# Introduction to

# Biological Databases

Valencia, September 2011

Marta Bleda

Ignacio Medina (*Nacho*)

imedina@cipf.es

http://bioinfo.cipf.es/imedina

*Bioinformatics and Genomics Department*

*Centro de Investigacion Principe Felipe (CIPF)*

*Valencia, Spain*

# Index

- Introduction
- Sequence Databases
- Functional Annotation Databases
- Protein Databases
- Variation Databases
- Genome Databases and Browsers
- Exercises

# Introduction

- Last years has been an exponential increase in the number of biological databases and in their content.

- ***Nucleic Acids Research* online Molecular Biology Database Collection** is a public repository that lists principal *biological databases*

- Updated every year. The Nov-2010 update includes **1330** databases !!

**http://www3.oup.co.uk/nar/database/c/**

# Introduction

- These databases contain:

  - Data and results from experiments with microarrays, *NGS*, ...

  - Genes, transcripts and *EST* sequences

  - DNA variation and frequencies (*SNP*, mutations, indels, …)

  - Protein sequences, structures and variations

  - Functional information about what a gene/protein is doing in the cell

  - User interface to search, navigate and explore the genomes

# Sequence Databases
## Genome Reference Consortium (*GRC*)

The ***GRC*** is a collaborative effort and only works with input from the larger scientific community

We strive to work closely with external groups to gather all relevant data

The ***GRC*** is now working to create ***assemblies*** that better represent this ***diversity*** and provide more robust substrates for genome analysis
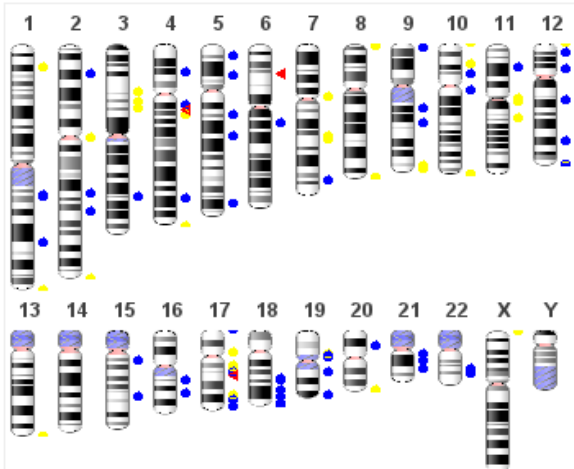
http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/human/index.shtml

# Sequence Databases
## European Bioinformatics Institute (*EBI*)

Mission:
• To provide freely available data and bioinformatics services to all facets of the scientific community in ways that promote scientific progress

• To contribute to the advancement of biology through basic investigator-driven research in bioinformatics

• To provide advanced bioinformatics training to scientists at all levels, from PhD students to independent investigators

• To help disseminate cutting-edge technologies to industry

Financiado por el EMBL, por tanto con dinero Europeo

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

# Sequence Databases
## Nat. Center for Biotech. Information (*NCBI*)

http://www.ncbi.nlm.nih.gov/guide/

Conjunto de herramientas y bases de datos para el estudio y análisis genómico y biomédico

Financiado por USA, en cierta forma compite con el EBI en objetivos y recursos

# Sequence Databases
## 1000 Genomes project

*1000 Genomes Project* is the first project to sequence the genomes of a large number of people, to provide a comprehensive resource on human genetic variation

http://www.1000genomes.org/

# Functional Annotation DDBB
## Overview

Some of the biological databases contains *Functional Information* of the genes and sequences

*Homo sapiens*  *Mus musculus*  *Rattus norvegicus*  *Gallus gallus*  *Danio rerio*  *Drosophila melanogaster*  *Caenorhabditis elegans*  *Saccharmoyces cerevisae*  *Arabidopsis thaliana*

UniProt/Swiss-Prot

UniProtKB/TrEMBL

Ensembl IDs

EntrezGene

Affymetrix

Agilent

**Gene IDs**

HGNC symbol

EMBL acc

RefSeq

PDB

Protein Id

IPI....

## Functional databases

**KEGG pathways**

**Reactome**

**Gene Ontology**

Biological Process
Molecular Function
Cellular Component

**Regulatory elements**

**miRNA**

**CisRed**

**Transcription Factor**

**Binding Sites**

**Biocarta pathways**

**Keywords Swissprot**

**InterPro Motifs**

**Gene Expression in tissues**

**Bioentities from literature:**

**Diseases terms**
**Chemical terms**

# Functional Annotation DDBB
## Gene Ontology (GO terms)

- The *Gene Ontology* project provides a **controlled vocabulary** to describe gene and gene product attributes in any organism

- Lastest version has **33808** terms (March, *2011*)

- The controlled vocabularies of terms are structured

**http://www.geneontology.org/**

# Functional Annotation DDBB
## Gene Ontology (GO terms)

**The three categories of GO**

**Molecular Function**

the tasks performed by individual gene products; examples are *transcription factor* and *DNA helicase*

**Biological Process**

broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

**Cellular Component**

subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *origin recognition complex*

⊟GO:0003673 : Gene Ontology (65883)
  ⊟ⓟ GO:0008150 : biological process (44405)
    ⊞ⓘ GO:0007610 : behavior (357)
    • ⓘ GO:0000004 : biological process unknown (7877)
    ⊟ⓘ GO:0009987 : cellular process (32672)
      ⊞ⓘ GO:0007154 : cell communication (5384)
      ⊞ⓘ GO:0008219 : cell death (744)
      ⊞ⓘ GO:0030154 : cell differentiation (464)
      ⊞ⓘ GO:0008151 : cell growth and/or maintenance (28802)
      ⊞ⓘ GO:0006928 : cell motility (911)
      ⊞ⓘ GO:0006944 : membrane fusion (257)
    ⊞ⓘ GO:0016265 : death (793)
    ⊞ⓘ GO:0007275 : development (4615)
    ⊞ⓘ GO:0008371 : obsolete (1581)
    ⊞ⓘ GO:0007582 : physiological processes (31124)
    ⊞ⓘ GO:0016032 : viral life cycle (115)
  ⊞ⓟ GO:0005575 : cellular component (32869)
  ⊞ⓟ GO:0003674 : molecular function (53910)

# Functional Annotation DDBB
## Gene Ontology (GO terms)

**GO is a DAG**
**(Directed Acyclic Graph)**

More general information

Levels

More detailed information

**terms are structured**

biological process
78842 genes

physiological process
55602 genes

cellular process
29557 genes

cell growth and/or maintenance
21215 genes

transport
11722 genes

secretory pathway
4505 genes

vesicle-mediated transport
1525 genes

intracellular transport
2255 genes

Golgi vesicle transport
442 genes

ER to Golgi transport
190 genes

Annotations are given to the most specific (low) level.

True path rule:
Annotation at a term implies annotation to all its parent terms

Annotation is given with an Evidence Code:
 EXP: inferred from Experiment
 IDA: inferred by direct assay
 TAS: traceable author statement
 ISS: inferred by sequence similarity
 IEA: electronic annotation

# Functional Annotation DDBB
## Gene Ontology (GO terms)

- AmiGO provides a web interface to search and browse the ontology and annotation data

  **http://amigo.geneontology.org/cgi-bin/amigo/go.cgi**

- QuickGO (EBI) provides also a web interface

  **http://www.ebi.ac.uk/ego**

# Functional Annotation DDBB
## GO Slim

- ***GO slims*** are cut-down versions of the GO ontologies *containing a **subset*** of the terms in the whole GO. They give a broad overview of the ontology content without the detail of the specific fine grained terms

http://www.geneontology.org/GO.slims.shtml

# Functional Annotation DDBB
## Kyoto Encyclopedia of Genes and Genomes (KEGG)

**1. Metabolism**

**1.1 Carbohydrate Metabolism**
Glycolysis / Gluconeogenesis
Citrate cycle (TCA cycle)
Pentose phosphate pathway
Pentose and glucuronate interconversions
Fructose and mannose metabolism
Galactose metabolism
Ascorbate and aldarate metabolism
Starch and sucrose metabolism
Amino sugar and nucleotide sugar metabolism
Pyruvate metabolism
Glyoxylate and dicarboxylate metabolism
Propanoate metabolism
Butanoate metabolism
C5-Branched dibasic acid metabolism
Inositol phosphate metabolism

**1.2 Energy Metabolism**
Oxidative phosphorylation
Photosynthesis
Photosynthesis - antenna proteins
Carbon fixation in photosynthetic organisms
Reductive carboxylate cycle in photosynthetic bacteria
Methane metabolism
Nitrogen metabolism
Sulfur metabolism

**1.3 Lipid Metabolism**
Fatty acid biosynthesis
Fatty acid elongation in mitochondria
Fatty acid metabolism
Synthesis and degradation of ketone bodies
Steroid biosynthesis
Primary bile acid biosynthesis
Secondary bile acid biosynthesis
Steroid hormone biosynthesis
Glycerolipid metabolism
Glycerophospholipid metabolism
Ether lipid metabolism
Sphingolipid metabolism
Arachidonic acid metabolism
Linoleic acid metabolism
alpha-Linolenic acid metabolism
Biosynthesis of unsaturated fatty acids

**1.4 Nucleotide Metabolism**
Purine metabolism
Pyrimidine metabolism

**1.5 Amino Acid Metabolism**
Alanine, aspartate and glutamate metabolism
Glycine, serine and threonine metabolism
Cysteine and methionine metabolism
Valine, leucine and isoleucine degradation
Valine, leucine and isoleucine biosynthesis
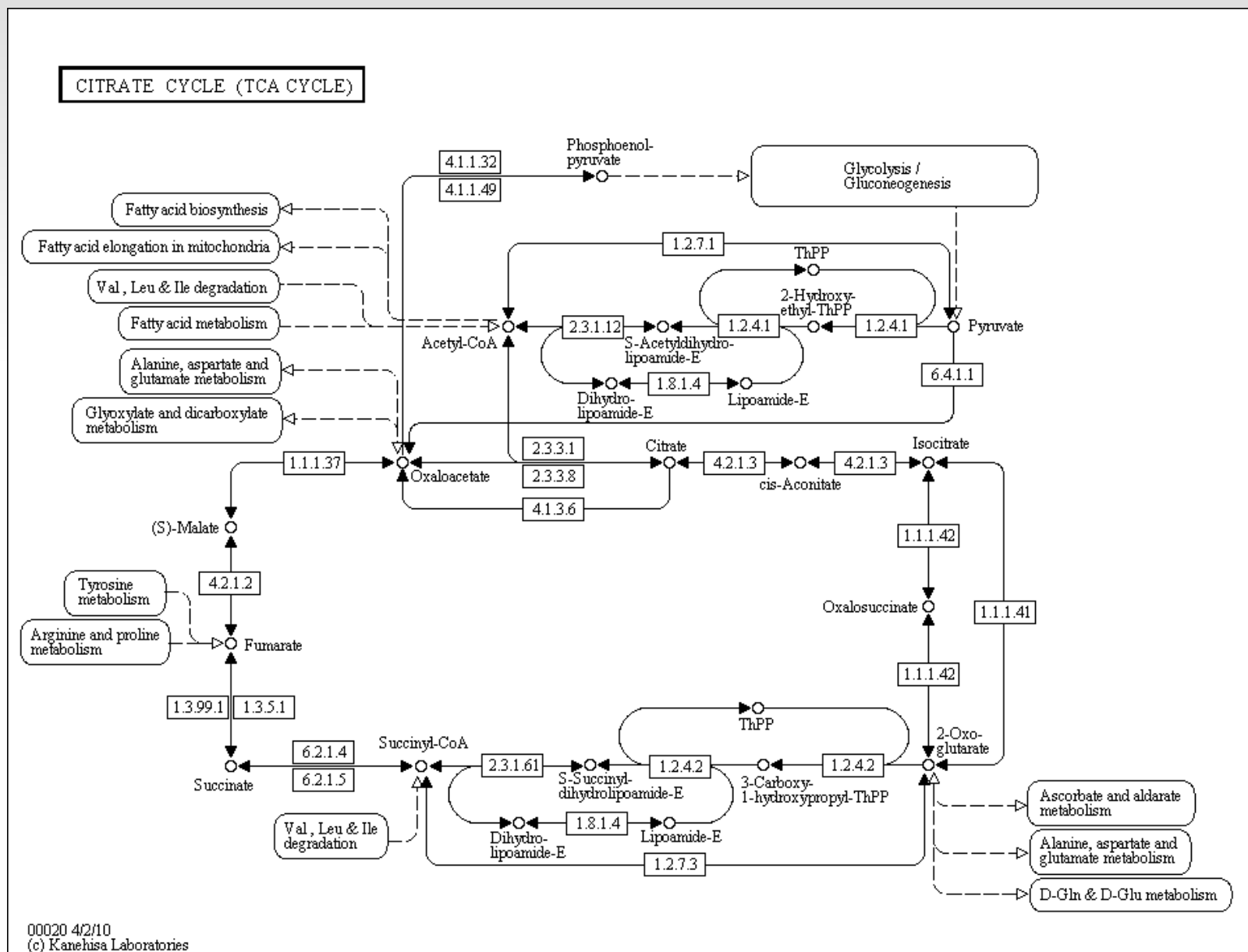Lysine biosynthesis

## KEGG pathways

**KEGG Databases** as of 2011/3/24

| | | |
|---|---|---|
| KEGG PATHWAY | Pathway maps, reference (total) | 389 (134,354) |
| KEGG BRITE | Functional hierarchies, reference (total) | 98 (37,769) |
| KEGG MODULE | KEGG modules, reference (total) | 0 (79,118) |
| KEGG DISEASE | Human diseases | 375 |
| KEGG DRUG | Drugs | 9,332 |
| KEGG EDRUG | Crude drugs and other natural products | 834 |
| KEGG ORTHOLOGY | KEGG Orthology (KO) groups | 14,360 |
| KEGG GENOME | KEGG Organisms | 1,558 |
| KEGG GENES | Genes in high-quality genomes (140 eukaryotes, 1205 bacteria, 97 archaea) | 6,359,583 |

## http://www.genome.jp/kegg/
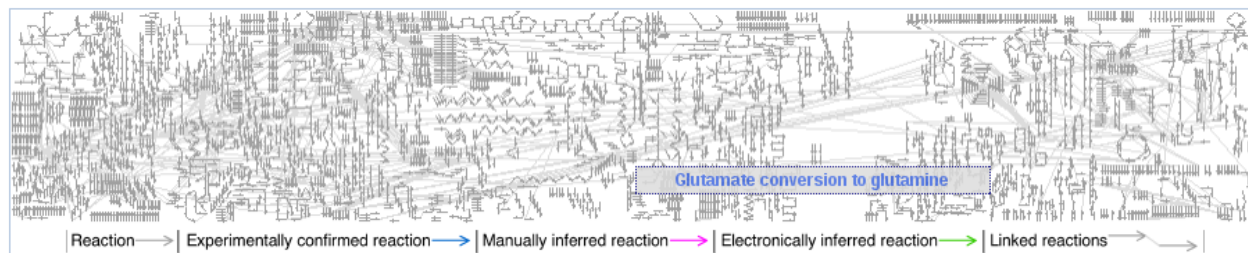
# Functional Annotation DDBB
## KEGG

# Functional Annotation DDBB
## Reactome

- It is a free, online, open-source, curated pathway database encompassing many areas of human biology. Information is authored by expert biological researchers
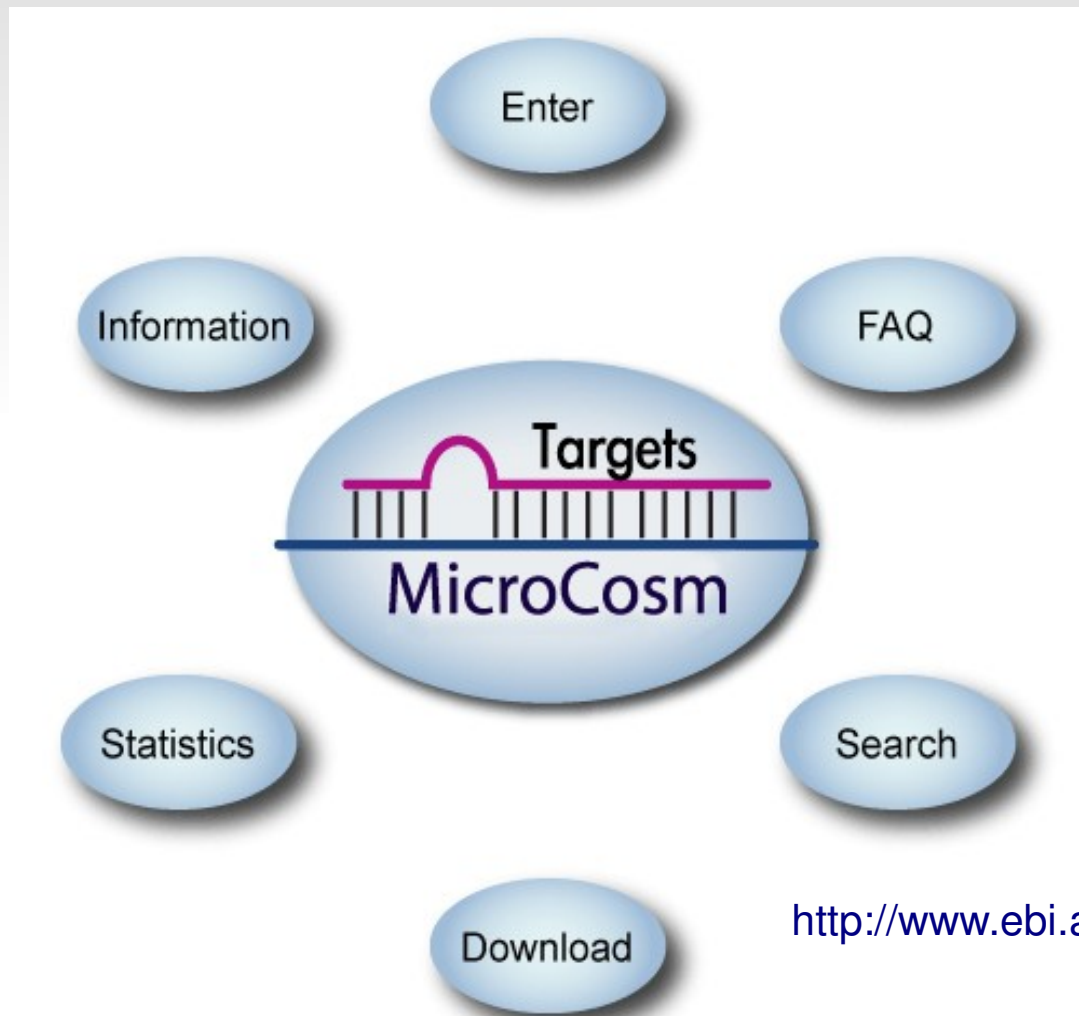
http://www.reactome.org/



| Apoptosis | Axon guidance | Biological oxidations | Botulinum neurotoxicity |
|---|---|---|---|
| Cell junction organization | Cell Cycle Checkpoints | Cell Cycle, Mitotic | DNA Repair |
| DNA Replication | Diabetes pathways | Electron Transport Chain | Gap junction trafficking and regulation |
| Gene Expression | Hemostasis | HIV Infection | Influenza Infection |
| Integration of energy metabolism | Integrin cell surface interactions | Metabolism of lipids and lipoproteins | Membrane Trafficking |
| Metabolism of amino acids and derivatives | Metabolism of carbohydrates | Metabolism of nitric oxide | Metabolism of nucleotides |
| Metabolism of polyamines | Metabolism of porphyrins | Metabolism of proteins | Metabolism of RNA |
| Metabolism of vitamins and cofactors | Muscle contraction | mRNA Processing | Myogenesis |
| Pyruvate metabolism and Citric Acid (TCA) cycle | Regulation of beta-cell development | Regulatory RNA pathways | Signaling by BMP |
| Signaling by EGFR | Signaling by FGFR | Signaling by GPCR | Signaling by PDGF |
| Signaling in Immune system | Signaling by Insulin receptor | Signalling by NGF | Signaling by Notch |
| Opioid Signalling | Signaling by Rho GTPases | Signaling by TGF beta | Signaling by VEGF |
| Signaling by Wnt | Synaptic Transmission | Telomere Maintenance | Transcription |
| Transmembrane transport of small molecules | | | |

# Functional Annotation DDBB
## MicroRNA



- Involved in gene regulation
- Last versions has 15172 entries (Release 16, Sept 2010)
- The *target database* contains computationally predicted targets for microRNAs across many species

http://www.mirbase.org/

http://www.ebi.ac.uk/enright-srv/microcosm/htdocs/targets/v5/

# Functional Annotation DDBB
## Jaspar TFBS

- The JASPAR database contains a curated, non-redundant set of profiles, derived from published collections of experimentally defined transcription factor binding sites for eukaryotes

- The prime difference to similar resources (TRANSFAC, etc) consist of the open data acess, non-redundancy and quality

http://jaspar.genereg.net/

# Functional Annotation DDBB
## ORegAnno

- It's an open database for the curation of known regulatory elements from scientific literature (*TFBS*)

- Annotation is collected from users worldwide for various biological assays



REGULATORY HAPLOTYPE: 7 entries.
REGULATORY REGION: 37520 entries.
TRANSCRIPTION FACTOR BINDING SITE: 14608 entries.
REGULATORY POLYMORPHISM: 175 entries.

http://www.oreganno.org/oregano/

# Functional Annotation DDBB
## String

Database for known and predicted protein-protein interactions (direct and indirect associations)
Cover four sources of annotations: Genomic association (prokaryotes), high-throughput experiments (e.g. y2h), conserved co-expression, previous knowledge (text-mining).

A combined score is calculated for every association based on benchmarks of the different types of associations against a common reference set.

http://string.embl.de/

Combined score (all sources)

**Your Input:**

P53   Cellular tumor antigen p53 (Tumor suppressor p53) (Phosphoprotein p53) (Antigen NY-CO-13) (393 aa)
(Homo sapiens)

**Predicted Functional Partners:**

| | | Score |
|---|---|---|
| HDM2 | Ubiquitin-protein ligase E3 Mdm2 (EC 6.3.2.-) (p53-binding protein Mdm2) (Oncoprotein Mdm2) (Double minute 2 protein) (Hdm2) (491 aa) | 0.999 |
| ATM | Serine-protein kinase ATM (EC 2.7.11.1) (Ataxia telangiectasia mutated) (A-T, mutated) (3056 aa) | 0.999 |
| p300 | Histone acetyltransferase p300 (EC 2.3.1.48) (E1A-associated protein p300) (2414 aa) | 0.999 |
| CHEK1 | Serine/threonine-protein kinase Chk1 (EC 2.7.11.1) (476 aa) | 0.999 |
| CDKN2A | Cyclin-dependent kinase inhibitor 2A, isoform 4 (p14ARF) (p19ARF) (173 aa) | 0.999 |
| MDM4 | Mdm4 protein (p53-binding protein Mdm4) (Mdm2-like p53-binding protein) (Mdmx protein) (Double minute 4 protein) (490 aa) | 0.999 |
| RNF53 | Breast cancer type 1 susceptibility protein (RING finger protein 53) (1863 aa) | 0.999 |
| HYRC | DNA-dependent protein kinase catalytic subunit (EC 2.7.11.1) (DNA-PK catalytic subunit) (DNA-PKcs) (DNPK1) (p460) (4127 aa) | 0.999 |
| TP53BP1 | Tumor suppressor p53-binding protein 1 (p53-binding protein 1) (p53BP1) (53BP1) (1972 aa) | 0.999 |
| PLK3 | Serine/threonine-protein kinase PLK3 (EC 2.7.11.21) (Polo-like kinase 3) (PLK-3) (Cytokine-inducible serine/threonine-protein kinase) (FGF-inducible kinase) (Proliferation-related kinase) (646 aa) | 0.999 |

# Protein Databases
## *UniProt*, protein sequence and information

**UniProtKB/Swiss-Prot** contains
531473 sequence entries

http://www.uniprot.org/

# Protein Databases
## *InterPro*, protein annotation database

- A ***centralized database*** of protein families, domains, repeats and sites in which identifiable features found in known proteins can be applied to new protein sequences

**http://www.ebi.ac.uk/interpro/**

**Member database information**

| Signature Database | Version | Signatures* | Integrated Signatures** |
|---|---|---|---|
| GENE3D | 3.3.0 | 2386 | 1377 |
| HAMAP | 021210 | 1675 | 1429 |
| PANTHER | 7.0 | 80933 | 1777 |
| PIRSF | 2.74 | 3248 | 2791 |
| PRINTS | 41.1 | 2050 | 2009 |
| PROSITE patterns | 20.66 | 1308 | 1292 |
| PROSITE profiles | 20.66 | 901 | 877 |
| Pfam | 24.0 | 11912 | 11465 |
| PfamB | 24.0 | 142303 | 0 |
| ProDom | 2006.1 | 1894 | 1008 |
| SMART | 6.1 | 895 | 882 |
| SUPERFAMILY | 1.73 | 1774 | 1154 |
| TIGRFAMs | 9.0 | 3808 | 3796 |

**Contents of InterPro 31.0 (Feb 2011)**

| | |
|---|---|
| Active site | 97 |
| Binding site | 65 |
| Conserved site | 615 |
| Domain | 5936 |
| Family | 14194 |
| PTM | 16 |
| Repeat | 262 |

# Protein Databases
## *IntAct*, protein-protein interaction database

**IntAct** provides a freely available, open source database system and analysis tools for protein interaction data

All interactions are derived from literature curation or direct user submissions and are freely available

http://www.ebi.ac.uk/intact/main.xhtml

# Protein Databases
## Protein Data Bank (*PDB*)

The **PDB** archive contains information about experimentally-determined structures of proteins, nucleic acids, and complex assemblies

http://www.rcsb.org/pdb/home/home.do

# Variation Databases
## *dbSNP*, the repository of all the *SNPs*

**dbSNP**
**Short Genetic Variations**

NCBI

PubMed  Nucleotide  Protein  Genome  Structure  PopSet  Taxonomy  OMIM  Books  SNP

Search for SNP on NCBI Reference Assembly

Search Entrez  SNP ▼  for [          ]  Go

Have a question about dbSNP? Try searching the SNP FAQ Archive!
[          ]  Go

ANNOUNCEMENT
8/15/2011: dbSNP Build 134 Release

Please see the build announcement for more details
(http://www.ncbi.nlm.nih.gov/projects/SNP/docs/build134.txt)

**GENERAL**
RSS Feed
Contact Us
Site Map
dbSNP Homepage
Announcements
dbSNP Summary
FTP Download
**HUMAN VARIATION**
**SNP SUBMISSION**
**DOCUMENTATION**
**SEARCH**
**RELATED SITES**

### Search by IDs on All Assemblies
Note: rs# and ss# must be prefixed with "rs" or "ss", respectively (i.e. ss25)

ID: [          ]  Reference cluster ID(rs#) ▼
Search  Reset

### Submission Information
- By Submitter
- New Submitted Batches
- Method
- Population
- Publication

http://www.ncbi.nlm.nih.gov/projects/SNP/

**BUILD STATISTICS:**

| Organism | dbSNP Build | Genome Build | Number of Submissions (ss#'s) | Number of RefSNP Clusters (rs#'s) (# validated) | Number of (rs#'s) in gene | Number of (ss#'s) with genotype | Number of (ss#'s) with frequency |
|---|---|---|---|---|---|---|---|
| Homo sapiens | 134 | 37.2 | 179,506,198 | 41,365,915 (6,961,883) | 16,880,922 | 73,208,602 | 35,627,484 |
| Mus musculus | 132 | 37.1 | 26,991,031 | 15,522,011 (6,439,098) | 6,696,618 | | |
| Pongo abelii | 132 | | 10,225,850 | 10,065,309 (0) | | | |
| Pongo pygmaeus | 127 | | 7,854,083 | 7,854,081 (0) | | | |
| Rattus norvegicus | 130 | 4.1 | 6,472,989 | 119,436 (1,605) | 1,024,738 | | |
| Gallus gallus | 131 | 2.1 | 11,318,097 | 3,504,588 (3,269,983) | 1,452,147 | | 50 |
| Glycine max | 127 | | 6,378,350 | 6,352,034 (234) | | | |
| Phoenix dactylifera | 133 | | 3,518,029 | 3,429,753 (0) | | | |
| Zea mays | 128 | | 4,555,638 | 4,350,627 (80) | | | |
| Oryza sativa | 128 | 4.1 | 5,872,306 | 5,359,569 (21,773) | 1,897,895 | | |
| Ovis aries | 128 | | 2,899,286 | 2,899,215 (66) | | | 91 |
| Bos taurus | 131 | 4.1 | 4,931,454 | 2,210,557 (13,881) | 677,906 | | 446 |
| Canis familiaris | 131 | 2.1 | 3,527,071 | 3,258,962 (214,713) | 982,946 | | 17 |

# Variation Databases
## *HapMap*, human Haplotype Map

To develop a public resource that will help researchers find genes associated with human disease and response to pharmaceuticals

http://hapmap.ncbi.nlm.nih.gov/

# Variation Databases
## Mutations: OMIM, COSMIC, Mitelman, …

http://www.ncbi.nlm.nih.gov/omim

http://www.sanger.ac.uk/genetics/CGP/cosmic/

# Genome DDBB and Browsers
## *Ensembl*, the most used and reliable



http://www.ensembl.org/index.html

The **Ensembl** project produces genome databases for vertebrates and other eukaryotic species, and makes this information freely available online

# Genome DDBB and Browsers
## *UCSC*

# Ejercicios
## Introducción

- Estamos interesados en estudiar un gen llamado *BCL2*.

- Utilizando las bases de datos explicadas anteriormente vamos a buscar información acerca de:

    - Secuencia génica

    - Información funcional y reguladora

    - Variaciones conocidas

    - Proteinas

# Ejercicios

## Información sobre la secuencia génica de *BCL2*

- Desde la página de Ensembl (http://www.ensembl.org/) intenta responder a las siguientes preguntas:

    - Indica la localización del gen y en qué cadena se encuentra

    - ¿Para cuántos tránscritos codifica?

    - ¿Y para cuantas proteínas?

    - Encuentra su secuencia de DNA

    - Indica el número de exones que contiene el gen

# Ejercicios
## Información funcional y reguladora de *BCL2*

- Ayúdate de GO (http://www.geneontology.org/) y encuentra los términos "biological process" y "cellular components" (GO terms) relacionados con el gen.

- Utiliza MicroCosm (http://www.ebi.ac.uk/enright-srv/microcosm/) para determinar si existe algún microRNA que regule a este gen.

- Dirígete a OregAnno (http://www.oreganno.org/). ¿Existe algún factor de transcripción conocido que regule a este gen?

- Utiliza KEGG (http://www.genome.jp/kegg/) y Reactome ( http://www.reactome.org/) para determinar en qué rutas (pathways) podemos encontrar este gen involucrado.

# Ejercicios
## Información sobre variaciones en *BCL2*

- Consulta en dbSNP (http://www.ncbi.nlm.nih.gov/snp/) el número de SNPs localizados en la secuencia de nuestro gen.

- En OMIM (http://www.ncbi.nlm.nih.gov/omim/) podemos encontrar información médica relacionada con mutaciones en genes. ¿Existe alguna enfermedad relacionada con nuestro gen?

- UniProtKB (http://www.uniprot.org/) contiene anotaciones sobre las consecencias observadas al mutar determinados aminoácidos en las secuencias proteicas. Observa qué tipo de consecuencias pueden tener estas mutaciones sobre la proteína.

# Ejercicios
## Información sobre las <u>proteínas</u> producidas por *BCL2*

- Obtén la siguiente inforamción de la proteína BCL2:

  - Secuencia (*UniProt*)

  - ¿Qué dominos proteicos funcionales tienen (*interpro*)?

  - Estructura 3D (*PDB*)