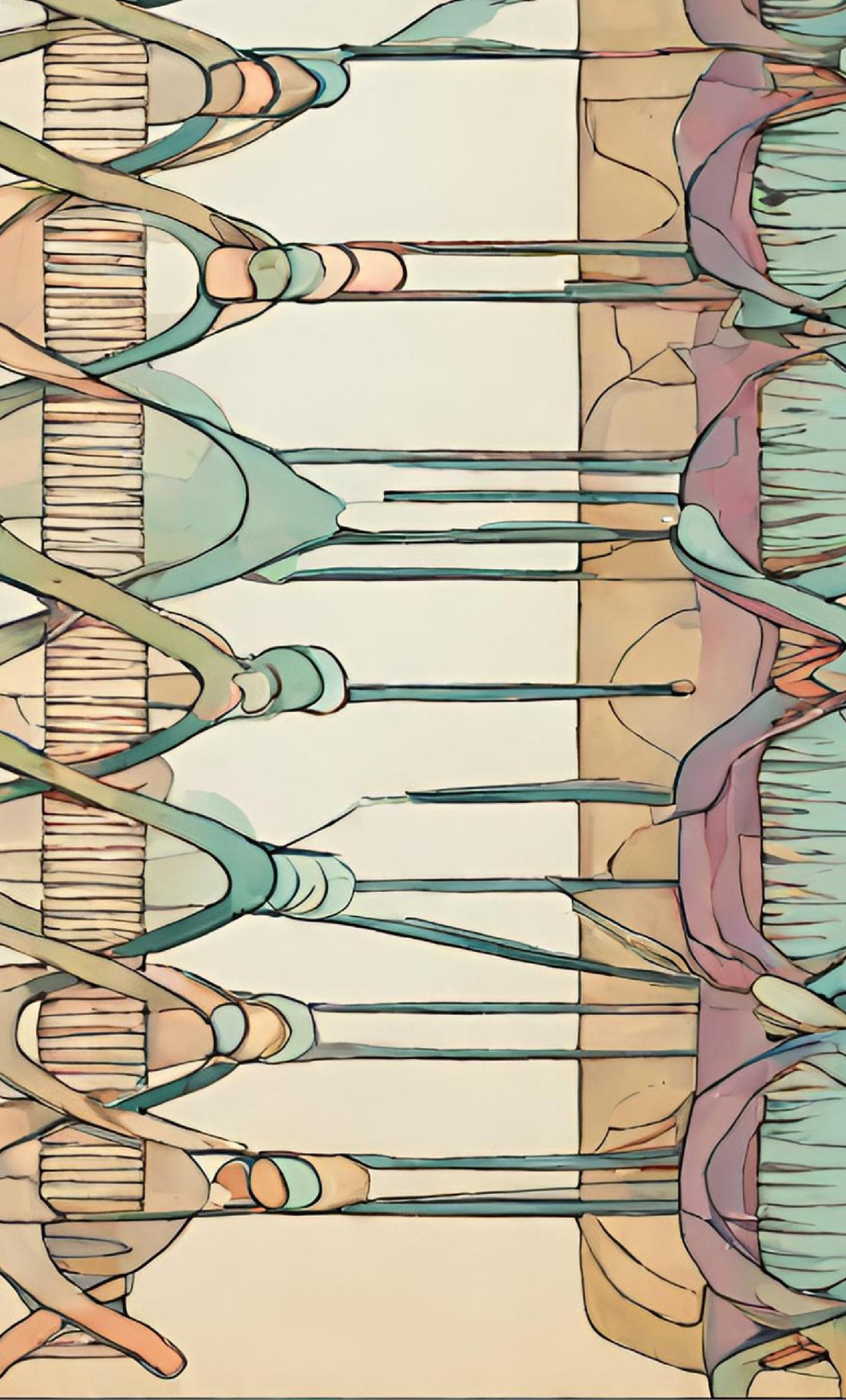


Design Credit Project :-

Topic:- Bioinformatics Tool for
Sequence Alignment and De Bruijn
Graph Analysis

Group members:-

Yashraj Chaturvedi(B22AI059)
Anuj Chincholikar(B22ES018)



CONTENTS

1 .Problem Statement overview

2. Substring Computation

3. De-Brujin, overlapping Graphs and Paths

**4. DNA Alignment and Needleman
Wunsch Algorithm**

5. Conclusion

Introduction

- **Brief overview of the problem statement:-**

- Sequence alignment using the Needleman-Wunsch algorithm and De Bruijn graphs construction for sequence assembly as well as analyzing genetic sequences along with substring computations for large dataset.

- **Objective of the Project:**

- Substring Computation: Extracts substrings.
 - De Bruijn Graph Construction, Visualizations, Paths
 - Needleman-Wunsch Algorithm: Alignment using dynamic programming.
 - Alignment Integration with De Bruijn Graph: Efficiently identifies overlaps.
 - Similarity Percentage Calculation: Assesses alignment quality.

Importance of problem statement and Challenges in Field

Integration of De Bruijn Graph and Sequence Alignment:

- For comparing genetic sequences
- Prediction of functional regions in genomic DNA and Nucleotides
- Identification of similarities and differences among bps
- Aids in Evolutionary ,phylogenetic and mutation studies
- Efficient representation sequence data
- Facilitate genome assembly from short reads

Challenges:

- Coping with diverse genetic sequences and variations.
- Managing high computational demands for processing large datasets.
- Uncertainties and errors in alignments due to sequencing.
- Handling complexities such as repetitive regions and gene duplications in genomes.
- Time-consuming DNA alignment on real genomic datasets.

Various Algorithms for Alignment and Assembly

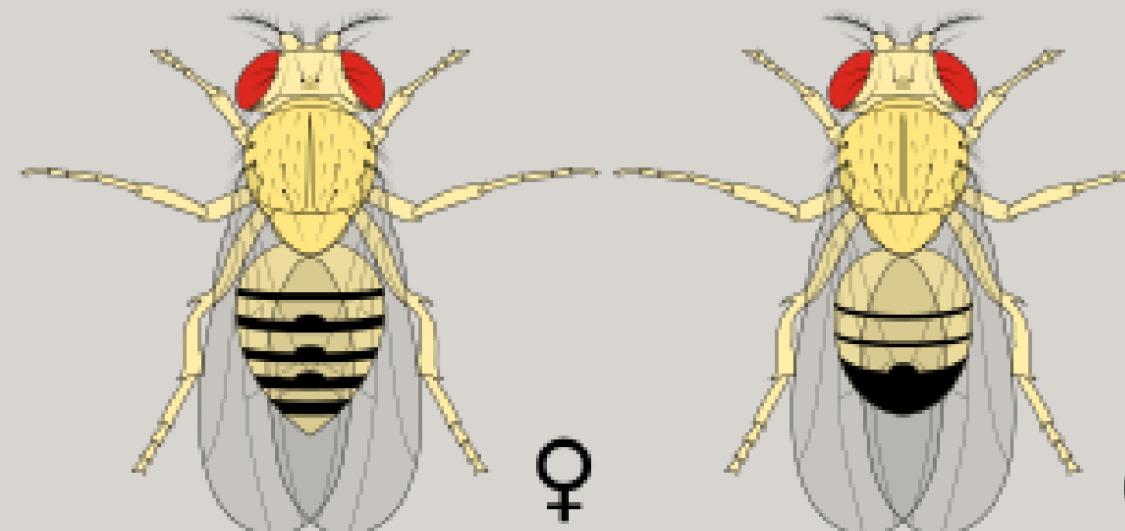
Genome Assembly:

- **Naive Algorithm:**
 - A basic approach concatenating overlapping reads without considering errors or overlaps, leading to inaccurate results, especially in repetitive regions.
- **Greedy Algorithm:**
 - Assembling reads by extending contigs based on immediate overlaps, efficient but may lead to suboptimal assemblies, particularly in complex genomic regions.
- **De-bruijn Graphs**

Genome Alignment:

- **Smith-Waterman:**
 - Local alignment via dynamic programming, finding similar regions.
- **Needleman-Wunsch:**
 - Global alignment via dynamic programming, aligning entire sequences.
- **BLAST**
 - Heuristic local alignment tool for comparing sequences.
- **Hidden Markov Models (HMMs)**
 - Statistical models for aligning sequences with structural similarities.

Input File of Drosophila Melanogaster DNA



ref. https://ftp.ensembl.org/pub/release-111/fasta/drosophila_melanogaster/d

- **Insects (*Drosophila melanogaster*)**: Approximately 180 million base pairs (Mb), equivalent to 180,000,000 bp.
 - **Humans (*Homo sapiens*)**: Approximately 3,200 million base pairs (Mb), equivalent to 3,200,000,000 bp.
 - **Assembly Types**: Primary, Nonchromosomal, Toplevel.
 - **Sequence Characteristics**: dna, dna_rm, dna_sm.
 - **Variability in Size and Timestamps**: Reflects data content, annotation levels, and currency.

← → C 🔍 https://ftp.ensembl.org/pub/release-111/fasta/drosophila_melanogaster/dna/

Index of /pub/release-111/fasta/drosophila_melanogaster/dna

	Name	Last modified	Size	Description
	Parent Directory			-
	CHECKSUMS	2023-10-19 16:11	2.3K	
	Drosophila melanogaster.BDGP6.46.dna.nonchromosomal.fa.gz	2023-10-04 12:42	1.5M	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.2L.fa.gz	2023-10-04 12:42	6.9M	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.2R.fa.gz	2023-10-04 12:42	7.3M	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.3L.fa.gz	2023-10-04 12:42	8.2M	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.3R.fa.gz	2023-10-04 12:42	9.3M	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.4.fa.gz	2023-10-04 12:42	395K	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.X.fa.gz	2023-10-04 12:42	6.8M	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.Y.fa.gz	2023-10-04 12:42	889K	
	Drosophila melanogaster.BDGP6.46.dna.primary_assembly.mitochondrion_genome.fa.gz	2023-10-04 12:41	5.2K	

Substring computation of large strings on *Drosophila melanogaster* DNA

- Approach Used:

- Read genetic sequences from a FASTA file.
- Iterate through the DNA sequence in steps equal to the desired substring length
- Write computed substrings to a CSV file for further analysis.
- Later this will be utilized in de-bruijn graph construction

Output of computed substring in csv file

```
→ Enter the genetic sequence: ATGCTAGCTGATGCTGATGCTAGCGTAGCGTACTATCATCTACGTGCGAGGACTGACTGCTGATGCTGAGCGAGCGCGCGCGCTTCCTCTCGCGCGC  
Enter the substring length: 5  
Substrings of length 5 in the genetic sequence:  
ATGCT  
TGCTA  
GCTAG  
CTAGC  
TAGCT  
AGCTG  
GCTGA  
CTGAT  
TGATG  
GATGC  
.....
```

Output for small protein sequence

1 to 10 of 19510 entries

Substring

AATGAATT
ATGAATTG
TGAATTGC
GAATTGCC
AATTGCCT
ATTGCCTG
TTGCCTGA
TGCCTGAT
GCCTGATA
CCTGATAA

Show per page

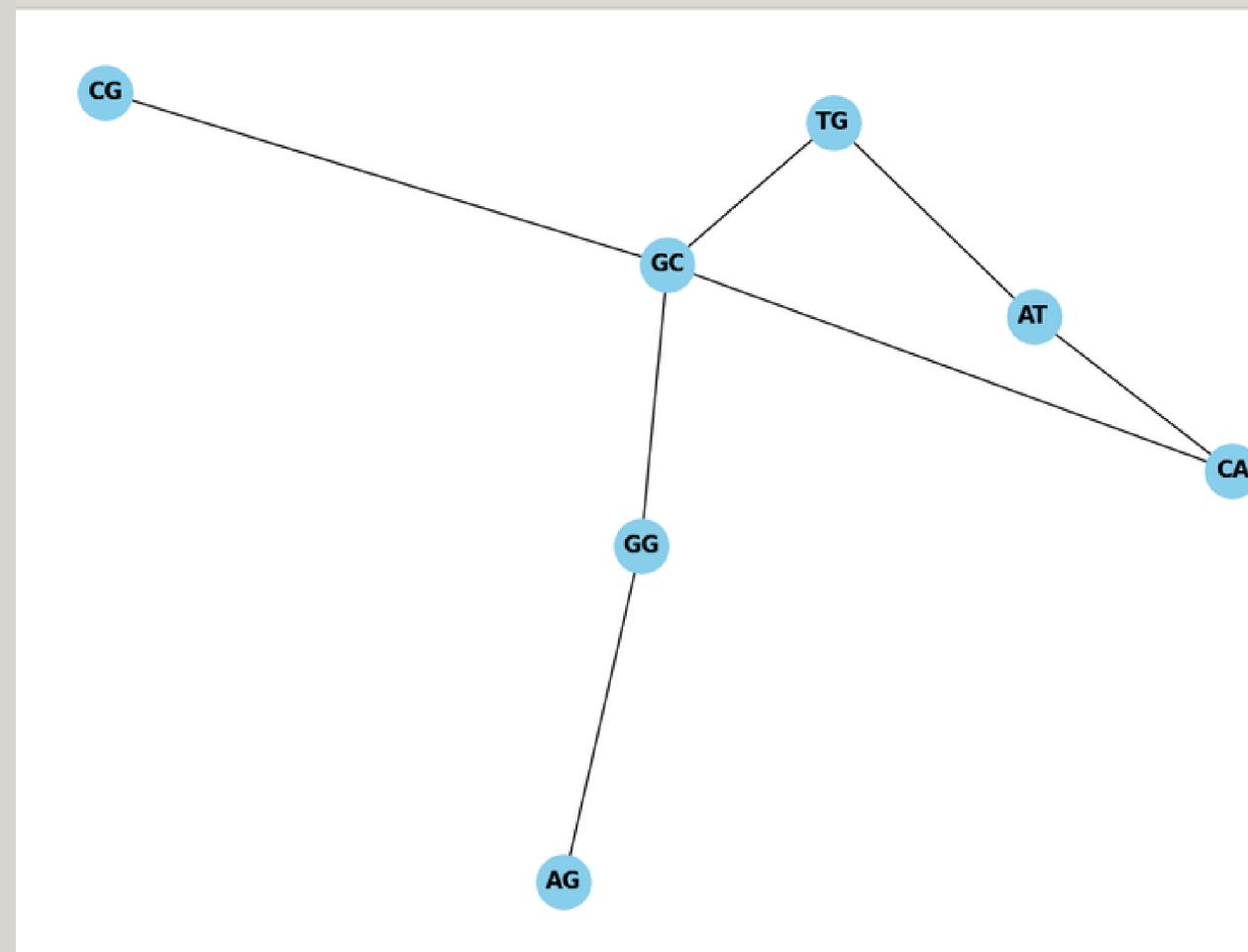
1 2 10 100 1000 1900 1950 1999

There are a total of 19510 substring computed of k-mer length 8 for drosophila

Definitions and Importance: Overlapping , De-bruijn Graphs , Paths

- Overlapping graphs show how short DNA reads overlap, forming a network.
- De Bruijn graphs (a special case of overlapping graph) represent overlaps between substrings of fixed length in sequences.
 - They aid in genome assembly and reconstructing DNA sequences from short reads.
- Eulerian paths traverse each edge of a graph exactly once
 - They are used in genome assembly algorithms like Eulerian assembly as well as sequence analysis
- Hamiltonian paths visit each vertex of a graph exactly once.
 - They are applied in various optimizations and finding sub-sequence of particular length but less common in genomics due to complexity.

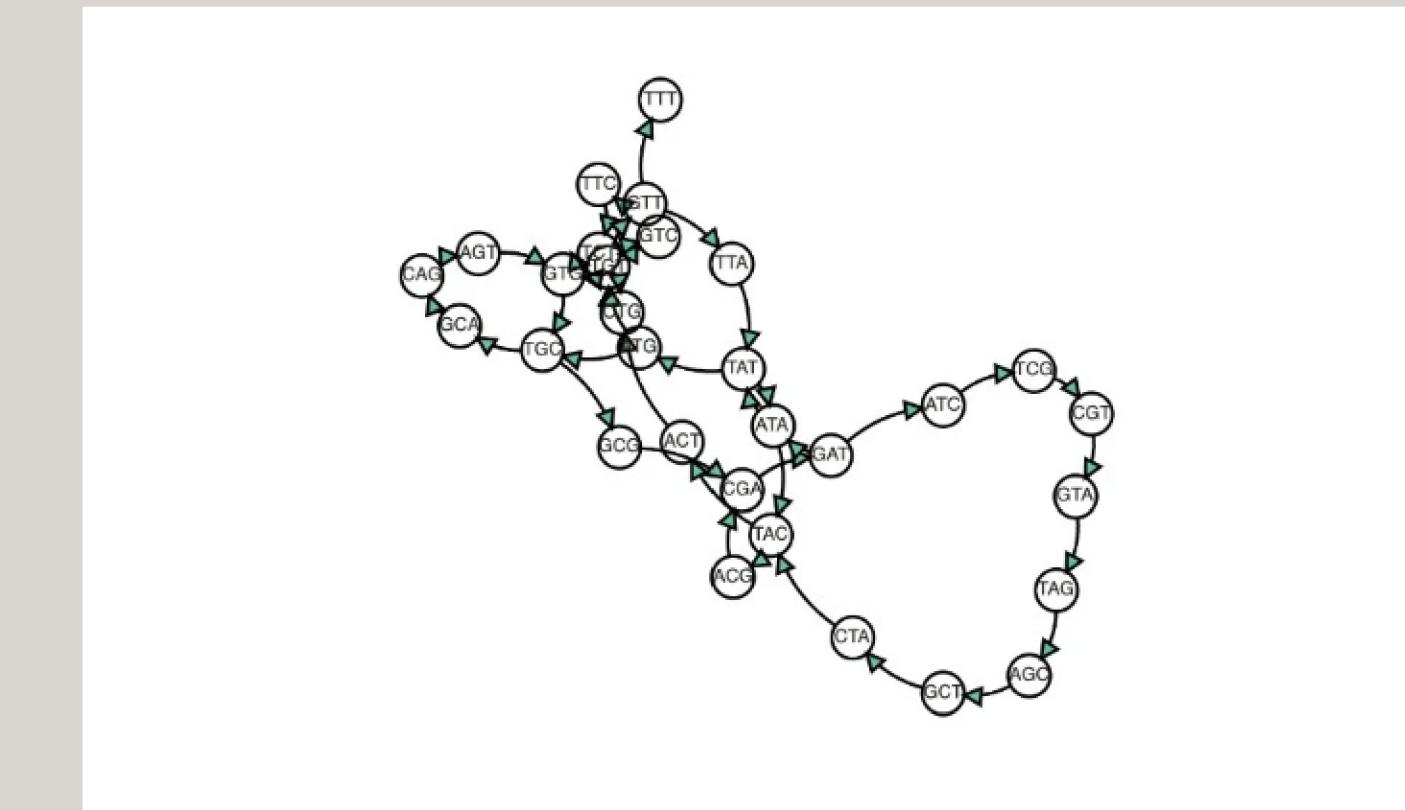
Visualizations of De-bruijn graphs obtained from small Reads



```
# Example reads
reads = ["ATGCG", "GCATG", "CATGC", "AGGCA", "GGCAT"]

# Length of k-mer
k = 3
```

```
ATCGATCGATCGATCGAATGTTCGATCGATCGATTACT
```



Output of Obtained Debrujin Edges,nodes, Eulerian path and hamiltonian paths for protein sequence.fasta file:-

sequence.fasta ×

```
1 >AAH22532.1 Prion protein [Homo sapiens]
2 MANLGCWMLVLFVATWSDLGLCKRPKPGGWNTGGSRYPGQGSPGGNRYPQPPQGGGGWGQPHGGGWGQPHG
3 GGWGQPHGGGWGQPHGGGWGQGGGTHSQWNKPSKPKNMKHMAGAAAAGAVVGGLGGYVLGSAMSRPIIH
4 FGSDYEDRYYRENMRYPNQVYYRPMDEYSNQNNFVHDCVNITIKQHTVTTTKGENFTETDVKMMERVV
5 EQMCITQYERESQAYYKRGSSMVLFSSPPVILLISFLIFLIVG
6
7
```

De Bruijn Edges:

```
[('L', 'F'), ('C', 'V'), ('S', 'S'), ('D', 'R'), ('G', 'L'), ('I', 'H'), ('T', 'G'), ('H', 'R'), ('F', 'T'), ('N', 'K'),
```

De Bruijn Nodes:

```
['C', 'R', 'Y', 'G', 'T', 'V', 'W', 'I', 'L', 'F', 'A', 'N', 'Q', 'E', 'H', 'M', 'S', 'P', 'K', 'D']
```

Eulerian Path:

```
['C', 'W', 'M', 'A', 'N', 'L', 'G', 'C', 'K', 'R', 'P', 'K', 'R', 'P', 'G', 'L', 'V', 'L', 'F', 'V', 'A', 'T', 'W', 'S',
```

Overlapping Graph:

```
( L , F ) ( L , C ) ( L , V ) ( L , G ) ( L , L ) ( L , I ) ( C , V ) ( C , I ) ( C , K ) ( C , W ) ( S , S ) ( S , F ) .
```

Hamiltonian Path:

```
['C', 'V', 'G', 'L', 'F', 'T', 'K', 'M', 'S', 'A', 'N', 'R', 'Y', 'P', 'I', 'H', 'D', 'E', 'Q', 'W']
```



DNA

ALIGNMENT

Needleman Wunsch Algorithm for global alignments on protein sequences

- **Needleman-Wunsch algorithm:** The Needleman_wunsch takes two input sequences seq1 and seq2, along with optional parameters for match, mismatch, and gap scores.
- **Initialization:** Creates score and traceback matrices.
- **Scoring Scheme:** Assigns scores for matches, mismatches, and gaps. We have taken match_score=1, mismatch_score=-1, gap_score=-2
- **Filling the Matrix:** Calculates scores for each cell using dynamic programming.
- **Traceback:** Reconstructs optimal alignment by following traceback path.
- **Alignment and Scoring:** Outputs aligned sequences and calculates alignment score.
- **Applications:** Used in DNA alignment to identify evolutionary relationships, functional elements, and mutation effects.
- **Implementation:** Reads sequences from a file, aligns pairs of sequences, calculates similarity percentage, constructs De Bruijn graphs, finds common paths, and visualizes graphs.
- **Importance:** Fundamental tool in bioinformatics for sequence analysis and similarity detection.

SEND
. ||
-AND
Score=2

Input File of 11 Protein Sequences of different organisms

ref. <https://www.ncbi.nlm.nih.gov/protein/>

Seq1 [organism=Carpodacus mexicanus] [clone=6b] actin (act) mRNA, partial
CTTTATCTAATCTTGGAGCATGAGCTGGCATAGTTGGAACCGCCCTCAGCCTCCTCATCCGTGCAGA
TTGGACAACCTGGAACTCTCTAGGAGACGACCAAATTACAATGTAATCGCACTGCCACGCCTTC
AATAATTTCTTATAGTAATACCAATCATGATCGGTGGTTGGAAACTGACTAGTCCCACTCATAA
GGGCCCGGCGACATAGCATTCCCCGTATAAACAAACATAAGCTCTGACTACTCCCCATCATTCT
TACTCTAGCATCCTCACAGTAGAAGCTGGAGCAGGAACAGGTGAACAGTATATCCCCCTCGCT
TAACCTAGCCATGCCGGTGCTTCAGTAGACCTAGCCATCTCTCCACTTAGCAGGTGTTCTC
ATCCTAGGTGCTATTAACTTATTACAACC GCCATCAACATAAAACCCCCAACCTCTCCAATACCA
CCCCCTATT CGTATGATCAGTC TTATTACCGCCGTC TTCTCCTACTCTCTCCAGTCCTCGCT
TGGCATTACTATACTACTAACAGACCGAAACCTAAACACTACGTTCTTGACCCAGCTGGAGGAGGAGA
CCAGTCCTGTACCAACACCTCTGATTCTCGGCCATCCAGAAGTCTATATCCTCATTTC
Seq2 [organism=uncultured bacillus sp.] [isolate=A2] corticotropin (CRH)
GTAGGTACCGCCCTAAGNCTCTAATCCGAGCAGAACTANGCCAACCGGAGCCCTCTGGGAGACGA
AAATCTACAACGTAGTCGTTACGGCCACGCCTCGTAATAATCTTTCTAGTAATGCCAATCATAC
CGGAGGATTGGAACTGACTAGTTCTCTAATGATTGGGGCCCCAGACATAGCATTCCCTCGAATAAA
AACATAAGCTTTGACTACTACCACCATCATTCTACTCCTAATAGCCTCTAACAGTAGAAGCAGGA
CCGGAACCGGATGAACCGTGTACCCACCACTAGCTGGAAACCTGGCCACGCCGGAGCCTAGTAGAC
AGCTATCTCTCCCTACACCTAGCAGGTATCTCATCCATCCTGGGGCAATTAACTCATTACAACAG
ATCAACATAAAACCACCGCCCTCTCACAATACCAACCCACTATTGTGTATCCGTCTTAATTAC
CCGTACTACTCCTACTATCTCTCCAGTACTAGCCGCCGTATCACCAGTCTACTCACAGACCGCAAC
AACACCAACCTTCTTGACCCAGCAGGAGGAGACCCAGTACTATACCAGCACCTATTCTGATTCT
GGACACCCAGAAGTCTACATCCTAATTCTC
Seq3 [organism=Phalaenopsis equestris var. leucaspis]
CTATACCTAATTTCGCGCATGAGCCGGAATGGTGGGTACCGCTCTAACGCTCCTCATCGAGCAGA
TAGGCCAACCGGAGCCCTCTGGGAGACGACCAAGTCTACAACGTGGTTGTACGGCCATGCCTTC
AATAATCTCTTATAGTTATGCCGATTATAATGGAGGATTGGAAACTGACTAGTCCCCCTAATAA
GGAGCCCCAGACATAGCATTCCCGGAATAAACAAACATAAGCTCTGACTACTCCACCATCTTCTC
CTCTCTTAGCATCCTCACAGTGGAAGCAGGCGTAGGTACAGGCTGAACAGTGTATCCCCACTAGCT
CAACCTAGCTCATGCCGGGCTCAGTCACCTCGCAATCTCTCCTACACCTAGCTGGTATTCCTC
ATCCTCGGAGCAATTAACTCATTACAACAGCAATTAACTGAAACCTCTGCCCTCTACAATACCA
CCCCACTATTGTGTATCAGTGTAAATTACTGCAGTCCTCCTTCTCCAGTTCTAGCT
AGGAATCACAATGCTCCTCACAGACCGCAACCTAACACCAATTCTGACCCCTGCCGGAGGAGGAGA
CCCGTCTATATCAACATCTCTGATTCTCGGCCACCCAGAAGTCTACATCCTAATTCTC
Seq4 [organism=uncultured archaeon]
ATGAGCTGGAATAGTAGGTACCGCCCTAACGCTCCTAACCTGAGCAGAGCTAGGCCAACCGGAGCC
CTGGGAGACGACCAAATCTACAACGTAGTCGNACGGCCATGCTTTGTAATAATCTCTCATAGCT
TGCCAACTCATATCGGAGGGTTGGAAACTGACTGGCCCCCTAACATAATTGGAGCTCCAGACATAGCT

Score and Traceback matrices

- **Score Matrix:**

- Represents scores of aligning subsequences.
- Initialized with scores based on a scoring scheme.
- Filled iteratively based on adjacent cell scores and character similarities.

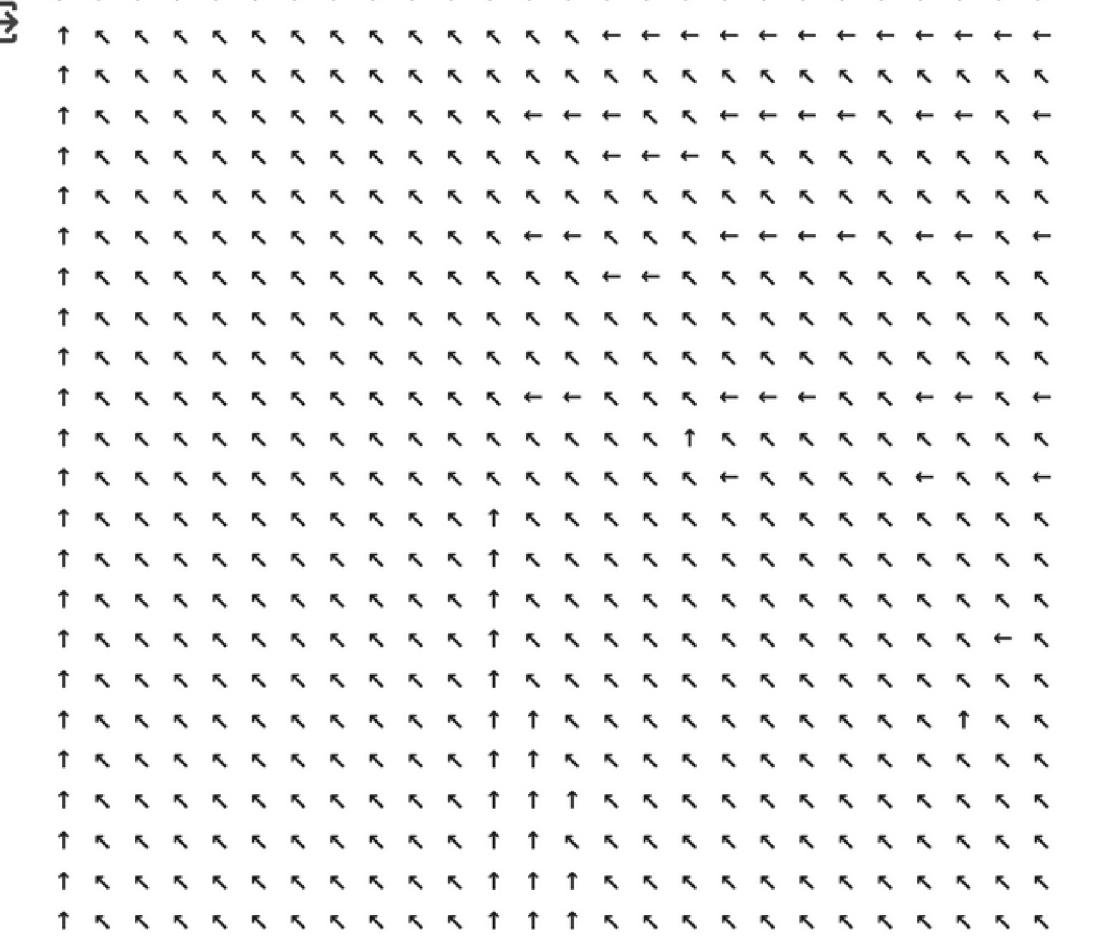
- **Trace-back Matrix:**

- Guides backtracking from bottom-right to top-left.
- Each cell indicates direction: diagonal (match/mismatch), up (gap in first sequence), or left (gap in second sequence).
- Used to reconstruct the optimal alignment.

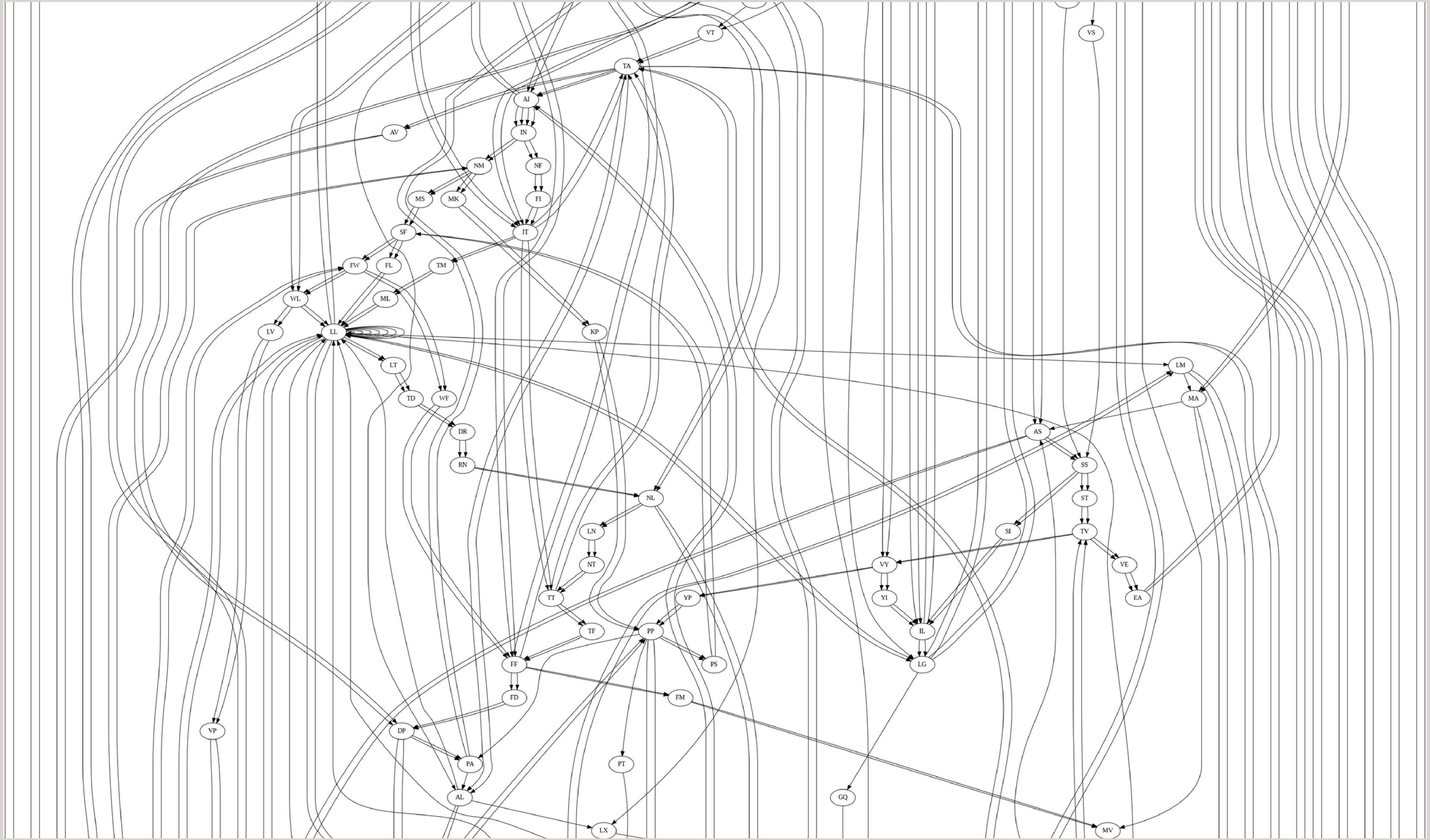
Score Matrix:

0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22	-24	-26	-28	-30	-32	-34	-36	-38	-40	-42	-44	-46	-48	-50
-2	-1	-3	-5	-7	-9	-11	-13	-15	-17	-19	-21	-23	-25	-27	-29	-31	-33	-35	-37	-39	-41	-43	-45	-47	
-4	-3	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22	-24	-26	-28	-30	-32	-34	-36	-38	-40	-42	-44	-46	-48
-6	-5	-4	-3	-5	-7	-9	-11	-13	-15	-17	-19	-21	-23	-25	-27	-29	-31	-33	-35	-37	-39	-41	-43		
-8	-7	-6	-5	-4	-6	-8	-10	-12	-14	-16	-18	-18	-20	-22	-24	-26	-28	-30	-32	-34	-36	-38	-40	-42	
-10	-9	-8	-7	-6	-5	-7	-9	-11	-13	-15	-17	-19	-19	-21	-23	-25	-27	-29	-31	-33	-35	-37	-39	-41	
-12	-11	-10	-9	-8	-7	-6	-5	-7	-9	-11	-13	-15	-15	-17	-19	-21	-23	-25	-27	-29	-31	-33	-35	-37	
-14	-13	-12	-11	-10	-9	-8	-7	-6	-5	-7	-9	-11	-13	-15	-15	-17	-19	-21	-23	-25	-27	-29	-31	-33	
-16	-15	-14	-13	-12	-11	-10	-9	-8	-7	-6	-5	-7	-9	-11	-13	-15	-16	-18	-20	-22	-24	-26	-28	-30	
-18	-17	-16	-15	-14	-13	-12	-11	-10	-9	-8	-7	-6	-5	-7	-9	-11	-13	-15	-17	-17	-19	-21	-23	-25	
-20	-19	-18	-17	-16	-15	-14	-13	-12	-11	-10	-10	-12	-14	-16	-18	-20	-22	-24	-26	-28	-30				
-22	-21	-20	-19	-18	-17	-16	-15	-14	-13	-12	-11	-11	-13	-15	-17	-18	-19	-21	-23	-25	-27	-29			
-24	-23	-22	-21	-20	-19	-18	-17	-16	-15	-14	-11	-12	-12	-14	-16	-18	-20	-22	-24	-26	-28	-30			
-26	-25	-24	-23	-22	-21	-20	-19	-18	-17	-16	-13	-12	-11	-13	-15	-17	-17	-19	-19	-21	-23	-25	-27		
-28	-27	-26	-25	-24	-23	-22	-21	-20	-19	-18	-15	-14	-13	-12	-14	-16	-18	-18	-20	-22	-24	-26			
-30	-29	-28	-27	-26	-25	-24	-23	-22	-21	-20	-17	-16	-15	-14	-13	-15	-17	-19	-19	-21	-23	-25			
-32	-31	-30	-29	-28	-27	-26	-25	-24	-23	-22	-19	-16	-17	-14	-15	-14	-16	-18	-20	-22	-22	-24			
-34	-33	-32	-31	-30	-29	-28	-27	-26	-25	-24	-21	-18	-17	-16	-15	-16	-15	-17	-19	-21	-21	-23			
-36	-35	-34	-33	-32	-31	-30	-29	-28	-27	-26	-23	-20	-17	-18	-17	-16	-17	-18	-16	-18	-20	-22	-23		
-38	-37	-36	-35	-34	-33	-32	-31	-30	-29	-28	-25	-22	-19	-18	-19	-18	-17	-18	-17	-19	-21	-23	-24		
-40	-39	-38	-37	-36	-35	-34	-33	-32	-31	-30	-27	-24	-21	-20	-19	-18	-19	-18	-19	-20	-22	-24			
-42	-41	-40	-39	-38	-37	-36	-35	-34	-33	-32	-29	-26	-23	-22	-21	-20	-21	-20	-19	-20	-21	-23			
-44	-43	-42	-41	-40	-39	-38	-37	-36	-35	-34	-31	-28	-25	-24	-23	-22	-21	-22	-21	-20	-21	-22			
-46	-45	-44	-43	-42	-41	-40	-39	-38	-37	-36	-33	-30	-27	-26	-25	-24	-23	-22	-23	-22	-21	-23			
-48	-47	-46	-45	-44	-43	-42	-41	-40	-39	-38	-35	-32	-29	-28	-27	-26	-25	-24	-23	-24	-23	-22			
-50	-49	-48	-47	-46	-45	-44	-43	-42	-41	-40	-37	-34	-31	-30	-29	-28	-27	-26	-25	-24	-25	-24			
-52	-51	-50	-49	-48	-47	-46	-45	-44	-43	-42	-39	-36	-33	-32	-31	-30	-29	-28	-27	-26	-25	-24			
-54	-53	-52	-51	-50	-49	-48	-47	-46	-45	-44	-41	-38	-35	-34	-33	-32	-31	-30	-29	-28	-27	-26			
-56	-55	-54	-53	-52	-51	-50	-49	-48	-47	-46	-43	-40	-37	-36	-35	-34	-33	-32	-31	-30	-29	-28			
-58	-57	-56	-55	-54	-53	-52	-51	-50	-49	-48	-45	-42	-39	-38	-37	-36	-35	-34	-33	-32	-31	-30			
-60	-59	-58	-57	-56	-55	-54	-53	-52	-51	-50	-47	-44	-41	-40	-39	-38	-37	-36	-35	-34	-33	-32			
-62	-61	-60	-59	-58	-57	-56	-55	-54	-53	-52	-49	-46	-43	-42	-41	-40	-39	-38	-37	-36	-35	-34			
-64	-63	-62	-61	-60	-59	-58	-57	-56	-55	-54	-51	-48	-45	-44	-43	-42	-41	-40	-39	-38	-37	-36			
-66	-65	-64	-63	-62	-61	-60	-59	-58	-57	-56	-53	-50	-47	-46	-45	-44	-43	-42	-41	-40	-39	-38			
-68	-67	-66	-65	-64	-63	-62	-61	-60	-59	-58	-55	-52	-49	-48	-47	-46	-45	-44	-43	-42	-41	-40			
-70	-69	-68	-67	-66	-65	-64	-63	-62	-61	-60	-57	-54	-51	-50	-49	-48	-47	-46	-45	-44	-43	-42			
-72	-71	-70	-69	-68	-67	-66	-65	-64	-63	-62	-59	-56	-53	-52	-51	-50	-49	-48	-47	-46	-45	-44			

- ▶ Traceback Matrix:



Visualization of Debruijn graphs on protein dataset using graphviz library



Output generated from Needleman Wunsch algorithm when we applied it on Protein.txt file

Alignment and De Bruijn graph for sequences 1 and 2:

Alignment for sequences 1 and 2:

LYLIFGAWAGMVGTLALSLLIRAEGLQPGTLLGDDQIYNVIVTAHAFVMIFFMVMPIMIGGFGNWLVPPLMIGAPDMAFPRMNNMSFWLLPPSFLLLASSTVEAGAGTGWTVYPPLAGNLAF
|||||.|||||||.|||.|||||||||.|||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
-----VGTALXLLIRAELEXQPGALLGDDQIYNVVVTAHAFVMIFFMVMPIMIGGFGNWLVPPLMIGAPDMAFPRMNNMSFWLLPPSFLLLASSTVEAGAGTGWTVYPPLAGNLAF

Score=184

similarity percentage from our model: 92.20779220779221

Similarity Percentage from library: 92.20779220779221

Common Paths: [('VGT', 'GTA'), ('GTA', 'TAL'), ('LSL', 'SLP'), ('LLI', 'LIR'), ('LIR', 'IRA'), ('IRA', 'RAE'), ('RAE', 'A')]

This is the image for alignment ,similarity percentage and common paths between sequence 1 and sequence 2 similarly there are total of 7 sequences so 7c2 pairs of comparisions.

De Bruijn graph:

LY → YL, YQ, YQ

YL → LI

LI → IF, IR, IT, IL, IR, IT, IL

IF → FG, FF, FS, FF, FS

FG → GA, GN, GH, GN, GH

GA → AW, AP, AG, AS, AI, AL, AP, AG, AS, AI

AW → WA

WA → AG

AG → GM, GA, GT, GN, GA, GV, GI, GG, GA, GT, GN, GA, GI, GI, GG

GM → MV

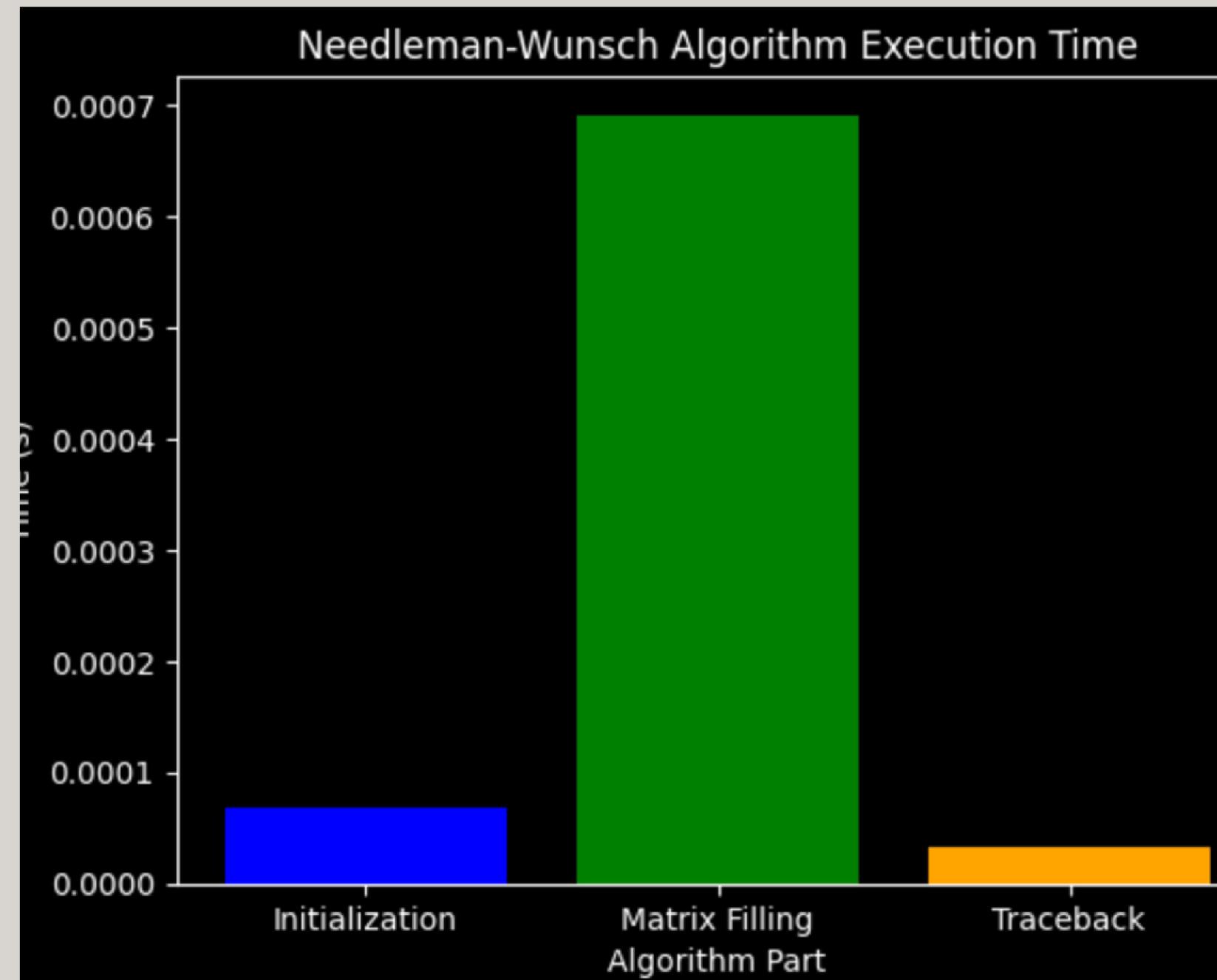
MV → VG, VM, VM

VG → GT, GT

We also computed cumulative error percentage in our model with respect to library function “pairwise2” in Bio , which came out to be 0.04409549123602203% in this case.

% error in our model 0.04409549123602203%

Time taken by each step of Needleman wunsch algorithm on protein data



It is clearly from graph that Matrix filling step in needleman wunsch algorithm takes the maximum time to compute relative to initialization and traceback.

We can optimize this by parallel processing.

Conclusion

In this project we started off with base idea of how to compute substrings on large genetic sequences ,then we moved over to concept of de-bruijin graphs , later on which we connected to overlapping graphs , we calculated eulerian and hamiltonian paths for these graphs and then started with concept of DNA Alignment wehere we did Needlenman Wunsch Algorithm which is a global alignment technique ,we saw similarity percentage between any 2 DNA sequences , made de-Brujin graph of those sequences to find common paths between both and lastly we saw time analysis of our algorithm on a Protein real time Dataset



Thank you!