# NGS read mapping : answers to questions

Céline Keime
keime@igbmc.fr

# Exercise 1
# 1. Alignment summary statistics

```
Reads:
        Input     :    1000000
        Mapped    :     983595 (98.4% of input)
          of these:      88965 ( 9.0%) have multiple alignments (434 have >20)
98.4% overall read mapping rate.
```

History

search datasets

**RNAseq1709**
7 shown

207.09 MB

7: TopHat 2 on data 2
and data 1:
accepted_hits

6: TopHat 2 on data 2
and data 1: splice
junctions

5: TopHat 2 on data 2
and data 1: deletions

4: TopHat 2 on data 2
and data 1: insertions

3: TopHat 2 on data 2
and data 1:
align_summary

1.1. 983,595 reads mapped onto hg38
1.2. 9% of these reads have multiple alignments

# Exercise 1
# 2. Splice junctions

| Chrom | Start | End | Name | Score | Strand | ThickStart | ThickEnd | ItemRGB | BlockCount | BlockSizes | BlockStarts |
|-------|-------|-----|------|-------|--------|------------|----------|---------|------------|------------|-------------|
| track name=junctions description="TopHat junctions" | | | | | | | | | | | (relative to *chromStart)* |
| chr1 | 15015 | 15822 | JUNC00000001 | 1 | − | 15015 | 15822 | 255,0,0 | 2 | 23,27 | |
| chr1 | 15931 | 16640 | JUNC00000002 | 1 | − | 15931 | 16640 | 255,0,0 | 2 | 16,34 | 0,675 |
| chr1 | 16751 | 16902 | JUNC00000003 | 4 | − | 16751 | 16902 | 255,0,0 | 2 | 14,45 | 0,106 |
| chr1 | 17359 | 17646 | JUNC00000004 | 1 | − | 17359 | 17646 | 255,0,0 | 2 | 9,41 | 0,246 |
| chr1 | 17730 | 17952 | JUNC00000005 | 1 | − | 17730 | 17952 | 255,0,0 | 2 | 12,38 | 0,184 |
| chr1 | 18322 | 24758 | JUNC00000006 | 3 | − | 18322 | 24758 | 255,0,0 | 2 | 44,21 | 0,6415 |
| chr1 | 30656 | 31014 | JUNC00000007 | 1 | + | 30656 | 31014 | 255,0,0 | 2 | 11,39 | 0,319 |
| chr1 | 164755 | 165897 | JUNC00000008 | 1 | − | 164755 | 165897 | 255,0,0 | 2 | 36,14 | 0,1128 |
| chr1 | 185544 | 186351 | JUNC00000009 | 1 | − | 185544 | 186351 | 255,0,0 | 2 | 15,35 | 0,772 |
| chr1 | 186453 | 187162 | JUNC00000010 | 1 | − | 186453 | 187162 | 255,0,0 | 2 | 16,34 | 0,675 |
| chr1 | 187273 | 187424 | JUNC00000011 | 4 | − | 187273 | 187424 | 255,0,0 | 2 | 14,45 | 0,106 |
| chr1 | 188254 | 188476 | JUNC00000012 | 1 | − | 188254 | 188476 | 255,0,0 | 2 | 12,38 | 0,184 |
| chr1 | 188891 | 195308 | JUNC00000013 | 2 | − | 188891 | 195308 | 255,0,0 | 2 | 11,46 | 0,6371 |
| chr1 | 495032 | 497141 | JUNC00000014 | 1 | − | 495032 | 497141 | 255,0,0 | 2 | 17,33 | 0,2076 |
| chr1 | 733347 | 735455 | JUNC00000015 | 1 | − | 733347 | 735455 | 255,0,0 | 2 | 17,33 | 0,2075 |
| chr1 | 756108 | 758983 | JUNC00000016 | 1 | − | 756108 | 758983 | 255,0,0 | 2 | 33,17 | 0,2858 |
| chr1 | 765211 | 766342 | JUNC00000017 | 1 | − | 765211 | 766342 | 255,0,0 | 2 | 36,14 | 0,1117 |
| chr1 | 805866 | 808598 | JUNC00000018 | 1 | − | 805866 | 808598 | 255,0,0 | 2 | 25,25 | 0,2707 |

History

search datasets

RNAseq1709
7 shown

207.09 MB

7: TopHat 2 on data 2
and data 1:
accepted_hits

6: TopHat 2 on data 2
and data 1: splice
junctions

5: TopHat 2 on data 2
and data 1: deletions

4: TopHat 2 on data 2
and data 1: insertions

3: TopHat 2 on data 2
and data 1:
align_summary

Number of alignments
spanning the junction

Each junction consists of 2 connected BED blocks →
Each block is as long as the maximal overhang of any
read spanning the junction

## 2.1. Splice junctions provided in a BED file

# Exercise 1
## 2.2. Splice junctions visualization
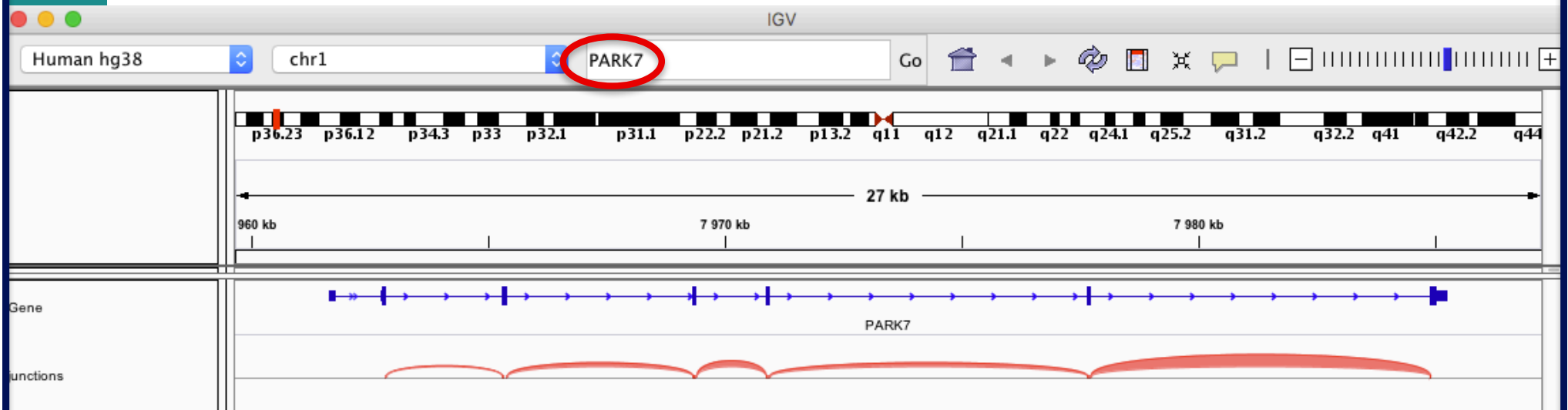
- **Galaxy**
  - Download splice junctions BED file



- **IGV**
  - Select the appropriate genome assembly (hg38)
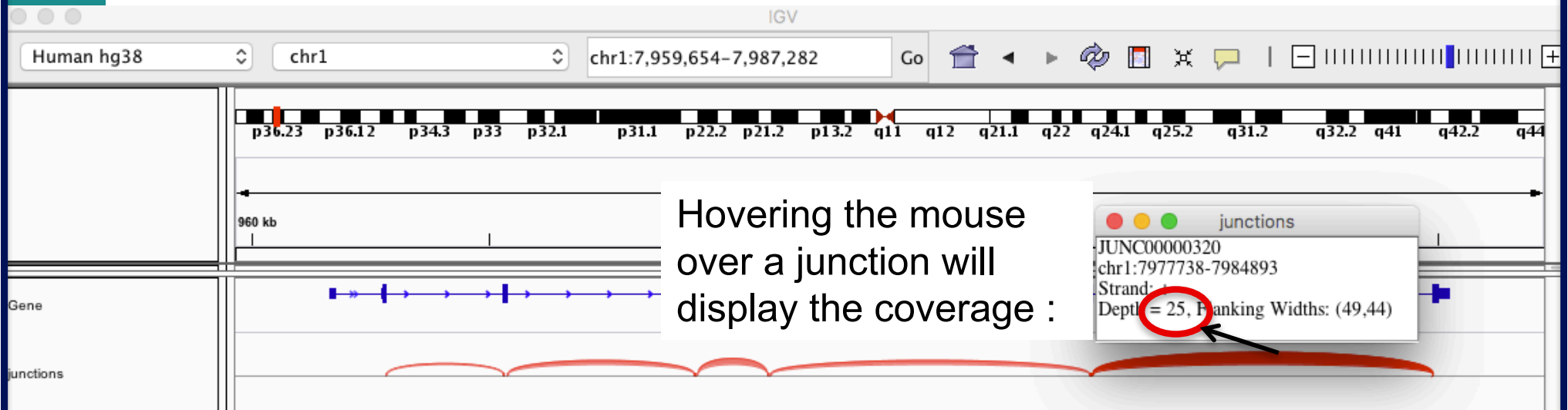  - File → Load from file and choose the downloaded BED file

# Exercise 1
# 2.3. Splice junctions visualization

# Exercise 1
# 2.3. Splice junctions visualization

- IGV



Hovering the mouse over a junction will display the coverage :

junctions

JUNC00000320
chr1:7977738-7984893
Strand:
Depth = 25, Flanking Widths: (49,44)

- BED file on Galaxy

chr1    7977689    7984937    JUNC00000320    (25)    +    7977689    7984937    255,0,0    2    49,44    0,7204

7977689+49 = 7977738 : junction start position
7984937-44 = 7984893 : junction end position

→ 25 alignments span the junction that joins the last 2 exons of *Park7* gene

# Exercise 1
# 3. Alignment visualization

- **Galaxy**

  3.1. Tophat2 provides an alignment in BAM format

  3.2. Download this file together with the corresponding index (in the same directory)
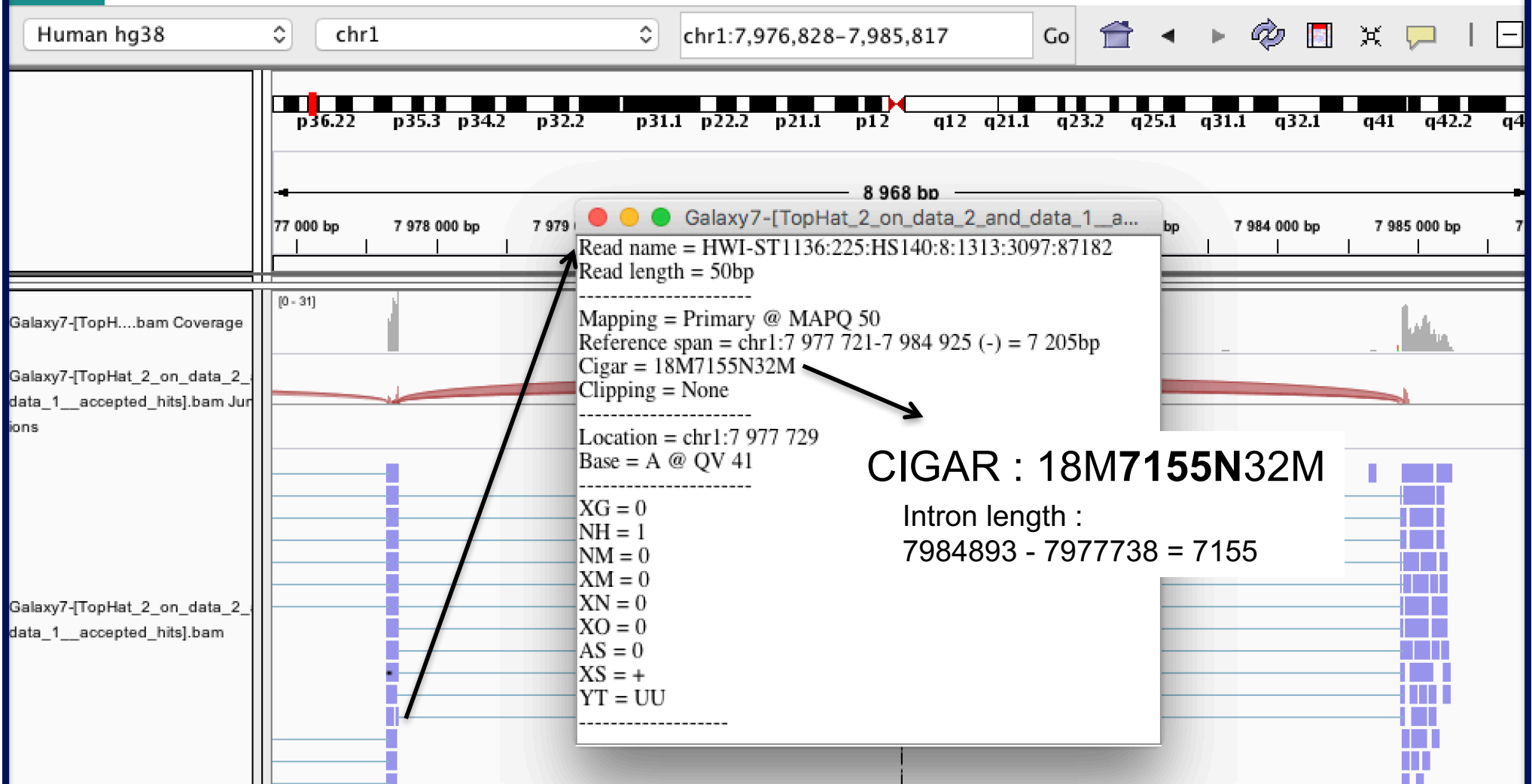
  

- **IGV**

  - File → Load from file and choose the downloaded BAM file

# Exercise 1
## 3.3. Reads aligned on a splice junction



Human hg38 | chr1 | chr1:7,976,828–7,985,817 | Go

p36.22  p35.3  p34.2  p32.2  p31.1  p22.2  p21.1  p12  q12  q21.1  q23.2  q25.1  q31.1  q32.1  q41  q42.2  q4

8 968 bp

77 000 bp    7 978 000 bp    7 979    Galaxy7-[TopHat_2_on_data_2_and_data_1_a...    bp    7 984 000 bp    7 985 000 bp

Galaxy7-[TopH....bam Coverage

Galaxy7-[TopHat_2_on_data_2_
data_1__accepted_hits].bam Jur
ions

[0 - 31]

Galaxy7-[TopHat_2_on_data_2_
data_1__accepted_hits].bam

Read name = HWI-ST1136:225:HS140:8:1313:3097:87182
Read length = 50bp
---------------------
Mapping = Primary @ MAPQ 50
Reference span = chr1:7 977 721-7 984 925 (-) = 7 205bp
Cigar = 18M7155N32M
Clipping = None
---------------------
Location = chr1:7 977 729
Base = A @ QV 41
---------------------
XG = 0
NH = 1
NM = 0
XM = 0
XN = 0
XO = 0
AS = 0
XS = +
YT = UU
---------------------

CIGAR : 18M**7155N**32M
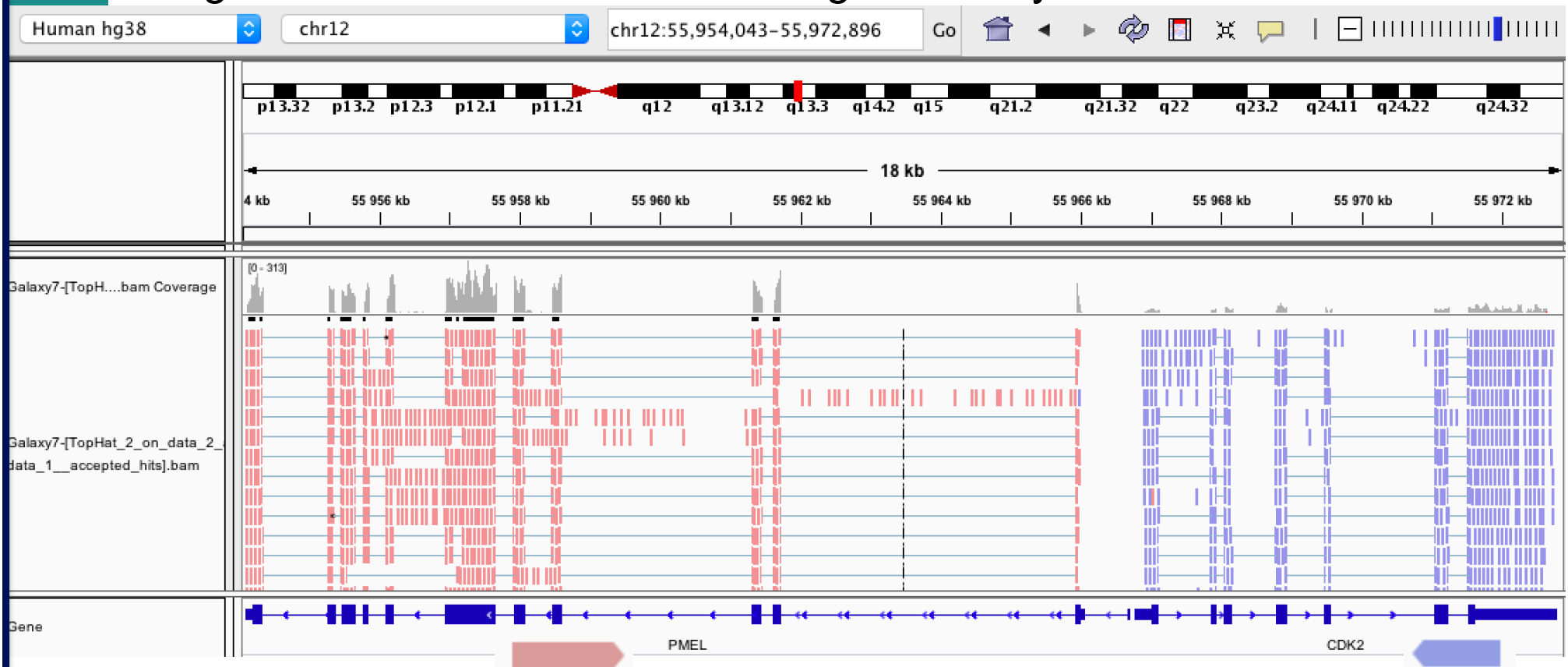
Intron length :
7984893 - 7977738 = 7155

# Exercise 1
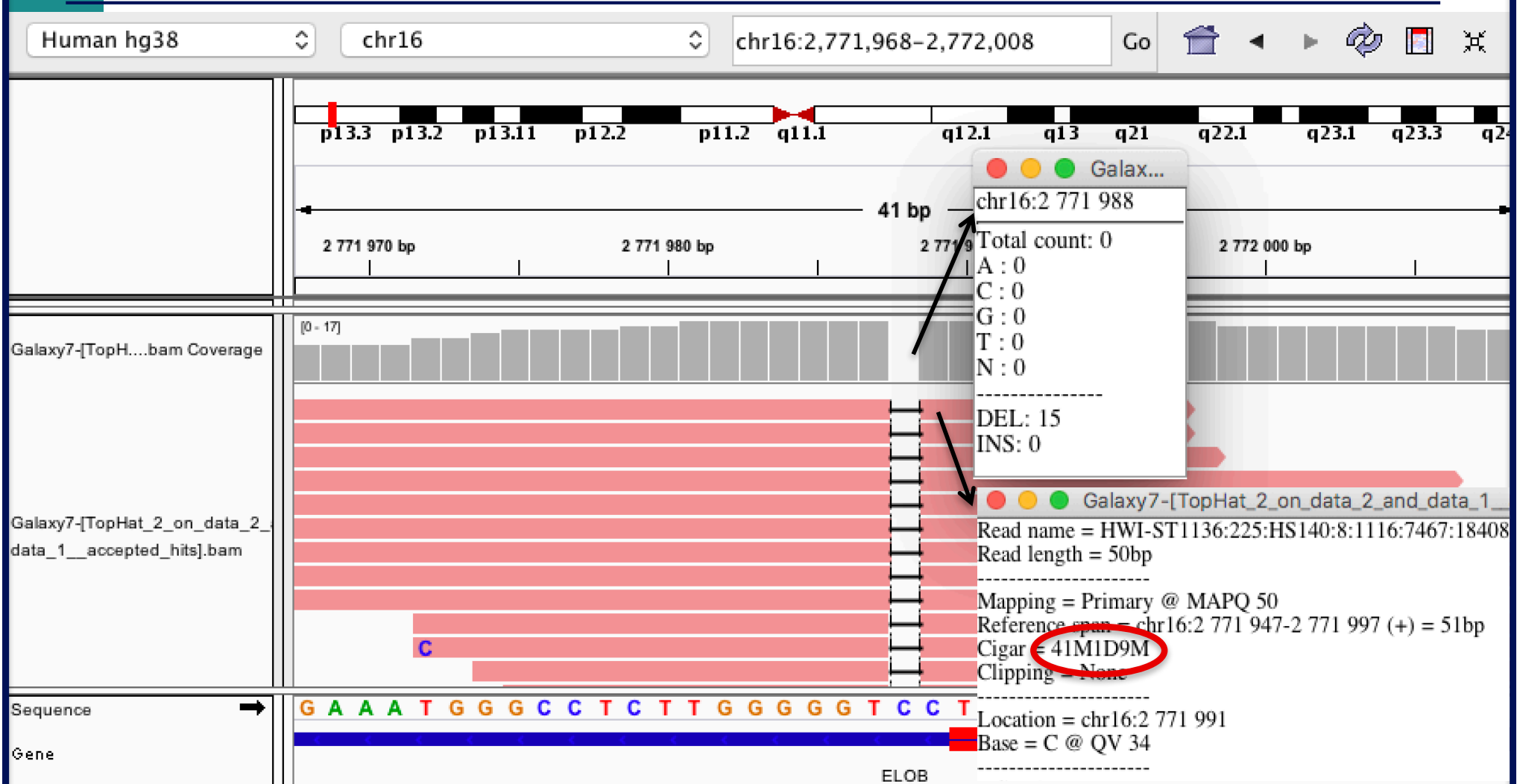## 3.4. Visualization of strand specificity

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



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

# Exercise 1
## 3.5. Visualization of a deletion



15 reads aligned with a deletion at position chr16:2,771,988

# Exercise 2 - Question 1
# Proportion of mapped reads in all samples

Galaxy : Shared Data → Data Libraries → CNRS training
RNAseq → alignment → align_summary :

| Name | Tophat2 on siLuc2: align_summary |
| --- | --- |

Reads:

Input : 43672265

Mapped : 42797297 (98.0% of input)

of these: 5829092 (13.6%) have multiple alignments (1132 have >20)

98.0% overall read mapping rate.

| Name | Tophat2 on siLuc3: align_summary |
| --- | --- |

Reads:

Input : 46565834

Mapped : 45633110 (98.0% of input)

of these: 6030755 (13.2%) have multiple alignments (861 have >20)

98.0% overall read mapping rate.

| Name | Tophat2 on siMitf3: align_summary |
| --- | --- |

Reads:

Input : 43985979

Mapped : 43048694 (97.9% of input)

of these: 5763991 (13.4%) have multiple alignments (765 have >20)

97.9% overall read mapping rate.

| Name | Tophat2 on siMitf4: align_summary |
| --- | --- |

Reads:

Input : 51348313

Mapped : 50317655 (98.0% of input)

of these: 6826164 (13.6%) have multiple alignments (643 have >20)

98.0% overall read mapping rate.

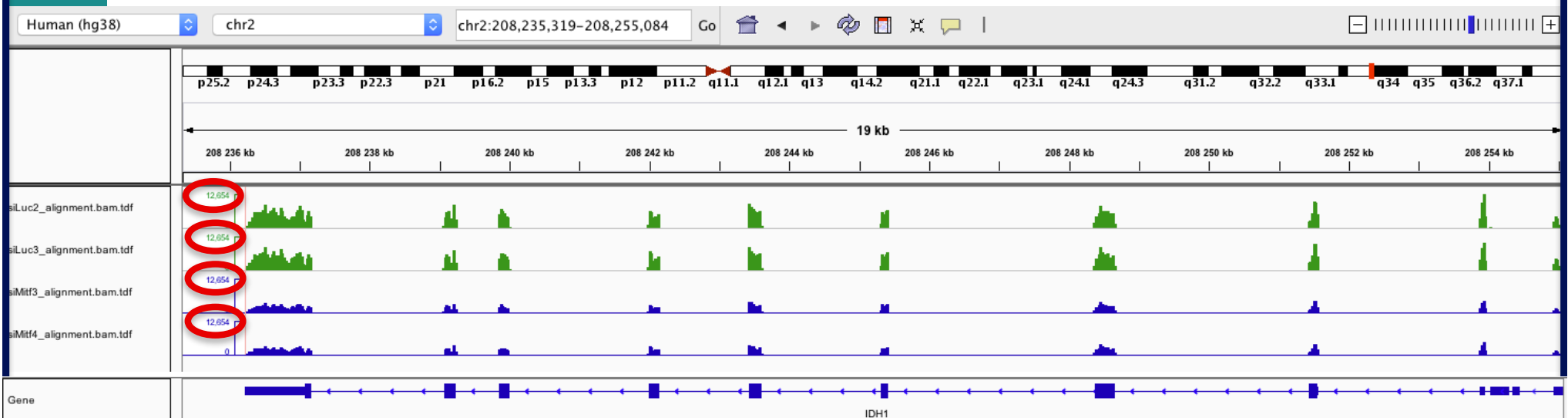→ This proportion is high and consistent across samples

# Exercise 2 – Question 2

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)
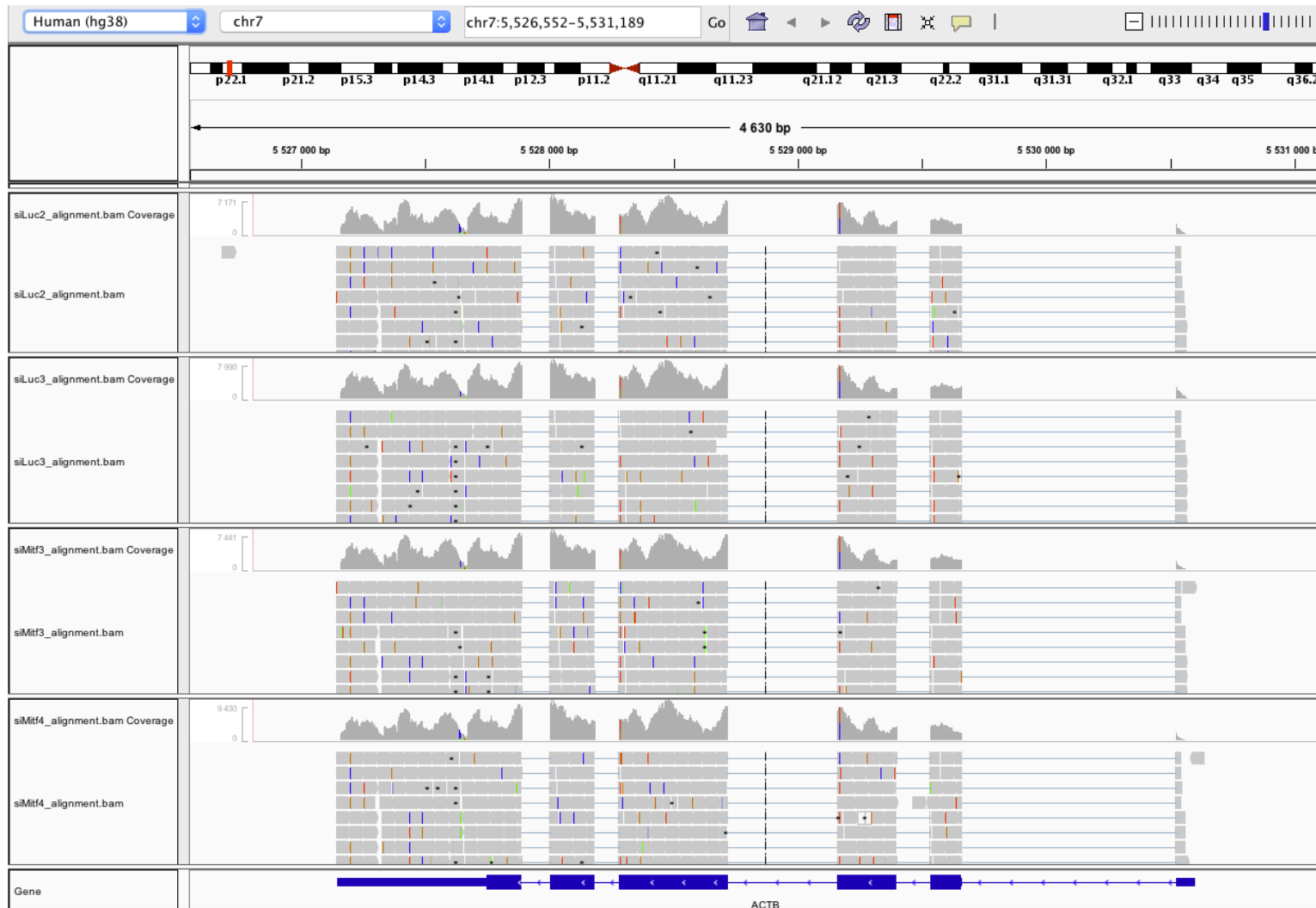
→ all tracks are on the same scale



*Idh1* is under-expressed in siMitf samples compared to siLuc ones
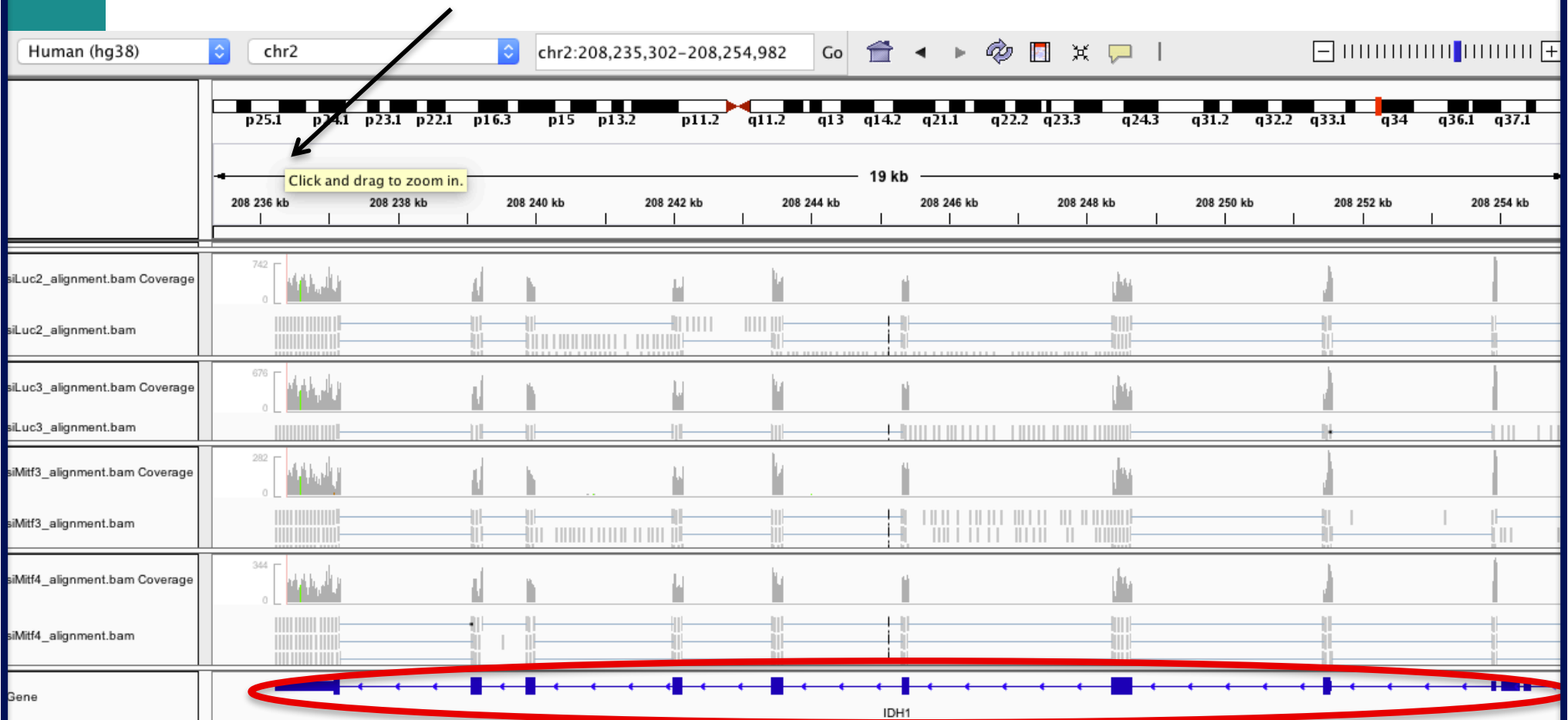
# Exercise 2 – Question 3
# Alignments visualization

IGV : File → Load from file and select the 4 BAM files
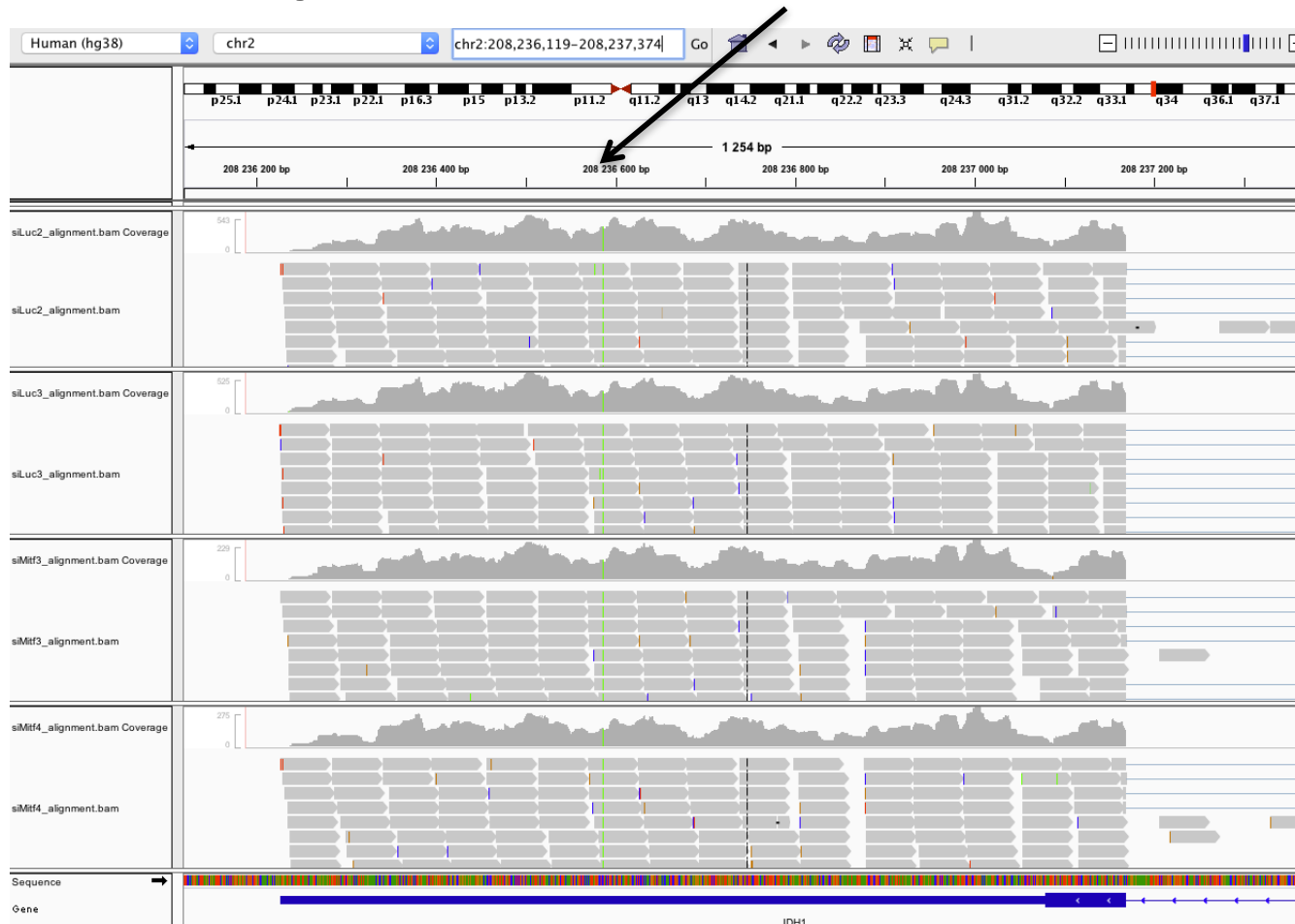
# Exercise 2 - Question 3

Click and drag to define a window around the last exon to zoom in



Arrows indicate annotated transcribed strand
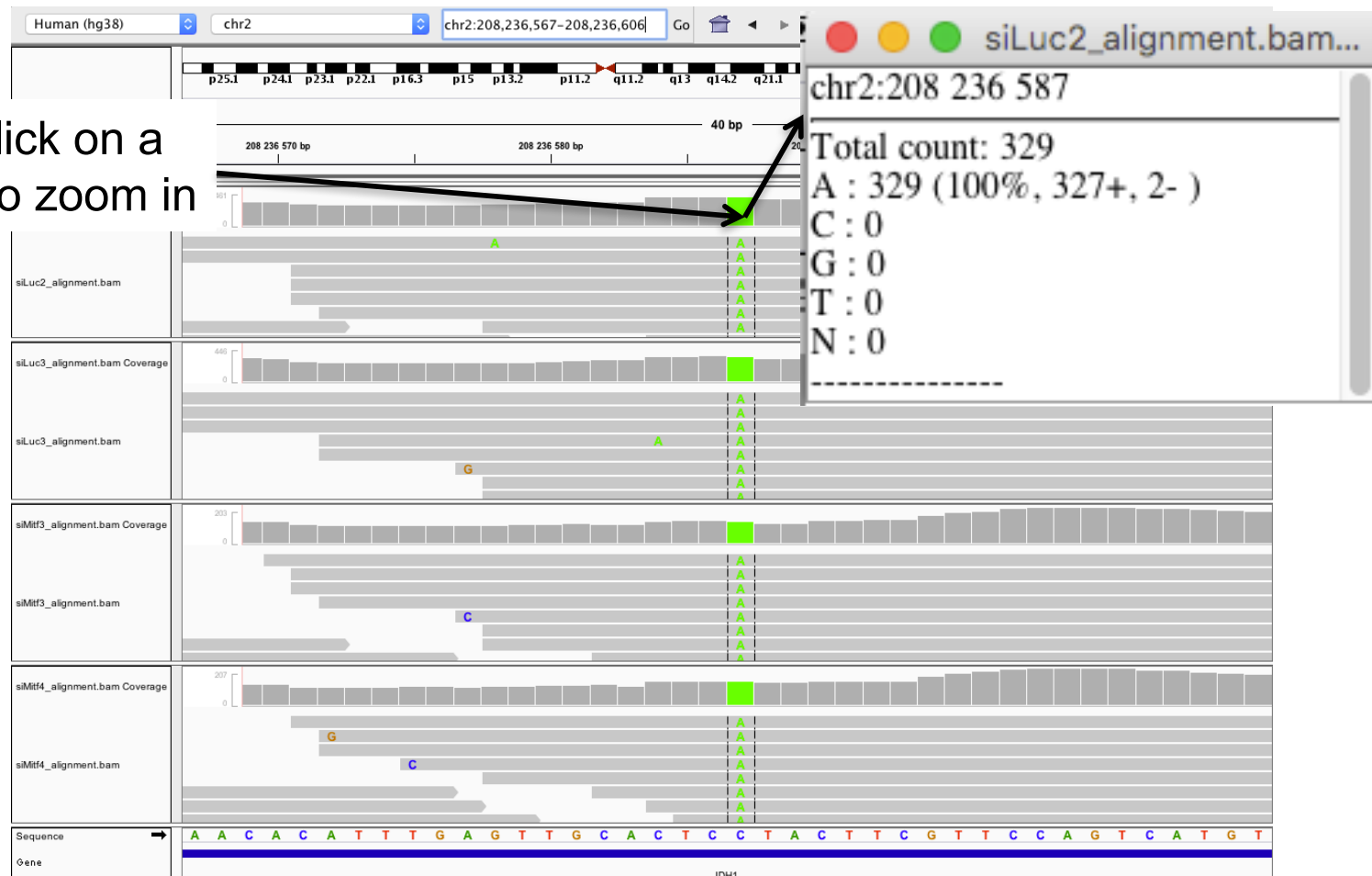
# Exercise 2 - Question 3

- You can see a nucleotide difference in green
- Click and drag to zoom in around this position

# Exercise 2 - Question 3

- In the location chr2:208,236,587 :
  - A in 100% of the RNA-seq reads, C in the reference genome
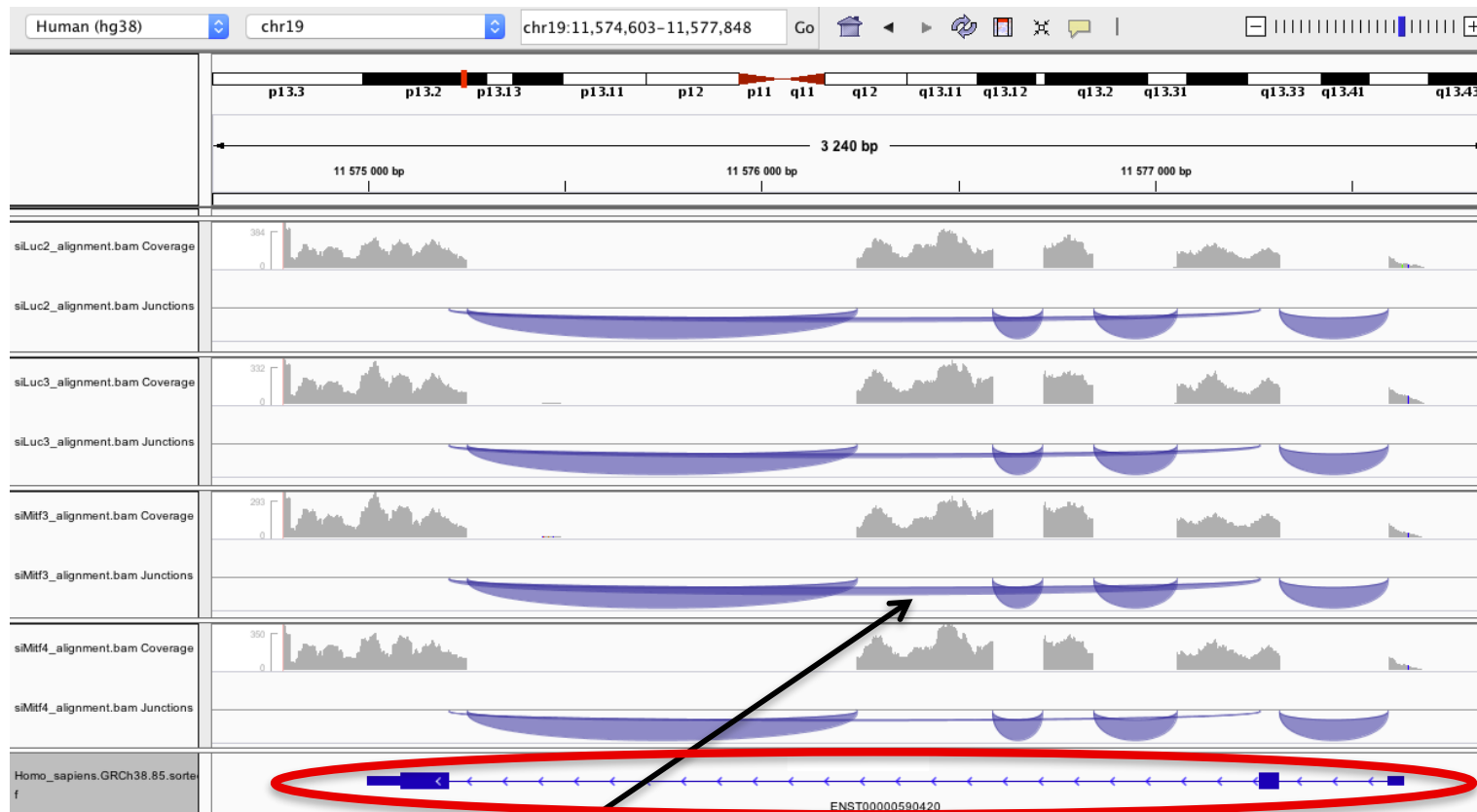
Double click on a position to zoom in

# Exercise 2 – Question 4



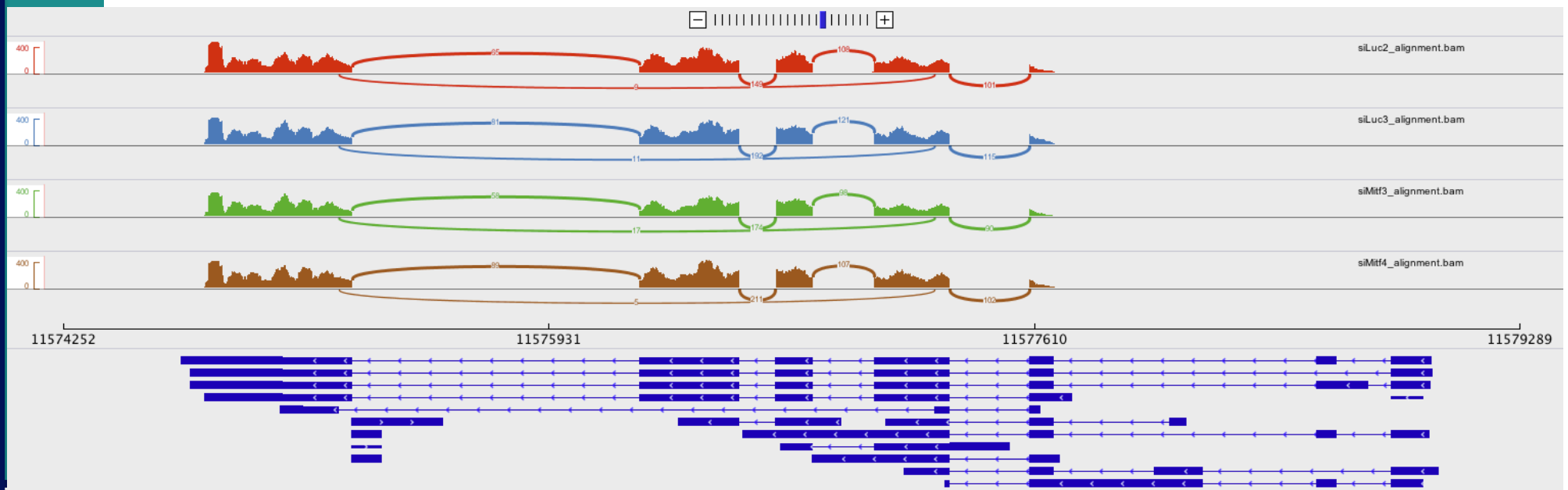This junction is not in Refseq annotations

# Exercise 2 – Question 4

- File → Load from file and select Ensembl annotations (Homo_sapiens.GRCh38.85.sorted.gtf)
- Right click on Ensembl annotations track and select Expanded



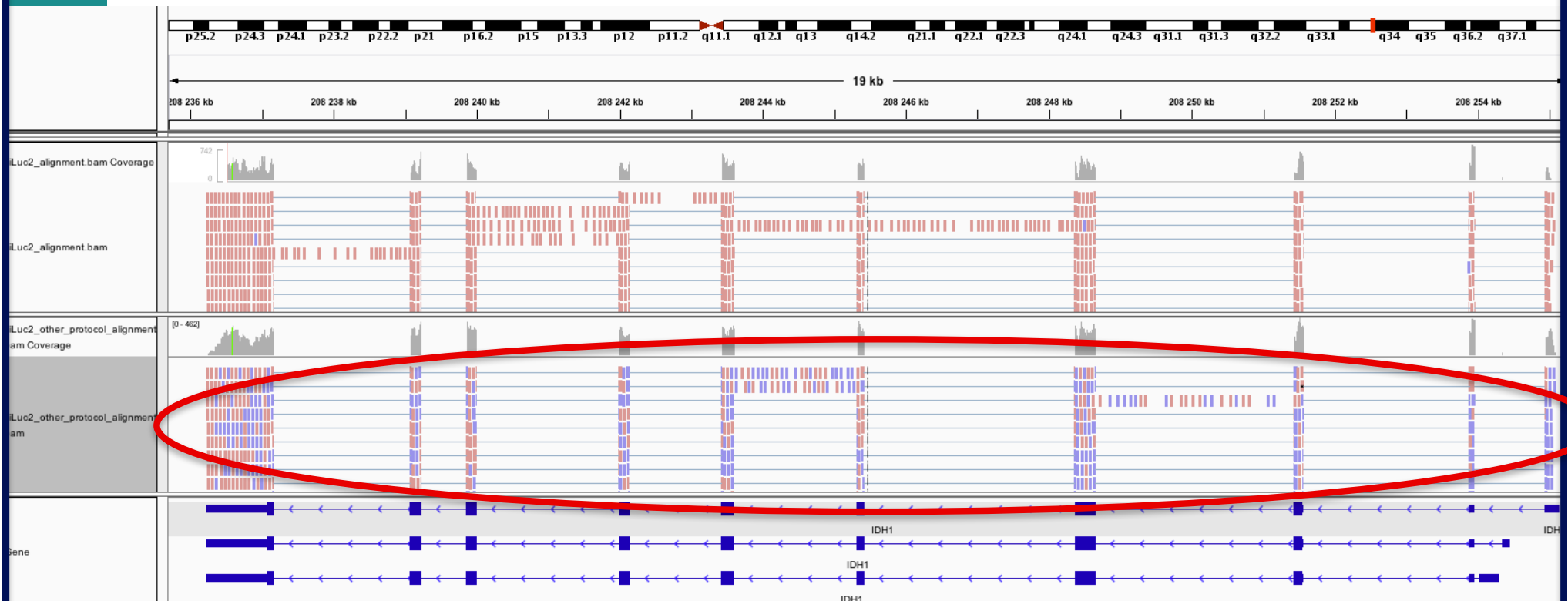This junction is in Ensembl annotations

# Exercise 2 – Question 4

- Sashimi plot



➔ Very useful to quickly screen differentially spliced exons
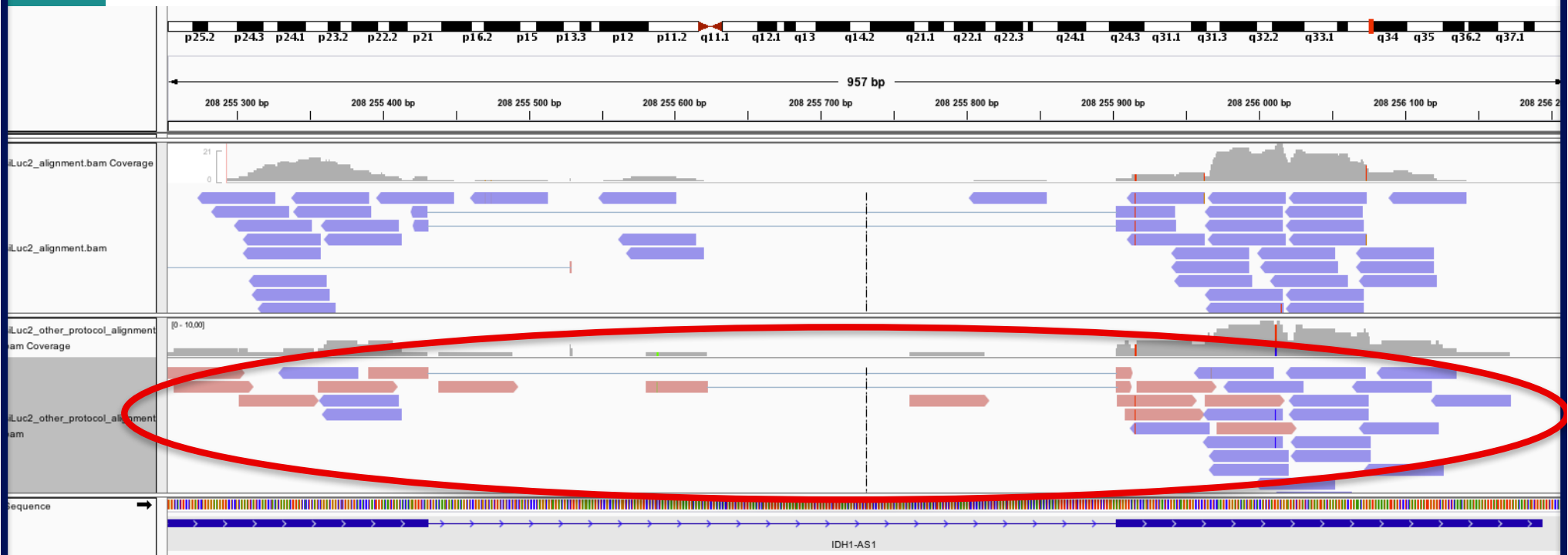   along genomic regions of interest (more accurate with paired-end data)

# Exercise 2 – Question 5

- File → load from file and select siLuc2_other_protocol_alignment.bam
- Right-click on BAM file → Color alignments by → read strand
- e.g. *Idh1* gene

# Exercise 2 – Question 5

- e.g. *Idh1-as1* gene



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