

Введение в анализ данных NGS

Анастасия Жарикова

10 ноября 2020

azharikova89@gmail.com

Для чего нужно секвенирование

- Эволюция
- Филогения
- Клиника
- Метагеномика
- Анализ транскриптомов
- Single cell (различные приложения)
-

ЭВОЛЮЦИЯ

Доместикация риса

1 MSGSSADPSP SASTAGAAVS PLALLRAHGH GHGHLTATPP SGATGPAPPP
51 PSPASGSAPR DYRKGNWTLH ETLILITANR LDDDRRAGVG GAAAGGGGAG
101 SPPTPRSAEQ RWKVENYCW KNGCLRSQNG CNDKWDNLLR DYKKVRDYES
151 RVA AAAAATGG AAAANSAPLP SYWTMERHER KDCNLPTNLA PEVDALSEV
201 LSRRAARRGG ATIAPTPPPP PLALPLPPPP PPSPPKPLVA QQQH HHHGHH
251 HHPPPPQPPP SSLQLPPAVV APPPASVSAE EEMSGSSESG EEEEGSGGEP
301 EAKRRRLSRL GSSVRSATV VARTLVACEE KRERRHRELL QLEERRLRLR
351 EERTEVRRQG FAGLIAAVNS LSSAIHALVS DHRSGDSSGR

sh4

Li et al., Science, 2006

Дикий рис

ААG

Лизин

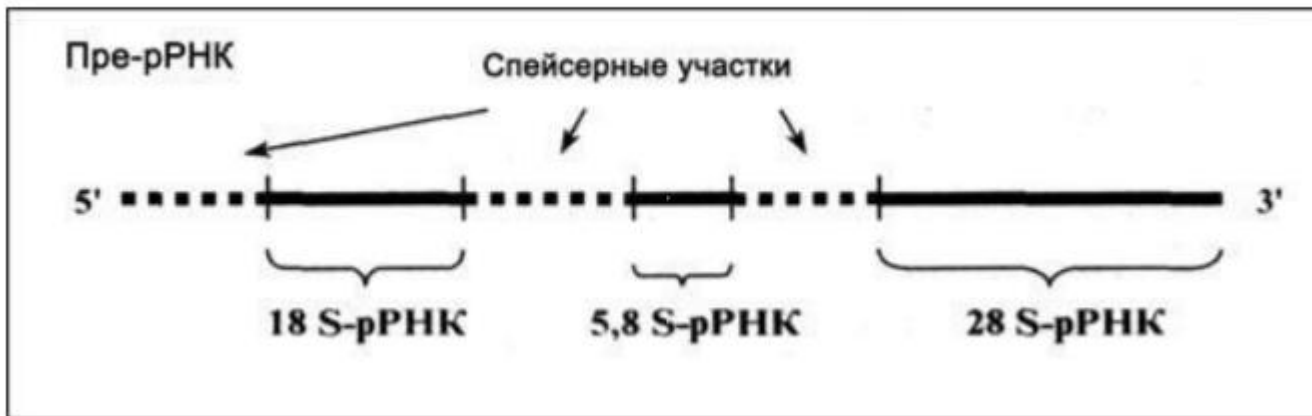
Культурный рис

ААТ

Аспарагин

Филогения

Транскрибируемые спейсеры



Спейсерные последовательности наиболее переменные с точки зрения эволюционной консервативности.

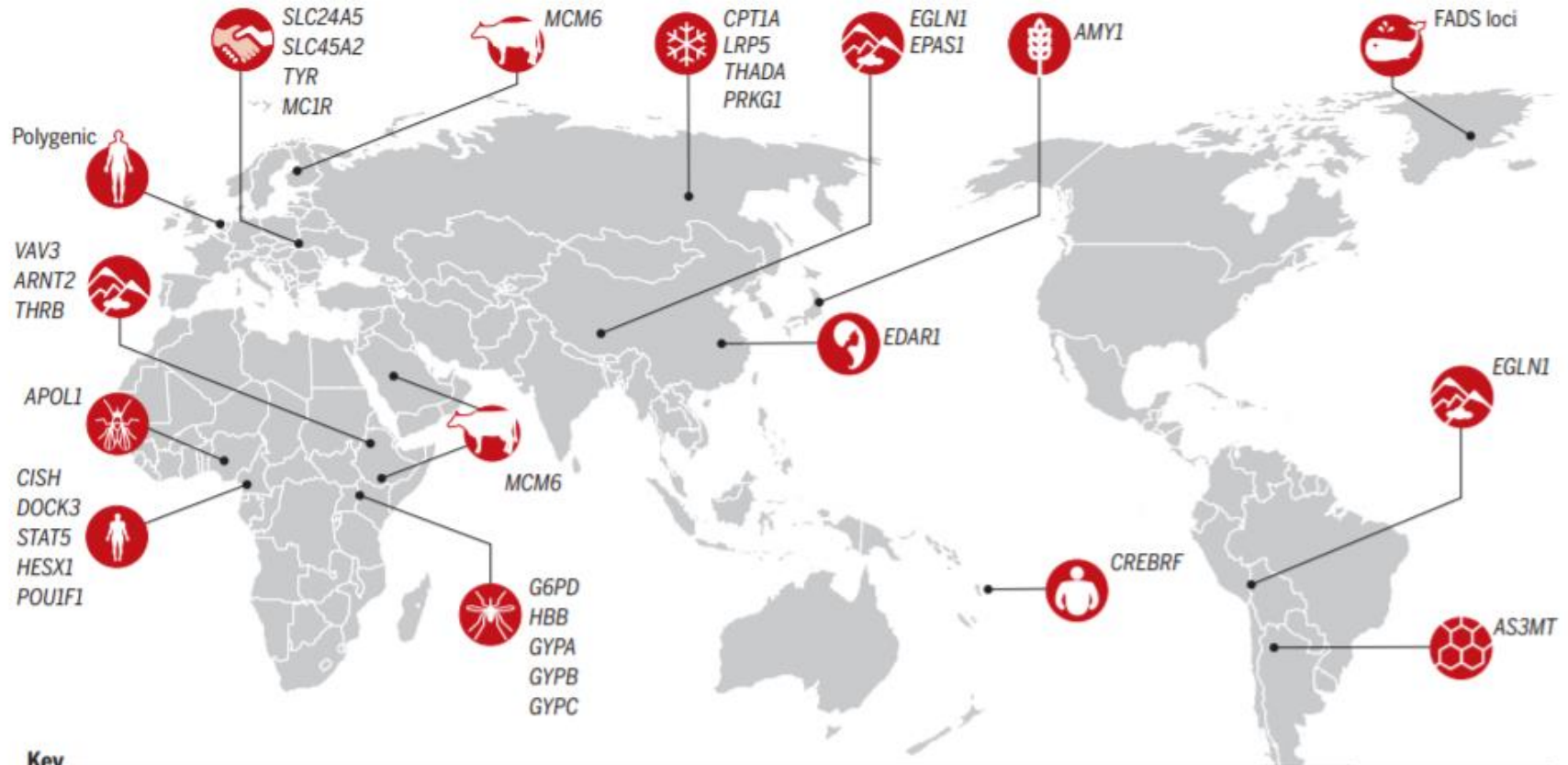
Секвенирование и анализ транскрибируемый спейсеров используется для изучения видового разнообразия и классификации близкородственных организмов

Популяционные и клинические исследования

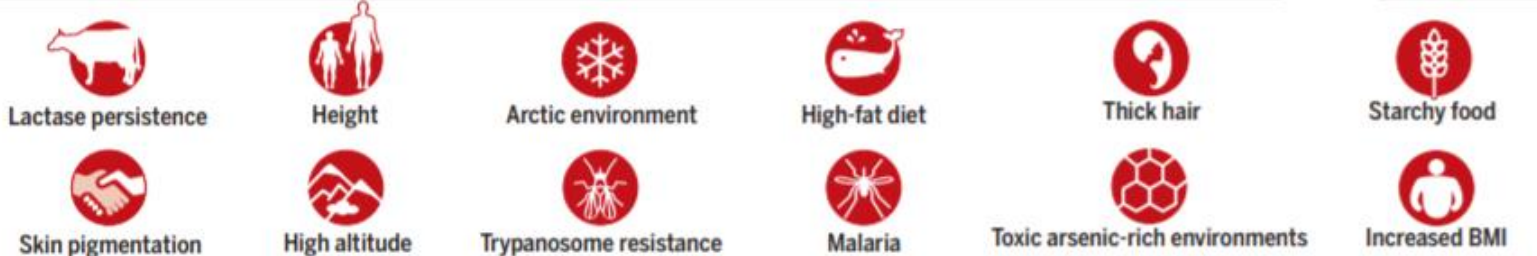
- 1000 геномов
 - Частоты snp в популяциях
- GWAS
 - Поиск полиморфизмов, ассоциированных с болезнями:
 - моногенные (муковисцидоз, ген CFTR)
 - полигенные (ишемическая болезнь сердца, шизофрения)
- Фармакогенетика
 - Индивидуальное лечение
 - Варфарин – предотвращает образование тромбов. Генетические факторы определяют до 53-54 % вариабельности дозы. Гены CYP2C9, CYP4F2, VKORC1.

Going global by adapting local: A review of recent human adaptation

Shaohua Fan,^{1*} Matthew E. B. Hansen,^{1*} Yancy Lo,^{1,2*} Sarah A. Tishkoff^{1,3†}



Key



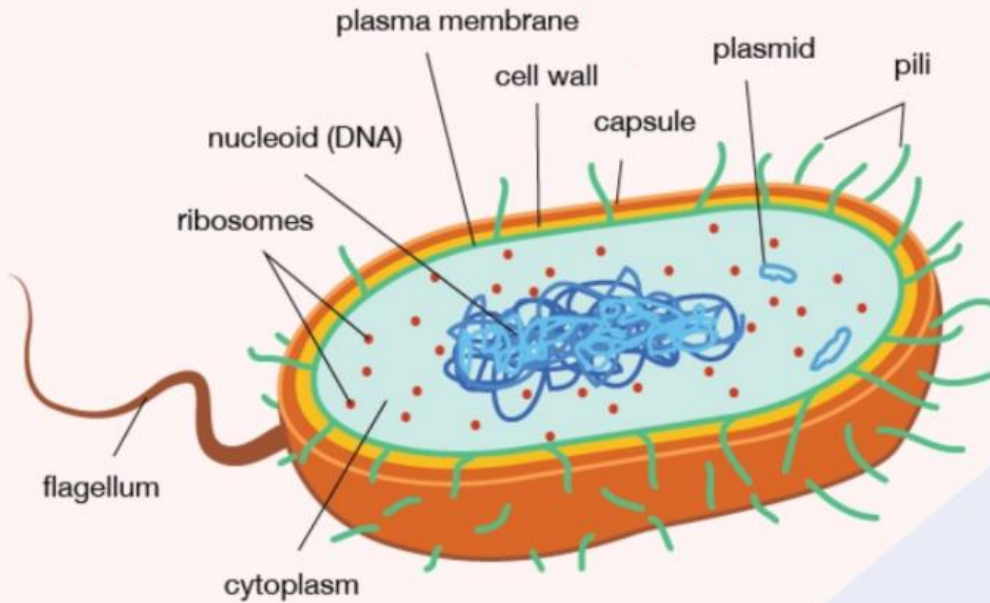
Что бывает

- DNA-seq
- RNA-seq
- Chip-seq
- HiC
- ATAC-seq
- DNase-seq
- eClip
- GRO-seq
- CRISPR-seq
- Ribo-seq
-

Что бывает

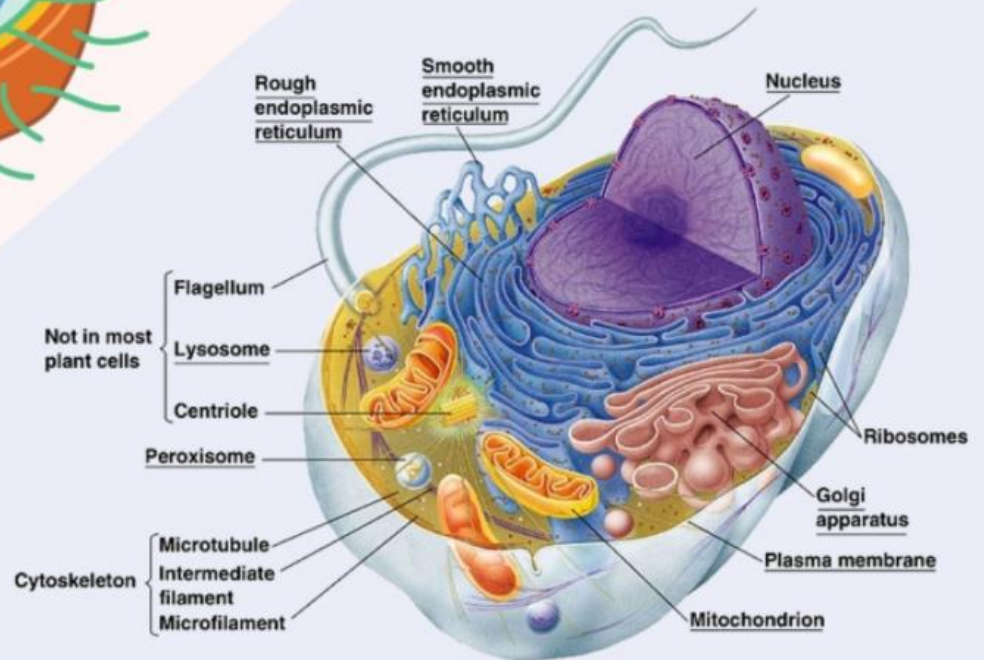
- **DNA-seq**
- RNA-seq
- Chip-seq
- HiC
- ATAC-seq
- DNase-seq
- eClip
- GRO-seq
- CRISPR-seq
- Ribo-seq
-

Где ДНК?

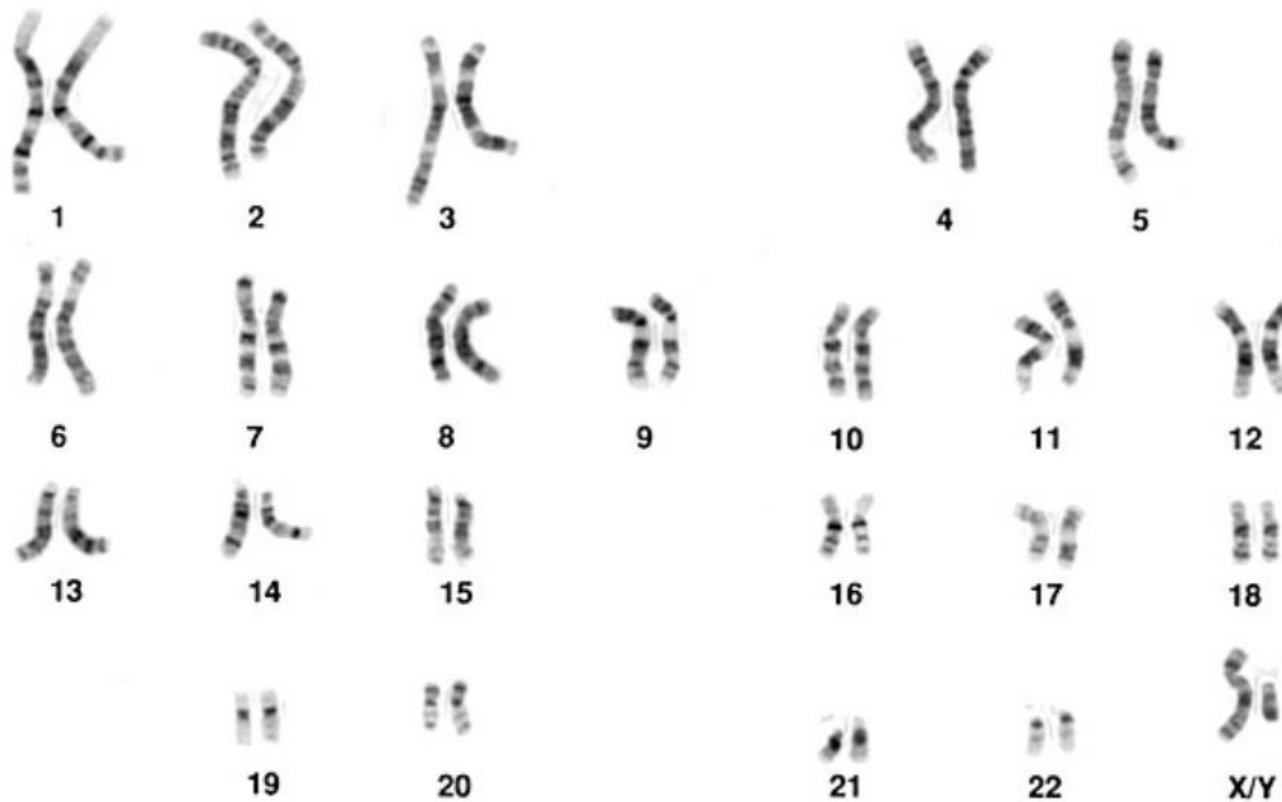


Prokaryotic Cell

Eukaryotic Cell



У человека 23 пары хромосом Много или мало?



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

Число хромосом у разных видов



Гиббоны - 44



Макака - 42



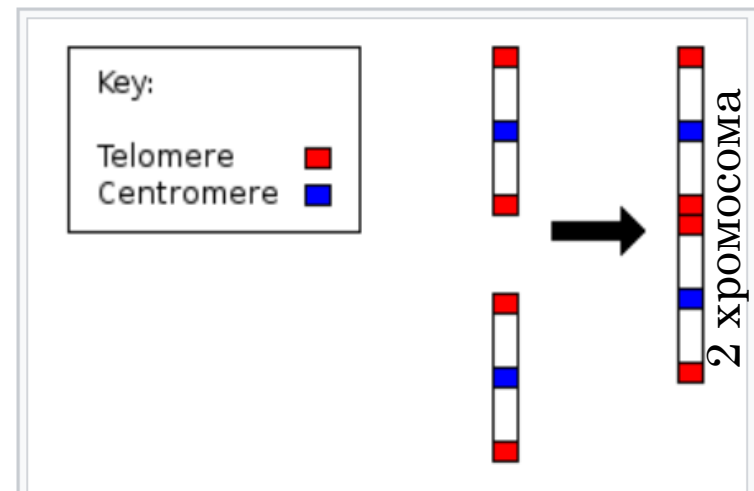
Капуцин - 54



48



46



После слияния двух хромосом остаются характерные следы: остатки теломер и рудиментарная центромера

Число хромосом у разных видов

Муравей (*Myrmica pilosula*) – 2

Плодовая мушка - 8

Арабидопсис – 10

Голубь – 16

Кошка – 38

Лиса - 34

Мышь - 40

Собака – 78

Утка – 80

Сазан - 104

Корова – 120

Рак (*Cambarus clarkii*) – 200

Хвощ – 216

Краб - 254

Бабочка – 380



Число хромосом у разных видов

| Подроды (в скобках) и виды муравьев | Число хромосом | |
|--|----------------|------------|
| | 2n (самцы) | 2n (самки) |
| Род <i>Formica</i> | | |
| <i>F. (Formica) aquilonia</i> - северный лесной муравей | 26 | 52 |
| <i>F. (Serviformica) cinerea</i> - серый лесной муравей | 27 | 54 |
| <i>F. (Serviformica) cunicularia</i> - пряткий муравей | 27 | 54 |
| <i>F. (Coptoformica) exsecta</i> - тонкоголовый муравей | 26 | 52 |
| <i>F. (Serviformica) fusca</i> - бурый лесной муравей | 27 | 54 |
| <i>F. (Formica) lugubris</i> - волосистый лесной муравей | 26 | 52 |
| <i>F. (Serviformica) picea</i> - черный блестящий муравей | 26 | 52 |
| <i>F. (Formica) polyctena</i> - малый лесной муравей | 26 | 52 |
| <i>F. (Formica) pratensis</i> - луговой муравей | 26 | 52 |
| <i>F. (Formica) rufa</i> - рыжий лесной муравей | 26 | 52 |
| <i>F. (Serviformica) rufibarbis</i> - краснощекый муравей | 27 | 54 |
| <i>F. (Raptiformica) sanguinea</i> - кровавый муравей-рабовладелец | 26 | 52 |
| <i>F. (Formica) truncorum</i> - красноголовый муравей | 26 | 52 |
| <i>F. (Serviformica) uralensis</i> - черноголовый муравей | 26 | 52 |
| Род <i>Lasius</i> | | |
| <i>L. (Lasius) alienus</i> - бледноногий муравей | 15 | 30 |
| <i>L. (Cautolasius) flavus</i> - желтый земляной муравей | 15 | 30 |
| <i>L. (Dendrolasius) fuliginosus</i> - пахучий муравей-древоточец | 14 | 28 |
| <i>L. (Lasius) niger</i> - черный садовый муравей | 15 | 30 |

Размер генома у разных видов



Количество белок кодирующих генов у разных видов

Картофель – 39 000

Человек ~ 20 000

Черви – 14 000

Мухи – 12 000

Грибы – 6 000

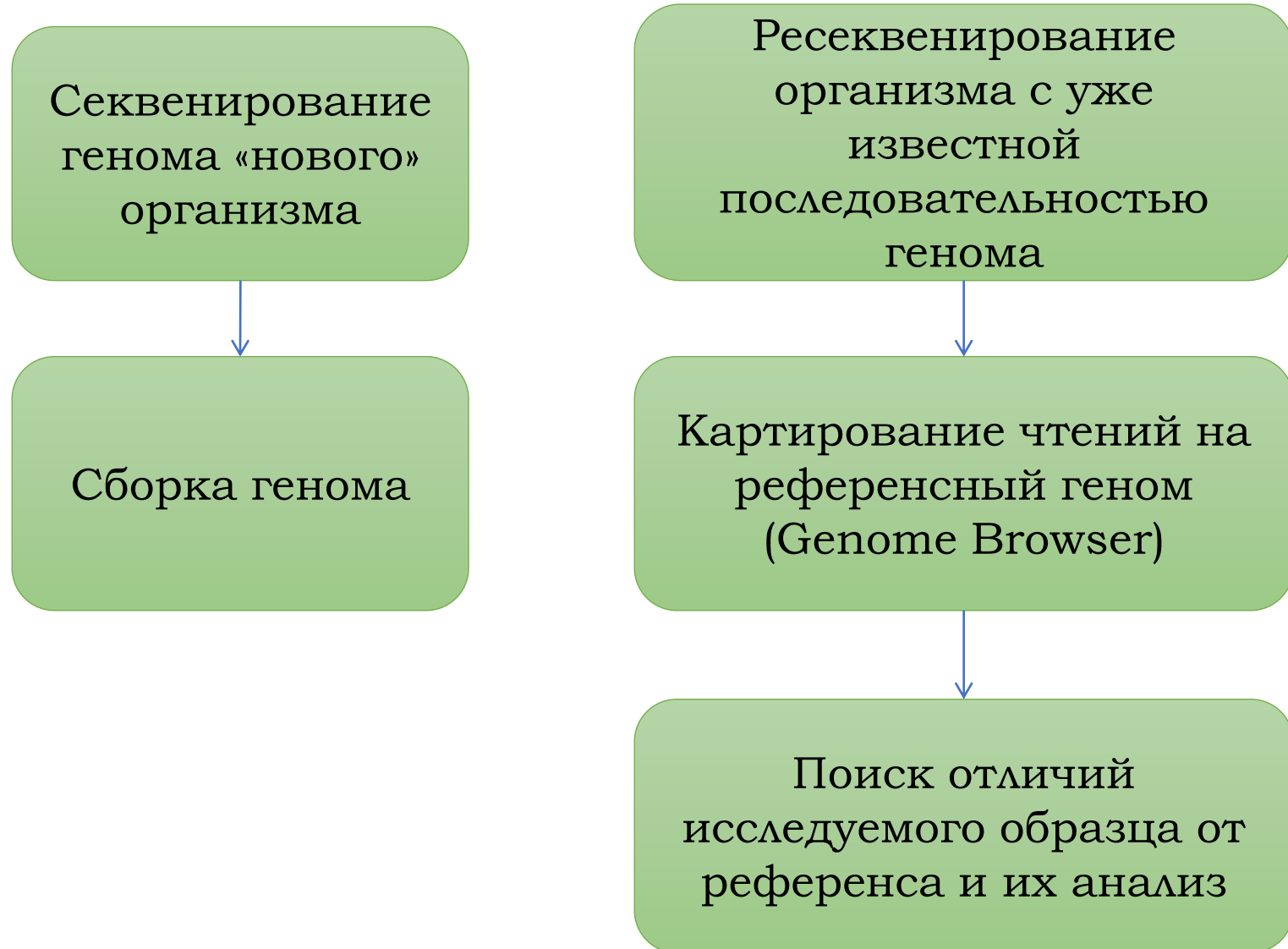
Бактерии – 2 000 – 4 000

Микоплазмы - 500

Вирус гриппа – 12

Какие еще гены бывают?

Секвенирование ДНК бывает ...



Возможности ресеквенирования

Можно ресеквенировать:

- ПОЛНЫЙ ГЕНОМ
- ЭКЗОМ (кодирующую часть генома)
- ОТДЕЛЬНЫЕ ТАРГЕТНЫЕ ГЕНЫ ИЛИ ОБЛАСТИ

!!!Выбор в зависимости от бюджета и целей исследования!!!

Экзомное ресеквенирование

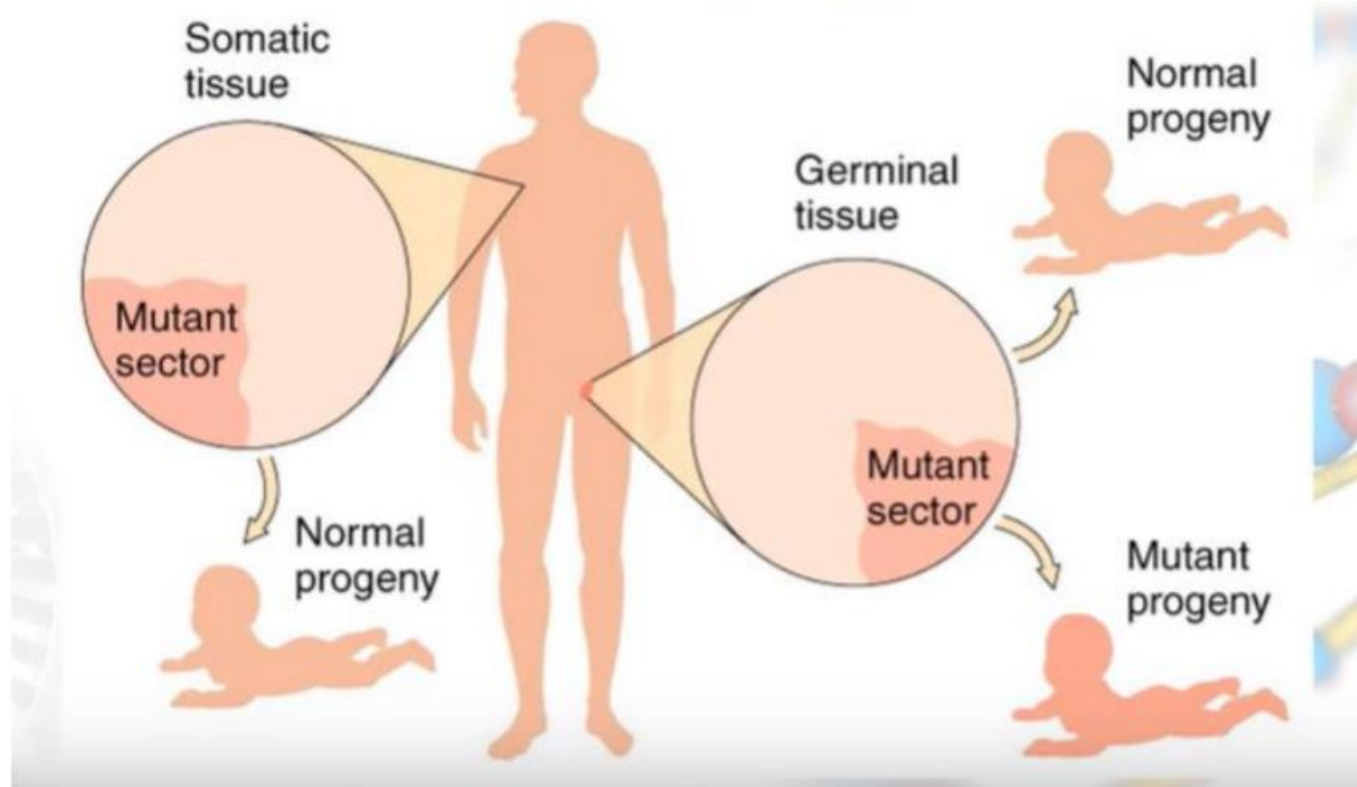
- «Плюсы»

- Небольшой объем кодирующих последовательностей – ниже цена
- Кодирующие последовательности лучше изучены
- Большое число болезнетворных мутаций находится в кодирующей последовательности (особенно менделевские заболевания)

- «Минусы»

- Нет информации о некодирующих участках
- Неравномерность покрытия экзонов

Соматические и герминальные мутации



<https://www.youtube.com/watch?v=3MiHSeH6yLg>

Какие бывают мутации

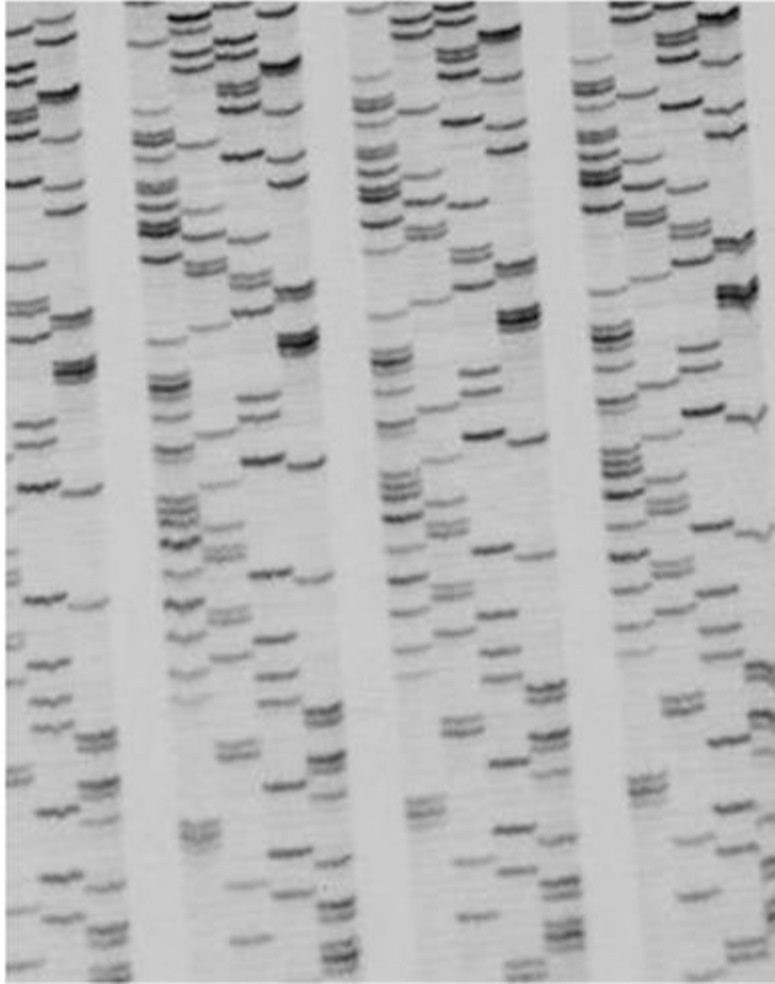
SNV – однонуклеотидные варианты, т.е. изменение одного нуклеотида

Короткие вставки и делеции (~ до 50 п.н.)

Структурные варианты: инверсии и транслокации; CNV

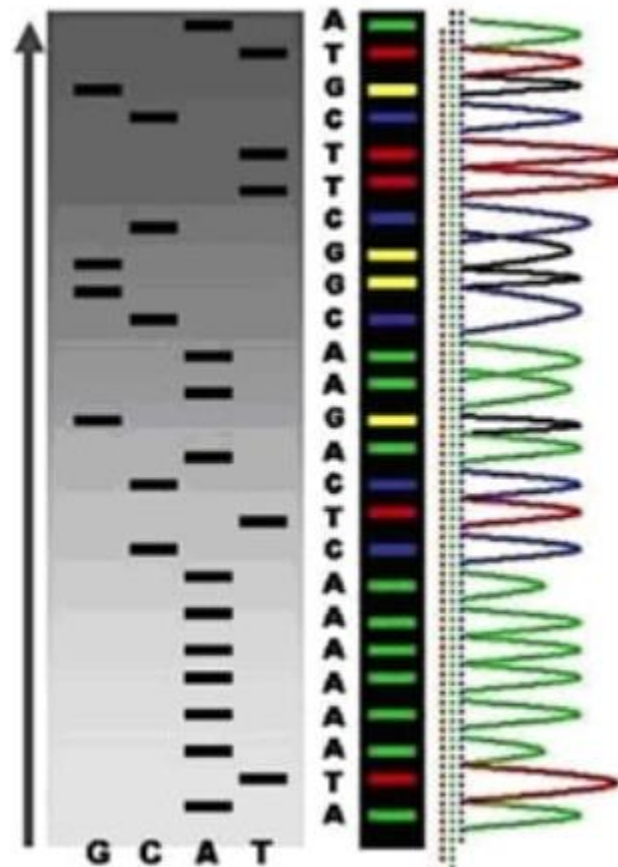
Анеуплоидии: нульсомии, моносомии, трисомии, полисомии

Полиплоидизация



Секвенирование ДНК

Метод «терминаторов»



~ 1000 п.н.

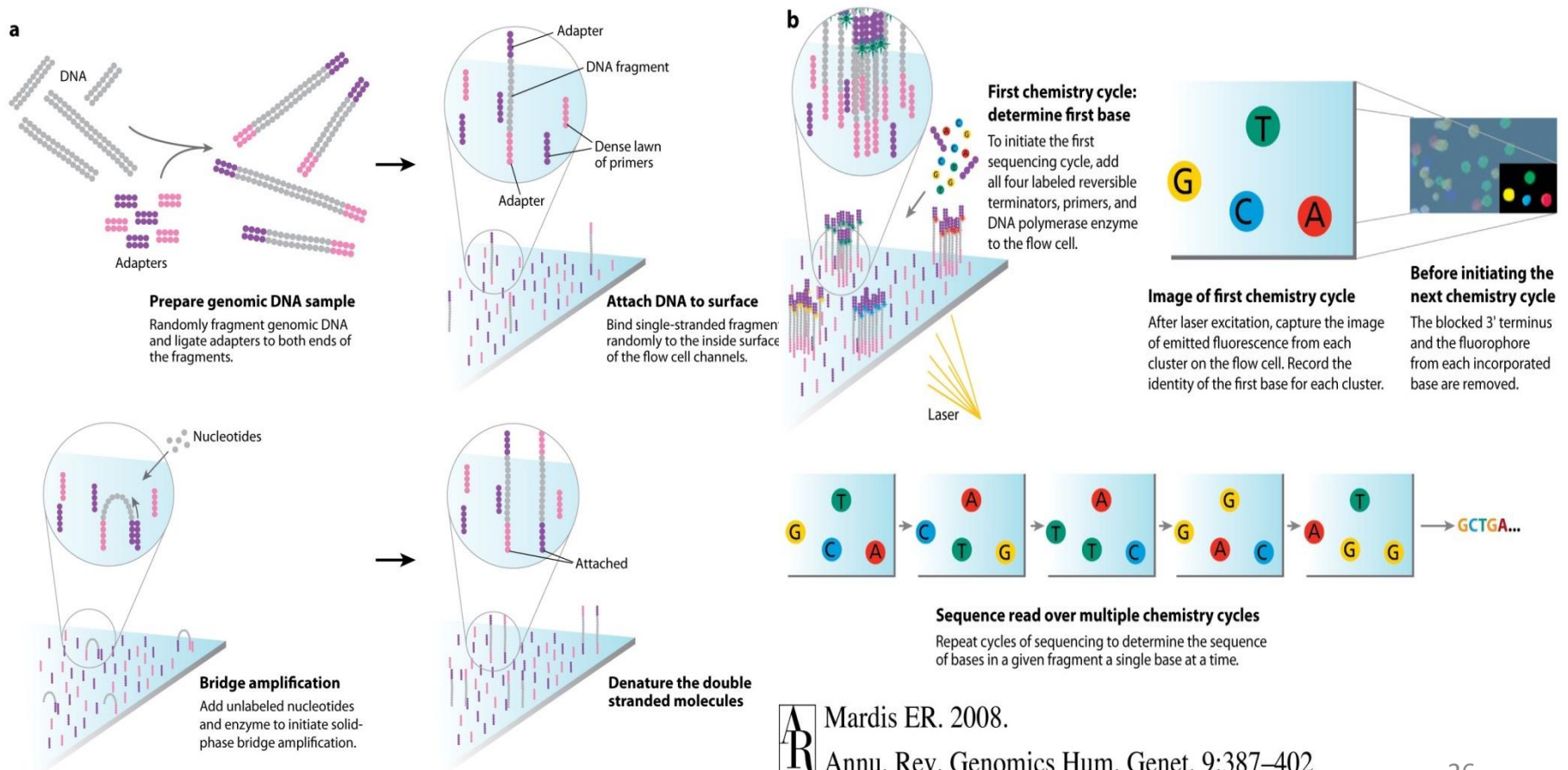
«Золотой стандарт»


Секвенирование второго поколения Next-generation sequencing (NGS)

“+” : одновременно идет секвенс большого количества разных фрагментов

“-” : прочтения длиной 75-150 нуклеотидов

Секвенирование второго поколения Next-generation sequencing (NGS) Illumina (но есть и другие приборы)



 Mardis ER. 2008.
Annu. Rev. Genomics Hum. Genet. 9:387–402

Парные и одноконцевые чтения



ATGCAGA????????????????CACTTTA

Для Illumina характерная длина чтения 100-200 п.н.

Что может пойти не так

Димеры адаптеров: адаптеры соединяются друг с другом без фрагмента ДНК между ними

Норма



Димер



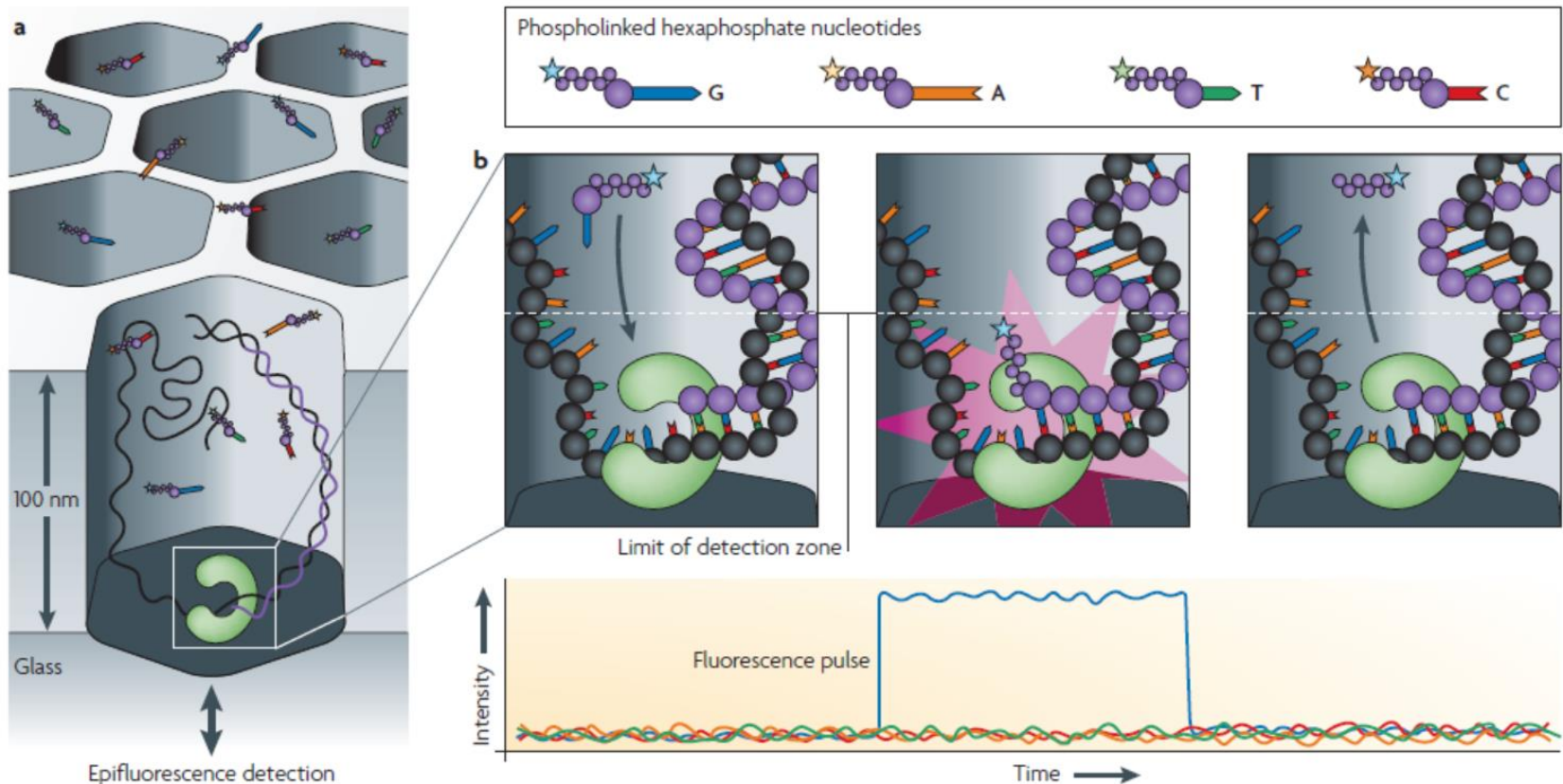
Фрагмент ДНК слишком короткий, чтение захватывает последовательности адаптеров



Одномолекулярное секвенирование в реальном времени

Pacific Biosciences

Pacific Biosciences — Real-time sequencing

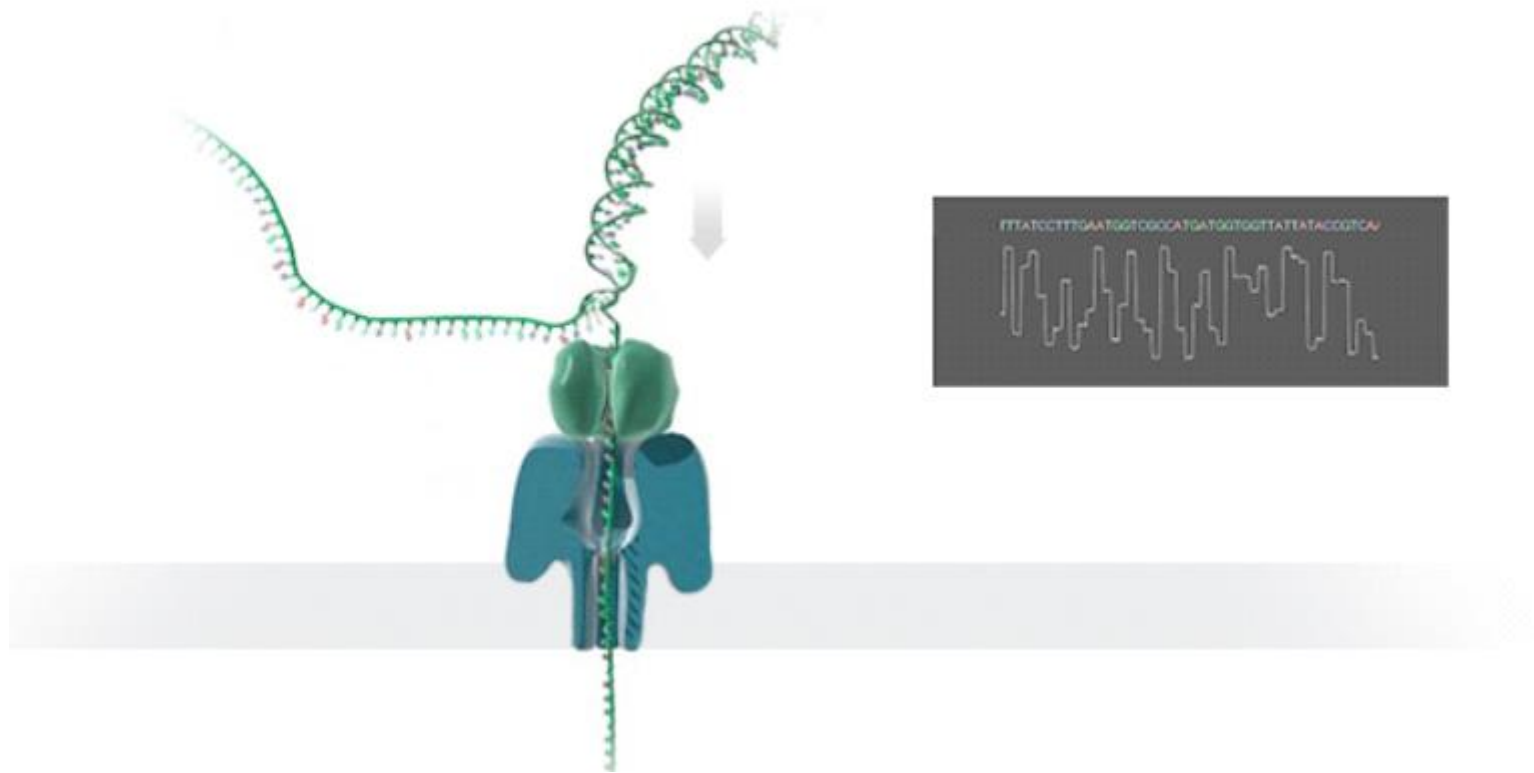


Одномолекулярное секвенирование в реальном времени Pacific Biosciences

“+” : длина прочтений 20000-60000
без амплификации
быстро

“-” : большой процент ошибок
цена

Нанопоровое секвенирование Oxford Nanopore



https://www.skygen.com/katalog/oborudovanie/oxford_nanopore_technologies/nanoporovyy_sekvenator_mini_on/

Нанопоровое секвенирование Oxford Nanopore

“+” : длина прочтений 20000-60000
без амплификации
быстро
компактность и мобильность

“-” : большой процент ошибок



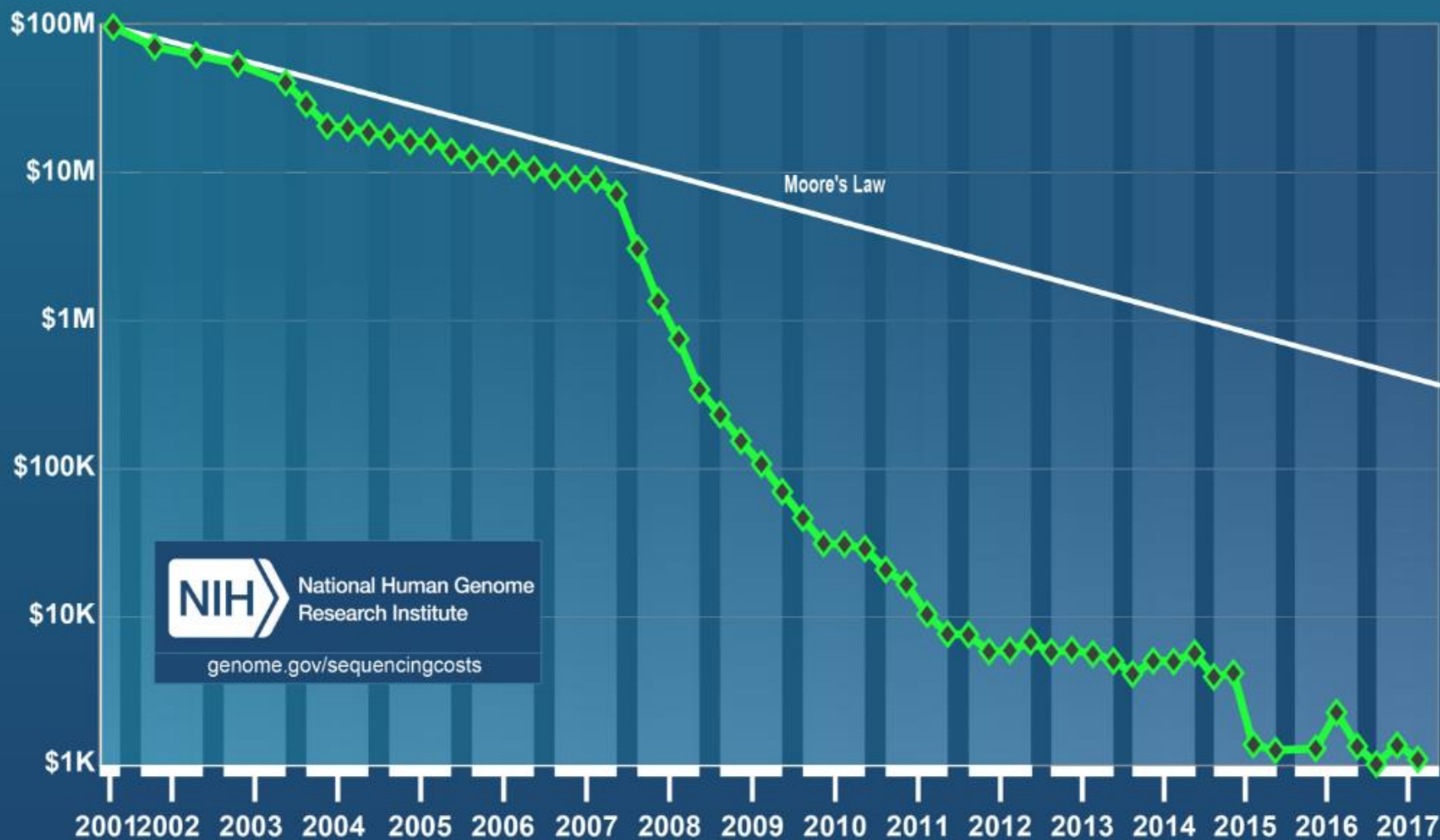
Что же выбрать?

Все зависит от задачи!

Комбинировать платформы

Увеличивать покрытие

Cost per Genome



ИСТОЧНИК: Source: [DNA Sequencing Costs: Data](http://genome.gov/sequencingcosts)

Откуда взять чтения?

<https://www.ncbi.nlm.nih.gov/sra>

SRX8794662: Whole exome seq of primary culture established from PDX tumor: Sample E9

1 ILLUMINA (Illumina HiSeq 4000) run: 31.4M spots, 6.3G bases, 2.3Gb downloads

Design: "Exom enrichment with Agilent SureSelect Human All Exon V6, based on UCSC hg19, GRCh37, February 2009"

Submitted by: NIH-phs002051

Study: DNA methylation in rhabdomyosarcoma PDX and PDX-derived primary cells

[PRJNA641459](#) • [SRP273116](#) • [All experiments](#) • [All runs](#)

[show Abstract](#)

Sample: Tumor DNA sample from N/A of a human participant in the dbGaP study "DNA Methylation in Rhabdomyosarcoma PDX and PDX-Derived Primary Cells"

[SAMN15468651](#) • [SRS7062698](#) • [All experiments](#) • [All runs](#)

Organism: [Homo sapiens](#)

Library:

Name: ON-2018/8626: E9

Instrument: Illumina HiSeq 4000

Strategy: WXS

Source: GENOMIC

Selection: PCR

Layout: PAIRED

The SRA run(s) below contain human sequence ([more...](#))

Runs: 1 run, 31.4M spots, 6.3G bases, 2.3Gb

| Run | # of Spots | # of Bases | Size | Published |
|-----------------------------|------------|------------|-------|------------|
| SRR12291396 | 31,383,123 | 6.3G | 2.3Gb | 2020-08-27 |

ID: 11430914

Важная информация!

<https://www.ncbi.nlm.nih.gov/sra>

SRX8794662: Whole exome seq of primary culture established from PDX tumor: Sample E9

1 ILLUMINA (Illumina HiSeq 4000) run: 31.4M spots, 6.3G bases, 2.3Gb downloads

Design: "Exom enrichment with Agilent SureSelect Human All Exon V6, based on UCSC hg19, GRCh37, February 2009"

Submitted by: NIH-phs002051

Study: DNA methylation in rhabdomyosarcoma PDX and PDX-derived primary cells

[PRJNA641459](#) • [SRP273116](#) • [All experiments](#) • [All runs](#)

[show Abstract](#)

Sample: Tumor DNA sample from N/A of a human participant in the dbGaP study "DNA Methylation in Rhabdomyosarcoma PDX and PDX-Derived Primary Cells"

[SAMN15468651](#) • [SRS7062698](#) • [All experiments](#) • [All runs](#)

Organism: [Homo sapiens](#)

Library:

Name: ON-2018/8626: E9

Instrument: Illumina HiSeq 4000

Strategy: WXS

Source: GENOMIC

Selection: PCR

Layout: PAIRED

The SRA run(s) below contain human sequence ([more...](#))

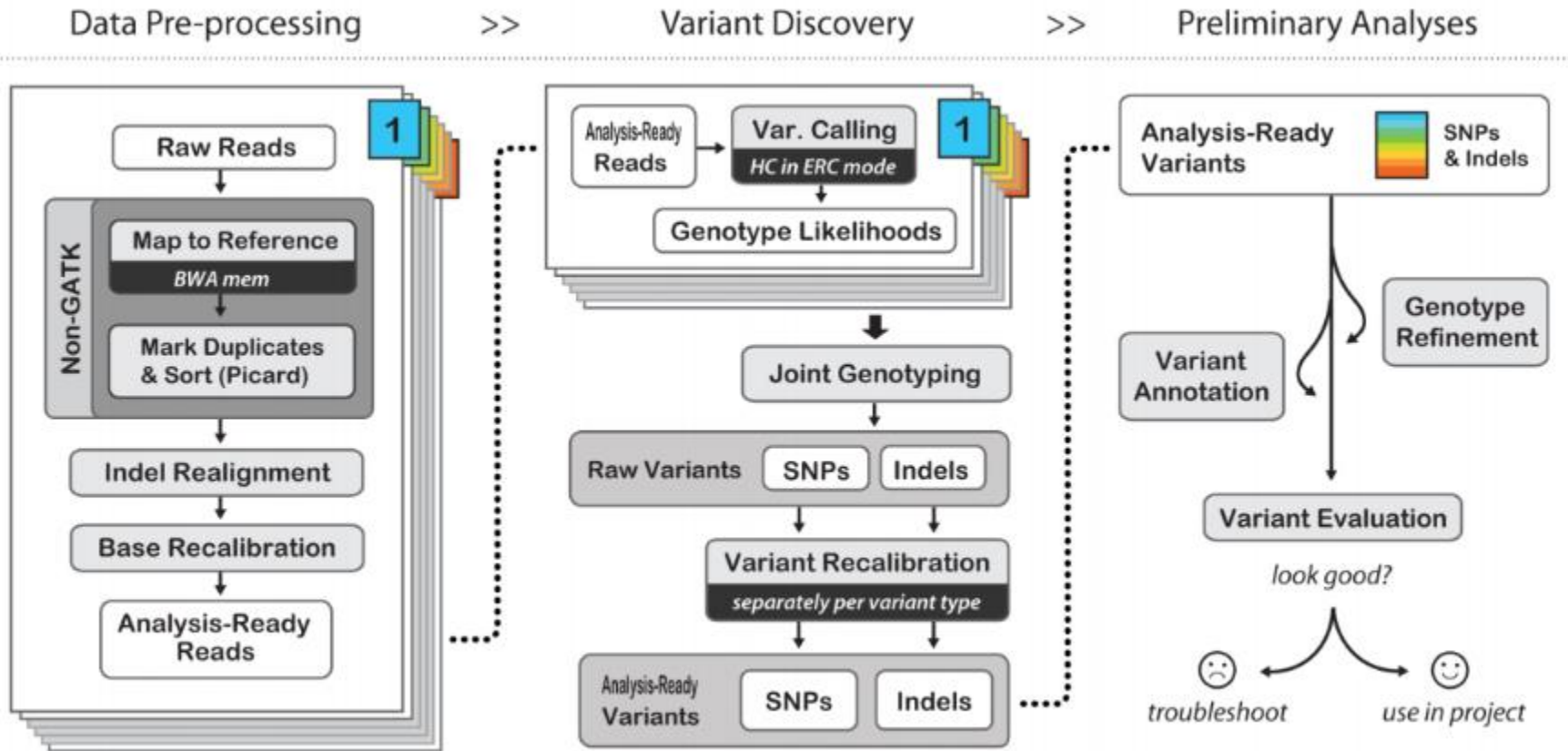
Runs: 1 run, 31.4M spots, 6.3G bases, 2.3Gb

| Run | # of Spots | # of Bases | Size | Published |
|-----------------------------|------------|------------|-------|------------|
| SRR12291396 | 31,383,123 | 6.3G | 2.3Gb | 2020-08-27 |

ID: 11430914

Sra toolkit - <https://www.ncbi.nlm.nih.gov/sra/docs/sradownload/>

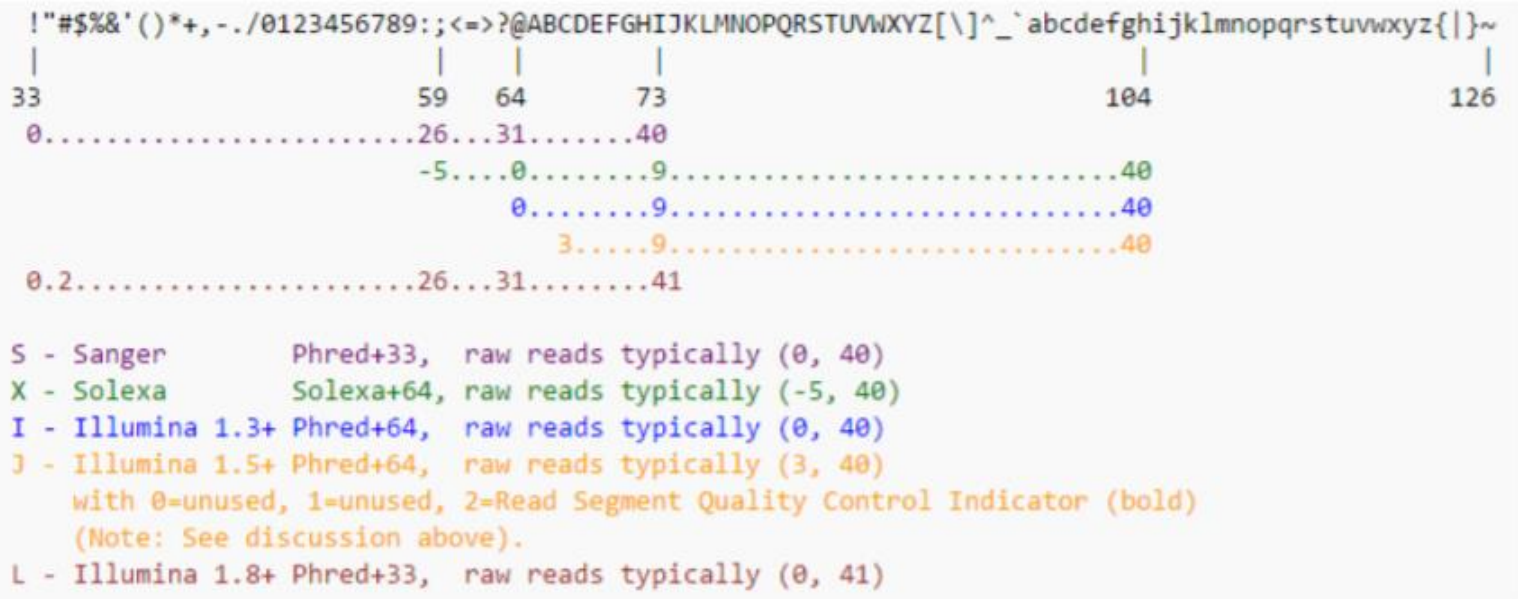
Обработка данных



Fastq формат

```

@HWI-ST992:147:D22HDACXX:3:1112:14175:15297 2:N:0:GGCTAC
Последовательность TAATGGCTTTTCCAAAACGCTCCACTCTTAAAGATGTGTATAAGAGACAGCAACAACAATTA
+
Качество 8??DDDBEDHHFHJJJJJJAFFGIIIIIGIGEEGIIIIHBFGEEGCGIJIFFIDIIJJIIII
    
```



Качество чтений

P – вероятность ошибки

Q – параметр качества. (Phred Quality Score)

$$Q = -10\log_{10}P$$

| Вероятность ошибки | Q |
|---------------------------------|----|
| 0.001 (<i>точность 99,9%</i>) | 30 |
| 0.01 (<i>точность 99%</i>) | 20 |
| 0.1 (<i>точность 90%</i>) | 10 |

Типичные значения Q от 1 до 40

Q>20 – «хорошее качество»

Пересчет качества в вероятность ошибки

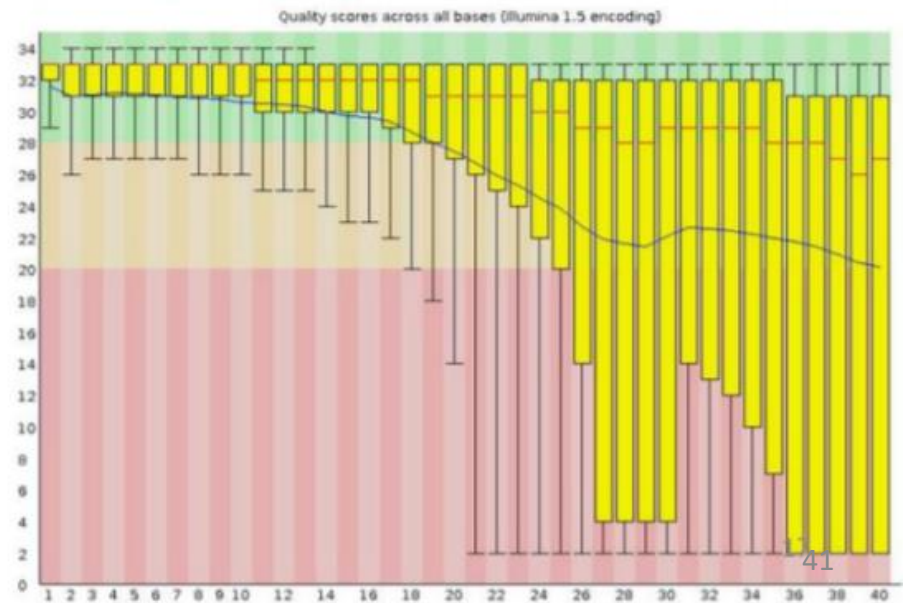
| Phred Quality Score | Символ | Вероятность ошибки | Точность |
|----------------------------|---------------|---------------------------|-----------------|
| 10 | + | 1/10 | 90% |
| 20 | 5 | 1/100 | 99% |
| 30 | ? | 1/1000 | 99,9% |
| 40 | | 1/10 000 | 99,99% |
| 50 | S | 1/100 000 | 99,999% |
| 60 |] | 1/1 000,000 | 99,9999% |

fastQC

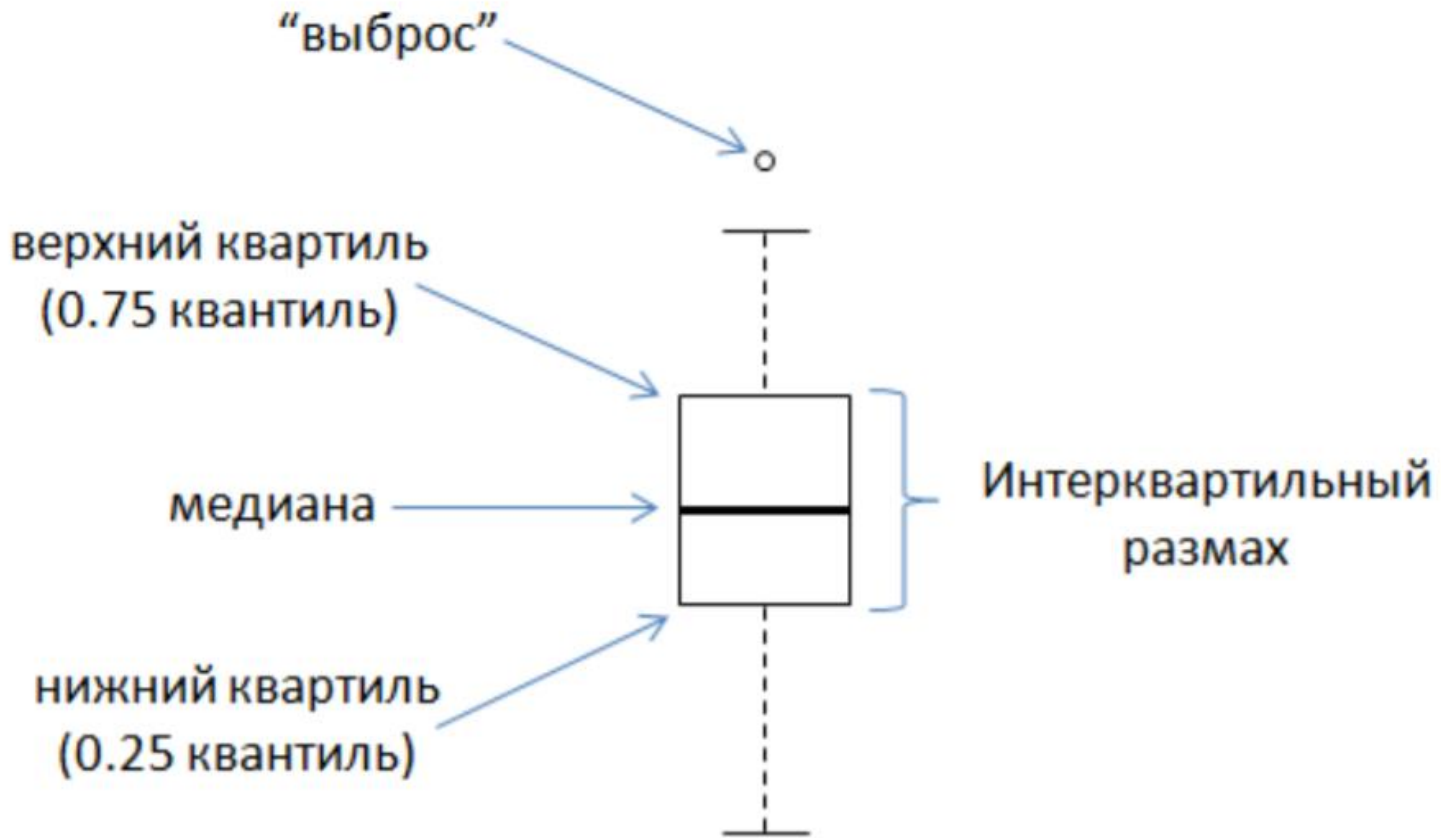
✔ Per base sequence quality



✘ Per base sequence quality



«Ящик с усами» / диаграмма размахов / boxplot



fastQC

<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Программа FasqQC стоит на kodoMo

Версию с графическим интерфейсом можно поставить на свой компьютер.

На сайте отличное руководство!

