

Основи секвенування нового покоління (NGS)

Анна Яручик

Експерт Бюро ВООЗ в Україні



План:

1. Терміни.
2. Огляд основних технологій NGS.
3. Аналіз даних NGS.
4. Бази даних.
5. Значення NGS та геномного нагляду в охороні здоров'я.



Терміни:



- **Принцип комплементарності:** чітка відповідність нуклеотидів між двома ланцюгами ДНК, які за допомогою водневих зв'язків формують парні комплекси: аденін в одному ланцюзі протилежний тиміну (або урацилу в РНК) в іншому ланцюзі, а цитозин в одному ланцюзі протилежний гуаніну в іншому.
- **Секвенування нового покоління (Next Generation Sequencing, NGS, масивне паралельне секвенування)** - це високопродуктивна технологія секвенування нуклеїнових кислот, яка дозволяє швидко та ефективно визначити послідовність багатьох молекул ДНК/РНК одночасно в одному об'ємі біохімічної реакції.
- **Зчитування (Read):** послідовність одного ДНК-фрагмента, яка секвенується.
- **Бібліотека (Library):** колекція ДНК-фрагментів з приєднаними адаптерами, підготовлена до секвенування.
- **Адаптери:** олігонуклеотиди з відомими послідовностями, які лігуються з ДНК-фрагментами для зв'язування їх з матрицею в інструменті NGS.
- **Лігування:** це з'єднання двох фрагментів нуклеїнової кислоти ферментом лігазою.
- **Індекс (баркод):** коротка послідовність олігонуклеотидів, прикріплених до зразка для його ідентифікації при змішуванні з іншими зразками. Сам індекс також секвенується, таким чином діючи як «штрих-код» для ідентифікації зразка.

Підходи секвенування

- **Таргетне секвенування:** секвенування лише обраних ділянок геному за допомогою праймерів або зондів, які специфічно зв'язуються з цільовими ділянками геному. Це дозволяє зменшити обсяг даних, отриманих в процесі секвенування, зменшуючи час і витрати на аналіз даних.
- **Секвенування методом дробовика (Shotgun sequencing):** стратегія секвенування, яка полягає в розбитті випадковим чином геному на велику кількість коротких фрагментів, що секвенуються окремо та потім збираються в одну послідовність. Метод використовується для дослідження нової або невідомої ДНК, а також для секвенування метагеномів - зразків з довкілля, що містять суміш багатьох організмів, таких як ґрунт, вода, кишкові мікроорганізми тощо.
- **Повногеномне секвенування (Whole Genome Sequencing):** секвенування повного геному цільового виду.
- **Повноекзомне секвенування (Whole Exome Sequencing):** секвенування лише екзонів (ділянок, що кодують білки) цільового виду.

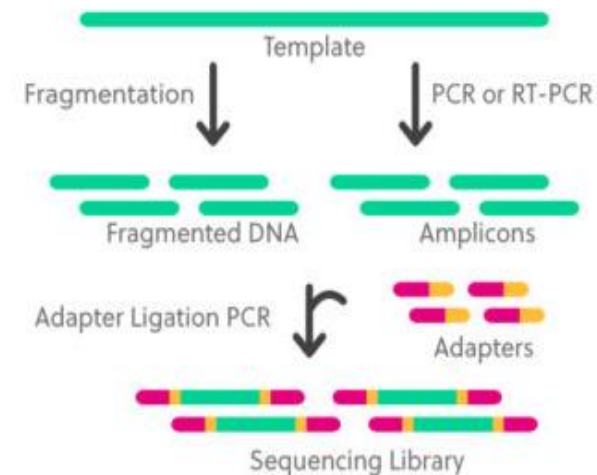
Основні етапи робочого процесу NGS

1. Екстракція ДНК/РНК зі зразків, оцінка кількості та якості.
* Зворотня транскрипція для РНК
2. Підготовка бібліотеки
3. Секвенування
4. Аналіз даних

STEP 1: Extraction



STEP 2: Library Prep



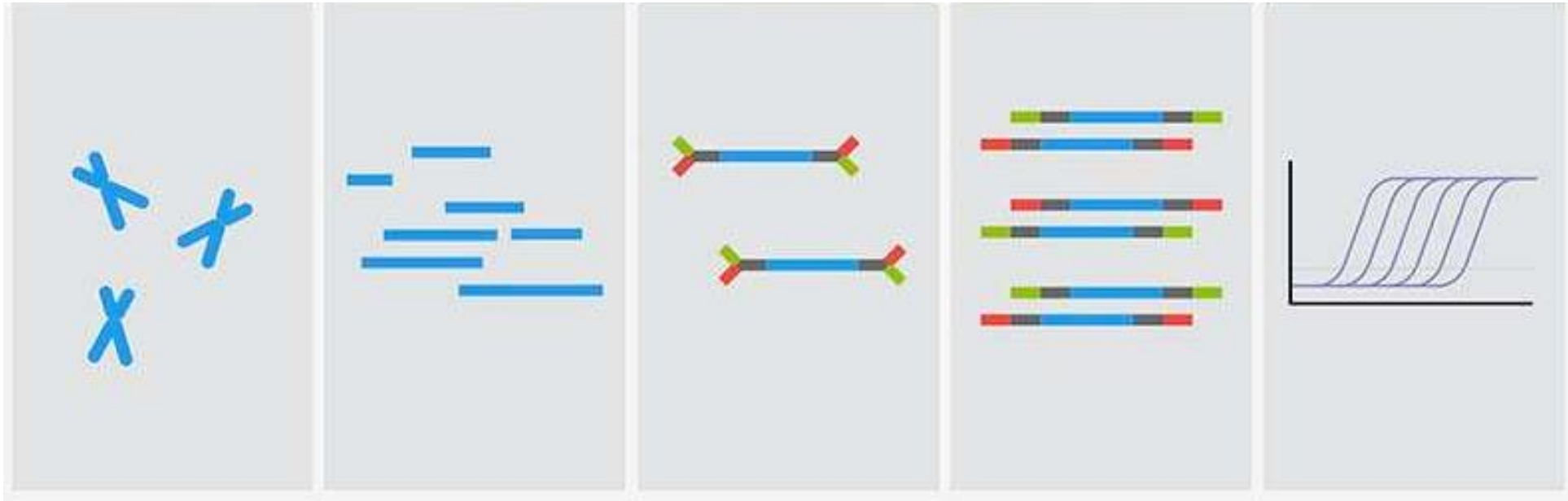
STEP 3: Sequencing



STEP 4: Analysis



Підготовка бібліотеки



Утворення
фрагментів ДНК

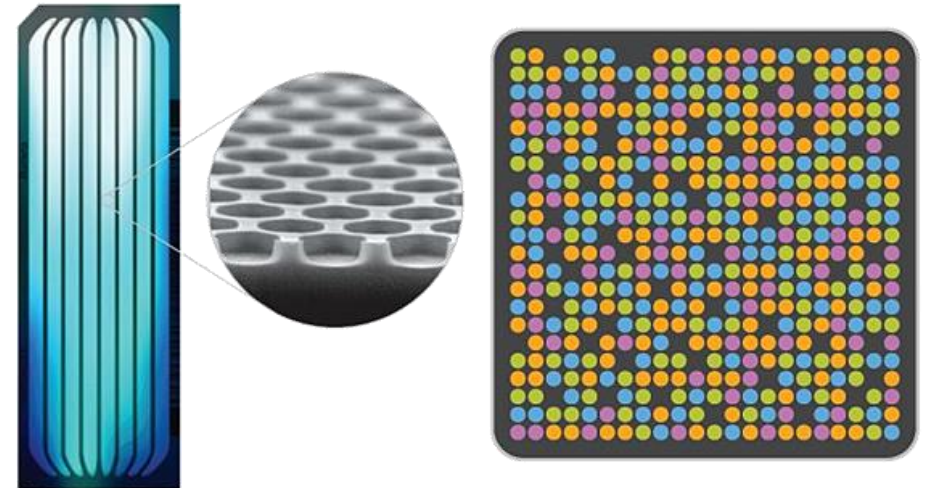
Приєднання адаптерів
та індексів

Ампліфікація та
відбір фрагментів
потрібної довжини

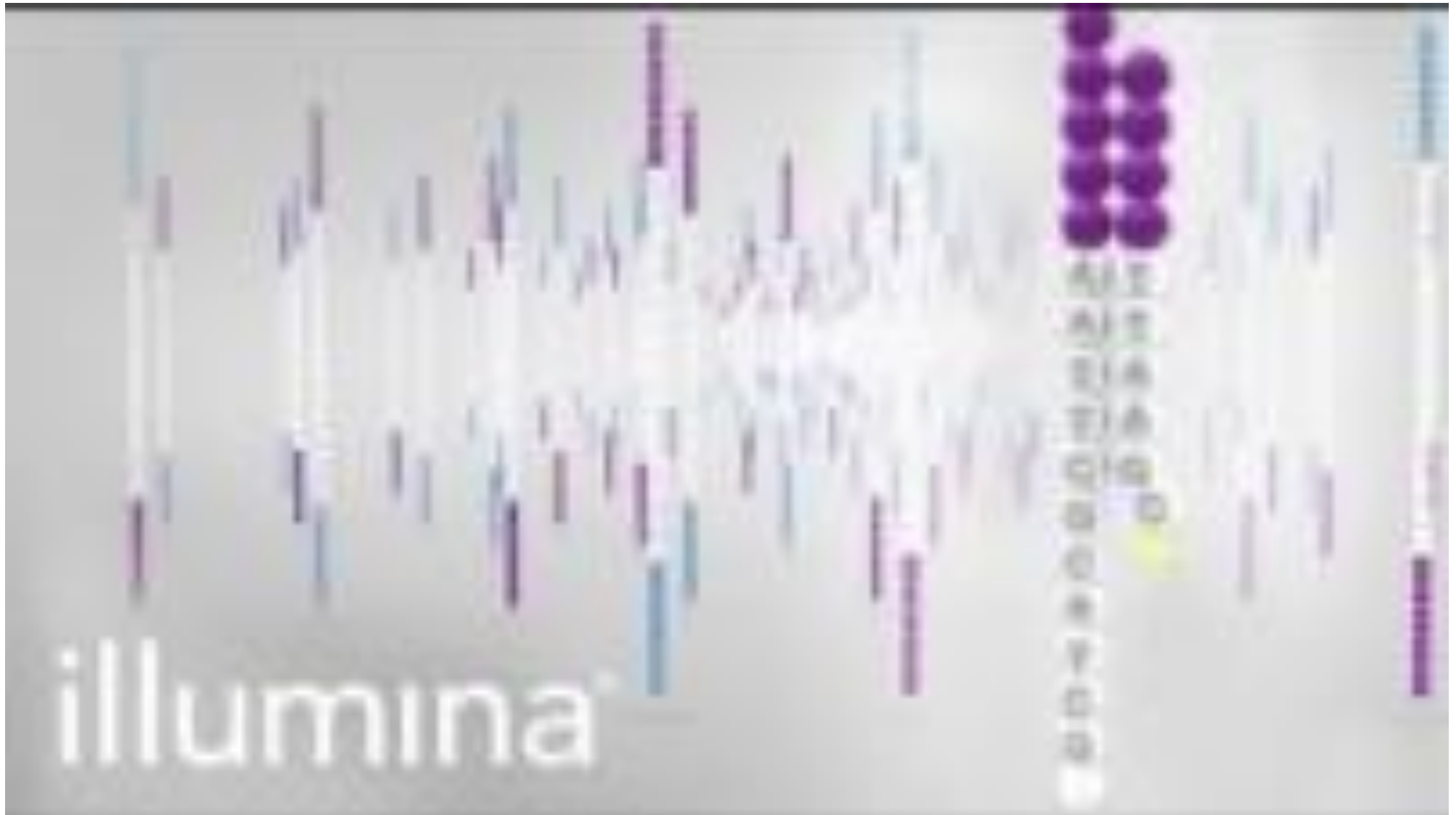
Оцінка кількості
бібліотеки

Секвенування Illumina

- Секвенування шляхом синтезу (SBS- Sequencing by synthesis)
- Технологія SBS використовує чотири флуоресцентно мічені нуклеотиди
- Генерування кластерів шляхом місткової ампліфікації
- Паралельне секвенування мільйонів ДНК-фрагментів на проточній лунці (flow cell)
- Короткі зчитування (до 300bp)
- Висока точність та низький рівень помилок



Секвенування Ілліміна



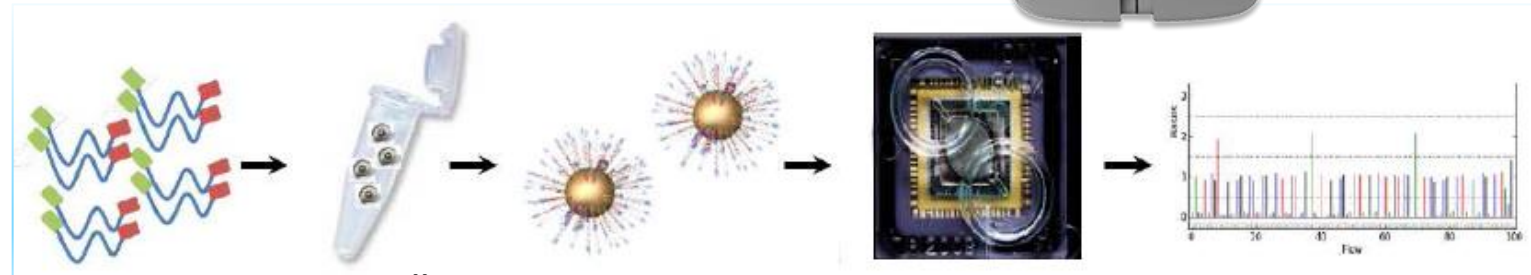
Відео:

<https://www.youtube.com/watch?v=fCd6B5HRaZ8>



European Region

Секвенування Ion Torrent



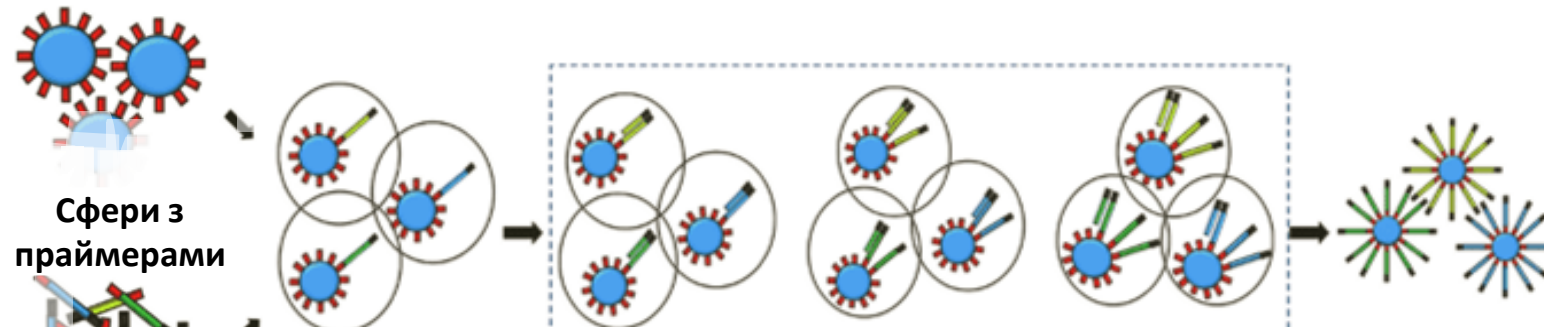
Приготування бібліотеки

Емульсійна ПЛР

Сфери ISP (Ion Sphere Particles) з копіями ДНК

Секвенування на чіпі

Аналіз даних



Сфери з праймерами

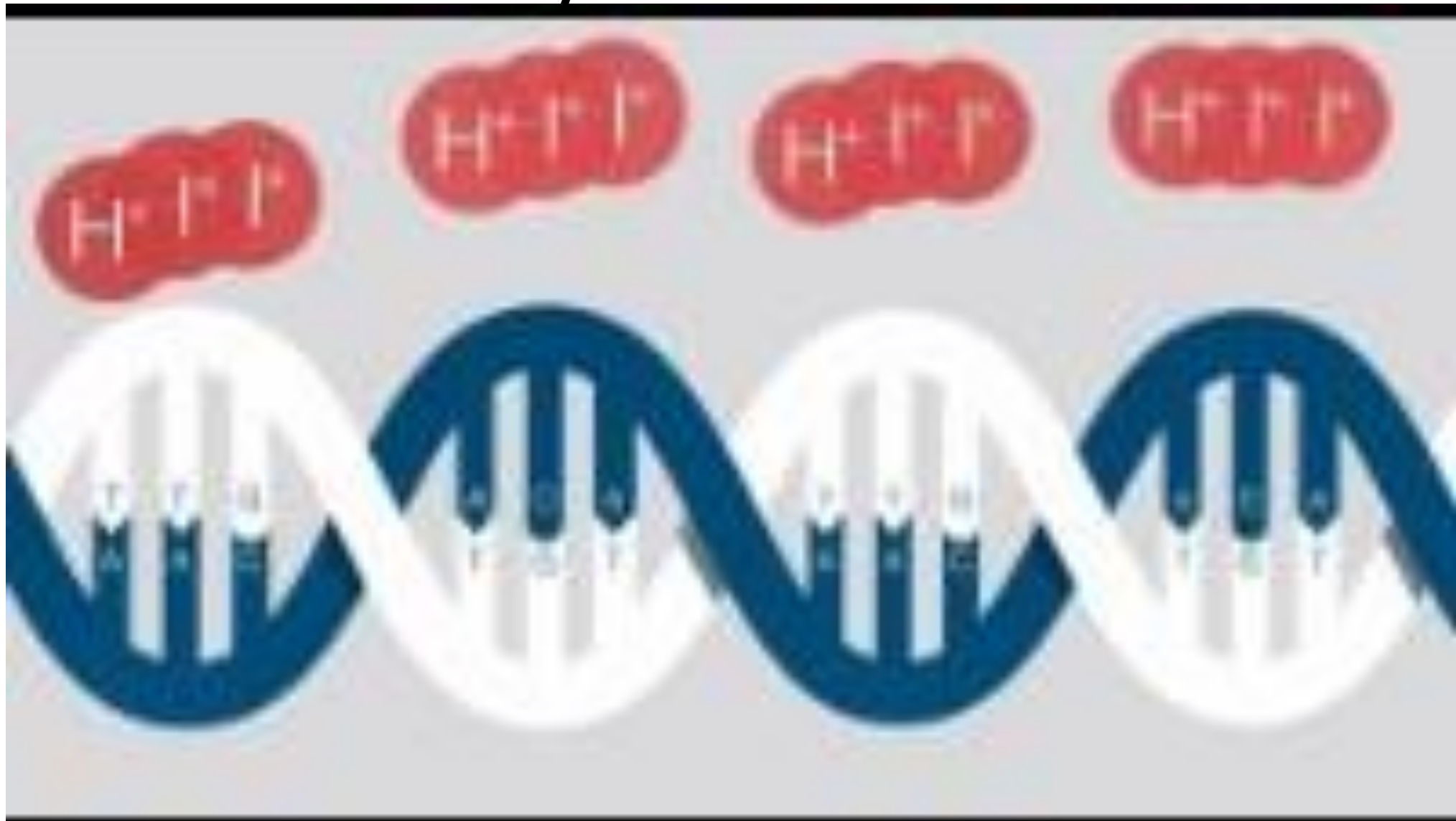
Бібліотека фрагментів ДНК

Емульсійна ПЛР

Клони ДНК

- Напівпровідникове секвенування на основі іонного потоку (ion semiconductor sequencing)
- Засноване на реєстрації йонів водню, виділених в середовище у ході синтезу полімеразою ДНК-фрагментів > зміна рН
- Клональна ампліфікація шляхом емульсійної ПЛР (ePCR): бібліотеки приєднуються до іонних сфер та шляхом ПЛР утворюється багато копій на поверхні сфери
- Короткі зчитування (200bp-600bp)
- Висока точність та низький рівень помилок

Секвенування Ion Torrent

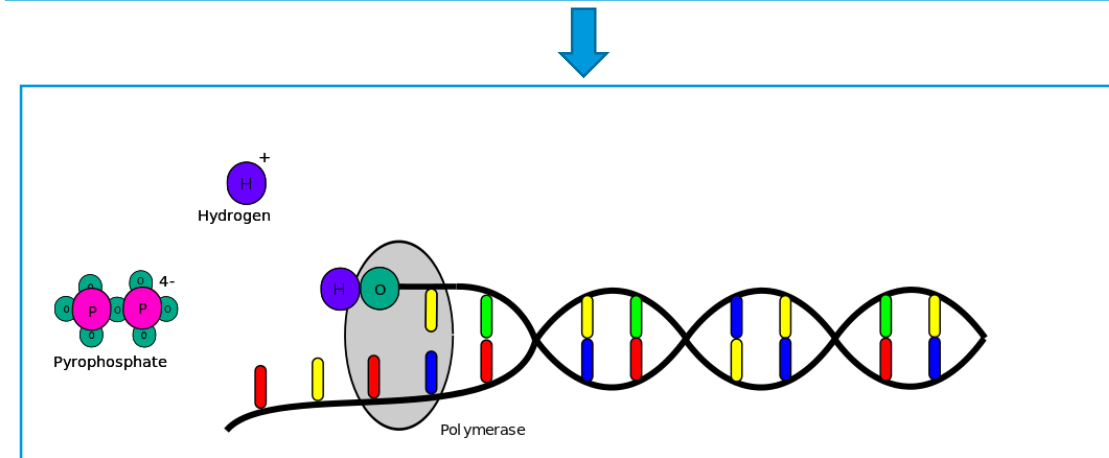
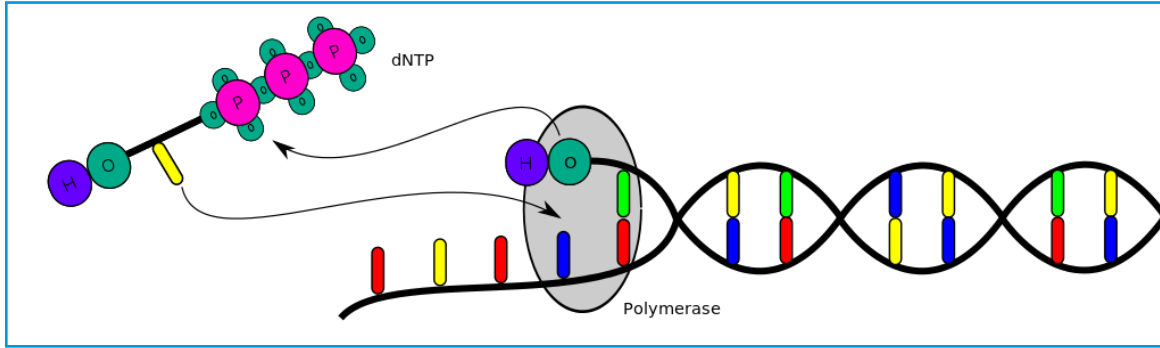


Відео:

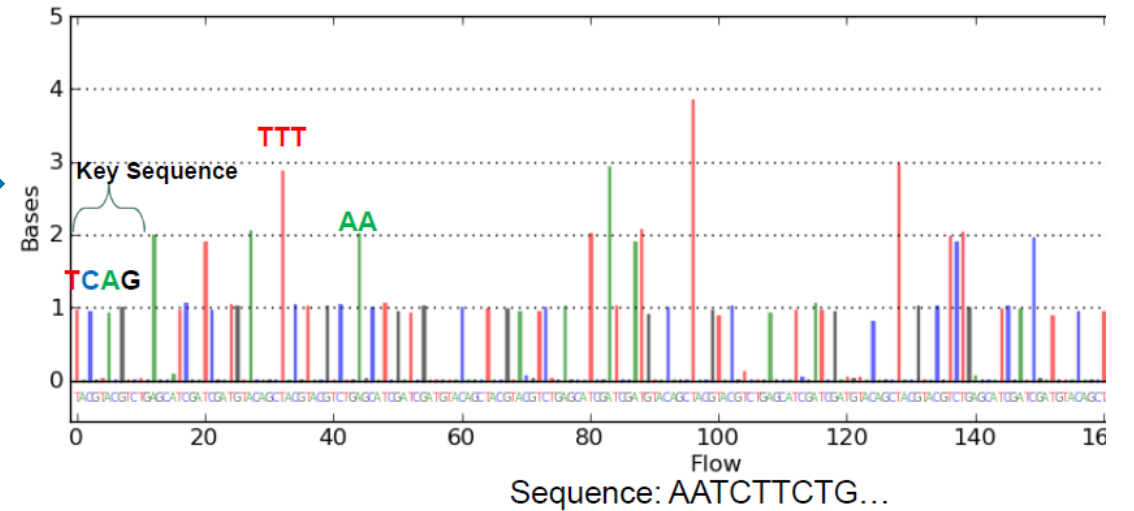
<https://youtu.be/zBPKj0mMcDg?t=25>

Секвенування Ion Torrent

Реакція приєднання полімеразою нуклеотиду до ланцюга ДНК, внаслідок якої вивільняється водень і пірофосфат

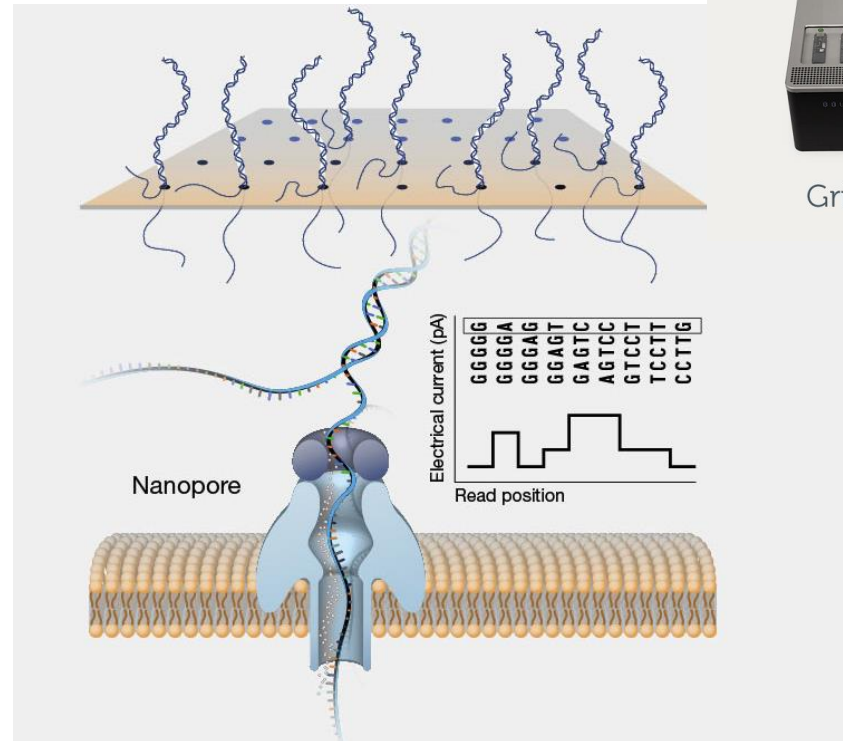


Детекція нуклеотидів в послідовності ДНК

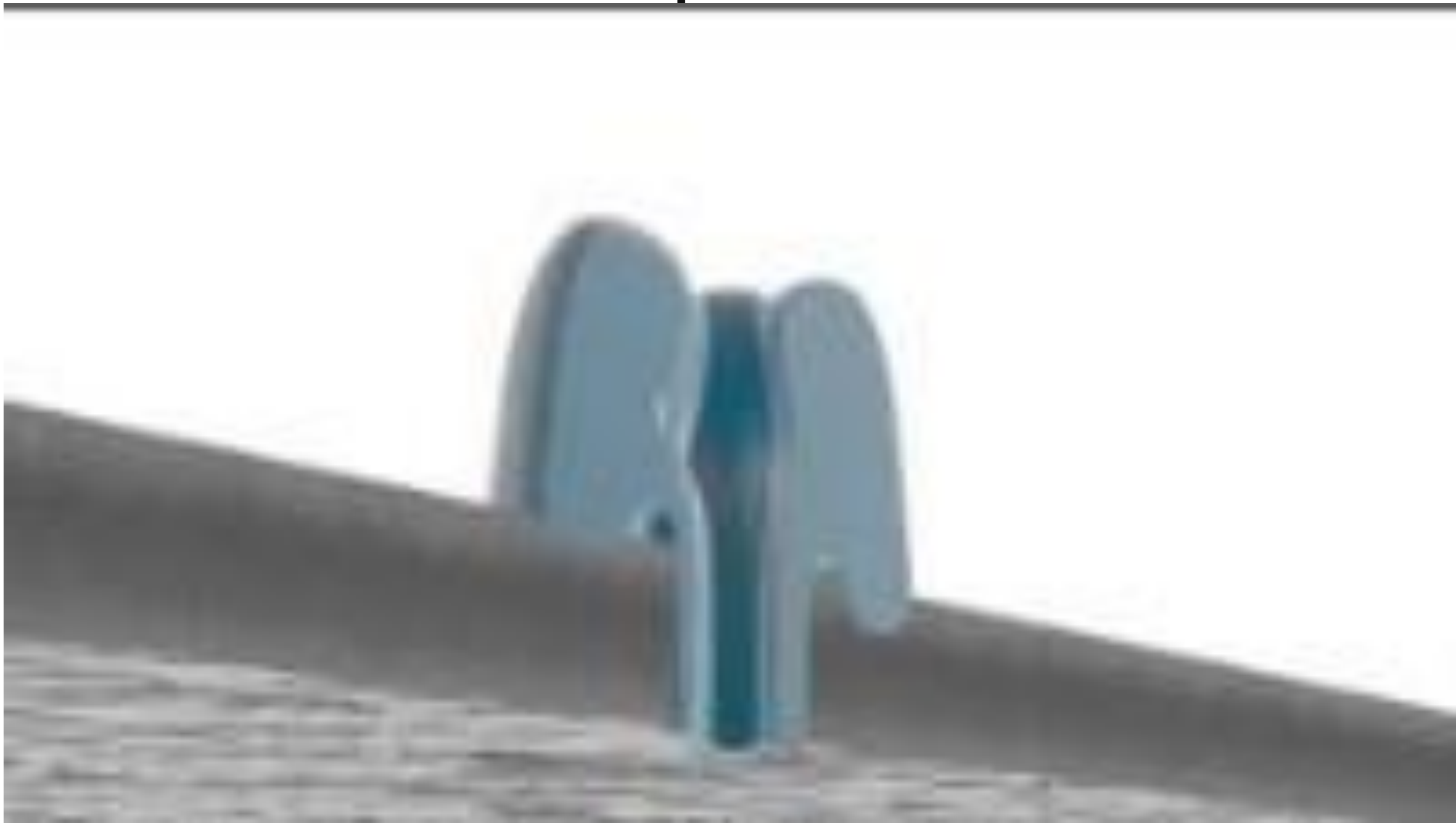


Oxford Nanopore

- Базується на використанні нанопор для вимірювання зміни струму, що протікає через пору під час проходження ДНК-фрагментів, що дає змогу визначити послідовність нуклеотидів.
- Можливість секвенувати одну молекулу ДНК або РНК без необхідності ампліфікації.
- Довгі зчитування (тисячі bp), що дає змогу вирішувати задачі, пов'язані з дослідженням структури та функцій геномів.
- Метод дозволяє проводити секвенування в режимі реального часу.
- Однією з особливостей є портативність та можливість використання в польових умовах.



Oxford Nanopore



Відео:

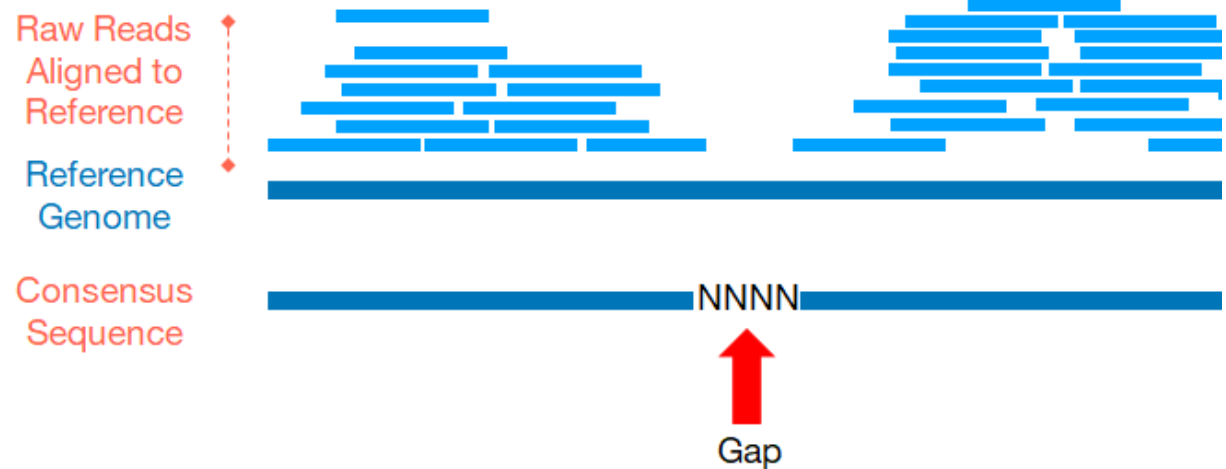
<https://www.youtube.com/watch?v=CGWZvHli3i0>



European Region

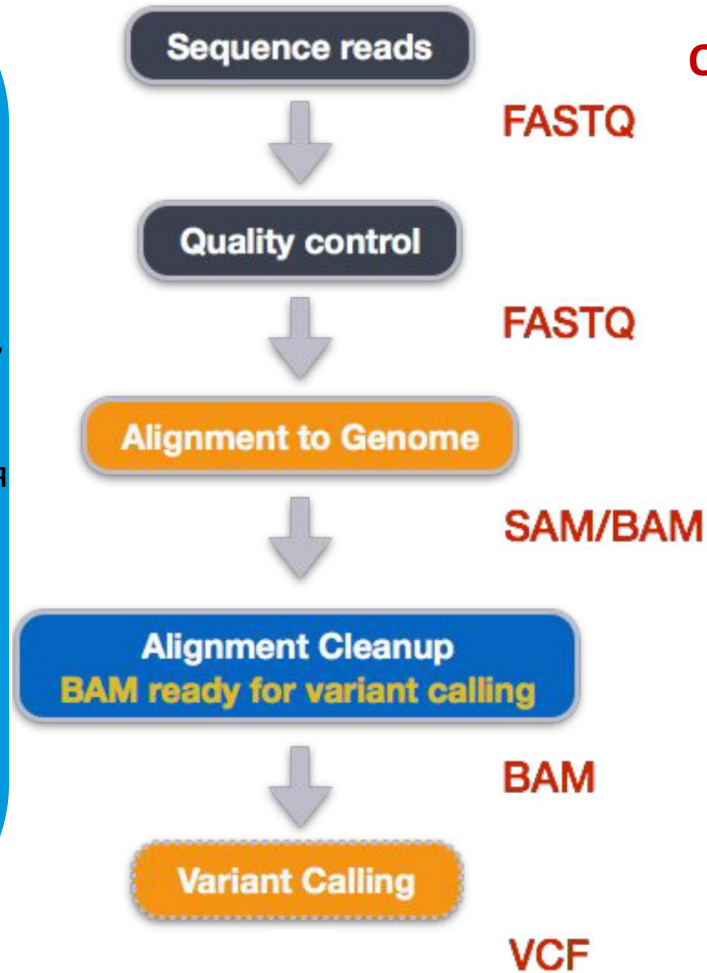
Аналіз даних NGS

- **Аналіз даних NGS** - це процес обробки та інтерпретації великих обсягів даних, отриманих з NGS-платформ, з використанням біоінформатичних інструментів з метою отримання консенсусної послідовності ДНК.
- **Консенсусна послідовність** - це узагальнена послідовність досліджуваної ДНК або РНК, яка представляє найбільш поширену основу в кожній позиції.
- **Референсна послідовність геному** - формально визнана офіційна послідовність відомого геному, зібрана вченими як репрезентативний приклад геному окремого виду, використовується як стандарт для порівняння досліджуваної ДНК цього виду.



Основні етапи аналізу даних

- 1. Визначення послідовності сирих даних (Basecalling)** - перетворення реєстрованих сигналів з секвенатора в послідовність нуклеотидів (наприклад, в форматі FASTQ).
- 2. Оцінка якості даних (Quality control)** - перевірка якості секвенування та обробки зчитувань, для визначення наявності помилок та можливих проблем, які потребують додаткової обробки.
- 3. Вирівнювання (Mapping/Alignment)** – вирівнювання усіх зчитувань та на референсний геном (якщо такий використовується) та збірка консенсусної послідовності.
- 4. Визначення варіантів (Variant calling)** - визначення відмінностей між секвенованою послідовністю та референсною послідовністю.
- 5. Інтерпретація результатів.**



Основні формати файлів даних

- **FASTQ file:** зберігає послідовності зчитувань, які отримані під час секвенування, та показники якості визначення кожного нуклеотиду.
- **SAM (Sequence Alignment/Map)** – стандартний текстовий формат файлу для зберігання результатів вирівнювання (Alignment) послідовностей.
- **BAM (Binary Alignment/Map)** - бінарний аналог формату SAM, який забезпечує швидше зберігання та обробку великих об'ємів даних.
- **VCF (Variant Call Format)** - формат файлу для зберігання інформації про варіанти (SNPs, InDels) геному порівняно з референсом.

Сирі дані: FASTQ file

1. унікальний ідентифікатор для кожного зчитування
2. рядок послідовності ДНК у вигляді тексту
3. рядок заголовка (+)
4. рядок символів, які кодують показник якості для кожної основи.

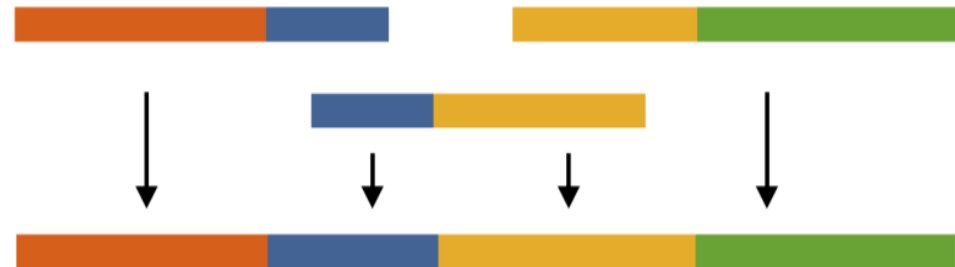
```
@K00359:71:HJJL7BBXX:3:1101:1996:1508 1:N:0:ATCACG
AAAATTCCAAGCTGGTTTCAACAGTACTTTGTTTCCAGAACAAAGAAATG
+
AAAFFJJJJJFJJ<J<FJJJJJJJJJJJJJJJJFJJFJJJJFFJJJJJJJ<
@K00359:71:HJJL7BBXX:3:1101:2240:1508 1:N:0:ATCACG
GTAAGGATGCGTAGGGATGGGAGGGCGATGAGGACTAGGATGATGGCGG
+
AAFFFJJJJJJJJF<J7JJFJJJJJJFFFJFJJJJJJJJJJJJJJJJJJJJ
@K00359:71:HJJL7BBXX:3:1101:2402:1508 1:N:0:ATCACG
GTCGACCATGTGGGCAGAACCTTGATGTTGGATTCCAGCAGGACCTGTCC
+
AAFFFJJJJJJJJJ<JJJJJJJJJJ<JFJJJJJJJJJJJJJJJJJJJJJJJJJJ
@K00359:71:HJJL7BBXX:3:1101:2463:1508 1:N:0:ATCACG
ATGTGGTGTATGCATCGGGTAGTCCGAGTAACGTCGGGGCATTCCGGAT
+
AAAFFJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
```


Вирівнювання послідовностей (Alignment)

- Вирівнювання на відому референсну послідовність



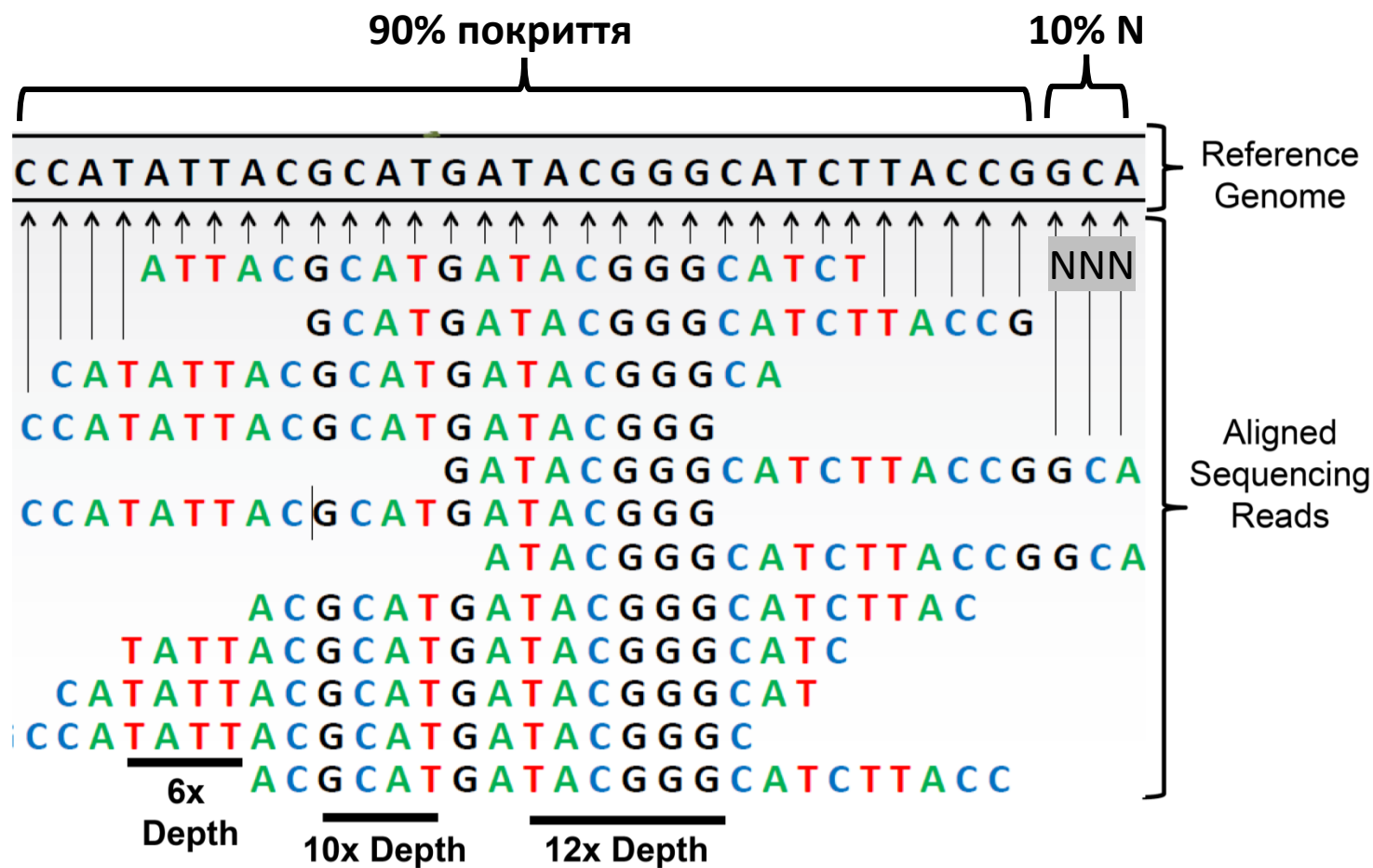
- De novo збірка геному: без попереднього знання послідовності шляхом виявлення взаємного перекриття зчитувань



Кількісні показники секвенування:

- **Довжина зчитування (ріда):** довжина секвенованого фрагмента в бібліотеці
- **Глибина покриття (Depth)** - кількість зчитувань вирівняних на окрему позицію геному (12x - 12-кратне покриття)
Чим більша глибина зчитування, тим вища впевненість у визначеному нуклеотиді.
- **Середнє покриття (Coverage):** середня кількість вирівняних зчитувань, які перекривають усі позиції на цільовому геномі.
- **Загальне покриття послідовності** - % геному, яка має достатню кількість рідів для аналізу
- **% N** - Відсоток основ у кожній позиції без визначеного нуклеотиду

Приклад: консенсусний геном охоплює 90%, тобто ми ідентифікували послідовність для 90% геному, але маємо 10% N з відсутніми даними.



Типи мутацій (варіантів)

Хромосомні



Точкові



Заміна: заміна одного нуклеотиду на інший

Делеція: втрата нуклеотиду

Інсерція: вставка додаткового нуклеотиду

Зсув рамки зчитування (Frameshift): коли додається або видаляється нуклеотид, змінюється вся послідовність кодонів, що йде за мутацією.

A T C C G A G T T

Стандартна послідовність

A T C C G C G T T

A T C ~~X~~ G A G T T

frameshift

A T C C T G A G T T

frameshift

IGV

(Integrative Genomics Viewer) – Інтегрований геномний переглядач

Sars-CoV-2 (ASM985889v3) NC_045512.2 NC_045512.2:22,769-22,880

113 bp

22 780 bp 22 800 bp 22 820 bp 22 840 bp

Click and drag to zoom in.

IonXpress_044_rawlib.realigned_processed.bam Coverage

IonXpress_044_rawlib.realigned_processed.bam

NC_045512.2:22 786

Total count: 351
A : 2 (1%, 0+, 2-)
C : 349 (99%, 54+, 295-)
G : 0
T : 0

IonXpr...

NC_045512.2:22 813

Total count: 5615
A : 0
C : 3 (0%, 3+, 0-)
G : 0
T : 5612 (100%, 3451+, 2161-)
N : 0

DEL: 2
INS: 8

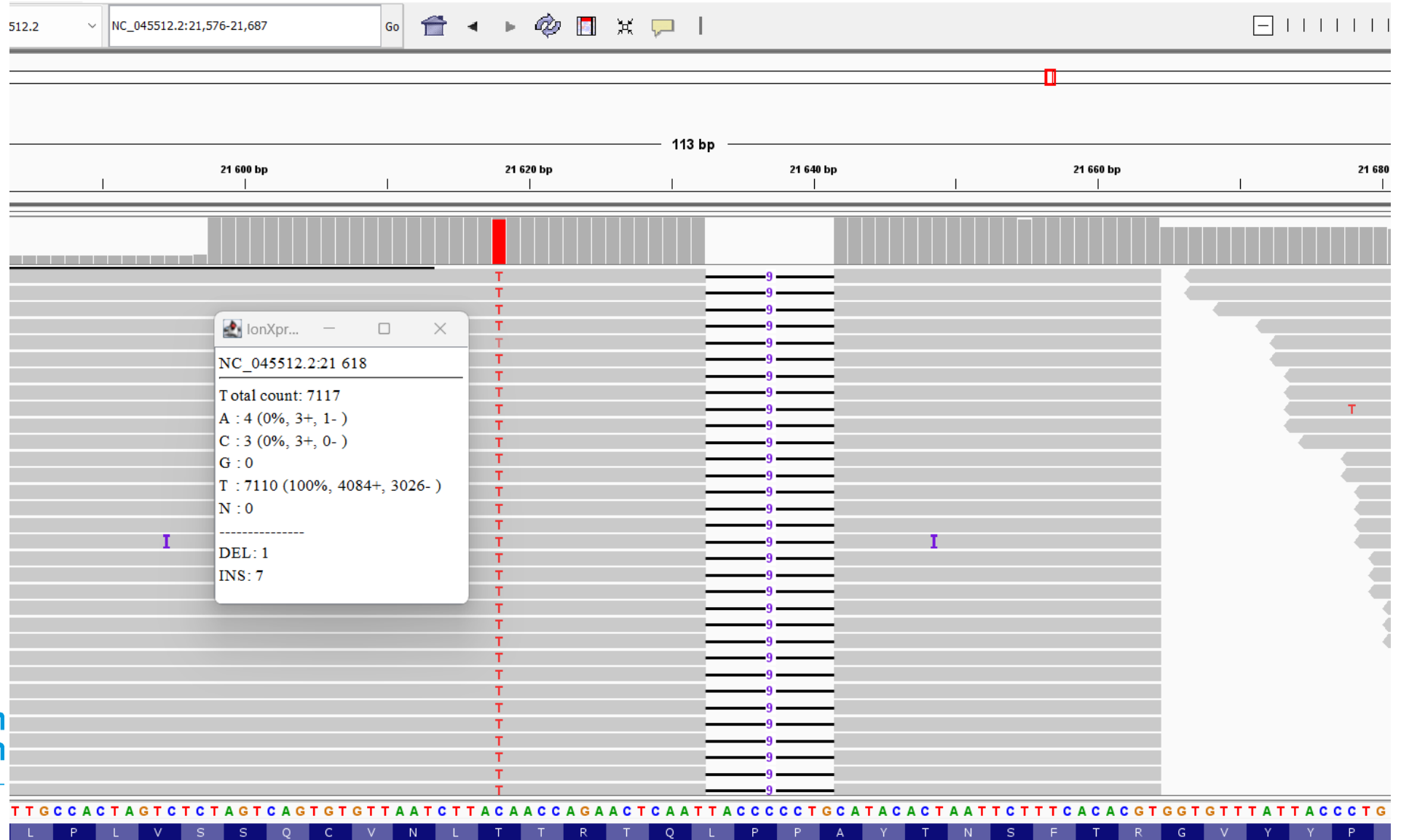
Sequence → AGAGGTTGATGAAGTCAAGACAAATCGCTCCAGGGCAAACCTGGAAAGATTGCTGATTATAATTATAAATTACAGATGATTT

Annotations R G D E V R Q I A P G Q T G K I A D Y N Y K L P D D F



IGV

(Integrative Genomics Viewer) – Інтегрований геномний переглядач



Формат FASTA

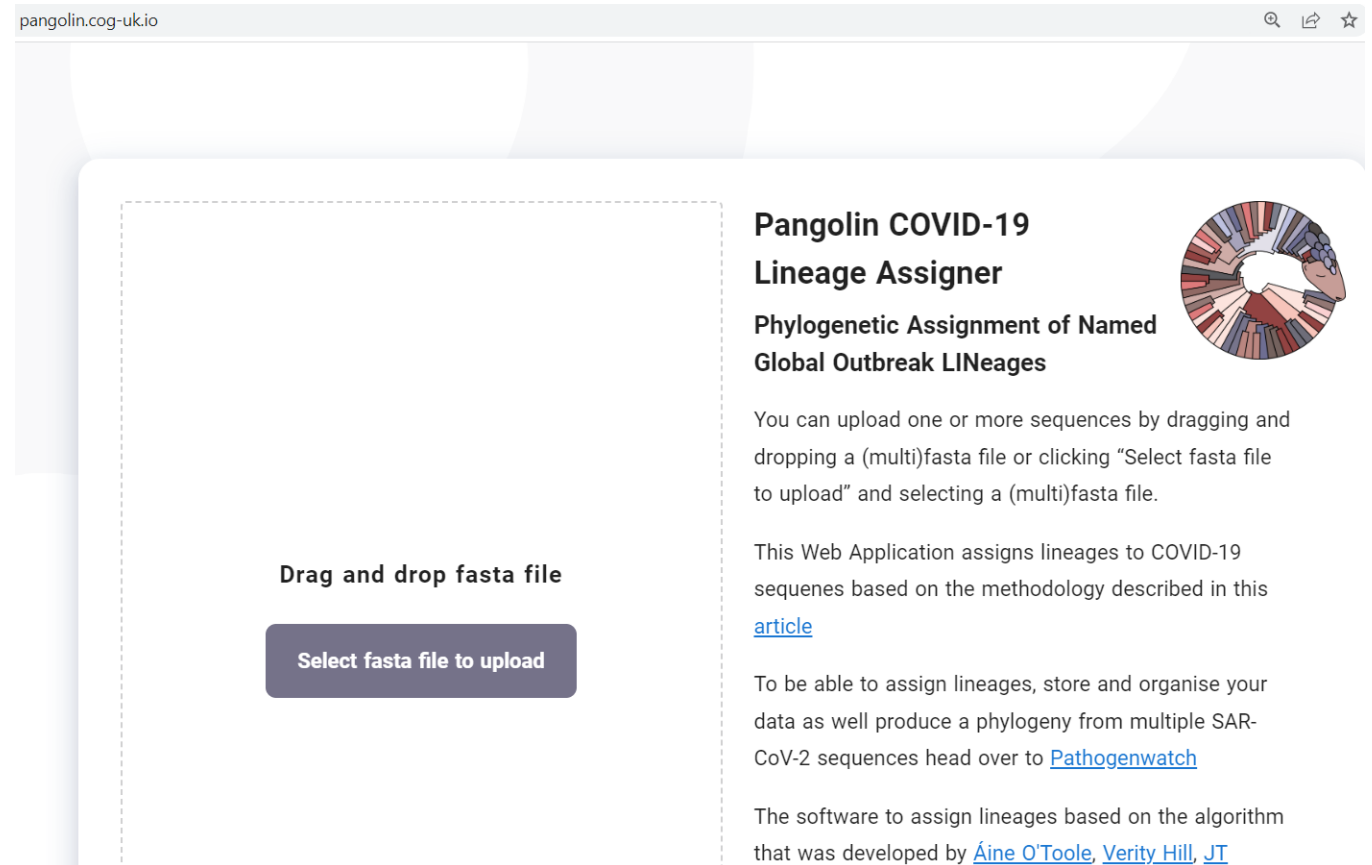
- простий текстовий формат для файлів послідовностей ДНК і білків.
- Кожній послідовності передуює рядок заголовка, який починається з «>» і за яким йде назва.
- Будь-який інший текст після першого рядка вважається частиною послідовності.
- Один файл може містити багато послідовностей

```
>hCoV-19/Ukraine/74373/2023
TTGATCTCTTGTAGATCTGTTCTCTAAACGAACCTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTA
TAATAACTAATTACTGTTCGTTGACAGGACACGAGTAACTCGTCTATCTTCTGCAGGCTGCTTACGGTTTCG
CAGCANNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNAGGTAAGATGGAGAGCCTTGTCCCTGGTTTCAACGAG
TTTGCCTGTTTTACAGGTTTCGCGACGTGCTCGTACGTGGCTTTGGAGACTCCGTGGAGGAGGTCTTATCAG
GATGGCACTTGTGGCTTAGTAGAAGTTGAAAAAGGCGTTTTGCCTCAACTTGAACAGCCCTATGTGTTTCAT
CTGCACCTCATGGTCATGTTATGGTTGAGCTGGTAGCAGAACTCGAAGGCATTTCAGTACGGTCGTAGTGGT
CCCTCATGTGGGCGAAATACCAGTGGCTTACCGCAAGGTTCTTCTTCGTAAGAACGGTAATAAAGGAGCTG
GATCTAAAGTCATTTGACTTAGGCGACGAGCTTGGCACTGATCCTTATGAAGATTTTCAAGAAAAGTGGAA
TTACCCGTGAACTCATGCGTGAGCTTAACGGAGGGGCATACACTCGCTATGTTCGATAACAACCTTCTGTGGC
GTGCATTAAGACCTTCTAGCACGTGCTGGTAAAGCTTCATGCACTTTGTCCGAACAACCTGGACTTTATTG
TGCTGCCGTGAACATGAGCATGAAATTGCTTGGTACACGGAACGTTCTGAAAAGAGCTATGAATTGCAGAC
CAAAGAAATTTGACACCTTCAATGGGGAATGTCCAAATTTTGTATTTCCCTTAAATTCATAATCAAGACT
GAAAAAGCTTGATGGCTTTATGGGTAGAATTCGATCTGTCTATCCAGTTGCGTCAACCAATGAATGCAACC
ATGAAGTGTGATCATTGTGGTGAACTTCATGGCAGACGGGCGATTTTGTAAAGCCACTTGCGAATTTTG
AAGAAGGTGCCACTACTTGTGGTTACTTACCCCAAATGCTGTTGTTAAAATTTATTGTCCAGCATGTCAC
GCATAGTCTTGCCGAATACCATAATGAATCTGGCTTGAAAACCATTCTTCGTAAGGGTGGTCGCACTATTG
TCTTATGTTGGTTGCCATAACAAGTGTGCCTATTGGGTTCCACGTGCTAGCGCTAACATAGGTTGTAACCA
GTTCCGAAGGTCTTAATGACAACCTTCTTGAAATACTCCAAAAAGAGAAAGTCAACATCAATATTGTTGGT
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNGCTTCCACAAGTGCTTTTGTGGAACTGTGAAAGGTTTGGATT
GTTGAATCCTGTGGTAATTTTAAAGTTACAAAAGGAAAAGCTAAAAAGGTGCCTGGAATATTGGTGAACA
TTTATGCATTTGCATCAGAGGCTGCTCGTGTGTACGATCAATTTTCTCCCGCACTCTTGAAACTGCTCAA
GAAGGCCGCTATAACAATACTAGATGGAATTTACAGTATTCCTGAGACTCATTGATGCTATGATGTTCA
```

Pangolin – визначення варіантів SARS-CoV-2

Варіанти ДНК ≠ варіанти SARS-CoV-2

- **Варіант ДНК/РНК (мутація)** - це відмінність в генетичній послідовності в порівнянні з референсною послідовністю.
- **Варіант SARS-CoV-2** - це вірусний геном, який відрізняється від інших вірусів однією або кількома мутаціями. Група вірусів з подібною послідовністю та загальними характеристиками (вірулентність, інфекційність тощо)




























The screenshot shows the Pangolin COVID-19 Lineage Assigner web application. The browser address bar displays 'pangolin.cog-uk.io'. The main content area features a large dashed box for file upload with the text 'Drag and drop fasta file' and a button labeled 'Select fasta file to upload'. To the right, the title 'Pangolin COVID-19 Lineage Assigner' is displayed above a circular phylogenetic tree icon. Below the title is the subtitle 'Phylogenetic Assignment of Named Global Outbreak LINEages'. The text explains that users can upload sequences by dragging and dropping a (multi)fasta file or clicking 'Select fasta file to upload' and selecting a (multi)fasta file. It also states that the application assigns lineages to COVID-19 sequences based on the methodology described in an [article](#). For more information, users are directed to [Pathogenwatch](#). The software is attributed to [Áine O'Toole](#), [Verity Hill](#), and [JT](#).

Pangolin – визначення варіантів SARS-CoV-2

← → ↻ 🏠 pangolin.cog-uk.io

 [Reset entries](#) [Upload another file](#)

File name	Sequence name	Lineage
— ANALYSED (Click tick icon for more info) 77 sequences 		
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73496/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73497/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73498/2023	BF.14   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73499/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73504/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73505/2023	BA.5.2   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73506/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73509/2023	BA.2   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73510/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73511/2023	BQ.1.10   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73539/2023	BA.2   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73558/2023	BA.5.2   
✓ 21_02_2023_torrent-server_61.fasta	hCoV-19/Ukraine/73559/2023	BA.5.2   

Бази даних

1. Послідовність геному

2. **Метадані** - набір даних, який дає інформацію про первинні дані
- додаткові дані про зразок, з якого була отримана послідовність ДНК/РНК:

- Дата забору
- Географічні дані
- Організм
- Вид біоматеріалу
- Стать і вік пацієнта
- Статус вакцинації пацієнта
- інші дані

! Метадані повинні бути стандартизовані між усіма лабораторіями



База даних GISAID

- це глобальна науково-дослідна ініціатива та основне джерело геномних даних і метаданих усіх вірусів грипу, респіраторно-синцитіального вірусу (RSV) і SARS-CoV-2

© 2008 - 2023 | Terms of Use | Privacy Notice | Co

You are logged in as **Anna Iaruchyk** - [log out](#)

Registered Users | EpiFlu™ | **EpiCoV™** | EpiRSV™ | EpiPox™ | My Profile

EpiCoV™ | Search | Downloads | Upload

Pandemic coronavirus causing COVID-19

A previously unknown human coronavirus (hCoV-19) was first detected in late 2019 in patients in the City of Wuhan, who suffered from respiratory illnesses including atypical pneumonia, an illness that has become known as coronavirus disease (COVID-19). The coronavirus originated from an animal host and is closely related to the virus responsible for the Severe Acute Respiratory Syndrome coronavirus (SARS).

On 10. January 2020, the first virus genomes and associated data were publicly shared via GISAID. The World Health Organization announced on 11. March 2020 the first coronavirus pandemic. As the pandemic progresses, scientists from around the globe are tracking the virus and its genome sequences to ensure optimal virus diagnostic tests, to track and trace the ongoing outbreak and to identify potential intervention options. Several analyses to assist with these efforts are offered here, including sequence alignments, diagnostic primer and probe coordinates, 3D protein models, drug targets, phylogenetic trees and many more.

[Search](#)

Audacity

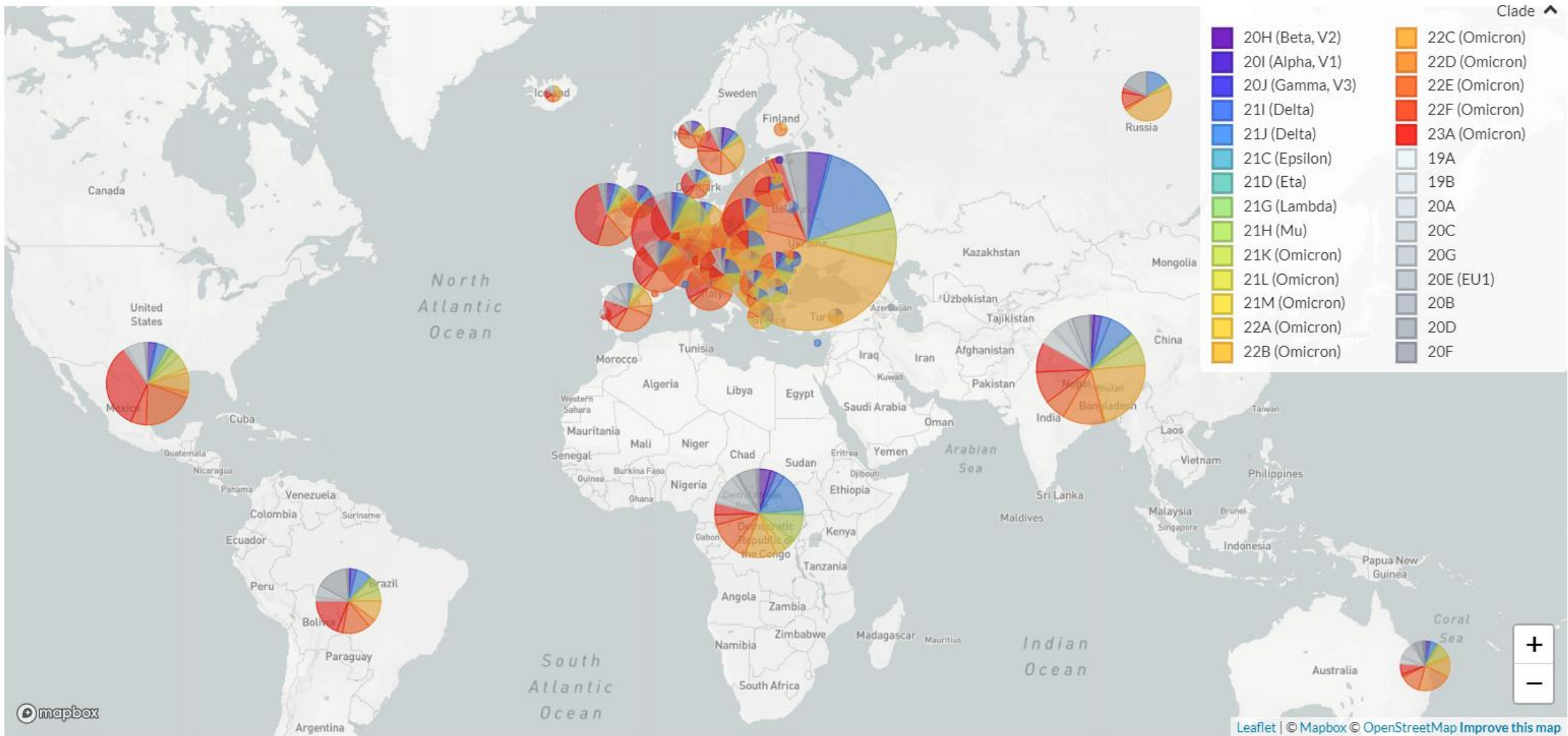
AudacityInstant

BLAST

CoVizu^e

Emerging Variants

Lineage Frequency



World Health Organization

European Region

<https://next.nextstrain.org>

Значення NGS та геномного нагляду в охороні здоров'я

- Дослідження походження інфекційних захворювань
- Виявлення спалахів/епідемій/пандемій інфекційних захворювань та управління ними
- Дослідження резистентності мікроорганізмів
- Розуміння джерел/резервуарів
- Виявлення нових патогенів
- Підтримка розробки діагностики, ліків і вакцин, а також інформування про реакцію на спалах



Дякую за увагу!

Анна Яручик

Експерт Бюро ВООЗ в Україні

iaruchyka@who.int



European Region

