# Topic 9.3: Genomes

# Genome Projects

Currently there are ~500,000 organisms with genomes sequenced or being sequenced.

Mostly – ~400,000 – bacteria

~20,000 viruses

~50,000 Eukarya

~200 animals

https://gold.jgi.doe.gov/index

# Comparative Genomics

By comparing genomes of related animals, we may be able to see what sequences are conserved

5-10% of our genome is conserved cf. other animals

Remember: <1% of our genome is "coding DNA"

(exons that code for protein)

→ a LOT of non-coding functional sequence information for chromosome packaging, segregation, and replication; non-coding RNA, and gene regulation

# Viewing Genomes

Sites that allow us to view various genomes:

- NCBI: Genome Data Viewer

- UC Santa Cruz Genome Browser

- JBrowse http://jbrowse.org/

- Ensemble's Genomes

The free Integrative Genomics Viewer (IGV) program allows us to access the same data but runs as a standalone program.

Other web sites:

The gnomAD web site allows us to look at variations in 120,000 genomes.

The Galaxy web site has one of the largest collection of algorithms for analyzing genomes

# Data in Genome Browsers

Besides the actual genomic sequences themselves, these genome browsers usually also have

- RefSeq genes (mRNA and protein)

- Ensembl genes

- Expressed sequence tags

- Gene predictions

- SNPs from a database like dbSNP

- Non-coding functional elements as found by projects like ENCoDE

# Human genome data in the UCSC Genome Browser

The University of California Santa Cruz genome browser was set up to make the human genome accessible to the general public

Go to:

http://genome.ucsc.edu/cgi-bin/hgTracks?db=hg38

# UCSC Genome Browser

We will be shown data for a a particular chromosome and at a particular location:

We cannot see indivdual nucleotides because there are a lot of them in a chromosome.

Notice all the tracks of data

# UCSC Genome Browser Tracks

The first track displays the location (in bases) in a particular chromosome and a scale

# UCSC Genome Browser Tracks

The first track displays the location (in bases) in a particular chromosome and a scale

# UCSC Genome Browser Tracks

There are usually tracks displaying information about genes in this area of the chromosome, including, if available, curated RefSeq data



In this case, we see that this particular gene is on the complementary strand

Clicking on the gene name will take us to a page displaying more details about the gene

# UCSC Genome Browser Tracks

Clicking on the gene name will take us to a page displaying more details about the gene

## Human Gene POMC (ENST00000380794.5) from GENCODE V41

**Description:** ACTH stimulates the adrenal glands to release cortisol. (from UniProt P01189)

**RefSeq Summary (NM_001035256):** This gene encodes a preproprotein that undergoes extensive, tissue-specific, post-translational processing via cleavage by subtilisin-like enzymes known as prohormone convertases. There are eight potential cleavage sites within the preproprotein and, depending on tissue type and the available convertases, processing may yield as many as ten biologically active peptides involved in diverse cellular functions. The encoded protein is synthesized mainly in corticotroph cells of the anterior pituitary where four cleavage sites are used; adrenocorticotrophin, essential for normal steroidogenesis and the maintenance of normal adrenal weight, and lipotropin beta are the major end products. In other tissues, including the hypothalamus, placenta, and epithelium, all cleavage sites may be used, giving rise to peptides with roles in pain and energy homeostasis, melanocyte stimulation, and immune modulation. These include several distinct melanotropins, lipotropins, and endorphins that are contained within the adrenocorticotrophin and beta-lipotropin peptides. The antimicrobial melanotropin alpha peptide exhibits antibacterial and antifungal activity. Mutations in this gene have been associated with early onset obesity, adrenal insufficiency, and red hair pigmentation. Alternatively spliced transcript variants encoding the same protein have been described. [provided by RefSeq, Jan 2016].

**Gencode Transcript:** ENST00000380794.5
**Gencode Gene:** ENSG00000115138.11
**Transcript (Including UTRs)**
   **Position:** hg38 chr2:25,160,853-25,168,690 **Size:** 7,838 **Total Exon Count:** 4 **Strand:** -
**Coding Region**
   **Position:** hg38 chr2:25,161,081-25,164,772 **Size:** 3,692 **Coding Exon Count:** 2

# UCSC Genome Browser Tracks

There may be an OMIM track displaying information about allelic variations



The bars represent variations that may be important in determining phenotypes or inherited diseases

# UCSC Genome Browser Tracks

A Conservation track displays best guesses for aligning this stretch of human genomic DNA with those of other vertebrates:

# UCSC Genome Browser Tracks

The last of the default tracks is often the RepeatMasker track showing any detected repeat sequences, which are mostly in non-coding parts of the genome.



Repeat sequences are genomic DNA with short and long interspersed nuclear elements (SINE and LINE), long terminal repeat elements (LTR), and other "low complexity" areas with short repeating sequences.

# UCSC Genome Browser Tracks

RepeatMasker track: Click on the "RepeatMasker" link to reveal computationally detected repeat sequences.



Here we can get more details of repeat sequences - genomic DNA with

- short and long interspersed nuclear elements (SINE and LINE)

- long terminal repeat elements (LTR)

- other "low complexity" areas with short repeating sequences.

When performing sequence alignments, we can disable these repeat sequences from appearing in our alignments.

# Zebra Mussel data in the NCBI Genome Data Viewer

Go to: https://www.ncbi.nlm.nih.gov/genome/gdv/

# Zebra Mussel data in the NCBI Genome Data Viewer

Go to: https://www.ncbi.nlm.nih.gov/genome/gdv/

# Zebra Mussel data in the NCBI Genome Data Viewer

## In the "Search organisms" box, enter Zebra Mussel

# Zebra Mussel data in the NCBI Genome Data Viewer

It will show up as "Dreissena polymorpha"

# Zebra Mussel data in the NCBI Genome Data Viewer

Genome Data Viewer will show us some of the details of Chromosome 1:



16 chromosomes plus a mitochondrial genome

Chromosome 1 is selected.
The view defaults to all 211 million bases.

Use the zoom controls to select less than about a million bases

# Zebra Mussel data in the NCBI Genome Data Viewer



We can see exons and introns

A lot of the genes are marked as "mRNA hypothetical" because they were detected using computational methods

At this resolution, we will see a lot more tracks

# Zebra Mussel data in the NCBI Genome Data Viewer

We can download the nucleotides in this view as a Fasta file. Click on Download → Download Fasta → Fasta (Visible Range)

# Zebra Mussel data in the NCBI Genome Data Viewer

We will see a window asking us where to save a file with a name that contains the chromosome number and the range of nucleotides.



We can open up the file using a text editor like Notepad or Notepad++

# Viewing gene promoter data in the Ensemble site

Go to:  http://ensembl.org/index.html

# Promoter data in Ensemble

Go to:

http://ensembl.org/index.html

Click "Human"

# Promoter data in Ensemble

A lot of options are presented



Click "BRCA2" or search for a particular gene

# Promoter data in Ensemble

Click the first search result

We will get a page with
a lot of data:

Click on "Region in detail"

# Promoter data in Ensemble

We will get a page with a view of the chromosome:



Chromosome 13: 32,313,383-32,401,971

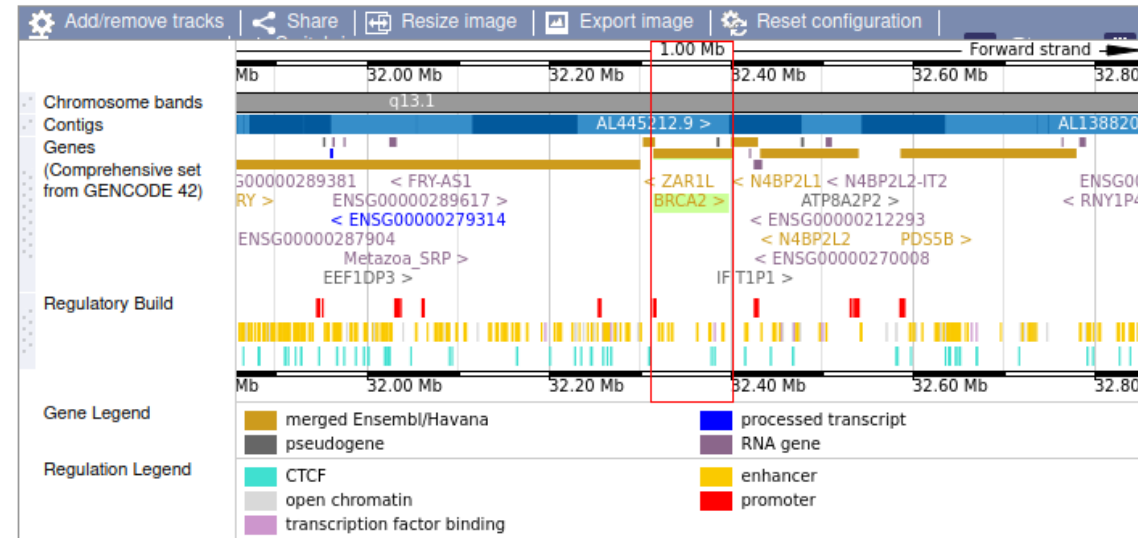A red box shows the part of the chromosome for which we will get a "region in detail" view:

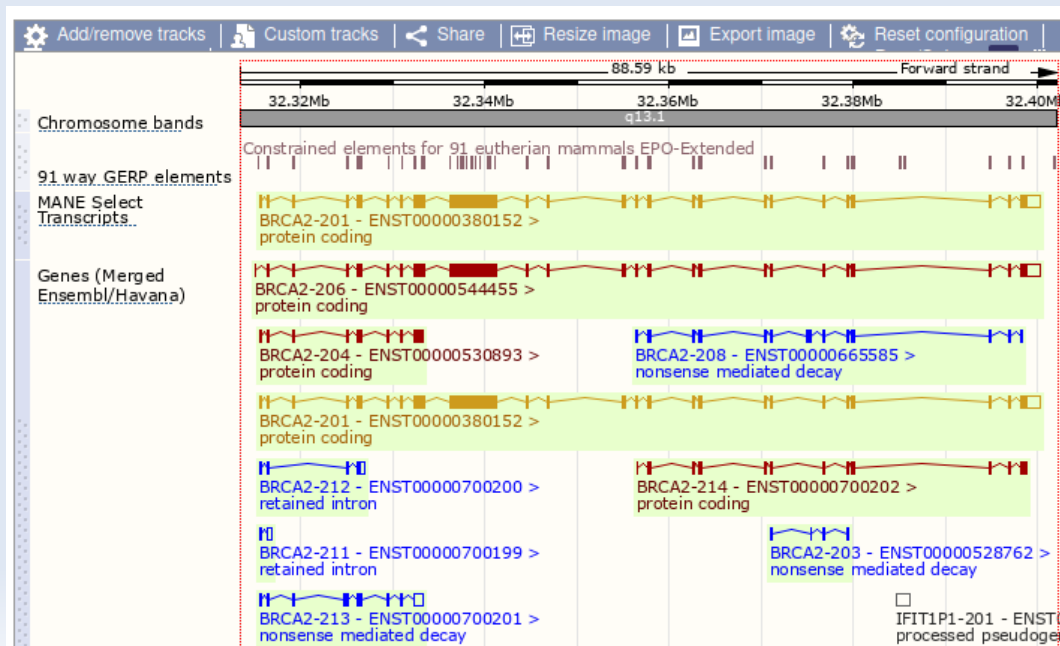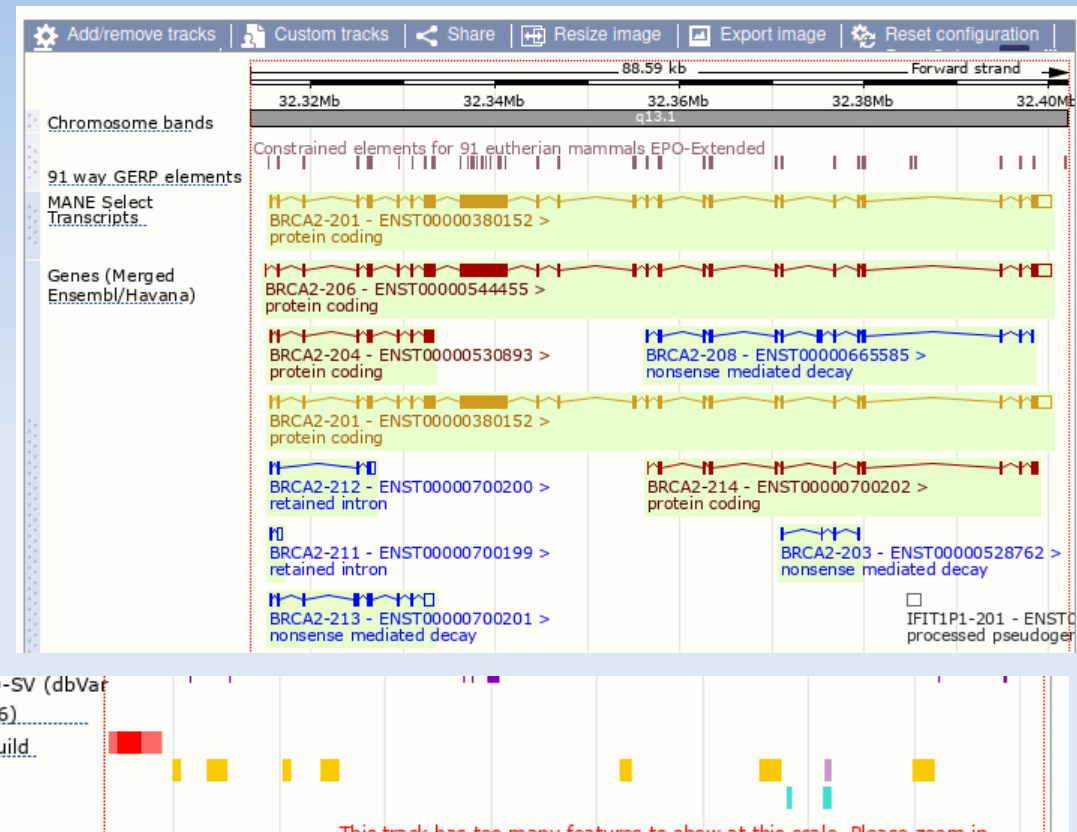# Promoter data in Ensemble

Below the Region
in detail view



we see more details of the gene:

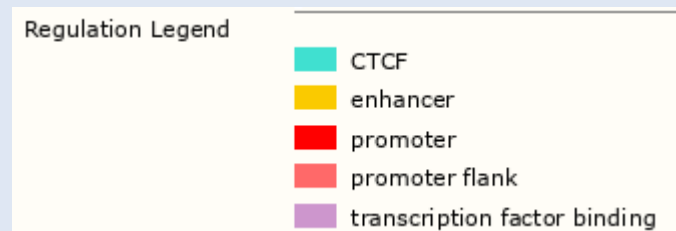# Promoter data in Ensemble

And below the gene view

In the Regulatory Build track, we see

a promoter region and promoter flanks

# Promoter data in Ensemble

## Click the promoter

In the Regulatory Build track



And in the pop-up window, choose the link for "Bounds" for a closer view



Zoom in for a closer view



and eventually, we will be able to see nucleotides: