

# Александр Предеус

Институт Биоинформатики

intro@bioinformatics:~\$ █

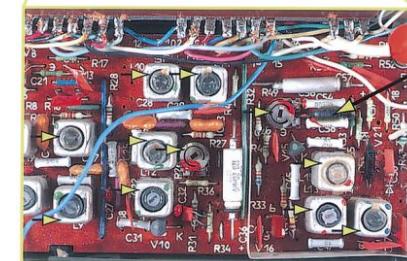
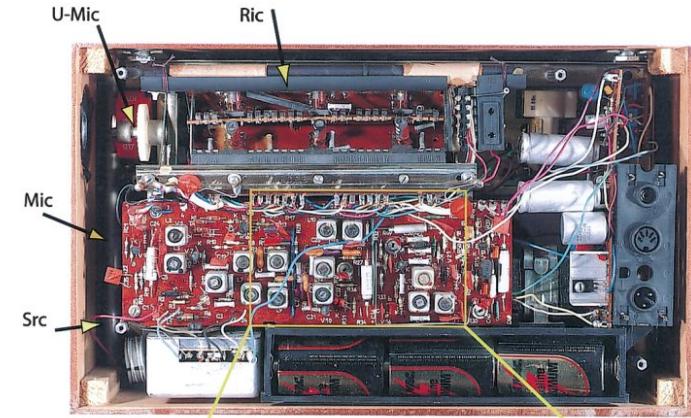
# Биология и математика

- нормальное развитие естественной науки подразумевает увеличение количественных моделей
- описательная стадия → научные модели → инженерия
- биология избегала этого очень долго!



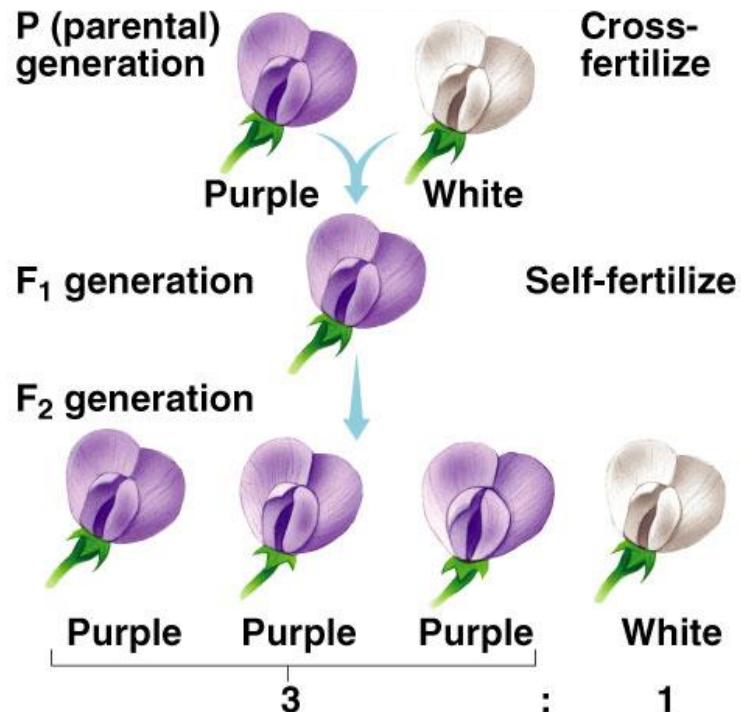
# Способен ли биолог починить радио?

- “Can biologist fix a radio - Or, what I have learnt while studying apoptosis”, by Yury Lazebnik



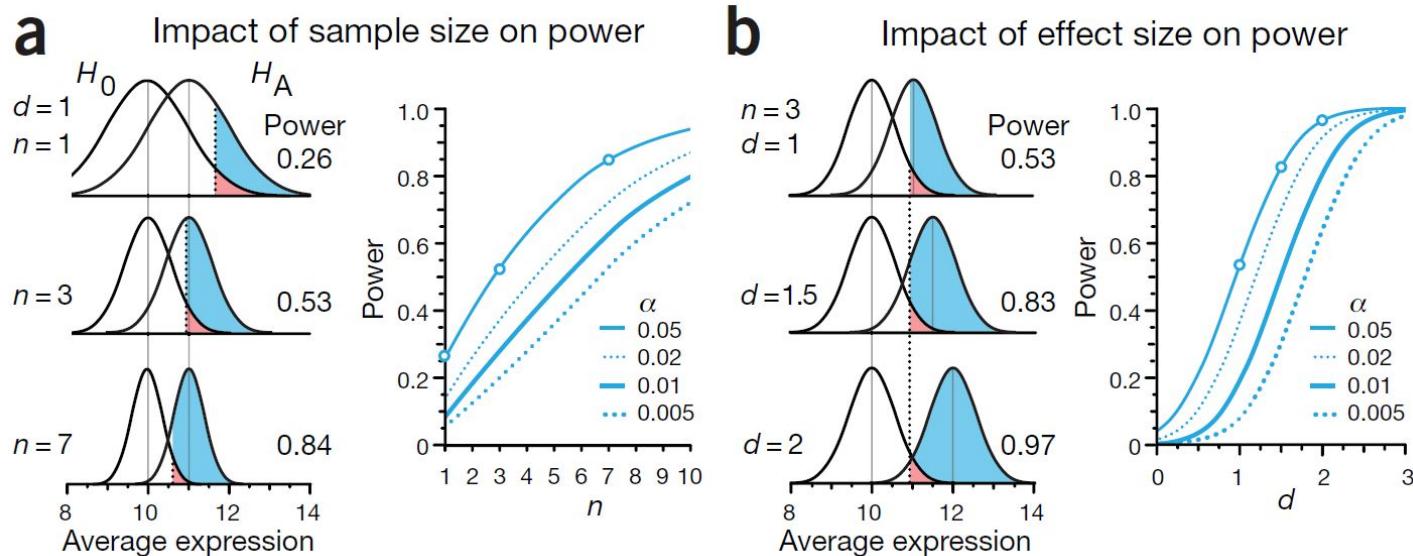
# Эволюция математики в биологии

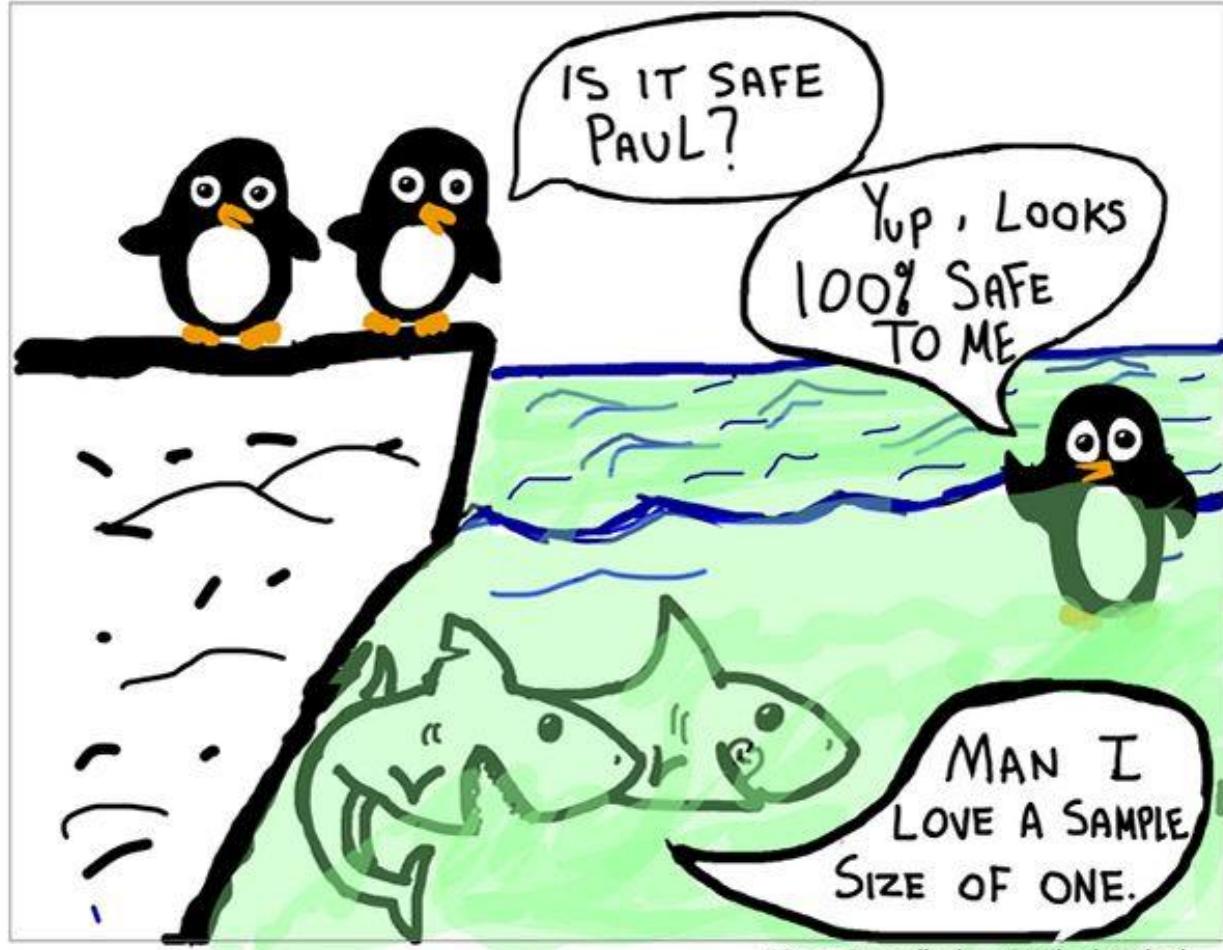
- любая описательная наука требует знания статистики
- статистический анализ - мощнейший инструмент во всех науках, где построение точных моделей затруднено



# Значимость результатов

- размер выборки и эффекта имеет решающее значение в интерпретации любого количественного эксперимента

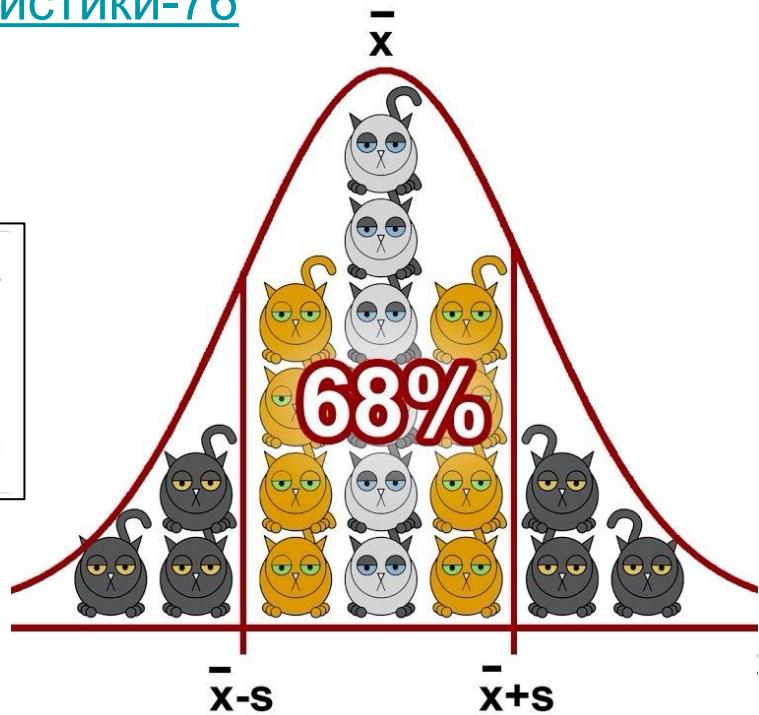
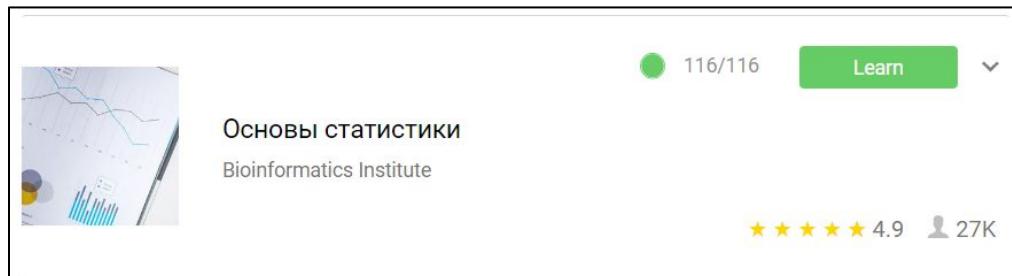




#DataDoodle by @WholeWhale

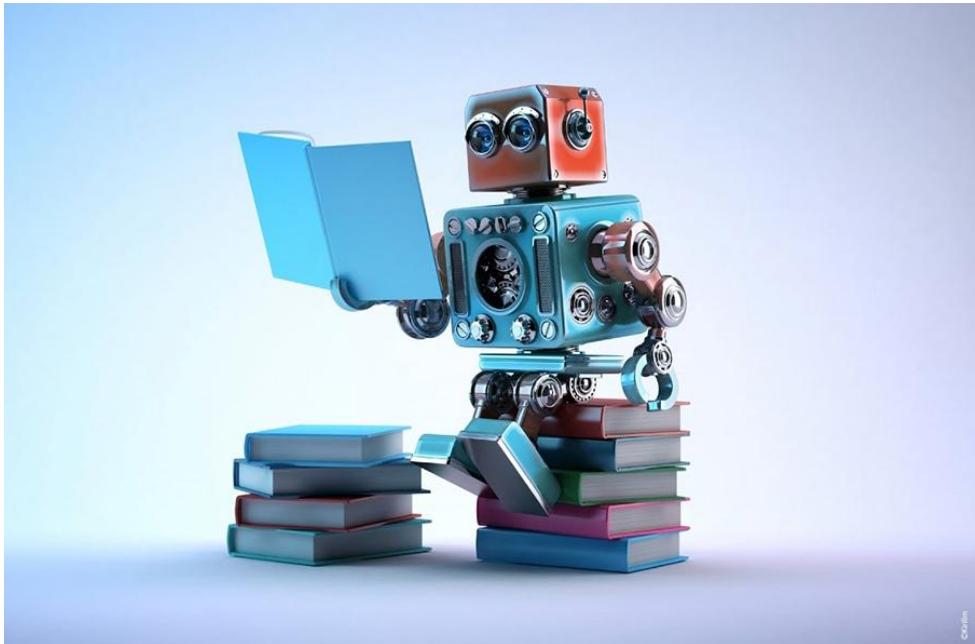
# Статистика и котики

- мало кто любит статистику, зато все любят котиков (с)
- <https://stepik.org/course/Основы-статистики-76>



# Машинное обучение

- классификация
- регрессия
- кластеризация
- уменьшение размерности
- нейросети
- поиск аномалий



# Эволюция “биоинформатики”

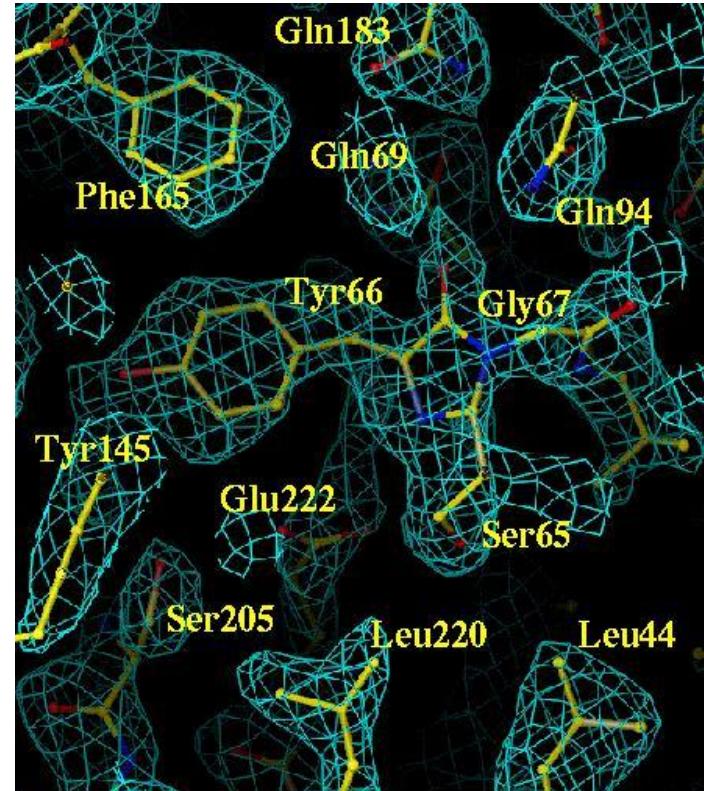
- математическая биология
  - дифференциальные уравнения, кинетические модели, и т.д.
- структурная биоинформатика
  - ака биофизика
  - рост связан с бумом кристаллографии протеинов в 90-е
- геномная биоинформатика
  - рост связан с бумом секвенирования
- системная биология
  - графы (сети), байесовские модели, анализ потоков, и т.д.
  - попытки интегрировать различные данные

# Структурная биоинформатика

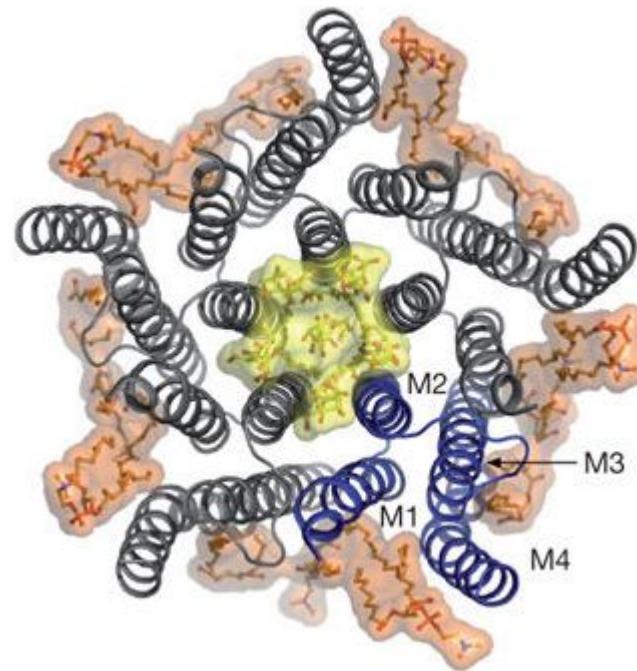
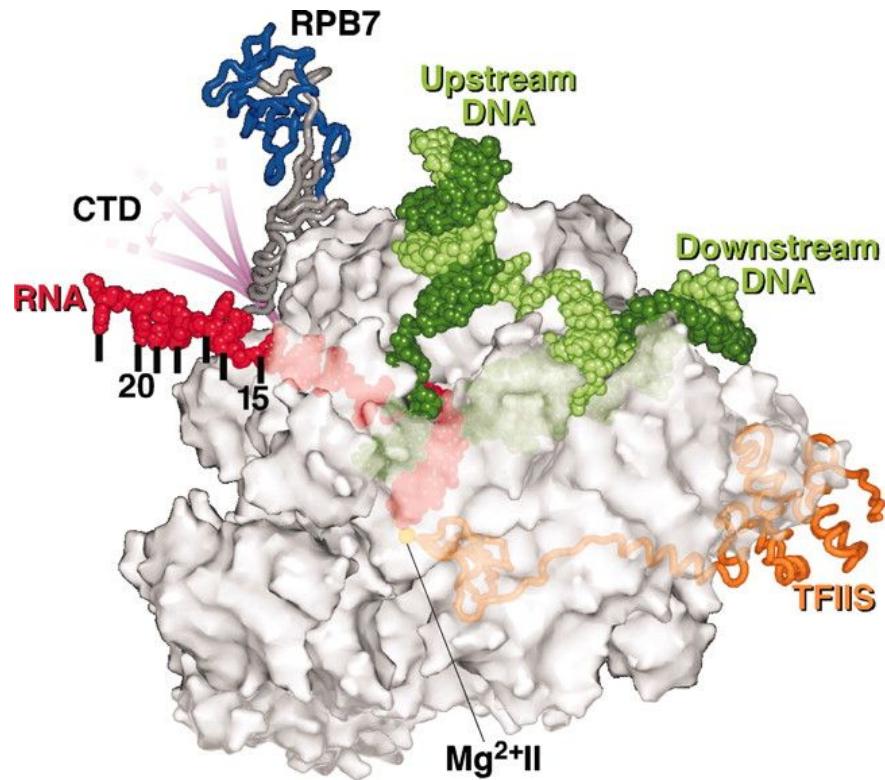
- предсказание и фолдинг неизвестных структур белков
- симуляция поведения белков и комплексов белок-ДНК
- “крупнозернистые” симуляции больших биологических объектов (органелл, вирусов, клеток, участков мембраны)
- докинг малых молекул (лекарств)

# Кристаллография протеинов

- монокристалл протеина облучается рентгеновскими лучами
- дифракционная картина позволяет восстановить позиции индивидуальных атомов с высокой достоверностью

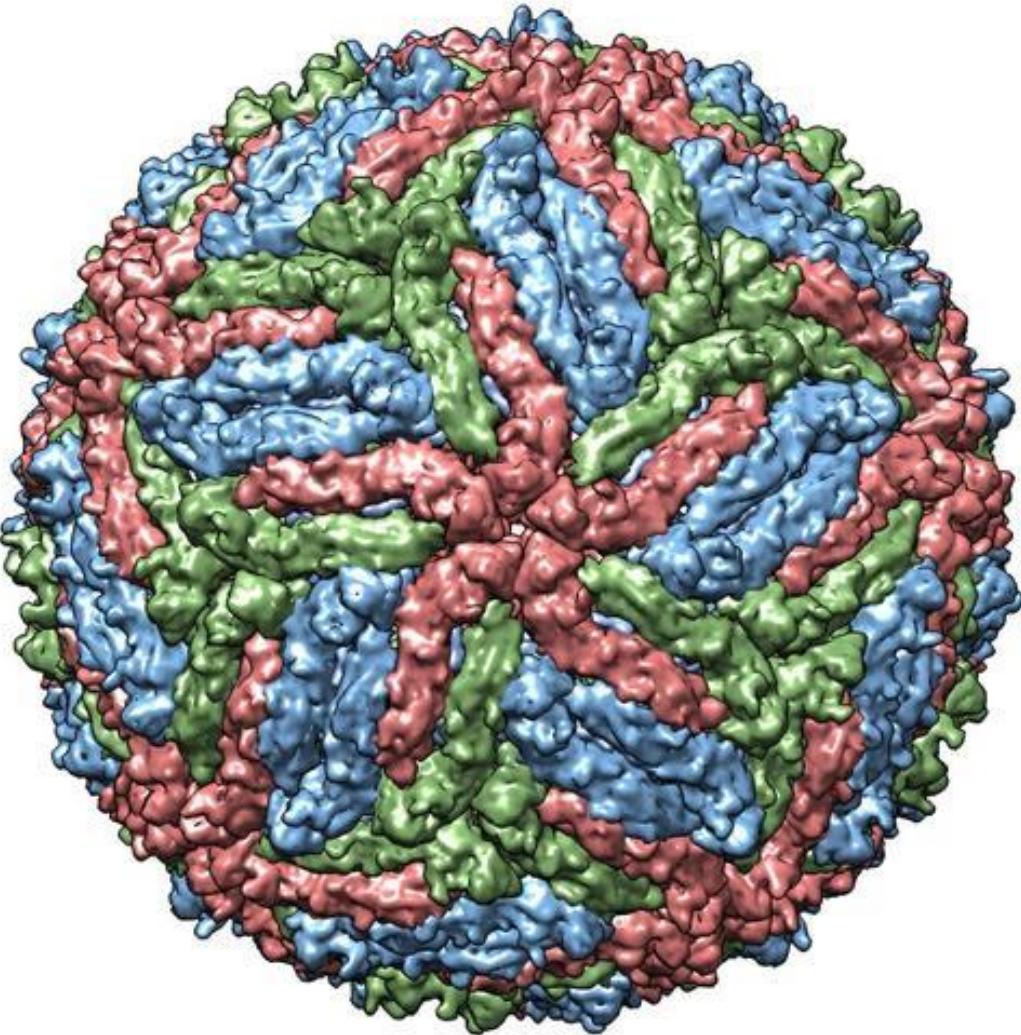


# Структуры, потрясающие воображение



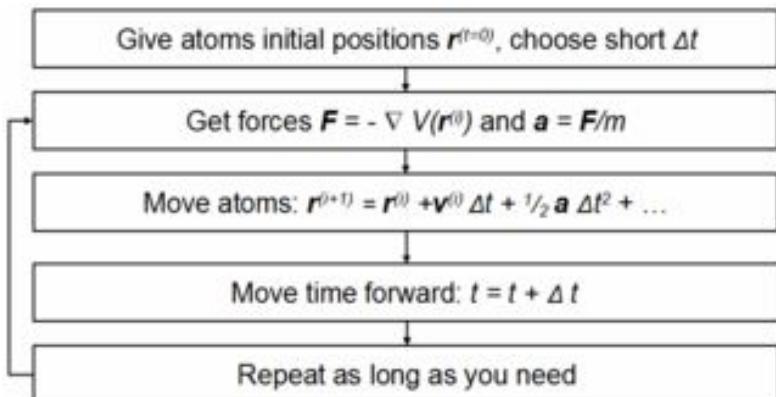
# Cryo-EM

- не требует кристаллизации!
- разрешение становится все лучше
- незаменим для определения структур огромного размера
- например, вируса Зика

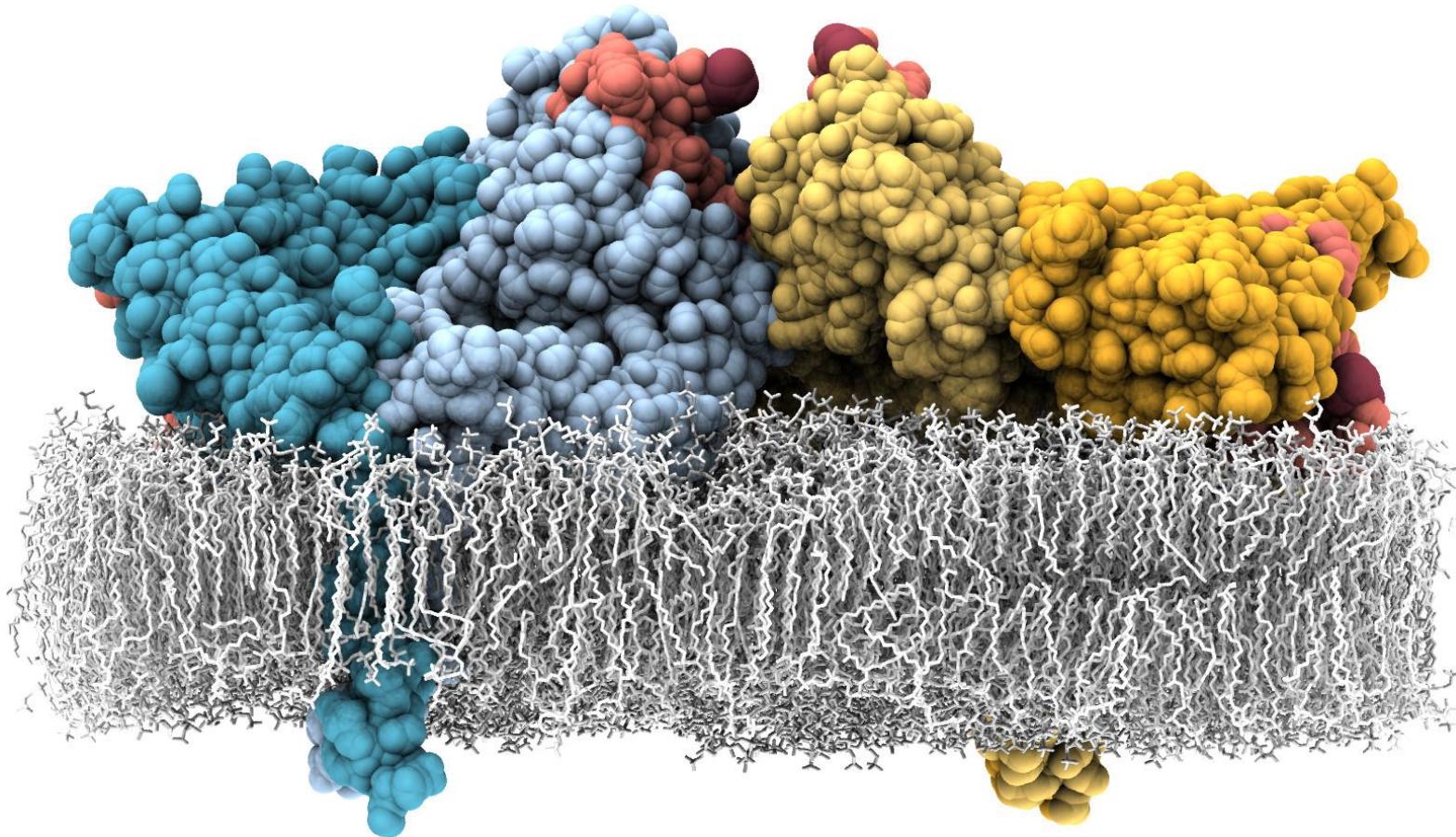


# Молекулярная динамика

- решение  
ニュтоновских  
дифференциальных  
уравнений движения

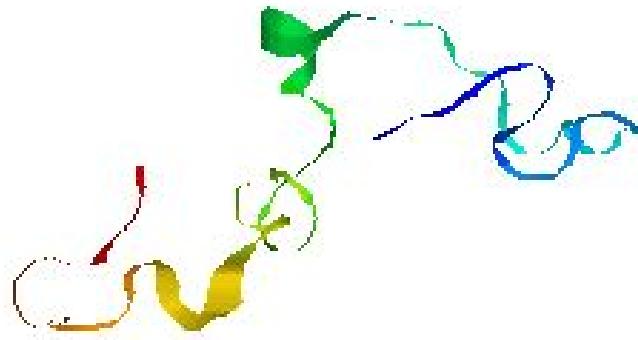


# Симуляции миллионов атомов



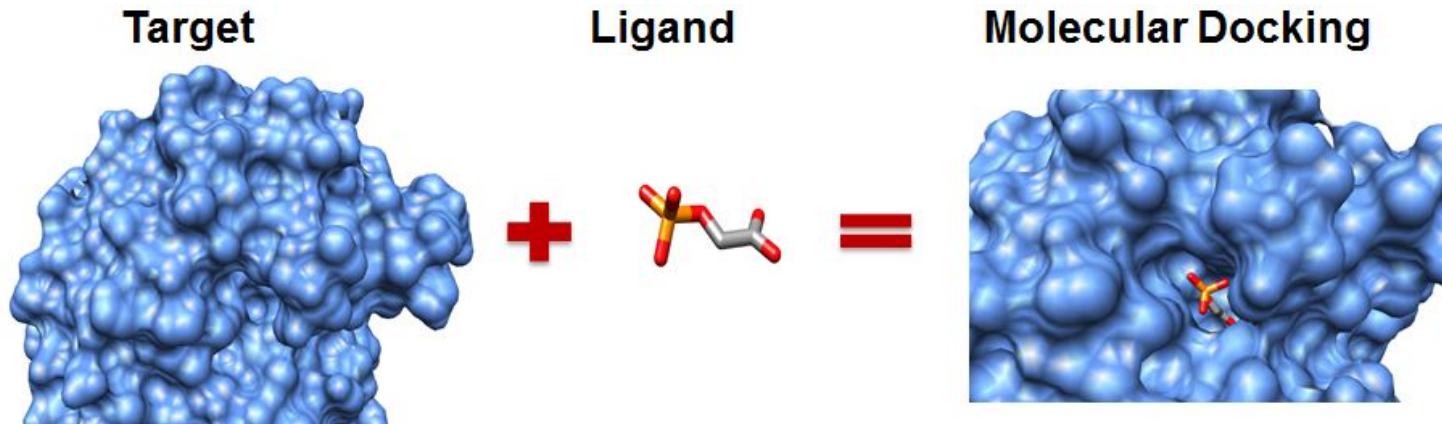
# Фолдинг протеинов

- отличается от молекулярной динамики целью и подходами
- молекулярная динамика может использоваться для дальнейшего улучшения структуры



# Докинг малых молекул

- большинство современных лекарств ингибируют ферменты, блокируя активный центр
- предсказание активных молекул может сильно снизить цену разработки новых лекарств



# Геномная биоинформатика

- прогресс связан с прогрессом в секвенировании
- секвенирование - определение последовательности (ДНК, РНК, протеина)
- задачи:
  - поиск нуклеотидных и протеиновых последовательностей в базах данных (*blast*)
  - сборка геномов *de novo* (*De Bruijn graphs*)
  - выравнивание на собранные геномы и транскриптомы (*FM-index, Burrows-Wheeler transform*)

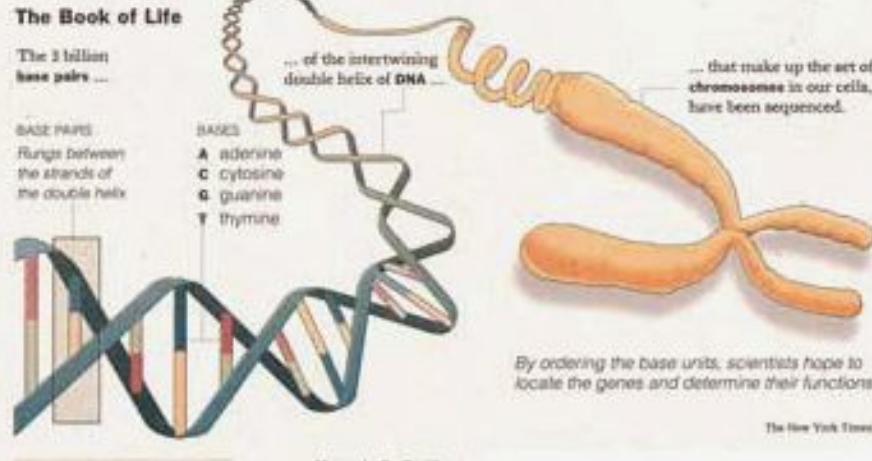
# Проект “геном человека”

ws  
Print"

# The New York Times

No. 51,432 Copyright © 2000 The New York Times TUESDAY, JUNE 27, 2000 Printed in Arizona ONE DOLLAR

## *tic Code of Human Life Is Cracked by Scientists*



### National Edition

Arizona and New Mexico: M  
cloudy in New Mexico; thunder  
in the mountains. Partly sunny  
where. Highs 80 mountains, over  
deserts. Weather map is on Page

### A SHARED SUCCESS

2 Rivals' Announcement Marks New Medical Era, Risks and All

By NICHOLAS WADE

WASHINGTON, June 26 — An achievement that represents a milestone of human self-knowledge, rival groups of scientists said today that they had deciphered the human genetic script, the set of instructions that defines the human organism.

The New York Times

# Университеты vs. Celera Genomics



Francis Collins



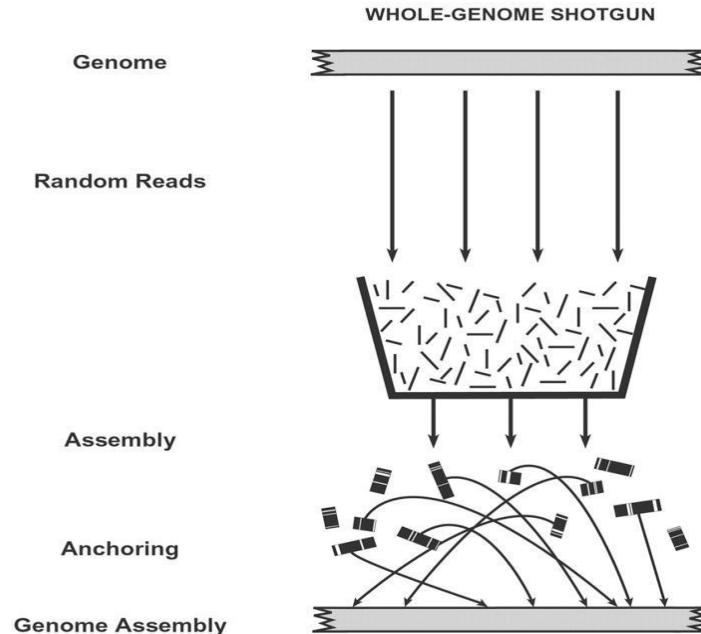
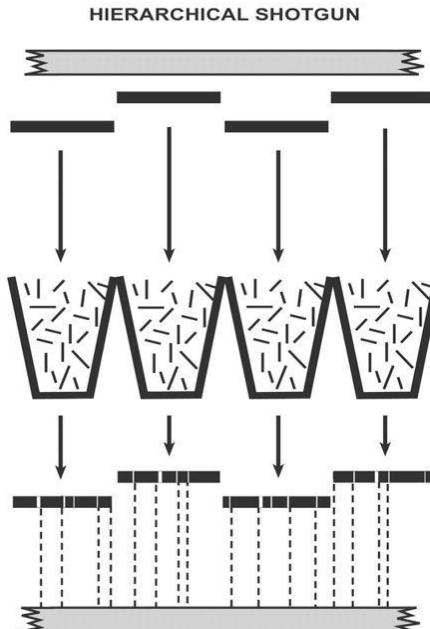
Eric Lander



PHOTO: ANDREW HARRER/BLOOMBERG VIA GETTY IMAGES

Craig Venter

# Соревнование стратегий

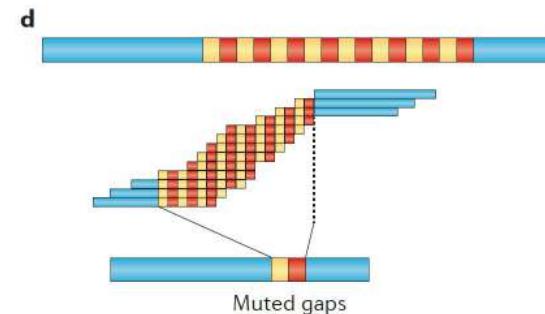
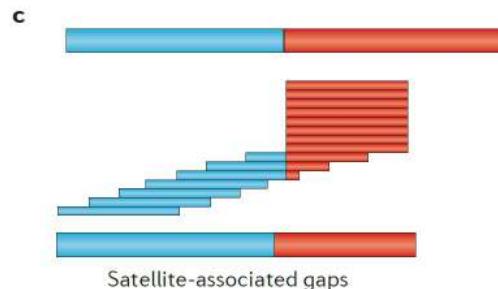
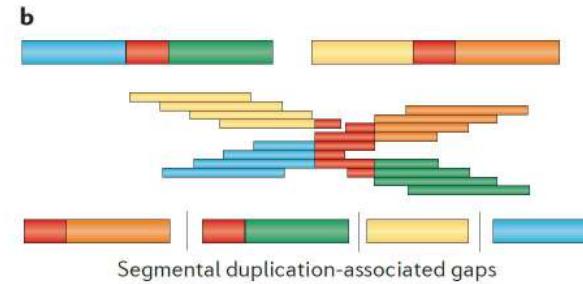
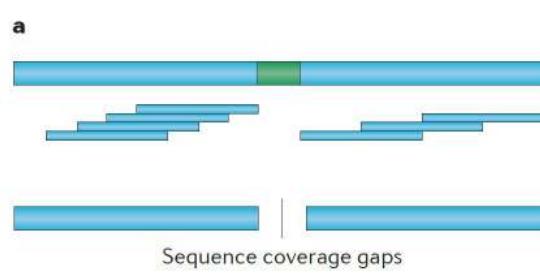


**Public (Universities)**  
1990-2001 (2003)  
3 billion dollars

**Celera Corporation**  
1999-2001 (2003)  
300 million dollars

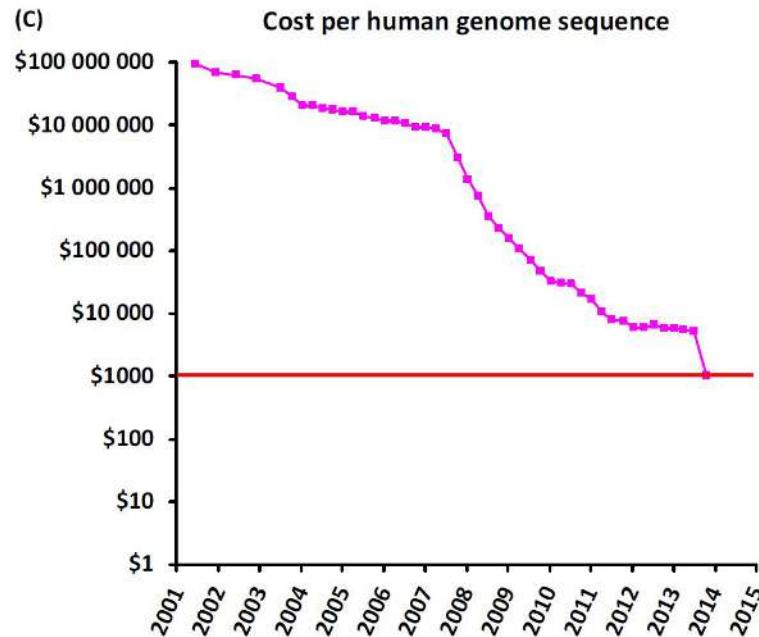
# Возможно ли?

- Сомнение в возможности сборки коротких прочтений в геном
- Трудности – повторы, плохое покрытие



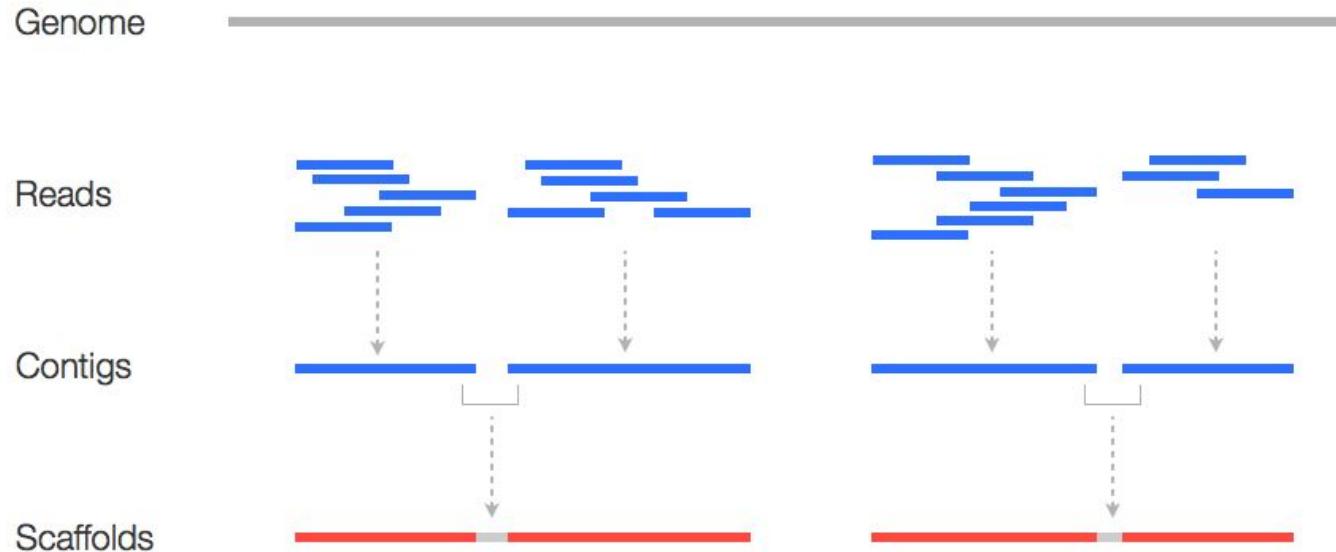
# Секвенирование: революция 2000-х

- Секвенирование генома человека упало в цене от 100 миллионов до ~1000 \$



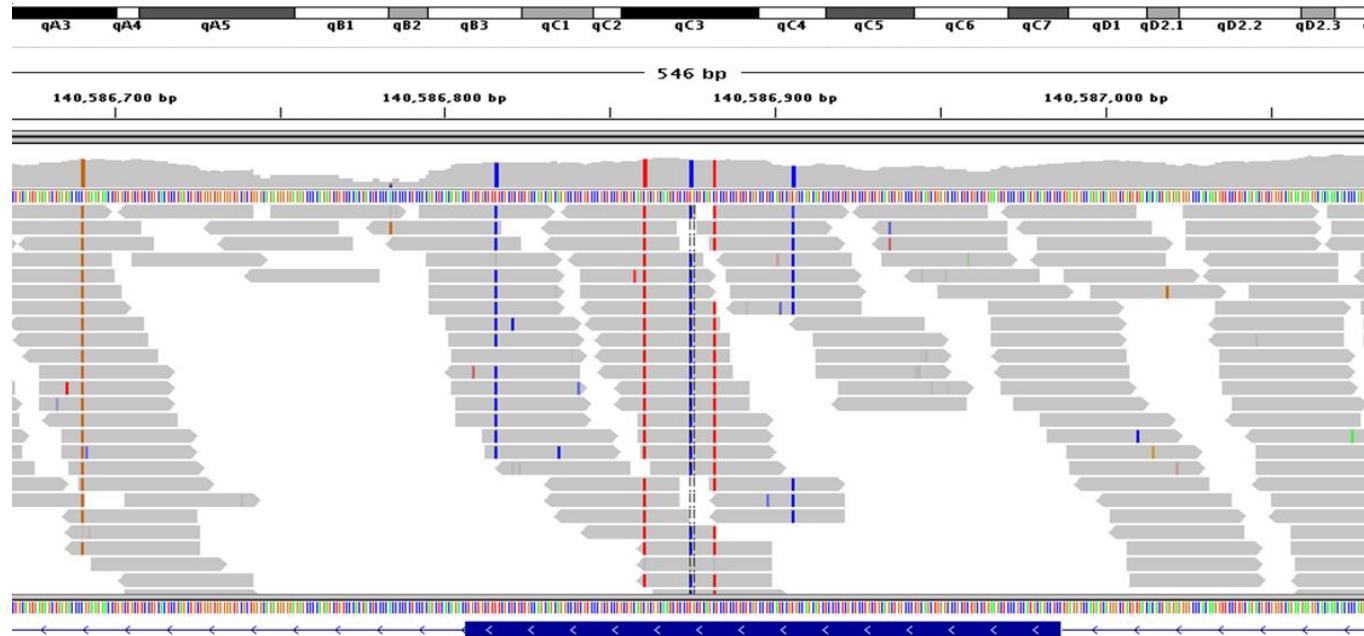
# Тип 1: неизученные виды (сборка)

- Сборка геномов (de-novo genome assembly)
- Сборка транскриптомов (de-novo transcriptome assembly)



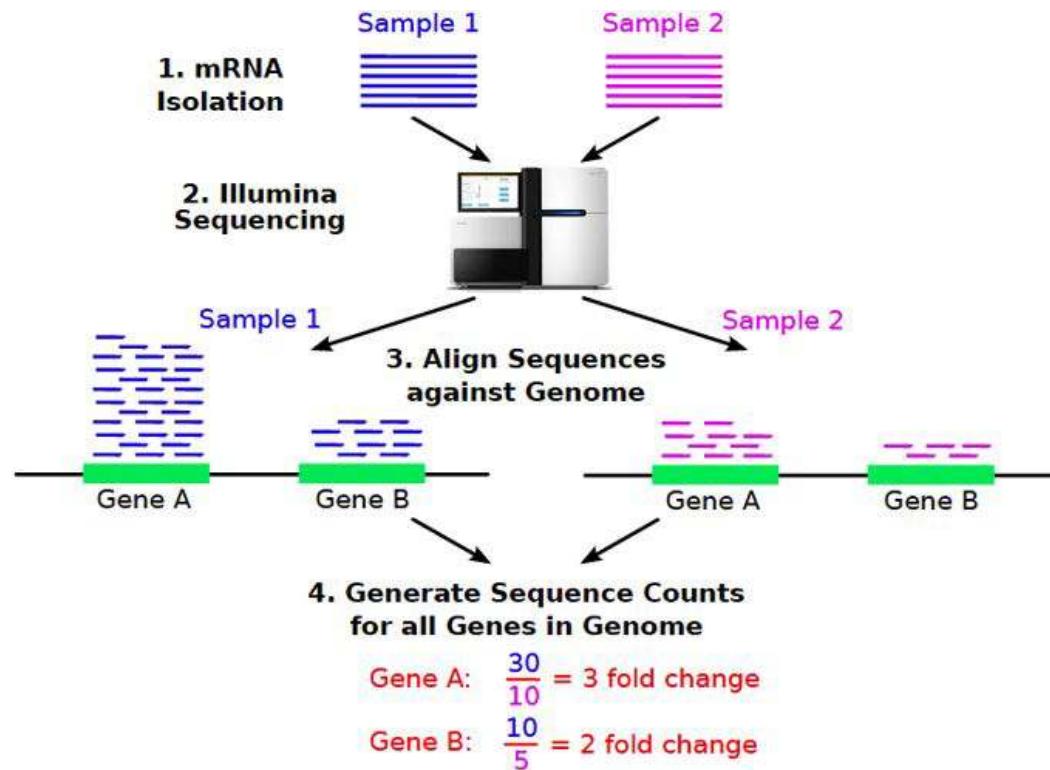
# Тип 2: Уточнение последовательности (выравнивание)

- Ресеквенирование геномов
- Ресеквенирование частей генома



# Тип 3: Количественный анализ (выравнивание)

- Транскриптом (RNA-seq)
- Позиционный (ChIP-seq, genome footprinting)



# Что (и кого) мы секвенируем?

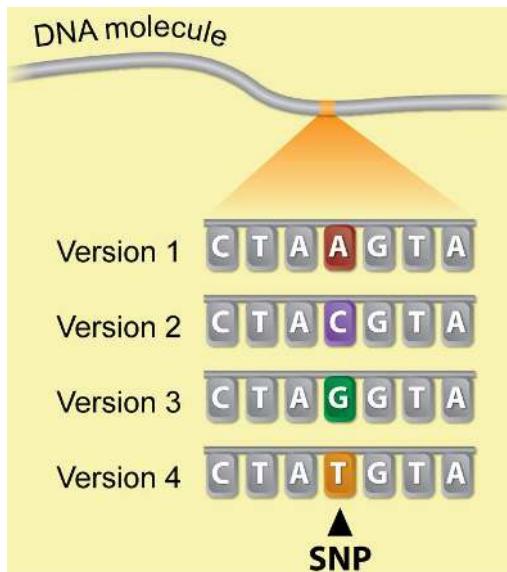
- Плазмиды и конструкты
- Вирусы
  - новые
  - количество (вирусная нагрузка)
  - варианты (штаммы)
- Прокариоты
  - новые
  - варианты (штаммы)
- Эукариоты
  - новые
  - модельные животные
- Люди
  - пациенты
  - исследование здоровых волонтеров

# Примеры применения NGS

- исследование вариации в геноме человека
- определение присутствующих организмов (“баркодинг”)
- анализ дифференциальной экспрессии - сравнить те же ткани у больных и здоровых людей
- поиск участков ДНК, активных в определенных тканях и условиях

# Вариативность в геномах

- SNP (SNV) - Single Nucleotide Polymorphism (Variant)
- Может быть гетеро/гомозиготный

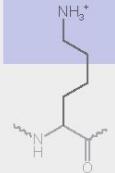
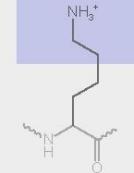
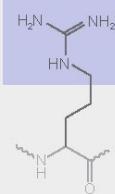
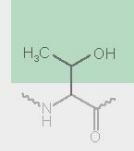
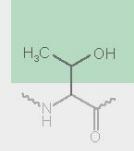


**SNP**

A C G T G T C	<b>G</b>	G T C T T A	Maternal chrom.
A C G T G T C	<b>A</b>	G T C T T A	Paternal chrom.
A C G T G T C	<b>G</b>	G T C T T A	Maternal chrom.
A C G T G T C	<b>G</b>	G T C T T A	Paternal chrom.
A C G T G T C	<b>A</b>	G T C T T A	Maternal chrom.
A C G T G T C	<b>A</b>	G T C T T A	Paternal chrom.

# Типы однонуклеотидных замен

- Эффект замены радикально зависит от ее типа

		Point mutations		
		No mutation	Silent	Nonsense
				Missense
				conservative non-conservative
DNA level	TTC	TTT	ATC	TCC
mRNA level	AAG	AAA	UAG	AGG
protein level	Lys	Lys	STOP	Arg
				
				
				basic polar

# Вариативность в геномах

- “Инделы” - инсерции и делеции
- Крупные структурные вариации

## Indel examples

### wild-type sequence

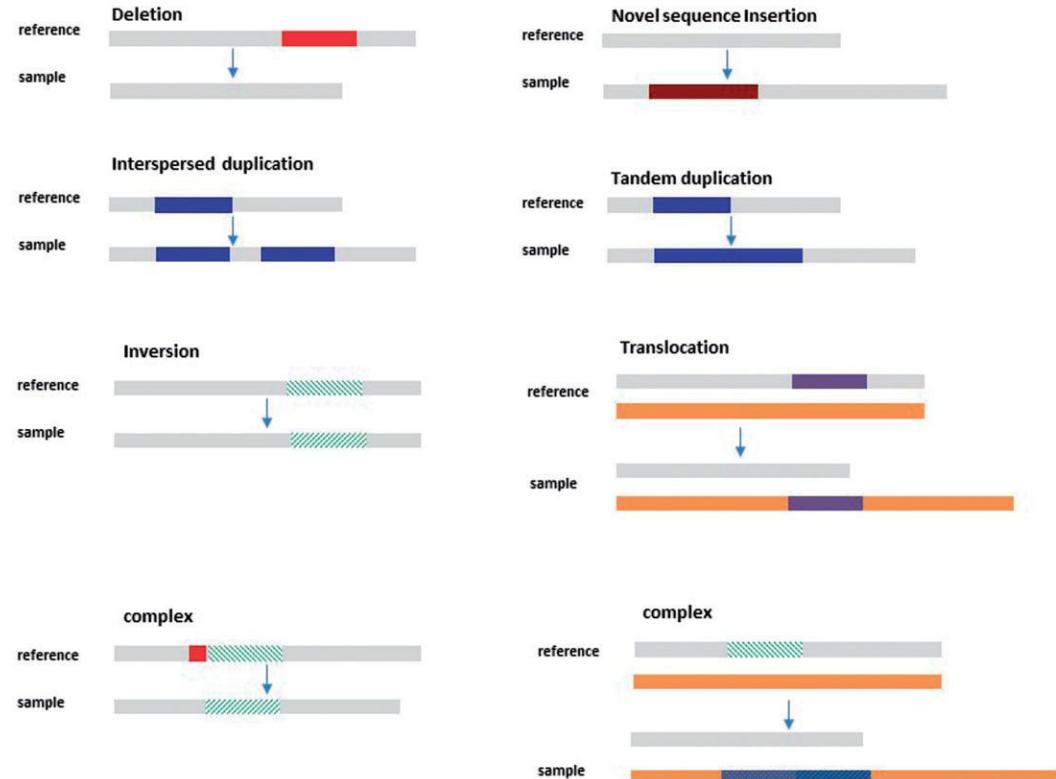
ATCTTCAGCCATAAAAGATGAAGTT

### 3 bp deletion

ATCTTCAGCCAAAGATGAAGTT

### 4 bp insertion (orange)

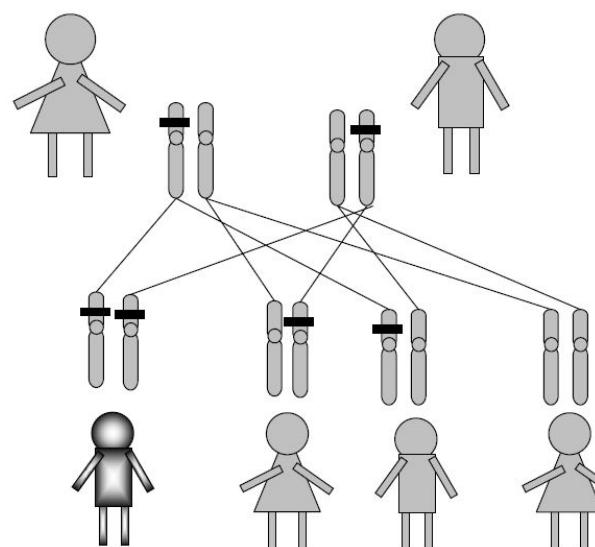
ATCTTCAGCCATATGTGAAAGATGAAGTT



# Медицинский NGS

- Около 7000 редких заболеваний, 80% генетических
- Общая частота - около 10%
- Мендельевские болезни - от 2 до 5%

recurrence risk: 25%, horizontal transmission,  
both sex are similarly affected



■ = carrier  
■ = symptoms  
A = dominant allele (wt)  
a = recessive allele (mut)

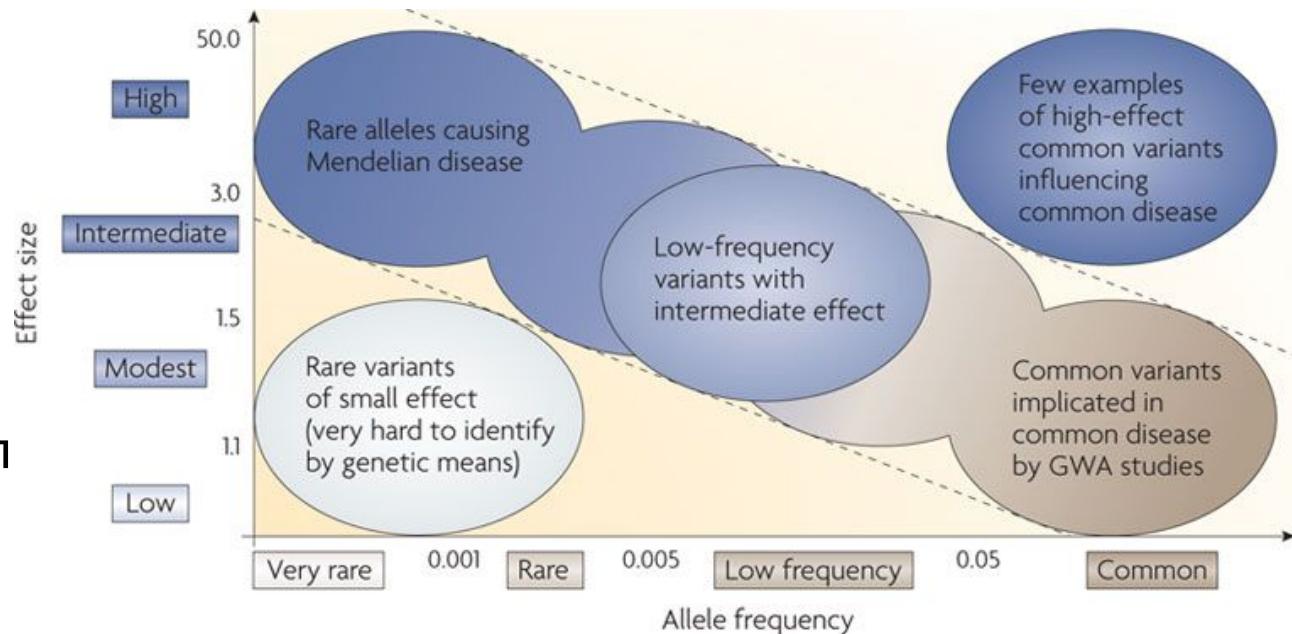
	A	a
A	AA	Aa
a	Aa	aa

genotype: 1:2:1  
phenotype: 1:3

Cornelia Schubert, Institute of Human Genetics, Goettingen

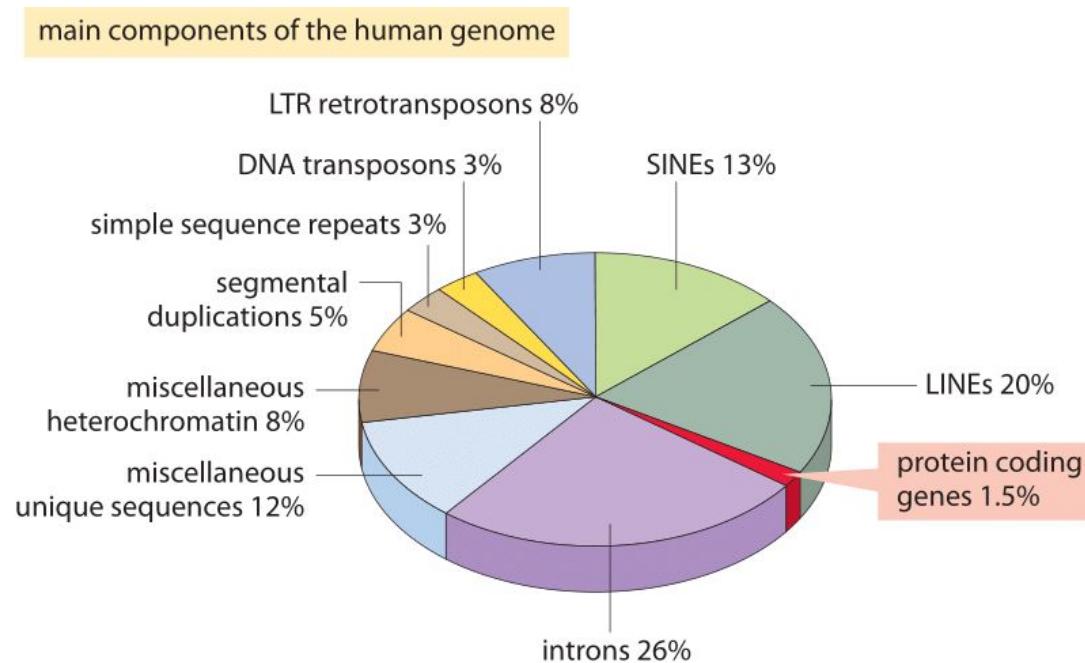
# Частота и влияние вариантов

- Большинство сильно-патогенных вариантов - редкие
- Эффект частых вариантов, как правило, очень мал



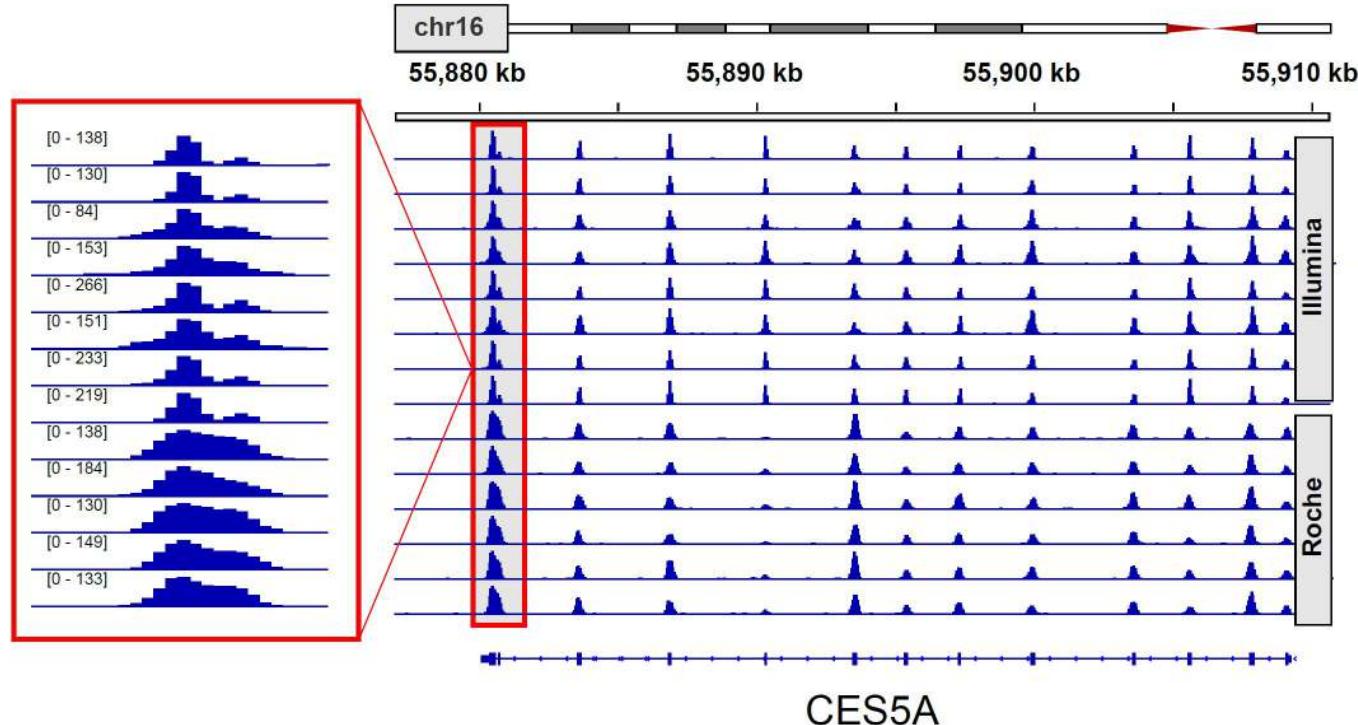
# Гены в геноме человека

- кодирующие последовательности составляют всего 1.5% от всего генома (45 Мб из 3 Гб)



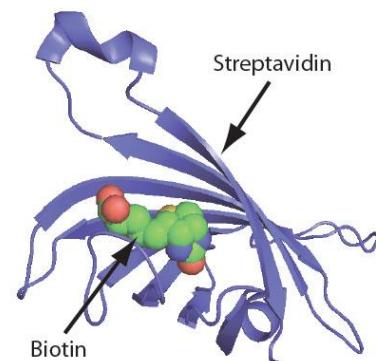
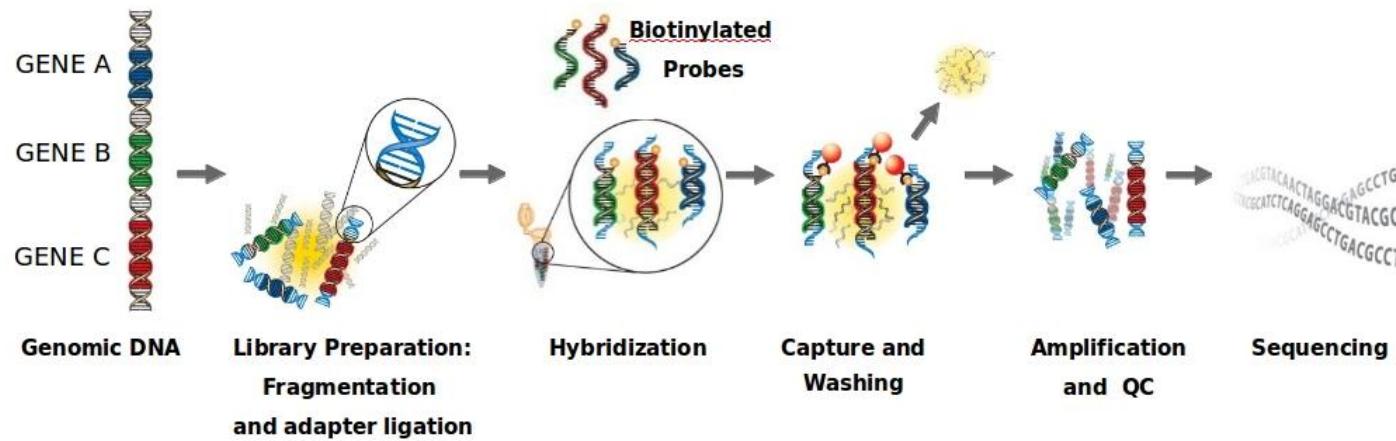
# Что такое экзом?

- Фокус только на кодирующих регионах (экзонах)



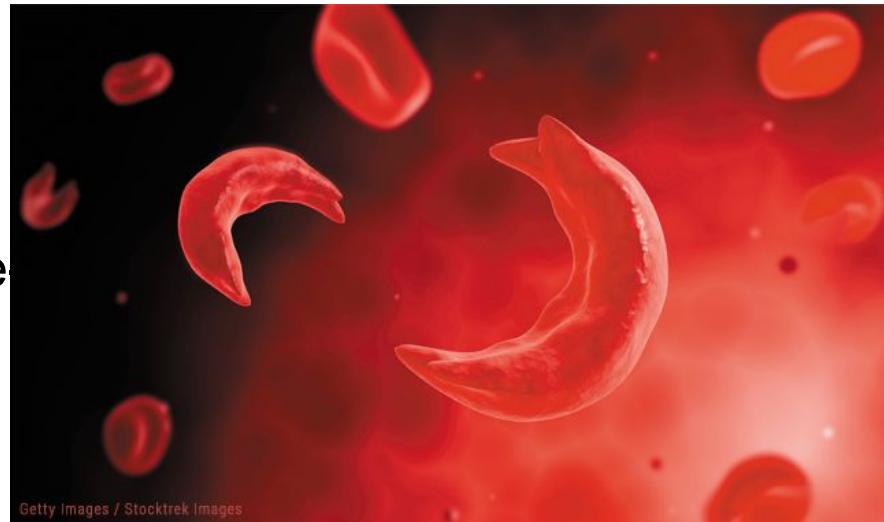
# Обогащение ДНК

- Выделяются только нужные нам фрагменты
- Набор фрагментов может быть практически любым!

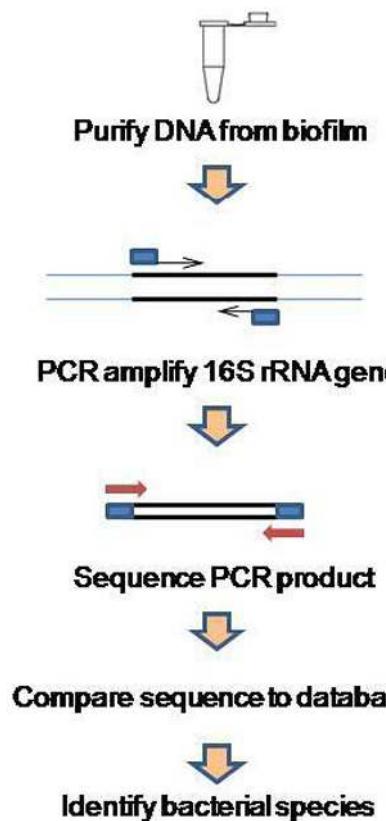


# Ресеквенирование ДНК

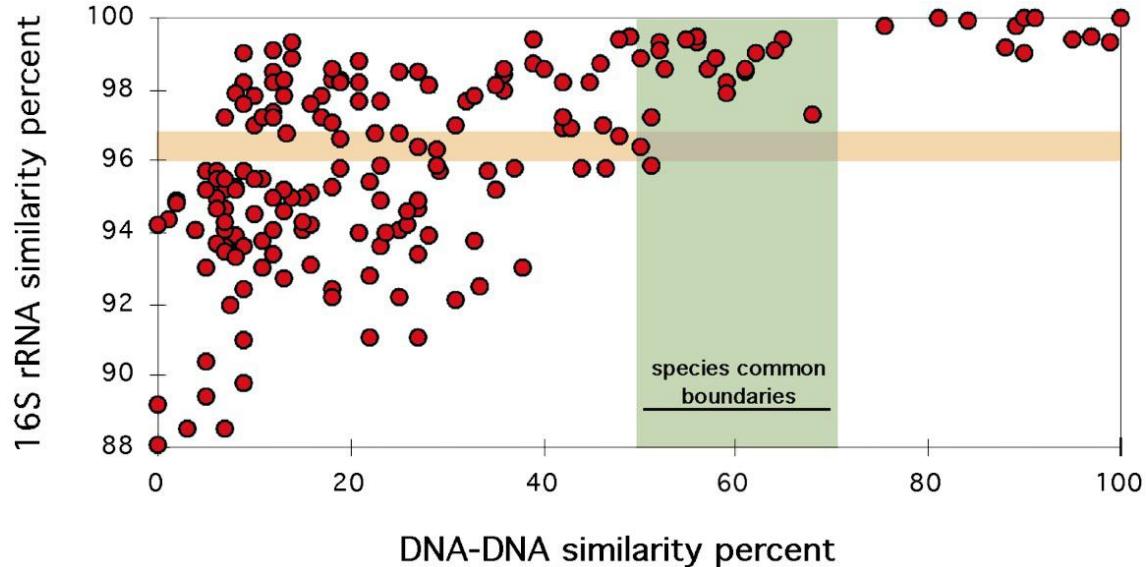
- Таргетная панель (1-500 генов)
- Клинический экзом (5000 важных генов) - CES
- Полный экзом (20000 генов) - WES
- Полный геном - WGS
- Являются необходимостью и факто стандартом при генетических (наследуемых) заболеваниях



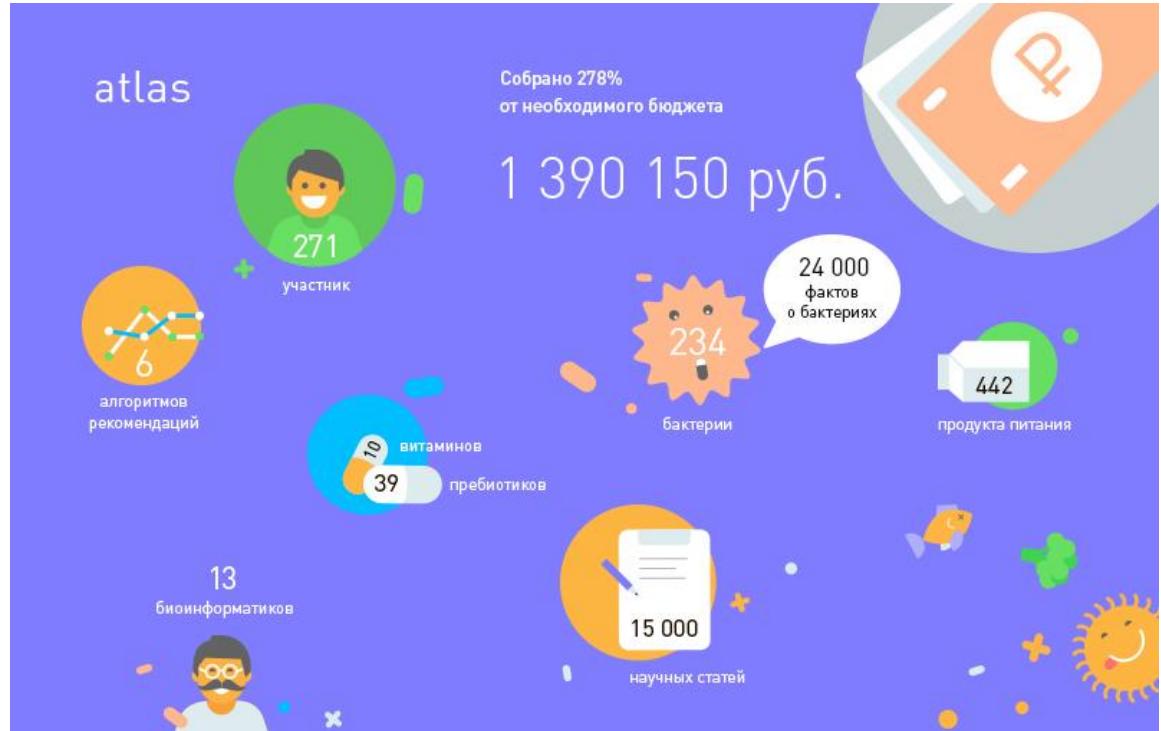
# 16S-секвенирование



- Секвенирование ампликона

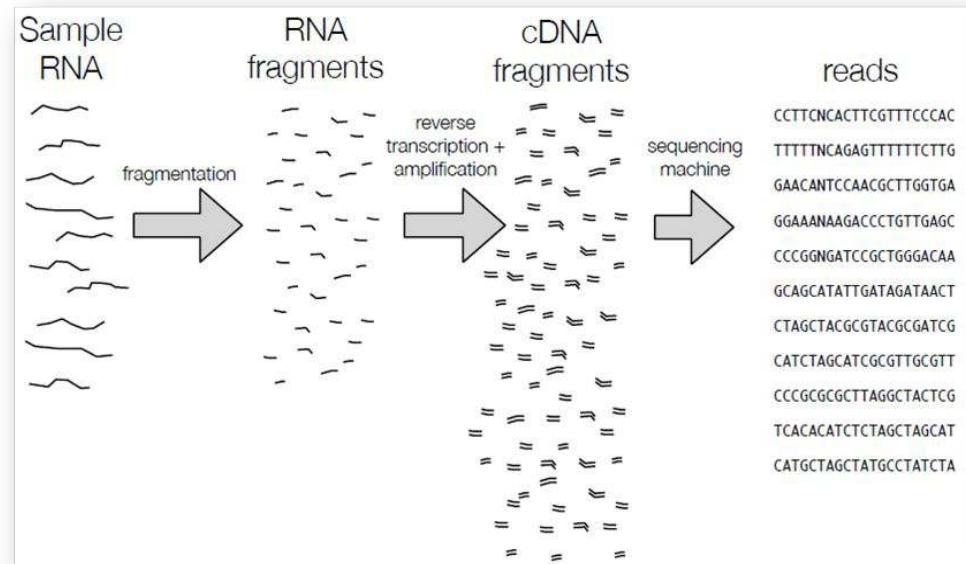


# Oh My Gut!



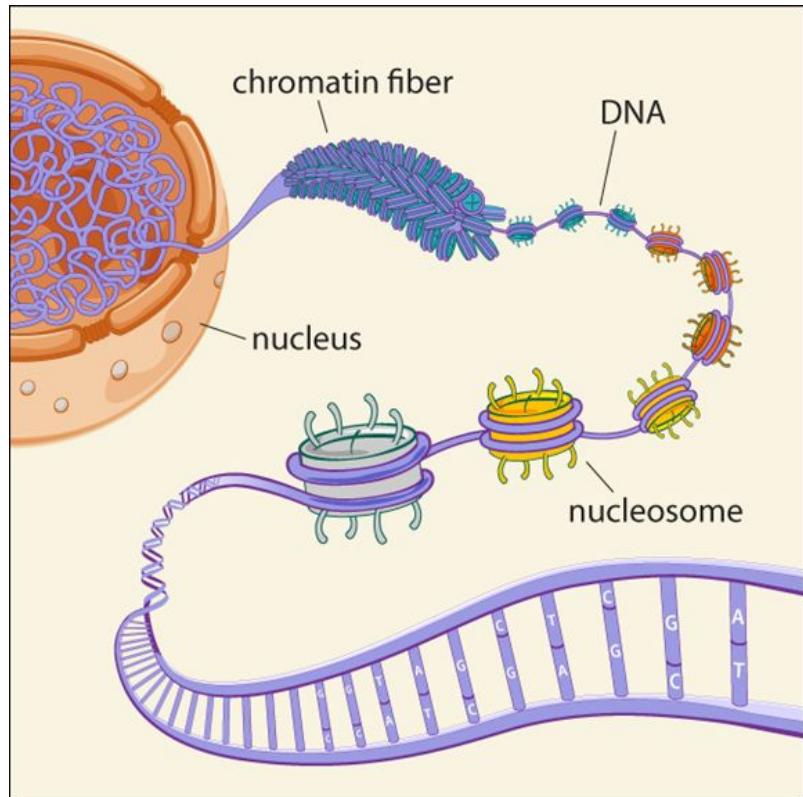
# Секвенирование транскриптома

- De novo - если сборка генома слишком трудна, или нужна разница между тканями
- На хорошо аннотированных организмах - связь с фенотипом или условиями
- Цели:
  - Количественный анализ экспрессии
  - Поиск новых транскриптов



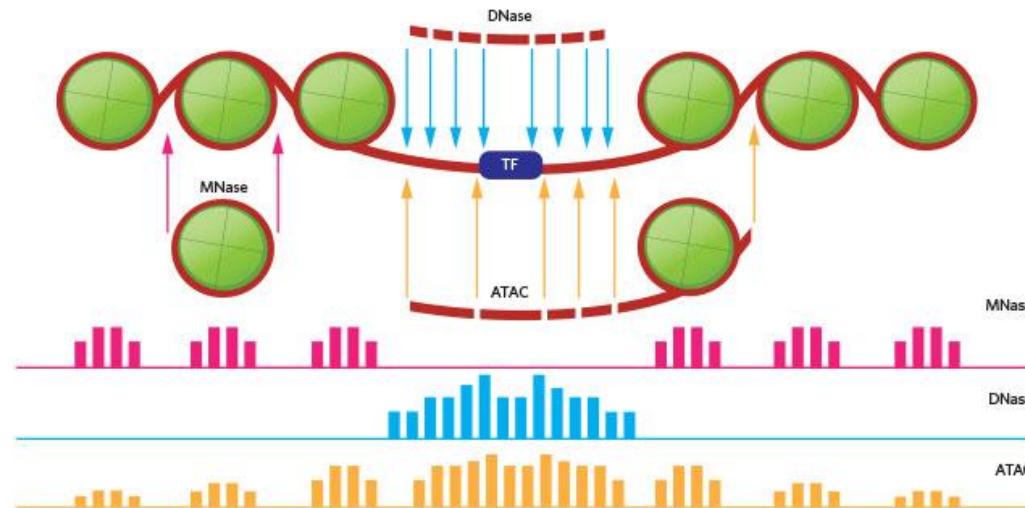
# Хроматин в эукариотах

- Три уровня упаковки:
  - "бусы на нитке" - нуклеосомы
  - 30 нм фибрила
  - Плотно упакованная структура - гетерохроматин
- Нуклеосомы содержат гистонные протеины
- Длина ДНК hs - около 2 м, размер клетки - порядка микрометров



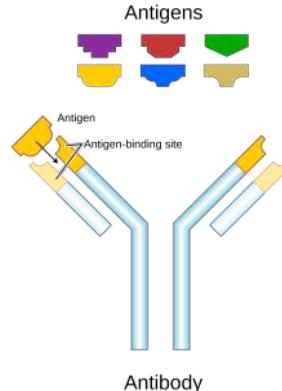
# Genome footprinting

- Смысл: поймать наиболее “развернутые” участки хроматина
- DNase-seq, FAIRE-seq, ATAC-seq



# ChIP-Seq

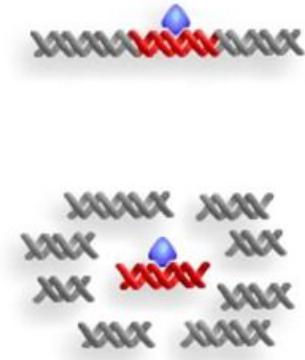
- Смысл: поймать участки ДНК, связанные с определенным протеином
- Делается с помощью селективных антител



Add formaldehyde to **chemically bind** proteins and DNA



**Break** DNA into small fragments (usually ultrasound)



Precipitate fragments with the target protein using **selective antibody**



Amplify the DNA using PCR and **sequence** it



# Куда податься биоинформатику?

- биоинформатика - крайне быстро растущий рынок
- в том числе и в России
- большинство потребителей услуг - в медицине и науке

FRONT ALL RANDOM AGGRESSION FUNNY VIDEOS NEWS PICS GAMING WORLDNEWS GIFS MOVIES TODAYILEARNED AWW SHOWERTHOUGHTS MILDLYINTERESTING



## BIOINFORMATICS

COMMENTS

This is an archived post. You won't be able to vote or comment.

question Is bioinformatics a viable career right now? (self/bioinformatics)  
23 submitted 1 year ago by [deleted]

Deciding if i should study this or not. I have a Masters Degree in Biology. I want to find a way to convert it to a useful career in Canada. Potato English

40 comments share save hide report

all 40 comments sorted by: best ▾  
[-] compbioguy 10 points 1 year ago  
Very. My lab trainees gets jobs easily in both academics and industry  
permalink embed save give gold

## SHARE

This special feature is brought to you by the [Science/AAAS Custom Publishing Office](#)

## An Explosion Of Bioinformatics Careers

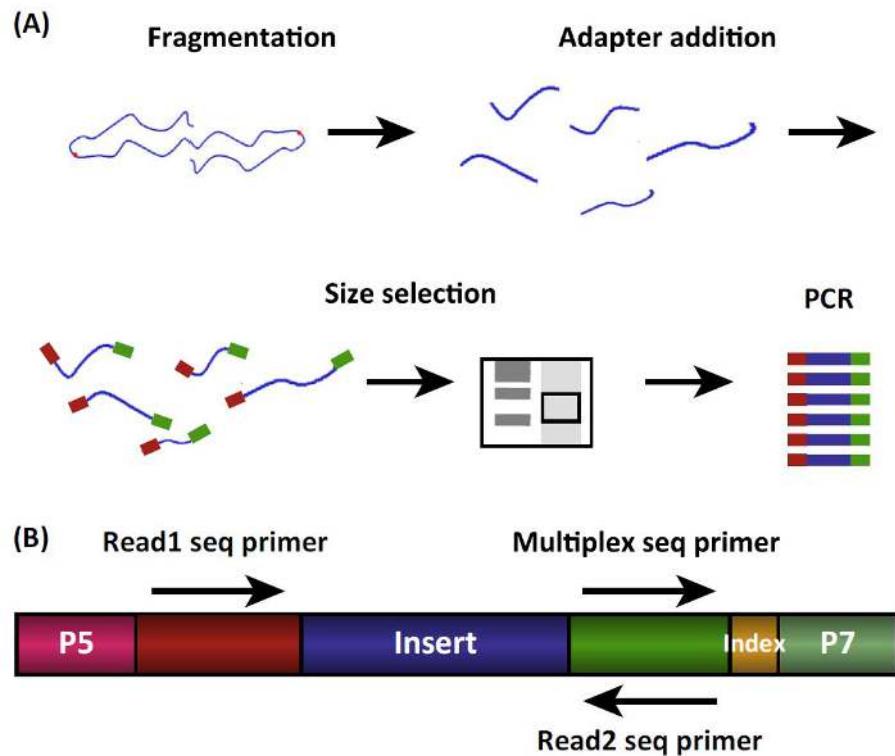
By Alaina G. Levine | Jun. 13, 2014, 2:00 PM



Big data is pouring out of life sciences research, creating ample opportunities for scientists with computer science expertise.

# Типичный NGS-эксперимент

- фрагментация ДНК
- лигация адаптеров
- выбор фрагментов нужного размера
- амплификация
- секвенирование



# Биоинформатика IRL

- Практически работа биоинформатика ~ NGS
- Работа с NGS ~ quality control (QC)
- Плюс, грамотное применение подходящих пайплайнов

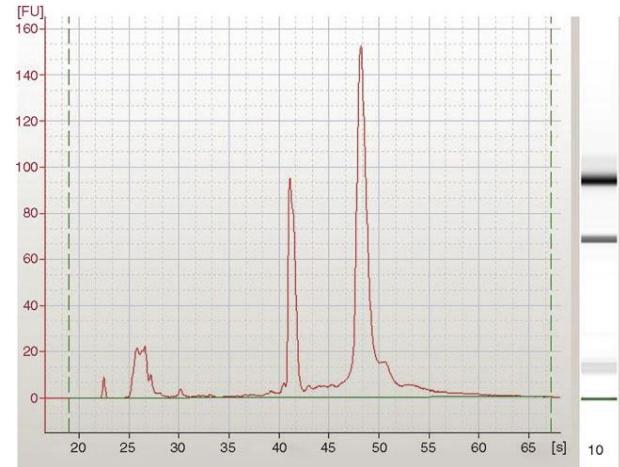
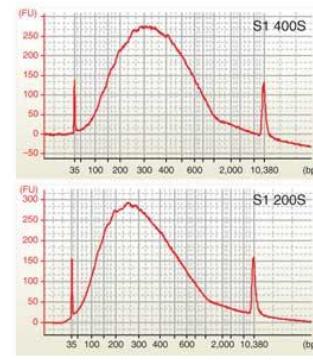
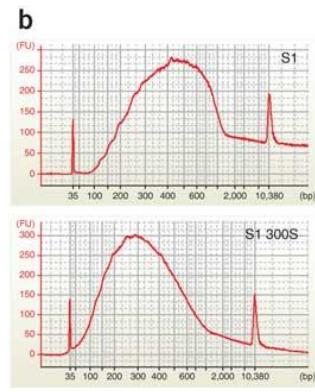
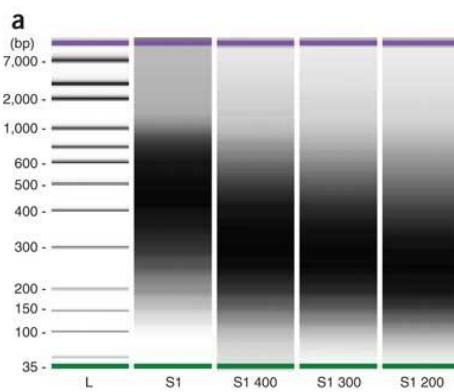


# Постстадийный QC

- QC эксперимента и приготовления библиотеки
  - Биологические наблюдения (живые клетки, гомогенность)
  - Качество экстракции ДНК/РНК, ее состояние
  - Распределение фрагментов по размеру - Bioanalyzer
  - PCR на нужные фрагменты (до секвенирования)
- QC секвенатора
  - На Illumina - плотность кластеров, basecalls, состояние прибора
- QC полученных прочтений
- QC после сборки и выравнивания

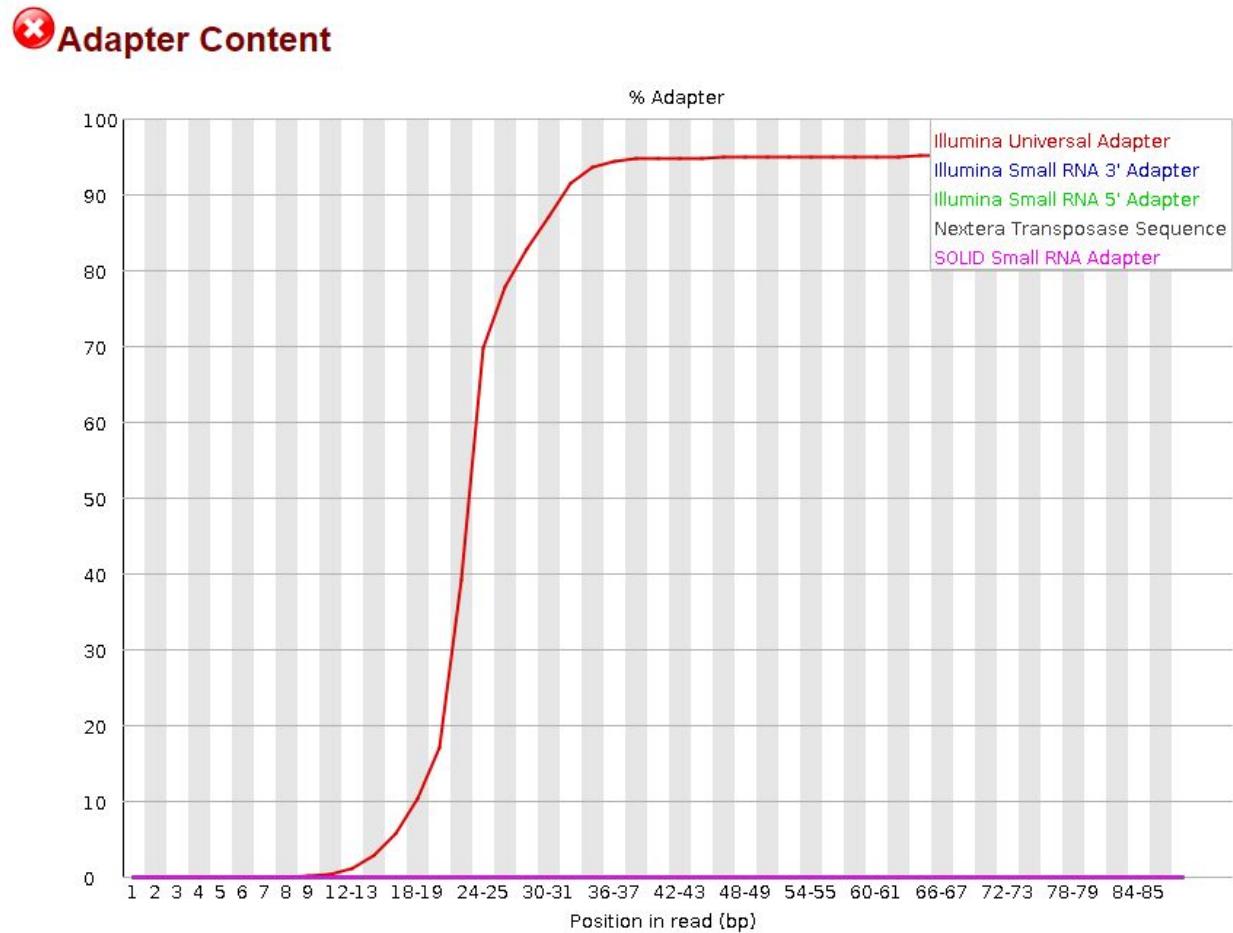
# Типичные ранние проблемы

- РНК или ДНК деградировал
- Чрезмерная/недостаточная соникация или дробление энзимом
- Ошибки в количествах реагентов (плохие реагенты)
- Библиотека “низкой сложности” (мало уникальных фрагментов)



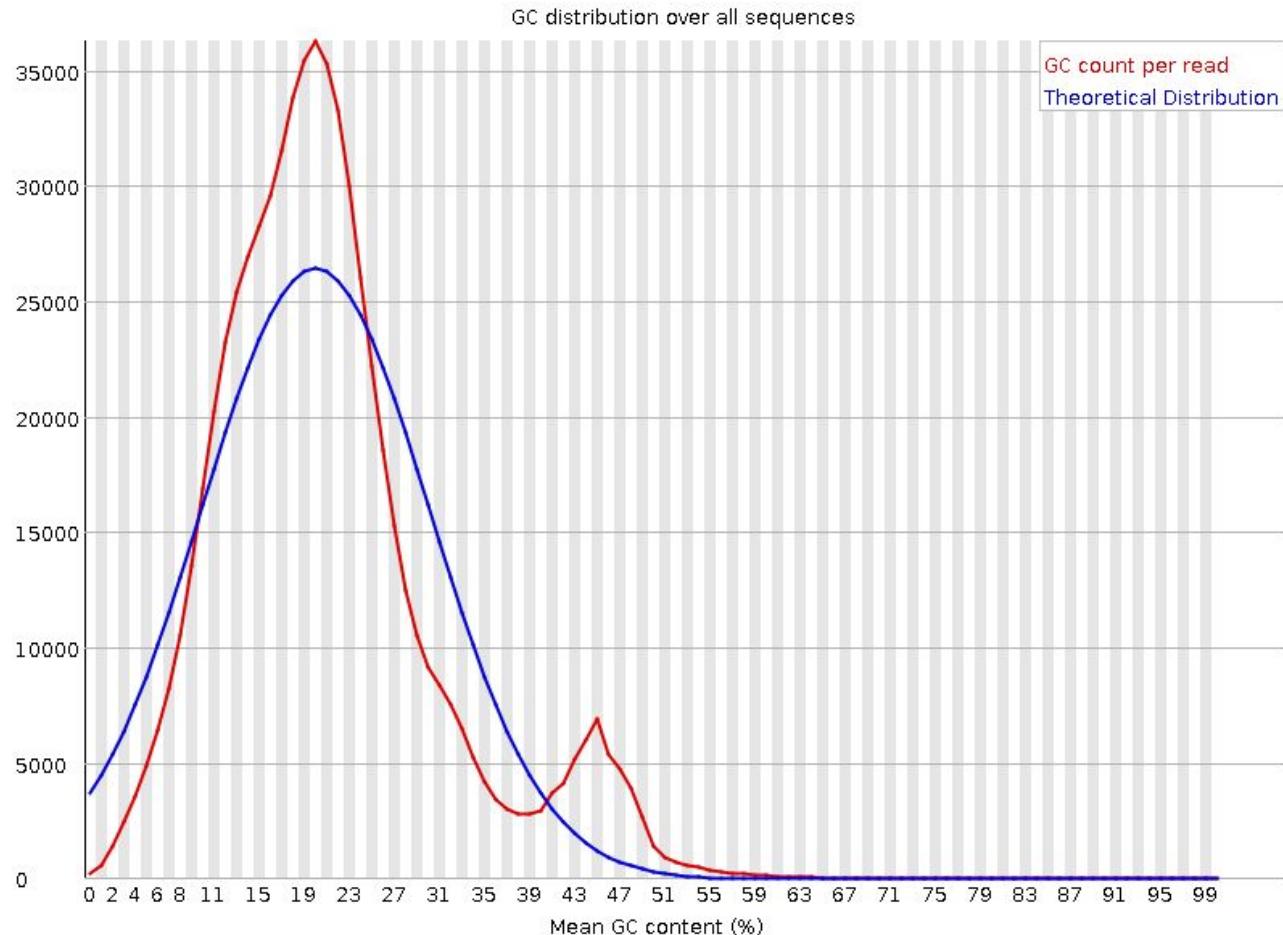
# Пример 1

- Слишком короткие фрагменты
- Секвенирование доходит до конца и идет в адаптер



# Пример 2

- загрязнение  
бактериальным  
геномом



# Примеры жизненных неудач

- QCFail: <https://sequencing.qcfail.com/>



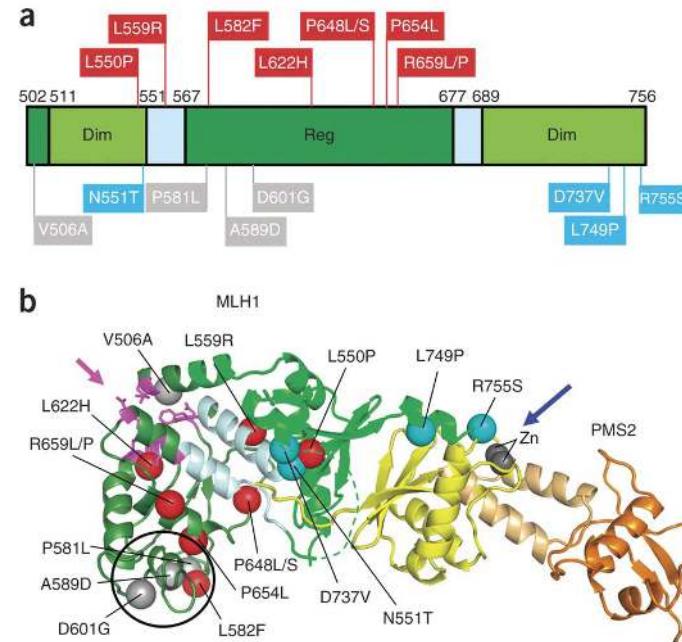
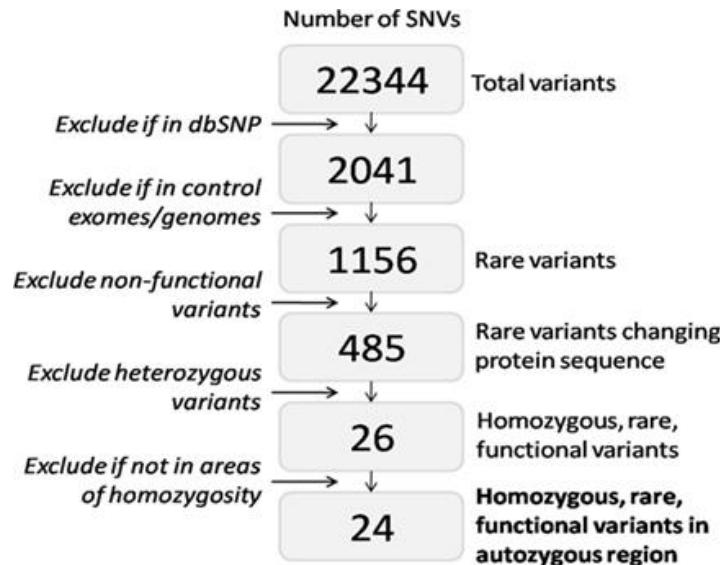
# Непонятные прочтения

- BLAST - понять, какой биологический вид
- Методика - чтобы понять, откуда они взялись

The screenshot shows the NCBI BLAST search interface. At the top, it displays the NIH logo, U.S. National Library of Medicine, and NCBI National Center for Biotechnology Information. Below this, the title "BLAST® > blastn suite" is shown. The main form is titled "Enter Query Sequence" and contains a text input field for "Enter accession number(s), gi(s), or FASTA sequence(s)". To the right of this field are "Clear" and "Query subrange" buttons, along with "From" and "To" input fields. Below the query input, there are fields for "Or, upload file" (with a "Choose File" button showing "No file chosen") and "Job Title" (with a text input field). A checkbox for "Align two or more sequences" is also present. The next section, "Choose Search Set", includes a "Database" dropdown set to "Human genomic + transcript" (with options for "Mouse genomic + transcript" and "Others (nr etc.)" and "Nucleotide collection (nr/nt)" selected), an "Organism" dropdown for "Optional", and an "Exclude" section with checkboxes for "Models (XM/XP)", "Uncultured/environmental sample sequences", and "Sequences from type material". An "Entrez Query" field is also available. The "Program Selection" section allows optimization for "Highly similar sequences (megablast)" (selected), "More dissimilar sequences (discontiguous megablast)", or "Somewhat similar sequences (blastn)". A "Choose a BLAST algorithm" link is provided. At the bottom, a large blue "BLAST" button is followed by a link to "Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences)". There is also a checkbox for "Show results in a new window" and a link to "+Algorithm parameters".

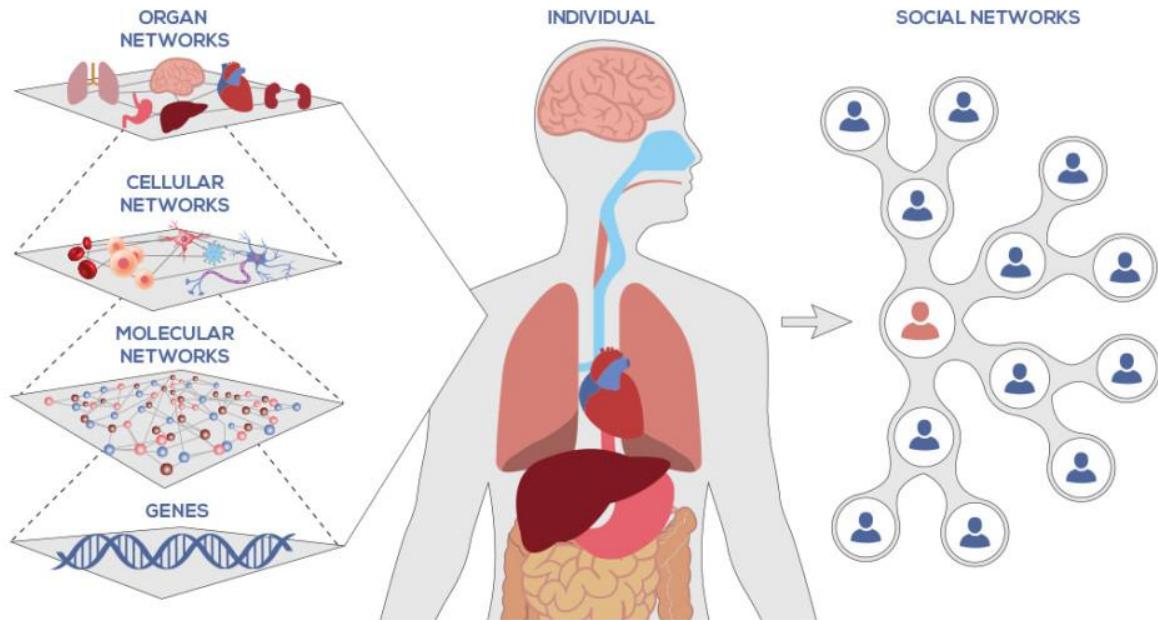
# Чем занимается мед-биоинформатик?

- общение с врачом-генетиком
- контроль качества образцов, анализ покрытия, кандидаты
- анализ и приоритезация мутаций



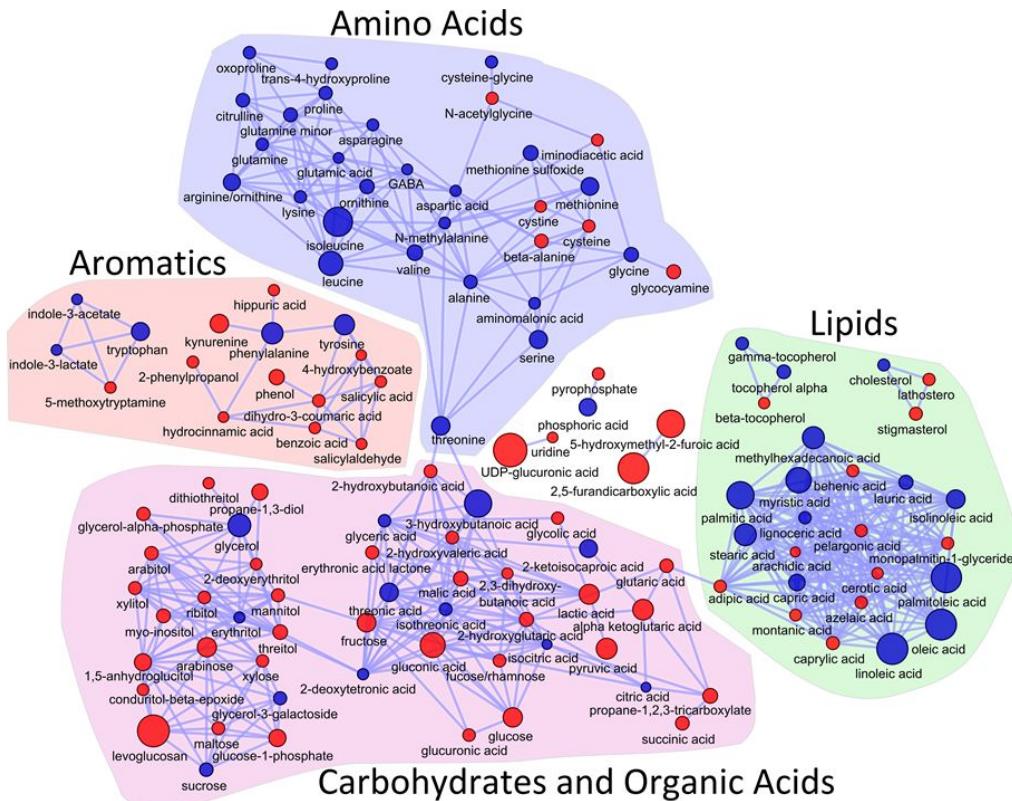
# Системная биология

- попытки объединить разрозненные данные в целостную картину
- попытки построить более точные математические модели



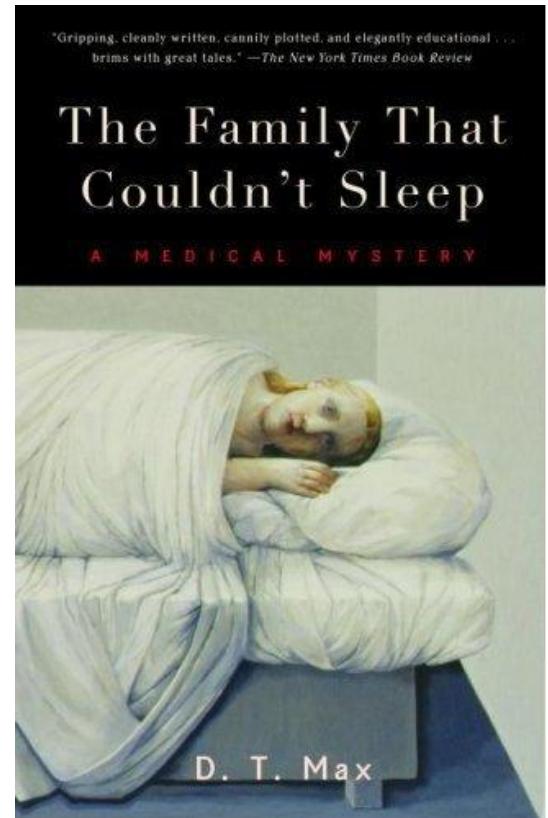
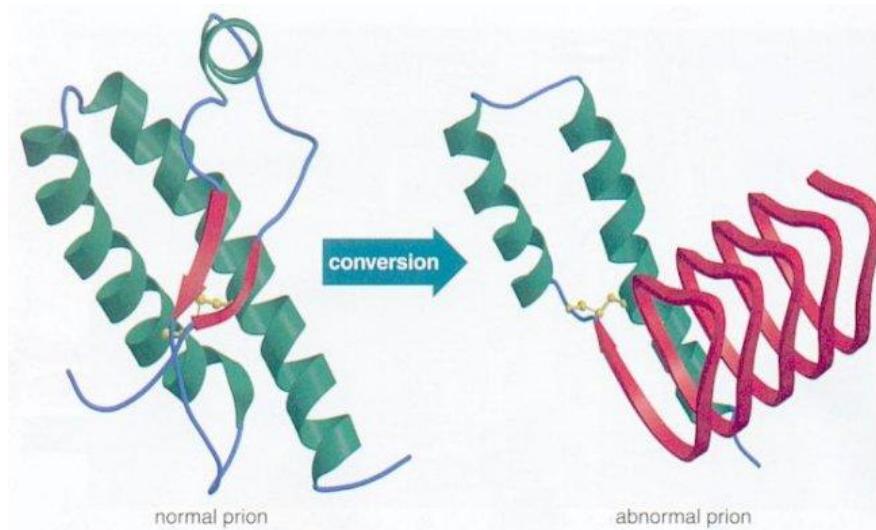
# Интеграция омиксных технологий

- геномика
  - транскриптомика
  - протеомика
  - метаболомика
    - лиpidомика
    - гликомика



# Fatal Familial Insomnia (FFI)

- аутосомно-доминантное прионное заболевание
- известно всего 25 семей
- знаменитая семья из Венето, страдающая от заболевания уже более 200 лет!



# История другой семьи

- Мать умерла от деменции в 52 года
- Дочь оказалась носителем мутации
- Оба бросили свои карьеры (в 2013-м году) и были приняты в аспирантуру Broad Institute
- cureFFI.org

CureFFI.org

About Archives Contact

## Does this mean I'll definitely get the disease?

Jan 20, 2016 • ericminikel • Cambridge, MA

I am excited to announce today the publication of our new study, "Quantifying prion disease penetrance using large population control cohorts" [full text, PDF, perspective piece by Robert Green]. I wanted to write a blog post to announce it and to offer up the comments section (bottom of this post) as a place for post-publication peer review. I also wanted to take this opportunity to tell the history of this study, a profoundly personal endeavor that has grown and taken shape over the entire four year course of my re-training as a scientist.



# На данный момент

- Организовали небольшую лабораторию в Broad Institute
- Собрали деньги на клинические испытания
- Успели еще очень много всякого!



Спасибо за внимание!