

Functional Annotation

Lesson 3

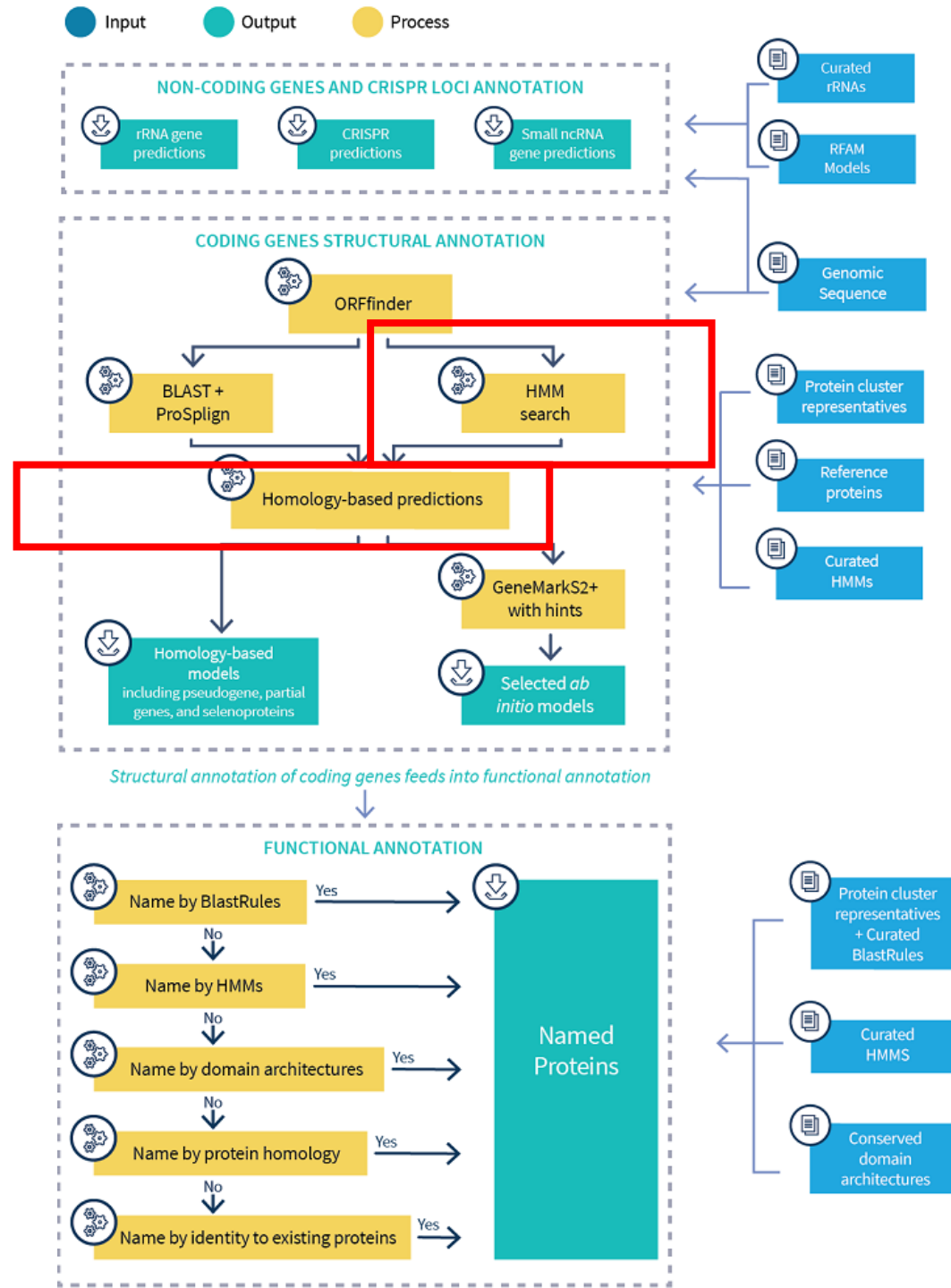
Inferring function from position-sensitive models

- HMMER



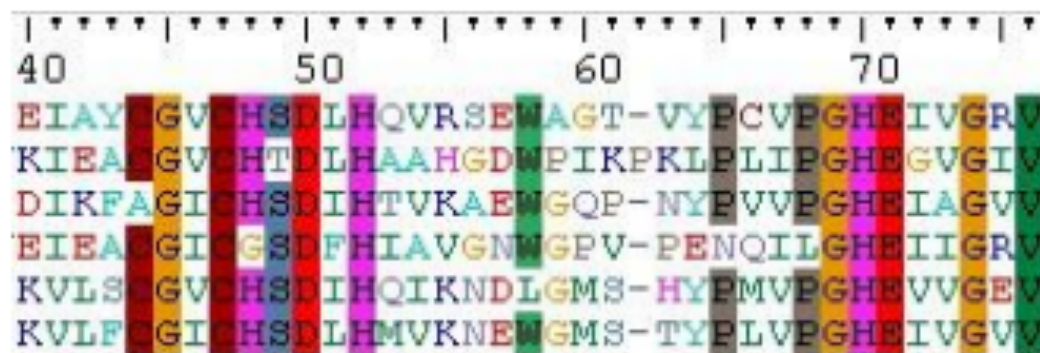
NCBI PGAP

Prokaryotic Genome Annotation Pipeline



The basis for the prediction of features is nearly always a sequence alignment

Based on experimentally verified sequence annotations, a multiple sequence alignment is constructed



Different methods exist to capture the information gained from this multiple sequence alignment

PROBABILISTIC APPROACHES to Homology

- BLAST is computationally greedy, alternative annotation options can use a probabilistic approach
- Probabilistic approach – incorporate random variables and probability distributions into the model of an event or phenomenon
- hidden Markov model - HMM
- HMMs require building a profile based on a training data set
- Like other “machine learning” approaches.
- More Sensitive

HMM starts with a Multiple Sequence Alignment

CLUSTAL 2.0.12 multiple sequence alignment

```

sp|088479|F0S_MESAU      MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNAQDFCTDLVSSANF 60
sp|Q56TN0|F0S_PHORO      MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNAQDFCADL SVSSANF 60
sp|077628|F0S_BOVIN      MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNAQDYCTDLAVSSANF 60
sp|Q8HZP6|F0S_FELCA      MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNAQDFCTDLAVSSANF 60
sp|P01100|F0S_HUMAN      MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNAQDFCTDLAVSSANF 60
sp|P12841|F0S_RAT        MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNTQDFCADL SVSSANF 60
sp|P01102|F0S_MSVFB      MMFSGFNADYEASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNTQDFCADL SVSSANF 60
sp|P11939|F0S_CHICK      MMYQGFAGEYEAASSRCSSASPAAGDSL YYHSPADSFSSMGSPVNSQDFCTDLAVSSANF 60
sp|P53539|F0S_B_HUMAN    -MFQAFP GDYDSGS-RCSS-SPSAESQ--YLSSVDSFGSPPTAAASQE-CAGLGEMPGSF 54
sp|Q9TUB3|F0S_B_CANFA    -MFQAFP GDYDSGS-RCSS-SPSAESQ--YLSSVDSFGSPPTAAASQE-CAGLGEMPGSF 54
sp|P13346|F0S_B_MOUSE    -MFQAFP GDYDSGS-RCSS-SPSAESQ--YLSSVDSFGSPPTAAASQE-CAGLGEMPGSF 54

```

```

sp|088479|F0S_MESAU      IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGVPTPS-----TGAYSR 108
sp|Q56TN0|F0S_PHORO      IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGVPTPS-----TGAYSR 108
sp|077628|F0S_BOVIN      IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGVPTPS-----AGAYSR 108
sp|Q8HZP6|F0S_FELCA      IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGVPTPS-----AGAYSR 108
sp|P01100|F0S_HUMAN      IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGVPTPS-----AGAYSR 108
sp|P12841|F0S_RAT        IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGLPTPS-----TGAYAR 108
sp|P01102|F0S_MSVFB      IPTVTAI STSPDLQWL VQPTLVSSV PS-----QTRAPHYGLPTQS-----AGAYAR 108
sp|P11939|F0S_CHICK      VPTVTAI STSPDLQWL VQPTLVSSV PS-----QNRG-HPYGVPA PAP-PAAYSR 108
sp|P53539|F0S_B_HUMAN    VPTVTAI TTSQLQWL VQPTLVSSV PS-----QSGQGLASQPPVDPYDMPGTSYSTPGMSGYSS 114
sp|Q9TUB3|F0S_B_CANFA    VPTVTAI TTSQLQWL VQPTLVSSV PS-----QSGQGLASQPPVDPYDMPGTSYSTPGMSGYSS 114
sp|P13346|F0S_B_MOUSE    VPTVTAI TTSQLQWL VQPTLVSSV PS-----QSGQGLASQPPVDPYDMPGTSYSTGLSAYST 114

```

```

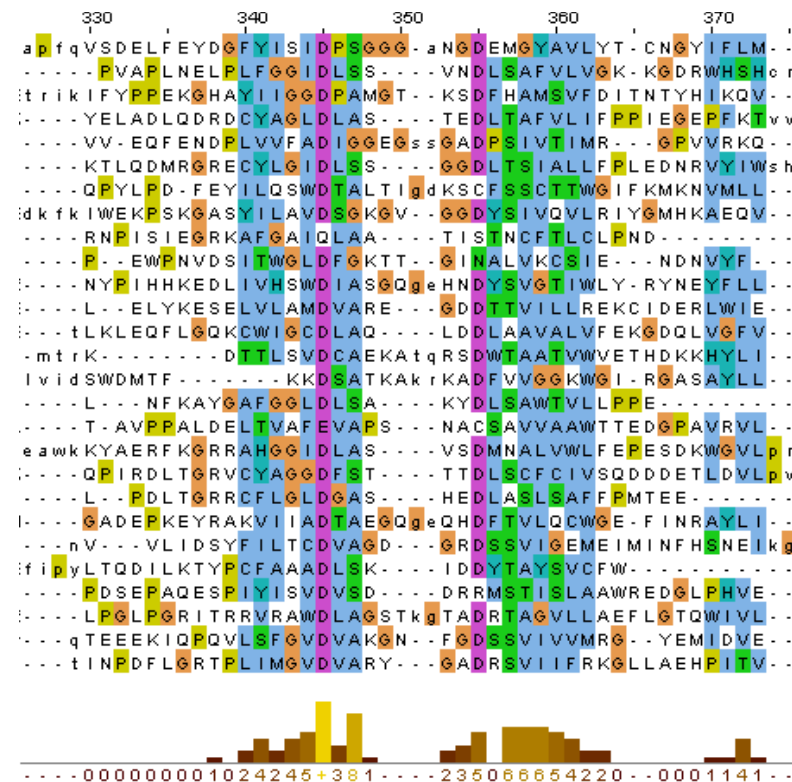
sp|088479|F0S_MESAU      -----AGMVKTVSGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 56
sp|Q56TN0|F0S_PHORO      -----AGMVKTVSGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 56
sp|077628|F0S_BOVIN      -----AGMVKTMGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 56
sp|Q8HZP6|F0S_FELCA      -----AGVVKTVTAGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 57
sp|P01100|F0S_HUMAN      -----AGVVKTMGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 56
sp|P12841|F0S_RAT        -----AGVVKTMGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 56
sp|P01102|F0S_MSVFB      -----AEMVKTVSGG---RAQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 56
sp|P11939|F0S_CHICK      -----PAVLK-APGG---RGQSI GRRGKVEQSP EEEEEKRRIRRRNKMAAAKCRN 55
sp|P53539|F0S_B_HUMAN    GGASGSGGPSTSGTSGP GPAPARARPRRPREETLTP EEEEEKRRVRRRNKLA AAKCRN 74
sp|Q9TUB3|F0S_B_CANFA    GGASGSGGPSTSGTSGP GPAPARARLRPRPREETLTP EEEEEKRRVRRRNKLA AAKCRN 74
sp|P13346|F0S_B_MOUSE    GGASGSGGPSTSTTSGP VARPAPARPRRPREETLTP EEEEEKRRVRRRNKLA AAKCRN 74

```

```

sp|088479|F0S_MESAU      RRRELDTLQ AETDQLEDEK SALQTEIANLLKEKEKLEFILA AHRPA CKIPDDL GFPEEM 216
sp|Q56TN0|F0S_PHORO      RRRELDTLQ AETDQLEDEK SALQTEIANLLKEKEKLEFILA AHRPA CKIPDDL GFPEEM 216
sp|077628|F0S_BOVIN      RRRELDTLQ AETDQLEDEK SALQTEIANLLKEKEKLEFILA AHRPA CKIPDDL GFPEEM 216
sp|Q8HZP6|F0S_FELCA      RRRELDTLQ AETDQLEDEK SALQTEIANLLKEKEKLEFILA AHRPA CKIPDDL GFPEEM 217
sp|P01100|F0S_HUMAN      RRRRLDTLQ AETDQLEDEK SALQTEIANLLKEKEKLEFILA AHRPA CKIPDDL GFPEEM 216

```



- Use structural and mechanistic information (e.g., catalytic sites)
- More sensitive than pairwise – detect distant relationships

Protein sequences can consist of structurally different parts

EXAMPLE

Domain

part of the tertiary structure of a protein that can exist, function and evolve independently of the rest, linked to a certain biological function

ATP-binding domain

Motif

part (not necessarily contiguous) of the primary structure of a protein that corresponds to the signature of a biological function. Can be associated with a domain.

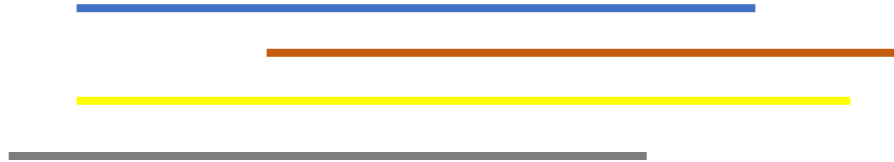
Heme motif CHXXH

Feature

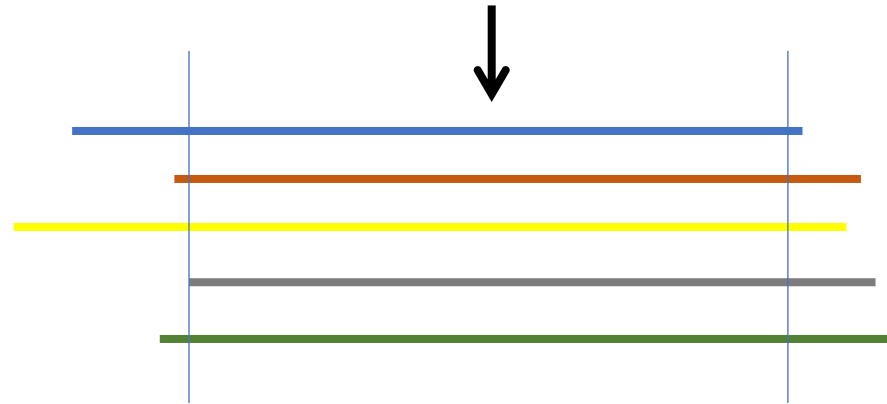
part of the sequence for which some annotation has been added. Some features correspond to domain or motif assignments.

MSA to HMM or PSSM profile

Collect “seed” proteins



Generate & Trim
Alignment



Region of good alignment and closest similarity

Generate Profile with
HMM or PSSM


*Compute statistical probabilities for amino
acid patterns in the seed*

Search New Model
against all proteins

*Choose “noise” and “**trusted**” cutoff scores based
on “known” versus “unknown” protein scores*

Frequency matrices or profiles include the chance of observing the residues


For every position of a motif, a list of all amino acids is made with their frequency. Position-specific weight/scoring matrix or profile. More sensitive way.

		Profile							
		Position:	1.	2.	3.	4.	5.	6.	
123456	ATPKAE		A	0.625	0	0	1/8	6/8	3/8
	KKPKAA		D	0	0	0	0	0	1/8
	AKPKAK		E	0	0	0	0	0	1/8
	TKPKPA		K	0.25	6/8	0	7/8	0	2/8
	AKPKT-		L	0	1/8	0	0	0	0
	AKPAAK		P	0	0	1	0	1/8	0
	KLPKAD		T	1/8	1/8	0	0	1/8	0
	AKPKAA		-	0	0	0	0	0	1/8
Consensus:	AKPKA-	Sum	1	1	1	1	1	1	
? Query: AKPKTE									
? Query: KKPETE									
? Query: TLPATE									

Example: <http://expasy.org/prosite/PS51092>

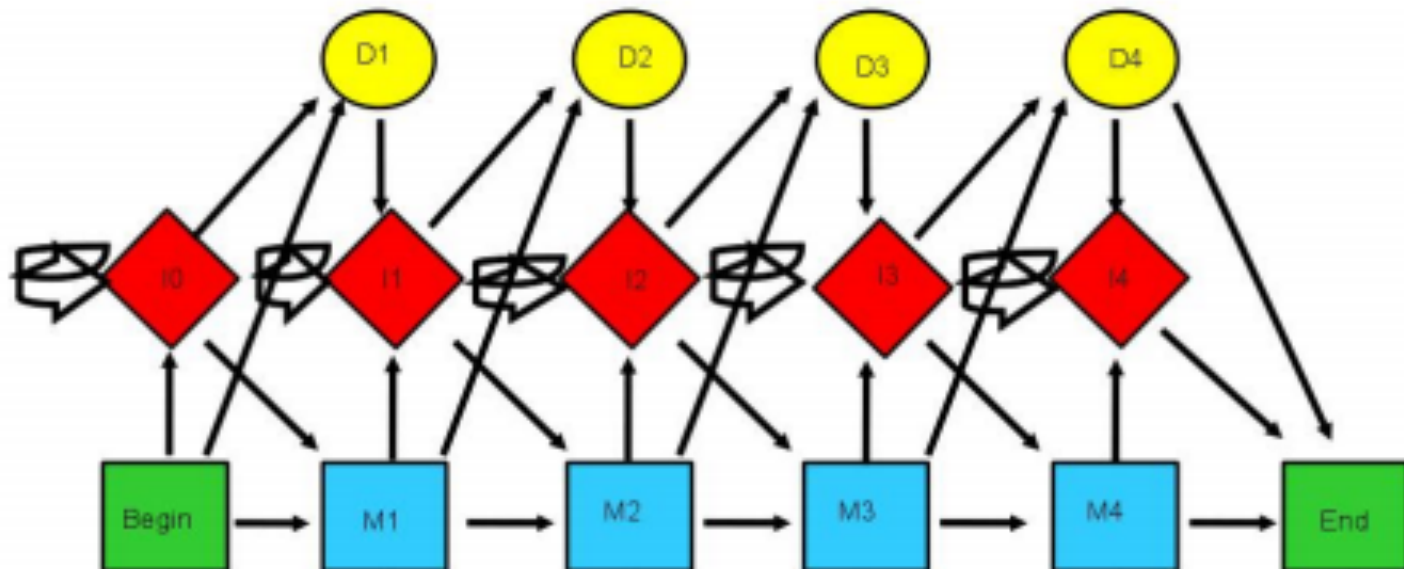
<http://prosite.expasy.org/prosuser.html#meth2>

How good a sequence matches a profile is reported with a score

		PSWM: scores							
	123456	Position:	1.	2.	3.	4.	5.	6.	
	ATPKAE		A	2.377	-2.358	-2.358	0.257	2.631	1.676
	KKPKAA		D	-2.358	-2.358	-2.358	-2.358	-2.358	0.257
	AKPKAK		E	-2.358	-2.358	-2.358	-2.358	-2.358	0.257
	TKPKPA		K	1.134	2.631	-2.358	2.847	-2.358	1.134
	AKPKT-		L	-2.358	0.257	-2.358	-2.358	-2.358	-2.358
	AKPAAK		P	-2.358	-2.358	0.257	-2.358	0.257	-2.358
	KLPKAD		T	0.257	0.257	-2.358	-2.358	0.257	-2.358
	AKPKAA								
Consensus:			AKPKA-						
? Query: AKPKTE			Score = 11.4						
? Query: KKPETE			Score = 5.0						
? Query: TLPATE			Score = 4.3						

A hidden Markov Model takes also into account the gaps in an alignment

The schematic representation of a HMM



HMMER

biosequence analysis using profile hidden Markov models

<http://www.myoops.org/twocw/mit/NR/rdonlyres/Electrical-Engineering>

MSA to HMM profile using HMMER

Input: Query Sequence Set

...SKEAEYLVKQLNTVME...
...SKEAKYLIQQLDTVMK...
...SKERYAAISMFMK...
...AKEGEYLYSNMLNAVMK...

Multiple Alignment

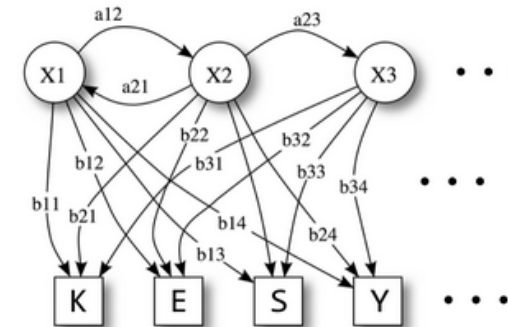
...SKEAEYLVK-QLNTVME...
...SKEAKYLIQ-QLDTVMK...
...SKERYAA----ISMFMK...
...AKEGEYLYSNMLNAVMK...



Input: Target Sequence Set

...CMSDKPDLSEVETFDKSKLTIQQEKEYNQRS...
...SCALEEHV**SKEAEYLVKMLNAVMKV**TGSFDP...
...DRSQNPPQSKGCCFVTFYTRKAALQAQNALH...
...KMPKDKERSLNPAAAQRKLDKQKSLKKGKAE...
...

hmmsearch



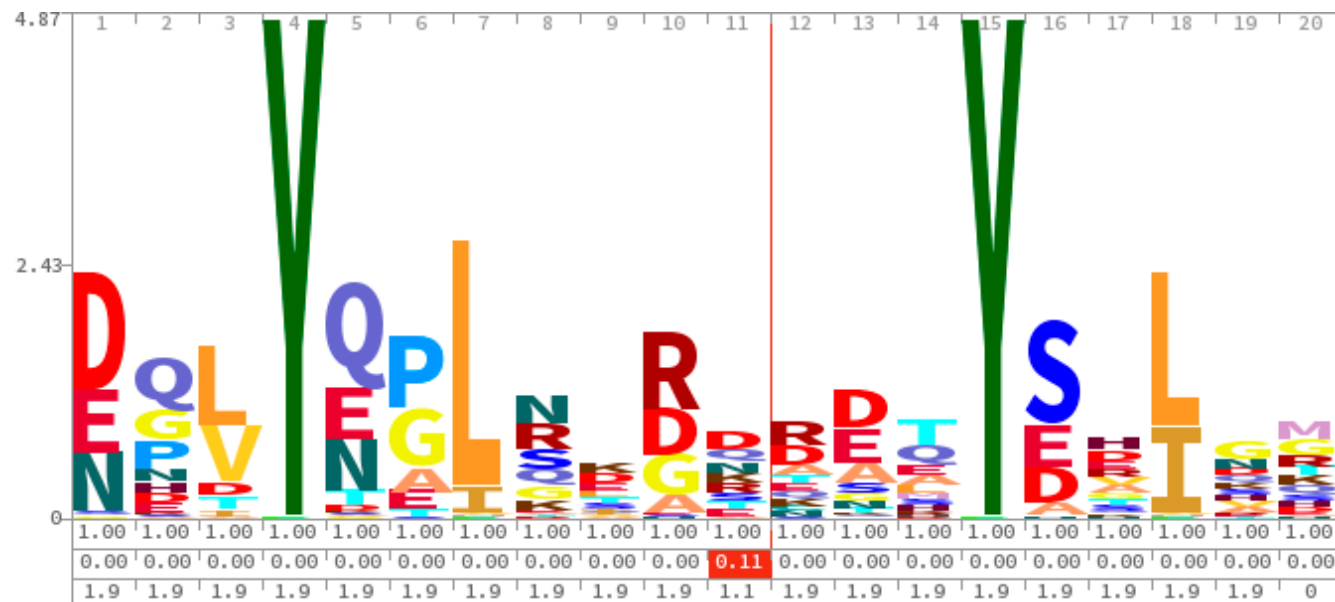
HMM Profile

SKEAEYLVKMLNAVMKV

Output: Resulting Match

HMMER

HMMER is used for searching sequence databases for sequence homologs and for making sequence alignments.



Program	Task
hmmalign	build an MSA from sequences
hmmbuild	convert MSA to HMM
phmmer	single protein vs protein DB (like-BLAST)
hmmsearch	protein sequence vs HMM database
hmmsearch	HMM database vs protein sequence
hmmconvert	convert HMM profile formats
hmmfetch	grab a single HMM from a database
hmmcompress	compress HMM into binary format for hmmsearch

Difference between hmmscan and hmmsearch

- hmmscan and hmmsearch are doing exactly the same compute
 - comparing one profile to one sequence at a time
 - bit score results are identical
 - save both in tabular output files
- hmmscan needs to read in all HMM profiles in a database for each query protein being searched → compute time \ll I/O time
 - hmmpress solves some of this issue
- hmmsearch reads in the HMM profile database so the search is now per query protein

PSI-BLAST & RPS-BLAST (NCBI)

PSI-BLAST

Position-Specific Iterated BLAST

- finds sequences significantly similar to the query in a database search
- Significant matches are used to make an alignment
- MSA used to build a Position-Specific Score Matrix (PSSM) for the query
- The new PSSM is searched against the same database again to pull in more significant hits based on conserved features
- Can be used to further refine the scoring model

RPS-BLAST

Reverse Position-Specific BLAST

- query sequence to search a database of pre-calculated PSSMs
 - reports significant hits in a single pass
- The PSSM has changed from "query" to "subject"
- Commonly used to search the NCBI CDD (Conserved Domain Database)

Lesson 4

Reference databases – series of lectures and examples about databases with a variety of functions

- **Pfam**
- **KEGG**
- **eggNOG**

Reference databases – lectures and examples about databases for specific functions

- **CAZy**
- **anitSMASH**