

# **GS01 0163**

## **Analysis of Microarray Data**

Keith Baggerly and Kevin Coombes  
Section of Bioinformatics

Department of Biostatistics and Applied Mathematics  
UT M. D. Anderson Cancer Center

`kabagg@mdanderson.org`  
`kcoombes@mdanderson.org`

5 September 2006

# Lecture 3: Linking Numbers to Biology

- So, why are we here?
- Why do we care?
- Affymetrix source for annotations
- List of Affymetrix annotations
- Updating the annotations in dChip
- What is GeneOntology?
- Using GeneOntology in dChip
- GoMiner

## So, why are we here?

We want to learn about Gene Annotations.

Microarrays are *designed*, which means that someone first chooses a set of genes of interest, selects probe sequences to target those genes, and then places those sequences on a microarray. In order to interpret (and possibly to analyze) the data produced from a microarray experiment, you need to refer to the accompanying annotations, which describe both the probes and the targeted genes.

## Things Change

One might naively think that gene annotations are static; meaning that they are produced when the microarray is designed and never change again. Let me disabuse you of that notion immediately. It is true that the biological sequences of the probes that were placed on the array do not change. However, our knowledge of the human genome continues to evolve, and thus our opinion about exactly what genes are targeted by those sequences must be continually updated.

For Affymetrix microarrays, the company maintains a web site that always contains their latest opinion on the nature and identity of the targeted genes.

## Why Do We Care?

Recall from the last lecture that we compared microarray data from samples of acute lymphocytic leukemia (ALL) patients and mixed-lineage leukemia (MLL) patients. Using the criteria that the lower bound of fold change (LBFC) should be at least 1.2-fold and the mean difference in expression should be greater than 100, we found a list of 610 probe sets that were differentially expressed.

It is considered bad form to just hand the biologists a list of 610 genes and wish them good luck as they go on their way. They typically want to know: do these genes reflect particular biological functions that are different between the two groups of samples, or do they identify specific biological pathways or networks that are perturbed?

# List of Differentially Expressed Genes

	A	B	AA	AB	AU	AV	AW	AX	AY
12	probe set	gene	baseline mean	baseline	experiment mean	experiment	fold change	lower bound	upper bound
13	37680_at	A kinase (PRKA) anchor protein (gravin) 12	2973.7	560.63	148.29	24.19	-20.05	-12.93	-30.28
14	1325_at	MAD, mothers against decapentaplegic homolog 1	7759.92	1390.4	595.18	64.06	-13.04	-8.89	-18.03
15	37280_at	MAD, mothers against decapentaplegic homolog 1	9124.17	1538.9	702.89	37.85	-12.98	-9.29	-16.88
16	37908_at	guanine nucleotide binding protein 11	2160.91	565.93	226.99	58.16	-9.52	-4.92	-18.23
17	34194_at	Homo sapiens mRNA; cDNA DKFZp564B076 (from	962.11	296.29	107.48	34.97	-8.95	-3.95	-21.14
18	753_at	nidogen 2 (osteonidogen)	2558.48	890.45	304.16	22.09	-8.41	-3.58	-13.49
19	1992_at	fragile histidine triad gene	1742.98	252.98	209.02	29.64	-8.34	-5.92	-11.72
20	1488_at	protein tyrosine phosphatase, receptor type, K	4128.67	1140	572.2	38.89	-7.22	-3.91	-10.70
21	1077_at	recombination activating gene 1	6927.92	1443.9	1021.43	204.85	-6.78	-4.09	-11.13
22	33910_at	Homo sapiens mRNA; cDNA DKFZp564P116 (from	460.85	209.6	72.66	7.64	-6.34	-1.59	-11.49
23	34800_at	leucine-rich repeats and immunoglobulin-like domain	5255.48	907	899.41	189.08	-5.84	-3.74	-9.53
24	35614_at	transcription factor-like 5 (basic helix-loop-helix)	7264.11	1378.1	1248.25	122.02	-5.82	-3.9	-8.05
25	41266_at	integrin, alpha 6	7923.59	1222.5	1445.79	200.87	-5.48	-3.84	-7.73
26	37343_at	inositol 1,4,5-triphosphate receptor, type 3	5231.99	747.28	966.99	97.72	-5.41	-3.98	-7.15
27	31892_at	protein tyrosine phosphatase, receptor type, M	801.09	336.26	150.51	9.57	-5.32	-1.64	-9.12
28	35669_at	KIAA0633 protein	1738.34	360.27	343.94	22.32	-5.05	-3.3	-6.93
29	38578_at	tumor necrosis factor receptor superfamily, member	4038.17	674.75	847.39	129.09	-4.77	-3.23	-6.94
30	37780_at	piccolo (presynaptic cytomatrix protein)	2856.4	830.13	601.56	40.43	-4.75	-2.46	-7.15
31	40570_at	forkhead box O1A (rhabdomyosarcoma)	10218.69	1178.1	2227.99	482.41	-4.59	-3.16	-7.34
32	39878_at	protocadherin 9	12518.61	2120.5	2816.54	552.51	-4.44	-2.89	-7.03
33	307_at	arachidonate 5-lipoxygenase	6743.7	992.9	1521.71	136.37	-4.43	-3.26	-5.80
34	38408_at	transmembrane 4 superfamily member 2	6543.7	1009.8	1489.02	230.77	-4.39	-3.04	-6.36

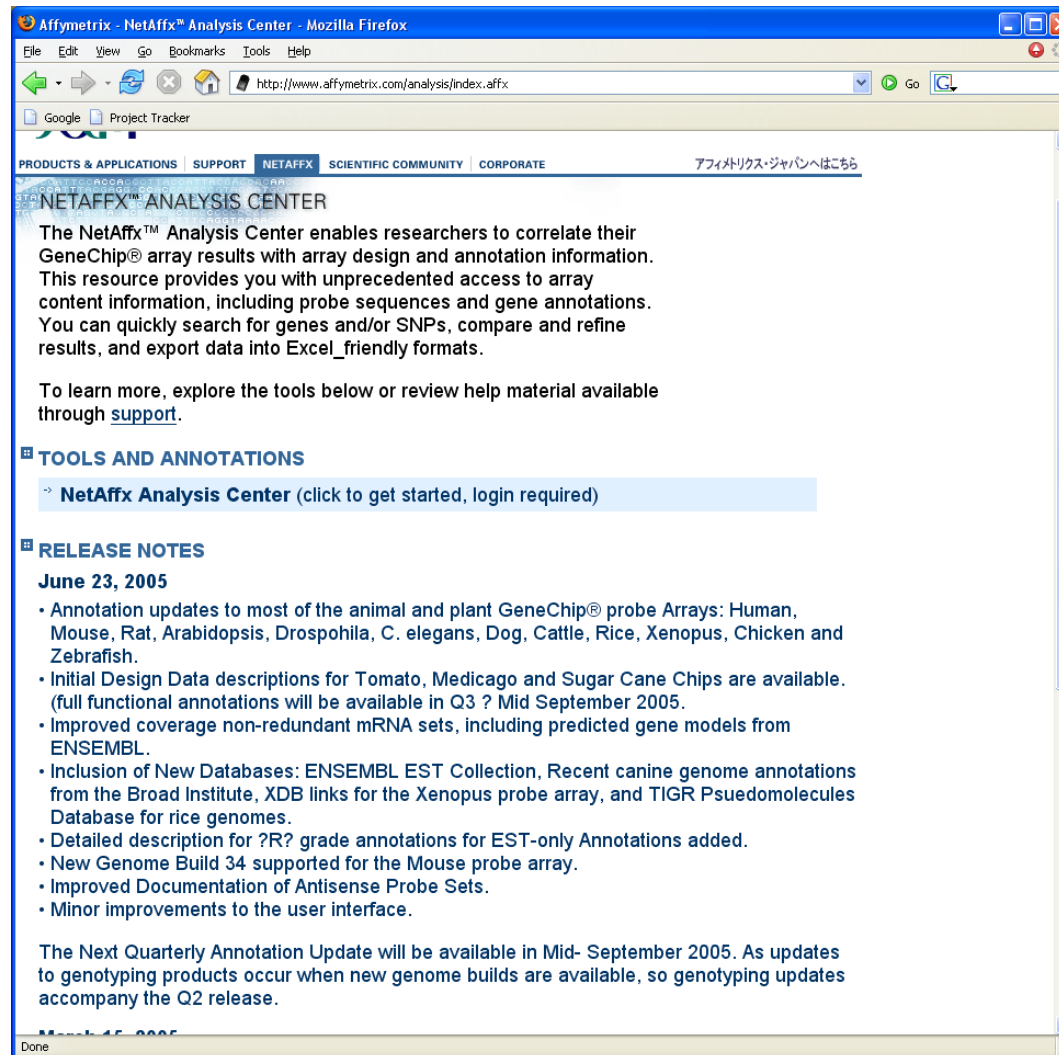
# Affymetrix Web Site

`http://www.affymetrix.com`



# NETAFFX

Annotations are updated quarterly...





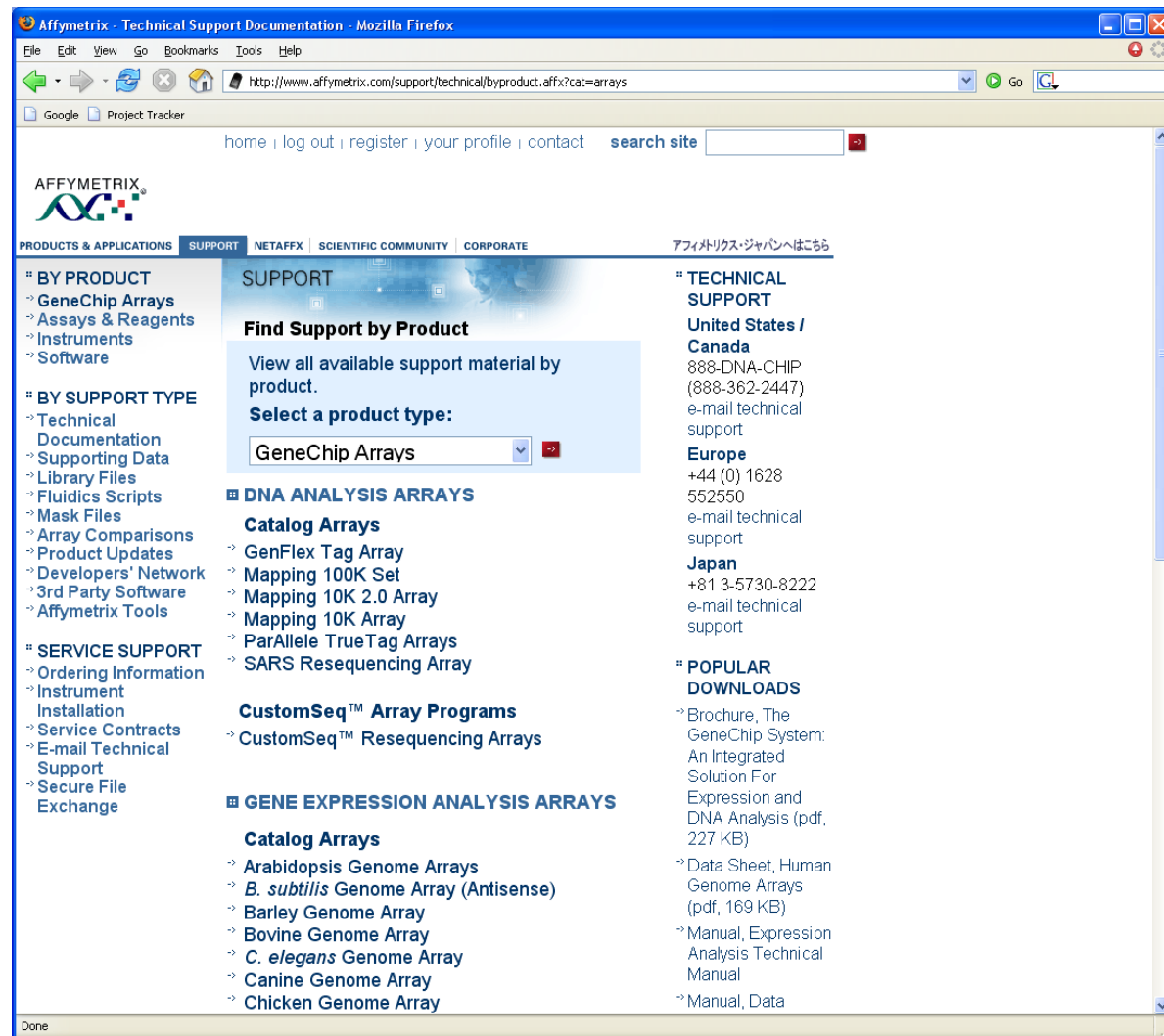
# Affymetrix Support

Go to the Affymetrix support page to get the full annotations.



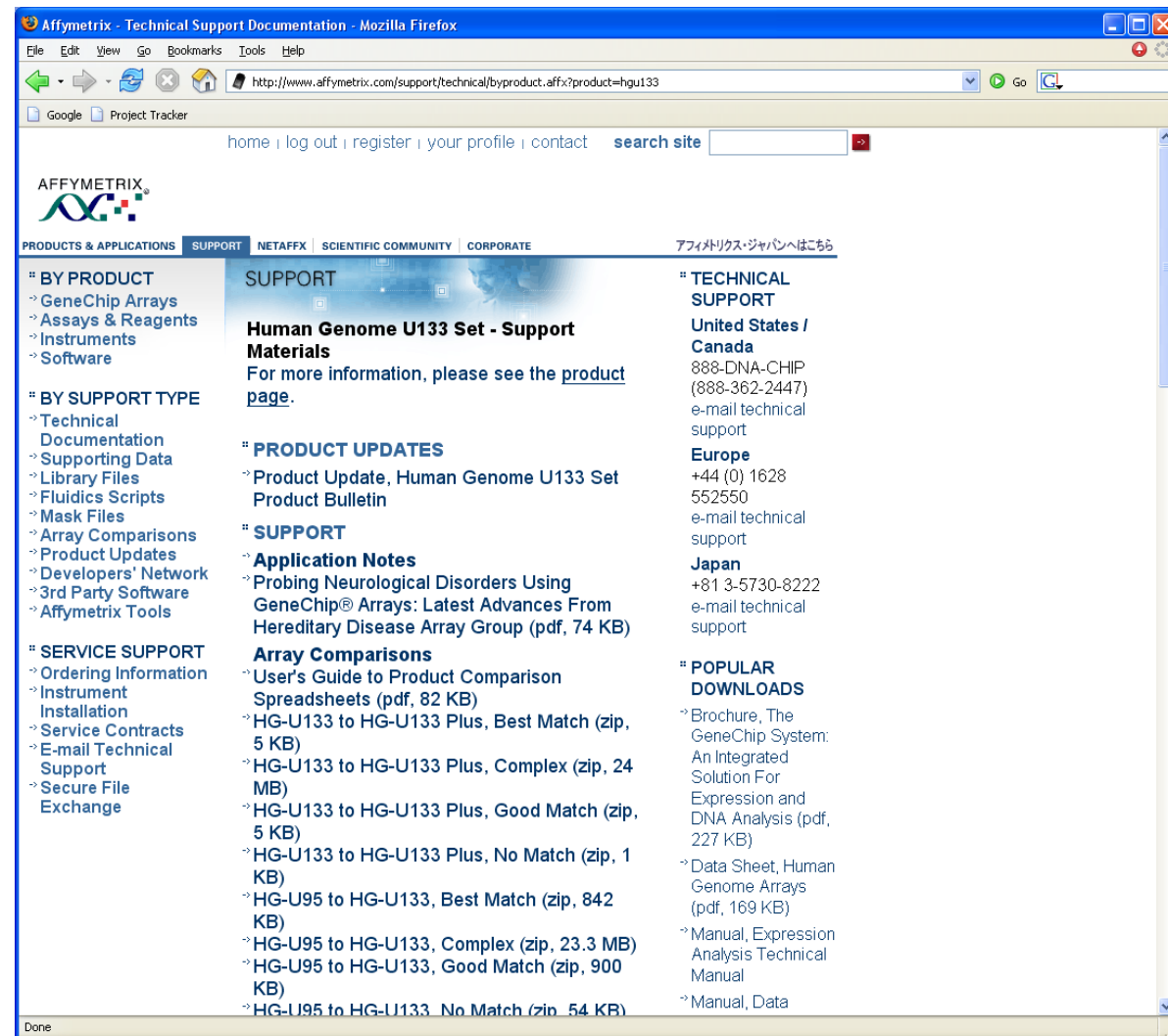
# Support By Product

Follow the “support by product” link to “GeneChip Arrays”.



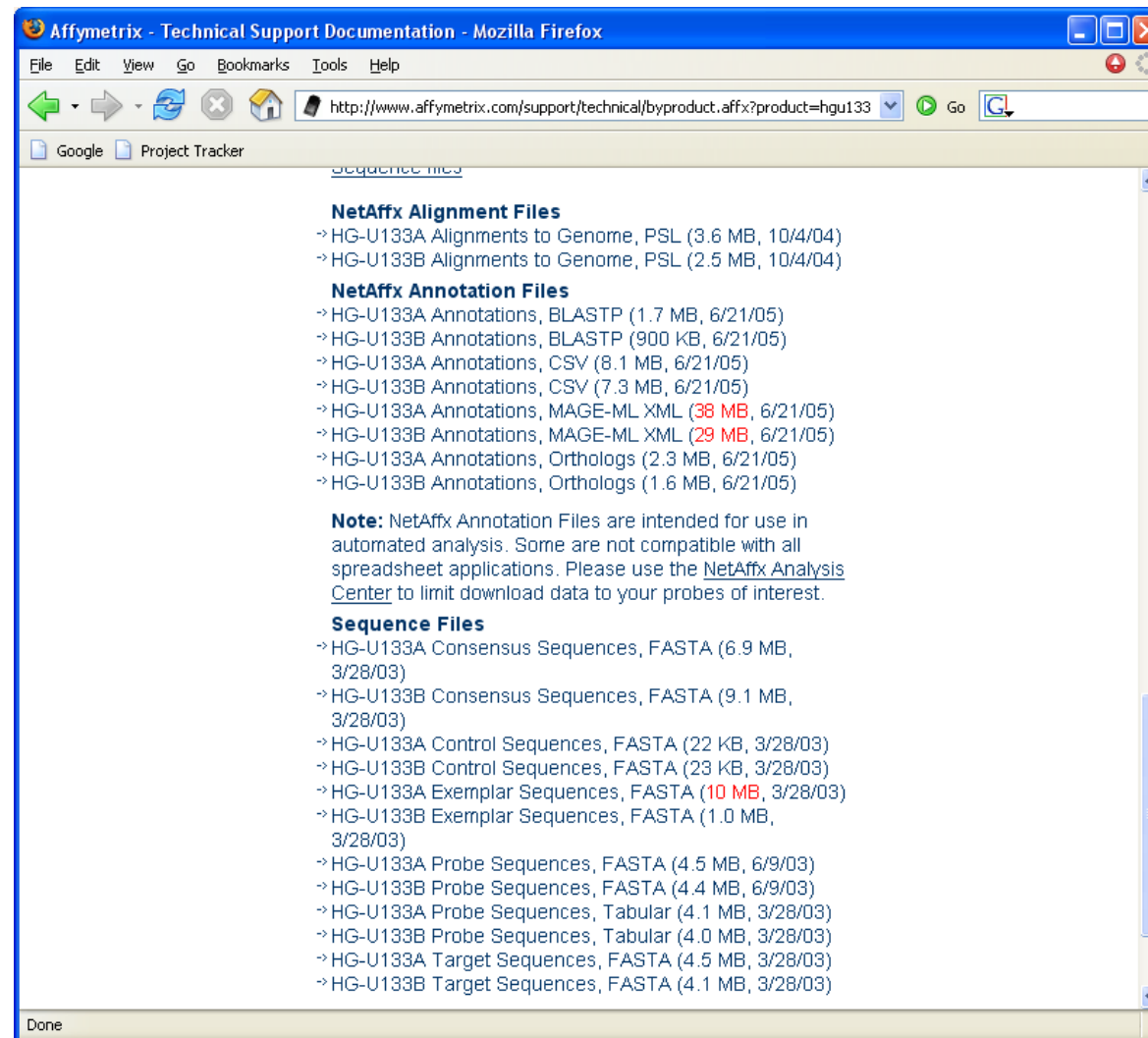
# Affymetrix Annotations for HU133

Scroll down to “Human Genome Arrays”; select “HG-U133 Set”



# Affymetrix Annotations for HU133

Scroll to get a list of available files.



# Affymetrix Main Annotation Files

There are three primary annotation files:

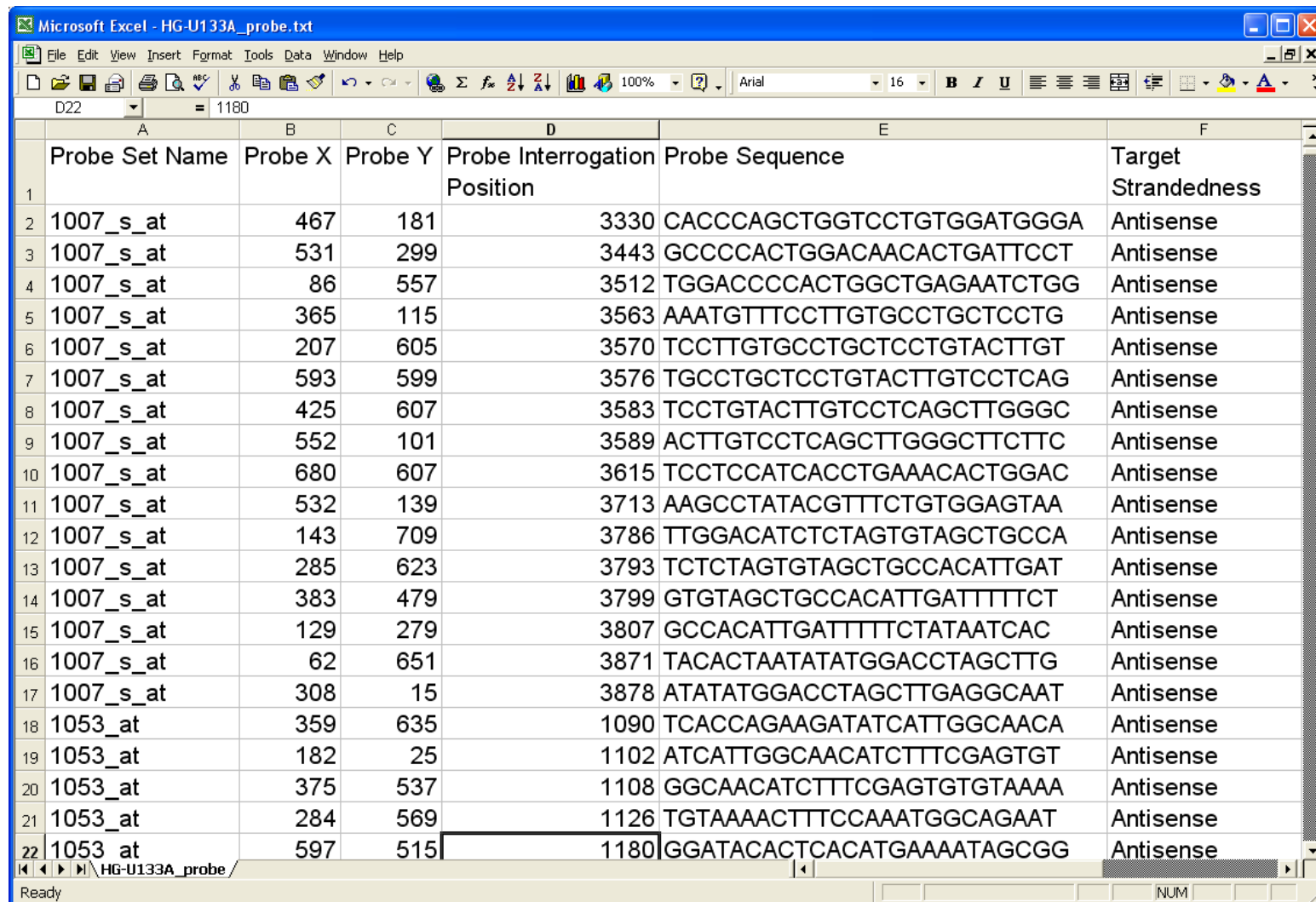
**Probe Sequence File:** Contains a complete listing of all the probes (25-mers) and probe sets on the microarray. (In tab-separated values format, the zipped file is 4.1MB; unzipped, it is 14.4MB.)

**Alignment File:** Contains mappings of targets and probes to the human genome. (In PSL format, the zipped file is 3.6MB; unzipped, it is 12.7MB.)

**Annotation File:** Contains the updated annotations of all the genes targeted by the microarray. (In comma-separated-value format, the zipped file is 6.1MB; unzipped, it is 47.9MB.)

# Affymetrix Probe Sequences for HU133

The HG-U133A\_probe\_tab file lists the probe sequences.



	A	B	C	D	E	F
	Probe Set Name	Probe X	Probe Y	Probe Interrogation Position	Probe Sequence	Target Strandedness
1						
2	1007_s_at	467	181	3330	CACCCAGCTGGTCCTGTGGATGGGA	Antisense
3	1007_s_at	531	299	3443	GCCCCACTGGACAACACTGATTCCT	Antisense
4	1007_s_at	86	557	3512	TGGACCCCACTGGCTGAGAATCTGG	Antisense
5	1007_s_at	365	115	3563	AAATGTTTCCTTGTGCCTGCTCCTG	Antisense
6	1007_s_at	207	605	3570	TCCTTGTGCCTGCTCCTGTACTTGT	Antisense
7	1007_s_at	593	599	3576	TGCCTGCTCCTGTACTTGTCTCAG	Antisense
8	1007_s_at	425	607	3583	TCCTGTACTTGTCTCAGCTTGGGC	Antisense
9	1007_s_at	552	101	3589	ACTTGTCTCAGCTTGGGCTTCTTC	Antisense
10	1007_s_at	680	607	3615	TCCTCCATCACCTGAAACACTGGAC	Antisense
11	1007_s_at	532	139	3713	AAGCCTATACGTTTCTGTGGAGTAA	Antisense
12	1007_s_at	143	709	3786	TTGGACATCTCTAGTGTAGCTGCCA	Antisense
13	1007_s_at	285	623	3793	TCTCTAGTGTAGCTGCCACATTGAT	Antisense
14	1007_s_at	383	479	3799	GTGTAGCTGCCACATTGATTTTCT	Antisense
15	1007_s_at	129	279	3807	GCCACATTGATTTTCTATAATCAC	Antisense
16	1007_s_at	62	651	3871	TACACTAATATATGGACCTAGCTTG	Antisense
17	1007_s_at	308	15	3878	ATATATGGACCTAGCTTGAGGCAAT	Antisense
18	1053_at	359	635	1090	TCACCAGAAGATATCATTGGCAACA	Antisense
19	1053_at	182	25	1102	ATCATTGGCAACATCTTTCGAGTGT	Antisense
20	1053_at	375	537	1108	GGCAACATCTTTCGAGTGTGTAAAA	Antisense
21	1053_at	284	569	1126	TGTAACACTTTCCAAATGGCAGAAT	Antisense
22	1053_at	597	515	1180	GGATACACTCACATGAAAATAGCGG	Antisense

# Affymetrix Alignment Information for HU133

The HG-U133A.link.psl file aligns the probes to the genome.

```

emacs@BSFC2-COOMBES
File Edit Options Buffers Tools Help
# Filename: HG-U133A.link.psl
# Array: HG-U133A
# Genome: Human_May_2004
# Date: Fri Sep  3 16:41:38 PDT 2004.
#
# This PSL file contains alignments of Affymetrix consensus and exemplar
# sequences to a genome, and the mapping of Affymetrix probe sets,
# poly-A sites, and poly-A stacks onto those sequences.
#
# Feature types are segregated in blocks, each preceded by a track line.
#
# For terms of use, see http://www.affymetrix.com/site/terms.affx
#
track name="HG-U133A netaffx consensus" description="Consensus Sequences"
0      16      192      1      2      5      2      174      -      HG
-U133A:217524_X_AT      416      51      265      chr10      0      104279249
104279632      4      53,26,123,7,      151,206,232,358,      104279
249,104279302,104279500,104279625,
0      17      244      10      4      27      4      33      +      HG
-U133A:217625_X_AT      477      15      313      chr14      0      74280760
74281064      5      24,44,62,8,133,      15,43,90,171,180,      742807
60,74280795,74280841,74280909,74280931,
0      18      388      0      3      14      8      501      -      HG
-U133A:208120_X_AT      849      393      813      chr14      0      22175357
22176264      9      179,27,8,8,9,73,24,61,17,      36,215,245,261
,272,281,354,378,439,      22175357,22175690,22175718,22175737,22175751,221
75927,22176078,22176185,22176247,
0      21      256      1      2      25      2      7      -      HG
-U133A:217650_X_AT      671      222      525      chr14      0      93075015
93075300      3      166,52,60,      146,324,389,      93075015,93075
182,93075240,
0      22      294      0      4      302      5      937835      +      HG
-U133A:217501_AT      811      34      652      chr16      0      20791663
21729814      6      114,30,55,45,7,65,      34,148,186,244,293,587
,20791663,21055794,21433872,21433929,21549291,21729749,
--(Unix)-- HG-U133A.link.psl Tue Jul 19 2:42PM (Text Fill)--L??--Top-----
Mark set

```



# What annotations does Affymetrix supply?

As noted earlier, HG-U133A\_annot.csv contains 47.9MB worth of annotation information. What occupies all that space?

Probe Set ID	GeneChip Array	Species	Scientific Name	Annotation Date	Sequence Type	Sequence Source	Transcript ID (Array D)	Target Description	Representation
1007_s_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	Affymetrix Propriet	U48705mRNA	U48705 /FEU487	
1053_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	GenBank	M87338	M87338 /FIM873	
117_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	Affymetrix Propriet	X51757cgs	X51757 /FE X517	
121_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	GenBank	X69699	X69699 /FE X696	
1255_g_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	Affymetrix Propriet	L36861expanded_cc	L36861 /FEL368	
1294_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	GenBank	L13852	L13852 /FEL138	
1316_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	Affymetrix Propriet	X55005mRNA	X55005 /FE X550	
1320_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	Affymetrix Propriet	X79510cgs	X79510 /FE X795	
1405_i_at	Human Genome U133A	Homo sapiens		20-Jun-05	Exemplar sequence	GenBank	M21121	M21121 /FIM21	
22275	AFFX-r2-Hs28	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-r2-Hs28SrRNA	M11167.1 HAF	
22276	AFFX-r2-Hs28	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-r2-Hs28SrRNA	M11167.1 HAF	
22277	AFFX-r2-P1-cre	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-r2-P1-cre-3	Bacterioph	AF
22278	AFFX-r2-P1-cre	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-r2-P1-cre-5	Bacterioph	AF
22279	AFFX-ThrX-3	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-ThrX-3	B. subtilis	/CAF
22280	AFFX-ThrX-5	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-ThrX-5	B. subtilis	/CAF
22281	AFFX-ThrX-M	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-ThrX-M	B. subtilis	/CAF
22282	AFFX-TrpnX-3	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-TrpnX-3	B. subtilis	/CAF
22283	AFFX-TrpnX-5	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-TrpnX-5	B. subtilis	/CAF
22284	AFFX-TrpnX-M	Human Genome U133A	Homo sapiens	20-Jun-05	Control sequence	Affymetrix Propriet	AFFX-TrpnX-M	B. subtilis	/CAF

First, we note that the file seems to contain redundant copies of lots of information. Second, it has information on 22,283 probe sets, one per line, in 43 columns.



## Description of annotation columns

**Probe Set ID.** The unique identifier that describes an Affymetrix probe set. Also used in CEL files and CDF files.

**GeneChip Array.** The chip type on which the probe set appears. The same entry is repeated for all probe sets.

**Species Scientific Name.** The scientific name of the species whose gene sequences are on the array. The same information is repeated for all probe sets.

**Annotation Date.** The date when the annotations were last updated. The same information is repeated for all probe sets.

**Sequence Type.** The kind of sequence used in the design of the array: can be “Consensus”, “Control”, or “Exemplar”.

**Sequence Source.** Where did the design sequence come from? Usually “GenBank”, but rarely (only 81 times on the HG-U133A) from “Affymetrix Proprietary Database”.

**Transcript ID(Array Design).** An identifier into one of several unspecified databases indicating the designed target sequence.

**Target Description.** Long text string describing the target, formed by combining several other fields.

**Representative Public ID.** For non-control sequences, a GenBank identifier.

**Archival UniGene Cluster.** The UniGene cluster identifier from the sequence at the time the array was designed (in this case, from UniGene build 133).

**UniGene ID.** UniGene cluster identifier from the build of UniGene current at the time the annotations were updated.

**Genome Version.** The build of the human genome used for sequence alignments. The same information is repeated for all probe sets.

**Alignments.** Location of the target sequence along the human genome, in base pairs along the chromosome.

**Gene Title.** Official gene title (either from UniGene or Entrez Gene).

**Gene Symbol.** Official gene symbol (either from UniGene or Entrez Gene).

**Chromosomal Location.** Location of the gene in terms of

cytogenetic bands; e.g., 16p12.

**Unigene Cluster Type.** Either absent if not present in this build of UniGene (indicated by “—”), “est”, “full length”, or “est /// full length”.

**Ensembl.** The unique identifier of the target sequence in the Ensembl database.

**Entrez Gene.** The unique identifier of the target sequence in Entrez Gene (formerly LocusLink). Sequences with these identifiers tend to be better understood and more reliable than genes without them. The identifiers refer to genetic loci that have been mapped explicitly because of their connection to specific diseases or biological processes.

**SwissProt.** The SwissProt identifier of the protein product

produced by the gene corresponding to the target sequence.

**EC.** Yet another database identifier.

**OMIM.** The unique identifier associated to the target sequence gene in the Online Mendelian Inheritance in Man (OMIM) database, describing the ways in which the gene is known to be associated with genetic diseases.

**RefSeq Protein ID.** The GenBank identifier of the consensus sequence for the protein produced by the target sequence.

**RefSeq Transcript ID.** The GenBank identifiers of the consensus sequences for the mRNA's produced by the target gene. (Alternative splicing accounts for multiples.) In many cases, this coincides with the “Representative Public ID”.

**FlyBase.** Corresponding identifier in the drosophila database.

**AGI.** Unknown.

**WormBase.** Corresponding identifier in the *C. elegans* database.

**MGI Name.** Probably the identifier in the mouse database.

**RGD Name.** Probably the identifier in the rat database.

**SGD accession number.** The identifier in the saccharomyces database.

**Gene Ontology Biological Process.** List of identifiers for annotations of the target gene into the “biological process” section of GeneOntology. More about this later.

**Gene Ontology Cellular Component.** Similar.

**Gene Ontology Molecular Function.** Similar.

**Pathway.** List of pathways that the target sequence is involved in.

**Protein Families.** Families to which the protein belongs.

**Protein Domains.** Domains included in the protein.

**InterPro.** Another protein database.

**Trans Membrane.** Description of trans-membrane part of the protein, if known or if applicable.

**QTL.** Unknown.

**Annotation Description.** Text description of how the probe set was annotated.

**Annotation Transcript Cluster.** Unclear.

**Transcript Assignments.** Very long description of the annotations.

**Annotation Notes.** Additional comments.



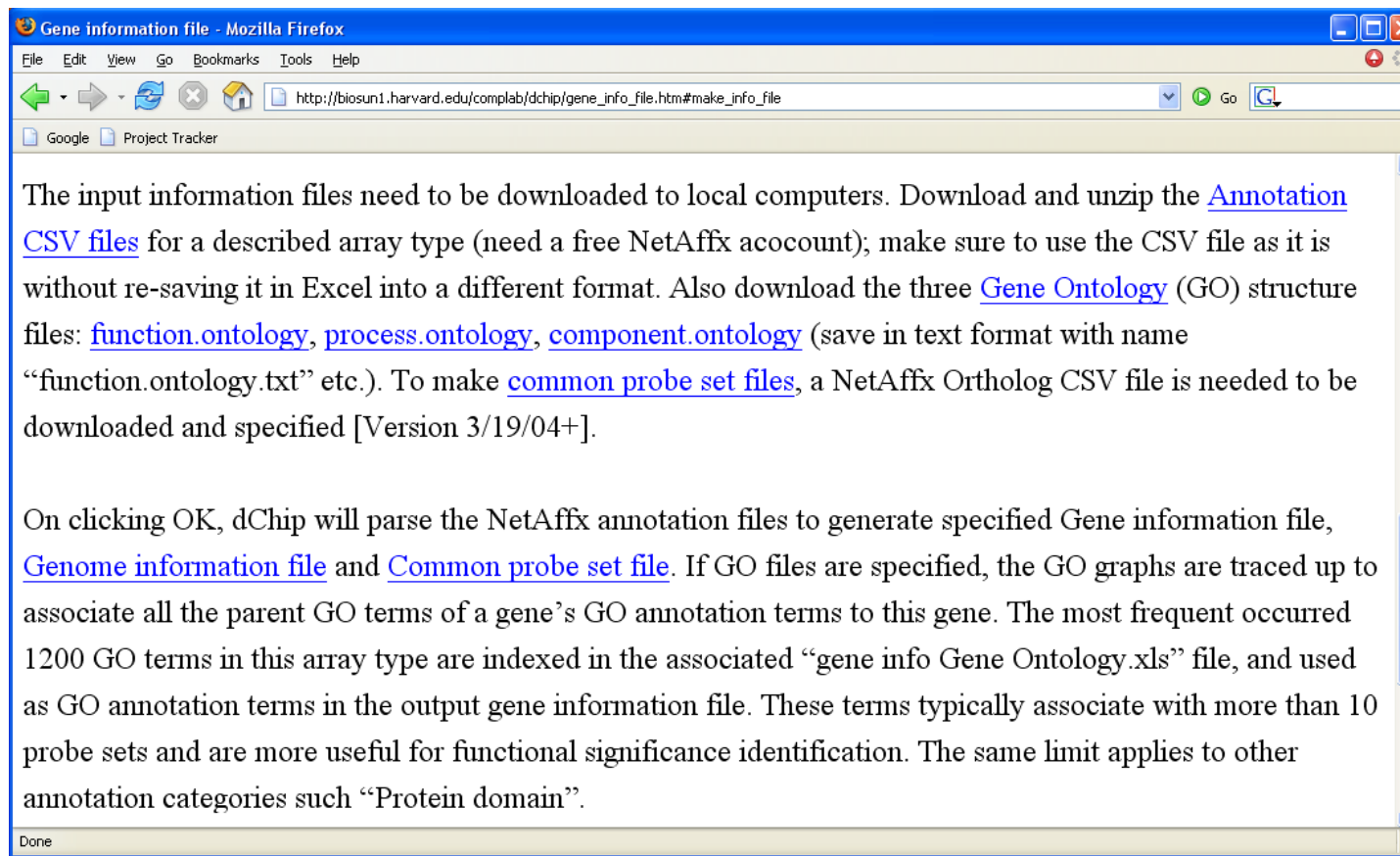
## Updating annotations in dChip

In order for dChip (or any other Affymetrix microarray analysis package) to use the updated annotations, you have to tell the software package where to get the information.

In the case of dChip, their online manual page tells you how to build new gene information and genome information files.

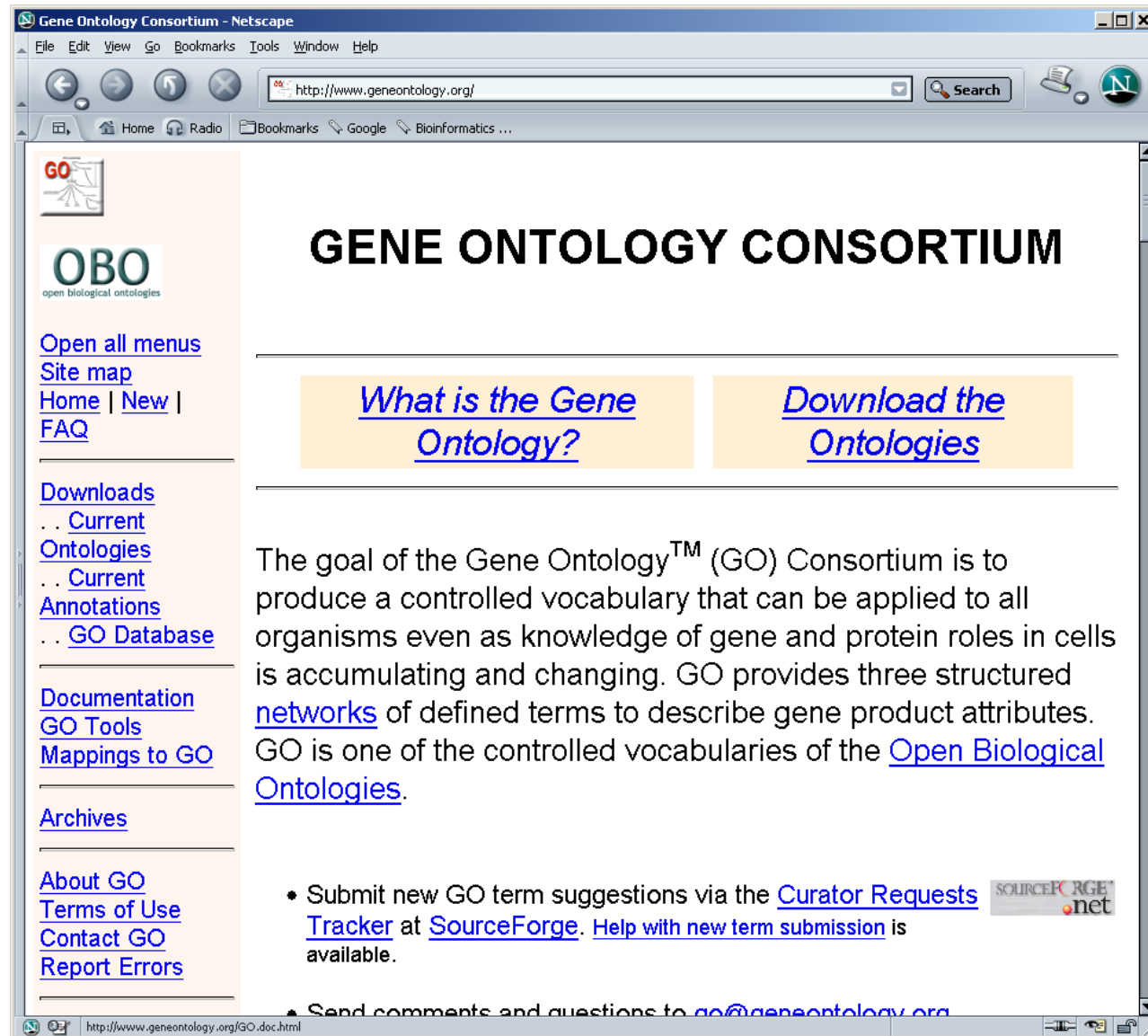
For many common chip types, the dChip web site contains up-to-date copies of these files. It's still useful to see where the data comes from how and how you can update your own versions.

# dChip Manual on Gene Information



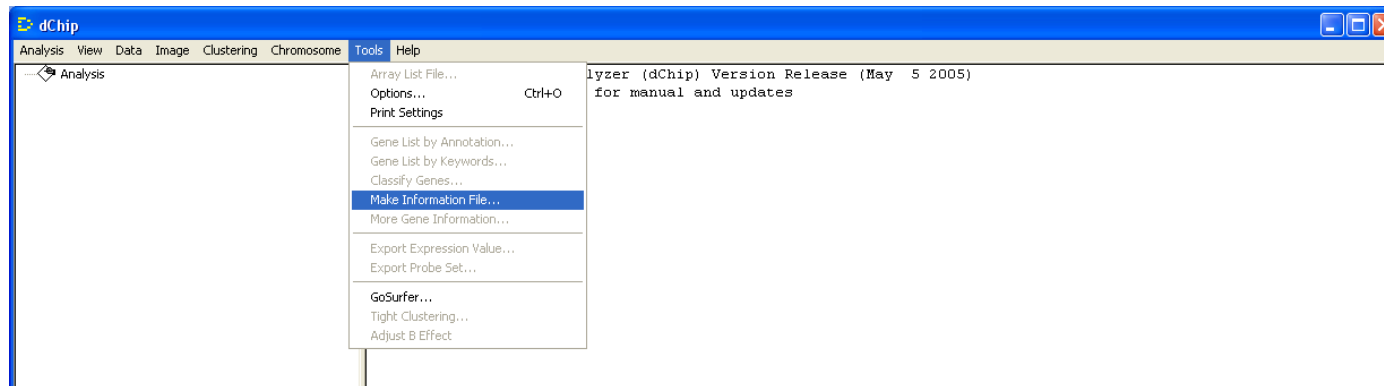
Requires the annotation CSV files from Affymetrix, along with three Gene Ontology files, which you can get from dChip or from the primary source.

# http://www.geneontology.org



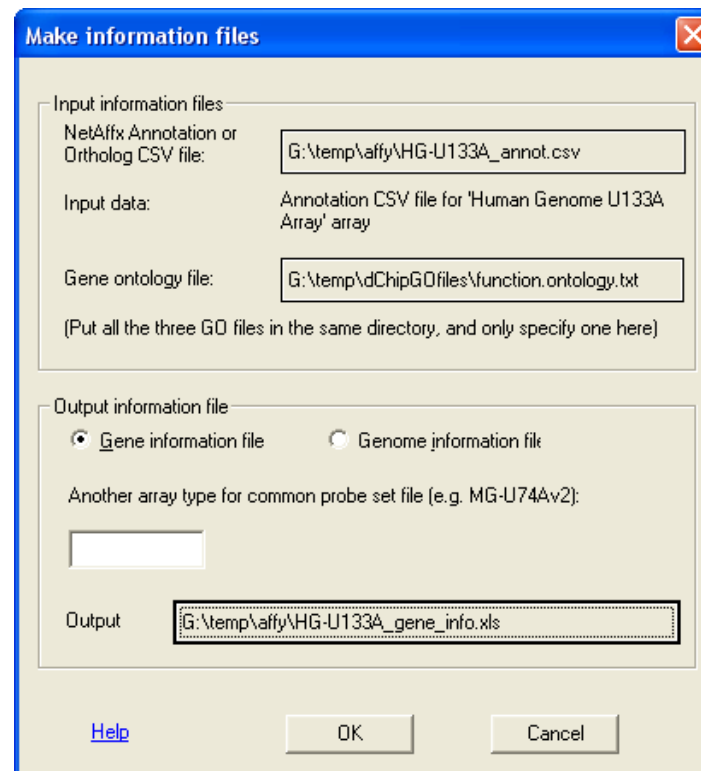
# Making the Gene Information file

1. Get and unzip the file containing the updated annotation CSV file from Affymetrix.
2. Get the three updated text files from GeneOntology.
3. Rename the three GeneOntology files, adding the “.txt” extension.
4. Use “Tools” – > “Make information file” in dChip.

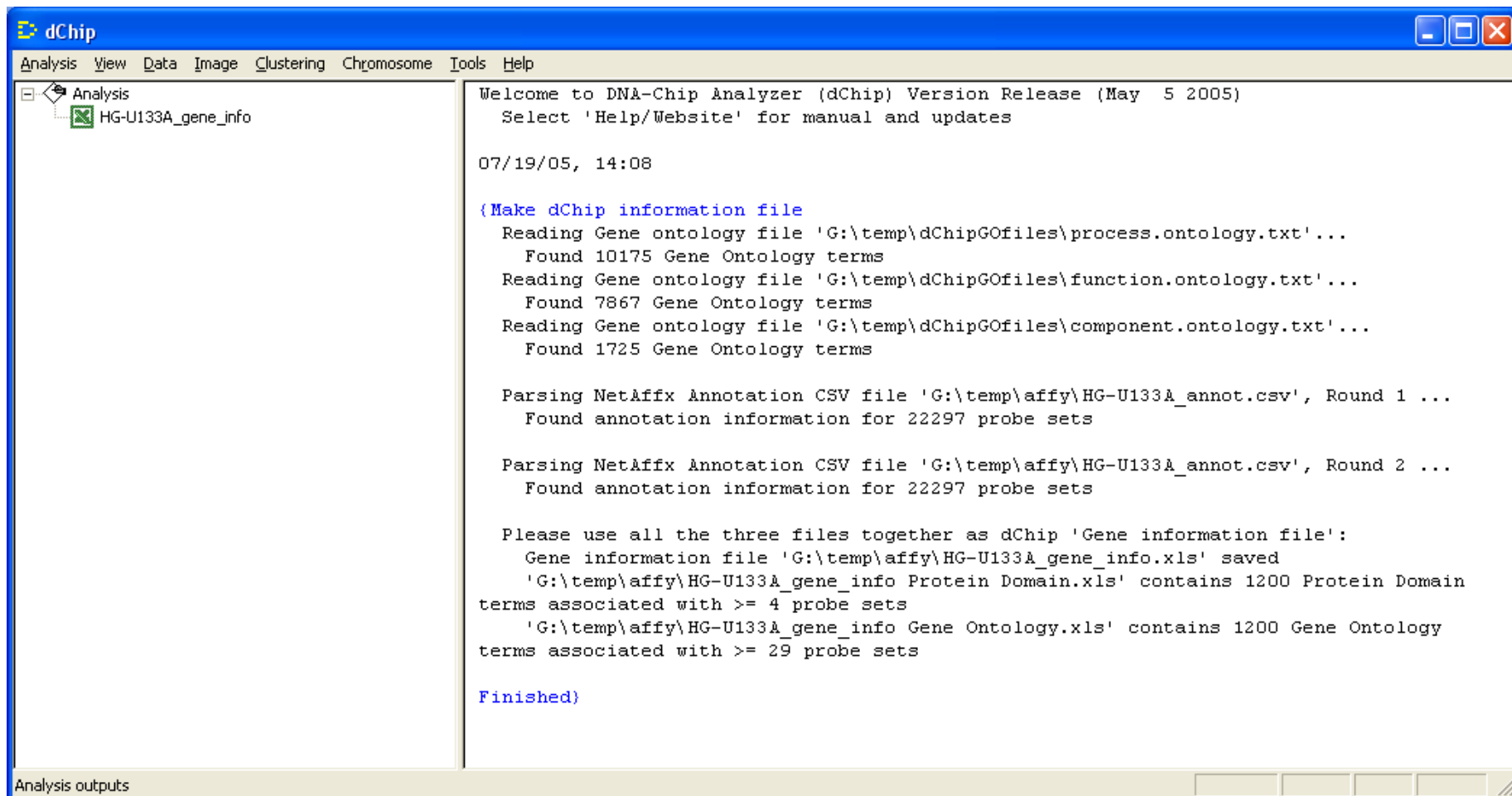


## Making the Gene Information file

Specify the locations of the CSV file and the GeneOntology files. Also say where you want the output sent. Note that I edited the default output file name to (i) start with the standard chip name and (2) use the underscore character as a separator.



# The Gene Information file



```
Welcome to DNA-Chip Analyzer (dChip) Version Release (May 5 2005)
Select 'Help/Website' for manual and updates

07/19/05, 14:08

{Make dChip information file
  Reading Gene ontology file 'G:\temp\dChipGOfiles\process.ontology.txt'...
  Found 10175 Gene Ontology terms
  Reading Gene ontology file 'G:\temp\dChipGOfiles\function.ontology.txt'...
  Found 7867 Gene Ontology terms
  Reading Gene ontology file 'G:\temp\dChipGOfiles\component.ontology.txt'...
  Found 1725 Gene Ontology terms

  Parsing NetAffx Annotation CSV file 'G:\temp\affy\HG-U133A_annot.csv', Round 1 ...
  Found annotation information for 22297 probe sets

  Parsing NetAffx Annotation CSV file 'G:\temp\affy\HG-U133A_annot.csv', Round 2 ...
  Found annotation information for 22297 probe sets

  Please use all the three files together as dChip 'Gene information file':
  Gene information file 'G:\temp\affy\HG-U133A_gene_info.xls' saved
  'G:\temp\affy\HG-U133A_gene_info Protein Domain.xls' contains 1200 Protein Domain
  terms associated with >= 4 probe sets
  'G:\temp\affy\HG-U133A_gene_info Gene Ontology.xls' contains 1200 Gene Ontology
  terms associated with >= 29 probe sets

Finished}
```

This step produces the three dChip annotation files that were described in Lecture 2.

# Making the Genome Information file

Using the same input files, you can also use dChip to create a “Genome information file”, which maps genes to specific positions along the genome.

The screenshot shows a Windows-style dialog box titled "Make information files". It has a blue title bar with a close button (X) in the top right corner. The dialog is divided into two main sections: "Input information files" and "Output information file".

**Input information files:**

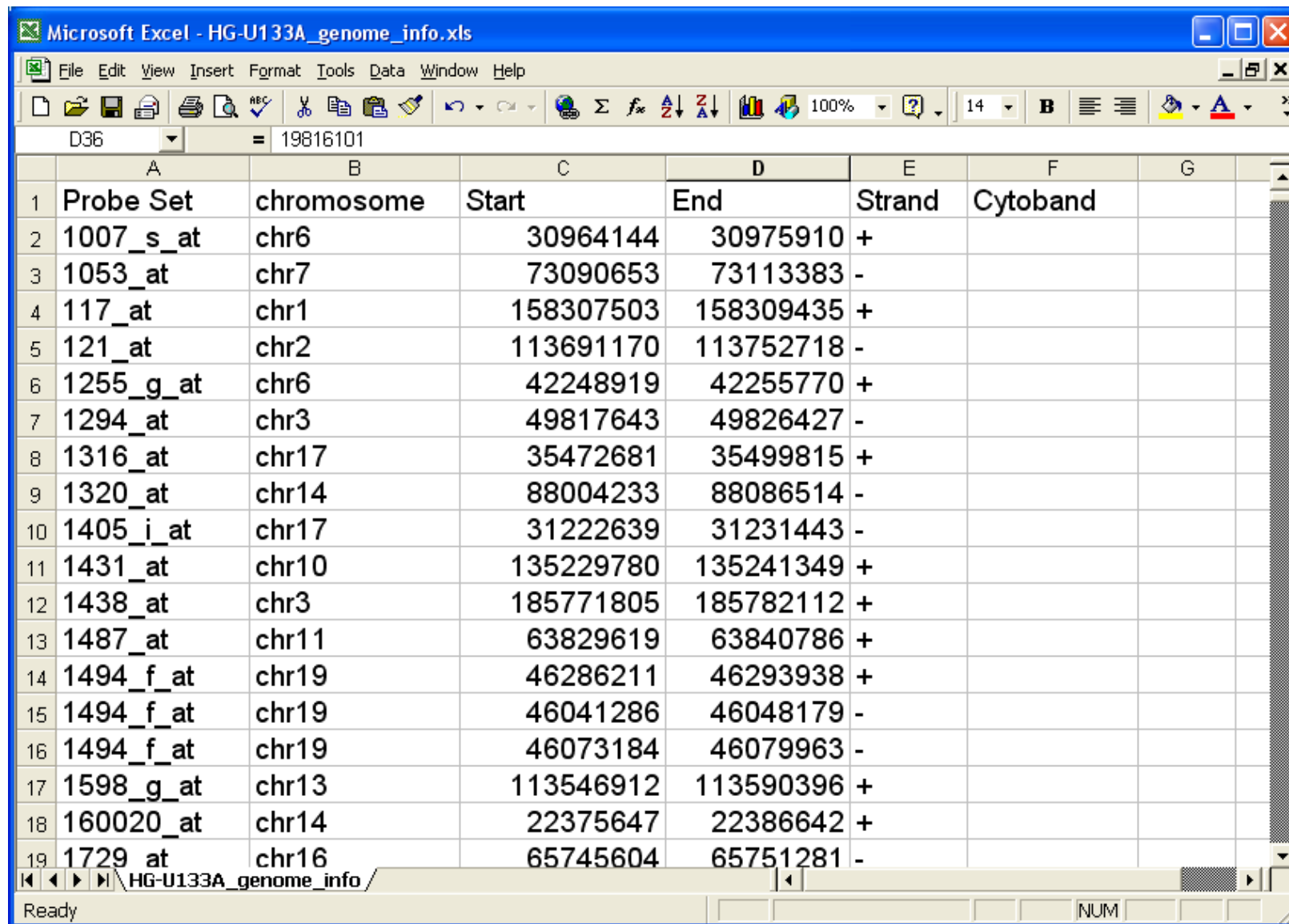
- NetAffx Annotation or Ortholog CSV file:** A text box containing the path "G:\temp\affy\HG-U133A\_annot.csv".
- Input data:** A label with the text "Annotation CSV file for 'Human Genome U133A Array' array".
- Gene ontology file:** A text box containing the path "G:\temp\dChipGOfiles\function.ontology.txt".
- Below these fields is a note: "(Put all the three GO files in the same directory, and only specify one here)".

**Output information file:**

- Two radio buttons are present: "Gene information file" (unselected) and "Genome information file" (selected).
- Below the radio buttons is a label: "Another array type for common probe set file (e.g. MG-U74Av2):" followed by an empty text box.
- Output:** A text box containing the path "G:\temp\affy\HG-U133A\_genome\_info.xls".

At the bottom of the dialog, there are three buttons: "Help" (with a blue underline), "OK", and "Cancel".

# The Genome Information file



Microsoft Excel - HG-U133A\_genome\_info.xls

File Edit View Insert Format Tools Data Window Help

D36 = 19816101

	A	B	C	D	E	F	G
	Probe Set	chromosome	Start	End	Strand	Cytoband	
1	1007_s_at	chr6	30964144	30975910	+		
2	1053_at	chr7	73090653	73113383	-		
3	117_at	chr1	158307503	158309435	+		
4	121_at	chr2	113691170	113752718	-		
5	1255_g_at	chr6	42248919	42255770	+		
6	1294_at	chr3	49817643	49826427	-		
7	1316_at	chr17	35472681	35499815	+		
8	1320_at	chr14	88004233	88086514	-		
9	1405_i_at	chr17	31222639	31231443	-		
10	1431_at	chr10	135229780	135241349	+		
11	1438_at	chr3	185771805	185782112	+		
12	1487_at	chr11	63829619	63840786	+		
13	1494_f_at	chr19	46286211	46293938	+		
14	1494_f_at	chr19	46041286	46048179	-		
15	1494_f_at	chr19	46073184	46079963	-		
16	1598_g_at	chr13	113546912	113590396	+		
17	160020_at	chr14	22375647	22386642	+		
18	1729_at	chr16	65745604	65751281	-		

Ready

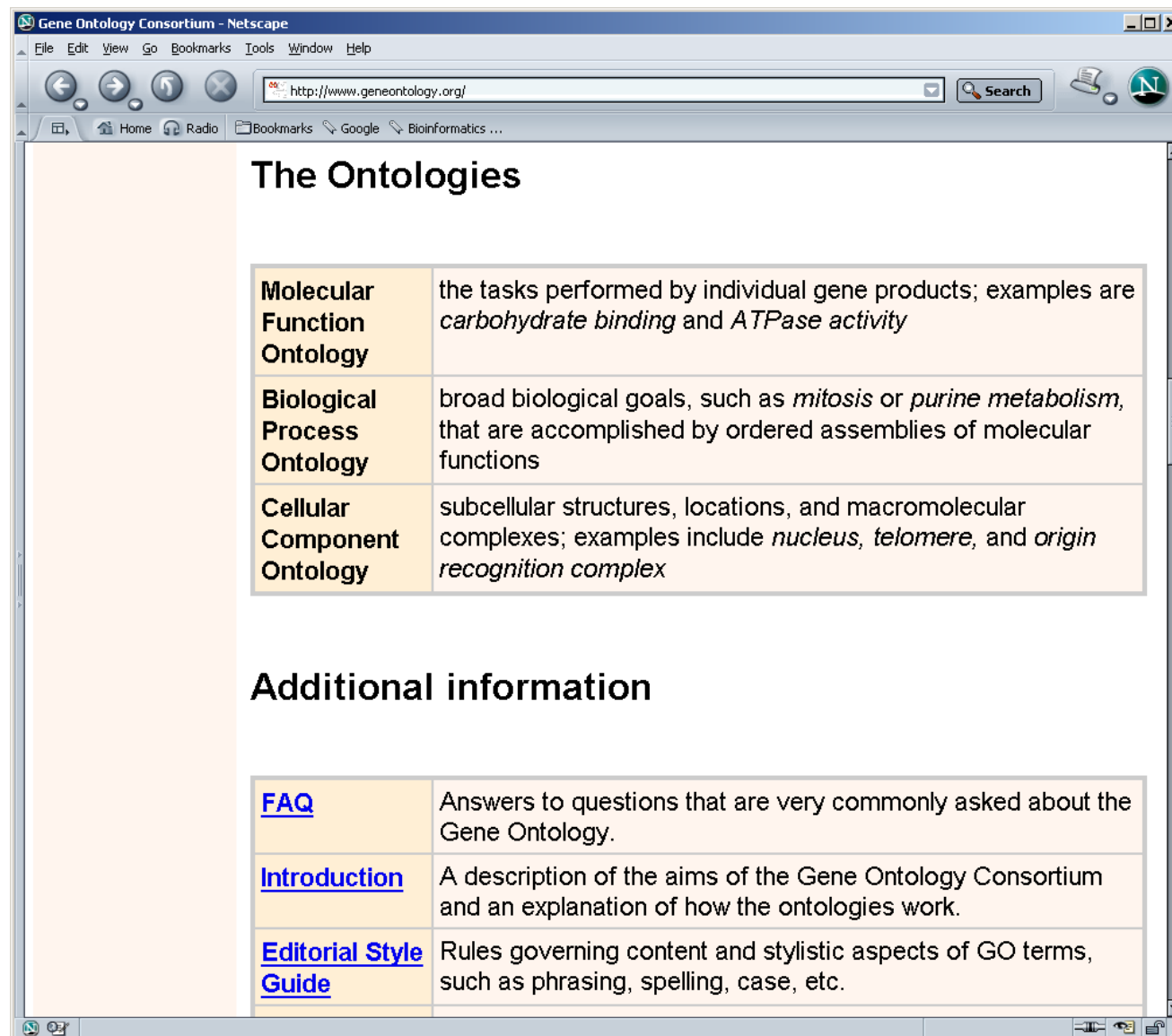


# What is GeneOntology?

GeneOntology uses controlled vocabularies to create a directed acyclic graph (DAG; a generalized tree) that describes the kinds of functions or properties that a gene might have. There are two parts to GeneOntology:

- Annotations, maintained in databases like Entrez Gene, that describe which genes actually have which functions.
- The DAG, maintained by the GeneOntology Consortium, that describes functions and relations between them:
  1. Biological process (what)
  2. Molecular function (how)
  3. Cellular component (where)

# GeneOntology: The top level



The screenshot shows a Netscape browser window titled "Gene Ontology Consortium - Netscape". The address bar displays "http://www.geneontology.org/". The page content is organized into two main sections: "The Ontologies" and "Additional information".

**The Ontologies**

<b>Molecular Function Ontology</b>	the tasks performed by individual gene products; examples are <i>carbohydrate binding</i> and <i>ATPase activity</i>
<b>Biological Process Ontology</b>	broad biological goals, such as <i>mitosis</i> or <i>purine metabolism</i> , that are accomplished by ordered assemblies of molecular functions
<b>Cellular Component Ontology</b>	subcellular structures, locations, and macromolecular complexes; examples include <i>nucleus</i> , <i>telomere</i> , and <i>origin recognition complex</i>

**Additional information**

<a href="#">FAQ</a>	Answers to questions that are very commonly asked about the Gene Ontology.
<a href="#">Introduction</a>	A description of the aims of the Gene Ontology Consortium and an explanation of how the ontologies work.
<a href="#">Editorial Style Guide</a>	Rules governing content and stylistic aspects of GO terms, such as phrasing, spelling, case, etc.

## GeneOntology annotations in Entrez Gene

You can find the GeneOntology annotations for individual genes in Entrez Gene. For genes with known functions, the Entrez Gene page will contain a section titled “GeneOntology”, which contains a list of the known functions for that gene.

Every GO annotation asserts that a specific gene has a specific function. As part of the design of GO, each assertion is itself annotated to explain the kinds of evidence the assertion is based on, as well as the organization or individual that supplied the annotation.

# GO annotations of the androgen receptor

**General gene information**

**GeneOntology**  
Provided by [GOA](#)

Category	Term	Evidence	PubMed
Function	<a href="#">androgen binding</a>	NAS	<a href="#">PubMed</a>
	<a href="#">androgen receptor activity</a>	IEA	<a href="#">PubMed</a>
	<a href="#">androgen receptor activity</a>	NAS	<a href="#">PubMed</a>
	<a href="#">metal ion binding</a>	IEA	
	<a href="#">protein dimerization activity</a>	NAS	<a href="#">PubMed</a>
	<a href="#">receptor activity</a>	IEA	
	<a href="#">steroid binding</a>	IEA	
	<a href="#">transcription factor activity</a>	IEA	
Process	<a href="#">cell proliferation</a>	NAS	<a href="#">PubMed</a>
	<a href="#">cell-cell signaling</a>	TAS	<a href="#">PubMed</a>
	<a href="#">prostate gland development</a>	NAS	<a href="#">PubMed</a>
	<a href="#">regulation of transcription, DNA-dependent</a>	IEA	
	<a href="#">sex differentiation</a>	NAS	<a href="#">PubMed</a>
	<a href="#">signal transduction</a>	TAS	<a href="#">PubMed</a>
	<a href="#">transcription</a>	IEA	
	<a href="#">transport</a>	TAS	<a href="#">PubMed</a>
Component	<a href="#">nucleus</a>	NR	
	<a href="#">nucleus</a>	IEA	

# http://www.ebi.ac.uk/GOA/

The screenshot shows the Gene Ontology Annotation (GOA) website at EBI, viewed in Netscape. The browser window title is "Gene Ontology Annotation @ EBI - Netscape". The address bar shows "http://www.ebi.ac.uk/GOA/". The page features a navigation menu with links: EBI Home, About EBI, Research, Services, Toolbox, Databases, Downloads, and Submissions. The main content area is titled "GOA DATABASE".

**GOA @EBI**

GOA is a project run by the European Bioinformatics Institute that aims to provide assignments of gene products to the Gene Ontology (GO) resource.

The goal of the Gene Ontology Consortium is to produce a dynamic controlled vocabulary that can be applied to all organisms, even while knowledge of gene and protein roles in cells is still accumulating and changing. In the GOA project, this vocabulary will be applied to a non-redundant set of proteins described in the UniProt Resource (Swiss-Prot/TrEMBL/PIR-PSD) and Ensembl databases that collectively provide complete proteomes for Homo sapiens and other organisms.

In the first stage of this project, GO assignments have been applied to a data set representing the human proteome by a combination of electronic mappings and manual curation. Subsequently, GO assignments for all complete and incomplete proteomes that exist in UniProt have been

**Access GOA**

Search GOA via SRS

Go

**GO**

The EBI's Gene Ontology Consortium pages. GO is an international consortium of scientists with the editorial office based at the EBI.

**UniProt/Swiss-Prot**

# GO browsing

The screenshot shows a Netscape browser window titled "AmiGO! Your friend in the Gene Ontology". The address bar contains the URL `http://www.godatabase.org/cgi-bin/amigo/go.cgi?view=details&depth=1&query=GO:0005497`. The page content is as follows:

**AmiGO**

**androgen binding**

**Accession:** GO:0005497  
**Aspect:** function  
**Synonyms:** None  
**Definition:**  
Interacting selectively with any androgen, male sex hormones.

**Term Lineage**

- all : all ( 153306 )
- ① GO:0003674 : molecular\_function ( 103037 )
- ① GO:0005488 : binding ( 29138 )
- ① GO:0042562 : hormone binding ( 36 )
- ① **GO:0005497 : androgen binding ( 7 )**
- ① GO:0005496 : steroid binding ( 83 )
- ① **GO:0005497 : androgen binding ( 7 )**

[Graphical View](#)

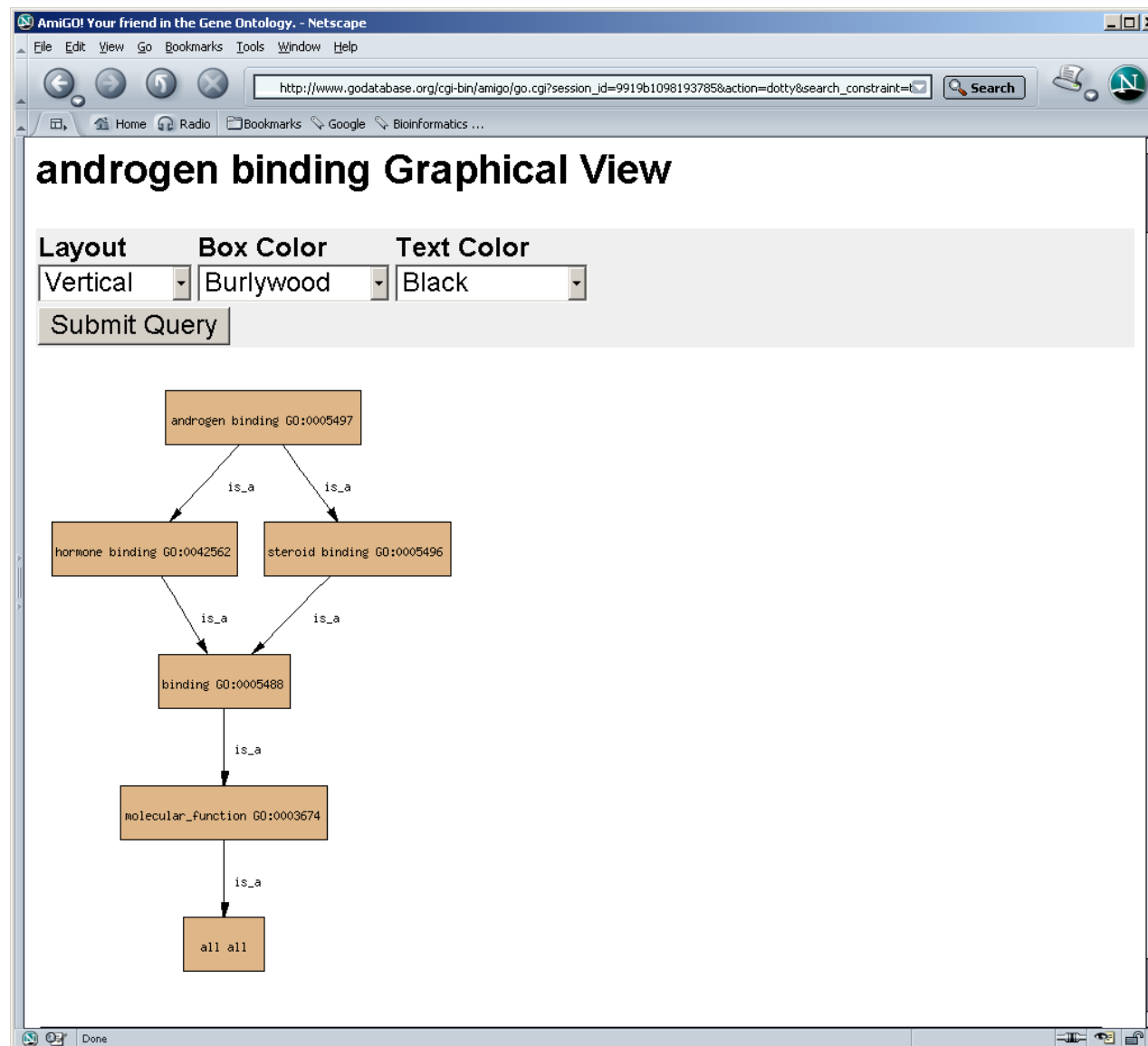
**External References**

None.

**Direct Gene Product Associations** Get ALL associations here:

Direct Associations

# GO browsing



## Edges are relationships

Edges in the DAG represent two kinds of relationships:

**is\_a** : Used when the child node is a special case of the parent node. For example, `hormone binding` is\_a kind of `binding`.

**part\_of** : Used when the child node is a component of the parent node. For example, a `membrane` is part\_of a `cell`

Genes may be annotated into different levels of the hierarchy, depending on how detailed the evidence is. In general, a gene not only has the function corresponding to the node with direct annotation, but also has every property at parent nodes up through the hierarchy.



# GO annotations of the androgen receptor

Gene - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=gene&cmd=Retrieve&dopt=Graphics&list\_uids=367

Google Project Tracker

Gene information file Gene

► General gene information ?

**GeneOntology**

Provided by [GOA](#)

Function	Evidence
<a href="#">androgen binding</a>	NAS <a href="#">PubMed</a>
<a href="#">androgen receptor activity</a>	IEA <a href="#">PubMed</a>
<a href="#">androgen receptor activity</a>	NAS <a href="#">PubMed</a>
<a href="#">metal ion binding</a>	IEA
<a href="#">protein dimerization activity</a>	NAS <a href="#">PubMed</a>
<a href="#">receptor activity</a>	IEA
<a href="#">steroid binding</a>	IEA
<a href="#">transcription factor activity</a>	IEA

Process	Evidence
<a href="#">cell proliferation</a>	NAS <a href="#">PubMed</a>
<a href="#">cell-cell signaling</a>	TAS <a href="#">PubMed</a>
<a href="#">prostate gland development</a>	NAS <a href="#">PubMed</a>
<a href="#">regulation of transcription, DNA-dependent</a>	IEA
<a href="#">sex differentiation</a>	NAS <a href="#">PubMed</a>
<a href="#">signal transduction</a>	TAS <a href="#">PubMed</a>
<a href="#">transcription</a>	IEA
<a href="#">transport</a>	TAS <a href="#">PubMed</a>

Component	Evidence
<a href="#">nucleus</a>	NR
<a href="#">nucleus</a>	IEA

Done

## GeneOntology: Evidence Codes

**IDA** : inferred from direct assay; indicates that the annotation is based on a paper describing an experiment that directly tested this function for this gene

**TAS** : traceable author statement; based on a review article or textbook that includes references to the original experiments

**IMP** : inferred from mutant phenotype; based on experiments involving mutations, knockouts, antisense, etc.

**IPI** : inferred from physical interaction; based on assays (like co-immunoprecipitation) that demonstrate physical interactions between the gene in question and other gene products

**IGI** : inferred from genetic interaction; based on experiments (such as synthetic lethals, suppressors, functional complementation) that show a genetic interaction between the gene in question and another gene

**ISS** : inferred from sequence or structure similarity; based on BLAST results that have been reviewed for accuracy by a curator

**IEP** : inferred from expression pattern; based on Northern, Westerns, or microarray experiments that reveal information about the timing or location of expression

**NAS** : non-traceable author statement; statements in papers (abstract, introduction, discussion) that a curator cannot trace to another publication

**IEA** : inferred from electronic annotation; based on sequence similarity searches or database records that have not been reviewed by a curator

**IC** : inferred by curator; even though no direct evidence is available, the property can reasonably be inferred by the curator. For example, it is reasonable to infer from direct evidence of “transcription factor activity” that the gene product is found in the nucleus

**ND** : no biological data available; only used for annotations to “unknown”

**NR** : not recorded; used only for annotations created before curators started adding evidence codes

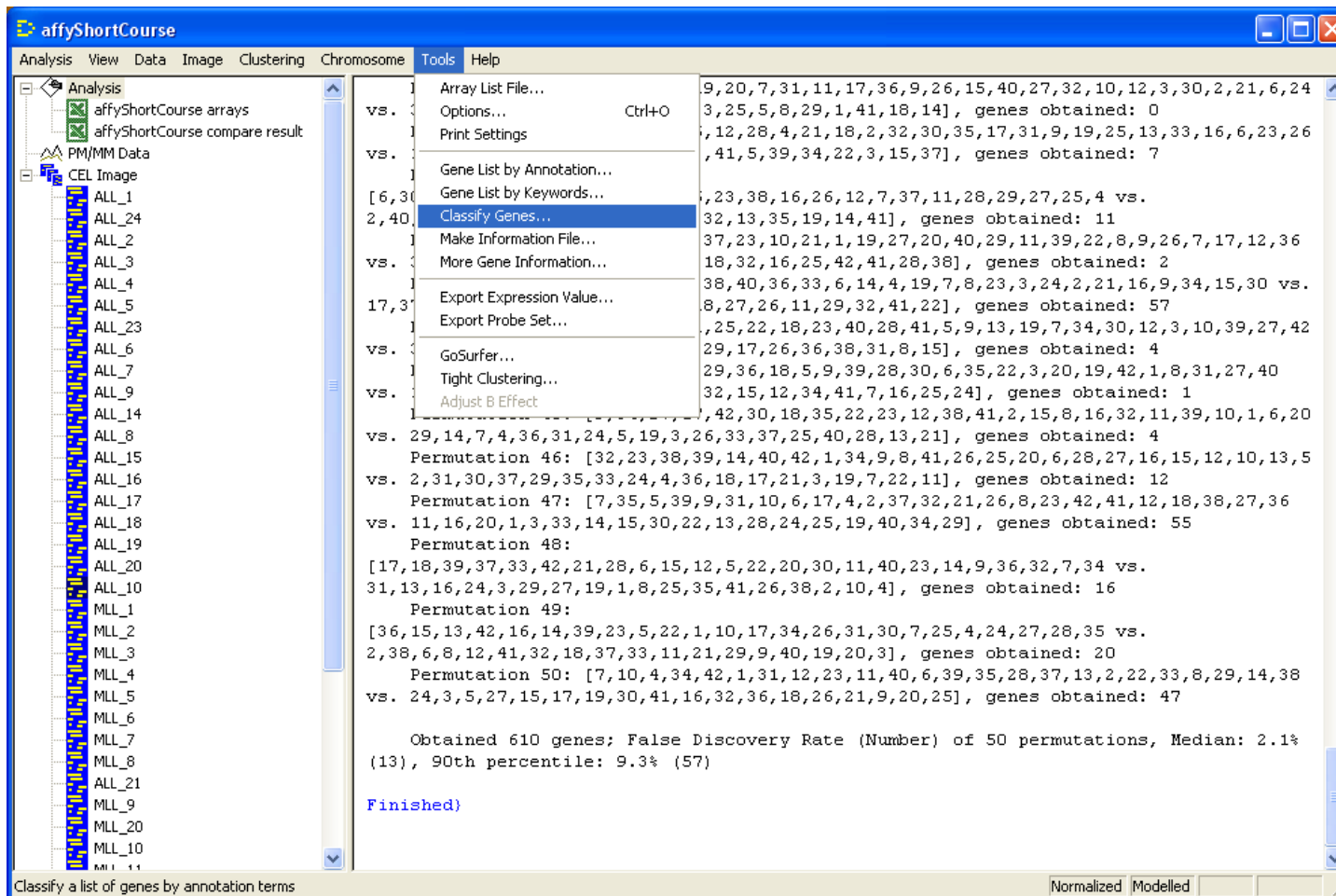
## Quality of evidence

The evidence codes fall into a rough hierarchy indicating how strongly the annotation of function should be believed.

1. IDA, TAS
2. IMP, IPI, IGI
3. ISS, IEP
4. NAS
5. IEA
6. IC

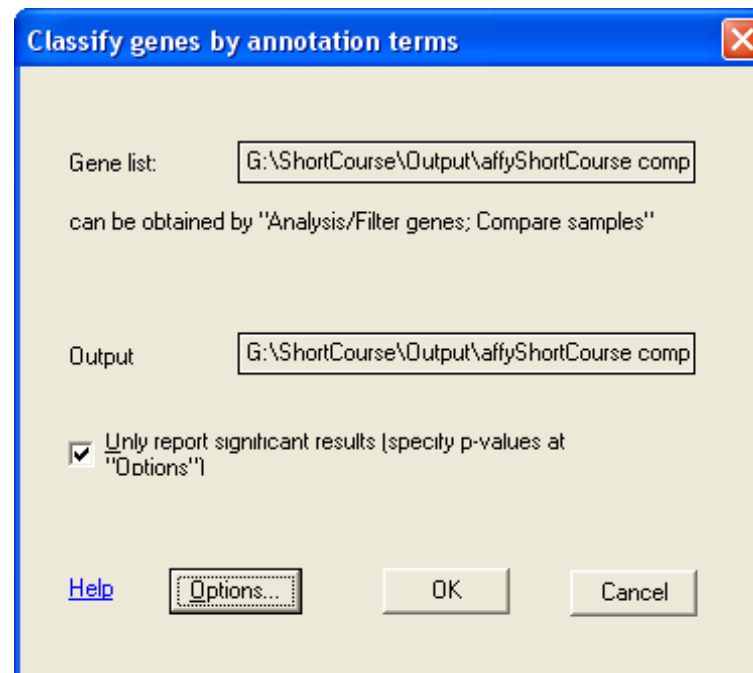
# Using GeneOntology in dChip

After running a sample comparison to find interesting genes, use the menu item “Tools” – > “Classify genes”.



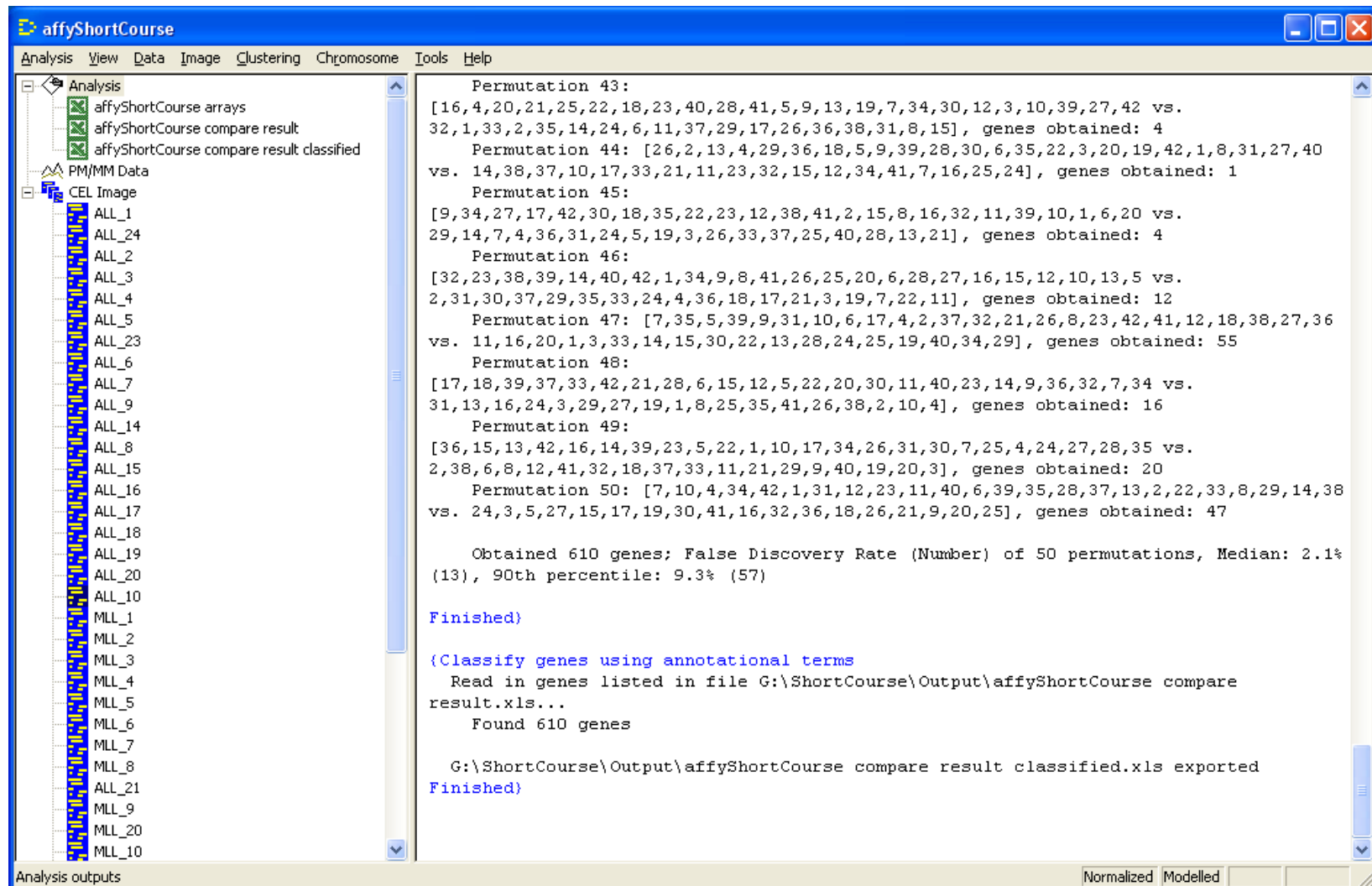
# Using GeneOntology in dChip

For the gene list file, select the “compare result” file produced previously. It’s also a good idea to check the box telling dChip only to report the groups with significant p-values.



# Using GeneOntology in dChip

The results are available in a few seconds.





# What do the results look like?

Microsoft Excel - affyShortCourse compare result classified.xls

File Edit View Insert Format Tools Data Window Help

A1 = probe set

	A	B	AA	AB	AU	AV	AW	BA
	probe set	gene	baseline mean	baseline	experiment mean	experiment	fold change	filtered
3	Found 21 Gene Ontology "protein tyrosine kinase" genes in a list with 391 annotated genes (all: 157/7685, PValue: 0.000042) *****							
4	40936_at	cysteine-rich motor neuron 1	7994	564	5144	612	-1.55	*
5	1485_at	EphA7	243	28	133	14	-1.83	*
6	2057_g_at	fibroblast growth factor receptor 1 (fms-re	5421	430	2717	100	-2	*
7	1964_g_at	fms-related tyrosine kinase 1 (vascular en	1555	167	982	51	-1.58	*
8	1545_g_at	fms-related tyrosine kinase 1 (vascular en	745	85	471	16	-1.58	*
9	34583_at	fms-related tyrosine kinase 3	9522	1513	16788	784	1.76	*
10	1065_at	fms-related tyrosine kinase 3	8414	1696	15615	933	1.86	*
11	40480_s_at	FYN oncogene related to SRC, FGR, YES	5038	514	3304	326	-1.52	*
12	34877_at	Janus kinase 1 (a protein tyrosine kinase)	15776	843	10823	834	-1.46	*
13	41594_at	Janus kinase 1 (a protein tyrosine kinase)	6687	345	4360	301	-1.53	*
14	1457_at	Janus kinase 1 (a protein tyrosine kinase)	3098	197	1886	177	-1.64	*
15	33238_at	lymphocyte-specific protein tyrosine kinas	3794	572	1936	258	-1.96	*
16	1988_at	platelet-derived growth factor receptor, alp	14547	602	10367	351	-1.4	*
17	36117_at	PTK2 protein tyrosine kinase 2	3730	242	2613	117	-1.43	*
18	37756_at	RYK receptor-like tyrosine kinase	1155	129	399	48	-2.89	*
19	539_at	RYK receptor-like tyrosine kinase	2294	107	1665	48	-1.38	*
20	572_at	TTK protein kinase	1309	128	792	76	-1.65	*
21	1674_at	v-yes-1 Yamaguchi sarcoma viral oncogen	1438	283	496	32	-2.9	*
22	32616_at	v-yes-1 Yamaguchi sarcoma viral related c	3247	219	4842	498	1.49	*
23	2024_s_at	v-yes-1 Yamaguchi sarcoma viral related c	1913	141	2960	322	1.55	*
24	1402_at	v-yes-1 Yamaguchi sarcoma viral related c	4141	289	6292	581	1.52	*
26	Found 12 Gene Ontology "protein tyrosine phosphatase" genes in a list with 391 annotated genes (all: 81/7685, PValue: 0.000740) *							
27	32916_at	protein tyrosine phosphatase, receptor typ	6814	927	3050	454	-2.23	*
28	31892_at	protein tyrosine phosphatase, receptor typ	801	336	151	10	-5.32	*

affyShortCourse compare result

Ready

## Interpreting the Results

Each group of entries in the results file is introduced by a line like:

```
Found 21 Gene Ontology "protein tyrosine kinase"  
genes in a list with 391 annotated genes (all:  
157/7685, PValue: 0.000042) *****
```

The part within quotation marks is the name of the GeneOntology category that was found to be significantly overrepresented among the differentially expressed genes.

## The numbers tell us:

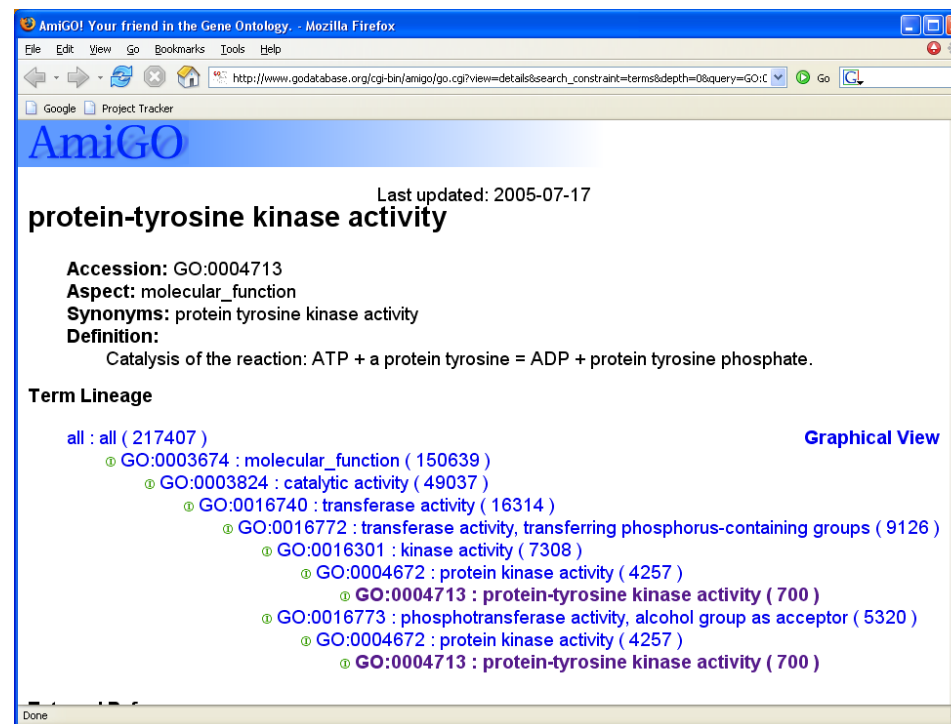
1. There were 7685 probe sets on the array with some kind of GeneOntology annotation.
2. There were 391 probe sets selected as differentially expressed that had some kind of GeneOntology annotation.
3. Of all the annotated probe sets, 157 had the “protein trosine kinase” function.
4. Of the selected annotated probe sets, 21 had the “protein tyrosine kinase” function.

The p-value comes from modeling the data using a hypergeometric distribution, which means it is the same value produced by Fisher’s Exact Test on a  $2 \times 2$  contingency table.

# What's wrong with the results?

First, the p-values have not been adjusted for multiple testing.

Second, we cannot tell if the software has accounted for the fact that the GeneOntology categories form a DAG. In particular, a gene with “protein tyrosine kinase” activity also inherits every annotation above it in the DAG.



# What's wrong with the results?

Third, by working with probe sets instead of genes, the counts are wrong.

Microsoft Excel - affyShortCourse compare result classified.xls

1	probe set	gene	AA	AB	AU	AV	AW	BA
			baseline mean	baseline	experiment mean	experiment	fold change	filtered
3	Found 21 Gene Ontology "protein tyrosine kinase" genes in a list with 391 annotated genes (all: 157/7685, PValue: 0.000042) *****							
4	40936_at	cysteine-rich motor neuron 1	7994	564	5144	612	-1.55 *	
5	1485_at	EphA7	243	28	133	14	-1.83 *	
6	2057_g_at	fibroblast growth factor receptor 1 (fms-re	5421	430	2717	100	-2 *	
7	1964_g_at	fms-related tyrosine kinase 1 (vascular en	1555	167	982	51	-1.58 *	
8	1545_g_at	fms-related tyrosine kinase 1 (vascular en	745	85	471	16	-1.58 *	
9	34583_at	fms-related tyrosine kinase 3	9522	1513	16788	784	1.76 *	
10	1065_at	fms-related tyrosine kinase 3	8414	1696	15615	933	1.86 *	
11	40480_s_at	FYN oncogene related to SRC, FGR, YES	5038	514	3304	326	-1.52 *	
12	34877_at	Janus kinase 1 (a protein tyrosine kinase)	15776	843	10823	834	-1.46 *	
13	41594_at	Janus kinase 1 (a protein tyrosine kinase)	6687	345	4360	301	-1.53 *	
14	1457_at	Janus kinase 1 (a protein tyrosine kinase)	3098	197	1886	177	-1.64 *	
15	33238_at	lymphocyte-specific protein tyrosine kinas	3794	572	1936	258	-1.96 *	
16	1988_at	platelet-derived growth factor receptor, alp	14547	602	10367	351	-1.4 *	
17	36117_at	PTK2 protein tyrosine kinase 2	3730	242	2613	117	-1.43 *	
18	37756_at	RYK receptor-like tyrosine kinase	1155	129	399	48	-2.89 *	
19	539_at	RYK receptor-like tyrosine kinase	2294	107	1665	48	-1.38 *	
20	572_at	TTK protein kinase	1309	128	792	76	-1.65 *	
21	1674_at	v-yes-1 Yamaguchi sarcoma viral oncogen	1438	283	496	32	-2.9 *	
22	32616_at	v-yes-1 Yamaguchi sarcoma viral related c	3247	219	4842	498	1.49 *	
23	2024_s_at	v-yes-1 Yamaguchi sarcoma viral related c	1913	141	2960	322	1.55 *	
24	1402_at	v-yes-1 Yamaguchi sarcoma viral related c	4141	289	6292	581	1.52 *	
26	Found 12 Gene Ontology "protein tyrosine phosphatase" genes in a list with 391 annotated genes (all: 81/7685, PValue: 0.000740) *							
27	32916_at	protein tyrosine phosphatase, receptor typ	6814	927	3050	454	-2.23 *	
28	31892_at	protein tyrosine phosphatase, receptor typ	801	336	151	10	-5.32 *	

# What alternatives are there?



The screenshot shows the GoMiner Home Page in a Netscape browser window. The browser's address bar displays <http://discover.nci.nih.gov/gominer/>. The page features a blue header with the GoMiner logo (Build:138) and the text "Genomics and Bioinformatics Group LMP, NCI, NIH" and "Medical Informatics & Bioimaging Lab, BME GA Tech/ Emory University". Below the header is a navigation bar with buttons for Home, Getting Started, Requirements, Installation, Command Line, Database, FAQ, Citing, Mirrors, and Credits. The main content area explains that GoMiner is a tool for biological interpretation of 'omic' data, including gene expression microarrays. It states that omic experiments often generate lists of dozens or hundreds of genes that differ in expression between samples, raising the question: *What does it all mean biologically?* To answer this, GoMiner leverages the [Gene Ontology \(GO\)](#) to identify biological processes, functions, and components. Instead of analyzing microarray results with a gene-by-gene approach, GoMiner classifies genes into biologically coherent categories and assesses these categories. The insights gained through GoMiner can generate hypotheses to guide additional research. A section titled "To get started using GoMiner" provides a list of steps: read the [Instructions](#), verify that the environment satisfies the system requirements, download the GoMiner program file, [gominer.jar](#), and read the [Quick Start](#) and try out the [sample files](#).

GoMiner Home Page - Netscape

File Edit View Go Bookmarks Tools Window Help

<http://discover.nci.nih.gov/gominer/> Search

Home Radio Bookmarks Google Bioinformatics ...

**GoMINER** Build:138

Genomics and Bioinformatics Group  
LMP, NCI, NIH

Medical Informatics & Bioimaging Lab,  
BME GA Tech/  
Emory University

Home Getting Started Requirements Installation Command Line Database FAQ Citing Mirrors Credits

GoMiner is a tool for biological interpretation of 'omic' data – including data from gene expression microarrays. Omic experiments often generate lists of dozens or hundreds of genes that differ in expression between samples, raising the question

*What does it all mean biologically?*

To answer this question, GoMiner leverages the [Gene Ontology \(GO\)](#) to identify the biological processes, functions and components represented in these lists. Instead of analyzing microarray results with a gene-by-gene approach, GoMiner classifies the genes into biologically coherent categories and assesses these categories. The insights gained through GoMiner can generate hypotheses to guide additional research.

**To get started using GoMiner**

- Read the [Instructions](#), and verify that your environment satisfies the system requirements
- Download the GoMiner program file, [gominer.jar](#)
- Read the [Quick Start](#) and try out the [sample files](#)

<http://discover.nci.nih.gov/gominer/index.jsp>

# http://discover.nci.nih.gov/gominer



# GoMiner: Getting Started

You need a machine with

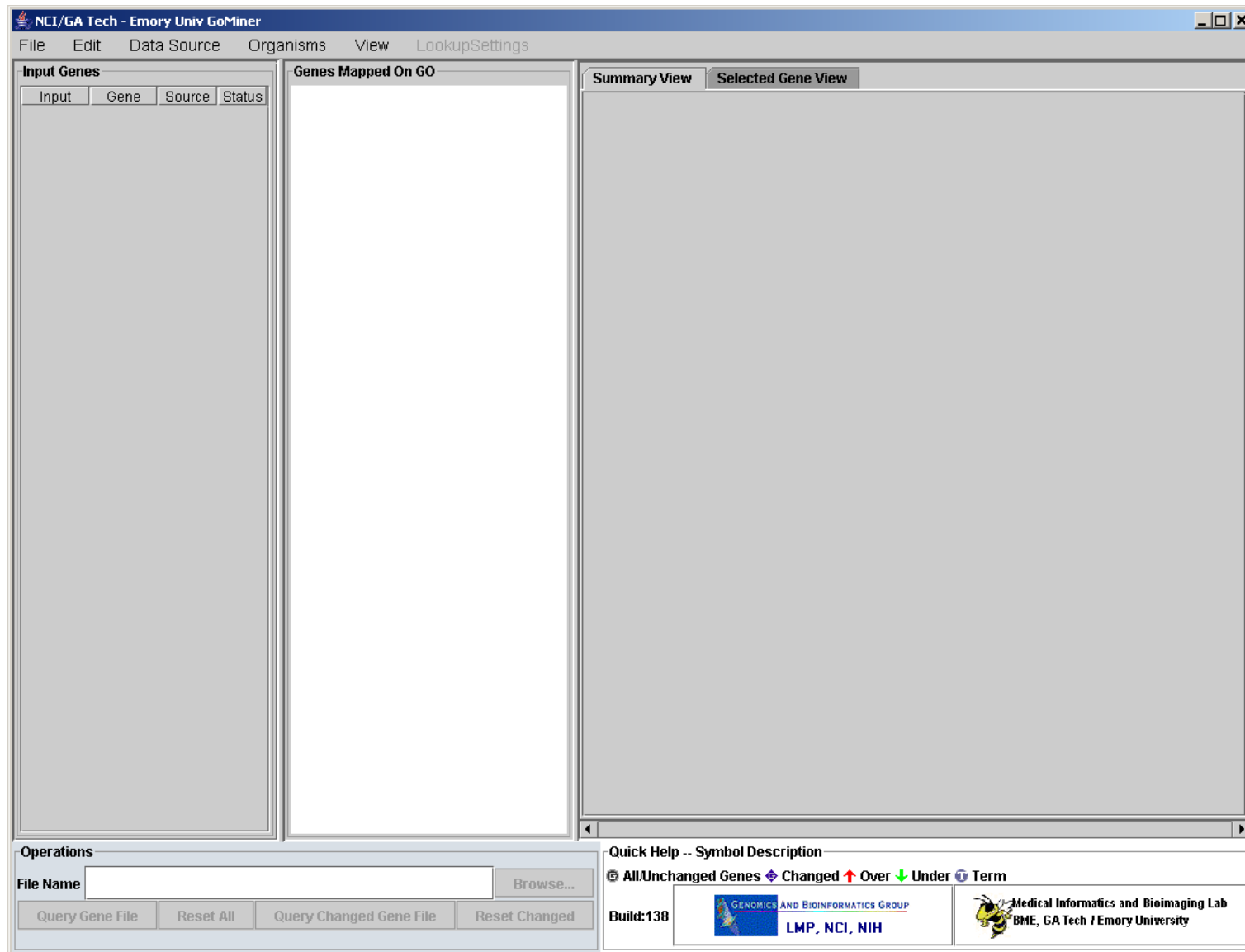
- Java 1.3 or higher
- Windows 98 or higher, Mac OS X or higher, Solaris, Linux, or FreeBSD
- High-speed internet access

Download the GoMiner Java code, install it, and double-click on it to start the program.

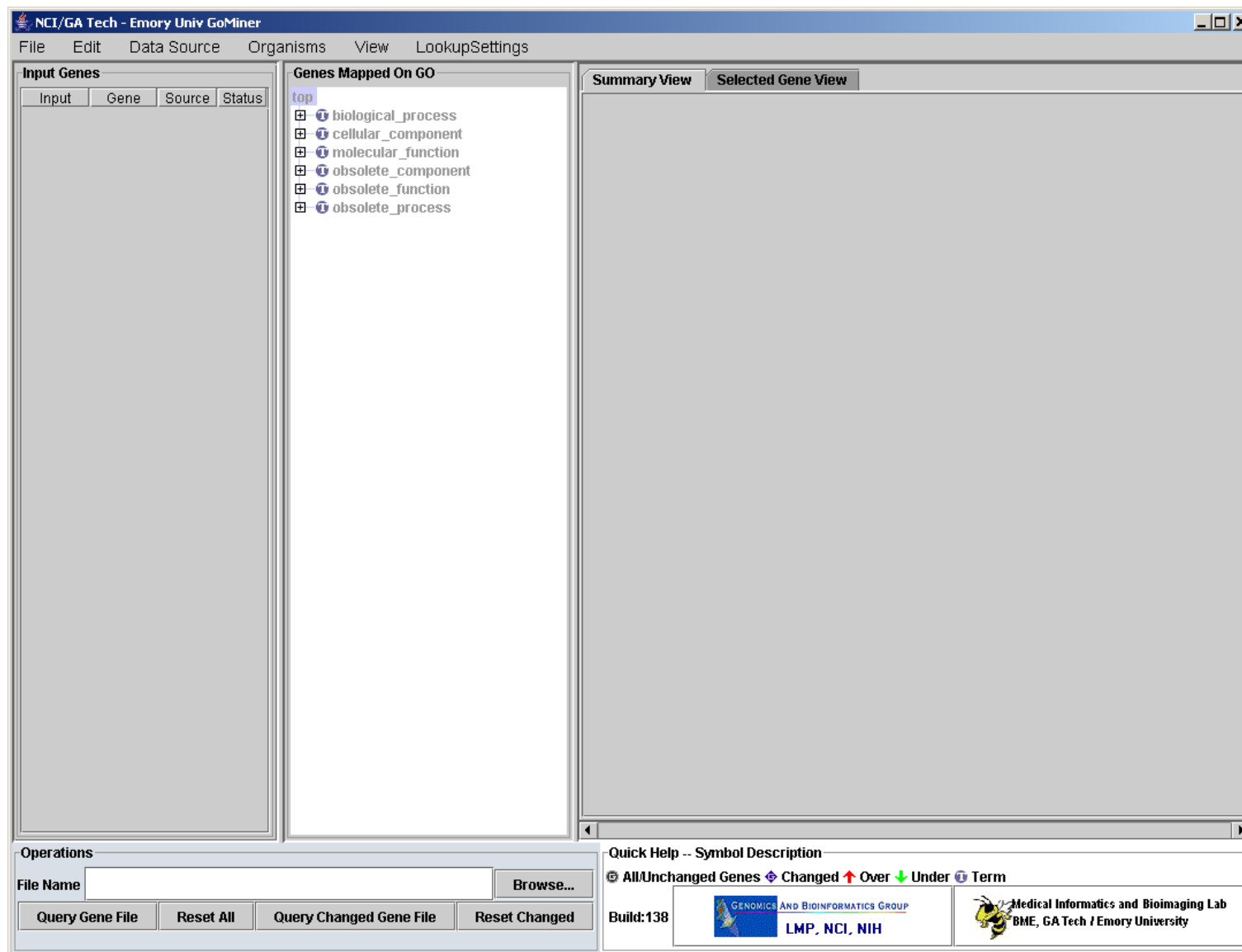
Then go to “File” — > “Load GO Terms” and click “OK”. Wait a few minutes while the program loads the GeneOntology information from the NCI.



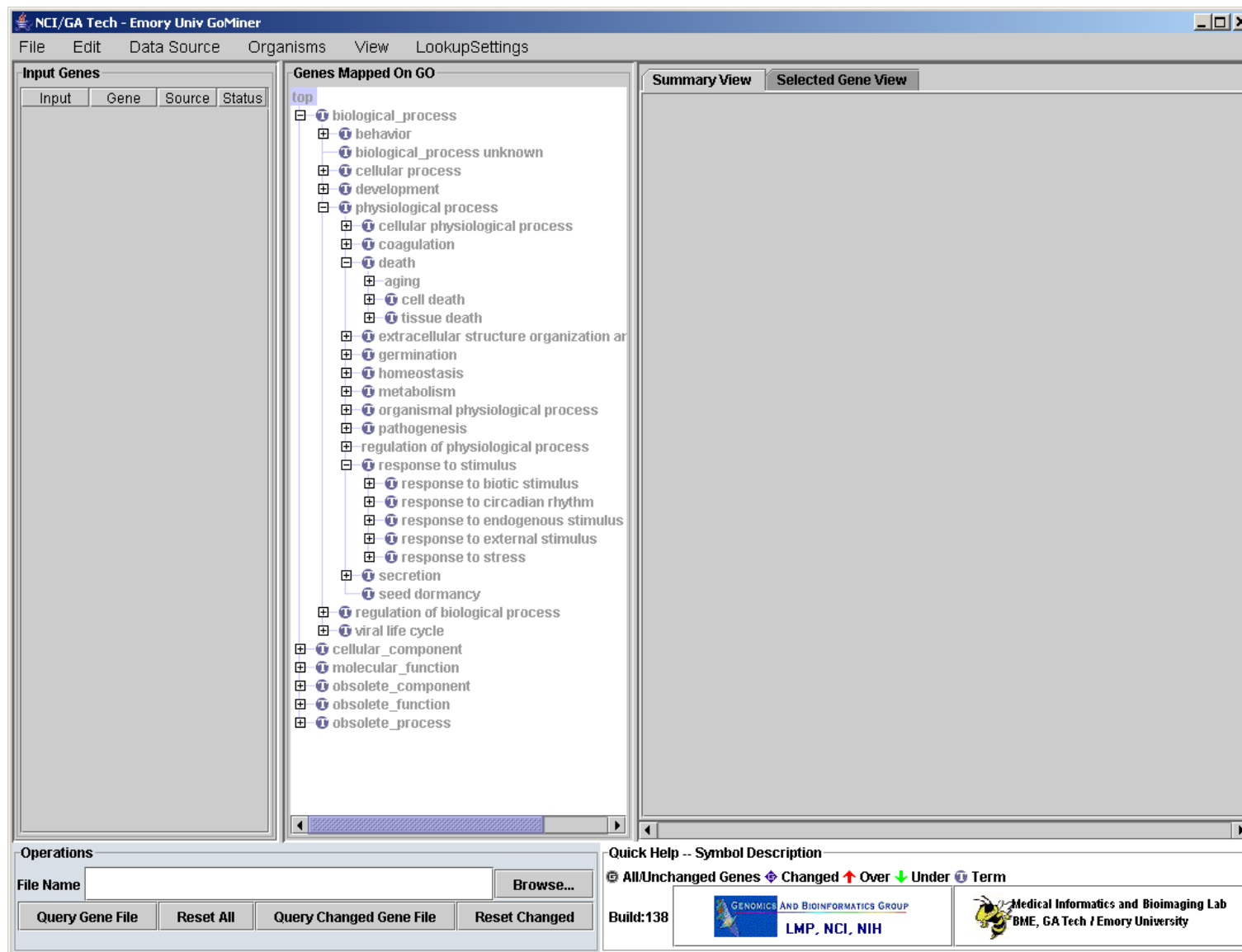
# GoMiner Start



# GoMiner: GO terms loaded



# GoMiner as GO browser



## Getting array data into GoMiner

1. Go to “Data Source” and select “UniProt (Hs)” to restrict to human gene annotations
2. Need a file containing a list of all genes in the experiment, one HUGO symbol per line. Use the “Browse” button, and then click “Query Gene File” to load this information. This may take some time...
3. Need a file containing a list of genes that changed. Can be one HUGO symbol per line. Optionally, you can include a second column with 1 (overexpressed) or -1 (underexpressed). Use “Browse” and “Query Changed Gene File” to load this data.

Note: GeneLink or Source can convert from various gene ids to HUGO symbols.

# GoMiner with array gene list loaded

NCI/GA Tech - Emory Univ GoMiner

File Edit Data Source Organisms View LookupSettings

**Input Genes**

Input	Gene	Source	Status
YWHAE	143E_...	UniProt	⊗
SFN	143S_...	UniProt	⊗
PPP2R...	2A5A_...	UniProt	⊗
PPP2R...	2A5B_...	UniProt	⊗
PPP2R...	2A5D_...	UniProt	⊗
PPP2R...	2A5E_...	UniProt	⊗
PPP2R...	2A5G_...	UniProt	⊗
PPP2R...	2AAA_...	UniProt	⊗
PPP2R...	2AAB_...	UniProt	⊗
PPP2R...	2ABA_...	UniProt	⊗
PPP2R...	2ABB_...	UniProt	⊗
HLA-DMA	2DMA_...	UniProt	⊗
HLA-D...	2DMB_...	UniProt	⊗
HLA-DOA	2DOA_...	UniProt	⊗
HLA-DRA	2DRA_...	UniProt	⊗
SH3BP2	3BP2_...	UniProt	⊗
SLC3A2	4F2_H_...	UniProt	⊗
A2M	A2MG_...	UniProt	⊗
ACTN1	AAC1_...	UniProt	⊗
PRKAB1	AAKB_...	UniProt	⊗
PRKAG1	AAKG_...	UniProt	⊗
ATBF1	ABF1_...	UniProt	⊗
ABL1	ABL1_...	UniProt	⊗
ABL2	ABL2_...	UniProt	⊗
ABR	ABR_H_...	UniProt	⊗
ACY1	ACY1_...	UniProt	⊗
ADAM17	AD17_...	UniProt	⊗
ADA	ADA_H_...	UniProt	⊗
ADD3	ADDG_...	UniProt	⊗
ADH6	ADH6_...	UniProt	⊗
ADK	ADK_H_...	UniProt	⊗
AOX1	ADO_H_...	UniProt	⊗
ADSS	ADSS_...	UniProt	⊗
SLC25A5	ADT2_...	UniProt	⊗
MLLT2	AF4_H_...	UniProt	⊗
GLA	AGAL_...	UniProt	⊗
ANGPT1	AGP1_...	UniProt	⊗
ANGPT2	AGP2_...	UniProt	⊗
AHR	AHR_H_...	UniProt	⊗

**Genes Mapped On GO**

top (1299)

- biological\_process (1245)
  - behavior (8)
  - biological\_process unknown (27)
  - cellular process (847)
  - development (220)
  - physiological process (1139)
    - IGF2\_HUMAN (IGF2) - (UniProt)
    - IGFA\_HUMAN (IGF1) - (UniProt)
    - O43200 (TSHR) - (UniProt)
    - PGH1\_HUMAN (PTGS1) - (UniProt)
    - PGH2\_HUMAN (PTGS2) - (UniProt)
    - REL1\_HUMAN (RLN1) - (UniProt)
    - RLF\_HUMAN (RLF) - (UniProt)
  - cellular physiological process (568)
  - coagulation (16)
  - death (123)
    - aging (2)
      - cell death (122)
        - cell aging (1)
        - cytolysis (3)
        - programmed cell death (120)
          - apoptosis (120)
          - regulation of programmed cell de...
      - extracellular structure organization and b...
    - homeostasis (13)
    - metabolism (823)
    - organismal physiological process (254)
    - pathogenesis (3)
    - regulation of physiological process (239)
    - response to stimulus (359)
    - secretion (2)
    - regulation of biological process (403)
    - viral life cycle (8)
  - cellular\_component (1070)
  - molecular\_function (1215)
    - obsolete\_component
    - obsolete\_function
    - obsolete\_process

**Summary View** **Selected Gene View**

Category Name	P-Chng	P-Undr	P-Ovr	Tot	Chng	Undr	Ovr	Category ID
ATP-dependent hel...	1.0000	1.0000	1.0000	11	0	0	0	GO:00080...
transcription elong...	1.0000	1.0000	1.0000	1	0	0	0	GO:00080...
protein C-terminus ...	1.0000	1.0000	1.0000	2	0	0	0	GO:00080...
microtubule binding	1.0000	1.0000	1.0000	3	0	0	0	GO:00080...
regulation of heart r...	1.0000	1.0000	1.0000	1	0	0	0	GO:00080...
circulation	1.0000	1.0000	1.0000	9	0	0	0	GO:00080...
beta-catenin binding	1.0000	1.0000	1.0000	1	0	0	0	GO:00080...
chemokine activity	1.0000	1.0000	1.0000	18	0	0	0	GO:00080...
oligopeptide transp...	1.0000	1.0000	1.0000	1	0	0	0	GO:00151...
peptide transporter ...	1.0000	1.0000	1.0000	2	0	0	0	GO:00151...
L-amino acid trans...	1.0000	1.0000	1.0000	1	0	0	0	GO:00151...
acidic amino acid tr...	1.0000	1.0000	1.0000	1	0	0	0	GO:00151...
amino acid transpo...	1.0000	1.0000	1.0000	2	0	0	0	GO:00151...
hexose transporter ...	1.0000	1.0000	1.0000	2	0	0	0	GO:00151...
monosaccharide tr...	1.0000	1.0000	1.0000	2	0	0	0	GO:00151...
carbohydrate trans...	1.0000	1.0000	1.0000	3	0	0	0	GO:00151...
nitric oxide metabol...	1.0000	1.0000	1.0000	4	0	0	0	GO:00462...
sodium ion transpo...	1.0000	1.0000	1.0000	1	0	0	0	GO:00150...
hydrogen ion trans...	1.0000	1.0000	1.0000	3	0	0	0	GO:00150...
monovalent inorga...	1.0000	1.0000	1.0000	3	0	0	0	GO:00150...
ion transporter activ...	1.0000	1.0000	1.0000	4	0	0	0	GO:00150...
protein phosphatas...	1.0000	1.0000	1.0000	3	0	0	0	GO:00150...
thrombin receptor a...	1.0000	1.0000	1.0000	1	0	0	0	GO:00150...
glutathione disulfid...	1.0000	1.0000	1.0000	1	0	0	0	GO:00150...
peptide disulfide ox...	1.0000	1.0000	1.0000	1	0	0	0	GO:00150...
disulfide oxidoredu...	1.0000	1.0000	1.0000	4	0	0	0	GO:00150...
protein transport	1.0000	1.0000	1.0000	26	0	0	0	GO:00150...
Cajal body	1.0000	1.0000	1.0000	1	0	0	0	GO:00150...
coreceptor activity	1.0000	1.0000	1.0000	6	0	0	0	GO:00150...
glucuronosyltransf...	1.0000	1.0000	1.0000	2	0	0	0	GO:00150...
nuclear organizatio...	1.0000	1.0000	1.0000	25	0	0	0	GO:00069...
organelle organizat...	1.0000	1.0000	1.0000	32	0	0	0	GO:00069...
unfolded protein re...	1.0000	1.0000	1.0000	1	0	0	0	GO:00069...
alcohol catabolism	1.0000	1.0000	1.0000	4	0	0	0	GO:00461...
response to unfold...	1.0000	1.0000	1.0000	5	0	0	0	GO:00069...
ER-nuclear signali...	1.0000	1.0000	1.0000	1	0	0	0	GO:00069...
response to lipid hy...	1.0000	1.0000	1.0000	1	0	0	0	GO:00069...
phenol metabolism	1.0000	1.0000	1.0000	2	0	0	0	GO:00189...

**Operations**

File Name: C:\Source\GoMinerExample\total.gene Browse...

Query Gene File Reset All Query Changed Gene File Reset Changed

**Quick Help -- Symbol Description**

⊗ All/Unchanged Genes ⊕ Changed ↑ Over ↓ Under ⓘ Term

Build:138

GENOMICS AND BIOINFORMATICS GROUP  
LMP, NCI, NIH

Medical Informatics and Bioimaging Lab  
BME, GA Tech / Emory University

# GoMiner with changed gene list loaded

NCI/GA Tech - Emory Univ GoMiner

File Edit Data Source Organisms View LookupSettings

**Input Genes**

Input	Gene	Source	Status
YVHAE	143E_...	UniProt	⊖
SFN	143S_...	UniProt	⊖
PPP2R...	2A5A_...	UniProt	⊖
PPP2R...	2A5B_...	UniProt	⊖
PPP2R...	2A5D_...	UniProt	⊖
PPP2R...	2A5E_...	UniProt	⊖
PPP2R...	2A5G_...	UniProt	⊖
PPP2R...	2AAA_...	UniProt	⊖
PPP2R...	2AAB_...	UniProt	⊖
PPP2R...	2ABA_...	UniProt	⊖
PPP2R...	2ABB_...	UniProt	⊖
HLA-DMA	2DMA_...	UniProt	⊖
HLA-D...	2DMB_...	UniProt	⊖
HLA-DOA	2DOA_...	UniProt	⊖
HLA-DRA	2DRA_...	UniProt	⬇
SH3BP2	3BP2_...	UniProt	⊖
SLC3A2	4F2_H_...	UniProt	⊖
A2M	A2MG_...	UniProt	⊖
ACTN1	AAC1_...	UniProt	⊖
PRKAB1	AAKB_...	UniProt	⊖
PRKAG1	AAKG_...	UniProt	⊖
ATBF1	ABF1_...	UniProt	⊖
ABL1	ABL1_...	UniProt	⊖
ABL2	ABL2_...	UniProt	⊖
ABR	ABR_H_...	UniProt	⊖
ACY1	ACY1_...	UniProt	⊖
ADAM17	AD17_...	UniProt	⊖
ADA	ADA_H_...	UniProt	⊖
ADD3	ADDG_...	UniProt	⊖
ADH6	ADH6_...	UniProt	⊖
ADK	ADK_H_...	UniProt	⊖
AOX1	ADO_H_...	UniProt	⊖
ADSS	ADSS_...	UniProt	⊖
SLC25A5	ADT2_...	UniProt	⊖
MLLT2	AF4_H_...	UniProt	⊖
GLA	AGAL_...	UniProt	⊖
ANGPT1	AGP1_...	UniProt	⊖
ANGPT2	AGP2_...	UniProt	⊖
AHR	AHR_H_...	UniProt	⬆

**Genes Mapped On GO**

p (1299 1.00 p=1.00 1.00 p=1.00 1.00 p=1.00)

- biological\_process (1245 1.03 p=0.17 1.01 p=0.48 1.02 p=0.17)
  - biological\_process unknown (27 1.30 p=0.46 0.53 p=0.86 0.88 p=0.1)
  - cellular\_process (847 0.99 p=0.58 0.97 p=0.69 0.98 p=0.67)
  - development (220 0.96 p=0.62 1.12 p=0.35 1.04 p=0.43)
  - physiological\_process (1139 1.08 p=0.04 1.05 p=0.11 1.06 p=0.01)
    - cellular\_physiological\_process (568 1.02 p=0.48 0.99 p=0.57 1.0)
    - coagulation (16 1.10 p=0.61 0.90 p=0.69 0.99 p=0.62)
    - death (123 1.28 p=0.26 1.17 p=0.34 1.22 p=0.20)
      - cell\_death (122 1.29 p=0.25 1.18 p=0.33 1.23 p=0.19)
        - cytolysis (3 0.00 p=1.00 4.81 p=0.19 2.64 p=0.33)
        - programmed\_cell\_death (120 1.32 p=0.24 1.08 p=0.45 1.1)
        - apoptosis (120 1.32 p=0.24 1.08 p=0.45 1.19 p=0.24)
          - regulation\_of\_programmed\_cell\_death (77 1.14 p=0.45 0.9)
    - homeostasis (13 0.00 p=1.00 5.55 p=0.00 3.05 p=0.02)
    - metabolism (823 0.90 p=0.91 1.14 p=0.04 1.03 p=0.33)
    - organismal\_physiological\_process (254 1.87 p=0.00 0.91 p=0.71)
    - regulation\_of\_physiological\_process (239 1.10 p=0.38 1.39 p=0.05 1)
    - response\_to\_stimulus (359 1.47 p=0.01 1.09 p=0.34 1.26 p=0.02)
  - regulation\_of\_biological\_process (403 1.18 p=0.18 1.25 p=0.06 1.22 p=0.06)
  - viral\_life\_cycle (8 2.19 p=0.38 1.80 p=0.44 1.98 p=0.27)
  - cellular\_component (1070 0.97 p=0.78 0.97 p=0.78 0.97 p=0.84)
  - molecular\_function (1215 0.95 p=0.96 0.91 p=1.00 0.93 p=1.00)
  - obsolete\_component
  - obsolete\_function
  - obsolete\_process

**Summary View** **Selected Gene View**

Category Name	P-Chng	P-Undr	P-Ovr	Tot	Chng	Undr	Ovr	Category ID
cytoplasmic sequ...	0.0002	0.0178	0.0260	4	4	2	2	GO:00429...
negative regulation ...	0.0002	0.0178	0.0260	4	4	2	2	GO:00429...
transcription factor...	0.0002	0.0178	0.0260	4	4	2	2	GO:00429...
regulation of transc...	0.0002	0.0178	0.0260	4	4	2	2	GO:00429...
regulation of protei...	0.0002	0.0178	0.0260	4	4	2	2	GO:00423...
regulation of nucleo...	0.0002	0.0178	0.0260	4	4	2	2	GO:00468...
chemokine activity	0.0008	0.0782	0.0060	18	8	3	5	GO:00080...
G-protein-coupled r...	0.0008	0.0782	0.0060	18	8	3	5	GO:00016...
chemokine recepto...	0.0008	0.0782	0.0060	18	8	3	5	GO:00423...
chemotaxis	0.0012	0.0547	0.0112	37	12	5	7	GO:00069...
taxis	0.0012	0.0547	0.0112	37	12	5	7	GO:00423...
response to wound...	0.0015	0.0227	0.0296	75	19	9	10	GO:00096...
response to chemi...	0.0018	0.0814	0.0097	54	15	6	9	GO:00422...
response to pathog...	0.0030	0.2972	0.0055	6	4	1	3	GO:00096...
regulation of transp...	0.0030	0.0414	0.0593	6	4	2	2	GO:00510...
immune response	0.0033	0.0002	0.4695	207	39	24	15	GO:00069...
response to pest, p...	0.0036	0.0178	0.0743	123	26	13	13	GO:00096...
extracellular space	0.0038	0.0039	0.2217	47	13	8	5	GO:00056...
protein threonine/tyr...	0.0063	0.0558	0.0794	7	4	2	2	GO:00047...
MAP kinase kinase ...	0.0063	0.0558	0.0794	7	4	2	2	GO:00047...
response to pathog...	0.0063	0.3374	0.0092	7	4	1	3	GO:00428...
antigen processing	0.0070	0.0001	1.0000	15	6	6	0	GO:00303...
antigen presentation	0.0070	0.0001	1.0000	15	6	6	0	GO:00198...
MHC class II recept...	0.0074	0.0024	0.5475	11	5	4	1	GO:00450...
response to extern...	0.0075	0.0400	0.0743	123	25	12	13	GO:00096...
defense response	0.0088	0.0008	0.4993	225	40	24	16	GO:00069...
response to biotic s...	0.0089	0.0013	0.4397	246	43	25	18	GO:00096...
inflammatory respo...	0.0096	0.1695	0.0232	52	13	5	8	GO:00069...
innate immune res...	0.0096	0.1695	0.0232	52	13	5	8	GO:00450...
physiological proce...	0.0099	0.0372	0.1127	1139	153	70	83	GO:00075...
metal ion homeost...	0.0114	1.0000	0.0008	12	5	0	5	GO:00068...
cell ion homeostasis	0.0114	1.0000	0.0008	12	5	0	5	GO:00068...
di-, tri-valent inorga...	0.0114	1.0000	0.0008	12	5	0	5	GO:00300...
cation homeostasis	0.0114	1.0000	0.0008	12	5	0	5	GO:00300...
ion homeostasis	0.0114	1.0000	0.0008	12	5	0	5	GO:00508...
response to abiotic ...	0.0119	0.1597	0.0309	65	15	6	9	GO:00096...
transforming growt...	0.0159	1.0000	0.0048	2	2	0	2	GO:00306...
NF-kappaB-nucleu...	0.0159	0.1107	0.1338	2	2	1	1	GO:00423...

**Operations**

File Name: C:\Source\GoMinerExample\undr.over.2col Browse...

Query Gene File Reset All Query Changed Gene File Reset Changed

**Quick Help -- Symbol Description**

⊖ All/Unchanged Genes ⊕ Changed ⬆ Over ⬇ Under ⊕ Term

Build:138

GENOMICS AND BIOINFORMATICS GROUP  
LMP, NCI, NIH

Medical Informatics and Bioimaging Lab  
BME, GA Tech / Emory University

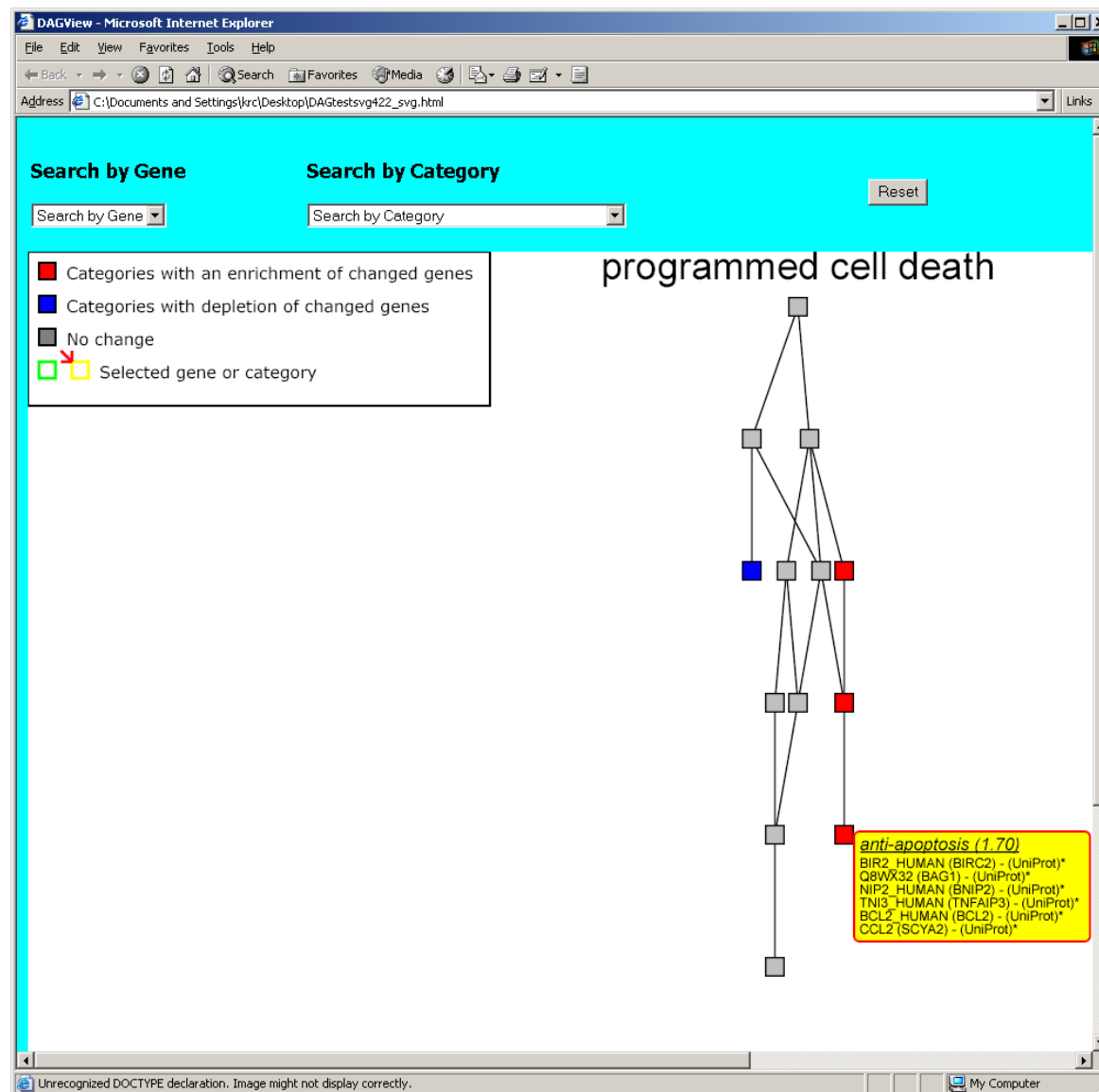
# GoMiner subgraphs

The screenshot displays the NCI/GA Tech - Emory Univ GoMiner software interface. The main window is divided into several panels:

- Input Genes:** A table listing genes and their sources. The table has columns: Input, Gene, Source, and Status. Genes listed include YVHAE, SFN, PPP2R, HLA-DMA, HLA-D, HLA-DOA, HLA-DRA, SH3BP2, SLC3A2, A2M, ACTN1, PRKAB1, PRKAG1, ATBF1, ABL1, ABL2, ABR, ACY1, ADAM17, ADA, ADD3, ADH6, ADK, AOX1, ADSS, SLC25A5, MLLT2, GLA, ANGPT1, ANGPT2, and AHR.
- Genes Mapped On GO:** A list of GO terms and their associated p-values. The list includes terms like biological\_process, cellular\_process, development, physiological\_process, cellular\_physiological\_process, coagulation, death, cell\_death, cytolysis, programmed\_cell\_death, apoptosis, and others, each with a p-value.
- Summary View:** A panel showing the selected gene view and a summary of the results. It includes a list of genes and their associated p-values.
- Operations:** A panel at the bottom left with buttons for "Query Gene File", "Reset All", "Query Changed Gene File", and "Reset Changed". It also includes a "File Name" field and a "Browse..." button.
- Quick Help -- Symbol Description:** A panel at the bottom right with a legend for symbols: All/Unchanged Genes (blue circle), Changed (red circle), Over (green circle), Under (blue circle), and Term (blue circle). It also includes a "Build:138" label and logos for the GENOMICS AND BIOINFORMATICS GROUP and the Medical Informatics and Bioimaging Lab.

A context menu is visible over the "Genes Mapped On GO" panel, showing options: "Export summary data to text file", "DAG of changed genes", "Export DAG of changed genes to file", and "Export Genes By Category".

# GoMiner subgraphs





## Intepreting GoMiner results

Enrichment is computed as

$$\frac{\text{changed genes in category} / \text{total genes in category}}{\text{changed genes on array} / \text{all genes on array}}$$

Statistical evidence of enrichment is based on a Fisher exact test.

## Intepreting GoMiner results

The p-values from the Fisher test are not corrected for multiple testing, but they should be since one is potentially looking at all GO categories. The categories are not independent, so it is not clear exactly how one should correct for multiple testing.

If one filters the gene list from the array before testing differential expression (for example, by removing low expressing or low variance genes), should those genes be included in the “query gene file” for the experiment?

The Fisher exact test is not completely appropriate, since genes can have multiple overlapping annotations into the GO DAG.

No existing test exploits the quality of evidence for the GO annotations.