# SCHOOL OF COMPUTING
## Faculty of Engineering

## Semester II 2021/2022

**Subject**      : Bioinformatics I (SECB2103)
**Section**      : 01 – Dr Haslina Hashim
**Topic**        : Lab 02 – Access to Biological Databases

<u>**Name**</u>          : Gui Yu Xuan
<u>**Matric ID**</u>     : A20EC0039

# 1. The purpose of this problem is to introduce you to using Entrez and related NCBI resources.

a) How many human proteins are bigger than 300,000 daltons?
<u>There are 2744 human proteins are bigger than 300,000 daltons.</u>

b) What is the longest human protein?
<u>Titin is the longest human protein.</u>

There are several different ways to solve these questions.

(1) Try to first limit your search to human by using TaxBrowser. From the home page of NCBI select the alphabetical list of resources and find the Taxonomy Browser and the entry for human. Then follow the link to Entrez Protein, where all the results will be limited to human.

(2) Enter a command in the format xxxxxx:yyyyyy[molwt] to restrict the output to a certain number of Daltons; for example, 002000:010000[molwt] will select proteins of molecular weight 2,000 to 10,000.

(3) As a different approach, search 30000:50000[Sequence Length]

(4) You can read more about titin (NP_596869.4), the longest human protein, in an NCBI newsletter (WebLink 2.73 ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt ). While the average protein has a length of several hundred amino acids, incredibly titin is 34,423 amino acids in length.

(5) Explore additional ways to limit Entrez searches by using an NCBI Handbook chapter (http://www.ncbi.nlm.nih.gov/books/NBK44864/)

*Write a report of how you would answer the question in this week's session.*

*Include screen captures as well as detailed explanation and observations.*

*All screen captures should be labelled as figures and numbered accordingly.*

*All tables should be labelled as tables and numbered accordingly.*

Hint: What is the Latin name for humans? Homo sapiens.
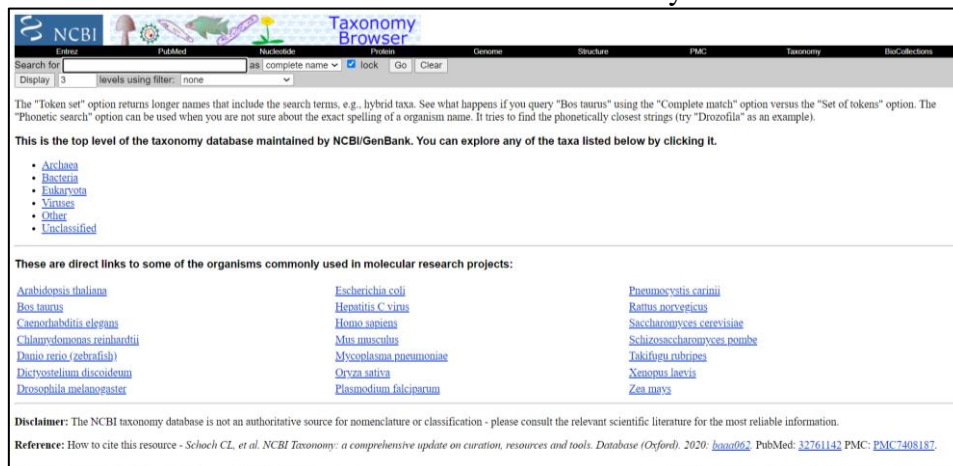
Here is the screen shot of the Taxonomy Browser.



*Figure 1*

After typing "human" in the searching bar, then select the "home sapiens(human).



*Figure 2*

At the right side of the website, click the Subtree links that under the Entrez Records.



*Figure 3*

At the left side of the website, select the molecular weight and enter the range from 300,000 to 6,000,000.
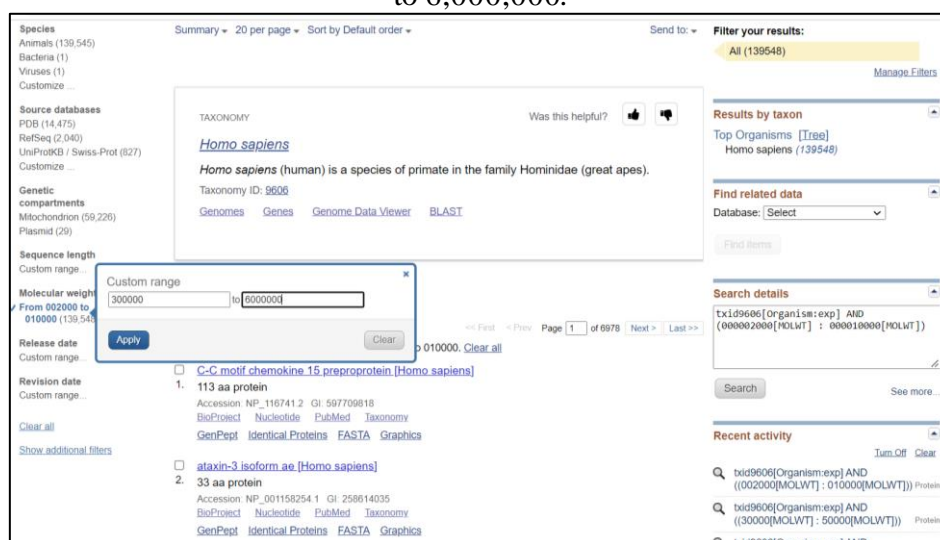


*Figure 4*

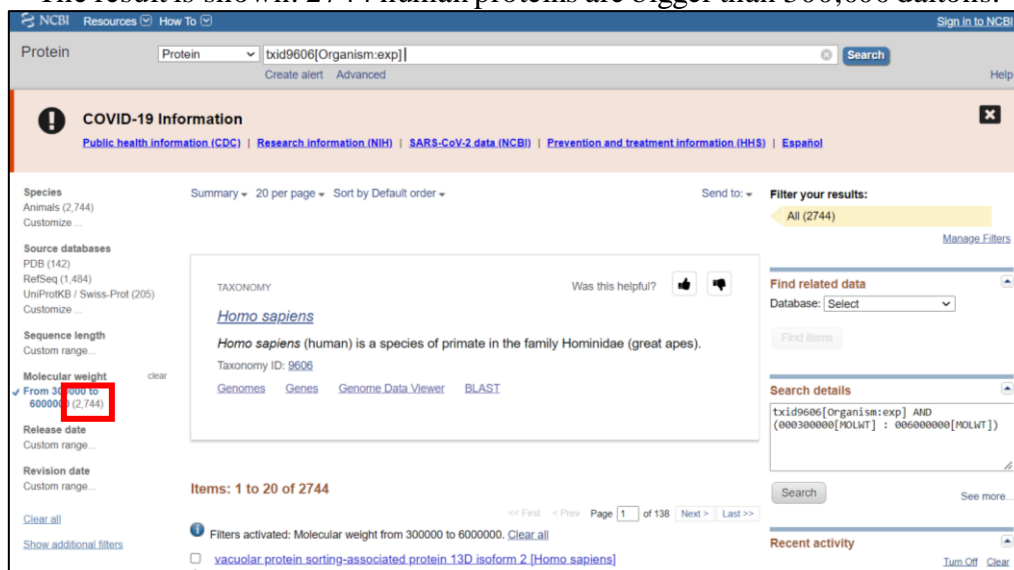The result is shown. 2744 human proteins are bigger than 300,000 daltons.



*Figure 5*

To see the largest protein, the molecular weight needs to change from 3,900,000 to 6,000,000. Then, you may see the result where titin is the longest human protein.
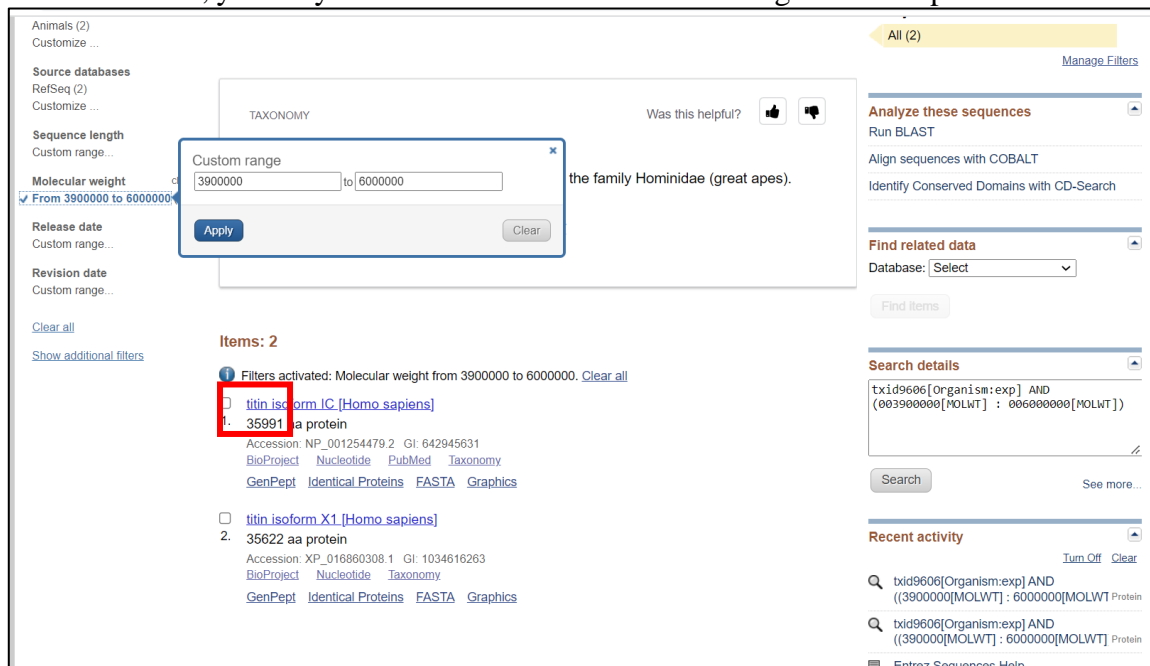


*Figure 6*

# 2. The purpose of this problem is to obtain information from the NCBI website.

The RefSeq accession number of human beta globin protein is NP_000509. Go to NCBI

(http://www.ncbi.nlm.nih.gov/).

a) What is the RefSeq accession number of beta globin protein from the chimpanzee (Pan troglodytes)?
XP_508242.1.

(1) There are several different ways to solve this. Try typing chimpanzee globin into the home page of NCBI; or use the Taxonomy Browser to find chimpanzee Entrez Gene entries.

(2) HomoloGene (http://www.ncbi.nlm.nih.gov/homologene )(WebLink 2.38) is a great resource to learn about sets of related eukaryotic proteins. Use HomoloGene to find a set of beta globins including chimpanzee.

4

Entering chimpanzee globin into the search bar of the NCBI home page. Then the results in both Gene and HomoloGene will be seen.
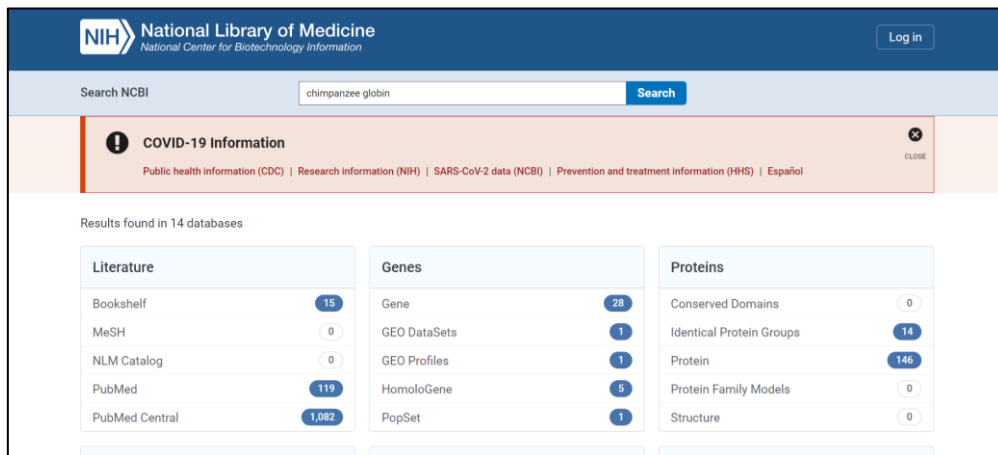


*Figure 7*

Then, by clicking the Gene, the gene results sort the HBB to the top. You may see the RefSeq accession NC_036890., this corresponds to the entire chromosome 11 of the chimp.
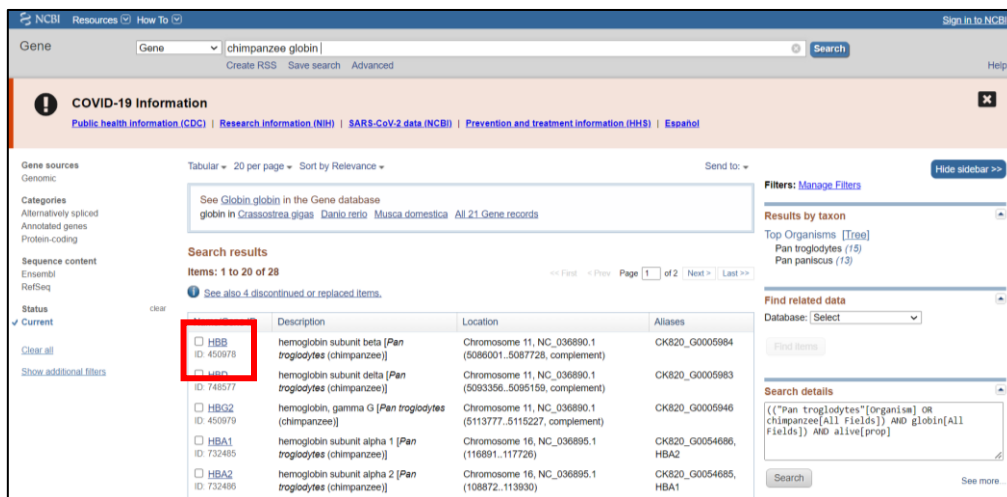


*Figure 8*

By following the HomoloGene, you may select the hemoglobin beta, and it will lead you to the accession for the chimp HBB protein.
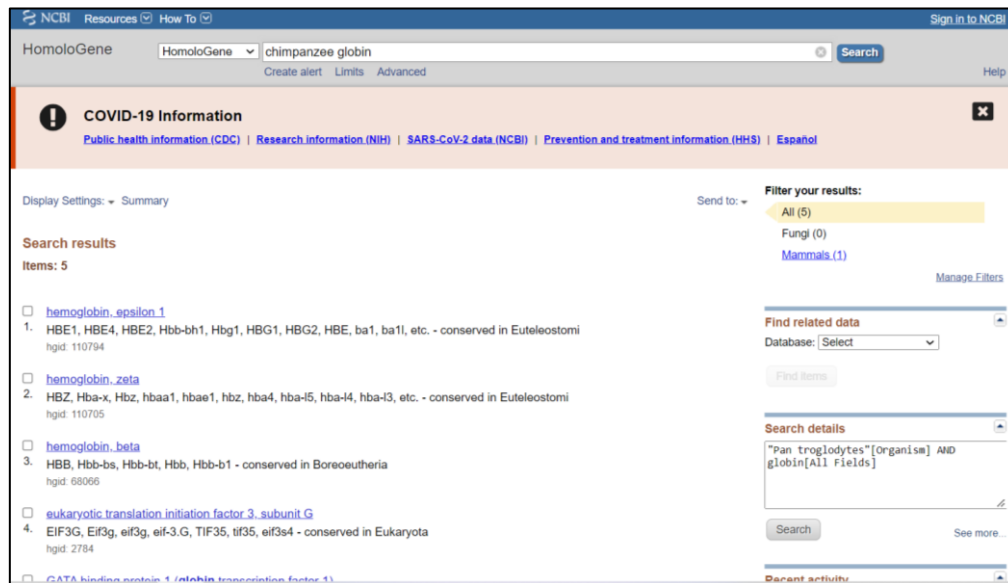


*Figure 9*

Hence, you may see the result of the RefSeq accession number of beta globin protein from the chimpanzee (Pan troglodytes).



*Figure 10*

## 3. The purpose of this exercise is to become familiar with the EBI website and how to use it to access information.

(1) Visit the site (http://www.ebi.ac.uk/ )(WebLink 2.5). Enter hemoglobin beta in the main query box (alternatively use the query human hemoglobin beta).

(2) Inspect the results. Explore the various links to information about pathways, genomes, nucleotide and protein sequences, structures, protein families, and more.

Here is the EBI home page. Typing the hemoglobin beta at the search bar.



*Figure 11*

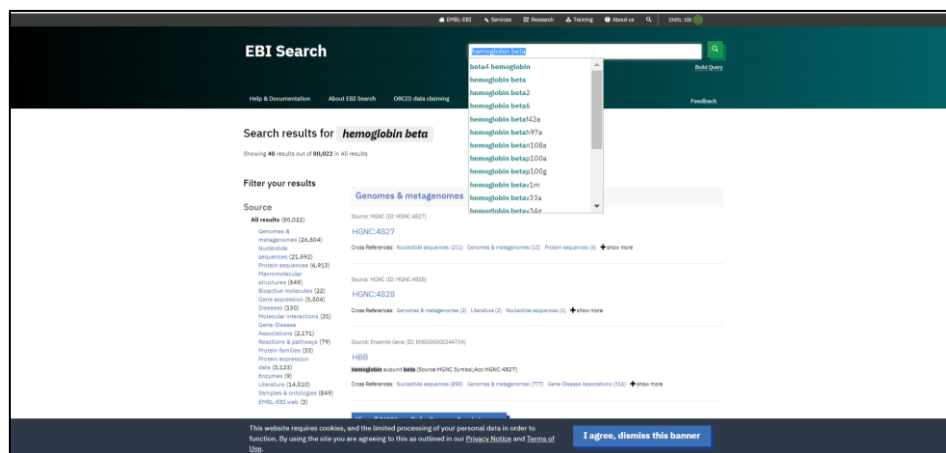Here is the search results of the hemoglobin beta.



*Figure 12*

## 4. Accessing information from BioMart: the beta globin locus

(1) Go to www.ensembl.org and follow the link to BioMart.

(2) First choose a database; we will select Ensembl Genes 71.

(3) Choose a dataset: Homo sapiens genes (GRCh37.p10). Note the other available datasets.

(4) Choose a filter. Here the options include region, gene, transcript event, expression, multispecies comparisons, protein domains, and variation. Select "region", chromosome 11, and enter 5240000 for the Gene Start (bp) and 5300000 (bp) for the Gene End. (Note that this region spans 60 kilobases and corresponds to chr11:5,240,001-5,300,000.)

(5) Choose attributes. Select the following features. Under "Gene" select Ensembl Gene ID and %GC content; under "External" select the external references CCDS ID, HGNC symbol (this is the official gene symbol) and HGNC ID(s).

(6) At the top left select "Count." Currently there are 8 genes matching these criteria.

(7) To view these results press "Results. "Note that you can export your results in several formats (including a comma separated values or CSV file) that can be further manipulated (e.g., converted to a BED file).

Here is the Ensembl home page. At the home page, select the BioMart.



*Figure 13*

Then, choose a Ensembl Genes 105 database.



*Figure 14*

Then, select a Human genes (GRCh38.p13) dataset.



*Figure 15*

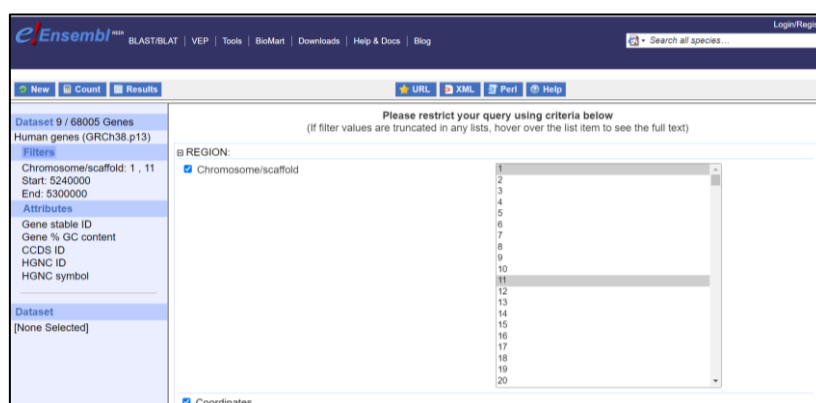After that, select the filters. Choose a region and fill the gene start and gene end.



*Figure 16*

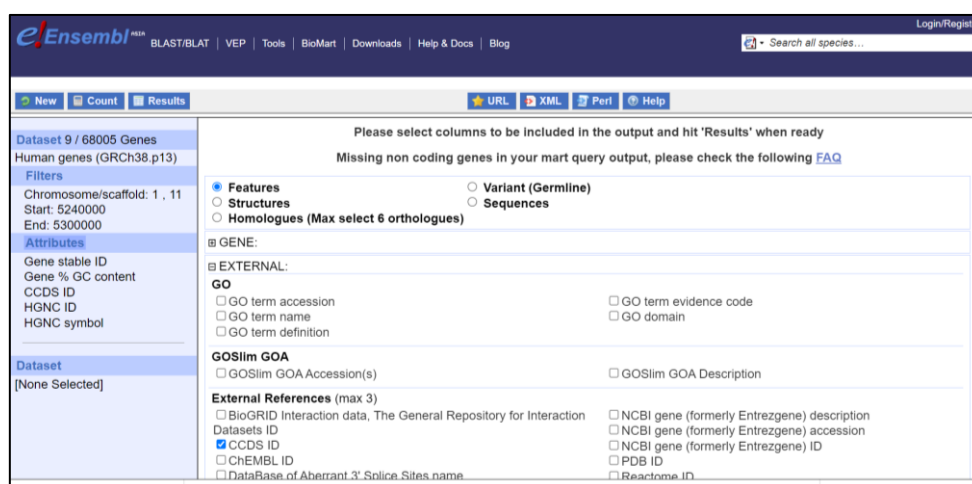Then, select the attributes. Choose gene and external references.
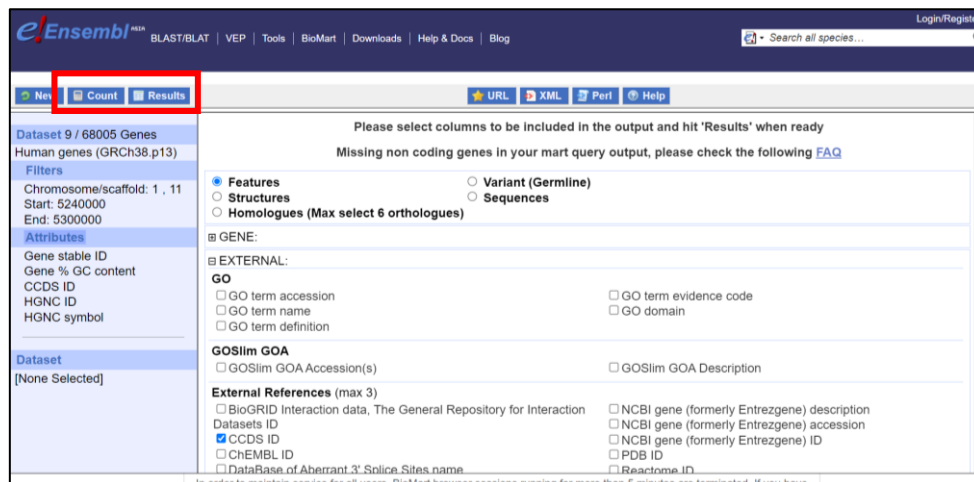


*Figure 17*

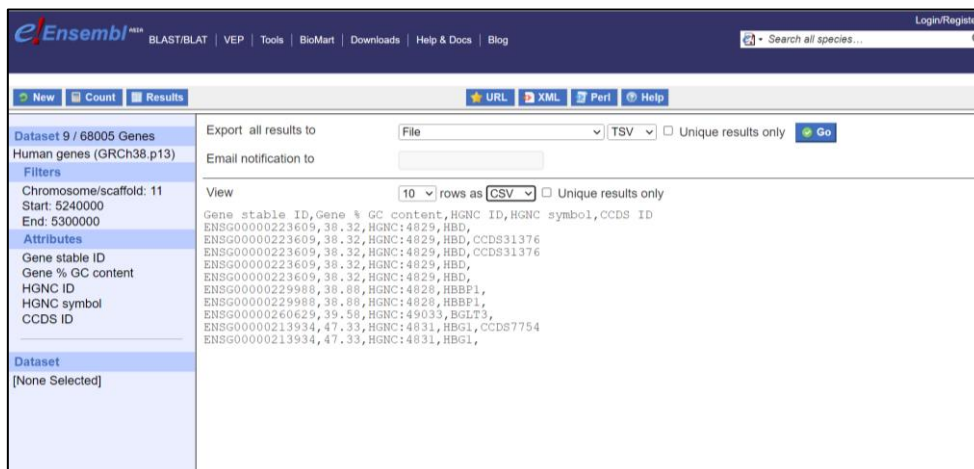Click "count" and "results".



*Figure 18*

The result will be shown.



*Figure 19*

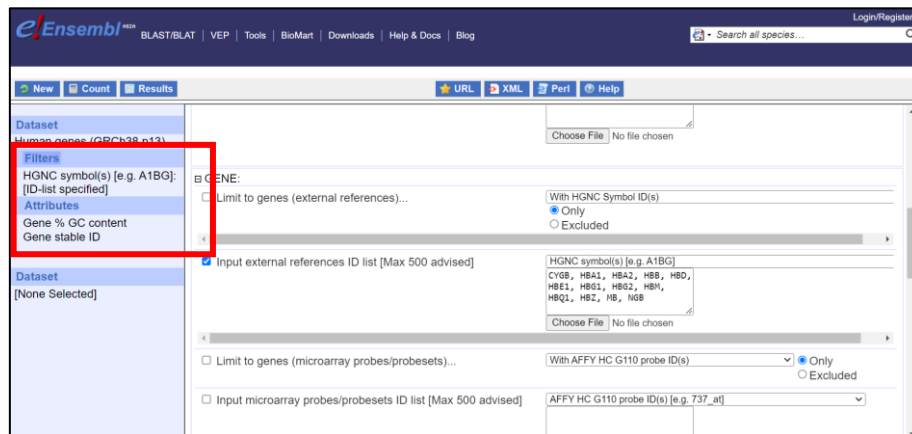## 5. BioMart: working with lists. The goal of this exercise is to access information in BioMart by uploading a text file listing gene identifier of interest.

Follow the steps from problem 2-4, but for the filter set choose Gene (instead of Region), select ID list limit and adjust the pulldown menu to HGNC symbol, then browse for a text file having a list of gene symbols.

See Web Document 2.5 ([WebDocument_2-5_13_humanGlobins_HGNCsymbols.txt](#)) for a text file listing official HGNC symbols for 13 human globin genes (*CYGB*, *HBA1*, *HBA2*, *HBB*, *HBD*, *HBE1*, *HBG1*, *HBG2*, *HBM*, *HBQ1*, *HBZ*, *MB*, *NGB*).

You could also enter these gene symbols manually.

For attributes choose any set of features that is different than in problem 2-4, so that you can further explore BioMart resources.

Change the filters and attributes.



*Figure 20*

The results will be shown after click on the results.



*Figure 21*