

Data mining with Ensembl Biomart

Stéphanie Le Gras
(slegras@igbmc.fr)

Guidelines

- Genome data
- Genome browsers
- Getting access to genomic data: Ensembl/BioMart

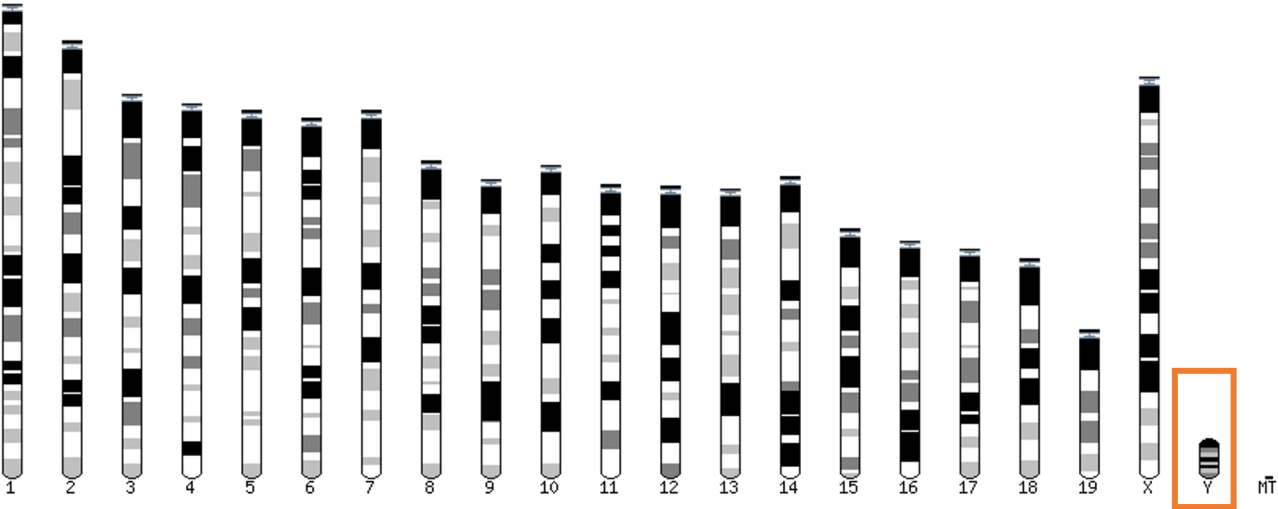
Genome builds

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)

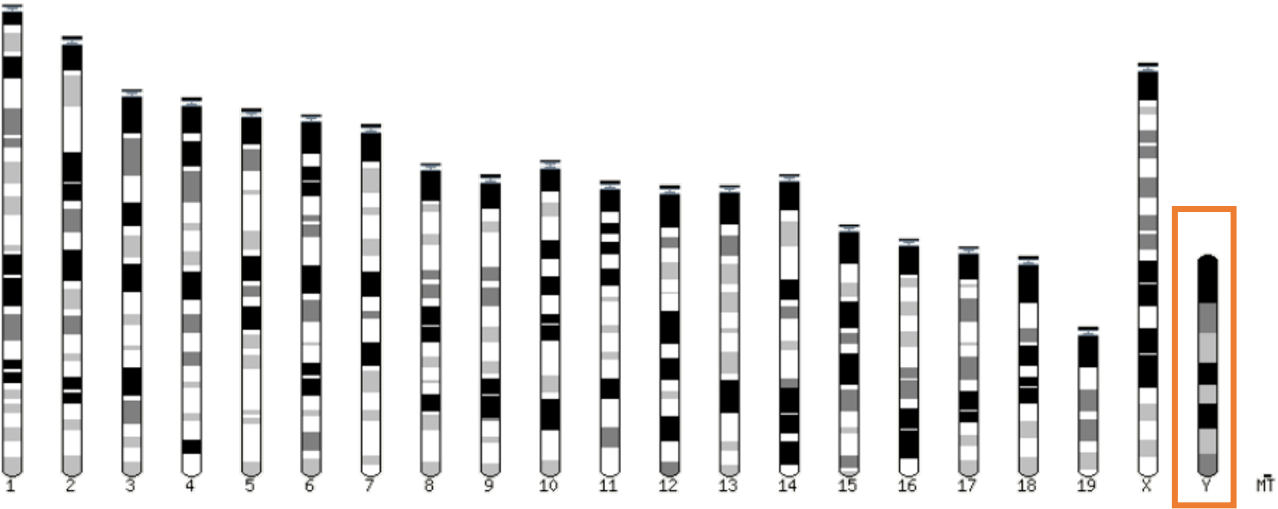
Source: <https://genome.ucsc.edu/FAQ/FAQreleases.html>

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



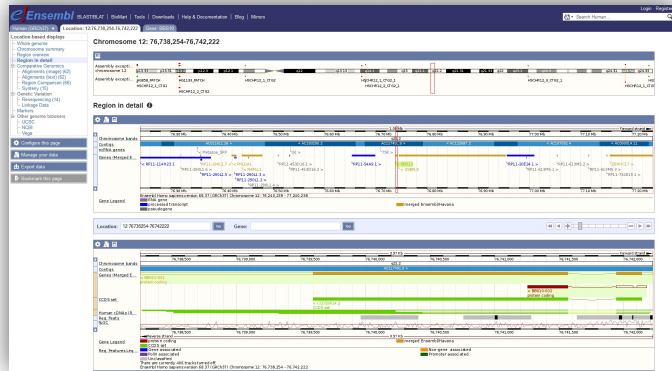
Genome browsers

Genome Browsers

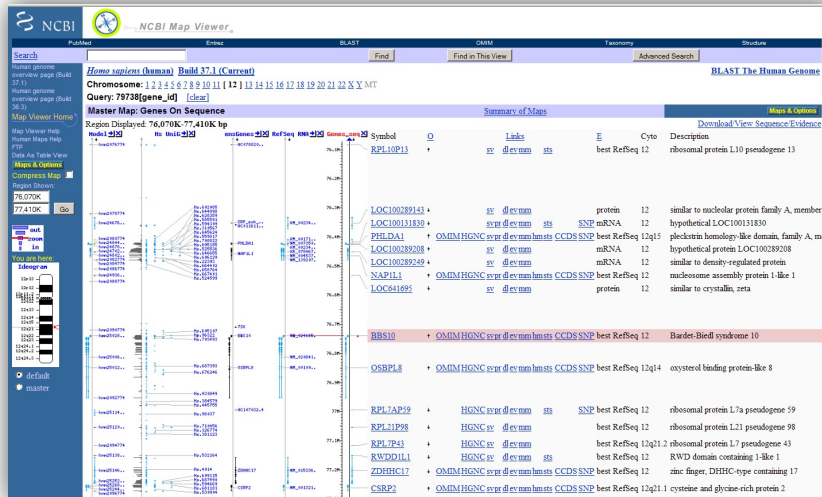
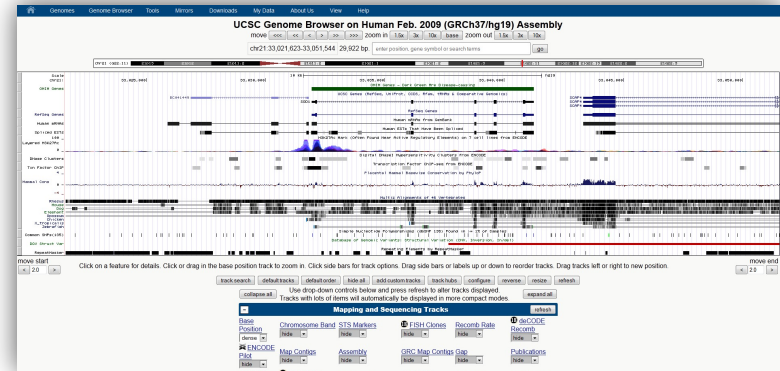
- Graphical interface to display genomic data
- Visualize and browse entire genomes with annotated data
 - Gene prediction and structure
 - Proteins,
 - Expression,
 - Regulation,
 - Variation,
 - Comparative analysis...

There are Genome Browsers...

EBI - Ensembl

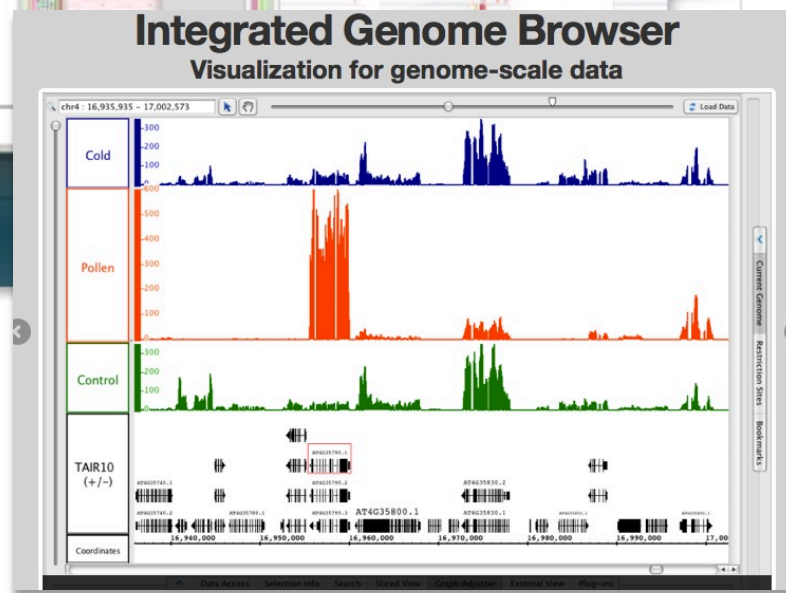
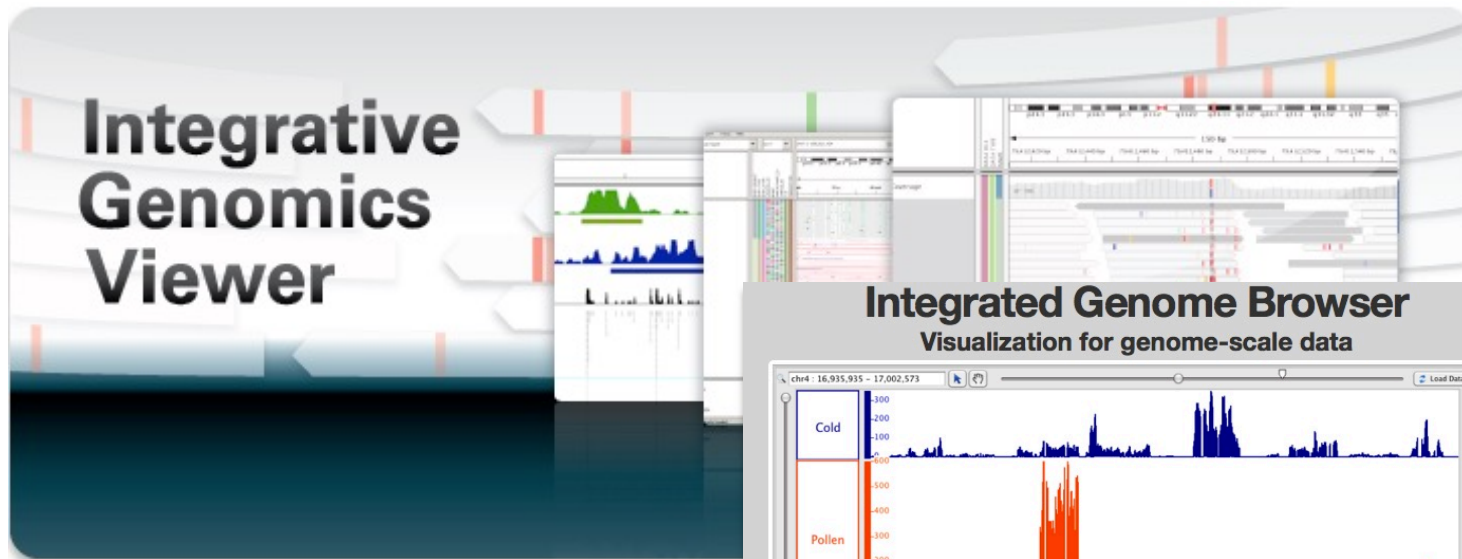


UCSC – Genome Browser



NCBI – Genome Data Viewer

And Genome browsers...



Getting access to genomic data: ENSEMBL/BIOmart




Access Ensembl's data




Web site

The screenshot shows the Ensembl website homepage. At the top, there is a navigation bar with links for BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, and Mirrors. A search bar is located in the top right corner. Below the navigation bar, there is a search box with the text "Search: [All species] for []" and a search button. The main content area is divided into several sections: "Browse a Genome" with a list of popular genomes (Human, Mouse, Zebrafish), "What's New in Ensembl Release 83 (December 2015)" with a list of updates, "Latest blog posts" with a list of recent posts, and "Tweets by @ensembl" with a list of tweets. There are also several interactive widgets for "Bill using Human GRCh37", "Variant Effect Predictor", "Gene expression in different tissues", "Find SNPs and other variants for my gene", "Retrieve gene sequence", "Compare genes across species", "Use my own data in Ensembl", and "ENCODE data in Ensembl".

Mining tool: BioMart

The screenshot shows the Ensembl BioMart interface. At the top, there is a navigation bar with links for BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, and Mirrors. A search bar is located in the top right corner. Below the navigation bar, there is a search box with the text "Search: [All species] for []" and a search button. The main content area is divided into several sections: "Dataset" with a dropdown menu for "CHOOSE DATABASE", "Filters (filtering and inputs)", "Attributes (desired output)", and "Results". There are also links for "BioMart tutorial", "YouTube", and "YouKu".


-  User friendly
-  Straightforward
-  Only one request at once

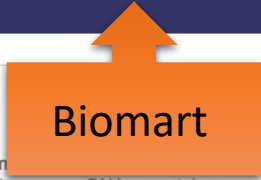
-  Get answer to complex query
-  Very fast
-  Need training

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!

BioMart/Ensembl

 [BLAST/BLAT](#) | [VEP](#) | [Tools](#) | [BioMart](#) | [Downloads](#) | [Help & Docs](#) | [Blog](#) [Login/Register](#)



Tools [All tools](#)

BioMart > Export custom datasets from Ensembl with this data-mining tool

Biomart

Variant Effect Predictor > Analyse your own variants and predict the functional consequences of known and unknown variants

Ensembl is a genome browser for vertebrate genomes that supports research in comparative genomics, evolution, sequence variation and transcriptional regulation. Ensembl annotate genes, computes multiple alignments, predicts regulatory function and collects disease data. Ensembl tools include BLAST, BLAT, BioMart and the Variant Effect Predictor (VEP) for all supported species.

Ensembl Release 95 (January 2019)

- New regulatory build for human, incorporating new data from ENCODE
- Update to GENCODE M20 for mouse
- New genomes: donkey, polar bear, black bear, red fox, koala, dingo, tuatara, painted turtle and desert tortoise
- Updated genomes for chicken, cow and horse
- New protein structure variation view

[More release news](#) on our blog

Search

All species

e.g. [BRCA2](#) or [rat 5:62797383-63627669](#) or [rs699](#) or [coronary heart disease](#)

All genomes

- [View full list of all Ensembl species](#)

Favourite genomes

 **Human**
GRCh38.p12

Still using GRCh37?

Other news from our blog

- 01 Mar 2019: [Getting to know us: Guy from Ensembl Plants](#)
- 27 Feb 2019: [Job: Ensembl Infrastructure Project Leader](#)
- 27 Feb 2019: [Custom data upload: creating URLs for large](#)

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

Example: Step 1 (Select datasets)

Dataset
[None selected]

Ensembl Genes 101

- ✓ - CHOOSE DATASET -
- Chicken genes (GRCg6a)
- Human genes (GRCh38.p13)**
- Mouse genes (GRCm38.p6)
- Rat genes (Rnor_6.0)
- Zebrafish genes (GRCz11)

- Abingdon island giant tortoise genes (ASM359739v1)
- Agassiz's desert tortoise genes (ASM289641v1)
- Algerian mouse genes (SPRET_EiJ_v1)
- Alpaca genes (vicPac1)
- Alpine marmot genes (marMar2.1)
- Amazon molly genes (Poecilia_formosa-5.1.2)
- American beaver genes (C.can_genome_v1.0)
- American bison genes (Bison_UMD1.0)
- American black bear genes (ASM334442v1)
- American mink genes (NNQGG.v01)
- Angola colobus genes (Cang.pa_1.0)
- Anole lizard genes (AnoCar2.0)
- Arabian camel genes (CamDro2)
- Argentine black and white tegu genes (HLtupMer3)
- Armadillo genes (Dasnov3.0)
- Asian bonytongue genes (fSciFor1.1)
- Atlantic herring genes (Ch_v2.0.2)
- Atlantic salmon genes (ICSASG_v2)
- Australian saltwater crocodile genes (CroPor_comp1)
- Ballan wrasse genes (BallGen_V1)
- Barramundi perch genes (ASB_HGAPassembly_v1)
- Bengalese finch genes (LonStrDom1)

First choose database and dataset

Example: Step 2 (Filter)

Ensembl BLAST/BLAT | VEP | Tools | BioMart | Downloads | Help & Docs | Blog

New | Count | Results | URL | XML | Perl | Help

Dataset
Human genes (GRCh38.p13)

Filters

Chromosome/scaffold: 1
Start: 78895
End: 224561

Attributes
Gene stable ID
Gene stable ID version
Transcript stable ID
Transcript stable ID version

Dataset
[None Selected]

8
9
10
11
12
13
14
15
16
17
18
19
20

Coordinates
Start: 78895
End: 224561

Karyotype band
Band Start
Band End

Marker

Limit to chromosome 1

Limit to given coordinates

Example: Step 3 (Count results)

Compute result count

Database 12 / 67130 Genes
Human genes (GRCh38.p13)

Filters
Chromosome/scaffold: 1
Start: 78895
End: 224561

Attributes
Gene stable ID
Gene stable ID version
Transcript stable ID
Transcript stable ID version

Dataset
[None Selected]

Please restrict your query using criteria below
(If filter values are truncated in any lists, hover over the list item 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

Coordinates

Example: Step 4 (Select attributes)

e!Ensembl BLAST/BLAT | VEP | Tools | BioMart | Downloads | Help & Docs | Blog Login/Register

Search all species...

New Count Results URL XML Perl Help

Dataset 12 / 67130 Genes
Human genes (GRCh38.p13)

Filters
Chromosome/scaffold: 1
Start: 78895
End: 224561

Attributes
Gene stable ID
Gene stable ID version

Dataset
[None Selected]

Please select columns to be included in the output and hit 'Results' when ready

Missing non coding genes in your mart query output, please check the following [FAQ](#)

Features **Variant (Germline)**
 Structures **Variant (Somatic)**
 Homologues (Max select 6 orthologues) **Sequences**

GENE:

Ensembl

- Gene stable ID
- Gene stable ID version
- Transcript stable ID
- Transcript stable ID version
- Protein stable ID
- Protein stable ID version
- Exon stable ID
- Gene description
- Chromosome/scaffold name
- Gene start (bp)
- Gene end (bp)
- Strand
- Karyotype band

- GENCODE basic annotation
- APPRIS annotation
- RefSeq match transcript
- Gene name
- Source of gene name
- Transcript name
- Source of transcript name
- Transcript count
- Gene % GC content
- Gene type
- Transcript type
- Source (gene)
- Source (transcript)

Select attributes to be output

Example: Step 5 (get results)

e!Ensembl BLAST/BLAT | VEP | Tools | BioMart | Downloads | Help & Docs | Blog Login/Register

🔍 Search all species...

New **Count** **Results** **URL** **XML** **Perl** **Help**

Dataset 12 / 67130 Genes
Human genes (GRCh38.p13)

Filters
Chromosome/scaffold: 1
Start: 78895
End: 224561

Attributes
Gene stable ID
Gene stable ID version

Dataset
[None Selected]

Export all results to Unique results only

Email notification to


View rows as Unique results only

Gene stable ID	Gene stable ID version
ENSG00000238009	ENSG00000238009.6
ENSG00000239945	ENSG00000239945.1
ENSG00000233750	ENSG00000233750.3
ENSG00000268903	ENSG00000268903.1
ENSG00000269981	ENSG00000269981.1
ENSG00000239906	ENSG00000239906.1
ENSG00000241860	ENSG00000241860.7
ENSG00000222623	ENSG00000222623.1
ENSG00000241599	ENSG00000241599.1
ENSG00000279928	ENSG00000279928.2


Start using Ensembl/BioMart


- We are going to use Ensembl/Biomart for **Ensembl v95**.
- On the main page of Ensembl website, click on **View in archive site**


-- Select a species --

 **Pig breeds**
Pig reference genome and 12 additional breeds

[View full list of all species](#)

 **Human**
GRCh38.p13
[Still using GRCh37??](#)

 **Mouse**
GRCm38.p6

 **Zebrafish**
GRCz11

- 28 Aug 2020: [Cool stuff the Ensembl VEP can do: annotating SARS-CoV-2 variants](#)
- 27 Aug 2020: [Job: Computational Data Analyst](#)

Compare genes across species

Find SNPs and other variants for my gene

```
GTGATACATTCC
CRTRAAAGTCTT
CTTCTAAATTCT
GRAACATTTTCC
```

Gene expression in different tissues



Retrieve gene sequence

```
GCCAGACTTCGGGTGG:
GGGCTGTGGGGGAGC:
GGGCTCTGGCTGGCCT:
AGGGACAGATTTGTGA:
CAGCTCTGGAGCGTTT:
CCCAGTCCAGCTGGCG:
```

Find a Data Display

Use my own data in Ensembl



Ensembl creates, integrates and distributes reference datasets and analysis tools that enable genomics. We are based at [EMBL-EBI](#) and our software and data are freely available.

Our [acknowledgements page](#) includes a list of current and previous funding bodies. [How to cite Ensembl](#) in your own publications.

Start using Ensembl/BioMart

- Choose Ensembl 95: Jan 2019 (GRCh38.p12)

View in archive site

Search

- Help topics
 - Frequently Asked Questions
 - Video Tutorials
 - Glossary

Contact HelpDesk

The following archives are available for this page:

- [Ensembl GRCh37](#): Full Feb 2014 archive with BLAST, VEP and BioMart
- [Ensembl 103: Feb 2021](#) (GRCh38.p13) - patched/updated gene set Aug 2020
- [Ensembl 102: Nov 2020](#) (GRCh38.p13) - patched/updated gene set Sep 2020
- [Ensembl 101: Aug 2020](#) (GRCh38.p13) - patched/updated gene set Mar 2020
- [Ensembl 100: Apr 2020](#) (GRCh38.p13) - patched/updated gene set Jun 2019
- [Ensembl 99: Jan 2020](#) (GRCh38.p13) - patched/updated gene set Aug 2019
- [Ensembl 98: Sep 2019](#) (GRCh38.p13) - patched/updated gene set Jun 2019
- [Ensembl 97: Jul 2019](#) (GRCh38.p12) - patched/updated gene set Mar 2019
- [Ensembl 96: Apr 2019](#) (GRCh38.p12) - patched/updated gene set Nov 2018
- [Ensembl 95: Jan 2019](#) (GRCh38.p12)
- [Ensembl 94: Oct 2018](#) (GRCh38.p12) - patched/updated gene set Jul 2018
- [Ensembl 93: Jul 2018](#) (GRCh38.p12)
- [Ensembl 92: Apr 2018](#) (GRCh38.p12) - patched/updated gene set Jan 2018
- [Ensembl 91: Dec 2017](#) (GRCh38.p10)
- [Ensembl 90: Aug 2017](#) (GRCh38.p10) - patched/updated gene set Jun 2017
- [Ensembl 89: May 2017](#) (GRCh38.p10) - patched/updated gene set Jan 2017
- [Ensembl 88: Mar 2017](#) (GRCh38.p10)
- [Ensembl 87: Dec 2016](#) (GRCh38.p7)
- [Ensembl 86: Oct 2016](#) (GRCh38.p7)
- [Ensembl 85: Jul 2016](#) (GRCh38.p7) - patched/updated gene set Jun 2016
- [Ensembl 80: May 2015](#) (GRCh38.p2) - patched/updated gene set Jan 2015
- [Ensembl 77: Oct 2014](#) (GRCh38) - patched/updated gene set Aug 2014
- [Ensembl 75: Feb 2014](#) (GRCh37.p13) - patched/updated gene set Sep 2013
- [Ensembl 67: May 2012](#) (GRCh37.p7) - patched/updated gene set Feb 2012
- [Ensembl 54: May 2009](#) (NCBI 36) - patched/updated gene set Oct 2008

[More information about the Ensembl archives](#)

Start using Ensembl/BioMart

- Click on Biomart (top menu)

The screenshot shows the Ensembl BioMart website. At the top, there is a navigation bar with the Ensembl logo, 'BioMart' label, and links for 'Downloads', 'Help & Docs', and 'Blog'. A search bar on the right contains the text 'Search all species...'. Below the navigation bar, there is a 'Tools' section with a 'BioMart >' link and a description: 'Export custom datasets from Ensembl with this data-mining tool'. A search box is present with a dropdown menu set to 'All species' and a 'Go' button. Below the search box, there are examples of search terms: 'e.g. BRCA2 or rat 5:62797383-63627669 or rs699 or coronary heart disease'. On the right side, there is a section titled 'Ensembl Archive Release 95 (January 2019)' with a list of updates: 'New regulatory build for human, incorporating new data from ENCODE', 'Update to GENCODE M20 for mouse', 'New genomes: donkey, polar bear, black bear, red fox, koala, dingo, tuatara, painted turtle and desert tortoise', 'Updated genomes for chicken, cow and horse', and 'New protein structure variation view'. Below this, there is a link to 'More release news'. At the bottom right, there is a section for 'Other news from our blog' with a link to a post from 27 Mar 2020 about annotating structural variants. On the left side, there is a section for 'All genomes' with a dropdown menu set to '-- Select a species --' and a link to 'View full list of all Ensembl species'. Next to it is a 'Favourite genomes' section with a profile picture and the text 'Human GRCh38.p12'.

- CHOOSE DATABASE : select “Ensembl Genes 95”
- CHOOSE DATASET : select “Human genes (GRCh38.p12)”

Exercise 1: get annotations of a gene (1/2)

- 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 **v95**, for Human GRCh38.p12
 - Click on Filters :
 - Expand the GENE section
 - Select « Input external references ID list »
 - Select Gene Name(s) in the drop down menu
 - Enter IDH1 in the text box
 - Click on Attributes :
 - Select “Features” (top panel, selected by default)
 - Select Gene stable ID, Transcript stable ID, Gene Name
 - Click on Results

Exercise 1: get annotations of a gene (2/2)

- 2. Extract all exon sequences of the IDH1 gene in fasta format. Headers will contain the Gene names, transcript stable IDs and Exon stable IDs.
- 3. Extract all coding sequences of the IDH1 gene in fasta format. Headers will contain the transcript stable IDs and Exon stable IDs.
- 4. Retrieve GO-terms associated to the IDH1 gene (select GO Term Name, GO domain and GO Term Accession along with Gene stable ID, Transcript stable ID and 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 stable ID, Transcript stable ID and 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 you extract from Ensembl/BioMart**

Exercise 2: get annotations for a set of genes

siMitfvssiLuc.up.txt

mart_export.txt (from Ensembl/Biomart)

Gene stable ID	siLuc2	siLuc...
ENSG00000018408	4685 ...	
ENSG00000081189	1716 ...	
ENSG00000106772	3063 ...	
ENSG00000124942	309 ...	
ENSG00000142871	243 ...	
ENSG00000143341	3760 ...	
ENSG00000154556	352 ...	
ENSG00000185565	679 ...	
ENSG00000163328	136 ...	
ENSG00000064042	1160 ...	
ENSG00000114423	2293 ...	

Gene stable ID	Gene name	Chro...
ENSG00000000971	CFH 1	19665187...
ENSG00000001461	NIPAL3	1 2441...
ENSG00000124942	AHNAK	11 624...
ENSG00000002330	BAD 11	642698...
ENSG00000002549	LAP3 4	175771...
ENSG00000002586	CD99 X	269113...
ENSG00000002834	LASP117	3886...
ENSG00000002919	SNX11	17 4810...
ENSG00000003137	CYP26B1	2 7212...
ENSG00000003436	TFPI 2	187464...
ENSG00000018408	WWTR1	3 1495...




Result file

Gene stable ID	siLuc2	siLuc3	...	Gene name	Chro...
ENSG00000124942	309	...	AHNAK	11	624...
ENSG00000018408	4685	...	WWTR1	3	1495...

Exercise 2: get annotations for a set of genes

If you encountered any trouble with the generation of the dataset

- go to GalaxEast (<http://use.galaxeast.fr>)
 - go to Shared Data/ Data Libraries / NGS data analysis training / RNAseq / statistical_analysis.
 - Import the dataset SARTools_DESeq2_tables to your history.
1. Click on  to display the content of the dataset [SARTools DESeq2 table](#) and download the file siMitfvssiLuc.up.txt (click right, save ...)
 2. Open the file siMitfvssiLuc.up.txt and change the name of the column which contains “Id” to “Gene stable ID” (first word, first line). Save the change.
 3. Use the file siMitfvssiLuc.up.txt to extract gene annotations for those genes. Annotation to extract are : gene stable IDs, Chromosome/scaffold name, Gene start, Gene end, strand, Gene name, Gene type. Save the results to a compressed TSV file. (don't close the Ensembl/Biomart window once done)
 - Tip: columns are in the same order as columns are selected
 4. Upload the file siMitfvssiLuc.up.txt and the annotation file (mart_export.txt.gz) 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 (**siMitfvssiLuc.up.txt** **then** **mart_export.txt**) based on the “Gene stable IDs” field.
 - Gene stable 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” (<- this exact case) in the dataset siMitfvssiLuc.up.annot.txt.
 - Tip 1: Search “lincRNA” in the column containing Gene types
 - Tip 2: c3 refers to column 3 of a dataset.
 - Tip3 : look at examples below the form to help you find the correct syntax

Exercise 2: get annotations for a set of genes

- Bonus question: go back to Ensembl/BioMart. You want to extract sequences of all promoters of the up-regulated genes (the ones from the file siMitfvssiLuc.up.txt) to run a *de novo* motif discovery and search for over represented nucleotide sequence. Retrieve the 200nt upstream of these genes. Header should contain Gene stable ID, Transcript stable ID, Gene name and Gene description.

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 stable ID, Chromosome/scaffold name, Gene Start (bp), Gene End (bp), strand and Gene Name