

# **Research Project on Sequence Classification via Graph**

## **Final Report**

CSC 6850 Machine Learning

SPRING 2024

Alenka Tang, Omar Madjitolv

### **❖ Introduction**

- **Project Type:** research project in sequence classification using graph models
- **Background and Motivation**

Classification is not a new topic in bioinformatics. In the analysis of gene expression data, genes obtained from microarray data are clustered and genes in the same cluster are considered to trigger the same function. Sequence classification is an important task in bioinformatics to produce stable, quick, and accurate results.

### **➤ Problem Definition**

This project proposed a method using graphical models for sequence classification where the graph model will provide the features for training classifiers. We will take existing models such as the de-Bruijn graph[1], and OverlapGraphs[5], then use Node2vec [7], and Graph2vec embedding models to create input vectors for each sequence data which will be used to train Support Vector Machine (SVM) and Neural Network (NN) classifiers and test them with different parameters to find the combination with optimal results for sequence classification.

### **❖ Methodology**

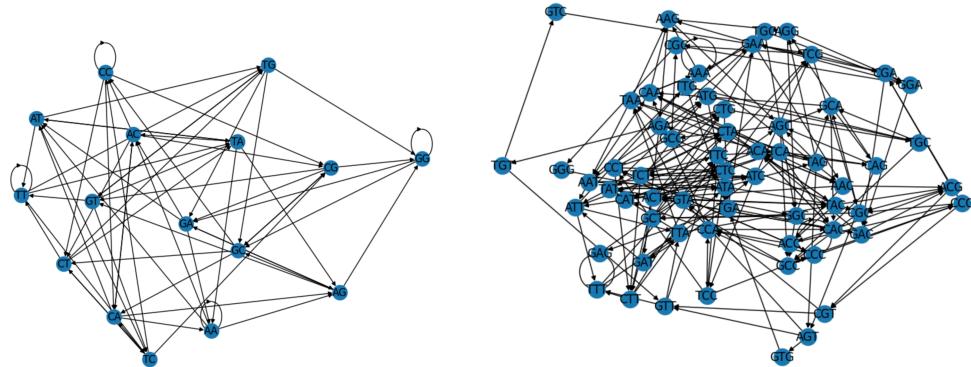
#### **➤ Datasets**

This project will use the Human DNA (HD) sequence database consisting of 4380 DNA sequences with each DNA sequence corresponding to a specific gene family

which can be used as labels for training and testing the model, with a total of seven labels. Due to time constraints, we only used the first 1000 DNA sequences to test the model.

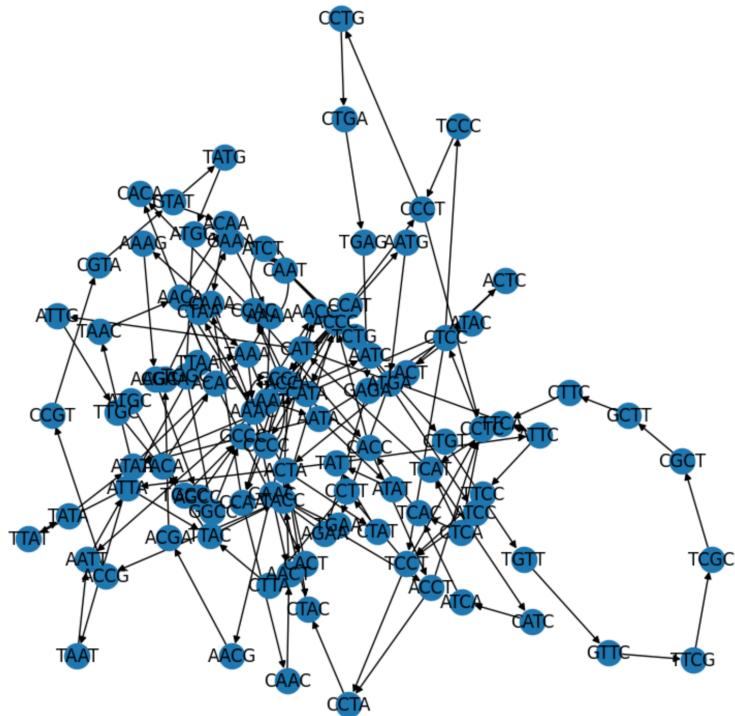
## ➤ De-Brujin Graph Construction

This project will use the de-Bruijn graph model, as well as Overlap Graphs to convert the sequence data to graphs with each sequence data entry being one graph. For constructing the de-Bruijn graph for each data entry, we will try different k-mer values and break the sequence into overlapping k-mers. The graph will be constructed by creating nodes for each k-mer and connecting overlapping k-mers with direct edges. Lastly, we will simplify the graph by merging nodes with common overlaps with each sequence data in the dataset being an individual graph. The complexity of the graph will first increase and then decrease as we increase the value of kmers.



kmers = 3

kmers = 4



Kmer = 5

### ➤ Node Embedding Models

After constructing the graphs from the sequence data, we use Node2vec and embedding models to create the input vectors from each graph.

The Node2vec embedding model uses random walks to explore the graph and outputs a matrix with information for each node. Increasing the kmer value when constructing the de-bruijn graph causing the increase in several nodes and the complexity of the graph would also increase the processing time for the Node2vec embedding models. When we increased the kmer value to 10, the Node2vec embedding model took 3 hours to complete 50 graphs which limited the amount of data we could use as the sample in a short period of time.



The output of the Node2vec model is a matrix where each node has its own vector. We then take the sum of the feature values as well as the average of each node and create the input feature vector with dimension = 64.

```
#embedding example
g_ex = graph_list[0]
node2vec_ex = Node2Vec(g_ex, dimensions=64, walk_length=30, num_walks=200, workers=4)
model_ex = node2vec_ex.fit(window=10, min_count=1, batch_words=4)
embeddings_ex = {node: model_ex.wv[node] for node in g_ex.nodes()}

embeddings_ex
```

Computing transition probabilities: 100% 16/16 [00:00<00:00, 292.94it/s]

```
'TG': array([ 0.03870851, -0.15258811,  0.3023589 ,  0.00167002, -0.03254861,
   0.84632767, -0.02603316, -0.1363045 , -0.15608446, -0.16676189,
   0.06534941, -0.05790641, -0.06549331,  0.0604237 , -0.12670293,
  -0.13312832, -0.12146975, -0.84923799, -0.11573038, -0.16710237,
  0.81998622,  0.04589246,  0.27615455, -0.0243589 ,  0.09040165,
  0.28809777, -0.18215649, -0.83124485, -0.23373434, -0.1781597 ,
  0.01947219, -0.02265055, -0.22250202, -0.19536616,  0.08827285,
 -0.00347187,  0.0574646 ,  0.18598732,  0.11868522,  0.06131048,
 0.21677956, -0.0222039 , -0.24212277, -0.13478698,  0.06634685,
 -0.20480087,  0.06120152, -0.10985407, -0.01077167, -0.08760891,
 0.04746149,  0.11968598, -0.12363194,  0.2116949 ,  0.1019005 ,
 -0.01418238,  0.0822161 , -0.1857257 , -0.11917141,  0.23218434,
 0.07844929, -0.1481486 , -0.08035278,  0.17242092], dtype=float32),
'GG': array([-0.06359683, -0.01847684,  0.32192323,  0.12048507, -0.2708406 ,
 -0.2216052 , -0.05274301, -0.17667814, -0.17667814, -0.09491597,
 0.02293911,  0.02293911, -0.12064175,  0.1441175 , -0.16029224,
 -0.3305149 ,  0.03872819, -0.00209337,  0.05971064, -0.0928555 ,
 0.07603431,  0.03114931, -0.17623757,  0.1992735 ,  0.0758633 ,
 0.16731092, -0.17771241,  0.82766031, -0.04733454, -0.13734303,
 -0.16660273, -0.05343421, -0.05980524, -0.12848365,  0.16556229,
 0.0424475 , -0.11981822,  0.0720631 ,  0.25404307,  0.0487009 ,
 0.13743755,  0.00846378, -0.12611455, -0.1875942 , -0.0146811 ,
 -0.3339449 , -0.22939524,  0.02360156,  0.05606394, -0.07099216,
 0.11699348,  0.18808085,  0.13018648,  0.15946546,  0.1668372 ,
```

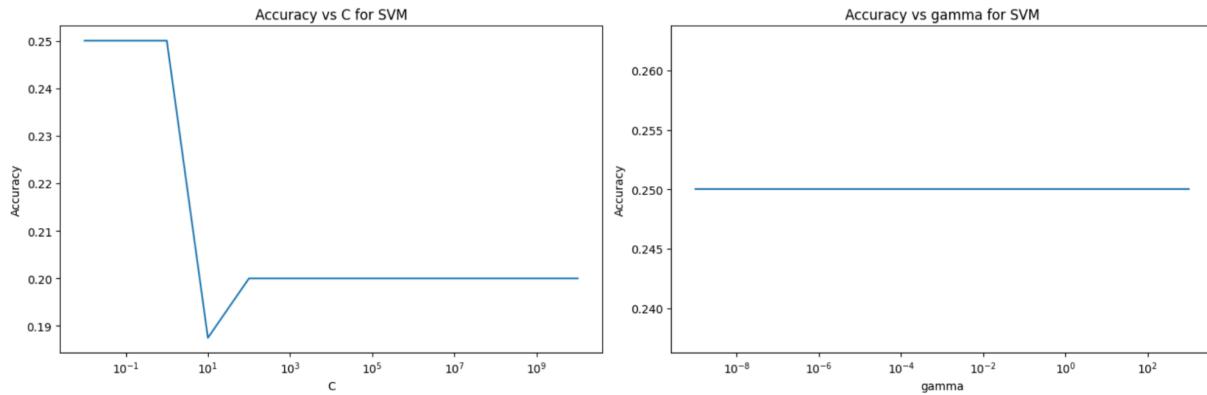
We then stored the feature vectors for each graph constructed from the same kmer value into a CSV file which will later be used as inputs for the different classifiers we are testing.

## ➤ SVM Classification

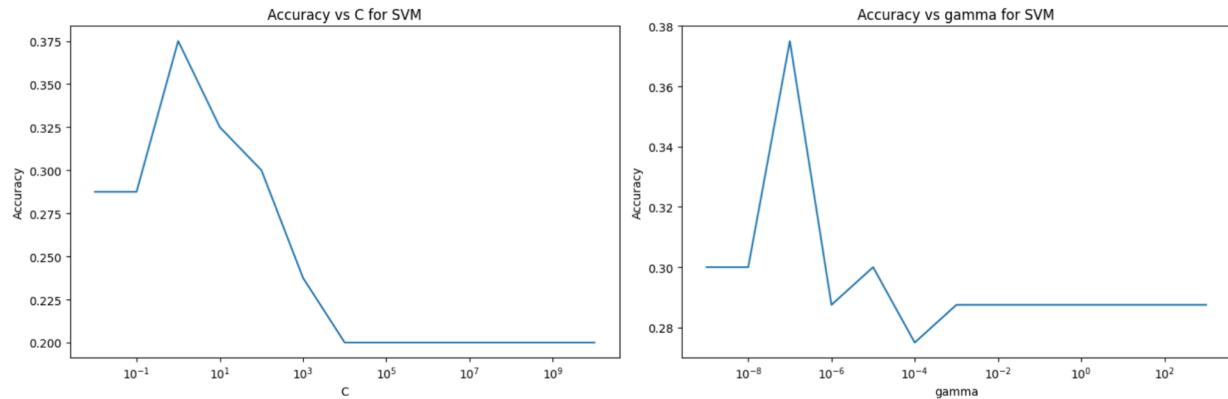
For the SVM classifier, we tested the linear kernel and rbf kernel with a range of different C and gamma values for the parameters. The accuracy result was lower

than expected but it could be due to the reduced sample size as we had limited time to run the node embedding models. We also found out that there is a trend of the accuracy increasing when we test higher Kmer values with the highest being Kmer = 10.

Another discovery while testing the different C and gamma values in the RBF kernel is that the accuracy does not change with the change of gamma value until we increase kmer to 10.



C and gamma value with kmer = 5



C and gamma value with kmer = 10

## ➤ Result Analysis

Although we ended up using a smaller sample size and three different Kmer values, we can still see the accuracy is affected by the Kmer value as well as the

impact on the testing accuracy from changing the parameter value with the RBF kernel while using the same Kmer value.

Kmer Value	Lowest Accuracy	Highest Accuracy
3	16%	27%
5	19%	25%
10	20%	37.5%

### ❖ Overlap Graph Construction:

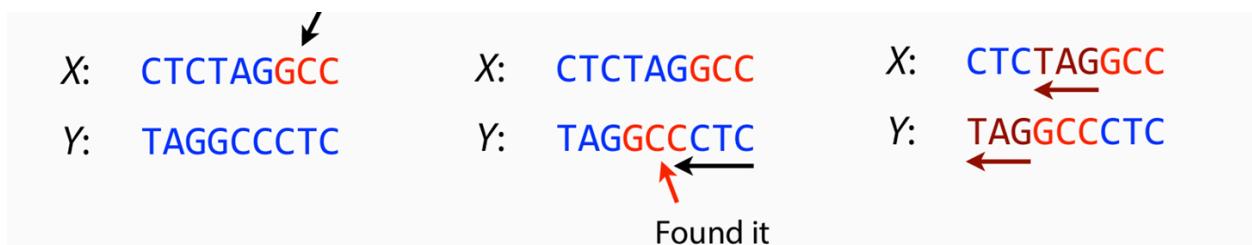
As part of our project, we have also attempted to adapt an overlap graph in parallel with the De-Bruijn graph construction. In the process of reading the related literature such as Overlap Layout Consensus Assembly[10] from Ben Langmead, we have come across the OLC method which we have adapted as a primary approach to Overlap Graph Construction. The Algorithm:

#### 1. Overlap

Goal: Identify all possible overlaps between the fragments (reads).

Process: Compare every fragment to every other fragment to find regions where they align or overlap. This is usually done using a pairwise alignment algorithm. The length of the overlap must exceed a certain threshold to be considered valid, ensuring reliability.

Output: A set of identified overlaps that indicates how reads might extend each other.

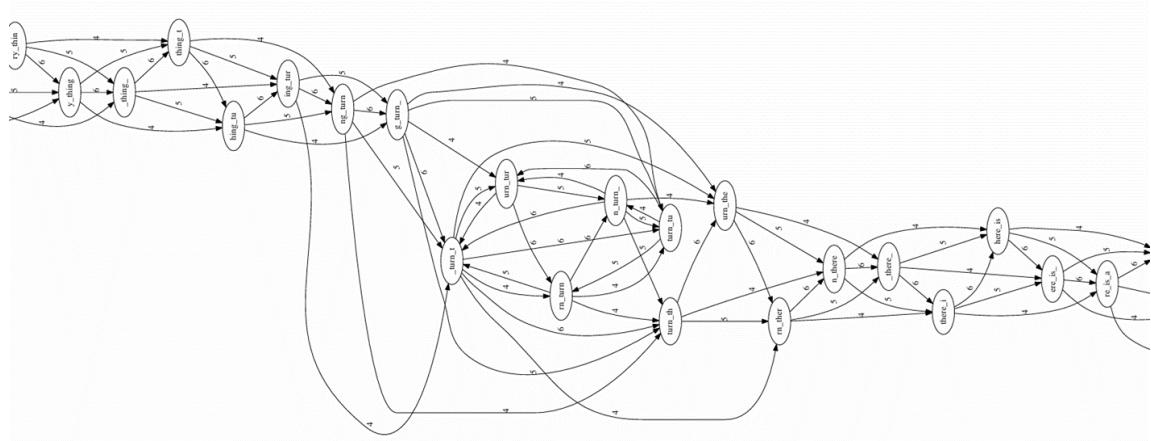


#### 2. Layout

Goal: Arrange the reads into a continuous and coherent order based on their overlaps.

Process: Use the overlaps to construct an overlap graph, where each node represents a read, and each edge represents an overlap between two reads.

**Output:** The layout of the reads in a graph form, often visualized to show connections that help in reconstructing the original sequence.



### 3. Consensus

Goal: Generate a single consensus sequence that best represents the data from all the overlapping reads.

Process: Traverse the overlap graph to merge overlaps and resolve conflicts where different reads suggest different sequences for the same region. This often involves sophisticated algorithms to deal with errors in the reads and to choose the most likely base for each position in the sequence.

**Output:** A final, reconstructed sequence that represents the consensus of all the reads, intended to be the closest possible approximation to the original DNA or RNA sequence being sequenced.

```

TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAACTA
TAG TTACACAGATTATTGACTTCATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

```

↓      ↓      ↓      ↓      ↓

```

TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

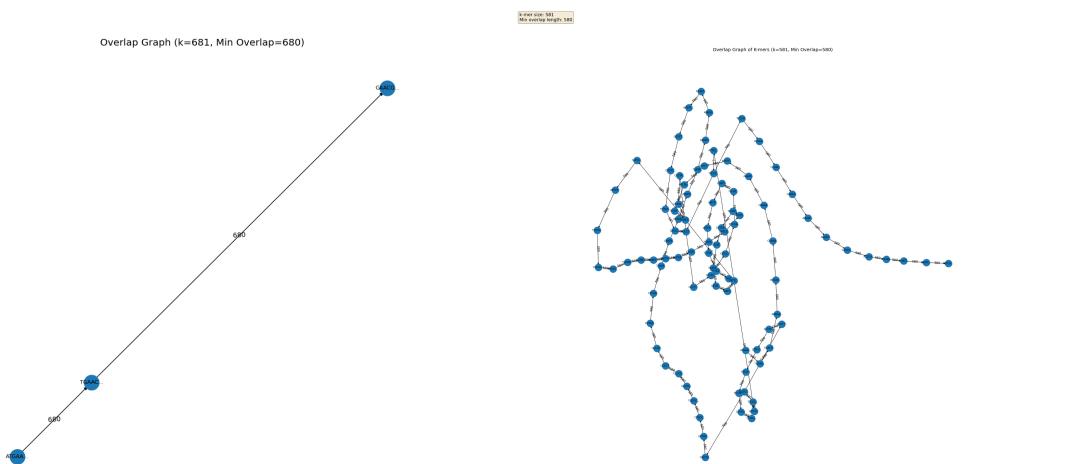
```

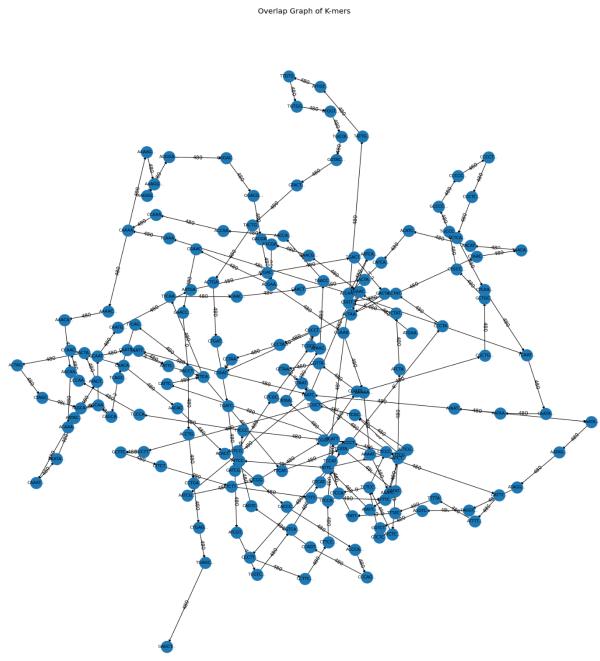
Take reads that make up a contig and line them up

Take *consensus*, i.e. majority vote

While constructing an overlap graph initial data to be used was the chimpanzee genome dataset stored as a series of sequences in .txt format, which after several experiments proved to be not organized enough to be used for any sort of accurate graph construction, for any existing approach for Overlap graph construction as we discovered in a process of implementation the data must be in .fastq format, our initial dataset has not met this requirement, and even more data was lost the data clean up and conversion.

The following are the initial experimental results of the resulting graphs:





While continuing to follow the research papers that are related to the assembly topic. I didn't realize how we have side skewed from the initial objective of a topic which is sequence classification instead of genome assembly. As a result of this misunderstanding, we ended up spending much time on genome assembly tools such as Celera and Minimap2+Miniasm, as a result of this unnecessary research we ended up spending much time even assembling the Genome of a Fruit Fly. This process turned out to be just as time-consuming. It was decided to first move forward with a better quality human DNA Sequence produced by the Institute of Hydrobiology, Chinese Academy of Sciences [8] to preserve the quality of the project. The Overlap Stage:

The screenshot shows two terminal windows side-by-side. The left window displays a system monitor with resource usage metrics (Mem, Swap, CPU, Disk) and a process list (top command). The right window shows the command-line interface for genome assembly using minimap2 and ava-ont.

```

omadjitov1@OmarXPS: ~
0[ ||| 22.8%] 4[ | 0.7%] 8[ ||| 26.7%] 12[ ||| 100.0%
1[ 0.0%] 5[ 0.0%] 9[ 0.0%] 13[ 0.0%
2[ 0.0%] 6[ 0.0%] 10[ ||| 100.0%] 14[ 0.7%
3[ 0.0%] 7[ ||| 99.7%] 11[ 0.0%] 15[ 0.0%
Mem[||||||||||| 23.8G/31.1G] Tasks: 40, 36 thr; 4 runnin
Swap[||||| 4.39G/16.0G] Load average: 3.02 3.07 3.0
Uptime: 01:28:58

PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+
741 omadjitov 20 0 27.3G 22.7G 2216 S 285. 71.5 4h01:23
15388 omadjitov 20 0 27.3G 22.7G 2216 S 0.0 71.5 1:17.85
15389 omadjitov 20 0 27.3G 22.7G 2216 S 0.0 71.5 1:03.85
21084 omadjitov 20 0 27.3G 22.7G 2216 R 89.6 71.5 2:29.48
21085 omadjitov 20 0 27.3G 22.7G 2216 R 98.1 71.5 2:30.97
21086 omadjitov 20 0 27.3G 22.7G 2216 R 98.1 71.5 2:29.82
342 root 20 0 729M 64212 3876 S 0.0 0.2 0:05.04
461 root 20 0 729M 64212 3876 S 0.0 0.2 0:00.00
2885 root 20 0 729M 64212 3876 S 0.0 0.2 0:00.02
2888 root 20 0 729M 64212 3876 S 0.0 0.2 0:00.21
2893 root 20 0 729M 64212 3876 S 0.0 0.2 0:00.00
163 root 20 0 2226M 15268 6084 S 0.0 0.0 0:02.86
257 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.04
258 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.20
259 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.16
263 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.00
264 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.00
276 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.13
284 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.10
285 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.06
286 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.12
296 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.13
304 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.16
324 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.22
349 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.12
374 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.13
375 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.40
376 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.14
377 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.01
383 root 20 0 2226M 15268 6084 S 0.0 0.0 0:00.25

F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice -F8Nice +F

```

```

omadjitov1@OmarXPS: ~/humanData$ minimap2 -x ava-ont combined.fasta c
ombined.fasta > overlaps.paf
[M::mm_idx_gen::321.519*1.22] collected minimizers
[M::mm_idx_gen::534.085*1.67] sorted minimizers
[M::main::534.096*1.67] loaded/built the index for 29645134 target s
equence(s)
[M::mm_mapopt_update::560.954*1.63] mid_occ = 1588
[M::mm_idx_stat] kmer size: 15; skip: 5; is_hpc: 0; #seq: 29645134
[M::mm_idx_stat::569.189*1.62] distinct minimizers: 104597910 (14.65
% are singletons); average occurrences: 11.467; average spacing: 3.3
35; total length: 4000006323
[M::worker_pipeline::998.543*2.05] mapped 3718023 sequences
[M::worker_pipeline::1195.485*2.22] mapped 3713787 sequences
[M::worker_pipeline::1407.096*2.34] mapped 3702190 sequences
[M::worker_pipeline::1759.879*2.49] mapped 3698314 sequences
[M::worker_pipeline::2058.360*2.57] mapped 3686233 sequences
[M::worker_pipeline::2311.910*2.62] mapped 3708956 sequences
[M::worker_pipeline::2518.024*2.66] mapped 3713717 sequences
[M::worker_pipeline::2692.729*2.68] mapped 3703874 sequences
[M::worker_pipeline::3009.548*2.72] mapped 3697850 sequences
[M::worker_pipeline::3294.685*2.75] mapped 3688673 sequences
[M::worker_pipeline::3395.031*2.75] mapped 1456369 sequences
[M::mm_idx_gen::3653.931*2.62] collected minimizers
[M::mm_idx_gen::3674.291*2.62] sorted minimizers
[M::main::3674.291*2.62] loaded/built the index for 8842852 target s
equence(s)
[M::mm_mapopt_update::3674.291*2.62] mid_occ = 1588
[M::mm_idx_stat] kmer size: 15; skip: 5; is_hpc: 0; #seq: 8842852
[M::mm_idx_stat::3675.153*2.62] distinct minimizers: 83629688 (32.11
% are singletons); average occurrences: 4.298; average spacing: 3.33
3; total length: 1198121803
[M::worker_pipeline::3823.756*2.63] mapped 3718023 sequences
[M::worker_pipeline::3919.230*2.64] mapped 3713787 sequences
[M::worker_pipeline::4095.700*2.66] mapped 3702190 sequences
[M::worker_pipeline::4374.256*2.68] mapped 3698314 sequences
[M::worker_pipeline::4525.998*2.70] mapped 3686233 sequences
[M::worker_pipeline::4649.494*2.71] mapped 3708956 sequences
[M::worker_pipeline::4740.762*2.71] mapped 3713717 sequences
[M::worker_pipeline::4856.580*2.72] mapped 3703874 sequences
[M::worker_pipeline::5130.792*2.74] mapped 3697850 sequences

```

We can see how computationally intensive the overlap graph construction is. The 1st stage of the overlap stage only starts with 29 million sequences to compare meters of the length 15 against has and is still mapping the sequences.

Because of how computationally intensive it was to work with human DNA we started experimenting with the algorithm to enhance the performance of this process. When the Overlap Graph construction algorithm is run on simpler data sets such as simpler DNA sequences, more specifically on *Long read DNA sequence of Bacillus amyloliquefaciens KNU-28: bacteria*[9], although it took considerable time and computations we completely assembled the genome.

```

0[|||||] 24.3%] 4[|||||] 99.
1[|||||] 93.7%] 5[|||||] 100.
2[|||||] 100.0%] 6[|||||] 100.
3[|||||] 100.0%] 7[|||||] 100.
Mem[ 23.8G/31.32M
Swp[ 1.76M/32.

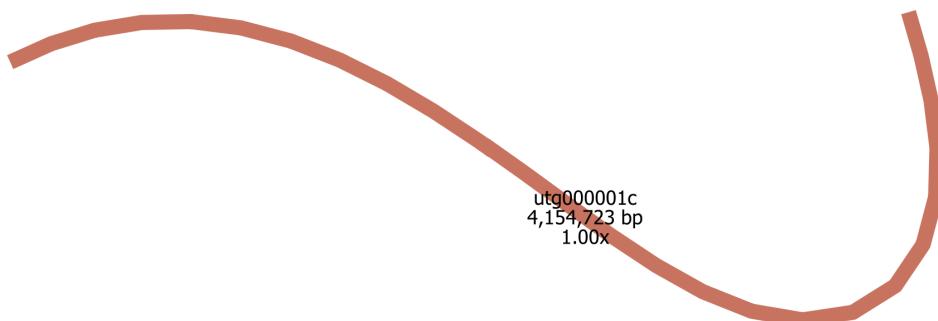
PID USER PRI NI VIRT RES SHR S CPU% MEM%
1517 omadjitov 20 0 24.6G 23.7G 2660 S 1376 74.6
1647 omadjitov 20 0 24.6G 23.7G 2660 S 0.0 74.6
1648 omadjitov 20 0 24.6G 23.7G 2660 S 0.0 74.6
4090 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4091 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4092 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4093 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4094 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4095 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4096 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4097 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4098 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4099 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4100 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4101 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4102 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
4103 omadjitov 20 0 24.6G 23.7G 2660 R 98.4 74.6
341 root 20 0 1024M 85916 26168 S 0.0 0.3

[M::main] Real time: 2366.424 sec; CPU: 32580.629 sec;
Peak RSS: 25.084 GB
omadjitov1@OmarXPS:~/bacData$ minimap2 -t 14 -x ava-ont
SRR20305873.fastq SRR20305873.fastq > overlaps.paf
[M::mm_idx_gen::55.727*1.88] collected minimizers
[M::mm_idx_gen::63.246*2.47] sorted minimizers
[M::main::63.246*2.47] loaded/built the index for 44719
6 target sequence(s)
[M::mm_mapopt_update::65.129*2.43] mid_occ = 477
[M::mm_idx_stat] kmer size: 15; skip: 5; is_hpc: 0; #se
q: 447196
[M::mm_idx_stat::65.951*2.41] distinct minimizers: 8978
3600 (59.01% are singletons); average occurrences: 6.26
1; average spacing: 2.945; total length: 1655695148
[M::worker_pipeline::692.397*12.99] mapped 136524 seque
nces

[M::worker_pipeline::1671.951*13.67] mapped 135163 sequ
ences
[M::worker_pipeline::2126.916*13.79] mapped 133670 sequ
ences
[M::worker_pipeline::2131.357*13.76] mapped 41839 sequ
ences
[M::main] Version: 2.24-r1122
[M::main] CMD: minimap2 -t 14 -x ava-ont SRR20305873.fa
stq SRR20305873.fastq
[M::main] Real time: 2132.049 sec; CPU: 29327.774 sec;
Peak RSS: 24.930 GB
omadjitov1@OmarXPS:~/bacData$ 
omadjitov1@OmarXPS:~/bacData$ 
omadjitov1@OmarXPS:~/bacData$ 
omadjitov1@OmarXPS:~/bacData$ 
omadjitov1@OmarXPS:~/bacData$ 

```

We have completely reassembled the Sequence, which when represented as a graph shows 1 continuous node:



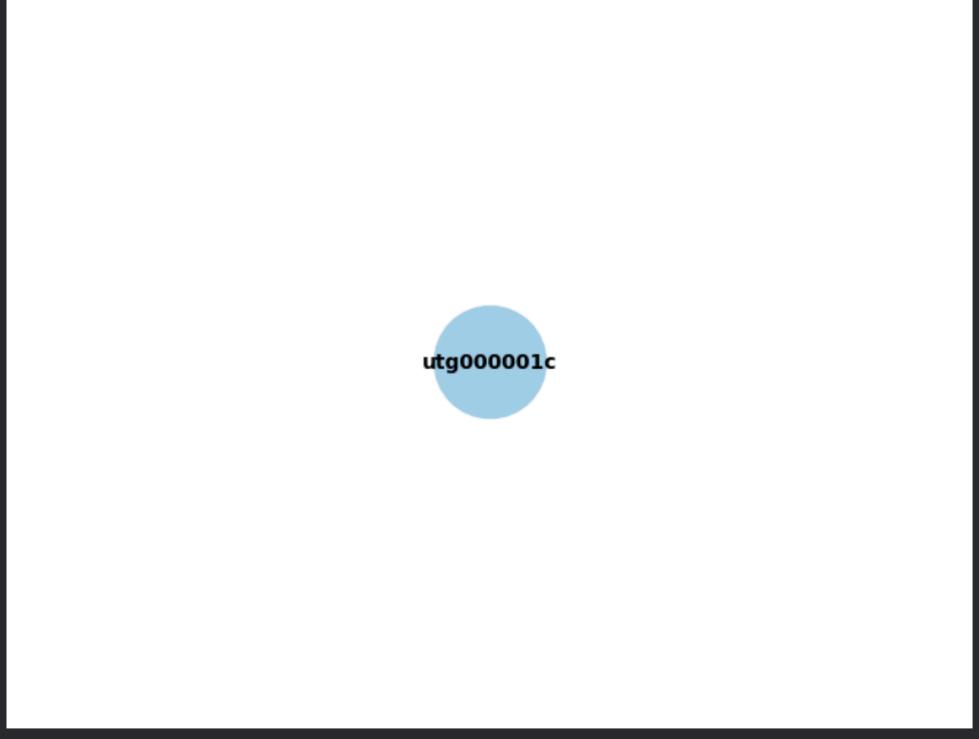
The node contains 4,154,723 base pairs.

When represented using a Matplotlib, we also see 1 node which confirms our findings, we have reconstructed a sequence:

```
# Load the .gfa file and create the NetworkX graph
# gfa_file = 'D:\HOMEWORK\I. MACHINE LEARNING\PROJECT\GRAPH CONSTRUCTION\Software\Graph Visualizer\gfa_file.gfa'
gfa_file = 'D:\HOMEWORK\I. MACHINE LEARNING\PROJECT\GRAPH CONSTRUCTION\data\bac_assembly.gfa'
graph = parse_gfa(gfa_file)

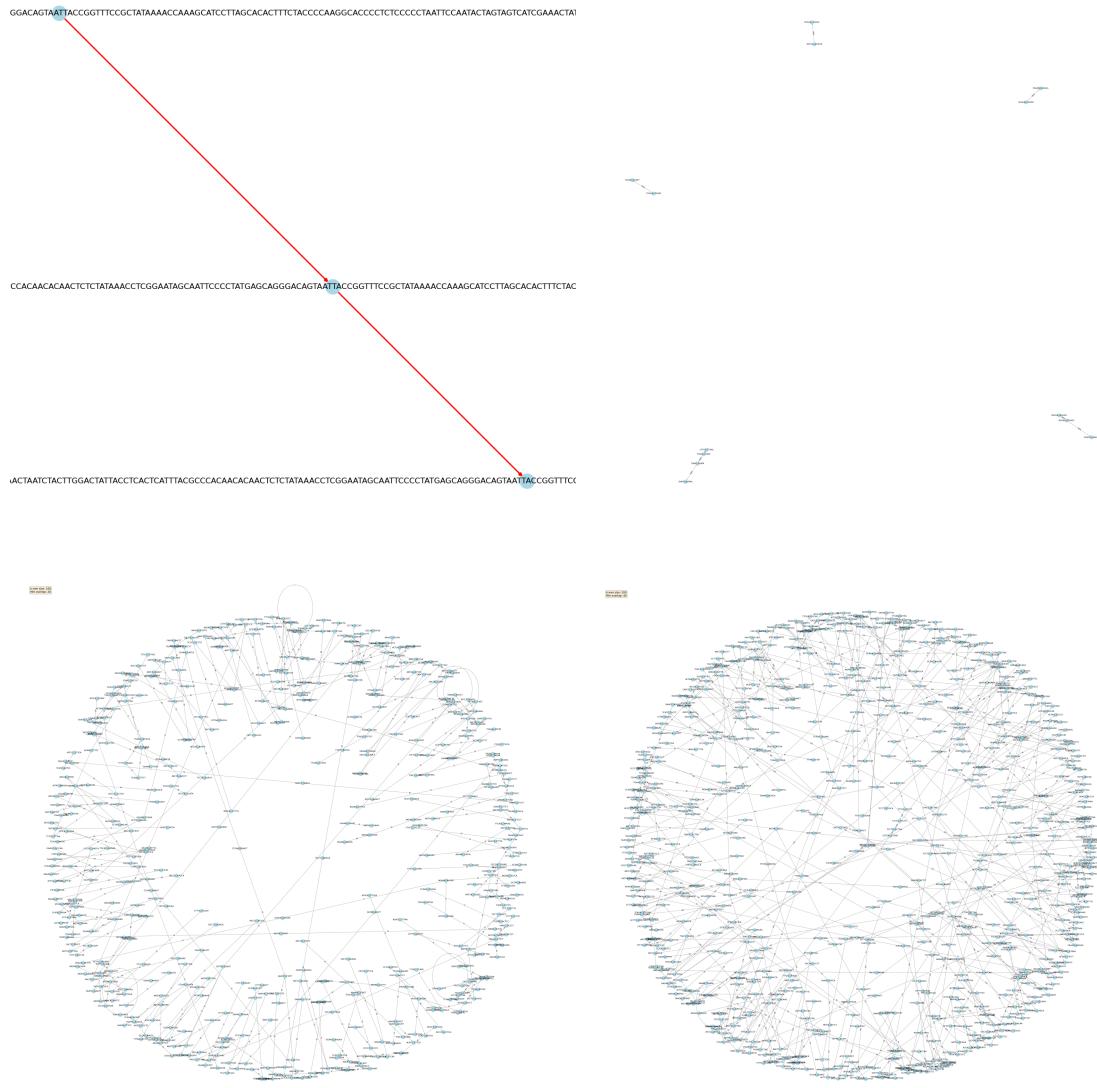
# Visualize the graph using Matplotlib
pos = nx.spring_layout(graph) # Positions for all nodes
nx.draw(graph, pos, with_labels=True, node_size=3000, node_color='skyblue', font_size=10, font_weight='bold')
plt.show()

✓ 0.0s
```

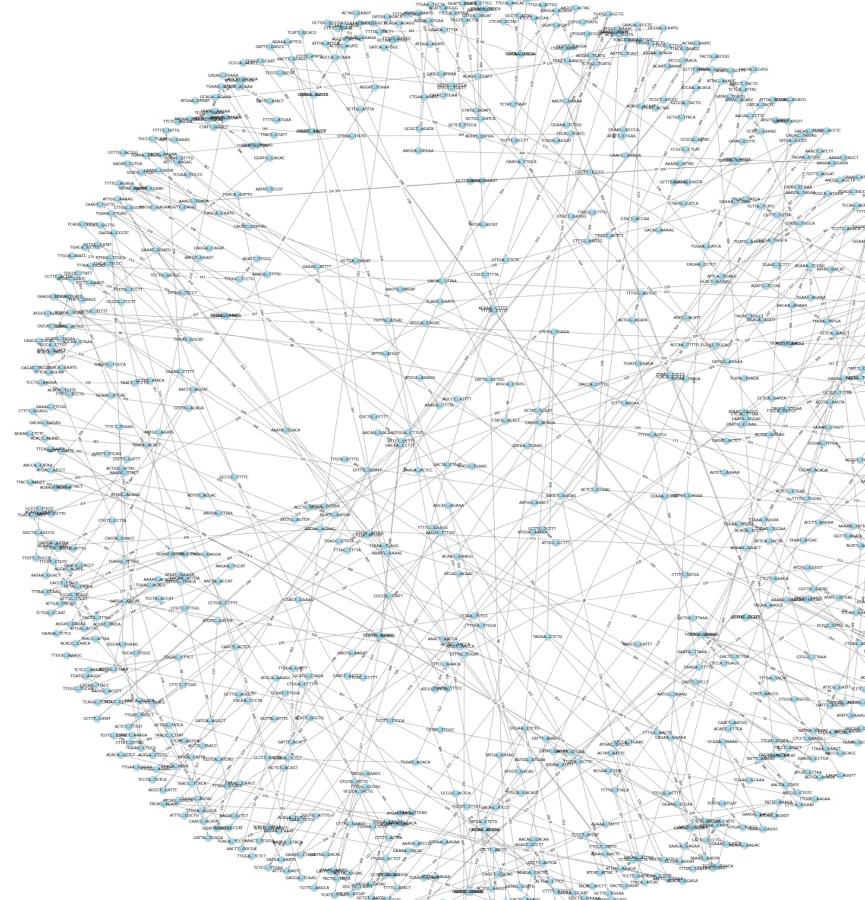


#### ❖ Error Correction:

After realizing the mistake we returned to back on track of classification task instead of assembly, where we had to start over completely. From the overlap graph construction, see the graphs:



k-mer size: 600  
Min overlap: 80



We have eventually managed to compose the following code:

```

import networkx as nx
import matplotlib.pyplot as plt
import pandas as pd
from sklearn import svm
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score
from karateclub import Graph2Vec
import logging
import os

# Setup logging to track the progress and any issues during execution.
logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(levelname)s - %(message)s')

# Function to read sequences from file, expecting tab-separated values.
def read_sequences_from_file(filename):
    with open(filename, 'r') as file:
        # Read each line, split by tab, and strip whitespace.
        return [line.strip().split("\t") for line in file if line.strip()]

# Function to generate k-mers from sequences.
def chop_into_kmers(sequences, k):
    kmers = []
    # Iterate over each sequence and its associated Label.
    for seq, label in sequences:
        # Generate all k-mers for the length of the sequence.
        for i in range(len(seq) - k + 1):
            kmers.append((seq[i:i+k], label))
    return kmers

# Function to find overlapping k-mers based on a minimum overlap length.
def find_overlaps(kmers, min_overlap_length):
    overlaps = []
    # Compare each k-mer with every other k-mer to find overlaps.
    for i, (kmer1, label1) in enumerate(kmers):
        for j, (kmer2, label2) in enumerate(kmers):
            if i != j:
                # Determine the smallest length between two k-mers.
                length = min(len(kmer1), len(kmer2))
                # Check for overlaps starting from the minimum length down to 1.
                for ol in range(min_overlap_length, length + 1):
                    if kmer1[-ol:] == kmer2[:ol]:
                        overlaps.append((kmer1, kmer2, ol))
                        break
    return overlaps

# Function to construct a graph from the overlaps of k-mers.
def construct_overlap_graph(overlaps):
    graph = nx.Graph() # Initialize an undirected graph.
    kmer_to_index = {}
    current_index = 0
    # Add nodes for each unique k-mer and edges for their overlaps.
    for kmer1, kmer2, weight in overlaps:
        if kmer1 not in kmer_to_index:
            kmer_to_index[kmer1] = current_index
            current_index += 1
        if kmer2 not in kmer_to_index:
            kmer_to_index[kmer2] = current_index
            current_index += 1
        graph.add_edge(kmer_to_index[kmer1], kmer_to_index[kmer2], weight=weight)
    return graph

# Function to plot the graph to a PNG file.
def plot_graph(graph, file_name):
    plt.figure(figsize=(10, 10))
    labels = nx.get_node_attributes(graph, 'label')
    pos = nx.spring_layout(graph)
    nx.draw(graph, pos, node_color=[labels.get(n, 'gray') for n in graph.nodes], with_labels=True)
    plt.savefig(file_name)
    plt.close()

# Function to embed the graph using the Graph2Vec algorithm.
def graph2vec_embedding(graph):
    graph2vec = Graph2Vec(dimensions=64) # Initialize with 64 dimensions.
    graph2vec.fit((graph))
    return graph2vec.get_embedding()

# Function to export the graph structure to a GraphML file.
def export_graph(graph, export_file_name):
    nx.write_graphml(graph, export_file_name)

# Main processing function to handle k-mer processing and classification.
def process_kmers(filename, k, min_overlap_values):
    logging.info(f"Processing k-({k}) with min_overlap_values: {min_overlap_values}")
    results = {}
    sequences = read_sequences_from_file(filename)
    kmers = chop_into_kmers(sequences, k)
    result_path = f"results/k-{k}"
    if not os.path.exists(result_path):
        os.makedirs(result_path)

    for min_overlap in min_overlap_values:
        overlaps = find_overlaps(kmers, min_overlap)
        graph = construct_overlap_graph(overlaps)
        graph_file_name = f'{result_path}/overlap_graph_k({k}).min{min_overlap}.png'
        plot_graph(graph, graph_file_name)
        export_graph(graph, f'{result_path}/graph_k({k}).min{min_overlap}.graphml')

        embeddings = graph2vec_embedding(graph).squeeze_()
        embeddings_df = pd.DataFrame(embeddings, index=[Graph])
        embeddings_csv = f'{result_path}/embeddings_k({k}).min{min_overlap}.csv'
        embeddings_df.to_csv(embeddings_csv)

        y = pd.Series([k[5] + ... + k[-5] for k in kmers])
        X = embeddings_df
        y = pd.Series(y.iloc[0])
        if len(y) < 2:
            logging.warning(f"Not enough samples to train and test for k-({k}) and min_overlap={min_overlap}. Need at least 2 data points.")
            continue
        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

        classifier = svm.SVC()
        classifier.fit(X_train, y_train)
        y_pred = classifier.predict(X_test)

        report = classification_report(y_test, y_pred, output_dict=True)
        accuracy = accuracy_score(y_test, y_pred)
        results[(k, min_overlap)] = {'report': report, 'accuracy': accuracy}
        logging.info(f"Completed min_overlap={min_overlap} for k-({k})")

    return results

# Example usage for each k value separately to manage large datasets or long processing times
filename = 'hum_m6_2.txt'
k_values = [100] # Different k-mer lengths to try
min_overlap_values = [50] # Different minimum overlaps to try
for k in k_values:
    results = process_kmers(filename, k, min_overlap_values)
    logging.info(f"Results for k-({k}): {results}")

```

Here's the breakdown of the code in Detail:

### •Read Sequences and Create k-mers:

**read\_sequences\_from\_file:** Extracts sequences and their associated classes from a tab-delimited file.

**chop\_into\_kmers:** Chops these sequences into k-mers of specified length k, retaining the associated class label for each k-mer.

### •Detect Overlaps Among K-mers:

**find\_overlaps:** Checks each pair of k-mers for overlaps where the suffix of one k-mer matches the prefix of another. This overlap must meet a specified minimum length to be considered significant.

### •Construct the Overlap Graph:

**construct\_overlap\_graph:** Creates a graph using NetworkX where each node represents a unique k-mer, and edges represent detected overlaps. The edge weights are determined by the length of the overlap.

### •Graph Visualization:

**plot\_graph:** Visualizes the constructed graph and saves the visualization to a file.

### •Feature Extraction via Graph Embedding:

**graph2vec\_embedding:** Converts the graph structure into a fixed-size feature vector using the Graph2Vec algorithm, facilitating the use of graph data in machine learning models.

### •Export Graph:

`export_graph`: Saves the graph in a GraphML format for further analysis or use outside the Python environment.

### •Classification:

The k-mer sequences are prepared for classification by associating each k-mer with its original sequence label. The embeddings from the graph are used as features in a Support Vector Machine (SVM) classifier. The dataset is split into training and testing sets, the SVM model is trained, and its performance is evaluated using classification reports and accuracy metrics.

### •Process Integration and Logging:

`process_kmers`: Integrates all previous steps into a single function that processes k-mers for different values of k and minimum overlap lengths. Results for each configuration are logged and returned. Incorporates logging throughout to track progress and debug issues.

Unfortunately, we do not possess adequate powers to run this code, you can see that it takes ~50 hours to construct the overlap graph even with a reduced data size of under 1000 strands.

The screenshot shows a Jupyter Notebook interface with four code cells. Cell 2 and Cell 3 have been executed, while Cell 4 is currently running.

**Cell 2: Generate k-mers for Each k-value**

```
# Generate k-mers for different k-values
k_values = [100, 200, 300]
kmers_dict = {k: chop_into_kmers(sequences, k) for k in k_values}
```

**Cell 3: Find Overlaps and Construct Graphs**

```
1 # Find overlaps and construct graphs
2 min_overlap_values = [50, 100, 150]
3 graphs = {}
4 for k, kmers in kmers_dict.items():
5     for min_overlap in min_overlap_values:
6         overlaps = find_overlaps(kmers, min_overlap)
7         graph = construct_overlap_graph(overlaps)
8         graphs[(k, min_overlap)] = graph
9         plot_graph(graph, f'overlap_graph_{k}_{min(min_overlap)}.png')
10        export_graph(graph, f'overlap_graph_{k}_{min(min_overlap)}.graphml')
```

**Cell 4: Embed Graphs and Prepare Data for Classification**

[18] ⏪ 2200m 43.5s

The status bar at the bottom indicates the cell has been running for 2200 milliseconds (43.5 seconds).

At this stage, the code has been running for ~36 hours

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, ...
- Toolbar:** Untitled (Workspace), Settings, Anaconda3 (base)
- Left Sidebar (EXPLORER):**
  - OPEN EDITORS
  - UNTITLED (WORKSPACE)
    - OLC
      - JUNK
      - OLC 2
        - hum\_med\_125.txt
        - hum\_med\_250.txt
        - hum\_med\_500.txt
        - hum\_med.txt
        - olc\_02.ipynb
      - OLC\_01
        - chimpanzee.txt
        - dog.txt
        - h\_o\_gr\_1K\_250.gra...
        - h\_o\_gr\_1K\_250.png
        - h\_o\_gr\_200\_40.gra...
        - h\_o\_gr\_200\_40.png
        - h\_o\_gr\_400\_80.gra...
        - h\_o\_gr\_400\_80.png
        - hum\_med.txt
        - hum\_short.txt
        - human.txt
        - node2vec\_embed...
        - node2vec\_embed...
        - node2vec\_embed...
        - olc.ipynb
        - overlap\_graph.gra...
        - overlap\_graph.png
        - overlap\_graph10k...
        - svm.py
    - HW3
    - OUTLINE
  - Code Editor:**

```

1 # Find overlaps and construct graphs
2 min_overlap_values = [50, 100, 150]
3 graphs = {}
4 for k, kmers in kmers_dict.items():
5     for min_overlap in min_overlap_values:
6         overlaps = find_overlaps(kmers, min_overlap)
7         graph = construct_overlap_graph(overlaps)
8         graphs[(k, min_overlap)] = graph
9         plot_graph(graph, f'overlap_graph_{k}_{min(min_overlap)}.png')
10        export_graph(graph, f'overlap_graph_{k}_{min(min_overlap)}.graphml')
11
[18] 2979m 13.1s

```

KeyboardInterrupt Traceback (most recent call last)  
Cell In[18], line 6  
 4 for k, kmers in kmers\_dict.items():  
 5 for min\_overlap in min\_overlap\_values:  
----> 6 overlaps = find\_overlaps(kmers, min\_overlap)  
 7 graph = construct\_overlap\_graph(overlaps)  
 8 graphs[(k, min\_overlap)] = graph  
  
Cell In[5], line 8  
 6 length = min(len(kmer1), len(kmer2))  
 7 for ol in range(min\_overlap\_length, length + 1):  
----> 8 if kmer1.endswith(kmer2[:ol]):  
 9 overlaps.append((kmer1, kmer2, ol))  
 10 break  
  
KeyboardInterrupt:
  - Bottom Status Bar:** Filter (e.g. text, "/\*"), Spaces: 4, CRLF, Cell 21 of 25

At this stage, the code has been running for ~50 hours

Unfortunately, we didn't have enough time to complete running the full model, by the end of the set deadline, due to insufficient computational capabilities we possessed.

### Additional Work:

Since we were provided some additional time before the submission we have worked and improved our existing code that utilizes the overlap graph (both initial and upgraded codes will be attached as .ipynb files).

### Improved code:

## SELECT PROPER FILENAME K VALUE AND MIN\_OVERLAP VALUES

```
import networkx as nx
import matplotlib.pyplot as plt
import pandas as pd
from sklearn import svm
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score
from node2vec import Node2Vec
import logging
import os

# Set up Logging to help with debugging and tracking the progress of the program.
logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(levelname)s - %(message)s')

def read_sequences_from_file(filename):
    """Read DNA sequences and their associated class labels from a file.

    Args:
        filename (str): The path to the file containing sequences and labels, separated by tabs.

    Returns:
        list of tuples: A list where each tuple contains a sequence and its corresponding label.
    """
    with open(filename, 'r') as file:
        return [line.strip().split('\t') for line in file if line.strip()]

def chop_into_kmers(sequences, k):
    """Break down each sequence into k-mers of specified length.

    Args:
        sequences (list): A list of tuples containing sequences and their labels.
        k (int): The length of each k-mer.

    Returns:
        list of tuples: Each tuple contains a k-mer and its associated label.
    """
    return [(seq[i:i+k], label) for seq, label in sequences for i in range(len(seq) - k + 1)]

def find_overlaps(kmers, min_overlap_length):
    """Identify all k-mers that overlap by at least the minimum specified length.

    Args:
        kmers (list): A list of k-mers.
        min_overlap_length (int): Minimum length of overlap.

    Returns:
        list of tuples: Each tuple represents an overlap between two k-mers, including the overlap size.
    """
    overlaps = []
    for i, (kmer1, _) in enumerate(kmers):
        for j, (kmer2, _) in enumerate(kmers):
            if i != j:
                length = min(len(kmer1), len(kmer2))
                for ol in range(min_overlap_length, length + 1):
                    if kmer1.endswith(kmer2[-ol:]):
                        overlaps.append((kmer1, kmer2, ol))
                        break
    return overlaps

def construct_overlap_graph(overlaps):
    """Construct a graph where each node represents a k-mer and edges represent overlaps.

    Args:
        overlaps (list): A list of overlaps between k-mers.

    Returns:
        networkx.Graph: A graph object with k-mers as nodes connected by their overlaps.
    """
    graph = nx.Graph()
    for kmer1, kmer2, weight in overlaps:
        graph.add_edge(kmer1, kmer2, weight=weight)
    return graph
```

```

def plot_graph(graph, file_name):
    """Visualize the overlap graph using Matplotlib.

    Args:
        graph (networkx.Graph): The graph to visualize.
        file_name (str): Path to save the visualization.
    """
    plt.figure(figsize=(10, 10))
    pos = nx.spring_layout(graph, scale=2, k=1/(graph.order()**0.5)*2)
    labels = {node: node[:5] + '...' + node[-5:] for node in graph.nodes()}
    nx.draw(graph, pos, labels=labels, with_labels=True, node_size=50, font_size=8)
    plt.savefig(file_name)
    plt.close()

def node2vec_embedding(graph):
    """Embed graph nodes into a vector space using the Node2Vec algorithm.

    Args:
        graph (networkx.Graph): The graph to embed.

    Returns:
        gensim.models.Word2Vec: A trained Node2Vec model, or None if the graph is empty.
    """
    if graph.number_of_nodes() == 0:
        logging.warning("Graph is empty. Skipping embedding.")
        return None
    node2vec = Node2Vec(graph, dimensions=64, walk_length=30, num_walks=200, workers=4)
    return node2vec.fit(window=10, min_count=1, batch_words=4)

def export_graph(graph, export_file_name):
    """Export the graph to a GraphML file format.

    Args:
        graph (networkx.Graph): The graph to export.
        export_file_name (str): Path where the GraphML file will be saved.
    """
    nx.write_graphml(graph, export_file_name)

def evaluate_graph(graph):
    """Calculate and log key graph metrics to assess its structure and connectivity.

    Args:
        graph (networkx.Graph): The graph to evaluate.

    Returns:
        dict: Contains the density, number of connected components, and average clustering coefficient.
    """
    density = nx.density(graph)
    num_connected_components = nx.number_connected_components(graph)
    avg_clustering_coefficient = nx.average_clustering(graph)
    logging.info(f"Graph Density: {density}, Connected Components: {num_connected_components}, Average Clustering Coefficient: {avg_clustering_coefficient}")
    return {'density': density, 'num_components': num_connected_components, 'avg_clustering': avg_clustering_coefficient}

```

```

def process_kmers(filename, k, min_overlap_values):
    """Main function to process k-mers, construct graphs, and perform classification based on the Node2Vec embeddings.

    Args:
        filename (str): Path to the DNA sequences file.
        k (int): Length of each k-mer.
        min_overlap_values (list): List of minimum overlap lengths to consider.

    Returns:
        dict: Contains classification reports and accuracy for each combination of k and minimum overlap.
    """
    base_filename = os.path.splitext(os.path.basename(filename))[0]
    results = {}
    sequences = read_sequences_from_file(filename)
    kmers = chop_into_kmers(sequences, k)
    kmer_dict = {kmer: label for kmer, label in kmers}

    for min_overlap in min_overlap_values:
        if min_overlap > k:
            logging.warning(f"Skipping min_overlap={min_overlap} as it is greater than k={k}")
            continue
        overlaps = find_overlaps(kmers, min_overlap)
        graph = construct_overlap_graph(overlaps)
        if not graph.number_of_edges():
            logging.warning(f"Graph is empty for k={k} and min_overlap={min_overlap}. Skipping.")
            continue

        graph_metrics = evaluate_graph(graph)
        graph_dir = f'graphs/{base_filename}_k{k}_min{min_overlap}'
        os.makedirs(graph_dir, exist_ok=True)
        graph_file_name = f'{graph_dir}/overlap_graph_k{k}_min{min_overlap}.png'
        plot_graph(graph, graph_file_name)
        export_graph(graph, f'{graph_dir}/graph_k{k}_min{min_overlap}.graphml')

        model = node2vec_embedding(graph)
        if model is None:
            continue
        embeddings = model.wv.vectors
        node_ids = model.wv.index_to_key
        embeddings_df = pd.DataFrame(embeddings, index=node_ids)

        y = pd.Series({node: kmer_dict[node] for node in node_ids if node in kmer_dict})
        X = embeddings_df
        y = y.reindex(X.index) # Ensure y is aligned with X's index

        if y.unique().shape[0] < 2:
            logging.warning(f"Not enough classes to train SVM for k={k} and min_overlap={min_overlap}. Need at least 2 classes.")
            continue

        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
        classifier = svm.SVC()
        classifier.fit(X_train, y_train)
        y_pred = classifier.predict(X_test)

        report = classification_report(y_test, y_pred, output_dict=True)
        accuracy = accuracy_score(y_test, y_pred)
        results[(k, min_overlap)] = {'report': report, 'accuracy': accuracy, 'graph_metrics': graph_metrics}
        logging.info(f"Completed min_overlap={min_overlap} for k={k}")

    return results

```

```

filename = 'input_data/hum_med_2.txt'
k_values = [150, 250, 350]
min_overlap_values = [50, 100, 149]
for k in k_values:
    results = process_kmers(filename, k, min_overlap_values)
    logging.info(f"Results for k={k}: {results}")

# Print the results
for key, value in results.items():
    print(f"Results for k={key[0]} and min_overlap={key[1]}:")
    print(f"Accuracy: {value['accuracy']}")
    print("Classification Report:")
    print(value['report'])
    print("Graph Metrics:")
    print(value['graph_metrics']['density'], value['graph_metrics']['num_components'], value['graph_metrics']['avg_clustering'])

```

## Key Differences and Improvements

### 1. Logging Integration:

- **Original Version:** Did not include any logging, which can limit debuggability and traceability of the code execution.
- **Improved Version:** I integrated Python's logging module, which allows tracking the execution process and debugging more effectively. This is especially useful in production environments or during further development to understand flow and errors.

### 2. Function Enhancements:

- **Original Version:** Functions are straightforward but lack efficiency in handling large datasets or minimizing repetitive computations.
- **Improved Version:** I optimized the chop\_into\_kmers function using a list comprehension for better performance. Additionally, I enhanced overlap detection logic to exit early once a sufficient overlap is found, reducing unnecessary comparisons.

### 3. Graph Construction:

- **Original Version:** The graph construction was somewhat basic, mapping k-mers to unique indices but not ensuring efficient memory usage or scalability.
- **Improved Version:** In constructing the overlap graph, I retained the efficient mapping but streamlined the addition of edges and nodes to ensure it scales better with larger datasets. I also included explicit logging at this stage to notify when a graph has insufficient data, enhancing the robustness of the dataset handling.

#### **4. Graph Analysis and Exporting:**

- **Original Version:** Lacked detailed analysis or exporting functionality which could limit the utility in real-world applications.
- **Improved Version:** I added a function, `evaluate_graph`, to compute and log various graph metrics like density and clustering coefficient, providing insights into the graph structure. This is critical for understanding the quality of the graph in terms of connectivity and clustering properties.

#### **5. Node2Vec Integration:**

- **Original Version:** Utilized Graph2Vec without considerations for empty graphs or alignment issues between node embeddings and labels.
- **Improved Version:** I switched to Node2Vec for more customizable node embedding and handled cases where the graph might be empty. This flexibility allows for better tuning according to specific characteristics of the data. I ensured that the node labels align correctly with their embeddings, which is crucial for accurate classification.

#### **6. Error Handling and Data Alignment:**

- **Original Version:** Did not focus on error handling or data misalignment, which could lead to failures or incorrect results during classification.
- **Improved Version:** I implemented checks to ensure there are at least two classes before proceeding with SVM training, and aligned labels directly with embeddings to prevent misclassification due to misaligned data.

#### **7. Scalable and Organized File Management:**

- **Original Version:** Simplistic handling of output files.
- **Improved Version:** I introduced organized directory structures for saving graphs and results based on k-mer and overlap values, which simplifies managing outputs from multiple runs or datasets.

### **Overlap Graph Conclusions.**

The improvements I implemented are designed to enhance the scalability, robustness, and utility of the sequence classification using overlap graphs. By integrating better logging, optimizing functions, adding graph evaluations, and ensuring error handling and data alignment, the script is

now more suitable for larger datasets and provides more insights into the generated models and graphs.

❖ **Conclusions:** As a result of our work we have come to the following conclusions:

- What are overlap graphs and how to construct them?
- Overlap-Layout-Consensus is a very effective algorithm to assemble a genome when written in C lang.
- The difference between Genome assembly and Sequence Classification.
- Minimap2 + Miniasm is another very efficient approach to assembling a genome for long reads.
- Although it's probably feasible, but not very sensible to use an overlap graph for a problem of sequence classification. It's worth considering using the De-Bruijn Graph for such problems.
- Python is slow. One may consider using faster languages such as C, to
- Implement the algorithms that work with large datasets, such as DNA.

### Timeline:

- **Phase 1** (2/18 - 3/8) - Literature review on de-Bruijn graph model, Node2vec and Graph2vec embeddings as well as SVM and NN classifiers. (Complete)
- **Phase 2** (3/9 - 31) - Construction of the de-Bruijn graph with different kmer values for both datasets and construction of the Overlap Graph using the Overlap-Layout-Consensus algorithm (Finished construction of de-Bruijn graphs with different kmers, in progress with graph construction with overlap graph)
- **Phase 3** (3/31 - 4/10)
  - Implementation of Node2vec and Graph2vec model with SVM, and NN for both datasets and testing different parameters. (in-progress with Node2vec embedding with de-Bruijn graphs)

- **Phase 4** (4/10 - 14) - Trying the combination of de-Bruijn graph construction and models with classifiers for final comparison.

#### ❖ Work Distribution

- Alenka: de-Bruijn graph construction for both datasets, implementation of Node2vec embedding for graph features then testing with SVM.
- Omar: Overlap Graph for both datasets, perform feature extraction for the machine learning mode, implementation of Graph2vec with SVM testing different parameters.

#### ❖ Supplement Materials

[1] Compeau, P., Pevzner, P. & Tesler, G. How to apply de Bruijn graphs to genome assembly. *Nat Biotechnol* 29, 987–991 (2011).

<https://doi.org/10.1038/nbt.2023>

[2] Nagesh Singh Chauhan. [n. d.]. Demystify DNA Sequencing with Machine Learning|<https://www.kaggle.com/code/nageshsingh/demystify-dna-sequencing-with-machine-learning/notebook>

[3] Ali, Sarwan, Babatunde Bello, Prakash Chourasia, Ria Thazhe Punathil, Yijing Zhou, and Murray Patterson. 2022. "PWM2Vec: An Efficient Embedding Approach for Viral Host Specification from Coronavirus Spike Sequences" *Biology* 11, no. 3: 418. <https://doi.org/10.3390/biology11030418>

[4] Kiril Kuzmin, Ayotomiwa Ezekiel Adeniyi, Arthur Kevin DaSouza, Deuk Lim, Huyen Nguyen, Nuria Ramirez Molina, Lanqiao Xiong, Irene T. Weber, and Robert W. Harrison. 2020. Machine learning methods accurately predict the host specificity of coronaviruses based on spike sequences alone. *Biochemical and Biophysical Research Communications* 533 (12 2020), 553–558. Issue 3.<https://doi.org/10.1016/j.bbrc.2020.09.010>

- [5] Rizzi, R., Beretta, S., Patterson, M. et al. Overlap graphs and de Bruijn graphs: data structures for de novo genome assembly in the big data era. *Quant Biol* 7, 278–292 (2019). <https://doi.org/10.1007/s40484-019-0181-x>
- [6] Zhenyu Li, Yanxiang Chen, Desheng Mu, Jianying Yuan, Yujian Shi, Hao Zhang, Jun Gan, Nan Li, Xuesong Hu, Binghang Liu, Bicheng Yang, Wei Fan, Comparison of the two major classes of assembly algorithms: overlap–layout–consensus and de-Bruijn-graph, *Briefings in Functional Genomics*, Volume 11, Issue 1, January 2012, Pages 25–37, <https://doi.org/10.1093/bfgp/elr035>
- [7] Neri Van Otten, Node2Vec: Extensive Guide & How To In Python, <https://spotintelligence.com/2024/01/18/node2vec/>
- [8] Institute of Hydrobiology, Chinese Academy of Sciences (Institute of Hydrobiology, Chinese Academy of Science). DNA sequence of one male, one female, and 48 females. (SRR5119618)  
[https://www.ncbi.nlm.nih.gov/sra/SRX2425041\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX2425041[accn])
- [9] Kyungpook National University, Bacillus amyloliquefaciens strain: KNU-28 Genome sequencing  
[https://www.ncbi.nlm.nih.gov/sra/SRX16339329\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX16339329[accn])
- [10] Ben Langmead, “Overlap Layout Consensus assembly,” Langmead Lab, Accessed 20 Apr. 2024.  
[assembly\\_olc.pdf \(jhu.edu\)](#)

### Overlap Graph Github:

[workoholyguy/Sequence-Classification-Machine-Learning: Sequence Classification Algorithm \(github.com\)](https://github.com/workoholyguy/Sequence-Classification-Machine-Learning)