

Identificació de seqüències i comprensió d'anotacions

Bases de dades en Biologia

En la actualitat es disposa d'una gran quantitat d'informació sobre seqüència de proteïnes i àcids nucleics. Aquesta informació es troba emmagatzemada a bases de dades de tipus divers:

Bases de dades primàries

Son aquelles que es fan servir com a repositoris de la informació de seqüència obtinguda directament de les fonts, sense processar, o amb un processament només tècnic.

Les més rellevants són:

GenBank (<https://www.ncbi.nlm.nih.gov/genbank/>): Seqüències de DNA, proteïna, genomes complets.

RefSeq (<https://www.ncbi.nlm.nih.gov/refseq/>): Seqüències de referència

ENA (European Nucleotide Archive, <https://www.ebi.ac.uk/ena/>): Seqüències de DNA

Uniprot (<https://www.uniprot.org/>): Seqüències de proteïna. Conté dues seccions: trEMBL: Seqüències obtingudes automàticament; i SwissProt: seqüències de proteïna amb anotacions funcionals acurades.

Protein Data Bank (<https://www.rcsb.org/>): Estructures tridimensionals de macromolècules

Bases de dades secundàries

Aquestes bases de dades inclouen informació addicional que pot provenir d'anotació manual (com SwissProt), o de resultats d'anàlisi bioinformàtica.

Algunes de les més rellevants serien:

SwissProt: La secció anotada de Uniprot

Ensembl (<http://ensembl.org>): Lloc de referència per a informació sobre gens, transcrits, variants de seqüència.

Interpro (<https://www.ebi.ac.uk/interpro/>): Classificació de famílies de proteïna amb força informació funcional.

Comparació de seqüències

La operació més rellevant que es realitza en bioinformàtica és la comparació de seqüències. La semblança de seqüència entre dos gens o dues proteïnes ens permet inferir que les propietats de les mateixes seran semblants. En particular, fem servir la comparació de seqüències per:

- Identificar noves seqüències a partir de la semblança amb les seqüències conegudes
- Avaluar el grau d'homologia, que a la seva vegada ens permet inferir la semblança de seqüència i estructura
- Identificar patrons relacionats amb propietats funcionals com l'activitat enzimàtica.
- Predir estructures tridimensionals de proteïnes
- Predir lloc d'interacció proteïna-DNA

L'eina més habitual per identificar seqüències i localitzar homòlegs és BLAST (<https://blast.ncbi.nlm.nih.gov/>). Blast permet comparar seqüències de DNA o proteïna contra diverses bases de dades.

- Bases de dades no reduntants (nr). Les més extensives, contenen un representant de tots les seqüències conegudes
- Seqüències de referència: Compostes de seqüències seleccionades com a referència (correctament anotades i contrastades).
- SwissProt: Seqüències de proteïna anotades
- PDB: Seqüències de proteïna amb informació estructural.