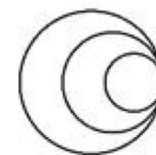




H3ABioNet

Pan African Bioinformatics Network for H3Africa



**wellcome
connecting
science**

Next Generation Sequencing Bioinformatics Course 2021

Module 6 – Pathogen variant calling Session 1: Alignment to reference



H3ABioNet

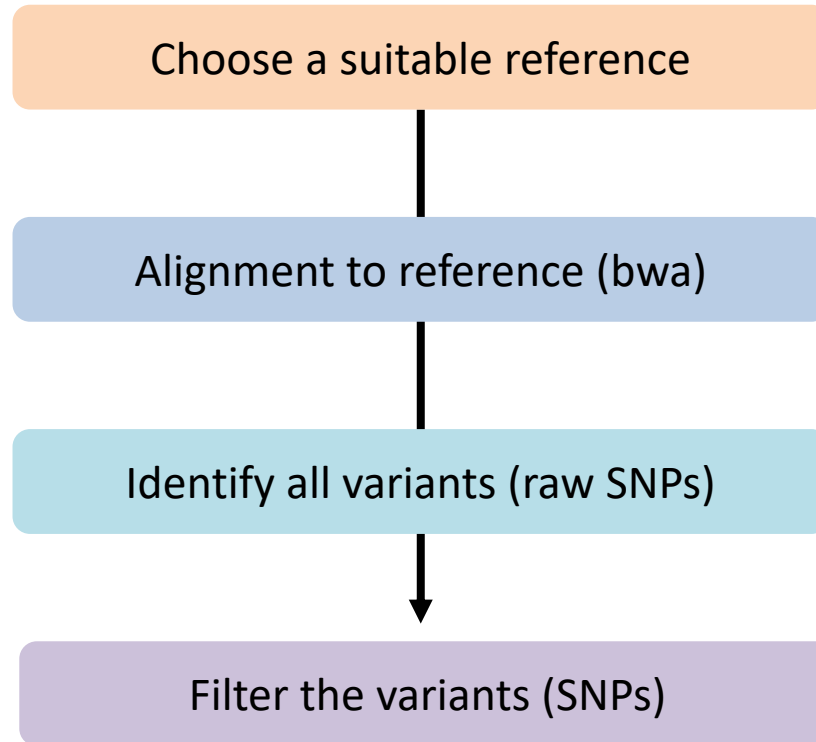
Pan African Bioinformatics Network for H3Africa

WELLCOME GENOME CAMPUS
CONNECTING
SCIENCE
ADVANCED
COURSES+
SCIENTIFIC
CONFERENCES



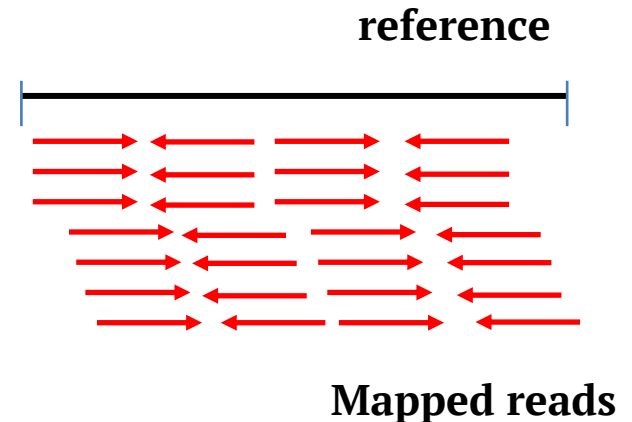
NGS Bioinformatics Course Africa 2021
Trainer name

Detection of variants (Variant calling)

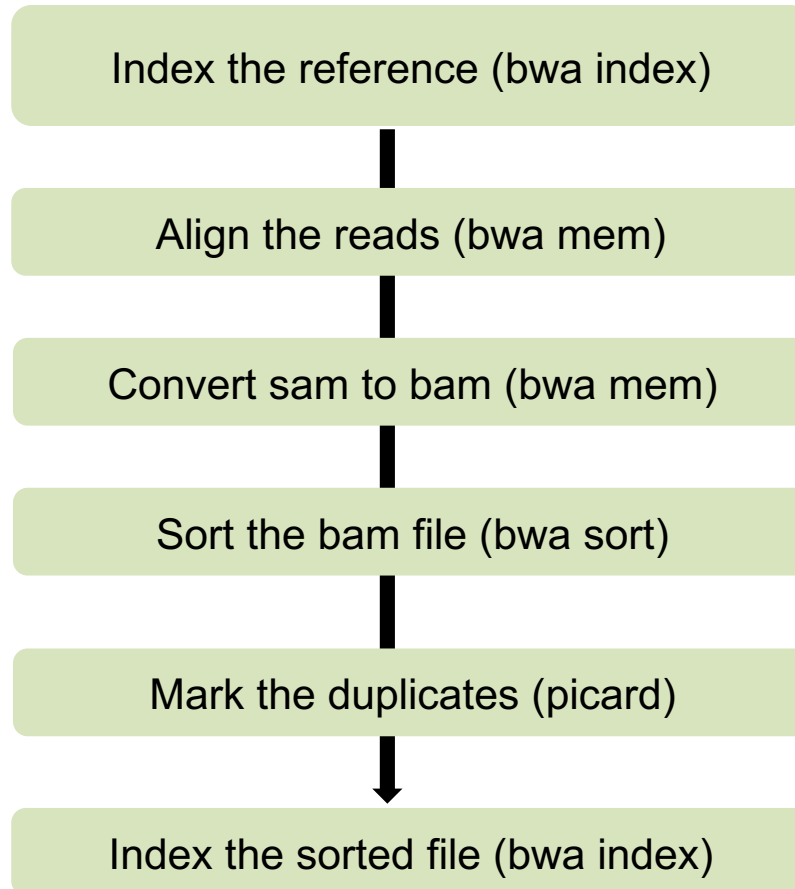


Selecting suitable reference

- Best source: NCBI reference database or prior published studies
- Complete genome sequence available
 - NCBI has versions : select the latest
- Genome annotation
- From the same species or close relative
- Present case : *M. tuberculosis* H37Rv (NC_000962.3)



Step1: Alignment to reference



Duplicates

- 1) `picard MarkDuplicates I=sorted.bam O=markdup.bam M=metrics.txt`
- 2) `grep -A 2 “^## METRICS” metrics.txt`

```
1 ## METRICS CLASS      picard.sam.DuplicationMetrics
2 LIBRARY              UNPAIRED_READS_EXAMINED
  READ_PAIRS_EXAMINED   SECONDARY_OR_SUPPLEMENTARY_RDS
  UNMAPPED_READS        UNPAIRED_READ_DUPLICATES
  READ_PAIR_DUPLICATES  READ_PAIR_OPTICAL_DUPLICATES
  PERCENT_DUPLICATION   ESTIMATED_LIBRARY_SIZE
3 Unknown Library      231556 2524511          30570 289914
  120543 300973 49885 0.13682          11354842|
```



Mapping statistics

1) `samtools stats in.bam>bamstat.txt`

2) `grep “^SN” bamstats.txt >stats.txt`

```
1 SN raw total sequences: 5570492
2 SN filtered sequences: 0
3 SN sequences: 5570492
4 SN is sorted: 1
5 SN 1st fragments: 2785246
6 SN last fragments: 2785246
7 SN reads mapped: 5280578
8 SN reads mapped and paired: 5049022 # paired-end technology bit set + both mates mapped
9 SN reads unmapped: 289914
10 SN reads properly paired: 4995486 # proper-pair bit set
11 SN reads paired: 5570492 # paired-end technology bit set
12 SN reads duplicated: 722489 # PCR or optical duplicate bit set
13 SN reads MQ0: 69525 # mapped and MQ=0
14 SN reads QC failed: 0
15 SN non-primary alignments: 0
16 SN total length: 835573800 # ignores clipping
17 SN total first fragment length: 417786900 # ignores clipping
18 SN total last fragment length: 417786900 # ignores clipping
19 SN bases mapped: 792086700 # ignores clipping
20 SN bases mapped (cigar): 776568548 # more accurate
21 SN bases trimmed: 0
22 SN bases duplicated: 108373350
23 SN mismatches: 5722273 # from NM fields
24 SN error rate: 7.368664e-03 # mismatches / bases mapped (cigar)
25 SN average length: 150
26 SN average first fragment length: 150
27 SN average last fragment length: 150
28 SN maximum length: 150
29 SN maximum first fragment length: 150
30 SN maximum last fragment length: 150
31 SN average quality: 35.9
32 SN insert size average: 350.1
33 SN insert size standard deviation: 109.3
34 SN inward oriented pairs: 2483610
35 SN outward oriented pairs: 34508
36 SN pairs with other orientation: 19681
37 SN pairs on different chromosomes: 0
38 SN percentage of properly paired reads (%): 89.7
```



Thank you



H3ABioNet

Pan African Bioinformatics Network for H3Africa

WELLCOME GENOME CAMPUS
CONNECTING
SCIENCE
ADVANCED
COURSES+
SCIENTIFIC
CONFERENCES



Next Generation Sequencing Bioinformatics
Trainer Name: Narender Kumar