



# Data formats

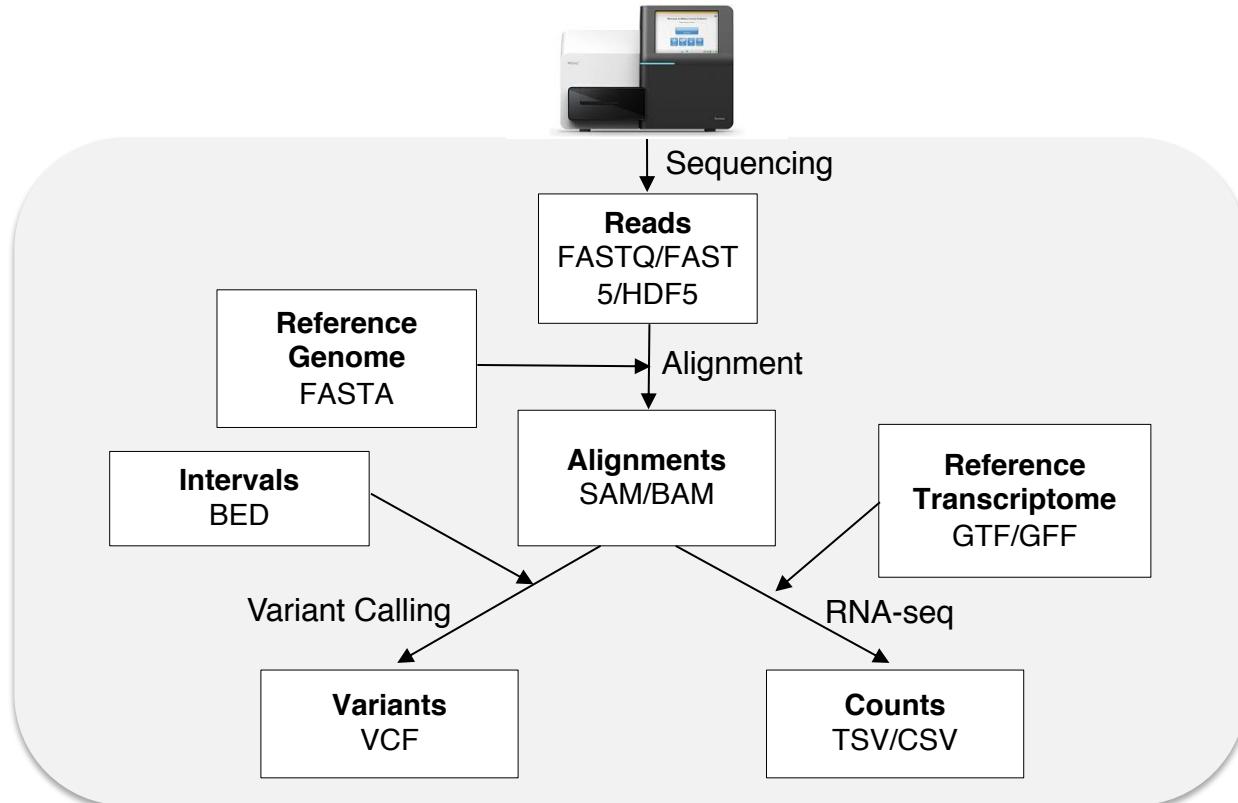
## Precision Oncology Course



CNIO BIOINFORMATICS UNIT

**Coral Fustero Torre**  
Bioinformatics Unit,  
Structural Biology Programme.  
[cfustero@cnio.es](mailto:cfustero@cnio.es) | [bioinformatics.cnio.es](http://bioinformatics.cnio.es)

# Formats outline



# Data formats

## Reference Genome: FASTA format

- Typical extensions: .fasta, .fas, .fa, .fna, .fsa
- Each sequence is composed by at least two consecutive lines:
  - ">" Sequence name and optional description (space separated)
  - Line(s) with the whole sequence

The diagram illustrates the FASTA format. The top section shows a single DNA sequence named "DNA\_SEQUENCE\_1" with a length of 60 characters. The sequence starts with "NNNNNNCT" and ends with "CCGTGAAGATGGAGCCATATTCC". A horizontal arrow below the sequence indicates its total length of 60 characters. The bottom section shows multiple sequences, labeled "DNA\_SEQUENCE\_1", "DNA\_SEQUENCE\_2", and "DNA\_SEQUENCE\_3", each starting with a greater than symbol and followed by their respective DNA sequences.

```
>DNA_SEQUENCE_1
NNNNNNCTGGGGGACAGAACCCATGGTGGCCCCGGCTCTCCCCAGTATCCAGTCCT
CCGTGAAGATGGAGCCATATTCC
60 chars

>DNA_SEQUENCE_1
NNNNNNCTGGGGGACAGAACCCATGGTGGCCCCGGCTCTCCCCAGTATCCAGTCCT
>DNA_SEQUENCE_2
GGGGGACAGAACCCATGGTGGCCCCGGCTCTCCCCAGTATCCAGTCCT
>DNA_SEQUENCE_3
CTCCTCCCCAGTATCCAGTCCTGGGGGACAGAACCCATGGTGGCCCCGGCTCTCCCCAGTATCCCA
```

We can have multiple sequences in the same file: **multifasta**

# Data formats

## Nucleotides codes: IUPAC

IUPAC nucleotide code	Base
A	Adenine
C	Cytosine
G	Guanine
T (or U)	Thymine (or Uracil)
R	A or G
Y	C or T
S	G or C
W	A or T
K	G or T
M	A or C
B	C or G or T
D	A or G or T
H	A or C or T
V	A or C or G
N	any base
. or -	gap

# Data formats

## Reads: FASTQ

- Typical extensions: .fq, .fastq
- Each read is composed by 4 lines:
  - "@" Read name and optional description (space separated)
  - Sequence
  - "+" (optionally: repeat the read name)
  - Base Quality Score

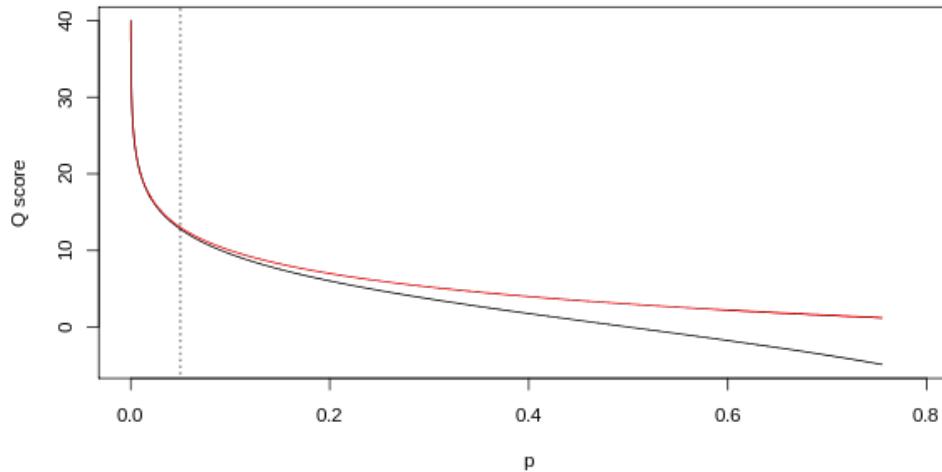
The diagram illustrates the structure of a FASTQ read. It consists of four lines of text:

- readname**: @HWI-EAS209\_0006\_FC706VJ:5:58:5894:21141 ATCACG
- sequence**: TTAATTGGTAAATAAATCTCCTAACAGCTTAGATNTTACCTNNNNNNNNNTAGTTCTTGAGATTGTTGGGGGAGACATTTGTGATTGCCCTGAT
- comment**: + ! "#\$%& ' ()\*+, - ./0123456789: ; <=> ? @ABCDEFGHIJKLMNPQRSTUVWXYZ[\]^\_ ` abcdefghijklmnopqrstuvwxyz{| }))))))))
- Quality**: (This label points to the last line of the sequence, which contains a series of ASCII characters representing base qualities.)

# Data formats

## Reads: FASTQ – Quality Score

- Phred quality scores  $Q_{Phred}$  are defined as a property which is logarithmically related to the base-calling **error** probabilities  $p$   
 $Q_{Phred} = -10 \log_{10}(p)$
- The score is written as the character whose ASCII code is:  $Q_{Phred} + 33$
- The higher the the  $Q_{Phred}$  , the lower the probability that the base calling is erroneous



# Data formats

## Reads: FASTQ – ASCII code

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0 000	NUL	(null)	32	20 040	&#32;	Space		64	40 100	&#64;	Ø		96	60 140	&#96;	`	
1	1 001	SOH	(start of heading)	33	21 041	&#33;	!	!	65	41 101	&#65;	A		97	61 141	&#97;	a	
2	2 002	STX	(start of text)	34	22 042	&#34;	"	"	66	42 102	&#66;	B		98	62 142	&#98;	b	
3	3 003	ETX	(end of text)	35	23 043	&#35;	#	#	67	43 103	&#67;	C		99	63 143	&#99;	c	
4	4 004	EOT	(end of transmission)	36	24 044	&#36;	\$	\$	68	44 104	&#68;	D		100	64 144	&#100;	d	
5	5 005	ENQ	(enquiry)	37	25 045	&#37;	¤	¤	69	45 105	&#69;	E		101	65 145	&#101;	e	
6	6 006	ACK	(acknowledge)	38	26 046	&#38;	¤	¤	70	46 106	&#70;	F		102	66 146	&#102;	f	
7	7 007	BEL	(bell)	39	27 047	&#39;	‘	‘	71	47 107	&#71;	G		103	67 147	&#103;	g	
8	8 010	BS	(backspace)	40	28 050	&#40;	(	(	72	48 110	&#72;	H		104	68 150	&#104;	h	
9	9 011	TAB	(horizontal tab)	41	29 051	&#41;	)	)	73	49 111	&#73;	I		105	69 151	&#105;	i	
10	A 012	LF	(NL line feed, new line)	42	2A 052	&#42;	*	*	74	4A 112	&#74;	J		106	6A 152	&#106;	j	
11	B 013	VT	(vertical tab)	43	2B 053	&#43;	+	+	75	4B 113	&#75;	K		107	6B 153	&#107;	k	
12	C 014	FF	(NP form feed, new page)	44	2C 054	&#44;	,	,	76	4C 114	&#76;	L		108	6C 154	&#108;	l	
13	D 015	CR	(carriage return)	45	2D 055	&#45;	-	-	77	4D 115	&#77;	M		109	6D 155	&#109;	m	
14	E 016	SO	(shift out)	46	2E 056	&#46;	.	.	78	4E 116	&#78;	N		110	6E 156	&#110;	n	
15	F 017	SI	(shift in)	47	2F 057	&#47;	/	/	79	4F 117	&#79;	O		111	6F 157	&#111;	o	
16	10 020	DLE	(data link escape)	48	30 060	&#48;	Ø	Ø	80	50 120	&#80;	P		112	70 160	&#112;	p	
17	11 021	DC1	(device control 1)	49	31 061	&#49;	1	1	81	51 121	&#81;	Q		113	71 161	&#113;	q	
18	12 022	DC2	(device control 2)	50	32 062	&#50;	2	2	82	52 122	&#82;	R		114	72 162	&#114;	r	
19	13 023	DC3	(device control 3)	51	33 063	&#51;	3	3	83	53 123	&#83;	S		115	73 163	&#115;	s	
20	14 024	DC4	(device control 4)	52	34 064	&#52;	4	4	84	54 124	&#84;	T		116	74 164	&#116;	t	
21	15 025	NAK	(negative acknowledge)	53	35 065	&#53;	5	5	85	55 125	&#85;	U		117	75 165	&#117;	u	
22	16 026	SYN	(synchronous idle)	54	36 066	&#54;	6	6	86	56 126	&#86;	V		118	76 166	&#118;	v	
23	17 027	ETB	(end of trans. block)	55	37 067	&#55;	7	7	87	57 127	&#87;	W		119	77 167	&#119;	w	
24	18 030	CAN	(cancel)	56	38 070	&#56;	8	8	88	58 130	&#88;	X		120	78 170	&#120;	x	
25	19 031	EM	(end of medium)	57	39 071	&#57;	9	9	89	59 131	&#89;	Y		121	79 171	&#121;	y	
26	1A 032	SUB	(substitute)	58	3A 072	&#58;	:	:	90	5A 132	&#90;	Z		122	7A 172	&#122;	z	
27	1B 033	ESC	(escape)	59	3B 073	&#59;	:	:	91	5B 133	&#91;	[		123	7B 173	&#123;	{	
28	1C 034	FS	(file separator)	60	3C 074	&#60;	<	<	92	5C 134	&#92;	\		124	7C 174	&#124;		
29	1D 035	GS	(group separator)	61	3D 075	&#61;	=	=	93	5D 135	&#93;	]		125	7D 175	&#125;	}	
30	1E 036	RS	(record separator)	62	3E 076	&#62;	>	>	94	5E 136	&#94;	^		126	7E 176	&#126;	~	
31	1F 037	US	(unit separator)	63	3F 077	&#63;	?	?	95	5F 137	&#95;	_		127	7F 177	&#127;	DEL	

# Data formats

## Reads: FASTQ – Single-end or Paired-end?

One unique sample can have 1 or 2 files:

- If **single-end** seq → 1 file  
filename: \*.fastq
- If **paired-end** seq → 2 files  
filenames: \*\_R1.fastq AND \*\_R2.fastq

\_R1

```
@SN1083:379:H8VA1ADXX:2:1101:1248:2144 1:N:0:12
CCTAAATGGTGCATGCTAGGAGGCCGTGCCCTTCTTGAAAAGTTGTATGTGAA
+
BBBFFFFFFFBFFFIIIIIFI<FFIIIIIFIIFBFIIIIIIIIFFIIIIIFI
```

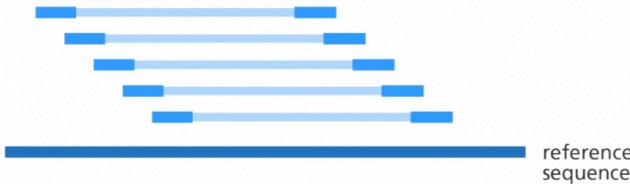
\_R2

```
@SN1083:379:H8VA1ADXX:2:1101:1248:2144 2:N:0:12
CATTTTCGACGTTGTTAATAAGCTCTGCGTACTTGCAAGCTATCTGCGCGAACG
+
BBBFFFFFFFIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIFFF
```

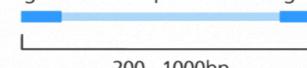
### Single-end reads



### Paired-end reads



sequenced fragment    unknown sequence    sequenced fragment



# Exercise: FastQC

1. Download data: [https://fundacioncnio-my.sharepoint.com/:u/g/personal/cfustero\\_cnio\\_es/Eb1S\\_Y-D8SpBspZK3MFZQiUBcbncd8hul6\\_c0YPOihnwQ?e=AQeHGu](https://fundacioncnio-my.sharepoint.com/:u/g/personal/cfustero_cnio_es/Eb1S_Y-D8SpBspZK3MFZQiUBcbncd8hul6_c0YPOihnwQ?e=AQeHGu)

2. Go to the terminal and open the file

```
$ more WEx_Normal_R1.fastq
```

3. How would you detect the number of sequences present in the file?

## Useful commands

**grep**: searches plain-text data sets for lines that match a regular expression

**wc**: counts words

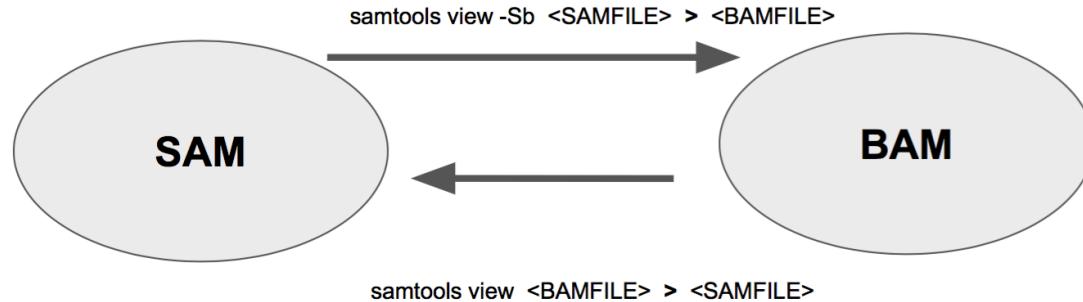
**wc -l**: counts lines

**|**: lets you send the output of one command to another.

# Data formats

## Alignment: SAM or BAM

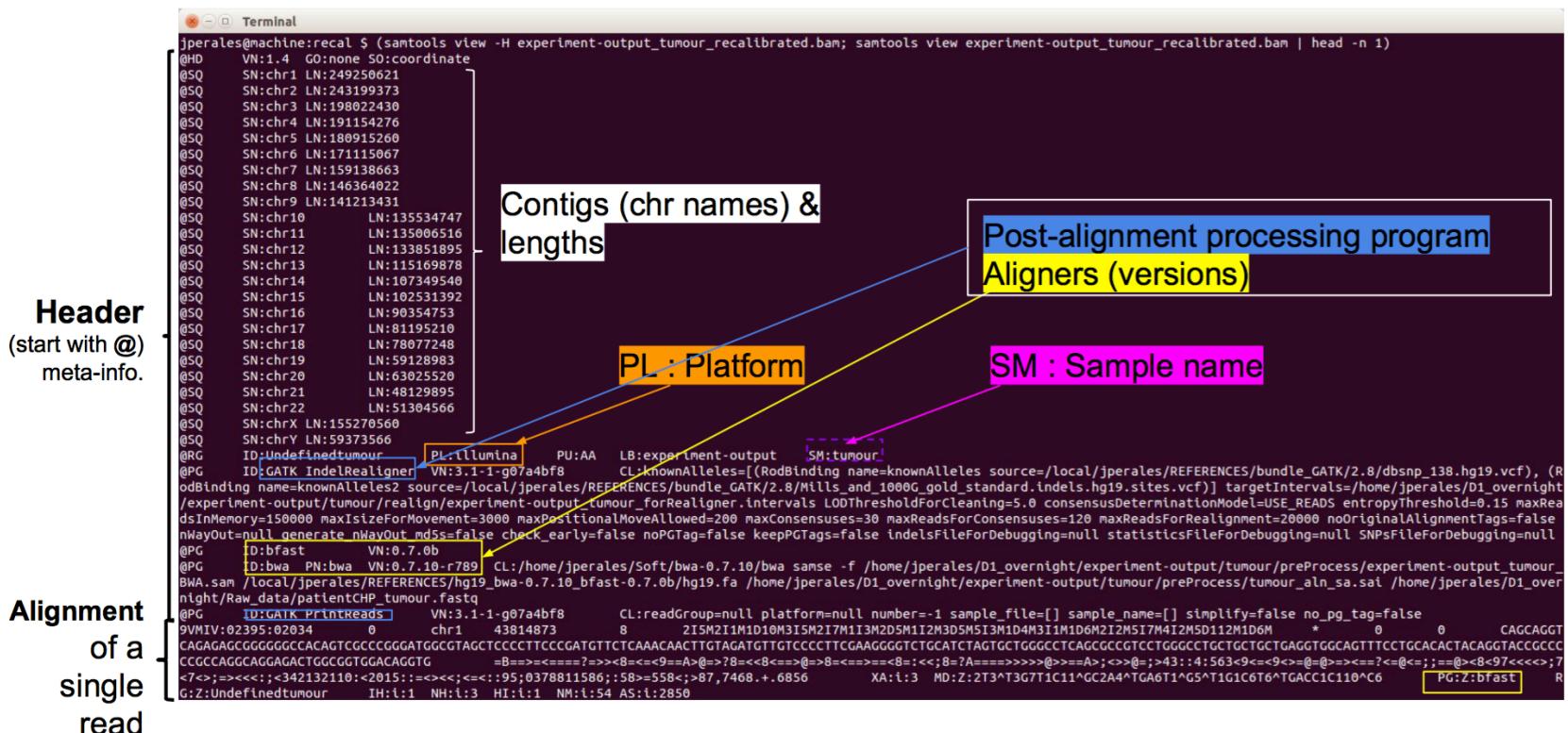
- **SAM** is the human readable text format (.sam extension)
- **BAM** is the binary, machine efficient format (.bam extension)
- Both contains exactly the same information and are interconvertible (samtools)



File specifications <https://samtools.github.io/hts-specs/SAMv1.pdf>

# Data formats

## Alignment: SAM or BAM - Header



# Data formats

## Alignment: SAM or BAM - Alignments

#1 ReadName	#2 99	#3 chr10	#4 2	#5 30	#6 <b>3MD2M1I1M</b>	#7 =	#8 14	#9 20	#10 CATCTG	#11 jjjjjjj	#12 z:Aligner
----------------	----------	-------------	---------	----------	------------------------	---------	----------	----------	---------------	----------------	------------------



If single-end:

7. reference sequence name of the alignment of the next read in sequence
8. position in the alignment of the next read in sequence
9. number of bases covered by reads from the same fragment. Plus/minus means the current read is the leftmost/rightmost read

# Data formats

## Alignment: SAM or BAM - Alignments

#1 ReadName	#2 99	#3 chr10	#4 2	#5 30	#6 3MD2M1I1M	#7 =	#8 14	#9 20	#10 CATCTG	#11 iiiiiji	#12 z:Aligner
<b>SAM FLAGS</b>											

FIELDS		
#Col	Field	Description
1.	QNAME	read name
2.	FLAG	bitwise FLAG* (unmapped, pair unmapped, properly mapped, ...)
3.	RNAME	Reference sequence name (e.g. chr1).
4.	POS	1-based leftmost position.
5.	MAPQ	Mapping Quality (Phred-scaled). Scale 0 to 255.
6.	CIGAR	extended CIGAR string.
7.	MRNM	Paired-end: Mate Reference sequence Name (= if same as RNAME).
8.	MPOS	Paired-end: 1-based Mate position.
9.	TLEN	Paired-end: Insert size
10.	SEQ	Read sequence
11.	QUAL	Base Quality Score from the Read sequence.
12.	OPT	Optional Tags.

Decimal	Description of read
1	Read paired
2	Read mapped in proper pair
4	Read unmapped
8	Mate unmapped
16	Read reverse strand
32	Mate reverse strand
64	First in pair
128	Second in pair
256	Not primary alignment
512	Read fails platform/vendor quality checks
1024	Read is PCR or optical duplicate
2048	Supplementary alignment

One of the reads is unmapped:  
[73](#), [133](#), [89](#), [121](#), [165](#), [181](#), [101](#), [117](#),  
[153](#), [185](#), [69](#), [137](#)

Both reads are unmapped:  
[77](#), [141](#)

Mapped within the insert size and in correct orientation:  
[99](#), [147](#), [83](#), [163](#)

Mapped within the insert size but in wrong orientation:  
[67](#), [131](#), [115](#), [179](#)

Mapped uniquely, but with wrong insert size:  
[81](#), [161](#), [97](#), [145](#), [65](#), [129](#), [113](#), [177](#)

# Data formats

## Alignment: SAM or BAM - CIGAR

### Concise Idiosyncratic Gapped Alignment Report

- Compressed representation of an alignment
- Format:** A CIGAR string is made up of <integer><op> pairs
- Where "op" is an operation specified as a single character, usually an uppercase letter (see table)

RefPos:	1 2 3 4 5 6 7 8 9
Reference:	C C A T A C T - G A
Read:	C A T - C T A G
POS:	2
CIGAR:	3M1D2M1I1M

Op	Description
M	alignment match (can be a sequence match or mismatch)
I	insertion to the reference
D	deletion from the reference
N	skipped region from the reference
S	soft clipping (clipped sequences present in SEQ)
H	hard clipping (clipped sequences NOT present in SEQ)
P	padding (silent deletion from padded reference)
=	sequence match
X	sequence mismatch

Reference sequence with aligned reads	CIGAR string	Explanation
C T 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 A A G G A T A * C T G G A T A A * G G A T A T G T T A [REDACTED] T G C T A	1M2I4M1D3M 5M1P1I4M 5M1N5M 3S8M 3H8M	Insertion & Deletion Padding & Insertion Spliced read Soft clipping Hard clipping
a a a C A T G T T A G A A A C A T G T T A G		

# Data formats

## Intervals: BED

- Typical extension: .bed
  - The first three are required BED fields, the rest are optional
1. **chrom** - The name of the chromosome (e.g. chr3, chrY, chr2)
  2. **chromStart** - The starting position of the feature in the chromosome  
    0-based (The first base in a chromosome is numbered 0)
  3. **chromEnd** - The ending position of the feature in the chromosome or scaffold

Additionally, 9 optional fields:

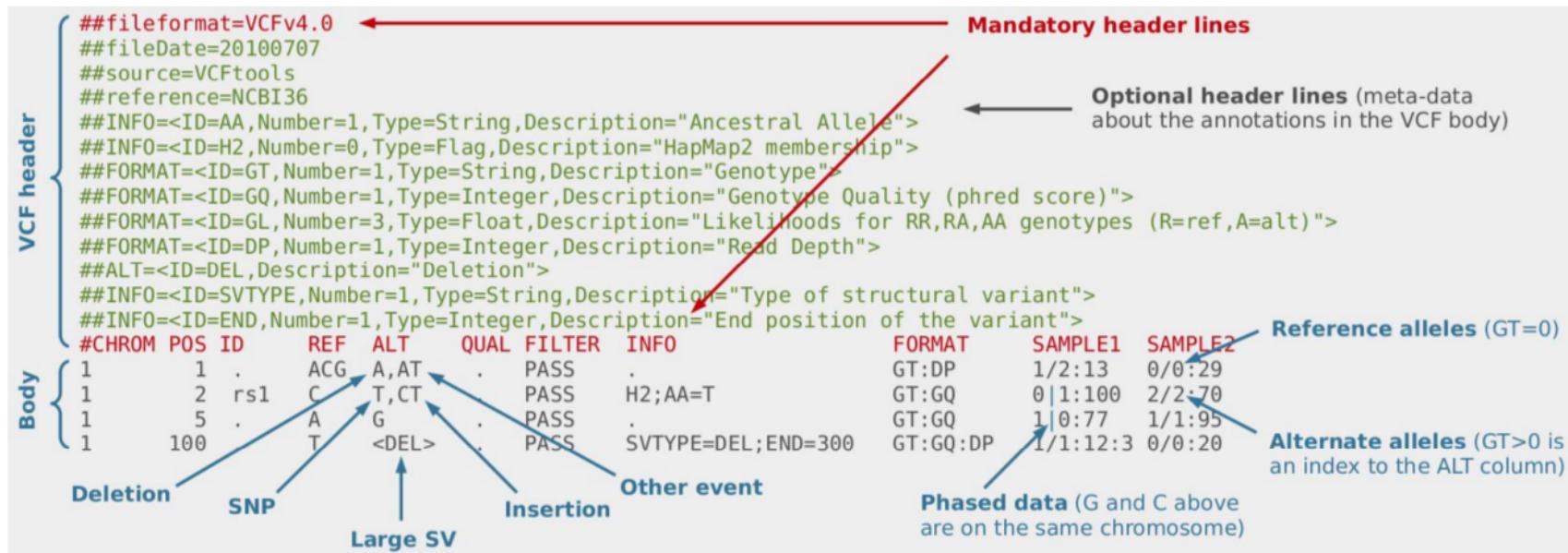
4. name - Defines the name of the BED line
5. Score (. or a number between 0 and 1000)
6. strand (+ forward, - reverse)
7. thickStart
8. thickEnd
9. itemRgb (255,0,0)
10. blockCount
11. blockSizes
12. blockStarts

```
track name="CHP2_Designed" description="Amplicon_Insert_CHP2" type=bedDetail ionVersion=4.0
chr1 43814968    43815086    CHP2_MPL_1    .    GENE_ID=MPL
chr1 115252185   115252269   CHP2_NRAS_3   .    GENE_ID=NRAS
chr1 115256504   115256584   CHP2_NRAS_2   .    GENE_ID=NRAS
chr1 115258689   115258774   CHP2_NRAS_1   .    GENE_ID=NRAS
chr2 29432572    29432680    CHP2_ALK_2    .    GENE_ID=ALK
chr2 29443607    29443729    CHP2_ALK_1    .    GENE_ID=ALK
chr2 209113103   209113206   CHP2_IDH1_1   .    GENE_ID=IDH1
chr2 212288904   212288990   CHP2_ERBB4_8   .    GENE_ID=ERBB4
chr2 212530051   212530180   CHP2_ERBB4_7   .    GENE_ID=ERBB4
```

# Data formats

## Variants: vcf

- Typical extension: .vcf (.bcf binary counterpart)
- Not all records in a VCF are true calls, the FILTER column specifies those which passed the calling
- QUAL is the score assigned to a given call. The greater QUAL is, the more reliable is. It is in log-scale



# Data formats

## Reference transcriptome: GTF/GFF

- Typical extensions: .gtf, .gff
- General Feature Format / Gene Transfer Format
- Annotation file for features
- One line per feature. 9 columns. Tab separated
- Used in: Reference transcriptomes (RNAseq) or to upload features to Genomic Browsers

chr1	unknown	exon	3214482	3216968	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	stop_codon	3216022	3216024	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	CDS	3210025	3210908	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	CDS	3421702	3421901	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	exon	3421702	3421901	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	CDS	3670552	3671348	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	exon	3670552	3671498	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	start_codon	3671340	3671348	.	.	.	gene_id "Xkr4"; gene_name "Xkr4"; p_id "P14345"; transcript_id "NM_001011874"; tss_id "TSS25485";
chr1	unknown	exon	4290846	4293012	.	.	.	gene_id "Rp1"; gene_name "Rp1"; p_id "P16188"; transcript_id "NM_001195662"; tss_id "TSS5760";
chr1	unknown	stop_codon	4292981	4292983	.	.	.	gene_id "Rp1"; gene_name "Rp1"; p_id "P16188"; transcript_id "NM_001195662"; tss_id "TSS5760";

seqname  
source  
feature  
start (1-based)  
end (1-based)  
strand  
frame  
score  
attribute

'0' first base of the feature is the first base of a codon  
'1' second base of the feature is the first base of a codon  
'2' third base of the feature is the first base of a codon

# Data formats

## Counts: CSV/TSV

### Comma-separated Text

```
Gene,Sample1,Sample2,Sample3  
DDR1,0,13,21  
TP53, 12,1,2  
KRAS, 50, 5,5
```

### Tab-separated Text

```
SRA_SAMPLE      SAMPLE_NAME  
SRS1307733  GSM2069823  
SRS1307731  GSM2069824
```



Excel

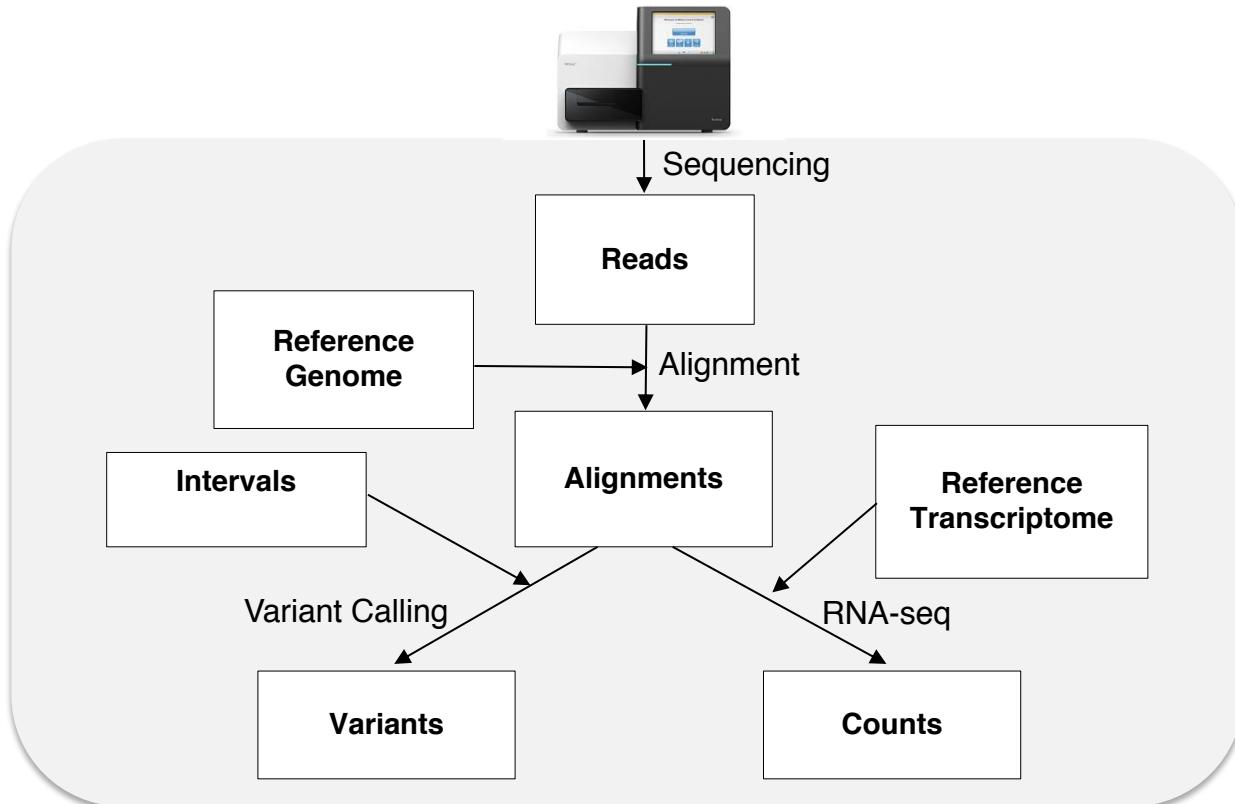
No format specifications.

There are no limits in terms of the matrix dimension (rows X columns).

However, there are good practices in Bioinformatics:

1. Use intuitive column names without spaces or rare chars. Instead of those, use "." or "\_" .
2. Data dimension: higher number of rows than columns.
3. Rows → individual observations or records. E.g. genomic positions.
4. Columns→ Individual variables for each individual observation. E.g. Chromosome, position, Score, Gene mutated, etc.
5. Do NOT mix data from different observations

# Formats outline



# Data formats cheat sheet

Format	Uses	Example	File type	Software Management	File Extension
Fasta	<a href="#"><b>Human genome</b></a> Define biological sequences (DNA, RNA, cDNA, proteins).	human_genome.fa	Plain text	samtools, picard-tools	.fa; .fasta
FastQ	<a href="#"><b>Raw sequencing data</b></a> Single-end sequencing → 1 file Paired-end sequencing → 2 files (R1 and R2 for each end, respectively)	DNAseq_raw_data.fastq (DNAseq_R1.fastq and DNAseq_R2.fastq)	Plain text	samtools, picard-tools <a href="#"><b>Aligners</b></a>	.fq; .fastq
SAM	Define read alignments. Store alignment meta-info (reference, methods, one- or multi-sample).	mapped_reads.sam	Plain text	samtools, picard-tools	.sam
BAM	<a href="#"><b>VISUALIZE ALIGNMENTS (IGV)</b></a> The same as SAM, but compressed and indexed. Also to store UNMAPPED reads (compressed).	mapped_reads.bam unmapped_reads.bam	Binary	samtools, bcftools, picard-tools, <a href="#"><b>IGV</b></a> (Integrative Genome Viewer)	.bam
VCF	<a href="#"><b>SNV &amp; Indels calls</b></a> Indicates genomic variations. Store Variant calling meta-info (reference, methods, one- or multi-sample).	point_variants.vcf	Plain text	bcftools, Unix	.vcf
BED	<a href="#"><b>Intervals</b></a> Delimit genomic regions (i.e. intervals) w or w/o annotations.	targeted_regions.bed intervals.bed	Plain text	bedtools, Unix <a href="#"><b>GATK, picard-tools</b></a>	.bed
TSV or CSV	Create data matrix (rows X Columns)	annotated_variants.tsv	Plain text	Unix, Microsoft Excel, OpenOffice	.tsv; .csv; .txt



# Thanks!

Credits for many class materials to:

Héctor Tejero: [htejero@cnio.es](mailto:htejero@cnio.es)

Elena Piñeiro: [epineiro@cnio.es](mailto:epineiro@cnio.es)

Javier Perales-Patón: [jperales@cnio.es](mailto:jperales@cnio.es)