

**INTERNATIONAL ORGANISATION FOR STANDARDISATION  
ORGANISATION INTERNATIONALE DE NORMALISATION  
ISO/IEC JTC 1/SC 29/WG 11  
CODING OF MOVING PICTURES AND AUDIO**

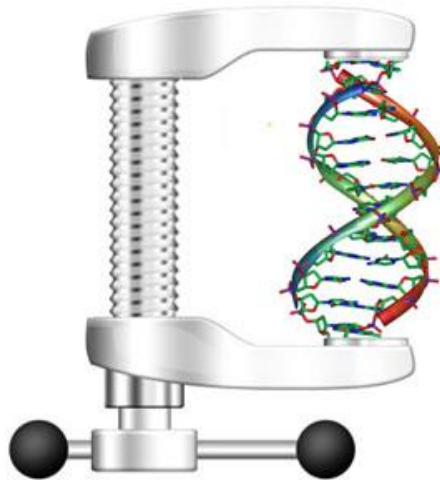
**ISO/IEC JTC 1/SC 29/WG 11 N15527  
Warsaw, CH – June 2015**

<b>Source:</b>	Requirements
<b>Authors:</b>	Claudio Alberti, Marco Mattavelli (EPFL)
<b>Title:</b>	Genome Compression 101 - Tutorial on Genome Compression and Storage

## **1 Introduction**

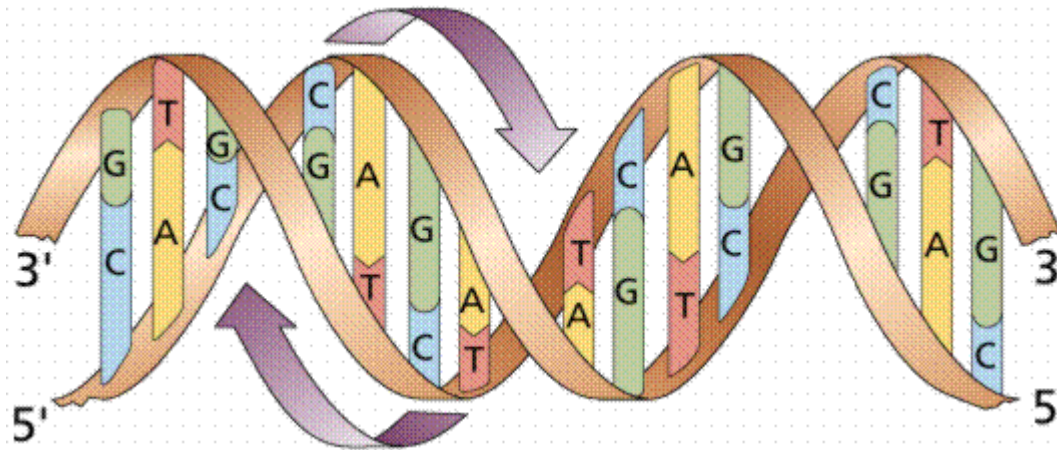
This slides set was presented during the 111<sup>th</sup> MPEG meeting held in Geneva in February 2015 and published at the 112<sup>th</sup> MPEG meeting in Warsaw (June 2015). The aim is to provide an introduction to the current status of genomic information presentation and storage focusing on the sequencing and alignment stages of the genome information life cycle.

# Genome Information Compression



# The Human Genome

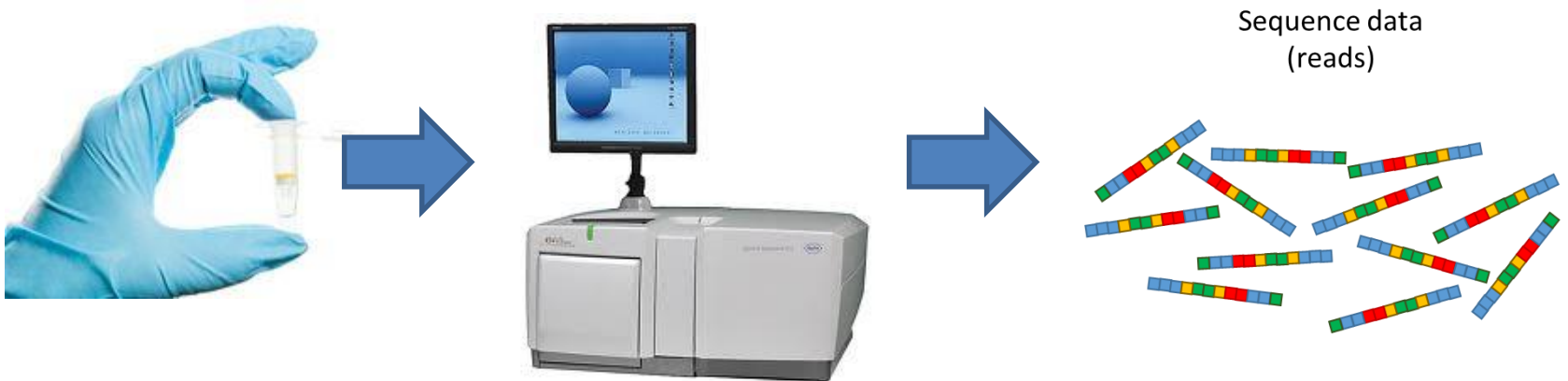
- 4 possible symbols (**bases**): A,C,G,T
- 3.2 billions base pairs for the human genome
  - $3.2 \times 2 \text{ bit} = 6.4 \text{ billion bits} / 8 = 800 \text{ Mbytes}$



- So why do we need compression?

# Genome sequencing

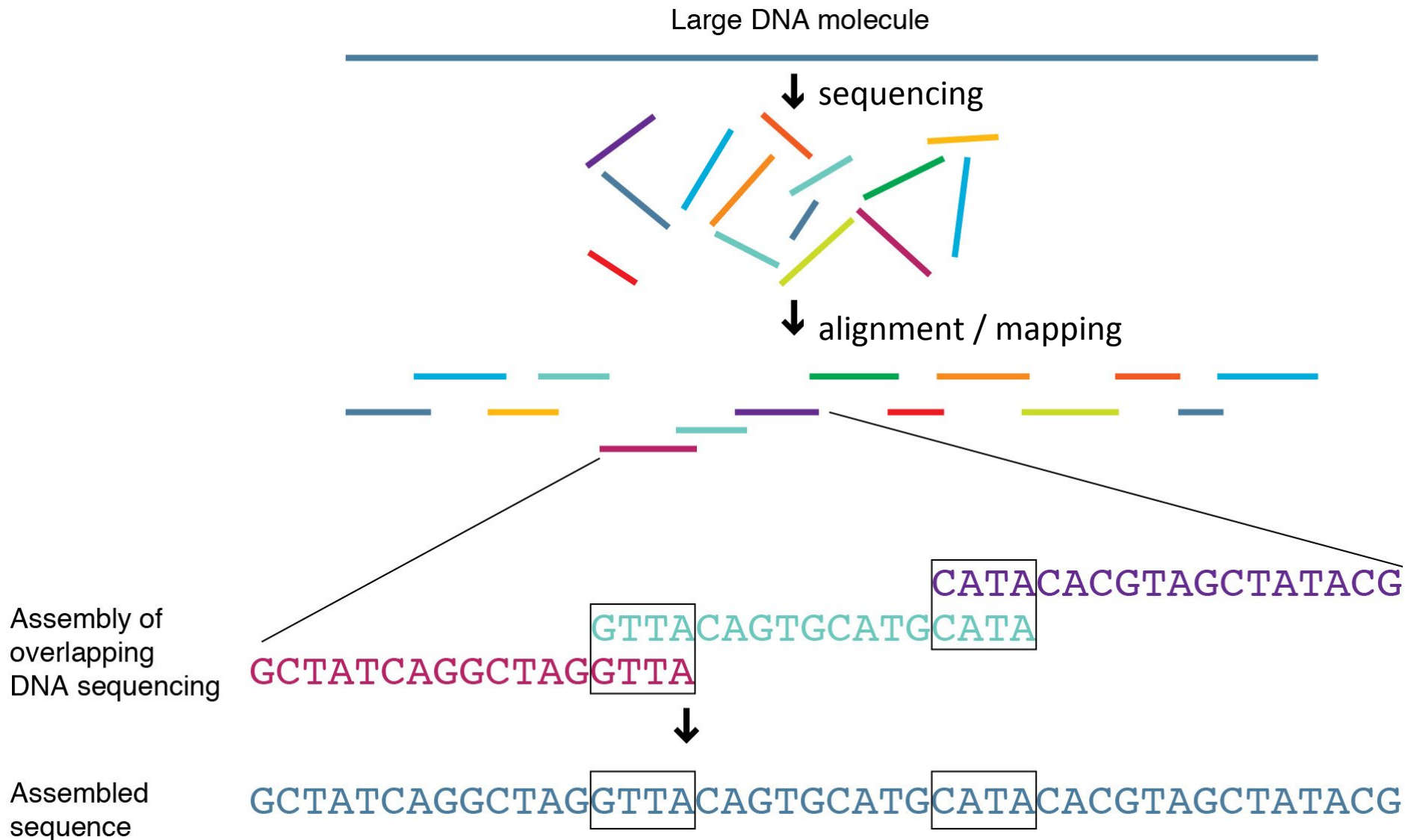
- How do we obtain a genome representation from biological samples?
- Current technology provides random fragments of genome data called “reads”
- This process is called **genome sequencing**



# Short and long reads

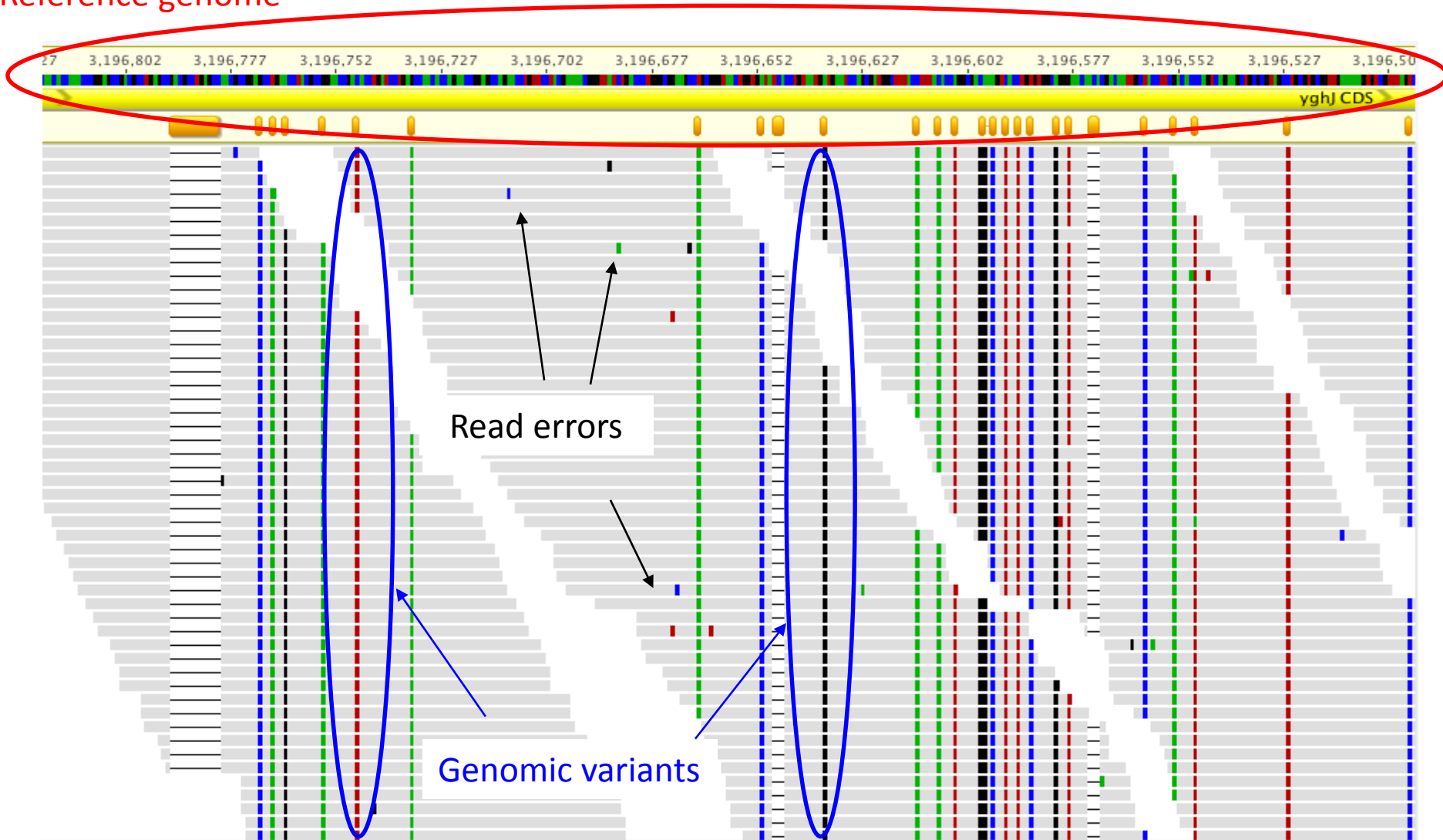
- Different sequencing technologies provide
  - Read lengths from about 100 to 20,000+ bases
  - Single or coupled (paired) reads
  - Different accuracy levels (from 60% to 99.9%)
- Shorter reads are
  - more accurate (up to 99.9%)
  - produced in much larger volumes (currently up to 600 billion bases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

# From sequencing to assembly

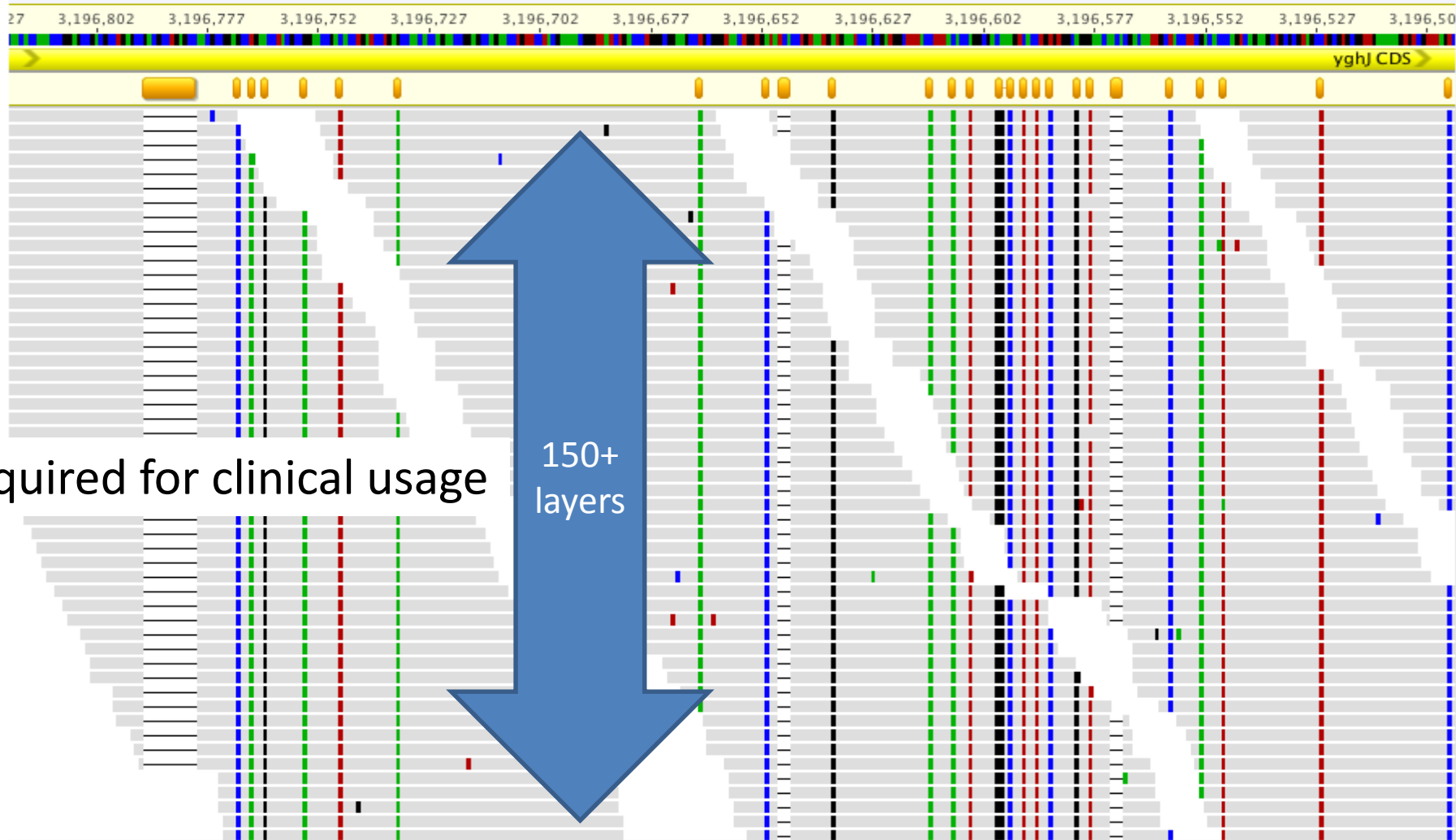


# Mapped reads

Reference genome



# Coverage



Required for clinical usage

150+  
layers



# 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	TGCTGG	GCTG								
+330322 SRR006330.3559	ACAGG*CGGG*CACC	TGCTGG	GCTGCAACAAA								
-330374 SRR006330.4248	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC						
-330389 SRR006330.3045	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTT			
+330414 SRR006330.1334	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330422 SRR006330.2510	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330426 SRR006330.2564	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330435 SRR006330.3770	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAA				
+330440 SRR006465.7152	ACAGG*CGGG*CACC	TGCTGGG									
-330452 SRR006330.1830	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330454 SRR006330.3586	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330469 SRR006330.2574	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT		
-330480 SRR006330.7000	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330481 SR	Is this reading	€CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG	
+330488 SR	noise or a	€CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG	
-330491 SR	mutation?	€CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAA			
+330500 SR		€CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT	
+330503 SR		€CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*CGTGCTT	
-330510 SRR006330.0270	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330511 SRR006330.1472	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
+330512 SRR006465.7531	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
-330520 SRR006330.4700	ACAGG*CGGG*CACC	TT	CTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAA	CAAAC	TTTTTGT	*CGTGCTT
+330527 SRR006332.7495	AC										
+330528 SRR006332.4165	AC										
-330529 SRR006332.4389	AC										
-330530 SRR006330.4357	ACAGG*CGGG*CACC	TT	CTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAA	CAAAC	TTTTTGT	*CGTGCTT
-330532 SRR006332.4943	ACAGG*C										
-330533 SRR006330.2871	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	ATGGCAAT	TAAC	TTTTTGT	*GTGCTG		
>	CONSENSUS	***	ACAGG*CGGG*CACC	TGCTGG	GCTGCAAA	AAACT	GATCTTGTCACTGGC	GTG	GCAATCAAAC	TTTTTGT	*GTGCTG

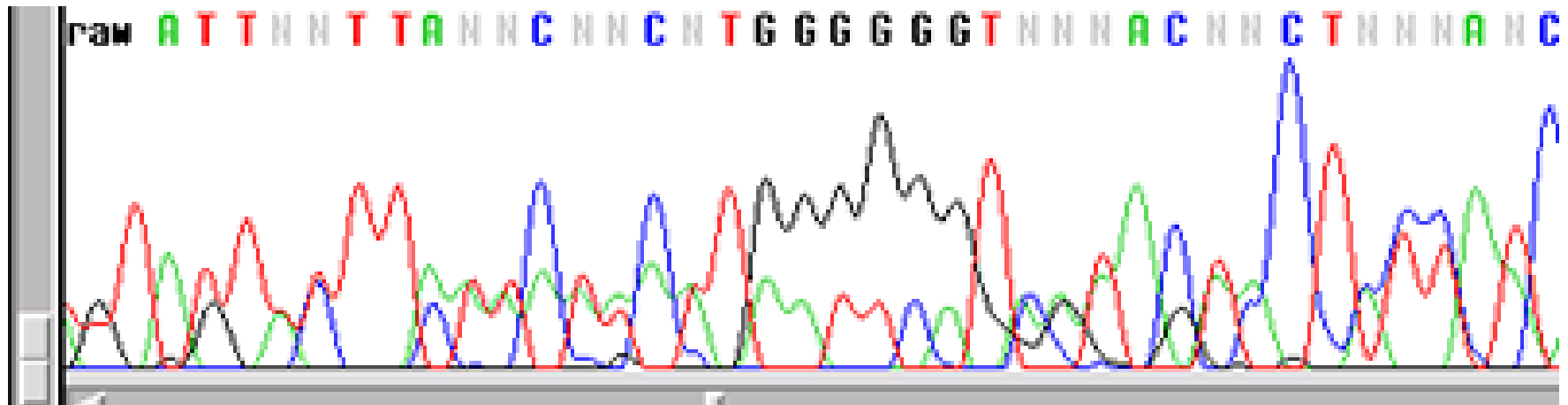
Tag type:SRMr Direction:- Comment:"Strong Repeat Marker base found by MIRA"

# 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

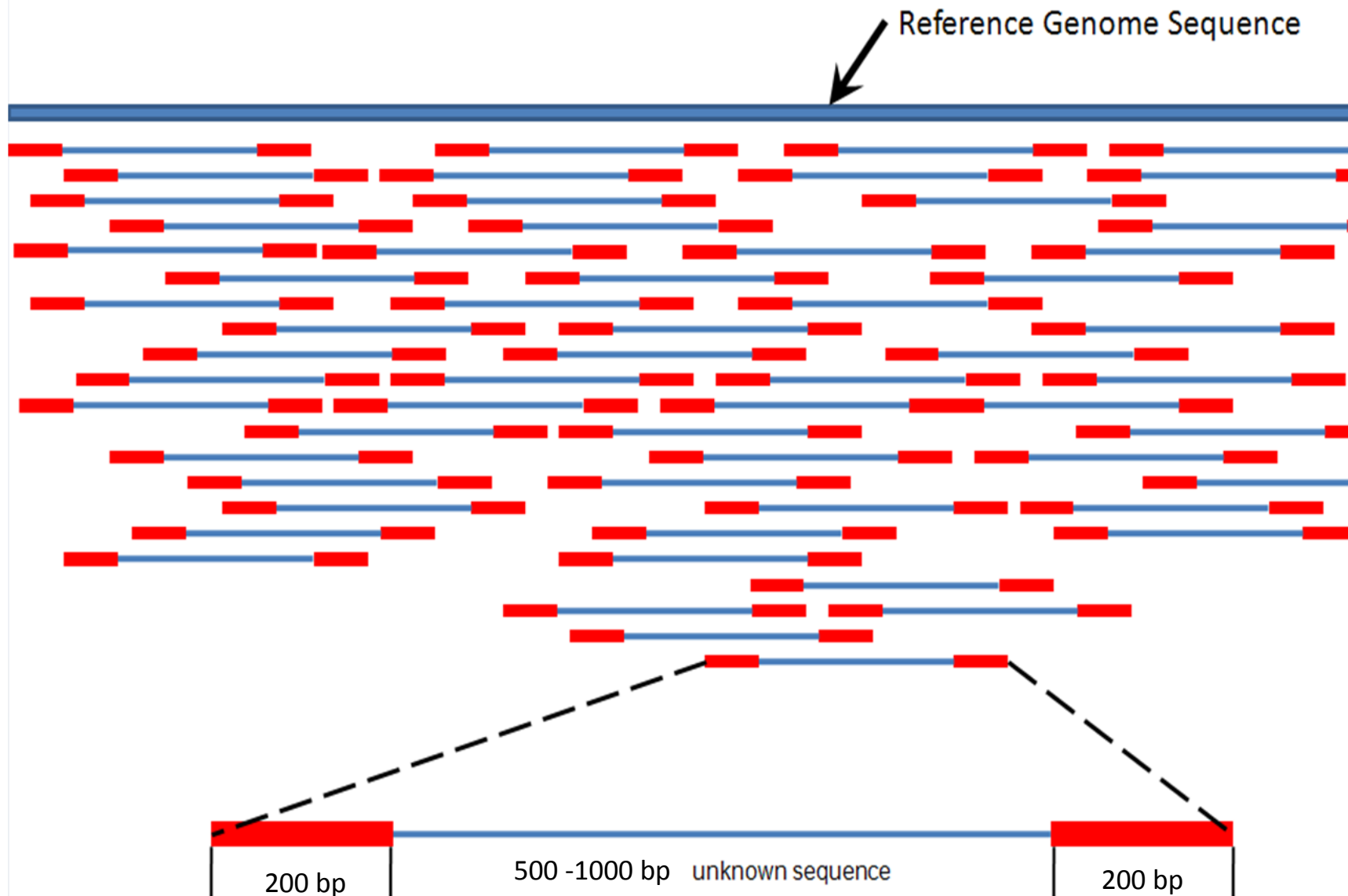


- The level of confidence is expressed as “quality score” in a range that is machine dependent represented as ASCII character.

# 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

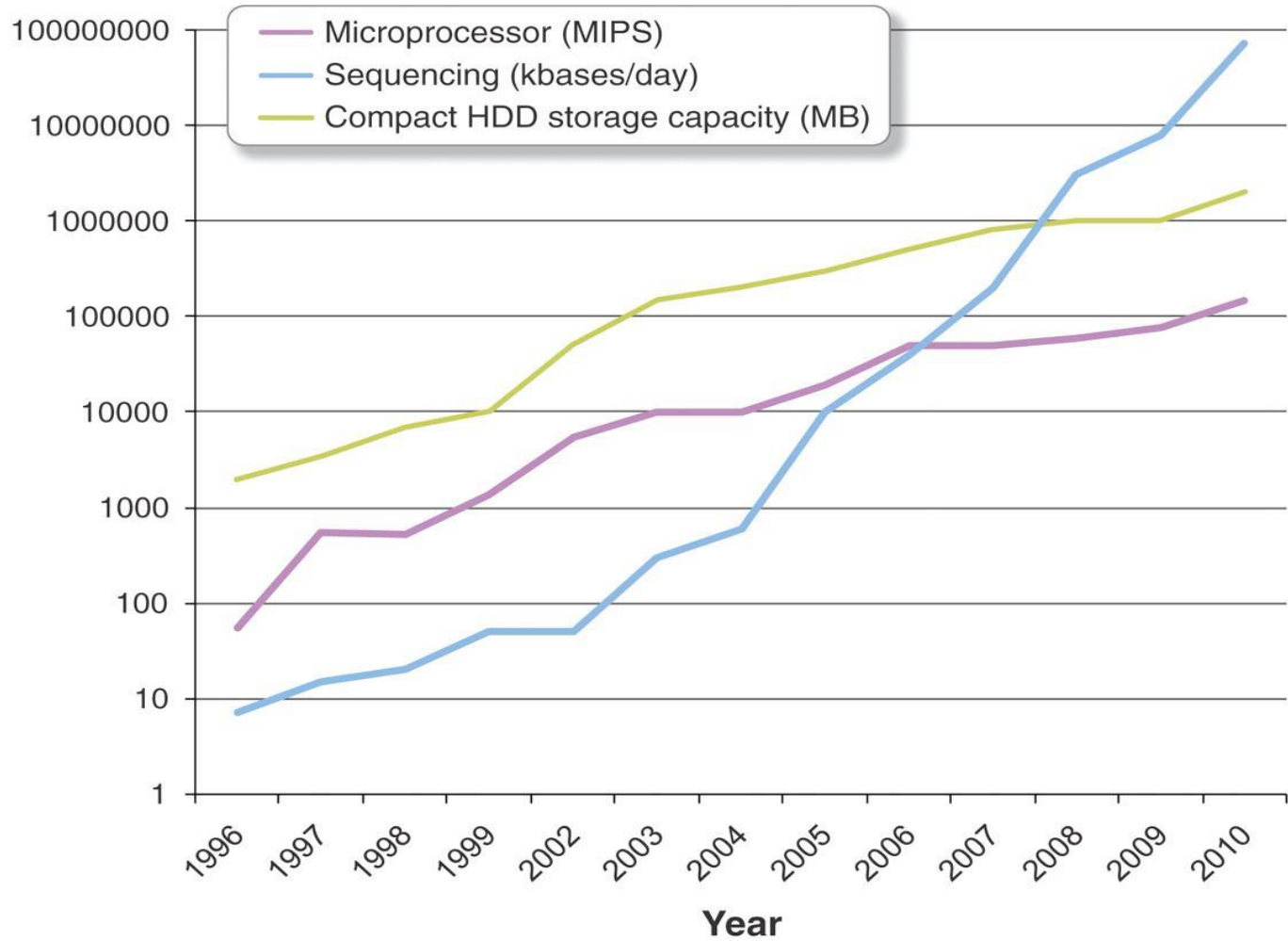


# Data volumes

- 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 = 3.2GB x 200 x 2 x labelling  $\approx$  1.5 TB
  - Labelling  $\approx$  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

AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTGTGAAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAA

+

CCCCFFFFFHHHHHIJIIJIIJIIJIIJFHHGGIFHHHHIIJJIJJIHIJJIJLI

# Example

One read:

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

AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGTGCCTATAGTTCCAAGTTGTAAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAA  
+

CCCCFFFFHHHHHIIJIIJIIJIIJFHHGGIFHHHIIJJJIJIIHIJJFIJJIIJJJJHIIIIIGIIIIIIIFHGGGHFFFFFF?BBEECDD?CCCECDD?AC;5@

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

ATCAGATGTATAATTTGCAAAATAGTTTCTCTCATTCTTTTTTTTTTTTTTTTTTTTATAGACAGGGTCTCACTGTATTGCCCAGGCTGGAGTGCAGTGGTGCAATC  
+

CCCCFFFFHHHHHJJJJJJJIIJJJIJJJIJJJIJJJIJJJIJJJJJJJHFDD@#####

# Example

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

```

• Chr1      39999220      39999240      39999260      39999280      39999300      39999320
• |
•
• AATGGTGCAGTCACAGCTGTCTACAAAAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCAAGACTGTGGTGAG
• .....AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGCCTATAGTTCCAAGTGTGTAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCA .....

```



# Example

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

- Chr1          39999220                  39999240                  39999260                  39999280                  39999300                  39999320  
39999340  
|  
|  
AATGGTGAGCTGATGGCACCCTGCACTCCAGCCTGGGCAATACAGTGAGACCCCTGTCTAAAAAAAAAAAAAAAAAGAATGAGAGAAACTATTTGCAAATTATACATCTGATAAGAGACCTGTAT CTA  
.....AAATATTTTTTAAAATTAGCCAGGTGTGGTGGTGTGCCTATAGTTCCAACTAGTTGTAAAGCTGAAACATAAGGACCACTTGGGTACAGGAGTTCCA .....  
|

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

- 39999340          39999360          39999380          39999400          39999420          39999440          39999460  
|  
|  
GTGGTGAGCTGATGGCACCCTGCACTCCAGCCTGGGCAATACAGTGAGACCCCTGTCTAAAAAAAAAAAAAAAAAGAATGAGAGAAACTATTTGCAAATTATACATCTGATAAGAGACCTGTAT CTA  
.....GTTTTGCCCCCTCACCCAGCCTGGGGAATAAAGGGGGCCCTTTCTAAAAAAAAAAAAAAAAAGAATGAGAGAAACTATTTGCAAATTATACATCTGAT .....  
|  
HS2000-1240\_45:1:1234:6966:12500/2  
ATCAGATGTATAATTTGCAAATAGTTTCTCTCATCTTTTTTTTTTTTTTTTTTAGACAGGGTCTCACTGTATTGCCAGGCTGGAGTGCAGTGGTGAATC

Note that the original sequence of the second read had to be reverse complemented to match the genome.

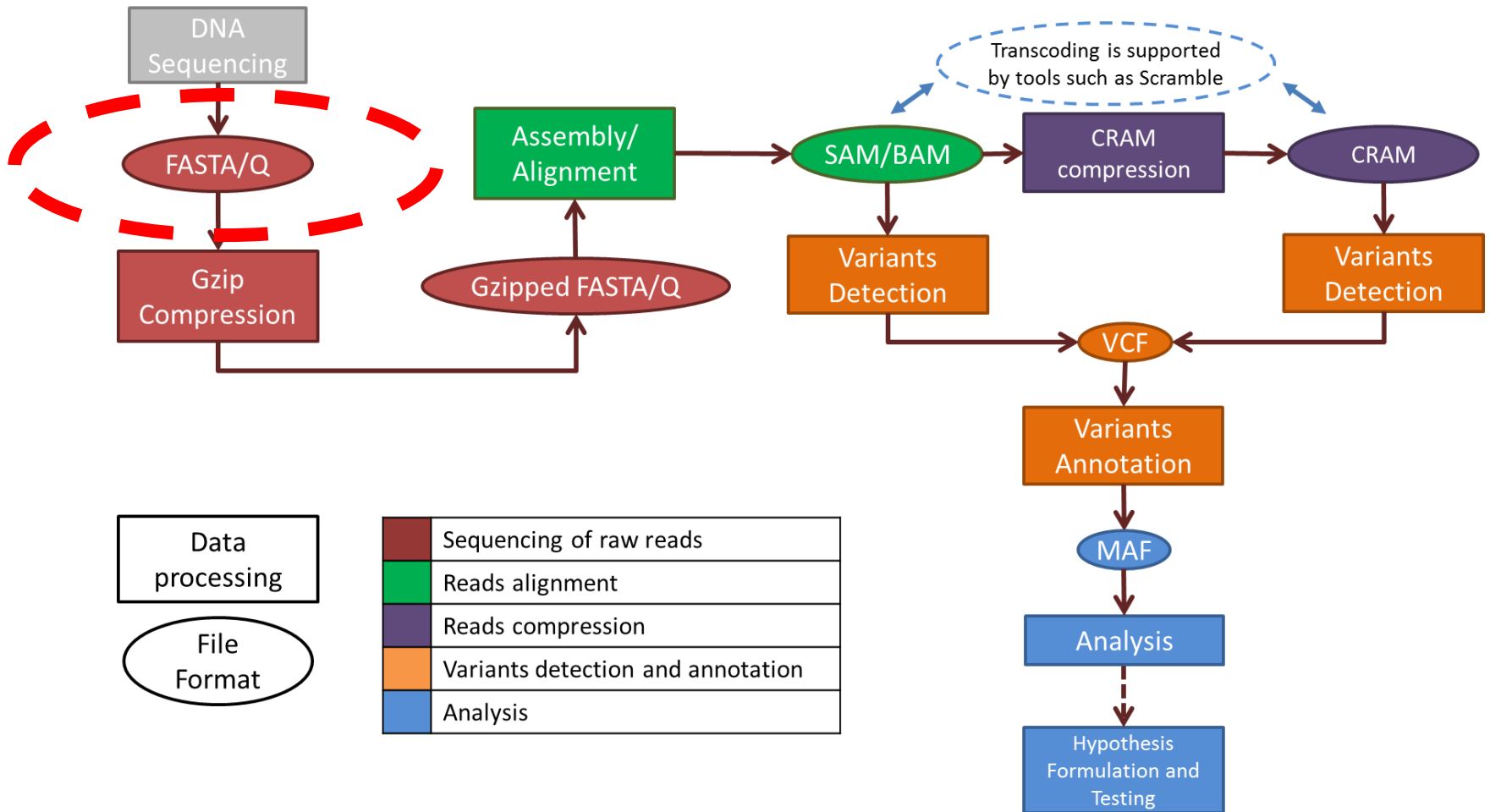








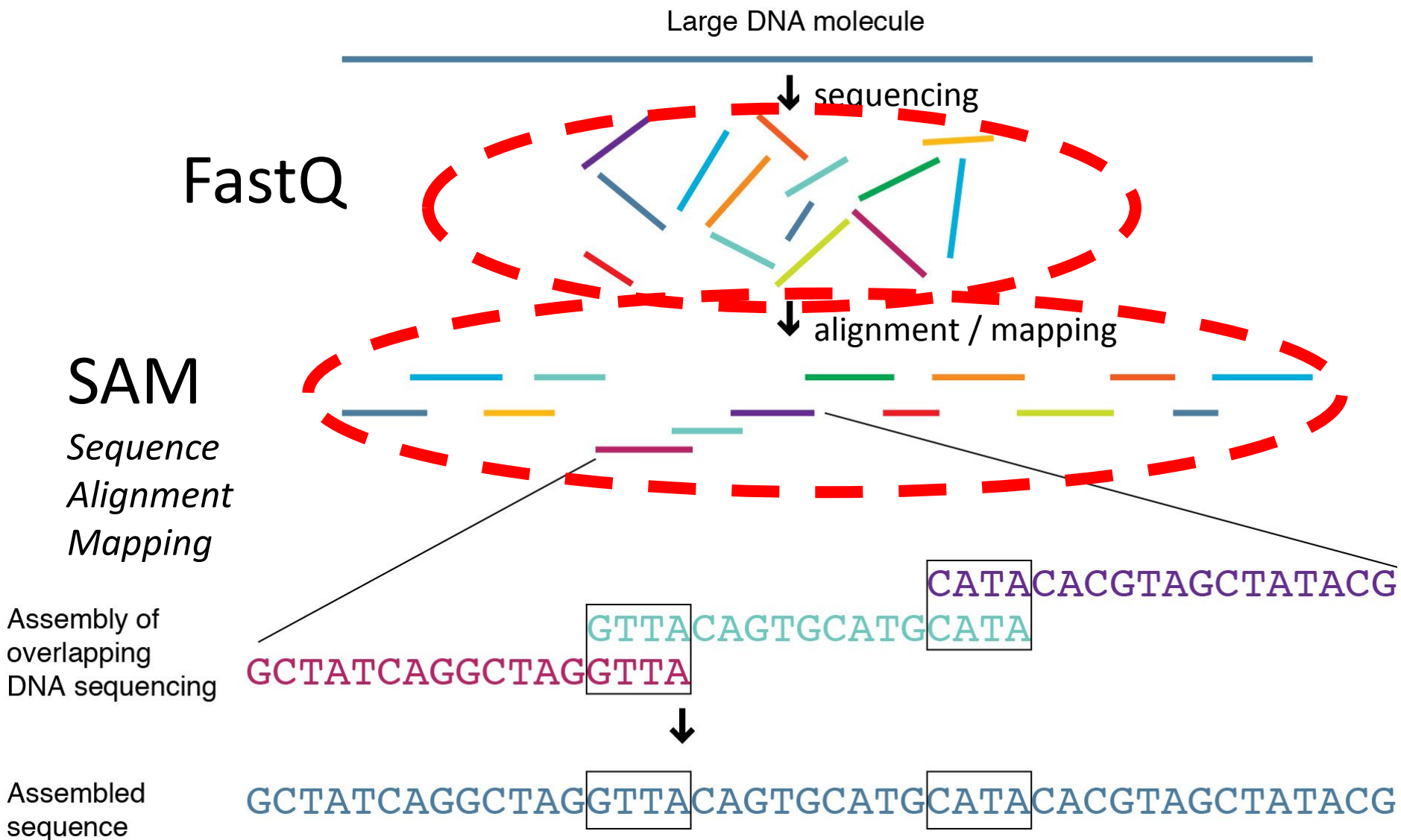
# Genome processing pipeline



# FastQ compression today

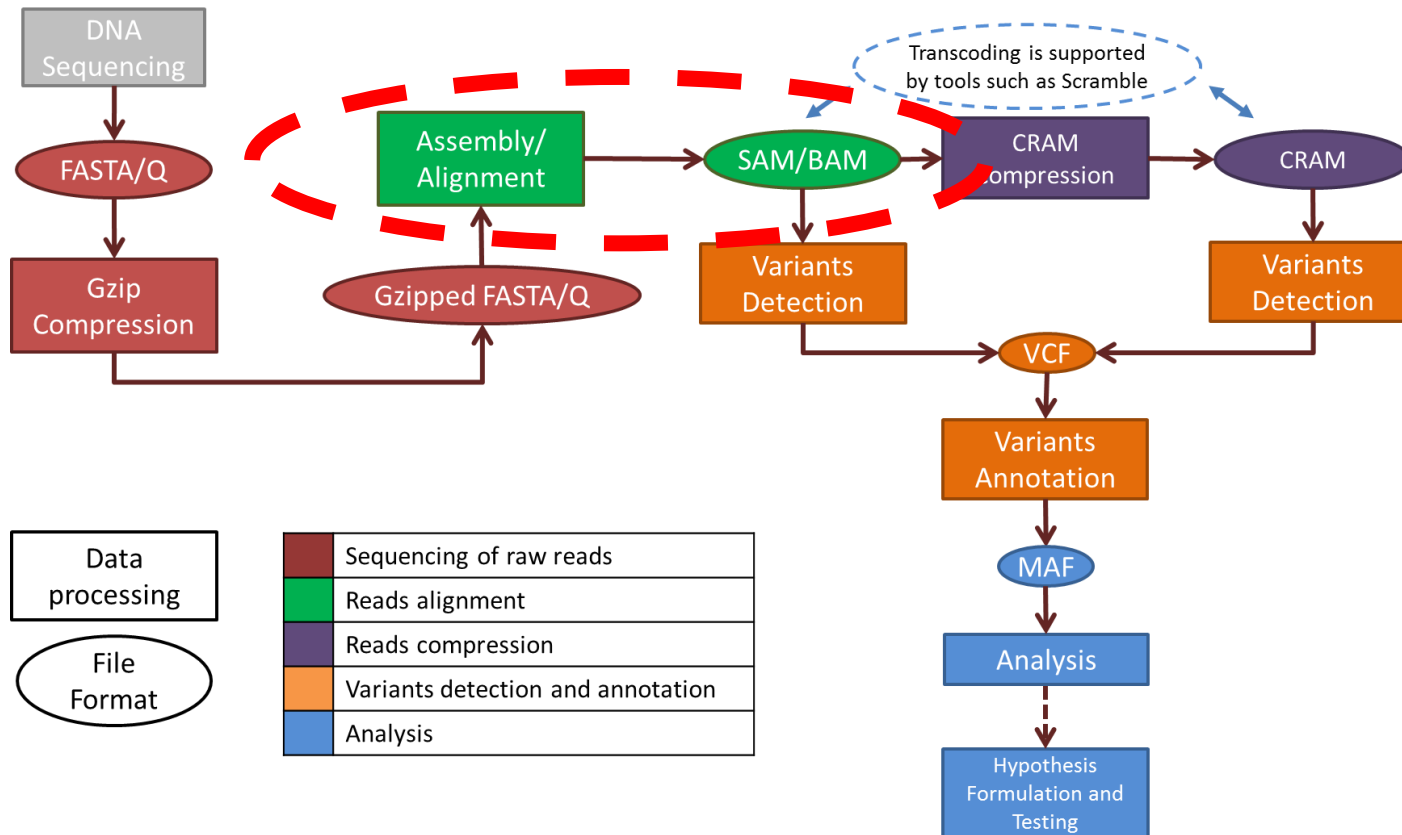
- Gzip of the entire txt file (sometimes split into several files)
- Compression ratio : 3 to 5
- According to the coverage 1 genome can take up to 2 TB

# From sequencing to assembly



# Alignment / mapping

- Raw data + alignment information



Coor 12345678901234 5678901234567890123456789012345  
ref 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 SM2I4M1D3M = 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:i:1

FastQ Headers

Quality scores if present

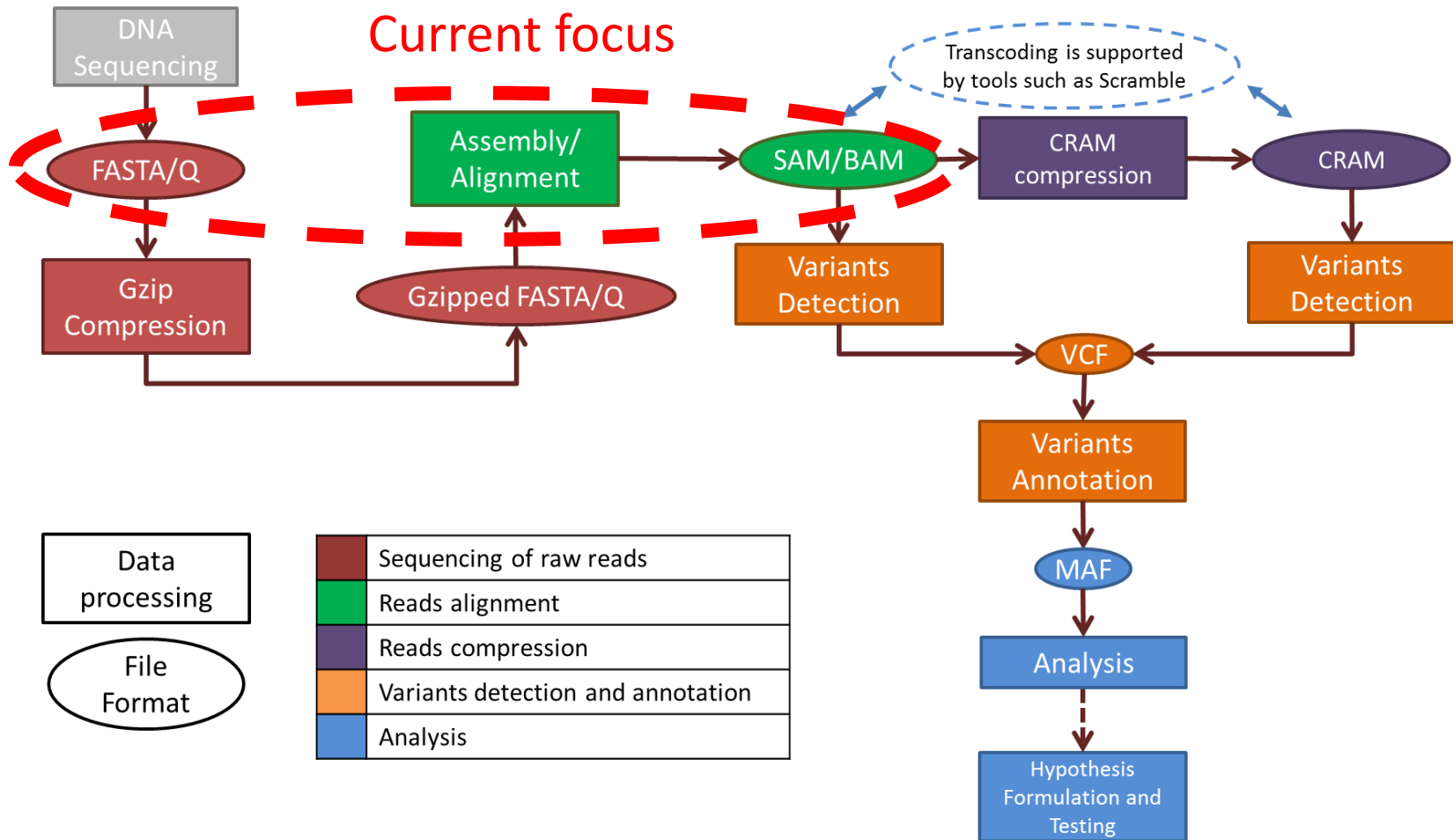
# Compressed SAM = BAM

- 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

