# Data mining with Ensembl Biomart

Stéphanie Le Gras
(slegras@igbmc.fr)

# Guidelines

- Genome data

- Genome browsers

- Getting access to genomic data: Ensembl/BioMart

# Genome Sequencing

Example: Human genome

- 2000: First draft of the human genome
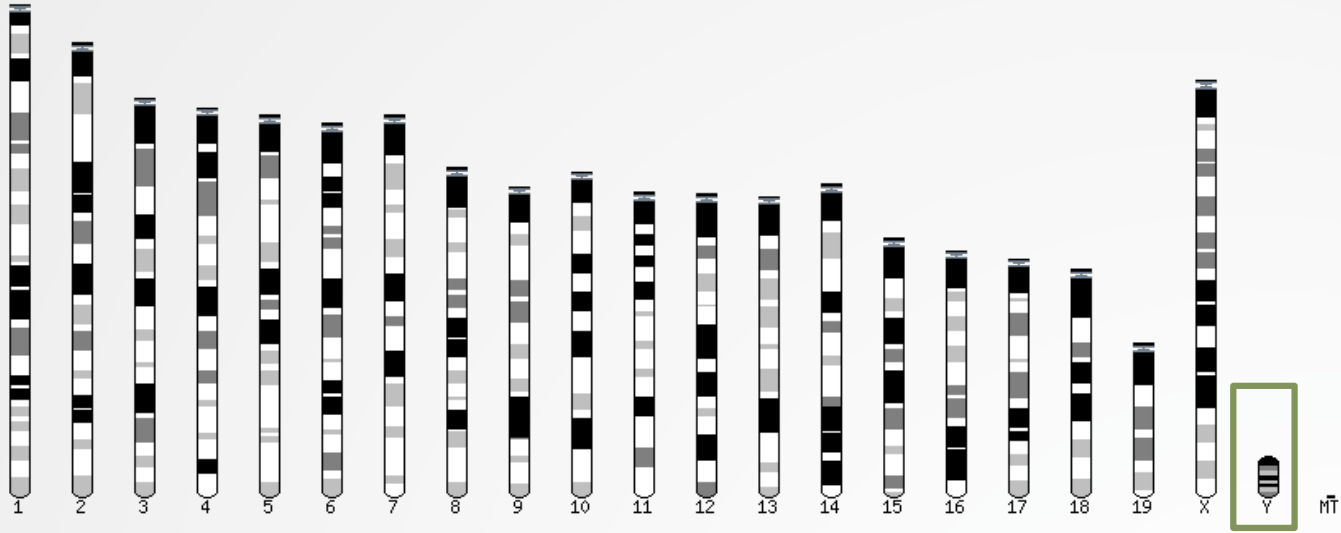
- 2003: Human genome sequencing complete

# Genome builds

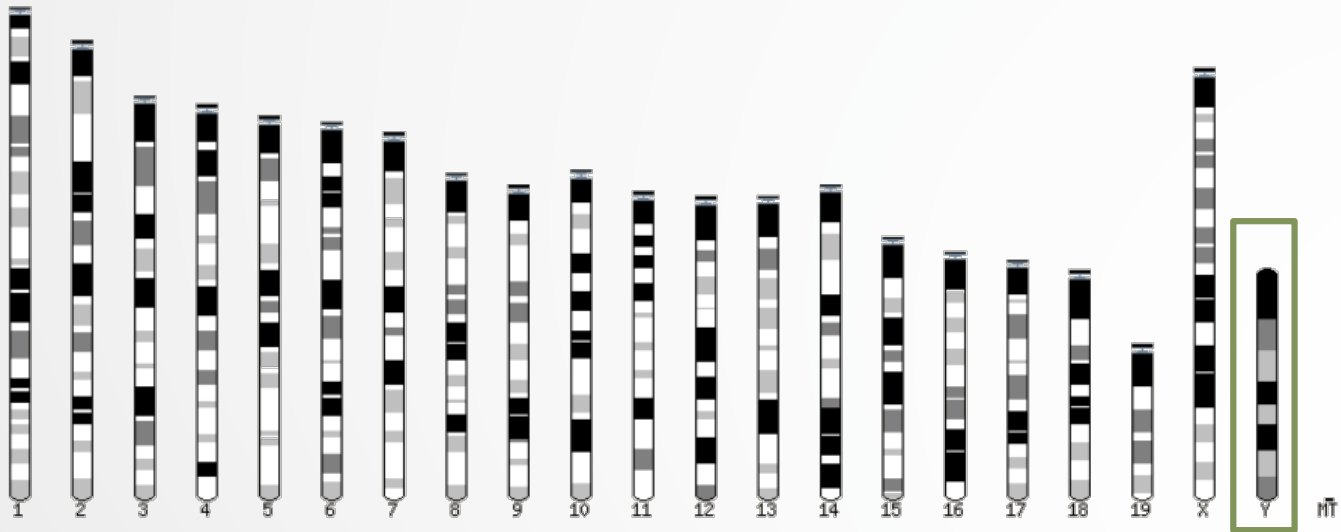| SPECIES | UCSC VERSION | RELEASE DATE | RELEASE NAME | STATUS |
|---|---|---|---|---|
| **MAMMALS** | | | | |
| Human | hg38 | Dec. 2013 | Genome Reference Consortium GRCh38 | Available |
| | hg19 | Feb. 2009 | Genome Reference Consortium GRCh37 | Available |
| | hg18 | Mar. 2006 | NCBI Build 36.1 | Available |
| | hg17 | May 2004 | NCBI Build 35 | Available |
| | hg16 | Jul. 2003 | NCBI Build 34 | Available |
| | hg15 | Apr. 2003 | NCBI Build 33 | Archived |
| | hg13 | Nov. 2002 | NCBI Build 31 | Archived |
| | hg12 | Jun. 2002 | NCBI Build 30 | Archived |
| | hg11 | Apr. 2002 | NCBI Build 29 | Archived (data only) |
| | hg10 | Dec. 2001 | NCBI Build 28 | Archived (data only) |
| | hg8 | Aug. 2001 | UCSC-assembled | Archived (data only) |
| | hg7 | Apr. 2001 | UCSC-assembled | Archived (data only) |
| | hg6 | Dec. 2000 | UCSC-assembled | Archived (data only) |
| | hg5 | Oct. 2000 | UCSC-assembled | Archived (data only) |
| | hg4 | Sep. 2000 | UCSC-assembled | Archived (data only) |
| | hg3 | Jul. 2000 | UCSC-assembled | Archived (data only) |
| | hg2 | Jun. 2000 | UCSC-assembled | Archived (data only) |
| | hg1 | May 2000 | UCSC-assembled | Archived (data only) |

Source: https://genome.ucsc.edu/FAQ/FAQreleases.html
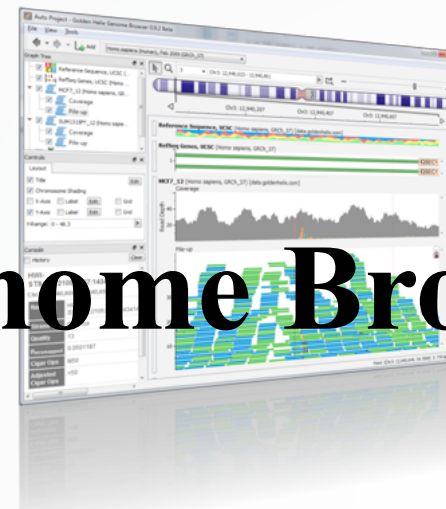
# Genome builds



mm9

mm10

# Get access to genomic data

- Need a way to gather all genomic information in one place

- Availability of the data

- Accessibility to the data

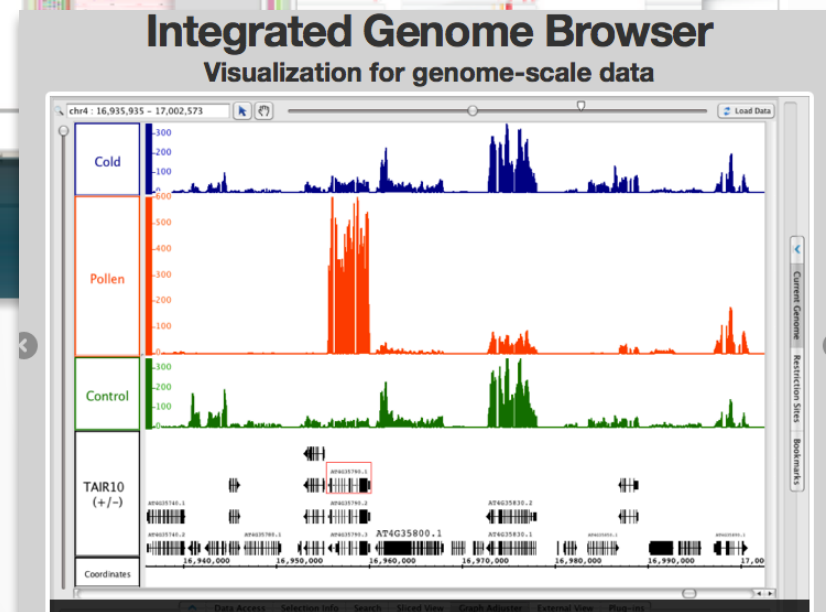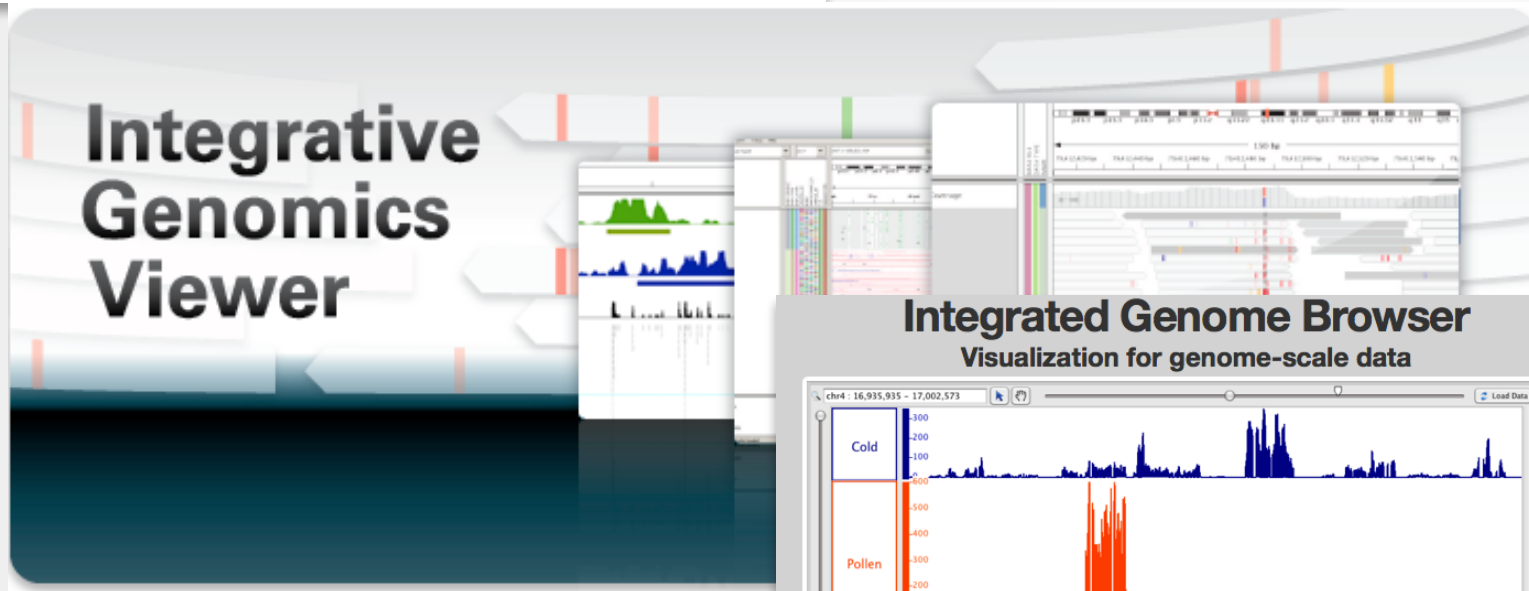**Genome Browser**

# Genome browsers

# Genome Browsers

- Graphical interface to display genomic data

- Visualize and browse entire genomes with annotated data
  - Gene prediction and structure
  - Proteins,
  - Expression,
  - Regulation,
  - Variation,
  - Comparative analysis...

# There are Genome Browsers...

## EBI - Ensembl



## UCSC – Genome Browser





## NCBI – Map Viewer

# And Genome browsers...

# Getting access to genomic data: ENSEMBL/BIOmart

# Access Ensembl's data

## Web site



## Mining tool: BioMart



🙂 User friendly
🙂 Straightforward
🙁 Only one request at once

🙂 Get answer to complex query
🙂 Very fast
🙁 Need training

# BioMart

- http://www.biomart.org/

- Joint development between EBI and Cold Spring Harbor Laboratory (CSHL)

- Open source project

- BioMart can access diverse databases from a single interface

- It is search engine that can find multiple terms and put them into a table format

- No programming required!

# Many uses of BioMart

# BioMart/Ensembl



- Get access to :
  - Genomic annotation (genes, SNPs)
  - Functional annotation
  - Expression data

15

# Example: Step 1 (Select datasets)

# Example: Step 2 (Filter)



e!Ensembl  BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors

Register

New   Count   Results              URL   XML   Perl   Help

**Dataset**
Human genes (GRCh38.p7)

**Filters**
Chromosome/scaffold: 1
Gene Start (bp): 78895
Gene End (bp): 10000000

**Attributes**
Gene ID
Transcript ID

**Dataset**
[None Selected]

**Please restrict your query using criteria belo
(If filter values are truncated in any lists, hover over the list it       to see the full text)

⊟ REGION:

☑ Chromosome/scaffold

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
```

Limit to chromosome 1

☑ Base pair
Gene Start (bp)         78895
Gene End (bp)           04561

Limit to given coordinates

☐ Band
Band Start
Band End

17

# Example: Step 3 (Count results)

# Example: Step 4 (Select attributes)

# Example: Step 4 (get results)

# Exercise 1: get annotations of a gene

- 1. Using Ensembl/BioMart, retrieve all transcripts IDs and the gene ID of IDH1 gene (human). How many transcripts the gene IDH1 has?
  - Use Ensembl Gene **v85**, for Human GRCh38.p7
  - Click on Filters :
    - Expand the GENE section
    - Select « Input external references ID list »
    - Select HGNC symbol(s) in the drop down menu
    - Enter IDH1 in the text box
  - Click on Attributes :
    - Select "Features" (top panel, selected by default)
    - Select Gene ID, Transcript ID, Associated Gene Name

- 2. Extract all exon sequences of the IDH1 gene in fasta format. Headers will contain the Associated gene names, transcript IDs and Exon IDs.

- 3. Extract all coding sequences of the IDH1 gene in fasta format. Headers will contain the transcript IDs and Exon IDs.

- 4. Retrieve GO-terms associated to the IDH1 gene (select GO Term Name, GO domain and GO Term Accession along with Gene ID, Transcript ID and Associated Gene Name)

- 5. Retrieve the germline variations found in this gene. Annotations to be found (Variant Name, Variant Alleles, Minor allele frequency, Chromosome/scaffold name, Chromosome/scaffold position start (bp), Chromosome/scaffold position end (bp), Variant Consequence along with Gene ID, Transcript ID and Associated Gene Name)

# Exercise 2: get annotations for a set of genes

- Annotate the file siMitfvssiLuc.up.txt you have generated using SARTools with gene annotations extracted from Ensembl/BioMart
  - If you encountered any trouble with the generation of the dataset
    - go to GalaxEast (http://use.galaxeast.fr)
    - go to Shared Data/ Data Libraries / CNRS training / RNAseq / statistical_analysis.
    - Import the dataset SARTools_DESeq2_tables to your history.
    - Click on 👁 to display the content of the dataset and download the file siMitfvssiLuc.up.txt (click right, save …)

- 1. Open the file siMitfvssiLuc.up.txt and change the name of the column which contains "Id" to "**Gene ID**". Save the change.

- 2. Use the file siMitfvssiLuc.up.txt to extract gene annotations for those genes. Annotation to extract are : gene IDs, chromosome, start of gene, end of gene, strand, associated gene name, gene type. Save the results to a compressed TSV file. (don't close the Ensembl/Biomart window once done)

- 3. Upload the file siMitfvssiLuc.up.txt and the annotation file you obtained from Ensembl/BioMart to GalaxEast into your current history "RNA-seq data analysis".
  - Type: tabular
  - Genome: hg38

# Exercise 2: get annotations for a set of genes

- 4. Use the tool "Join two Datasets" to merge the two datasets based on the Gene IDs.
  - Gene IDs are used as unique identifiers common to the two datasets. For a given gene, data spread in the two files are going to be merged in the same line in the newly generated file.

- 5. rename the generated dataset in 4. to siMitfvssiLuc.up.annot.txt

- 6. Is there lncRNAs in the upregulated genes? Use the tool "Filter data on any column using simple expressions" to search for "lincRNA" in the dataset siMitfvssiLuc.up.annot.txt

- 7. Go back to Ensembl/BioMart. You want to run a *de novo* motif discovery on all promoters of the up-regulated genes (the ones from the file siMitfvssiLuc.up.txt). Extract the promoter sequences of all up-regulated genes: retrieve the 2kb upstream of the transcripts of these genes.

## Exercise 3: get annotations in the genome

- 1. How many genes are located in the genomic region: **2:208226227-208276270**

- 2. Extract the coordinates of all human genes located on chromosomes (exclude scaffolds). Information to extract for each gene: Gene ID, Chromosome/scaffold name, Gene Start (bp), Gene End (bp), strand and associated Gene Name