# Finding a Motif in DNA

Orr Shomroni

29 September, 2022

## Problem

Given two strings $s$ and $t$, $t$ is a **substring** of $s$ if $t$ is contained as a contiguous collection of symbols in $s$ (as a result, $t$ must be no longer than $s$).

The **position** of a symbol in a string is the total number of symbols found to its left, including itself (e.g., the positions of all occurrences of 'U' in "AUGCUUCAGAAAGGUCUUACG" are 2, 5, 6, 15, 17, and 18). The symbol at position $i$ of $s$ is denoted by $s[i]$.

A substring of $s$ can be represented as $s[j:k]$, where $j$ and $k$ represent the starting and ending positions of the substring in $s$; for example, if $s =$ "AUGCUUCAGAAAGGUCUUACG", then $s[2:5] =$ "UGCU".

The **location** of a substring $s[j:k]$ is its beginning **position** $j$; note that $t$ will have multiple locations in $s$ if it occurs more than once as a substring of $s$ (see the Sample below).

**Given**: Two DNA strings $s$ and $t$ (each of length at most 1 kbp).

**Return**: All locations of $t$ as a substring of $s$.

## Sample Dataset

```
GATATATGCATATACTT
```

```
ATAT
```

## Sample Output

```
2 4 10
```

```python
f=open("rosalind_subs.txt",'r')
s=f.readlines()[0].replace("\n","")
f=open("rosalind_subs.txt",'r')
t=f.readlines()[1].replace("\n","")
if len(s)<len(t):
  print("Error: length of substring ",t," is longer than ",s)
  exit()

motifLoci=[]
for i in range(0,len(s)-len(t)):
  tmpText=s[i:i+len(t)]
  if tmpText==t:
    motifLoci.append(i+1)

print("String: "+s)
```

## String: CATCTTGGCATCTTGGTTATGACTCATCTTGCATCTTGTGACATCATCTTGCCATCTTGGCATCTTGCATCTTGACATCTTGACATCTTGAGA

```
print("Substring: "+t)
```

## Substring: CATCTTGCA

```
print("Locations of motifs: "+str(motifLoci))
```

## Locations of motifs: [25, 61, 127, 159, 184, 191, 198, 205, 305, 312, 374, 381, 507, 514, 672, 687,