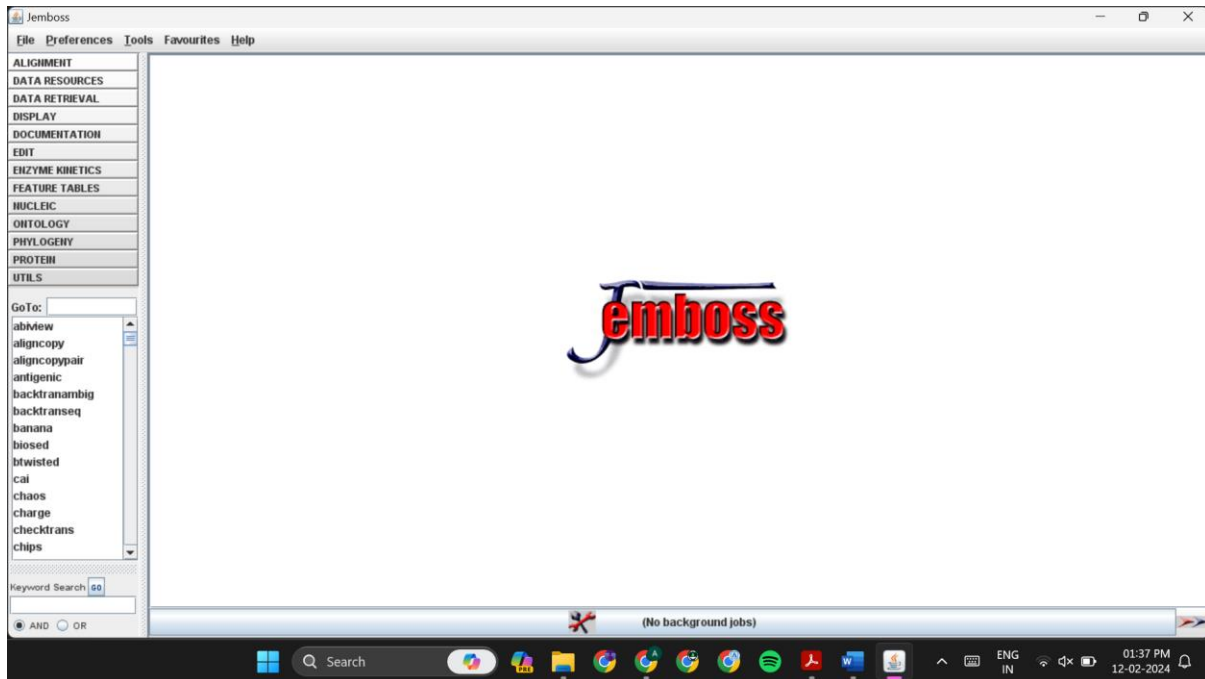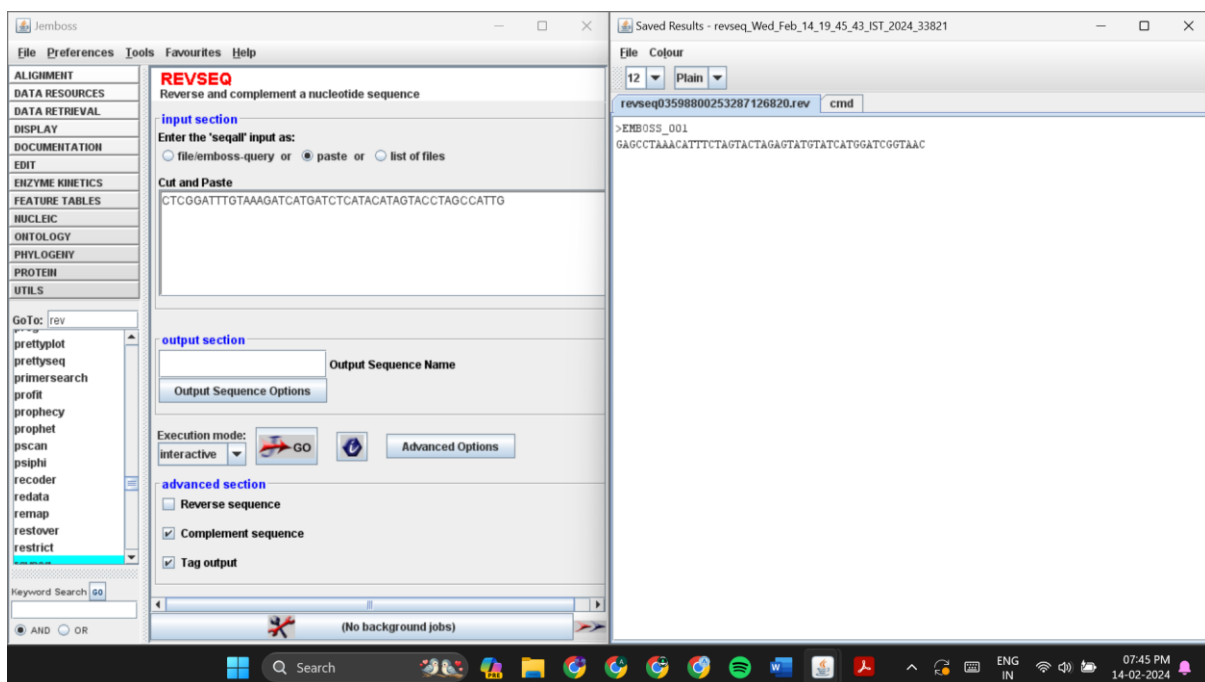Anirudh Rao

BE21B004

# BT3040 – Bioinformatics

## Practical 1
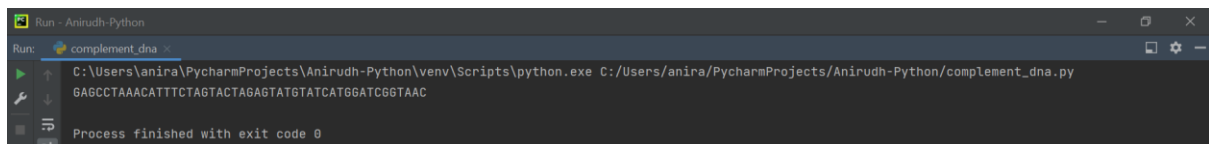
**1.**



**2.**

**3.**

The Python code to find the complementary strand for the DNA sequence given in Q2 can be found below:

```python
def complement(sequence):

    """
    Returns the complementary sequence (in 3' to 5' direction) of a given
DNA sequence (in 5' to 3' direction)

    """

    complement = {"A": "T", "T": "A", "G": "C", "C": "G"}
    complementary = ""
    for base in sequence:
        complementary += complement[base]
    return complementary


test_seq = "CTCGGATTTGTAAAGATCATGATCTCATACATAGTACCTAGCCATTG"
ans = complement(test_seq)
print(ans)
```
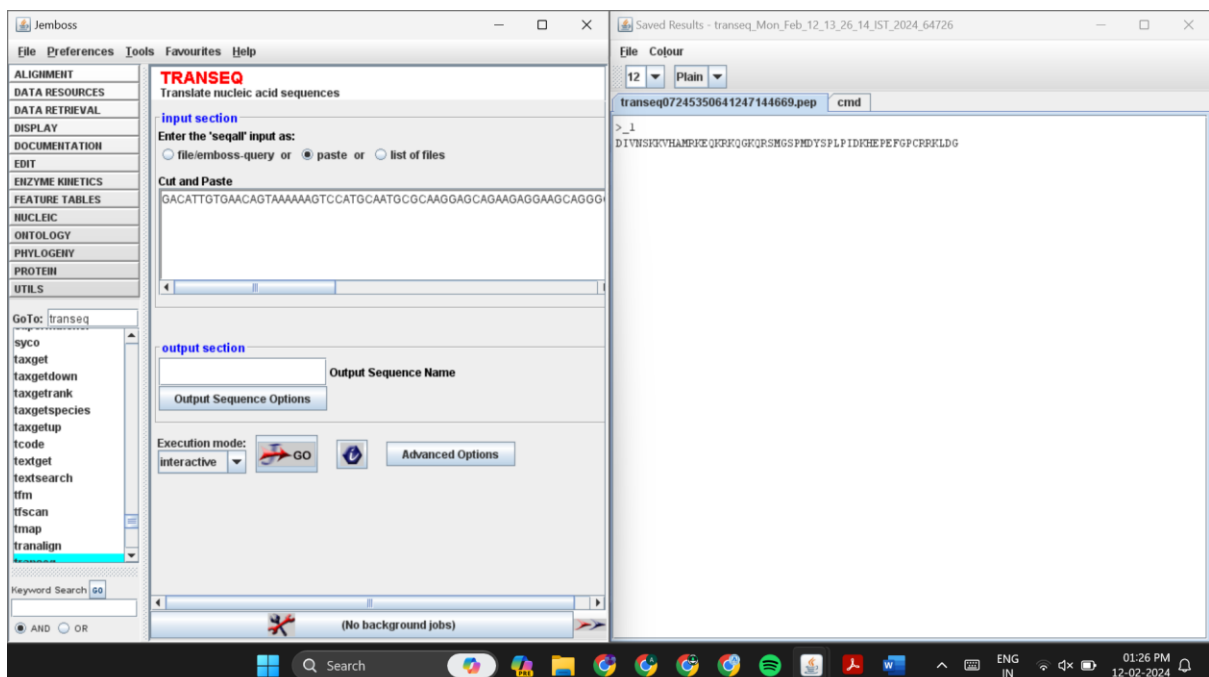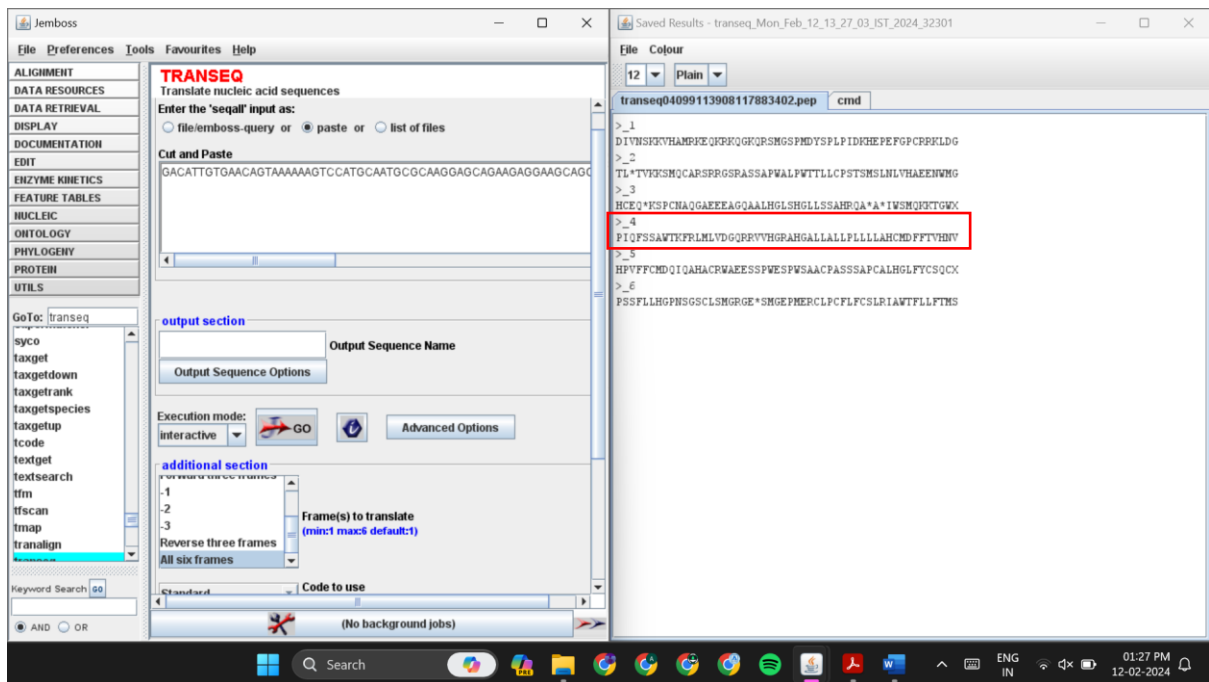
The output looks like:



Please note that the output is in the 3' to 5' direction.

**4. (i)**

## 4. (ii)



The required reading frame is **reading frame 4**, which is the same as the **reading frame –1**.

## 5.

The Python code to find the translated protein sequence for the DNA sequence given in Q4 can be found below:

```python
def translate(sequence):

    """
    Returns the translated protein sequence of a given DNA sequence (in 5'
to 3' direction)

    """

    codon_table = {
    "GCT": "A", "GCC": "A", "GCA": "A", "GCG": "A",
    "TGT": "C", "TGC": "C",
    "GAT": "D", "GAC": "D",
    "GAA": "E", "GAG": "E",
    "TTT": "F", "TTC": "F",
    "GGT": "G", "GGC": "G", "GGA": "G", "GGG": "G",
    "CAT": "H", "CAC": "H",
    "ATA": "I", "ATT": "I", "ATC": "I",
    "AAA": "K", "AAG": "K",
    "TTA": "L", "TTG": "L", "CTT": "L", "CTC": "L", "CTA": "L", "CTG": "L",
    "ATG": "M",
    "AAT": "N", "AAC": "N",
    "CCT": "P", "CCC": "P", "CCA": "P", "CCG": "P",
    "CAA": "Q", "CAG": "Q",
    "CGT": "R", "CGC": "R", "CGA": "R", "CGG": "R", "AGA": "R", "AGG": "R",
    "TCT": "S", "TCC": "S", "TCA": "S", "TCG": "S", "AGT": "S", "AGC": "S",
    "ACT": "T", "ACC": "T", "ACA": "T", "ACG": "T",
```

```
    "GTT": "V", "GTC": "V", "GTA": "V", "GTG": "V",
    "TGG": "W",
    "TAT": "Y", "TAC": "Y",
    "TAA": "*", "TAG": "*", "TGA": "*"}

    protein = ""
    for i in range(0, len(sequence) - 3 + 1, 3):
        codon = sequence[i:i+3]
        protein += codon_table[codon]

    return protein


test_seq = \
"GACATTGTGAACAGTAAAAAAGTCCATGCAATGCGCAAGGAGCAGAAGAGGAAGCAGGGCAAGCAGCGCTCCAT
GGGCTCTCCCATGGACTACTCTCCTCTGCCCATCGACAAGCATGAGCCTGAATTTGGTCCATGCAGAAGAAAACT
GGATGGG"
ans = translate(test_seq)
print(ans)
```
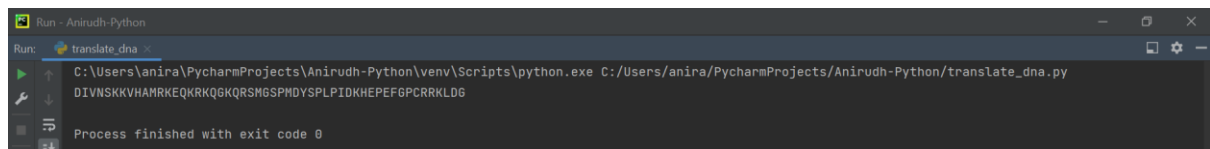
The output looks like:



## 6.

The Python code to find the strings 'AAG', 'GTC', 'GAG, 'ACTA', and 'ATAT' in the DNA sequence given in Q4 can be found below:

```
def find_match(sequence, string):

    match_count = 0
    match_positions = []
    for i in range(len(sequence) - len(string) + 1):
        if sequence[i:i+len(string)] == string:
            match_count += 1
            match_positions.append(str(i))
    positions = ", ".join(match_positions)
    print(f"Enter the string: {string}")
    print(f"Total matches: {match_count}")
    print(f"Positions of matches: {positions}\n")


test_seq = \
"GACATTGTGAACAGTAAAAAAGTCCATGCAATGCGCAAGGAGCAGAAGAGGAAGCAGGGCAAGCAGCGCTCCAT
GGGCTCTCCCATGGACTACTCTCCTCTGCCCATCGACAAGCATGAGCCTGAATTTGGTCCATGCAGAAGAAAACT
GGATGGG"
test_strings = ["AAG", "GTC", "GAG", "ACTA", "ATAT"]
for test_string in test_strings:
    find_match(test_seq, test_string)
```

The output looks like:



## 7.

DAN is used to calculate melting temperature while BANANA is used to plot the bending and curvature of DNA.

**8.**

The Python code to compute the average base stacking energy for the DNA sequence given in Q2 can be found below:

```python
def find_base_stack_energy(sequence):

    energies = {"AA": -4, "AT": -7, "AC": -5, "AG": -11, "TA": -7, "TT": -2, "TC": -3, "TG": -4, "CA": -9, "CT": -5, "CC": -6, "CG": -7, "GA": -9, "GT": -6, "GC": -4, "GG": -11}
    total_energy = 0
    for i in range(len(sequence)-2+1):
        dinucleotide = sequence[i:i+2]
        total_energy += energies[dinucleotide]
    avg_base_stack_energy = total_energy/(len(sequence)-2+1)
    return avg_base_stack_energy


test_seq = "CTCGGATTTGTAAAGATCATGATCTCATACATAGTACCTAGCCATTG"
print(find_base_stack_energy(test_seq))
```

The output looks like:

## 9. (i)



## 9. (ii)



The melting temperature of the 1st sequence is much lower compared to the 2nd sequence. This is because the 1st sequence is AT rich while the 2nd sequence is GC rich. GC rich sequences will have a higher melting temperature due to the greater energy required to break the three hydrogen bonds between G and C rather than the two hydrogen bonds between A and T.

**10.**

https://www.iitm.ac.in/bioinfo/cgi-bin/SBFE/seq_dna2.py

**Your input seq is:**

AAATGGCCCTAA

## Nucleotide Content:

|  | Nucleotide content in % |
|---|---|
| **AT_content** | 58.333333 |
| **Adenine_content** | 41.666667 |
| **Cytosine_content** | 25.000000 |
| **GC_content** | 41.666667 |
| **Guanine_content** | 16.666667 |
| **Keto_GT_content** | 33.333333 |
| **Purine_AG_content** | 58.333333 |
| **Thymine_content** | 16.666667 |

Download Now!