# Obesity Induced Changes in Gene Expression

This project will focus on analyzing simulated RNAseq data in order to evaluate differential gene expression between control and treatment samples. Your goal is to combine your knowledge of Unix, Python/R, and bioinformatic tools *muscle* and *hmmer* to "replicate" a recent study of gene expression in kidney tissues of normal (control) and obese mice (Kuhns & Pluznick 2017). RNAseq involves high-througput sequencing of cDNA libraries constructed from all expressed RNA in the target tissue. Sequences are filtered and aligned to a transcriptome or genome assembly with established bioinformatic pipelines, including tools such as *Bowtie* or *Tophat*. Expression levels of each transcript in each tissue sample are then quantified/converted to counts using statistical models implemented in programs such as *HTSeq* and differential expression assessed with various programs/packages (e.g. *edgeR* or *DESeq*). In this project, you will approximate this process by building protein models of 6 target transcripts, searching simulated transcriptome data for "hits," and counting those hits as a proxy for gene expression.

Begin by inspecting your files, you should have six to work with:

1.*uniquetranscripts.fasta*: A fasta file including the sequences of 6 transcripts that showed differential expression between kidney tissues of normal (control) and obese mice in Kuhns & Pluznick 2017.

2.*codonmap.txt*: A tab delimited text file with amino acid abbreviation in the first column and associated three base pair codons in the second.

3&4.*control1.fasta* & *control2.fasta*: Transcript sequences produced in a "new *RNAseq* experiment" with kidney tissue from two normal (control) male mice.

5&6.*obese1.fasta* & *obese2.fasta*: Transcript sequences produced in a "new *RNAseq* experiment" with kidney tissue from two diet-induced obese male mice.

## Use BLAST to identify the genes encoded by the 6 differentially expressed transcripts listed in *uniquetranscripts.fasta.*

Go to https://blast.ncbi.nlm.nih.gov/Blast.cgi and select *nucleotide BLAST*, because your data are transcriptome (cDNA) sequences. Upload *uniquetranscripts.fasta* via the *choose file* button shown below.



For this assignment, search the *Nucleotide collection (nr/nt)* Database.

Also, use the default optimization for *Highly similar sequences (megablast)*. And click *BLAST* to initiate the search.



When the searches are complete, your browser will redirect to the results page. The drop down menu next to *Results for:* will allow you to toggle between results for the 6 transcripts in your fasta file.



Scroll down to the *Descriptions* heading for your table of hits, which you should **save to your repo**. The best hits, listed first, will have extremely low E-values and high (~100%) Identity. By selecting *All* (top left of the table) the *Download* button becomes available and you can choose the download format from the dropdown menu. **Using Unix commands, make a single table that includes the top hit for each transcript.**
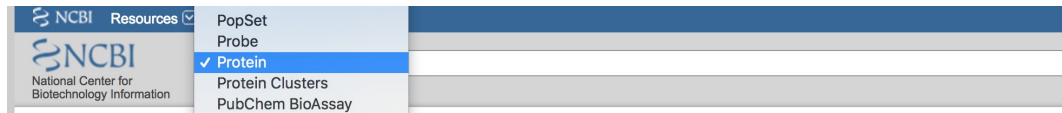
# Search the NCBI protein database for amino acid sequences corresponding to these 6 transcripts.

Point your browser to the NCBI homepage https://www.ncbi.nlm.nih.gov/ and choose the *Protein* database.



Based on the name of the protein listed in the top hit for each transcript in the previous step, search the NCBI protein data base for sequences. By selecting *Animals* under *Species* at the top left of the page, the option to filter by taxon should appear at the top right.



Depending on how conserved a protein is, the protein sequence can vary substantially between distantly related organisms (view taxa as a *TREE* to display relationships). This could change the efficacy of the HMM model you will build below. Choose 10-20 protein sequences from mice (*Mus musculus*) and other closely related organisms (if available) by checking the boxes to the left of the listed matches. Download your selected sequences by clicking on the arrow next to *Send to:*, in the upper right corner, selecting *File* and choosing *FASTA* from the dropdown menu. **You should save one fasta file of protein sequences per identified transcript (6 total).**

## Translate the 4 provided files of "RNAseq data."

**Write a python or R script to translate nucleotide sequences into amino acid sequences.** Your script should read the codon map (*codonmap.txt*) and the nucleotide fasta file you want to translate and write a corresponding fasta file containing the translated amino acid sequences. (If you want to try out something new, look up how to pass arguments to a python or R script.) I won't give too many suggestions here, because we want you to decompose the problem as a group. But, here are a few hints for finding the first open reading frame in each sequence. All of the correct start codons are the first in frame, so you could count nucleotides by 3's from the beginning of the sequence or you could use a regular expression. However, beware of greedy regex!! You'll want to stop at the first stop codon in frame (not the last!). **Apply this script to the 4 files of "RNAseq data" (*control1.fasta, control2.fasta, obese1.fasta & obese2.fasta*).**

## Build a Hidden Markov Model for each of the 6 transcript proteins and search the 4 translated "RNAseq files."

For each of the 6 transcripts make a muscle alignment from your downloaded protein sequences. Using these alignments, construct 6 HMM protein models using hmmbuild. Finally, search all 4 translated "RNAseq files" for each of the 6 HMM protein models using hmmsearch. **Use a bash script to loop over the 6 transcript files and 4 "RNAseq files," executing these commands.**

## Graph the "expression levels" of each protein in each of the "RNAseq files."

Based on the counts of hmm hits for each transcript (our measure of RNA expression) in each "RNAseq file," make a **graphical comparison of expression levels** across the 2 normal and 2 obese mice. **Qualitatively compare these results to those reported in Kuhns & Pluznick 2017.**

## Further Exploration

1. For a few (2-3) of your 6 transcripts, return to the original *BLAST* search and change the *Optimize for* option. Qualitatively, how do *discontinuous megablast* and *blastn* change your table of *BLAST* hits? This may be easier to explore, if you also restrict the *Database* option to either *Human* or *Mouse*. (You do not need to repeat any other steps here.)

2. For a few (2-3) of your 6 transcripts, return to the NCBI protein search and explore the effects of phylogenetic relatedness of your amino acid sequences on the performance of your HMM model. For example, what would happen if you built your HMM protein model using more distantly related mammals (e.g. primates)? Would you still get the same quality hits if your HMM protein model was based on non-mammalian sequences? Pick one of the "RNAseq files" to search in order to test your hypotheses. Compare e-values among HMMs built from differing taxa.