

# What are the inputs to a bioinformatics analysis?

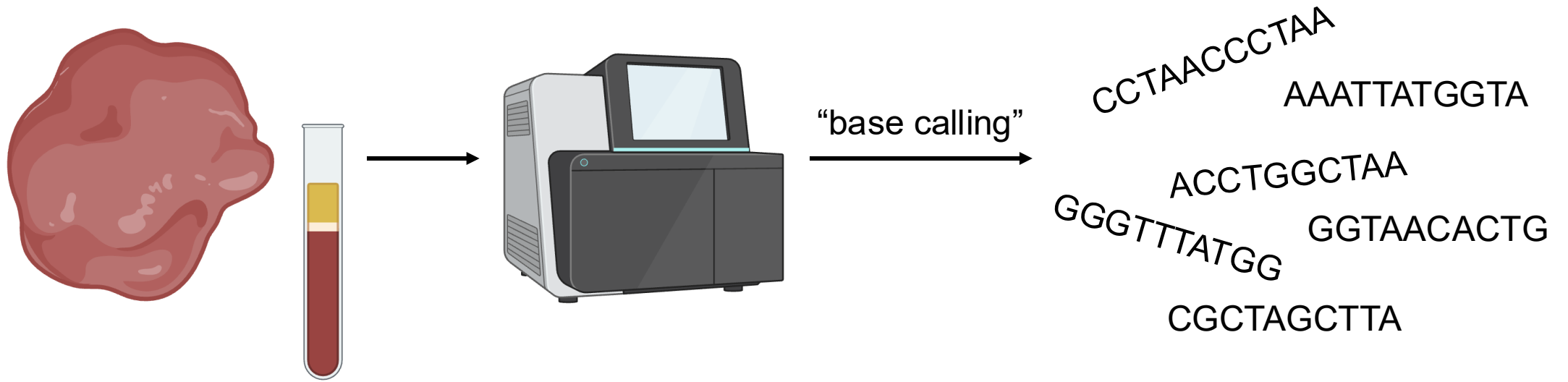
---

Clinical Cancer Genomics  
12 May 2025



# Learning objectives and lecture agenda

- Prepare for the analysis of DNA-sequencing data during lectures and exercises on Wednesday
- Learn about the components of a bioinformatics analysis other than the code and software – **what are the inputs?**
  - Become familiar with some of the main file types
- Learn the basics of the Integrative Genomics Viewer (IGV)



Overview of Illumina sequencing: <https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Sequencing instruments like the Illumina sequencer produce millions of reads that are stored in **FASTQ** files: 1 file for single-end, 2 files for paired-end sequencing.

### Example

#### Paired-end sequencing of a DNA fragment.



In paired-end sequencing, a DNA fragment is sequenced from its both ends resulting in two sequencing reads: *read1* and *read2*. These are stored in 2 FASTQ files.

**Note** that the DNA fragment to be sequenced is also referred to as an *insert*.

**Remember** that an Illumina sequencer sequences each of the strands separately and one after another. The illustration above is only meant to show from where each of the reads is originating.

This is *read1* obtained from the DNA fragment above (stored in the 1st FASTQ file).

```
@HWI-ST898:563:C4CK5ACXX:5:1212:14521:85648/1
```

```
AGGTTAATGCTTTCTCTCTA.....
```

```
+
```

```
=?>=@AADCC@@BBBBBB@.....
```



$$Q = -10 \cdot \log_{10}(P_{\text{err}})$$

---

This is *read2* obtained from the same DNA fragment (stored in the 2nd FASTQ file).

```
@HWI-ST898:563:C4CK5ACXX:5:1212:14521:85648/2
```

```
TTTCTCCTAAGATCTCAGTG.....
```

```
+
```

```
>><@BBBC@@C@@BBBBC?C.....
```

$$Q = -10 \cdot \log_{10}(P_{\text{err}})$$

Probability of a base being wrong = 0.001 ( $P_{\text{err}}$ )

$$-10 \cdot \log_{10}(0.001) = 30$$

Symbol	Phred Quality Score	Probability of Incorrect Base Call
<	27	0.002
=	28	0.001
>	29	0.001
?	30	0.001
@	31	0.0008
A	32	0.0006

# Quality control

A typical step after obtaining FASTQ files is to perform **quality control** using tools such as [FastQC](#).

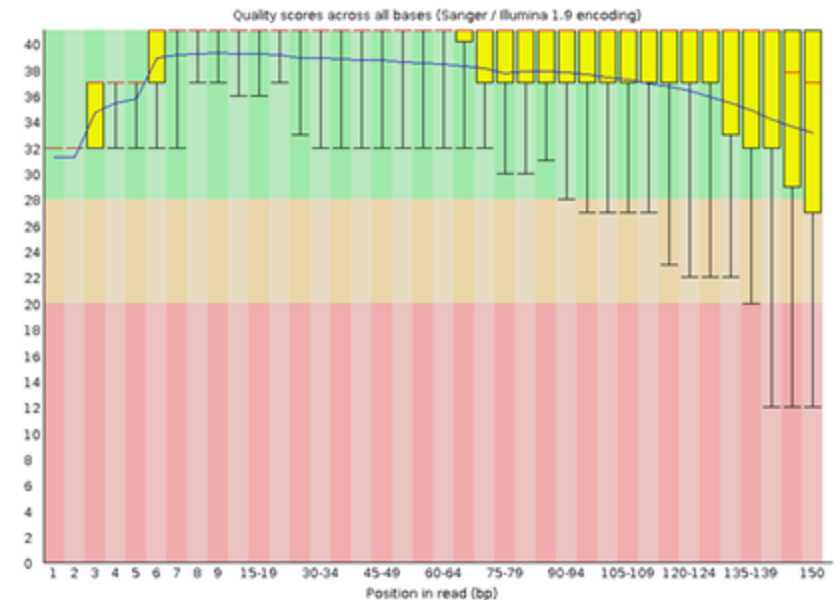
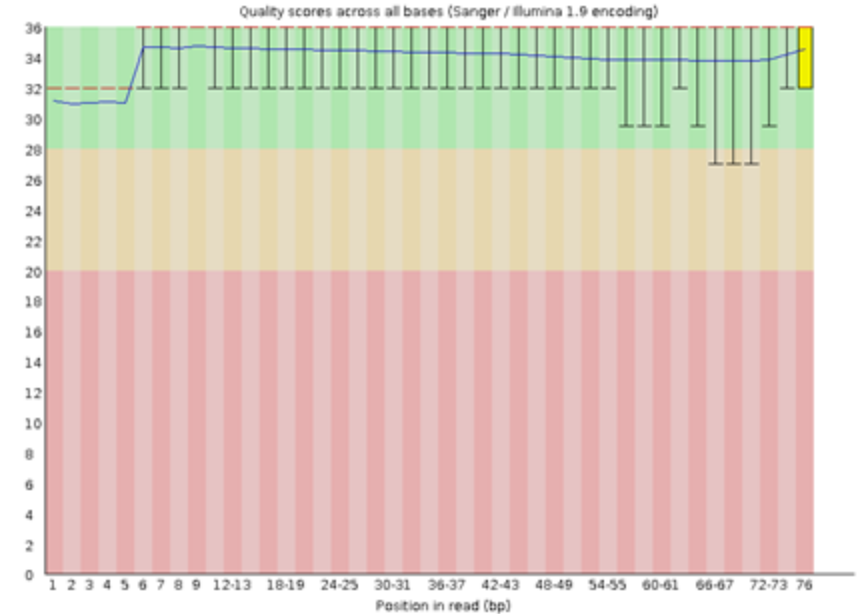
This analysis tells us if we need to take further actions before proceeding with downstream analyses.

The figure at the bottom right shows closer to the end of reads, quality score drops.

Trimming low-scoring end bases can enhance downstream analysis outcomes:

- e.g. using bioinformatics tools **Trimmomatic**, **Cutadapt**, **TrimGalore**.

✓ Per base sequence quality





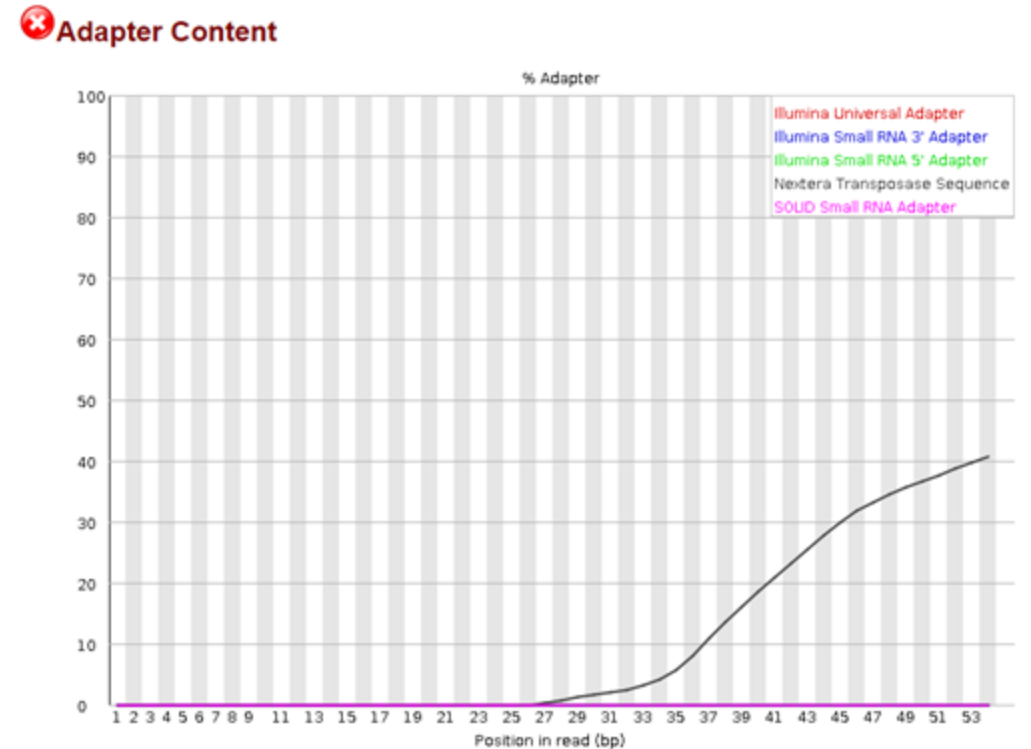
The top right figure shows **adapter contamination** closer to the end of reads.

This is caused by **adapter read-through**.

- In paired-end sequencing, this may occur when the DNA fragment length is shorter than the read length.

Trimming the adapter sequences enhances downstream analysis outcomes:

- e.g. using bioinformatics tools **Trimmomatic**, **Cutadapt**, **TrimGalore**.

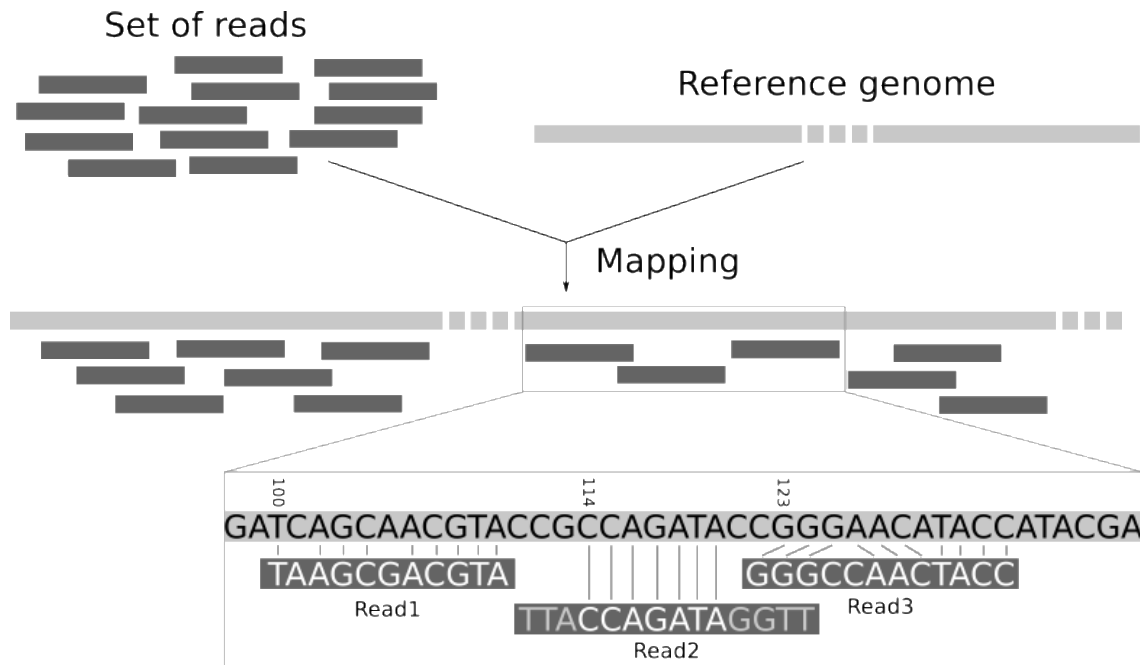


#### Example

##### Adapter read-through

```
ADAPTER-----read1----->
                                     ADAPTER
                                     ↓↓↓↓↓↓
5' NNNNNNAGGTTAATGCTTTCTCT.....CACTGAGATCTTAGGAGAAANNNNNN 3'
   |||
3' NNNNNNTCCAATTACGAAAGAGA.....GTGACTCTAGAATCCTCTTTNNNNNN 5'
<-----read2-----
↓↓↓↓↓
```

# Sequence alignment



To sequence DNA, we have to fragment it into smaller pieces and therefore lose the location information (i.e. where in the genome they belong).

**Sequence alignment or mapping** is the process of finding the genomic locations of the sequencing reads obtained from the sequencing of DNA fragments.

## Example



### Sequence alignment.

We have a reference genome, and a set of sequencing reads. We want to align these reads to the genomic region they have originated from.

```
CTAGAGCGTGGCCCGGAGCTGCCCTTTCCTCTTCGGTGAAGTTTTTAAAAGCTGCTGCGA # reference
```

```
CCGGAGCTGCCCTTTCCTCTTCGGTGA
```

```
CCTTACCTCTTCGGTGAAGTTTTTAAA
```

```
GAGCTGCCCTTTCCTCTTCGGTGAAGT
```

```
TGCCCTTACCTCTTCGGTGAAGTTTTT
```

```
CCCGGAGCTGCCCTTTCCTCTTCGGTG
```

```
TACCTCTTCGGTGAAGTTTTTAAAAGC
```

## After alignment.

↓

000000000111111112222222223333333334444444445555555556 # position-10s

12345678901234567890123456789012345678901234567890 # position-1s

CTAGAGCGTGGCCCGGAGCTGCCCTTTCCTCTTCGGTGAAGTTTTTAAAAGCTGCTGCGA # reference

CCGGAGCTGCCCTTTCCTCTTCGGTGA

CCTTACCTCTTCGGTGAAGTTTTTAAA

GAGCTGCCCTTTCCTCTTCGGTGAAGT

TGCCCTTACCTCTTCGGTGAAGTTTTT

CCCGGAGCTGCCCTTTCCTCTTCGGTG

TACCTCTTCGGTGAAGTTTTTAAAAGC

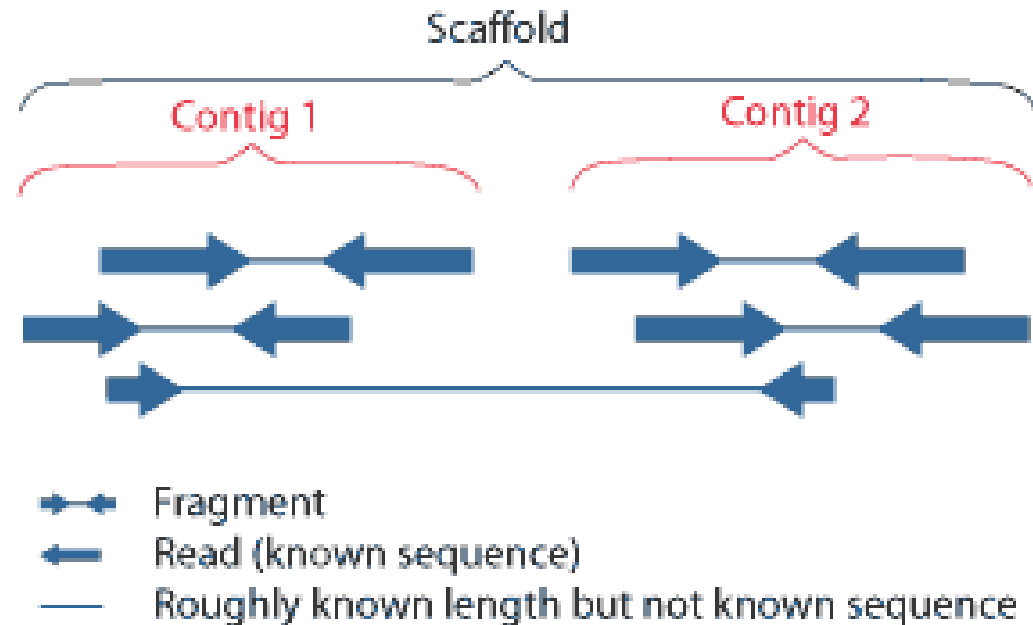
↑

🧩 In a real situation, the **reference is billions of nucleotides long** and we have **hundreds of millions of sequencing reads** each are e.g. **100 nucleotides long**.

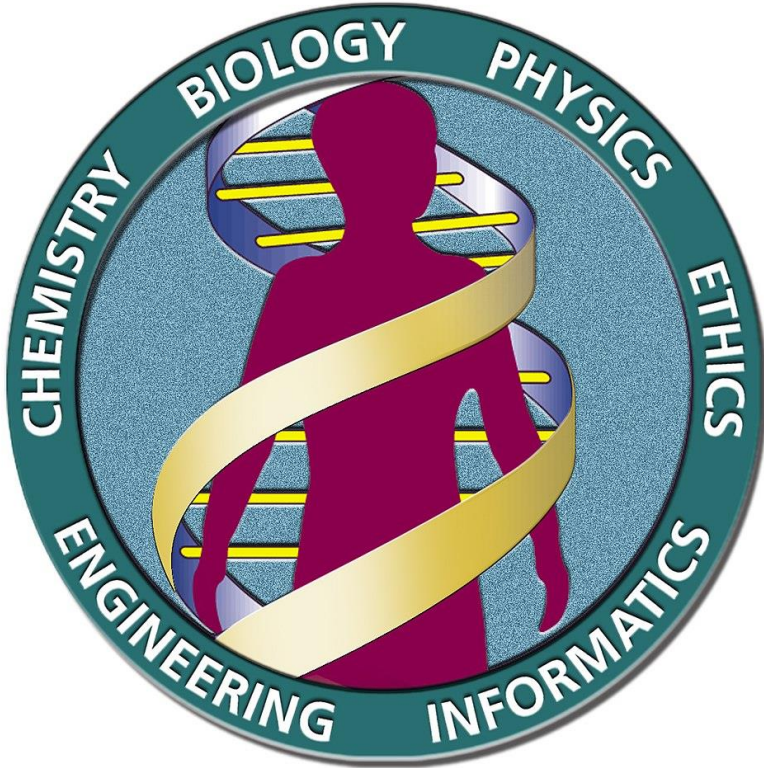
? Incidentally, what is happening at the position **27** marked by the arrow?

# We need a reference human genome

- An attempt at a complete representation of the nucleotide sequence of an individual genome.
- When we sequence new samples, we can map the reads to the reference rather than doing a new assembly each time.



# We need a reference human genome



Human Genome  
Project  
1990 - 2003

- The Human Genome Project genome sequence accounted for over 90% of the human genome in 2003
  - ... but it contained some **gaps** and **errors**
- Why are there gaps?
  - **Difficult-to-sequence regions**, e.g. repetitive regions
  - **Sequencing errors** and **low sequence read coverage** at certain loci make it difficult to reliably reconstruct the genome, thus resulting in gaps
- Efforts by the **Genome Reference Consortium (GRC)** and recently by the **Telomere to Telomere (T2T)** consortium have sought to fill in the remaining gaps.
- There are further efforts by the **Human Pangenome Reference Consortium (HPRC)** and others to "*better represent the human diversity*" in the human genome sequence.

## Major releases:

2006 → hg18 / NCBI36

2009 → hg19 / NCBI37 / GRCh37

2013 → hg38 / GRCh38

2022 → T2T-CHM13

## Minor releases or patches:

GRCh38.p13

Need to convert genome coordinates from one assembly to another (e.g. from hg19 to hg38)? → LiftOver



# The FASTA file

```
1 >gi|568336023|gb|CM000663.2| Homo sapiens chromosome 1, GRCh38 reference primary assembly
2 CCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTA
3 ACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTA
4 TAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACC
5 CTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACC
6 CAACCCCAACCCCAACCCCAACCCCAACCCCAACCCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAAC
7 CCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAAC
8 TAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAAC
9 TCTGACCTGAGGAGAACTGTGCTCCGCCTTCAGAGTACCACCGAAATCTGTGCAGAGGACAACGCAGCTC
10 CGCCCTCGCGGTGCTCTCCGGGTCTGTGCTGAGGAGAACGCAACTCCGCCGTTGCAAAGGCGCGCCGCGC
11 CGGCGCAGGCGCAGAGAGGCGCGCCGCGCCGGCGCAGGCGCAGAGAGGCGCGCCGCGCCGGCGCAGGCGC
12 AGAGAGGCGCGCCGCGCCGGCGCAGGCGCAGAGAGGCGCGCCGCGCCGGCGCAGGCGCAGAGAGGCGCGC
13 CGCGCCGGCGCAGGCGCAGACACATGCTAGCGCGTCGGGGTGGAGGCGTGGCGCAGGCGCAGAGAGGCGC
14 GCCGCGCCGGCGCAGGCGCAGAGACACATGCTACCGCGTCCAGGGGTGGAGGCGTGGCGCAGGCGCAGAG
15 AGGCGCACCGCGCCGGCGCAGGCGCAGAGACACATGCTAGCGCGTCCAGGGGTGGAGGCGTGGCGCAGGCG
16 GCAGAGACGCAAGCCTACGGGCGGGGGTGGGGGGGCGTGTGTTGCAGGAGCAAAGTCGCACGGCGCCGG
17 GCTGGGGCGGGGGGAGGGTGGCGCCGTGCACGCGCAGAACTCACGTACGGTGGCGCGGCGCAGAGACG
18 GGTAGAACCTCAGTAATCCGAAAAGCCGGGATCGACCGCCCCTTGCTTGCAGCCGGGCACTACAGGACCC
19 GCTTGCTCACGGTGCTGTGCCAGGGCGCCCCCTGCTGGCGACTAGGGCAACTGCAGGGCTCTCTTGCTTA
20 GAGTGGTGGCCAGCGCCCCCTGCTGGCGCCGGGGCACTGCAGGGCCCTCTTGCTTACTGTATAGTGGTGG
```



# Aligned sequences

The alignment results to a reference genome are stored in a standard file format called **Sequence Alignment/Map (SAM)**.

- This file is text-based, human readable, and takes up a lot of space.
- The **binary alignment map (BAM)** is a binary file that stores similar information in a compressed form and thus uses less space, but it is not human-readable.



# Common genomic data file formats

- **FASTQ** – raw sequence data (with qualities)
- **FASTA** – sequences (DNA, RNA, protein)
- **SAM/BAM** – aligned sequence data
- **BED** – other genome features
- **GTF** – gene/transcript annotations
- **VCF** – variant calls (individual, multi-individual)
- **MAF** – aggregated variant information (project, population)
- Plus, many other (potentially custom) data formats output by specialized tools

Post-alignment, sequencing reads can be summarized over/within genomic intervals.

**Browser Extensible Data (BED) file format** - Genomic intervals and information.

1	chr7	127471196	127472363
2	chr7	127472363	127473530
3	chr7	127473530	127474697
4	chr7	127474697	127475864
5	chr7	127475864	127477031
6	chr7	127477031	127478198
7	chr7	127478198	127479365
8	chr7	127479365	127480532
9	chr7	127480532	127481699



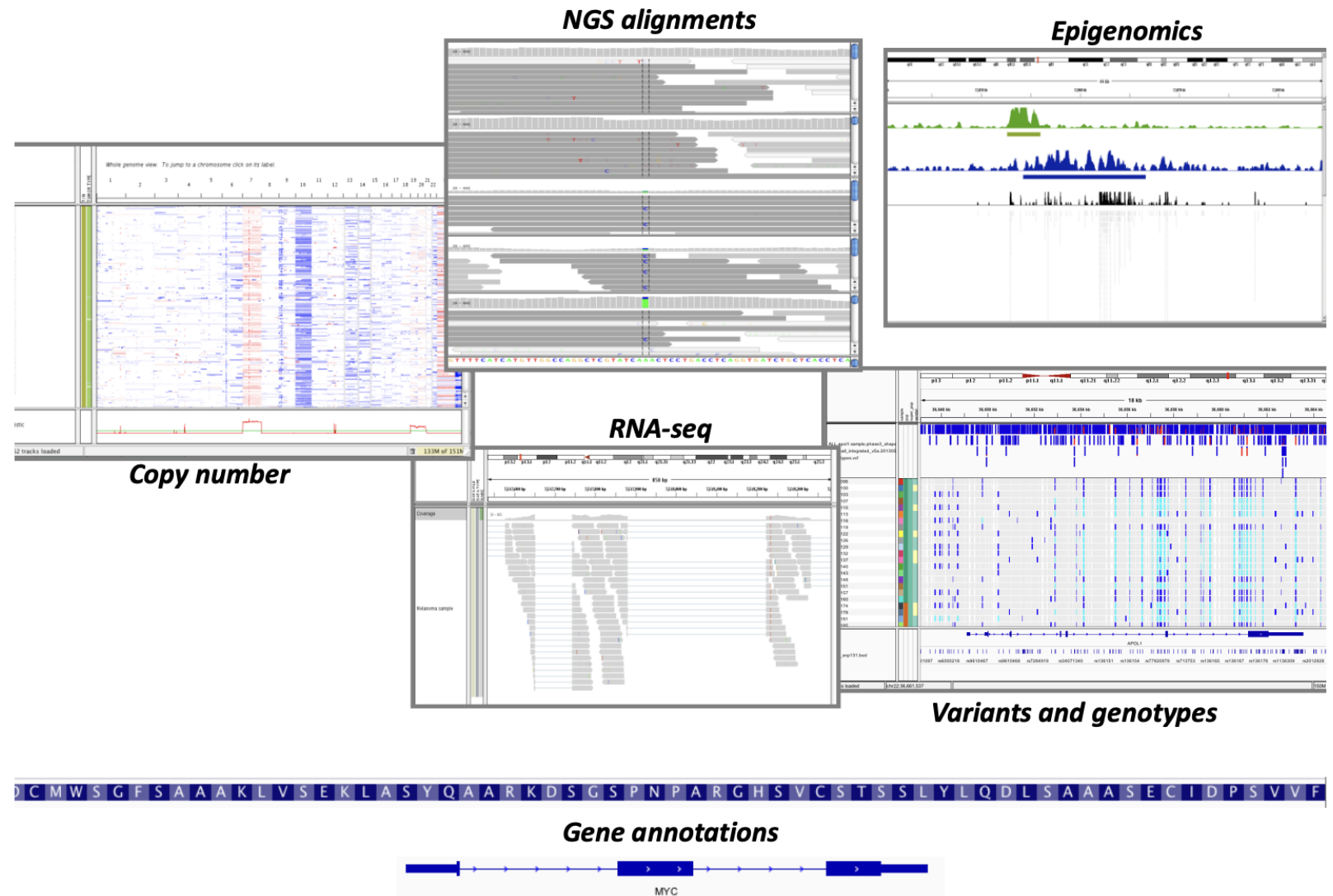
We may also need some annotations of our genome  
... for example, where are certain genomic features located?

## Gene Transfer Format (GTF) – Information on gene structure

<u>Col 1</u>	<u>Col 2</u>	<u>Col 3</u>	<u>Col 4</u>	<u>Col 5</u>	<u>Col 6</u>	<u>Col 7</u>	<u>Col 8</u>	<u>Col 9</u>
chr21	HAVANA	transcript	10862622	10863067	.	+	.	gene_id "ENSG00000169..
chr21	HAVANA	exon	10862622	10862667	.	+	.	gene_id "ENSG00000169..
chr21	HAVANA	CDS	10862622	10862667	.	+	0	gene_id "ENSG00000169..
chr21	HAVANA	start_codon	10862622	10862624	.	+	0	gene_id "ENSG00000169..
chr21	HAVANA	exon	10862751	10863067	.	+	.	gene_id "ENSG00000169..
chr21	HAVANA	CDS	10862751	10863064	.	+	2	gene_id "ENSG00000169..
chr21	HAVANA	stop_codon	10863065	10863067	.	+	0	gene_id "ENSG00000169..
chr21	HAVANA	UTR	10863065	10863067	.	+	.	gene_id "ENSG00000169..

# Integrative Genomics Viewer (IGV)

A visualization tool to simultaneously integrate and analyze multiple types of genomic data.



## IGV supports many different file formats

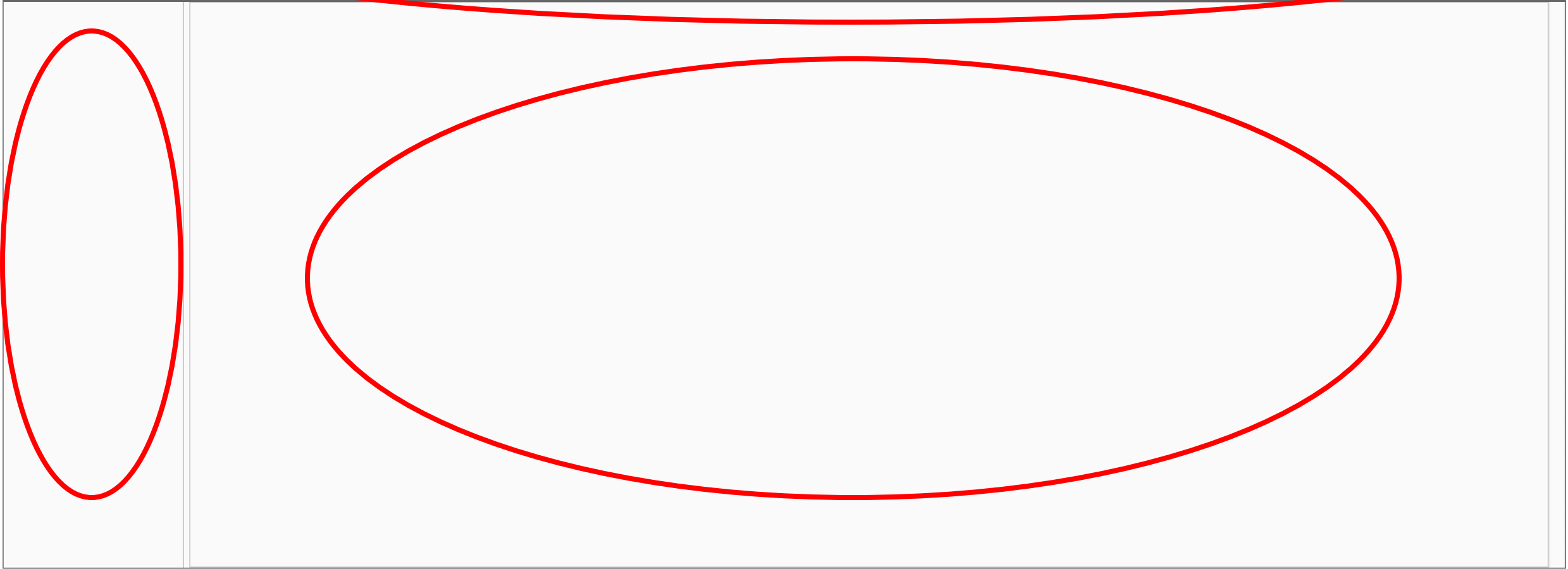
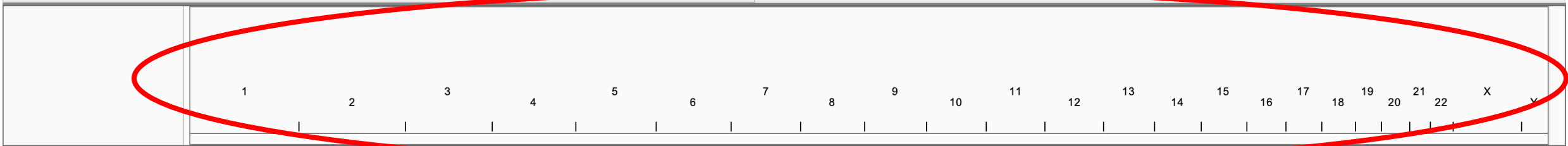
- [BAM](#)
- [BED](#)
- [BedGraph](#)
- [bigBed](#)
- [bigWig](#)
- [Birdsuite Files](#)
- [broadPeak](#)
- [CBS](#)
- [CN](#)
- [Custom File Formats](#)
- [Cytoband](#)
- [FASTA](#)
- [GCT](#)
- [genePred](#)
- [GFF/GTF](#)
- [GISTIC](#)
- [Goby](#)
- [GWAS](#)
- [IGV](#)
- [LOH](#)
- [MAF \(Multiple Alignment Format\)](#)
- [MAF \(Mutation Annotation Format\)](#)
- [Merged BAM File](#)
- [MUT](#)
- [narrowPeak](#)
- [PSL](#)
- [SAM](#)
- [Sample Info \(Attributes\) file](#)
- [SEG](#)
- [SNP](#)
- [TAB](#)
- [TDF](#)
- [Track Line](#)
- [Type Line](#)
- [VCF](#)
- [WIG](#)
- [chrom.sizes](#)

For more info see: [www.broadinstitute.org/igv/FileFormats](http://www.broadinstitute.org/igv/FileFormats)

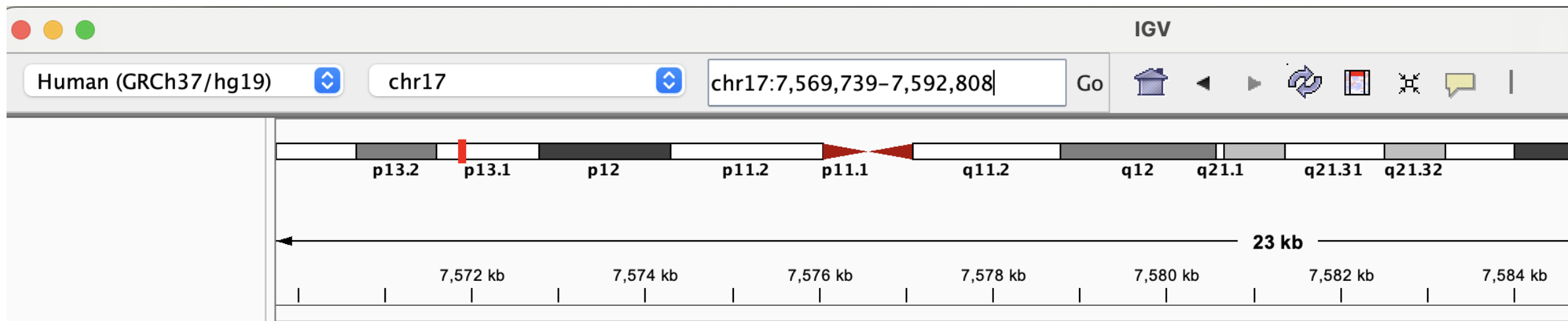
IGV

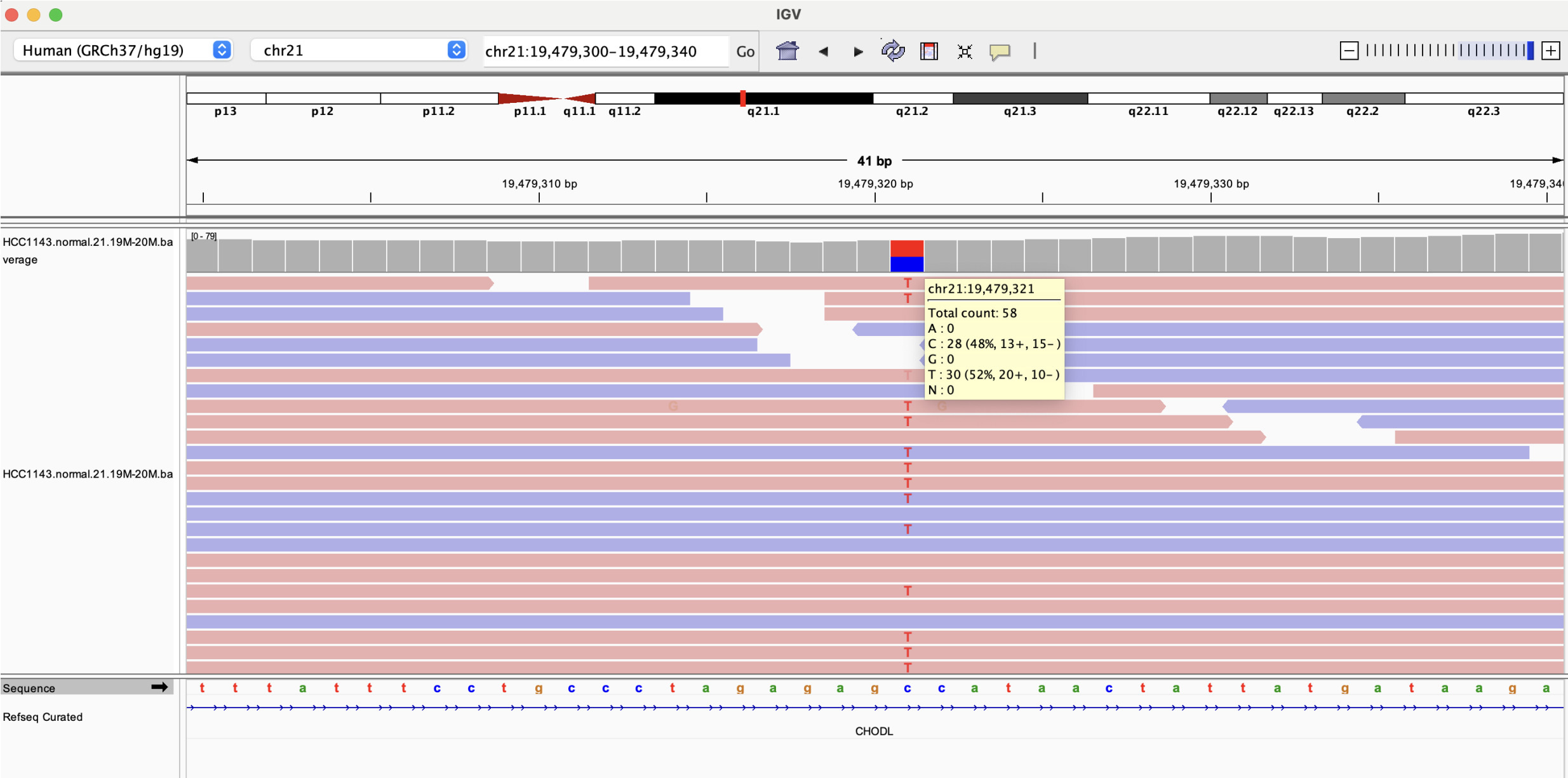
Human (GRCh37/hg19) All Go

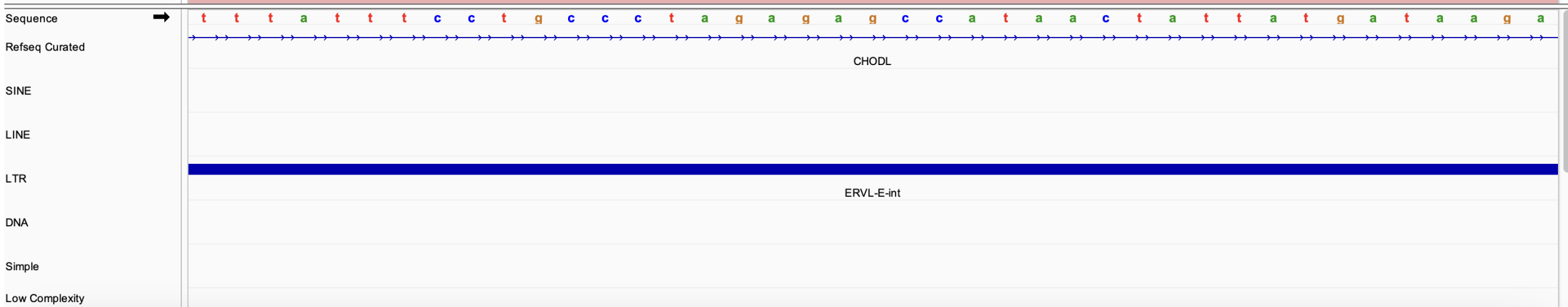
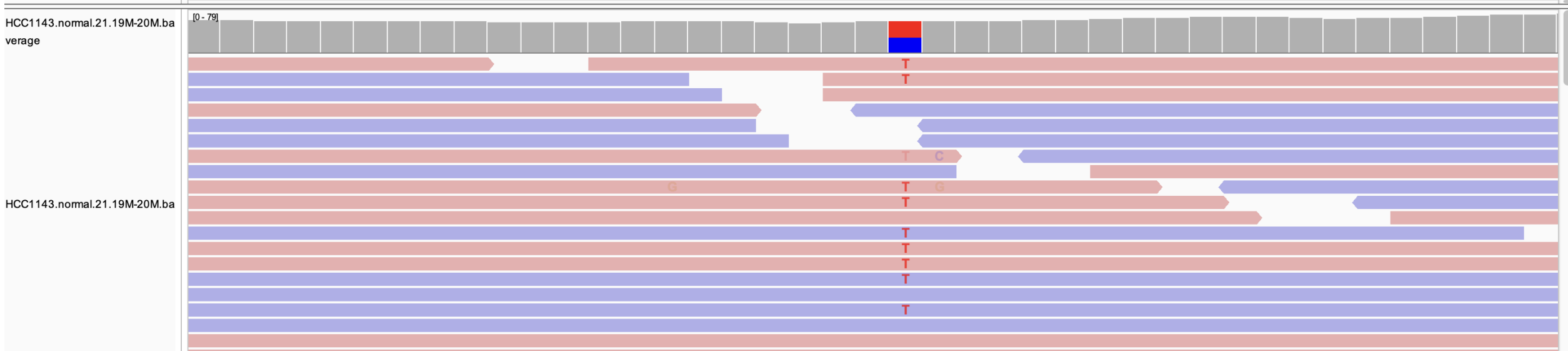
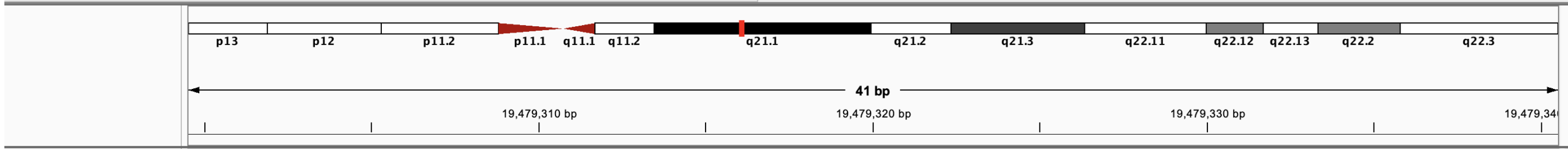
[-] [Zoom Bar] [+]











**... now all we need is some tools and software!**

- In **Exercise Set 1**, we will be inspecting some of the files and file types discussed in this lecture
- We will also work on becoming familiar with the **Integrative Genomics Viewer (IGV)** software, which we will use many times during later Exercise Sets
- Goal: prepare for analysis of DNA-sequencing data in Exercise Set 2 (starting Wednesday)

## **Exercise Set 1 overview**

1. Overview of the computing environment
2. Parts of a bioinformatics analysis: reference genome, annotation files, and raw data files
3. Introduction to IGV

Optional:

- Programming in R for Bioinformatics



**Karolinska  
Institutet**