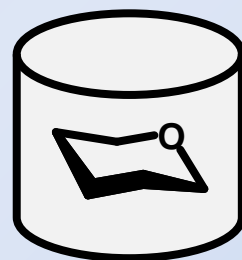




**Philip Toukach**

# **Carbohydrate Databases**



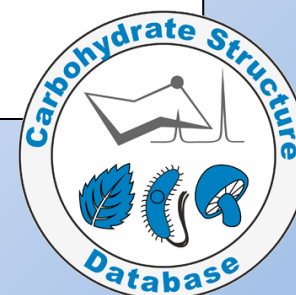
<http://toukach.ru/rus/glyco-db.htm>

# Carbohydrate databases

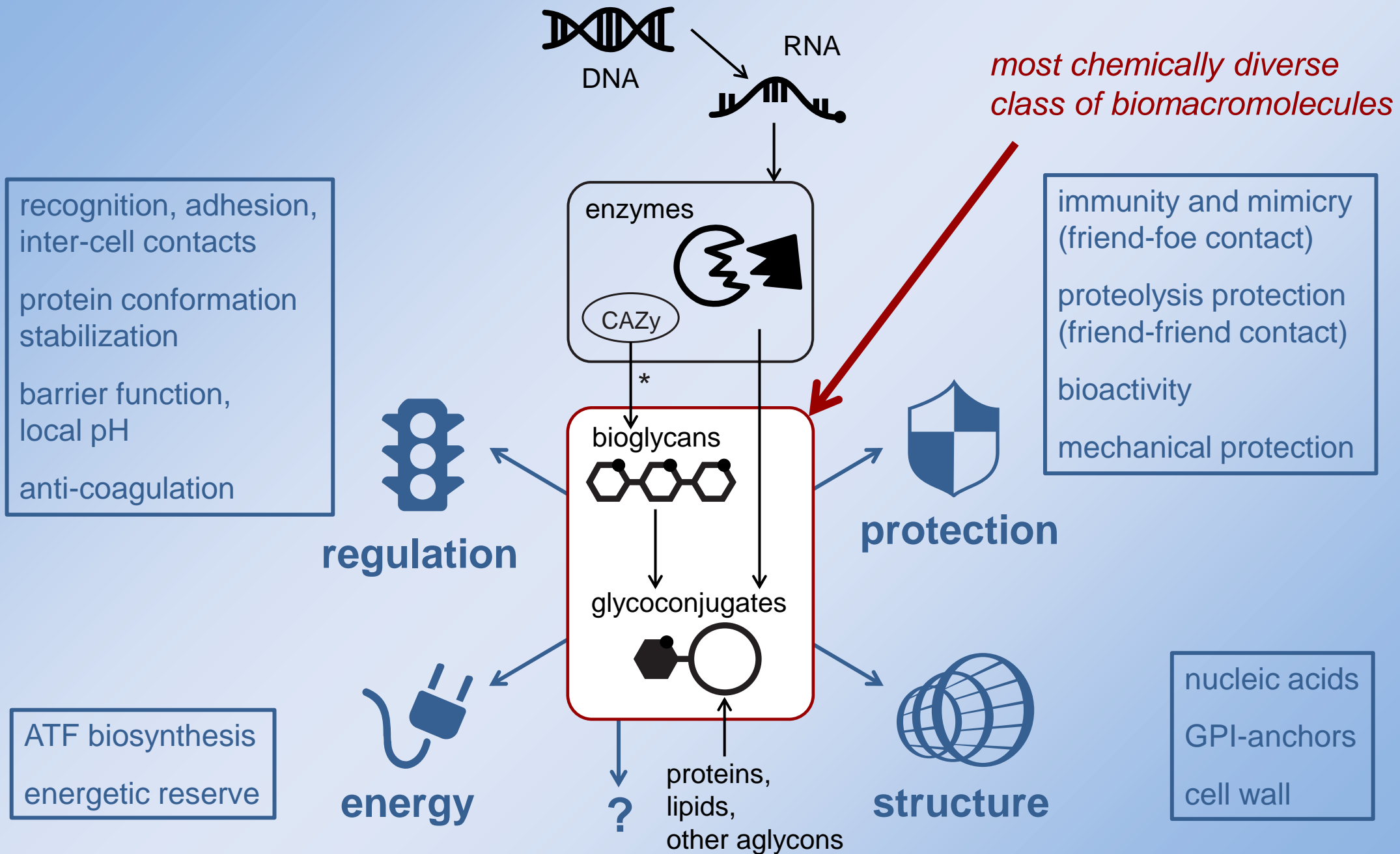
- **What are glycans and what are databases?**
- **What do we expect from them?**
- **Why do not we get it?**

• **What is an ideal database?**

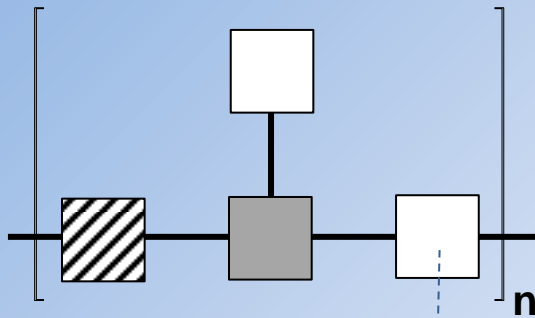
• **What do we have already?**



# Glycans in cells

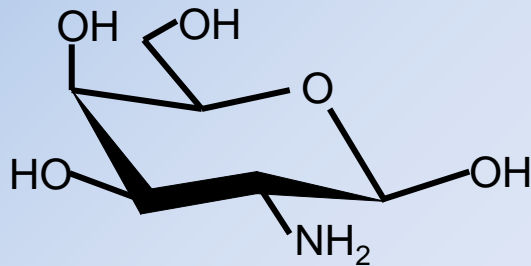


# Glycan structure

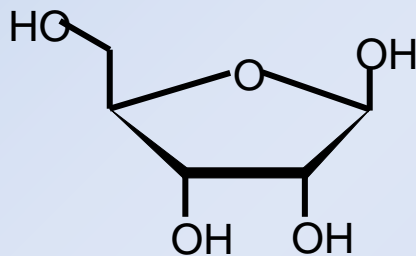


## Complete structure

- monomeric residues, incl. non-sugars
- topology & sequence
- linkage positions
- side chain stoichiometry
- polymer size and frame positioning



*aldo-pyranose* example ( $\beta$ -D-GalpN)



*aldo-furanose* example ( $\beta$ -D-Ribf)

## Residue structure

- carbon skeleton size (4-10)
- stereo pattern (monomer identity)
- ring form (*p/f/a*, *aldo/keto*)
- anomeric configuration ( $\alpha/\beta$ )
- absolute configuration
- modifications (-NH<sub>2</sub>, -COOH, deoxy)

# Glycomics in biology

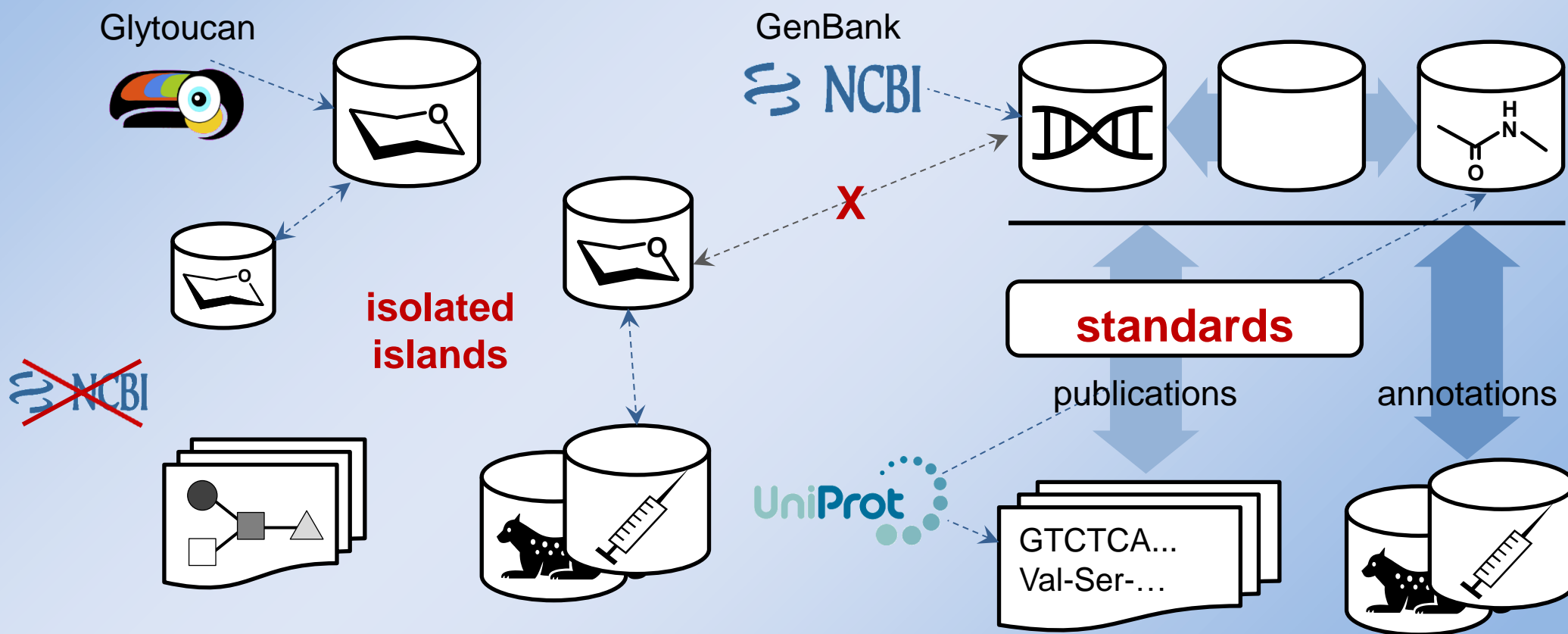
- Structure, diversity, conformation of carbohydrates
- Taxonomy and classification of microorganisms
- Glyco-epitopes and immunospecificity of strains
- Explanation of antigen-antibody interaction
- Carbohydrate vaccines and immunostimulators
- Correlation of microorganism bioactivity to its glycome
- Carbohydrate biosynthesis and turn-over

# Glycomics vs. genomics, proteomics

6

*as compared to other -omics:*

- similar information scope (>100 000 known structures)
- greater chemical diversity
- poorer IT involvement (databases, services)
- less standardized



# Why do we need a database?

7

- **Easy access to knowledge and research automation**

Which natural structures look like a model of interest? Which of their fragments are specific to a genus of interest? Where were they published, and in association to which taxa, diseases, and organs? Which enzymes synthesize them and how was it proven? Which glyco-epitopes induce immune response in an organism of interest?

- **Simulation of molecular properties**

Molecular geometry and energy, MS and NMR spectra, bioactivity, ...

- **Structure prediction from experimental data**

- **Prediction of taxon properties**

Glycome-based clustering, search for similarities and differences in taxa, chemotaxonomic classification

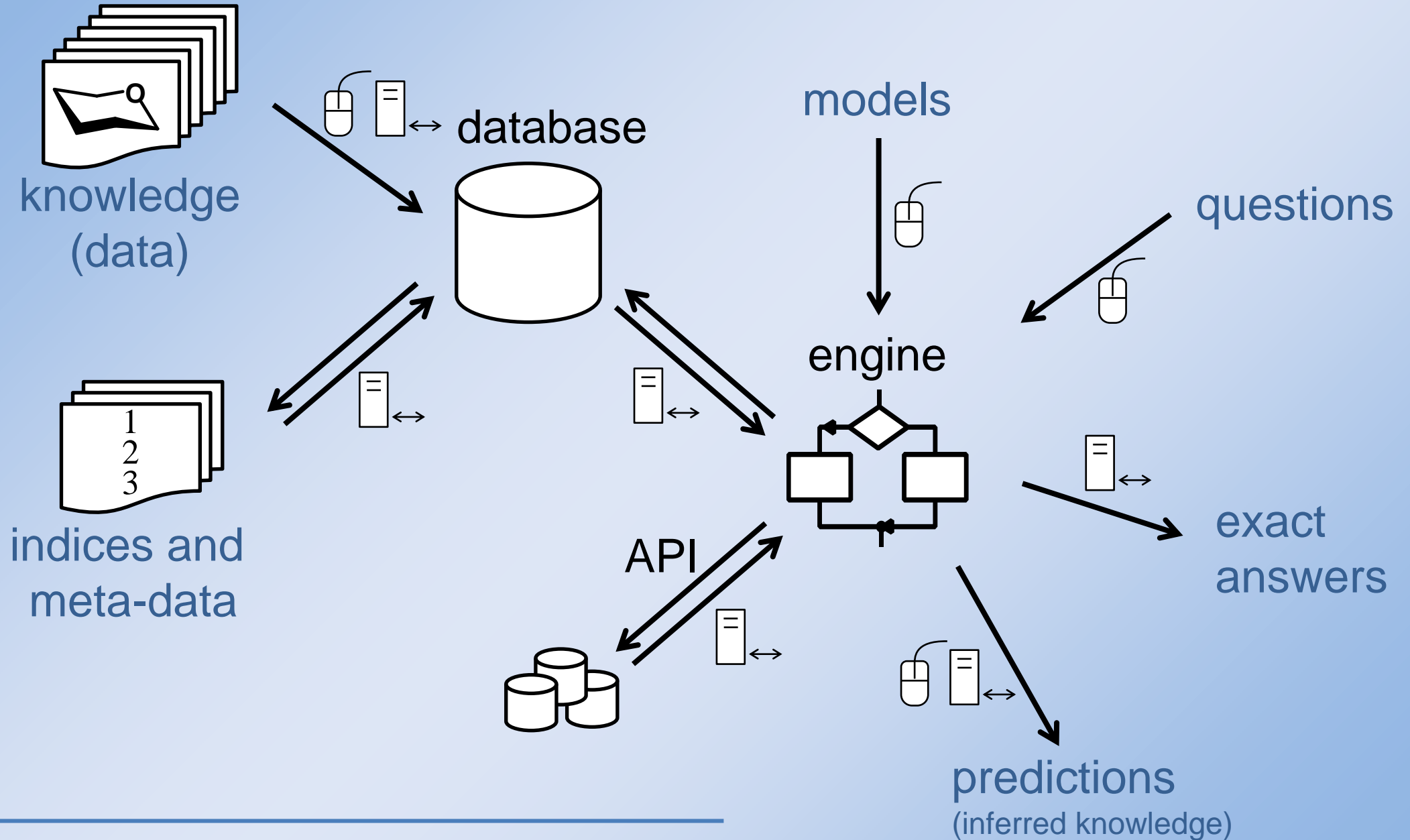
- **Molecule identification and visualization** (in publications as well)

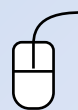
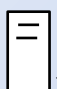
# Glycoinformatic challenges

- Diversity and heterogeneity of objects
- Ambiguous structural description
- Difficult input and visualization of complex structures
- Project isolation and lack of standards
- Incompleteness of data and poor data quality
- Resource-greedy algorithms
- Lack of systematic view from users and developers  
(no commonly-agreed services; incompatibility of initiatives)



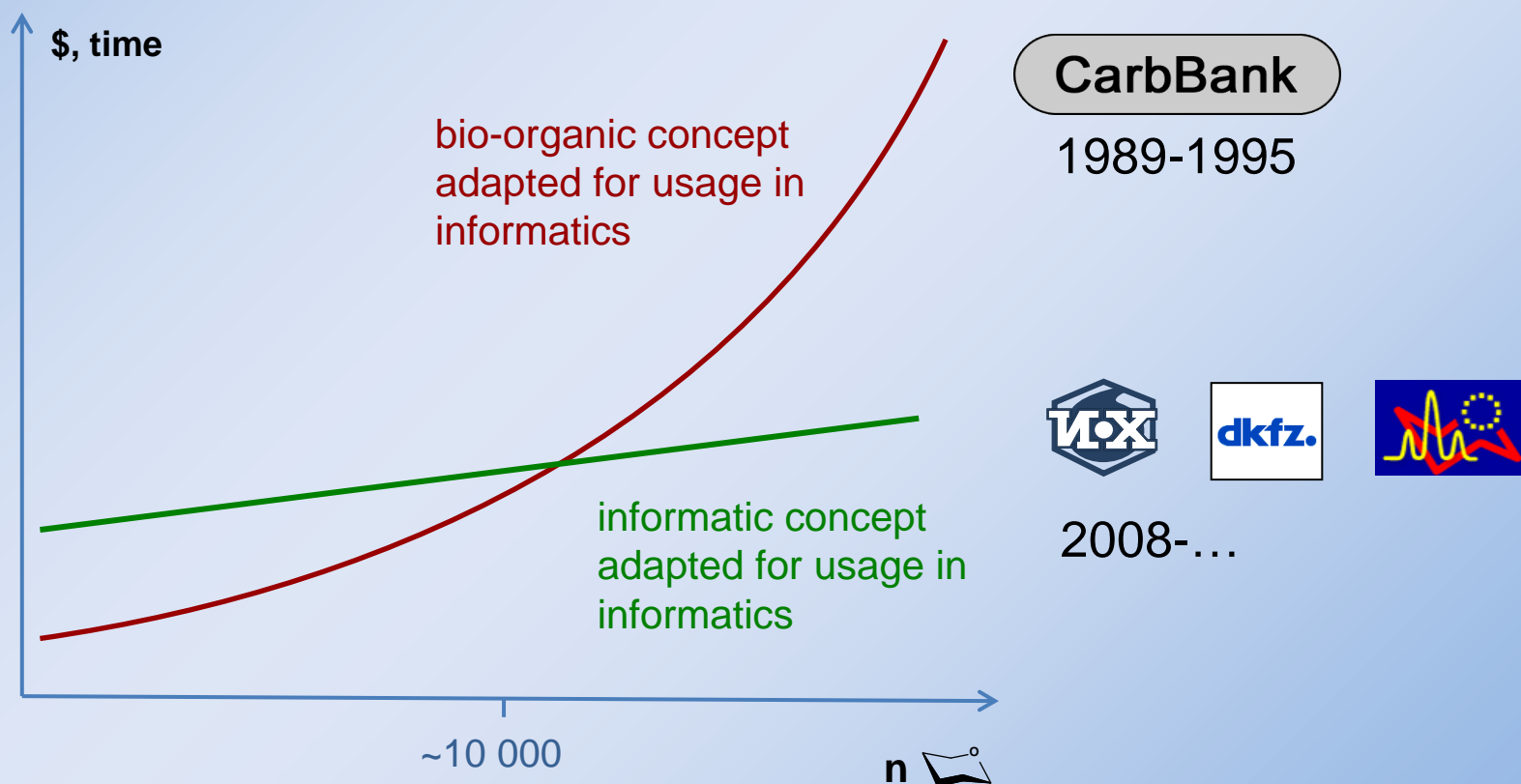
# What is a database?



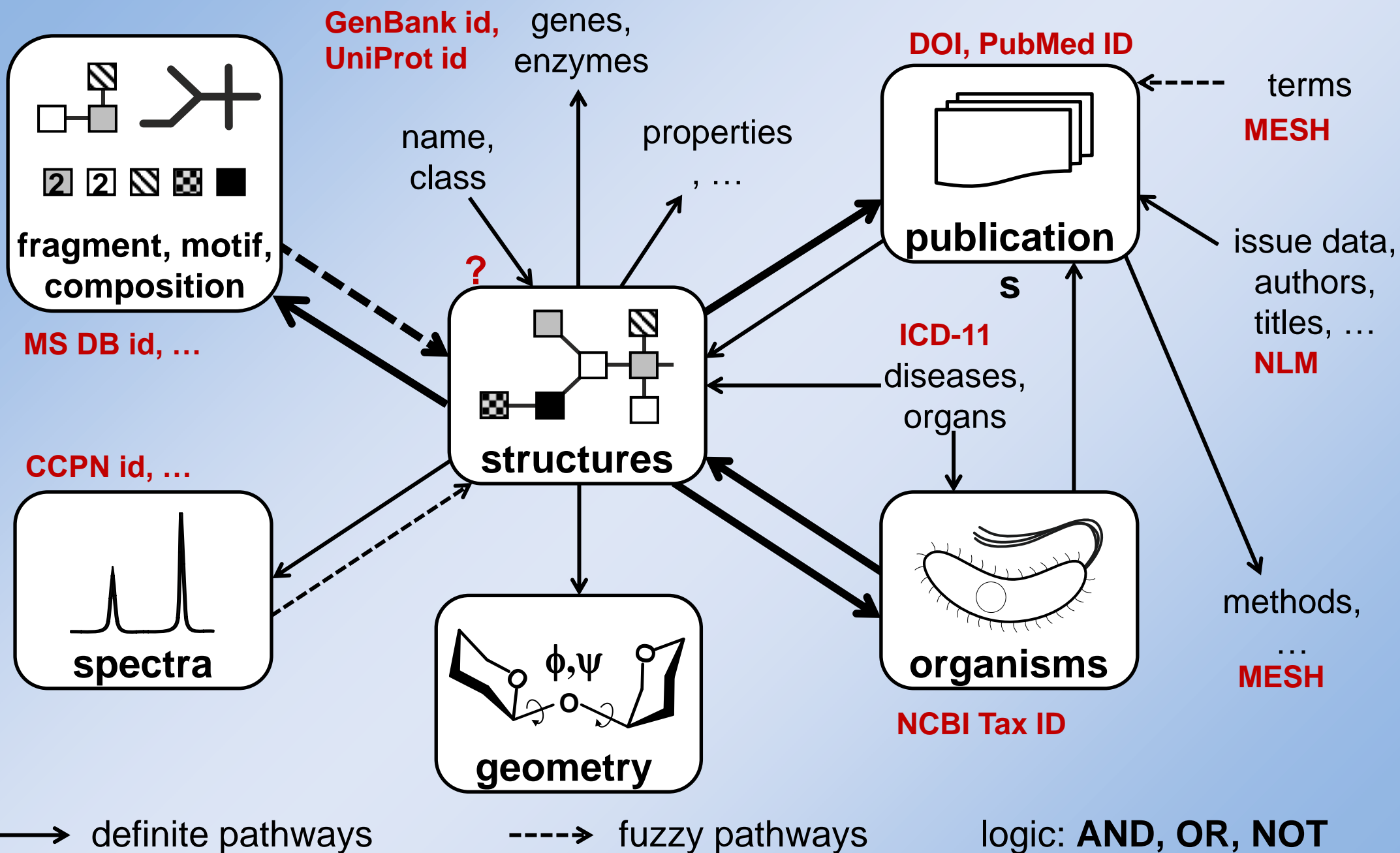
 human-manned       ↔ automated

# Development approaches

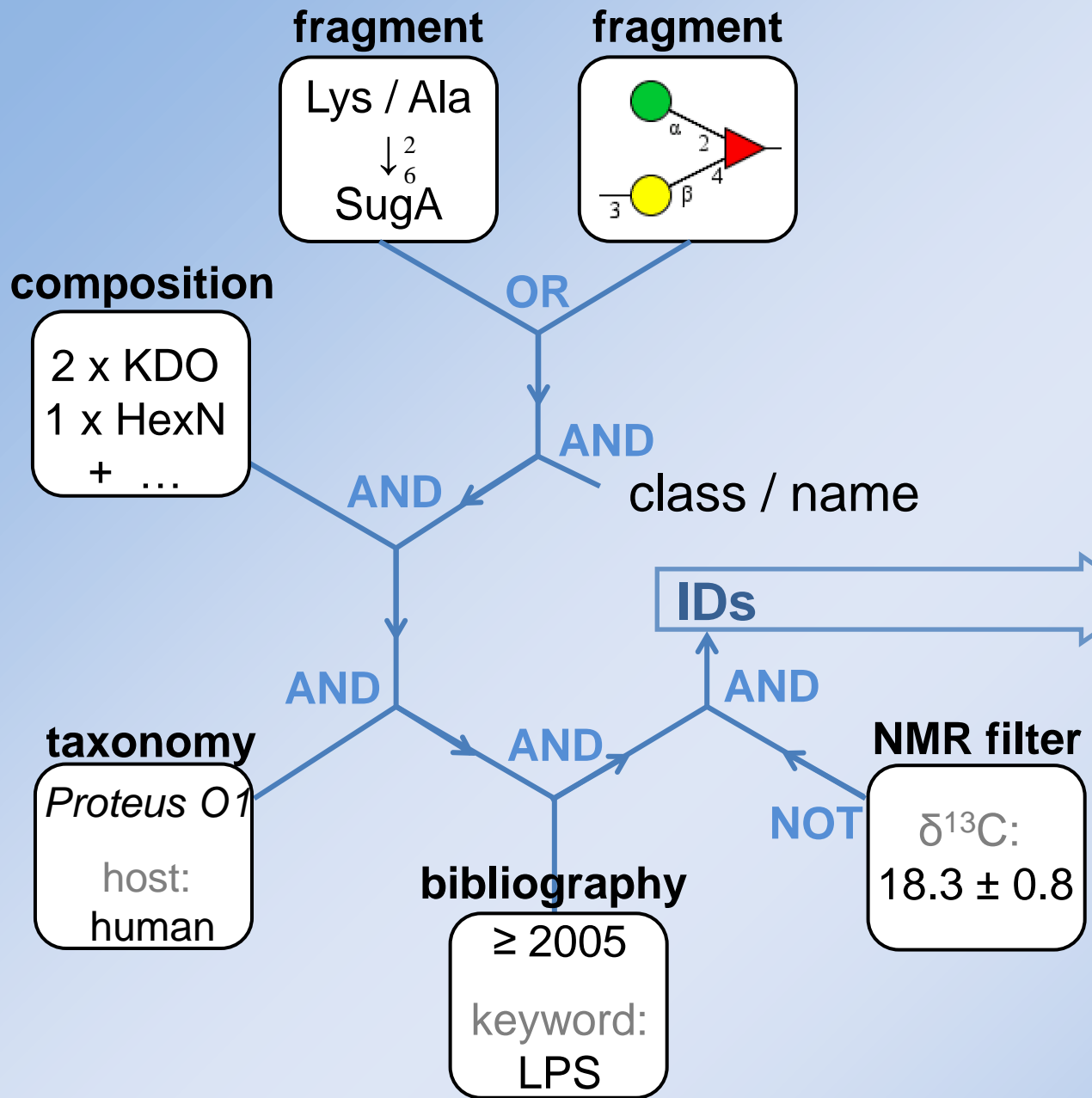
- A database and a platform should follow the rules of informatics
- These rules should be adapted and specificated for carbohydrates



# Typical queries



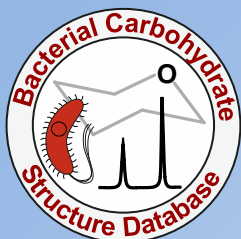
# CSDb: complex query



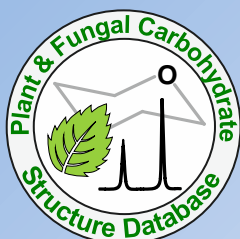
Data arranged by compound, publication, organism, etc.

The screenshot shows the search results for compound ID 10502. It displays two chemically distinct structures for selection. The first structure is a repeating unit of a polymer chemical repeating unit, identified as O-polysaccharide, O-antigen. The second structure is a ball-and-stick model of the same compound. The search results also include a list of publications, including one from 2012 (DOI: 10.1016/j.carres.2012.04.006) and another from 2013 (DOI: 10.1016/j.carres.2013.04.006). The NMR spectrum for the compound is also shown, with a peak at 18.3 ± 0.8 ppm. The spectrum also has 2 signals at unknown positions (not plotted).

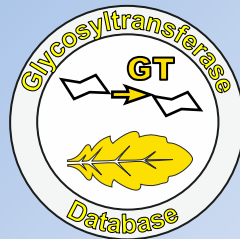
# CSDB: Carbohydrate Structure Database<sup>13</sup>



since  
2005



since  
2012

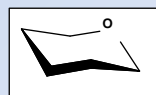


since  
2017

## Database set + platform for services



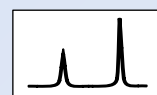
Zelinsky Institute  
Moscow, Russia



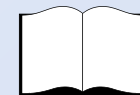
25K



13K



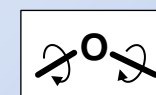
14K



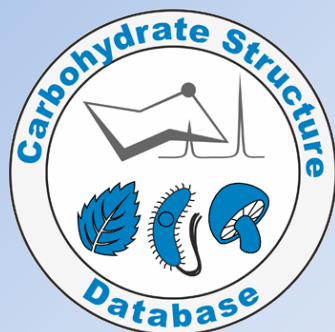
10K



2K






3K



CSDB



- regular updates
- extensible architecture
- data analysis tools
- curated content (15% = Carbbank, 85% = literature)
- complete coverage (bacteria and fungi)
- integration with other DBs


# Structural databases




**CarbBank** 23   complete to 1995    architecture, % errors








**GlycomeDB** 39   meta-repository  poor annotations   no aglycons





**GLYTOUCAN** 99  meta-repository  poor annotations   no aglycons


**CFG glycan**   mammals, > 6






SweetDB, SugaBase 

**GIYCO-SCIENCES.DE**    25 / 19


**GLYCAN**        11

**Eurocarb DB**    architecture  model only

BCSDB  PFCSDB **13 / 5** (bacteria, archaea) **8 / 3** (fungi, plants)

**Carbohydrate Structure Database**   complete on prokaryots  curated  

**GlycoSuite**      mammals+... **10 / 1**  complete to 2005

**JCGGDB**  **> 70**  database collection  poor annotations

**UniCarbKB**   **4 / 1**     curated

**nibrT**  **0.7 O- & N-**    

**Glycoconjugate Data Bank**  **44**   

**EcoDAB** **0.2 E. coli**  

**GlycoBase**  **0.3 animals**  

# Special content



human glyco genes



~0.2

N- & O-glycan MS<sup>2,3,4</sup>



~0.2

glycochemical reactions



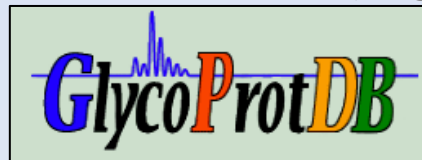
~3.0 (4.4 str)

conjugates & aglycones



>70

N-glycoproteins  
*C. elegans* + mice



~2.5

protocols for  
synthesis & analysis



~0.2 (0.5 sub)

glyco-epitopes  
& antibodies



~0.2 (0.6 ABs)

GlyTOUcan  
(id repository)



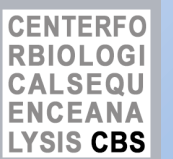
~99

binding to pathogens



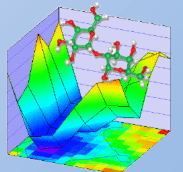
~0.9

O-glycBase,  
O- & C-glycoproteins



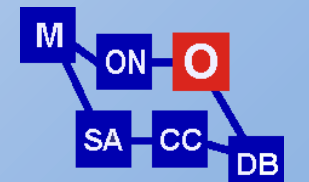
~0.2

GlycoMaps,  
computed conformational maps



~2.6

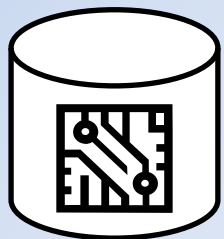
MSDB  
monosaccharides & notations



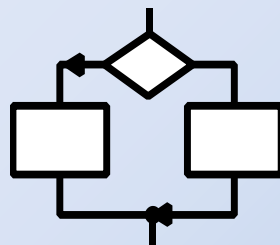
~0.8

# Evaluation criteria

- **Completeness** (+ chosen domain)
- **Data quality** (% of errors, human-readability)
- **Functionality** (data and index types, query processing)
- **Integration** (supported data formats, import & export, API, RDF)
- **Interface** (user-friendliness, stability, performance)



architecture



control scripts

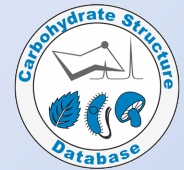


data



# Architecture, functionality

17



- Relational database
- Indexation + standard indices (DOI, TaxID, ICD-11, PMID, Genbank, ...)
- Structures, taxonomy, bibliography (different entries and data types)
- Human-readable dump (organization of data upload)
- Controlled term vocabularies (monomers; MSDB)
- ~~Free text~~ ↗
- Connection table

## minimum

structure,  
taxonomy,  
bibliography,  
cross-DB references

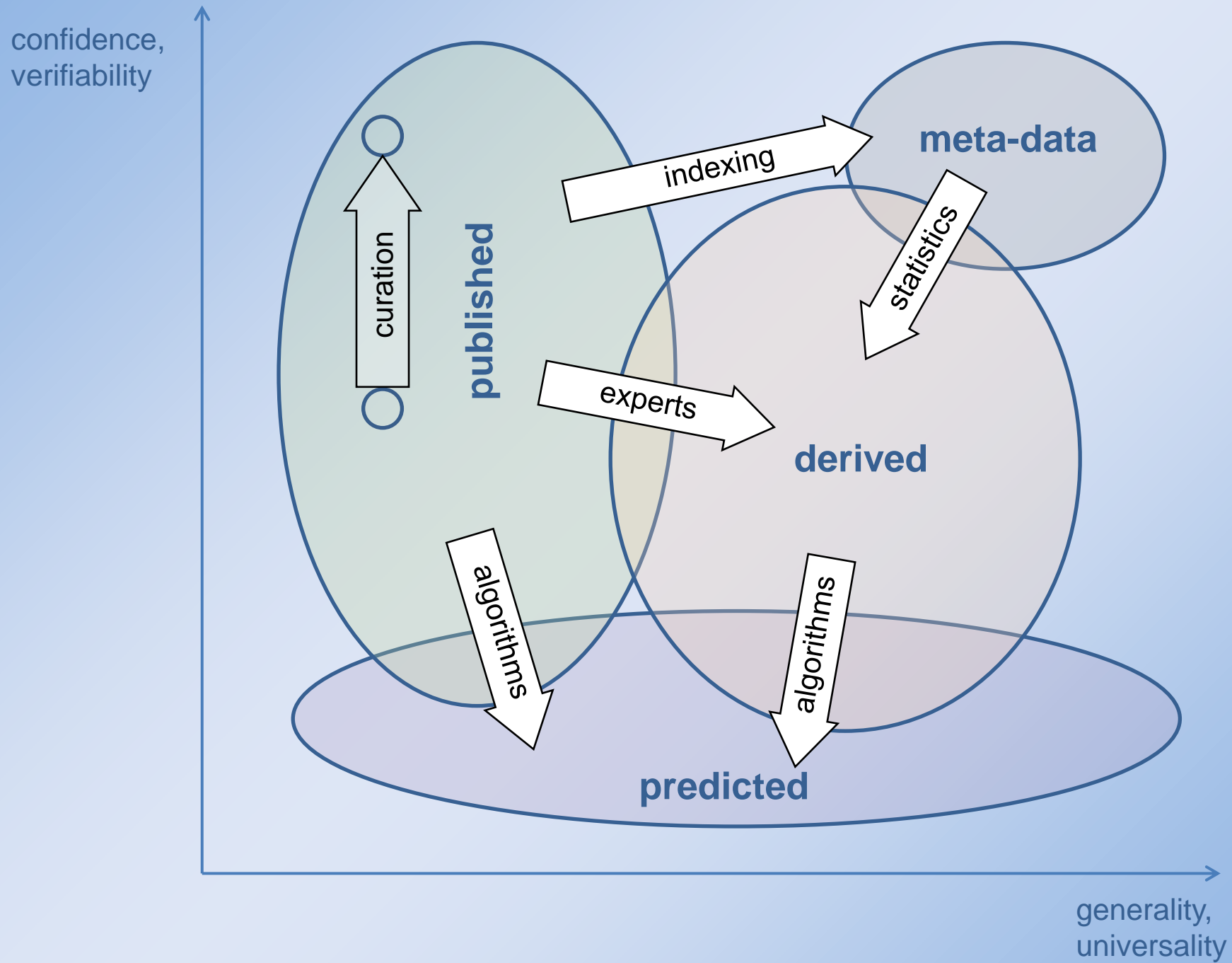
## optional

trivial names,  
NMR and MS spectra,  
spectroscopic conditions,  
bio-activity,  
genes, enzymes,  
conformation

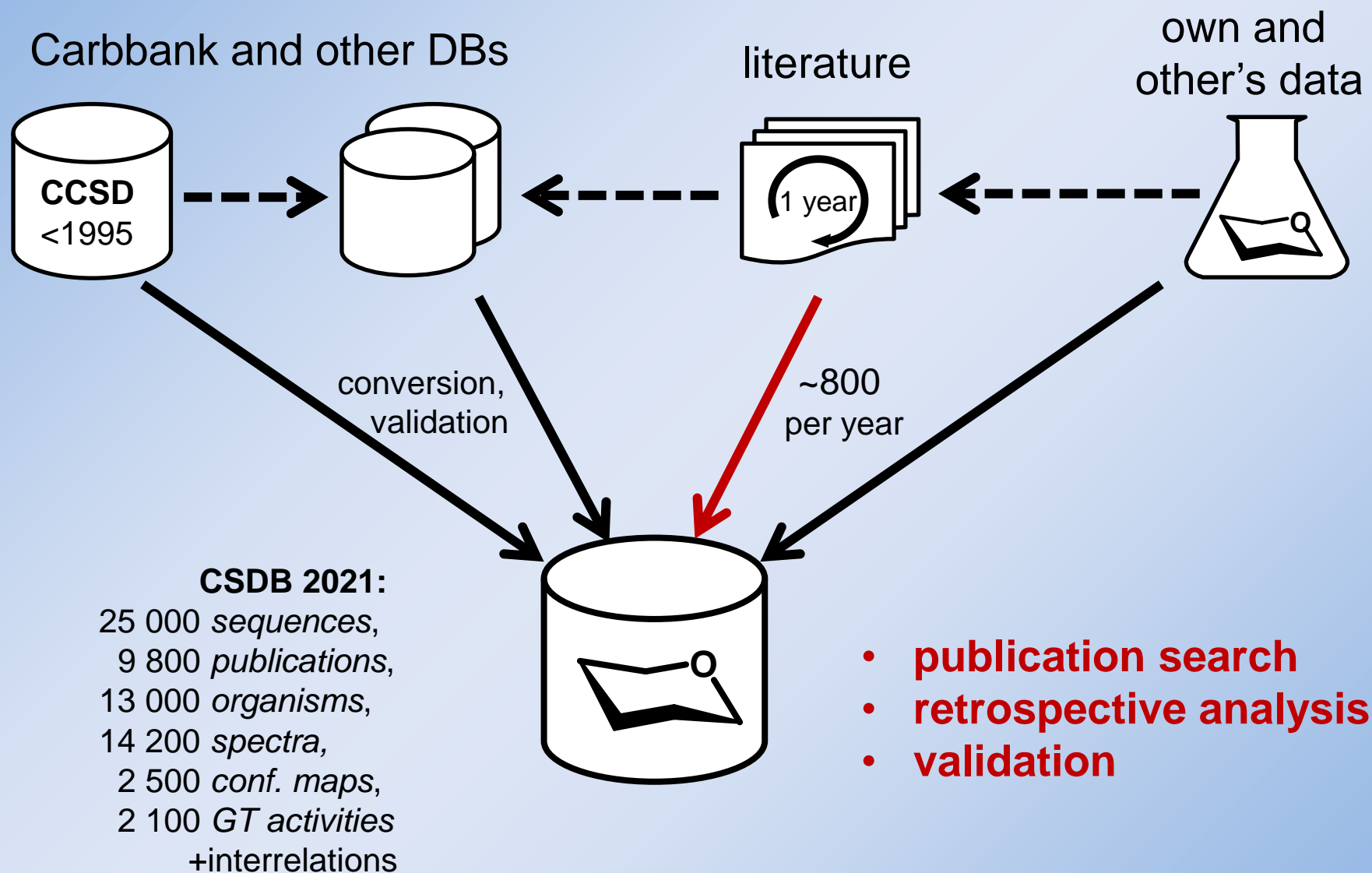
diseases,  
organs, tissues,  
genotype, life stage

keywords,  
abstracts,  
affiliations

# Data levels



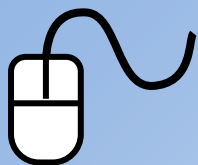
# Data sources



**complete coverage:**

negative search result = still valuable scientific information

# Data quality



from annotators

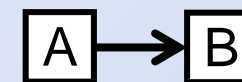


from databases

found &  
corrected



from publications



in scripts

## errors, inconsistency

### correctable

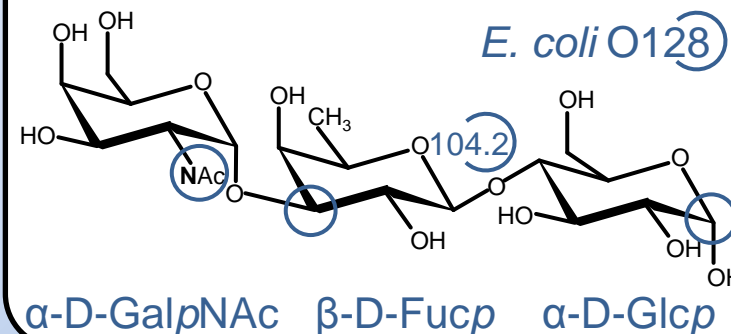
2dGlc → araHex,  
α-Rib-ol → Rib-ol,  
D-Kdo → Kdo,  
1-methyl → 1-Me,  
n.m.r. → NMR,  
taxid 583 → Proteus,  
...

### detectable

Glc(1-2)GlcN,  
anhydro-Kdo,  
D-manHep,  
Galp5N,  
Ac(1-2)[Glc(1-2)]Gal,  
*Escherichia sapiens*,  
*Dev Food Sci* 2012,  
#Ac : 23ppm, 65 ppm,  
D-Gcl, ...

### undetectable

*E. coli* O127:  
aDGalpN(1-4)bDFucp(1-4)bDGlc  
Fuc C1 103.2 ppm



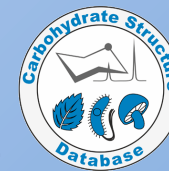
CarbBank



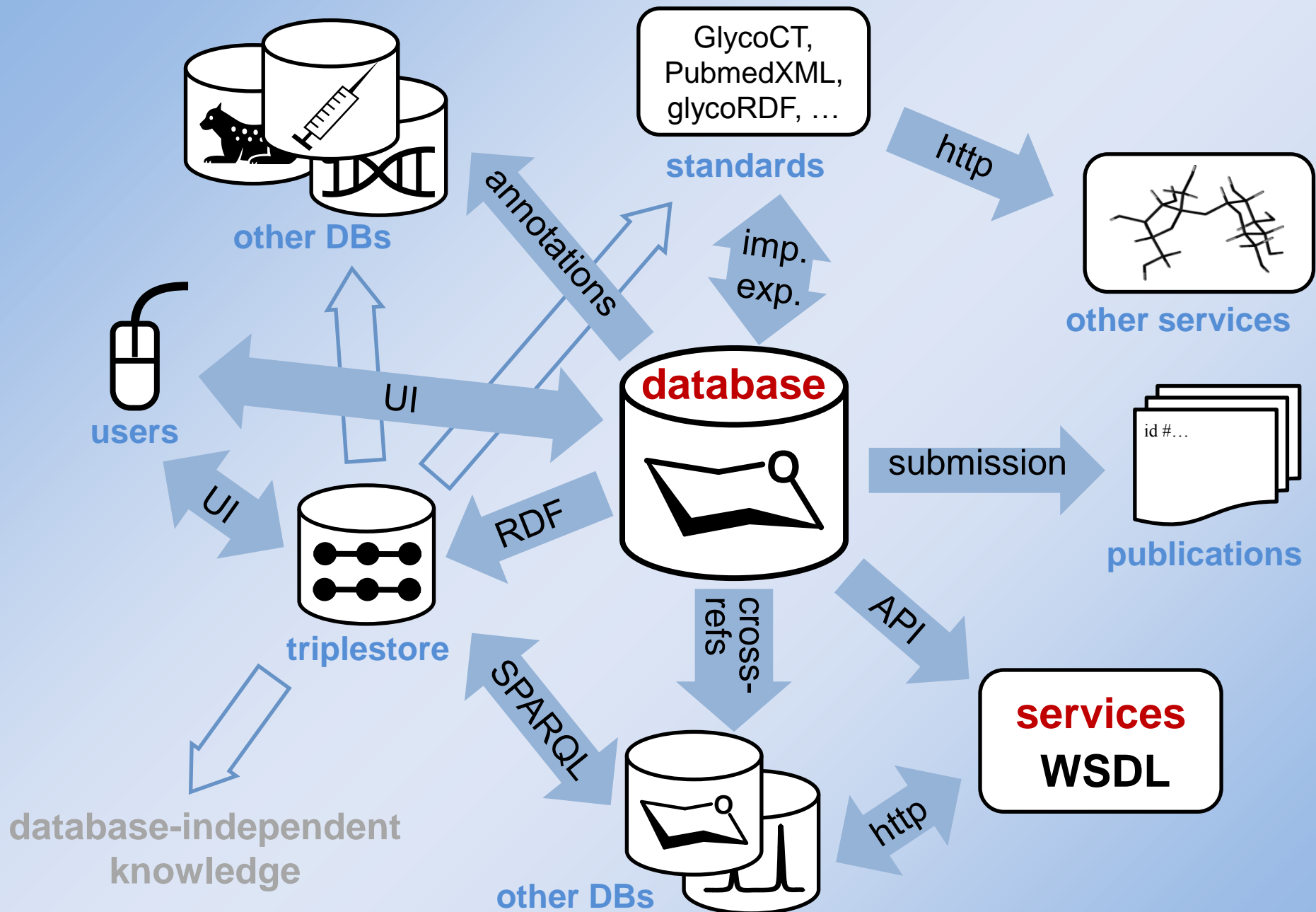
>50% (incorrect, missing, falsely present  
structures, strains, annotations)

↳ other DBs










😊 <10%



# Ideal integration



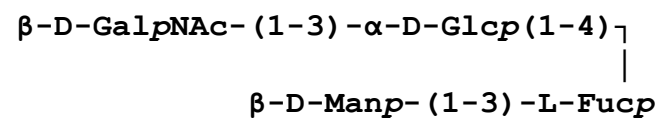
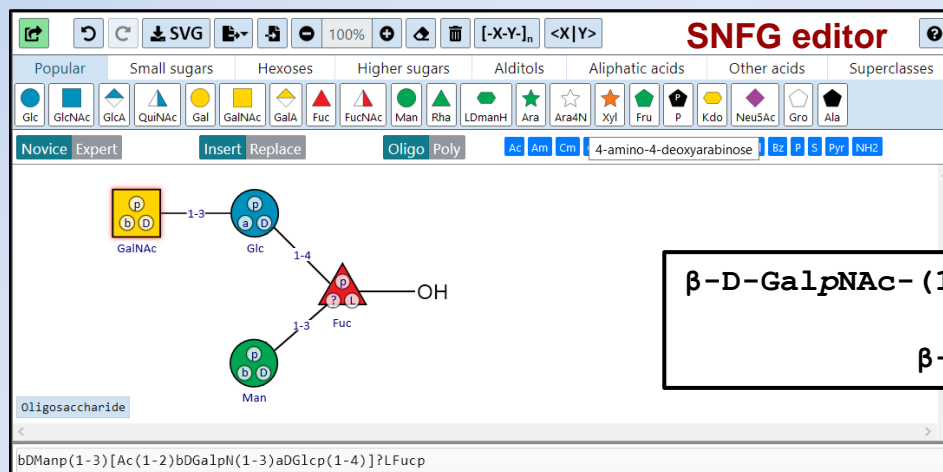
# Interfaces

- 
 ↔ Data conversion ↔ **other formats**
-  ↔ Automated web-services (WSDL)
- 
 ↔ Import, export
-  Documentation, user's HELP
-  User-friendliness, performance
-  Links to **other projects** (queries, indices, data)
-  Structure **input** & output

SNFG,  
WURCS,  
GlycoCT,  
SMILES,  
MOL, PDB,  
Glydell,  
LinUCS,  
Sweet-DB  
GLYCAM,  
GlycoRDF,  
DCI XML, ...

NCBI PubMed,  
NCBI Taxonomy,  
DOI,  
Uniprot, Genbank,  
Glycosciences.DE,  
MonosaccharideDB,  
Glytoucan, ICD-11

wizard,  
graphic builder,  
library,  
CSDB Linear,  
GlycoCT



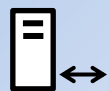
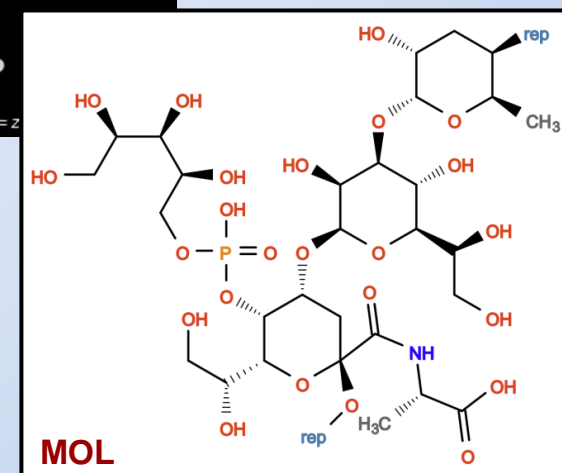
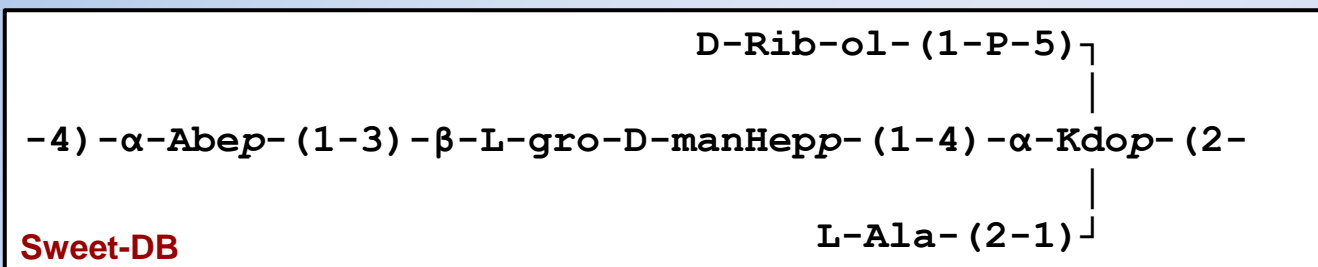
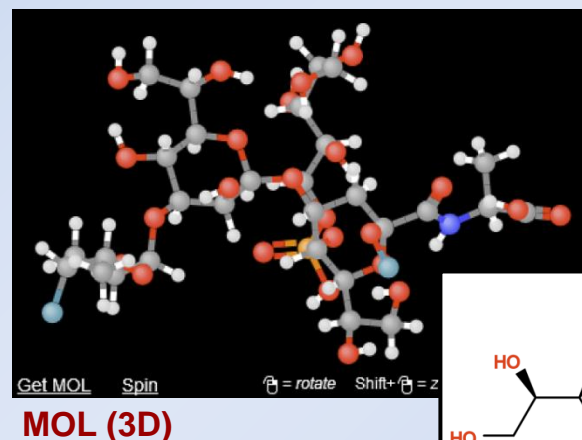
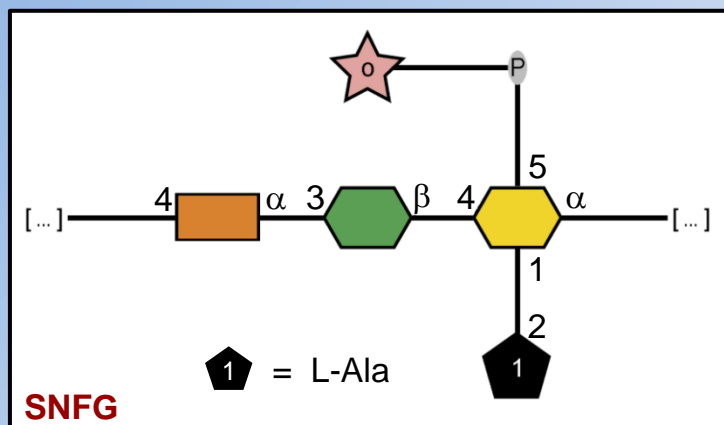
**bDManp(1-3)[Ac(1-2)DGalpN(1-?)aDG1cp(1-4)]?LFucp**

CSDB

# Structure output



## Visualization in human-readable formats:



## Export in machine-readable formats:

```
[*]O[C@]1(C(=O)N[C@@H](C)C(=O)O)[C@@H](O[C@@H]2O[C@H]([C@@H](O)CO)[C@@H](O)[C@H](O[C@H]3O[C@H](C)[C@H]([*])C[C@H]3O)[C@@H]2O)[C@@H](OP(=O)(O)OC[C@H](O)[C@H](O)[C@H](O)CO)[C@@H]([C@H](O)CO)O1
```

**SMILES**

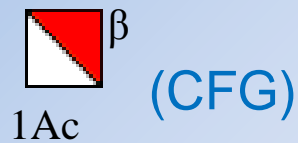
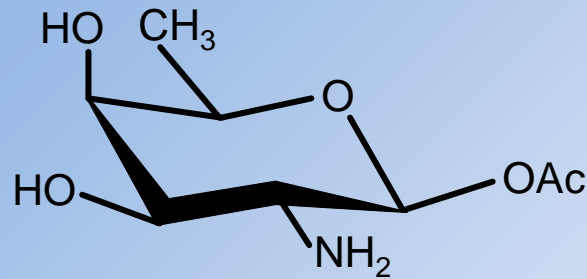
```
2.0/5,5,5/[Aad1122h-2a_2-6][h222h][a11221h-1b_1-5][a2d12m-1a_1-5][A1m_2*N]/1-2-3-4-5/a1-e2_a2-d4~_a4-c1_a5-b1*OPO*/3O/3=O_c3-d1
```

**WURCS**

```
-4) aXAbep (1-3) bXLDmanHepp (1-4) [xDRib-ol (1-P-5) , xLAla? (2-1) ] aXKdop (2-
```

**CSDB**

# Nomenclature fuzziness



bDFucpN(1-1)Ac (CSDB)

← unambiguously maps to structure but, nevertheless, is human-readable

D-FucpN-β1OAc

beta-fucosamine acetate

1-acetoxy-beta-D-fucopyranosamine

2-deoxy-2-amino-β-D-fucopyranosyl acetate (IUPAC)

β-D-fucosamine acetic ester

β-6-deoxy-D-galactosamine acetate

b-dgal-HEX|1:5|2-amino|1-acetate (GlycoCT)

β-D-фукозамин-1-О-ацетат (another human language)

(2S,3R,4R,5R,6R)-3-amino-4,5-dihydroxy-6-methyltetrahydro-2H-pyran-2-yl acetate (IUPAC)

N[C@H]([C@H]([C@H]([C@@H](C)O1)O)O)[C@@H]1OC(C)=O (SMILES)

1S/C8H15NO5/c1-3-6(11)7(12)5(9)8(13-3)14-4(2)10/h3,5-8,11-12H,9H2,1-2H3/t3-,5-,6+,7-,8+/m1/s1 (InChI)



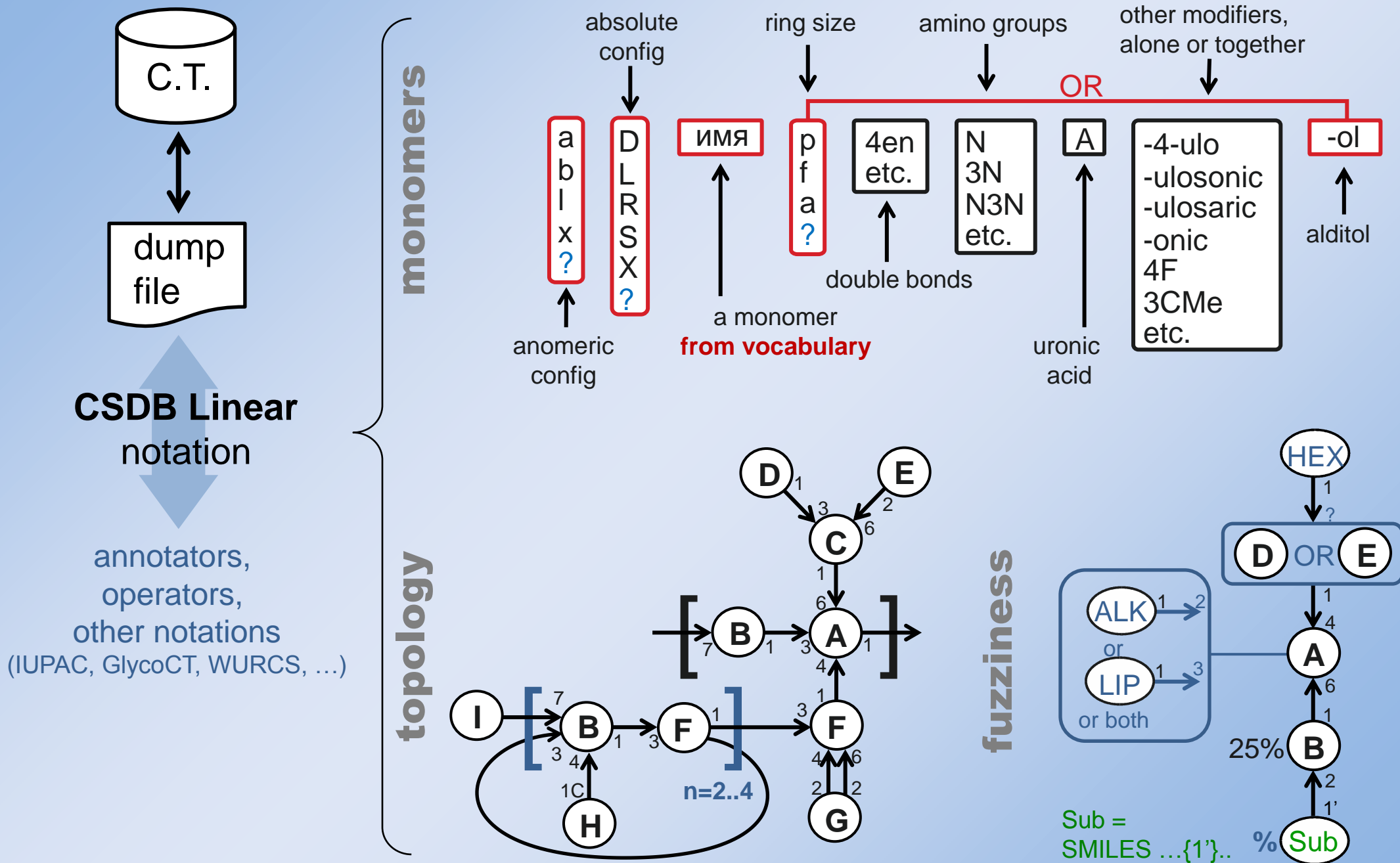
# Existing molecular description, MOL? 25

- A lot of deep stones in translation of coordinates to residue-based notation (to publish like  $\alpha$ -D-Galp-(1-3)- $\beta$ -D-Glcp)
- Difficult translation from a residue-based notation to coordinates  
=> annotation laboriousness
- Cannot describe underdetermined structures
- Data are not visually linked to knowledge  
↓
- Non-human-readable → difficult to curate → data errors
- Atom coordinates are not primary data  
but incomplete MOL (without 3D) has the same format as 3D MOL
- Large in storage and network transfer  
(and cannot be used as an URL part)

52.0606	<del>6.3910</del>	<del>-0.1606</del>	<del>C</del>	0	0	0
51.8591	<del>8.6986</del>	<del>-0.1875</del>	<del>N</del>	0	0	0
52.9844	<del>9.0584</del>	<del>0.7259</del>	<del>C</del>	0	0	1
53.8550	<del>8.9929</del>	<del>0.0662</del>	<del>H</del>	0	0	0
52.9684	<del>10.5530</del>	<del>1.3121</del>	<del>C</del>	0	0	2
52.2705	<del>11.0903</del>	<del>8.6993</del>	<del>H</del>	0	0	0
1117	1	0	0	0	0	
1118	1	0	0	0	0	
1119	1	0	0	0	0	

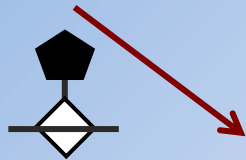
**atoms, coordinates, connectivity**

# Structural features and notation



# Structure abstraction levels

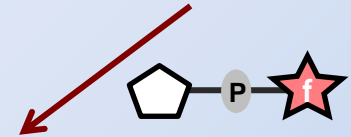
-3) [xLAla (2-6) , Ac (1-2) ]bDGalpNA (1-  
*exact fragment and its linkage*



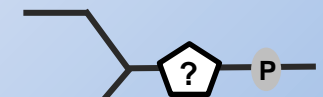
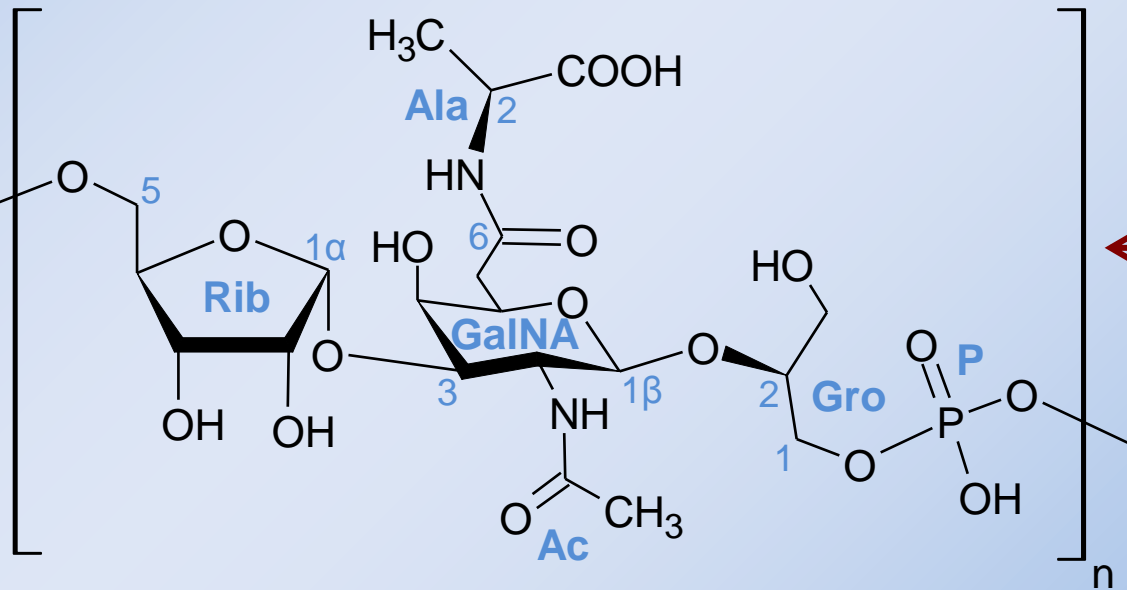
L-Ala (2-6) +

-5) α-D-Ribf(1-3) β-D-GalpNAcA (1-2) D-Gro (1-P-

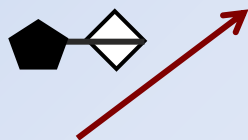
xDGro (1-P-5) aDRibf  
*unspecified topology and capping bonds*



-?) ?Dhex (1-  
*only residue class and 1-linkage*



*approximate motif*




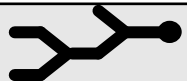












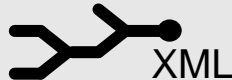





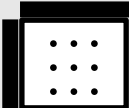





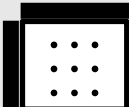





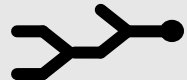





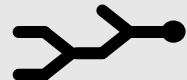





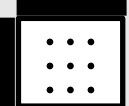











x?Ala (2-?) ?DGal?NA

*unspecified configurations, substitution positions, ringsize, N-acetylation*



HEX, xDRib?, PEP  
*partial composition*

# Notation comparison

	<i>approach</i>	complete	unambi- guous	human control	parseable	fuzziness support
IUPAC 						
IUPAC extended (SweetDB, Carbbank) 	pseudo- graphics					
Glyde I 	 XML				 URL	
WURCS (JCGGDB, ChEBI, PDB)					 URL	
GlycoCT (Glycome-DB)						
LinearCode (CFG)					 URL	
LinUCS (GlycoSCIENCES)					 URL	
KCF (KEGG)						
CSDB linear (CSDB)					 URL	

 compatibility

worse     better



- Biosynthesis pathway analysis (glycosyltransferase database)
- Oligosaccharide conformation maps
- NMR simulation and spectrum assignment ( $^{13}\text{C}$ ,  $^1\text{H}$ , 2D)
- Structure prediction from NMR and other data
- Glycome-based taxon clustering
- Feature distribution in structures and taxa
- Monomer classification

# Glycosyltransferases

Criteria:  
(in any combination)

- IDs in databases
- enzyme name / group
- gene name / cluster
- CAZy family
- organism (species, strain)
- synthesized linkage
- donor (or its fragment)
- acceptor (or its fragment)
- object cellular role
- trustworthiness level

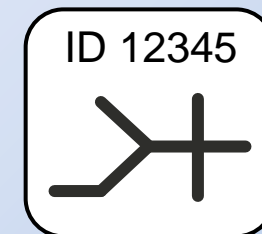
**71B1**  
 Uniprot  
 # Q9LSY9.1  
 Genbank  
 # 821729

IDs

Object:

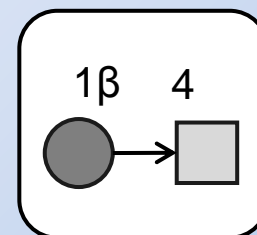


organism,  
organ, tissue

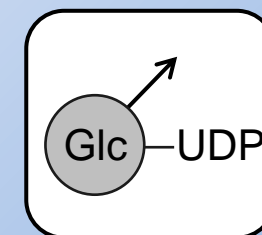


full product  
structure

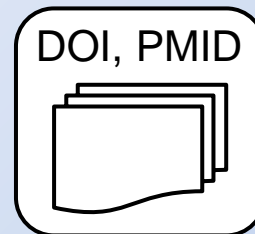
Activity:



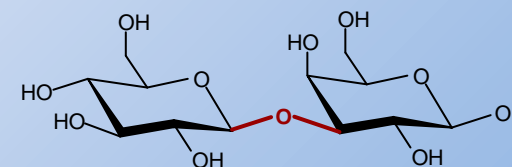
synthesized  
fragment



donor and  
acceptor



references

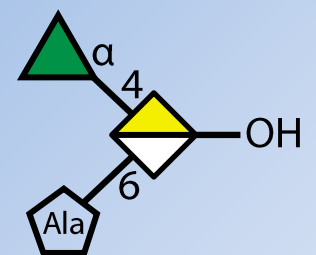


## CSDB GT

# Conformation analysis

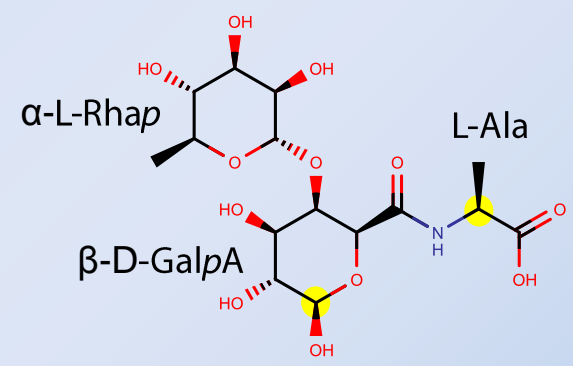
aLRhap(1-4)[x?Ala?(2-6)]?DGalpA

structure,  
incl. incomplete

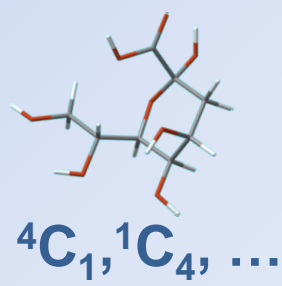
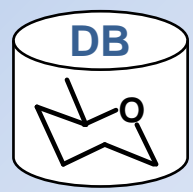


other variants  
( $\alpha$ -GalA, D-Ala, etc.)

SMILES

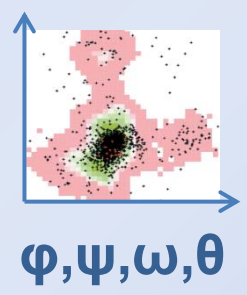
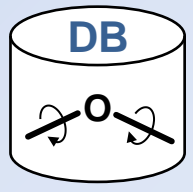


populated ring conformers  
~1000 residues

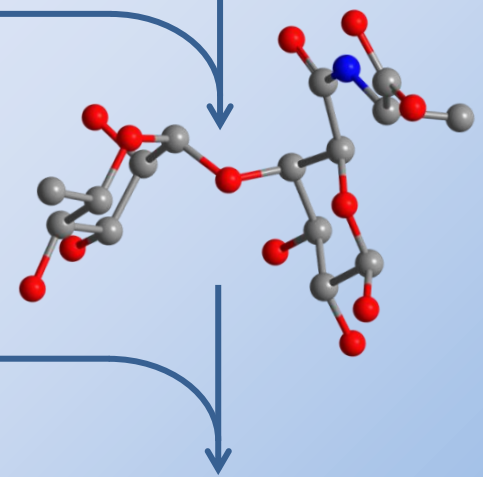


«chairification»

populated bond torsions

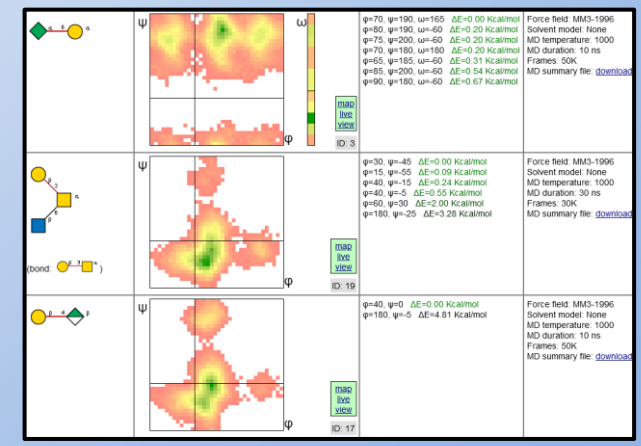


picking of minima  
MM-relaxation

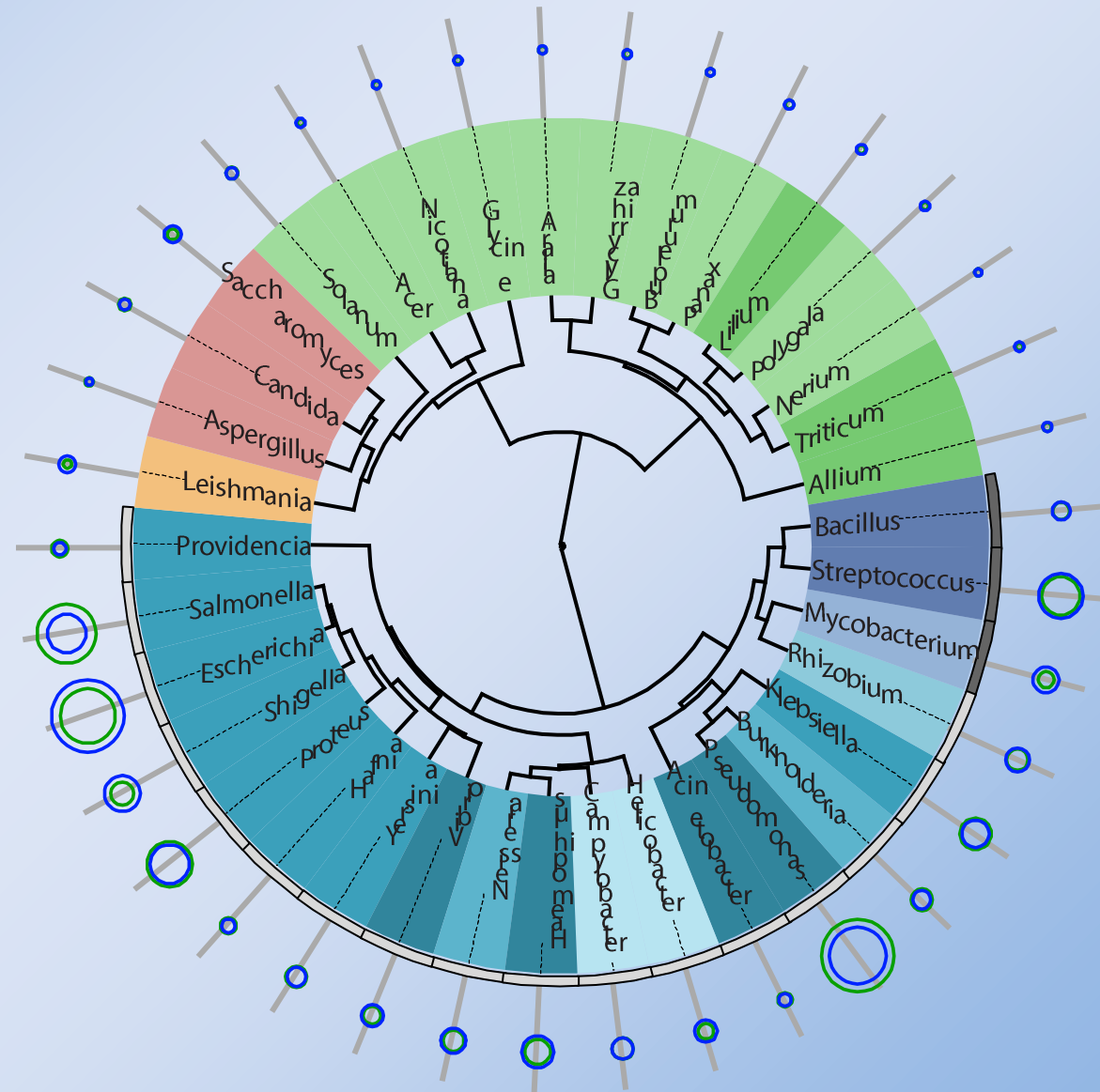
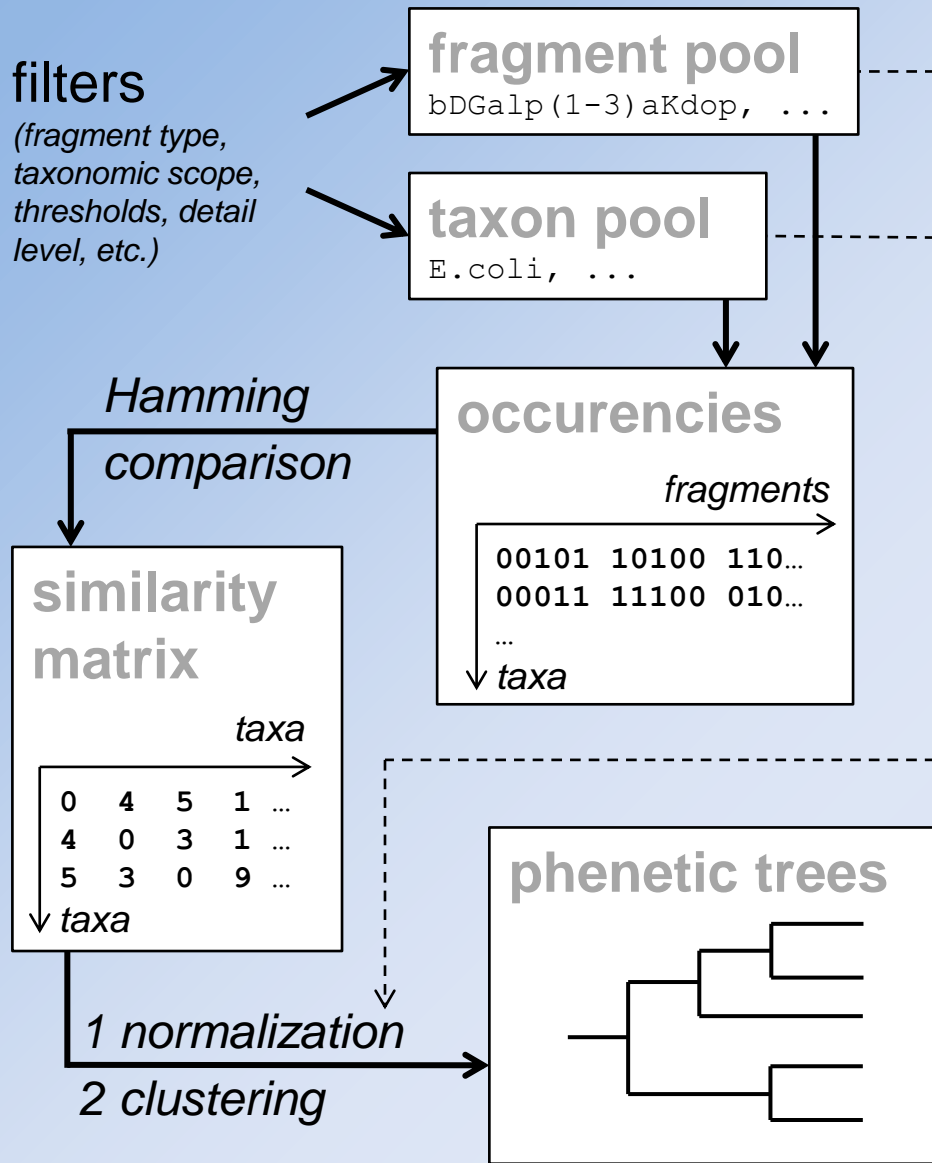


mol. dynamics  
300K, 100ns, H<sub>2</sub>O

conformers  
+ energies



# Taxon clustering



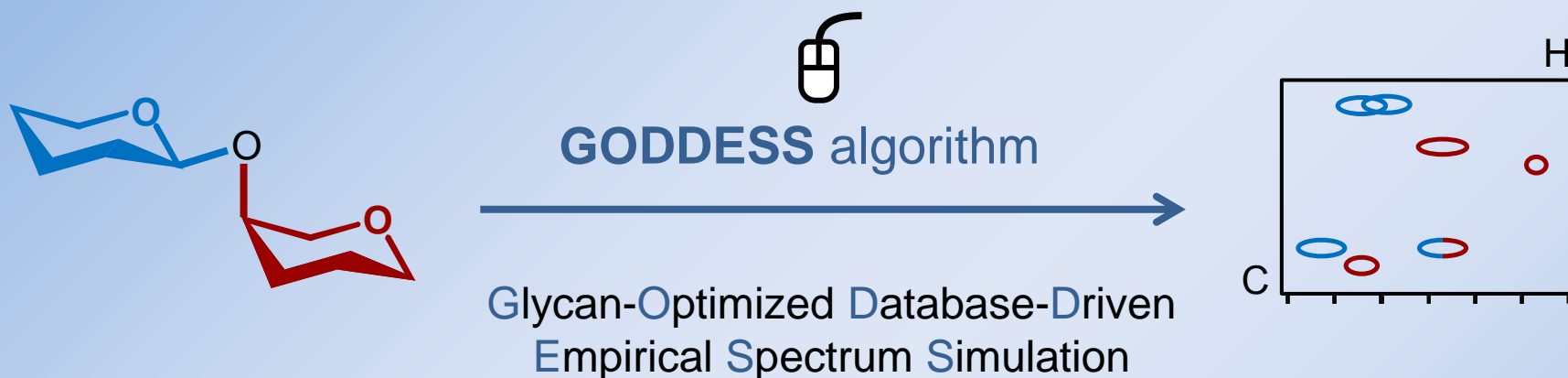
assigned structures and organisms

bacteria: Gram+ Gram-



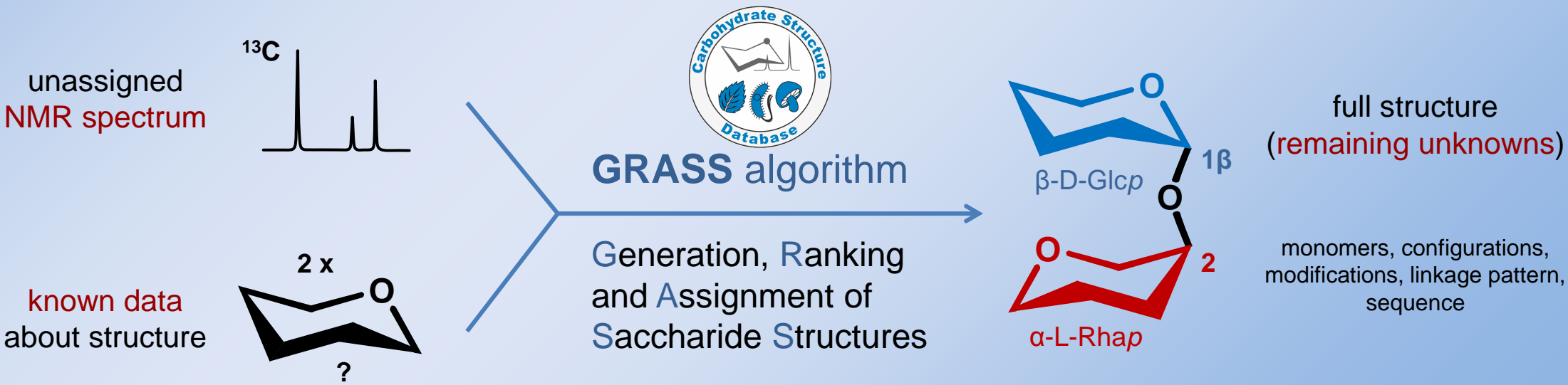
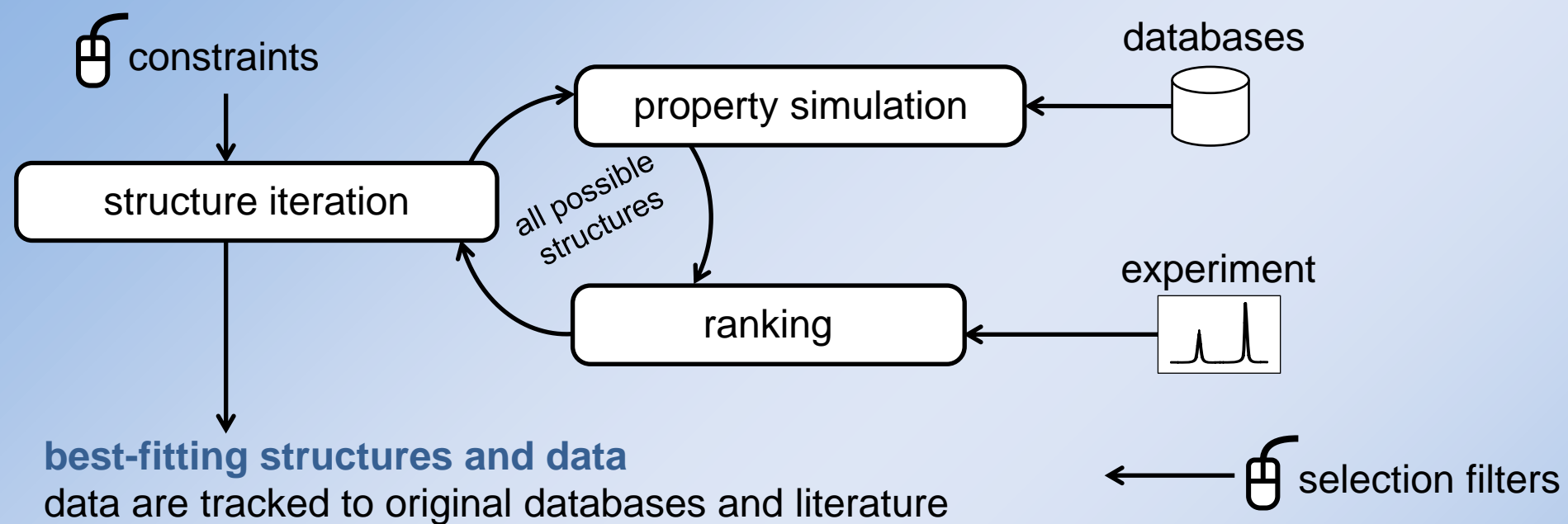
# NMR simulation

NMR is a main methods for sequence analysis in glycobiology




- Chemical shift simulation ( $^{13}\text{C}$  ~ 0.7 ppm,  $^1\text{H}$  ~ 0.06 ppm)
- Support for structure elucidation
- Signal assignment and hypothesis validation
- Structure iteration and simulation-to-experiment comparison
- Verification of molecular geometry models

# Structure iterator



# Perspective

● done in CSDB ● to do ● almost done

- Recognized human-readable language (SNFG, CSDB Linear, ...)
- Cross-project data access (GlycoRDF, GlycoCoO, central triplestore)
- Cross-project services  
(structure input & output , conformational calculations, spectra simulation, ...)
- Recognized indices  
(Glytoucan ID, MSDB, PMID, DOI, TaxID, ICD-11, PDB id, Genbank, ...)
- Standard models and protocols (API, WSDL, SPARQL, ...)
- Ideological replacement of Carbbank 
- Requirement to include IDs in publications (Glytoucan ID?)  
(who will remove unpublished / erroneous data?)

# Links and further reading

36



<http://glytoucan.org>

J. Abrahams et al. **Recent advances in glycoinformatic platforms for glycomics and glycoproteomics** (2020) *Curr Opin Struct Biol* **62**, 59-69. doi: [10.1016/j.sbi.2019.11.009](https://doi.org/10.1016/j.sbi.2019.11.009)

K. Aoki-Kinoshita **A practical guide to using glycomic databases** (2017) *Springer*. doi: [10.1007/978-4-431-56454-6](https://doi.org/10.1007/978-4-431-56454-6)



[http://jcgddb.jp/index\\_en.html](http://jcgddb.jp/index_en.html)

T. Lütteke **The use of glyco-informatics in glycochemistry** (2012) *Beilstein J Org Chem* **8**, 915-929. doi: [10.3762/bjoc.8.104](https://doi.org/10.3762/bjoc.8.104)



<http://glycosciences.de>



<http://www.genome.jp/kegg/glycan/>

Ph. Toukach, K. Egorova **Carbohydrate Structure Database merged from bacterial, plant and fungal parts** (2016) *Nucl Acid Res* **44**, D1229–D1236. doi: [10.1093/nar/gkv840](https://doi.org/10.1093/nar/gkv840)



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

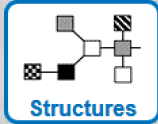
<http://csdb.glycoscience.ru>

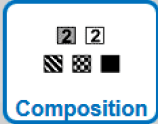



<http://toukach.ru/rus/glyco-db.htm>


# CSDB on the Internet

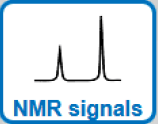
**Database search**

  
Structures

  
Composition

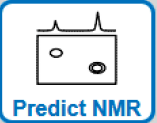
  
Organisms

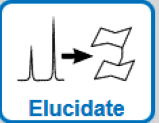
  
Publications

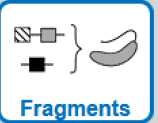
  
NMR signals

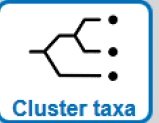
Additional operations are available from the [left menu](#). If you don't see it [click here](#)


**Useful tools**


  
Predict NMR

  
Elucidate

  
Fragments

  
Cluster taxa

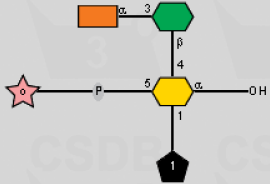
  
GT activities

  
Examples

**NMR spectrum simulation**

Please, select how to input a structure:

- [Input using Structure Wizard](#)
- [Select from library](#)
- [Draw in Glycan Builder](#)
- [Convert from GlycoCT](#)
- [Use expert form \(field below\)](#)




1 = L-Ala

**Structure in CSDB encoding:**

(this field is editable) [Help on structure encoding](#)

Nucleus:   More parameters...

Solvent:   Coverage



Prokaryotes » Plants » Fungi

7005 publications (1941-2017):  
18923 compounds from  
8859 organisms  
last update: 2017 Jun 2

**Search**

- CSDB IDs
- (Sub)structure
- Composition
- Taxonomy
- Bibliography
- NMR signals

**Help**

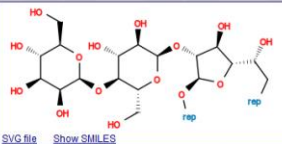
**Extras**

- NMR simulation
- Elucidation from NMR
- Monomer namespace
- Fragment abundance
- Coverage stats
- Taxon clustering
- Submit record
- Translate structure
- Feedback

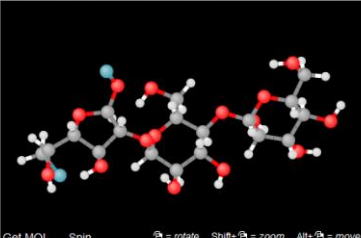
**Maintenance**

Related record ID(s): [101](#)  
NCBI Taxonomy refs (TaxIDs): [64489](#)

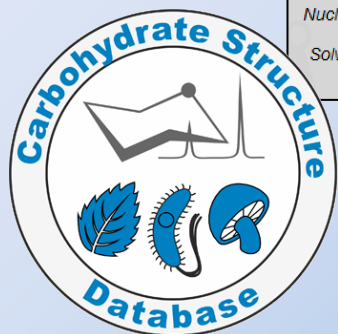
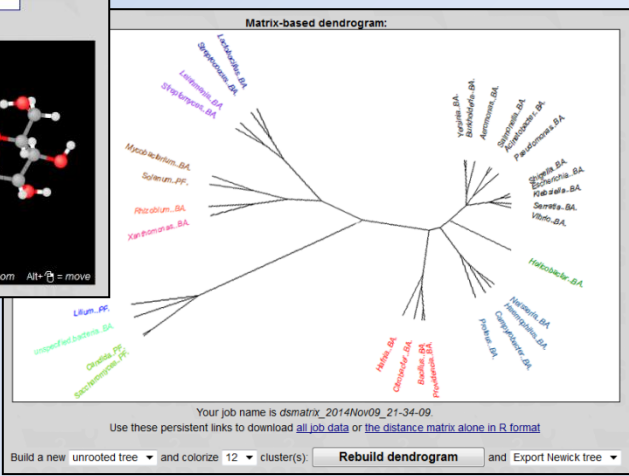
There is only one chemically distinct structure:



[SVG file](#) [Show SMILES](#)



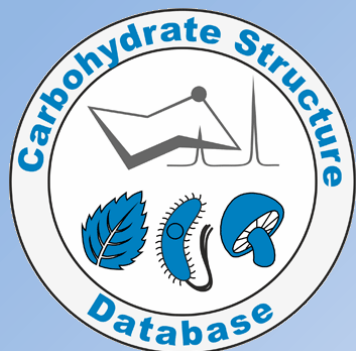
[Get MOL](#) [Spin](#) [rotate](#) [Shift+zoom](#) [Alt+move](#)



<http://csdb.glycoscience.ru>

- free access
- detailed manuals
- problem solution examples

# Credits



## Carbohydrate Structure Database

*curated content, close-to-full coverage*



Zelinsky Institute  
Moscow, Russia

programming

literature processing & verification

general support, data collection

integration, ontology

conformation analysis

ideas, R&D, notation, programming,

interface, supervision

partners

  Roman Kapaev, Andrei Bochkov, Ivan Chernyshov, ...

 Ksenia Egorova, Nadezhda Kalinchuk, Kirill Kazantsev, ...

 Yuriy Knirel

  René Ranzinger, Kiyoko Aoki-Kinoshita, Thomas Lütteke, ...

 Victor Stroylov, Sofya Scherbinina, ...

 Philip Toukach



Russian  
Foundation for  
Basic Research

2005-2007,  
2012-2020 (x3)



International Science  
and Technology  
Center

2004-2005



Russian  
Federation  
President Grants

2005-2007



Deutsches Krebs-  
ForschungsZentrum

2007-2010 (x4)



Russian Science  
Support Agency

2008-2009



Russian Science  
Foundation

2018-2020

# **Supplementary**

(used for discussion)

# Examples of queries to CSDB

- Study how an introduction of the amino group will affect the NMR chemical shifts of the lactose fragment
- Find bacterial glycans containing a galacturonic acid residue and at least one more hexose, published after 2005 in relation to antigens
- Find all compounds extracted from the plants of the genus *Solanum* which contain a solanidine constituent
- Find all carbohydrate structures having a signal close to 34 ppm in the  $^{13}\text{C}$  NMR spectrum, except those containing any octose
- Find all papers by Knirel or Shashkov AS on bacterial structures containing quinovose-4-amine amidated by any N-acetylated amino acid
- Find all bacterial nonose monosaccharide structures (monomers or homopolymers)
- Predict  $^{13}\text{C}$  NMR spectrum of 3-O- $\alpha$ -abequosyl-6-deoxy- $\beta$ -D-mannoheptopyranosyl-(D-ribitol-1)-phosphate in water solution and explore the credibility of chemical shifts simulated with lowest reported trustworthiness
- Rank structural hypotheses for an unelucidated oligomer conforming to an experimental  $^{13}\text{C}$  NMR spectrum and containing bacillosamine, galacturonic acid and lysine residues
- Study monomeric composition of two fungal species, *Aspergillus oryzae* and *Aspergillus fumigatus*, and reveal which monomers occupy termini of side chains
- Find which dimeric fragments (including sugars, aglycons and other residues) of higher plant carbohydrates are specific to lupins
- Study coverage statistics of *Proteobacteria*



# Resource Description Framework

**RDF is a model to store data as *object-predicate-subject* triples.**

- 😊 allows federated queries with minimal knowledge of source DB formats
- 😞 needs a triplestore and agreed ontology

**Question:** find carrier protein data for any given glycan in JCGGDB.

**Problem:** JCGGDB does not have links to protein databases.

## **Preamble:**

JCGGDB entries have links to GlycomeDB IDs.

Both GlycomeDB and UniCarbKB have structures in GlycoCT format.

UniCarbKB entries have links to UniProt IDs.

## **Solution** (*9-line SPARQL script*):

Map JCGGDB IDs to UniCarbKB IDs using GlycomeDB and retrieve the UniProt IDs from UniCarbKB for each JCGGDB ID.

## **Need:**

ontology standard → data exported to RDF → a triplestore → SPARQL endpoint

**GlycoRDF is a first formal carbohydrate ontology (OWL),**  
GlycoCoO is its extension for glycoconjugates.

# Structure input and output

CSDB/SNFG structure editor

Popular Small sugars Hexoses Higher sugars Alditols Aliphatic acids Other acids Superclasses

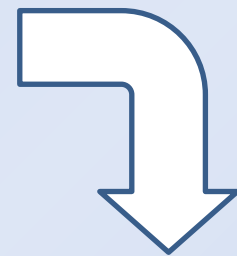
Glc GlcNAc GlcA QuiNAc Gal GalNAc GalA Fuc FucNAc Man Rha LDmanH Ara Ara4N Xyl Fru P Kdo Neu5Ac Gro Ala

Novice Expert Insert Replace Oligo Poly Ac Am Cm Cho Fo Me Et Pr EtN Allyl Bz P S Pyr NH2

search residues search modifications

Chemical repeating unit; n=10

-3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]?DGlcp(1-?)[Ac(1-2)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(



Previews Refresh

RES  
1r:r1  
REP  
REP1:6o(3+1)2d=-1--1  
RES  
2b:b-dgal-HEX-1:4  
3b:x-dglc-HEX-1:5

Subst-(7-3)-D-Rib-ol-(1--P--4)--+  
-3)-a-L-Fucp-(1-6)-D-Glcp-(1-?)b-D-GalfNac-(1-

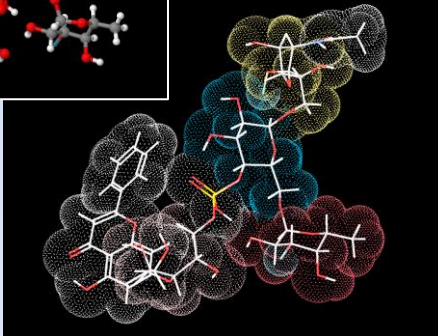
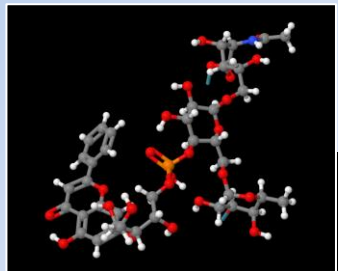
Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23

REMARK 300 USING GENERATED NAME Fda FOR RESIDUE SUB:(Ac,2)bdGalFN

REMARK 300 USING GENERATED NAME Fda FOR RESIDUE xSubst

AUTHOR GENERATED BY OPENSABEL 2.4.1, RESIDUE ASSIGNMENT BY CARBOHYDRATE STRUCT

CONFND	CSDB linear =	1 R	2 O1	3 C1	4 O4	5 O4	6 O5	7 O5	8 O6	9 O6	10 O3	11 O3	12 O2	13 N2	14 H1
	-3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]adGlcp(1-6)[Ac(1-2)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23	1	2	3	4	5	6	7	8	9	10	11	12	13	14
		0.156	1.111	2.494	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		1.377	1.487	3.097	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		1.643	2.889	2.961	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		1.951	3.280	1.601	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.374	3.123	1.402	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.933	4.392	2.736	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.890	5.497	1.642	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.185	4.772	-0.562	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.501	3.900	-1.665	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		4.005	2.825	2.771	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		4.252	1.419	2.915	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		2.912	3.229	3.759	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.019	4.689	4.105	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		0.780	3.467	3.305	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		2.260	0.749	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		4.211	0.481	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		6.274	1.168	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		5.771	-0.887	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		4.869	-0.431	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.336	2.931	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		1.146	2.209	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		2.647	4.689	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		5.266	3.368	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		5.202	5.040	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		6.697	5.164	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		4.543	5.689	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		6.973	6.220	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		7.056	4.765	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		7.165	4.614	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		2.820	-1.808	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.154	-2.751	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.270	-4.096	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
		3.728	-4.973	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000



There are 3 chemically distinct structures. Please, select:

1. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]?DGlcp(1-3)[Ac(1-2)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23
2. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]?DGlcp(1-5)[Ac(1-2)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23
3. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]?DGlcp(1-6)[Ac(1-2)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23

SMILES  
[\*]O[C@@H]1O[C@@H]([C@H](O)COC2O[C@H](CO[C@@H]3O[C@@H](C)[C@@H](O)[C@@H]([\*])[C@@H]3O)[C@@H](OP(=O)(O)OC[C@H](O)[C@@H](O)C3cc(O)c4c(=O)cc(-c5ccccc5)oc4c3)[C@H](O)CO)[C@H](O)[C@@H]2O)[C@@H](O)[C@@H]1NC(C)=O

There are 2 sterically distinct structures. Please, select:

1. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]adGlcp(1-6)[Ac(1-2)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23
2. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]bDGalfN(1- // Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23

Files: MOL, PDB, Glycam

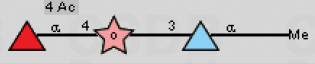
3D Shift+Z = zoom Shift+X = pan Alt+R = rotate Ctrl+M = menu

# (Sub)structure search

### Structure wizard

Topology: 3 residues (linear: A->B->C) (A)→(B)→(C)

Structure:



Residue (A):

()

[aLFucp](#)

substitutes  of Residue B

is terminal

add substitution  
 add substituent  at   
 add substituent  
 add substituent  
 add substituent

Residue (B):

()

[DRib-ol](#)

substitutes  of Residue C

add substitution  
 add substituent  
 add substituent  
 add substituent

Residue (C):

()

[a?6dTal?](#)

has aglycon:

add substitution  
 add substituent  
 add substituent  
 add substituent

**Structure in CSDB encoding:**


[Return the structure to the search page and close this window](#)

[Home](#) [Help](#)


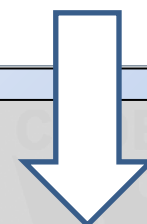


### Glycan Builder

File Edit Structure View Help




Linkage    Chirality  Ring

### Search for (sub)structure

Please, select how to input structure:

- [Input using Structure Wizard](#)
- [Select from library](#)
- [Draw in Glycan Builder](#)
- [Convert from GlycoCT](#)
- [Copy from the previous query \(aLFucp3N\)](#)
- [Use expert form \(field below\)](#)



**Structural fragment in CSDB encoding:**

*(this field is editable)* [Help on structure encoding](#)

Only those containing text:   in aglycons, aliases or linear code  in trivial names

**Search scope:**

Search the whole database  Search in the result of the previous query (logical AND)  
 Combine with the result of the previous query (logical OR)  
 Negate search (find results NOT matching current query)

Treat search term as a   
 Search for molecule types:   
 Search for structures with published NMR data only  
 Restrict compound class:   
 Restrict taxonomical domain:

Previous results: 122 structures: [<ID list>](#)

& display  records per page.

[Predict NMR](#) [Sweet 3D model](#) [Home](#) [Help](#) [HELP !!!](#)

# Organism search

Found **12** organisms. Displayed organisms from **1** to **12**  
[Expand all organisms](#)   [Show all as text \(SweetDB notation\)](#)

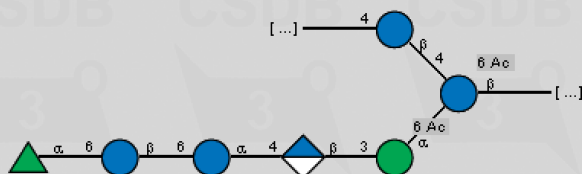
1. (Organism ID: 1005)

**[Acetobacter xylinum](#)**  
 (Ancestor NCBI TaxID 28448, [species name lookup](#))

Later renamed to: [Komagataeibacter xylinus](#)  
 Taxonomic group: bacteria  
 Phylum: Proteobacteria

The following compound(s) are assigned to this organism:

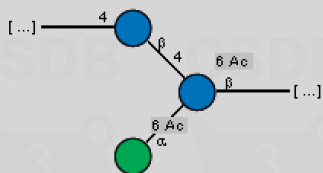
- Compound ID: 1717



[Show legend](#)  
[Show as text](#)

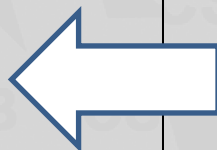
Carbohydrate Research 2004, "Synergistic interactions between the genetically modified bacterial polysaccharide P2 and carob or konjac mannan"  
[CSDB ID 9262](#) (all data & tools)

- Compound ID: 1720



[Show legend](#)  
[Show as text](#)

Carbohydrate Research 2004, "Synergistic interactions between the genetically modified bacterial polysaccharide P2 and carob or konjac mannan"  
[CSDB ID 9414](#) (all data & tools)



## Search for organism

**Display domains:**    bacteria    archaea    protista    algae    fungi    plants    animals

---

**Genus:**    
 Acetobacter  
 Acholeplasma  
 Acidithiobacillus  
 Acinetobacter  
 Acremonium  
 Actinobacillus  
 Actinobaculum

**Species:**    
 sp.  
 diazotrophicus  
 methanolicus  
 tropicalis  
 xylinum

**Strain / subspecies:**    
 B42  
 CKE5  
 CKEP  
 CR1/4  
 IFO 13693

Specify:

---

**Search scope:**

Search the whole database       Search among HOST organisms  
 Search in the result of the previous query (logical AND)       Use NCBI taxID  
 Combine with the result of the previous query (logical OR)       Include subtaxons  
 Negate search (find results NOT matching current query)

Previous results: 6 structures: [<ID list>](#)

& display  records per page.

[List of organisms](#)      [Home](#)      [Help](#)

---

**Process taxonomy in NCBI Taxonomy DB** (fields are editable):

Genus:     Species:

# Bibliography search

Found **3** publications. Displayed publications from **1** to **3**

[Expand all publications](#)   [Show all as text \(SweetDB notation\)](#)

1. (Article ID: 1525)

Knirel YA, Lindner B, Vinogradov EV, Shaikhutdinova RZ, Senchenkova SN, Kocha  
**Cold temperature-induced modifications to the composition and structure of Yersinia pestis**

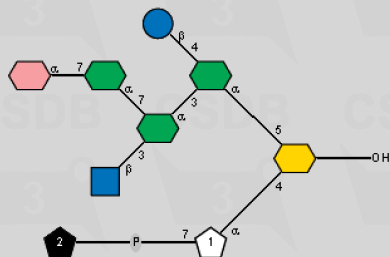
*Carbohydrate Research* **340(9)** (2005) 1625-1630

Following a report of variations in the lipopolysaccharide (LPS) structure of *Y. pestis* at 6 degrees C and flea (25 degrees C) temperatures, a number of changes to the LPS of the bacterium was identified. The LPS of the bacterium was cultivated at a temperature of winter-hibernating rodents (6 degrees C) differs from the LPS of the bacterium cultivated at 25 degrees C. The LPS of the bacterium known Y. pestis KM218 differs from the LPS of the bacterium cultivated at 25 degrees C in the following: (i) replacement of terminal galactose with terminal mannose; (ii) phosphorylation of terminal oct-2-ulosonic acid with phosphoethanolamine; (iii) the absence of glycine; lipid A differs in the lack of any 4-amino-4-deoxyoctanoic acid; (iv) the absence of a fatty acid(s). The data obtained suggest that cold temperature-induced modifications to the LPS of Y. pestis may be a mechanism of control of the synthesis of Y. pestis LPS.

*Lipopolysaccharide, structure, core, modification, agent, composition, Yersinia pestis, Plague*

The publication contains the following compound(s):

• Compound ID: 4209



1 = a-Kop  
2 = EtN

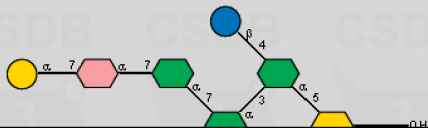
[Show legend](#)

[Show as text](#)

*Yersinia pestis* KM218

[CSDB ID 10076](#) (all data & tools)

• Compound ID: 4210



## Search for bibliography

**Authors:** "Knirel YA" OR Toukach

[Help on author/keyword query syntax](#)

ä ö ü á é í ó ç š

start with: tou

**Title:** pestis OR plague\*

(content of title) [Help on title/abstract query syntax](#)

search also in abstract

**Keywords:** structure? OR composition?

(content of keyword section) [Help on author/keyword query syntax](#)

search also in title

**Journal:** Carbohydrate Letters  
Carbohydrate Polymers  
Carbohydrate Research  
Cell  
Cell Chemical Biology  
Cell and Tissue Research

**Year:** 1983  
1984  
1985  
1986  
1987  
1988  
1989

**Vol:** \*  
**Page:** \*

### Search scope:

- Search the whole database  
 Search in the result of the previous query (logical AND)  
 Combine with the result of the previous query (logical OR)  
 Negate search (find results NOT matching current query)
- Publications with structure elucidation only  
 Restrict taxonomical domain: All domains

Previous results: 3 publications: 2953,201,1525

& display 30 records per page.

[PubMed XML](#)

[Home](#)

[Help](#)

## Author index:

[Toubetto K](#)  
[Toukach FV](#)

[Toussaint A](#)

The listed author names start with 'Tou'.  
Click an author name to copy it to the author field in the caller form.

[Close this window](#)

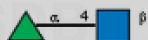
# Conformation search

## Search for disaccharide conformation maps

Use the following criteria alone or in any combination to search for conformation maps.

Conformation ID:  ANY

Model:       (only those components are listed for which conformation maps are stored)

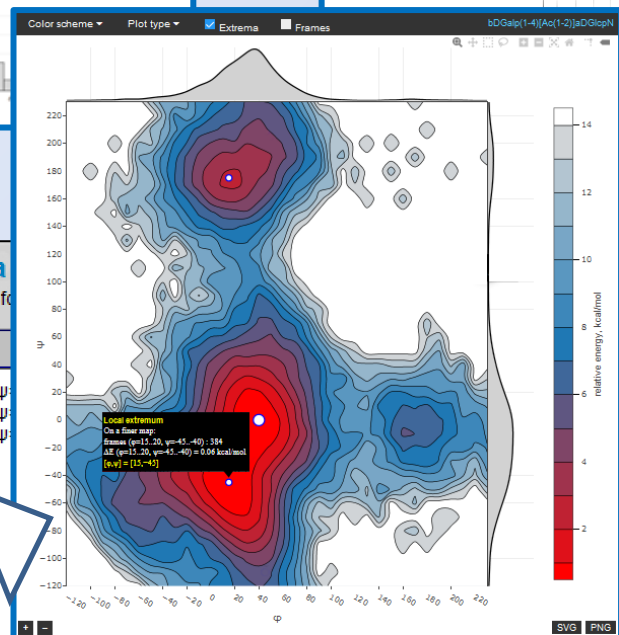
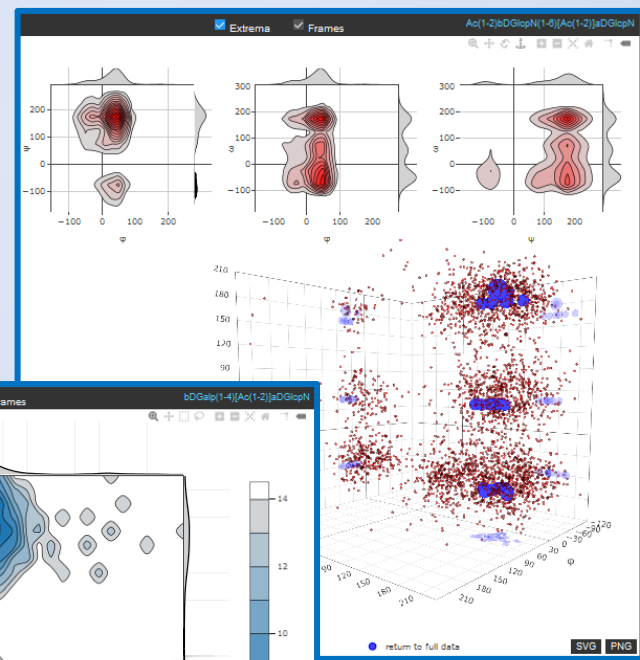
or type dimeric fragment in CSDB encoding  

Force field:  MM3-2000

Temperature:  1000


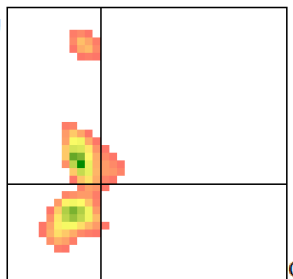
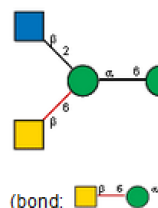

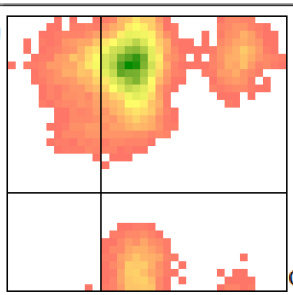
Solvent model:  any

[Home](#) [Help](#)



## CSDB conformation data

12 conformation maps have been found

Model structure	Conformation map	
		<p>φ=-30, ψ=...</p> <p>φ=-35, ψ=...</p> <p>φ=-25, ψ=...</p> <p>map live view</p> <p>ID: 1610</p>
 (bond:  )		<p>φ=25, ψ=165, ω=-75 ΔE=0.00 Kcal/mol</p> <p>φ=45, ψ=170, ω=-60 ΔE=0.00 Kcal/mol</p> <p>φ=45, ψ=195, ω=-75 ΔE=0.38 Kcal/mol</p> <p>φ=50, ψ=185, ω=-60 ΔE=0.72 Kcal/mol</p> <p>φ=25, ψ=165, ω=180 ΔE=0.72 Kcal/mol</p> <p>φ=30, ψ=155, ω=-60 ΔE=0.72 Kcal/mol</p> <p>φ=35, ψ=180, ω=180 ΔE=0.85 Kcal/mol</p> <p>φ=45, ψ=215, ω=-60 ΔE=0.85 Kcal/mol</p> <p>φ=40, ψ=170, ω=165 ΔE=0.99 Kcal/mol</p> <p>φ=20, ψ=185, ω=165 ΔE=0.99 Kcal/mol</p> <p>φ=55, ψ=195, ω=165 ΔE=0.99 Kcal/mol</p> <p>φ=55, ψ=155, ω=-60 ΔE=1.13 Kcal/mol</p> <p>φ=45, ψ=185, ω=165 ΔE=1.13 Kcal/mol</p> <p>φ=30, ω=145, ω=-60 ΔE=1.13 Kcal/mol</p> <p>Force field: MM3-1996</p> <p>Solvent model: None</p> <p>MD temperature: 1000</p> <p>MD duration: 30 ns</p> <p>Frames: 30K</p> <p>MD summary file: <a href="#">download</a></p> <p>map live view</p> <p>ID: 907</p>

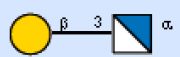
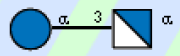
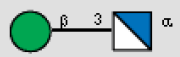
3D

2D

# Glycosyltransferase search

## CSDB glycosyltransferase search

42 glycosyltransferase activities have been identified in the CSDB.  
Please note that GTR database covers only two species: *Escherichia coli* and *Yersinia enterocolitica*.

Enzyme	Gene	Activity
Name: WbbD UniProt ID: <a href="#">Q03084*</a>	?	Synthesized dimer: bDGalp(1-3)aDGlcPn  Donor (ID 19342): <a href="#">DGalp(1-P-P-5)nucU</a> Acceptor (ID 19715): <a href="#">Ph(1-11)[Ac(1-2)aDGlcPn(1-P-1)]Subst // Subst = undecan-1,11-diol</a> Status: evidence <i>in vitro</i> <a href="#">?</a> Confirmation methods: <i>in vitro</i> (crude extract) ID: 2053
Name: WbbG UniProt ID: <a href="#">Q0H8C8*</a>		Synthesized dimer: aDGlcP(1-3)aDGlcPn  Status: indirect evidence <i>in vivo</i> <a href="#">?</a> Confirmation methods: mutation (knockout) Notes: Repeating unit of the O148 antigen. ID: 2151
Name: WbaD UniProt ID: <a href="#">Q1L815*</a>	Name: wbaD GenBank ID: <a href="#">7156002*</a>	Synthesized dimer: bDManp(1-3)aDGlcPn  Donor (ID 19855): <a href="#">DManp(1-P-P-5)nucG</a>


### CSDB glycosyltransferase search



Use the following conditions alone or in any combination to search for glycosyltransferases. Any field may be left blank for no restrictions.

**GT names and IDs:** Type enzyme name, e.g. "Orf10". Wildcards (\* and ?) are supported.  
 Enzyme name:

**Organism:** Select species  Type strain/serogroup

**Molecule role:** Filter by target structure

**Synthesized bond:** Type dimeric fragment in CSDB encoding or use tools  
 [Use Wizard](#) 

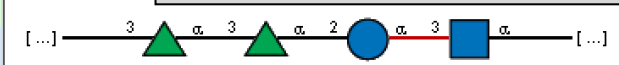
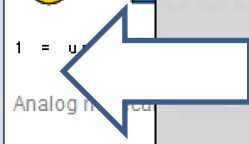
**Donor & acceptor:** Type donor CSDB encoding or use tools  
 [Use Wizard](#)   
Type acceptor CSDB encoding or use tools  
 [Use Wizard](#) 

Treat donor/acceptor as fragments

**Confirmation status:** Filter results to those

[Search!](#)

[Home](#)   [Help](#)   [HELP !!!](#) [?](#)



CSDB ID(s): [11572](#), [21578](#), [23062](#), [26257](#), [27289](#)



Molecule role: O-antigen

Organism (ID 1863): *Escherichia coli* O77

Full structure (ID 4600):



**Zhou et al. 2016**  
 DOI: [10.1016/j.carres.2016.02.007](#)

Wang et al. 2007  
 DOI: [10.1099/mic.0.2007](#)

# Taxon clusterization

**Scope settings**

Limit taxonomical scope to: **phylum**

Display groups:  bacteria  archaea  protista  algae  fungi  plants  animals

**Phylum:** (select multiple with CTRL key)

- (unspecified bacteria)
- (unspecified protista)
- Actinobacteria
- Bacteroidetes
- Chlamydiae
- Chloroflexi**
- Crenarchaeota
- Cyanobacteria

**General settings**

species Rank of taxons to compare (should be lower than selected scope). [Specify exact species \(all\)](#)

50 **Taxon population threshold.** Minimal number of structures\* assigned to a taxon or its subtaxons, to include this taxon in calculation (affects selection of taxons). Check to use this filter.

15 % **Normalized taxon population threshold.** Minimal part of structures\* assigned to a taxon or its subtaxons, to include this taxon in calculation (affects selection of taxons). Normalized by the total number of structures\* in the database. Check to use this filter.

50 **Structure abundance threshold.** Minimal number of structures\* in which a fragment should be contained to be qualified as 'present in biota' (affects selection of fragments)

60 **Fragment abundance threshold.** Minimal number of instances\* in which a fragment should be present to be qualified as 'present in biota' (affects selection of fragments)

2 **Fragment presence threshold.** Minimal number of instances\* in which a fragment should be present in organisms of a taxon to be qualified as present in this taxon (affects occurrence codes and thus, taxon dissimilarity)

**two residues** **Type of fragments to analyze** (dimeric or monomeric)

**only polymers** **Type of structures to analyze.** Only structures of this type are considered in fragment analysis and where marked by (\*). 'Optimized' = only polymers from bacteria, archaea and fungi, and only mono/oligomers from plants.

**R-project** **Format of the dissimilarity matrix**

**Fragment pool generation settings**

- Combine anomeric forms.** All sugar residues will be treated as 'any anomer'
- Exclude underdetermined residues.** Residues with unknown anomeric, absolute or ringsize configuration will be omitted from analysis.
- Exclude monovalent residues.** Residues like Me, Ac, etc. will be omitted from analysis. Please note, that Ac in N-acetylated amnosugars is a separate residue.
- Exclude superclasses.** Fragments with residues represented by aliases and superclasses will be omitted from analysis.
- Differentiate aliases.** Residue aliases (used for atypical residues) will be differentiated by actual residue names, otherwise they are combined under an alias name.
- Sugars only.** Fragments with non-sugar residues (including monovalent residues, like N-acetyls) will be omitted from analysis.
- Exclude aglycons.** Fragments with atypical residues at non-reducing ends will be omitted from analysis.
- Differentiate location.** The same fragments at different locations (inline, terminal, reducing) will be treated as different.
- Strict comparizon** of fragments. Unknown configurations and ringsizes are always unequal to those known (otherwise a fuzzy comparizon is performed).

**Distance matrix based on fragment presence**

The analysis was performed over all cellular organisms

Prepared 20 monomers  
Prepared 32 genera  
Generated occurrence bit-codes. [Show](#)  
Generated dissimilarity matrix. [Show](#)

**Calculation parameters:**

Hamming mode: YES  
Fragment size: monomer  
Fragment abundance filter: instance threshold: 550  
Fragment abundance filter: structures threshold: 500  
Fragment presence threshold: 2  
Differentiate structures of this type: any  
Filter: differentiate monomer positions (inline/terminal/reducing) in structures: NO  
Matrix data format: R

**Coverage data on used taxons:**

(taxons, number of organisms in a taxon, number of structures assigned to these organisms)

Acinetobacter (BA)	68	140
Aeromonas (BA)	64	122
Bacillus (BA)	95	234
Burkholderia (BA)	36	219
Solanum (PF)	46	127

**Matrix-based dendrogram:**

Your job name is `dsmatrix_2014Nov09_21-34-09`

Use these persistent links to download [all job data](#) or [the distance matrix alone in R format](#)

Build a new **unrooted tree** and colorize **12** cluster(s): **Rebuild dendrogram** and **Export Newick tree**



# Structure prediction

### Structure generation constraints:

The structure contains **6 residue(s)**: [Add residue](#)

$\alpha/\beta$	D/L	Residue	Ring form
1. ?	D	galact-2N-uronic acid	pyranose
2.		acetic acid	
3.	D	show all residues	
4.		phosphoric acid	
5. $\alpha$	?	any octose	pyranose
6.	L	alanine	

**Allowed linkages:**

	C1	C2	C3	C4	C5	C6	C7+
D-GalpNA	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	None
Ac	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	None
D-Rib-ol	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	None
P	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	None
$\alpha$ -Octp	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	None
L-Ala	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	None

**Advanced options:** [Hide](#)

Min in	Max in	Location	Ac at N	Acceptors	Remove
1	2	any	demanded	any	<input checked="" type="checkbox"/>
?	?	any		any	<input checked="" type="checkbox"/>
?	?	terminal		3, 5	<input checked="" type="checkbox"/>
?	?	reducing			<input checked="" type="checkbox"/>
?	?	any	forbidden	1	<input checked="" type="checkbox"/>

**Search depth:** Widespread structures only

**Scope:**  oligomers  polymers   $\beta$ -anomers = 1

**Advanced scope:** CH<sub>2</sub> carbons: ?  no furanoses

### Top 15 matches:

#Rank	Structure	Experimental spectrum	Simulated spectrum	Comments
#1.				$\Delta \sim 0.94$ ppm Corr = 1.000 RMS dev = 1.46 ppm Trust = 46%
#2.				$\Delta \sim 0.95$ ppm Corr = 1.000 RMS dev = 1.46 ppm Trust = 46%
#15.				$\Delta \sim 1.42$ ppm Corr = 0.999 RMS dev = 1.99 ppm Trust = 49%

### Find best matching structures:

Experimental <sup>13</sup>C NMR spectrum in water (24 signals of 24 expected):

17.4 22.9 34.7 50.5 52.4 63.9 64.9 66.2 68.3 70.6 72.3 72.4 72.7 73.3 73.6 76.5 78.6 78.8 99.2 102.6 171.2 175.2 176.0 176.5

$\pm$  2 signals

Find 15 best-fitting structures

Save generated structures

[Go!](#)

E-mail for results: [why?](#)  
user@gmail.com

