# BIOINFORMATICS LAB – AP BIOLOGY

Bioinformatics is the science of collecting and analyzing complex biological data. Bioinformatics combines computer science, statistics and biology to allow scientists to do many things:

- Build phylogenetic trees
- Study mutations found in human cancers
- Investigate where genes are expressed (transcribed and translated) in an organism
- Examine interactions between genes
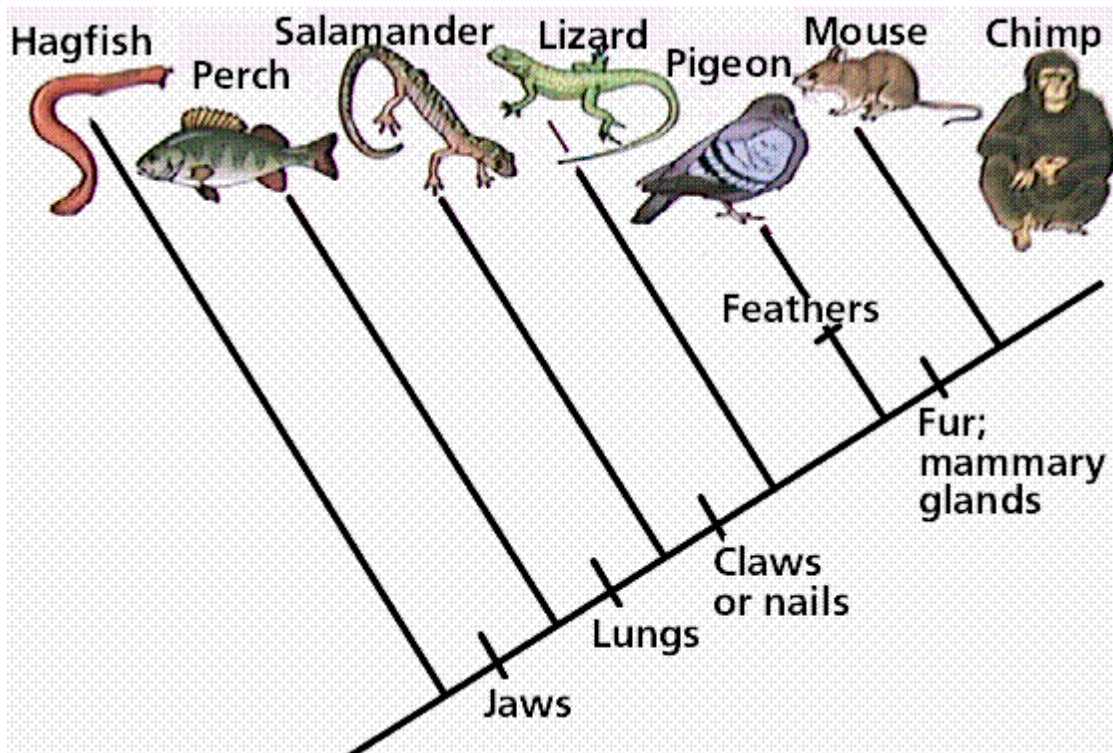- Analyze circadian rhythms in expression of genes

And much, much more.

We will start by using DNA sequences isolated from a fossil to propose placement of that fossil on a phylogenetic tree, then we will explore some of the other free resources available online to do bioinformatics.


## BLAST – Basic Local Alignment Search Tool

One of the ways to determine the degree of relatedness between organisms is to compare their DNA sequences.  The more recently organisms shared a common ancestor, the more closely related they will be, and they will have a greater homology (similarity) in their DNA sequences.  In bioinformatics, a **sequence alignment** is a way of arranging the **sequences** of **DNA**, RNA, or protein to identify regions of similarity.  BLAST compares the sequence you are interested in (the "query" sequence) to the millions of known DNA sequences available, looking for the greatest degree of homology.  Looking at these homologous sequences (called "hits") can be useful in determining the function of a gene (if a gene is involved in glycolysis in mice, it is very likely that the human version of that gene is also involved in glycolysis) or how closely related two species are (the more DNA in common, the more closely related).  You will be using BLAST to propose a possible placement of a fossil on a cladogram (aka phylogenetic tree).


Here is an example of a simple cladogram:

This cladogram includes synapomorphies, shared derived characteristics. From the cladogram, we can say that all of the organisms except the hagfish have the derived characteristic of jaws, while only mice and chimps share the derived characteristics of fur and mammary glands.

While the previous cladogram could tell us about the characteristics that organisms share, we can also use shared characteristics to construct cladograms. Below is a table showing some of the characteristics of major plant groups:

| Physical Characteristics of Major Plant Groups | | | |
|---|---|---|---|
| Plant Group | Vascular Tissue | Flowers | Seeds |
| Mosses | 0 | 0 | 0 |
| Conifers | + | 0 | + |
| Angiosperms | + | + | + |
| Ferns | + | 0 | 0 |

The more shared characteristics two organisms have, the more closely related they are. Using the data above, propose an arrangement of mosses, conifers, angiosperms and ferns on the cladogram below.

GAPDH (glyceraldehydes 3-phosphate dehydrogenase) is an enzyme that catalyzes the sixth step in glycolysis, an important reaction that produces molecules used in cellular respiration. The following data table shows the percentage similarity of this gene and the protein it expresses in humans versus other species. For example, according to the table, the GAPDH gene in chimpanzees is 99.6% identical to the gene found in humans, while the primary sequence of the corresponding protein is identical.

| Percentage Similarity of the GADPH gene and protein with *Homo sapiens* | | |
|---|---|---|
| Organism | Gene Percentage Similarity with Homo sapiens | Protein Percentage Similarity with Homo sapiens |
| *Pan troglodytes* | 99.6 | 100 |
| *Canis lupus familiaris* | 91.3 | 95.2 |
| *Drosophila melanogaster* | 72.4 | 76.7 |
| *Caenorhabditis elegans* | 68.2 | 74.3 |

Based on the data, which species is most closely related to humans?  Least closely?
Why are the percentage similarities higher for the proteins than the genes?

The Liaoning Province in northeastern China is a rich source of dinosaur fossils.  In the 1990s, scientists found the following fossil:

Based on the morphology (shape) of the fossil in the picture, make a hypothesis as to where the organism in the fossil should be placed on the cladogram below. (Draw a line labeled "fossil specimen" on the cladogram). In Darwin's day, scientists had to rely on the morphology of the fossil and information gleaned from the surrounding rock layers when classifying a fossil Now, with the advent of biotechnology, scientists can sometimes isolate intact DNA from fossils, sequence the DNA, and use this information to refine the fossil's classification.

Scientists were able to isolate and sequence DNA from this fossil. Your job is to take the fossil DNA sequences (found on the class webpage), run them through the BLAST tool, and use your results to propose a placement of that fossil on a cladogram.

**PART ONE**

1. Go to the BLAST homepage (https://blast.ncbi.nlm.nih.gov/Blast.cgi). Click on the Nucleotide Blast box.



2. Your screen should look something like this:



Query box   In the query box, copy and paste your first DNA sequence.

3. For "Job Title" enter "Gene 1" or "Gene 2" or etc…..
4. For "Database" use the pull down menu to choose "Nucleotide Collection."  Leave "Organism" black since we want to search all organisms.
5. Under "Program Selection" optimize for "highly similar sequences."
6. Click the "BLAST" button at the bottom of the screen, then **be patient**!  Depending on the speed of your internet connection and how busy the server is, it may take a few minutes to see your results.

Your results will look something like this:



Here is some of the information this screen is telling you:

Query ID:  this is a unique search identifier that indexes your specific search

Molecule type: nucleic acid (since we were looking at DNA)

Query length: 5575 (the number of nucleotides in your query sequence – gene 1 in this example)

7. Now look at the colored bars.  Red indicates a strong match with a high level of homology.  Each bar represents a "hit," or a sequence that the BLAST program found in the database that has a high level of homology with your query sequence.  The bars are arranged in order with the strongest hit (with the highest degree of homology) on top.  Click on the strongest hit (top red bar)

You should see something like this:

8. The white box tells you that this sequence came from the organism *Gallus gallus*  (google it if you don't know what organism that is) and the DNA sequence in that organism codes for the protein collagen.   Now click on the word "Alignment" in the white box
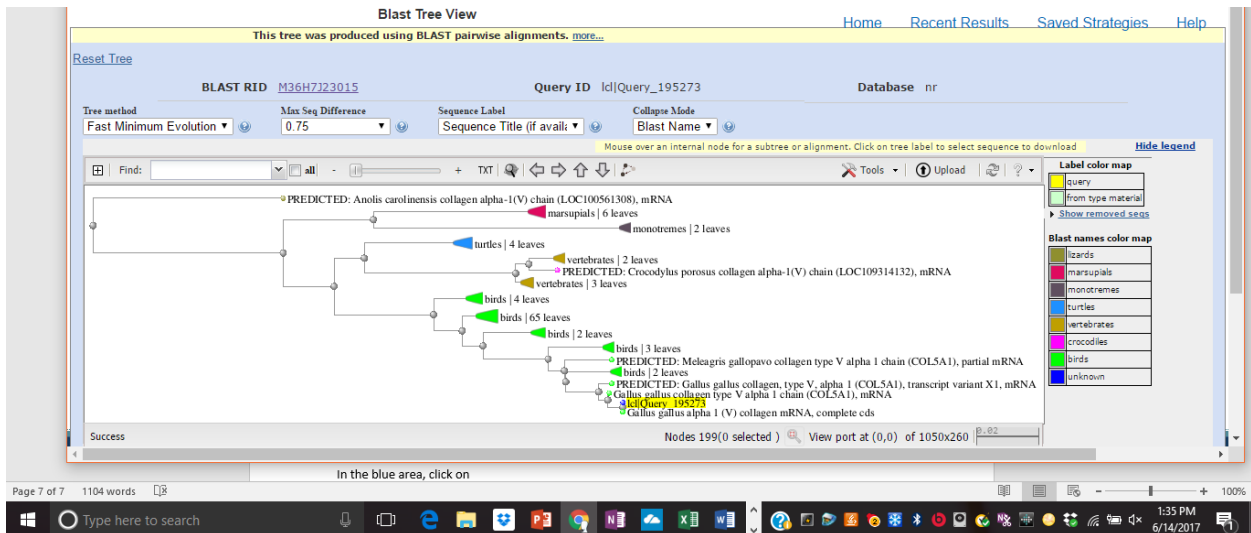


9. This is NOT a double stranded sequence of DNA.  How do you know this?

The query line shows the DNA sequence you copied into BLAST.  The Sbjct line show the aligned DNA sequence that BLAST found in the database.  Each vertical line represents a match between the query and subject sequences.  Score tells you how strong the homology is, anything over 200 is considered a

strong match. Record the score in your notes. Now use the back arrow to return to the screen with the red bars.



10. In the blue area, click on "distance tree of results." This will create for you a proposed cladogram with your query sequence highlighted in yellow and placed on the cladogram.



11. Based on the identity of the species *Gallus gallus* and the distance tree of results, do you want to revise your initial placement of the fossil on the cladogram? Why or why not?

Repeat steps #1-11 for the remaining gene sequences obtained from the fossil. Where do you think the fossil should be placed on the initial cladogram and why?

**PART TWO – ONLINE MENDELIAN INHERITANCE IN MAN (OMIM)**

1. Google "OMIM" or click on this link https://www.omim.org/



In the search box enter the name of a disease or gene you are interested in.  I'm using breast cancer

Here is what you might see:



Note that the third result is for one of the most well know breast cancer genes, BRCA1.  Click on one of your search results to find more information about your gene.  I clicked on BRCA1 and saw this:

On the left hand side of the screen in blue, you can find a wealth of information about the gene, such as its structure and function.

2.  On the right hand side of the screen, click on "genome" and then NCBI Map viewer to see what chromosome(s) the gene resides on. Record the location of your gene.

**PART THREE – BioGPS – Bio Gene Portal System**

3.  On the right hand side of the screen, click on "Gene Info" and then "BioGPS."  This will open a new window that will show you in which tissues your gene of interest is most expressed in (highest levels of mRNA).  If there is little difference in the expression of your gene among different tissues, it is probably necessary in most, if not all, cells and is called a "housekeeping gene."  You can use the slider to zoom in and out. Record where your gene is most strongly expressed.



**PART FOUR – KEGG – KYOTO ENCYCLOPEDIA OF GENES AND GENOMES**

4.  On the right side of the screen, click on "KEGG"  This will give you information on the metabolic pathways the gene is part of, diseases related to this gene, the structure of the gene, the amino acid sequence of the protein coded for by the gene, and the nucleotide sequence of the gene.

Metabolic pathways   Diseases

Scroll down to structure and click on "JMOL" to see the 3D structure of the protein coded for by the gene.



Estimate from the picture how many alpha helices and beta sheets are present in the protein.

## PART FIVE – CIRCADB – CIRCADIAN RHYTHMS

1. Go to http://circadb.hogeneschlab.org/

**2.** Enter the name of your gene

**3.** Examine the results and determine if the expression of your gene of interest  is subject to circadian rhythms.

**PART SIX – MGI – Mouse Genome Informatics**

4. Go to http://www.informatics.jax.org/vocab/gene_ontology



5. In the search box enter the name of your gene.
6. This screen should give you some ideas as to the function of your gene

Brian - Google Drive × | BRCA1-BARD1 complex × | CIRCA: Circadian gene e ×

www.informatics.jax.org/vocab/gene_ontology/GO:0031436

MGI

About Help FAQ

Keywords, Symbols, or IDs | Quick Search

Home | Genes | Phenotypes | Human Disease | Expression | Recombinases | Function | Strains / SNPs | Homology | Pathways | Tumors

Search ▾ | Download ▾ | More Resources ▾ | Submit Data | Find Mice (IMSR) | ⚙ Analysis Tools | Contact Us | Browsers

## Gene Ontology Browser
Molecular Function | Biological Process | Cellular Component

### GO Search

| brca1 | | Clear |

5 terms, sorted by best match

**BRCA1**-BARD1 complex
**BRCA1**-A complex
**BRCA1**-B complex
**BRCA1**-C complex
**BRCA1**-Rad51 complex

### GO Term Detail

Term: **BRCA1-BARD1 complex**
Definition: A heterodimeric complex comprising BRCA1 and BARD1, which possesses ubiquitin ligase activity and is involved in genome maintenance, possibly by functioning in surveillance for DNA damage.
Parent Terms: *is-a* nuclear ubiquitin ligase complex
Category: Cellular Component

### GO Tree View

- nuclear stress granule
- ▶ nuclear transcription factor complex
- ▼ nuclear ubiquitin ligase complex
  - anaphase-promoting complex
  - Asi complex
  - **BRCA1-BARD1 complex** (2 genes, 3 annotations)
  - CLRC ubiquitin ligase complex
  - nuclear SCF ubiquitin ligase complex
  - PRC1 complex
  - SUMO-targeted ubiquitin ligase complex

Contributing Projects:

Mouse Genome Database (MGD), Gene Expression Database (GXD), Mouse Tumor Biology (MTB), Gene Ontology (GO), MouseCyc