# IDC306 - Assignment 2

Rishi Vora                                                                                  MS21113

**Q1. Using the function, write the complement of the sequence in the fasta file.**

```Shell
#!/bin/bash

complement() {
    local seq="$1"
    local comp_seq=""
    for (( i=0; i<${#seq}; i++ )); do
        ch="${seq:$i:1}"
        case "$ch" in
            A) comp_seq+="T" ;;
            T) comp_seq+="A" ;;
            C) comp_seq+="G" ;;
            G) comp_seq+="C" ;;
            *) comp_seq+="$char" ;;
        esac
    done
    echo "$comp_seq"
}

input_file="./fasta_file.txt"
output_file="./fasta_file_complement.txt"

if ! [ -f $input_file ]; then
    echo "$input_file does not exist. Exiting..."
    exit 1
fi

>$output_file

while read -r line; do
    if [[ $line == ">"* ]]; then
        echo "$line" >> "$output_file"
    else
        complement "$line" >> "$output_file"
    fi
done < "$input_file"

echo "Complement sequence written to $output_file"
```

**Q2. Write a function to perform:**

**a. Find the composition of the DNA sequence.**

**b. Report the number of ORFs of length more than 20 codons in the (+) strand and (-) strand**

```bash
1    #!/bin/bash                                          Shell
2
3    composition () {
4        local seq=$1
5        seqlen=${#seq}
6
7        declare -A counts=([A]=0 [T]=0 [G]=0 [C]=0)
8
9        for i in `seq 0 $((seqlen-1))`; do
10           ch=${seq:$i:1}
11           ((counts[$ch]+=1))
12       done
13
14       echo -e "Composition of $seqid:\nA: ${counts[A]}, T: ${counts[T]}, G:
         ${counts[G]}, C: ${counts[C]}"
15   }
16
17   reverse () {
18       seq=$1
19       len=${#seq}
20       revd=""
21
22       for (( i=$len-1; i>=0; i-- )); do
23           revd=$revd${seq:$i:1}
24       done
25
26       echo $revd
27   }
28
29   complement () {
30       seq=$1
31       seqlen=${#seq}
32       comp=""
33
34       declare -A complements=([A]=T [T]=A [G]=C [C]=G)
35
36       for i in `seq 0 $((seqlen-1))`; do
37           ch=${seq:$i:1}
38           comp=$comp${complements[$ch]}
39       done
40
41       echo $comp
```

```bash
42  }
43
44  orf_finder () {
45      start_codons=("ATG" "GTG")
46      stop_codons=("TAG" "TAA" "TGA")
47
48      seq=$1
49      seqlen=${#seq}
50      direction=$2
51
52      for frame in 0 1 2; do
53          i=$frame
54
55          while [[ $i -lt $((seqlen-2)) ]]; do
56              start_codon="${seq:$i:3}"
57
58              if [[ ${start_codons[@]} =~ $start_codon ]]; then
59
60                  for (( j=$i; j<=$seqlen-2; j+=3)); do
61                      stop_codon="${seq:$j:3}"
62
63                      if [[ ${stop_codons[@]} =~ $stop_codon ]]; then
64                          orf="${seq:$i:$((j - i + 3))}"
65                          if [[ ${#orf} -le $((20*3)) ]]; then continue; fi
66
67                          echo -e "\nORF found in $seqid in frame $((frame+1))
                            from $((i+1)) to $((j+3)) in the ($direction) strand:
                            \n$orf"
68                          i=$j
69                          break
70                      fi
71                  done
72              fi
73              i=$((i+3))
74          done
75      done
76  }
77
78  file=./fasta_file.txt
79
80  if ! [ -f $file ]; then
81      echo "$file does not exist. Exiting..."
82      exit 1
83  fi
84
85  declare -A sequences
```

```
86
87  while read line; do
88      if [[ ${line:0:1} == ">" ]]; then
89          seqid=${line:1:14}
90      else
91          sequences[$seqid]=${sequences[$seqid]}$line
92      fi
93  done < $file
94
95  echo -e "There are ${#sequences[@]} sequences in $file\n"
96
97  for seqid in ${!sequences[@]}; do
98      echo -e "\n$seqid"
99      composition ${sequences[$seqid]}
100 done
101
102 for seqid in ${!sequences[@]}; do
103     sequence=${sequences[$seqid]}
104
105     orf_finder $sequence "+"
106
107     revd=$(reverse $sequence)
108     comp=$(complement $revd)
109     orf_finder $comp "-"
110 done
```