# BIOL550 - Assignment #2 (5 questions, 45 points)

**Always read through assignment/project/exam instructions carefully; they are provided for a reason, and a lot of time can be saved by taking an extra second to read through them!**

In this assignment, you will be provided with five (5) different biology-related files. For each file, you will be tasked with designing a regular expression that returns desired lines. Please feel free to utilize Regexer (https://regexr.com/) to aid in the design of your regular expressions. A 'desired output' file containing the correctly extracted lines is provided for each input file.

To check the validity and correctness of your regular expression, you can perform the following command lines:

```
## This extracts the lines that match your regular expression
grep -P 'regular expression' Q#_input.txt > Q#_output.txt
```

followed by:

```
## This checks if the files have any differences
diff -s Q#_output.txt Q#_desired_output.txt
```

***What to do*:**

1. Download the A2_F23_files.tar.gz from Blackboard to your local computer.
2. Upload A2_F23_files.tar.gz from your local computer to your ~/Assignment/Assignment_2 folder on the class server
3. Untar the A2_F23_files.tar.gz file archive and work from your account on the class server.
4. Write your final regular expression for each question in a single text file named assignment2_[lastname]_[firstname].txt. **Do not** use a word processor to save your regular expression, use a code editor instead, to prevent corrupting your regular expressions!
5. Submit the text file containing your regular expression to Blackboard.

   ```
   ## Example of format
   # Question 1
   ^\d+\w{3}\s+$
   # Question 2
   \d\d\d\w+\s*
   ```
   **Please reach out for any clarification needed.**

## Question 1 (5 points)

FASTA (for FAST-All), has become one of the standards in which sequence information, either nucleotide or amino-acid, is stored. The format is made up of two components: the header and the sequence. The header component is a single line that starts with a '>' and is followed by a string of characters containing the sequence ID and potentially some other metadata. The sequence component has variable number of lines (it depends on the length of the sequence and text-wrap formatting used).
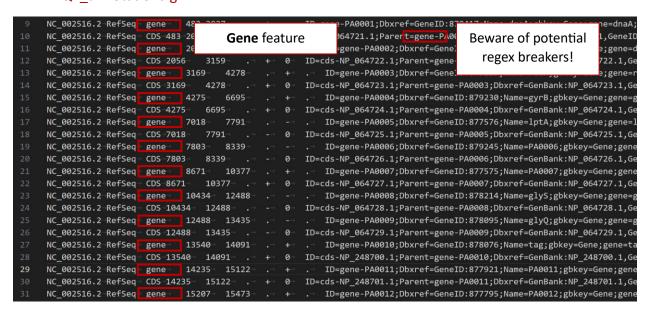
1. Create a regular expression that returns only the **FASTA header** lines of the file Q1_fasta.fna.

## Question 2 (5 points)

Genome annotations can be stored in various file formats; for this question, the annotations are in gene feature format (GFF). The GFF format stores gene location, source database, feature type, as well as several additional pieces of information needed to completely define genomic annotations, with one feature entry per line.
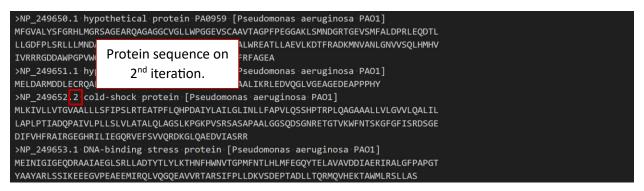
1. Create a regular expression that returns only lines with the **gene feature-type** from Q2_annotations.gff.

## Question 3 (10 points)

Sometimes newer evidence comes forward that may result in the need to update the sequences present in FASTA files. When these changes are made, the version number of the protein sequence is modified to delineate past analysis from future analysis.

1. Create a regular expression the grabs FASTA headers for proteins whose sequence is not on its first iteration (**iteration ≥ 2**) from Q3_protein.faa.

```
>NP_249650.1 hypothetical protein PA0959 [Pseudomonas aeruginosa PAO1]
MFGVALYSFGRHLMGRSAGEARQAGAGGCVGLLWPGGEVSCAAVTAGPFPEGGAKLSMNDGRTGEVSMFALDPRLEQDTL
LLGDFPLSRLLLMND            ALWREATLLAEVLKDTFRADKMNVANLGNVVSQLHMHV
IVRRRGDDAWPGPVW    Protein sequence on    FRFAGEA
>NP_249651.1 hy       2nd iteration.       Pseudomonas aeruginosa PAO1]
MELDARMDDLECRQA            AALIKRLEDVQGLVGEAGEDEAPPPHY
>NP_249652.2 cold-shock protein [Pseudomonas aeruginosa PAO1]
MLKIVLLVTGVAALLLSFIPSLRTEATPFLQHPDAIYLAILGLINLLFAPVLQSSHPTRPLQAGAAALLVLGVVLQALIL
LAPLPTIADQPAIVLPLLSLVLATALQLAGSLKPGKPVSRSASAPAALGGSQDSGNRETGTVKWFNTSKGFGFISRDSGE
DIFVHFRAIRGEGHRILIEGQRVEFSVVQRDKGLQAEDVIASRR
>NP_249653.1 DNA-binding stress protein [Pseudomonas aeruginosa PAO1]
MEINIGIGEQDRAAIAEGLSRLLADTYTLYLKTHNFHWNVTGPMFNTLHLMFEGQYTELAVAVDDIAERIRALGFPAPGT
YAAYARLSSIKEEEGVPEAEEMIRQLVQGQEAVVRTARSIFPLLDKVSDEPTADLLTQRMQVHEKTAWMLRSLLAS
```

## Question 4 (10 points)

A dump file (.dmp), named because they are typically created by "dumping" a bunch of data into a file, can be used to store a variety of information. It's important to note that just because the data is "dumped" does not mean there isn't a defined structure to the file! In the NCBI names.dmp, the data stored is used to connect the taxonomic ID of an organism to its more descriptive name.

1. Create a regular expression that returns all lines that have a **scientific name** related to **Arabidopsis** in Q4_names.dmp.

## Question 5 (15 points)

The PDB format (.pdb) contains information related to the orientation of biological molecules in 3D space, most commonly proteins (the P in PDB stands for proteins!). In its barebones form, the PDB file contains 3D coordinates for almost each atom in the structure. Additionally, PDB files will almost always contain a metadata block at the top of the file, describing what the molecules in the structure are, how they were obtained, and the individuals who helped obtain it. When it comes to text parsing, the PDB format is one of the most interesting, as the file does not delimit data by the traditional method (tab, comma, semicolon), rather each column is assigned to a particular piece of data depending on what kind of line is being read (For more details: https://www.cgl.ucsf.edu/chimera/docs/UsersGuide/tutorials/pdbintro.html)!

1. Grab all lines that denote a protein (ATOM) alpha carbon (CA) in a tyrosine residue (TYR) for chains B, D, or E in Q5_structure.pdb.



**Alpha carbon** entries.

**Chain** entries.

**Atom** entries.

**TYR residue** entries.