

Relational DBs and Repositories

DAVID L. TABB, PH.D.

Overview

- Storing complex structures in files
- Scaling storage via relational databases
- Repositories for biomolecular data:
 - Sequences and variants
 - Transcriptomes
 - Structures
 - Proteomes

JSON: streamlined messages among tools

```
{  
  "contentManagementSystems" : [  
    {  
      "name": "WordPress",  
      "percentMarketShare": 58.9  
    },  
    {  
      "name": "Joomla",  
      "percentMarketShare": 6.1  
    },  
    {  
      "name": "Drupal",  
      "percentMarketShare": 4.9  
    }  
  ]  
}
```

← Array of objects

↑ Name: Value pair

- Web services frequently communicate data by JSON.
- JavaScript Object Notation began at Netscape, but many languages now support it.

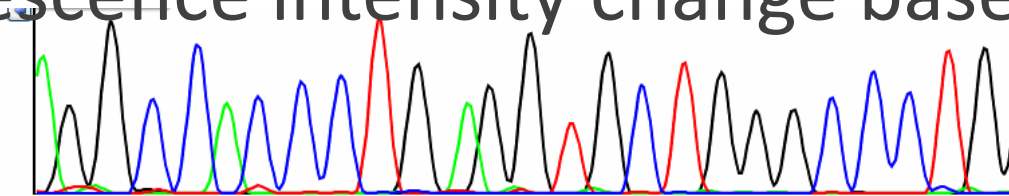
XML labels contents through markup

```
<?xml version="1.0" encoding="UTF-8"?>
<CATALOG>
  <CD>
    <TITLE>The very best of</TITLE>
    <ARTIST>Cat Stevens</ARTIST>
    <COUNTRY>UK</COUNTRY>
    <COMPANY>Island</COMPANY>
    <PRICE>8.90</PRICE>
    <YEAR>1990</YEAR>
  </CD>
  <CD>
    <TITLE>Unchain my heart</TITLE>
    <ARTIST>Joe Cocker</ARTIST>
    <COUNTRY>USA</COUNTRY>
    <COMPANY>EMI</COMPANY>
    <PRICE>8.20</PRICE>
    <YEAR>1987</YEAR>
  </CD>
</CATALOG>
```

- XML encloses data in tags that label content.
- This markup type is also used in HTML, a sibling.
- A separate DTD or schema defines allowable elements.

High-throughput Sequencing

- Electropherogram: a trace showing fluorescence intensity change base by base



- FASTQ (text): sequences of reads with per-base quality assessments

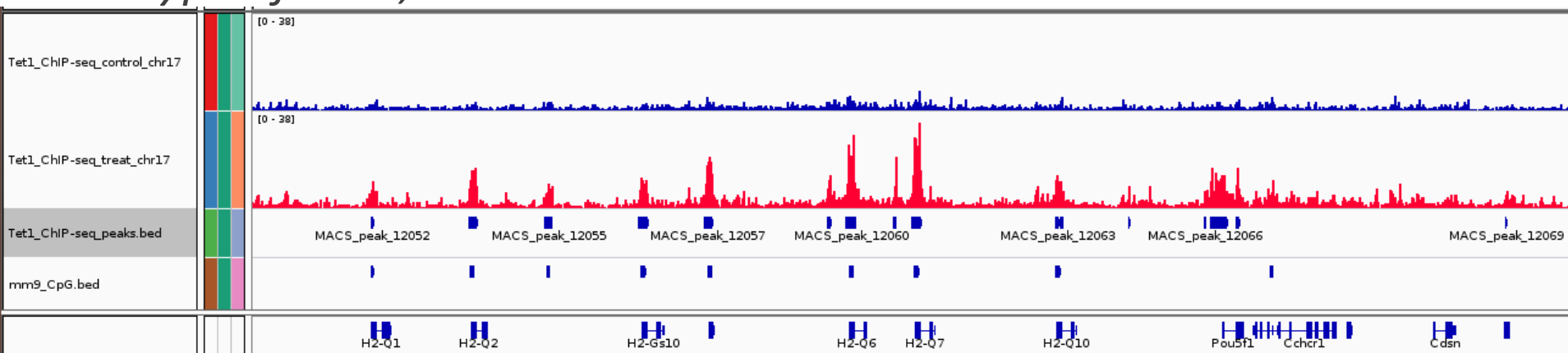
Modern sequencers routinely discard electropherograms due to size

<http://www.gendx.com/>

Reporting alignment products

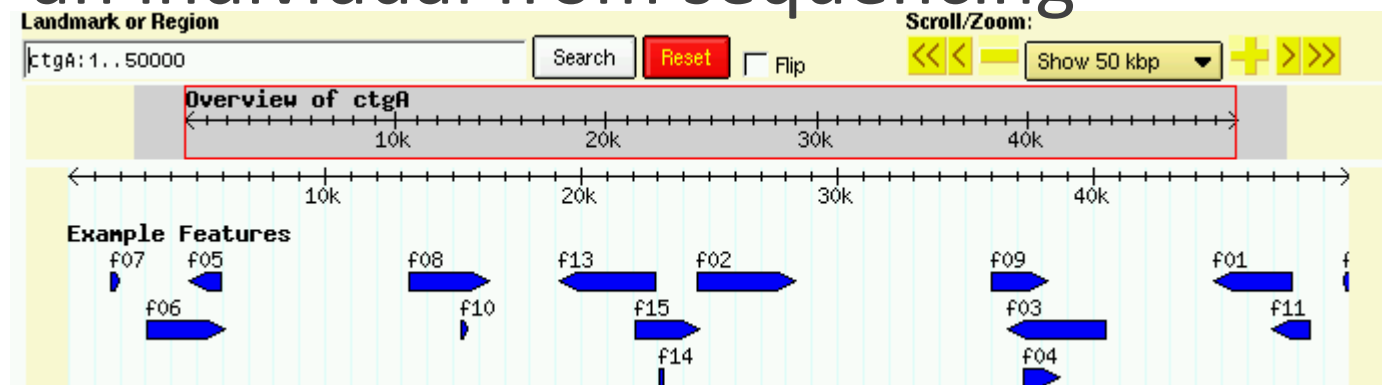
- SAM(text)/BAM(binary): Sequence Alignment Maps align reads to reference annotation
- BED (text): Browser Extensible Data defines genomic features for visualization

BAM and BED represent different emphases for the same type of data; converters exist.



Highly distilled information

- GFF (text): General Feature Format supplies coordinates of genetic features
- VCF (text)/BCF (binary): Variant Call Format details variant nucleotides (mutations or SNPs) for an individual from sequencing



<https://github.com/The-Sequence-Ontology/Specifications/blob/master/gff3.md>

<https://samtools.github.io/hts-specs/>

<http://chlamycyc.mpimp-golm.mpg.de/gbrowse/tutorial/dbgff/tutorial.html>

RDMS and SQL history

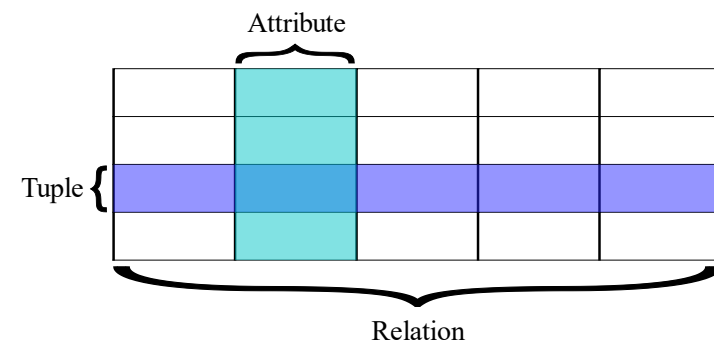
- Relational database management systems (RDMS) were defined in 1970 by E.F. Codd.
- IBM created the SEQUEL (later SQL) language to support RDMS manipulation.
- Essentially any modern database can support SQL queries: MS Access, SQLite, MySQL, PostgreSQL, Oracle...



<https://www.w3schools.com/sql/>

https://docs.oracle.com/cd/B12037_01/server.101/b10759/toc.htm

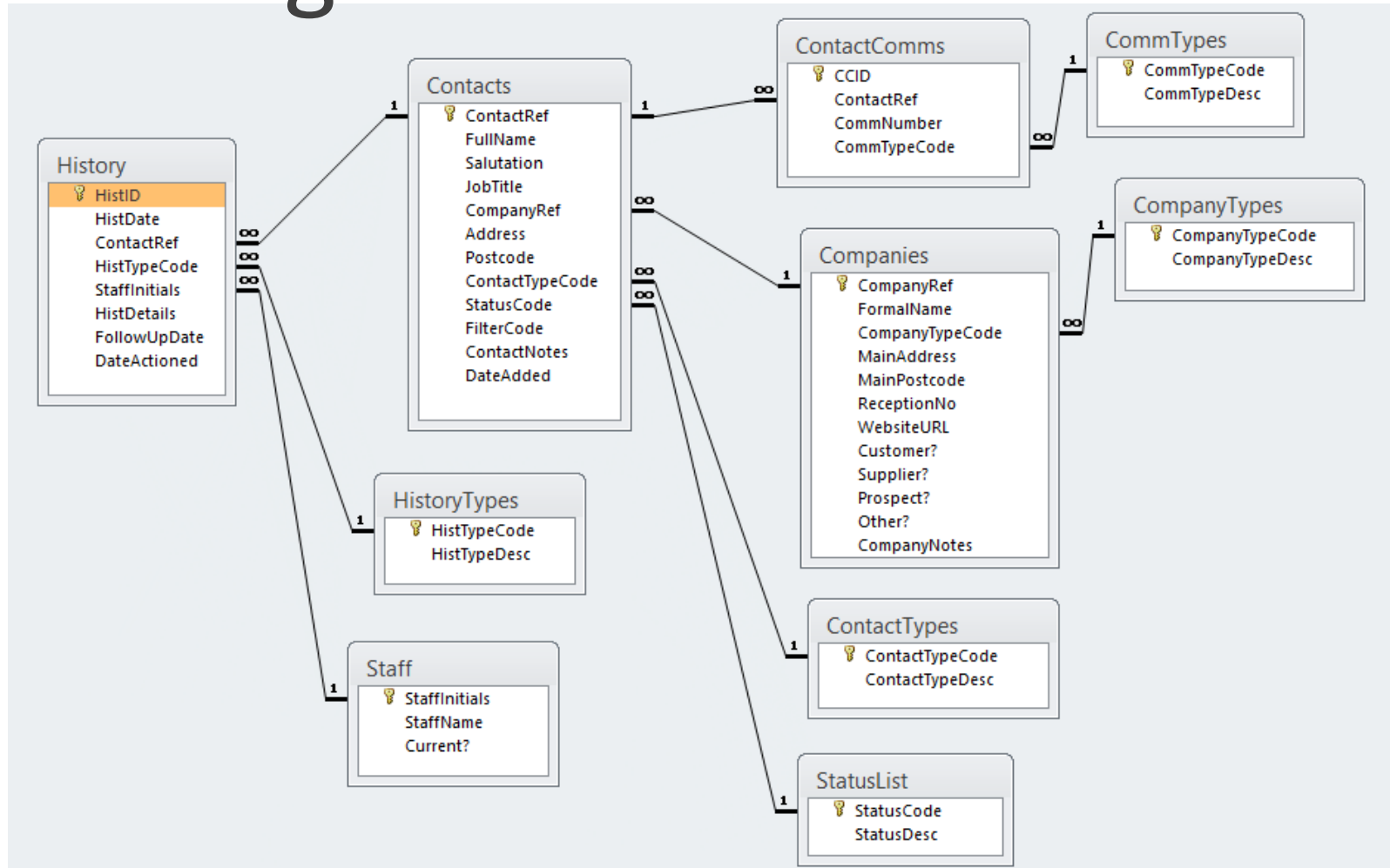
Keys ID each *tuple*



- A single record or row from a table is frequently called a *tuple*.
- Some tables have a column (attribute) that is unique for each row. This is a *primary key*.
- Keys make it possible to relate tables to each other very rapidly in a *join*.



Joins connect information among tables



Relational databases

model complex structures

- Index: a lookup table of a unique value for each record that enables rapid retrieval.
- Inner Join: combine data from two tables, aligning on a field in common to both.
- Query: select fields of data from a body of data (whether single table or a “join”).

```
SELECT ID, NAME, AMOUNT, DATE FROM  
CUSTOMERS INNER JOIN ORDERS ON  
CUSTOMERS.ID = ORDERS.CUSTOMER_ID;
```

The case for normalization

- Update anomaly: adding a record requires re-entry of extant data for consistency
- Insertion anomaly: we may not have all the fields we need when adding a record
- Deletion anomaly: when we remove a record, we lose more data than intended

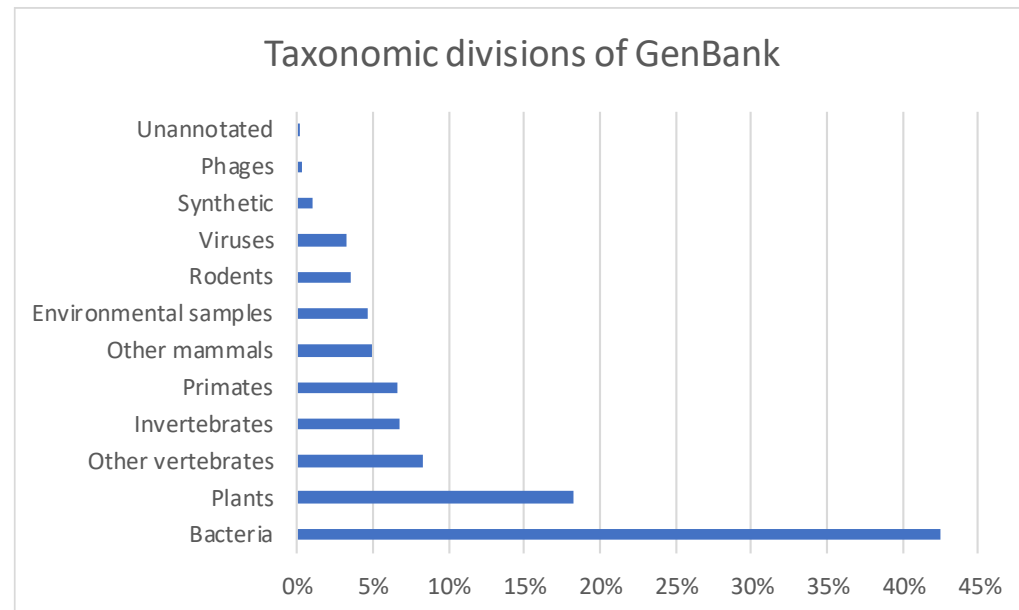
[http://www.studytonight.com/
dbms/database-normalization.php](http://www.studytonight.com/dbms/database-normalization.php)

ID	Name	City	Subject
401	Adam	Durbanville	Bio
402	Alex	Bellville	Maths
403	Stuart	Cape Town	Maths
404	Adam	Durbanville	Physics

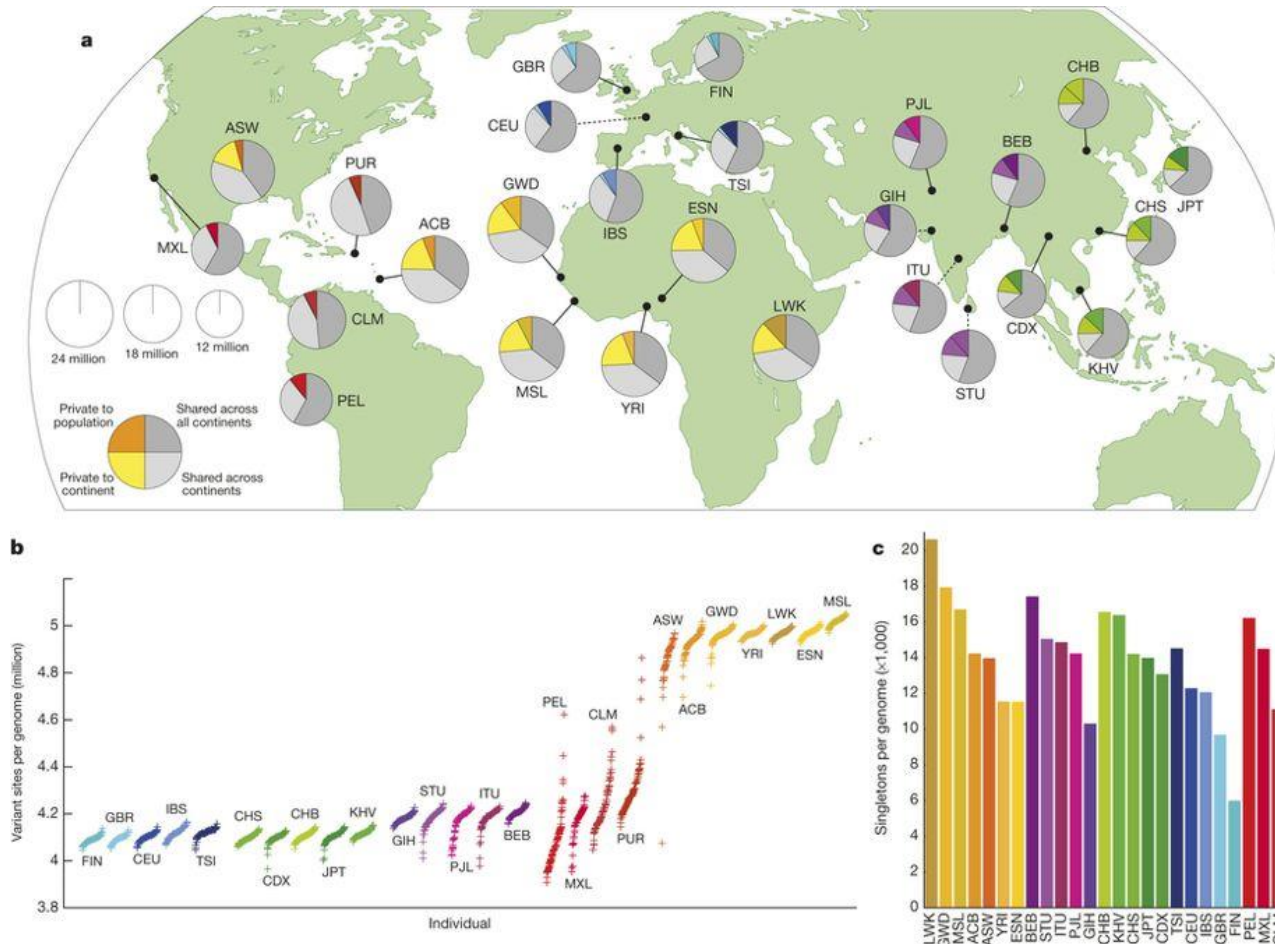
Repositories house vast collections of sequences.

- GenBank / Euro Nucleotide Archive / DDBJ Synchronized DNA repositories at NCBI, European Molecular Biology Lab, and DNA Data Bank of Japan

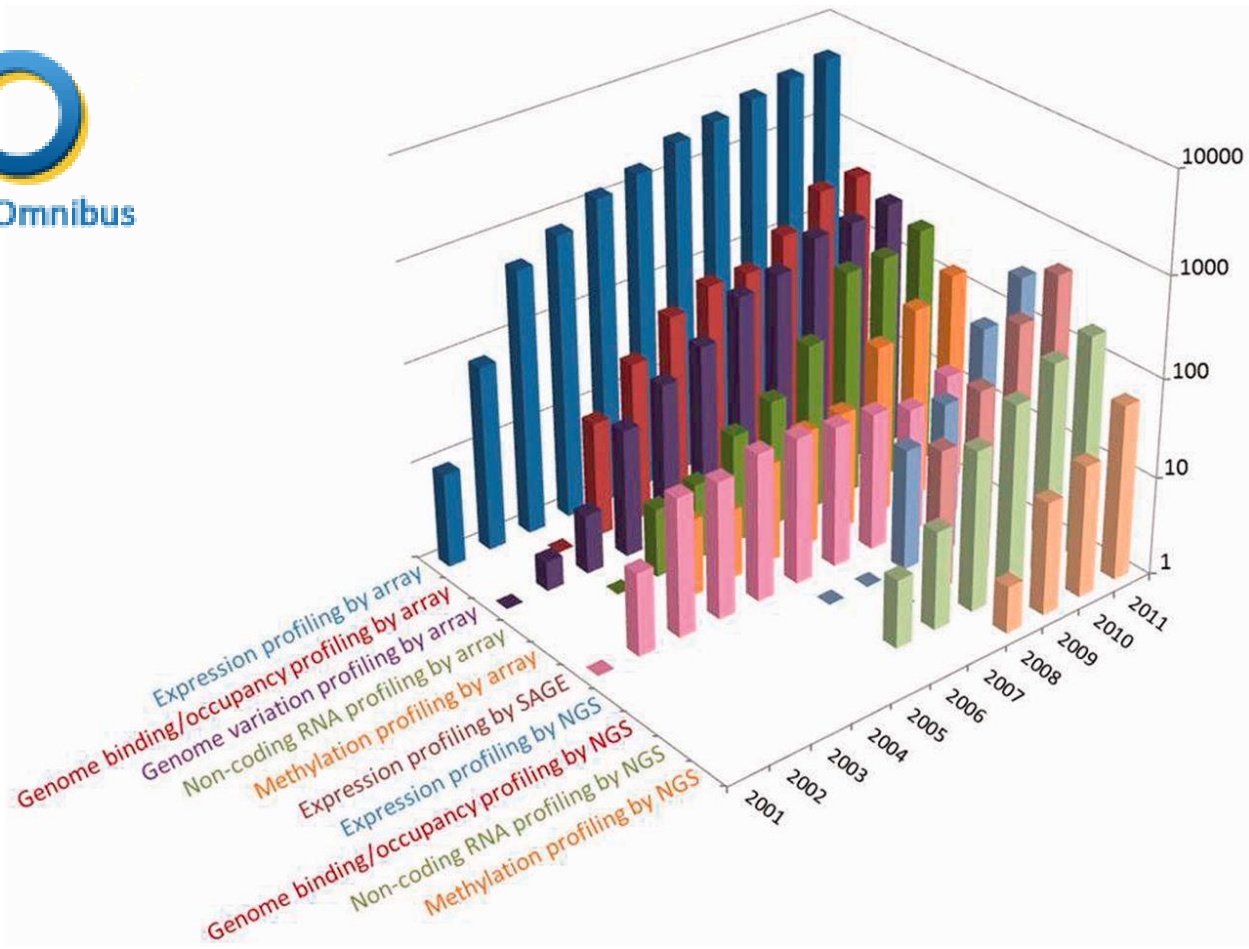
- #1 species:
Homo sapiens
19 752 523 722 bp



“1000 Genomes” project data: www.internationalgenome.org



Gene Expression Omnibus: Functional genomics galore!

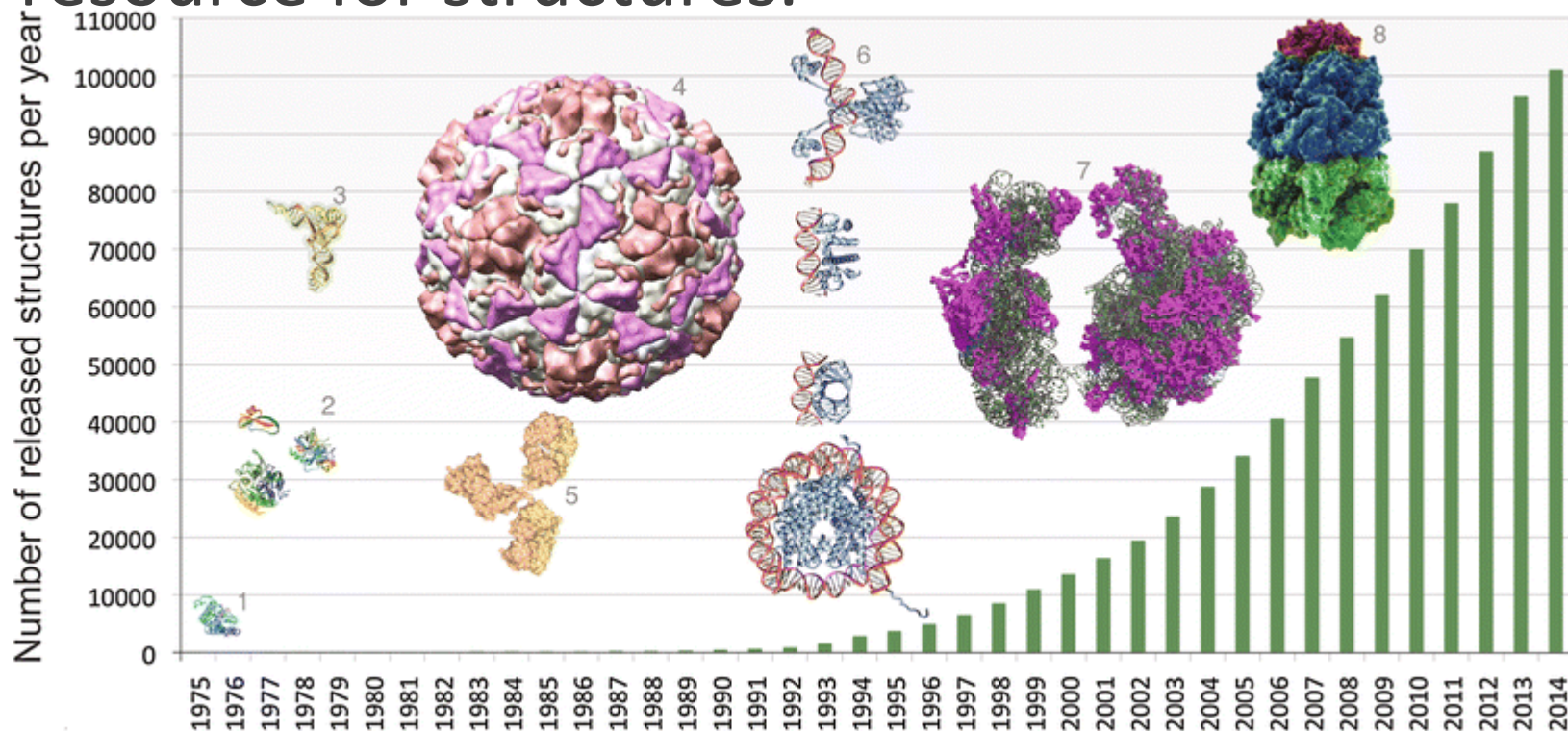


Protein structures since 1971



The Protein Data Bank is an international resource for structures.

<http://www.wwpdb.org/>



Proteomics data from via OmicsDI



The screenshot shows a web browser window with the URL [https://www.omicsdi.org/search?q=Bacillus%20licheniformis-AND-omics_type:"Proteomics"](https://www.omicsdi.org/search?q=Bacillus%20licheniformis-AND-omics_type%3A%22Proteomics%22). The page displays search results for "Bacillus licheniformis AND omics". The left sidebar contains filters for "Organisms" (listing Bacillus licheniformis (2), Corylus avellana (1), Schistosoma japonicum (1), Bos taurus (1), Saccharomyces cerevisiae (1), and Homo sapiens (1)), "Repository" (listing Pride (3)), and "Tissue". The main content area shows two results. The first result is "Bacillus licheniformis exponential phase" with a green gear icon indicating 45 datasets. It describes the cytoplasmic, membrane, and exo-proteome of cells growing in minimal medium. The second result is "Gastrointestinal digestion of hazelnut allergens on molecular level" with a green gear icon indicating 25 datasets. It describes the elucidation of degradation kinetics and resistant immunoactive peptides using mass spectrometry. The page includes navigation links (Home, Browse, Submit Data, API, Database, Help, Login) and a search bar.

Search results for **Proteomics** (3)

Organisms

Find your Organisms

- Bacillus licheniformis (2)
- Corylus avellana (1)
- Schistosoma japonicum (1)
- Bos taurus (1)
- Saccharomyces cerevisiae (1)
- Homo sapiens (1)

Repository

Find your Repository

- Pride (3)

Tissue

Sort by: Relevance Page size 10

« Previous 1 Next »

Bacillus licheniformis exponential phase 45

Cytoplasmic, membrane and exo-proteome of cells growing in minimal medium

ORGANISM(S): Bacillus licheniformis

2015-04-21 | PXD000791 | Pride

- cytoplasmic proteome
- membrane proteome
- Bacillus licheniformis
- exoproteome

Cite

Gastrointestinal digestion of hazelnut allergens on molecular level: 25

Elucidation of degradation kinetics and resistant immunoactive peptides using mass spectrometry

Allergy to hazelnut seeds ranks among the most prevalent food allergies in Europe. The aim of this study was to elucidate the gastrointestinal digestion of hazelnut

Closing thoughts

- Spreadsheets are fine for 2D tables of numbers, but knowing about databases empowers you to handle more complexity.
- Supporting your experimental data with those published to repositories makes your manuscript stronger.