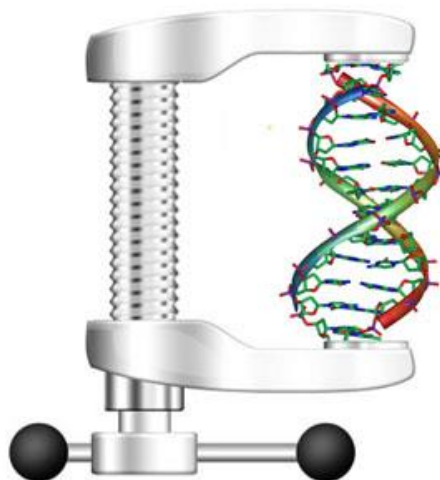
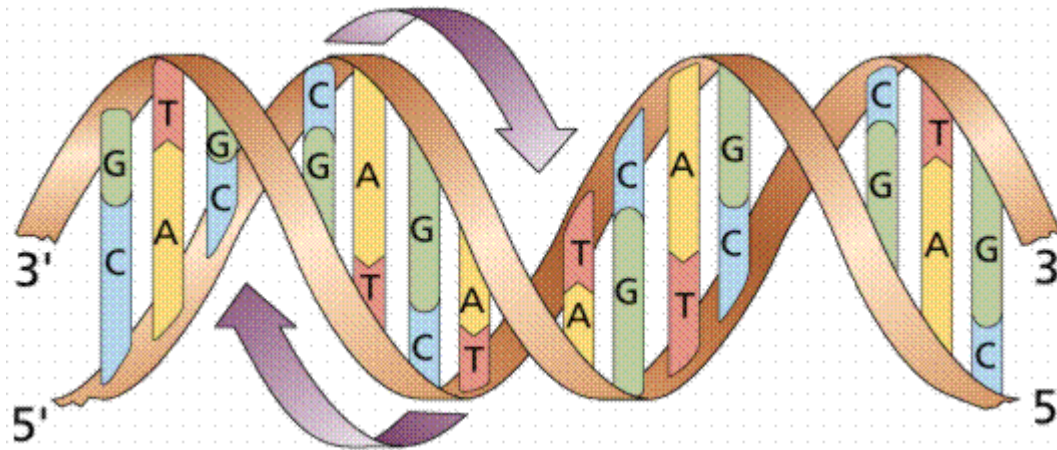


Genome Information Compression



The Human Genome

- 4 possible symbols (**bases**): A,C,G,T
- 3.2 billions base pairs for the human genome
 - $3.2 \times 2 \text{ bit} = 6.4 \text{ billion bits} / 8 = 800 \text{ Mbytes}$



- So why do we need compression?

Genome sequencing

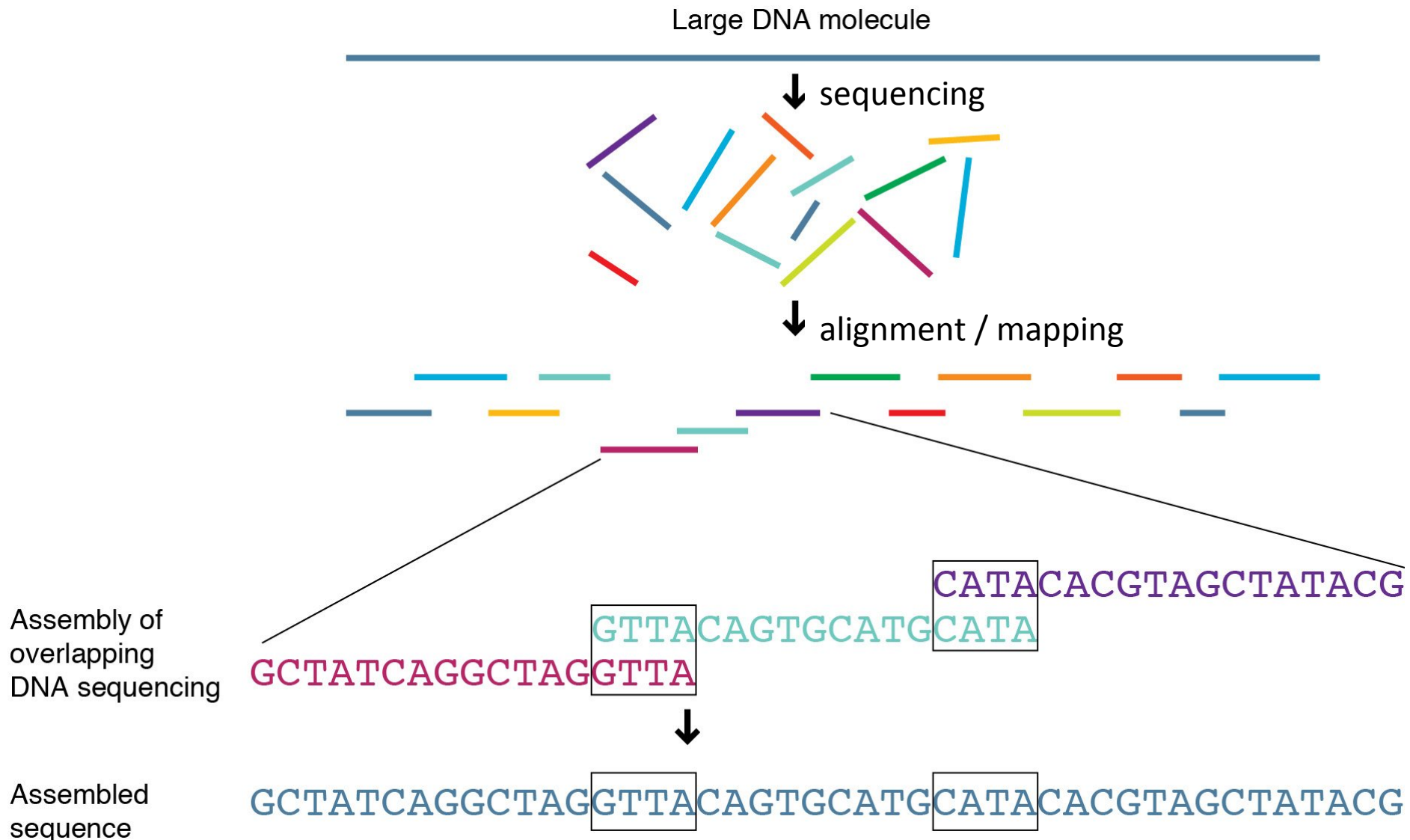
- How do we obtain a genome representation from biological samples?
- Current technology provides random fragments of genome data called “reads”
- This process is called **genome sequencing**



Short and long reads

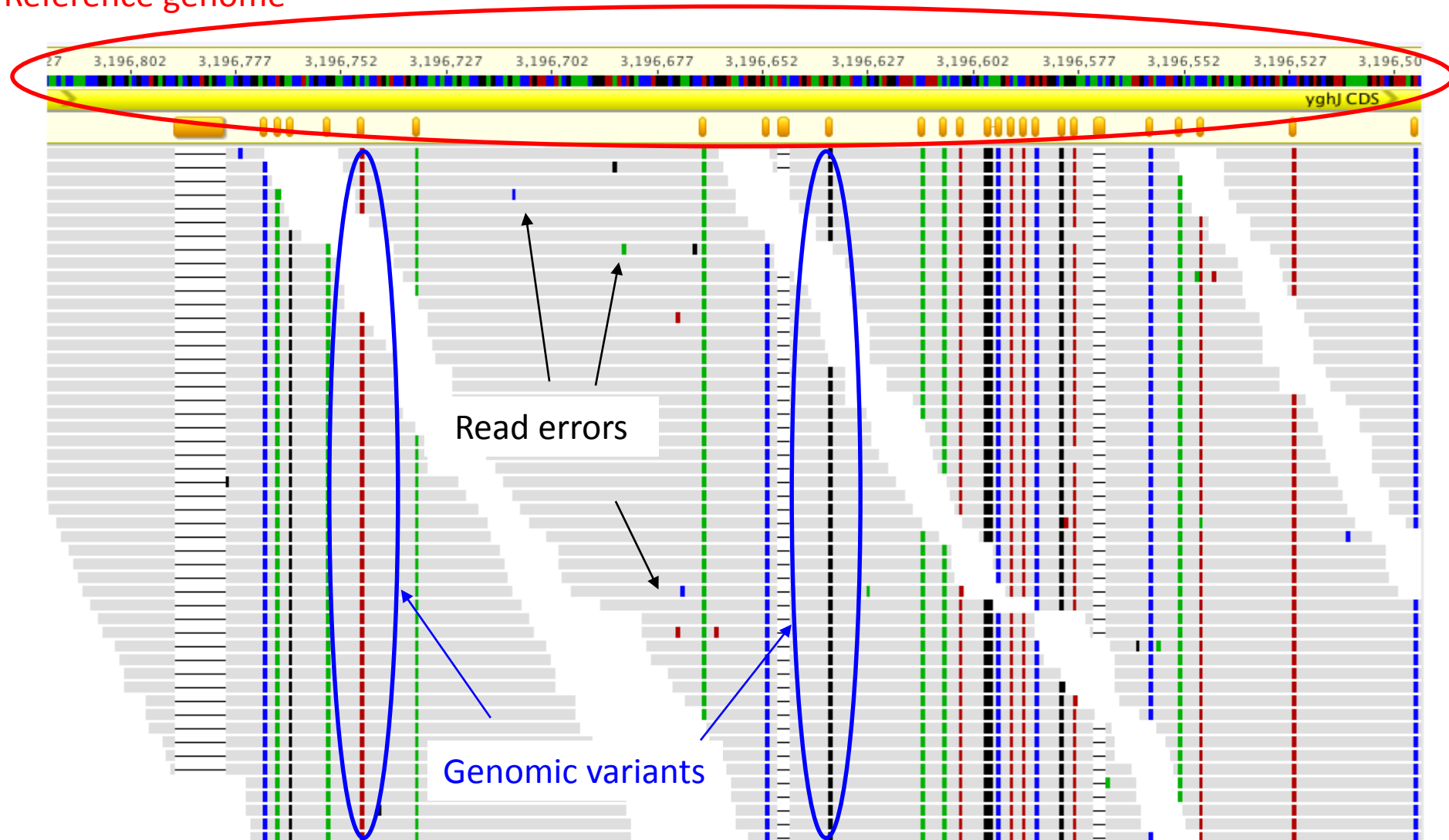
- Different sequencing technologies provide
 - Read lengths from about 100 to 20,000+ bases
 - Single or coupled (paired) reads
 - Different accuracy levels (from 60% to 99.9%)
- Shorter reads are
 - more accurate (up to 99.9%)
 - produced in much larger volumes (currently up to 600 billion bases per single run)

From sequencing to assembly

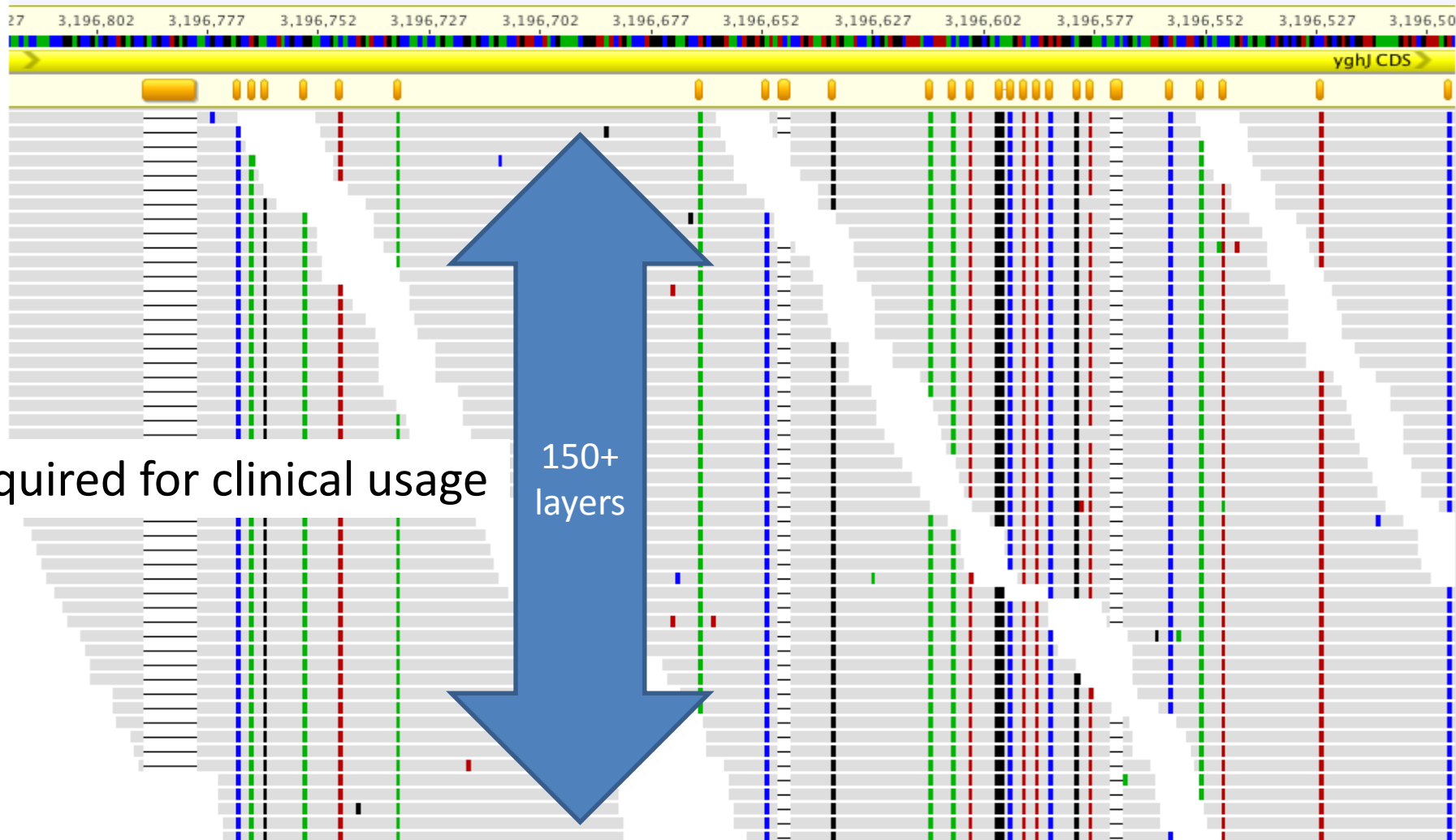


Mapped reads

Reference genome



Coverage



All reads must be preserved

Contig Editor: +225151 SRR006330.470516

Cons 2 Qual 0 Insert Edit Modes >> Cutoffs Undo Next Search Commands >> Settings >> Quit Help >>

<< < > >>

	00100	100110	100120	100130	100140	100150	100160	100170	1	
+330308 SRR006330.2238	ACAGG*CGGG*CACC	TTGCTGGCTG								
+330322 SRR006330.3559	ACAGG*CGGG*CACC	TTGCTGGGCTGCAACAAA								
-330374 SRR006330.4248	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC						
-330389 SRR006330.3045	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTT			
+330414 SRR006330.1334	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330422 SRR006330.2510	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330426 SRR006330.2564	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330435 SRR006330.3770	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAA				
+330440 SRR006465.7152	ACAGG*CGGG*CACC	TTGCTGGG								
-330452 SRR006330.1830	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330454 SRR006330.3586	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330469 SRR006330.2574	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330480 SRR006330.7000	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330481 SR	Is this reading	*CACCCTGCTGTGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
+330488 SR	noise or a	*CACCCTGCTGTGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330491 SR	mutation?	*CACCCTGCTGTGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAA				
+330500 SR		*CACCCTGCTGTGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
+330503 SR		*CACCCTGCTGTGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330510 SRR006330.1472	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330511 SRR006330.1472	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
+330512 SRR006465.7531	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330520 SRR006330.4700	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCA	CAAAC	TTTTTGT	*CGTGCTT	
+330527 SRR006332.7495	AC									
+330528 SRR006332.4165	AC									
-330529 SRR006332.4389	AC									
-330530 SRR006330.4357	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCA	CAAAC	TTTTTGT	*CGTGCTT	
-330532 SRR006332.4943	ACAGG*C									
-330533 SRR006330.2871	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	AATCTTGCTACTGGC	ATGGGCAAT	TAAC	TTTTTGT	*GTGCTG		
>	CONSENSUS	---	ACAGG*CGGG*CACC	TTGCTGGGCTGCAAA	AAACT	GATCTTGCTACTGGC	GTG	GCAATCAAAC	TTTTTGT	*GTGCTG

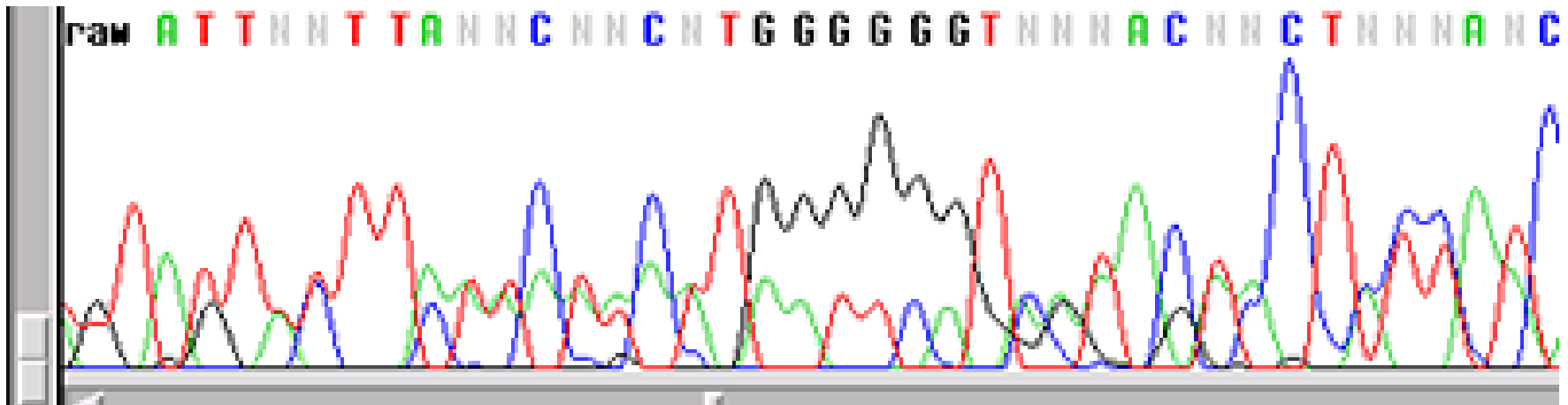
Tag type:SRMr Direction:- Comment:"Strong Repeat Marker base found by MIRA"

Data volumes

- 3.2 billion per genome
 - X 200+ for clinical usage (to get at least 150 layers)
- Additional information
 - 1 quality score per each base (up to 96 possible levels)

Reads quality scores

- Each base call in a read has a level of confidence

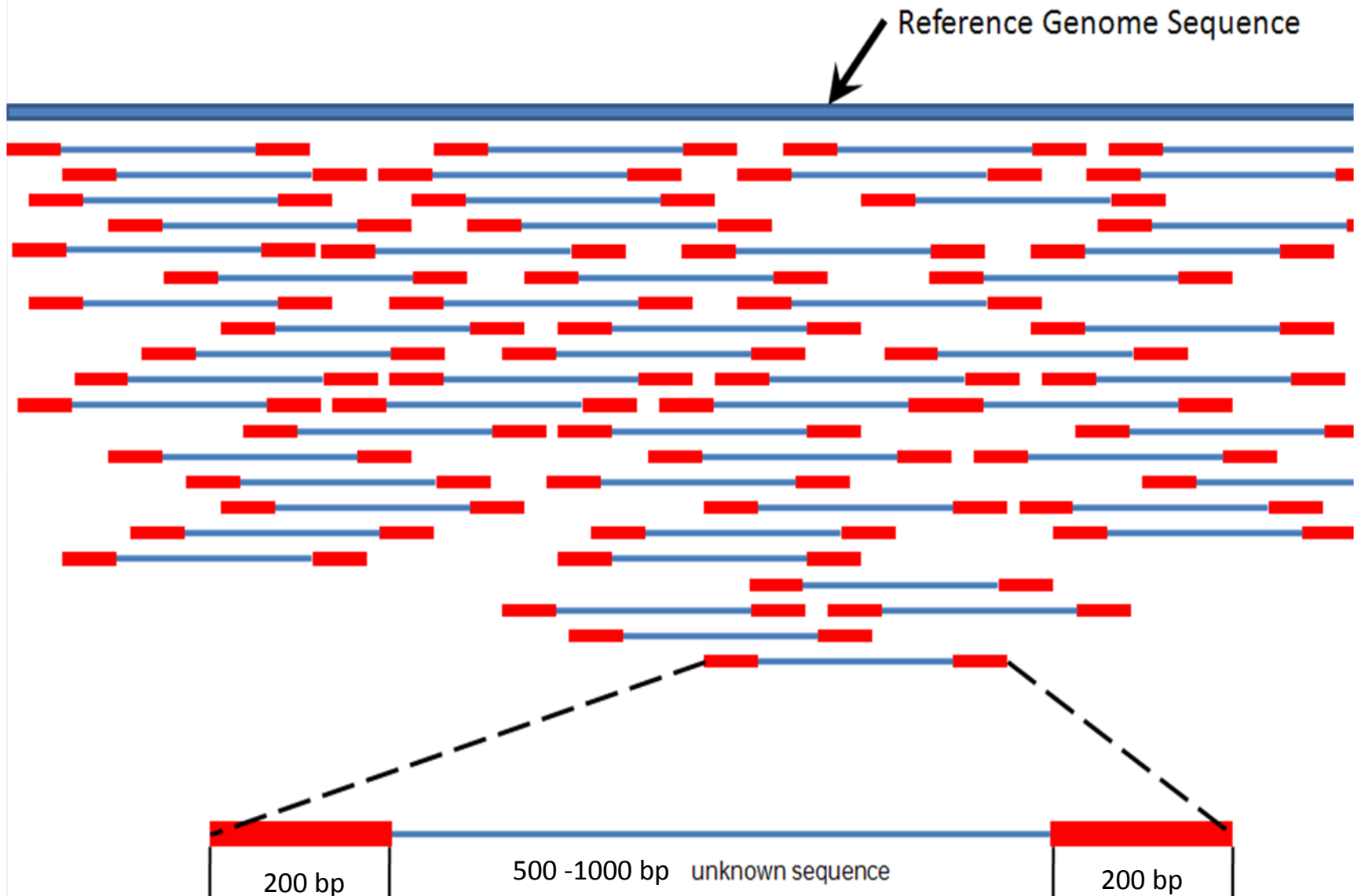


- The level of confidence is expressed as “quality score” in a range that is machine dependent represented as ASCII character.

Data volumes

- 3.2 billion per genome
 - X 200+ for clinical usage (to get at least 150 layers)
- Additional information
 - 1 quality score (ASCII char) per each base
 - pairing information for coupled reads (labelling)

Paired reads

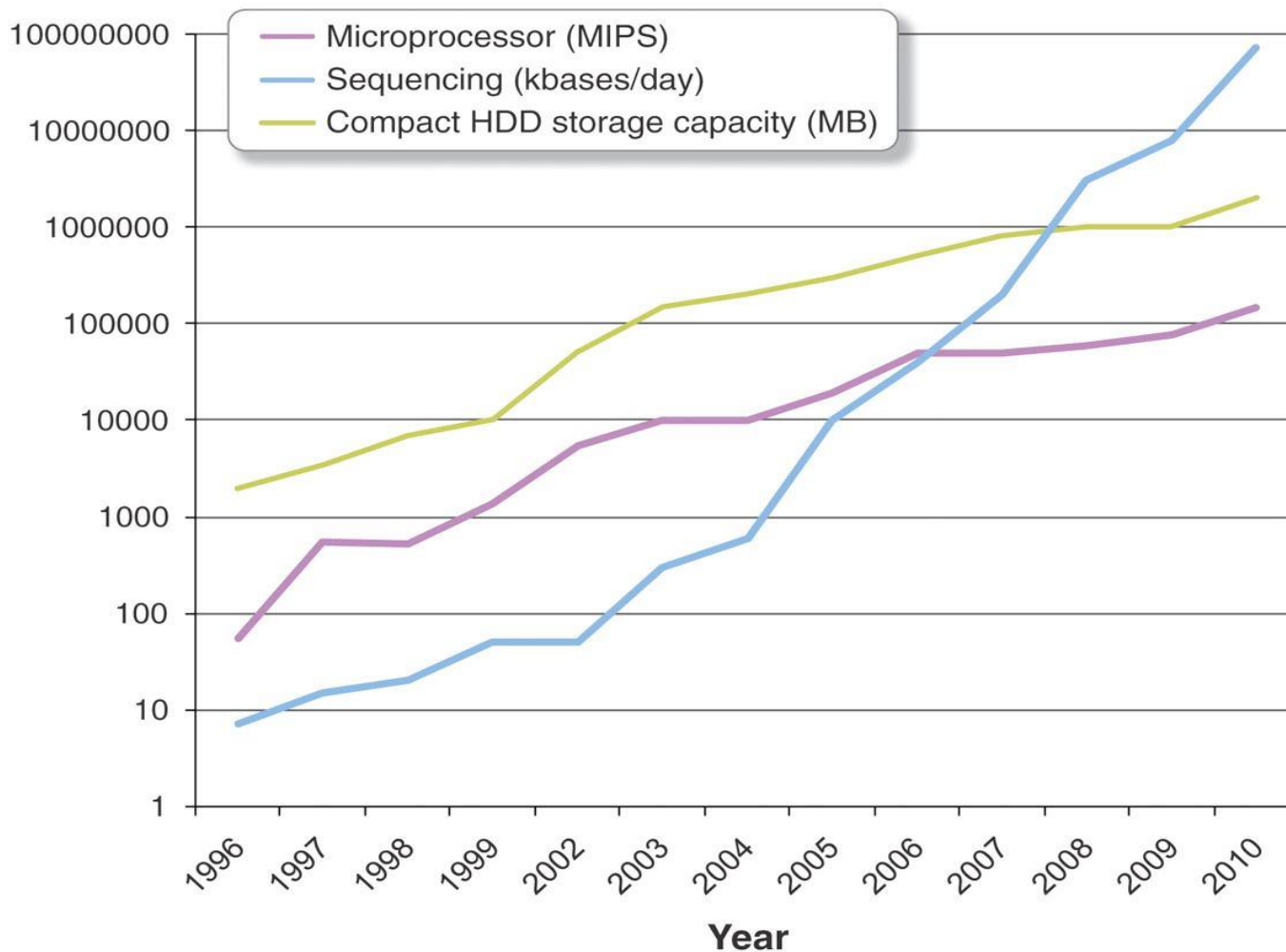


Data volumes

- 3.2 billion per genome
 - X 200+ for clinical usage (to get at least 150 layers)
- Additional information
 - 1 quality score per each base (up to 96 possible levels)
 - pairing information for coupled reads (labelling)
- Total = $3.2\text{GB} \times 200 \times 2 \times \text{labelling} \approx 1.5 \text{ TB}$
 - Labelling ≈ 1.15

Sequencing Progress vs Compute and Storage

Moore's and Kryder's Laws fall far behind



Raw data format

FASTQ	Field	FASTA
@HWUSI-EAS100R:6:73:941:197 3#0/1	<i>Header (Unique ID plus other information). Only the first character is standard.</i>	>HWUSI-EAS100R:6:73:941:197 3#0/1
GATTTGGGGT.....	<i>Nucleotides sequence</i>	GATTTGGGGT.....
+SRR001666.1 071112_SLXA-EAS1_s_7	<i>Optional description. Only the first character is standard. This is becoming obsolete</i>	Not present
!"*((((**+)	<i>Quality scores</i>	Not present

Example

One read:

HS2000-1240_45:1:1234:6966:12500

AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAA

+

CCCCFFFFFHHHHHHIJIJIJIJIJJFHGGIFHHHHIJJJIJIHIJJFIJJIIJJJJHIIIGIIIIIIIFHGGGHFFFFFF?BBEECDD?CCCECDD?AC;5@

Example

One read:

HS2000-1240_45:1:1234:6966:12500/1

AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTGTAAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAA

+

CCCCFFFFFHHHHHIJIIJIIJIIJFHHGGIFHHHHIIJJJIJJIHIIJFIJJIIJJJHIIIIIGIIIIIIIFHGGGHFFFFFF?BBEECDD?CCCECDD?AC;5@

HS2000-1240_45:1:1234:6966:12500/2

ATCAGATGTATAATTTGCAAATAGTTTCTCTCATTCTTTTTTTTTTTTTTTTTTTAGACAGGGTCTCACTGTATTGCCCAGGCTGGAGTGCAGTGGTGCAATC

+

CCCCFFFFFHHHHHJJJJJJJJIIJJJIJJJIJJJIJJJIJJJJJJJHFD@#####

Example

- A read aligned onto a section of the chromosome 1 of the reference genome.

- ```

• Chr1 39999220 39999240 39999260 39999280 39999300 39999320
• 39999340
•
• |
•
• AATGGTGCAGTCACAGCTGTCTACAAAAATATTTTTTAAATATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCAGTTGGGTACAGGAGTTCCAAGACTGTGGTGAG
•
•AAATATTTTTTAAATATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCAGTTGGGTACAGGAGTTCCAA

```

# Example

- A read aligned onto a section of the chromosome 1 of the reference genome.

- Chr1 39999220 39999240 39999260 39999280 39999300 39999320  
39999340  
• . | . | . | . | . | . | . | . | . | . | . | . | .  
• |  
•  
• AATGGTGCAGTCACAGCTGTCTACAAAAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAAGACTGTGGTGAG  
• .....AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAA.....

- A second read, known to be associated with the first read.

- [illegible]

- A read aligned onto a section of the chromosome 1 of the reference genome.

- A second read, known to be associated with the first read.

EPFL SCI-STI-MM

# Example

- [illegible]

Note that the sequence obtained may not reflect exactly the genomic sequence.

- A read aligned onto a section of the chromosome 1 of the reference genome.

[illegible]

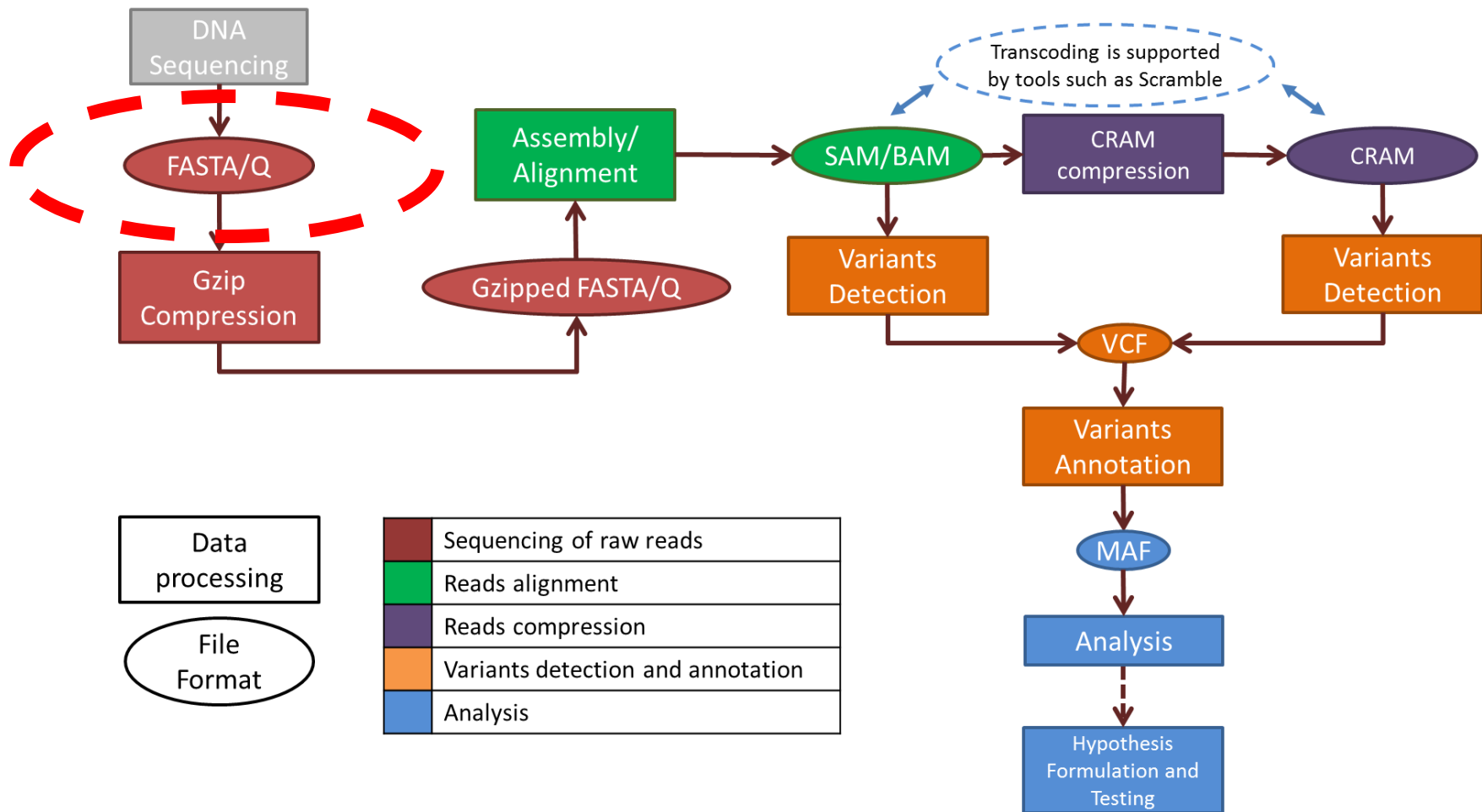
- A second read, known to be associated with the first read.

[illegible]

Here, the quality of the sequence obtained, associated to each nucleotide has been flagged as poor in the region where it does not match well the reference genome.



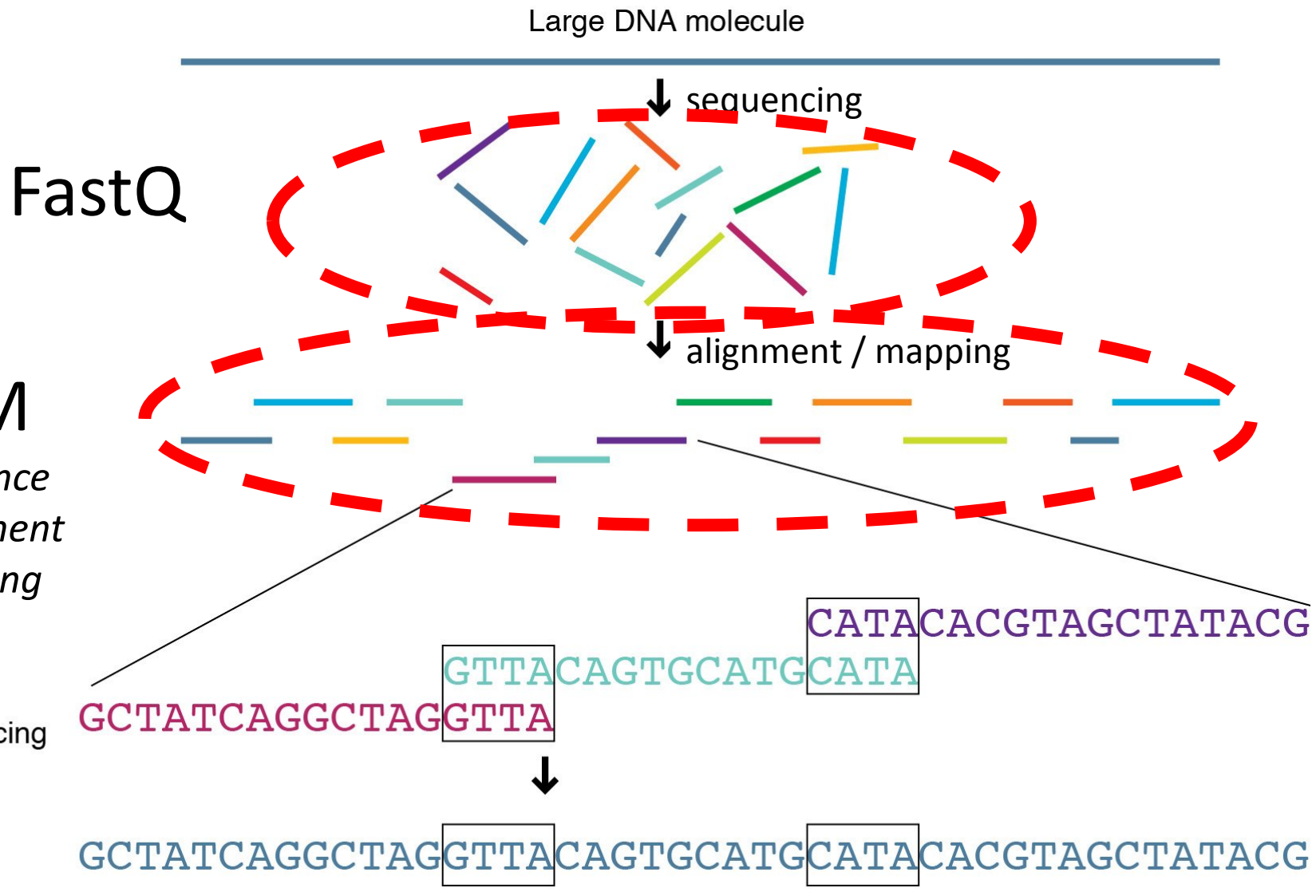
# Genome processing pipeline



# FastQ compression today

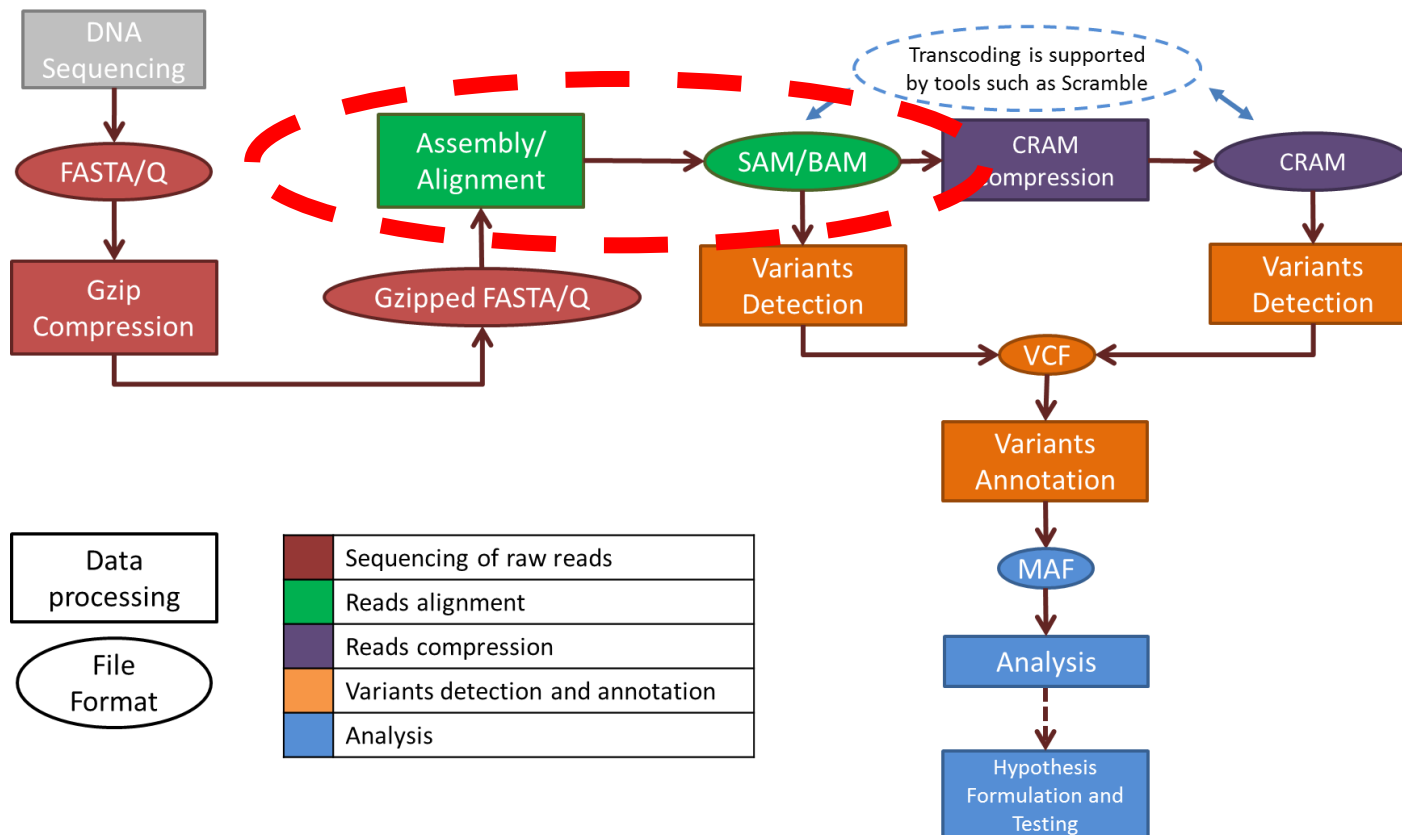
- Gzip of the entire txt file (sometimes split into several files)
- Compression ratio : 3 to 5
- According to the coverage 1 genome can take up to 2 TB

# From sequencing to assembly



# Alignment / mapping

- Raw data + alignment information



Coord  
ref 12345678901234 5678901234567890123456789012345

AGCATGTTAGATAA\*\*GATAGCTGTGCTAGTAGGCAGTCAGCGGCAT

Position index  
reference genome

+r001/1  
+r002  
+r003  
+r004  
-r003  
-r001/2

TTAGATAAAGGATA\*CTG  
aaaAGATAA\*GGATA  
gcctaAGCTAA  
ATAGCT.....TCAGC  
ttagctTAGGC  
CAGCGGCAT

FastQ reads

Paired reads

FastQ headers

Positions

Indels

Base sequences

SAM

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

r001 163 ref 7 30 6M2I4M1D3M = 37 39 TTAGATAAAGGATACTC \*  
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;  
r004 0 ref 16 30 6M14N5M \* 0 0 ATAGCTTCAGC \*  
r003 2064 ref 29 17 6H5M \* 0 0 TAGGC \* SA:Z:ref,9,+,5S6M,30,1;  
r001 83 ref 37 30 9M = 7 -39 CAGCGGCAT \* NM:1:1

FastQ Headers

Quality scores if present

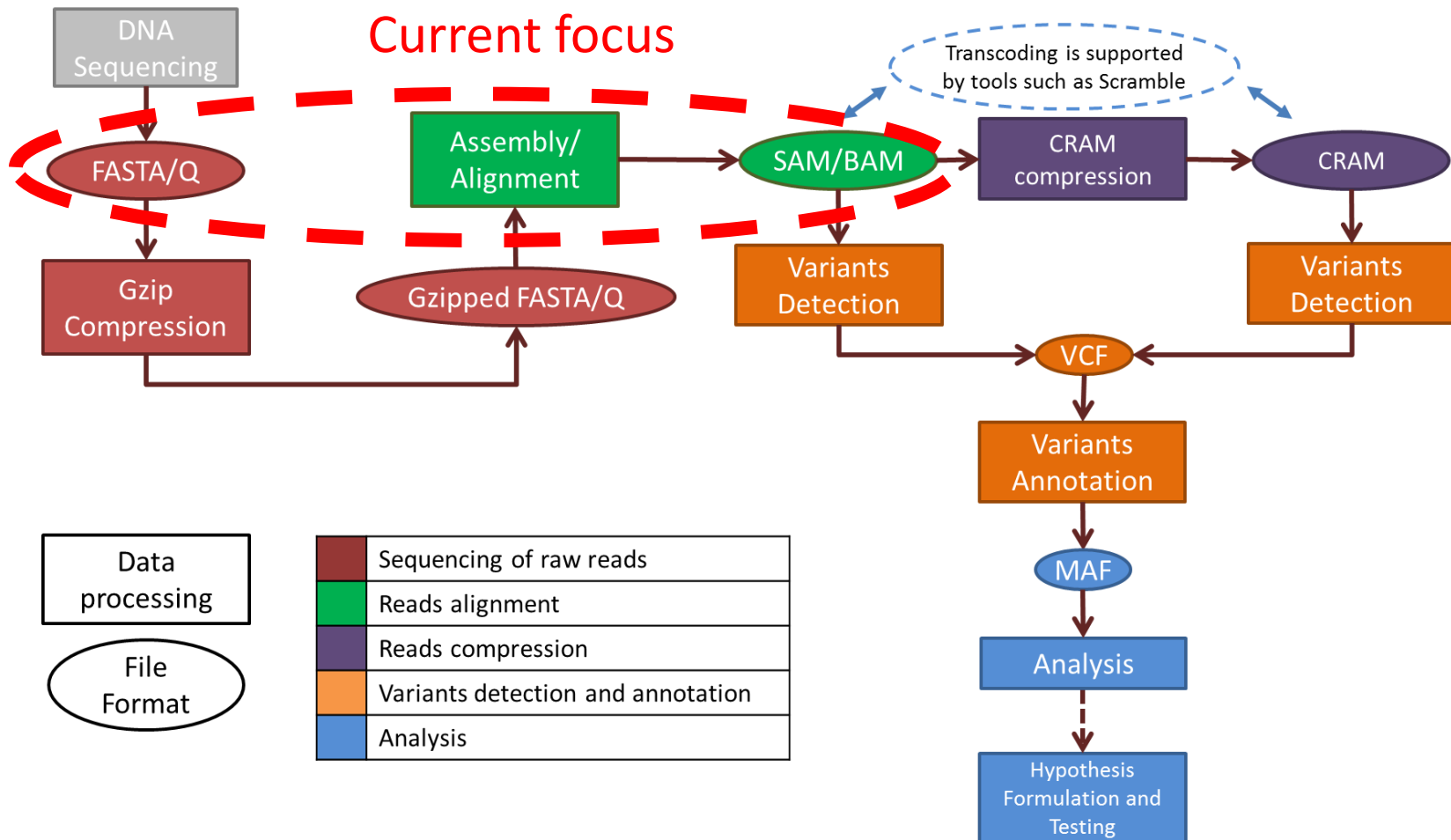
# Compressed SAM = BAM

- BAM = Block based zipped SAM
- Indexable for random access
- Compression ratio over textual SAM  $\sim$  4 to 6
- Example (coverage 200x)
  - Fastq = 1.5 TB
    - Fastq.gz = 370 GB
    - Fastq.bzip = 255 GB
    - quip = 205 GB (best compression tool for FastQ)
  - SAM =  $\sim$ 3 TB
  - BAM = 500 GB

# SAM/BAM view demo

- Human sample from the MPEG dataset
  - /human/illumina/LowCoverage/NA21144.chrom11
- Chromosome 11
  - Reads length: 100 bases
  - No. of reads: ~10.1 millions
  - Unaligned data: 2.2 GB
  - Aligned SAM: 4.5 GB
  - Compressed BAM: 1 GB

# Raw sequence data + Aligned data



# Thank you

