# INTRODUCTION TO UNIX-LINUX – ALIGNMENT/MAPPING

*Klebsiella Workshop*

*Sep 2024*

ASLM — AFRICAN SOCIETY FOR LABORATORY MEDICINE

UNIVERSITY *of the* WESTERN CAPE

SANBI — South African National Bioinformatics Institute

BIG DATA INSTITUTE / UNIVERSITY OF OXFORD

African Union

AFRICA CDC — Centres for Disease Control and Prevention — Safeguarding Africa's Health

PUBLIC HEALTH ALLIANCE FOR GENOMIC EPIDEMIOLOGY

# Introduction to Alignment / Mapping

- A sequence alignment is a way of arranging the primary sequence/reads of DNA, RNA or Proteins to identify regions of similarity that may be a consequences of FUNCTIONAL, STRUCTURAL and EVOLUTIONARY relationship between the sequences

- Sequence alignment is the procedure of comparing two (pair-wise alignment) or more (multiple sequences) by searching for a series of individual characters/patterns that are in the same order in the sequences.
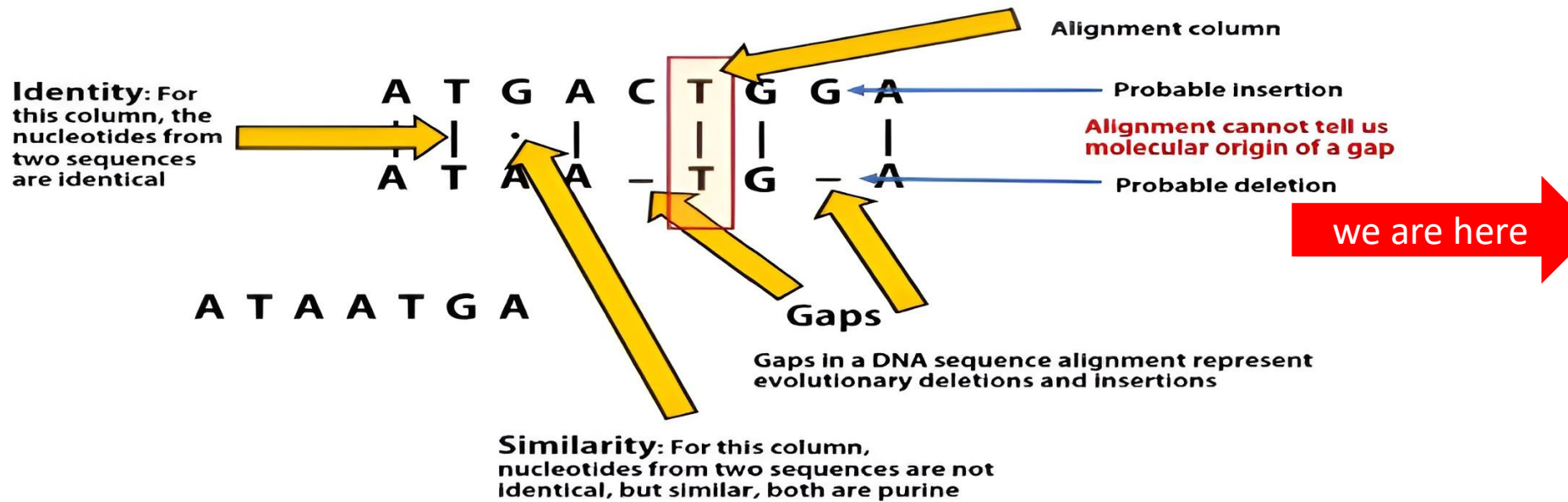
we are here

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Introduction to Alignment / Mapping cont.

Identity: For this column, the nucleotides from two sequences are identical

Alignment column

Probable insertion

Alignment cannot tell us molecular origin of a gap

Probable deletion

Gaps

Gaps in a DNA sequence alignment represent evolutionary deletions and insertions

Similarity: For this column, nucleotides from two sequences are not identical, but similar, both are purine

- Known sequence is the REFERENCE SEQUENCE, and unknown is the QUERY SEQUENCE.
- Align QUERY SEQUENCE to CONTIGS

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# What is Alignment / Mapping



- Short reads must be combined into larger fragments
- **Mapping:** use a reference genome as a guide
- **Alignment:** referred to as alignment of query sequences to contigs. Also have gene alignment
- **De Novo assembly:** without a reference genome
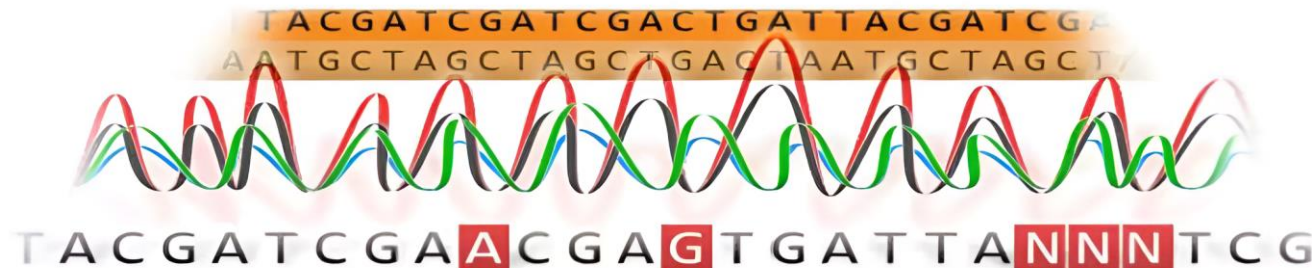
Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Why is Alignment needed?

- To compare sequences for similarities and differences
  - Mutations
- Often, we are looking for similarities
- Homology: similarity due to descent from a common ancestor
- For traceability
- To infer differences and similarities in structures, function and sequences



**Quality control**

**Adaptor / Trimming**

we are here →

**Alignment /mapping**

**Variant calling**

**Tree Building**

# Mapping in Bioinformatics



## Example NGS pipeline

**Sequence reads**

FASTQ

**Quality control**

FASTQ

**Alignment to Genome**

SAM/BAM

**Alignment Cleanup**
**BAM ready for variant calling**

BAM

**Variant Calling**

VCF

A high level view of a typical NGS bioinformatics workflow

we are here

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Mapping in Bioinformatics

## Sequence alignment

- Determine position of short read on the reference genome

```
Reference: . . . A A - C G C C T T . . .              | = match
.                  | : - : | | | | |                  : = mismatch
Read:              A G G G C C T T                     - = gap
```

- Read could align to multiple places



we are here

- How to handle multi-mapped reads? Depends on tool:
    - Map to best region (but what is "best"? And what about ties?)
    - Map to all regions
    - Map to one region randomly
    - Discard read
- How do we determine *best* region?
    - Assign **alignment score** to every mapping

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Mapping in Bioinformatics

## Sequence Alignment

```
Reference: AAA CAGTGA GAA
Observed:  AAA TCTCT  GAA
```

| Alignment | | Tool | Variant calls |
|---|---|---|---|
| `AAA-CAGTGAGAA`<br>`\|\|\|-\|--\|::\|\|\|`<br>`AAATC--TCTGAA` | Maybe like this? | Novoalign | `ins T`<br>`del AG`<br>`sub GA -> CT` |
| `AAACAGTGAGAA`<br>`\|\|\|-::\|::\|\|\|`<br>`AAA-TCTCTGAA` | Or this? | Ssaha2 | `del C`<br>`sub AG -> TC`<br>`sub GA -> CT` |
| `AAACAGTGAGAA`<br>`\|\|\|:-:\|::\|\|\|`<br>`AAAT-CTCTGAA` | Or..? | BWA | `snp C -> T`<br>`del A`<br>`snp G -> C`<br>`sub GA -> CT` |
| `AAACAGTCA-----GAA`<br>`\|\|\|----------\|\|\|`<br>`AAA------TCTCTGAA` | What about this? | Complete Genomics | `del CAGTGA`<br>`ins TCTCT` |

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Mapping in Bioinformatics

## Mapping tools

| Mapping tool | Uses | Characteristics |
|---|---|---|
| HISAT2 | DNA/RNA | Short reads. Based on GCSA. Reference. |
| RNASTAR | RNA | Short reads. Extremely fast. High sensitive and accuracy. Based on Maximal Mappable Prefixes (MMPs). Reference. |
| BWA-MEM2 | DNA | Short reads. Twice as faster as BWA-MEM. Memory efficient. Based on Burrows-Wheeler. Reference. |
| Minimap2 | DNA/RNA | Long reads (PacBio and ONT). Extremely fast. Based on DALIGN and MHAP. Reference. |
| Bismark | DNA/RNA | Short reads. Bisulfite treated sequencing. Based on GCSA. Reference. |
| BBMap | DNA/RNA | Short and long reads (PacBio and ONT). Memory demanding. Reference. |
| Whisper 2 | DNA | Short reads. Indel sensitive. Variant-calling oriented. Reference. |
| S-conLSH | DNA | Long reads (ONT). High sensitivity and accuracy. Reference. |

we are here

- There are many more tools, the mapping tool is chosen based on which best fits your data

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Mapping in Bioinformatics



**SAM/BAM file format**

**SAM:** **S**equence **A**lignment **M**ap

**BAM:** Binary (compressed) SAM; not human-readable

we are here

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Tool: SNIPPY

- An alignment and variant calling tool – all-in-one
- Snippy finds SNPs between a haploid reference genome and your NGS sequence reads. It will find both substitutions (SNPS) and insertions/deletions (indels).
- Input: NGS Reads in fastq format (SE or PE) & a Reference file in either fasta or genbank format
- Various parameters can be defined "options"
  - `--contigs` allows you to call SNPs from contigs rather than reads.
  - `--rgid` will set the Read Group (`RG`) ID (`ID`) and Sample (`SM`) in the BAM and VCF file.
- Simplicity, input → output

**Quality control**

**Adaptor / Trimming**

we are here ➔ **Alignment /mapping**

we are here ➔ **Variant calling**

**Tree Building**

# How SNIPPY works

**Output Files**

| Extension | Description |
|-----------|-------------|
| .tab | A simple tab-separated summary of all the variants |
| .csv | A comma-separated version of the .tab file |
| .html | A HTML version of the .tab file |
| .vcf | The final annotated variants in VCF format |
| .bed | The variants in BED format |
| .gff | The variants in GFF3 format |
| .bam | The alignments in BAM format. Includes unmapped, multimapping reads. Excludes duplicates. |
| .bam.bai | Index for the .bam file |
| .log | A log file with the commands run and their outputs |
| .aligned.fa | A version of the reference but with `-` at position with `depth=0` and `N` for `0 < depth < --mincov` (does not have variants) |
| .consensus.fa | A version of the reference genome with *all* variants instantiated |
| .consensus.subs.fa | A version of the reference genome with *only substitution* variants instantiated |
| .raw.vcf | The unfiltered variant calls from Freebayes |
| .filt.vcf | The filtered variant calls from Freebayes |
| .vcf.gz | Compressed .vcf file via BGZIP |
| .vcf.gz.csi | Index for the .vcf.gz via `bcftools index` ) |

- SNIPPY has the following incorporated:
  - BWA-MEM – maps individual sequence
  - SAM tools – converts SAM→ BAM
  - GNU parallel – executing job in parallel (isolates)
  - Freebayes – variant calling program
  - VCFlib – parsing and manipulating VCF files
  - VCFtools – output .vcf file
  - SNPEFF – Prediction tool, annotates and predicts the effects of genetic variants

Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

Tree Building

# Installation - SNIPPY

- Installation - CLI:
  - `cd $HOME`
  - `git clone https://github.com/tseemann/snippy.git`
  - `$HOME/snippy/bin/snippy –help`

- Installation – Bioconda:
  - `conda install -c conda-forge -c bioconda -c defaults snippy`

Quality control

Adaptor / Trimming

Alignment /mapping

**we are here**

Variant calling

**we are here**

Tree Building

# Usage - SNIPPY

- Usage:
  - a reference genome in FASTA or GENBANK format (can be in multiple contigs)
  - sequence read file(s) in FASTQ or FASTA format (can be .gz compressed) format
  - a folder to put the results in

- https://github.com/tseemann/snippy

Quality control

Adaptor / Trimming

we are here →

Alignment /mapping

we are here →

Variant calling

Tree Building

# VCF - Output



Quality control

Adaptor / Trimming

Alignment /mapping

Variant calling

we are here

Tree Building