

# Biomedical Data Bases

Saul Pierotti

December 29, 2019

## Introduction

- A database is an organized collection of data
- A DBMS is a piece of software that allows DB implementation and data mining
  - It allows data storing, retrieval, to perform backups, to maintain security
- Curation is a synonym of manual annotation
- In the last few years there has been a big effort to standardize DBs
- Even in the best DBs there are errors (!)
- Integration among DBs happens with corss-linking and with merging of DBs
  - Uniprot is an example
- Not all the structures available are on PDB (!)
  - Pharmaceutical companies have their own provate DBs
- Some DB terms
  - A record is a DB entry containing different fields
  - An accession key is a unique identifier of a record
  - A table is a DB file containing many records
    - \* Different tables can be connected by the same accession key for some records
  - A query is a data request submitted to a DB
- A query can be organized with boolean operator, in order to retrieve the desired result
- The schema of a DB is the logical structure of the data
  - A schema can be written in a flatfile, XML, ecc.
- The instance is the set of actual data
- The schema allows the interpretation of the instance
- Pubmed is based in Betsheda, at the National Library of Medicine
- Many DBs have built-in a way to retrieve data from the command line in a structured way
  - Entrez Direct allows to retrieve results from Entrez in the Unix command line
- The serch engine of PubMed uses automatic term mapping: if I do not specify the search filed, it tries to guess the correct one for each term in the query
- It first looks for the query in the MeSH Translation Table, then for the Journal translation table and then for an Author Index
- MeSH (Medical Subject Headings) is a list of common search terms that are classified in the correct context
  - If I search MeSH directly for a term, it gives me a list of possible meanings
  - Topics in MeSH are organized in a hierarchical wayH
  - I can limit my search to some topics
- Some of the main servers that we will use
  - NCBI (USA) hosts PubMed
  - SIB (Switzerland) hosts the tool ExPASy
    - \* In ExPASy, we can find SwissModel, UniProt, T-Coffee, LALIGN
    - \* T-Cofee is a really powerful alignment tool from Prof.Cedric Notredame, who maybe will give us an advanced lecture
  - EBI-EMBL (Germany and England)

- \* It started in Heidelberg and then was transferred to Hinxton
  - \* They provide a lot of training schemes, and freely available data
  - \* Ensemble is a genome browser from EBI
  - \* It hosts also InterPro
- Elixir is an european platform for bioinformatics data resources and tools
  - Casadio is on the board of directors of Elixir

## Uniprot

- Uniprot is composed of entries from Swiss-Prot, PIR and TrEMBL, and from external sources
  - PIR was founded by Margaret Dayhoff, the creator of PAM matrices and of the first protein sequence atlas
  - SwissProt was part of SIB and it is the most well-curated part
  - TrEMBL are automatic translations of EMBL DNA sequences
  - It became a central hub for protein information
  - It offers extensive cross-linking
- Uniprot is itself composed of different parts
  - The UniProt Archive (UniParc) is a comprehensive and non-redundant database that contains most of the publicly available protein sequences in the world. Proteins may exist in different source databases and in multiple copies in the same database. UniParc avoided such redundancy by storing each unique sequence only once and giving it a stable and unique identifier (UPI) making it possible to identify the same protein from different source databases. A UPI is never removed, changed or reassigned. UniParc contains only protein sequences. All other information about the protein must be retrieved from the source databases using the database cross-references. UniParc tracks sequence changes in the source databases and archives the history of all changes. UniParc has combined many databases into one at the sequence level and searching UniParc is equivalent to searching many databases simultaneously.
  - The UniProt Knowledgebase (UniProtKB) is the central hub for the collection of functional information on proteins, with accurate, consistent and rich annotation. In addition to capturing the core data mandatory for each UniProtKB entry (mainly, the amino acid sequence, protein name or description, taxonomic data and citation information), as much annotation information as possible is added. This includes widely accepted biological ontologies, classifications and cross-references, and clear indications of the quality of annotation in the form of evidence attribution of experimental and computational data.
  - The UniProt Reference Clusters (UniRef) provide clustered sets of sequences from the UniProt Knowledgebase (including isoforms) and selected UniParc records in order to obtain complete coverage of the sequence space at several resolutions while hiding redundant sequences (but not their descriptions) from view. Unlike in UniParc, sequence fragments are merged in UniRef: The UniRef100 database combines identical sequences and sub-fragments with 11 or more residues from any organism into a single UniRef entry, displaying the sequence of a representative protein, the accession numbers of all the merged entries and links to the corresponding UniProtKB and UniParc records.
- When we publish something using a database, we need to specify the release that we used (!)
- SwissProt is composed of proteins manually curated entries from TrEMBL
  - The first step of curation is to assess if there is evidence at the protein level
  - GO-terms (and other ontologies), references, isoforms, PTMs (post translational modifications), polymorphisms, crosslinks are added
  - If available, the structure of the protein is reported, modelled or experimental
  - The correct nomenclature is established and reported
- If different isoforms of the protein are identified, the curator searches for the reason of the difference and documents it
- Different isoforms can be due to alternative splicing, natural variations, alternative initiation sites, erroneous prediction, ecc.

- Gene ontology is a classification of a sequence in terms of biological process, molecular function and cellular component
- The TrEMBL part of UniProtKB is mainly redundant, while SwissProt is not-redundant
- How to recognize a SwissProt entry from a TrEMBL entry in Uniprot
  - The accession key for TrEMBL proteins is an alphanumeric string, and the entry name repeats the accession key
  - The accession key for SwissProt proteins is an alphanumeric string, and the entry name is not related to the accession key
  - In the result page, reviewed proteins have a golden badge, non reviewed proteins have a blue badge
- In Uniprot we can customize the result list to show the information we need
- Entries in UniProt are huge (!)
- Each entry has an annotation score in 1 to 5
- Functional information is a bit sparse in the entry, it can be in many different fields
- There is also a comment field in the entry
- We did an exercise about how to do an advanced search on UniProt and download the result as an Excel file, that I cannot reproduce here
  - We basically have to use boolean terms for the search, select the needed columns for display and then download in the required file format
- To retrieve the IDs for other databases, we use the Retrieve/ID mapping section
  - Many databases have a mapping feature, but UniProt is the best place for doing ID mapping (!)
  - UniProt mapping can be done also from the Unix command line
- In order to facilitate the retrieval of information from an entry, UniProt created a visual interface called feature viewer
  - This is required because of the high information density of UniProt entries, which makes it difficult to read all the information in a textual format
- Automatic annotation in TrEMBL uses data from SwissProt to infer information
  - It is composed of 2 main pipelines: Unirule and SAAS
- Unirule uses expert-curated rules
  - For each automatic annotation, we can see the set of rules that generated it
- SAAS uses machine learning for finding the best set of rules, based on sequence features
- HAMAP is used as a partially automated rule source in Unirule
  - It is effective when the protein is part of a well defined family
  - It is used in place of manual annotation when it gives a result at least as good as the manual one
  - It claims to be more error-free than manual annotation
- UniParc is the biggest non-redundant protein sequence database, it contains sequences from various sources
  - It is not annotated
- UniRef provides clusters of sequences from UniProtKB and selected UniParc entry
  - It can be accessed from any entry going to the “similar sequences” link
  - UniRef100 combines identical sequences with at least 11 residues
  - UniRef90 clusters UniRef100 sequences that share at least 90% sequence identity and 80% overlapping with the seed
  - UniRef50 in the same way clusters sequences from UniRef90
  - It is useful because I can use the seed as a representative of the cluster, and so reduce the number of entries that I analyse
    - \* It can make BLAST searches faster (!)
  - The seed is the longest or best annotated sequence
    - \* It is considered annotation score, if it is manually annotated, and the source organism
    - \* If the source organism is a model organism it is preferred
- All the changes that an entry witnessed are stored in UniProt in the “history” link
- In UniProt I can also find complete proteomes
  - They are a collection of all the entries from an organism with a completely sequenced genome
  - Reference proteomes are the ones with best annotated entries and usually belong to model organisms
  - They evolve with time, when new proteins are added

- They involve a lot of manual curation
- There are 25000 genes, 100000 transcripts and 1000000 proteins in humans
- Uniprot can be accessed programmatically through a python library

## Errors in databases

- For a long time sequences have been annotated on the basis of structural and sequence similarity
- This approach has its limits, since not always structural or sequence similarity correspond to similar function
- Burkard Rost published a paper where he says that enzyme function is less conserved than expected
- The current level of misannotation is unknown, and seems to be increasing
- The level of misannotation in enzyme superfamilies that comprise families with different functions is particularly high
- The analysis of the conservation of catalytic residues can improve annotation accuracy

## Protein classification databases

- Classifying proteins reduces the amount of data we need to deal with, and allows us to make predictions about protein function
- Interpro and Pfam are databases of protein classification based on sequence
- Proteins can be classified in superfamilies, families, subfamilies
  - The deeper we go, the more we can infer functional features for the members
- A domain is an independently-folding unit of protein sequence which harbours a function
- Sequence features are active sites, binding sites and other characteristics short sequences found inside domains
- A protein signature is a mathematical model build from multiple sequence alignments
  - Motives, fingerprint, profiles, HMM are signatures
- A single motif, also called pattern, can be represented by a regular expression
  - Prosite is a DB of patterns
- A fingerprint is composed of multiple motives in a specific arrangement
  - PRINTS is a DB for fingerprints
  - Fingerprints are useful for differentiating subfamilies
- A profile is built from a multiple global alignment and consists of a position-specific scoring matrix (PSSM)
- HMMs are complex algorithms capable of modelling the probabilities of residue change, insertion and deletion
  - They are a really powerful protein signature
- The general workflow of building a signature is to do a MSA, generate a draft model, use it against the all Uniprot and develop a mature model
- InterPro is a meta-DB that gathers data about patterns, profiles, fingerprints, HMMs
  - It is an hub of protein classification and annotation
  - It also gathers data from UniProt and PDB
  - It is cited by uniprot in each entry
- Pfam families are made from a seed alignment manually curated, and then are refined by searching the whole uniprot
  - The mature MSA is automatic, while the seed is curated manually
  - A Pfam family is represented by an HMM
  - Superfamilies in Pfam are called clans
  - From 2016 Pfam started to use reference proteomes instead of the whole Uniprot, to decrease redundancy
- The number of protein folds is limited, and classifying protein based on their fold allows to understand their functional and structural relationships
- Topology schemes are useful for identifying the fold of a protein

- We can represent the topology of a protein by drawing its secondary structure as a flat diagram
- Structural classification DBs are SCOP and CATH/Gene3D
- It seems that we have discovered all fold types
  - No new folds have been reported in SCOP and CATH since 2010
- SCOP in itself is dead, but its legacy has been taken by SCOPe
  - The first level classification is alpha, alpha/beta, alpha+beta, all beta, small proteins
  - The classification is class>fold>superfamily>family>sequence
- SCOP2 is a different DB, made from the same people that made SCOP, that uses a different classification system based on networks
- CATH classifies structures from the PDB while Gene3D predicts the location of functional domains in available sequences (Uniprot) using informations from CATH
  - The first level classification is Mainly alpha - Alpha-beta - Mainly beta
  - The classification is class>architecture>topology or fold>homologous superfamily>sequence family
- ECOD is another DB that claims to classify many more domains than CATH and SCOP (it makes a finer classification)
  - Its classification is based on evolutionary relationships
  - It is a newcomer, we will see how it will go

## X-ray crystallography

- Bragg's law:  $2d \sin \theta = n\lambda$ 
  - $2d \sin \theta$  is the excess distance traveled by light when hitting a diffraction plane with distance  $d$  with an incidence angle  $\theta$
  - $n\lambda$  is an integer number of wavelengths
  - If the additional distance traveled is an integer number of wavelengths, I have constructive interference and therefore a signal
  - The diffracted beam is deviated by  $2\theta$  from the incident beam
- The distance of a diffraction spot  $r$  from the expected incidence point of the primary beam on a detection screen placed at distance  $A$  from the crystal is  $A \tan 2\theta$
- Synchrotron produce X-rays with really high brilliance, so I have an high signal-to-noise ratio
  - They allow to use smaller crystals
- The inverse Fourier transform is equivalent to passing the radiation through a lens
  - It does a Fourier synthesis
  - We do not have materials that can focus X-rays
- The X-rays are scattered by the electron cloud and never reach the nuclei, this is why we get electron densities
- We use crystals instead of single molecules because it amplifies the signal and reduces the noise
- The reciprocal space of a Fourier transform is a graph were each frequency is reported with a corresponding amplitude
  - I always get pairs of frequencies, for example 3 and -3, because I cannot discriminate the direction of a wave
- A 2d wave can be represented in the reciprocal space as points on the plane, were the combinations of  $x$  and  $y$  frequency is represented
- The coordinates in the reciprocal space are called  $h$  and  $k$
- When I sum more 2d waves, in the reciprocal space I just put all the points deriving from them together
- The phase of the different frequencies refers to how they allign together
- If I apply a low-pass filter to the reciprocal space, I get back a less detailed image
  - High frequencies are responsible for fine details
- On the opposite, with an high-pass filter I get only the details but not the bulk image
- From a 3d Fourier transform I get a 3d reciprocal space with coordinates  $hkl$
- From rotating the crystal, I get a series of 2d images that allow to reconstruct the 3d reciprocal space
- The space groups refers to the rotational symmetries of the unit cell
  - It is important for the crystallographer because they determine the rotations that I have to apply

to my crystal during data collection

## NMR

- The conditions in which it is performed are more similar to those encountered in vivo
- It uses the same principle as MRI
- It allows to determine also the structure of mobile regions
  - This is also used for studying protein folding
- It cannot be used for big proteins (>80kDa)
- The wavelength is on the microwave spectrum
- An NMR structure is an average of many models
- Sometimes also the single models in the bundle are deposited

## CryoEM

- The resolution has increased enormously in the last 20 years
  - We went from 25 Å in 2005 to 2 Å in 2015
- Lenses for an electron beam are electromagnets
- What allowed a great increase in resolution is cooling the sample, so that it does not get damaged by the electron beam
- We started also to use very low intensity of the beam in order to avoid damage, but this lowers the signal quality
  - In order to increase quality we take many images and average them
- The cooling is done very fast so that water solidifies in an amorphous form
- The microscope sees many different molecules, all oriented randomly
- An algorithm clusters the images by orientation and averages the ones in the same orientation and conformation, so to increase quality
- From the averaged images, I can obtain a 3d structure
- In the last years a series of small improvements increased significantly resolution
- CryoEM gives you an electron density (!)
- Size does not matter (!)

## PDB file

- The B factor is the mean square deviation of the position, so it is measured in Å<sup>2</sup>
- There are several parsing tools for PDB files
- Now data are deposited also in xml format