



# **Integrated Access to Biological Data**

**Budva, 11th of June 2006**

**Ainhoa Lorente**

robotiker  
tecnalia

- Introduction and objectives
- Biological Data Repositories
- Biological ontologies
- Ontology merging and mapping
- Database annotation
- Example of an ontology merging/mapping
- Conclusions
- Questions and answers

# I. Introduction and objectives

To handle biological data repositories by means of semantic and artificial intelligence technologies.

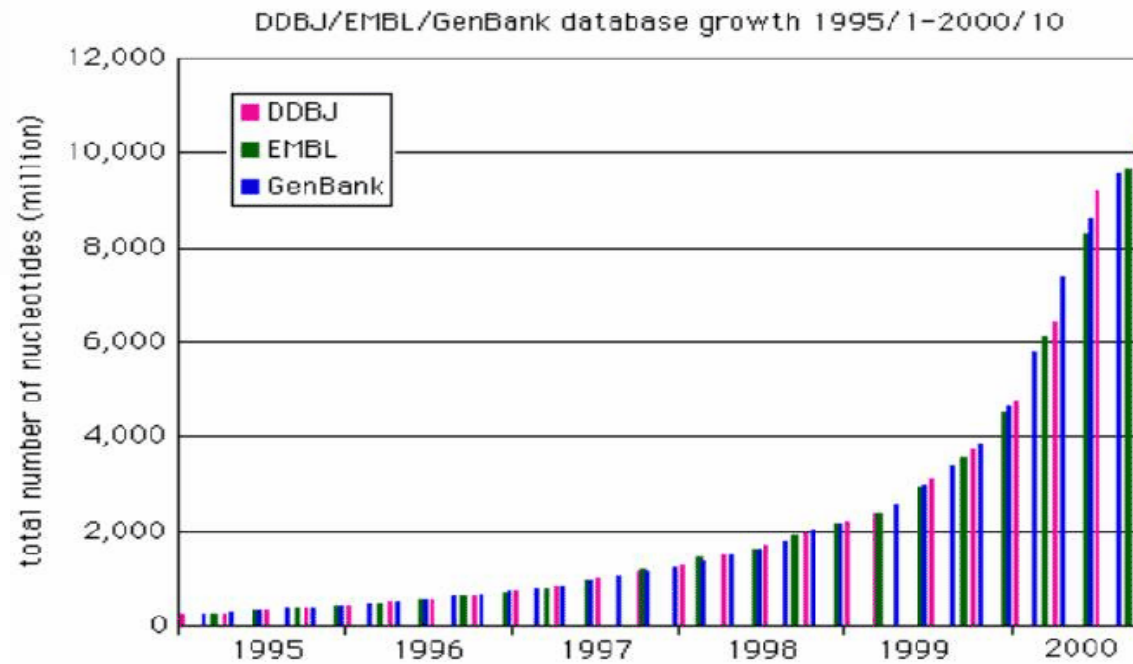
- **Starting point:**
  - ▶ Human genome sequencing has given rise to a great number of biological data repositories.
- **Problem addressed:**
  - ▶ Quantity and heterogeneity
- **Our aim:**
  - ▶ To provide an unified access point to diverse biological data repositories.
- **Our challenge:**
  - ▶ To change the existing vision of ontologies in biology:
    - ▶ **Up to now:** As mere guides for data structure.
    - ▶ **From now on:** As integrated modelling of the biological data by combining or associating ontologies.

## II. Biological Data Repositories (i)

### Most important categories

- Nucleotides Sequences → DNA
- Amino acid Sequences → Proteins
- Gene expression
- Scientific literature
- Corporate databases
- Health cards

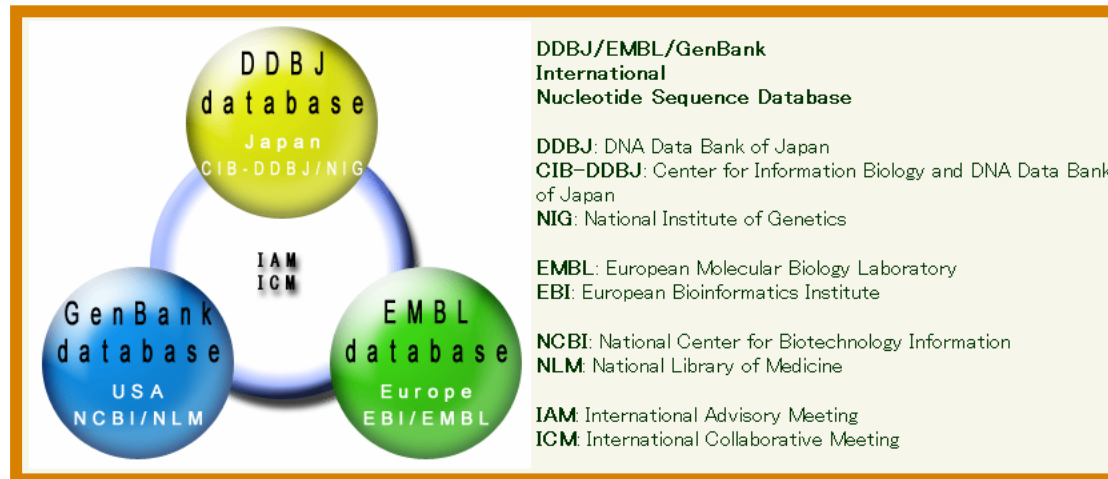
## II. Biological Data Repositories (ii)



<http://www.math.tau.ac.il/~rshamir/algmb/00/scribe00/html/lec05/node3.html>

## II. Biological Data Repositories (iii)

### THE MOST IMPORTANT DNA DATABASES



## II. Biological Data Repositories (iv)

### ■ The most important Protein Databases

- ▶ **SwissProt:** (<http://us.expasy.org/sprot/>).
- ▶ **PIR:** Protein Information Resource.  
(<http://pir.georgetown.edu/>).
- ▶ **PDB:** Protein Data Bank (<http://www.rcsb.org/pdb/>).

### ■ Gene Expression

- ▶ **GDX**
- ▶ **ExpressDB** (<http://arep.med.harvard.edu/ExpressDB>).

### ■ Scientific Literature

- ▶ **MEDLINE**
- ▶ **PubMed**
- ▶ **UpToDate**

### ■ Corporate Databases

### ■ Health Cards:

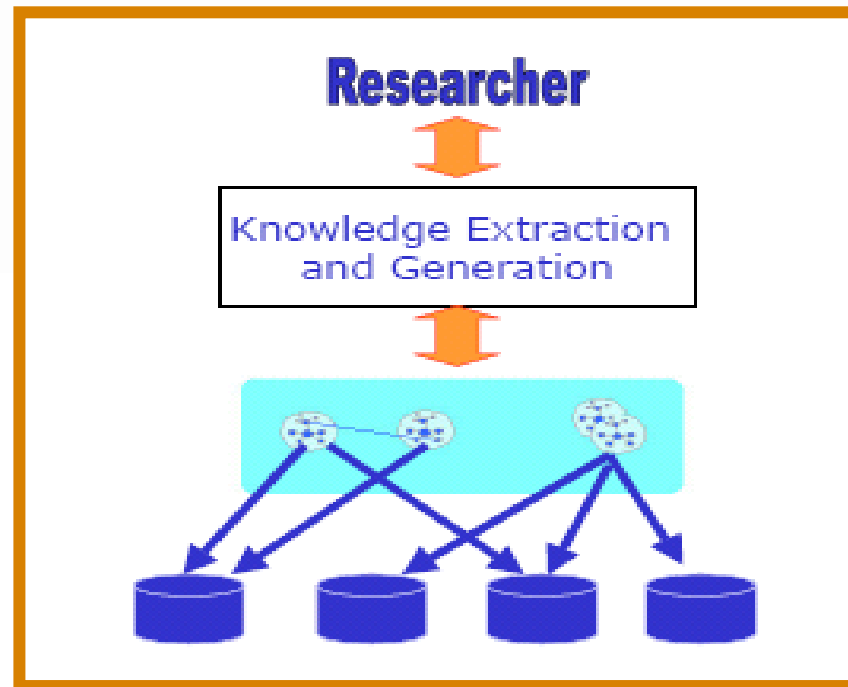
#### ▶ **Future Use:**

- ▶ To facilitate the medical care attention
- ▶ To increase patients mobility comfort
- ▶ Administrative tasks
- ▶ Emergency health cards
- ▶ Specific care records
- ▶ Patients general medical records
- ▶ To match genetic patient data with biological databases

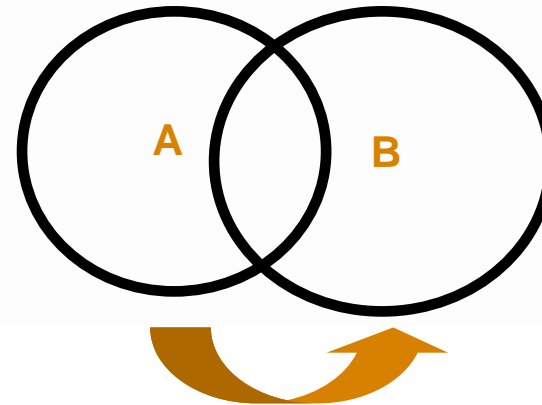
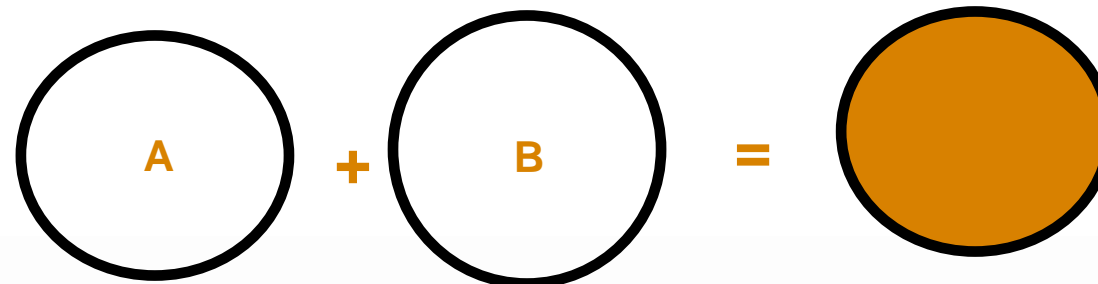
#### ▶ **Challenges:**

- ▶ To unify the data to store
- ▶ Unification of the media
- ▶ Unification of medical nomenclature



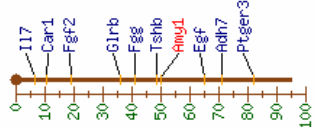
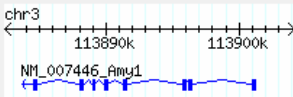


- **Gene Ontology**
- **The Microarray Gene Expression Data (MGED)**
- **UMLS**

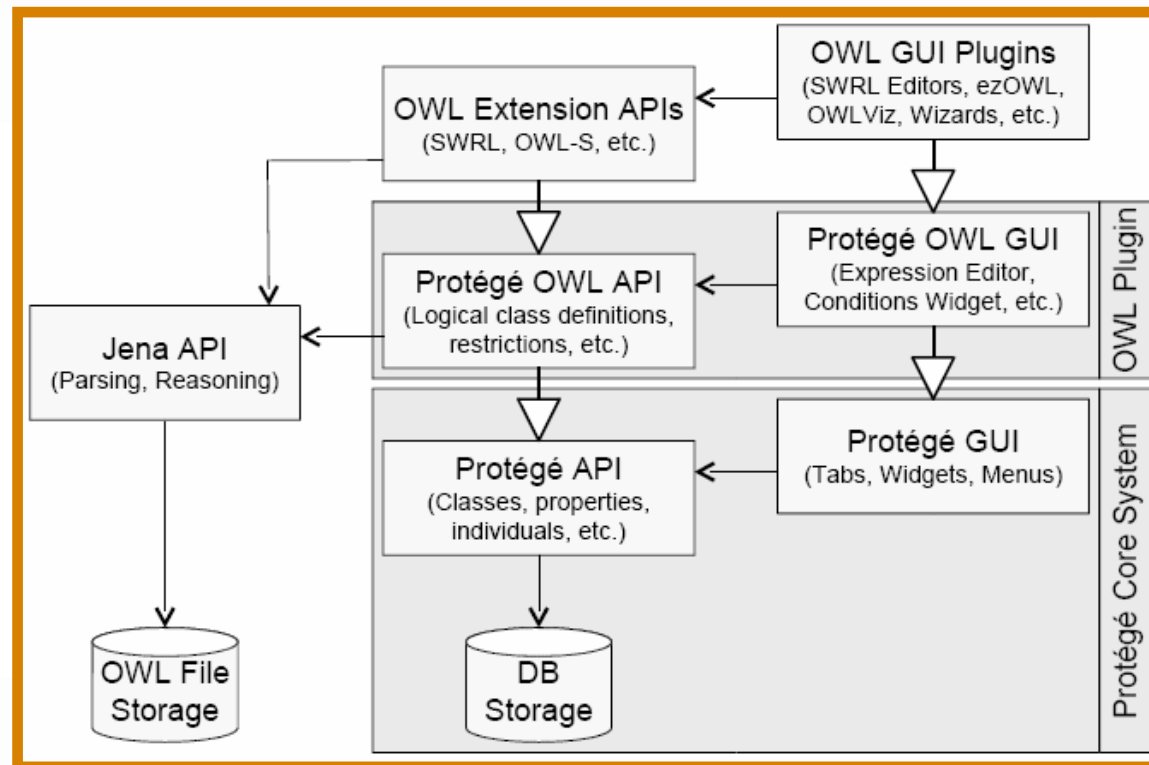
■ **Ontology Merging:**■ **Ontology Mapping:**

## Gene ontology annotations

?
**Gene Detail**
Your Input Welcome

<b>Symbol Name ID</b>	<b>Amy1</b> amylase 1, salivary MGI:88019	<a href="#">Nomenclature History</a>																
<b>Synonyms</b>	Amy-1																	
<b>Genetic Map</b>	Chromosome 3 50.0 cM <a href="#">Detailed Genetic Map ± 1 cM</a> Mapping data( <a href="#">65</a> )																	
<b>Sequence Map</b>	113882378-113904173 bp, - strand (From Ensembl annotation of NCBI Build 33) <a href="#">Ensembl ContigView</a>   <a href="#">UCSC Browser</a>   <a href="#">NCBI Map Viewer</a>	 <a href="#">MGI Mouse GBrowse</a>																
<b>Mammalian orthology</b>	human; cattle; rat ( <a href="#">Mammalian Orthology</a> ) Comparative Map ( <a href="#">Mouse/Human Amy1 ± 2 cM</a> )																	
<b>Sequences</b>	<table border="0" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left;">Representative Sequences</th> <th style="text-align: left;">Length</th> <th style="text-align: left;">Strain/Species</th> <th style="text-align: left;">Flank</th> </tr> </thead> <tbody> <tr> <td><input type="checkbox"/> genomic ENSMUSG00000045402 <a href="#">Ensembl Gene Model</a>   <a href="#">MGI Sequence Detail</a></td> <td>21796</td> <td>C57BL/6J</td> <td>± <input type="text" value="0"/> kb</td> </tr> <tr> <td><input type="checkbox"/> transcript NM_007446 <a href="#">RefSeq</a>   <a href="#">MGI Sequence Detail</a></td> <td>1771</td> <td>-</td> <td></td> </tr> <tr> <td><input type="checkbox"/> polypeptide P00687 <a href="#">SWISS-PROT</a>   <a href="#">EBI</a>   <a href="#">MGI Sequence Detail</a></td> <td>511</td> <td>Not Applicable</td> <td></td> </tr> </tbody> </table> <p>For the selected sequences <input type="text" value="download in FASTA format"/> <input type="button" value="Go"/></p> <p>All sequences(<a href="#">18</a>)</p>		Representative Sequences	Length	Strain/Species	Flank	<input type="checkbox"/> genomic ENSMUSG00000045402 <a href="#">Ensembl Gene Model</a>   <a href="#">MGI Sequence Detail</a>	21796	C57BL/6J	± <input type="text" value="0"/> kb	<input type="checkbox"/> transcript NM_007446 <a href="#">RefSeq</a>   <a href="#">MGI Sequence Detail</a>	1771	-		<input type="checkbox"/> polypeptide P00687 <a href="#">SWISS-PROT</a>   <a href="#">EBI</a>   <a href="#">MGI Sequence Detail</a>	511	Not Applicable	
Representative Sequences	Length	Strain/Species	Flank															
<input type="checkbox"/> genomic ENSMUSG00000045402 <a href="#">Ensembl Gene Model</a>   <a href="#">MGI Sequence Detail</a>	21796	C57BL/6J	± <input type="text" value="0"/> kb															
<input type="checkbox"/> transcript NM_007446 <a href="#">RefSeq</a>   <a href="#">MGI Sequence Detail</a>	1771	-																
<input type="checkbox"/> polypeptide P00687 <a href="#">SWISS-PROT</a>   <a href="#">EBI</a>   <a href="#">MGI Sequence Detail</a>	511	Not Applicable																

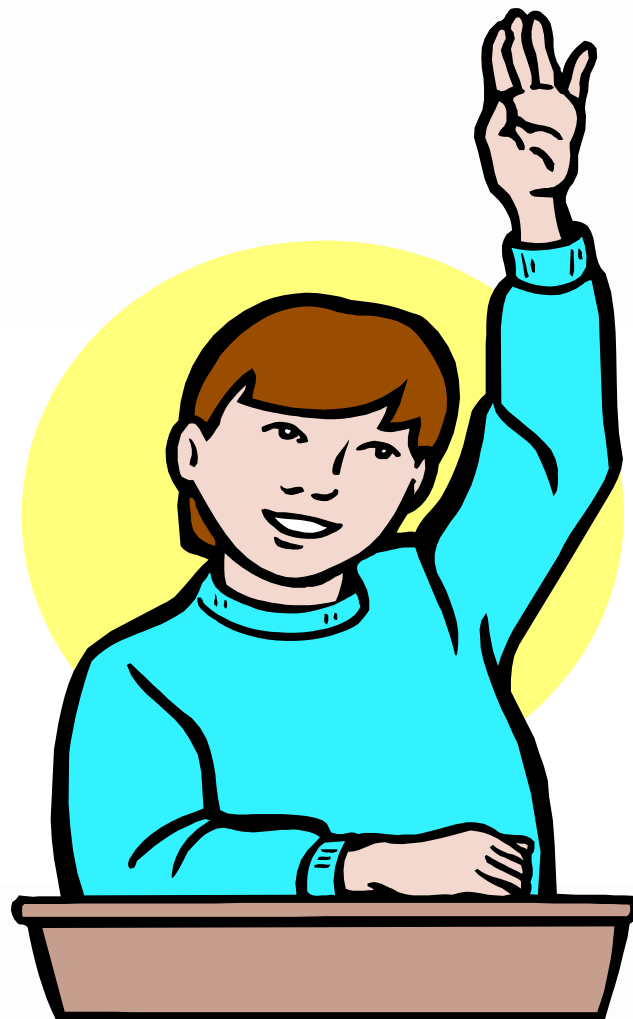
### Protege Architecture with OWL Plugin



## VII. Conclusions

- The application of semantic web technologies to the biological domain is rather limited because:
  - ▶ the semantic web technologies
  - ▶ the tools needed to implement themare still under development.
- Future advances could be applied to solve the immediate scientific needs:
  - ▶ Data aggregation and interoperability.
  - ▶ Unique entry point for data and processes.
  - ▶ Agreement in terminology.
  - ▶ Syntax and semantics related to biological data.
  - ▶ Semantic data annotation to turn human- understandable data into machine-understandable data.
  - ▶ Inference languages to extract and generate knowledge from aggregated data.

## VIII. Questions and answers





**More information:  
allorete@robotiker.es**

robotiker  
tecnalia