. However, P(RH)=0.24, P(QD)=0.24, while P(QH)=0.16

# Markov chains: Positional dependencies



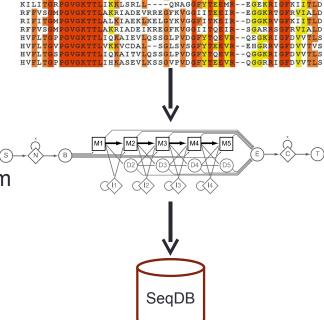
The connectivity or **topology** of a Markov chain can easily be designed to capture dependencies and variable length motifs.



Recall that a PSSM for this motif would give the sequences **WEIRD** and **WEIRH** equally good scores even though the **RH** and **QR** combinations were not observed

# Use of HMMER

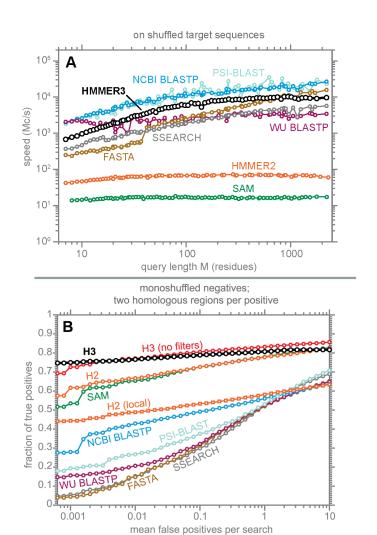
- Widely used by protein family databases
  - Use 'seed' alignments
- Until 2010
  - Computationally expensive
  - Restricted to HMMs constructed from multiple sequence alignments
- Command line application



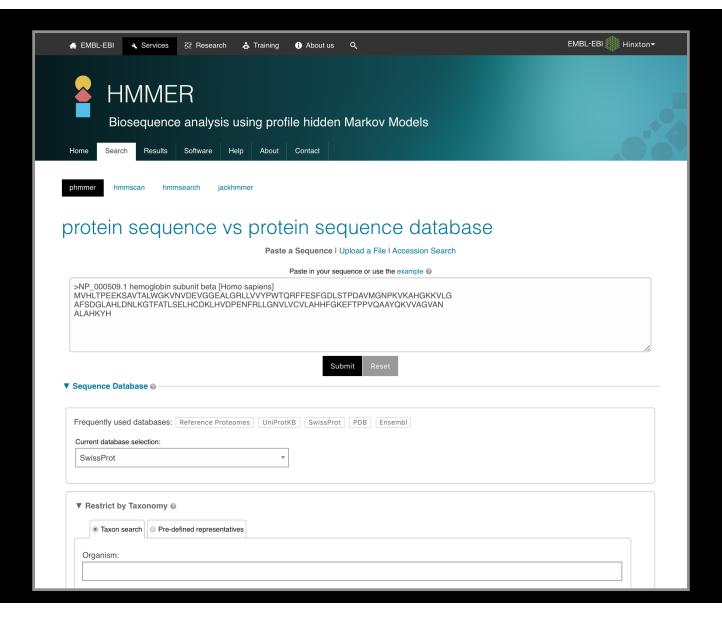


# **HMMER vs BLAST**

	HMMER	BLAST							
Program	PHMMER	B LA STP							
Query	Single sequence								
Targe <b>t</b> Databas <b>e</b>	Sequence database								
Program	HMM SCA <b>N</b>	RP SB LA ST							
Query	Single sequence								
Targe <b>t</b> Databas <b>e</b>	Profile HMM database, e.g. Pfam	PSSM database, e.g. CDD							
Program	HM M SE A RCH	P SI-B LA ST							
Query	Profile HM <b>M</b>	PSSM							
Target Databas <b>e</b>	Sequence database								
Program	JA CKHM MER	P SI-B LA ST							
Query	Single sequence								
Target Databas <b>e</b>	Sequence database								



Modified from: S. R. Eddy PLoS Comp. Biol., 7:e1002195, 2011.



Signif	Significant Query Matches (12) in swissprot (v.2018_11)  Customise Customise								
	Target	Description	Species	Cross-references	E-value				
>	HBB_HUMAN₽	Hemoglobin subunit beta	Homo sapiens <i>₫</i>	m s m li 🕸 🗈	6.8e-99				
>	HBD_HUMAN ₪	Hemoglobin subunit delta	Homo sapiens <i>₫</i>	m s m li 🕸 🕩 🗈	1.6e-91				
>	HBE_HUMANØ	Hemoglobin subunit epsilon	Homo sapiens <i>₫</i>	m s m li 🕸 🕩 🗈	1.5e-74				
>	HBG2_HUMAN₽	Hemoglobin subunit gamma-2	Homo sapiens <i>₫</i>	m s m li 🕸 🕞	8.8e-73				
>	HBG1_HUMAN₽	Hemoglobin subunit gamma-1	Homo sapiens <i>₫</i>	m s m s s	6.2e-72				
>	HBA_HUMAN ₪	Hemoglobin subunit alpha	Homo sapiens₫	m s m li 🕸 🕞	3.8e-29				
>	HBAZ_HUMAN®	Hemoglobin subunit zeta	Homo sapiens <i>₫</i>	m s m li 🕸 🕩 🖹	4.5e-23				
>	HBAT_HUMAN ₪	Hemoglobin subunit theta-1	Homo sapiens <i>₫</i>	m s m	5.2e-22				
>	HBM_HUMAN®	Hemoglobin subunit mu	Homo sapiens₫	m s m	3.4e-19				
>	CYGB_HUMAN®	Cytoglobin	Homo sapiens <i>₫</i>	m s m li 🕸 🗈	3.1e-14				
>	MYG_HUMAN®	Myoglobin	Homo sapiens ₪	m s m is * ·	2.3e-06				
>	NGB_HUMAN₽	Neuroglobin	Homo sapiens <i>₫</i>	m s m li 🕸 🗈	0.0017				
(show	all) alignments	Your search took: 0.06 search	ecs	showir	ng rows 1 - 12 of 12				

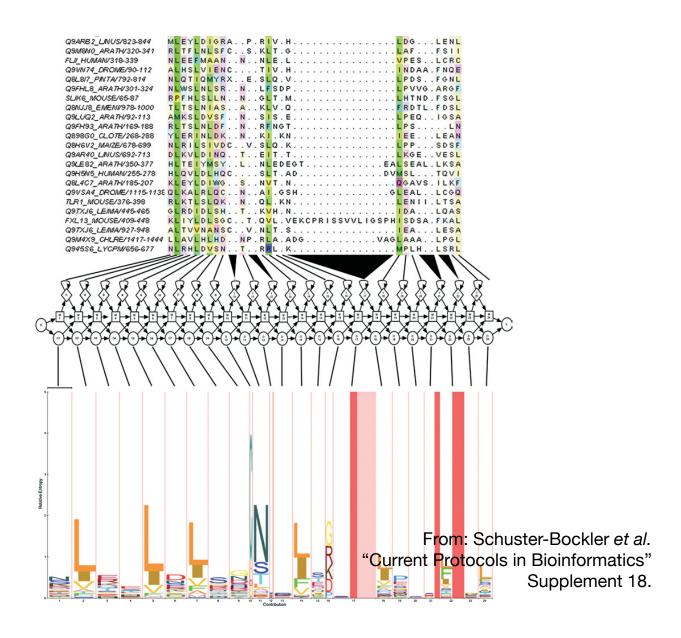
#### **PFAM**: Protein Family Database of Profile HMMs

Comprehensive compilation of both multiple sequence alignments and profile HMMs of protein families.

http://pfam.sanger.ac.uk/

PFAM consists of two databases:

- Pfam-A is a manually curated collection of protein families in the form of multiple sequence alignments and profile HMMs. HMMER software is used to perform searches.
- **Pfam-B** contains additional protein sequences that are automatically aligned. Pfam-B serves as a useful supplement that makes the database more comprehensive.
- Pfam-A also contains higher-level groupings of related families, known as clans



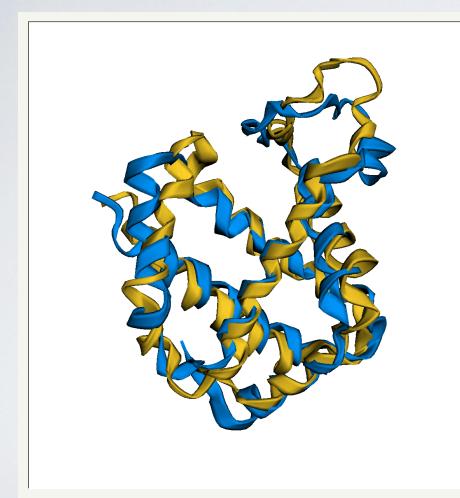
# YOUR TURN!

 There are four required and one optional hands-on sections including:

1.	Limits of using BLAST	[~10 mins]
2.	Using PSI-BLAST	[~30 mins]
3.	Examining conservation patterns	[~20 mins]
	— BREAK [15 mins]—	
4.	[Optional] Using HMMER	[~10 mins]

5. Divergence of protein sequence and structure [~25 mins]

- Please do answer the last review question (Q20).
- We encourage <u>discussion</u> at your **Table** and on **Piazza**!



#### ALIGNMENT

CONTACT MAP

Align	2hb	sB.po	db 146	6 wit	h 4m	pmB.p	db 1	.48									
Twists	s 0	ini-	len 13	36 in	i-rm	sd 3.	05 d	pt-eq	u 143	opt-	rmsd 2	2.65	chain	-rmsd	3.05		
Score	318	.72 a	align-	-len	150	gaps	7 (4	.67%)									
P-valı	ue 3	.26e-	-14 A	fp–nu	m 14	073 I	dent	ity 2	0.67%	Simi	larity	40.	00%				
Block	0	afp 1	17 sc	ore 3	18.7	2 rms	d 3	.05 g	ар 9	(0.06	%)						
					:		:		:		:		:		:		:
Chain	1	2	HLTP	VEKSA'	VTAL	WGKVN	IVD	EVGGE	ALGRL	LVVYP	WTQRFI	FESFG	-DLST	PDAVM	<b>GNPKVI</b>	KAHGK	(VL
			111	1111	1111	11111	. 11	.11111	11111	11111	11111	11111	1111	11111	11111	111111	111
Chain	2	2	ERP	-EPEL	IRQS	WRAVS	RSPL	.EHGTV	LFARL	FALEP	DLLPLI	FQYNC	RQFSS	PEDCL	SSPEFI	DHIR	(VM
Chain	1	69	GAFSI	DGLAH	LDNL	KGTFA	TLSE	LHCD-	-KLHV	DPENF	RLLGN	/LVCV	LAHHF	GKEFT	PPVQA	AYQKV۱	/AG
			11111	11111	1111	11111	.1111	1111	1111	11111	11111	11111	11111	11111	11111	111111	111
Chain	2	70	LVID	AAVTN	VEDL	SSLEE	YLAS	LGRKH	RAVGV	KLSSF	STVGES	SLLYM	LEKCL	GPAFT	PATRA	AWSQLY	/GA
Chain	1	137	VANAI	LAHKY	Н												
			11111	11111	1												
Chain	2	140	VVQA	<b>ISRGW</b>	D												

# Summary

- Find a gene project: You can start working on this now. Submit your responses to Q1-Q4 to get feedback.
- PSI-BLAST algorithm: Application of iterative position specific scoring matrices (PSSMs) to improve BLAST sensitivity
- Hidden Markov models (HMMs): More versatile probabilistic model for detection of remote similarities
- Structure comparisons as gold standards: Structure is more conserved than sequence

# **Homework:** DataCamp!

Install R and RStudio (see website)

Complete the Introduction to R course on DataCamp (Check Piazza for your DataCamp invite and sign up with your UCSD email (i.e. first part of your email address) please.

Let me know **NOW** if you don't have access to DataCamp!