

Tema 4. Bases de datos secundarias

Actualizado en: 13/02/2023

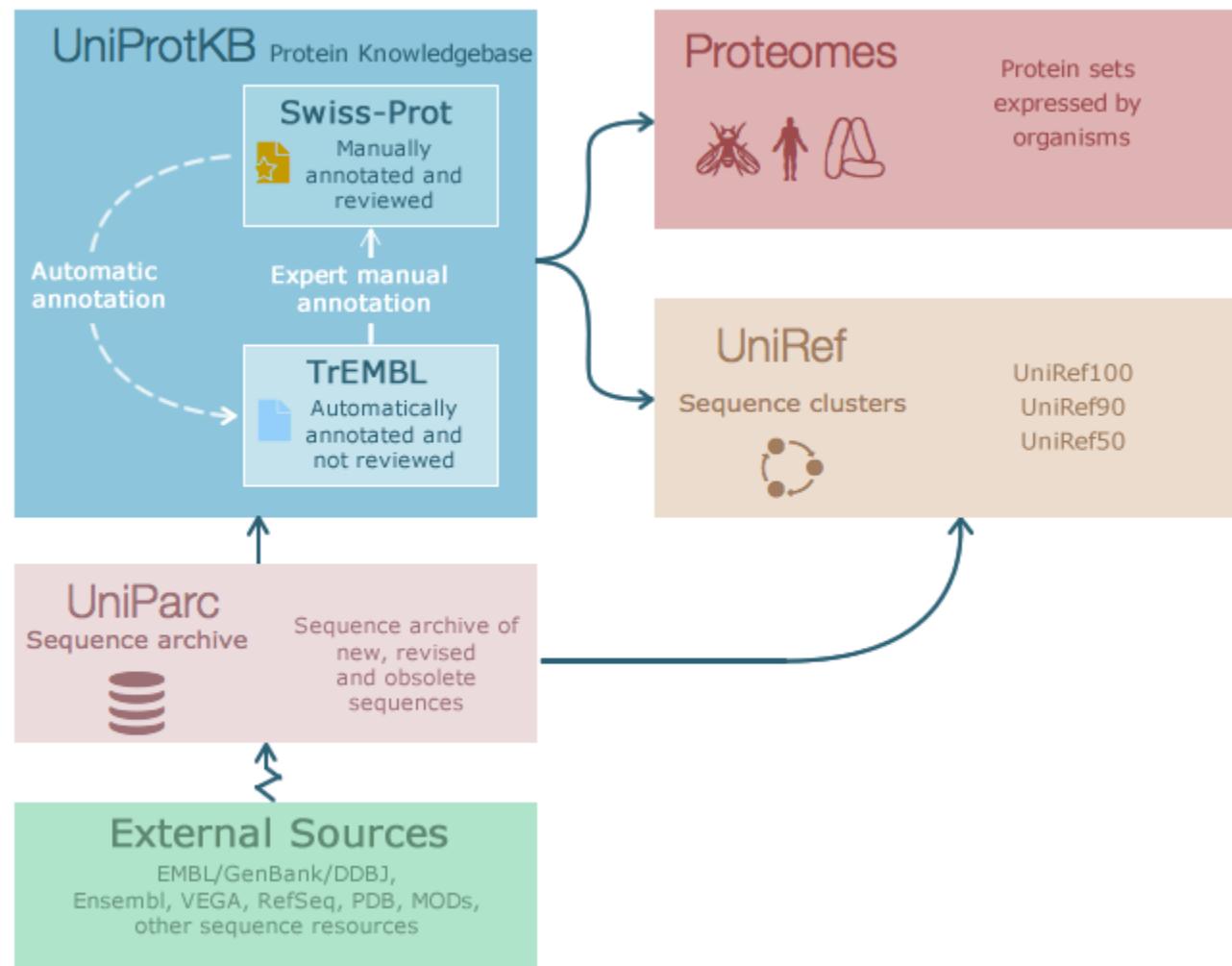
Objetivos

- UniProt
- Interpro
- ENSEMBL
- OrthoDB

UniProt

“The Universal Protein Resource (UniProt) is a comprehensive resource for protein sequence and annotation data. The UniProt databases are the UniProt Knowledgebase (UniProtKB), the UniProt Reference Clusters (UniRef), and the UniProt Archive (UniParc). The UniProt consortium and host institutions EMBL-EBI, SIB and PIR are committed to the long-term preservation of the UniProt databases.”

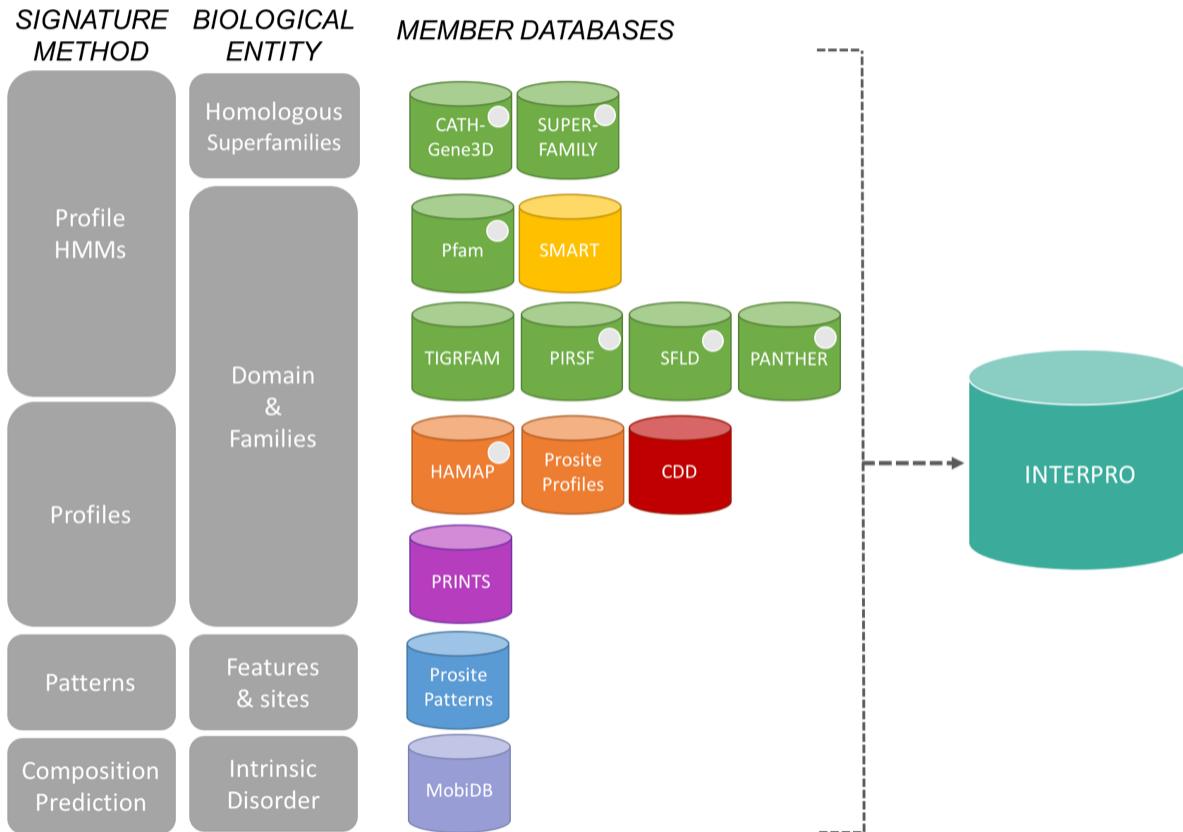
UniProt



UniProt. Cuestionario.

- ¿En qué sección de UniProt buscarías información sobre la función de una proteína?
- ¿Qué sección es menos *redundante*?
- ¿Por qué coexisten TrEMBL y Swiss-Prot?
- ¿De dónde salen los datos primarios?

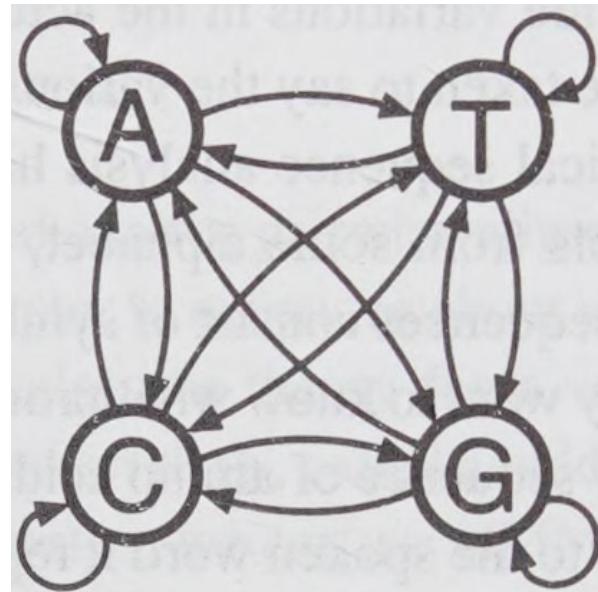
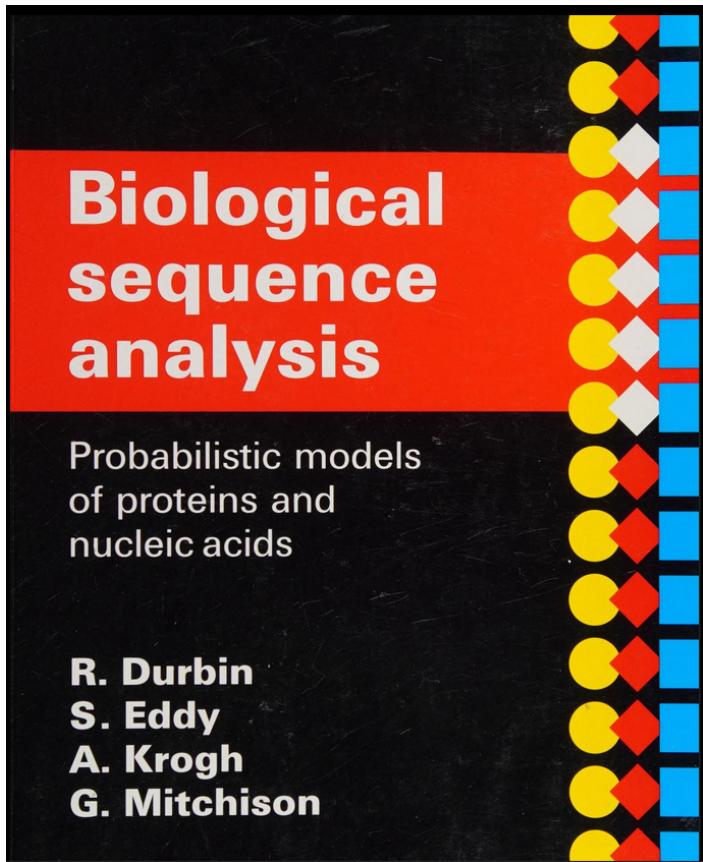
Interpro



Interpro. Tipos de registros

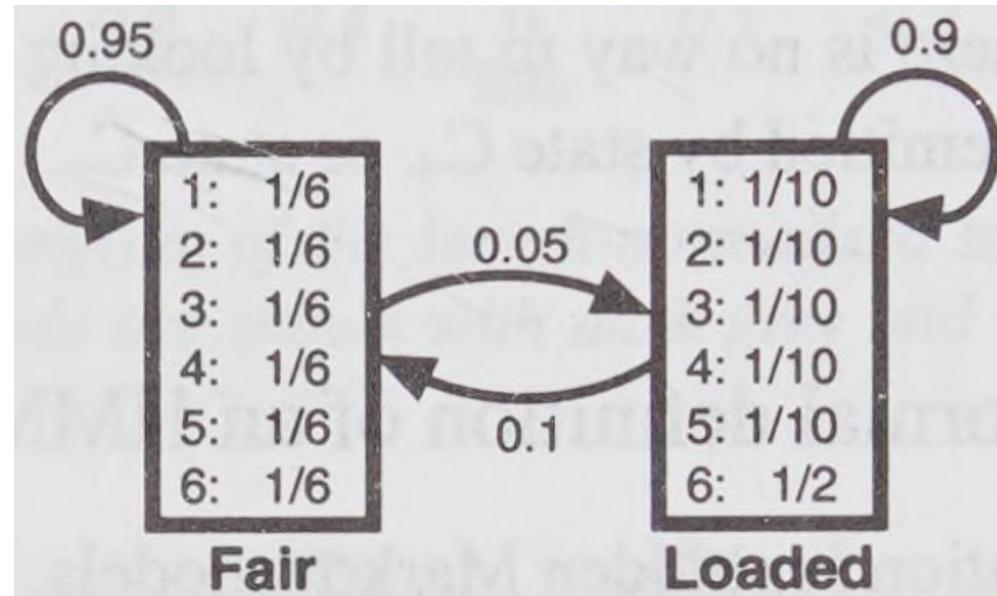
- **Familia:** Proteínas homólogas con secuencias, estructuras y funciones relacionadas.
- **Dominio:** Unidad funcional y estructural de las proteínas.
- **Sitio:** Motivos cortos: sitios de unión, sitios activos...
- **Repetición:** Secuencia de <50 aa. repetida muchas veces en una proteína.
- **Superfamilia:** Proteínas homólogas, aunque las secuencias no se parezcan.
- **No integrados:** Registros de las bases de datos participantes, no incluídos en Interpro.

Modelo de Markov



$$a_{st} = P(x_i = t | x_{i-1} = s)$$

Modelo de Markov Oculto (HMM)



Ejemplo del casino con dos dados, uno justo y otro *cargado*, con probabilidades de *emisión* diferentes. No sabemos en cada momento con qué dado se está jugando.

Modelo de Markov Oculto (HMM)

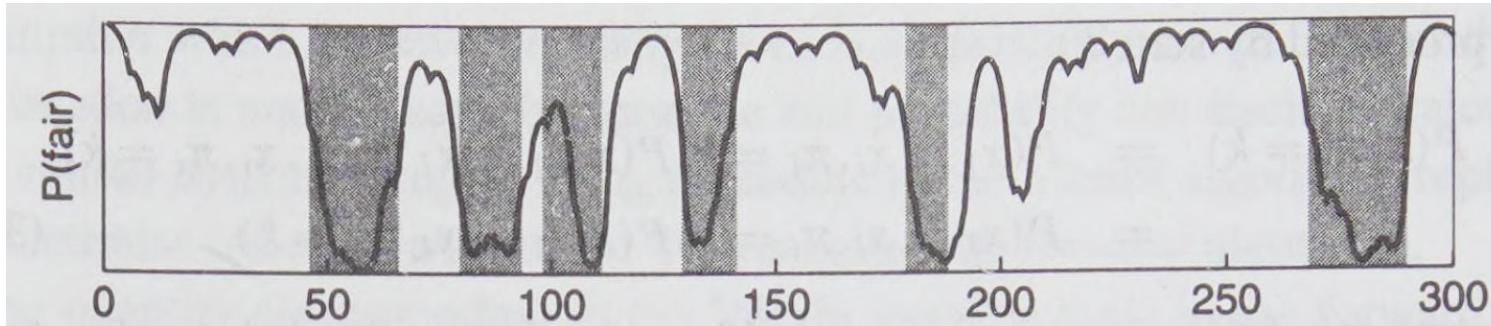


Figure 3.6 The posterior probability of being in the state corresponding to the fair die in the casino example. The x axis shows the number of the roll. The shaded areas show when the roll was generated by the loaded die.

A partir de los datos visibles, podemos ajustar los parámetros del modelo de Markov y predecir los estados ocultos a lo largo de la cadena de eventos.

Modelo de Markov Ocultos (HMM)

Start with a multiple sequence alignment



Insertions / deletions can be modelled



Occupancy and amino acid frequency at each position in the alignment are encoded



Profile created

seq1 A C G - L D
seq2 S C G -- E
seq3 N C G g F D
seq4 T C G - W Q

deletion

1 2 3 - 4 5

insertion

W

F

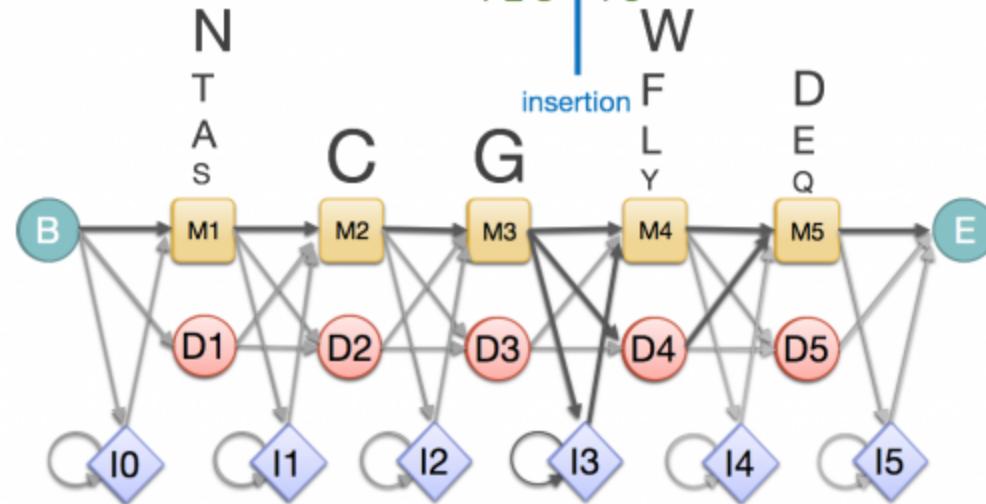
L

Y

D

E

Q



Secuencias genómicas



Predicción automática de genes



Base de datos MySQL



Análisis y visualización

ENSEMBL

- Servicio del EMBL-EBI.
- Originalmente para genoma humano.
- Incluye genomas de 311 especies.
- La base de datos y sus herramientas estan disponibles.

ENSEMBL. Acceso a los datos

- Consulta y visualización.
- Opciones de descarga de resultados de búsqueda.
- Descarga de tablas por ftp.
- Perl API.
- REST API.

ENSEMBL

Location-based displays

- Whole genome
- Chromosome summary
- Region overview
- Region in detail**
- Comparative Genomics
 - Synteny
 - Alignments (image)
 - Alignments (text)
 - Region Comparison
- Genetic Variation
 - Variant table
 - Resequencing
 - Linkage Data
 - Markers
- Other genome browsers
 - UCSC
 - NCBI
 - Ensembl GRCh37

Configure this page

Custom tracks

Export data

Share this page

Bookmark this page

Human (GRCh38.p13) ▾

Location: 17:63,992,802-64,038,237

Chromosome 17: 63,992,802-64,038,237

Assembly exceptions

Chr. 17

Assembly exceptions

Region in detail

Chromosome bands

Contigs

Genes (Comprehensive set from GENCODE 36)

Regulatory Build

Gene Legend

Regulation Legend

Location: 17:63992802-64038237

Go

Gene: Go

Chromosome bands

45.44 kb

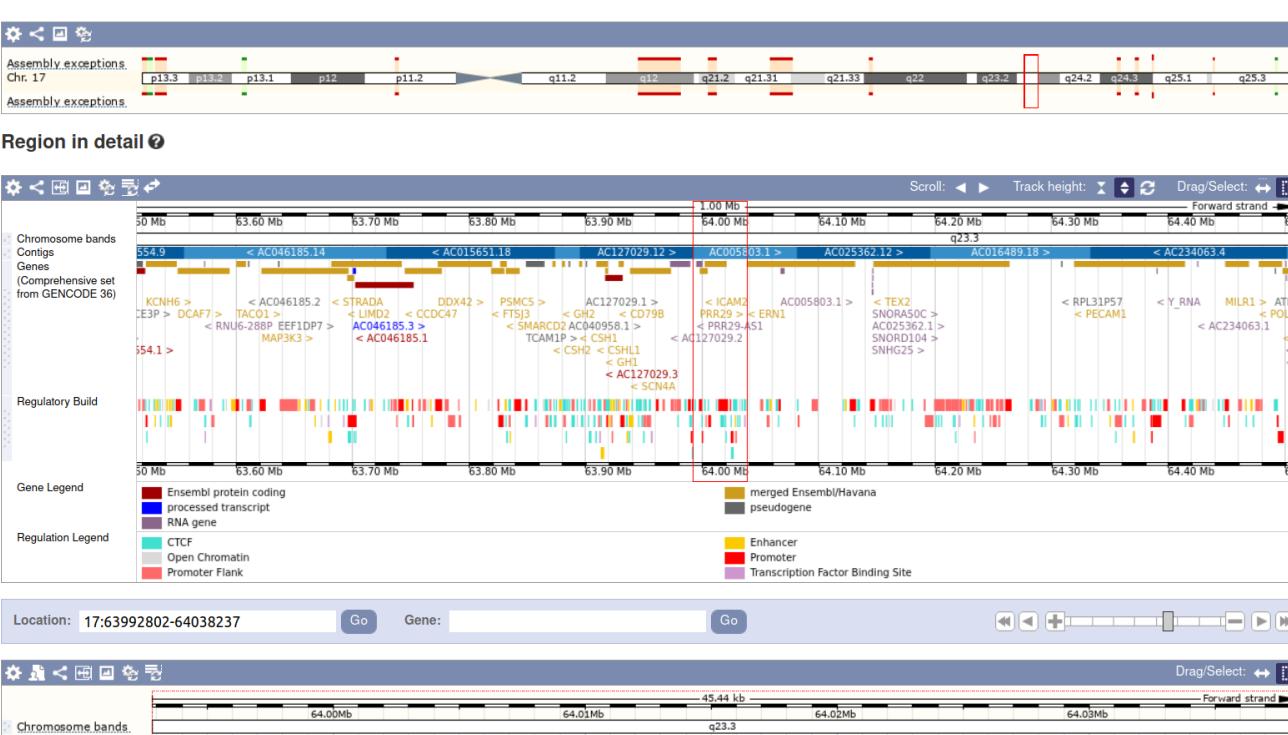
64.00Mb 64.01Mb 64.02Mb 64.03Mb

q23.3

Forward strand

Search all species...

Login/Register 



OrthoDB

The image shows two screenshots of the OrthoDB v11 website. Screenshot A (left) shows the homepage with a search bar containing 'hsp'. A green oval highlights the search bar with the text 'Keyword search with autocomplete'. Another green oval highlights the 'Species to display' section with the text 'Level-of-orthology and species to expand'. A third green oval highlights the 'Search or browse covered species' section at the bottom left. Screenshot B (right) shows a search result for 'hsp70' with 53 groups found. A green oval highlights the 'Drag & drop bookmarklet to the toolbar' option. A green arrow points from the 'Bookmark OrthoDB' link in screenshot B to the 'Toolbar' icon in screenshot A. A green oval highlights the Sankey diagram in the center of the page with the text 'Sankey diagram facilitates navigation by levels-of-orthology'.

Fuente: [Kuznetsov, D. et al. 2023](#)

A Se puede limitar la búsqueda por grupo taxonómico. B El grupo de ortología dispone de un diagrama interactivo.

OrthoDB

The image shows two screenshots of the OrthoDB web interface, labeled C and D.

Screenshot C: A gene search interface. The search bar contains "5141_1:000ecf". Below the search bar, there is a green callout box with the text "Search with an identifier for a gene-centric view". The main content area displays "Gene Info" for *Neurospora crassa*. It includes details such as UniProt ID (GE21DRAFT_9193), function (Endoplasmic reticulum chaperone BiP), and various identifiers like Entrez, UniProt, STRING, and GO terms. A second green callout box highlights the "Orthologs in example species" section, which lists orthologs from *Schizosaccharomyces pombe* (pomb) and *Rattus norvegicus*.

Screenshot D: A search results page for the same gene identifier. The search bar again contains "5141_1:000ecf". The main content area displays "Gene Info" for *Neurospora crassa*, identical to screenshot C. Below this, there is a "Get Ortholog Groups" section with a search bar for "search species (by name)". It shows ortholog groups for *Schizosaccharomyces pombe* (group 985719at4890) and *Saccharomyces cerevisiae* (YJM1573) at the Ascomycota level.

Fuente: [Kuznetsov, D.](#)

et al. 2023

C La visualización del gen permite acceder a genes ortólogos. D Se puede buscar el ortólogo de una especie concreta.

OrthoDB

- Los datos primarios són las proteínas de UniProtKB.
- Utiliza el [programa orthologer](#) para agrupar proteínas ortólogas.
- El programa [BUSCO](#) utiliza a su vez OrthoDB para determinar qué genes suelen estar en copia única en la mayor parte de organismos de un grupo.