**ASSIGNMENT #2 – Regular Expressions – BIOL550 (5 Questions; 35 PTS)**

*Regular expressions* (*aka* regexes) are motifs/patterns designed to match specific strings of characters in text files. These matches can then be captured or excluded/ignored depending on the desired objective.

Regular expressions can vary in complexity, from very simple to highly complex. When designing a regular expression, I recommend starting with something broad and simple that matches all the characters that you are looking for. This often leads to capturing extra characters and/or patterns that were not desired. If so, you can modify/complexify your expression as needed make it more specific and remove undesired matches.

Files commonly used in bioinformatics can be found as binary files or as human-readable ASCII text files. Text files in bioinformatics include formats commonly found in programming (*e.g.* .json, .xml, .tsv) and custom formats designed specifically to store biological data (*e.g.* .fastq, .gbk, .vcf). These texts files can be further designed for human readability or to facilitate parsing with regular expressions and/or scripting/programming languages.

In this assignment, we will practice how to use regular expressions by parsing several text files often encountered in bioinformatics. We'll learn more about these text files in the second half of the semester. ***Feel free to use online regular expression tools like Regexr (https://regexr.com/) to help you design your regular expressions.***

The input files for this assignment are available Blackboard as A2_F22_files.tar.gz. Note that while we will be using files that are relevant to bioinformatics in the questions below, the goal here is to test your understanding of regular expressions, not biology/bioinformatics. ***If you are not in biology or lack the relevant background pertaining to these files, feel free to ask questions about them.***

*What to do:*

      a)   Download the file A2_F22_files.tar.gz from Blackboard to your local computer

      b)   In your account on the class server, create a folder named: ~/Assignment/Assignment_2

      c)   Upload A2_F22_files.tar.gz from your local computer to your ~/Assignment/Assignment_2 folder on the class server

      d)   Untar the A2_F22_files.tar.gz file archive and work from your account on the class server.

      e)   You can test your Perl Compliant Regular Expressions (PCRE) directly on the input files with grep -P ## See question 1 below for an example of how to do so.

      f)   Write the regular expressions you designed for each question in a single text file that includes your name (*e.g. **Pombert_JF_regexes.txt***). Make sure to use regular expressions, *i.e.* **DO NOT** use the output examples as input files for grep. Use a code editor to save your regexes; Microsoft Word will likely corrupt your regexes!

      g)   Submit the text file containing your regular expressions to Blackboard.
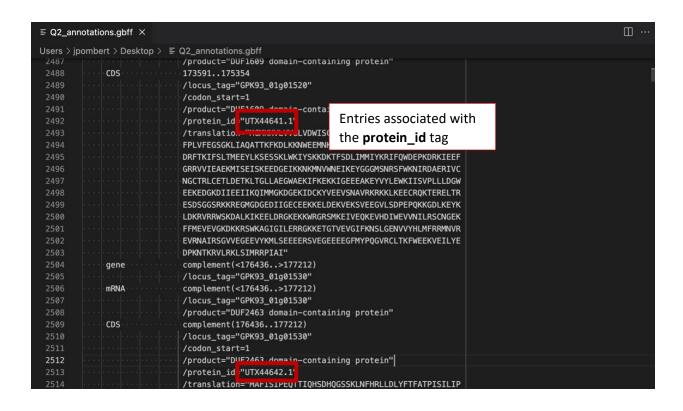
**7 PTS. QUESTION 1**   ## FASTA is a standard file format for biological sequences, it can contain one
## or more sequences. Each sequence is identified by a header starting with >.
## This header is followed by one or more line(s) containing the corresponding
## sequence: https://www.ncbi.nlm.nih.gov/genbank/fastaformat/

1. Create a regular expression that will return only the **sequence lines**; *i.e.* lines that do not matches the fasta headers from the file Q1_multi.fasta.
2. Test it with grep on Q1_multi.fasta. Write the output in Q1_myHeaders.txt; e.g.
   grep -P 'regular expression' Q1_multi.fasta > Q1_myHeaders.txt
3. You can compare your output to that of Q1_desired_output.txt with **wc -l** and **diff**:
   wc -l Q1_myHeaders.txt Desired_outputs/Q1_desired_output.txt
   diff -s Q1_myHeaders.txt Desired_outputs/Q1_desired_output.txt
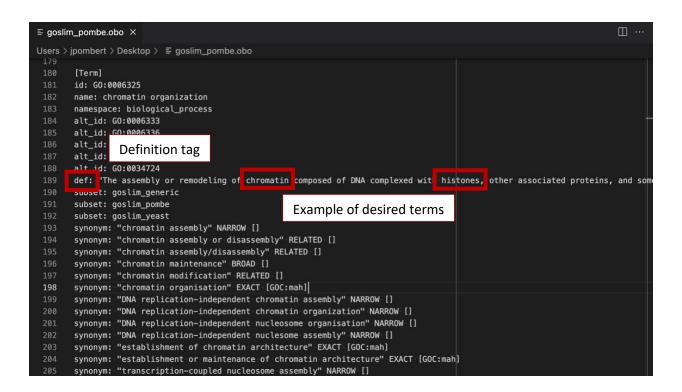
**7 PTS. QUESTION 2**

## In genomics, several formats are available to store biological annotations.
## These annotations include where genes are located and information about
## the products they encode (if known). Some of the most common genome
## annotation formats include those from the NCBI GenBank (.gb/.gbff) and
## EMBL (.embl) databases as well as the Genome File Format (.gff/.gff3)
## version 3

1. Create a regular expression that matches all entries associated with the **protein_id** tags from the file Q2_annotations.gbff and returns the corresponding lines with grep (see Q2_desired_output.txt for example).
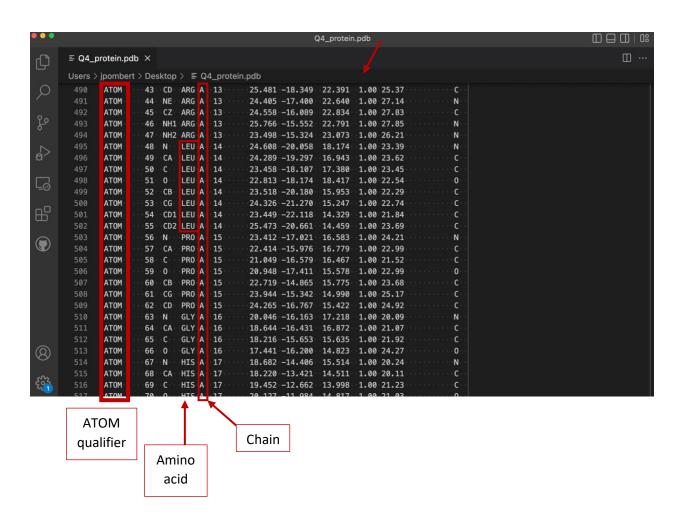2. Test it with grep on Q2_annotations.gbff. Write the output in Q2_myProtein_IDs.txt.

**7 PTS. QUESTION 3**

1. Create a regular expression that matches all lines with the **def:** tag and which contains at least one of the following terms: **chromatin**, **telomere**, or **histone** from Q3_goslim_pombe.obo (see Q3_desired_output.txt for desired output).
   ## pombe = *Schizosaccharomyces pombe* (fission yeast); an important model organism in
   ## biology
2. Test your regular expression with grep on Q3_ goslim_pombe.obo. Write the output in Q3_myGOs.txt.

## Protein 3D structures are often stored using the Protein Data Bank (PDB)
## eponym PDB format (https://www.rcsb.org/). This legacy format was later
## expanded into the PDBx/mmCIF (macromolecular Crystallographic
## Information File) format to accommodate larger molecules. Both file formats
## can be visualized using various protein 3D structure viewers like ChimeraX
## (https://www.cgl.ucsf.edu/chimerax/)

1. Create a regular expression that matches only lines starting with the **ATOM** qualifier and which contain the leucine amino acid (**LEU**) located on chain **A** in Q4_protein.pdb. See lines from Q4_desired_output.txt for the desired outcome.
2. Test your regular expression with grep on Q4_protein.pdb. Write the output in Q4_myLeucine_residues.txt.

**7 PTS. QUESTION 5**

1. Create a regular expression that matches all lines with **homozygous transitions (HOM=1)** between **pyrimidines (C -> T, T -> C)** from Q5_Streptococcus.vcf (see Q5_desired_output.txt for desired output). ## In this exercise, a dataset from a reference *Streptococcus* genome was mapped against a new clinical isolate...

2. Test it with grep and write the output in Q5_mySNPs.txt.