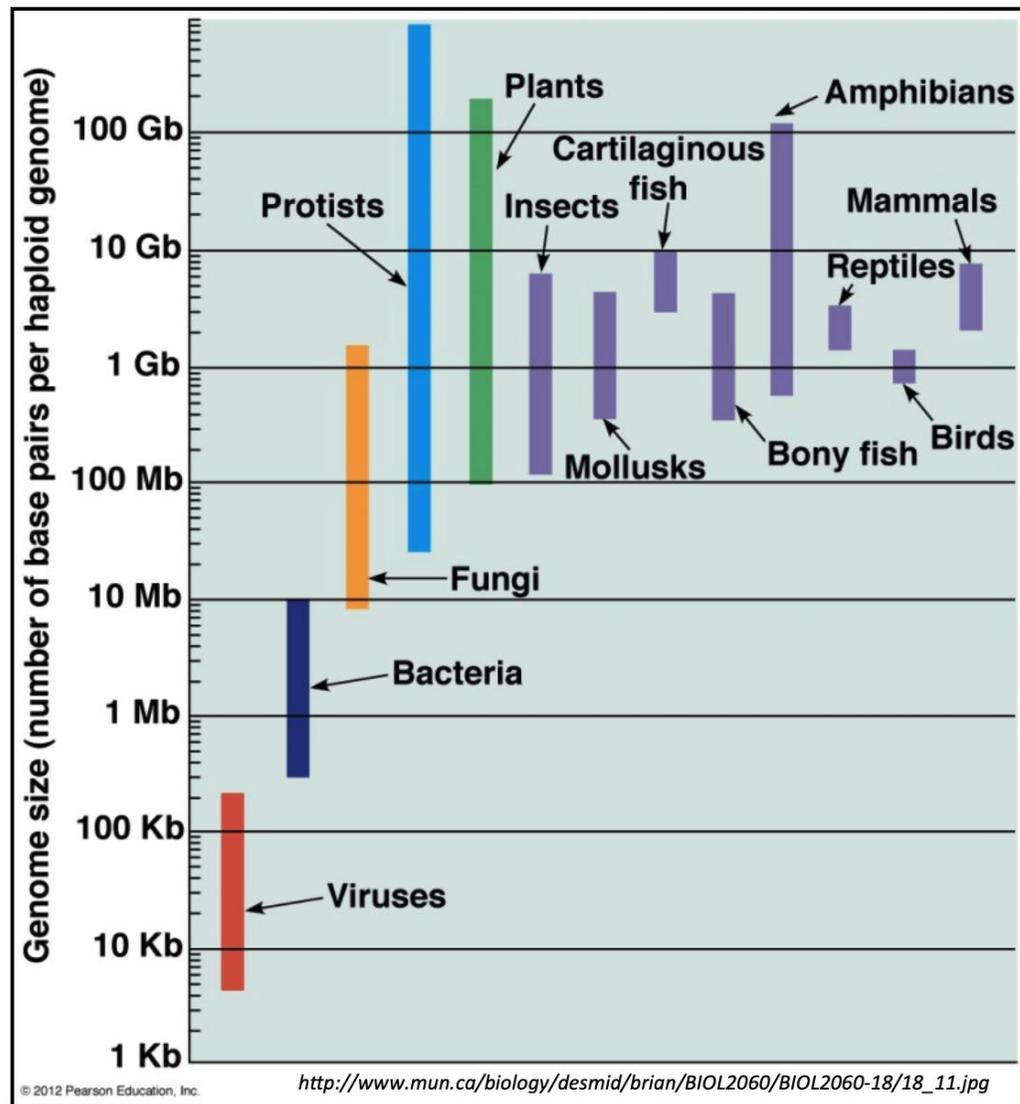


СБОРКА

**ГЕРАСИМОВ ЕВГЕНИЙ,
2020**

КАКИМИ БЫВАЮТ ГЕНОМЫ?



Малые геномы:

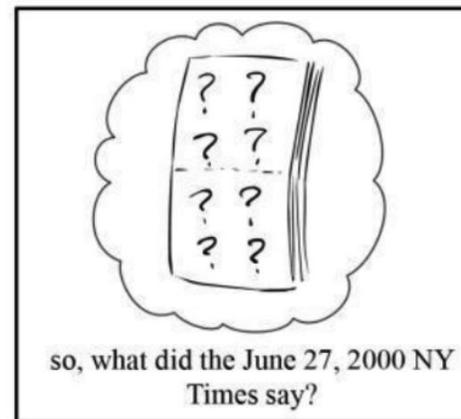
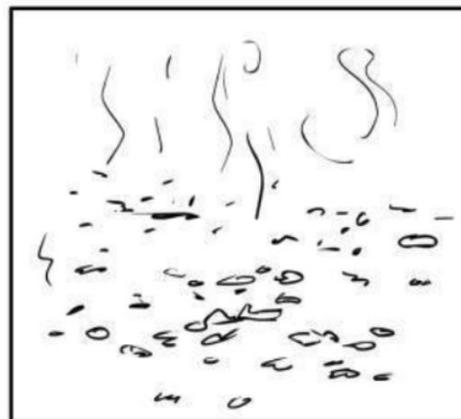
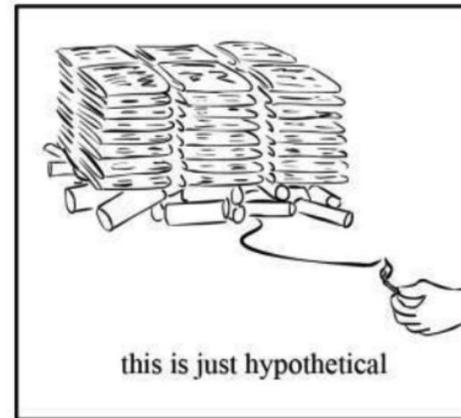
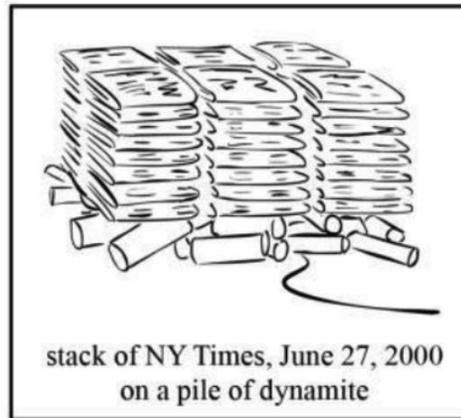
- вирусы
- бактерии
- органеллы

Несколько
примеров

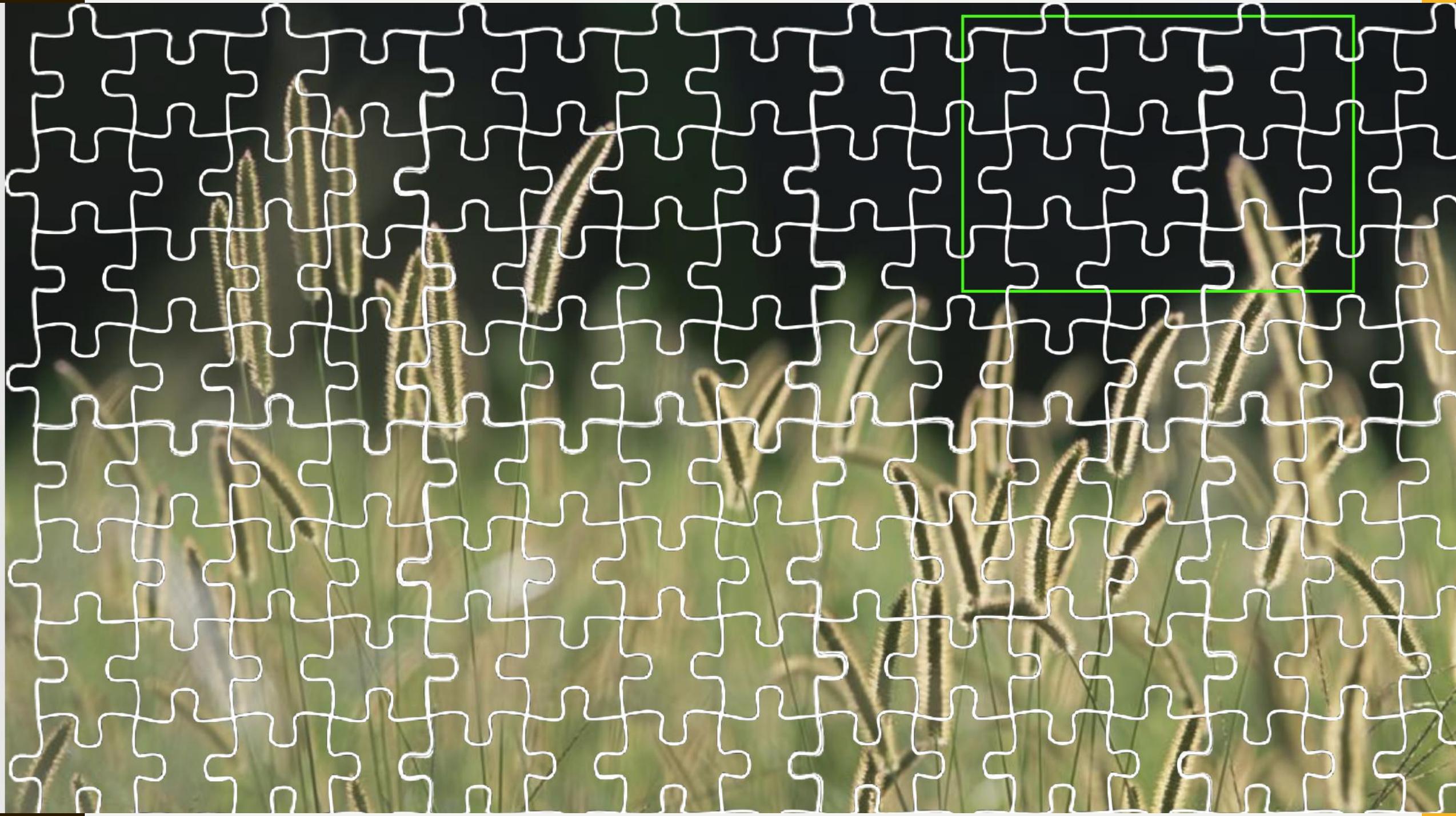
- человек (3.2 Гб)²³
- *E. coli* (4.6 Мб)¹
- дрозофила (139.5 Мб)⁴
- арабидопсис (135 Мб)⁵
- *Paris japonica* (150 Гб)⁴⁰

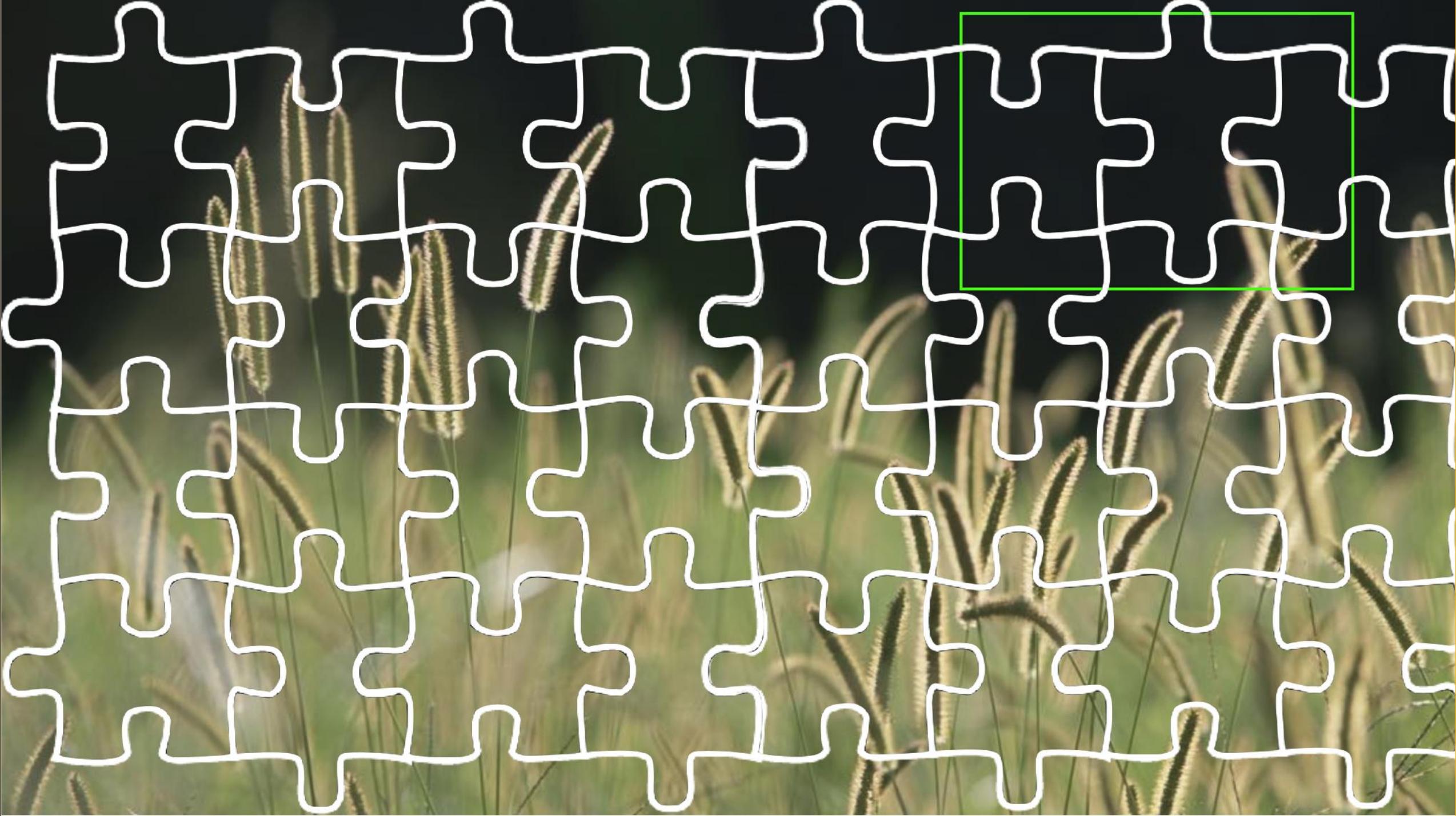
(Вороний глаз японский)

СУТЬ СБОРКИ DE NOVO



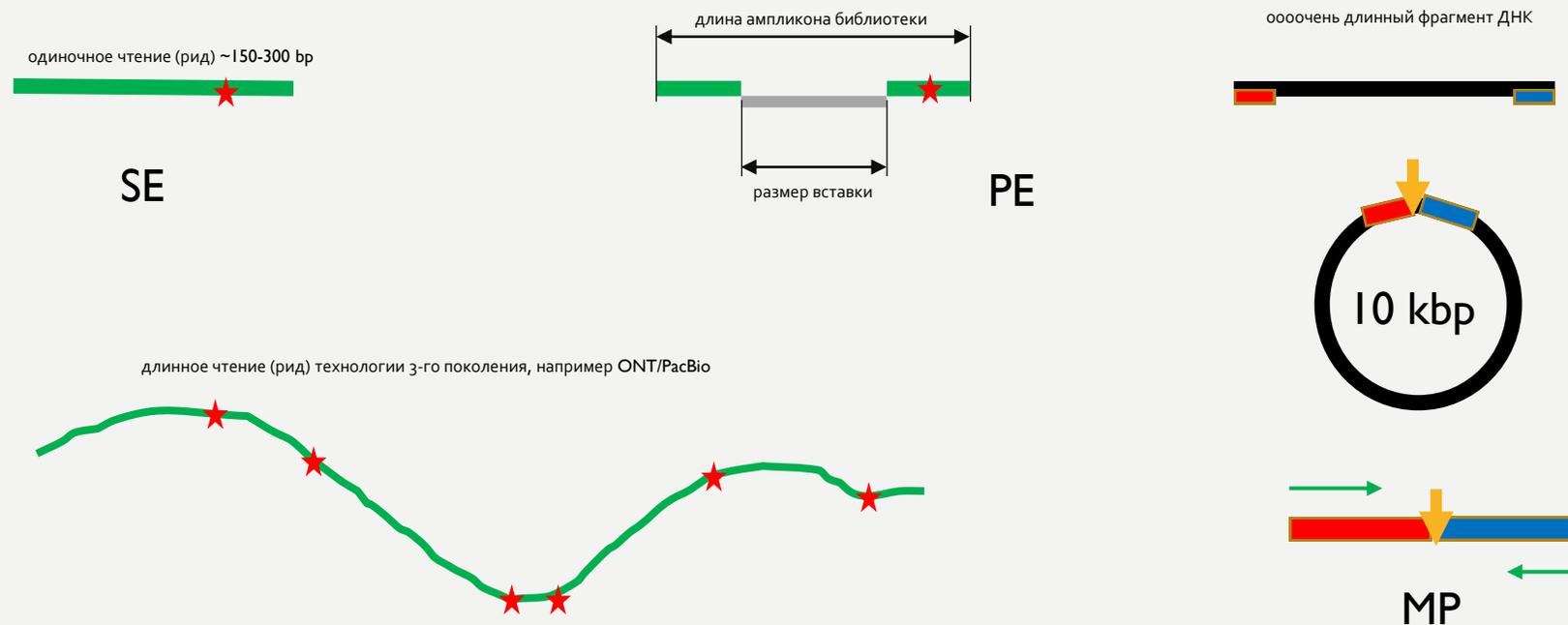




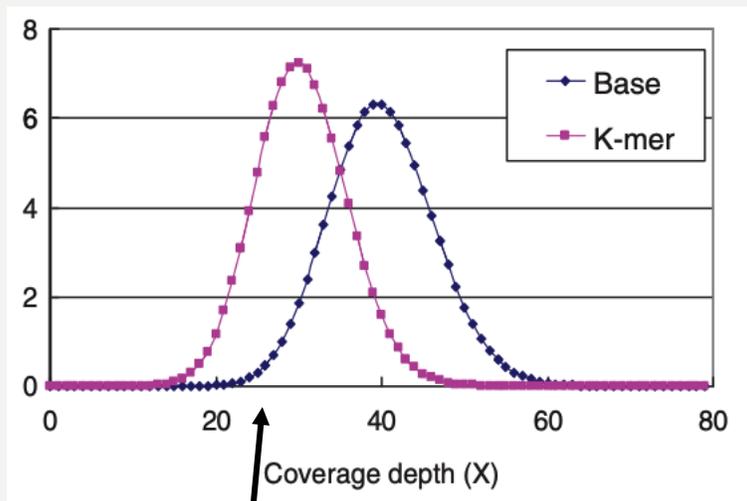


ВХОДНЫЕ ДАННЫЕ

- Чтения одной из платформ NGS, обычно Illumina или PacBio (реже 454 или ONT).



ПОКРЫТИЕ



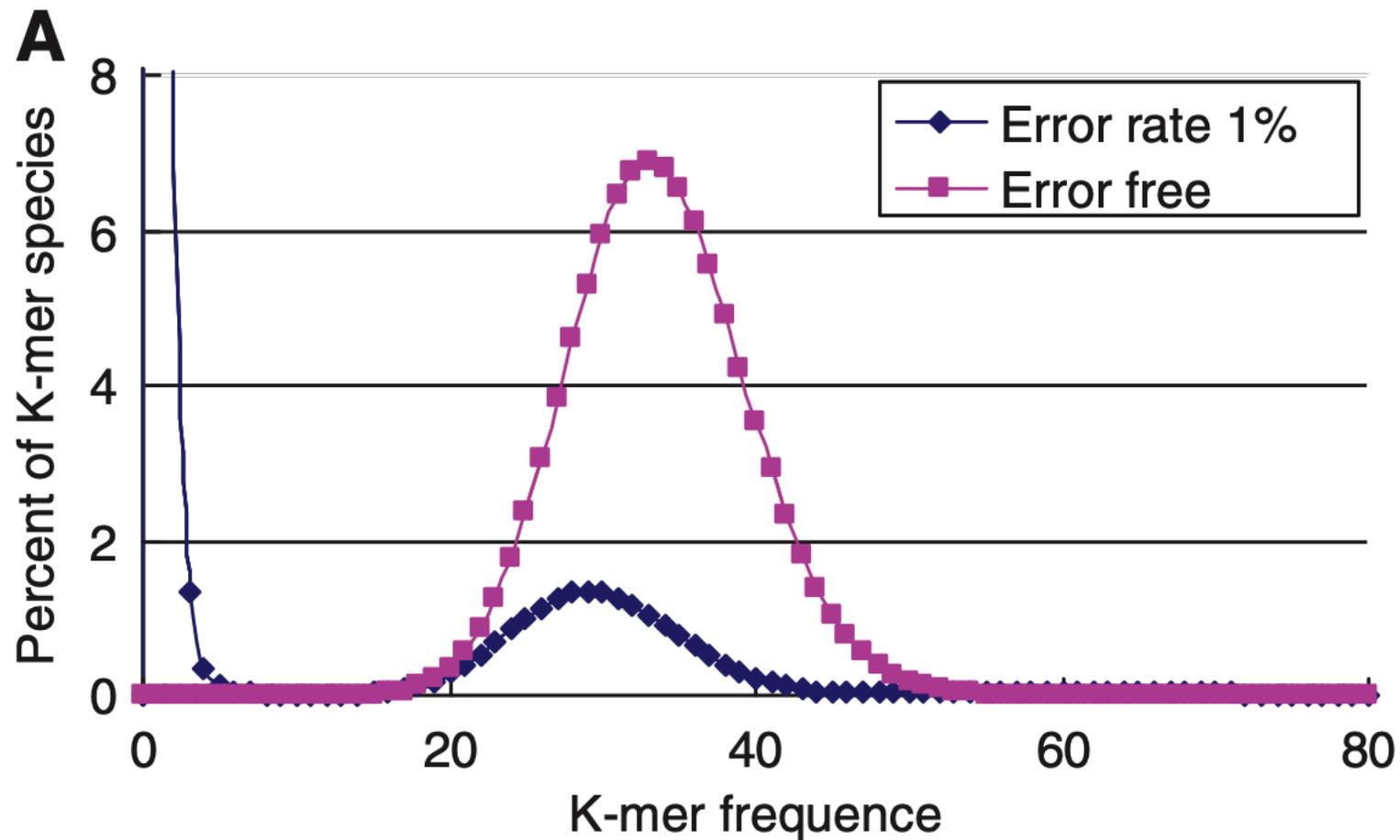
$d_k = ?$

N – число ридов
 L – длина рида
 G – длина генома
 K – длина k-мера
 d_b – среднее покрытие нуклеотида
 d_k – среднее покрытие (частота) k-мера
 n_b – число прочитанных нуклеотидов
 n_k – число прочитанных k-меров

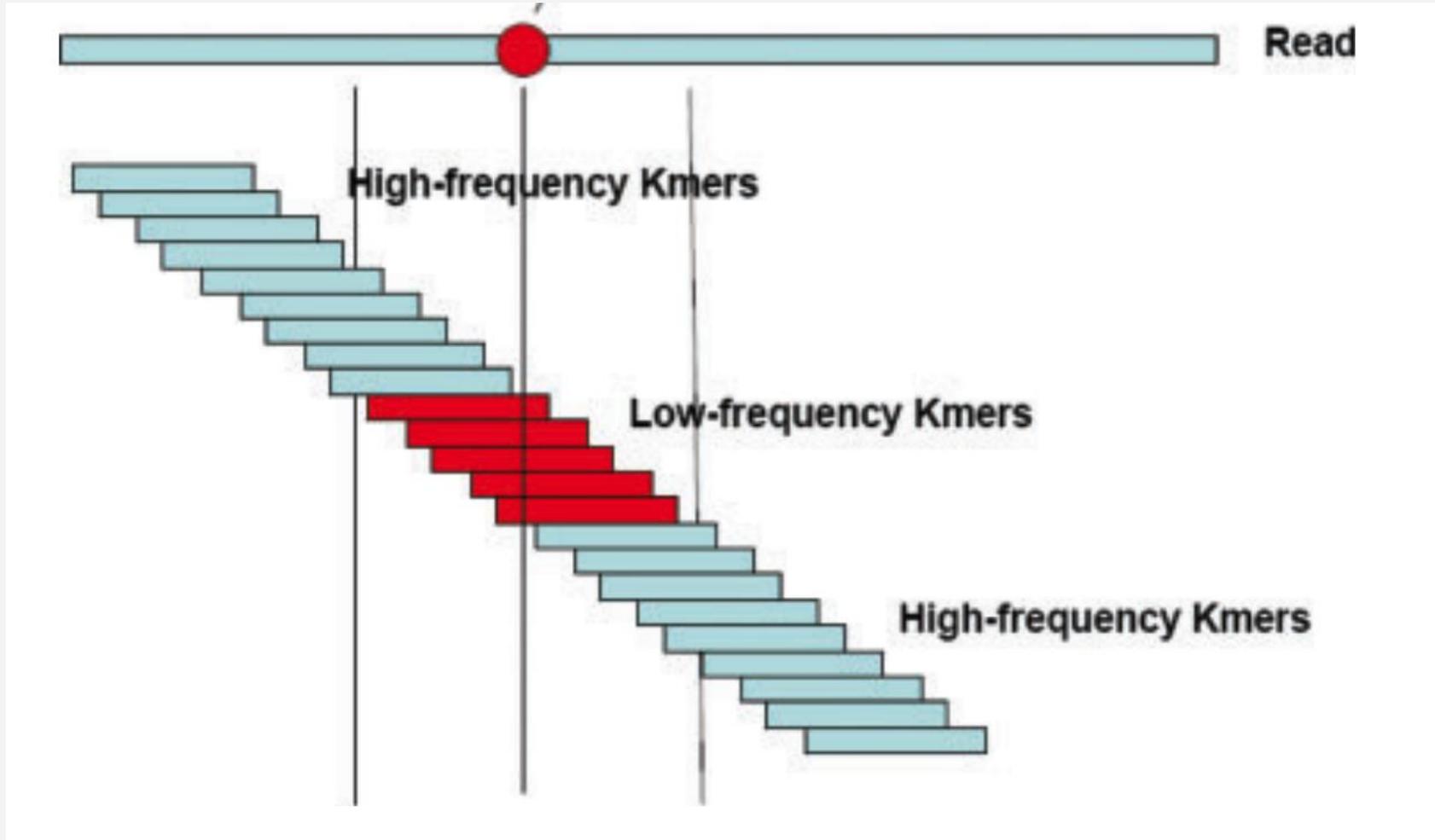
- $n_b = N * L$
 $n_k = N * (L - K + 1)$
- $d_b = n_b / G$
 $d_k = n_k / G; \quad G = n_k / d_k$
- $d_k = d_b * n_k / n_b = d_b * (L - K + 1) / L$

$d_b = ?$

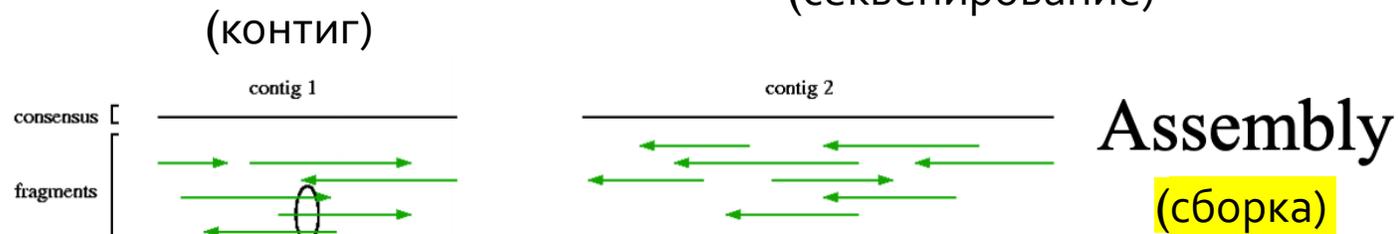
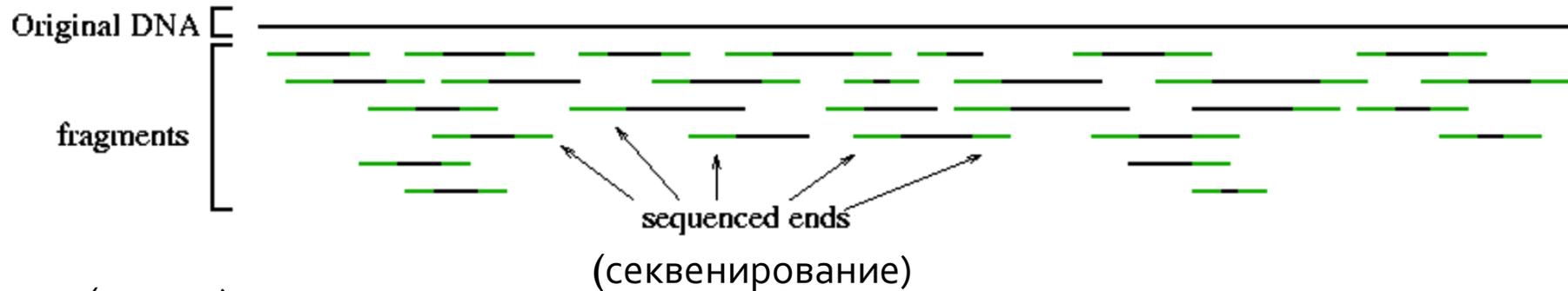
ОШИБКИ В РИДАХ



ИДЕЯ ДЛЯ КОРРЕКЦИИ ОШИБОК



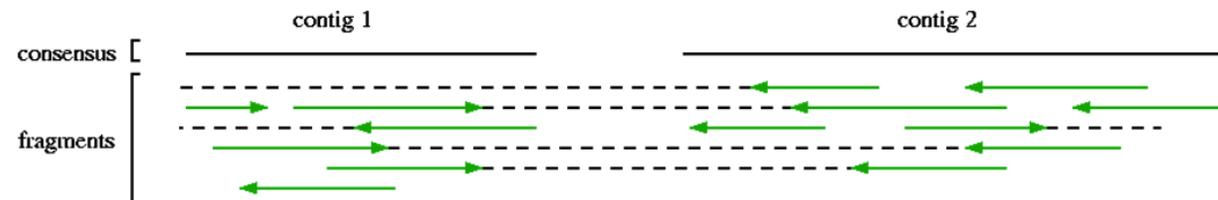
СБОРКА И СКАФФОЛДИНГ



```
AAAАСТСГССТСГСТТАТСААССГАТСССССГСТАССТТСТАСАГССАТСАТТ  
AAAАСТСГССТСГСТТАТСААССГАТСССССГСТАССТТСТАСАГССАТСАТТ  
AAAАСТСГССТСГСТТАТСААССГАТСССССГСТАССТТСТАСАГССАТСАТТ
```

Scaffolding
(скаффолдинг)

(скаффолд из 2 КОНТИГОВ)



OVERLAP-CONSENSUS-LAYOUT

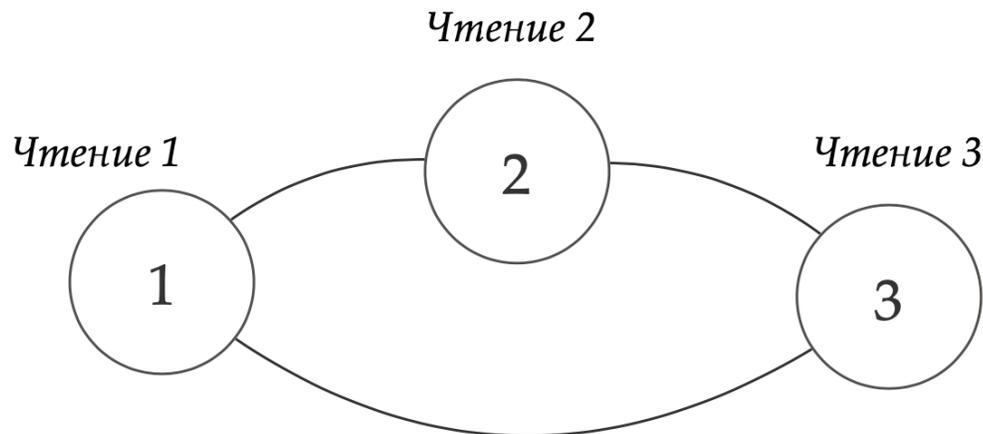
ATGCTA
ATGC
TGCT
GCTA



.fastq
ATGC (1)
TGCT (2)
GCTA (3)

OLC assembler

1-2
ATGC
TGCT
1-3
ATGC
GCTA
GCTA
TGCT
3-2



ПРИНЦИП OLC

Найти все возможные перекрытия между рядами (сложно).

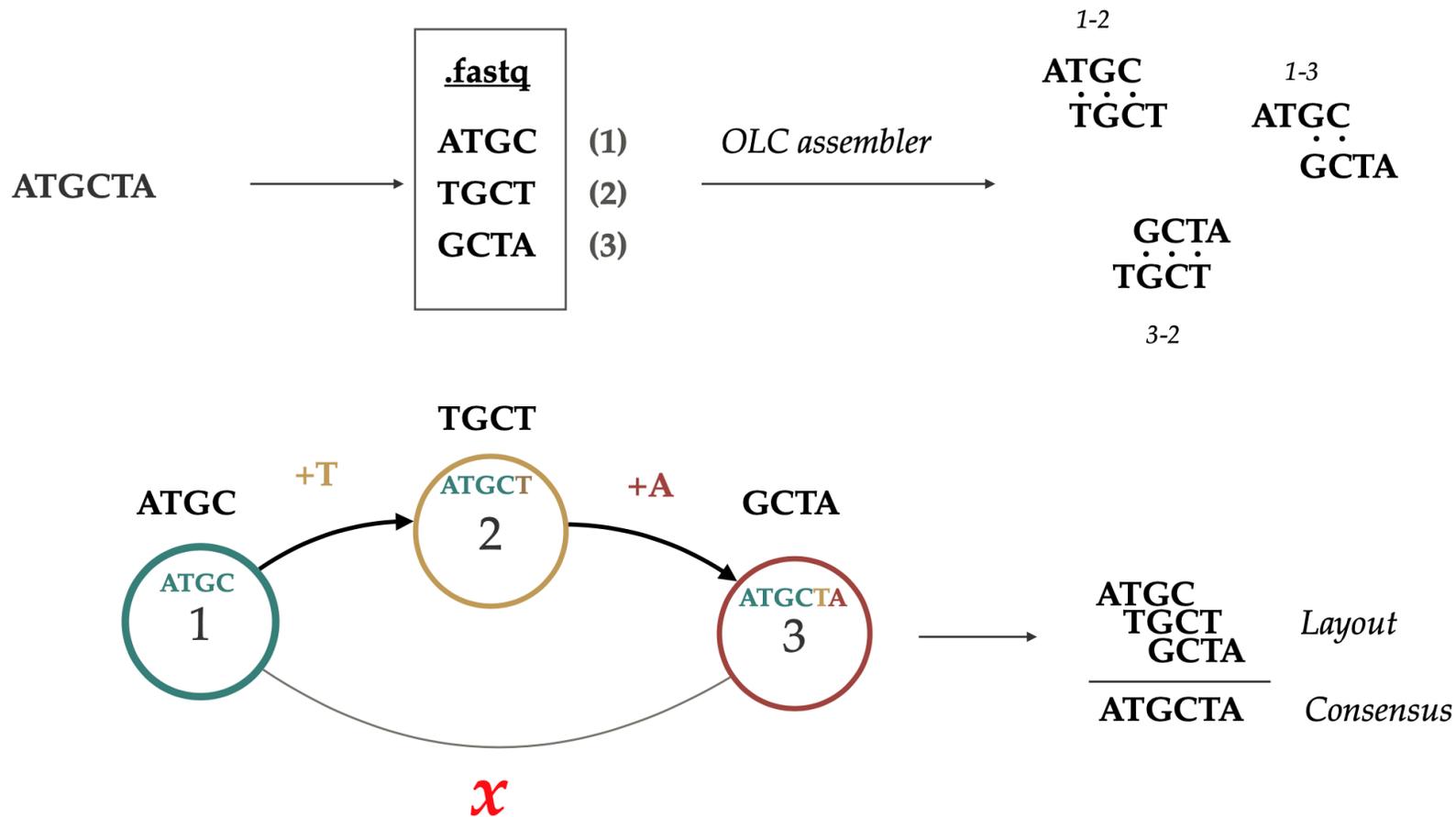
На данном этапе используются суффиксные деревья, динамическое программирование.

Допускается выравнивание, содержащее определенное число замен и гэпов (один из рядов может содержать ошибки)

Строится таблица (граф) перекрытий, в процессе выбираются только надежные перекрытия (обычно устанавливается критерий минимальной длины, % идентичности)

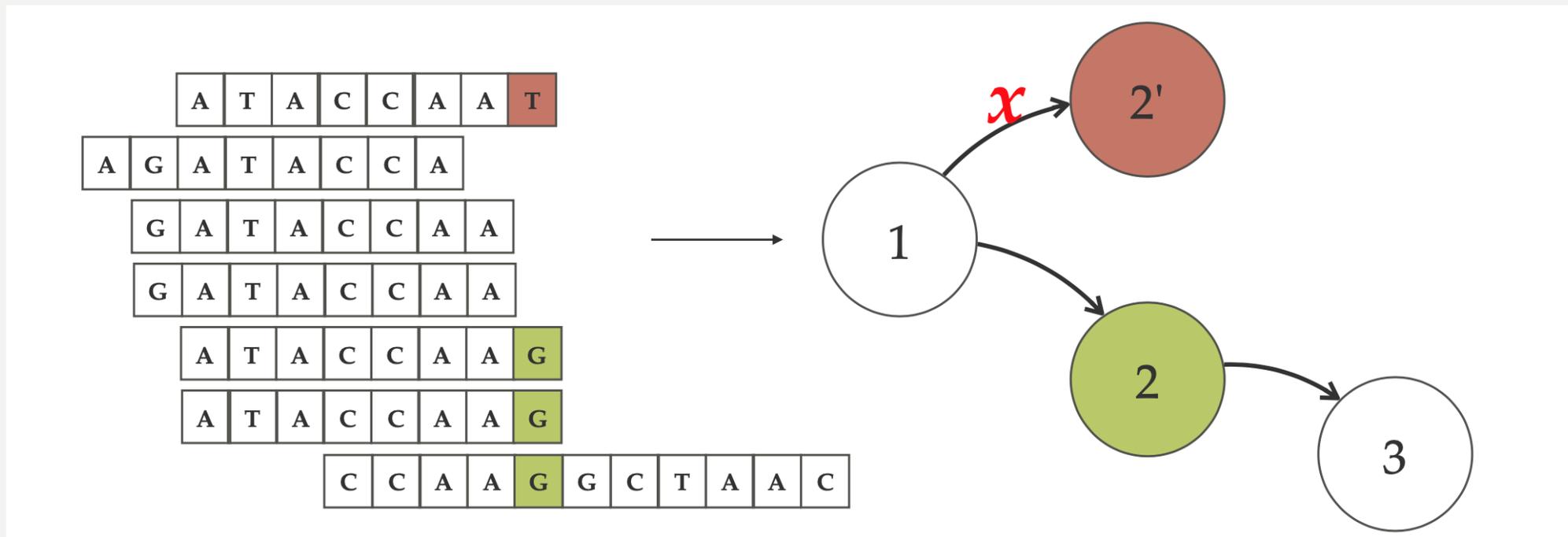
Ряды выравниваются друг с другом согласно графу перекрытий, образуя консенсусную последовательность

OVERLAP-CONSENSUS-LAYOUT

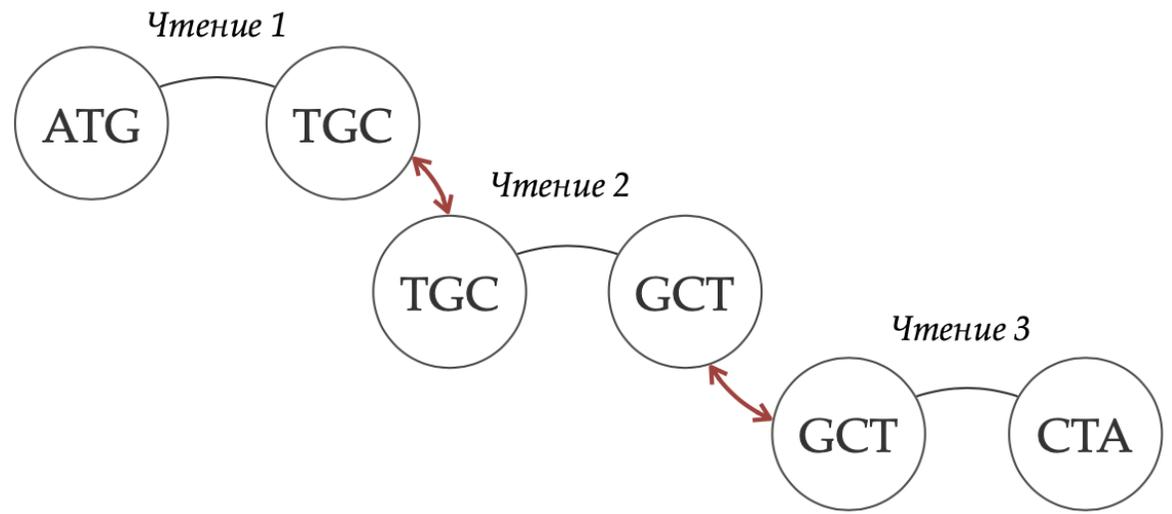
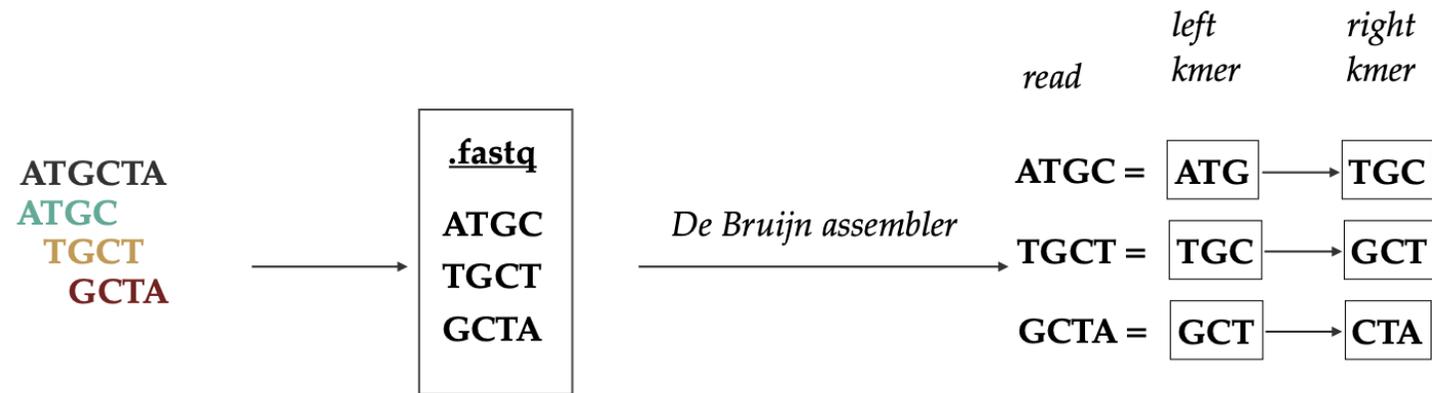


Минимальное перекрытие: 3

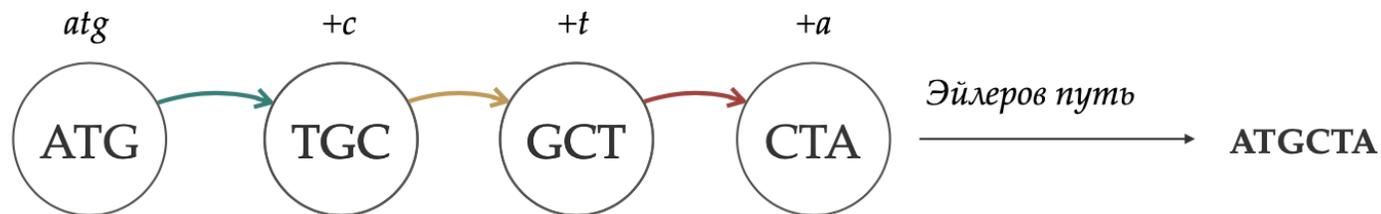
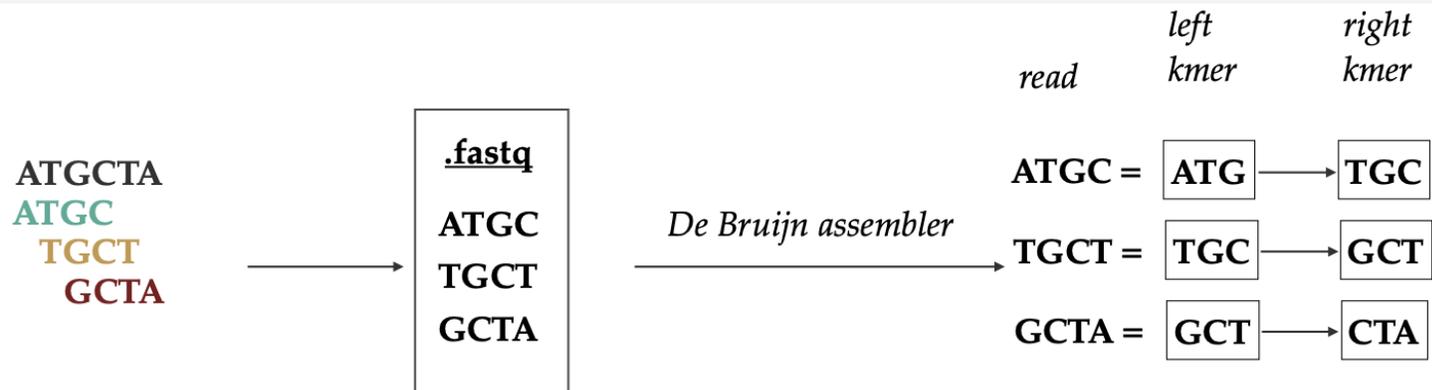
РАБОТА С LAYOUT



ГРАФЫ ДЕ БРЁЙНА (DBG)



ГРАФЫ ДЕ БРЁЙНА (DBG)

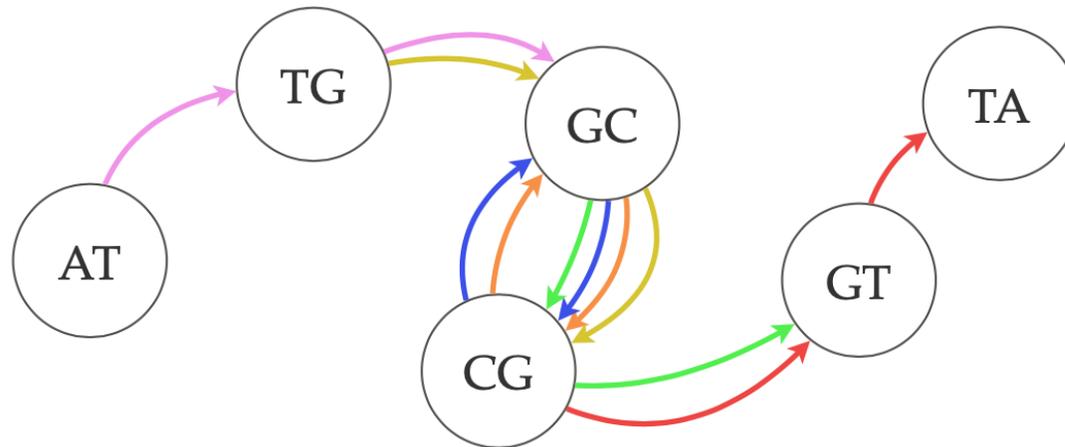


ГРАФЫ ДЕ БРЁЙНА (DBG)

ATGCGCGTA
ATGC
TGCG
GCGC
CGCG
GCGT
CGTA

`.fastq`
ATGC
TGCG
GCGC
CGCG
GCGT
CGTA

ATGC = AT → TG → GC
TGCG = TG → GC → CG
GCGC = GC → CG → GC
CGCG = CG → GC → CG
GCGT = GC → CG → GT
CGTA = CG → GT → TA

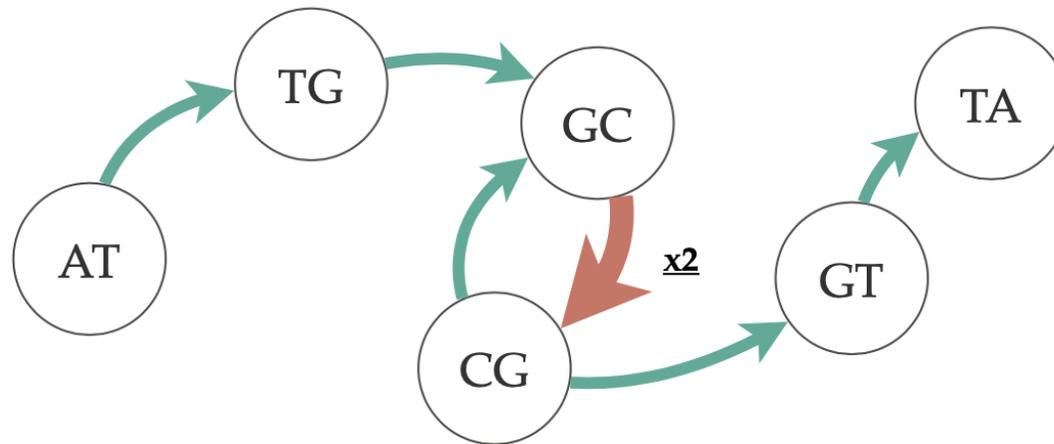


ГРАФЫ ДЕ БРЁЙНА (DBG)

1 2 3 4 4 4 3 2 1
ATGCGCGTA
ATGC
TGGC
GCGC
CGCG
GCGT
CGTA

.fastq
ATGC
TGGC
GCGC
CGCG
GCGT
CGTA

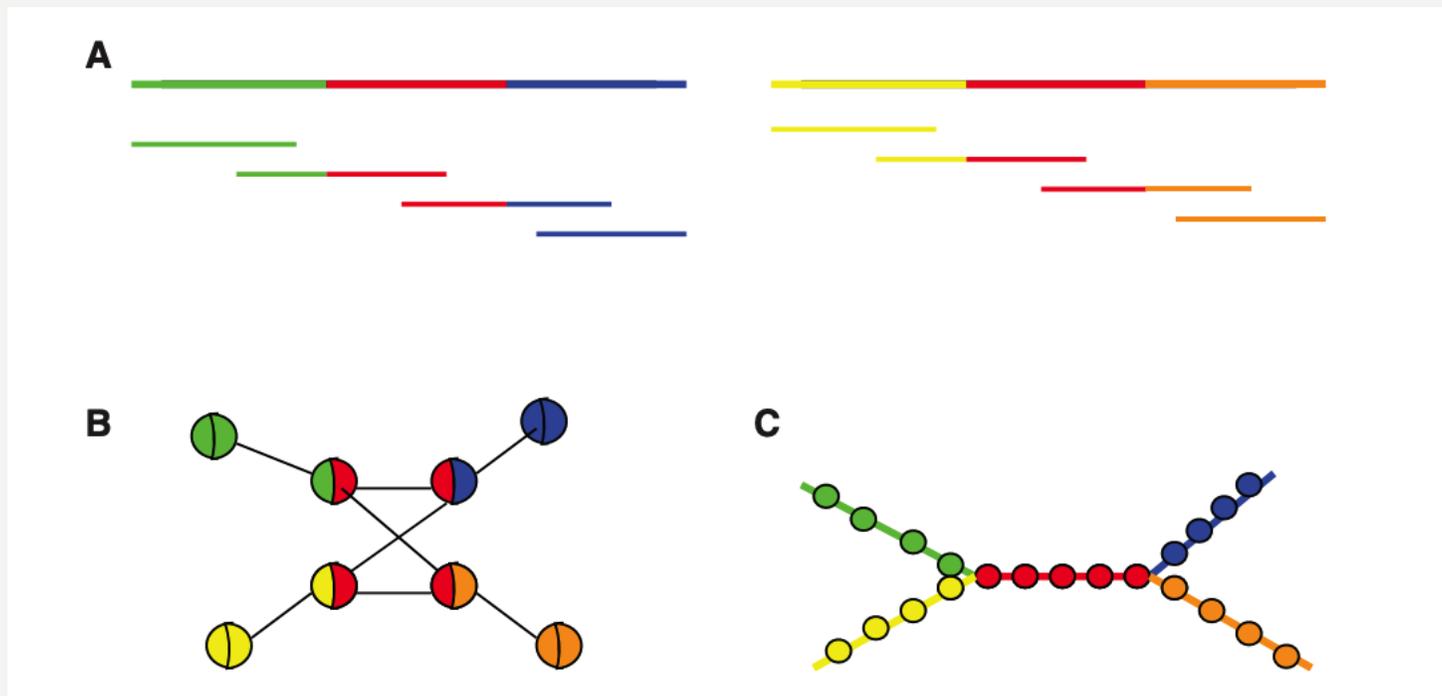
1
2 ← среднее покрытие нуклеотида
3
4



OLC И DBG

OLC:
рид – узел, перекрытие - ребро

DBG:
рид – ребро или ребра, kmer - узел



РАБОТА С ПОВТОРАМИ

Кроме ошибок в данных и естественной гетерогенности (например, гетерозиготные состояния аллеля в диплоидном геноме или естественный полиморфизм в популяции), повторы – основной источник сложности при сборке

Существуют методы маскирования повторов до сборки

AGCTTTTTTGCA

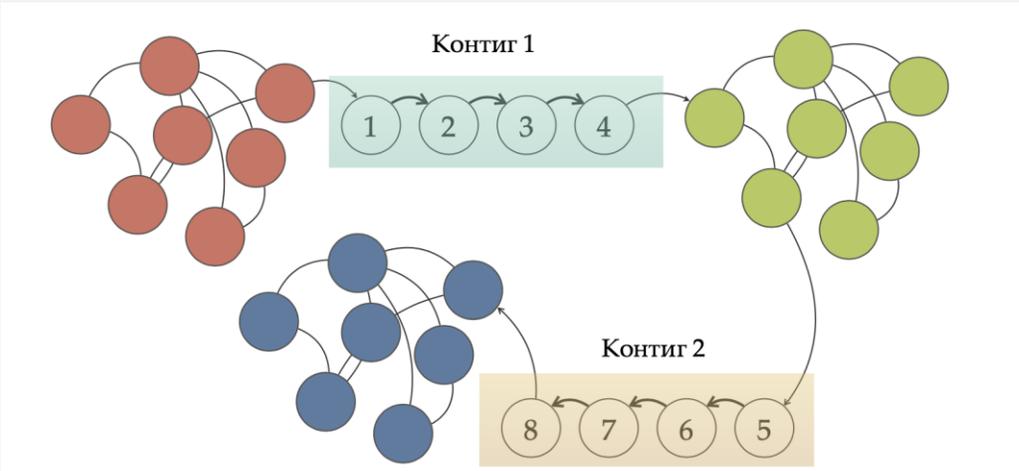
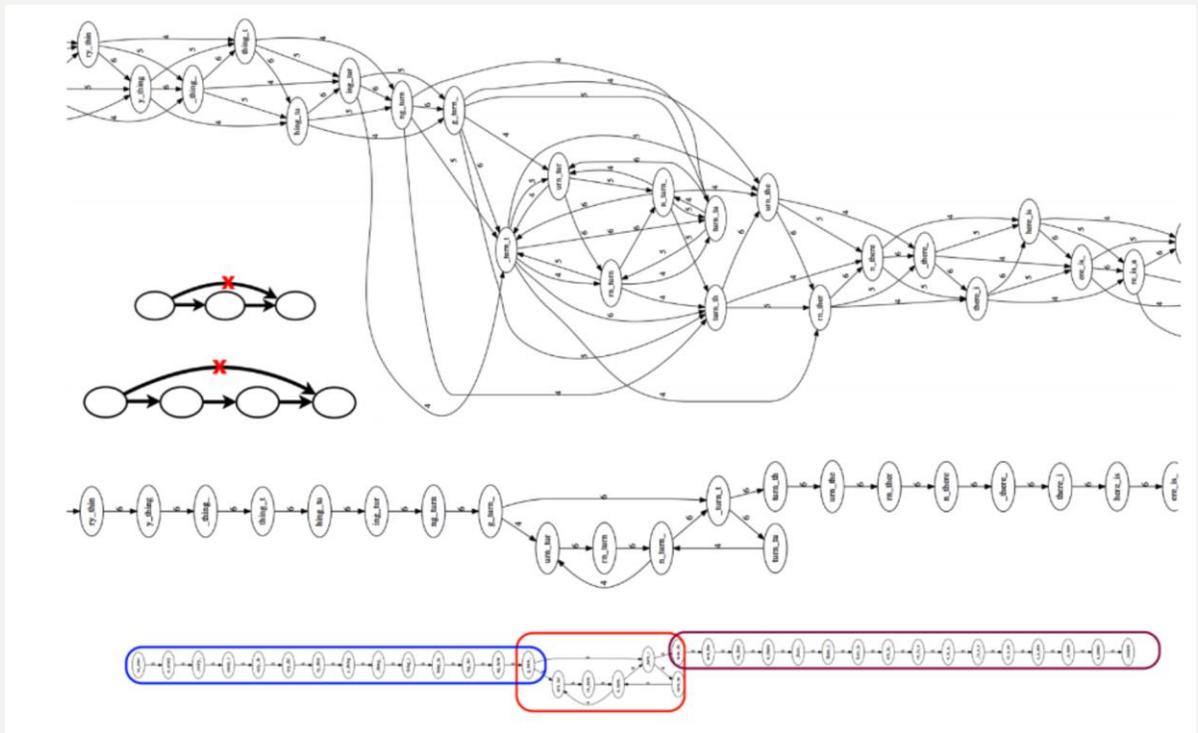
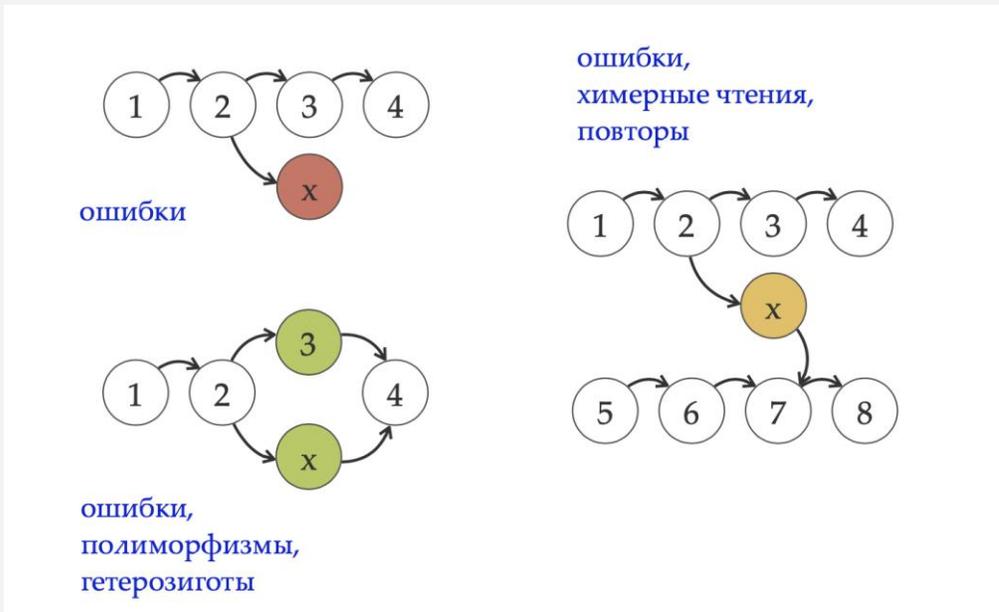
AGC → TTT → GCA

DBG: повторы не сильно усложняют граф, так как риды, приходящие из повтора, дают уже существующие в графе kmers

AGCTTTTTTGCA
CTTT
TTTTT
TTTTT
TTGCA

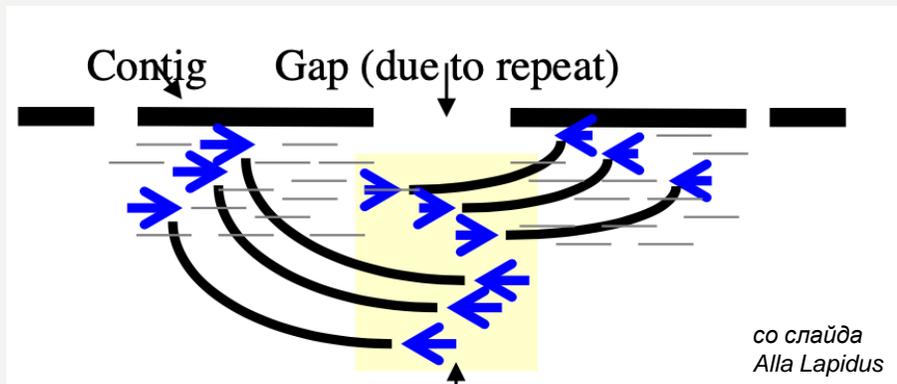
OLC: повторы усложняют граф, так как каждый рид дает в графе один новый узел; повторы увеличивают сложность вычисления, так как риды повторов имеют множество перекрытий (создают много ребер)

РАБОТА С ГРАФОМ: КОНТИГИ

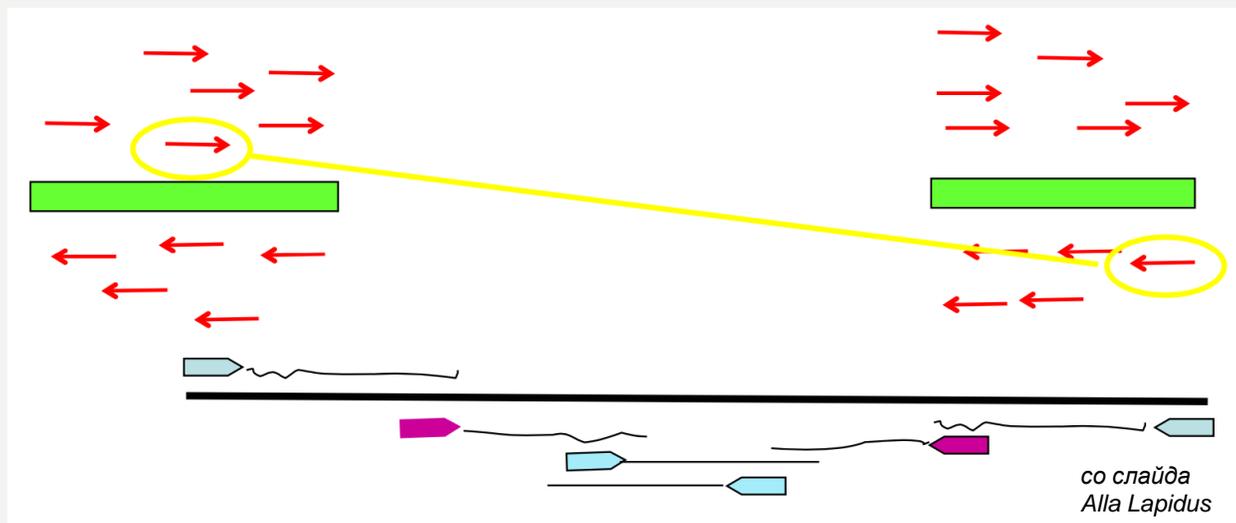


УЛУЧШЕНИЕ СБОРКИ

Локальная пересборка



Локальное досеквенирование



ДРУГИЕ ВИДЫ И СПОСОБЫ СБОРКИ

- Оптическое картирование
- Генетические карты
- Карты контактов HiC
- По имеющемуся референсу
- Транскриптомная сборка

КАЖДЫЙ БИОИНФОРМАТИК РАЗ В ЖИЗНИ...

Название	Алгоритм	Технологии	Авторы	Представлен
ABYSS	De Bruijn	Solexa, SOLiD	Simpson, J. et al.	2008
ALLPATHS-LG	De Bruijn	Solexa, SOLiD	Gnerre, S. et al.	2011
Celera WGA Assembler / CABOG	OLC	Sanger, 454, Illumina	Myers, G. et al.; Miller G. et al.	2004
CLC Genomics Workbench	String Graph	Sanger, 454, Solexa, SOLiD	CLC bio	2008
Edena	OLC	Illumina	D. Hernandez et al.	2008
Euler	De Bruijn	Sanger, 454 (,Solexa ?)	Pevzner, P. et al.	2001
Euler-sr	De Bruijn	454, Solexa	Chaisson, MJ. et al.	2008
IDBA	De Bruijn	Sanger,454,Solexa	Yu Peng, Henry C. M. Leung, Siu-Ming Yiu, Francis Y. L. Chin	2010
MIRA	OLC	Sanger, 454, Solexa	Chevreur, B.	1998
Newbler	String Graph	454, Sanger, Solexa, Ion	454/Roche	2009
PCAP	OLC	Sanger, 454	Huang et al.	2003
SGA	String Graph	Illumina, Ion Torrent	Simpson, J. et al.	2011
SOPRA		Illumina, SOLiD, Sanger, 454	Dayarian, A. et al.	2010
SOAPdenovo	De Bruijn	Solexa	Li, R. et al.	2009
SPAdes	De Bruijn	Illumina, Solexa	Bankevich, A et al.	2012
Velvet	De Bruijn	Sanger, 454, Solexa, SOLiD	Zerbino, D. et al.	2007