

# 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](http://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/>

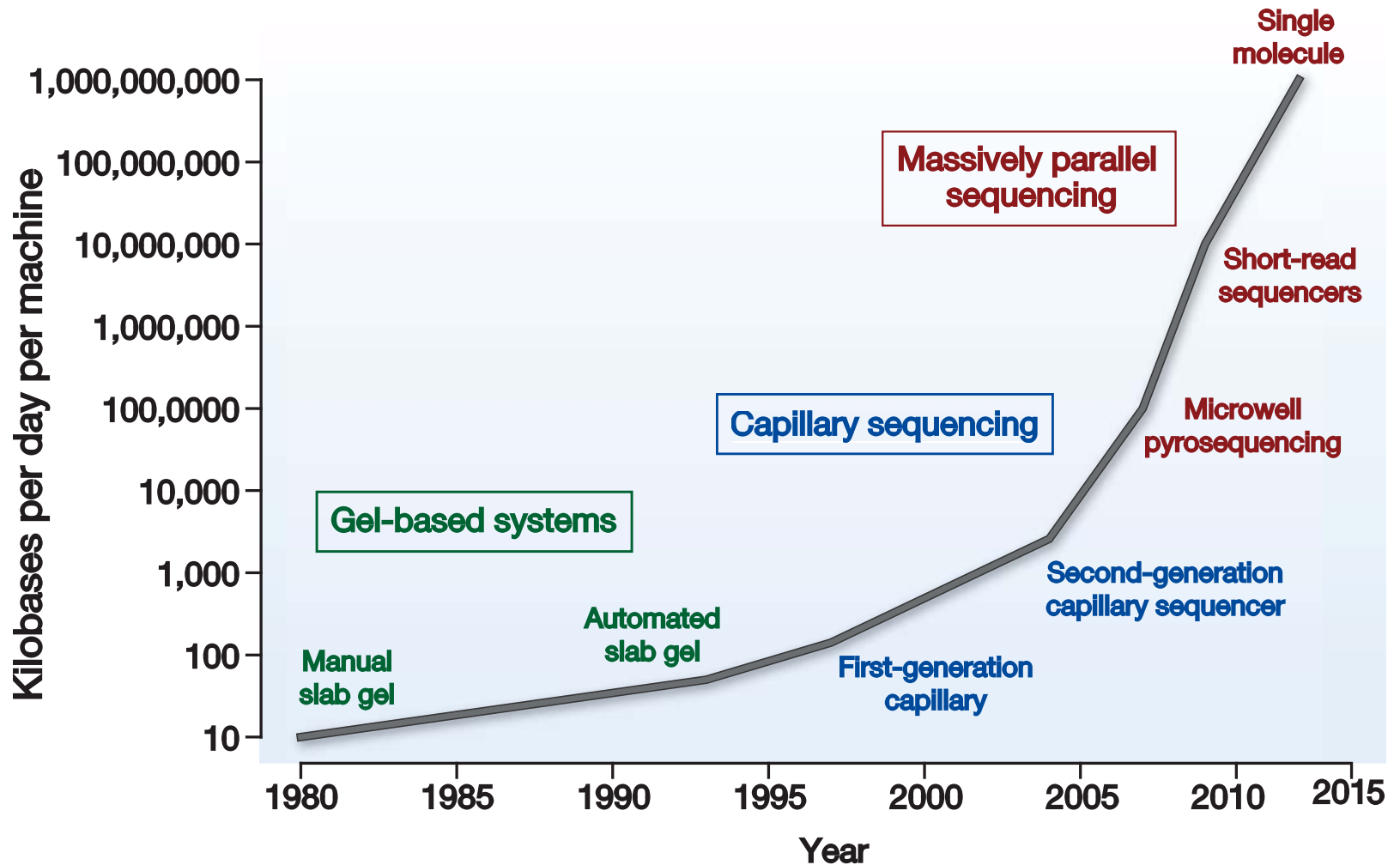
# Next Generation Sequencing

---

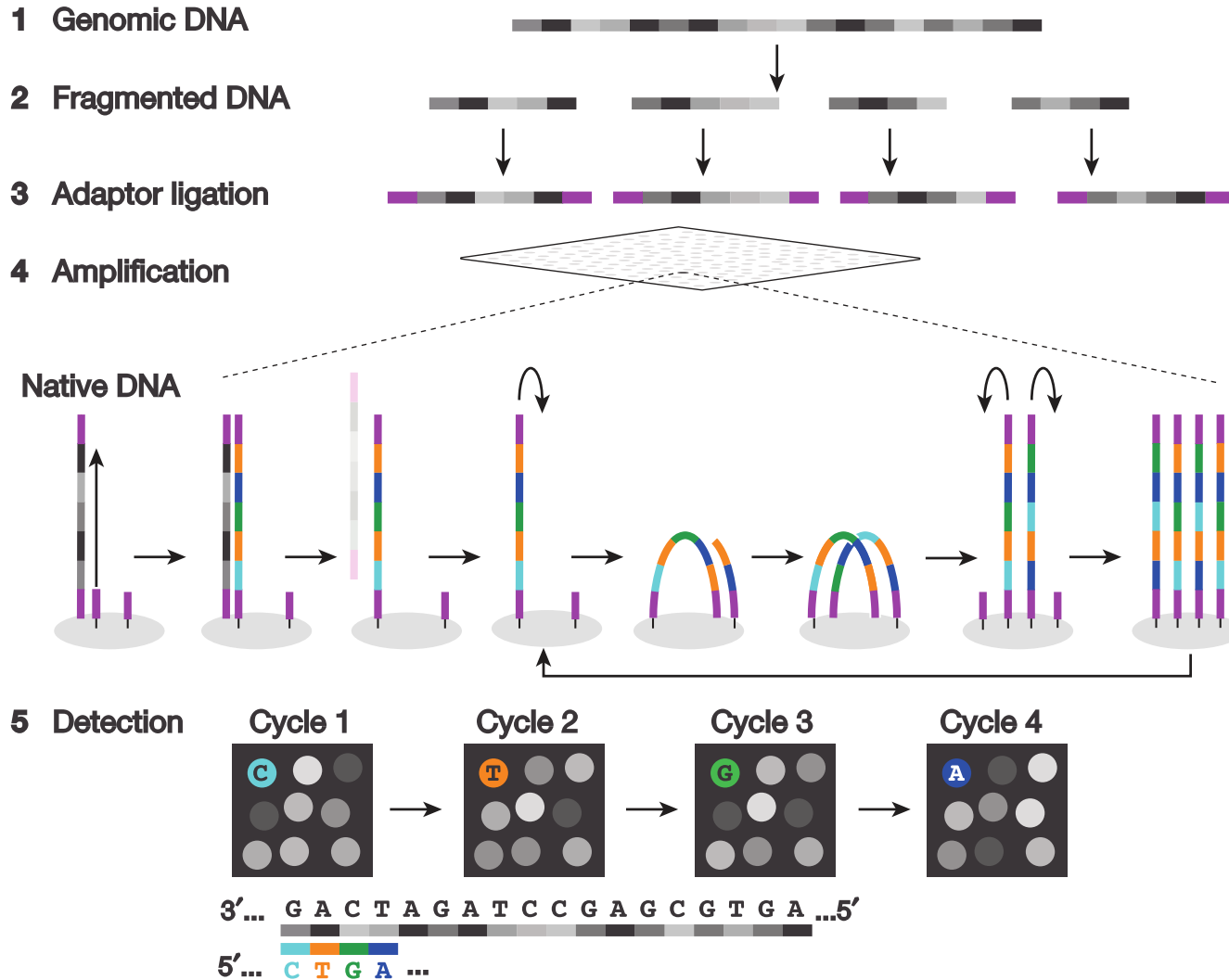
- ▶ NGS take advantage of **parallelization**
  - ▶ reads millions/billions of reads for a time
  - ▶ shorter reads (50-100 bps)
  - ▶ higher error rates (0.1-1%)
- ▶ commercial products:
  - ▶ 454
  - ▶ SOLiD
  - ▶ **Solexa (Illumina)**



# Next Generation Sequencing



# Next Generation Sequencing



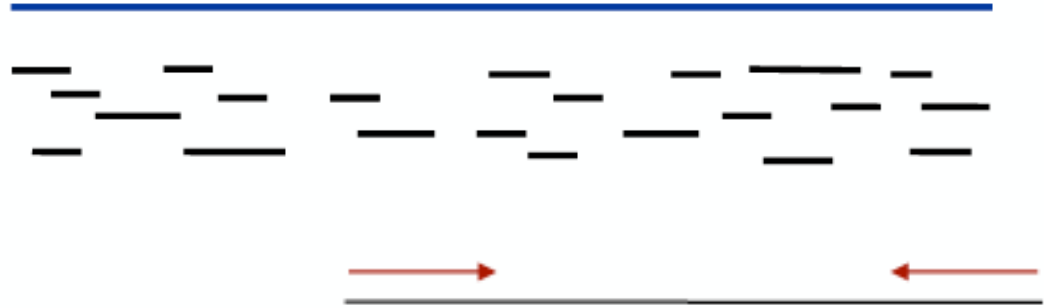


# Read Types

Fragment DNA:



Single end



Paired end  
Ins: 200-800 bp

# Read Types

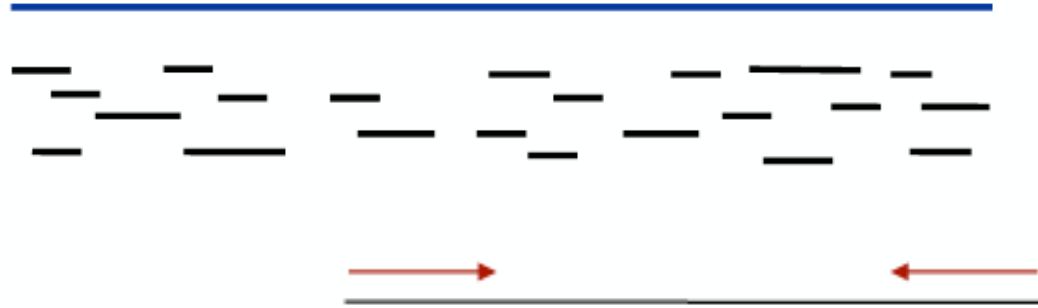
Fragment DNA:



Single end

Advantages:

- cheaper
- compatible with protocols producing small fragments (Ribo-seq, mirRNA-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  
MSLSLWQQCLARLQDELPATEFSMWIRPLQAELSDNTLALYAPNRFVLDW  
VRDKYLEALRDLLALQEKLVTIDNIQKTVAEYYKIKVADLLSKRRRSRVARP  
RQMAMALAKELLHAVGNGIMARKPNAKVVYMHSERFVQDMVKALQNNAI  
EEFKRYYSVDALLIDDFSLPEIGDAFGGRDHTTVLHACRKIEQLREESH  
KEDFSNLIRTLSS

# FASTA files

---

Start symbol

Sequence ID  
(no spaces)

Sequence description  
(spaces allowed)

```
>dnaA chromosomal replication initiator protein DnaA
MSLSLWQQCLARLQDELPATEFSMWIRPLQAELSDNTLALYAPNRFVLDW
VRDKYLEALRDLLALQEKLVTIDNIQKTVAEYYKIKVADLLSKRRRSVARP
RQMAMALAKELLHAVGNGIMARKPNAKVVYMHSERFVQDMVKALQNNAI
EEFKRYYSVDALLIDDFSLPEIGDAFGGRDHTTVLHACRKIEQLREESH
KEDFSNLIRTLSS
```

The sequence

# FASTQ files

---

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

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

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

This number ( $Q$ ) can be  
converted to  $P$

Phred-scale

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

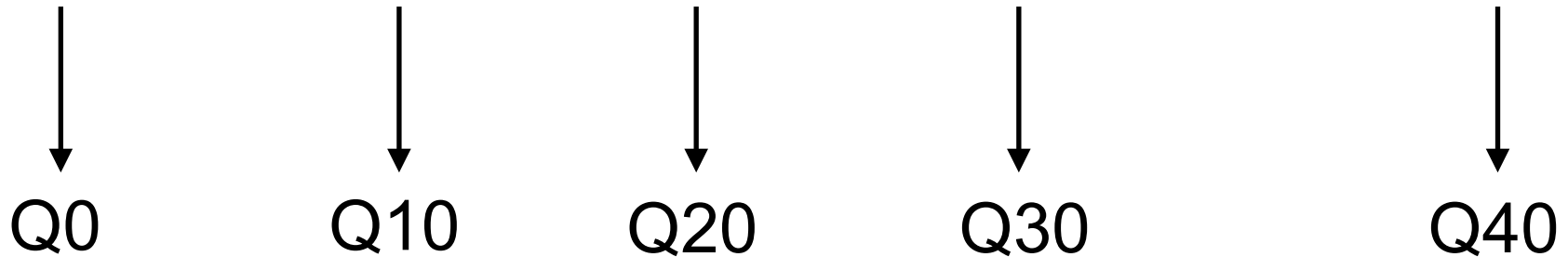
$$P = 10^{(-Q/10)}$$

# FASTQ files

---

Uses letters/symbols to represent numbers:

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



*bad*

*maybe*

*ok*

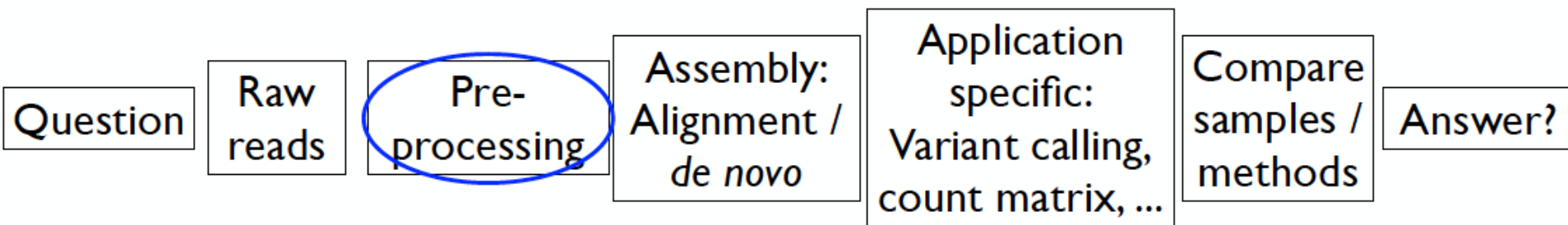
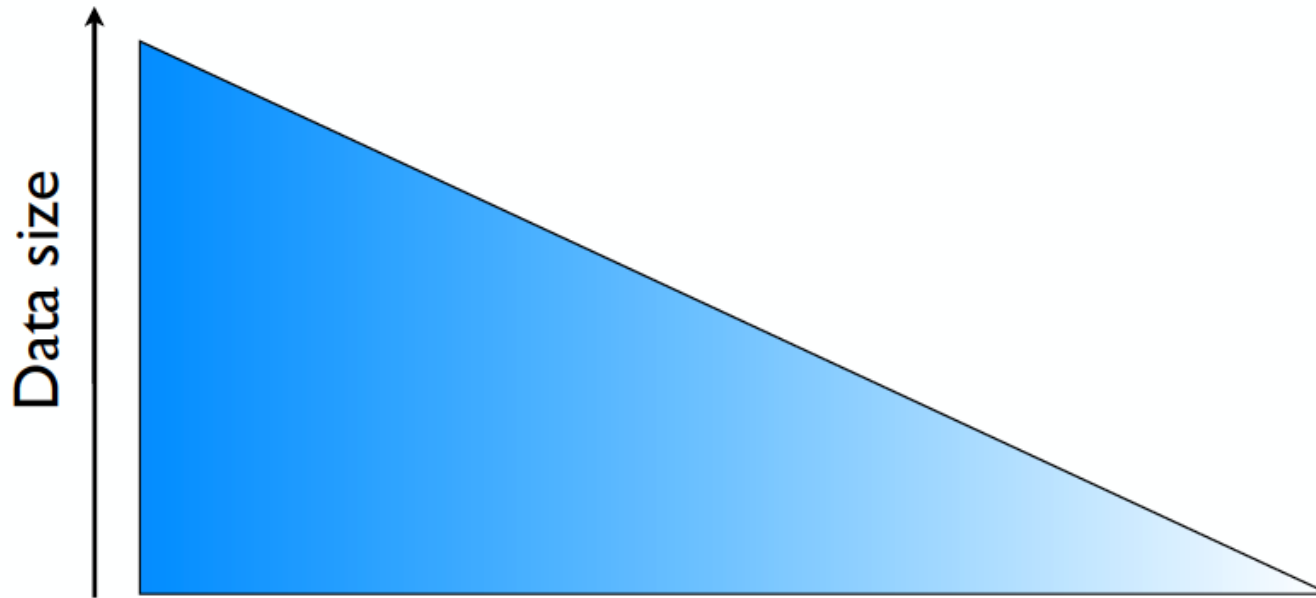
*good*

*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), ...**

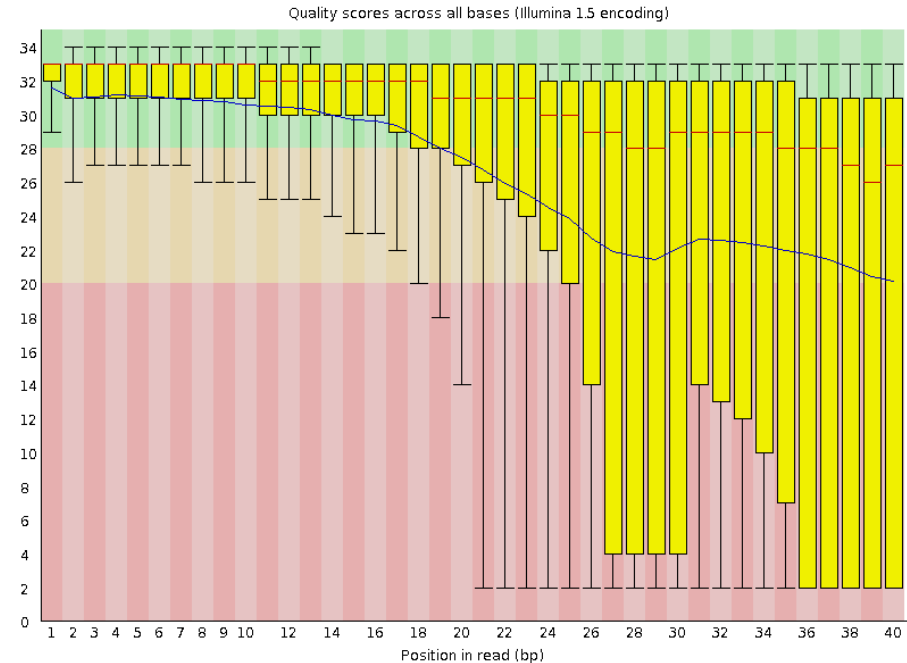
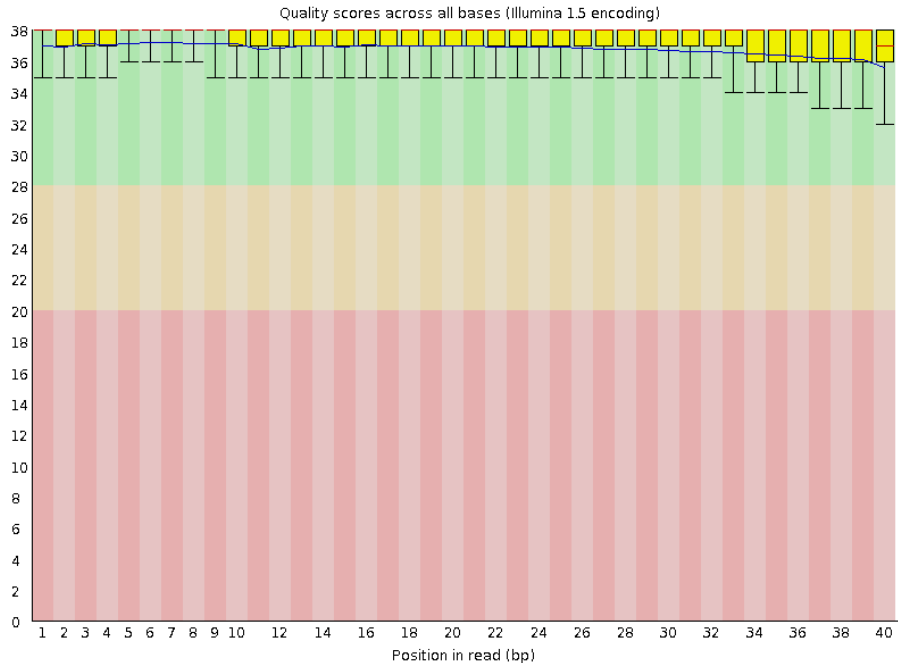
# Quality Control

---

- **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](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](http://www.bioinformatics.babraham.ac.uk/projects/fastqc/bad_sequence_fastqc.html)

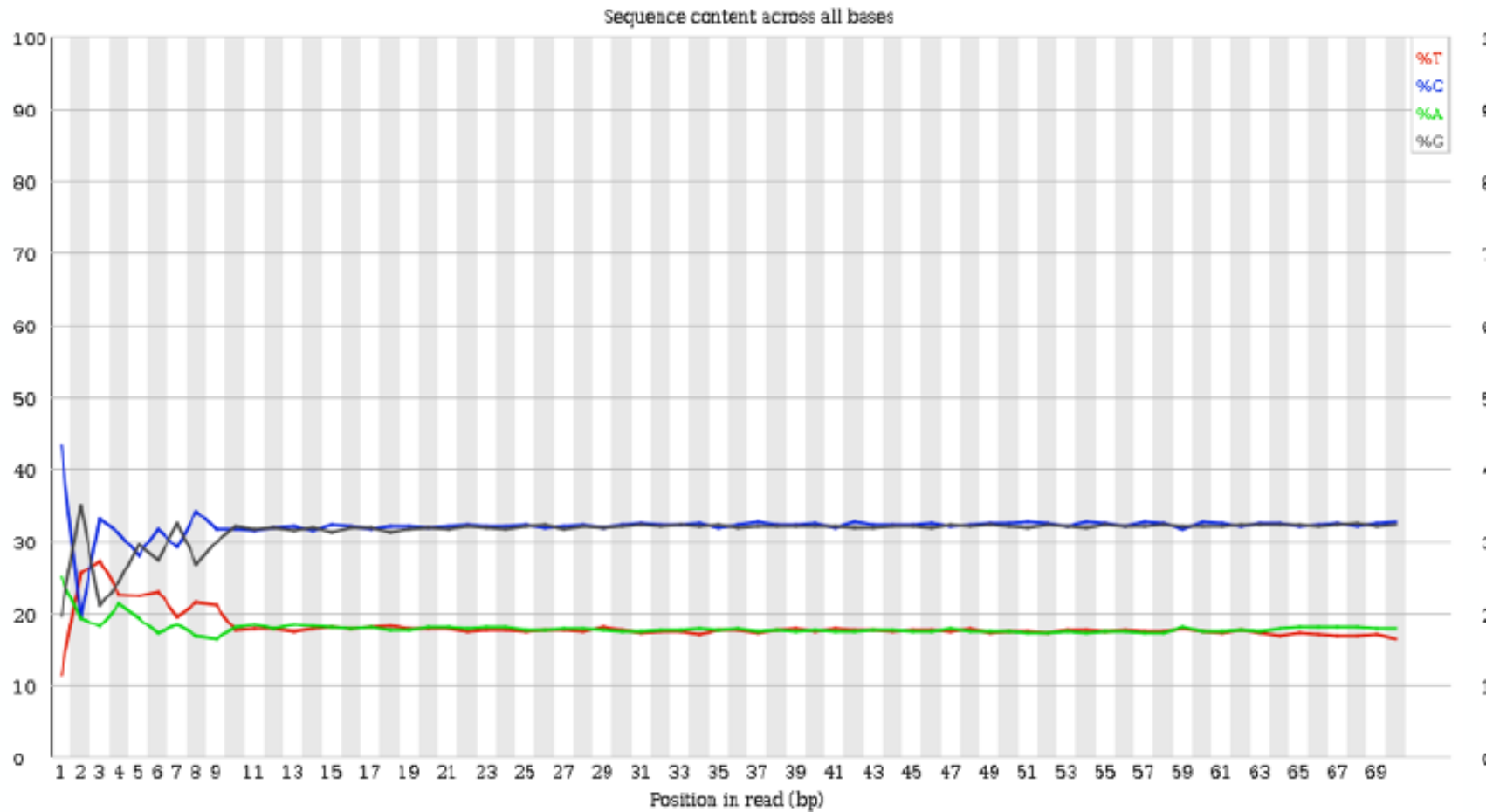
# Quality Control

Sequencing quality decreases with size.



Solution: trim end of reads with low quality

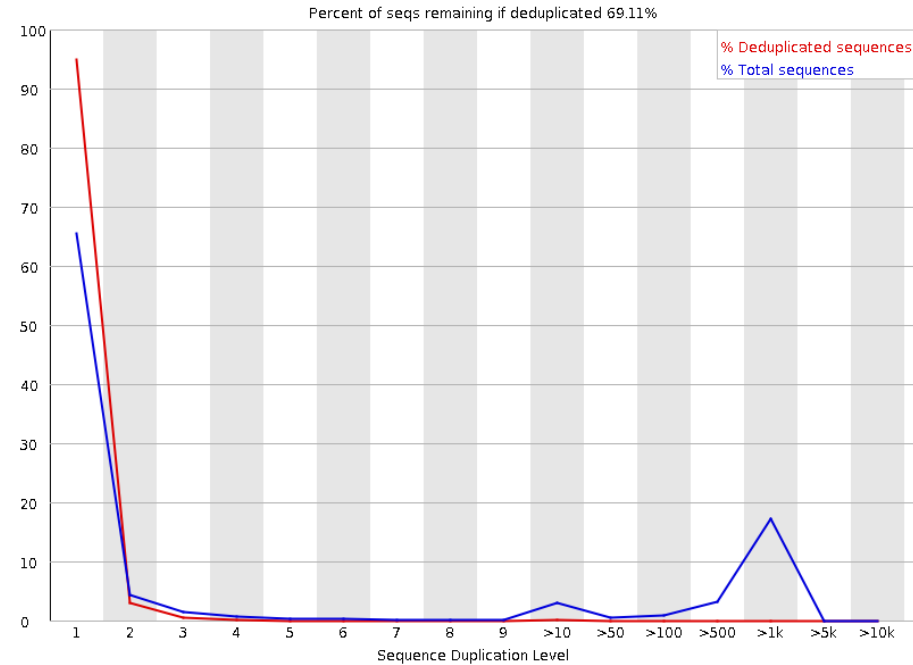
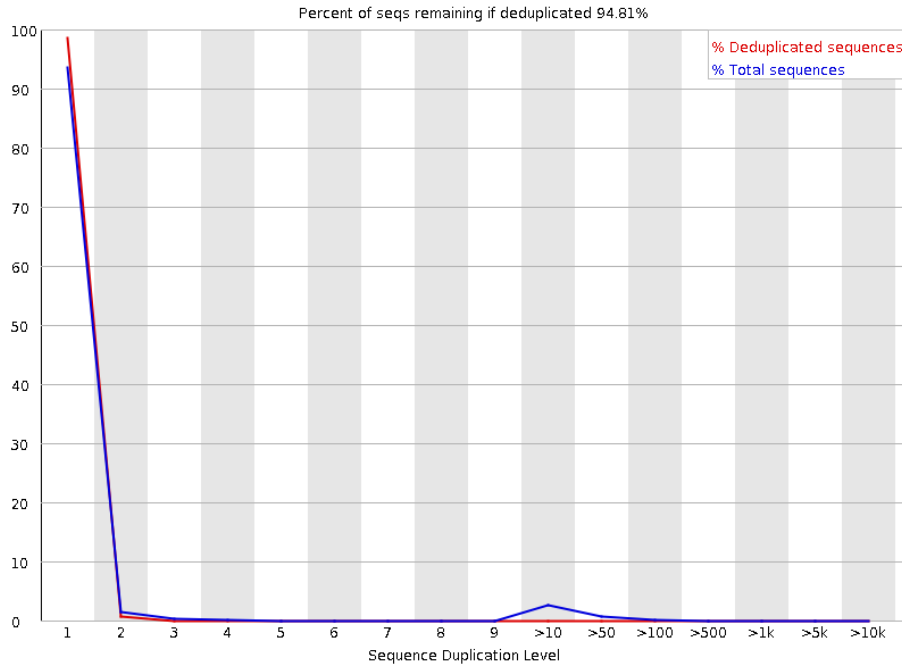
# Read position sequence bias



- Trim read starts

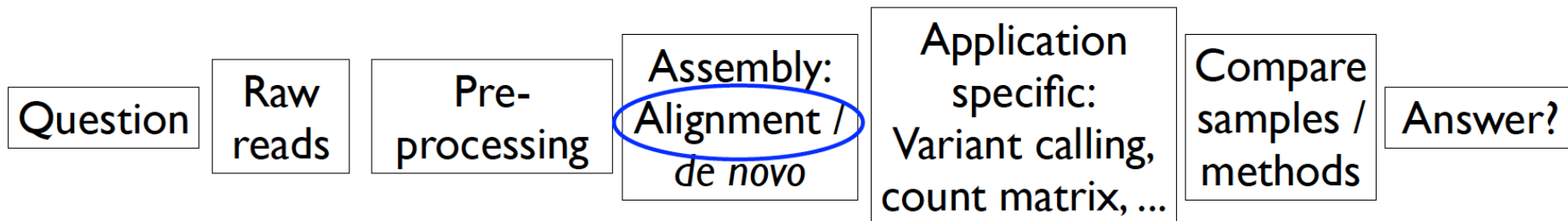
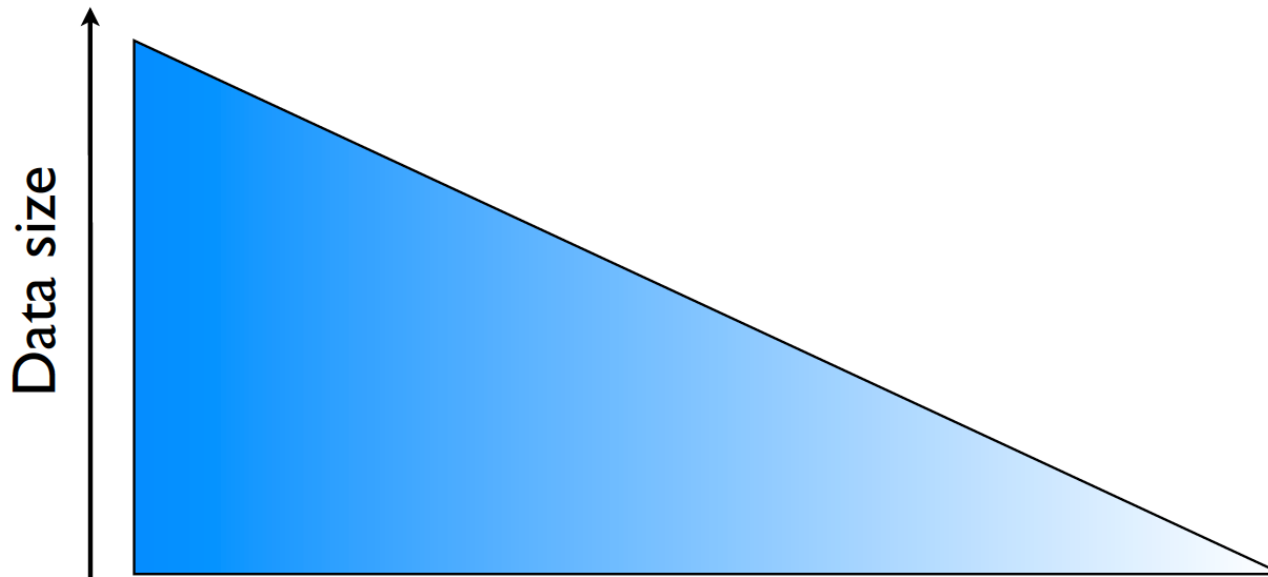
# Quality Control

## Sequence duplication levels



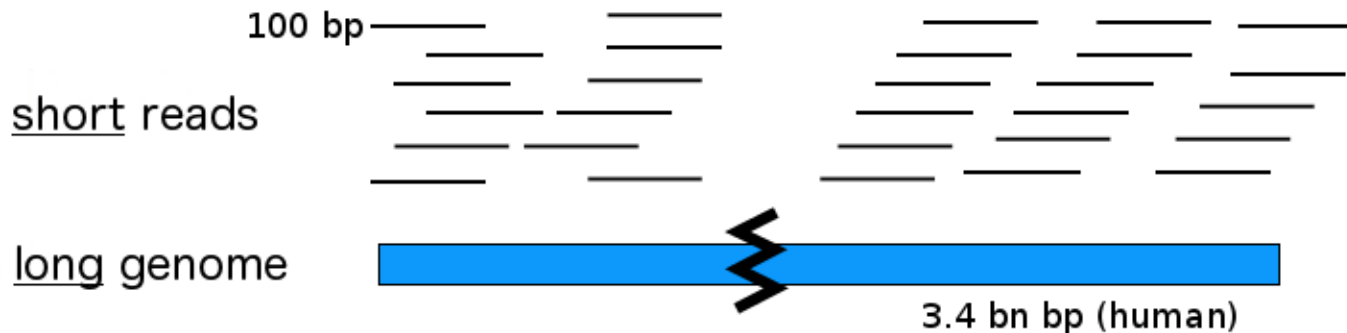
- **Solution: remove duplicates**

# (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

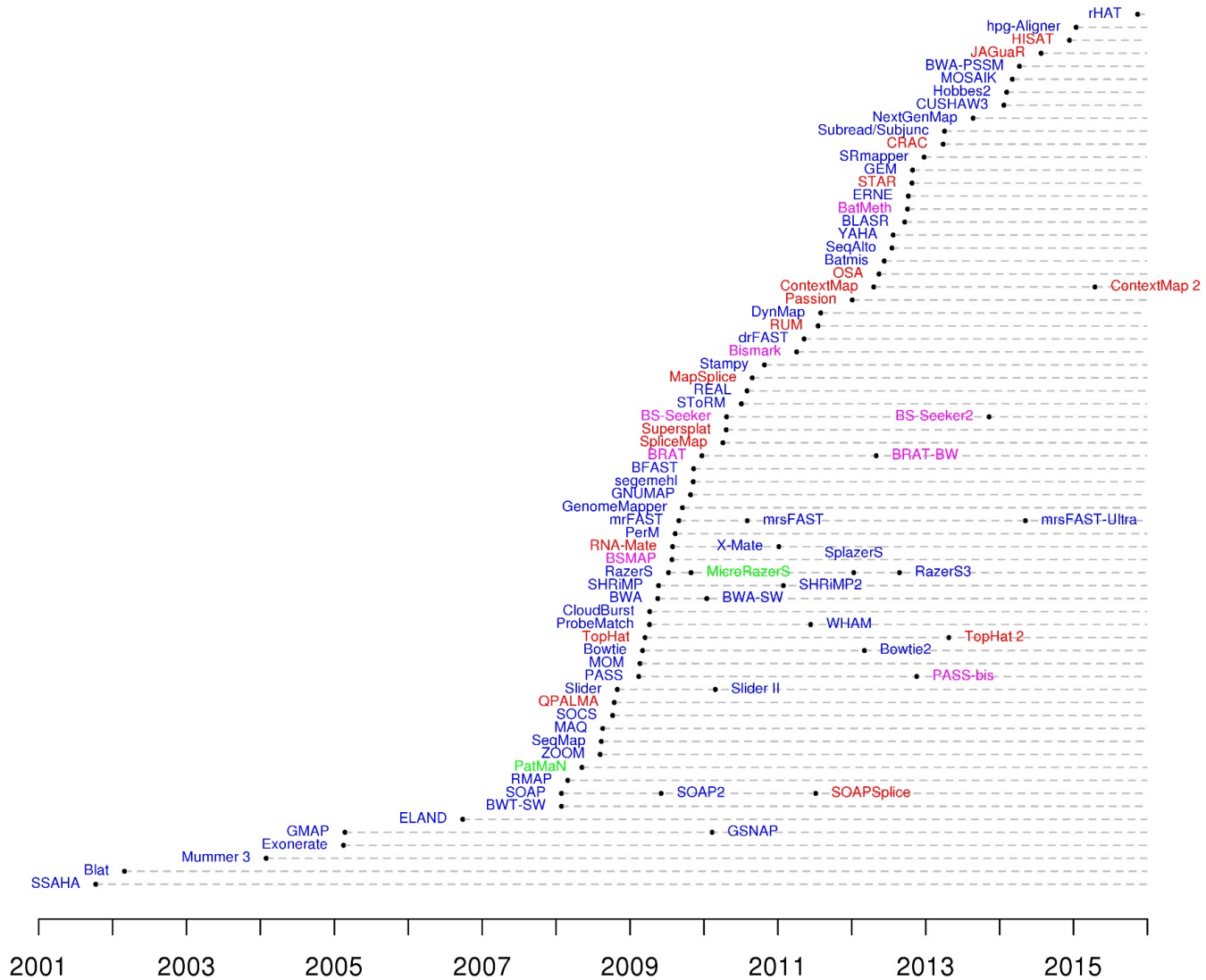
**Solution:**

**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



# Reference based aligners - Overview

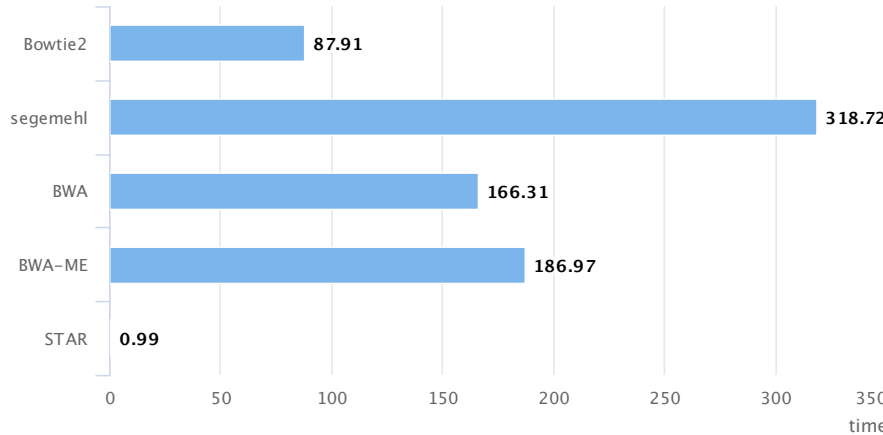
	<i>Time</i>	<i>Precision</i>	<i>Pairs</i>	<i>GAPS</i>	<i>Phred</i>	<i>Memory</i>	<i>Application (Comments)</i>
<b>BOWTIE</b>	+		+	-	-	<b>5GB</b>	<b>General</b> <i>(max. 3 missmatches)</i>
<b>BWA</b>	+		+	+	+	<b>8GB</b>	<b>General</b> <i>(max of 200bps reads)</i>
<b>NOVOALIGN</b>		+	+	+	+	<b>8GB</b>	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

non comprehensive list

# Alignment Tools

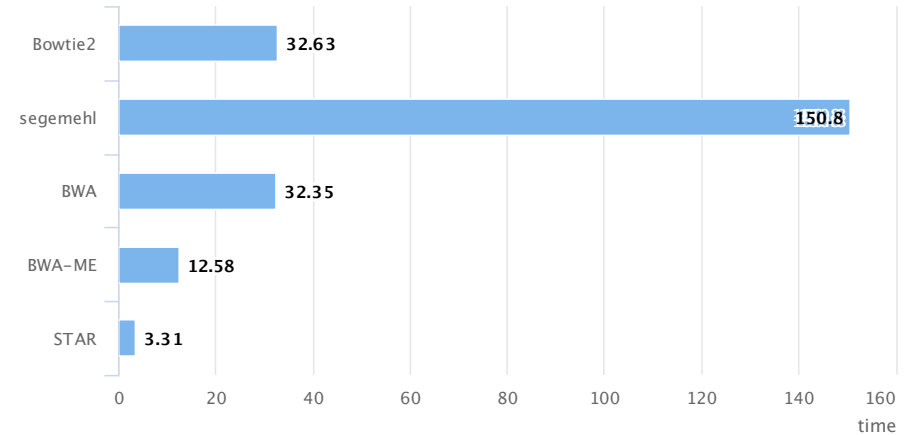
User time [s]

DNA-Seq



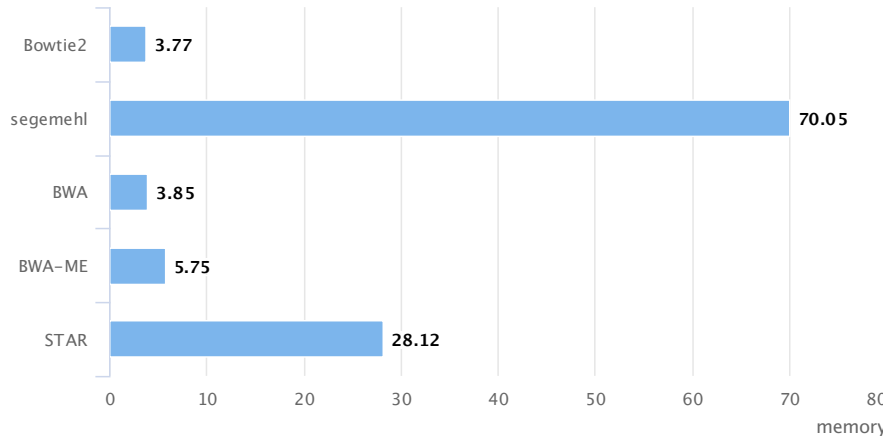
User time [s] (mRNA-Seq)

mRNA-Seq



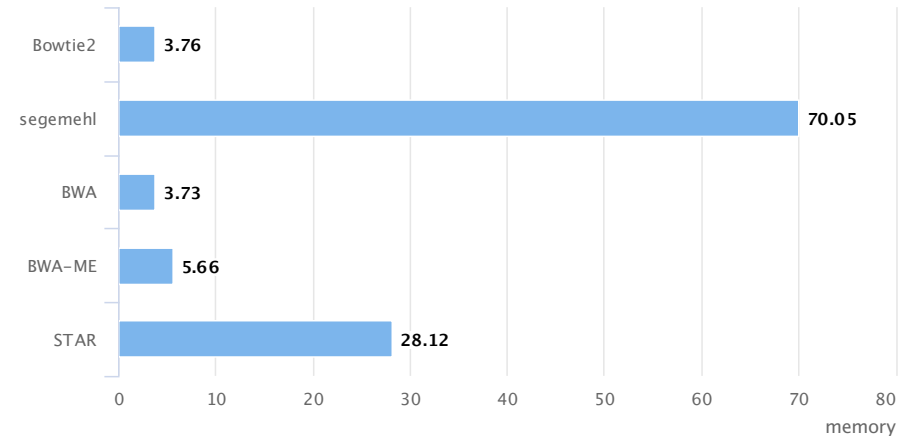
Memory consumption [GB]

DNA-Seq



Memory consumption [GB]

mRNA-Seq



# SAM Files

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

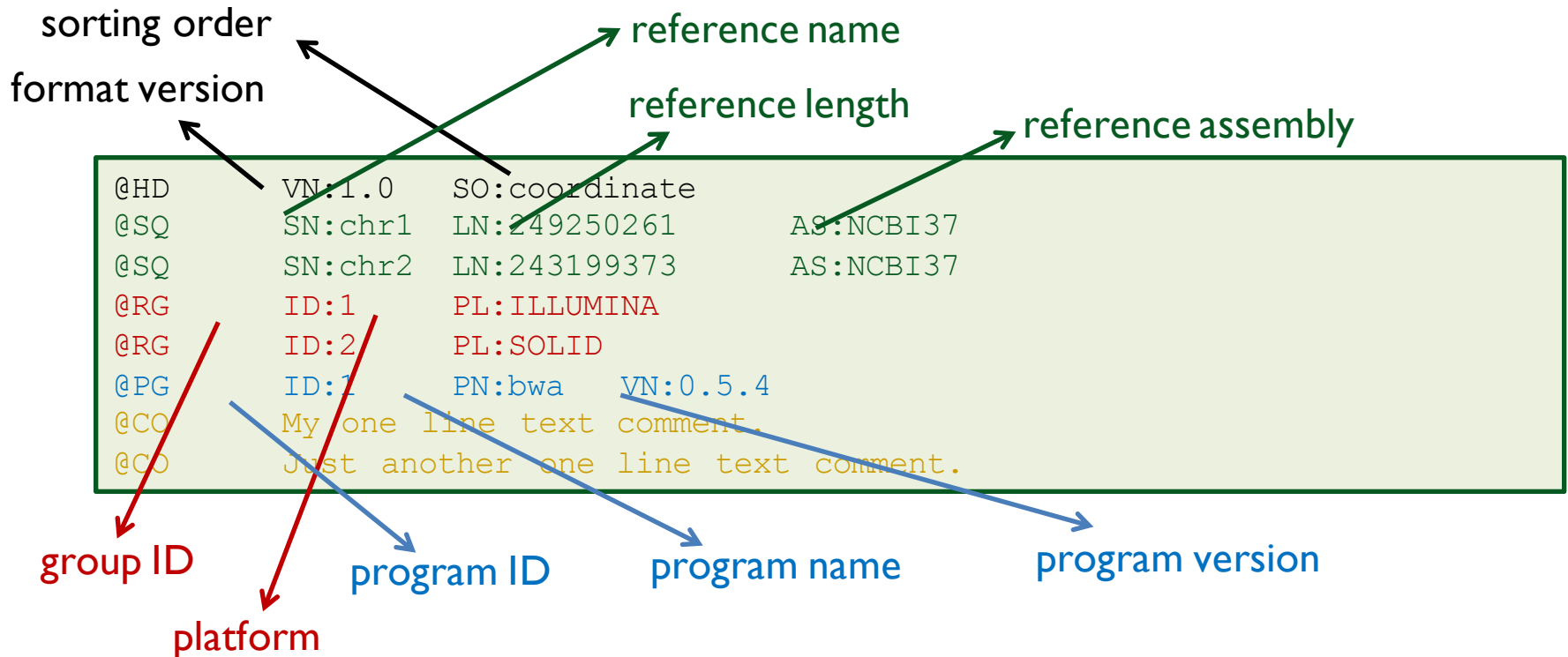
Header

```
@HD VN:1.5 S0: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 - Header

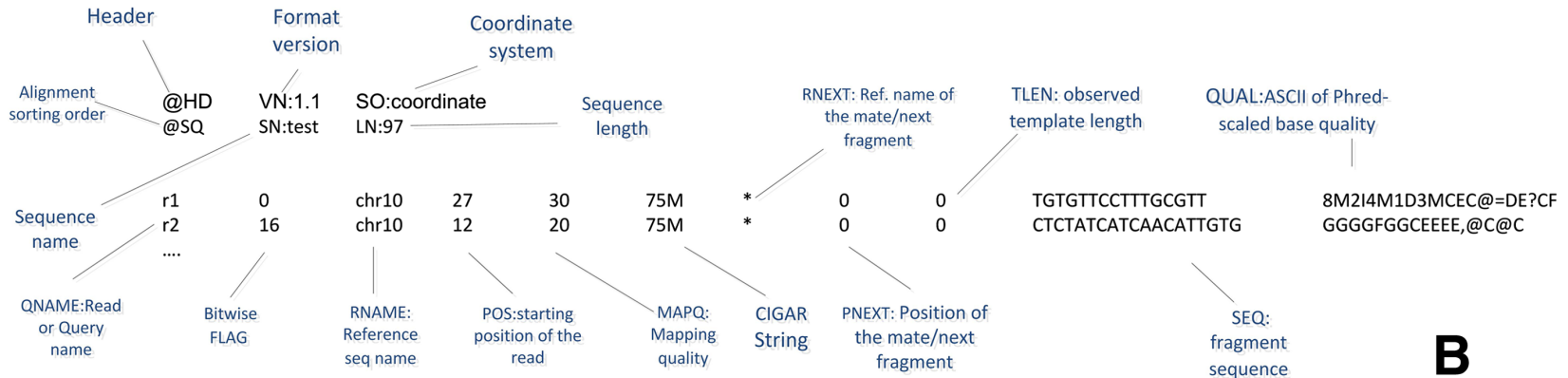
- ▶@HD – Header line.
- ▶@SQ – Reference genome information.
- ▶@RG – Read group information.
- ▶@PG – Program (software) information.
- ▶@CO – Commentary line.



# SAM Files- alignment section

Coordinates 123456789...  
 Reference AAATGAATAATCTCTATCATCAACATTGTGTTCCCTTTGCGTTTTAACCTTTCCT  
 Reads r1 TGTGTTCCCTTTGCGTTI  
 r2 CTCTATCATCAACATTGTG

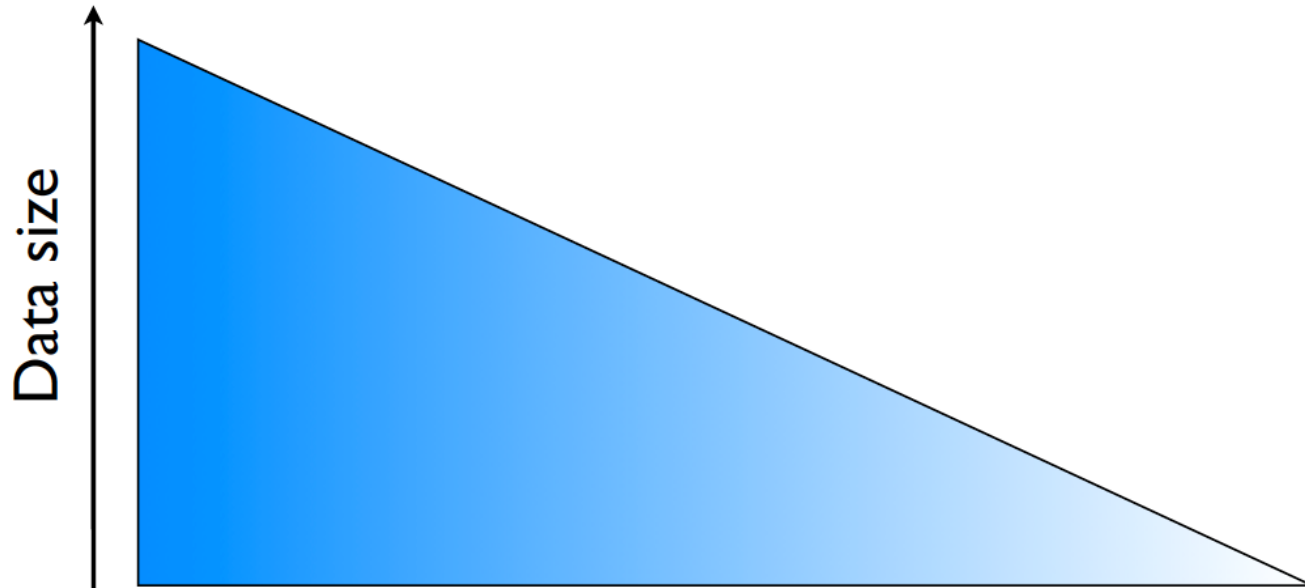
**A**



**B**

- 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 - Peak Calling



Question

Raw reads

Pre-processing

Assembly:  
Alignment /  
*de novo*

Application

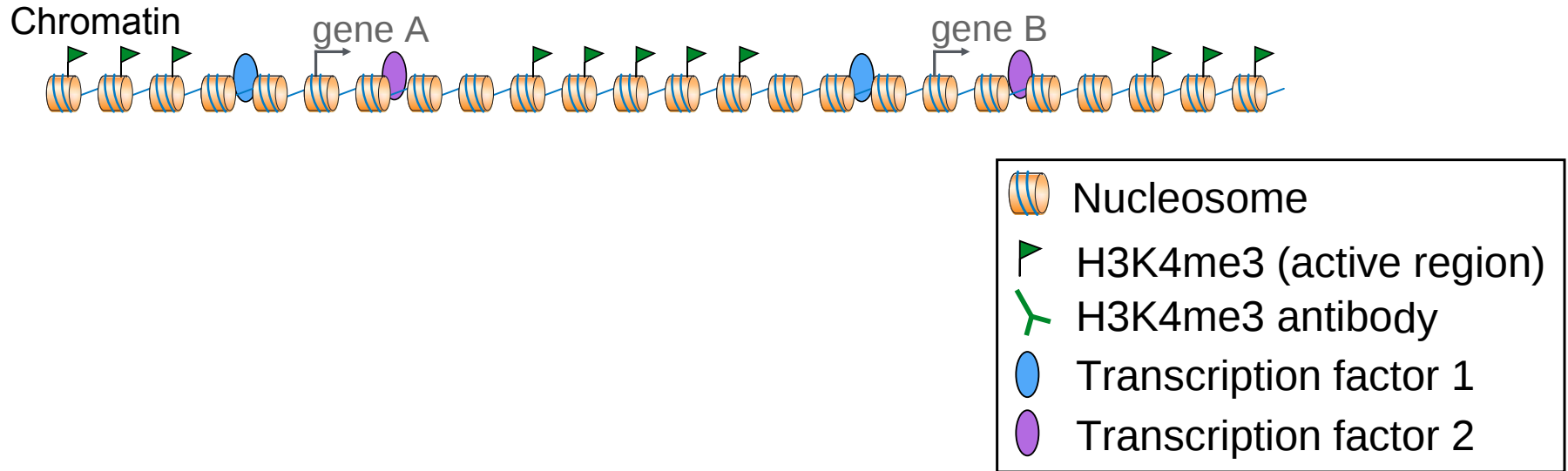
Peak Calling

Compare samples /  
methods

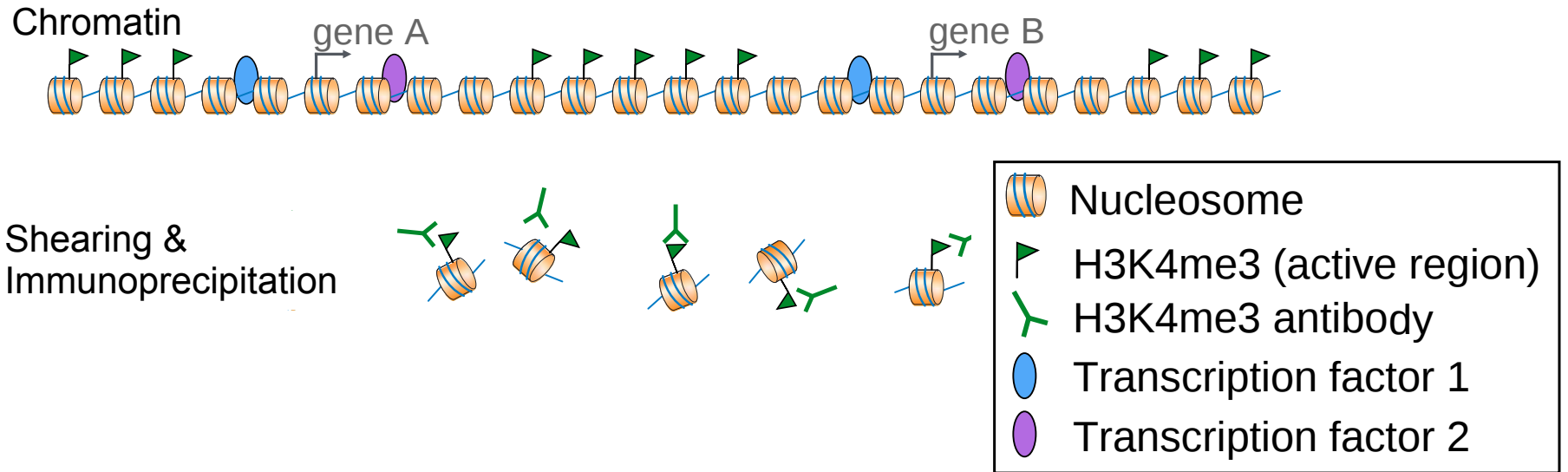
Answer?



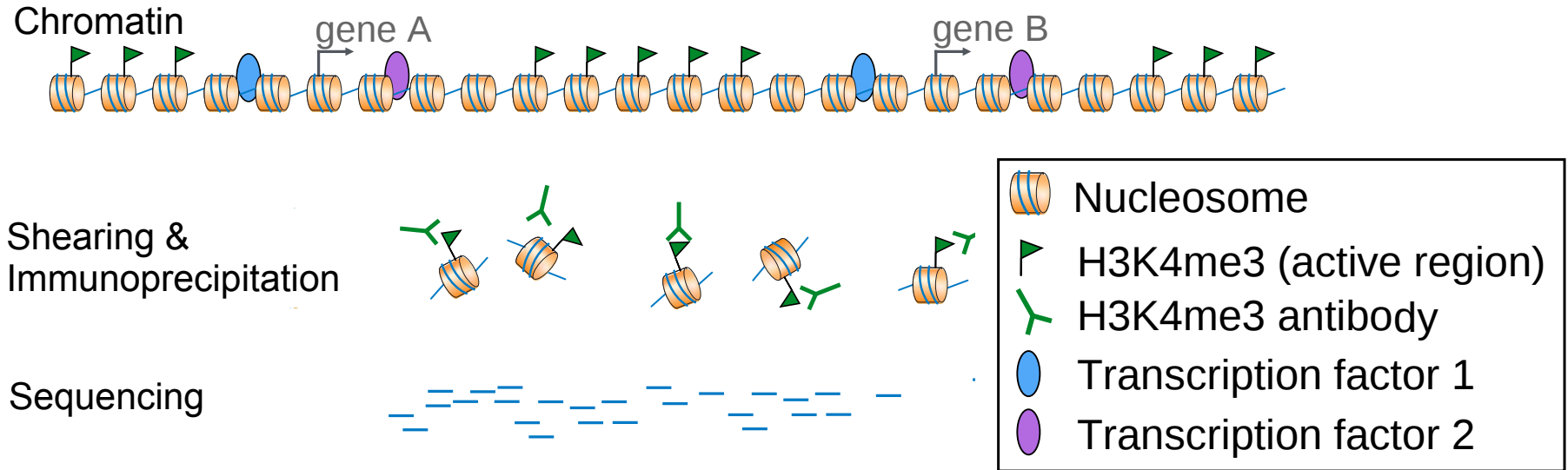
# DNA - Protein interactions with ChIP-Seq



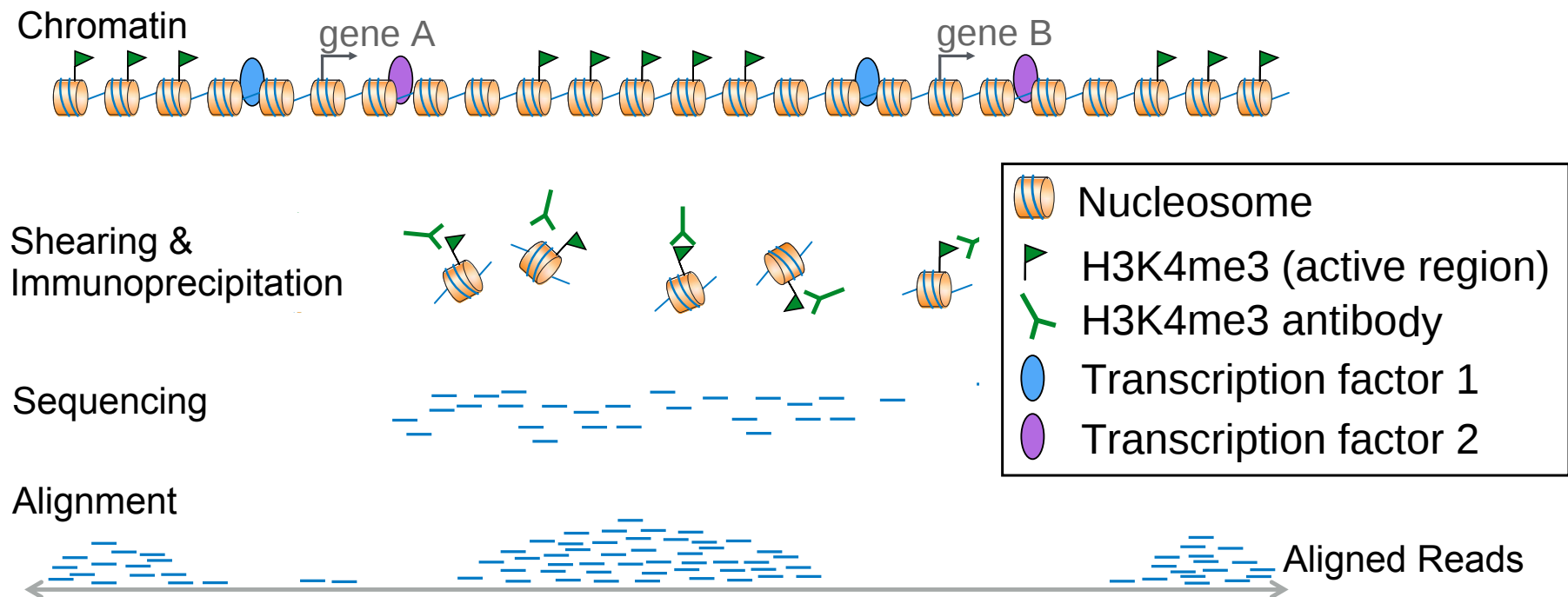
# DNA - Protein interactions with ChIP-Seq



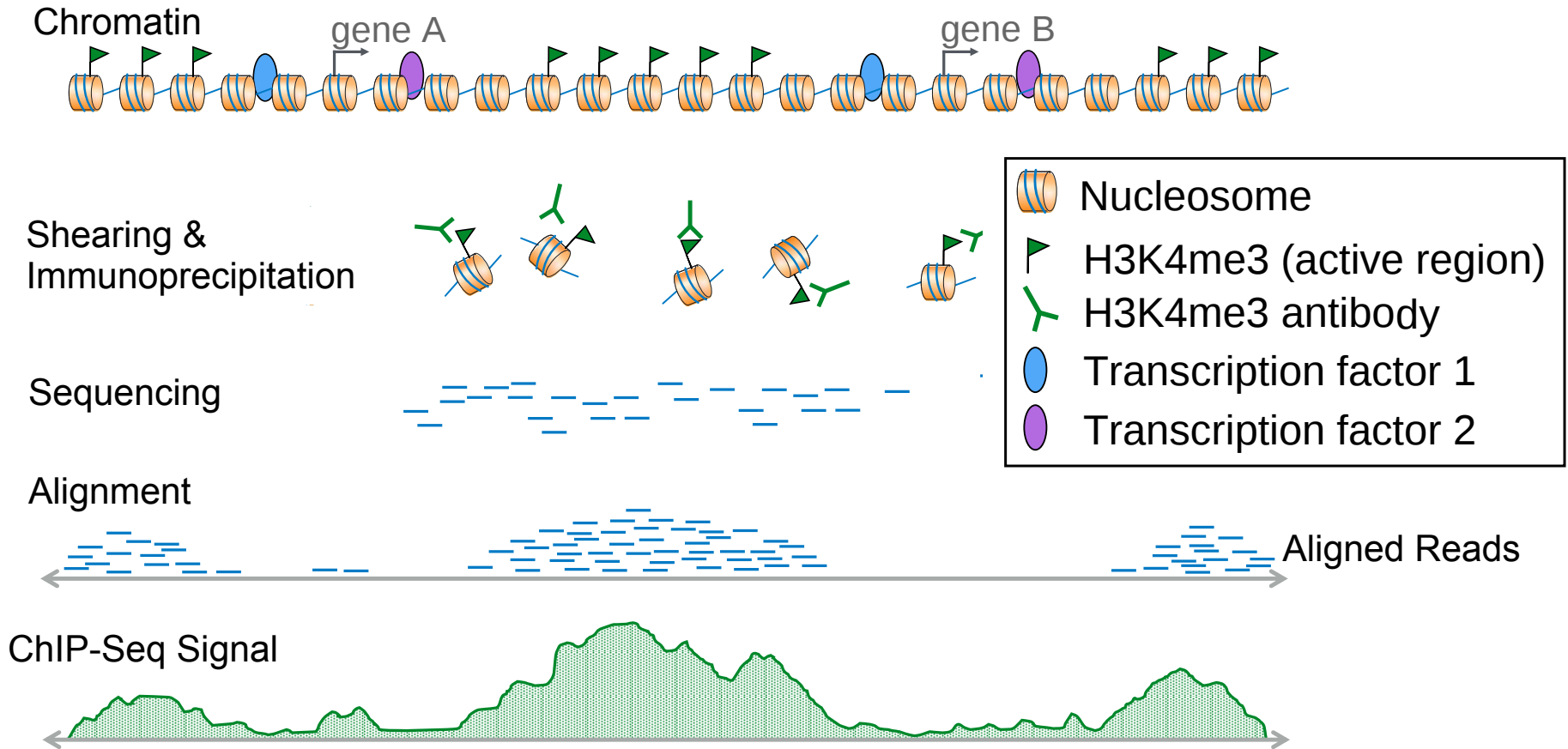
# DNA - Protein interactions with ChIP-Seq



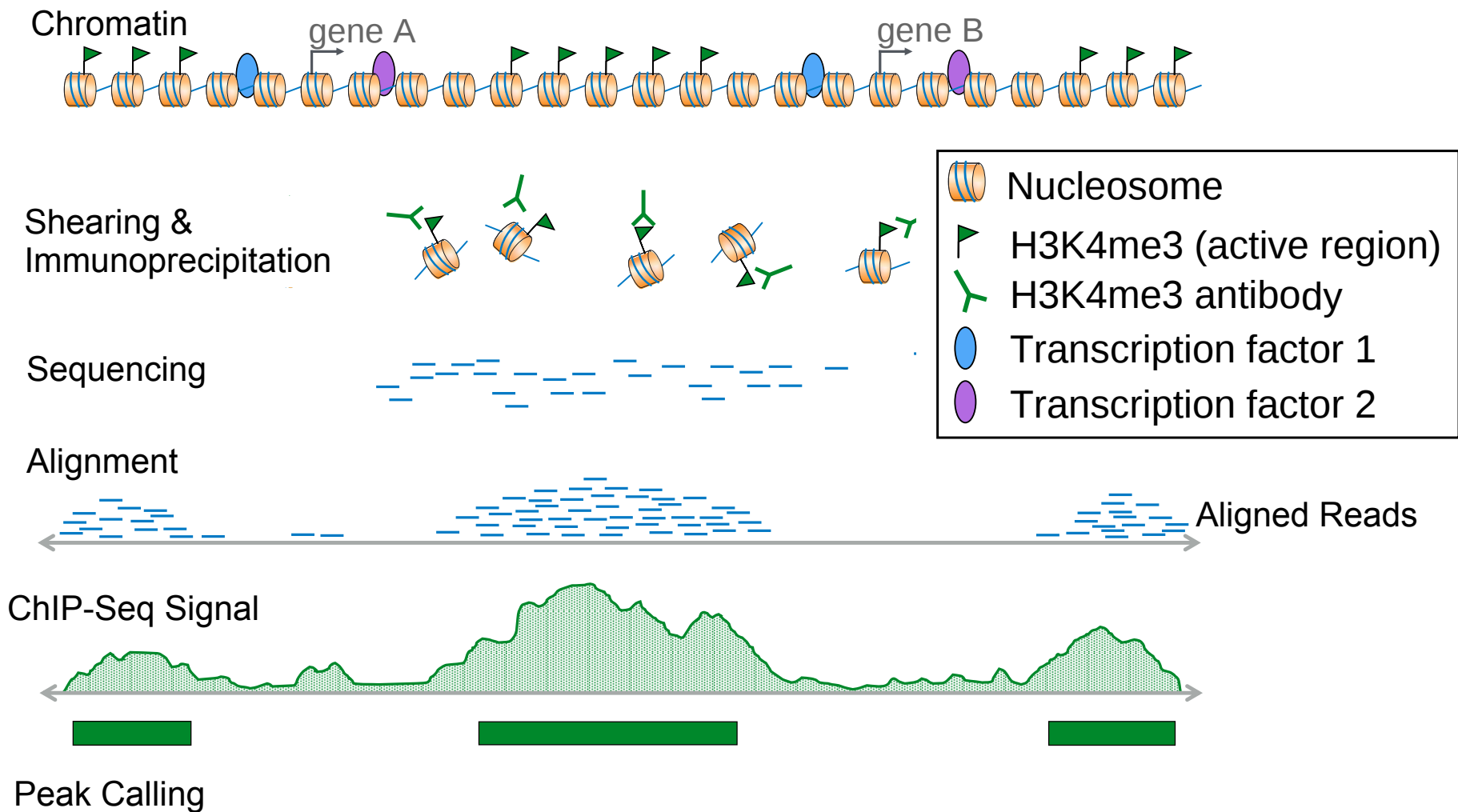
# DNA - Protein interactions with ChIP-Seq



# DNA - Protein interactions with ChIP-Seq

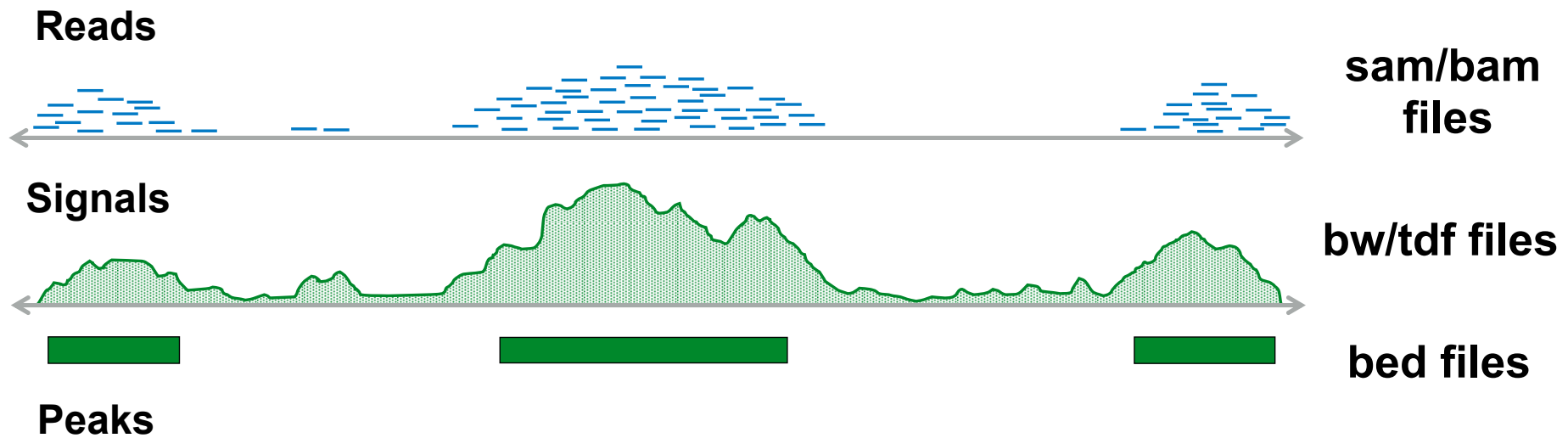


# DNA - Protein interactions with ChIP-Seq



# ChIP-Seq - Data Files

---



# Peaks - Bed files

- Peaks / genomic regions are stored in bed files

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



# Genomic Signals - WIG/TDF Files

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

- Example of WIG file

Header                      Chromosome

```
variableStep    chrom=chr1
140000    30.5
140100    25.1
141200    14
142000    -32.8
```

Genomic Coordinate                      Signal

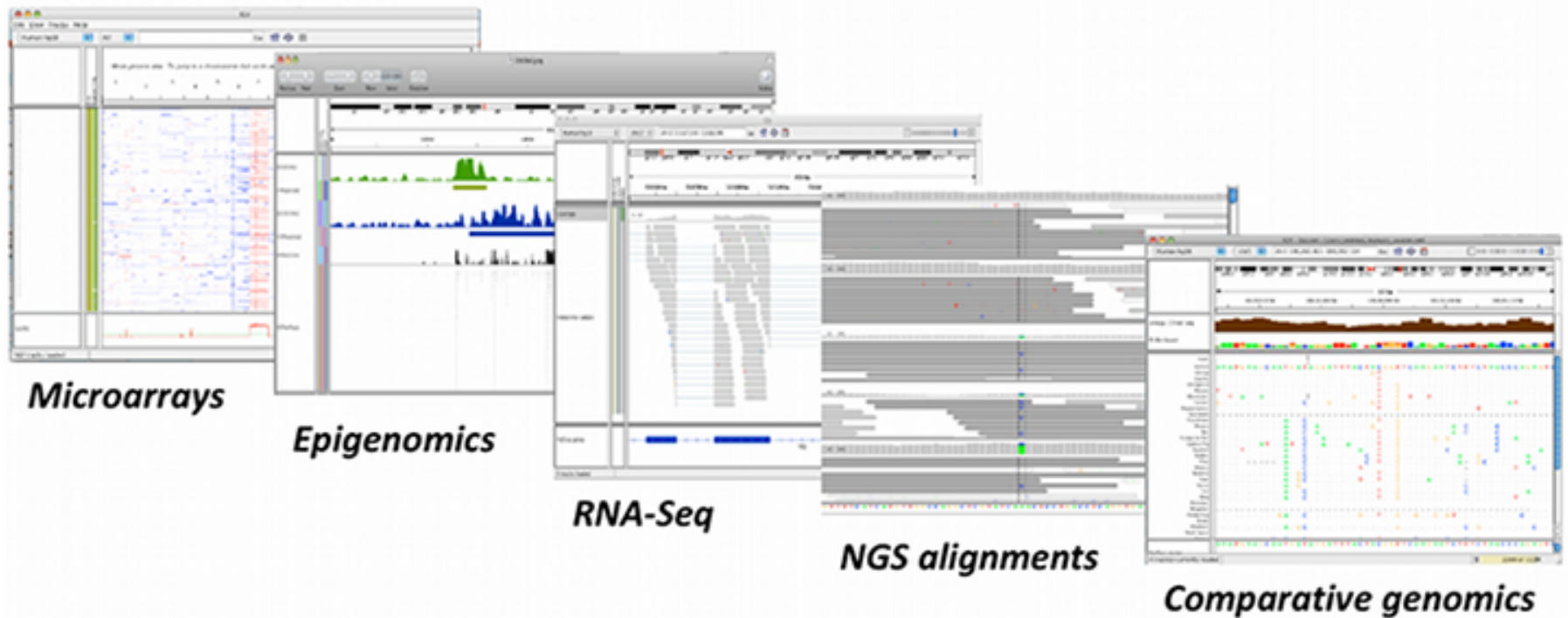
- In practice, we use binary version of WIG (BIGWIG) or TDF files

# 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>

**Integrative Genomics Viewer**

**Home**

**Integrative Genomics Viewer**

**Overview**

The **Integrative Genomics Viewer (IGV)** is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets. It supports a wide variety of data types, including array-based and next-generation sequence data, and genomic annotations.

**Downloads**

Download the IGV desktop application and igvtools.

**Citing IGV**

To cite your use of IGV in your publication:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P.

**Funding**

Development of IGV has been supported by funding from the National Cancer Institute (NCI) of the National Institutes of Health (NIH).

**Search website**

search

© 2013-2018  
Broad Institute, and  
the Regents of the  
University of California

**Navigation Menu:** Home, Downloads, Documents, Hosted Genomes, FAQ, IGV User Guide, File Formats, Release Notes, Credits, Contact.

# Launch IGV

- Home
- Downloads
- Documents

- Hosted Genomes
- FAQ
- IGV User Guide
- File Formats
- Release Notes
- Credits
- Contact

Search website

search


© 2013-2018  
Broad Institute, and  
the Regents of the  
University of California

**NOTE: IGV 2.4.x releases require [Java 8](#). For Java 10 see the [development snapshot build](#).**

## Install IGV


### Download IGV Mac App

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



### Download IGV on Windows

Download and unzip the Archive, then double-click the *igv.bat* file to run IGV. See [readme.txt](#) to run IGV from the command line



**For high DPI screens:** Use Java 10 and the [development snapshot build of IGV](#).

Keynote      Download IGV to run on Linux / MacOS command line

Institute for  
Computational Genomics  
01011011010  
10100100101



# Launch IGV

IGV\_2.4.14 File Genomes View Tracks Regions Tools GenomeSpace Help

IGV

Human hg19 All Go

Select genome from the drop-down menu

1 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 X Y

RefSeq Genes

2 tracks loaded 413M of 1,373M

# Launch IGV

The screenshot displays the IGV\_2.4.14 application window. The 'File' menu is open, listing several options. A green arrow points to the 'Load from ENCODE (2012)...' option. A yellow callout box with the text 'Load data from the ENCODE project' is positioned over the arrow. The main interface shows a genome browser with a track for 'Human hg19' and a 'RefSeq Genes' track at the bottom. The status bar at the bottom indicates '2 tracks loaded' and '501M of 1.373M'.

- File
- Genomes
- View
- Tracks
- Regions
- Tools
- GenomeSpace
- Help

IGV

Human hg19

Go

3 4 5 7 9 11 12 13 14 15 16 17 18 19 20 21 22 X Y

Load data from the ENCODE project

RefSeq Genes

2 tracks loaded

501M of 1.373M



# Launch IGV

The screenshot shows the IGV interface with a search dialog box. The dialog box title is 'Encode Production Data'. The search filter is 'CTCF GM12878'. A table of results is displayed below the filter. A yellow highlight and a green arrow point to the search filter. The table has columns for cell, dataType, antibody, view, replicate, type, lab, and hub.

cell	dataType	antibody	view	replicate	type	lab	hub
GM12878	ChipSeq	CTCF	Peaks		narrowPeak	Broad	Data
GM12878	ChipSeq	CTCF_(SC-15...	Peaks		narrowPeak	Stanford	Data
GM12878	ChipSeq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChipSeq						Data
GM12878	ChipSeq						Data
GM12878	ChipSeq						Data
GM12878	ChipSeq	CTCF	Signal		bigWig	Broad	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_Overlap...		bigWig	UT-A	Data
GM12878	ChipSeq	CTCF	Peaks		narrowPeak	UT-A	Data
GM12878	ChipSeq	CTCF	Signal		bigWig	UT-A	Data
GM12878	ChipSeq	CTCF_(SC-15...	Alignments	1	bam	Stanford	Data
GM12878	ChipSeq	CTCF_(SC-15...	Alignments	2	bam	Stanford	Data
GM12878	ChipSeq	CTCF_(SC-15...	Peaks		narrowPeak	Stanford	Data
GM12878	ChipSeq	CTCF_(SC-15...	Signal		bigWig	Stanford	Data
GM12878	ChipSeq	CTCF	Alignments	1	bam	UW	Data
GM12878	ChipSeq	CTCF	Alignments	2	bam	UW	Data
GM12878	ChipSeq	CTCF	Hotspots	1	broadPeak	UW	Data
GM12878	ChipSeq	CTCF	Hotspots	2	broadPeak	UW	Data
GM12878	ChipSeq	CTCF	Peaks	1	narrowPeak	UW	Data
GM12878	ChipSeq	CTCF	Peaks	2	narrowPeak	UW	Data
GM12878	ChipSeq	CTCF	RawSignal	1	bigWig	UW	Data
GM12878	ChipSeq	CTCF	RawSignal	2	bigWig	UW	Data
GM12878	ChipSeq	CTCF	Peaks		bigBed	Broad	analysis

# Launch IGV

IGV

Human hg19

Encode Production Data

Filter:  50 rows

<input type="checkbox"/>	cell	dataType	antibody	view	replicate	type	lab	hub
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks		narrowPeak	Broad	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF_(SC-15...	Peaks		narrowPeak	Stanford	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks		narrowPeak	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks		narrowPeak	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	1	bam	Broad	
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	2	bam	Broad	
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks		broadPeak	Broad	
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Signal		bigWig	Broad	
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	1	bam	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	2	bam	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	3	bam	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Base_Overlap...		bigWig	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks		narrowPeak	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Signal		bigWig	UT-A	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF_(SC-15...	Alignments	1	bam	Stanford	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF_(SC-15...	Alignments	2	bam	Stanford	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF_(SC-15...	Peaks		narrowPeak	Stanford	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF_(SC-15...	Signal		bigWig	Stanford	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	1	bam	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Alignments	2	bam	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Hotspots	1	broadPeak	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Hotspots	2	broadPeak	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks	1	narrowPeak	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks	2	narrowPeak	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	RawSignal	1	bigWig	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	RawSignal	2	bigWig	UW	Data
<input type="checkbox"/>	GM12878	ChipSeq	CTCF	Peaks		bigBed	Broad	analysis

sort by lab

Load Cancel

RefSeq Genes

# Launch IGV

IGV

Human hg19 All

Encode Production Data

Filter: CTCF GM12878 50 rows

cell	dataType	antibody	view	replicate	type	lab	hub
<input checked="" type="checkbox"/>	GM12878	ChIPSeq	CTCF	Peaks	narrowPeak	Broad	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Alignments	bam	Broad	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Alignments	bam	Broad	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Peaks	broadPeak	Broad	Data
<input checked="" type="checkbox"/>	GM12878	ChIPSeq	CTCF	Signal	bigWig	Broad	
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Peaks	bigBed	Broad	
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Peaks	bigBed	Broad	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Signal	bigWig	Broad	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Peaks	bigBed	Broad	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Signal	bigWig	Broad	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_SC-15...	Peaks	narrowPeak	Stanford	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_SC-15...	Alignments	bam	Stanford	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_SC-15...	Alignments	bam	Stanford	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_SC-15...	Peaks	narrowPeak	Stanford	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_SC-15...	Signal	bigWig	Stanford	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_C-20	Peaks	bigBed	Stanford	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_C-20	Peaks	bigBed	Stanford	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_C-20	Signal	bigWig	Stanford	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_C-20	Peaks	bigBed	Stanford	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_C-20	Peaks	bigBed	Stanford	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF_C-20	Signal	bigWig	Stanford	analysis
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Peaks	narrowPeak	UT-A	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Alignments	bam	UT-A	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Alignments	bam	UT-A	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Alignments	bam	UT-A	Data
<input type="checkbox"/>	GM12878	ChIPSeq	CTCF	Base_Overlap...	bigWig	UT-A	Data

sort by lab

Load Cancel

RefSeq Genes

# Launch IGV - Example on CTCF ChIP-seq



# 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

---

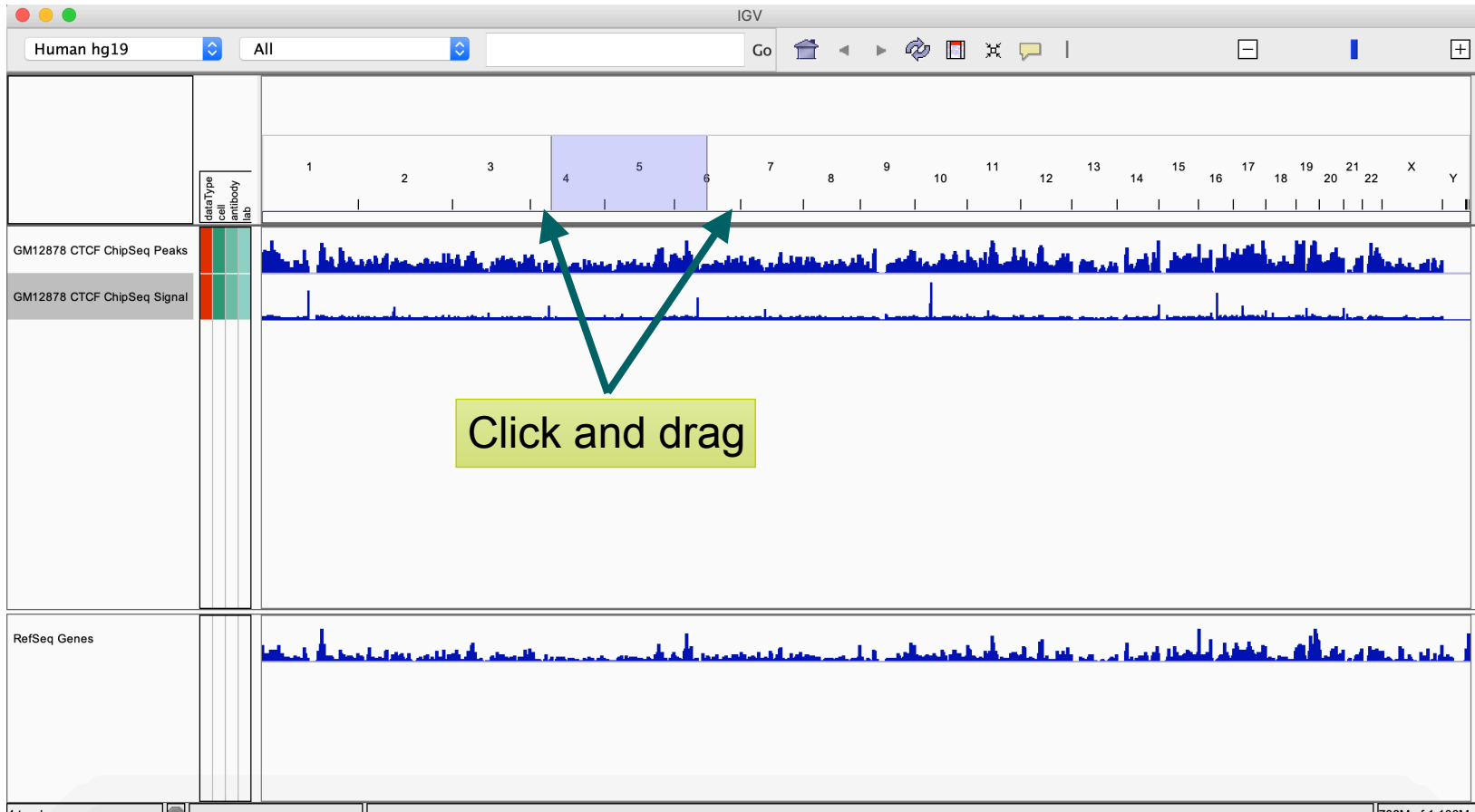
Note: for large files use indexed formats

# Hands on

---

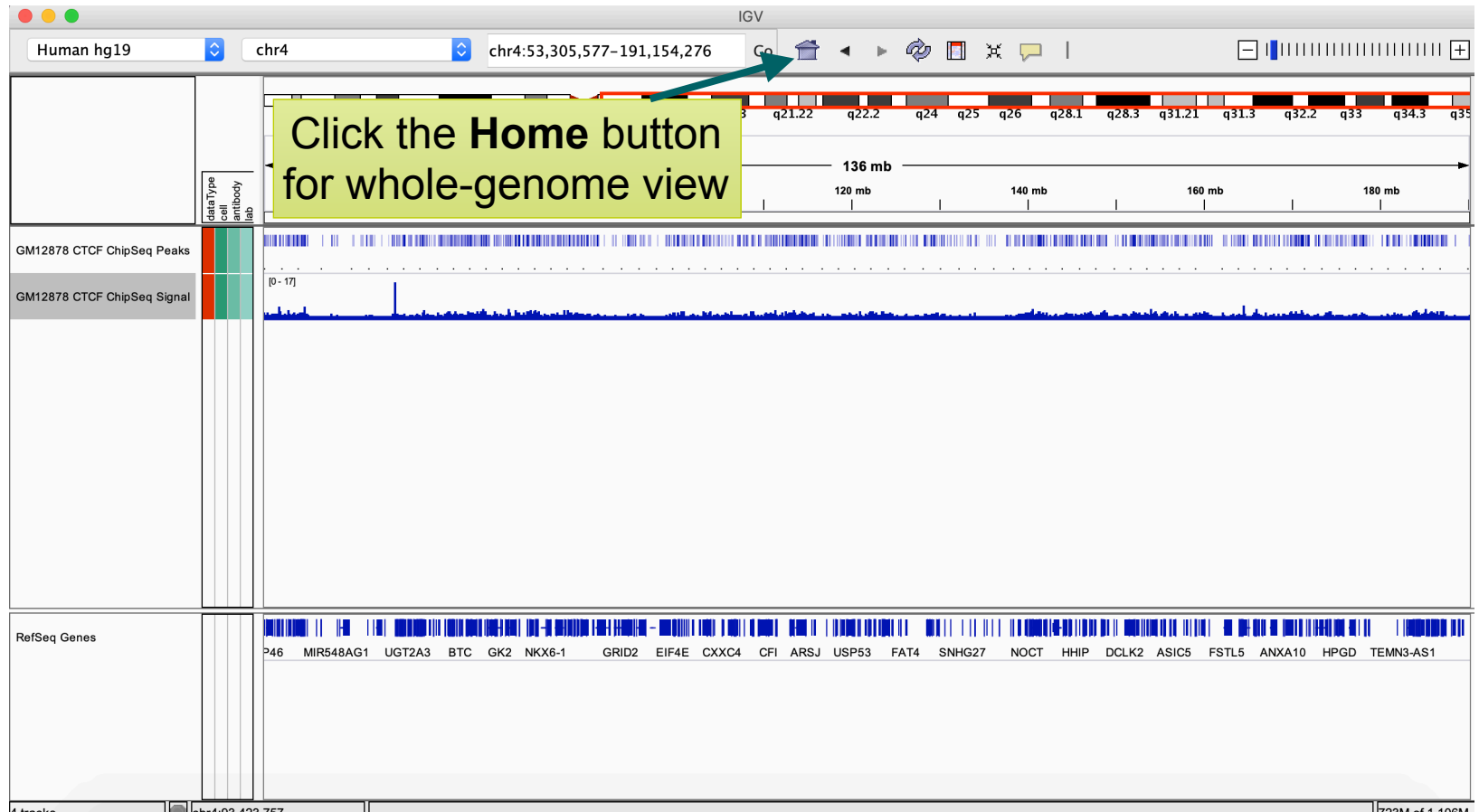
- Lunch IGV on your computer
- Choose human genome hg19
- Load data from ENCODE project
  - ChIP-seq of factor CTCF for GM12878 cell type
  - Peaks and signal

# Navigate





# Navigate



# Navigate

IGV 2.4.14 File Genomes View Tracks Regions Tools GenomeSpace Help

Human hg19 chr5 chr5:137,799,181-137,807,004 Go

p15.32 p15.1 p14.2 p13.3 p13.1 p11 q11.2 q12.2 q13.2 q14.1 q14.3 q15 q21.2 q22.1 q23.1 q23.3 q31.2 q32 q33.2 q34 q35.1 q35

7,765 bp

137,800,000 bp 137,801,000 bp 137,802,000 bp 137,803,000 bp 137,804,000 bp 137,805,000 bp 137,806,000 bp 137,807,000 bp

GM12878 CTCF ChIP-seq Peaks

GM12878 CTCF ChIP-seq Signal

RefSeq Genes

EGR1

Type gene name or other annotation into RefSeq Genes and click **Go**

# Navigate

IGV 2.4.14 File Genomes View Tracks Regions Tools GenomeSpace Help

Human hg19 chr5 chr5:137,799,181-137,807,004 Go

7,765 bp

137,800,000 bp 137,801,000 bp 137,802,000 bp 137,803,000 bp 137,804,000 bp 137,805,000 bp 137,806,000 bp

GM12878 CTCF ChIPSeq Signal

- Rename Track...
- Change Track Color (Positive Values)...
- Change Track Color (Negative Values)...
- Change Track Height...
- Change Font Size...

Type of Graph

- Heatmap
- Bar Chart
- Points
- Line Plot

Windowing Function

- Minimum
- Mean
- Maximum
- None

RefSeq Genes

- Set Data Range...
- Set Heatmap Scale...
- Log scale
- Autoscale
- Show Data Range
- Save image...
- Export track names...

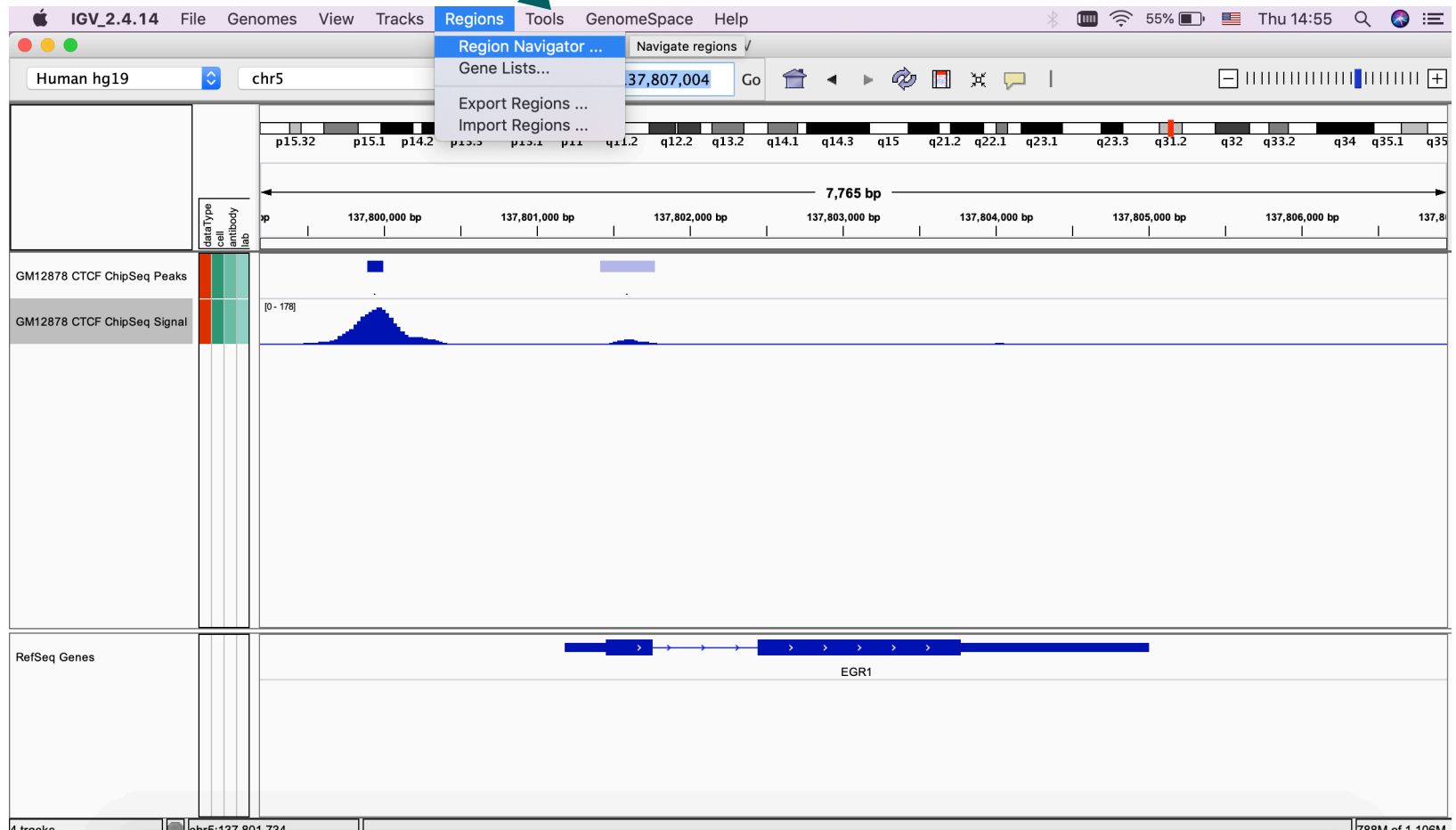
Remove Track

Adjust the scale

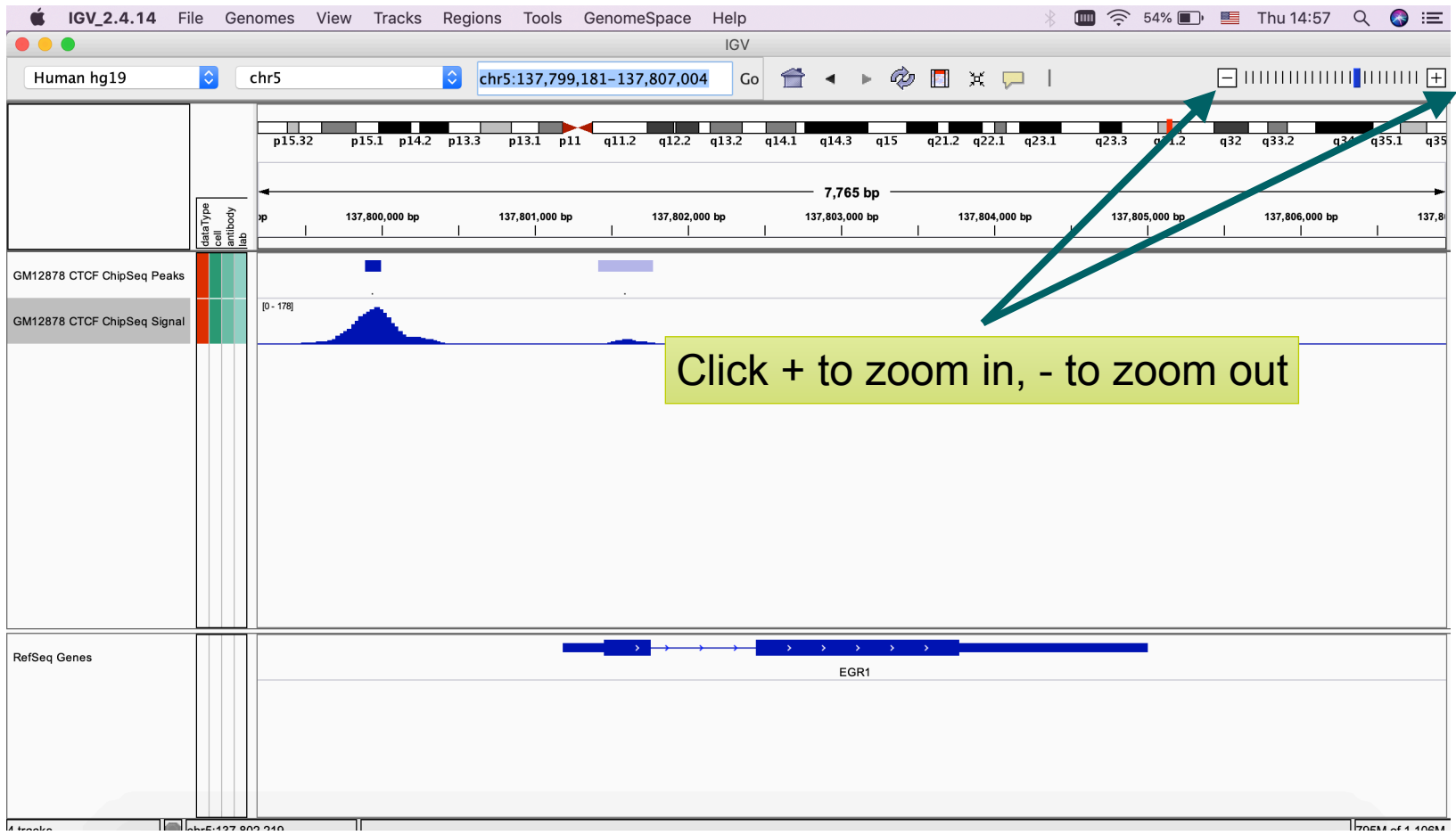
EGR1

# Navigate

Add interesting regions



# Navigate



# Viewing multiple regions



# Change the track color

The screenshot displays the IGV 2.4.14 application window. The top menu bar includes 'File', 'Genomes', 'View', 'Tracks', 'Regions', 'Tools', 'GenomeSpace', and 'Help'. The address bar shows 'Human hg19', 'chr10', and a genomic region 'chr10:64,569,756-64,580,927'. The main view is split into two panels: the left panel shows chromosome 5 with a region 'chr5:137799182-137807004' highlighted, and the right panel shows chromosome 10 with a region 'chr10:64569757-64580927' highlighted. A context menu is open over the 'GM12878 CTCF ChIPSeq Signal' track in the left panel. The menu options are:

- Rename Track...
- Change Track Color (Positive Values)...**
- Change Track Color (Negative Values)...
- Change Track Height...
- Change Font Size...
- Type of Graph**
  - Heatmap
  - Bar Chart
  - Points
  - Line Plot
- Windowing Function**
  - Minimum
  - Mean
  - Maximum
  - None
- Set Data Range...
- Set Heatmap Scale...
  - Log scale
- Autoscale
- Show Data Range
- Sort frames
- Save image...
- Export track names...
- Remove Track

The right panel shows a ChIP-seq signal plot with a red peak. Below the tracks, gene models for 'EGR1' and 'EGR2' are visible. The status bar at the bottom right indicates '104288 of 4.16784'.

# 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

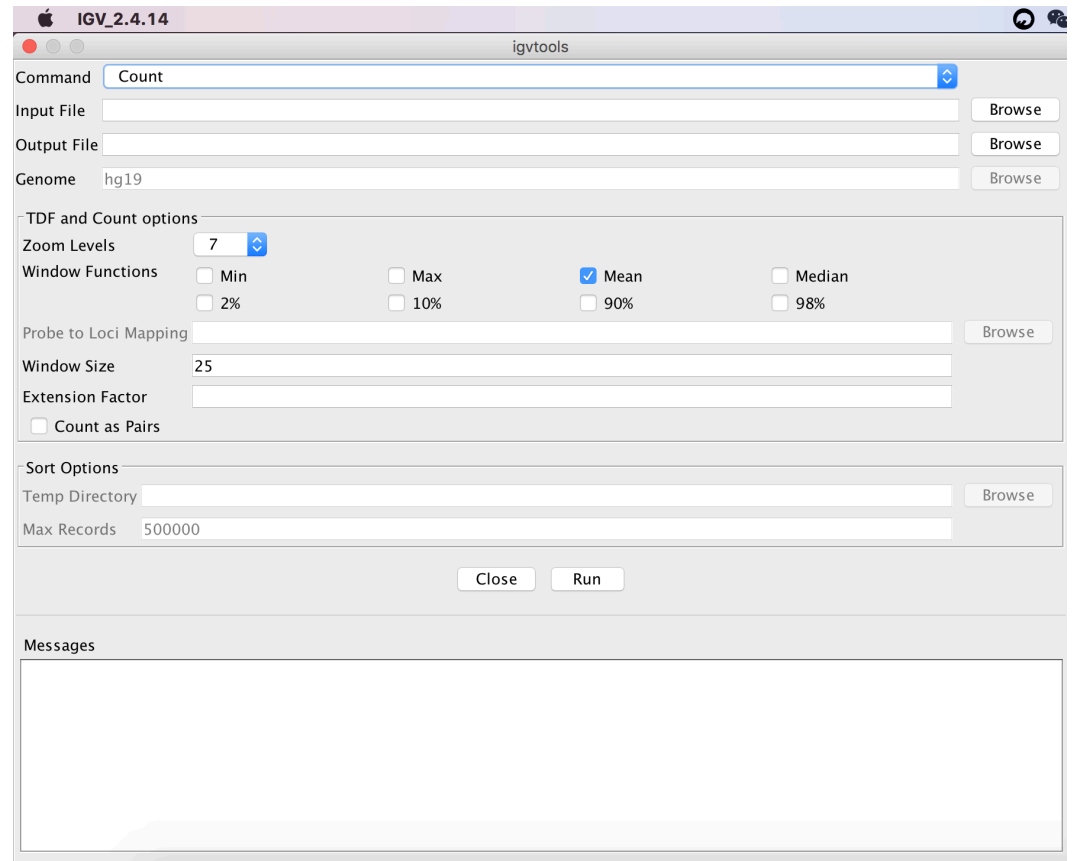
---

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

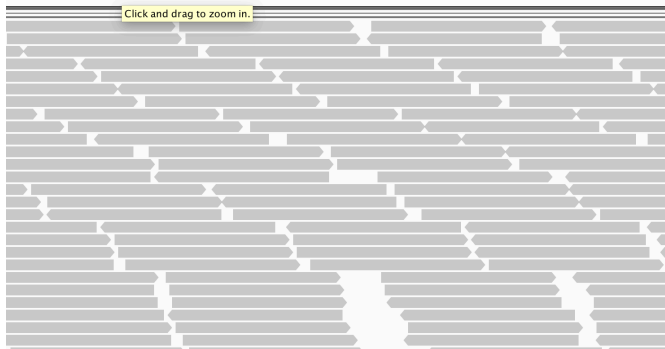
TDF files have the following advantages:

- 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

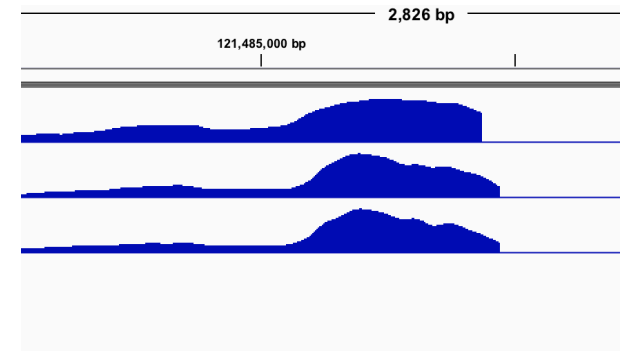
---

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



**Alignment**

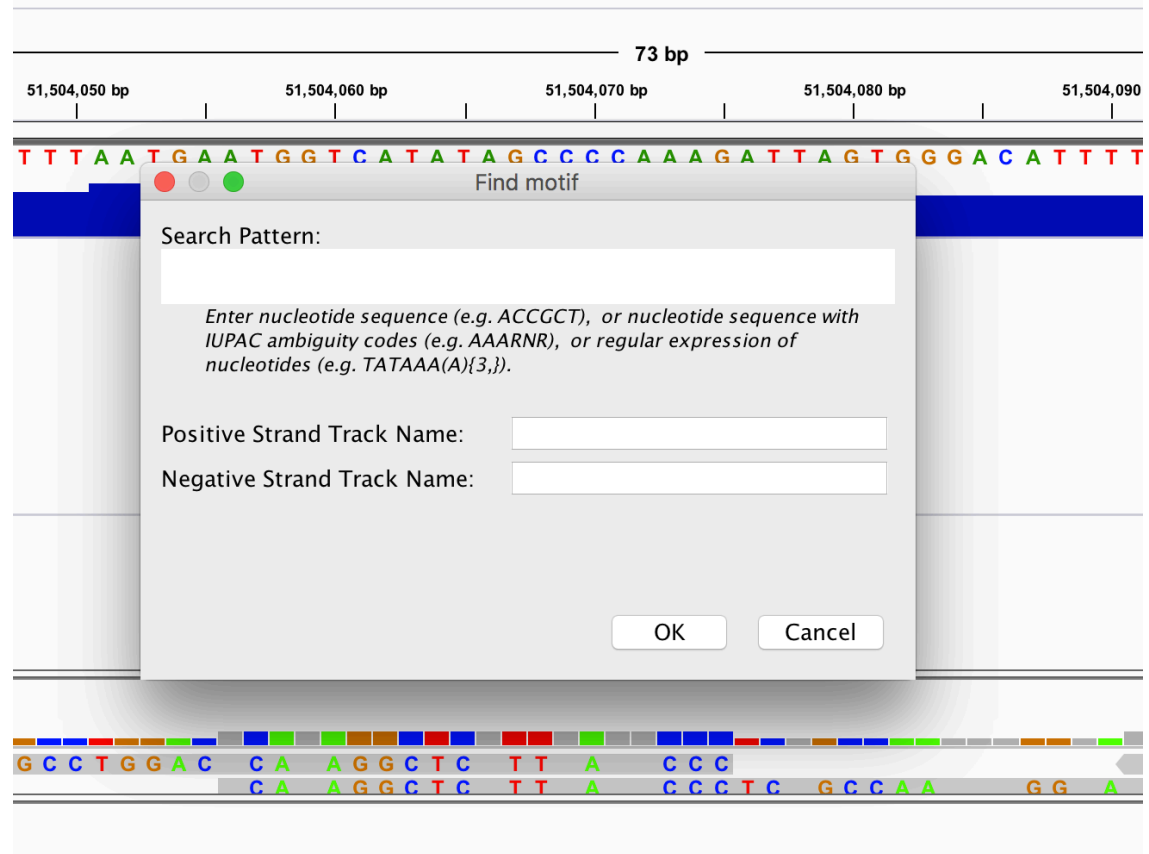
**count**



**Read Density**

# Find Motif

Lunched from the  
IGV user interface  
*Tools > Find Motif*



# hands on

---

download the BAM file from:

[https://costalab.ukaachen.de/open\\_data/](https://costalab.ukaachen.de/open_data/)

[Bioinformatic Analysis in R 2018/BIAR\\_D5/practice/igv\\_data.zip](https://costalab.ukaachen.de/open_data/Bioinformatic_Analysis_in_R_2018/BIAR_D5/practice/igv_data.zip)

Files included in the zip file:

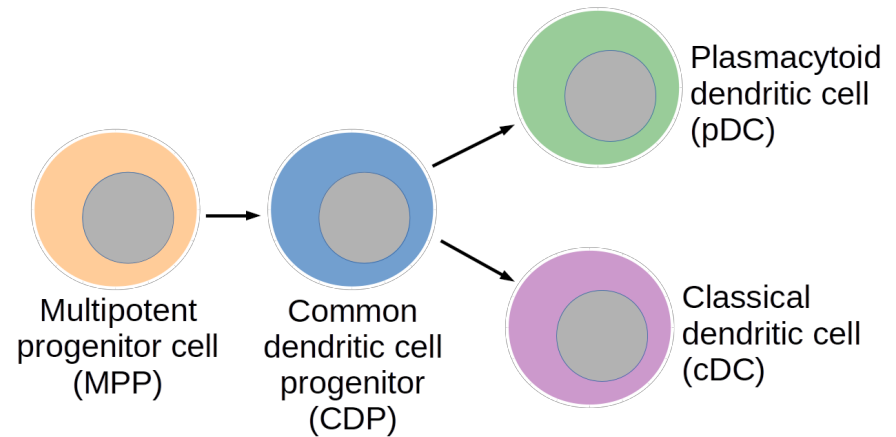
- cDC\_H3K4me1\_chr8.bam
- cDC\_H3K4me1\_chr8\_peak.bed
- cDC\_PU1\_chr8.bam
- cDC\_PU1\_chr8\_peaks.bed
- script.zsh (only for reference)

Convert BAM into TDF using *count*

Explore the data

# Example: Epigenetic Changes in Cell Differentiation

*invitro* system for mimicking dendritic cell differentiation



PU.1

Find binding sites of the transcription factor PU.1



PU.1 Motif

H3K4me3

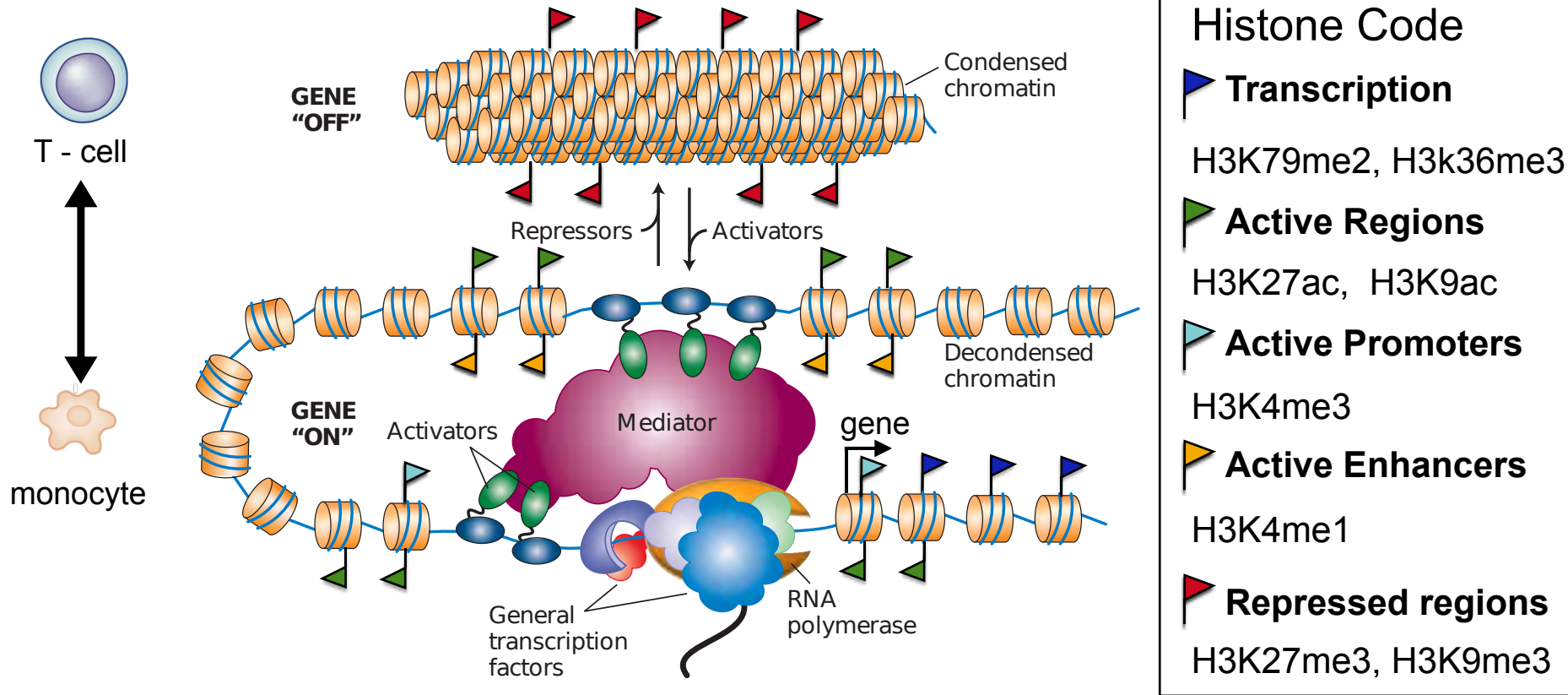
Find histone modifications associated to PU.1 binding sites

H3K4me1

H3K27me3

ChIP-Seq

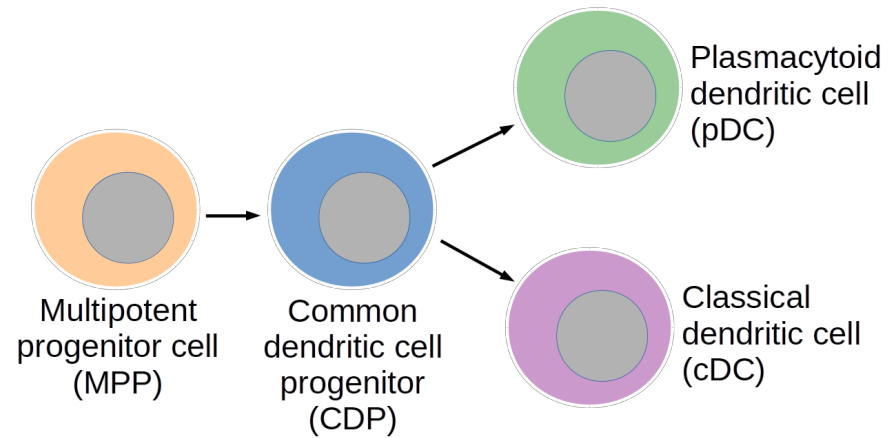
# Chromatin, Histone Code and TF binding





# Example: Epigenetic Changes in Cell Differentiation

*invitro* system for mimicking dendritic cell differentiation



PU.1

Find binding sites of the transcription factor PU.1



PU.1 Motif

ChIP-Seq

H3K4me3

Active Promoters

H3K4me1

Active Enhancers

H3K27me3

Repressed regions

Modifications  
PU.1 binding sites



[www.costalab.org](http://www.costalab.org)

Institute for  
Computational Genomics  
01011011010  
1010010010

**RWTH**AACHEN  
UNIVERSITY