# Analysis of Mnx1 protein, its homologous sequences, phylogenetic tree and mutations

G. Görkem Köse - 25359
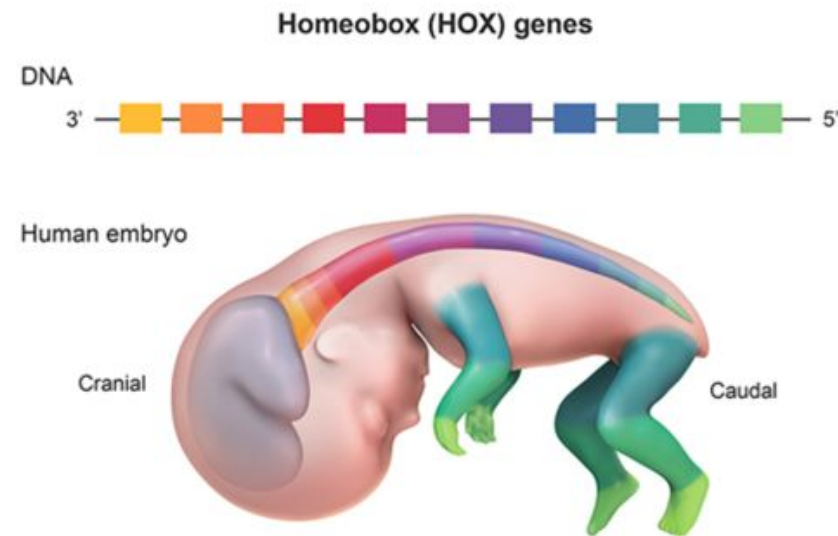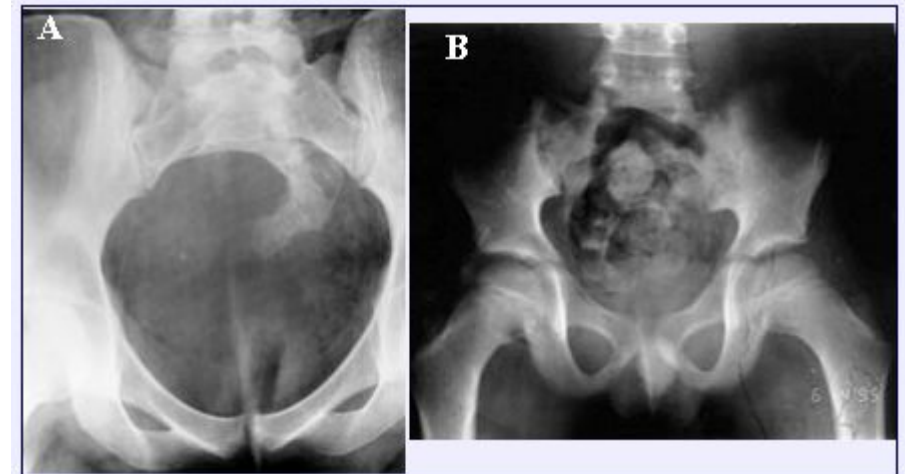
Ahmet Ölçüm - 25915

Sayed Damon Sadraije Najafi - 26260
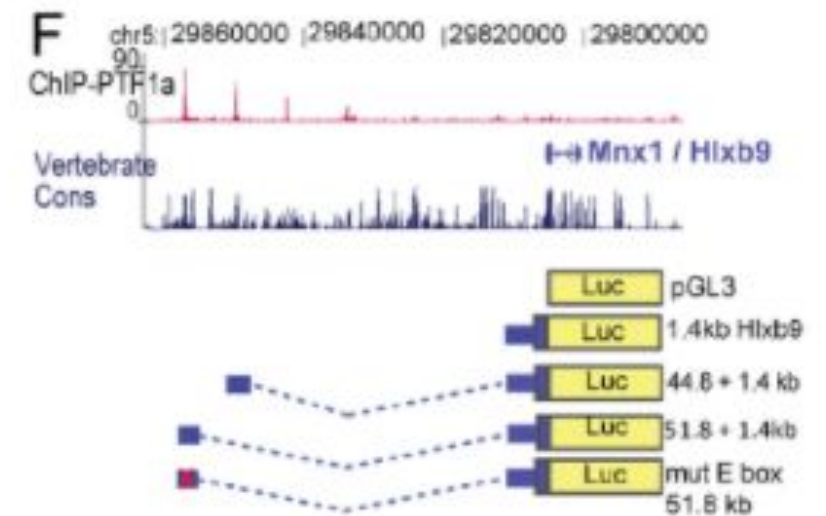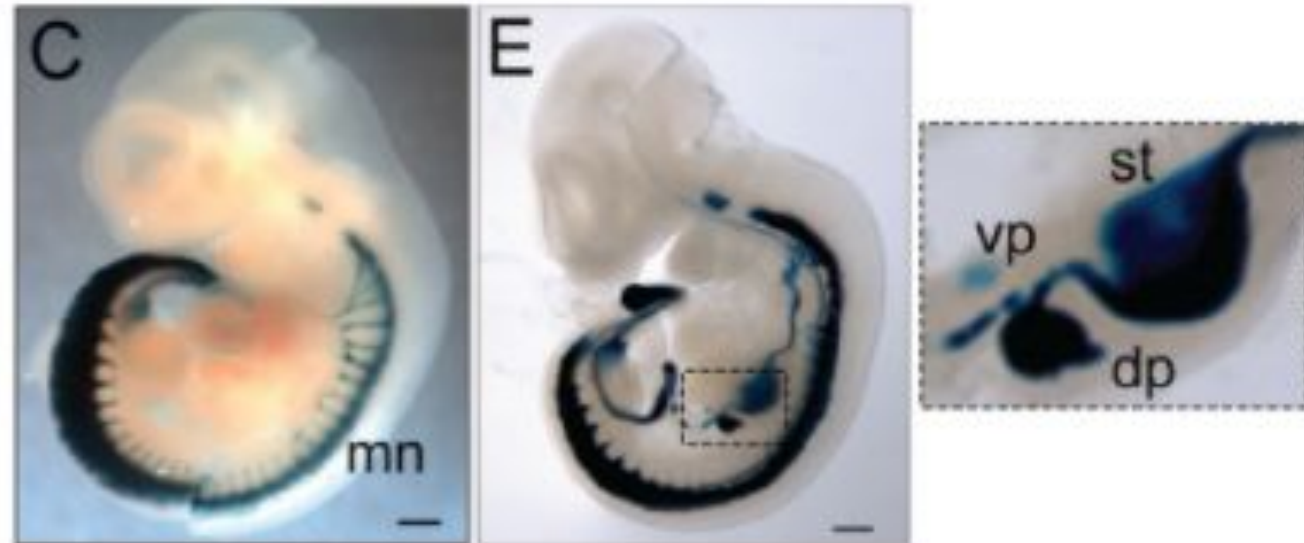
Emir Can Ülkü - 26740

# Currarino Syndrome

- Teratoma
- Hamartoma
- Neurenteric cyst
- Anterior meningocele
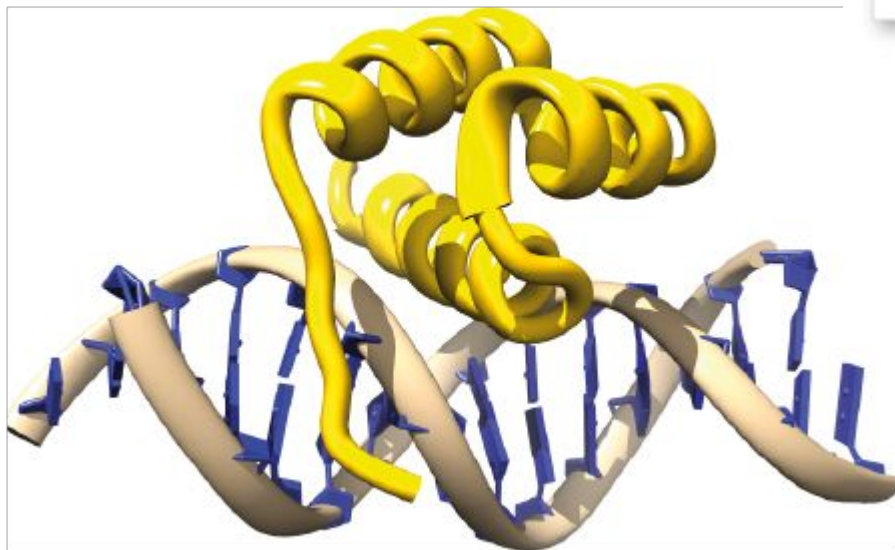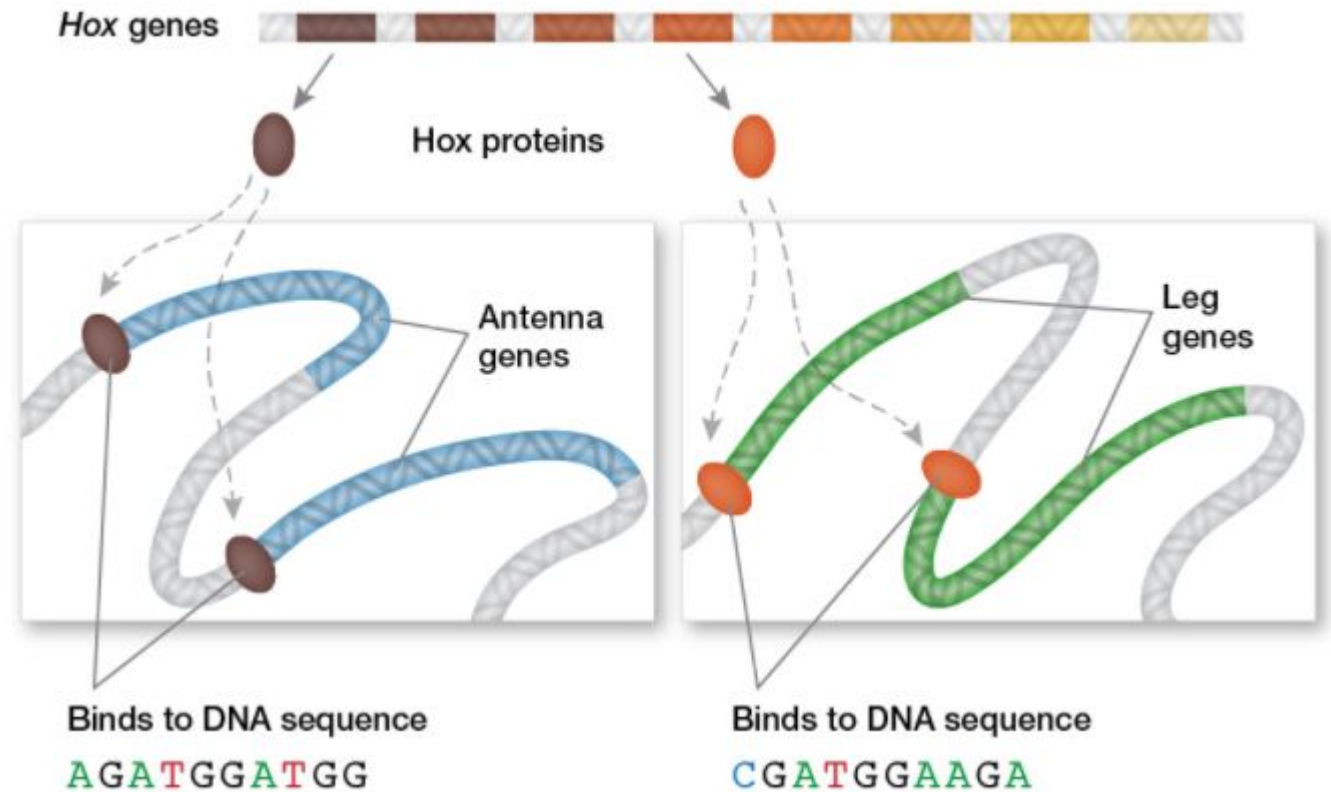


Homeobox (HOX) genes

# Mnx1 Protein

- Nuclear Protein
- Homeodomain
- Hox Protein Like
- Pancreas
- Motor-neuron

# Homeodomain Proteins

- Gene expression regulation
- Developmental Proteins



Hox genes

Hox proteins

Antenna genes

Binds to DNA sequence
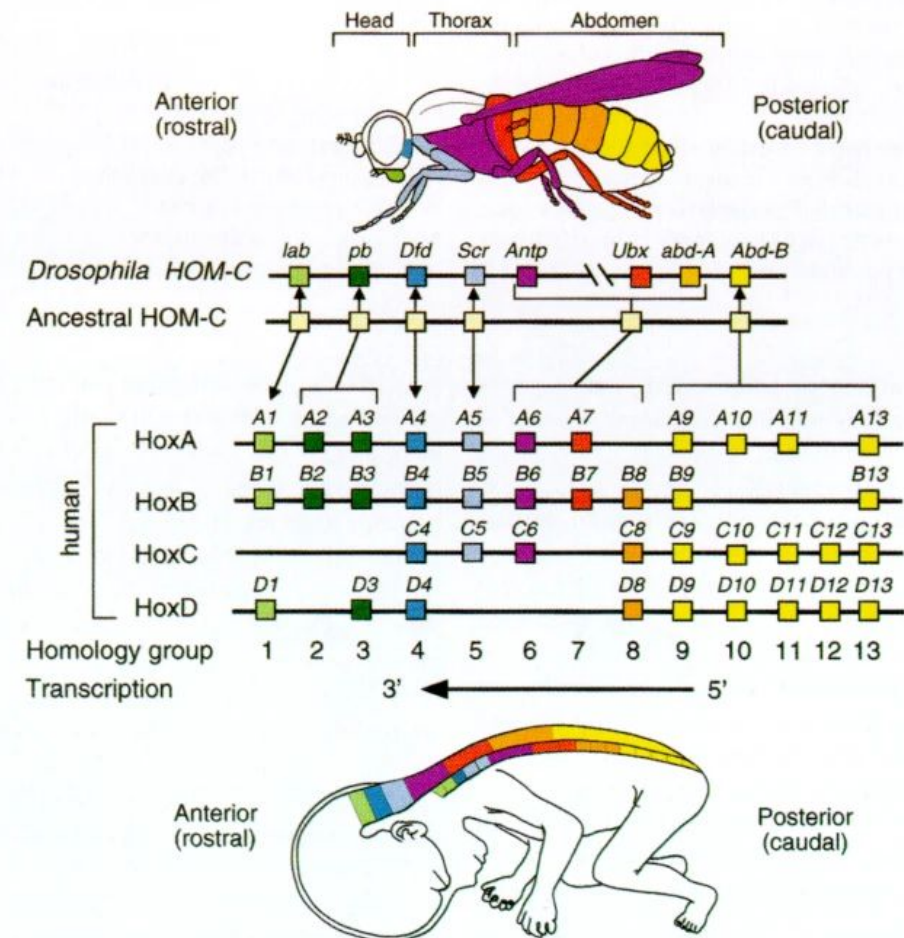AGATGGATGG

Leg genes

Binds to DNA sequence
CGATGGAAGA

# Homeodomain Proteins

- Gene expression regulation
- Developmental Proteins



Top: (Left) Normal fruitfly; (Right) Fruitfly with mutation in antennapedia gene Bottom: (Left) Normal fruitfly; (Right) Fruitfly with a homeotic mutation that gives it two thoraxes. Bottom images courtesy of the Archives, California Institute of Technology.

# Materials & Methods

1. Research
2. Getting the protein sequence for Mnx1 ( UniProt)
3. Finding homologous proteins (BlastP)
4. Alignment of the homologous proteins (MEGA11)
5. Building Trees (MEGA11)
6. Rerooting the Phylogenetic Tree (FigTree)
7. Pruning the Clade (Python, ete3 library)
8. Calculating conservation scores (Python)
9. Retrieving mutations (Clinical papers & gnomAD)
10. Mapping mutation occurring sites to the aligned sequences (Python)
11. Classification
12. Statistical tests to assess the effect of allele frequencies on mutation type (Python, scipy library)

# Research

- From rarediseases.info.nih.gov website, Currarino Triad Syndrome is chosen.
- The protein which cause Currarino Triad Syndrome when a mutation occurs is identified.
- The protein sequence for Mnx1 is retrieved from UniProt.

# Finding homologous sequences

- BlastP

# Multiple sequence alignment

- MEGA11, Muscle (MUltiple Sequence Comparison by Log-Expectation)

# Phylogenetic tree construction

- MEGA11, Neighbor-Joining method

# Rerooting the phylogenetic tree

- FigTree

# Pruning the clade

- Python3, Ete3 library

# Calculating conservation scores

```python
seqDict = fastareader(filename)
sequences = list(seqDict.values())
aaList = ["A", "R", "N" , "D" , "C" , "Q" , "E" , "G" , "H" , "I" , "L" , "K" , "M" , "F" , "P" , "S" , "T" , "W" ,"Y" ,"V"]

consensus = ""
consensus_aminoacid_score = {}
for pos in range(len(sequences[0])):
    aa_percent_dict = dict.fromkeys(aaList, 0)
    aminoacids_onsamepos = ""
    for seq in sequences:
        aminoacids_onsamepos += seq[pos]
    for aa in aaList:
        count_aa = aminoacids_onsamepos.count(aa)
        aa_percent_dict[aa] = count_aa / len(sequences)
    aa_percent_dict = dict(sorted(aa_percent_dict.items(), key=lambda x: x[1], reverse=True))
    consensus_aa = list(aa_percent_dict.keys())[0]
    max_score = list(aa_percent_dict.values())[0]
    consensus += consensus_aa
    consensus_aminoacid_score[pos] = {consensus_aa:max_score}

conservation_scores_file = open("../conservation/conservation_scores_pruned_s=500.tsv", "w")
conservation_scores_file.write("Position"+"\t"+"Consensus Aminoacid"+"\t" +"Conservation Score" + "\n")
```

# 3 different conservation score calculation sets

- First method:
  - 500 aligned sequences
- Second method:
  - Already aligned sequences that are in the focus clade
- Third method:
  - Realigned sequences that are in the focus clade

# Conservation score histograms for different methods



First method

Second method

Third method

# Conservation scores per position

First
method



Second
method



Third
method

# Retrieving mutations

- Known pathogenic mutations are retrieved from clinical papers about Currarino Syndrome. (17 such mutations)
- The mutations for which the clinical significance is not know are retrieved from gnomAD. (173 such mutations)

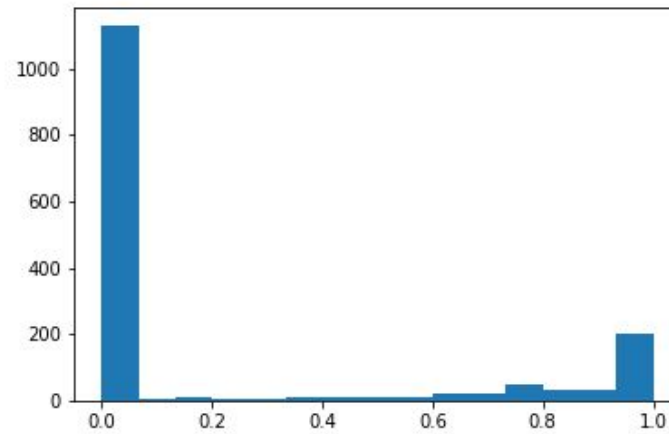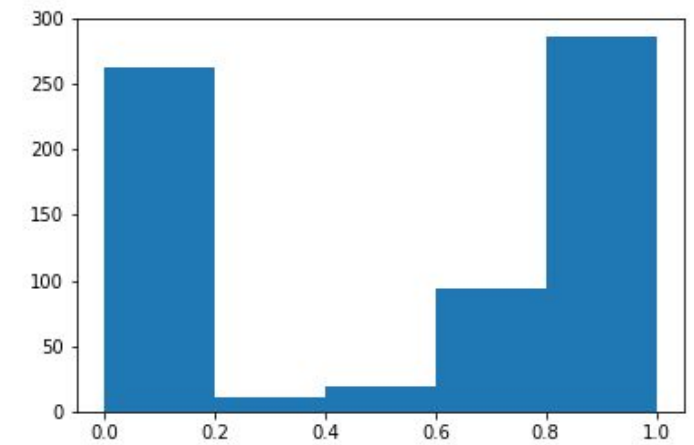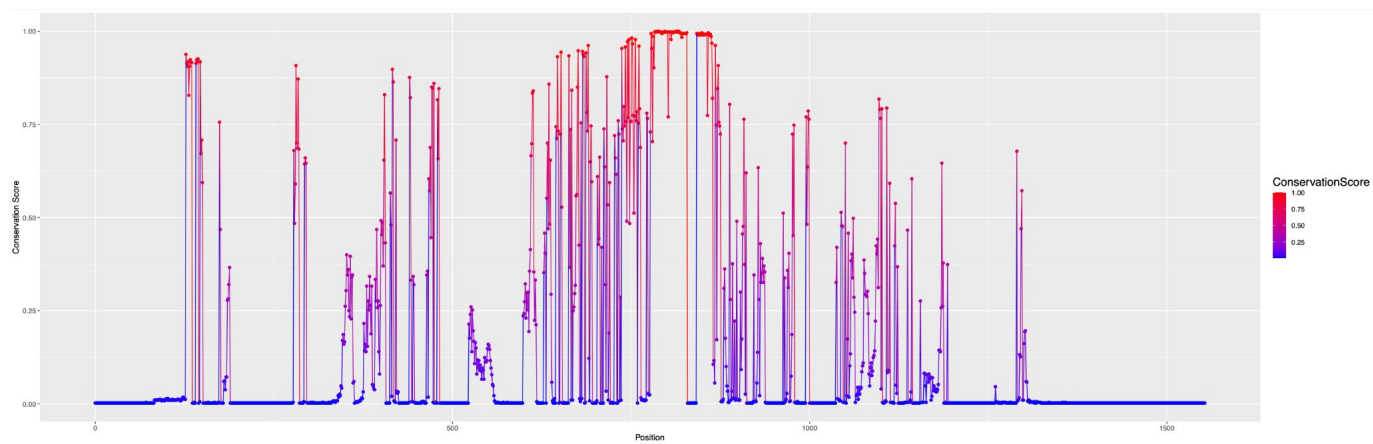| Variant ID | Source | HGVS Consequence | VEP Annotation | LoF Curation | Clinical Significance | Flags | Allele Count | Allele Number | Allele Frequency | H |
|---|---|---|---|---|---|---|---|---|---|---|
| 7-156798251-G-C | E | p.Pro390Arg | missense | | | | 1 | 166492 | 6.01e-6 | |
| 7-156798251-G-A | E G | p.Pro390Leu | missense | | | | 3 | 197408 | 1.52e-5 | |
| 7-156798254-G-A | E | p.Ser389Leu | missense | | | | 11 | 172340 | 6.38e-5 | |
| 7-156798258-C-A | E G | p.Asp388Tyr | missense | | | | 3 | 212878 | 1.41e-5 | |
| 7-156798259-G-C | E | p.Asp387Glu | missense | | | | 1 | 185122 | 5.40e-6 | |
| 7-156798266-G-A | E | p.Ser385Leu | missense | | | | 1 | 202342 | 4.94e-6 | |
| 7-156798267-A-G | E | p.Ser385Pro | missense | | | | 1 | 204042 | 4.90e-6 | |
| 7-156798269-G-A | E | p.Ser384Phe | missense | | | | 1 | 205914 | 4.86e-6 | |
| 7-156798269-G-T | E | p.Ser384Tyr | missense | | | | 3 | 205914 | 1.46e-5 | |
| 7-156798281-G-A | E | p.Ser380Phe | missense | | | | 1 | 218876 | 4.57e-6 | |
| 7-156798282-A-C | E | p.Ser380Ala | missense | | | | 1 | 219394 | 4.56e-6 | |
| 7-156798282-A-G | E | p.Ser380Pro | missense | | | | 1 | 219394 | 4.56e-6 | |
| 7-156798284-G-A | E | p.Ala379Val | missense | | | | 6 | 220344 | 2.72e-5 | |
| 7-156798294-C-T | E | p.Val376Ile | missense | | | | 1 | 226638 | 4.41e-6 | |
| 7-156798294-C-G | E | p.Val376Leu | missense | | | | 1 | 226638 | 4.41e-6 | |
| 7-156798299-G-A | E | p.Ala374Val | missense | | | | 1 | 230024 | 4.35e-6 | |
| 7-156798302-C-T | E | p.Gly373Asp | missense | | | | 1 | 231712 | 4.32e-6 | |
| 7-156798303-C-G | E | p.Gly373Arg | missense | | | | 1 | 232518 | 4.30e-6 | |
| 7-156798309-T-G | E | p.Ser371Arg | missense | | | | 1 | 236242 | 4.23e-6 | |
| 7-156798321-GGT... | E | p.Asp364_Asp366del | inframe deletion | | | | 1 | 238208 | 4.20e-6 | |

# Mapping mutation occurring sites to aligned sequences

- Original sequence:
  - MEKSKNFRIDALLAVDP…

- Aligned sequence:
  - …--------------MEKSKNFRI-----DA…

- For a mutation such as M1A, the position is not 1 in the aligned sequence. Therefore, the actual positions and the aligned positions are mapped.

# Classification

- For each conservation score table, the conservation scores at the positions for the known pathogenic mutations are averaged and set as a threshold for classification.

- ```
  For each unknown-type mutation:

        if the conservation score is above the threshold:

            classify as pathogenic

        else

            classify as neutral
  ```

# Classification results

**First method:**

- Classification threshold: 0.68
- Number of pathogenic mutations: 77
- Number of neutral mutations: 96

**Second method:**

- Classification threshold: 0.89
- Number of pathogenic mutations: 104
- Number of neutral mutations: 69

**Third method:**

- Classification threshold: 0.91
- Number of pathogenic mutations: 102
- Number of neutral mutations: 71

# Statistical tests to assess the effect of allele frequencies on mutation type

- Pathogenic and neutral mutations are separated into two samples, and an independent t-test is conducted on the allele frequencies.

```python
from scipy import stats as st
import pandas as pd

df = pd.read_csv("../classification/classification_pruned_realigned_s=500.csv")

a = df.loc[df['isPathogenic'] == True, 'Allele_frequency'].to_numpy()
b = df.loc[df['isPathogenic'] == False, 'Allele_frequency'].to_numpy()
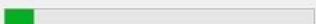
pvalue = st.ttest_ind(a=a, b=b, equal_var = True).pvalue
```

- p-values for:
  - First method → 0.18
  - Second method → 0.08
  - Third method → 0.35

# Issues

- Why we did not use the sequence files of 100 and 250?
- Why we did not use the sequence files of 1000 and 5000?

# Discussion

- Number of known pathogenic mutations should increase in order to classify the unknown mutations accurately.
- As the genetic relevance of the sequences increases in a set of sequences, conservation scores also increase.
- Given the same conservation score calculation method:
  - as classification threshold increases $\rightarrow$ sensitivity decreases.
  - as classification threshold decreases $\rightarrow$ specificity decreases.
- Number of mutations which identified as pathogenic and the classification thresholds are very close for the second and third method.
- With the observation of p-values, there is no statistical significance of the allele frequencies.

# References

- Thompson, Nancy & Gésina, Emilie & Scheinert, Peter & Bucher, Philipp & Grapin-Botton, Anne. (2012). RNA Profiling and Chromatin Immunoprecipitation-Sequencing Reveal that PTF1a Stabilizes Pancreas Progenitor Identity via the Control of MNX1/HLXB9 and a Network of Other Transcription Factors. Molecular and cellular biology. 32. 1189-99. 10.1128/MCB.06318-11.

- Truscott M., Nepveu A. (2005) Homeodomains. In: Encyclopedic Reference of Genomics and Proteomics in Molecular Medicine. Springer, Berlin, Heidelberg . https://doi.org/10.1007/3-540-29623-9_2740

- Currarino syndrome Lynch, SA Review on Currarino syndrome, with data on clinics, and the genes involved. Jean-Loup Huret (Editor-in-Chief) INIST-CNRS (Publisher) 2008

- Mark, M., Rijli, F. & Chambon, P. Homeobox Genes in Embryogenesis and Pathogenesis. *Pediatr Res* 42, 421–429 (1997). https://doi.org/10.1203/00006450-199710000-00001

- Genetic Science Learning Center. (2016, March 1) Homeotic Genes and Body Patterns. Retrieved December 25, 2021, from https://learn.genetics.utah.edu/content/basics/hoxgenes/

- Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, Chen X, Taipale J, Hughes TR, Weirauch MT.(2018) The Human Transcription Factors. Cell. 172(4):650-665. doi: 10.1016/j.cell.2018.01.029. Review.