# Biological Data Resources

Benjamin Hitz, PhD
4/1/2015 CBPSS Panel
http://cherylab.stanford.edu

- The original model organism database for the first eukaryotic genome (*Saccharomyces cerevisiae)* to be sequenced

- Stores genes and biological annotations extracted from scientific literature *via* manual curation

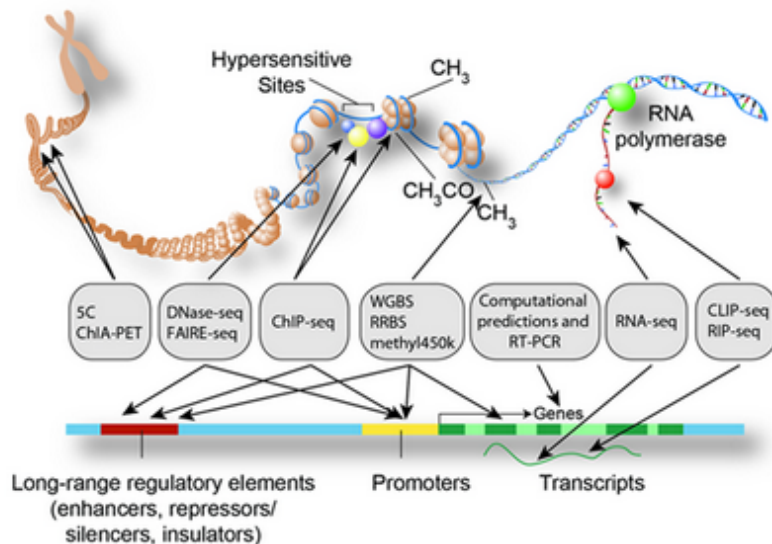http://www.yeastgenome.org/ … since 1993

# SGD Annotations

GO: FUS3 encodes a MAP protein kinase as shown by direct enzyme assay in Bao et al. (2004).

Phenotype: A null mutant of VAC14 has abnormal vacuolar morphology in S288C strain background as shown in Alghamdi et al. (2013).

| Bioentity (what is it) | Bioconcept (what it does) | Reference (who said it) | Experiment (how is it known) | Strain (details) |
|---|---|---|---|---|
| FUS3 | GO: MAP kinase activity | PMID:15620357 | Direct enzyme assay | |
| VAC14 | Phenotype: Vacuolar morphology | PMID: 23389034 | Classical genetics | S288C |

Additional properties (e.g., allele, conditions, qualifier) can be attached to any annotation.

# ENCODE: Encyclopedia of DNA Elements



The ENCODE (Encyclopedia of DNA Elements) Consortium is an international collaboration of research groups funded by the National Human Genome Research Institute (NHGRI). The goal of ENCODE is to build a comprehensive parts list of functional elements in the human genome, including elements that act at the protein and RNA levels, and regulatory elements that control cells and circumstances in which a gene is active.

*Image credits: Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)*

## Data

To find and download ENCODE Consortium data:
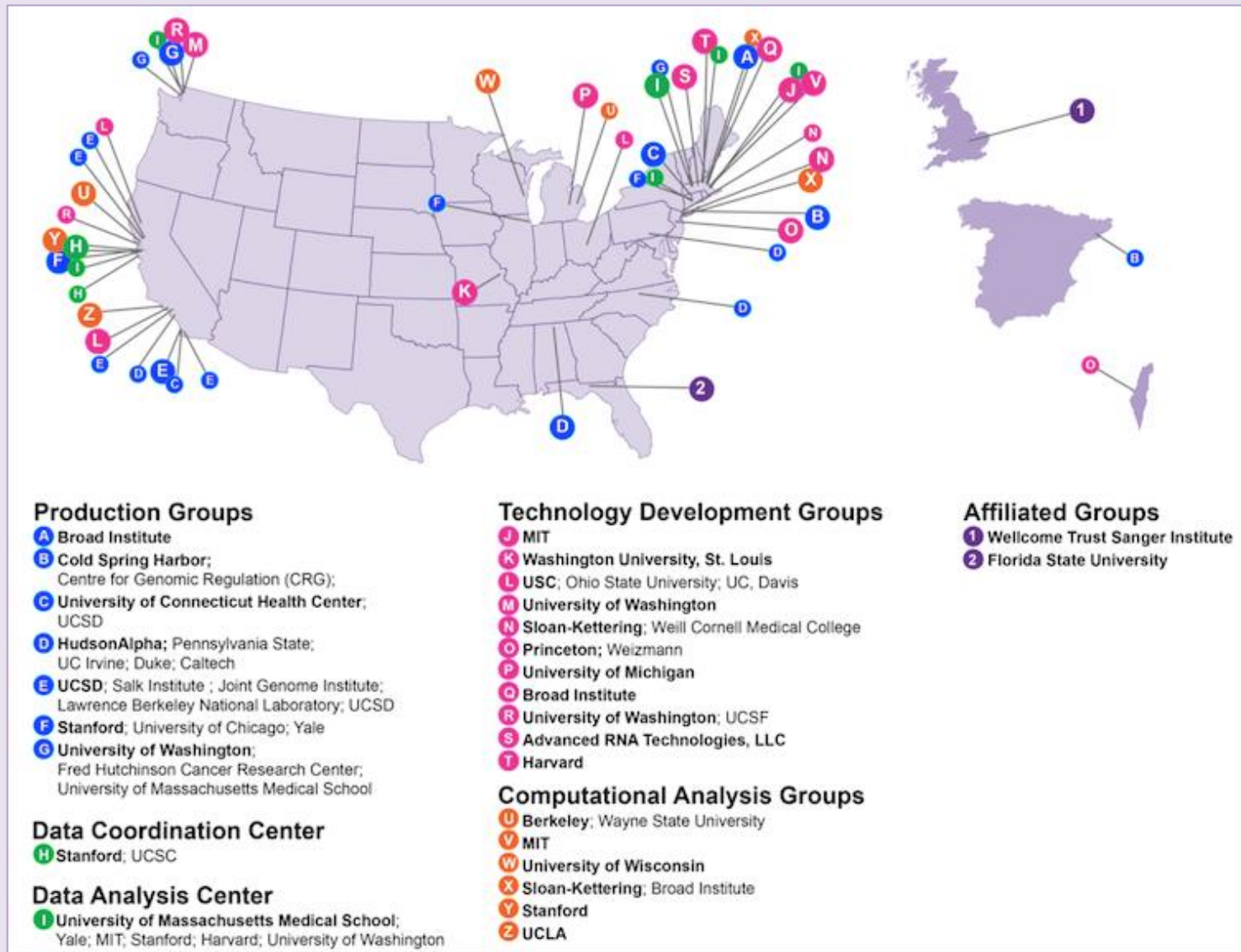
- Click the Data toolbar above and browse data

## News
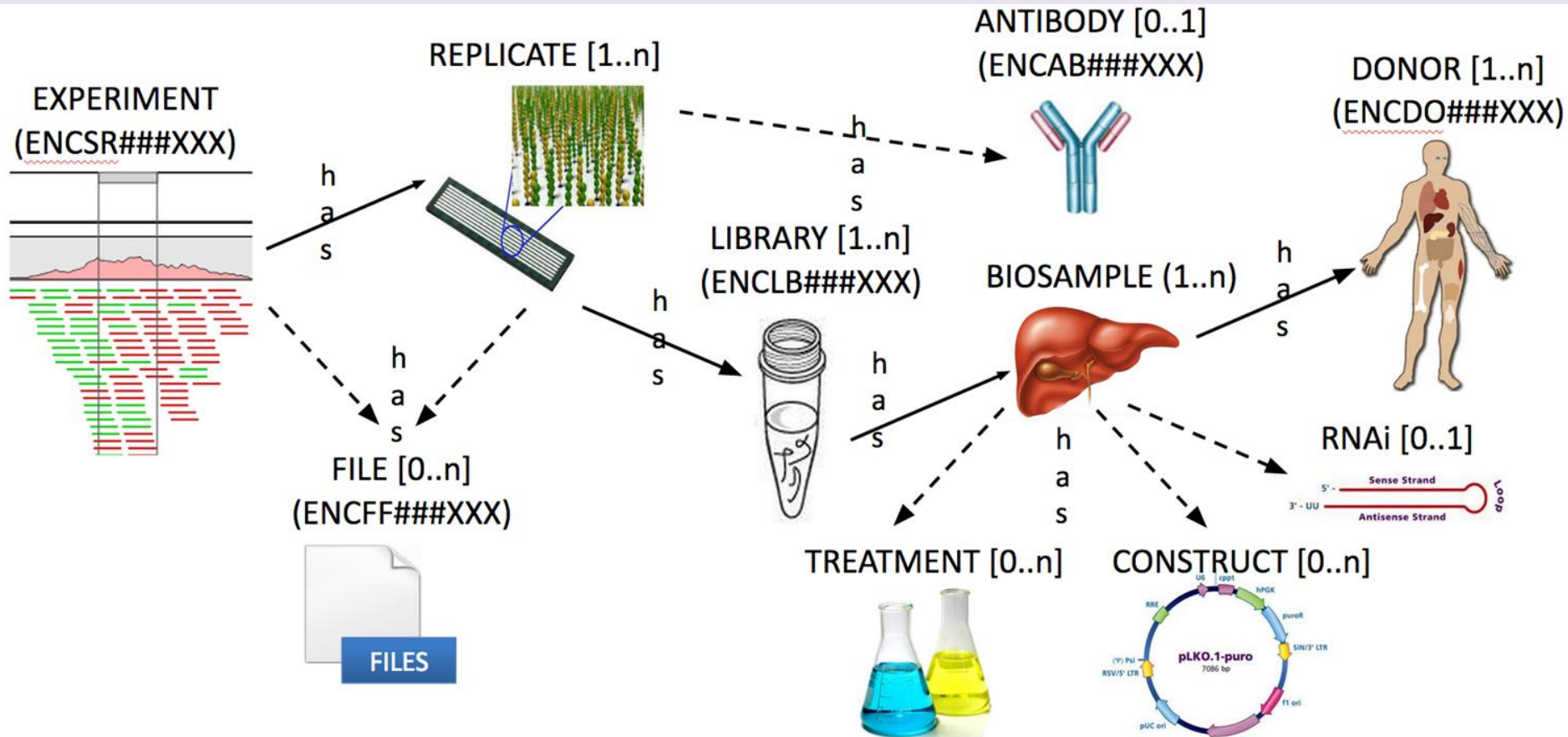
**Sept 12, 2014**: Data release: 23 human and 5 mouse datasets. [read more]

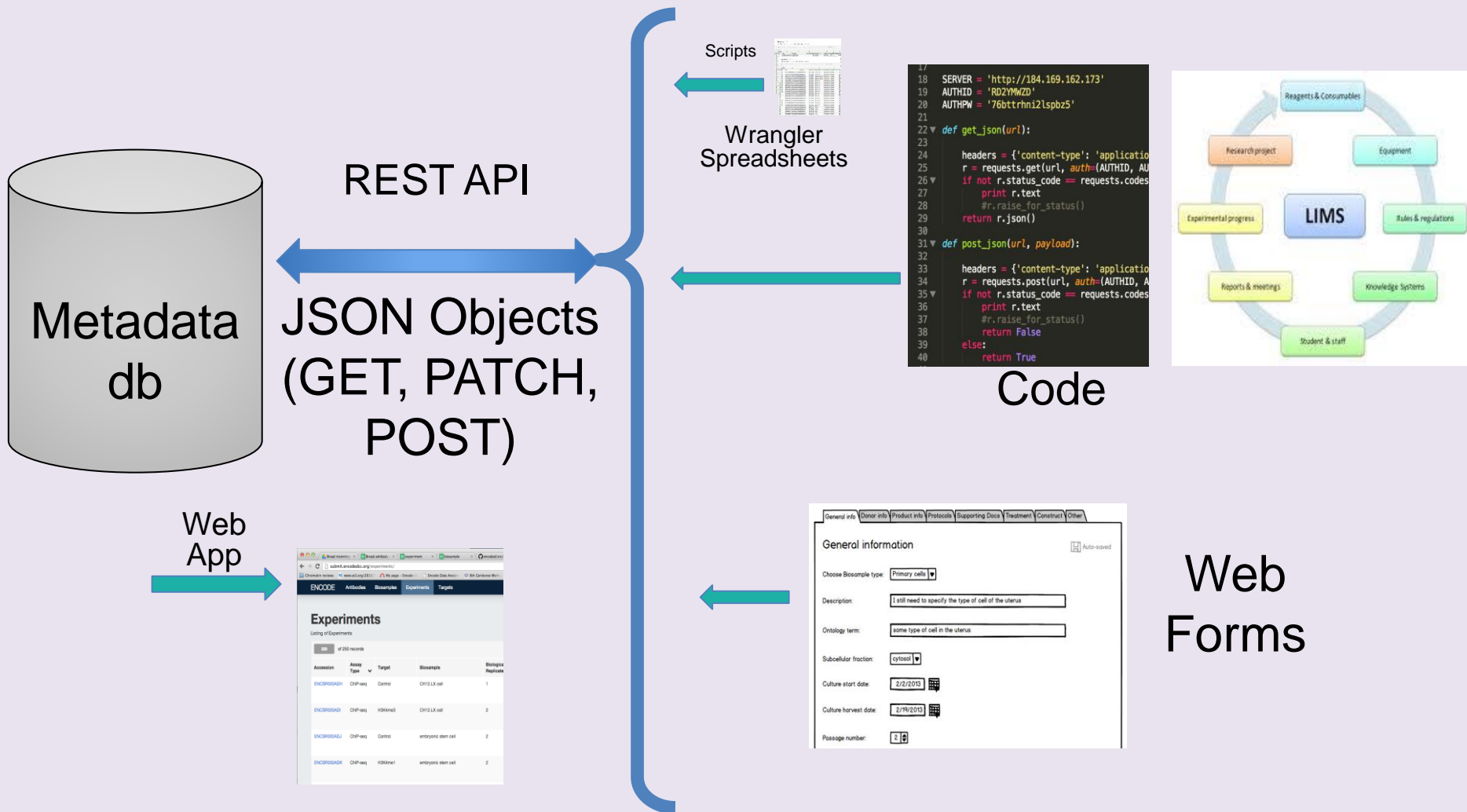**August 28, 2014**: modENCODE and ENCODE comparison papers
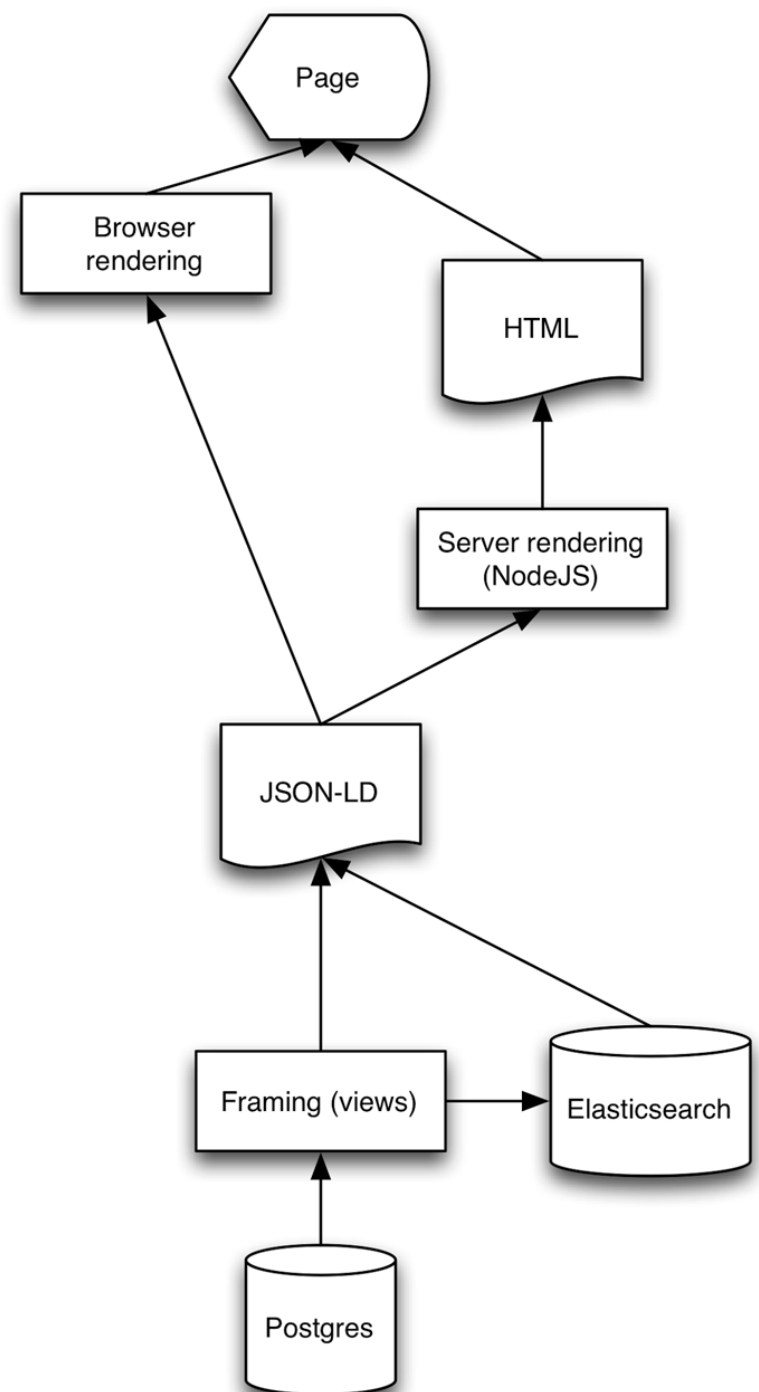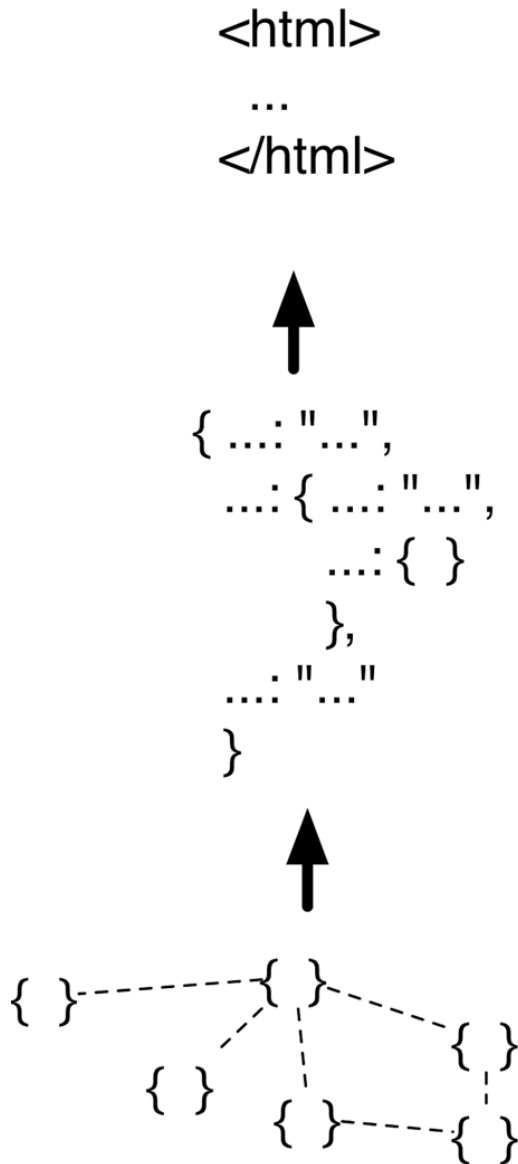
# What is the ENCODE Consortium?



**Production Groups**
- **A** Broad Institute
- **B** Cold Spring Harbor;
  Centre for Genomic Regulation (CRG);
- **C** University of Connecticut Health Center;
  UCSD
- **D** HudsonAlpha; Pennsylvania State;
  UC Irvine; Duke; Caltech
- **E** UCSD; Salk Institute ; Joint Genome Institute;
  Lawrence Berkeley National Laboratory; UCSD
- **F** Stanford; University of Chicago; Yale
- **G** University of Washington;
  Fred Hutchinson Cancer Research Center;
  University of Massachusetts Medical School

**Data Coordination Center**
- **H** Stanford; UCSC

**Data Analysis Center**
- **I** University of Massachusetts Medical School;
  Yale; MIT; Stanford; Harvard; University of Washington

**Technology Development Groups**
- **J** MIT
- **K** Washington University, St. Louis
- **L** USC; Ohio State University; UC, Davis
- **M** University of Washington
- **N** Sloan-Kettering; Weill Cornell Medical College
- **O** Princeton; Weizmann
- **P** University of Michigan
- **Q** Broad Institute
- **R** University of Washington; UCSF
- **S** Advanced RNA Technologies, LLC
- **T** Harvard

**Computational Analysis Groups**
- **U** Berkeley; Wayne State University
- **V** MIT
- **W** University of Wisconsin
- **X** Sloan-Kettering; Broad Institute
- **Y** Stanford
- **Z** UCLA

**Affiliated Groups**
- **1** Wellcome Trust Sanger Institute
- **2** Florida State University

*Image credit: NHGRI*

# Metadata model

# Metadata submission



Scripts

Wrangler
Spreadsheets

REST API

JSON Objects
(GET, PATCH,
POST)

Metadata
db

Code

Web
App

Web
Forms

# Software Stack

```
<html>
...
</html>
```

↑

```
{ ...: "...",
  ...: { ...: "...",
         ...: { }
       },
  ...: "..."
}
```

↑

```
{ } ----- { }  ----- { }
       { }   ----- { } ----- { }
```

Page

Browser rendering

HTML

Server rendering (NodeJS)

JSON-LD

Framing (views) → Elasticsearch

Postgres

# Find Common Biosamples Between Two Consortia



356 terms

http://genome.ucsc.edu/ENCODE/cellTypes.html

314 terms

GEO characteristics: common_name, tissue_type, cell_type, lines

# Labs were internally consistent

… but only 3 biosample names match exactly between projects

# 45 Biosamples in Common Between Current ENCODE & REMC



31 UBERON

11 CL

2 EFO

1 NTR

439 terms

154 terms

# An ontology is a set of words...

mitochondrion

mitochondrial
chromosome

nucleus

cell

chromosome

# An ontology is a set of words...
*.. with different types of relationships to each other.*

# Why use ontologies?

- Consistency of language and identifiers facilitates identification of data programmatically.  Alternative spellings & phrases are synonyms.  Independent of a particular data model.

$$F \neq f \neq Female \neq female$$

- Biological concepts are defined to provide scope

  *Mitochondria: A semiautonomous, self replicating organelle that occurs in varying numbers, shapes, and sizes in the cytoplasm of virtually all eukaryotic cells. It is notably the site of tissue respiration.*

- Relationships between terms can be computed to provide additional annotation details for grouping, searching, or analysis

# Challenge: Find all heart-related tissues?

**Showing 16 of 16**

## Organism
| | |
|---|---|
| *Homo sapiens* | 10 |
| *Mus musculus* | 6 |

## Biosample status
| | |
|---|---|
| in progress | 7 |
| released | 7 |
| deleted | 2 |

## Biosample type
| | |
|---|---|
| tissue | 8 |
| primary cell | 7 |
| in vitro differentiated cells | 1 |

## Organ
| | |
|---|---|
| heart | 14 |

## Sex
| | |
|---|---|
| male | 8 |
| unknown | 6 |
| female | 2 |

## Life stage
| | |
|---|---|
| fetal | 7 |
| adult | 6 |
| unknown | 3 |

## Source
| | |
|---|---|
| John Stamatoyannopoulos | 7 |
| BDRL | 6 |

---

### heart (*Mus musculus*, adult 8 week)
**Type**: tissue
**Source**: John Stamatoyannopoulos

Biosample
ENCBS536YRO
**deleted**

---

### heart (*Homo sapiens*, fetal 80 day)
**Type**: primary cell
**Source**: BDRL

Biosample
ENCBS913ULP
**in progress**

---

### heart (*Homo sapiens*, fetal 76 day)
**Type**: primary cell
**Source**: BDRL

Biosample
ENCBS953MIB
**in progress**

---

### heart (*Mus musculus*, adult 8 week)
**Type**: tissue
**Source**: John Stamatoyannopoulos

Biosample
ENCBS331ENC
**released**

---

### cardiac fibroblast (*Homo sapiens*)
**Type**: primary cell
**Source**: ScienCell

Biosample
ENCBS307AAA
**released**

---

### heart (*Mus musculus*, adult 8 week)
**Type**: tissue
**Source**: John Stamatoyannopoulos

Biosample
ENCBS846GWQ
**released**

---

### heart (*Mus musculus*, adult 8 week)

Biosample

**A**

ENCODE  Data ▾  Methods ▾  About ENCODE ▾  Help ▾          Search ENCODE  🔍   Sign In

**Assay**
ChIP-seq ............................... 2392
RNA-seq ............................... 655
DNase-seq ............................ 265
RNA profiling by array assay .... 180
shRNA knockdown followed by .. 167
RNA-seq
+ See more...

**Experiment status**
released ............................... 4400
revoked ..................................... 4

**Genome assembly**
hg19 ..................................... 2542
mm9 ....................................... 560
dm3 ....................................... 108

**Organism**
*Homo sapiens* ..................... 3389
*Mus musculus* ........................ 879
*Drosophila melanogaster* ...... 108

**Biosample type**
immortalized cell line ............ 2530
primary cell ............................ 767
tissue ..................................... 700
stem cell ................................ 208
in vitro differentiated cells ...... 122
+ See more...

**Organ**
brain ...................................... 200
skin of body ........................... 165
blood vessel ........................... 109
lung ......................................... 89
liver ......................................... 78
+ See more...

**Biosample treatment**
ethanol ..................................... 54
17β-estradiol ........................... 36
dimethyl sulfoxide ................... 35
dexamethasone ...................... 28
all-trans-retinoic acid .............. 21
+ See more...

**Available data**
fastq ..................................... 3890
bam ...................................... 3045
bigWig ................................... 3012
bed_narrowPeak ................... 1316
broadPeak ............................ 1295
+ See more...

Showing 25 of 4404          Filter to 500 to visualize ⧉   View All

ChIP-seq of MEL cell line (*Mus musculus*)                 Experiment
  Target: H3K36me3                                          ENCSR972MRN
  Lab: Michael Snyder, Stanford                             released
  Project: ENCODE

ChIP-seq of K562 (*Homo sapiens*, adult 53 year)           Experiment
  Target: Control                                           ENCSR996HUG
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of K562 (*Homo sapiens*, adult 53 year)           Experiment
  Target: Control                                           ENCSR599SPN
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of K562 (*Homo sapiens*, adult 53 year)           Experiment
  Target: Control                                           ENCSR824WND
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of K562 (*Homo sapiens*, adult 53 year)           Experiment
  Target: Control                                           ENCSR896JFO
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of K562 (*Homo sapiens*, adult 53 year)           Experiment
  Target: Control                                           ENCSR830KKO
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of GM12878 (*Homo sapiens*, adult 53 year)        Experiment
  Target: Control                                           ENCSR360XQV
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of GM12878 (*Homo sapiens*, adult 53 year)        Experiment
  Target: Control                                           ENCSR237VSG
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of HepG2 (*Homo sapiens*, child 15 year)          Experiment
  Target: Control                                           ENCSR655LCA
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

ChIP-seq of HepG2 (*Homo sapiens*, child 15 year)          Experiment
  Target: Control                                           ENCSR649FPL
  Lab: Richard Myers, HAIB                                  released
  Project: ENCODE

shRNA knockdown followed by RNA-seq of HepG2 (*Homo sapiens*, child 15 year)   Experiment

**B**

**Assay**
ChIP-seq ............................... 106
RNA-seq ................................. 59
DNase-seq ............................. 11
whole genome bisulfite sequencing 9
RAMPAGE ................................. 6
+ See more...

**Experiment status**
released ................................. 200

**Genome assembly**
mm9 ......................................... 45
hg19 ......................................... 16

**Organism**
*Mus musculus* ........................ 136
*Homo sapiens* .......................... 64

Showing 25 of 200          Visualize ▾   View All

RNA-seq of forebrain (*Mus musculus*, embryonic 11.5 day)    Experiment
  Lab: Barbara Wold, Caltech                                 ENCSR160IIN
  Project: ENCODE                                            released

RNA-seq of hindbrain (*Mus musculus*, embryonic 11.5 day)    Experiment
  Lab: Barbara Wold, Caltech                                 ENCSR760TOE
  Project: ENCODE                                            released

RNA-seq of midbrain (*Mus musculus*, embryonic 11.5 day)     Experiment
  Lab: Barbara Wold, Caltech                                 ENCSR307BCA
  Project: ENCODE                                            released

RNA-seq of hindbrain (*Mus musculus*, postnatal 0 day)       Experiment
  Lab: Barbara Wold, Caltech                                 ENCSR861FGB

**Organ**
brain                              200 ⊗
skin of body                       165
blood vessel                       109
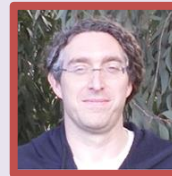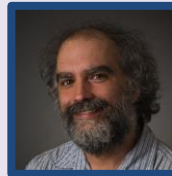lung                                89
liver                               78
heart                               75
kidney                              67
muscle organ                       41
mammary gland                      36
extraembryonic structure           33
bone element                       32
eye                                25
gonad                              22

**Biosample treatment**

| | |
|---|---|
| ethanol | 54 |
| 17β-estradiol | 36 |
| dimethyl sulfoxide | 35 |
| dexamethasone | 28 |
| all-trans-retinoic acid | 21 |

+ See more...

**Available data**

| | |
|---|---|
| fastq | 3890 |
| bam | 3051 |
| bigWig | 3012 |
| bed_narrowPeak | 1316 |
| broadPeak | 1295 |

+ See more...

B

**Assay**

| | |
|---|---|
| ChIP-seq | 106 |
| RNA-seq | 59 |
| DNase-seq | 11 |
| whole genome bisulfite sequencing | 9 |
| RAMPAGE | 6 |

+ See more...

**Experiment status**

| | |
|---|---|
| released | 200 |

**Genome assembly**

| | |
|---|---|
| mm9 | 45 |
| hg19 | 16 |

**Organism**

| | |
|---|---|
| *Mus musculus* | 136 |
| *Homo sapiens* | 64 |

**Biosample type**

| | |
|---|---|
| tissue | 160 |
| primary cell | 40 |

**Organism**

| | |
|---|---|
| brain | 200 ⊗ |
| skin of body | 165 |
| blood vessel | 109 |
| lung | 89 |
| liver | 78 |
| heart | 75 |
| kidney | 67 |
| muscle organ | 41 |
| mammary gland | 36 |
| extraembryonic structure | 33 |
| bone element | 32 |
| eye | 25 |
| gonad | 22 |
| stomach | 22 |
| placenta | 19 |
| bronchus | 18 |
| small intestine | 18 |
| mouth | 16 |
| spleen | 14 |
| thymus | 13 |
| large intestine | 11 |
| prostate gland | 11 |
| esophagus | 10 |
| pancreas | 8 |
| spinal cord | 8 |
| urinary bladder | 6 |
| adrenal gland | 5 |
| thyroid gland | 3 |
| tongue | 3 |
| trachea | 3 |
| lymphatic vessel | 1 |

- See fewer

**Showing 25 of 200**  Visualize  View All

RNA-seq of forebrain (*Mus musculus*, embryonic 11.5 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR160IIN
released

RNA-seq of hindbrain (*Mus musculus*, embryonic 11.5 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR760TOE
released

RNA-seq of midbrain (*Mus musculus*, embryonic 11.5 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR307BCA
released

RNA-seq of hindbrain (*Mus musculus*, postnatal 0 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR861FGB
released

RNA-seq of forebrain (*Mus musculus*, postnatal 0 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR527RFK
released

RNA-seq of midbrain (*Mus musculus*, postnatal 0 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR026ZRP
released

RNA-seq of hindbrain (*Mus musculus*, postnatal 0 day)
Lab: Thomas Gingeras, CSHL
Project: ENCODE
Experiment
ENCSR749BAG
released

RNA-seq of midbrain (*Mus musculus*, postnatal 0 day)
Lab: Thomas Gingeras, CSHL
Project: ENCODE
Experiment
ENCSR255SDF
released

RNA-seq of forebrain (*Mus musculus*, postnatal 0 day)
Lab: Thomas Gingeras, CSHL
Project: ENCODE
Experiment
ENCSR723SZV
released

RNA-seq of forebrain (*Mus musculus*, embryonic 11.5 day)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR000OXO
released

RNA-seq of Purkinje cell (*Homo sapiens*, adult 20 year)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR417EDR
released

RNA-seq of Purkinje cell (*Homo sapiens*, adult 20 year)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
ENCSR672QFA
released

RNA-seq of Purkinje cell (*Homo sapiens*, adult 20 year)
Lab: Barbara Wold, Caltech
Project: ENCODE
Experiment
released

RNA-seq of Purkinje cell (*Homo sapiens*, adult 20 year)
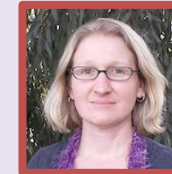Experiment

# ENCODE DCC

https://www.encodeproject.org/
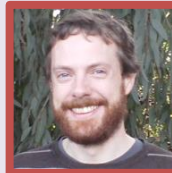
Eurie Hong, Mike Cherry (PI), Jim Kent (co-PI), Ben Hitz

Zhiping Weng, ENCODE DAC

**Data Wranglers**

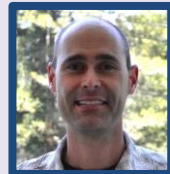Esther Chan, Jean Davidson, Venkat Malladi, Cricket Sloan, J. Seth Strattan

**Software Engineers**

Funding Source: NHGRI

Nikhil Podduturi, Laurence Rowe, Forrest Tanaka

**QA, administration, biocuration**

Brian Lee, Stuart Miyasato, Matt Simison, Zhenhua Wang, Marcus Ho

*Malladi et al. ; Database, 2015, 1–11 doi: 10.1093/database/bav010*

@encodedcc     encode-help@lists.stanford.edu     https://github.com/ENCODE-DCC/

# Using ontologies for metadata annotation

1. **Uber Anatomy ontology** (UBERON; http://uberon.org/)
   - tissues: heart, blood, brain

2. **Cell Ontology** (CL; http://cellontology.org/)
   - primary cell types: hepatocyte, cardiomyocyte

3. **Experimental Factor Ontology** (EFO; http://www.ebi.ac.uk/efo/)
   - immortalized cell lines: K562, HepG2, MCF-7

4. **Ontology for Biomedical Investigations** (OBI; http://obi-ontology.org/page/Main_Page)
   - experimental assays: RNA-seq, CLIP-seq, ChIP-seq, etc

5. **Chemical Entities of Biological Interest** (ChEBI; http://www.ebi.ac.uk/chebi/)
   - chemical treatments: estradiol, ethanol, etc

6. **Sequence Ontology** (SO; http://www.sequenceontology.org/)
   - nucleic acid being sequenced: microRNA, poly-A+ mRNA, etc

7. **Gene Ontology** (GO; http://www.geneontology.org/)
   - group gene products that are targets of ChIP-seq or RNAi experiments