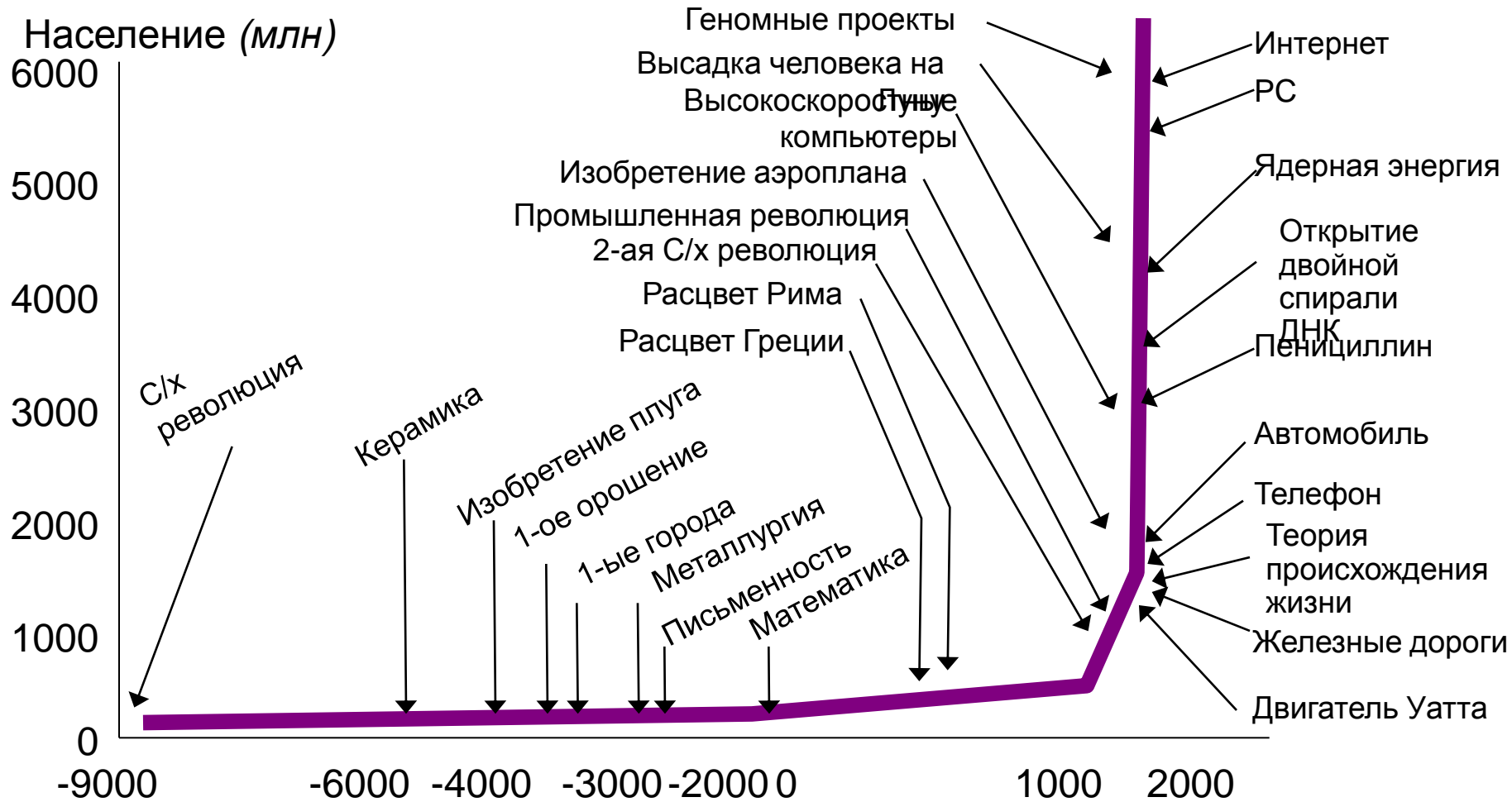# Биоинформатика

Коротков Евгений Вадимович
Институт Биоинженерии, ФИЦ Биотехнологии РАН

bioinf@yandex.ru

# ИСТОРИЯ ТЕХНОЛОГИЙ



Население *(млн)*

6000
5000
4000
3000
2000
1000
0

-9000  -6000  -4000  -3000 -2000  0  1000  2000

Геномные проекты
Высадка человека на Луну
Высокоскоростные компьютеры
Изобретение аэроплана
Промышленная революция
2-ая С/х революция
Расцвет Рима
Расцвет Греции

Интернет
PC
Ядерная энергия
Открытие двойной спирали ДНК
Пенициллин
Автомобиль
Телефон
Теория происхождения жизни
Железные дороги
Двигатель Уатта

С/х революция
Керамика
Изобретение плуга
1-ое орошение
1-ые города
Металлургия
Письменность
Математика

# Чтение и анализ генетических текстов
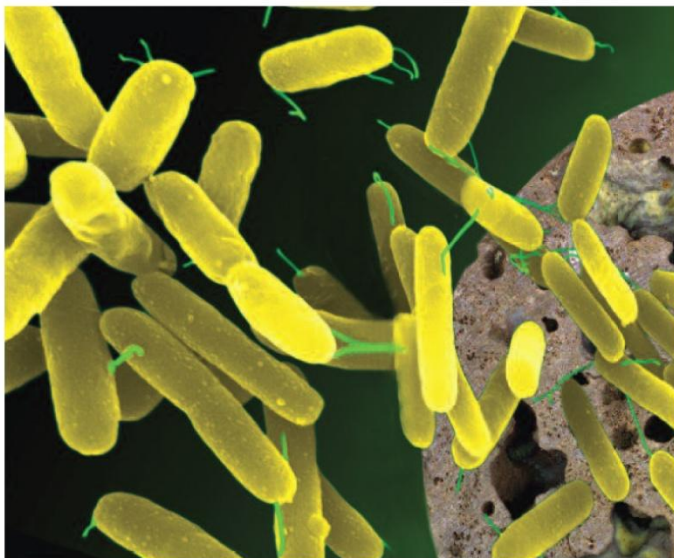
# Манипулирование известными и создание новых, ранее не существовавших в природе, генетических текстов

# Создание органов и организмов с рукотворными генетическими программами

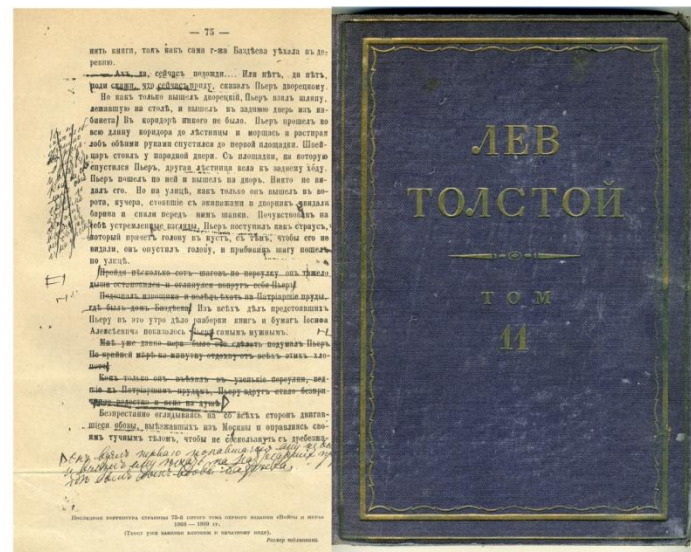число букв
в геноме микроба

число букв в романе
Л.Н.Толстого "Война и Мир"
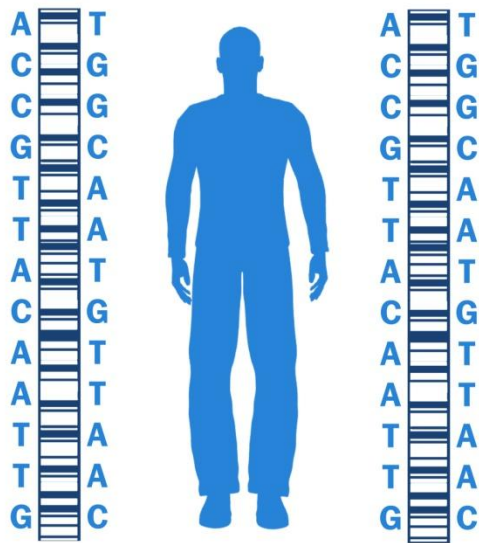




Геном бактерии
2 500 000 букв

роман Л.Н.Толстого
"Война и Мир"
2 500 000 букв
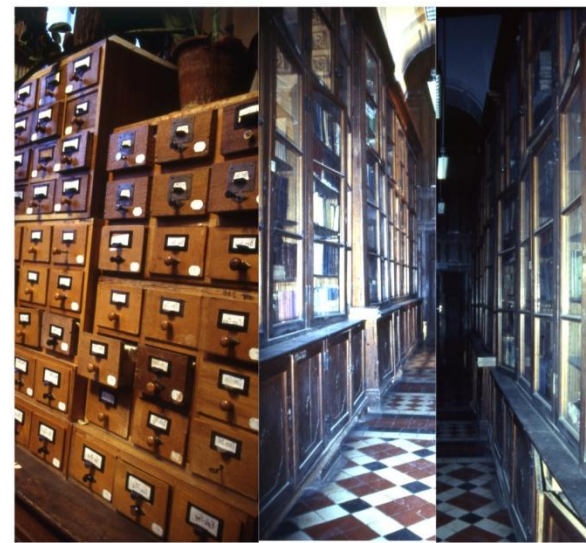
число букв
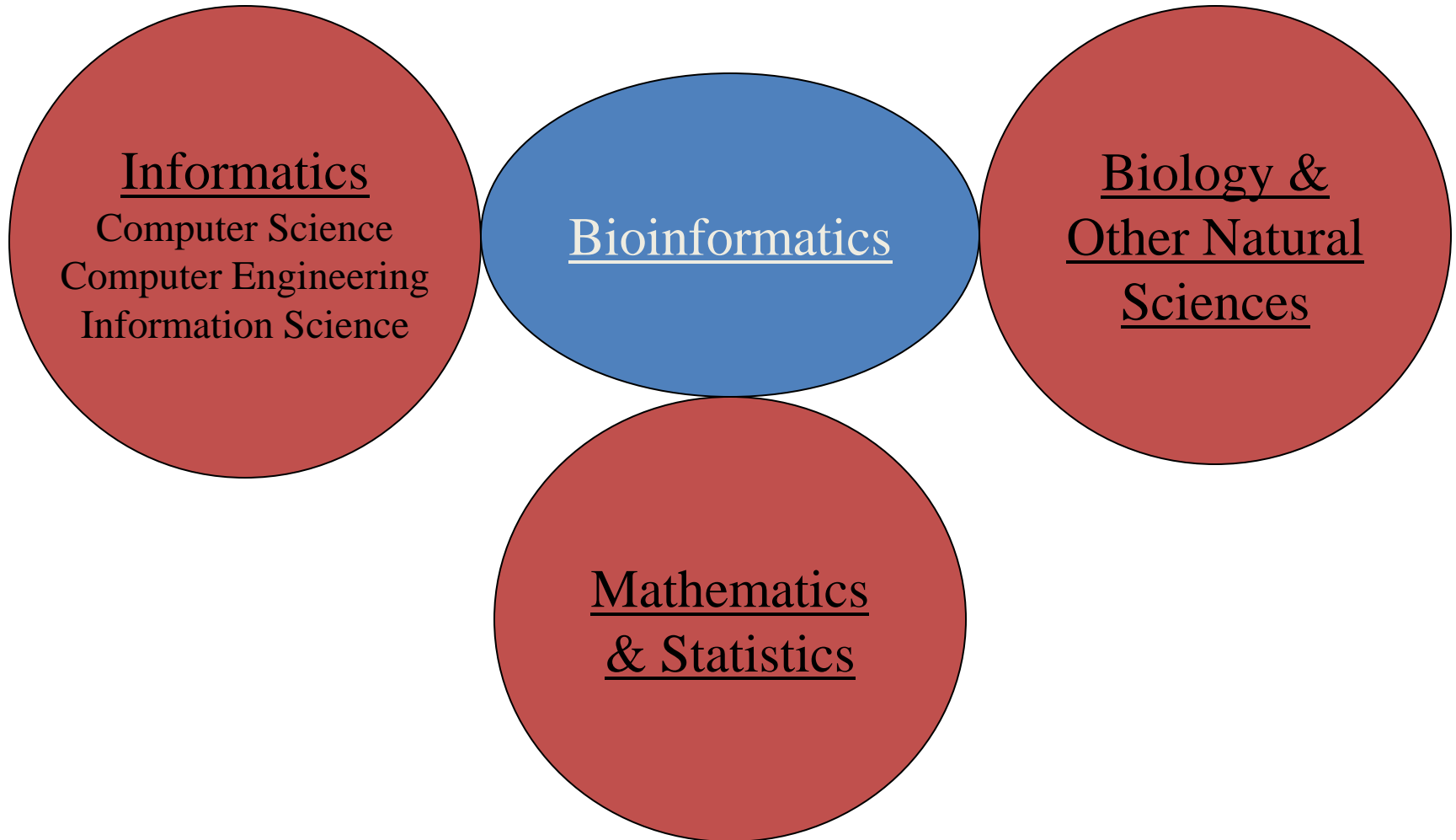в индивидуальном
геноме человека

Геном человека
6 000 000 000 букв

число букв в книгах
библиотеки Л.Н.Толстого
в Ясной Поляне

библиотека Л.Н.Толстого
в Ясной Поляне
6 000 000 000 букв

- В мире идет процесс накопления генетической информации: данные биобанков удваиваются в объеме, примерно, каждые 7 месяцев. Лаборатории по расшифровке ДНК получают несколько петабайт секвенированных данных в год (1 терабайт содержит около 1 трлн субъединиц ДНК).

- Во многих странах создание банков биологической информации выливается в проекты национального масштаба, они становятся системообразующими для мировой науки.

# What is Bioinformatics?

**Informatics**
Computer Science
Computer Engineering
Information Science

**Bioinformatics**

**Biology & Other Natural Sciences**

**Mathematics & Statistics**

# Bioinformatics Related Fields

- Computational biology
- Computational molecular biology
- Biomolecular informatics
- Computational genomics
- …

# Составные части биоинформатики

- 1D и 3D биология
- Разработка биологических баз данных
- Генетические сети и их использование
- Геномика
- Протеомика

Рентгеноструктурный анализ (РСА) макромолекул

Индикаторы качества модели макромолекулы, построенной по данным РСА

Алгоритмы вычисления поверхности макромолекулы

Алгоритмы нахождения гидрофобного ядра молекулы белка

Алгоритмы нахождения структурных доменов белков

Пространственное выравнивание структур белков

Структурные классификации доменов SCOP и CATH

Молекулярная динамика

# Biological Data

- Genomes
  - DNA Sequences of A, T, C, G
  - Annotated with function, "interesting" features
- Proteins
  - Amino Acid Sequences
    - Sequences of 20 letters
  - Annotated with structure, function, etc.

# Biological Data

- Gene Expression
  - Dynamic behavior of genes
- Protein Expression
  - Dynamic behavior of proteins
- Structural Features
  - RNA and proteins
- …

# Biological Data
## Sus scrofa agouti-related protein gene

```
   1 ggcacattct cctgttgagc caggctatgc tgaccacaat gttgctgagc tgtgccctac
  61 tgctggcaat gcccaccatg ctgggggccc agataggctt ggcccccctg gagggtatcg
 121 gaaggcttga ccaagccttg ttcccagaac tccaaggtca gtgcgggcag gagtgggttg
 181 ggtggggctt ggacatcctc tggccacaaa gtattctgct tgtatgagcc ctttcttccc
 241 cttcccaatc ccaggcctgg gaggtgggtg ttttgtgcat gggtggttct gccctcacat
 301 catctgtccc agatctaggc ctgcagcccc cactgaagag gacaactgca gaacgggcag
 361 aagaggctct gctgcagcag gccgaggcca aggccttggc agaggtaaca gctcaggaa
 421 agggctgagg ccacaagtct tgagtgggtg tgtcaagcat caacctctat ctgtgcttgg
 481 agttgccact gtggtacaac gggattggcg gtgtcttggg agcgctggga cgtggtttca
 541 tccccggcca gcacaagtgg gttaaggatc tggccttgcc atcccttcag cttaggctga
 601 gactgtggct tggagctgat ctctgaccgg aagctccata tgctctgggg tgaccaaaaa
 661 tggaaaaaca aacatacaaa acacctctac ctgcacttcc tgaccccctc acccgggggcg
 721 acactgcaga ccatcccgtt cacgctccac ttccatcctg ccttgatctg gcgcattcca
 781 tgaatgtgct tttggaagtc cttgtttccc aacccttgta ggtgctagat cctgaaggac
 841 gcaaggcacg ctccccacgt cgctgcgtaa ggctgcacga atcctgtctg ggacaccagg
 901 taccatgctg cgacccatgt gctacatgct actgccgttt cttcaacgcc ttctgctact
 961 gccgcaagct gggtactgcc acgaacccct gcagccgcac ctagctggcc agccaatgtc
1021 gtcg
```

# Пионеры биоинформатики

## Лайнус Полинг

- **Анализ аминокислотных последовательностей глобинов нескольких позвоночных**

- **Гипотеза молекулярных часов**

Zuckerkandl, E., and L. Pauling. **1962**. Molecular disease, evolution, and genic heterogeneity. Horizons in Biochemistry, Academic Press, New York, 189-225.

Zuckerkandl, E., and L. Pauling. **1965**. Evolutionary divergence and convergence in proteins. Evolving Genes and Proteins, Academic Press, New York, 97-166.

# Пионеры биоинформатики

## Маргарет Дейхофф

- **Однобуквенный код аминокислот A,C,D,E,F,G,H…**

- **Матрицы аминокислотных замен PAM (Point Accepted Mutation)**

**Атлас последовательностей белков и их структур (1965)**

# Секвенирование

✓ 1977 г. Maxam-Gilbert and Sanger Sequencing

✓ 2005 г. Next-Generation Sequencing

➤ Virus – 3222 (Bacteriophage phiX 174, 5386 пн – 1977 г.)

➤ Bacteria – 2289 (*Haemophilus influenza,* 1.8 x $10^6$ пн – 1995 г.)

➤ Eukarya – 168 (*S. cerevisiae* 1.2 x $10^7$ пн – 1995 г; *H. sapien*, 3 x $10^9$ пн -2001 г.)

➤ Archaea – 152 (*Methanococcus jannaschi ,* 1.7 x $10^6$ nt – 1996 г.)

1953 : Discovery of DNA structure by Watson and Crick

1973 : First sequence of 24 bp published
1977 : Sanger sequencing method published
1980 : Nobel Prize Wally Gilbert and Fred Sanger
1982 : Genbank started
1983 : Development of PCR
1987 : 1st automated sequencer : Applied Biosystems Prism 373

1996 : Capillary sequencer : ABI 310
1998 : Genome of Caenorhabditis elegans sequenced
2000 : Human genome sequenced
2005 : 1st 454 Life Sciences Next Generation Sequencing system : GS 20 System
2006 : 1st Solexa Next Generation Sequencer : Genome Analyzer
2007 : 1st Applied Biosystems Next Generation Sequencer : SOLiD
2009 : 1st Helicos single molecule sequencer : Helicos Genetic Analyser System
2011 : 1st Ion Torrent Next Generation Sequencer : PGM
2011 : 1st Pacific Biosciences single molecule sequencer : PacBio RS Systems
2012 : Oxford Nanopore Technologies demonstrates ultra long single molecule reads

# Стоимость секвенирования с развитием технологий NGS



Cost per Raw Megabase of DNA Sequence

| Date | Cost per Mb | Cost per Genome |
|---|---|---|
| Sep-01 | $5,292.39 | $95,263,072 |
| Sep-02 | $3,413.80 | $61,448,422 |
| Oct-03 | $2,230.98 | $40,157,554 |
| Oct-04 | $1,028.85 | $18,519,312 |
| Oct-05 | $766.73 | $13,801,124 |
| Oct-06 | $581.92 | $10,474,556 |
| Oct-07 | $397.09 | $7,147,571 |
| Oct-08 | $3.81 | $342,502 |
| Oct-09 | $0.78 | $70,333 |
| Oct-10 | $0.32 | $29,092 |
| Oct-11 | $0.09 | $7,743 |
| Oct-12 | $0.07 | $6,618 |
| Jan-13 | $0.06 | $5,671 |

# Общий принцип

ДНК нарезается на фрагменты определенной длины

К ним лигируются адаптеры

Амплификация каждого отдельного фрагмента в изолированных от других условиях

Анализ последовательности амплифицированных клонов ДНК

# A schematic of sequencing

# Laser Dye Based Sequencing

# Four-Color Sequencing

# Automated Base Calling

# A Biology Lab?

# Human Genome Project



Minimal set of overlapping BACs selected from physical

BAC shotgun reads
4x - prefinish

Sequence assembly

Whole- genome shotgun reads, including next-gen sequence

Combine overlapping whole-genome and BAC-derived reads

Assemble clone sequences to represent chromosomes and annotate using Ensembl

C:\Users\Korostin\Desktop\DN/

25

# Human Genome Sequencing

# Hierarchical shotgun sequencing

Genomic DNA

BAC library

Organized mapped large clone contigs

BAC to be sequenced

Shotgun clones

Shotgun sequence

```
...ACCGTAAATGGGCTGATCATGCTTAAA
        TGATCATGCTTAAACCCTGTGCATCCTACTG...
```

Assembly ...ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

27

A group of scientists were asked to describe Bioinformatics...

# Fundamental Problems in Bioinformatics

- Pairwise Sequence Alignment
- Multiple Sequence Alignment
- Phylogenetic Analysis
- Sequence Based Database Searches
- Gene Prediction
- Structure Prediction (RNA and Protein)
- Protein Classification
- Gene Expression
- Genetic nets

# Какие бывают выравнивания?

Выравнивания

парные

множественные

глобальные      локальные      глобальные      локальны

# Почему нам интересует парное сходство последовательностей?

Функцию, структуру и многие свойства белка/ДНК определяет последовательность

Родственные белки имеют похожие свойства. Молекулы, похожие по свойствам, похожи по последовательностям.

Свойсва можно предсказать, если мы найдем последовательности похожие на данную.

# Pairwise Sequence Alignment

- Given two DNA or AA sequences, find the best way to "line them up"
  - Biology allows for variation
  - Gaps, mismatches, etc..

```
HEAGAWGHEE
PAWHEAE
```

```
HEAGAWGHE-E
 |   |  || |
P-A--W-HEAE
```

```
HEAGAWGHE-E
  ||  || |
--P-AW-HEAE
```

# Парное выравнивание: вес

Две последовательности:

```
>P1
ALGTEEIC
>P2
ALGTIAA
```

Параметры:

- матрица замен
- штрафы за пропуски

Алгоритм
Нидельмана – Вунша

Алгоритм
Смита – Ватермана

```
P1 ALGTEEIC-
P2 ALGT--IAA
```

Оптимальное полное выравнивание

```
P1 ALGT
P2 ALGT
```

Оптимальное частичное выравнивание

\>sp|P69905|HBA_HUMAN Hemoglobin subunit alpha OS=Homo sapiens GN=HBA1
PE=1 SV=2
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHGKKVA
DALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLA
SVSTVLTSKYR


Бэта-2 субъединицей гемоглобина Rattus norvegicus (Серая крыса).
hemoglobin subunit beta-2 [Rattus norvegicus]

Score = 115 bits (288), Expect = 1e-24, Method: Compositional matrix adjust.
Identities = 63/145 (44%), Positives = 87/145 (60%), Gaps = 8/145 (5%)

```
Query 3   LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLSHGSA-----QV 56
          L+ A+K V WGKV +A GAEAL R+ + +P T+ YF F DLS SA QV
Sbjct 4   LTDAEKATVSGLWGKV--NADNVGAEALGRLLVVYPWTQRYFSKFGDLSSASAIMGNPQV 61

Query 57  KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA 116
          K HGKKV +A + + H+D++ + LS+LH KL VDP NF+LL + +++ L HL
Sbjct 62  KAHGKKVINAFNDGLKHLDNLKGTFAHLSELHCDKLHVDPENFRLLGNMIVIVLGHHLGK 121

Query 117 EFTPAVHASLDKFLASVSTVLTSKY 141
          EFTP A+ K +A V++ L KY
Sbjct 122 EFTPCAQAAFQKVVAGVASALAHKY 1
```

# Sequence Based Database Searches

- Keyword
  - Find all sequences named "cytochrome c"

- Sequence
  - Find all sequences similar to `HEAGAWGHEE`
  - Remember, there are gigabytes to search, and I'm not about to wait two days for an answer!

- BLAST, FASTA, …

# Multiple Sequence Alignment

- Extend pairwise problem to multiple sequences

# Для чего строят множественные выравнивания?



позволяет найти общее

мотивы, паттерны, профили

поиск активного центра

предсказание 3D-структуры

позволяет оценить эволюционные отношения

реконструкция эволюции

**Построение множественных выравниваний — необходимый этап решения многих задач молекулярной биологии**

# Множественное выравнивание

Можно определить вес (хотя ситуация со штрафом за пропуски сложнее)

Но не существует приемлемого алгоритма, гарантирующего нахождение оптимального по данному весу выравнивания

Аналог алгоритма Нидельмана – Вунша имеет приемлемое время работы лишь для очень малого числа последовательностей (до 4–5)

# Программы множественного выравнивания

- **ClustalW** — к настоящему времени явно устарела, но по-прежнему очень популярна
- **Muscle** — пожалуй, на текущий момент программа первого выбора
- **MAFFT** — тоже очень популярная программа
- **DiAlign**
- **T-Coffee**
- **Kalign**
- **ProbCons**

Всё это программы **полного** выравнивания

Единственная популярная программа частичного множественного выравнивания — **MEME** (ищет блоки, то есть выравнивания без пропусков)

# Алгоритм ClustalW – пример эвристического прогрессивного алгоритма



Руководящее дерево

Очевидные недостатки :

- результат зависит от порядка выравниваний;

- «один раз гэп – навсегда гэп»

# Phylogenetic Analysis

- Study relationships between organisms
  - Characteristic similarity
  - Sequence similarity
  - Whole genome comparison
  - …

# Phylogenetic Analysis

**Ортологи** — последовательности, возникшие из одного общего предшественника в процессе видообразования. Ортологи, как правило, имеют одну и ту же функцию

**Паралоги** — последовательности, возникшие из одного общего предшественника в результате дупликации одного гена в одном организме. Паралоги, как правило, имеют разные функции.

# Gene Prediction

- Does the following sequence contain a gene?

```
TTGTAATCTCCTCTGTGACTATAATGACTAGTCTCAGGCCTGCCTTCCCCAGAAACCTCTCTTTTGGCTATTTCTCTTTC
TAGTTCTCTGTTTAAACAAAATTTATTCTATATATCTATCTATCTGTCTATCTATCTATCTATCTATCTATCTATCTATC
TATCTATCTATCTATCATCTACTTATCATCTGTCTAGCCATTTGAAGCATCTTTGTGTTTTAGGTCCTGTTAGATTCTCC
TTTCAGCCAGTGGAGGATCTGGACAGAGCTATTTCTTAGCTTCCCCTAAGCCATGTTGTTAGAACGAATCCCCCACACCT
CCTCTGAGTGCTACGTCTCCGTCAAGAATTATGTATGTGGGATCCAGATGGCCCAGTGGATAAAACTGCAAGTGTCATGA
CCATGACCTGACTTCAAGGGATTGTGTAGAAAGGGAGTTATCACAGTGTGAGGGACAGGGCTAAGGACACTAACCCGTAT
GTTGAGGGGCACAGACGCTAGCAACAACAGTGAAGTGTTTAAAAAGGCAAAAATCATGTTTCTAGAAGTCAGGAAGAGCC
TAACTTGTGGACAAGGACCAACAGGCAGCAGTTGTAATGGGGCAGGGCAGAGGGAGAGCGGACACGCAGCTTTTGGCATC
AAACACACCCAGAGTGTGGATAGAGAGTAGGGAAATACTCTAGTCTCTGGCTAGGATACTCCCCTCTCTTTTTGACATTT
CTCATTGGCAGCCCCAAGTGGTCACTGGAGAGCCAGGAAGCCTAAAGGACACAGTTAGTAGCAGCCAGCTCCTTTGGTGG
AATTTTGGGGACATGGTGGGGTGACTTGGCTCTATCCAGGCCAGGGCTGGGTGTGAGTATACACTTAGTGACTGGCCTTC
```

- How many introns? Exons? Promoters? Other features?

# Genome annotation

# Structure Prediction (RNA, Protein)

- From sequence, predict 2 and 3D structures.

# Protein Classification

- From sequence, identify characteristics of a protein
  - Active sites
  - Families (e.g. globin)
  - Blocks
  - Domains
  - Folds
  - Motifs
  - Etc.

# Protein engineering

# Gene Expression

- Study of gene activity under experimental conditions
  - Large scale studies with microarrays

**Фрагмент одной из карт метаболических путей.**

**Современная биология стала источником огромных объемов экспериментальной информации, осмысливание которых невозможно без использования эффективных информационных технологий и методов математического моделирования**

# Компьютерная технология формализованного описания, конструирования и визуализации генных сетей



Gene network "Macrophage activation" (E. Nedosekina, E. Ananko)

Subscheme "Jak-Stat signal transduction payhway" E. Nedosekina, E. Ananko)

51

# МЕТАБОЛИЧЕСКИЕ ПУТИ – ОБЯЗАТЕЛЬНЫЕ ЭЛЕМЕНТЫ ГЕННЫХ СЕТЕЙ. Адипоцит: мевалонатный путь биосинтеза холестерина в клетке.

# Интеграция генных сетей при противовоспалительном ответе

# Соотношение метаболической и регуляторной компонент цикла трикарбоновых кислот E. Coli K-12:

**Исполняющая компонента (метаболизм)**

**Регуляторная компонента (управление метаболизмом)**



**139** процессов



**1882** процессов

■ - ПРОЦЕСС

→ - участие в процессе с ненулевой стехиометрией

⇢ - участие в процессе с нулевой стехиометрией

**Полный граф метаболической компоненты E. COLI K-12: 3973 процесса**



*Нижние оценки сложности модели (без детального учета этапов матричного биосинтеза):*
*~ 60 000 – 100 000 процессов*

*Более детальная модель:*
*~ 1 000 000 процессов*

*Портретная модель:*
*не менее 10 000 000 процессов*

# Первый "банк данных"



*1965 -1978*

Атлас белковых последовательностей и их структур

Первая версия атласа содержала описание **65 !** последовательностей белков

# Genome Sizes

| Species | Genome Size |
|---|---|
| Bacteriophage MS2 | 3569 bp |
| Esherichia coli | 4.7 million bp |
| Human | 3.3 billion bp |

# Nucleotide Sequence Databases

- **3 main databases**
  - **EMBL:** www.ebi.ac.uk/embl
  - **GenBank:** www.ncbi.nlm.nih.gov/GenBank
  - **DDBJ:** www.ddbj.nig.ac.jp

The 3 databases are synchronized on a daily basis, and the accession numbers are consistent.

There are no legal restriction in the usage of these databases. However, there are some patented sequences in the database

# Protein Sequence Databases



**Swiss-Prot**
Protein knowledgebase
**TrEMBL**
Computer-annotated supplement to Swiss-Prot

**UniProt**
*the universal protein resource*

The UniProt Knowledgebase consists of:

- **UniProtKB/Swiss-Prot**; a curated protein sequence database which strives to provide a high level of annotation (such as the description of the function of a protein, its domains structure, post-translational modifications, variants, etc.), a minimal level of redundancy and high level of integration with other databases [More details / References / Linking to Swiss-Prot / User manual / Recent changes / Disclaimer].
- **UniProtKB/TrEMBL**; a computer-annotated supplement of Swiss-Prot that contains all the translations of EMBL nucleotide sequence entries not yet integrated in Swiss-Prot.

These databases are developed by the Swiss-Prot groups at SIB and at EBI.

UniProt Knowledgebase Release 6.5 consists of:
UniProtKB/Swiss-Prot Release 48.5 of 22-Nov-2005: 199607 entries (More statistics)
UniProtKB/TrEMBL Release 31.5 of 22-Nov-2005: 2406391 entries (More statistics)

> *Swiss-Prot headlines*
Keyword hierarchies and categories (Read more...)

The SWISS-PROT database has **some legal restrictions**: the entries are copyrighted, but freely accessible by academic researchers.
Commercial companies must buy a license fee from SIB.

58

*http://www.expasy.ch/sprot/*

# Анализ белковых последовательностей: Swiss-Prot

Swiss-Prot – одна из первых баз данных белковых последовательностей, "gold standard" белковой аннотации.

Аннотация выполнена вручную группой профессиональных экспертов на основе экспериментальной информации, описанной в научных статьях.

Организована в 1986 году – SIB+EBI+PIR+GU = prof. Amos Bairoch

На сегодняшний день – 556568 последовательностей

# UniProt DB

UniProt = Swiss-Prot + TrEMBL (Translated EMBL sequence database)

TrEMBL – 107 427635  sequences

# Поиск белка в Swiss-Prot (по названию)

# Advances search



**UniProt Knowledgebase (Swiss-Prot and TrEMBL) Advanced Search**

This search program uses SRS to perform queries. Simpler forms are available to search by description or by full text. Available connectors within a field are "&" (and), "|" (or) a "!" (but not). You can prefix your search terms by ! to specify "not" (this is not possible in SRS). Example queries:

- To retrieve all AP1 complex proteins from mouse (AP1S1, AP1G1, etc. but not MIAP1, IQGAP1, ...), specify *Gene Name: ap1\**, *Organism: Mus*, and deselect *"Append a prefix \* to query terms"*.
- To retrieve the three human beta-adrenergic receptor proteins in UniProtKB/Swiss-Prot, but not the beta-adrenergic receptor kinases, specify *Description: beta&adrenergic&receptor!kinase*, *Organism: Homo sapiens*, and select *"Append and prefix \* to query terms"*.

Search ☑UniProtKB/Swiss-Prot ☐UniProtKB/TrEMBL

|  | Description | ATPase |
| AND | Gene name | |
| AND | Organism | | or choose from the list: |

☑ Append and prefix * to query terms

detailed ▾ view of 100 ▾ results

[ Отправить запрос ] [ Сброс ]

This tool can be used to create links to UniProtKB by using the URL of the results page.

*The gory details:*

- **The description line** is indexed as a series of words. If no wildcard (*) is present at the start of the query, it will only match entries where the query is the start of the

62

# Biomolecule Structure Database

- PDB: http://www.rcsb.org
- SCOP: http://scop.berkeley.edu
- CATH: http://biochem.ucl.ac.uk/bsm/CATH
- ASTRAL: http://astral.berkeley.edu
- Interfaces to PDB:
  - **PDB at a glance** http://cmm.info.nih.gov/modeling/pdb_at_a_glance.html
  - **Molecules to go** http://molbio.info.nih.gov/cgi-bin/pdb/
  - **EBI interface**: http://www.ebi.ac.uk/msd/
  - **PDBSum:** http://www.ebi.ac.uk/thornton-srv/databases/pdbsum

# Serine-threonine and tyrosine protein kinases

# Data flow in ASTRAL

The **ASTRAL** compendium provides databases and tools useful for analyzing protein structures and their sequences

# Поиск литературы: PubMed

PubMed is a service of the **U.S. National Library of Medicine** that includes over 18 million citations from MEDLINE and other life science journals for biomedical articles back to the 1950s. PubMed includes links to full text articles and other related resources.

URLs: **www.pubmed.gov**

   **www.ncbi.nlm.nih.gov**

# Поиск по названию белка

# Как это выглядит

# Как получить статью

# Другие виды поиска

По любым ключевым словам или их сочетаниям (AND – необязательно)

По автору (лучше с инициалами!)

По названию статьи

По журналу

По аффилиации авторов

Только в аннотациях

По PMID

По дате – год, либо год/месяц

По словосочетанию – взять в кавычки

# Ген-ориентированные базы данных и геномные браузеры

Что такое ген-ориентированные базы данных?

Самые простые примеры таких БД

Примеры геном-ориентированных баз данных и геномные браузеры

Human Genome Browser

# Что такое ген-ориентированные базы данных?

- Единица исследования – ген (а не экспериментальная последовательность)
- Призваны снабжать информацией по конкретному гену, а не "последовательностям, относящимся ко данному конкретному гену" – интегрируют все такие части в единое целое за Вас

# Первый пример – Gene Entrez (бывший LocusLink) в NCBI

- Единица – генетический локус – конкретное место на хромосоме, кодирующее данный белок и/или соответствующее данному гену

# DUT ген человека

# Продолжение записи:

**Bibliography**
- **Related Articles in PubMed**
- **GeneRIFs: Gene References Into Function**

**Interactions**

**General gene information**
- **Markers**
- **Genotypes**
- **Pathways**
- **Homology**

**GeneOntology**

**General protein information (Names, ECs, ACs)**

**NCBI Reference Sequences (RefSeq)**
- **mRNAs and proteins**
- **Reference assembly + Alternate assembly: Genomic**

**Related Sequences (links between ACs of different types)**

**Additional Links (OMIM, PharmGKB, HRDP, UniGene)**

# Геномные базы данных

Объект – полный геном

Возможность одновременно изучать все гены одного генома

Сравнение друг с другом целых геномов – сравнительная геномика (comparative genomics)

Интеграция всей доступной информации о данном геноме

Основная информация о генах, но в геномном контексте

Геномные браузеры – графическое представление всей интегрированной информации

NCBI -> Genomic Biology (http://www.ncbi.nlm.nih.gov/Genomes/)

# MapViewer

# Sequence Viewer

# Human

Два основных браузера:

Ensembl ([http://www.ensembl.org](http://www.ensembl.org)) – EBI & Sanger Institute, использует свои IDs, 35 эукариотических видов

Human Genome Browser ([http://genome.ucsc.edu/](http://genome.ucsc.edu/)) – UCSC, USA использует GenBank IDs, 41 эукариотический вид

# Human Genome Browser

# DUT gene (dUTPAse)

# Как это выглядит?

# Different perspectives on Bioinformatics

- Bioinformatics is a tool
  - Biologists, biochemists, medical professionals, etc.
  - Obtain meaningful and understandable results
- Bioinformatics is a discipline
  - Informaticians, mathematicians, statisticians, etc.
  - Generate meaningful and understandable results

# Summary

- Bioinformatics is truly interdisciplinary
  - Biology (natural sciences), informatics, mathematics & statistics
- Databases
  - Large, semistructured, incomplete, inaccurate
- Wide-range of problems
  - Solutions employ knowledge from sciences with algorithms and models from informatics, mathematics, and statistics

- Веб-страница для курса
- http://intbio.org/bioinf2018

- Связь с лекторами
- bioinf2018@googlegroups.com
- 

  Онлайн опрос
- https://goo.gl/forms/0RDO3xnIqiotvFYz1