# Gene Expression Matrix

## ENCODE AWG, Nov 13 2015

Sarah Djebali

Anna Vlasova

# Gene Expression Matrix

- STAR-RSEM pipeline, gencode v.19

- TPM and FPKM values for each replicate

- Tabular-separated (tsv) and Json formats

- long RNA-seq (>200nt)


- Total number of experiments =163

- Total number of bioreplicates = 320

- Total number of genes = 58,540

- Labs: Thomas Gingeras, Barbara Wold, Brenton Graveley

http://genome.crg.es/~sdjebali/STAR-RSEM/geneid_genename_with_tpmallrep_fpkmallrep.tsv.gz
http://genome.crg.es/~sdjebali/STAR-RSEM/geneid_genename_with_tpmallrep_fpkmallrep.json.g

# Json format

{ "gene_name": "SEC62", "ensembl_id": "ENSG00000008952.12",
  "expression_values": [
   { "dataset": "ENCSR000AAA",
    "rep1_tpm": 10.70, "rep2_tpm": 4.78, "rep1_fpkm": 26.43, "rep2_fpkm": 21.64 },
      { "dataset": "ENCSR000AAB",
  "rep1_tpm": 10.51, "rep2_tpm": 1.40, "rep1_fpkm": 24.57, "rep2_fpkm": 9.69 },
{ "dataset": "ENCSR000AAC",
 "rep1_tpm": 3.40, "rep2_tpm": 10.81, "rep1_fpkm": 17.47, "rep2_fpkm": 38.94 },
{ "dataset": "ENCSR000AAD",
 "rep1_tpm": 1.04, "rep2_tpm": 3.11, "rep1_fpkm": 9.34, "rep2_fpkm": 15.16 },
 { "dataset": "ENCSR000AAE",
 "rep1_tpm": 7.32, "rep2_tpm": 2.44, "rep1_fpkm": 26.35, "rep2_fpkm": 13.35 }....

# Gene Expression Matrix

| Lab | # Experiments/ Bioreplicates | Fraction (experiments) | | | Preparation (experiments) | | |
|---|---|---|---|---|---|---|---|
| | | Whole cell | Nucleus | Cytosol | Total | polyA+ | nonPolyA+ |
| Barbara Wold | 15 / 30 | 15 | 0 | 0 | 13 | 2 | 0 |
| Brenton Graveley | 9 / 18 | 9 | 0 | 0 | 9 | 0 | 0 |
| Thomas Gingeras | 139 / 272 | 131 | 8 | 8 | 107 | 21 | 11 |
| **All** | **163 / 320** | **155** | **8** | **8** | **129** | **23** | **11** |

**Grouping by biosample_type**

primary cells 81
immortalized cell lines 47
tissues 23
in vitro differentiated cells 9
induced pluripotent stem cell lines 2
stem cells 1

**Experiments can also be summarized**

- organ_slims
- system_slims
- develomental_slims

# Gene Expression Matrix

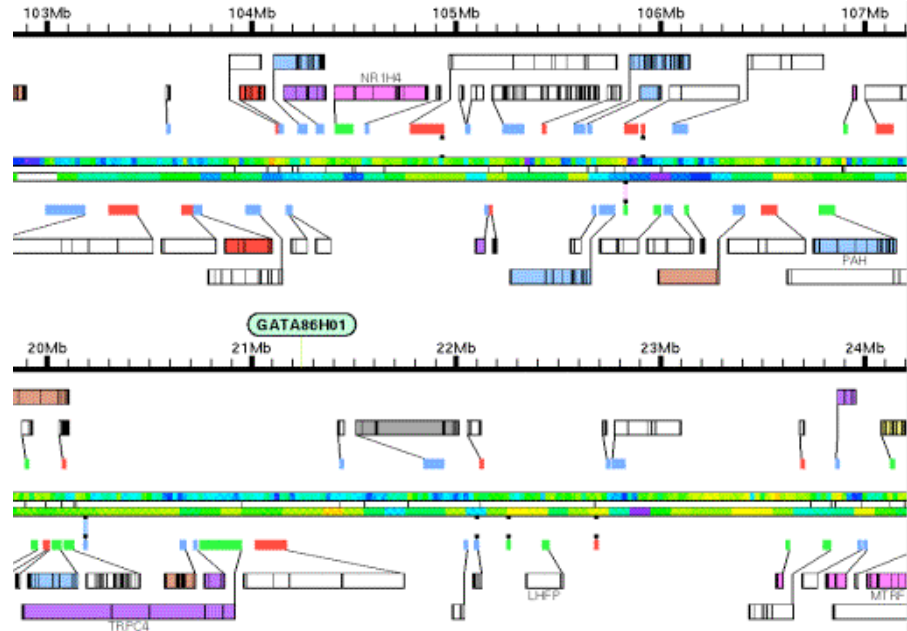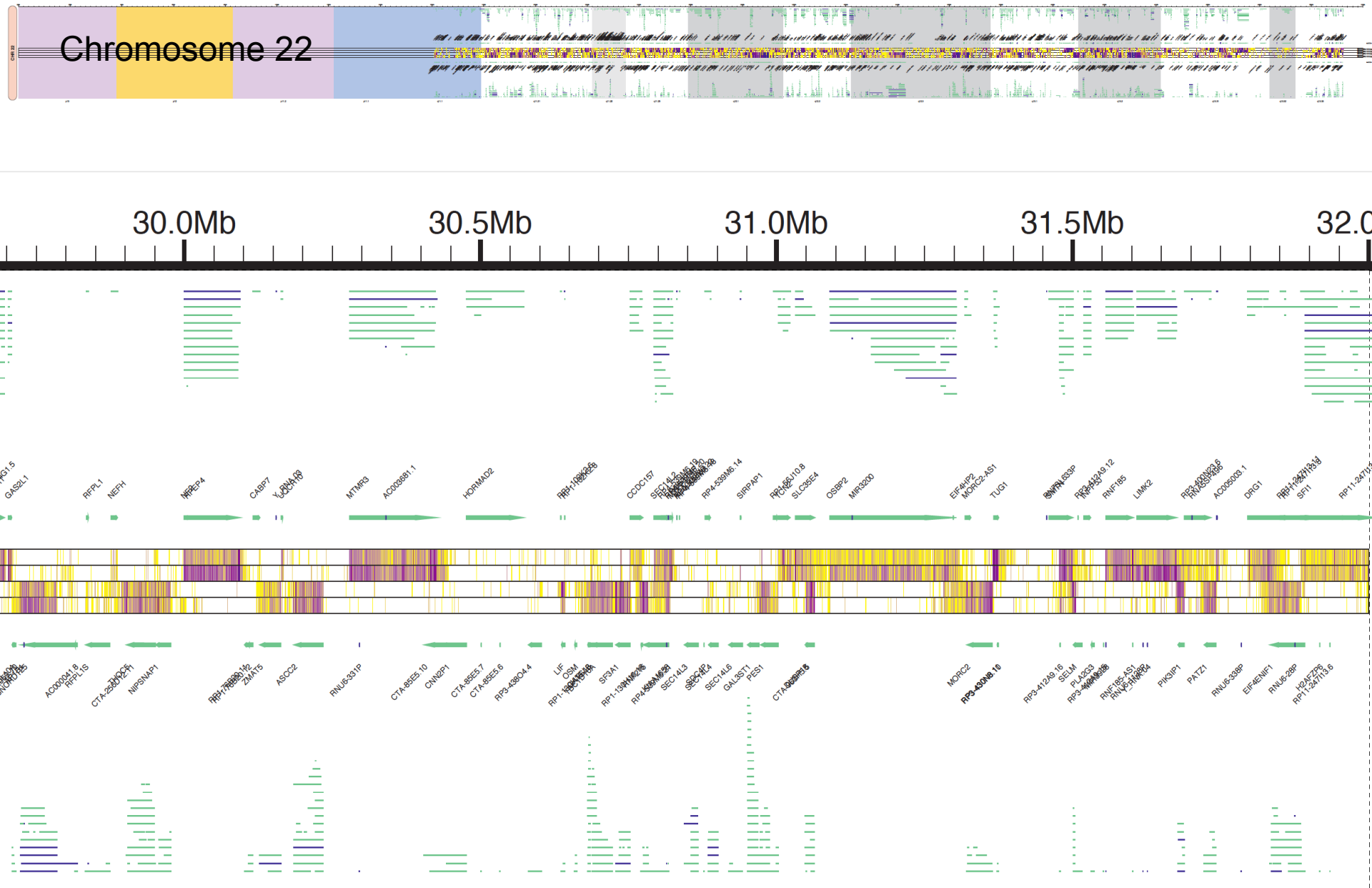ANNOTATION OF THE CELERA HUMAN GENOME ASSEMBLY

# ANNOTATION OF THE CELERA HUMAN GENOME ASSEMBLY

Chromosome 22

30.0Mb  30.5Mb  31.0Mb  31.5Mb  32.0

q12.2

# Gene based expression plots
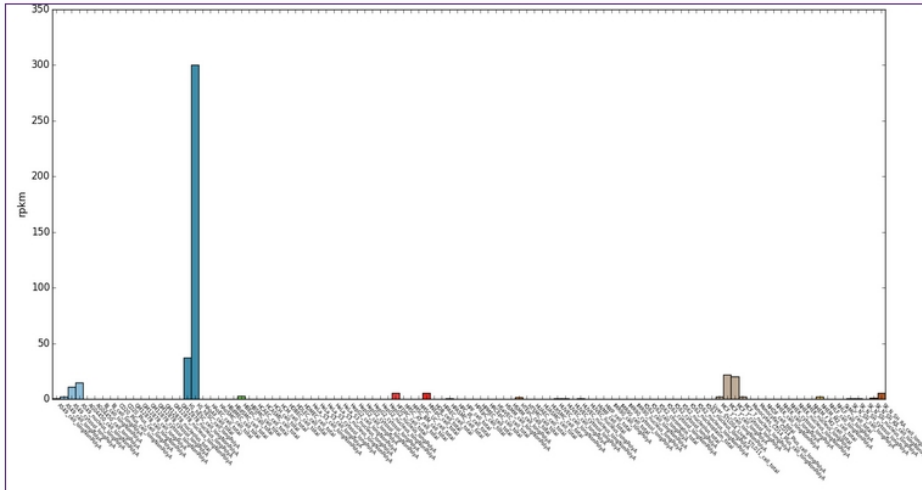
# GTEx

# ENCODE browser

Welcome to YUE Lab
Computational and Functional Genomics/Epigenomics

ABOUT | MOUSE | HUMAN | DOWNLOAD | LINKS | CONTACT

**Gene Expression Across Multiple Tissues/Cell Lines.**

Human (hg19)

Gene **SOX2** [NM_003106, ENSG00000181449, ENST00000325404]

Save as CSV

| | |
|---|---|
| A549_cell_longNonPolyA | 0.495852 |
| A549_cell_longPolyA | 2.30451 |
| A549_cytosol_longPolyA | 10.8769 |
| A549_nucleus_longPolyA | 14.639 |
| AG04450_cell_longNonPolyA | 0 |
| AG04450_cell_longPolyA | 0 |
| BJ_cell_longNonPolyA | 0 |
| BJ_cell_longPolyA | 0 |
| CD20_Plus_cell_longNonPolyA | 0 |
| CD20_Plus_cell_longPolyA | 0 |
| GM12878_cell_longNonPolyA | 0 |
| GM12878_cell_longPolyA | 0.045262 |
| GM12878_cytosol_longNonPolyA | 0 |
| GM12878_cytosol_longPolyA | 0.028214 |
| GM12878_nucleolus_total | 0.095877 |
| GM12878_nucleus_longNonPolyA | 0.113634 |
| GM12878_nucleus_longPolyA | 0.05033 |
| H1_hESC_cell_longNonPolyA | 37.3365 |

http://promoter.bx.psu.edu/ENCODE/get_human_expr.php?assembly=hg19&gene=Sox2

# BioGPS

# Expression Atlas

https://www.ebi.ac.uk/gxa/experiments/E-MTAB-2706

# All samples. Absolute TPMs
# **Sorted**: by sample+bioreplicate name. **Colored**: by organ

All samples. Absolute TPMs
**Sorted**: by sample+bioreplicate name. **Colored**: by organ

ENSG00000000003.10, TSPAN6

# Gingeras samples. Log10 TPMs
**Sorted**: by sample+bioreplicate name.
**Colored**: by biosample type. **Faceted**: by preparation protocol



ENSG00000000003.10, TSPAN6

dataset.biosample_type
- **immortalized cell line**
- **in vitro differentiated cells**
- **induced pluripotent stem cell line**
- **primary cell**
- **stem cell**
- **tissue**

# All samples. Log10 TPMs
**Sorted**: by organ. **Grouped**: by organ. **Colored**: by organ

All samples. Log10 TPMs
**Sorted**: by median gene expression in organs.
**Grouped**: by organ. **Colored**: by organ

GTEX portal

# All samples Absolute TPMs
**Sorted**: by meidan gene expression in system.
**Grouped**: by system. **Colored**: by system

GTEX portal

All samples. Log10 TPMs
**Sorted**: by median gene expression by organ.
**Grouped**: by organ. **Colored**: by organ.
**Faceted**: by extraction protocol

All samples. Log10 TPMs
**Sorted**: by median gene expression by organ.
**Grouped**: by organ. **Colored**: by biosample type.
**Faceted**: by RNA fraction

ENSG0000000003.10

Gene based sunburst plots

extraction
- polyA−
- polyA+
- totalRNA

ENSG00000000003.10

organ_slims
- blood vessel
- blood vessel_heart
- blood vessel_skin of body
- blood vessel_urinary bladder
- bone element
- brain
- bronchus
- extraembryonic structure
- extraembryonic structure_placenta
- eye
- heart
- kidney
- liver
- lung
- lymphatic vessel_skin of body
- mammary gland
- mouth_tongue
- muscle organ
- na
- skin of body
- spinal cord
- stomach
- thyroid gland
- trachea
- urinary bladder

# GENOOM
## ENCODE PROJECT

LOREM LOREM LOREM

## ON BODY GENE EXPRESSION
CH1

CH2
CH3
CH1

45,889,387-80,087,217 34,197,831 bp.

Internal yugular vein
Subclavian artery

Pulmonary Trunck
Ascending aorta
Pulmonary veins

LOAD DATA

SAVE

SEND

NOTES

COLOR RANGE

STATISTICS

CH1    45,889,387-80,087,217 34,197,831 bp.

45,898,698    45,898,698

OneBigRobot

# GENOOM
ENCODE PROJECT

LOREM    LOREM    LOREM

## KARIOTYPE DATA NAME

Lorem Ipsum sid amet
dolor est.

⚙ CH1

## PROXIMITY HIERACY

CH1 45,889,387-80,087,217 34,197,831 bp.

⚙

CODING
NON CODING
SELECTED

## CONSOLE →

LOAD DATA   SAVE   SEND   NOTES

COLOR RANGE   STATISTICS

## 🖌 COLOR RANGE

SEQUENTIAL

DIVERGING

QUALITATIVE

#FF55CC

86%
48%
30%

## CRITERIA  Lorem ipsum sid dolor amet Lorem ipsum sid dolor amet

CH1   45,889,387-80,087,217 34,197,831 bp.

🔖 45,898,698          🔖 45,898,698

OneBigRobot

# Ben Fry

## Genetics projects



### The Genetic Code

A redesign of traditional diagrams of the genetic code to clarify and highlight patterns in the data. Includes an interactive version that depicts how the code works. (2001, updated 24 February 2008)

### Isometric Haplotype Blocks

Combination of several different representations of haplotype data into a single interactive tool. (2001, updated February 2004)
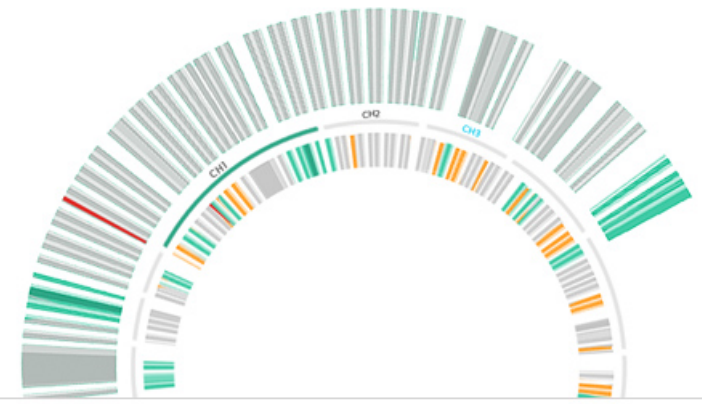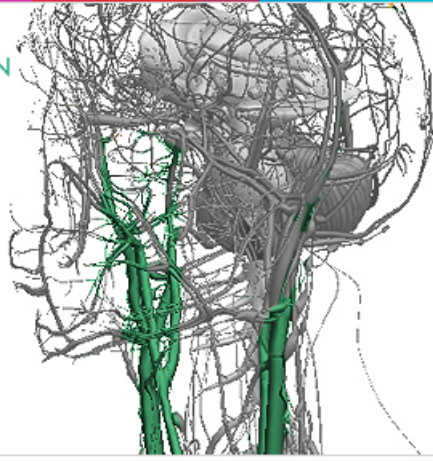
### Genome Valence

A later adaptation of the Valence project that visualizes biological data, and was created for the 2002 Whitney Biennial. (March 2002, updated November 2003)

### Handheld Genome Browser

The biologist's calculator: a genome browser that runs on a handheld device. (2001)

## Tools

### Bifurcator, 2004

This tool creates a bifurcation plot suitable for publication from a set of haplotype data. Given a set of SNPs that define a "core" region, the program creates an image of how individual genotypes differentiate from that point.

### Microarray Clustering with CAST, December 2000

Implementation of the CAST algorithm to cluster microarray data, developed for a class project

## Illustrations

### Aligning Humans + Mammals, December 2007

Sequences of human DNA aligned with about a dozen other mammals, created as an illustration for *Seed Magazine*.

### Nature HapMap Cover, October 2005

Cover for the journal Nature, announcing the completion of the first phase of the HapMap project.

### Humans vs. Chimps, October 2005

An illustration of how the gene FOXP2, believed to be connected to language acquisition, differs in humans versus chimps.

### Axonometric Introns & Exons, January 2003

Large format (9 x 18 feet) print of all the known and predicted genes in the human genome

### Chromosome 21, January 2003

Installation depicting thirteen million letters (one quarter) of human chromosome 21, colored by their use

### Chromosome 14, May 2001

Poster depicting all the genes of chromosome 14 in the human genome

### Chromosome 22, April 2001

The A, C, G, and T letters of human chromosome 22 shown in a three pixel font

# acknowledgements

- Sarah Djebali

- Anna Vlasova

- Julien Lagarde

- Didac Santesmasses (sunburst plots)

- Josep F. Abril (genome plots, University of Barcelona)

- Griselda Serra (human body browser, OneBigRobot)