



СБМ 

НАУЧНО-ИССЛЕДОВАТЕЛЬСКИЙ ИНСТИТУТ
СИСТЕМНОЙ БИОЛОГИИ И МЕДИЦИНЫ
РОСПОТРЕБНАДЗОРА



ВВЕДЕНИЕ В БИОИНФОРМАТИКУ



*Андрей Евгеньевич САМОЙЛОВ
н.с. лаб. математической биологии
и биоинформатики
НИИ СБМ Роспотребнадзора*

ПЛАН ЛЕКЦИИ

- Что такое биоинформатика и зачем она нужна
- Основные понятия и форматы данных
- SARS-CoV-2 и как его секвенировать
- Оценка качества сборки генома
- Заключительные мысли

ЧТО ТАКОЕ БИОИНФОРМАТИКА?

БИОИНФОРМАТИКА...

- решает с вычислительной точки зрения крупномасштабные биологические проблемы (анализ “больших данных”).
- изучает и разрабатывает компьютерные методы в биологии.
- получает, анализирует, хранит, организует и визуализирует биологические данные.
- объединяет общую и молекулярную биологию, кибернетику, генетику, химию, компьютерные науки, математику, статистику.

КАК ВЫГЛЯДЯТ БИОИНФОРМАТИЧЕСКИЕ ИНСТРУМЕНТЫ

BLAST® » blastn suite

blastn blastp blastx tblastn tblastx **Standard Nucleotide BLAST**

BLASTn programs search nucleotide databases using a nucleotide query. more...

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) Clear

Query subrange From To

Or, upload file Выберите файл Файл не выбран

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database Standard databases (nr etc.): rRNA/ITS databases Genomic + transcript databases Betacoronavirus

Nucleotide collection (nr/nt)

Organism Enter organism name or id—completions will be suggested exclude

Exclude Models (XM/XP) Uncultured/environmental sample sequences

Limit to Sequences from type material

Entrez Query Create custom database

для обычного человека

```
(blast) samoilov_ae@node2:~$ blastn -h
USAGE
blastn [-h] [-help] [-import_search_strategy filename]
[-export_search_strategy filename] [-task task_name] [-db database_name]
[-dbsize num_letters] [-glist filename] [-seqidlist filename]
[-negative_glist filename] [-negative_seqidlist filename]
[-taxids taxids] [-negative_taxids taxids] [-taxidlist filename]
[-negative_taxidlist filename] [-entrez_query entrez_query]
[-db_soft_mask filtering_algorithm] [-db_hard_mask filtering_algorithm]
[-subject subject_input_file] [-subject_loc range] [-query input_file]
[-out output_file] [-evaluate evaluate] [-word_size int_value]
[-gapopen open_penalty] [-gapextend extend_penalty]
[-perc_identity float_value] [-qcov_hsp_perc float_value]
[-max_hsps int_value] [-xdrop_ungap float_value] [-xdrop_gap float_value]
[-xdrop_gap_final float_value] [-searchsp int_value]
[-sum_stats bool_value] [-penalty penalty] [-reward reward] [-no_greedy]
[-min_raw_gapped_score int_value] [-template_type type]
[-template_length int_value] [-dust DUST_options]
[-filtering_db filtering_database]
[-window_masker_taxid window_masker_taxid]
[-window_masker_db window_masker_db] [-soft_masking soft_masking]
[-ungapped] [-culling_limit int_value] [-best_hit_overhang float_value]
[-best_hit_score_edge float_value] [-subject_besthit]
[-window_size int_value] [-off_diagonal_range int_value]
[-use_index boolean] [-index_name string] [-lcase_masking]
[-query_loc range] [-strand strand] [-parse_deflines] [-outfmt format]
[-show_gis] [-num_descriptions int_value] [-num_alignments int_value]
[-line_length line_length] [-html] [-sorthits sort_hits]
[-sorthsps sort_hsps] [-max_target_seqs num_sequences]
[-num_threads int_value] [-mt_mode int_value] [-remote] [-version]
```

DESCRIPTION
Nucleotide-Nucleotide BLAST 2.13.0+

Use '-help' to print detailed descriptions of command line arguments
(blast) samoilov_ae@node2:~\$

для биоинформатика

ЗАЧЕМ ВАМ РАЗБИРАТЬСЯ В БИОИНФОРМАТИКЕ?

1

Чтобы заниматься
биоинформатикой

2

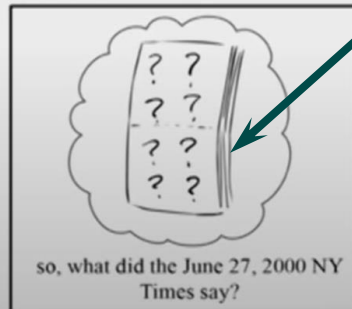
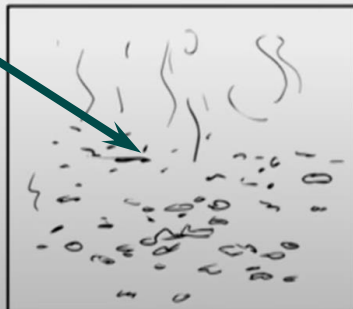
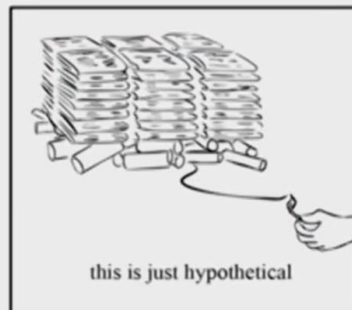
Чтобы понимать
биоинформатиков

2

Так сказала
А.Ю. Попова
в интервью

ЧТО ТАКОЕ СБОРКА ГЕНОМА?

The Newspaper Problem



fastq

fasta

ФОРМАТЫ ДАННЫХ: fasta

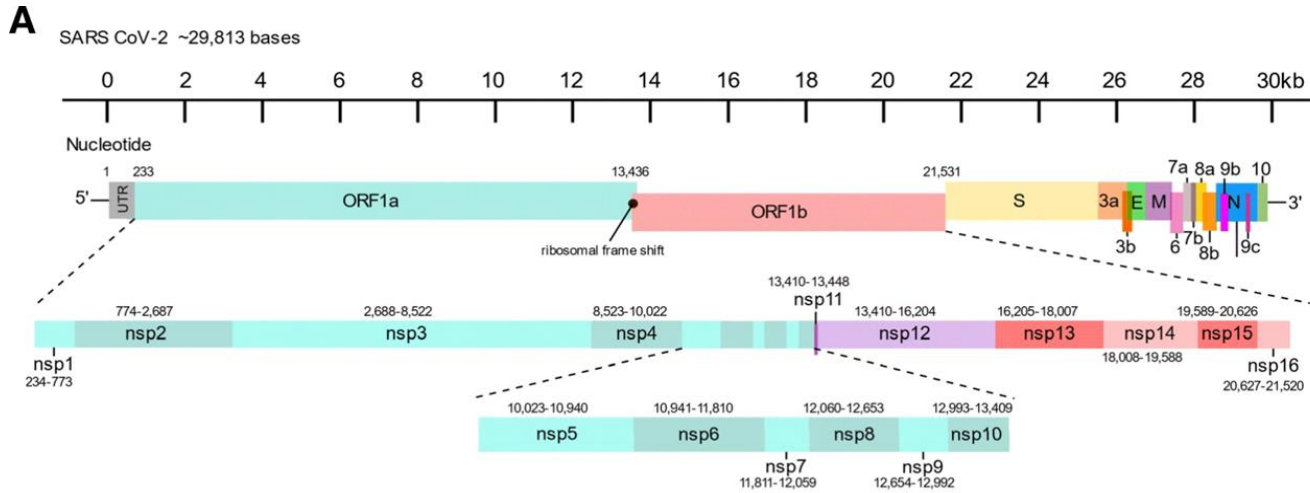
>NC_045512.2 Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1, complete genome

```
ATTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAACTTTTCGATCTCTTGTAGATCTGTTCTCTAAA  
CGAACTTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCACTCACGCAGTATAATTAATAAC  
TAATTA CTGTCGTTGACAGGACACGAGTAACTCGTCTATCTTCTGCAGGCTGCTTACGGTTTCGTCCGTG  
TTGCAGCCGATCATCAGCACATCTAGGTTTTCGTCCGGGTGTGACCGAAAGGTAAGATGGAGAGCCTTGTC  
CCTGGTTTTCAACGAGAAAACACACGTCCA ACTCAGTTTGCCTGTTTTACAGGTTTCGCGACGTGCTCGTAC  
GTGGCTTTGGAGACTCCGTGGAGGAGGTCTTATCAGAGGCACGTCAACATCTTAAAGATGGCACTTGTGG  
CTTAGTAGAAGTTGAAAAAGGCGTTTTTGCCTCAACTTGAACAGCCCTATGTGTT CATCAAACGTTCCGGAT  
GCTCGAACTGCACCTCATGGTCATGTTATGGTTGAGCTGGTAGCAGAACTCGAAGGCATT CAGTACGGTC  
GTAGTGGTGAGACACTTGGTGTCTTGTCCCTCATGTGGGCGAAATACCAGTGGCTTACCGCAAGGTTCT  
TCTTCGTAAGAACGGTAATAAAGGAGCTGGTGGCCATAGTTACGGCGCCGATCTAAAGTCATTTGACTTA  
GGCGACGAGCTTGGCACTGATCCTTATGAAGATTTTCAAGAAA ACTGGAACACTAAACATAGCAGTGGTG  
TTACCCGTGAACTCATGCGTGAGCTTAACGGAGGGGCATACACTCGCTATGTGATAACA ACTTCTGTGG  
CCCTGATGGCTACCCTCTTGAGTGCATTAAGACCTTCTAGCACGTGCTGGTAAAGCTTCATGCACTTTG  
TCCGAACA ACTGGACTTTATTGACACTAAGAGGGGTGTATACTGCTGCCGTGAACATGAGCATGAAATTG
```


1. Сборка *de novo*
 - a. не требуется референсная последовательность
 - b. подходит всегда, но результаты сложнее анализировать

2. Сборка с референсом
 - a. нужно знать, что именно собираем, и иметь достаточно близкий референсный геном
 - b. результат - готовый геном, который можно анализировать сразу

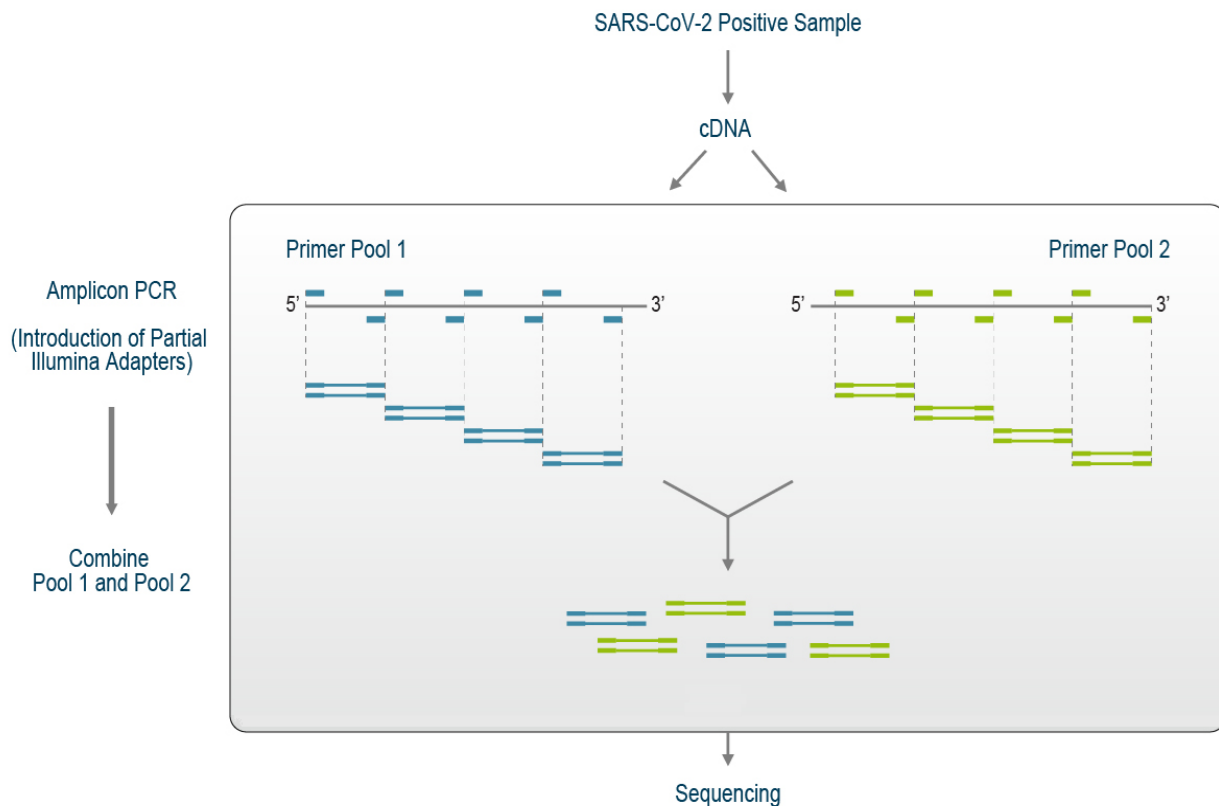
СТРУКТУРА ГЕНОМА SARS-CoV-2



C Percent identity matrix – full genome

1: SARS-CoV-2	100.00	86.85	81.25	81.58	79.34	78.40	80.09	96.75
2: SARS-CoV	86.85	100.00	78.31	77.09	75.59	74.81	76.30	86.56
3: MERS-CoV	81.25	78.31	100.00	79.39	77.76	77.06	78.20	80.81
4: HKU1	81.58	77.09	79.39	100.00	83.89	79.68	83.86	80.40
5: OC43	79.34	75.59	77.76	83.89	100.00	77.50	77.97	78.27
6: 229E	78.40	74.81	77.06	79.68	77.50	100.00	82.86	77.58
7: NL63	80.09	76.30	78.20	83.86	77.97	82.86	100.00	79.03
8: RaTG13	96.75	86.56	80.81	80.40	78.27	77.58	79.03	100.00

ПРОТОКОЛЫ ARTIC И MIDNIGHT



Без амплификации доля генетического материала SARS-CoV-2 в образце может составлять около 1% или существенно меньше!

ПРОТОКОЛЫ ARTIC И MIDNIGHT

ARTIC

- 99 пар праймеров, разбитых на 2 пула
- характерный размер ампликона - 300 п.о.

Midnight

- 29 пар праймеров, разбитых на 2 пула
- характерный размер ампликона - 1200 п.о.



ПРОТОКОЛЫ ARTIC И MIDNIGHT: НЕДОСТАТКИ

Так как почти все их используют, то у всех одни и те же проблемы:

- появление мутации в месте посадки праймера \Rightarrow амплификация не происходит \Rightarrow часть генома невозможно отсеквенировать

или

- появление мутации в месте посадки праймера \Rightarrow не можем её найти, так как мешают праймеры (если вдруг не удалось от них избавиться в процессе обработки)

*В протоколе ARTIC праймеры занимают около 15% генома.
По указанным причинам до 30% геномов варианта Дельта
и 60% геномов линии BA.2 (Омикрон)
содержат меньше мутаций, чем должны!*

КРИТЕРИИ КАЧЕСТВА ГЕНОМНОЙ СБОРКИ

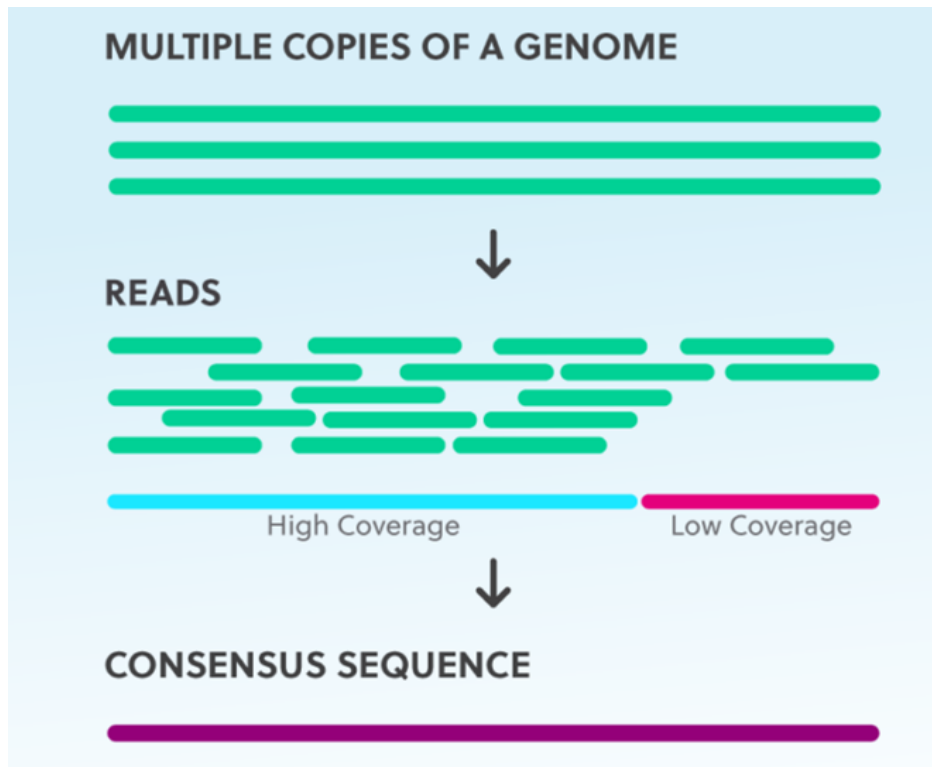
1. Глубина покрытия

- a. покрытие отдельного нуклеотида - сколько раз мы прочитали этот нуклеотид
- b. покрытие генома - сколько в среднем прочтений приходится на каждый нуклеотид генома

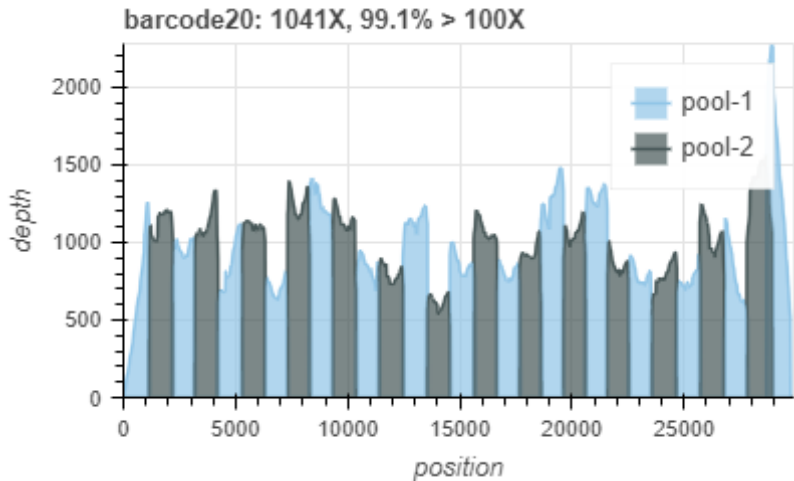
Для обработки данных SARS-CoV-2 платформы Oxford Nanopore по умолчанию достаточным считается покрытие 20, если покрытие меньше, ставится буква N.

2. Ширина покрытия - доля длины генома с достаточным покрытием.

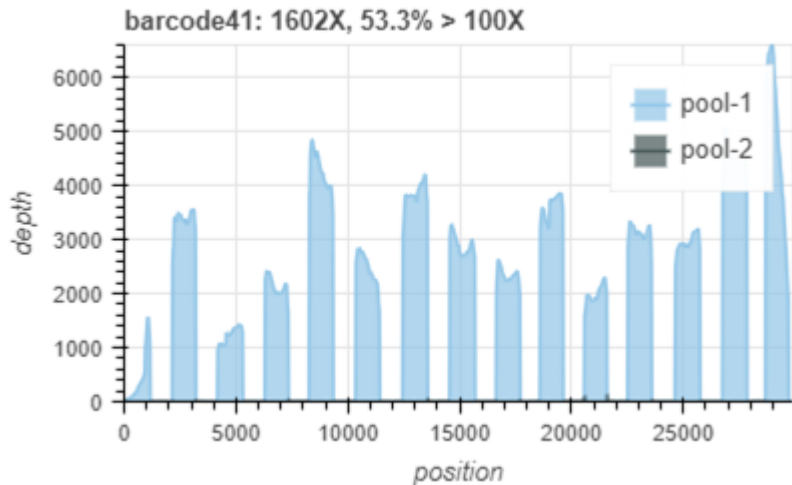
ПОКРЫТИЕ КАК ВАЖНЫЙ КРИТЕРИЙ СБОРКИ



ОДНОРОДНОСТЬ ПОКРЫТИЯ - ЭТО ВАЖНО!

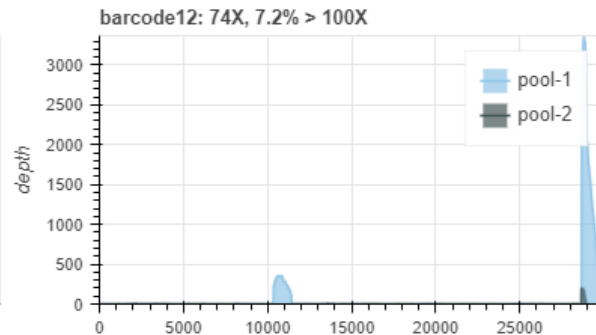
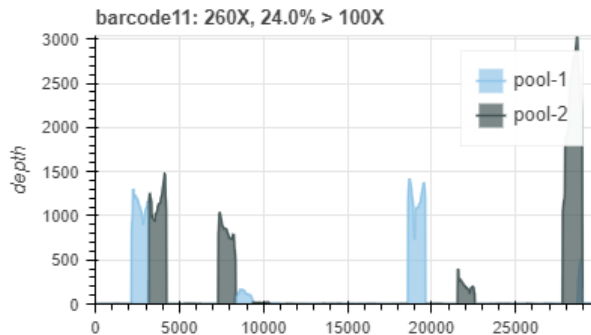
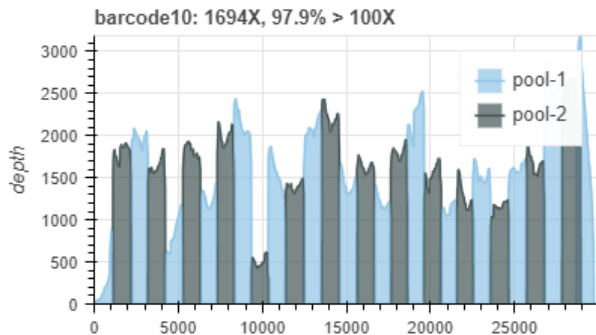
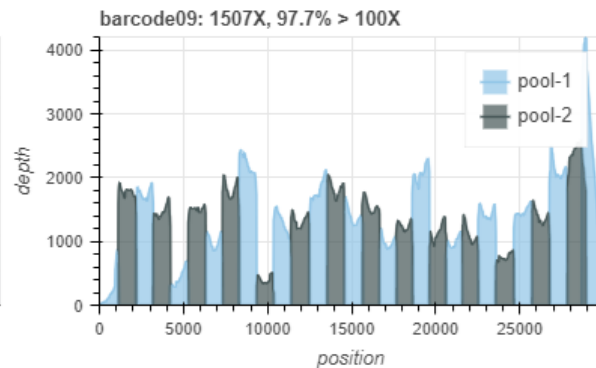
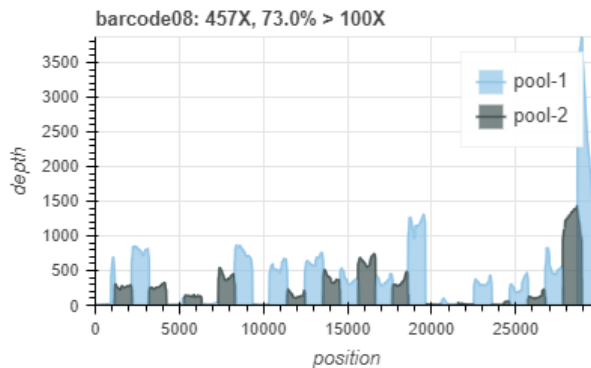
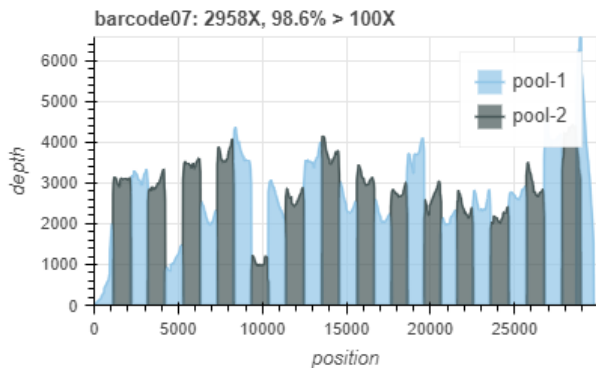


Однородное покрытие



Неоднородное покрытие
(нет амплификации в пуле № 2)

ГРУСТНАЯ (НО ТИПИЧНАЯ) СИТУАЦИЯ В ЖИЗНИ...



КАК ЕЩЁ МОЖНО ПРИМЕНИТЬ ВЫСОКОПРОИЗВОДИТЕЛЬНОЕ СЕКВЕНИРОВАНИЕ И БИОИНФОРМАТИКУ?

- Диагностика инфекционных заболеваний
- Исследование внутрибольничных инфекций
- Поиск генов антибиотикорезистентности
- Валидация бактериальных и вирусных коллекций

... и многое другое!

ГДЕ УЧИТЬСЯ БИОИНФОРМАТИКЕ?

Bioinformatics Institute

Подписаться 1M подписчиков

[Bioinformatics Institute](#) is the first non-governmental non-commercial bioinformatics research and educational institution in Russia.

Bioinformatics Institute's mission is to educate the new generation of highly qualified specialists in the field of bioinformatics and to popularize Bioinformatics in Russia.

The Institute's programs and events are open to students and professionals with Mathematics, Software Engineering, or Biology background. Apart from several offline educational initiatives, we develop a series of open online courses.

🔍 Быстрый поиск курсов (по названию или ID)

Программирование на Python

Bioinformatics Institute



★ 4.7 👤 678K 🕒 19 ч 🗨️

Бесплатно

Основы статистики

Bioinformatics Institute



★ 4.9 👤 238K 🕒 7 ч 🗨️

Бесплатно

Python: основы и применение

Bioinformatics Institute



Введение в Linux

Bioinformatics Institute



<https://stepik.org/org/bioinf>

WWW.SYSBIOMED.RU



СБМ 

НАУЧНО-ИССЛЕДОВАТЕЛЬСКИЙ ИНСТИТУТ
СИСТЕМНОЙ БИОЛОГИИ И МЕДИЦИНЫ
РОСПОТРЕБНАДЗОРА

WWW.SYSBIOMED.RU

Спасибо за внимание!