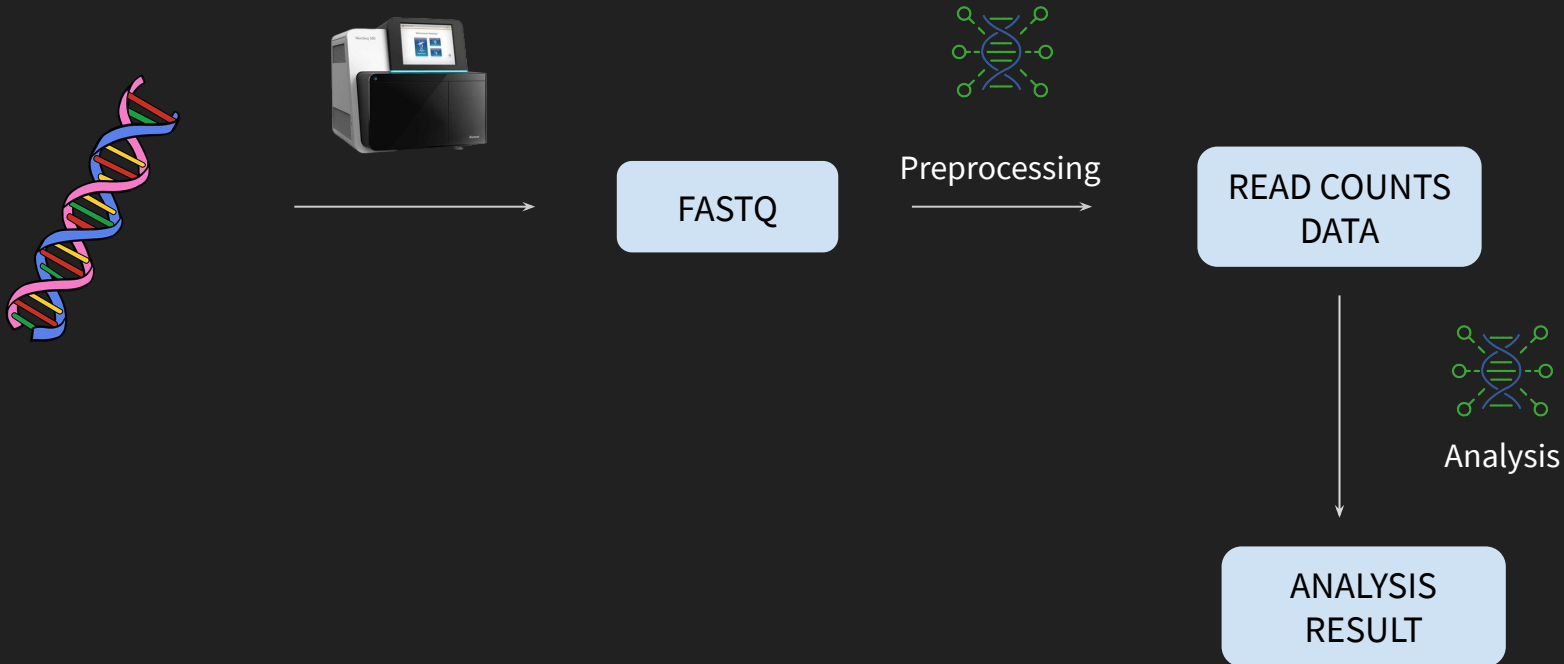# Introduction to Bioinformatics
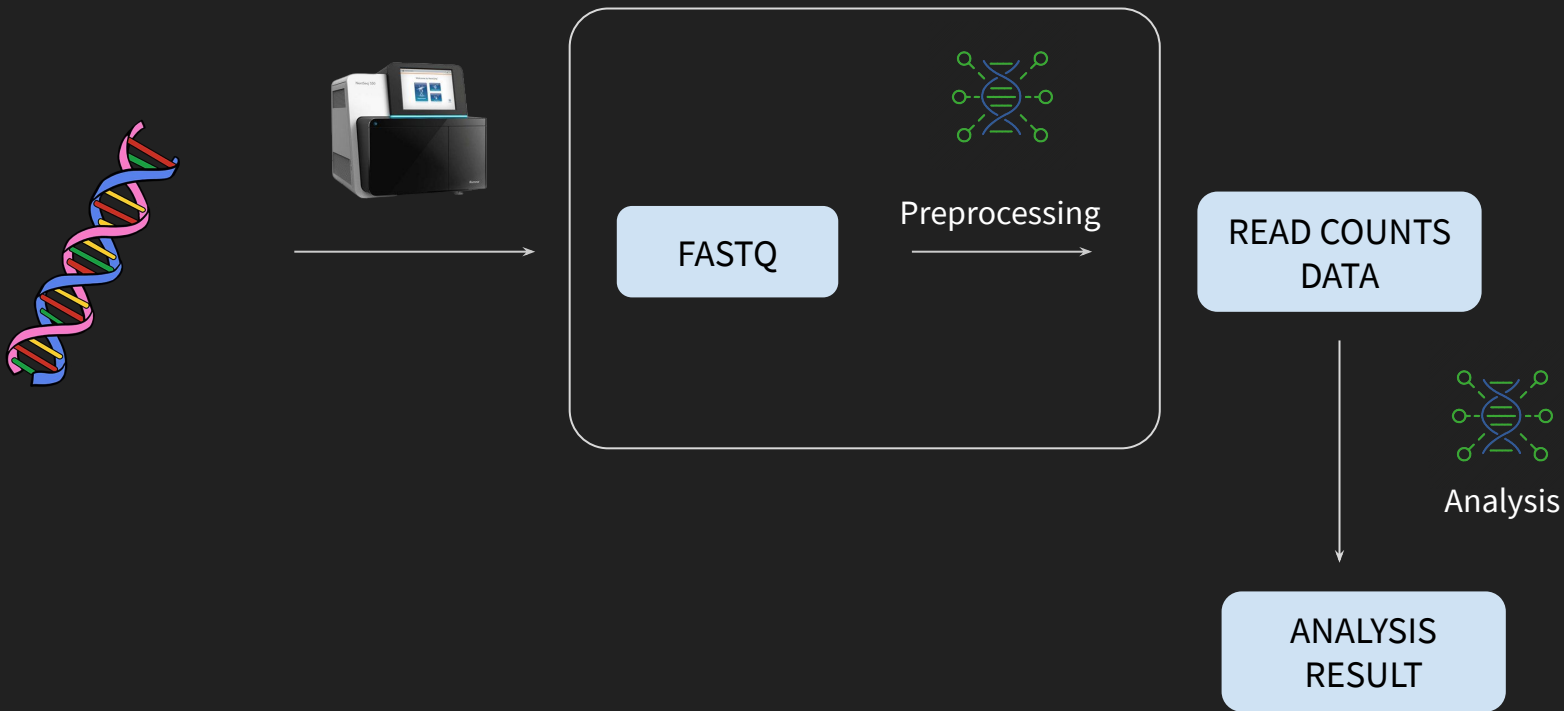
More on UNIX and Bioinformatics Data Formats

# A Quick Glimpse
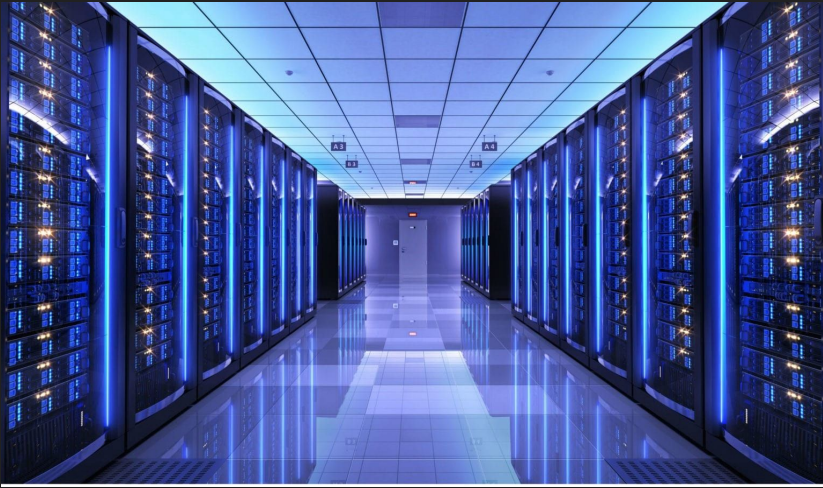
NGS (Bisulfite Sequencing)

# A Quick Glimpse

NGS (Bisulfite Sequencing)

# Why do we need to learn UNIX?



- Most of the software in the preprocessing part for NGS analysis is only available and can be run in UNIX system only
- Most high performance computing is using UNIX

# What are we doing and what are we processing?

- During the preprocessing part we will preprocess the output from the sequencing machine
- The output is mostly looking like DNA sequence since we're sequencing DNA
- Our job here is to convert this data so it'd be usable for analysis

# Preprocessing Pipeline using MethylSeq

## Pipeline Summary

The pipeline allows you to choose between running either Bismark or bwa-meth / MethylDackel. Choose between workflows by using `--aligner bismark` (default, uses bowtie2 for alignment), `--aligner bismark_hisat` or `--aligner bwameth` .

| Step | Bismark workflow | bwa-meth workflow |
| --- | --- | --- |
| Generate Reference Genome Index *(optional)* | Bismark | bwa-meth |
| Raw data QC | FastQC | FastQC |
| Adapter sequence trimming | Trim Galore! | Trim Galore! |
| Align Reads | Bismark | bwa-meth |
| Deduplicate Alignments | Bismark | Picard MarkDuplicates |
| Extract methylation calls | Bismark | MethylDackel |
| Sample report | Bismark | - |
| Summary Report | Bismark | - |
| Alignment QC | Qualimap | Qualimap |
| Sample complexity | Preseq | Preseq |
| Project Report | MultiQC | MultiQC |

# File Types

- Plain text file formats
    - Information often structured into lines and columns
    - Human-readable
    - Easy to process


- Binary file formats
    - Not human-readable
    - Require special software for processing
    - Efficient storage
    - (significant) reduction to file size when compared to a plain text counterpart (e.g. 75 % space saved)

# Common File Formats that You Will Encounter

- **FASTA** - Simple collections of named DNA/protein sequences (text)
- **FASTQ** - Extension of FASTA format, contains additional quality information. Widely used for storing unaligned sequencing reads (text)
- **SAM/BAM** - Alignments of sequencing reads to a reference genome (text/binary)
- **BED** - Region-based genome annotation information (e.g. a list of genes and their genomic locations).
- **GFF/GTF** - gene-centric annotations (text)
- **VCF** - variant call format, to store information about genomic variants (text)
- **CSV/TSV** - Usually stores read counts/expression information per sample

# FASTA format

The nucleic acid codes that can be found in FASTA file:

A --> adenosine

T --> thymidine

C --> cytidine

S --> G C (strong)

G --> guanine

W --> A T (weak)

B --> G T C

U --> uridine

N --> A G C T

R --> G A (purine)

Y --> T C (pyrimidine)

Example of fasta format http://hgdownload.cse.ucsc.edu/goldenPath/hg38/chromosomes/

# Quick UNIX Check

- **How long is chrY?**
  - $ grep -v ">" hg38.chrY.fa | grep -o "[ATCGatcg]" | wc -l 26415043
- **How many adenosines are there? $**
  - $ grep -v ">" hg38.chrY.fa | grep -o -i "A" | wc -l 7886192

# FASTQ format

- Nearly all sequencing technologies produce sequencing reads in FASTQ format
  - **Sequence ID** @SEQ_ID
  - **Sequence** GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACA GTTT
  - **Separator** +
  - **Quality scores** !''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>

# FASTQ Quality Scores (Phred Scores)

- PHRED Base quality (Q) – integer value derived from the estimated probability (P) of the corresponding base being determined wrong
  - $Q = -10 * \log 10(Perr)$ (rounded to nearest integer)

- PHRED Base quality (Q) – integer value derived from the estimated probability (P) of the corresponding base being determined wrong A higher quality score is better (>=20 is considered "good")
  - Score of 10 means 10% of probability of it's being error
  - Score of 20 means 1%
  - Score of 30 means 0.1% etc

# FastQC Helps Quality Control



Quality scores across all bases (Illumina >v1.3 encoding)

**More information on interpreting:**

https://hbctraining.github.io/Intro-to-rnaseq-hpc-salmon/lessons/qc_fastqc_assessment.html

# Sequence Alignment Map (SAM)

- Intended for storing read alignments against reference sequences
- Has a binary version with good software support (BAM format)

- The SAM format consists of two sections:
  - Header section Used to describe source of data, reference sequence, method of alignment, etc.
  - Alignment section Used to describe the read, quality of the read, and nature alignment of the read to a region of the genome

# Sequence Alignment Map (SAM)



Example SAM/BAM header section (abbreviated)



Example SAM/BAM alignment section (only 10 alignments shown)

# SAM/BAM Header Section

- Used to describe source of data, reference sequence, method of alignment, etc.
- Each section begins with '@' followed by a two-letter record type code. These are followed by two-letter tags and values, example:
    - @HD The header line
    - VN: format version
    - SN: reference sequence name
    - LN: reference sequence length
    - SP: species

# SAM/BAM Alignment Section

| Col | Field | Type | Regexp/Range | Brief description |
|---|---|---|---|---|
| 1 | QNAME | String | [!-?A-~]{1,254} | Query template NAME |
| 2 | FLAG | Int | $[0, 2^{16}-1]$ | bitwise FLAG |
| 3 | RNAME | String | \*\|[!-()+-<>-~][!-~]* | Reference sequence NAME |
| 4 | POS | Int | $[0, 2^{31}-1]$ | 1-based leftmost mapping POSition |
| 5 | MAPQ | Int | $[0, 2^{8}-1]$ | MAPping Quality |
| 6 | CIGAR | String | \*\|([0-9]+[MIDNSHPX=])+ | CIGAR string |
| 7 | RNEXT | String | \*\|=\|[!-()+-<>-~][!-~]* | Ref. name of the mate/next read |
| 8 | PNEXT | Int | $[0, 2^{31}-1]$ | Position of the mate/next read |
| 9 | TLEN | Int | $[-2^{31}+1, 2^{31}-1]$ | observed Template LENgth |
| 10 | SEQ | String | \*\|[A-Za-z=.]+ | segment SEQuence |
| 11 | QUAL | String | [!-~]+ | ASCII of Phred-scaled base QUALity- |

# Tools to work with BAM/SAM

- **samtools** - view, sort, index, QC, stats on SAM/BAM files, and more
- **sambamba** - view, sort, index, merge, stats, mark duplicates. fast laternative to samtools
- **picard** - QC, validation, duplicates removal and many more utility tools

# BED File Formats

- Text-based, tab-separated list of genomic regions
- Each region is specified by a reference sequence and the start and end positions on it
- Optionally, each region can have additional properties defined – E.g. strand, name, score, color
- Intended for visualizing genomic annotations in IGV, UCSC Genome Browser (context of expression, regulation, variation, conservation, . . . )

# BED File Formats

- 3 mandatory columns (must be in correct order)
  - "chrom" – chromosome
  - "chromStart" – the first base of the region with respect to the chromosome (counting starts from 0)
  - "chromEnd" – the first base after the region with respect to the chromosome [chromStart, chromEnd) allows easy region-length calculation
  - Optional fields: "name", "score", "strand", other annotation columns

# Example of BED File Formats

chr1 115263684 115263685 rs10489525 0 +
chr12 97434219 97434220 rs6538761 0 +
chr14 102360744 102360745 rs7142002 0 +
 chr16 84213683 84213684 rs4150167 0 -
chr2 206086170 206086171 rs4675502 0 +
chr20 14747470 14747471 rs4141463 0 +

# BED File Formats

- 9 additional optional fields, their order is binding (unlike with SAM format).
- All regions must have the same optional fields
- Most important optional fields:
  - "name" – name of the region
  - "score" – score value between 0 and 1000 (read-count, transformed p-value, "quality", . . . ) Can be interpreted as shades of grey during visualization
  - "strand" – either "+" or "-" (not "1"/"-1") BED12 format specification available

# Tools to work with BED File Formats

- **bedtools** - universal tools for manipulating genomic regions
- **bedops** - complementary to bedtools, providing additional functionality and speedup

# Genomic Data Resources

- **GEO**: Gene Expression Omnibus.
  - Host array- and sequencing-based data.
- **ArrayExpress**: European version of GEO.
  - Better curated than GEO but has less data.
- **SRA**: Sequence Read Archive. Designed for hosting large scale high-throughput sequencing data, e.g., high speed file transfer. Data are required to be deposited in one of the databases when paper is accepted

# Sequence Read Archive

- The NCBI database which stores sequence data obtained from next generation sequence (NGS) technology
- Archives raw NGS data for various organisms from several platforms (FASTQ files) Serves as a starting point for "secondary analyses"
- Provides access to data from human clinical samples to authorized users who agree to the datasets' privacy and usage mandates
- Search metadata to locate the sequence reads for download and further downstream analyses

# Getting data from SRA

- The NCBI sratoolkit provides two command line tools to allow local BLAST searches against specific sra files directly
  - fastq-dump: Convert SRA data into fastq format
  - prefetch: Allows command-line downloading of SRA, dbGaP, and ADSP data
  - sam-dump: Convert SRA data to sam format
- .sra files are NOT FASTQ files - need to further convert them using sratoolkit