# Learning Outcomes - Day 1

- Learning to identify and download genomic data from NCBI.

– Understand the differences between SRA, SRX, PRJ, SAM and SRR.

– Identify the different IDs in NCBI website.
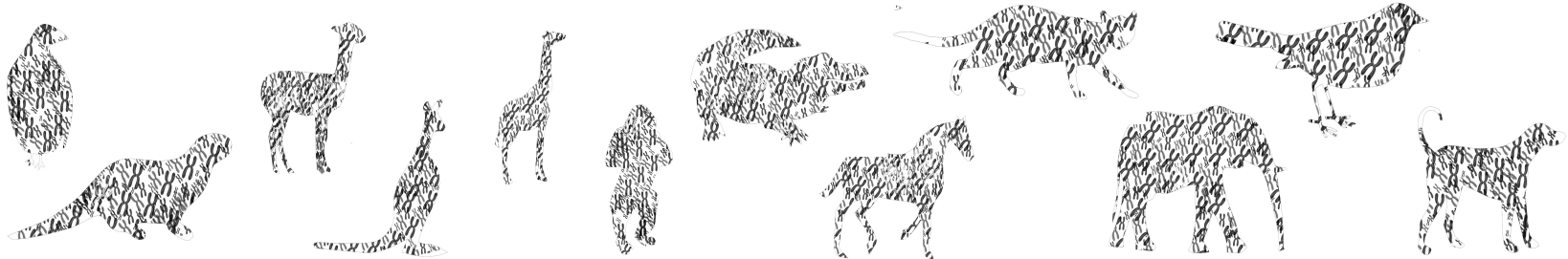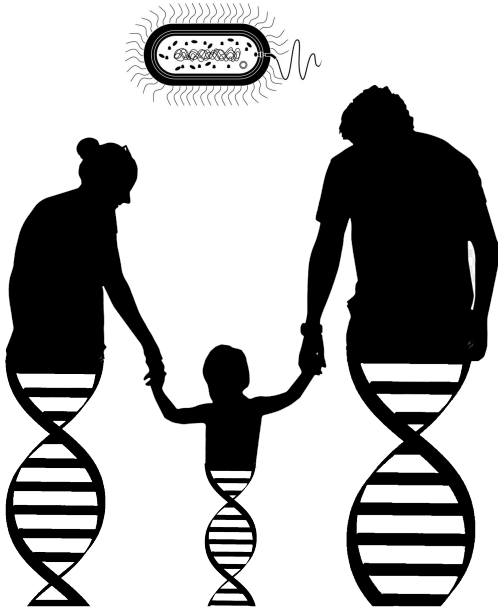
# Gene sequencing

The Mendelian gene is a basic unit of heredity.

The molecular gene is a sequence of nucleotides in DNA, that is transcribed to produce a functional RNA and protein.
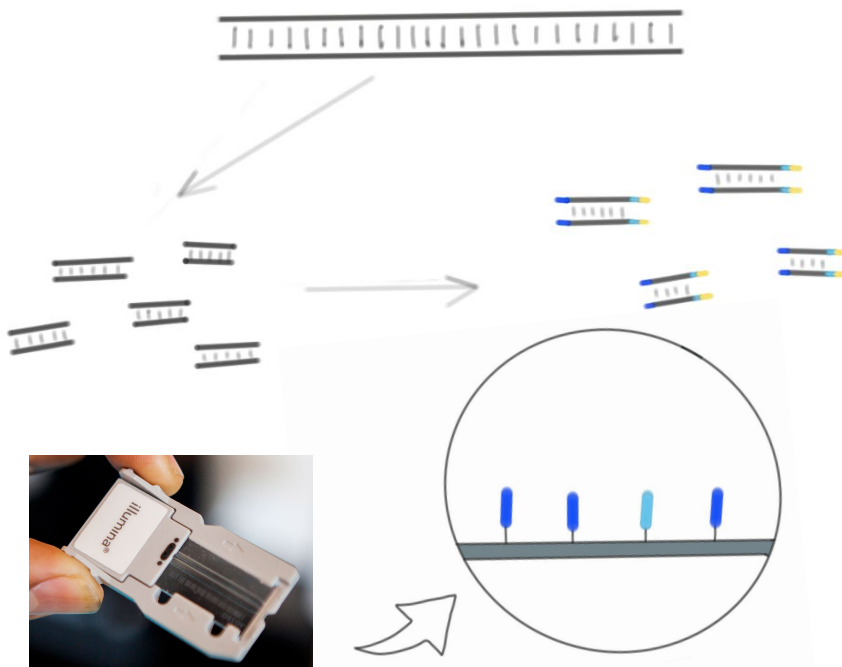
Gene detection in humans: Disease detection (i.e. cancer)

Gene detection in other animals: Disease detection, breeding.

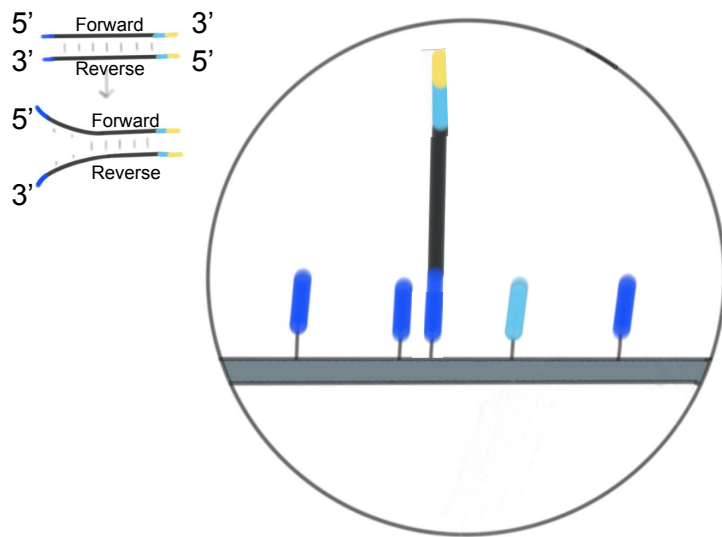Gene detection in pathogens → AMR, virulence, other specific genes

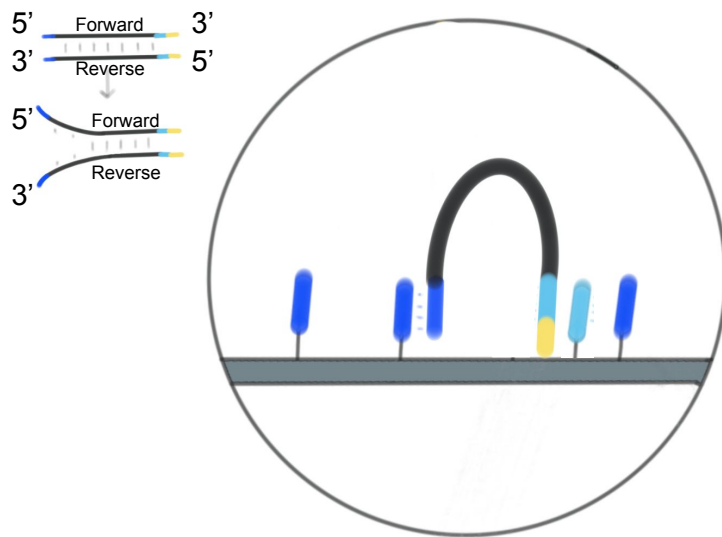DNA can be cut in small pieces and sequenced based on a **reference genome**.



- **Library preparation**: Long strands of DNA are cut in small fragments, and specific sequences of DNA called adapters are added at each end of each fragment, including an index to identify the sample.

- **Sequencing by synthesis (Illumina)**: A glass surface called flow cell contains small fragments of DNA called oligonucleotides. These match the adapters in the library.

DNA can be cut in small pieces and sequenced based on a **reference genome**.
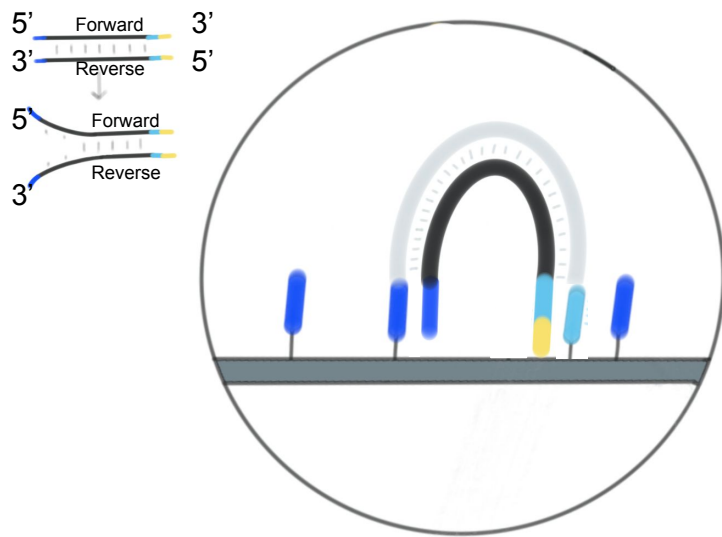


- One of your sequences (forward sense) attaches to one of the oligos by the adapter and it gets copied. This copy is the continuation of the oligo but it is the opposite sense (reverse). The forward sequences remaining are washed away.
- **Clonal amplification via PCR**: The reverse sequence bends and attaches to the next oligo, and it is copied. The result is the original sense of the sequence (forward sense). The reverse sense sequences are discarded.

5

# Next Generation Sequencing (NGS)

DNA can be cut in small pieces and sequenced based on a **reference genome**.
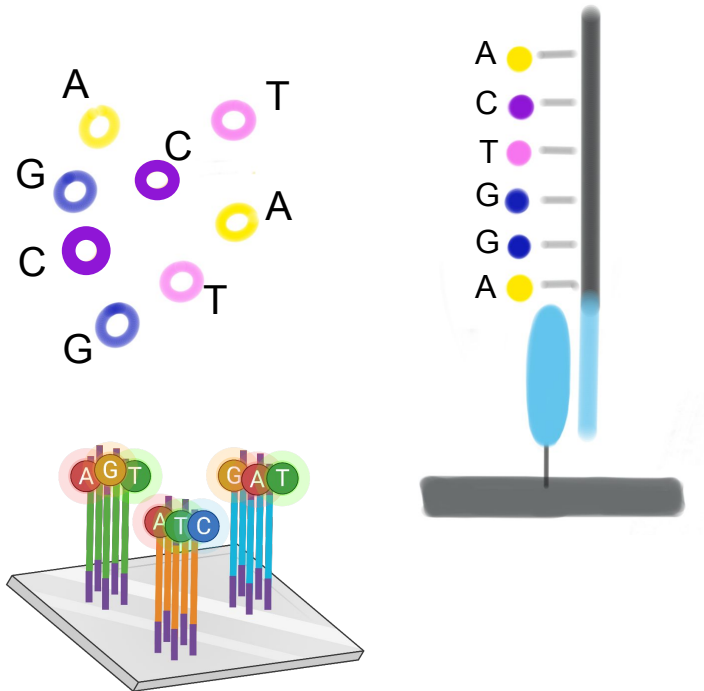


- One of your sequences (forward sense) attaches to one of the oligos by the adapter and it gets copied. This copy is the continuation of the oligo but it is the opposite sense (reverse). The forward sequences remaining are washed away.
- **Clonal amplification via PCR**: The reverse sequence bends and attaches to the next oligo, and it is copied. The result is the original sense of the sequence (forward sense). The reverse sense sequences are discarded.

6

NC STATE UNIVERSITY

DNA can be cut in small pieces and sequenced based on a **reference genome**.
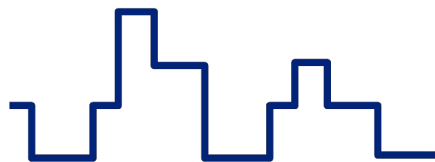


- One of your sequences (forward sense) attaches to one of the oligos by the adapter and it gets copied. This copy is the continuation of the oligo but it is the opposite sense (reverse). The forward sequences remaining are washed away.
- **Clonal amplification via PCR**: The reverse sequence bends and attaches to the next oligo, and it is copied. The result is the original sense of the sequence (forward sense). The reverse sense sequences are discarded.

7

NC STATE UNIVERSITY

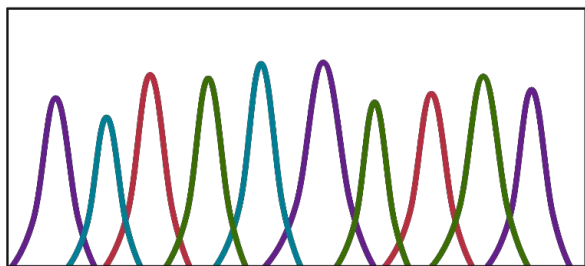DNA can be cut in small pieces and sequenced based on a **reference genome**.



- Sequencing primers bind to the forward strand. In the sequencing medium, there are fluorescent nucleotides with one specific fluorescent tag, polymerase enzymes and a terminator.
- Only one nucleotide per cycle is attached to the sequence, and its color is detected and stored by the instrument. Finally, the index gets sequenced. **This is called single-end.**
- If a second index is sequenced, along with the reverse strand of the library, it is known as **paired-end sequencing**.

8

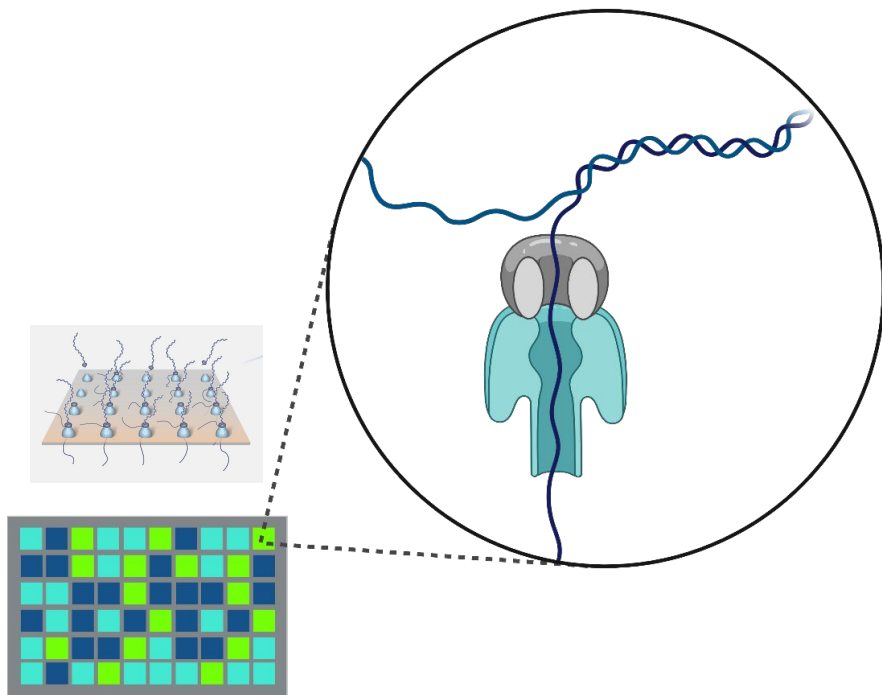DNA can be cut in small pieces and sequenced based on a **reference genome**.
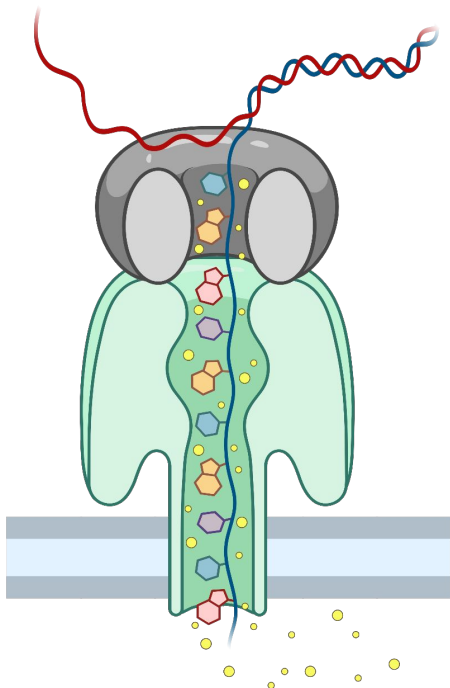
AGTCCCTGAATCGA

- After several filters within the instrument, the resulting sequences are **demultiplexed**, which means sort reads by sample based on the attached indexes.
- After demultiplexing, the sequences are **mapped to a reference genome.**

- This type of sequencing is **short read sequencing,** with fragments between 75-300bp

NC STATE UNIVERSITY

Long-read sequencing. Oxford Nanopore Technologies (ONT)



- Similar than with Illumina, libraries are prepared where an adapter is attached to each sequence.
- ONT sequencer flow cells are made of an electric-resistant bilayer created by a synthetic polymer with an array of nanopores (~1.8nm).
- A potential is applied across the membrane, resulting in a current flowing only through the nanopores.

# Nanopore sequencing

Long-read sequencing. Oxford Nanopore Technologies (ONT)



- Single molecules flowing through the nanopores cause characteristic disruptions in the electric current across the membrane.
- Measuring those disruptions, the molecules can be easily identified.
- ONT analyzes intact DNA strands that pass through the nanopores and analyzed in real time.
- By preparing the DNA to have a hairpin structure, ONT can produce pair-end data in one continuous read.
- Long-read sequencing still has higher error rate than short-read sequencing.

FASTQ format is a text-based format for storing both a biological sequence (usually nucleotide sequence) and its corresponding quality scores.



FASTQ File Format Analysis

# Genomic data and where to find it

# SRAs and NCBI Entrez

Sequence Read Archive (**SRA**) data, available through multiple cloud providers and NCBI servers, is the largest publicly available repository of high throughput sequencing data.

Entrez is a molecular biology database system that provides integrated access to nucleotide and protein sequence data, gene-centered and genomic mapping information, 3D structure data, PubMed MEDLINE, and more.

Entrez covers over 20 databases including the complete protein sequence data from PIR-International, PRF, Swiss-Prot, and PDB and nucleotide sequence data from GenBank that includes information from EMBL and DDBJ.

**NCBI**

**Entrez Molecular Sequence Database System**

| PubMed | Entrez | BLAST | OMIM | Taxonomy | Structure |

**Sequence Read Archive (SRA):** repository that stores raw sequencing data. It includes the actual sequences obtained from the high-throughput sequencing machines before any processing or analysis.

**PRJ:** Refers to a genomic **project**. It represents a collection of related experiments and data associated with a particular scientific investigation.

**SRX**: Refers to a **sequencing experiment.** It is an organized set of runs, which are the actual sequencing data generated for a sample or set of samples. SRX[...], ERX[...], DRX[...].

**SAM**: Represents an individual **biological sample** that is part of a project. It is a specific instance of genetic material taken from an organism.

**SRR:** Represents a **specific run** of sequencing data. It is the raw data generated in a single sequencing run for a particular sample.

# NCBI Data Acronyms

# NCBI Data Acronyms



PRJ 1
(Project → Shelf)

PRJ 2

PRJ 3

SRA
(Sequence Read Archive → Library)

SRX
(Experiment → book)

SAM
(Sample → Chapter)

SAM 1 SAM 2

SRR
(Run → Paragraph)

SRR 1

SRR 2

# NCBI Data Acronyms



PRJ 1
(Project → Shelf)

PRJ 2

PRJ 3

SRA
(Sequence Read Archive → Library)

SRX
(Experiment → book)

SAM
(Sample → Chapter)

SAM 1          SAM 2

SRR          SRR 1
(Run → Paragraph)

SRR 2

Multiple experiments
(SRX) using the
same sample (SAM)

- Sequence details are important (i.e. single vs paired).

- Difference between long and short-read sequencing.

  - FASTq files contain important information.

- Numerous public databases for genomic data.

  - SRA → PRJ → SRX → SAM → SRR.