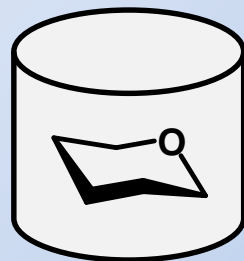


Филипп Тоукач

Информационные технологии в структурной гликохимии и гликобиологии

доклад для защиты диссертации д.х.н. (2018) с дополнениями (2020)



2005-н.в.

В ЭТОМ ДОКЛАДЕ:

введение • Зачем информатизировать гликомику?

глико-информатика

- Углеводные базы данных
- Правила гликоинформатики на примере CSDB

структуры

- Формализация углеводных структур
- Молекулярная геометрия углеводов

надстройки

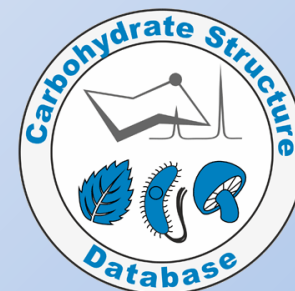
- Гликотаксономия
- Гликозилтрансферазы

глико-ЯМР

- Моделирование спектров ЯМР
- Предсказание первичной структуры

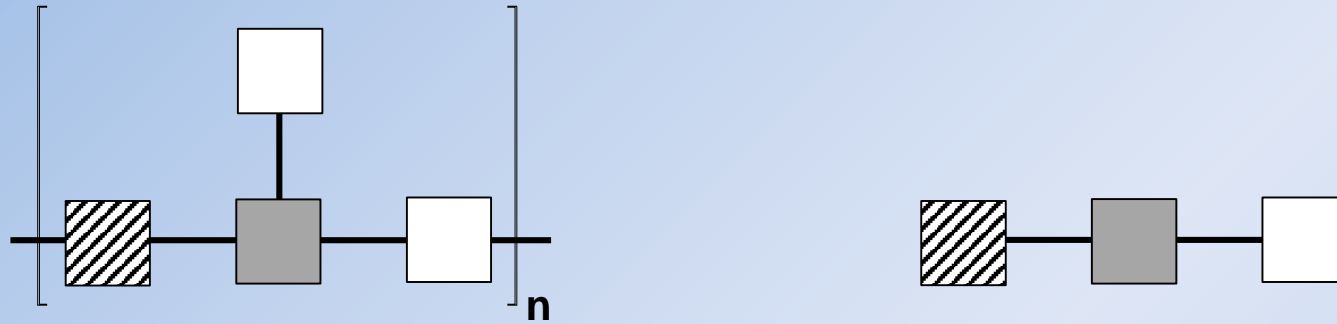
выводы

- Что сделано, опубликовано, внедрено?



<http://toukach.ru/rus/glycoinf.htm>

Гликохимия в биологии



- Структура, изомерия, конформация углеводов в клетках
- Таксономия и классификация микроорганизмов
- Гликоэпитопы и иммуноспецифичность организмов
- Гликосодержащие вакцины и лекарства
- Корреляция макросвойств организма с его углеводами
- Биосинтез и круговорот углеводов

Гликоинформатика

- **Роли углеводов:** гликозилирование белков, антигены микроорганизмов, межклеточные контакты, клеточная стенка, биоактивные гликозиды.
 - Интерес рос, объём информации увеличивался, но использовать её было проблематично.
- **Глобальная задача гликоинформатики** - обеспечить исследователей углеводов всей мощью информационных технологий.

>2000 г.

Частные задачи гликоинформатики

- **Легкий доступ к знаниям и автоматизация исследований**

Какие природные структуры похожи на заданные? Какие их фрагменты специфичны для заданных биологических видов? Где они опубликованы, в привязке к каким таксонам, болезням, и т.д.? Какие ферменты их синтезируют и с какой достоверностью это показано? На какие гликоэпитопы реагируют антитела?

- **Моделирование свойств молекул**

Молекулярная геометрия, спектры, биоактивность, ...

- **Предсказание структуры по наблюдаемым свойствам**

- **Предсказание свойств таксонов**

Кластеризация на основании гликомов, поиск схожести и различий таксонов, хемотаксономическая классификация

- **Идентификация и визуализация молекул** (в т. ч. в публикациях)

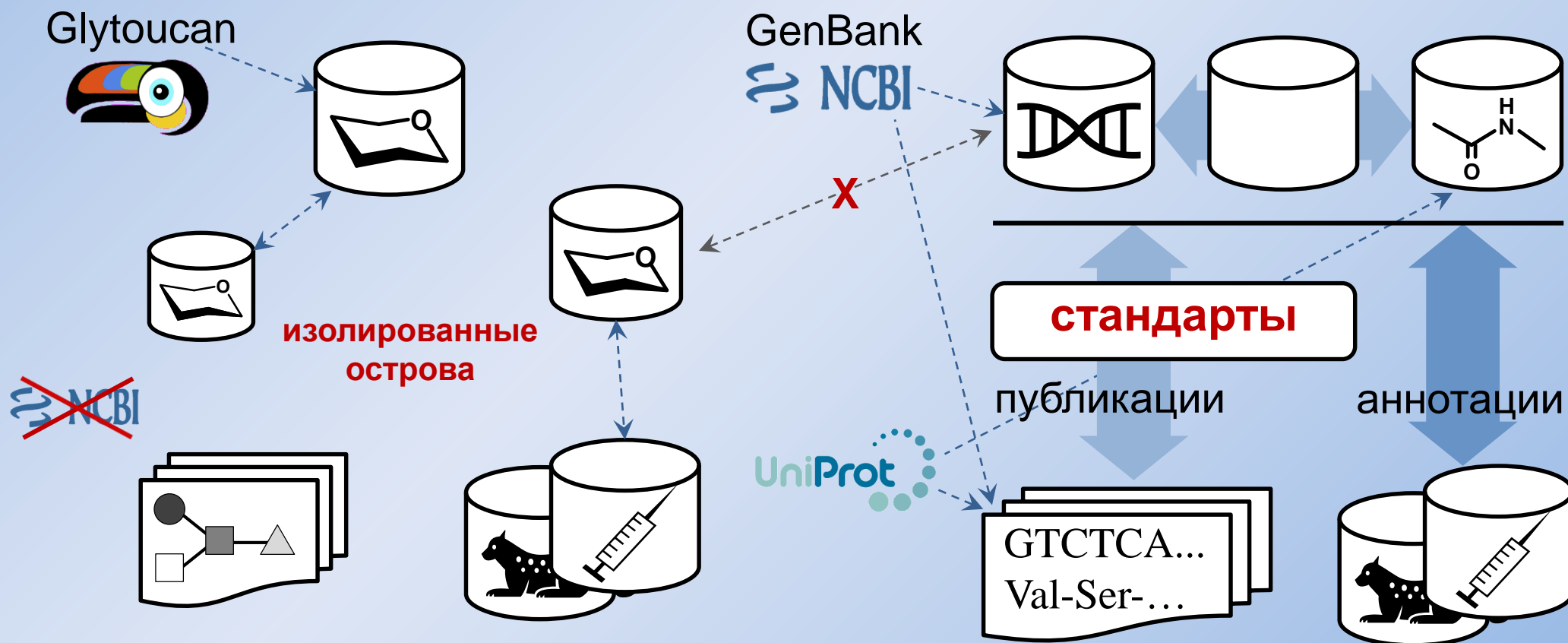
- **Методология обработки накопленных знаний**

Аннотирование публикаций, интеграция проектов, стандартизация знаний

Гликомика vs. геномика, протеомика

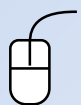
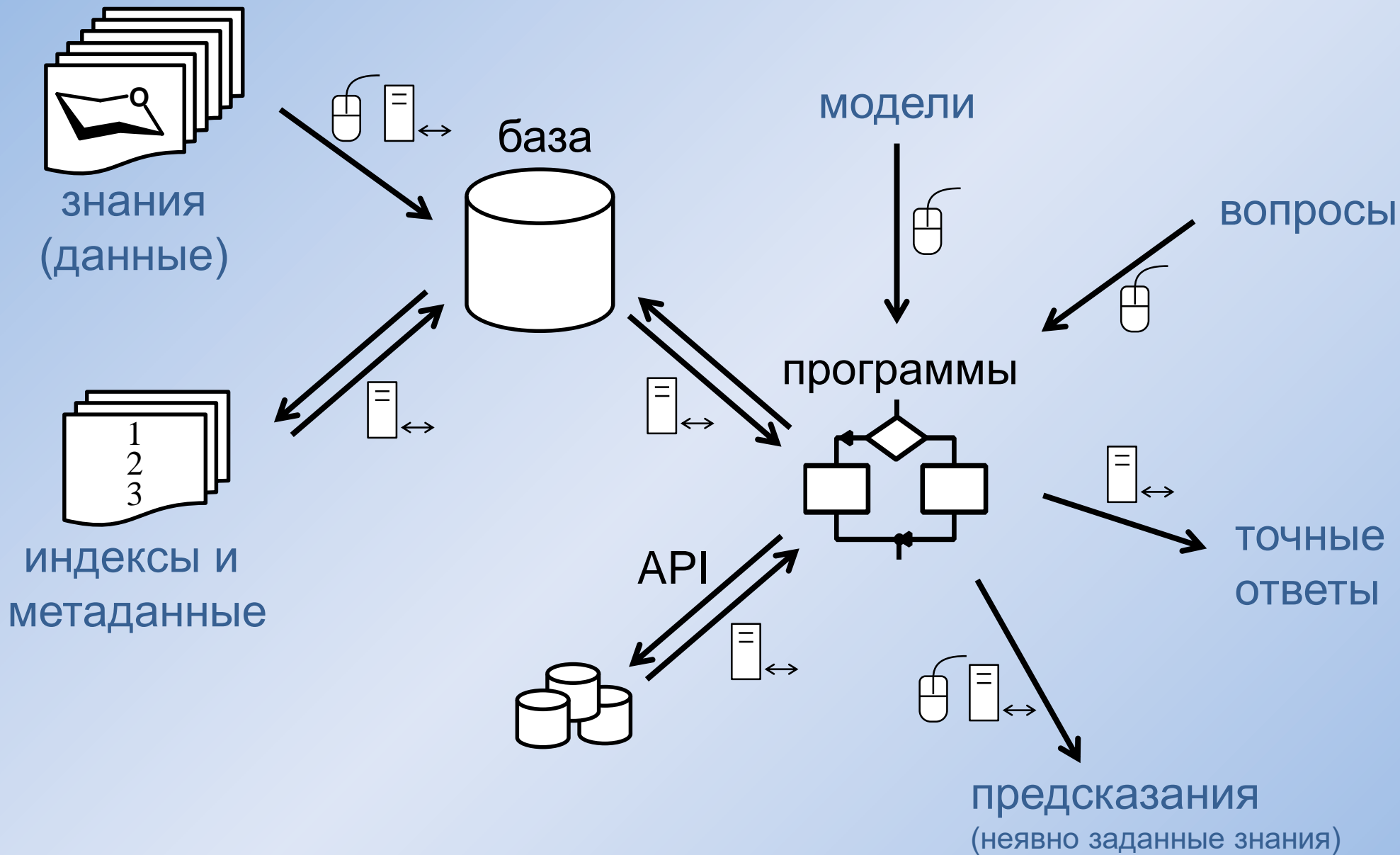
по сравнению с другими *-omics*:

- сходный объем информации (>100 000 известных структур)
- бóльшая химическая вариативность
- меньшее использование IT (базы данных, сервисы)
- меньшая стандартизация



- Вариативность и гетерогенность объектов
- Неоднозначное описание структуры
- Сложности с вводом и визуализацией больших структур
- Отсутствие стандартов
- Изолированность проектов
- Неполнота и низкое качество данных в базах
- Ресурсоемкие алгоритмы
- Нехватка системного видения у разработчиков и пользователей
(нет общепризнанных сервисов, инициативы несовместимы друг с другом)

Базы данных

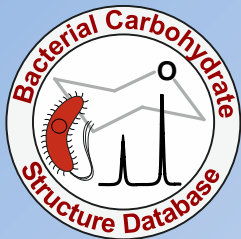


с участием человека

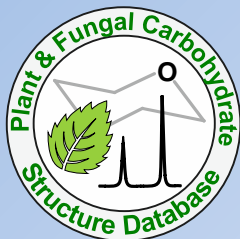


автоматически

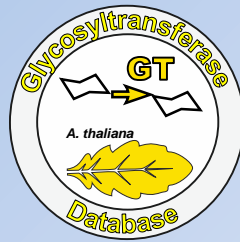
Carbohydrate Structure Database



с 2005



с 2012



с 2017

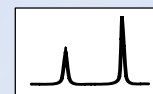
Комплекс баз данных +
Платформа для сервисов



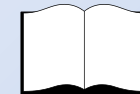
25K



13K



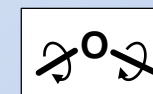
14K



10K



2K



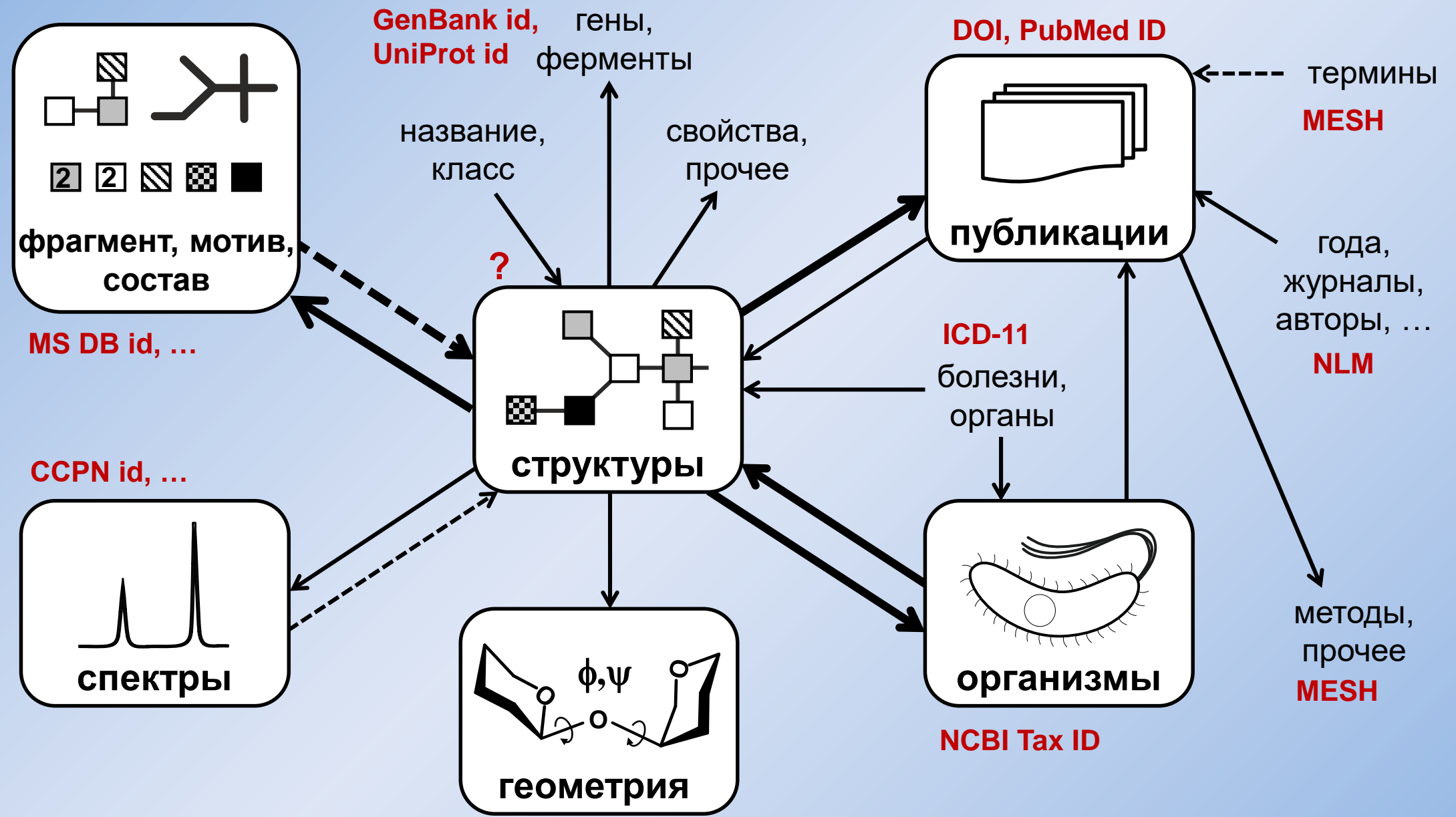
3K



CSDB

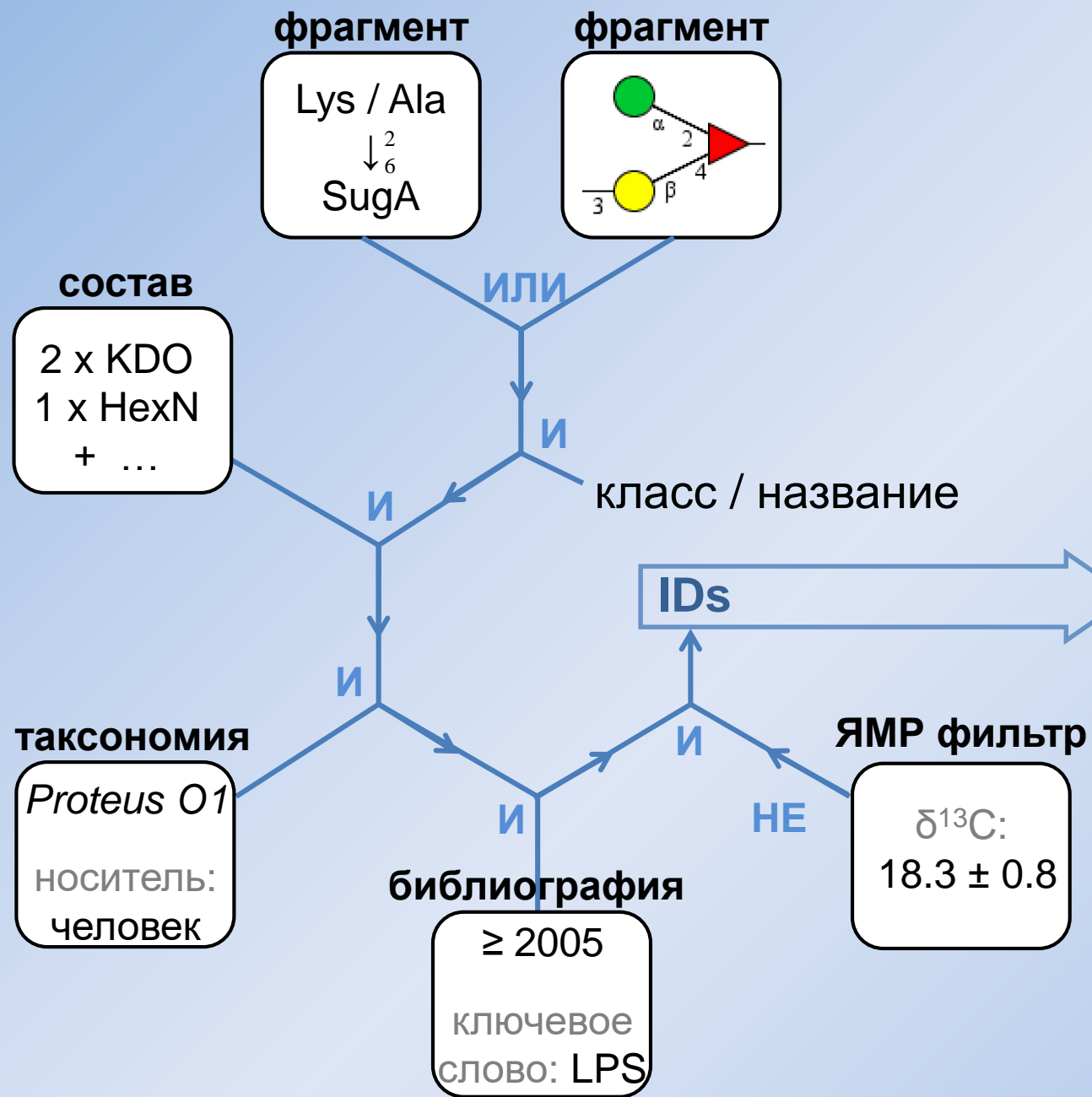
- регулярное пополнение
- расширяемая архитектура
- множество инструментов обработки данных
- проверка контента (15% = CarbBank, 85% = литература)
- полное покрытие (бактерии, грибы)
- интеграция с другими базами

Типичные запросы



→ однозначные переходы - - - - -> нечеткие переходы комбинации: **И**, **ИЛИ**, **НЕ**

CSDb: составной запрос



Данные сгруппированы по соединениям, публикациям, организмам и т.д.

Found 5 structures. Displayed structures from 1 to 5.

Expand all compounds Show all as text (Sweet)

1. Compound ID: 10502

Structure type: polymer chemical repeating unit
Compound class: O-polysaccharide, O-antigen

Structural formula & atomic coordinates
Sweet-II 3D model

The structure is contained in the following articles:

- Article ID: 4266
Boyko AS, Dmitrenko AS, Fedonenko YP
"O-polysaccharide of the lipopolysaccharide of the bacterium *Acetivibrio acotriose*" - Carbohydrate Research
- Article ID: 4833
Fedonenko YP, Boyko AS
"The review of the O-polysaccharides of the genus *Acetivibrio*" - Carbohydrate Research

Two types of neutral O-polysaccharide were isolated from the aqueous phenol-water extraction from the acetoacetate of the major O-polysaccharide was characterized by ¹H and ¹³C NMR spectroscopy. The D-rhamnose is indicated by italics.

Lipopolysaccharide, structural, O-polysaccharide, *Acetivibrio brasiliense*, O-antigen, D-Acofriose

NCBI PubMed ID: 22575749
Publication DOI: 10.1016/j.carres.2012.04.006
Journal NLM ID: 0043535
Publisher: Elsevier
Correspondence: room308@ibppm.sgu.ru
Institutions: Institute of Biochemistry and Physiology of Plants and Microorganisms, Russian Academy of Sciences
Methods: 1H NMR, 13C NMR, NMR-2D, methylation, chemical analysis, GLC, Smith degradation

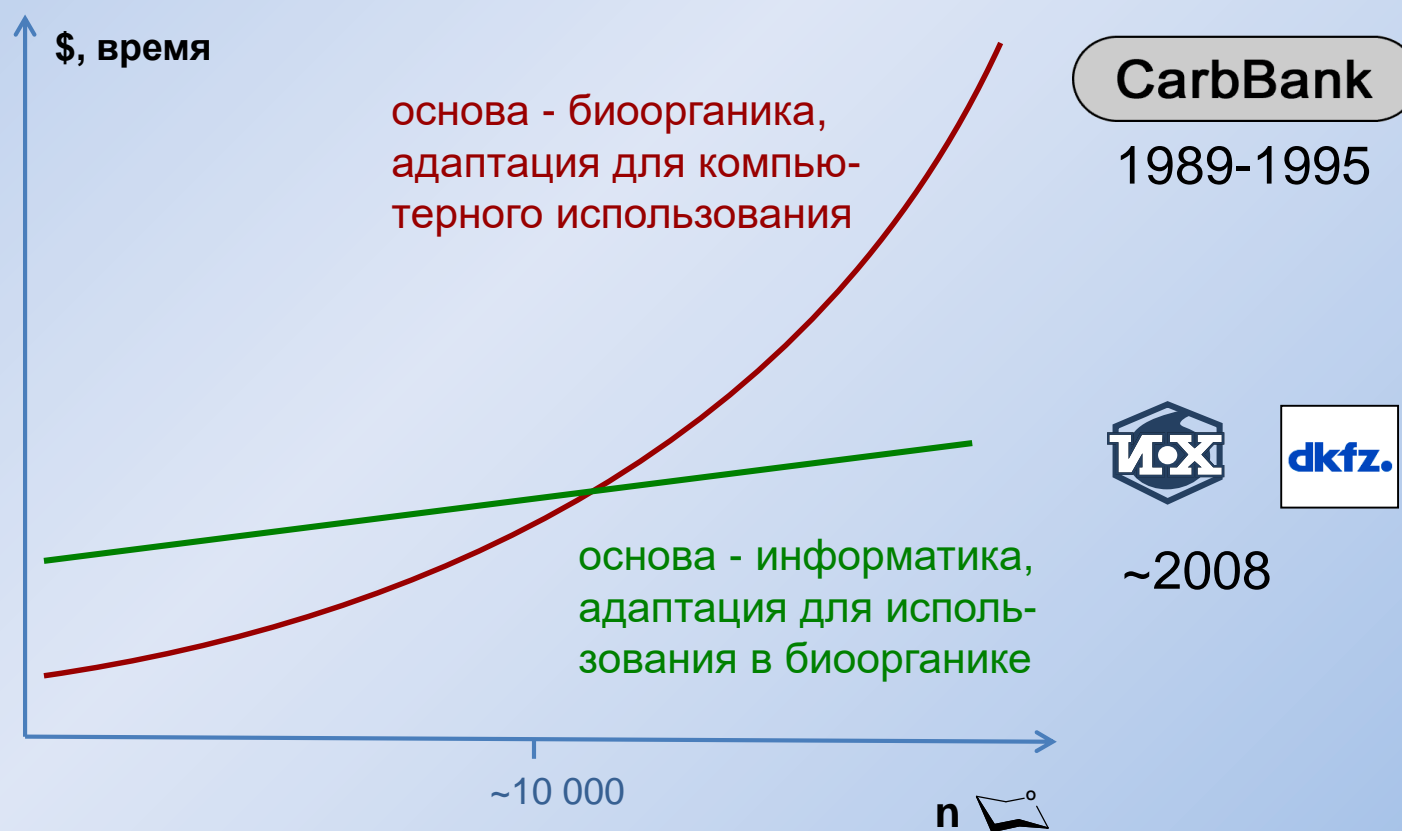
Acetivibrio brasiliense Jm6B2
CSDB ID 2831

Get MOL Spin rotate zoom move

The spectrum also has 2 signals at unknown positions (not plotted).

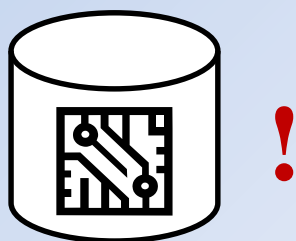
Организация разработки

- База данных и платформа должны быть построены по правилам информатики.
- Правила были конкретизированы и адаптированы для углеводов.

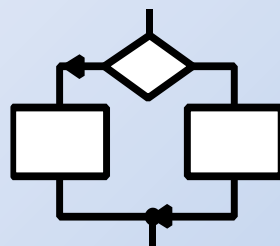


Критерии оценки

- **Функциональность** (типы данных и индексов, обработка запросов)
- **Полнота покрытия** (+ выбранный класс)
- **Качество данных** (% ошибок, прозрачность)
- **Интеграция** (поддержка форматов, импорт-экспорт, API, RDF)
- **Интерфейс** (простота, стабильность, производительность)



внутренняя
архитектура

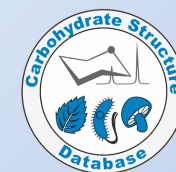


управляющие
программы



наполнение
данными

Архитектура



- Реляционная база данных
- Индексация + стандартные индексы (DOI, TaxID, ICD-11, PMID, ...)
- Структуры, таксономия, библиография – разные типы записей
- Человекочитаемый дамп (организация процесса наполнения)
- Контролируемые словари терминов (мономеров; MSDB)
- ~~Free text~~ ↗
- Таблица связности

минимум

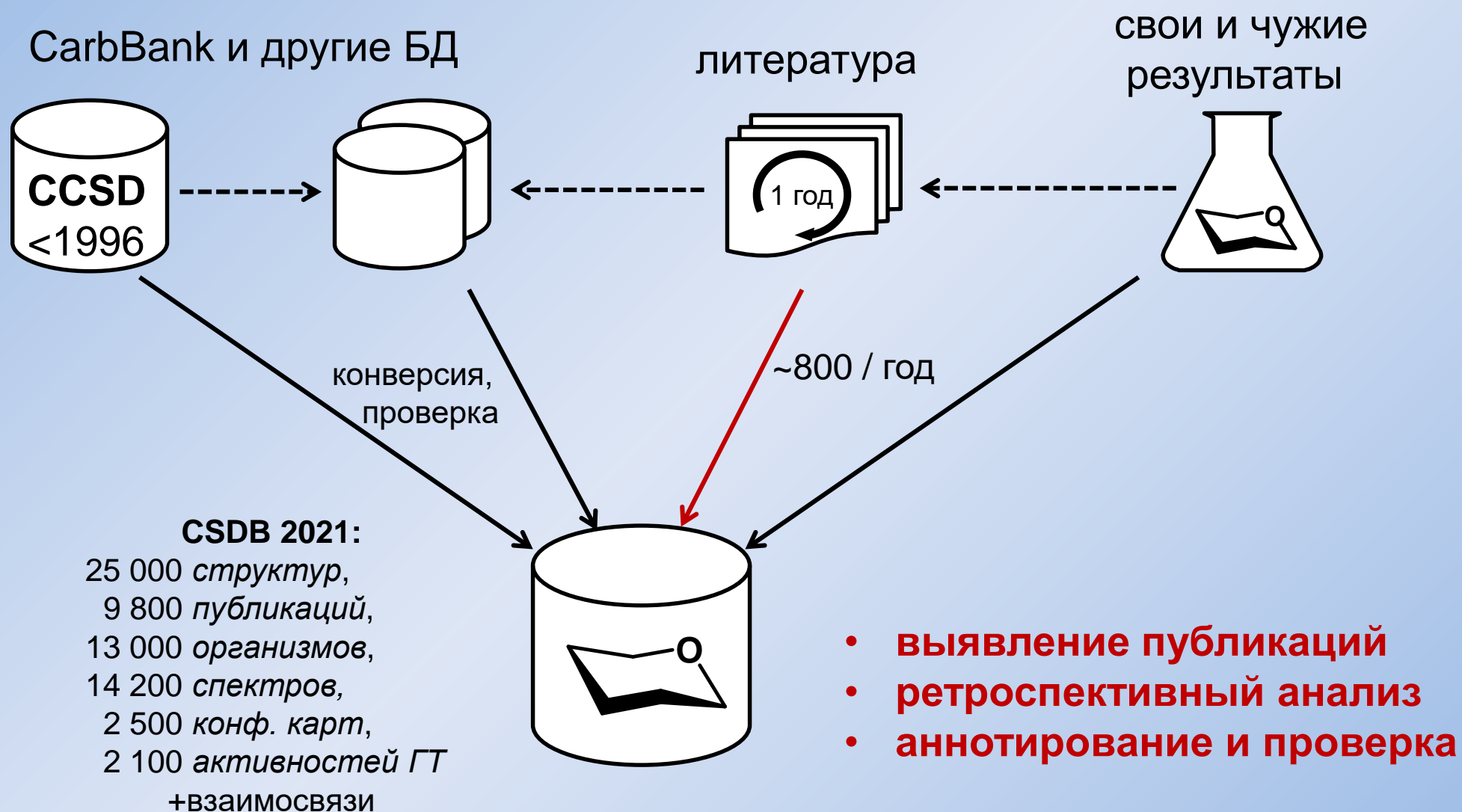
структура,
таксономия,
библиография,
ссылки на другие БД

дополнительно

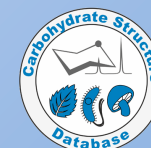
тривиальные названия,
спектры ЯМР, МС
условия съемки спектров,
биоактивность,
гены, ферменты,
конформация

болезни
органы, ткани
генотип, стадия
ключевые слова
рефераты
институты
методы

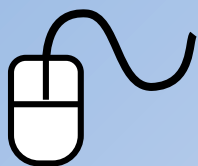
Источники данных



полное покрытие по прокариотам:
 отрицательный результат поиска =
 значимая научная информация



Качество данных



операторские

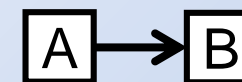


в базах

найдено &
исправлено



в статьях



в программах

ошибки, противоречия

исправляемые

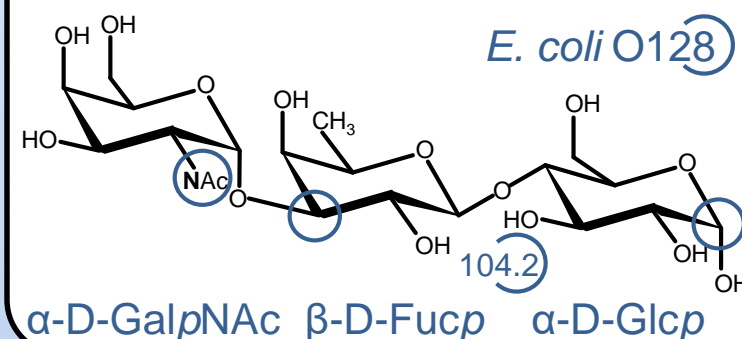
2dGlc → araHex,
α-Rib-ol → Rib-ol,
D-Kdo → Kdo,
1-methyl → 1-Me,
n.m.r. → NMR,
taxid 583 → Proteus,
...

выявляемые

Glc(1-2)Glc**N**,
anhydro-Kdo,
D-manHep,
Gal**p**5N,
Ac(1-2)[Glc(1-2)]Gal,
Escherichia sapiens,
Dev Food Sci 2012,
#Ac : 23 м.д., **65** м.д.
D-G**cl**, ...

невывявляемые

E. coli O127:
aDGalp**N**(1-4)bDFucp(1-4)**b**DGlc p
Glc C1 10**3.2** ppm



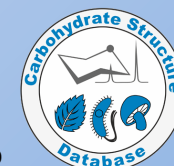
CarbBank



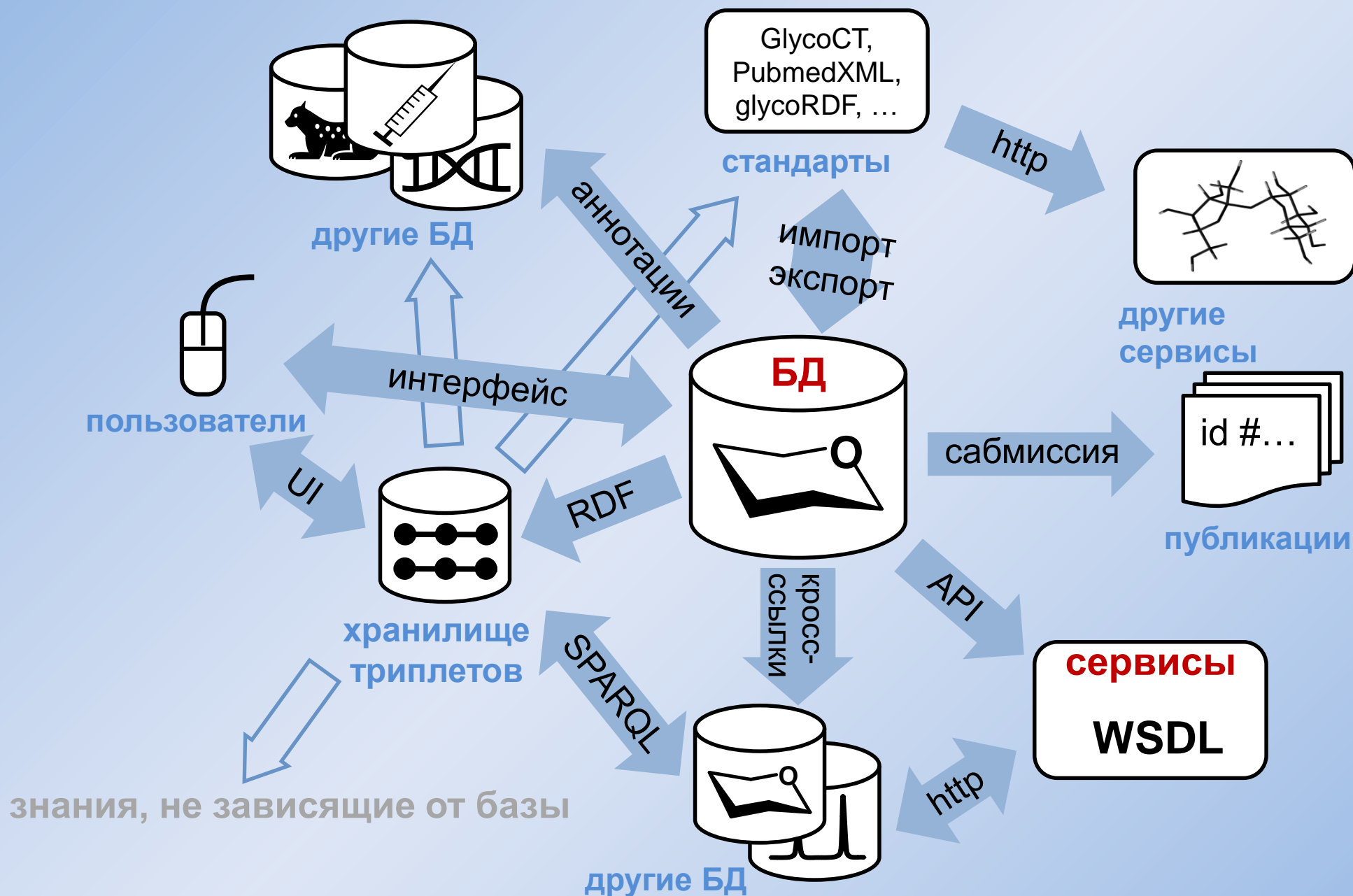
>50% (неправильные, отсутствующие, ложно присутствующие структуры, штаммы, аннотации)

↳ другие БД

😊 <10%



Идеальная интеграция



RDF – модель данных в виде триплетов *объект-предикат-субъект*.

- 😊 допускает распределенные запросы с минимальным знанием форматов баз
- ☹️ требует репозитория триплетов и согласованной онтологии

Задача: найти белок-носитель для произвольного гликана из JCGGDB.

Проблема: JCGGDB не связана ссылками с белковыми базами.

Преамбула:

Записи в JCGGDB имеют ссылки на идентификаторы в GlycomeDB.

Как GlycomeDB, так и UniCarbKB могут экспортировать структуры в формате GlycoCT.

Записи в UniCarbKB имеют ссылки в белковую базу UniProt.

Решение (*9-строчный скрипт на SPARQL*):

Сопоставить идентификаторы JCGGDB и UniCarbKB, используя GlycomeDB, и получить идентификаторы UniProt из UniCarbKB для каждого идентификатора JCGGDB.

Требуется:

стандартная онтология → экспорт данных в RDF → репозиторий триплетов → интерфейс SPARQL

GlycoRDF – первая формальная углеводная онтология (OWL)

Интерфейс

Конверсия данных ↔ другие форматы

Автоматические web-сервисы (WSDL)

Импорт, экспорт

Документация, HELP

Дружественность, быстрое действие

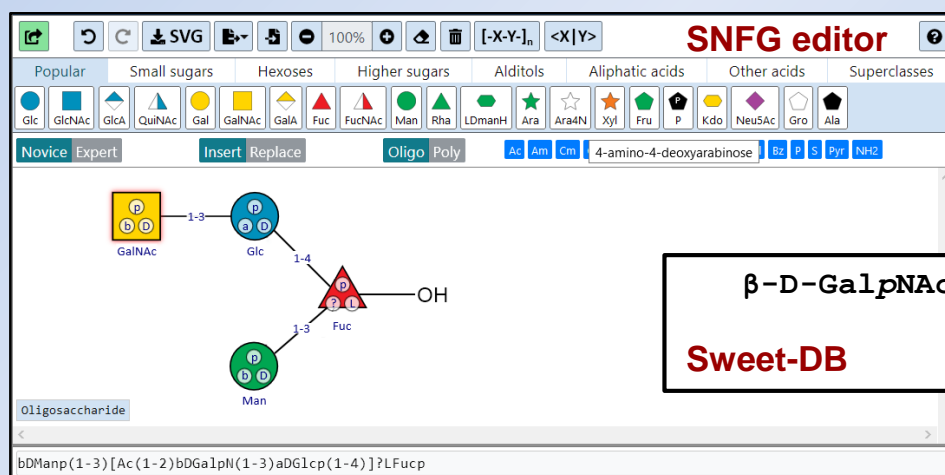
Ссылки на записи в других проектах (запросы, индексы, данные)

Ввод-вывод структур

SNFG,
WURCS,
GlycoCT,
SMILES,
MOL, PDB,
Glyde II,
LinUCS,
Sweet-DB
GLYCAM,
GlycoRDF,
PubMed XML,
DCI XML

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

веб-помощник,
граф. редактор,
библиотека структур,
CSDB Linear,
GlycoCT



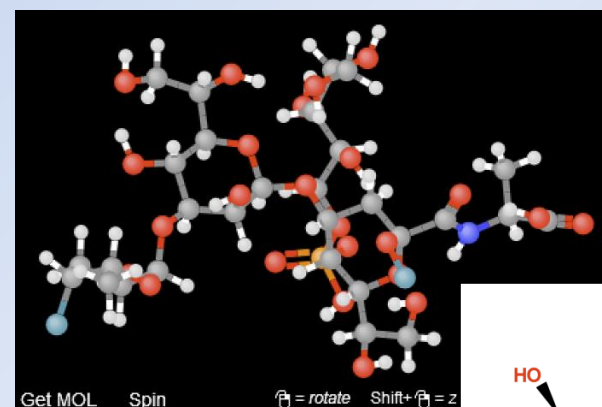
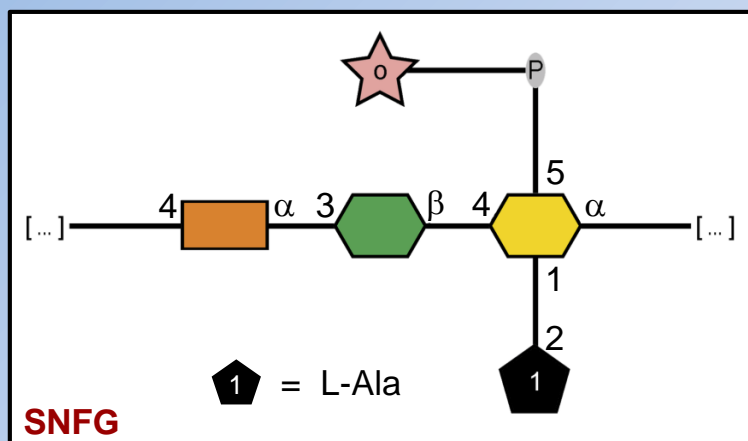
β -D-GalpNAc-(1-3)- α -D-G1cp(1-4)
Sweet-DB β -D-Manp-(1-3)-L-Fucp

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

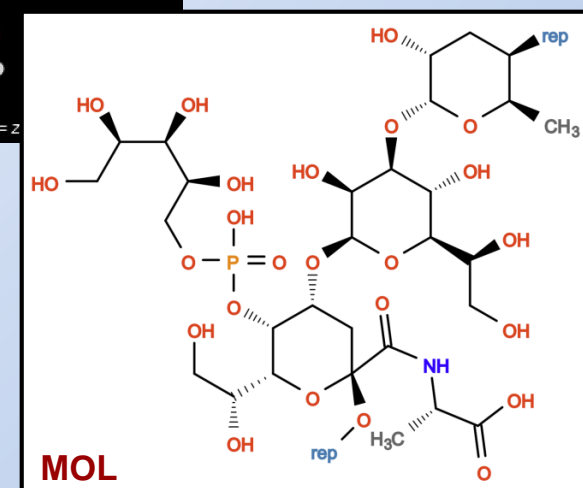
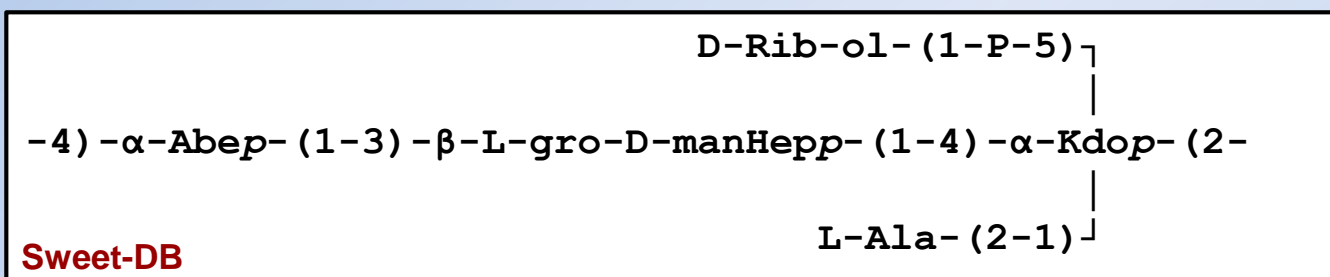
CSDB

Вывод структур

Визуализация в человекочитаемых форматах:



MOL (3D)



Экспорт в машиночитаемых форматах:

[*]O[C@]1(C(=O)N[C@@H](C)C(=O)O)C[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

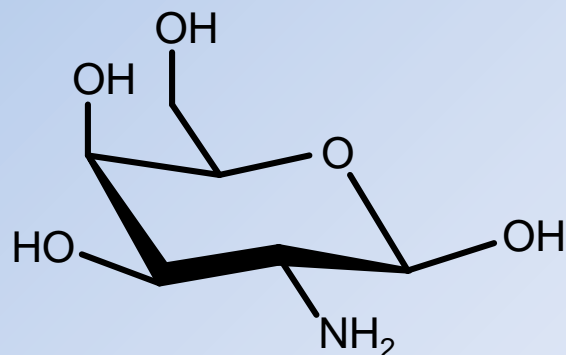
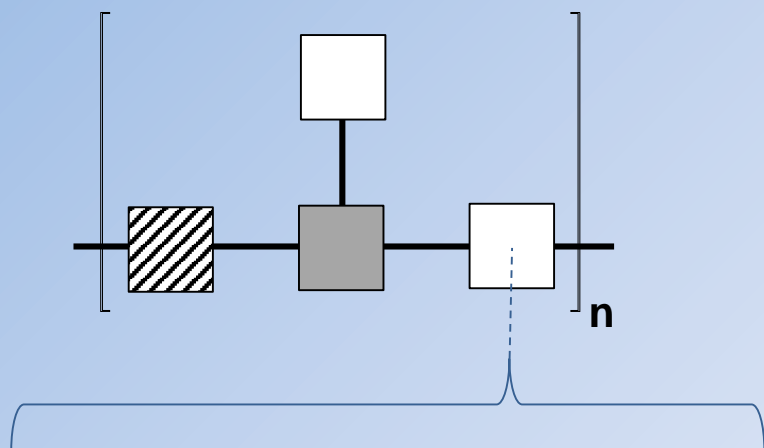
SMILES

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

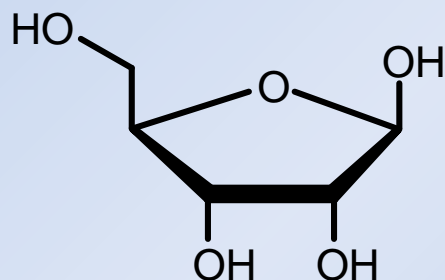
WURCS

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

CSDB Linear



пример альдо-пиранозы (β -D-GalpN)



пример альдо-фуранозы (β -D-Ribf)

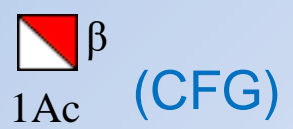
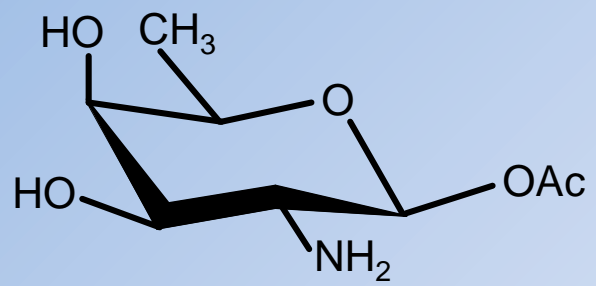
Полная структура

- мономерный состав, в т.ч. неуглеводный
- топология и последовательность
- позиции замещения
- стехиометрия боковых цепей
- число, границы и типы повторов

Структура остатка

- размер скелета (4-10)
- стереохимия всех центров (=мономер)
- способ циклизации (*p/f/a*, *aldo/keto*)
- аномерная форма (α/β)
- абсолютная конфигурация (D/L)
- модификации ($-\text{NH}_2$, $-\text{COOH}$, deoxy)

Неоднозначность номенклатуры



bDFucpN (1-1) Ac (CSDB) ←

однозначно соотнесено со структурой, но при этом понятно людям

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 (на естественных языках)

(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)

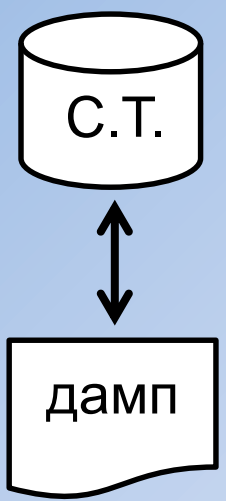
Почему не аналог MOL?

- Очень сложный перевод поатомного описания в семантическое (в структуру типа α -D-Galp-(1-3)- β -D-Glcp)
- Сложный перевод семантического описания в поатомное => трудоемкость аннотирования статей
- Нельзя описать структуры с неопределенностями
- Данные визуально не сопоставлены со знаниями
-  Нечеловекочитаемый → трудно курировать → ошибки в данных
- Координаты атомов не являются первичными данными но неполный MOL (без 3D) может быть воспринят как 3D MOL
- Громоздко для хранения и передачи в сети (и не передается как параметр в URL)

52.0606	6.3910	-0.1606	O	0	0	0
51.8591	8.6986	-0.1875	N	0	0	0
52.9844	9.0584	0.7259	C	0	0	1
53.8550	8.9929	0.0662	H	0	0	0
52.9684	10.5530	1.3121	C	0	0	2
52.2705	11.0903	0.6993	H	0	0	0
1117	1	0	0	0	0	
1118	1	0	0	0	0	
1119	1	0	0	0	0	

атомы, координаты, связность

Особенности структур и язык

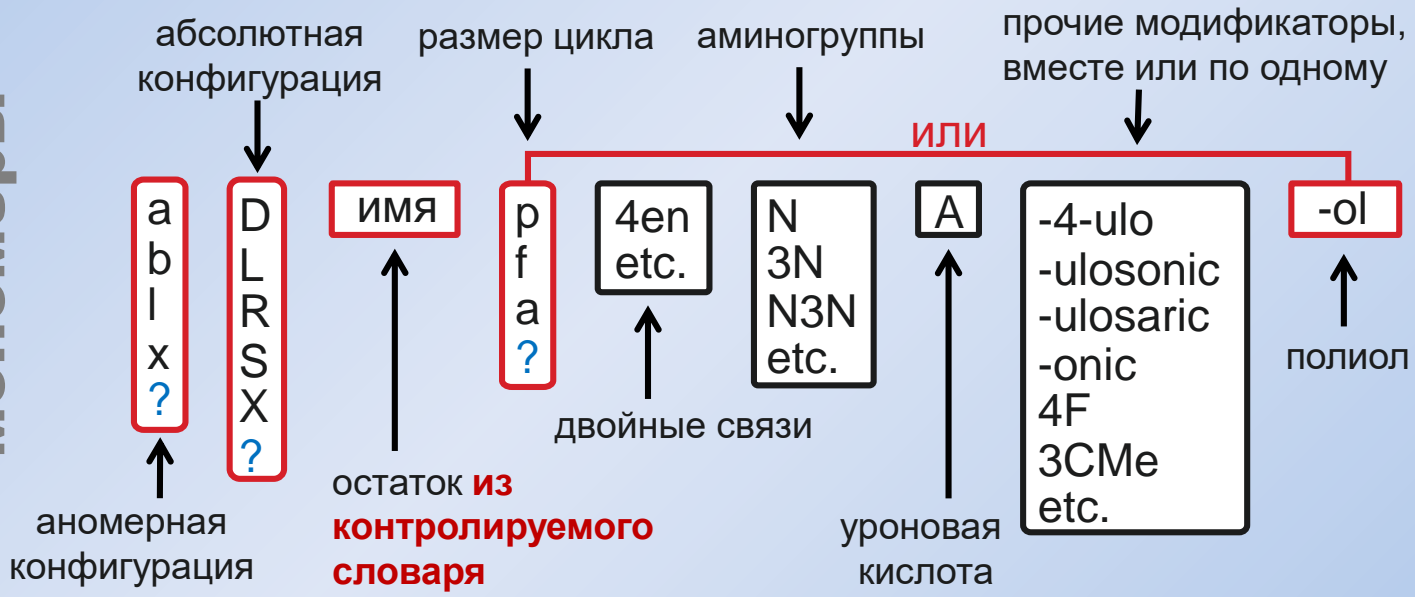


углеводный язык CSDB Linear

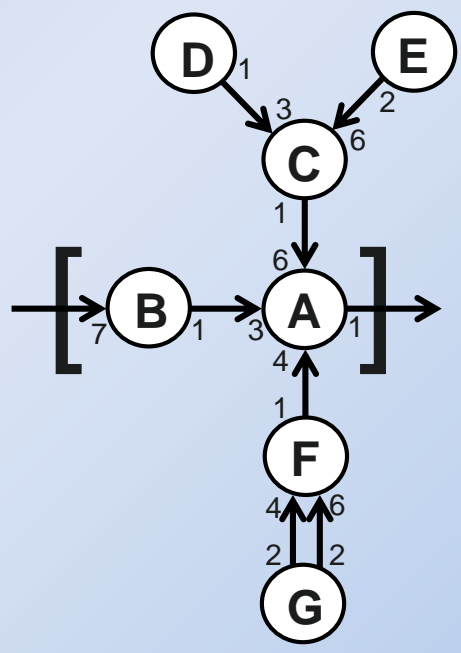
аннотаторы, операторы,
другие языки
(IUPAC, GlycoCT, WURCS, ...)

- полнота
- однозначность
- человекочитаемость
- машиночитаемость
- неточные структуры

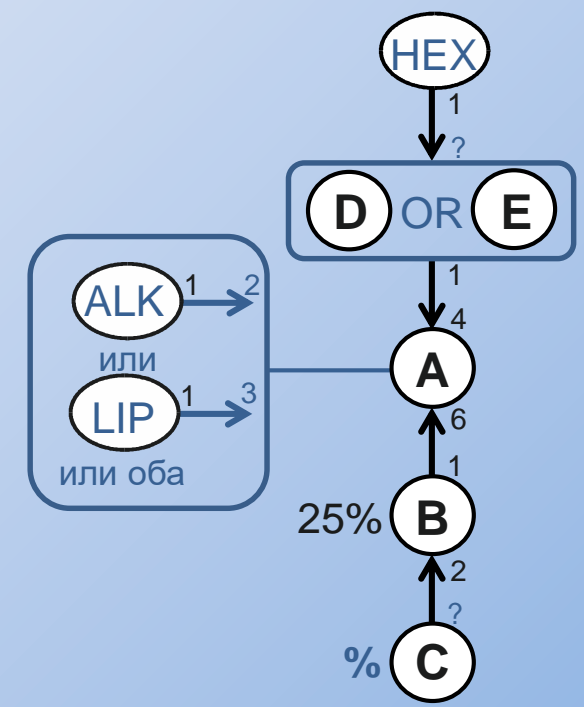
МОНОМЕРЫ



ТОПОЛОГИЯ



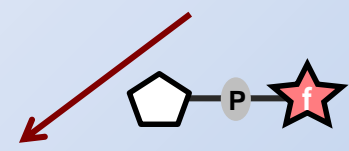
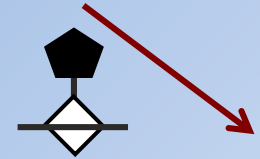
НЕТОЧНОСТИ



Уровни абстракции

-3) [xLAla (2-6) ,Ac (1-2)]bDGalpNA (1-
точный фрагмент и его связи

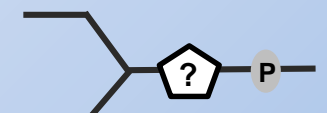
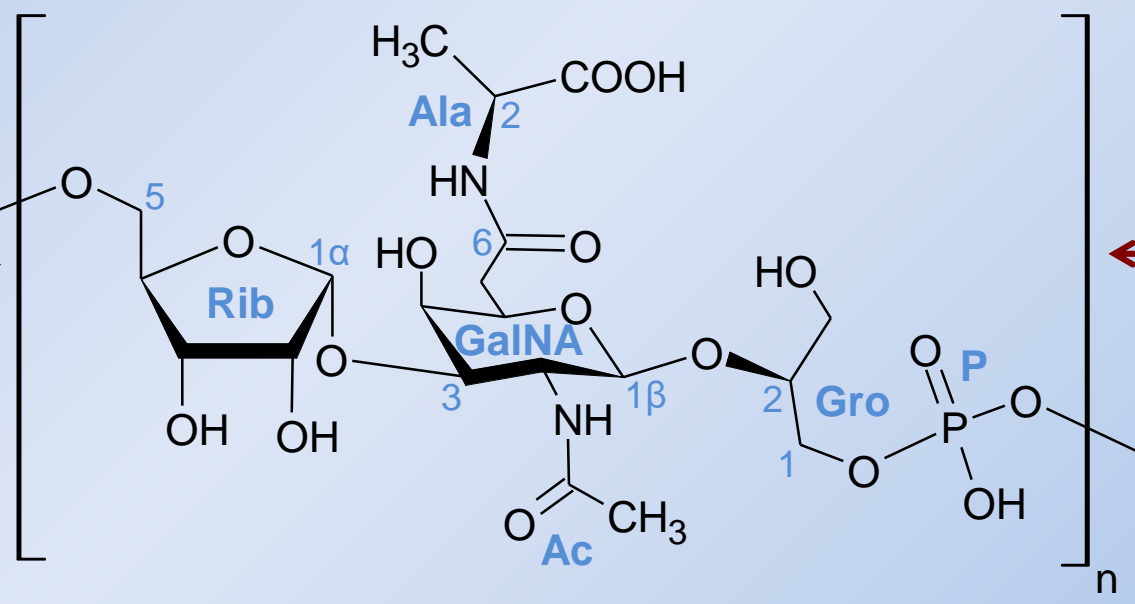
xDGro (1-P-5) aDRibf
не указана топология и места присоединения



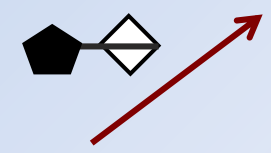
L-Ala (2-6) +

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

-?) ?Dhex (1-
указаны только класс и 1-связь



приблизительный мотив



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


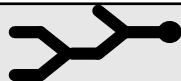








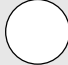


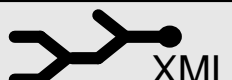





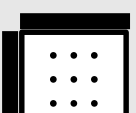




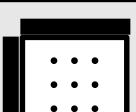

























HEX, xDRib?, PEP

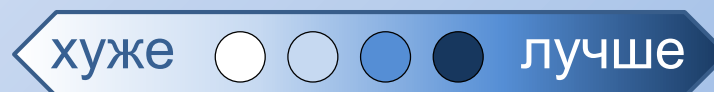
не указаны конфигурации, позиции замещения, размер цикла, N-ацетилирование

частичный состав

Сравнение языков

	<i>подход</i>	полнота	однозначность	контроль	парсинг	неточные структуры
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	

 СОВМЕСТИМОСТЬ

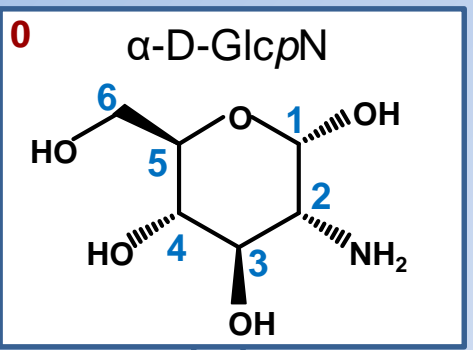


От семантики к атомам

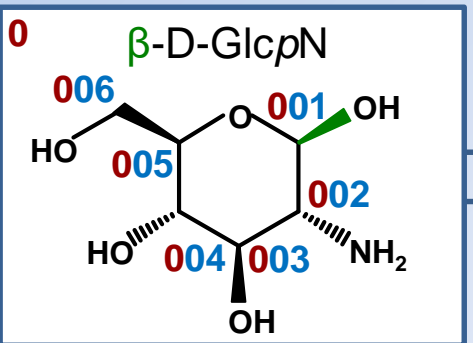
анализ CSDB Linear или другой нотации



SMILES прототипов

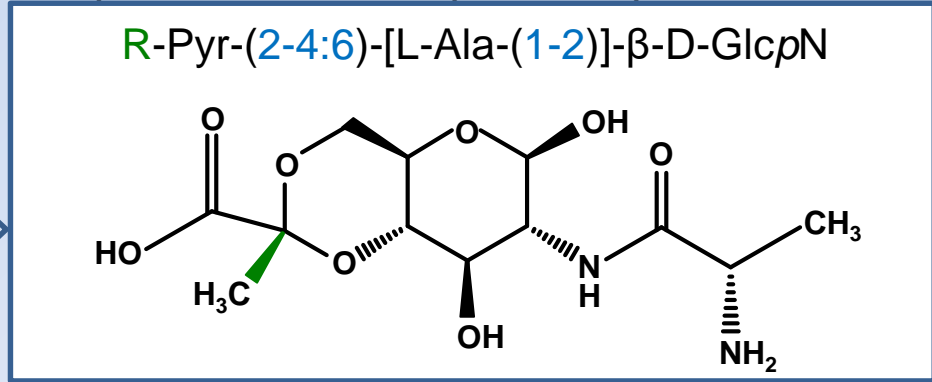


стереоконфигурации + метки



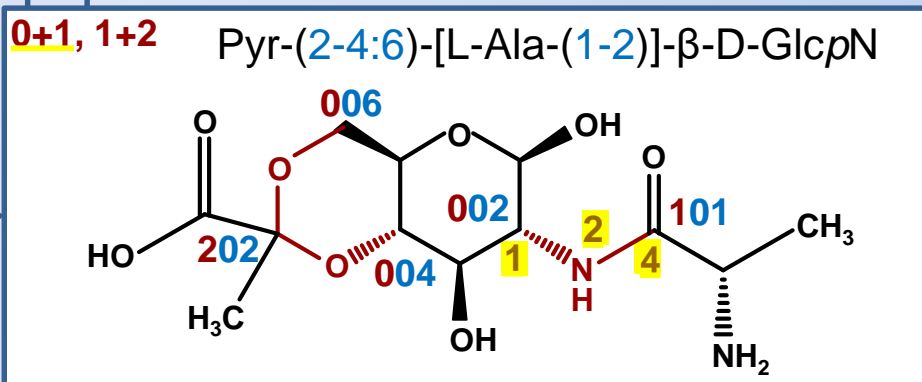
```
[006CH2](O)[005C@@H](O1)
[004C@@H](O)[003C@H](O)
[002C@@H](O)[001C@H](O)1
```

сборка + новые стереоцентры



SMILES: C[C@H](N)C(=O)N[C@H]1[C@@H](O)O[C@@H]2CO[C@@](C)(C(=O)O)O[C@H]2[C@@H]1O

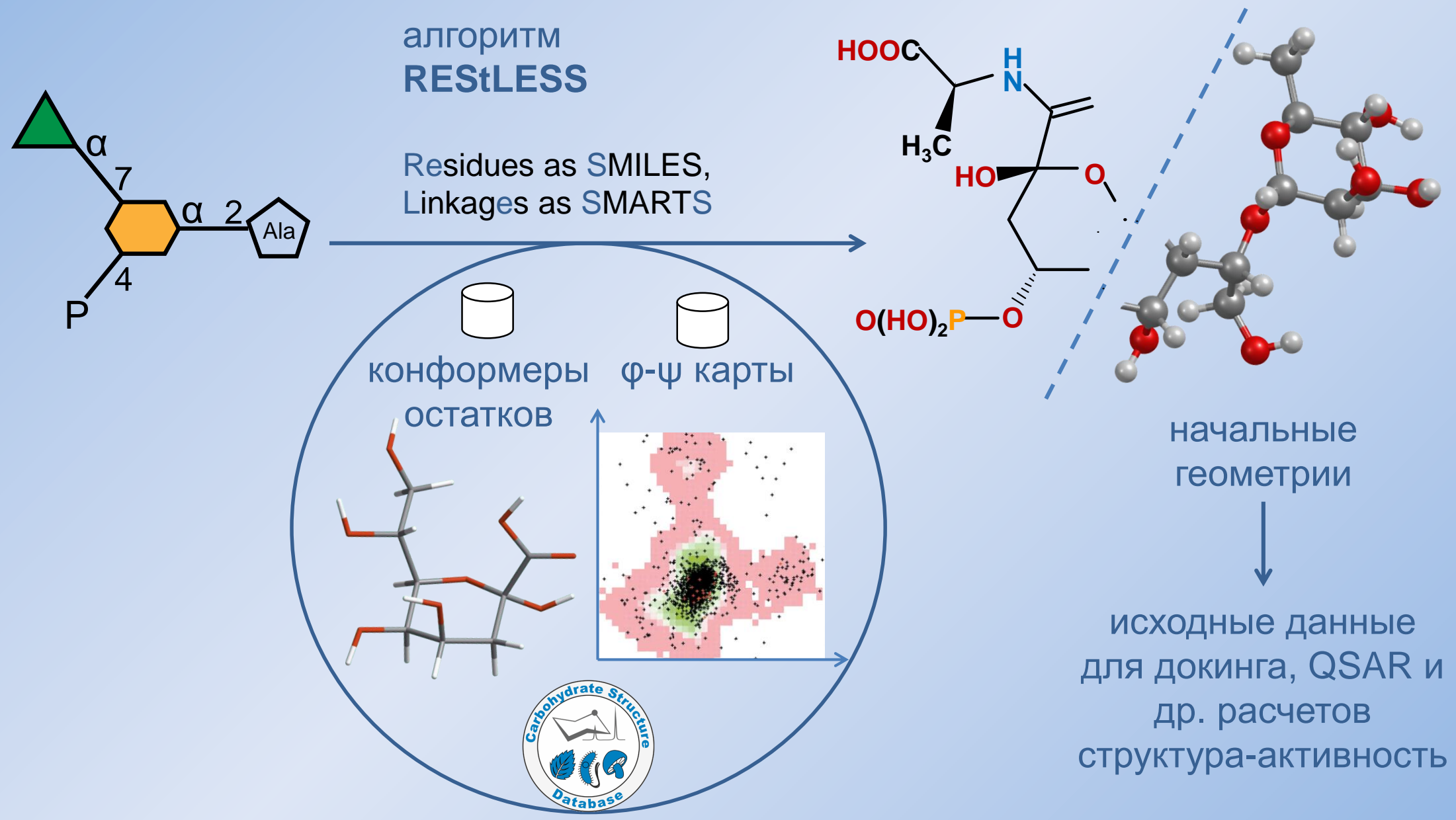
SMARTS связей между остатками



```
[202*:1]=[O:2].[O;!H0:3][006*:4]>>[*:2][*:1][*:3][*:4]
[202*:1][O;!H0:2].[O;!H0:3][004*:4]>>[*:1][*:2][*:4]
[002*:1][N;!H0:2].[O;!H0:3][101*:4]>>[*:1][*:2][*:4]
```

Любые операции на атомарном уровне + Внешние химические программы

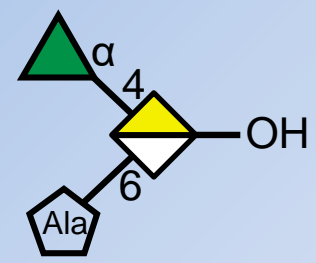
Молекулярная геометрия



Моделирование конформеров

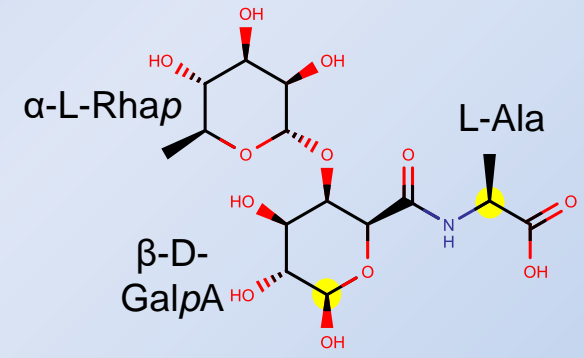
aLRhap(1-4)[x?Ala?(2-6)]?DGalpA

структуры,
в т.ч. неполные

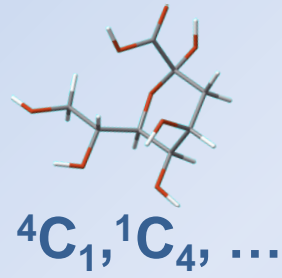


другие варианты
(α-GalA, D-Ala, и т. д.)

SMILES

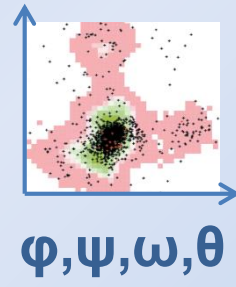


выгодные
конформеры
~1000 остатков

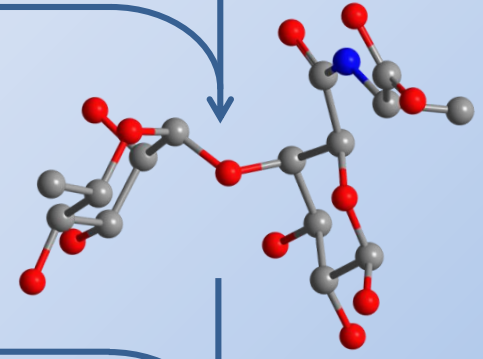


«креслификация»

заселенные
состояния
мостиков

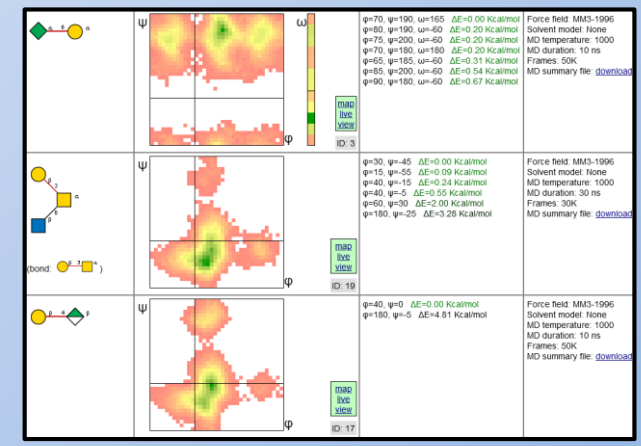


выбор минимумов
ММ-релаксация

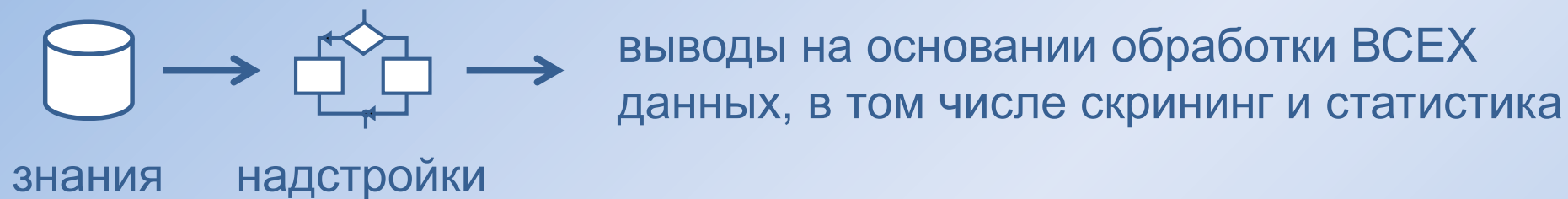


мол. динамика
300К, 100нс, H₂O

конформеры
+ их энергии

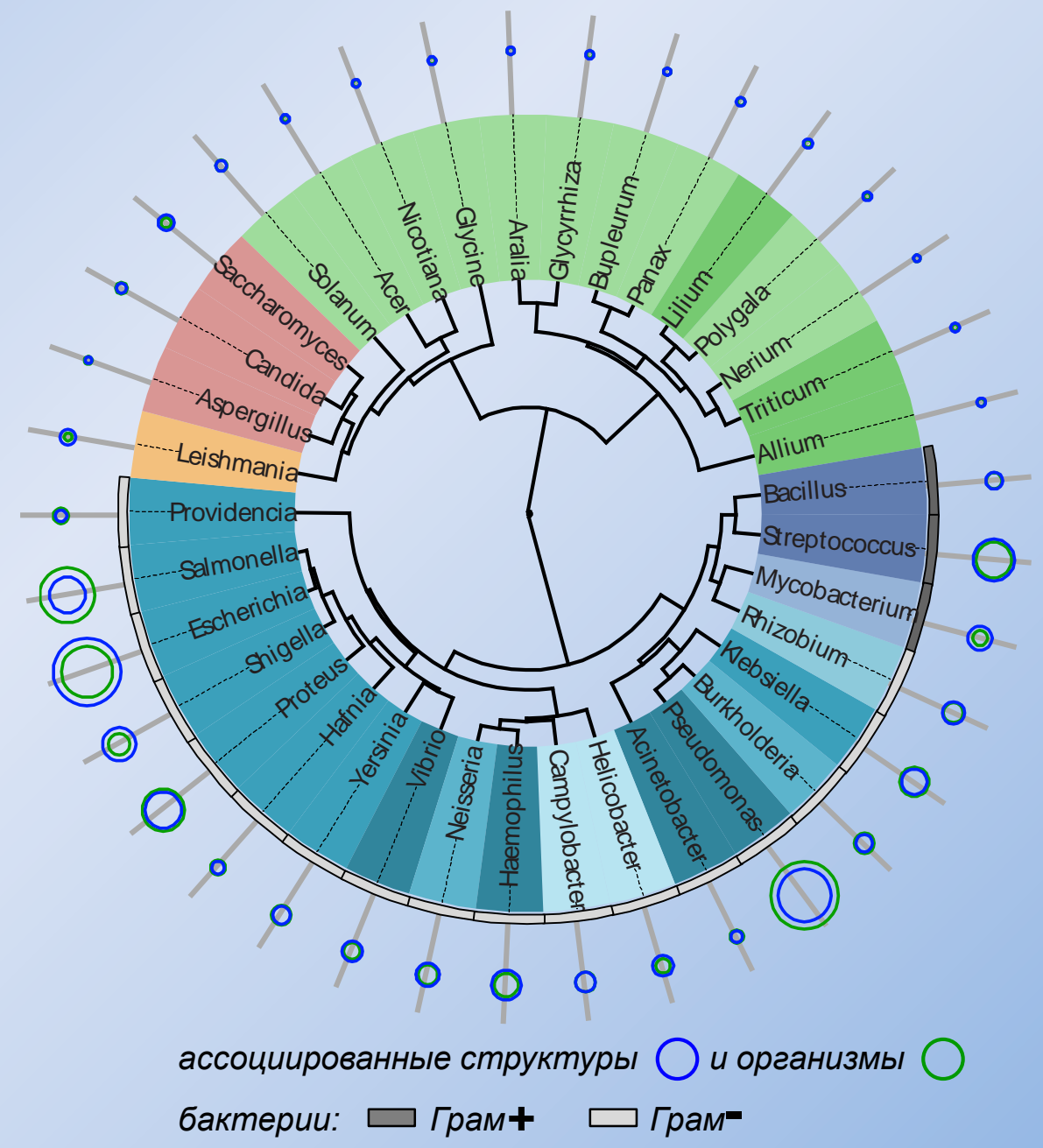
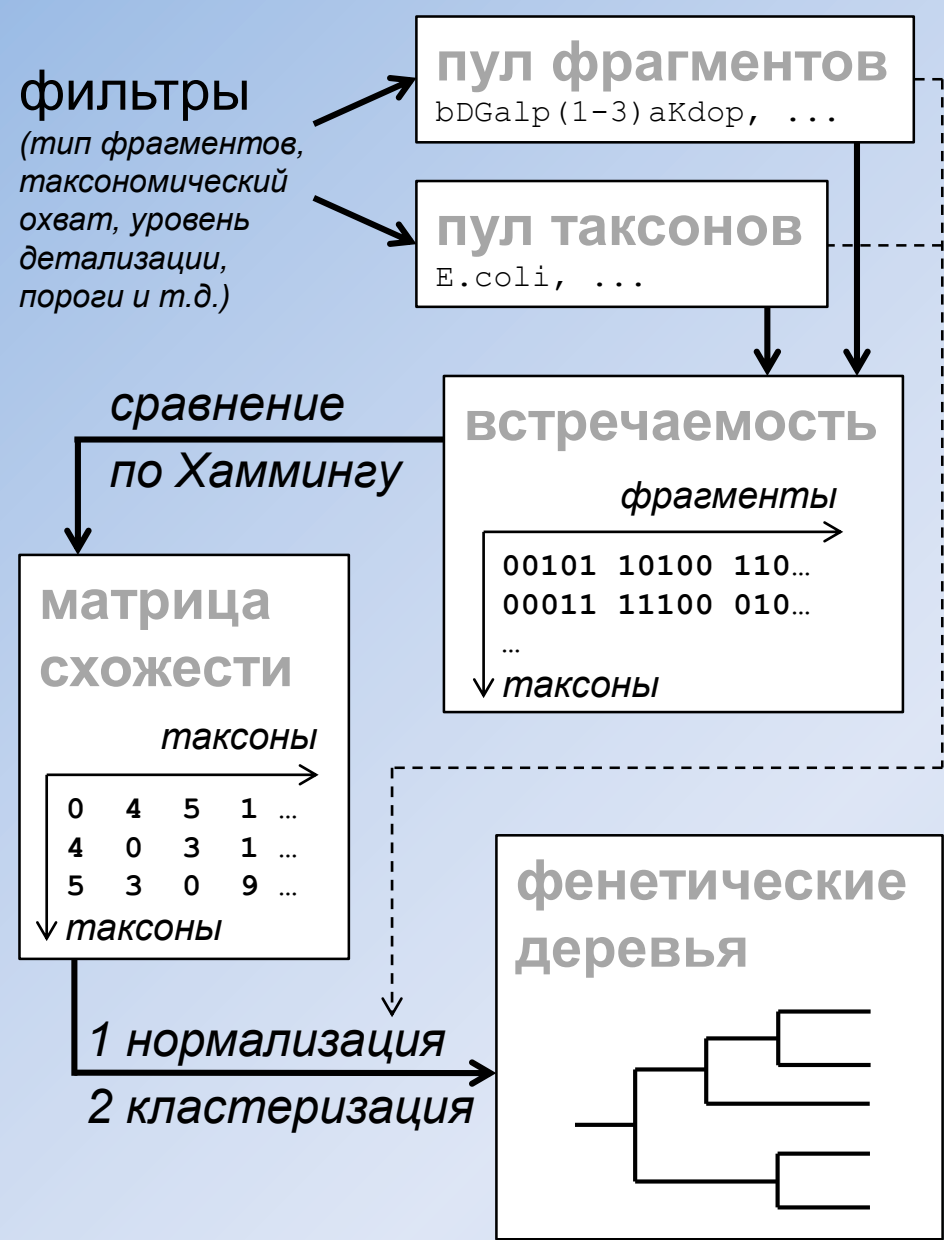


Надстройки



- Анализ путей биосинтеза (база гликозилтрансфераз)
- Конформационные карты олигосахаридов
- Предсказание и отнесение спектров ЯМР ^{13}C , ^1H , 2D
- Предсказание структуры по спектрам и другим данным
- Кластеризация таксонов на основании их гликомов
- Распределение фрагментов по таксонам и положению в структурах
- Классификация мономеров

Кластеризация таксонов



Гликозилтрансферазы

Критерии:
(в любом сочетании)

- идентификаторы в базах
- название фермента / группа
- название гена / кластер
- семейство CAZy
- организм (вид, штамм)
- синтезируемая связь
- донор (или его фрагмент)
- акцептор (или его фрагмент)
- роль объекта в клетке
- уровень достоверности

71B1
Uniprot
Q9LSY9.1
Genbank
821729

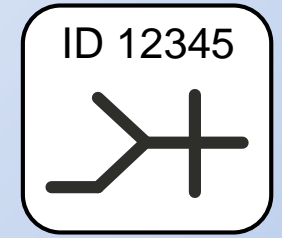
идентификаторы



Объект:

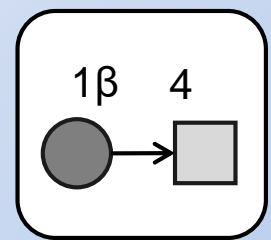


организм,
орган, ткань

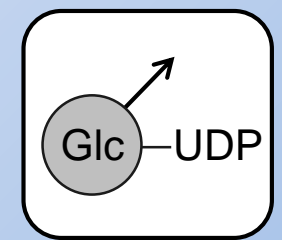


полная
структура

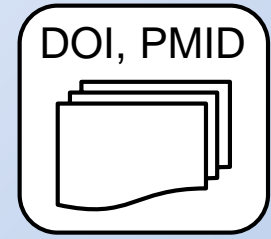
Активность:



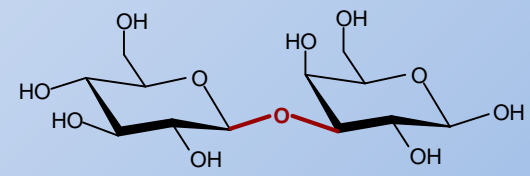
синтезируемый
фрагмент



донор и
акцептор



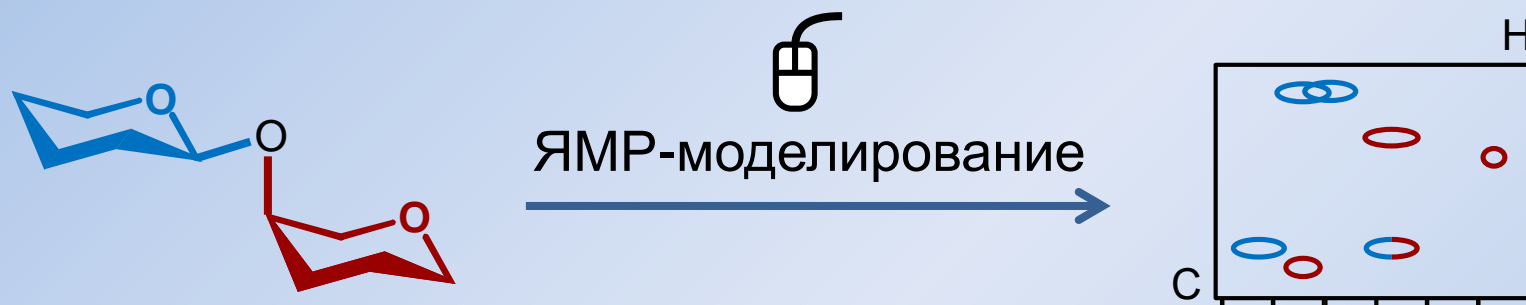
публикации



CSDB GT

ЯМР-моделирование

ЯМР – основной метод анализа первичной структуры в гликобиологии



- Помощь эксперту в установлении структуры
- Отнесение сигналов и проверка гипотез
- Перебор структур и сравнение симуляции с экспериментом
- Верификация расчетной геометрии молекул

Статистически

химические сдвиги ^{13}C и ^1H

- ☺ Существует база (CSDB)
- ☺ Отслеживается до источников
- ☹ Медленно (~минуты)

адаптация для углеводов:
GODDESS

Эмпирически

химические сдвиги и КССВ

- ☺ Очень быстро (~миллисекунды)
- ☹ Требуется модель
- ☹ Нужны специальные базы

адаптация для углеводов:
GODDESS, BIOPSEL, CASPER

Квантово-механически

геометрия + все параметры ЯМР

- ☺ Не зависит от базы
- ☹ Низкая точность (>3 м.д.)
- ☹ Очень медленно (~месяцы)

адаптация для углеводов: не нужна

Иначе

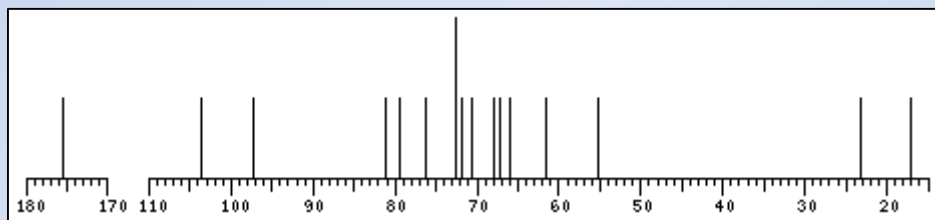
(нейронная сеть, регрессия, ММ-QM)

- ☹ Интегрально малопригодно

адаптация для углеводов: нет

- Использует регулярно обновляемую базу CSDB (>9000 спектров)
- Обобщает химическое окружение предсказываемого атома, пока не будет найдено достаточно похожих фрагментов в базе
- Работает на уровне групп атомов и специальных дескрипторов
- Оценивает достоверность и позволяет отследить источники
- Смешивает эмпирические и статистические результаты

применимо к предсказанию любых атомарных свойств, представленных в базе



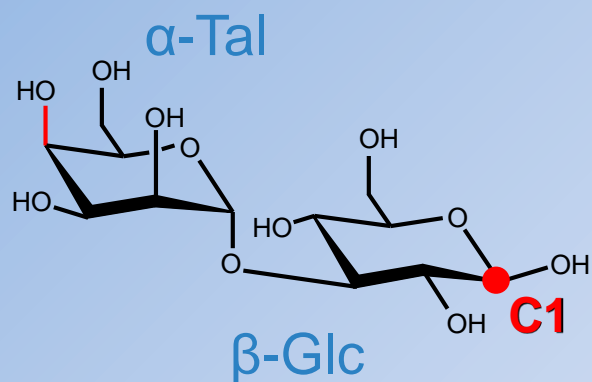
β -D-GlcpNAc-(1-3)- α -D-Fucp-(1-P-3)-D-Gro

¹³C NMR data: Solvent: Water (H or D) Recalc

Linkage	Residue	C1	C2	C3	C4	C5	C6	Accuracy
0	xDGro <i>trustworthiness-></i> 0 <i>NMR references-></i> 19	67.3 0	72.6 0.58	66.0 2.93				1.17
3	xXP?	-						
3,0	aDFucp <i>trustworthiness-></i> 2.83 <i>NMR references-></i> 3	97.2 2.83	72.5 2.62	79.5 3.8	70.6 1.61	68.0 3.39	17.2 3.63	2.98
3,0,3	bDGlcP <i>trustworthiness-></i> 3 <i>NMR references-></i> 1	103.7 3	55.1 3	81.2 3	71.9 3	76.2 3	61.5 3	3
3,0,3,2	Ac <i>trustworthiness-></i> 3.65 <i>NMR references-></i> 19	175.4 3.65	23.3 3.71					3.68

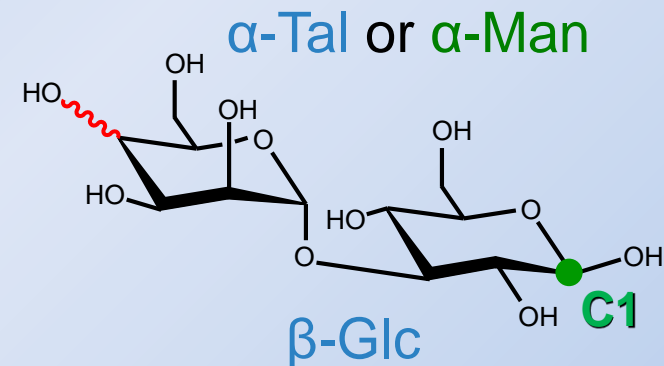
Обобщение структуры

фрагмент структуры



пример обобщения
конфигурация C4 донора

обобщенные фрагменты



типы заместителей
(X = OH, NH₂, H, ...)

R/S всех
центров

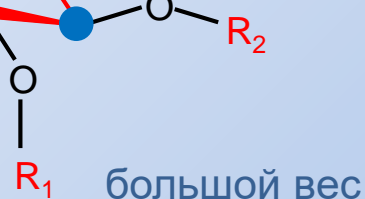
циклизация
(p, f, a)

обобщаемые
дескрипторы

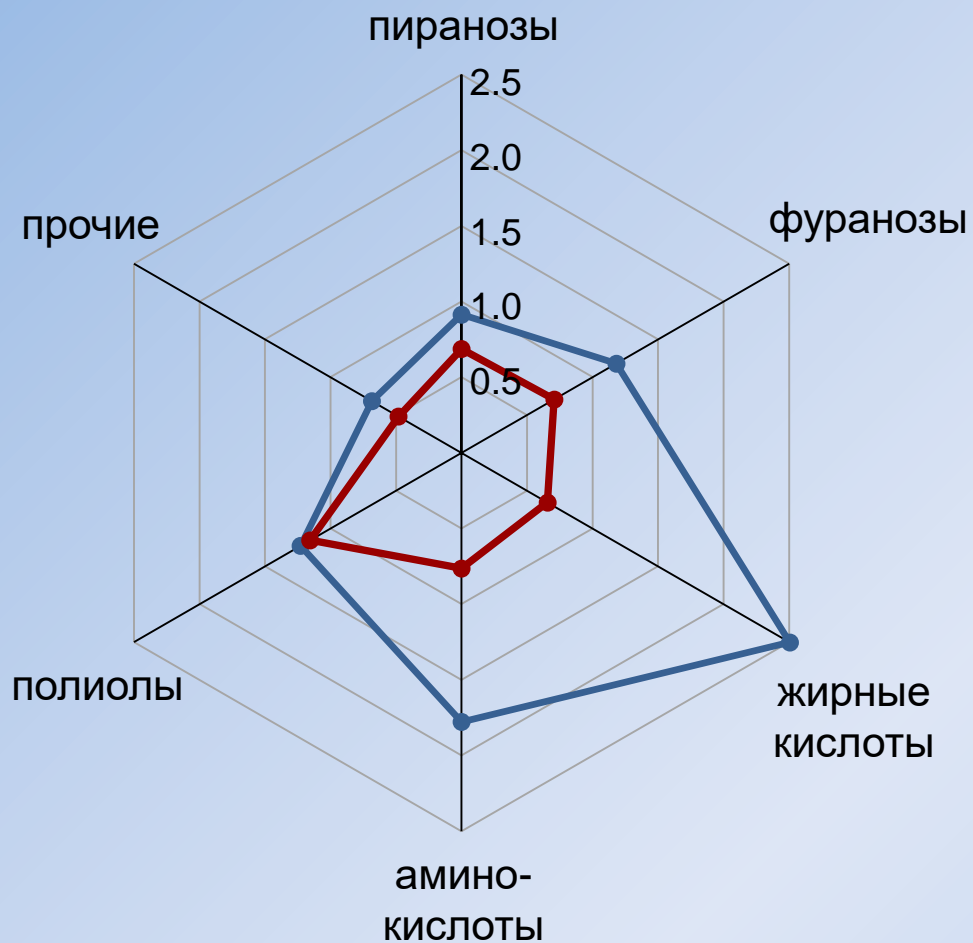
малый вес
для ● (далеко)

средний вес

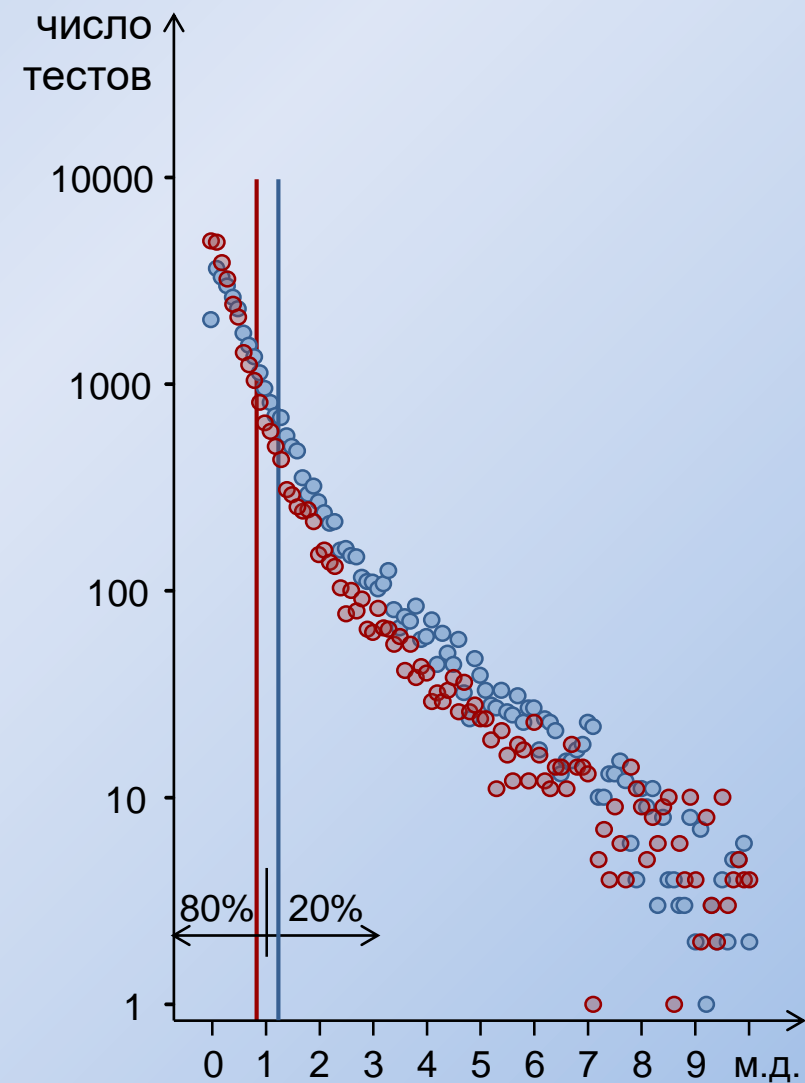
соседние фрагменты
(R_n) и их положения



каждый олигомерный фрагмент обобщается от малых весов к большим,
пока в базе не найдется достаточно обобщенных фрагментов,
затем данные из базы усредняются с учетом выбросов



средняя точность по классам, м.д.

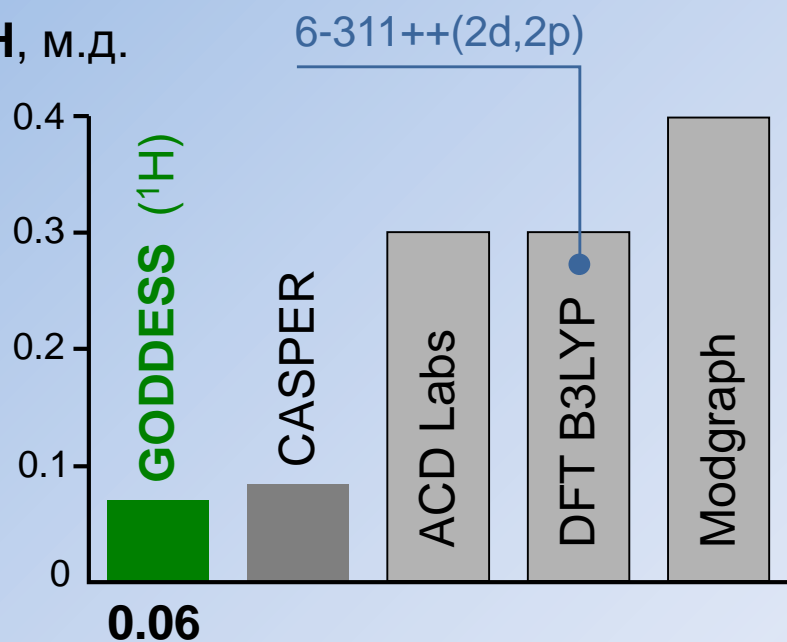
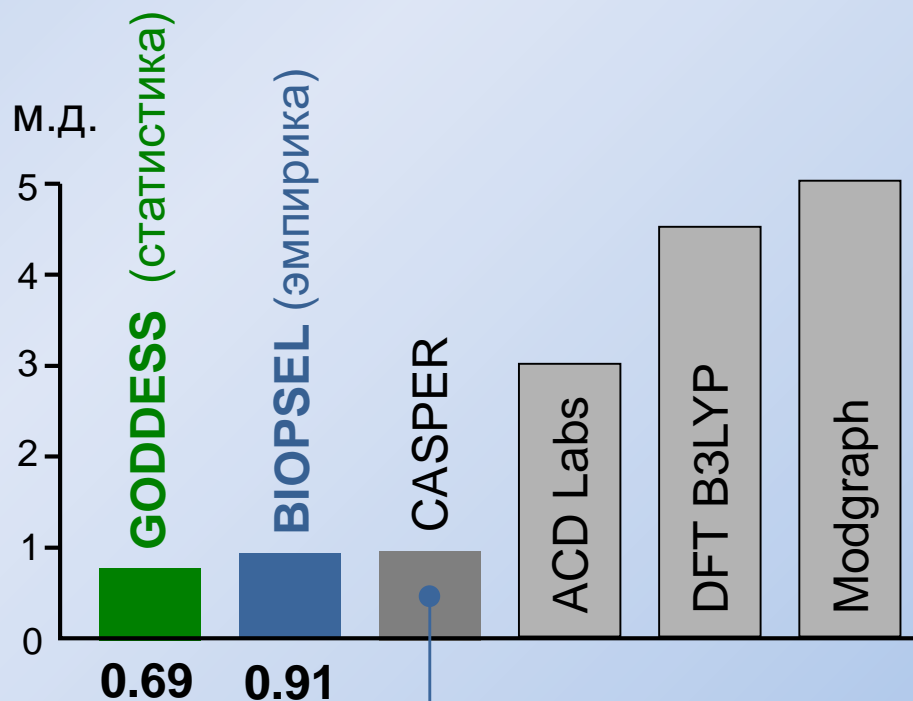


ошибка моделирования
химического сдвига

^{13}C - эмпирика

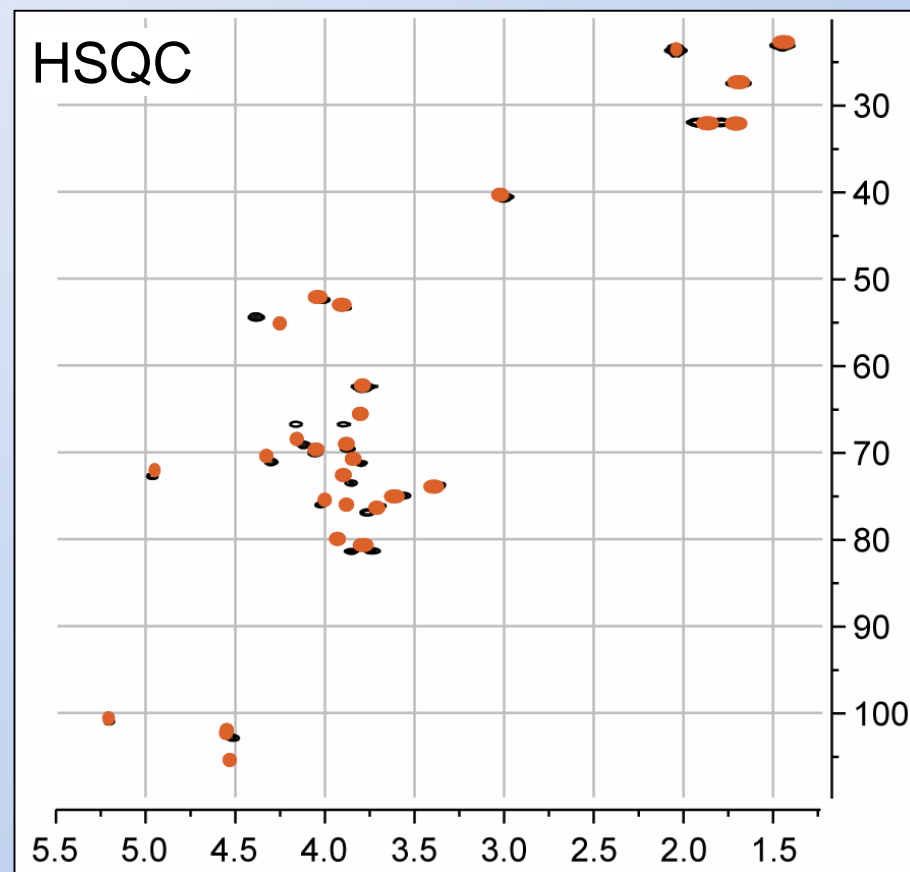
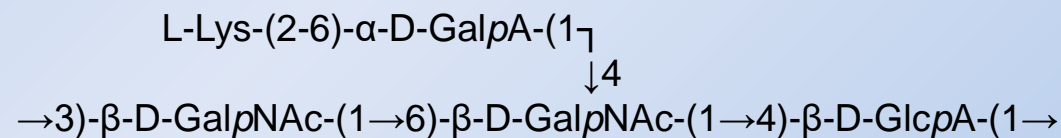
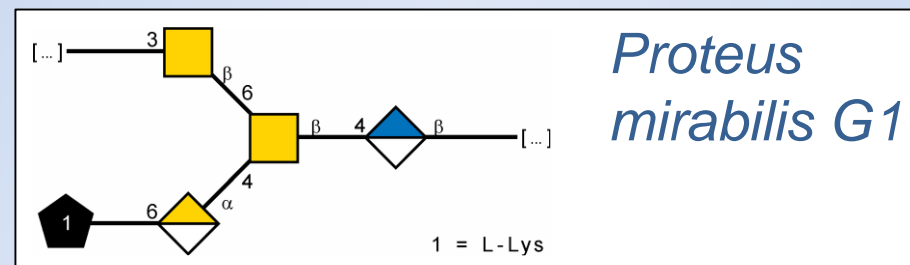
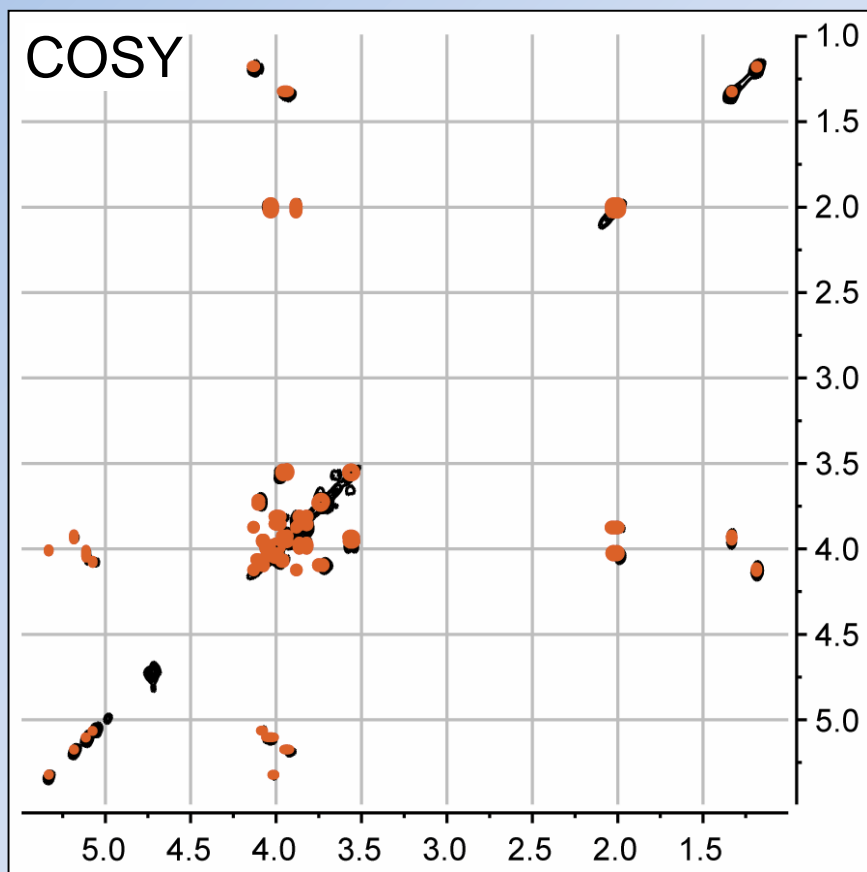
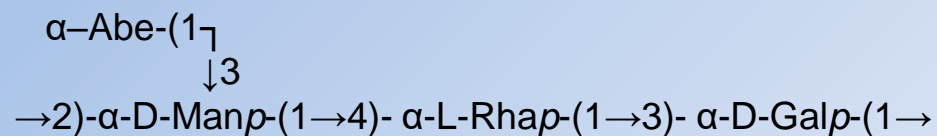
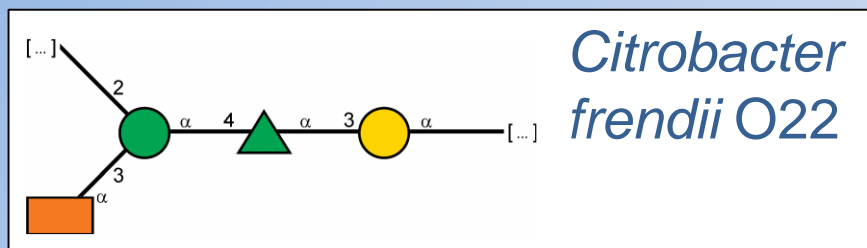
^{13}C - статистика

Сравнение подходов

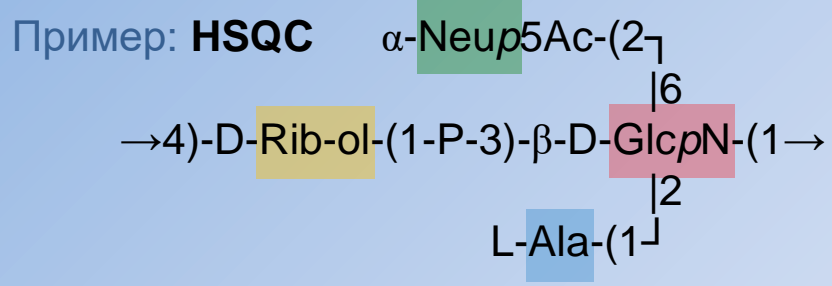
 $\Delta^1\text{H}$, м.д. $\Delta^{13}\text{C}$, м.д.

только гликозидные связи и
распространенные мономеры

сравнение «предсказание-эксперимент» проведено на
~32000 химических сдвигов из всех классов биогликанов



Визуализация результатов



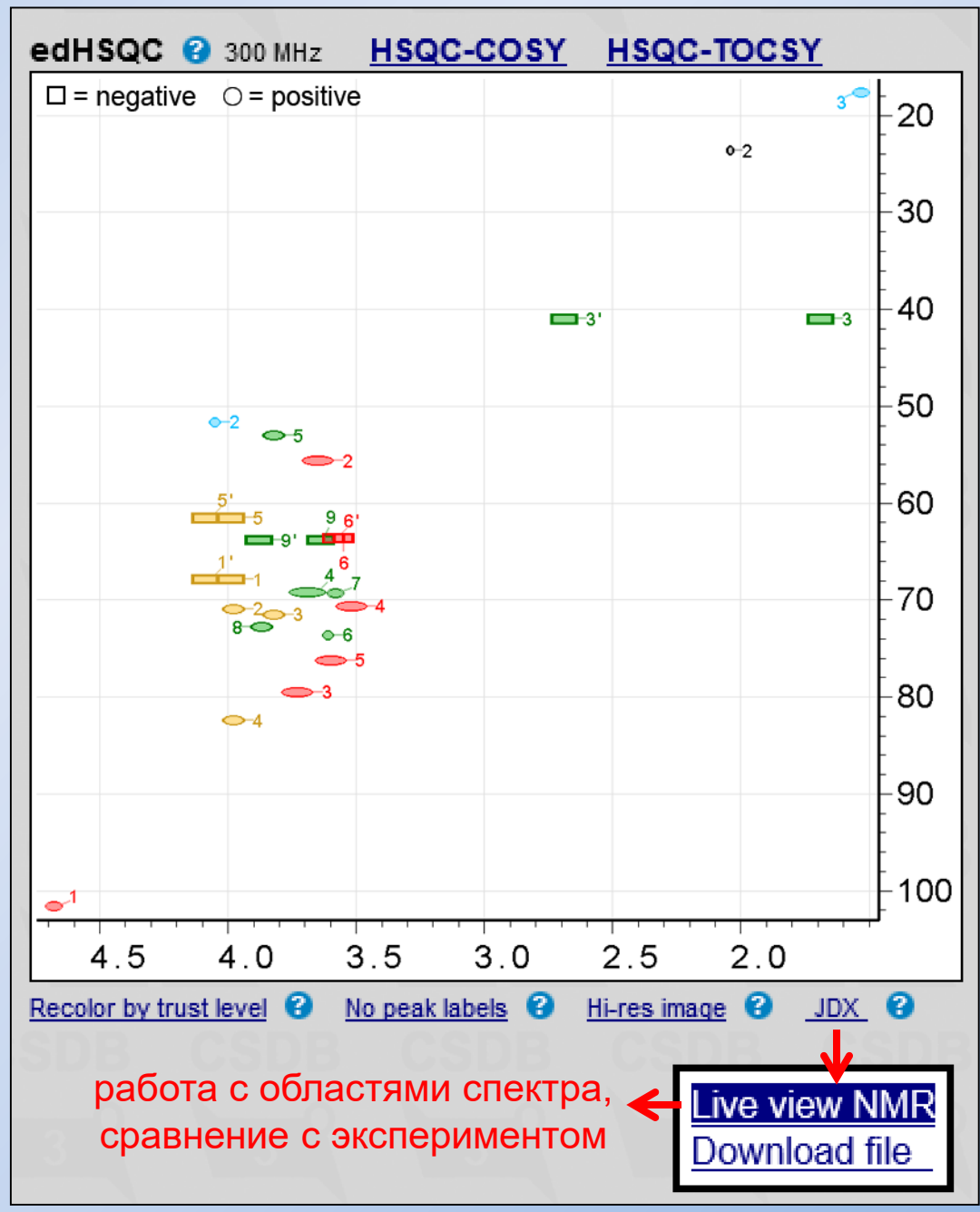
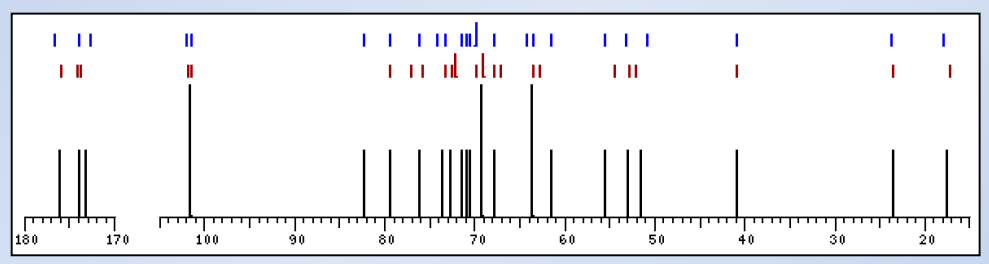
Поддерживается:

- 1D:** ¹³C BB, таблицы отнесения ¹H и ¹³C
- ¹H-¹H:** COSY, TOCSY, DQF COSY, COSY RCT
- ¹H-¹³C:** edHSQC, HMBC, HSQC-COSY, HSQC-TOCSY

эмпирическая оценка KCCB
+ частота спектрометра → ширина кросс-пигов

Скоро: NOESY / ROESY

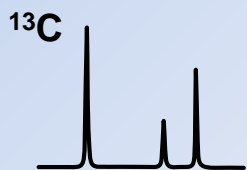
На выходе: таблицы, изображения, «живая» работа в браузере, экспорт в CSV и Jcamp-DX



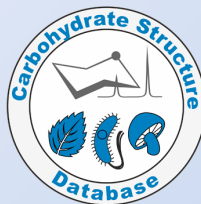
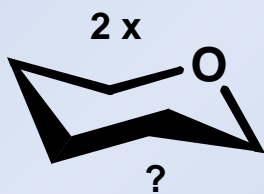
Итератор структур



неотнесенный
спектр ЯМР

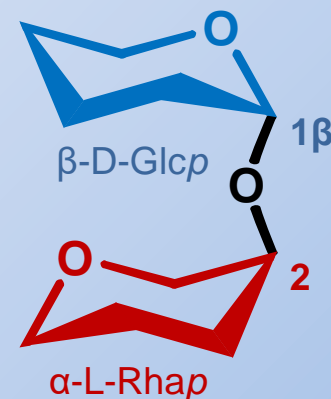


известные
данные
о структуре



алгоритм GRASS

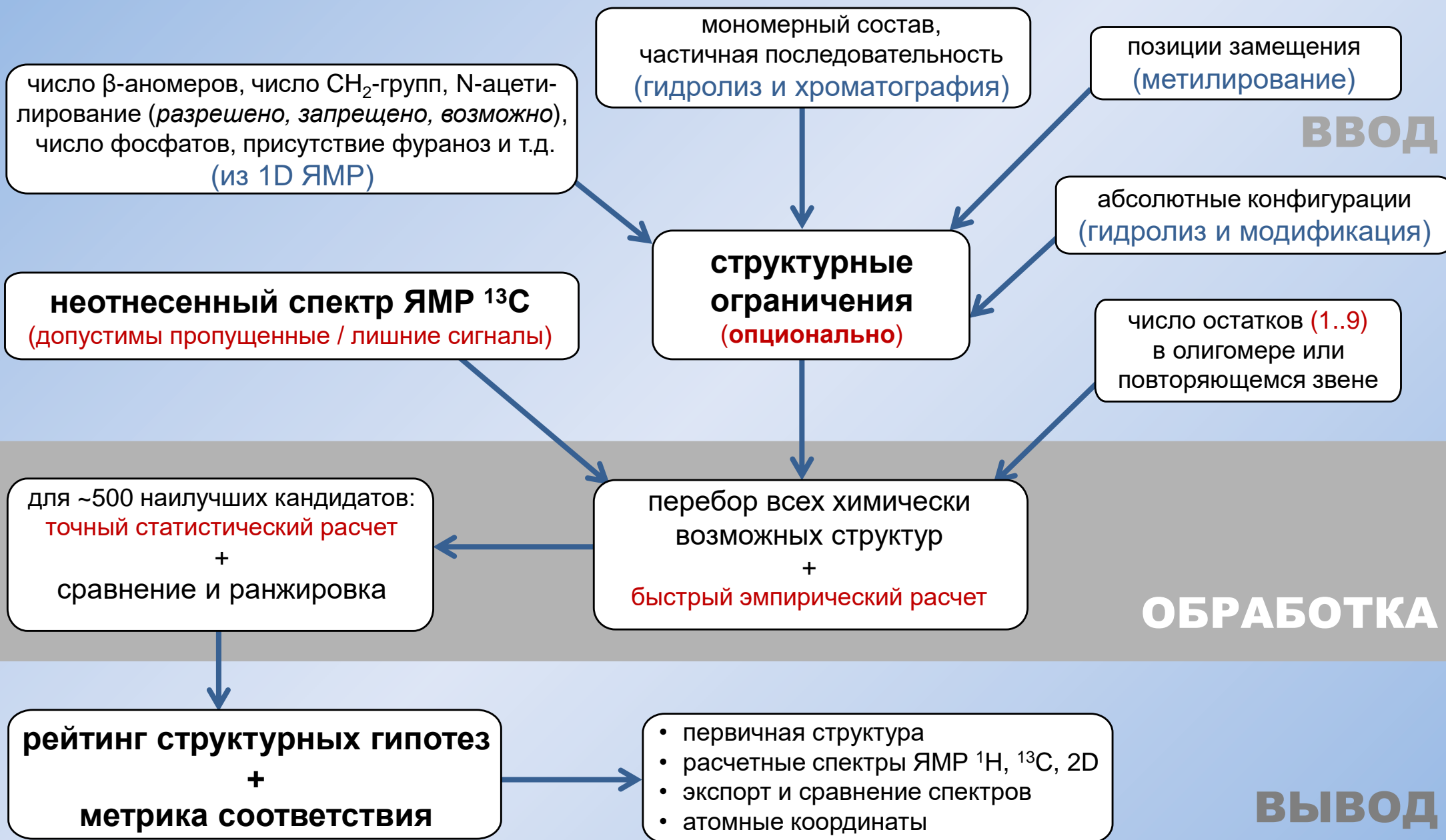
Generation, Ranking
and Assignment of
Saccharide Structures



полная структура
(оставшиеся
неизвестные)

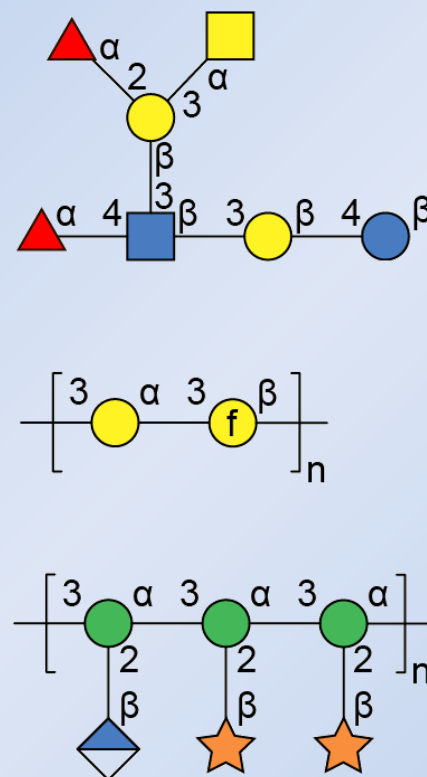
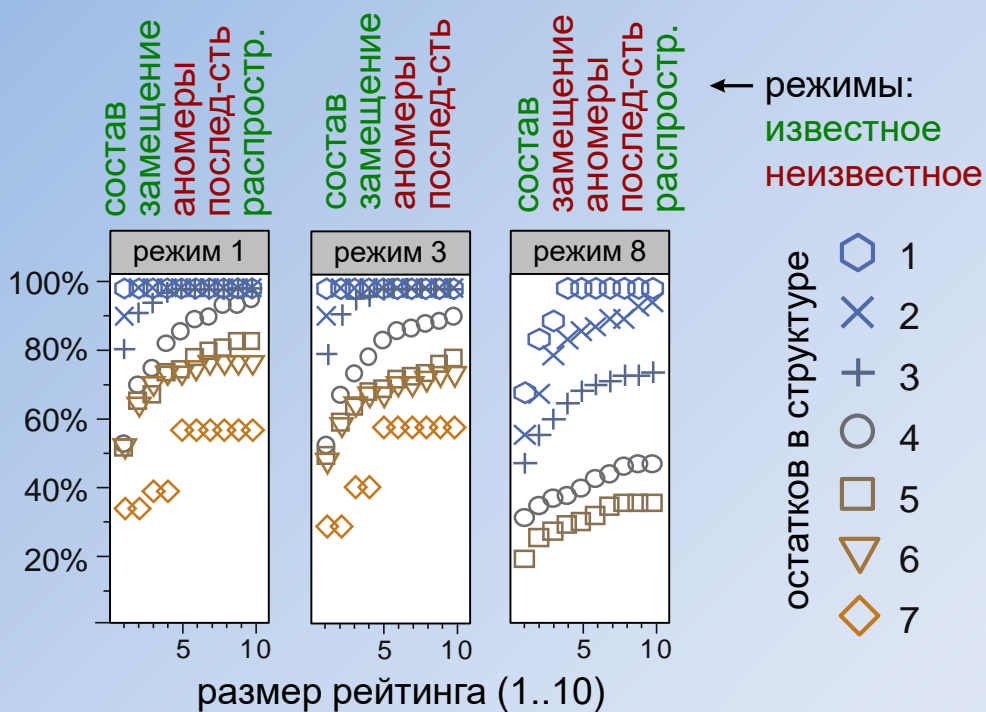
мономеры, конфигурации,
модификации, позиции
замещения,
последовательность

Выбор структурных гипотез

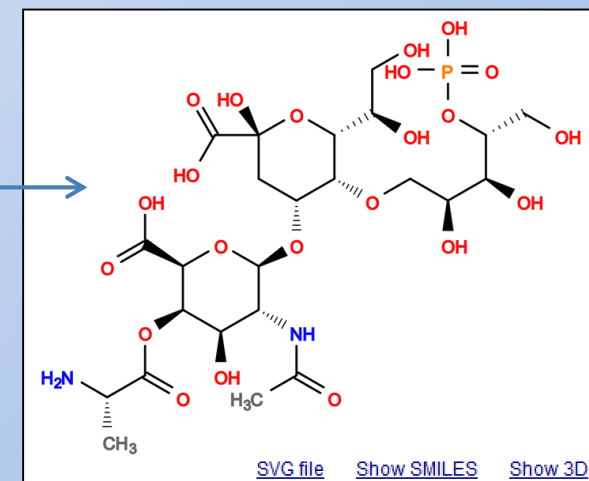
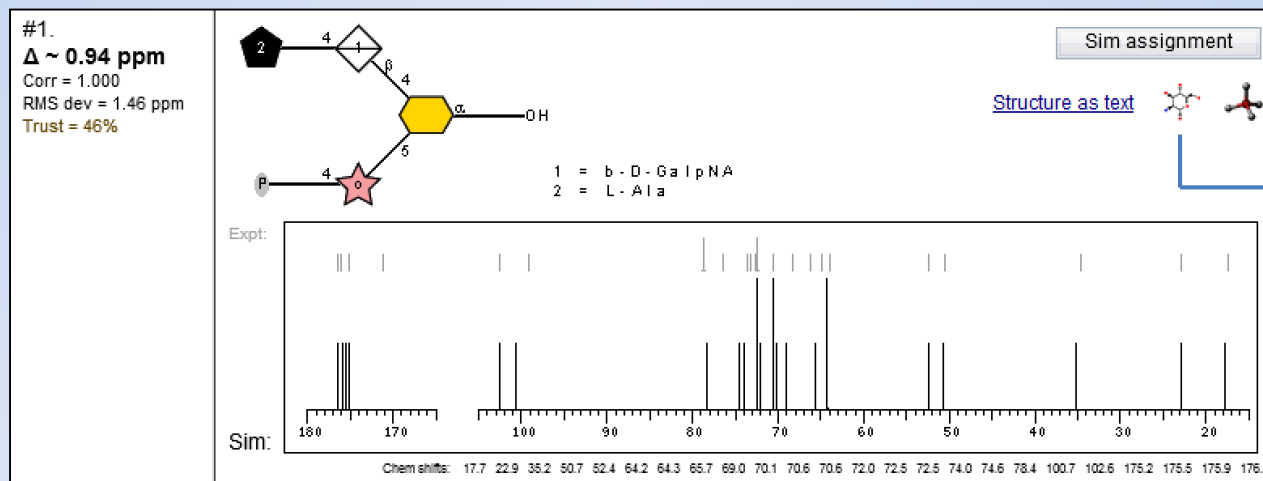


Результаты ранжирования

Статистика по 556 структурам



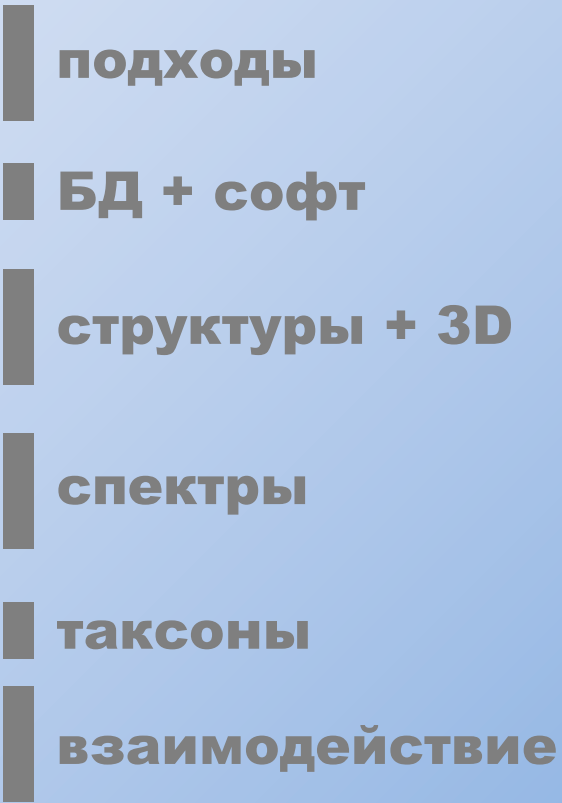
Метрики
 +
 ВЫВОД



Что сделано?

- **Цель**: оптимизация и автоматизация структурно-функциональных исследований углеводов; информатизация гликохимии и гликобиологии.
- **Средства**: разработка и апробация идеологии, моделей, стандартов, инструментов работы со структурами, спектрами, таксономией.

Что сделано?

- Цель: оптимизация и автоматизация структурно-функциональных исследований углеводов; информатизация гликохимии и гликобиологии.
 - Средства: разработка и апробация идеологии, моделей, стандартов, инструментов работы со структурами, спектрами, таксономией.
 - Разработано и объединено в систему:
 - правила обработки информации в гликохимии;
 - онтология в гликомике;
 - база данных и платформа CSDB;
 - углеводный язык и его связь с атомарными моделями;
 - быстрое получение молекулярной геометрии углеводов;
 - ЯМР-моделирование углеводов;
 - сравнение и предсказание структур по спектрам;
 - кластерный анализ химического разнообразия гликомов;
 - интеграция с конкурентами;
 - web-портал.
- 
- подходы
 - БД + софт
 - структуры + 3D
 - спектры
 - таксоны
 - взаимодействие

Что сделано?

- Цель: оптимизация и автоматизация структурно-функциональных исследований углеводов; информатизация гликохимии и гликобиологии.
- Средства: разработка и апробация идеологии, моделей, стандартов, инструментов работы со структурами, спектрами, таксономией.

- Итог:

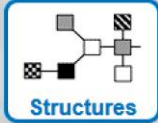



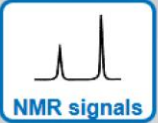
Новые гликохимические инструменты созданы, верифицированы на модельных системах и использованы для реальных исследований.

Заложен фундамент для статистических и прямых расчетов корреляции структура-свойство (*востребовано в направленном синтезе*).

Преобразилась молодая область знания – гликоинформатика, задан и обеспечен мировой вектор ее развития.



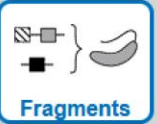



Сайт в Интернете

Database search

 Structures
  Composition
  Organisms
  Publications
  NMR signals

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

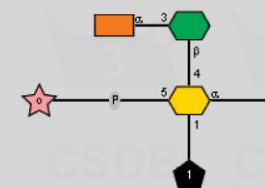
Useful tools

 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 \(field below\)](#)



Structure in CSDB encoding:

aXAbep(1-3)bXLDmanHepp(1-4)[xDRib-ol(1-P-5),xLAla?(2-1)]aXKdop
(this field is editable) [Help on structure encoding](#)

Nucleus: 1H/13C (2D) More parameters...

Solvent: Water (H or D) Coverage

Carbohydrate Structure Database

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

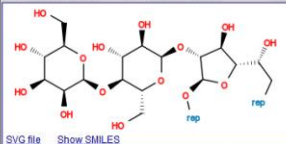
Extras

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

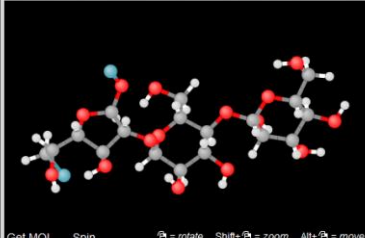
Maintenance

Related record ID(s): 101
NCBI Taxonomy refs (TaxID): 64489

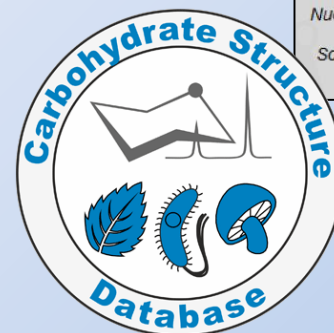
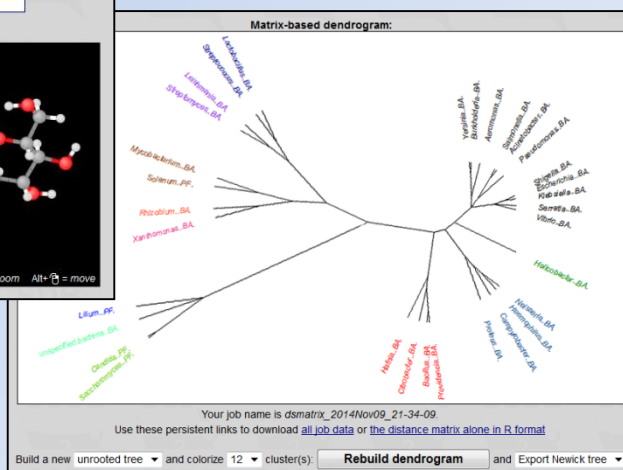
There is only one chemically distinct structure:



SVG file Show SMILES



Get MOL Spin rotate Shift+ zoom Alt+ move



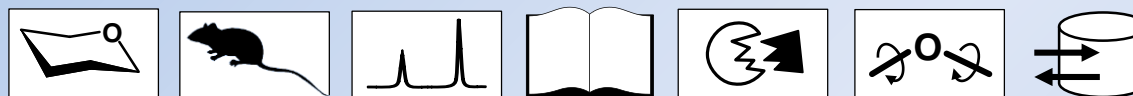
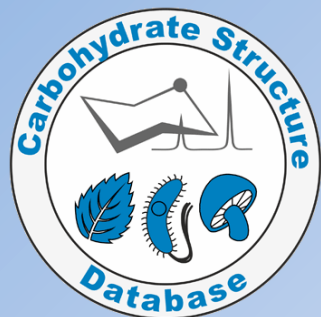
<http://csdb.glycoscience.ru>

- свободный доступ
- подробная документация
- примеры решения задач









Участники

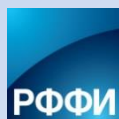
Carbohydrate Structure Database

проверяемый контент, полное покрытие



программирование
аннотирование и проверка данных
общая поддержка и сбор данных
интеграция, онтология
конформационный анализ
идеи, язык, архитектура, интерфейс,
программирование, координация
партнеры

-   Роман Капаев, Андрей Бочков, Иван Чернышов, ...
-  Ксения Егорова, Надежда Калинчук, Кирилл Казанцев, ...
-  Юрий Книрель
-   Рене Ранцингер, Кийоко Аоки-Киношита, Томас Люттеке, ...
-  Виктор Стройлов, Софья Щербинина, ...
-  Филипп Тоукач



Российский Фонд
Фундаментальных
Исследований

Книрель (1 грант)
Тоукач (3 гранта)



Международный
Научно-Технический
Центр

Книрель (1 грант)



Комиссия по
грантам при
президенте РФ

Тоукач (1 грант)



Немецкий Центр
Исследования
Рака

Тоукач (4 гранта)



Фонд Содействия
Отечественной
Науке

Тоукач (1 грант)



Российский
Научный
Фонд

Тоукач (2 гранта)

Публикации

- По теме доклада: 34 статьи + 3 монографии (суммарный импакт-фактор: ~170)
- Использование наработок в исследованиях других авторов: ~600 цитирований
- Каждый год (кроме роботов): ~2500 поисковых запросов, ~1000 запросов на расчет
- Избранные публикации:

Toukach PhV, Joshi H, Ranzinger R, Knirel YuA, von der Lieth C-W **Sharing of worldwide distributed carbohydrate resources: online connection of the Bacterial Carbohydrate Structure DataBase and GLYCOSCIENCES.de** *Nucl Acid Res* 2007, 35:D280-D286

Toukach PhV **Bacterial Carbohydrate Structure Database 3: Principles and Realization** *J Chem Inf Model* 2011, 51:159-170

Egorova KS, Toukach PhV **Critical analysis of CCSD data quality** *J Chem Inf Model* 2012, 52:2812-2814

Toukach PhV, Ananikov VP **Recent advances in computational predictions of NMR parameters for structure elucidation of carbohydrates: methods and limitations** *Chem Soc Rev* 2013, 42: 8376-8415

Капаев RR, Egorova KS, Toukach PhV **Carbohydrate structure generalization scheme for database-driven simulation of experimental observables, such as NMR chemical shifts** *J Chem Inf Model* 2014, 54:2594-2611

Egorova KS, Kondakova AN, Toukach PhV **CSDB: tools for statistical analysis of bacterial, plant and fungal glycomes** *Database* 2015, ID bav073

Toukach PhV, Egorova KS **Bacterial, Plant, and Fungal CSDB: daily Usage.** в **Glycoinformatics**, ред. Lütteke T, Frank M. Springer, 2015, гл. 5: 55-85

Капаев RR, Toukach PhV **Improved carbohydrate structure generalization scheme for ¹H and ¹³C NMR simulations** *Anal Chem* 2015, 87:7006-7010

Ranzinger R *et al.* (13 авторов вкл. Toukach PhV) **GlycoRDF: an ontology to standardize Glycomics data in RDF** *Bioinformatics* 2015, 31:919-925

Toukach PhV, Egorova KS **Carbohydrate Structure Database merged from multiple taxonomic domains** *Nucl Acid Res* 2016, 44:D1229-D1236

Капаев RR, Toukach PhV **Simulation of 2D NMR spectra of carbohydrates using GODESS Software** *J Chem Inf Model* 2016, 56:1100-1104

Egorova KS, Toukach PhV **CSDB_GT: a new curated database on glycosyltransferases** *Glycobiology* 2017, 27:285-290

Капаев RR, Toukach PhV **GRASS: semi-automated NMR-based structure elucidation of saccharides** *Bioinformatics* 2018, 34:957-963

Chernyshov IYu, Toukach PhV **REStLESS: automated translation of glycan sequences from residue-based notation to SMILES and atomic coordinates** *Bioinformatics* 2018, 34:2679-2681

Egorova KS, Toukach PhV **Glycoinformatics: bridging isolated islands in the sea of data** *Angewandte Chemie* 2018, 57:14986-14990

Toukach PhV, Egorova KS **New features of CSDB Linear, as compared to other carbohydrate notations** *J Chem Inf Model* 2020, 60:1276-1289

Stroylov VS, Panova MP, Toukach PhV **Comparison of methods for bulk simulation of glycosidic bond conformations** *Int J Mol Sci* 2020, 21: ID 7626

Scherbinina SI, Toukach PhV **Three-dimensional structures of carbohydrates and where to find them** *Int J Mol Sci* 2020, 21: ID 7702

Дополнения

(демонстрируются в рамках ответов на вопросы)

Сокращения на слайдах

BIOPSEL	BIOPolymer Structure ELucidation	MESH	Medical Subject Headings
CASPER	Computer Assisted SPectrum Evaluation of Regular polysaccharides	MSDB	MonoSaccharide DataBase
CFG	Consortium for Functional Glycomics	NCBI	National Center for Biotechnology Information
COSY	Correlation SpectroscopY	OWL	Ontology Web Language
CSDB	Carbohydrate Structure DataBase	PDB	Protein Data Bank
CSV	Comma-Separated Values	PMID	PubMed IDentifier
DFT B3LYP	Density Functional Theory: Becke 3-parameter Lee-Yang-Parr	RDF	Resource Description Framework
DOI	Digital object Identifier	REStLESS	ResiduEs as SMILES, LinkagEs as SMARTS
GODDESS	Glyco-Optimized Database-Driven Empirical Spectrum Simulation	SMARTS	SMiles ARbitrary Target Specification
GRASS	Generation, Ranking and Assignment of Saccharide Structures	SMILES	Simplified Molecular-Input Line-Entry System
GT	GlycosylTransferase	SNFG	Symbolic Notation For Glycans
HMBC	Heteronuclear Multiple Bond Correlation	SPARQL	SPARQL Protocol and RDF Query Language
HOSE	Hierarchical Organization of Spherical Environment	TOCSY	TOTal Correlation SpectroscopY
HSQC	Heteronuclear Single Quantum Coherence	WSDL	Web Service Description Language
ICD	International Classification of Diseases	WURCS	Web-3 Unique Representation of Carbohydrate Structures
IUPAC	International Union of Pure and Applied Chemistry		

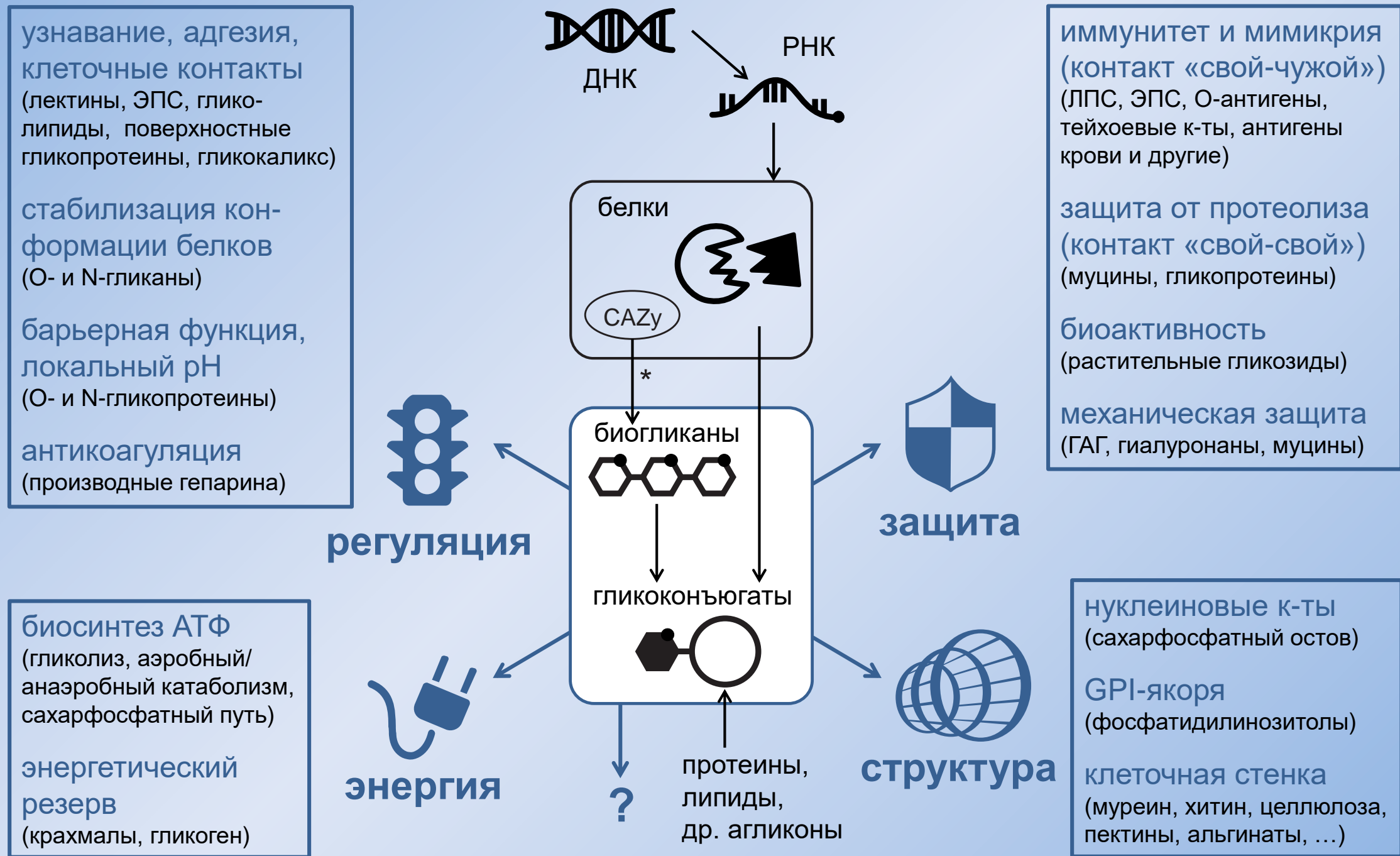
● СТАНДАРТНЫЕ

● НОВЫЕ

Выводы - кратко

- Аннотации ~10 тыс. публикаций (~20 тыс. структур). Свободно доступная база данных Carbohydrate Structure Database (CSDB), оснащённая многочисленными видами поиска, представления и анализа данных.
- Инструмент автоматического моделирования структуры биогликанов - фундамент для расчётов корреляции «структура – свойство» в химии углеводов.
- Аннотации ~9 тыс. спектров ЯМР биогликанов. Исследование взаимосвязи «структура – спектр». Метод обобщения атомного окружения в углеводных структурах и моделирования ХС с точностью 0.06 (^1H) и 0.69 (^{13}C) м.д.
- Программа генерирования и ранжирования структурных гипотез для установления первичной структуры по экспериментальным данным.
- Углеводный язык CSDB Linear. Семантическое описание углеводов, используемое в публикациях, связано с поатомным описанием, используемым в химических расчётах.
- Стандарты визуализации углеводных структур, их одно- и двумерных спектров ЯМР. SNFG рекомендована ведущими углеводными журналами.
- База данных активностей гликозилтрансфераз.
- Анализ распространённости структурных особенностей и характерных признаков углеводов в различных таксономических группах. Альтернативные углеводные «деревья жизни».
- Правила гликоинформатики и углеводная онтология GlycoRDF. Автоматическое взаимодействие между проектами - получение знаний, неявно содержащихся в нескольких базах разного типа.
- Выявление и исправление ошибок в ~340 опубликованных структурах / спектрах и ~2 тыс. ошибок в базах данных.

Роли углеводов



Структурные базы

CarbBank 23 полная до 1996 ORIG архитектура, % ошибок

GlycomeDB 99 мета-репозиторий неполные аннотации нет агликонов

GLYTOUCAN

CFG glycan млекопитающие, > 6

SweetDB, SugaBase

GI-CO-SCIENCES.DE 25 / 19

GLYCAN 11

Eurocarb DB архитектура только модель

BCSDB PFCSDB **13 / 5** (бактерии, археи) **7 / 2.5** (грибы, растения)

Carbohydrate Structure Database ORIG полная по прокариотам курируемая

GlycoSuite ORIG млекопитающие+... **10 / 1** полная до 2005

JCGGDB > 70 коллекция баз слабо аннотированы

UniCarbKB 4 / 1 курируемая

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

Glycoconjugate Data Bank 44

EcoDAB **0.2 E. coli** ORIG

GlycoBase **0.3 животные** ORIG

Специальные базы



углеводные гены человека



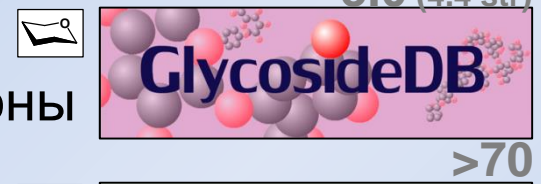
MS^{2,3,4} N- и O-гликанов



химические реакции



конъюгаты & агликаны



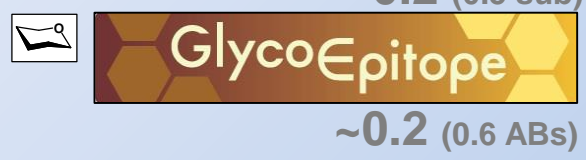
N-гликопротеины
C. elegans + мышь



методики синтеза
и анализа



гликоэпитопы
и антитела



GlyTOUcan,
репозиторий идентификаторов



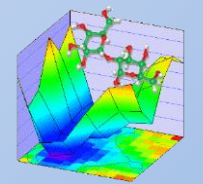
адгезия к патогенам



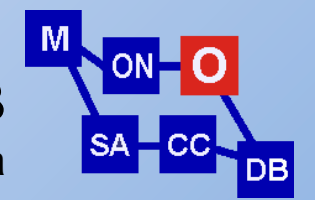
O-glycBase,
O- и C-гликопротеины



GlycoMaps,
расчетные конформационные карты



MSDB
моносахариды и номенклатура



~0.8

Примеры вопросов к CSDB

- Как введение аминогруппы влияет на химические сдвиги ЯМР в лактозном фрагменте?
- Какие структуры О-антигенов, содержащие галактуроновую кислоту и еще как минимум одну гексозу, были опубликованы после 2005-го года?
- Какие гликозиды, выделенные из растений рода паслёна, содержат агликон соланидин?
- Какие углеводы, кроме октозосодержащих, имеют сигнал ЯМР ^{13}C около 34 м.д. ?
- Какие бактериальные структуры, опубликованные А.С. Шашковым или Ю.А. Книрелем, содержат хиновоз-4-амин, амидированный любой N-ацетилированной аминокислотой?
- Гомополимеры каких нонапираноз встречаются в бактериях?
- Каков ожидаемый спектр ЯМР ^{13}C 3-О- α -абеквозил-6-деокси- β -D-манногептопиранозил-(D-рибитол-1)-фосфата в воде и на основании каких источников предсказаны химические сдвиги, для которых указана наименьшая достоверность?
- Какова наиболее вероятная последовательность остатков бациллозамина, глюкуроновой кислоты и лизина в олигомере с указанным экспериментальным спектром ЯМР ^{13}C ?
- Какие моносахариды склонны занимать концевые позиции в гликанах Аспергилла дымящего и Аспергилла кодзи?
- Какие димерные фрагменты (включая дисахариды) гликанов высших растений специфичны для рода люпинов?
- Для скольких гликанов протеобактерий опубликованы спектры ЯМР?

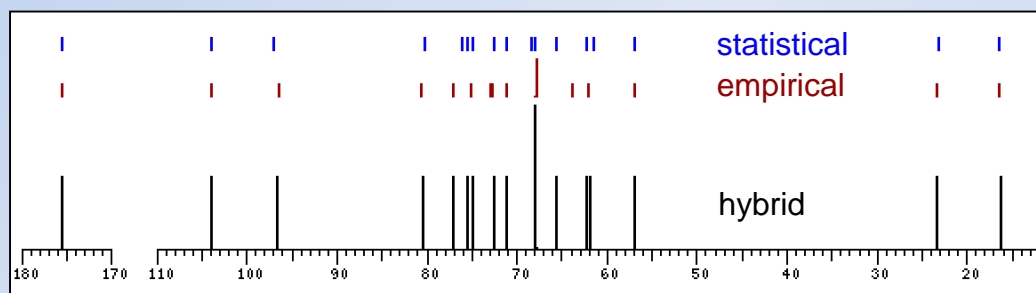
Эмпирическое моделирование

- Инкрементная схема со стерической коррекцией (BIOPSEL)
- Учитывает структурное окружение (9-13 дескрипторов)
- Использует специальные базы хим. сдвигов и эффектов

(440 мономеров, 3300 димеров и тримеров, 300 эмпирических эффектов)

- Поддерживает большинство углеводных структур
- Оценивает степень достоверности предсказания

результаты смешиваются
со статистической симуляцией
с учетом достоверности



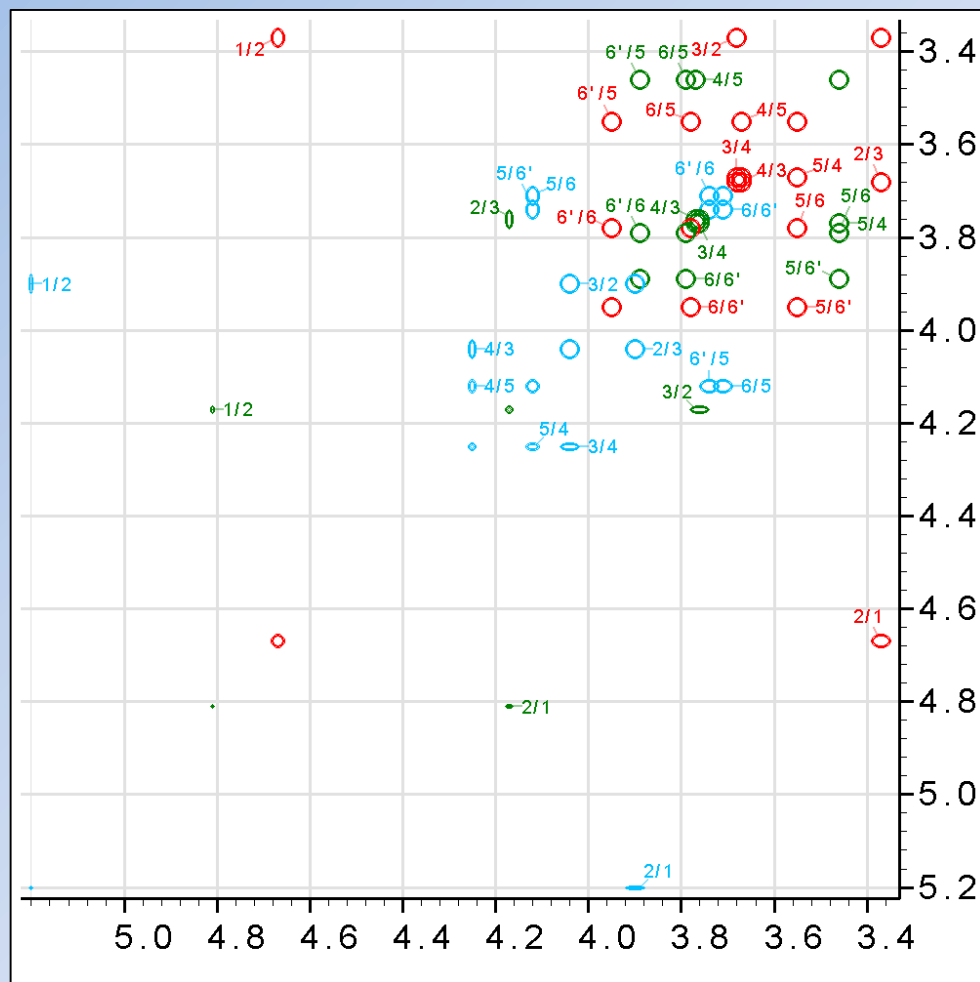
β -D-GlcpNAc-(1-3)- α -D-Fucp-(1-P-3)-D-Gro

¹³C NMR data (in D₂O):

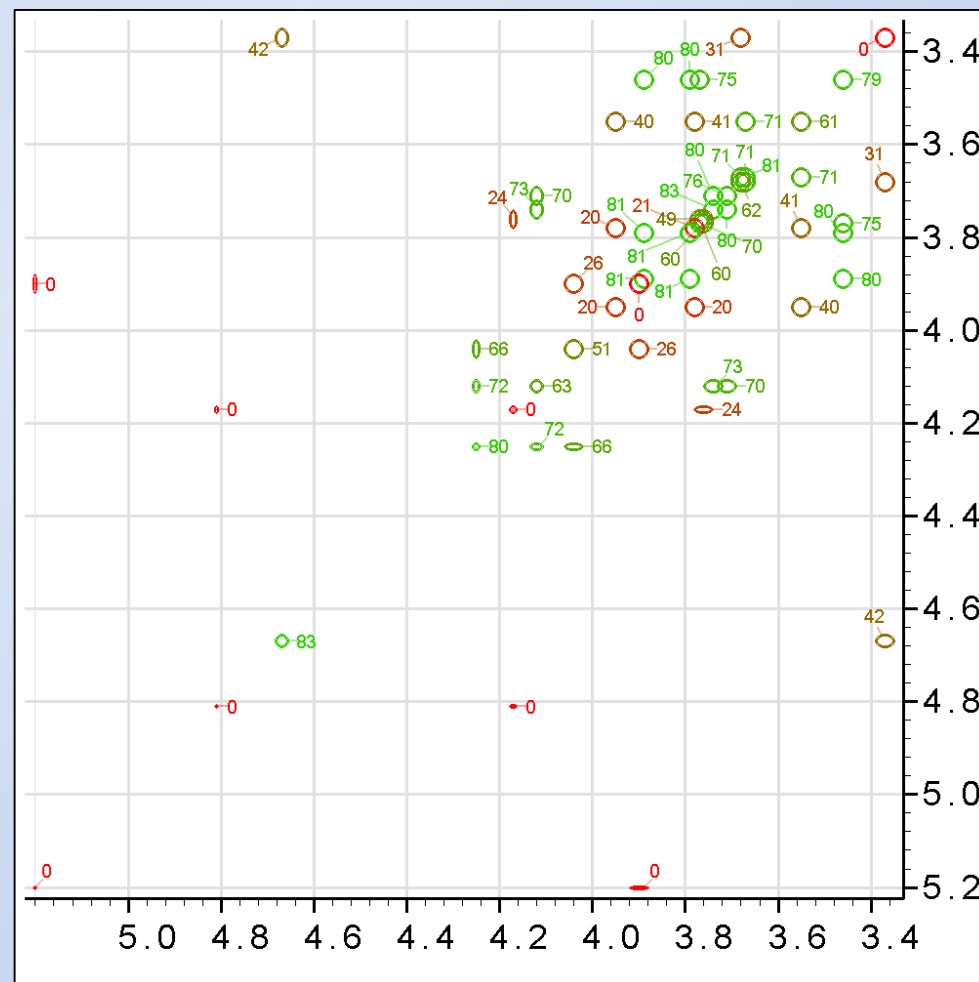
Linkage	Residue	C1	C2	C3	C4	C5	C6	C7	C8	Accuracy
0	xDGro unsubstituted-> effect->	64.0 64.0	73.0 73.5 -0.5	68.0 64.0 +4.0						0.5
3	xXP? unsubstituted-> effect->									2
3,0	aDFucp unsubstituted-> effect->	96.5 93.3 +3.2	67.9 69.2 -1.3	80.6 70.4 +10.2	72.7 73.0 -0.3	67.8 67.4 +0.4	16.5 16.7 -0.2			4
3,0,3	bDGlcNAc unsubstituted-> effect->	104.0 96.2 +7.8	57.0 58.0 -1.0	75.1 75.1	71.2 71.2	77.2 77.2	62.1 62.1	175.6 175.6	23.5 23.5	4

пример: ^1H - ^1H COSY

$\rightarrow 3$ - α -D-Galp-(1 \rightarrow 3)- β -D-Manp-(1 \rightarrow 4)- β -D-Glcp-(1 \rightarrow



принадлежность остатку



оценка достоверности

Оценка достоверности

0%



100%



- большой общий вес генерализаций
- в базе мало структур
- противоречивые данные ЯМР

- малый общий вес генерализаций
- в базе много структур
- данные ЯМР согласуются
- статистические и эмпирические предсказания близки

+ ожидаемая погрешность каждого сигнала (*м.д.*)
предсказывается из ХС и степени достоверности
(регрессией)

ССЫЛКИ



<http://glytoucan.org>

Abrahams J.L. 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)



http://jcgddb.jp/index_en.html

K.F. 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://glycosciences.de>

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



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

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://csdb.glycoscience.ru>



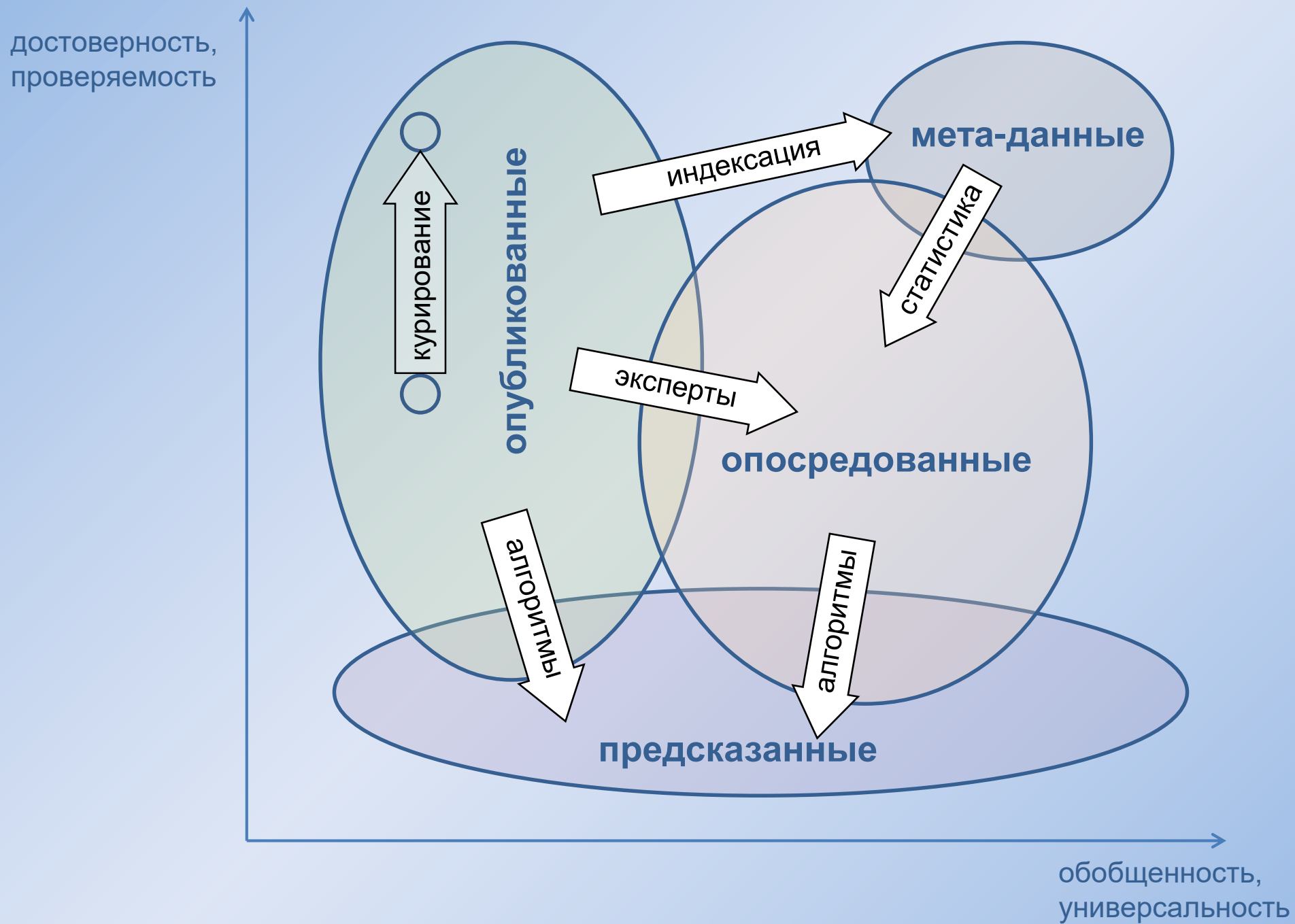
<http://toukach.ru/rus/glycoinf.htm>

Дальнейшее развитие

● сделано в CSDB ● предстоит сделать ● близко к завершению

- **Стандартный человекочитаемый язык** (SNFG, CSDB Linear, ...)
- **Расширение онтологий** (интеграция через GlycoRDF и GlycoCoO)
- **Кросс-проектные сервисы**
(ввод-вывод структур, конформационные расчеты, предсказание свойств)
- **Стандартные индексы**
(Glytoucan ID, MSDB, PMID, DOI, TaxID, ICD-11, PDB id, Genbank, ...)
- **Стандартные протоколы** (API, WSDL, SPARQL, REST, ...)
- **Идеологическая замена CarbBank** 
- **Требование включать ID в публикации** (Glytoucan ID?)
(кто будет чистить базы от неопубликованных/ошибочных данных?)

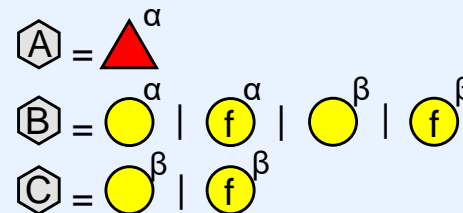
Уровни данных



1. Наборы остатков

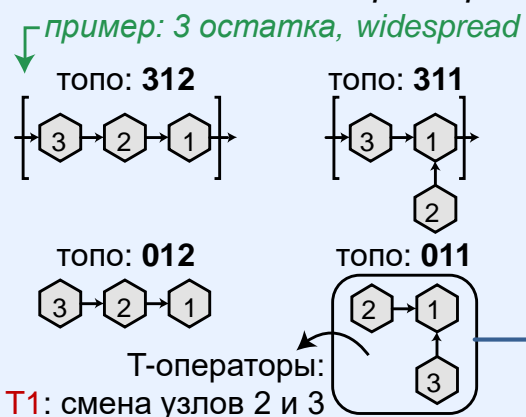
фильтры: «нет фураноз»,
widespread:
редкие остатки
(если есть обычные)

	набор	α/β	D/L	имя	цикл
пример	A	α	L	Fuc	p
	B	?	D	Gal	?
	C	β	D	Gal	?



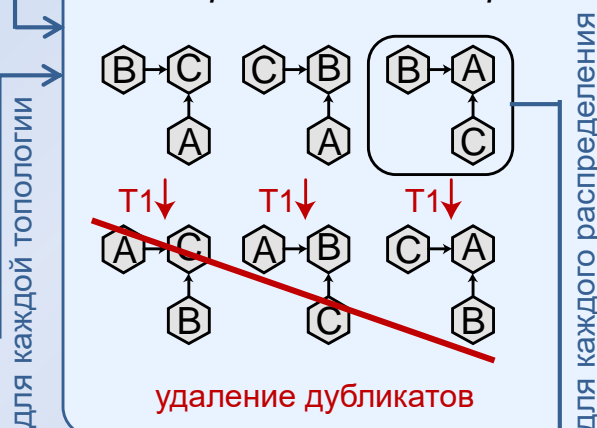
2. Топологии и T-операторы

фильтры: поли/олиго
widespread:
суперразветвленные
большие боковые цепи



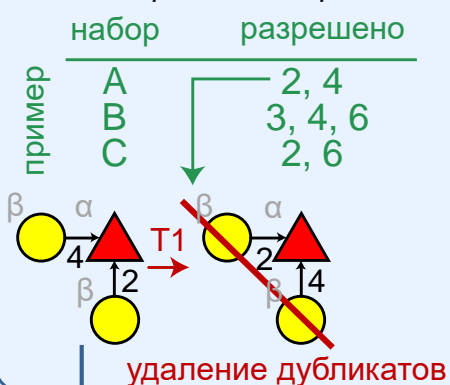
3. Распределения наборов

фильтры: разрешенные акцепторы
min/max заместителей
расположение в структуре
widespread:
суперразветвленные



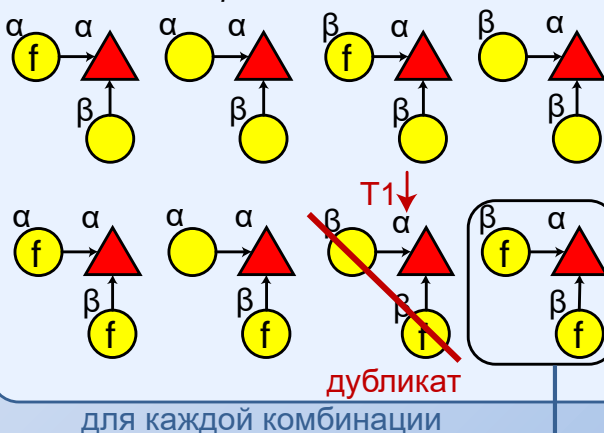
5. Позиции замещения

фильтры: widespread:
(1-1) связи и т.д.



4. Комбинации остатков

фильтры: число сигналов в спектре
тип N-ацетилирования
число β -сахаров
число CH_2 групп
widespread:
суперразветвленные



для каждой структуры → симуляция спектра ЯМР ^{13}C и ранжирование

Схема CSDB

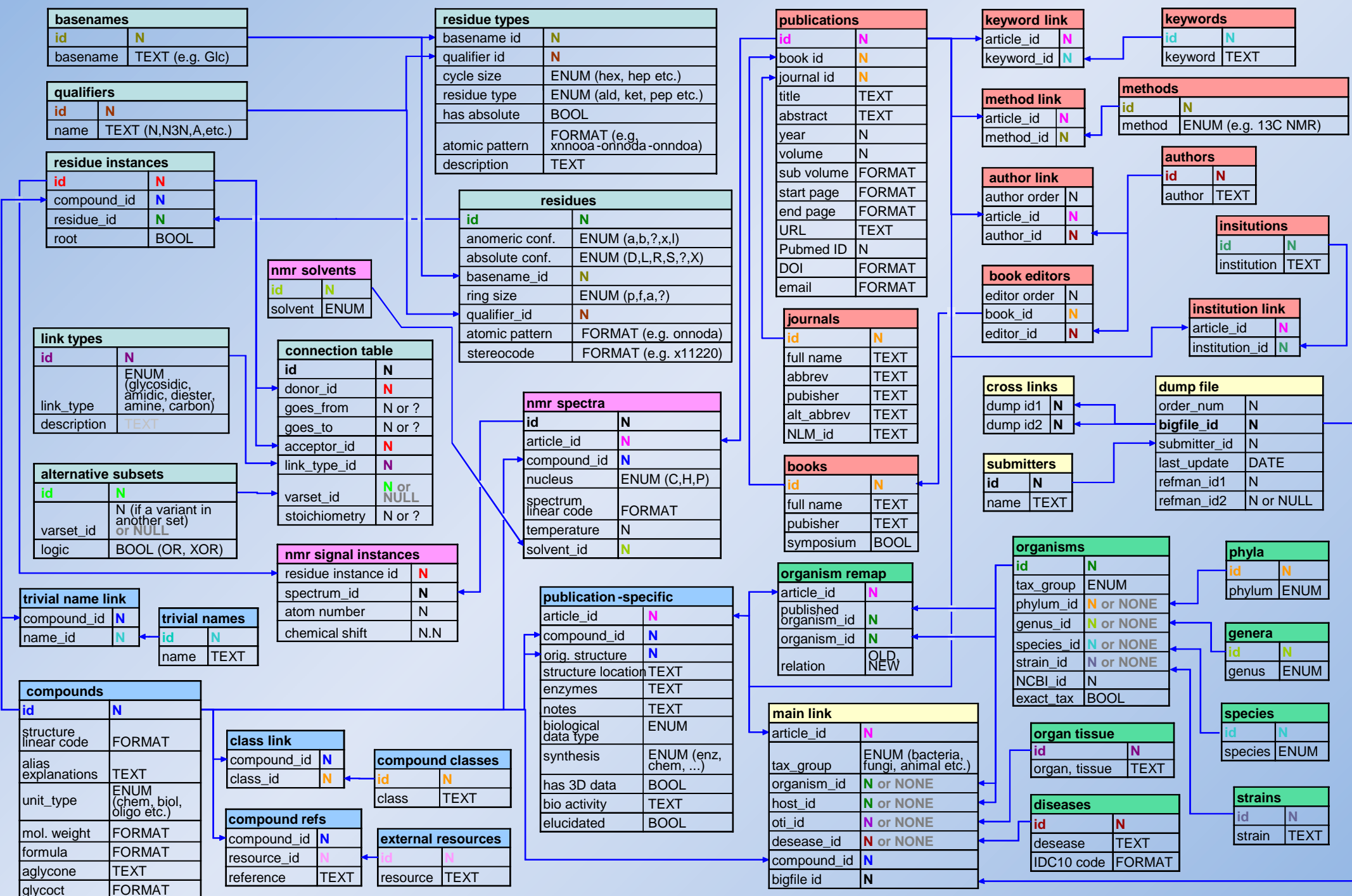
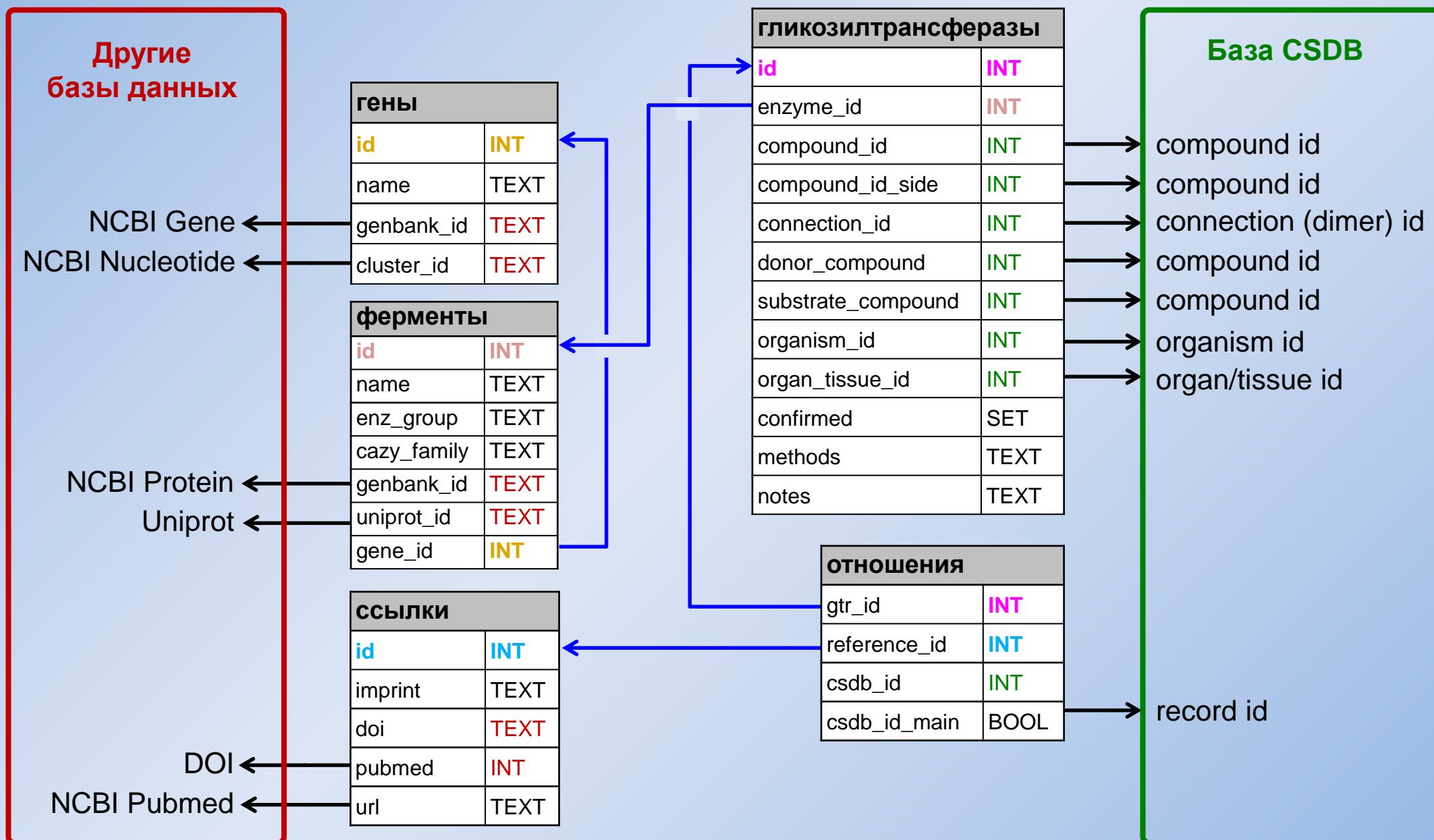


Схема CSDB GT



Экспорт в RDF

Структура

Ссылки (ресурс, ID, ссылка)

Записи структуры (GlycoCT, GlydeII, BCSDB, ... , SweetDB)

Состав (ссылки на MSDB, *non-glyco residues?*)

Эксп. молекулярный вес (+комментарии, напр. 1234 [M+])

Тип звена (*monosaccharide or derivative, oligomer, chemical repeat, biological repeat, cyclic repeat, homopolymer repeat*)

Степень полимеризации

Тривиальное название

Структурный класс

Спектры ЯМР, отнесение, растворитель, температура

Атомные дескрипторы и координаты (ссылка на файл)

Изображения структуры (ссылка на файл)

Структура – Публикация

Методы (ссылки на MeSH)

Положение структуры в статье

Ассоциированные публикации

Ссылки (ресурс, ID, ссылка)

PMID, DOI, ISSN, ISBN, NLM ID, www-link

Название ресурса (*журнала или книги*),

Издательство, Язык

Тип ресурса (*журнал, симпозиум, глава*)

Год, том, выпуск, страницы, глава (Dublin Core)

Институты авторов

Ключевые слова

Ассоциированные организмы

Ссылки (ресурс, ID, ссылка) (*UniProt и др.*)

Домен, Тип

Таксон (род, вид, штамм/серогруппа, ген. линия)

Хост-организм (ссылка на запись RDF)

Орган, ткань (+ссылка на хост)

Структура – Организм

Болезнь (*IDC code*)

Стадия жизни

Вовлеченные белки (ссылки)

и их гены (ссылки на GenBank)

```
http://csdb.glycoscience.ru/integration/rdf.php
?id=<####> &mode=<structure, publication, organism, spectrum, relation, record>
&format=<turtle, rdfxml, rdfjson, ntriples> &clean=<0,1>
```

Ввод и вывод структуры

CSDB/SNFG structure editor

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

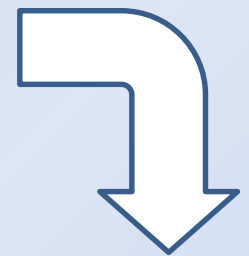
Glc GlcNAc GlcA QuiNAc Gal GalNAc GalA Fuc FucNAc Man Rha LDmanH Ara Ara4N Xyl Fru P Kdo Neu5Ac Gro Ala

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-

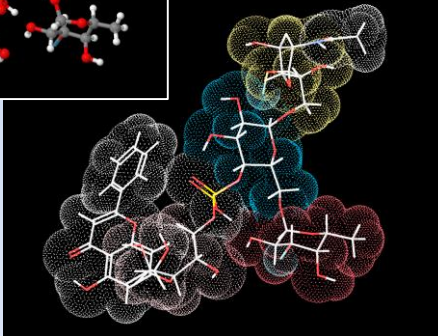
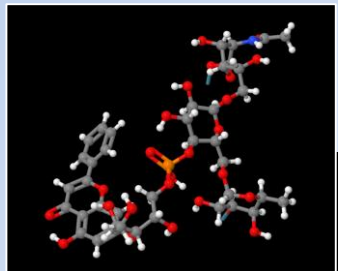
Subst = chrysin = SMILES O=c2cc(c1ccccc1)oc3c{7}c(0)cc(0)c23

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

REMARK 300 USING GENERATED NAME Fdb FOR RESIDUE xSubst

AUTHOR GENERATED BY OPENSABEL 2.4.1, RESIDUE ASSIGNMENT BY CARBOHYDRATE STRUCT

COMPND	CSDB linear =								
HEXTAM	1 R	Pda	A	1	0.156	1.111	2.494	1.00	0.00
HEXTAM	2 C1	Fda	A	1	1.377	1.487	3.097	1.00	0.00
HEXTAM	3 C1	Pda	A	1	1.643	2.889	2.961	1.00	0.00
HEXTAM	4 O4	Pda	A	1	1.951	3.280	1.601	1.00	0.00
HEXTAM	5 O4	Fda	A	1	3.374	3.123	1.402	1.00	0.00
HEXTAM	6 O5	Fda	A	1	3.933	4.392	3.736	1.00	0.00
HEXTAM	7 O5	Pda	A	1	3.890	5.497	1.642	1.00	0.00
HEXTAM	8 O6	Pda	A	1	3.185	4.772	-0.562	1.00	0.00
HEXTAM	9 O6	Fda	A	1	3.501	3.900	-1.665	1.00	0.00
HEXTAM	10 O3	Pda	A	1	4.005	2.825	2.771	1.00	0.00
HEXTAM	11 O3	Fda	A	1	4.252	1.419	2.915	1.00	0.00
HEXTAM	12 O2	Fda	A	1	2.912	3.229	3.759	1.00	0.00
HEXTAM	13 N2	Fda	A	1	3.019	4.649	4.105	1.00	0.00
HEXTAM	14 H1	Fda	A	1	0.780	3.467	3.305	1.00	0.00
					2.260	0.749	1.00	0.00	
					4.211	0.481	1.00	0.00	
					6.274	1.168	1.00	0.00	
					5.771	-0.887	1.00	0.00	
					4.869	-0.431	1.00	0.00	
					3.336	2.931	1.00	0.00	
					1.146	2.209	1.00	0.00	
					2.647	4.689	1.00	0.00	
					5.266	3.368	1.00	0.00	
					5.202	5.040	1.00	0.00	
					6.697	5.164	1.00	0.00	
					4.543	5.689	1.00	0.00	
					6.973	6.220	1.00	0.00	
					7.056	4.765	1.00	0.00	
					7.165	4.614	1.00	0.00	
					2.820	-1.898	1.00	0.00	
					3.154	-2.751	1.00	0.00	
					3.270	-4.096	1.00	0.00	
					3.738	-4.973	1.00	0.00	



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
2. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]?DGlcp(1-5)[Ac(1-2)]bDGalfN(1- // Subst
3. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]?DGlcp(1-6)[Ac(1-2)]bDGalfN(1- // Subst

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](Oc3cc(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
2. -3)aLFucp(1-6)[Subst(7-3)xDRib-ol(1-P-4)]bDGlcp(1-6)[Ac(1-2)]bDGalfN(1- // Subst

Files: MOL, PDB, Glycam

3D Shift+ zoom Shift+ pan Alt+ rotate Ctrl+ menu

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?](#)

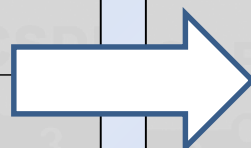
has aglycon:

add substitution
 add substituent
 add substituent
 add substituent

Structure in CSDB encoding:

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

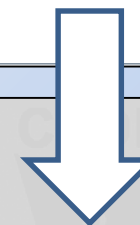
[Home](#) [Help](#)



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 !!!](#)

Поиск по составу

Prokaryotes + Plants + Fungi

7005 publications (1941-2017):
18923 compounds from
8859 organisms
last update: 2017 Feb 13

Search

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

Help

- About
- Basic usage
- Statistical tools
- NMR tools
- Usage examples
- Advanced features
- Structure encoding
- Database docs
- Credits

Extras

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

Maintenance

Found 6 structures. Displayed structures from 1 to 6

Expand all compounds Show all as text (SweetDB notation)

Requested composition: 2 Gal, 1 HEX, 1 P

1. **Compound ID: 6420**

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

Structure type: polymer biological repeating unit ; n~27
Aglycon: core part of the molecule
Trivial name: P-saccharide core
Compound class: lipophosphoglycan

The structure is contained in the following publication(s):

- Article ID: 2648
McConville MJ, Thomas-Oates JE, Ferguson MAJ, Homans SW
"Structure of the lipophosphoglycan from *Leishmania major*" -
Journal of Biological Chemistry **265** (1990) 19611-19623

Leishmania major
[CSDB ID 235965](#) (all data & tools)

Expand this compound

2. **Compound ID: 6361**

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

Search for residue composition

Complete structural composition (3 units) :

1. phosphoric acid x 1 --> 1 x P
2. galactose x 2 --> 2 x Gal
3. any hexose x 1 --> 1 x HEX

[Add unit \(+\)](#) [Remove unit \(-\)](#)

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)

Search for molecule types: All molecule types

- Search for complete composition (not a fragment)
- Restrict compound class: -select class-
- Restrict taxonomical domain: All domains

Previous results: 2 structures: 4104,8808

[Go!](#) & display 30 records per page.

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

Поиск по таксономии

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

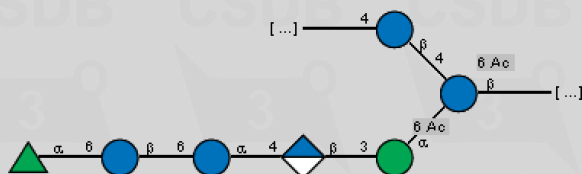
1. (Organism ID: 1005)

[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:

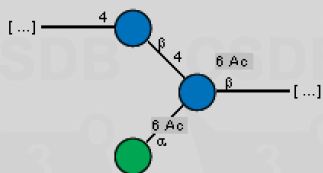
- Compound ID: 1717



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

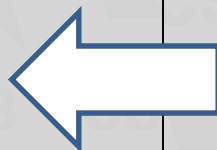
Carbohydrate Research 2004, "Synergistic interactions between the genetically modified bacterial polysaccharide P2 and carob or konjac mannan"
[CSDB ID 9262](#) (all data & tools)

- 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:

Species:

Strain / subspecies:

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:

Поиск по библиографии

Found 3 publications. Displayed publications from 1 to 3

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

1. (Article ID: 1525)

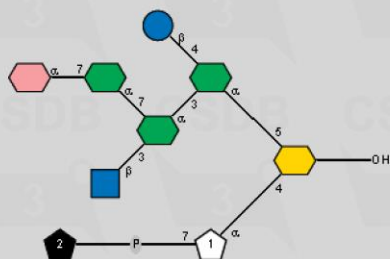
Knirel YA, Lindner B, Vinogradov EV, Shaikhutdinova RZ, Senchenkova SN, Kocha
Cold temperature-induced modifications to the composition and structure of Yersinia pestis lipopolysaccharide
Carbohydrate Research **340(9)** (2005) 1625-1630

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. LPS of a new type was isolated from *Y. pestis* KM218 that the latter differs in: (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-deoxy-4-phosphoryl groups; (iv) the presence of a 4-amino-4-deoxy-4-phosphoryl group; (v) the absence of a fatty acid(s). The data obtained suggest that cold temperature-induced modifications 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

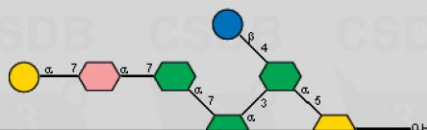


1 = α-Kop
 2 = EtN

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

Yersinia pestis KM218
[CSDB ID 10076](#) (all data & tools)

• Compound ID: 4210



Search for bibliography

Authors: start with:
[Help on author/keyword query syntax](#) [ä ö ü á é í ó ç š](#)

Title: search also in abstract
 (content of title) [Help on title/abstract query syntax](#)

Keywords: search also in title
 (content of keyword section) [Help on author/keyword query syntax](#)

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

Year:
 1984
1985
 1986
 1987
 1988
 1989

Vol:
Page:

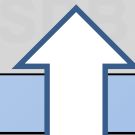
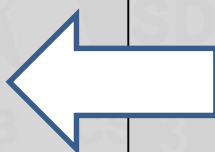
Search scope:

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:

& display records per page.

[PubMed XML](#) [Home](#) [Help](#)



Author index:

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

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

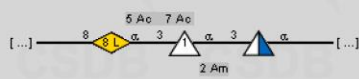
[Close this window](#)

Поиск сигналов ЯМР

Found **2** structures. Displayed structures from **1** to **2**

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

1. Compound ID: 150 (similarity: 2.5)



1 = a-L-FucpN [Show legend](#) [Show as text](#)

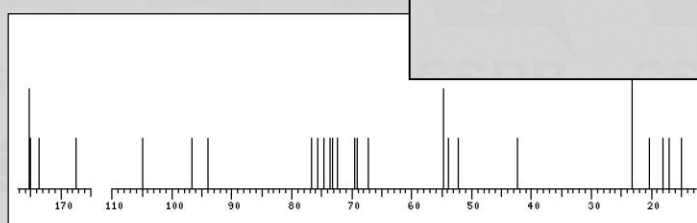
Structure type: polymer biological repeating unit

The average similarity of its ¹³C NMR spectra with the search term (signals in bold) is **2.5** ([help on similarity values](#))

¹³C NMR spectra assigned to the structure:

- in [Article ID 3480](#):
NMR conditions: in D2O at 313 K
¹³C NMR data:

Linkage	Residue	C1	C2	C3	C4	C5	C6	C7	C8	C9
3,3,5	Ac	175.3	23.1							
3,3,7	Ac	175.3	23.1							
3,3	aX8eLegp	173.7	105.0	42.4	69.2	53.8	73.6	54.7	73.2	14.9
3,2	Am	167.5	20.4							
3	aLFucpN	96.7	52.3	75.6	72.4	67.3	17.1			
2	Ac	175.0	23.3							
	aDQuipN	94.1	54.6	76.7	74.6	69.6	18.0			



The structure is contained in the following publication(s):

- Article ID: 31
Bystrova OV, Lindner B, Moll H, Kocharova NA, Knirel YA, Zähringer U, Pier GB "**Structure of the lipopolysaccharide of Pseudomonas aeruginosa O-12 with a randomly O-acetylated core region**" - *Carbohydrate Research* **338(18)** (2003) 1895-1905
Pseudomonas aeruginosa O12
[CSDB ID 1824](#) (all data & tools)
- Article ID: 3480
King JD, Mulrooney EF, Vinogradov E, Kneidinger B, Mead K, Lam JS "**IfnA from Pseudomonas aeruginosa O12 and wbuX from Escherichia coli O145 encode membrane-associated proteins and are required for expression of 2,6-dideoxy-2-acetamido-L-galactose in lipopolysaccharide O antigen**" - *Journal of Bacteriology* **190(5)** (2008) 1671-1679
Pseudomonas aeruginosa O12
[CSDB ID 23007](#) (all data & tools)

[Expand this compound](#)

2. Compound ID: 11058 (similarity: 1.364)

Search for NMR signals

Nucleus: Carbon Proton

Threshold: [Threshold explanation](#)

Chemical shifts:

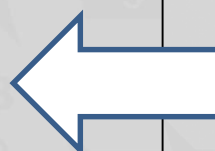
Search scope:

- Search the whole database
- Search in the result of the previous query (logical AND)
- Combine with the result of the previous query (logical OR) Signals within a single residue
- Negate search (find results NOT matching current query)

& display records per page.

Previous results: 2 structures: 150,11058

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



Поиск по ID и полная запись

2. (BCSDB ID: 50001)

Egorova KS, Toukach FV
This is a test record for BCSDB 3 debug N2
 Carbohydrate Research 489 (1944) 9726

/Variants 0/ is:
 Me-3) -
 OR (inclusively)
 Ac-2) -

a-Neup5 (10%) Ac- (2-6) -+
 |
 S-Pyr- (2-4:2-3) -a-Fucp- (1--P--4) --b-D-GlcpN
 |
 /Variants 0/-+

Deliriumus trementii O67 PCM2005
 (Ancestor NCBI TaxID 374, [species name lookup](#))

Taxonomic group: bacteria / Pseudophylum, Newphylum (Phylum: Pseudophylum, Newphylum)
 Host organism: Worra worra, Homo sapiens
 Organ / tissue: [brain](#), [ass](#)
 Associated disease: [googlomania](#)

NCBI PubMed ID: [123456789](#)
 Journal NLM ID: [0043535](#)
 Publisher: Elsevier
 Institutions: N.D. Zelinsky Institute of Organic Chemistry, Hirzfeld Institute

The structure of the O-specific polysaccharide of Deliriumus trementii O66 has been elucidated using 2d-NMR approach, including... The studies of biology the effect of this glycopolymer on sexual behaviour of Worra-worra's bla-bla-bla...

structure, antigen, polysaccharide, worra, Deliriumus trementii

Structure type: oligomer ; 2000-3000
 Location inside paper: HPLC fraction 8
 Trivial name: khrenbiose
 Compound class: X-antigen, Y-antigen

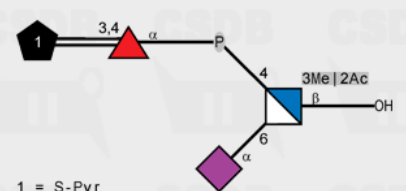
Methods: 1D-NMR, 2D-NMR, 13C-NMR
 Biological activity: causes suppression of Worra-worra's natural instincts
 Enzymes that release or process the structure: CoQ-III
 Biosynthesis and genetic data: biochemical data
 Comments, role: structure was revised (see RR 1234), absolute configurations are not determined
 3D data: conformation data, computer modelling, dynamics

NCBI Taxonomy refs (TaxIDs): [6661313](#), [374](#)
 Reference(s) to other database(s): SpecDB:01.9, CA:123, CA-RN:456, patent:USA#123, ProtDB:ABC456

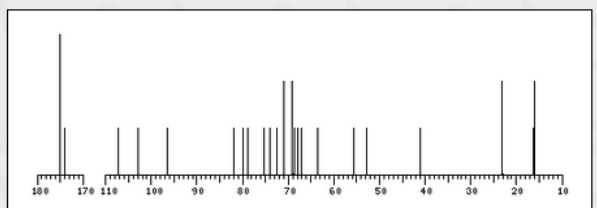
NMR conditions: in D2O at 308 K

¹³C NMR data:

Linkage	Residue	C1	C2	C3	C4	C5	C6	C7	C8	C9
6,5	10%Ac	175.0	23.0							
6	aXNeup	174.1	107.2	41.0	69.0	52.9	73.9	69.2	72.6	63.6
4,0,4	xSPyr?	175.1	80.0	16.0						
4,0	a?Fucp	96.5	68.6	71.0	81.9	67.8	15.9			
4	P									
0	bDGlcpN	103.0	55.6	78.9	70.9	75.4	67.1			
2	Ac	175.0	23.2							
3	Me	16.3								



1 = S-Pyr



Search for CSDB IDs

Scope: Record IDs Structure IDs Publication IDs Organism IDs
 RDF only: Source IDs Spectrum IDs Relation IDs

CSDb record IDs:
 You can use commas to separate IDs and hyphens to specify the range, e.g. 100-150,160-165,170

& display records per page.

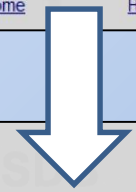
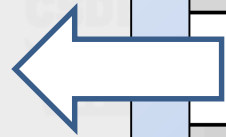
[Make RDF in](#) [DCI XML](#) [Home](#) [Help](#)

Found **88** records. Displayed records from **31** to **60**
[Previous 30 record\(s\)](#) [Next 30 record\(s\)](#)

[Expand all records](#) [Show all as text \(SweetDB notation\)](#) [Report data error](#)

< ПОЛНЫЕ ЗАПИСИ ЗДЕСЬ >

[Previous 30 record\(s\)](#) [Next 30 record\(s\)](#)
 Resort records by:
[New query](#) [Home](#) [Help](#)



Поиск конформаций

Search for disaccharide conformation maps

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

Conformation ID:

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

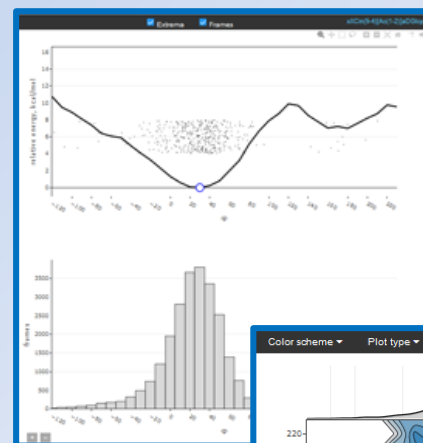
or type dimeric fragment in CSDB encoding

Force field:

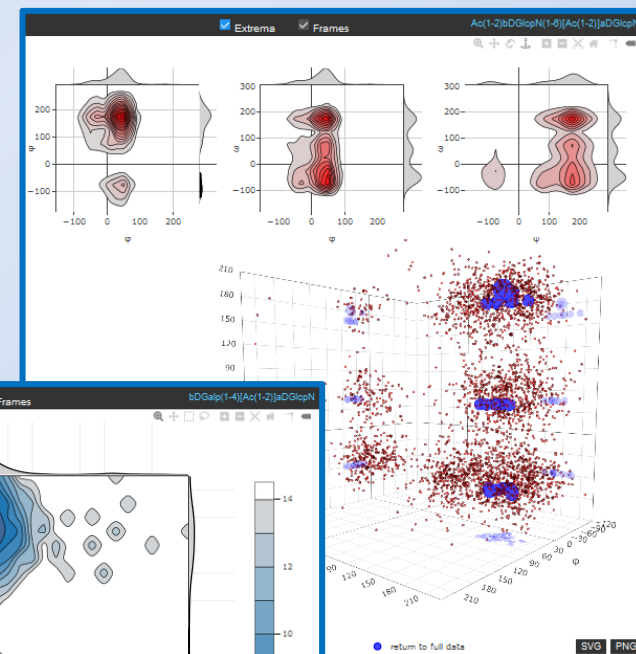
Temperature:

Solvent model:

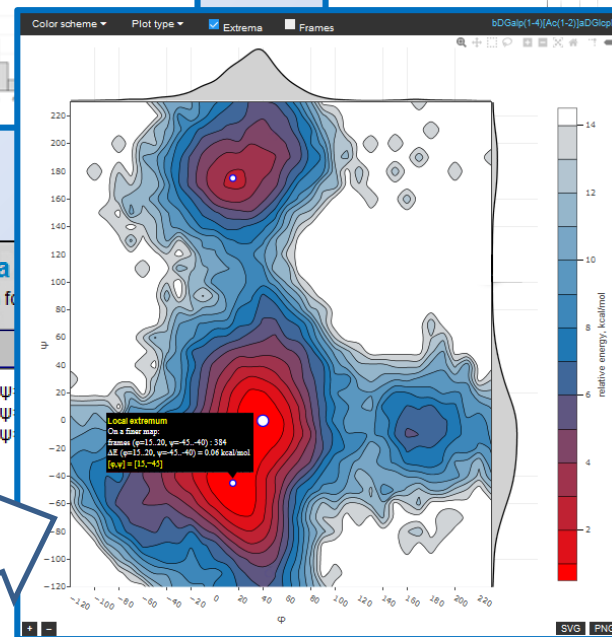
[Home](#) [Help](#)



1D



3D



2D

CSDB conformation data

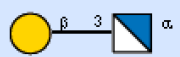
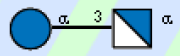
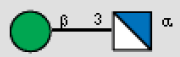
12 conformation maps have been found

Model structure	Conformation map	
		<p>$\varphi = -30, \psi = -35, \omega = -25$</p> <p>map live view</p> <p>ID: 1610</p>
		<p> $\varphi = 25, \psi = 165, \omega = -75 \Delta E = 0.00 \text{ Kcal/mol}$ $\varphi = 45, \psi = 170, \omega = -60 \Delta E = 0.00 \text{ Kcal/mol}$ $\varphi = 45, \psi = 195, \omega = -75 \Delta E = 0.38 \text{ Kcal/mol}$ $\varphi = 50, \psi = 185, \omega = -60 \Delta E = 0.72 \text{ Kcal/mol}$ $\varphi = 25, \psi = 165, \omega = 180 \Delta E = 0.72 \text{ Kcal/mol}$ $\varphi = 30, \psi = 155, \omega = -60 \Delta E = 0.72 \text{ Kcal/mol}$ $\varphi = 35, \psi = 180, \omega = 180 \Delta E = 0.85 \text{ Kcal/mol}$ $\varphi = 45, \psi = 215, \omega = -60 \Delta E = 0.85 \text{ Kcal/mol}$ $\varphi = 40, \psi = 170, \omega = 165 \Delta E = 0.99 \text{ Kcal/mol}$ $\varphi = 20, \psi = 185, \omega = 165 \Delta E = 0.99 \text{ Kcal/mol}$ $\varphi = 55, \psi = 195, \omega = 165 \Delta E = 0.99 \text{ Kcal/mol}$ $\varphi = 55, \psi = 155, \omega = -60 \Delta E = 1.13 \text{ Kcal/mol}$ $\varphi = 45, \psi = 185, \omega = 165 \Delta E = 1.13 \text{ Kcal/mol}$ $\varphi = 30, \psi = 145, \omega = -60 \Delta E = 1.13 \text{ Kcal/mol}$ </p> <p>Force field: MM3-1996 Solvent model: None MD temperature: 1000 MD duration: 30 ns Frames: 30K MD summary file: download</p> <p>map live view</p> <p>ID: 907</p>

Поиск гликозилтрансфераз

CSDB glycosyltransferase search

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

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


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](#) 

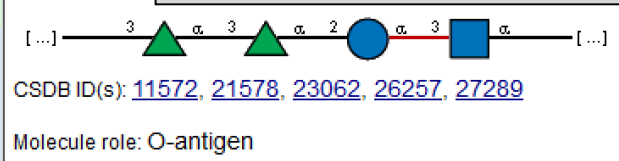
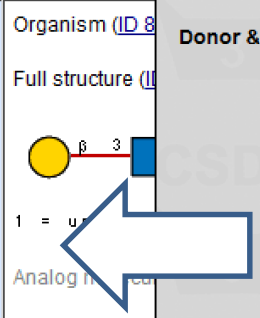
Donor & acceptor: Type donor CSDB encoding or use tools
 [Use Wizard](#) 
Type acceptor CSDB encoding or use tools
 [Use Wizard](#) 

Treat donor/acceptor as fragments

Confirmation status: Filter results to those

[Search!](#)

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



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

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

Фрагменты и таксоны

CSDB dimer abundance

The table lists **17** dimeric fragments present in **8** saccharides associated with **2** organisms from: *Herpetomonas muscarum*, *Herpetomonas samuelpeessoai* (species).

Anomeric forms were combined. Dimers containing monovalent constituents were excluded. Residues with undefined configurations or ringsizes are greyed. Superclasses are in blue. To re-sort the list click the according column name.

Position	Donor	Linkage	Acceptor	Abundance	Compound IDs	Abundance in selected species
terminal	DManp	1-2	DManp	18 (17%)	5555 , 5556 , 5557 , 5559 , 5560 , 5848	<i>Herpetomonas samuelpeessoai</i> : 12 (67%) <i>Herpetomonas muscarum</i> : 6 (33%)
di-branched	DManp	1-3	DManp	15 (14%)	5345 , 5555 , 5556 , 5557 , 5558 , 5559 , 5560 , 5848	<i>Herpetomonas samuelpeessoai</i> : 11 (73%) <i>Herpetomonas muscarum</i> : 4 (27%)
inline linear	DManp	1-2	DManp	3 (3%)		<i>Herpetomonas muscarum</i> : 3 (43%)
			DManp	5 (4.7%)	5556 , 5557 , 5559 , 5560	<i>Herpetomonas samuelpeessoai</i> : 4 (80%) <i>Herpetomonas muscarum</i> : 1 (20%)
tri-branched, reducing end	DGlcPn	1-4	DGlcN	5 (4.7%)	5345 , 5557 , 5558 , 5559 , 5560	<i>Herpetomonas samuelpeessoai</i> : 5 (100%)
di-branched, reducing end	DGlcPn	aglycone	diphosphodolichol	4 (3.7%)	5555 , 5556	<i>Herpetomonas samuelpeessoai</i> : 2 (50%)
tri-branched	DGlcPn	1-4	DGlcPn	4 (3.7%)	5555 , 5556	
tri-branched	DManp	1-3	DGlcPn	2 (1.9%)	5848	
penta-branched	DGlcPn	1-4	DGlcPn	2 (1.9%)		
inline linear	DManp	1-3	DManp	1 (0.9%)	5557	
terminal	DGlcP	1-2	DManp	1 (0.9%)	5558	
<i>Total</i>				<i>107 (100%)</i>		

[Export TSV](#)

[Monomers](#)

[Home](#)

Monomer and dimer abundance

This page will generate a pool of monomers or dimers abundant in glycans from the selected taxonomic group(s). First, please select a taxonomic rank of a desired group: species

Display groups: bacteria archaea protista algae fungi plants animals

Genus:
(select only one)

- Herbaspirillum
- Herbidospora
- Herpetomonas**
- Histophilus
- Hyphomonas
- Idiomarina
- Inquilinus
- Kaistella
- Kineosporia
- Kingella
- Klebsiella
- Kocuria

Species:
(select multiple with CTRL key)

- sp. (unassigned)
- muscarum**
- samuelpeessoai**

Select all

- Combine anomeric forms
 - Include undefined configs
 - Include ONLY saccharides
 - Include monovalent residues
 - Include aglycons in oligomers
 - Include aliases
 - Explain 'Subst' aliases
-
- Distinguish inline / terminal / reducing
 - Distinguish branching degree
 - ...and ignore monovalent substituents

Display only those fragments that are unique for this species in all biota

[Display monomers](#)

[Display dimers](#)

[Monomer namespace](#)

[Home](#)

[Help](#)

Кластеризация таксонов

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**

Предсказание структуры

Structure generation constraints:

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

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

Allowed linkages: C1 C2 C3 C4 C5 C6 C7+

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"/>
?	?	reducing			<input checked="" type="checkbox"/>
?	?	any	forbidden	1	<input checked="" type="checkbox"/>

Search depth: Widespread structures only

Scope: oligomers polymers Δ

Advanced scope: β -anomers: = 1 CH₂ 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%				1 = b-D-GalpNA 2 = L-Ala
#2. $\Delta \sim 0.95$ ppm Corr = 1.000 RMS dev = 1.46 ppm Trust = 46%				1 = b-D-GalpNA 2 = L-Ala
#15. $\Delta \sim 1.42$ ppm Corr = 0.999 RMS dev = 1.99 ppm Trust = 49%				1 = b-D-GalpNA 2 = L-Ala

Find best matching structures:

Experimental ¹³C 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

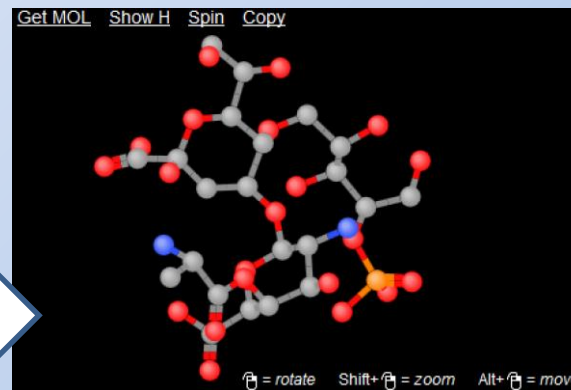
\pm 2 signals

Find 15 best-fitting structures

Save generated structures

[Go!](#)

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



Соответствие специальности

Методологические компьютерные работы традиционно защищаются по специальности, для которой предназначены создаваемые инструменты.

Формулировки, охватываемые темой работы, в паспорте специальности **02.00.10** «биоорганическая химия» и их расшифровки в программе специальности :

- «Изучение структуры и функций биомолекул [...] физико-химическими методами; структурно-функциональные исследования полисахаридов и смешанных биополимеров».
- «Моно-, олиго-, полисахариды: номенклатура, стереохимия, конформация. Методы изучения строения».
- «Соединения из микроорганизмов, грибов, водорослей, растений, [...] липополисахариды бактерий».
- «Спектроскопия ЯМР, [...] связь параметров спектров ЯМР с химической и пространственной структурой биомолекул, [...] двумерная спектроскопия ЯМР».
- «Компьютерное моделирование молекулярной механики биомолекул».
- формулировки, связанные с компьютерным моделированием и расчётами геометрии *белков*. На момент составления паспорта специальности моделирование *гликанов* ещё не было популярно.

Кроме этого, работа расширяет методологию биоорганической химии на новые области, не указанные в описании специальности.