



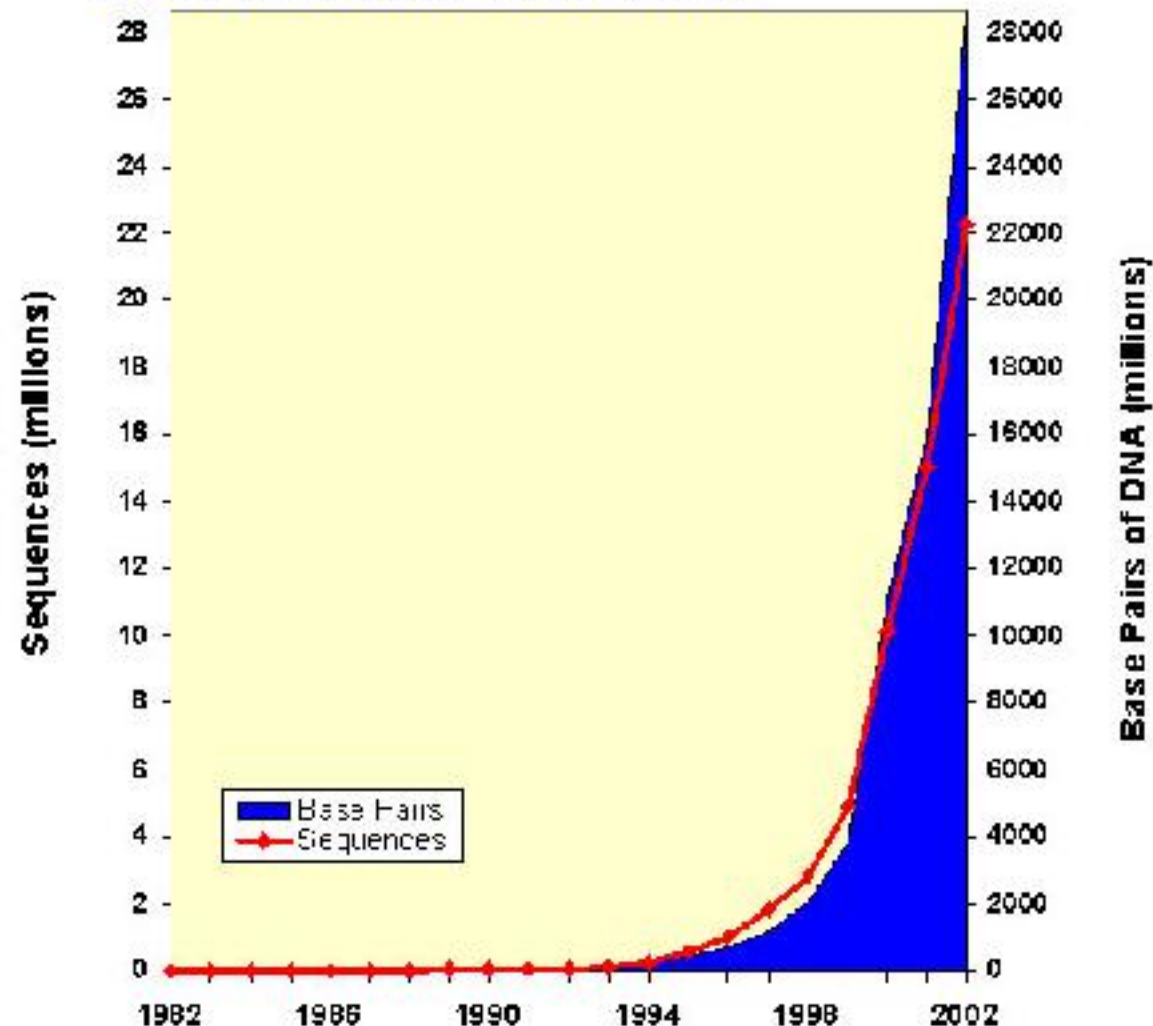
**Операции с последовательностями,  
доступ к базам и форматы данных. Базы  
данных последовательностей ДНК EMBL  
и GenBank**

*Орлов Юрий Львович*

# Международные проекты геномных исследований

Стремительно растут темпы исследований по секвенированию геномной ДНК (Benson et al., 2000; Wheeler et al., 2000). На сентябрь 2003 г. доступны 139 полных геномов прокариот, включая 16 видов археобактерий и 123 вида бактерий.

## Рост объема GenBank



**GenBank за 2002 год:**

**28,507,990,166 п.о.**

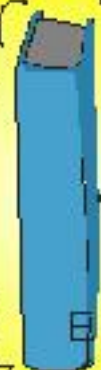
**22,318,883  
последовательностей**

Полностью секвенированы эукариотические геномы мышевидного салата *Arabidopsis thaliana*, червя *Caenorhabditis elegans*, плодовой мушки *Drosophila melanogaster*, дрожжей *Saccharomyces cerevisiae* и *Schizosaccharomyces pombe*, некоторых внутриклеточных паразитических организмов (*Plasmodium falciparum*, *Encephalitozoon cuniculi*).



# Библиотеки геномных последовательностей

**1 том**



E. Coli  
~  $10^6$  п.о.

**10 томов**



Дрожжи  
~  $10^7$  п.о.

**100 томов**



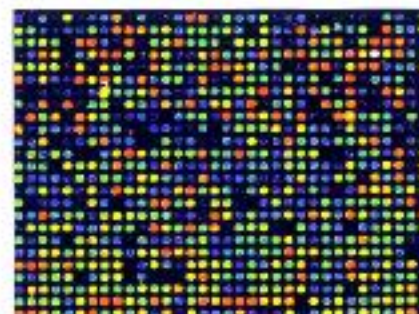
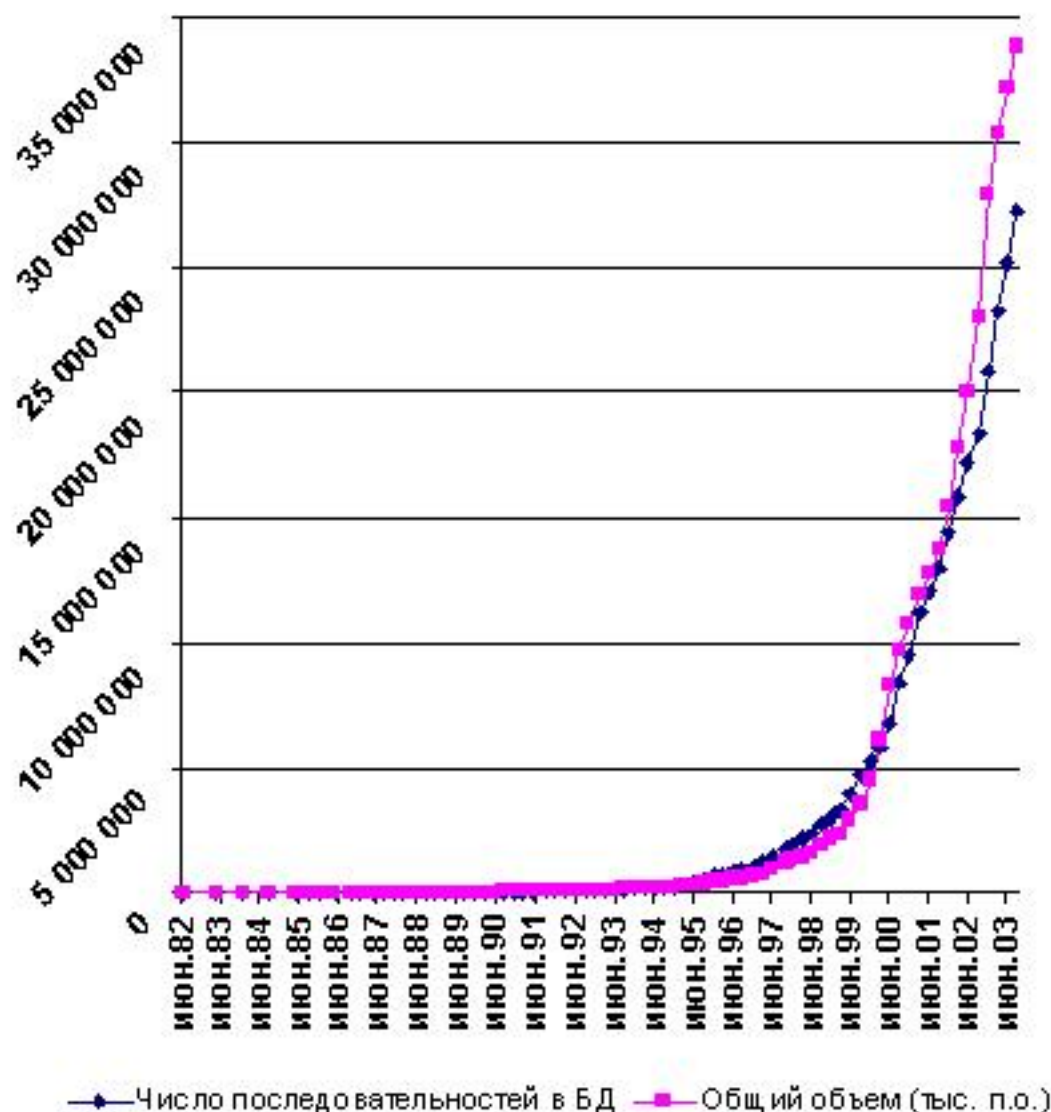
Дрозофила  
~  $10^8$  п.о.

**1000 томов**



Человек  
~  $3 \times 10^9$  п.о.

# Рост числа последовательностей в банке данных EMBL с 1982 по июнь 2003 г.

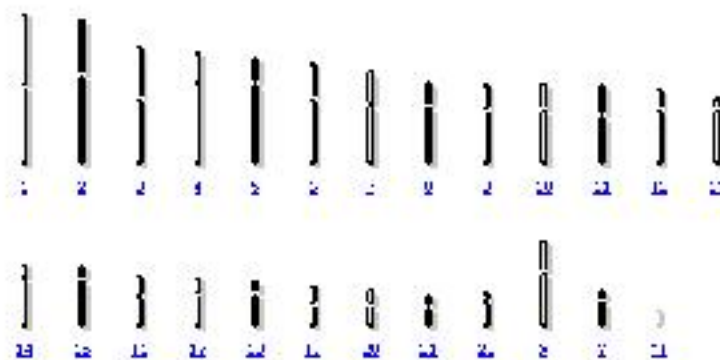


2000 г. – секвенирование  
генома человека

постгеномная эпоха

2003 г. – завершение чернового  
секвенирования

(6 государств, Россия с 1989 г.)



## Дальнейшая структура лекции и курса

Лекция - Базы и банки данных.

Курс - Компьютерная геномика – задачи анализа генетических макромолекул

### Структура лекции

- Интернет. Технические компьютерные средства поддержки БД.
- Классификация баз данных. История развития БД. Репозитории и банки данных.
- Структура карточек GenBank/EMBL
- Поиск информации о нуклеотидных последовательностях. ENTREZ
- Запросы к БД. Перекрестные ссылки (линки).
- Поиск научных статей PubMed

## Техническая основа работы с базами данных. Поиск информации в Интернет.

Пример:

<http://www.bionet.nsc.ru/SRCG/index.html>

**http://**тип документа (http = гипертекстовый документ)

[www.bionet.nsc.ru](http://www.bionet.nsc.ru) адрес сервера в сети Internet

**/SRCG/**каталог сервера в котором находится HTML-документ **index.html** имя файла в котором находится HTML-документ

Сервис в Интернет: почта E-mail; FTP, Archie; Gopher; WWW, JAVA.

Просмотр гипертекстовых файлов (броузеры): программа "Netscape Communicator" "Internet Explorer"

Поиск молекулярно-биологической информации проводится по тому же принципу, что поиск любой другой информации в Интернет с использованием поисковых систем.

# Сервер ИЦиГ СО РАН. Программа Netscape Navigator

The screenshot shows the Netscape Navigator browser window with the address bar set to <http://www.bionet.nsc.ru/>. The website header features the logo of the Institute of Cytology and Genetics (ИЦиГ) and the text "ИНСТИТУТ ЦИТОЛОГИИ И ГЕНЕТИКИ" and "Сибирское Отделение Российской Академии Наук". Navigation links for "По-русски" and "In English" are present. The main content area includes a large blue banner with microscopic images and a "События" (Events) section with news items from September 2003. A sidebar on the left contains navigation links for "Институт", "Научная деятельность", "Кольцо сайтов", and "Информационные ресурсы". A search box is located in the bottom right of the page content.

ИЦиГ СО РАН

ИНСТИТУТ ЦИТОЛОГИИ И ГЕНЕТИКИ  
Сибирское Отделение Российской Академии Наук

По-русски In English

События

Новости

- Архив новостей
- Доска объявлений

Сентябрь 2003

Союз Научной Молодежи Новосибирского Научного Центра (СНМ ННЦ СО РАН) проводит электронный опрос молодежи: "Молодые ученые: поддержка, роль, проблемы и новые возможности". [подробнее >>>](#)

Сентябрь 2003

Региональный Общественный Фонд содействия отечественной науке объявляет ОТКРЫТЫЙ КОНКУРС для ученых Российской академии наук на основании грантов Фонда в области естественных и гуманитарных наук. [подробнее >>>](#)

Вход для сотрудников

Поиск на сервере

«Компьютерная геномика» НГУ, Лекция 1, 2003



# Тот же сайт просматривается с помощью Internet Explorer

Институт цитологии и генетики - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites History

Address <http://www.bionet.nsc.ru/> Go Links

**ИНСТИТУТ ЦИТОЛОГИИ И ГЕНЕТИКИ**  
*Сибирское Отделение Российской Академии Наук*

По-русски In English

**События**

**Новости**

- Архив новостей
- Доска объявлений

**Сентябрь 2003** Совет Научной Молодежи Новосибирского Научного Центра (СНМ ННЦ) СО РАН проводит электронный [опрос молодых ученых](#) "Молодые ученые: лидерская роль, проблемы и новые возможности". [подробнее >>>](#)

**Сентябрь 2003** Региональный Общественный Фонд содействия отечественной науке объявляет [ОТКРЫТЫЙ КОНКУРС для ученых Российской академии наук](#) на соискание грантов Фонда в области естественных и гуманитарных наук. [подробнее >>>](#)

**Вход для сотрудников**

**Поиск на сервере**

Done Internet

Start

19:23

Умение работать с Интернетом – необходимое практическое требование, без которого современная научная работа с информацией, с базами данных просто невозможна.

### Список полезных сайтов:

Учебник по Интернет (Microsoft)	<a href="http://noms.microsoft.com/intl/ru/tutorials/default.htm">http://noms.microsoft.com/intl/ru/tutorials/default.htm</a>
Средства поиска в Интернет (собраны вместе)	<a href="http://cuwww.unige.ch/mats-index.html">http://cuwww.unige.ch/mats-index.html</a>
Список ссылок на поисковые системы	<a href="http://www.nisp.ru/ssarch/">http://www.nisp.ru/ssarch/</a>
Средства поиска (Microsoft)	<a href="http://home.microsoft.com/intl/ru/access/all/ruone.asp">http://home.microsoft.com/intl/ru/access/all/ruone.asp</a>
<b>Поисковые сервера</b>	
Yahoo!	<a href="http://www.yahoo.com/">http://www.yahoo.com/</a>
AltaVista	<a href="http://www.altavista.com/">http://www.altavista.com/</a>
100hot	<a href="http://www.100hot.com">http://www.100hot.com</a>
InfoSeek	<a href="http://www.infoseek.com/">http://www.infoseek.com/</a>
Excite	<a href="http://www.excite.com/">http://www.excite.com/</a>
Lycos	<a href="http://www.lycos.com/">http://www.lycos.com/</a>
WhoWhere	<a href="http://www.whowhere.com/">http://www.whowhere.com/</a>
<b>Российские поисковые системы</b>	
Rambler	<a href="http://www.rambler.ru/">http://www.rambler.ru/</a>
List.ru	<a href="http://www.list.ru/">http://www.list.ru/</a>
Яндекс	<a href="http://www.yandex.ru">http://www.yandex.ru</a>
@Rus	<a href="http://www.eurus.ru/">http://www.eurus.ru/</a>
Апорт	<a href="http://www.cport.ru/">http://www.cport.ru/</a>
Weblist	<a href="http://weblist.ru/">http://weblist.ru/</a>
<b>Сервера бесплатной электронной почты</b>	
Mail.Ru	<a href="http://mai.ru">http://mai.ru</a>

# Базы данных последовательностей ДНК. Исторический обзор.

Первые последовательности, собранные в базах данных были аминокислотными последовательностями.

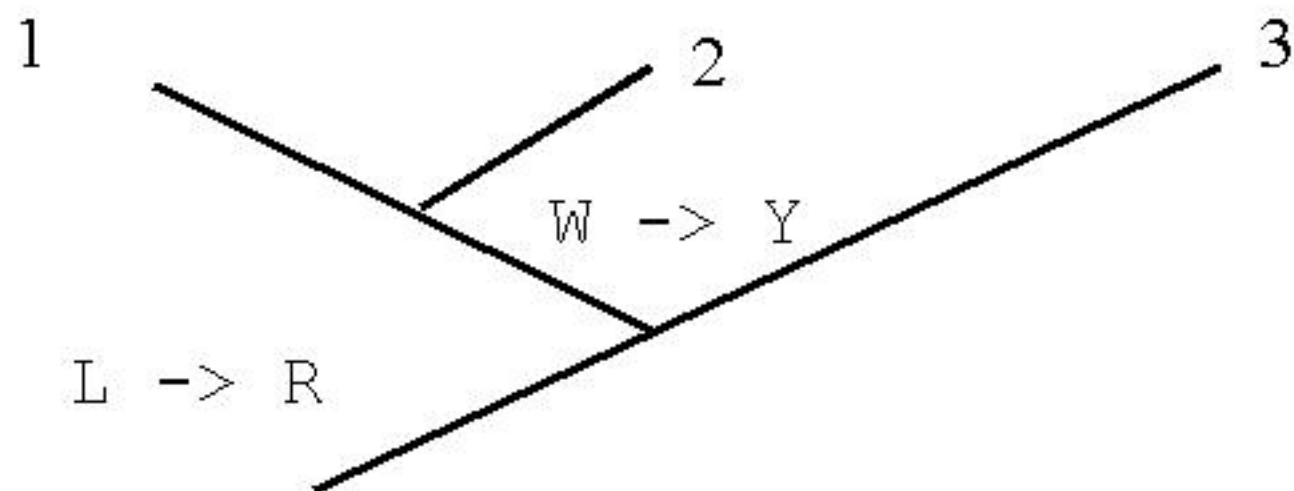
Методы секвенирования белков (Sanger and Tuppy, 1951)  
Margaret Dayhoff (1925-1983)

National Biomedical Research Foundation (NBRF) (Вашингтон, США)  
**1960-е Protein Sequence Atlas**  
коллекция последовательностей, которая перешла в базу данных  
Protein Informational Resource (PIR ранее Protein Identification Resource)  
<http://watson.gmu.edu:8080/pirwww/index.html>

**1980-е PIR-International Protein Sequence Database**  
<http://www-nbrf.georgetown.edu/pir>

Сотрудничество международных центров NBRF, MIPS (Munich Center for Protein Sequences),  
JPIID (Japan International Protein Information Database)  
Данные были организованы в семейства и суперсемейства на основе схожести последовательностей.





Organism\_1 AWTVASAVRTSI  
 Organism\_2 AYTVAAAVRTSI  
 Organism\_3 AWTVAAAVLTSI

### Таблицы частот замен.

Были отобраны белки с не более 15% различий, чтобы наблюдаемые замены аминокислот отражали только одну замену, а не две.

Метод предсказания филогенетических отношений и вероятных аминокислотных замен в процессе эволюции родственных белковых последовательностей.

Матрицы аминокислотных замен PAM – матрицы Дайхоф (Dayhoff) MDM (Mutation Data Matrix), или PAM (Percent Accepted Mutation)

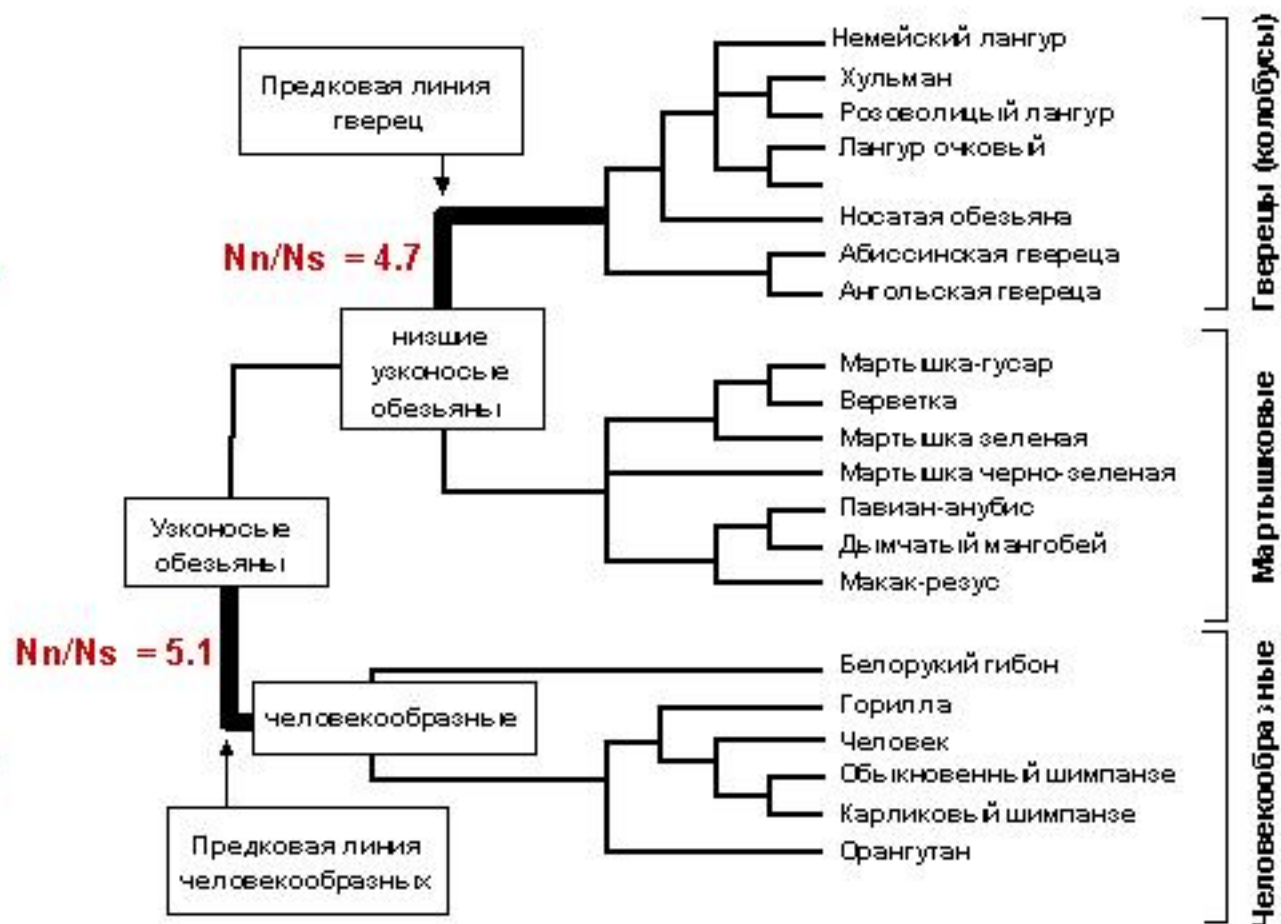
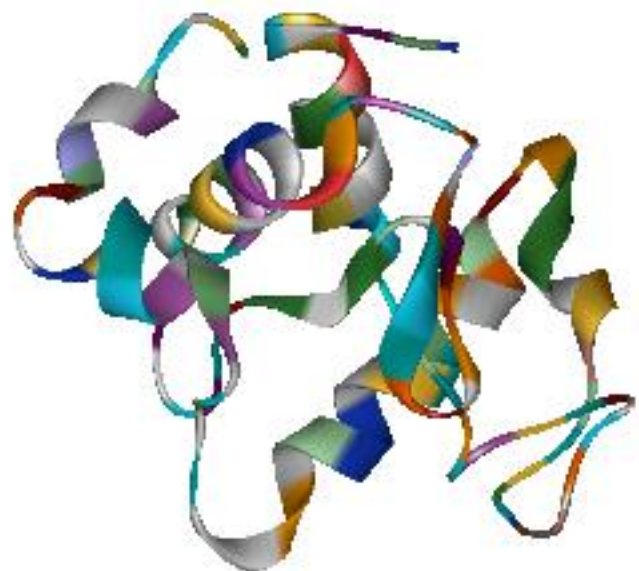
# Научные задачи, решаемые с помощью анализа родственных белков

## Молекулярный филогенетический анализ: адаптивная эволюция генов, кодирующих лизоцимы обезьян

Критерий адаптивной эволюции (Кишура, 1984):

$$Nn/Ns > 1$$

Здесь  $Nn$  и  $Ns$  - количество несинонимических замен, фиксировавшихся в исследуемом эволюционном маршруте



# Базы данных последовательностей ДНК

Walter Goad (1925-2000)

LANL (Los Alamos National Laboratory) Нью-Мексико, США

EMBL (European Molecular Biology Laboratory) Гейдельберг, Германия

1979 - прототип GenBank разрабатывался в LANL 1982-1992,

далее в NCBI (National Center for Biotechnology Information)

[www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)

1980 - EMBL Data Library ([www.ebi.ac.uk](http://www.ebi.ac.uk))

1984 – DDBJ ([www.ddbj.nig.ac.jp](http://www.ddbj.nig.ac.jp))

1983 - ГЕНЭКСПРЕСС и Программа "Геном человека" СССР (А.А.Баев)

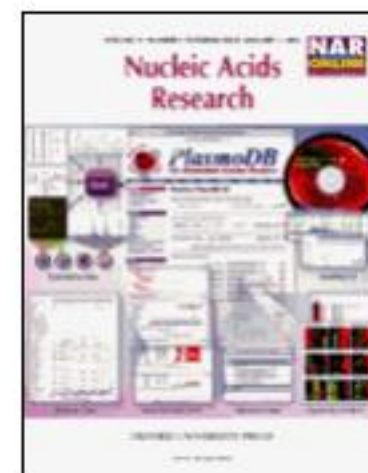
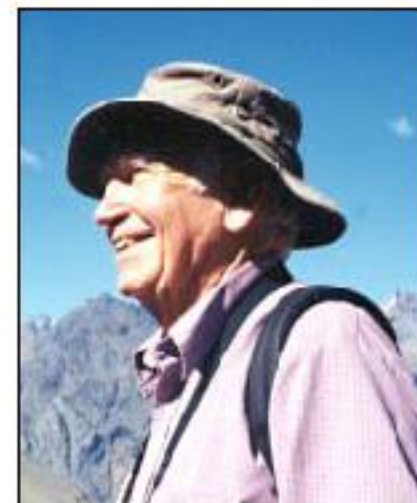
Сейчас эти центры образуют международный консорциум

International Nucleotide Sequence Database Collaboration

[www.ncbi.nlm.nih.gov/collab](http://www.ncbi.nlm.nih.gov/collab)

Описание баз данных (более 300 молбиолбаз данных)

NAR - ежегодный первый выпуск (N1, январь) :



## Nucleic Acids Research

<http://nar.oupjournals.org/content/vol31/issue1/>

OXFORD  
Journals online

# Программы анализа последовательностей

Методы секвенирования ДНК

Maxam and Gilbert, 1977 - Sanger et al., 1977

Секвенирование

Программы анализа последовательностей ДНК, РНК и белков.

1982, 1984 - NAR - первые спец.

выпуски по программам

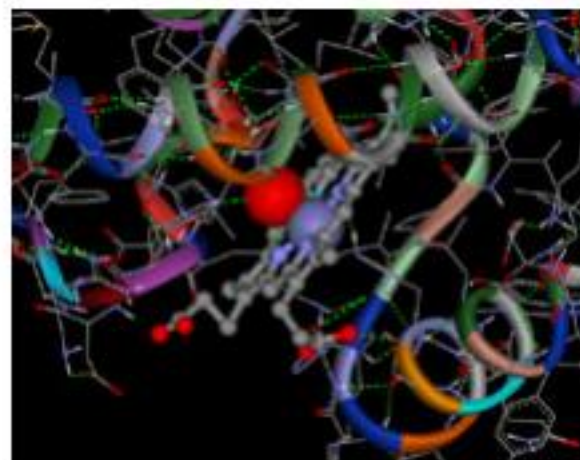
GCG на компьютерах VAX (Ун-т Висконсин)

[www.gcg.com/](http://www.gcg.com/)

Intelligenetics, DNASTar

PHRED, PHRAP ([www.codoncode.com](http://www.codoncode.com))

коммерческие и некоммерческие пакеты



Методы, алгоритмы и современные программные реализации будут подробно рассмотрены в дальнейших разделах курса

# Получение последовательностей из открытых баз данных

Важнейший шаг в обеспечении технического доступа к базам данных – разработка Web-страниц.

GENINFO – D.Benson, D.Lipman

Система ENTREZ (позднее включившая MedLine, NLM)

<http://www.ncbi.nlm.nih.gov/entrez/>



Новый релиз выходит каждые два месяца. Образован международный консорциум (International Nucleotide Sequence Database Collaboration), который состоит из Банка данных ДНК Японии (DNA DataBank of Japan, DDBJ), Европейской Молекулярно-биологической Лаборатории (the European Molecular Biology Laboratory, EMBL), и GenBank США (NCBI). Эти три организации обмениваются данными на ежедневной основе.



# Коллекции молекулярно-биологических баз данных.

Andreas D. Baxevanis. The Molecular Biology Database  
Collection: 2003 update. Nucleic Acids Research, 2003, Vol. 31,  
No. 1 1-12.

## Основные репозитарии:

(<http://nar.oupjournals.org/cgi/content/full/31/1/1>)

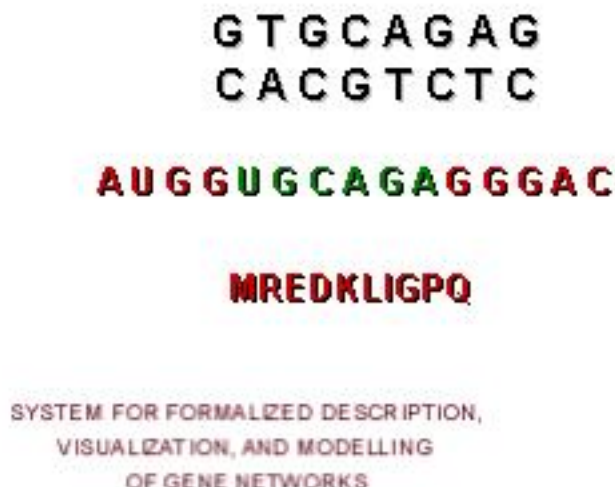
DNA Data Bank of Japan (DDBJ)	<a href="http://www.ddbj.nig.ac.jp">http://www.ddbj.nig.ac.jp</a>	последовательности ДНК и белков, консорциум
EMBL Nucleotide Sequence Database	<a href="http://www.ebi.ac.uk/embl.html">http://www.ebi.ac.uk/embl.html</a>	-//-
GenBank	<a href="http://www.ncbi.nlm.nih.gov/">http://www.ncbi.nlm.nih.gov/</a>	-//-
NCBI Reference Sequence Project	<a href="http://www.ncbi.nlm.nih.gov/RefSeq/">http://www.ncbi.nlm.nih.gov/RefSeq/</a>	Природные биологические молекулы
Ensembl	<a href="http://www.ensembl.org/">http://www.ensembl.org/</a>	Аннотированная информация, Геномы эукариот
UCSC Genome Browser	<a href="http://genome.ucsc.edu/">http://genome.ucsc.edu/</a>	Аннотация геномов
STACK	<a href="http://www.sanbi.ac.za/Dbases.html">http://www.sanbi.ac.za/Dbases.html</a>	Ген-ориентированные кластеры
TIGR Gene Indices	<a href="http://www.tigr.org/tdb/tgi.shtml">http://www.tigr.org/tdb/tgi.shtml</a>	-//-
UniGene	<a href="http://www.ncbi.nlm.nih.gov/UniGene/">http://www.ncbi.nlm.nih.gov/UniGene/</a>	-//-

Ресурсы по молекулярной биологии могут быть классифицированы как **репозитарии, банки данных, базы данных, информационные сайты, электронные библиотеки и Интернет-доступное программное обеспечение.**

Каталог ресурсов: <http://www.mgs.bionet.nsc.ru/mgs/links/links.html>

Специализированные базы данных можно подразделить в соответствии с иерархической организацией хранения и передачи наследственной информации:

- Уровень ДНК
- Уровень РНК
- Уровень белка
- Генные сети



**Эл. каталог ресурсов ИЦиГ СО РАН: <http://www.mgs.bionet.nsc.ru/mgs/links/links.html>**

[\[Links of IC&G SB RAS, Novosibirsk\]](#) [\[Priority links\]](#) [\[WWW-Biosciences\]](#) [\[GENOME\]](#) [\[NCBI National Center of Biotechnological Informatic\]](#) [\[Data bases\]](#) [\[Bio Tools\]](#) [\[EBI European Bioinformatics Institute\]](#) [\[BCM Baylor College of Medicine\]](#) [\[GDB Genome Database\]](#) [\[PDB Protein Data Bank Brookhaven National Lab.\]](#) [\[Patents\]](#) [\[PROGR\]](#) [\[France\]](#) [\[SilverPlatter\]](#)

## **Links of IC&G SB RAS, Novosibirsk:**

[INSTITUTE OF CYTOLOGY AND GENETICS, The Siberian Division of the Russian Academy of Sciences](#)

[The First International conference on Bioinformatics of Genome Regulation and Structure](#)

[Vogis Courier \(Russian only\)](#)

[Lab. of Experimental Modelling of Evolutionary Processes](#)

[Meiosis Laboratory Home Page](#)

[Sector of Molecular Evolution](#)

[X inactivation](#)

[Methods of Genetic Analysis Lab](#)

[The Laboratory of Animal Molecular Genetics](#)

## **Priority links:**

[Novosibirsk Institute of Bioorganic Chemistry](#)

[Institute of Protein Research, RAS](#)

[Engelhardt Institute of Molecular Biology \(EIMB\)](#)

[A National Laboratory for Computational Science and Engineering \(SDSC\)](#)

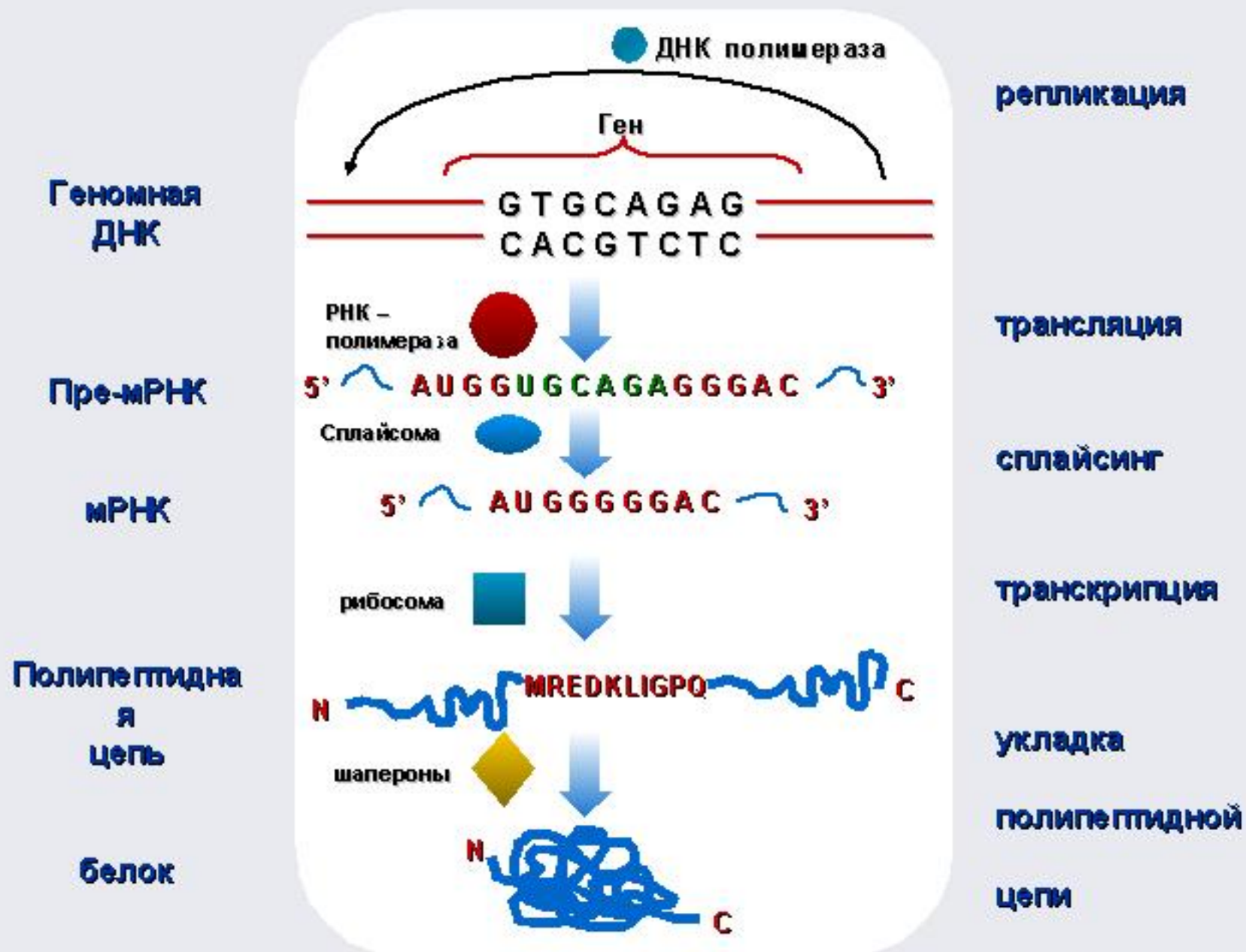
[Molecular Bioinformatics of Gene Regulation, \(Transfac DATA Base\)](#)

[Institute of Advanced Biomedical Technologies \(ITBA\)](#)

Полный каталог содержит  
более 1000 эл.ресурсов по  
молекулярной биологии

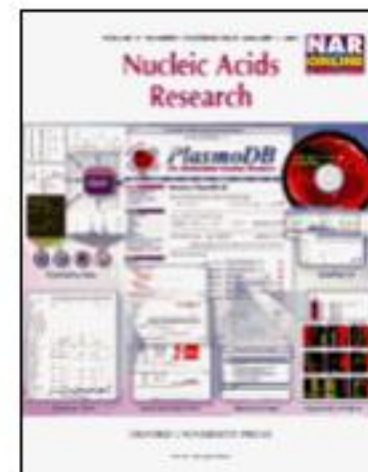
# Базы данных и электронные ресурсы по биоинформатике

## Фундаментальные генетические процессы



## ОСНОВНЫЕ НАПРАВЛЕНИЯ (ГРУППЫ) БАЗ ДАННЫХ ПО МОЛЕКУЛЯРНОЙ БИОЛОГИИ (Более 300 баз данных)

- Сравнительная геномика (Comparative Genomics)
- Экспрессия генов с помощью микрочипов
- Структура и регуляция генов
- Генетическое и физическое картирование
- Геномные базы данных.
- Межмолекулярные взаимодействия
- Метаболические пути и клеточная регуляция
- Мутации
- Белки
- РНК
- Структуры (Пространственные)
- Трансгенные растения и животные
- Другие направления (Литературные базы данных, таксономия)



Andreas D. Baxevanis. The Molecular Biology Database Collection: 2003 update. Nucleic Acids Research, 2003, Vol.31, No. 1 1-12.

<http://nar.oupjournals.org/cgi/content/full/31/1/1>

# ОСНОВНЫЕ НАПРАВЛЕНИЯ БАЗ ДАННЫХ ПО МОЛЕКУЛЯРНОЙ БИОЛОГИИ (NAR, 2003)

## Сравнительная геномика (Comparative Genomics)

Clusters of Orthologous Groups (COG) <http://www.ncbi.nlm.nih.gov/COG>

содержит филогенетическую классификацию белков из 43 полных геномов



## Экспрессия генов с помощью микрочипов (Gene Expression) microarray

## Структура и регуляция генов (Gene Identification and Structure)

SNP Consortium database <http://snp.cshl.org> SNP Consortium data

TRRD <http://www.bionet.nsc.ru/trrd/> Transcription regulatory regions of eukaryotic genes

TRANSCRIPTION REGULATORY  
REGIONS DATABASE



## Генетическое и физическое картирование (Genetic and Physical Maps)

## Геномные базы данных (Genomic Databases)

ACeDB information <http://www.acedb.org/> *C. elegans*, *S. pombe*, and human sequences and genomic information

GeneCards <http://bioinfo.weizmann.ac.il/cards/> Integrated database of human genes, maps, proteins and diseases

GOLD <http://igweb.integratedgenomics.com/GOLD/> Information regarding complete and ongoing genome projects



## Межмолекулярные взаимодействия (Intermolecular Interactions)

BIND <http://bind.ca> Molecular interactions, complexes and pathways

## ГЕННЫЕ СЕТИ. Метаболические пути и клеточная регуляция (Metabolic Pathways and Cellular Regulation)

EcoCyc <http://ecocyc.org/> *Escherichia coli* K-12 genome, metabolic pathways, transporters and gene regulation

EpoDB <http://www.cbil.upenn.edu/EpoDB/> Genes expressed during human erythropoiesis

Kyoto Encyclopedia of Genes and Genomes (KEGG) <http://www.genome.ad.jp/kegg> Metabolic and regulatory pathways

RegulonDB [http://www.cifn.unam.mx/Computational\\_Genomics/regulondb/](http://www.cifn.unam.mx/Computational_Genomics/regulondb/) *Escherichia coli* transcriptional regulation and operon organization



## Мутации

### Mutation Databases

Human Gene Mutation Database (HGMD) <http://www.hgmd.org> Known (published) gene lesions underlying human inherited disease

«Компьютерная геномика»

НГУ, Лекция 1, 2003

# ОСНОВНЫЕ НАПРАВЛЕНИЯ БАЗ ДАННЫХ ПО МОЛЕКУЛЯРНОЙ БИОЛОГИИ (NAR, 2003)

## БЕЛКИ (Protein Databases)

**Kabat Database** <http://immuno.bme.nyu.edu/> Sequences of proteins of immunological interest

**PIR-NREF** <http://pir.georgetown.edu/pirwww/pirref.shtml> Non-redundant reference database with comprehensive protein sequences

**SWISS-PROT/TrEMBL** <http://www.expasy.ch/sprot> Curated protein sequences

**TRANSFAC** <http://transfac.gbf.de/TRANSFAC/index.html> Transcription factors and binding sites

## Мотивы в белках (Protein Sequence Motifs)

**CluSTr** <http://www.ebi.ac.uk/clustr/> Automatic classification of SWISS-PROT+TrEMBL proteins

**Pfam** <http://www.sanger.ac.uk/Software/Pfam/> Multiple sequence alignments and hidden Markov models of common protein domains

**PROSITE** <http://www.expasy.org/prosite> Biologically significant protein patterns and profiles

## Протеомные ресурсы (Proteome Resources)

**Proteome Analysis Database** <http://www.ebi.ac.uk/proteome/> Online application of InterPro and CluSTr for the functional classification of proteins in whole genomes

## Системы поиска информации (Retrieval Systems and Database Structure)

**TESS** <http://www.cbil.upenn.edu/teess> Transcription element search system

## РНК (RNA Sequences)

**ACTIVITY** <http://util.bionet.nsc.ru/databases/activity.html> Functional DNA/RNA site activity

**RNA Modification Database** <http://medlib.med.utah.edu/RNAmods/> Naturally modified nucleosides in RNA

**UTRdb/UTRsite** <http://bighost.area.ba.cnr.it/sr5/> 5'- and 3'-UTRs of eukaryotic mRNAs and relevant functional patterns

## Пространственные структуры (Structure)

**HSSP** <http://www.sander.ebi.ac.uk/hssp/> Structural families and alignments; structurally-conserved regions and domain architecture

**PDB** <http://www.pdb.org/> Structure data determined by X-ray crystallography and NMR

**SCOP** <http://scop.mrc-lmb.cam.ac.uk/scop> Familial and structural protein relationships

## Трансгенные растения и животные (Transgenics)

Другие направления (Varied Biomedical Content)

**PubMed** <http://www.ncbi.nlm.nih.gov/PubMed/> MEDLINE and Pre-MEDLINE citations

**Tree of Life** <http://phylogeny.arizona.edu/tree/phylogeny.html> Information on phylogeny and biodiversity

## Поиск литературной информации.

### База MEDLINE

Определение базы данных. Поле, вход, запись.

Адрес базы данных MEDLINE : <http://www.ncbi.nlm.nih.gov/PubMed/>

### Простой поиск.

Формат запроса: word1 (AND word2 OR word3 ....)

### Сложный поиск.

Формат запроса: word1[FieldCode] AND word2[FieldCode]

Общий формат операторов запроса.

*Логические операторы*

| - "ИЛИ" (OR).  
& - "И" (AND).  
! - "И НЕ" (AND NOT, BUT NOT).

Пример действия логических операторов:

The image shows three Venn diagrams illustrating logical operators. Each diagram consists of two overlapping circles labeled 'A' and 'B'.  
1. OR ('|'): Both circles A and B are filled with red. Below it is the text 'OR ("|")'.  
2. AND ('&'): Only the overlapping area between circles A and B is filled with red. Below it is the text 'AND ("&")'.  
3. BUT NOT ('!'): Only circle A is filled with red, while circle B is empty. Below it is the text 'BUT NOT ("!")'.

Сохранение  
найденной  
информации

(принцип  
Cut&Paste)



Ключевые слова структурированы в специальном словаре для поиска литературной (медицинской и др.) информации



MEDICAL SUBJECT HEADINGS



Contact NLM | Site Index | Search Our Web Site | NLM Home

Health Information | Library Services | Research Programs | New & Noteworthy | General Information

Каталог доступен на сервере NCBI  
<http://www.nlm.nih.gov/mesh/>



## MeSH Browser

- [Online searching](#) of MeSH vocabulary
- [About the MeSH Browser](#)



## All About MeSH

- [MeSH Fact Sheet](#)
- [Presentations and papers](#) by MeSH staff.
- [Online introductory material](#) to the *Annotated MeSH*
- [Information from Previous Years](#)
- [Suggestions](#) for authors' keywords

## Obtaining MeSH

- [Download](#) electronic copies.
- Ordering information for [Printed](#) versions.



## What's New



## MeSH Staff

- [Biographies and email](#).
- [Publications and presentations](#).

## Related Efforts

- [Congenital Abnormalities associated with Mental Retardation \(MCA/MR\)](#).
- [Unified Medical Language System \(UMLS®\)](#)
- [NLM Classification](#)



## MeSH Suggestions

- [Send](#) MeSH vocabulary suggestions.

## MEDLINE/PubMed

Ссылки и резюме статей из 4600 биомедицинских журналов

[What's New](#) | [MeSH Fact Sheet](#) | [MeSH Browser](#) | [MeSH Files](#) | [MeSH Publications](#) | [MeSH Staff](#) | [MeSH Suggestions](#)

# Понятие карточки банка данных. Пример EMBL

ID - identification	(begins each entry; 1 per entry)
AC - accession number	(>=1 per entry)
SV - sequence version	(1 per entry)
DT - date	(2 per entry)
DE - description	(>=1 per entry)
KW - keyword	(>=1 per entry)
OS - organism species	(>=1 per entry)
OC - organism classification	(>=1 per entry)
OG - organelle	(0 or 1 per entry)
RN - reference number	(>=1 per entry)
RC - reference comment	(>=0 per entry)
RP - reference positions	(>=1 per entry)
RX - reference cross-reference	(>=0 per entry)
RG - reference group	(>=0 per entry)
RA - reference author(s)	(>=0 per entry)
RT - reference title	(>=0 per entry)
RL - reference location	(>=1 per entry)
DR - database cross-reference	(>=0 per entry)
CC - comments or notes	(>=0 per entry)
AH - assembly header	(0 or 1 per entry)
AS - assembly information	(0 or >=1 per entry)
FH - feature table header	(0 or 2 per entry)
FT - feature table data	(>=0 per entry)
XX - spacer line	(many per entry)
SQ - sequence header	(1 per entry)
CO - contig/construct line	(0 or >=1 per entry)
bb - (blanks) sequence data	(>=1 per entry)
// - termination line	(ends each entry; 1 per entry)

Текстовый  
файл, длина  
строки  
ограничена,  
Ключевые слова  
жестко заданы

[http://www.ebi.ac.uk/embl/Documentation/User\\_manual/usrman.html](http://www.ebi.ac.uk/embl/Documentation/User_manual/usrman.html)

# Пример карточки EMBL, содержащий ген протеин-киназы дрозофилы

```
ID  DMCAGKCA  standard; DNA; INV; 3955 BP.
XX
AC  M18555; J03504;
XX
SV  M18555.1
XX
DT  16 JUL 1988 (Rel. 16, Created)
FF  06-JUL-1988 (Rel. 20, Last updated, Version 1)
XX
DE  D.melanogaster cAMP dependent protein kinase catalytic subunit,
FE  complete cos.
XX
FU  cAMP dependent protein kinase; kinase; protein kinase.
XX
OS  Drosophila melanogaster (fruit fly)
OC  Eukaryota; Metazoa; Arthropoda; Tracheata; Hexapoda; Insecta; Pterygota;
OC  Neoptera; Endopterygota; Diptera; Brachycera; Muscivora; Ephydroidea;
OC  Drosophilidae; Drosophila.
XX
RN  [1]
RP  1-3959
RX  MEDLINE; 88115281.
RI  Foster J.L., Higgins G.C., Jackson M.P.:
RT  "Cloning, sequence, and expression of the Drosophila cAMP-dependent protein
RT  kinase catalytic subunit gene":
RI  J. Biol. Chem. 263:1676-1681(1988).
XX
DR  FLYBASE; FEgnC00C273; Ika C1.
FE  SWISS-PROT; P12370; KAPC_FROFL.
XX
CC  Draft entry and computer readable copy of sequence [1] kindly
CC  provided by J.L.Foster (17-MAR-1988).
XX
FH  Key          Location/Qualifiers
FH
FT  source       1..3959
FT              /db_xref="taxon:7227"
```

Идентификатор

Код доступа

Описание  
последовательности

Литературная ссылка

Разметка



Кроме форматов EMBL/GenBank, предназначенных для хранения информации существуют форматы для компьютерной обработки и анализа последовательностей :

NBRF-PIR

GCG

Plain/ASCII Staden

Genetic Data Environment (GDE)

Fasta/Pearson format

Intelligenetics

PIR/CODATA

ASN.1 Abstract Syntax Notation

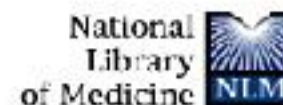
NEXUS

Fasta

```
>seq1  
aggctgct agct agct  
>seq2  
aactaact
```

NBRF

```
>DL;seq1  
seq1., 16 bases, 268 checksum  
aggctgctag ctagct*  
>DL;seq2  
aactaact
```



PubMed

Nucleotide

Protein

Genome

Structure

PMC

Taxonomy

Search: PubMed

for

Go

Clear

Limits

Preview/index

History

Clipboard

Details

About Entrez

Feed Versions

Entrez PubMed

Overview

Help | FAQ

Tutorials

New/Noteworthy

E-Utilities

PubMed Services

Journals Database

MeSH Database

Single Citation

Matcher

Batch Citation Matcher

Clinical Queries

LinkOut

Caddy

Related Resources

Order Documents

NLM Gateway

TOXNET

Consumer Health

Clinical Alerts

Clinical Trials.gov

PubMed Central

- Enter one or more search terms, or click [Preview/index](#) for advanced searching
- Enter [author names](#) as smith j. Initials are optional
- Enter [journal titles](#) in full or as MEDLINE abbreviations. Use the [Journals Database](#) to find j

PubMed is a service of the National Library of Medicine that provides access to over 12 million MEDLINE and additional life science journal articles, including full-text articles.

Entrez – система навигации по базе данных GenBank и электронной библиотеке PubMed

### Bookshelf Additions



Now available: new and updated material in *The NCBI Handbook* and *Genes and Disease*.

### New PubMed Features



The Summary page displays a new icon link for free full-text articles.

New data and additional search options, including an [e-mail](#) selection, have been added to PubMed. See [New/Noteworthy](#).

### Severe Acute Respiratory Syndrome

Citation matches about [Severe Acute Respiratory Syndrome \(SARS\)](#) are provided during this time of peak interest to facilitate searching this topic.

# Entrez – электронная библиотека PubMed

The screenshot shows the PubMed website interface in Microsoft Internet Explorer. The browser's address bar displays the URL: <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?CMD=Search&DB=PubMed>. The page features the NCBI logo on the left and the National Library of Medicine (NLM) logo on the right. The main navigation bar includes links for PubMed, Nucleotide, Protein, Gene, Genome, Literature, PDB, Anatomy, and PIV. The search bar contains the text "Tuberc" and "for Monoc". Below the search bar, there are options for "Units", "Preview/Index", "History", "Clipboard", and "Details". The results display section shows "Disp by Summary", "Show 20", "Print", and "Send to Text". The results list includes several articles, with the first one being "Characterization and analysis of posttranslational modifications of the human large cytoplasmic ribosomal subunit proteins by mass spectrometry and Edman sequencing." A red arrow points to the search bar.

# Просмотр найденной записи – резюме статьи

The screenshot shows the PubMed website interface. At the top, there are logos for NCBI, PubMed, and the National Library of Medicine (NLM). Below the logos is a navigation bar with tabs for 'PubMed', 'Nucleotide', 'Protein', 'Structure', 'FYD', 'Taxonomy', 'Cyt', and 'Euk'. A search bar is present with the text 'Search PubMed for' and buttons for 'Go' and 'Clear'. Below the search bar are links for 'Limits', 'Previous/Next', 'History', 'Clipboard', and 'Details'. A secondary navigation bar includes 'Inquiry', 'Abstract', 'Show: 20', 'Sort', 'View: 1', and 'Text'. The main content area displays a search result for a protein chemistry article. The article title is 'Characterization and analysis of posttranslational modifications of the human large cytoplasmic ribosomal subunit proteins by mass spectrometry and Edman sequencing.' The authors listed are Odintsova TI, Muller EC, Ivanov AV, Igorov IA, Bismert R, Vladimirov SF, Kostka S, Otto A, Wittmann-Liebold B, and Karpova GG. The affiliation is the Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russian Federation. The abstract text describes the isolation and analysis of 60S ribosomal proteins from human placenta, identifying 22 proteins with various post-translational modifications. The PMID is 12582925, and it is noted as 'PubMed - in process'.

PubMed  
Nucleotide  
Protein  
Structure  
FYD  
Taxonomy  
Cyt  
Euk

Search PubMed for [Go] [Clear]

Limits Previous/Next History Clipboard Details

Inquiry Abstract Show: 20 Sort View: 1 Text

1 | J Protein Chem. 2002 Apr;21(3):245-58. [Full Text](#) [Link](#)

**Characterization and analysis of posttranslational modifications of the human large cytoplasmic ribosomal subunit proteins by mass spectrometry and Edman sequencing.**

Odintsova TI, Muller EC, Ivanov AV, Igorov IA, Bismert R, Vladimirov SF, Kostka S, Otto A, Wittmann-Liebold B, Karpova GG.

Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russian Federation.

The 60S ribosomal proteins were isolated from ribosomes of human placenta and separated by reversed phase HPLC. The fractions obtained were subjected to trypsin and Glu-C digestion and analyzed by mass fingerprinting (MALDI-TOF), MS/MS (ESI), and Edman sequencing. Forty-nine basic subunit proteins were found, 22 of which allowed mass fingerprinting in accordance with the SwissProt database (June 2002) masses (proteins L8, L7, L9, L13, L15, L17, L18, L21, L22, L24, L26, L27, L30, L32, L34, L35, L36, L37, L37A, L38, L39, L41). Eleven (proteins L7, L10A, L11, L12, L13A, L23, L23A, L27A, L30, L29, and P0) resulted in mass changes that are consistent with N-terminal loss of methionine, acetylation, formation of all-glycine, or hydroxylation. A loss of methionine without acetylation was found for protein L8 and L17. For many proteins (L3, L4, L5, L7A, L10, L14, L15, L31, and L40), the molecular masses could not be determined. Proteins P1 and protein L3-L2c were not identified by the methods applied.

PMID: 12582925 [PubMed - in process]

Пример – статья найденная по ключевому слову Ivanov (автор)



# Поиск по ключевым словам в базе данных GenBank.

## Запрос «альбумин».

The screenshot shows a Microsoft Internet Explorer browser window displaying the NCBI Nucleotide search results for the query 'albumin'. The search was performed on the URL <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?CMD=Search&DB=nucleotide>. The search results are displayed in a table with columns for accession number, description, and links. The results are as follows:

Accession Number	Description	Links
U197042	Rattus norvegicus albumin (Alb), mRNA g U197042 c U197042.1 U197042.1	Links
U197873	Rattus norvegicus group specific component (Gc), mRNA g U197873 c U197873.1 U197873.1	Links
U197895	Home regions cDNA FL167013, clone U197895, highly similar to Syrian hamster precursor g U197895 c U197895.1 U197895.1	Links
U197512	Mus musculus alpha-fetoprotein (Afp), mRNA g U197512 c U197512.1 U197512.1	Links
U197517	Mus musculus plasminogen activator, urokinase (Uro), mRNA g U197517 c U197517.1 U197517.1	Links
U197518		Links

The browser interface includes a search bar with the query 'albumin', a 'Go' button, and a 'Search' button. The search results are displayed in a table with columns for 'Accession', 'Description', and 'Links'. The search results are sorted by 'Relevance' and show 1 to 20 of 6511 results. The current page is 1 of 328.

# Пример найденной карточки GenBank



Pubmed Nucleotide Protein Genome Structure PMC Таксоны

Search Nucleotide for [ ] Go Clear

Limits Preview/Index History Clipboard

Display default Show 20 Send File Get Subsequence Features

1: [NM\\_134326](#) *Rattus norvegicus* [g:19705430]

LOCUS NM\_134326 1956 bp mRNA linear ROD. IC-SFP-2003

DEFINITION *Rattus norvegicus* albumin (Alb), mRNA.

ACCESSION NM\_134326

VERSION NM\_134326.1 GI: 9705430

KEYWORDS .

SOURCE *Rattus norvegicus* (Norway rat)

ORGANISM [Rattus norvegicus](#)  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae;  
*Rattus*.

REFERENCE 1 (bases 1 to 1956)

AUTHORS Kuot, S., Papet, I., Béchereau, J., Denis, P., Buffiere, C., Gimonet, J.,  
Gimont, P., Elyousfi, M., Fieulle, E. and Oled, C.

TITLE Increased albumin plasma efflux contributes to hypoalbuminemia only  
during early phase of sepsis in rats

JOURNAL Am. J. Physiol. Regul. Integr. Comp. Physiol. 284 (3): R707-R713  
(2003)

MEDLINE [22458063](#)

PUBMED [12571074](#)

REMARK GeneRIF: an accelerated plasma efflux of albumin contributes to  
hypoalbuminemia only during the early period of sepsis

REFERENCE 2 (bases 1 to 1956)

AUTHORS Forlak, J., Dangero, H. and Thum, T.

TITLE ArcoRox 1254 modulates gene expression of nuclear transcription  
factors: implications for albumin gene transcription and protein  
synthesis in rat hepatocyte cultured

Описание  
полей в  
GenBank:  
  
Аналогичная  
структура  
  
Более длинные  
ключевые  
слова (11)

```

      /rntse="reading frame (prn pep13F) "
BASE COUNT      545 a      492 c      470 g      441 t
ORIGIN
      atgaagtggg taacatttctt cccccccccn nncatccccg gttctgcertt ttctagggggt
  61 ctgtttccgcc caqaacqaca caaqaqtqaa atcccacata qgttcaagca ctcaqcaqaa
 121 cagcacttes asggcotagt cctgattgpc btttcczagt acctccagaa atgcccstat
 181 caagagrata tcaaat.t.gtr. gtaggaagna acaaacctt.g caaaaaaarg tgtgctgar
 241 caqaaccccc aaactqtqa caaqtccabt cacactctct tccqacacca ctcatqccc
 301 atcccasagc ctgctgaaaa ctacggtgaa ctggtctgact gctgtgcasa caasgagccc
 361 caaagaaagc agt.gtt.t.ctt. gtagcaaaag gat.gaa.aat. caaacct.gtr. aacct.tccag
 421 aqcccqcaqg ctqacqccat qtqacactcc ttccaaqaaq accctaccag ctctccqqa
 481 caetacttgc atgasgttgc caggazacat ccttatttct acccccaga acctccctac
 541 tatgctgaga aat.aaanga ggtt.ctgarn taggtctgca cacagcttga caaagcagtr
 601 tqcctqccac ccaaccttqa tqccqtqaaa caqaaazcaz tqctcccacc tqccctcac
 661 aggatgaagt gctccagbat gcagazatbt gtagazagag ccttccaaag ctgggcagta
 721 gttgctatga gctagagat. cctcaat.gtr. gagtt.ctgag aaat.caccaa at.tggcaaa
 781 caactccccc aaaccacaaa qqaqtctctt cacqccazcc tqttcaatq ccccqatqac
 841 agggcagAAC ctgccaagtA catgtgtgag aaccagzcca cctctctccag caaacctccag
 901 ccttgetgtg ataagccagt gctccagaaa ccczagctc tcccagagc agaacatgac
 961 caectccctg ccgatctgcc ctcaabagst gctgactttg tccaggaaca ggcagctgtg
1021 sagsactatg ctgaggecaa ggatgtcttc ctgggcactt ttttgtatga atattccage
1081 aggcaccccc attactccyl yltccctycty ctgaaatly ctcaagaaala tgaagcaca
1141 ctggagcagt gctgtgetga agggctctct cctgctgett accgcccagt gctgcccagc
1201 ttccagccc ctgtcgaaga acctaaagaa ttggtcaaaa ctaactgga ccttaccgag
1261 aagctggag agtatggat tccaaacccy ctctctyctt galacaccca caaagcactt
1321 caggtgtcgc cccccaectt cgtggaggca gcaagaaacz tgggacgagt gggcaaccag
1381 tgtgtccccc tteetgaagc tcagagactg cctgtgtggy aagactactt gctctgcacc
1441 ctgaacccy ctgtgtctct ycatgaaay acctcaatga gcaagaaat taccaaatgt
1501 tgcagctgggt cctcgttggc aagaccggca tgtttctctg cctgcccagt tgcagagacc
1561 taagtccccc aagagttaa agtgtgagaz ttcaactctt actctgactt ctgcaccctc
1621 ccagaccagc agaaacagat aaayaaayca aagctctcty ctgaatlyct caaacacaa
1681 ccccagccc ccagcagaca gctgaaagz gtgctgggtg acttegcaca ctccgtggac
1741 aagtgttgc aggtgctga caaggat aac tgettgcaca ctgaggggccc caaccttgtt
1801 ctcaagaaac aagaaacctt agctcaaaa calcaaaccc auctcaagctt accctgagaa
1861 caacagccat gcagcctzag gactctctct ttctgttggg gcaaacccca caccctcagg
1921 aaccacact cctctgaaca tttgacttct tttctc

```

Продолжение  
карточки GenBank.

Описание  
последовательности с  
11-й позиции,  
нумерация,

Ключевое слово  
ORIGIN



PubMed Nucleotide Protein Genome Structure PMC  
Search: PubMed for [ ] Go Clear  
Limits Preview/Index History Clipboard Details

- About Entrez
- Entrez Overview
- Help/FAQ
- New/Noteworthy
- E-Utilities

- PubMed
- Protein
- Nucleotide
- Structure
- Genome
- Books
- 3D Domains
- Domains
- GEC
- GEC Transcripts
- Journals

- Enter one or more search terms, or click [Preview/Index](#) for advanced searching.
- Enter [author names](#) as sm.ty. Initials are optional.
- Enter [journal titles](#) in full or as MEDLINE abbreviations. Use the [Journals Database](#) to find journal titles.

PubMed, a service of the National Library of Medicine, provides access to over 12 million MEDLINE citations back to the mid-1960s and additional life science journals. PubMed includes links to many sites providing full text articles and other related resources.

### Bookshelf Additions



Now available: new and updated material in [The NCBI Handbook](#) and [Genes and Disease](#).

### New PubMed Features

The Summary page displays a new icon link for free full-text articles.

New data and additional search options, including an [EMAIL](#) selection, have been added to PubMed. See [New/Noteworthy](#).

### Severe Acute Respiratory Syndrome

Citations to articles about [Severe Acute Respiratory Syndrome \(SARS\)](#) are provided during this time of peak interest to facilitate searching this topic.

- PubMed Services
- Journals Database
- MESH Database
- Single Citation Matcher
- Batch Citation Matcher
- Clinical Queries
- LinkOut
- Cut by
- Related Resources
- Order Documents
- NLM Gateway
- OXN=1
- Consumer Health
- Clinical Alerts
- Clinical Trials.gov
- PubMed Central



Entrez Genomes - Microsoft Internet Explorer

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Genomes

Summary

Accession	Organism	Genome	Genome	Genome
1: G1_001141	Saccharomyces cerevisiae chromosome III, complete chromosome sequence [27]	g[17910050] (G1_001141.1)	Link	Link
2: G1_001307	Yeast genome assembly (complete genome) [1333]	g[17910052] (G1_001307.1)	Link	Link
3: G1_004976	Leishmania major chromosome I, complete sequence [7352]	g[15310099] (G1_004976.1)	Link	Link
4: G1_001227	AF1494 genome, complete genome [6744]	g[2132113] (G1_001227.1)	Link	Link
5: G1_004295	Comma erythrox virus V2, complete genome [6710]	g[1544915] (G1_004295.1)	Link	Link
6: G1_001557	Melanoplax segungensis complete genome [14129]	g[953129] (G1_001557.1)	Link	Link
7: G1_001973	Lucilia hirsuta complete genome [2721]	g[953044] (G1_001973.1)	Link	Link

Entrez Genomes  
Help

Submitting  
Genome Project  
Genome sequence

Microbial Genomes  
Complete Genomes  
List of projects  
PDB neighbors

Genomic BLAST  
Microbial  
Eukaryotic

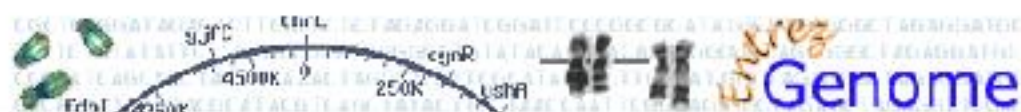
Archaea  
Genome  
Plasmids  
Unfinished

Bacteria  
Genome  
Plasmids  
Unfinished

Eukaryotes  
Genome  
Plasmids  
Organelles

VIROSES

Рассмотрим классификация геномов  
более подробно



CLAST PubMed Nucleotide Protein Structure

Search Genome for  Go Clear

Limits Index History Clipboard

### Complete Archaea Taxonomy / List 36

<a href="#">Halobacterium salinarum</a>	NC_000854	1699995 bp	Apr 7 2003
<a href="#">Halobacterium salinarum</a>	NC_000917	2178400 bp	Dec 17 1997
<a href="#">Halobacterium salinarum</a>	NC_002124	1765 bp	Dec 4 1989
<a href="#">Halobacterium salinarum</a>	NC_002121	1796 bp	Apr 2 1988
<a href="#">Halobacterium salinarum</a>	NC_003158	16541 bp	Oct 16 2001
<a href="#">Halobacterium sp. NRC-1</a>	NC_002807	2014239 bp	Oct 9 2001
<a href="#">Halobacterium sp. NRC-1</a>	NC_002808	365425 bp	Jul 14 2000
<a href="#">Halobacterium sp. NRC-1</a>	NC_001889	191346 bp	Jan 30 1993
<a href="#">Halobacterium salinarum</a>	NC_004531	3918 bp	Jan 15 2003
<a href="#">Methanohalobium salinarum</a>	NC_000909	1664970 bp	Sep 10 2001
<a href="#">Methanohalobium salinarum</a>	NC_001732	38407 bp	Feb 14 2002
<a href="#">Methanohalobium salinarum</a>	NC_001733	16550 bp	Feb 15 2002
<a href="#">Methanohalobium salinarum</a>	NC_001811	8285 bp	May 7 1997
<a href="#">Methanopyrus kandleri A719</a>	NC_003551	1694999 bp	Feb 4 2003
<a href="#">Methanopyrus kandleri A719</a>	NC_002007	5467 bp	Jan 2 1997



Search for

As

 lock



Display

Levels using Filter:

## Archaeoglobus fulgidus DSM 4304

*Taxonomy ID:* 224325

*Rank:* no rank

*Genetic code:* [Translation table 11](#)

*Other names:*

**Archaeoglobus fulgidus str. DSM 4304**[synonym]

[Lineage \(full\)](#)

[cellular organisms](#); [Archaea](#); [Euryarchaeota](#); [Archaeoglobi](#);

[Archaeoglobales](#); [Archaeoglobaceae](#); [Archaeoglobus](#); [Archaeoglobus](#)

[fulgidus](#)

Entrez records	
Database name	Direct links
Nucleotide	<a href="#">174</a>
Protein	<a href="#">4827</a>
Genome	<a href="#">1</a>
Taxonomy	<a href="#">1</a>

**Comments and References:**

Более подробная информация о полностью секвенированном геноме археобактерии *A. fulgidus*



[Klenk HP et al. \(1997\)](#)

Klenk, H.P., Clayton, R.A., Tomb, J.-F., White, O., Nelson, K.E., Ketchum, K.A., Dodson, R.J.,





PubMed Nucleotide Protein Genome Structure PDB Taxonomy

Search Nucleotide for

Limits Preview/index History Clipboard

Display default Show 1 Send to File

NC\_000917 Archaeoglobus fulgidus [gi:11497621]

LOCUS NC\_000917 2178400 bp DNA circular 307 18-JUN-2003

DEFINITION Archaeoglobus fulgidus DSM 4304, complete genome.

ACCESSION NC\_000917

VERSION NC\_000917.1 GI:11497621

KEYWORDS

SOURCE Archaeoglobus fulgidus DSM 4304  
 ORGANISM [Archaeoglobus fulgidus DSM 4304](#)  
 Archaea: Euryarchaeota: Archaeoglobi: Archaeoglobales:  
 Archaeoglobaceae; Archaeoglobus.

REFERENCE 1 (bases 1 to 2178400)  
 AUTHORS Klenk, H.P., Clayton, R.A., Tomb, J.-F., White, O., Nelson, K.E., Ketchum, K.A., Dodson, R.C., Gwinn, M., Hickox, E.K., Giersevan, J.O., Richardson, D.L., Kerlavage, A.P., Graham, J.E., Kyropides, M.C., Fleischmann, R.D., Quackenbush, J., Lee, N.H., Sutton, G.G., Gill, S., Kirkness, E.F., Dougherty, B.A., McKerney, K., Adams, M.D., Loftus, B., Peterson, S., Welch, C.L., McNeil, L.K., Badger, J.H., Glodek, A., Shou, L., Overbeek, R., Gocayne, J.D., Urdeman, J.F., McDonald, L., Utterback, I., Cotton, M.P., Spriggs, T., Artach, I., Kaine, B.L., Sykes, S.M., Sadow, P.W., D'Andrea, K.P., Bowman, C., Fujita, C., Garland, S.A., Mazon, T.M., Olsen, G.J., Fraser, C.M., Smith, H.O., Woese, C.R. and Venter, J.C.

TITLE The complete genome sequence of the hyperthermophilic, sulphate-reducing archaeon Archaeoglobus fulgidus

JOURNAL Nature 390 (6658), 361-370 (1997)

REPLINE [93043:43](#)  
[9389475](#)

Карточка GenBank  
 генома  
 археобактерии  
*A.fulgidus*

# Рассмотрим карточку GenBank NC\_000917 для *A. fulgidus*

LOCUS NC\_000917 2173403 bp DNA circular BCI 13-JUN-2003  
DEFINITION Archaeoglobus fulgidus DSM 4304, complete genome.  
ACCESSION NC\_000917  
VERSION NC\_000917.1 GI:11497621  
KEYWORDS .  
SOURCE Archaeoglobus fulgidus DSM 4304  
ORGANISM Archaeoglobus fulgidus DSM 4304  
Archaea; Euryarchaeota; Archaeoglobi; Archaeoglobales;  
Archaeoglobaceae; Archaeoglobus.  
REFERENCE 1 (bases 1 to 2173403)  
AUTHORS Klenk, H.P., Clayton, R.A., Tomb, C.-F., White, O., Nelson, K.E.,  
Ketchum, K.A., Dodson, R.J., Swinn, M., Hickey, E.K., Peterson, J.D.,  
Richardson, D.L., Karlavage, A.K., Graham, D.E., Kyprides, M.C.,  
Fleischmann, R.D., Coakley, J., Lee, N.H., Sutton, G.G., Gill, S.,  
Kirkness, E.F., Dougherty, B.A., McKenney, K., Adams, M.D., Loftus, B.,  
Peterson, S., Reich, C.L., McNeil, L.K., Badger, J.H., Glodek, A.,  
Zhou, L., Overbeek, R., Gocayne, J.D., Weidman, J.F., McDonald, L.,  
Utterback, T., Cotton, M.D., Spriggs, T., Artiach, P., Raine, D.P.,  
Sykes, S.M., Sadow, P.W., D'Andrea, E.P., Bowran, C., Fujii, C.,  
Garland, S.L., Mason, T.W., Olsen, G.J., Fraser, T.W., Smith, H.O.,  
Woese, C.R. and Venter, C.C.  
TITLE The complete genome sequence of the hyperthermophilic,  
sulphate-reducing archaeon Archaeoglobus fulgidus  
JOURNAL Nature 393 (6650): 364-370 (1997)  
MEDLINE [38049343](#)  
PIRMBD [3389475](#)  
COMMENT PROVISIONAL [REFSEQ](#): This record has not yet been subject to final  
NCBI review. The reference sequence was derived from [AE000782](#).  
FEATURES  
    Location/Qualifiers  
    source 1..2173403  
          /organism="Archaeoglobus fulgidus DSM 4304"  
          /mol\_type="genomic DNA"  
          /db\_xref="taxon:224035"  
    gene complement(436..786)

# Рассмотрим далее поля карточки GenBank

```
FEATURES             Location/Qualifiers
     source            1..2178400
                     /organism="Archaeoglobus fulgidus DSM 4304"
                     /mol_type="genomic DNA"
                     /db_xref="taxon:321325"
     gene              complement(406..786)
                     /locus_tag="AF0001"
     CDS               complement(406..786)
                     /locus_tag="AF0001"
                     /note="hypothetical protein; identified by GeneMark;
putative"
                     /codon_start=1
                     /transl_table=11
                     /product="hypothetical protein"
                     /protein_id="NP_068342.1"
                     /db_xref="GI:11497622"
                     /translation="MQLSIPFWSDFNSAFEEFVKLFLALSIPFWSDFNSISVVISLSM
RIFPQSHFGLISTEERKGLFGSNPRFQSHFGLISIDTLFEDLRNRLADPQSHFGLIST
PRCLPCDPSRRLFCSHFGLISTIR"
     gene              complement(3222..3749)
                     /locus_tag="AF0002"
     CDS               complement(3222..3749)
                     /locus_tag="AF0002"
                     /note="hypothetical protein; identified by GeneMark;
putative"
                     /codon_start=1
                     /transl_table=11
                     /product="hypothetical protein"
                     /protein_id="NP_068343.1"
                     /db_xref="GI:11497622"
                     /translation="MKAASYGVPFQSHFGLISTIRGNDGLPRVVYFQSHFGLISTVK
IVEVLEKDEALSIPFWSDFNCHLQNDQRIRLAPFQSHFGLISTCTLQTSCTNANFLSI
PFWSDFMLELLEKLEKLEMLFQSHFGLISTRNRMVWGRSELQPAIFQSHFGLISTQSSQ
LIRPVKMTFNPILV"
     gene              1320..1781
```

# Продолжение карточки GenBank

```
g2nc      complement(2177643..2178146)
         /locus_tag="AF2436"
CDS      complement(2177643..2178146)
         /locus_tag="AF2436"
         /note="identified by sequence similarity; putative"
         /codon_start=1
         /transl_table=11
         /product="conserved hypothetical protein"
         /protein_id="NP_071258.1"
         /db_xref="GI:1150012"
         /db_xref="COG:C0G1468"
         /translation="MVEGELFVRGTEVRYFVCKTKLULFERNIAEHEESDSVKLGK
LVNRQHFSDPKVRIQRVALDIVRECELEVEVEKEDRMEKADFYQLAYVLYVLSK
HGVPKRSPLSYPKSRKNVSVLETFENLLVRLKSLLEELKLIKSSSEPKPKKSYCTKCA
YYELCF3"
```

BASE COUNT 562096 a 527303 c 531003 g 55801 t

CP131M

```
1  gaaaatggll  caatcgaaa  ttgagtagaa  ggalaaaagt  gcatgcta  lalaatgag
61  atgcacttcc  gaaccctcgc  ggaagtatat  caatgacagc  agccttcaga  aacccttaca
121  ttggaaatag  agggaaaatt  actgatggtt  gaastcagac  caaastggga  tggaaagacc
181  ctttcagccc  tagtctgagt  gtcagytlla  ctcttgltga  aatcagacca  aaatgggall
241  gacaggtttg  ttacgggctt  tgatttgctc  ctccgggctt  ctggttgcaa  tcagaccaca
301  atgggattga  aagtaaaaga  gttcaacctt  ctcactgggt  caactgctt  gttgaaacca
361  gacaaaaaty  gaatlyaaa  gtabllgaa  caalyaaaay  aaacttacc  tctctctgca
421  aatcagacca  caatgggatt  gaaagagtct  tcgggatggg  caaccaggga  gacaccgagc
481  cgttgaaatc  agaccaaaat  gggattgaaa  gtcagcaagg  ctactcggga  gatcctcgaa
541  gagggtatca  gttgaaatca  gaccaaaaty  ggalgaaaay  gaggcaatgc  tggcaaaagc
601  agcgcctcgt  tcttcggttg  caatcagacc  caactgggat  caaacgacag  tgcacctgct
661  taagaaaatg  ctgacagaaa  ttgagttgaa  atcagaccac  aatgggattg  aagagcggcg
721  gaaagatllt  agcaatlltt  caaaagctga  gttgaaatca  gaccaaaaty  ggalgaaaay
781  ttgcacttcc  ctccgcccga  ttggtctcgt  cggcaggeat  gttgcaatea  gaccacaatg
841  ggattgaaag  cgggtctctt  actactcat  ccgsgaagtc  agactcggct  tggtcgaaat
901  caaaccaaaa  tgggallgaa  agagcaagtc  gttgaaatga  gtaglcaaaa  caaalctgca
961  gttgcaatca  gaacaacatg  ggaatgcaag  ttctccacc  caactaccac  tccgcccggc
1021  aaaaagttag  aatcagacca  aatgggatt  caaagagcgc  gttcagaccac  gtcgtcggtc
1081  tctctctgca  agtctgaaat  caaaccaaaa  tgggallgaa  agtctgaaat  gtaglcaaaa
```

# Окончание карточки GenBank

• • •

```
2176561 ccttcgaaaa gggagtcttg aacccaatte aeatagegpc tcaaaaactt zttgacctg
2176621 tcccccetes ccaegtbccc gtggtcggea octatgaggt acacootcac zcccacatca
2176661 cqaacqqlqa qlatllctca acqjccqlaa actlqcljag aayllcllay zallcqaycc
2176741 tccacaggcg ctggtaagac accctccctg cttagctccct gtgcttgacg zccttcacca
2176801 gectctegct gaagtgtctg acgaacttc ttttccctc ctcgctcagc aasateccgc
2176861 egagcctges gtogaagtct gactccgtaa cctcccggtt gtcaccacag togaacacca
2176921 qcctctcaac aallacjquc llqaaaquc ccccaajlc qajjqcvaqa ctcaatctc
2176981 tctcagaggg ctgctggagg tagcttattg cgggatctag ctgggtctgg taatctcag
2177041 aaagsacggc ccagtaaacg agggcaattc csaagctgat cacggcatte actteacttt
2177101 caggaggcac tctgcccctc ttcccgaaac ccaaatcttt gagtatttca togaatccag
2177161 tqlaclaquc qltllcjqca qcaqctctc ccccaajac ctcaaccatq ctlllaqcay
2177221 agctcagatg ctgcagagat ggcagcagct ctgagctatc ctgctcagcc tctctcagaa
2177281 cctttgcaac attgagaata gcctccctca csaagctctt ggcacagctc agcctctctt
2177341 cgtctcgaag gtgatgtca gcctgcctca aacacacctc acccgaaacg agcctctctc
2177401 tgggatagaa tgcagccata tagtaccctg aacggtcga gaagtccatg caactcctt
2177461 cctttgcaag aagctctatt gccttcgaag ttctgcctac agagccsaga zcctagattg
2177521 cgtagattga gttgacagga atggccctc tcccatctc atctctaac taatggctg
2177581 tctcctgccc tctcagcttc ccctccagaca ccaggtatga atctctctt ctcatggccc
2177641 gaccacgaaa accacagctc gtactagggc cacttcgctc agtagctctt ctctcagcc
2177701 tctggcatcg atgagctctt gactgcctta atctctaac gaatgcttt caactcaac
2177761 agtccattcc cgtccagctc caactctaac ttttctctc tctccgata zcttatectc
2177821 cccctcqucc taacacjqlc ctllcjqayq taclaaajql ayjqcvaay ctcaaaajay
2177881 tggccctctc ccactctgtc tgactctctg acctcaaac cctccagctc ctccccccgc
2177941 cccacgatgt ccagagccsc cctgccaate ccacactctt tggctctct zctgaaatgc
2178001 tggctgtgaa ccagccttcc cagcttcaac gcctcctctt cagctctcat zgcactgctg
2178061 cjqjqcvaaaa qcacajqlc llqlllqcaa acccaajayl ayctcaactc jqltccccct
2178121 acaactagct ctccctcaac aactatgatc acaatctcag tctgcttctca taataaacaa
2178181 tggcattctc tgggaaata gtctggattc gcaaacctcg gtcctctgag tgtttgaatg
2178241 tgttcgaaaa acctacaacg ggtctgcaca gtaattctac ttgagcttta attggcaga
2178301 actatqllc atatllajll qllllqaaa llcqlcljag llqllctca llcqaajql
2178361 aaagctctgag gaaaagattt taagcaaaa atgtatctag
```

Accession: [NC\\_000517](#)  
Total Bases: 2178400 bp  
Completed: Dec 17, 1997.

**Feature table:**  
[Protein coding genes](#)  
[Structural RNAs](#)

**BLAST protein homologs:**

- [CCGs](#) (Clusters of Orthologous Groups)
- [3D Structure](#) (Sequences with known structure)
- [TaxMap](#) (Sequences grouped by superkingdom)
- [TaxPlot](#) (3-way genome comparison)
- [CCD](#) (Conserved Domain Database)

Contributor: [TIGR](#)

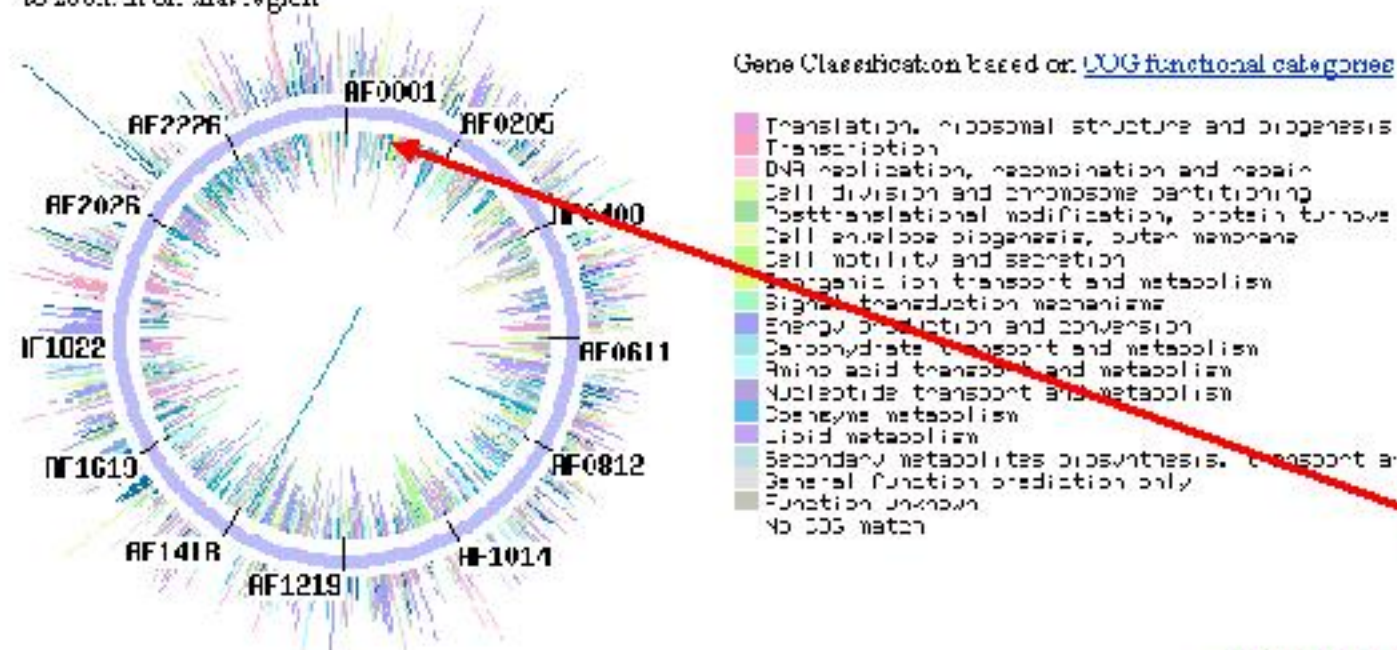
Download chromosome sequence data from [NCEI FTP site](#)

BLAST your query sequence against the genome

Start from:   Search for gene:

**Protein coding genes distribution map**

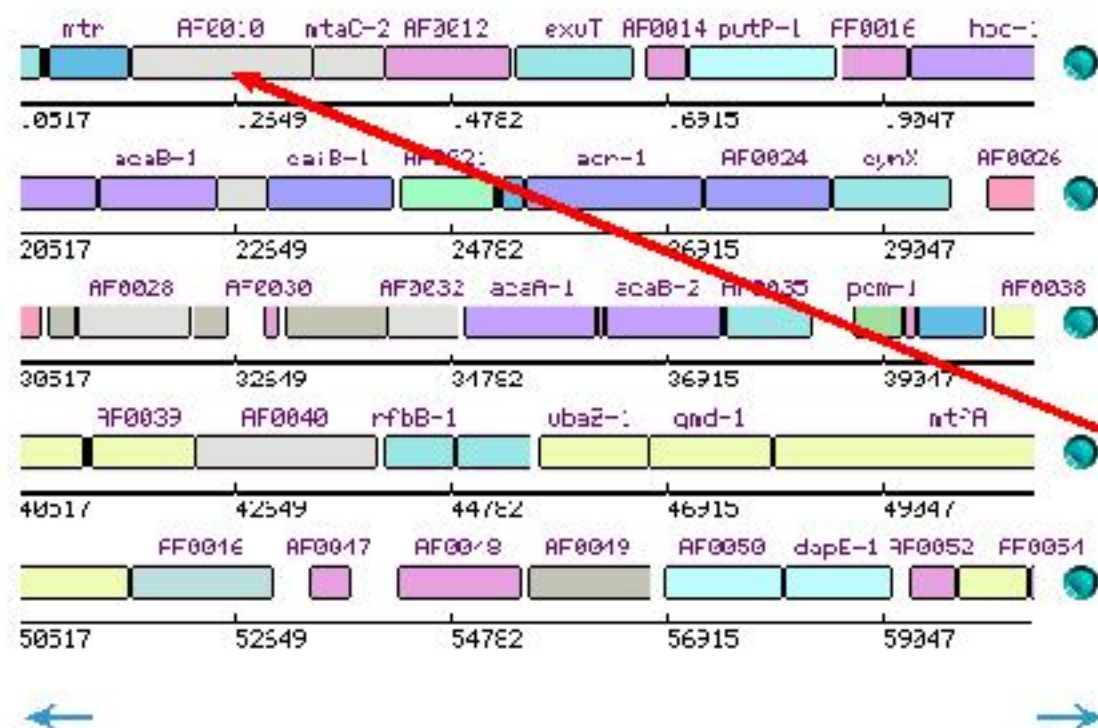
To see map locations of genes, click on a region in the map, to zoom in on that region



Рассмотрим  
графическое  
представление  
генов в геноме  
*A.fulgidus*

Рассмотрим  
район более  
подробно

Click on the rectangle to get BLAST neighbors for the gene of interest or click on the overview below to see a distant region



Указано расположение генов на последовательности

Рассмотрим ген AF0010 более подробно

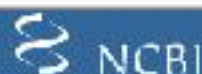
- Translation, ribosome structure and biogenesis
- Transcription
- DNA replication, recombination and repair
- Cell division and chromosome partitioning
- Posttranslational modification, protein turnover
- Cell envelope biogenesis, outer membrane
- Cell motility and secretion
- Inorganic ion transport and metabolism
- Signal transduction mechanisms
- Energy production and conversion
- Carbohydrate transport and metabolism
- Amino acid transport and metabolism
- Nucleotide transport and metabolism
- Coenzyme metabolism
- Lipid metabolism
- Secondary metabolites biosynthesis, transport and catabolism
- General function prediction only
- Function unknown
- No DCG match



Рассматриваемый район на карте генома

# Поиск гомологии интегрирован в банк данных GenBank NCBI

## Результаты поиска гомологичных генов для AF0010, *A. fulgidus*



BLAST	Protein	Structure	PubMed	Taxonomy
Genome	Nucleotide	3D-Domains	Books	Jobs

Query: [U197031](#) hypothetical protein [Archaeoglobus fulgidus DSM 1801]

Mapping: [gi:2650649,7463924](#)

COG2871 assigned by Computer (7 best hits)

[Best hits](#)
[Common tree](#)
[CytochromeP450](#)
[3D structures](#)
[BLAST search](#)
[C lists](#)

95 BLAST hits to 64 unique species. Sort by taxonomy, proximity

[12](#) Archaea [82](#) Bacteria [0](#) Metazoa [0](#) Fungi [1](#) Plants [0](#) Viruses [0](#) Other Eukaryotes

Accession:  Out-Of:

597 hits

	SCORE	E	ALIGNMENT	GI	PROTEIN DESCRIPTION
	1033	2	SP_000072 23111500	23111500	hypothetical protein [Thermotoga sibiricum DSM 1801]
	<a href="#">1027</a>	2	<a href="#">XP_000076</a> 23_111510	23_111510	hypothetical protein [Thermotoga sibiricum DSM 1801]
	<a href="#">924</a>	2	<a href="#">XP_000074</a> 23_11158	23_11158	hypothetical protein [Thermotoga sibiricum DSM 1801]
	631	2	U022473 20197572	20197572	mkkA [Mycobacterium tuberculosis]
	618	2	U022473 20515057	20515057	unclassified hypothetical protein [Thermotoga sibiricum DSM 1801]
	<a href="#">800</a>	10	<a href="#">U022473</a> 2647034	2647034	<i>A. fulgidus</i> predicted coding region 57050 [Archaeoglobus fulgidus DSM 1801]
	857	2	U022473 3171208	3171208	prokaryotic electron transfer protein [Korarchaeum leontii]
	855	2	U022473 15075010	15075010	COG2871B PROTEIN [Simulium venustum]
	841	2	XP_000073 23_12951	23_12951	hypothetical protein [Thermotoga sibiricum DSM 1801]
	<a href="#">872</a>	2	<a href="#">XP_000074</a> 23_13051	23_13051	hypothetical protein [Thermotoga sibiricum DSM 1801]
	873	2	XP_000073 23_12958	23_12958	hypothetical protein [Thermotoga sibiricum DSM 1801]
	745	2	U022473 23737473	23737473	hypothetical protein [Thermotoga sibiricum DSM 1801]
	<a href="#">453</a>	2	<a href="#">XP_000076</a> 33055585	33055585	hypothetical protein [Gallus gallus domesticus]
	316	4	U022473 9911555	9911555	hypothetical protein (nucleic domain) [Korarchaeum leontii DSM 1801]
	313	4	U022473 20005000	20005000	unclassified protein [Mycobacterium tuberculosis]
	253	4	XP_000073 23052581	23052581	hypothetical protein [Thermotoga sibiricum DSM 1801]



# Геномы эукариот в GenBank

## Eukaryote Genomes List

- [5] *Anopheles gambiae*  
chromosomes X, 2, 3
  - [5] *Arabidopsis thaliana*  
chromosomes I, II, III, IV, V
  - [5] *Caenorhabditis elegans*  
chromosomes I, II, III, IV, V, X
  - [5] *Drosophila melanogaster*  
chromosomes 1, 2, 3, 4, Y
  - [11] *Eurephiala formosa cuniculi genome*  
chromosomes I, II, III, IV, V, VI, VII, VIII, IX, X, XI
  - [3] *Gemmatortum nucleomorph genome*  
chromosomes 1, 2, 3
  - [16] *Neohelminthosoma caryocarpum*  
chromosomes I, II, III, IV, V, VI, VII, VIII, IX, X, XI, XII, XIII, XIV, XV, XVI
  - [14] *Plasmodium falciparum*  
chromosomes 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14
  - [3] *Schizosaccharomyces pombe*  
chromosomes I, II, III
- Maps -> See genomes in Map View:**
- [24] *Homo sapiens*  
chromosomes: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, X, Y
  - [21] *Macaca mulatta*  
chromosomes: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, X, Y
  - [12] *Rattus norvegicus (rat)*  
chromosomes: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17

Интегрированная  
система представления  
данных

Работа с GenBank  
Хромосома I человека

MapViewer Home

MapViewer Help  
Human Maps Help  
FTP

Data Available View  
**Maps & Options**  
Compress Map

Region Shown:  
  
 Gc

Legend:  
out  
in

[Homo sapiens Map View](#) [build](#) [BLAST The Human Genome](#)  
**33**

Chromosome [1] 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18  
 19 20 21 22 X Y

Master Map: Genes On Sequence **Maps & Options**

Total Genes On Chromosome: **3232** [153 not localized]  
 Region Displayed: **0-245M bp** [Download Sequence](#) [View Evidence](#)  
 Genes Labeled: **20** Total Genes in Region: **3079**

6H9...	STS Contig	Genes_seq	Symbol	LinkOut	Cyto	Descript
			<a href="#">DNF5</a>	<a href="#">OMIM sv pr dl ev mm hn</a>	C 1p36.1-p36.2	deleted i
			<a href="#">FAD2</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p35.2-p35.1	peptidyl
			<a href="#">LOC284632</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p36.11	hypotact
			<a href="#">EPB41</a>	<a href="#">OMIM sv pr dl ev mm hn</a>	C 1p33 p32	erythrocy
			<a href="#">MGC4796</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p34.2	hypotact
			<a href="#">KIAA0467</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p34.1	KIAA04
			<a href="#">FLT12499</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p32.3	hypotact
			<a href="#">FIN1</a>	<a href="#">OMIM sv pr dl ev mm hn</a>	C 1p31	protein C
			<a href="#">KIAA1107</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p22.1	KIAA11
			<a href="#">FLJ14743</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1p13.1	hypotact
			<a href="#">HSD3B2</a>	<a href="#">OMIM sv pr dl ev mm hn</a>	C 1p13.1	hydroxy-
			<a href="#">HHCDS7</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1q21	TcD37 b
			<a href="#">IL6R</a>	<a href="#">OMIM sv pr dl ev mm hn</a>	C 1q21	interleuk-
			<a href="#">IRTA1</a>	<a href="#">OMIM sv pr dl ev mm hn</a>	C 1q21	immunog
			<a href="#">DKFZP586I151</a>	- <a href="#">sv pr cl ev mm hn</a>	C 1q23.1	DKFZP-

The screenshot shows the EMBL Nucleotide Sequence Database website. The browser title is "The EMBL Nucleotide Sequence Database - Microsoft Internet Explorer". The address bar shows "http://www.ebi.ac.uk/embl/". The website header includes the EMBL-EBI logo and the text "European Bioinformatics Institute". A search bar is visible with the text "Nucleotide sequences" and a "Go" button. Below the header is a navigation menu with links: EBI Home, About EBI, Research, Services, Toolbox, Databases, Downloads, and Submissions. The main content area is titled "EMBL Nucleotide Sequence Database" and contains a description of the database, a search box, and several informational boxes for EMBL-EBI, TPA, and NCBI. A sidebar on the left contains a list of links: Index, Access, Documentation, News, Submission, General info, and Contact. At the bottom of the main content area, there is a table with two columns: "Link" and "Explanation".

**EMBL Nucleotide Sequence Database**

The EMBL Nucleotide Sequence Database (also known as EMBL-EBI) constitutes Europe's primary nucleotide sequence resource. Main sources of DNA and RNA sequences include submissions from medical researchers, genome sequencing projects and patent applications.

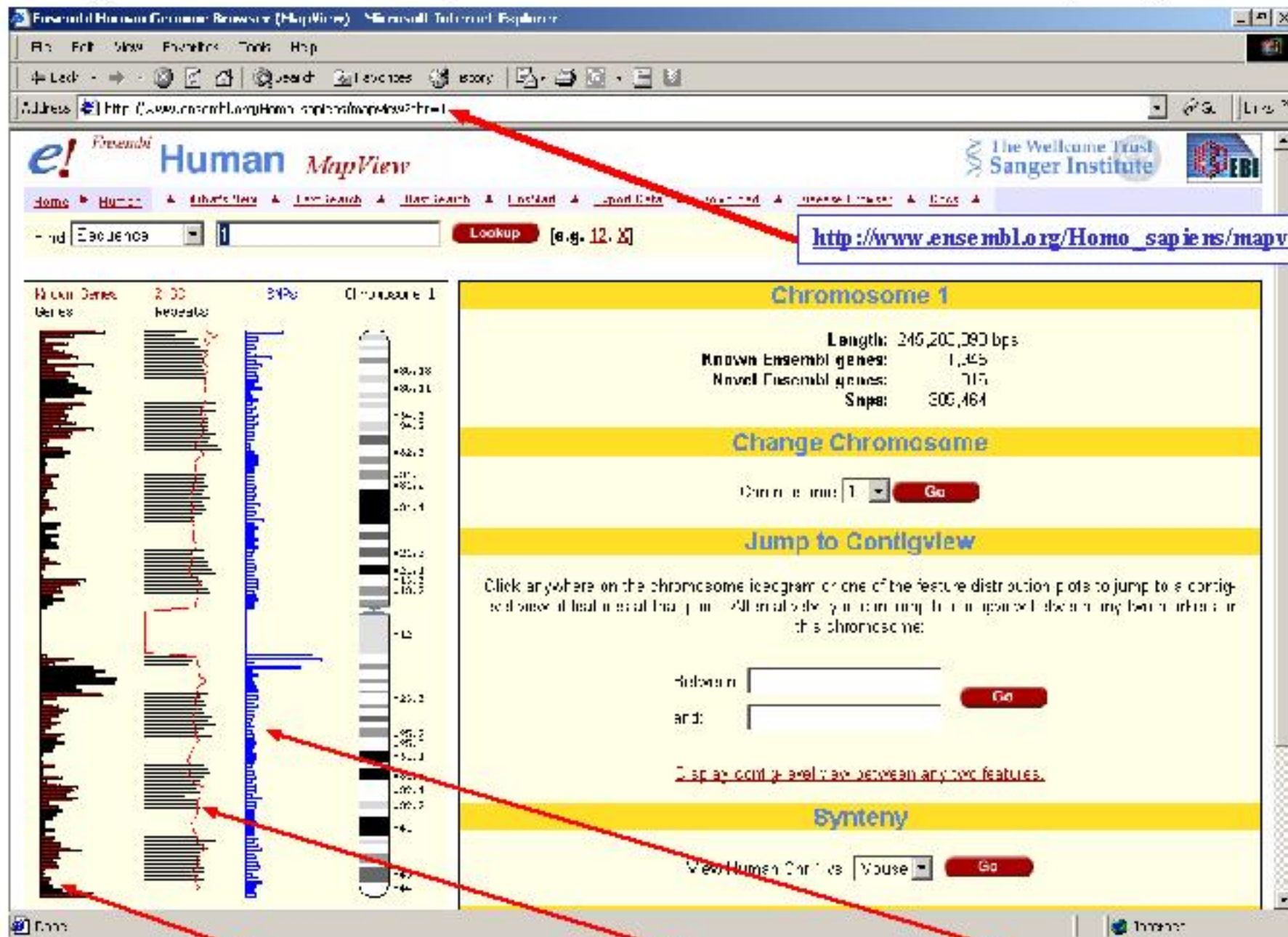
The database is produced in a multinational collaboration with GenBank (USA) and the DNA Data Bank of Japan (DDBJ). Each of the three groups collects a portion of the total sequence data each worldwide, and all new and updated database entries are exchanged between the groups on a daily basis. The current database release (Release 75, September 2002) with accompanying Release 75 manual and associated information for EBI users. A complete database entry is shown [here](#).

A publication in [Nucleic Acids Res.](#) 2002, 30, 2161-17-22 provides further information and details.

The EMBL Nucleotide Sequence Database group is headed by [Rolf Appelber](#).

Link	Explanation
<a href="#">Access</a>	Completed Genomes Webserver, database queries (RTSBI) and FTP archives (EMBL release 75 increments)

# Представление хромосомы 1 человека на сервере EBI



Информация о плотности генов на хромосоме, GC составе и нуклеотидном полиморфизме

# База данных GeneCards (Вейцмановский институт, Израиль)

**GeneCard for gene *TP53***  
**GC17M008311**

Approved HUGO Gene Nomenclature Catabase symbol  
*TP53* (tumor protein p53 (L-Human syndrome))

**Aliases and Additional Descriptions**  
(According to HUGO, UniProt, SWISS-PROT, and/or RefSeq)

- *p53*
- *TRP53*
- *p13*
- tumor protein p53 (L-Human syndrome)
- Cellular tumor antigen p53 (tumor suppressor p53) (oncoprotein p53) (Antigen Y-CC-13)

Previous CC identifier: GC17PC06016

**Chromosomal Location**  
(According to GeneCards and/or HUGO and/or LocusLink and/or UniProt)

Chromosomes: **17** [General gene cistries](#)

Centromeric region: **17p13.1** Telomeric region: **17p13.1**

Gene locations: [Detailed description of GeneCards](#) [Detailed description of GeneCards](#) [Detailed description of GeneCards](#)

**Genomic Views**  
(According to UCSC and Ensembl)

GeneCard Initiative for GC17M008311: [About GC Descriptions](#)

Start: 8,311,492 bp from pter

End: 8,330,670 bp from oter

Основная идея базы GeneCards: анализ и биомедицинских знаний в Интернете и представление их для пользователя.



## Пример карточки гена p53

### OMIM ID: 191170

search databases for OMIM named disorders.

- [Colorectal cancer](#)
- [Li-Fraumeni syndrome](#)

### SWISS-PROT: P53 | UMAP

- **Disease:** TP53 is found in increased amounts in a wide variety of transformed cells. TP53 is frequently mutated or inactivated in about 80% of cancers.
- **Disease:** defects in tp53 are **also the cause of germline cancers such as li-fraumeni syndrome (lfs)** [mim: 191170]. LFS is an abnormal dominant familial cancer syndrome that in its classic form is defined by the existence of both a proband with a sarcoma and two other first-degree relatives with a cancer by age 45 years. In these families the affected relatives develop a diverse set of malignancies including breast carcinoma, sarcoma, and brain tumors at unusually early ages.
- **Disease:** Variant A-a-143 is temperature sensitive. At 37.5 degrees Celsius it possesses strong DNA binding ability, but at 37.5 degrees Celsius its transcriptional activities are greatly reduced.
- **Disease:** Defects in TP53 are also the cause of Janett's adenocarcinoma (JA). JA is a condition in which the normally stratified squamous epithelium of the lower esophagus is replaced by a metaplastic columnar epithelium. The condition develops as a complication in approximately 10% of patients with chronic gastroesophageal reflux disease and predisposes to the development of esophageal adenocarcinoma.
- **Disease:** Defects in TP53 are the cause of head and neck squamous carcinoma (HNSC) and oral squamous cell carcinoma (OSCC). Cigarette smoke is a primary mutagenic agent in cancer of the aerodigestive tract.

**Disorders & Mutations**  
(in which this Gene is  
Involved, According to  
[OMIM](#), [SWISS-PROT](#),  
[GenAtlas](#), [GeneTests](#),  
[HGMD](#), [BCGD](#) and/or  
[IQDE](#).)

### GeneTests: TP53

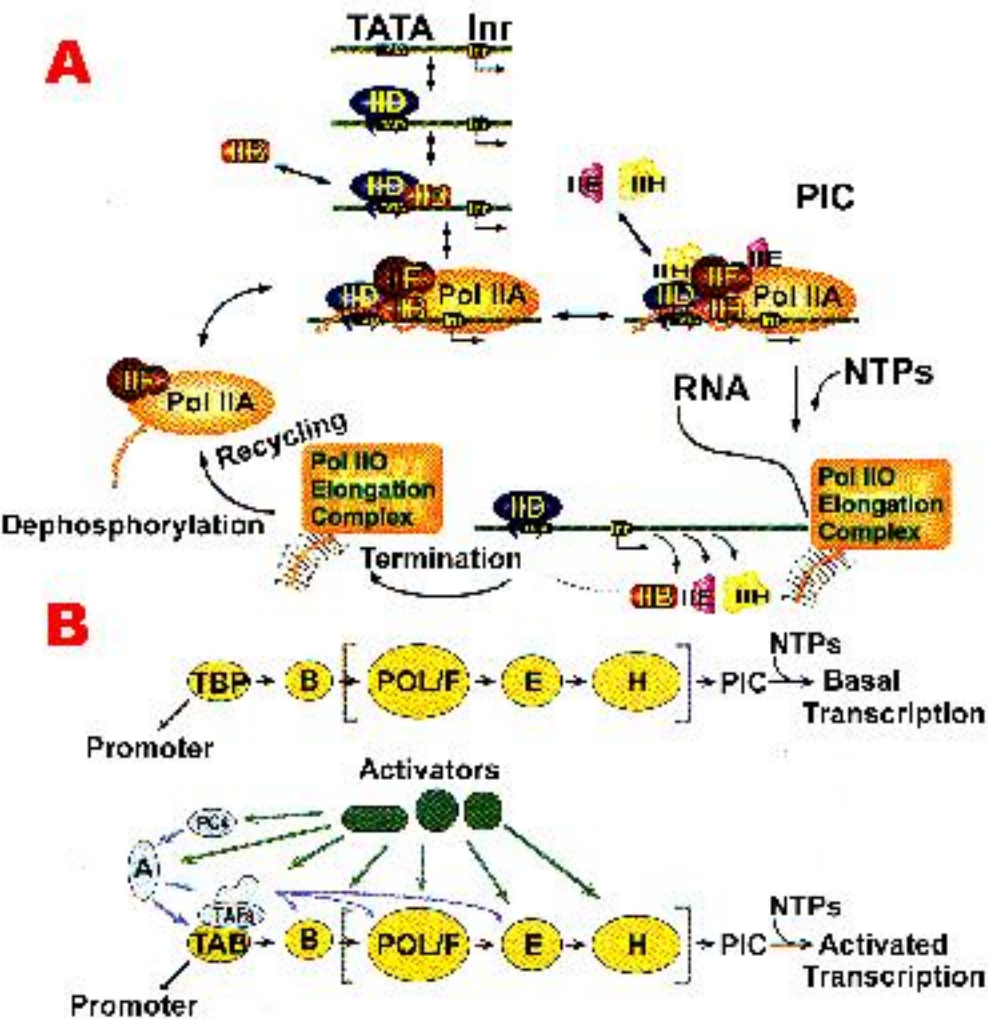
- [Li-Fraumeni Syndrome](#)

**Human Gene Mutation Database** entry for TP53

**Breast Cancer Gene Database** entry for TP53

## Задачи исследования регуляции генной экспрессии

Следует отметить новые подходы к анализу текстов, связанные с моделированием генных сетей - комплексов взаимодействующих макромолекул, включающих гены и их продукты - РНК, белки и метаболиты. Это новое направление биоинформатики оперирует в качественно иных терминах, связанных с биофизическим и биохимическим моделированием и не затрагивает напрямую анализ самих последовательностей.



Эти задачи будут освещены в отдельном курсе

**Регуляторные геномные последовательности (РГП)**



## Список рекомендуемой литературы

**Математические методы для анализа последовательностей ДНК. (Под ред. М.С.Уотермена, Перевод с англ. под ред. П.А.Певзнера), Москва, «Мир», 1999.**

**Франк-Каменецкий М.Д. под ред. (1990) Компьютерный анализ генетических текстов Москва, Наука, 1990, 267 с.**

**Кель А.Э., Колчанов Н.А., Соловьев В.В. Математическое моделирование в молекулярной биологии и генетике. Теория мутационного процесса: делеции и дупликации. Новосибирск, НГУ, учебное пособие, 1989, 86 с.**

**David W. Mount «Bioinformatics. Sequence and genome analysis» CSHL Press, New York, 2001.**

**BLAST By Joseph Bedell, Ian Korf, Mark Yandell. Publisher : O'Reilly, July 2003 (ISBN : 0-596-00299-8 ) 360 pp.**

**Durbin R., Eddy S.R., Krogh A., Mitchson G. Biological sequence analysis. 1998, Cambridge: Cambridge University Press, 356 p.**

**Sequence-Evolution-Function. Computational Approaches in comparative Genomics. (Eds Eugene V. Koonin, Michael Y. Galperin) Kluwer, 2002, 488 pp.**

## Ресурсы:

### ДНК

NCBI (Natl Center Biotech Information) - GenBank <http://www.ncbi.nlm.nih.gov/>

EBI (European Bioinformatics Institute) - EMBL <http://www.ebi.ac.uk/>

ENTREZ <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi>

PUBMED <http://www.ncbi.nlm.nih.gov/pubmed>

DNA Search <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Nucleotide>

GenBank <http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html>

### Выравнивание

BLAST (Basic Local Alignment Sequence Tool) <http://www.ncbi.nlm.nih.gov/blast>

BLAST 2 SEQUENCES <http://www.ncbi.nlm.nih.gov/blast/bl2seq/bl2.html>

### Белки

Expasy (Expert Protein Analysis System, SwissProt, TrEMBL) <http://www.expasy.ch/>

PDB - Protein 3D Structure database <http://www.rcsb.org/pdb/>

Номенклатура IUPAC <http://www.chem.qmul.ac.uk/iupac/misc/naabb.html>

## Рекомендуемый ресурс на русском языке

# molbiol.ru

методы, информация и программы для молекулярных биологов

[Войти](#) | [Регистрация](#)

ПОРТАЛ | О ПРОЕКТЕ | **СПРАВОЧНИК** | МЕТОДЫ | РАСТВОРЫ | РАСЧЁТЫ | ЛИТЕРАТУРА | RSS | Текст | ОБУЧЕНИЕ | WEB-РЕСУРСЫ  
ФИРМЫ | БИОПАТ | КАРТА САЙТА | ПОИСК | COFFEE BREAK | РАБОТЫ И УСЛУГИ | БИРЖА ТРУДА | ФОРУМ