



"I hope you've got a lot of disk space, Ted. I think I accidentally just faxed you the entire Internet."

## Introducción a las Bases de datos biológicas

Febrero 2010

---

---

---

---

---

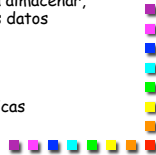
---

---

---

### Porque las bases de datos ?

- Crecimiento exponencial de los datos biológicos
- Datos (secuencias, 3D estructuras, análisis gel 2D, MS análisis....) no son publicados en revistas, pero si en bases de datos
- Son usadas en investigación biológica, como lo eran la revistas científicas !
- Biólogos dependen de los computadores para almacenar, organizar, buscar, manipular, y recuperar los datos
- Libre Acceso es clave
- Base de todas las herramientas bioinformáticas



---

---

---

---

---

---

---

---

### Que es una base de datos ?

- Una colección
  - estructurada
  - De fácil búsqueda (indexada) → tabla de contenido
  - Actualizada periódicamente (release) → Nuevas dediciones
  - Referencias cruzadas ([hipervínculos](#)) → vínculos con otras DB
- Incluye la herramientas (software) para acceso, actualización, inserción, borrado.... en la DB
- Manejo del almacenamiento de los datos: Texto plano (flat files), tablas vinculadas (bases de datos relacionales)

---

---

---

---

---

---

---

---

## DB: Texto plano « flat file »

Base de datos de estudiantes:  
(texto plano, 3 entradas)

```
Código: 183023
Nombre: Julián
Apellido: Pulecio
Cursos : 19003-01, 21001-01
Email: jpul@ibun.unal.edu.co
//
//
Código: 183024
Nombre: Sonia
Apellido : Cuartas
Cursos : 19003-01, 17001-01
Email: soniacol@hotmail.com
//
//
Código: 183025
Nombre: Jaime
Apellido : Moreno
Cursos : 19003-01
Email: pm186111@ibun.unal.edu.co
//
```

- Fácil de manejar: todas las entradas se pueden ver al tiempo!

---

---

---

---

---

---

---

---

## Bases de datos « relacionales »

Curso	Nom_Curso
19003-01	Bioinformática
17001-01	Bioquímica Avanzada
21001-01	Análisis Molecular

Alumno	Código
Gutiérrez	183023
Cuartas	183024
Moreno	182425

Curso	Código
19003-03	183023
19003-03	183024
19003-03	182425
17001-01	183024
21001-01	183023

Fácil: manejo y selección de la salida

---

---

---

---

---

---

---

---

## Algunas estadísticas

- Más de 2000 bases de datos diferentes
- Generalmente accesibles a través de WEB (
  - Google: <http://www.google.com/>
  - Biohunt: <http://www.expasy.org/BioHunt/>
  - Amos' links: [www.expasy.ch/alinks.html](http://www.expasy.ch/alinks.html)
- Tamaño variable: <100Kb a >150Gb
  - DNA: > 180 Gb
  - Proteínas: 7 Gb
  - Estructuras 3D : 15 Gb
  - Otras: Pequeñas

---

---

---

---

---

---

---

---

### Clases de DB para las ciencias de la vida

- Secuencias (DNA, proteínas) -> DB primarias
- Genómicas
- Dominios/familias proteicos -> DB secundarias
- Mutación/polimorfismo
- Proteómica (2D gel, MS)
- 3D estructura -> DB de estructuras
- Metabolismo
- Bibliografía
- Otras

---

---

---

---

---

---

---

---

### Contenido Ideal mínimo de una DB de « secuencias »

- Número de acceso (AC)
- Secuencias !!
- Referencias
- Datos taxonómicos
- ANOTACIONES/CURADURIA
- Palabras claves
- Referencias cruzadas
- Documentación

---

---

---

---

---

---

---

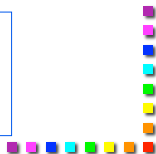
---

### Ejemplo...Formato de entrada para secuencias

\_entrada SWISS-PROT, en formato fasta :

```
>sp|P01588|EPO_HUMAN_ERYTHROPOIETIN_PRECURSOR - Homo sapiens (Human).
MGVHECPAWLWLLLSLGLPLGLPVLGAPPLICDSRVLEKLEAKKAE
NITTCACAEHCNSLNENITVPTKVNFKYAMKREVEQQQAVVWQGLALSEA
VLRGQALLVNSQWPELQLRVDKAVSGLRSLTLLRALGAQKEAISPFD
AASAAFLRTTADTFRKLRFVYSNFLRGLKLYTGEACRTGDR
```

Formatos empleados (texto plano):  
fasta  
GCG  
NBRF/PIR  
MSF....  
Formatos estandarizados ?



---

---

---

---

---

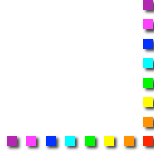
---

---

---

## Bases de datos anotadas - curadas?

- contienen la secuencia, comentarios, referencias de la literatura, notas sobre experimentos
- Derivadas de la integración de las herramientas de cómputo y conocimiento biológico
  - por ejemplo, genes conocidos y predichos
- Registros añadidos solo después de verificar su precisión y las anotaciones
- Ejemplo :  
SWISS-PROT, OMIM, RefSeq, LocusLink




---

---

---

---

---

---

---

---

---

---

## Ejemplo: Base de datos de secuencia

### SWISS-PROT Flat file

```

ID     EPO_HUMAN          STANDARD;          PRT;   193 AA.
AC     F01588; Q9UHA0; Q9UEZ5; Q9UDE0;
DT     21-JUL-1986 (Rel. 01, Created)
DT     21-JUL-1986 (Rel. 01, Last sequence update)
DT     20-AUG-2001 (Rel. 40, Last annotation update)
DE     Erythropoietin precursor.
OS     EPO.
OX     Homo sapiens (Human).
OC     Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC     Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.
OX     NCBI_TaxID=9606;
RN     [1]
RE     SEQUENCE FROM N.A.
RX     MEDLINE=85137899; PubMed=3838366;
RA     Jacobs K., Shoemaker C., Rudersdorf R., Neill S.D., Kaufman R.J.,
RA     Wilson A., Szebera J., Jones S.S., Hewick R., Frisoe R.F.,
RA     Kawakita M., Shimizu T., Miyake T.;
RT     "Isolation and characterization of genomic and cDNA clones of human
RT     erythropoietin."
RL     Nature 313:806-810(1985).
    
```

#### Taxonomía

#### Referencia




---

---

---

---

---

---

---

---

---

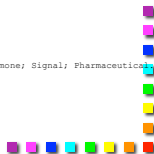
---

## Ejemplo: Base de datos de secuencia

### SWISS-PROT Flat file

```

---
CC     -!- FUNCTION: ERYTHROPOIETIN IS THE PRINCIPAL HORMONE INVOLVED IN THE
CC     REGULATION OF ERYTHROCYTE DIFFERENTIATION AND THE MAINTENANCE OF A
CC     PHYSIOLOGICAL LEVEL OF CIRCULATING ERYTHROCYTE MASS.
Anotaciones CC     -!- SUBCELLULAR LOCATION: SECRETED.
CC     -!- TISSUE SPECIFICITY: PRODUCED BY KIDNEY OR LIVER OF ADULT MAMMALS
CC     AND BY LIVER OF FETAL OR NEONATAL MAMMALS.
CC     -!- PHARMACEUTICAL: Available under the names Epoen (Amgen) and
CC     Procrit (Ortho Biotech).
---
DR     EMBL: X02158; CAA26095.1; -.
DR     EMBL: X02157; CAA26094.1; -.
Referencias DR     EMBL: M11319; AAA2480.1; -.
Cruzadas   DR     EMBL: AF033356; AAC78791.1; -.
DR     EMBL: AF202308; AAF23132.1; -.
DR     EMBL: AF202306; AAF23132.1; JOINED.
---
Palabras claves KW     Erythrocyte maturation; Glycoprotein; Hormone; Signal; Pharmacological
    
```




---

---

---

---

---

---

---

---

---

---

## Ejemplo: Base de datos de secuencia

```

FT SIGNAL 1 27
FT CHAIN 28 193 ERYTHROPOIETIN.
FT PROPEP 190 193 MAY BE REMOVED IN PROCESSED PROTEIN.
FT DISULFID 34 188
FT DISULFID 56 60
FT CARBOHYD 51 51 N-LINKED (GLCNAC...).
FT CARBOHYD 65 65 N-LINKED (GLCNAC...).
FT CARBOHYD 110 110 N-LINKED (GLCNAC...).
FT CARBOHYD 153 153 O-LINKED (GALNAC...).
FT VARIANT 131 132 SL -> RF (IN AN HEPATOCELLULAR
FT CARCINOMA).
FT VARIANT 149 149 /FTID=VAR_009870,
FT E -> Q (IN AN HEPATOCELLULAR CARCINOMA).
FT /FTID=VAR_009871.
FT CONFLICT 40 40 E -> Q (IN REF. 1; CAA26095).
FT CONFLICT 85 85 Q -> QQ (IN REF. 5).
FT CONFLICT 140 140 G -> R (IN REF. 1; CAA26095).
**
** ***** INTERNAL SECTION *****
**CL 7q22;
SQ
SEQUENCE 193 AA; 21306 MW; C91F0E4C26A52033 CRC64;
MOVHSCFQAWL WLLSLLSLLP LGLVPLGAP SLICDSRVLE WLLKSDAE NITTCQAEHC
SLNENIVFD TRVNFYAKR MEVQQQVEV WQLALLSEA VLKQDALLVN SQQWPLQLL
HVMKAVSLR SLTLLLRALQ AQKEAISFPD AASAAFLRTI TADTFRKLFK VYSNPLRGLK
KLYTGEACTG GER
//

```

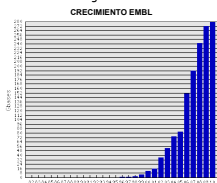
## Bases de Datos de Secuencias de nucleótidos

### EMBL/GenBank/DDJB

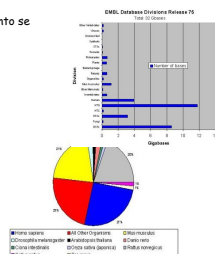
- Contienen principalmente la misma información en el plazo de 2-3 días (pocas diferencias en el formato y sintaxis)
  - Contienen las secuencias (genes únicos, ESTs, genomas completos, etc.) derivadas de:
    - Los proyectos genoma y centro de secuenciación
    - Científicos Individuales
    - Oficinas de patentes (es decir oficina de patentes europea, EPO)
- Datos No-confidenciales son intercambiados diariamente  
 Actualmente: 20 secuencias x10<sup>6</sup>, sobre 30 x10<sup>9</sup> bp;  
 Estadísticas: <http://www3.ebi.ac.uk/Services/DBStats/>  
 secuencias > 73'000 especies diferentes;

## Impresionante incremento en secuencias de nucleótidos

- Datos en EMBL... su gran incremento se debió al surgimiento de PCR...



Febrero 2010 . 279,125,380,402 bases en 181,171,193 registros.  
 Fuente: <http://www3.ebi.ac.uk/Services/DBStats/>

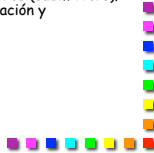


1980: 80 genes totalmente secuenciados!

## EMBL/GenBank/DDBJ



- Longitud Heterogenea de secuencias: genomas fragmentos...
- Tamaños de Secuencia:
  - max 300'000 pb / entrada (! secuencias genómicas, solapadas)
  - min 10 bp / entrada
- Archivo: nada sale!!! -> muy redundante!
- lleno de errores: en las secuencias, en anotaciones, en atribución de CDS...
- Ninguna consistencia en las anotaciones; la mayoría de las anotaciones son hechas por los investigadores (submitters); heterogeneidad en la calidad y en la realización y actualización de la información




---

---

---

---

---

---

---

---

---

---

---

---

## EMBL/GenBank/DDJB



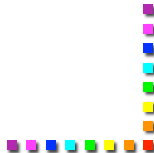
- Se puede encontrar información inesperada en estas bd:

```

FT source      1.124
FT             /db_xref="taxon:4097"
FT             /organelle="plastid:chloroplast"
FT             /organism="Nicotiana tabacum"
FT             /isolate="Cuban cahibo cigar, gift from President Fidel Castro"
    
```

```

O:
FT source      1.17084
FT             /chromosome="complete mitochondrial genome"
FT             /db_xref="taxon:9267"
FT             /organelle="mitochondrion"
FT             /organism="Didelphis virginiana"
FT             /dev_stage="adult"
FT             /isolate="fresh road killed individual"
FT             /tissue_type="liver"
    
```




---

---

---

---

---

---

---

---

---

---

---

---

## Entrada de EMBL : ejemplo

```

ID L1000 standard_Olea_F01_716 bp
AC B44811.578372
XX
CC B44811.1
DT 10-APR-1993 (Rel 1) (Created)
UT 10-JUN-1993 (Rel 5) (Last updated: Version 6)
OR Listeria innocens, sod gene for superoxide dismutase
OX and gene, superoxide dismutase
OC
OS Listeria innocens
OC Neisseria-Fishelson, Bacillus-Clottidina group
OC Bacillus-Staphylococcus group: Listeria
XX
CC
XX
CC (1)
CC BC021096.1(10197)
CC GenBank: BC021096.1
CC Cloning of a superoxide dismutase gene from Listeria innocens by
CC functional complementation in Escherichia coli and characterization of the
CC gene product.
CC Mol. Gen. Genet. 231: 91-92 (1992).
XX
CC L1000
CC Entry 1:
XX
CC
CC Submitted (21-APR-1993) to the EMBL/GenBank/CCDS databases.
CC L1000: 10000001.100000000. Submitted Sequence. Recombinant in
CC Accession: B7834882.2
CC
CC S055-P807: P24763_S00X_L151F
XX
CC Key Location/Qualifiers
XX
CC source          1..716
CC             /db_xref="taxon:5381"
CC             /organism="Listeria innocens"
CC             /protein_coding="yes"
CC             /catalytic_activity="ATC 13119"
CC             /EC="1.11.1.1"
CC             /termination_codon="TAG"
CC             /translation_start="1"
CC             /translation_end="716"
CC             /CDR
CC             /db_xref="GenBank:BC021096.1"
CC             /db_xref="EMBL:B7834882.2"
CC             /db_xref="S055-P807:P24763"
CC             /catalytic_activity="ATC 13119"
    
```

### Palabras claves

### Taxonomía

### Referencias

### Referencias cruzadas




---

---

---

---

---

---

---

---

---

---

---

---

## Entrada de EMBL : (Cont.)

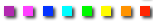
```

HL Submitted (21-APR-1992) to the EMBL GenBank/DBJ databases.
L1 seq1: Institut für Biologie, Universität Würzburg, Plattenstr. 4a
L1 Rahland, 97080 Würzburg, FRG
L1
L1 SP188-PR07: F28763, SGCN_L351Y.
L1
L1 Key Location/Qualifiers
L1
L1 FT source 1..756
L1 FT /db_xref="taxon:1628"
L1 FT /organism="Escherichia faecalis"
L1 FT /strain="ATCC 19113"
L1 RBS 96..100
L1 FT /gene="tnpA"
L1 FT terminator 725..748
L1 FT /gene="tnpA"
L1 FT CDS 157..794
L1 FT /db_xref="SP188-PR07: F28763"
L1 FT /db_xref="F01043: tnpA1"
L1 FT /gene="tnpA"
L1 FT /protein_id="tnpA1"
L1 FT /product="transposase"
L1 FT /protein_id="tnpA1"
L1 FT /translation="MTYSYDLYDTTALRHHSPWMEHTFTSNNIVYLYLKAIVLQ
L1 SLLTILPSSDDEINTEEDYVWVYDGRDGIWYTFVLSVFVQVDFQALPQVSLK
L1 RYDFRIFDQKDLKALASDFQDAF: FRSQSLITLTLKQDFLSDGATFTD
L1 DYVHYATLFLQSRKSTIDTFDFYIHWGRNDRDPAK"
L1
L1 CC
L1 SO
Sequence 756 bp. 247 A. 136 C. 151 G. 222 T. 0 other.
Sequence 756 bp. 247 A. 136 C. 151 G. 222 T. 0 other:
60
GGATTTC TGGATGAT GATGATGG AAGAGATC GACTTTTG GGCTACAA 60
TACCAAAAT TACCTATAC ATATGATCT TGGAGTGG ATTGATAT GGAACAAG 120
AAATCACT AAGACAGA AAGAGATC TATGTATG AATGATAG AGATGTCA 180
GGGACAGG AACTTGGG AAAAGTGG GAGGATAG TGTATATC AATAGACT 240
GCGACAAA TGTGTTGG AGTATGCG GAGATGTA GCAATCTA 300
TCTGTTTG GCTTGGCC AATTGGTAT GATGTACG CTTATCACT AAGAGTGG 360
AAAGAGAG AATGGAGG AATGATAA TGGAGAAA AATGATAG GAGGATGG 420
CGCTATTG ATTGATGG CTTGAGTA GGTGGTGA ATGTAGCA AAGATGCT 480
TGCCTATG ACCAAGTC TCACTTGG GAGGATAG CTCGATTC TGTGTGAT 540
ATTTGGAG AGTATATA ATTTAGCT CAGAGCTG GCGTATAG CATTGGAG 600
TTTGGATG AATTAAGT GATGACCA AATAGACT TGGAGGCG AATATAAT 660
TGGAGGCT CAATAAGT GATCTATA TTTCTA

```

Anotación

Secuencia

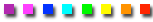


## Entrada de GenBank : ejemplo

```

□ 1: A3416340. Escherichia coli . [g:16304811]
LOCUS EC0416340 1470 bp DNA linear BCT 07-JUL-2002
DEFINITION Escherichia coli partial tnpA gene for transposase and blaCTX-M-1
gene for CTX-M-1 beta-lactamase.
ACCESSION A3416340
VERSION A3416340.1 GI:16304811
KEYWORDS blaCTX-M-1 gene; CTX-M-1 beta-lactamase; tnpA gene; Transposase.
SOURCE Escherichia coli
ORGANISM Escherichia coli
Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
Enterobacteriaceae; Escherichia.
REFERENCE 1
AUTHORS Shinoh,M., Chu,V.T., Lambert,T., Doney,J.L., Hermann,J.L.,
Ould-Hocine,Z., Verdet,C., Deisze,F., Philippon,A. and Arlet,G.
TITLE Diversity of CTX-M beta-lactamases and their promoter regions from
Enterobacteriaceae isolated in three Parisian hospitals
JOURNAL FEBS Microbiol. Lett. 209 (2), 161-168 (2002)
MEDLINE 2300747
PubMed 12007950
REFERENCE 2 (bases 1 to 1470)
AUTHORS Arlet,G.J.
TITLE Direct Submission
JOURNAL Submitted (08-OCT-2001) Arlet G.J., Bacteriology, UFR
Saint-Antoine, 27, rue de Chaligny, Paris 75571 cedex 12, FRANCE

```



## Entrada de GenBank : ejemplo

```

FEATURES             Location/Qualifiers
     source            1..1470
                     /organism="Escherichia coli"
                     /mol_type="genomic DNA"
                     /strain="tnc1"
                     /db_xref="taxon:562"
                     /country="France:Paris"
     repeat_region     1..318
                     /insertion_seq="ISKpi1"
     gene              1..307
                     /gene="tnpA"
     CDS                1..307
                     /gene="tnpA"
                     /function="transposition"
                     /codon_start=2
                     /transl_table=1
                     /product="transposase"
                     /protein_id="tnc09165.1"
                     /db_xref="GI:16304811"
                     /db_xref="SF16348.1:Q93328"
                     /translation="MDVVYKSNVYKRTYKAVVSHLLRFFVHANEVVFQNNMLYVLYL
1STZKGLQVGVVYV...
BASE COUNT           429 a 100 c 376 g 365 t
ORIGIN
1 tggaaagtgt gtagaatgct aaactatcat caaagaagcc aaatagaca tggaggtagg
61 tcaatctctg cttaagcaat tttggggaga tgaagctgtg ttcataatga tgaagctctc
121 atacaactca ttttggagt tgaagctgct tttcctgagc tctcagaatc aatgagtaga

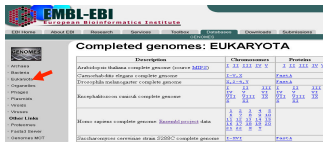
```



## Divisiones EMBL

- EMBL ha sido dividida en sub-bases de datos para permitir el fácil mantenimiento y búsqueda de datos
  - fun, hum, inv, mam, org, phg, pln, pro, rod, syn, unc, vrl, vrt
  - est, gss, htg, htc, sts, patent

<http://www.ebi.ac.uk/genomes/>



Organismo	Chromosomas	Proyecto
Arabis thaliana - complete genome (Arabid)	5, 10, 11, 2, 3, 4, 5	1998-2000
Canis familiaris - complete genome	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100	2000-2001
Drosophila melanogaster - complete genome	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100	2000-2001
Macaca mulatta - complete genome (RhesusMacaca)	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100	2001-2002
Mus musculus - complete genome (MusMus)	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100	2001-2002
Danio rerio - complete genome (Danio)	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100	2001-2002

## DB genómicas

- Contienen información sobre genes, su localización (mapping), nomenclatura y links a DB de secuencia; **usualmente no contiene secuencias**;
- Existen para los organismos más importantes en investigación de las ciencias de la vida;
- Ejemplos: MIM, GDB (human), MGD (mouse), FlyBase (Drosophila), SGD (yeast), MaizeDB (maize), SubtiList (B.subtilis), etc.;
- Formato: generalmente relacional (Oracle, SyBase or AceDb).

## MIM

- OMIM™: Online Mendelian Inheritance in Man
- Catalogo de genes humanos y desordenes genéticos
- Contiene un sumario de literatura, fotos, y referencias. Contiene numerosos links a información de secuencia y artículos.



## MIM:

\*133170 ERYTHROPOIETIN; EPO

Alternative titles; symbols

EP

TABLE OF CONTENTS

TEXT  
REFERENCES  
SEE ALSO  
CONTRIBUTORS  
CREATION DATE  
EDIT HISTORY

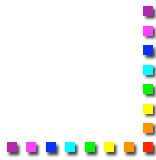
Database Links

Gene Map Locus: 7q21

Note: pressing the symbol will find the citations in MEDLINE whose text most closely matches the text of the preceding MIM paragraph, using the Entrez MEDLINE neighboring function.

TEXT

Human erythropoietin is an acidic glycoprotein hormone with molecular weight 34,000. As the prime regulator of red cell production, its major functions are to promote erythroid differentiation and to initiate hemoglobin synthesis. Sherwood and Shouval (1986) described a human renal carcinoma cell line that continuously produces erythropoietin. Eschbach et al. (1987) demonstrated the effectiveness of recombinant human erythropoietin in treating the anemia of end-stage renal disease. Lee-Rueng (1984) cloned human erythropoietin cDNA in *E. coli*. McDonald et al. (1986) and Shoemaker and Mitscock (1986) cloned the mouse gene and the latter workers showed that coding DNA and amino acid sequences are about 80% conserved between man and mouse. This is a much higher order of conservation than for various interferons, interleukin-2, and GM-CSF.



---

---

---

---

---

---

---

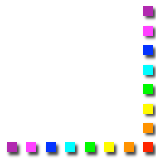
---

---

---

## Proyectos DB de Secuencia y genómica

- Ensembl
- TIGR



---

---

---

---

---

---

---

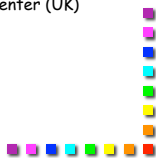
---

---

---

## Ensembl: automatic annotation of eukaryotic genomes

- Contiene todas las secuencias de DNA del genoma humano normalmente disponibles para el dominio público.
- Anotación automatizada: usando diferentes herramientas de software, identificando dentro de las secuencias de DNA:
  - Genes (conocidos o predichos)
  - Polimorfismos de un solo nucleótido (SNPs)
  - Repeticiones
  - Homologías
- Creado y mantenido por el EBI y el Sanger Center (UK)



---

---

---

---

---

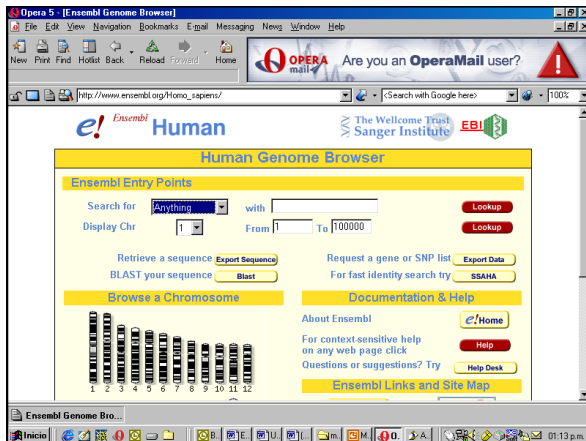
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

---

---

**Ensembl: [www.ensembl.org](http://www.ensembl.org)**

Con Ensembl usted puede ...

- Buscar secuencias de DNA en el genoma humano
- Visualizar los mapas cromosómicos
- Encontrar genes, SNPs y sitios comunes con el genoma de ratón
- Mirar proteínas y familias de proteínas

• Ensembl provee:

- Identificación del 90% de los genes humanos conocidos
- Predicción de 10,000 genes adicionales, todos con evidencia soportada

---

---

---

---

---

---

---

---

---

---

---

---

**BD de secuencias de proteína**

- **SWISS-PROT**: creada en 1986 (A.Bairoch)  
<http://www.expasy.org/sprot/>
- **TrEMBL**: creada en 1996; complementa a SWISS-PROT; derivada de las traducciones de CDS EMBL automatizadas (versión «proteómica» de EMBL)
- **PIR-PSD**: Protein Information Resources  
<http://pir.georgetown.edu/>

---

---

---

---

---

---

---

---

---

---

---

---

## BD de secuencias de proteína

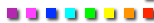
■ **PRF**: Protein Research Foundation (Japan): Peptide/Protein Sequence Database (PRF/SEQDB)

<http://www.prf.or.jp/en/index.html>

■ **GenPept**: producida a partir del análisis del release de GenBank correspondiente para regiones codificantes traducidas.

■ Muchas bases de datos de proteína especializada para familias específicas o grupos de proteínas.

■ Ejemplos: **YPD** (proteínas de levadura), **AMSDb** (peptidos antibacterianos), **GPCRDB** (receptores 7 TM), **IMGT** (sistema inmune) etc.




---

---

---

---

---

---

---

---

---

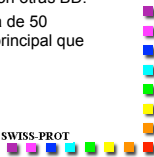
---

## Swissprot

<http://www.expasy.ch/sprot/>

**SWISS-PROT** es una base de datos de secuencias de proteína,

- Colaboración entre el SIB (CH) y EMBL/EBI, 1986
- Curada y Anotada (descripción de la función, estructura de los dominios, modificaciones post- traduccionales, variantes)
- Mínima redundancia y alta Integración con otras BD.
- Release semanal; disponible en cerca de 50 servidores a través del mundo, la fuente principal que es ExPASy




---

---

---

---

---

---

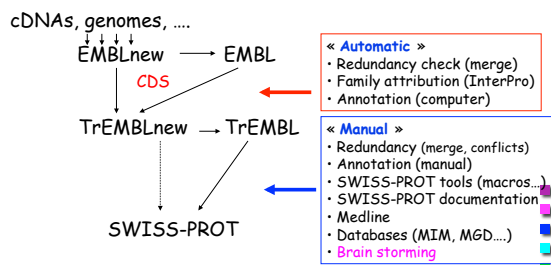
---

---

---

---

## Historia de una entrada en SWISS-PROT



Once in SWISS-PROT, the entry is no more in TrEMBL, but still in EMBL (archive)

CDS: proposed and submitted at EMBL by authors or by genome projects (can be experimentally proved or derived from gene prediction programs). TrEMBL does not take CDS into account. Also, genome projects use gene prediction programs: only take CDS already annotated in the EMBL entry.




---

---

---

---

---

---

---

---

---

---





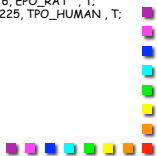
## Prosite (pattern): ejemplo

```

ID EPO_TPO; PATTERN.
AC P500817;
DT OCT-1993 (CREATED); NOV-1995 (DATA UPDATE); JUL-1998 (INFO UPDATE).
DE Enkephalin / Humoboposiaa ciguatera.
PA P-x(4)-C-D-x-R-[LIVM(2)-x-[KR]-x(14)-C.
NR /RELEASE=30,00000.
NR /TOTAL=14(14); /POSITIVE=14(14); /UNKNOWN=0(0); /FALSE_POS=0(0);
NR /FALSE_NEG=0; /PARTIAL=1;
CC /TAXO-RANGE=?EPP; /MAX-REPEAT=1;
CC /SITE=3 disulfide; /SITE=11 disulfide;
DR P48617, EPO_BOVIN , T; P33707, EPO_CANFA , T; P33708, EPO_FELCA , T;
DR P01588, EPO_HUMAN , T; P07865, EPO_MACFA , T; Q28513, EPO_MACMU , T;
DR P07321, EPO_MOUSE , T; P49157, EPO_PIG , T; P29676, EPO_RAT , T;
DR P33709, EPO_SHEEP , T; P42705, TPO_CANFA , T; P40225, TPO_HUMAN , T;
DR P40226, TPO_MOUSE , T; P49745, TPO_RAT , T;
DR P42706, TPO_PIG , P;
DO PD000644;
//
    
```

Diagnostic performance

List of matches



## Prosite (profile): ejemplo

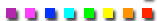
```

PROSITE: P55097
ID: BTL_MATRIX.
AC: P55097.
DT: DEC-1999 (CREATED); DEC-1999 (DATA UPDATE); DEC-1999 (INFO UPDATE).
DE: BTL_MATRIX.
MA: /GENERAL_SPEC ALPHABET=ABCDEFGHIJKLMNPQRSTVWYZ; LENGTH=67.
MA: /DISJOINT DEFINITION PROTECT NONE; NO GAP.
MA: /NORMALIZATION MODE=L; FUNCTION LINEAR; RI: 9751; S2: 0.0068002; TEXT=L; LOGE.
MA: /CUT-OFF LEVEL=3; SCORE=363; N: SCORE=5; MODE=1; TEXT=Y.
MA: /CUT-OFF LEVEL=1; SCORE=357; N: SCORE=5; MODE=1; TEXT=Y.
MA: /DEFAULT D: -20; I: -20; E: -50; EI: -50; MI: -105; MD: -105; IM: -105; DM: -105; MM: -1; WD: -2.
MA: /Z RND: B1=305; B0=305.
MA: /M SV: C: M=-30 28 14 9 10 20 14 19 15 17 14 8 19 14 25 0 0 9 32 17 12.
MA: /M SV: D: M=-16 8 1 23 10 14 11 1 23 0 27 23 21 11 0 8 2 26 38 10 7.
MA: /M SV: E: M=-2 23 8 28 24 1 24 25 16 20 7 6 20 25 23 30 30 4 24 23 9 24.
MA: /M SV: F: M=-12 19 10 19 10 27 20 19 16 8 23 11 17 11 10 10 0 0 0 0 4 10.
MA: /M SV: G: M=-11 30 22 23 24 15 32 23 25 29 30 17 26 27 25 22 24 9 16 17 3 24.
MA: /M SV: H: M=-11 16 13 10 20 20 23 14 4 2 2 20 19 4 7 4 2 8 25 9 9.
MA: /M SV: I: M=-1 25 3 29 25 2 29 26 17 22 10 7 23 25 23 22 11 3 24 27 10 25.
MA: /M SV: J: M=-8 7 28 8 7 25 6 7 27 0 23 17 8 13 0 3 3 4 23 27 17 3.
MA: /Z I: I: MD=0; IM=0; DM=15; WD=15.
MA: /M SV: K: M=-8 27 8 3 27 22 7 30 8 26 19 10 14 8 9 2 9 24 28 21 6.
MA: /M SV: L: M=-7 4 23 4 7 23 13 2 21 10 18 9 3 12 7 9 4 4 26 25 12 6.
MA: /M SV: M: M=-8 4 21 8 15 21 7 7 1 10 5 3 14 0 1 2 2 4 26 9 1.
MA: /M SV: N: M=-15 28 22 24 26 31 31 11 18 18 26 9 2 27 27 21 20 9 14 4 13 26.
MA: /M SV: O: M=-13 9 24 20 3 11 21 7 17 7 18 4 4 8 2 9 9 16 20 1 2.
MA: /M SV: P: M=-15 28 22 24 26 31 31 11 18 18 26 9 2 27 27 21 20 9 14 4 13 26.
MA: /M SV: Q: M=-15 22 2 1 20 16 6 26 8 21 5 15 15 6 1 2 11 26 32 7 0.
MA: /M SV: R: M=-12 5 29 5 25 18 8 26 34 24 9 1 14 8 24 8 2 17 20 10 5.
MA: /M SV: S: M=-4 12 16 16 16 16 16 16 16 16 16 16 16 16 16 16 16 16 16 16 16 16.
MA: /M SV: T: M=-7 26 19 31 28 7 32 24 27 23 14 11 22 25 23 23 10 28 19 3 26.
MA: /M SV: U: M=-10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10.
MA: /M SV: V: M=-8 0 12 8 18 4 16 15 10 18 12 2 14 8 13 18 11 5 32 19 8.
    
```

•Tabla de puntaje específica de la posición de los residuos.

•Cuales son conservados o degenerados

•Posiciones que toleran inserciones



## ProDom

Consiste en una compilación automatizada del alineamiento de dominios homólogos

## PRINTS

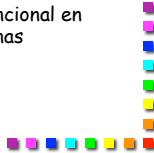
- Compendio de fingerprints de motivos proteicos
- La mayoría de familias de proteína poseen varios motivos conservados
- Fingerprint:** grupo de motivos (simple o compuesto, por ejemplo multidominios) = característica de los miembros de una familia
- Un miembro de la familia exhibe todos los elementos del fingerprint, mientras que los miembros de la subfamilia pueden poseer solamente una parte



## BD de Familia/dominio de Proteínas:BD compuestas

Ejemplo: **InterPro**

- Reune a PROSITE, PRINTS, Pfam, ProDom y SMART en un recurso integrado de familias de proteínas protein families, dominios y sitios funcionales;
- Single set of «documents» linked to the various methods;
- Utilizado para mejorar la anotación funcional en SWISS-PROT (clasificación de proteínas desconocidas)



---

---

---

---

---

---

---

---

## InterPro: ejemplo

**IPRO00323**

Name: Erythropoietin/hrombopoietin

Type: Family

Abstract

Erythropoietin, a plasma glycoprotein, is the primary physiological mediator of erythropoiesis [1]. It is involved in the regulation of the level of peripheral erythrocytes by stimulating the differentiation of erythroid progenitor cells. Found in the spleen and bone marrow, into mature erythrocytes [2]. It is primarily produced in adult kidneys and fetal liver, acting by attachment to specific binding sites on erythroid progenitor cells, stimulating their differentiation [3]. Severe kidney dysfunction causes reduction in the plasma levels of erythropoietin, resulting in chronic anaemia - injection of purified erythropoietin into the blood stream can help to relieve this type of anaemia. Levels of erythropoietin in plasma fluctuate with varying oxygen tension of the blood, but undergoes and prostaglandins also modulate the levels to some extent [3]. Erythropoietin glycoprotein sequences are well conserved, a consequence of which is that the hormones are cross-reactive among mammals, i.e. that from one species, say human, can stimulate erythropoiesis in other species, say mouse or rat [4].

Thrombopoietin (TPO), a glycoprotein, is the mammalian hormone which functions as a megakaryocytic lineage specific growth and differentiation factor of facting the proliferation and maturation from their committed progenitor cells acting at a late stage of megakaryocyte development. It acts as a circulating regulator of platelet numbers.

Examples:1  
F23708  
F23709  
F49740  
view matches for the examples

Publications

1. Shoemaker C.B., Mitschke L.D. 849-858 (1986)
2. Takeuchi M., Takasaki S., Miyazaki H., Kato T., Hoshi S., Kuchibe N., Kobata A. J. Biol. Chem. 263: 3457-3463 (1988)
3. Liu F.K., Liu C.H., Liu P.H., Browne J.K., Egan J.C., Smalling R., Fox G.M., Chen K.K., Castro M., Suggs S. Gene 44: 201-209 (1994)
4. Nagao M., Suga H., Okano M., Masuda S., Nishita H., Dura K., Sasaki R. Nucleotide sequence of rat erythropoietin. 1171-99-302 (1992)

Children

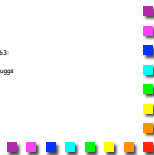
IPRO00303

Signatures

PROSITE PS00867 EPO\_TPO  
Pfam PF00758 EPO\_TPO

Matches

Table: [Genes/Prot](#)



---

---

---

---

---

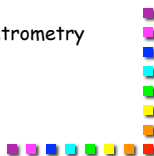
---

---

---

## BD de Proteomica

- Contain informations obtained by 2D-PAGE: master images of the gels and description of identified proteins
- Examples: SWISS-2DPAGE, ECO2DBASE, Maize-2DPAGE, Sub2D, Cyano2DBase, etc.
- Format: composed of image and text files
- Most 2D-PAGE databases are "federated" and use SWISS-PROT as a master index
- There is currently no protein Mass Spectrometry (MS) database (not for long...)



---

---

---

---

---

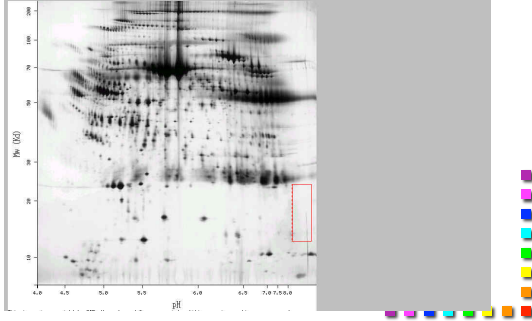
---

---

---

This protein does not exist in the current release of SWISS-2DPAGE.

### EPO\_HUMAN (human plasma)



---

---

---

---

---

---

---

---

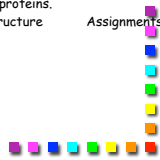
---

---

### BD 7: estructura 3D

- Contain the spatial coordinates of macromolecules whose 3D structure has been obtained by X-ray or NMR studies
- Proteins represent more than 90% of available structures (others are DNA, RNA, sugars, virus, complex protein/DNA...)
- PDB (Protein Data Bank), SCOP (structural classification of proteins (according to the secondary structures)), BMRB (BioMagResBank; RMN results)
- DSSP: Database of Secondary Structure Assignments.  
HSSP: Homology-derived secondary structure of proteins.  
FSSP: Fold Classification based on Structure-Structure
- Future: Homology-derived 3D structure db.

Assignments.



---

---

---

---

---

---

---

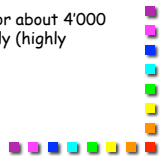
---

---

---

### PDB: Protein Data Bank

- Managed by Research Collaboratory for Structural Bioinformatics (RCSB) (USA).
- Contains macromolecular structure data on proteins, nucleic acids, protein-nucleic acid complexes, and viruses.
- Specialized programs allow the visualization of the corresponding 3D structure.
- Currently there are ~16'000 structure data for about 4'000 different molecules, but far less protein family (highly redundant)!



---

---

---

---

---

---

---

---

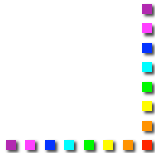
---

---



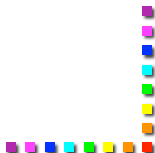
## PDB: example

```
HEADER LYASE(OXO-ACID) 01-OCT-91 1ZCA 1ZCA 2
COMPND CARBONIC ANHYDRASE I1 (CARBONATE DEHYDRATASE) (HCA I1) 1ZCA 3
COMPND 2 (E.C.4.1.1) MUTANT WITH VAL 121 REPLACED BY ALA (V121A) 1ZCA 4
SOURCE HUMAN (HOMO SAPIENS) RECOMBINANT PROTEIN 1ZCA 5
AUTHOR S.KNAIB,D.W.CHRISTIANSON 1ZCA 6
REVDAT 1 15-OCT-92 1ZCA 0 1ZCA 7
JRNL AUTH 5 KNAIB,T.L.CALDERONE,D.W.CHRISTIANSON,C.A.FIERKE 1ZCA 8
JRNL TITL ALTERING THE MOUTH OF A HYDROPHOBIC POCKET. 1ZCA 9
JRNL TITL 2 STRUCTURE AND KINETICS OF HUMAN CARBONIC ANHYDRASE 1ZCA 10
JRNL TITL 3 /115 MUTANTS AT RESIDUE VAL-121 1ZCA 11
JRNL REF 3 BIOLOGCHEM 7 266 17300 1991 1ZCA 12
JRNL REF1 ASTM 78CH43 US ISSN 0021-9258 071 1ZCA 13
REMARK 1 1ZCA 14
REMARK 2 1ZCA 15
REMARK 2 RESOLUTION: 2.4 ANGSTROMS. 1ZCA 16
REMARK 3 1ZCA 17
REMARK 3 REFINEMENT. 1ZCA 18
REMARK 3 PROGRAM PROLSQ 1ZCA 19
REMARK 3 AUTHORS HENDRIKSSON,KONNERT 1ZCA 20
REMARK 3 R VALUE 0.170 1ZCA 21
REMARK 3 RMSD BOND DISTANCES 0.021 ANGSTROMS 1ZCA 22
REMARK 3 RMSD BOND ANGLES 1.3 DEGREES 1ZCA 23
REMARK 4 1ZCA 24
REMARK 4 N-TERMINAL RESIDUES SER 2, HIS 3, HIS 4 AND C-TERMINAL 1ZCA 25
REMARK 4 RESIDUE LYS 260 WERE NOT LOCATED IN THE DENSITY MAPS AND, 1ZCA 26
REMARK 4 THEREFORE, NO COORDINATES ARE INCLUDED FOR THESE RESIDUES. 1ZCA 27
```



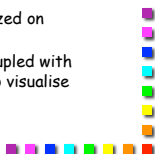
## PDB (cont.)

```
SHEET 3 SIDPHE 66 PHE 70 1 O ASN 67 N LEU 60 1ZCA 68
SHEET 4 SIDTR 88 TRP 97 1 O PHE 93 N VAL 68 1ZCA 69
SHEET 5 SIDALA 116 ASN 124 1 O HIS 119 N HIS 94 1ZCA 70
SHEET 6 SIDLEU 141 VAL 150 1 O LEU 144 N LEU 100 1ZCA 71
SHEET 7 SIDVAL 207 LEU 210 1 O ILE 210 N GLY 145 1ZCA 72
SHEET 8 SIDTR 195 GLY 196 1 O TRP 192 N VAL 211 1ZCA 73
SHEET 9 SIDLYS 237 ALA 198 1 O LYS 227 N THR 193 1ZCA 74
SHEET 10 SIDLYS 39 TRP 40 1 O LYS 39 N ALA 258 1ZCA 75
TURN 1 71 GLN 126 VAL 31 TYR 118 (CIS-PRO) 1ZCA 76
TURN 2 72 GLY 81 LEU 84 TYR 119 (TRANS) (GLY 82) 1ZCA 77
TURN 3 73 ALA 134 GLN 137 TYR 118 (CIS) 1ZCA 78
TURN 4 74 GLN 137 GLY 140 TYR 119 (ASP) 1ZCA 79
TURN 5 75 THR 200 LEU 203 TYR 119 (CIS-PRO) 1ZCA 80
TURN 6 76 GLY 221 GLY 226 TYR 119 (GLY) 1ZCA 81
CRYST1 42.700 42.700 73.000 90.00 104.60 90.00 P 21 2 1ZCA 82
ORF1K1 1.000000 0.000000 0.000000 0.00000 1ZCA 83
ORF2K2 0.000000 1.000000 0.000000 0.00000 1ZCA 84
ORF3K3 0.000000 0.000000 1.000000 0.00000 1ZCA 85
SCALE1 0.023419 0.020000 0.006800 0.00000 1ZCA 86
SCALE2 0.000000 0.023419 0.000000 0.00000 1ZCA 87
SCALE3 0.000000 0.000000 0.016456 0.00000 1ZCA 88
ATOM 1 N TRP 5 8397 -0.701 88.778 1.00 13.37 1ZCA 89
ATOM 2 CA TRP 5 1743 -1.668 11.585 1.00 13.42 1ZCA 90
ATOM 3 C TRP 5 6.786 -2.502 10.667 1.00 13.47 1ZCA 91
ATOM 4 O TRP 5 6.482 -2.085 9.667 1.00 13.57 1ZCA 92
ATOM 5 CB TRP 5 6.997 -0.917 12.645 1.00 13.34 1ZCA 93
ATOM 6 CG TRP 5 5.794 -0.209 12.215 1.00 13.40 1ZCA 94
ATOM 7 CD TRP 5 5.681 1.084 11.797 1.00 13.29 1ZCA 95
ATOM 8 CE TRP 5 4.417 -0.667 12.211 1.00 13.34 1ZCA 96
ATOM 9 NE1 TRP 5 4.388 1.418 11.515 1.00 13.30 1ZCA 97
ATOM 10 CE2 TRP 5 3.988 0.375 11.797 1.00 13.35 1ZCA 98
ATOM 11 CE3 TRP 5 3.817 -1.877 12.645 1.00 13.39 1ZCA 99
ATOM 12 CZ2 TRP 5 2.216 0.208 11.656 1.00 13.39 1ZCA 100
ATOM 13 CZ1 TRP 5 2.465 -0.241 12.504 1.00 13.33 1ZCA 101
ATOM 14 CH2 TRP 5 1.654 -1.001 12.009 1.00 13.34 1ZCA 102
```



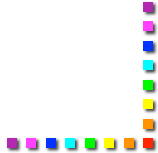
## BD 8: Metabólicas

- Contain informations that describe **enzymes**, biochemical reactions and **metabolic pathways**;
- ENZYME and BRENDA: **nomenclature databases** that store informations on enzyme names and reactions;
- Metabolic databases**: EcoCyc (specialized on Escherichia coli), KEGG, EMP/WIT;  
Usually these databases are tightly coupled with query software that allows the user to visualise reaction schemes.



## BD 9: Bibliográfica

- Bibliographic reference databases contain citations and abstract informations of published life science articles;
- Example: Medline
- Other more specialized databases also exist (example: Agricola).



---

---

---

---

---

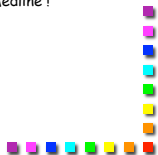
---

---

---

## Medline

- MEDLINE covers the fields of **medicine**, nursing, dentistry, veterinary medicine, the health care system, and the preclinical sciences
- more than 4,000 biomedical journals published in the United States and 70 other countries
- Contains over 10 million citations since 1966 until now
- Contains links to biological db and to some journals
- New records are added to PreMEDLINE daily!
  - Many papers not dealing with human are not in Medline !
  - Before 1970, keeps only the first 10 authors !
  - Not all journals have citations since 1966 !



---

---

---

---

---

---

---

---

## Medline/Pubmed

- the National Center for Biotechnology Information (NCBI)
- PubMed provides **access** to **bibliographic** information such as MEDLINE, PreMEDLINE, HealthSTAR, and to integrated molecular biology databases (composite db)
- PMID: 10923642 (PubMed ID), UI: 20378145 (Medline ID)



---

---

---

---

---

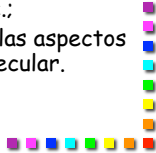
---

---

---

## Otras DB

- Hay muchas bases de datos que no pueden ser clasificadas;
- Ejemplos: ReBase (enzimas restricción), TRANSFAC (factores de transcripción), CarbBank, GlycoSuiteDB (Azúcares ligados), Proteína-proteína interacción db (DIR, ProNet, Interact), Proteasa db (MEROPS), patentes en biotecnología db, etc.;
- Como también otros aspectos de los aspectos de macromoléculas y biología molecular.



---

---

---

---

---

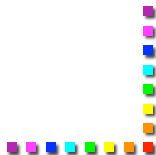
---

---

---

## Proliferación de DB

- Cual es la mejor DB para análisis de secuencia?
- Cual tiene la mejor calidad de datos ?
- Cual es la más completa ?
- Cual es la más actualizada ?
- Cual es la menos redundante ?
- Cual es la más indexada (permite búsquedas complejas) ?
- Cual es la que responde más rápido ?
- .....??????



---

---

---

---

---

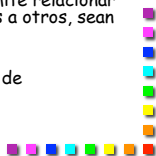
---

---

---

## Herramientas de búsqueda en DB

- **Sequence Retrieval System** (SRS, Europe) permite a cualquier db flat-file ser indexada con cualquier otra DB; permite formular consultas a través de un rango amplio de DBs usando una interfase simple, sin importar la estructura de los datos y los lenguajes de consulta...
- **Entrez** (USA): menos flexible que SRS pero explota el concepto de « neighbouring », el cual permite relacionar artículos de diferentes db para ser unidos a otros, sean o no referencias cruzadas específicas.
- **ATLAS**: Especifica para db de secuencias de macromoléculas (i.e. NRL-3D)
- ....



---

---

---

---

---

---

---

---

## SRS

Top Page Query Form Query Manager View Manager Databases Help

Search **SPTTR**

Do Query Reset Combine searches with AND Append wildcard \* to words.

Info AllText albumin

Info AllText

Info AllText

Info AllText

Include fields in output: ID, AccNumber, CreationDate, LastChangeDate, SubmissionDate, Description, GeneName

Entry List in chunks of: 30

Sequence Format: default

Use view: SequenceSimple

Retrieve set of: entry

Alternative Query Form: Separate multiple values by & (and), | (or), ! (and not)

---

---

---

---

---

---

---

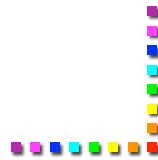
---

---

---

## Entrez-protein

- NCBI: <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein>
- Reune muchos recursos, incluyendo SWISS-PROT, PIR, PRF, PDB, y traducciones de regiones codificantes anotadas en el GenBank (« Genpept ») and RefSeq.
  - PRF: Protein Research Foundation (Japan): DB Peptido/Secuencias de proteínas (PRF/SEQDB)
  - PDB: Protein Data Bank (3D structure)
  - RefSeq: NCBI Reference Sequence project
  - PIR - International Protein Sequence Database
- Proteínas y DNA secuencias



---

---

---

---

---

---

---

---

---

---

Entrez Protein

NCBI Protein

Search Protein for albumin

Display Summary Save Text Details Add to Clipboard

Show 20 Items 1-20 of 644 Page 1 of 33 Select page: 1 2 3 4 5 6 7 8 9 10 >>

1 [CAA76847](#) PubMed, Related Sequences, Nucleotide, Taxonomy  
bovine serum albumin [Bos taurus]  
gi|3336842|emb|CAA76847.1|3336842

2 [AAH11965](#) Related Sequences, Nucleotide, Taxonomy  
D site of albumin promoter (albumin D-box) binding protein [Homo sapiens]  
gi|15080432|gb|AAH11965.1|AAH11965.1|5080432

3 [Q39837](#) PubMed, Related Sequences, Taxonomy  
ALBUMIN 1 PRECURSOR (PA1) [CONTAINS PA1A, LEGNUSLIN (PA1B)]  
gi|4916530|sp|Q39837|ALB1\_SOYBN|4916530

4 [P36233](#) PubMed, Related Sequences, Taxonomy  
SPARC PRECURSOR (SECRETED PROTEIN ACIDIC AND RICH IN CYSTEINE) (OSTEONECTIN) (ON) (BASEMENT MEMBRANE PROTEIN) BM-40  
gi|13959711|sp|P36233|SPARC\_RABIT|13959711

---

---

---

---

---

---

---

---

---

---