



**Mansoura University**  
**Faculty of Computers and Information**  
**Department of Computer Science**  
**Second Semester: 2020-2021**



# **[MED-145] Genomics: Genome Indexing & Reads Mapping**

## **Introduction to SAM/BAM formats**

**Grade: Third Year (Medical Informatics Program)**

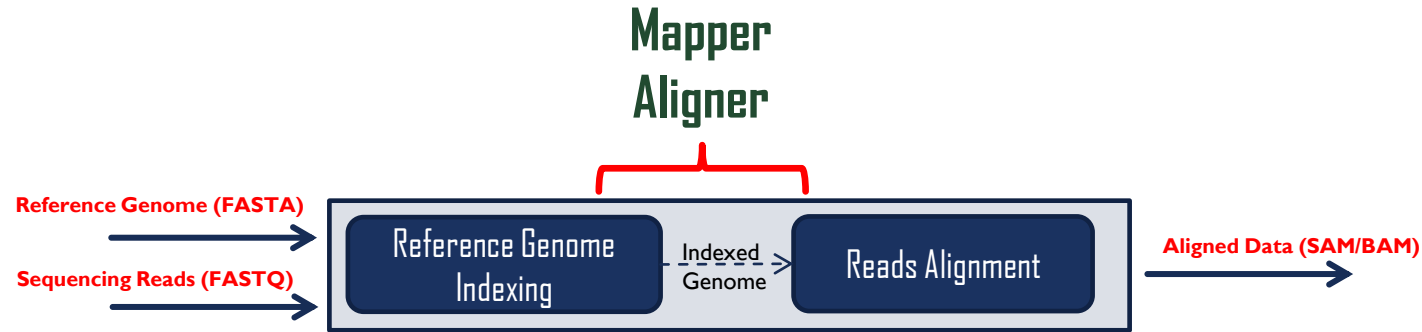
**Sara El-Metwally, Ph.D.**

**Faculty of Computers and Information,**

**Mansoura University,**

**Egypt.**

# TYPICAL MAPPING/ALIGNMENT WORKFLOW



## Genome Indexing and Mapping Approaches

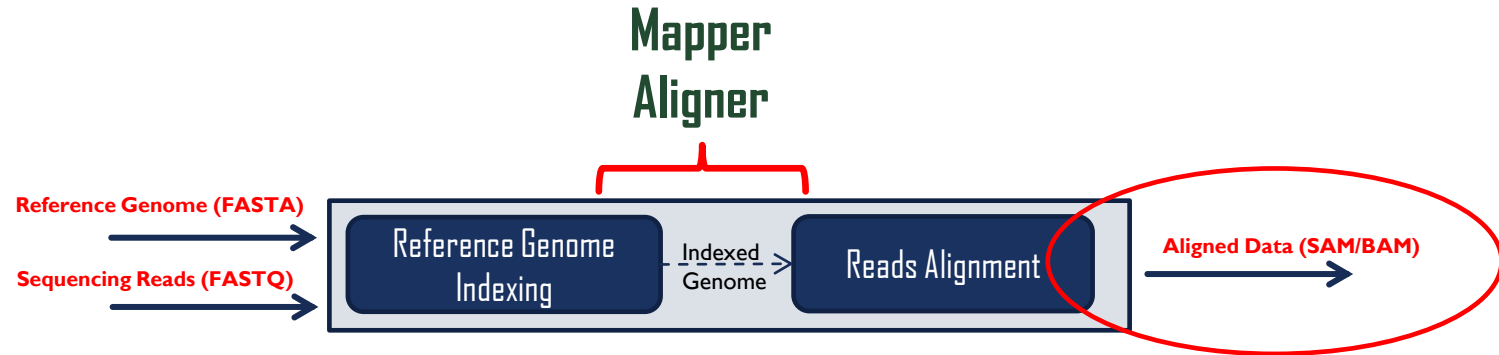


hash-based



Burrows-Wheeler

# TYPICAL MAPPING/ALIGNMENT WORKFLOW



## Genome Indexing and Mapping Approaches



hash-based

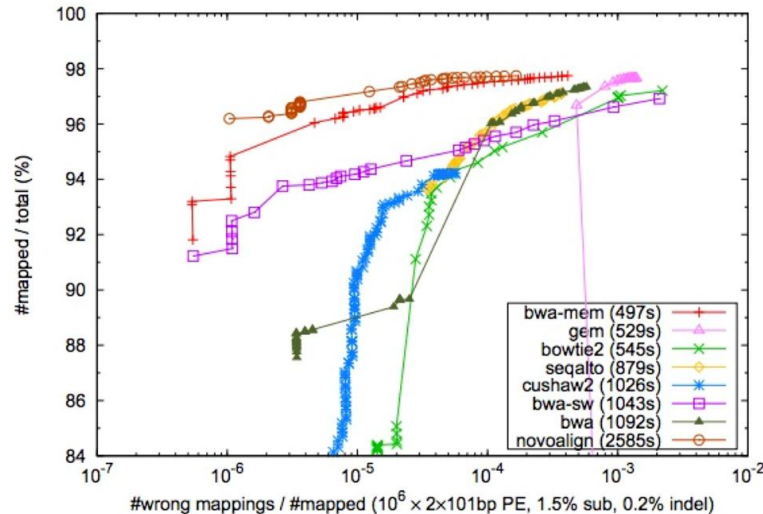


Burrows-Wheeler

# Sequence alignment software

<u>Aligner</u>	<u>Approach</u>	<u>Applications</u>	<u>Availability</u>
BWA-mem	Burrows-Wheeler	DNA, SE, PE, SV	open-source
Bowtie2	Burrows-Wheeler	DNA, SE, PE, SV	open-source
Novoalign	hash-based	DNA, SE, PE	free for academic use
TopHat	Burrows-Wheeler	RNA-seq	open-source
STAR	hash-based (reads)	RNA-seq	open-source
GSNAP	hash-based (reads)	RNA-seq	open-source

# BWA-MEM: never "published" ; widely used.



**Fig. 1.** Percent mapped reads as a function of the false alignment rate under different mapping quality cutoff. Alignments with mapping quality 3 or lower are excluded. An alignment is *wrong* if after correcting clipping, its start position is within 20bp from the simulated position.  $10^6$  pairs of 101bp reads are simulated from the human reference genome using wgsim (<http://bit.ly/wgsim2>) with 1.5% substitution errors and 0.2% indel variants. The insert size follows a normal distribution  $N(500, 50^2)$ . The reads are aligned back to the genome either as single end (SE; top panel) or as paired end (PE; bottom panel). GEM is configured to allow up to 5 gaps and to output suboptimal alignments (option '-e5 -m5 -s1' for SE and '-e5 -m5 -s1 -pb' for PE). GEM does not compute mapping quality. Its mapping quality is estimated with a BWA-like algorithm with suboptimal alignments available. Other mappers are run with the default setting except for specifying the insert size distribution. The run time in seconds on a single CPU core is shown in the parentheses.

## Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM


Heng Li

Broad Institute of Harvard and MIT, 7 Cambridge Center, Cambridge, MA 02142, USA

<https://arxiv.org/pdf/1303.3997v2.pdf>

# BWA-MEM

Unaligned  
Sample Data  
In FASTQ (SE or PE)



```
@seq1
ATTCGAAACA...
+
DDED88(999...
@seq2
CCCCGTTTCA...
+
AAC887BBAC...
```

Reference genome (FASTA)

```
>chr1
TACCTCCAGGGGGCATCCTCCCCCAATTCTG
AAACACAATCGTAGCCCCTGGCACTACCTATG
TGTGTCAATTTCGGAGAGAGAGATTACAGAA
AAAAAAGTCTGGACTCAACTAGGATACACACA
TTCGGCTACAGATACCAAAAAAAAAAAAAAAAA
AAATTTTCACCATTGAGGCACCACCTTCTCGT
CGCTGCGTCGCTCTGCTCGCTTCGGCTAAAAA
TTCGCGCAATACATTTCGGCTACAGATACCAAA
```

↓

BWA MEM

Aligned  
Sample Data in  
SAM format

```
seq1      99      1
3666901    60
149M      =
3666935    185
ATTCGAAACA...DDED88(999
MC:Z:151M  MD:Z:149
RG:Z:15-0017315_1
NM:i:0      MQ:i:60
AS:i:149    XS:i:44
147         1
3666935    60
151M      =
3666901    -185
CCCCGTTTCA...AAC887BBAC...
MC:Z:149M  MD:Z:151
RG:Z:15-0017315_1
NM:i:0      MQ:i:60
AS:i:151    XS:i:59

seq2
```

# BWA-MEM workflow

*This takes a long time, but  
you do it once*

Create BWT of reference genome.

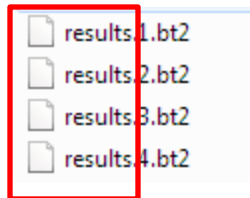
```
$ bwa index grch38.fa
```



*Output is in SAM format.  
Use multiple threads if you  
have a computer with  
multiple CPUs.*

Align paired-end FASTQ  
to BWT index.

```
$ bwa mem -t 16 grch38.fa 1.fq 2.fq > sample.sam
```



**Same prefix: grch38.fa**

# SAM format: a text-based standard(!) for representing sequence alignments

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 25 no. 16 2009, pages 2078–2079  
doi:10.1093/bioinformatics/btp352

*Sequence analysis*

## The Sequence Alignment/Map format and SAMtools

Heng Li<sup>1,†</sup>, Bob Handsaker<sup>2,†</sup>, Alec Wysoker<sup>2</sup>, Tim Fennell<sup>2</sup>, Jue Ruan<sup>3</sup>, Nils Homer<sup>4</sup>, Gabor Marth<sup>5</sup>, Goncalo Abecasis<sup>6</sup>, Richard Durbin<sup>1,\*</sup> and 1000 Genome Project Data Processing Subgroup<sup>7</sup>

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK, <sup>2</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02141, USA, <sup>3</sup>Beijing Institute of Genomics, Chinese Academy of Science, Beijing 100029, China, <sup>4</sup>Department of Computer Science, University of California Los Angeles, Los Angeles, CA 90095, <sup>5</sup>Department of Biology, Boston College, Chestnut Hill, MA 02467, <sup>6</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA and <sup>7</sup><http://1000genomes.org>

Received on April 28, 2009; revised on May 28, 2009; accepted on May 30, 2009

Advance Access publication June 8, 2009

Associate Editor: Alfonso Valencia

**Table 1.** Mandatory fields in the SAM format

No.	Name	Description
1	QNAME	Query NAME of the read or the read pair
2	FLAG	Bitwise FLAG (pairing, strand, mate strand, etc.)
3	RNAME	Reference sequence NAME
4	POS	1-Based leftmost POSition of clipped alignment
5	MAPQ	MAPping Quality (Phred-scaled)
6	CIGAR	Extended CIGAR string (operations: MIDNSHP)
7	MRNM	Mate Reference NaMe ('=' if same as RNAME)
8	MPOS	1-Based leftmost Mate POSition
9	ISIZE	Inferred Insert SIZE
10	SEQ	Query SEQUENCE on the same strand as the reference
11	QUAL	Query QUALity (ASCII-33=Phred base quality)



# SAM format overview

- In the dark ages, sequence aligners used disparate output formats. **Pain.**
- 1000 Genomes Project sought to standardize. **Standards are good.**
- The result is imperfect, but it's a **huge** improvement.
- **Strengths of the SAM and BAM formats**
  - Compressed: less disk hungry
  - Indexed: fast viewing, slicing, etc.
  - Single-end and paired-end
  - Relatively simple to produce
  - *Good toolkits available*

# SAM/BAM/CRAM

- **SAM**: Sequence Alignment Map format is a tab delimited text file describing mapping information for short read data.
- With the millions of reads generated from NGS machines, the SAM files become really big.
- **BAM** files contain the same information as SAM files but are encoded in condensed computer readable binary format to save disk space.
- Use [samtools](#) to convert **SAM** to **BAM**.

# SAM/BAM

- There are two types of BAM files: **unsorted or sorted**.
- **Sorted BAM file**: the alignments are sorted left-to-right along the reference genome.
- Most alignment tools output SAM (not BAM), and the alignments come out in an arbitrary order -- not sorted.
- An authoritative and complete document describing the SAM and BAM formats is the SAM specification.

# SAM/BAM

The header contains general information about the alignment, such as the reference name, reference length, the program used for the alignment, etc.

```
@HD VN:1.0 SO:coordinate
@SQ SN:chr20 LN:64444167
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping_tophat/L6 18 GTGAAA L007 R1 001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
CCGTGTTTAAAGGTGGATGCGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C BBDCCDDCCDDDDDDDDDDDDDDCCDCDDDDDDDDDDDDDDDDDDDDDDDBHFFFFDC@@
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0
TGCTGGATCATCTGGTTAGTGGCTTCTGACTCAGAGGACCTTCGTCCCCTGGGGCAGTGGACCTTCCAGTGATTCCCCTGACATAAGGGGCATGGACGA
G DCDDDDDEDDDDDDDDDDDDDDCCDDDDDDDEEC>DFFFEJJJJJIGJJJJIHGBHHGJIJJJJJJGJJJJJJJJJJHJJJJJJHHHHHHFFFFFCCC
AS:i:-16 XM:i:3 X0:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0
GGCTTTATTGGTAAAAAAGGAATAGCAGATTTAATCAGAAATTCACCTGGCCCAGCAGCACCAACCAGAAAGAAGGGAAGAAGACAGGAAAAAACCA
C DDDDDDDDDCCDDDDDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJJIJJJJIIIGGFJJJIHIIIIJJJJJJJIGHHFAHGFIHJHFGGHFFFDDBB
AS:i:-11 XM:i:2 X0:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0
0 GTGGCTCTTCACAGGAATGTTGAGGATGACATCCATGTCTGGGGTGCATTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
```

# Sequence Alignment/Map Format Specification

The SAM/BAM Format Specification Working Group

7 Jan 2021

The master version of this document can be found at <https://github.com/samtools/hts-specs>.  
This printing is version 981fe0f from that repository, last modified on the date shown above.

@HD

## 1 The SAM Format Specification

SAM stands for Sequence Alignment/Map format. It is a TAB-delimited text format consisting of a header section, which is optional, and an alignment section. If present, the header must be prior to the alignments. Header lines start with '@', while alignment lines do not. Each alignment line has 11 mandatory fields for essential alignment information such as mapping position, and variable number of optional fields for flexible or aligner specific information.

This specification is for version 1.6 of the SAM and BAM formats. Each SAM and BAM file may optionally specify the version being used via the @HD VN tag. For full version history see Appendix B.

Unless explicitly specified elsewhere, all fields are encoded using 7-bit US-ASCII<sup>1</sup> in using the POSIX / C locale. Regular expressions listed use the POSIX / IEEE Std 1003.1 extended syntax.

0 - 9										10-19										20-29										30-39											
1	2	3	4	5	6	7	8	9	0	1	2	3	4			5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0
A	G	C	A	T	G	T	T	A	G	A	T	A	A	*	*	G	A	T	A	G	C	T	G	T	G	C	T	A	G	G	C	A	G	T	C	A	G	C	G	C	C

Reference Genome

R1

A G C T G T G C T

R3

T G T G C T A G G

R2

G T T A G A T A A

R4

A A C C G A T A G

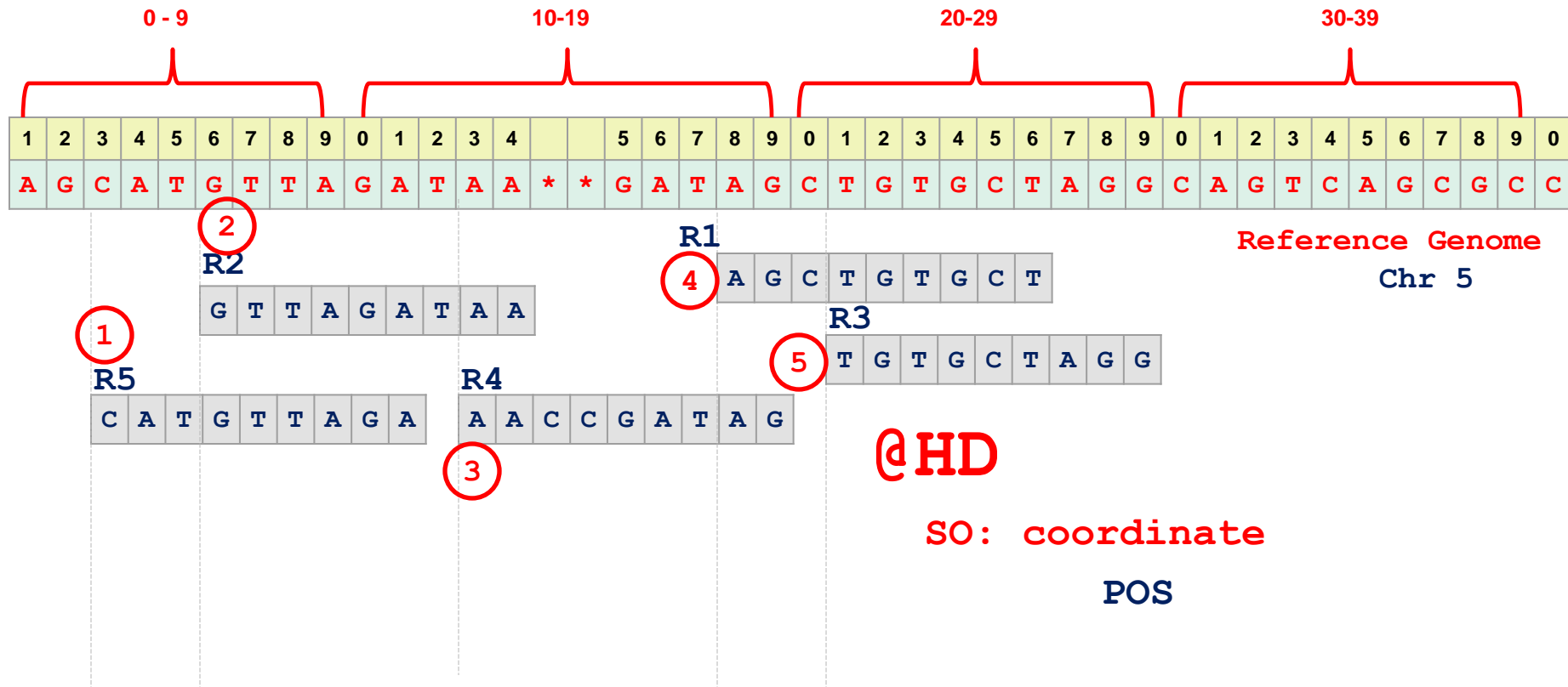
R5

C A T G T T A G A

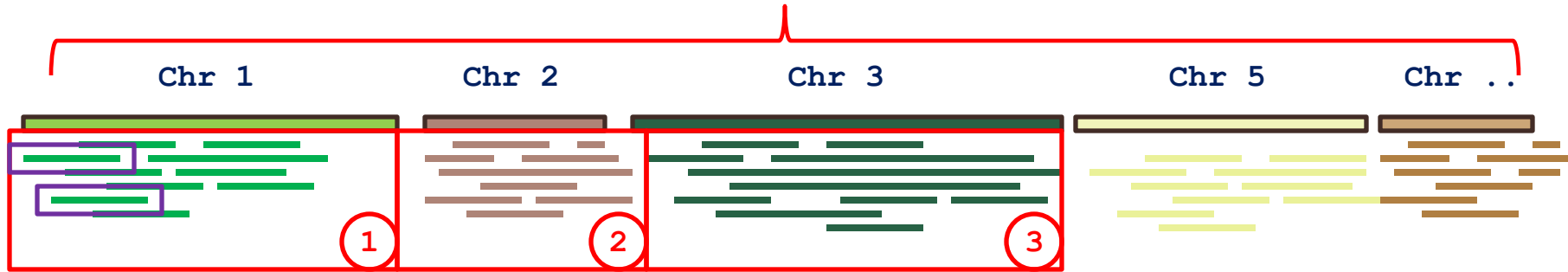
@HD

SO

- unknown
- unsorted
- queryname
- coordinate



# Reference Genome



@HD

SO: coordinate

RNAME

POS



# BWA-MEM

**@SQ**

Reference genome (FASTA)

```
>chr1 SN
TACCTCCAGGGGGCATCCTCCCCCAATTCG
AAACACAATCGTAGCCCCTGGCACTACCTATG
TGTGTCAATTCGGAGAGAGAGATTACAGAA
AAAAAAGTCTGGACTCAACTAGGATACACACA
TTCGGCTACAGATACCAAAAAAAAAAAAAA
AAATTTTCACCATTGAGGCACCACCTTCTCGT
CGCTGCGTCGCTCTGCTCGCTTCGGCTAAAAA
TTCGCGCAATACATTTCGGCTACAGATACCAA
```

**LN**

Aligned  
Sample Data in  
SAM format

Unaligned  
Sample Data  
In FASTQ (SE or PE)



```
@seq1
ATTCGAAACA...
+
DDED88(999...
```

**@PG** ↓  
**BWA MEM**  
**ID, VN**

```
$ bwa mem -t 16 grch38.fa 1.fq 2.fq > sample.sam
```

**CL**

```
seq1      99      1
3666901   60
149M      =
3666935   185
ATTCGAAACA...DDED88(999
MC:Z:151M MD:Z:149
RG:Z:15-0017315_1
NM:i:0     MQ:i:60
AS:i:149   XS:i:44
```

# SAM/BAM

```
@HD      VN:1.0  SO:coordinate
@SQ      SN:chr20      LN:64444167
@PG      ID:TopHat      VN:2.0.14      CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping tophat/L6 18 GTGAAA L007 R1 001.fastq
```

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDDCCDDCCDDDDDDDDDDDDDDCCCDDBC?DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
      AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961      16      chr20      193953      50      100M      *      0      0
      TGCTGGATCATCTGGTTAGTGGCTTCTGACTCAGAGGACCTTCGTCCCCTGGGGCAGTGGACCTTCCAGTGATTCCCCTGACATAAGGGGCATGGACGA
G      DCDDDDDEDDDDDDDDDDDDDDCCDDDDDDDEEC>DFFFEJJJJJIGJJJJIHGBHHGJIJJJJJJGJJJJJJHJJJJJJHHHHHHFFFFFCCC
      AS:i:-16      XM:i:3      XO:i:0      XG:i:0      MD:Z:60G16T18T3      NM:i:3      NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030      16      chr20      270877      50      100M      *      0      0
      GGCTTTATTGGTAAAAAAGGAATAGCAGATTTAATCAGAAATTTCCACCTGGCCCAGCAGCACCAACCAGAAAAGAAGGGAAGAAGACAGGAAAAAACCA
C      DDDDDDDDDCCDDDDDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJJIJJJJIIIGGFJJJIHIIIIJJJJJJIGHHFAHGFIHJHFGGHFFFD@BB
      AS:i:-11      XM:i:2      XO:i:0      XG:i:0      MD:Z:0A85G13      NM:i:2      NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699      0      chr20      271218      50      50M4700N50M      *      0
      0      GTGGCTCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGGTGCATTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
```

Every line after the header represents a single read, with 11 mandatory tab separated fields of information.

# What critical information do we need for sequence alignments?

# SAM format overview

Col #	Name	Meaning	Example
1	QNAME	Read or Pair name	HWI-ST156_1:278:1:1058:4544:0
2	FLAG	Bitwise FLAG	<i>Much more soon!</i>
3	RNAME	Reference sequence name	chr1
4	POS	1-based alignment start coordinate	8,724,005
5	MAPQ	Mapping quality	60
6	CIGAR	Extended CIGAR string	<i>Much more soon!</i>
7	MRNM	If paired, the mate's reference seq.	chr1
8	MPOS	If paired, the mate's alignment start	8,724,505
9	ISIZE	If paired, the insert size	562
10	SEQ	The sequence of the query/mate	ACAAATTCAG...
11	QUAL	The quality string for the query/mate	HHH\$^^%\$\$\$...
12	OPT	Optional Tags	XA:i:2, MD:Z:OT34G15

<http://samtools.sourceforge.net/samtools.shtml>

Col #	Name	Meaning	Example
1	QNAME	Read or Pair name	HWI-ST156_1:278:1:1058:4544:0
2	FLAG	Bitwise FLAG	<i>Much more soon!</i>
3	RNAME	Reference sequence name	chr1
4	POS	1-based alignment start coordinate	8,724,005
5	MAPQ	Mapping quality	60
6	CIGAR	Extended CIGAR string	<i>Much more soon!</i>
7	MRNM	If paired, the mate's reference seq.	chr1
8	MPOS	If paired, the mate's alignment start	8,724,505
9	ISIZE	If paired, the insert size	562
10	SEQ	The sequence of the query/mate	ACAAATTCAG...
11	QUAL	The quality string for the query/mate	HHH\$^~%\$\$\$...
12	OPT	Optional Tags	XA:i:2,MD:Z:OT34G15

```

arq5x@beast:~/beast_data/Pat — ssh — 101x8
arq5x@beas.../Pat — bash  arq5x@beas.../Pat — bash  arq5x@beas...a/Pat — ssh  java
[arq5x@beast Pat]$ samtools view 1094PC0005.possrt.conc.bam | head -1
1 HWI-ST156_1:278:66:2461:8880:0 2 99 3 chr10 4 50133 0 51M = 50214 132 TAATT
GACGCGCTGTTACGCCCTTTGAGTTCGGTTGAGTTTTGGTTGGAG GFDGDFFFDBEEEE:DDDDDFDFEFGGEEGFAFDGFFB?BBABBB@EBDA? X
T:A:R NM:i:0 SM:i:0 AM:i:0 X0:i:2 X1:i:0 XM:i:0 X0:i:0 XG:i:0 MD:Z:S1 XA:Z:chr10,+50133,51M
,0;

```



# SAM/BAM

qname: sequence name.

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCDDCCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
      AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

Field 1, qname is the name of the read. Read names often contain information about:

- 1.The scientific study for which the read was sequenced.
- 2.The sequencing instrument, and the exact part of the sequencing instrument, where the DNA was sequenced.

flag is a bit field encoding some yes/no pieces of information about whether and how the read aligned

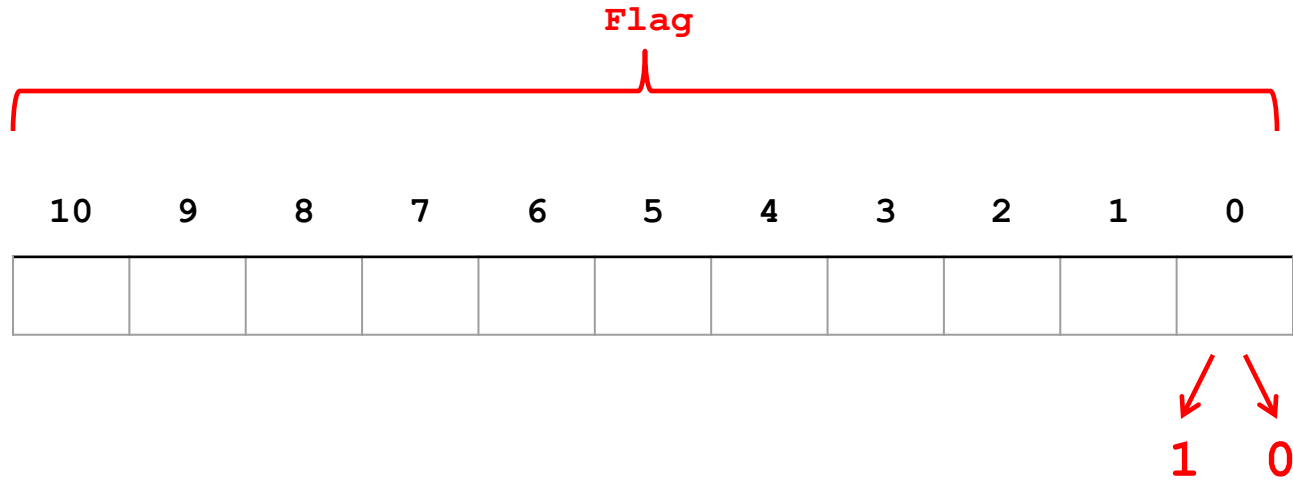
# SAM/BAM

flag

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
CCGTGTTTAAAGGTGGATGCGGTACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C BBDCCDDCCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCDDDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
```

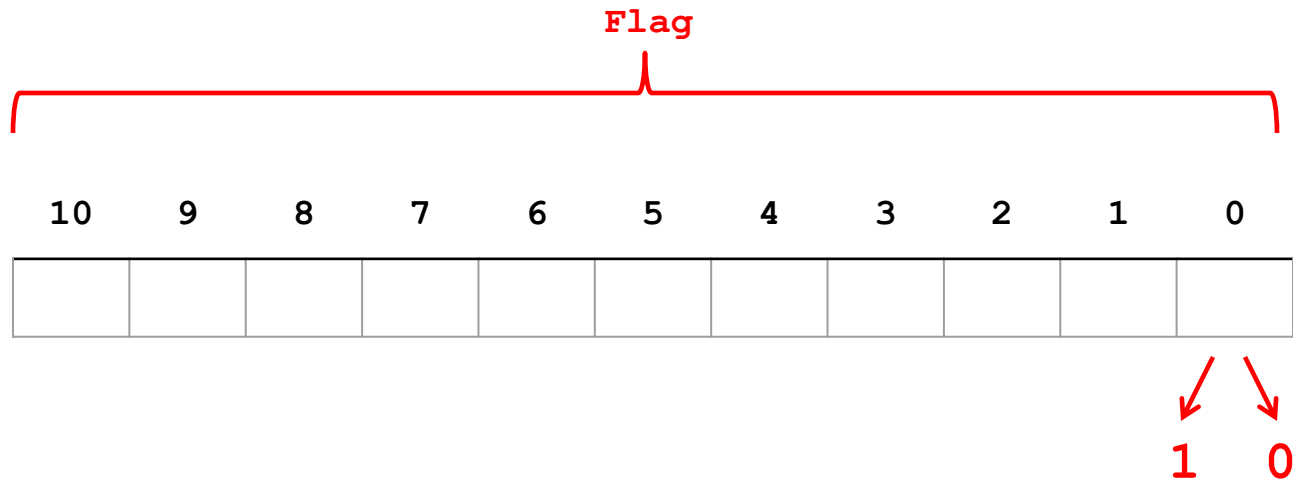
base2	base10	base16	Meaning	Applies to:
00000000001	1	0x0001	The read originated from a paired sequencing molecule	Both
00000000010	2	0x0002	The read is mapped in a proper pair	Pairs only
00000000100	4	0x0004	The query sequence itself is unmapped	Both
00000001000	8	0x0008	The query's mate is unmapped	Pairs only
00000010000	16	0x0010	Strand of the query (0 for forward; 1 for reverse strand)	Both
00000100000	32	0x0020	Strand of the query's mate	Pairs only
00001000000	64	0x0040	The query is the first read in the pair	Pairs only
00010000000	128	0x0080	The read is the second read in the pair	Pairs only
00100000000	256	0x0100	The alignment is not primary	Both
01000000000	512	0x0200	The read fails platform/vendor quality checks	Both
10000000000	1024	0x0400	The read is either a PCR duplicate or an optical duplicate	Both





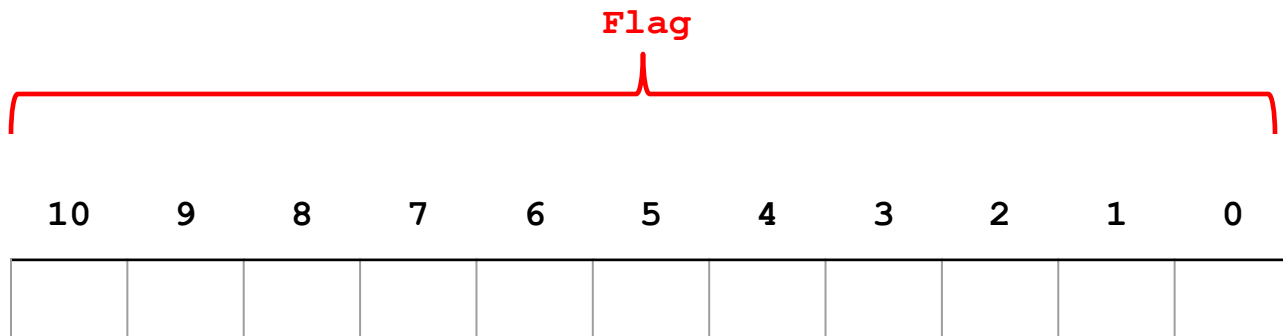
Single-end sequencing read





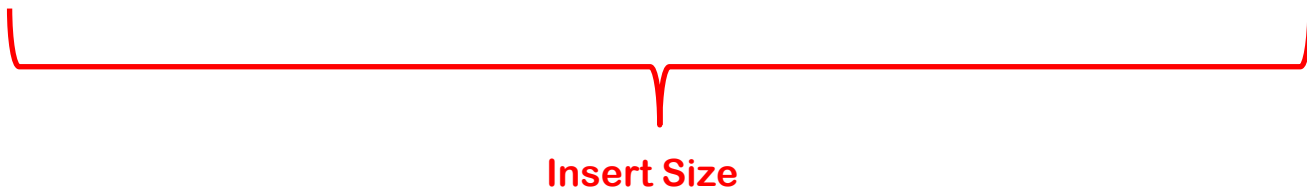
paired-end sequencing read

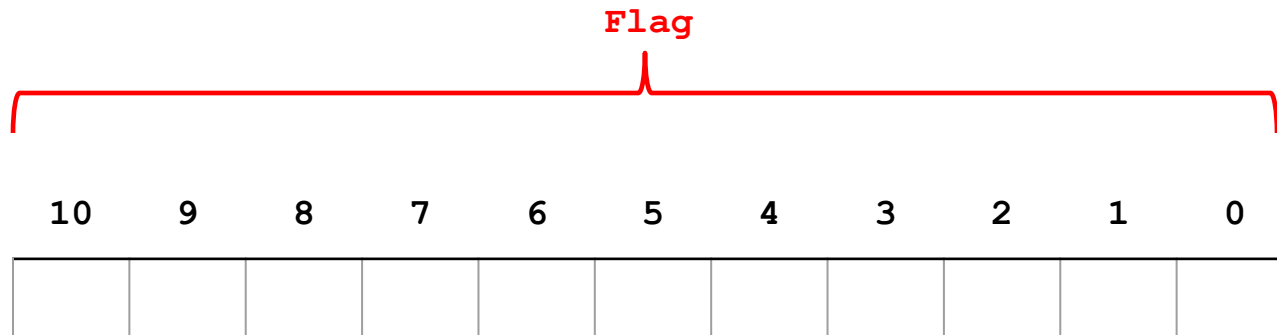




Read mapped in a proper pair

1 0

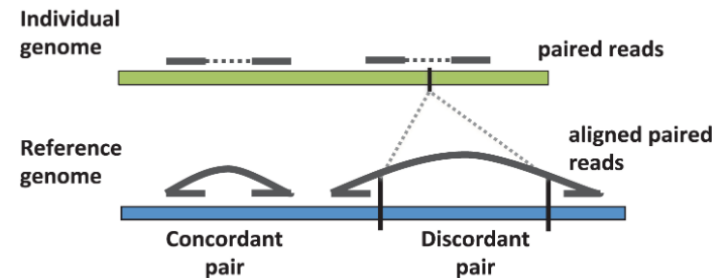




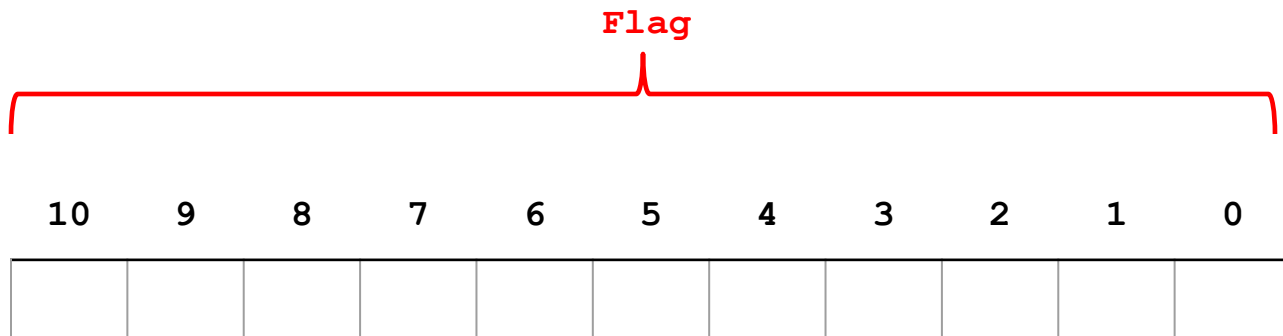
Read mapped concordantly

AGGG**T**TTGGTTCGTTTTAGGGTTTGGTTCGTTGA**G**GCTTTAG  
 ||| ||||| ||||| |||||  
 AGGG**A**TTGGTTCGTT \_\_\_\_\_ GGTTTCGTTGA**A**GCTT

1 0



doi: <https://doi.org/10.1371/journal.pcbi.1002821.g005>

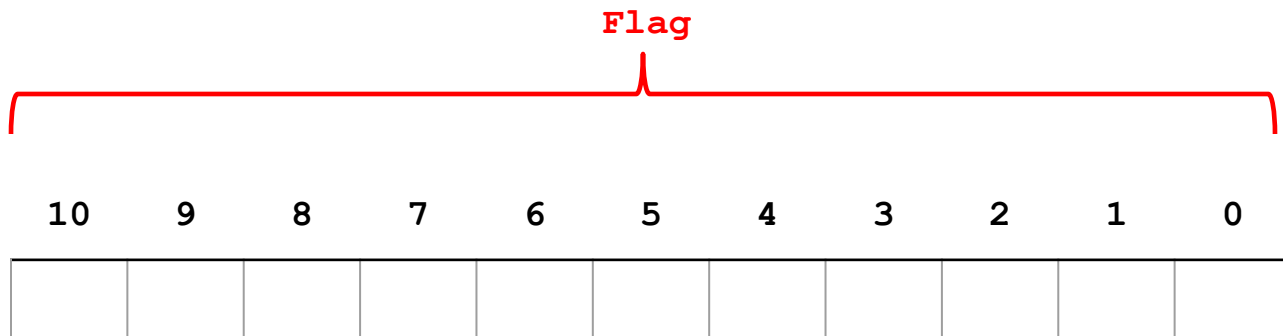


1 0

Read fails to align (**unmapped query**)

Read aligned successfully (**mapped query**)

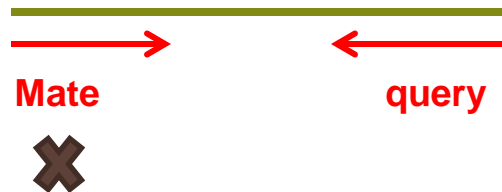
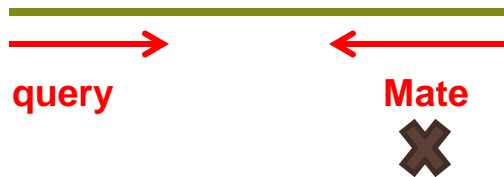


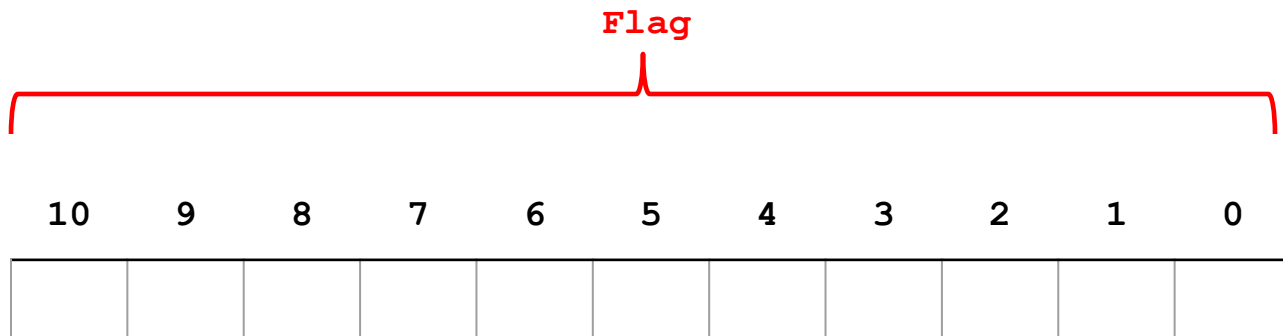


1 0

Read mate fails to align (**unmapped query mate**)

Read mate aligned successfully (**mapped query mate**)





1 0

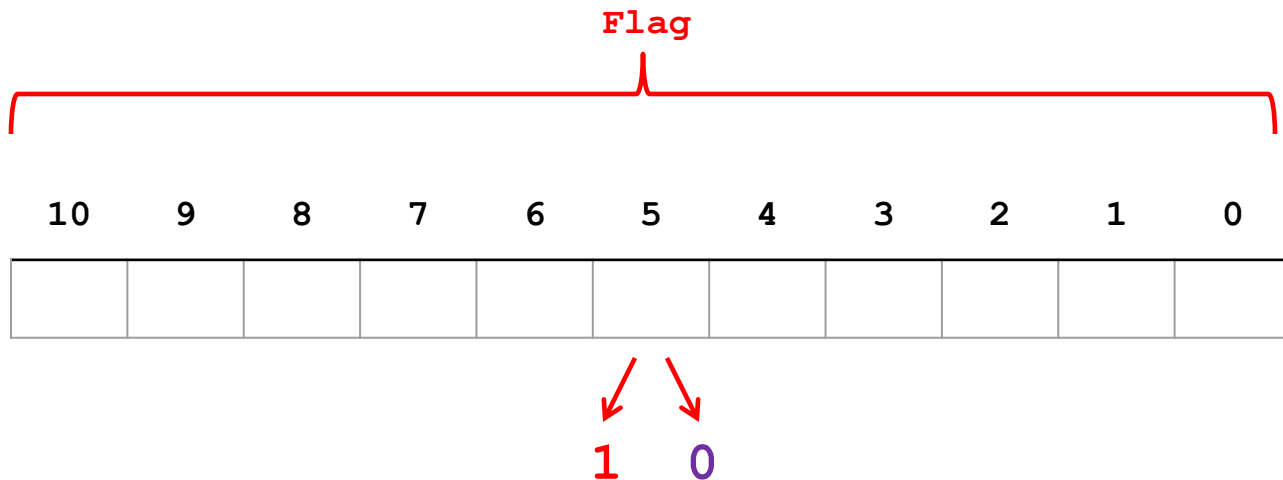
Strand of the query (0: forward direction,+, Watson strand )

Strand of the query (1: reverse complement,-, Crick strand )



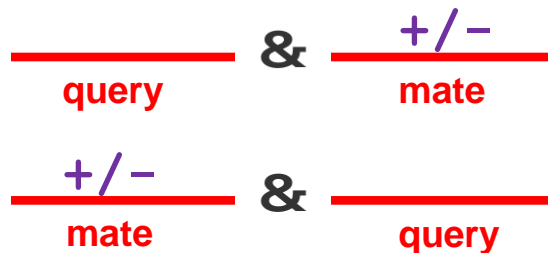
Query/Read



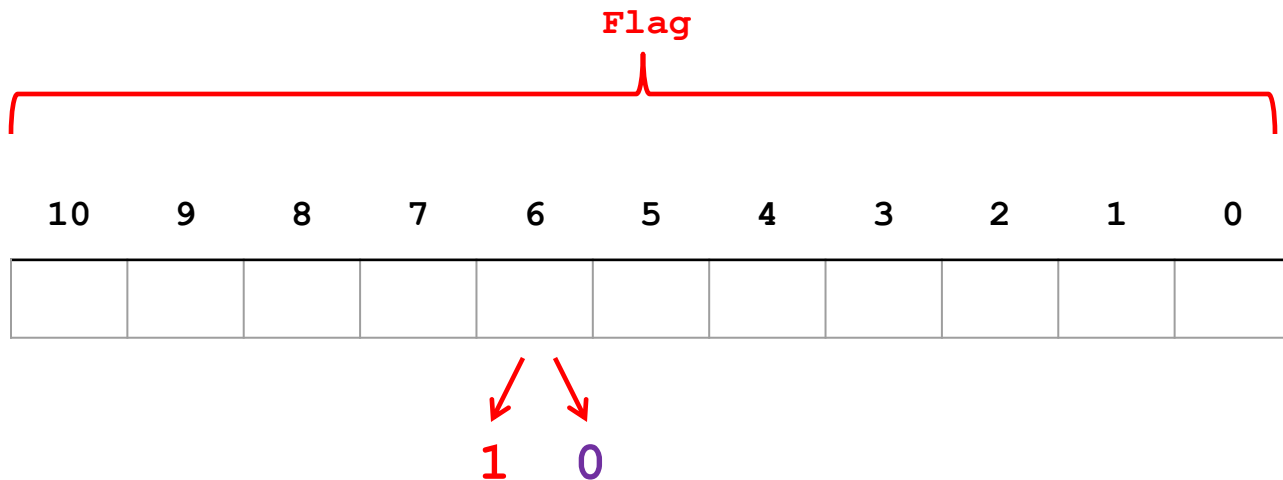


Strand of the query's mate (0: forward direction,+, Watson strand )

Strand of the query's mate (1: reverse complement,-, Crick strand )

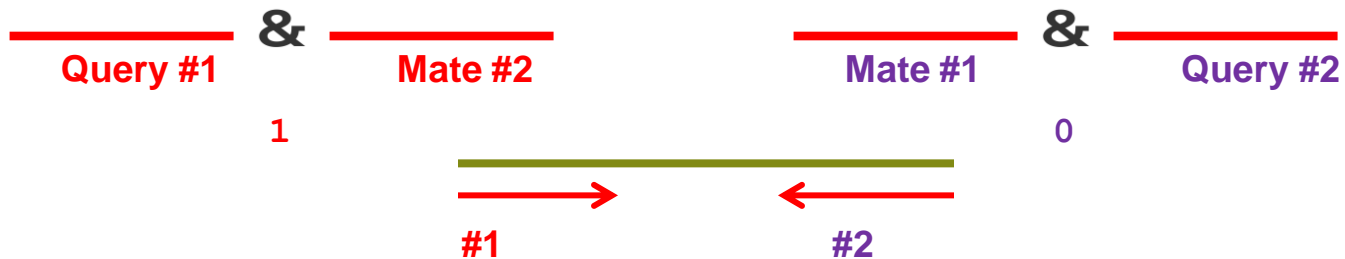


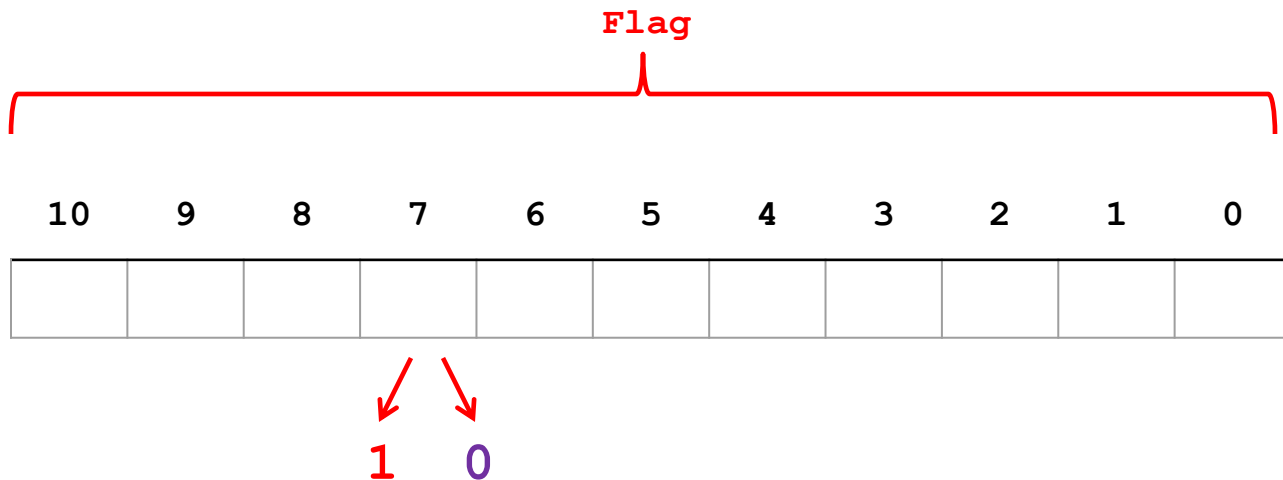




the query is the second read in the pair (0: #2 )

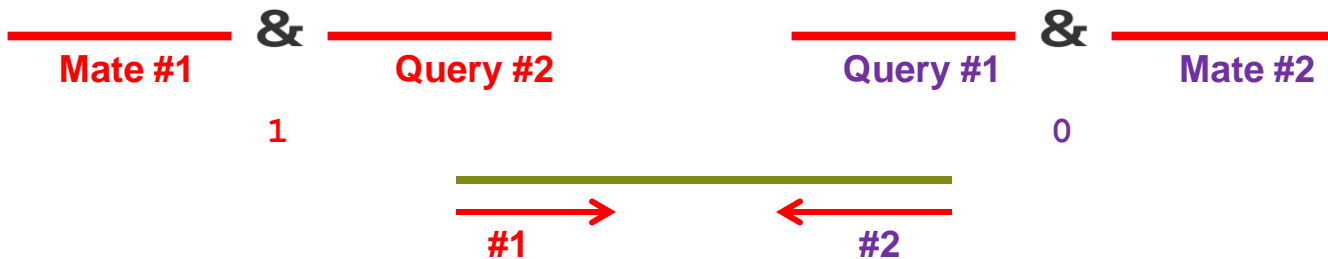
the query is the first read in the pair (1: #1 )





the query is the first read in the pair (0: #1 )

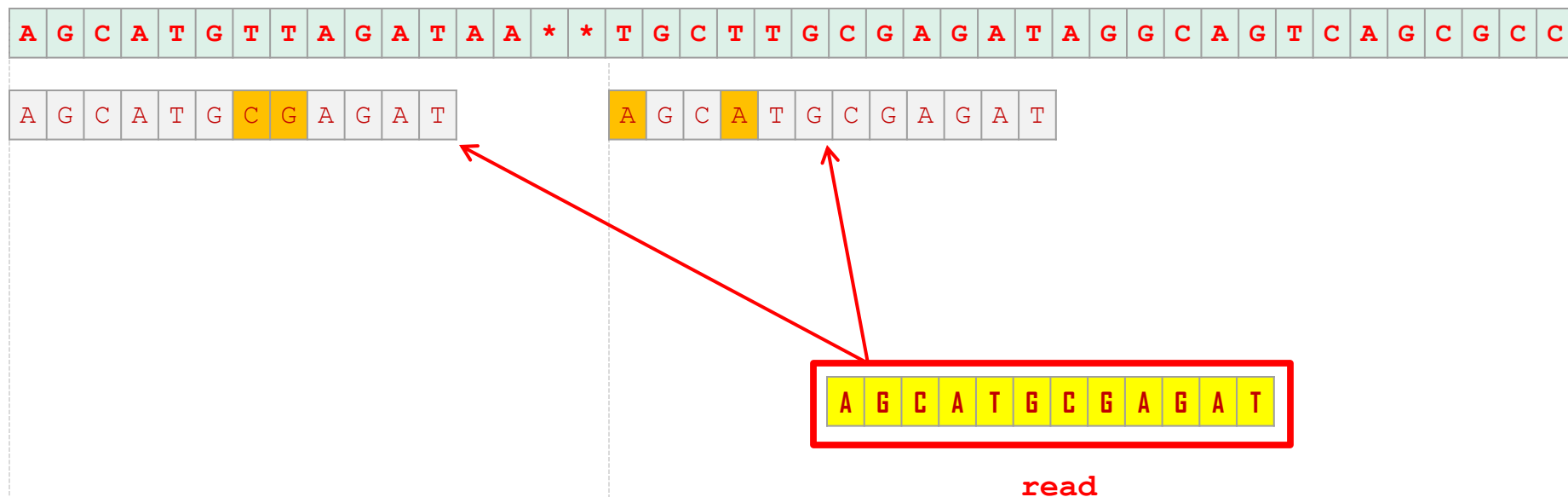
the query is the second read in the pair (1: #2 )

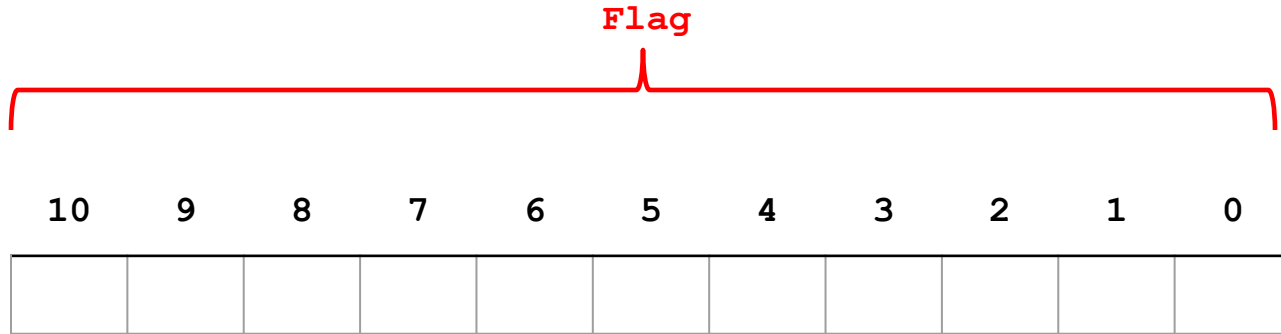


# Note

Multiple mapping positions

Which one is Primary ? Secondary?

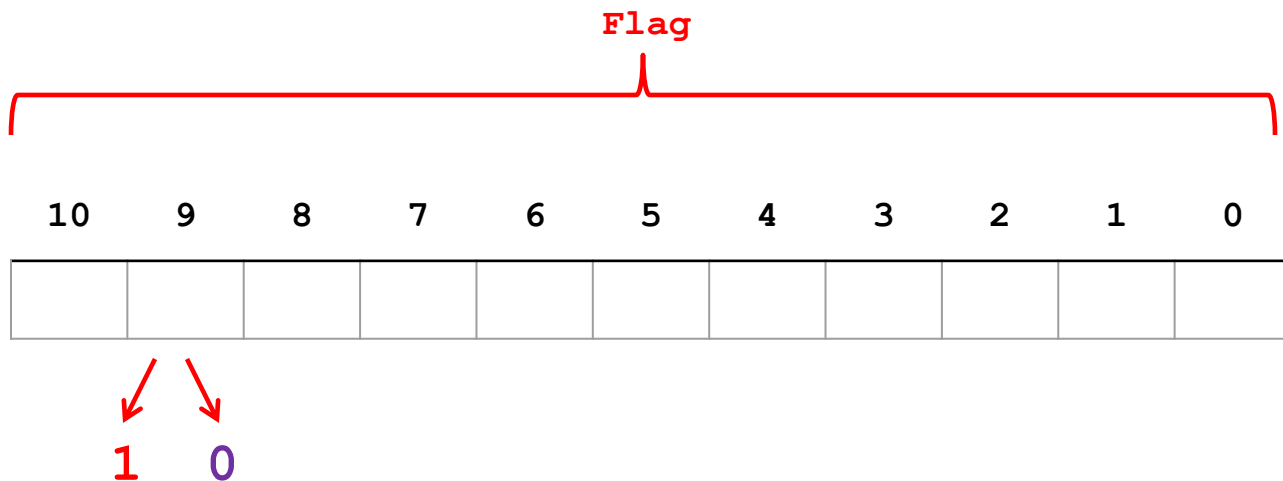




1 0

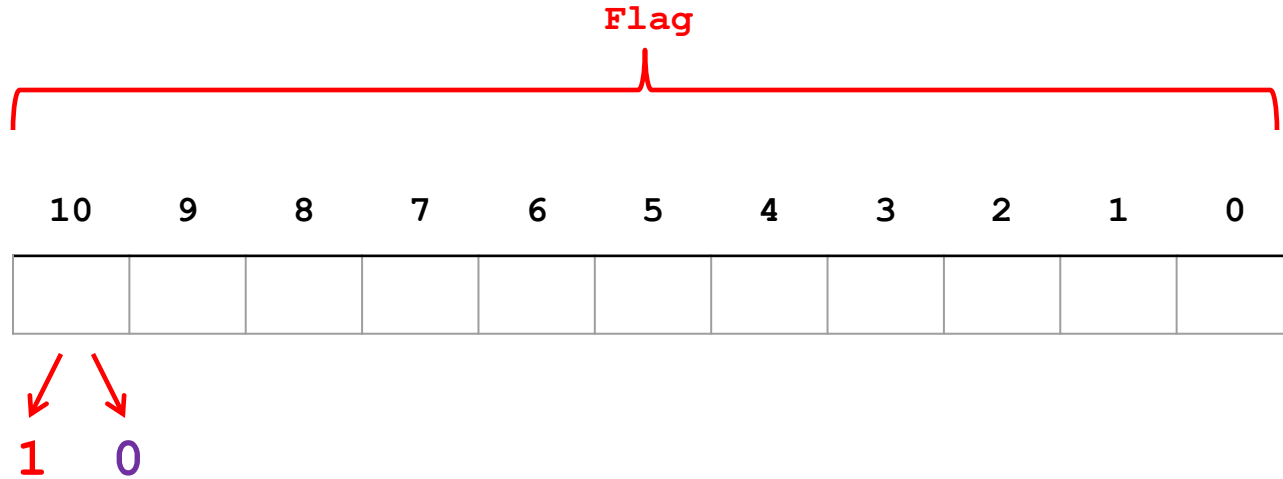
The alignment is primary(0: Primary Alignment )

The alignment is not primary (1: Secondary Alignment )



The read passing filters, such as platform/vendor quality controls (0)

The read not passing filters, such as platform/vendor quality controls (1)



Read is not PCR or optical duplicate(0 )

Read is PCR or optical duplicate (1)

# Note

A	G	C	A	T	G	T	T	A	G	A	T	A	A	*	*	G	A	T	A	G	C	T	G	T	G	C	T	A	G	G	C	A	G	T	C	A	G	C	G	C	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Linear alignment

A	A	C	C	G	A	T	A	G	C	T	G	T	G	C	T	A	G	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

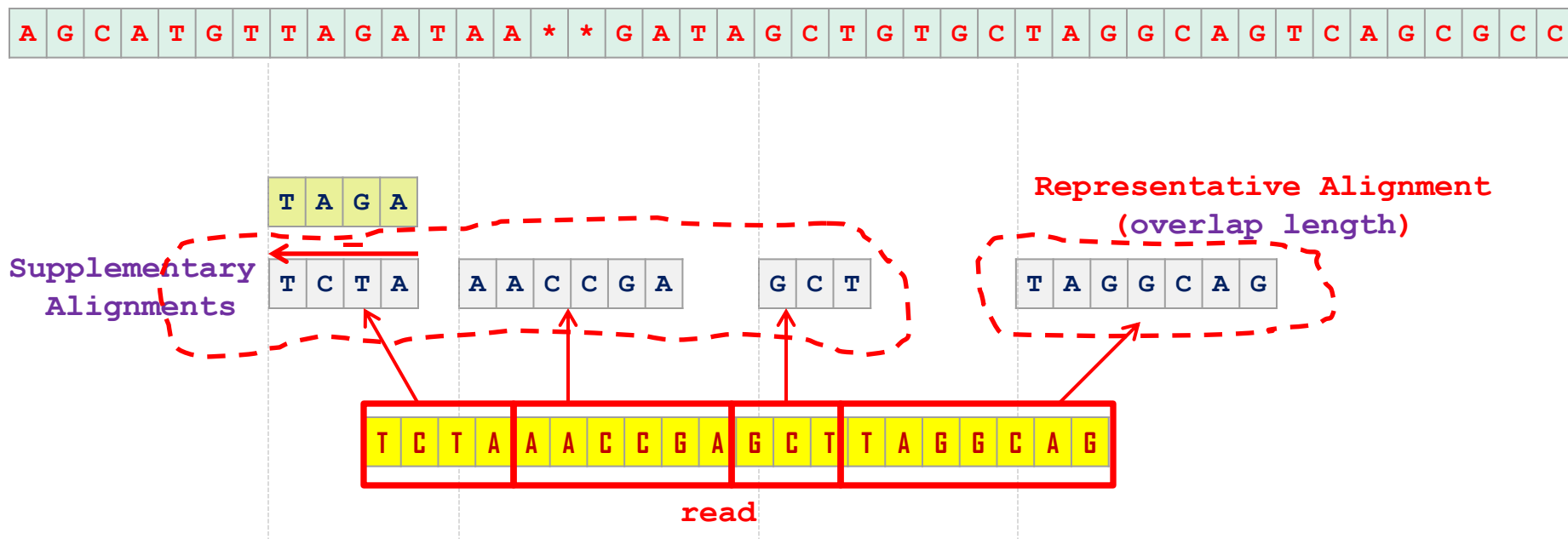


A	A	C	C	G	A	T	A	G	C	T	G	T	G	C	T	A	G	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

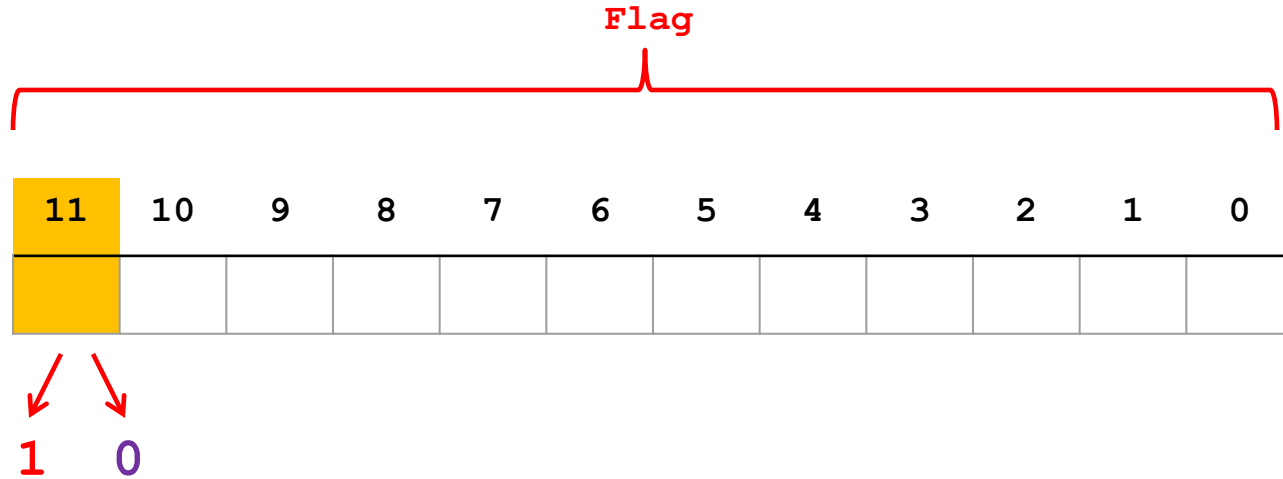
read

# Note

## Chimeric/Non-linear alignment

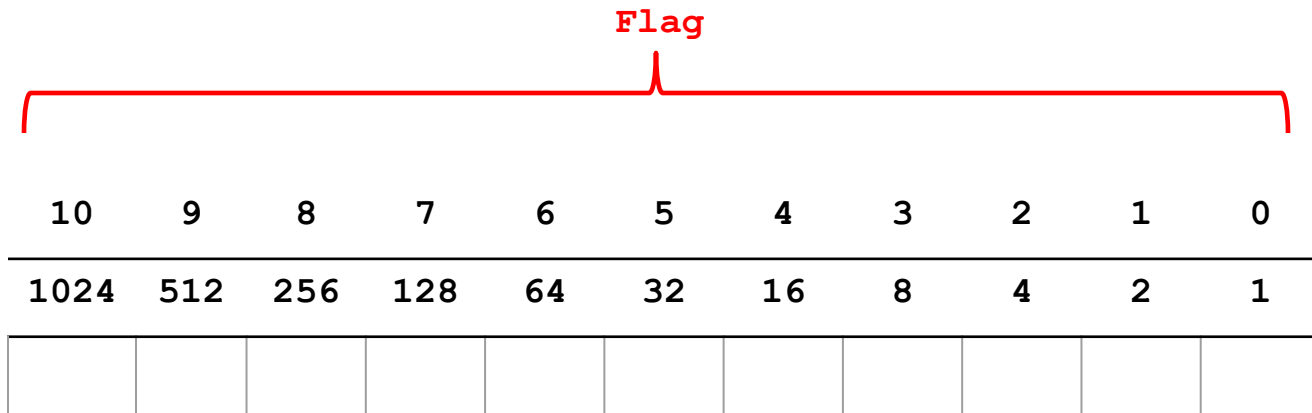






Not supplementary alignment (0)

Supplementary alignment (1)



ST-E00223:32:H5J57CCXX:4:1220:14651:8868	99	1	10086
--	----	---	-------

Flag										
10	9	8	7	6	5	4	3	2	1	0
1024	512	256	128	64	32	16	8	4	2	1
0	0	0	0	1	1	0	0	0	1	1

Read is aligned to the Watson strand, +

ST-E00223:32:H5J57CCXX:4:1220:14651:8868

Flag



99

1

10086

- ☐ Read is a paired-end.
- ☐ Pairs are aligned concordantly (proper pairs).
- ☐ Read's mate is aligned to the Crick strand, - .
- ☐ This is the first end ( read #1).



# SAM/BAM

## Flags

The `flags` field is a bitfield. Individual bits correspond to certain yes/no properties of the alignment. Here are the most relevant ones:

- Bit 0 (least significant): 1 if read is paired-end, 0 otherwise
- Bit 1: for paired-end reads only: 1 if the pair aligns concordantly, 0 otherwise
- Bit 2: 1 if read failed to align, 0 otherwise
- Bit 3: for paired-end reads only: 1 if the other end failed to align, 0 otherwise
- Bit 4: 1 if read aligned to Crick strand, 0 if Watson strand
- Bit 5: for paired-end reads only: 1 if the other end aligned to Crick strand, 0 if Watson strand
- Bit 6: for paired-end reads only: 1 if this is the first (#1) end, 0 if this is the second (#2) end
- Bit 7: for paired-end reads only: 0 if this is the first (#1) end, 1 if this is the second (#2) end

There are a few more that are used less often; see the [SAM specification] for details.

# SAM/BAM

POS

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCCDDCCDDDDCDCCCCDBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

**POS** is the 1-based offset into the reference sequence where the read aligned

# SAM/BAM

MAPQ

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCCDDCCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

**MAPQ:** For an aligned read, this is a confidence value; high when we're very confident we've found the correct alignment, low when we're not confident.

$$\text{MAPQ} = -10 * \log_{10}(P_{\text{map\_loc\_wrong}})$$

$(P_{\text{map\_loc\_wrong}})$	$\log_{10}(P_{\text{map\_loc\_wrong}})$	MAPQ
1	0	0
0.1	-1	10
0.01	-2	20
0.001	-3	30
0.0001	-4	40

# SAM/BAM

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCCDDCCDDDDDCDDDDDDCDCCDBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

The cigar string : encode the details of the alignment.

# Note

Mapping Position



A G C A T G T T A G A T A A G G G A T A G C T G T G C T A G G C A G T C A G C G C C

A A C C G A T A G C T A G C C T A A A C

CIGAR String

20M

Extended  
CIGAR String

2=2X7=3X3=2X1=



# Note

Mapping Position



A G C A T G T T A G A T A A - - G A T A G C T G T G C T A G G C A G T C A G C G C C

A A C C G A T A G C T A G C C T A A A C

CIGAR String

2M2I16M

Extended  
CIGAR String

2=2I7=3X3=2X1=

# Note

Mapping Position



A G C A T G T T A G A T A A - - G A T A G C T G T G C T A G G C A G T C A G C G C C

A A C C G A T A G C T A - - C T A A A C

CIGAR String

2M2I8M2D6M

Extended  
CIGAR String

2=2I7=1X2D3=2X1=

# Note

Mapping Position



A	G	C	A	T	G	T	T	A	G	A	T	A	A	G	G	G	A	T	A	G	C	T	G	T	G	C	T	A	G	G	C	A	G	T	C	A	G	C	G	C	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

T	T	C	C	G	A	T	A	G	C	T	A	G	C	C	T	G	C	C	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

# Note

Start Position

End Position

A G C A T G T T A G A T A A G G G A T A G C T G T G C T A G G C A G T C A G C G C C

T T C C G A T A G C T A G C C T G C C A

CIGAR String

4S13M3S

SEQ

TTCCGATAGCTAGCCTGCCA

# Note

Start Position

End Position

A G C A T G T T A G A T A A G G G A T A G C T G T G C T A G G C A G T C A G C G C C

T T C C G A T A G C T A G C C T G C C A

CIGAR String

4H13M3H

SEQ

GATAGCTAGCCTG

# Note

Mapping Position



A G C A T G T T A G A T A A G G G A T A G C T G T G C T A G G C A G T C A G C G C C

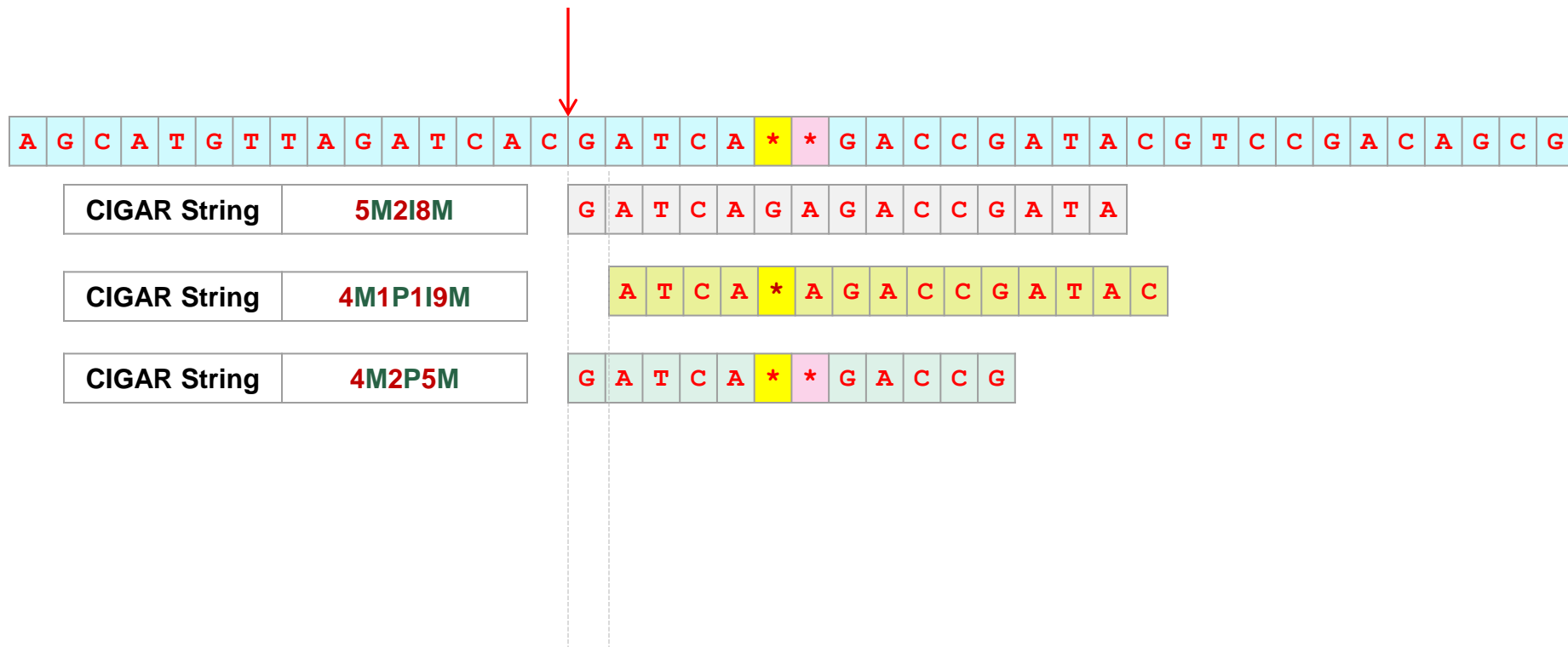
A A G G G . . . . . T A G G C

CIGAR String

5M10N5M

# Note

## Mapping Position



# SAM/BAM

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      cigar      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDDCCDDCCDDDDDCDDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDBDHFFFFDC@@
AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

The cigar string : encode the details of the alignment.

Operation	Meaning
M	Match*
D	Deletion w.r.t. reference
I	Insertion w.r.t. reference
N	Split or spliced alignment
S	Soft-clipping
H	Hard-clipping
P	Padding

Reference:  
Experimental:

ACCTGTC -- TACCTTACG  
ACCT - TCCATAC TTTATC

4M 1D 2M 2I 7M 2S

CIGAR string:

4M1D2M2I7M2S

LENGTH/OPERATION



# SAM/BAM

REF: AGCTAGCATCGTGTGCGCCCGTCTAGCATACGCATGATCGACTGTCAGCTAGTCAGACTAGTC

Read: GTGTAACCC.....TCAGAATA

Operation	Meaning
=	Exact match
X	Mismatch
D	Deletion w.r.t. reference
I	Insertion w.r.t. reference
N	Split or spliced alignment
S	Soft-clipping
H	Hard-clipping
P	Padding

The CIGAR for this alignment is :  
**9M32N8M.**

# SAM/BAM

The extended CIGAR string: M become = and X

Operation	Meaning
=	Exact match
X	Mismatch
D	Deletion w.r.t. reference
I	Insertion w.r.t. reference
N	Split or spliced alignment
S	Soft-clipping
H	Hard-clipping
P	Padding

Reference:

Experimental:

ACCTGTC--TACCTTACG

ACCT-TCCATACTTTATC

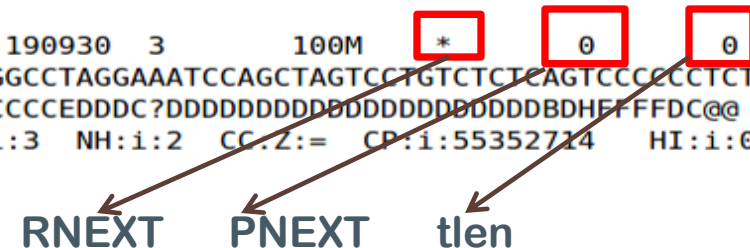


4= 1D 2= 2I 3= 1X 3= 2S

CIGAR string: 4=1D2=2I3=1X3=2S

# SAM/BAM

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTAAAGGTGGATGCGGTACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCCDDCCDDDDCDCCCCBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDDDDBDHFFDC@@
      AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```



- ✓ **rnext** only relevant for paired-end reads; name of the reference sequence where other end aligned.
- ✓ **pnext** only relevant for paired-end reads; 1-based offset into the reference sequence where other end aligned.
- ✓ **Tlen/size** only relevant for paired-end reads; insert length inferred from alignment.

# SAM/BAM

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
      CCGTGTTTTAAAGGTGGATGCGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCCDDCCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCCEDDDC?DDDDDDDDDDDDDDDDDDDDBDHEFFDC@@
      AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

RNEXT      PNEXT      tlen

= : RNAME equals RNEXT

\*: information is not available

```
my_read 99      chr2:172936693-172938111      129      60      100M      =      429      400      CTAAC TAGCCTGGGAAA
my_read 147     chr2:172936693-172938111      429      60      100M      =      129      -400     TCGAGCTCTGCATTC
```

# SAM/BAM

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
CCCTGTTTAAAGCTCCATCCCGCTCAGCTTCCGAGCTAGGCTTAGCGATTCTTACTTCCGCTAGCAAATCCAGCTACTCCTCTCTCTCAGTCCCCCTCT
CBBDCDDCCDDDDDCDDDDDDDCDCCCDRC?DDDDDDDDDDDDDDDDDDCCDCCDDDDDDDDDDDDCCCCEDDDDC?DDDDDDDDDDDDDDDDDDDDDDDDRDHEEEEDC@@
AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20T15A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55552714      HI:i:0
```

seq

qual

extras

tab-separated "extra" fields, usually optional and aligner-specific but often very important!

# SAM/BAM

```
HWI-ST1145:74:C101DACXX:7:1102:4284:73714      16      chr20      190930      3      100M      *      0      0
CCGTGTTTAAAGGTGGATGCGGTACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCCTGTCTCTCAGTCCCCCTCT
C      BBDCEDDCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDHFFFFDC@@
AS:i:-15      XM:i:3      XO:i:0      XG:i:0      MD:Z:55C20C13A9      NM:i:3      NH:i:2      CC:Z:=      CP:i:55352714      HI:i:0
```

Alignment Score  $x?$   
TAG:TYPE:VALUE reserved fields for end users

String for mismatching positions.

[http://chagall.med.cornell.edu/galaxy/references/SAM\\_BAM\\_Specification.pdf](http://chagall.med.cornell.edu/galaxy/references/SAM_BAM_Specification.pdf)

<https://samtools.github.io/hts-specs/SAMtags.pdf>

<https://github.com/vsbuffalo/devnotes/wiki/The-MD-Tag-in-BAM-Files>

# SAM/BAM

- MD: String for mismatching positions.
- The MD field aims to achieve SNP/indel calling without looking at the reference.
- The MD field ought to match the CIGAR string.

**MD: Z: 10A5^AC6**



**Thank you!**