

P2: Parsing

Write a python script that parses a GenBank file and outputs a FASTA file and an ordered table with some statistics (specifications below).

Input: a file containing multiple sequences in GenBank format. Location: <http://www.bioinformatics.nl/courses/BIF-30806/docs/argonaut.gb>

Tasks:

- Parse the GenBank file, collect accession numbers, organism names, and sequences
- Calculate length and GC content for each sequence
- Order the sequences from high to low GC content
- Produce two output files (specified below)

Output (2 files):

- A. FASTA file with the sequences, ordered from high to low GC content. Labels should be accession numbered followed by organism name.
- B. A tab-delimited file with the following columns:
 1. Accession number
 2. Organism name
 3. GC content (printed as percentage with two decimals)
 4. Sequence length

The lines should be ordered from high to low GC content (same order as in output A)

Create a python script, containing your **name** and **student number**, that performs the described task. **Turn in your python script on BlackBoard (under P2).**

Example output 1 (the data in this example is made up):

```
>NM_002022392 Solanum lycopersicum
ATGTCGTATAAACCAAGCTCAGAAATAGCTTTCCGGTTATGGAGGGTTGG...
>XM_001635194 Chlamydomonas moewusii
CTTAATTACATATTAATGTTCTGTACCAAGCGGCTTGTGCGGGCA...
...
```

Example output 2 (the data in this example is made up):

```
NM_002022392    Solanum lycopersicum    56.53    2579
XM_001635194    Chlamydomonas moewusii    40.50    2681
...
```