

Орг. вопросы

- Сайт <http://intbio.org/bioinf2019-2020>

ВВЕДЕНИЕ В БИОИНФОРМАТИКУ

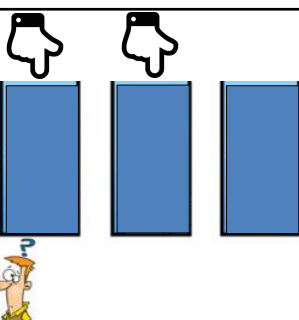
Лекции №2-3

Понятие информации. Способы измерения информации. ДНК как цифровой носитель информации. Системы счисления. Теория информации. Информационная энтропия. Сжатие информации. Теорема Котельникова. Теорема Шеннона-Хартли. Шифрование информации. Хранение информации. Источники больших данных в биомедицине. Проблемы передачи больших данных.

Алексей Константинович Шайтан, к.ф.-м.н.

Сайт курса: <http://intbio.org/bioinf2019-2020>

13 сентября 2019



- В результате игрок получает **0.67 бита информации**
- Если бы ведущий открыл дверь в самом начале – только **0.58 бита информации**

https://en.wikipedia.org/wiki/Monty_Hall_problem

Природа информации

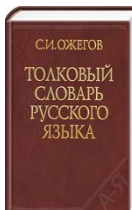
Информация

ИНФОРМАЦИЯ, 1. Сведения об окружающем мире и протекающих в нем процессах, воспринимаемые человеком или специальным устройством.



H.H. Моисеев

... универсального определения информации не только нет, но и быть не может из-за широты этого понятия.



Norbert Wiener

Information is information, not matter or energy.

http://www.aselibr.ru/press_center/journal/2007/number_3/number_3_6/number_3_657/

Информация



iPhone 7 Plus Серебристый Ёмкость

Теперь выберите ёмкость.

32 Гб¹
52 990.00 руб.

Доставка: на складе

128 Гб¹
60 990.00 руб.

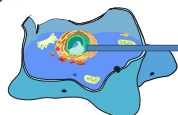
Доставка: на складе

Живые системы и цифровые технологии

Аналоговые технологии

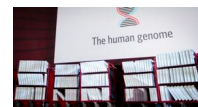
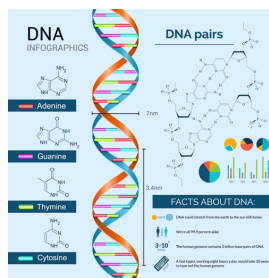


Цифровые технологии



ДНК – цифровой код

Генетический код



Геном человека
3 млрд нуклеотидов
262 тыс страниц
3 Гбайта



Геном кишечной палочки
4.6 млн нуклеотидов
400 страниц
4.6 Мбайт

1 bit

0 1 0 1 1 0 1 1 1 0

1 byte = 8 bits

Двоичный код

ATGC

Четвертичный код

Код может быть **позиционный** или **непозиционный**

Позиционный код активно используется в **системах счисления**

Соответствие цифр некоторых систем счисления

Основа системы счисления	2	8	10	16
0	0	0	0	0
1	1	1	1	1
10	2	2	2	2
11	3	3	3	3
100	4	4	4	4
101	5	5	5	5
110	6	6	6	6
111	7	7	7	7
1000	10	8	8	8
1001	11	9	9	9
1010	12	10	A	A
1011	13	11	B	B
1100	14	12	C	C
1101	15	13	D	D
1110	16	14	E	E
1111	17	15	F	F

Целое число без знака x в b -ичной системе счисления представляется в виде конечной линейной комбинации степеней числа b^k :


$$x = \sum_{k=0}^{n-1} a_k b^k, \text{ где } a_k \text{ — это целые числа, называемые цифрами, удовлетворяющие неравенству } 0 \leq a_k \leq b - 1.$$

TIME

SCIENCE

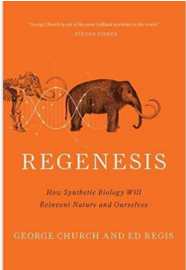
The First Book To Be Encoded in DNA

Two Harvard scientists have produced 70 billion copies of a book in DNA code --and it's smaller than the size of your thumbnail.



Lisa Poole / AP FILE

In his lab at the Harvard Medical School in Boston, George Church,



REGENESIS

How Synthetic Biology Will Redefine Nature and Ourselves

GEORGE CHURCH AND ED REGIS

Измерение информации

1 bit

0 1 0 1 1 0 1 1 1 0

1 byte = 8 bits

Измерения в байтах				Decimal	Binary
		ГОСТ 8.417—2002	Приставки СИ		
Название	Обозначение	Степень	Название	Степень	
байт	Б	10 ⁰	-	10 ⁰	0000
килобайт	кбайт	10 ³	кило-	10 ³	0001
мегабайт	Мбайт	10 ⁶	мега-	10 ⁶	0010
гигабайт	Гбайт	10 ⁹	гига-	10 ⁹	0011
терабайт	Тбайт	10 ¹²	тера-	10 ¹²	0100
петабайт	Пбайт	10 ¹⁵	пета-	10 ¹⁵	0101
эксабайт	Эбайт	10 ¹⁸	экса-	10 ¹⁸	0110
зеттабайт	Збайт	10 ²¹	зетта-	10 ²¹	0111
иоттабайт	Ибайт	10 ²⁴	иотта-	10 ²⁴	1000
					1001

Вероятность того, что машина за этой дверью: $\frac{1}{3}$ $\frac{1}{3}$ $\frac{1}{3}$



- Сколько информации нужно, чтобы закодировать положение машины?
1.5849625007211563... бит

Теория информации

“the father of [information theory](#)”



Claude Elwood Shannon
(April 30, 1916 – February 24, 2001)

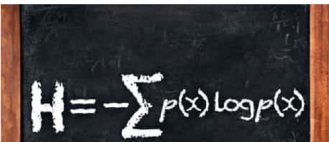
https://youtu.be/z2Whi_nL-x8

Теория информации

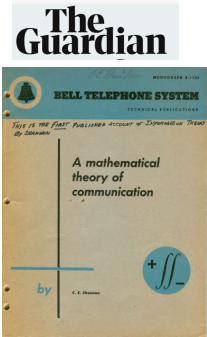
Science: A short history of equations

Without Claude Shannon's information theory there would have been no internet

It showed how to make communications faster and take up less space on a hard disk, making the internet possible



Введена мера информации(!)
кг, метр, секунда + БИТ



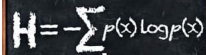
1948

<https://www.khanacademy.org/computing/computer-science/information-theory/info-theory/v/intro-information-theory>

Информационная энтропия

[Клод Шеннон](#) предположил, что прирост информации равен утраченной неопределённости, и задал требования к её измерению:

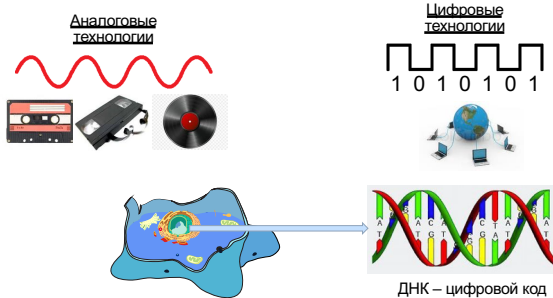
- мера должна быть непрерывной; то есть изменение значения величины вероятности на малую величину должно вызывать малое результирующее изменение функции;
- в случае, когда все варианты (буквы в приведённом примере) равновероятны, увеличение количества вариантов (букв) должно всегда увеличивать значение функции;
- должна быть возможность сделать выбор (в нашем примере букв) в два шага, в которых значение функции конечного результата должно являться суммой функций промежуточных результатов.



Живые системы и цифровые технологии - аналогии

От программирования компьютеров к программированию живых систем

Живые системы и цифровые технологии



Генетический код

Геном человека
3 млрд нуклеотидов
262 тыс страниц
3 Гбайта

Геном кишечной палочки
4.6 млн нуклеотидов
400 страниц
4.6 Мбайт

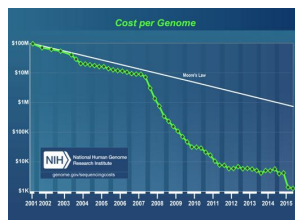
Секвенирование ДНК (метод Сэнгера)

Sanger Sequencing

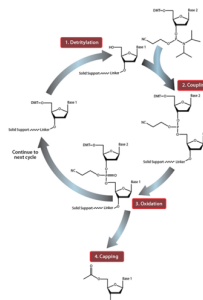
Sanger sequencing uses ddNTPs (dideoxynucleotide triphosphates) which do not have a free 3' OH mixed in with dNTPs. Whenever the DNA polymerase incorporates a ddNTP it won't be able to add any other nucleotides. Then gel electrophoresis is used to separate the DNA fragments.

<https://www.youtube.com/watch?v=593zWZNwbJI>

Технологии чтения и записи ДНК

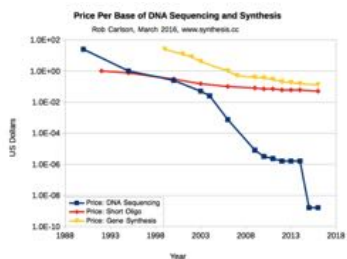


Синтез ДНК

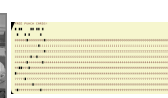


<https://www.sigmaaldrich.com/technical-documents/articles/biology/dna-oligonucleotide-synthesis.html>

Удешевление синтеза ДНК



От программирования компьютеров к программированию жизни



Ядро ОС Linux
~100 МБ



????

Аналогии компьютеров и живых клеток

- Software boots into hardware
- Genetic code boots into cells




2010 *Mycoplasma mycoides* JCVI-syn1.0

2016 Syn3.0

Human	~20000-25000
E. coli genome (K12 strain)	~4800
Syn 1.0	301
<i>Mycoplasma genitalium</i> *	525
Syn 3.0	473



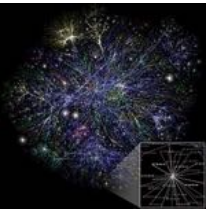


Craig Venter

Первый искусственный геном (~\$40 млн)

Минимальный геном

От программирования компьютеров к программированию жизни

~100 млрд. (10^{11}) нейронов
Каждый нейрон имеет
~7000 синаптических связей

23 млрд. устройств
подключенных к интернету

Сложные инженерные системы






~2500 лет
д.н.э.

Boeing 747
6 млн частей

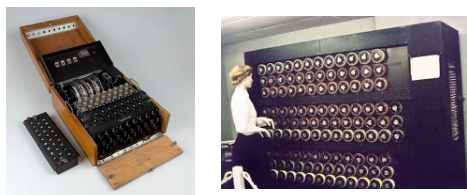
Boeing 777
Первый самолет спроектированный полностью на компьютере

Intel Xeon Phi
CPU
8*10⁹
транзисторов

Передача информации

Передача информации

- Комплексная область: **теоретические, практические, физические** аспекты
- Вопросы сжатия данных
- Вопросы надежности
- Вопросы **шифрования и защиты данных** (особенно в медицине и биологии)



https://en.wikipedia.org/wiki/The_Imitation_Game

Аналоговый сигнал



Цифровой сигнал



Как конвертировать?

Передача информации

Связь частоты сигнала и пропускной способности

Владимир Александрович Котельников



1908 - 2005

Harry Nyquist

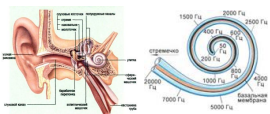


1889 - 1976

Теорема Котельникова-(Найквиста-Шенона)

«любую функцию $F(t)$, состоящую из частот от 0 до f , можно непрерывно передавать с любой точностью при помощи чисел, следующих друг за другом через $1/(2f)$ секунд»

https://ru.wikipedia.org/wiki/Теорема_Котельникова



44.1 кГц – частота дискретизации при записи звука

ВВЕДЕНИЕ В РЯДЫ ФУРЬЕ



Говоря шире, теорема Котельникова утверждает, что непрерывный сигнал $x(t)$ можно представить в виде интерполяционного ряда:

$$x(t) = \sum_{k=-\infty}^{\infty} x(k\Delta) \operatorname{sinc} \left[\frac{\pi}{\Delta} (t - k\Delta) \right],$$


где $\operatorname{sinc}(x) = \sin(x)/x$ – функция sinc. Интервал дискретизации удовлетворяет ограничениям

$0 < \Delta \leq \frac{1}{2f_c}$. Мгновенные значения данного ряда есть дискретные отсчёты сигнала $x(k\Delta)$.

Передача информации

Связь частоты сигнала и пропускной способности

Ralph Hartley



$$C = B \log_2 \left(1 + \frac{S}{N} \right),$$

где

- C — пропускная способность канала, бит/с;
- B — полоса пропускания канала, Гц;
- S — полная мощность сигнала над полосой пропускания, Вт или В²;
- N — полная шумовая мощность над полосой пропускания, Вт или В²;
- S/N — отношение мощности сигнала к шуму (SNR).

1888 – 1970


Теорема Шеннона-Хартли

https://ru.wikipedia.org/wiki/Теорема_Шеннона_—_Хартли

Передача информации




Оптоволокно



Антенны КНЧ

Рис. 1. Полосы пропускания линий связи и популярные частотные диапазоны

https://ru.wikipedia.org/wiki/Связь_с_поворотными_порками

Каналы связи



Карта подводных кабелей

<https://habrhabr.ru/company/rookwell/blog/305634/>

Шифрование информации



https://en.wikipedia.org/wiki/Public-key_cryptography

Криптосистемы с открытым ключом



Необратимая Хэш функция

```
mbptb:~ alexsha$ md5 -s 'Hello world!!!'
MD5 ("Hello world!!!") = 87ee732d831690f45b8606b1547bd09e
```

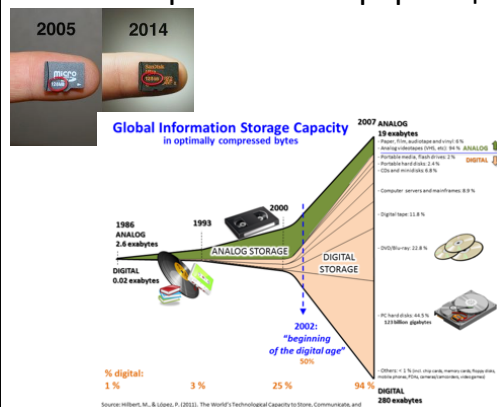
https://en.wikipedia.org/wiki/Public-key_cryptography

Хранение информации

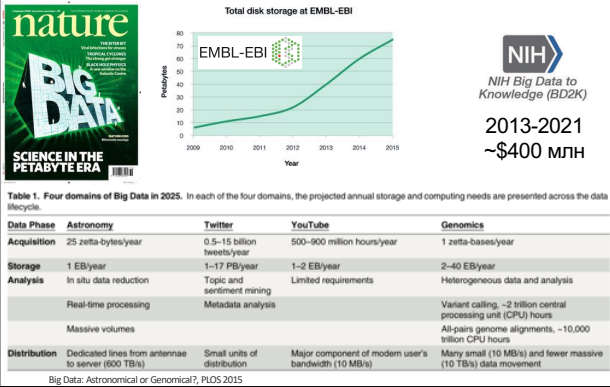
Хранение информации



Хранение информации



Большие данные в биомедицине



Источники больших данных в биомедицине

- **Омиксные технологии**
 - Секвенирование, геномика, транскриптомика, протеомика, метаболомика и т.д.
 - Коннектом мозга
- **Медицинская информация**
 - Электронные медицинские карты, результаты клинических исследований и т.д.
 - Медицинские изображения, МРТ и т.д.
- **Структурная биология и моделирование**
 - Данные с лазеров на свободных электронах (XFEL)
 - Моделирование структуры и динамики белков.

Данные секвенирования, пример Геномы раковых опухолей



Геном человека ~ 3.3 Gb
x100 секвенирование ~300Gb

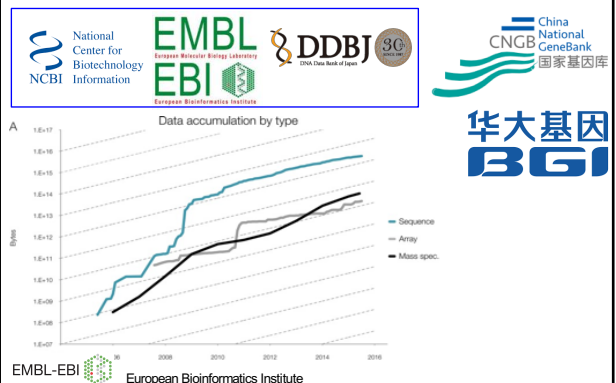


25000 образцов опухолей

Международный проект, данные распределены по миру

ID	Donor	Repository	Project	Study	Data Type	Storage	Size	URL
119995	00217962	PCAWG - London, PCAWG - Barcelona, Collaboratory - Toronto, USA - Houston	BRCA1	Aligned Reads	WGS	SAW	128.72 GB	
119994	00217962	PCAWG - London, PCAWG - Barcelona, Collaboratory - Toronto, USA - Houston	BRCA2	Aligned Reads	WGS	SAW	107.27 GB	
119994	0048280	PCAWG - London, PCAWG - Barcelona, USA - Houston, Collaboratory - Toronto, USA - Virginia	OVAR2	Aligned Reads	WGS	SAW	134.05 GB	
119993	0048280	PCAWG - London, PCAWG - Barcelona, Collaboratory - Toronto, USA - Virginia, USA - Houston	OVAR1	Aligned Reads	WGS	SAW	101.45 GB	
119995	00222828	PCAWG - Chicago (TCGA), RDC - Chicago	DLBC15	Aligned Reads	WGS	SAW	202.36 GB	
119995	00222828	PCAWG - Chicago (TCGA), RDC - Chicago	DLBC15	Aligned Reads	WGS	SAW	100.71 GB	

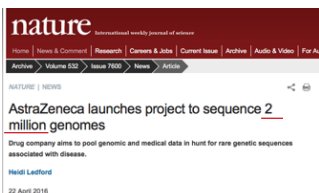
Централизованные репозитории омиксных данных



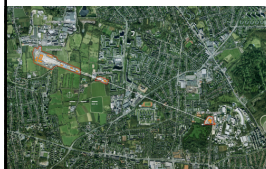
Genomes en masse



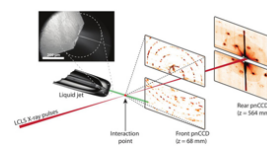
5 years ~ 100 000 genomes



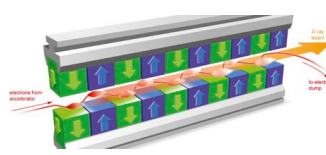
Структурная биология и моделирование



European XFEL, Hamburg



27000 импульсов в секунду



Передача больших данных

Выделенные научно-образовательные сети 100Gbit/s



Figure 1: Asia Pacific Ring (APR)



Программные решения



Базы данных

Базы данных

- Реляционные базы данных, объектно-ориентированные, RDF
- Системы управления базами данных СУБД
- Языки и стандарты SQL, SPARQL, RDF



Реляционные базы данных

Клиенты				Товары	
Id_кл	Фамилия	Имя	Отчество	Id_тов	Название
15	Иванов	Иван	Иванович	1	Шоф
16	Петров	Петр	Петрович	2	Стул
17	Николаев	Николай	Николаевич	3	Стол

Заказы				
Id_зак	Клиент	Товар	Дата	Количество
1	15	1	15.09.2003	1
2	17	1	17.09.2003	2
3	15	2	20.09.2003	12

SQL

Целостность данных

Транзакции

Соответствие требованиям ACID

Атомарность

Изолированность

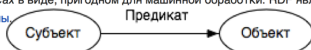
Надежность

<https://aws.amazon.com/ru/relational-database/>

Resource Description Framework

Язык SPARQL

Resource Description Framework (RDF, «среда описания ресурса»^[1]) — это разработанная консорциумом Всемирной паутины модель для представления данных, в особенности — метаданных^[2]. RDF представляет утверждения о ресурсах в виде, пригодном для машинной обработки. RDF является частью концепции семантической паутины.



UniProt Downloads

SPARQL Downloads

Your SPARQL query

Add common prefixes

```

1 PREFIX up: <http://purl.uniprot.org/owa/>
2 PREFIX biochem: <http://purl.uniprot.org/taxonomy/>
3 PREFIX owl: <http://www.w3.org/2002/07/owl#>
4 SELECT ?protein ?text
5 WHERE
6 {
7   ?protein a up:Protein .
8   ?protein up:organism taxon:9606 .
9   ?protein up:annotation ?annotation .
10  ?annotation a up:annotation:SequenceAnnotation .
11  ?annotation up:comment ?text .
12  FILTER (CONTAINS(?text, 'loss of function'))
13 }
14
  
```

Select all human UniProt entries with a sequence variant that leads to a 'loss of function'

https://ru.wikipedia.org/wiki/Resource_Description_Framework/

Биологические базы данных

Бiology is a data-intensive science!

- Нужно уметь хранить данные
- Нужно уметь обрабатывать данные
- Нужно уметь обмениваться данными
- Данные должны быть максимально открыты и доступны научному сообществу.
- Data provenance ("происхождение данных")

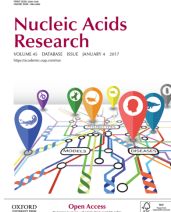
Data provenance [\[edit \]](#)

Scientific research is generally held to be of good provenance when it is documented in detail sufficient to allow reproducibility.^{[27][28]} Scientific workflow systems assist scientists and programmers with tracking their data through all transformations, analyses, and interpretations. Data sets are reliable when the process used to create them are

- Кризис воспроизводимости результатов в науке!?

Базы данных для биологии

- На данный момент количество не возможно сосчитать – очень много – важно не запутаться и не потеряться при их использовании
- Надежные источники информации о базах данных – научные журналы



Annual Database Issue – информация о ~200 БД каждый год.

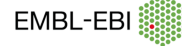
Базы данных для биологии

- Бесплатные vs Платные (по подписке)
- Свободно доступные vs Ограниченно доступные
- Большие ресурсы (NCBI, EBI/EMBL, etc.) интегрирующие многие базы данных - поддерживаются государством
- Коллаборации между университетами (напр. PDB)
- Коммерческие компании
- Локальные базы данных, поддерживаемые силами научных групп
- База данных vs Web Server – граница размыта.
- Хорошие БД - информационные ресурсы с возможностями сложного поиска и моделирования.

Крупные центры биологических БД



- Bethesda, MD USA
- Более 60 БД включая PubMed, GenBank, DBGap, SRA



- European Bioinformatics Institute, Cambridge, UK + Switzerland

Что храниться?

- БД статей, абстрактов, патентов
- Последовательности ДНК
- Последовательности белков
- 3D структуры молекул
- Геномы
- Данные экспрессии
- Сырые данные с секвенаторов
- Информация о химических соединениях и их активности
- Информация о болезнях, информация о пациентах
- Информация о видах живых организмов
- Информация о метаболических и сигнальных путях
- Информация о взаимодействии молекул
- Много производной информации: базы гомологичных последовательностей, аннотация отдельных классов белков и т.д.

План

- **Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)**
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

Реферативные базы данных

Clinical/Biomedical

PubMed – US National Library of Medicine database (Medline); refers to >25M articles from 5600 biomedical journals, 1940s to present, with some older items, in medicine, nursing, dentistry, veterinary medicine, allied health & pre-clinical sciences
 - bibliographic database with author-provided abstracts, added indexing terms from **MeSH** (Medical Subject Headings) thesaurus, & links to other resources

www.pubmed.gov



FREE

Реферативные базы данных

Clinical/Biomedical

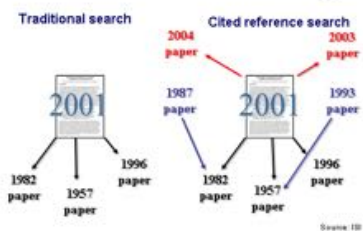
Embase – European based, includes all of Medline (database behind PubMed) and more; > 29M records, >8,500 journals, 1940s to present; includes coverage of more basic science journals & pre-clinical topics - especially useful for drug pipeline information, biotechnology, medical devices, conference coverage, toxicology, health policy/management, & alternative/complementary medicine
 Emtree thesaurus includes almost twice as many terms as PubMed

<https://www.elsevier.com/solutions/embase-biomedical-research>

ELSEVIER PAID

Реферативные базы данных

Cited Reference Searching



Source: IIR

Реферативные базы данных

Общенаучные базы данных цитирований

Web of Science - covers >12,000 journals from 1900 to present; useful for cited reference, **conference information & affiliations** (institutions)

<https://webofknowledge.com/>



Scopus – covers >18,500 journals from 1823 to present, complete citation counts for indexed articles 1996 to present; a general science database, not a specialized database – useful for cited reference, **conference information & affiliations** (institutions)

<https://www.scopus.com/>



Реферативные базы данных

Общенаучные базы данных цитирований

Elibrary.ru/РИНЦ

Реферативные базы данных

Общенаучные базы данных цитирований



Articles Case law

Recommended articles

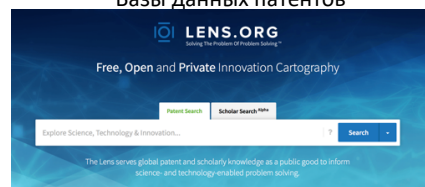
Exploring DNA dynamics within oligonucleosomes with coarse-grained simulations: SIRAH force field extension for protein-DNA complexes
A Brandner, A Schüller, F Melo, S Pantano - Biochemical and biophysical research ..., 2017

Базы данных диссертаций

Open DOAR <http://www.opendoar.org/index.html> ;
 OpenThesis <http://www.openthesis.org/> ;
 BASE – Bielefeld Academic Search Engine -
<http://www.base-search.net/>
 > refine search result > document type > theses

ProQuest Dissertations & Theses
 Database <http://www.proquest.com/products-services/pgdt.html> - from 1743
 to present; some fulltext since 1990; fee with some
 free search capability

Базы данных патентов



http://www.lens.org/lens/biological_search – ПОИСК
 ДНК последовательностей



Search and read the full text of patents from around the world.

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- **Базы данных последовательностей ДНК**
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

Базы данных нуклеотидных последовательностей

Нуклеотидные БД – это хранилища, принимающие данные от научного сообщества и представляющие их широкой общественности. Различные БД отличаются по источнику последовательностей, их надежности, широте аннотирования и т.д. В идеале БД должна содержать все известные последовательности.

The **International Nucleotide Sequence Database Collaboration** – совместный проект EMBL-Bank в Европейском Институте Биоинформатики (EBI), японского банка данных ДНК (DDBJ) в Центре Информационной Биологии (CIB) и GenBank в Национальном Центре Биотехнологической Информации (NCBI).



72

База данных GenBank

Открытая БД нуклеотидных последовательностей, учреждена в 1982 г.
 2017: > 300 000 организмов, ~ 203 млн. последовательностей,
 ~ 240 млрд. пар оснований

База данных GenBank. Структура файла

```

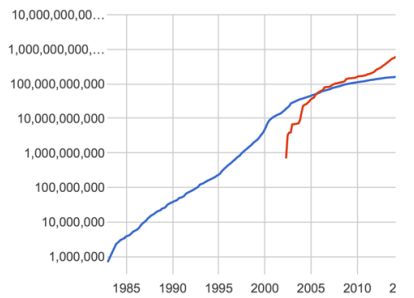
...
FT      /translation="MGQPGNGSAFLFLAPNGSHAPDHDVTTQQRDEVVVVGGMVMSLIVL
FT      AIVFGNVLVITAIKFERLQVTNYFITSLACADLVMLAVVPPGAAHILMKMNTFGNF
FT      WCFEFTSIDVLCVTASIEITLCVAVDRYFAITSPFKYQSLTKKARVILMVMIVSGL
FT      TSFLPIQMHWRATHQEAINCYANETCCDFFTNQAYAIASSIVSFYPLVMVFVYSRV
FT      FQEAQRQLQIKDSEGRFHVQNLQVEQDGRGTGHGLRBSKFLCKEKALKTLGIMGT
FT      FTLCWLPFFIVNIHVHVIQDNLIRKEVILLNWIGYVNSGFNPLYCRSPDFRIAFQELL
FT      CLRRSSLKAYNGYSSNGNTGEGSGYHVEQEKENKLLCEDLPGTDFVGHQGTVPDNI
FT      DSQGRNCSTNDSLK
FT      variation      46
FT      /gene="ADRB2"
FT      /replace="ax
FT      /note="Arg16 to Gly polymorphism"
XX
...
    
```

База данных GenBank. Структура файла

```

...
SQ Sequence 1242 BP; 275 A; 331 C; 326 G; 310 T; 0 other;
atggggcaac ccgggaacgg cagcgccttc ttgctggcac ccaatggaag ccatgcgccg 60
gaccacgacg tcaacgacga aaggacgag gttgggtgg tggcatggg catcgtcatg 120
tccttcacg cctcggccat cgtgtttggc aatgtgctgg tcatcacagc cattgccaaag 180
ttcagacgtc tgcagacggt caccactac ttcactact cactggcctg tgcgtgatcg 240
gtcatgggcc tggcagtggt gccctttggg gccgccata ttcttatgaa aatgtggact 300
tttggcaact tctgttgoga gttttggact tccattgatg tgcgttgctg caccggcacc 360
attgagaccc tgtcgtgat cgaatggat cgtactcttg ccattacttc acctttcaag 420
taccagacc tgcagccaa gaataaggcc cgggtgatca tctgtatggt gtgattgtg 480
tcaggcccta cctccttctt gccacttcag atgactcgtt accggccac ccaccaggaa 540
gccatcaact gctatccaa tgagactcgc tgtactcttc tcaagaacca agcctatgcc 600
attgcctctt ccatcgtgtc ctctacggt cccctgggtga tcaatgtctt cgtactactc 660
agggtctttc aggagcccaa aaggcagctc cagaagattg acaaatctga gggcgccttc 720
catgtccaga accctagcca gttggagcag gatggggcga cggggcatgg actccgcaga 780
tcttccaagt tctcctgaa ggagacaaa gccctcaaga cgttaggcat catcatgggc 840
acttcaacc tctcgtggct gccctctctc atcgttaaca ttgtgatgt gatccaggat 900
aaacctatcc gtaaggaaat ttacatcttc ctaaatgga tagcctatgt caattctggt 960
ttcaatccc tttctactg ccggagccca gatttcaggg ttgctctcca ggaactctg 1020
tgcctggcga ggtctctctt gaaggcctat ggaatgget actccagcaa cggcaacaca 1080
ggggagcaga gttgatatca cgtggaacag gaaagaata ataaactcgt gttgtgaagc 1140
ctcccaggca cyyaagaact tgtgggacat caaagtgact gctctagoga taaacttgat 1200
tcacaaggga gyaattgtag tacaatgac tctactcgtg aa 1242
//
    
```

Bases



GenBank and WGS Statistics

<https://www.ncbi.nlm.nih.gov/genbank/statistics/>

Genbank – is an archive! Contains everything.

Nicotiana tabacum chloroplast JLA region, sequence 2

GenBank: Z71230.1

[FASTA](#) [Graphics](#)

```

FEATURES             Location/Qualifiers
     source            1..124
                        /organism="Nicotiana tabacum"
                        /organelle="plastid:chloroplast"
                        /mol_type="genomic DNA"
                        /isolate="Cuban cahibo cigar, gift from President Fidel
                        Castro"
                        /db_xref="taxon:4097"

```

RefSeq – is a reference sequence database!

RefSeq – is a reference sequence database!

Если нужен список последовательностей всех генов человека – это вопрос к RefSeq, а не GenBank!

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- **Базы данных последовательностей белков**
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

GenBank/RefSeq is nucleotide centric, but

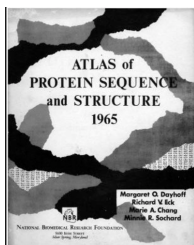
```

...
...
FT      /translation="MGQPGNGSAFLLPNGSHAPDHDVTQQRDEVVWVGMGIVMSLIVL
FT      AIVFGNVLVITAIKFERLQVTINYFITSLACADLVMLAVVFGAAHILMKWTFGNF
FT      WCFPWSIDVLCVTSASIELTLCVIAVDRYFAITSPFKYQSLLTNNKARVILMVMVIVSGL
FT      TSFLPIQMHWYRATHQEAINCYANETCCDFFTNGAYAIASSIVSYVPLVIMVFVSRV
FT      FQEAKRQLQKIDRSEGRFHVQNLQVQDGRTHGLRRSSKFLREHKALKTGIIMGT
FT      FTLCWLPFFIIVHVIQDNLIRREVIILLNIGVYNSGFNPLIYCRSPDFRIAFQELL
FT      CLRSSLKAYGNGCYSSNGNTGEQSGYHVEQEKENKLLCEDLPGETDFVGHQGTVPSPDNI
FT      DSQGRNCSTNDSLL*
FT      variation      46
FT      /gene="ADRB2*"
FT      /replace="ae
FT      /note="Arg16 to Gly polymorphism*
XX
...

```

Protein sequences are annotate within GB records ⁸⁰

Protein Centric Sequence Databases



Margaret Oakley Dayhoff
1925-1983

Margaret Dayhoff, a founder of the field of bioinformatics

Invented one-letter amino acid code, substitution matrices, etc.

https://en.wikipedia.org/wiki/Margaret_Oakley_Dayhoff

Protein Centric Sequence Databases



<http://pir.georgetown.edu>

In 2002, PIR along with its international partners, EBI (European Bioinformatics Institute) and SIB (Swiss Institute of Bioinformatics), were awarded a grant from NIH to create UniProt, a single worldwide database of protein sequence and function, by unifying the PIR-PSD, Swiss-Prot, and TrEMBL databases. As of 2010, PIR offers a wide variety of resources mainly oriented to assist the propagation and standardization of protein annotation: PIRSF,^[9] iProClass, and iProLINK.

The Protein Ontology (PRO) is another popular database released by the Protein Information Resource.^{[9][10]}

Белковые базы данных



UniProt – наиболее всеобъемлющий каталог информации о белках, объединяющий в себе данные из UniProtKB/Swiss-Prot, UniProtKB/TrEMBL и PIR.

83

Белковые базы данных



Качественно аннотированную информацию о белках
нужно искать в UniProtKB

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- **Базы данных 3D структур**
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

Структурные базы данных



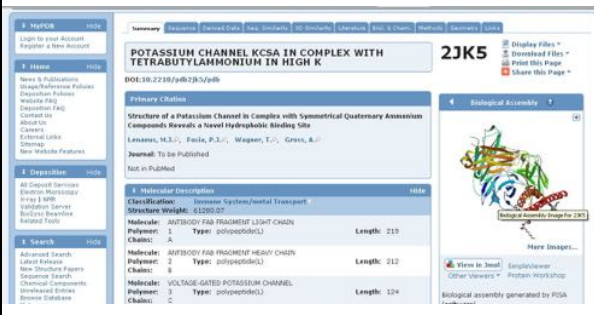
PDB – содержит информацию об экспериментально определенных структурах белков, нуклеиновых кислот и различных комплексов.

Структурные базы данных



Структурные базы данных

POTASSIUM CHANNEL KCSA IN COMPLEX WITH TETRABUTYLAMMONIUM IN HIGH K



Molecular Description	
Classification:	Ionocyte System/Animal Transport
Structure Weight:	61280.07
Molecule:	ANTIBODY FAB FRAGMENT LIGHT CHAIN
Polymer:	1 Type: polypeptide(L) Length: 219
Chain:	A
Molecule:	ANTIBODY FAB FRAGMENT HEAVY CHAIN
Polymer:	2 Type: polypeptide(L) Length: 212
Chain:	B
Molecule:	VOLTAGE-GATED POTASSIUM CHANNEL
Polymer:	3 Type: polypeptide(L) Length: 124
Chain:	C

Структурные базы данных

POTASSIUM CHANNEL KCSA IN COMPLEX WITH TETRABUTYLAMMONIUM IN HIGH K

Chain A: (polymer 1) [beta] [sheet] [alpha] [loop]

Description: ANTIBODY FAB FRAGMENT LIGHT CHAIN

Chain Type: poly(peptide)

Length: 219 residues

37% secondary structure (beta) [preference]

57% helix (4 helices); 12 residues

51% beta sheet (22 beta strands); 112 residues

Sequence Defaults

Residues: 1 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150 160 170 180 190 200 210 219

Sequence: MQQLGQPGALILVLPASVLLSEKASCTITFDWIHWKQDFKGLKELICETLPSKRAV
 NENIQKKAFLTADKSSSTAMQLSSELSSEDSAVVYKARERGDGFYFVWAGGTTVTYSKAK
 TTPPSVYPLFGSAAGENSMTYLGELKGLFPEYFVYVTVNNSLSSGVHTEFPAVLSGLY
 LSSLSVTFPSHWSPFSLVENVHPASSTAVDARLVFRQ

База данных PDB. Структура файла

```
HEADER      IMMUNE SYSTEM/METAL TRANSPORT              15-AUG-08   2JK5
TITLE      POTASSIUM CHANNEL KCSA IN COMPLEX WITH TETRABUTYLAMMONIUM
COMPND     2 IN HIGH K
COMPND     MOL_ID: 1;
COMPND     2 MOLECULE: ANTIBODY FAB FRAGMENT LIGHT CHAIN;
COMPND     3 CHAIN: A;
COMPND     4 ENGINEERED: YES;
COMPND     5 MOL_ID: 2;
...
KEYWDS     IMMUNE SYSTEM METAL TRANSPORT COMPLEX, QUATERNARY AMMONIUM,
...
EXPDTA     X-RAY DIFFRACTION
AUTHOR     M. J. LENAUS, P. J. FOZIA, T. WAGNER, A. GROSS
REVDAT     1 17-NOV-09 2JK5 0
JRNLS      AUTH M. J. LENAUS, P. J. FOZIA, T. WAGNER, A. GROSS
JRNLS      TITL STRUCTURE OF A POTASSIUM CHANNEL IN COMPLEX WITH
JRNLS      TITL 2 SYMMETRICAL QUATERNARY AMMONIUM COMPOUNDS REVEALS
JRNLS      TITL 3 A NOVEL HYDROPHOBIC BINDING SITE
JRNLS      REF TO BE PUBLISHED
JRNLS      REFW
REMARK     2
REMARK     2 RESOLUTION.      2.4  ANGSTROMS.
REMARK     3
REMARK     3 REFINEMENT.
REMARK     3 PROGRAM      : REFMAC 5.5.0051
...

```

Структурные базы данных

NDB – основана в 1992 г. для сбора и распространения информации о структуре нуклеиновых кислот. Формат хранения данных идентичен PDB.

WELCOME TO THE NUCLEIC ACID DATABASE
 a repository of three-dimensional structure information about nucleic acids

Number of Released Structures: 4623 Structures
 Last Update: 13-Oct-2010

Search the NDB by ID
 Enter an NDB ID or PDB ID [Search]

Search for Released Structures

Nucleic Acids Highlight

Структурные базы данных

EMBL-EBI

EMD-1367

Title: Three-dimensional structure of a voltage-gated potassium channel at 2.5 nm resolution.

Authors: Olga Sokolova, Ludmila Kolmakova-Partensky and Nikolay Grigorieff

Sample: Shaker-B channel

Aggregation state: Single particle (25 angstroms resolution)

Latest update: 2011-05-26

Summary

Experimental details

Visualization

Map information

Downloads

Status: Released

Deposition date: 2007-05-24

Header release date: 2007-05-30

Map release date: 2007-05-30

Primary citation: Sokolova O, Kolmakova-Partensky L, Grigorieff N. Three-dimensional structure of a voltage-gated potassium channel at 2.5 nm resolution. *STRUCTURE* (2007) 9, pp 215-220 [PubMed 15209500]

Sample: Shaker-B channel

Resolution: 25 Å (determined by FSC at 0.5 cut-off)

Filed PDB: PDB Authors

PubMed Status

Doyle, D.A., Cabral, J.M., Pfoetzner, 1188 R.A., Hood, A., Gubbis, J.M., Cohen, S.L., 952989 Released

Chak, S.L., Mackinnon, R.

Kreusch, A., Pfaffinger, P.J., Stevens, 1460 C.F., Choe, S. 9692076 Released

Структурные базы данных

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- **Базы данных хим. соединений**
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

Базы данных химических соединений

Базы данных химических соединений

Selected	Compound	Count
#BioAssays, Active		29
#BioAssays, Tested		60
#Protein 3D Structures		6
#Crystal Structure Of		1
#BioMedical Annotation		
Pharmacological Actions		69
#Synthetomatics		37
Biocystems		16
#Depositor Category		
Biological Properties		120
Chemical Vendors		65
Journal Publishers		87
NH Molecular Libraries		52

Базы данных химических соединений

Chemical Abstract Service – в регистре содержится 130 млн соединений (2018)

Базы данных химических соединений

Базы данных углеводов

For 2017:

7005 publications for 18924 compounds from 8859 organisms



99

Структурные базы данных

7009 структур липидов и сходных соединений – не поддерживается в настоящее время ☹

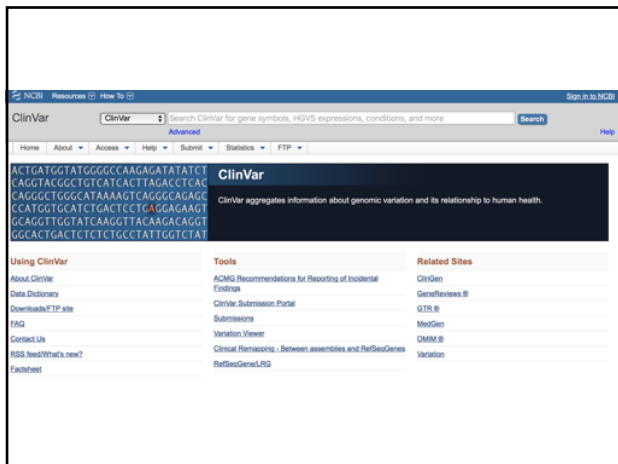
100

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- **Базы данных геномов и аннотаций**
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

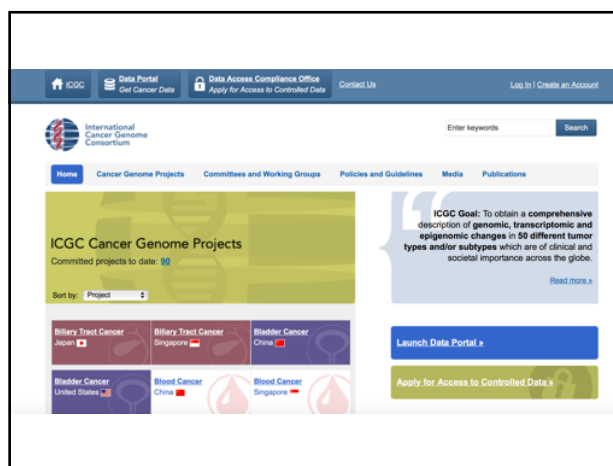
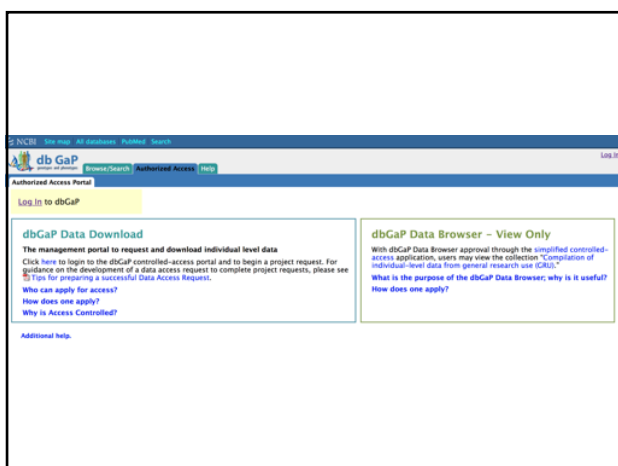
План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- **Базы данных вариаций генома**
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии



План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- **Базы данных геном-фенотип**
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии



План

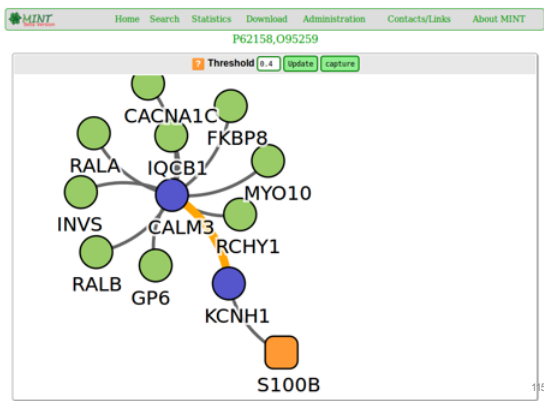
- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- **Базы данных взаимодействий**
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

Базы данных взаимодействий



114

Базы данных взаимодействий



115

Базы данных взаимодействий

Object	Name	UniProt ID	Sequence	FC5	Mass	Pub. Date	Activity
Lentiv guspol...	e-KT3.1.1 CMTX, CMTX L4	P13687	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	1804-102	4395.05	1998	Shaker-Kv1.1,Kv1.2
Lentiv guspol...	e-KT3.1.2 CMTX L4-q2, CMTX...	P49628	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	31.48	4335.08	1999	Kv1.4,Kv4.1
Meusobfus int...	e-KT3.1.3 BTK, Bcr/Abt...	P10483	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		4230.02	1990	Kv4.1
Meusobfus int...	e-KT3.1.4 BcrTK1, BcrTK...	Q9196	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	1893	4368.04	1997	Kv1.4,Kv1.2,Kv4.1
Meusobfus int...	e-KT3.1.4 BcrTK2, BcrTK...	Q9195	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	2847	4178.97	1997	Kv1.4,Kv1.2,Kv1.1,Kv1.3
Parabulbus int...	e-KT3.1.50 PBT3, Parabul...	P91117	EYVNERDQDQ...YLAIVLNTIK...Q...Y...F...		4274.23	2002	Kv1.1,Kv1.2,Kv1.1,Kv1.2
Centonides int...	e-KT3.1.151 BcrTK, BcrTK...	P10532	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		4095.98	2003	Kv4.1
Meusobfus int...	e-KT3.1.152 BcrTK, BcrTK...	C91432	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		4221.18	2005	Kv1.1
Meusobfus int...	e-KT3.1.171 BcrTK, BcrTK...	C91437	ZFTVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		4353.21	2005	Kv1.1
Centonides int...	e-KT3.2.1 NTK, NTK...	P08815	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	153M	4195.06	1992	Shaker-Kv1.1,Kv1.2
Centonides int...	e-KT3.2.2 NTK, NTK...	P40755	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	1M7X	4178.11	1993	Shaker-Kv1.1,Kv1.1
Centonides int...	e-KT3.2.3 CMTX L1, CMTX...	P49629	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		4391.06	1994	Kv1
Centonides int...	e-KT3.2.4 NTK-2, NTK-2...	Q97421	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		4383.05	1996	Kv1
Centonides int...	e-KT3.2.5 NTK-2, NTK-2...	P55847	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...	2M7Y	4231.25	1998	Kv1.1,Kv1.2,Kv1.3,Kv1.4
Centonides int...	e-KT3.2.7 CMTX L2, CMTX...	P49630	TETVQSTTAKQ...YLAIVLNTIK...Q...Y...F...		3909.82	1994	Kv1

116

Базы данных химических соединений

DRUGBANK

WHAT ARE YOU LOOKING FOR?

Drugs
Targets
Pathways
Indications

DRUGBANK

117

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- **Базы данных сигнальных/метаболических путей**
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

KEGG PATHWAY Database

Wiring diagrams of molecular interactions, reactions and relations

Menu [PATHWAY](#) [BRITE](#) [MODULE](#) [KO](#) [GENES](#) [LIGAND](#) [NETWORK](#) [DISEASE](#) [DRUG](#) [DBGET](#)

Select prefix Organism Enter keywords [Help](#)

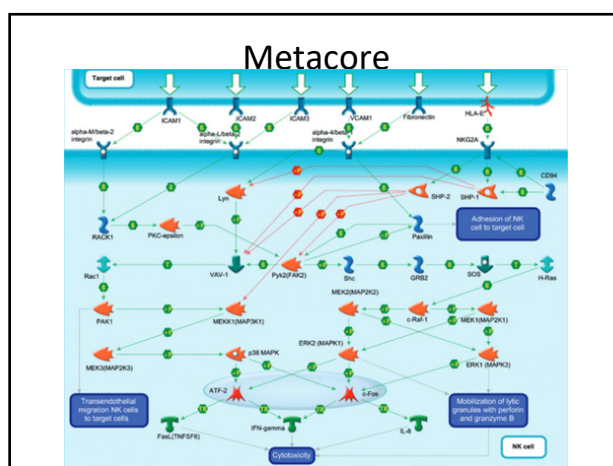
[[New pathway maps](#) | [Update history](#)]

Pathway Maps

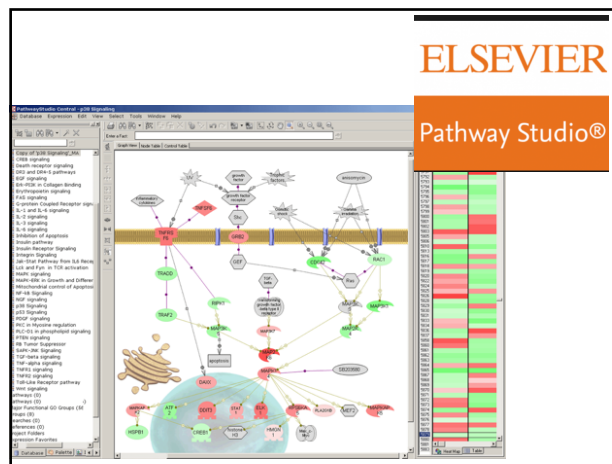
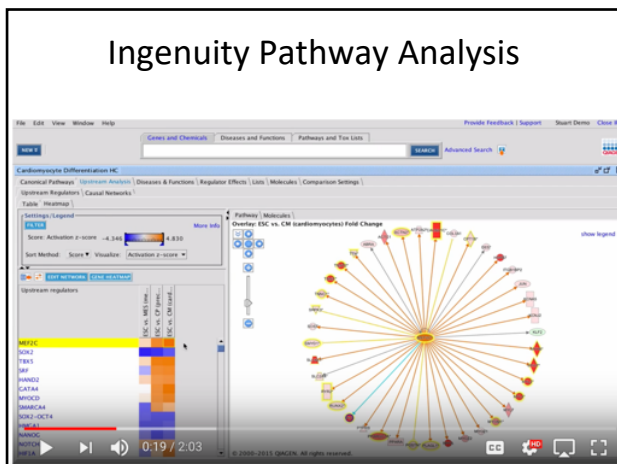
KEGG PATHWAY is a collection of manually drawn pathway maps representing our knowledge on the molecular interaction, reaction and relation networks for:

1. **Metabolism**
Global/overview Carbohydrate Energy Lipid Nucleotide Amino acid Other amino Glycan Cofactor/vitamin Terpenoid/PK Other secondary metabolite Xenobiotics Chemical structure
2. Genetic Information Processing
3. Environmental Information Processing
4. Cellular Processes
5. Organismal Systems
6. Human Diseases
7. Drug Development

KEGG PATHWAY is a reference database for [Pathway Mapping](#).



Ingenuity Pathway Analysis

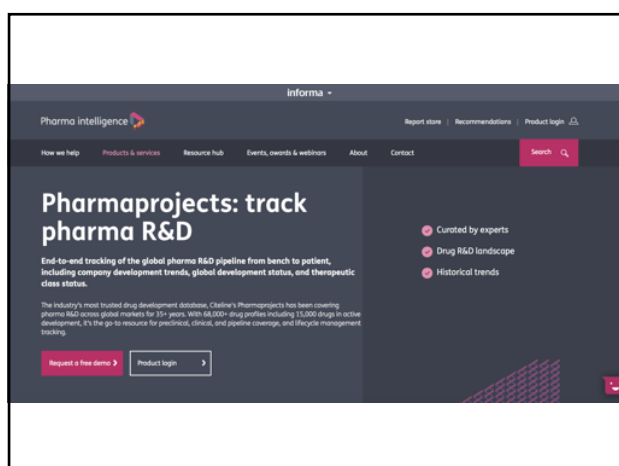
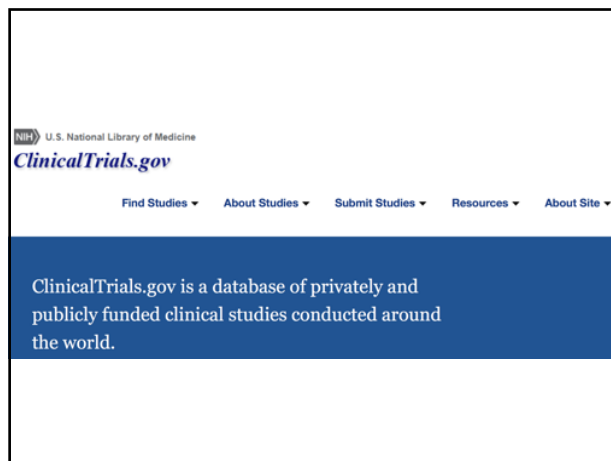


План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- **Базы данных секвенирования**
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- **Базы данных клинических исследований и лекарств**
- Базы данных по экспрессии генов/гистологии
- Базы данных по таксономии



План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- **Базы данных по экспрессии генов/гистологии**
- Базы данные по таксономии

Базы данных экспрессии генов

Базы данных экспрессии генов

THE HUMAN PROTEIN ATLAS

MENU HELP NEWS

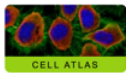
SEARCH

e.g. RBM3, insulin, CD36

Search Fields »



TISSUE ATLAS



CELL ATLAS



PATHOLOGY ATLAS

131

План

- Библиографические/реферативные базы данных литературных источников (статьи, тезисы, патенты, материалы конференций и т.д.)
- Базы данных последовательностей ДНК
- Базы данных последовательностей белков
- Базы данных 3D структур
- Базы данных хим. соединений
- Базы данных геномов и аннотаций
- Базы данных вариаций генома
- Базы данных геном-фенотип
- Базы данных взаимодействий
- Базы данных сигнальных путей
- Базы данных секвенирования
- Базы данных заболеваний и медицинской информации
- Базы данных по экспрессии генов/гистологии
- **Базы данных по таксономии**

Таксономические базы данных

Taxonomy Browser – знаменитая таксономическая БД, имеющая иерархическую структуру, основанную на анализе последовательностей и призванная упорядочить классификацию организмов, для которых известна хотя бы одна последовательность ДНК или белка.



Видовые базы данных

Содержат таксономическую, библиографическую, географическую, визуальную и прочую информацию



Видовые базы данных https://plant.depo.msu.ru

