

Bioinformatics Analysis in R

Day 5 Next Generation Sequencing (NGS) Data Analysis and Visualization

Ivan G. Costa, Zhijian Li

Institute for Computational Genomics
RWTH University Hospital
www.costalab.org

Outline

- Introduction to NGS data analysis pipeline
 - Quality check
 - Alignment
 - Higher level analysis (peak calling)
 - File formats
- Visualization of NGS data using IGV
 - RNA-seq, ChIP-seq, ...
 - IGV tools
- Practice

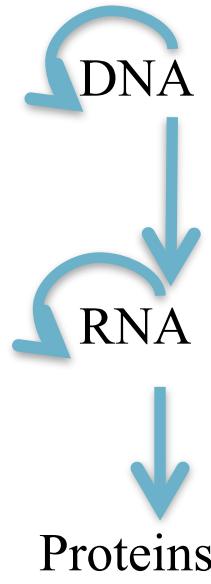
Bioinformatics Analysis in R

Next Generation Sequencing

Sequencing

- Read the bases of a DNA/RNA sequence
- Applications
 - Sequence DNA of known or unknown organism
 - Detect variants on patients
 - Sequence the RNA of a cell
 - Detect location of proteins interacting with DNA
- Problem
 - Only short DNA sequences (< 1000 bps) can be read
- Solution
 - Bioinformatics

Information Level vs. NGS



DNA Sequencing

- > detection of genetic variants
- > de-novo reconstructions of genomes

RNA Sequencing

- > quantification of RNA in a cell
- > de-novo identification of RNAs

Detection of Interactions:

- ChIP Sequencing -> a protein with DNA
- CLIP Sequencing -> a protein with RNA
- ChIRP Sequencing -> a RNA with DNA
- ...

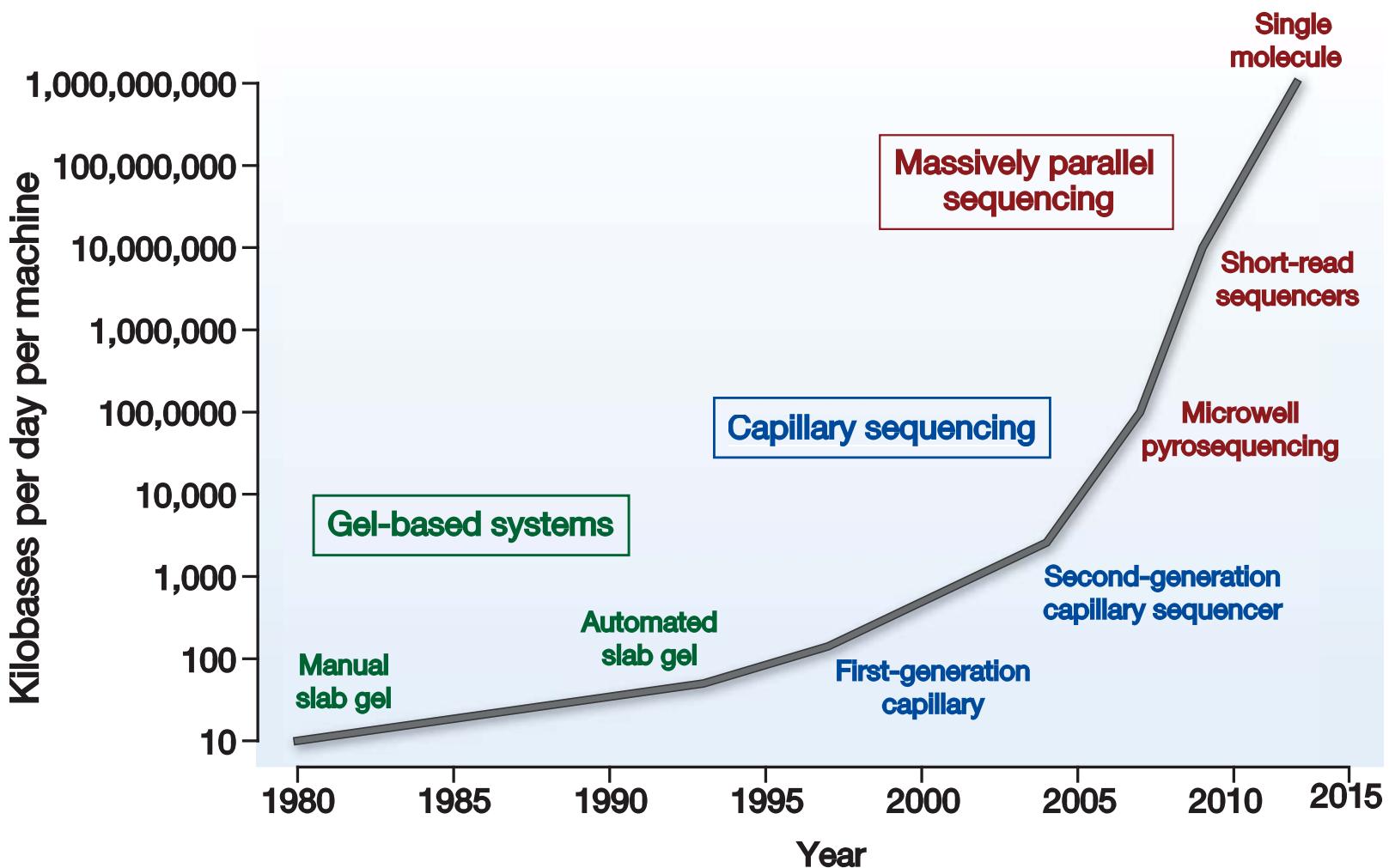
See here for a comprehensive list of Seq essays (>50)
<https://liorpachter.wordpress.com/seq/>

Next Generation Sequencing

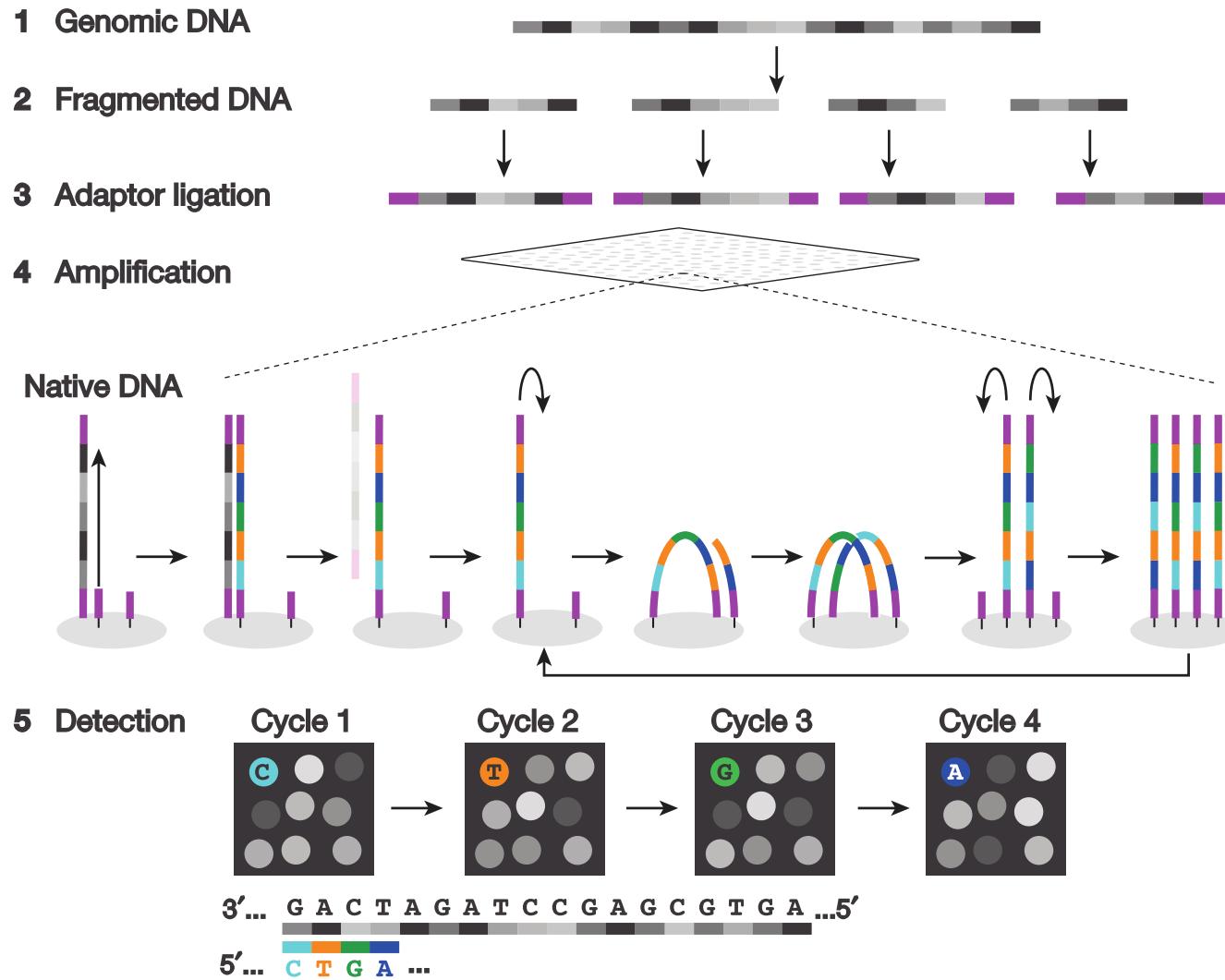
- ▶ NGS take advantage of **parallelization**
 - ▶ reads millions/billions of reads for a time
 - ▶ shorter reads (50-100 bps)
 - ▶ higher error rates (0.1-1%)
- ▶ commercial products:
 - ▶ 454
 - ▶ SOLiD
 - ▶ **Solexa (Illumina)**



Next Generation Sequencing

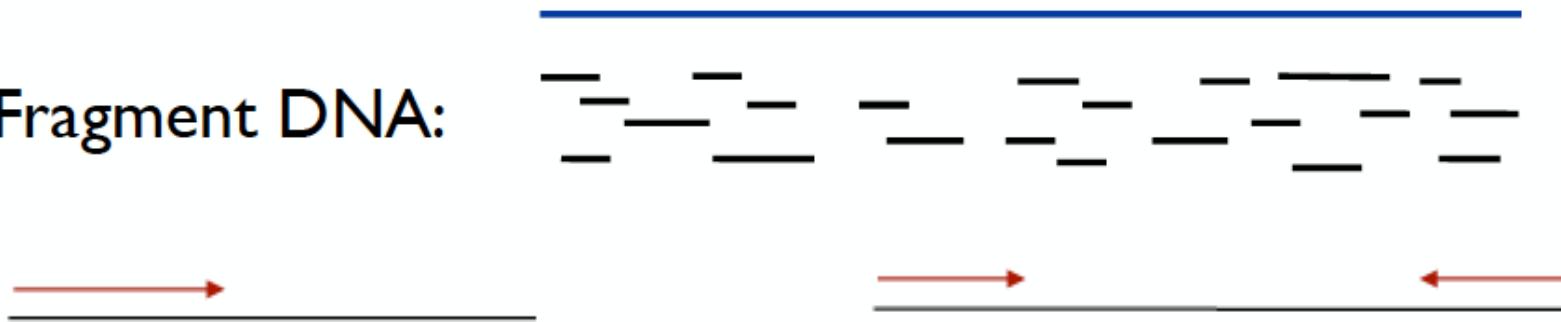


Next Generation Sequencing



Read Types

Fragment DNA:

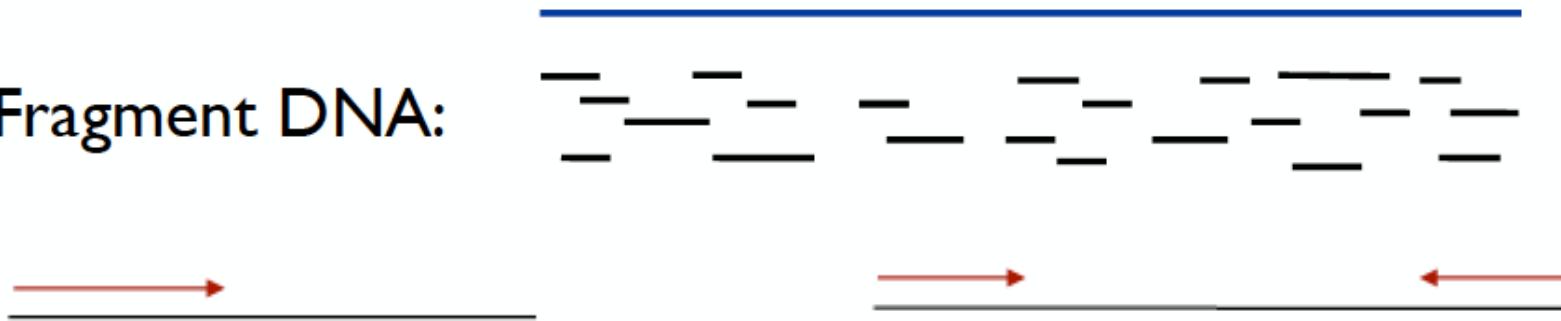


Single end

Paired end
Ins: 200-800 bp

Read Types

Fragment DNA:



Single end

Advantages:

- cheaper
- compatible with protocols producing small fragments (Ribo-seq, miRNA-seq)

Paired end Ins: 200-800 bp

Advantages:

- easier to align
- helps detection of variants (DNA), exon pairs (RNA)

FASTA files

```
>dnaA chromosomal replication initiator protein DnaA
MSLSLWQQCLARLQDEL PATEFSMWIRPLQAELSDNTLALYAPNRFVLDW
VRDKYLEALRDLLALQEKLVTIDNIQKTVAEYYKIKVADLLSKRRSRSVARP
RQMAMALAKELLHAVGNGIMARKPNAKV VYMHSERFVQDMVKALQNNAI
EEFKRYYRSVDALLIDDFSLPEIGDAFGGRDHTTVLHACRKIEQLREESHD
KEDFSNLIRTLSS
```

FASTA files

Start symbol

Sequence ID
(no spaces)

Sequence description
(spaces allowed)

> dnaA chromosomal replication initiator protein DnaA

MSLSLWQQCLARLQDEL PATEFS M WIRPL QAE LSDN T LALYAP NRFV LDW
VRDKY LEAL RDLL ALQ EKL VTI DNI QKTVA EYY KIKVAD LL SKRR SR SVA RP
RQMAMALAKELLHAVGNGIMARKPNAKV VYMHSERFVQDMVKALQNNAI
EEFKRYYRSVDALLIDDFSLPEIGDAFGGRDHTTVLHACRKIEQLREESHD
KEDFSNLIRTLSS



The sequence

FASTQ files

Header
Sequence
Qualities
(prob. that base call is wrong)

```
@ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1
ATTCCCAGGCCTTTCCAGGCCTGCCTGCTCGAGC
+
BAAAGECEE<EEDFEDF3DBDBB=A+==>9>>88?
```

One character encodes a number
using ascii table (0-255)

Phred-scale

$$Q = -10 * \log_{10} P$$

This number (Q) can be
converted to P

$$P = 10^{(-Q/10)}$$

FASTQ files

Uses letters/symbols to represent numbers:

!"#\$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJ

↓
Q0

↓
Q10

↓
Q20

↓
Q30

↓
Q40

bad

maybe

ok

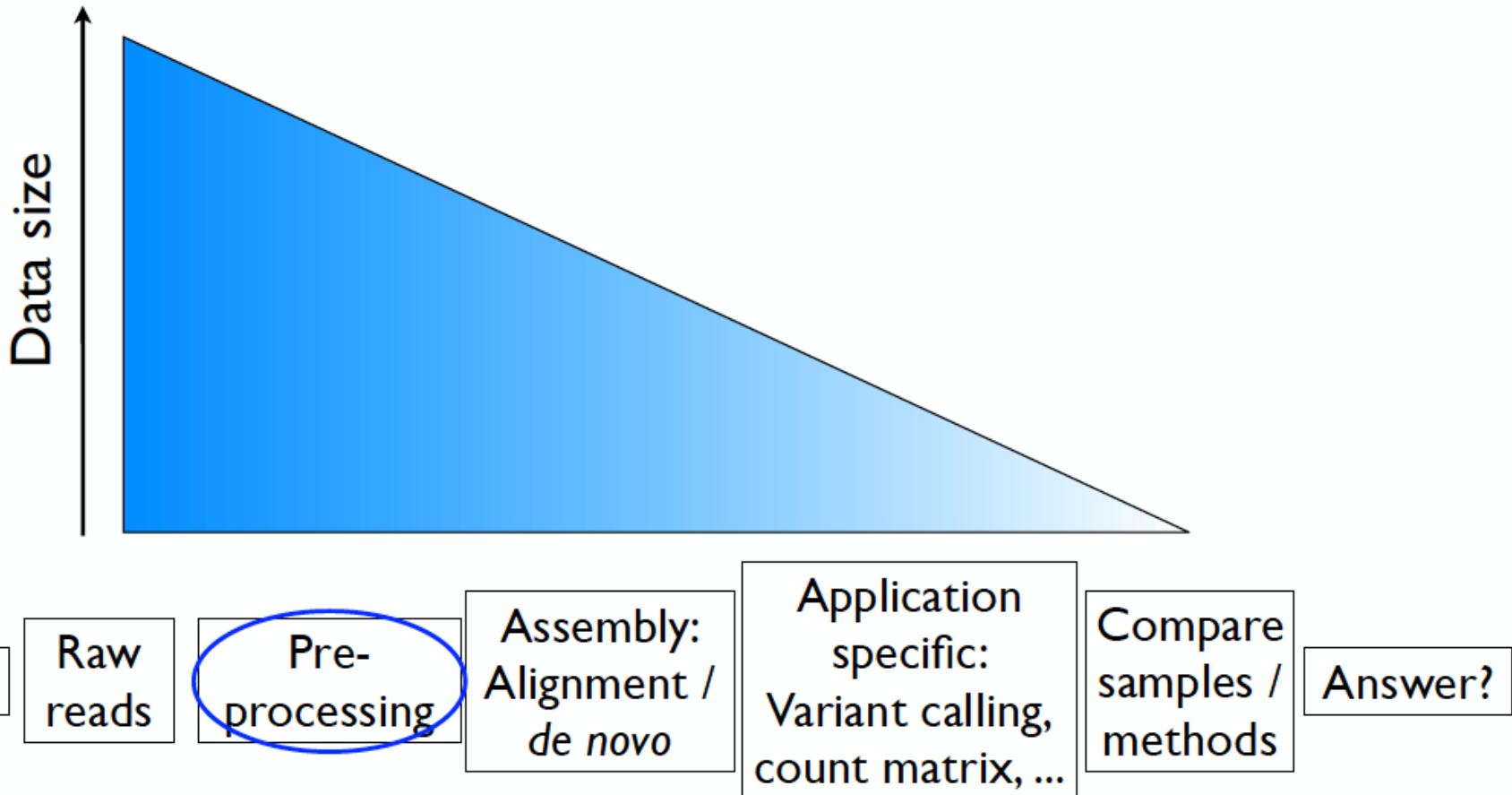
good

excellent

Bioinformatics Analysis in R

Next Generation Sequencing Data Analysis

Pre-processing



Pre-processing

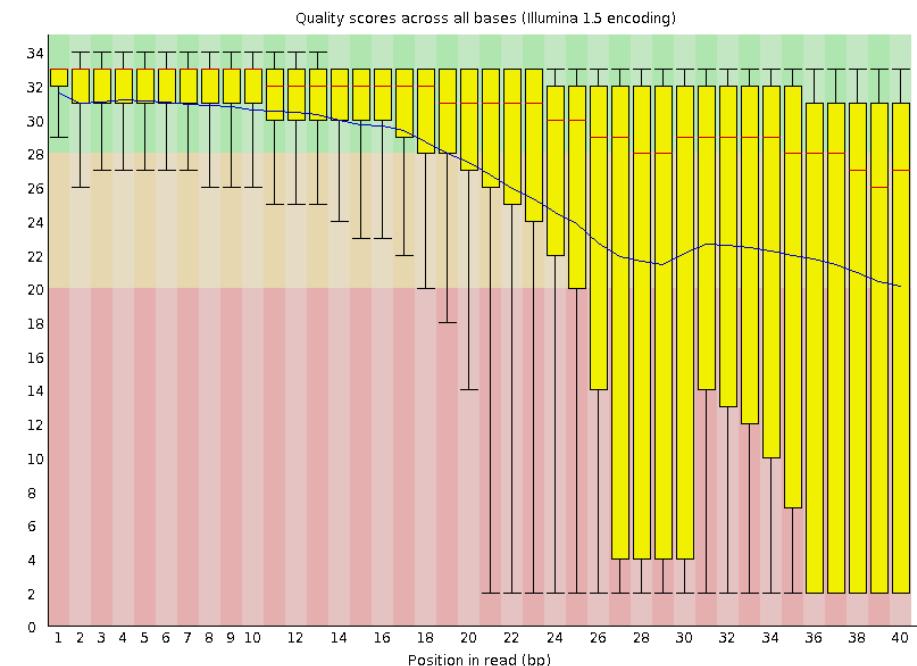
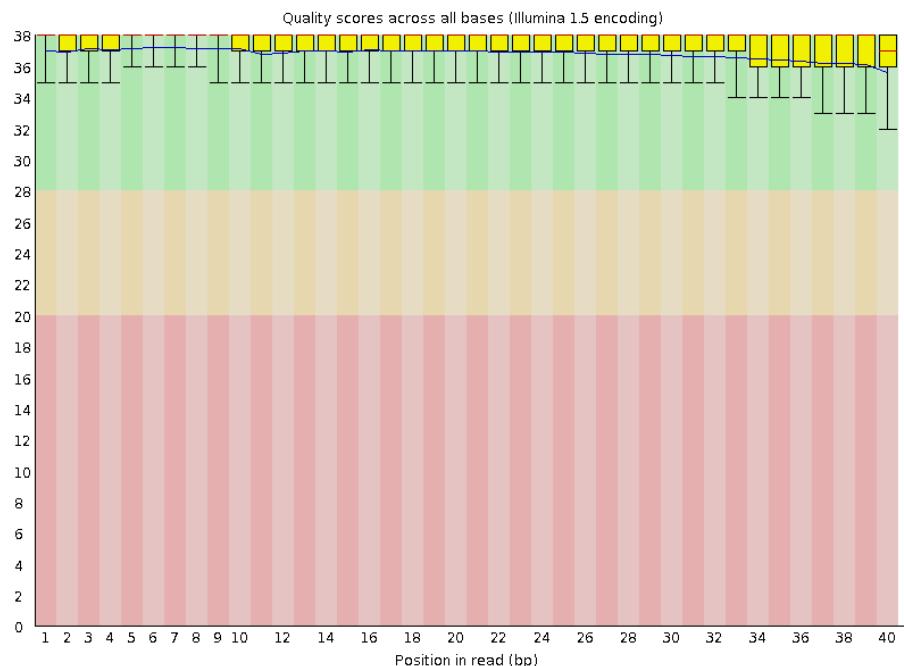
- Sequencing and sample preparation introduce errors
 - Errors in start/end of reads
 - Bases bias on read positions
 - Presence of adapter sequences
 - Fragment duplication from PCR
 - ...
- Tools: FastQC (for checking), Trimmomatic (for trimming), ...

Quality Control

- **FastQC (usually provided by NGS core facilities)**
 - tool to analyse quality of reads from sequencing.
 - indicate problems in library preparation or sequencing steps.
- Example of good sequences:
 - http://www.bioinformatics.babraham.ac.uk/projects/fastqc/good_sequence_short_fastqc.html
- or bad sequences:
 - http://www.bioinformatics.babraham.ac.uk/projects/fastqc/bad_sequence_fastqc.html

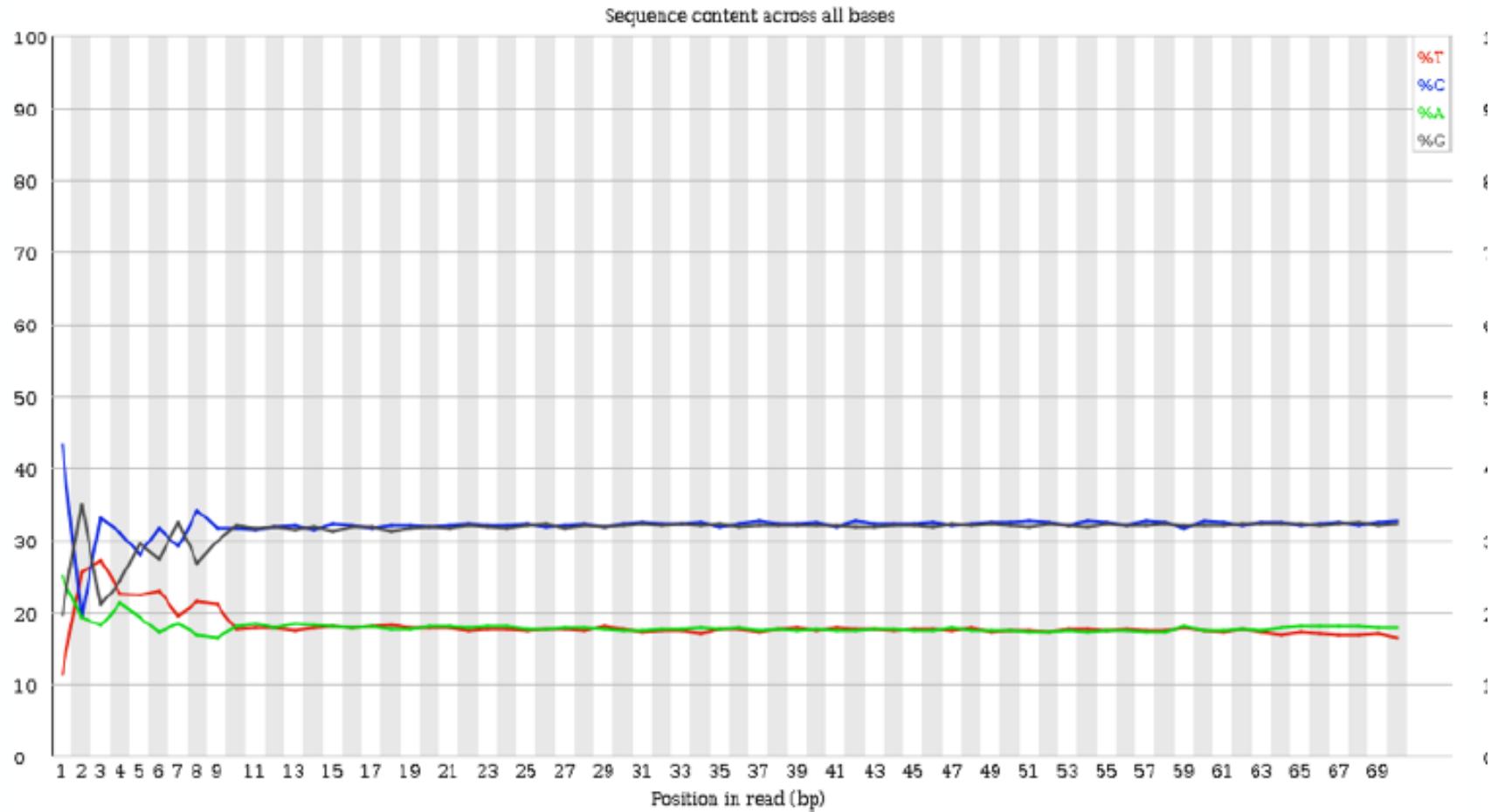
Quality Control

Sequencing quality decreases with size.



Solution: trim end of reads with low quality

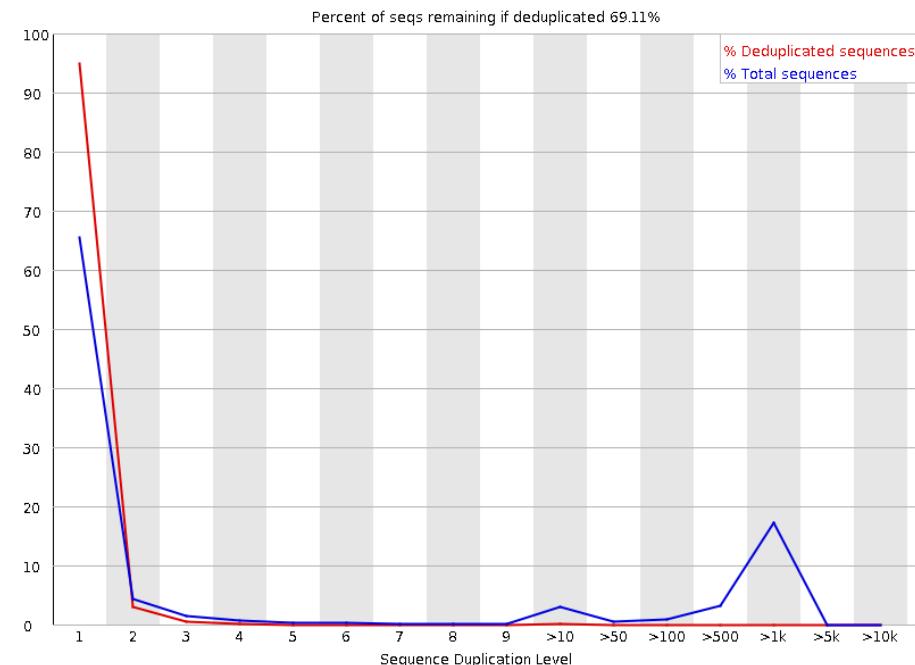
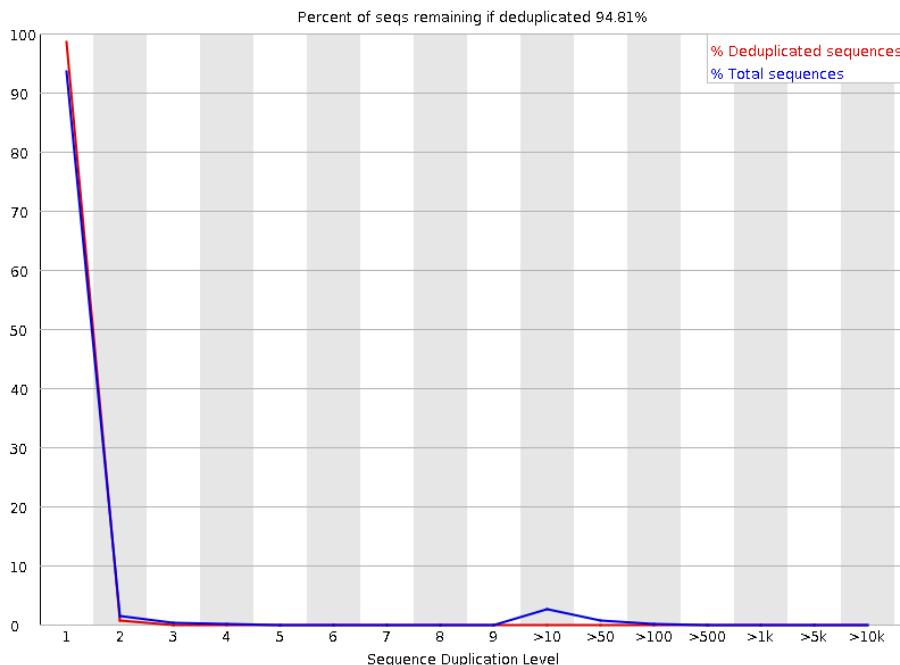
Read position sequence bias



- Trim read starts

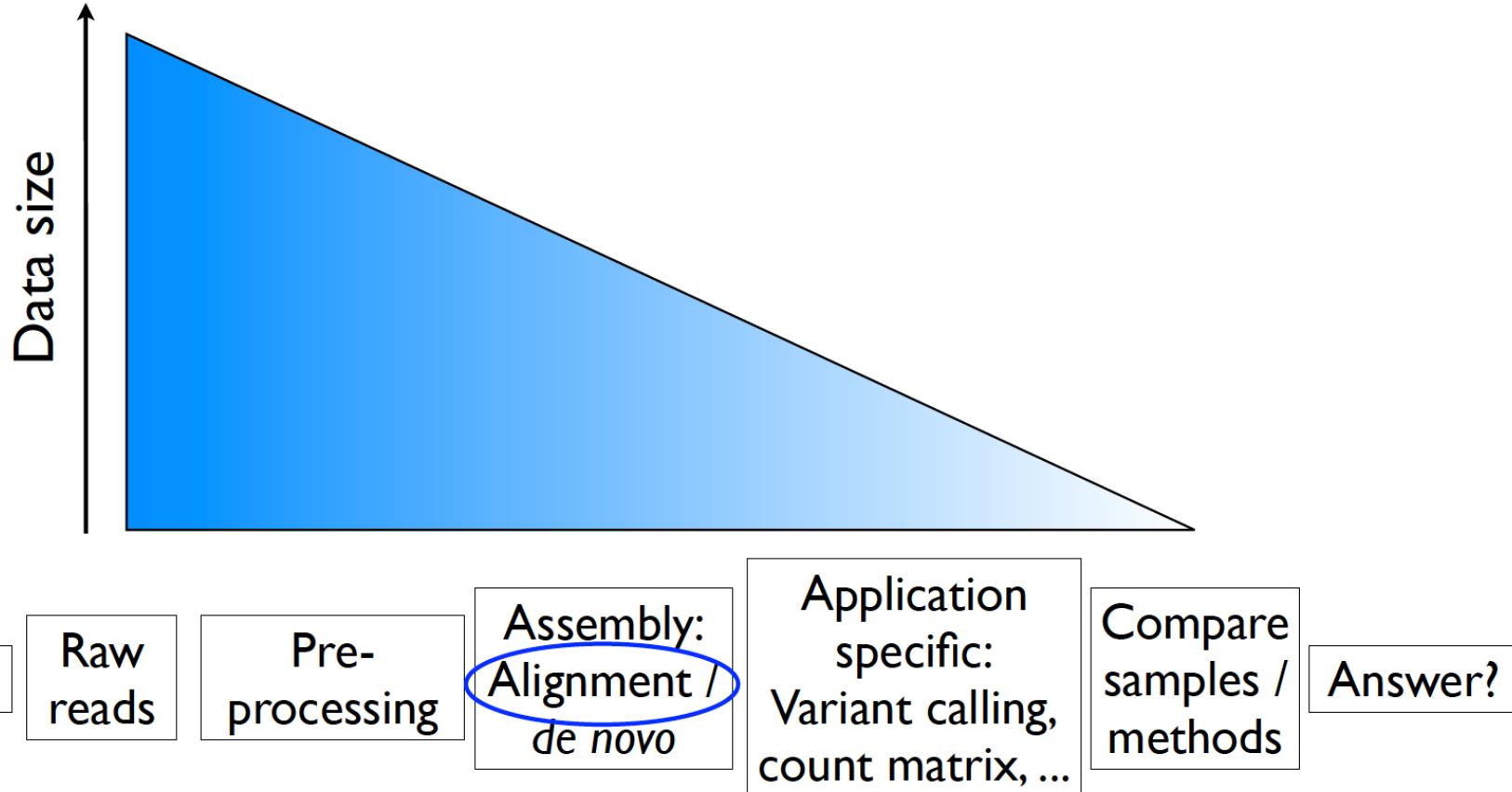
Quality Control

Sequence duplication levels



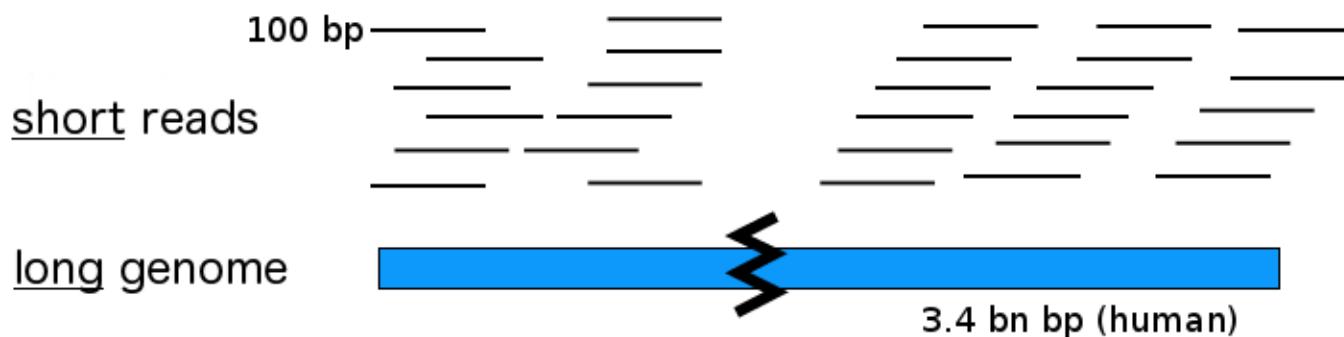
- Solution: remove duplicates

(Short reads) Alignment



Short Read Alignment

- Query
 - sequenced reads in FASTQ format
 - huge number of them, 1M ~ 100M
 - short read length, ~100 bp
- Reference
 - human genome in FASTQ format
 - total size ~3 billion bps
- Lots of short vs. a few longs
 - BLAST would take several years to run.



Pitfalls

- **(Unknown) divergence of sample and reference genome**
- **Poor genome reference quality**
- **Repeats in the genome (larger than read size)**
- **Recombinations**
- **Sequencing/read errors**

Algorithms - Alignment

Short read alignment is a special problem

- reference sequence (genome) is large and fixed
- query sequence (reads) are short and many

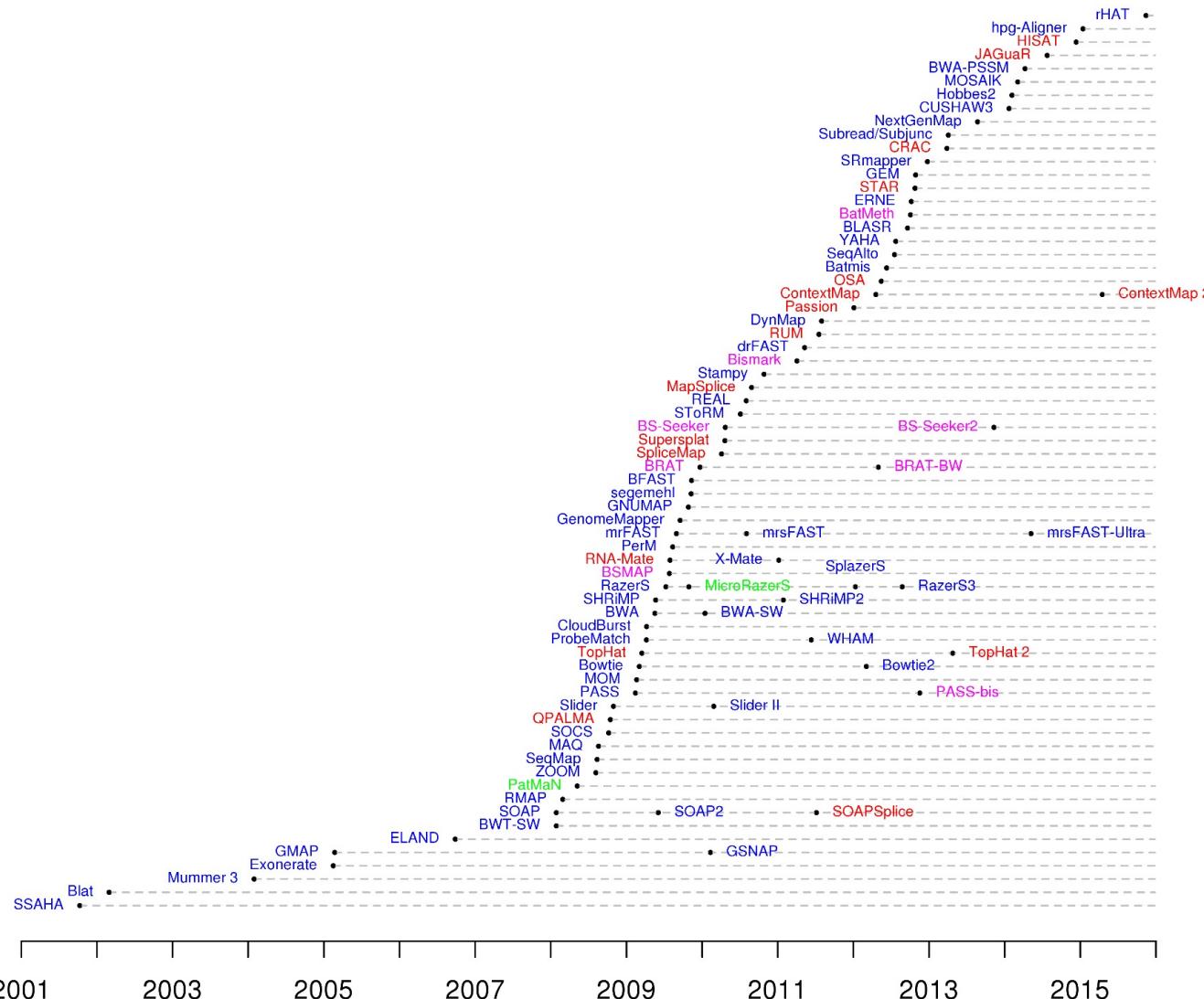
Solution:

1. Pre-process the genome finding all exact alignments for small sequences (>14bps) (index)

- k-mer hash table (>10GB)
- compressed suffix trees (> 4GB)

2. Break your read in small pieces (>14bps) and extend your alignment on all candidate positions using dynamic programming.

Alignment Tools



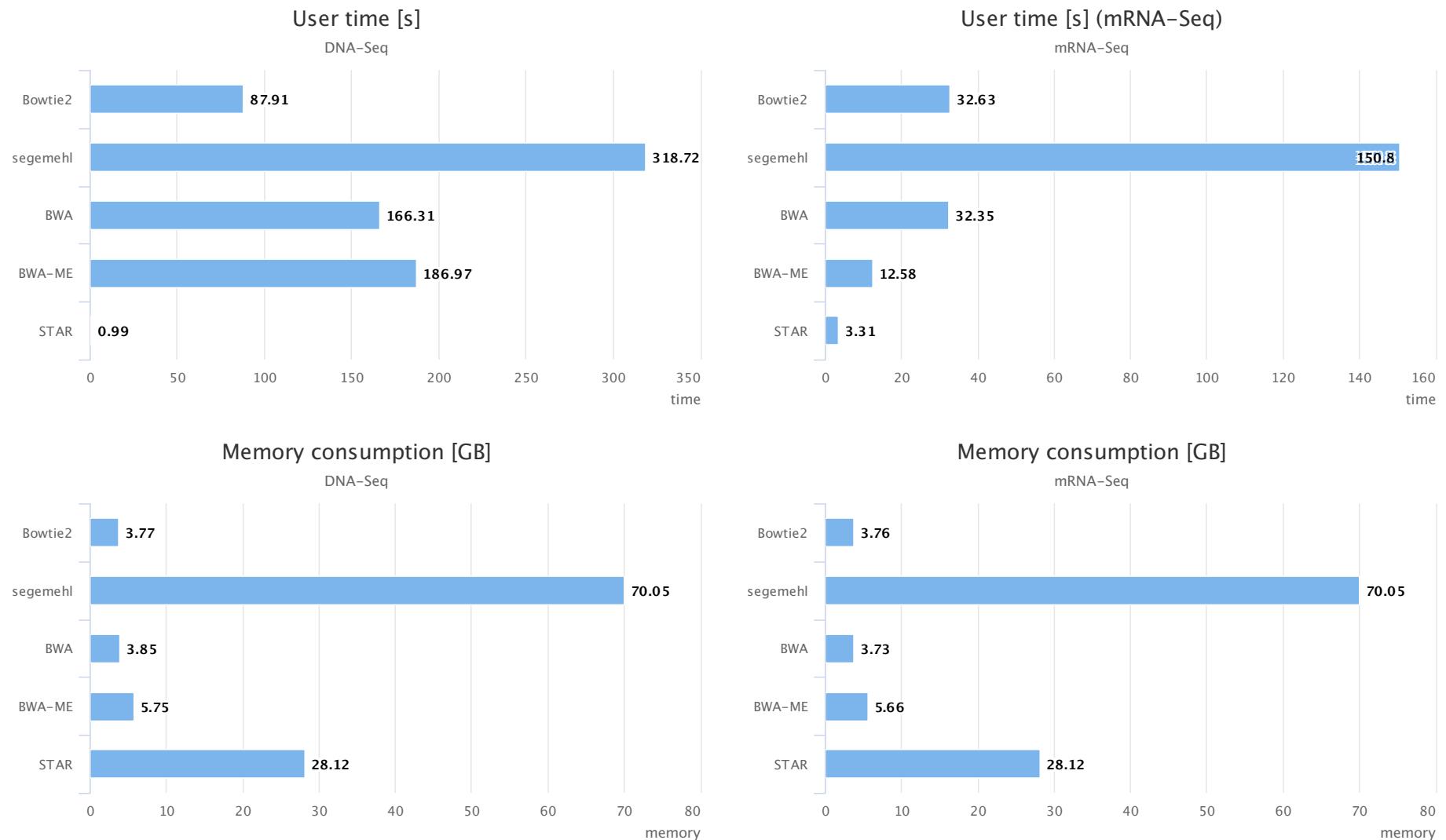
<https://www.ecseq.com/support/ngs/what-is-the-best/ngs-alignment-software>

Reference based aligners - Overview

	Time	Precision	Pairs	GAPS	Phred	Memory	Application (Comments)
BOWTIE	+	+	-	-	5GB	General <i>(max. 3 missmatches)</i>	
BWA	+	+	+	+	8GB	General <i>(max of 200bps reads)</i>	
NOVOALIGN		+	+	+	+	8GB	General <i>(commercial license)</i>
STAR	+	+	-	+	32GB	RNA-Seq <i>(allow split-maps)</i>	
BISMARCK	+	+	+	+	+	10GB	Bisulfite/reduced sequencing

non comprehensive list

Alignment Tools



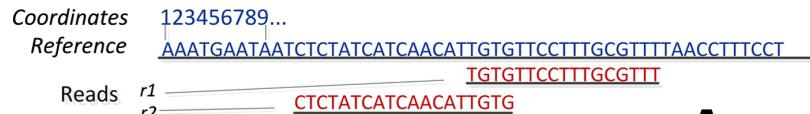
SAM Files

- Store alignment results as text-based file
- Consists of a header and an alignment section

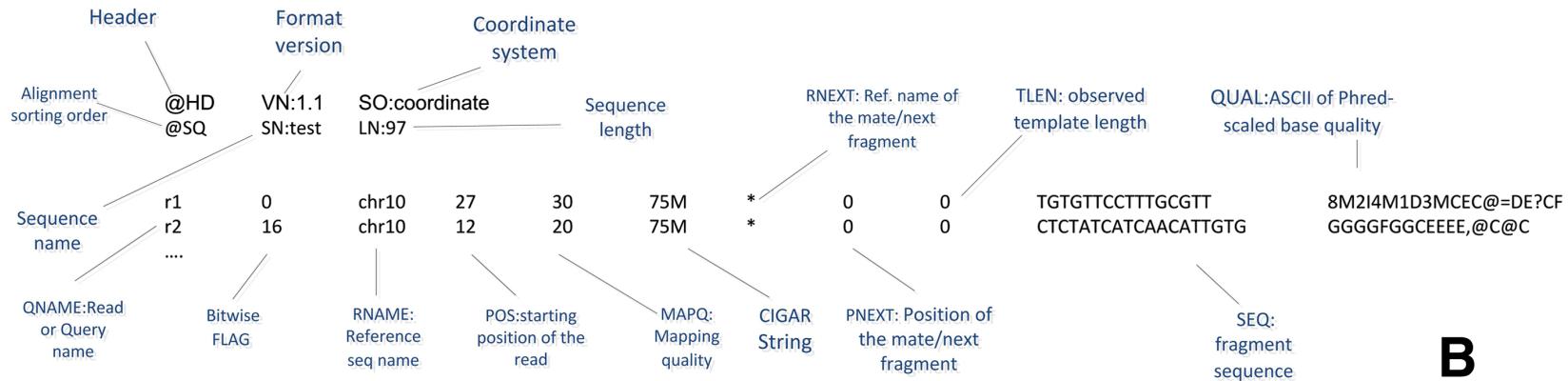
Header											
@HD VN:1.5 SO:coordinate											
@SQ SN:ref LN:45											
r001	99	ref	7	30	8M2I4M1D3M	=	37	39	TTAGATAAAGGATACTG	*	
r002	0	ref	9	30	3S6M1P1I4M	*	0	0	AAAAGATAAGGATA	*	
r003	0	ref	9	30	5S6M	*	0	0	GCCTAAGCTAA	*	SA:Z:ref,29,-,6H5M,17,0;
r004	0	ref	16	30	6M14N5M	*	0	0	ATAGCTTCAGC	*	
r003	2064	ref	29	17	6H5M	*	0	0	TAGGC	*	SA:Z:ref,9,+,5S6M,30,1;
r001	147	ref	37	30	9M	=	7	-39	CAGCGGCAT	*	NM:i:1

Alignment

SAM Files- alignment section

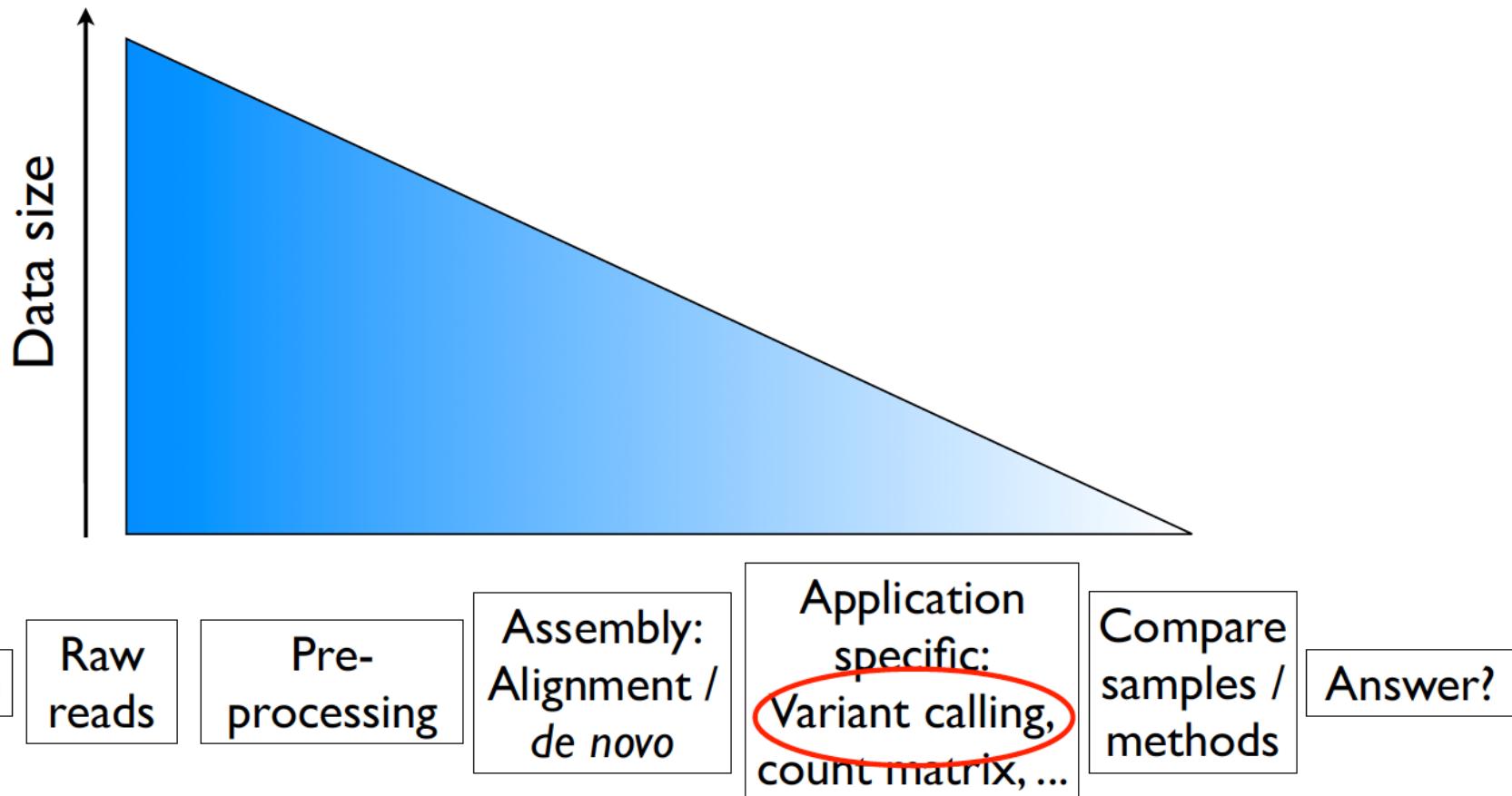


A



- SAM files will have large size (1-10 Gbs)
 - Usually a experiment has dozens of such files
- BAM files (zipped version of SAM) is more common and reduces the size by 30-50%. This file can be opened in genome browsers if a index file is also given.

Applications - Variant Calling



Example Application

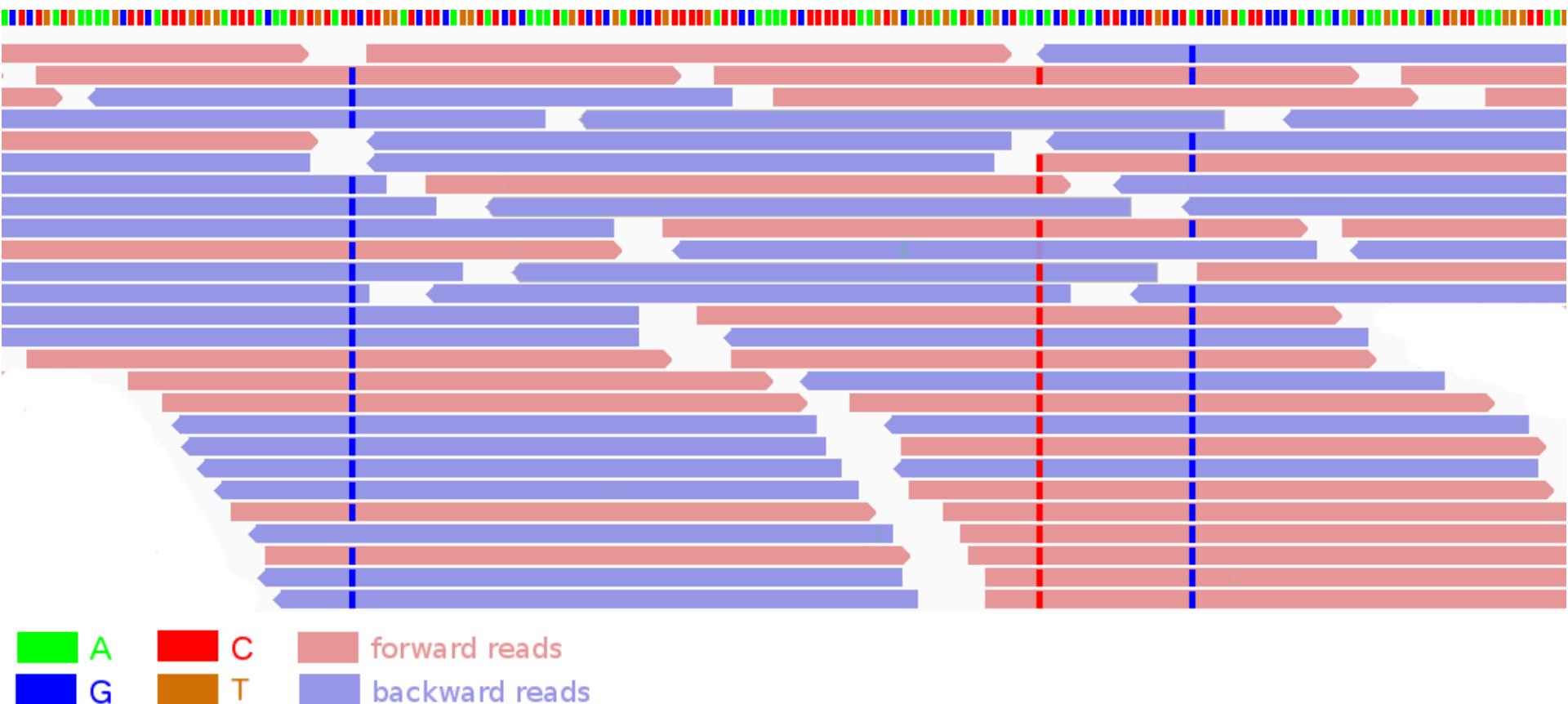
SNP Detection - Individual DNA Sequencing

---TCGTCGTGGTTGAACGTACCGTACCGTTCCCTGAGGCTTAT---

TCCTCGTGGTTGAACGGAA
CGTCGTGGTTGAACGGGAC
CGTGGTTGAACGGACGTA
GTGGTTGAACGGACGTACAG
GTTGAACGGACGTACCGTTCCCTG
TGAACGGACGTACCGTTCC
GAACGGACGTACCGTTCCCTGAGGC
ACGGACGTACAGTTCCCT
GGACGTACCGTTCC
GTACCCTGAG--TTA
TCCCTGAG--TTA
CCTGAGGCTTAT

Alignment of short reads on DNA

SNP Calling Example - Genome Browser

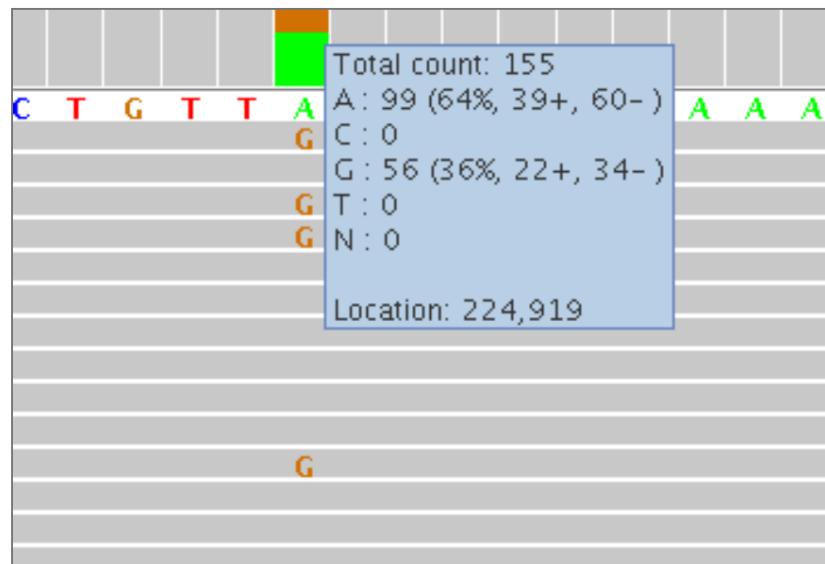
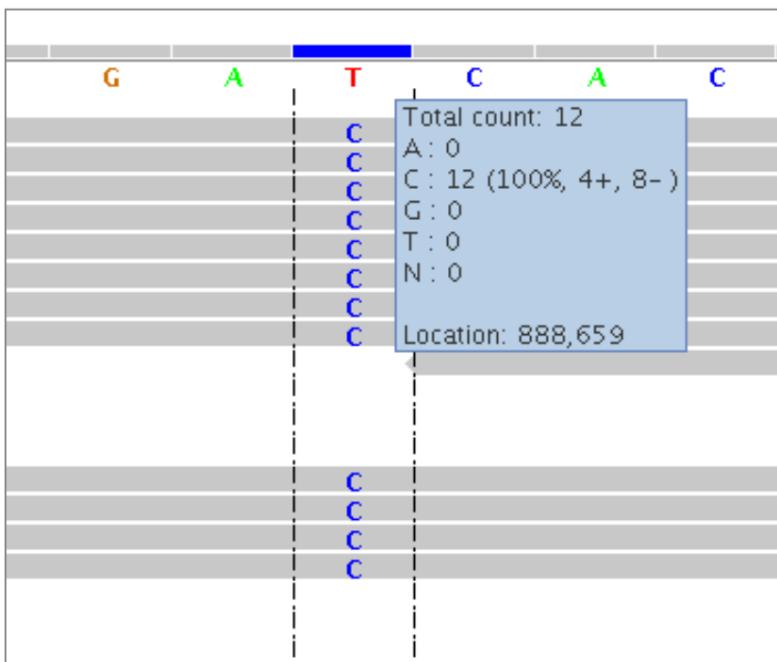


IGV only show miss-matches!

Simple SNP Calling Method

◆ Aims

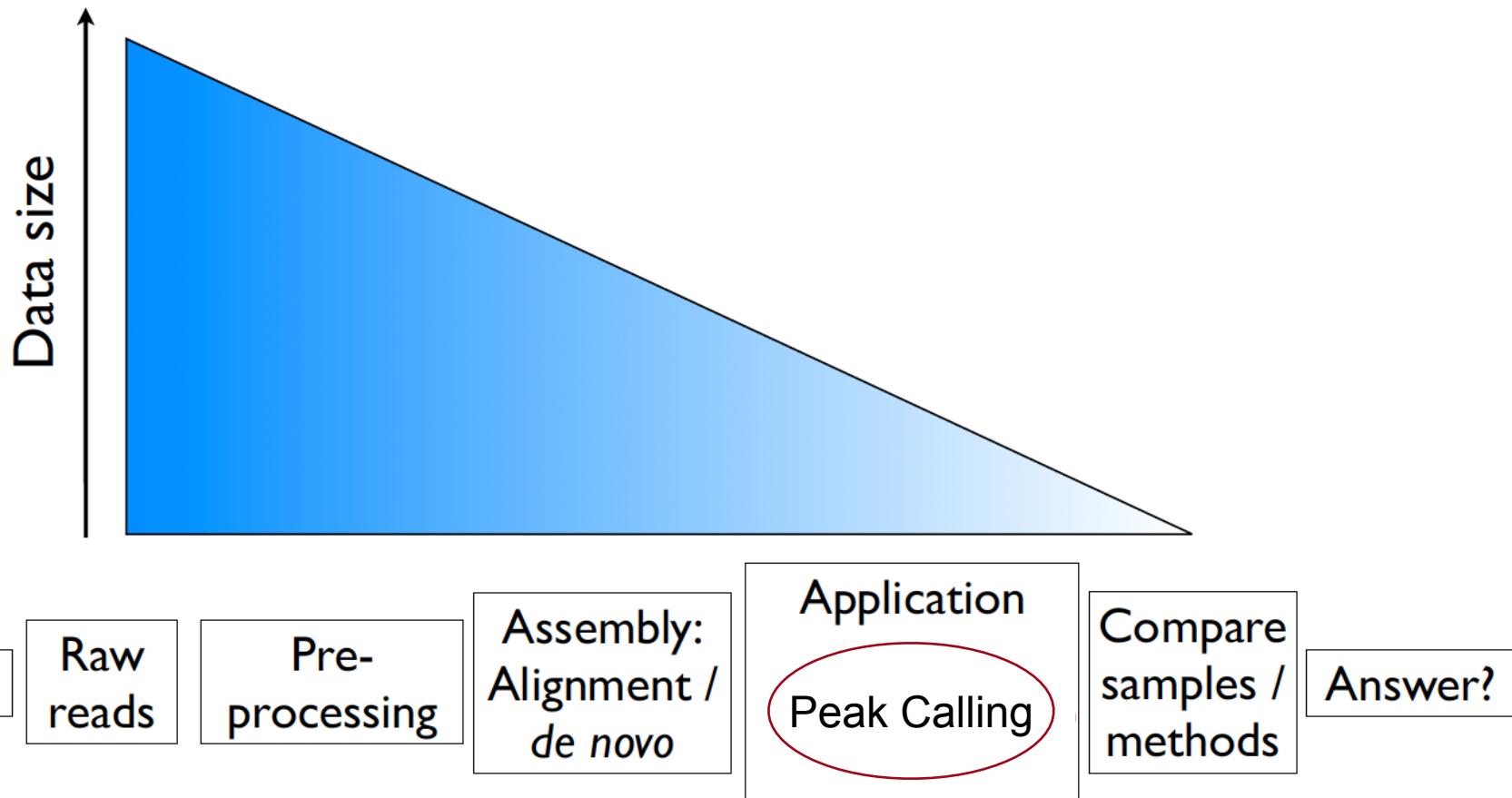
- ◆ Variant calling: Identify polymorphic sites => sites that differs from the reference.
- ◆ Genotyping: Determine the genotype for a certain individual at such sites.



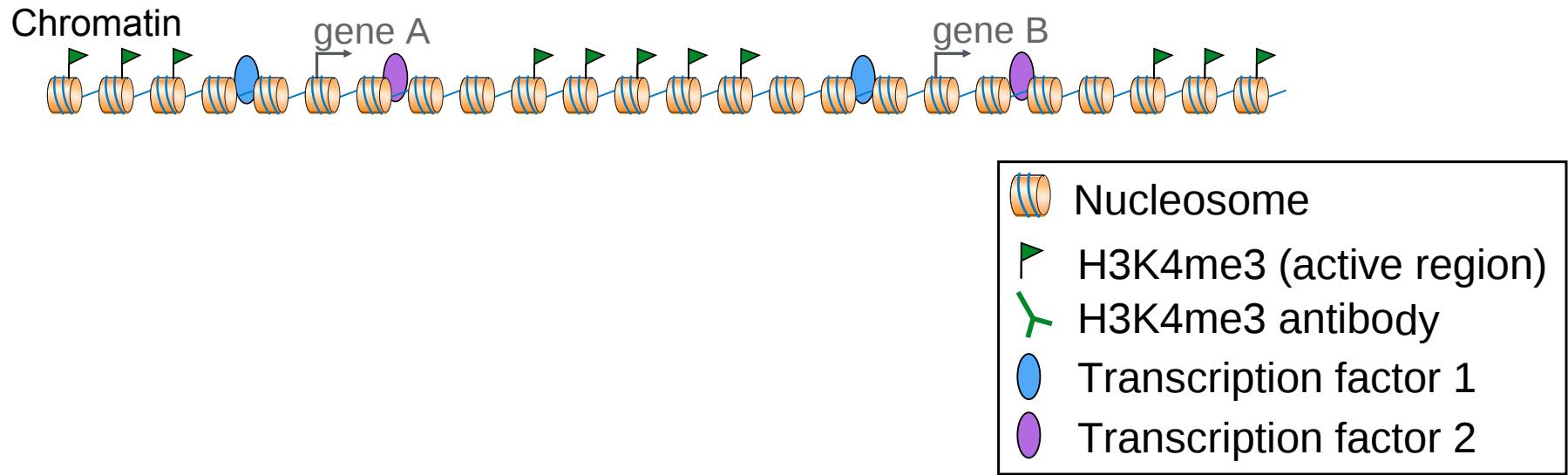
◆ Early methods

- ◆ Works by simply counting the alleles at each site, and then identifying a variant by use of simple cutoff rules.

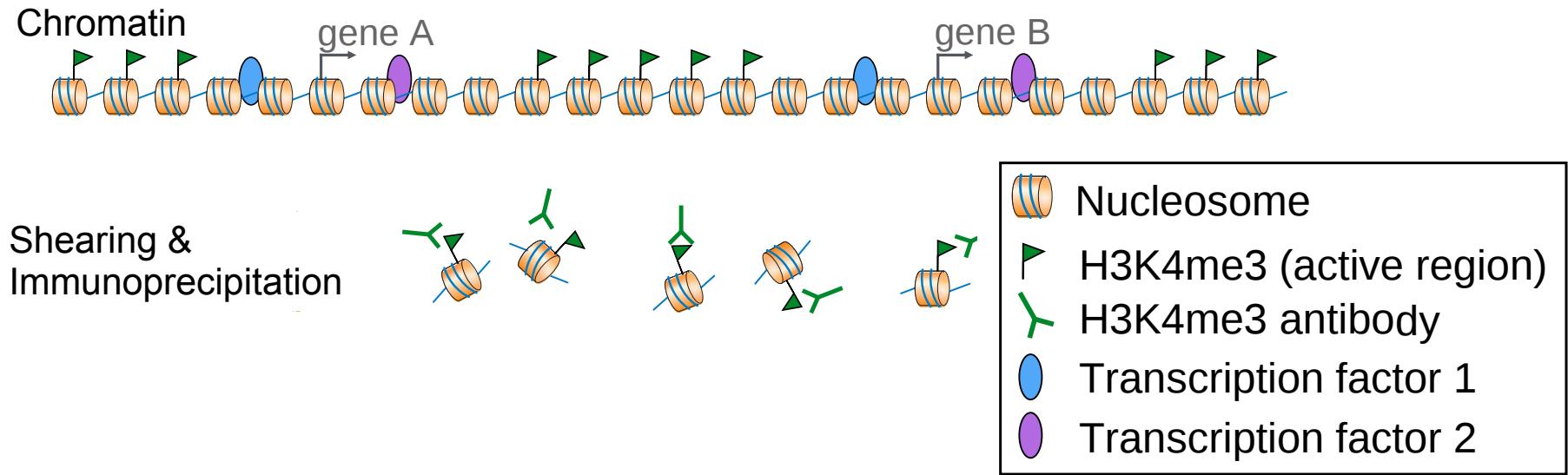
Applications - Peak Calling



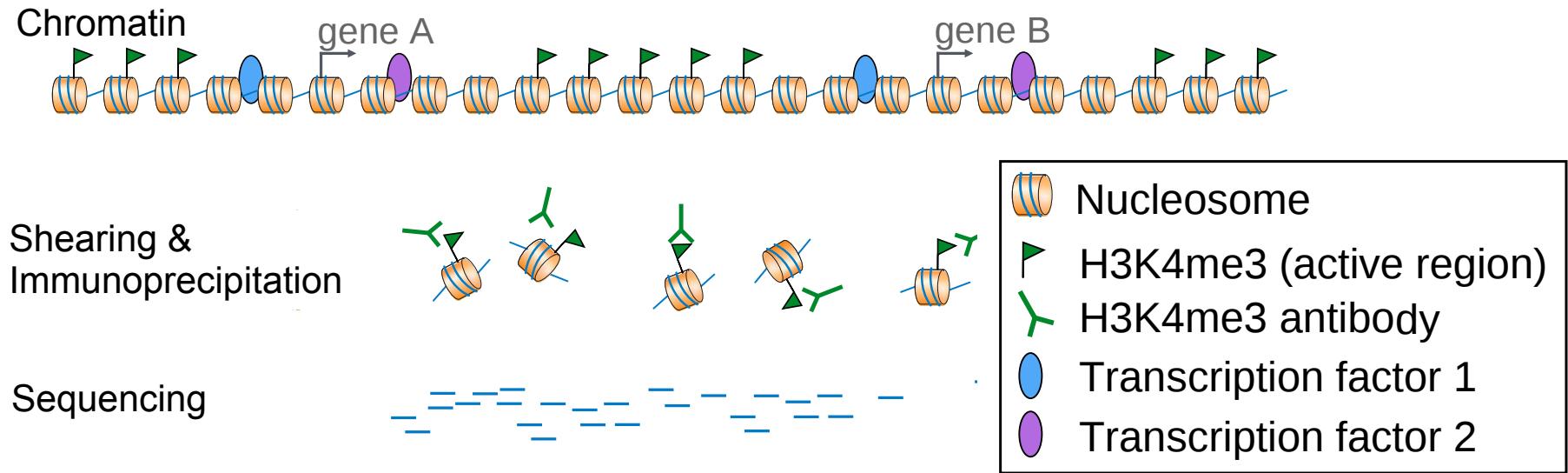
DNA - Protein interactions with ChIP-Seq



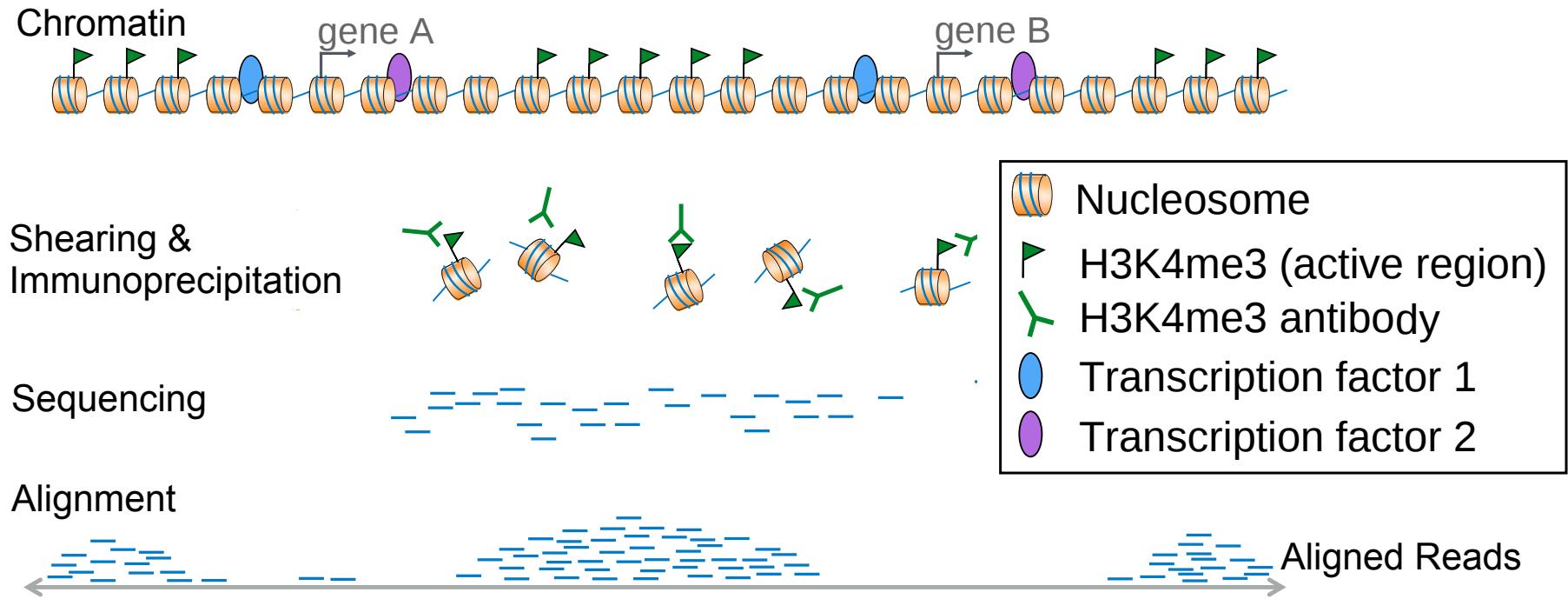
DNA - Protein interactions with ChIP-Seq



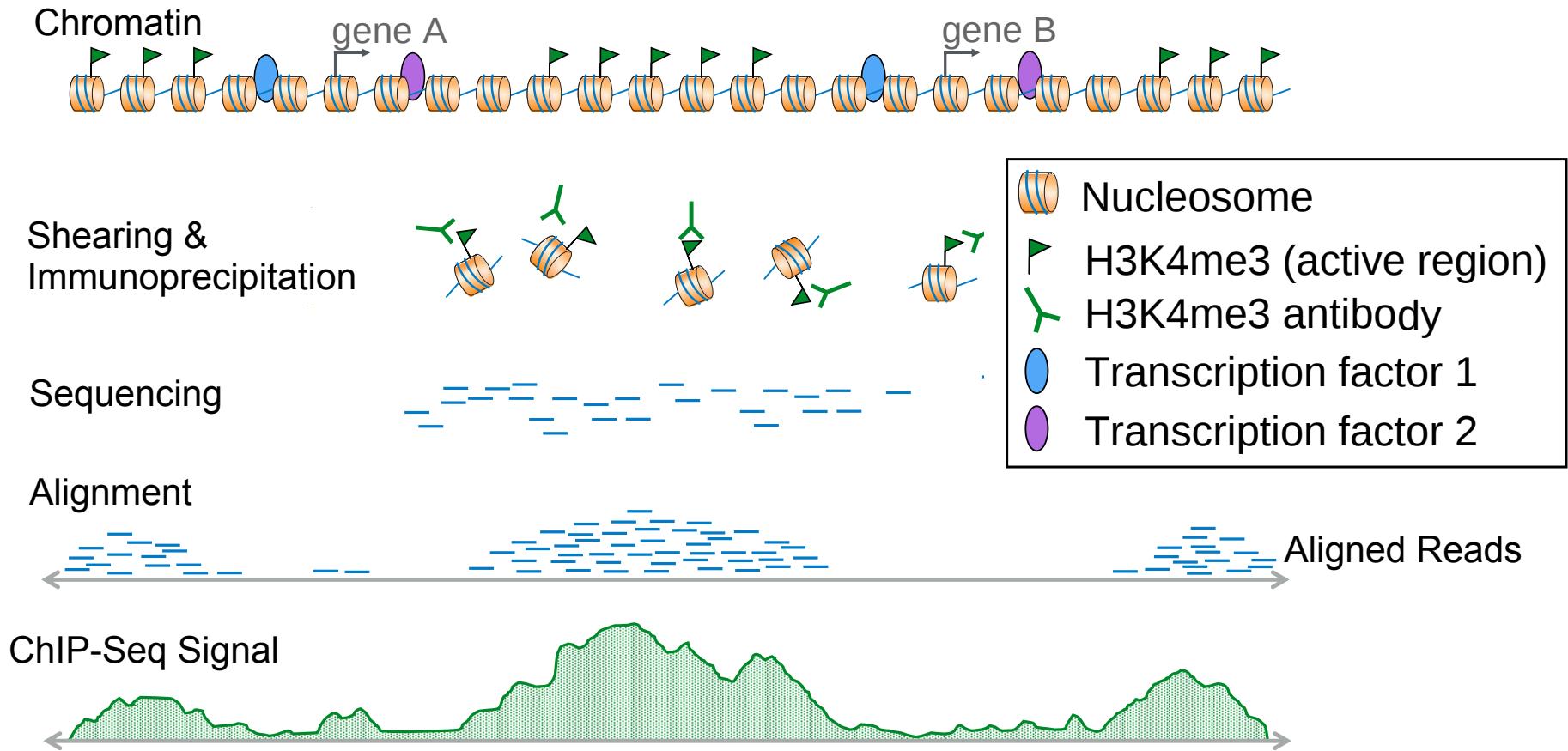
DNA - Protein interactions with ChIP-Seq



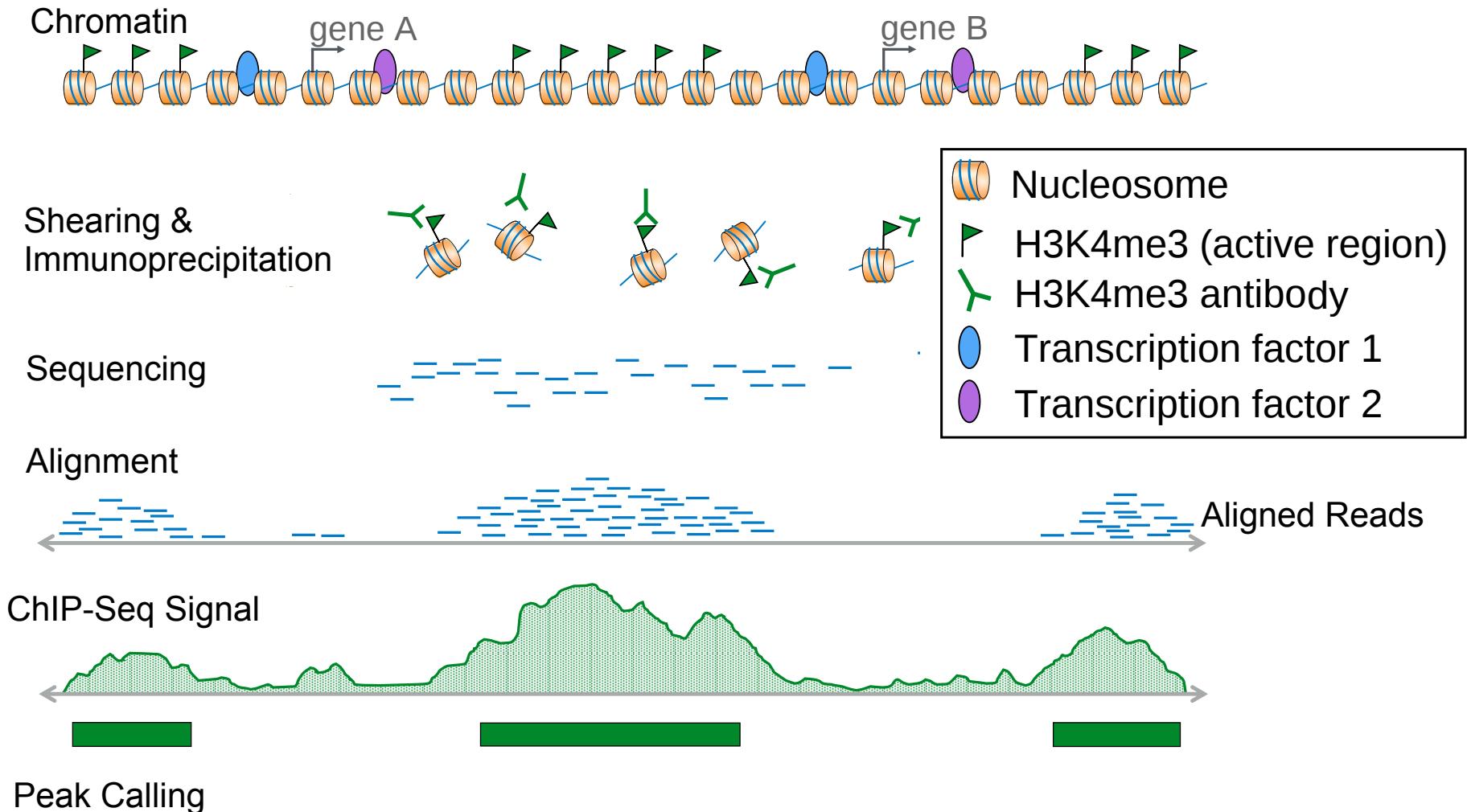
DNA - Protein interactions with ChIP-Seq



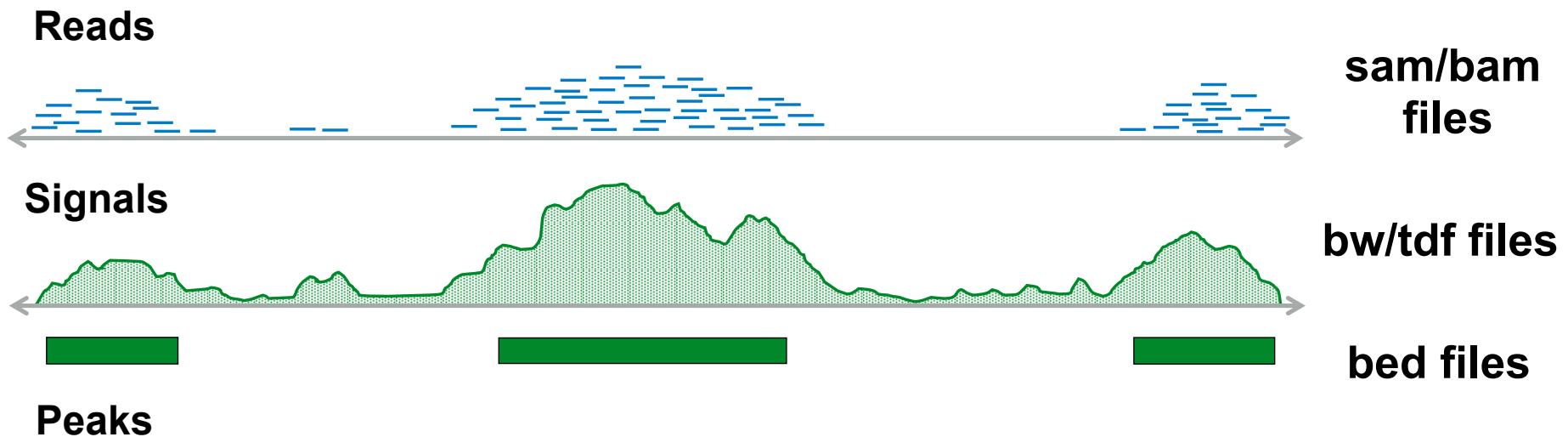
DNA - Protein interactions with ChIP-Seq



DNA - Protein interactions with ChIP-Seq



ChIP-Seq - Data Files



Peaks - Bed files

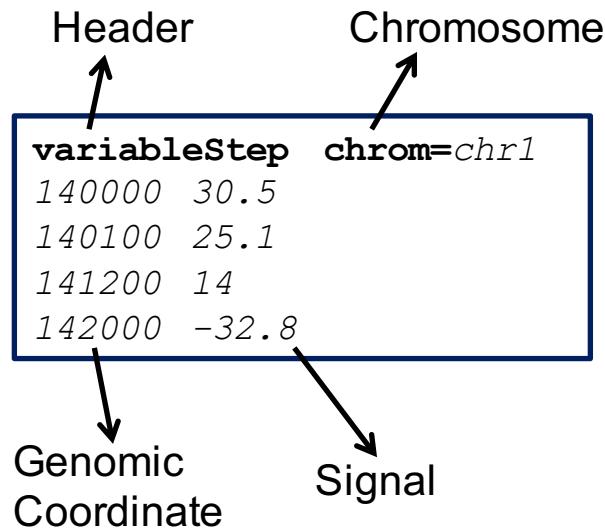
- Peaks / genomic regions are stored in bed files

Chromosome	Start	End	Name	Score	Strand
<i>chr7</i>	127471196	127472363	<i>Peak1</i>	0	+
<i>chr7</i>	127472363	127473530	<i>Peak2</i>	0	+
<i>chr7</i>	127473530	127474697	<i>Peak3</i>	0	+
<i>chr7</i>	127474697	127475864	<i>Peak4</i>	0	+
<i>chr7</i>	127475864	127477031	<i>Peak5</i>	0	-
<i>chr7</i>	127477031	127478198	<i>Peak6</i>	0	-
<i>chr7</i>	127478198	127479365	<i>Peak7</i>	0	-
<i>chr7</i>	127479365	127480532	<i>Peak8</i>	0	+
<i>chr7</i>	127480532	127481699	<i>Peak9</i>	0	-

Genomic Signals - WIG/TDF Files

- Files containing smoothed counts of reads for ChIP-seq, RNA-seq, or similar protocols.

- Example of WIG file



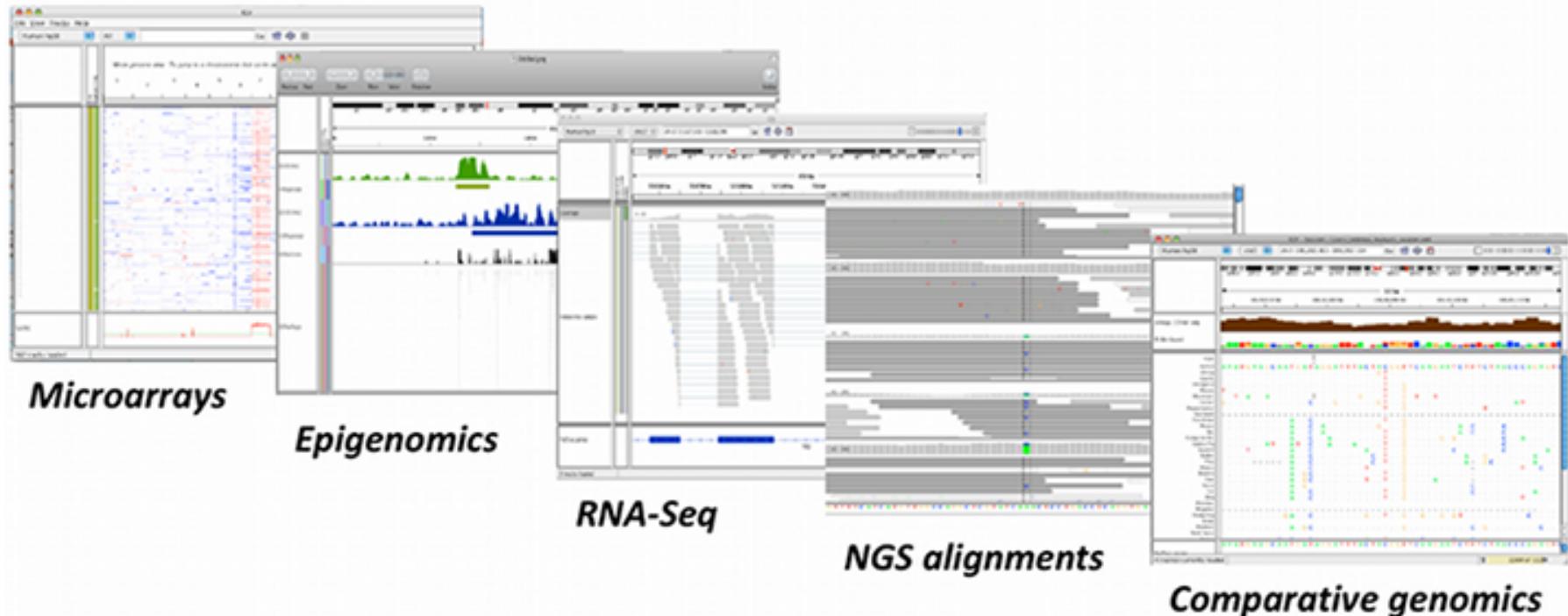
- In practice, we use binary version of WIG (BIGWIG) or TDF files

Bioinformatics Analysis in R

Next Generation Sequencing Data Visualization

IGV (Integrative Genome Viewer)

- Desktop application for the **visual interactive** exploration of **integrated** genomic datasets



Advantages

- A high-performance visualization tool
- Allows us interactively explore large, integrated dataset
- Supports a wide variety of data types, including microarray and next-generation sequencing data
- **FREE**

Launch IGV

<http://software.broadinstitute.org/software/igv/home>

The screenshot shows the "Home" page of the IGV website. On the left, there is a sidebar with a red circle highlighting the "Downloads" and "Documents" menu items. The main content area features a large image of the IGV software interface displaying genomic tracks. Below the image, there are two sections: "Overview" and "Citing IGV". The "Overview" section contains text about the tool's purpose and capabilities, along with a small icon of a computer monitor showing genomic data. The "Citing IGV" section provides instructions for citation and lists several academic publications. At the bottom of the page, there are links for "Download IGV" and "Funding".

Launch IGV

Downloads | Integrative Genomics Viewer

Home > Downloads

Downloads

Did you know that there is also an **IGV web application** that runs only in a web browser, does not use Java, and requires no downloads? See <https://igv.org/app>. Click on the [Help](#) link in the app for more information about using IGV-Web.

Install IGV 2.7.x

See the [Release Notes](#) for what's new in each release.

IGV Mac App

Download and unzip the Mac App Archive, then double-click the IGV application to run it. You can move the app to the *Applications* folder, or anywhere else.

MacOS Catalina users: We sign our Mac App as a trusted Apple developer, but it is not yet notarized by Apple (a new requirement in Catalina). To run it, right-click on the downloaded IGV app; select "Open" from the menu; and click the "Open" button in the window that pops up. After that, double-clicking on the app will work.

IGV for Windows

Download and run the installer. An IGV shortcut will be created on the Desktop; double-click it to run the application.

IGV for Linux

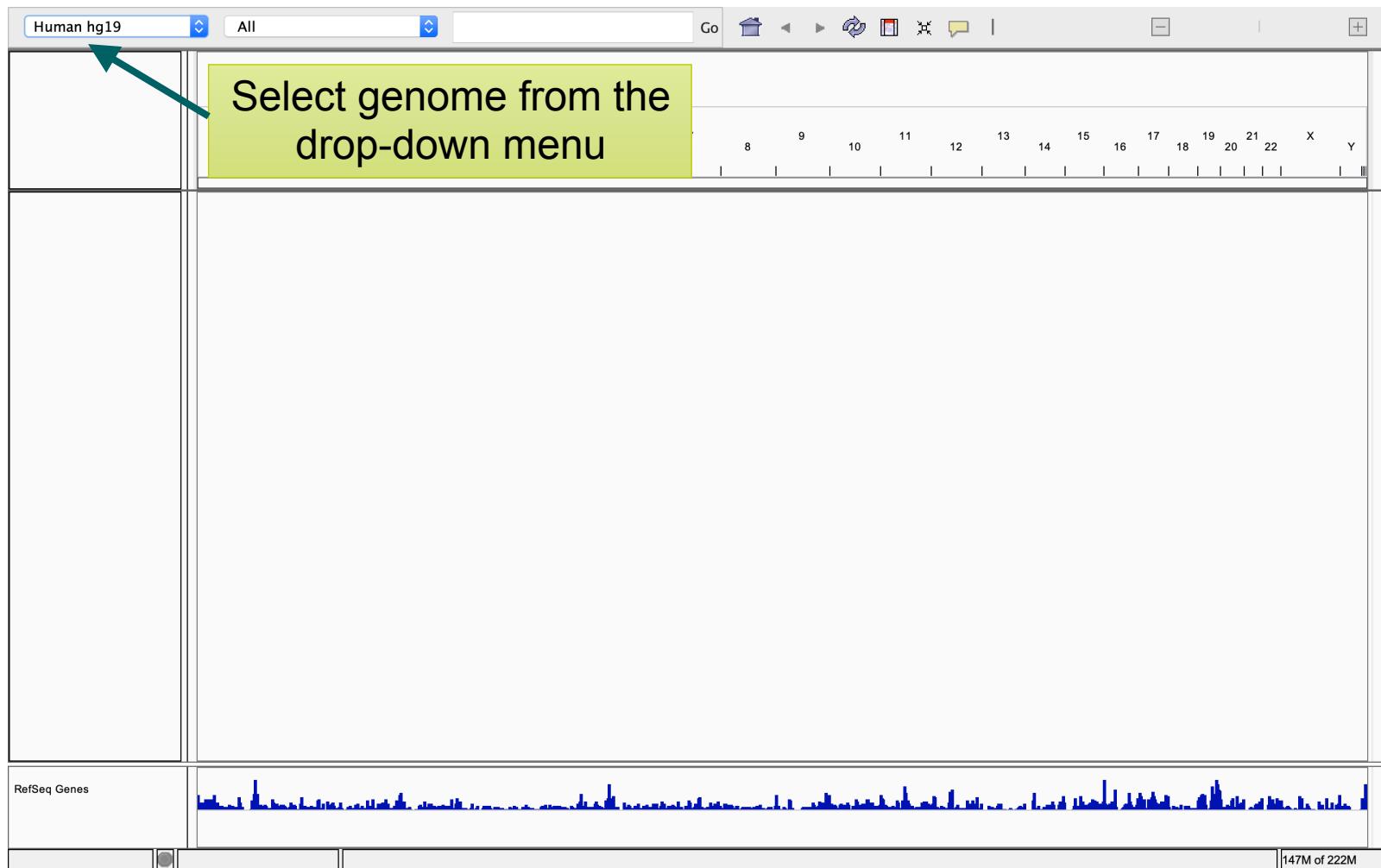
Download and unzip the Archive. See the downloaded *readme.txt* for further instructions.

IGV and igvtools to run on the command line (all platforms)

Download and unzip the Archive. **Requires Java 11.** See the downloaded *readme.txt* and *igvtools_readme.txt* for further instructions.

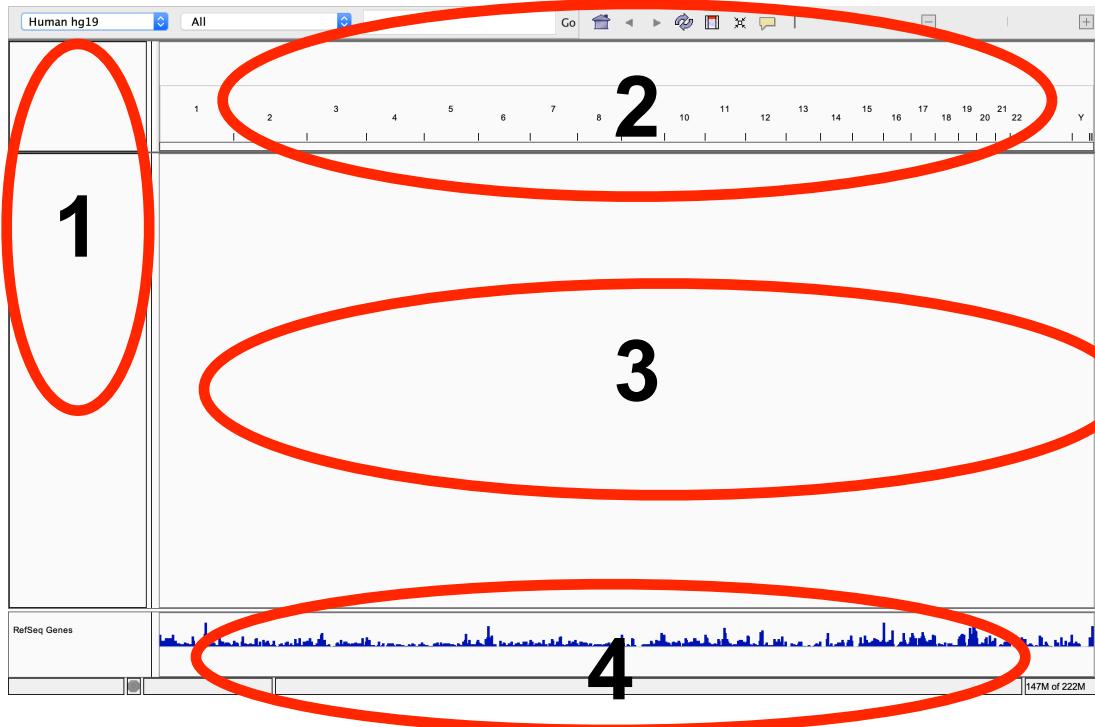
Other IGV Versions

Launch IGV



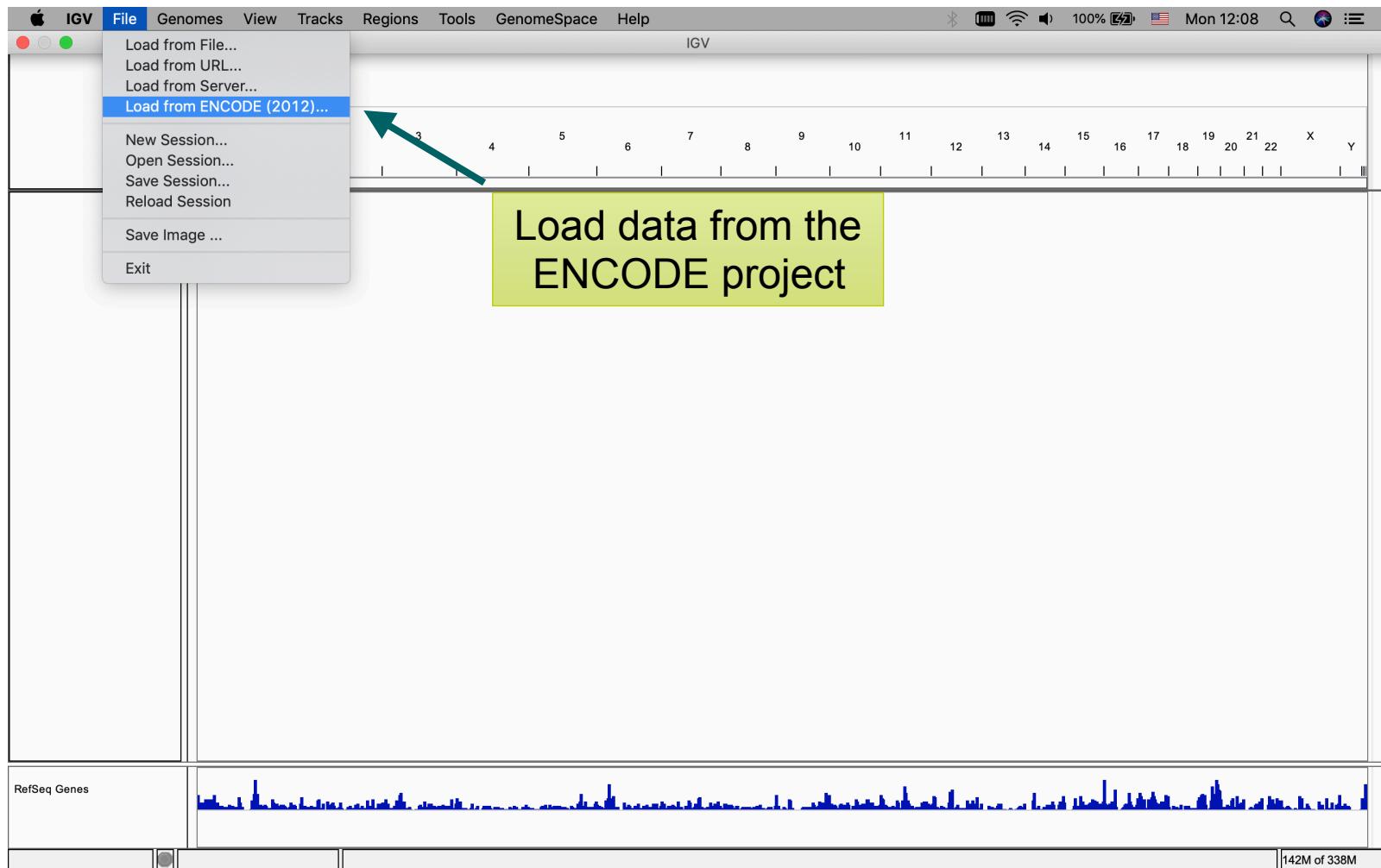
Launch IGV

Genome navigation



Genome locations and sequence

Viewing NGS Data



Viewing NGS Data

The screenshot shows a software interface for viewing NGS data. At the top, there are dropdown menus for "Human hg19" and "All". Below the menu bar is a toolbar with various icons. A central feature is a modal dialog titled "Encode Production Data" with the subtitle "50 rows". The dialog contains a table with columns: cell, dataType, antibody, view, replicate, type, lab, hub. A search bar at the top of the dialog is filled with "ctcf gm12878". A green arrow points from a text box labeled "search CTCF and GM12878" to the search bar. The table data includes multiple entries for GM12878 cells with CTCF antibody, showing various data types like Peaks, Alignments, and Signal, and replicates ranging from 1 to 3. The bottom right of the dialog has "Load" and "Cancel" buttons. Below the dialog, the main window shows a "RefSeq Genes" track with blue peaks representing signal across chromosomes 15 through 22. A status bar at the bottom right indicates "205M of 338M".

cell	dataType	antibody	view	replicate	type	lab	hub
GM12878	ChIP-Seq	CTCF	Peaks		narrowPeak	Broad	Data
GM12878	ChIP-Seq	CTCF_(SC...)					
GM12878	ChIP-Seq	CTCF					
GM12878	ChIP-Seq	CTCF					
GM12878	ChIP-Seq	CTCF					
GM12878	ChIP-Seq	CTCF	Alignments	2	bam	Broad	Data
GM12878	ChIP-Seq	CTCF	Peaks		broadPeak	Broad	Data
GM12878	ChIP-Seq	CTCF	Signal		bigWig	Broad	Data
GM12878	ChIP-Seq	CTCF	Alignments	1	bam	UT-A	Data
GM12878	ChIP-Seq	CTCF	Alignments	2	bam	UT-A	Data
GM12878	ChIP-Seq	CTCF	Alignments	3	bam	UT-A	Data
GM12878	ChIP-Seq	CTCF	Base_Ove...		bigWig	UT-A	Data
GM12878	ChIP-Seq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChIP-Seq	CTCF	Signal		bigWig	UT-A	Data
GM12878	ChIP-Seq	CTCF_(SC...)	Alignments	1	bam	Stanford	Data
GM12878	ChIP-Seq	CTCF_(SC...)	Alignments	2	bam	Stanford	Data
GM12878	ChIP-Seq	CTCF_(SC...)	Peaks		narrowPeak	Stanford	Data
GM12878	ChIP-Seq	CTCF_(SC...)	Signal		bigWig	Stanford	Data
GM12878	ChIP-Seq	CTCF	Alignments	1	bam	UW	Data
GM12878	ChIP-Seq	CTCF	Alignments	2	bam	UW	Data
GM12878	ChIP-Seq	CTCF	Hotspots	1	broadPeak	UW	Data
GM12878	ChIP-Seq	CTCF	Hotspots	2	broadPeak	UW	Data
GM12878	ChIP-Seq	CTCF	Peaks	1	narrowPeak	UW	Data
GM12878	ChIP-Seq	CTCF	Peaks	2	narrowPeak	UW	Data
GM12878	ChIP-Seq	CTCF	RawSignal	1	bigWig	UW	Data
GM12878	ChIP-Seq	CTCF	RawSignal	2	bigWig	UW	Data
GM12878	ChIP-Seq	CTCF	Peaks		bigBed	Broad	analysis
GM12878	ChIP-Seq	CTCF	Peaks		bigBed	UT-A	analysis
GM12878	ChIP-Seq	CTCF_(C...	Peaks		bigBed	Stanford	analysis
GM12878	ChIP-Seq	CTCF	Peaks		bigBed	UW	analysis

Viewing NGS Data

The screenshot shows a software interface for viewing NGS data. At the top, there are dropdown menus for "Human hg19" and "All". Below the menu bar is a toolbar with various icons. A central feature is a modal dialog titled "Encode Production Data" containing a table with 50 rows of data. The table has columns for cell, dataType, antibody, view, replicate, type, lab, and hub. A green arrow points to the "lab" column header, and a yellow box highlights the word "sort by lab". The "hub" column contains entries like "Data", "Data", "Data", "Data", "Data", "analysis", "analysis", "analysis", "Data", "Data", etc. At the bottom of the dialog are "Load" and "Cancel" buttons. In the background, there's a genomic track showing RefSeq Genes with blue peaks, and a status bar at the bottom right indicates "205M of 299M".

cell	dataType	antibody	view	replicate	type	lab	hub
GM12878	ChIPSeq	CTCF	Peaks		narrowPeak	Broad	Data
GM12878	ChIPSeq	CTCF	Alignments	1	bam	Broad	Data
GM12878	ChIPSeq	CTCF	Alignments	2	bam	Broad	Data
GM12878	ChIPSeq	CTCF	Peaks		broadPeak	Broad	Data
GM12878	ChIPSeq	CTCF	Signal		bigWig	Broad	Data
GM12878	ChIPSeq	CTCF	Peaks		bigBed	Broad	Data
GM12878	ChIPSeq	CTCF	Peaks		bigBed	Broad	Data
GM12878	ChIPSeq	CTCF	Signal		bigWig	Broad	Data
GM12878	ChIPSeq	CTCF_(SC...	Peaks		narrowPeak	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...	Alignments	1	bam	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...	Alignments	2	bam	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...	Peaks		narrowPeak	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...	Signal		bigWig	Stanford	Data
GM12878	ChIPSeq	CTCF_(C-...	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...	Peaks		bigWig	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...	Signal		bigWig	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...	Signal		bigWig	Stanford	analysis
GM12878	ChIPSeq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChIPSeq	CTCF	Alignments	1	bam	UT-A	Data
GM12878	ChIPSeq	CTCF	Alignments	2	bam	UT-A	Data
GM12878	ChIPSeq	CTCF	Alignments	3	bam	UT-A	Data
GM12878	ChIPSeq	CTCF	Base_Ove...		bigWig	UT-A	Data
GM12878	ChIPSeq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChIPSeq	CTCF	Signal		bigWig	UT-A	Data
GM12878	ChIPSeq	CTCF	Peaks		bigBed	UT-A	Data
GM12878	ChIPSeq	CTCF	Peaks				analysis

Viewing NGS Data

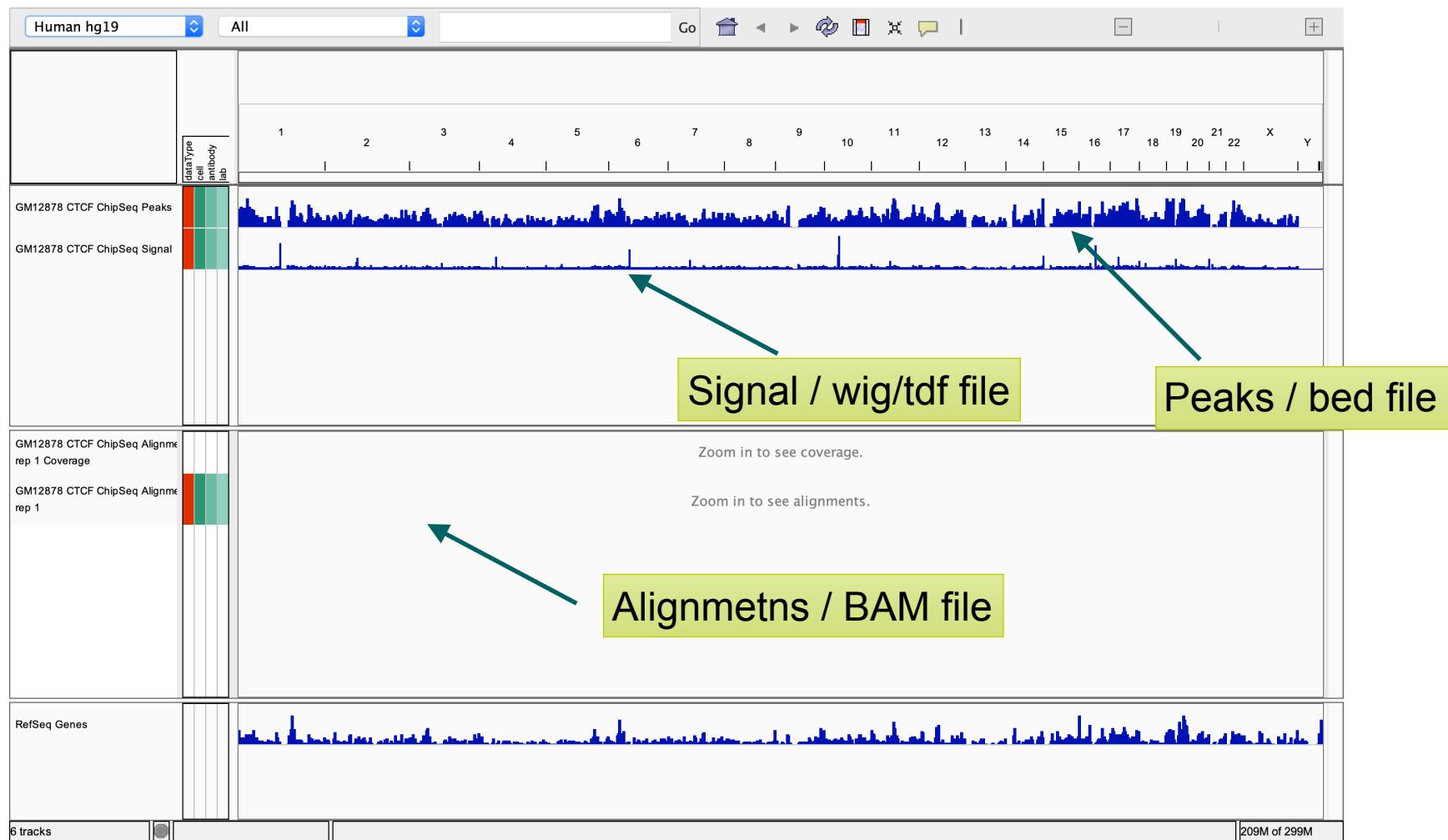
select alignments, signal and peaks

load data

The screenshot shows a software interface for viewing NGS data. On the left, there's a tree view of tracks: "Human hg19" (selected), "GM12878 CTCF ChIPSeq Peaks", "GM12878 CTCF ChIPSeq Signal", "GM12878 CTCF ChIPSeq Alignments rep 1 Coverage", "GM12878 CTCF ChIPSeq Alignments rep 1", and "RefSeq Genes". Below the tracks are status bars: "6 tracks" and "202M of 299M". A central window is titled "Encode Production Data" with a filter bar containing "ctcf gm12878". It lists 50 rows of data with columns: cell, dataType, antibody, view, replicate, type, lab, hub, and a detailed description. Several rows for "GM12878" are selected. Arrows point from the text "select alignments, signal and peaks" to the "alignments" rows in the table, and from "load data" to the "Load" button at the bottom of the dialog. The right side of the interface shows genomic tracks for chromosomes 13 through Y.

cell	dataType	antibody	view	replicate	type	lab	hub
GM12878	ChIPSeq	CTCF	Peaks		narrowPeak	Broad	Data
GM12878	ChIPSeq	CTCF	Alignments	1	bam	Broad	Data
GM12878	ChIPSeq	CTCF	Alignments	2	bam	Broad	Data
GM12878	ChIPSeq	CTCF	Peaks		broadPeak	Broad	Data
GM12878	ChIPSeq	CTCF	Signal		bigWig	Broad	Data
GM12878	ChIPSeq	CTCF	Peaks		bigBed	Broad	analysis
GM12878	ChIPSeq	CTCF	Peaks		bigBed	Broad	analysis
GM12878	ChIPSeq	CTCF	Peaks		bigWig	Broad	analysis
GM12878	ChIPSeq	CTCF	Peaks		bigBed	Broad	analysis
GM12878	ChIPSeq	CTCF	Peaks		bigWig	Broad	analysis
GM12878	ChIPSeq	CTCF_(SC...)	Peaks		narrowPeak	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...)	Alignments	1	bam	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...)	Alignments	2	bam	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...)	Peaks		narrowPeak	Stanford	Data
GM12878	ChIPSeq	CTCF_(SC...)	Signal		bigWig	Stanford	Data
GM12878	ChIPSeq	CTCF_(C-...)	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...)	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...)	Signal		bigWig	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...)	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...)	Peaks		bigBed	Stanford	analysis
GM12878	ChIPSeq	CTCF_(C-...)	Signal		bigWig	Stanford	analysis
GM12878	ChIPSeq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChIPSeq	CTCF	Alignments	1	bam	UT-A	Data
GM12878	ChIPSeq	CTCF	Alignments	2	bam	UT-A	Data
GM12878	ChIPSeq	CTCF	Alignments	3	bam	UT-A	Data
GM12878	ChIPSeq	CTCF	Base_Ove...		bigWig	UT-A	Data
GM12878	ChIPSeq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChIPSeq	CTCF	Signal		bigWig	UT-A	Data
GM12878	ChIPSeq	CTCF	Peaks		bigBed	UT-A	analysis

Viewing NGS Data



File Formats

- The **file format** defines the track type
- The track type determines the display options
- IGV supports many file formats
 - **BAM**
 - **BED**
 - BedGraph
 - bigBed
 - bigWig
 - Birdsuite Files
 - broadPeaks
 - CBS
 - CN
 - Cufflinks Files
 - FASTA
 - GCT
 - genePred
 - GFF
 - GISTIC
 - Goby
 - GWAS
 - IGV
 - LOH
 - MAF
 - MUT
 - narrowPeaks
 - PSL
 - RES
 - SAM
 - SEG
 - SNP
 - TAB
 - **TDF**
 - TrackLine
 - TypeLine
 - VCF
 - WIG

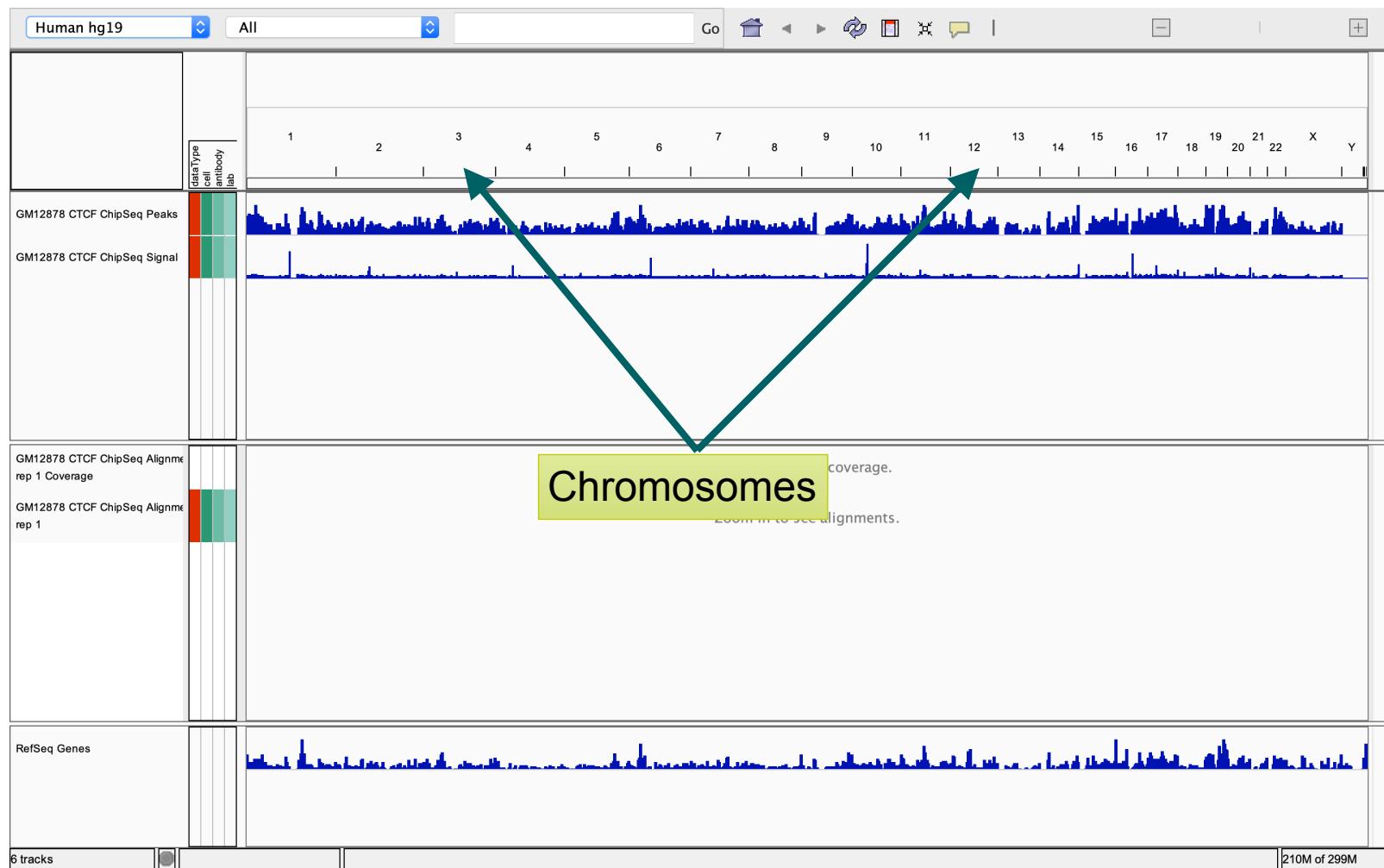
File Formats

Note: for large files use indexed formats

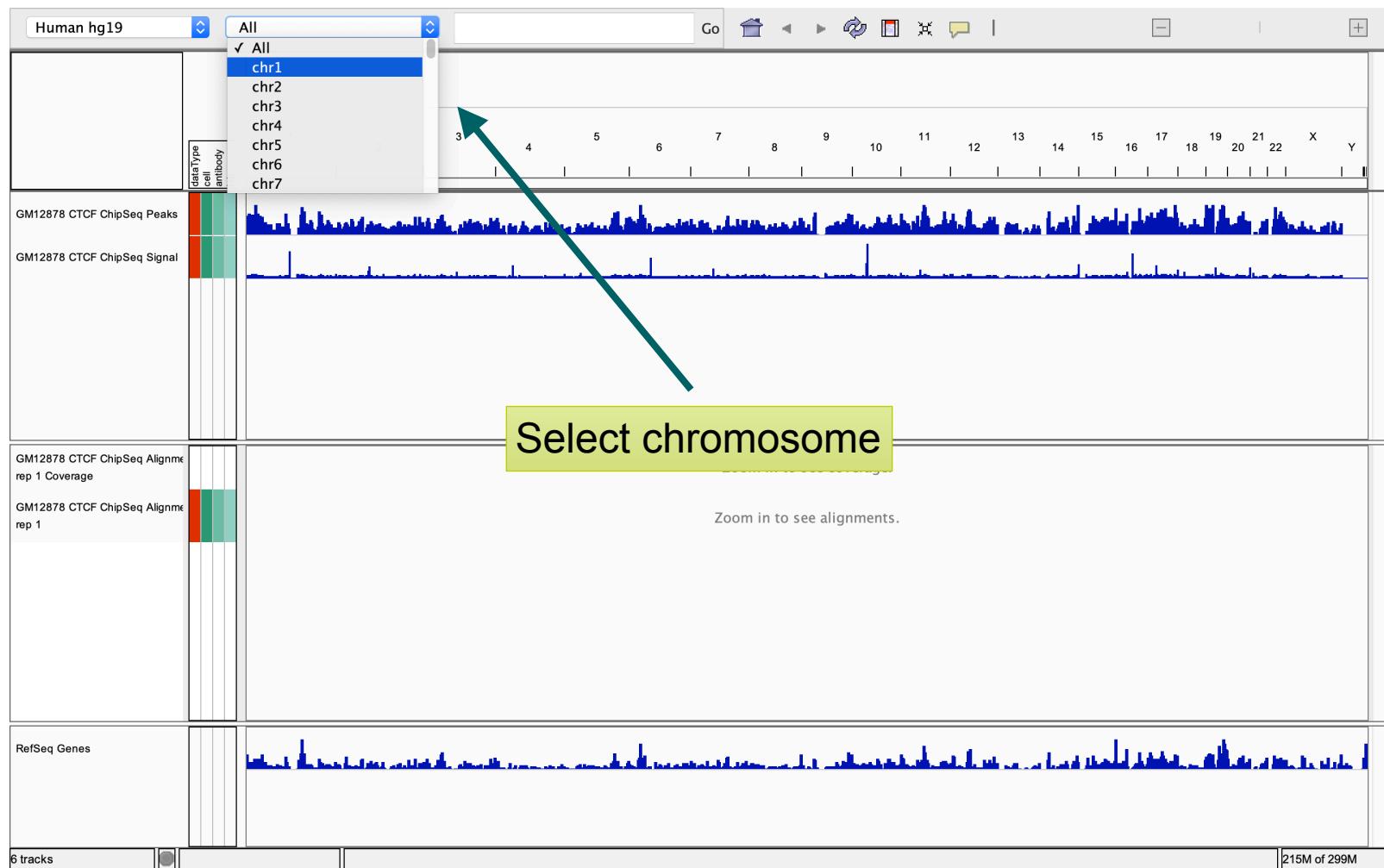
Hands on

- Launch IGV on your computer
- Choose human genome hg19
- Load data from ENCODE project
 - ChIP-seq of factor CTCF for GM12878 cell type
 - Alignments, peaks and signal

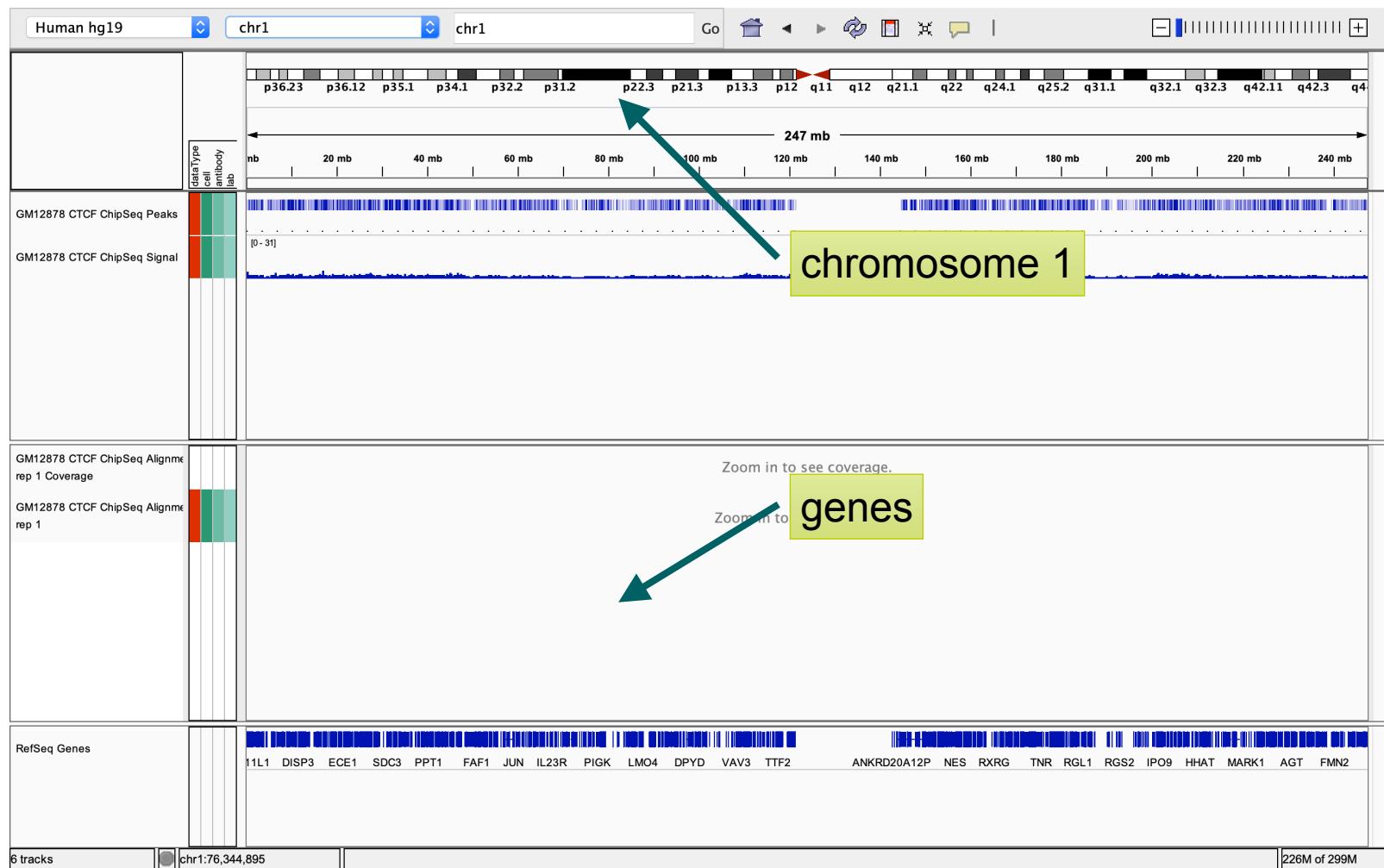
Navigation



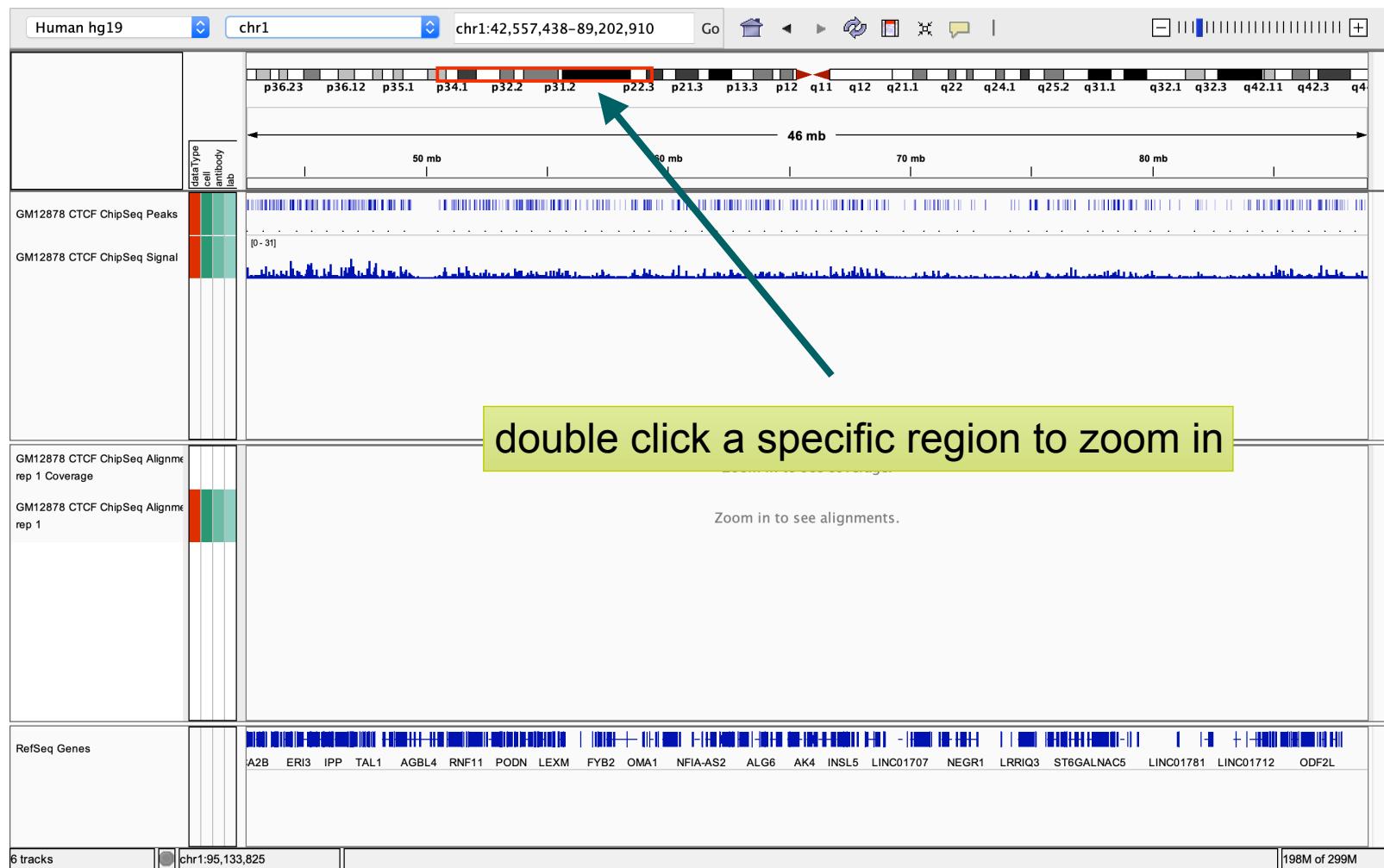
Navigation



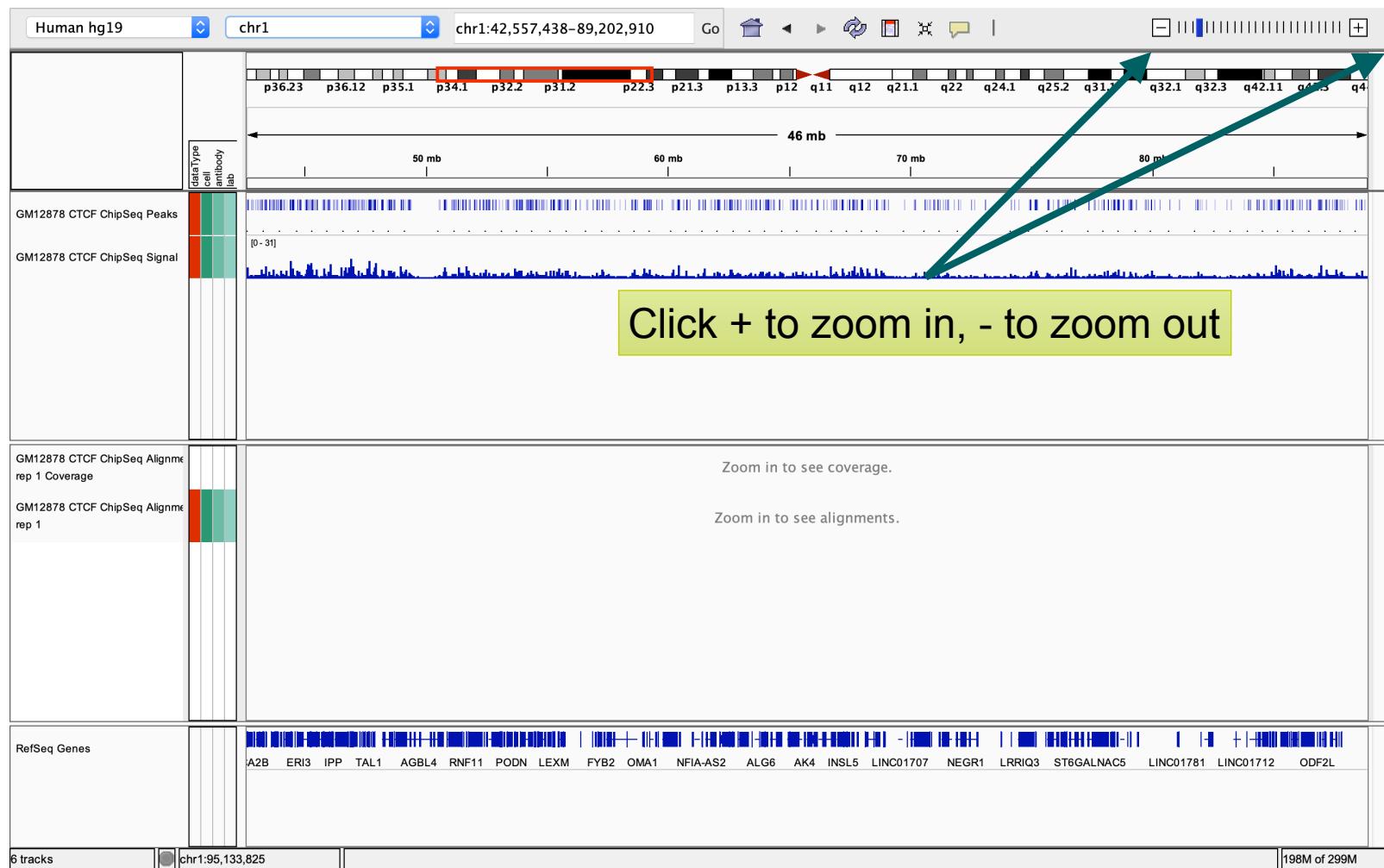
Navigation



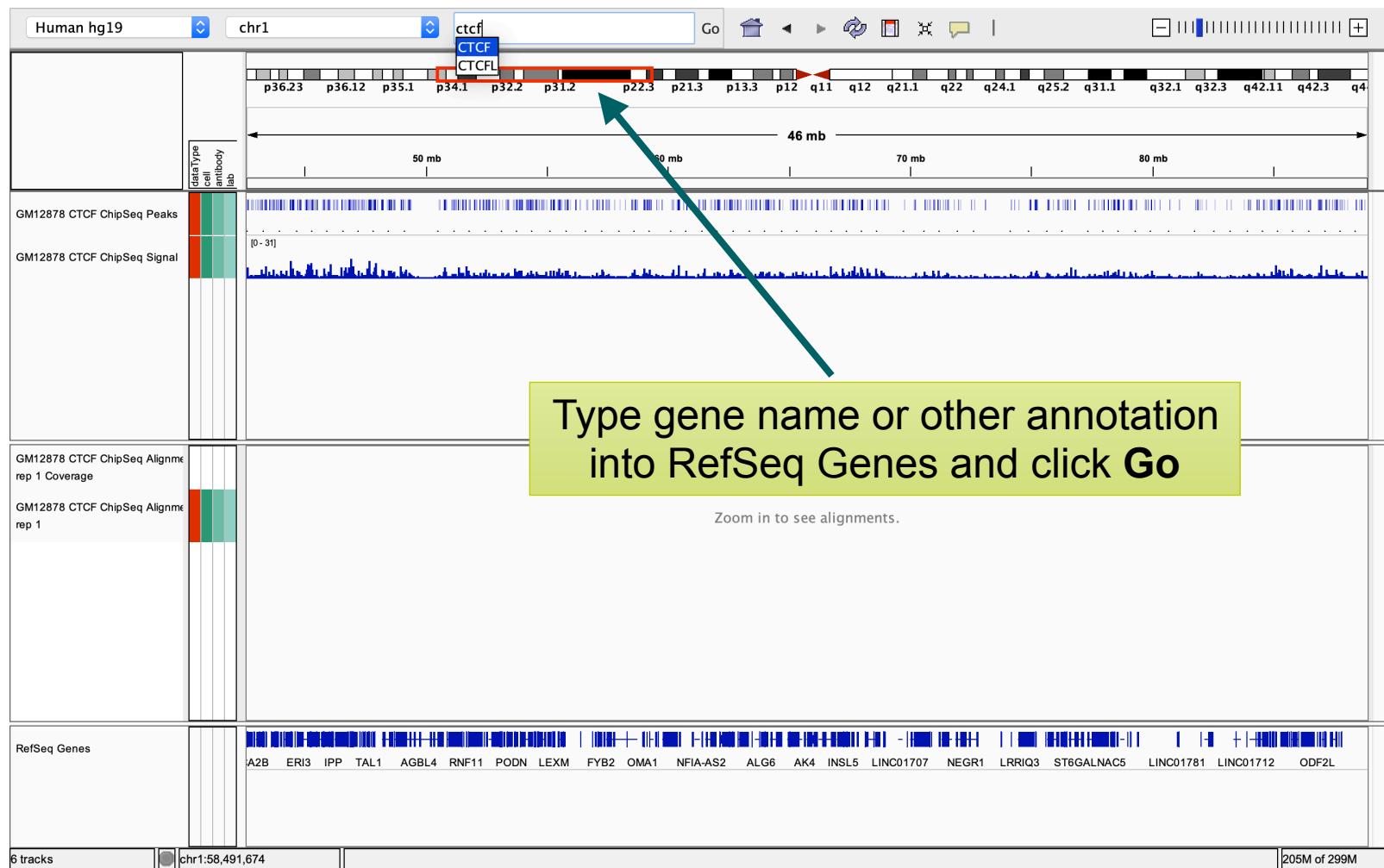
Navigation



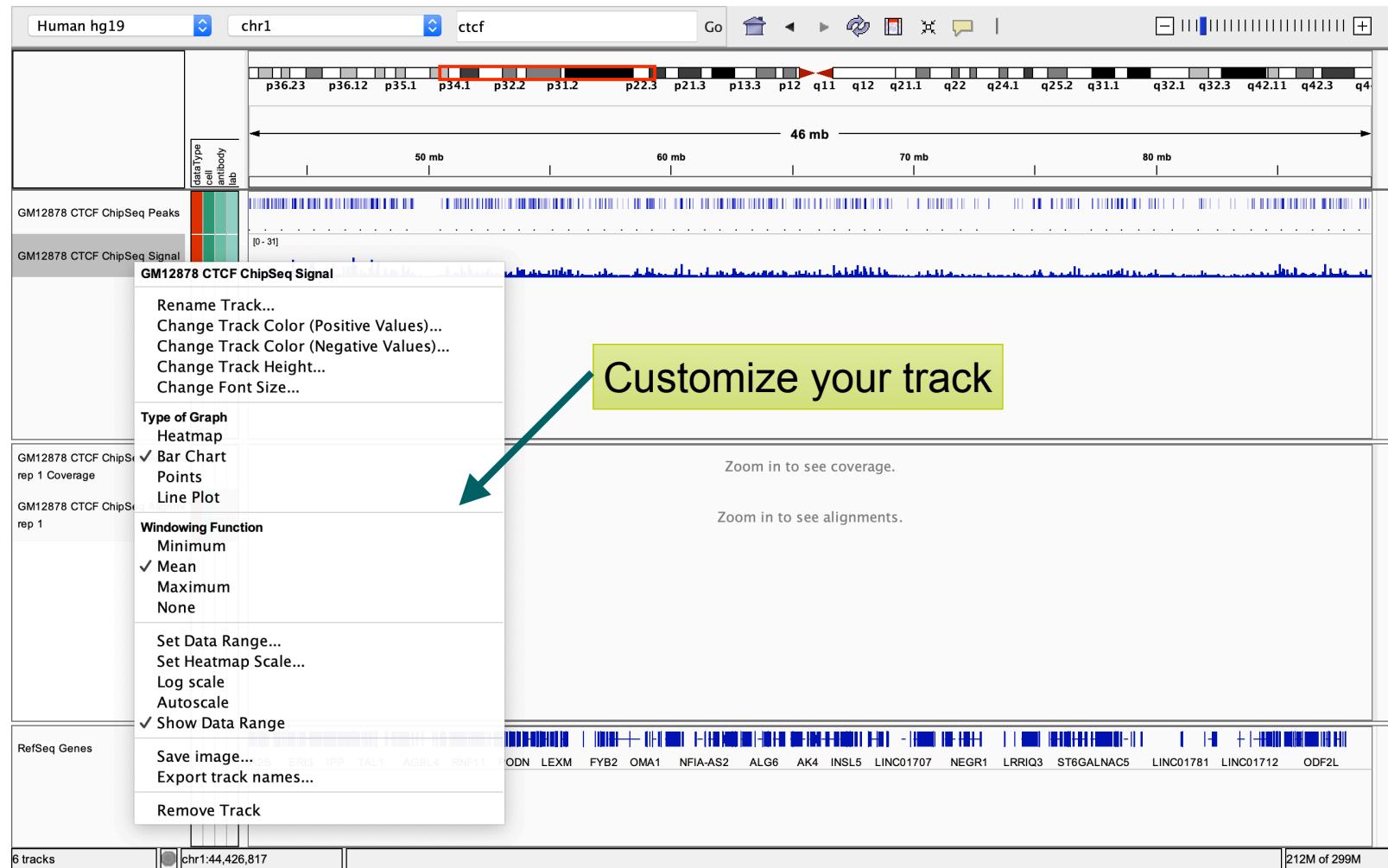
Navigation



Navigation

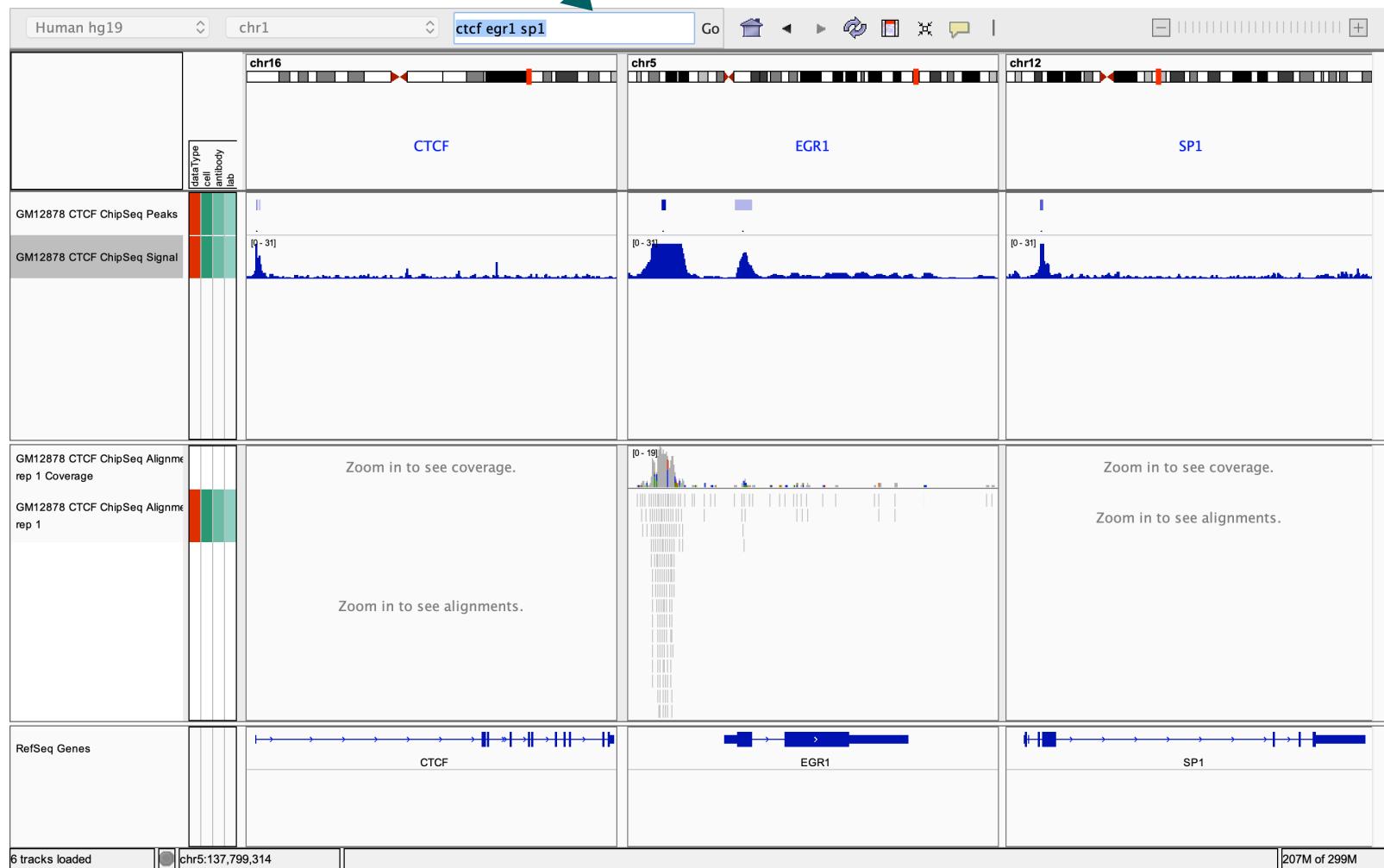


Navigation

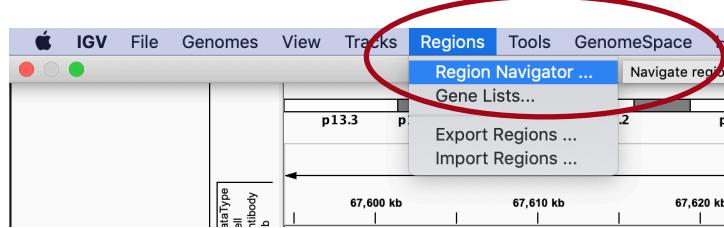


Viewing multiple regions

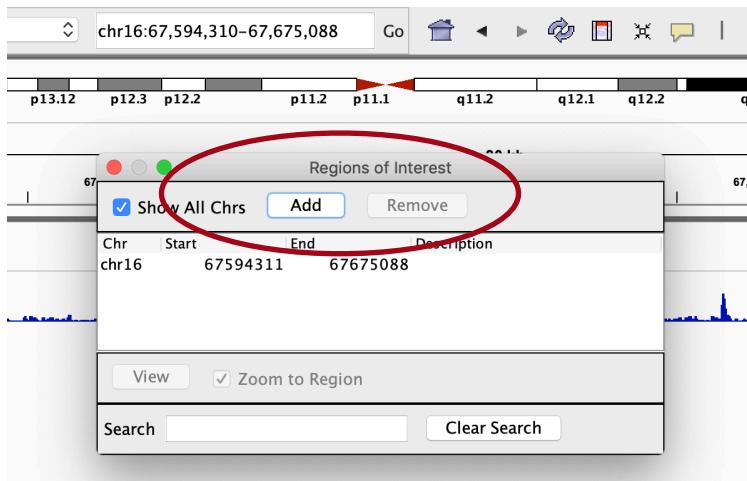
Input multiple regions



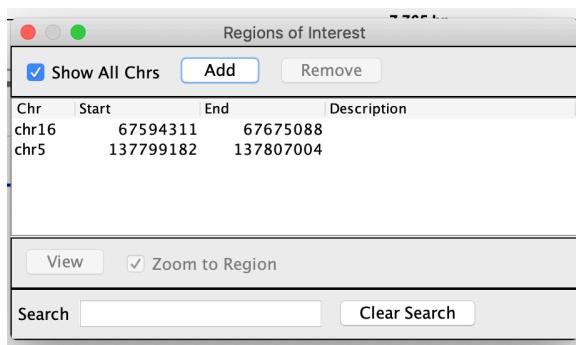
Viewing multiple regions



Click Regions > Region Navigator

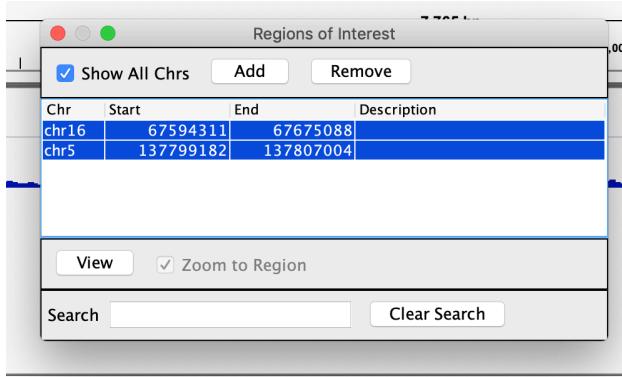


Click add



Repeat above step to add more regions

Viewing multiple regions



Select two or more regions
click **View**

Hands on

- Change the color of CTCF track
- Find more than two interesting regions
- View the multiple regions
- Load more dataset from ENCODE project
 - H3K4me1 of GM12878
 - H3K4me3 of GM12878
- Use different color for these three tracks

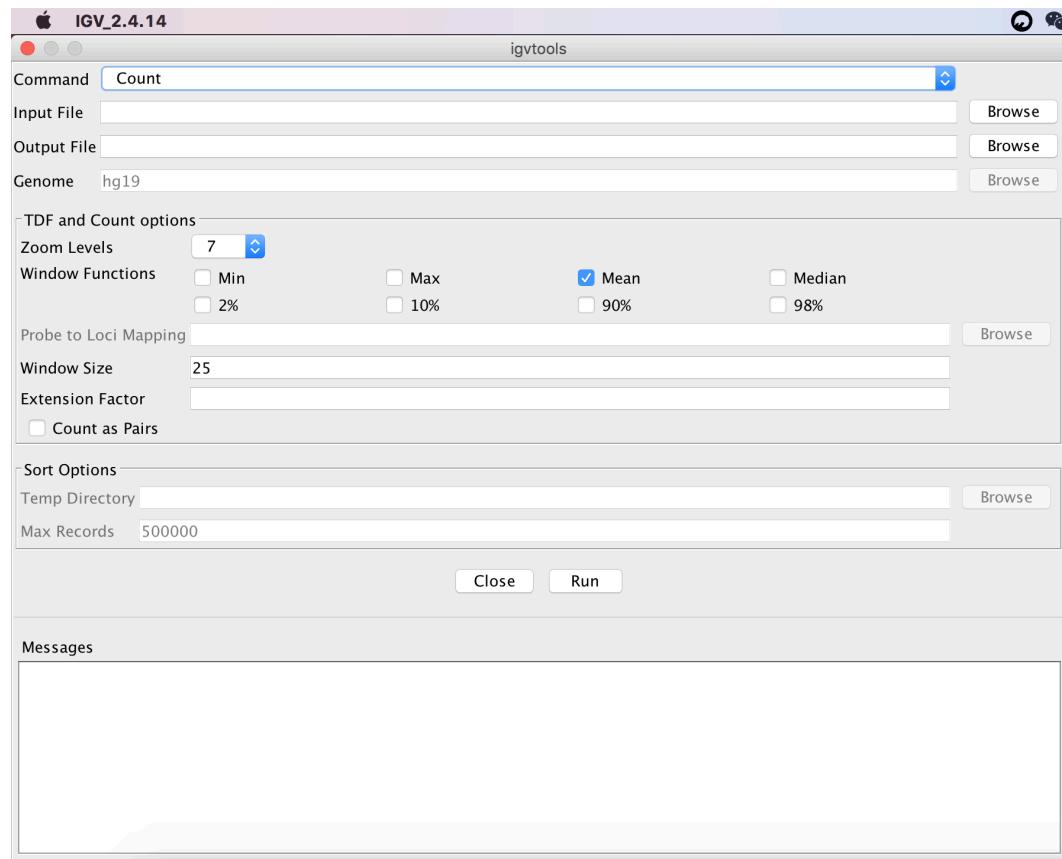
IGV Tools

A set of utilities for preparing files for efficient display

- toTDF
 - Converts sorted data file to binary file (TDF).
- counts
 - Computes average alignment or feature over a window size across the genome
- sort
 - Sorts file by genomic position
- index
 - Creates an index file for alignment or feature file

IGV Tools

Can be launched from
the IGV user interface
Tools > Run igvtools...



toTDF

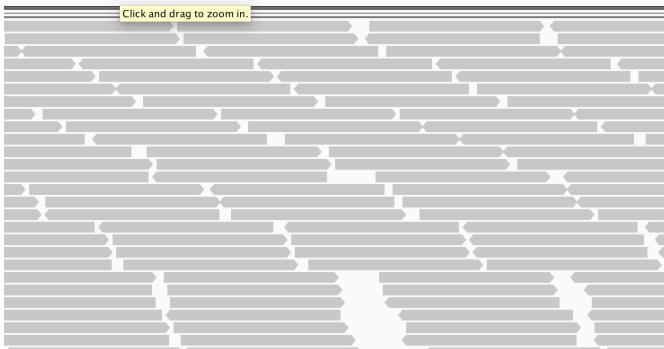
The **toTDF** utility converts large data files into tiled data format (.tdf) files

TDF files have the following advantages:

- Data is indexed for efficient retrieval
- Data is preprocessed for zoomed out views
- TDF files are web friendly - large data can be shared over the web.

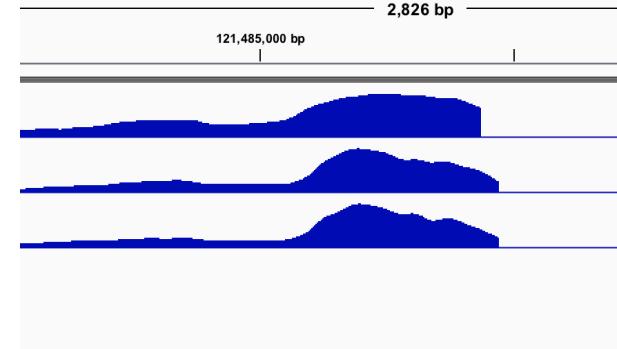
count

The **count command** is used to transform alignment files to read density TDF files, e.g. for ChIP-seq, RNA-seq and similar alignment counting experiments



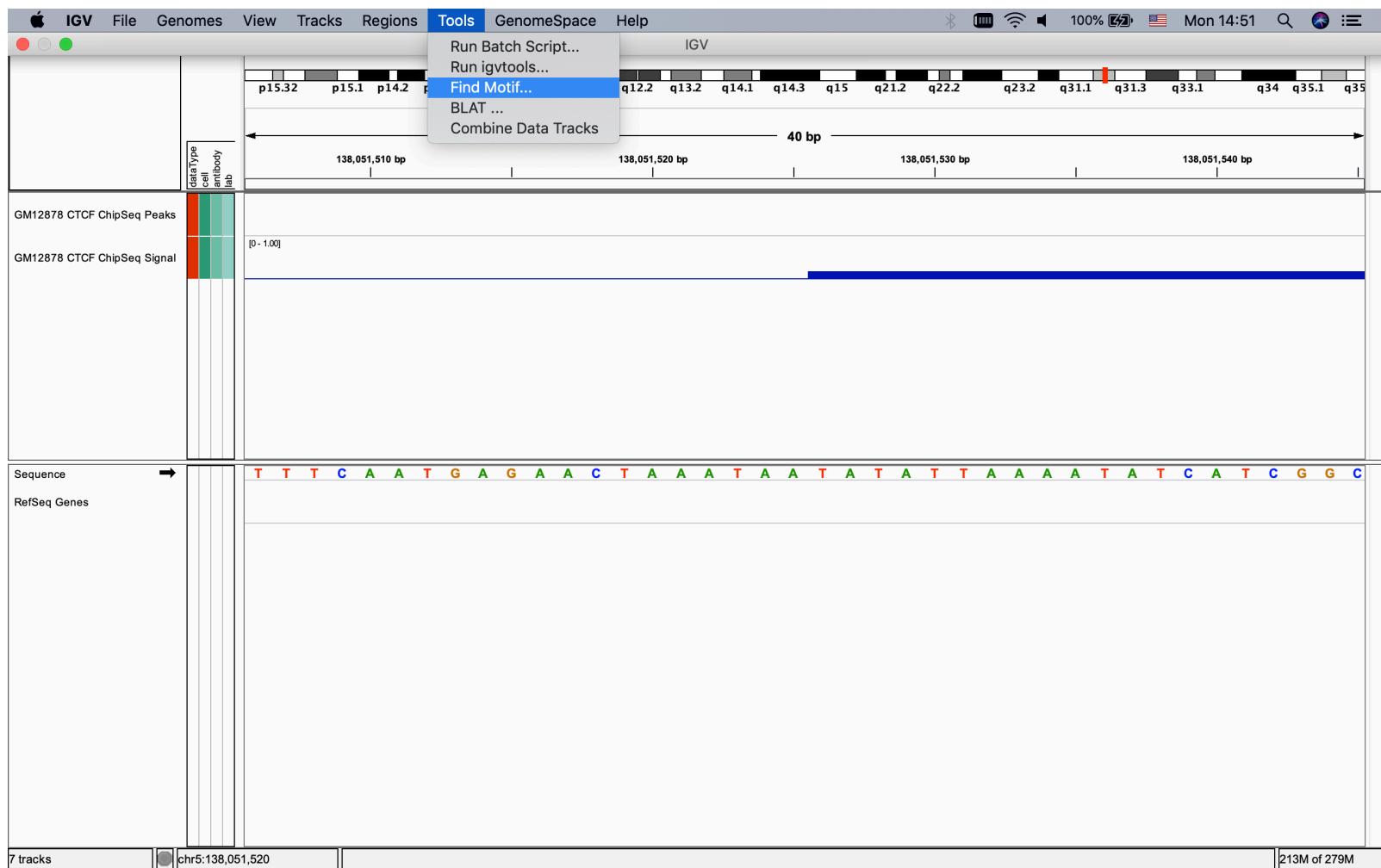
count

Alignment

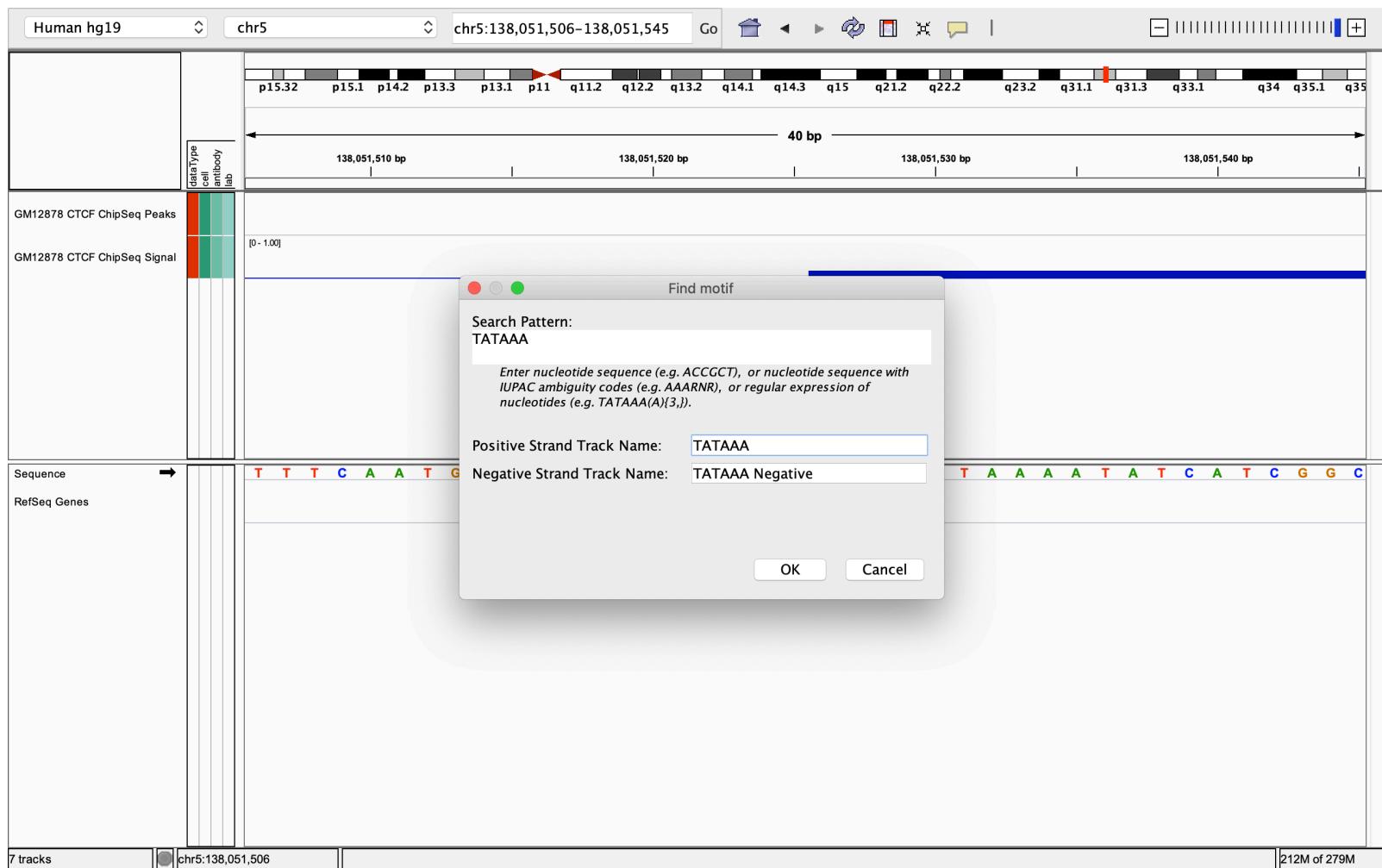


Read Density

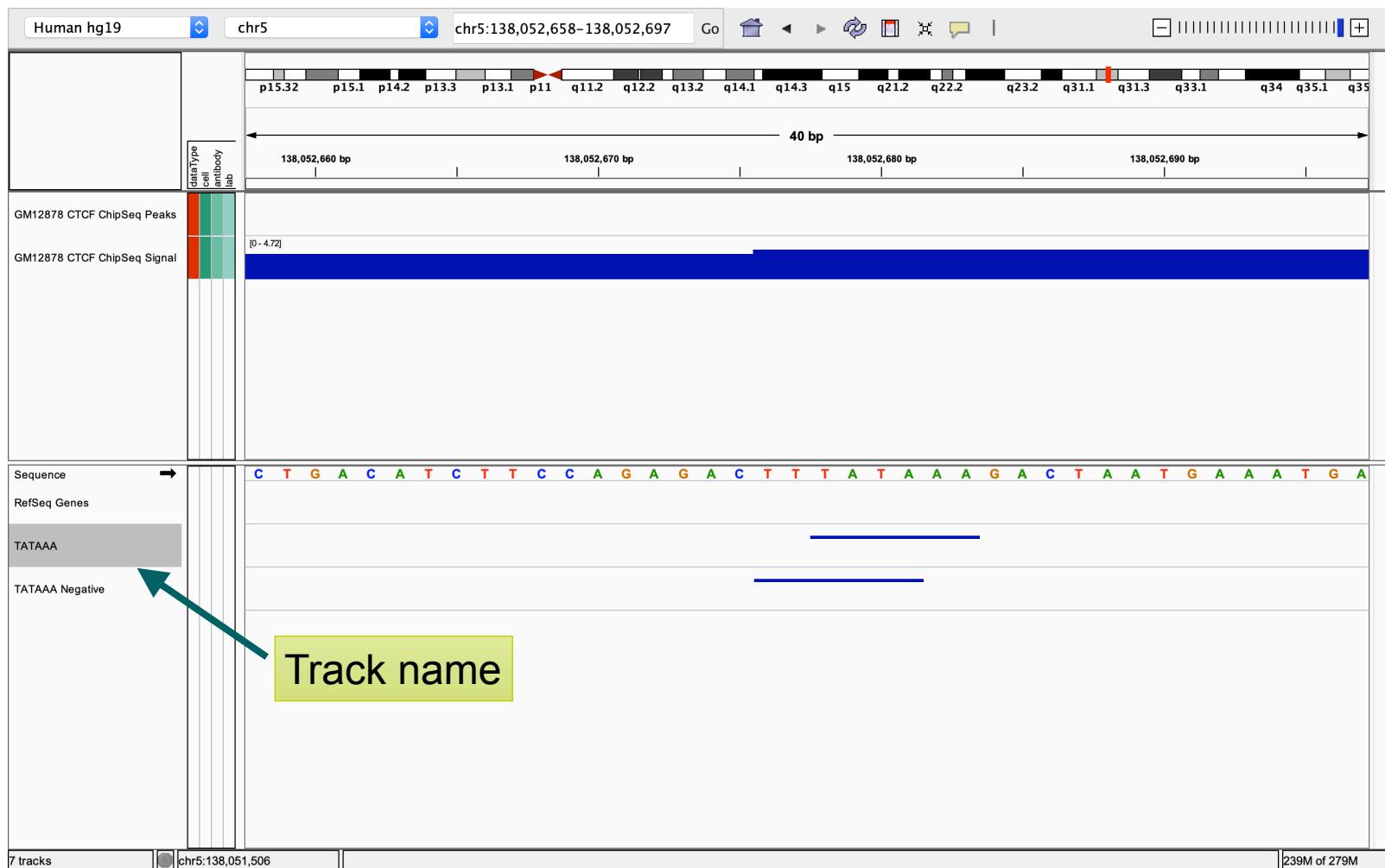
Find Motif



Find Motif

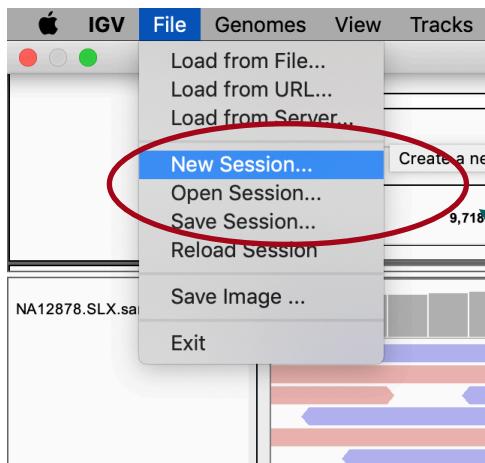


Find Motif



Viewing SNPs

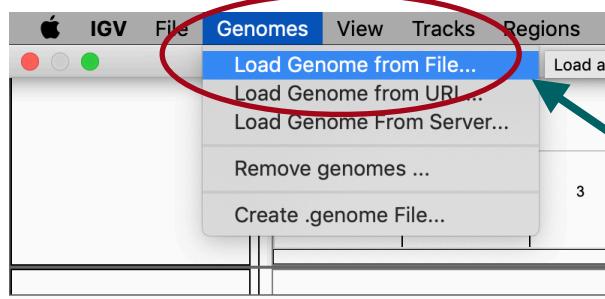
- Download data:
 - [https://costalab.ukaachen.de/open_data/
Bioinformatics Analysis in R 2019/BIAR_D5/
igvData.tar.gz](https://costalab.ukaachen.de/open_data/Bioinformatics_Analysis_in_R_2019/BIAR_D5/igvData.tar.gz)
- Create a new session



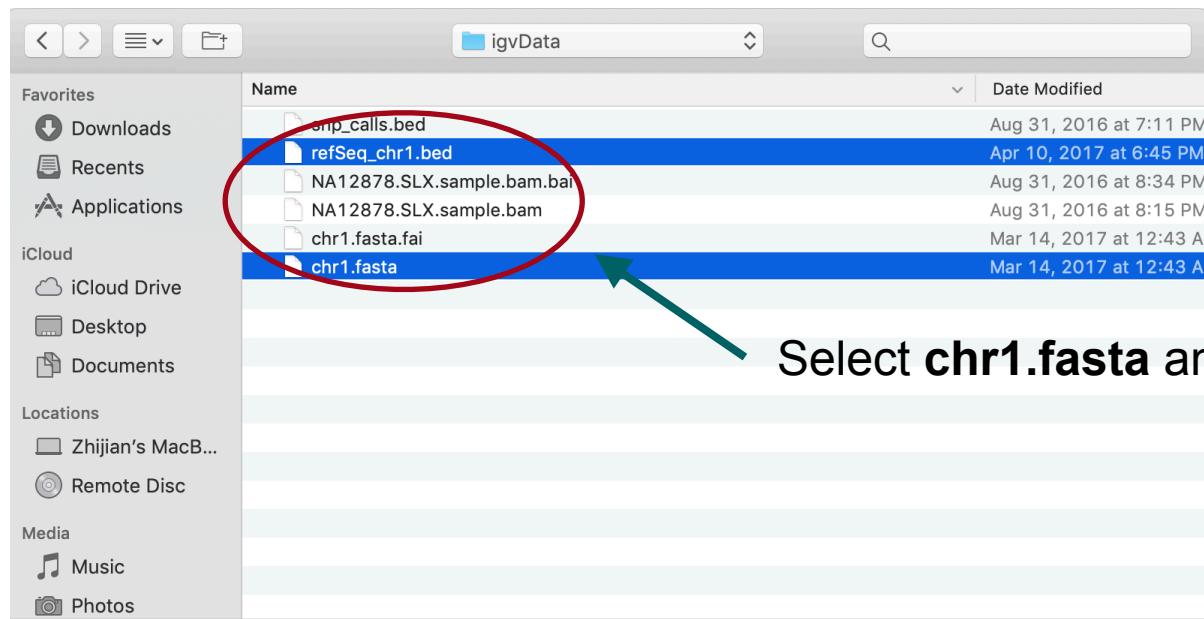
Click **New Session**

Viewing SNPs

- Load reference genome:



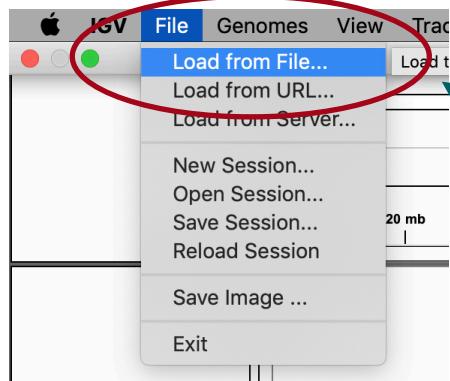
Click **Genomes > Load Genome from File**



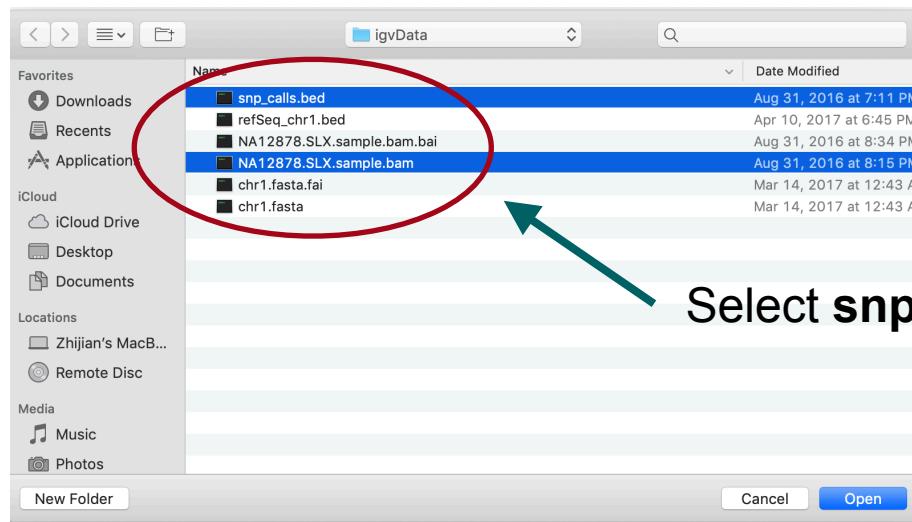
Select **chr1.fasta** and **refSeq_chr1.bed**

Viewing SNPs

- Load sequencing data:

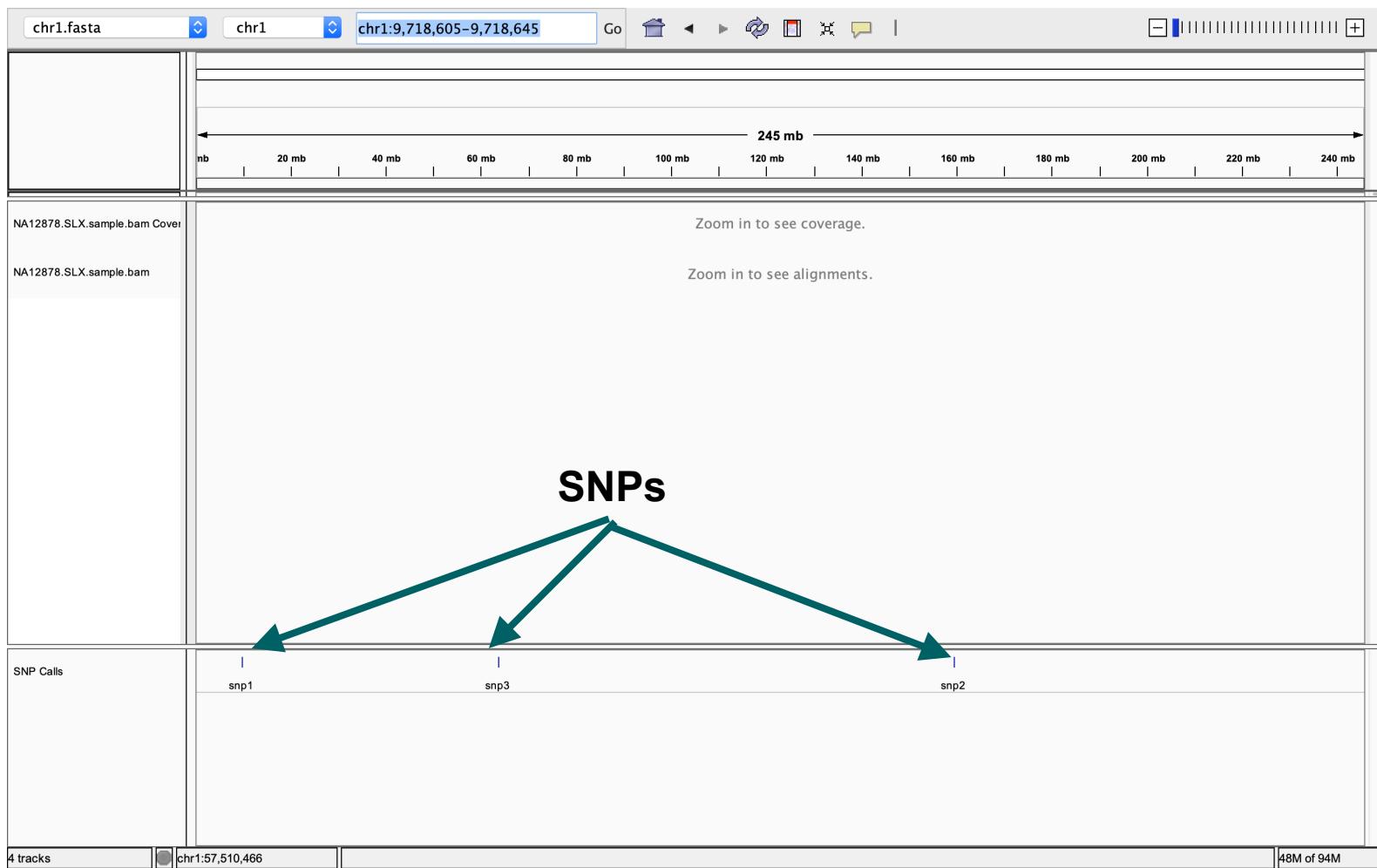


Click File > Load Genome from File

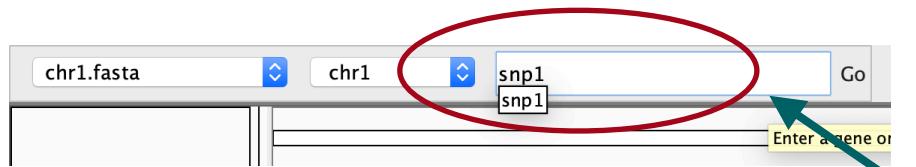


Select.snp_call,bed and NA12878.SLX.sample.bam

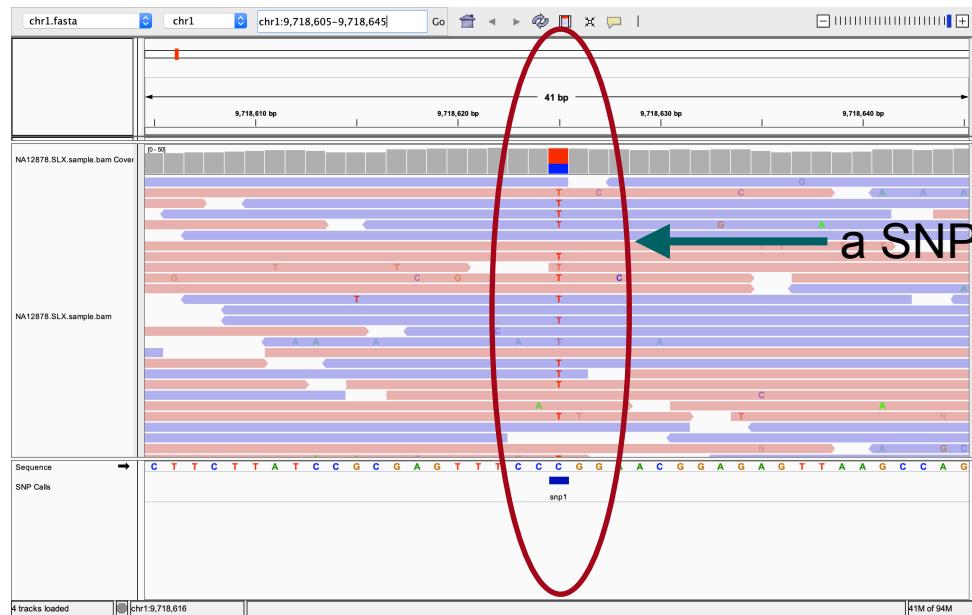
Viewing SNPs



Viewing SNPs

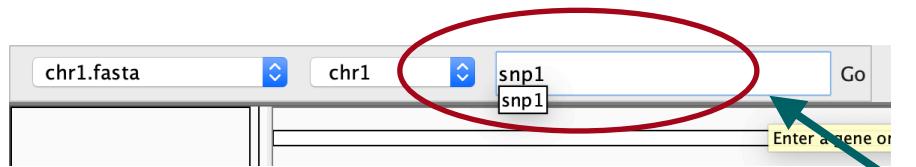


Input SNP1 and Click Go

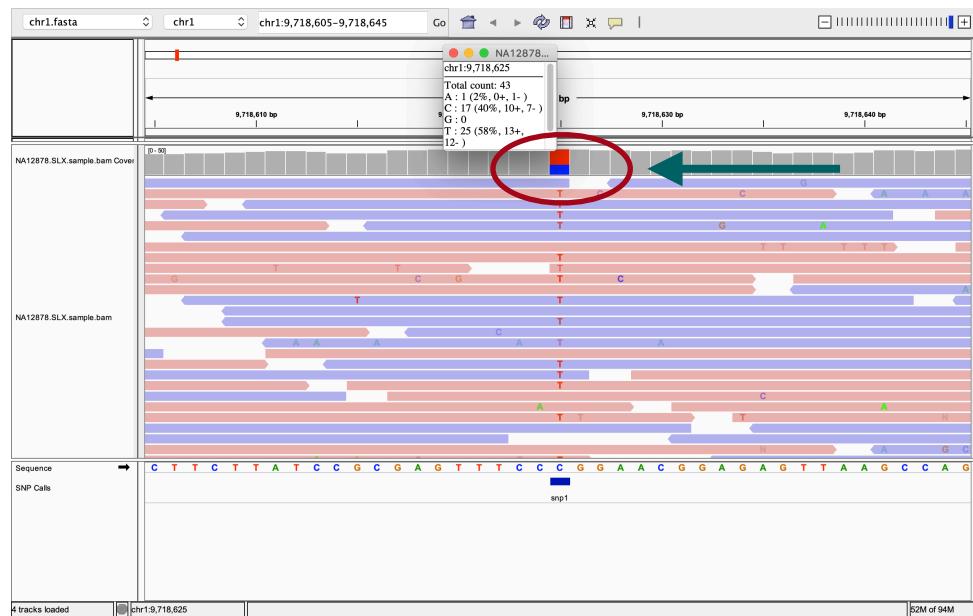


a SNP with C > T

Viewing SNPs



Input.snp1 and Click Go



Click here to see statistics

Viewing RNA-seq data

- Create a new session

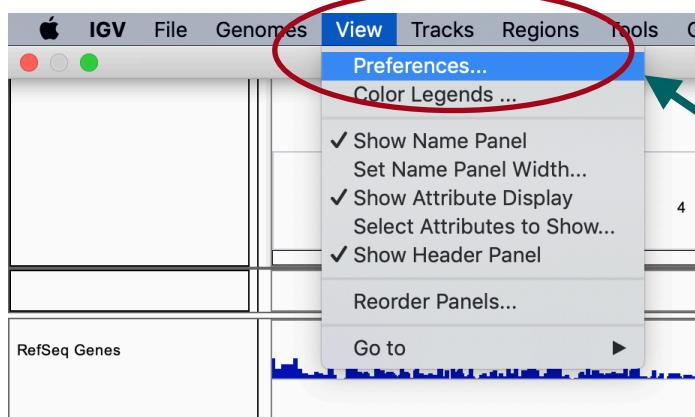


- Choose hg19 as reference genome

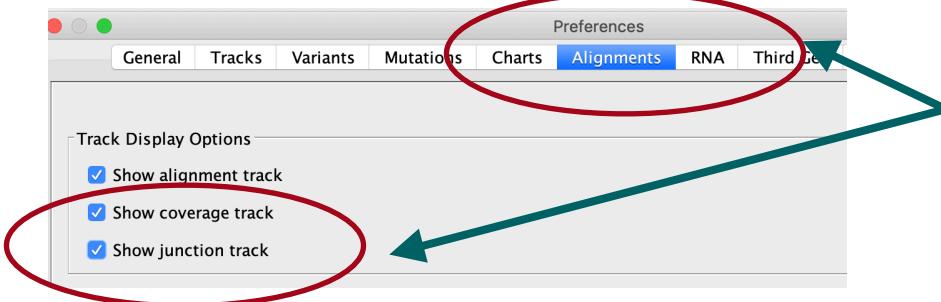


Viewing RNA-seq data

- Set preference for viewing RNA-seq data



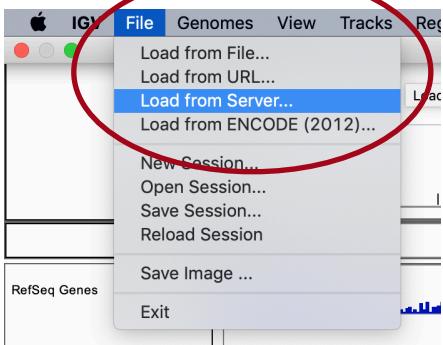
Click **View > Preference**



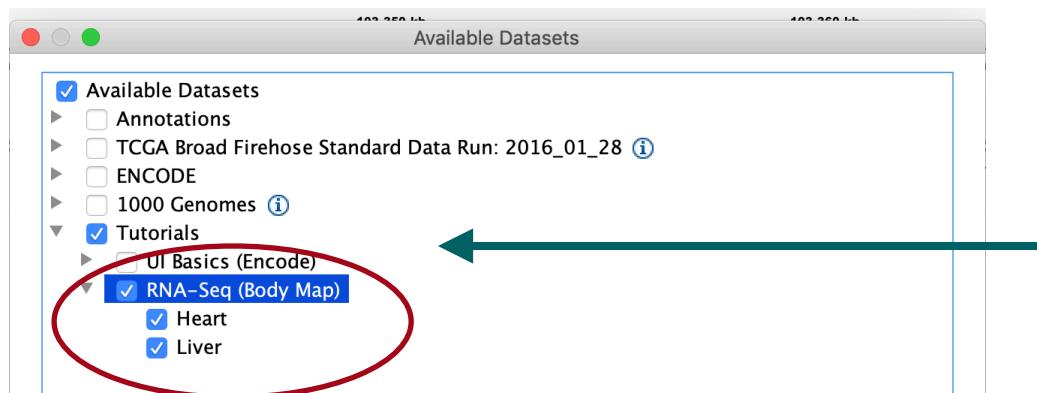
Select **Alignments tab**
Check **junction track**

Viewing RNA-seq data

- Load data



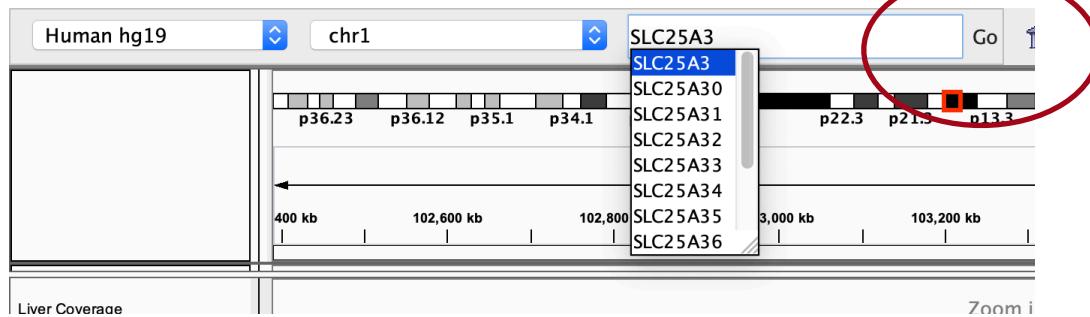
Click File > Load from Server



Open Tutorial
Click on RNA-seq > OK

Viewing RNA-seq data

- Jump to your favorite gene

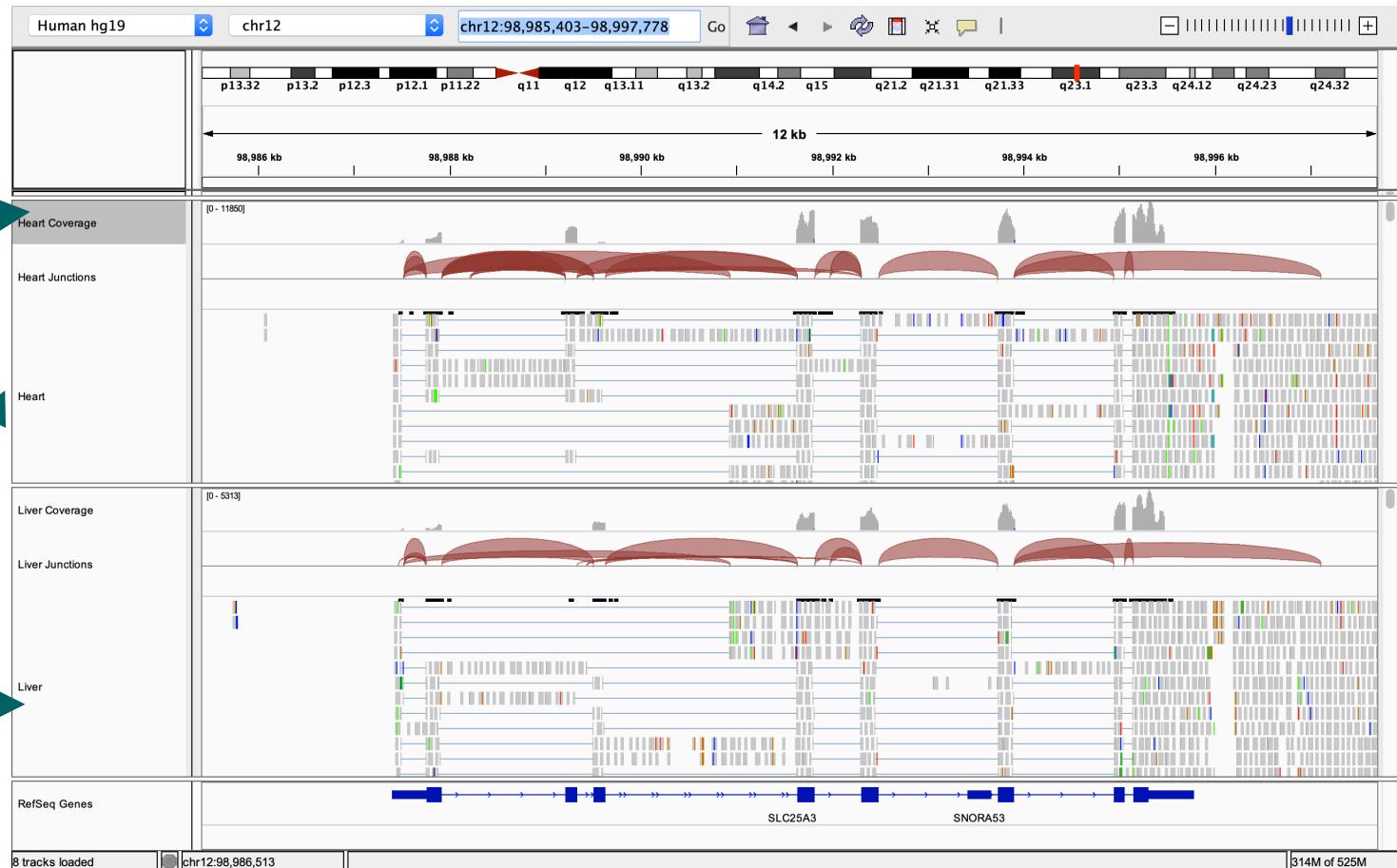


Viewing RNA-seq data

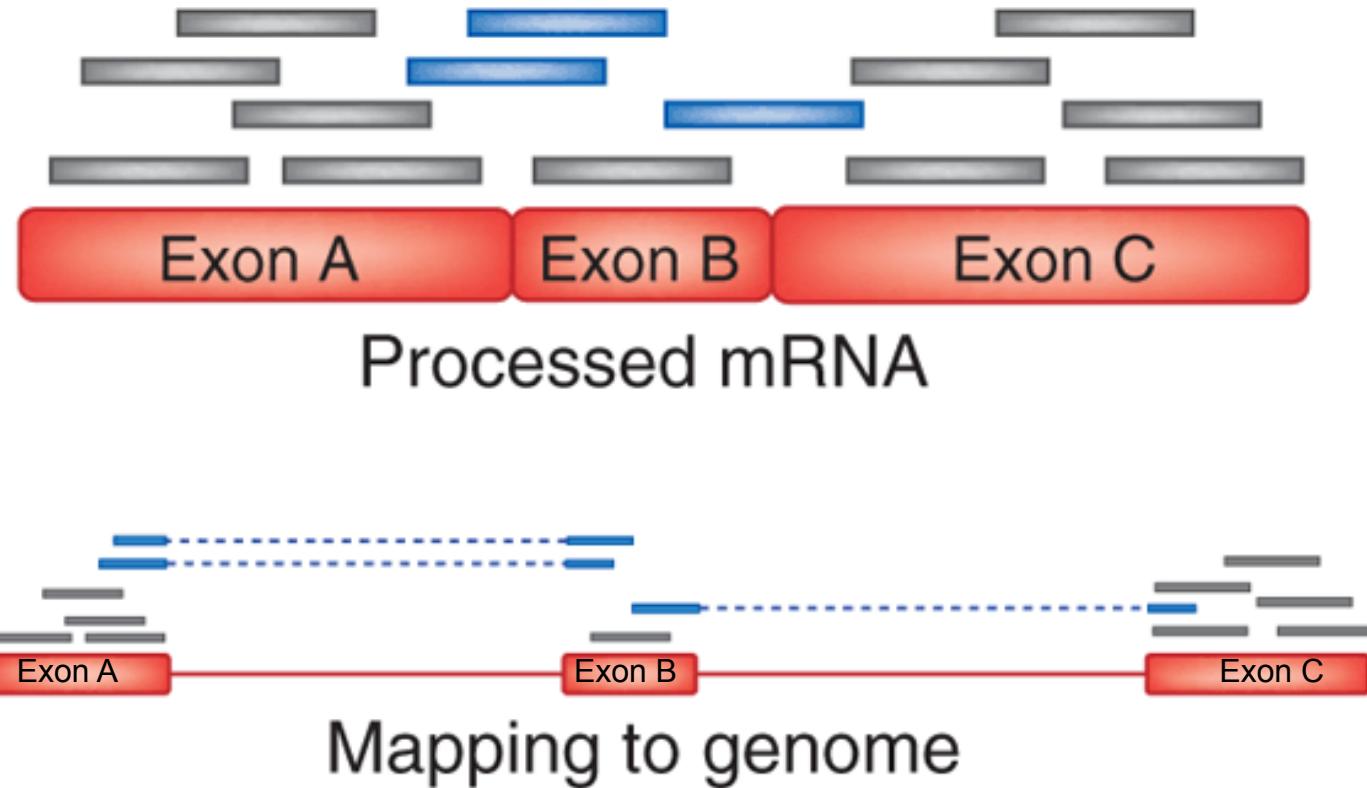
Coverage

Junction

Alignment



Split Read Mapping (RNA-Seq)



- Reads spanning distinct exons indicate junctions

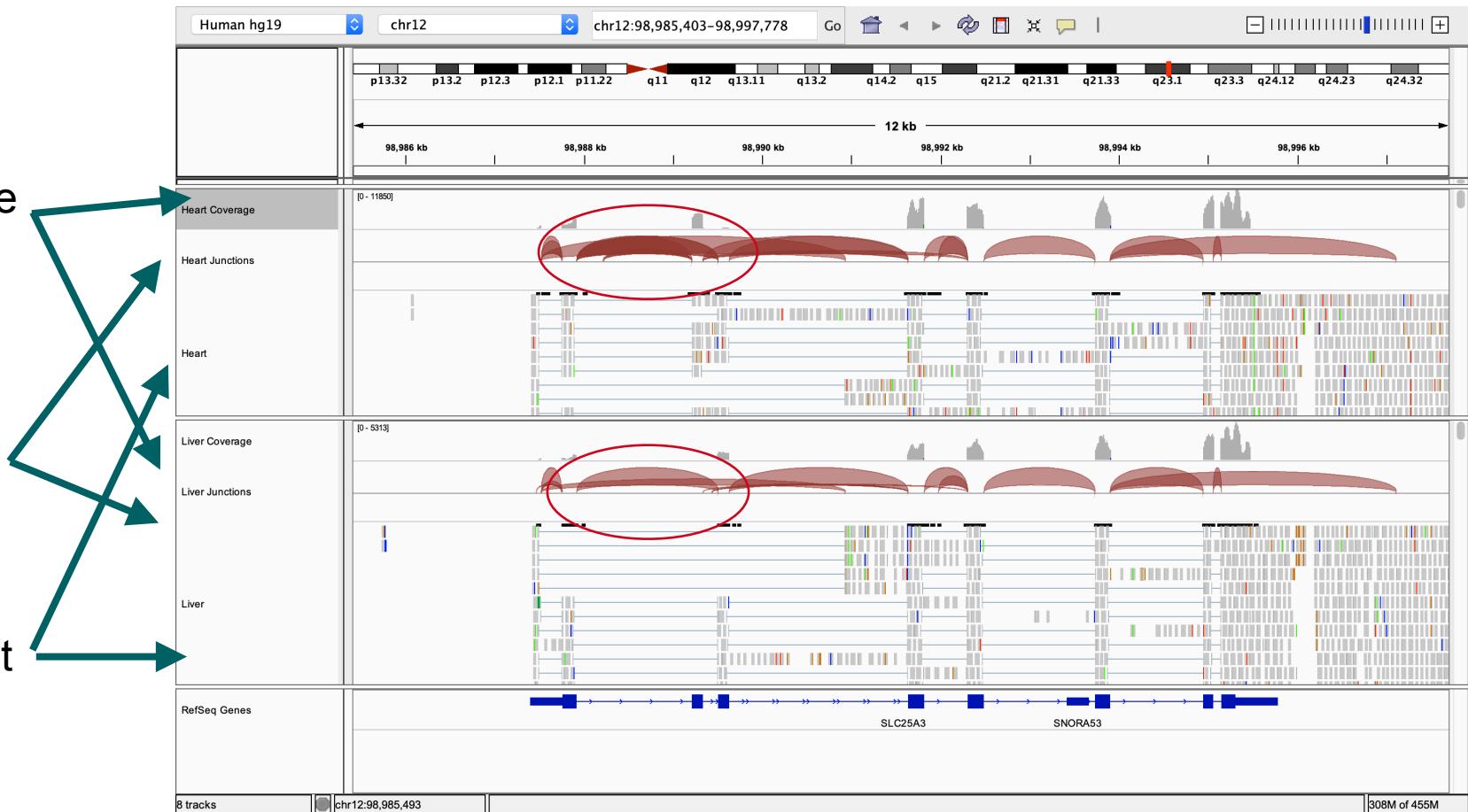
Viewing RNA-seq data

- Example of alternative splicing

Coverage

Junction

Alignment



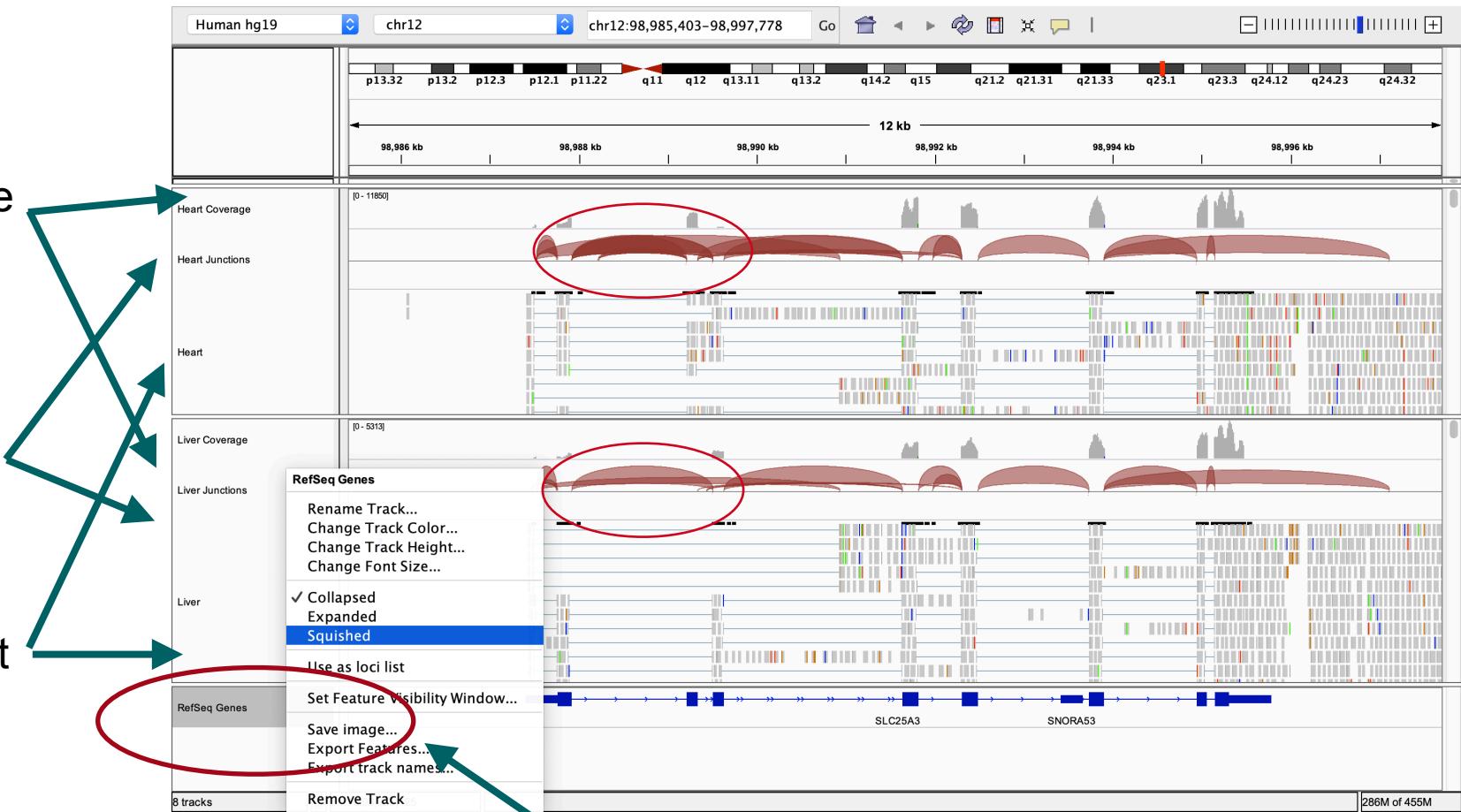
Viewing RNA-seq data

- Example of alternative splicing

Coverage

Junction

Alignment



Right click and select Squished

Viewing RNA-seq data

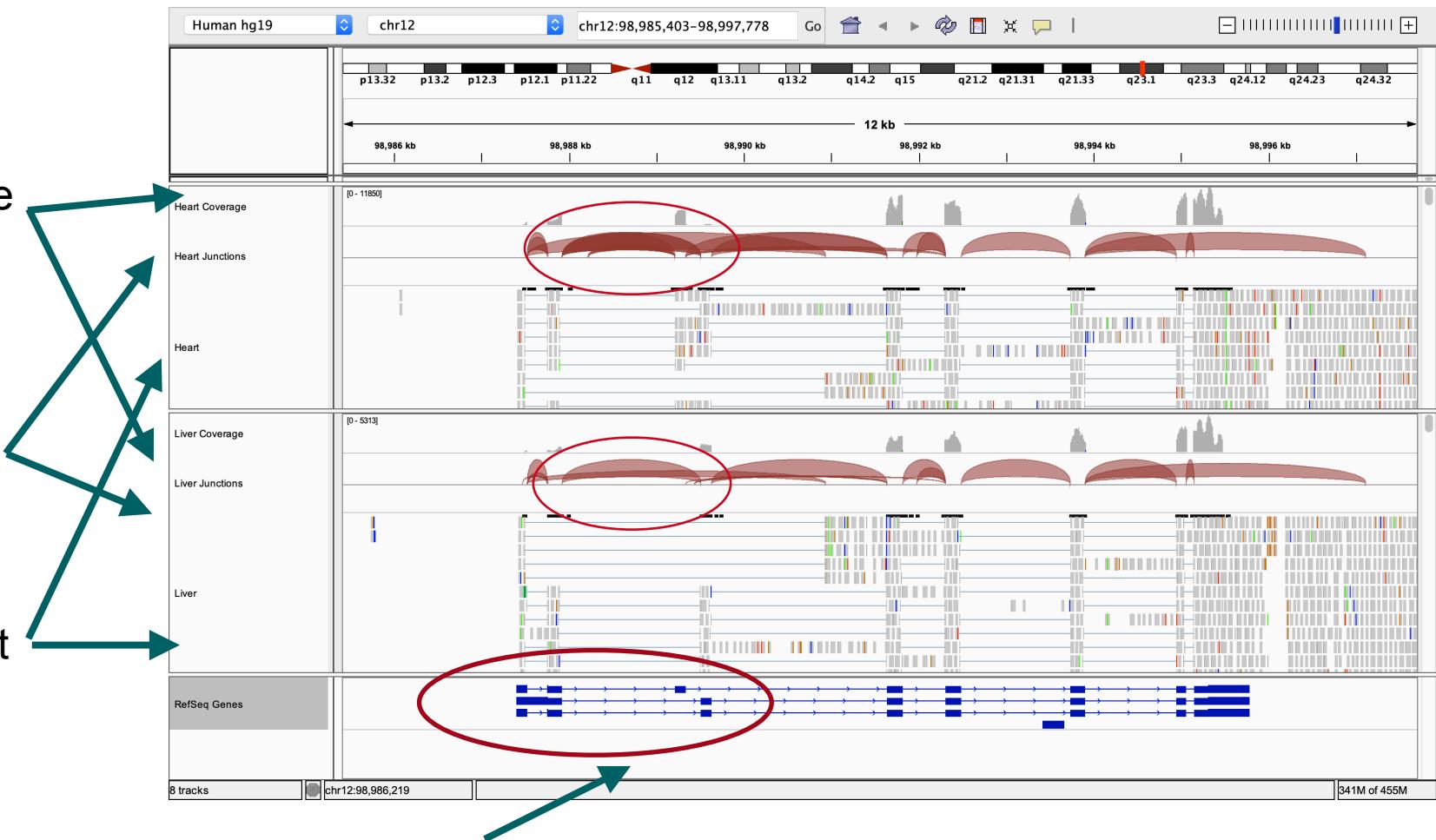
- Example of alternative splicing

Coverage

Junction

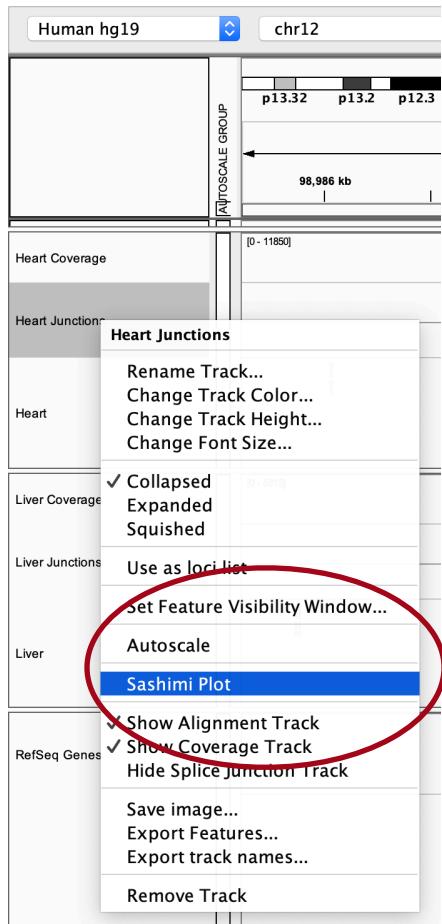
Alignment

Alternative splicing



Viewing RNA-seq data

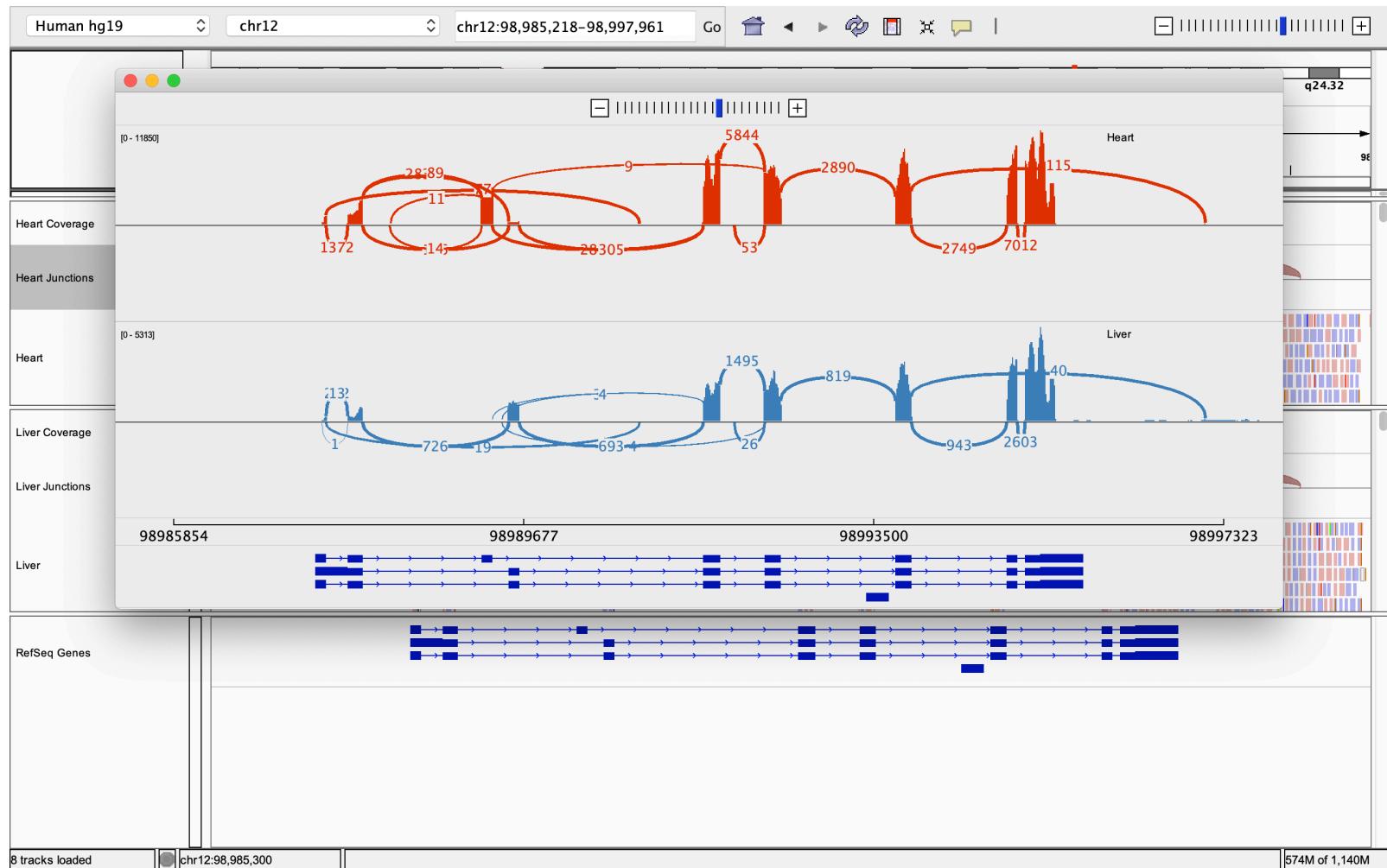
- Example of alternative splicing



Right click Junction > Sashimi plot

Viewing RNA-seq data

- Example of alternative splicing





www.costalab.org

Institute for
Computational Genomics
01011011010
10100100101



RWTHAACHEN
UNIVERSITY