# EBI web resources I: databases and tools
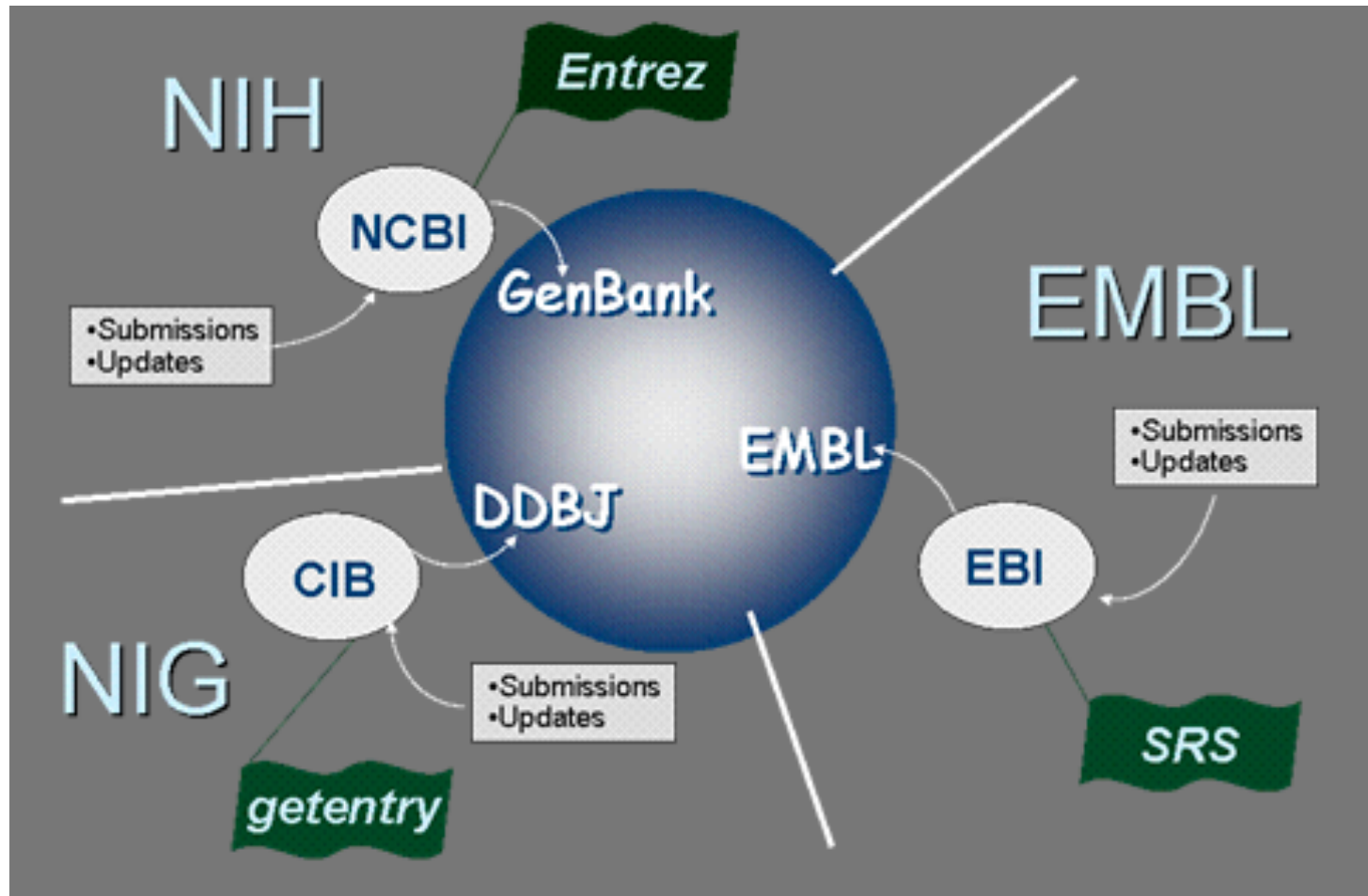
Yanbin Yin

Fall 2015

# Outline

- Intro to EBI

- Databases and web tools
  - UniProt
  - Gene Ontology

- Hands on Practice

MOST MATERIALS ARE FROM: http://www.ebi.ac.uk/training/online/course-list

# Three international nucleotide sequence databases

# The European Bioinformatics Institute (EBI)



Created in 1992 as part of European Molecular Biology Laboratory (EMBL)

EMBL was created in 1974 and is a molecular biology research institution supported by 20 European countries and Australia



Wellcome Trust Genome Campus, Hinxton,Cambridge, UK
Neighbor of Wellcome Trust Sanger Institute

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

EMBL-EBI

Services | Research | Training | About us

# The European Bioinformatics Institute

Part of the European Molecular Biology Laboratory

EMBL-EBI provides freely available data from life science experiments, performs basic research in computational biology and offers an extensive user training programme, supporting researchers in academia and industry.

## Find a gene, protein or chemical:

[ Search ]

Examples: blast, keratin, bfl1...

## News from EMBL-EBI

**The new, improved human genome**

Ensembl has incorporated a vast amount of knowledge into a fully annotated reference human genome, GRCh38, providing a solid foundation for future genomics research.

**Global Alliance** for Genomics & Health

**New Genomics API from the Global Alliance for Genomics and Health**

New software allows researchers to share anonymised genetic data seamlessly across platforms.

**Marmoset genome sheds light on chimeral twins**

Initial analyses of the marmoset genome provide insight into this tiny primate's reproductive system, which is well adapted to multiple births. The marmoset sequence is freely available in the

European Molecular Biology Laboratory

Visit **EMBL**.org

EMBL
40 YEARS : 1974-2014

## Popular

| | |
|---|---|
| Services | Jobs |
| Research | Visit us |
| Training | EMBL |
| News | Contacts |

## Events

📅 1 day course in metabolomics and bioinformatics for Nutritionists (London, UK)
Sep 23 2014
Registration deadline: Sep 16 2014

📅 diXa Open Meeting - 29-30 September 2014
Sep 29 2014 -Sep 30 2014
Registration deadline: Sep 12 2014

See all courses and conferences
See other events at EMBL-EBI

5

# Research groups in EBI

| | Group/team leader | Area of research | |
|---|---|---|---|
| **Genomes** | Ewan Birney | Algorithmic methods for genome analysis | <span style="color:red">InterPro</span> |
| | Paul Flicek | Vertebrate genomics | |
| | Nick Goldman | Evolutionary tools for sequence analysis | |
| **Transcriptomes** | Alvis Brazma | Functional genomics | |
| | Anton Enright | Functional genomics and analysis of small RNA function | <span style="color:red">miRBase</span> |
| | John Marioni | Computational and evolutionary genomics | |
| | Oliver Stegle | Statistical genomics and systems genetics | |
| **Proteins** | Janet Thornton | Computational biology of proteins: structure, function and evolution | |
| | Rolf Apweiler | Protein sequence analysis and functional annotation | <span style="color:red">UniProt</span> |
| | Gerard Kleywegt | Structural validation of proteins; protein-ligand interactions | |
| **Pathways and systems** | Nicolas Le Novère | Computational systems neurobiology | |
| | Nick Luscombe | Genomics and regulatory systems | |
| | Paul Bertone | Pluripotency, reprogramming and differentiation | |
| | Julio Saez-Rodriguez | Systems biomedicine | |
| **Literature** | Dietrich Rebholz-Schuhmann | Literature analysis and semantic data integration in life science research | |
| **Chemistry** | Christoph Steinbeck | Cheminformatics and metabolism | |
| | John Overington | Chemogenomics and drug discovery | |

6

# Major databases in EBI

GenBank ——————— EMBL-Bank (DNA and RNA sequences)
Genome MapView ——————— Ensembl (genomes)
GEO ——————— ArrayExpress(microarray-based gene-expression data)
nr (GenPept) ——————— UniProt (protein sequences)
CDD ——————— InterPro(protein families, domains and motifs)
MMDB ——————— PDBe (macromolecular structures)

Others, such as
IntAct (protein–protein interactions)
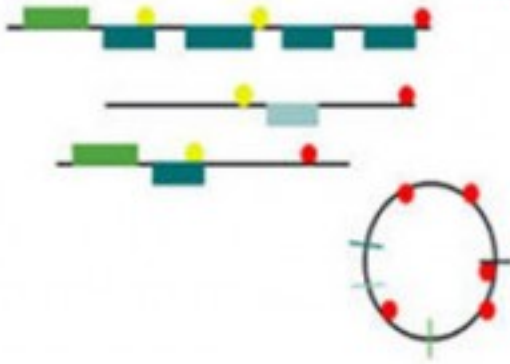Reactome (pathways)
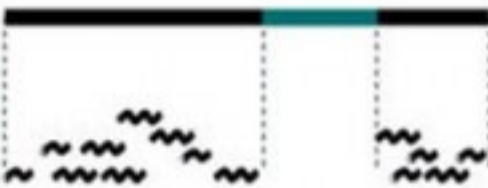ChEBI (small molecules)
IntEnz (enzyme classification)
GO (gene ontology)

Swiss Institute of Bioinformatics
Sanger Institute

chromatograms

Sequence might first enter ENA as SRA (Sequence Read Archive) fragmented sequence reads; it might be re-submitted as assembled WGS (Whole Genome Shotgun) sequence overlap contigs; it might be re-submitted again with further assembly as CON (Constructed) sequence entries, with the older WGS entries being consigned to the Sequence Version Archive

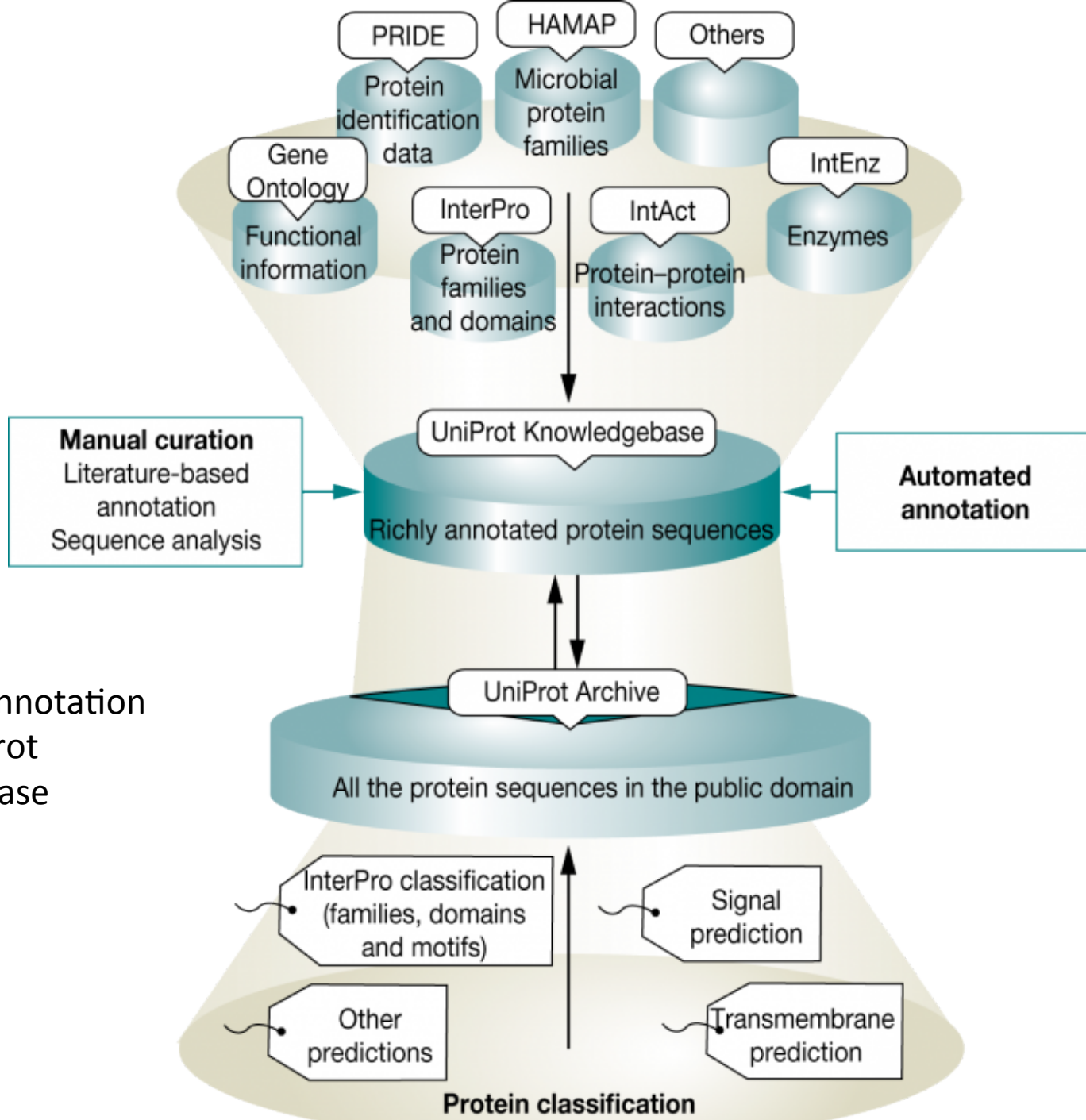Data is first split into classes, then it is split into intersecting slices by taxonomy

# UniProt



http://www.uniprot.org/help/uniparc

11

Sources of annotation for the UniProt Knowledgebase

Curation generation
http://cys.bios.niu.edu/yyin/teach/PBB/Bioinformatics%20Curation%20generation.pdf

Life as a **Scientific Curator**
http://www.ebi.ac.uk/about/jobs/career-profiles/scientific-curator

Scientific Database Curator job : Cambridge, United Kingdom
http://www.nature.com/naturejobs/science/jobs/444213-scientific-database-curator

# Hands on practice 1: UniProt

www.uniprot.org

http://www.uniprot.org/docs/uniprot_flyer.pdf
http://www.uniprot.org/help/about

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein

**UniProtKB**
UniProt Knowledgebase

Swiss-Prot
(549,215)
Manually annotated and reviewed.

TrEMBL
(50,825,784)
Automatically annotated and not reviewed.

**UniRef**
Sequence clusters

**UniParc**
Sequence archive

**Proteomes**

Supporting data

Literature citations

Cross-ref. databases

Taxonomy

Diseases
XXX

Subcellular locations

Keywords

**News**

Forthcoming
Planned chan

UniProt releas
Life (and dea
variation files

UniProt releas
Pseudo-allerg
access to Uni
of human var

News arch

## Getting started

### Text search
Our basic text search allows you to search all the resources available

### BLAST
Find regions of similarity between your sequences

## UniProt data

### Download latest release
Get the UniProt data

### Statistics
View Swiss-Prot and TrEMBL statistics

### How to cite us

## Protein sp

15

We are going to do ID mapping

16

# Upload Lists

## 1. Provide your identifiers

http://cys.bios.niu.edu/yyin/teach/PBB/at-id.txt

```
At1g24735
At3g61990
At3g62000
At1g67990
At1g67980
At4g26220
At1g15950
At1g80820
At1g76470
At2g02400
```

**OR** upload your own file: [ Choose File ] No file chosen

## 2. Select options

Choose TAIR here and UniProtKB here

From

To

[ TAIR ⬍ ]   [ UniProtKB ⬍ ]   [ Go ]

These are UniProt IDs



BLAST    Align    Retrieve/ID Mapping                                    Help    Contact

Show help for UniProtKB

**Results**

61 out of 61 TAIR identifiers were successfully mapped to 79 UniProtKB IDs.

🛒 Basket ▾

Filter by[i]

✏ Columns   🔧 BLAST   ≡ Align   ⬇ Download   🛒 Add to basket

◀ 1 to 25 of 79 ▶   Show  25 ▾

⭐ Reviewed (40)
Swiss-Prot

📄 Unreviewed (39)
TrEMBL

Popular
organisms
A. thaliana (79)

View by

Taxonomy
Keywords
Gene Ontology
Enzyme class
Pathway

UniRef

Your results in sequence
clusters with identity of:
100%, 90% or 50%

| ☐ | Your list:...APA6( | Entry | Entry name ⬍ | | Protein names ⬍ ⏩ | Gene names ⬍ | Organism ⬍ | Length ⬍ | ✏ |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | AT2G37040 | P35510 | PAL1_ARATH | ⭐ | **Phenylalanine ammonia-lyase 1** | **PAL1**, At2g37040, T1J8.22 | Arabidopsis thaliana (Mouse-ear cress) | 725 | |
| ☐ | AT3G53260 | P45724 | PAL2_ARATH | ⭐ | **Phenylalanine ammonia-lyase 2** | **PAL2**, At3g53260, T4D2.190 | Arabidopsis thaliana (Mouse-ear cress) | 717 | |
| ☐ | AT5G04230 | F4JW69 | F4JW69_ARATH | 📄 | **Phenylalanine ammonia-lyase** | **PAL3**, At5g04230 | Arabidopsis thaliana (Mouse-ear cress) | 698 | |
| ☐ | AT5G04230 | P45725 | PAL3_ARATH | ⭐ | **Phenylalanine ammonia-lyase 3** | **PAL3**, At5g04230, F21E1_150 | Arabidopsis thaliana (Mouse-ear cress) | 694 | |
| ☐ | AT3G10340 | Q9SS45 | PAL4_ARATH | ⭐ | **Phenylalanine ammonia-lyase 4** | **PAL4**, At3g10340, | Arabidopsis thaliana | 707 | |

18

Select the PAL proteins and align them

# Results

61 out of 61 TAIR identifiers were successfully mapped to 79 UniProtKB IDs.

🛒 Basket ▾

**Filter by**[i]

🔖 Reviewed (40)
Swiss-Prot

☐ Unreviewed (39)
TrEMBL

**Popular organisms**
A. thaliana (79)

**View by**

Taxonomy
Keywords
Gene Ontology
Enzyme class
Pathway

**UniRef**

Your results in sequence clusters with identity of: 100%, 90% or 50%

✎ Columns    🔍 BLAST    ☰ Align    ⬇ Download    🛒 Add to basket

◀ 1 to 25 of 79 ▶    Show 25 ⇕

| ☐ | Your list:...APA6C | Entry | Entry name ⬍ | | Protein names ⬍ ⟫ | Gene names ⬍ | Organism ⬍ | Length ⬍ | ✎ |
|---|---|---|---|---|---|---|---|---|---|
| ☑ | AT2G37040 | P35510 | PAL1_ARATH | 🔖 | Phenylalanine ammonia-lyase 1 | PAL1, At2g37040, T1J8.22 | Arabidopsis thaliana (Mouse-ear cress) | 725 | |
| ☑ | AT3G53260 | P45724 | PAL2_ARATH | 🔖 | Phenylalanine ammonia-lyase 2 | PAL2, At3g53260, T4D2.190 | Arabidopsis thaliana (Mouse-ear cress) | 717 | |
| ☑ | AT5G04230 | F4JW69 | F4JW69_ARATH | 📄 | Phenylalanine ammonia-lyase | PAL3, At5g04230 | Arabidopsis thaliana (Mouse-ear cress) | 698 | |
| ☑ | AT5G04230 | P45725 | PAL3_ARATH | 🔖 | Phenylalanine ammonia-lyase 3 | PAL3, At5g04230, F21E1_150 | Arabidopsis thaliana (Mouse-ear cress) | 694 | |
| ☑ | AT3G10340 | Q9SS45 | PAL4_ARATH | 🔖 | Phenylalanine ammonia-lyase 4 | PAL4, At3g10340, F14P13.6 | Arabidopsis thaliana (Mouse-ear | 707 | |

Clustal omega program will be called to align the selected protein seqs
May take 1 min to finish

19

This is the MSA result page

Toggle these options on will add colors in the alignment

☑ ALIGNMENT

☐ TREE

☐ RESULT INFO

# Highlight

**Annotation**
- ☐ Sequence conflict
- ☐ Binding site
- ☐ Chain
- ☐ Active site
- ☐ Modified residue
- ☐ Cross-link

**Amino acid properties**
- ☐ Similarity
- ☐ Hydrophobic
- ☐ Negative
- ☐ Positive
- ☐ Aliphatic
- ☐ Tiny
- ☐ Aromatic
- ☐ Charged
- ☐ Small
- ☐ Polar
- ☐ Big
- ☐ Serine Threonine

## Alignment

🖨 How to print an alignment in color

```
P35510  PAL1_ARATH   1   MEINGAHKSNGGGVDAMLCGGDIKTKNMVI--NAEDPLNWGAAAEQMKGSHLDEVKRMVA    58
P45724  PAL2_ARATH   1   ----------MDQIEAMLCGGGEKTKVAVTTKTLADPLNWGLAADQMKGSHLDEVKKMVE    50
F4JW69  F4JW69_ARATH 1   --------------------MEFR---QPNATALSDPLNWNVAAEALKGSHLEEVKKMVK    37
P45725  PAL3_ARATH   1   --------------------MEFR---QPNATALSDPLNWNVAAEALKGSHLEEVKKMVK    37
Q9SS45  PAL4_ARATH   1   --------------------MELCNQNNHITAVSGDPLNWNATAEALKGSHLDEVKRMVK    40
                                             *****   :*: :*****:***:**

P35510  PAL1_ARATH   59  EFRKPVVNLGGETLTIGQVAAISTIGNSVKVELSETARAGVNASSDWVMESMNKGTDSYG   118
P45724  PAL2_ARATH   51  EYRRPVVNLGGETLTIGQVAAISTVGGSVKVELAETSRAGVKASSDWVMESMNKGTDSYG   110
F4JW69  F4JW69_ARATH 38  DYRKGTVQLGGETLTIGQVAAVAS--GGPTVELSEEARGGVKASSDWVMESMNRDTDTYG    95
P45725  PAL3_ARATH   38  DYRKGTVQLGGETLTIGQVAAVAS--GGPTVELSEEARGGVKASSDWVMESMNRDTDTYG    95
Q9SS45  PAL4_ARATH   41  EYRKEAVKLGGETLTIGQVAAVARGGGGSTVELAEEARAGVKASSEWVMESMNRGTDSYG   100
                         ::*: .*:**************::       . .***:* :*.**:***:******: **:**

P35510  PAL1_ARATH   119 VTTGFGATSHRRTKNGVALQKELIRFLNAGIFGSTK---ETSHTLPHSATRAAMLVRINT   175
P45724  PAL2_ARATH   111 VTTGFGATSHRRTKNGTALQTELIRFLNAGIFGNTK---ETCHTLPQSATRAAMLVRVNT   167
F4JW69  F4JW69_ARATH 96  ITTGFGSSSRRRTDQGAALQKELIRYLNAGIFATGNEDDDRSNTLPRPATRAAMLIRVNT   155
P45725  PAL3_ARATH   96  ITTGFGSSSRRRTDQGAALQKELIRYLNAGIFATGNEDDDRSNTLPRPATRAAMLIRVNT   155
Q9SS45  PAL4_ARATH   101 VTTGFGATSHRRTKQGGALQNELIRFLNAGIFGPGAG--DTSHTLPKPTTRAAMLVRVNT   158
                         :*****::*:***.:* ***.****:*****.       :  ..***: :******:*:**

P35510  PAL1_ARATH   176 LLQGFSGIRFEILEAITSFLNNNITPSLPLRGTITASGDLVPLSYIAGLLTGRPNSKATG   235
P45724  PAL2_ARATH   168 LLQGYSGIRFEILEAITSLLNHNISPSLPLRGTITASGDLVPLSYIAGLLTGRPNSKATG   227
F4JW69  F4JW69_ARATH 156 LLQGYSGIRFEILEAITTLLNCKITPLLPLRGTITASGDLVPLSYIAGFLIGRPNSRSVG   215
P45725  PAL3_ARATH   156 LLQGYSGIRFEILEAITTLLNCKITPLLPLRGTITASGDLVPLSYIAGFLIGRPNSRSVG   215
Q9SS45  PAL4_ARATH   159 LLQGYSGIRFEILEAITKLLNHEITPCLPLRGTITASGDLVPLSYIAGLLTGRPNSKAVG   218
                         ****:*************.:** :*:* ***********************:* *****::.*

P35510  PAL1_ARATH   236 PNGEALTAEEAFKLAGISSGFFDLQPKEGLALVNGTAVGSGMASMVLFETNVLSVLAEIL   295
P45724  PAL2_ARATH   228 PDGESLTAKEAFEKAGISTGFFDLQPKEGLALVNGTAVGSGMASMVLFEANVQAVLAEVL   287
F4JW69  F4JW69_ARATH 216 PSGEILTALEAFKLAGVS-SFFELRPKEGLALVNGTAVGSALASTVLYDANILVVFSEVA   274
P45725  PAL3_ARATH   216 PSGEILTALEAFKLAGVS-SFFELRPKEGLALVNGTAVGSALASTVLYDANILVVFSEVA   274
Q9SS45  PAL4_ARATH   219 PSGETLTASEAFKLAGVS-SFFELQPKEGLALVNGTAVGSGLASTVLFDANILAVLSEVM   277
                         *.** *** ***: **:* .**:*:****************.:** **:::*:  *::*:
```

20

Go back to the protein list page
Selecting one protein will enable the BLAST button



Choose advanced will allow to change BLAST parameters

Here you can make changes

We are going to search UniProt proteomes for human protein set

Click on Advanced you will see a pop-out window

UniProt

Proteomes ▾ [                                  ] Advanced ▾ [🔍]

BLAST   Align   Upload Lists                                              Help   Contact

Welcome to the new UniProt website! We hope you enjoy the new design. If you're not quite ready yet, you can still go back to the old site.

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

| UniProtKB | UniRef | UniParc | Proteomes | News |
|---|---|---|---|---|
| Swiss-Prot (546,238) ⭐ Manually annotated and reviewed. | Sequence clusters | Sequence archive | | Ubiquitin caught at its own game \| New human variant types available on the FTP site UniProt release 2014_08 |
| | Supporting data | | | Lark or owl? PER3 is the answer \| Cross-references to CCDS and GeneReviews \| UniParc cross-references with protein and gene |
| TrEMBL (82,1... Aut annotat rev | Literature citations | Taxonomy | Subcellular locations | |

Searching in **Proteomes**                                              ✖

Term
[Organism [OS]    ▾]  [Human [9606]                          ]   🗑

Term
[AND ▾]  [All              ▾]  [                              ]   🗑  ➕

Here you can specify search terms

🔍

23

Click here to get help

BLAST  Align  Upload Lists                                                                        Help  Contact

Show help for Proteomes

## Results

○ Repeat search in UniProtKB (140,991)                                          🛒 Basket ▾

### Filter by

⬇ Download                                                        1 to **1** of **1**   Show 25 ▾

📕 1 Reference proteomes

### Map To

| ☐ | Proteome ID | Organism ◆ | Last modified ◆ | Protein Count |
|---|---|---|---|---|
| ☐ | UP000005640 | 📕 **Homo sapiens (Human)** | 2014-07-09 | 68049 |

UniProtKB

### Demo

1 to **1** of **1**   Show 25 ▾

▶ Help video

Click here to open a new page

## Results                                                                        🛒 Basket ▾

### Filter by

📊 Reviewed (20,187)
Swiss-Prot

📊 Unreviewed (47,862)
TrEMBL

### Popular organisms

Human (68,049) ✖

### Proteomes

UP000005640 (68,049) ✖

### View by

Taxonomy
Keywords
Gene Ontology
Enzyme class
Pathway

### UniRef

Your results in sequence clusters with identity of:
100%, 90% or 50%

### Demo

▶ Help video

✏ Columns  🔍 BLAST  ≡ Align  ⬇ Download  🛒 Add to basket        ◀ 1 to 25 of 68,049 ▶  Show 25 ▾

| ☐ | Entry | Entry name ◆ | Protein names ◆ | Gene names ◆ | Organism ◆ | Length ◆ | ✏ |
|---|---|---|---|---|---|---|---|
| ☐ | P31946 | 1433B_HUMAN | 14-3-3 protein beta/alpha | **YWHAB** | Homo sapiens (Human) | 246 | |
| ☐ | P62258 | 1433E_HUMAN | 14-3-3 protein epsilon | **YWHAE** | Homo sapiens (Human) | 255 | |
| ☐ | Q04917 | 1433F_HUMAN | 14-3-3 protein eta | **YWHAH**, YWHA1 | Homo sapiens (Human) | 246 | |
| ☐ | P61981 | 1433G_HUMAN | 14-3-3 protein gamma | **YWHAG** | Homo sapiens (Human) | 247 | |
| ☐ | P31947 | 1433S_HUMAN | 14-3-3 protein sigma | **SFN**, HME1 | Homo sapiens (Human) | 248 | |
| ☐ | P27348 | 1433T_HUMAN | 14-3-3 protein theta | **YWHAQ** | Homo sapiens (Human) | 245 | |
| ☐ | P63104 | 1433Z_HUMAN | 14-3-3 protein zeta/delta | **YWHAZ** | Homo sapiens (Human) | 245 | |
| ☐ | P30443 | 1A01_HUMAN | HLA class I histocompatibility anti... | **HLA-A**, HLAA | Homo sapiens (Human) | 365 | |
| ☐ | P01892 | 1A02_HUMAN | HLA class I histocompatibility anti... | **HLA-A**, HLAA | Homo sapiens (Human) | 365 | |
| ☐ | P04439 | 1A03_HUMAN | HLA class I histocompatibility anti... | **HLA-A**, HLAA | Homo sapiens (Human) | 365 | |
| ☐ | P13746 | 1A11_HUMAN | HLA class I histocompatibility anti... | **HLA-A**, HLAA | Homo sapiens (Human) | 365 | |
| ☐ | Q96QU6 | 1A1L1_HUMAN | 1-aminocyclopropane-1-carboxylate s... | **ACCS**, PHACS | Homo sapiens (Human) | 501 | |

24

# Gene Ontology

http://geneontology.org/page/documentation

The Gene Ontology (GO) project is a collaborative effort to address the need for consistent descriptions of gene products in different databases

The project began as a collaboration between three model organism databases, FlyBase (*Drosophila*), the *Saccharomyces* Genome Database (SGD) and the Mouse Genome Database (MGD), in 1998

Three structured controlled vocabularies (ontologies) that describe gene products in terms of their associated biological processes, cellular components and molecular functions in a species-independent manner.

There are three separate aspects to this effort:

1, the development and maintenance of the ontologies themselves;
2, the annotation of gene products, which entails making associations between the ontologies and the genes and gene products in the collaborating databases; and
3, development of tools that facilitate the creation, maintenance and use of ontologies.

# The scope of GO

Gene Ontology covers three domains:

**cellular component**, the parts of a cell or its extracellular environment;

**molecular function**, the elemental activities of a gene product at the molecular level, such as binding or catalysis;

**biological process**, operations or sets of molecular events with a defined beginning and end, pertinent to the functioning of integrated living units: cells, tissues, organs, and organisms

GO is not a database of gene sequences, nor a catalog of gene products. Rather, GO describes how gene products behave in a cellular context.

GO is not a dictated standard, mandating nomenclature across databases. Groups participate because of self-interest, and cooperate to arrive at a consensus.

GO is not a way to unify biological databases (i.e. GO is not a 'federated solution'). Sharing vocabulary is a step towards unification, but is not, in itself, sufficient.

The structure of GO can be described in terms of a graph, where each GO term is a node, and the relationships between the terms are edges between the nodes. GO is loosely hierarchical, with 'child' terms being more specialized than their 'parent' terms, but unlike a strict hierarchy, a term may have more than one parent term



http://geneontology.org/page/ontology-structure

```
id: GO:0000016
name: lactase activity namespace: molecular_function
def: "Catalysis of the reaction: lactose + H2O = D-glucose + D-galactose." [EC:3.2.1.108]
synonym: "lactase-phlorizin hydrolase activity" BROAD [EC:3.2.1.108]
synonym: "lactose galactohydrolase activity" EXACT [EC:3.2.1.108]
xref: EC:3.2.1.108
xref: MetaCyc:LACTASE-RXN
xref: Reactome:20536
is_a: GO:0004553 ! hydrolase activity, hydrolyzing O-glycosyl compounds
```

## What can I do with GO?

### What can I do with GO?

One of the most popular uses of GO is to find significant shared GO terms (or parents of those GO terms) that are annotated to genes in a particular query set (e.g. a set of genes that are overexpressed in a microarray experiment). This process helps you to find out what those genes may have in common and is known as a **GO enrichment analysis**.

GO is also used for purposes as diverse as:

- integrating proteomic information from different organisms;
- assigning functions to protein domains;
- finding functional similarities in genes that are overexpressed or underexpressed in diseases and as we age;
- analysing groups of genes that are co-expressed during development;
- developing automated ways of deriving information about gene function from the literature;
- verifying models of genetic, metabolic and product interaction networks.

The GO tools web page lists the tools that you can use to analyse the data from GO.

http://www.ebi.ac.uk/training/online/course/go-quick-tour/what-can-i-do-go

Enrichment analysis: use statistical test e.g. Fisher exact test

Example: in human genome background (20,000 gene total), 40 genes are involved in p53 signaling pathway. A given gene list has found that 3 out of 300 belong to p53 signaling pathway. Then we ask the question if 3/300 is more than random chance comparing to the human background of 40/20000



http://david.abcc.ncifcrf.gov/helps/functional_annotation.html#E4

# UniProt-GO annotation (GOA)

# UniProt-GOA format

The *reference* used to make the annotation (e.g. a journal article)
An *evidence code* denoting the type of evidence upon which the annotation is based
The date and the creator of the annotation

```
Gene product: Actin, alpha cardiac muscle 1, UniProtKB:P68032
GO term: heart contraction ; GO:0060047 (biological process)
Evidence code: Inferred from Mutant Phenotype (IMP) Reference: PMID 17611253
Assigned by: UniProtKB, June 6, 2008
```

# The idea of GO annotation for new sequences

If you have a new genome/transcriptome sequenced, how do you perform a GO annotation for it?

1. Find a closet model organism which has been annotated by GO
2. BLAST your data against this closest organism
3. Transfer the GO annotation of the best match to your query sequences

For instance, if we want to annotate fern transcriptome with GO function descriptions ….

1. Find Arabidopsis UniProt protein dataset
2. Find the Arabidopsis GOA association file
3. BLASTx fern reads (or assembled UniGenes) against the UniProt set
4. Analyze BLAST result to link fern reads GO terms

# Hands on practice 2: GO annotation

http://geneontology.org/

# http://amigo1.geneontology.org/cgi-bin/amigo/blast.cgi

*the Gene Ontology*

| Search | Browse | BLAST | Homolog Annotations | Tools & Resources | Help |

Search GO [                    ]  ⊙ terms  ○ genes or proteins  ☐ exact match  [Submit]

## BLAST Search

The sequence search is performed using either BLASTP or BLASTX (from the WU-BLAST package), depending on the type of the input sequence.

## BLAST Query

### Enter your query ❓

Enter a UniProtKB accession **or** upload a text file of queries **or** paste in FASTA sequence(s)

UniProtKB accession: [                ]

Text file (maximum file size 500K): [ Choose File ] No file chosen

FASTA sequence(s):
Sequences should be separated with an empty line.

```
>AT5G22740.1|AT5G22740.1|cs1A
MDGVSPKFVLPETFDGVRMEITGQLGMIWELVKAPVIVPLLQLAVYICLL
MSVMLLCERVYMGIVIVLVKLFWKKPDKRYKFEPIHDDEELGSSNFPVVL
VQIPMFNEREVYKLSIGAACGLSWPSDRLVIQVLDDSTDPTVKQMVEVEC
QRWASKGINIRYQIRENRVGYKAGALKEGLKRSYVKHCEYVVIFDADFQP
EPDFLRRSIPFLMHNPNIALVQARWRFVNSDECLLTRMQEMSLDYHFTVE
QEVGSSTHAFFGFNGTAGIWRIAAINEAGGWKDRTTVEDMDLAVRASLRG
WKFLYLGDLQVKSLPSTFRAFRFQQHRWSCGPANLFRKMVMELIVRNKKV
```

Get an example protein sequence file from http://cys.bios.niu.edu/yyin/teach/PBB/csl-pr.fa

35

# BLAST Query Submission

## Success!

Your job has been successfully submitted to the BLAST queue.

Please be patient as your job may take several minutes to complete. This page will automatically refresh with the BLAST results when the job is done.

Try retrieving your job now

## Query Summary

Your job contains 2 sequences.

Parameters
Threshold: 0.1
Maximum number of alignments shown: 50
BLAST filter: on

AmiGO version: 1.8

Try AmiGO Labs

## High Scoring Gene Products

Select all | Clear all | Perform an action with this page's selected gene products... ▼ | Go!

| | Symbol, full name | Information | | P value |
|---|---|---|---|---|
| ☐ | CSLA02<br>cellulose synthase-like A02 | BLAST match ⬇ view associations ➡ BLAST with CSLA02 ➡ | **protein** from *Arabidopsis thaliana* | 3.6e-295 |
| ☐ | ATCSLA09 | BLAST match ⬇ view associations ➡ BLAST with ATCSLA09 ➡ | **protein** from *Arabidopsis thaliana* | 4.5e-217 |
| ☐ | CSLA03<br>cellulose synthase-like A3 | BLAST match ⬇ view associations ➡ BLAST with CSLA03 ➡ | **protein** from *Arabidopsis thaliana* | 4.1e-191 |
| ☐ | ATCSLA15 | BLAST match ⬇ view associations ➡ BLAST with ATCSLA15 ➡ | **protein** from *Arabidopsis thaliana* | 2.9e-183 |
| ☐ | CSLA07<br>cellulose synthase like | BLAST match ⬇ view associations ➡ BLAST with CSLA07 ➡ | **protein** from *Arabidopsis thaliana* | 1.2e-182 |
| ☐ | CSLA10<br>cellulose synthase-like A10 | BLAST match ⬇ view associations ➡ BLAST with CSLA10 ➡ | **protein** from *Arabidopsis thaliana* | 3.3e-173 |
| ☐ | CSLA01<br>cellulose synthase-like A01 | BLAST match ⬇ view associations ➡ BLAST with CSLA01 ➡ | **protein** from *Arabidopsis thaliana* | 4.0e-170 |
| ☐ | CSLA14<br>cellulose synthase like A14 | BLAST match ⬇ view associations ➡ BLAST with CSLA14 ➡ | **protein** from *Arabidopsis thaliana* | 7.6e-167 |

6

This is easy. Now let's try to get a list of differentially expressed genes and then find what's common in this list of genes in terms of functions.

We're gonna use NCBI GEO website to get the gene list and then feed the gene list to GO enrichment analysis tools

Go to NCBI home page, search GEO DataSets with keyword "liver cancer", and hit search

Top hits are always GEO DataSets, let's choose the 3rd one, hit Analyze DataSet

Choose "Compare 2 sets of samples"

Choose "Value means difference"
Choose "8+ fold"
Choose "higher"

Then go to Step 2

Select to choose group A: three samples for COP 1 depletion and Huh7 cell line

Group B: three samples for negative control and Huh7 cell line

Hit ok, and go to Step 3

DataSet Record GDS4831: (Expression Profiles) (Data Analysis Tools) (Sample Subsets)

| Title: | COP1 depletion effect on hepatocellular carcinoma cell lines |
|---|---|
| Summary: | Analysis of Huh7, HepG2, and Hep3B hepatocellular carcinoma (HCC) cells depleted for the ubiquitin modulator COP1. COP1 regulates p53 activity by ubiquitination. p53 is wild type in HepG2, mutated in Huh7, and lacking in Hep3B. Results provide insight into the role of COP1 in HCC pathoge |
| Organism: | Homo sapiens |
| Platform: | GPL6883: Illumina HumanRef-8 v3.0 expression beadchip |
| Citation: | Lee YH, Andersen JB, Song HT, Judge AD et al. Definition of ubiquitination modulator COP1 as a novel therap 1;70(21):8264-9. PMID: 20959491 |
| Reference Series: | GSE21955 |
| Value type: | count |

Sample count: 22
Series published: 20

Click on accessions to select samples individually, click on colored blocks and then on blinking arrows to select groups of samples.

Ok
Reset
Cancel

| Samples, Group A | Factors | | Samples, Group B |
|---|---|---|---|
| | protocol | cell line | |
| GSM545954 | COP1 depletion | Huh7 | GSM545954 |
| GSM545955 | | | GSM545955 |
| GSM545956 | | | GSM545956 |
| GSM545960 | | HepG2 | GSM545960 |
| GSM545961 | | | GSM545961 |
| GSM545962 | | | GSM545962 |
| GSM545963 | | | GSM545963 |
| GSM545968 | | Hep3B | GSM545968 |
| GSM545969 | | | GSM545969 |
| GSM545970 | | | GSM545970 |
| GSM545971 | | | GSM545971 |
| GSM545957 | negative control | Huh7 | GSM545957 |
| GSM545958 | | | GSM545958 |
| GSM545959 | | | GSM545959 |
| GSM545964 | | HepG2 | GSM545964 |
| GSM545965 | | | GSM545965 |
| GSM545966 | | | GSM545966 |
| GSM545967 | | | GSM545967 |
| GSM545972 | | Hep3B | GSM545972 |
| GSM545973 | | | GSM545973 |
| GSM545974 | | | GSM545974 |
| GSM545975 | | | GSM545975 |

Data Analysis

Find genes

Compare 2 sets of samples [?]

Cluster heatmaps

Experiment design and value distribution

**Step 1:** Select test and significance level

Value means difference ⇕   Ā vs B̄: 8+ fold ⇕   higher ⇕

**Step 2:** Select which Samples to put in Group A and Group B

**Step 3:** Query Group A vs. B

40

Total 256 gene profiles are found with 8+ fold higher expression in COP 1 depletion than in negative control in Huh7 cell line

To get the list of genes, choose Gene database and hit Find items

Display Settings: ☑ Summary, 20 per page, Sorted by Default order

Send to: ☑

Filters: Manage Filters

**Results: 1 to 20 of 256**

<< First < Prev Page 1 of 13 Next > Last >>

☐ UBE2G2 - COP1 depletion effect on hepatocellular carcinoma cell lines
1. Annotation: UBE2G2, ubiquitin-conjugating enzyme E2G 2
Organism: Homo sapiens
Reporter: GPL6883, ILMN_2297824 (ID_REF), GDS4831, NM_182688
DataSet type: Expression profiling by array, count, 22 samples
ID: 104862213
GEO DataSets    Gene    UniGene    Profile neighbors    Chromosome neighbors    Homologene neighbors

☐ MS4A6A - COP1 depletion effect on hepatocellular carcinoma cell lines
2. Annotation: MS4A6A, membrane-spanning 4-domains, subfamily A, member 6A
Organism: Homo sapiens
Reporter: GPL6883, ILMN_2359800 (ID_REF), GDS4831, NM_152851
DataSet type: Expression profiling by array, count, 22 samples
ID: 104862285
GEO DataSets    Gene    UniGene    Profile neighbors    Chromosome neighbors    Homologene neighbors

☐ VIP - COP1 depletion effect on hepatocellular carcinoma cell lines
3. Annotation: VIP, vasoactive intestinal peptide
Organism: Homo sapiens
Reporter: GPL6883, ILMN_1794638 (ID_REF), GDS4831, NM_194435
DataSet type: Expression profiling by array, count, 22 samples
ID: 104862502
GEO DataSets    Gene    UniGene    Chromosome neighbors    Homologene neighbors

☐ HIST1H2BO - COP1 depletion effect on hepatocellular carcinoma cell lines
4. Annotation: HIST1H2BO, histone cluster 1, H2bo

**Profile data**

Download profile data

**Profile pathways**

Find pathways

**Find related data**

Database: Gene

based on gene annotation from platform pr
accessio

Find items

**Recent activity**

Turn Off

COP1 depletion effect on hepatocellula
carcinoma cell lines                    GDS

(GDS4831[ACCN]) AND GDS[filter] (1
GDS

41

Total 225 genes correspond to 256 gene profiles
To download the list of Gene IDs, hit Send to, choose UI list as format and hit Create file



A file named "gene_result.txt" will be automatically downloaded to your local computer
Find out where it is downloaded to, open it using notepad++

View the file using notepad++

Next we will use DAVID to perform function enrichment analysis

# The **D**atabase for **A**nnotation, **V**isualization and **I**ntegrated**D**iscovery (**DAVID** )



Hit start analysis

Upload the list of Gene IDs

Select ENTREZ_GENE_ID

Click on Gene list

This allows you to view functional annotation from various resources including GO

Check the submitted gene list

If you have clicked on Functional Annotation tool, you are at this page

Uncheck this

**Upload** **List** **Background**

**Gene List Manager**

Select to limit annotations by one or more species   Help

- Use All Species –
Homo sapiens(225)
Unknown(1)

Select Species

**List Manager**   Help

gene_result (6)

Select List to:

Use     Rename

Remove   Combine

Show Gene List

View Unmapped Ids

**Annotation Summary Results**

Help and Tool Manual

**Current Gene List: gene_result (6)**        **225 DAVID IDs**
**Current Background: Homo sapiens**        **Check Defaults** ☐     Clear All

⊞ **Disease** (0 selected)
⊞ **Functional_Categories** (0 selected)
⊞ **Gene_Ontology** (0 selected)
⊞ **General Annotations** (0 selected)
⊞ **Literature** (0 selected)
⊞ **Main_Accessions** (0 selected)
⊞ **Pathways** (0 selected)
⊞ **Protein_Domains** (0 selected)
⊞ **Protein_Interactions** (0 selected)
⊞ **Tissue_Expression** (0 selected)

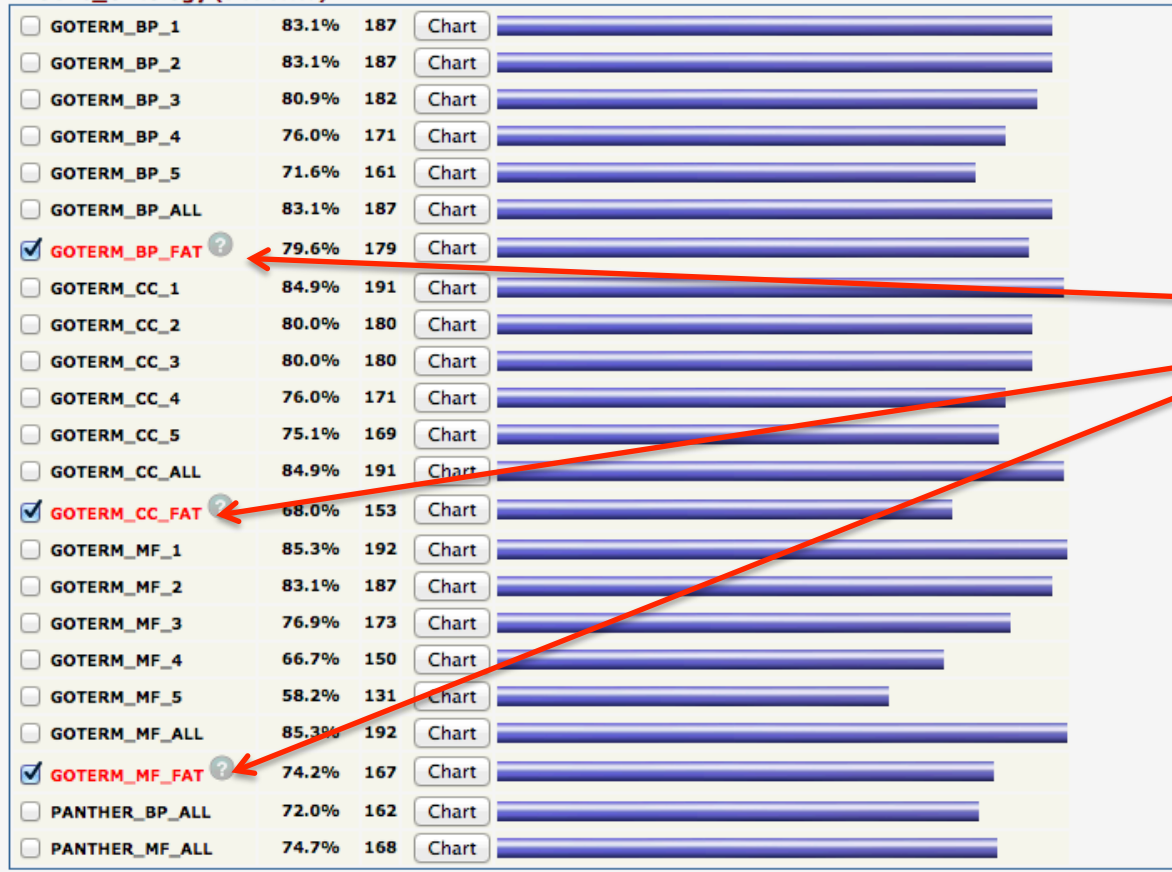***Red annotation categories denote DAVID defined defaults***

**Combined View for Selected Annotation**

Functional Annotation Clustering

Functional Annotation Chart

Functional Annotation Table

All these can be changed by users (to show or not to show and show what)

**Current Background: Homo sapiens**     Check Defaults ☐     Clear All

⊞ **Disease** (0 selected)
⊞ **Functional_Categories** (0 selected)
⊟ **Gene_Ontology** (3 selected)

| | | | | |
|---|---|---|---|---|
| ☐ GOTERM_BP_1 | 83.1% | 187 | Chart | |
| ☐ GOTERM_BP_2 | 83.1% | 187 | Chart | |
| ☐ GOTERM_BP_3 | 80.9% | 182 | Chart | |
| ☐ GOTERM_BP_4 | 76.0% | 171 | Chart | |
| ☐ GOTERM_BP_5 | 71.6% | 161 | Chart | |
| ☐ GOTERM_BP_ALL | 83.1% | 187 | Chart | |
| ☑ GOTERM_BP_FAT ⊙ | 79.6% | 179 | Chart | |
| ☐ GOTERM_CC_1 | 84.9% | 191 | Chart | |
| ☐ GOTERM_CC_2 | 80.0% | 180 | Chart | |
| ☐ GOTERM_CC_3 | 80.0% | 180 | Chart | |
| ☐ GOTERM_CC_4 | 76.0% | 171 | Chart | |
| ☐ GOTERM_CC_5 | 75.1% | 169 | Chart | |
| ☐ GOTERM_CC_ALL | 84.9% | 191 | Chart | |
| ☑ GOTERM_CC_FAT ⊙ | 68.0% | 153 | Chart | |
| ☐ GOTERM_MF_1 | 85.3% | 192 | Chart | |
| ☐ GOTERM_MF_2 | 83.1% | 187 | Chart | |
| ☐ GOTERM_MF_3 | 76.9% | 173 | Chart | |
| ☐ GOTERM_MF_4 | 66.7% | 150 | Chart | |
| ☐ GOTERM_MF_5 | 58.2% | 131 | Chart | |
| ☐ GOTERM_MF_ALL | 85.3% | 192 | Chart | |
| ☑ GOTERM_MF_FAT ⊙ | 74.2% | 167 | Chart | |
| ☐ PANTHER_BP_ALL | 72.0% | 162 | Chart | |
| ☐ PANTHER_MF_ALL | 74.7% | 168 | Chart | |

⊞ **General Annotations** (0 selected)
⊞ **Literature** (0 selected)
⊞ **Main_Accessions** (0 selected)
⊞ **Pathways** (0 selected)
⊞ **Protein_Domains** (0 selected)
⊞ **Protein_Interactions** (0 selected)
⊞ **Tissue_Expression** (0 selected)

***Red annotation categories denote DAVID defined defaults***

**Combined View for Selected Annotation**

Functional Annotation Clustering
Functional Annotation Chart
Functional Annotation Table

Select just GO

Click here will open a new window to show the 225 differentially expressed genes are enriched in what GO

48

# Functional Annotation Chart

**Current Gene List: gene_result (6)**
**Current Background: Homo sapiens**
**225 DAVID IDs**

⊞ Options

Genes are enriched in what GO categories (compared to the genome background)?

[ Rerun Using Options ]  [ Create Sublist ]

50 chart records

🖫 Download File

| Sublist | Category | Term | RT | Genes | Count | % | P-Value | Benjamini |
|---------|----------|------|----|-------|-------|---|---------|-----------|
| ☐ | GOTERM_BP_FAT | integrin-mediated signaling pathway | RT | | 7 | 3.1 | 3.1E-4 | 3.6E-1 |
| ☐ | GOTERM_CC_FAT | plasma membrane | RT | | 65 | 28.9 | 6.2E-4 | 1.4E-1 |
| ☐ | GOTERM_CC_FAT | integral to plasma membrane | RT | | 28 | 12.4 | 7.5E-4 | 8.8E-2 |
| ☐ | GOTERM_CC_FAT | intrinsic to plasma membrane | RT | | 28 | 12.4 | 1.0E-3 | 8.3E-2 |
| ☐ | GOTERM_BP_FAT | cell surface receptor linked signal transduction | RT | | 38 | 16.9 | 5.8E-3 | 9.8E-1 |
| ☐ | GOTERM_BP_FAT | G-protein coupled receptor protein signaling pathway | RT | | 26 | 11.6 | 6.5E-3 | 9.6E-1 |
| ☐ | GOTERM_CC_FAT | nucleosome | RT | | 5 | 2.2 | 6.6E-3 | 3.4E-1 |
| ☐ | GOTERM_BP_FAT | positive regulation of protein kinase activity | RT | | 9 | 4.0 | 9.4E-3 | 9.7E-1 |
| ☐ | GOTERM_BP_FAT | positive regulation of kinase activity | RT | | 9 | 4.0 | 1.1E-2 | 9.6E-1 |
| ☐ | GOTERM_CC_FAT | integral to membrane | RT | | 78 | 34.7 | 1.3E-2 | 4.8E-1 |
| ☐ | GOTERM_BP_FAT | cell activation | RT | | 10 | 4.4 | 1.4E-2 | 9.6E-1 |
| ☐ | GOTERM_BP_FAT | positive regulation of transferase activity | RT | | 9 | 4.0 | 1.4E-2 | 9.5E-1 |
| ☐ | GOTERM_BP_FAT | leukocyte activation | RT | | 9 | 4.0 | 1.5E-2 | 9.3E-1 |
| ☐ | GOTERM_CC_FAT | plasma membrane part | RT | | 38 | 16.9 | 1.6E-2 | 4.8E-1 |
| ☐ | GOTERM_BP_FAT | positive regulation of epithelial cell proliferation | RT | | 4 | 1.8 | 1.7E-2 | 9.3E-1 |
| ☐ | GOTERM_BP_FAT | activation of protein kinase activity | RT | | 6 | 2.7 | 1.7E-2 | 9.2E-1 |
| ☐ | GOTERM_BP_FAT | DNA packaging | RT | | 6 | 2.7 | 1.9E-2 | 9.2E-1 |
| ☐ | GOTERM_CC_FAT | protein-DNA complex | RT | | 5 | 2.2 | 1.9E-2 | 5.0E-1 |
| ☐ | GOTERM_CC_FAT | intrinsic to membrane | RT | | 79 | 35.1 | 2.2E-2 | 5.0E-1 |
| ☐ | GOTERM_BP_FAT | heart development | RT | | 8 | 3.6 | 2.4E-2 | 9.4E-1 |
| ☐ | GOTERM_BP_FAT | nucleosome assembly | RT | | 5 | 2.2 | 2.5E-2 | 9.4E-1 |
| ☐ | GOTERM_BP_FAT | chromatin assembly | RT | | 5 | 2.2 | 2.8E-2 | 9.4E-1 |
| ☐ | GOTERM_BP_FAT | locomotory behavior | RT | | 9 | 4.0 | 2.9E-2 | 9.4E-1 |
| ☐ | GOTERM_BP_FAT | leukocyte differentiation | RT | | 6 | 2.7 | 2.9E-2 | 9.3E-1 |

49

**Next lecture:** `EBI web resources II (ENSEMBL and InterPro)`