

SUBSTITUTION MODELS AND MEGA_(COOL) PHYLOGENETICS



DISTANCES BETWEEN SEQUENCES

AGATGCTAGCATCGACTAGCATCAGCTGACCCCGCGCGCGCAT

AGATAA TAGCATCGACTAGCATCAGCGGACCCCGCGCGCGCAC

How far are these a part?



P-DISTANCE

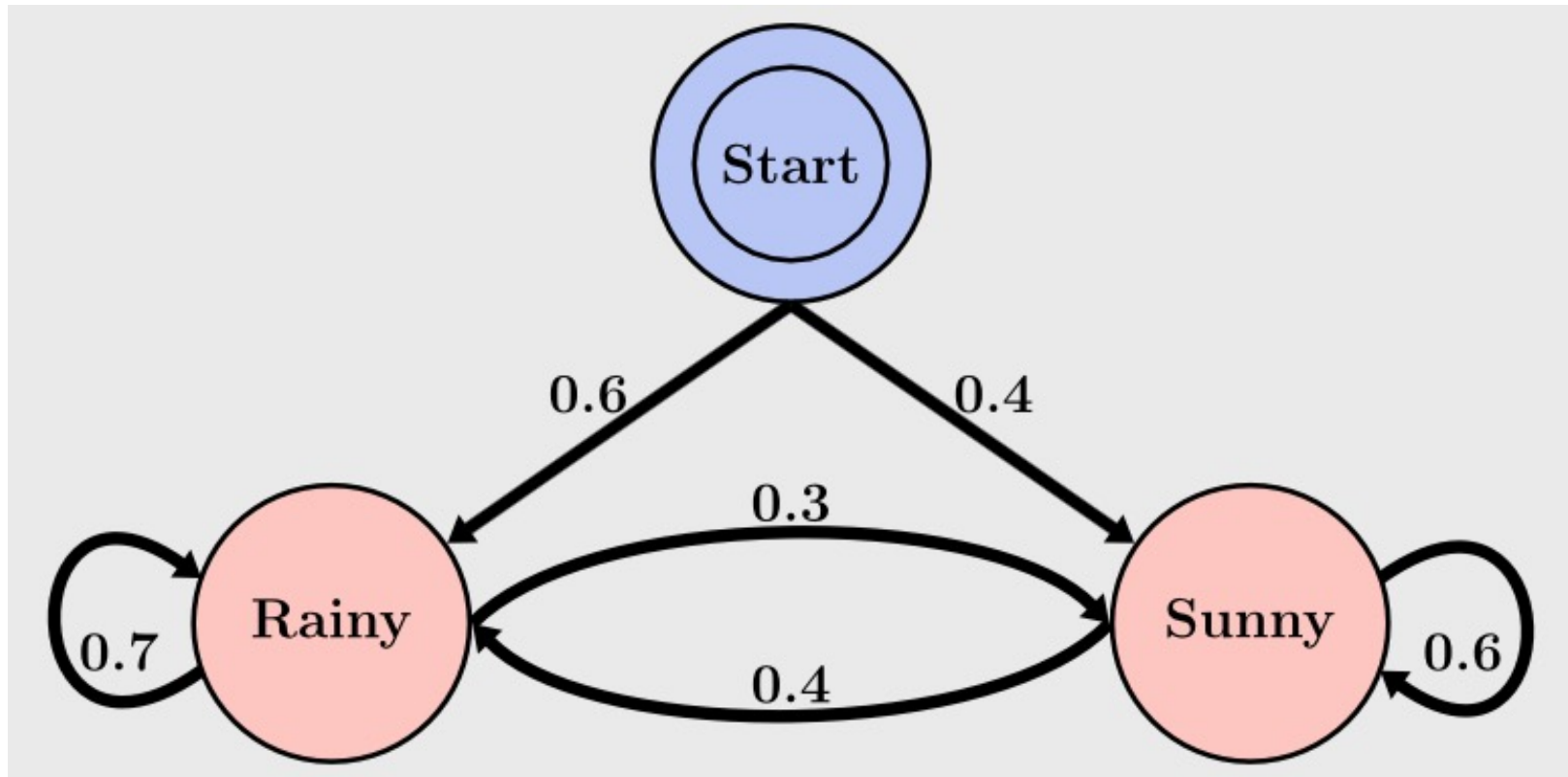
The proportion of differences between the two sequences

Differences / Total number of nucleotides compared

The sequences before were 43 long, and there were 4 differences

**What would be their p-distance?
When is this appropriate? When is it not?**

CONTINUOUS-TIME MARKOV CHAINS



BACK TO THE DISTANCES AND SEQUENCES

AGATGCTAGCATCGACTAGCATCAGCTGACCCCGCGCGCGCAT

AGAT**AA**TAGCATCGACTAGCATCAGC**G**GACCCCGCGCGCGCAC**C**

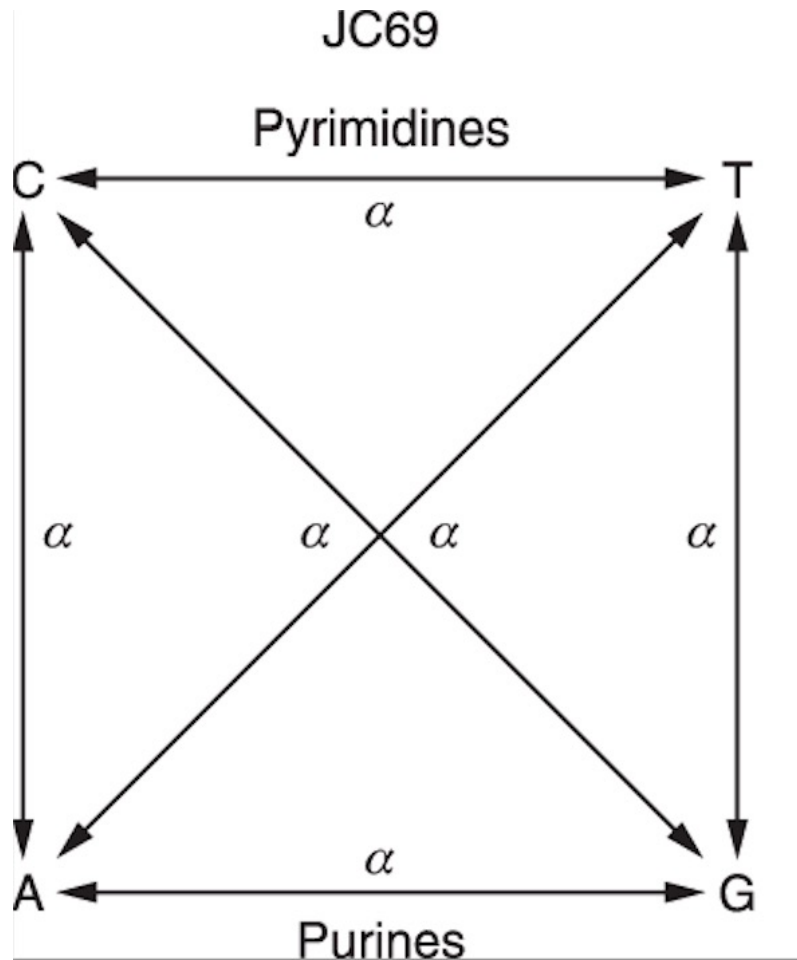
We want to quantify the unobserved, which is difficult.

We can do so by treating this as a Markov chain, but it will require us to assume somethings:

- Each site in the sequence are independent

Not always true but we will assume more unbelievable stuff later so never mind.

JUKES CANTOR (69)



$$P(t) = \{p_{ij}(t)\} \approx I + Qt = \begin{bmatrix} 1 - 3\lambda t & \lambda t & \lambda t & \lambda t \\ \lambda t & 1 - 3\lambda t & \lambda t & \lambda t \\ \lambda t & \lambda t & 1 - 3\lambda t & \lambda t \\ \lambda t & \lambda t & \lambda t & 1 - 3\lambda t \end{bmatrix}$$

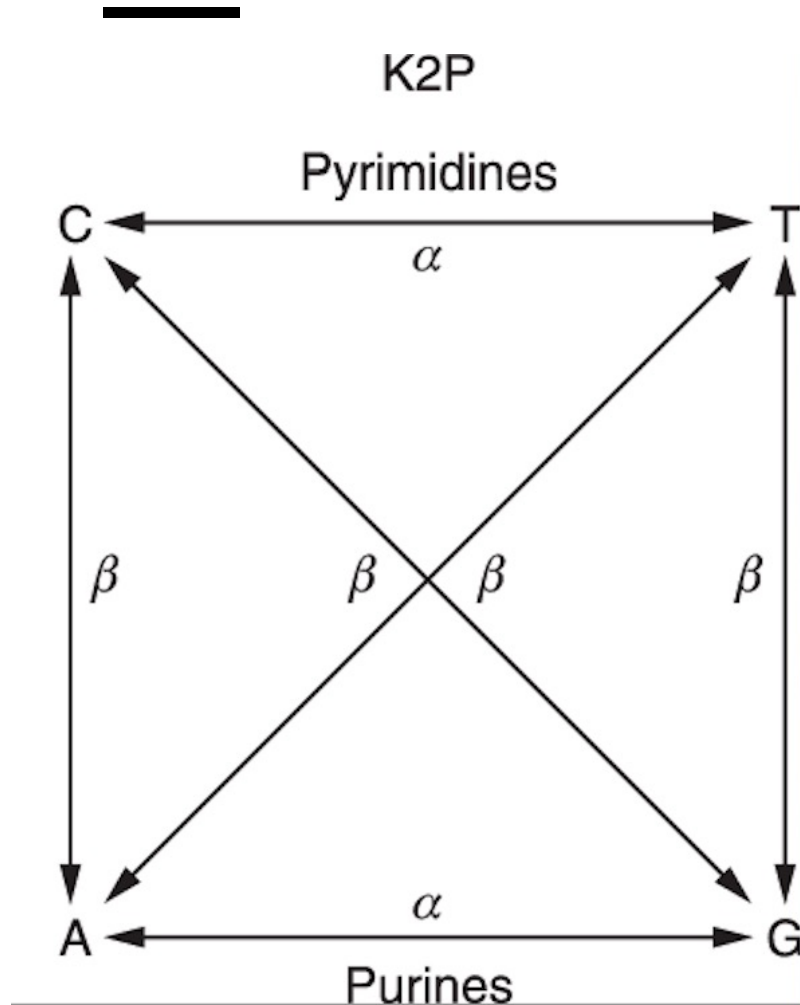
$$P(t) = e^{Qt} = \begin{bmatrix} p_0(t) & p_1(t) & p_1(t) & p_1(t) \\ p_1(t) & p_0(t) & p_1(t) & p_1(t) \\ p_1(t) & p_1(t) & p_0(t) & p_1(t) \\ p_1(t) & p_1(t) & p_1(t) & p_0(t) \end{bmatrix}, \text{ with } \begin{cases} p_0(t) = \frac{1}{4} + \frac{3}{4}e^{-4\lambda t}, \\ p_1(t) = \frac{1}{4} - \frac{1}{4}e^{-4\lambda t}. \end{cases} \quad (1.4)$$

$$\hat{d} = -\frac{3}{4} \log \left(1 - \frac{4}{3}\hat{p} \right), \quad (1.7)$$

What would the distance be between the two seqs now?

N=43
Differences=4

KIMURA (80) AKA K80 AKA AKA AKA K2P



$$Q = \begin{bmatrix} -(\alpha + 2\beta) & \alpha & \beta & \beta \\ \alpha & -(\alpha + 2\beta) & \beta & \beta \\ \beta & \beta & -(\alpha + 2\beta) & \alpha \\ \beta & \beta & \alpha & -(\alpha + 2\beta) \end{bmatrix}.$$

$$P(t) = e^{Qt} = \begin{bmatrix} p_0(t) & p_1(t) & p_2(t) & p_2(t) \\ p_1(t) & p_0(t) & p_2(t) & p_2(t) \\ p_2(t) & p_2(t) & p_0(t) & p_1(t) \\ p_2(t) & p_2(t) & p_1(t) & p_0(t) \end{bmatrix}, \quad (1.10)$$

where the three distinct elements of the matrix are

$$\begin{aligned} p_0(t) &= \frac{1}{4} + \frac{1}{4}e^{-4\beta t} + \frac{1}{2}e^{-2(\alpha+\beta)t} = \frac{1}{4} + \frac{1}{4}e^{-4\hat{d}/(\kappa+2)} + \frac{1}{2}e^{-2\hat{d}(\kappa+1)/(\kappa+2)}, \\ p_1(t) &= \frac{1}{4} + \frac{1}{4}e^{-4\beta t} - \frac{1}{2}e^{-2(\alpha+\beta)t} = \frac{1}{4} + \frac{1}{4}e^{-4\hat{d}/(\kappa+2)} - \frac{1}{2}e^{-2\hat{d}(\kappa+1)/(\kappa+2)}, \\ p_2(t) &= \frac{1}{4} - \frac{1}{4}e^{-4\beta t} = \frac{1}{4} - \frac{1}{4}e^{-4\hat{d}/(\kappa+2)} \end{aligned} \quad (1.11)$$

$$\begin{aligned} \hat{d} &= -\frac{1}{2} \log(1 - 2S - V) - \frac{1}{4} \log(1 - 2V), \\ \hat{\kappa} &= \frac{2 \times \log(1 - 2S - V)}{\log(1 - 2V)} - 1 \end{aligned} \quad (1.12)$$

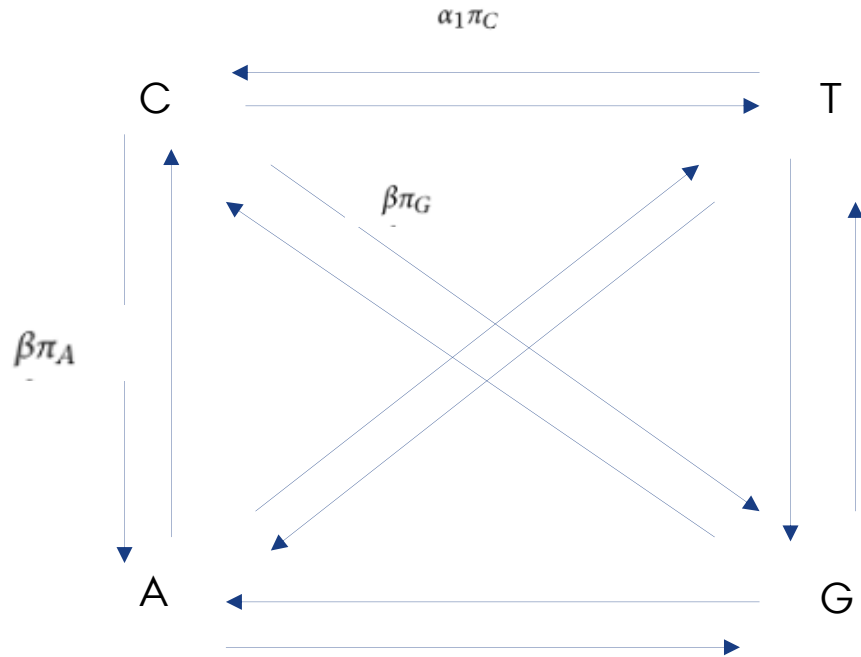
What would the distance be between the two seqs now?

N=43

Differences=4 (S=2, V=2)

S = transitions

V = transversions



$$Q = \begin{bmatrix} -(\alpha_1 \pi_C + \beta \pi_R) & \alpha_1 \pi_C & \beta \pi_A & \beta \pi_G \\ \alpha_1 \pi_T & -(\alpha_1 \pi_T + \beta \pi_R) & \beta \pi_A & \beta \pi_G \\ \beta \pi_T & \beta \pi_C & -(\alpha_2 \pi_G + \beta \pi_Y) & \alpha_2 \pi_G \\ \beta \pi_T & \beta \pi_C & \alpha_2 \pi_A & -(\alpha_2 \pi_A + \beta \pi_Y) \end{bmatrix}. \quad (1.16)$$

$$\hat{d} = \frac{2\pi_T \pi_C}{\pi_Y} (a_1 - \pi_R b) + \frac{2\pi_A \pi_G}{\pi_R} (a_2 - \pi_Y b) + 2\pi_Y \pi_R b,$$

What would the distance be between the two seqs now?

If you want you can do this one at home..

TN93, THE COMPLICATED CASE

Talk to your group for 2 minutes:

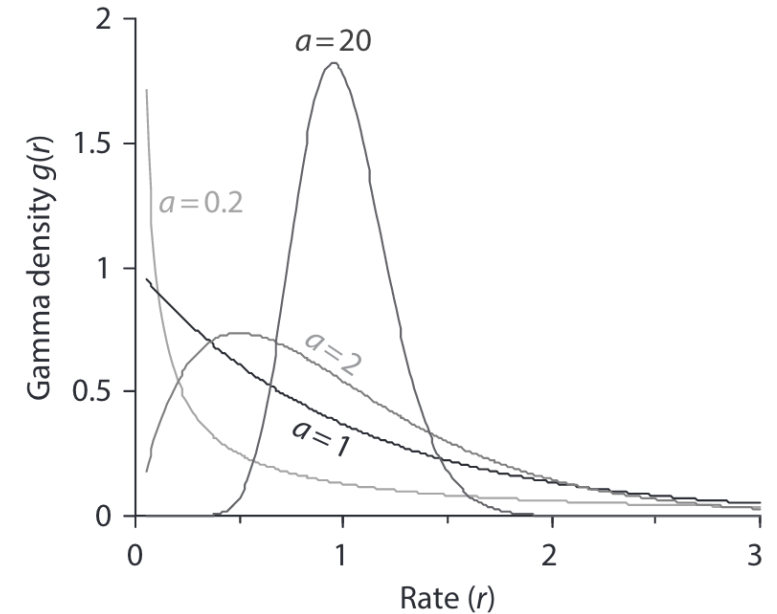
What would we have to change to make TN93 -> JC69?

What would we have to change to make TN93 -> K2P

(Fun) fact this is a nested model, which they'll cover in Datascience at some point.

EXTENTIONS OF THE MODELS

Not all sites are equally likely to change, we can adjust for this using a gamma distribution.



MAXIMUM LIKELIHOOD ESTIMATION (MLE)

Estimate the parameters for your model by picking the parameters that maximize the likelihood of seeing the data

It can get awefully complicated, but it is indeed very handy

EXERCISES

Now you have ~45 minutes to make the exercises for today

Ask questions if you have any

EVALUATION OF THE EXERCISES

The Sequence of the Human Genome

J. CRAIG VENTER, MARK D. ADAMS, EUGENE W. MYERS, PETER W. LI, RICHARD J. MURAL, GRANGER G. SUTTON, HAMILTON O. SMITH, MARK YANDELL, CHERYL A. EVANS, [...],

AND XIAOHONG ZHU

+264 authors

[Authors Info & Affiliations](#)

SCIENCE • 16 Feb 2001 • Vol 291, Issue 5507 • pp. 1304-1351 • [DOI: 10.1126/science.1058040](https://doi.org/10.1126/science.1058040)



1. Understanding PubMed and GenBank

GenBank is a useful database that contains DNA, RNA and protein sequences publicly available that we will access through MEGAX to download sequences of interest. You can read more about it [here](#).

In order to use GenBank efficiently, as other databases such as [PubMed](#) used to search for papers, it is a good idea to use specific search fields. These can be specified as e.g. [Author], [pdat] (publication time) and [Title] and use logical operators to combine search terms, e.g. AND, OR, NOT. For example, try to search "Zhao[Author] AND Wu F[Author] AND Yu[Author] AND coronavirus[Title] AND 2020[pdat]" on [PubMed](#). Which paper comes up?

By searching papers in [PubMed](#) and [GenBank](#) (and Google for sure), answer the following questions:

1. Find and download the paper of the first sequence of the human genome by the International Human Genome Sequencing Consortium or the one assembled by Craig Venter et al in 2001. Where was it published? How can the sequence be obtained?
2. Find the paper on the high coverage archaic Denisovan sequence published by Svante Pääbo's group published in 2012 and answer the same questions. Who were the Denisovans?
3. How many sequences from the dolphin (you might want to use the latin name *Delphinidae*) can you find? You can choose any other animal!
4. Find the Taxonomic position of Dolphins and find the number of DNA sequences deposited in GenBank for some of the dolphin species
5. Find sequences from FOXP2. What can you learn about this protein?

TAXONOMY

[Delphinidae](#)

Marine dolphins (*Delphinidae*) is a family of whale in pigs, camels etc.).

Taxonomy ID: [9726](#)

Taxonomy browser

Genomes

Items: 1 to 20 of 557525

The Sequence of the Human Genome

J. CRAIG VENTER, MARK D. ADAMS, EUGENE W. MYERS, PETER W. LI, RICHARD J. MURAL, GRANGER G. SUTTON, HAMILTON O. SMITH, MARK YANDELL, CHERYL A. EVANS, AND XIAOHONG ZHU +264 authors [Authors Info & Affiliations](#)

SCIENCE · 16 Feb 2001 · Vol 291, Issue 5507 · pp. 1304-1351 · DOI: 10.1126/science.1058040

15,734 9,244

This article has a correction.
Please see: [Corrections and Clarifications - 8 June 2001](#)

This article has a correction.
Please see: [Corrections and Clarifications - 22 February 2002](#)

Abstract

A 2.91-billion base pair (bp) consensus sequence of the euchromatic portion of the human genome was generated by the whole-genome shotgun sequencing method.

A High-Coverage Genome Sequence from an Archaic Denisovan Individual

MATTHIAS MEYER, MARTIN KIRCHER, MARIE-THERES GANSAUGE, HENG LI, FERNANDO RACIMO, SWAPAN MALLICK, JOSHUA G. SCHRAIBER, FLORA JAY, KAY PRÜFER, [...] AND SVANTE PÄÄBO +24 authors [Authors Info & Affiliations](#)

SCIENCE · 30 Aug 2012 · Vol 338, Issue 6104 · pp. 222-226 · DOI: 10.1126/science.1224344

4,744 1,202





Taxonomy Browser

Selected taxa	
Delphinidae (marine dolphins) Enter one or more taxonomic names	
Taxonomic name	Genomes
▼ <i>Eukaryota</i> (eukaryotes)	34,056
▼ <i>Metazoa</i> (animals)	12,392
▼ <i>Chordata</i> (chordates)	6,426
▼ <i>Mammalia</i> (mammals)	2,820
▼ <i>Artiodactyla</i> (even-toed ungulates)	415
▼ <i>Delphinidae</i> (marine dolphins)	26
➤ <i>Cephalorhynchus</i>	0
➤ <i>Delphinus</i>	3
➤ <i>Feresa</i>	0
➤ <i>Grampus</i>	2
➤ <i>Lagenodelphis</i>	0
➤ <i>Lagenorhynchus</i>	3
➤ <i>Lissodelphis</i>	0
➤ <i>Orcaella</i>	0
➤ <i>Orcinus</i>	3
➤ <i>Peponocephala</i>	0
➤ <i>Globicephala</i> (pilot whales)	1
➤ <i>Pseudorca</i>	0
➤ <i>Sotalia</i>	0
➤ <i>Sousa</i>	2
➤ <i>Stenella</i>	2
➤ <i>Steno</i>	1
➤ <i>Tursiops</i>	9



Example 2.1:

Launch the *Alignment Explorer* by selecting the **Align | Edit/Build Alignment** on the launch bar of the main *MEGA* window.

Select **Create New Alignment** and click **Ok**. A dialog will appear asking “Are you building a DNA or Protein sequence alignment?” Click the button labeled “**DNA**”.

From the *Alignment Explorer* main menu, select **Data | Open | Retrieve sequences from File**. Select the “hsp20.fas” file from the MEG/Examples directory.

Aligning Sequences by *ClustalW*

You can create a multiple sequence alignment in *MEGA* using either the *ClustalW* or Muscle algorithms. Here we align a set of sequences using the *ClustalW* option.

Example 2.2:

Open the alignment file (using the instructions above) hsp20.fas.

Select the **Edit | Select All** menu command to select all sites for every sequence in the data set.

Select **Alignment | Align by ClustalW** from the main menu to align the selected sequences data using the *ClustalW* algorithm. Click the “**Ok**” button to accept the default settings for *ClustalW*.

Once the alignment is complete, save the current *alignment session* by selecting **Data | Save Session** from the main menu. Give the file an appropriate name, such as “hsp20_Test.mas”. This will allow the current *alignment session* to be restored for future editing.

Exit the *Alignment Explorer* by selecting **Data | Exit Aln Explorer** from the main menu.

The screenshot displays the MEGA (Molecular Evolutionary Genetics Analysis) software interface. The main window shows a list of sequences and a DNA sequence alignment view. The 'ClustalW' button in the top toolbar is circled in red. A dialog box titled 'ClustalW Options' is open, showing the 'Alignment' section. The 'Pairwise Alignment' and 'Multiple Alignment' sections both have 'Gap Opening Penalty' set to 15.00 and 'Gap Extension Penalty' set to 6.66. The 'Matrix' section is collapsed. At the bottom of the dialog box, the 'OK' button is circled in red.

BLAST

BLAST (basic local alignment search tool) is an algorithm and program for comparing primary biological sequence information, such as the amino-acid sequences of proteins or the nucleotides of DNA and/or RNA sequences. A BLAST search enables a researcher to compare a subject protein or nucleotide sequence (called a query) with a library or database of sequences, and identify library sequences that resemble the query sequence above a certain threshold. [Wikipedia](#).

The screenshot displays the BLAST web interface. At the top, there is a toolbar with various icons for file operations and sequence manipulation. Below the toolbar, a text input field contains the sequence "BLAST selected sequence(s)". To the right of this field are two tabs: "DNA Sequences" and "Translated Protein Sequences". Below the tabs, there is a table with two columns: "Species/Abbrv" and "Sequence". The table contains two rows of data. The first row is labeled "1. Sequence 1" and the second row is labeled "2. NG_007491.3 Homo sapiens forkhead box P2 (FOXP2) RefSeqGene on chromoso". The sequence "A" is highlighted in yellow. Below the table, there is a blue banner with the NIH logo and the text "National Library of Medicine National Center for Biotechnology Information". To the right of the banner is a "Log in" button. Below the banner, there is a navigation bar with the text "BLAST®" and links to "Home", "Recent Results", "Saved Strategies", and "Help". Below the navigation bar, there is a section titled "Format Request". Under this section, there is a table with the following data:

Request ID	F14EZJWP016
Status	Searching
Time since submission	00:00:00

Below the table, there is a message: "This page will be automatically updated in 1 seconds until search is done".

EVALUATION



AARHUS
UNIVERSITY
DEPARTMENT OF MOLECULAR BIOLOGY AND GENETICS

TA
29 AUGUST 2023

BJARKE MEYER PEDERSEN
PH.D STUDENT



—

JOIN THE BARC FACEBOOK TODAY





AARHUS
UNIVERSITY