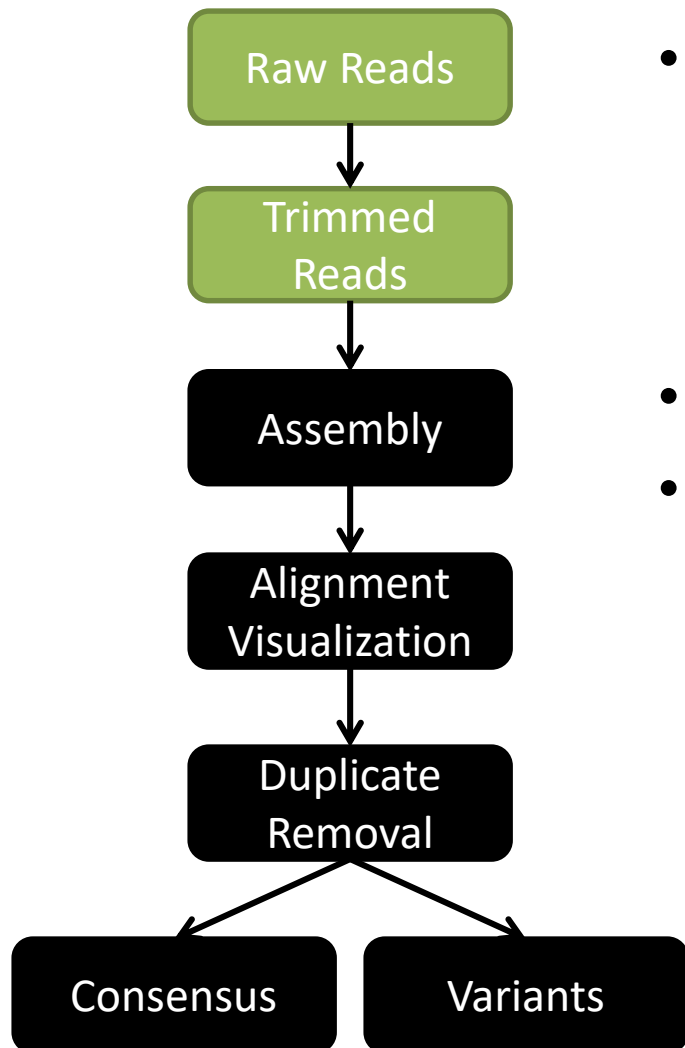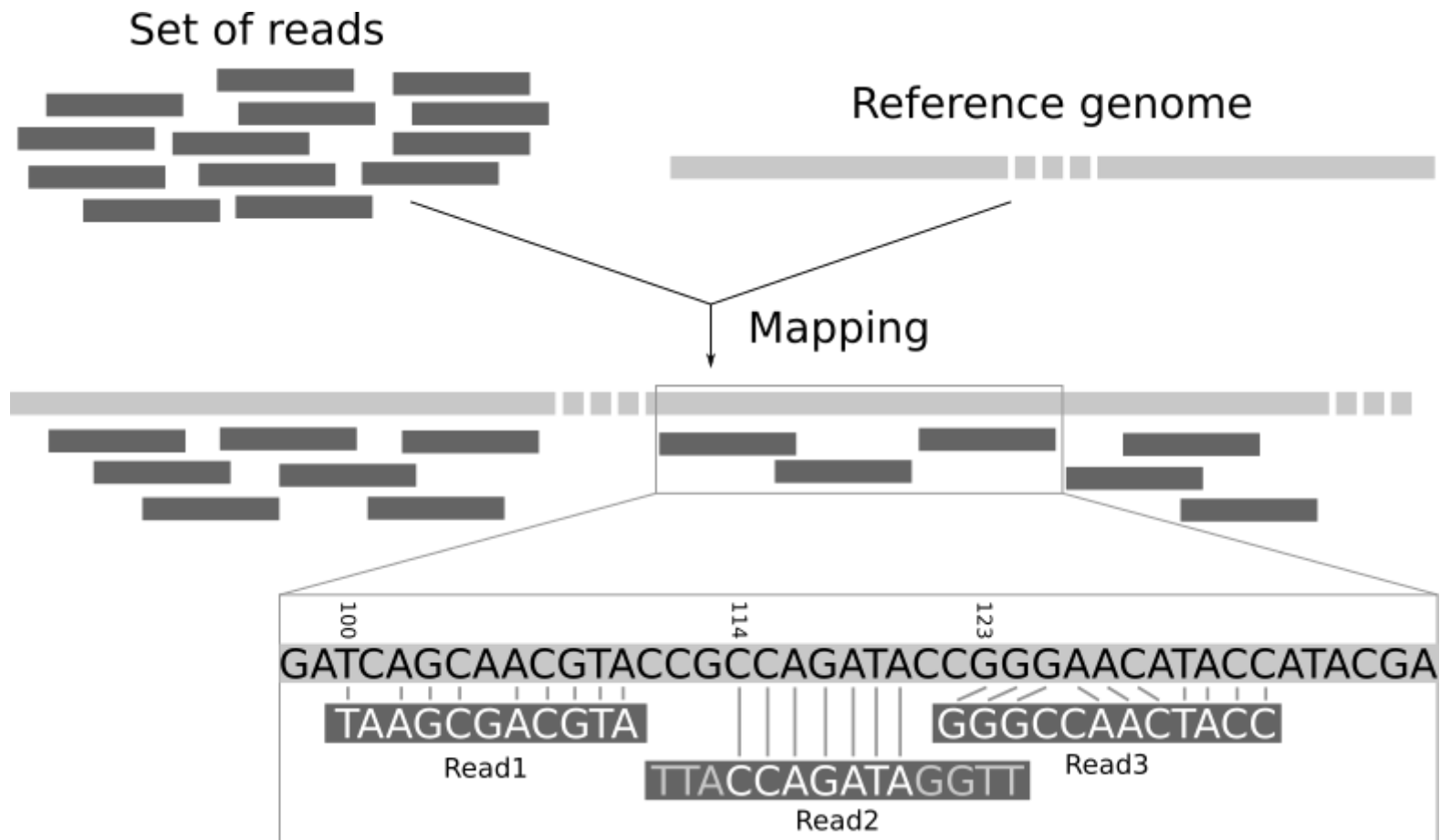# Reference Based Assembly

# Why do we assemble reads?

- Second Generation sequencing fragments genomic material which is then sequenced to give the reads

- We assemble these reads to identify the exact position of the genome they come from and their base-to-base correspondence.

- The two main methods of assembly are

  ➢ Reference Bases Alignment

  ➢ Denovo Assembly

# Reference

Raw Reads

↓

Trimmed Reads

↓

Assembly

↓

Alignment Visualization

↓

Duplicate Removal

Consensus     Variants

- All we need for reference alignment are the trimmed reads and a reference sequence.

- What is a suitable reference?
- Can you map any reads to any reference?

# Perfect Alignment

Pos:   1             10                20                30

Ref:   ACGTAGCATTGCTAGCGGTTAACCGTAGCT

```
       CGTAG     TGCTA     GTTAA CGTAG
       ACGTA             TAGCG           TAGCT
                ATTGC          GGTTA
         TAGCA       CTAGC       TAACC
       ACGTA                GCGGT        TAGCT
```

# Viral Sequencing

- Viruses exist as a quasispecies distribution which introduces variants

- Variants can also be introduced by PCR errors and poor-quality sequencing

- Variants can be identified from read alignments

```
Reference  CCGTTAGAGTTACAATTCGA
Read 2         TTAGAGTAACAA
Read 3     CCGTTAGAGTTA
Read 4               TTACAATTCGA
Read 5          GAGTAACAA
Read 6         TTAGAGTAACAAT
```

# Viral Alignment

Pos:   1              10              20             30

Ref:   ACGTAGCATTGCTAGCGGTTAACCGTAGCT

    CGTAG     TGATA     GTTAA  CGTAG

    ACGTA        TAGCG          T-GCA

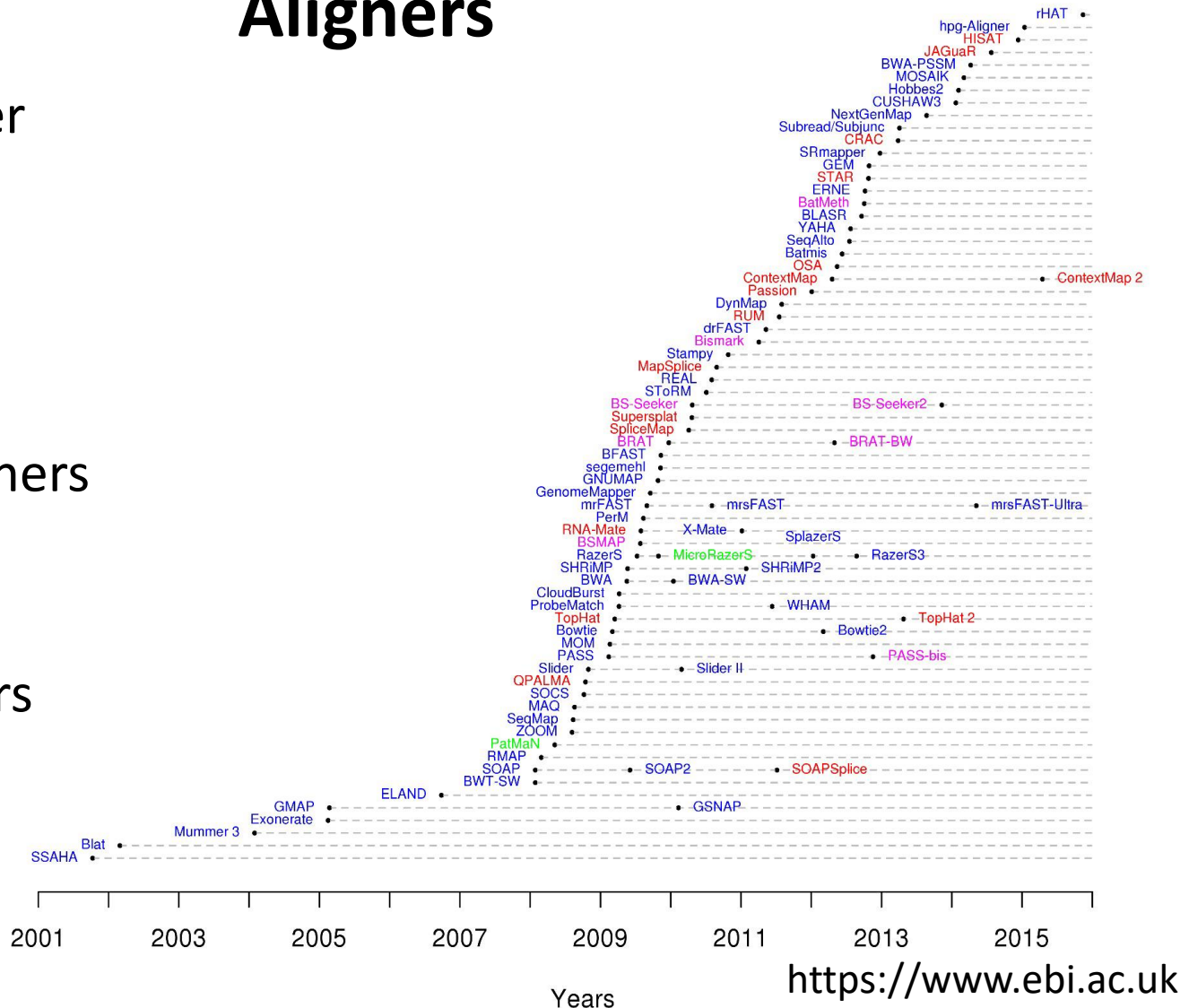         ATTGC      GCTTA

      TTGCA     CTAGC    TAACC

    ACGTA          GCGGT      TAGCT

CIGAR: Pos4 1M1X3M

CIGAR: Pos 26 1M1D2M1X

# Aligners

- Short Read Aligner
  - BWA
  - Bowtie
  - Tanoti
  - NovoAlign

- Splice Aware Aligners
  - Tophat
  - BBMap

- Long Read Aligners
  - Minimap
  - LAST



https://www.ebi.ac.uk

# Which Aligner?

- Hash based Aligners

- Burrows Wheeler Aligners

- Global vs Local Aligners

- Effect on Consensus Calling

- Effect on Variant Calling

**Always good to keep aligners consistent across any experiment**

# Which Reference?

- Choosing the wrong species

- Choosing the wrong genotypes (divergent sequences)

- If you do not have information on correct reference may have to map to panel of references

- If species not known Denovo Assembly

# Identifying incorrect Reference

| | Bit | Description |
|---|---|---|
| 1 | 0x1 | template having multiple segments in sequencing |
| 2 | 0x2 | each segment properly aligned according to the aligner |
| 4 | 0x4 | segment unmapped |
| 8 | 0x8 | next segment in the template unmapped |
| 16 | 0x10 | SEQ being reverse complemented |
| 32 | 0x20 | SEQ of the next segment in the template being reverse complemented |
| 64 | 0x40 | the first segment in the template |
| 128 | 0x80 | the last segment in the template |
| 256 | 0x100 | secondary alignment |
| 512 | 0x200 | not passing quality controls |
| 1024 | 0x400 | PCR or optical duplicate |
| 2048 | 0x800 | supplementary alignment |

| QNAME | FLAG | RNAME | POS | MAPQ | CIGAR | RNEXT | PNEXT | TLEN | SEQ | QUALITY |
|---|---|---|---|---|---|---|---|---|---|---|
| ReadN | 4 | * | 0 | 0 | * | * | 0 | 0 | ACGTAG | IHGFFF |
| ReadN2 | 4 | * | 0 | 0 | * | * | 0 | 0 | GGGGGGG | IIIHHGG |

- Looking at assembly stats
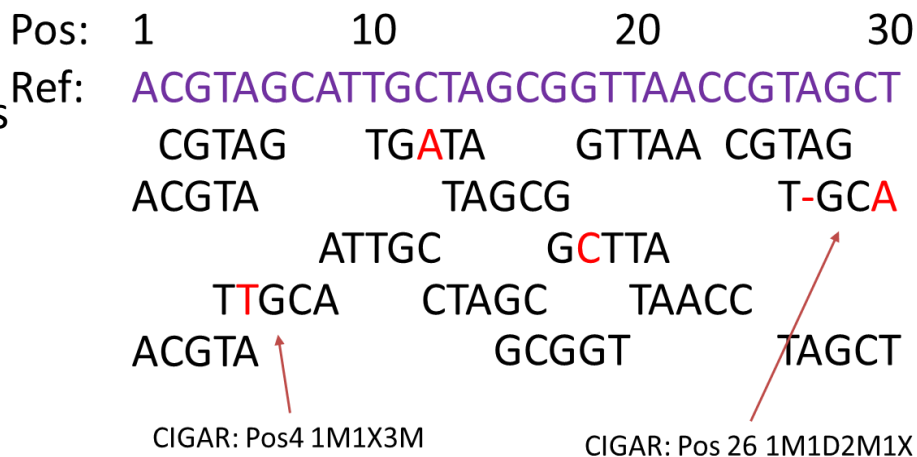- Visually inspecting the alignment

# Sequence Alignment

- Usually output as SAM file
  - Sequence Alignment Map contains coordinates and quality of mapping.

- Converted to BAM file
  - Binary Alignment Map is compressed

- Files are sorted and indexed
  - Makes it faster to access and work with the data

Pos:    1           10          20          30
Ref:    ACGTAGCATTGCTAGCGGTTAACCGTAGCT
        CGTAG       TGATA       GTTAA  CGTAG
        ACGTA           TAGCG              T-GCA
            ATTGC       GCTTA
            TTGCA       CTAGC       TAACC
        ACGTA                   GCGGT       TAGCT

CIGAR: Pos4 1M1X3M          CIGAR: Pos 26 1M1D2M1X

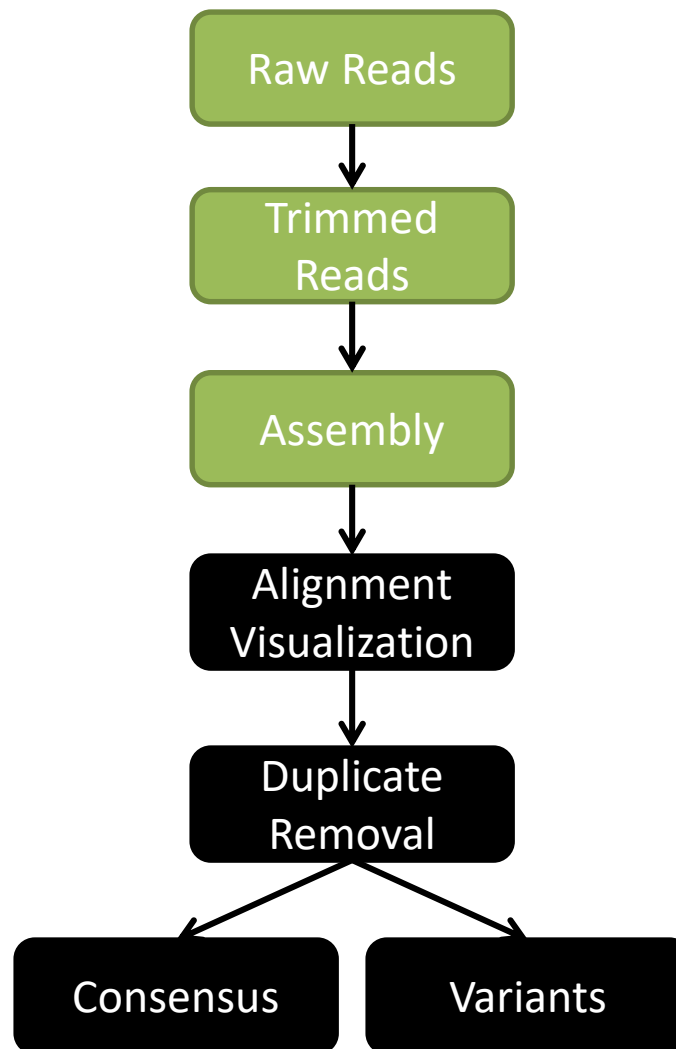Pos:    1           10          20          30
Ref:    ACGTAGCATTGCTAGCGGTTAACCGTAGCT
        ACGTA   ATTGC       GCTTA   CGTAG
        ACGTA       TGATA       GTTAA   T-GCA
            CGTAG           CTAGC       TAACC  TAGCT
            TTGCA           TAGCG
                            GCGGT

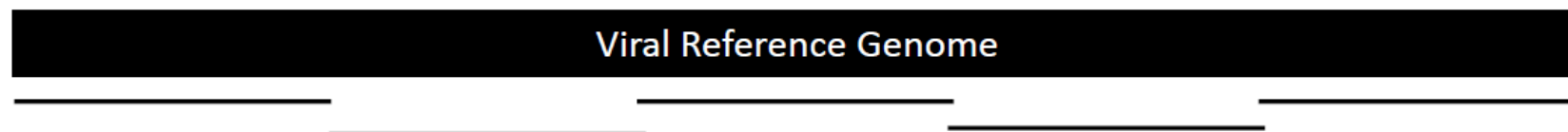CIGAR: Pos4 1M1X3M          CIGAR: Pos 26 1M1D2M1X

# Alignment Stats/Visualization

- Summary stats give an idea of the mapping
  - Mapped vs Unmapped reads
  - Average depth of coverage
  - Breadth of Coverage
- Coverage Plots
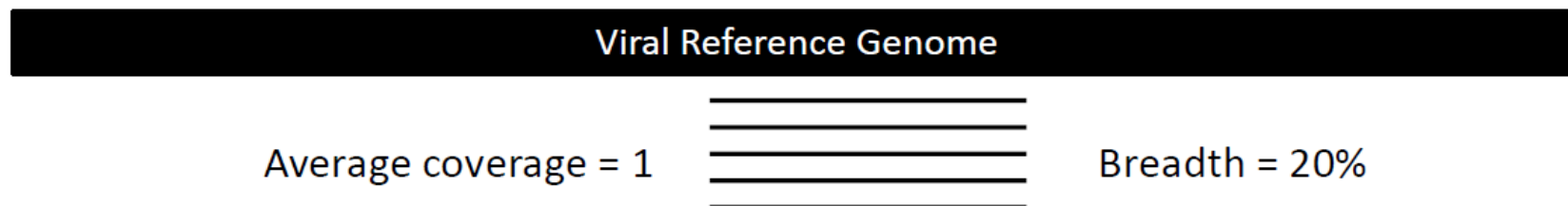- Visualization of complete alignment

# Coverage

```
                   1                   2                   3
Pos:   1234567890123456789012345678901234
Ref:   ACGGTGACACGTAGCAGTACGCGGGTTACACAGA
       ACGGCGA          CAGTTCG        AC-CAGA
           AGACGTA          GCGGGTT
               GTAGCAGT          TTACACAG
         GCGACAC          TCGCGGG
       CGGCGAC          AGTTCGC      TACACAT
           ACG-AGC          GGGGTAC
Cov:   1223334333333323333334333334443331
```

# Coverage Depth vs Breadth



Viral Reference Genome

Average coverage = 1                    Breadth = 100%
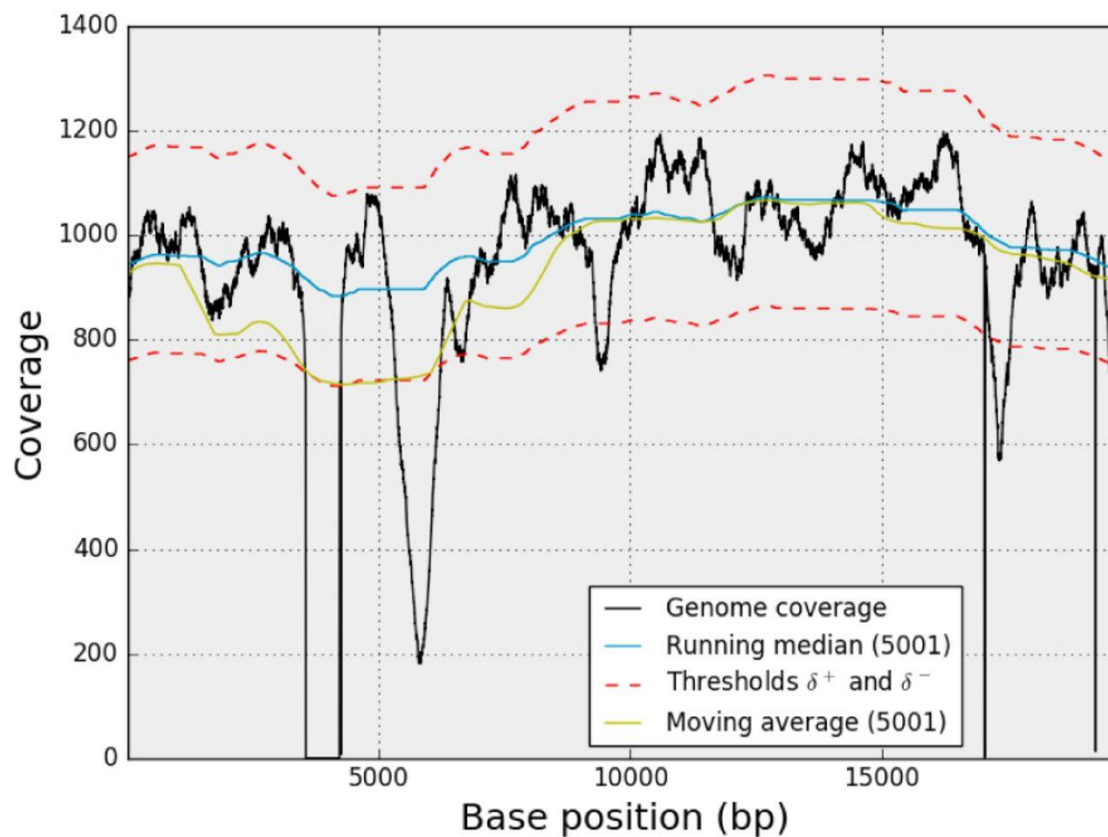
Viral Reference Genome

Average coverage = 1                    Breadth = 20%

Mode, Median, Quartiles would be different

# Coverage Plot

# Visualization - Tablet



BAM File + Reference File

https://ics.hutton.ac.uk/tablet/
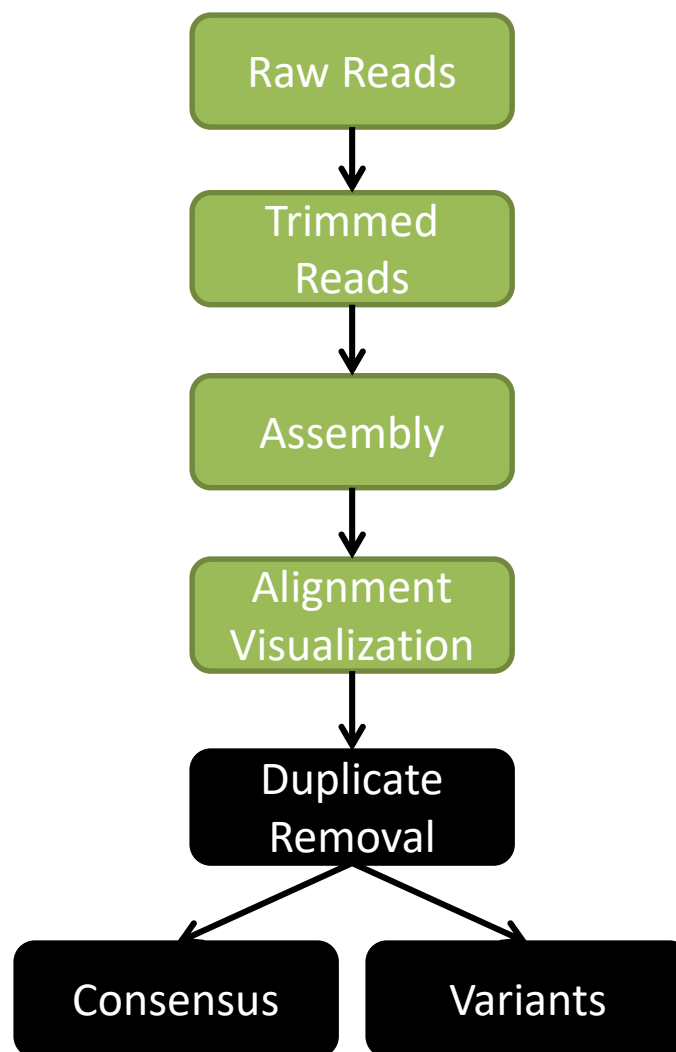
# Duplicate Removal

- Types of Duplicates
    - Optical Duplicates
    - PCR Duplicates
- Relationship of starting genomic material, sequencing depth and PCR Duplicates
- Why remove PCR Duplicates
    - Effect on Consensus Calling
    - Effect on Variant Calling

- **Note:** Only remove duplicates if library preparation has PCR steps

# Duplicate Removal

```
                                         *
TTTCATACTAACTAGCCTGCGGTCTGTGTTTCCCGACTTCTGAGTCATGGGGTTTCAATGCCTATAGATTC
                 ...........................C.
                 ............................
                      ..............T...............
                      ....................C.........
                          ...............................
                           ...............................
                              ...............................
                              ......C........................
                               ..C.............................
```

```
                                         *
TTTCATACTAACTAGCCTGCGGTCTGTGTTTCCCGACTTCTGAGTCATGGGGTTTCAATGCCTATAGATTC
                      ...............................
                      ..............T...............
                      ....................C.........
                      ....................C.........
                      ....................C.........
                      ....................C.........
                               ...............................
                               ...............................
                               ...............................
```