**Data Upload Exercises**

These exercises are designed to teach you how to upload your own data files to the website (http://www.ensembl.org).  A range of upload options are explored, including GFF, BigWig, BAM and BED file formats.  Sample upload files are provided.

Note: the answers to these exercises were composed in release version 68: http://e68.ensembl.org/index.html.

Please report any discrepancies with more current versions to helpdesk@ensembl.org.

_____

# CUSTOM ANNOTATION
_____

### Exercise 1 – Attaching a GFF file

Have a look at the following file:

http://www.ebi.ac.uk/~bert/n-scan_genes.gff

It contains annotations for three transcripts of the human *HFE* gene (ENSG00000010704) generated by the N-SCAN gene structure prediction software, as shown on the UCSC Genome Browser (http://genome.ucsc.edu/cgi-bin/hgTracks?position=chr6:26087509-26095469&knownGene=pack&nscanGene=pack).

The file is in GFF (General Feature Format) format:

http://www.ensembl.org/info/website/upload/gff.html

Attach the file to Ensembl and have a look at the result.
_____

*Answer*

🖰 Go to the Ensembl homepage (http://www.ensembl.org/).
🖰 Click on the picture of human (Donatello's St. George) or the word 'Human' next to it.
🖰 Click [Manage your data] in the side menu.
🖰 Click on 'Attach Remote File'.
🖰 Enter the URL of the file in the 'File URL' text box.
🖰 Select 'Data format: GFF'.
🖰 Enter 'N-SCAN genes' in the 'Name for this track' text box.
🖰 Click [Next>].

🖰 Click on 'Go to first region with data: 6:26037670-26137670'.

A new track named 'N-SCAN genes' should now have been added to the 'Region in detail' page.

You may want to turn off all tracks that you added to the display in the previous exercises.

🖰 Click [Configure this page] in the side menu.
🖰 Click [Reset configuration].

To display the names of the N-SCAN genes.

🖰 Click on 'Your data'.
🖰 Select 'N-SCAN genes – Labels'.
🖰 Click (✓).

Note that, at the moment, the CDS information in the GFF file is not taken into account in Ensembl and thus no distinction between the UTRs and CDS of the transcripts can be seen.

────────────────────────────────────────────────

**Exercise 2 – Attaching a BigWig file**

The *BCL11A* (B-cell CLL/lymphoma 11A (zinc finger protein)) gene functions as a myeloid and B-cell proto-oncogene.

The files

http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeCaltechRnaSeq/wgEncodeCaltechRnaSeqGm12878R2x75Th1014Il200SigRep1V4.bigWig

and

http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeCaltechRnaSeq/wgEncodeCaltechRnaSeqK562R2x75Th1014Il200SigRep1V4.bigWig

contain RNA-Seq data for the GM12878 and K562 cell lines, respectively.

The files are in BigWig format:

https://cgwb.nci.nih.gov/goldenPath/help/bigWig.html

Attach both files to Ensembl and have a look at the result. Is the *BCL11A* gene expressed in both cell lines?

────────────────────────────────────────────────

*Answer*

🖱 Go to the Ensembl homepage (http://www.ensembl.org).
🖱 Select 'Search: Human' and type 'bcl11a' in the 'for' text box.
🖱 Click [Go].
🖱 Click on 'Gene' on the page with search results.
🖱 Click on 'Human'.
🖱 Click on '2:60678302-60780702:-1'.

You may want to turn off all tracks that you added to the display in the previous exercises.

🖱 Click [Configure this page] in the side menu.
🖱 Click [Reset configuration].
🖱 Click (✓).

🖱 Click [Manage your data] in the side menu.
🖱 Click on 'Attach Remote File'.
🖱 Enter the URL of the first file in the 'File URL' text box.
🖱 Select 'Data format: BigWig'.
🖱 Enter 'GM12878_RNAseq' in the 'Name for this track' text box.
🖱 Click [Next>].
🖱 Click [Save].
🖱 Repeat for the second file.
🖱 Click (✓).

The *BCL11A* gene is expressed in the GM12878 cell line, while there is virtually no expression in the K562 cell line. Note that the vertical scale differs between the two attached RNA-Seq tracks.

_____

**Exercise 3 – Attaching a BAM file**

The following file contains alignments to the GRCh37 assembly of low coverage Illumina sequencing reads of chromosome 20 of individual HG00096 from the 'British from England and Scotland, UK' cohort (http://ccr.coriell.org/Sections/Search/Sample_Detail.aspx?Ref=HG00096&PgId=166):

http://www.ebi.ac.uk/~bert/HG00096.chrom20.ILLUMINA.bwa.GBR.low_coverage.20100901.bam

The file is in BAM format. BAM is the compressed binary version of the SAM (Sequence Alignment/Map) format, a compact and indexable representation of nucleotide sequence alignments:

http://samtools.sourceforge.net/SAM1.pdf

To display these data in Ensembl also the .bam.bai index file is needed:

http://www.ebi.ac.uk/~bert/HG00096.chrom20.ILLUMINA.bwa.GBR.low_coverage.20100901.bam.bai

The .bam.bai file should be placed in the same directory as the .bam file.

Attach the file to Ensembl and have a look at the result. Can you find any individual reads containing a nucleotide that differs from the sequence of the reference genome? And a position where individual HG00096 differs from the reference genome or where individual HG00096 is heterozygous?

─────────────────────────────────────────────────────────

*Answer*

⼿ Go to the Ensembl homepage (http://www.ensembl.org/).
⼿ Click on the picture of human (Donatello's St. George) or the word 'Human' next to it.
⼿ Click on 'Sample entry points – Karyotype' in the side menu.
⼿ Click on chromosome 20 in the karyotype.
⼿ Click on 'Jump to location View' in the pop-up menu.

You may want to turn off all tracks that you added to the display in the previous exercises.

⼿ Click [Configure this page] in the side menu.
⼿ Click [Reset configuration].
⼿ Click (✓).

⼿ Click [Manage your data] in the side menu.
⼿ Click on 'Attach Remote File'.
⼿ Enter the URL of the file in the 'File URL' text box.
⼿ Select 'Data format: BAM'.
⼿ Enter 'HG00096' in the 'Name for this track' text box.
⼿ Click [Next>].
⼿ Click (✓).

A new track named 'HG00096' should have been added to the 'Region in detail' page.

⼿ Zoom in to see the actual reads.

Individual reads are shown in grey, with the consensus sequence shown above the reads in colour.

Nucleotides that differ from the sequence of the reference genome are shown in red:

http://www.ensembl.org/Homo_sapiens/Location/View?db=core&r=20:4486120 7-44861246

An example of a position where individual HG00096 is heterozygous:

http://www.ensembl.org/Homo_sapiens/Location/View?db=core&r=20:4485470 6-44854746

_____


## Exercise 4 – Creating an annotated karyotype

This is a list of all human caspase genes:

*CASP1*, *CASP2*, *CASP3*, *CASP4*, *CASP5*, *CASP6*, *CASP7*, *CASP8*, *CASP9*, *CASP10*, *CASP12*, *CASP14*

Create a figure of the human karyotype showing the genomic position of the caspase genes.

_____


### *Answer*

⌐ Go to the Ensembl homepage (http://www.ensembl.org/).
⌐ Click on the picture of human (Donatello's St. George) or the word 'Human' next to it.
⌐ Click [Manage your data] in the side menu.
⌐ Click on 'Features on Karyotype'.
⌐ Enter the list of caspase genes in the 'ID(s)' text box.
⌐ Click [Show features].

The positions of the caspase genes should now be shown in the karyotype by red triangles. Note that some of the genes (on chromosome 2 and 11) are so close to each other that they cannot be shown by separate triangles.

_____


## Exercise 5 – Creating and uploading a BED file

Create a small text file containing some annotation in BED format (http://www.ensembl.org/info/website/upload/bed.html) and upload it to Ensembl.

Note that BED offers the simplest format, with only three required fields, i.e. chromosome, start and end.

_____

*Answer*

🖰 Create a text file with your annotation in for example Notepad or TextEdit and save it on your computer.
🖰 Go to the Ensembl homepage (http://www.ensembl.org/).
🖰 Select your favourite species.
🖰 Click [Manage your data] in the side menu.
🖰 Click on 'Upload Data'.
🖰 Enter the name for your track in the 'Name for this upload (optional)' text box.
🖰 Select 'Data format: BED'.
🖰 Click [Choose File] behind 'Upload file:'.
🖰 Select the text file you just created.
🖰 Click [Upload].
🖰 Click 'Go to first region with data:'.

Your data should now be shown as a new track on the 'Region in detail' page.

_____

### Exercise 6 – Removing custom annotation

Remove your attached and uploaded annotations.

_____

*Answer*

🖰 Go to the Ensembl homepage (http://www.ensembl.org/).
🖰 Click on the picture of human (Donatello's St. George) or the word 'Human' next to it.
🖰 Click [Manage your data] in the side menu.
🖰 Click for each added data set on the trash can icon.
🖰 Click (✓).

Your annotations should be removed now.

_____