Barry Grant UC San Diego http://thegrantlab.org/bimm143

BIMM 143 Hands-on Lab Session Class 03





Class 3: Hands-on section

http://thegrantlab.org/bimm143/

• • • • •				
Schedule · BIMM 143				
IIII UCSanDiego	Home	Gmail 2	Gcal GitHub	BIMM143 Nee appr algo
BIMM 143		3	Tue 10/05/21	Proj (Par Sequ due
A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of		*	Tue 10/05/21	Opt Dete BLAS Prot
Biological Sciences, UCSD I. Overview Schedule Computer Setup		4	Thu 10/07/21	Bioi Why RStu inter docu
Learning Goals				Data

bioboot.github.io

BGGN213 GDrive Atmosphere CloudLaunch BIMM194 Blink News + + + edleman-Wunsch, Smith-Waterman and BLAST heuristic proaches, Hands on with dot plots, Needleman-Wunsch and BLAST gorithms highlighting their utility and limitations.

Ç

oject: Find a gene project assignment

art 1) Principles of database searching, due in 2 weeks. (Part 2) quence analysis, structure analysis and general data analysis with R le at the end of the quarter.

otional: Advanced sequence alignment and database searching

etecting remote sequence similarity, Database searching beyond

AST, Substitution matrices, Using PSI-BLAST, Profiles and HMMs,

otein structure comparisons as a gold standard.

oinformatics data analysis with **R**

ny do we use R for bioinformatics? R language basics and the tudio IDE, Major R data structures and functions, Using R eractively from the RStudio console. Introducing Rmarkdown cuments.

ata exploration and visualization in **R**



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

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- environment that we have covered to date in class.

The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project)

The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R

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"
- environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project)
 - Monday of week 4 (Apr 21th, 04/21/25).
 - is due 12pm Monday of week 10 (Jun 2nd, 06/02/25).

The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R

description) and the example report for format and content guidance.

Your responses to questions Q1-Q4 are due 12pm San Diego time on

The complete assignment, including responses to all questions,

- A total of 20% of the course grade will be assigned based on the "find-a-gene project assignment"
- environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project
 - Monday of week 4 (Apr 21th, 04/21/25).
 - is due 12pm Monday of week 10 (Jun 2nd, 06/02/25).

The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R

description) and the example report for format and content guidance.

Your responses to questions Q1-Q4 are due 12pm San Diego time on

The complete assignment, including responses to all questions,

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 *#-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 <u>not</u> 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.



\bullet	l				
Schedule · BIMM 143					
	Home	Gmail	Gcal	GitHub	BIMM1

UCSanDiego

BIMM 143

A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of Biological Sciences, UCSD **I**.

Overview

Schedule

Computer Setup

Learning Goals

Assignments & Grading

Ethics Code

3: (Project) Find a Gene Assignment Part 1

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.

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.

- San Diego time.

Videos:

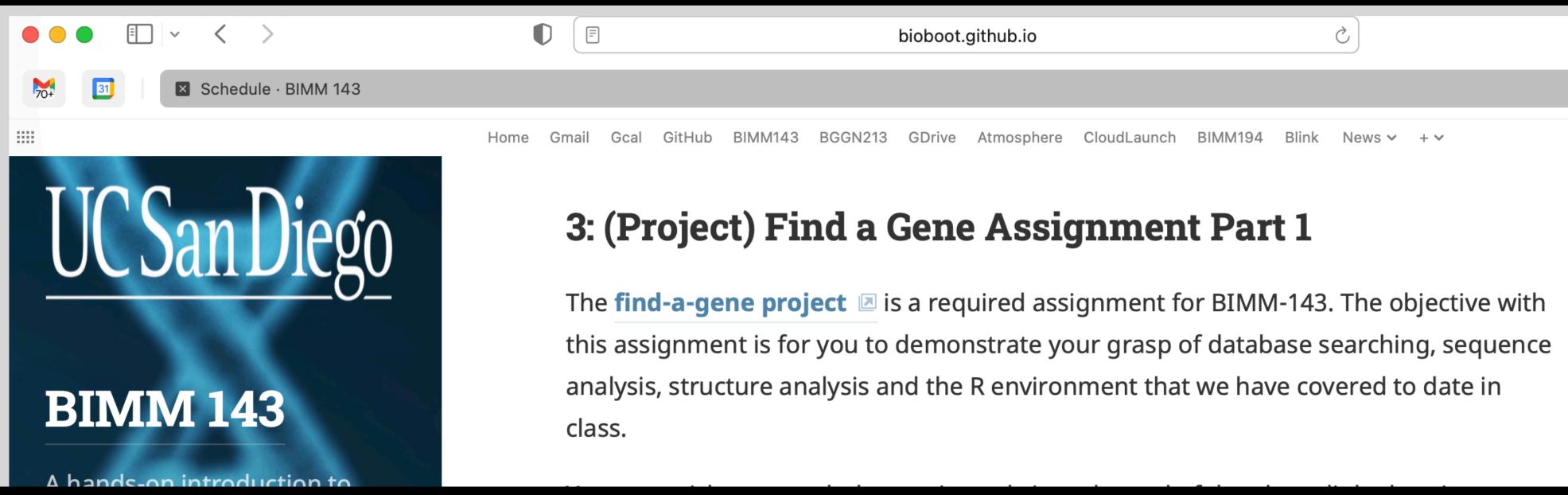
43 BGGN213 GDrive Atmosphere CloudLaunch BIMM194 Blink News 🗸 + 🗸

• Your responses to questions Q1-Q4 are due **Tuesday Oct 19th** (10/19/21) at 12pm

• 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.

• 3.1 - Project introduction I 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 (Apr 21th, 04/21/25).

The complete assignment, including responses to all questions, is due 12pm Monday of week 10 (Jun 2nd, 06/02/25).

Class 3: Hands-on section

http://thegrantlab.org/bimm143/

• • • • •	l			
Schedule · BIMM 143				
	Home	Gmail 2	Gcal GitHub 09/30/21	BIMM143
UCSanDiego			09/30/21	appi algo
BIMM 143		3	Tue 10/05/21	Proj (Par Sequ due
A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of		*	Tue 10/05/21	Opt Dete BLA Prot
Biological Sciences, UCSD I. Overview Schedule Computer Setup		4	Thu 10/07/21	Bioi Why RStu inter docu
Learning Goals				Data

bioboot.github.io

BGGN213 GDrive Atmosphere CloudLaunch BIMM194 Blink News + + + edleman-Wunsch, Smith-Waterman and BLAST heuristic proaches, Hands on with dot plots, Needleman-Wunsch and BLAST gorithms highlighting their utility and limitations.

Ç

oject: Find a gene project assignment

art 1) Principles of database searching, due in 2 weeks. (Part 2) quence analysis, structure analysis and general data analysis with R ie at the end of the quarter.

otional: Advanced sequence alignment and database searching

etecting remote sequence similarity, Database searching beyond

AST, Substitution matrices, Using PSI-BLAST, Profiles and HMMs,

otein structure comparisons as a gold standard.

oinformatics data analysis with **R**

ny do we use R for bioinformatics? R language basics and the tudio IDE, Major R data structures and functions, Using R eractively from the RStudio console. Introducing Rmarkdown cuments.

ata exploration and visualization in **R**



► Details:

Match Score	Mismatch Score	Gap Sco
Sequence 2	GTCGACGC	
Sequence 1	GATTAC	

Compute Optimal Alignment

٢

1

G

Α

Τ

Τ

Α

С

-1

Clear Path

٢

		G	Т	С	G	A	С	G	С	
	0	-2	-4	-6	-8	-10	-12	-14	-16	
1	-2	1	↓ -1	↓ -3	× ♦ -5	-6 + 1	from Di (Due to en G & C	a match		<u>re from Upper cell</u> -2 (The Gap score) =
	-4	↑ -1	0	↓ -2	× ↓ -4	-3 + -2	from Sid 2 (The Ga			ning (max) score is -5
	-6	↑ -3	0	-1	× ♦ -3	-5 • -5	-5	← -7	← -9	
	-8	↑ -5	► ↑ -2	-1	× -2	× ↓ -4	× ♦ -6	-6	↓ -8	
	-10	↑ -7	↑ -4	-3	-2	-1	↓ -3	4 -5	↓ -7	
	-12	↑ -9	↑ -6	-3	× ↑ -4	× ↑ -3	0	↓ -2	★ -4	

▼ Reference:

See the lecture and hands-on session for class 2 for a full discussion of Global, Local, and various Heuristic approaches to biomolecular sequence alignment. Barry J Grant.

ore ٢ Custom Path

-2

G A C G C G Т С **G A T T A C - -**Score = -4

<u>NW App Link</u>





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

Key Question:

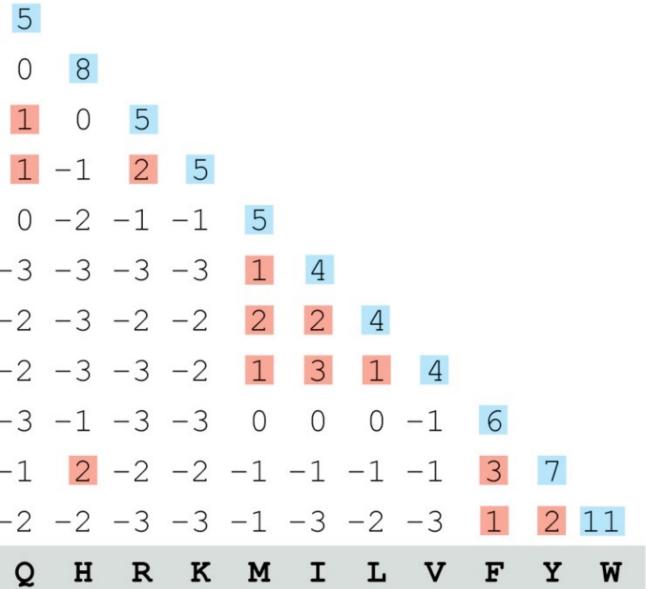
▼ <u>Algorithm paramet</u>	
General Parame	
Max target sequences	Select the maximu
Short queries	 Automatically
Expect threshold	10
Word size	3 💌 📀
Max matches in a query range	0
Scoring Parame	eters
Matrix	BLOSUM62
Gap Costs	Existence: 11 E
Compositional adjustments	Conditional con
Filters and Masl	king
Filter	Low complexit
Mask	Mask for looke Mask lower ca
BLAST	Search databas



By default BLASTp match scores come from the BLOSUM62 matrix

С	9									
S	-1	4					BI	ock	(S <u>S</u>	Su
т	-1	1	5				ob	ser	ve	d 1
Ρ	-3	-1	-1	7			se	que	enc	e
A	0	1	0	-1	4					
G	-3	0	-2	-2	0	6				
N	-3	1	0	-2	-2	0	6			
D	-3	0	-1	-1	-2	-1	1	6		
Е	-4	0	-1	-1	-1	-2	0	2	5	
Q	-3	0	-1	-1	-1	-2	0	0	2	
H	-3	-1	-2	-2	-2	-2	1	-1	0	
R	-3	-1	-1	-2	-1	-2	0	-2	0	
K	-3	0	-1	-1	-1	-2	0	-1	1	
М	-1	-1	-1	-2	-1	-3	-2	-3	-2	
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-
v	-1	-2	0	-2	0	-3	-3	-3	-2	-
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	_
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-
	С	S	т	Ρ	A	G	N	D	Е	

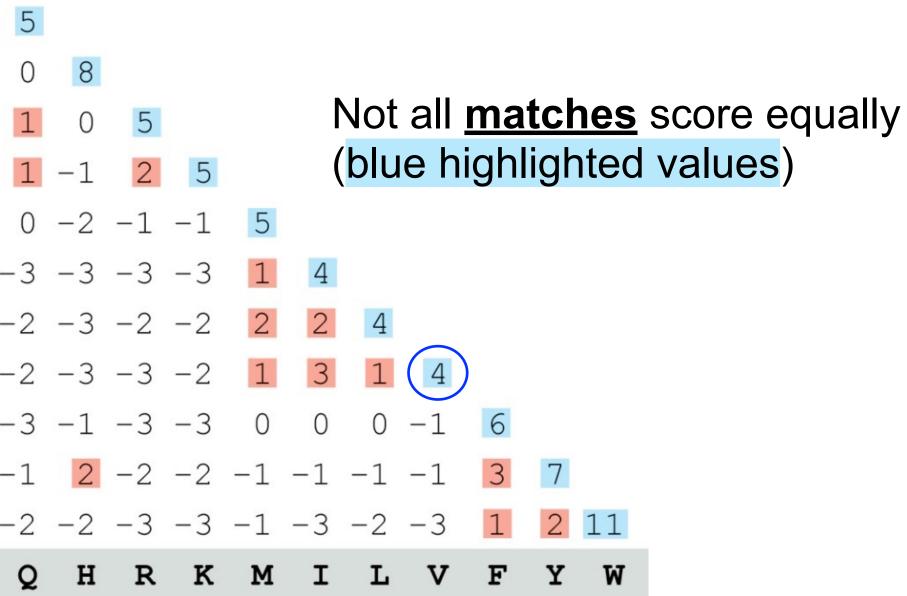
ubstitution <u>Matrix</u>. Scores obtained from frequencies of substitutions in blocks of aligned swith no more than 62% identity.



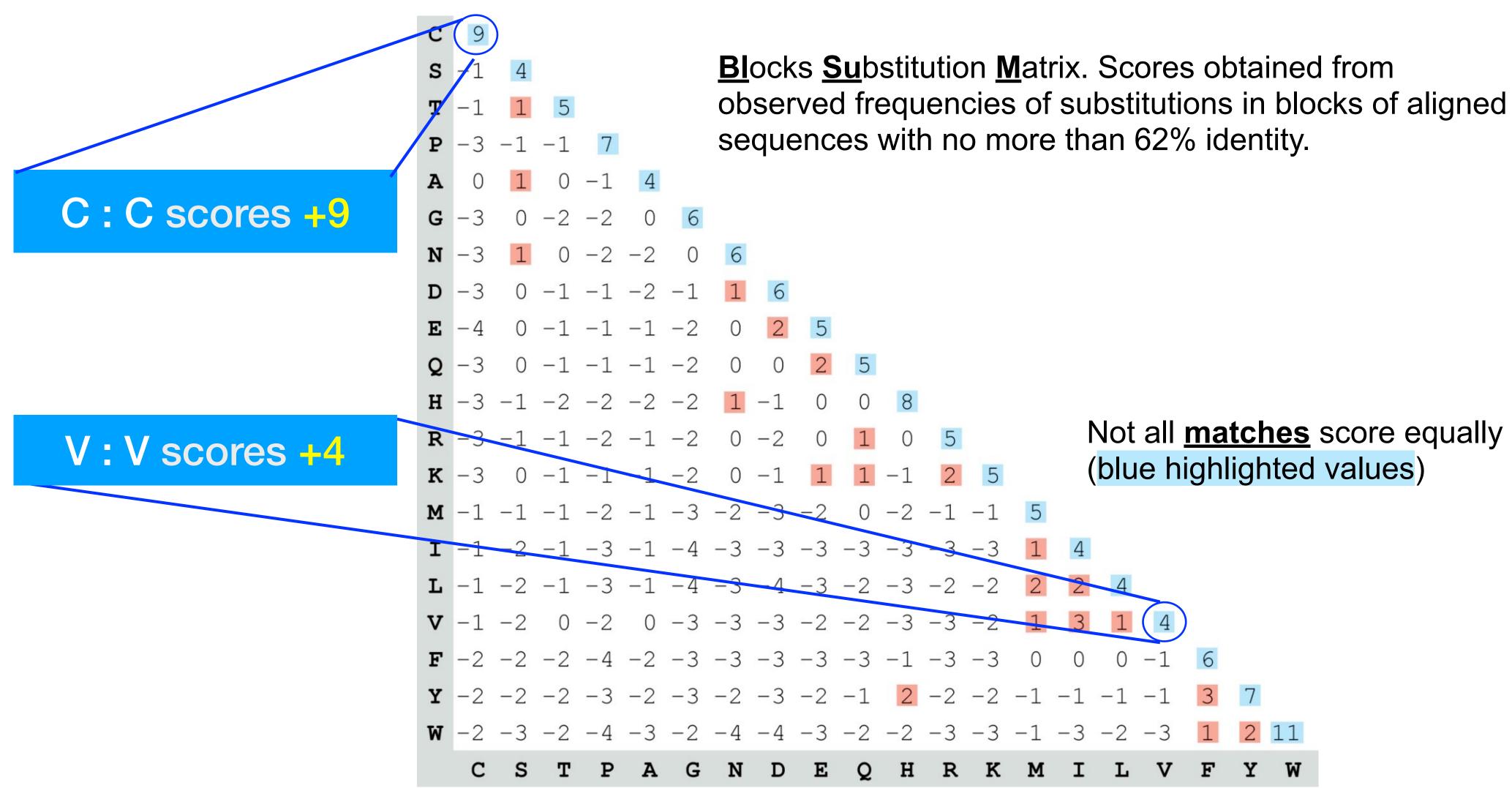
By default BLASTp match scores come from the BLOSUM62 matrix

С	9)								_
S	-1	4					BI	ock	(S <u>S</u>	Su
т	-1	1	5				ob	ser	ve	d .
Ρ	-3	-1	-1	7			se	que	enc	e
A	0	1	0	-1	4					
G	-3	0	-2	-2	0	6				
N	-3	1	0	-2	-2	0	6			
D	-3	0	-1	-1	-2	-1	1	6		
Е	-4	0	-1	-1	-1	-2	0	2	5	
Q	-3	0	-1	-1	-1	-2	0	0	2	
H	-3	-1	-2	-2	-2	-2	1	-1	0	
R	-3	-1	-1	-2	-1	-2	0	-2	0	
K	-3	0	-1	-1	-1	-2	0	-1	1	
М	-1	-1	-1	-2	-1	-3	-2	-3	-2	
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	_
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	_
v	-1	-2	0	-2	0	-3	-3	-3	-2	_
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	_
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	_
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-
	С	S	т	Ρ	A	G	N	D	Е	

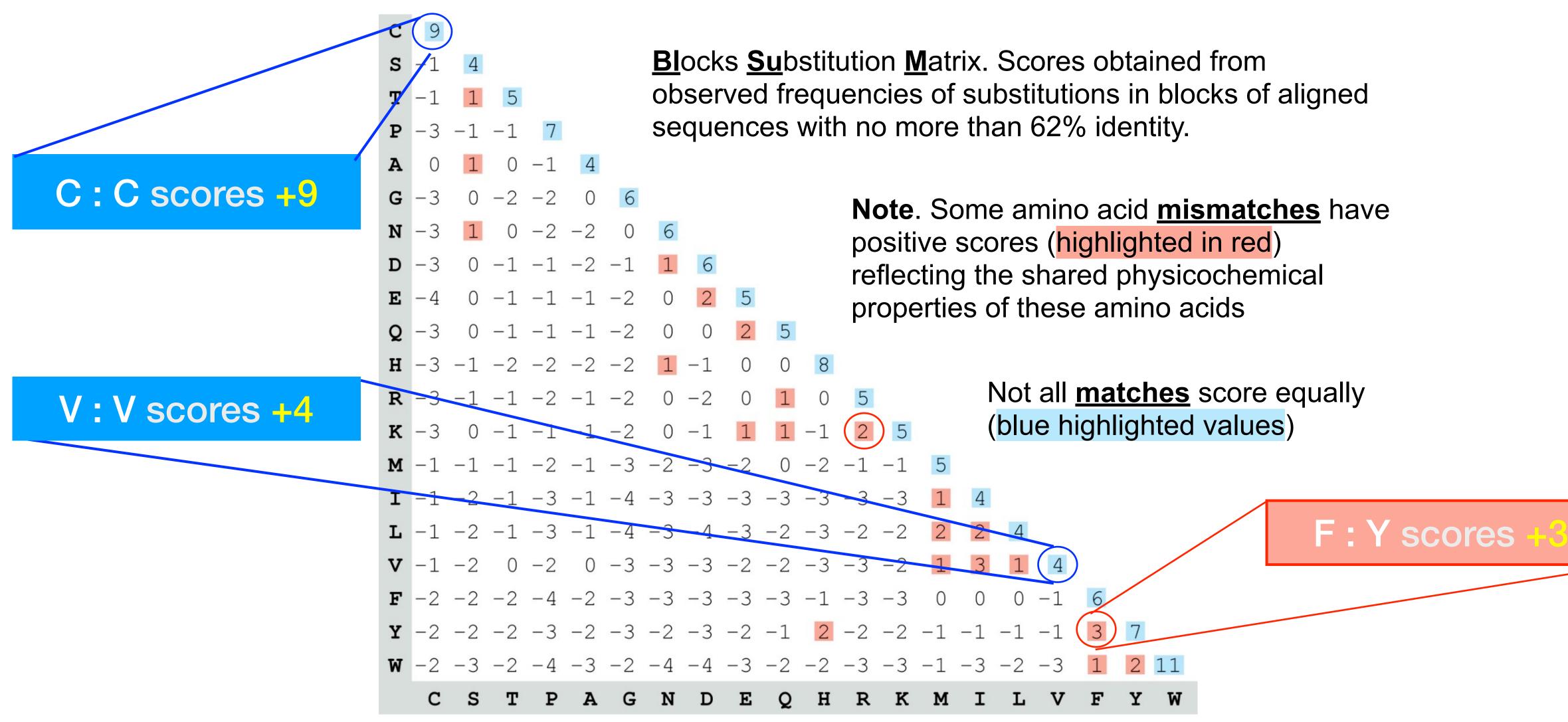
ubstitution <u>Matrix</u>. Scores obtained from frequencies of substitutions in blocks of aligned swith no more than 62% identity.



By default BLASTp match scores come from the BLOSUM62 matrix

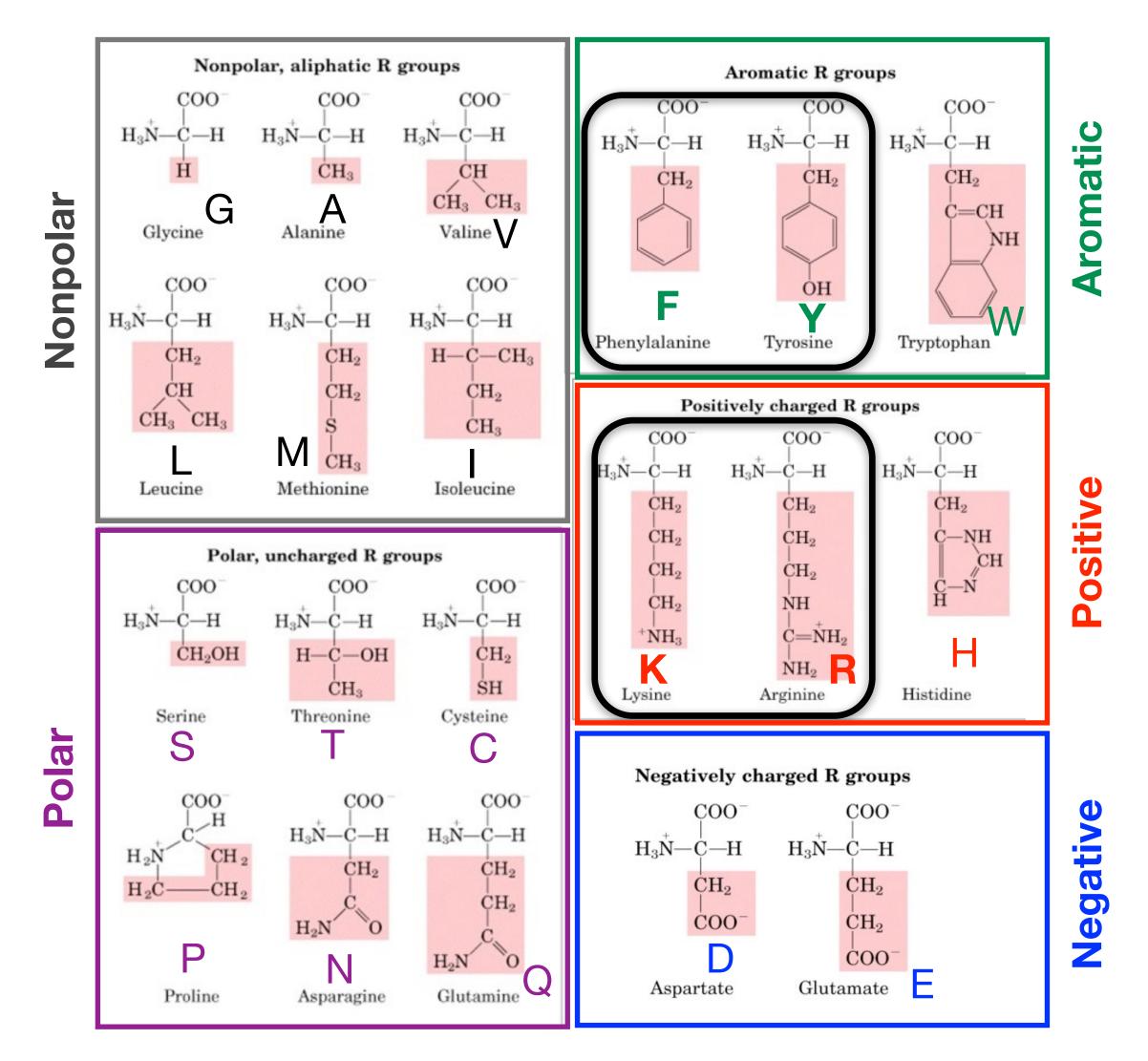


By default BLASTp match scores come from the BLOSUM62 matrix

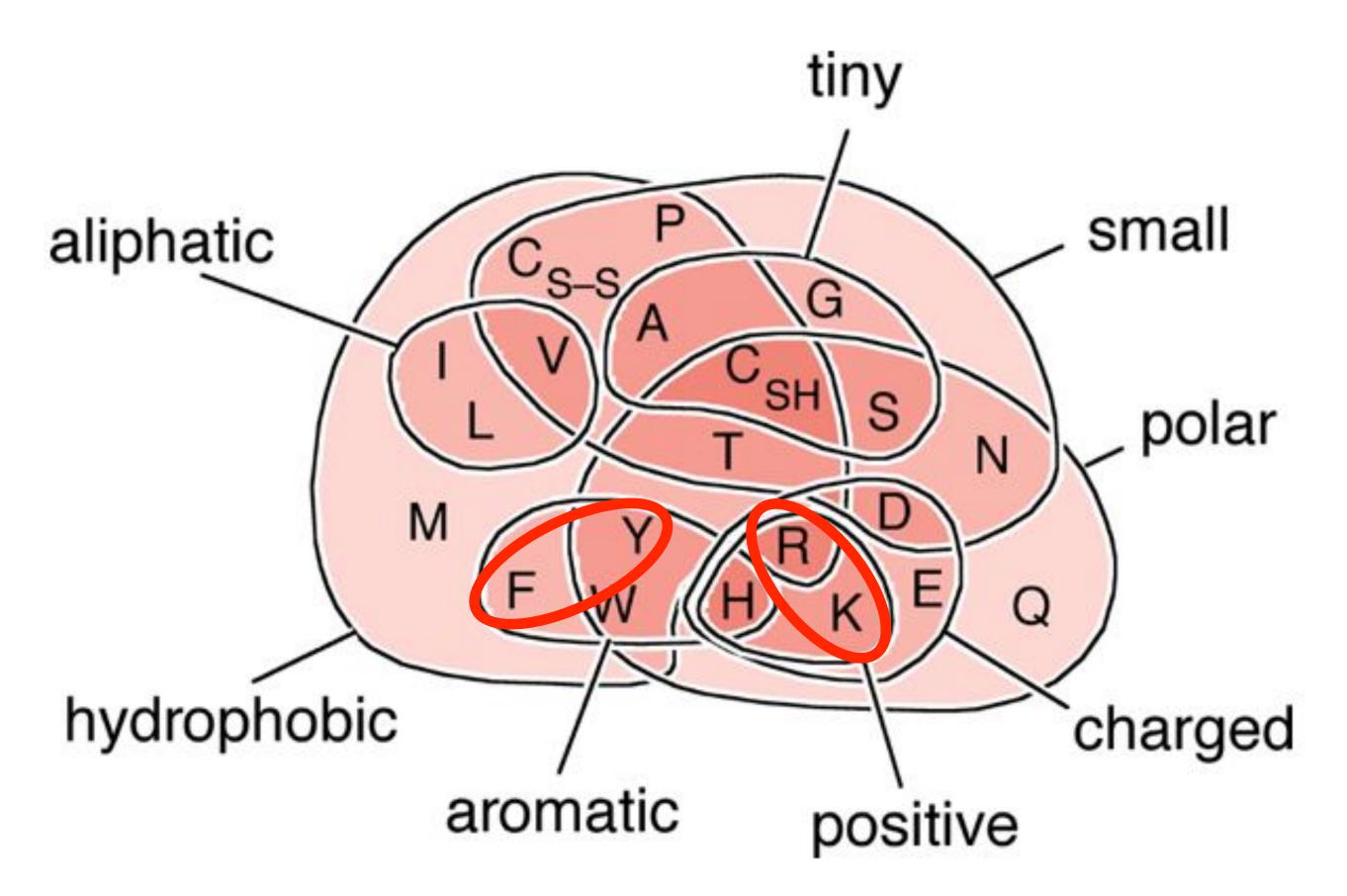




Protein scoring matrices reflect the properties of amino acids



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
 - 2. Using PSI-BLAST
 - 3. Examining conservation patterns

--- BREAK [15 mins]---

- 4. [Optional] Using HMMER
- 5. Divergence of protein sequence and structure
- Please do answer the last review question (Q20).
- We encourage <u>discussion</u> at your Table and on Piazza!

[~10 mins] [~30 mins] [~20 mins]

[~10 mins] [~25 mins]

YOUR TURN!

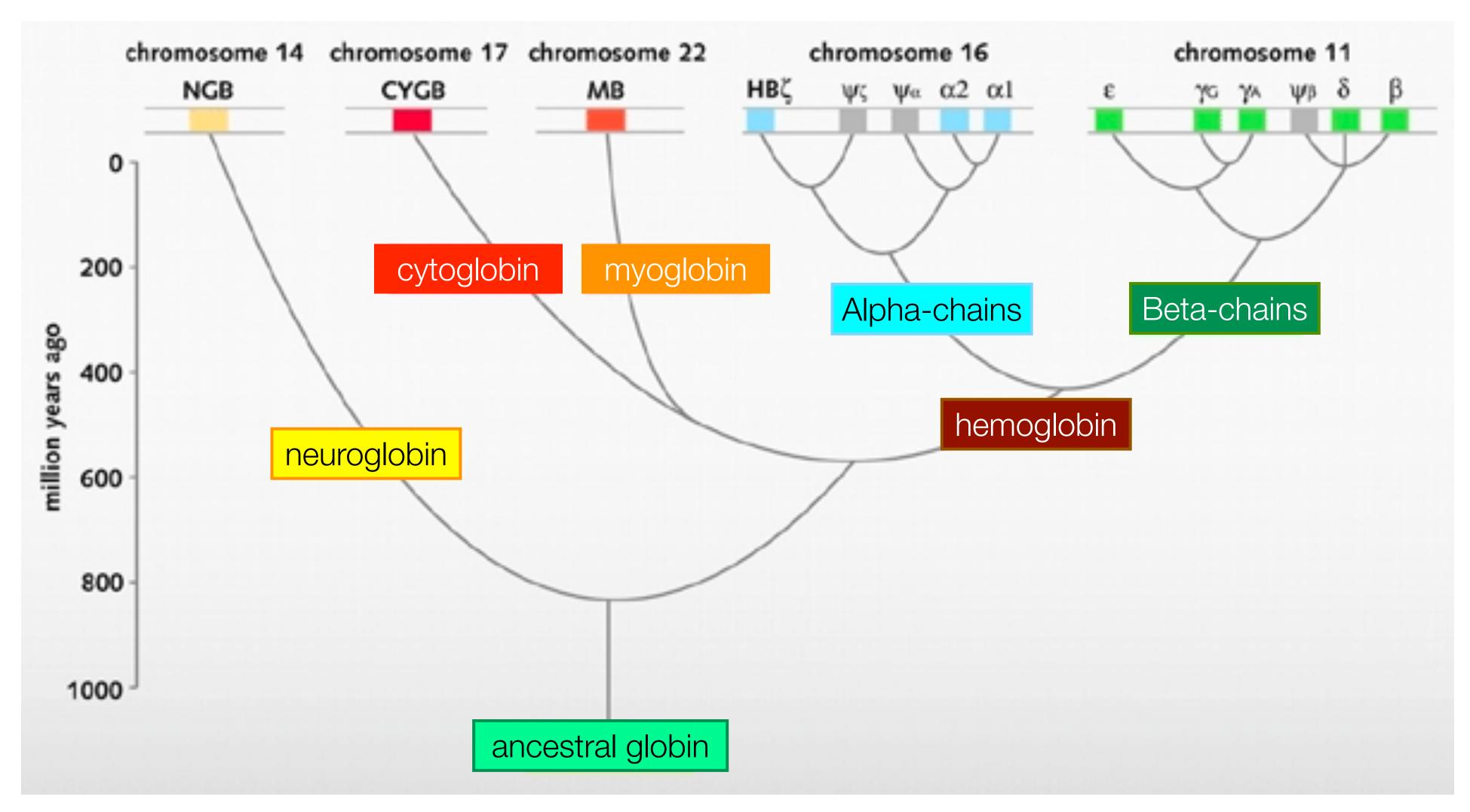
- There are four required and one optional hands-on sections including:
 - **1. Limits of using BLAST**
 - 2. Using PSI-BLAST
 - 3. Examining conservation patterns

--- BREAK [15 mins]---

- 4. [Optional] Using HMMER
- 5. Divergence of protein sequence and structure
- Please do answer the last review question (Q20).
- We encourage <u>discussion</u> at your Table and on Piazza!

[~10 mins] [~30 mins] [~20 mins]

[~10 mins] [~25 mins]



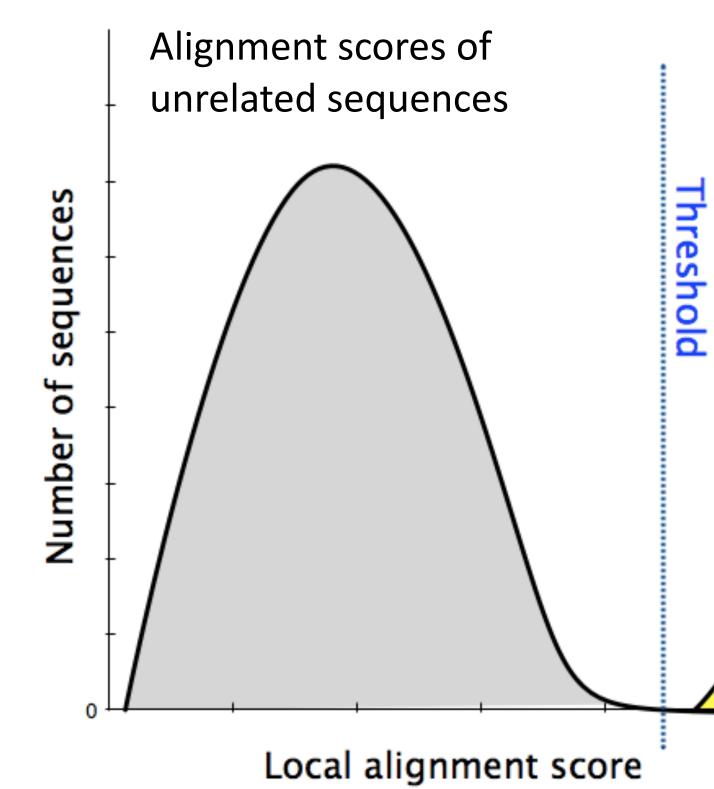
An evolutionary model of human globins.

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?

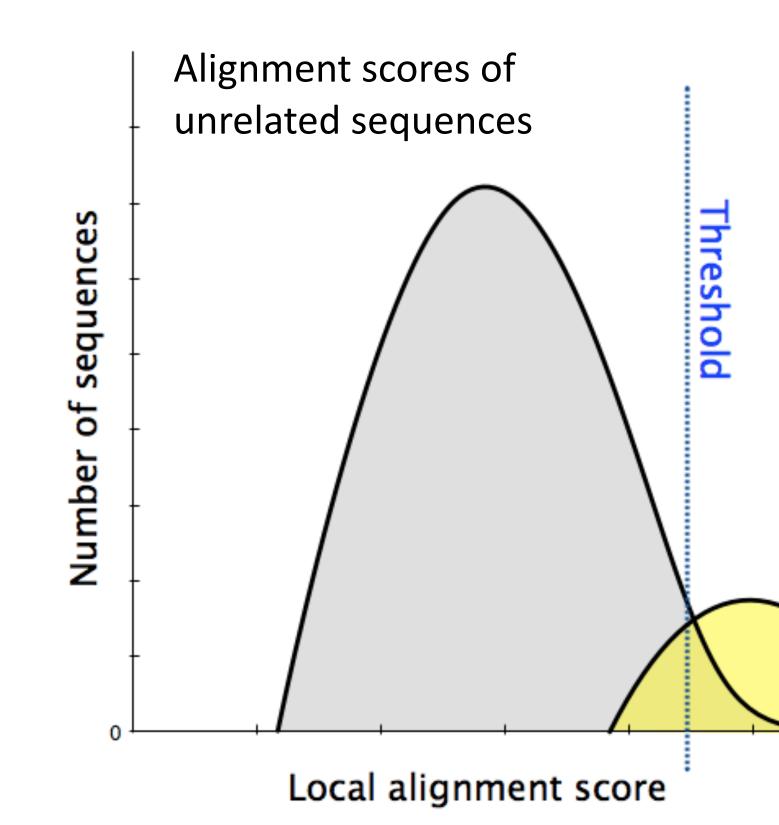
(yellow) from all unrelated sequences (gray)



• Ideally, a threshold separates all query related sequences

Alignment scores of related sequences

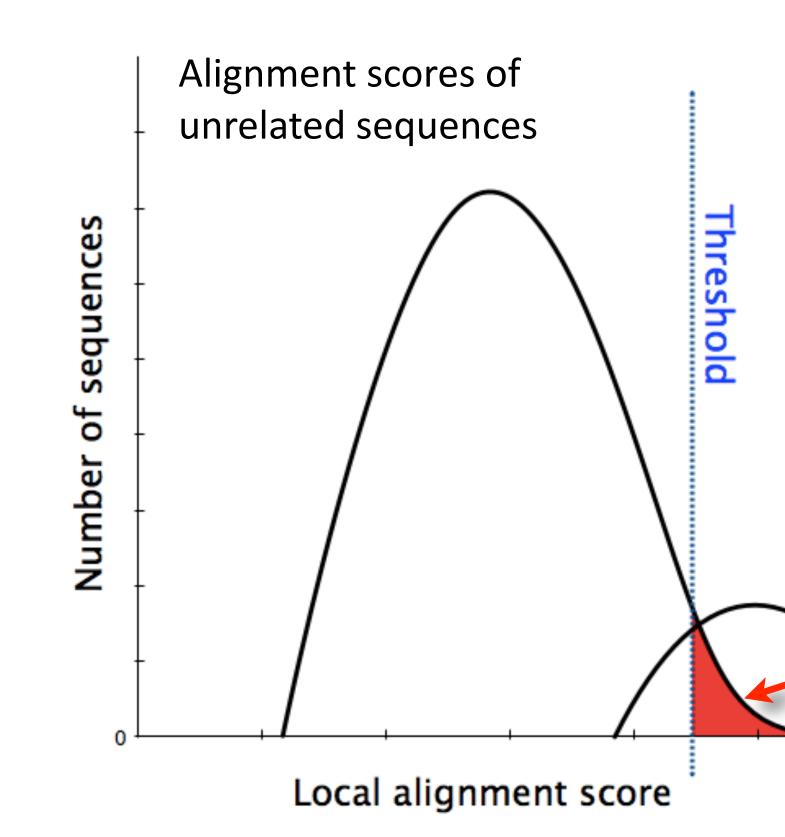
 Unfortunately, often both score distributions overlap unrelated



— The E value describes the expected number of hits with a score above the threshold if the query and database are

> Alignment scores of related sequences

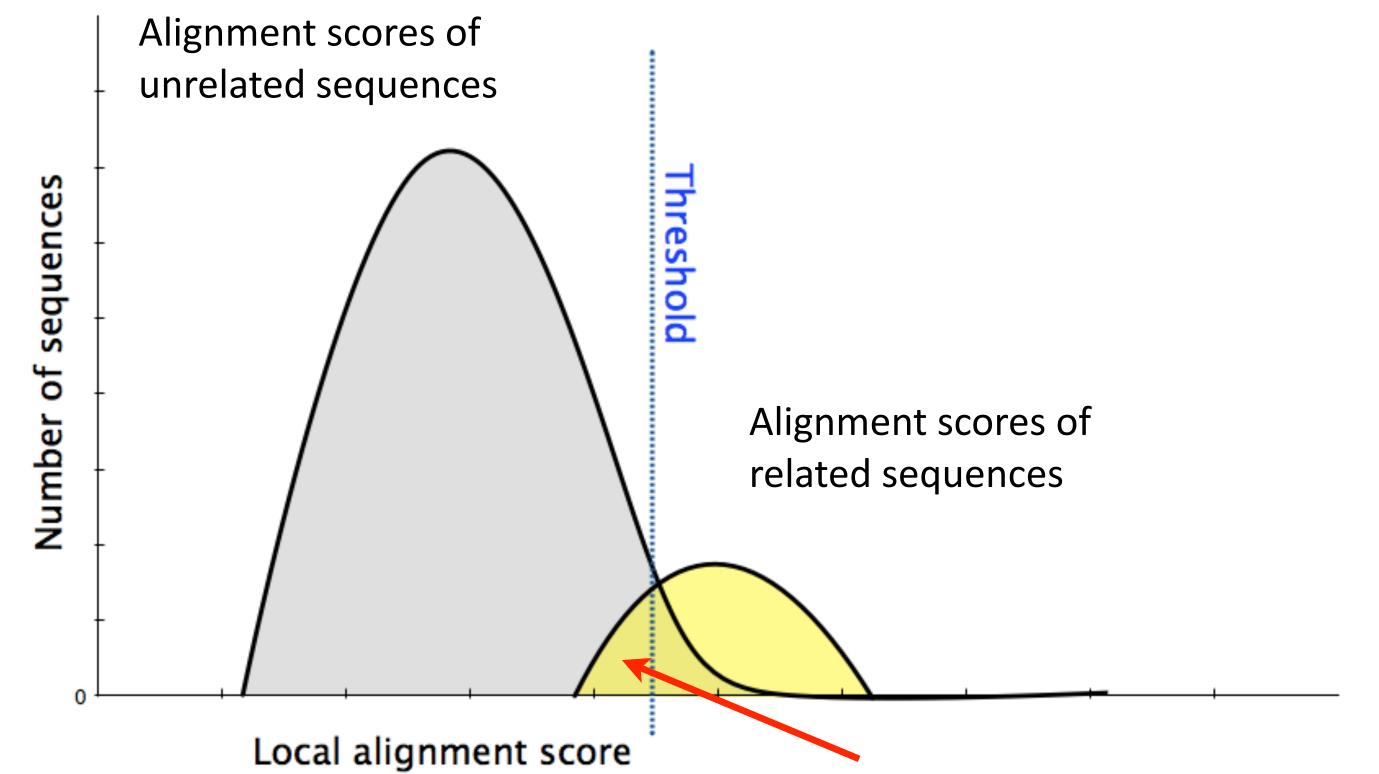
 Unfortunately, often both score distributions overlap unrelated



— The E value describes the expected number of hits with a score above the threshold if the query and database are

> The E-value provides an estimate of the number of false positive hits!

but not reported because of our E-value cutoff? – Lets change the cutoff and see…



• Maybe myoglobin, cytoglobin, neuroglobin etc. are found

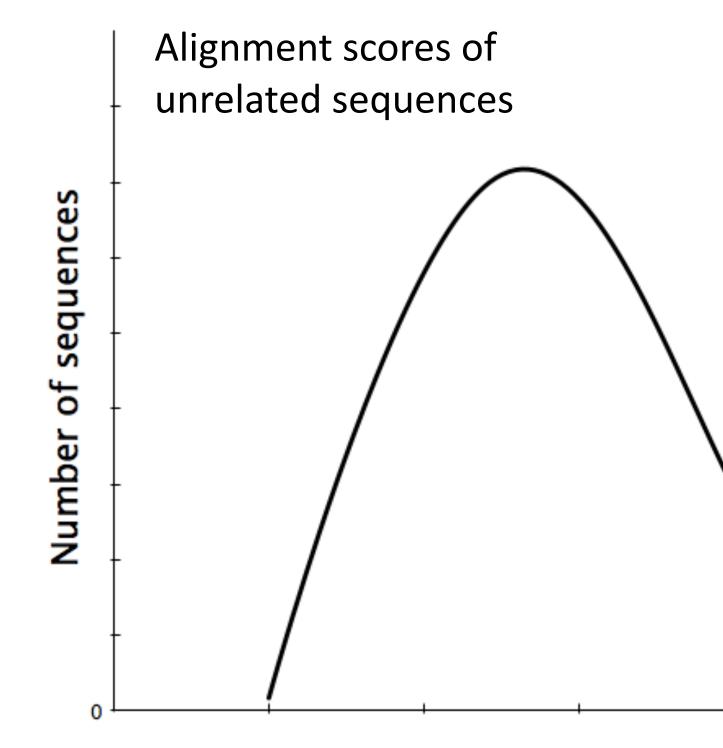
Description

hemoglobin subunit beta

hemoglobin subunit delta

hemoglobin subunit alpha

probable ATP-dependent RNA helicase



Local alignment score

Max score	Query cover	E value	Max ident	Accession
284	100%	0	100%	NP_000510.1
240	100%	0	75.5%	NP_005321.1
114	97%	0	43.45%	NP_000508.1
42.7	10%	0.93	32%	XP_011530405.1

A score of 42.7 or better is expected to occur by chance 93 in 100 times (Evalue = 0.93)

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

42.7

YOUR TURN!

- There are four required and one optional hands-on sections including:
 - 1. Limits of using BLAST
 - 2. Using PSI-BLAST
 - 3. Examining conservation patterns

--- BREAK [15 mins]---

- 4. [Optional] Using HMMER
- 5. Divergence of protein sequence and structure
- Please do answer the last review question (Q20).
- We encourage <u>discussion</u> at your Table and on Piazza!

[~10 mins] [~30 mins] [~20 mins]

[~10 mins] [~25 mins]

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).

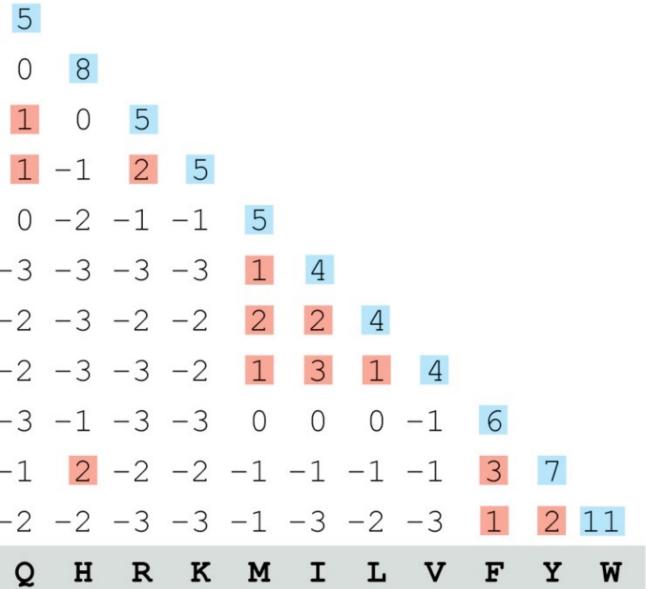
▼ <u>Algorithm paramet</u>	
General Parame	
Max target sequences	Select the maximu
Short queries	 Automatically
Expect threshold	10
Word size	3 💌 📀
Max matches in a query range	0
Scoring Parame	eters
Matrix	BLOSUM62
Gap Costs	Existence: 11 E
Compositional adjustments	Conditional con
Filters and Masl	king
Filter	Low complexit
Mask	Mask for looke Mask lower ca
BLAST	Search databas



By default BLASTp match scores come from the BLOSUM62 matrix

С	9									
S	-1	4					BI	ock	(S <u>S</u>	Su
т	-1	1	5				ob	ser	ve	d 1
Ρ	-3	-1	-1	7			se	que	enc	e
A	0	1	0	-1	4					
G	-3	0	-2	-2	0	6				
N	-3	1	0	-2	-2	0	6			
D	-3	0	-1	-1	-2	-1	1	6		
Е	-4	0	-1	-1	-1	-2	0	2	5	
Q	-3	0	-1	-1	-1	-2	0	0	2	
H	-3	-1	-2	-2	-2	-2	1	-1	0	
R	-3	-1	-1	-2	-1	-2	0	-2	0	
K	-3	0	-1	-1	-1	-2	0	-1	1	
М	-1	-1	-1	-2	-1	-3	-2	-3	-2	
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-
v	-1	-2	0	-2	0	-3	-3	-3	-2	-
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	_
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-
	С	S	т	Ρ	A	G	N	D	Е	

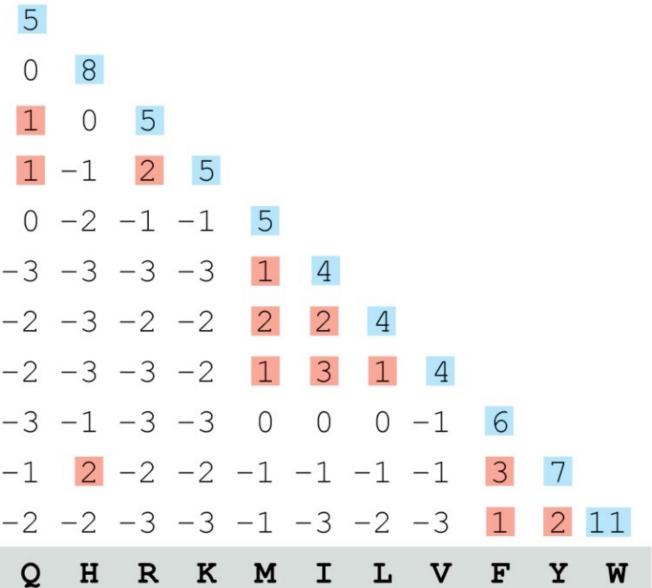
ubstitution <u>Matrix</u>. Scores obtained from frequencies of substitutions in blocks of aligned swith no more than 62% identity.



By default BLASTp match scores come from the BLOSUM62 matrix

C	W -2	Y -2	F -2	v -1	L -1	I -1	M -1	K -3	R -3	H -3	Q -3	E -4		D -3	N -3 D -3		G -3 N -3	G -3 N -3	 A 0 G -3 N -3 	 P - 3 A 0 G - 3 N - 3
c s	2 -3	2 -2	2 -2	-2	-2	-2	- 1	3 0	3 -1	3 -1		3 0		ł 0	3 0 1 0	3 1 3 0 4 0	3 0 3 1 3 0 4 0) 1 3 0 3 1 3 0 4 0	3 -1) 1 3 0 3 1 3 0 4 0	1 3 -1 1
т	-2	-2	-2	0	-1	-1	-1	-1	-1	-2		-1	-1 -1		-1 -1	-1 -1	0 -1 -1	0 -2 0 -1 -1	0 -2 0 -1 -1	5 -1 0 -2 0 -1 -1
Р	-4	-3	-4	-2	-3	-3	-2	-1	-2	-2	-	-1		-1	-1 -1	-2 -1 -1	-2 -2 -1 -1	-1 -2 -1 -1	-1 -2 -1 -1	-1 -2 -1 -1
Α	-3	-2	-2	0	-1	-1	-1	-1	-1	-2	T	_1		-1	-2 -1	-2 -2 -1	0 -2 -2 -1	0 -2 -2 -1	0 -2 -2 -1	0 -2 -2 -1
G	-2	-3	-3	-3	-4	-4	-3	-2	-2	-2	-2	- 2		-2	-1 -2	-1 -2	0 -1 -2	0 -1 -2	0 -1 -2	0 -1 -2
N	-4	-2	-3	-3	-3	-3	-2	0	0	1	0	0		0	0	1 0	1 0	1 0	1 0	1 0
D	-4	-3	-3	-3	-4	-3	-3	-1	-2	-1	0	0	2		2	2	2	2	2	2
Е	-3	-2	-3	-2	-3	-3	-2	1	0	0	2	2	_	_	_	_	_	_	_	_
	-	_	-	_	_	_			1											

Note. All matches of Alanine for Alanine score +4 regardless of their position or context in the molecule.



PSI-BLAST: Position specific iterated BLAST

customized to your query

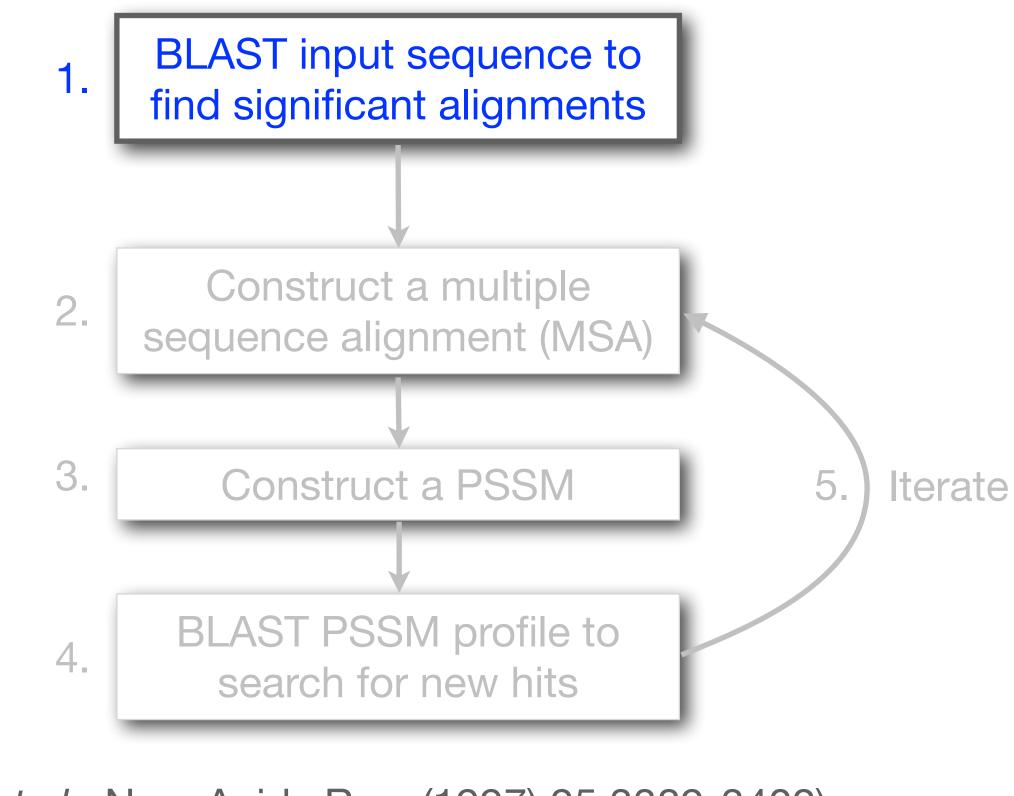
• The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a <u>scoring matrix that is</u>

- customized to your query

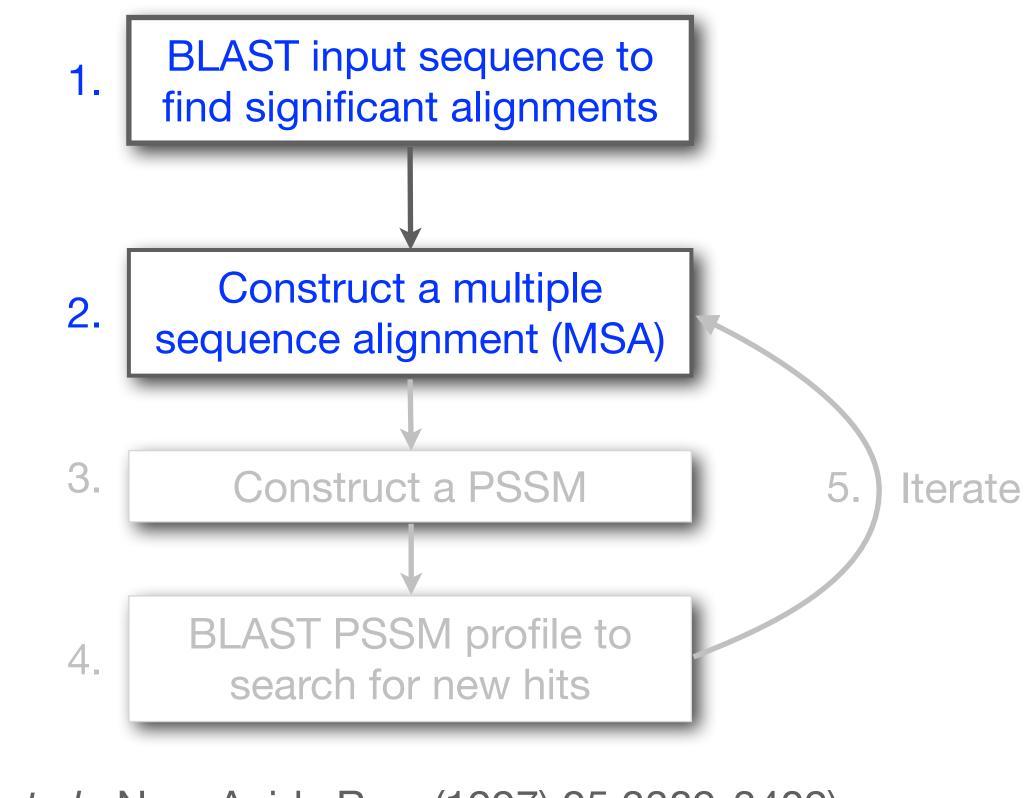
• The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a <u>scoring matrix that is</u>

– PSI-BLAST constructs a multiple sequence alignment from the results of a first round BLAST search and then creates a "profile" or specialized position-specific scoring matrix (PSSM) for subsequent search rounds

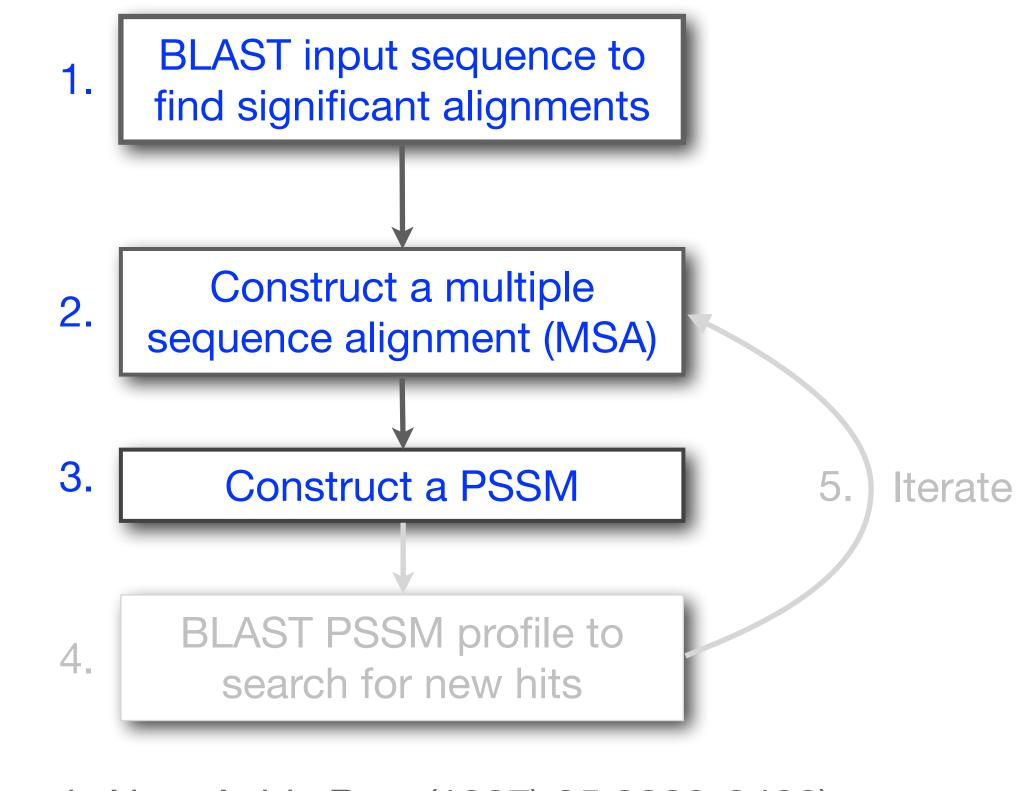
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



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



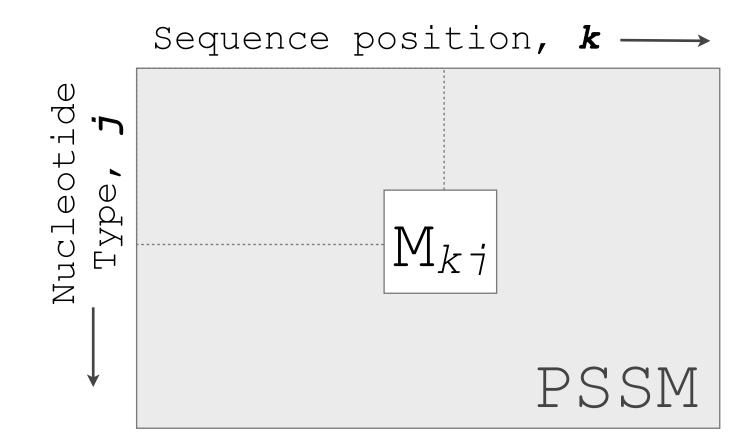
What is a **PSSM**?

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).



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

 M_{kj} score for the *j*th nucleotide at position k

- \mathbf{p}_{kj} probability of nucleotide *j* at position *k*
- **p**^j "background" probability of nucleotide *j*

See Gibskov et al. (1987) PNAS 84, 4355

CCAAATTAGGAAA CCTATTAAGAAAA AAA<mark>TTAGG</mark>AAA AAATTCGGATA TTTCCA ATTTAG' AATTAGGAAA AATTGGC AA ירחרוי **CCAATTTTCAAAA**

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).

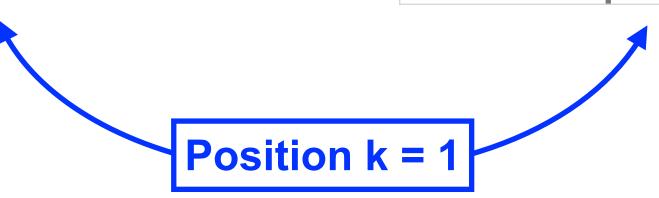
CC	AAA <mark>TT</mark> A <mark>GG</mark> AAA
CC	TATTAAGAAAA
CC	AAA <mark>TT</mark> AGGAAA
CC	AAA <mark>TT<mark>C</mark>GGA<mark>T</mark>A</mark>
CC	CATTTCGAAAA
CC	T <mark>ATTTAGTAT</mark> A
CC	AAA <mark>TT</mark> AGGAAA
CC	AAA <mark>TT</mark> GG <mark>C</mark> AAA
TC	T <mark>ATTTTGGAAA</mark>
CC	AA <mark>TTTT</mark> CAAAA

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:													

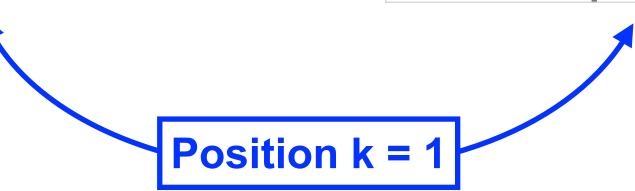
CCAAATTAGGAAA CCAAATTAGGAAA CCAAATTCGGATA CCCATTTCGAAAA CCTATTTAGTATA CCAAATTGGGAAA TCTATTTTGGAAA

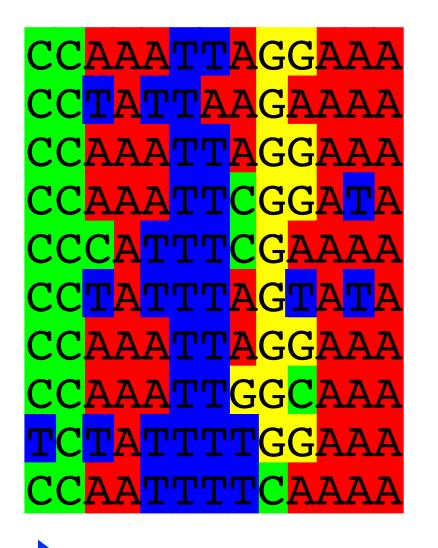
Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A :													
C:													
G:													
T:													



CCAAA TTAGGAAA CCAAA TTAGGAAA CCAAA TTCGGAT CCCATTCGAAA CCCATTCGAAA CCAAA TTAGGAAA CCAAA TTGGGAAA TCTATTTTGGAAA

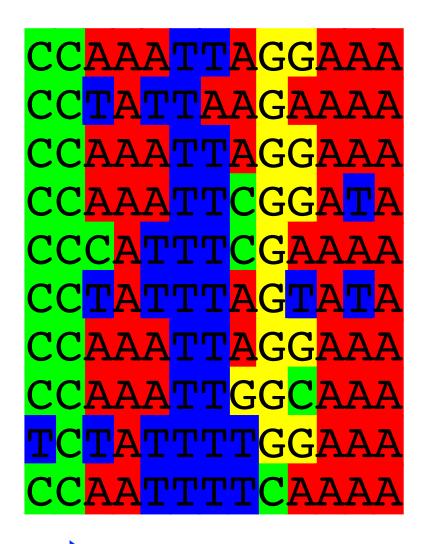
Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A :	0												
C:	9												
G:	0												
T:	1												





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												

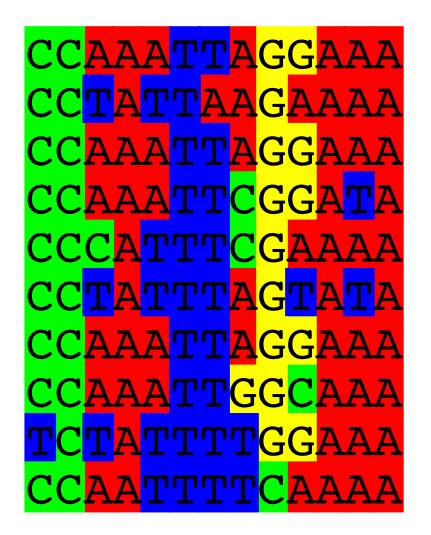




Alignment Counts matrix:

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											

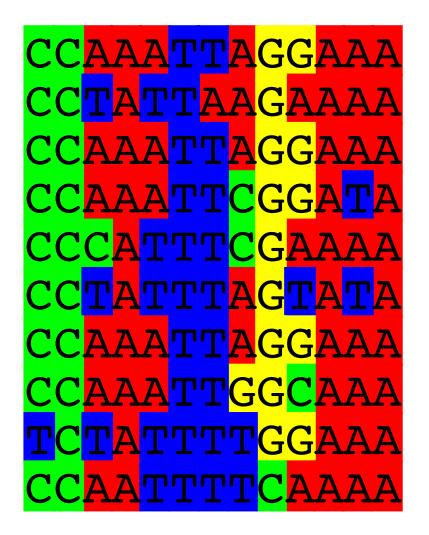
Position k = 2



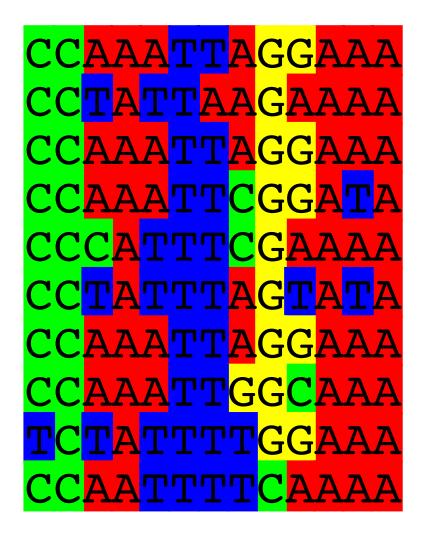
Alignment Counts matrix:

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										

Position k = 3



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	А	Α	[AT]	Т	Т	Α	G	G	Α	Α	Α



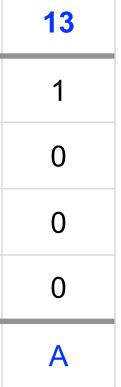
Alignment Counts matrix:

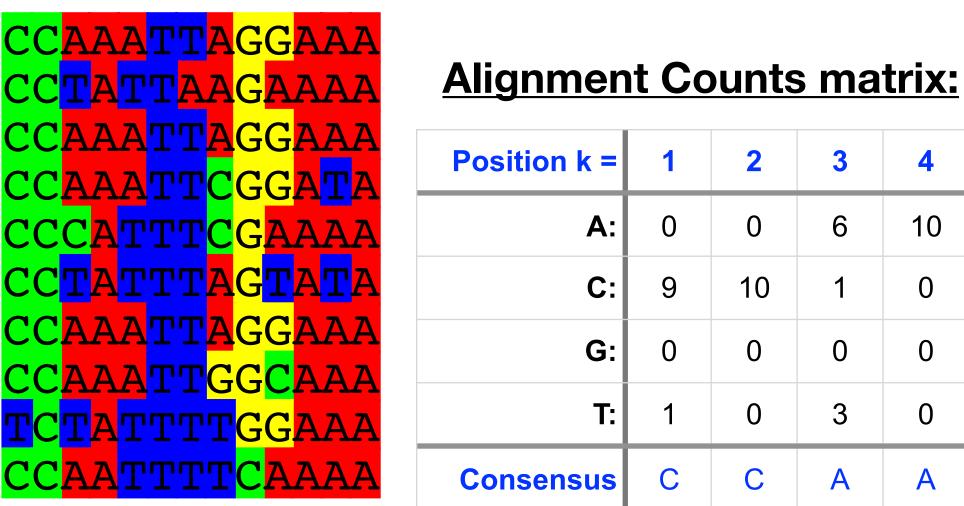
Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A :	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0
Consensus	С	С	Α	Α	[AT]	Т	Т	Α	G	G	Α	Α	Α

Average Profile (Frequency) matrix:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	
A :	0	0	0.6	1	0.5	0	0.1	0.5	0	0.3	1	0.8	
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	А	А	[AT]	Т	Т	А	G	G	Α	А	

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





Or the "score (M_{kj}) matrix" = PS<u>S</u>M

- C_{kj} Number of *j*th type nucleotide at position k
- Total number of aligned sequences Ζ
- "background" probability of nucleotide j $\mathbf{p}_{\mathtt{j}}$
- probability of nucleotide j at position k \mathbf{p}_{kj}

3	4	5	6	7	8	9	10	11	12	13
6	10	5	0	1	5	0	3	10	8	10
1	0	0	0	0	2	1	1	0	0	0
0	0	0	0	0	1	9	5	0	0	0
3	0	5	10	9	2	0	1	0	2	0
А	Α	[AT]	Т	Т	Α	G	G	Α	Α	Α

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

Adapted from Hertz and Stormo, Bioinformatics 15:563-577

		Align	ment Ma	atrix	C _k	j	
nome)	Position A: C: G: T:	k = 1 0 9 0	2 0 10 0 0	1 0	4 10 0 0 0	
		С:	• 0 • 1 , j=A ¹ :	10	1	0	p_j
s):			=1, <i>j</i> =C:		١		1
.4		k =	:1, <i>j</i> =Т:	$M_{\scriptscriptstyle kj}$ =	= log(<u>C</u> _{kj} +	$\frac{p_j}{p_j}$
			Co	mputi	ing t	he D	NA
12 8 0 0 2	13 10 0 0	Align Position k A: C: G: T:	ment Ma = 1 0 9 0 1	trix: 2 0 10 0 0	3 6 1 0 3	4 10 0 0 0	5 0 5 5
		T:	1	0	3	0	

ethod

$$\frac{5}{5}, \frac{6}{0}, \frac{7}{1}, \frac{8}{5}, \frac{9}{10}, \frac{11}{11}, \frac{12}{13}, \frac{13}{10}, \frac{5}{0}, \frac{1}{0}, \frac{1}{0}, \frac{5}{0}, \frac{1}{0}, \frac{1}{0}, \frac{5}{0}, \frac{1}{0}, \frac{$$

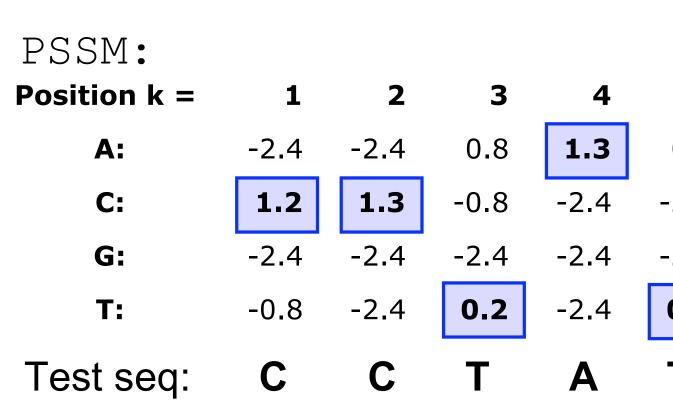
$$\frac{d^2 (Z+1)}{d_j} = \log \left(\frac{9 + 0.25 / 10 + 1}{p_{kj0}} \right) \frac{(E_{kj}^1 + p_j)}{p_j} / (Z+1) = \log \left(\frac{1 + 0.25 / 10 + 1}{0.25} \right) = \log \left(\frac{1 + 0.25 / 10 + 1}{0.25} \right) = -0.8$$

IA Sequence Profile (PSSM)

6	7	8	9	10	11	12	13
0	1	5	0	3	10	8	10
0	0	2	1	1	0	0	0
0	0	1	9	5	0	0	0
10	9	2	0	1	0	2	0

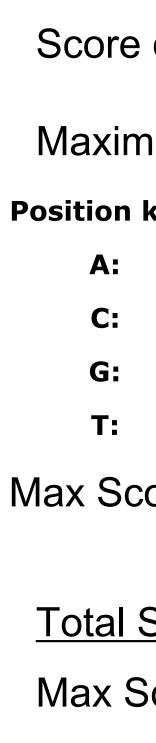
Scoring a test sequence

Query Sequence CCTATTTAGGATA



Query Score = 1.2 + 1.3 + 0.2 + 1.3 + 0.6 + 1.3 + 1.2 + 0.6 + 1.2 + 0.6 + 1.3 + -0.2 + 1.3 = 11.9

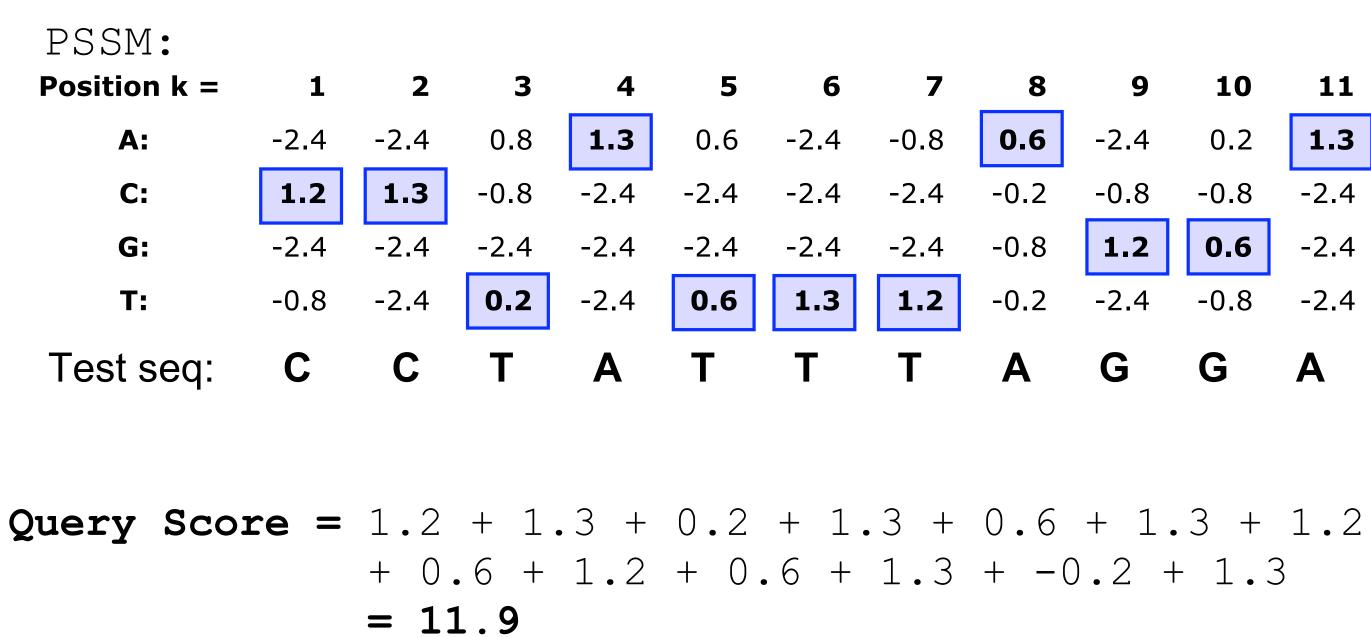
5	6	7	8	9	10	11	12	13
0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
							-2.4	
-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4
0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4
Т	Т	Т	Α	G	G	Α	т	Α



Max S

Scoring a test sequence

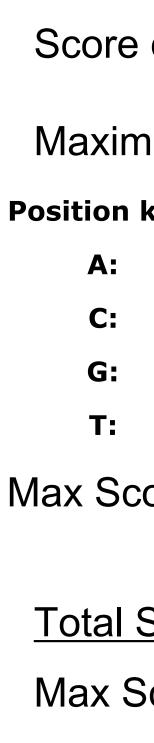
Query Sequence **CCTATTTAGGATA**



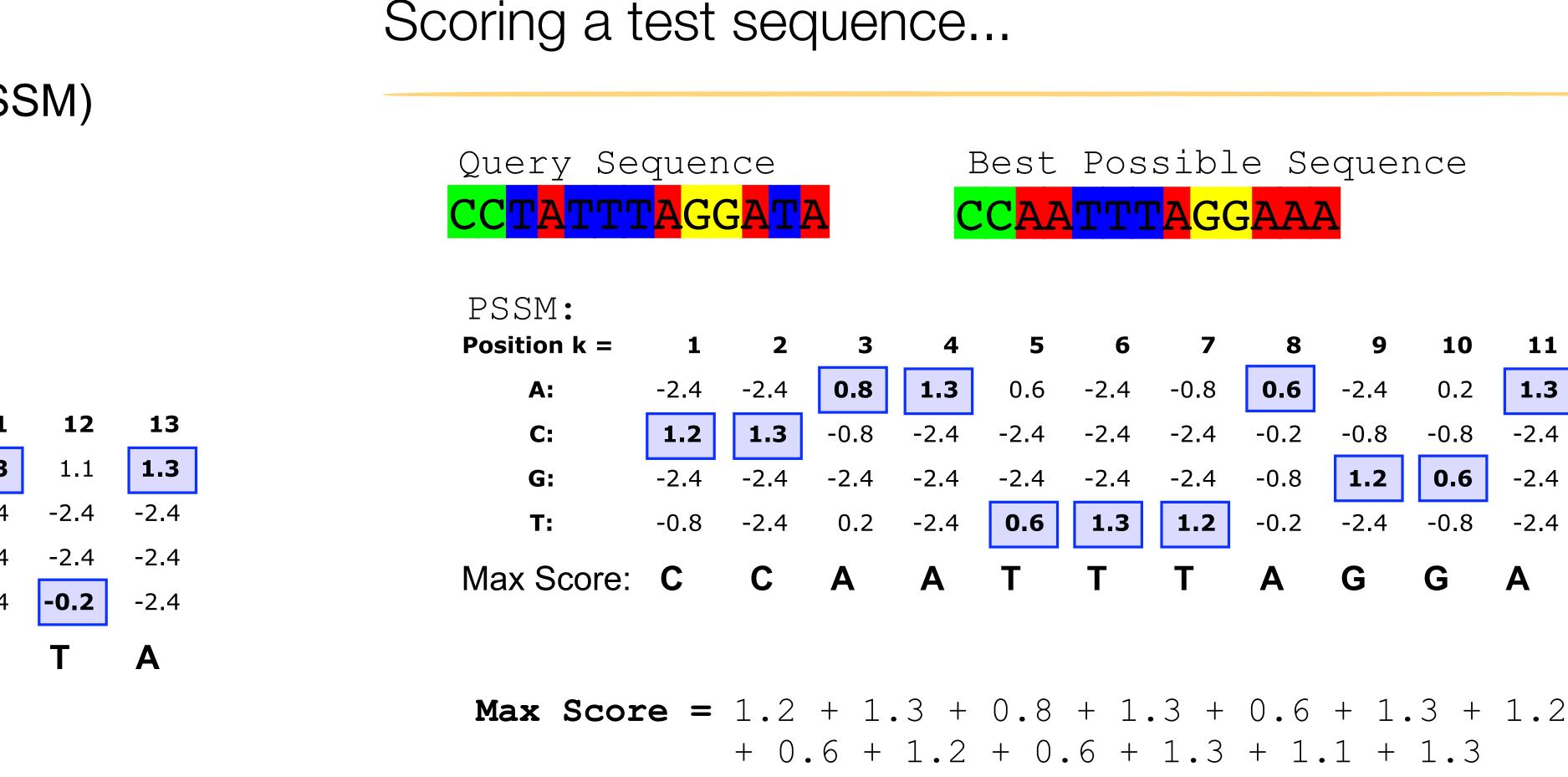
Q. Does the query sequence match the DNA sequence profile?

				9				
0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4
0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4
				G				

+ 0.6 + 1.2 + 0.6 + 1.3 + -0.2 + 1.3



Max S



+ -0.2 + 1.3

A. Following method in Harbison *et al.* (2004) Nature 431:99-104

= 13.8

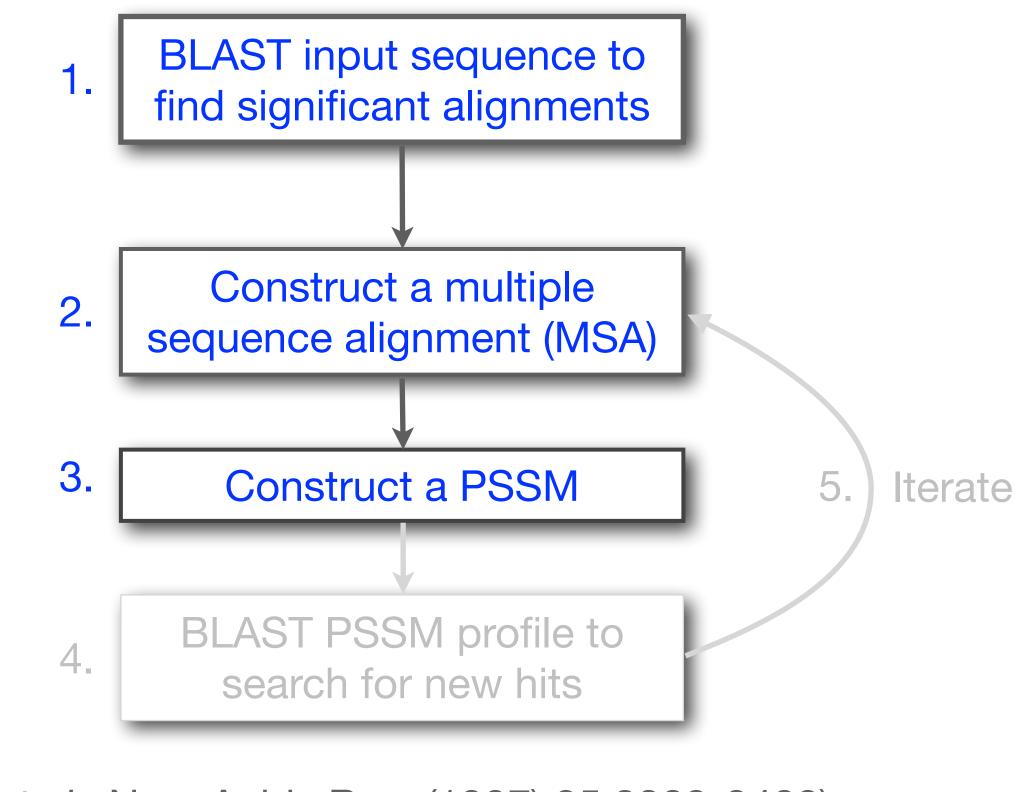
Heuristic threshold for match = $60\% \times Max$ Score = $(0.6 \times 13.8 = 8.28);$ 11.9 > 8.28; Therefore our query is a potential TFBS!

Best Possible Sequence **CCAATTTAGGAAA**

			8					
0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
			-0.2					
-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4
0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4
			Α					
78	+ 1	3 +	06	+ 1	3 +	1 2		

+ 0.6 + 1.2 + 0.6 + 1.3 + 1.1 + 1.3

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

FTVDENGQMSATAKGRVRLFNNU
FSVDEKGHMSATAKGRVRLLSNU
FSVDEKGHMSATAKGRVRLLSNU
MSATAKGRVRLLNNU
FKIEDNGKTTATAKGRVRILDKI
FSVDESGKVTATAHGRVIILNNU
FSVDGSGKVTATAQGRVIILNNU
FTIHEDGAMTATAKGRVIILNNU
FKVEEDGTMTATAIGRVIILNNU
FKVQEDGTMTATATGRVIILNNU
FSVDGSGKMTATAOGRVIILNNI

FSVDGSGKMTATAQGRVIILNN YTVEEDGTMTASSKGRVKLFGF

730496	66
200679	63
206589	34
2136812	2
132408	65
267584	44
267585	44
3777608	63
5687453	60
10697027	
13645517	1
13925316	
131649	65

- WDVCADMIGSFTDTEDPAKFKMKYWGVASFLQKGNDDH 125
- WEVCADMVGTFTDTEDPAKFKMKYWGVASFLQRGNDDH 122
- WEVCADMVGTFTDTEDPAKFKMKYWGVASFLQRGNDDH 93
- WDVCADMVGTFTDTEDPAKFKMKYWGVASFLQKGNDDH 53
- LELCANMVGTFIETNDPAKYRMKYHGALAILERGLDDH 124
- WEMCANMFGTFEDTPDPAKFKMRYWGAASYLQTGNDDH 103
- WEMCANMFGTFEDTPDPAKFKMRYWGAAAYLQSGNDDH 103
- WEMCADMMATFETTPDPAKFRMRYWGAASYLQTGNDDH 122
- WEMCANMFGTFEDTEDPAKFKMKYWGAAAYLQTGYDDH 119
- WEMCANMFGTFEDTEEPARFKMKYWGAAAYLQTGYDDH 140
 - MVGTFTDTEDPAKFKMKYWGVASFLQKGNDDH 32
- FSVDGSGKMTATAQGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAAAYLQSGNDDH 97
- YTVEEDGTMTASSKGRVKLFGFWVICADMAAQYTDPTTPAKMYMTYQGLASYLSSGGDNY 126

N,M,L,Y,G

GLB2_LUMTE/31-141
GLB2_TYLHE/32-143
GLB3_LAMSP/30-141
GLB_TUBTU/29-139
GLB4_LUMTE/36-146
GLB4_TYLHE/33-143
GLB3_TYLHE/33-143
GLB1_TYLHE/30-137
MYG_CYPCA/23-136
MYG_ALLMI/27-143
MYG_GALGA/23-138
MYG_HETPO/23-138
HBAM_LITCT/26-129
HBA3_PLEWA/27-137
HBAZ_CAPHI/27-137
HBB_HETPO/25-136
HBB_SQUAC/25-137
HBB2_XENLA/26-142
HBB1_CYGMA/26-142
HBB_LEPPA/25-142
HBB_ALLMI/25-141
HBB0_MOUSE/26-142
HBBN_AMMLE/20-136
HBB_LITCT/19-135
HBB1_XENBO/25-141
HBA_LEPPA/26-138
HBA1_TORMA/26-136
HBA_HETPO/33-143
HBA_SQUAC/26-136
HBAD_ERYML/26-136

		Е	Е	L	D	H	L	Q	۷	Q	Η	Е				G	R	K	T	Ρ	D
		А	Q	L	Е	н	L	R	Q	Q	н	Т				Κ	L	G	Т	т	G
		т	Q	L	А	н	L	А	S	Q	н	S				S	R	G	۷	s	А
		А	Q	L	А	н	L	к	S	Q	н	А				Е	R	Ν	Т	к	А
		S	н	L	G	н	L	А	D	Q	н	1			Q	R	Κ	G	۷	т	к
		S	L	Т	D	н	L	А	Е	Q	н	к			А	R	А	G	F	к	т
		Е	Е	L	к	н	L	А	R	Q	н	R			_						
		_				-			_			κ		-		_	-	_		-	-
		_							_			А			Ν						
												А			L						
		_										Т		-	т	_			-	-	-
												1								_	_
												Т			Ν						_
												A			Y						
												A			Y			_			
	-											A		_	E				-		
	-	_										Y	·	-	E		_		-	_	-
	-				Q								·		Т			_			
	-											S	·		E						
	•			_	A				-	_			·	-	E		_		-	_	-
•	•				A								·	-	E		-		-	_	-
•		E											·		D						
•						_			_			C		_	B	-	_				
·	•											S		_	G						
•	•	_			_							S	·		E						
•	•				н	_							·		R						
•	-											G		·	K	_		_	-	_	
•	•	_										G		·						D	
·	•											G			T						
·	•	G	1	L	5	Q	L	5	D	L	Η	А	·	·	Y	N	L	R	V	D	Ρ
			[HF	۶Ķ		A	D]			▲ H										

♥

♥

	➡ ➡									
GLB2_LUMTE/31-141	E E L D H L Q V Q H E G R K I P D)								
GLB2_TYLHE/32-143	A Q L E H L R Q Q H I K L G I T G	3								
GLB3_LAMSP/30-141	T Q L A H L A S Q H S S R G V S A	1								
GLB_TUBTU/29-139	A Q L A H L K S Q H A E R N I K A	1								
GLB4_LUMTE/36-146	S H L G H L A D Q <mark>H</mark> I Q R K G V T K	<								
GLB4_TYLHE/33-143	S L I D H L A E Q H K A R A G F K T	Г								
GLB3_TYLHE/33-143	E E L K H L A R Q H R E R S G V K A	٩.								
GLB1_TYLHE/30-137	Q A L A H Y A A F H K Q F G T I F	>								
MYG_CYPCA/23-136	A I L K P L A T T H A N T H K I A L									
MYG_ALLMI/27-143	E V L K P L A K S H A L E H K I P V	/								
MYG_GALGA/23-138	Q P V K A <mark>L</mark> A A T <mark>H</mark> I T T H K I P P	>								
MYG_HETPO/23-138	T N V K E L A D T <mark>H</mark> I N K <u>H</u> K I P P	>								
HBAM_LITCT/26-129		5								
HBA3_PLEWA/27-137	Q A L S K L S D L <mark>H</mark> A Y N L R V D P	>								
HBAZ_CAPHI/27-137	S A L S K L S E L <mark>H</mark> A Y V L R V D P									
HBB_HETPO/25-136	S Q F T D <mark>L</mark> S K K <mark>H</mark> A E E L H V D V									
HBB_SQUAC/25-137	P H F V E L S K K H Y E E L H V D P									
HBB2_XENLA/26-142	S S L Q Q L S K I <mark>H</mark> A T E L F V D P	>								
HBB1_CYGMA/26-142		>								
HBB_LEPPA/25-142	G H L A N L S H L <mark>H</mark> S E K L H V D P	>								
HBB_ALLMI/25-141	G H F A N L S K L <mark>H</mark> C E K F H V D P									
HBB0_MOUSE/26-142	E T F A H L S E L H C D K L H A D F	>								
HBBN_AMMLE/20-136	G A F A S <mark>L</mark> S Z L <mark>H</mark> C B K L H V B P									
HBB_LITCT/19-135	A Y Y A K <mark>L</mark> S E R <mark>H</mark> S G E L H V D P	>								
HBB1_XENBO/25-141	G Y Y A Q L S K Y H S E T L H V D P	>								
HBA_LEPPA/26-138	S C L H T L S E K H A R E L M V D P									
HBA1_TORMA/26-136	H H L N K L A E K H G K G L L V D P									
HBA_HETPO/33-143	THLHKLATF <mark>H</mark> G SELKVDP									
HBA_SQUAC/26-136	G H L D P L A V L H G T T L C V D P									
HBAD_ERYML/26-136	G T L S Q L S D L <mark>H</mark> A Y N L R V D P	>								
	[HPKSAD] H									

20 amino acids

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

. . .

	A R N D C Q	E G H	I L K M	F P S T W Y V
1 M	-1 -2 -2 -3 -2 -1	-2 -3 -2	1 2 -2 6	0 -3 -2 -1 -2 -1 1
2 K	-1 1 0 1 -4 2	4 -2 0	-3 -3 3 -2	-4 -1 0 -1 -3 -2 -3
3 W	-3 -3 -4 -5 -3 -2	-3 -3 -3	-3 -2 -3 -2	1 -4 -3 -3 12 2 -3
4 V	0 -3 -3 -4 -1 -3	-3 -4 -4	3 1 -3 1	-1 -3 -2 0 -3 -1 4
5 W	-3 -3 -4 -5 -3 -2	-3 -3 -3	-3 -2 -3 -2	1 -4 -3 -3 12 2 -3
6 A	5 -2 -2 -2 -1 -1	-1 0 -2	-2 -2 -1 -1	-3 -1 1 0 -3 -2 0
7 L	-2 -2 -4 -4 -1 -2	-3 -4 -3	2 4 -3 2	
8 L	-1 -3 -3 -4 -1 -3		_	
9 L	-1 -3 -4 -4 -1 -2		2 4 -3 2	
10 L	-2 -2 -4 -4 -1 -2		2 4 -3 2	
11 A	5 -2 -2 -2 -1 -1		-2 -2 -1 -1	
12 A	5 -2 -2 -2 -1 -1			
13 W	-2 -3 -4 -4 -2 -2			
14 A				-3 -1 1 -1 -3 -3 -1
15 A	2 -1 0 -1 -2 2		-3 -3 0 -2	
16 A	4 -2 -1 -2 -1 -1	-1 3 -2	-2 -2 -1 -1	-3 -1 1 0 -3 -2 -1
•••		0 0 1		
37 S			-2 -3 0 -2	
38 H				-4 -2 0 -2 -3 -3 -4
39 T			-1 -1 -1 -1	
40 W	-3 -3 -4 -5 -3 -2			
41 H				
42 A	4 -2 -2 -2 -1 -1	-I U -2	-2 -2 -1 -1	-3 -1 1 0 -3 -2 0

20 amino acids

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

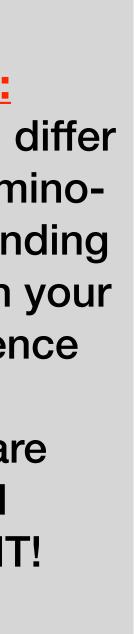
. . .

	A R N	D C Q	E G (H) I	L K M	F P S	T W Y V
1 M	-1 -2 -2	-3 -2 -1	-2 -3 -z 1	2 -2 6	0 -3 -2 -	1 -2 -1 1
2 K	-1 1 0	1 -4 2	4 -2 0 -3	-3 3 -2	-4 -1 0 -	1 -3 -2 -3
3 W	-3 -3 -4	-5 -3 -2	-3 -3 -3 -3	-2 -3 -2	1 -4 -3 -	3 12 2 -3
4 V	0 -3 -3	-4 -1 -3	-3 -4 -4 3			0 -3 -1 4
5 W	-3 -3 -4	-5 -3 -2	-3 -3 -3 -3	-2 -3 -2	1 -4 -3 -	-3 12 2 -3
6 A	5 -2 -2	-2 -1 -1	-1 0 -2 -2	-2 -1 -1	-3 -1 1	0 -3 -2 0
7 L	-2 -2 -4	-4 -1 -2	-3 -4 -3 2	4 -3 2		1 -2 -1 1
8 L	-1 -3 -3	-4 -1 -3	-3 -4 -3 2		3 -3 -2 -	1 -2 0 3
9 L	-1 -3 -4		-3 -4 -3 2	4 -3 2		1 -2 -1 2
10 L	-2 -2 -4					1 -2 -1 1
11 A	J Z Z	-2 -1 -1		-2 -1 -1	•	0 -3 -2 0
12 A	5 -2 -2		-1 0 -2 -2			0 -3 -2 0
13 W		-4 -2 -2				
14 A			-2 4 -2 -2			
15 A	_		0 2 -1 -3			
16 A	4 -2 -1	-2 -1 -1	-1 3 -2 -2	-2 -1 -1	-3 -1 1	0 -3 -2 -1
•••	0 1 0	1 1 0	0 0 1 0	• • • •	0 1 1	1 0 0 0
37 S 38 H	2 - 1 0		$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	-3 0 -2	-3 -1 4	1 - 3 - 2 - 2
38 H 39 T						
	• – •		-1 -2 -2 -1			
40 W 41 H	-3 -3 -4	-5 -3 -2	-3 -3 -3 -3 -3 -2 -3 14 -2	-2 -3 -2	$\begin{array}{cccc} 1 & -4 & -3 & -\\ 2 & 2 & 2 \end{array}$	-3 -3 -3 -3 -3 -3 -3 -3
41 H 42 A			-2 -3 14 -2 -1 0 -2 -2			
42 A	4 - 2 - 2	-7 -1 -1	-1 0 -2 -2	-2 -1 -1	-2 -T T	0 -3 -2 0

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!

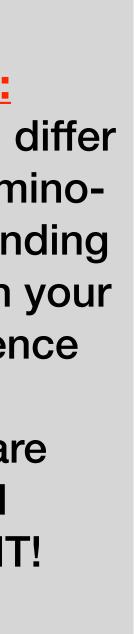


			20 amino acids	
	1 M 2 K 3 W 4 V 5 W 6 A	0 -3 -3 -4 -1 -3	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
All the amino	7 L 8 L 0 T		that is more sensitiv	
acids from position 1 to the end of your query seq.	9 L 10 L 11 A 12 A 13 W 14 A 15 A 16 A 	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
	37 S 38 H 39 T 40 ₩ 41 H 42 A	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$

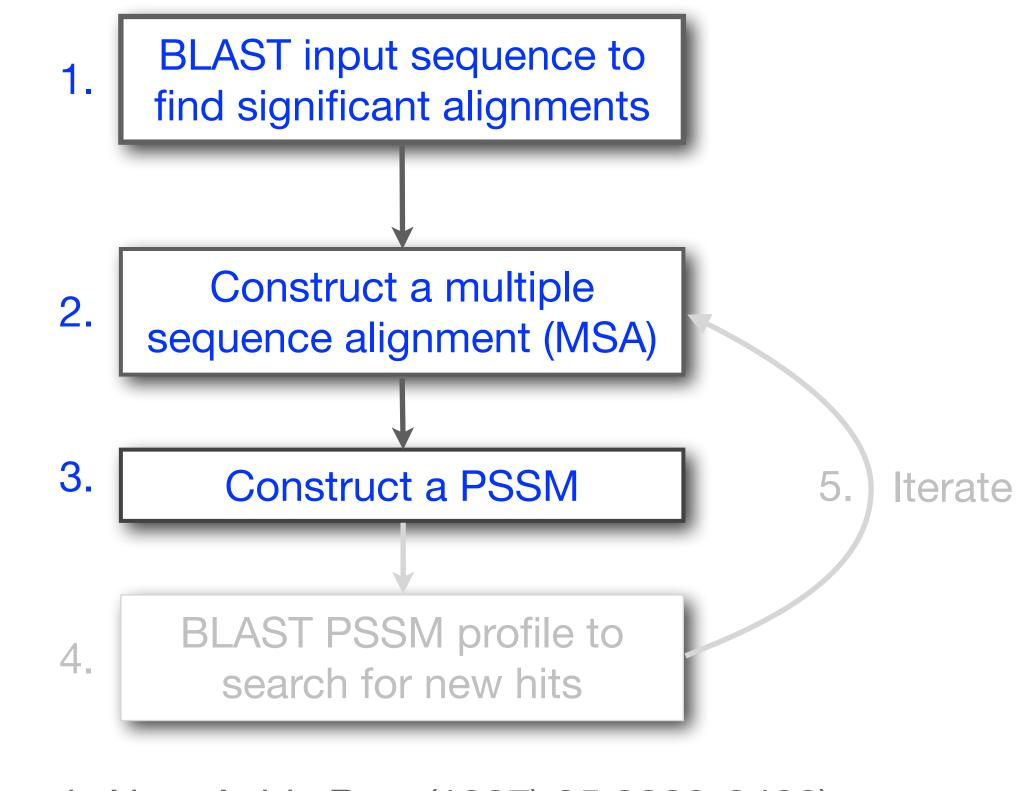
Key Point:

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

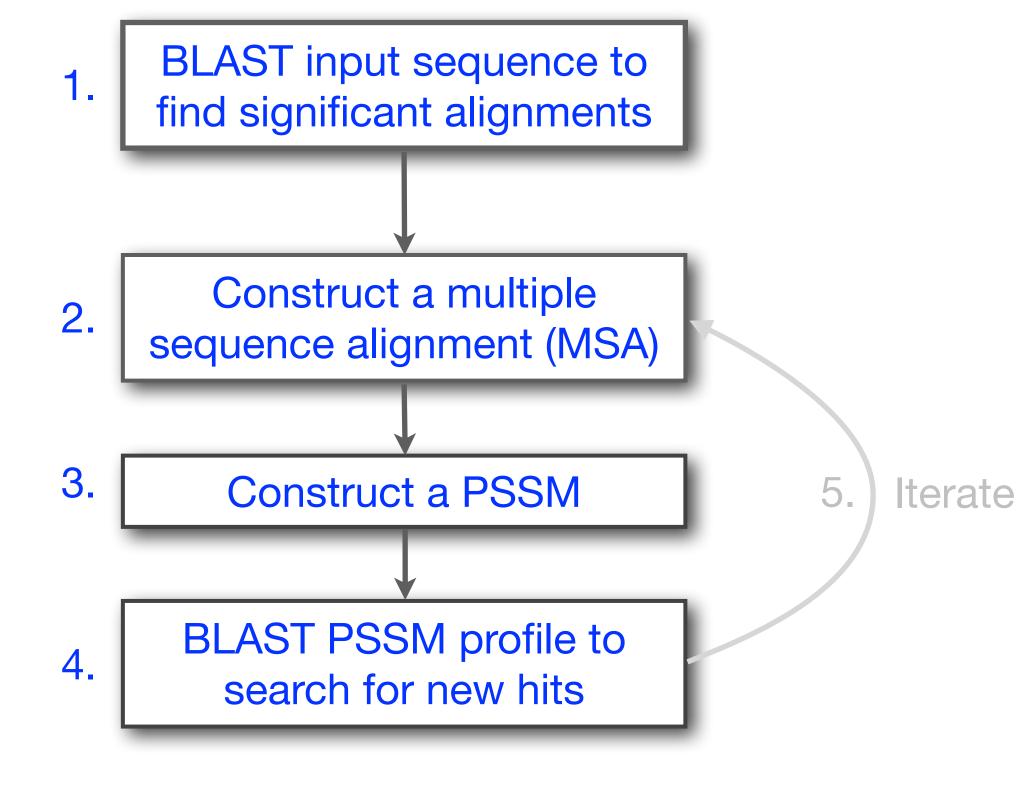
> - i.e. sores are POSITION **DEPENDENT!**



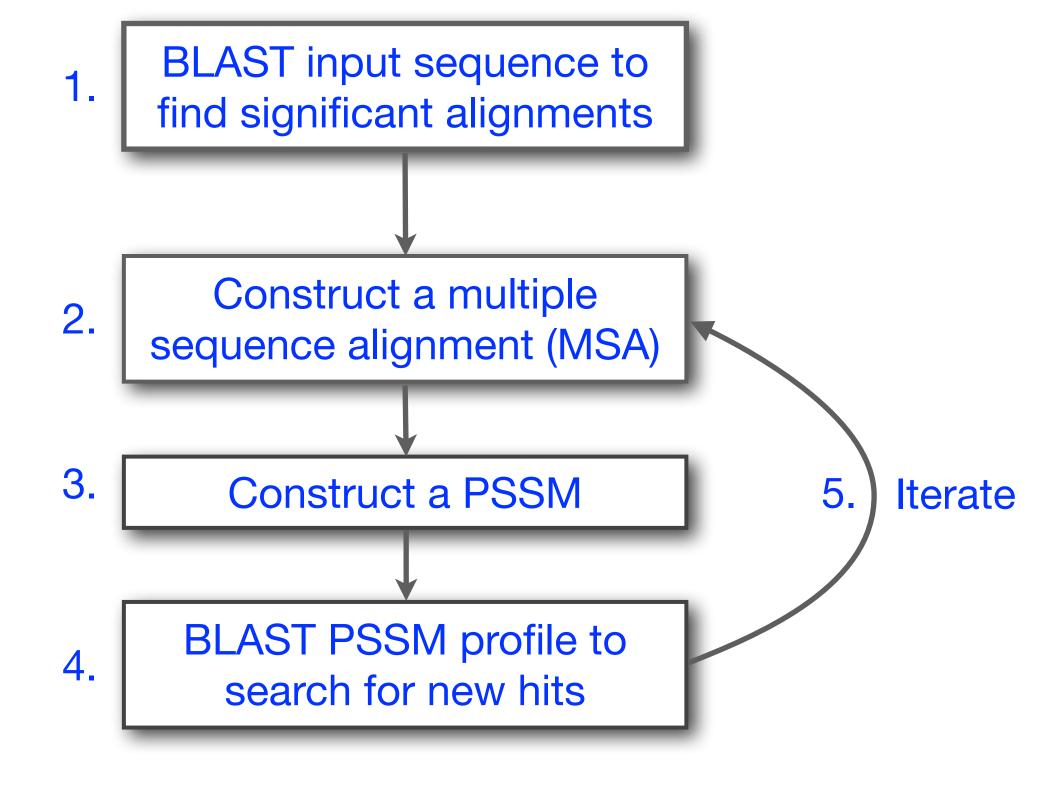
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



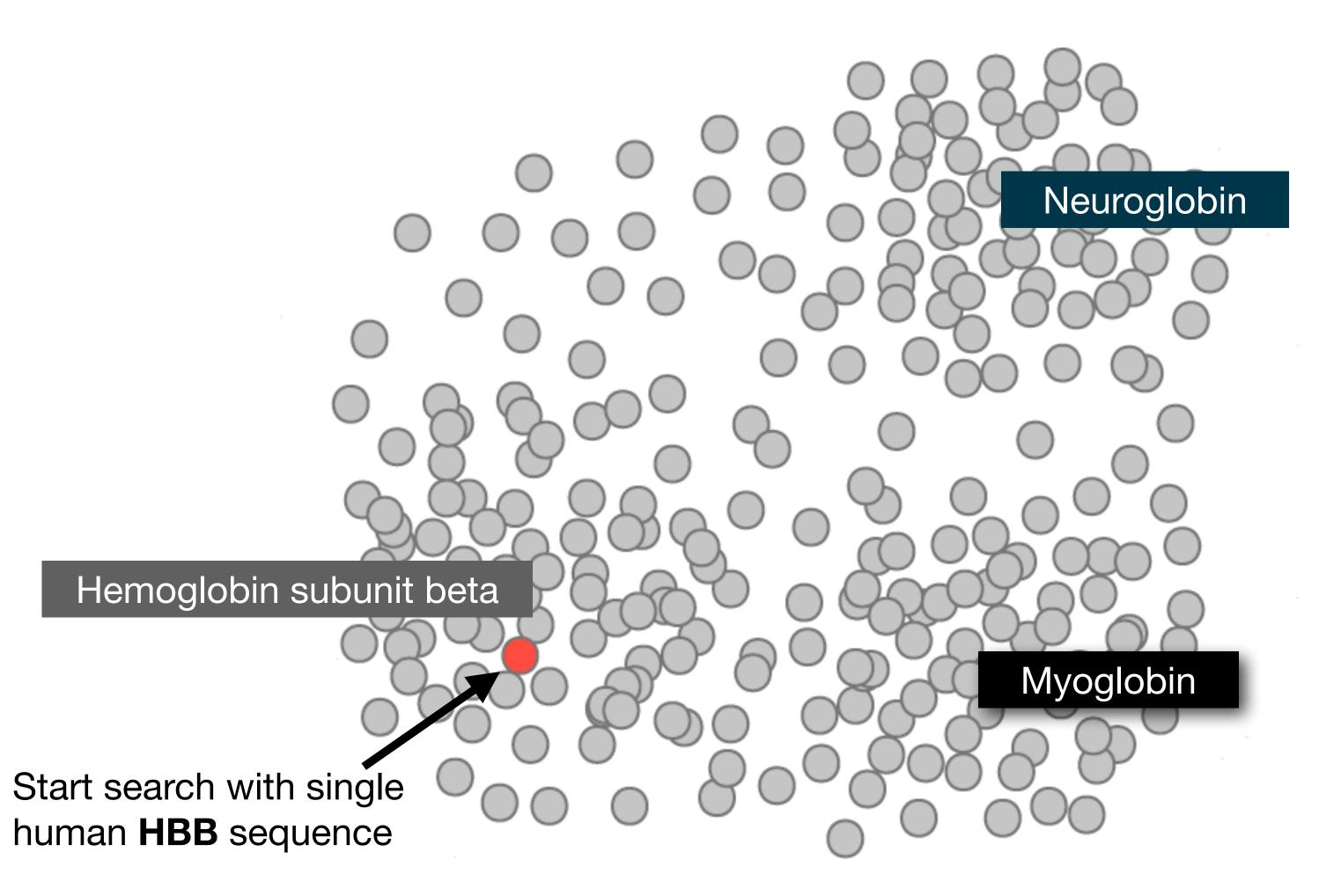
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



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



Hemoglobin subunit beta



 \bigcirc

 \bigcirc

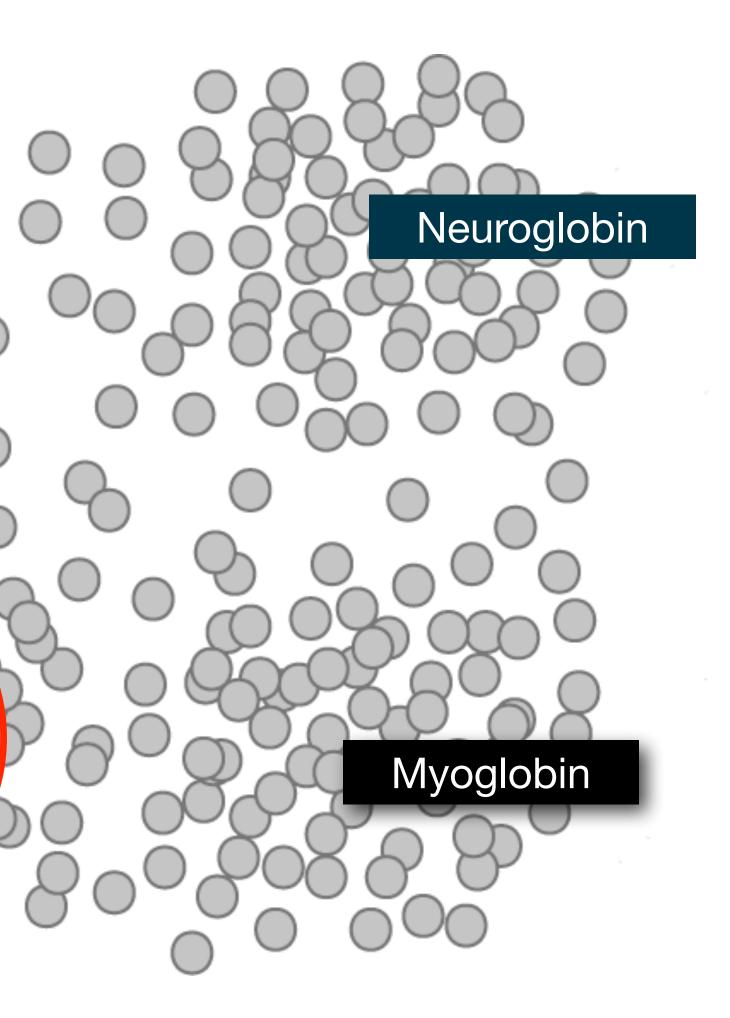
 \bigcirc

Hemoglobin subunit beta

 \bigcirc

200

Result of initial BLASTp search



 \bigcirc

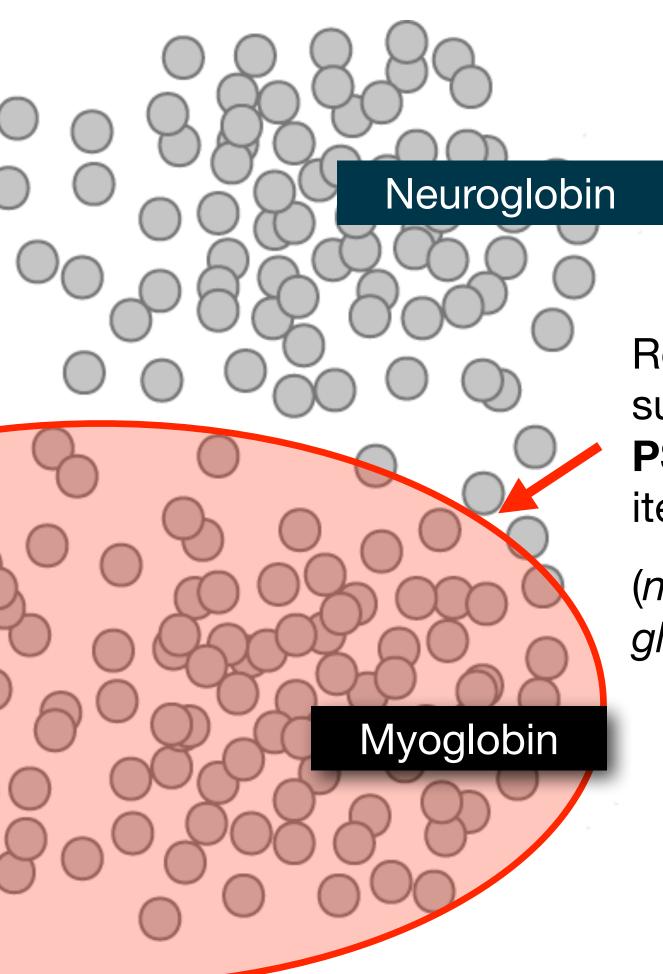
 \bigcirc

Hemoglobin subunit beta

 \bigcirc

 \bigcirc

200



 \bigcirc

Result of subsequent **PSI-BLAST** iteration

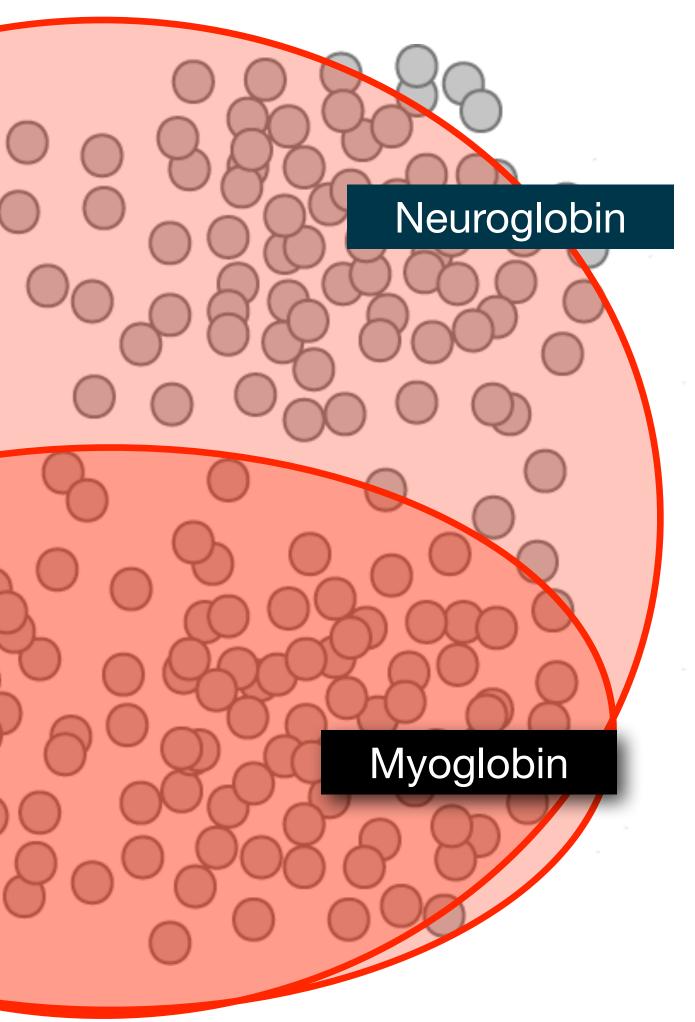
(note, more globin hits!)

TAF6-like RNA polymerase II?

(

 \bigcirc

Hemoglobin subunit beta



Result of later **PSI-BLAST** iteration

(note, potential "corruption"!)

Description

hemoglobin subunit beta [Homo sapiens]
hemoglobin subunit delta [Homo sapiens]
hemoglobin subunit epsilon [Homo sapiens]
hemoglobin subunit gamma-2 [Homo sapiens]
hemoglobin subunit alpha [Homo sapiens]
hemoglobin subunit zeta [Homo sapiens]

Max score	Total score	Query cover	E value	Ident	Accession
301	301	100%	2e-106	100%	<u>NP_000509.1</u>
284	284	100%	7e-100	93%	<u>NP_000510.1</u>
240	240	100%	2e-82	76%	<u>NP_005321.1</u>
235	235	100%	2e-80	73%	<u>NP_000175.1</u>
232	232	100%	3e-79	73%	<u>NP_000550.2</u>
114	114	97%	7e-33	43%	<u>NP_000508.1</u>
100	100	97%	3e-27	36%	<u>NP_005323.1</u>

Description

hemoglobin subunit beta [Homo sapiens] hemoglobin subunit delta [Homo sapiens] hemoglobin subunit epsilon [Homo sapiens] hemoglobin subunit gamma-2 [Homo sapiens] hemoglobin subunit gamma-1 [Homo sapiens] hemoglobin subunit alpha [Homo sapiens] hemoglobin subunit zeta [Homo sapiens]

myoglobin [Homo sapiens]

neuroglobin [Homo sapiens]

Max score	Total score	Query cover	E value	Ident	Accession
301	301	100%	2e-106	100%	<u>NP_000509.1</u>
284	284	100%	7e-100	93%	<u>NP_000510.1</u>
240	240	100%	2e-82	76%	<u>NP_005321.1</u>
235	235	100%	2e-80	73%	<u>NP_000175.1</u>
232	232	100%	3e-79	73%	<u>NP_000550.2</u>
114	114	97%	7e-33	43%	<u>NP_000508.1</u>
100	100	97%	3e-27	36%	<u>NP_005323.1</u>
80.5	80.5	97%	2e-19	26%	<u>NP_005359.1</u>
54.7	54.7	92%	2e-09	23%	<u>NP_067080.1</u>

2

1

New relevant globins found only by PSI-BLAST

Description

hemoglobin subunit beta [Homo sapiens]
hemoglobin subunit delta [Homo sapiens]
hemoglobin subunit epsilon [Homo sapiens]
hemoglobin subunit gamma-2 [Homo sapiens]
hemoglobin subunit alpha [Homo sapiens]
hemoglobin subunit zeta [Homo sapiens]

myoglobin [Homo sapiens]

neuroglobin [Homo sapiens]

myoglobin [Homo sapiens]

hemoglobin subunit alpha [Homo sapiens]

hemoglobin subunit mu [Homo sapiens]

hemoglobin subunit theta-1 [Homo sapiens]

neuroglobin [Homo sapiens]

PREDICTED: cytoglobin isoform X2 [Homo sapiens]

PREDICTED: microtubule cross-linking factor 1 isoform X1 [H

PREDICTED: microtubule cross-linking factor 1 isoform X4 [H

Inclusion of irrelevant hits can lead to PSSM corruption

Ma sco			E value	Ident	Accession	
30	01 301	100%	2e-106	100%	<u>NP_000509.1</u>	
28	284	100%	7e-100	93%	<u>NP_000510.1</u>	
24	0 240	100%	2e-82	76%	<u>NP_005321.1</u>	
23	5 235	100%	2e-80	73%	<u>NP_000175.1</u>	
23	232	100%	3e-79	73%	<u>NP_000550.2</u>	
11	4 114	97%	7e-33	43%	<u>NP_000508.1</u>	
10	00 100	97%	3e-27	36%	<u>NP_005323.1</u>	
8	80.5 80.	.5 97%	2e-19	26%	<u>NP_005359.1</u>	
5	64.7 54.	.7 92%	2e-09	23%	<u>NP_067080.1</u>	
15	59 159	97%	3e-50	26%	NP_005359.1	
15	51 151	97%	3e-47	42%	<u>NP_000508.1</u>	
14	147	97%	6e-46	35%	<u>NP_001003938.1</u>	
14	147	97%	2e-45	37%	<u>NP_005322.1</u>	
13	34 134	92%	3e-40	23%	<u>NP_067080.1</u>	
11	5 115	66%	3e-33	25%	<u>KP_016879605.1</u>	
<u>no sapie</u> 46	.3 46.3	27%	7e-06	39% 2	<u> (P_011523942.1</u>	
<u>no sapie</u> 46	.3 46.3	27%	7e-06	39% >	<u>KP_005258156.1</u>	

?

YOUR TURN!

- There are four required and one optional hands-on sections including:
 - 1. Limits of using BLAST
 - 2. Using PSI-BLAST
 - **3. Examining conservation patterns**

--- BREAK [15 mins]---

- 4. [Optional] Using HMMER
- 5. Divergence of protein sequence and structure
- Please do answer the last review question (Q20).
- We encourage <u>discussion</u> at your Table and on Piazza!

[~10 mins] [~30 mins] [~20 mins]

[~10 mins] [~25 mins]

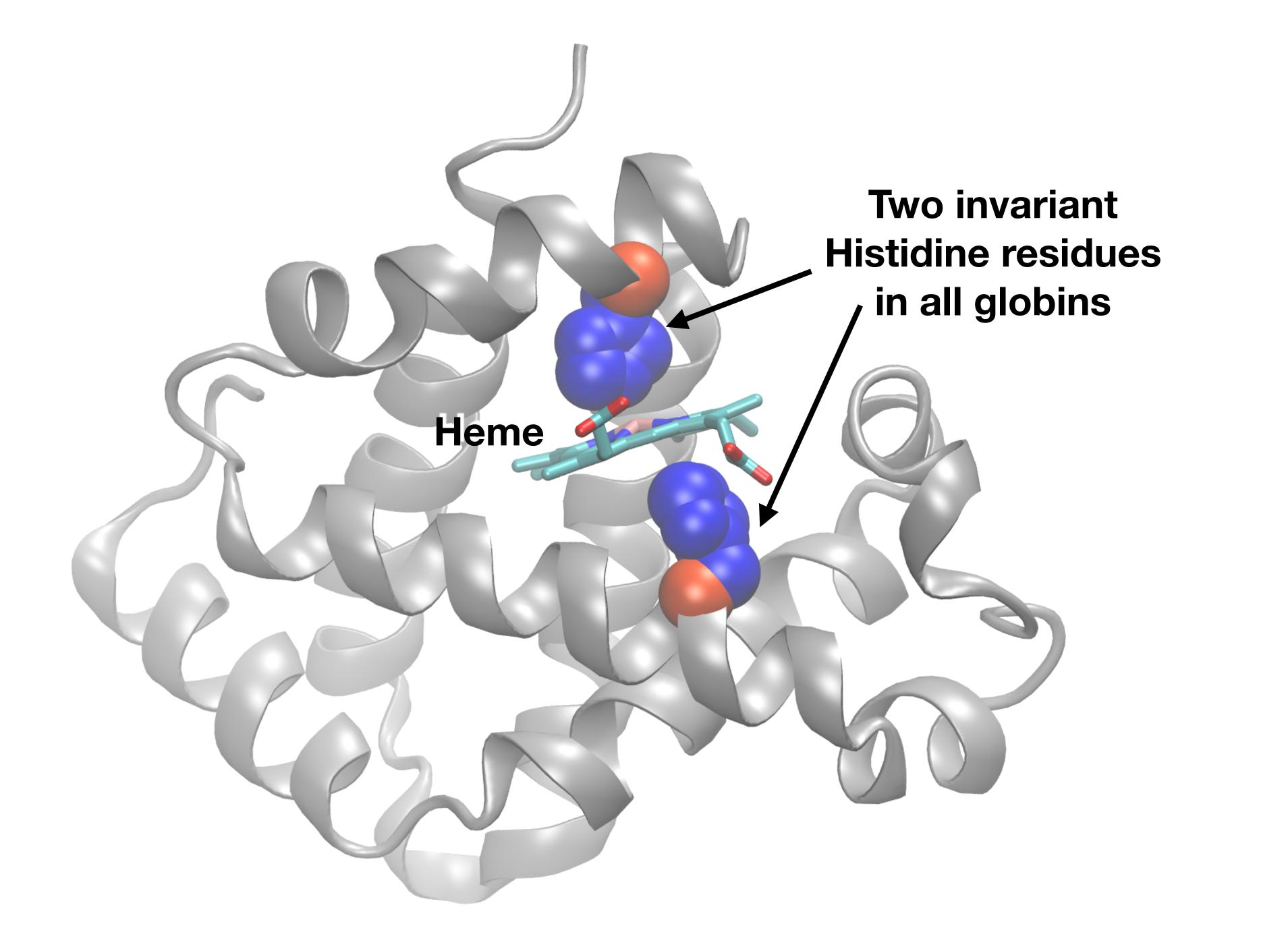
72	MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFE-SFGDLSTPDA VM-GNPKVKAHGKKVLGAF
72	MVHLTPEEKTAVNALWGKVNVDAVGGEALGRLLVVYPWTQRFFE-SFGDLSSPDA VM-GNPKVKAHGKKVLGAF
72	MGHFTEEDKATITSLWGKVNVEDAGGETLGRLLVVYPWTQRFFD-SFGNLSSASA IM-GNPKVKA<mark>H</mark>GKKVLTSL
72	MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFE-SFGDLSTPDA VM-GNPKVKAHGKKVLGAF
72	MVHFTAEEKAAVTSLWSKMNVEEAGGEALGRLLVVYPWTQRFFD-SFGNLSSPSA IL-GNPKVKA<mark>H</mark>GKKVLTSF
72	MGHFTEEDKATITSLWGKVNVEDAGGETLGRLLVVYPWTQRFFD-SFGNLSSASA IM-GNPKVKA<mark>H</mark>GKKVLTSL
67	-MSLTKTERTIIVSMWAKISTQADTIGTETLERLFLSHPQTKTYFP-HFDLHpGSAQLRAHGSKVVAAV
67	-MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFP-HFDLShGSAQVKGHGKKVADAL
89	[15]SEELSEAERKAVQAMWARLYANCEDVGVAILVRFFVNFPSAKQYFS-QFKHMEDPLEME-RSPQLRKHACRVMGAL
66	MLSAQERAQIAQVWDLIAGHEAQFGAELLLRLFTVYPSTKVYFP-HLSACQ-DATQLLSHGQRMLAAV
67	-MALSAEDRALVRALWKKLGSNVGVYTTEALERTFLAFPATKTYFS-HLDLSpGSSQVRAHGQKVADAL
89	[15]SEELSEAERKAVQAMWARLYANCEDVGVAILVRFFVNFPSAKQYFS-QFKHMEDPLEME-RSPQLRKHACRVMGAL
24	MEDPLEME-RSPQLRKHACRVMGAL
73	-MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFD-KFKHLKSEDE MK-ASEDLKK<mark>H</mark>GATVLTAL
72	MERPEPELIRQSWRAVSRSPLEHGTVLFARLFALEPDLLPLFQyNCRQFSSPEDCL-SSPEFLDHIRKVMLVI
18	MK-ASEDLKKHGATVLTAL

	73	SDGLAHLDNLKGTFATLSELHCDKL
	73	SDGLAHLDNLKGTFSQLSELHCDKL
	73	GDAIKHLDDLKGTFAQLSELHCDKL
	73	SDGLAHLDNLKGTFATLSELHCDKL
	73	GDAIKNMDNLKPAFAKLSELHCDKL
	73	GDATKHLDDLKGTFAQLSELHCDKL
	68	GDAVKSIDDIGGALSKLSELHAYIL
	68	TNAVAHVDDMPNALSALSDLHAHKL
1	90	NTVVENLHDPDKVssvLALVGKAHALKH
1	67	GAAVQHVDNLRAALSPLADLHALVL
	68	SLAVERLDDLPHALSALSHLHACQL
	90	NTVVENLHDPDKVssvLALVGKAHALKH
1	25	NTVVENLHDPDKVssvLALVGKAHALKH
1	74	GGILKKKGHHEAEIKPLAQSHATKH
	73	DAAVTNVEDLSSLeeyLASLGRKHRA-V
1	19	GGILKKKGHHEAEIKPLAQSHATKH

✓Query_73613	1
✓ <u>NP_000510.1</u>	1
✓ <u>NP_000175.1</u>	1
✓ <u>NP_000509.1</u>	1
✓ <u>NP_005321.1</u>	1
<u>VP_000550.2</u>	1
✓ <u>NP_005323.1</u>	1
<u>VNP_000508.1</u>	1
✓ <u>XP_005257062.1</u>	1
✓ <u>NP_001003938.1</u>	1
✓ <u>NP_005322.1</u>	1
✓ <u>NP_599030.1</u>	1
✓ <u>XP_016879605.1</u>	1
<u>VNP_001349775.1</u>	1
<u>VNP_067080.1</u>	1
<u>VNP_001369741.1</u>	1

✓Query_73613
✓ <u>NP_000510.1</u>
✓ <u>NP_000175.1</u>
<u>VNP_000509.1</u>
<u>VNP_005321.1</u>
<u>VNP_000550.2</u>
<u>VNP_005323.1</u>
<u>VNP_000508.1</u>
✓ <u>XP_005257062.1</u>
✓ <u>NP_001003938.1</u>
<u>VNP_005322.1</u>
<u>VNP_599030.1</u>
✓ <u>XP_016879605.1</u>
✓ <u>NP_001349775.1</u>
<u>VNP_067080.1</u>
✓ <u>NP_001369741.1</u>

147 LHVDPENFRLLGNVLVCVLAHHFGKE**F**TPPVQAAYQKVVAGVANALAHKYH 147 LHVDPENFRLLGNVLVCVLARNFGKEFTPQMQAAYQKVVAGVANALAHKYH LHVDPENFKLLGNVLVTVLAIHFGKEFTPEVQASWQKMVTGVASALSSRYH 147 147 LHVDPENFRLLGNVLVCVLAHHFGKE**F**TPPVQAAYQKVVAGVANALAHKYH LHVDPENFKLLGNVMVIILATHFGKEFTPEVQAAWQKLVSAVAIALAHKYH 147 147 LHVDPENFKLLGNVLVTVLAIHFGKEFTPEVQASWQKMVTAVASALSSRYH 142 LRVDPVNFKLLSHCLLVTLAARFPADFTAEAHAAWDKFLSVVSSVLTEKYR 142 LRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR 202 HKVEPVYFKILSGVILEVVAEEFASDFPPETQRAWAKLRGLIYSHVTAAYK[35] 141LRVDPANFPLLIQCFHVVLASHLQDEFTVQMQAAWDKFLTGVAVVLTEKYR 142 LRVDPASFQLLGHCLLVTLARHYPGDFSPALQASLDKFLSHVISALVSEYR HKVEPVYFKILSGVILEVVAEEFASDFPPETQRAWAKLRGLIYSHVTAAYK[23] 190 HKVEPVYFKILSGVILEVVAEEFASDFPPETQRAWAKLRGLIYSHVTAAYK[35] 137 HKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYK[6] 154 VGVKLSSFSTVGESLLYMLEKCLGPAFTPATRAAWSQLYGAVVQAMSRGWD[2] 151 HKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYK[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
 - 2. Using PSI-BLAST
 - 3. Examining conservation patterns --- BREAK [15 mins]---

4. [Optional] Using HMMER

- 5. Divergence of protein sequence and structure
- Please do answer the last review question (Q20).
- We encourage <u>discussion</u> at your Table and on Piazza!

[~10 mins] [~30 mins] [~20 mins]

[~10 mins] [~25 mins]

Problems with PSSMs: Positional dependencies

Do not capture positional dependencies



WEIRD WEIQH WEIQH WEIQH

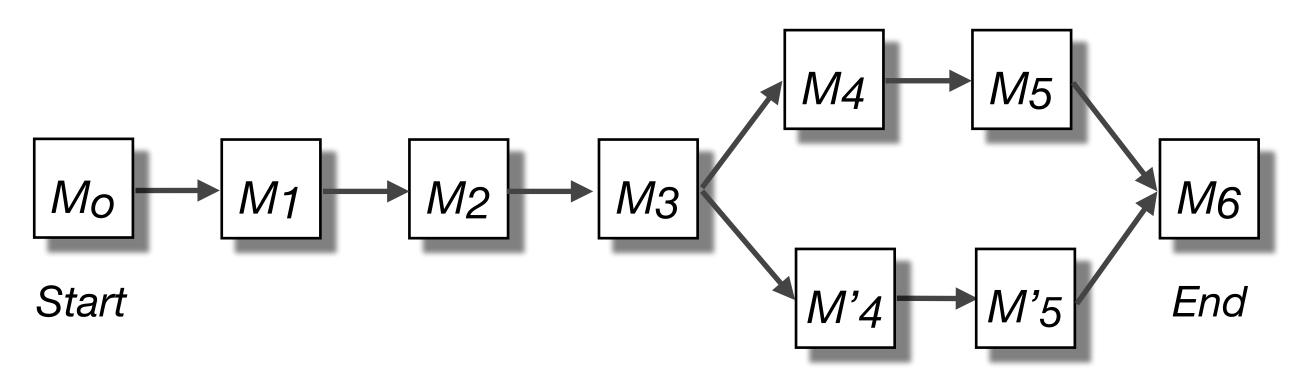
Note: We <u>never</u> see **QD** or **RH**, we only see **RD** and **QH**. However, P(RH)=0.24, P(QD)=0.24, while P(QH)=0.16

		_		
)				0.6
-				
4				0.4
		I		
5			0.4	
2			0.6	
V	Ι			

Markov chains: Positional dependencies

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



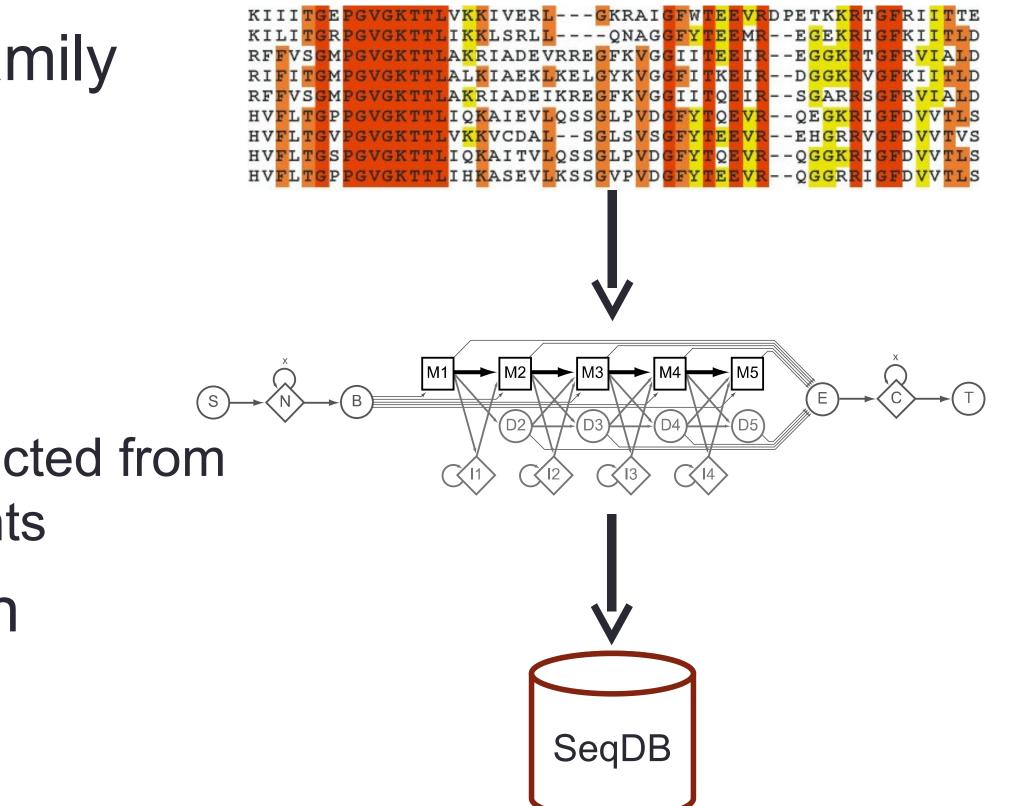


Recall that a PSSM for this motif would give the sequences **WEIRD** and **WEIRH** equally good scores even though the **RH** and **QR** combinations were not observed



Use of HMMER

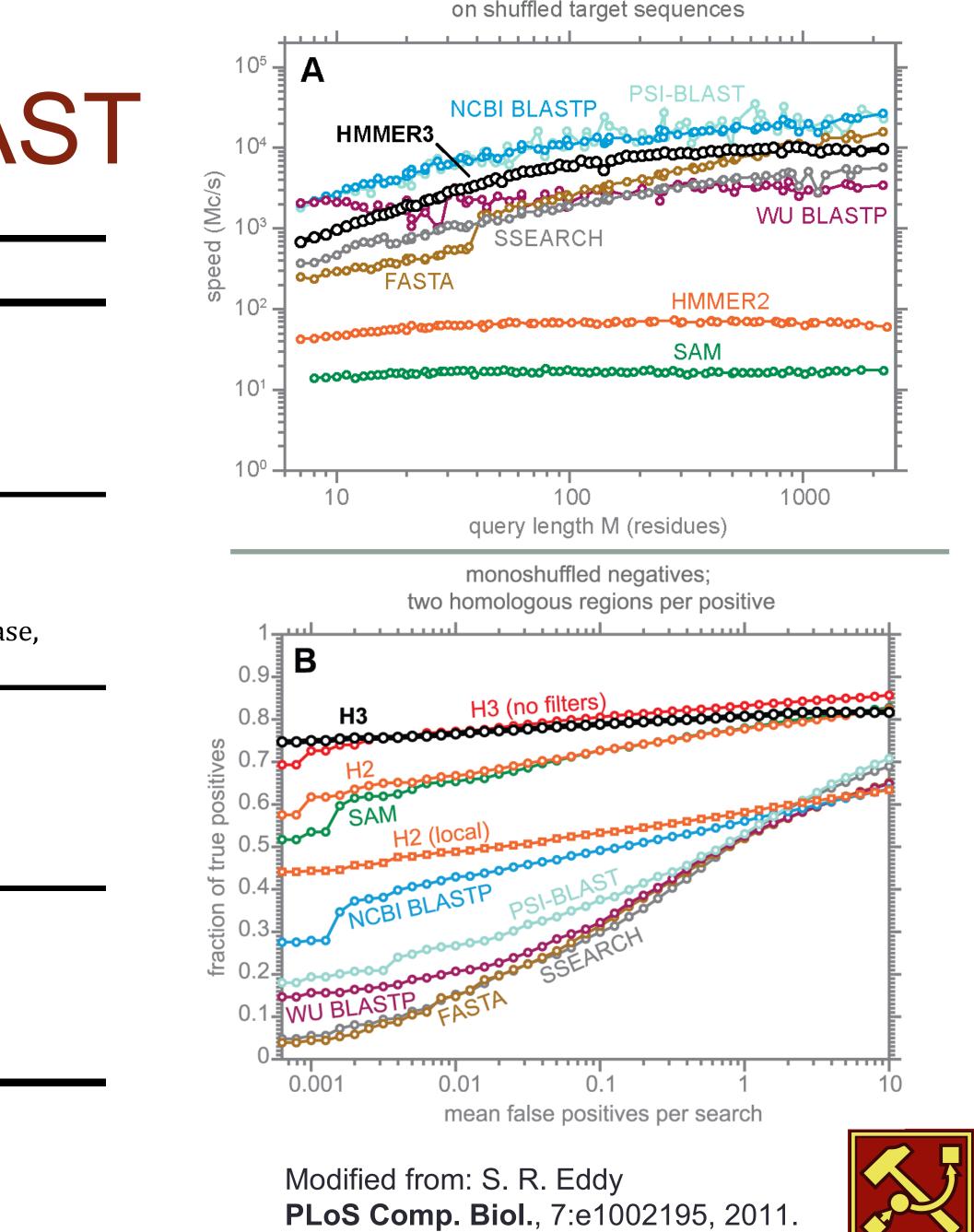
- Widely used by protein family databases
 - Use 'seed' alignments
- Until 2010
 - Computationally expensive
 - Restricted to HMMs constructed from multiple sequence alignments
- Command line application





HMMER vs BLAST

	HMMER	BLAST		
Program	PHMMER	B LA STP		
Query	Single sec	luence		
Target Database	Sequence database			
Program	HMMSCAN	RP SB LA ST		
Query	Single sec	luence		
Target Database	Profile HMM database, e.g. Pfam	PSSM databas e.g. CDD		
Program	HMMSEARCH	P SI-B LA ST		
Query	Profile HMM	PSSM		
Target Database	Sequence d	atabase		
Program	JACKHMMER	P SI-B LA ST		
Query	Single sec	luence		
Target Database	Sequence database			



HMMER Biosequence analysis using profile hidden N Home Search Results Software Help About Contact pmmer hmmscan hmmsearch jackhmmer Crotecin sequence analysis using profile hidden N Paste in your seque Paste a Sequence I Up Paste in your seque NP_000509.1 hemoglobin subunit beta [Homo sapiens] WHLTPEEKSAVTALWGKNNVDEVGGEALGRLUVYPWTORFFESFGDLSTE AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEF ALAHKYH Subm Sequence Database @ Frequently used databases: Reference Proteomes UniProtKB SwissProt Current database selection: SwissProt Restrict by Taxonomy @ Taxon search @ Pre-defined representatives Organism:	EMBL-EBI Services		rch 🎄	Training	i About	us C
phmmer hmmscan hmmsearch jackhmmer protection sequence use protection sequence use Paste a Sequence use Paste in your seque Paste in your seque >NP_000509.1 hemoglobin subunit beta [Homo sapiens] MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTF AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEF Submer Sequence Database @			sis usir	ng prof	ile hidd	en M
Protein sequence vs protein seque Paste a Sequence lup Paste in your seque >NP_000509.1 hemoglobin subunit beta [Homo sapiens] MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTORFFESFGDLSTFA AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFALAHKYH Subm Sequence Database @ Frequently used databases: Reference Proteomes UniProtKB SwissProt * Restrict by Taxonomy @ @ Taxon search	Home Search Results	Software	Help	About	Contact	
Paste a Sequence I Up Paste in your seque >NP_000509.1 hemoglobin subunit beta [Homo sapiens] MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTF AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEF ALAHKYH Sequence Database Frequently used databases: Reference Proteomes UniProtKB SwissProt Current database selection: SwissProt * Restrict by Taxonomy Taxon search Pre-defined representatives 	phmmer hmmscan h	mmsearch	jackhmme	r		
Paste in your seque >NP_000509.1 hemoglobin subunit beta [Homo sapiens] MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTFAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFALAHKYH Subm Sequence Database Frequently used databases: Reference Proteomes UniProtKB SwissProt Restrict by Taxonomy Image: Taxon search Pre-defined representatives	protein sequ	lence	VS	orote	ein se	equ
>NP_000509.1 hemoglobin subunit beta [Homo sapiens] MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTF AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEF ALAHKYH Subm Sequence Database @ Frequently used databases: Reference Proteomes UniProtKB SwissProt V Restrict by Taxonomy @ @ Taxon search Pre-defined representatives				Paste	a Sequenc	e I Uplo
MVHLTPEEKSAVTALWGKVNVDEVGGEÄLGRLLVVYPWTQRFFESFGDLSTF AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEF ALAHKYH Sequence Database @ Frequently used databases: Reference Proteomes UniProtKB SwissProt Current database selection: SwissProt Restrict by Taxonomy @ Taxon search Pre-defined representatives					Paste in you	r sequer
Sequence Database Frequently used databases: Reference Proteomes UniProtKB SwissProt Current database selection: SwissProt Restrict by Taxonomy Taxon search Pre-defined representatives	MVHLTPEEKSAVTALWGK AFSDGLAHLDNLKGTFAT	(VNVDEVGGB	EALGRLL	VVYPWTC		
Current database selection: SwissProt Restrict by Taxonomy Pre-defined representatives	Sequence Database @—					Submit
 Restrict by Taxonomy @ Taxon search Pre-defined representatives 		S: Reference F	Proteomes	UniProt	KB Swiss	Prot
Taxon search Pre-defined representatives	SwissProt			•		
	▼ Restrict by Taxonomy	0				
Organism:	Taxon search Pre	e-defined repres	entatives			
	Organism:					

en Markov Models

equence database

e I Upload a File I Accession Search

r sequence or use the example 📀

DLSTPDAVMGNPKVKAHGKKVLG GKEFTPPVQAAYQKVVAGVAN

Submit Reset

Signif	Significant Query Matches (12) in swissprot (v.2018_11) Customise Customise							
	Target	Description	Species	Oross-references	E-value			
>	HBB_HUMAN	Hemoglobin subunit beta	Homo sapiens 🗗		6.8e-99			
>	HBD_HUMAN ₪	Hemoglobin subunit delta	Homo sapiens 🗗		1.6e-91			
>	HBE_HUMAN ₪	Hemoglobin subunit epsilon	Homo sapiens 🗗		1.5e-74			
>	HBG2_HUMANI	Hemoglobin subunit gamma-2	Homo sapiens 🗗		8.8e-73			
>	HBG1_HUMAN☞	Hemoglobin subunit gamma-1	Homo sapiens 🗗		6.2e-72			
>	HBA_HUMAN ₪	Hemoglobin subunit alpha	Homo sapiens 🗗		3.8e-29			
>	HBAZ_HUMAN፼	Hemoglobin subunit zeta	Homo sapiens 🗗		4.5e-23			
>	HBAT_HUMAN ₪	Hemoglobin subunit theta-1	Homo sapiens 🗗		5.2e-22			
>	HBM_HUMAN ₪	Hemoglobin subunit mu	Homo sapiens 🗗		3.4e-19			
>	CYGB_HUMAN ₪	Cytoglobin	Homo sapiens 🗗		3.1e-14			
>	MYG_HUMAN ₪	Myoglobin	Homo sapiens 🗗		2.3e-06			
>	NGB_HUMAN፼	Neuroglobin	Homo sapiens 🗗		0.0017			
(show	(show all) alignments Your search took: 0.06 secs showing rows 1 - 12 of 12							

<u>Local Link</u>

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:

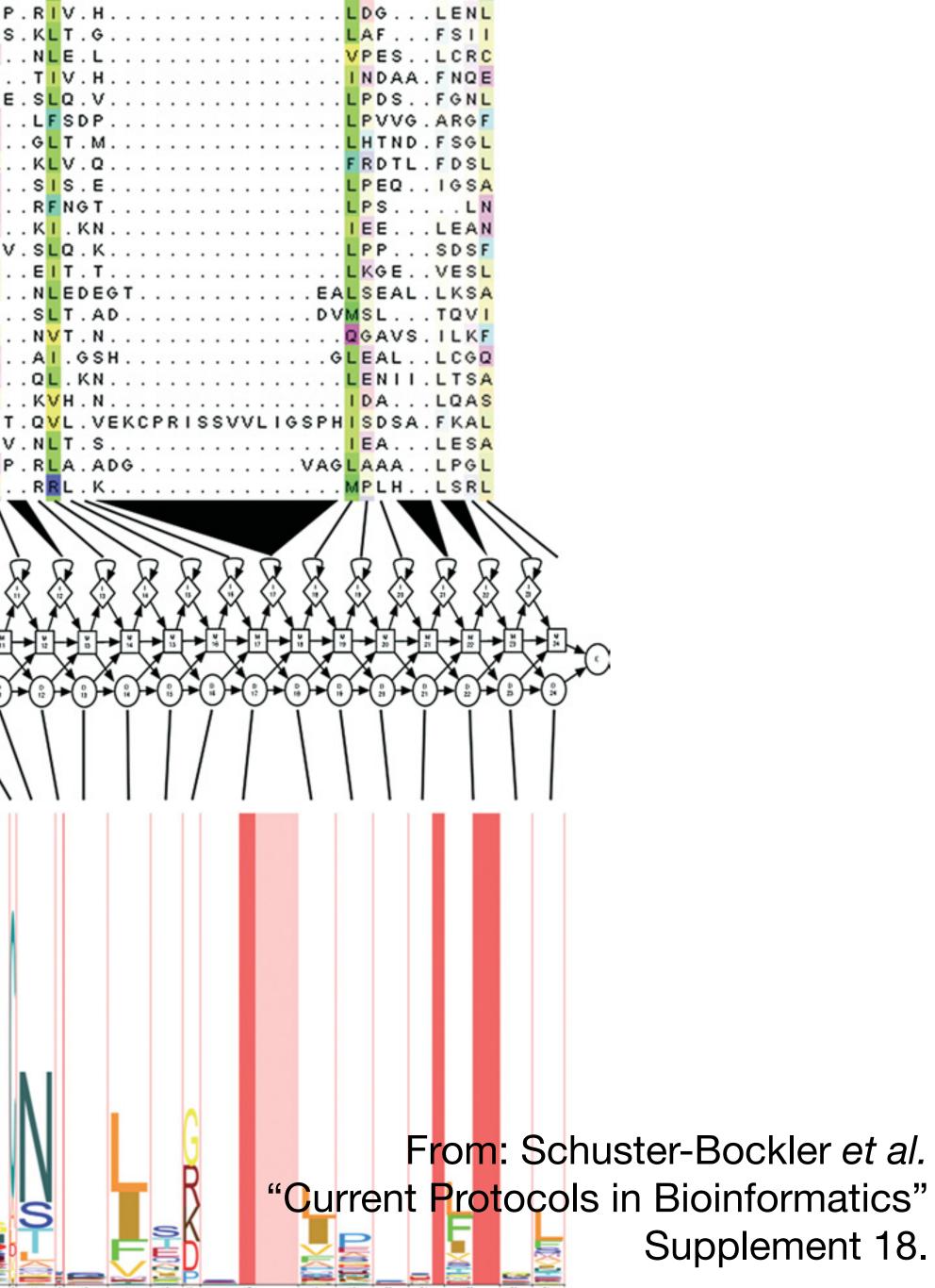
- lacksquaresearches.
- •
- •

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

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**

Q9M8N0_ARATH/320-341 FLII_HUMAN/318-339 Q9VN74_DROME/90-112 Q8L8/7_PINTA/792-814 Q9FHL8_ARATH/301-324 SLIK6_MOUSE/65-87 Q8NJJ8_EMENI/978-1000 Q9LUQ2_ARATH/92-113 Q9FH93_ARATH/169-188 Q898G0_CLOTE/268-288 Q8H6V2_MAIZE/678-699 Q9AR40_LINUS/692-713 Q9LE82_ARATH/350-377 Q9H5N5_HUMAN/255-278 Q8L4C7_ARATH/185-207 Q9VSA4_DROME/1115-1138 TLR1_MOUSE/376-398 Q9TXJ6_LEIMA/445-465 FXL13_MOUSE/409-448 Q9TXJ6_LEIMA/927-948 Q9M4X9_CHLRE/1417-1444 Q945S6_LYCPIN/656-677	RLKTLSLQKN. GLRDIDLSHT. KLIYLDLSGCT ALTVVNANSCV
Retative Entropy	



YOUR TURN!

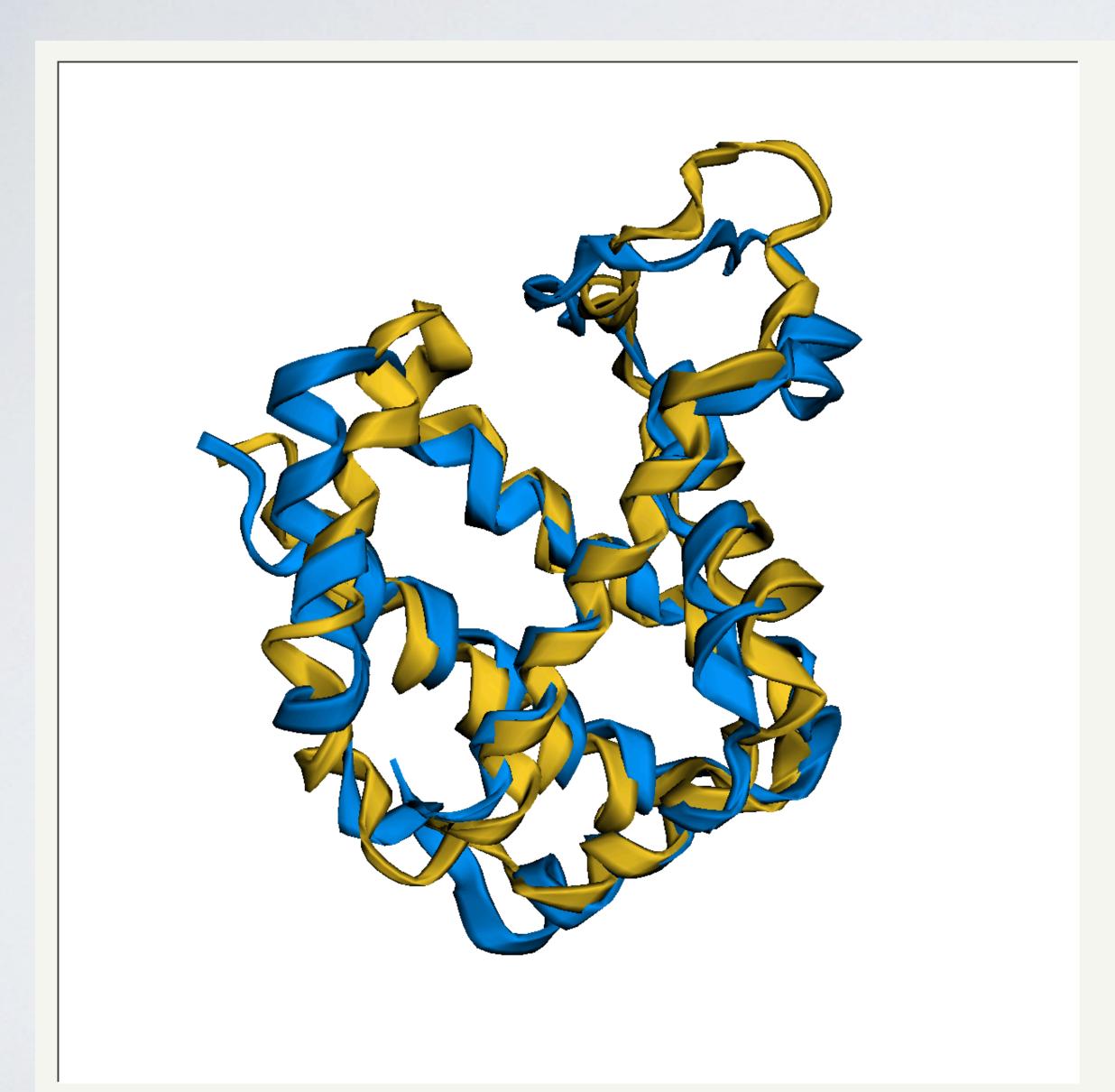
- There are four required and one optional hands-on sections including:
 - 1. Limits of using BLAST
 - 2. Using PSI-BLAST
 - 3. Examining conservation patterns

--- BREAK [15 mins]---

- 4. [Optional] Using HMMER
- 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!

[~10 mins] [~30 mins] [~20 mins]

[~10 mins]



ALIGNM	IENT	CONTACT MAP						
Align 2	nbsB.p	db 146 with 4m	pmB.pdb 148					
Twists () ini-	len 136 ini-rm	sd 3.05 opt-eq	u 143 opt-r	msd 2.65 c	hain-rmsd	3.05	
Score 31	L8.72	align–len 150	gaps 7 (4.67%)					
P-value	3 . 26e	-14 Afp-num 14	073 Identity 2	0.67% Simil	arity 40.0	0%		
Block () afp	17 score 318.7	2 rmsd 3.05 g	ap 9 (0.06%)			
		. :	. : .	: .	: .	: .	: .	:
Chain 1	2	HLTPVEKSAVTAL	WGKVNVDEVGGE	ALGRLLVVYPW	TQRFFESFG-	DLSTPDAVMG	NPKVKAHG	KKVL
		111 11111111	11111 1111111	11111111111	111111111	1111111111	11111111	.1111
Chain 2	2	ERPEPELIRQS	WRAVSRSPLEHGTV	LFARLFALEPD		OFSSPEDCLS	SPEFLDHI	RKVM
Chain 1	69	GAFSDGLAHLDNL	KGTEATL SELHCD-	-KI HVDPENER		AHHEGKEETP	PVOAAYOK	VVAG
			11111111111111					
Chain 2	70		SSLEEYLASLGRKH					
	70		SSEETERSEONAT		I VOLULE IIIL			
Chain 1	127	VANALAHKYH						
	137							
		1111111111						
Chain 2	140	VVQAMSRGWD						



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 <u>NOW</u> if you don't have access to DataCamp!