# NGS read mapping : answers to questions

Céline Keime
keime@igbmc.fr

# Exercise 1
# 1. Log file

Proportion of uniquely mapped reads :

# Exercise 1
## 2. Alignment file

- **Galaxy**
    - STAR provides an alignment in BAM format
    - Download this file together with the corresponding index (in the same directory)

- **IGV**
    - File → Load from file and choose the downloaded BAM file
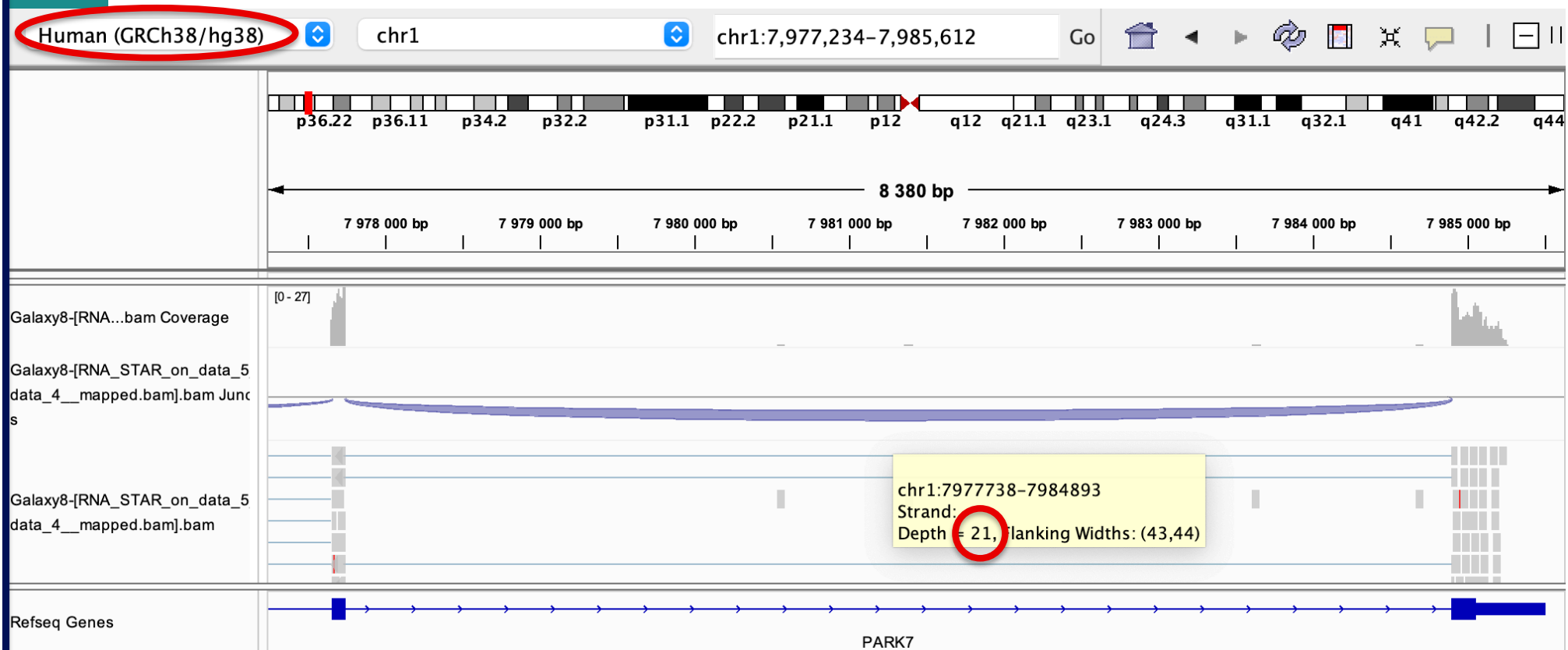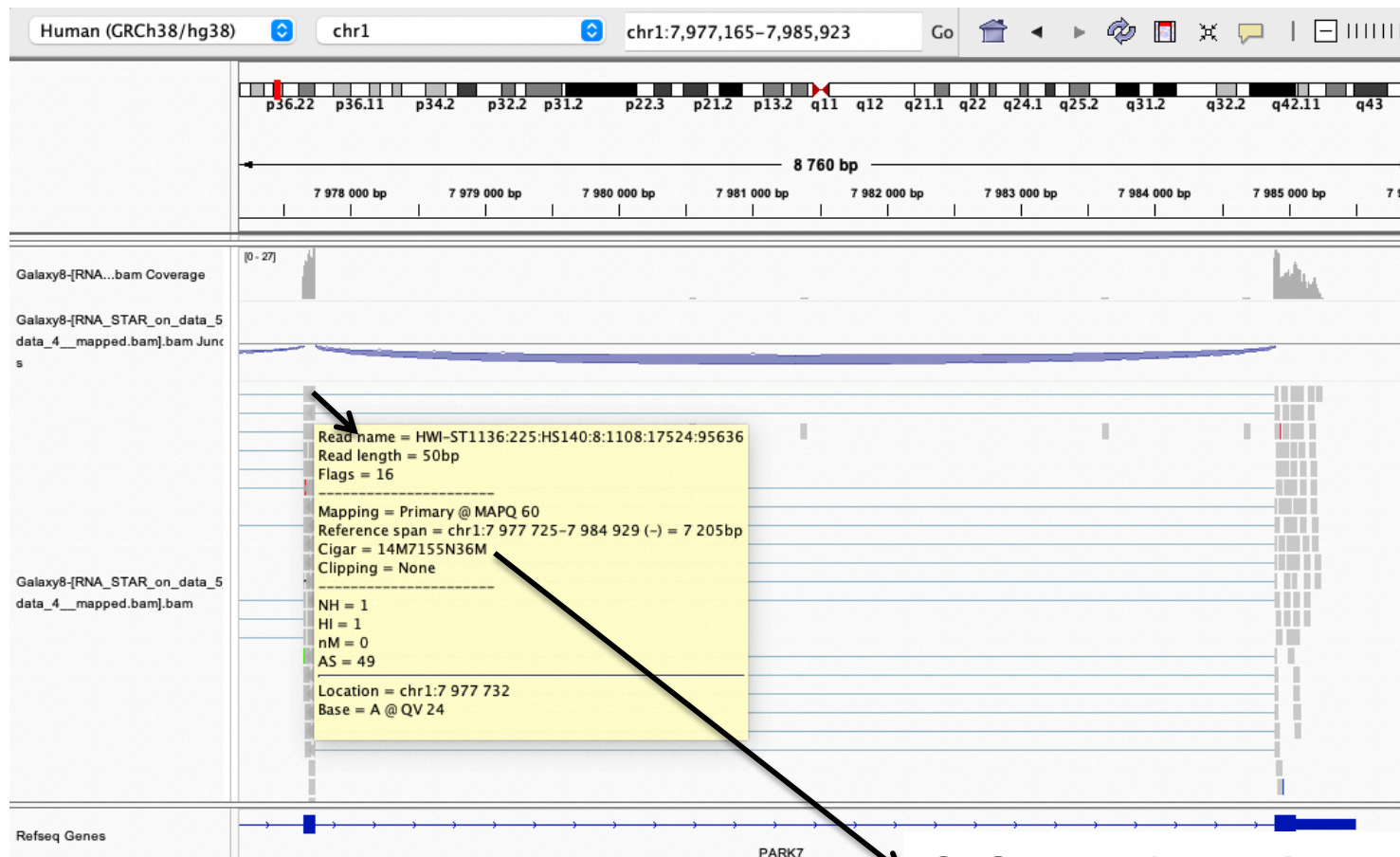
# Exercise 1
## 2. Splice junction



→ 21 reads span the junction that joins the last 2 exons of *Park7* gene

# Exercise 1
## 2. Splice junction
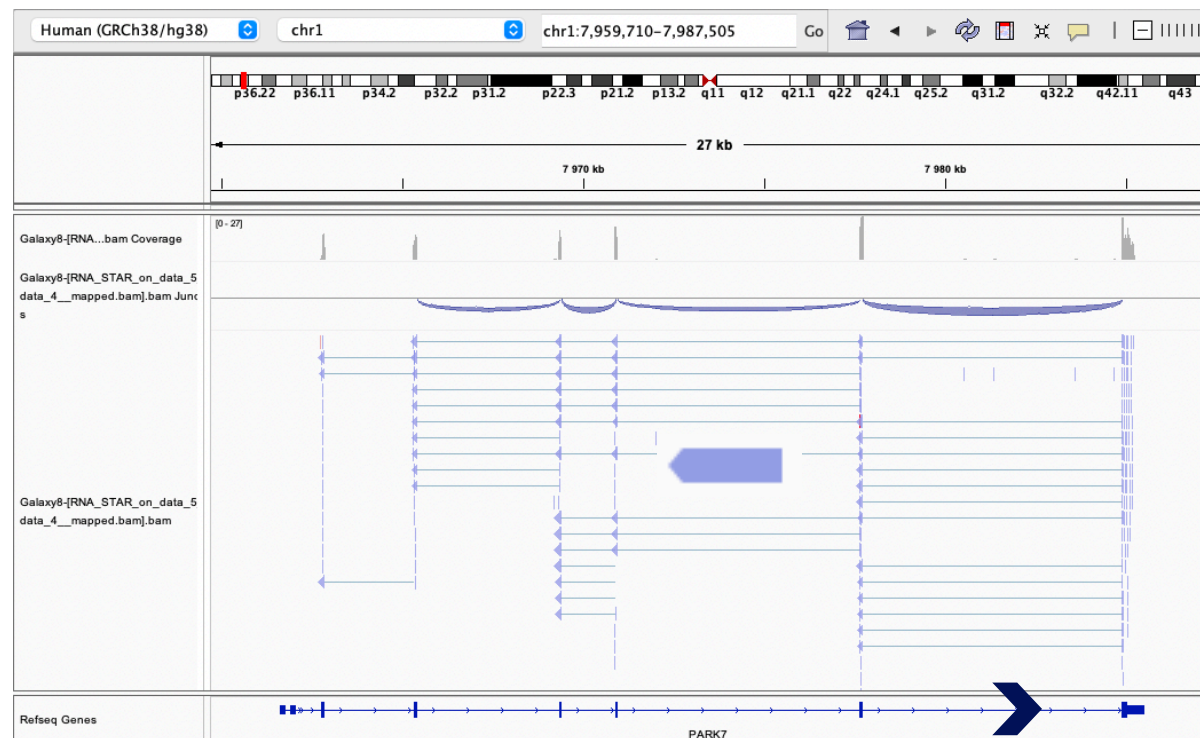


CIGAR : 14M**7155N**36M

Intron length :
7984893 - 7977738 = 7155

# Exercise 1
# 2. Strand specificity

Right click on BAM file → Color alignments by → read strand

*Park7 :*



The library has been prepared with a directional mRNAseq protocol which retains strand information :
reads are in the opposite direction as the transcribed strand

# Exercise 1
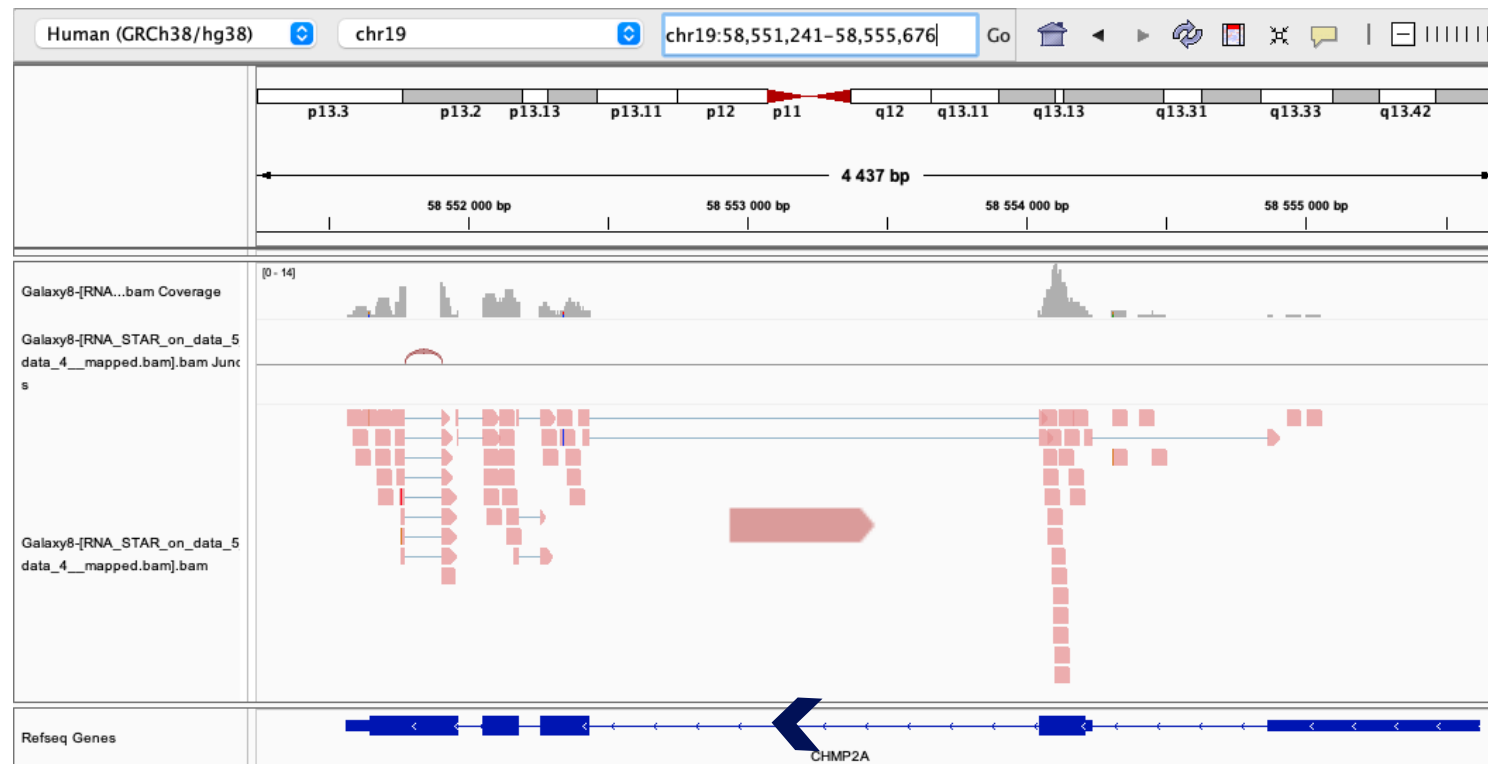## 2. Strand specificity

*Chmp2a :*



The library has been prepared with a directional mRNAseq protocol which retains strand information :
reads are in the opposite direction as the transcribed strand

# Exercise 1
# 2. Multiple mapped reads

Right click on BAM file → Color alignments by → tag → NH



Number of reported alignments
→ see NH tag in pop-up windows to visualize
color-coding (that can be different from this one) :  ▇ 1   ▇ 2   ▇ 3

There are multiple aligned reads on this gene

# Exercise 2 - Question 1
# Proportion of uniquely mapped reads

Galaxy : "NGS data analysis training Strasbourg" history



→ This proportion is consistent across samples

# Exercise 2 – Question 2
## *Idh1* gene expression

IGV : File → Load from file and select the 4 tdf files

Select all tdf tracks → Right-click → Group Autoscale :

→ IGV automatically adjusts the Y scale to the data range currently in view (this scaling continually adjusts as you move)

→ all tracks are on the same scale

Search for *Idh1*



*Idh1* is under-expressed in siMitf samples compared to siLuc ones

# Exercise 2 – Question 3

- File → new session
- File → load from files and load the 4 BAM files
- Search for *EEF2*

# Exercise 2 – Question 3

Exon numbers are provided on annotation track

Click and drag on a region to zoom in

# Exercise 2 – Question 3

- *Eef2* exon 11
  - chr19:3,979,410 : G in ~100% of the reads, A in the genome

# Exercise 2 – Question 3

- *Eef2* exon 13
  - chr19:3,977,488 : G in ~100% of the reads, A in the genome

# Exercise 2 – Question 4

- Position chr4:6707961 :
  - Deletion vs reference genome

# Exercise 2 – Question 5

- Region chr20:44,935,294-44,939,521 :
  - Right-click on Refseq Genes track → select Expanded
    to see all annotated isoforms

# Exercise 2 – Question 5

- Region chr20:44,935,294-44,939,521 :



We detect an isoform without this exon in siMitf samples

**IGV is only a visualization tool**
**In-depth analysis using paired-end data with more coverage is needed**

# Exercise 2 – Question 5

- If you would like to display Ensembl annotations, you can add this track
    - File → Load from file
    - Select Homo_sapiens.GRCh38.105.chr.sorted.gtf available in RNAseq/annotations folder

# Exercise 2 – Question 5

- You can save your IGV session
  - To save the current state of your IGV session to a named session file
  - File → Save Session
  - Data files must stay at the same location
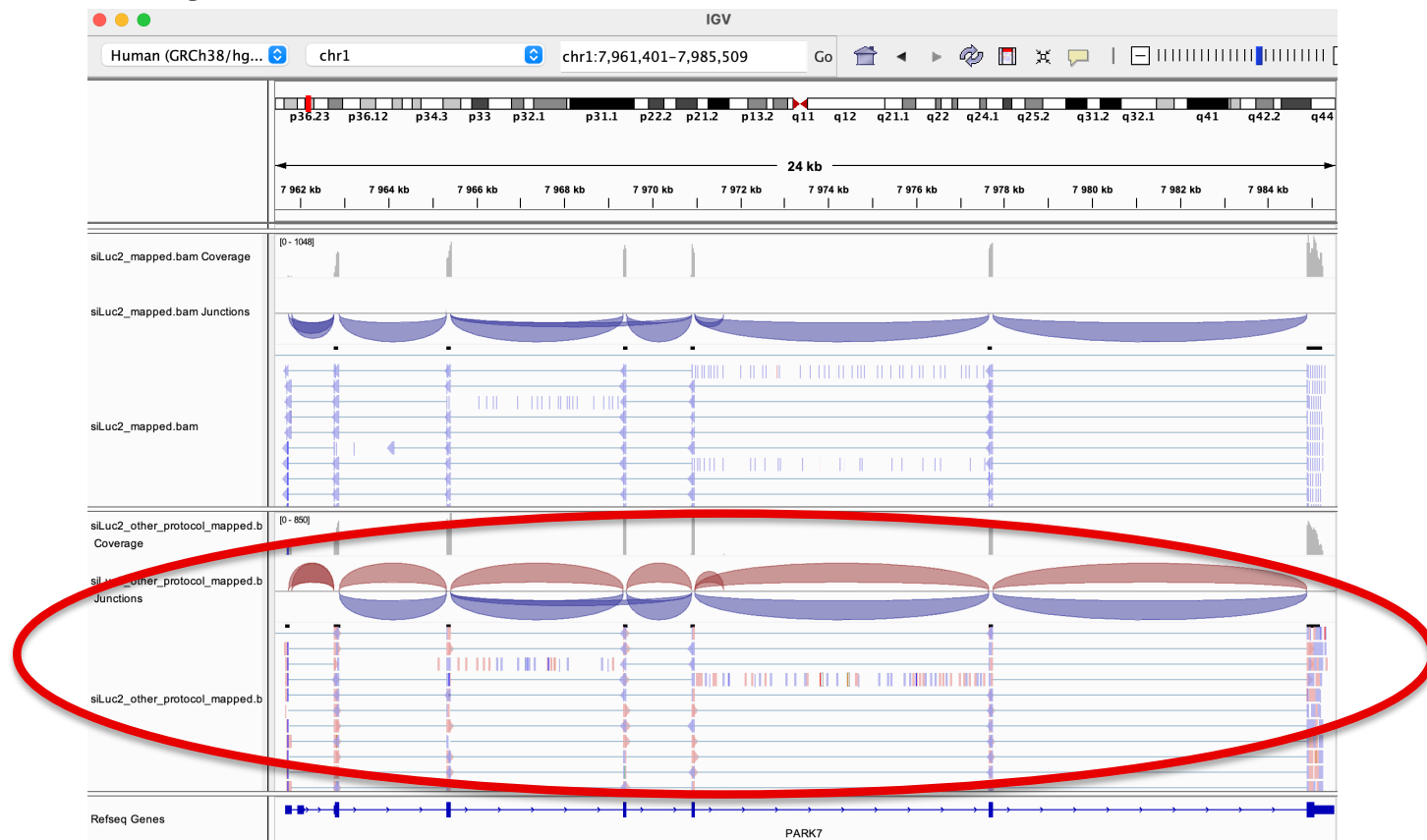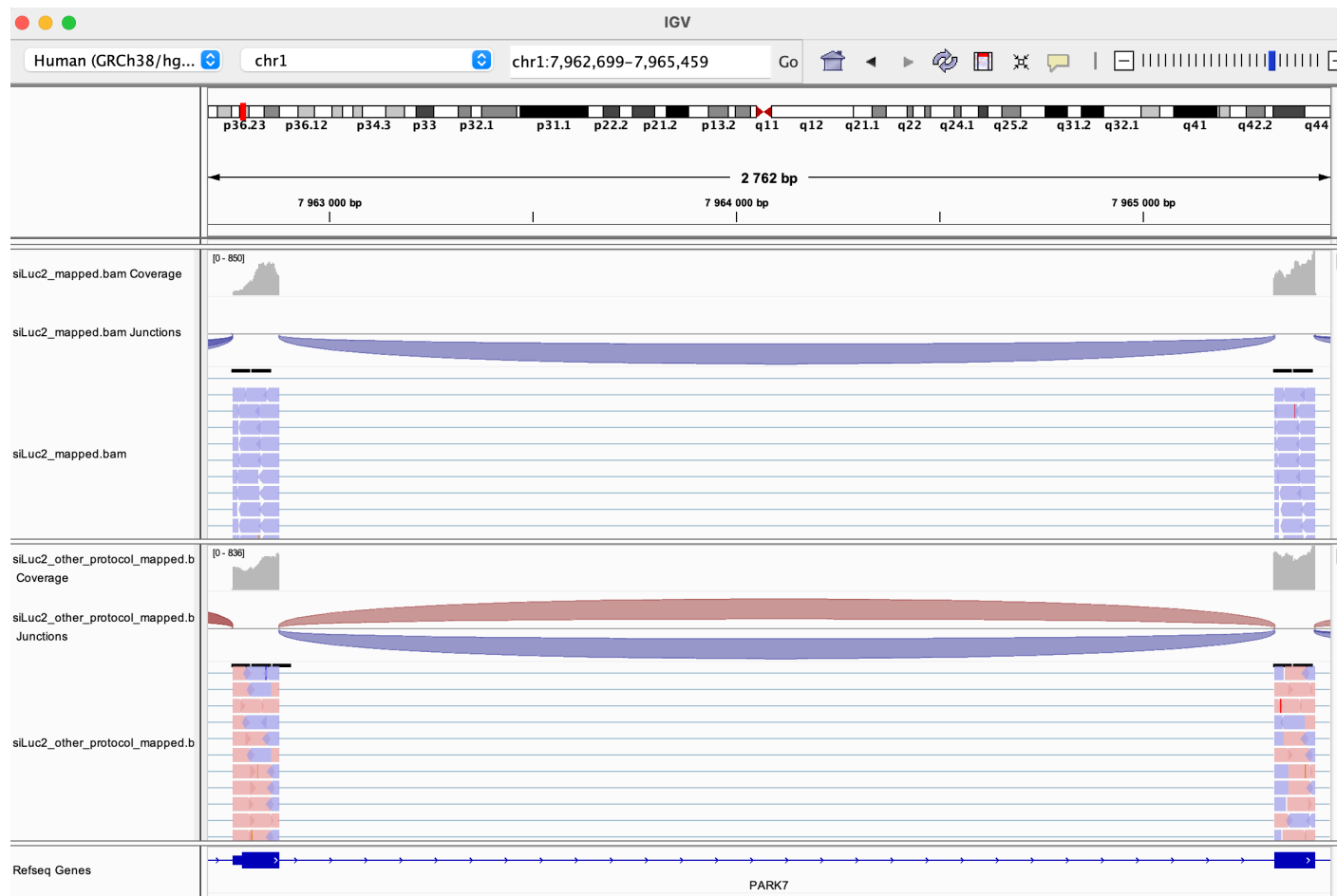- Use File → Open session to restore a saved session

# Exercise 2 – Question 6

- Remove siLuc3 and siMitf3/4 tracks (Right click on tracks → Remove track)
- File → load from file and select siLuc2_other_protocol_alignment.bam
- Right-click on BAM file → Color alignments by → read strand
- e.g. *Park7* gene

# Exercise 2 – Question 6



This protocol is not directional (it does not preserve strand information)