



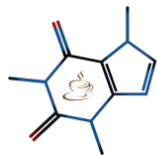
Challenges in Building NoSQL Databases for Natural Products Research

Maria Sorokina, Christoph Steinbeck

Friedrich-Schiller University Jena, Germany



About me



Cheminformatics and Computational Metabolomics
Friedrich-Schiller-University, Jena, Germany



FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA



Chem- and bioinformatician

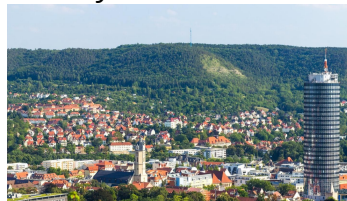
Did my studies in bioinformatics and statistics at Paris-Saclay University, France

Doctoral work at the Genoscope, Evry (Paris region), France on metabolic networks representation and finding new metabolic pathways

Now senior postdoctoral researcher at the Friedrich-Schiller University, in Jena, Germany:

- Natural Products cheminformatics (databases)
- Research Data Management for the ChemBioSys CRC
- Omics for marine diatoms

Steinbeck Lab:
<https://cheminf.uni-jena.de>

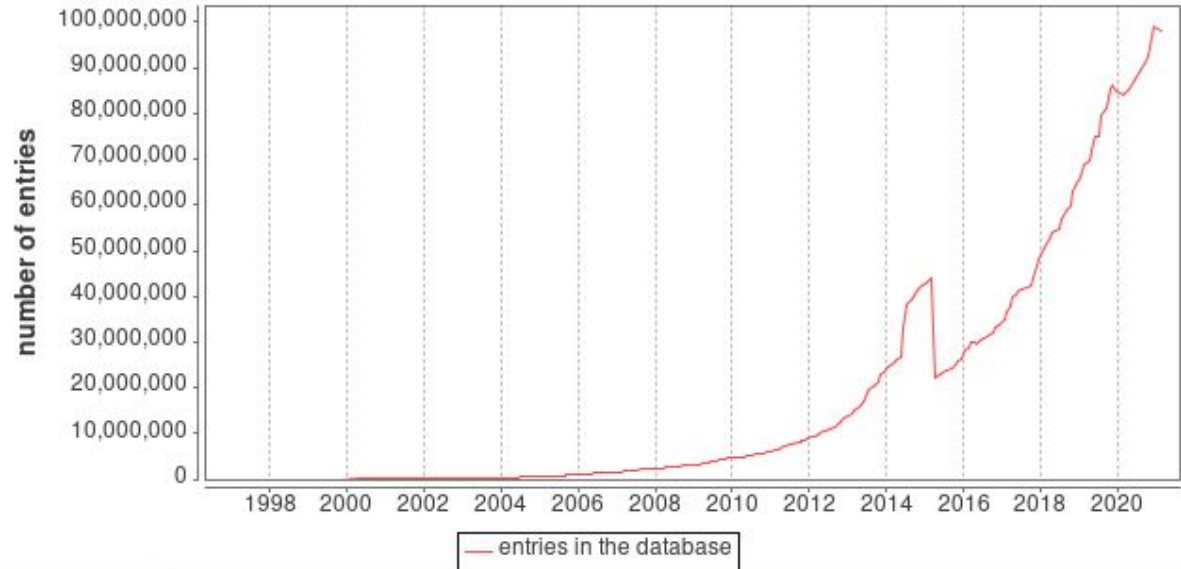


Why data storage is a challenge

Constantly increasing amount of complex data

-----> increasing need to store **increasingly complex** and **increasingly connected** data efficiently

**Number of entries in
UniProt/TrEMBL**

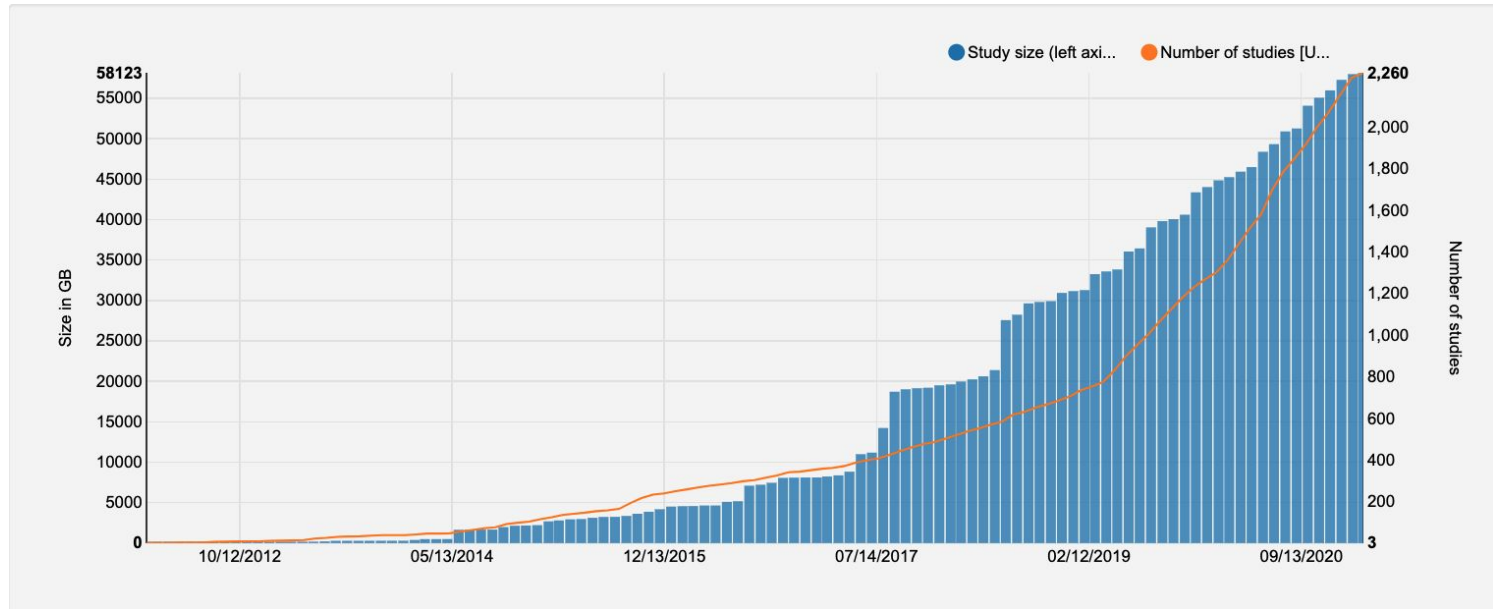


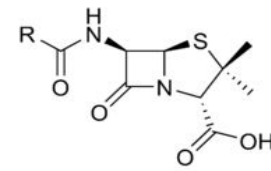
Why data storage is a challenge

Constantly increasing amount of complex data

-----> increasing need to store **increasingly complex** and **increasingly connected** data efficiently

Number of entries in MetaboLights



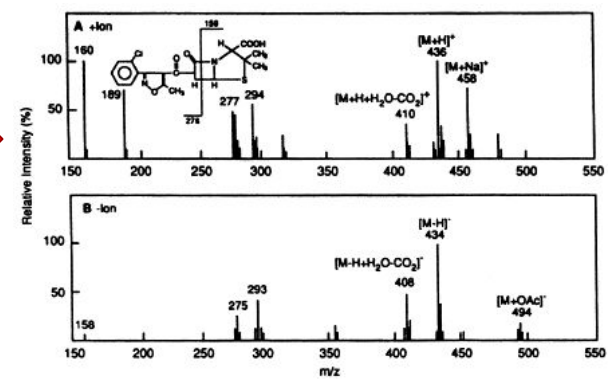
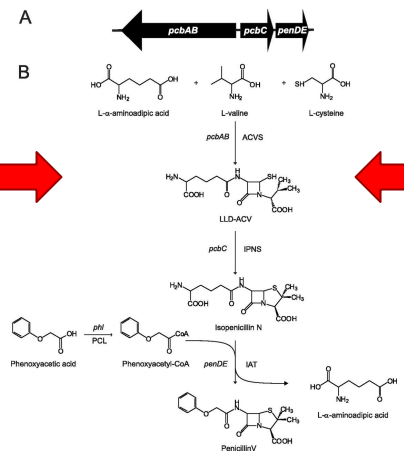
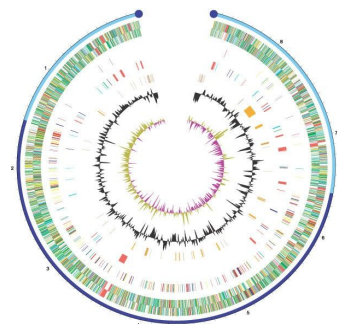


Penicillin

Why biological data storage is a challenge



Penicillium chrysogenum



Organism information & taxonomy

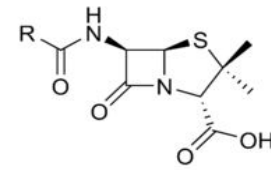
Assembled genome & gene functional annotations

Gene clusters & biosynthetic pathways

Spectral information

Sources: Applied and environmental microbiology, Vol 87 Issue 6, <https://aem.asm.org/content/78/19/7107/F1> [https://doi.org/10.1016/S0021-9673\(98\)00281-7](https://doi.org/10.1016/S0021-9673(98)00281-7)

Why biological data storage is a challenge

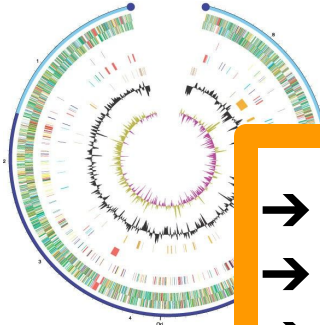


Penicillin

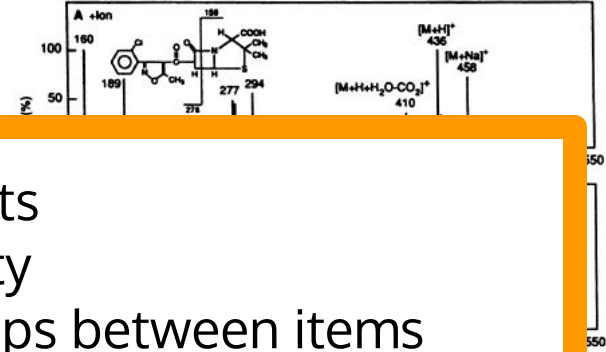
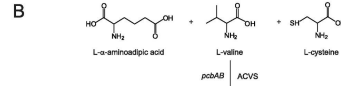


Penicillium chrysogenum

Organism information & taxonomy



Assembled genome & gene functional annotations



- ➔ Multiple data formats
- ➔ High data complexity
- ➔ Complex relationships between items
- ➔ Need an efficient way of modelling and storing data

Short introduction to database management systems

SQL: Structured Query Language - main query language for relational databases

Short introduction to database management systems

SQL: Structured Query Language - main query language for relational databases

Relational Database Model

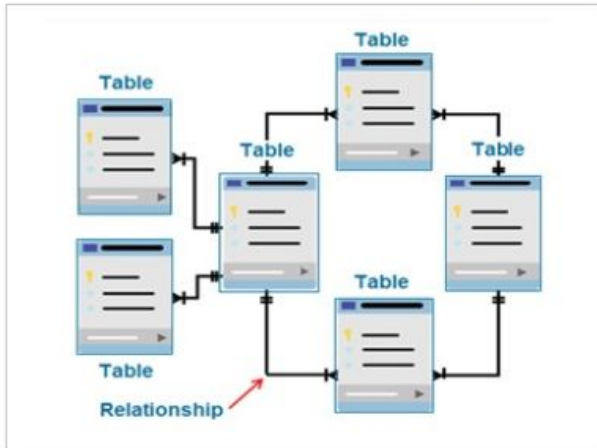


Table = Relation = Entity = Concept = Object

PK - Primary Key And FK - Foreign Key

www.learncomputerscienceonline.com

RDBMS - Relational Table Example

StuID	StuName	StuAge	StuClass	StuSection
1001	Alex	15	10	B
1002	Maria	14	11	A
1003	Maya	14	9	A
1004	Bob	16	11	C
1005	Newton	14	10	D
1006	Sanjay	15	10	B

Table / Entity / Relation

A Table Represents a Database Entity

Table Row is referred as Records Or Tuple

A Table Column Represents an Attribute

www.learncomputerscienceonline.com



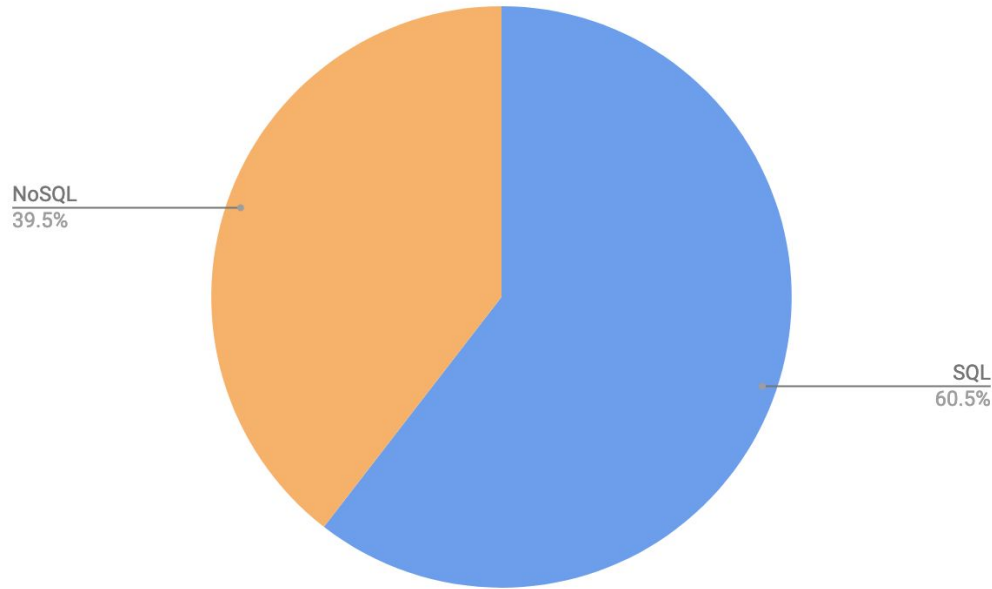
- create meaningful information by joining the tables
- joining tables allows to understand the *relationships* between the data, or how the tables connect
- good balance between flexibility and efficiency
- indexing

BUT

- table schema is expensive to change
- not terrific with parallel writing to the same table
- limited data formats and organization

Short introduction to database management systems

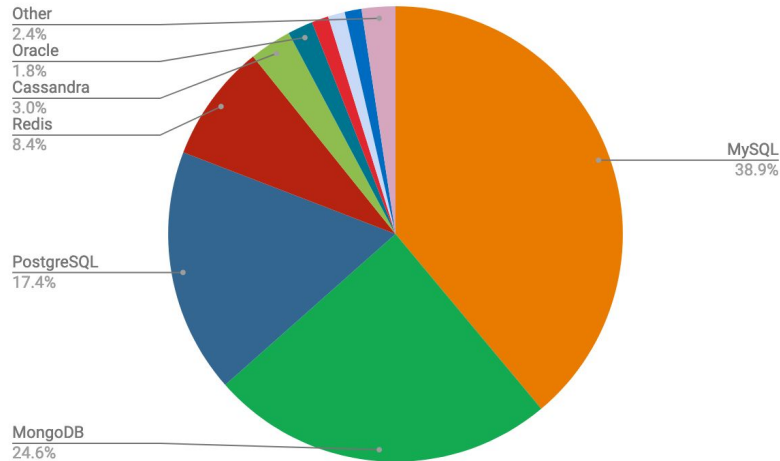
In 2019:



Short introduction to database management systems

SQL: Structured Query Language

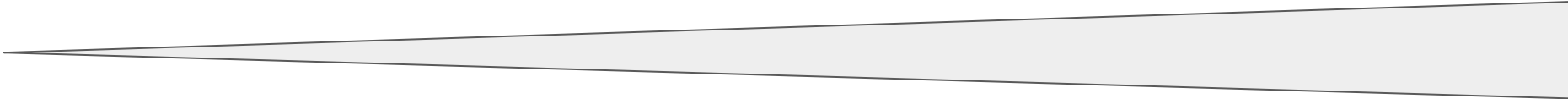
noSQL: “not only SQL” rather than “not SQL”



Short introduction to database management systems

SQL: Structured Query Language

noSQL: “not only SQL” rather than “not SQL”



Less relational

More relational

<p><u>Key-value DBs</u> Redis, Voldemort, Dynamo</p>	<p><u>Column-oriented DBs</u> (db-dependent, e.g. CQL) - CassandraDB, Google’s Big Table</p>	<p><u>Document DBs</u> (db-dependent, e.g. mongo query language) - MongoDB, CouchDB</p>	<p><u>Relational DBs:</u> (SQL) - MySQL, PostgreSQL, MariaDB, Oracle</p>	<p><u>Graph DBs</u> (db-dependent, e.g. Cypher) - Neo4j, OrientDB</p>
---	---	--	---	--

noSQL database:



mongoDB®

- Document database
 - Mostly open (they make money with the cloud for companies)
 - Mongo Query Language
 - Since 2018 main storage engine: WiredTiger - extremely effective for Big Data
 - Each databases is composed of “collections” (equivalent of tables)
 - Each collection is composed of “documents”
 - Documents are in JSON format
-
- Does NOT have in-built functions (yet) for chemical data (as it has for geo data)

noSQL database:



mongoDB®

A storage engine is a software module that a database management system uses to **Create, Read, Update and Delete** data from a database.

- Document database
 - Mostly open (they make money)
 - Mongo Query Language
 - Since 2018 main storage engine: WiredTiger - extremely effective for Big Data
 - Each database is composed of “collections” (equivalent of tables)
 - Each collection is composed of “documents”
 - Documents are in JSON format
-
- Does NOT have in-built functions (yet) for chemical data (as it has for geo data)

noSQL database:



mongoDB®

- Document database
 - Mostly open (they make money with the cloud for companies)
 - Mongo Query Language
 - Since 2018 main storage engine: WiredTiger - extremely effective for Big Data
 - Each databases is composed of “collections” (equivalent of tables)
 - Each collection is composed of “documents”
 - Documents are in JSON format
-
- Does NOT have in-built functions (yet) for chemical data (as it has for geo data)

Modelling data for NP: the COCONUT example

COCONUT: ColleCtion Of uNique natUral productS

<https://coconut.naturalproducts.net>

In the database, 3 main collections:

- sourceNaturalProduct
- uniqueNaturalProduct
- fragment

Contains data from **56** different public data sources (version from october 2020) and **401,624** unique “flat” molecules



Modelling data for NP: the COCONUT example

Why MongoDB?

- Structure flexibility: easy to add more complex fields if needed
- Good with big complex data
- Fuzzy text search auto-enabled
- Multiple in-build search functions (like selecting on-bits in fingerprints, which facilitates structure search)

Modelling data for NP: the COCONUT example

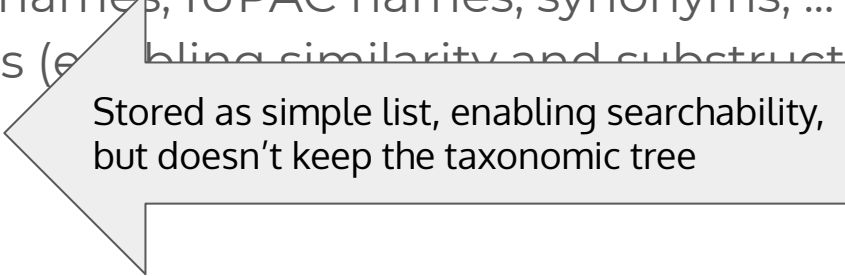
For each unique flat NP, a number of annotations is stored:

- Available 3D structures
- Molecular descriptors
- Traditional names, IUPAC names, synonyms, ...
- Fingerprints (enabling similarity and substructure search)
- Taxonomy
- Literature
- Geography
- Chemical classification
- Source databases (and x-references to them)

Modelling data for NP: the COCONUT example

For each unique flat NP, a number of annotations is stored:

- Available 3D structures
- Molecular descriptors
- Traditional names, IUPAC names, synonyms, ...
- Fingerprints (enabling similarity and substructure search)
- Taxonomy
- Literature
- Geography
- Chemical classification
- Source databases (and x-references to them)



Stored as simple list, enabling searchability,
but doesn't keep the taxonomic tree

Modelling data for NP: the LOTUS example

LOTUS (natural prOducTs occUrrences databaSe) - extra-curated NP database

<https://lotus.naturalproducts.net>

Concept:

each NP is associated with at least one literature reference and one organism

- COCONUTs' simple list of taxonomies and list of references is not good enough
- more complex data model needed

Modelling data for NP: the LOTUS example

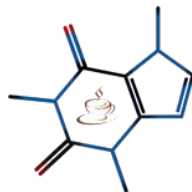
LotusUniqueNaturalProduct

- id
- structure representations
- taxonomyObject: {
 [{doi: '[{ taxonomyDatabase1:[{organism1,organism2,...}], {..}]}]
- ...

Take home messages

- ★ Natural products data quantity and complexity is constantly increasing
- ★ Proper data organisation is essential at the earliest stages
- ★ MongoDB and other document-based DBMs are extremely efficient and allow without too much damage a late data model replanification

Acknowledgements

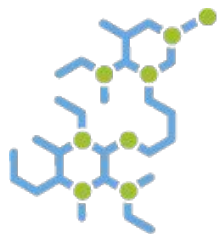


Cheminformatics and Computational Metabolomics

Friedrich-Schiller-University, Jena, Germany

Chris Steinbeck and the wonderful Caffeine group (cheminf.uni-jena.de)

ChemBioSys CRC



ChemBioSys

COLLABORATIVE RESEARCH CENTER 1127
CHEMICAL MEDIATORS IN COMPLEX BIOSYSTEMS



And the organisers of this workshop!

1. COCONUT

Stereochemistry in COCONUT:

423706 molecules with no the stereochemistry or where it was removed

- 50% only one stereocenter (but can be more in nature!)
- 23% more than one stereocenter (from the same or different DBs)
- 15.6% - truly no stereocenters
- 11.4% - have at least one stereocenter but info missing in the source database

Missing stereochemistry is a problem, as it has it's importance for molecular function.
But this can be solved mainly experimentally...