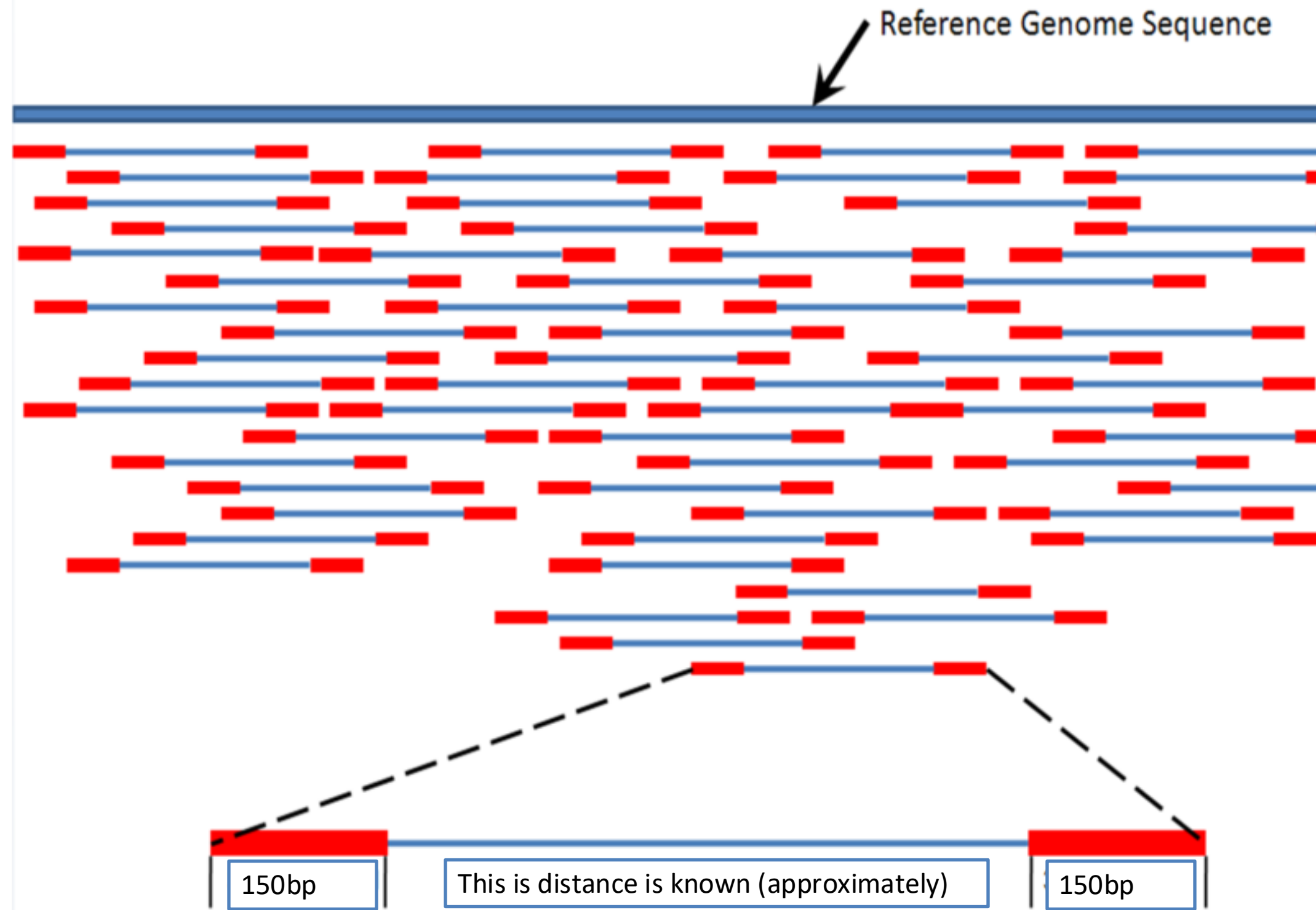# Genômica Computacional
## Analise de variantes

Professor: Ricardo A. Vialle

CS31 - Genômica Computacional

17 de Julho de 2025
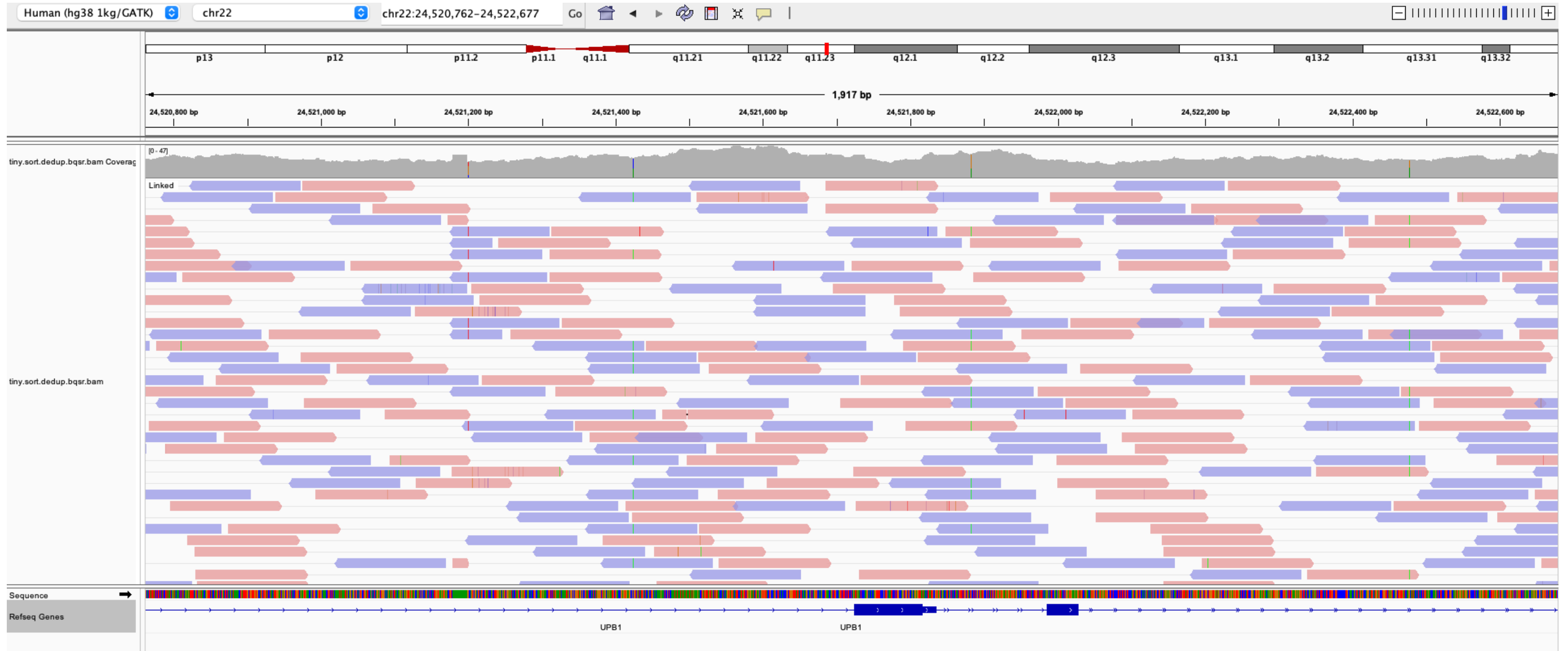
# Read mapping

# Read mapping

# IGV browser

# Types of genetic variation

- **99% of DNA is shared** between two individuals
- Variation in the remainder explains all our **predisposition** differences
- **Remaining** phenotypic variation: environmental/stochastic differences

| Name | Example | Frequency in one genome |
|---|---|---|
| Single nucleotide polymorphisms (**SNPs**) | GAGGAGAACG[C/G]AACTCCGCCG | 1 per 1,000 bp |
| Insertions/deletions (**indels**) | CACTATTC[C/CTATGG]TGTCTAA | 1 per 10,000 bp |
| Short tandem repeats (**STRs**) | ACGGCAGTCGTCGTCGTCACCGTAT | 1 per 10,000 bp |
| Structural variants (**SVs**) / Copy Number Variants (**CNVs**) | Large (median 5,000 bp) deletions, duplications, inversions | 1 per 1,000,000 bp |

# VCF and BCF file format

# SAM and BAM file format



Source: https://www.researchgate.net/figure/A-sample-of-the-SAM-file_fig3_309134977

# SAM and BAM file format (FLAG)

**Bitwise Flags**

| Integer | Binary | Description (Paired Read Interpretation) |
|---------|--------|------------------------------------------|
| 1 | 000000000001 | template having multiple templates in sequencing (read is paired) |
| 2 | 000000000010 | each segment properly aligned according to the aligner (read mapped in proper pair) |
| 4 | 000000000100 | segment unmapped (read1 unmapped) |
| 8 | 000000001000 | next segment in the template unmapped (read2 unmapped) |
| 16 | 000000010000 | SEQ being reverse complemented (read1 reverse complemented) |
| 32 | 000000100000 | SEQ of the next segment in the template being reverse complemented (read2 reverse complemented) |
| 64 | 000001000000 | the first segment in the template (is read1) |
| 128 | 000010000000 | the last segment in the template (is read2) |
| 256 | 000100000000 | not primary alignment |
| 512 | 001000000000 | alignment fails quality checks |
| 1024 | 010000000000 | PCR or optical duplicate |
| 2048 | 100000000000 | supplementary alignment (e.g. aligner specific, could be a portion of a split read or a tied region) |

Use *samtools flagstat* to summarize these metrics

https://en.wikipedia.org/wiki/SAM_(file_format)

# SAM and BAM file format



Source: https://www.researchgate.net/figure/A-sample-of-the-SAM-file_fig3_309134977

# CIGAR

For example:

```
RefPos:      1   2   3   4   5   6   7   8   9  10  11  12  13  14  15  16  17  18  19
Reference:   C   C   A   T   A   C   T   G   A   A   C   T   G   A   C   T   A   A   C
Read: ACTAGAATGGCT
```

Aligning these two:

```
RefPos:      1   2   3   4   5   6   7       8   9  10  11  12  13  14  15  16  17  18  19
Reference:   C   C   A   T   A   C   T       G   A   A   C   T   G   A   C   T   A   A   C
Read:                        A   C   T   A   G   A   A       T   G   G   C   T
```

With the alignment above, you get:

```
POS: 5
CIGAR: 3M1I3M1D5M
```

# SAM and BAM file format



```
@HD VN:1.5 SO:coordinate                                            Header
@SQ SN:ref LN:45                                                    section

r001    99 ref   7 30 8M2I4M1D3M = 37   39 TTAGATAAAGGATACTG *
r002     0 ref   9 30 3S6M1P1I4M * 0     0 AAAAGATAAGGATA    *
r003     0 ref   9 30 5S6M       * 0     0 GCCTAAGCTAA       * SA:Z:ref,29,-,6H5M,17,0;   Alignment
r004     0 ref  16 30 6M14N5M    * 0     0 ATAGCTTCAGC       *                            section
r003  2064 ref  29 17 6H5M       * 0     0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001   147 ref  37 30 9M           = 7 -39 CAGCGGCAT         * NM:i:1
```

**Optional fields** in the format of TAG:TYPE:VALUE

**QUAL:** read quality; * meaning such information is not available

**SEQ:** read sequence

**TLEN:** the number of bases covered by the reads from the same fragment. Plus/minus means the current read is the leftmost/rightmost read. E.g. compare first and last lines.

**PNEXT:** Position of the primary alignment of the NEXT read in the template. Set as 0 when the information is unavailable. It corresponds to POS column.

**RNEXT:** reference sequence name of the primary alignment of the NEXT read. For paired-end sequencing, NEXT read is the paired read, corresponding to the RNAME column.

**CIGAR:** summary of alignment, e.g. insertion, deletion

**MAPQ:** mapping quality

**POS:** 1-based position

**RNAME:** reference sequence name, e.g. chromosome/transcript id

**FLAG:** indicates alignment information about the read, e.g. paired, aligned, etc.

**QNAME:** query template name, aka. read ID

# SAM and BAM file format

```
@HD VN:1.5 SO:coordinate                                                    Header
@SQ SN:ref LN:45                                                            section

r001    99 ref  7 30 8M2I4M1D3M = 37  39 TTAGATAAAGGATACTG *
r002     0 ref  9 30 3S6M1P1I4M * 0    0 AAAAGATAAGGATA     *
r003     0 ref  9 30 5S6M       * 0    0 GCCTAAGCTAA        * SA:Z:ref,29,-,6H5M,17,0;   Alignment
r004     0 ref 16 30 6M14N5M    * 0    0 ATAGCTTCAGC        *                            section
r003  2064 ref 29 17 6H5M       * 0    0 TAGGC              * SA:Z:ref,9,+,5S6M,30,1;
r001   147 ref 37 30 9M         = 7  -39 CAGCGGCAT          * NM:i:1
```

Optional fields in the format of TAG:TYPE:VALUE

QUAL: read quality; * meaning such information is not available

SEQ: read sequence

**TLEN**: the number of bases covered by the reads from the same fragment. Plus/minus means the current read is the leftmost/rightmost read. E.g. compare first and last lines.

PNEXT: Position of the primary alignment of the NEXT read in the template. Set as 0 when the information is unavailable. It corresponds to POS column.

RNEXT: reference sequence name of the primary alignment of the NEXT read. For paired-end sequencing, NEXT read is the paired read, corresponding to the RNAME column.

CIGAR: summary of alignment, e.g. insertion, deletion

MAPQ: mapping quality

POS: 1-based position

RNAME: reference sequence name, e.g. chromosome/transcript id

FLAG: indicates alignment information about the read, e.g. paired, aligned, etc.

QNAME: query template name, aka. read ID

TLEN

This is distance is known (approximately)

150bp                    150bp

# SAM and BAM file format

| Tag ⬍ | Type ⬍ | Description | [hide] ⬍ |
|-------|--------|-------------|----------|
| AM | i | The smallest template-independent mapping quality in the template | |
| AS | i | Alignment score generated by aligner | |
| BC | Z | Barcode sequence identifying the sample | |
| BQ | Z | Offset to base alignment quality (BAQ) | |
| BZ | Z | Phred quality of the unique molecular barcode bases in the OX tag | |
| CB | Z | Cell identifier | |
| CC | Z | Reference name of the next hit | |
| CG | B,I | BAM only: CIGAR in BAM's binary encoding if (and only if) it consists of >65535 operators | |
| CM | i | Edit distance between the color sequence and the color reference (see also NM) | |
| CO | Z | Free-text comments | |
| CP | i | Leftmost coordinate of the next hit | |
| CQ | Z | Color read base qualities | |

Header
section

* 
* 
* SA:Z:ref,29,-,6H5M,17,0;    Alignment
*                              section
* SA:Z:ref,9,+,5S6M,30,1;
* NM:i:1

Optional fields in the format of TAG:TYPE:VALUE

QUAL: read quality; * meaning such information is not available

...ce

...overed by the reads from the same fragment. Plus/minus
...leftmost/rightmost read. E.g. compare first and last lines.

...gnment of the NEXT read in the template. Set as 0 when the
...information is unavailable. It corresponds to POS column.

RNEXT: reference sequence name of the primary alignment of the NEXT read. For paired-end
sequencing, NEXT read is the paired read, corresponding to the RNAME column.

CIGAR: summary of alignment, e.g. insertion, deletion

MAPQ: mapping quality

POS: 1-based position

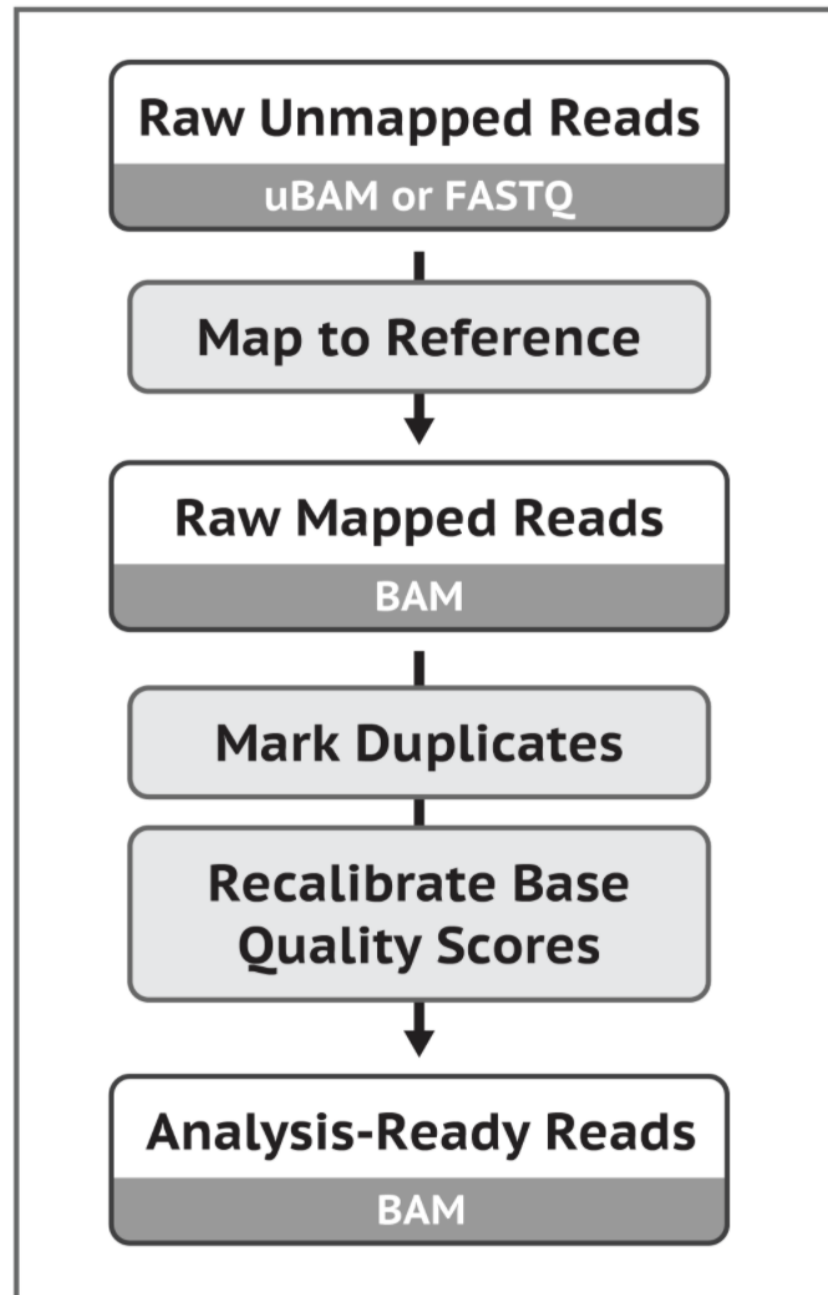RNAME: reference sequence name, e.g. chromosome/transcript id

FLAG: indicates alignment information about the read, e.g. paired, aligned, etc.

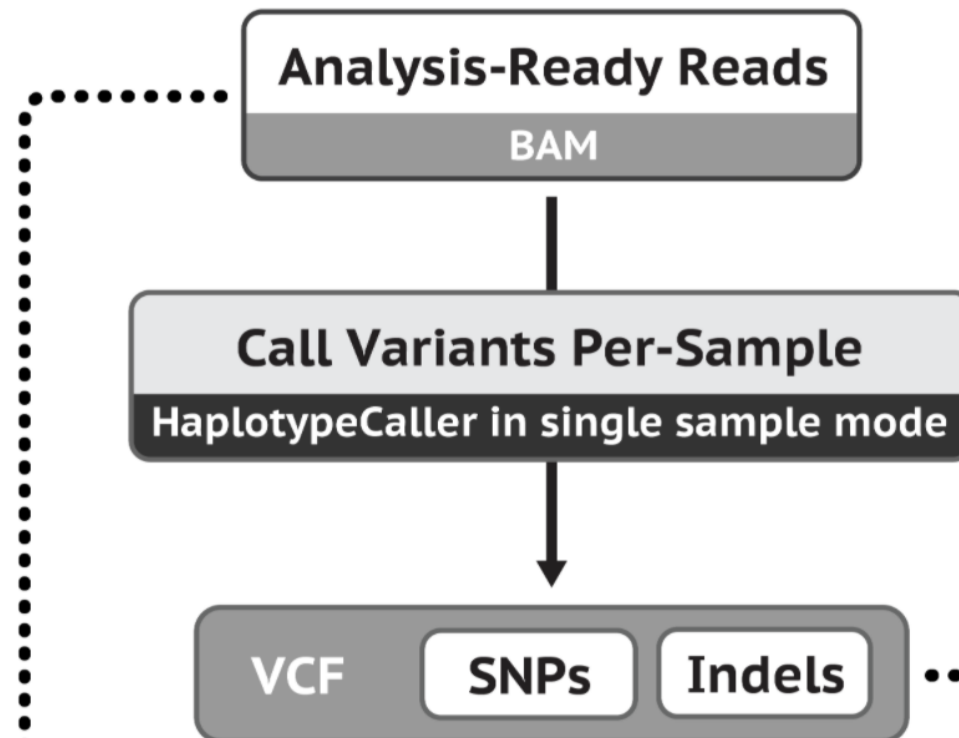QNAME: query template name, aka. read ID