

# Gene Annotation

Contributions from Carlo Colantuoni and Robert Gentleman

---

---

---

---

---

---

---

## What We Are Going To Cover

Cells, Genes, Transcripts → Genomics Experiments

Sequence Knowledge Behind Genomics Experiments

Annotation of Genes in Genomics Experiments

---

---

---

---

---

---

---

## Biological Setup

Every cell in the human body contains the entire human genome: 3.3 Gb or ~30K genes.

The investigation of gene expression is meaningful because different cells, in different environments, doing different jobs express different genes.

Tasks necessary for gene expression analysis:

Define what a gene is.

Identify genes in a sea of genomic DNA where <3% of DNA is contained in genes.

Design and implement probes that will effectively assay expression of ALL (most? many?) genes simultaneously. Cross-reference these probes.

---

---

---

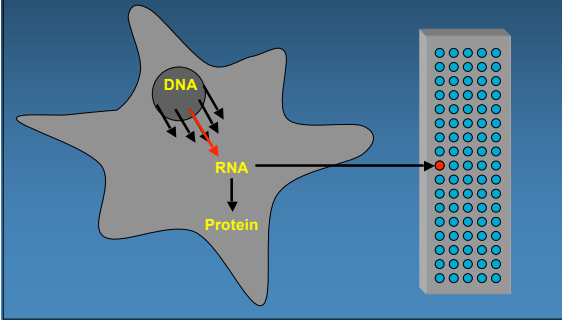
---

---

---

---

## Cellular Biology, Gene Expression, and Microarray Analysis




---

---

---

---

---

---

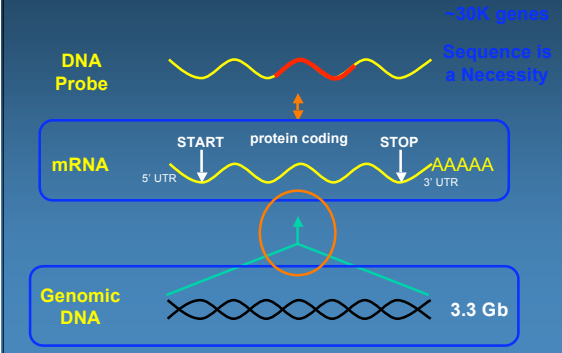
---

---

---

---

**Gene:** Protein coding unit of genomic DNA with an mRNA intermediate.




---

---

---

---

---

---

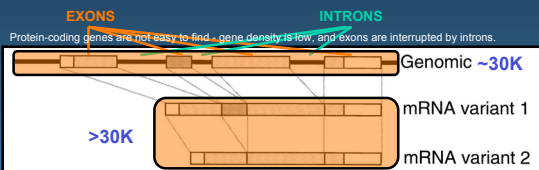
---

---

---

---

## From Genomic DNA to mRNA Transcripts



Alternative splicing  
Alternative start & stop sites in same RNA molecule  
RNA editing & SNPs

Transcript coverage  
Homology to other transcripts  
Hybridization dynamics  
3' bias

---

---

---

---

---

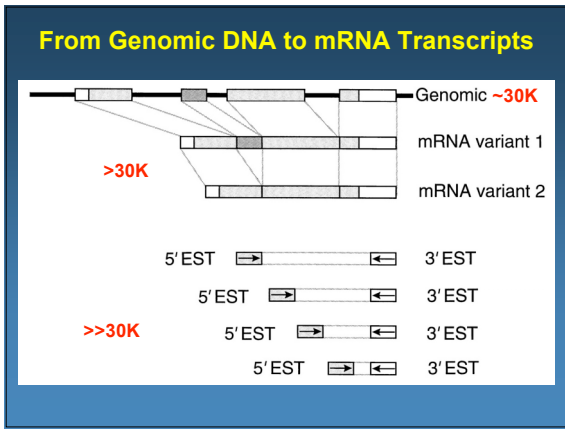
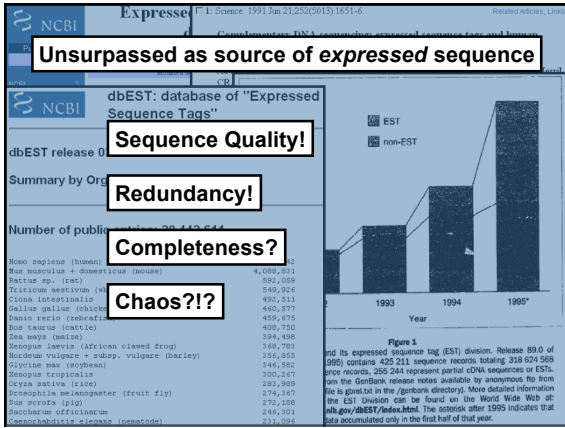
---

---

---

---

---



NCBI UniGene		CBI																																																																					
Query	Protein	Gene	Protein																																																																				
<p><b>COMT: Catechol-O-methyltransferase</b></p> <p>SELECTED MODEL ORGANISM PROTEIN SIMILARITIES</p> <table border="1"> <tr> <td>H. sapiens</td> <td>pir_A33459 - A38459 catechol O-methyltransferase (see EPI-EST)</td> <td>100 % / 264 aa</td> <td></td> </tr> <tr> <td>M. musculus</td> <td>gi_088527 - COMT_MOUSE Catechol O-methyltransferase, membrane-bound form (see EPI-EST)</td> <td>89 % / 260 aa</td> <td></td> </tr> <tr> <td>R. norvegicus</td> <td>pir_S22020 - S22090 catechol O-methyltransferase (see EPI-EST)</td> <td>79 % / 259 aa</td> <td></td> </tr> <tr> <td>A. thaliana</td> <td>pir_147051 - 147061 hypothetical protein F21F14.160... Arabidopsis thaliana (see EPI-EST)</td> <td>29 % / 113 aa</td> <td></td> </tr> </table> <p>MAPPING INFORMATION</p> <p>GeneMap: 22q11.21-q11.23 (M)</p> <p>UniSTS entry: chr22:3827791 (Map View)</p> <p>UniSTS entry: chr22:3827791 (Map View)</p> <p>UniSTS entry: chr22:3827791 (Map View)</p> <p>UniSTS entry: chr22:3827791 (Map View)</p> <p>UniSTS entry: chr22:3827791 (Map View)</p> <p>EXPRESSION INFORMATION</p> <p>Note: Highly represented (2.1 pct) in library 8332-GR0163</p> <p>cDNA sources: lymphoma, cell line, lobular carcinoma in situ, epithelium (cell line), melanocyte, malignant melanoma, sarcoma, to lymph node, amylaric, oligo-dendrocytes</p>				H. sapiens	pir_A33459 - A38459 catechol O-methyltransferase (see EPI-EST)	100 % / 264 aa		M. musculus	gi_088527 - COMT_MOUSE Catechol O-methyltransferase, membrane-bound form (see EPI-EST)	89 % / 260 aa		R. norvegicus	pir_S22020 - S22090 catechol O-methyltransferase (see EPI-EST)	79 % / 259 aa		A. thaliana	pir_147051 - 147061 hypothetical protein F21F14.160... Arabidopsis thaliana (see EPI-EST)	29 % / 113 aa																																																					
H. sapiens	pir_A33459 - A38459 catechol O-methyltransferase (see EPI-EST)	100 % / 264 aa																																																																					
M. musculus	gi_088527 - COMT_MOUSE Catechol O-methyltransferase, membrane-bound form (see EPI-EST)	89 % / 260 aa																																																																					
R. norvegicus	pir_S22020 - S22090 catechol O-methyltransferase (see EPI-EST)	79 % / 259 aa																																																																					
A. thaliana	pir_147051 - 147061 hypothetical protein F21F14.160... Arabidopsis thaliana (see EPI-EST)	29 % / 113 aa																																																																					
<p><b>mRNA SEQUENCES (13)</b></p> <table border="1"> <tr> <td>BT007125.1</td> <td>Homo sapiens catechol-O-methyltransferase mRNA, complete cds</td> <td>P</td> <td></td> </tr> <tr> <td>AK130031.1</td> <td>Homo sapiens cDNA FLJ25521 fs, clone KDN08050, highly similar to Catechol O-methyltransferase, membrane-bound form (EC 2.1.1.6)</td> <td>P</td> <td></td> </tr> <tr> <td>AK129492.1</td> <td>Homo sapiens cDNA FLJ25981 fs, clone AD50616</td> <td>P</td> <td></td> </tr> <tr> <td>NM_007310.1</td> <td>Homo sapiens catechol-O-methyltransferase (COMT), transcript variant S-COMT, mRNA</td> <td>P</td> <td></td> </tr> <tr> <td>NM_000754.2</td> <td>Homo sapiens catechol-O-methyltransferase (COMT), transcript variant MB-COMT, mRNA</td> <td>P</td> <td></td> </tr> <tr> <td>AL350148.1</td> <td>Homo sapiens mRNA, cDNA DKFZp547A166 from clone DKFZp547A166</td> <td>PA</td> <td></td> </tr> <tr> <td>BC000419.2</td> <td>Homo sapiens catechol-O-methyltransferase mRNA (cDNA clone MGC-8663 IMAGE 2964400), complete cds</td> <td>PA</td> <td></td> </tr> <tr> <td>BC011935.2</td> <td>Homo sapiens catechol-O-methyltransferase transcript variant MB-COMT, mRNA (cDNA clone MGC-20006 IMAGE 3503220), complete cds</td> <td>PA</td> <td></td> </tr> <tr> <td>BC005867.2</td> <td>Homo sapiens catechol-O-methyltransferase mRNA (cDNA clone MGC-4072 IMAGE 2964400), complete cds</td> <td>PA</td> <td></td> </tr> <tr> <td>M55212.1</td> <td>Homo sapiens catechol-O-methyltransferase (COMT) mRNA, complete cds</td> <td>P</td> <td></td> </tr> <tr> <td>M55213.1</td> <td>Homo sapiens catechol-O-methyltransferase (COMT) mRNA, 5' end</td> <td>P</td> <td></td> </tr> <tr> <td>M55225.1</td> <td>Homo sapiens catechol-O-methyltransferase (COMT) mRNA, complete cds</td> <td>P</td> <td></td> </tr> <tr> <td colspan="4"><b>EST SEQUENCES (10 of 84) (Show all ESTs)</b></td> </tr> <tr> <td>BE381620.1</td> <td>cDNA clone endonostriatum IMAGE3895343 adenocarcinoma cell line</td> <td></td> <td>5' read PM</td> </tr> <tr> <td>BQ226776.1</td> <td>cDNA clone embryonal IMAGE6045206 carcinoma, cell line</td> <td></td> <td>5' read PM</td> </tr> <tr> <td>BQ415574.1</td> <td>cDNA clone melanotic IMAGE4620406 melanoma</td> <td></td> <td>5' read PM</td> </tr> <tr> <td>BQ415500.1</td> <td>cDNA clone melanotic</td> <td></td> <td>5' read PM</td> </tr> </table>				BT007125.1	Homo sapiens catechol-O-methyltransferase mRNA, complete cds	P		AK130031.1	Homo sapiens cDNA FLJ25521 fs, clone KDN08050, highly similar to Catechol O-methyltransferase, membrane-bound form (EC 2.1.1.6)	P		AK129492.1	Homo sapiens cDNA FLJ25981 fs, clone AD50616	P		NM_007310.1	Homo sapiens catechol-O-methyltransferase (COMT), transcript variant S-COMT, mRNA	P		NM_000754.2	Homo sapiens catechol-O-methyltransferase (COMT), transcript variant MB-COMT, mRNA	P		AL350148.1	Homo sapiens mRNA, cDNA DKFZp547A166 from clone DKFZp547A166	PA		BC000419.2	Homo sapiens catechol-O-methyltransferase mRNA (cDNA clone MGC-8663 IMAGE 2964400), complete cds	PA		BC011935.2	Homo sapiens catechol-O-methyltransferase transcript variant MB-COMT, mRNA (cDNA clone MGC-20006 IMAGE 3503220), complete cds	PA		BC005867.2	Homo sapiens catechol-O-methyltransferase mRNA (cDNA clone MGC-4072 IMAGE 2964400), complete cds	PA		M55212.1	Homo sapiens catechol-O-methyltransferase (COMT) mRNA, complete cds	P		M55213.1	Homo sapiens catechol-O-methyltransferase (COMT) mRNA, 5' end	P		M55225.1	Homo sapiens catechol-O-methyltransferase (COMT) mRNA, complete cds	P		<b>EST SEQUENCES (10 of 84) (Show all ESTs)</b>				BE381620.1	cDNA clone endonostriatum IMAGE3895343 adenocarcinoma cell line		5' read PM	BQ226776.1	cDNA clone embryonal IMAGE6045206 carcinoma, cell line		5' read PM	BQ415574.1	cDNA clone melanotic IMAGE4620406 melanoma		5' read PM	BQ415500.1	cDNA clone melanotic		5' read PM
BT007125.1	Homo sapiens catechol-O-methyltransferase mRNA, complete cds	P																																																																					
AK130031.1	Homo sapiens cDNA FLJ25521 fs, clone KDN08050, highly similar to Catechol O-methyltransferase, membrane-bound form (EC 2.1.1.6)	P																																																																					
AK129492.1	Homo sapiens cDNA FLJ25981 fs, clone AD50616	P																																																																					
NM_007310.1	Homo sapiens catechol-O-methyltransferase (COMT), transcript variant S-COMT, mRNA	P																																																																					
NM_000754.2	Homo sapiens catechol-O-methyltransferase (COMT), transcript variant MB-COMT, mRNA	P																																																																					
AL350148.1	Homo sapiens mRNA, cDNA DKFZp547A166 from clone DKFZp547A166	PA																																																																					
BC000419.2	Homo sapiens catechol-O-methyltransferase mRNA (cDNA clone MGC-8663 IMAGE 2964400), complete cds	PA																																																																					
BC011935.2	Homo sapiens catechol-O-methyltransferase transcript variant MB-COMT, mRNA (cDNA clone MGC-20006 IMAGE 3503220), complete cds	PA																																																																					
BC005867.2	Homo sapiens catechol-O-methyltransferase mRNA (cDNA clone MGC-4072 IMAGE 2964400), complete cds	PA																																																																					
M55212.1	Homo sapiens catechol-O-methyltransferase (COMT) mRNA, complete cds	P																																																																					
M55213.1	Homo sapiens catechol-O-methyltransferase (COMT) mRNA, 5' end	P																																																																					
M55225.1	Homo sapiens catechol-O-methyltransferase (COMT) mRNA, complete cds	P																																																																					
<b>EST SEQUENCES (10 of 84) (Show all ESTs)</b>																																																																							
BE381620.1	cDNA clone endonostriatum IMAGE3895343 adenocarcinoma cell line		5' read PM																																																																				
BQ226776.1	cDNA clone embryonal IMAGE6045206 carcinoma, cell line		5' read PM																																																																				
BQ415574.1	cDNA clone melanotic IMAGE4620406 melanoma		5' read PM																																																																				
BQ415500.1	cDNA clone melanotic		5' read PM																																																																				





**http://www.ncbi.nlm.nih.gov/Entrez/**

---

---

---

---

---

---

---

---

---

---

### Functional Annotation of Lists of Genes

**KEGG**  
**PFAM**  
**SWISS-PROT**  
**GO**

**DRAGON**  
**DAVID**  
**BioConductor**

---

---

---

---

---

---

---

---

---

---

### Analysis of Functional Gene Groups

---

---

---

---

---

---

---

---

---

---







## WWW resources

- Nucleotide databases: e.g. GenBank.
- Gene databases: e.g. Entrez Gene, UniGene.
- Protein sequence and structure databases: e.g. SwissProt, Protein DataBank (PDB).
- Literature databases: e.g. PubMed, OMIM.
- Chromosome maps: e.g. NCBI Map Viewer.
- Pathways: e.g. KEGG.
- [Entrez](#) is a search and retrieval system that integrates information from databases at NCBI (National Center for Biotechnology Information).
- if you know of some we should be using – please let us know

---

---

---

---

---

---

---

---

## annotate: matching IDs

### Important tasks

- Associate manufacturers or in-house probe identifiers to other available identifiers.  
E.g.  
Affymetrix IDs → Entrez Gene IDs  
Affymetrix IDs → GenBank accession number.
- Associate probes with biological data such as chromosomal position, pathways.
- Associate probes with published literature data via PubMed (need PMID).

---

---

---

---

---

---

---

---

## annotate: matching IDs

Affymetrix identifier	"41046_s_at"
HGU95A chips	
Entrez Gene ID	"9203"
GenBank accession #	"X95808"
Gene symbol	"ZNF261"
PubMed, PMID	"10486218" "9205841" "8817323"
Chromosomal location	"X", "Xq13.1"

---

---

---

---

---

---

---

---

## Annotation data packages

- The Bioconductor project provides annotation data packages, that contain many different mappings to interesting data
  - Mappings between Affy IDs and other probe IDs: hgu95av2 for HGU95Av2 GeneChip series, also, hgu133a, hu6800, mgu74a, rgu34a, YG.
  - Affy CDF data packages.
  - Probe sequence data packages.
- These packages are updated and expanded regularly as new data become available.
- They can be downloaded from the Bioconductor website and also using `installDataPackage`.
- `DPEXplorer`: a widget for interacting with data packages.
- `AnnBuilder`: tools for building annotation data packages.

---

---

---

---

---

---

---

---

## annotate: matching IDs

- Much of what `annotate` does relies on matching symbols.
- This is basically the role of a hash table in most programming languages.
- In R, we rely on environments.
- The annotation data packages provide R environment objects containing key and value pairs for the mappings between two sets of probe identifiers.
- Keys can be accessed using the `Rs` function.
- Matching values in different environments can be accessed using the `get` or `multiget` functions.

---

---

---

---

---

---

---

---

## annotate: matching IDs

```
> library(hgu95av2)
> get("41046_s_at", env = hgu95av2ACCNUM)
[1] "X95808"
> get("41046_s_at", env = hgu95av2LOCUSID)
[1] "9203"
> get("41046_s_at", env = hgu95av2SYMBOL)
[1] "ZNF261"
> get("41046_s_at", env = hgu95av2GENENAME)
[1] "zinc finger protein 261"
> get("41046_s_at", env = hgu95av2SUMFUNC)
[1] "Contains a putative zinc-binding
motif (MYM)|Proteome"
> get("41046_s_at", env = hgu95av2UNIGENE)
[1] "Hs.9568"
```

---

---

---

---

---

---

---

---

## annotate: matching IDs

```
> get("41046_s_at", env = hgu95av2CHR)
[1] "X"
> get("41046_s_at", env = hgu95av2CHRLOC)
  X
-68692698
> get("41046_s_at", env = hgu95av2MAP)
[1] "Xq13.1"
> get("41046_s_at", env = hgu95av2PMID)
[1] "10486218" "9205841" "8817323"
> get("41046_s_at", env = hgu95av2GO)
  TAS      TAS      IEA
"GO:0003677" "GO:0007275" "GO:0016021"
```

---

---

---

---

---

---

---

---

## annotate: matching IDs

- Instead of relying on the general R functions for environments, new user-friendly functions have been written for accessing and working with specific identifiers.
- E.g. `getGO`, `getGODesc`, `getLL`, `getPMID`, `getSYMBOL`.

---

---

---

---

---

---

---

---

## annotate: matching IDs

```
> getSYMBOL("41046_s_at", data="hgu95av2")
41046_s_at
"ZNF261"
> gg<- getGO("41046_s_at", data="hgu95av2")
> getGODesc(gg[[1]], "MF")
$"GO:0003677"

"DNA binding activity"
> getLL("41046_s_at", data="hgu95av2")
41046_s_at
9203
> getPMID("41046_s_at", data="hgu95av2")
$"41046_s_at"
[1] 10486218 9205841 8817323
```

---

---

---

---

---

---

---

---

## annotate: querying databases

The `annotate` package provides tools for

- Searching and processing information from various WWW biological databases
  - GenBank,
  - LocusLink,
  - PubMed.
- Regular expression searching of PubMed abstracts.
- Generating nice HTML reports of analyses, with links to biological databases.

---

---

---

---

---

---

---

---

## annotate: WWW queries

- Functions for querying WWW databases from R rely on the `browseURL` function

```
browseURL("www.r-project.org")
```

**Other tools:** `HTMLPage` class, `getTDRows`, `getQueryLink`, `getQuery4UG`, `getQuery4LL`, `makeAnchor`.

- The `XML` package is used to parse query results.

---

---

---

---

---

---

---

---

## annotate: querying GenBank

[www.ncbi.nlm.nih.gov/Genbank/index.html](http://www.ncbi.nlm.nih.gov/Genbank/index.html)

- Given a vector of GenBank accession numbers or NCBI UIDs, the `genbank` function
  - opens a browser at the URLs for the corresponding GenBank queries;
  - returns an `XMLdoc` object with the same data.

```
genbank("X95808", disp="browser")
```

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?tool=bincom&cmd=Search&db=Nucleotide&term=X95808>

```
genbank(1430782, disp="data",  
       type="uid")
```

---

---

---

---

---

---

---

---

### annotate: querying LocusLink

[www.ncbi.nlm.nih.gov/LocusLink/](http://www.ncbi.nlm.nih.gov/LocusLink/)

- `locuslinkByID`: given one or more LocusIDs, the browser is opened at the URL corresponding to the first gene.

```
locuslinkByID("9203")  
http://www.ncbi.nlm.nih.gov/LocusLink/LocRpt.cgi?l=9203
```

- `locuslinkQuery`: given a search string, the results of the LocusLink query are displayed in the browser.

```
locuslinkQuery("zinc finger")  
http://www.ncbi.nlm.nih.gov/LocusLink/list.cgi?Q=zinc finger&ORG=Hs&V=0
```

- `getQuery4LL`.

---

---

---

---

---

---

---

---

### annotate: querying PubMed

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

- For any gene there is often a large amount of data available from PubMed.
- The `annotate` package provides the following tools for interacting with PubMed

- `pubMedAbst`: a class structure for PubMed abstracts in R.

- `pubmed`: the basic engine for talking to PubMed (`pmidQuery`).

---

---

---

---

---

---

---

---

### annotate: pubMedAbst class

Class structure for storing and processing PubMed abstracts in R

- `pmid`
- `authors`
- `abstText`
- `articleTitle`
- `journal`
- `pubDate`
- `abstUrl`

---

---

---

---

---

---

---

---

## annotate: high-level tools for querying PubMed

- `pm.getabst`: download the specified PubMed abstracts (stored in XML) and create a list of `pubMedAbst` objects.
- `pm.titles`: extract the titles from a list of PubMed abstracts.
- `pm.abstGrep`: regular expression matching on the abstracts.

---

---

---

---

---

---

---

---

## annotate: PubMed example

```
pmid <-get("41046_s_at", env=hgu95aPMID)
pubmed(pmid, disp="browser")

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?tool=bioconductor&cmd=Retrieve&db=PubMed&list\_uids=10486218%2c9205841%2c8817323

absts <- pm.getabst("41046_s_at", base="hgu95a")
pm.titles(absts)
pm.abstGrep("retardation",absts[[1]])
```

---

---

---

---

---

---

---

---

## annotate: PubMed HTML report

- The new function `pmAbst2HTML` takes a list of `pubMedAbst` objects and generates an HTML report with the titles of the abstracts and links to their full page on PubMed.

```
pmAbst2HTML(absts[[1]],
            filename="pm.html")
```

---

---

---

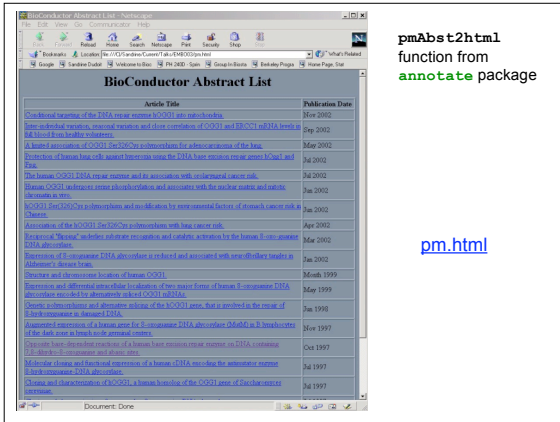
---

---

---

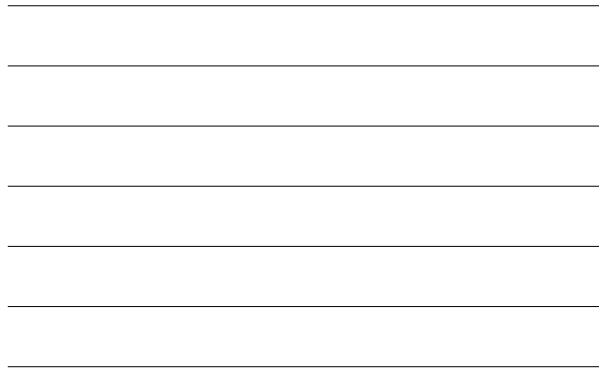
---

---



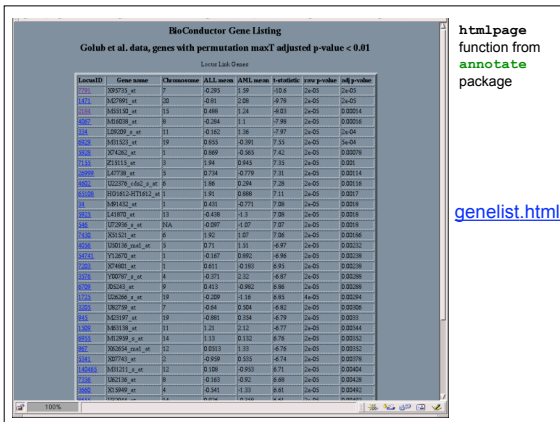
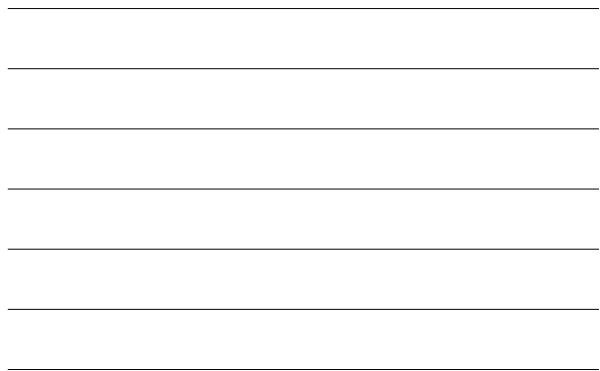
pmAbst2html  
function from  
annotate package

[pm.html](#)



## annotate: analysis reports

- A simple interface, [htmlpage](#), can be used to generate an HTML report of analysis results.
- The page consists of a table with one row per gene, with links to Entrez Gene, Affymetrix, SwissProt, UniGene, or OMIM.
- Entries can include various gene identifiers and statistics.



htmlpage  
function from  
annotate  
package

[genelist.html](#)



## annaffy

- Provides simplified mappings between Affymetrix IDs and annotation data
- Relies on chip-level annotation packages created by AnnBuilder
- Supplies functions to produce mappings for almost all environments in a given annotation package

---

---

---

---

---

---

---

---

## annaffy:Interactive

```
> symbol <- aafSymbol(probinds, "hgu95av2")
> getText(symbol)
[1] "COL11A2" "FLT3" "BDNF" "CD19" "GSTT2" "FGFR2" "IL18"
[8] "IFNB1" "RAB5B" "TAF11"

> gos <- aafGO(probinds, "hgu95av2")
> gos[[3]]
An object of class "aafGO"
[[1]]
An object of class "aafGOItem"
 @id  "GO:0007399"
 @name "neurogenesis"
 @type "Biological Process"
 @evid "TAS"
```

---

---

---

---

---

---

---

---

## annaffy:Interactive

```
> gbs <- aafGenBank(probinds, "hgu95av2")
> getURL(gbs[[3]])
[1]
"http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=search&db=nucleotide&term=M61176%5BACCN%5D&doptcmd=GenBank"
> browseURL(getURL(gbs[[3]]))

This will open a browser pointing to this particular GenBank ID
```

---

---

---

---

---

---

---

---



## annaffy:Non-interactive

- Primary function of annaffy is to produce very nice HTML or text tables
- These tables can contain:
  - Links to databases
  - Statistics
  - Expression measures
    - Color-coded to intensity for easy viewing

---

---

---

---

---

---

---

---

---

---

---

---

## annaffy:HTML Table

```
> aaf.handler()
[1] "Probe"           "Symbol"           "Description"
[4] "Function"        "Chromosome"       "Chromosome Location"
[7] "GenBank"         "LocusLink"        "Cytoband"
[10] "UniGene"         "PubMed"           "Gene Ontology"
[13] "Pathway"

> anntable <- aafTableAnn(probinds[1:10], "hgu95av2", aaf.handler())[c(1:3,
10)]
> sttable <- aafTable("t-stat" = rnorm(10), "p-value" = runif(10))
> exprtable <- aafTableInt(aafExpr, probeids = probinds[1:10])
> table <- merge(anntable, sttable)
> table <- merge(table, exprtable)
> saveHTML(table, "faketable.HTML", title="Some Fake Results")
```

---

---

---

---

---

---

---

---

---

---

---

---

Some Fake Results

Probe	Symbol	Description	Exp1	Exp2	Exp3	Exp4	Exp5	Exp6	Exp7	Exp8	Exp9	Exp10
1001_#	UCS11A2	UCS11A2, transmembrane alpha 2	0.26112	1.42334	0.479316	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1002_#	UC23	UC23, uncharacterized protein human 3	0.26112	-2.20305	0.43332	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1003_#	UC24P	UC24P, uncharacterized protein human 4	0.26112	-2.51622	0.46534	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1004_#	UC19	UC19, uncharacterized protein human	0.26112	-1.47451	0.409718	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1005_#	UC271	UC271, uncharacterized protein human	0.26112	0.74329	0.409797	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1006_#	UC272	UC272, uncharacterized protein human	0.26112	0.74329	0.409797	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1007_#	UC273	UC273, uncharacterized protein human	0.26112	0.74329	0.409797	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1008_#	UC274	UC274, uncharacterized protein human	0.26112	0.74329	0.409797	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1009_#	UC275	UC275, uncharacterized protein human	0.26112	0.74329	0.409797	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988
1010_#	UC276	UC276, uncharacterized protein human	0.26112	0.74329	0.409797	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988	0.9988

---

---

---

---

---

---

---

---

---

---

---

---

## Supplemental Slides

---

---

---

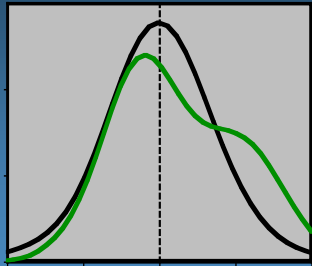
---

---

---

---

### Analysis of Functional Gene Groups



---

---

---

---

---

---

---

### Functional Gene/Protein Networks

**DIP**  
**BIND**  
**MINT**  
**HPRD**  
**PubGene**  
**Predicted Protein Interactions**

---

---

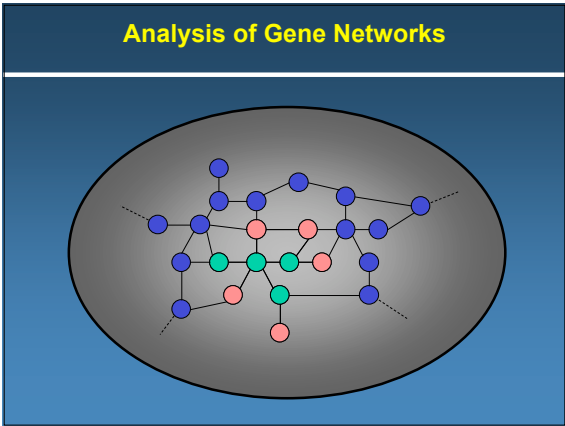
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

### Database of Interacting Proteins

DATABASE STATISTICS

Number of proteins	17043
Number of organisms	107
Number of interactions	44349
Number of distinct experiments describing an interaction	49304
Number of data sources (articles)	2694
Number of data sources (other)	34

ORGANISM	PROTEINS	INTERACTIONS	EXPERIMENTS	Details
<i>Drosophila melanogaster</i> (fruit fly)	7052	26996	21012	↗
<i>Saccharomyces cerevisiae</i> (baker's yeast)	4749	15656	19143	↗
<i>Caenorhabditis elegans</i>	2638	4030	4075	↗
<i>Helicobacter pylori</i>	710	1425	1425	[...]
<i>Homo sapiens</i> (Human)	897	1379	1998	↗
<i>Escherichia coli</i>	421	516	971	↗
<i>Mus musculus</i> (house mouse)	197	288	389	↗
<i>Rattus norvegicus</i> (Norway rat)	84	107	154	↗
Others (99)	300			

---

---

---

---

---

---

---

---

---

---

#### Blueprint

The Biomolecular Interaction Network Database (BIND) is a collection of records documenting molecular interactions. The contents of BIND include high-throughput data submissions and hand-curated information gathered from the scientific literature.

BIND is an interaction database with three classifications: interactions, complexes, and pathways. Interactions are formed from one or more interacting molecules, complexes are formed from two or more interacting molecules, and pathways are formed by a specific sequence of two or more interactions.

A BIND record represents an interaction between two or more objects that is believed to occur in a living organism. A biological object can be a protein, DNA, RNA, lipid, molecular complex, gene, photon or an unclassified biological entity. BIND records are created for interactions which have been shown experimentally and published in at least one peer-reviewed journal. A record also references any papers with experimental evidence that support or dispute the associated interaction.

Interactions are the basic units of BIND and can be linked together to form molecular complexes or pathways.

A molecular complex is a collection of two or more molecules that associate to form a function unit in a living organism. In BIND, these are represented as a molecular complex object, formed by linking two or more interaction records. Molecular complex records are supplemented with additional information such as complex topology and associated objects involved in the interaction.

A pathway is a collection of two or more interactions within a living organism. In BIND, pathway records that are formed by linking the pathway records are supplemented with additional information such as the pathway events associated with a particular disease.

The object-oriented design of BIND has allowed diversity of interactions in a format that is efficient and used by software engineers developing new to submissions from Blueprint curators and has been growing exponentially since its inception. Recent submissions from researchers in other use BIND to understand publicly available data.

#### Blueprint

Current BIND Database Statistics

Record Type	Count
Interactions (All)	77573
Interactions (Spoke Mode)	7624
Molecular Complexes	1522
Pathways	8
Organisms Represented	953
Sequences (GFs)	32995
Publications	9283

#### Accession Query

Results for NCBI Taxonomy ID(s): 9606

Search by:

Organism ID:

9606 is the Taxonomy ID for Homo Sapiens  
 Taxonomy ID 9606 is found in:  
[18136-20-0 Interaction Link](#) | [4 BIND Pathway Links](#) | [116 BIND Molecular Complex Links](#)

---

---

---

---

---

---

---

---

---

---



## Predicted Human Protein Interactions

Research

### A first-draft human protein-interaction map Ben Lehner and Andrew G Fraser

Address: The Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA, UK.

Correspondence: Andrew G Fraser. E-mail: agf@sanger.ac.uk

Published: 13 August 2004  
Genome **Biology** 2004, **5**:R63

Received: 7 May 2004  
Revised: 23 June 2004  
Accepted: 20 July 2004

The electronic version of this article is the complete one and can be found online at <http://genomebiology.com/2004/5/9/R63>

---

---

---

---

---

---

---

---

## Predicted Human Protein Interactions

Used high-throughput protein interaction experiments from fly, worm, and yeast to predict human protein interactions.

Human protein interaction is predicted if both proteins in an interaction pair from other organism have high sequence homology to human proteins.

>70K Hs interactions predicted  
>6K Hs genes

---

---

---

---

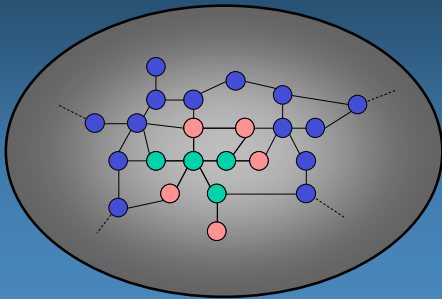
---

---

---

---

## Analysis of Gene Networks



---

---

---

---

---

---

---

---

### NCBI Web Links

<http://www.ncbi.nlm.nih.gov>  
<http://www.ncbi.nlm.nih.gov/Entrez/>  
<http://www.ncbi.nlm.nih.gov/Genbank/>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Nucleotide>  
<http://www.ncbi.nlm.nih.gov/dbEST/>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=gene>  
<http://www.ncbi.nlm.nih.gov/LocusLink/>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=unigene>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=homologene>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=OMIM>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed>  
<http://www.ncbi.nlm.nih.gov/PubMed/>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=cdd>  
<http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>  
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=snp>  
<http://www.ncbi.nlm.nih.gov/SNP/>  
<http://eutils.ncbi.nlm.nih.gov/entrez/query/static/advancedentrez.html>  
<http://www.ncbi.nlm.nih.gov/geo/>  
<http://www.ncbi.nlm.nih.gov/RefSeq/>

**FTP:**  
<ftp://ftp.ncbi.nlm.nih.gov/>  
<ftp://ftp.ncbi.nlm.nih.gov/repository/UniGene>  
<ftp://ftp.ncbi.nlm.nih.gov/pub/HomoloGene/>

---

---

---

---

---

---

---

---

---

---

---

---

<p><b>NUCLEOTIDE:</b></p> <p> <a href="http://genome.ucsc.edu/">http://genome.ucsc.edu/</a>  <a href="http://www.embl-heidelberg.de/">http://www.embl-heidelberg.de/</a>  <a href="http://www.ensembl.org/">http://www.ensembl.org/</a>  <a href="http://www.ebi.ac.uk/">http://www.ebi.ac.uk/</a>  <a href="http://www.gdb.org/">http://www.gdb.org/</a>  <a href="http://bioinfo.weizmann.ac.il/cards/index.html">http://bioinfo.weizmann.ac.il/cards/index.html</a>  <a href="http://www.gene.ac.uk/cgi-bin/medline/medline.cgi?db=pubmed">http://www.gene.ac.uk/cgi-bin/medline/medline.cgi?db=pubmed</a> </p> <p><b>PATHWAYS and NETWORKS:</b></p> <p> <a href="http://www.genome.ad.jp/kegg/">http://www.genome.ad.jp/kegg/</a>  <a href="ftp://ftp.genome.ad.jp/pub/kegg/">ftp://ftp.genome.ad.jp/pub/kegg/</a>  <a href="http://www.genome.ad.jp/anorftp/">(http://www.genome.ad.jp/anorftp/)</a>  <a href="http://dip.doe-mbi.ucla.edu">http://dip.doe-mbi.ucla.edu</a>  <a href="http://dip.doe-mbi.ucla.edu/dip/Download.cgi">http://dip.doe-mbi.ucla.edu/dip/Download.cgi</a>  <a href="http://www.blueprint.org/bind/">http://www.blueprint.org/bind/</a>  <a href="http://www.blueprint.org/bind/bind_downloads.html">http://www.blueprint.org/bind/bind_downloads.html</a>  <a href="http://160.80.34.4/mint/index.php">http://160.80.34.4/mint/index.php</a>  <a href="http://160.80.34.4/mint/release/main.php">http://160.80.34.4/mint/release/main.php</a>  <a href="http://www.hprd.org/">http://www.hprd.org/</a>  <a href="http://www.hprd.org/FAQ?selectedtab=DOWNLOAD+REQUESTS">http://www.hprd.org/FAQ?selectedtab=DOWNLOAD+REQUESTS</a>  <a href="http://www.pubgene.org/ (also .com)">http://www.pubgene.org/ (also .com)</a> </p>	<p><b>PROTEIN:</b></p> <p> <a href="http://us.expasy.org/">http://us.expasy.org/</a>  <a href="ftp://us.expasy.org/">ftp://us.expasy.org/</a>  <a href="http://www.sanger.ac.uk/Software/Pfam/">http://www.sanger.ac.uk/Software/Pfam/</a>  <a href="http://www.sanger.ac.uk/Software/Pfam/ftp.shtml">http://www.sanger.ac.uk/Software/Pfam/ftp.shtml</a>  <a href="http://smart.embl-heidelberg.de/">http://smart.embl-heidelberg.de/</a>  <a href="http://www.ebi.ac.uk/interpro/">http://www.ebi.ac.uk/interpro/</a>  <a href="http://us.expasy.org/prosite/">http://us.expasy.org/prosite/</a>  <a href="ftp://us.expasy.org/databases/prosite/">ftp://us.expasy.org/databases/prosite/</a> </p>
---	--

**More Web Links**

<http://www.bioconductor.org/>  
<http://apps1.niaid.nih.gov/david/>  
<http://www.geneontology.org/>  
<http://reelcover.nci.nih.gov/gomimer/index.jsp>  
<http://pubmatrix.grc.nia.nih.gov/>  
<http://pevsnerlab.kennedykrieger.org/dragon.htm>

---

---

---

---

---

---

---

---

---

---

---

---