

Alignment and Variant Calling

Introduction

Why use reference mapping as opposed to de novo assembly?

De novo assembly is computationally expensive and doesn't have the capacity for quality assessment.

Reference mapping is rapid, accurate and provides a mechanism for analysing sequence variation but relies on a reference genome for your organism.

Variants we can detect:

- single nucleotide polymorphisms (SNPs)
- small insertions/deletions (indels)
- structural variation (SVs)

Raw sequence data

Unique name for the read
(machine ID, flowcell, tile, barcode)

Quality scores for each base in the read

ACCCCCACAGTACTT
GTACTTATATACT ATTACACCTGAACCTA
GAACCTAAACCCC GAACCTAAACCCC ACCCCCCACAGTACTT
GTACTTATATACT CACCTGAACCTAAA
ATTACACCTGAACCTA CACCTGAACCTAAA GTACTTATATACT
ACCCCCACAGTACTT GAACCTAAACCCC ACCCCCCACAGTACTT

ACCCCCACAGTACTT

GTACTTATATACT ATTCACCTGAACCTA

GAACCTAAACCCC **GAACCTAAACCCC** **ACCCCCACAGTACTT**
GTACTTATATACT **CACCTGAACTAAA**

ATTCACCTGAACCTA CACCTGAACCTAAA GTACTTATATACT

ACCCCCACAGTACTT

GAACCTAAACCCC

ACCCCCACAGTACTT

GAACCTAAACCCCC ACCCCCCACAGTACTT

ATTCACCTGAAACCTAAACCCCCACAGTACTTATATACATAGTCATAATTACACTG

ACCCCACAGTACTT

GTACTTATATACT ATTACACCTGAAACCTA

GTACTTATATACT **GAACCTAAACCCC** CACCTGAAACCTAAA **GAACCTAAACCCC** ACCCCCACAGTACT

ATTCACCTGAACCTA CACCTGAACCTAAA GTACTTATATACTAC
ACCCCCACAGTACTT

GAACCTAAACCCCC **ACCCCCACAGTACTT**



ATTCACCTGAACCTAAACCCACAGTACTTATATACATAGTCATAATTACACTG
ATTCACCTGAACCT

TTCACCTGAACCTA
CACCTGAACCTAAA
GTTT

CTAAACCCCCACAGTA

AACCCCCACAGTACTTATA

CACAGTACTTATATA

CTTATATAACATAG

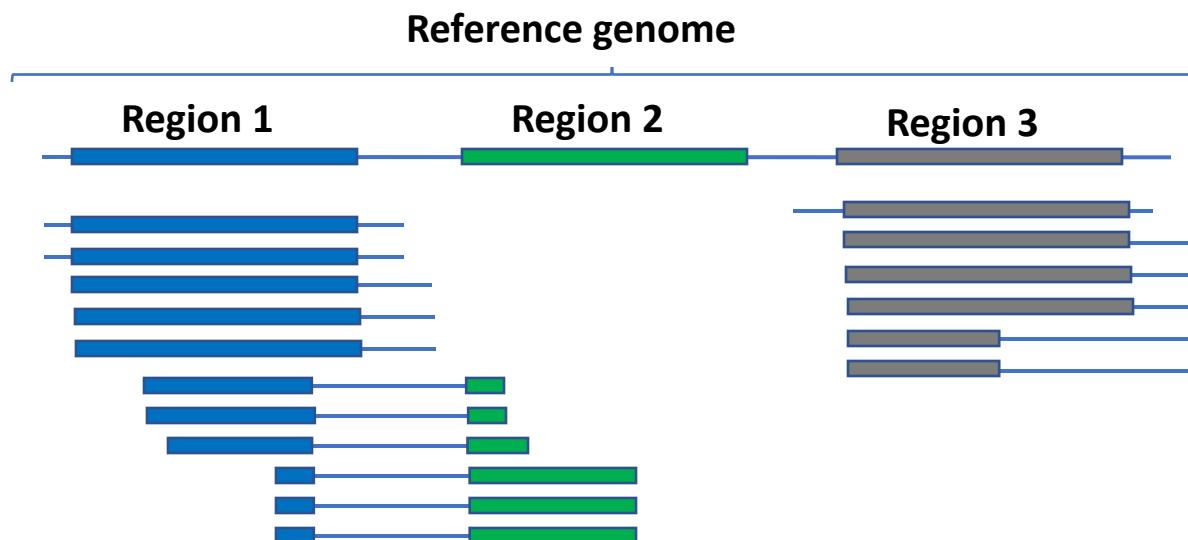
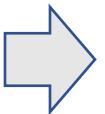
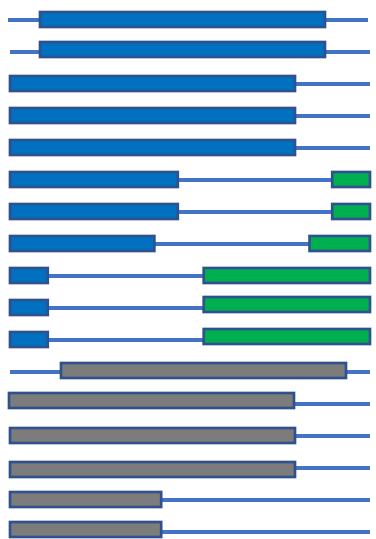
TATATACTAGTCA
-GATCTGCGA-

ACATAGTCATAAT
A G C G A T T

AGTCATAATTACAT

Mapping NGS reads

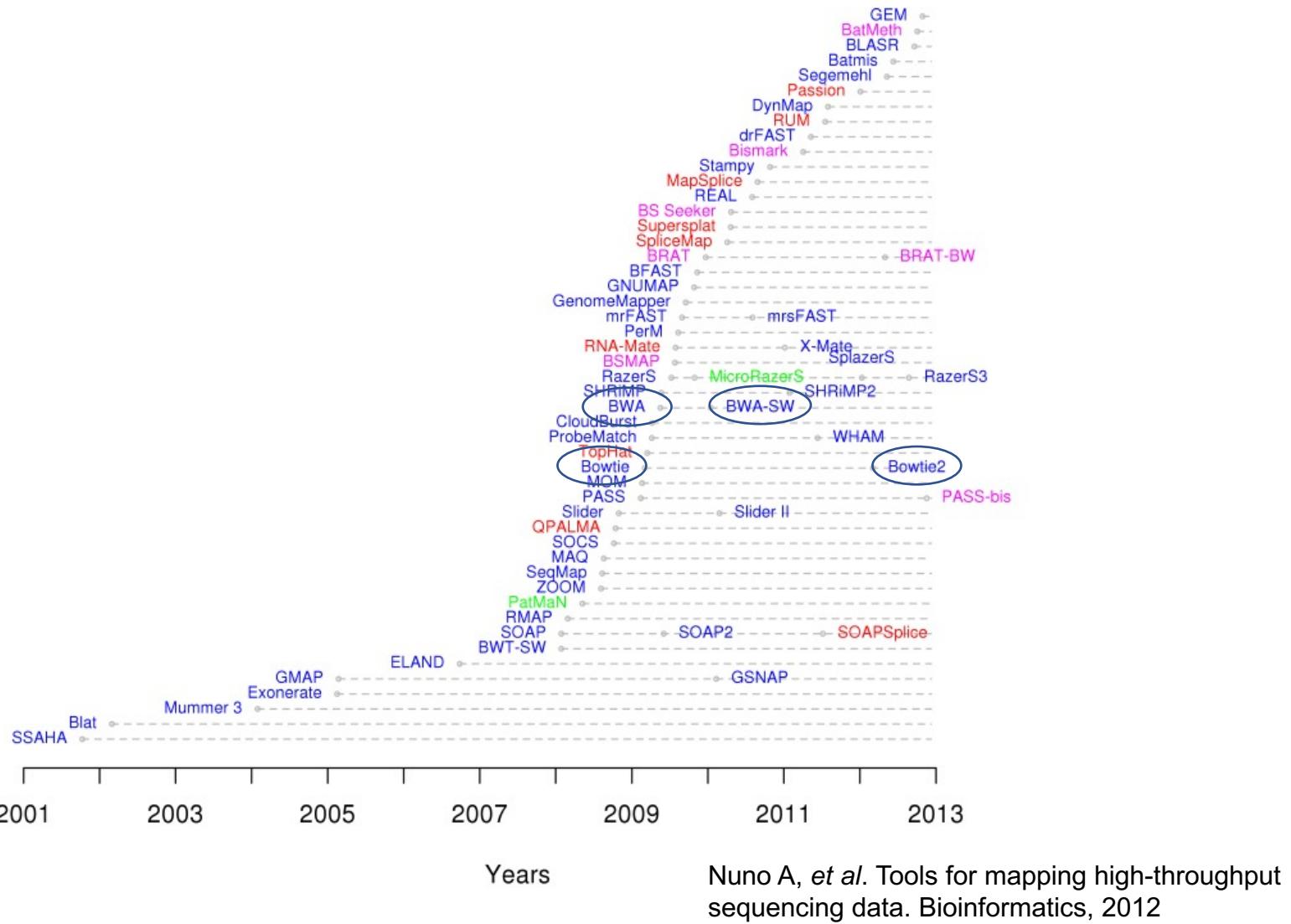
NGS generates a huge number of short reads



Alignment Algorithms

- Finding a match for your read in the genome can be slow
- Alignment algorithms make a special genome index for fast mapping
- Very fast for exact matches
- Slower for inexact matches with multiple possible locations
- Main alignment software you will use is BWA-MEM

Evolution of Mapping Algorithms

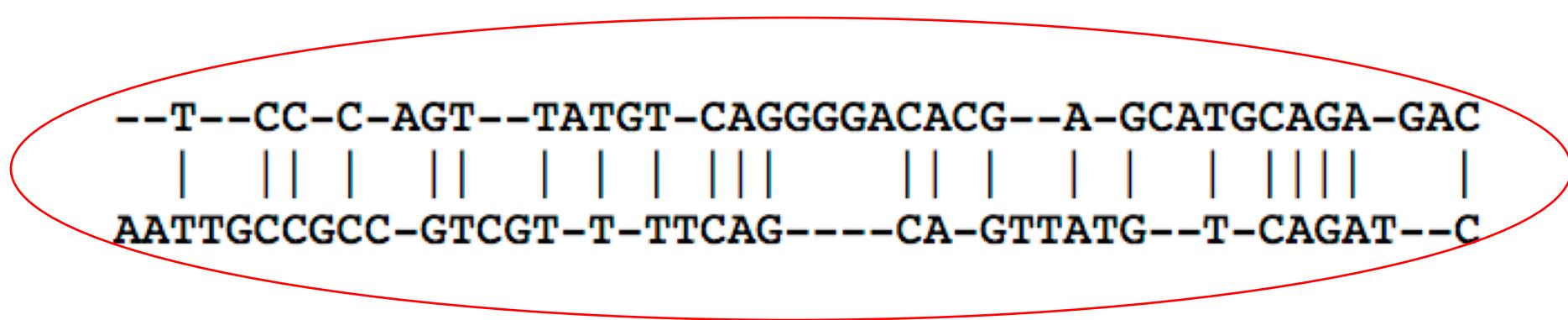


Mapping algorithms

- BWA-MEM & Bowtie2
- Perform gapped or local alignment of single or paired end reads
- Generally recommended for high quality queries as they are faster and more accurate than predecessors.

Global vs local alignment

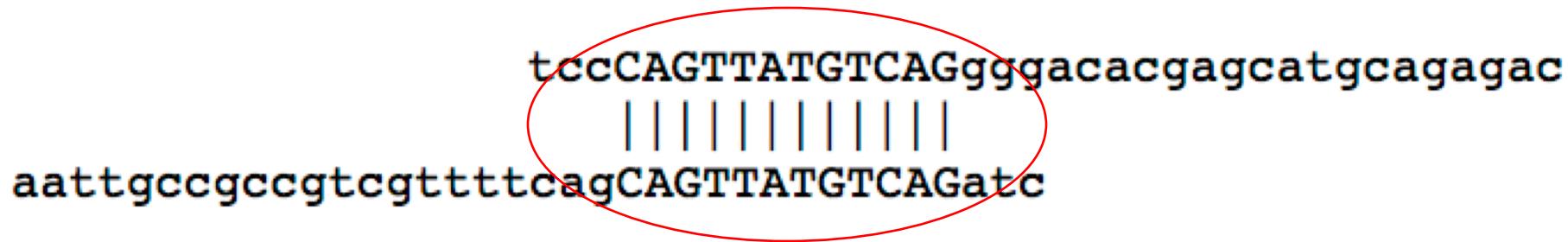
Global alignment:



A global alignment identifies the best match between 2 sequences from one end to the other. A good strategy for sequences of similar length and homology.

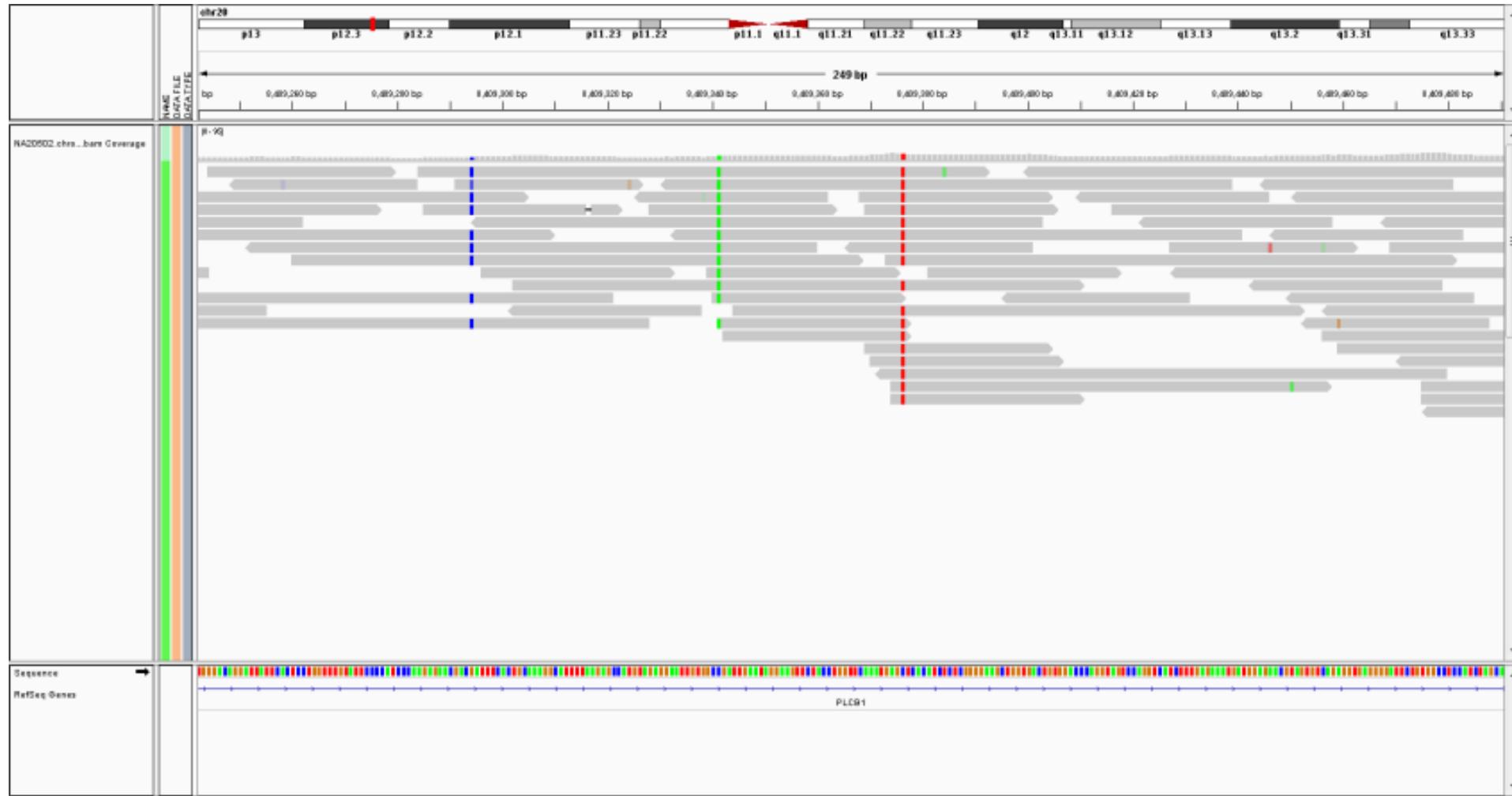
Global vs local alignment

Local alignment:

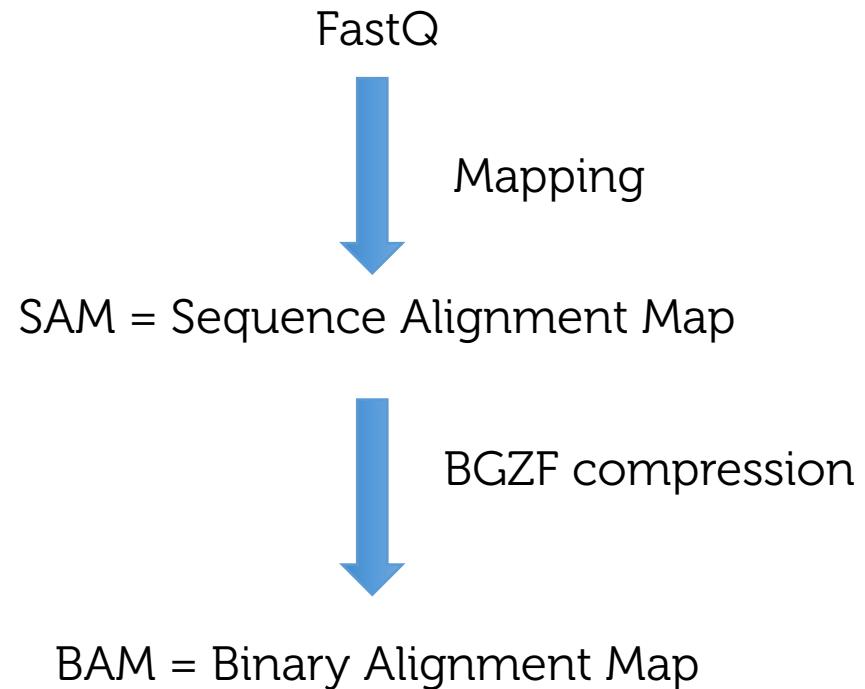


A local alignment identifies the best match between subsequences. A good strategy for short sequences with homology (e.g. Conserved motifs) but differ elsewhere. More sensitive for highly diverged sequences (e.g HIV)

Analysing data: alignment



WHAM, BAM?, SAM?



SAM Format

All short read aligners generate text based SAM output.

SAM files contain detailed information regarding the alignment results per read including mapping quality/Phred scores, mismatches/substitutions/indels (cigar line).

All this information is compressed into a BAM file and then indexed significantly reducing file size and allowing for rapid processing at the variant detection stage.

SAM Format

Unique name for the read
(machine ID, flowcell, tile, barcode)

M01481:15:00000000-A4YBL:1:1101:10709:6650

12594 355

ATTCCACTGTGTGCTTCAAATAAACTTCGCGTTGTAATTCCCTGACTTCACCGTCTGGAATCAGGTAGTCACATATCCCTCGCTTAAC

GGGACATTACAGTTATAAACAAATGTGGACAATGAAATTGTTAAGTCTT

ABBABFFFFFFGGGGGGGGHHHHHHGGGGGGHHHHHHHHHHHHHHHHHGHHGGFFHHHHFGGGHHHHHHGHHGGGHGHHHHGGGDGGGGHHGG

HGGHHGGHHHHHHHGGHHHHHHHEHHHHHHHGHHHHHHGHHHH

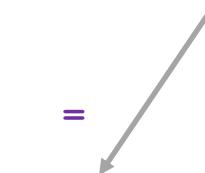
NM:i:0

AS:i:151

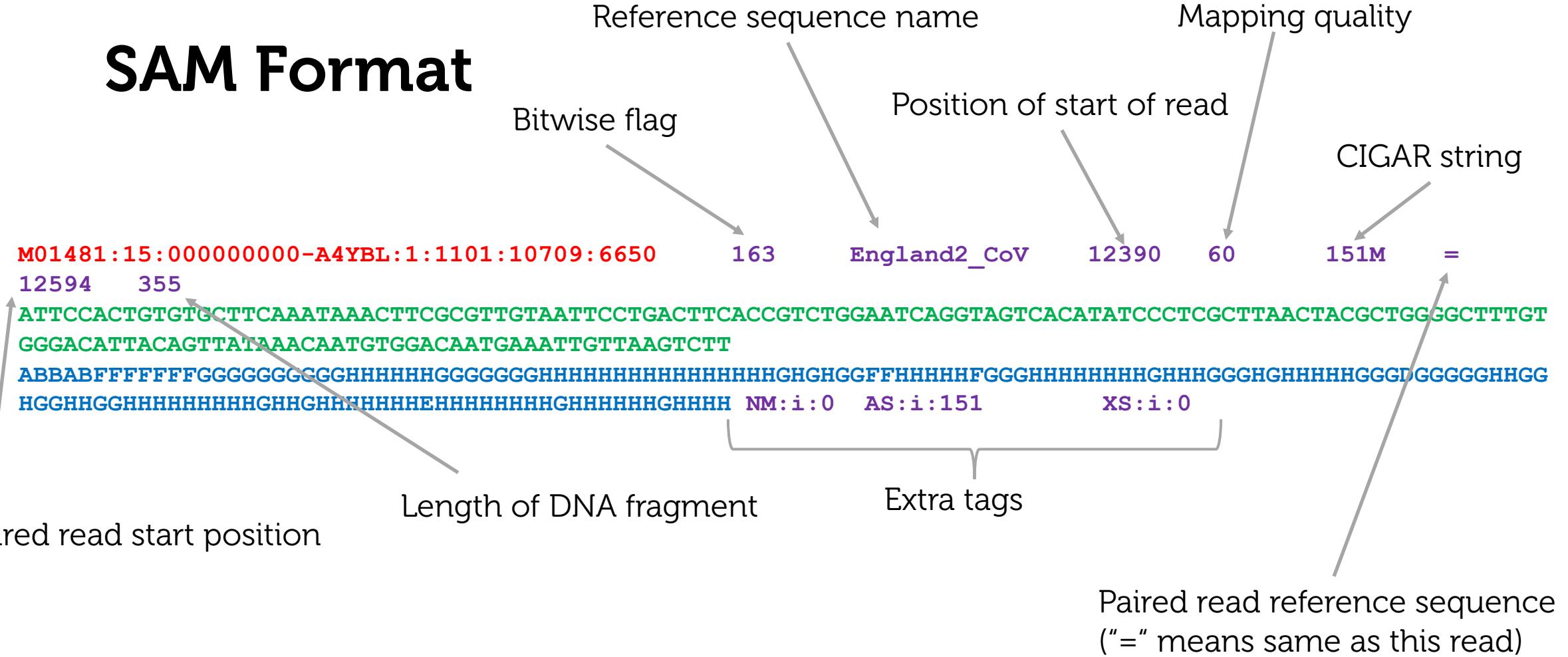
XS:i:0

Quality scores for each
base in the read

Sequence Read



SAM Format



Google 'SAM Format specification' or go
to <https://github.com/samtools/hts-specs>

Bitwise flag

numeric	binary	description
1	00000001	template has multiple fragments in sequencing
2	00000010	each fragment properly mapped according to aligner
4	00000100	fragment is unmapped
8	00001000	mate is unmapped
16	00010000	sequence is reverse complemented
32	00100000	sequence of mate is reversed
64	01000000	is first fragment in template
128	10000000	is second fragment in template

Decoding SAM flags

This utility makes it easy to identify what are the properties of a read based on its SAM flag value, or conversely, to find what the SAM Flag value would be for a given combination of properties.

To decode a given SAM flag value, just enter the number in the field below. The encoded properties will be listed under Summary below, to the right.

SAM Flag:

[Explain](#)

[Switch to mate](#)

Toggle first in pair / second in pair

Find SAM flag by property:

To find out what the SAM flag value would be for a given combination of properties, tick the boxes for those that you'd like to include. The flag value will be shown in the SAM Flag field above.

- read paired
- read mapped in proper pair
- read unmapped
- mate unmapped
- read reverse strand
- mate reverse strand
- first in pair
- second in pair
- not primary alignment
- read fails platform/vendor quality checks
- read is PCR or optical duplicate
- supplementary alignment

Summary:

read paired (0x1)
read mapped in proper pair (0x2)
mate reverse strand (0x20)
second in pair (0x80)

The CIGAR string

- Compress the read alignment into a short easy-to-parse format
- M – match
- I – insertion (base in the read, not the reference)
- D – deletion (base in the reference, not in the read)

Reference: CCCTACGTCCCAGTC-AC

CTACGTCCCAG	11M
TAC--CAC	3M2D3M
CCAGTCAAC	6M1I2M

The CIGAR line

```
M01481:15:00000000-A4YBL:1:1101:10709:6650      163      England2_Cov      12390    60  127M1D5I17M  =
12594    355
ATTCCACTGTGTGCTTCAAATAAACTTCGCGTTGTAATTCTGACTTCACCGTCTGGAATCAGGTAGTCACATATCCCTCGCTTAAC TACGCTGGGCTTG
TGGGACATTACAGTTATAACAATGTGGACAATGAAATTGTTAAGTCTT
ABBABFFFFFFGGGGGGGGHHHHHHGGGGGGHHHHHHHHHHHHHHHGHHGGFFHHHHHFGGGHHHHHHHGHHGGGHHHHHGGGDGGGGHHG
GHGGHHGGHHHHHHHHHGHHHHHHHEHHHHHHHGHHHHHHHGHHHH  NM:i:0  AS:i:151          XS:i:0
```

127 Matches, 1 Deletion, 5 Insertions, 17 Matches when compared to the reference sequence.

Error probabilities

M01481:15:000000000-A4YBL:1:1101:10709:6650 163 England2_CoV 12390 60 151M = 12594 355
ATCCACTGTGTGCTTCAAATAAACTCGCGTTGAATTCTGACTTCACCGTCTGGAATCAGGTAGTCACATATCCCTCGCTTAAC TACGCTGGGGCTT
TGTGGGACATTACAGTTATAACAAATGTGGACAATGAAATTGTTAAGTCTT
ABBABBBBBBGGGGGGGGHHHHHHGGGGGGHHHHHHHHHHHHHHHHHHHHGHGHGGFFHHHHHFGGGHHHHHHHGHHGGGHGHHHH
HGGGDGGGGGGHHGGHGGHHGGHHHHHHHGHHGGHHHHHHHEHHHHHHHHGHHHHHHHHGHHHHH NM:i:0 AS:i:151 XS:i:0

Phred quality scores are logarithmically linked to error probabilities

Phred quality score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%
50	1 in 100000	99.999%
60	1 in 1000000	99.9999%

Issues with alignment

- Repeat regions
 - Read may align equally well to multiple regions
 - Paired end reads have distance information which is also weighted
 - Alignment between the read and true source in the genome may have more differences than alignment with a repeat (read will be misplaced)
- Choice of reference genome
 - There may be many nucleotide differences between the reads and the most closely related reference genome (a significant problem in HIV and HCV viral NGS mapping)
 - Reads are easier to align with fewer variants
 - Leads to bias in alignment/variant calling
- Alignment of long reads with high error rate is difficult

Visualisation

- Many editors exist for visualising NGS data enabling you to view the read pileup.
- Tablet is a popular lightweight editor requiring the sorted BAM, the index file (BAI) and the reference genome in FASTA format.
- Variation can be graphically viewed across the assembled reads.

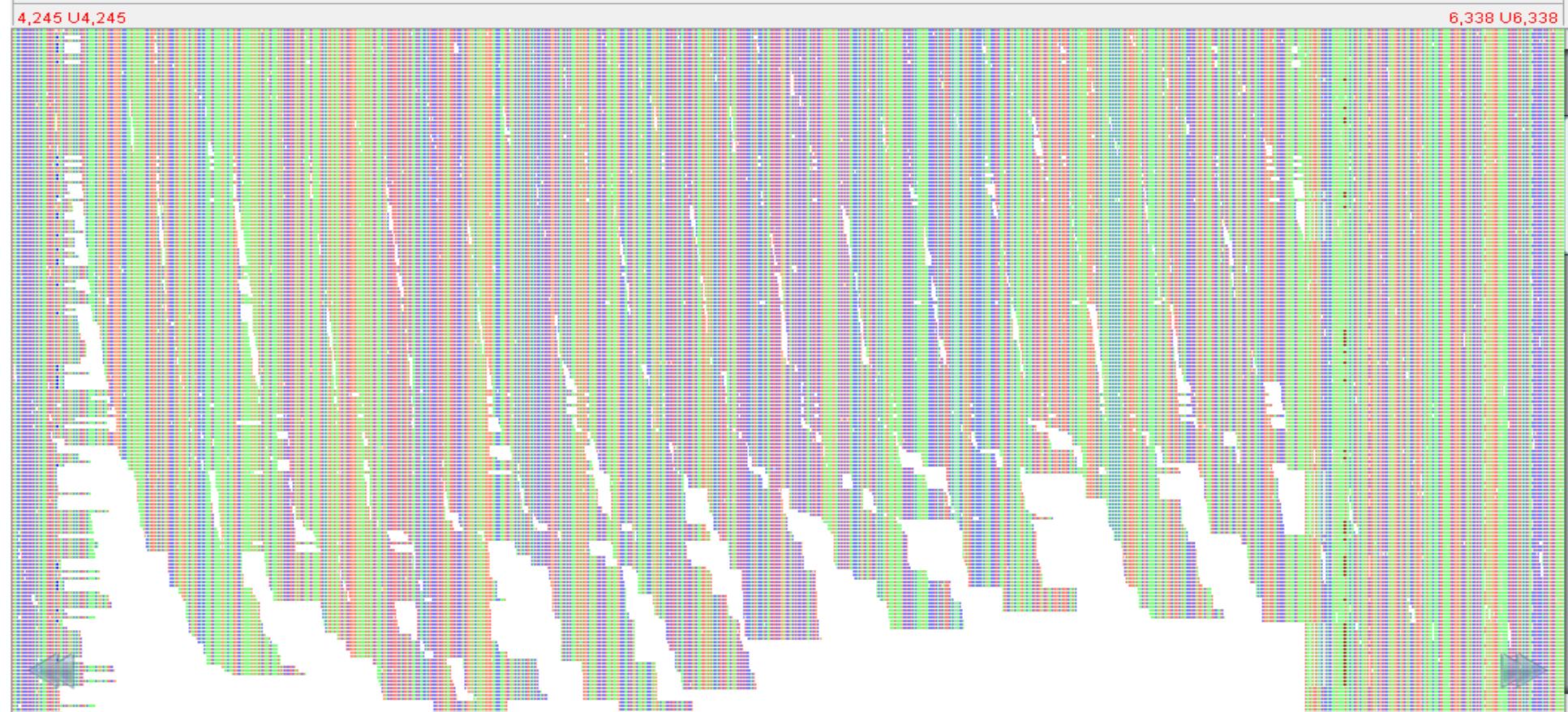


s Advanced

<input type="checkbox"/> Read Packing	<input type="checkbox"/> Zoom: <input type="range"/>	<input type="checkbox"/> Page Left	<input type="checkbox"/> Page Right	<input type="checkbox"/> Jump to Base
<input type="checkbox"/> Tag Variants	<input type="checkbox"/> Variants: <input type="range"/>	<input type="checkbox"/> Prev Feature	<input type="checkbox"/> Next Feature	<input type="checkbox"/> Read Info
<input type="checkbox"/> Read Colours		<input type="checkbox"/> Prev View	<input type="checkbox"/> Next View	<input type="checkbox"/> RS Off
Visual	Adjust		Navigate	<input type="checkbox"/> Show Cigar-I
				<input type="checkbox"/> Show Bases
				<input type="checkbox"/> RS Centre
				<input type="checkbox"/> Read Names
				<input type="checkbox"/> RS Custom
				Overlays



1 to 25,000 (25 Kb) 4,245 to 6,339 (2.1 Kb)



ACCCCCACAGTACTT

GTACTTATATA
ACCCCCACAGTACTT

GAACCTAAACCCC
GTACTTATATA
ATTCA
ACCCCCACAGTACTT

GAACCTAAACCCC
CACCTGAACCTAAA
CACCTGAACCTAAA
GTACTTATATA
GAACCTAAACCCC
ACCCCCACAGTACTT

↓

ATTCACCTGAACCTAAACCCCACAGTACTTTATACATAGTCATAATTACACTG
ATTCA
TTCACCTGAACCTA
CACCTGAACCTAAA
CTAAACCCCACAGTA
AACCCCCACAGTACTTATA
CACAGTACTTATATA
CTTATATA
TATATA
ACATAGTCATAAT
AGTCATAATTACA

SNP in our sample

ATTCACCTGAACCTAAACCCCACAGTACTTTATACATAGTCATAATTACACTG
ATTCACCTGAACCT
TTCACCTGAACCTA
CACCTGAACCTAAA
CTAAACCCACAGTA
AACCCCCACAGTACTT--A
CACAGTACTT--ATAC
CTT--ATACATAG
T--ATACATAGTCA
ACATAGTCATAAT
AGTCATAATTACAG

DELETION in our sample

ATTCACCTGAACCTAAACCCCACAGTACTTT--ATACATAGTCATAATTACACTG
ATTCACCTGAACCT
TTCACCTGAACCTA
CACCTGAACCTAAA
CTAAACCCACAGTA
AACCCCCACAGTACTTTCTA
CACAGTACTTTCTATAC
CTTTCAATACATAG
TTTCAATACATAGTCA
ACATAGTCATAAT
AGTCATAATTACAG

INSERTION in our sample

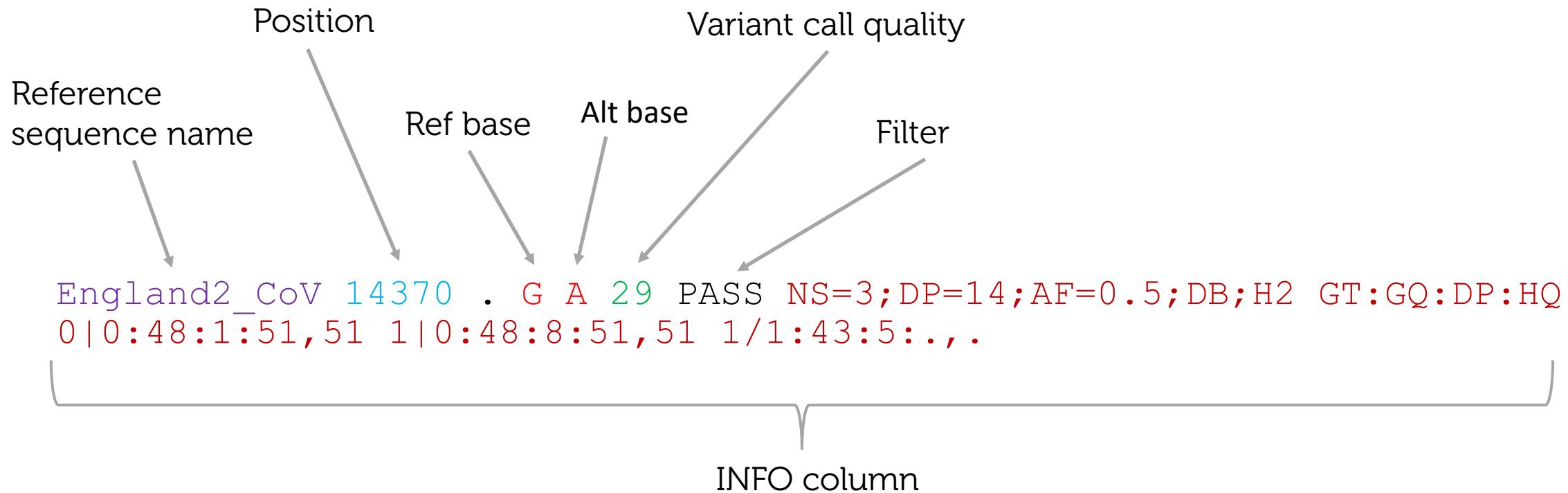
ATTCACCTGAACCTAAACCCCCACAGTACTTTATACATAGTCATAATTTACACTG
ATTCACCTGAACCT
TTCACCTGAACCTA
CACCCGAACCTAAA
CTAAACCCCCACAGTA
AACCCCCACAGTACTTTA
CACAGTACCTTATATAC
CTTATATACATAG
TTTATACATAGTCA
ACATAGTCATAAT
AGTCATAATTTACA

??? Is this a SNP?

Filtering variant calls

- How can we tell if the variant is real?
- Require a minimum **number** of reads with the variant
- Require a minimum **quality** of the reads with the variant
- Require a minimum **proportion** of reads with the variant

Variant Call Format



Variant Call Format (VCF) INFO column

AB=0; ABP=0; AC=1; AF=1; AN=1; AO=125; DP=125; DPB=125;
LEN=1; MQM=60; MQMR=0; NS=1; TYPE=snp

- 'NS' is number of samples
- 'DP' is depth of reads
- 'TYPE' is type of variant
- All columns are defined at the top of the file

Variant Call Format (VCF) FORMAT column

GT:DP:AD:RO:QR:AO:QA:GL

1:125:0,125:0:0:125:4434:-399.193,0

- 'GT' is genotype. 0 is ref, 1 is alt
- 'DP' is depth (125)
- 'GL' is genotype likelihood for ref and alt, highest is best

Variant Annotation

- What does this variant actually do?
- Requires a reference genome file with gene features

Position in nucleotide sequence Position in peptide sequence Effect
↓ ↓ ↓
54/1758 18/585 stop_gained c.54T>A p.Cys18* Saur_00043 mecR1 Methicillin resistance mecR1 protein

AA change

Gene locus tag
Gene name
Gene product

The diagram illustrates the mapping of variant annotations. It starts with three input fields: 'Position in nucleotide sequence' (54/1758), 'Position in peptide sequence' (18/585), and 'Effect' (stop_gained). An arrow points from each of these to a central summary row: 'c.54T>A p.Cys18*' (with a bracket over 'Cys18*'), 'Saur_00043', 'mecR1', and 'Methicillin resistance mecR1 protein'. From the 'Effect' field, another arrow points down to 'AA change' (c.54T>A p.Cys18*). Finally, arrows point from the summary row to three output fields: 'Gene locus tag' (Saur_00043), 'Gene name' (mecR1), and 'Gene product' (Methicillin resistance mecR1 protein).