

HIDDEN MARKOV MODELS (HMM) for detection of intrinsically disordered regions of proteins

Alina Skrylnik

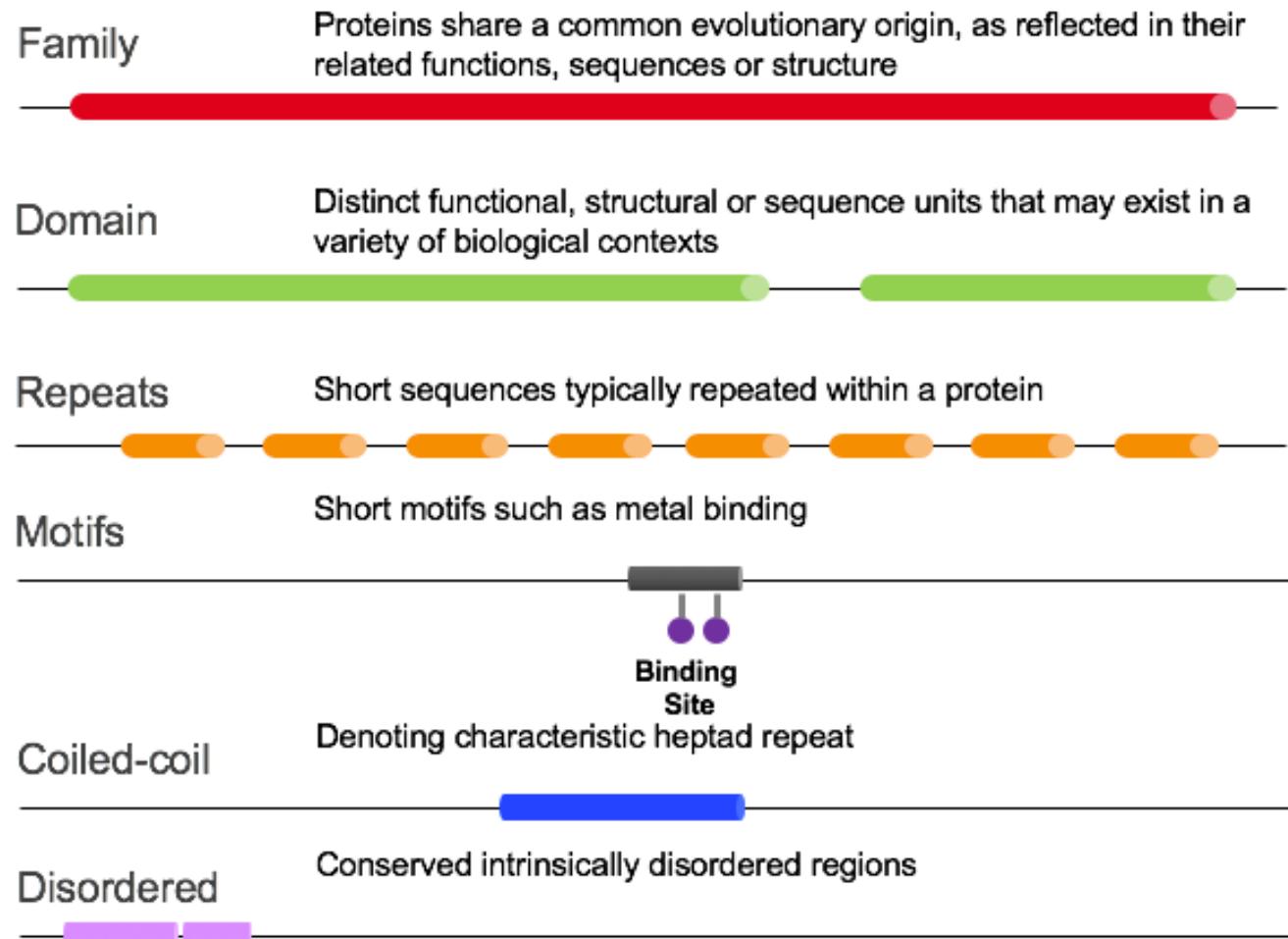
MSc. Data Science (Biological Data Analytics)

2023-2024

Domains are distinct functional and structural units in a protein.

Intrinsically disordered regions (**IDRs**) are considered to be a domain type.

The prediction of Pfam entries is based on the profile HMMs



Profile HMM (pHMM):

A probabilistic model that encapsulates the evolutionary changes that have occurred in a set of related sequences.

Start with a multiple sequence alignment



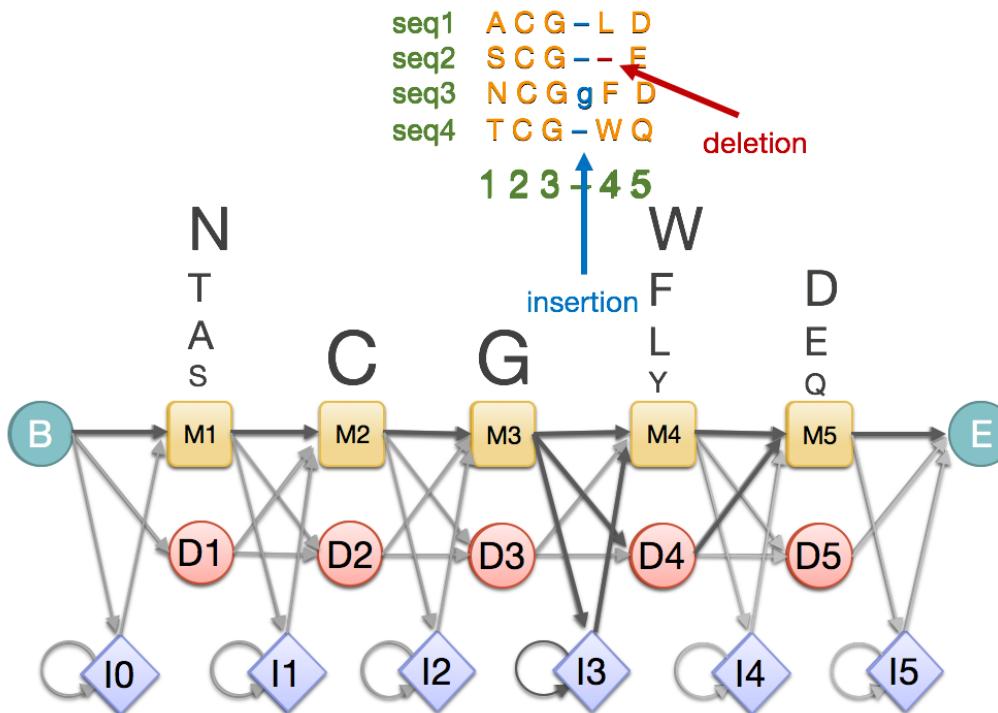
Insertions / deletions can be modelled



Occupancy and amino acid frequency at each position in the alignment are encoded



Profile created



Hidden states: IDRs

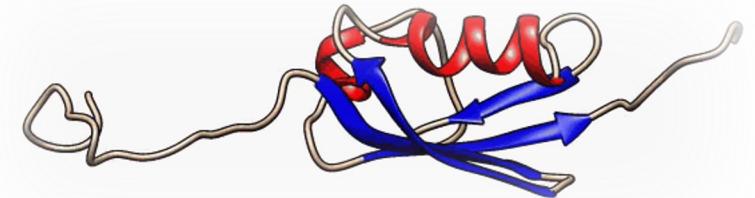
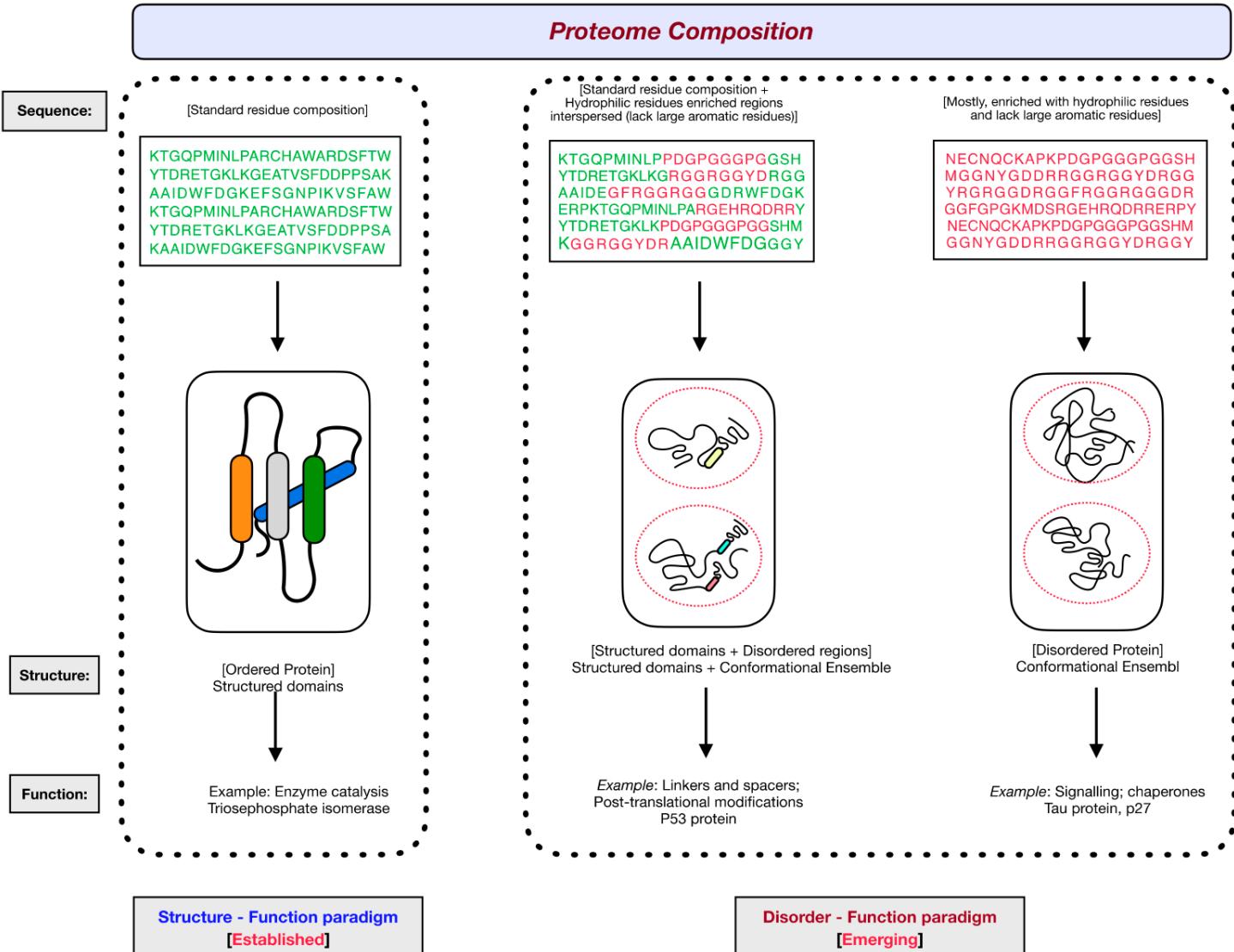
Observable states:
amino acid sequences

M: match state

I: insertion state

D: deletion state

Structure vs disordered paradigms

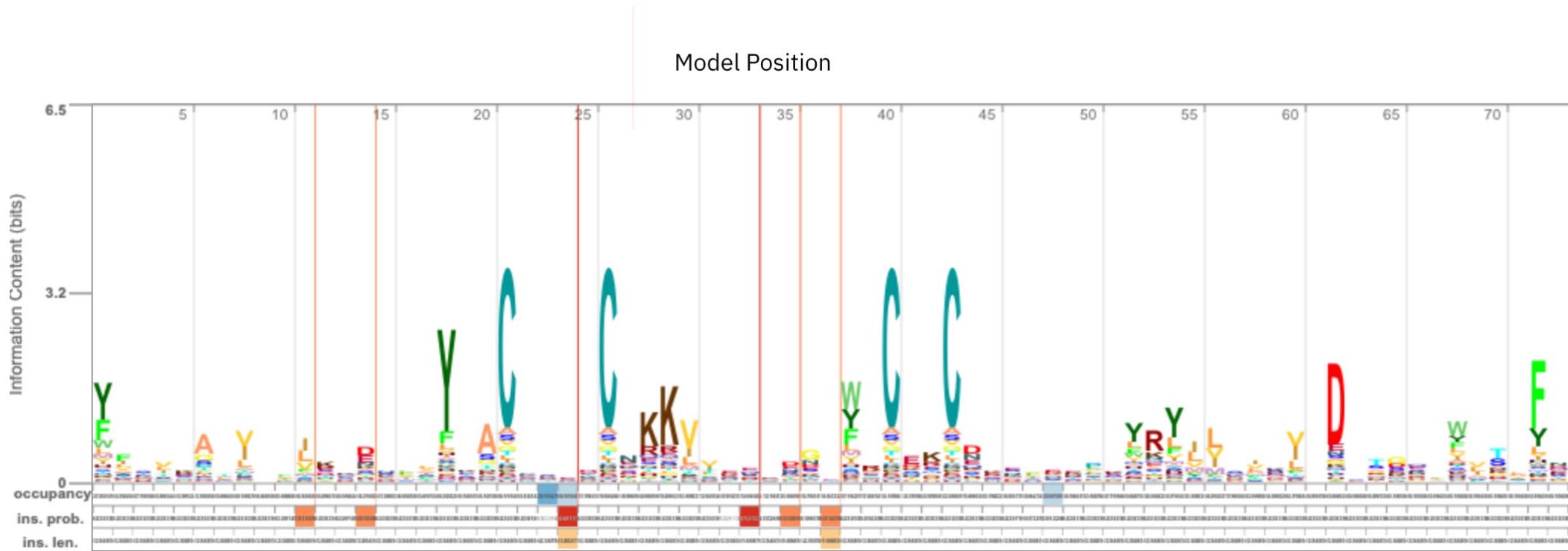


Intrinsically disordered regions (IDRs): polypeptide segments of the protein lacking 3D structure

Main functions: signalling and regulatory

Goal of the thesis:

analysis of HMMs developed for detecting
intrinsically disordered regions (IDRs) of the proteins

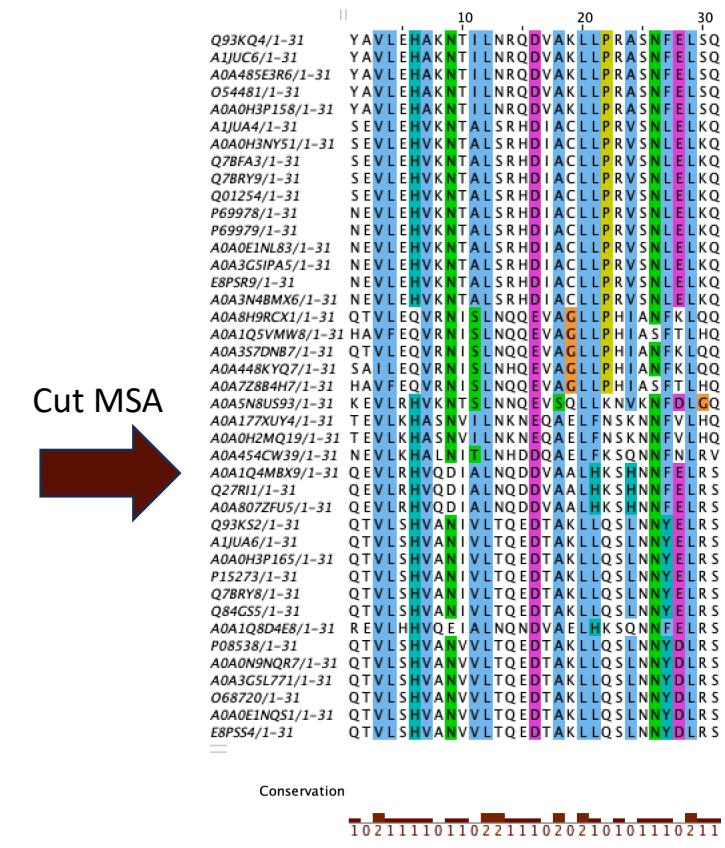
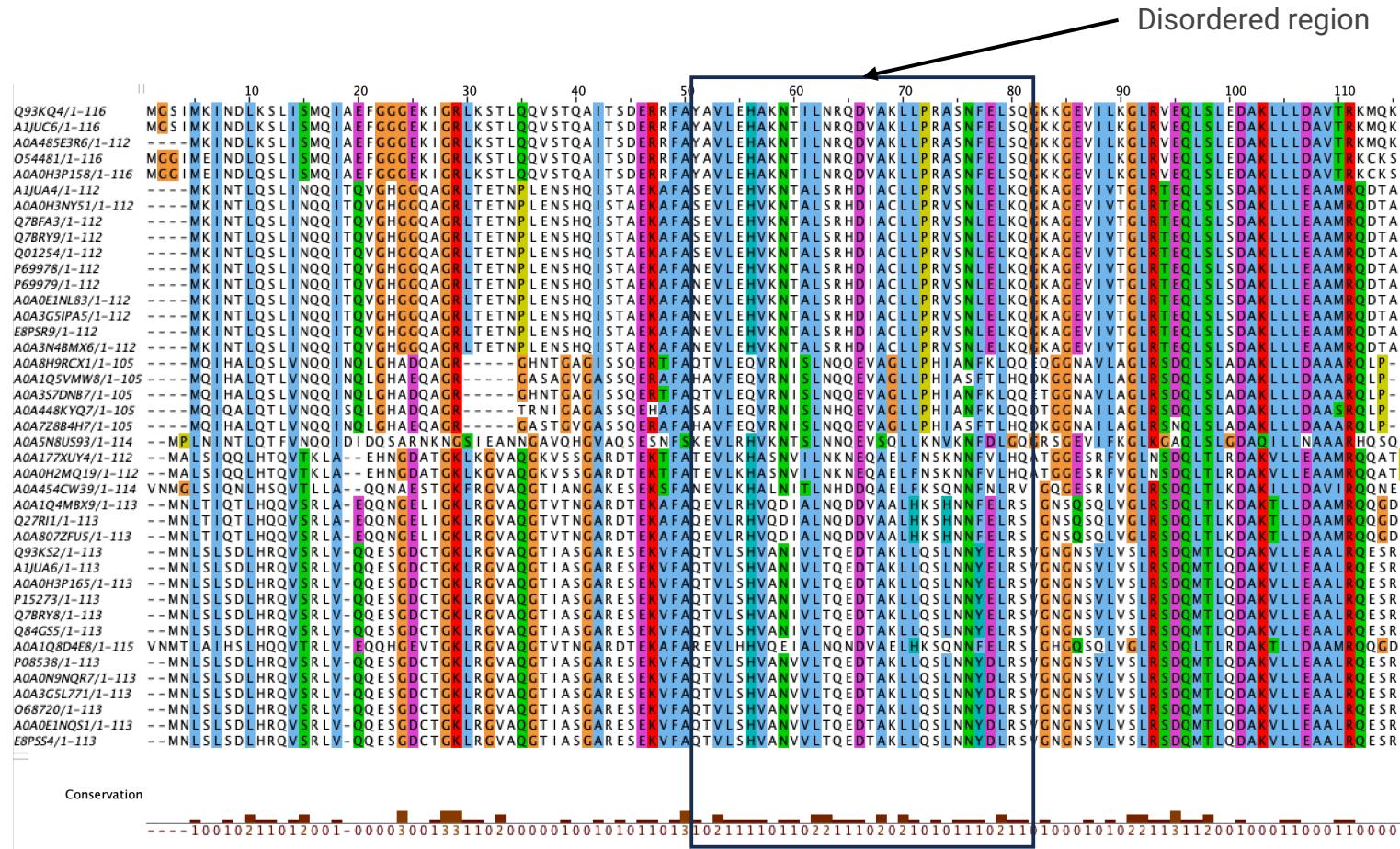


pHMM example: PF08646 - Replication factor-A C terminal domain

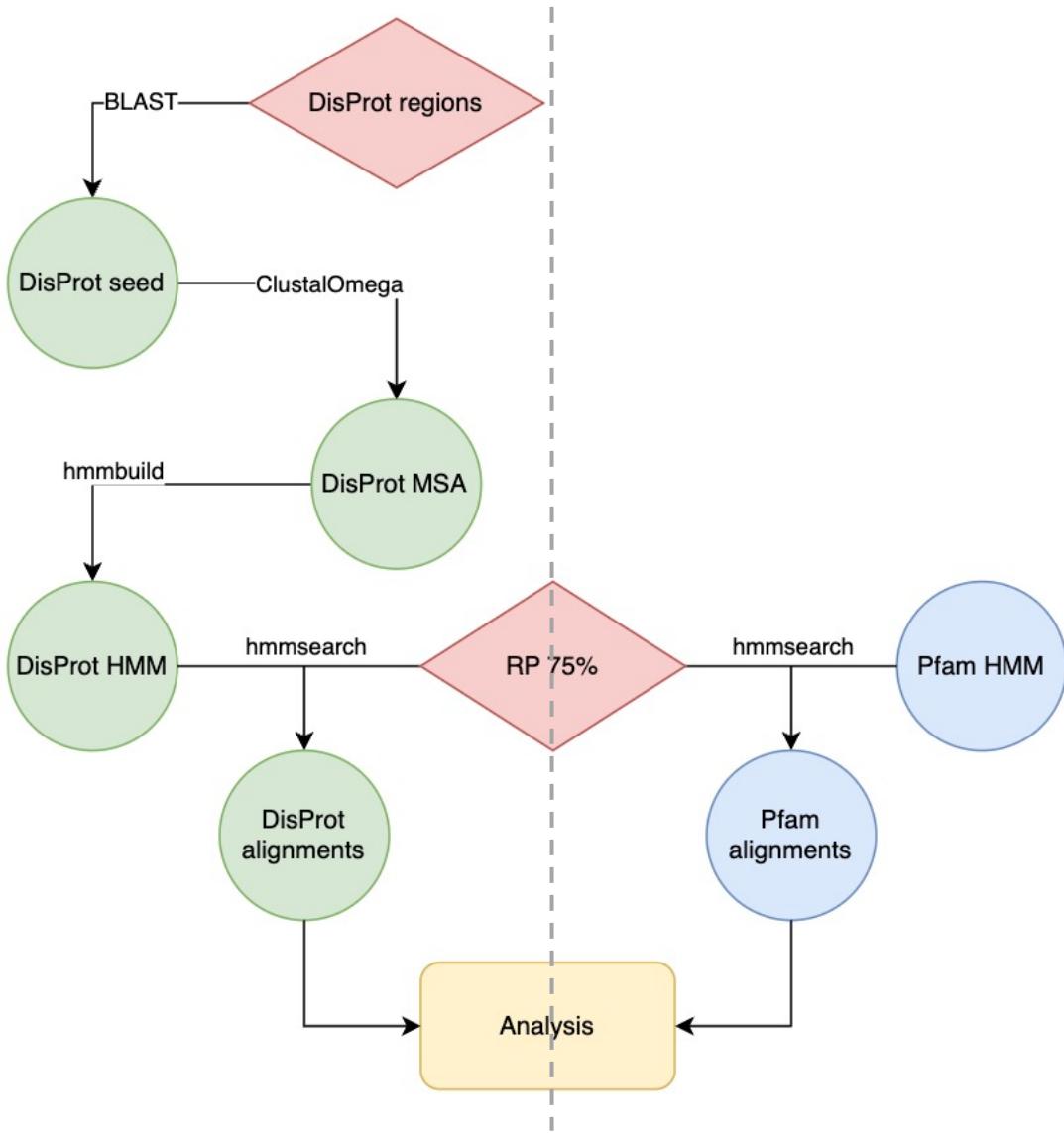
Multiple Sequence Alignment (MSA)

A technique used to arrange **3 or more** protein sequences, in a way that highlights both the similarities and differences among them.

- Helps in identifying **the least conserved regions**



Thesis workflow



Databases		
DisProt	Manually curated repository of IDPs and IDR (2272 proteins, 3151 regions)	Extract information about IDR positions
UniProt	Database of protein sequence and functional information	Obtain the seed alignments
Reference Proteome (75%)	Subset of proteomes that have been selected to provide a representative cross-section of the taxonomic diversity	Search for the proteins matching HMMs
Tools		
ClustalOmega	Program to build MSAs to generate alignments between three or more sequences	Build the MSAs
HHMMER	Software package for sequence analysis with the various tools to build and analyse HMMs	Build and search HMMs

HMM build and search

- 1 step: building the HMMs

Input: MSAs of IDRs

Output: HMMs

! Very fast – about 1-2 minutes for the whole dataset processing

- 2 step: searching the proteins based on HMMs in a proteome database (RP 75%)

Input: HMMs

Output: table of matched proteins

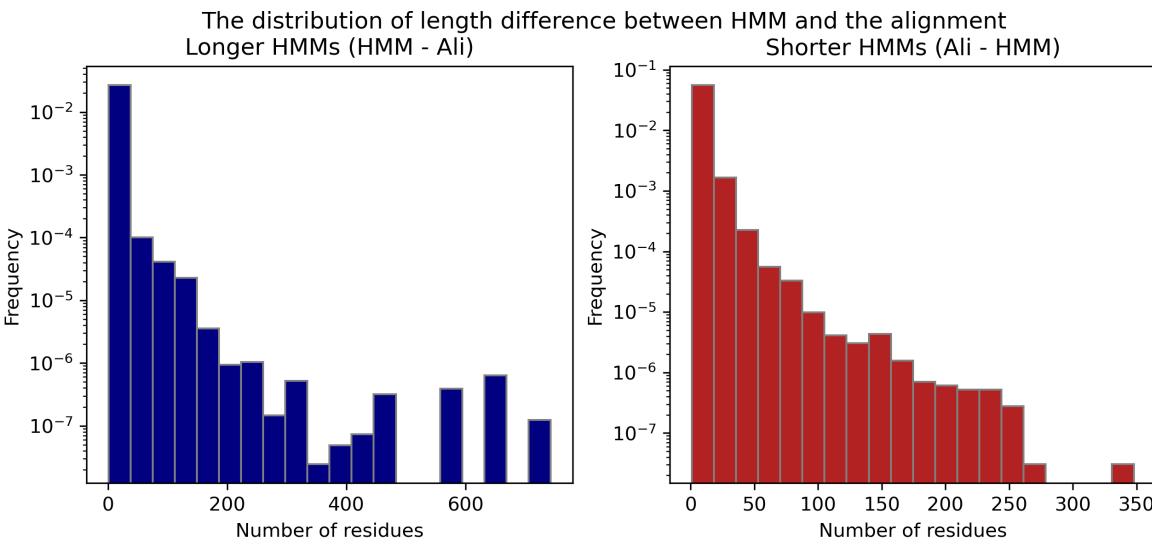
hmmbuild

HMM	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
COMPO	6.93836	6.68823	2.78359	1.56765	5.66943	2.59069	6.01310	6.05767	5.13250	5.33880	2.71543	6.23546	5.66889	1.98465	1.98371	2.65961	2.66563	5.72936	6.93832	2.70045
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
1	6.93836	6.68823	2.78359	1.56765	5.66943	2.59069	6.01310	6.05767	5.13250	5.33880	2.71543	6.23546	5.66889	1.98465	1.98371	2.65961	2.66563	5.72936	6.93832	2.70045
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
2	3.90441	5.75457	5.46612	5.58732	6.22794	4.48988	6.08677	5.69397	5.67298	5.98140	5.28159	6.33975	5.87445	5.09021	5.28159	5.74284	6.09654	4.48874	5.26850	3.88682
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
3	5.87485	6.77789	6.31463	6.35795	7.37367	5.79855	5.28869	5.98481	6.27452	6.38558	5.98140	6.29431	6.18429	6.11330	5.83993	6.29431	6.18429	5.77283	4.07738	3.90993
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
4	5.46200	7.11353	4.52584	0.85997	6.87713	5.21683	6.21336	6.86719	5.44425	6.19336	7.34117	5.16433	5.98552	5.56833	5.83970	5.41278	5.83666	6.44135	5.75335	6.77741
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
5	5.46108	6.90219	5.34496	5.38534	6.45842	5.37626	6.33943	6.66943	5.13181	5.92175	7.10076	5.70381	6.03086	5.05688	5.28247	5.56877	5.85329	6.29851	7.33411	6.44392
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
6	5.46108	6.90219	5.34496	5.38534	6.45842	5.37626	6.33943	6.66943	5.13181	5.92175	7.10076	5.70381	6.03086	5.05688	5.28247	5.56877	5.85329	6.29851	7.33411	6.44392
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
7	5.46200	7.11353	4.52584	0.85997	6.87713	5.21683	6.21336	6.86719	5.44425	6.19336	7.34117	5.16433	5.98552	5.56833	5.83970	5.41278	5.83666	6.44135	5.75335	6.77741
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
8	4.23715	6.88464	5.80839	5.81984	6.34928	5.75247	6.35653	5.90568	5.84033	5.59644	6.16134	5.60216	5.56413	6.11511	5.85383	4.46125	8.07963	5.29017	7.42638	6.58977
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
9	4.47134	7.25522	4.22347	2.58282	6.54552	4.37082	5.13198	6.24331	4.45863	5.63542	6.64556	3.83844	5.11848	4.34399	5.27949	4.24969	4.82266	5.72548	7.68694	6.85904
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
10	5.38193	6.76793	6.21331	6.29851	7.06383	0.03137	6.07645	7.19554	6.57116	6.51694	6.74063	6.39622	6.12961	6.87818	6.47528	5.62143	5.92216	5.62526	7.61280	7.26792
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
11	5.46108	6.90219	5.34496	5.38534	6.45842	5.37626	6.33943	6.66943	5.13181	5.92175	7.10076	5.70381	6.03086	5.05688	5.28247	5.56877	5.85329	6.29851	7.33411	6.44392
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
12	5.65795	6.94695	6.22418	5.71763	6.76765	5.53028	6.14445	6.68109	5.09617	5.70381	5.45217	0.05374	5.78842	5.91615	6.34709	7.31411	6.63805	7.08796	5.29017	7.42638
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
13	5.46108	6.90219	5.34496	5.38534	6.45842	5.37626	6.33943	6.66943	5.13181	5.92175	7.10076	5.70381	6.03086	5.05688	5.28247	5.56877	5.85329	6.29851	7.33411	6.44392
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
14	5.07225	6.78828	5.94856	5.04823	6.45413	5.15537	6.37870	6.73284	5.62427	6.12144	7.21838	0.06174	5.88206	5.85265	5.18592	5.55846	6.15910	7.38948	6.44416	
	2.68618	4.42225	2.77519	2.73123	3.46354	2.40513	3.72494	3.29354	2.67741	2.69355	4.24698	2.98347	2.73739	3.18146	2.89801	2.37887	2.77519	2.98518	4.58477	3.61503
	0.08416	5.48317	*	0.61958	0.77253	*	0.00000	*	*	*	*	*	*	*	*	*	*	*	*	*

hmmssearch

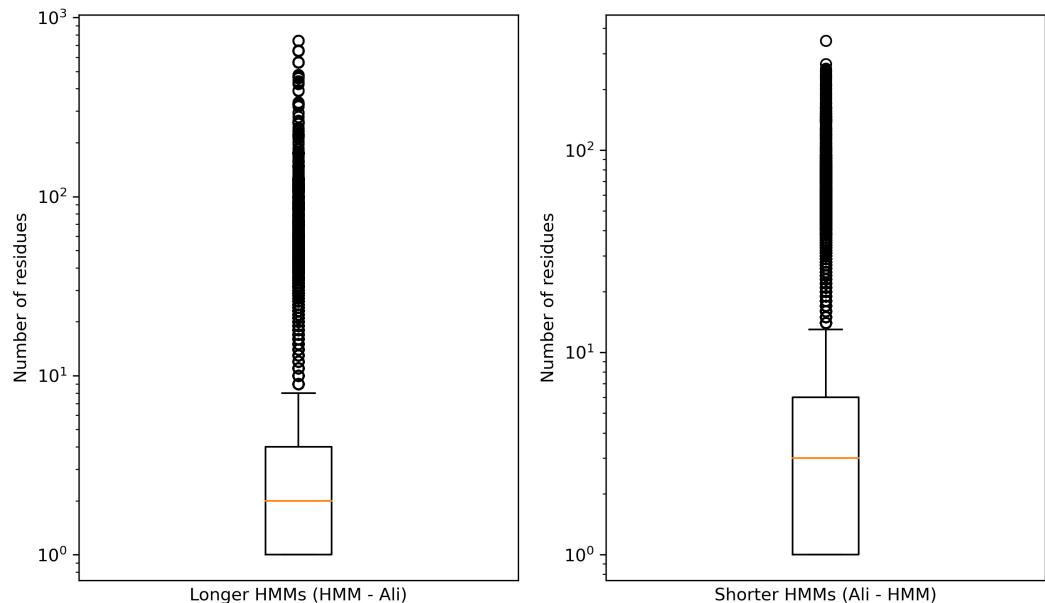
Query: A8AZZ3_146-195 [M=46]													
Scores for complete sequences (score includes all domains):													
--- full sequence ---			--- best 1 domain ---			#dom-							
E-value	score	bias	E-value	score	bias	exp	N	Sequence	Description				
5.4e-07	42.4	10.3	23	17.9	0.0	5.3	5	R0LE47	R0LE47_STRMT^ ^^ ^Ornithine carbamoyltransferase				
2.1e-06	40.5	20.6	2.1e-06	40.5	20.6	3.2	3	A8AZZ3	A8AZZ3_STRGC^ ^^ ^Amylase-binding protein AbpA {ECO:0000313 EMBL:ABV10035.1}				
5.7e-05	35.9	1.9	2.4e+04	8.3	0.2	8.3	8	A0A371S028	A0A371S028_9BACI^ ^^ ^Uncharacterized protein {ECO:0000313 EMBL:ABV10035.1}				
8.6e-05	35.3	0.1	1.5e+04	8.9	0.0	5.3	4	A0A158TMK1	A0A158TMK1_9CLOT^ ^^ ^Autolysin {ECO:0000313 EMBL:ABV10035.1}				
0.0011	31.7	14.9	0.0011	31.7	14.9	1.8	1	T0TW47	T0TW47_9STRE^ ^^ ^Uncharacterized protein {ECO:0000313 EMBL:ABV10035.1}				
0.0022	30.8	2.1	6.6e+04	6.9	0.1	6.5	7	U5MZS1	U5MZS1_CLOSA^ ^^ ^Surface protein PspC {ECO:0000313 EMBL:ABV10035.1}				
0.0064	29.3	6.6	0.0064	29.3	6.6	3.4	3	E7S8K8	E7S8K8_9STRE^ ^^ ^Uncharacterized protein {ECO:0000313 EMBL:ABV10035.1}				
0.0068	29.2	25.2	0.0068	29.2	25.2	4.6	4	A0A0F2CVN8	A0A0F2CVN8_STRCR^ ^^ ^Uncharacterized protein {ECO:0000313 EMBL:ABV10035.1}				
0.007	29.2	29.4	0.12	25.3	0.0	5.6	4	E1LN85	E1LN85_STRMT^ ^^ ^Cell wall binding repeat family				
Alignments for each domain:													
== domain 1 score: 0.3 bits; conditional E-value: 5.7													
A8AZZ3_146-195 22 eakkeaaaknagaK 35													
ea++at++n++at++													

Length comparison: HMM vs alignment



HMM length:
the length of the **query** sequence
Alignment length:
the length of the **target** sequence

The boxplots of length difference between HMM and the alignment



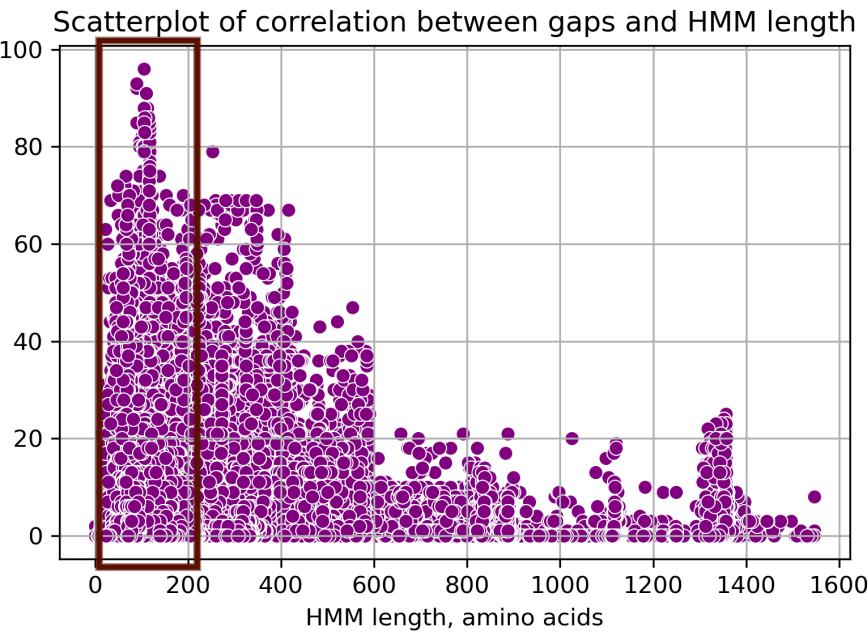
The number of pairwise alignments: **5,878,578**, from which HMMs:

- Match the target sequences length: 2,906,648 (**49.44%**)
- Longer HMMs: 1,097,774 (**36.94%**)
- Shorter HMMs: 1,874,156 (**63.06%**)

- **13.0%** – longer HMMs outliers (8 residues)
- **7.6%** – shorter HMMs outliers (13 residues)

The low proportion of outliers suggests that we can use **ali from/to** for further comparison with Pfam domains

Gaps analysis

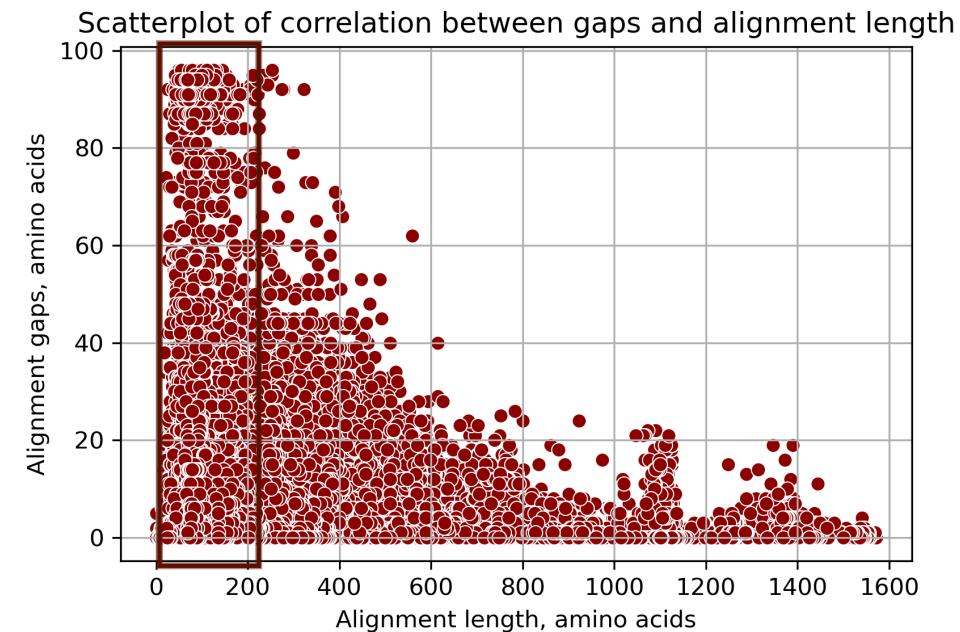


A8AZZ3_146-195	1	hgaqtgkaakasa...atkpeakeeaa....aknakagaKagqKALP	40		
	+g++++++	k+++	++ ++a+++++a	++n++++++	+++ ++
R0LE47	215	QGQSAKGWQKDAKgqwSYLKDAQGTKAtgw	lKDNGTWYYLNAEGVMQ	261	

Example of the **insertions** for A8AZZ3 protein: HMM

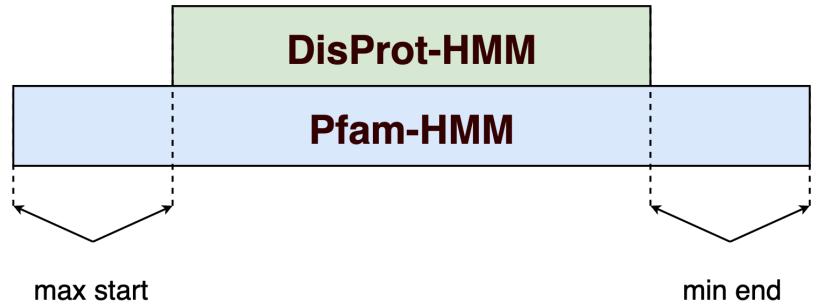
A8AZZ3_146-195	5	tgkaakasaatkpeakeeaaaknakagaKagqKALPKTsAVK	46
	t++a++a+	+ k e+++++a a n+ka +++++	KALPKTsAVK
T0TW47	169	TNEATNAA-K-KTEEGAKA-AQNGKASSAQAGKALPKTSAVK	207

Example of the **deletions** for A8AZZ3 protein: target sequence



The distribution of gaps across the sequences is similar for both HMM and target sequences. **Shorter sequences** tend to have **more gaps** relative to their length.

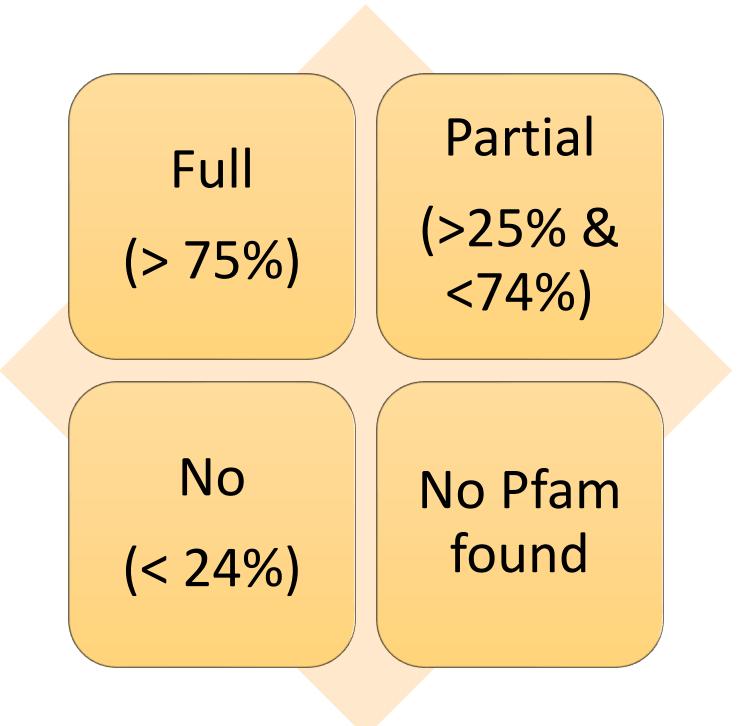
Overlap types



Overlap length: end – start

Total length: pfam length + disprot length – overlap

Overlap percentage: overlap/total length (%)



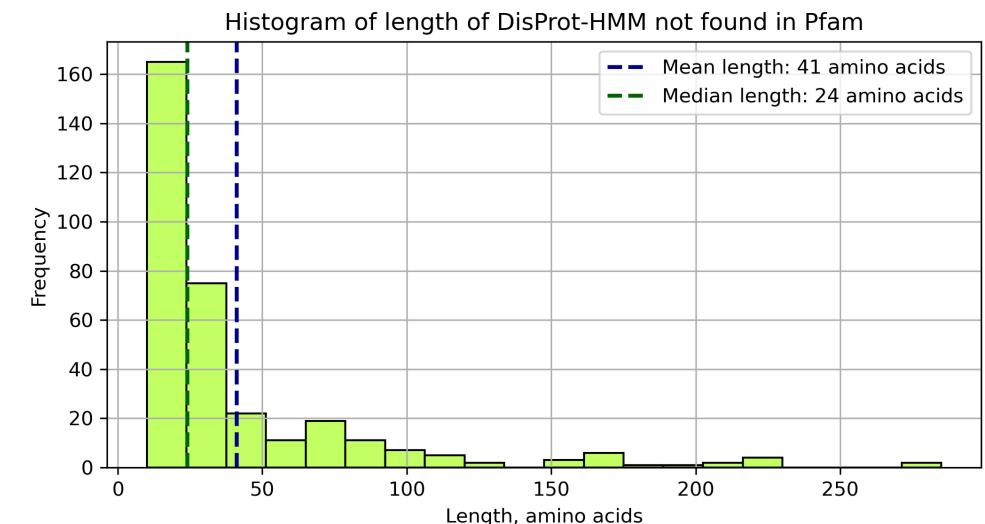
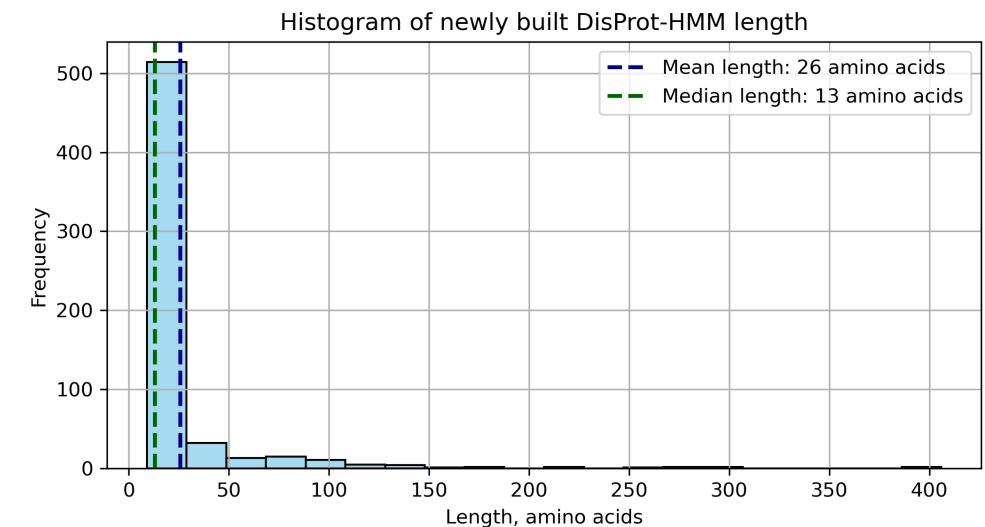
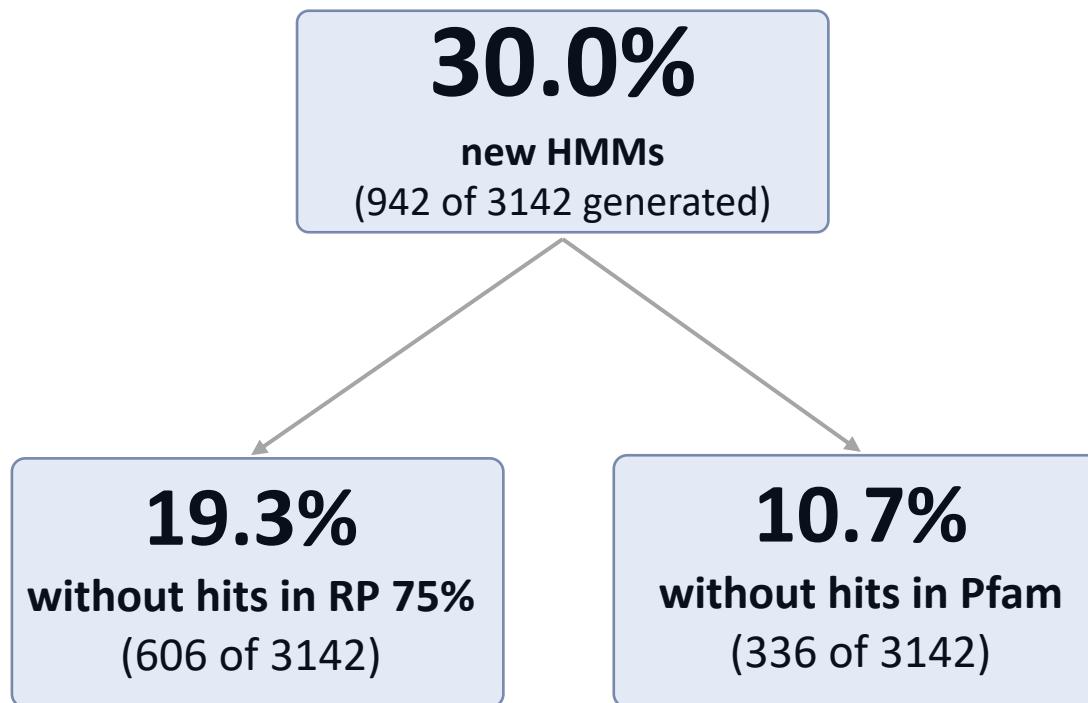
DisProt-HMM distribution:

- Full overlap: 371;
- Partial overlap: 891;
- No overlap: **2196**;
- Full and partial overlaps: 24;
- Full and no overlaps: 341;
- Partial and no overlaps: 766;
- All types of overlaps: 14.

Orphan HMMs

The new HMMs:

- 1) generated by hmmbuild but not found in RP 75%;
- 2) found in RP 75% but not in Pfam.



Summary statistics



Dataset	# of proteins	# of regions
Curated DisProt	2272	3151
UniProt BLAST	2271	3150
ClustalOmega	2268	3147
hmmbuild	2263	3142
hmmsearch	1971	2536
Pfam *	1734 (76.32%)	2200 (69.82%)

-1 instance: not found as a query sequence by BLAST
-3 instances: no subject sequences found
-5 instances: failed at cluster splitting in ClustalOmega
-606 instances: not found in RP 75%
-236 instances: not found in Pfam

* - the numbers are dependent on the specified thresholds and may vary accordingly

Results

- **70%** of DisProt-HMMs overlapped Pfam-HMMs, while **30%** of DisProt-HMMs are new;
- Among overlapping instances the larger amount of DisProt-HMMs corresponds to “no overlap” type, comprising **2196** models;
- Most gaps are found in **shorter** sequences, indicating that the HMM considers the IDRs to be shorter than those stated in the Curated DisProt;
- HMM without hits in RP 75% and Pfam are usually short in length.

Possible developments

- Different alignment format: .sto instead of .fasta;
- Other types of MSA: MAFFT, T-Coffee, etc;
- Regions selection: envelope regions instead of alignment;
- Take into account conditional and independent E-values.

Thank you for your attention!