Lessons 5 & 6. Content

#### 1. Alignment algorithms

2. Multiple alignments

# Alignment

- 1) Score for the correspondence of aa/bases
- 2) Penalty of insertions/deletions
- 3) Algorithm performing the alignment
- 4) Measure of the alignment significance

# Alignment

- 1) Score for the correspondence of aa/bases
- 2) Penalty of insertions/deletions
- 3) Algorithm performing the alignment
- 4) Measure of the alignment significance

LAMIAPRIMASEQCREATA-------MIAALTR----ASEQDAALLINEARE LAMIAPRIMASEQ-----CREATA --MIAAL---TRASEQDAALLINEARE



# Some definitions...

From the persian mathematician al-Kharezmi (al-Khawarizmi) from the IX century

Algorithm: Sequence of well-defined instructions (univocally interpretable) that allows to reach an outcome in a finite number of steps



# Some definitions...

From the persian mathematician al-Kharezmi (al-Khawarizmi) from the IX century

Algorithm: Sequence of well-defined instructions (univocally interpretable) that allows to reach an outcome in a finite number of steps

Step 1: Gather Your Ingredients for the Sandwich.

Step 2: Put Gloves on (optional).

Step 3: Pull Out Two Slices of Bread and put on plate.

Step 4: Open Peanut Butter and Jelly.

Step 5 Pick up butter knife.

- Step 6: Spread the Peanut Butter Onto One Slice of Bread.
- Step 7: Spread the Jelly Onto the Other Slice of Bread.
- Step 8: Combine the Two Slices.
- Step 9: Say: Enjoy



# An example...



The **binary search** (or half-interval search) is an algorithm that finds the position of a target value within a sorted array.

It compares the target value to the middle element of the array; if they are not equal, the search continues on the remaining half, again taking the middle element *etc*.

# Some definitions...

From the persian mathematician al-Kharezmi (al-Khawarizmi) from the IX century

**Algorithm:** Sequence of well-defined instructions (univocally interpretable) that allows to reach an outcome in a un risultato in a finite number of steps

**Programme:** description of an algorithm in a specific coding language (e.g. C, fortran, python, perl).

# Exact algorithms vs. heuristic algorithms

Type of algorithm	programmes	Pro	Cons
Exact	<ol> <li>Needlman Wunch (1970)</li> <li>Smith Waterman (1981)</li> </ol>	Sensitive	Slow
Heuristic	<ol> <li>BLAST –Altschul (1990)</li> <li>FASTA- Pearson (1985)</li> </ol>	Fast	Non sensitive

# Some definitions...

Heuristic method: in computer science, it is heuristic a method that can generally find a good solution to a problem, though it is not possible to prove that the found solution is the correct one

#### Exact algorithms: e.g. algorithm of Smith-Waterman

Exact, it guarantees to find out the **best** alignament(s) for a pair of sequences.

For 2 sequences: A of length n and B of length m, Smith-Waterman takes n\*m computational steps

In searching a database of sequences -

- In case the query sequence A is long n=200 nucleotides
- And we search for homologues sequences in the EST database containing, e.g., 23\*10<sup>6</sup> sequences, B<sub>i</sub>, each of length m=500.
- Number of computational steps:
   23\*10<sup>6</sup> \* 500 \* 200 ~ 10<sup>11</sup> total steps !

#### Exact algorithms: e.g. algorithm of Smith-Waterman

Exact, it garantees to find out the **best** alignament(s) for a pair of sequences.

For 2 sequences: A of length n and B of length m, Smith-Waterman takes  $n^*m$  computational steps \*10<sup>6</sup> \* 500 \* 200 ~ 10<sup>11</sup> passi totali !

How do we discard the irrelevant alignments?

The heuristic algorithms (BLAST, FASTA) can filter most of the irrelevant alignments.

2 sequences of length *n* and *m* with a 'sliding' algorithm would require  $n \ge m$  comparisons between positions: problem O(*nm*) ~ O( $n^2$ )\* (<u>quadratic size</u>)

Examp	ole <b>6</b> x <b>5</b> :		
1)	LLKKQW	2) I	LKKQW>
	LLKQW		LLKQW
3)	LLKKQW	4)	LLKKQW
	LLKQW		LLKQW
5)	L <mark>L</mark> KK <mark>QW</mark>	6)	<b>LLK</b> KQW
	LLK <mark>QW</mark>		LLKQW
7)	LKKQW	8)	LLKKQW
	L <mark>L</mark> KQW		LLKQW
9)	LLKKQW	10)	LLKKQW
	LLKQW		LLKQW
* If <i>n</i>	and <i>m</i> have the same of	order of magn	itude

#### What if we allow gaps?

In a sequence long *n* one can insert gaps in *n*-1 positions.

Just allowing 1-res gaps "-" "n" different sequences are obtained

Example n = 6:

- 1) LLKKQW 2) L-LKKQW
- 3) LL-KKQW 4) LLK-KQW

5) LLKK-QW 6) LLKKQ-W

#### What if we allow gaps?

In a sequence long *n* one can insert gaps in *n*-1 positions.

Just allowing 1-res gaps "-" "n" different sequences are obtained

Example  $\mathbf{n} = 6$ :

- 1) LLKKQW 2) L-LKKQW
- 3) LL-KKQW 4) LLK-KQW

5) LLKK-QW 6) LLKKQ-W

In case we allow a larger number of gaps (besides the 1-res) the number of possible sequences increases <u>exponentially</u> and so does the problem size

Exact alignment algorithms such as the Needlman-Wunch and Smith-Waterman are examples of Dynamic Programming

Breaking down the problem into simpler subproblems, then recursively finding the optimal solutions to the sub-problems

#### Pots of gold game: rules



Going from START to END without passing twice through the same point and without moving backwards, while collecting the max number of pots of gold

# Pots of gold game



# Pots of gold game



Choice of the best path from A to END does not depend on the path taking me from START to A

## Pots of gold game: the optimal path (solution)



## Pots of gold game: the optimal path (solution)





introduction of GAPs (both *opening*, and *extension* ones) in the alignment

e.g. How we align the following sequences:

#### **HEAGAWGHEE**

#### **PAWHEAE**

?

# An example of alignment algorithms: cumulative matrix

#### **HEAGAWGHEE** *vs* **PAWHEAE**

Step1: building the site-specific matrix (values from BLOSUM45)

	H	E	Α	G	Α	W	G	H	E	E
Ρ	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1
Α	-2	-1	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
н	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
Α	-2	-1	5	0	5	-3	0	-2	-1	-1
Ε	0	6	-1	-3	-1	-3	-3	0	6	6

#### **Alignment algorithm: cumulative matrix**



Step2: building the *cumulative* matrix, where each element represents the optimal score that can be obtained from start to that point







#### **Cumulative matrix**

		H	E	Α	G	Α	W	G	н	E	E
	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
Ρ	-8	-2	-9	-17	-25	-33	-41	-49	-57	-65	-73
Α	-16	-10	-3	-4	-12	-20	-28	-36	-44	-52	-60
W	-24	-18	-11	-6	-7	-15	-5	-13	-21	-29	-37
н	-32	-14	-18	-13	-8	-9	-13	-7	-3	-11	-19
Е	-40	-22	-8	-15	-15	-9	-12	-15	-7	3	-5
Α	-48	-30	-15	-3	-11	-10	-12	-12	-15	-5	2
Е	-56	-38	-23	-11	-6	-12	-13	-15	-12	-9	1

	<u> </u>	E	Α	G	Α	W	G	H	E	E	
Ρ	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1	
Α	-2	-1	5	0	5	-3	0	-2	-1	-1	<b>•</b> •
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3	Orig
н	10	0	-2	-2	-2	-3	-2	10	0	0	
E	0	6	-1	-3	-1	-3	-3	0	6	6	
Α	-2	-1	5	0	5	-3	0	-2	-1	-1	
E	0	6	-1	-3	-1	-3	-3	0	6	6	

Original site-specific matrix



Step3: backwards path through cells that allowed to obtain the best alignment scores

		Н	E	Α	G	Α	W	G	Н	E	<u> </u>
	0-	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
Ρ	-8	-2	-9	-17	-26	-33	-44	-50	-58	-65	-73
Α	-16	-10	-3	-4	-12	<u>-1</u> 5	-26	-34	-44	-53	-62
W	-24	-19	-13	-6	-7	-15	0	-11	-22	-33	-44
Н	-32	-14	-19	-15	-8	-9	-11	-2-	<b>— 0</b> -	-8	-16
Е	-40	-22	-8	-17	-18	-9	-12	-13	-2	6	4
Α	-48	-32	-17	-3	-11	-12	-12	-12	-12	-3	5
Е	-56	-40	-19	-12	-6	-12	-15	-15	-12	-5	3
	HE.	AG	AW	GH	E – ]	2		del	<b>etio</b>	n	
				•••		-		ins	ortic	h	·
	-P.	<b>A</b> -	– W	-H.	ĽA	<b>H.</b>		113		,,,	

# An example of alignment algorithms: cumulative matrices

#### **HEAGAWGHEE** *vs* **PAWHEAE**

Step1: building a site-specific matrix (values from a PAM or BLOSUM matrix)

- Step2: building a *cumulative* matrix, where each element represents the maximum score achievable to go from start to that point
- Step3: backwards path through the cells which allowed to obtain the best scores = optimal alignment

#### **Multiple alignment**

Luotvyygupunne atttlfcase a aydteunnu	MATHACUPTEP NP OF UNL GNOTEN FNMUN NMMUTOMOED I ISLUD OSL
LINTUYYGUPUN <mark>KE</mark> ATTTLI <mark>CASB</mark> A <mark>M</mark> AY <mark>DTEUHN</mark> U	wathacupternpercenturutenfrimmendruconned i isludosl
LUUTUYYGUPUM <mark>HE</mark> AATTLFCASDA <mark>R</mark> AYLTEUHNU	WATHACUPTIP NP OF UULENUTEN FINNUR NNNUE OMHED I ISL WDOSL
LWUTUYYGUPUW <mark>KE</mark> ATTTL FCABB A <mark>M</mark> AY <mark>STE AHE</mark> U	wath acopt up up gevolenot en fumur numor onner i i slud o sl
LWOTVYYGVPVWHE ANTTLYCASE ARAYITE IHRV	WATH ACUPTOP NP OF LUMGNUTEN FINNER NUMUE ON HED IISL ND 08L
LUOTUYYGUPUWKE ATTTLTCASE AMAYNTEUHNU	WARH ACOPTOP NP OF ULLEN OT EN FWHOR NNMUF OMNED I I SLOD Q SL
LUDTUYYGUPUMHE ATTTLFCASD A AYATEUHNU	WATHACUPTIP NP OF UUL GNUTEN FIMUN NNMUE OMOED I ISL ND OSL
LUOTUYYSUPUUNE ATTTLFC ASD AMAYDTE ANNU	WATHACOPTNPNP OF WILENUTEN FINIOR NUMVE ONHED I ISLUDOSL
LAGTOWY GOP ON TALTFLECASD AN AVET FRHOM	WATHACOPTIPNPOLLSLONOTENFIMOR NUMORONHEDOISLUDOSL
L GOT OVYGOPOGREAT TTLTCASH ARS VETEARN T	WATHACOPTIP NP GEIALENUTEN FINNOR MIMUE OMHED IISLOBOSI.
LUOT OYYGOPODERATTELY CASD ARSYL PEARNT	GATHACOPTIPRPETIMENOTIAFMORMMOLOMETITELO
E OOT OTTGOP TOWEATTTETCASE ANATORE AND T	WATH SCOPTIP AP OF THE ENOTES FAMILY MANUEL PRODUCTS
LWOIDTTGOPONIESTILLTCASDARS TILLOMMI	WATCACOPTOP SPORT PROVIDE STATES AND A STATE
INCIDENTIAL CONTRACTOR CAPTURE	MATTACOPIEP OF CHOILE FOR AND AND
CARTER AND A STOLE AT A STOLE	THE TREE DISTURDED FOR NOT THE ARE AND
OWNERS AND AND THAT THE ATTING	MET TOCLED NOT YOF THE AFT AND WINTED A TED TO THE FET ST
CARTINE SOFTING NASIPERCATER	MGT INCEPTION OF ITLN-MTEAFLAND NTWEE AVEL MARK FETST
OFUTUFYGIPAN NASIPLFCATERDT	MGT IOCLPDNDD YOF ITLN-OTE APT AWNNTOTED AUEDUMAL PETS I
OVUTVEYOUP AN MAT IPLECATED	WGTTOCLPDNDDYSELALN-UTESELAWENTUTEOAIEDUWOLFETSI
TUTUFYGIP AN NATUPLICATING	WGTVOCLPDNGDYTE I LM- ITE AFT AWDNTUT 00 AUDDUMEL FETS I
GUITVFYSVPVWWWSEVQAFCMTPTT	WATTNC IP DERD TREVPLN - ITEP TE AWAD NPLVAQAG SN INLL FE GTL
LYUTUFYGIPUMURSTUCAFCMTPNTNM	MATTNC IPDINDATEUPLN - ITE AFT AND NPLUNO AE SE INLL FE OTM
GYUTUFYGUPUNKE AKTHLICATONSSL	WUTTWC IP SLPEY THUT IP I MENT GL I ENGIVYO ANH AMG SML IT IL

Few definitions...

Pair-wise sequence alignment = Alignment of TWO similar sequences

Multiple sequence allignment (MSA) = Alignment of MANY similar sequences, generally coming from a search in databases

## **Exact algorithms for multiple alignments**

The necessary number of computational steps is in the order  $L^N$ , where L is the length and N the number of sequences to be aligned

e.g. for 4 sequences 100 aa/nt long

The number of computational steps would be 100<sup>4</sup> = 100 millions !

not viable

# **Approximate solution**

We perform a pairwise alignment(s) then align additional sequences to it(them)


## **Approximate solution**



A profile of the alignment can be built in the form of a PSSM (position-specific scoring matrix) but also by HMM

	А	S	D	Κ	L
	V	S	Е	R	F
А	$\frac{1}{2}(PAM(A,A) + PAM(A,V)) = \frac{1}{2}(2+0)$	<sup>1</sup> / <sub>2</sub> (1+1)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (-1-2)	<sup>1</sup> / <sub>2</sub> (-2-4)
G	1/2(1-1)	<sup>1</sup> / <sub>2</sub> (1+1)	1/2(1+0)	<sup>1</sup> / <sub>2</sub> (-2-3)	1/2(-4-5)
R	<sup>1</sup> / <sub>2</sub> (-2-2)	<sup>1</sup> / <sub>2</sub> (0+0)	1/2(-1-1)	<sup>1</sup> / <sub>2</sub> (3+6)	<sup>1</sup> / <sub>2</sub> (-3-4)
S	<sup>1</sup> / <sub>2</sub> (1+0)	1/2(3+3)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (-3-3)
G	1/2(1-1)	<sup>1</sup> / <sub>2</sub> (1+1)	1/2(1+0)	<sup>1</sup> / <sub>2</sub> (-2-3)	1/2(-4-5)
S	<sup>1</sup> / <sub>2</sub> (1+0)	1/2(3+3)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (-3-3)

Example of PSSM

## **Approximate solution**



The easiest way to go is building a PSSM by performing an arithmetic <u>average of the</u> <u>scores</u> for the alignment of ress of the 3<sup>rd</sup> to those of the aligned sequences

	А	S	D	Κ	L
	V	S	Е	R	F
Α	$\frac{1}{2}(PAM(A,A) + PAM(A,V)) = \frac{1}{2}(2+0)$	<sup>1</sup> / <sub>2</sub> (1+1)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (-1-2)	<sup>1</sup> / <sub>2</sub> (-2-4)
G	1/2(1-1)	<sup>1</sup> / <sub>2</sub> (1+1)	1/2(1+0)	1/2(-2-3)	1/2(-4-5)
R	1/2(-2-2)	1/2(0+0)	1/2(-1-1)	<sup>1</sup> / <sub>2</sub> (3+6)	<sup>1</sup> / <sub>2</sub> (-3-4)
S	<sup>1</sup> / <sub>2</sub> (1+0)	1/2(3+3)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (-3-3)
G	1/2(1-1)	<sup>1</sup> / <sub>2</sub> (1+1)	1/2(1+0)	1/2(-2-3)	1/2(-4-5)
S	<sup>1</sup> / <sub>2</sub> (1+0)	1/2(3+3)	<sup>1</sup> / <sub>2</sub> (0+0)	<sup>1</sup> / <sub>2</sub> (0+0)	1/2(-3-3)

Example of PSSM



	Α	S	D	K	L		A V	S S	D E	K R	L F
А	2	1	0	-1	-2	А	1	1	0	-1.5	-3
G	1	1	1	-2	-4	G	0	1	0.5	-2.5	-4.5
R	-2	0	-1	3	-3 -	R	-2	0	-1	4.5	-3.5
S	1	3	0	0	-3	S	0.5	3	0	0	-3
G	1	1	1	-2	-4	G	0	1	0.5	-2.5	-4.5
S	1	3	0	0	-3	S	0.5	3	0	0	-3

ASDKL VSERF A-GRSGS

# A classical approach: CLUSTALW

- Programme for MSAs based on a hierarchic approach
- Step1: pairwise alignment for all the input sequences (for N seq.: N(N-1)/2 pairwise alignments)
- Step2: building a guide tree, i.e. a hierarchy of sequences in order of their similarity (cluster analysis)
- Step3: building the multiple alignment based on the guide tree by first aligning the most similar pairs, then aligning the other sequences with those pairs until all have been aligned

#### Measured distance = % of different amino acids

	Seq1	Seq2	Seq3	Seq4
Seq1	0	5	11	14
Seq2		0	9	10
Seq3			0	7
Seq4				0

The measure of dissimilarity between sequences represents their evolutionary distance and can be obtained in several ways

The % of different amino acids is one possible way



#### Measured distance = % of different amino acids



	Cls1-2	Cls 3-4
Cls1-2	0	1/2[d(Cls 1-2),3] + 1/2[d(Cls 1-2),4] = 11
Cls 3-4		0





### Example of guide tree

The guide tree is used to guide the order of constructing the multiple alignment

![](_page_44_Figure_2.jpeg)

#### Another example of guide tree

How a guide tree, built on the basis of pair-wise alignments between all the sequences, guides the order of constructing the multiple alignment

![](_page_45_Figure_2.jpeg)

In principle, the pair-wise alignments can also be approximate, for instance based on the presence of k-mers (stretches of k residues) in common for two sequences

# **MSAs: current approaches**

- ClustalW has been retired and substituted as a web service by Clustal Omega
- Clustal Omega is a new MSA tool that uses seeded guide trees and HMM profile-profile techniques to generate alignments

A HMM is a **machine learning method**: it learns from known cases, assigns probabilities to events based on observations and makes predictions

HMMs can be designed for many tasks, in Bioinformatics and beyond

Markov Model means that there is a statistical **Markov chain**, i.e. each state at the step n depends on the state at the step n-x (if x=1, the Markov chain is of first order)

Hidden means that its **states are not observable** 

#### **HMM structure**

A HMM is characterized by:

Number of states (X), number of symbols (Y), distribution of transition probabilities between states (A), distribution of emission probabilities of symbols (B), initial state distribution ( $\pi$ ).

States can be "silent", in case they do not emit symbols

![](_page_48_Figure_4.jpeg)

Hidden Markov models are made of unobservable (hidden) states

Each state emits symbols from a fixed alphabet, e.g. ACGT, according to specific emission (or output) probabilities

Different (hidden) states are connected by precise transition probabilities

The sequence of states is a Markov chain: the choice of next element (state) depends on the actual one (1<sup>st</sup> order chain)

While states are "hidden", symbols (e.g. ACGT) are "observable"

Observations are probability functions of the "hidden" states

Hidden Markov models are made of unobservable (hidden) states

Each state emits symbols from a fixed alphabet, e.g. ACGT, according to specific emission (or output) probabilities

Different (hidden) states are connected by precise transition probabilities

The sequence of states is a Markov chain: the choice of next element (state) depends on the actual one

While states are "hidden", symbols (e.g. ACGT) are "observable"

We can think of a HMM as a generator of sequences with defined probabilities

Hidden Markov models are made of unobservable (hidden) states

Each state emits symbols from a fixed alphabet, e.g. ACGT, according to specific emission (or output) probabilities

Different (hidden) states are connected by precise transition probabilities

The sequence of states is a Markov chain: the choice of next element (state) depends on the actual one

While states are "hidden", symbols (e.g. ACGT) are "observable"

Basically we use a sequence of observations to estimate the sequence of hidden states

In this example of HMM the hidden states are the used dices (regular or modified) and the outcomes of rolls are the symbols (observable)

![](_page_52_Figure_2.jpeg)

Rigged (modified) dice: Same external features of the regular dices pair

Roll	abcdefghilmnop
Num1	64251324653255
Num2	43516241313243

?? Given the above sequence of symbols, what pair of dices have we thrown ??

#### Dices choice

![](_page_53_Figure_1.jpeg)

#### **One more example of HMM**

![](_page_54_Picture_1.jpeg)

![](_page_54_Picture_2.jpeg)

![](_page_54_Picture_3.jpeg)

![](_page_54_Picture_4.jpeg)

#### **Python representation of HMM parameters**

```
states = ('Rainy', 'Sunny')
observations = ('walk', 'shop', 'clean')
start_probability = { 'Rainy': 0.6, 'Sunny': 0.4}
transition_probability = {
    'Rainy' : {'Rainy': 0.7, 'Sunny': 0.3},
    'Sunny' : {'Rainy': 0.4, 'Sunny': 0.6},
    }
emission_probability = {
    'Rainy' : {'walk': 0.1, 'shop': 0.4, 'clean': 0.5},
    'Sunny' : {'walk': 0.6, 'shop': 0.3, 'clean': 0.1},
    }
```

![](_page_55_Figure_2.jpeg)

## Another example of HMM for a sequence alignment

![](_page_56_Figure_1.jpeg)

The hidden states here are: deletion, insertion and amino acid matches. Deletions are instances of "silent states"

## **Basic problems for HMMs**

Given the structure of a HMM (X,Y,A,B, $\pi$ )

Problem 1: how we calculate the probability of a sequence of observations (e.g. "LASD")  $O = O_1 O_2 O_3 ... O_n$ ? (forward-backward algorithm)

**Problem 2**: given a sequence of observations, how we choose an optimal sequence of states, which 'explains' the sequence of observations? (Viterbi algorithm).

**Problem 3**: how we adjust the parameters (transition probabilities) of the model to maximize the probability of a sequence of observations? (Baum-Welch algorithm) – this is how a HMM MSA is built

## Another example of HMM for a sequence alignment

![](_page_58_Figure_1.jpeg)

given path through the model will emit a sequence with an associated probability

With the forward-backward algorithm, we can calculate the probability of having a specific sequence, e.g. **PETS** (problem 1 - it will be the sum of all the paths emitting the sequence)

With the Viterbi algorithm, we can choose the optimal sequence of states (most probable *path*), which 'explains' the sequence (problem 2); this is analogous to the best-scoring aln in dynamic programming

## HMM profiles (or profile HMMs)

![](_page_59_Figure_1.jpeg)

https://www.ebi.ac.uk/training/online/courses/pfam-creating-protein-families/what-are-profile-hidden-markov-models-hmms/

# **Sequence logos**

Profiles of MSAs can be represented graphically in the form of sequence logos, easily showing the residue preference or conservation at particular positions, which point to a functional role

We have already encountered PSSMs (Position Specific Scoring Matrices), examples of scoring schemes of MSAs for searching for other similar sequences, represented as sequence logos

![](_page_60_Figure_3.jpeg)

Examples from Web Logo

## HMM profiles (or profile HMMs)

HMMs are commonly used to align a novel sequence to a HMM profile or to align HMM profiles one to each other

A HMM profile is a HMM model containing the information present in a multiple alignment

HMM profiles can be visualized using sequence logos to illustrate the emission probabilities for different residues types at each match state

They are more sophisticated versions of a PSSM, especially because they can treat INDELs in a position-dependent way

HMM profiles of protein families and subfamilies are reported in several databases (BLOCKS, Pfam etc.)

## HMM profiles (or profile HMMs)

HMMs are commonly used to align a novel sequence to a HMM profile or to align HMM profiles one to each other

A HMM profile is a HMM model containing the information present in a multiple alignment

HMM profiles can be visualized using sequence logos to illustrate the emission probabilities for different residues types at each match state

They are more sophisticated versions of a PSSM, especially because they can treat INDELs in a position-dependent way

![](_page_62_Figure_5.jpeg)

Example of a HMM profile for the **Toxin\_7 family**, from the **Pfam database** 

## **Alignment of HMM profiles**

The alignment of two HMM profiles is actually the alignment of two alignments; in it the gap scoring is position-dependent

In a possible approach, one multipe alignment is firstly reduced to a profile HMM, then a modification of the Viterbi algorithm is used to find the most probable set of paths which emit the other alignment (to get the overall probability for the alignment the probabilities for each sequence path must be multiplied)

**HHsearch** aligns two profile HMMs and is designed to identify very remote homologs; it also uses a variant of the Viterbi algorithm to find the alignment with the best score

![](_page_63_Figure_4.jpeg)

Simplified visualization of the alignment of two HMMs (from Pfam) using logos to illustrate the emission probabilities at each match state

![](_page_64_Picture_0.jpeg)

**Multiple Sequence Alignment** 

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three** or more sequences. For the alignment of two sequences please instead use our pairwise sequence alignment tools.

Important note: This tool can align up to 4000 sequences or a maximum file size of 4 MB.

STEP 1 - Enter your input sequences	
Enter or paste a set of	
PROTEIN	<b>▼</b>
sequences in any supported format:	
Or, upload a file: Choose File no file selected	Use a example sequence I Clear sequence I See more example input
Or, upload a file: Choose File no file selected	Use a <u>example sequence</u> I <u>Clear sequence</u> I <u>See more example input</u>

## https://www.ebi.ac.uk/Tools/msa/clustalo/

	106010206 mod ND 001110574 1	UDBUUMDUCDUCUODTI NRUMEOLOUCU VORDMOVODLOCLEDUUDI	205
gı	186910296 [ref NP_001119574.1]	VPEKKTPKSPVGVQPILNEHTFCAGMSKYQEDTCYGDAGSAFAVHDL	305
gı	4826/62 rer NP_005134.1	VPERKTPKSPVGVQPILNEHTFCAGMSKYQEDTCYGDAGSAFAVHDL	364
gı	21264363 rei NP_006601.2	YGGKDSCRGDSGGALVFLDS	64Z
gı	21264357 ref NP_001870.3	GGKDACAGDSGGPMVTLNR	655
gı	21264359 ref NP_624302.1	GRYSVTENMFCAGYYEGGKDTCLGDSGGAFVIFDD	673
gi	4502495 ref NP_001725.1	TADAEAYVFTPNMICAGGEKGMDSCKGDSGGAFAVQDP	641
gi	66347875 ref NP_001724.3	NRMDVFSQNMFCAGHPSLKQDACQGDSGGVFAVRDP	663
gi	289547636 ref NP_057630.2	QRPEVFSDNMFCVGDETQRHSVCQGDSGSVYVVWDN	445
gi	4758502 ref NP_004123.1	GQDTCQGDSGGPLTCE	516
gi	295054188 ref NP_001171131.1	GQDTCQGDSGGPLTCE	490
gi	4504383 ref NP_001519.1	KSDACQGDSGGPLACE	605
gi	14702169 ref NP_127509.1	TVTDNMLCAGDTRSGGPQANLHDACQGDSGGPLVCL	474
gi	4505861 ref NP_000921.1	TVTDNMLCAGDTRSGGPQANLHDACQGDSGGPLVCL	520
gi	4505863 ref NP_002649.1	KTDSCQGDSGGPLVCS	383
gi	222537759 ref NP_001138503.1	KTDSCQGDSGGPLVCS	366
gi	119392081 ref NP 000195.2	GSIDACKGDSGGPLVCMDA	534
gi	4503635 ref NP_000497.1	RITDNMFCAGYKPDEGKRGDACEGDSGGPFVMKSP	577
gi	4506115 ref NP 000303.1	DRQDACEGDSGGPMVASF	410
gi	10518503 ref NP 062562.1	GSKDSCKGDSGGPHATHY	390
ģi	4503645 ref NP 000122.1	GSKDSCKGDSGGPHATHY	412
ģi	4503625 ref NP 000495.1	KQEDACQGDSGGPHVTRF	427
ģi	4503649 ref NP 000124.1	GGRDSCQGDSGGPHVTEV	419
qi	32698940 ref NP 872365.1	GAFDTCRGDSGGPLMCYLP	277
ģi	148231605 ref NP 937828.2	GGRDACOGDSGGSLMCRNK	249
ģi	300244530 ref NP 003610.2	RFTGRMLCAGNLHEHKRVDSCOGDSGGPLMCERP	834
ai	205360943 ref NP 005647.3	GNVDSCOGDSGGPLVTSK	449
ai	227499990 ref NP 001128571.1	GNVDSCOGDSGGPLVTSK	486
ai	13173471 ref NP 076927.1	GGVDSCOGDSGGPLVCOE	409
qi	33667063 ref NP 892018.1	GGVDSCSGDAGGPLACREP	1018
qi	4505881 ref NP 000292.1	GGTDSCOGDSGGPLVCFE	768
ai	116292750 ref NP 005568.2	RGTDSCOGDSGGPLVCFE	1998
ai	58331209 ref NP 001962.3	GVRSGCOGDSGGPLHCLV	214
ai	62526043 ref NP 009203.2	GVISACNGDSGGPLNCOLE	225
ai	15559207 ref NP 254275.1	GVISSCNGDSGGPLNCOAS	225
ai	58331211 ref NP 056933.2	GVICTNGDSGGPLNCOAS	225
ai	110815798 ref NP 899234.2	GITEKMICAGEAA-SGEKDECOGDSGGPLVCRHE	772
g_ ai	4503137 ref NP 001898.1	SITESMICAGGAGASSCOGESGGPLVCOKG	223
91 ai	118498341 ref NP 001897.4	GVSSCMGDSGGPLVCQKD	222
91	118498350 ref NP 001020371 3		222
	· · · · · · · · · · · · · · · · · · ·		

![](_page_66_Picture_0.jpeg)

Tools > Multiple Sequence Alignment > Clustal Omega

#### Results for job clustalo-I20220829-154434-0539-1190104-p2m

Alignments	Result Summary	Guide Tree	Phylogenetic Tree	Results Viewers	Submission Details	
Download G	Guide Tree Data					
Phylogra	am					
Branch length:	💿 Cladogram 🔵	Real				
			gil gil gil gil gil gil gil gil gil gil	186910296IrefINP 4826762IrefINP 21264363IrefINP 21264357IrefINP 21264359IrefINP 4502495IrefINP 66347875IrefINP 289547636IrefINP 289547636IrefINP 295054188IrefINP 4504383IrefINP 14702169IrefINP 4505861IrefINP 222537759IrefIN	P_001119574.11 0. 005134.11 0.00432 _006601.21 0.3046 _001870.31 0.1394 _624302.11 0.1394 001725.11 0.35933 _001724.31 0.2638 P_057630.21 0.2638 004123.11 0 P_001171131.11 0 001519.11 0.35036 _127509.11 0.00387 002649.11 0.01448 P_001138503.11 0	00432277 2277 365 485 485 33 36 386 386 386 37597 7597 7597 928 .0144928

# **MSAs: other approaches**

DIALIGN is a local alignment method

It constructs pairwise ad multiple alignments by comparing whole ungapped segments several residues long

The alignment is then constructed from pairs of equal-length gap-free segments (diagonals)

Many diagonals will overlap and the program has to find a set of diagonals which can be combined into one consistent alignment

DIALIGN is suitable for sequences of moderate length

## **Multiple alignments**

LUNTVYYGUPUN <mark>KE</mark> ATTTLFCASDA	AYDTEVHNUMATH ACUPTOP #P	EVAL GNOTEN FUMAR	NNMVEQ <mark>MQED</mark> I ISLW <mark>EQSLE</mark> F
LINTUYYGUPUN <mark>KE</mark> ATTTLFCASBAR	AVDTEUHNUWATH ACUPTEP NP (	EVOLUNOTEN FRAME	NDMUEQMHED I ISLWEQSLEF
LUUTUYYGUPUNNE AATTLECASDAK	AYDTEUHNUMATHACUPTIPNP	BOOLEMOTER FRMME	NEW JEOMHED I ISLUDOSLEF
LOOT OVYGOP OWNE AT TTL FC ASD AF	AY STE AND OWATH ACOPTAP AP (	POOL PROTEREMMON	MNAOLOMELD IISLOLOSLER
LIGTOTYCOPOINT ATTICCASE A	AN TELEVISION TRACOPTOP APO	ELONGNOTEN FRMINE	
LINCI OFFOUR ATTILICASE A	AVATES HUS MATHACIZET IN HIS	PLAT GRATEN PUMUN	NUM ROMORD ITSLUDOST
LINTUTY GUP MINE ATTTL PC ASD A	AVOTE ANNOUNT HACUPTHP HP	ENAL ERATER FRAME	MANDED AND THE TANK OF STATE
LWOTUYYGUPUMIL AETTL FCASDAR	AYDTERHNUWATH ACUPTOP NP	ELSLONUTERFOMME	NNMVEOMHEDVISLWDOSL KF
LINTVYYGVPUN <mark>KE</mark> ATTTL FCASE A <mark>R</mark>	S <mark>vete ann iwath</mark> acop tep np o	E I ALERUTERFRMOR	nnnveo <mark>nned</mark> i islandslæf
LWOTVYYGUPUM <mark>KE</mark> ATTTLFCASD <mark>A</mark> K	SYEPEAHN IWATHACUPTOPNP	E IEMENUTEN FRMM	NNMVEONHED I I SLWDORL RF
LIGTUTY GUP INKE ATTTLFCASDAK	AVARE ANN IWATH ACUPT IP NPO	E IELENUTENFRMM	NNMUEOMHED I ISLUDOSLEF
LUOTOYYGOPOWNE ATTTLFCASDAN	SYETEVHN IWATHACOPTIPNP	E IELENUTEGEMMON	NEW OF OMHED I ISLUD OSLEF
LOOTOYYCOPOOR A POLYCASI AR	ANSTEANN IWAT GACOPT OP SPO	TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTT	NNAOD QAHED I I SLOU QSLEF
GYOTOT FOT AND AS IP TO AT AN	PTNCT INCEPTION	P TOT N-STP & PT AND	TITES & LETS TONL FETS INC.
OVUTUFYGIPAN MASIPLFCATER	DTMGT IOCLPDNDDY	TTLN-UTEAFDAND	STOTES A LEDITAL FETS INF
QVUTUFYGUPUM NAS IPLFCATEM	BTWGT INC PDWDDY	E ITLN-VTE AFD AM	NTUTEQ AUEDUMSL FETS INF
Q FUTUFYG IP AW MAS IPL FCATERS		e itln-steaptawn	ntoteq ao <mark>edownl p</mark> ets i <mark>k</mark> f
QYUTUFYGUP AN MAT IPL FC AT KM	<mark>DTWGTTOCLPDNDD</mark> Y	ELALN-UTESFLAME	NTOTEQ AIEBOWOLFETS INF
OYUTUFYGIPAN MATUPLICATINE		EILN-ITEAFLAWD	NTOTO AUDITON L FETS INF
QUITOFYSOPOUNNESUGAFCMTPTT	LWATTNC IP DIRD YI	FUPLN-ITEP FE AMAD	NPLUAGAGEN INLL FEOTLER
A YOTOFYGIPOUNATIONAFCMTPHT	ST IN STATE OF THE	TOPLN - ITE AND AND	NPLOROALSE INLLE TOTME

Two sequences whisper, many homologous sequences talk loud

A. Lesk

#### Why are multiple alignments so important?

Because

they allow to obtain accurate alignments outline positions subjected to evolutionary pressure provide relevant functional/structural insight

![](_page_69_Picture_3.jpeg)

Provide information on the evolutionary process

#### Why are multiple alignments so important?

(A) p110αTFILGIGDRHNSNIMVKDDG-QLFHIDFGHFLDHKKKKFGYKRERVPFVLT--QDFLIVI 142cAMP-kinaseQIVLTFEYLHSLDLIYRDLKPENLLIDQQGYIQVTDFGFAKRVKGRTWXLCGTPEYLAPE 179

#### pairwise alignment of the catalytic domains of PI3-kinase p110 $\alpha$ and a cAMP-dependent protein kinase

(B)	p110β	SYVLGIGDRHSD <mark>NI</mark> NVKKT <mark>G</mark> QLFHI <mark>DFG</mark> HILGNFKSKFGIKRERVPFILT	136
	p110δ	TYVLGIGDRHSDNIMIRES <mark>G</mark> QLFHI <mark>DFG</mark> HFLGNFKTKFGINRERVPFILT	136
	p110α	TFILGIGDRHN <mark>SNIMV</mark> KDD <mark>G</mark> QLFHI <mark>DFG</mark> HFLDHKKKKFGYKRERVPFVLT	135
	p110γ	TFVLGIGDRHNDNIMITET <mark>G</mark> NLFHIDFGHILGNYKSFLGINKERVPFVLT	135
	p110_dicti	TYVLGIG <mark>D</mark> RHN <mark>DNLMV</mark> TKG <mark>G</mark> RLFHI <mark>DFG</mark> HFLGNYKKKFGFKRERAPFVFT	135
	cAMP-kinase	QIVLTFEYLHSLDLIYR <mark>D</mark> LKP <mark>ENLLI</mark> DQQ <mark>G</mark> YIQVT <mark>DFG</mark> FAKRVKGRTWXLCGTPEYLA	177

#### their ClustalW alignment with other PI3-kinases

In the multiple alignment, the functionally important residues (highlighted in green) are correctly aligned

# **One more example: thioredoxins**

Involved in cell proliferation, blood coagulation, insulin degradation, enzymatic regulation etc.

Fold  $\alpha/\beta$ :  $\beta$  sheet of five strands flanked by  $\alpha$ -helices

![](_page_71_Picture_3.jpeg)
# **One more example: thioredoxins**

(a)

Racharichia coli Porphyra perperaa Thiohacillus ferroomidans Straptompos clavuldans Cyanidiaschynon merolas Wiman Eksus monkey Sheep Rabhit Chicken Dictyostelium discoideum Dictyostelium discoideum

Racherichia coli

Rheutis soukey

Wenner

Sheep

Rabbit

Chicken

Potphyra purpursa Thiobhc illus' ferrooni dans

Streptomyces clavuligerus Cyanidi dechymon merciae

Dict was tell as discoldean

Diction telium discoldeum

Drosophila melanorastar Casnothabditis eDepars

Ricinus communis

Neurospora crassa



### Color code:

Gly, Ala, Ser, Thr : small Cys, Val, Ile, Leu, Pro, Phe, Tyr, Met, Trp : hydrophobic Asn, Gln, His : polar Asp, Glu : negatively charged Lys, Arg : positively charged



Racherichia coli Porphyra purpurea Thioharillus ferrooridaus Streptonyces clavuligerus Cyanisti dechyzon merci se Numu Rheutis poskey Sheep Rathit Chicken Dictycstelium discoddeum Dictycstelium discoddeum Drosophila melanogastar Cassorhabditis elegans Ricinus communis Neurospora crassa







## Profile sequence logo representation

Racherichia coli Porphyra purpursa Thiobacillus ferroonidans Streptomess claveligerus Cyanidi aschynon as rol as Rusa Rheutis poskey Sheep Rathit Chicken Dictypertelium discoddeum Dictypertelium discoddeum Drosophila melanoraster Casadrhabditis elegans Ricinus communis Neurospora crassa

Recherichia coli Porphyr a perpersa Thioline illus" f arrowidans Streptomyces clavuligerus Cyaridi dechyzon merciae Numer Rheutis poskey Sheep Rabbit Chicken Dictyostelium discoideum Dictyostelium discoideum Drosophila melanoraster Caenothabditis elegans Ricinus communis Neurospora crassa

- 0400

566666



(a)

(a)

Racherichia coli Porphyra purpures Thiobacillus fermonidaus Streptonyces clavuligenus Cyanidi aschyzon zerol as Brazi Rheutis poskey Sheep Rathit Chicken Dictyortelium discoideum Dictyortelium discoideum Drosophila melanogaster Caesophabditis elegans Ricinus communis Neurospora crassa

Recherichia coli Porphyr a perpersa Thiolise iffus' f arroad dans Streptomyces clavuligerus Cyaridi dechyzon merciae Numer Rheutis poskey Sheep Rabbit Chicken Dictyostelium discoideum Dictyostelium discoideum Drosophila melanoraster Caenothabditis elegans Ricinus communis Neurospora crassa

- 0400 -

538355





(a)

Racherichia coli Porphyra perpena Porphyra perpena Thiobacillus ferroanidans Streptonyces claveligerus Cyaridi oschyzon zerolae Wean Rheutis poskey Sheep Rabbit Chicken Dictyortelium discoideum Dictyortelium discoideum Drosophila melanoraster Casadrhabditis elegans Ricinus communis Neurospora crassa

Recherichia coli Porphyr a perpersa Thioline iffus' f arroani daus Streptomore clavuligerus Cyaridi dechyzon merciae Numer Rheutis poskey Sheep Rabbit Chicken Dictyostelium discoideum Dictyostelium discoideum Drosophila melanoraster Caenothabditis elegans Ricinus communis Neurospora crassa

- 0400





### Surface β-strand



Racherichia coli Porphyra purpurea Thioharillus ferroaridans Streptomyces clavuligenus Cyanidi aschyzon zerol as Brazi Rheutis poskey Sheep Rathit Chicken Dictypertelium discoddeum Dictypertelium discoddeum Drosophila melanogaster Caesophabditis elegans Ricinus communis Neurospora crassa

Recherichia coli Porphyr a perpersa Thioline illus" f arrowidans Streptomyces clavuligerus Cyaridi dechyzon merciae Numer Rheutis poskey Sheep Rabbit Chicken Dictyostelium discoideum Dictyostelium discoideum Drosophila melanoraster Caenothabditis elegans Ricinus communis Neurospora crassa

- 0400

538355



100100

8288



Racherichia coli Porphyra purpures Thiobacillus fermonidaus Streptonyces clavuligenus Cyanidi aschyzon zerol as Brazi Rheutis poskey Sheep Rathit Chicken Dictyortelium discoideum Dictyortelium discoideum Drosophila melanogaster Caesophabditis elegans Ricinus communis

Recherichia coli Porphyr a perpersa Thioline illus" f arrowidans Streptomores clavuligerus Cyaridi dechyzon merdiae Numer Rheutis poskey Sheep Raibhit Chicken Dictyostelium discoideum Dictyostelium discoideum Drosophila melanoraster Caenothabditis elegans Ricinus communis Neurospora crassa

5666 666 6



8588

ophobic hydrophilic 6 5 amphiphatic helix 203060

(a)

#### Lessons 5 & 6. Content

1. Alignment algorithms. There are exact and heuristic ones. We choose one or the other depending on the wanted application

2. Multiple alignments. They contain precious information about the evolutionary path. Non-exact methods are used to obtain them. HMMs can be applied. They are extremely informative on the structure and function of corresponding proteins