



BIMM 143

Advanced Database Searching
Lecture 3

Barry Grant
UC San Diego

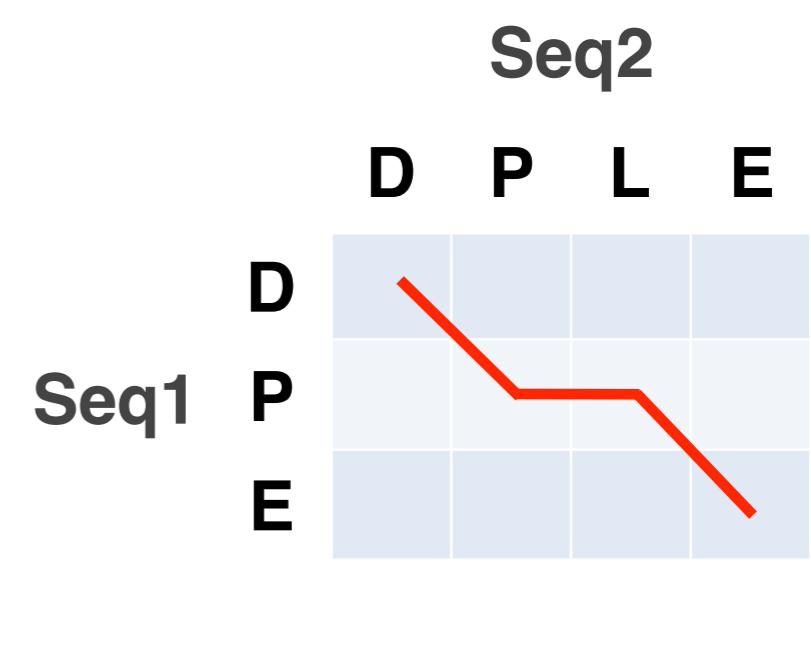
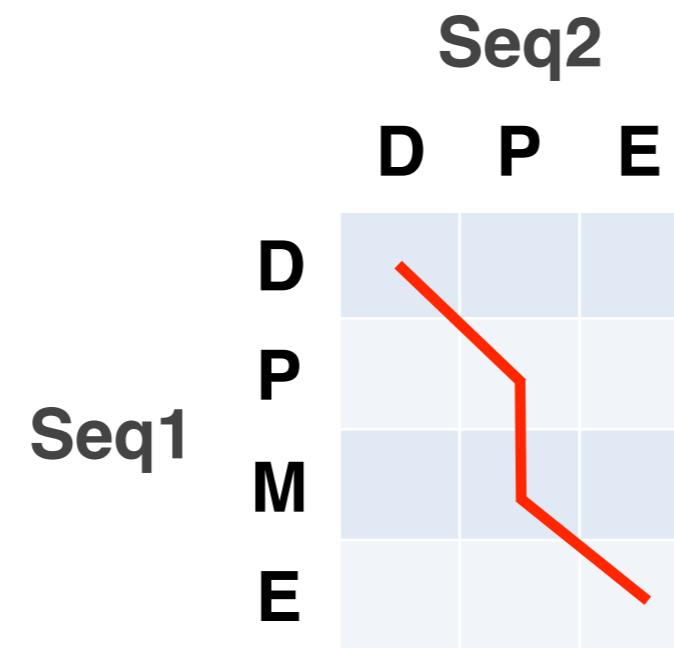
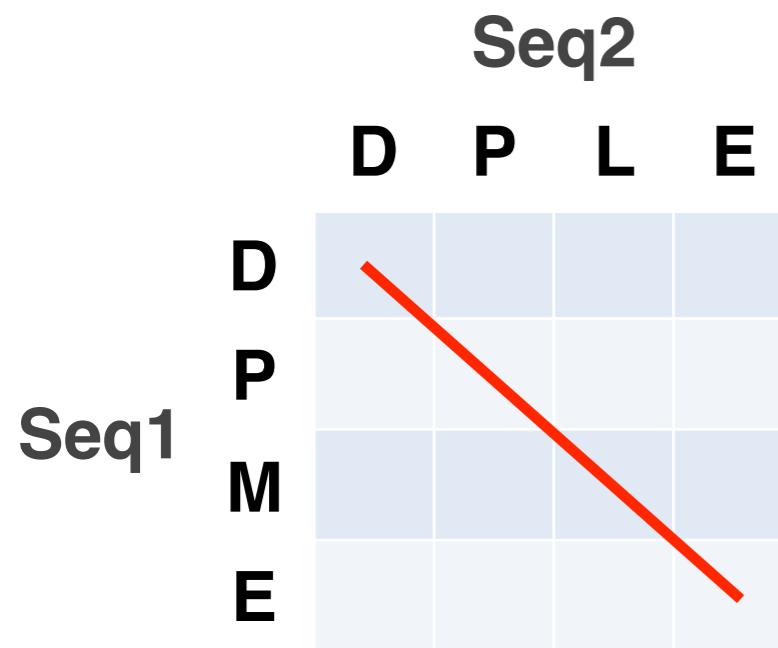
<http://thegrantlab.org/bimm143>

Recap From Last Time:

- Sequence alignment is a fundamental operation underlying much of bioinformatics.
- Introduced dot matrices, dynamic programming and the BLAST heuristic approaches.
 - ➔ *Key point:* Even when optimal solutions can be obtained they are not necessarily unique or reflective of the biologically correct alignment.
- Introduced classic global and local alignment algorithms (Needleman–Wunsch and Smith–Waterman) and their major application areas.
- Heuristic approaches are necessary for large database searches and many genomic applications.

[Feedback](#)

Muddy Point: Different paths represent different alignments



Seq1: D P L E
| | : |
Seq2: D P M E

Seq1: D P M E
| | | |
Seq2: D P - E

Seq1: D P - E
| | | |
Seq2: D P L E

(Mis)matches are represented by diagonal paths &
Indels with horizontal or vertical path segments

Todays Menu

- Sequence motifs and patterns: Simple approaches for finding functional cues from conservation patterns
- Sequence profiles and position specific scoring matrices (PSSMs): Building and searching with profiles, Their advantages and limitations
- PSI-BLAST algorithm: Application of iterative PSSM searching to improve BLAST sensitivity
- Hidden Markov models (HMMs): More versatile probabilistic model for detection of remote similarities

Side Note:

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

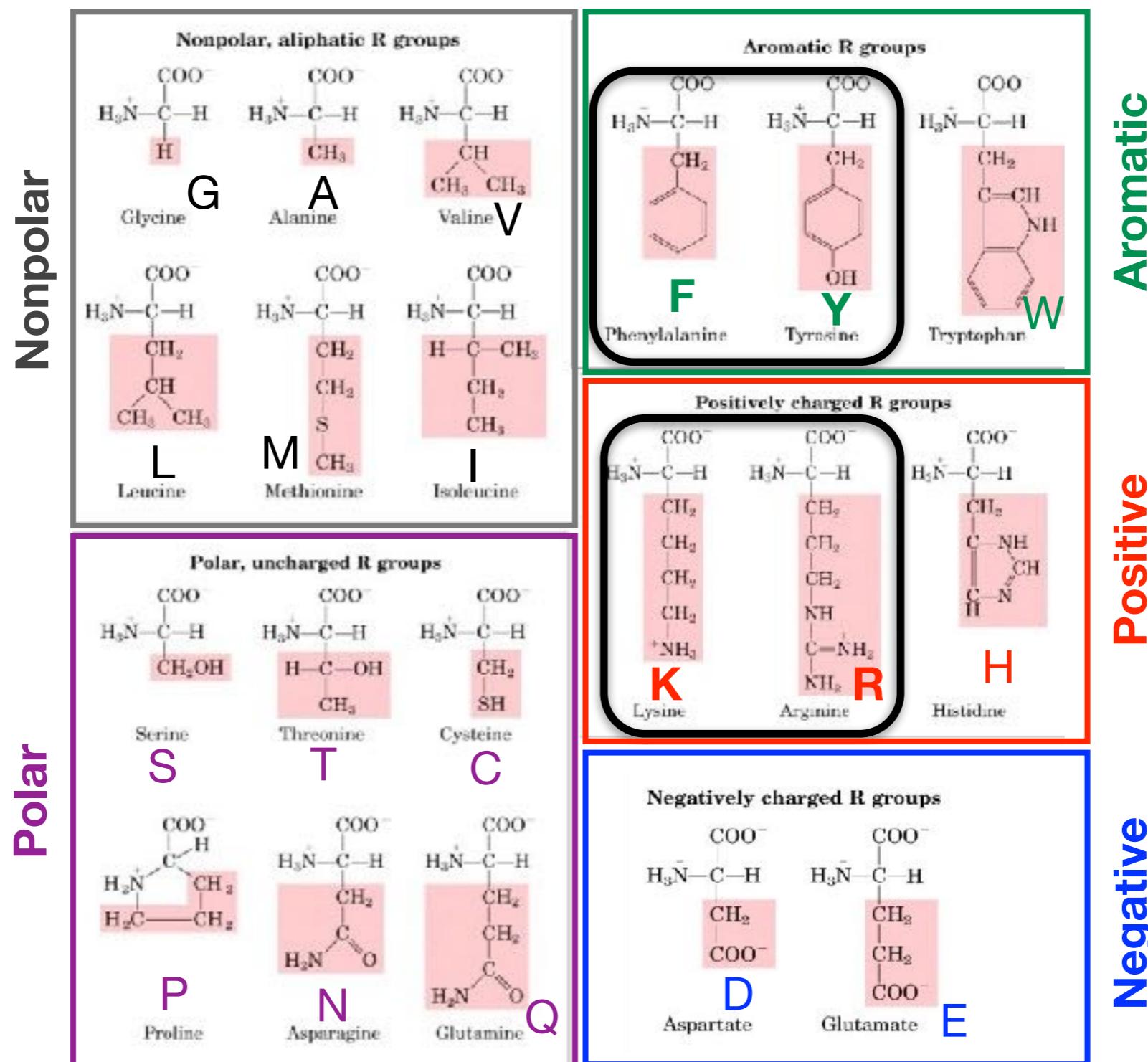
By default BLASTp match scores come from the BLOSUM62 matrix

Blocks Substitution Matrix. Scores obtained from observed frequencies of substitutions in blocks of aligned sequences with no more than 62% identity.

Note. Some amino acid mismatches have positive scores (highlighted in red) reflecting the shared physicochemical properties of these amino acids

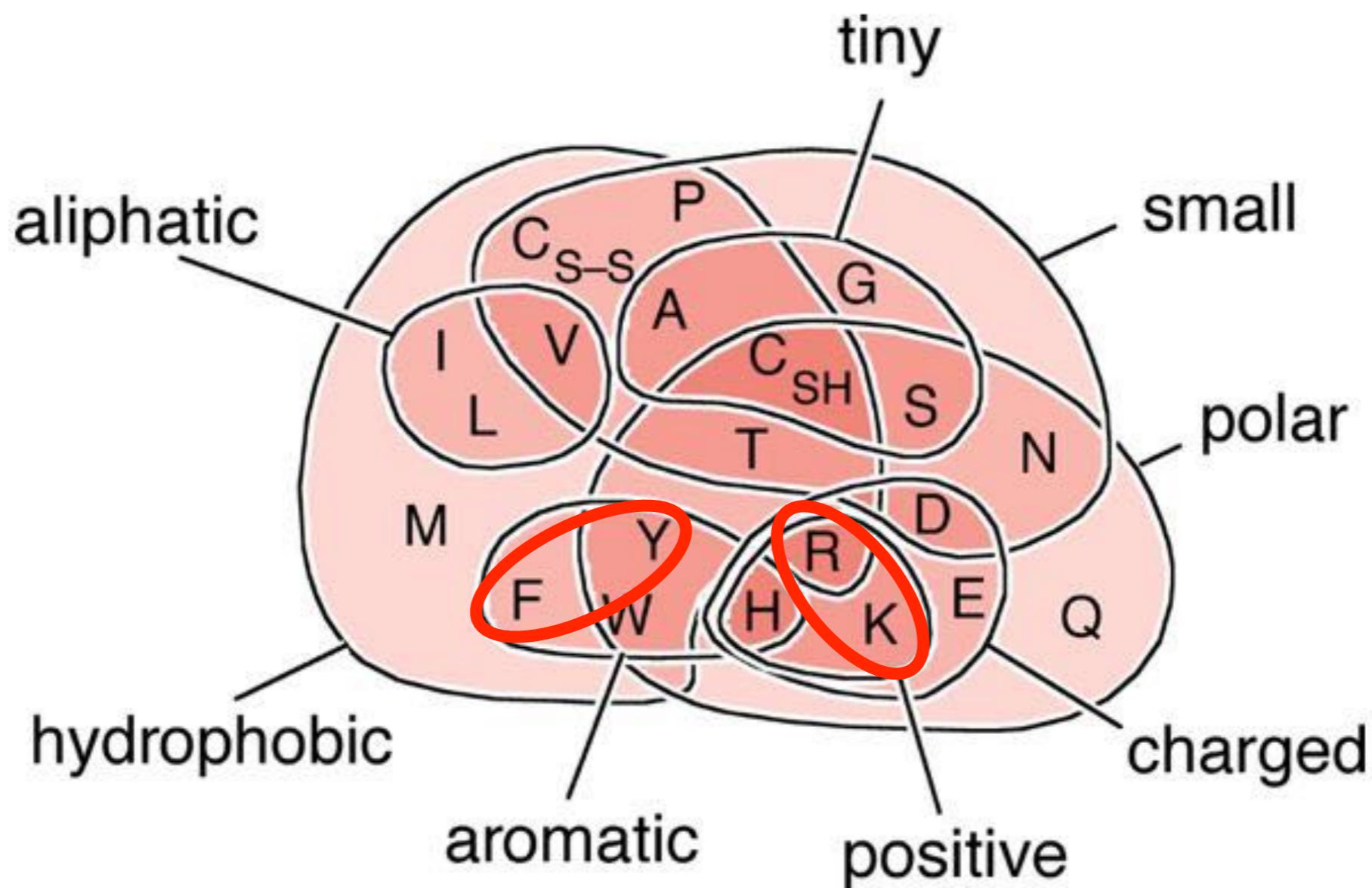
Not all matches score equally
(blue highlighted values)

Protein scoring matrices reflect the properties of amino acids



Protein scoring matrices reflect the properties of amino acids

"from Louis Taylor"



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...

Todays Menu

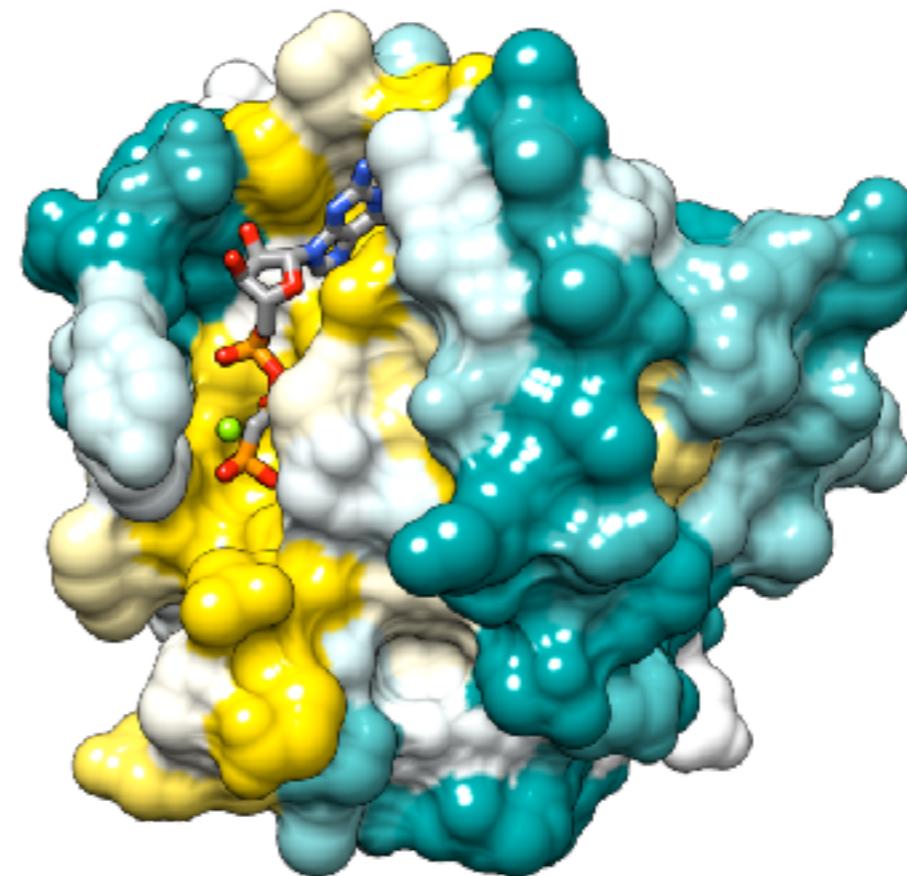
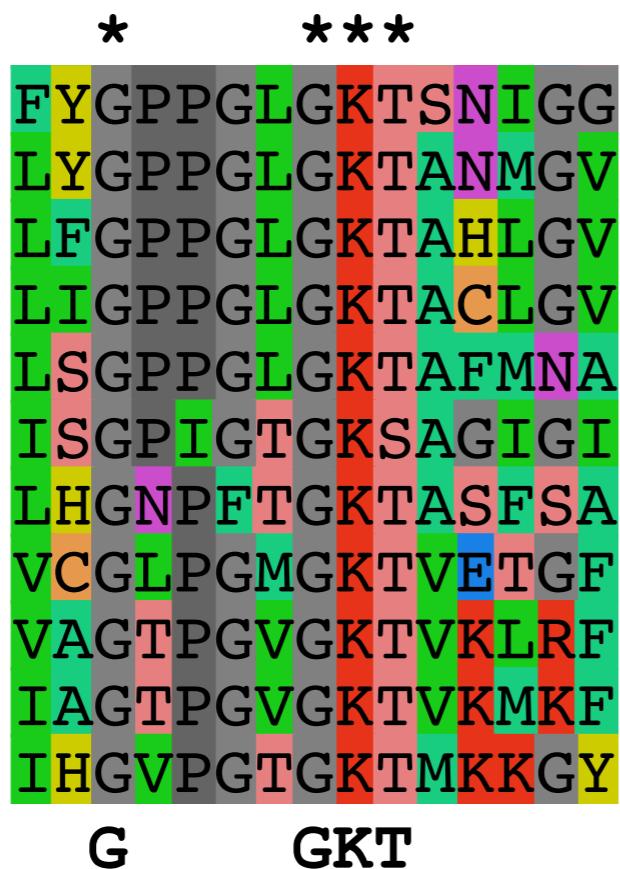
- Sequence motifs and patterns: Simple approaches for finding functional cues from conservation patterns
- Sequence profiles and position specific scoring matrices (PSSMs): Building and searching with profiles, Their advantages and limitations
- PSI-BLAST algorithm: Application of iterative PSSM searching to improve BLAST sensitivity
- Hidden Markov models (HMMs): More versatile probabilistic model for detection of remote similarities

Functional cues from conservation patterns

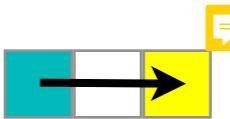
Within a protein or nucleic acid sequence there may be a small number of characteristic residues that occur consistently. These conserved “sequence fingerprints” (or **motifs**) usually contain functionally important elements

- E.g., the amino acids that are consistently found at enzyme active sites or the nucleotides that are associated with transcription factor binding sites.

ATP/GTP-binding proteins: G-x(4)-G-K-T



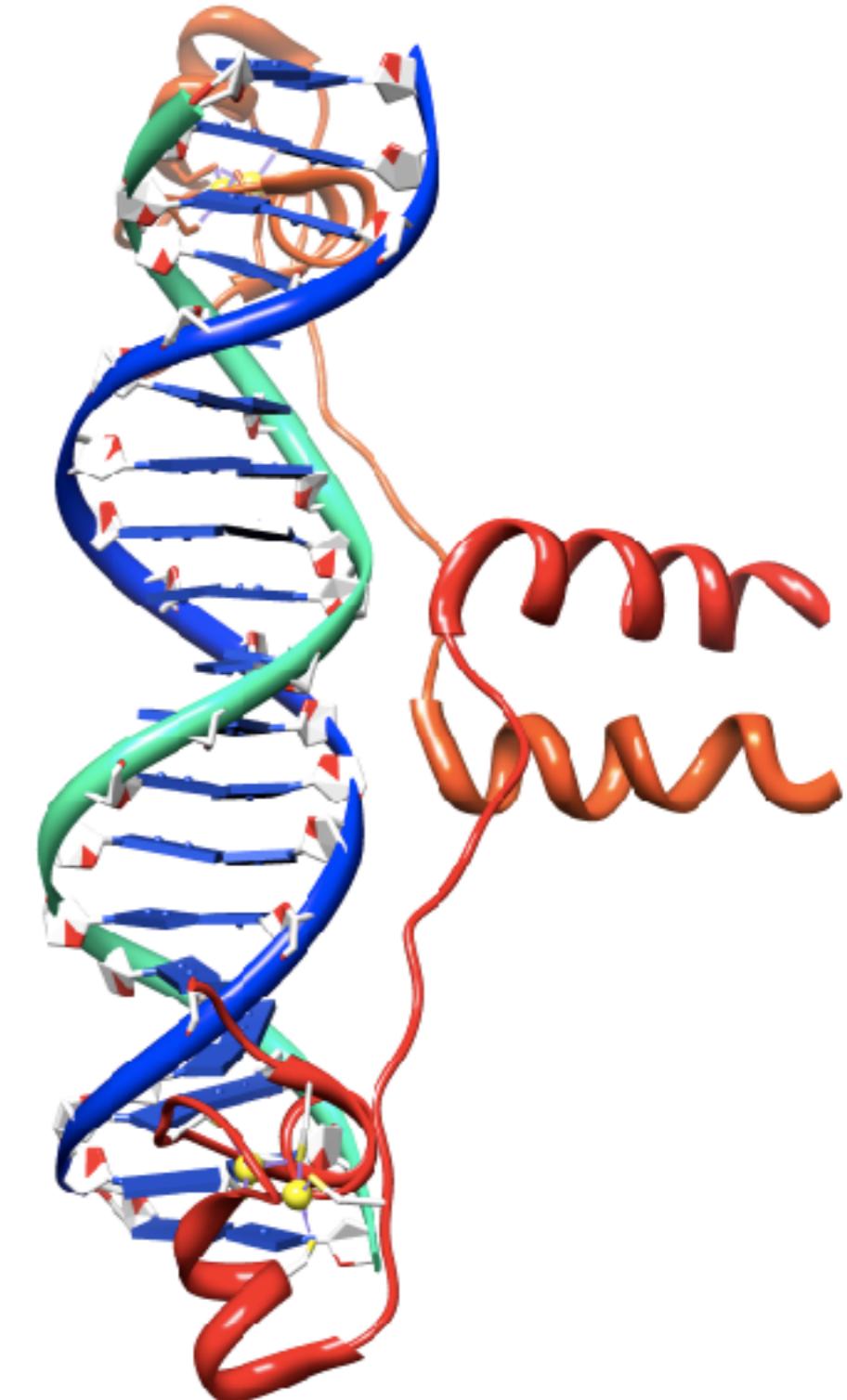
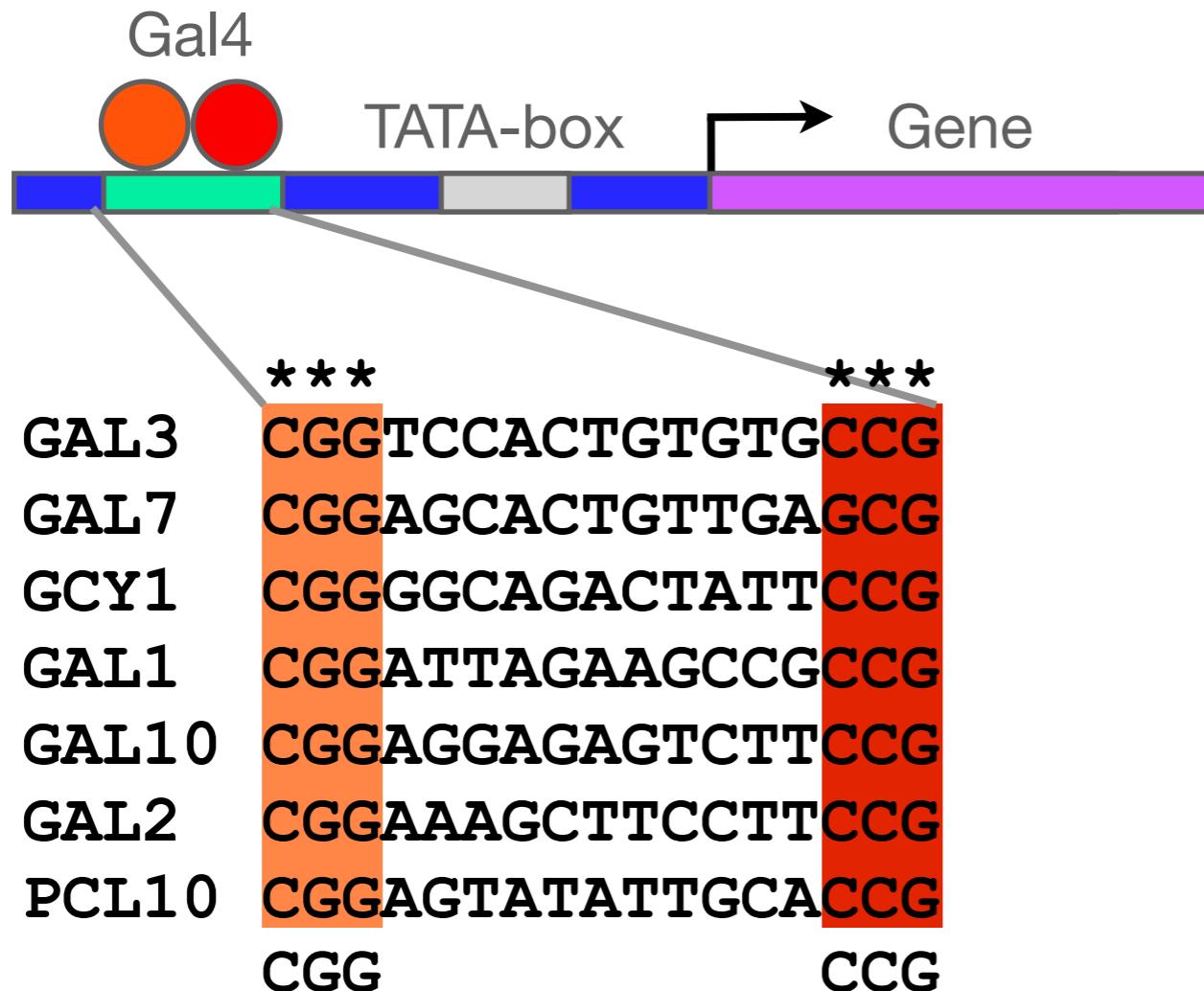
Conservation



Functional cues from conservation patterns...

Many DNA patterns are binding sites for
Transcription Factors.

- E.g., The Gal4 binding sequence
C-G-G-N (11) -C-C-G



Representing recurrent sequence patterns

Beyond knowledge of invariant residues we can define **position-based** representations that highlight the range of permissible residues per position.

- **Pattern:** Describes a motif using a qualitative consensus sequence (e.g., IUPAC or regular expression). N.B. Mismatches are not tolerated!



[LFI]-x-G-[PT]-P-G-x-G-K-[TS]-[AGSI]

- **Profile:** Describes a motif using quantitative information captured in a position specific scoring matrix (weight matrix).
Profiles quantify similarity and often span larger stretches of sequence.
- **Logos:** A useful visual representation of sequence motifs.

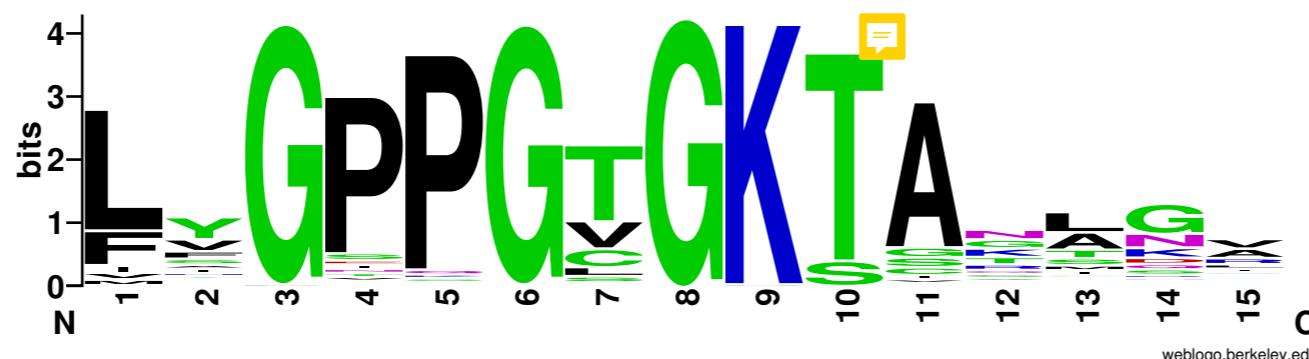
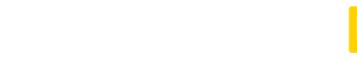


Image generated by:
weblogo.berkeley.edu



PROSITE is a protein pattern and profile database

Currently contains > 1790 patterns and profiles: <http://prosite.expasy.org/>

Example PROSITE patterns:

PS00087; SOD_CU_ZN_1

[GA]-[IMFAT]-H-[LIVF]-H-{S}-x-[GP]-[SDG]-x-[STAGDE]

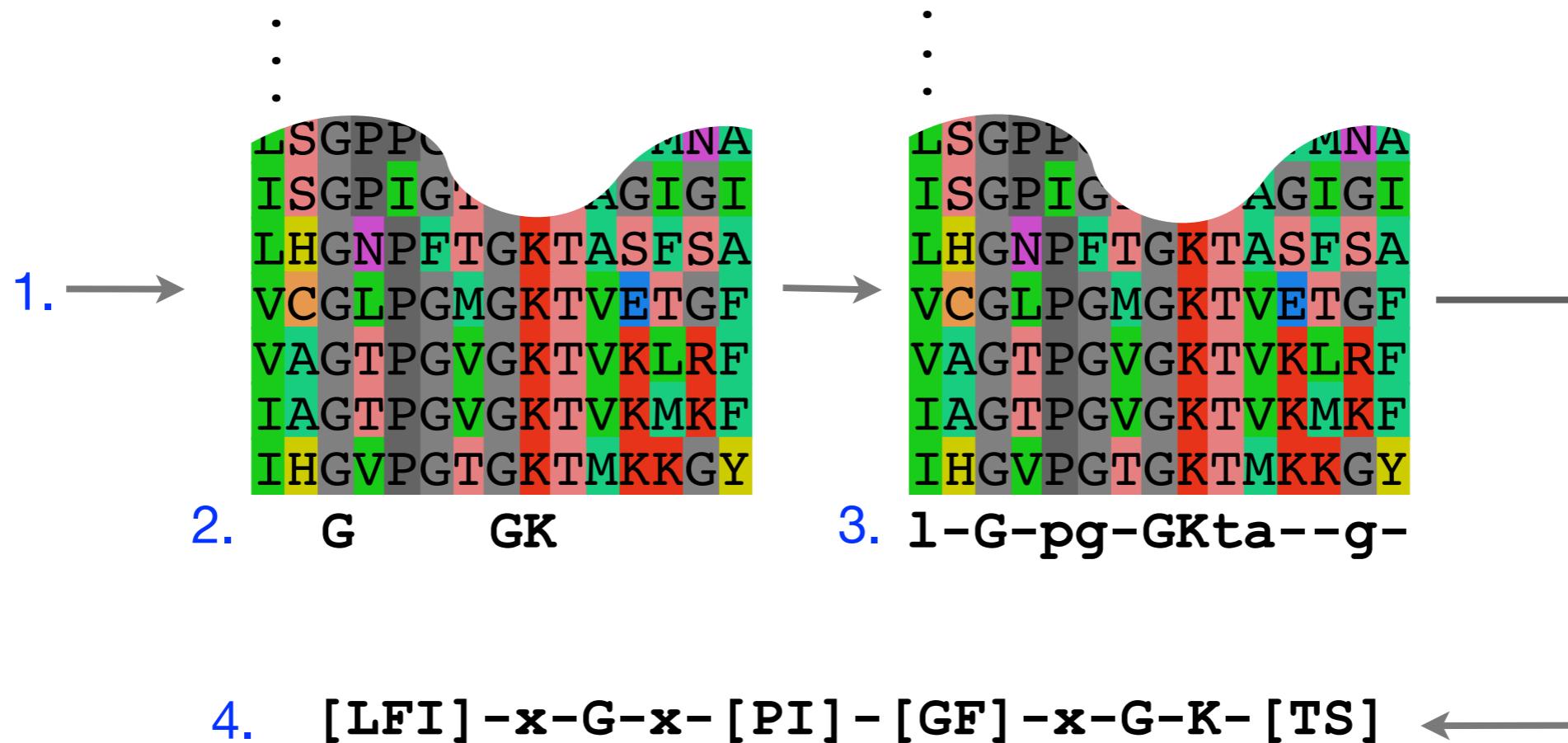
The two Histidines are copper ligands

- Each position in the pattern is separated with a hyphen
- x can match any residue
- [] are used to indicate ambiguous positions in the pattern
 - e.g., [SDG] means the pattern can match S, D, or G at this position
- { } are used to indicate residues that are not allowed at this position
 - e.g., {S} means NOT S (not Serine)
- () surround repeated residues, e.g., A(3) means AAA

Defining sequence patterns

There are four basic steps involved in defining a new PROSITE style pattern:

1. Construct a multiple sequence alignment (MSA)
2. Identify conserved residues
3. Create a core sequence pattern (i.e. *consensus sequence*)
4. Expand the pattern to improve **sensitivity** and **specificity** for detecting desired sequences - more on this shortly...



You want to be sensitive enough to find them all, but not specific enough to find none.

Pattern advantages and disadvantages

Advantages:

- Relatively straightforward to identify (exact pattern matching is fast)
- Patterns are intuitive to read and understand
- Databases with large numbers of protein (e.g., [PROSITE](#)) and DNA sequence (e.g., [JASPER](#) and [TRANSFAC](#)) patterns are available.

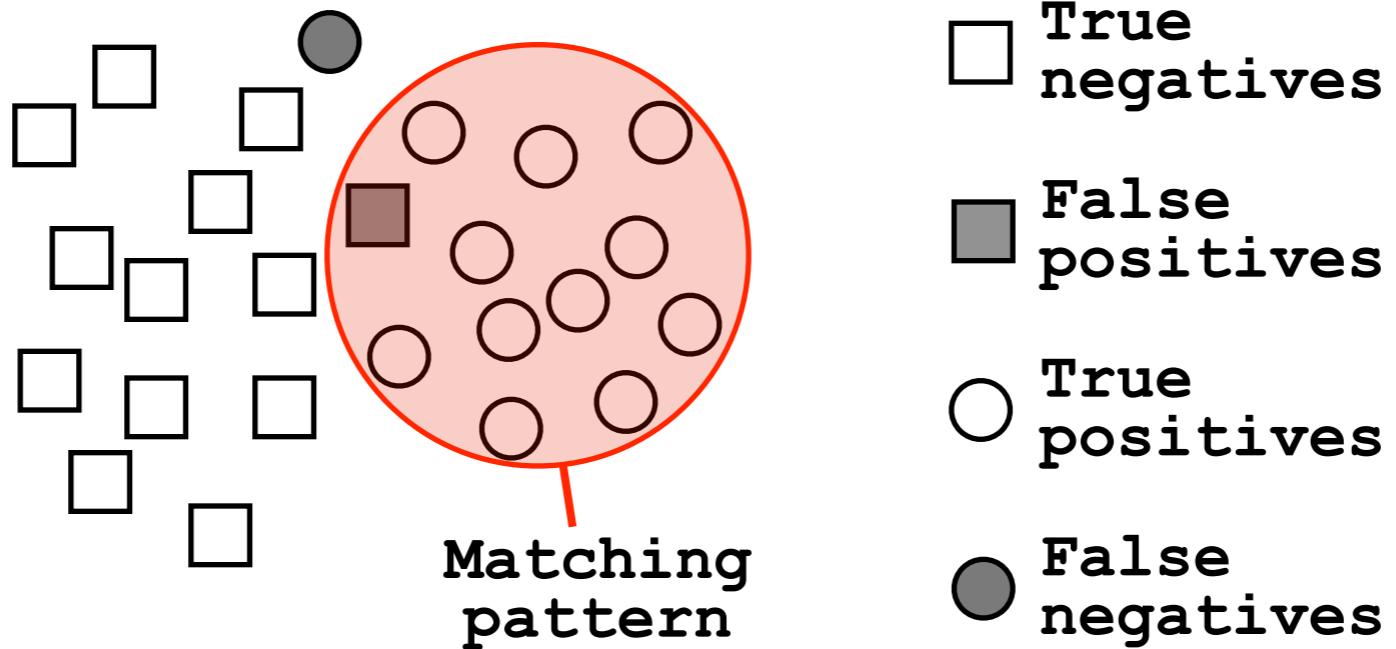
Disadvantages:

- Patterns are qualitative and deterministic (i.e., either matching or not!)
- We lose information about relative frequency of each residue at a position
E.g., [GAC] vs 0.6 G, 0.28 A, and 0.12 C
- Can be difficult to write complex motifs using regular expression notation
- Cannot represent subtle sequence motifs

Side note: pattern sensitivity, specificity, and PPV

In practice it is not always possible to define one single regular expression type pattern which matches all family sequences (*true positives*) while avoiding matches in unrelated sequences (*true negatives*).

circles = true
squares = false

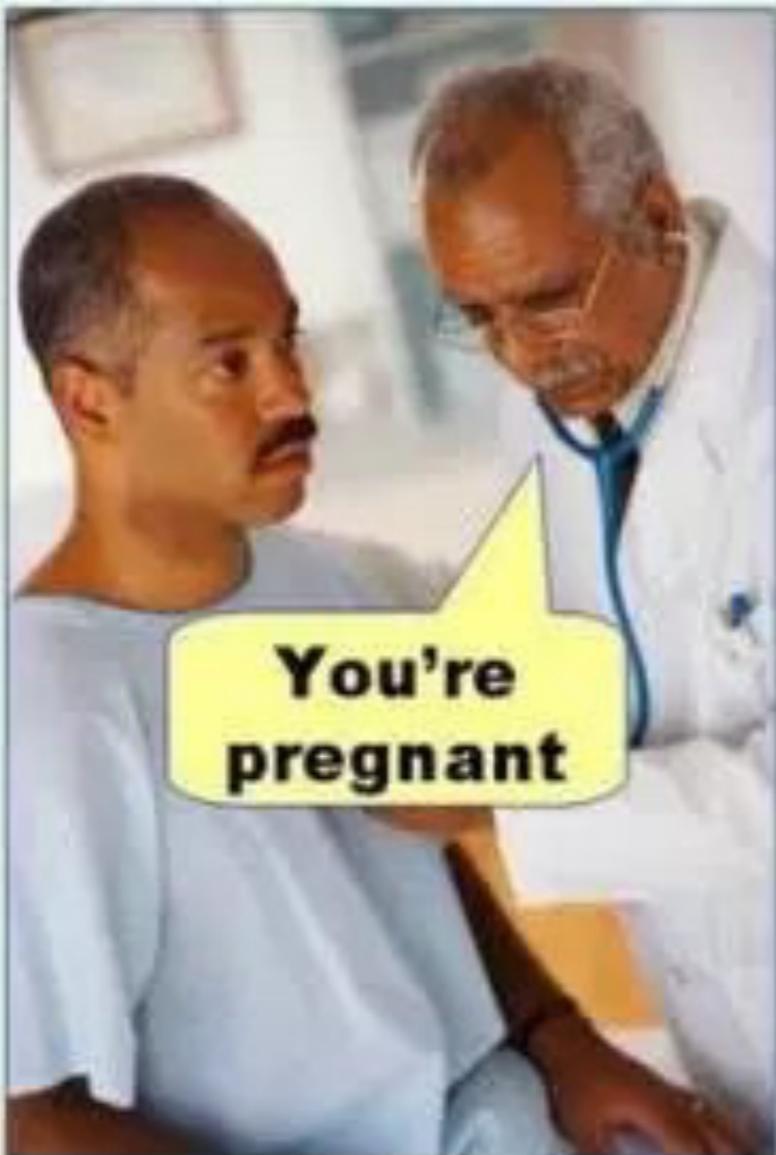


$$\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN})$$

$$\text{Specificity} = \text{TN} / (\text{TN} + \text{FP}) \quad \text{PPV} = \text{TP} / (\text{TP} + \text{FP})$$

The positive predictive value (or PPV) assesses how big a proportion of the sequences matching the pattern are actually in the family of interest.
(i.e., the probability that a positive result is truly positive!)

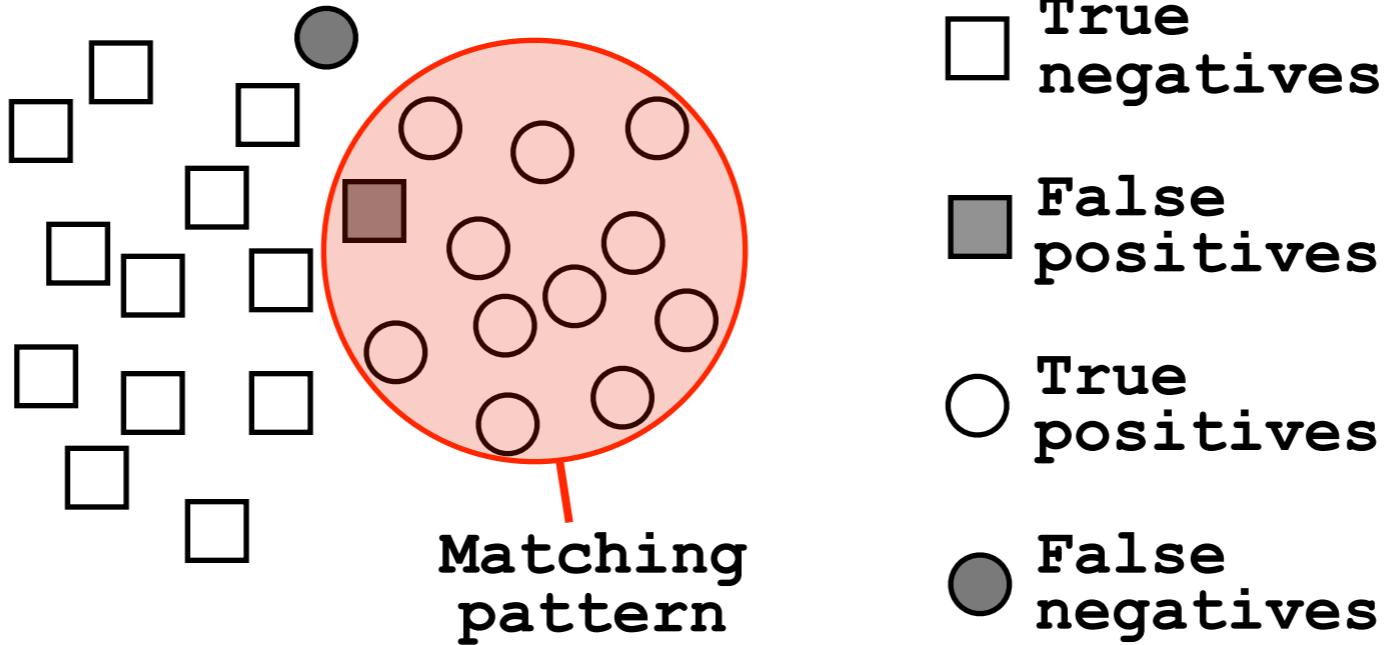
Type I error
(false positive)



Type II error
(false negative)



Side note: pattern sensitivity, specificity, and PPV



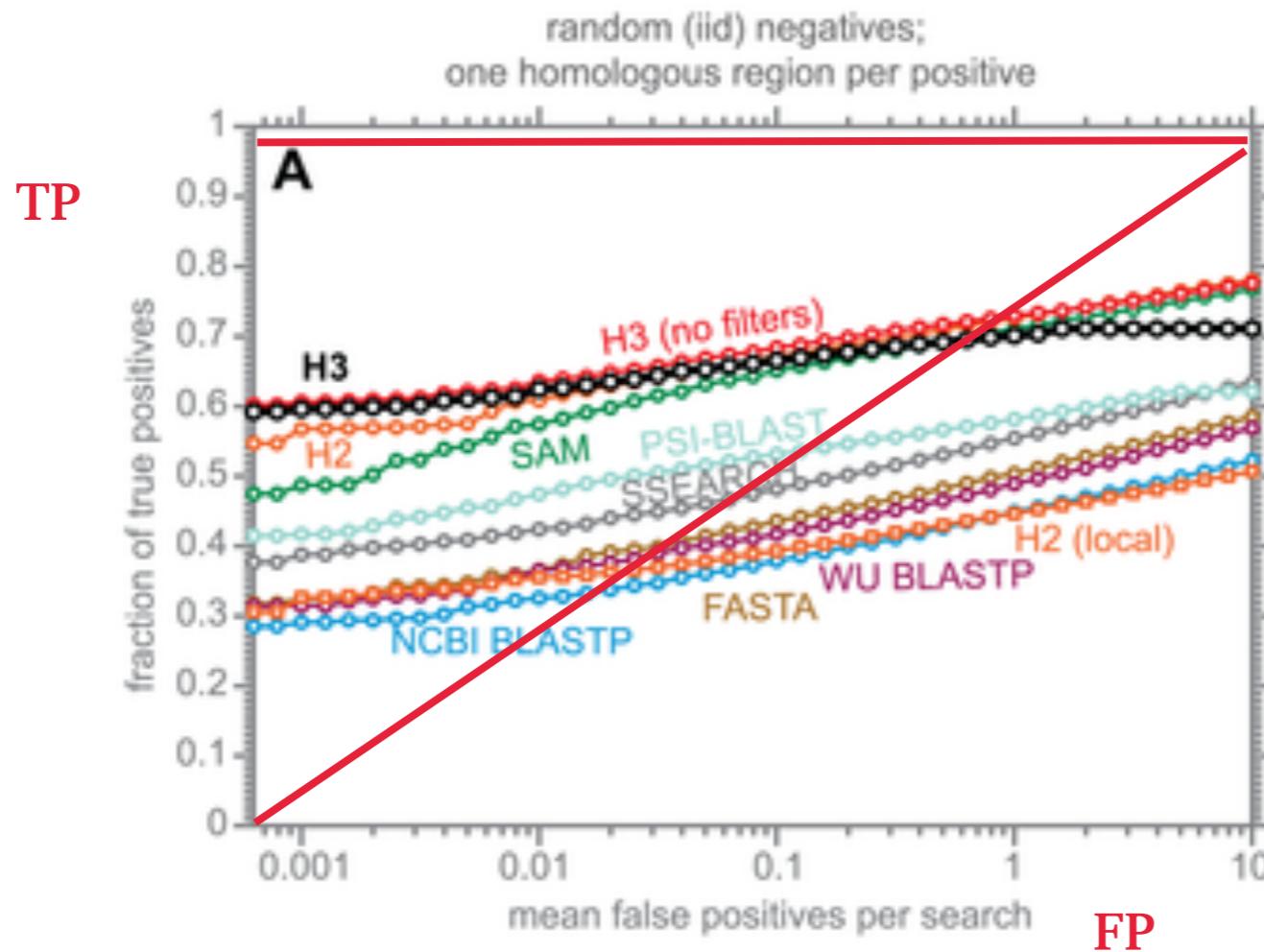
Sensitivity = $TP / (TP+FN)$ = Fraction of total circles we found
(i.e. things we want!)

Specificity = $TN / (TN+FP)$ = Fraction of total squares we missed
(i.e. things we don't want!)

PPV = $TP / (TP+FP)$ = Fraction of our highlighted matches that are actually circles
(i.e. proportion of the things we found that are what we want!)

ROC plot example

ROC plot of sequence searching performance...



More area under the curve makes it a better match (higher true positivity likelihood)

H3 (HMMER3) has a much higher search sensitivity and specificity than BLASTp

In each benchmark, true positive subsequences have been selected to be no more than 25% identical to any sequence in the query alignment ... (see paper for details).

See: Eddy (2011) PLoS Comp Biol 7(10): e1002195

Todays Menu

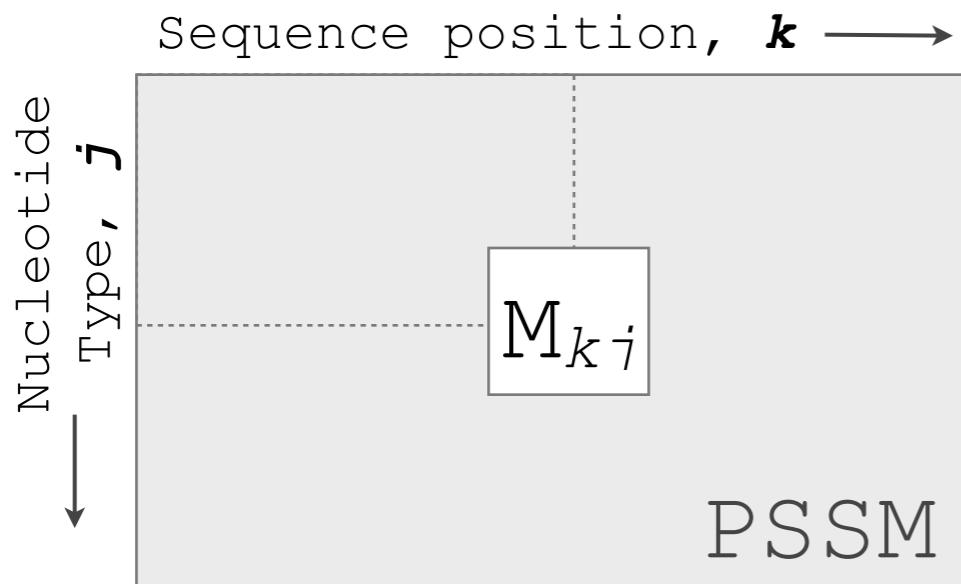
- Sequence motifs and patterns: Simple approaches for finding functional cues from conservation patterns
- Sequence profiles and position specific scoring matrices (PSSMs): Building and searching with profiles, Their advantages and limitations
- PSI-BLAST algorithm: Application of iterative PSSM searching to improve BLAST sensitivity
- Hidden Markov models (HMMs): More versatile probabilistic model for detection of remote similarities

Sequence profiles

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

Unlike deterministic patterns, profiles assign a score to a query sequence and are widely used for database searching.

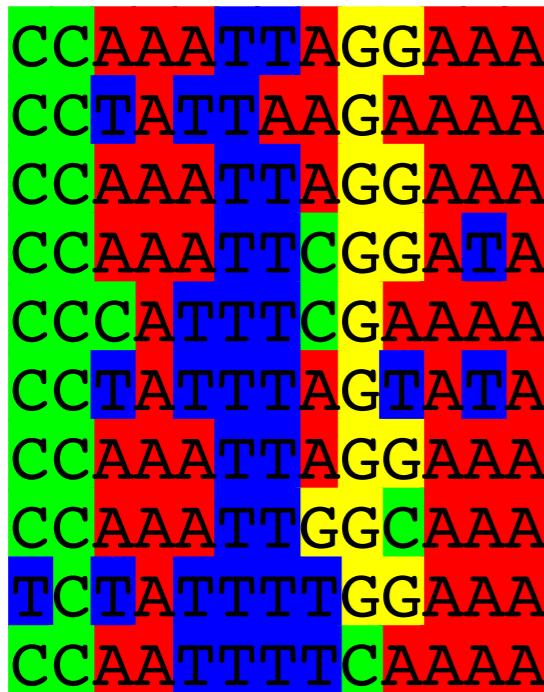
A simple PSSM has as many columns as there are positions in the alignment, and either 4 rows (one for each DNA nucleotide) or 20 rows (one for each amino acid).



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

M_{kj} score for the j th nucleotide at position k
 p_{kj} probability of nucleotide j at position k
 p_j “background” probability of nucleotide j

Computing a transcription factor bind site PSSM



10 sequences

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:	C	C	[ACT]	A	[AT]	T	T	N	G	N	A	[AT]	A

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

C_{kj} Number of jth type nucleotide at position k

Z Total number of aligned sequences

p_j “background” probability of nucleotide j

p_{kj} probability of nucleotide j at position k

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

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

Computing a transcription factor bind site PSSM...

Alignment Matrix: C_{kj}

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

$$k=1, j=A: M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right) = \log\left(\frac{0 + 0.25 / 10 + 1}{0.25}\right) = -2.4$$

$$k=1, j=C: M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right) = \log\left(\frac{9 + 0.25 / 10 + 1}{0.25}\right) = 1.2$$

$$k=1, j=T: M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right) = \log\left(\frac{1 + 0.25 / 10 + 1}{0.25}\right) = -0.8$$

PSSM: M_{kj}

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-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	-0.8	-2.4	-0.2	-2.4

Scoring a test sequence

Query Sequence
CCTATTAGGATA

PSSM:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-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	-0.8	-2.4	-0.2	-2.4
Test seq:	C	C	T	A	T	T	T	A	G	G	A	T	A

$$\begin{aligned}\text{Query Score} &= 1.2 + 1.3 + 0.2 + 1.3 + 0.6 + 1.3 + 1.2 \\ &\quad + 0.6 + 1.2 + 0.6 + 1.3 + -0.2 + 1.3 \\ &= \boxed{11.9}\end{aligned}$$

does this query sequence match the consensus motif enough to be a transcription factor? is 11.9 meaningful?

what's the best possible match?

Scoring a test sequence

Query Sequence
CCTATTAGGATA

PSSM:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-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	-0.8	-2.4	-0.2	-2.4
Test seq:	C	C	T	A	T	T	T	A	G	G	A	T	A

$$\begin{aligned}\text{Query Score} &= 1.2 + 1.3 + 0.2 + 1.3 + 0.6 + 1.3 + 1.2 \\ &\quad + 0.6 + 1.2 + 0.6 + 1.3 + -0.2 + 1.3 \\ &= 11.9\end{aligned}$$

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

Scoring a test sequence...

Query Sequence
CCTATTAGGATA

Best Possible Sequence
CCAATTAGGAAA

PSSM:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-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	-0.8	-2.4	-0.2	-2.4

Max Score: C C A A T T T A G G A A A

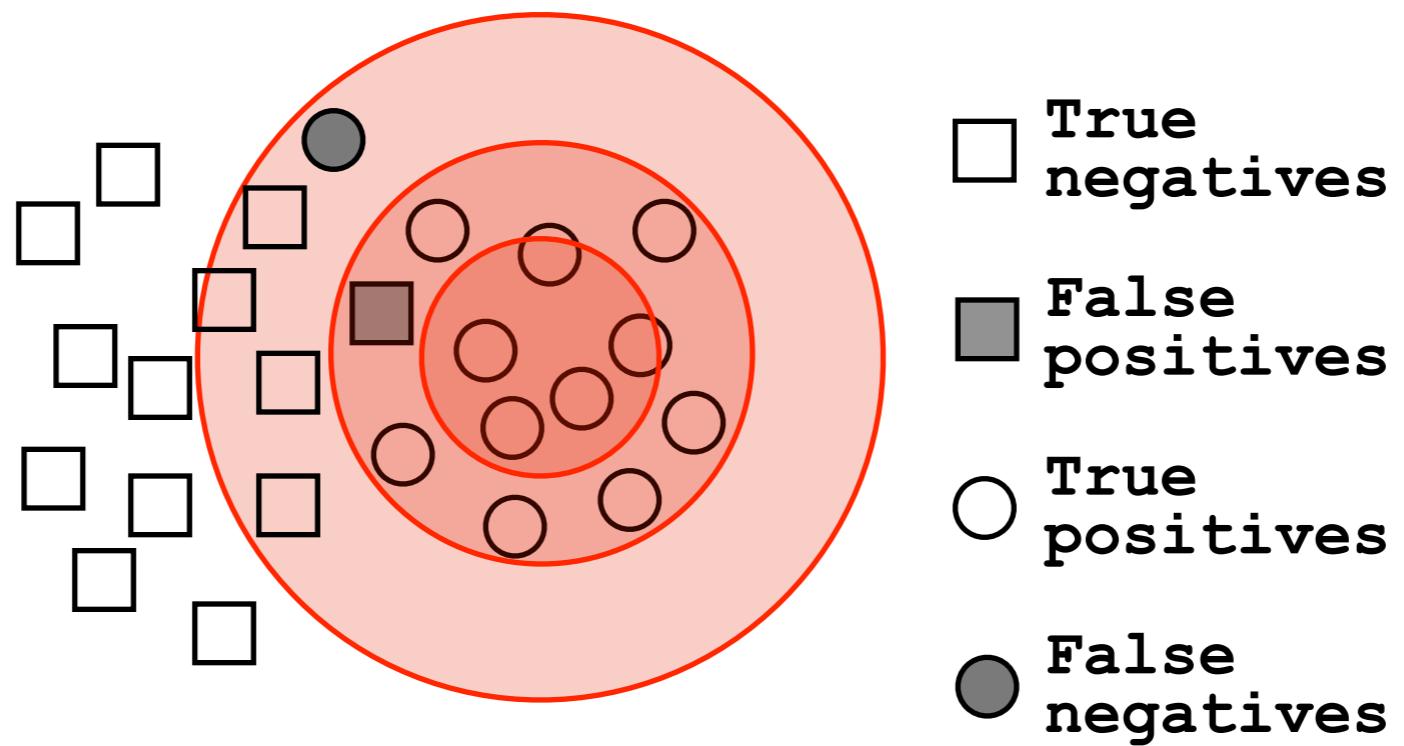
$$\begin{aligned}\text{Max Score} &= 1.2 + 1.3 + 0.8 + 1.3 + 0.6 + 1.3 + 1.2 \\ &\quad + 0.6 + 1.2 + 0.6 + 1.3 + 1.1 + 1.3 \\ &= 13.8\end{aligned}$$

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

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

Picking a threshold for PSSM matching

Again, you want to select a threshold that **minimizes FPs** (e.g., how many shuffled or random sequences does the PSSM match with that score) and **minimizes FNs** (e.g., how many of the ‘real’ sequences are missed with that score).



$$FP=0, FN=7, TP=5$$

$$5/(5+0) = 1$$

$$FP=1, FN=1, TP=11$$

$$11/(11+1) = 0.92$$

$$FP=5, FN=0, TP=12$$

$$12/(12+5) = 0.71$$

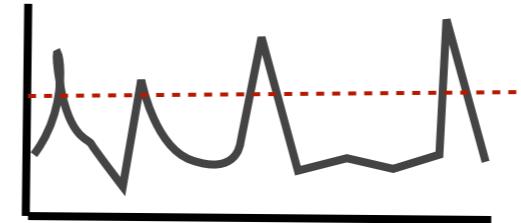
Q. Which threshold has the best PPV ($TP/(TP+FP)$) ?

Searching for PSSM matches

If we do not allow gaps (i.e., no insertions or deletions):

- Perform a linear scan, scoring the match to the PSSM at each position in the sequence - the “sliding window” method

GCAGGTATCCTATTAGCAATAGC....
↓↓↓↓↓
→



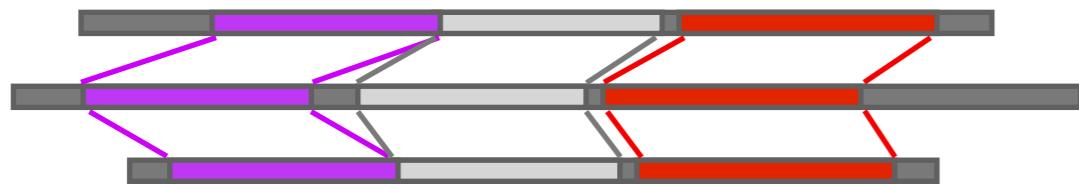
If we allow gaps:

- Can use dynamic programming to align the profile to the protein sequence(s) (with gap penalties)
We will discuss PSI-BLAST shortly...
see Mount, Bioinformatics: sequence and genome analysis (2004)
- Can use hidden Markov Model-based methods
We will cover HMMs in the next lecture...
see Durbin et al., Biological Sequence Analysis (1998)

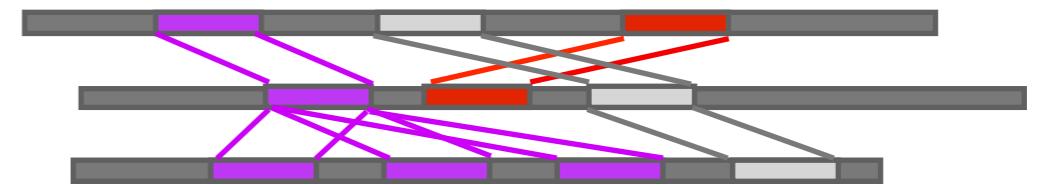
Side note: Building PSSMs from unaligned sequences

Patterns and profiles are most often built on the basis of known site equivalences (i.e. from a pre-calculated MSA).

However, a number of programs have been developed that employ local multiple alignments to search for common sequence elements in unaligned sequences.



Global similarity



Local non-consistent similarity

Gibbs *sampling* methods:

Motif Sampler - <http://bayesweb.wadsworth.org/gibbs/gibbs.html>

AlignAce - <http://atlas.med.harvard.edu/cgi-bin/alignace.pl>

Expectation maximization method:

MEME - <http://meme.sdsc.edu/>

See: Lawrence et al. (1993) Science. 262, 208-14

Side note: Profiles software and databases

Pftools is a package to build and search with profiles,
<http://www.isrec.isb-sib.ch/ftp-server/pftools/>

The package contains (among other programs):

- ▶ **pfmake** for building a profile starting from multiple alignments
- ▶ **pfsearch** to search a protein database with a profile
- ▶ **pfscan** to search a profile database with a protein

PRINTS database of PSSMs

<http://bioinf.man.ac.uk/dbbrowser/PRINTS>

Collection of conserved motifs used to characterize a protein

- ▶ Uses fingerprints (conserved motif groups).
- ▶ Very good to describe sub-families.

BLOCKS is another PSSMs database similar to prints

<http://www.blocks.fhcrc.org>

ProDom is collection of protein motifs obtained automatically using PSI-BLAST

<http://prodes.toulouse.inra.fr/prodom/doc/prodom.html>

Side note: Profiles software and databases...

InterPro is an attempt to group a number of protein domain databases.

<http://www.ebi.ac.uk/interpro>

It currently includes:

- ▶ Pfam
 - ▶ PROSITE
 - ▶ PRINTS
 - ▶ ProDom
 - ▶ SMART
 - ▶ TIGRFAMs
-
- InterPro tries to have and maintain a high quality of annotation
 - The database and a stand-alone package (**iprscan**) are available for UNIX platforms, see:

<ftp://ftp.ebi.ac.uk/pub/databases/interpro>

Todays Menu

- Sequence motifs and patterns: Simple approaches for finding functional cues from conservation patterns
- Sequence profiles and position specific scoring matrices (PSSMs): Building and searching with profiles, Their advantages and limitations
- PSI-BLAST algorithm: Application of iterative PSSM searching to improve BLAST sensitivity
- Hidden Markov models (HMMs): More versatile probabilistic model for detection of remote similarities

Your Turn!

Hands-on sections 1 & 2: Comparing methods and the trade-off between sensitivity, selectivity and performance

~50 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).

By default BLASTp match scores come from the BLOSUM62 matrix

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9																			
S	-1	4																		
T	-1	1	5																	
P	-3	-1	-1	7																
A	0	1	0	-1	4															
G	-3	0	-2	-2	0	5														
N	-3	1	0	-2	-2	0	6													
D	-3	0	-1	-1	-2	-1	1	6												
E	-4	0	-1	-1	-1	-2	0	2	5											
Q	-3	0	-1	-1	-1	-2	0	0	2	5										
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8									
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5								
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5						
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4					
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4				
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4			
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	3	7		
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11

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

PSI-BLAST: Position specific iterated BLAST

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

Inspect the blastp output to identify empirical “rules” regarding amino acids tolerated at each position

<u>730496</u>	66	FTVDENGQMSATAKGRVRLFNNWDVCADMIGSFTD TEDPAFKMKYWGVASFLQKGNDH	125
<u>200679</u>	63	FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTF TD TEDPAFKMKYWGVASFLQRGNDDH	122
<u>206589</u>	34	FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTF TD TEDPAFKMKYWGVASFLQRGNDDH	93
<u>2136812</u>	2	MSATAKGRVRLNNWDVCADMVGTF TD TEDPAFKMKYWGVASFLQKGNDH	53
<u>132408</u>	65	FKIEDNGKTTATAKGRVRILDKLELCANMVGTF IETNDPAKYRMKYHGALAILERGLDDH	124
<u>267584</u>	44	FSVDESGKVTATAHGRVIILNNWEMCANMF GTFEDTPDPAKFKMRYWGAASYLQTGNDDH	103
<u>267585</u>	44	FSVDGSGKVTATAQGRVIILNNWEMCANMF GTFEDTPDPAKFKMRYWGAASYLQSGNDDH	103
<u>8777608</u>	63	FTIHEDGAMTATAKGRVIILNNWEMCADMMATFETTPDPAKFRMRYWGAASYLQTGNDDH	122
<u>6687453</u>	60	FKVEEDGTMTATAIGRVIILNNWEMCANMF GTFEDTEDPAFKMKYWGAASYLQTGYDDH	119
<u>10697027</u>	81	FKVQEDGTMTATATGRVIILNNWEMCANMF GTFEDTEEPARFKMKYWGAASYLQTGYDDH	140
<u>13645517</u>	1	MVGTF TD TEDPAFKMKYWGVASFLQKGNDH	32
<u>13925316</u>	38	FSVDGSGKMTATAQGRVIILNNWEMCANMF GTFEDTPDPAKFKMRYWGAASYLQSGNDDH	97
<u>131649</u>	65	YTVEEDGTMTASSKGRVKLFGFUVVICADMAAQYTDPTTPAKMYMTYQGLASYLSSGGDNY	126

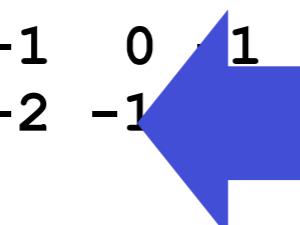
↑ ↑ ↑ ↑ ↑

R,I,K **C** **D,E,T** **K,R,T** **N,L,Y,G**

	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	3	1	2	2	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								1	-3	-3	12	2	-3		
4 V	0	-3	-3	-4	-1	-3	-3	-4								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	0	-3	-3	-1	-2	-1	1			
8 L	-1	-3	-3	-4	-1	-3	-3	-4	-3	2	2	-3	1	3	-3	-2	-1	-2	0	3			
9 L	-1	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	2			
10 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1			
11 A	5	-2	-2	-2	-1	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0			
12 A	5	-2	-2	-2	-1	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0			
13 W	-2		All the amino acids from position 1 to N (the end of your query protein)												-3	2	1	-3	-3	-2	7	0	0
14 A	3														-1	-2	-3	-1	1	-1	-3	-3	-1
15 A	2														0	-2	-3	-1	3	0	-3	-2	-2
16 A	4														-1	-1	-3	-1	1	0	-3	-2	-1
...																							
37 S	2	-1	0	-1	-1	0	0	0	-1	-2	-3	0	-2	-3	-1	4	1	-3	-2	-2			
38 G	0	-3	-1	-2	-3	-2	-2	6	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4			
39 T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-3	-2	0			
40 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-3	-2	1	-4	-3	-3	9	2	-3			
41 Y	-2	-2	-2	-3	-3	-2	-2	-3	2	-2	-1	-2	-1	3	-3	-2	-2	2	7	-1			
42 A	4	-2	-2	-2	-1	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0			

	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	0	-3	-3	-1	-2	-1	1		
8 L	-1	-3	-3	-4	-1	-3	-3	-4	-3	2	2	-3	1	3	-3	-2	-1	-2	0	3		
9 L	-1	-3	-4	-4	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	-2	-1	2			
10 L	-2	-2	-4	-4	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	-2	-1	1			
11 A	5	-2	-2	-2	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	1	0	-3	-2	0	
12 A	5	-2	-2	-2	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	1	0	-3	-2	0	
13 W	-2	-3	-4	-4	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-2	7	0	0			
14 A	3	-2	-1	-2	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	1	-1	-3	-3	-1	
15 A	2	-1	0	-1	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	3	0	-3	-2	-2	
16 A	4	-2	-1	-1	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	1	0	-3	-2	-1	
...																						
37 S	2	-1	0	-1	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	4	1	-3	-2	-2	
38 G	0	-3	-1	-2	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-2	0	-2	-3	-3	-4	
39 T	0	-1	0	-1	-1	-2	-2	-3	-2	-2	-1	-2	-2	-3	-3	-1	1	5	-3	-2	0	
40 W	-3	-3	-4	-5	-2	-1	-2	-3	-2	-2	-1	-2	-2	-3	-3	-4	-3	-3	9	2	-3	
41 Y	-2	-2	-2	-3	-1	-1	-1	-2	-2	-2	-1	-1	-2	-2	-3	-3	-2	-2	7	-1		
42 A	4	-2	-2	-2	-1	-1	-1	-2	-2	-2	-1	-1	-2	-2	-3	-3	-1	1	0	-3	-2	0

Note: A given amino acid (such as alanine) in your query protein can receive different scores for matching alanine depending on the position in the protein
 (BLOSUM SAA = +4)



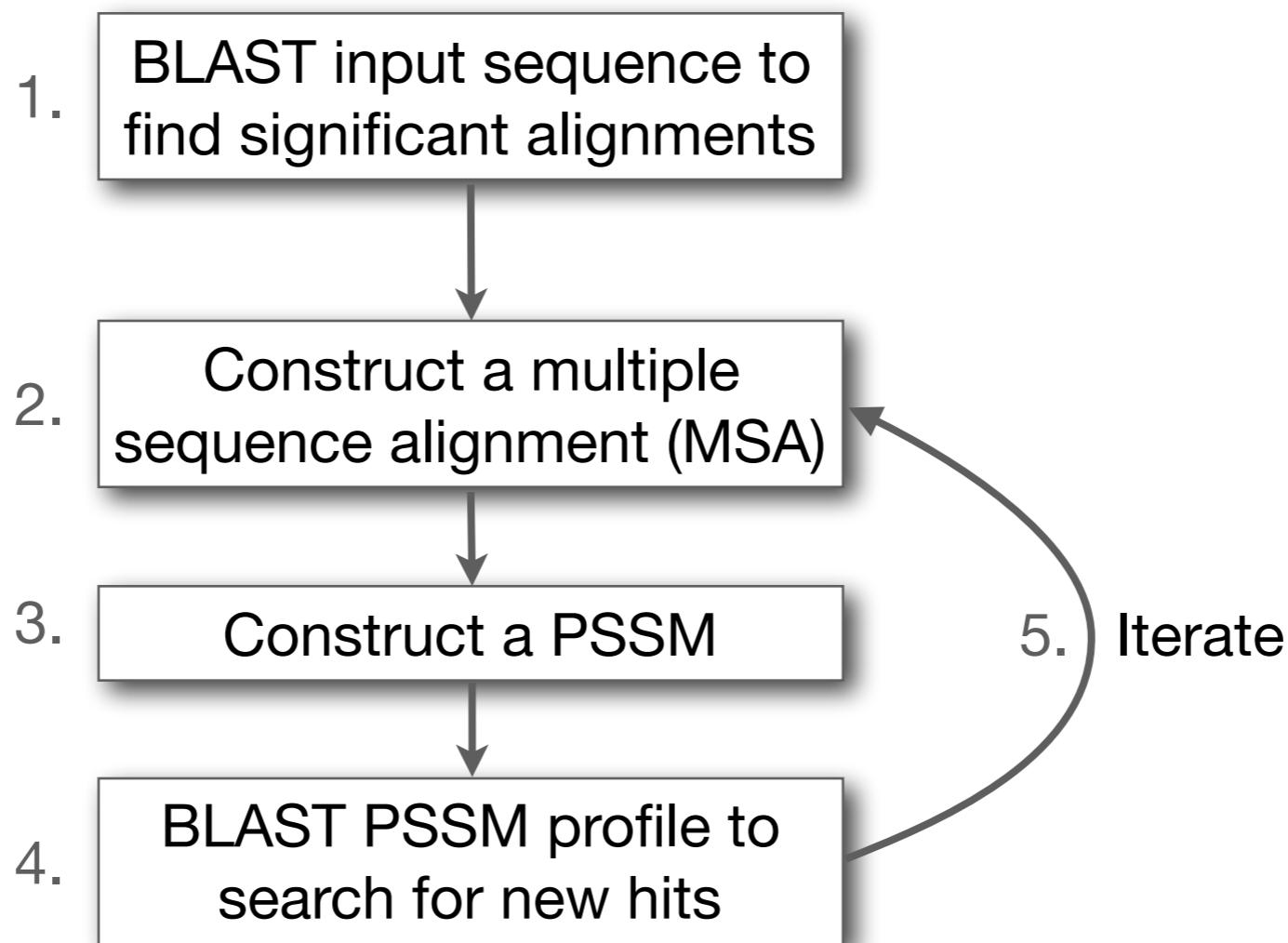
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
1	M																				
2	K																				
3	W																				
4	V																				
5	W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-2	-3	-2	1	-4	-3	-3	12	2	-3	
6	A	5	-2	-2	-2	-1	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0
7	L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3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

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

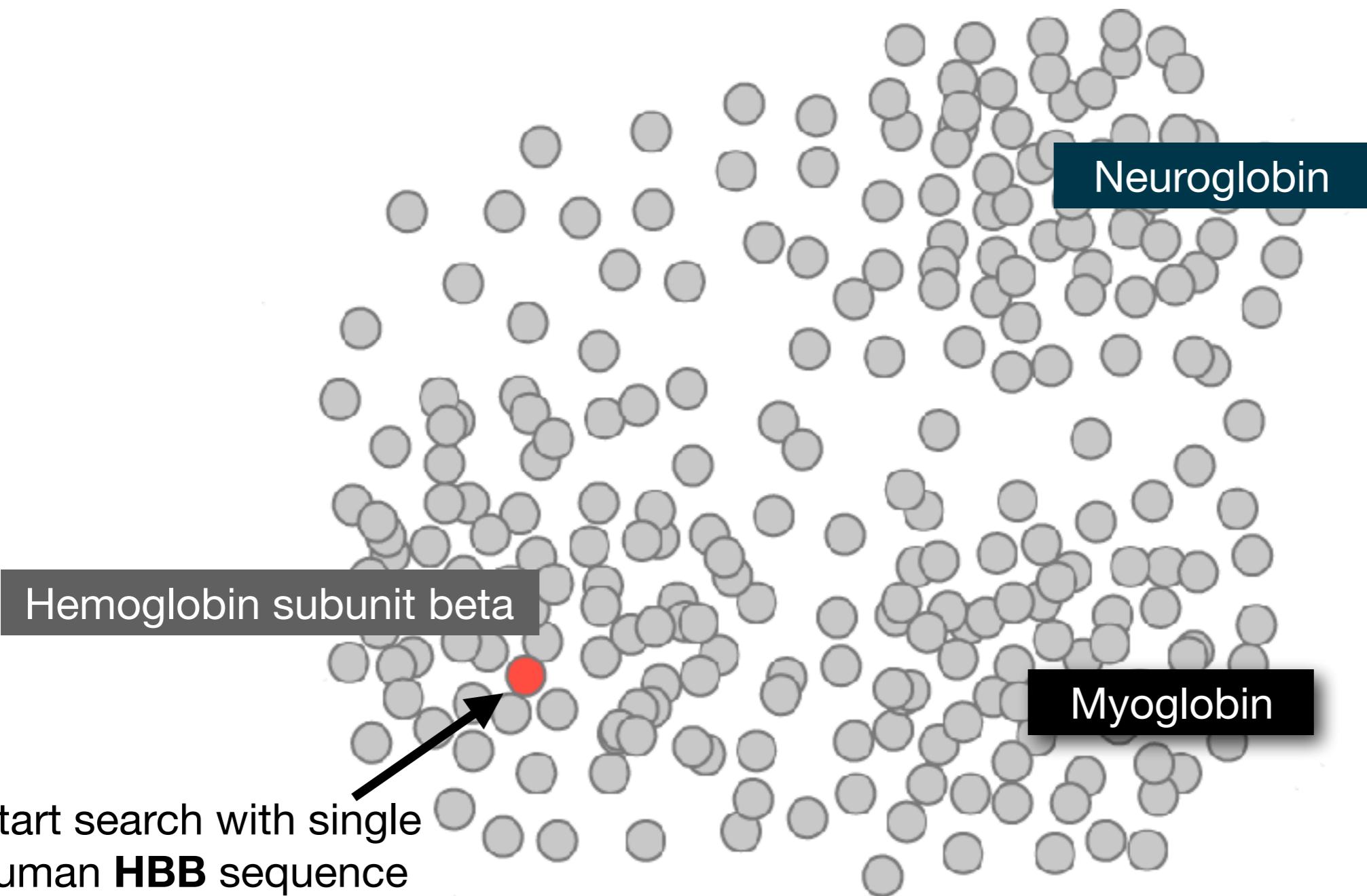
Note: A given amino acid (such as alanine) in your query protein can receive different scores for matching alanine depending on the position in the protein
 (BLOSUM SAA = +4)

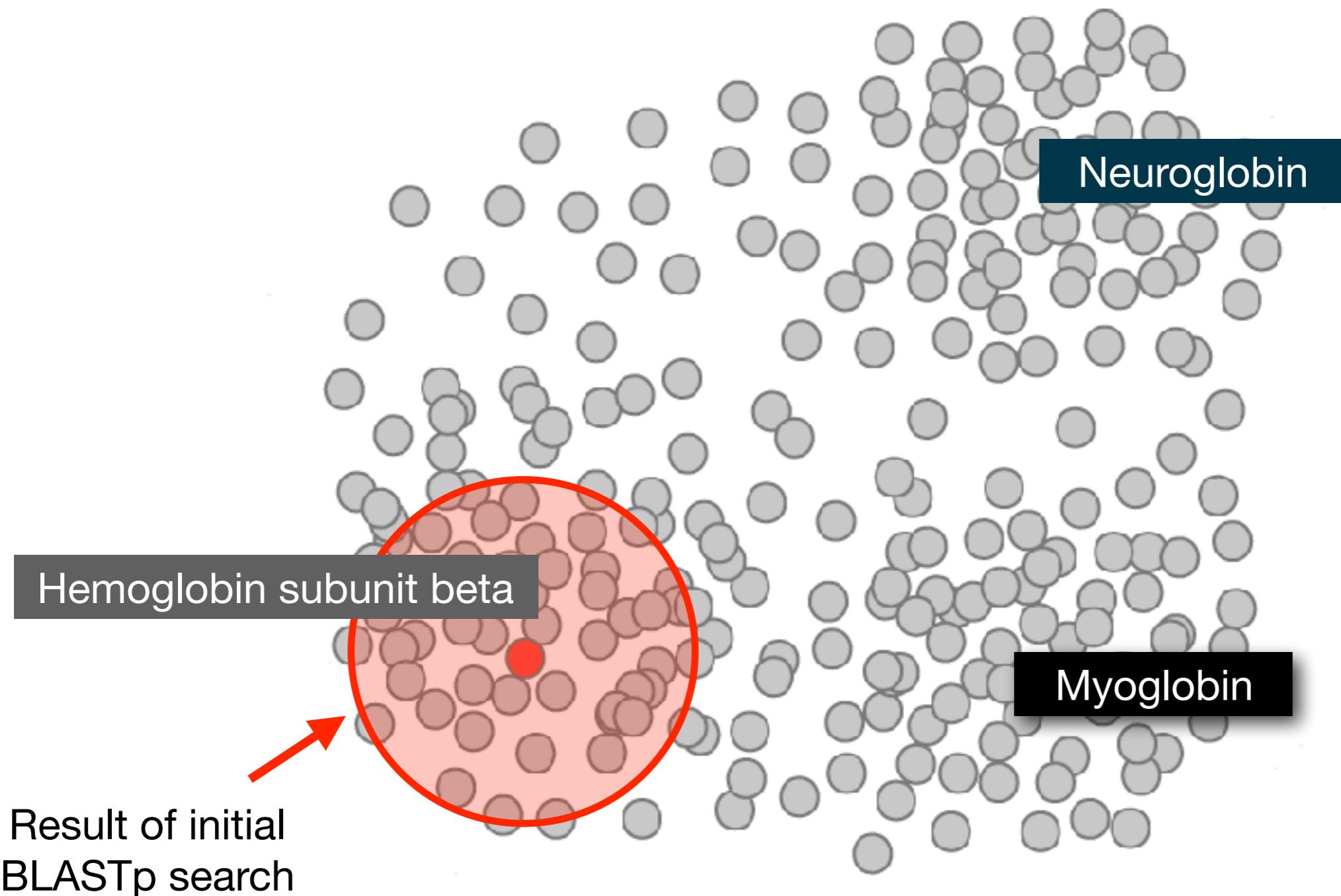
PSI-BLAST: Position-Specific Iterated BLAST

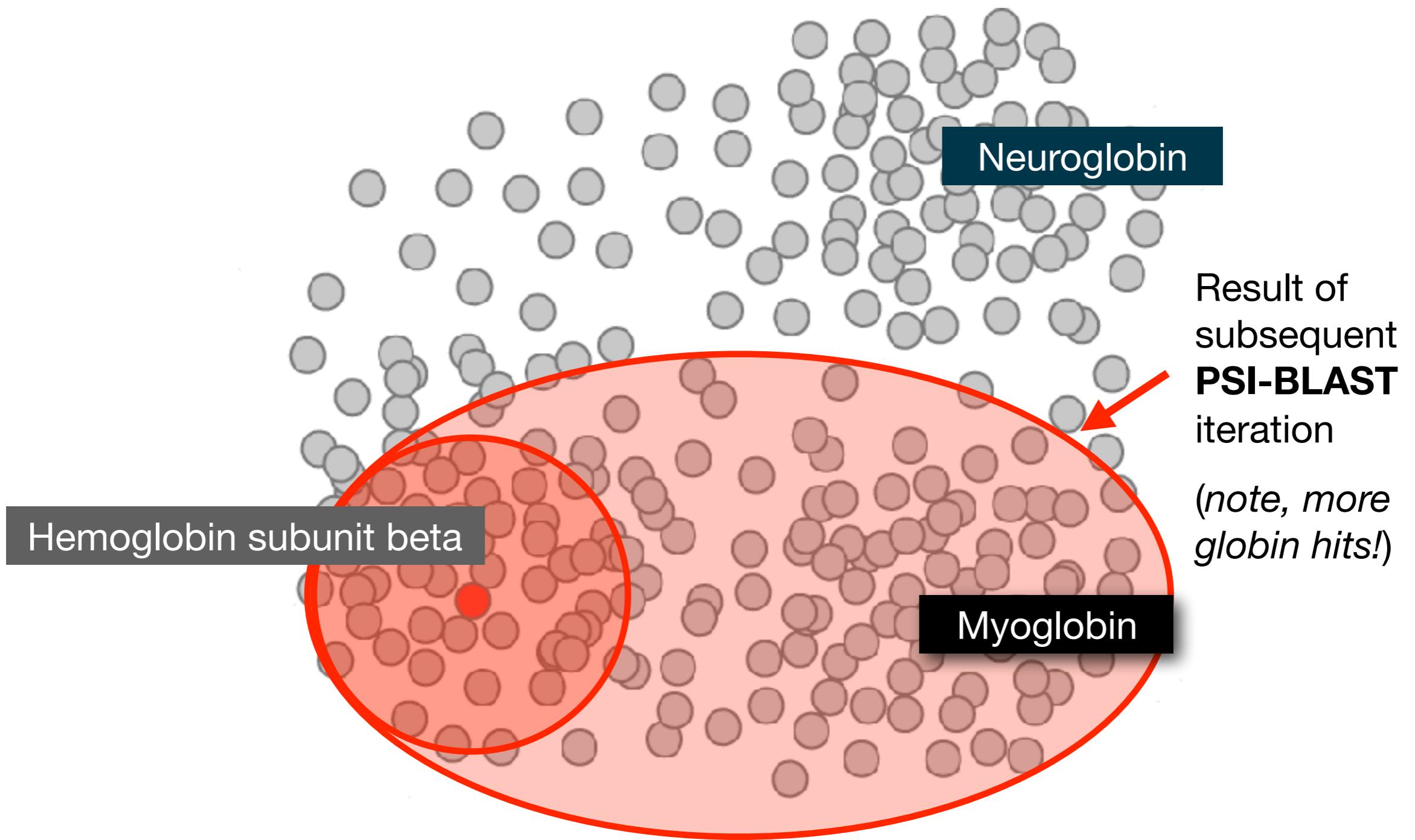
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



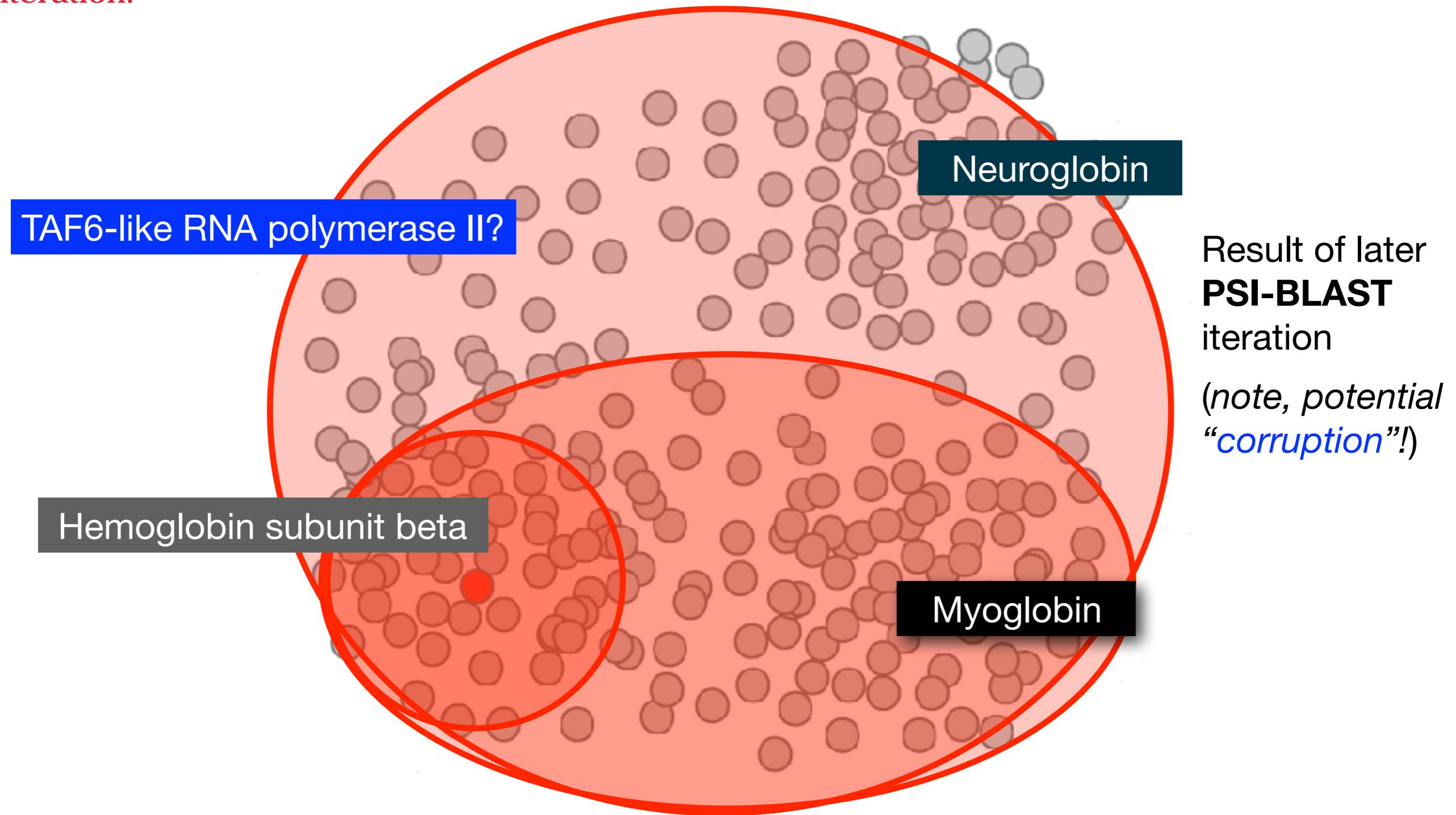
(see Altschul *et al.*, Nuc. Acids Res. (1997) 25:3389-3402)







PSSM can be corrupted, too, which is why PSI-BLAST asks you what you want to include in each iteration.



Description	Max score	Total score	Query cover	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

1

Description	Max score	Total score	Query cover	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

1

2

New relevant globins found only by PSI-BLAST

Description	Max score	Total score	Query cover	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 sapien	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

Description	Max score	Total score	Query cover	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

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

2

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 sapien	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

3

Score and E value depends on PSSM

PSI-BLAST is performed in five steps

- A normal blastp search uses a scoring matrix (e.g., BLOSUM62) to perform pairwise alignments of your query sequence (such as RBP) against the database. PSI-BLAST also begins with a protein query that is searched against a database of choice.
- PSI-BLAST constructs a multiple sequence alignment (MSA) from an initial blastp-like search. It then creates a **PSSM** based on that multiple alignment.
- This **PSSM** is then used as a query to search the database again.
- PSI-BLAST estimates the statistical significance of the database matches, essentially using the parameters we described for gapped alignments.
- The search process is continued iteratively,^{typically 3 to 5 times.} At each step a new PSSM is built.

PSI-BLAST returns dramatically more hits

You must decide how many iterations to perform and which sequences to include!

You can stop the search process at any point - typically whenever few new results are returned or when no new “sensible” results are found.

Iteration	Hits with $E < 0.005$	Hits with $E > 0.005$
1	34	61
2	314	79
3	416	57
4	432	50
5	432	50

Human retinol-binding protein 4 (RBP4; P02753) was used as a query in a PSI-BLAST search of the RefSeq database.

Example PSI-BLAST PSSM at iteration 3

The PSI-BLAST PSSM is essentially a query customized scoring matrix that is more sensitive than BLOSUM (e.g. BLOSUM S_{AA} = +4)

PSI-BLAST errors: the corruption problem

The main source of error in PSI-BLAST searches is the spurious amplification of sequences that are unrelated to the query.

Use biology knowledge to select only the relevant proteins.

There are three main approaches to stopping corruption of PSI-BLAST queries:

- Perform multi-domain splitting of your query sequence
 - If a query protein has several different domains PSI-BLAST may find database matches related to both individually. One should not conclude that these hits with different domains are related.
 - Often best to search using just one domain of interest.
- Inspect each PSI-BLAST iteration removing suspicious hits.
 - E.g., your query protein may have a generic coiled-coil domain, and this may cause other proteins sharing this motif (such as myosin) to score better than the inclusion threshold even though they are not related.
 - Use your biological knowledge!
- Lower the default expect level (e.g., $E = 0.005$ to $E = 0.0001$).
 - This may suppress appearance of FPs (but also TPs)

Profile advantages and disadvantages

Advantages:

- Quantitate with a good scoring system
- Weights sequences according to observed diversity
Profile is specific to input sequence set
- Very sensitive
Can detect weak similarity
- Relatively easy to compute
Automatic profile building tools available

Disadvantages:

- If a mistake enters the profile, you may end up with irrelevant data
The corruption problem!
- Ignores higher order dependencies between positions
i.e., correlations between the residue found at a given position and those found at other positions (e.g. salt-bridges, structural constraints on RNA etc...)
- Requires some expertise and oversight to use proficiently

Todays Menu

- Sequence motifs and patterns: Simple approaches for finding functional cues from conservation patterns
- Sequence profiles and position specific scoring matrices (PSSMs): Building and searching with profiles, Their advantages and limitations
- PSI-BLAST algorithm: Application of iterative PSSM searching to improve BLAST sensitivity
- Hidden Markov models (HMMs): More versatile probabilistic model for detection of remote similarities

Your Turn!

Hands-on sections 3 & 4:
Comparing methods and the trade-off
between sensitivity, selectivity and
performance

~30 mins

Problems with PSSMs: Positional dependencies

Do not capture positional dependencies

WEIRD
WEIRD
WEIQH
WEIRD
WEIQH

D					0.6
E		I			
H					0.4
I			I		
Q				0.4	
R				0.6	
W	I				

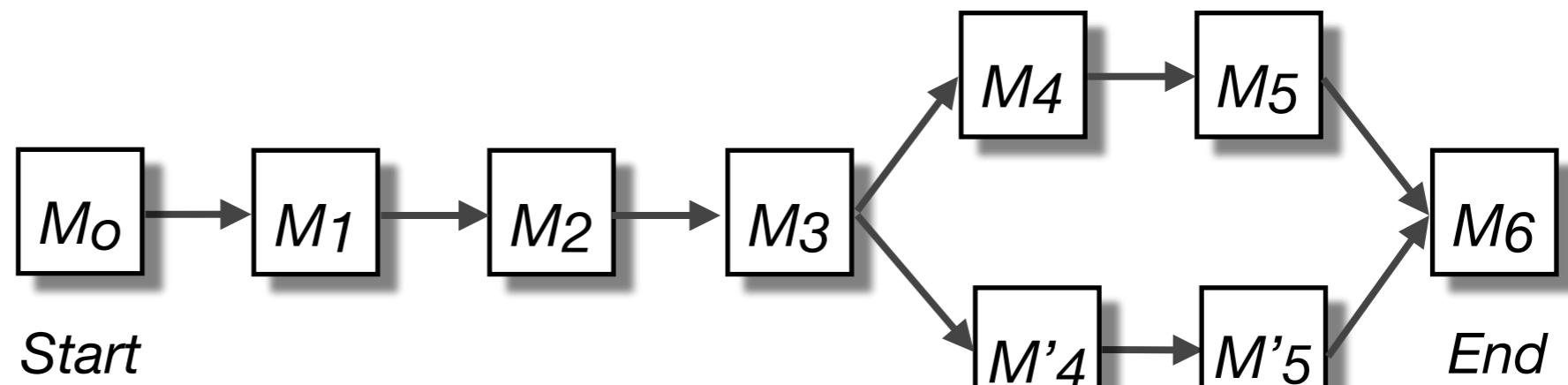
Note: We never see **QD** or **RH**, we only see **RD** and **QH**.
However, $P(RH)=0.24$, $P(QD)=0.24$, while $P(QH)=0.16$

Markov chains: Positional dependencies



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

WEIRD
WEIRD
WEIQH
WEIRD
WEIQH

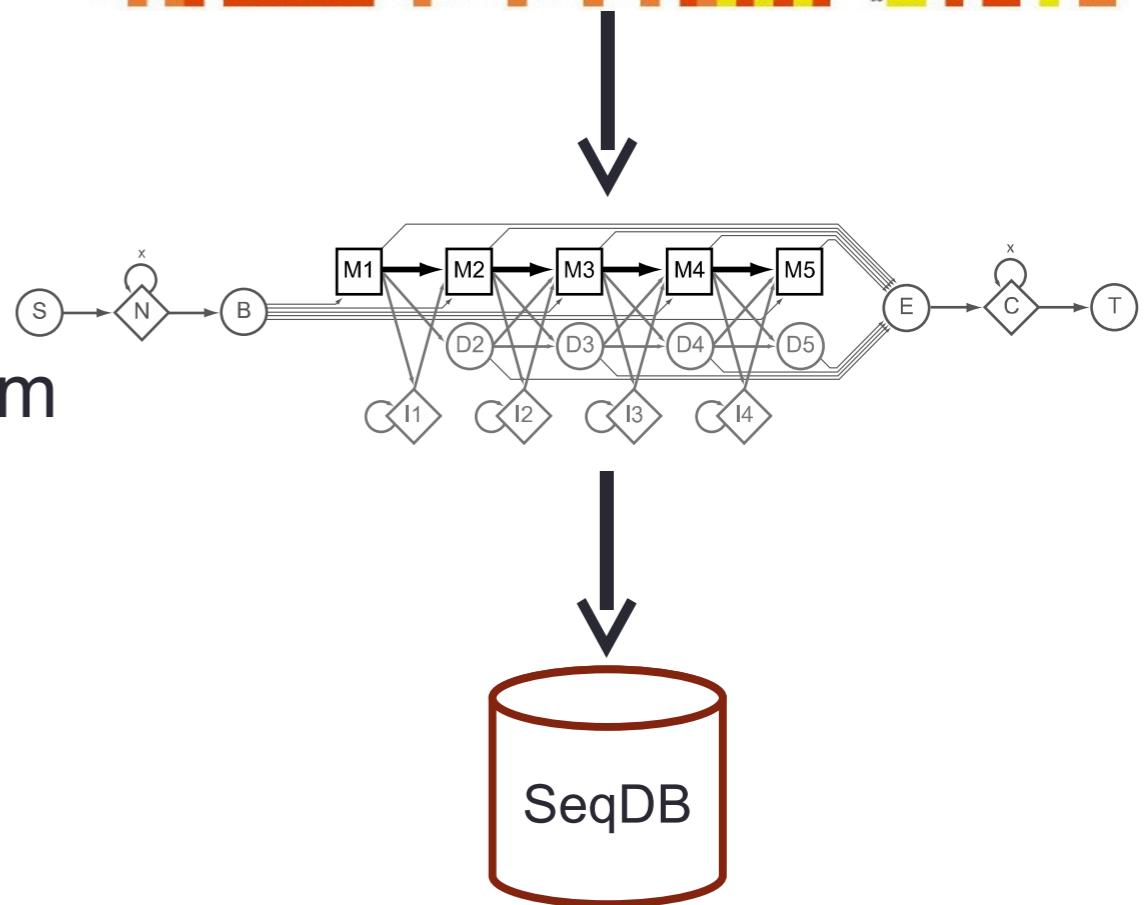


Recall that a PSSM for this motif would give the sequences **WEIRD** and **WEIRH** equally good scores even though the **RH** and **QR** combinations were not observed

Use of HMMER

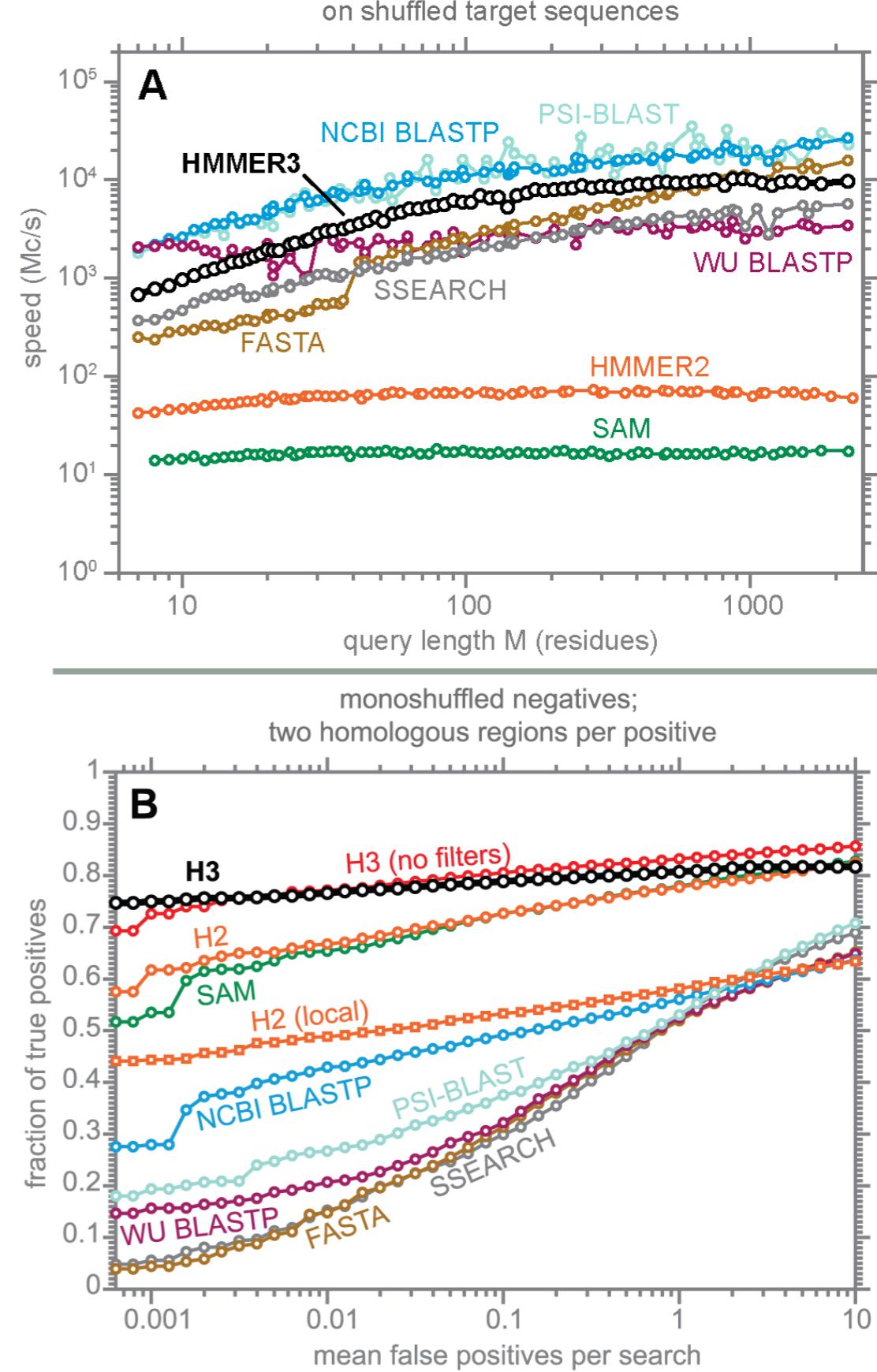
- Widely used by protein family databases
 - Use ‘seed’ alignments
- Until 2010
 - Computationally expensive
 - Restricted to HMMs constructed from multiple sequence alignments
- Command line application

A multiple sequence alignment of a protein family. The sequences are shown as horizontal lines of amino acid residues. Conserved positions are highlighted with orange boxes, while variable positions are in black. The alignment shows a highly conserved N-terminal domain followed by a more variable C-terminal domain.



HMMER vs BLAST

	HMMER	BLAST
Program	<i>PHMMER</i>	<i>BLASTP</i>
Query		Single sequence
Target Database		Sequence database
Program	<i>HMMSCAN</i>	<i>RPSBLAST</i>
Query		Single sequence
Target Database	Profile HMM database, e.g. Pfam	PSSM database, e.g. CDD
Program	<i>HMMSEARCH</i>	<i>PSI-BLAST</i>
Query	Profile HMM	PSSM
Target Database		Sequence database
Program	<i>JACKHMMER</i>	<i>PSI-BLAST</i>
Query		Single sequence
Target Database		Sequence database

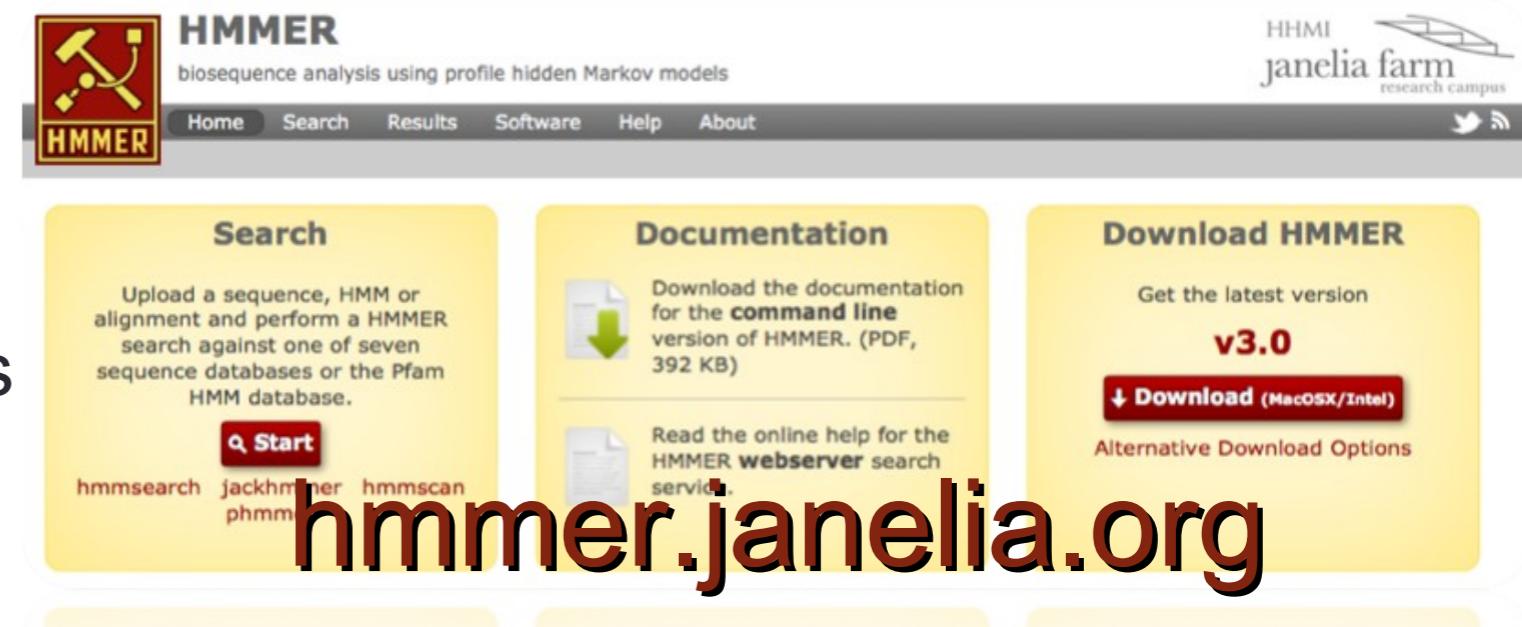


Modified from: S. R. Eddy
PLoS Comp. Biol., 7:e1002195, 2011.



Fast Web Searches

- Parallelized searches across compute farm
 - Average query returns ~1 sec
- Range of sequence databases
 - Large Comprehensive
 - Curated / Structure
 - Metagenomics
 - Representative Proteomes
- Family Annotations
 - Pfam
- Batch and RESTful API
 - Automatic and Human interface



Visualization of Results – By Score

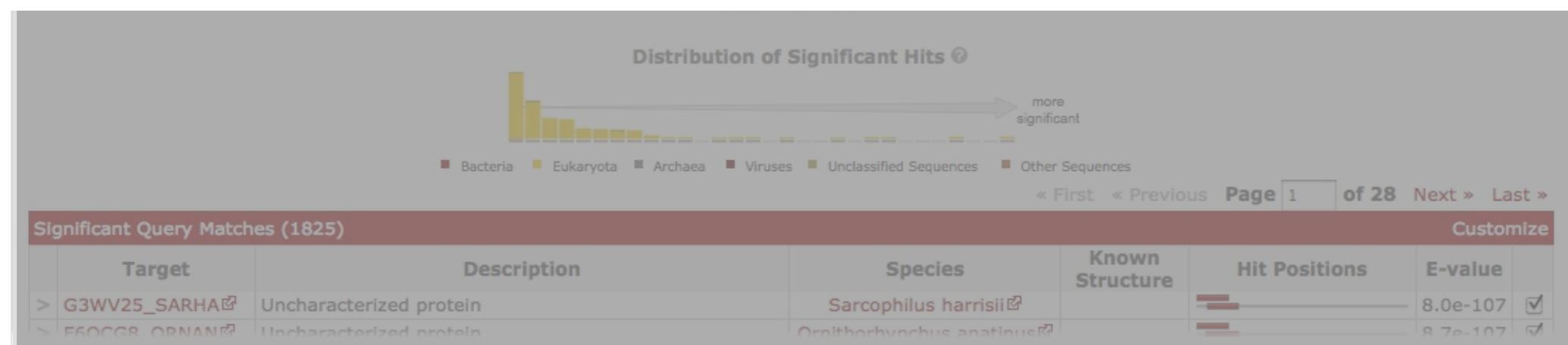
Pfam Domains

Distribution of Significant Hits



- Bacteria
- Eukaryota
- Archaea
- Viruses
- Unclassified Sequences
- Other Sequences

Distribution of Significant Hits 



The screenshot shows a search interface for significant query matches. At the top, there is a histogram of hit distribution and a legend for sequence types. Below the histogram, a table lists "Significant Query Matches (1825)". The table columns include Target, Description, Species, Known Structure, Hit Positions, and E-value. Two entries are shown: G3WV25_SARHA and F6OCCG8_ORNAN. Both targets are uncharacterized proteins from Sarcophilus harrisii and Ornithodoros anatinus respectively. The E-values are 8.0e-107 and 8.7e-107, both with checked boxes.

Significant Query Matches (1825)						Customize
Target	Description	Species	Known Structure	Hit Positions	E-value	
G3WV25_SARHA	Uncharacterized protein	Sarcophilus harrisii			8.0e-107	<input checked="" type="checkbox"/>
F6OCCG8_ORNAN	Uncharacterized protein	Ornithodoros anatinus			8.7e-107	<input checked="" type="checkbox"/>



Visualization of Results – By Score



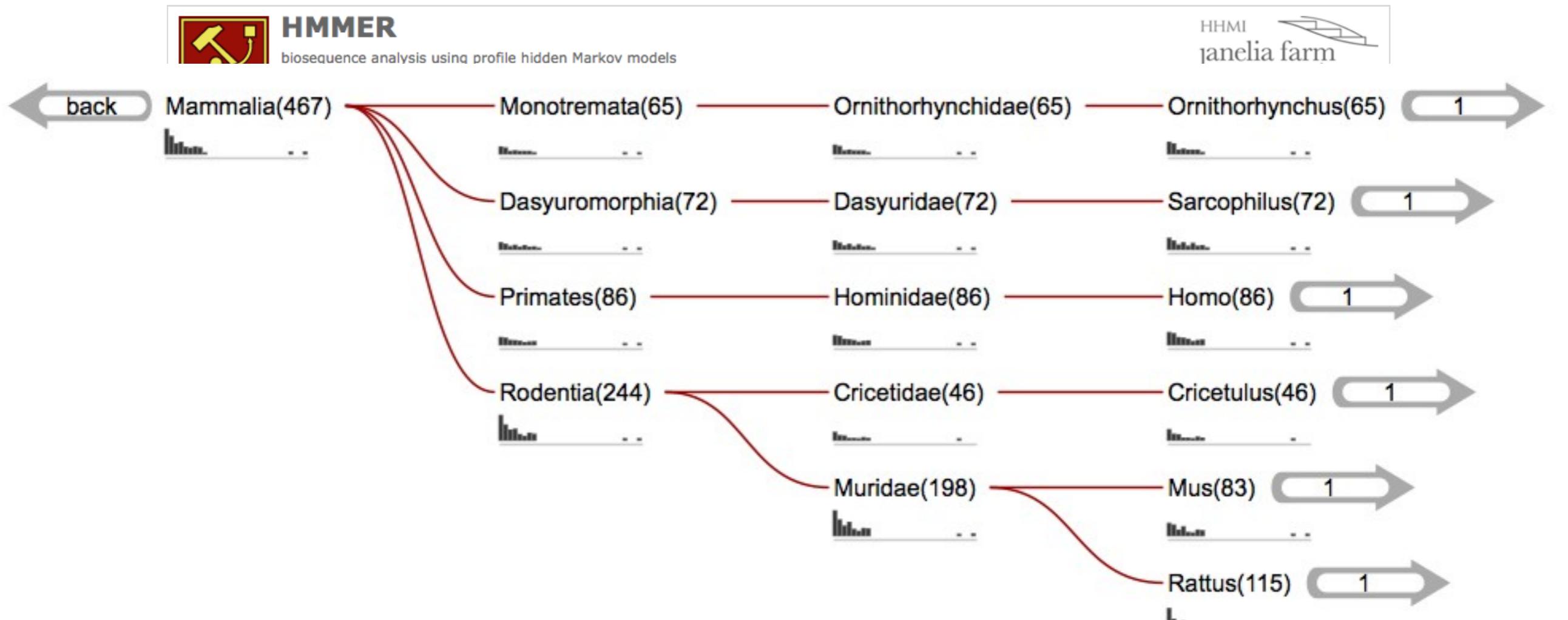
V	Q16IU8_AEDAE	SH2/SH3 adaptor protein	Aedes aegypti						2.5e-31	<input checked="" type="checkbox"/>		
Query	Target Envelope		Target Alignment		Bias	Accuracy	% Identity (count)	% Similarity (count)	Bit Score	E-value		
	start	end	start	end						Ind.	Cond.	
	7	62	4	81	9	63	0.02	0.81	36.4 (20)	50.9 (28)	19.5	0.2 0.00011

.....*.....*.....*.....*.....*.....*.....*.....*.....

Query	7	d p n l f v a l y d f v a s g d n t l s i t k g e k l r v l g y n h n g e w c e a q t . k n g q q w v p s n y i t	62
		d va yd+ a g l + k+e+ +l + w q n g+vpsny+	
Target	9	D V C Y V V A K Y D Y A A Q G A Q E L D L R K N E R Y L L D -- D S K H W W R V Q N s H N Q S G Y V P S N Y V K	63
PP		5566799*****987775..455677766516777*****96	



Visualization of Results – By Taxonomy



Species Distribution

Species	Count	View
Rattus norvegicus	115	Show
Homo sapiens	86	Show
Mus musculus	83	Show
Sarcophilus harrisii	72	Show
Ornithorhynchus anatinus	65	Show
Cricetus griseus	46	Show

Show All Visible

Search Details
Jump to threshold page



Visualization of Results – By Domain

The screenshot shows the HMMER web interface. At the top, there's a logo with a hammer and sickle, followed by the text "HMMER" and "biosequence analysis using profile hidden Markov models". Below this is a navigation bar with links for Home, Search, Results (which is highlighted), Software, Help, and About. The search term "jackhmmer" is entered in the search bar. On the right side, there's a logo for "HHMI janelia farm research campus" with social media links for Twitter and RSS, and a "Search Again" button.

Below the navigation, there are tabs for Score, Taxonomy, Domain (which is also highlighted), and Download. A section labeled "Iteration 1" is visible.

The main content area displays several domain architectures:

- 214 SEQUENCES** with domain architecture: **SH2**, example: [F6Q3Z0_CIOIN](#) [View Scores](#)
Show All
SH3_1 (PF00018.23)
Description: SH3 domain [Pfam] Coordinates: 88 - 135 (alignment region 88 - 135)
- 202 SEQUENCES** with domain architecture: [SH2](#)
with domain architecture: [SH2](#), example: [FYN_HUMAN](#) [View Scores](#)
Show All
SH3_1 (PF00018.23)
Description: SH3 domain [Pfam] Coordinates: 88 - 135 (alignment region 88 - 135)
- 57 SEQUENCES** with domain architecture: [SH2](#)
with domain architecture: [SH2](#), example: [D6W7G8_TRICA](#) [View Scores](#)
Show All
SH3_1 (PF00018.23)
Description: SH3 domain [Pfam] Coordinates: 88 - 135 (alignment region 88 - 135)
- 46 SEQUENCES** with domain architecture: **SH2, SOCS_box**, example: [B3F7U0_ANOGA](#) [View Scores](#)
Show All
SH3_1 (PF00018.23)
Description: SH3 domain [Pfam] Coordinates: 88 - 135 (alignment region 88 - 135)
- 42 SEQUENCES** with domain architecture: **SH3_1, SH2, SH3_1**, example: [A8XPY6_CAEBR](#) [View Scores](#)
Show All



PFAM: Protein Family Database of Profile HMMs

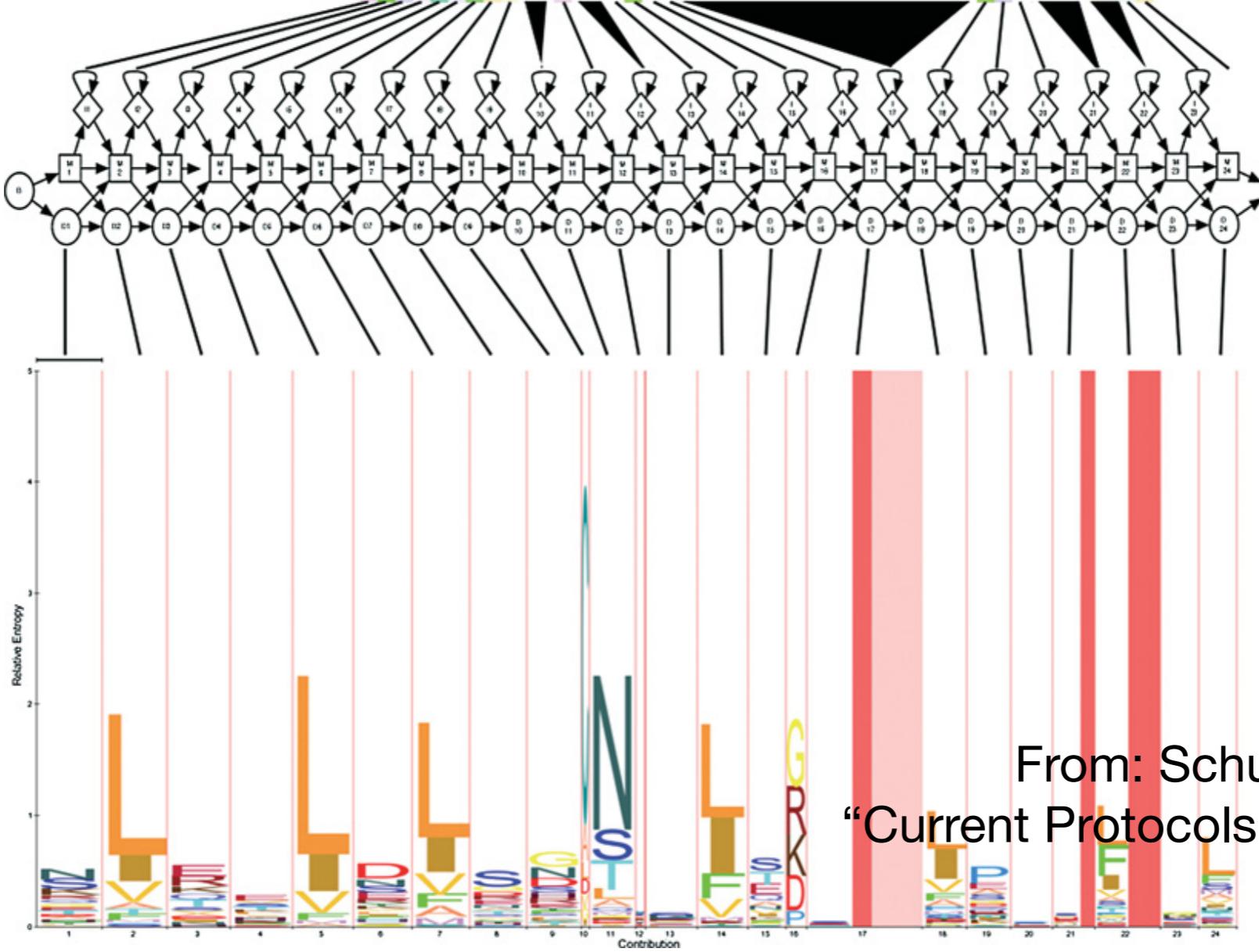
Comprehensive compilation of both multiple sequence alignments and profile HMMs of protein families.

<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**

Q9ARB2_LNUS/823-844	MLEYLDIGRA..P.RIV.H.....	LDG...LENL
Q9M8N0_ARATH/320-341	RLTFLNLNSFC..S.KLT.G.....	LAF...FSII
FLJ_HUMAN/318-339	NLEEFMAAN..N..NLE.L.....	VPES..LCRC
Q9VN74_DROME/90-112	ALHSLVIENC....TIV.H.....	INDAA.FNQE
Q8L8I7_PNTA/792-814	NLQTIQMYRX..E.SLQ.V.....	LPDS..FGNL
Q9FHL8_ARATH/301-324	NLWSLNLSR..N..LFSDP.....	LPVVG.ARGF
SLK6_MOUSE/65-87	RPFHLSSLN..N..GLT.M.....	LHTND.FSGL
Q8NJJ8_EMEN/978-1000	TLTSLNIAS..A..KLV.Q.....	FRDTL.FDSL
Q9LUQ2_ARATH/92-113	AMKSLDVSF..N..SIS.E.....	LPEQ..IGSA
Q9FH93_ARATH/169-188	RLTSLNLDF..N..RFNGT.....	LPS....LN
Q898G0_CLOTE/268-288	YLERINLDK..N..KI.KN.....	IEE...LEAN
Q8H6V2_MAIZE/678-699	NLRILSIVDC..V.SLQ.K.....	LPP...SDSF
Q9AR40_LNUS/692-713	DLKVLDINQ..T..EIT.T.....	LKGE..VESL
Q9LE82_ARATH/350-377	HLTEIYMSY..L..NLEDEGT..	EALSEAL.LKSA
Q9H5N5_HUMAN/255-278	HLQVLDLHQC....SLT.AD.....	DVMSL...TQVI
Q8L4C7_ARATH/185-207	KLEYLDIWG..S..NVT.N.....	QGAVS.ILKF
Q9VSA4_DROME/1115-1138	QLKALRLQC..N..AI.GSH..	GLEAL..LCGQ
TLR1_MOUSE/376-398	RLKTLSSLQK..N..QL.KN.....	LENII.LTSA
Q9TXJ6_LEIMA/445-465	GLRDIDLGH..T..Kvh.N.....	IDA...LQAS
FXL13_MOUSE/409-448	KLIYLDLSGC..T.QVL.VEKCPRISSVVLI	GSPHI.SDSA.FKAL
Q9TXJ6_LEIMA/927-948	ALTVVNANSC..V.NLT.S.....	IEA...LESA
Q9M4X9_CHLRE/1417-1444	LLAVLHLHD..NP.RLA.ADG.....	VAGLAAA..LPGL
Q945S6_LYCPMV/656-677	NLRHLDVSN..T..RRL.K.....	MPLH..LSRL



HMM limitations

HMMs are linear models and are thus **unable to capture higher order correlations** among positions (e.g. distant cysteins in a disulfide bridge, RNA secondary structure pairs, etc).

Another flaw of HMMs lies at the very heart of the mathematical theory behind these models. Namely, that the probability of a sequence can be found from the product of the probabilities of its individual residues.

This claim is only valid if the probability of a residue is independent of the probabilities of its neighbors. In biology, there are frequently **strong dependencies between these probabilities** (e.g. hydrophobic residues clustering at the core of protein domains).

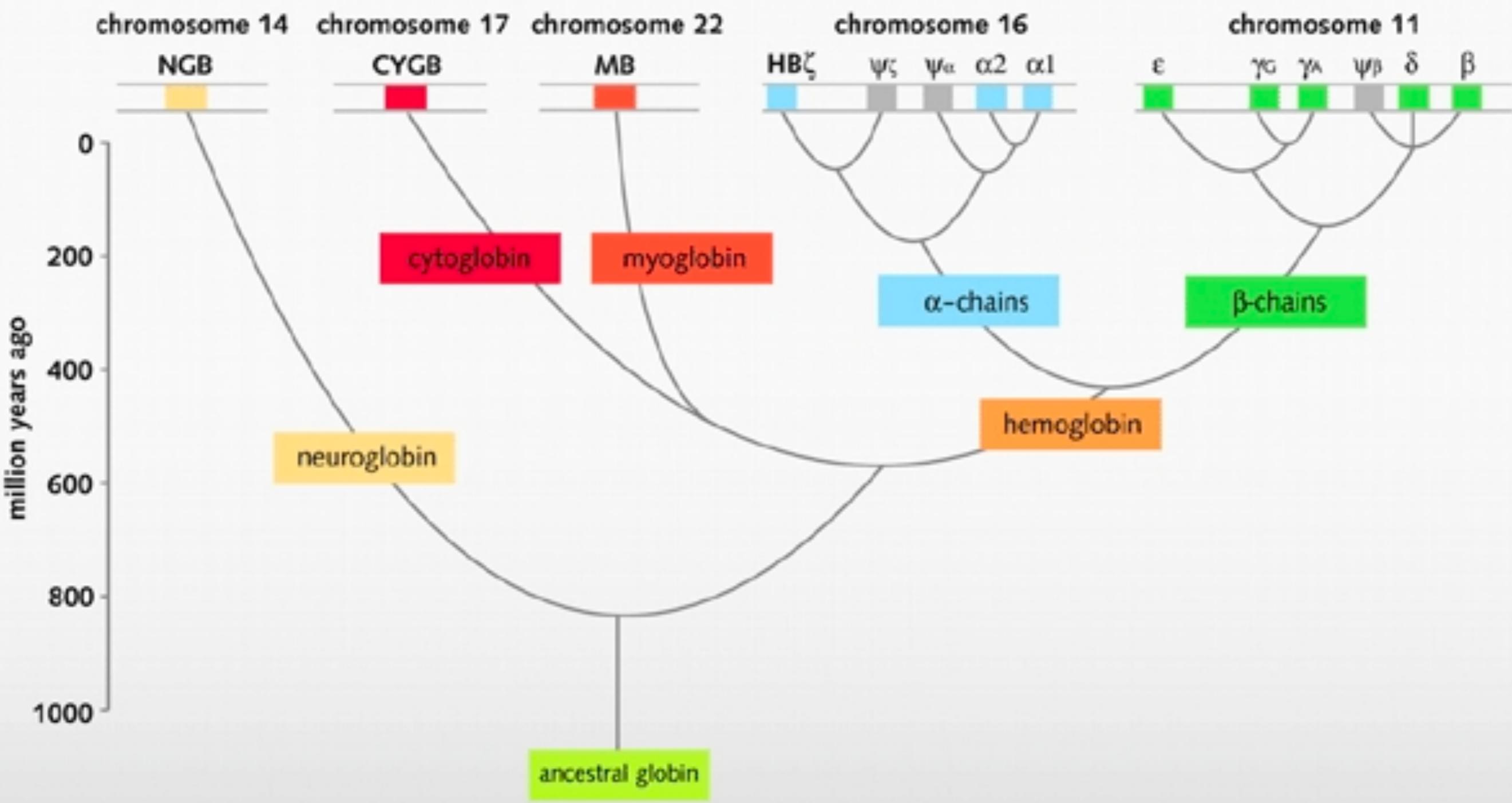
These biological realities have motivated research into new kinds of statistical models. These include hybrids of HMMs and neural nets, dynamic Bayesian nets, factorial HMMs, Boltzmann trees and stochastic context-free grammars.

See: Durbin et al. “Biological Sequence Analysis”



That's it!

Side Note: Human Globins



An evolutionary model of human globins.

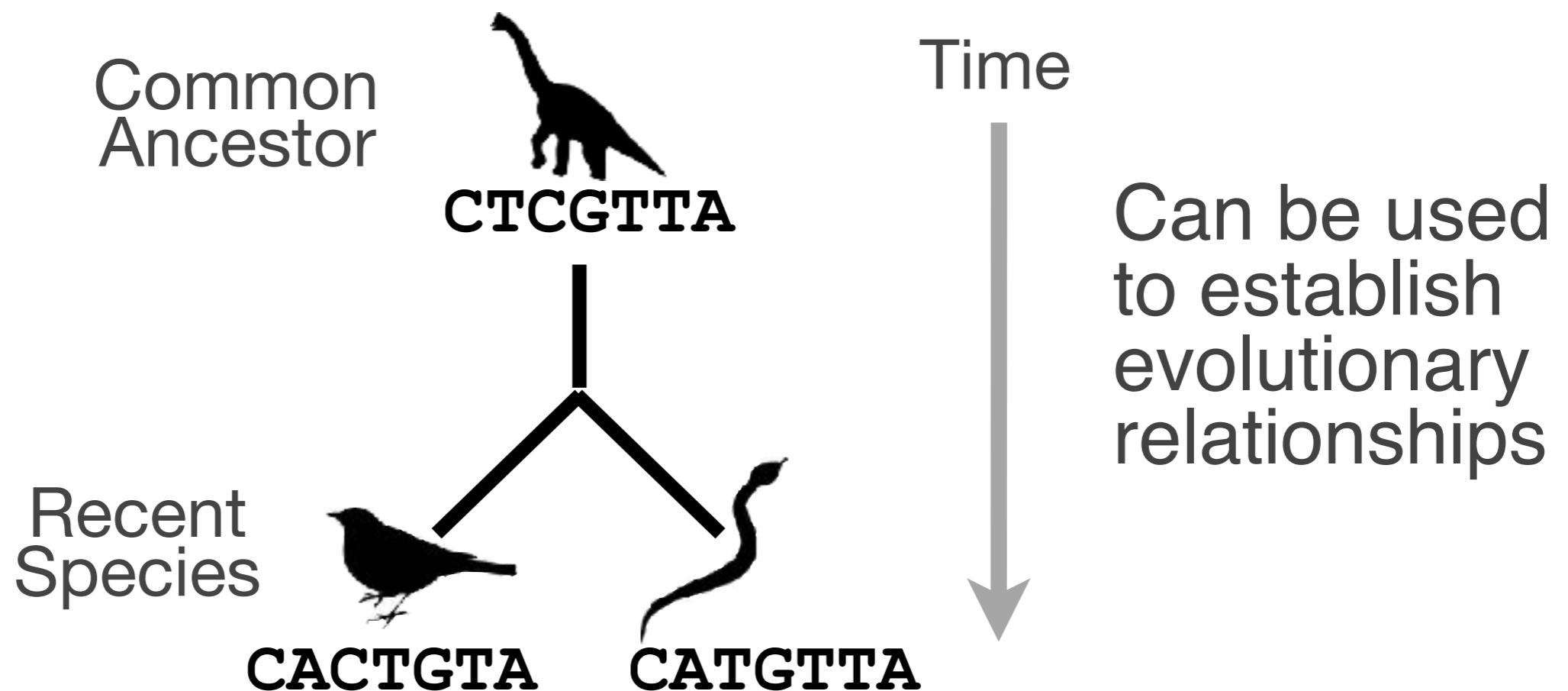
The different locations of globin genes in human chromosomes are reported at the top of the figure, distinguishing between the functional genes (in color) and the pseudogenes (in grey).

Side Note: Orthologs vs Paralogs

Sequence comparison is most informative when it detects homologs

Homologs are sequences that have common origins
i.e. they share a **common ancestor**

- They may or may not have common activity



Key terms

When we talk about related sequences we use specific terminology.

Homologous sequences may be either:

- **Orthologs or Paralogs**

(Note. these are all or nothing relationships!)

Any pair of sequences may share a certain level of:

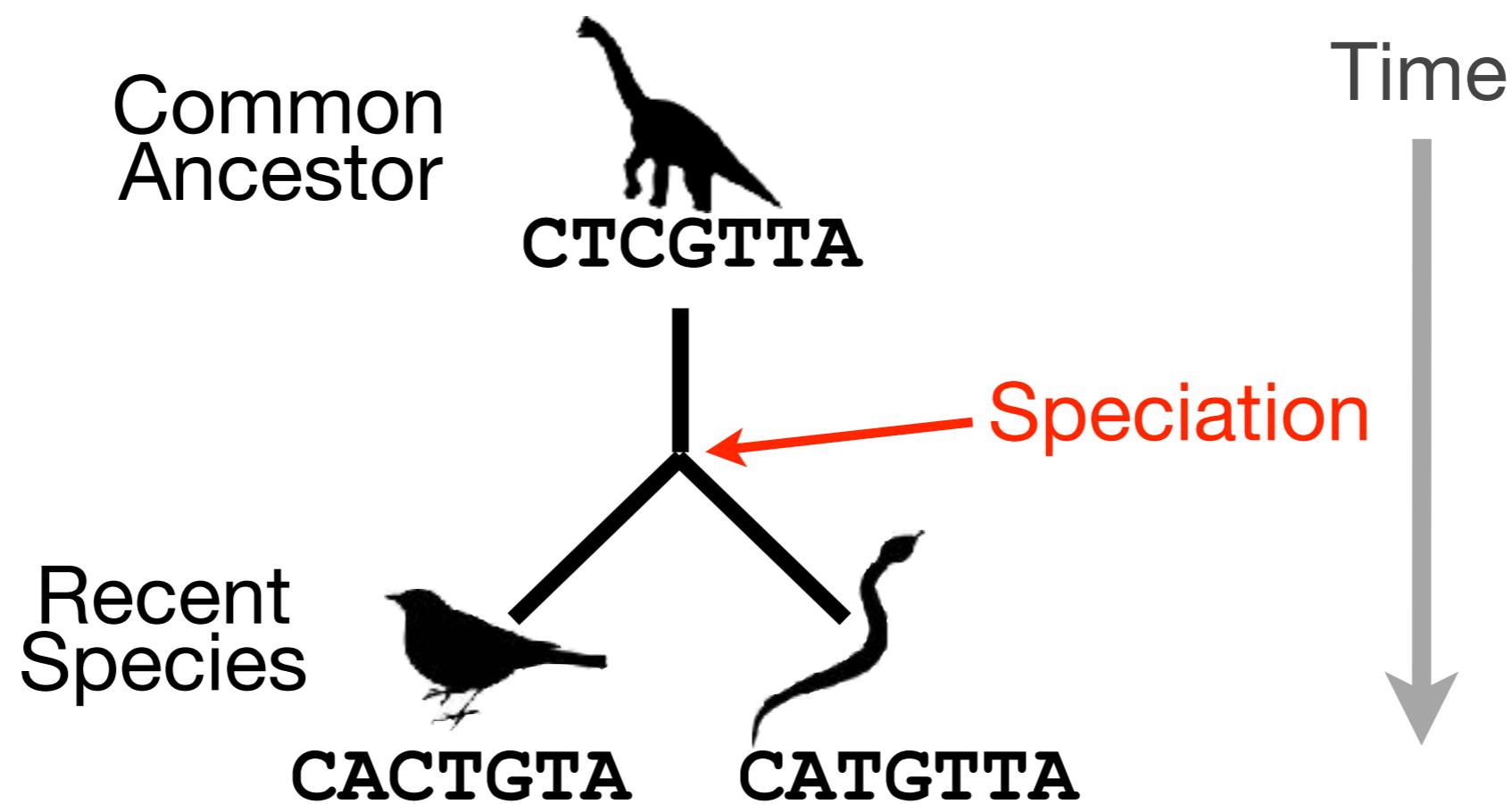
- **Identity and/or Similarity**

(Note. if these metrics are above a certain level we often infer homology)

Orthologs tend to have similar function

Orthologs: are homologs produced by speciation that have diverged due to divergence of the organisms they are associated with.

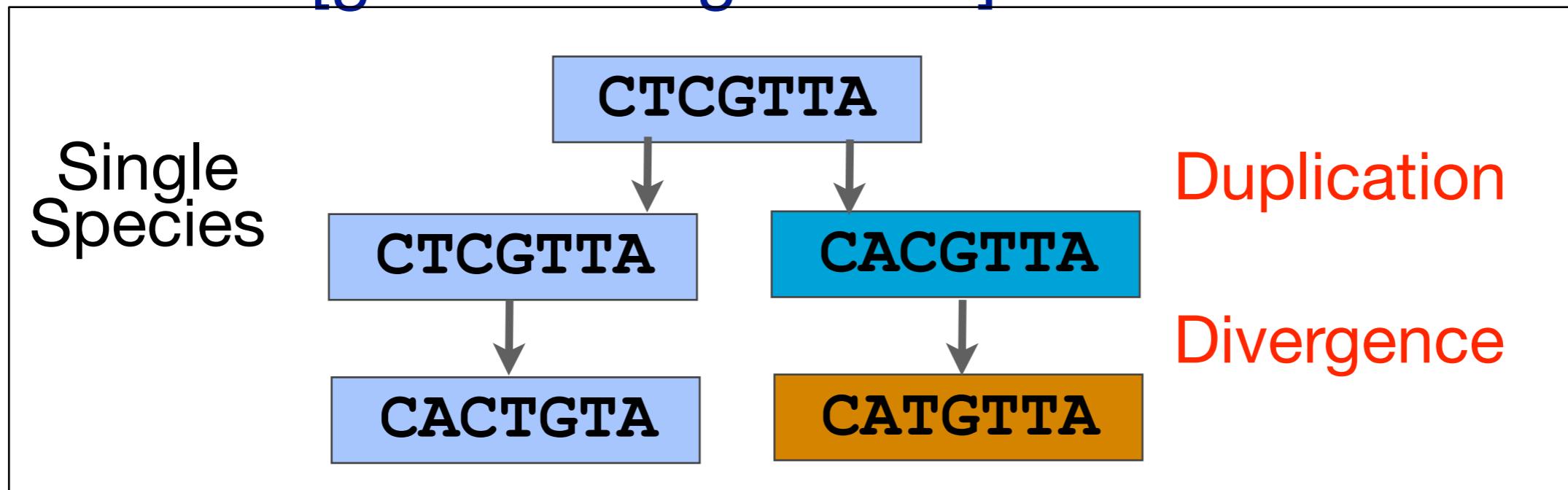
- Ortho = [greek: straight] ... implies direct descent



Paralogs tend to have slightly different functions

Paralogs: are homologs produced by **gene duplication**. They represent genes derived from a common ancestral gene that duplicated within an organism and then subsequently diverged by accumulated mutation.

- Para = [greek: along side of]



Orthologs vs Paralogs

- In practice, determining ortholog vs paralog can be a complex problem:
 - gene loss after duplication,
 - lack of knowledge of evolutionary history,
 - weak similarity because of evolutionary distance
- Homology does not necessarily imply exact same function
 - may have similar function at very crude level but play a different physiological role