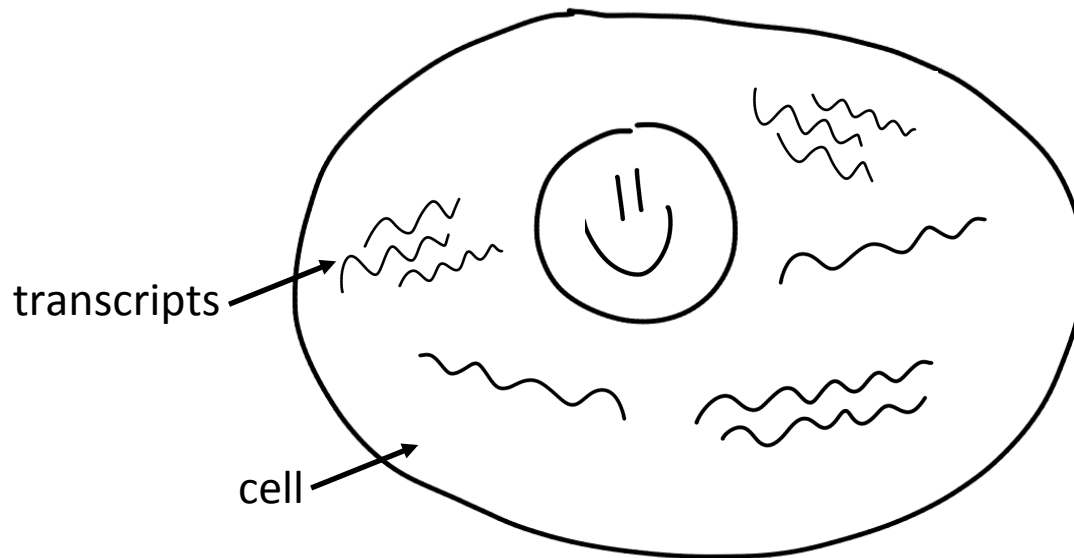


RNA-seq: What? Why? How?

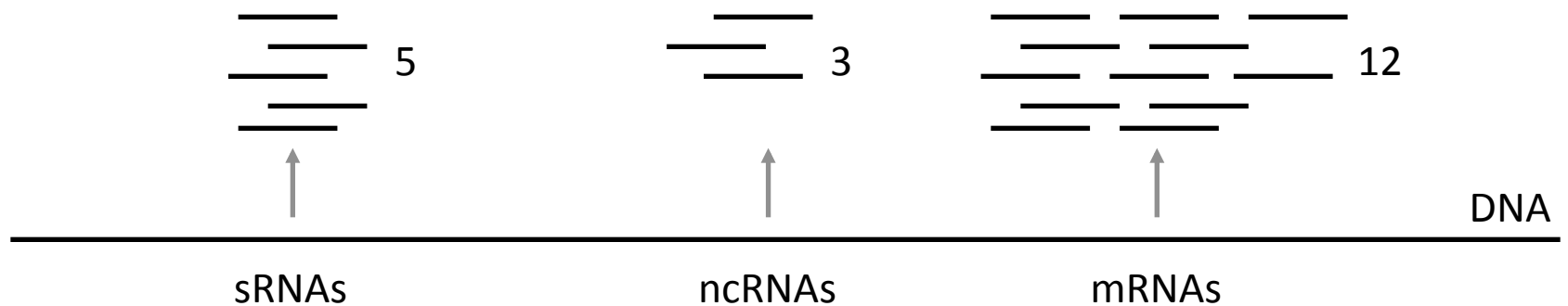
Dariya Sydykova

What is RNA-seq?

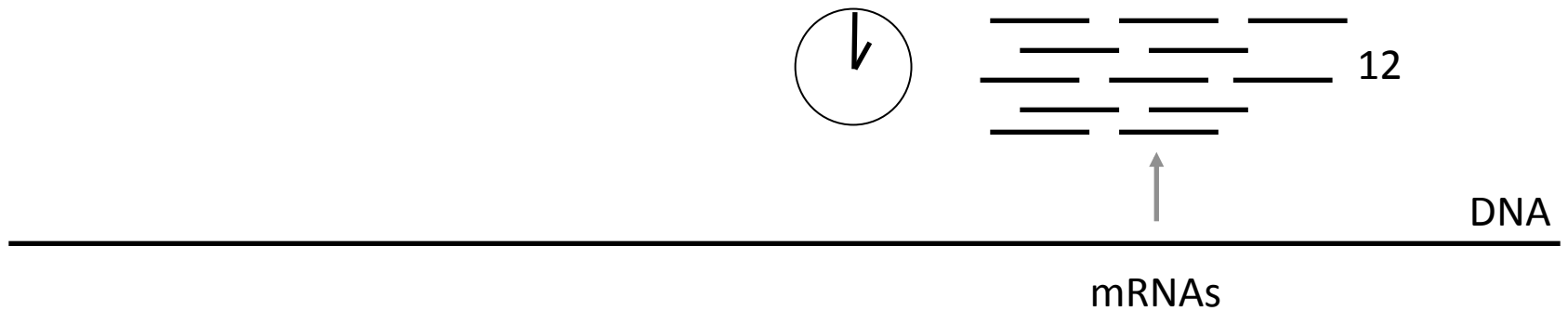
RNA-seq is a deep-sequencing technology that captures and quantifies the complete set of transcripts in a cell or the cell's transcriptome.



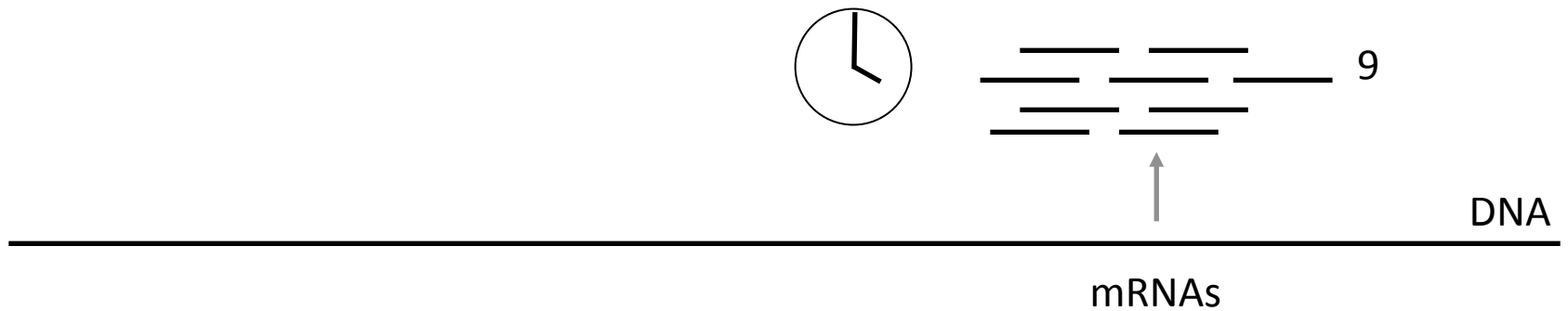
Why do we do RNA-seq?



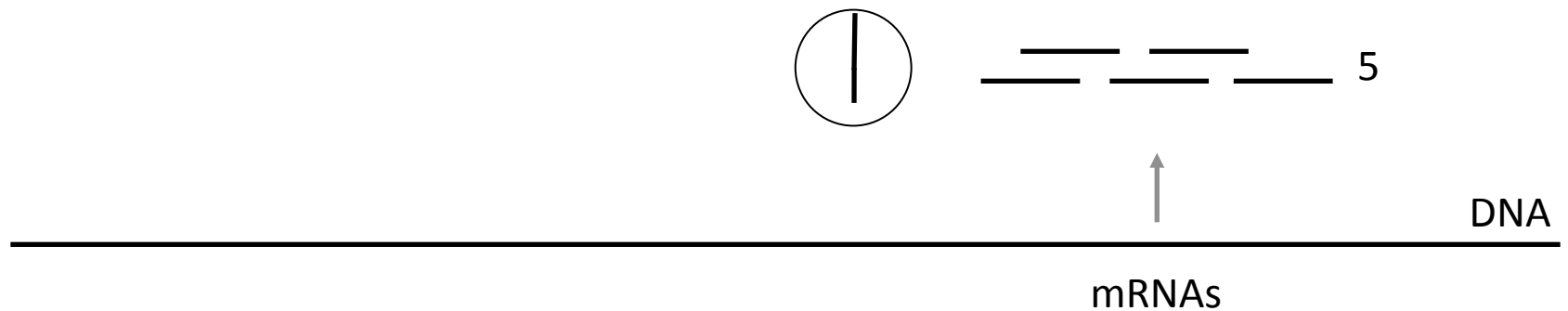
Why do we do RNA-seq?



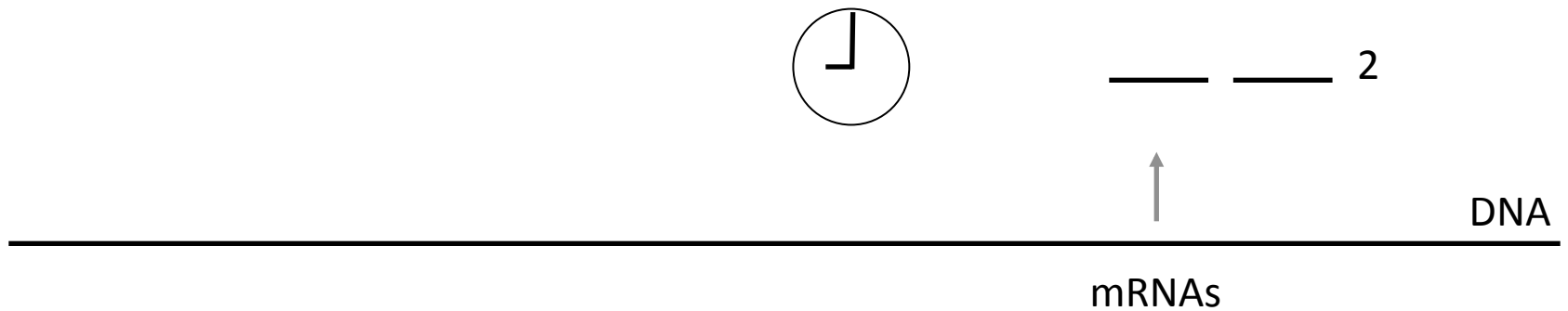
Why do we do RNA-seq?



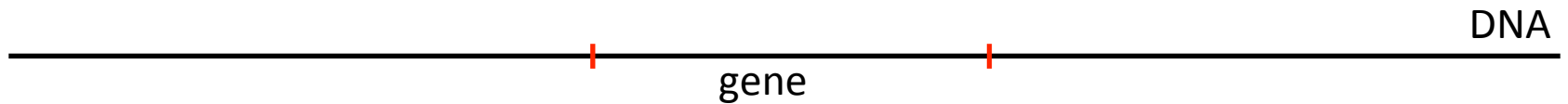
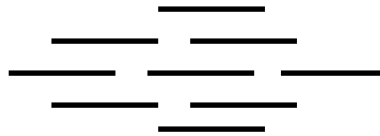
Why do we do RNA-seq?



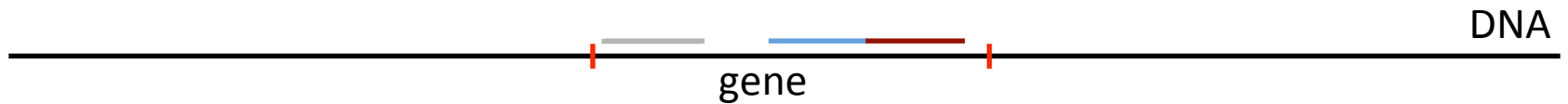
Why do we do RNA-seq?



Why do we do RNA-seq?

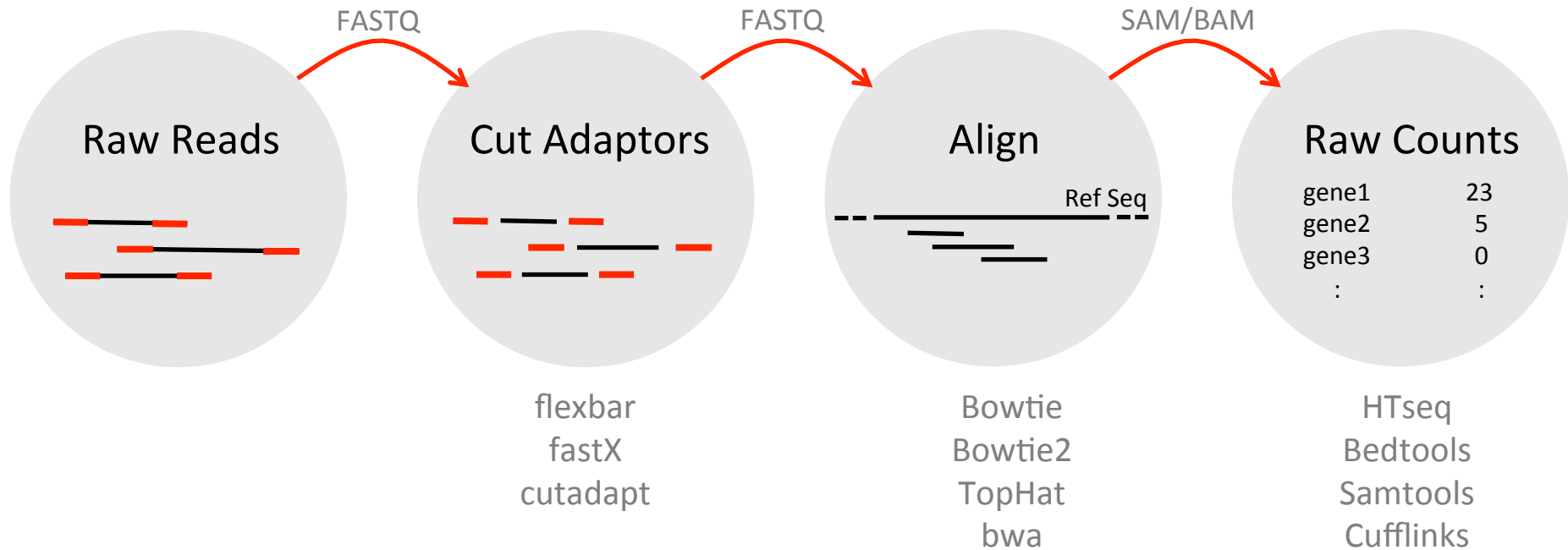


Why do we do RNA-seq?

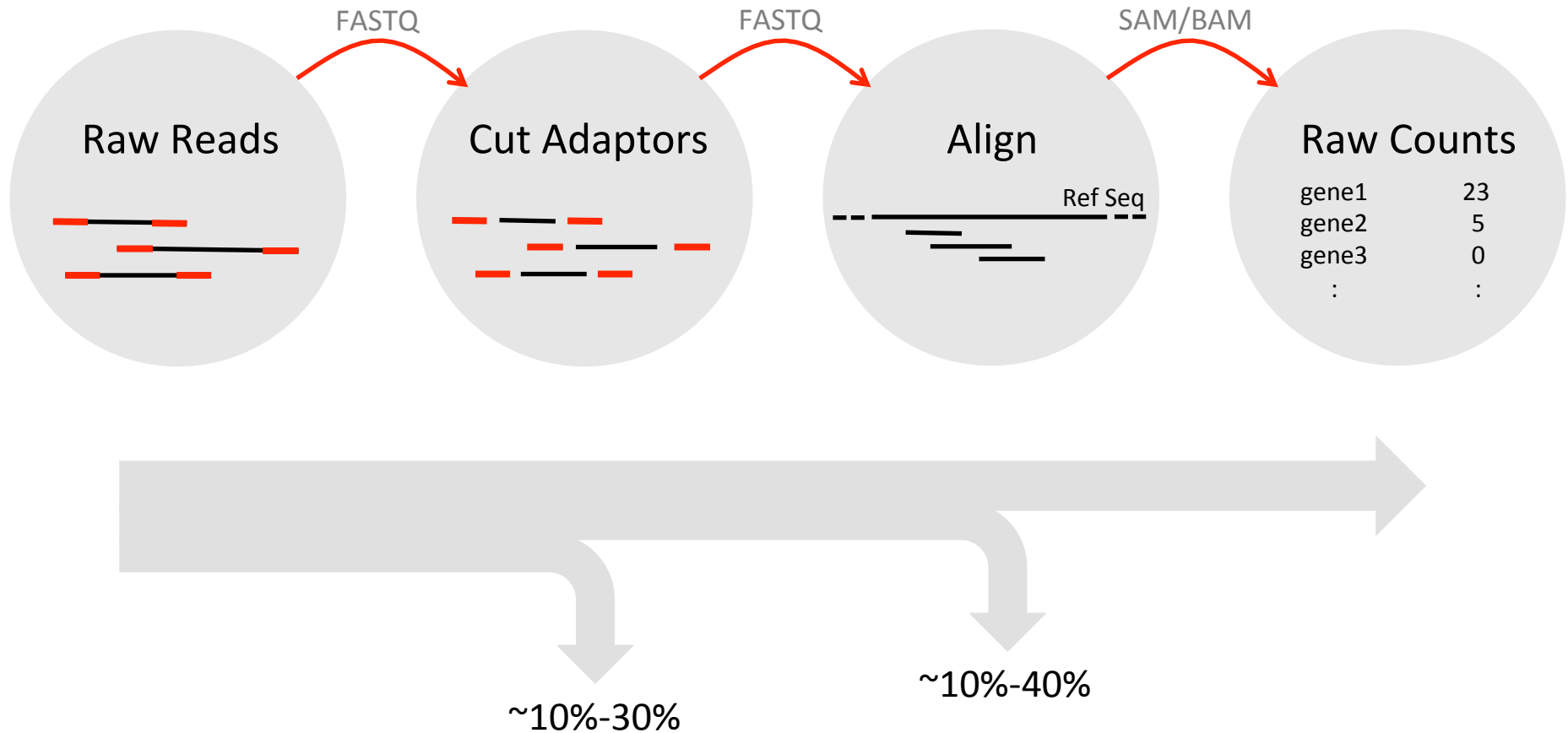


Next Gen Sequencing

RNA-seq pipeline



RNA-seq pipeline



FASTQ file



Read identifier:

Machine name run ID flowcell ID x y 1/2:Y/N:0:index seq

@HWI-ST1097:104:D13TNACXX:4:1101:1715:2142 1:N:0:CGATGT

seq → GCGTTGGTGGCATAGTGGTGAGCATAGCTGCCTTCCAAGCAGTTATGGGAG

+ ← GSAF/sequence description

→ =<@BDDD=A;+2C9F<CB?;CGGA<<ACEE*1?C:D>DE=FC*0BAG?DB6

└ Ascii-encoded base quality score

How to get the number of reads?

```
grep ^@ <file.fastq> | wc -l
```

```
wc -l <file.fastq>
```

```
Daria-Sydykovas-MacBook-Pro:Downloads DariaSydykova$ grep ^@ sample_1.fq | wc -l
750000
Daria-Sydykovas-MacBook-Pro:Downloads DariaSydykova$ wc -l sample_1.fq
3000000 sample_1.fq
```

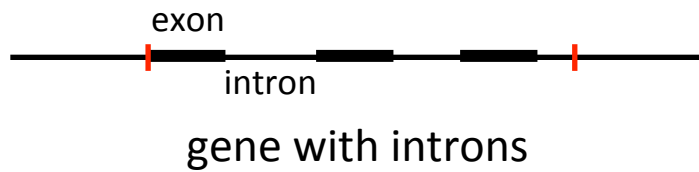
How to cut off the adaptors?

GCGTTGGTGGCATAGTGGTGAGCATAGCTGCCTTCCAAGCAGTTATGGGAG

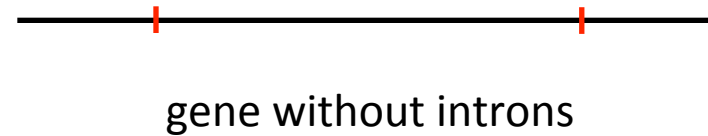
```
flexbar -n 1 -t <new_file_name> -r <reads1.fastq>  
-p <reads2.fastq> -a adaptors.fasta -f <fastq/fasta> > flexbar.out
```

Mapping

TopHat/Bowtie2

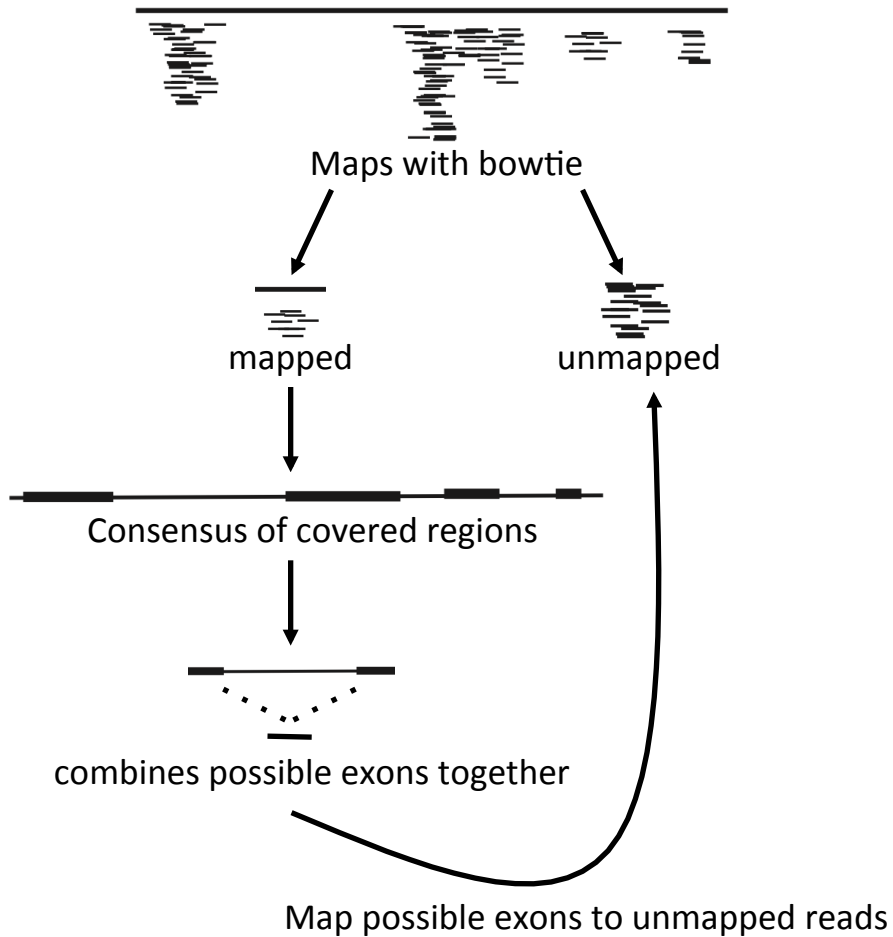


Bowtie/Bowtie2

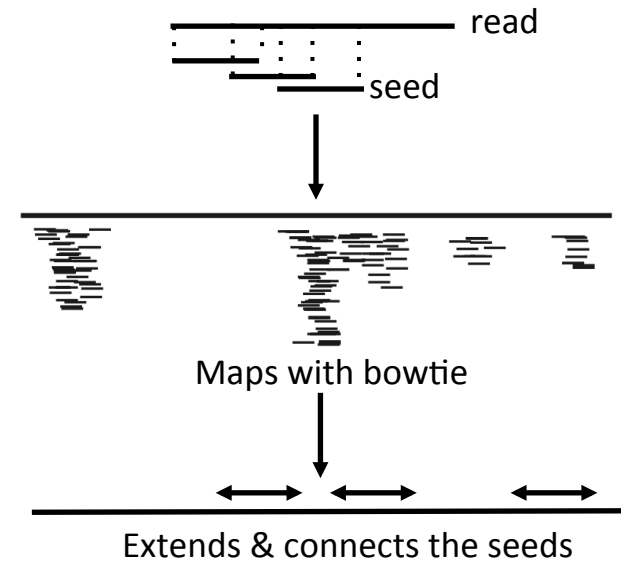


Mapping

TopHat



Bowtie2



Mapping

```
tophat -p 4 -G <gtf_ref_seq_file> -o <output.sam>  
--no-novel-juncs <reference_index> <reads1.fastq>  
<reads2.fastq>
```

```
bowtie2 -x <reference_index> -1 <reads1.fastq> -2  
<reads2.fastq> -S <output.sam>
```

GTF file

1	REL606	.	exon	190	255	.	+	.	transcript_id "ECB_00001";
2	REL606	.	gene	190	255	.	+	.	transcript_id "ECB_00001"; gene_id "ECB_00001"; gene_name "thrL";
3	REL606	.	exon	336	2798	.	+	.	transcript_id "ECB_00002";
4	REL606	.	gene	336	2798	.	+	.	transcript_id "ECB_00002"; gene_id "ECB_00002"; gene_name "thrA";
5	REL606	.	exon	2800	3732	.	+	.	transcript_id "ECB_00003";
6	REL606	.	gene	2800	3732	.	+	.	transcript_id "ECB_00003"; gene_id "ECB_00003"; gene_name "thrB";
7	REL606	.	exon	3733	5019	.	+	.	transcript_id "ECB_00004";
8	REL606	.	gene	3733	5019	.	+	.	transcript_id "ECB_00004"; gene_id "ECB_00004"; gene_name "thrC";
9	REL606	.	exon	5232	5528	.	+	.	transcript_id "ECB_00005";
10	REL606	.	gene	5232	5528	.	+	.	transcript_id "ECB_00005"; gene_id "ECB_00005"; gene_name "yaaX";
11	REL606	.	exon	5681	6457	.	-	.	transcript_id "ECB_00006";
12	REL606	.	gene	5681	6457	.	-	.	transcript_id "ECB_00006"; gene_id "ECB_00006"; gene_name "yaaA";
13	REL606	.	exon	6527	7957	.	-	.	transcript_id "ECB_00007";
14	REL606	.	gene	6527	7957	.	-	.	transcript_id "ECB_00007"; gene_id "ECB_00007"; gene_name "yaaJ";
15	REL606	.	exon	8236	9189	.	+	.	transcript_id "ECB_00008";
16	REL606	.	gene	8236	9189	.	+	.	transcript_id "ECB_00008"; gene_id "ECB_00008"; gene_name "talB";
17	REL606	.	exon	9304	9891	.	+	.	transcript_id "ECB_00009";
18	REL606	.	gene	9304	9891	.	+	.	transcript_id "ECB_00009"; gene_id "ECB_00009"; gene_name "mogA";
19	REL606	.	exon	9926	10492	.	-	.	transcript_id "ECB_00010";
20	REL606	.	gene	9926	10492	.	-	.	transcript_id "ECB_00010"; gene_id "ECB_00010"; gene_name "yaaH";
21	REL606	.	exon	10641	11354	.	-	.	transcript_id "ECB_00011";
22	REL606	.	gene	10641	11354	.	-	.	transcript_id "ECB_00011"; gene_id "ECB_00011"; gene_name "yaaW";
23	REL606	.	exon	10828	11313	.	+	.	transcript_id "ECB_00012";
24	REL606	.	gene	10828	11313	.	+	.	transcript_id "ECB_00012"; gene_id "ECB_00012"; gene_name "htgA";

Mapping Output: SAM/BAM file

Mate name	Mate position	Template length
--------------	------------------	--------------------

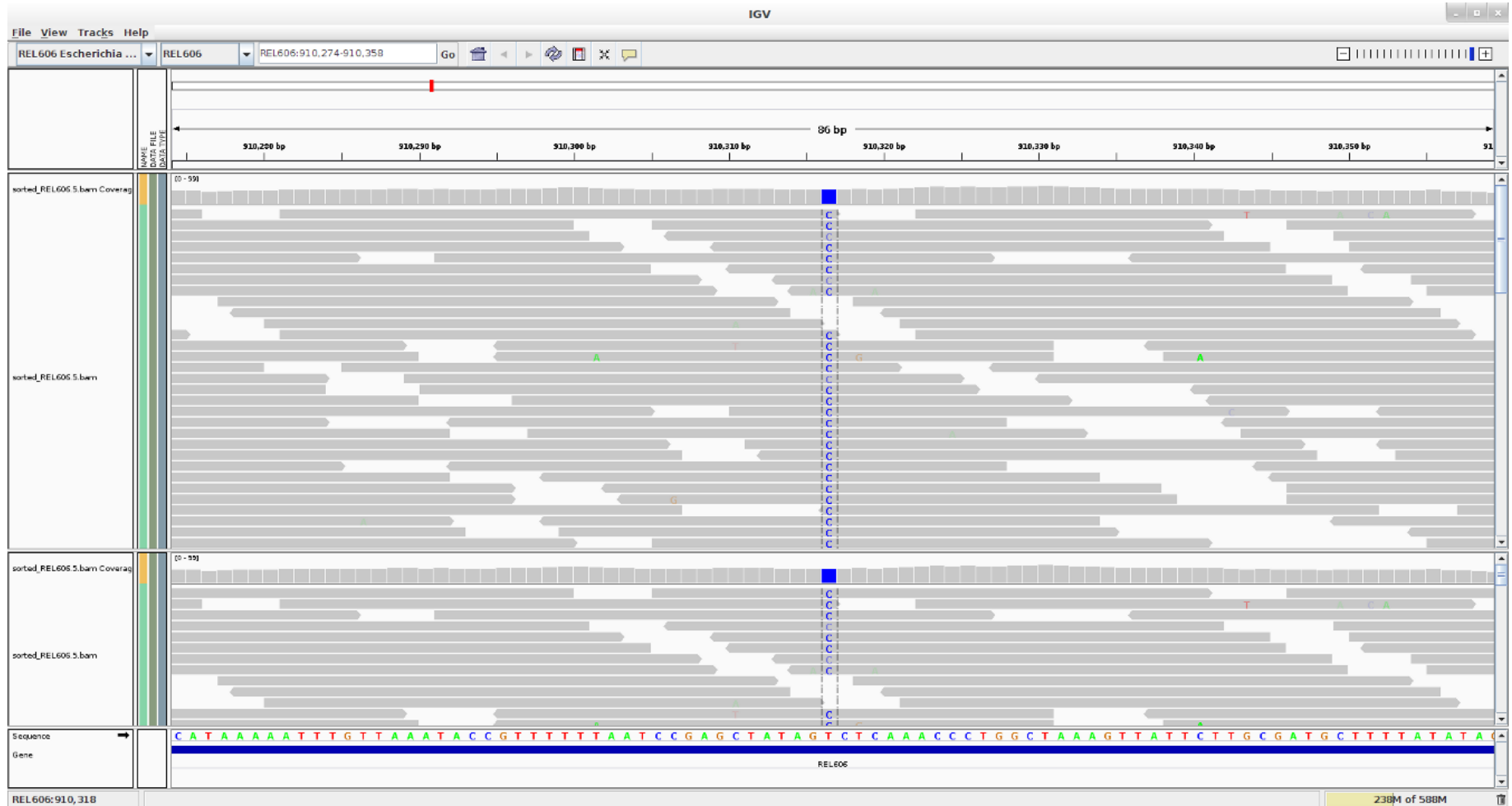
Read identifier	FLAG	Ref Seq	position	-10 log ₁₀ P(mapping pos is wrong)
HWI-ST1001:137:C12FPA	0	chr1	12805	1

42M4I5M * 0 0

TTGGATGCCCCTCCACACCCTCTTGATCTTCCCTGTGATGTCACCAATATG seq
CCCCFFFFHHGHJJJJHJJJJJJJJJJJJJJJJJJJJJJJJJJ AS:i:-28 XN:i:0 XM:i:2 XO:i:
1XG:i:4 NM:i:6 MD:Z:2C41C2 YT:Z:UU NH:i:3 CC:Z:chr15 CP:i:102518319
XS:A:+ HI:i:0

Visualizing SAM files

Integrative genomics viewer (IGV):



Visualizing SAM files

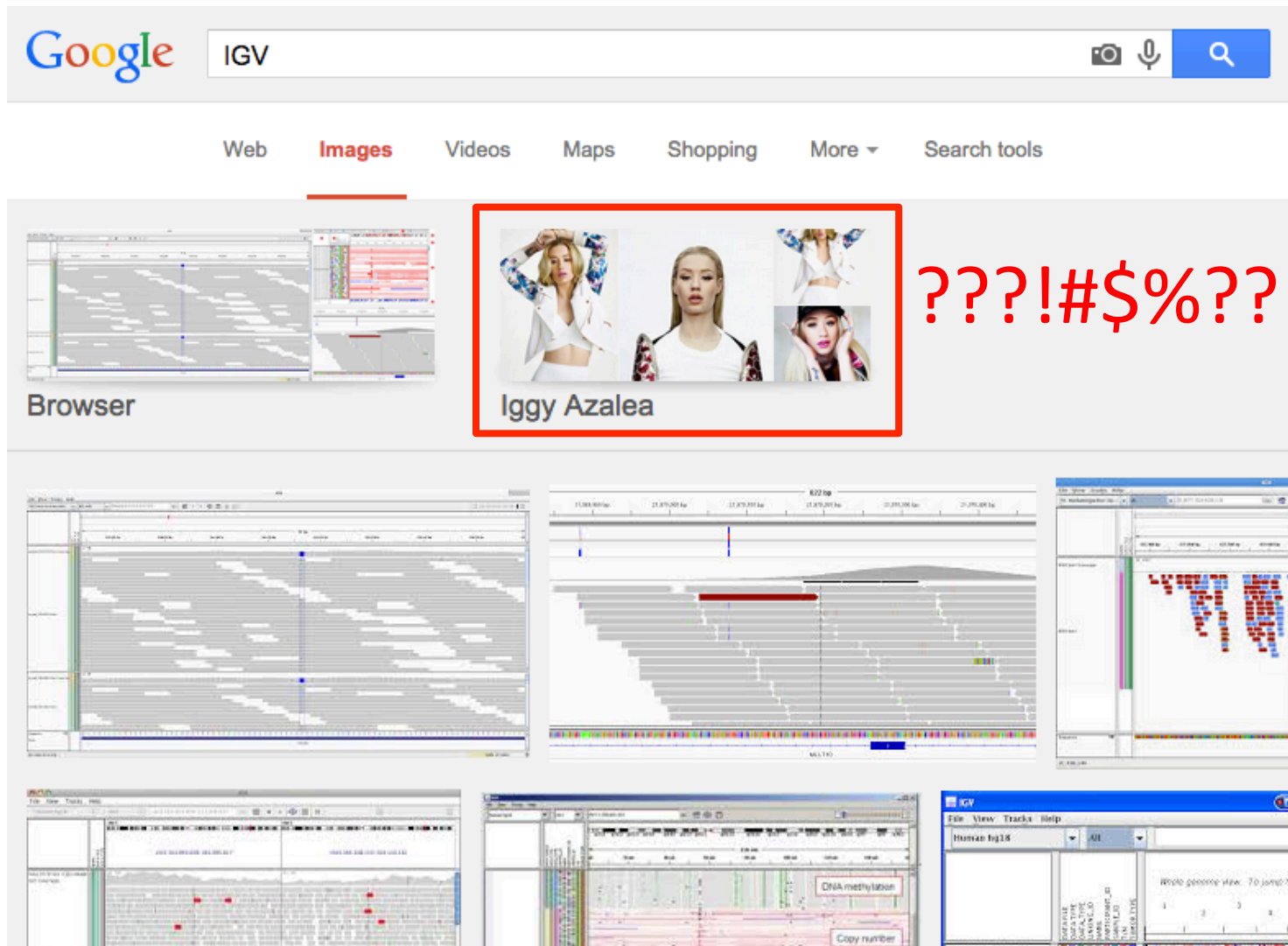
Google IGV

Web Images Videos Maps Shopping More Search tools

Browser

Iggy Azalea

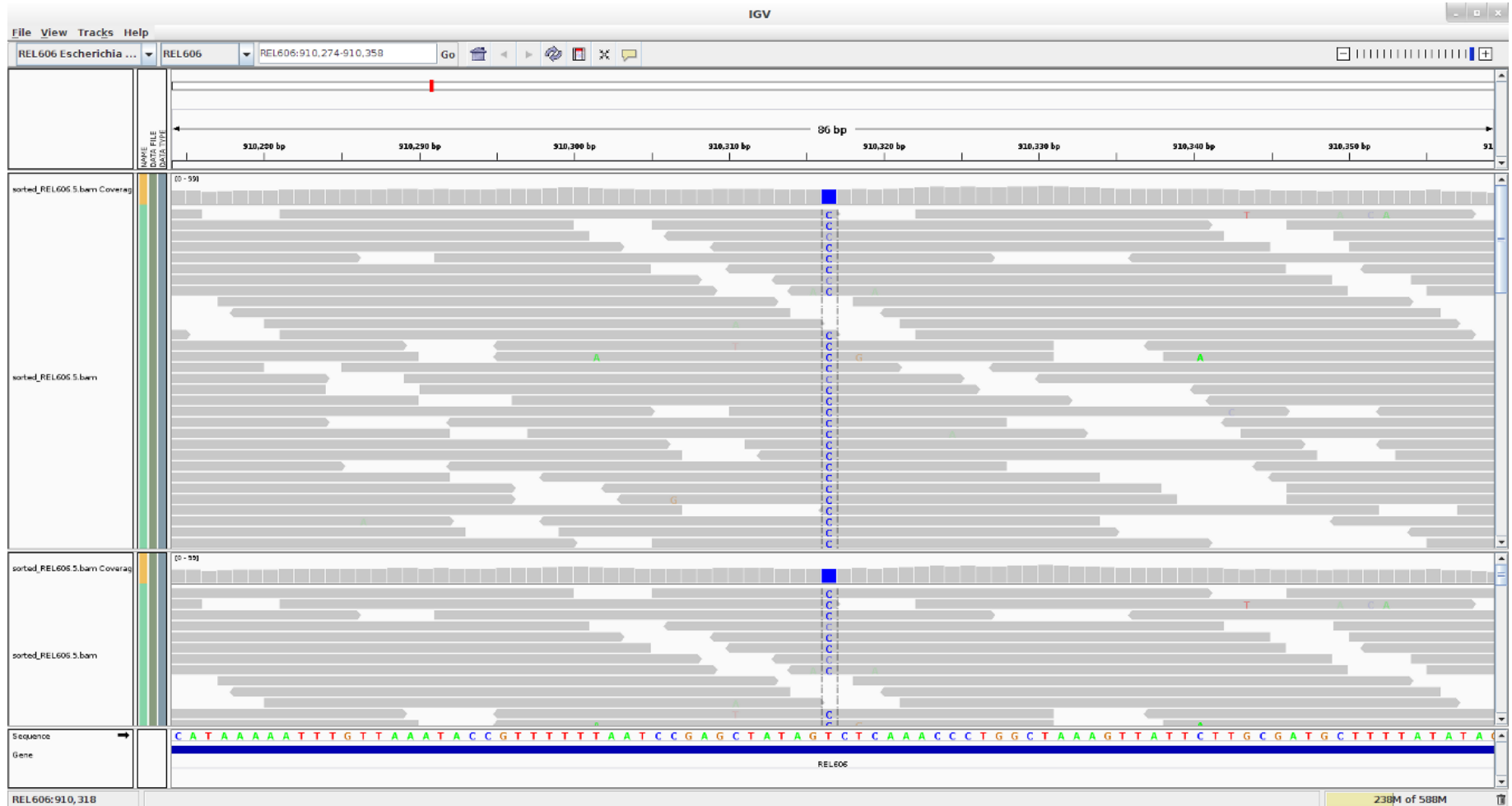
???!#\$%??!?



The image is a collage of screenshots from the IGV (Integrative Genomics Viewer) software. The top left shows a 'Browser' view with a track of genomic data. The top center features a red-bordered box containing three images of Iggy Azalea, with the text 'Iggy Azalea' below them. To the right of this box is the text '???!#\$%??!?' in red. The bottom section contains several smaller screenshots showing different views of genomic data, including tracks for 'DNA methylation' and 'Copy number'.

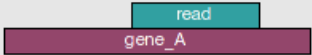
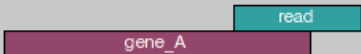


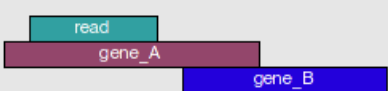
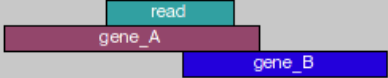
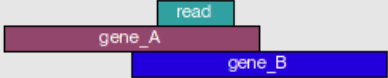
Visualizing SAM files

Integrative genomics viewer (IGV):



Quantifying Gene Counts

htseq-count -m <mode> -t <feature_type> -i <id_attribute> <output.sam>
<gtf_ref_seq_file>

	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

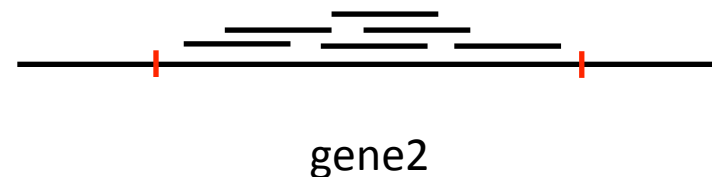
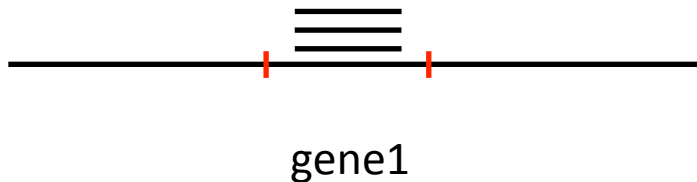
sorted
tools sort

Normalizing transcript counts

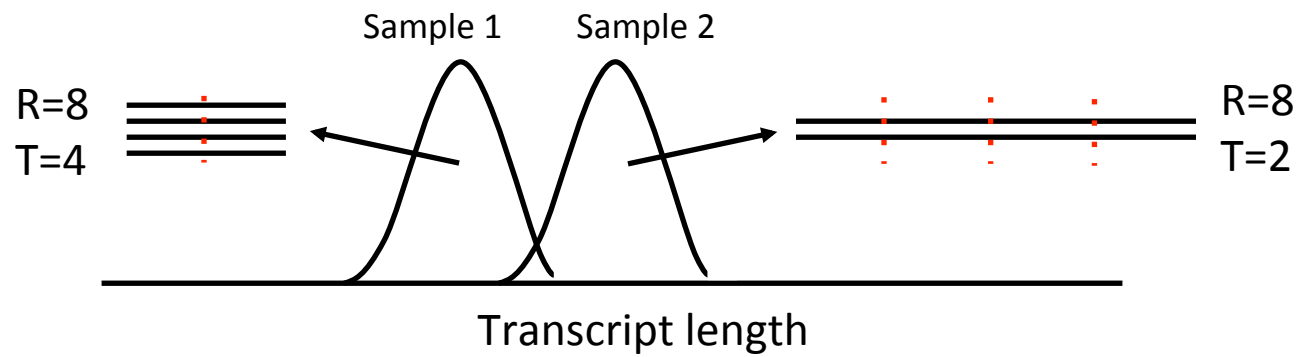
Reads (fragments) per kilobase of exon per million fragments mapped (RPKM or FPKM) normalizes for:

- Total number of reads
- Gene lengths

$$\text{RPKM} = \frac{\text{\# reads mapped to a gene}}{\text{exon length}/10^3 * \text{total \# reads mapped}/10^6}$$



Downside of RPKM



DESeq2

	sample1	sample2
gene1	4	2
gene2	10	5
Size Factor:	20	10

↓ Divide by size factor

	sample1	sample2
gene1	.2	.2
gene2	.5	.5

Questions?

TPM

Transcripts per million (TPM) normalizes:

- Total number of transcripts
- Gene lengths
- Read lengths/gene coverage

$$\text{TPM} = \frac{\text{\# reads mapped to a gene} * \text{mean read length}}{\text{exon length} * T/10^6}$$

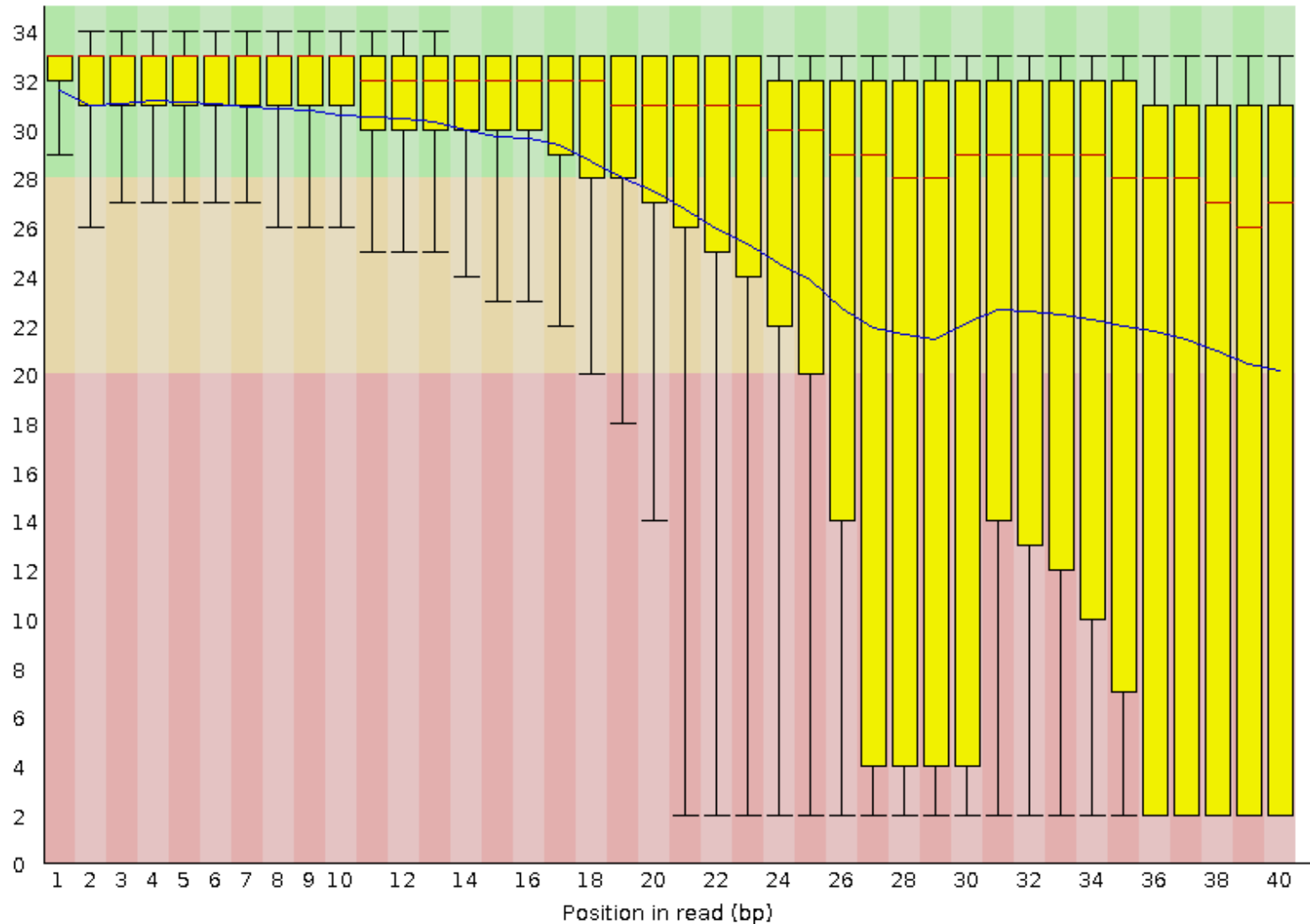
$$T = \sum \frac{\text{\# reads mapped to a gene} * \text{mean read length}}{\text{exon length}}$$

What is the sequence quality of your reads?

```
fastqc <file.fastq>
```

Bad sequence quality

Quality scores across all bases (Illumina 1.5 encoding)



Good sequence quality

