# Bash, File formats and Quality Control

**Lecture: Introduction to Linux systems, command line and file formats**
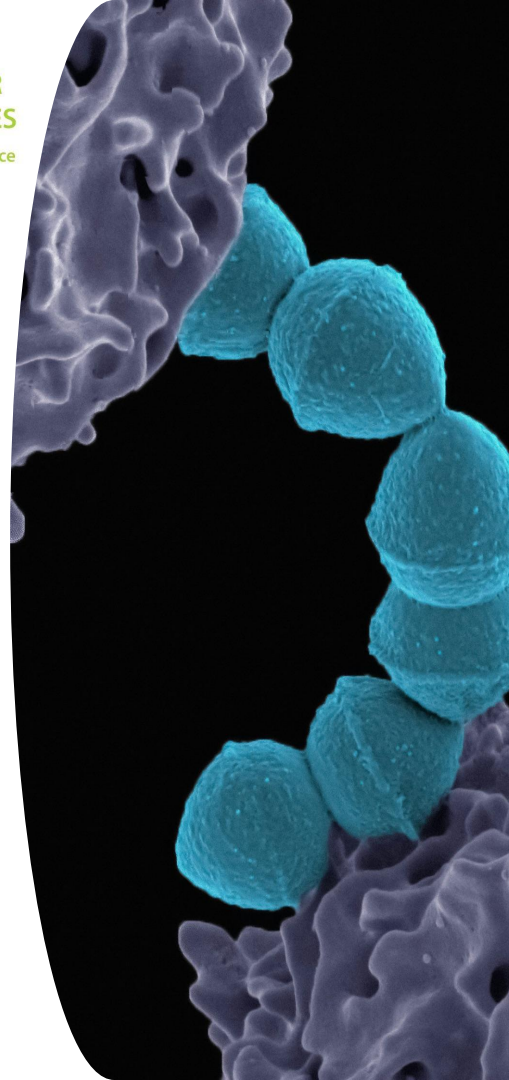
**Instructors:**

Rito Mikhari

Collins Kigen

AMR Virtual course (Africa and Asia)

**21 May 2024**

## Learning outcome

- Learn how to use the Unix command

- Understand the different file formats

## Overview
- Background: Operating systems, Unix

- Basic Unix commands

- Introduction to colab notebooks (demo)

- Introduction to the basic file formats and their functions in genomics

# Background: Operating systems

- An operating system (OS) is the software that functions as an interface between the computer user and its hardware

- The software enables and controls the set of programs that we install in our machines

- Some functions include: allocation of storage, processor management, security, error detection, file management

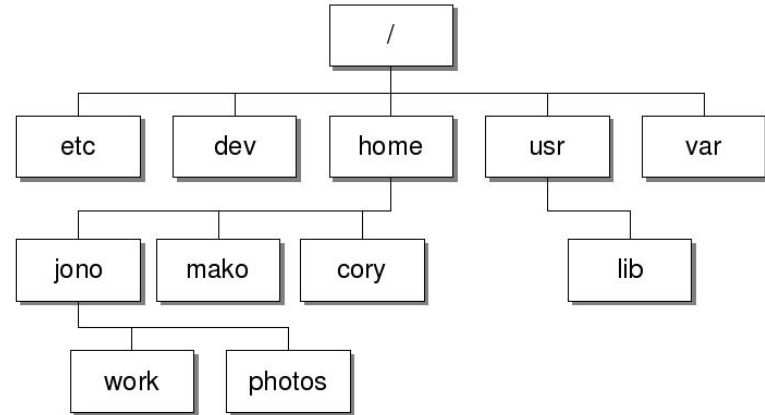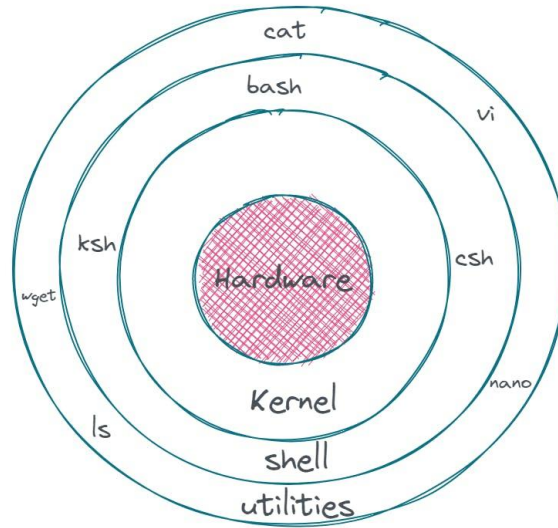- Examples of most used OS include: Windows, Linux, Unix, MAC OS, Android, etc.

# Background: Unix

- Unix was originally called AT&T Unix, and was first developed in 1969 at the Bell Lab research center

- A family of OS's that allow for multiple users and allow for multitasking

- Beneficial for working with large text files (or genomic data files) and accommodates several powerful yet flexible commands

- Ability to combine these multiple commands for different objectives

4

# Basic Unix: Basic terms

- Shell
  Command-line interpreter

- Terminal
  Tool used for shell commands

- Command prompt
  (base) ubuntu@student-1:~$

- Directories
  - Equivalent to folders with files

- Root
  - The mother directory where all directories stem from

- Syntax
  - Command [options] [arguments]
  - Different programs will have their own syntax with a number of options
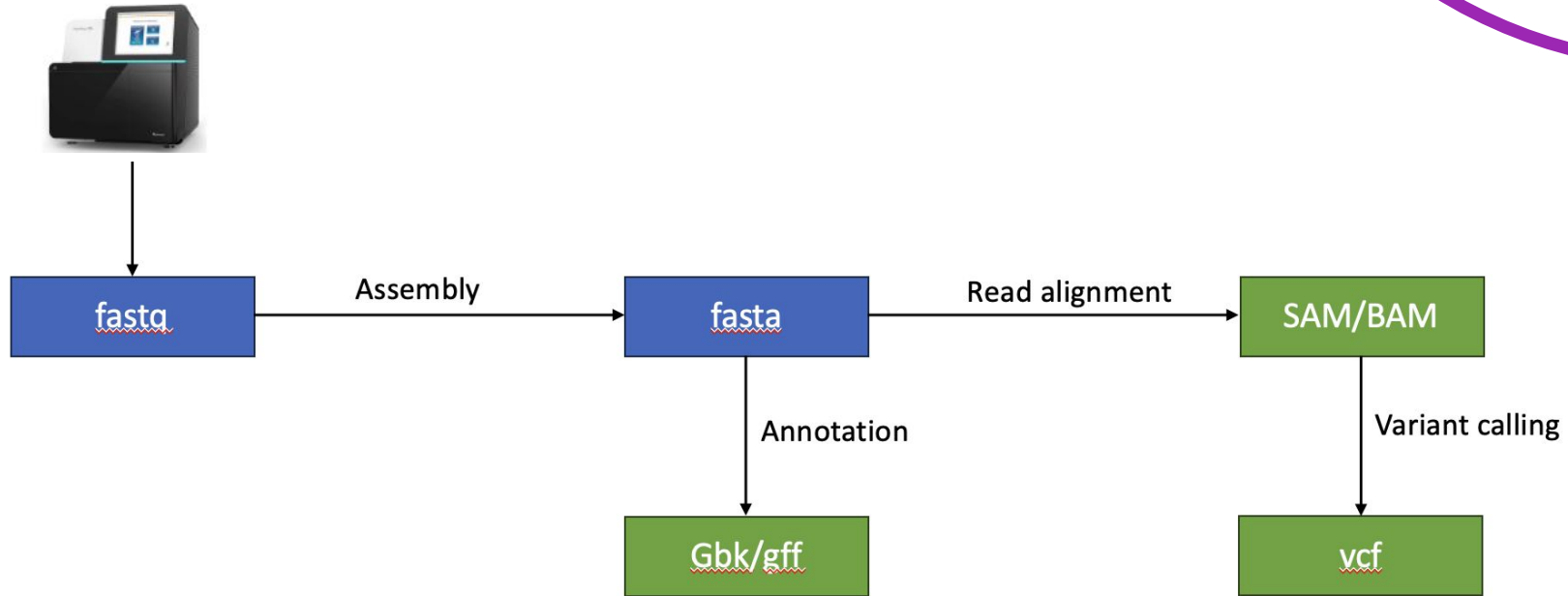
# Basic Unix: Commands

- Navigation through files and directories
  (*cd, ls, pwd*)

- File and directory management
  (mkdir, touch, nano, cat, less, rm, mv, cp, head, tail)

- Pipes
  (we, sort)

- Searching and filtering
  (*grep, find, locate*)

- Process control
  (*kill*)

- Tricks and shortcuts
  (tab complete, wildcard, arrows)



https://cdn.hostinger.com/tutorials/pdf/Linux-Commands-Cheat-Sheet.pdf?_gl=1*1qpygsr*_gcl_au*MTc3ODMzODAwMi4xNzE0OTM3MDUx&_ga=2.258160322.999754537.1714937051-998924758.1714937051

# Introduction to the basic file formats
# - Data transformation

# Introduction to the basic file formats

- Data is information contained in a file

- Different file formats carry data in a specific design and optimization for some programs to be able to read and display the information in an understandable way

- The basic progression of sequencing data and formats involved:

    - **Fastq** – raw sequence reads

    - **Fasta** – genome assembly

    - **.gff/.gbk** – annotation files

    - **SAM/BAM** – read alignment

    - **VCF** – variant calling

# Introduction to the basic file formats

## Fastq
- Standard format directly from the sequencing instrument

Identifier —— @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence —— TTGCCTGCCTATCATTTTAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+' sign —— +
Quality scores —— hhhhhhhhhhghhghhhhhfhhhhhffffffe'ee['X]b[d[ed'[Y[^Y
Identifier —— @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence —— GATTTGTATGAAAGTATACAACTAAAACTGCAGGTGGATCAGAGTAAGTC
'+' sign —— +
Quality scores —— hhhhgfhhcghghggfcffdhfehhhhcehdchhdhahehffffde'bVd

## Fasta
- Following genome assembly
- fasta format is much simpler
- Contains only an identifier and the sequence
- Can also be amino acids (20 letters)

>NG_008679.1:5001-38170 Homo sapiens paired box 6 (PAX6)
ACCCTCTTTTCTTATCATTGACATTTAAACTCTGGGGCAGGTCCTCGCGTAGAACGCGGCTGTCAGATCT
GCCACTTCCCCTGCCGAGCGGCGGTGAGAAGTGTGGGAACCGGCGCTGCCAGGCTCACCTGCCTCCCCGC
CCTCCGCTCCCAGGTAACCGCCCGGGCTCCGGCCCCGGCCCGGCTCGGGGCCCGCGGGGCCTCTCCGCTG
CCAGCGACTGCTGTCCCCAAATCAAAGCCCGCCCCAAGTGGCCCCGGGGCTTGATTTTTGCTTTTAAAAG
GAGGCATACAAAGATGGAAGCGAGTTACTGAGGGAGGGATAGGAAGGGGGGTGGAGGAGGGACTTGTCTT
TGCCGAGTGTGCTCTTCTGCAAAAGTAGCAAAATGTTCCACTCCTAAGAGTGGACTTCCAGTCCGGCCCT
GAGCTGGGAGTAGGGGGCGGGAGTCTGCTGCTGCTGTCTGCTAAAGCCACTCGCGACCGCGAAAAATGCA
GGAGGTGGGGACGCACTTTGCATCCAGACCTCCTCTGCATCGCAGTTCACGACATCCACGCTTGGGAAAG
TCCGTACCCGCGCCTGGAGCGCTTAAAGACACCCTGCCGCGGGTCGGGCGAGGTGCAGCAGAAGTTTCCC
GCGGTTGCAAAGTGCAGATGGCTGGACCGCAACAAAGTCTAGAGATGGGGTTCGTTTCTCAGAAAGACGC

# Introduction to the basic file formats

## Multifasta
- Concatenated fasta files usually following multisequence alignment

```
            >VIT_201s0011g03530.1
Header  ●
Sequence  ● AATTAAGCATAAATACTCACTCTTACCCCCTTATTTTCTTATCTCTCATCACTTTTGGTGCGAAG
            GACCATGAGAACAAGCTGCAATGGGTGTAGGGTTCTTCGCAAGGCATGCAGCCAAGACTGCATCA
Header  ●   >VIT_201s0011g03540.1
Sequence  ● CAGGTAGCGTGAAGTTAAACCCTAGCGCTTTAGACAAACAGCTGTAGTCACCGCCCACAAACACC
            AGCCTCTGAGACACCACCTCAAACCTTTCCACTTAAATACACATCCCTCACACCCTTTTCAATTC
Header  ●   >VIT_201s0011g03550.1
Sequence  ● CATGCAAAGCTGAACGCGATGCTGCTGTGATTGGTGGTAAGTGGTAGTTGAGTAAATTTGACAGTGAA
            GCCGAAATGGTAAAAGACTAAGGCTAGAAGTAGAATACCACTGTTCTTCTCATCACGTGGGCCCA
```

## Genbank (gbk)
- Genome annotation output
- Combination of computer and human-readable information
- Still need another application to visualize annotations

```
LOCUS       NC_002516              6264404 bp    dna     circular UNK
DEFINITION  Pseudomonas aeruginosa PAO1 chromosome, complete genome.
ACCESSION   NC_002516
FEATURES             Location/Qualifiers
     source          1..6264404
                     /mol_type="genomic DNA"
                     /db_xref="taxon:208964"
                     /strain="PAO1"
                     /organism="Pseudomonas aeruginosa PAO1 (Reference)"
     gene            483..2027
                     /locus_tag="PA0001"
                     /db_xref="GeneID:878417"
                     /name="dnaA"
     CDS             483..2027
                     /locus_tag="PA0001"
                     /db_xref="GeneID:878417"
                     /translation="MSVELWQQCVDLLRDELPSQQFNTWIRPLQVEAEGDELRVYAPN
                     RFVLDWVNEKYLGRLLELLGERGEGQLPALSLLIGSKRSRTPRAAIVPSQTHVAPPPP
                     VAPPPAPVQPVSAAPVVVPREELPPVTTAPSVSSDPYEPEEPSIDPLAAAMPAGAAPA
                     VRTERNVQVEGALKHTSYLNRTFTFENFVEGKSNQLARAAAWQVADNLKHGYNPLFLY
                     GGVGLGKTHLMHAVGNHLLKKNPNAKVVYLHSERFVADMVKALQLNAINEFKRFYRSV
                     DALLIDDIQFFARKERSQEEFFHTFNALLEGGQQVILTSDRYPKEIEGLEERLKSRFG
                     WGLTVAVEPPELETRVAILMKKAEQAKIELPHDAAFFIAQRIRSNVRELEGALKRVIA
                     HSHFMGRPITIELIRESLKDLLALQDKLVSIDNIQRTVAEYYKIKISDLLSKRRSRSV
                     ARPRQVAMALSKELTNHSLPEIGVAFGGRDHTTVLHACRKIAQLRESDADIREDYKNL
                     LRTLTT"
                     /product="chromosomal replication initiator protein DnaA"
```

# Introduction to the basic file formats

## Sequence alignment map (SAM)

- A tab-delimited, line oriented text format of information on alignments

- Header section with metadata in each column

- Alignment section with corresponding information on the alignment



https://samtools.github.io/hts-specs/SAMv1.pdf

# Introduction to the basic file formats

**Binary alignment map (BAM)**

- Compressed version of SAM
- Not human-readable
- Can be recognized by certain programs that allow you to visualize the alignment (e.g IGV, Artemis)
- Ready for variant calling

# Introduction to the basic file formats

**Variant call format (VCF)**

- Output from extraction of variance/variations between the query sequence and the reference sequence directly from the BAM file.

# Introduction to the basic file formats
# - Data transformation

Acknowledgements:

- Collins Kigen (Kenya)
- Jorge da Rocha (United kingdom)
- Progress Dube (Zimbabwe)
- Marcela Suarez Esquivel (Cosat Rica)