# Evolutionary Multiple Sequence Alignment
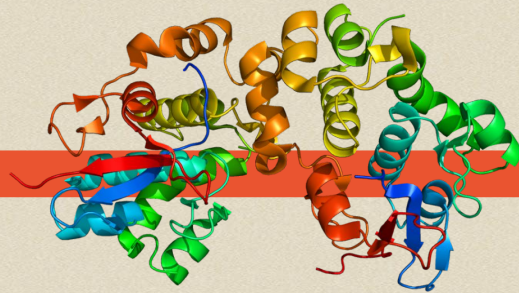
Sanjana Bhagavatula, John Drohan, Rachel Utama, Sreevatsa Nukala

DS3500 Final Project

# Multiple sequence alignment (MSA)

- **Alignment of 3 or more biological sequences of similar length**

| A | T | T | G | C | C | A | T | T |
|---|---|---|---|---|---|---|---|---|
| A | T | G | G | C | C | A | T | T |
| A | T | C | C | A | A | T | T | T | T |
| A | T | C | T | T | C | T | T |
| A | C | T | G | A | C | C |

| A | T | T | G | C | C | A | T | T | - | - |
|---|---|---|---|---|---|---|---|---|---|---|
| A | T | G | G | C | C | A | T | T | - | - |
| A | T | C | - | C | A | A | T | T | T | T |
| A | T | C | T | T | C | - | T | T | - | - |
| A | C | T | G | A | C | C | - | - | - | - |

- **Process of minimizing gaps and mismatches**

- **Detect regions of variability or conservation in a family of proteins**

- **May indicate functional, structural, or evolutionary relationships between biological sequences**

# Phylogenetic Tree

# Implementation

Evolutionary Multiple Sequence Alignment

# Read Fasta File

```
1   >BAA78379.1 P53 [Canis lupus familiaris]
2   MQEPQSELNIDPPLSQETFSELWNLLPENNVLSSELCPAVDELLLPESVVNWLDEDSDDAPRMPATSAPT
3   APGPAPSWPLSSSVPSPKTYPGTYGFRLGFLHSGTAKSVTWTYSPLLNKLFCQLAKTCPVQLWVSSPPPP
4   NTCVRAMAIYKKSEFVTEVVRRCPHHERCSDSSDGLAPPQHLIRVEGNLRAKYLDDRNTFRHSVVVPYEP
5   PEVGSDYTTIHYNYMCNSSCMGGMNRRPILTIITLEDSSGNVLGRNSFEVRVCACPGRDRRTEEENFHKK
6   GEPCPEPPPGSTKRALPPSTSSSPPQKKKPLDGEYFTLQIRGRERYEMFRNLNEALELKDAQSGKEPGGS
7   RAHSSHLKAKKGQSTSRHKKLMFKREGPDSD
8   >BAC16799.1 P53 [Homo sapiens]
9   MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEAPRMPEAA
10  PRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYSPALNKMFCQLAKT
11  CPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLAPPQHLIRVEGNLRVEYLDDRN
12  TFRHSVVVPYEPPEVGSDCTTIHYNYMCNSSCMGGMNRRPILTIITLEDSSGNLLGRNSFEVHVCACPGR
13  DRRTEEENLRKKGEPHHELPPGSTKRALSNNTSSSPQPKKKPLDGEYFTLQIRGRERFEMFRELNEALEL
14  KDAQAGKEPGGSRAHSSHLKSKKGQSTSRHKKLMFKTEGPDSD
15  >AAA39883.1 p53 [Mus musculus]
16  MTAMEESQSDISLELPLSQETFSGLWKLLPPEDILPSPHCMDDLLLPQDVEEFFEGPSEALRVSGAPAAQ
17  DPVTETPGPVAPAPATPWPLSSFVPSQKTYQGNYGFHLGFLQSGTAKSVMCTYSPPLNKLFFQLAKTCPV
18  QLWVSATPPAGSRVRAMAIYKKSQHMTEVVRRCPHHERCSDGDGLAPPQHLIRVEGNLYPEYLEDRQTFR
19  HSVVVPYEPPEAGSEYTTIHYKYMCNSSCMGGMNRRPILTIITLEDSSGNLLGRDSFEVRVCACPGRDRR
20  TEEENFRKKEVLCPELPPGSAKRALPTCTSASPPQKKKPLDGEYFTLKIRGRKRFEMFRELNEALELKDA
21  HATEESGDSRAHSSLQPRAFQALIKEESPNC
```

# Sum-Pairs Score Matrix

|       | S₁  | S₂  | S₃  | S₄  | S₅  |      |
|-------|-----|-----|-----|-----|-----|------|
| **S₁** | -   | 7   | -2  | 0   | -3  | **2** |
| **S₂** | 7   | -   | -2  | 0   | -4  | **1** |
| **S₃** | -2  | -2  | -   | 0   | -7  | **-11** |
| **S₄** | 0   | 0   | 0   | -   | -3  | **-3** |
| **S₅** | -3  | -4  | -7  | -3  | -   | **-17** |
|       | **2** | **1** | **-11** | **-3** | **-17** | |

S₁ is the sequence most similar to the rest, and below are the best alignments between S₁ and the rest of the sequences.

# Implementation

**Evo framework**

- **Fitness criteria**
    - **Sum-pairs scores**
    - **BLOSUM Matrix**
- **Modification Agents**
    - **Smith waterman algorithm**
- **Evolve**

- **Visualizations**
- **Alignment Library**



Matrix *T* looks like this, with the pink traceback:

Alignment:

G T T G
| | | |
G T T G

(Pink traceback)



BLOSUM 62 scoring matrix

(positive values are shaded)

The values for amino acid substitutions were obtained from Henikoff S & Henikoff JG (1992) Amino acid substitutions matrices from protein blocks. *Proc. Natl. Acad. Sci.* **89**: 10915-10919.
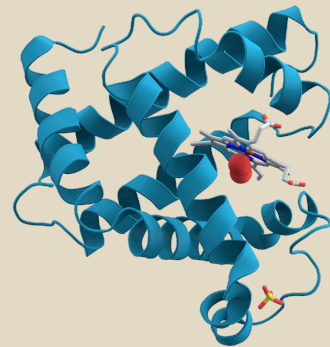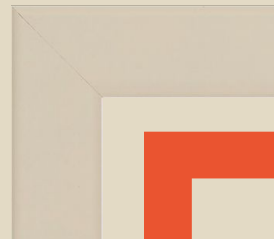
# Preliminary Results

# **Potential Next Steps**
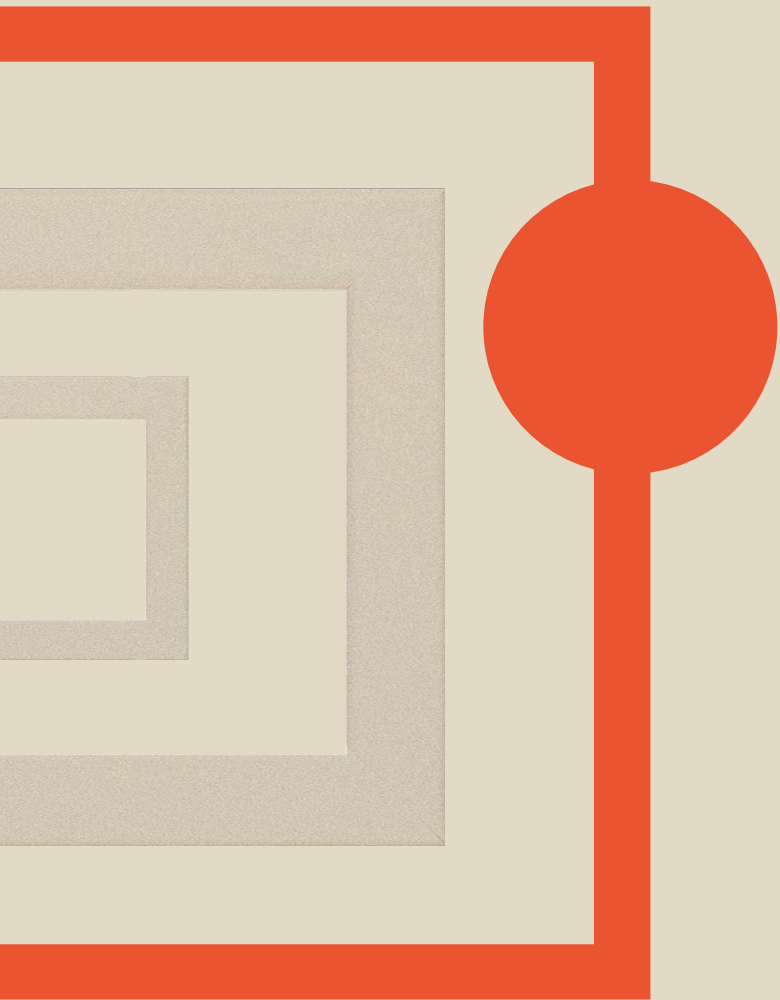


1) **Add more Agents**

    a) Modification Agents -

        i) Needleman-Wunsch

2) **Add more Fitness Criteria**

3) **Convert Amino acid to numerical representation**

    a) Enable faster run time

    b) More capability to compare amino acids on a spectrum

# References

1. https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-020-3526-6
2. https://www.ebi.ac.uk/Tools/msa/
3. https://www.cs.princeton.edu/~mona/Lecture/msa1.pdf
4. https://d1w9csuen3k837.cloudfront.net/Pictures/2000x2000fit/7/9/6/138796_Hydrophobic-and-polar-amino-acids.jpg
5. https://d1w9csuen3k837.cloudfront.net/Pictures/2000x2000fit/7/9/6/138796_Hydrophobic-and-polar-amino-acids.jpg

# Thank you!

Does anyone have any questions?