



# ВОЗМОЖНОСТИ СЕКВЕНЦИРОВАНИЯ

Анастасия Жарикова – 2022

[azharikova89@gmail.com](mailto:azharikova89@gmail.com)

# Что такое секвенирование?

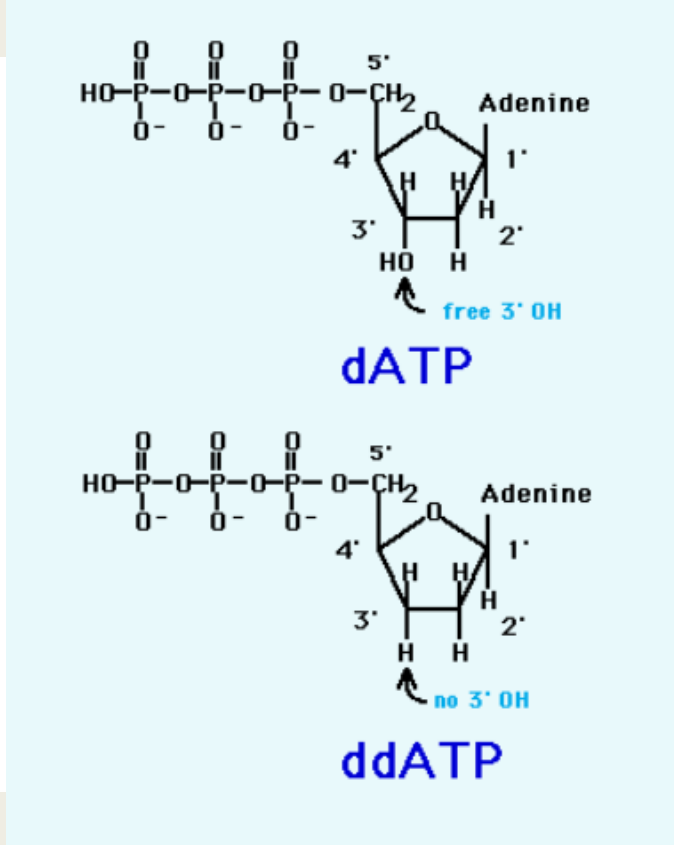
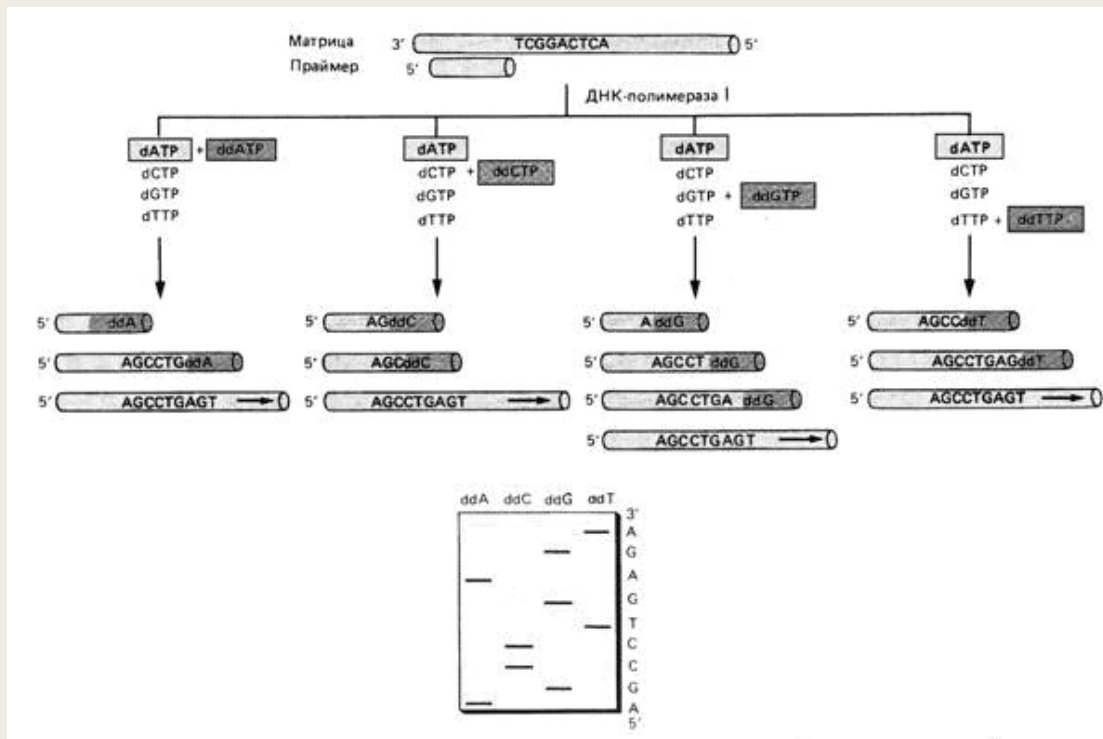
Определение последовательности макромолекулы

Будем говорить о нуклеиновых кислотах =>

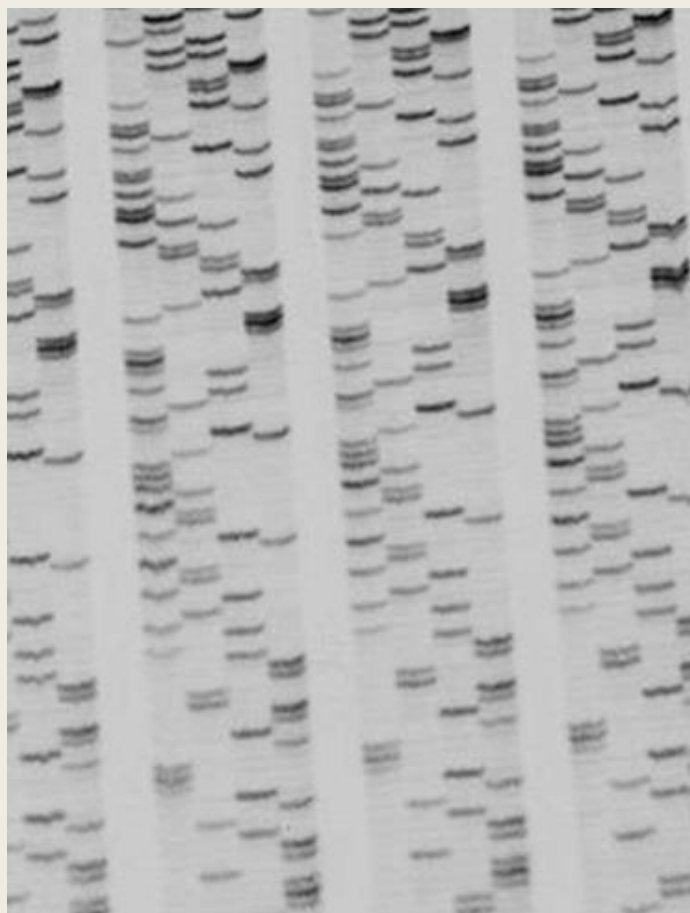
**РНК** (рибонуклеиновая кислота) или **ДНК**  
(дезоксирибонуклеиновая кислота) =>

Последовательность нуклеотидов

# Метод «терминаторов» секвенирование ДНК по Сэнгеру – 1977г

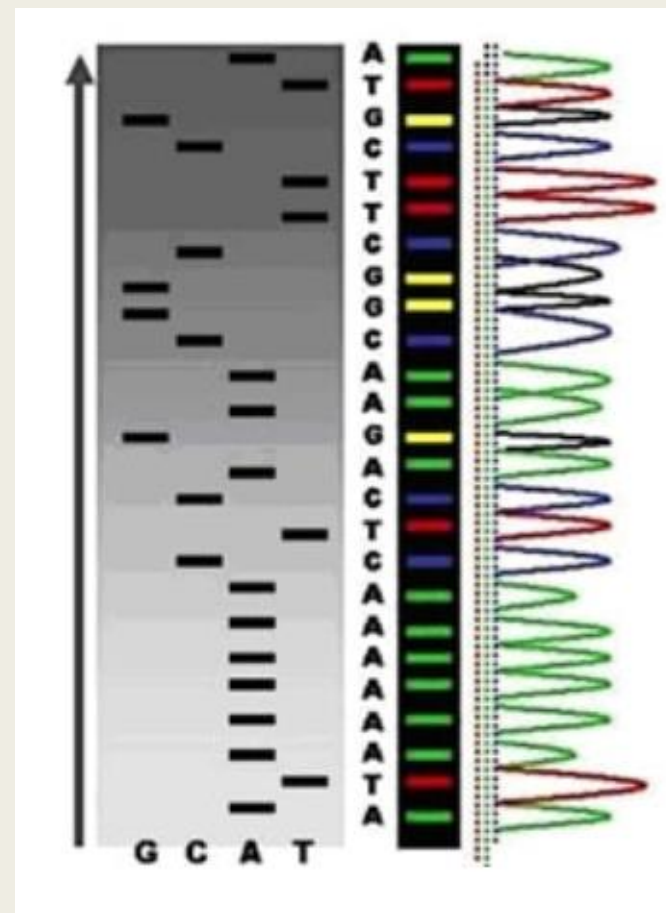


# Метод «терминаторов»



← Было

→ Стало

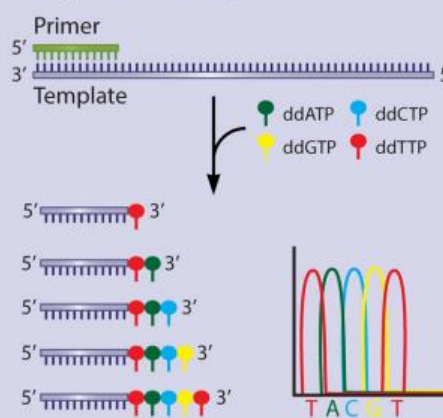


Секвенирование нуклеиновых кислот

«Золотой стандарт»  
~ 1000 п.н.

# Поколения подходов к секвенированию

### First Generation Shotgun Sequencing

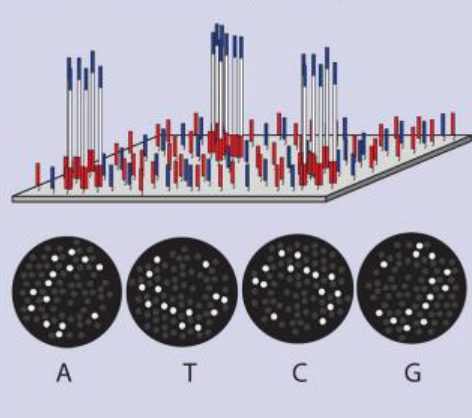


The diagram illustrates the first generation of sequencing. It shows a DNA template with a primer hybridized to it. The primer is extended using four dideoxynucleotides (ddNTPs): ddATP (green), ddCTP (blue), ddGTP (yellow), and ddTTP (red). The extension is stopped at various positions, creating fragments of different lengths. These fragments are then separated by size using gel electrophoresis, resulting in a chromatogram with four distinct peaks corresponding to the bases T, A, C, and G.

- Sequencing by synthesis
- High accuracy
- Long read lengths
- Relatively small amount of data generated

e.g., ABI capillary sequencer (ABI)

### Second Generation Massively Parallel Sequencing

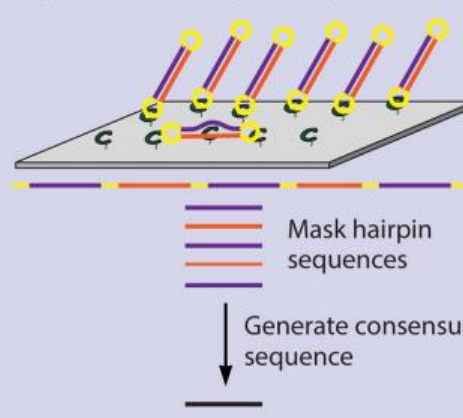


The diagram illustrates the second generation of sequencing. It shows a large number of DNA templates being amplified and sequenced simultaneously. The resulting data is represented by a grid of colored bars (red, blue, green, yellow) for each base. Below the grid are four circular images labeled A, T, C, and G, showing the distribution of reads for each base.

- Sequencing by synthesis
- Amplified templates are generated during sequencing, reducing the requirements for starting material
- High accuracy
- Short read lengths

e.g., MiSeq (Illumina), Ion Torrent (Thermo Fisher Scientific)

### Third Generation Single-molecule Sequencing



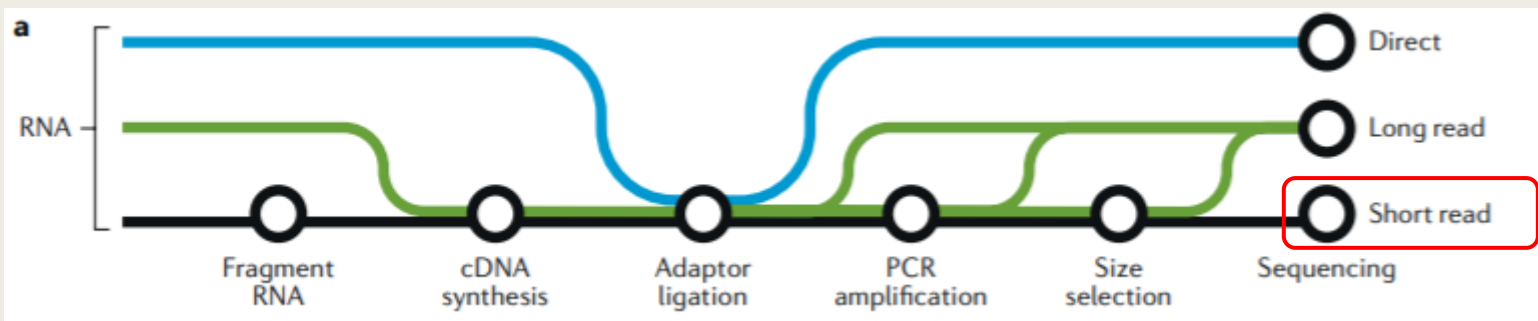
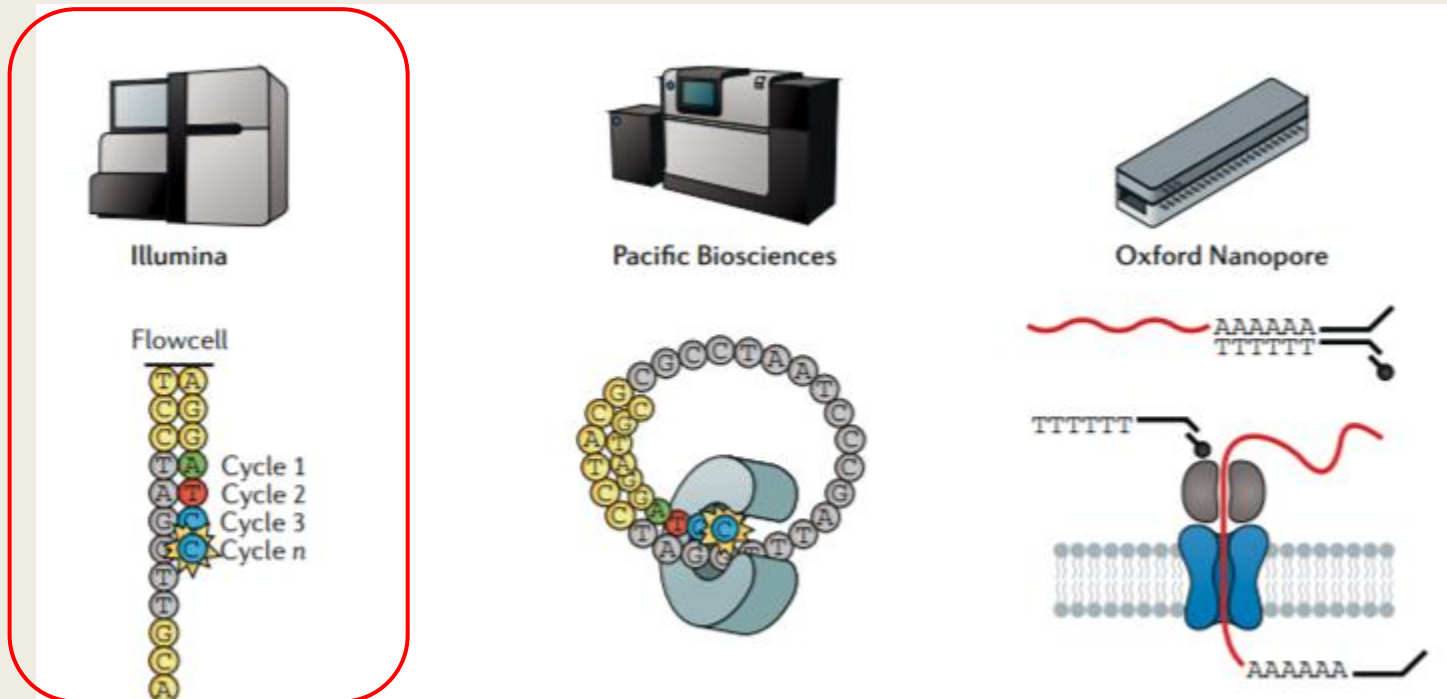
The diagram illustrates the third generation of sequencing. It shows single-molecule templates being sequenced. Mask hairpin sequences are used to facilitate the process. The resulting data is used to generate a consensus sequence.

- Single-molecule templates
- Low accuracy
- Long read lengths

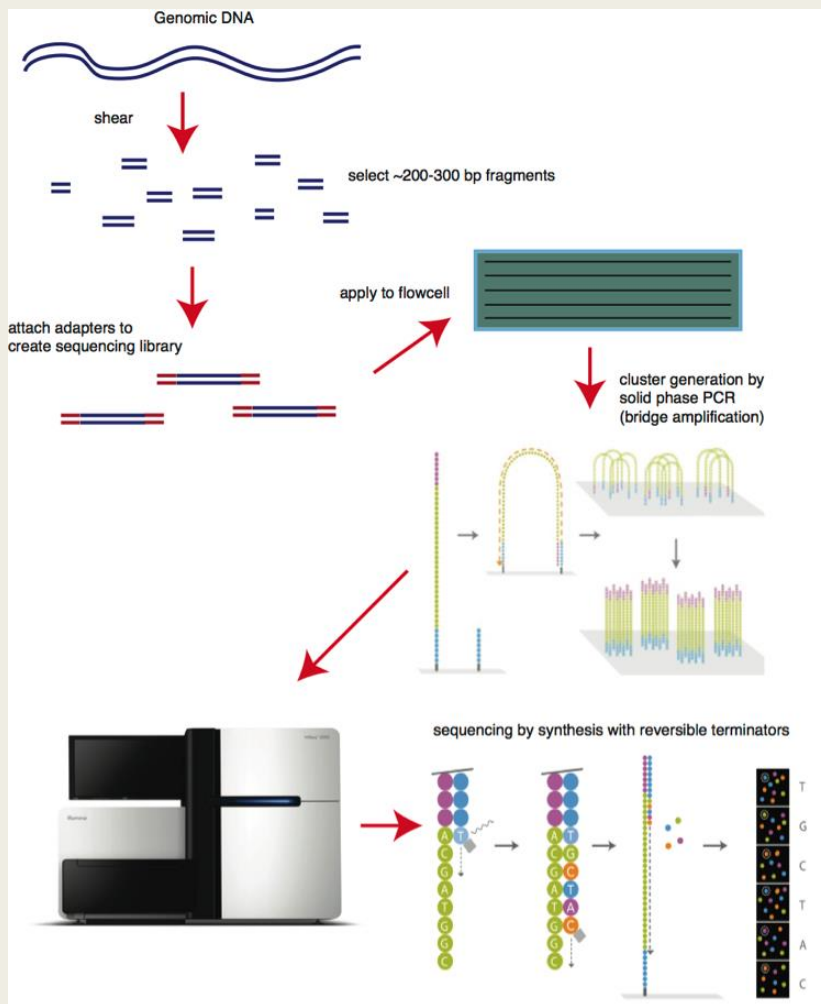
e.g., Single-Molecule Real-Time (SMRT) — Sequencing (Pacific Biosciences), MinION (Oxford Nanopore Technologies)

## Секвенирование нуклеиновых кислот

# Разнообразие платформ



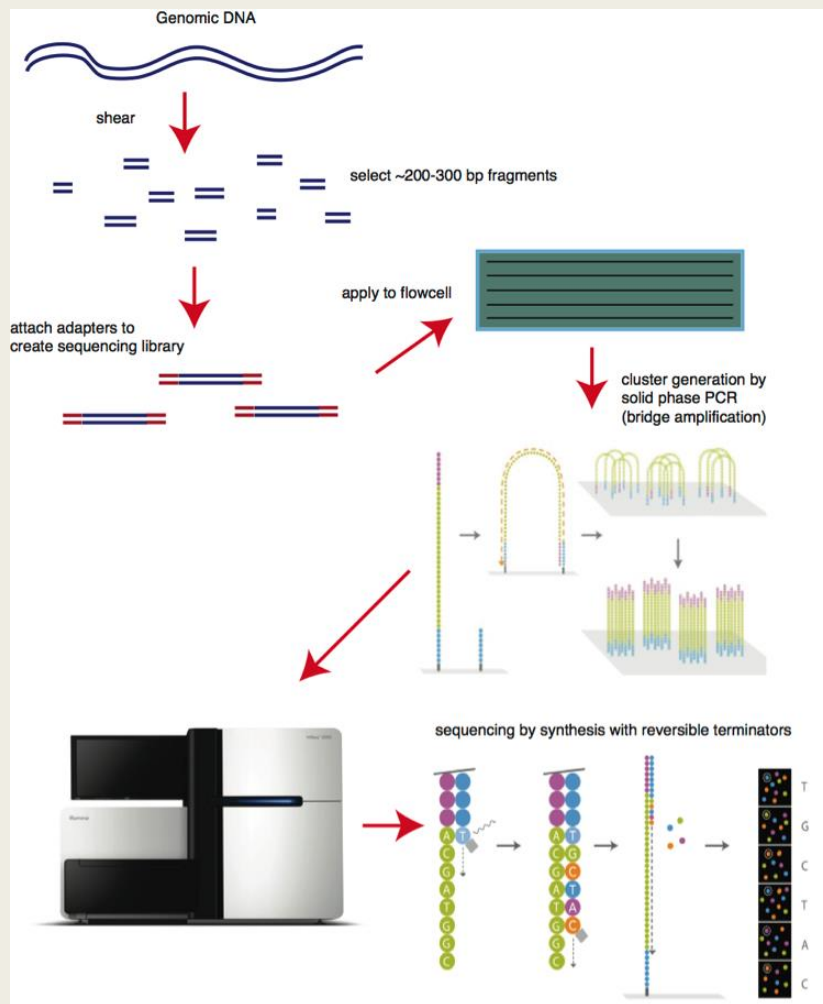
# Next generation sequencing – NGS - Illumina



Выделение ДНК  
Фрагментация  
Подготовка библиотеки  
Секвенирование

МУЛЬТИК

# Next generation sequencing – NGS - Illumina



Выделение ДНК  
Фрагментация  
Подготовка библиотеки  
Секвенирование

МУЛЬТИК

Умеем секвенировать только ДНК  
А что же с РНК?  
РНК => обратная транскрипция =>  
кДНК => секвенирование

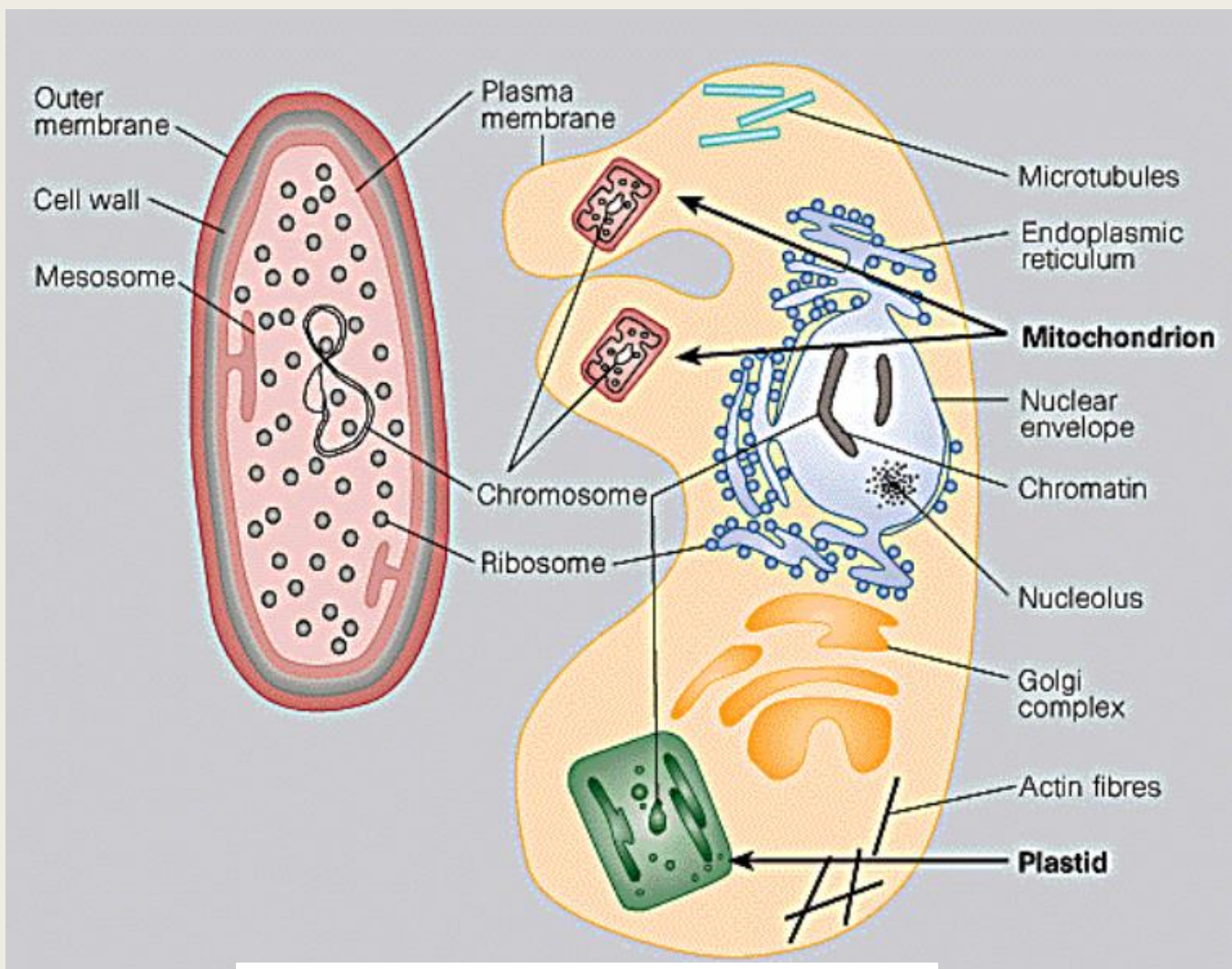
В итоге мы в любом случае  
секвенируем короткие фрагменты  
ДНК



# DNA-seq & RNA-seq

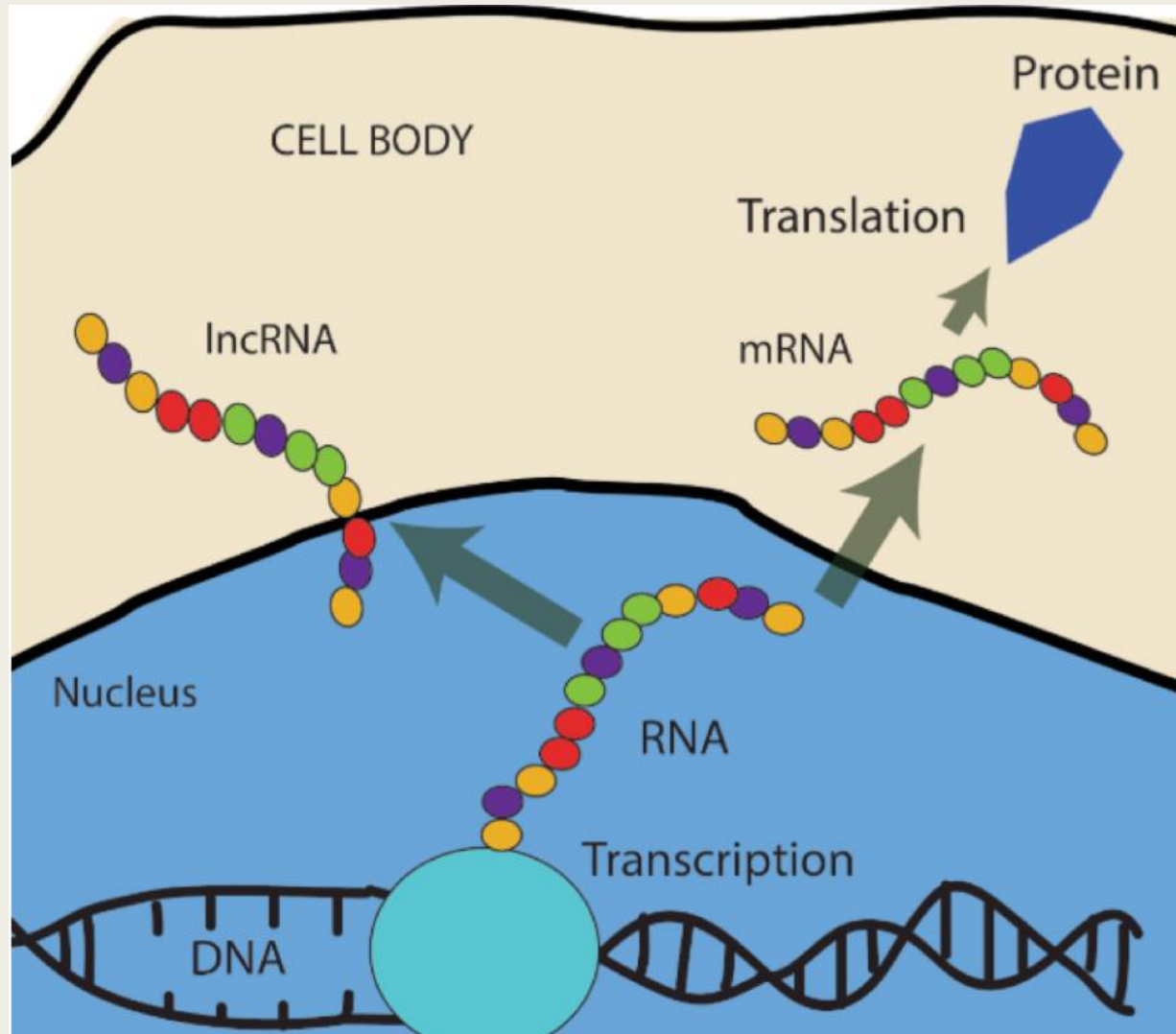
- Стремимся секвенировать
  - *Всю ДНК => геном*
  - *Конкретные области генома*
  - *Всю РНК => транскриптом*
  - *Отдельные фракции транскриптома*

# Где находится ДНК?

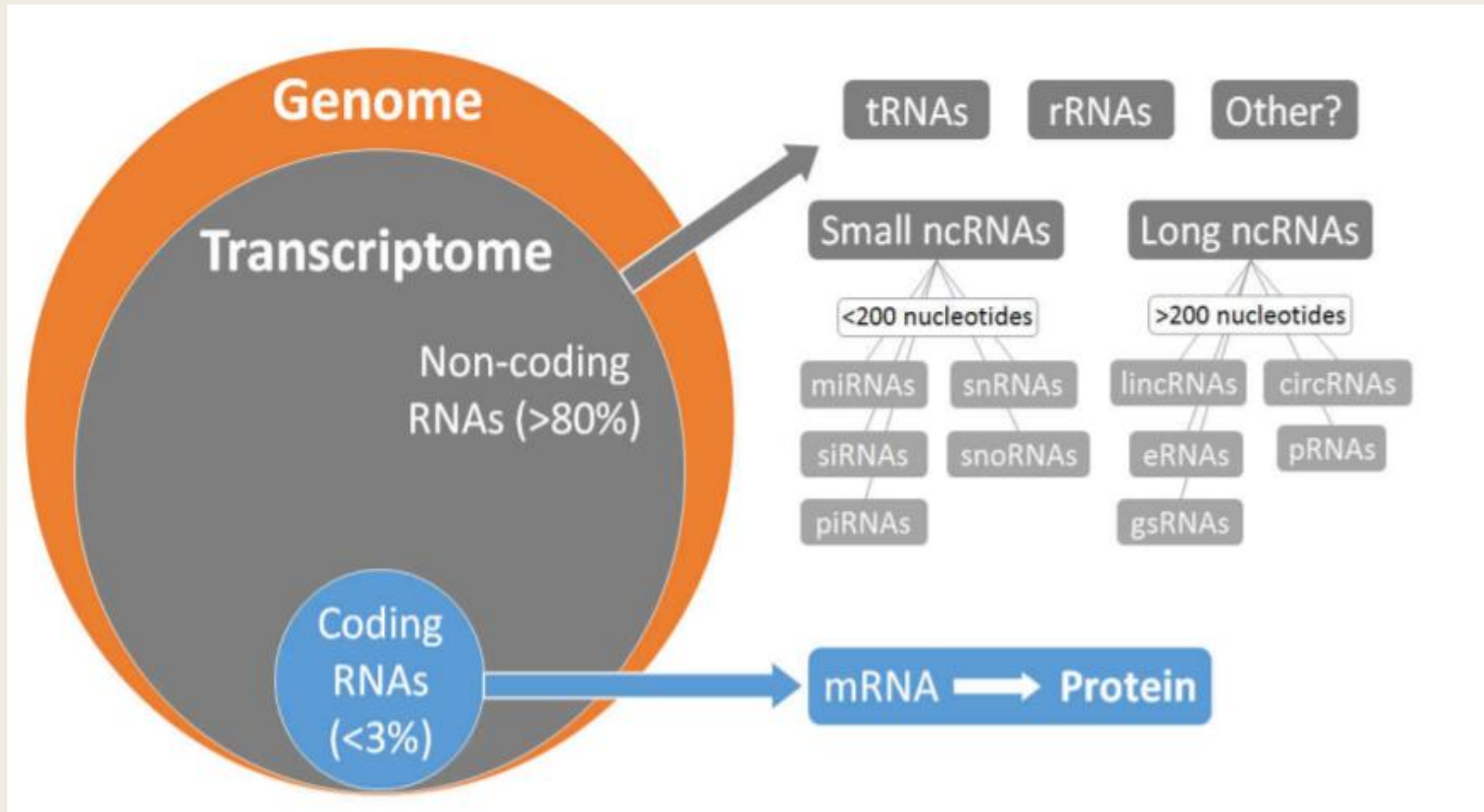


Typical prokaryotic (left) and eukaryotic (right) cells

# Откуда берется РНК?



# Какие бывают РНК?



# GENCODE

## разметка генов человека

Total No of Genes	61533
Protein-coding genes	19982
Long non-coding RNA genes	18811
Small non-coding RNA genes	7567
Pseudogenes	14763
- processed pseudogenes	10662
- unprocessed pseudogenes	3557
- unitary pseudogenes	243
- polymorphic pseudogenes	50
- pseudogenes	15
Immunoglobulin/T-cell receptor gene segments	
- protein coding segments	409
- pseudogenes	236

# 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

# Некодирующие РНК

## **Noncoding RNAs Databases: Current Status and Trends**

**Vinicius Maracaja-Coutinho, Alexandre Rossi Paschoal, José Carlos Caris-Maldonado, Pedro Vinícius Borges, Almir José Ferreira, and Alan Mitchell Durham**

### **Abstract**

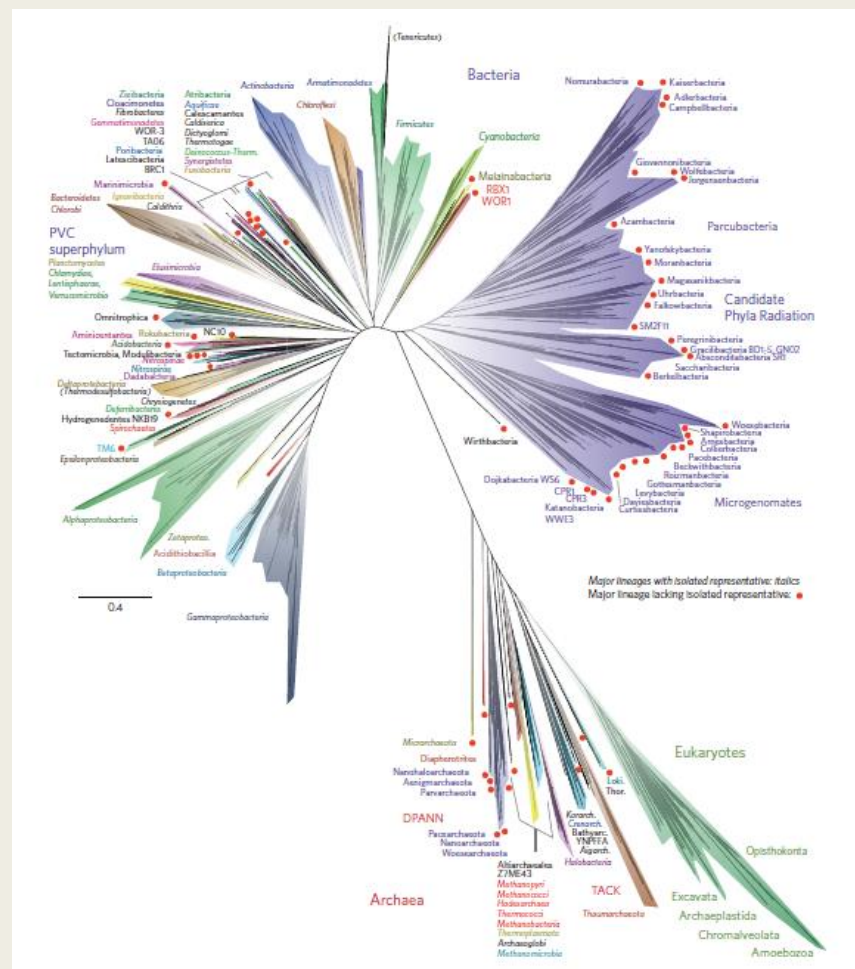
One of the most important resources for researchers of noncoding RNAs is the information available in public databases spread over the internet. However, the effective exploration of this data can represent a daunting task, given the large amount of databases available and the variety of stored data. This chapter describes a classification of databases based on information source, type of RNA, source organisms, data formats, and the mechanisms for information retrieval, detailing the relevance of each of these classifications and its usability by researchers. This classification is used to update a 2012 review, indexing now more than 229 public databases. This review will include an assessment of the new trends for ncRNA research based on the information that is being offered by the databases. Additionally, we will expand the previous analysis focusing on the usability and application of these databases in pathogen and disease research. Finally, this chapter will analyze how currently available database schemas can help the development of new and improved web resources.

# Геномы организмов в NCBI

- Eukaryotes (21680)
- Prokaryotes (390644)
- Viruses (47571)
- Plasmids (36213)
- Organelles (21972)

Задачи:

- Сборка геномов
- Сравнительная геномика



<https://www.nature.com/articles/nmicrobiol201648.pdf>



# Число хромосом у разных видов



Гиббоны - 44



Макака - 42



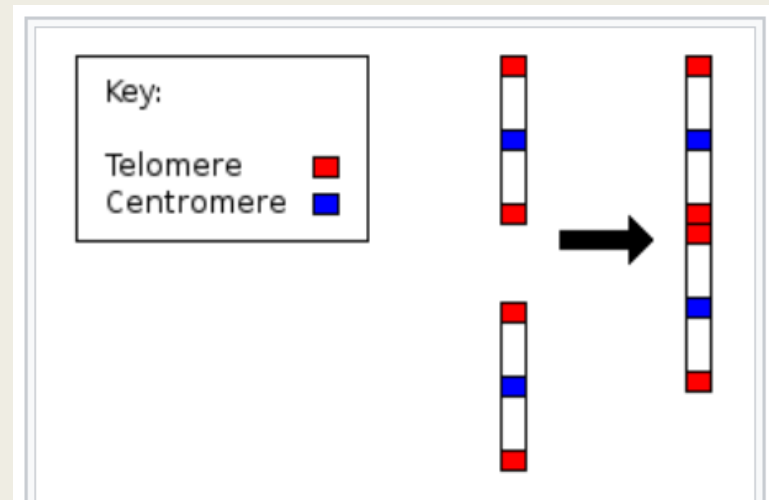
Капуцин - 54



48



46



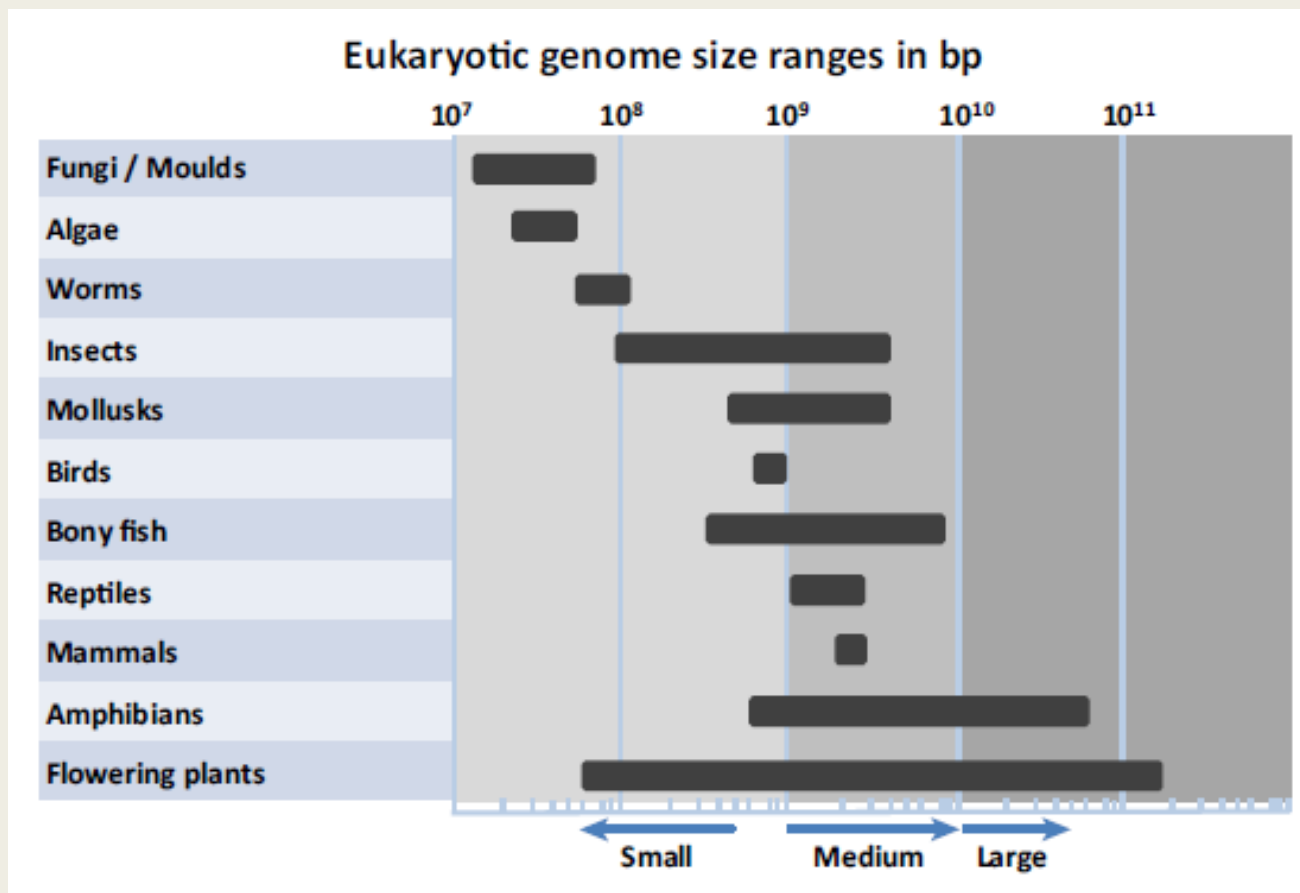
После слияния двух хромосом остаются характерные следы: остатки теломер и рудиментарная центромера

# Число хромосом у разных видов

- Муравей (*Myrmecia 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



# Размер генома



Размер генома человека ~ 3 млрд п.н.

База про размеры геномов: <https://www.genomesize.com/statistics.php>

# Доместикация риса

```
1  MSGSSADPSP  SASTAGAAVS  PLALLRAHGH  GHGHLTATPP  SGATGPAPPP
51  PSPASGSAPR  DYRKGNWTLH  ETLILITANR  LDDDRRAGVG  GAAAGGGGAG
101 SPPTPRSAEQ  RWKWVENYCW  KNGCLRSQNG  CNDKWDNLLR  DYKKVRDYES
151 RVAAAAATGG  AAAANSAPLP  SYWTMERHER  KDCNLPTNLA  PEVYDALSEV
201 LSRRAARRGG  ATIAPT PPPP  PLALPLPPPP  PPSPPKPLVA  QQQHHHHGHH
251 HHPPPQPPP  SSLQLPPAVV  APPPASVSAE  EEMSGSSESG  EEEEGSGGEP
301 EAKRRRLSRL  GSSVVR SATV  VARTLVACEE  KRERRHRELL  QLEERRRLRE
351 EERTEVRRQG  FAGLIAAVNS  LSSAIHALVS  DHRSGDSSGR
```

Дикий рис – AAG – лизин

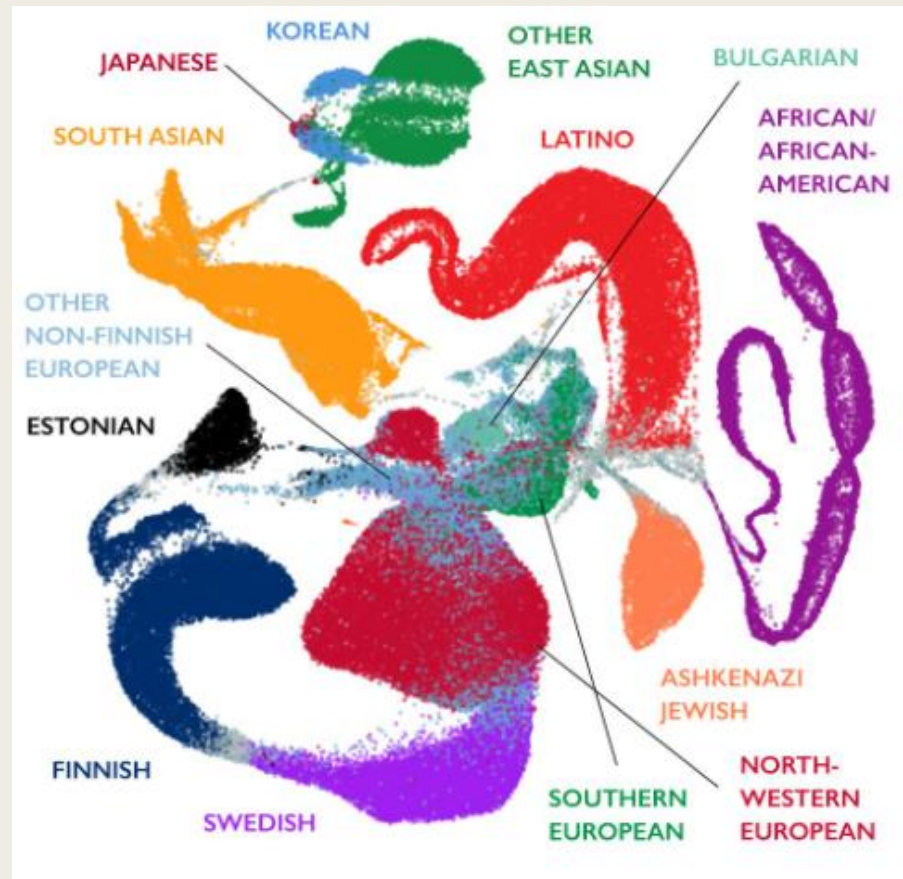
Культурный рис – AAT – аспарагин

# DNA-seq

- Изучение отличий образца от референсного генома
- Можно секвенировать
  - *Геном*
  - *Экзоны мРНК => экзом (< 2%)*
  - *Конкретные области или гены*

# Коллекции

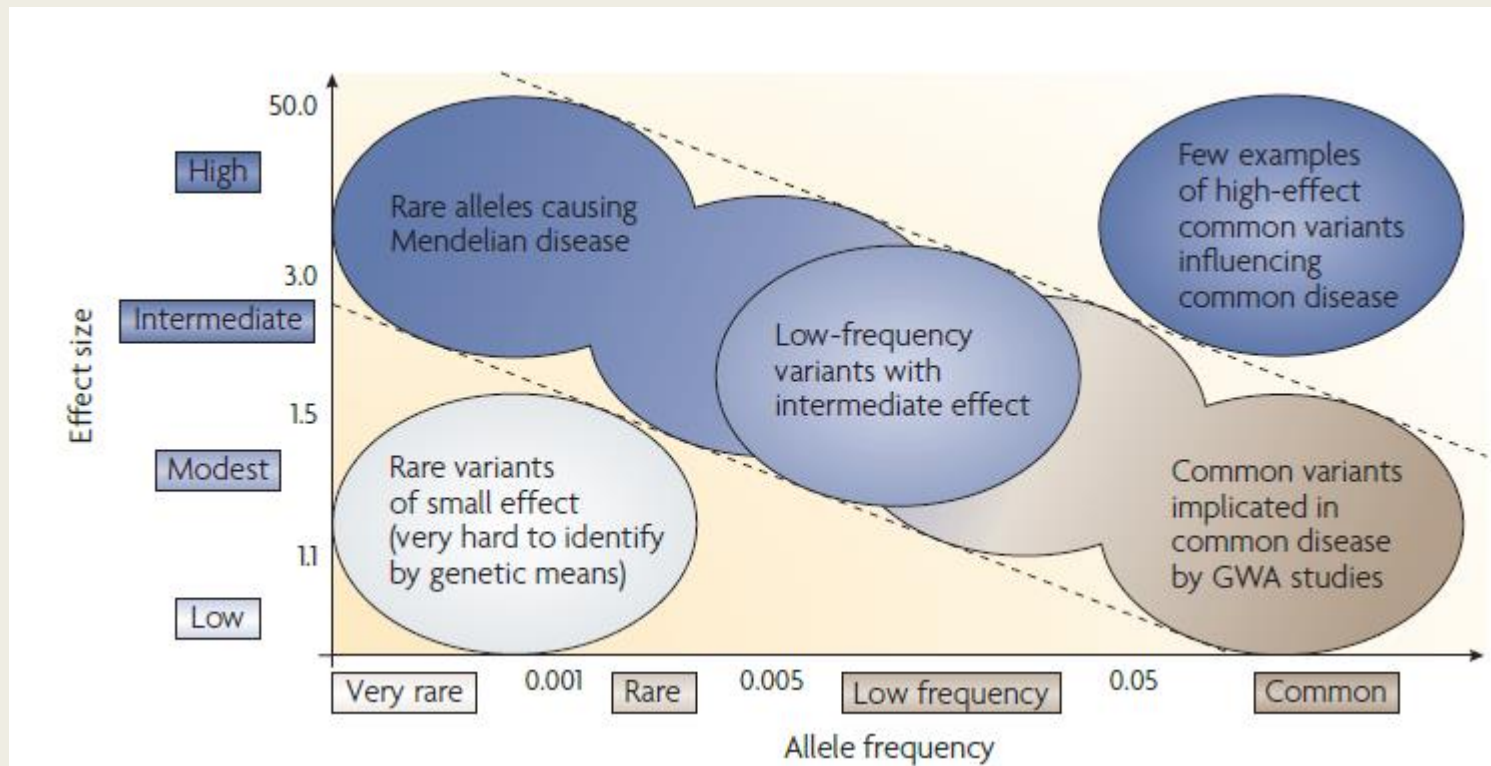
- [1000 genomes](#) project
- The genome aggregation database ([gnomAD](#))
- [100000 genomes](#)
- [UK Biobank](#)
- ...



# Зачем?

- Геногеография – изучает географическое распространение генетических признаков живых организмов
- Медицинская геномика – изучает влияние геномных нарушений на развитие наследственных заболеваний
- Фармакогенетика – изучает генетические особенности человека, влияющие на индивидуальную переносимость лекарственных средств
- ...

# Генетика и заболевания

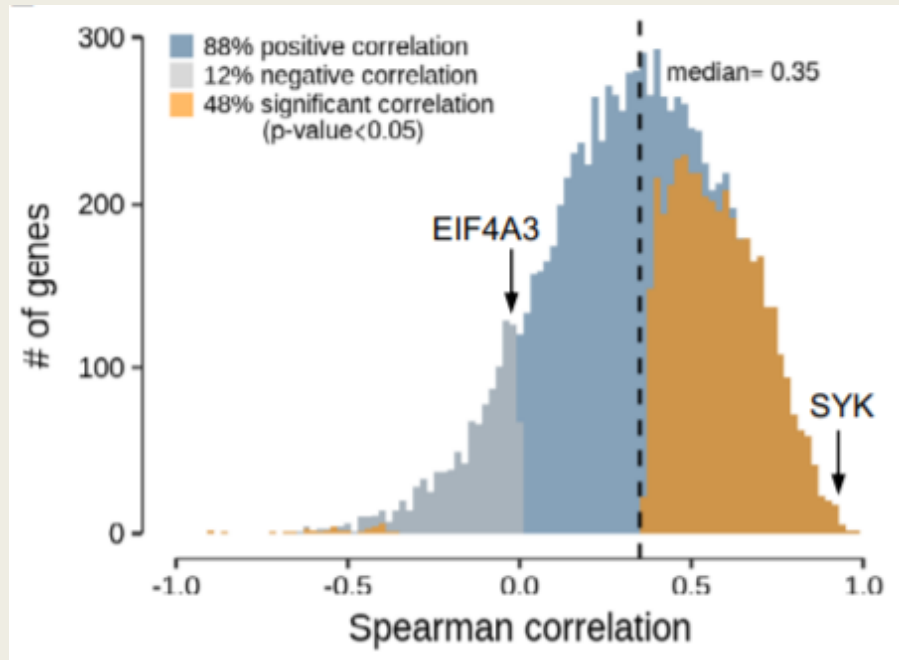




# RNA-seq - Задачи

- Экспрессия генов
- Дифференциальная экспрессия генов
- Сплайсинг
- Редактирование РНК
- Сборка транскриптома

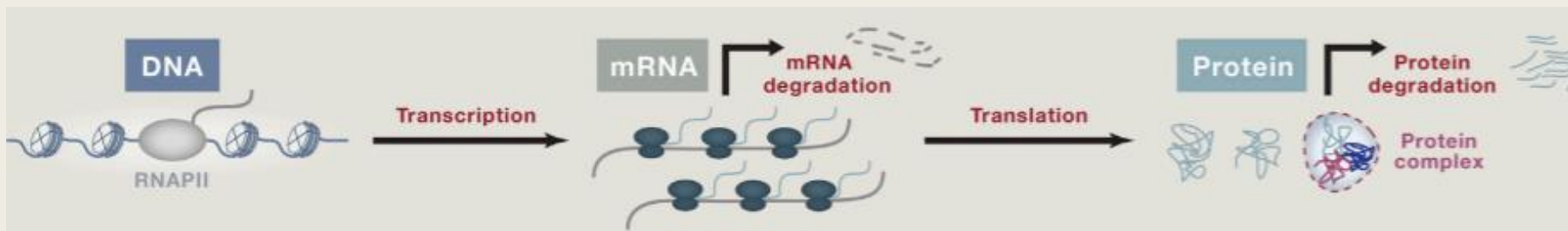
# MPHC



Транскрипционный профиль  
белок-кодирующих генов

Можем косвенно судить о  
концентрации белков

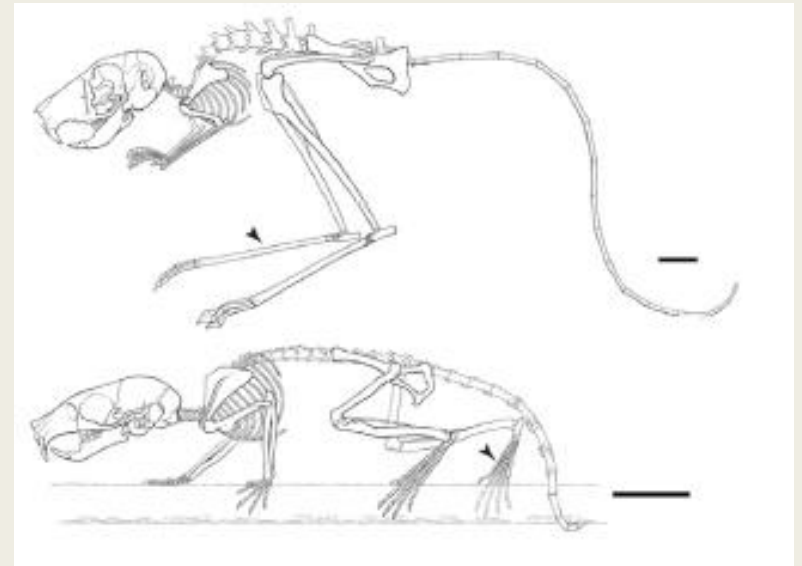
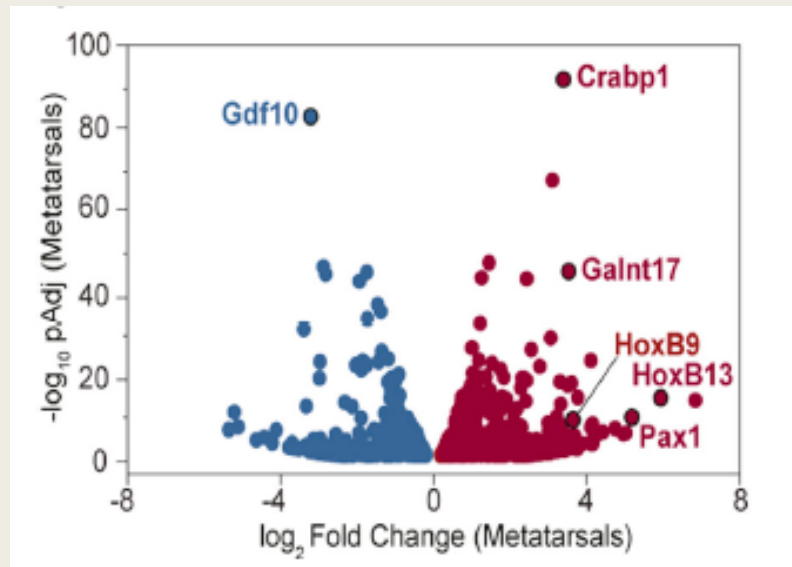
[A deep proteome and transcriptome abundance atlas of 29 healthy human tissues](#)



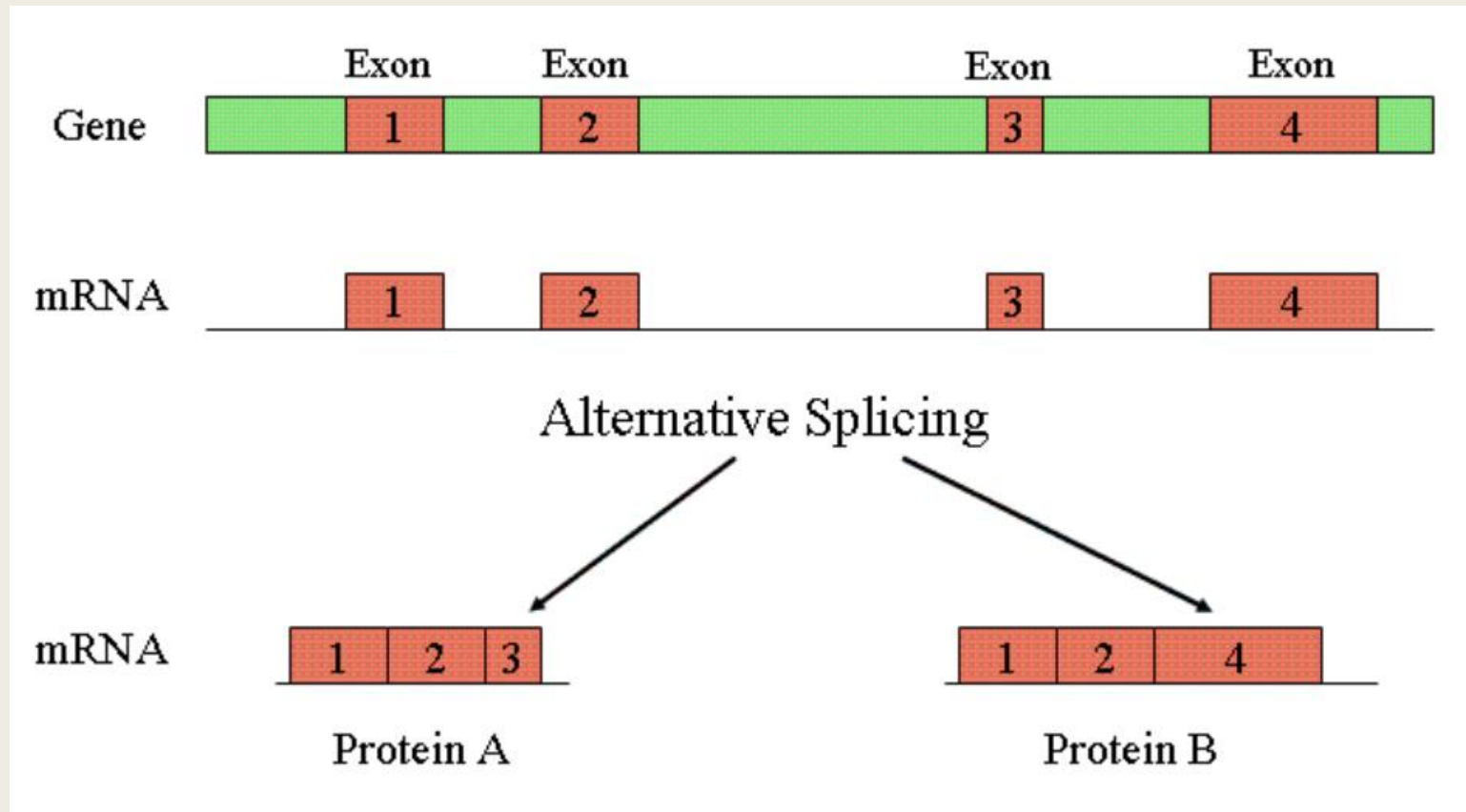
[On the Dependency of Cellular Protein Levels on mRNA Abundance](#)

# Дифференциальная экспрессия

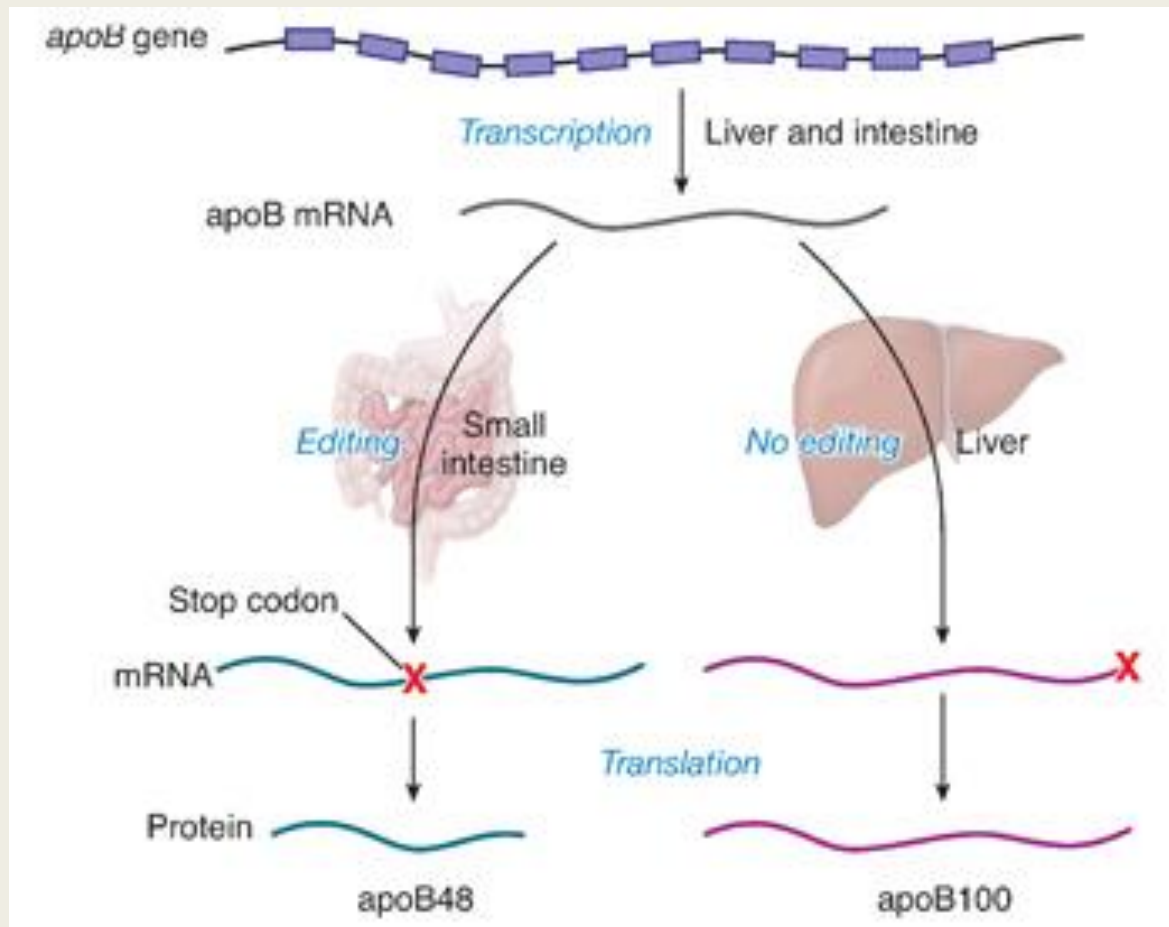
Interspecies transcriptomics identify genes that underlie disproportionate foot growth in jerboas



# Альтернативный сплайсинг



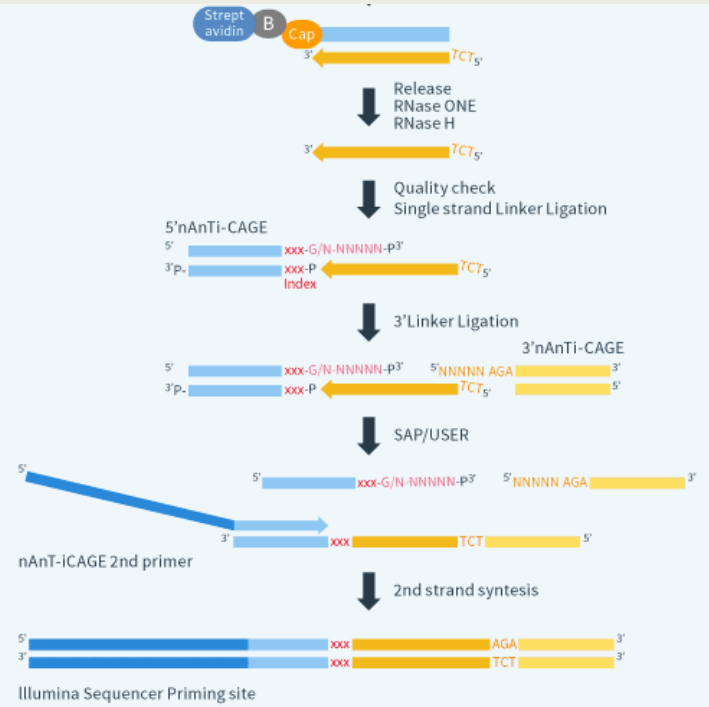
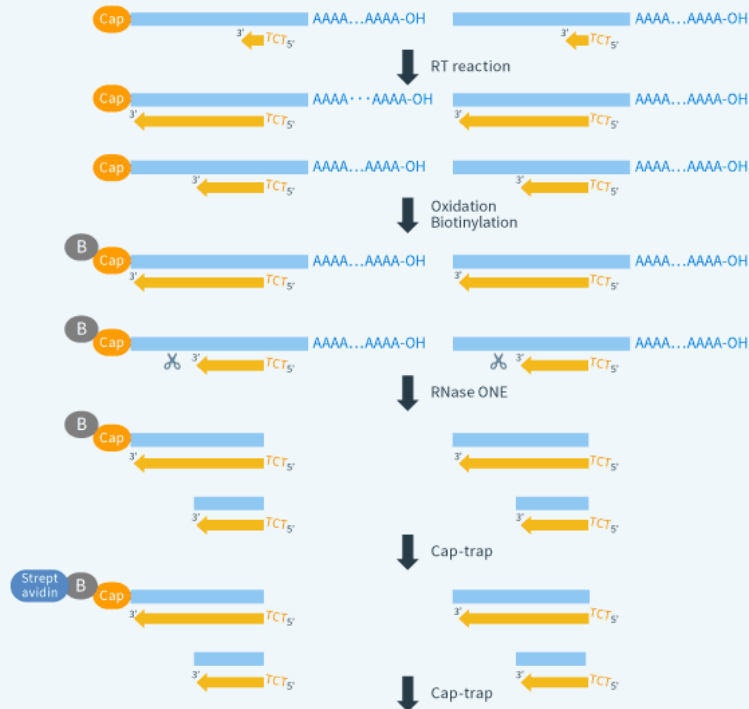
# Редактирование РНК



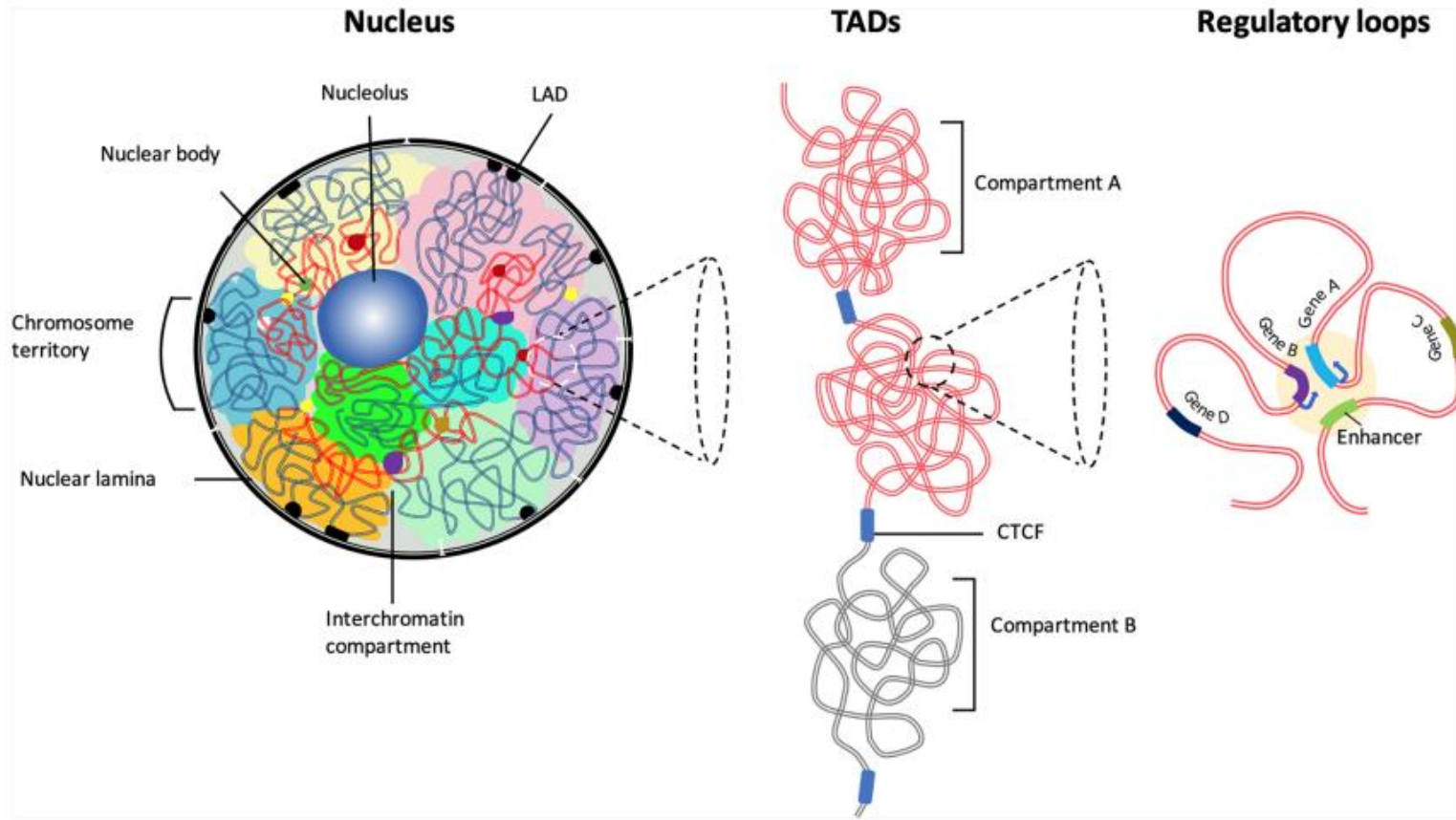
# CAGE

## уточнение старта транскрипции

### CAGE library preparation



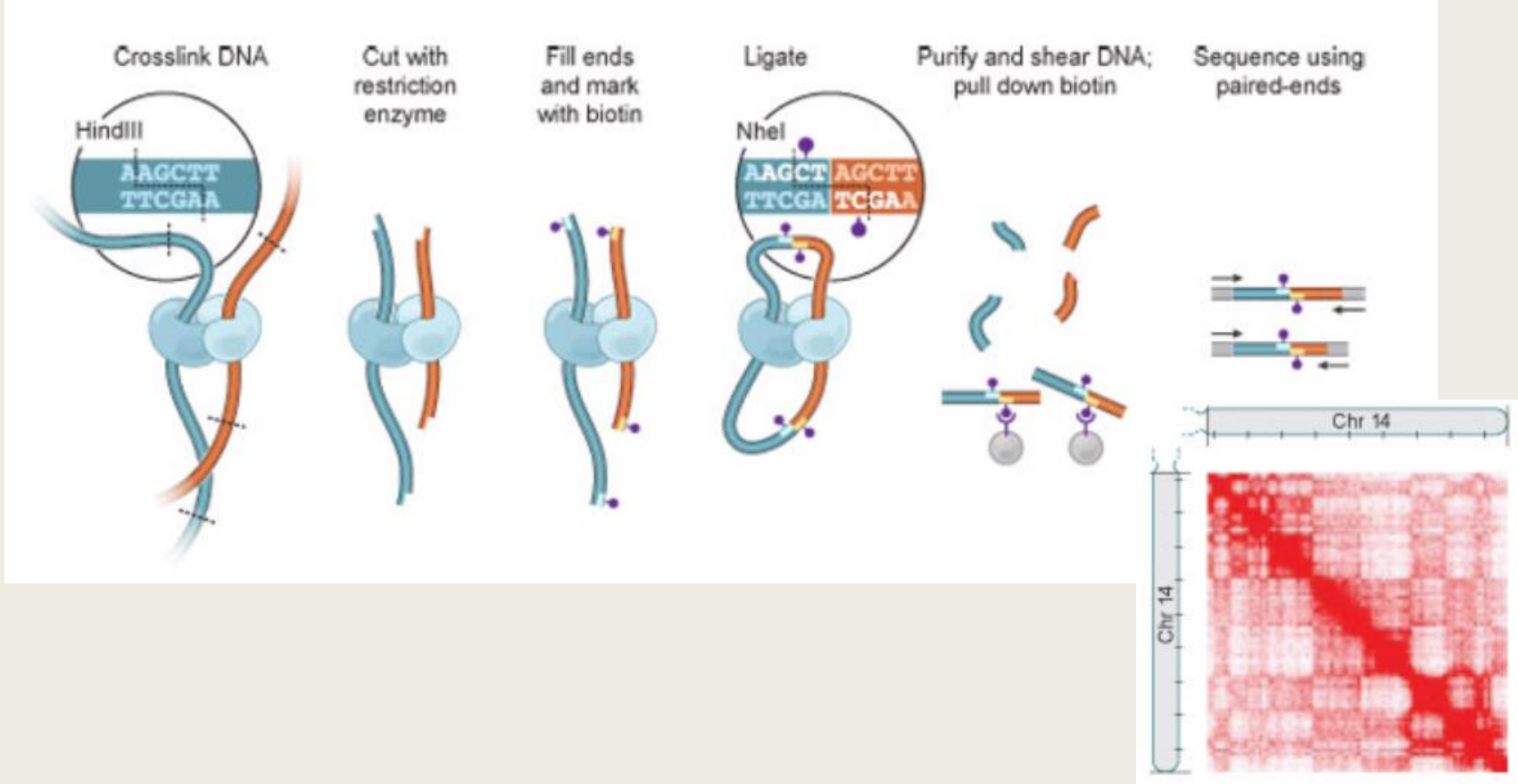
# Укладка ДНК в ядре



# Hi-C - ДНК+ДНК

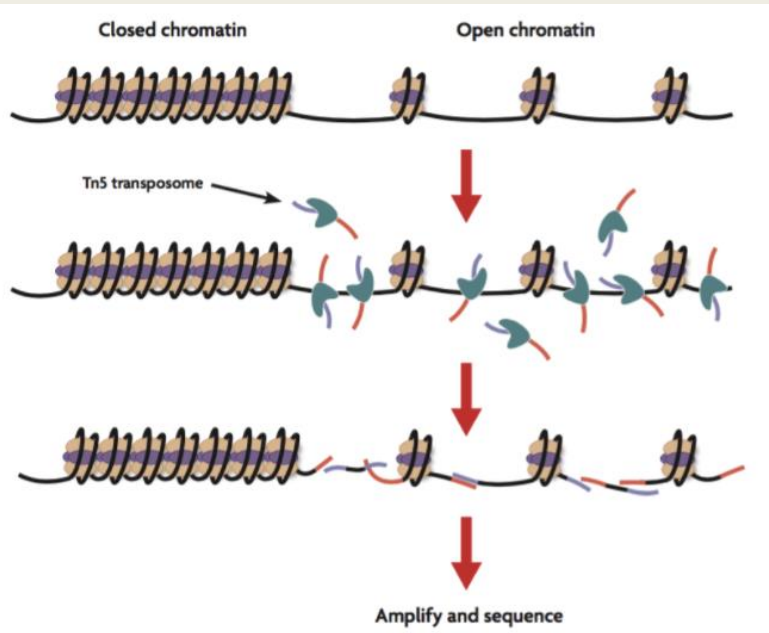
## определение 3D-структуры хроматина

<https://pubmed.ncbi.nlm.nih.gov/20461051/>



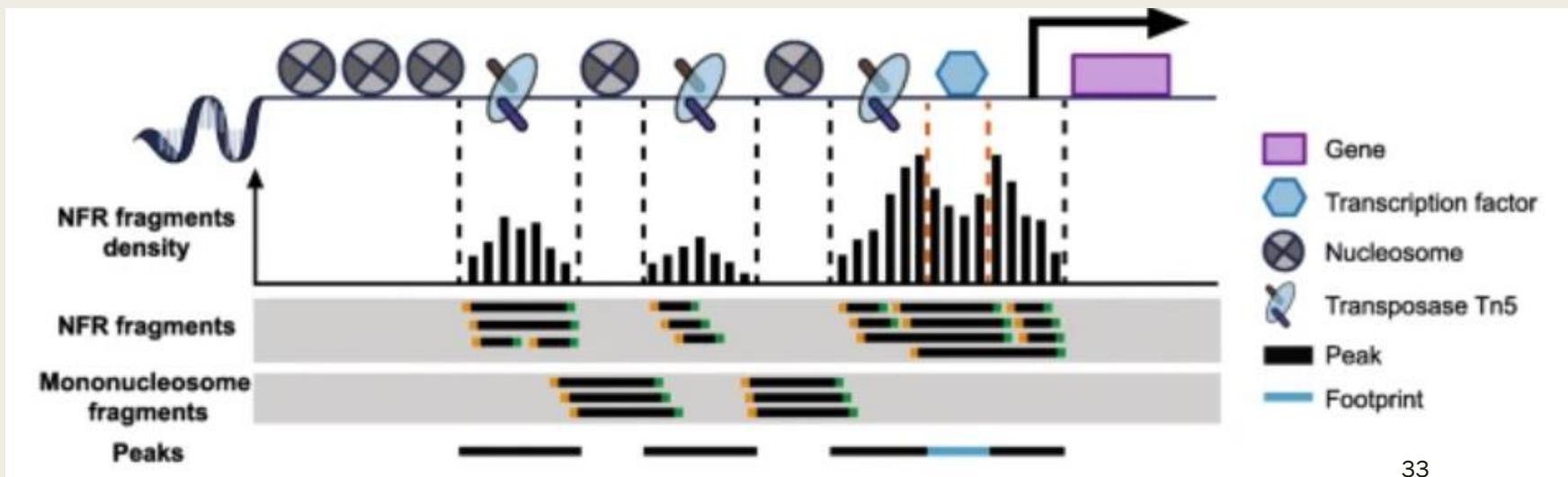


# ATAC-seq - доступный хроматин

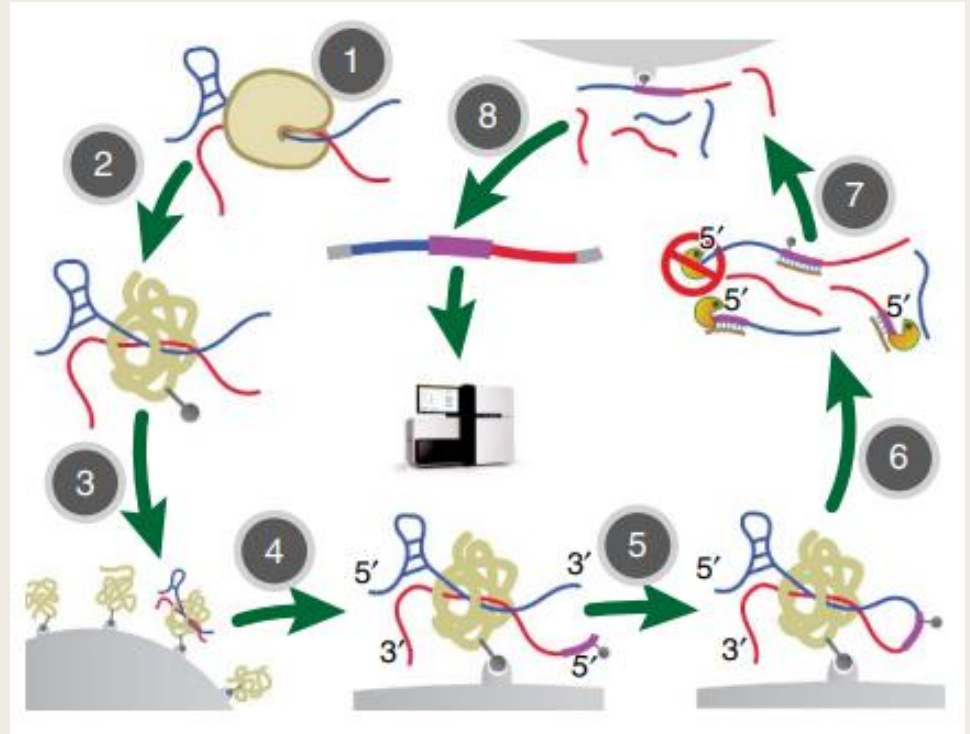


Высокоактивная транспозаза Tn5 вносит двуцепочечные разрывы в открытые участки хроматина и вставляет в области разрывов адаптеры для секвенирования

<https://www.activemotif.com/blog-atac-seq>

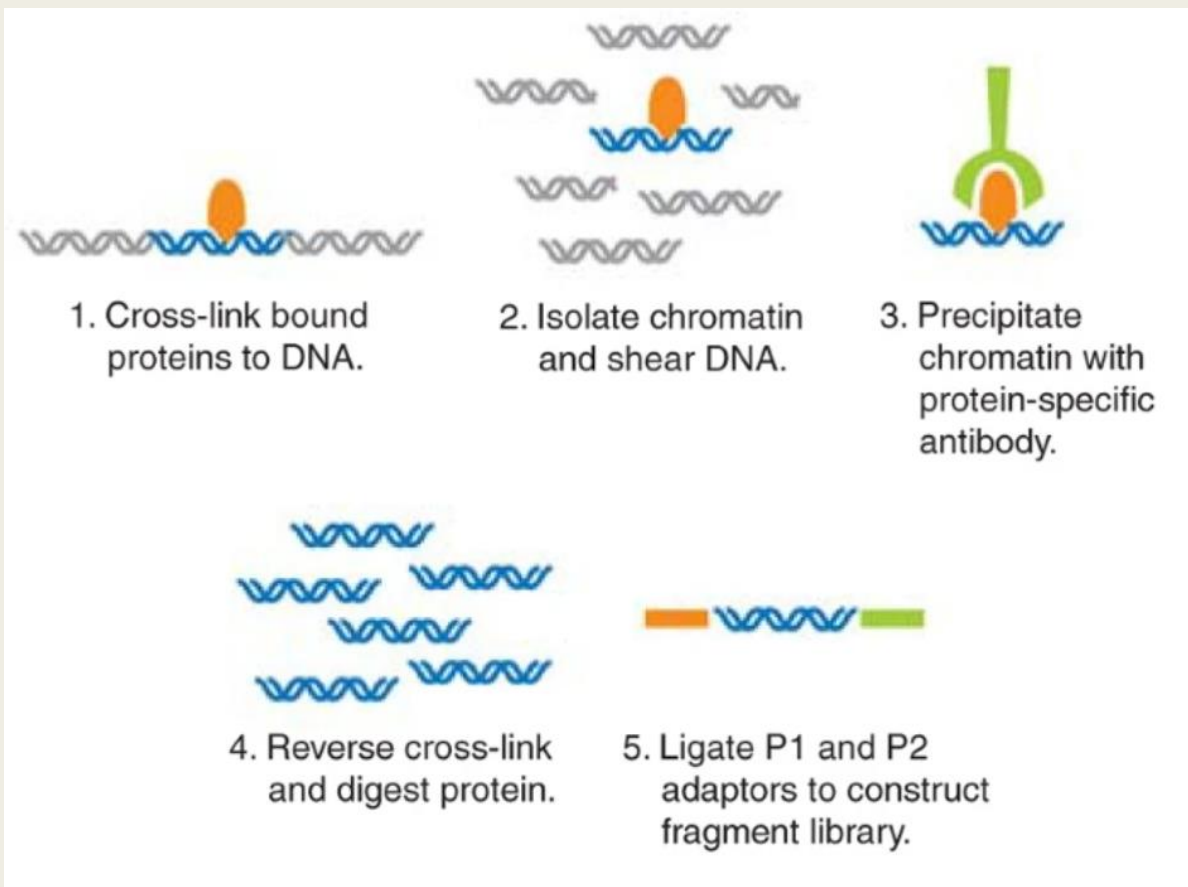


# MARIO – PHK+PHK



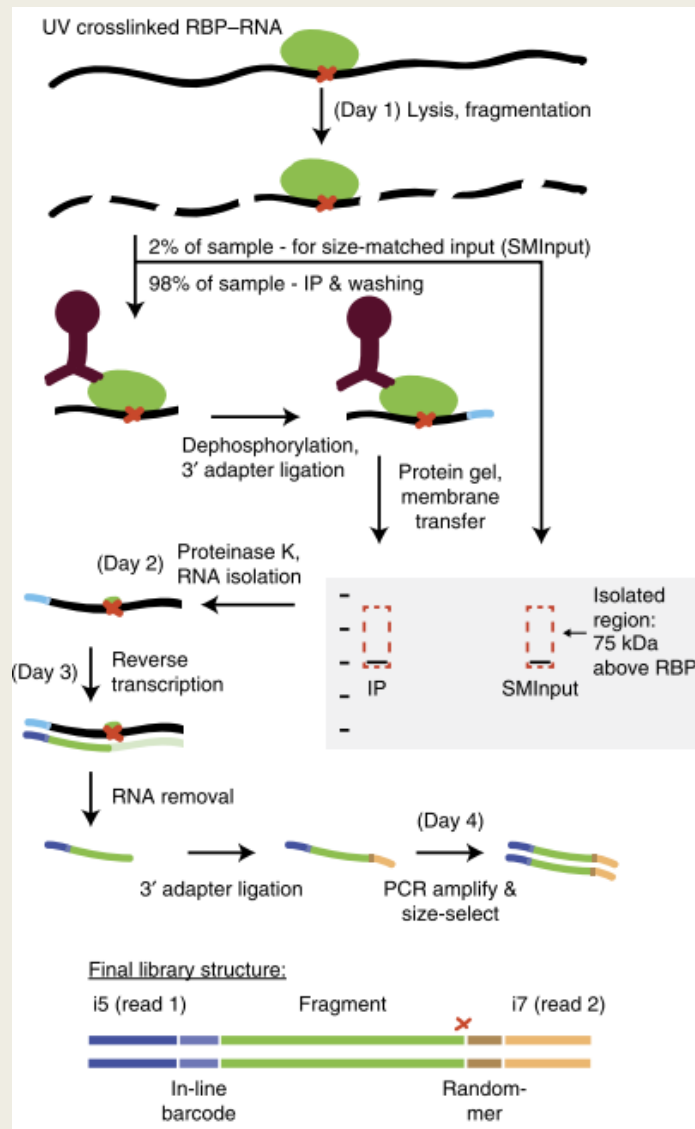
- (1) cross-linking RNAs to proteins;
- (2) RNA fragmentation, protein denaturing and biotinylation;
- (3) immobilization of RNA-binding proteins at low density;
- (4) ligation of a biotinylated RNA linker;
- (5) proximity ligation under a dilute condition;
- (6) RNA purification and RT;
- (7) biotin pull-down;
- (8) construction of sequencing library

# Chip-seq – ДНК+белок

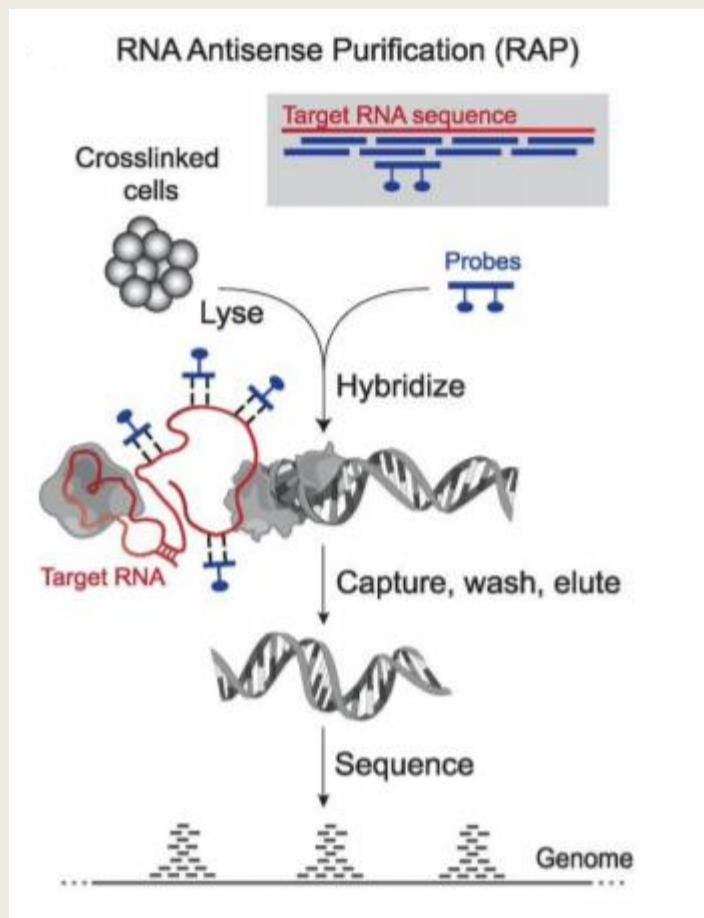


Антитела специфичны к одному белку => за один эксперимент устанавливаем сайты связывания по всему геному, но для одного белка

# eClip – РНК+белок

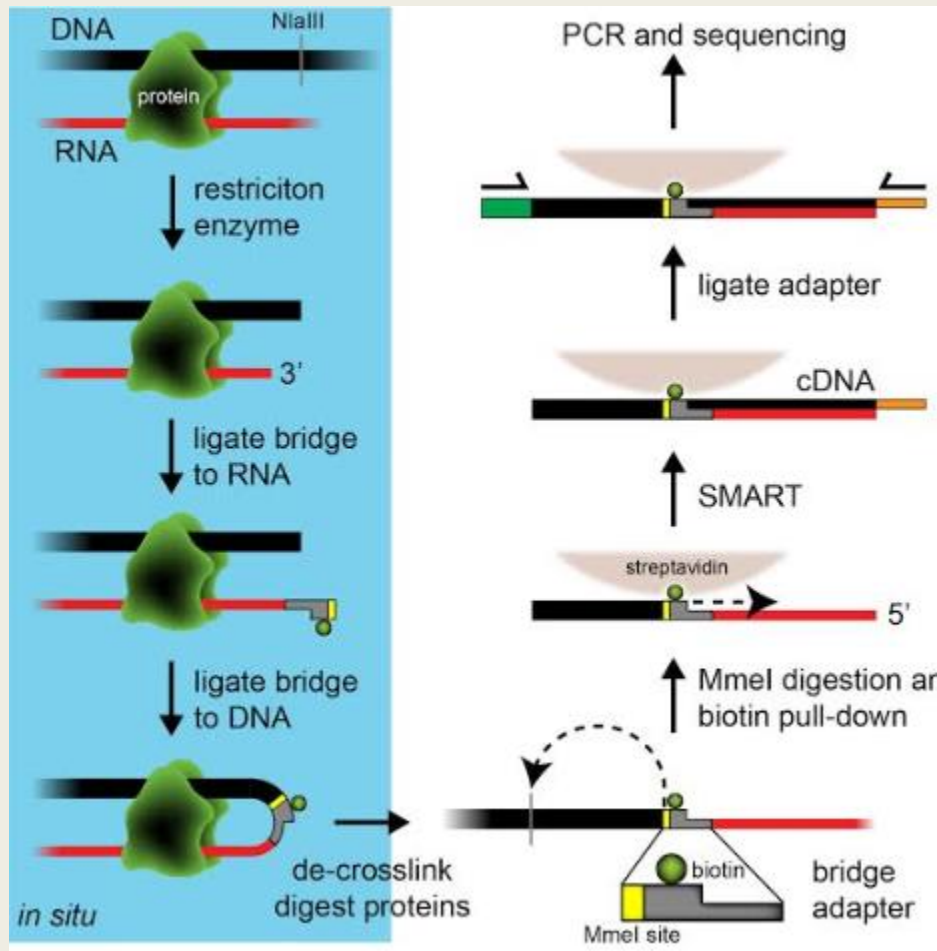


# RAP-seq – РНК+ДНК



Для одной заранее (!) известной РНК устанавливаем локусы взаимодействия с хроматином

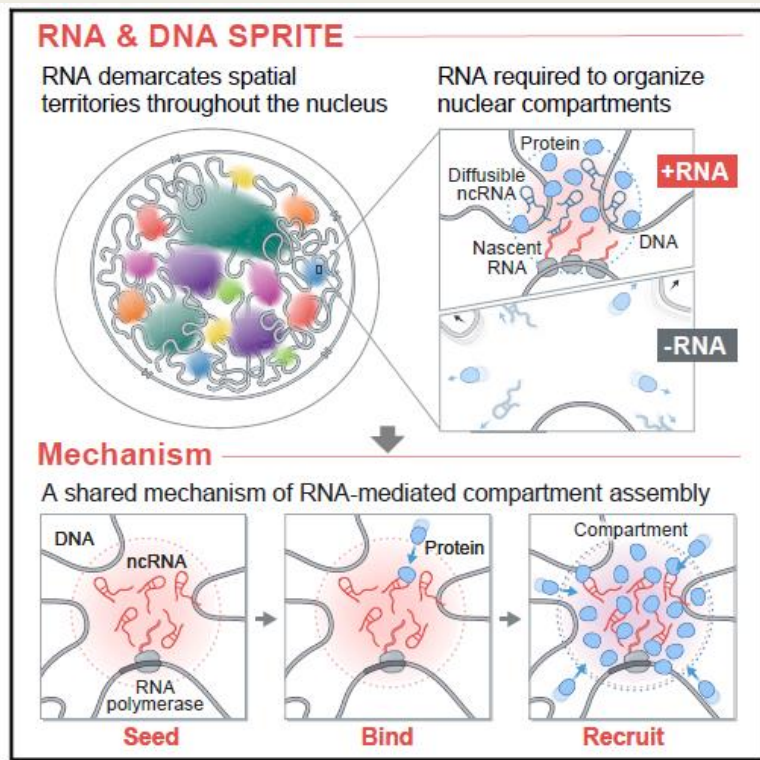
# Red-C – РНК+ДНК



Для всех РНК, независимо от типа и длины, устанавливаем все локусы взаимодействия с хроматином

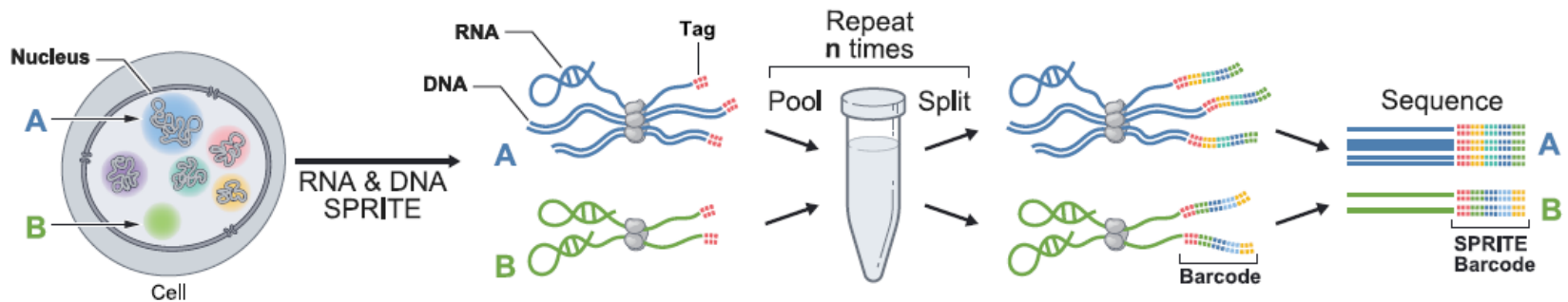
Есть другие методы: GRID-seq, RADICL-seq, ...

# RD-SPRITE – РНК+ДНК



Можно детектировать крупные макромолекулярные комплексы, ядерные структуры

Нет стадии лигирования близко расположенных молекул

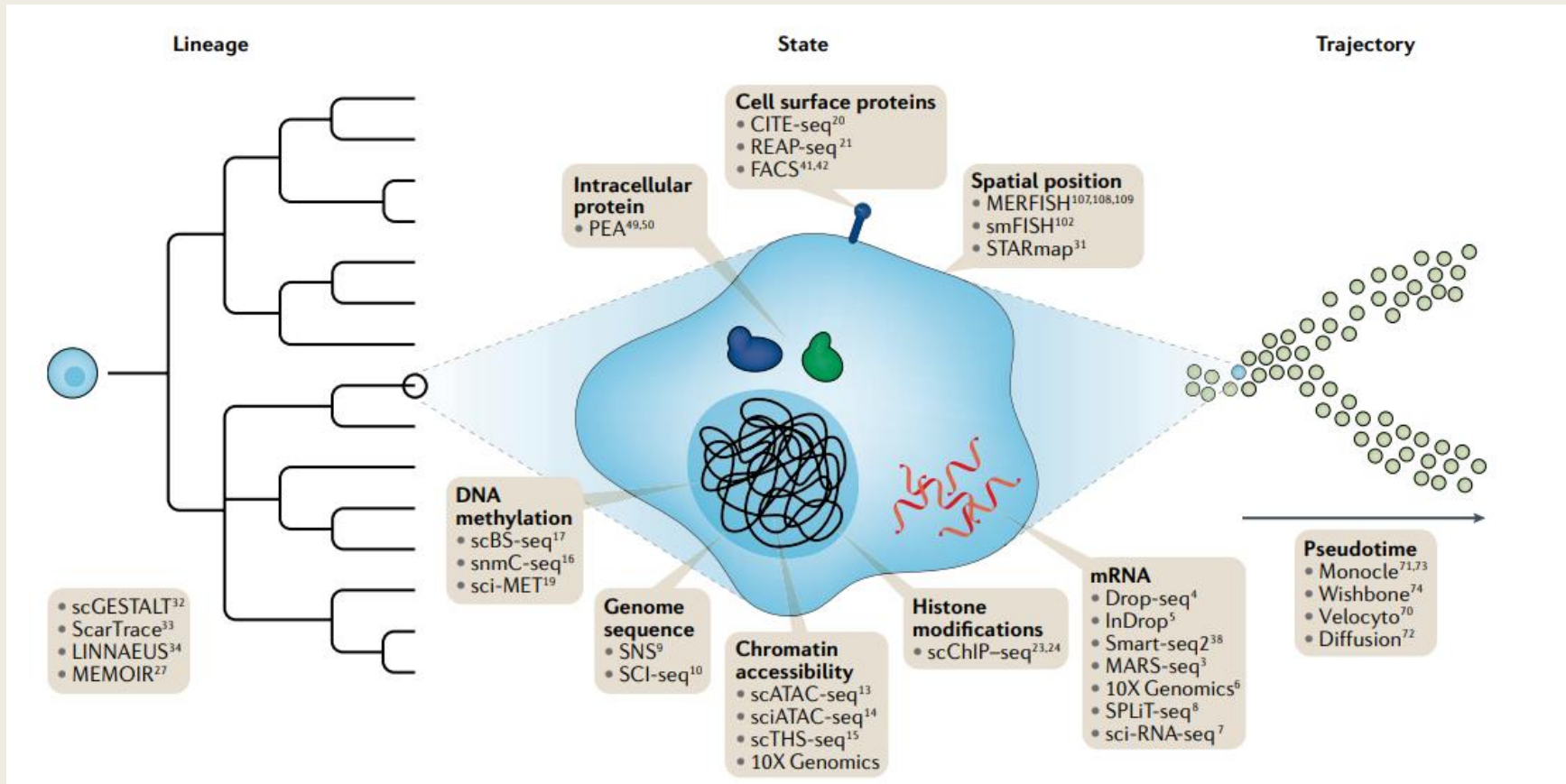


# enSEQlopedia

- <http://enseqlopedia.com/enseqlopedia/>
- Описание большого количества протоколов
- Есть далеко не все



# А еще есть single cell методы



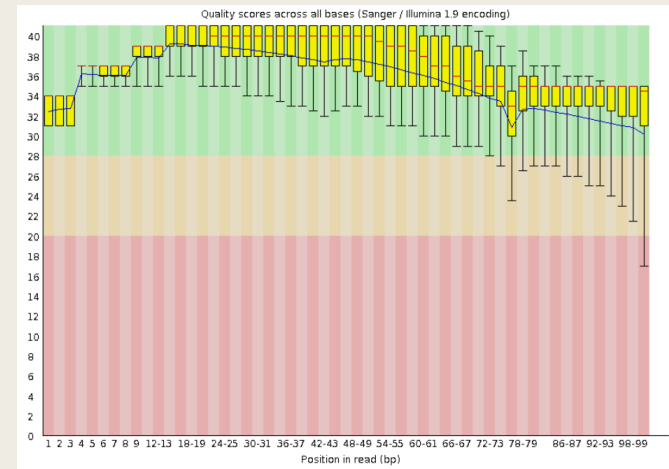
# Важно

- У всех методов есть свои ограничения, особенности, «подводные камни», источники «шума», вариации протокола
- Необходимо понимать природу данных (организм, ткань, клеточная линия, процесс выделения и обработки)
- Разобраться во всех этапах пробоподготовки и контролях

# Как обрабатывать данные?

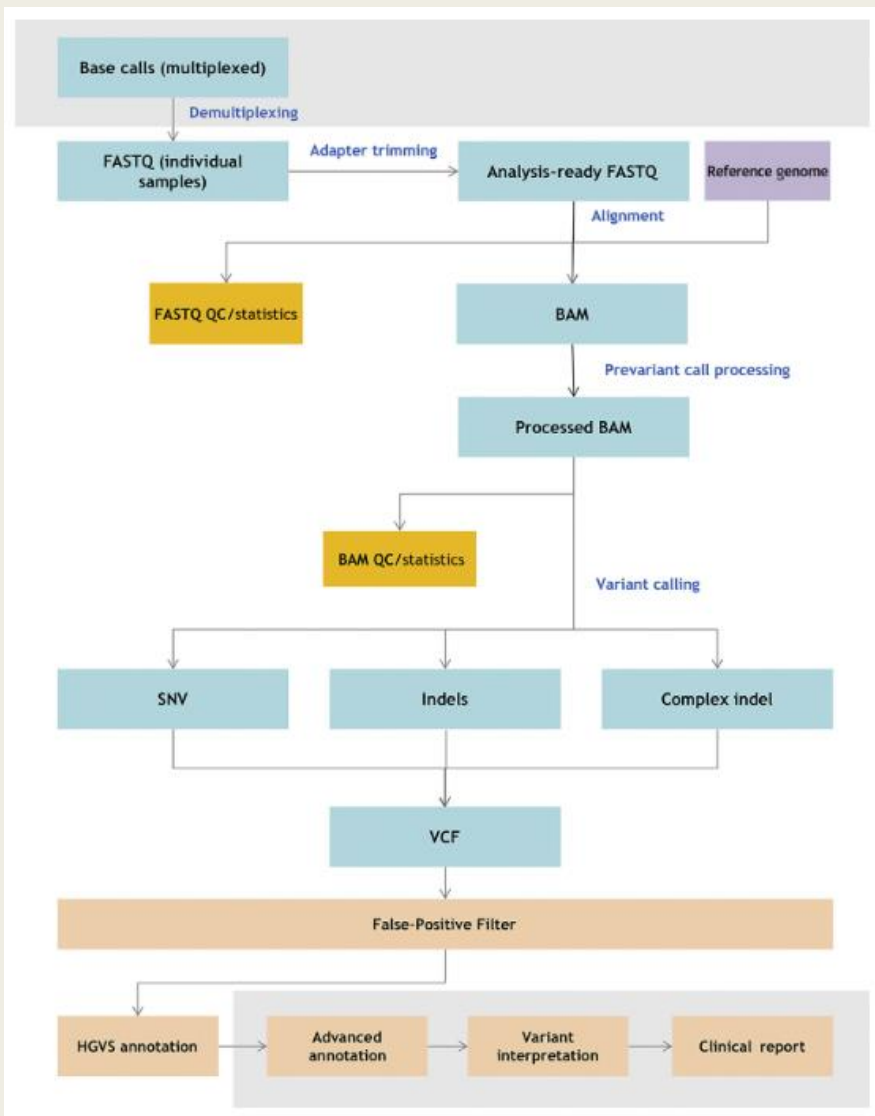
- Говорим только о Illumina
- В любом случае получаем чтения или «риды»
- Для каждого типа секвенирования разрабатывают свой подход биоинформатической обработки
- Есть общие шаги в обработке чтений

# Всегда нужно



- Проверить качество чтений
- Проверить наличие адаптеров или других технических последовательностей в чтениях
- *Картировать чтения на референсный геном или собрать референсный геном \ транскриптом*
- *Проверить качество сборки или картирования*

# Программный конвейер



Каждый блок – шаг обработки,  
отдельная программа

Не все эксперименты одинаковые

«Стандартный» протокол обработки  
может не подойти

# Benchmarks

Практически для каждого шага программного конвейера можно найти несколько похожих программ

Нужно быть в курсе актуального ПО

