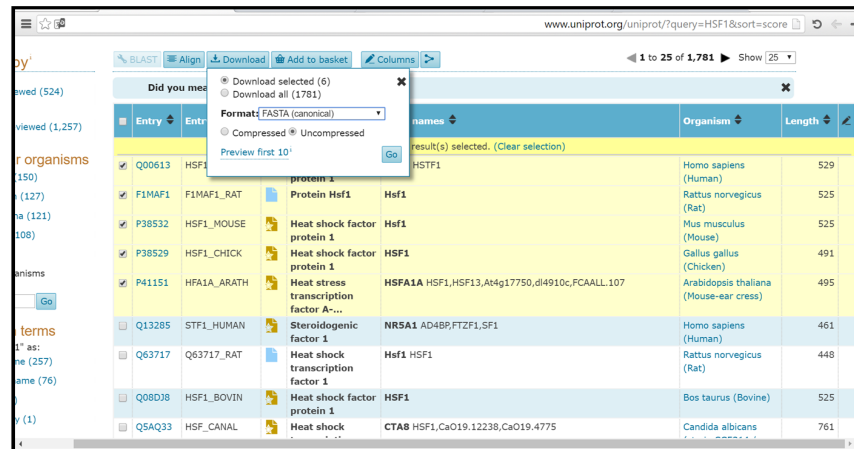# Introduction to Bioinformatics Online Course: IBT

**Multiple Sequence Alignment**

**Building Multiple Sequence Alignment**

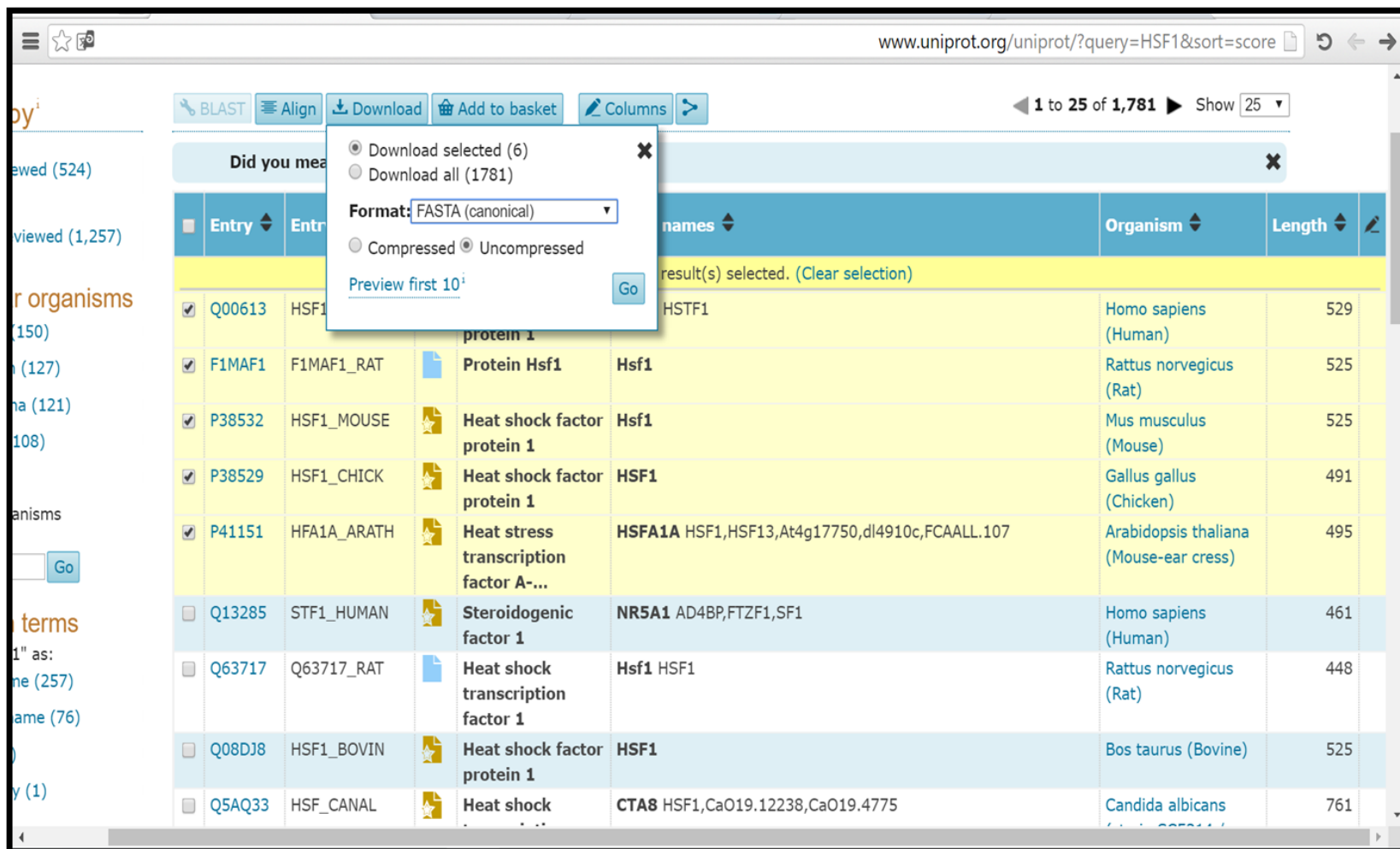**Lec2  Choosing the Right Sequences**

# Choosing the Right Sequences

"Before you build your alignment, you must **carefully select the sequences** you want to align. These sequences are members of the **same protein family**, and they all **share a common ancestor**. The family is usually **too large** to be entirely included in your multiple alignment, and **picking the right sequences is an art**."



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2ⁿᵈ Edn). *Wiley publishing, Inc*. 436 pp.

Introduction to Bioinformatics Online Course:IBT
Multiple Sequence Alignment| Prof. Ahmed M. Alzohairy

# Retrieve Sequences from Uniprot (www.uniprot.org

# Retrieve Sequences from NCBI (https://www.ncbi.nlm.nih.gov/protein)

# A Few Guidelines for Selecting Sequences

# Proteins or DNA

Use **proteins whenever possible**. You can turn them back into DNA after doing the multiple alignment.

If the sequences are **non-coding sequences**, you must **use DNA**

**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

H3ABioNet
Pan African Bioinformatics Network for H3Africa

# A Few Guidelines for Selecting Sequences

## Many sequences

Start with 10–15 sequences; **avoid** aligning **more than 50** sequences.



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

# A Few Guidelines for Selecting Sequences

## Very different sequences

Sequences that are **less than 30 percent** identical to more than half the other sequences in the set often **cause troubles**.



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

# Identical sequences

They **never help**. Unless you have a very good reason to do so, **avoid incorporating** into your multiple alignment any sequence that's more than **90 percent identical** to another sequence in the set.
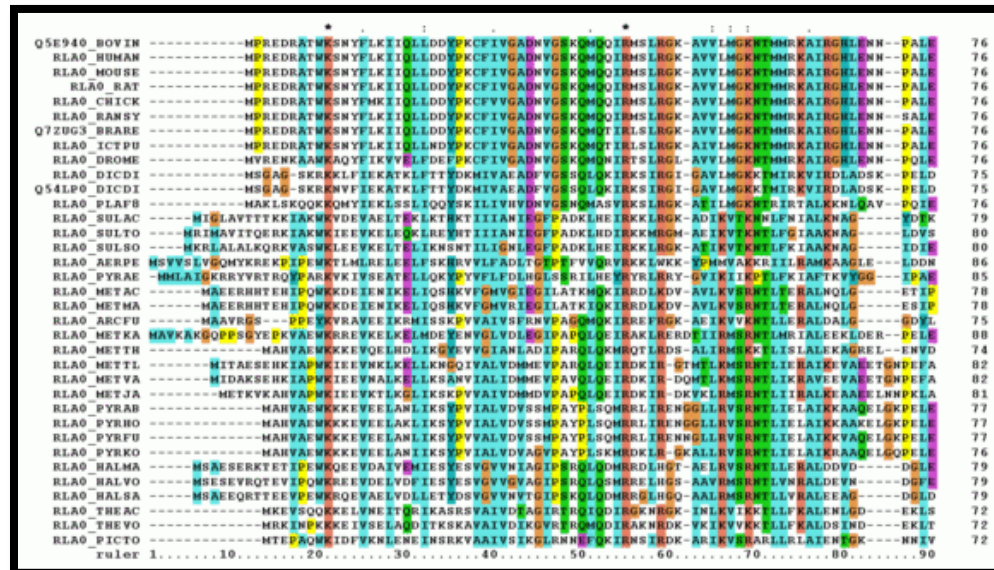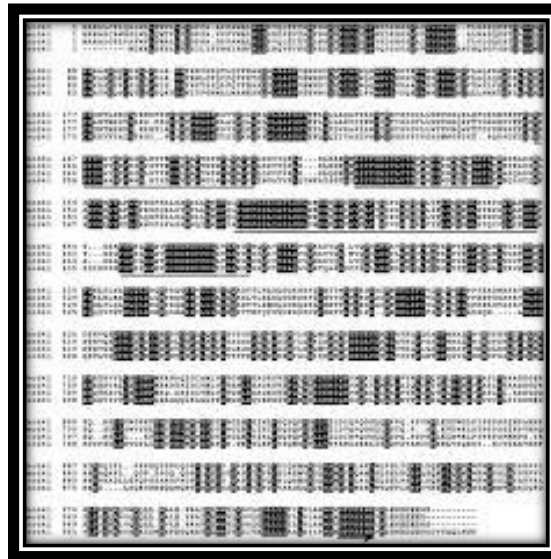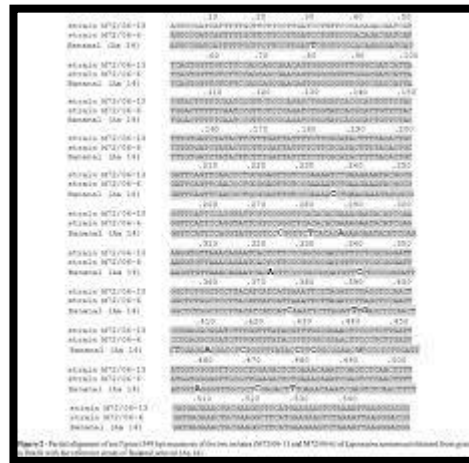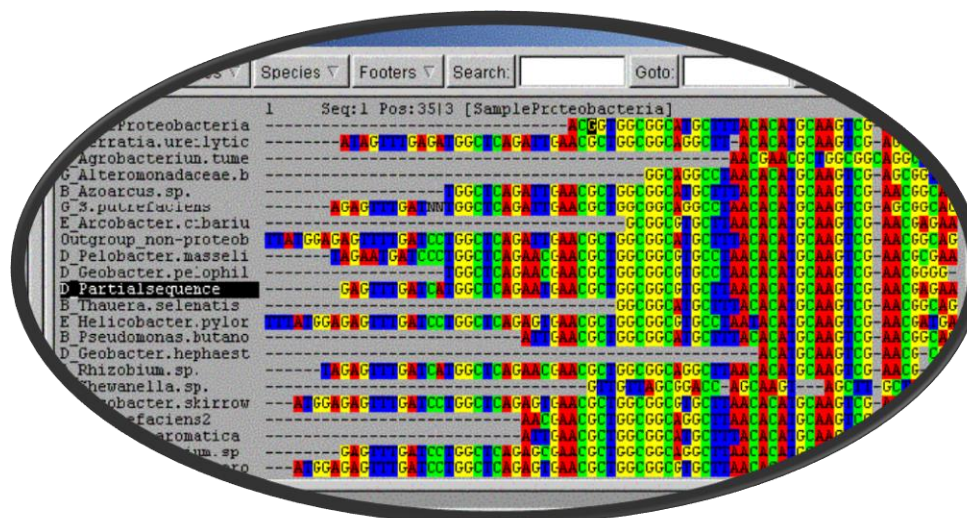


**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

# A Few Guidelines for Selecting Sequences

## Partial sequences

Multiple-sequence-alignment programs prefer sequences that are roughly the same length. Programs often have **difficulties comparing** items in a **mixture of complete sequences** and **shorter fragments**.
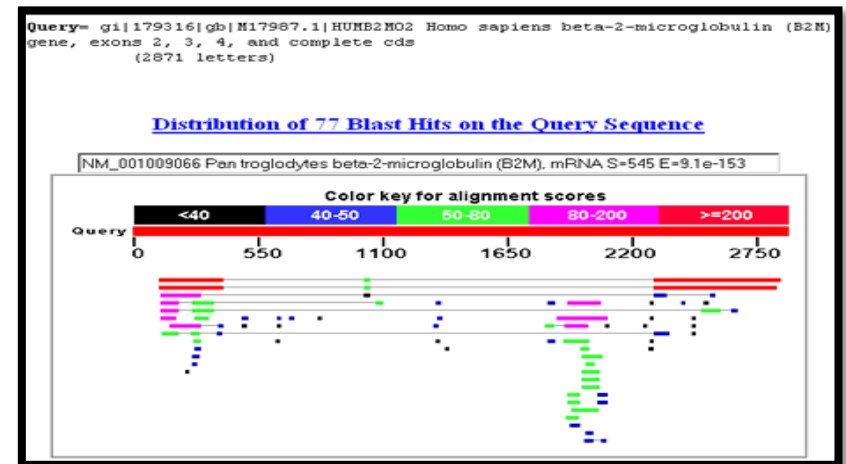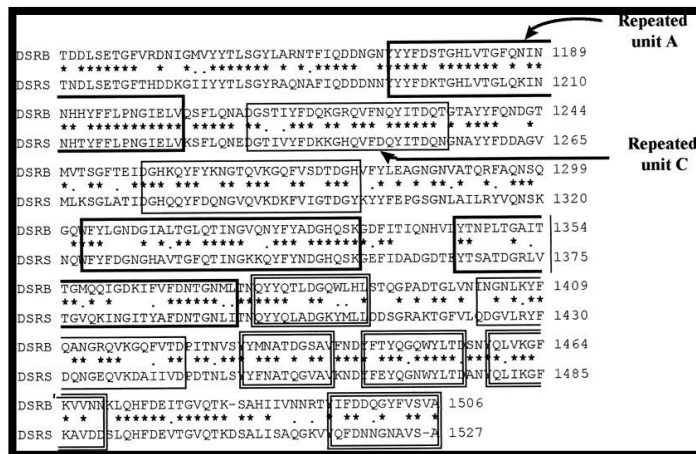
Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

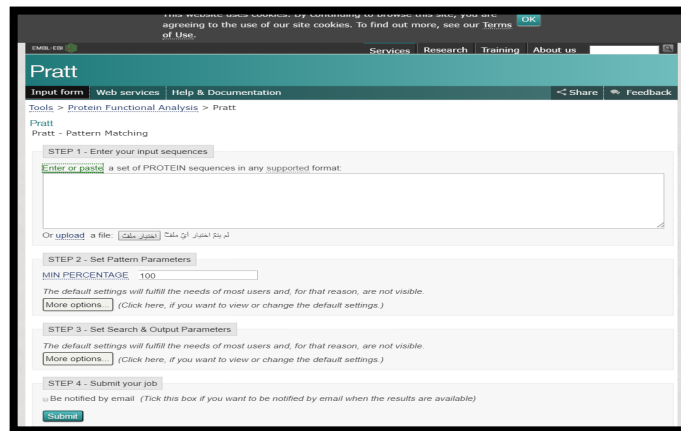# A Few Guidelines for Selecting Sequences

## Repeated domains

Sequences with repeated domains **cause trouble** for most multiple-alignment programs — especially if the **number of domains is different**. When this happens, you may be better off extracting the domains yourself with Dotlet or Lalign and making a multiple alignment of those segments.

**Claverie J, Notredame C** (2007). *Bioinformatics for Dummies* (2nd Edn). *Wiley publishing, Inc*. 436 pp.

# A Few Guidelines for Selecting Sequences

If you still **cannot generate a proper alignment** from sequences that you know are related, you could **use a local** multiple alignment method, such as the **Gibbs sampler**, or a pattern extraction motif, such as **Pratt**.



**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

H3ABioNet
Pan African Bioinformatics Network for H3Africa

Introduction to Bioinformatics Online Course:IBT
Multiple Sequence Alignment| Prof. Ahmed M. Alzohairy

For carrying out a **phylogenetic analysis** on a set of **coding DNA** sequences, do the following:

1.    **Translate** your **DNA** sequences **into proteins**.

2.    **Perform** the **multiple alignments** on the proteins.

3.    **Thread** the **DNA** back onto the protein multiple sequence alignment framework using **pal2nal** (coot.embl.de/pal2nal) or **Protogene** if you do not have the original DNA sequence (www.tcoffee.org).



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc.* 436 pp.
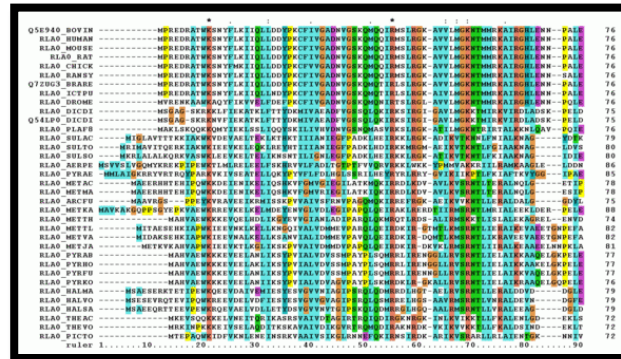
# Choosing the right number of sequences

**1- you should start with a relatively small number of sequences — between 10 and 15 sequences would be suitable for most cases.**

**2- After you get something interesting happening with this small set, you can always increase its size.**

**3- it's hard to see any reason for generating a multiple alignment with more than 50 sequences, unless you're interested in building some extensive phylogenetic tree.**

# Why you should not use too much sequences to align?

1- *Computing* big alignments is difficult.

2- *Building* big alignments is difficult

3- *Displaying* big alignments is difficult.

4-*Using* big alignments is difficult.

5-Making *accurate* big alignments is difficult



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

**Should you choose sequences that are very similar or very different?**

- Make the right **compromise** between **similarity** and **new information**
- An alignment that only contains **very similar** sequences brings **little information**.
- You can use it to **extrapolate annotations**, but you **can't do phylogeny**, **structure prediction**, **function perdiction,** or any of the other useful applications that we mentioned before
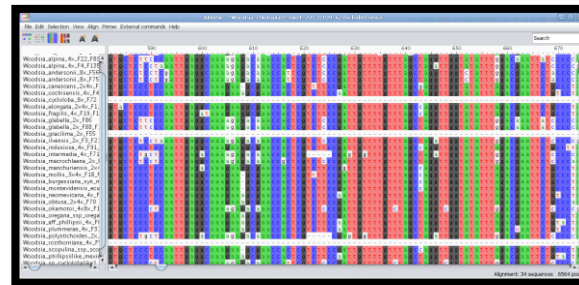
**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

# The general rule is that you want them to be **as distantly related as possible - without** requiring too many **gaps in** order to be properly **aligned**.

**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

**H3ABioNet**
Pan African Bioinformatics Network for H3Africa

Introduction to Bioinformatics Online Course:IBT
Multiple Sequence Alignment| Prof. Ahmed M. Alzohairy

**Two things multiple-sequence-alignment programs *really* don't like are**

## 1- Sequences that are very different from every other sequence in the group
## 2- Sequences that need long insertions/deletions to be properly aligned



**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2[nd] Edn). *Wiley publishing, Inc*. 436 pp.

H3ABioNet
Pan African Bioinformatics Network for H3Africa

Introduction to Bioinformatics Online Course:IBT
Multiple Sequence Alignment| Prof. Ahmed M. Alzohairy

# Gathering your sequences with online BLAST servers

**Characterized sequences:** Try to include sequences with **good annotations and experimental information** in your alignment because they bring biological information with them — and also **allow feature propagation**.

**Uncharacterized sequences:** including them in your multiple alignment is to **distinguish between the conserved positions** that cannot mutate and the other, less-important columns. They help in **getting some contrast** on your sequence of interest.



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

**Sequences** that are **so similar** to the query are **probably homologous**. We commonly refer to such sequences **as *hits*** or ***matches***.



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

# Choosing the Right Method of Multiple Sequence Alignment
# <u>ClustalOmega</u>

The **most commonly** used multiple sequence alignment package. 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



Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

H3ABioNet
Pan African Bioinformatics Network for H3Africa

# 2-Tcoffee

One of the latest multiple-sequence-alignment packages that you can use. With Tcoffee, you **can combine sequences** and **structures**, **evaluate an alignment**, or **merge** several alternative multiple alignments into a single unified result.



**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

Claverie J, Notredame C (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc.* 436 pp.

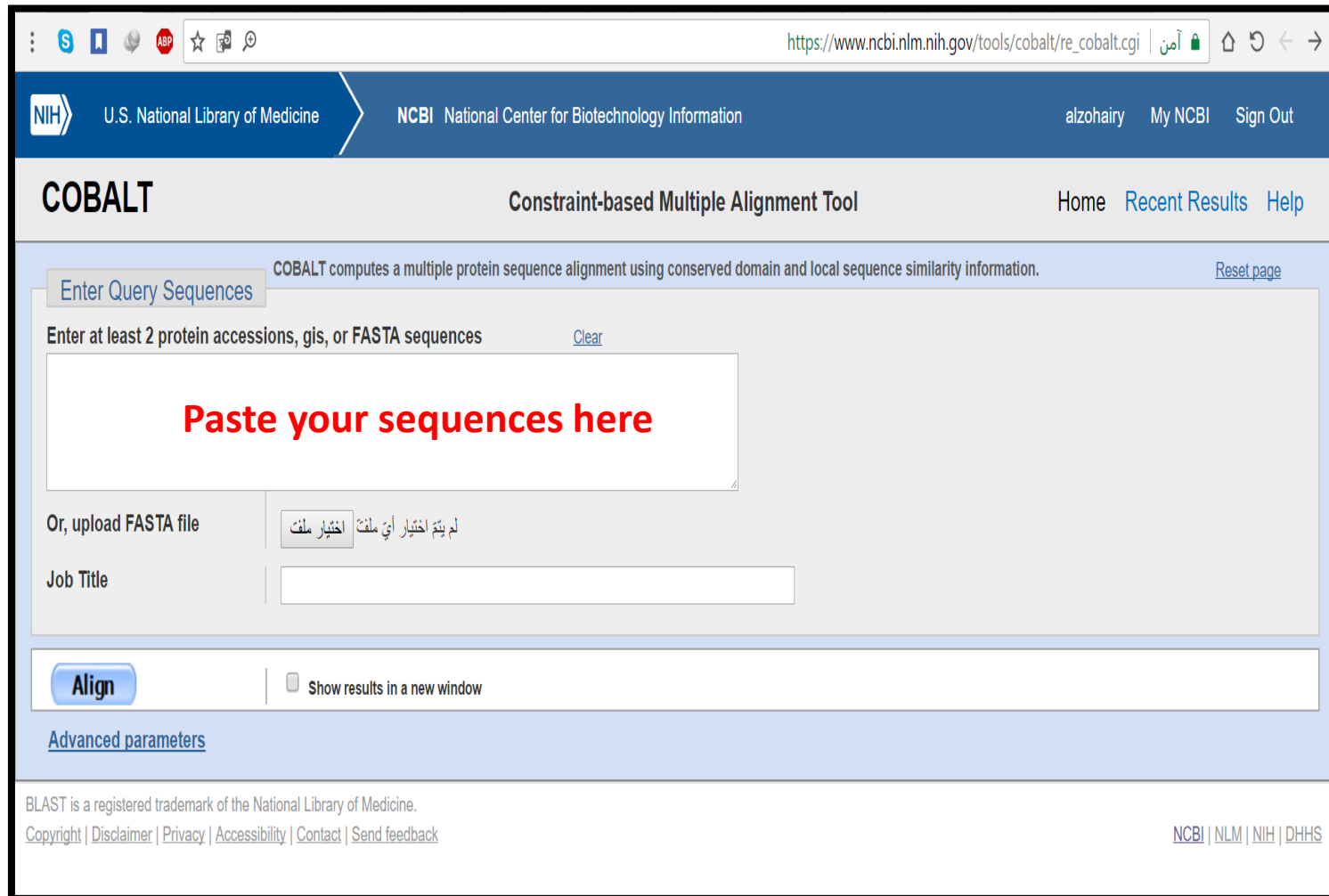# Choosing the Right Method of Multiple Sequence Alignment

## <u>3- MUSCLE</u>

One of **the fastest** alignment methods around for aligning large set of sequences

**Claverie J, Notredame C** (2007). Bioinformatics for Dummies (2nd Edn). *Wiley publishing, Inc*. 436 pp.

![H3ABioNet] **H3ABioNet**
**Pan African Bioinformatics Network for H3Africa**

Introduction to Bioinformatics Online Course:IBT
Multiple Sequence Alignment| Prof. Ahmed M. Alzohairy

# 4- COBALT

## Constraint-based Multiple Alignment Tool

H3ABioNet
Pan African Bioinformatics Network for H3Africa

By

Ahmed Mansour Alzohairy

**Department of Genetics, Zagazig University, Zagazig, Egypt**