

# Programming Homework 1 Instructions (Read First)

In lecture and in a practical, we saw an implementation of the naive exact matching algorithm:

```
1 def naive(p, t):
2     occurrences = []
3     for i in range(len(t) - len(p) + 1): # loop over alignments
4         match = True
5         for j in range(len(p)): # loop over characters
6             if t[i+j] != p[j]: # compare characters
7                 match = False
8                 break
9         if match:
10            occurrences.append(i) # all chars matched; record
11    return occurrences
12
```

...and we saw a function that takes a DNA string and returns its reverse complement:

```
1 def reverseComplement(s):
2     complement = {'A': 'T', 'C': 'G', 'G': 'C', 'T': 'A', 'N': 'N'}
3     t = ''
4     for base in s:
5         t = complement[base] + t
6     return t
7
```

...and we saw a function that parses a DNA reference genome from a file in the FASTA format.

```
1 def readGenome(filename):
2     genome = ''
3     with open(filename, 'r') as f:
4         for line in f:
5             # ignore header line with genome information
6             if not line[0] == '>':
7                 genome += line.rstrip()
8     return genome
9
```

...and we saw a function that parses the read and quality strings from a FASTQ file containing sequencing reads.

```
1 def readFastq(filename):
2     sequences = []
3     qualities = []
4     with open(filename) as fh:
5         while True:
6             fh.readline() # skip name line
7             seq = fh.readline().rstrip() # read base sequence
8             fh.readline() # skip placeholder line
9             qual = fh.readline().rstrip() # base quality line
10            if len(seq) == 0:
11                break
12            sequences.append(seq)
13            qualities.append(qual)
14    return sequences, qualities
15
```

First, implement a version of the naive exact matching algorithm that is *strand-aware*. That is, instead of looking only for occurrences of P in T, additionally look for occurrences of the *reverse complement* of P in T. If P is ACT, your function should find occurrences of both ACT and its reverse complement AGT in T.

If P and its reverse complement are identical (e.g. AACGTT), then a given match offset should be reported only once. So if your new function is called naive\_with\_rc, then the old naive function and your new naive\_with\_rc function should return the same results when P equals its reverse complement.

Hint: See [this notebook](#) for a few examples you can use to test your naive\_with\_rc function.

Next, download and parse the lambda virus genome, at: [https://d28rh4a8wq0iu5.cloudfront.net/ads1/data/lambda\\_virus.fa](https://d28rh4a8wq0iu5.cloudfront.net/ads1/data/lambda_virus.fa)