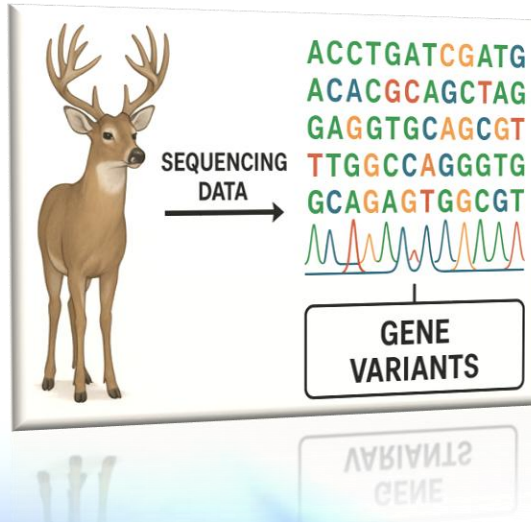


Clustering-Based Identification -  
Gene Variants in Deer Using Sequencing Data

# Introduction



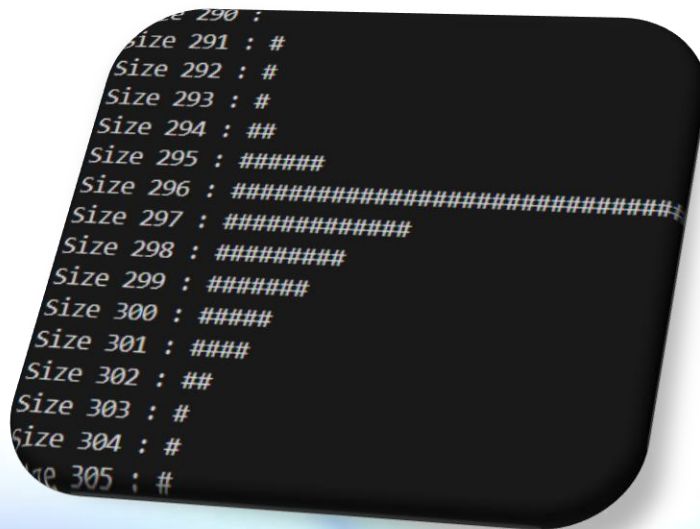
- analysis of gene "Major Histocompatibility Complex"
  - immune system in deer recognizes pathogens
- comparison of appearance in 41 individual deer
  - ~1000 sequence extractions per deer
  - FASTQ format



# Our program

- eliminate corrupted sequences
- correcting measurement errors
- distinguish multiple gene variants
  - $k$  between 2 and 6
  - try to isolate the main variants

# Design & Implementation

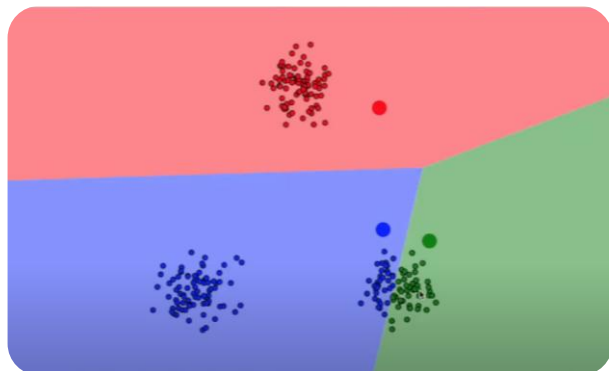


```
Size 290 :  
Size 291 : #  
Size 292 : #  
Size 293 : #  
Size 294 : ##  
Size 295 : #####  
Size 296 : #####  
Size 297 : #####  
Size 298 : #####  
Size 299 : #####  
Size 300 : #####  
Size 301 : ####  
Size 302 : ##  
Size 303 : #  
Size 304 : #  
Size 305 : #
```

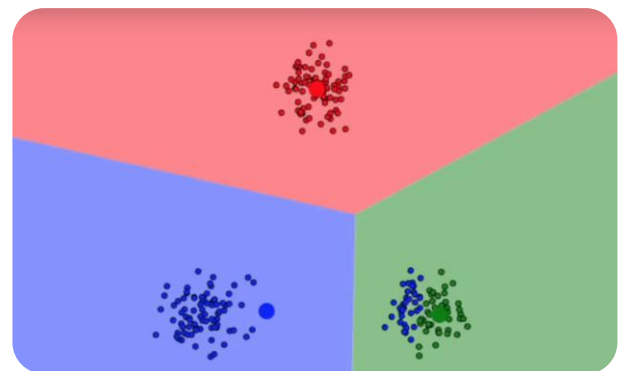
- extract FASTQ file into vector of string
- plot histogram to eliminate irrelevant measurements
- implement k-means clustering



# K-means clustering

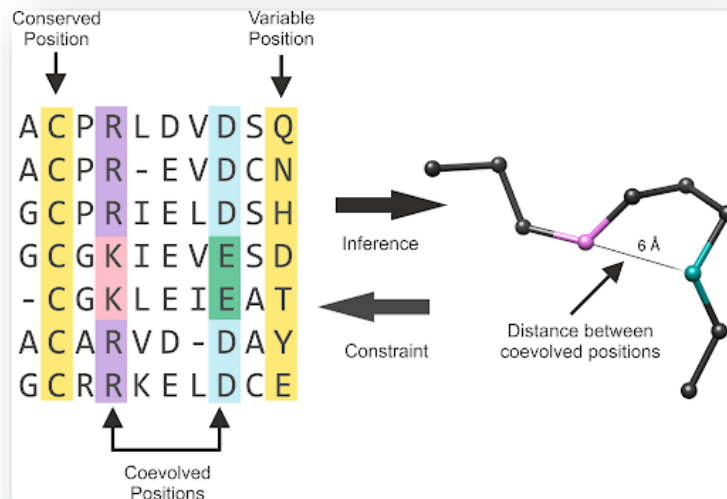


➤ distance metric



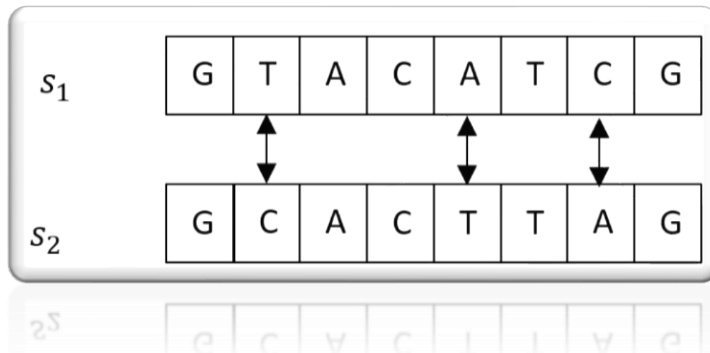
➤ calculating centroids

# SP $\Delta$ A for Alignment & Consensus



- MSA alignment
- graph based alignment
  - generate consensus (= centroid)

# Distance Metrics



➤ Levenshtein distance

➤ Hamming distance



# Evaluation Strategy

- number of 296 bp sequences per file
  - 200-3500 bp
- amount of clusters
  - 3, 6
- sample file name variations for some deer
  - differences not specified





# Observations

- 0 = identical
- digit = difference in bp
- - = too many differences

| Comparison for: J4_S_CE_IonXpress_021 (k3_I296) with Ground Truth J29B |           |           |           |
|--|-----------|-----------|-----------|
|  | Cluster 1 | Cluster 2 | Cluster 3 |
| cluster no 1 : 36  | -         | 0         | -         |
| cluster no 2 : 276   | 0         | -         | -         |
| cluster no 3 : 1   | -         | 1         | -         |
| TOTAL TIME : 22.5554 seconds   |           |           |           |



# Observations

- clustering
  - no difference between amount 3 or 6
  - longer runtime when higher amount
  - 1<sup>st</sup> cluster found, 2<sup>nd</sup> with errors, 3<sup>rd</sup> not found
- matches
  - most with 500+ sequences per file and J30B
  - rare matches in smaller files, none in largest
- similar filenames produce different results



# Conclusion

- SP0A-based program with Hamming distance overall produces realistic result
  - enough input needed,  $k=3$  sufficient, J30B more common, 3<sup>rd</sup> cluster not identifiable (confirmed by supervisor), similar filename  $\neq$  genetic similarity
  - refinements for larger datasets with dynamic  $k$  and better prefiltering
- 