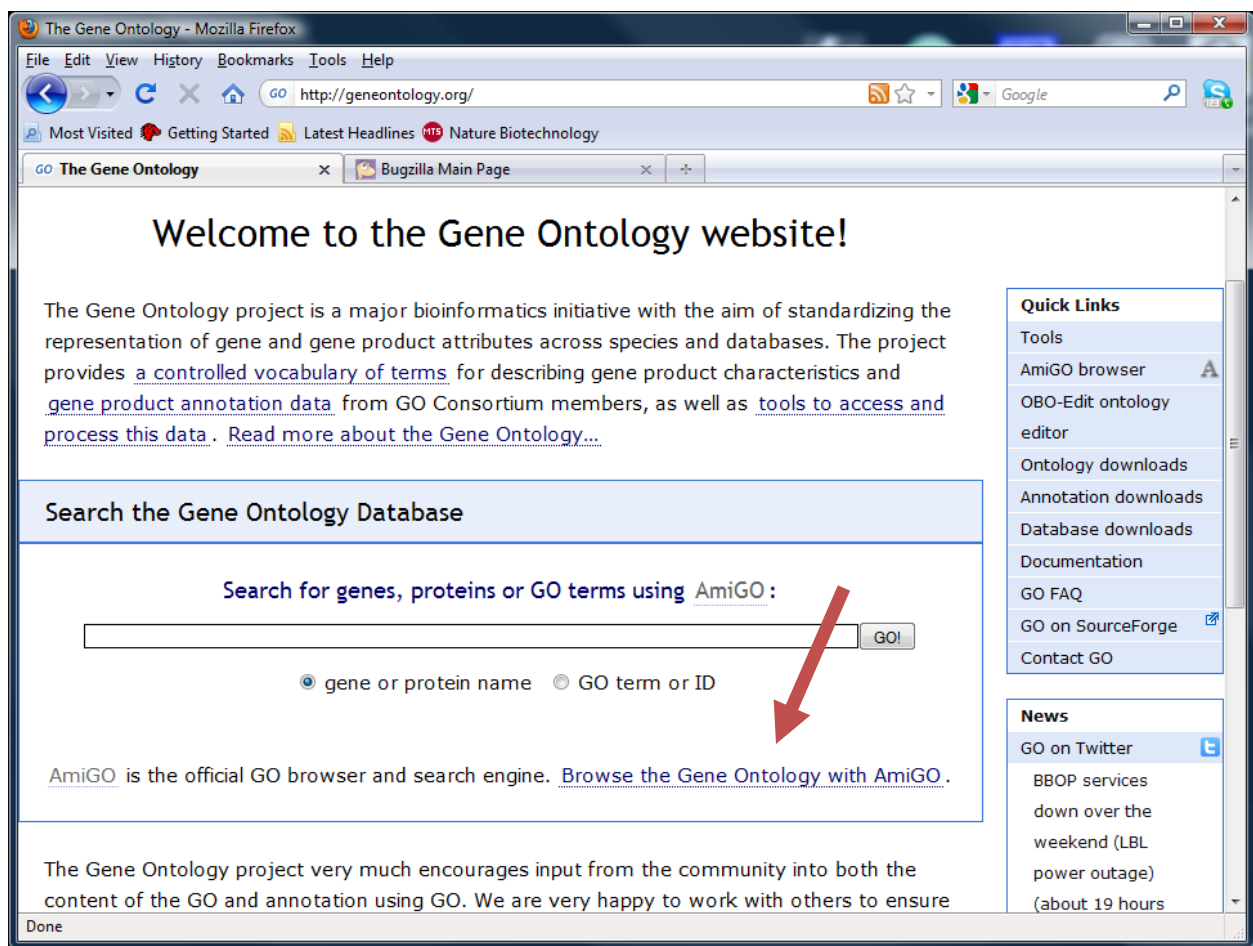# GO, DAVID and iRefScape revisited.

Ian Donaldson

MBV-INF 4410/9410

This exercise will revisit a number of applications introduced in the class. You will learn how to browse the Gene Ontology, select a set of genes that are annotated with some specific GO term and then you will analyze this gene list using DAVID and iRefScape.

Go to http://geneontology.org/



Click on "Browse the Gene Ontology with AmiGO".

Spend some time browsing through each of the three GO Ontologies.  Look for terms that you are familiar with and see how they relate to terms above (parent terms) and below (child terms) in the GO.



Try to make your way to the term that describes "nucleus" without directly searching for it. Instead, start by expanding the cellular_component "root node" (GO:0005755) by clicking on the + beside it.  Look for the next closest thing to "nucleus" in the expanded list and then click on that.  If you want to look at the definition of any given term, just click on it.

If you cant find the term entry for "nucleus", use the next page as a hint or search for nucleus in the "Search GO" box at the top of the interface.

Once you make it to the nucleus, keep navigating down the tree to chromatin assembly complex (GO:0005678). This term is actually used to refer to a number of complexes (try expanding the term). What are these complexes? Note that all of them are "leaf" nodes (you cant expand them any further).

Now go back and click on the 47 gene products that correspond to the "chromatin assembly complex".

Click on the "47 gene products" to view them.

Spend some time browsing the links from this page (or just hovering over them). Scroll down. Genes annotated with each of the three leaf terms are grouped separately on the page.

What kinds of evidence are provided for these assignments. What kinds might you be cautious of (less likely to believe) or more likely to believe? Who assigns these GO terms to the genes?

Try filtering the list using the filters at the top. Make selections from the menus (Gene Product Type, Data source, Species and Evidence code and then click "Set filters".

Next, try to export your results for all genes from all data sources from human (Homo sapiens) that have any evidence code. Set filters and then click on "Download all association information in gene association format".

Copy and paste this to an Excel spread sheet (or similar). Right-click, paste special, as text.





Select and copy the UniProt accessions from column B.

Go to http://david.abcc.ncifcrf.gov/home.jsp , click on Start Analysis and then paste the list of accessions into the query box like this:

The tell DAVID that you have entered UniProt Accessions. Like this

Then select "Gene List" under List type and click on Submit list:

Your list has been saved by DAVID as List_1 and it has automatically recognized "Homo sapiens" as the species from which the list is derived:

If you click on the "Background" tab, you will see that DAVID has also set "Homo sapiens" as the "background". Its important that you check these settings and change them if necessary. Discuss why.

You can now look at categories that are over-represented in your list. For example, click on the + beside Gene_Ontology and scroll down to GOTERM_CC_FAT like this:

If you then click on the "Chart" button, you see a new window open like this:

6 out of the 7 genes in your list (85.7%) are annotated with the GO term for chromatin assembly complex. The probability of randomly choosing 7 genes from the human genome where 6 of them all have this associated GO term is 1.3E.17. When you correct for multiple hypothesis testing (Benjamini), the probability is 2.3E-16 (still quite surprising – i.e. we would suspect that whatever "process" was used to pick out these 7 genes was not unrelated to this annotation). And of course, we know this to be the case.

Click on the "chromatin assembly complex" to see details about the term that is overrepresented.

Click on the blue bar underneath "Genes" to see the list of genes that had this annotation.

Go back to the "Annotation summary results" and explore whether other categories of annotation were over represented in this list.

There is a lot of other material you can explore on this site. Make a note of the Nature Protocols tutorial on use of DAVID for later.
http://www.nature.com/nprot/journal/v4/n1/pdf/nprot.2008.211.pdf

Next, we will explore known interactions between the genes in our list.

Start Cytoscape and the iRefScape plugin from Plugins menu?iRefScape 0.9. Then copy and paste the identifiers from column B of the Excel spreadsheet to the query box. Like this…



Make sure that the search type is set to UniProt Ac, the taxon is set to Any (or Homo Sapiens) and that "Use canonical expansion" is selected and that iterations are set to 0 (see red arrows above). Iterations of zero will only return interactions that occur between proteins in our query list. Then click on search and load.

The initial results will look like this

And you should be able to clean them up (rearrange them) to look like this.

Hint use Layout/align and distribute to align and stack the purple hexagons (nodes that represent complexes) and then iRefScape/ViewTools/Toggle selected multi-edges to hide multi-edges (representing multiple experiments that support the same protein-protein interaction).

From this view it should be apparent that

1) 4 of the nodes don't interact with any of the other proteins in the list. Although 2 of them are self interacting (loops)

2) 2 of the proteins (CAF1A and CAF1B) appear to be co-members of multiple complexes that are documented by multiple databases (these complexes are represented by the purple hexagons). By definition, these hexagons represent complex records with 3 or more proteins.

3) These same two proteins are reported to interact with one another by one database (HPRD). You will only se this interaction if you are using the version of iRefScape that contains HPRD data (you have to ask for this version).

4) If you want to explore the evidence for any interaction, click on the edge and look at details in the edge attribute browser. For example, the i.PMID feature lists the publication where the evidence for the interaction was found.

5) If you want to explore the other members of each of the complexes, click on the complex node ("pseudo-node") and then select iRefScape/Search tools/ retrieve interactions for selected nodes.

Go back to the Gene Ontology pages. Should more of these proteins have been involved in a single complex?

It is left as an exercise to see if you would get different results by querying iRefScape with proteins from different organisms (say Drosophila) that are annotated as belonging to the "chromatin assembly complex".

Next, we will move on to a technique that will identify proteins related to our initial query. We want to find proteins that interact with 2 or more of the proteins in our starting list (i.e. things that are associated with the "chromatin assembly complex".

Select all nodes in the current view.

Retrieve their neighbours using iRefIndex menu/Search tools/Retrieve interactions for selected nodes.

You will see (a ridiculously large network – 1106 nodes and 9006 edges) like this:



Select all nodes using control-A.

In the node attribute browser, select to view the i.query node feature and then sort on this feature by clicking on the i.query column heading. You should see all the nodes from your initial search at the top of the node attribute browser. Like this:



Left-click and drag over these table entries to select them with the mosue.

The right-click on one of the hi-lited entries and choose "Select from table".

After this operation, only the original nodes directly returned by the query will apeear in the node attribute browser.
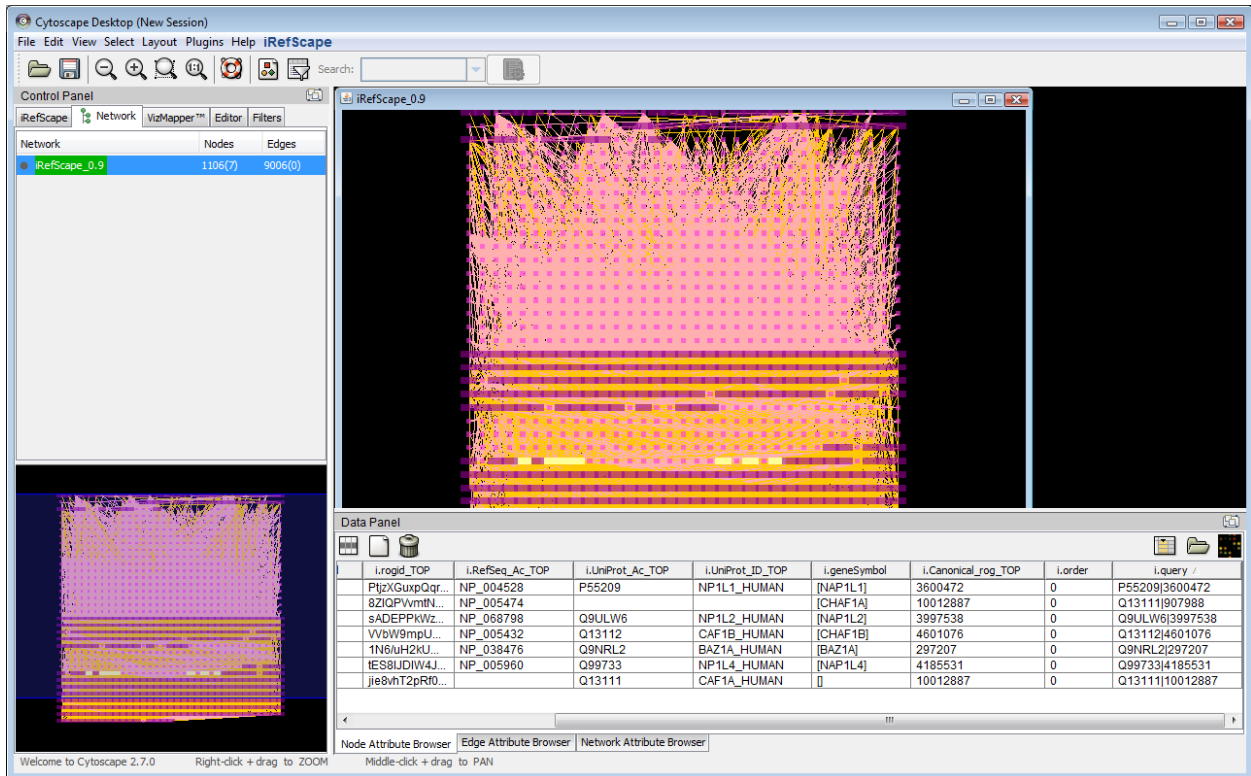
Then go to the iRefScape menu/View tools/Select between nodes.

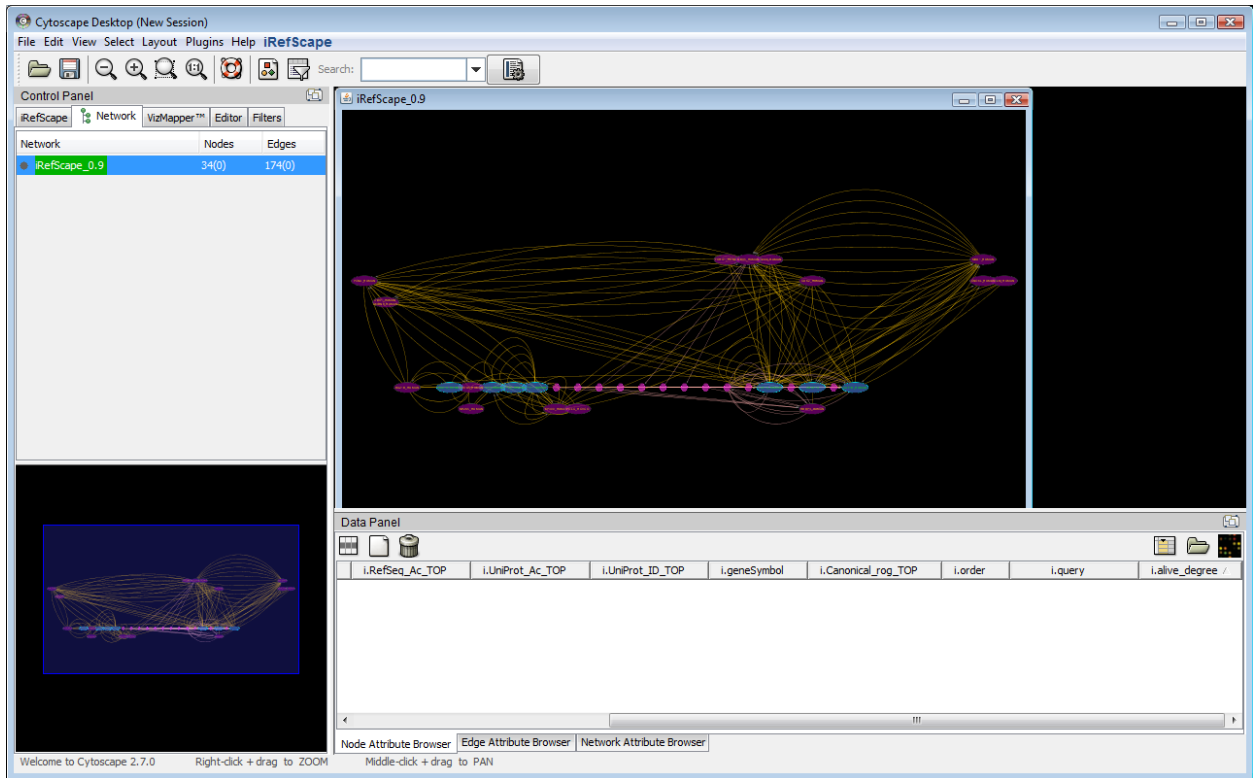Use the node attribute browser to select all nodes that are either

1) Returned by the original query (see i.query node feature) or
2) Have an i.alive_degree of 2 or more (see the i.alive_degree node feature)

This may require a live demo.  Ask if you have problems.

Then select all other nodes and delete them.  Use Select/Nodes/Invert node selection

And then hit delete.


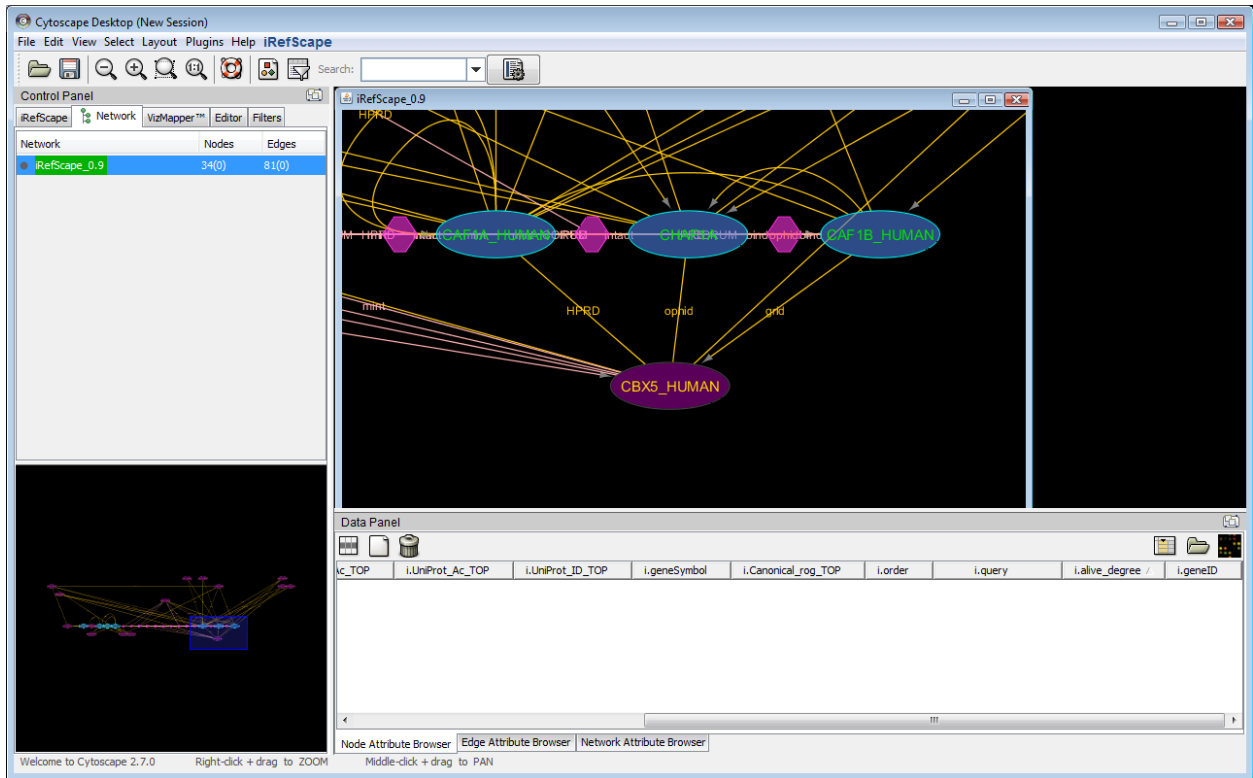You should see something like this.

These new "between nodes" are connected to two or more nodes from your original search list (seed list).  Guilt-by-association makes these nodes  candidates to look at for things that may be related to the "chromatin assembly complex".  The higher the i.alive_degree, the better.

For example, CBX5_HUMAN interacts with CAF1A, CAF1B and CHAF1A from our original search.

A brief review of the Entrez Gene record for CBX5 (Entrez Gene ID 23468) shows that there is evidence to support this connection. It is left as an exercise to review the papers that support evidence of an interaction between CBX5 and each of the three original query proteins.

If you have time left, you can use the techniques described in this tutorial to investigate another GO term and its associated genes that are of specific interest to you.