

# Коррекция ошибок в рядах

Алгоритмы в биоинформатике

Мелешко Дмитрий  
[meleshko.dmitrii@gmail.com](mailto:meleshko.dmitrii@gmail.com)

# Что было в прошлом модуле

- Сравнение двух последовательностей: расстояния, глобальное и локальное выравнивания, штрафы за гэпы, эффективное использование памяти
- Поиск специфичных участков генома при помощи HMM
- Сравнение многих последовательностей между собой и одной последовательности со многими
- Выравнивание на референсный геном (BWT, BWA)
- Перестройки в геноме, синтенные

# Что будет в этом модуле

- Секвенирование! NGS данные. Артефакты и важная информация.
- Сборка генома из коротких прочтений.
- Сборка многих геномов. Метагеном, гаплотипы и связанные задачи.
- Эволюция и ее параметры. Зачем нужны вероятностные модели?
- Вторичная структура РНК

# В этой лекции

- Методы секвенирования
- Секвенирование как случайный процесс
- Ошибки секвенирования
- Способы отбрасывать риды с ошибками
- Способы коррекции ошибок

# Технологии секвенирования

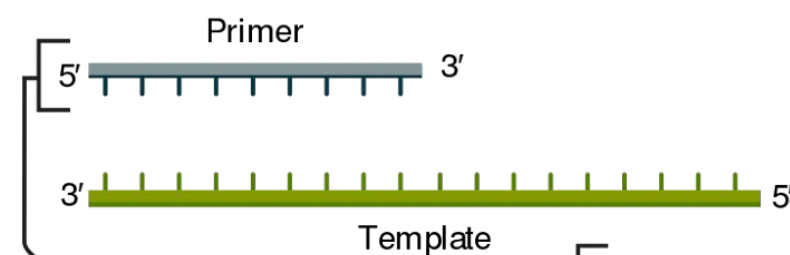
# Технологии секвенирования

- ● Sanger
- Illumina/MGI
- Nanopore
- PacBio
- 10X Genomics/stLFR

# Технологии: Sanger

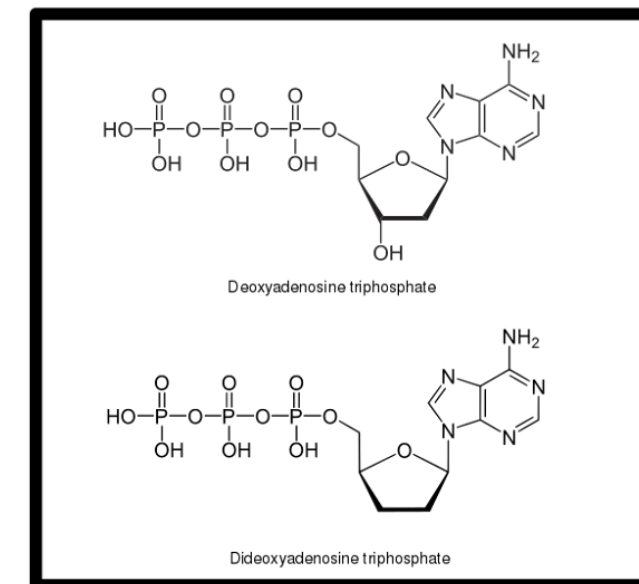
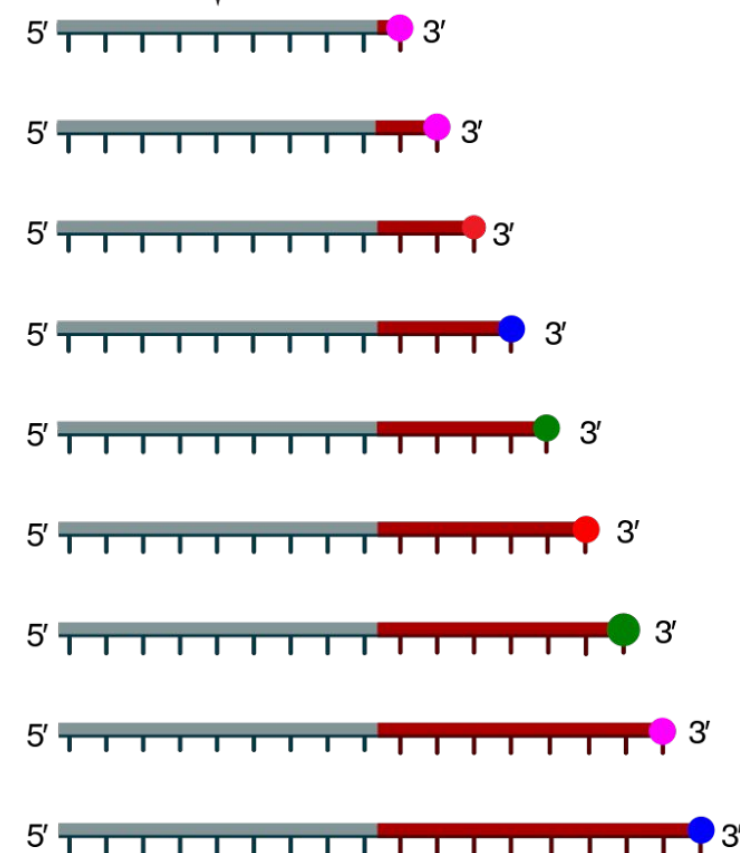
## ① Reaction mixture

- ▶ Primer and DNA template
- ▶ DNA polymerase
- ▶ ddNTPs with flourochromes ▶ dNTPs (dATP, dCTP, dGTP, and dTTP)

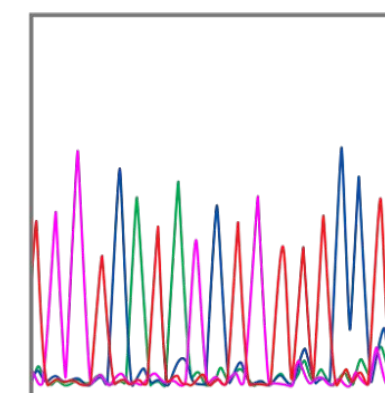
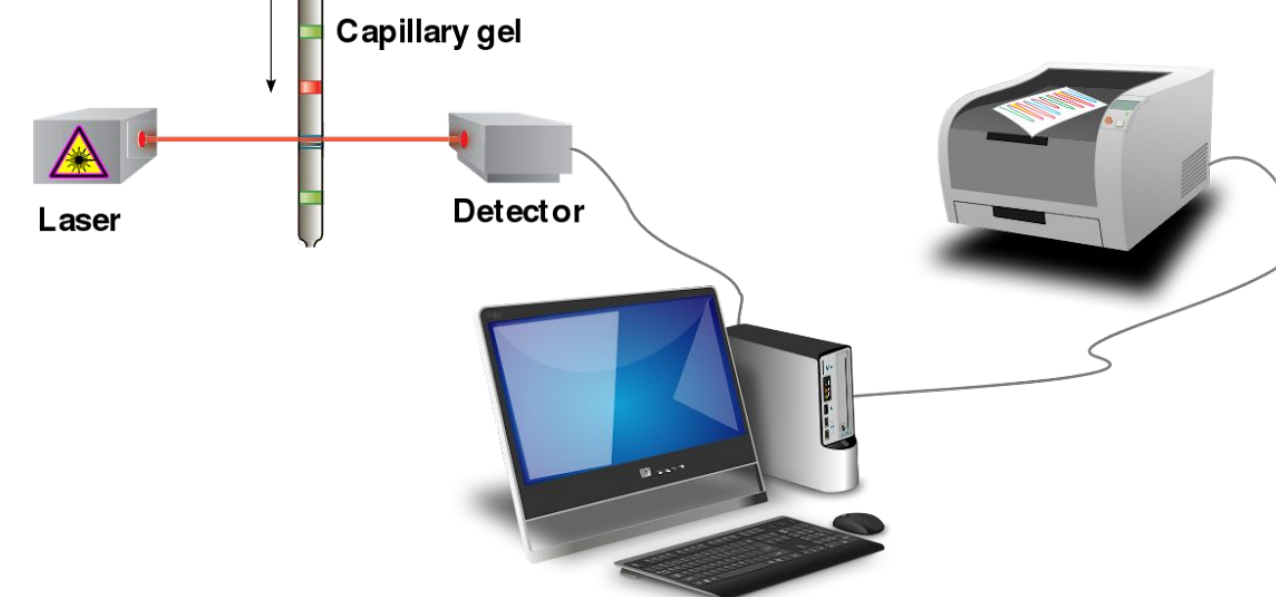


ddNTPs  
ddTTP ●  
ddCTP ●  
ddATP ●  
ddGTP ●

## ② Primer elongation and chain termination



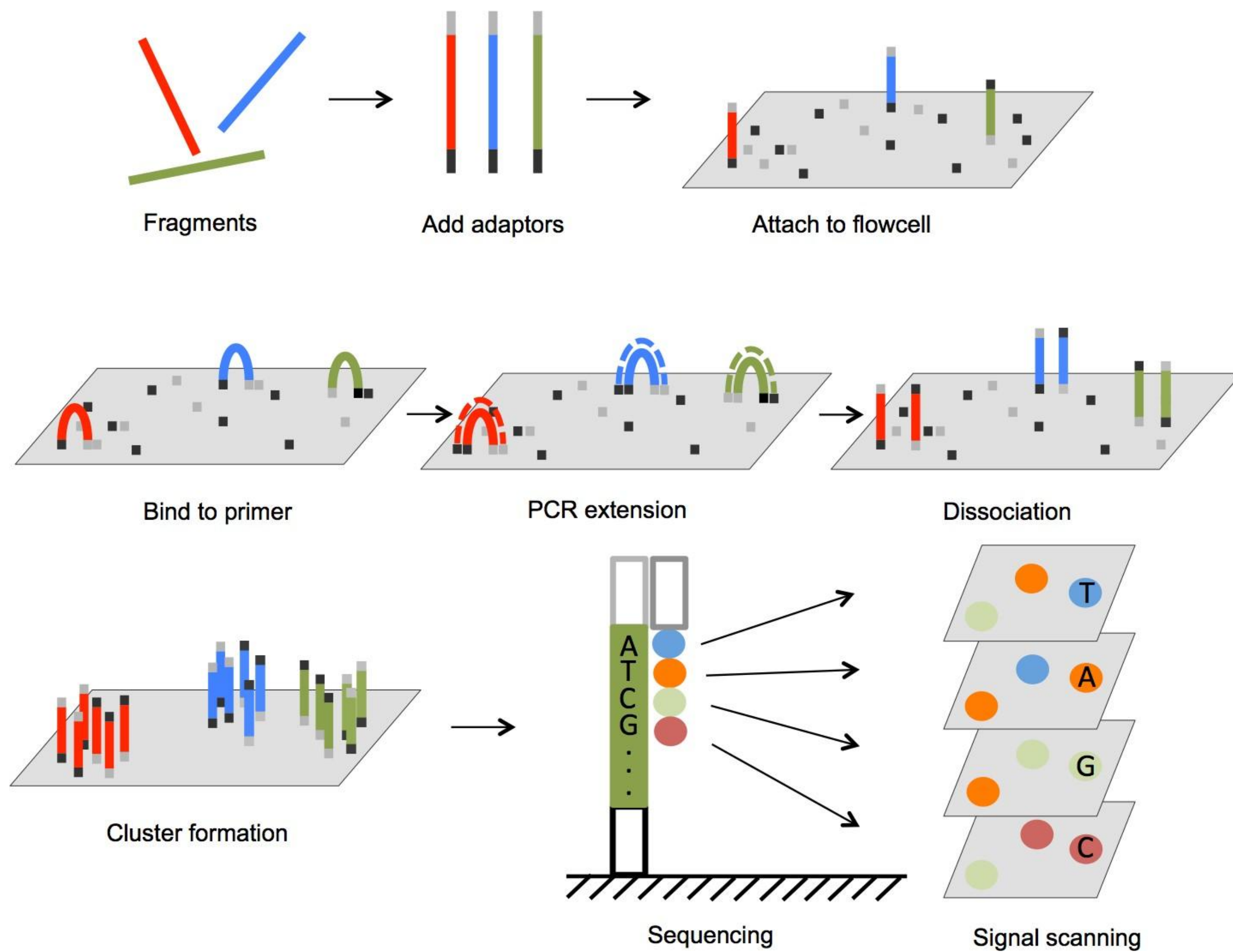
## ③ Capillary gel electrophoresis separation of DNA fragments



Chromatograph

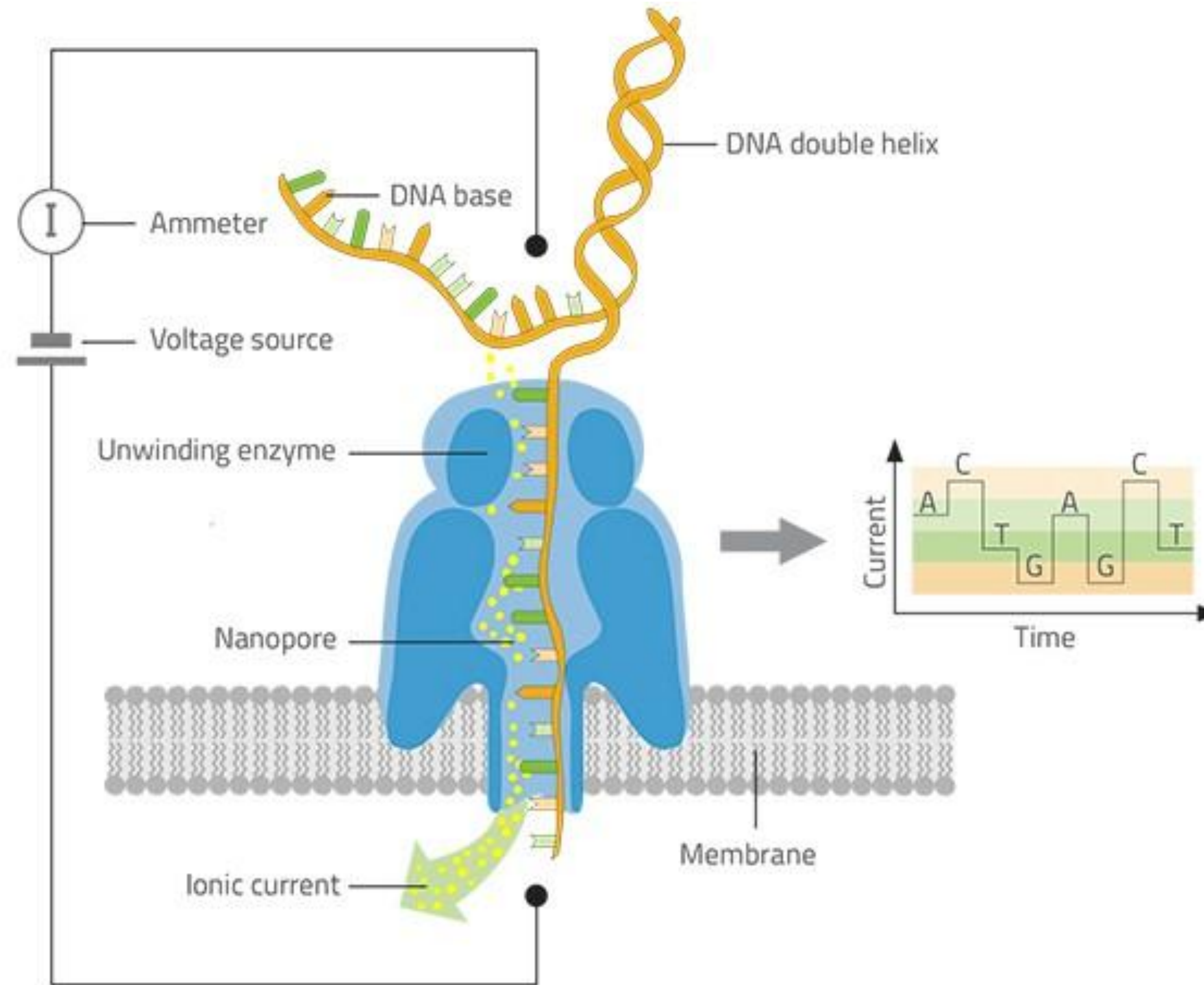
## ④ Laser detection of flourochromes and computational sequence analysis

# Технологии: Illumina

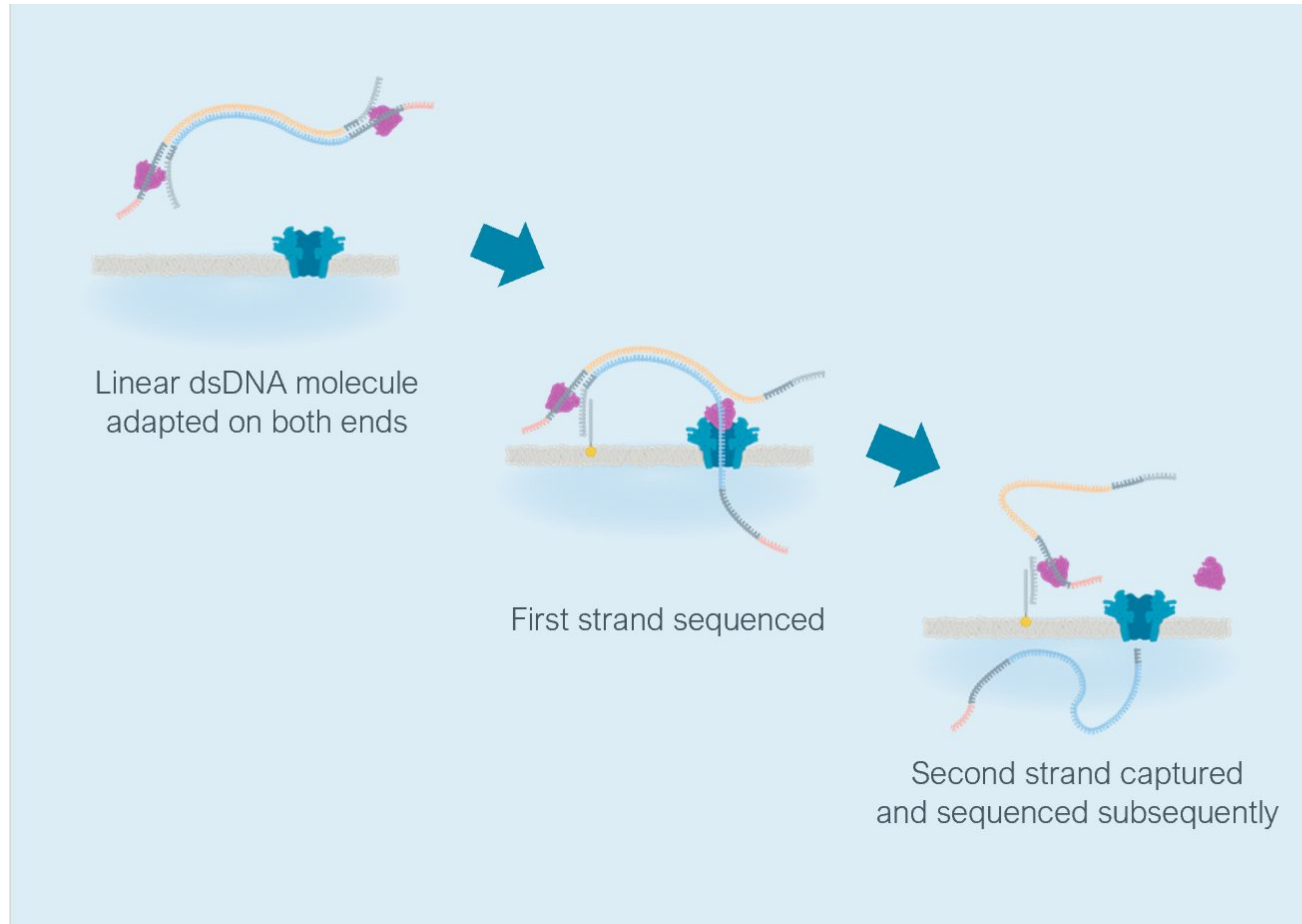




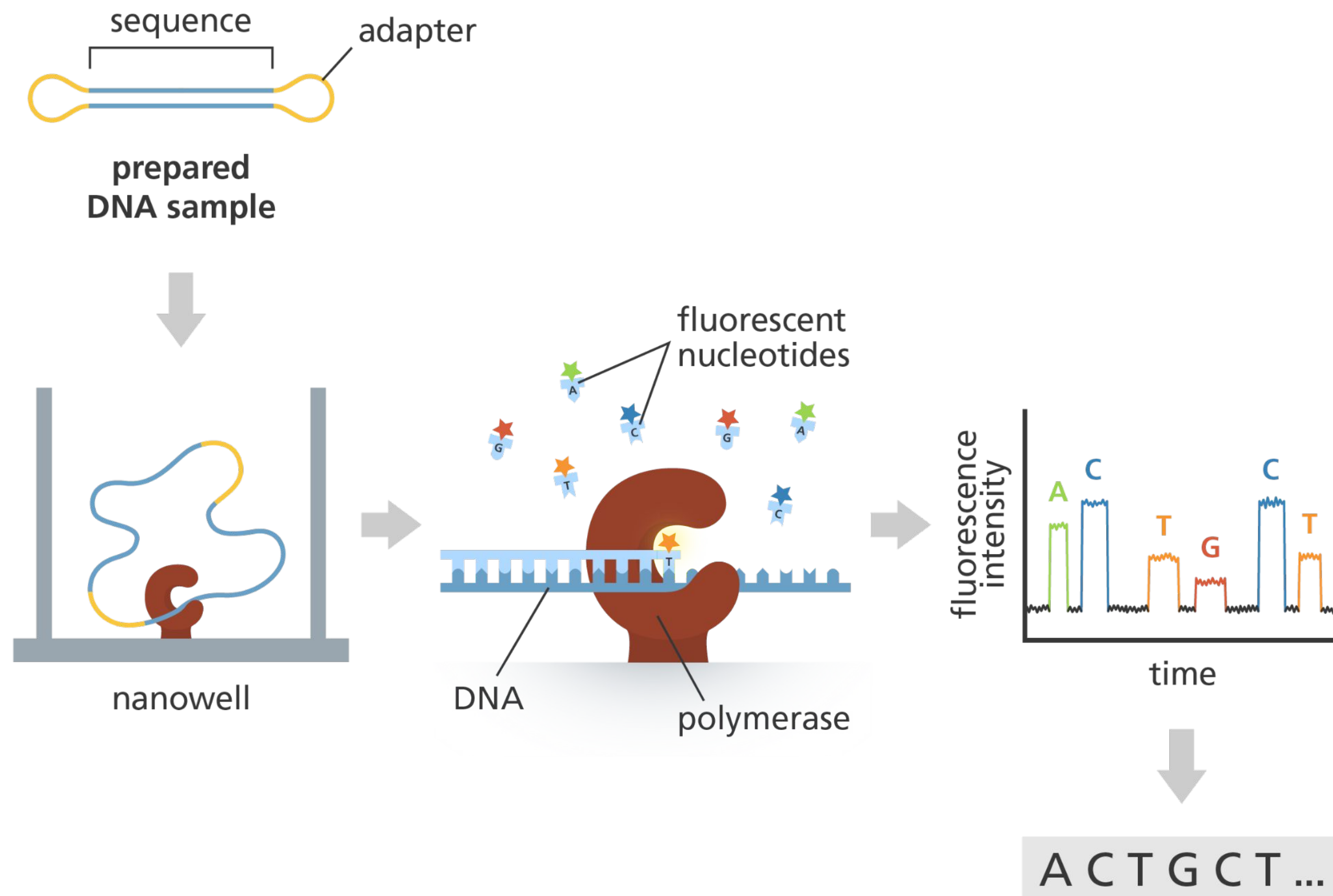
# Технологии: Nanopore



# Технологии: Nanopore Duplex



# Технологии: PacBio

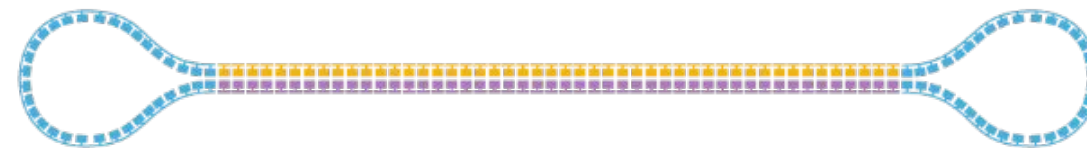


# Технологии: PacBio Hi-Fi

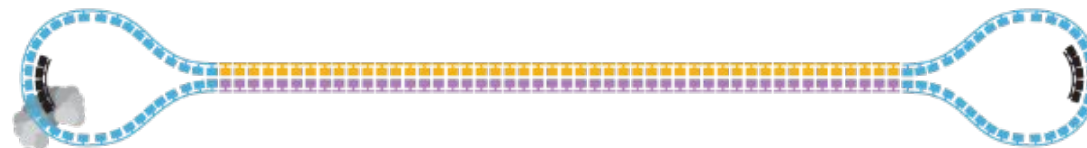
Start with high-quality double stranded DNA



Ligate SMRTbell adapters and size select



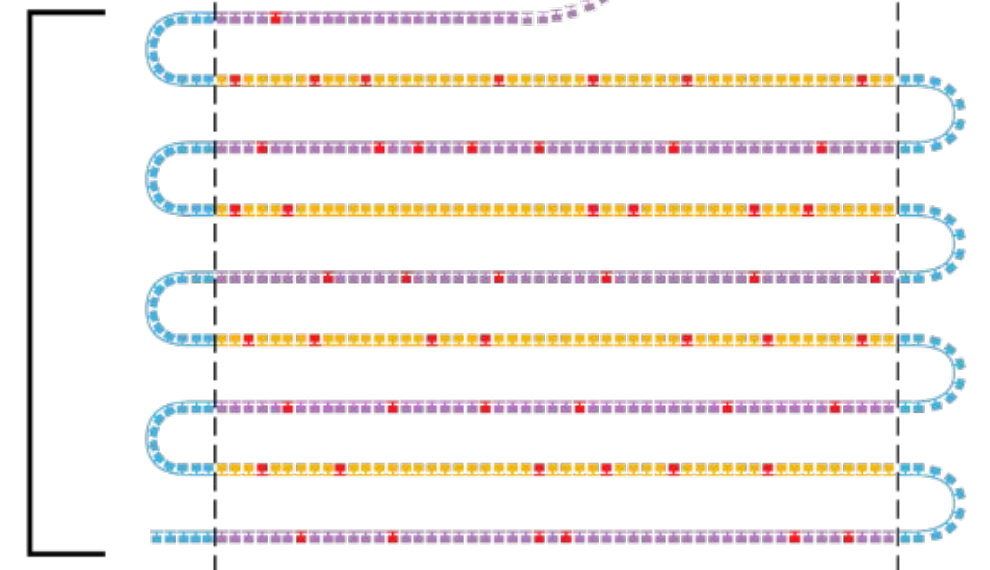
Anneal primers and bind DNA polymerase



Circularized DNA is sequenced in repeated passes



The polymerase reads are trimmed of adapters to yield subreads



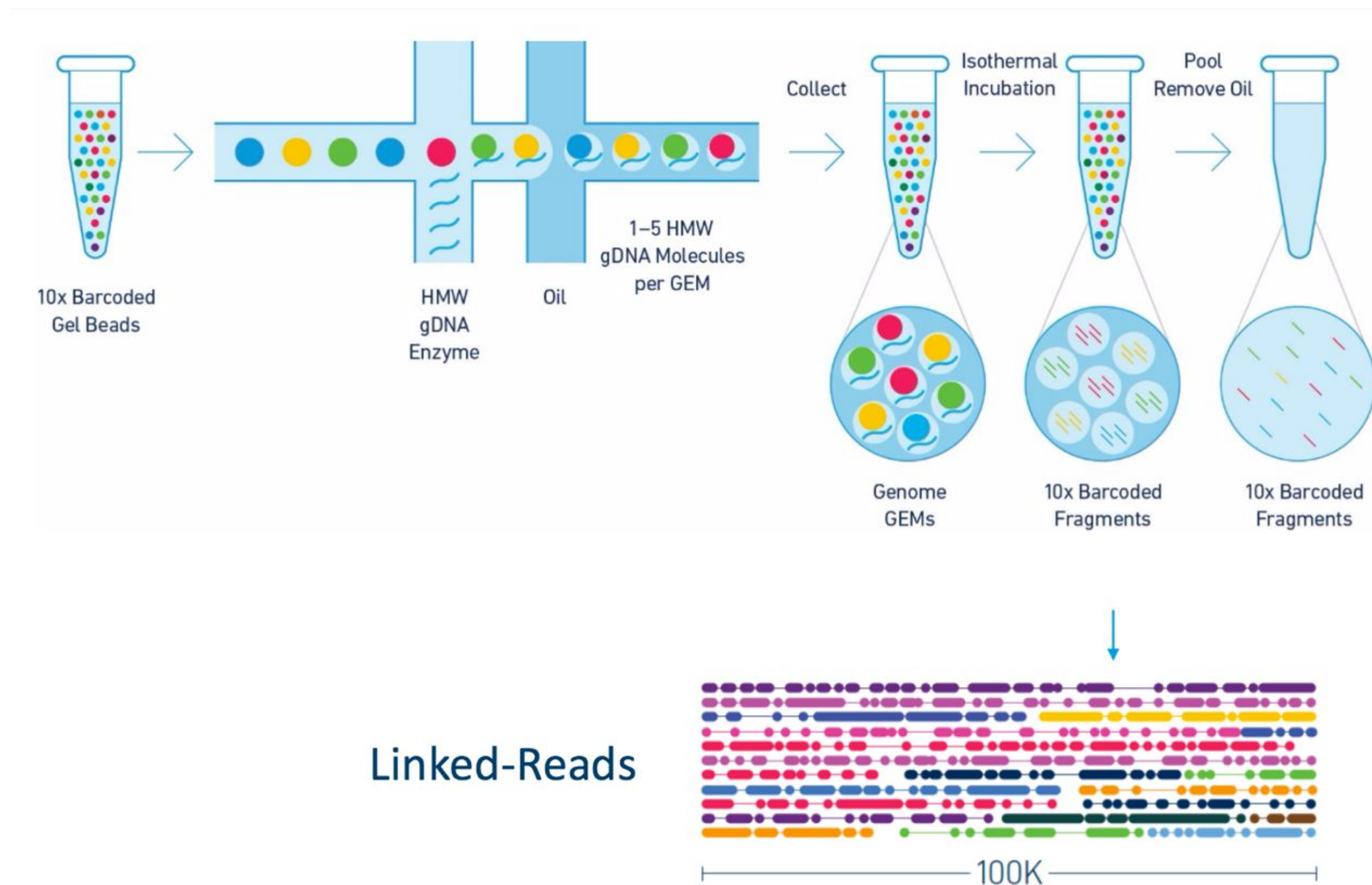
Consensus is called from subreads



**HiFi READ**  
(>99% accuracy)



# Технологии: 10X Genomics



# Что можно секвенировать?

-

# Что можно секвенировать?

- Whole Genome Sequencing

## Whole Genome Sequencing

30-60x Coverage



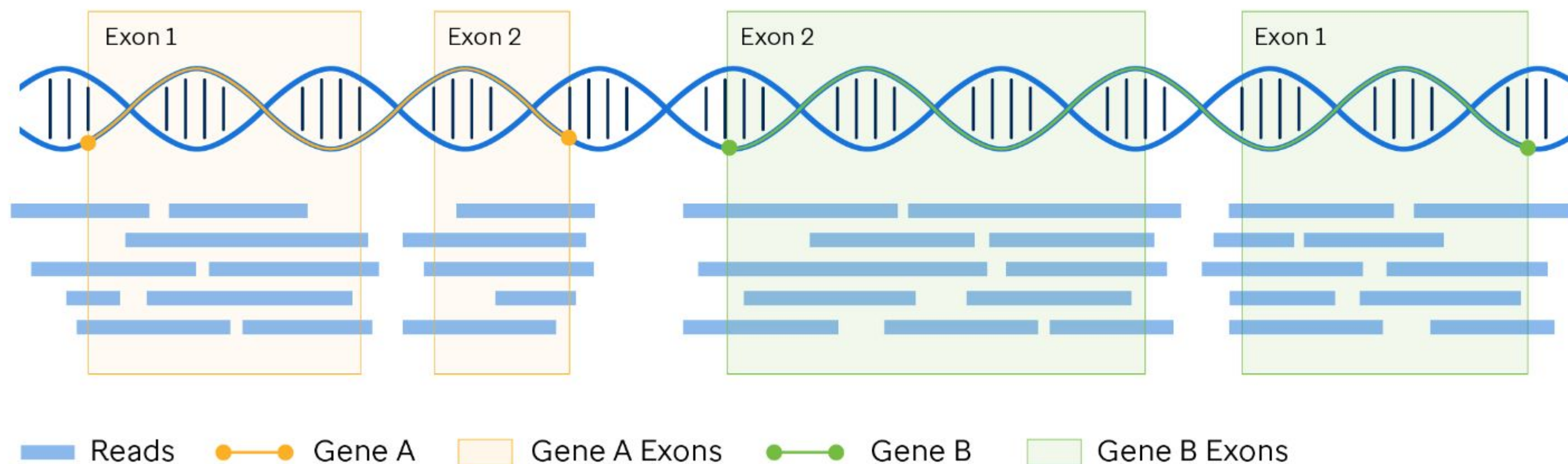
— Reads

# Что можно секвенировать?

- Whole Exome Sequencing

## Whole Exome Sequencing

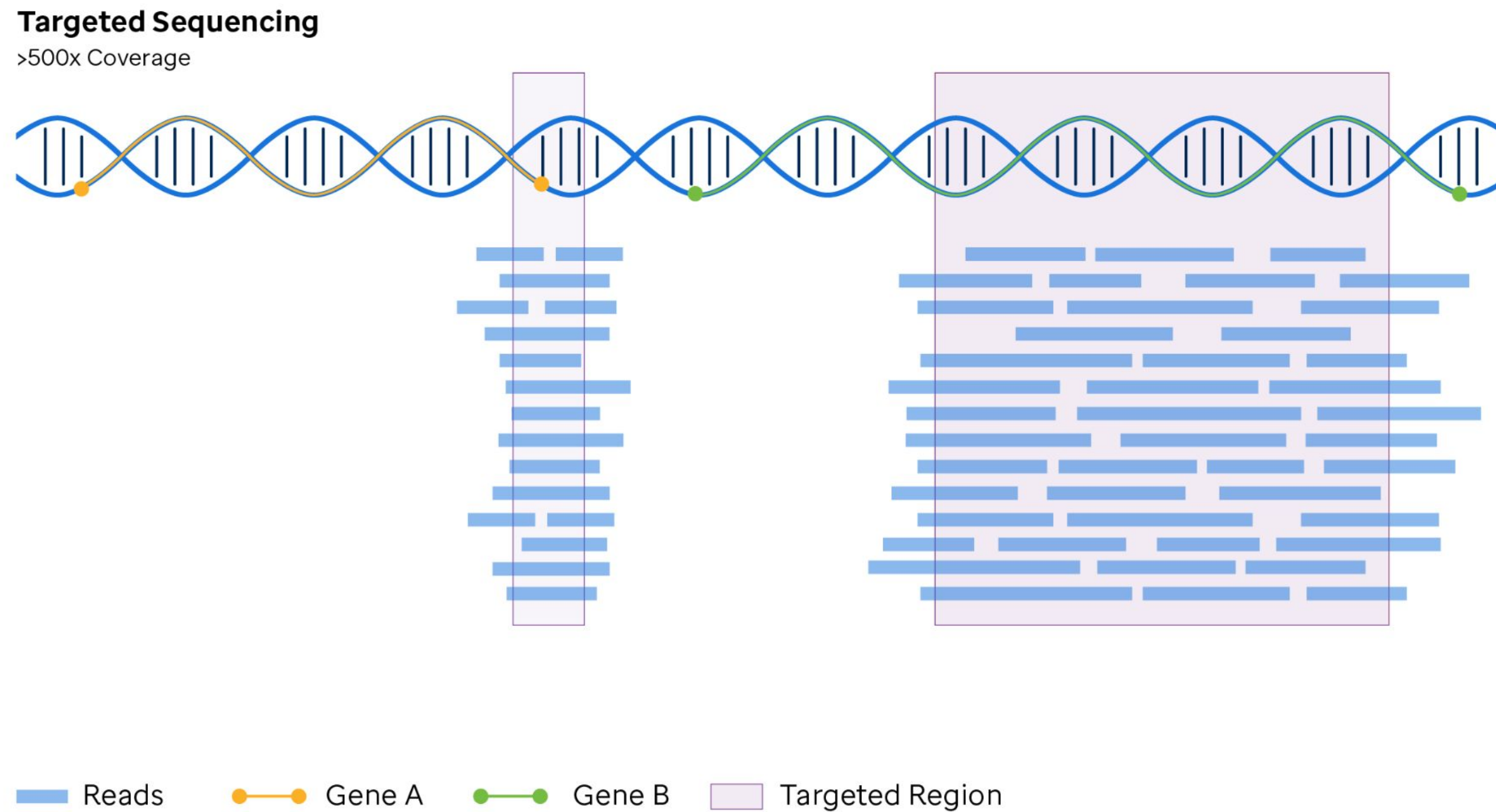
50-100x Coverage





# Что можно секвенировать?

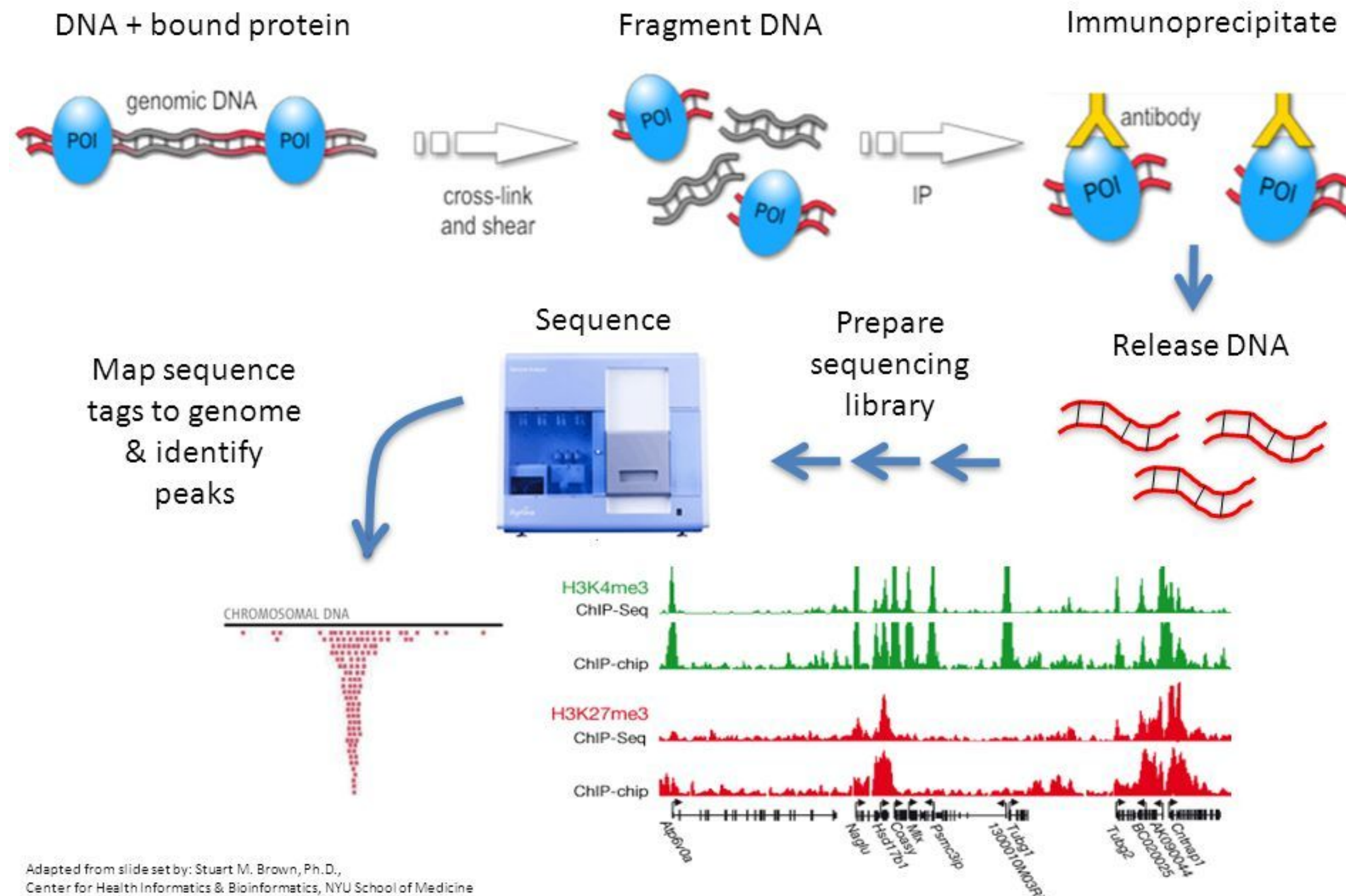
## ○ Targeted Sequencing



# Что можно секвенировать?

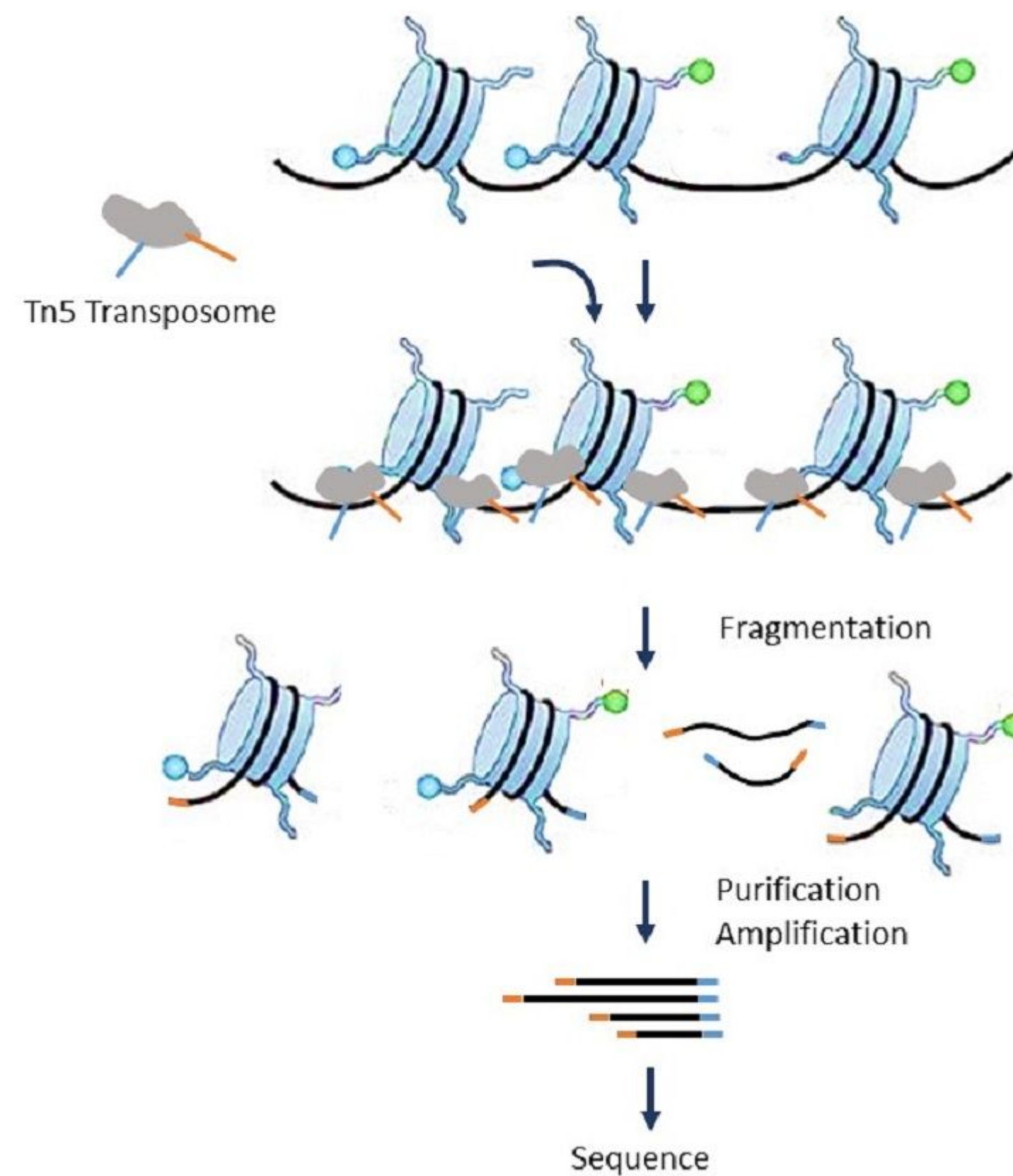
- ChIP-Seq

## ChIP-seq overview



# Что можно секвенировать?

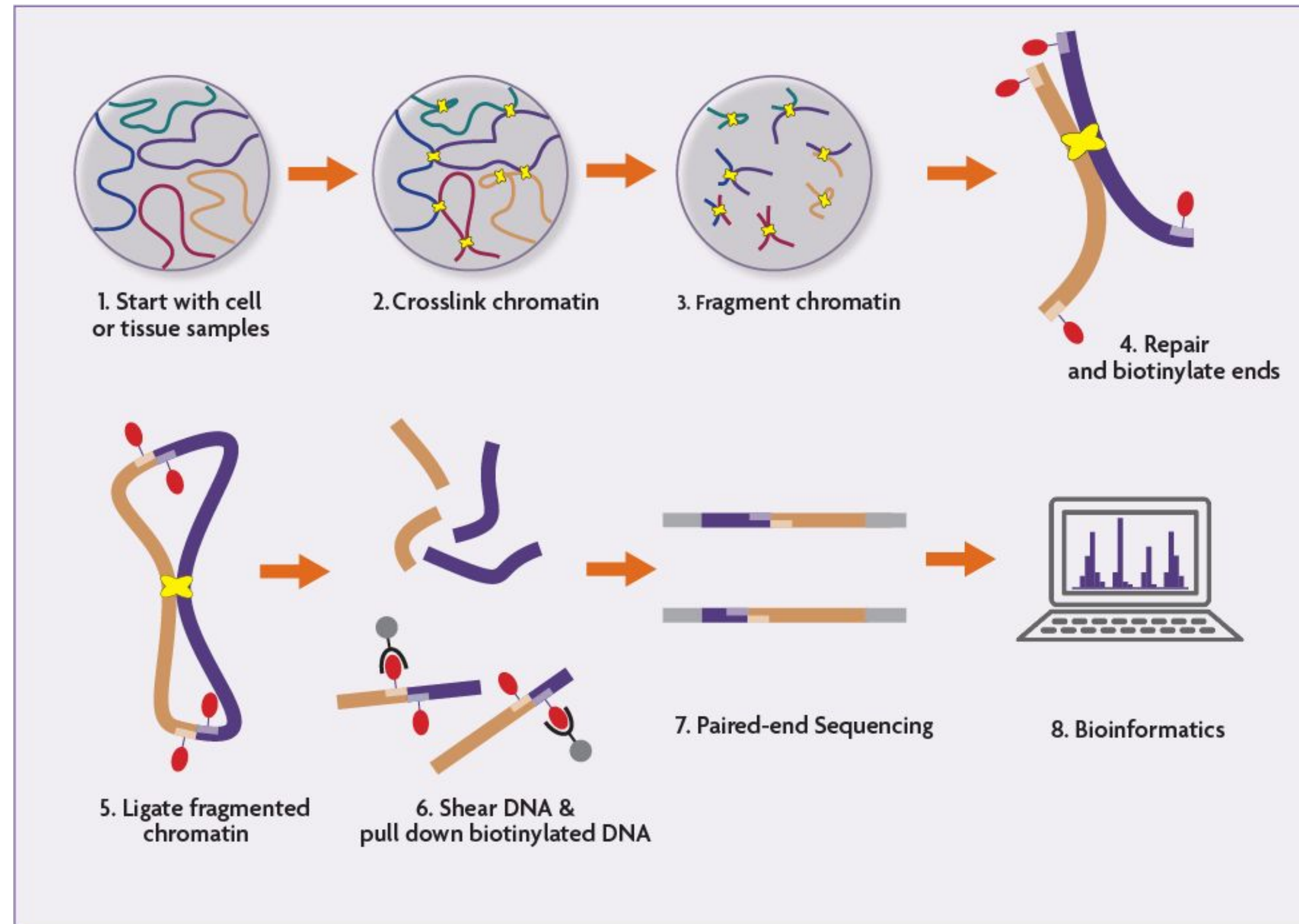
- ATAC sequencing





# Что можно секвенировать?

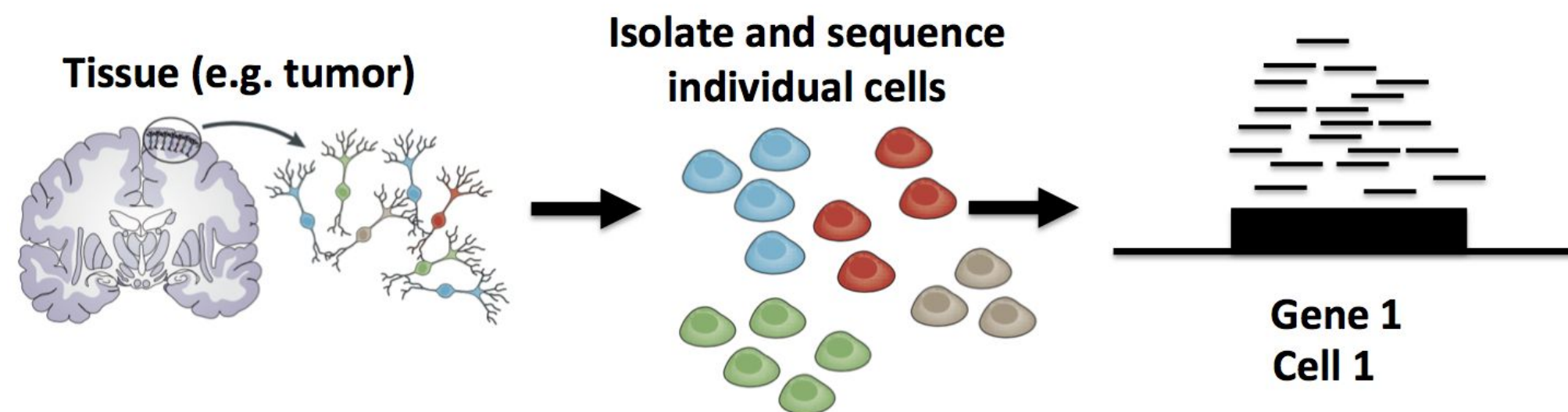
- Hi-C



# Что можно секвенировать?

- Single-Cell Sequencing

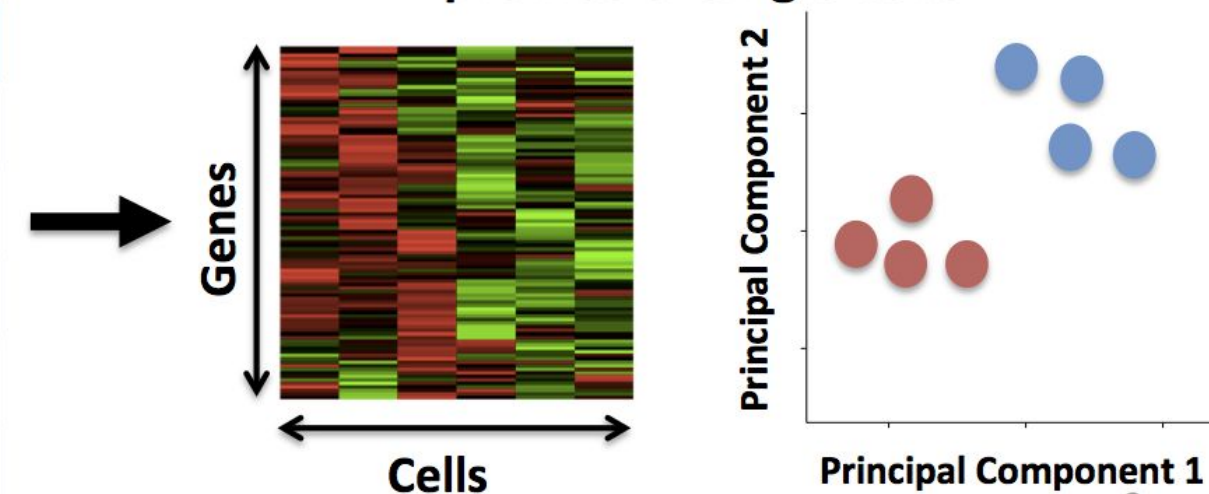
## Single-cell RNA-Seq (scRNA-Seq)



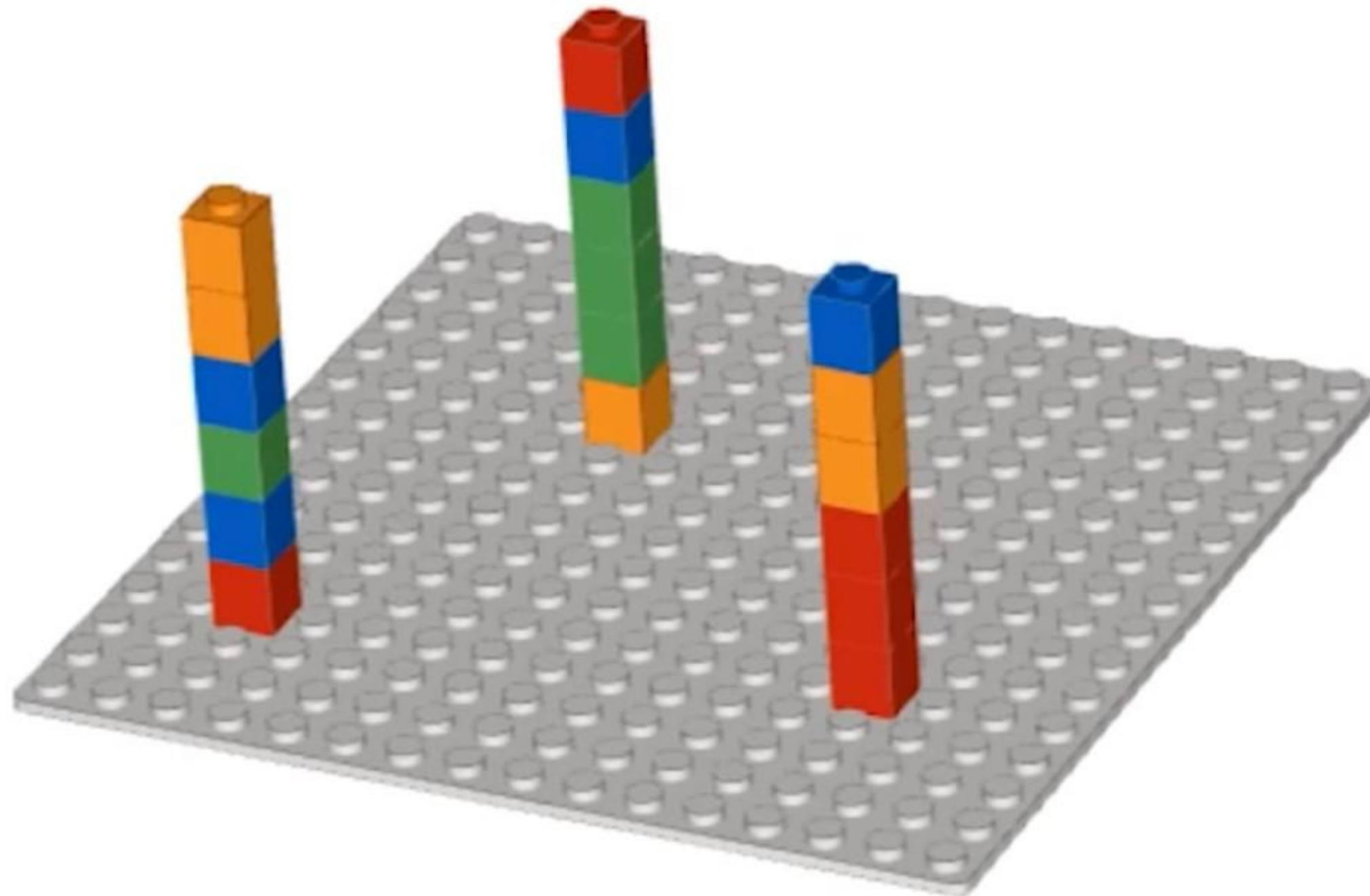
Read Counts

	Cell 1	Cell 2	...
Gene 1	18	0	
Gene 2	1010	506	
Gene 3	0	49	
Gene 4	22	0	
...			

Compare gene expression profiles of single cells

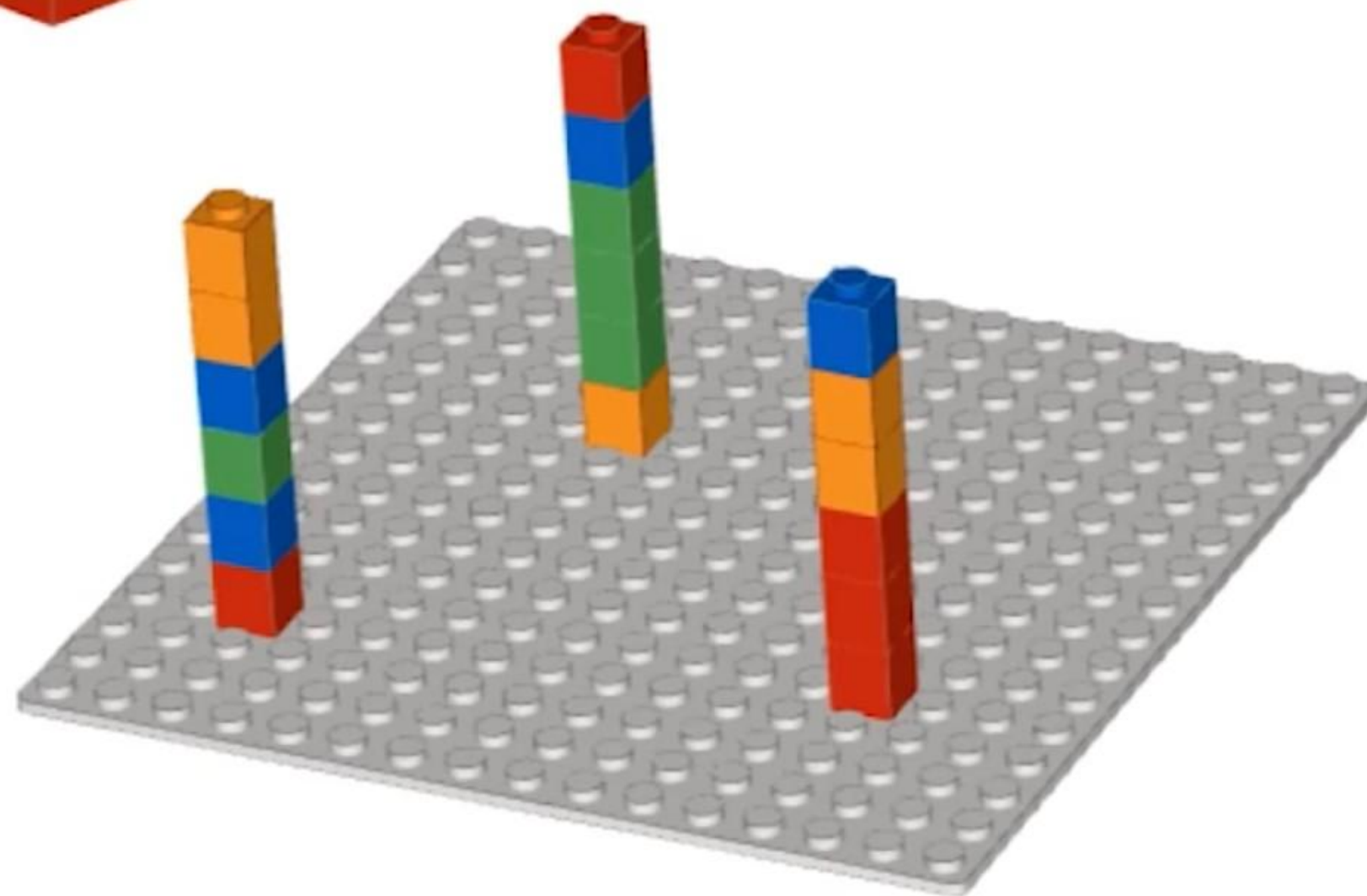
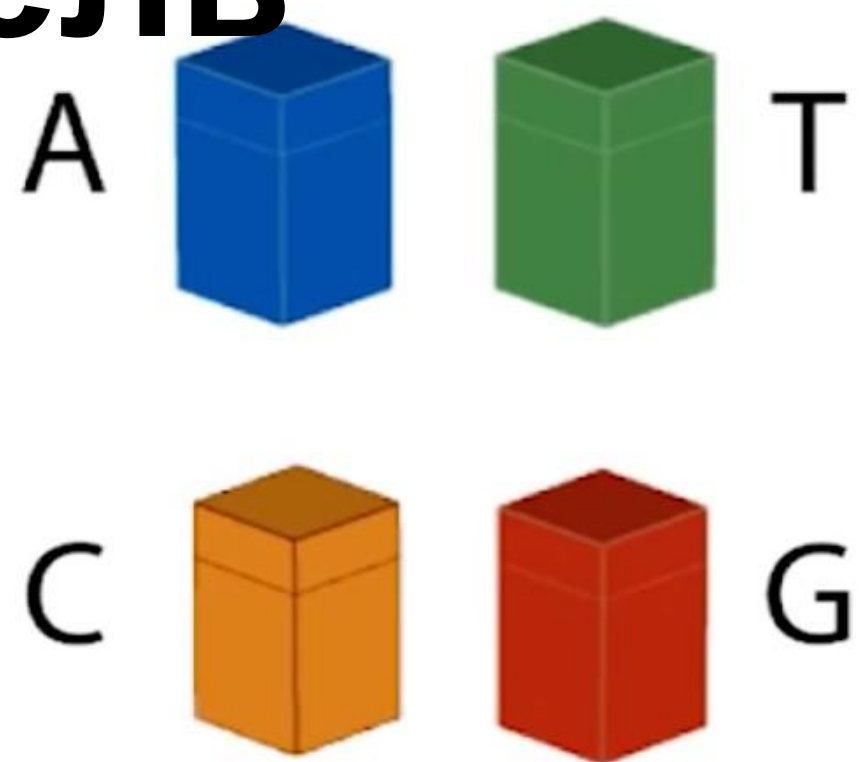


# Шуміна, модель

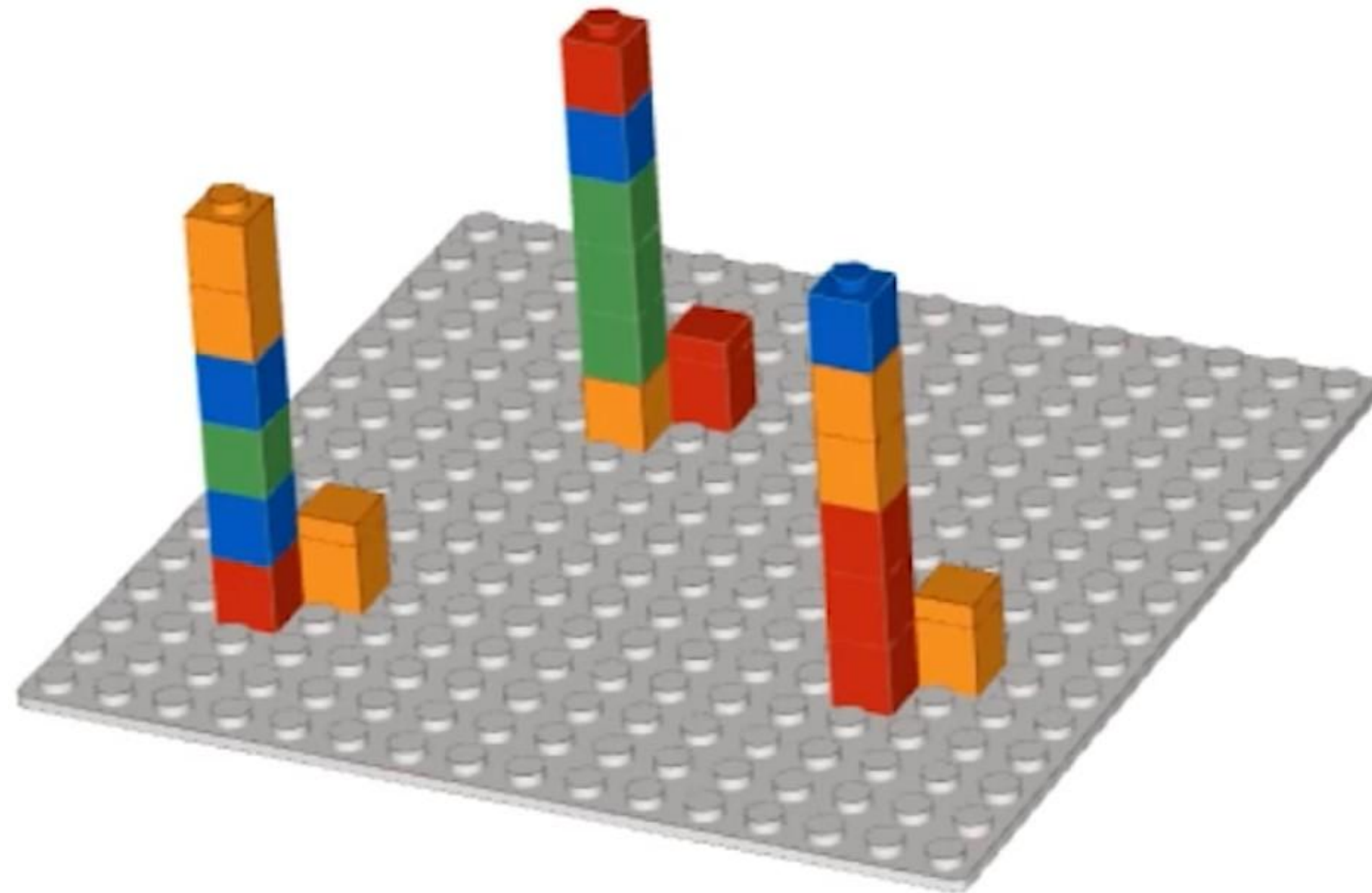




# Illumina, модель

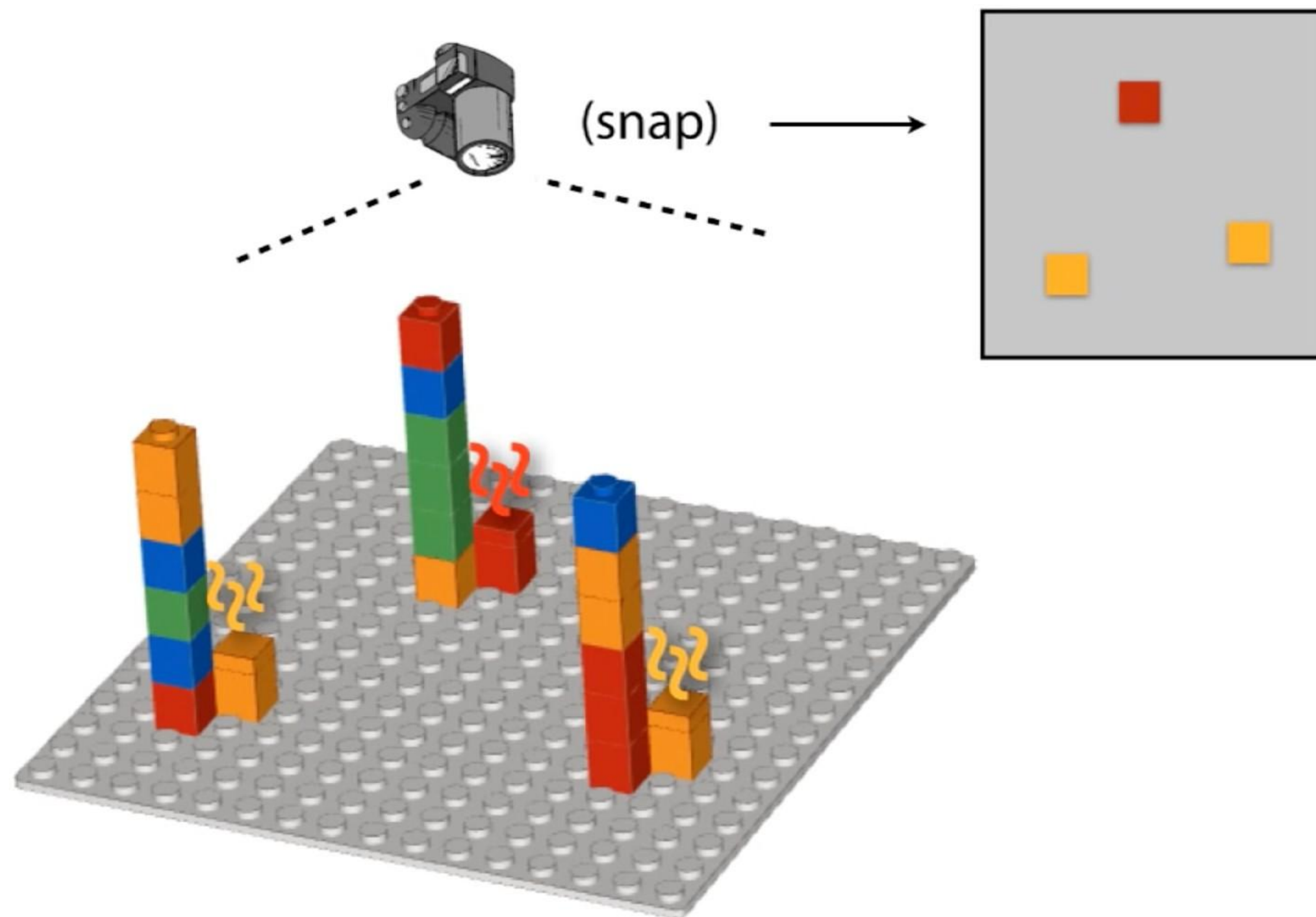


# Илүміна, модель

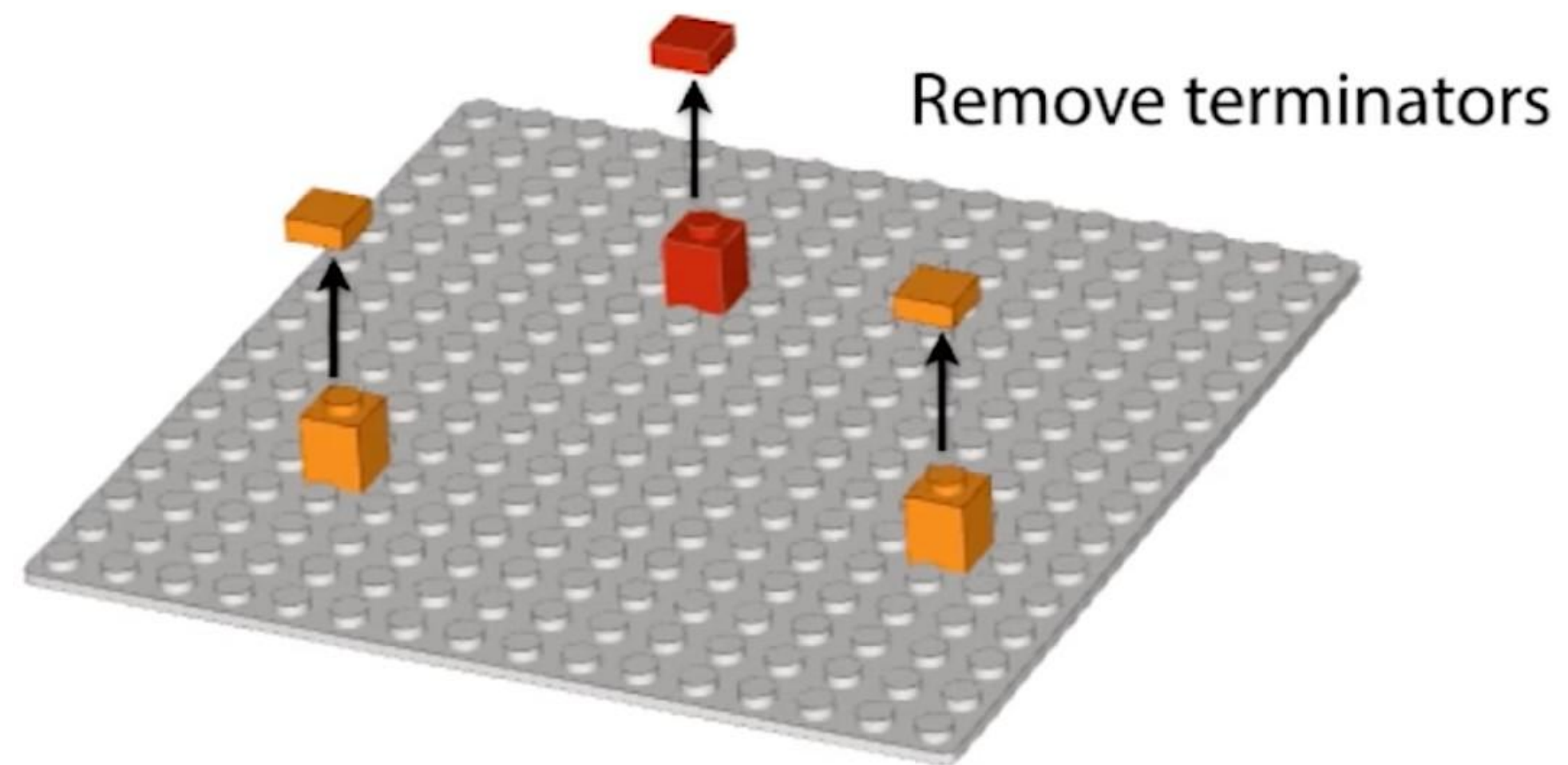




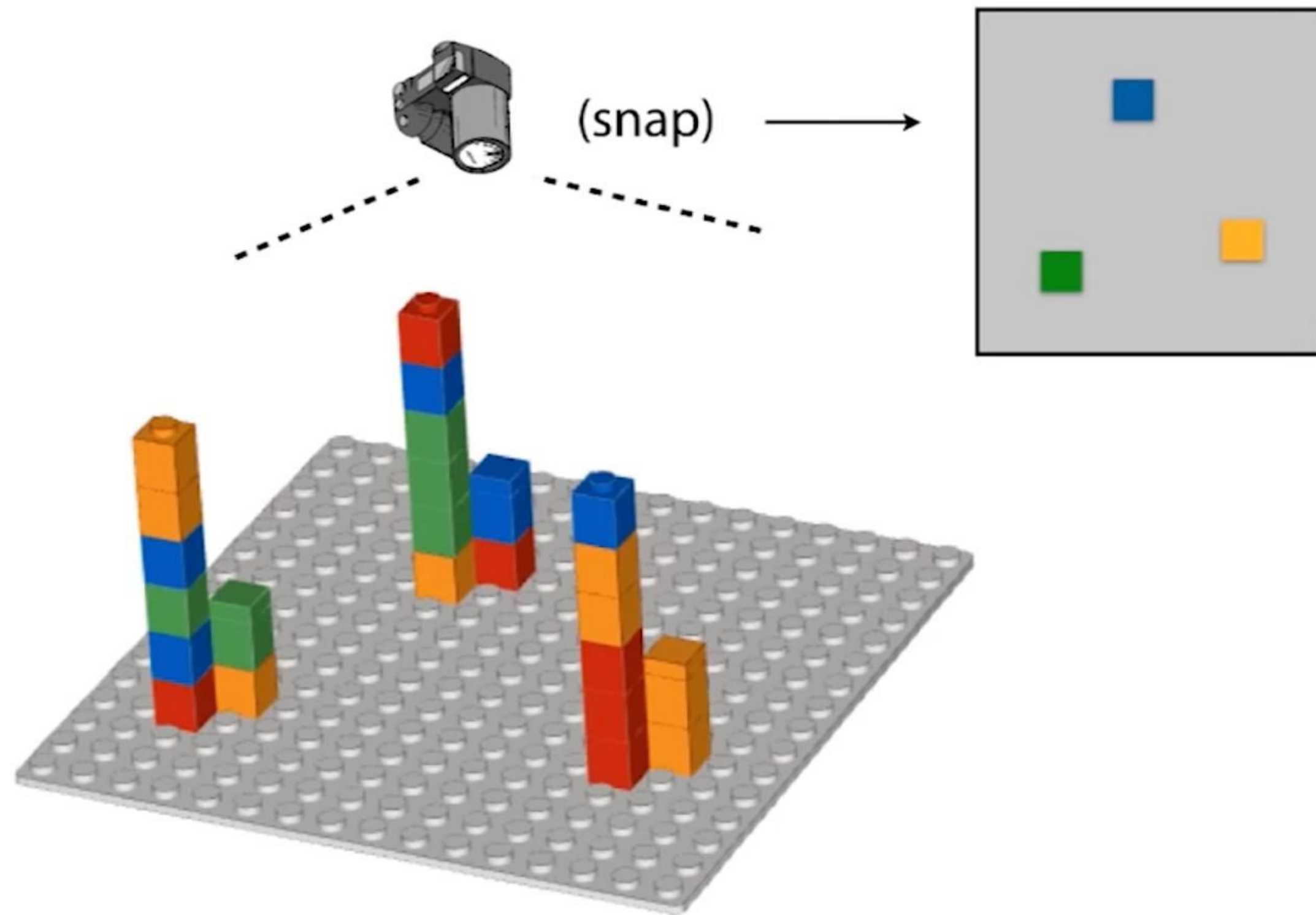
# Илмина, модель



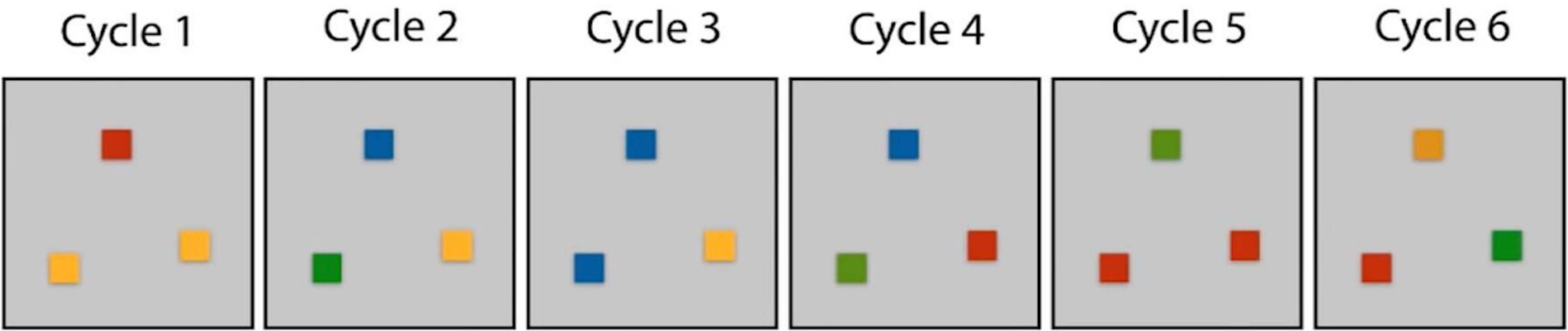
# Илліміна, модель



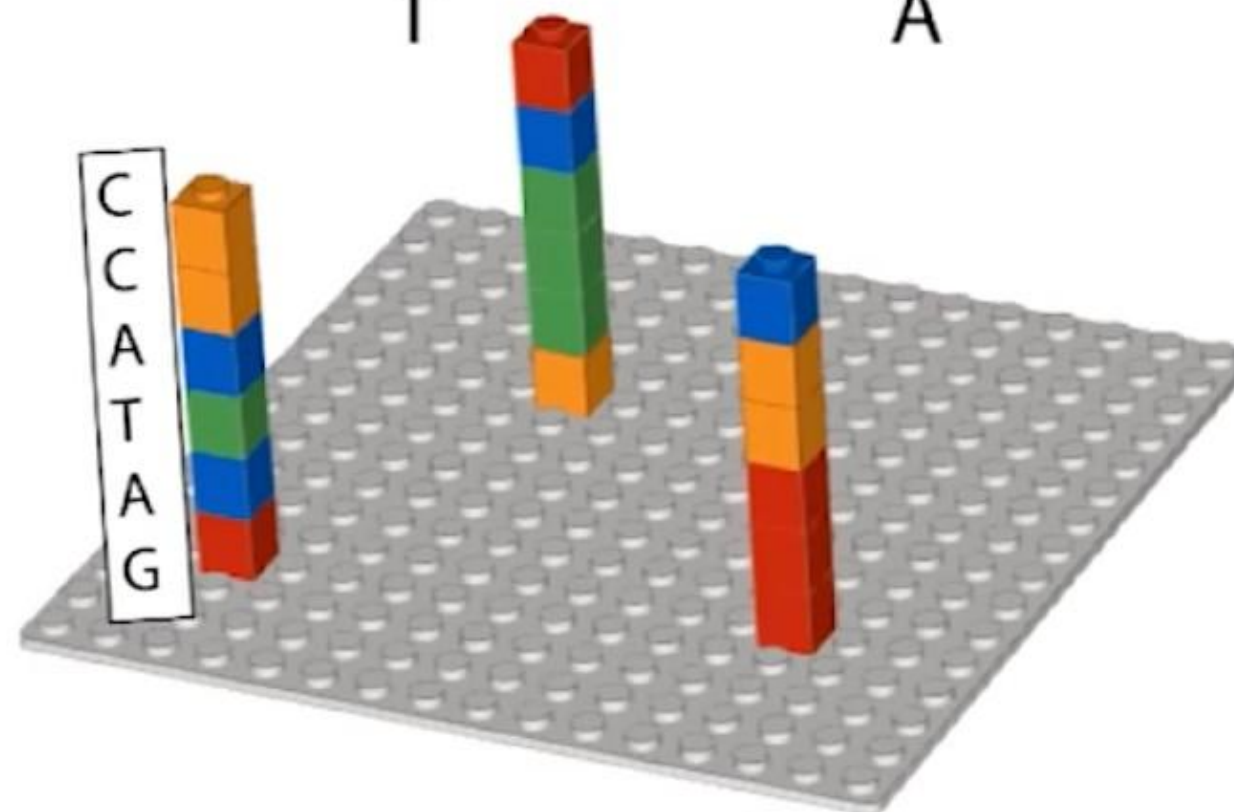
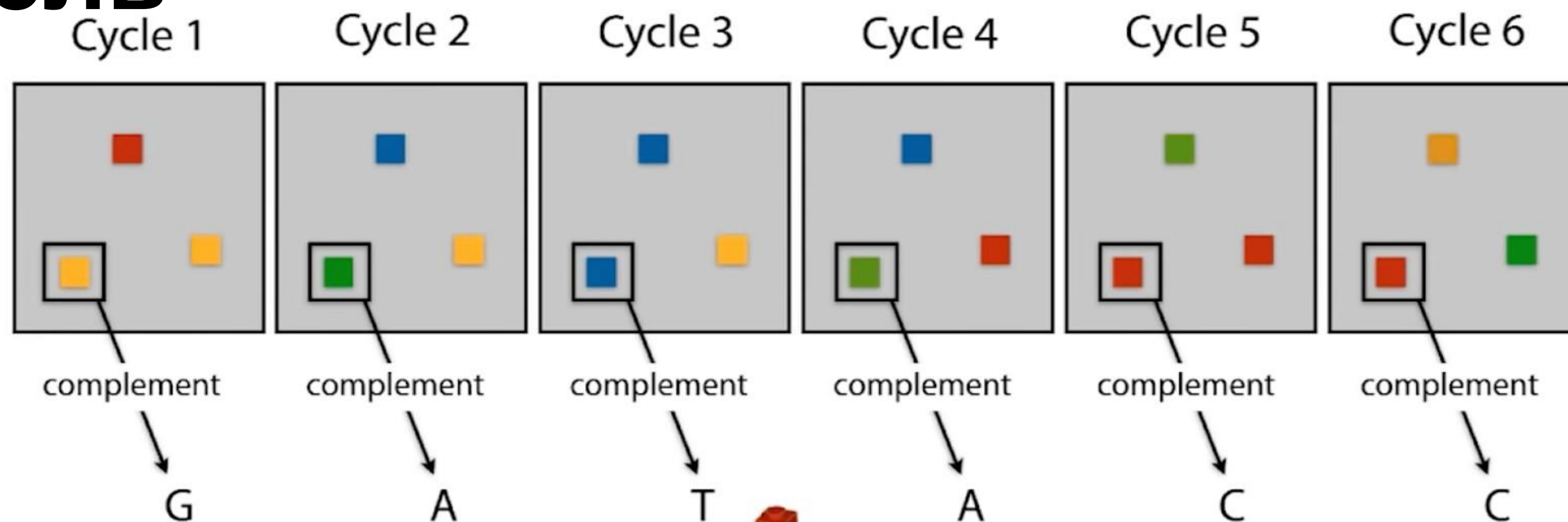
# Илмина, модель



# Иllumina, модель

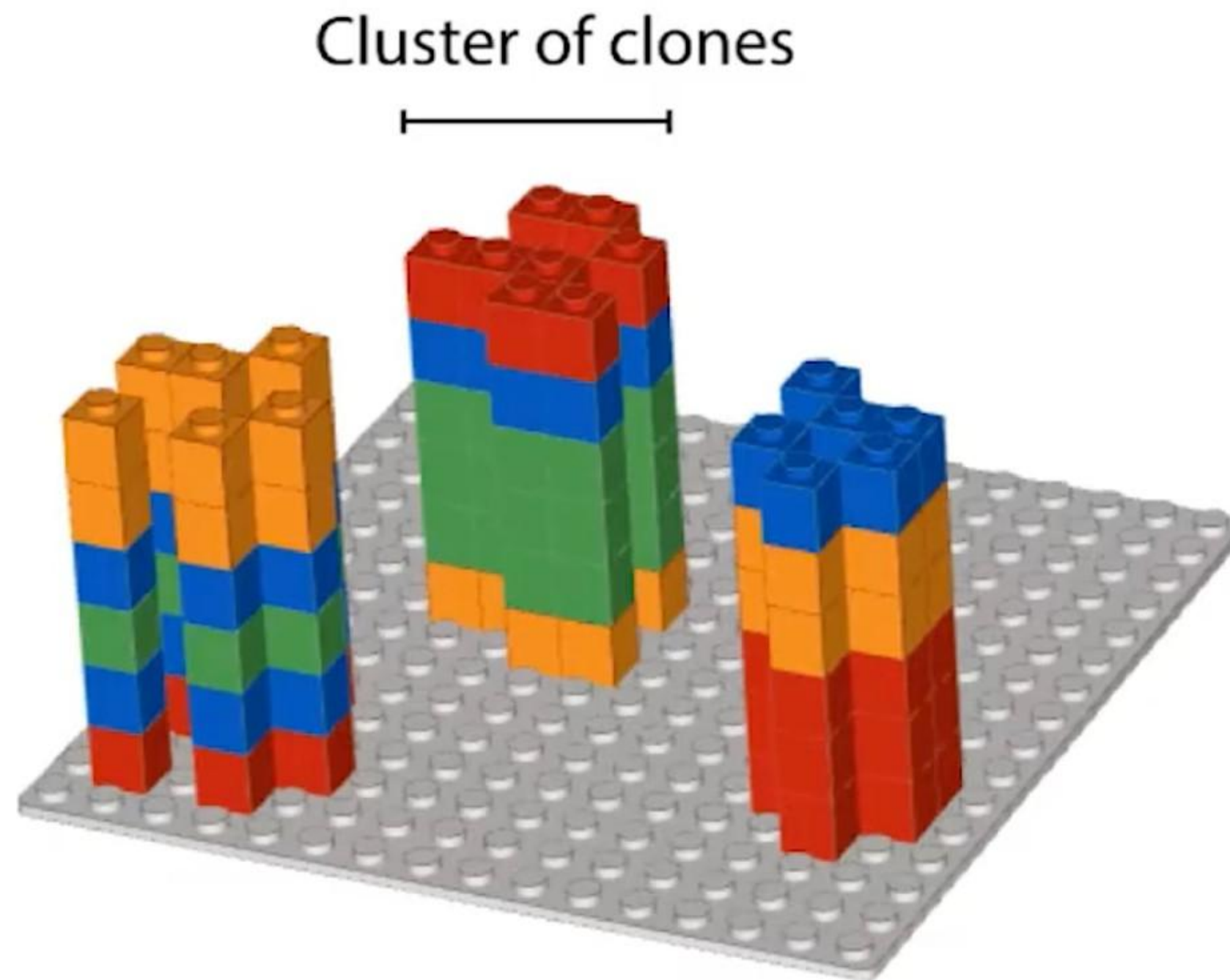


# Illumina, модель

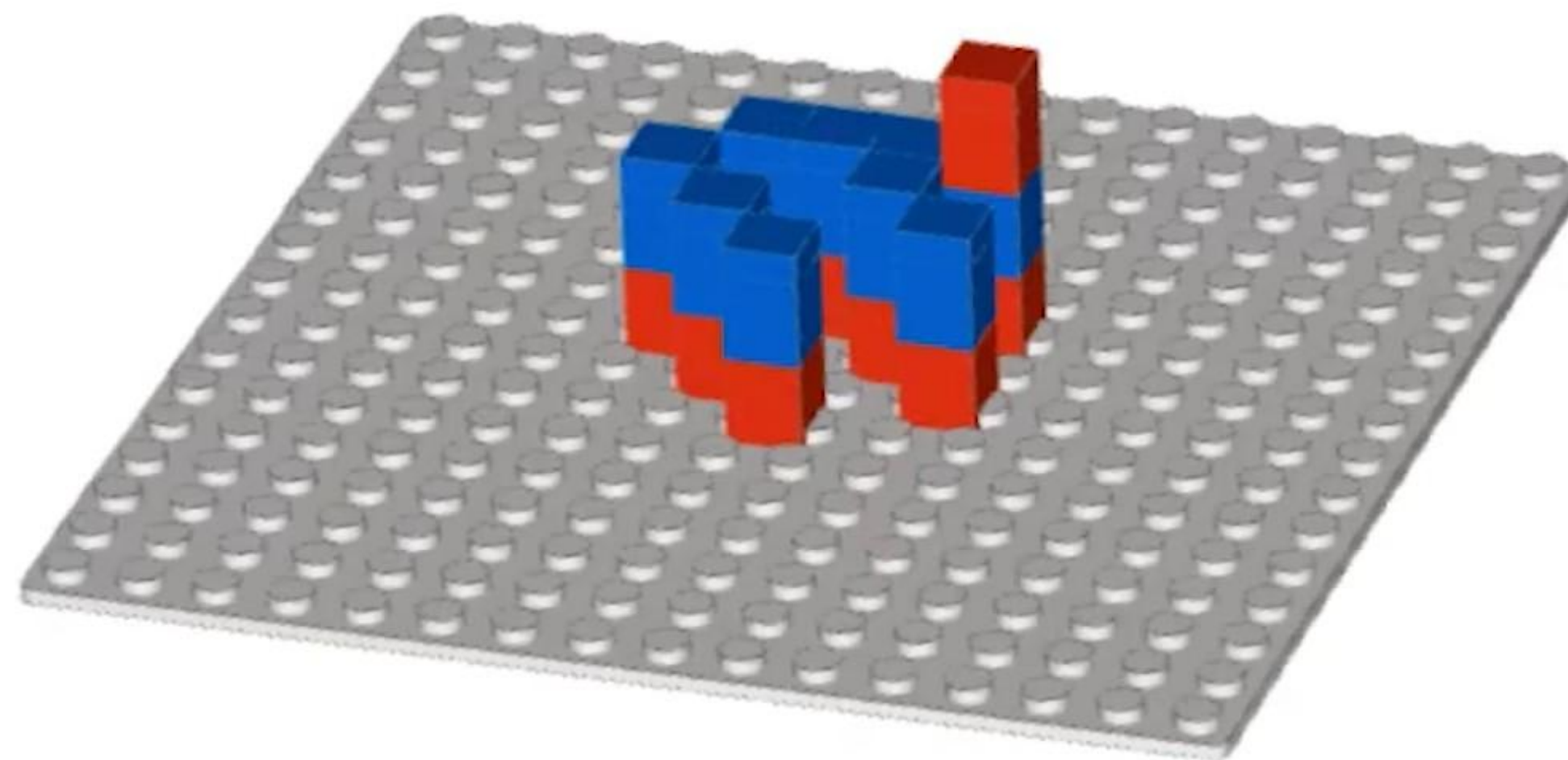




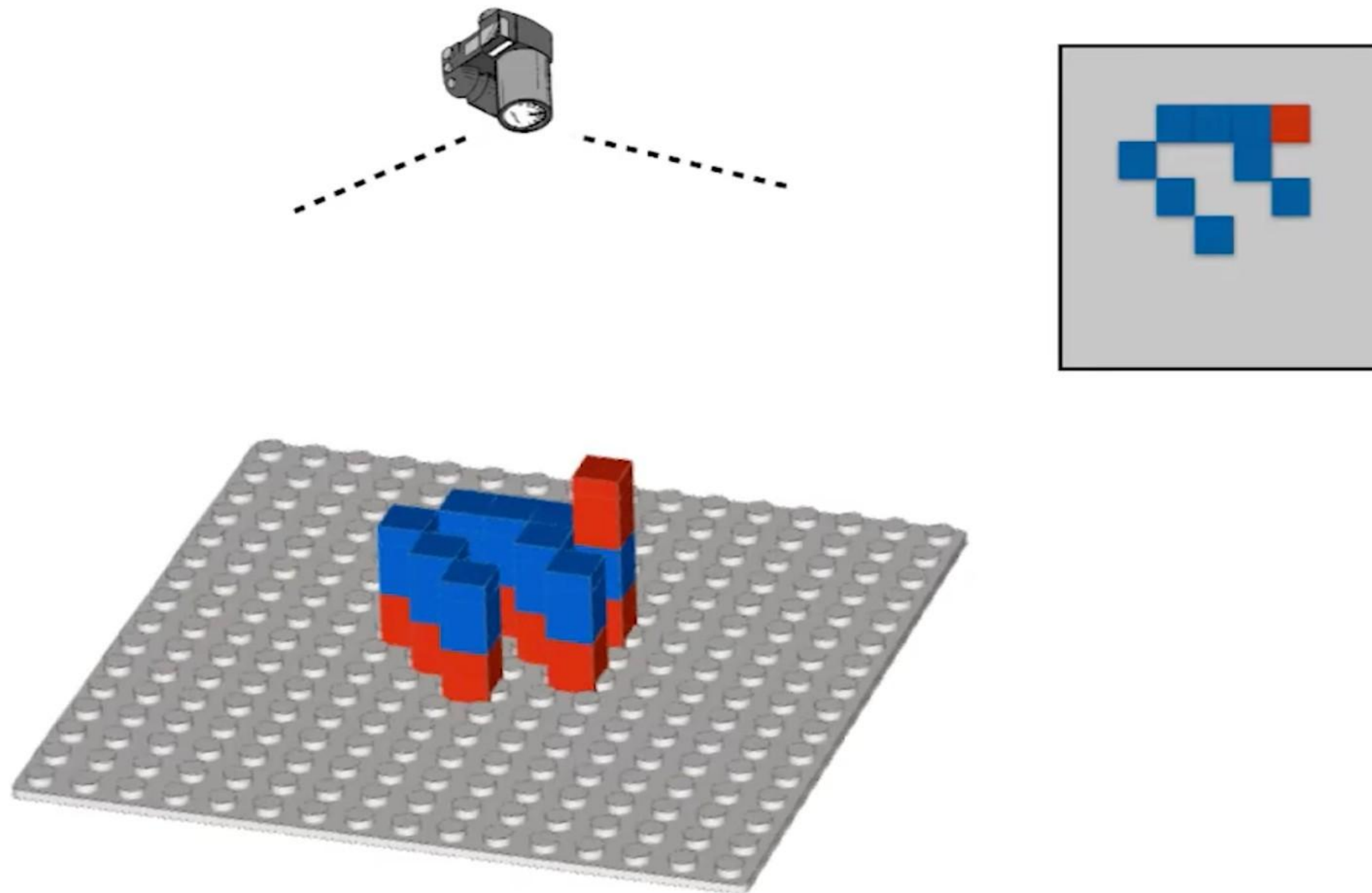
# Ошибки, на примере Illumina



# Ошибки, на примере Illumina

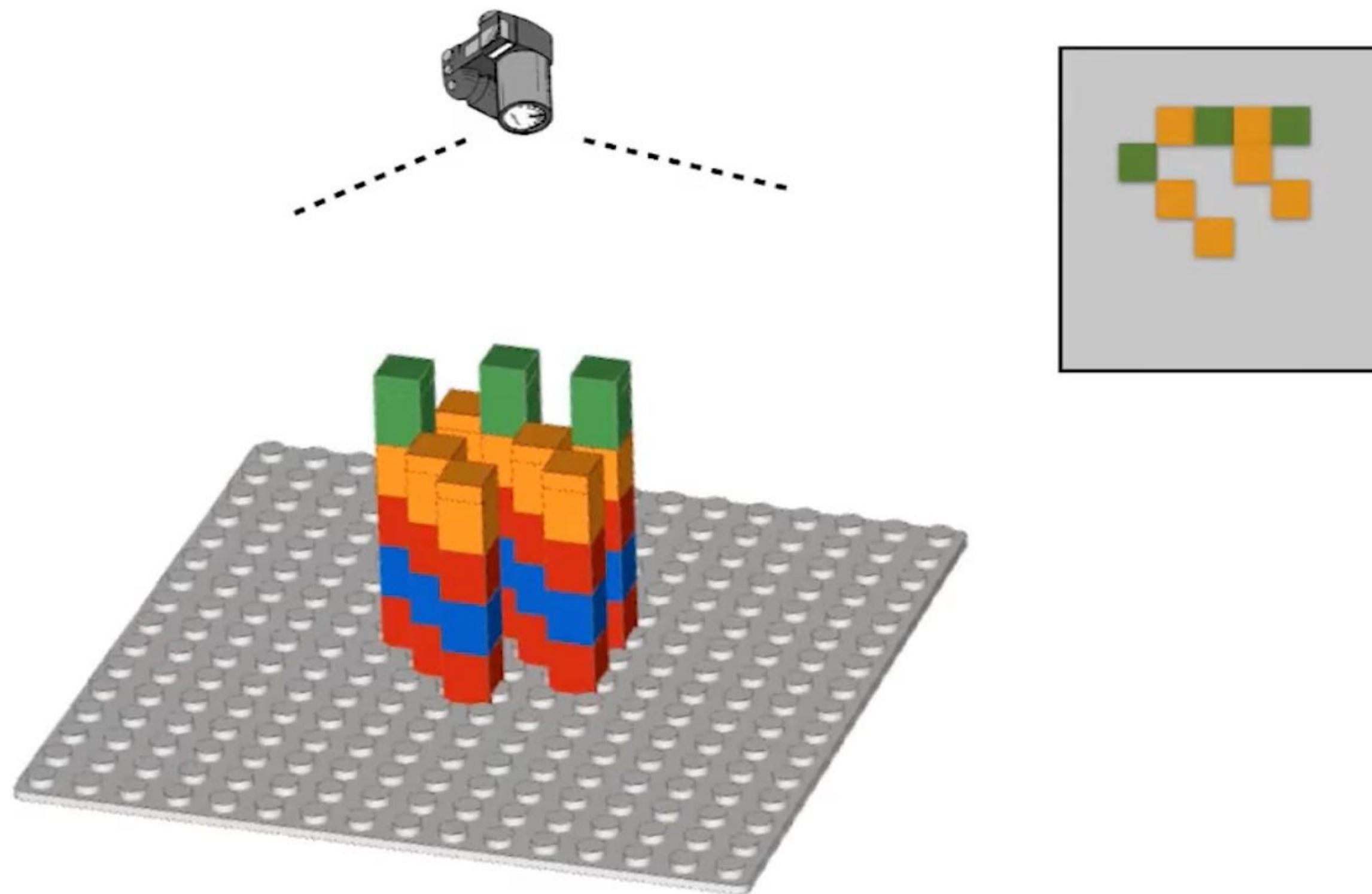


# Ошибки, на примере Illumina





# Ошибки, на примере Illumina



# Ошибки, на примере Illumina

$$Q = - 10 \log_{10}(p)$$

Качество прочтения



Вероятность ошибки



# Ошибки, на примере Illumina

$$Q = -10 \log_{10}(p)$$

Качество прочтения



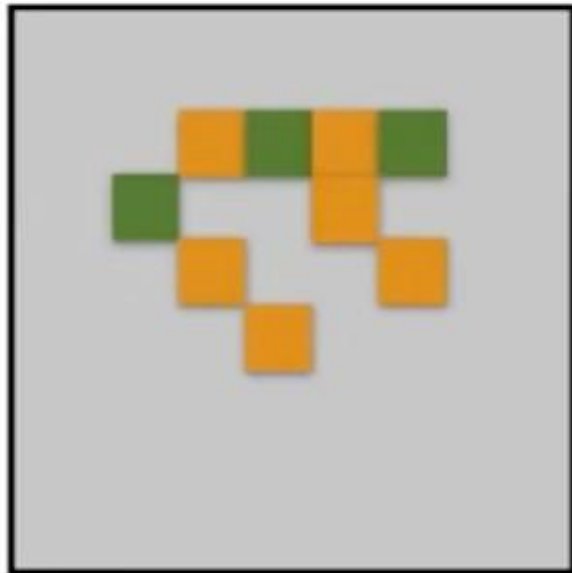
Вероятность ошибки

$Q = 10 \rightarrow 1$  к  $10$  что произошла ошибка

$Q = 20 \rightarrow 1$  к  $100$

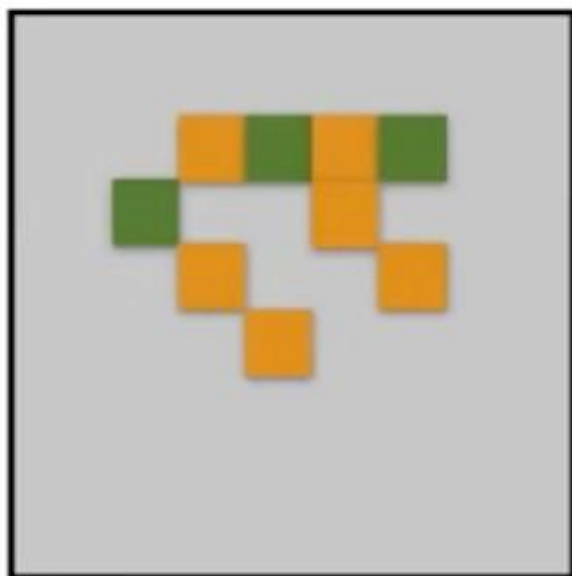
$Q = 30 \rightarrow 1$  к  $1000$

# Ошибки, на примере Illumina



$$p = \frac{3}{9} = \frac{1}{3}$$

# Ошибки, на примере Illumina



$$p = \frac{3}{9} = \frac{1}{3}$$

$$Q = -10 \log_{10} \left( \frac{1}{3} \right) = 4.77$$

# FASTQ

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!' '*(((***+))%%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65
```

# FASTQ

```
@SEQ_ID [идентификатор]  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT  
+  
!' '*(((***+))%%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65
```

# FASTQ

```
@SEQ_ID [идентификатор]  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT [последовательность]  
+  
!' '*(((***+))%%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65
```



# FASTQ

```
@SEQ_ID [идентификатор]  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT [последовательность]  
+ [необязательная строка]  
!' '*((( (**+)) %%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65
```

# FASTQ

```
@SEQ_ID [идентификатор]  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT [последовательность]  
+ [необязательная строка]  
!' '*((( (**+)) %%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65 [качество прочтения]
```

# FASTQ

```
@SEQ_ID [идентификатор]  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT [последовательность]  
+ [необязательная строка]  
!' '*((( (**+)) %%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65 [качество прочтения]
```

ASCII представление того самого  $Q$

Из качества в символ: `chr(Q + 33)`

Обратно: `ord(qual) - 33`

# Визуализация

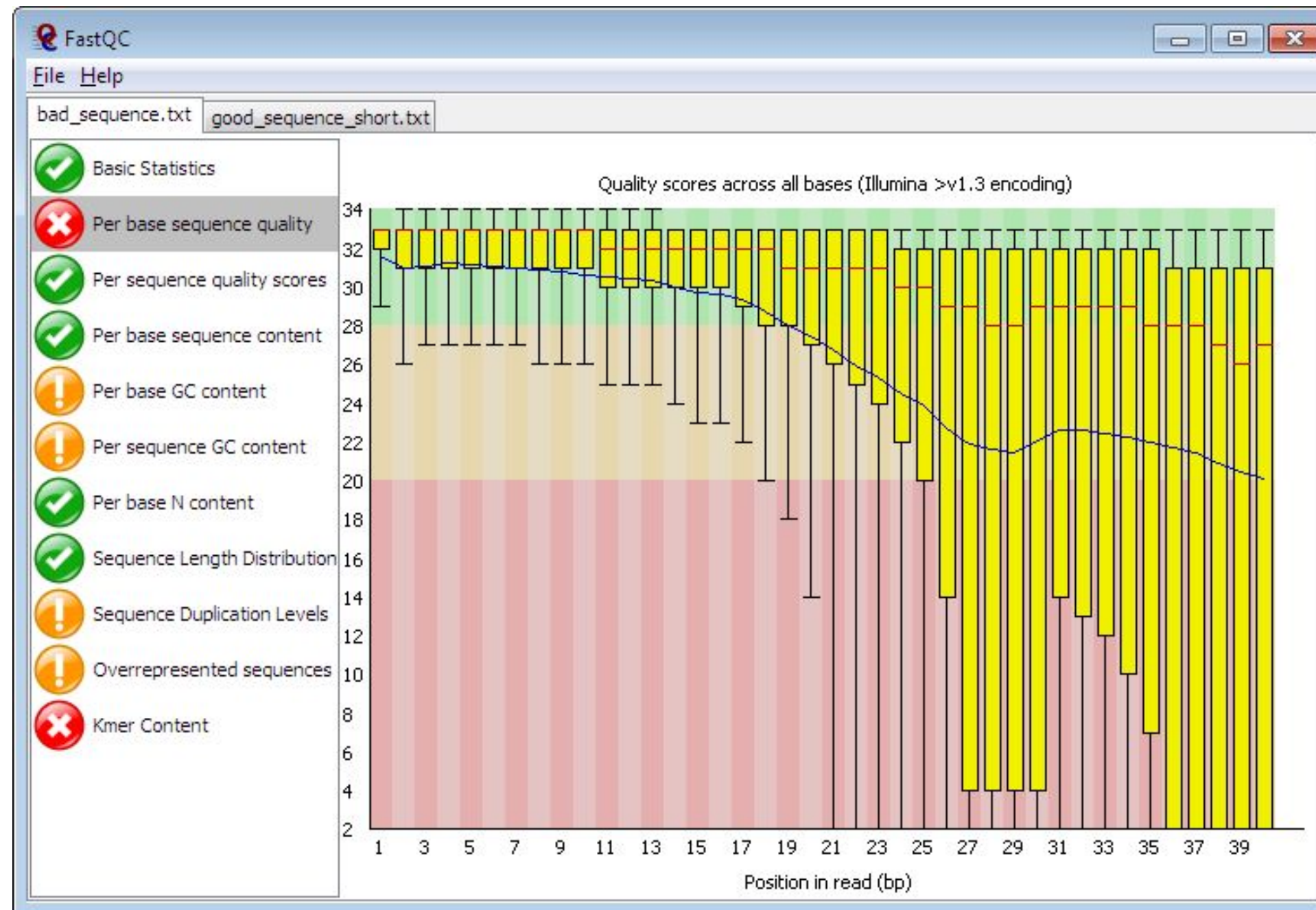
**FastQC** [<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>]

```
>>fastqc bad_sequence.txt good_sequence.txt
```

# Визуализация

**FastQC** [<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>]

```
>>fastqc bad_sequence.txt good_sequence.txt
```



# Удаление плохих ридов

**Trimmomatic** [<http://www.usadellab.org/cms/?page=trimmomatic>]

```
java -jar trimmomatic-0.35.jar SE -phred33 input.fq.gz output.fq.gz  
ILLUMINACLIP:TruSeq3-SE:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

# Удаление плохих ридов

**Trimmomatic** [<http://www.usadellab.org/cms/?page=trimmomatic>]

```
java -jar trimmomatic-0.35.jar SE -phred33 input.fq.gz output.fq.gz  
ILLUMINACLIP:TruSeq3-SE:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

Удаление адаптеров (ILLUMINACLIP:TruSeq3-SE:2:30:10)

# Удаление плохих ридов

**Trimmomatic** [<http://www.usadellab.org/cms/?page=trimmomatic>]

```
java -jar trimmomatic-0.35.jar SE -phred33 input.fq.gz output.fq.gz  
ILLUMINACLIP:TruSeq3-SE:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

Удаление адаптеров (ILLUMINACLIP:TruSeq3-SE:2:30:10)

Удаление низкокачественных вначале (с качеством хуже 3) (LEADING:3)



# Удаление плохих ридов

**Trimmomatic** [<http://www.usadellab.org/cms/?page=trimmomatic>]

```
java -jar trimmomatic-0.35.jar SE -phred33 input.fq.gz output.fq.gz  
ILLUMINACLIP:TruSeq3-SE:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

Удаление адаптеров (ILLUMINACLIP:TruSeq3-SE:2:30:10)

Удаление низкокачественных вначале (с качеством хуже 3) (LEADING:3)

Удаление низкокачественных в конце (с качеством хуже 3) (TRAILING:3)

# Удаление плохих ридов

**Trimmomatic** [<http://www.usadellab.org/cms/?page=trimmomatic>]

```
java -jar trimmomatic-0.35.jar SE -phred33 input.fq.gz output.fq.gz  
ILLUMINACLIP:TruSeq3-SE:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

Удаление адаптеров (ILLUMINACLIP:TruSeq3-SE:2:30:10)

Удаление низкокачественных вначале (с качеством хуже 3) (LEADING:3)

Удаление низкокачественных в конце (с качеством хуже 3) (TRAILING:3)

Сканировать окном в 4 нуклеотида, если среднее качество в окне ниже 15, то удалять (SLIDINGWINDOW:4:15)

Удалять риды короче 36 нуклеотидов (MINLEN:36)

# Исправление ошибок

Идея!

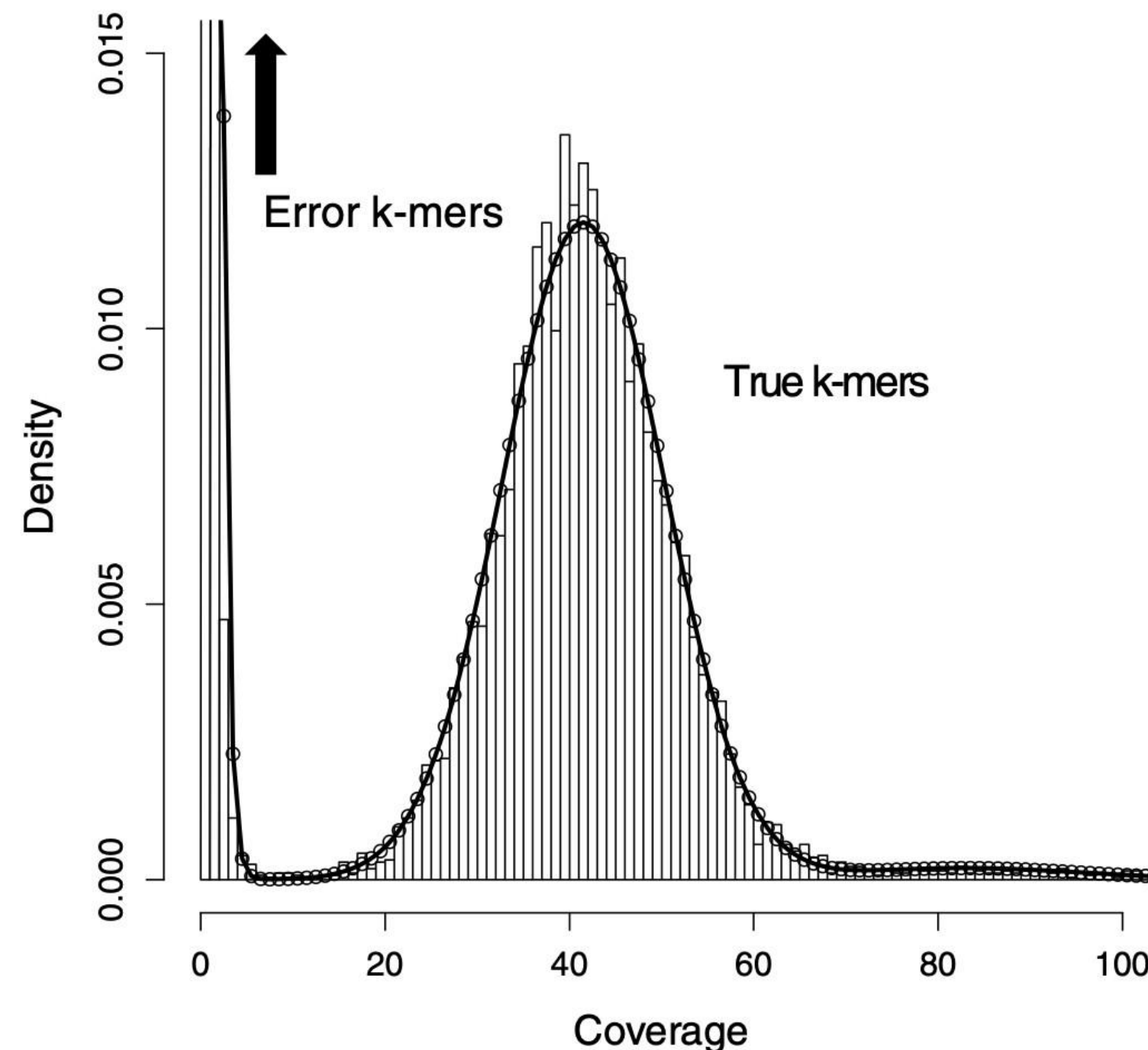
```
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGTGGTTCAAAGCAGTATCGATCAAATA -> GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATAAAATA -> GATTTGGGGTTCAAAGCAGTATCGATCAAATA
GATTTGGGGTTCAAAGCAGTATCGATCAAATA
```

# Исправление ошибок: Quake

1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
  - $k$ -мер: его количество не учитывает качество
  - Будем считать  $q$ -меры.  $k$ -мер: его количество \* произведение качеств

# Исправление ошибок: Quake

1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
  - $k$ -мер: его количество не учитывает качество
  - Будем считать  $q$ -меры.  $k$ -мер: его количество \* произведение качеств



# Исправление ошибок: Quake

1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
  - $k$ -мер: его количество не учитывает качество
  - Будем считать  $q$ -меры.  $k$ -мер: его количество \* произведение качеств
  - По распределению разделим  $k$ -меры на 2 кластера (ошибочные и правильные)
2. Все ошибочные  $k$ -меры – кандидаты на исправление в рядах  
Нуклеотиды которые наблюдаем в риде  $O = O_1, O_2, \dots, O_N$   
Те которые на самом деле в геноме  $A = A_1, A_2, \dots, A_N$   
Нам бы хотелось по наблюдаемым понять наиболее вероятные  $A_i$



# Исправление ошибок: Quake

1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
  - $k$ -мер: его количество не учитывает качество
  - Будем считать  $q$ -меры.  $k$ -мер: его количество \* произведение качеств
  - По распределению разделим  $k$ -меры на 2 кластера (ошибочные и правильные)
2. Все ошибочные  $k$ -меры – кандидаты на исправление в рядах  
Нуклеотиды которые наблюдаем в ряде  $O = O_1, O_2, \dots, O_N$   
Те которые на самом деле в геноме  $A = A_1, A_2, \dots, A_N$ 
  - Нам бы хотелось по наблюдаемым понять наиболее вероятные  $A_i$

$$P(A = a | O = o) = \prod_{i=1}^N \frac{P(O_i = o_i | A_i = a_i) P(A_i = a_i)}{P(O_i = o_i)}$$

# Исправление ошибок: Quake

1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
  - $k$ -мер: его количество не учитывает качество
  - Будем считать  $q$ -меры.  $k$ -мер: его количество \* произведение качеств
  - По распределению разделим  $k$ -меры на 2 кластера (ошибочные и правильные)
2. Все ошибочные  $k$ -меры – кандидаты на исправление в рядах  
Нуклеотиды которые наблюдаем в ряде  $O = O_1, O_2, \dots, O_N$   
Те которые на самом деле в геноме  $A = A_1, A_2, \dots, A_N$ 
  - Нам бы хотелось по наблюдаемым понять наиболее вероятные  $A_i$

- $$P(A = a | O = o) = \prod_{i=1}^N \frac{P(O_i = o_i | A_i = a_i) P(A_i = a_i)}{P(O_i = o_i)}$$

$$P(O_i = o_i | A_i = a_i) = \begin{cases} p_i & \text{if } o_i = a_i \\ (1 - p_i) E_{q_i}(a_i, o_i) & \text{otherwise} \end{cases} \quad \text{где} \quad p_i = 1 - 10^{-\frac{q_i}{10}}$$

# Исправление ошибок: Quake

1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
  - $k$ -мер: его количество не учитывает качество
  - Будем считать  $q$ -меры.  $k$ -мер: его количество \* произведение качеств
  - По распределению разделим  $k$ -меры на 2 кластера (ошибочные и правильные)
2. Все ошибочные  $k$ -меры – кандидаты на исправление в рядах  
Нуклеотиды которые аблюдаем в риде  $O = O_1, O_2, \dots, O_N$   
Те которые на самом деле в геноме  $A = A_1, A_2, \dots, A_N$ 
  - Нам бы хотелось по наблюдаемым понять наиболее вероятные  $A_i$

$$\circ P(A = a | O = o) = \prod_{i=1}^N \frac{P(O_i = o_i | A_i = a_i) P(A_i = a_i)}{P(O_i = o_i)}$$

$$P(O_i = o_i | A_i = a_i) = \begin{cases} p_i & \text{if } o_i = a_i \\ (1 - p_i) E_{q_i}(a_i, o_i) & \text{otherwise} \end{cases} \quad \text{где} \quad p_i = 1 - 10^{-\frac{q_i}{10}}$$

3.  
Коррекция

# Исправление ошибок: Quake

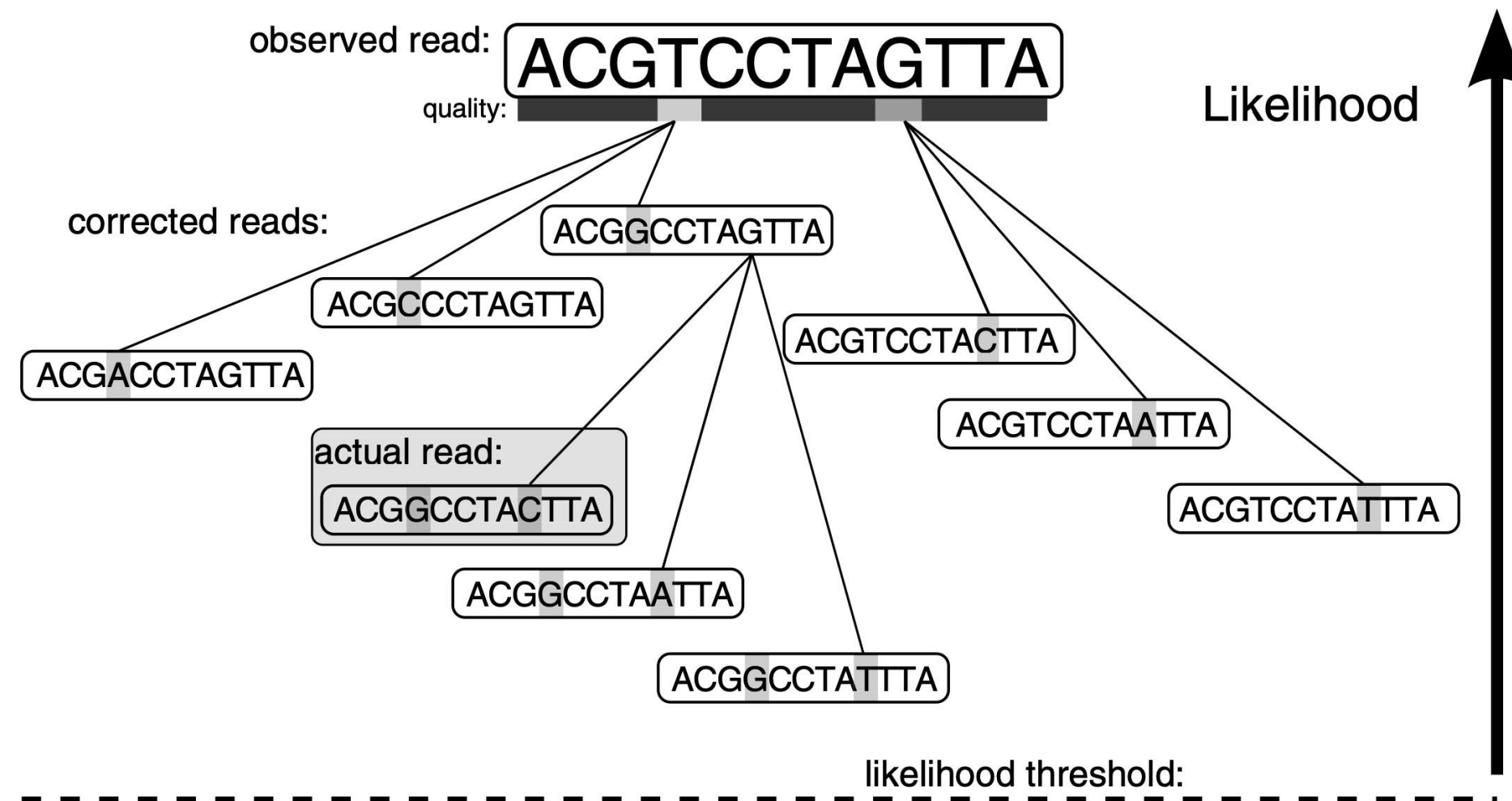
1. Нужно посчитать все  $k$ -меры и определить, какие ошибочные
2. Все ошибочные  $k$ -меры — кандидаты на исправление в рядах
3.
  - Коррекция
    - Хотим сделать в каждом ряде такие замены, которые переместят все  $k$ -меры из
    - кластера ошибочных в кластер правильных

```
1: function SEARCH( $R$ )
2:    $P$ .PUSH( $\{\}$ , 1)
3:   while ( $C, L$ )  $\leftarrow$   $P$ .POP() do
4:     if VALID( $R, C$ ) then
5:       return  $C$ 
6:     else
7:        $i \leftarrow$  lowest quality unconsidered position
8:       for  $nt \in [A, C, G, T]$  do
9:         if  $R[i] == nt$  then
10:            $C_{nt} = C$ 
11:         else
12:            $C_{nt} = C + (i, nt)$ 
13:            $L_{nt} \leftarrow$  LIKELIHOODRATIO( $R, C_{nt}$ )
14:           if  $L_{nt} > likelihood\_threshold$  then
15:              $P$ .PUSH( $C_{nt}, L_{nt}$ )
16:   return  $\{\}$ 
```

# Исправление ошибок: Quake

1. Нужно посчитать все k-меры и определить, какие ошибочные
2. Все ошибочные k-меры — кандидаты на исправление в рядах
3. Коррекция

Хотим сделать в каждом ряде такие замены, которые переместят все k-меры из о кластера ошибочных в кластер правильных



# Исправление ошибок: Quake

**Quake:** [<http://www.cbcb.umd.edu/software/quake/manual.html>]

```
>>quake.py -f [fastq file list] -k [k-mer size] -p 4
```



# Резюмируем

- Секвенирование — случайный процесс
- Могут происходить ошибки секвенирования
- Можно просто отбрасывать ряды с ошибками
- Но можно и исправлять!