# Project 1: Fun with DNA (REGEX Lookaround)!

In this project we find Opening reading frame or ORF from DNA sequences with the help of Python regex.

**DNA** is a sequence of bases, `A`, `C`, `G`, or `T`. They are translated into proteins 3-bases where each sequence is called a **codon**. There is a special start codon `ATG`, and three stop codons, `TGA`, `TAG`, and `TAA`. Example:

```
cgcgcATGcATGcgTGAcTAAcgTAGcgcgcgcgc
```

An opening reading frame or **ORF** consists of a **start codon**, followed by some more codons, and ending with a **stop codon**. The above example has overlapping ORFs.

- `ATGcATGcgTGA` and
- `ATGcgTGAcTAA`.

The following pattern only finds the first ORF (`atgcatgcgtga'`). Since it consumes the first ORF, it also consumes the beginning of the second ORF.

```
1   from re import *
2
3   dna = 'cgcgcATGcATGcgTGAcTAAcg
4   dna = dna.lower()
5   orfpat = r'(?x) ( atg  (?: (?!
6   print findall(orfpat,dna)
```

We want to find an ORF without consuming it, we can use a **positive lookahead** assertion (`(?= ( atg`). We put the whole ORF pattern inside the lookahead and find the two `atgcatgcgtga` and `atgcgtgactaa`.

```
1   from re import *
2
3   dna = 'cgcgcATGcATGcgTGAcTAAcg
```

```
4    dna = dna.lower()
5    orfpat = r'(?x) (?= ( atg  (?:
6    s = findall(orfpat,dna)
7    if s:
8        print ', '.join(s)
```

This project **adopts** and **simplifies** the Splitsvile examples (DNA) from Rex Dwyer's ipython notebook.