

# Using UCSC Tools for Browsing and Data Mining ENCODE Data

## Aims

- Learn to locate and display ENCODE data in the UCSC Genome Browser
- Learn to retrieve ENCODE data from the UCSC Genome Browser database using the Table Browser data retrieval tool

## Introduction

The University of California Santa Cruz (UCSC) Genome Browser at <http://genome.ucsc.edu> is a web-based set of tools providing access to a database of genome sequence and annotations for visualization, comparison and analysis by the scientific, medical and academic communities. The primary mission of the site is to provide timely and convenient open access to high-quality human genome sequence and annotations in a framework that enables easy exploration from genome-wide down to the base level. Annotation datasets, or 'tracks', on the human genome cover conservation and evolutionary comparisons, gene models, regulation, expression, epigenetics and tissue differentiation, variation, phenotype and disease associations. A substantial contributor to our mission has been participation in the ENCODE project as the designated data repository in the ENCODE Pilot (2003-2007) and as the Data Coordination Center (DCC) in the ENCODE whole-genome data production phase (2007-2011). All production ENCODE data is routed to UCSC for validation, quality review, database storage, visualization, and dissemination to other public databases. At this time over 2700 distinct ENCODE experiments have been processed by the DCC and made publicly available.

Other organisms represented at the site include 4 non-human primates, 14 other mammals including a marsupial and a monotreme, 10 non-mammalian vertebrates and 24 non-vertebrates. The Genome Browser hosts mapping and sequence annotation tracks that describe assembly, gap and GC content for all organisms in the browser database. Additionally, for most organisms we show alignments from RefSeq genes, mRNAs and ESTs from GenBank, and other gene or gene prediction tracks such as Ensembl Genes (6). For human and mouse assemblies, we also offer a locally generated UCSC Genes track based upon RefSeq, GenBank, CCDS and UniProt data. About half of the genomes hosted at UCSC include a multiple sequence alignment track and pairwise genomic alignments between assemblies to further comparative and evolutionary investigations. Expression, regulation, variation and phenotype tracks are available for many of the assemblies. We also support user data upload and visualization, and have recently introduced a data hub mechanism allowing visualization of user data hosted remotely.

## Worked Example 1: You will need to follow this through with the OpenHelix exercises at the end of this section.

Examining RNA expression in the vicinity of the TP53 gene

- 1 Browse to genome.ucsc.edu.

UCSC Genome Bioinformatics

Genomes **2** Blat Tables Gene Sorter PCR VisiGene Proteome Session FAQ Help

**Genome Browser**

**About the UCSC Genome Bioinformatics Site**

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to the [ENCODE](#) and [Neandertal](#) projects.

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the Center for Biomolecular Science and Engineering (CBSE) at the University of California Santa Cruz (UCSC). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

Home Genomes Blat Tables Gene Sorter PCR Session FAQ Help

**Human (*Homo sapiens*) Genome Browser Gateway**

The UCSC Genome Browser was created by the [Genome Bioinformatics Group of UC Santa Cruz](#).  
Software Copyright (c) The Regents of the University of California. All rights reserved.

clade	genome	assembly	position or s	term	gdc
Mammal	Human	Feb. 2009 (GRCh37/hg19)	chr17:7,566,934	5,649	tp53

**3** [Click here to reset](#) the browser user interface settings to their defaults.

track search add custom tracks track hubs configure tracks and display clear position

**4**

About the Human Feb. 2009 (GRCh37/hg19) assembly ([sequences](#))

The image shows a screenshot of the UCSC Genome Browser interface. At the top, there is a navigation bar with links like Home, Genomes, Blat, Tables, Gene Sorter, PCR, DNA, Convert, PDF, Session, Ensembl, NCBI, and Help. Below this, the main title reads "UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly". A search bar contains the coordinates "chr17:7,566,934-7,595,649" and a "zoom out (1.5x)" button is circled in red. The main content area displays several tracks: RefSeq Genes, Human mRNAs, Human ESTs, H3K27ac (ENCODE), DNase-seq (ENCODE), Phylo-P (Phylo-P), RepeatMasker, and Simple Nucleotide Polymorphisms (SNPs). The tracks show various genomic features and data points across a 18 kb region. At the bottom, there are controls for track search, default tracks, default order, hide all, and expand all, along with a refresh button.

The screenshot shows the UCSC Genome Browser interface for Human Feb. 2009 (GRCh37/hg19) Assembly. The main track shows a region on chromosome 17 (75,737,720-75,900,812 bp). Below the main track are several categories of tracks:

- Mapping and Sequencing Tracks:** Includes Base Position, Chromosome Band, STS Markers, FISH Clones, Recomb Rate, ENCODE Pilot, Map Contigs, Assembly, GRC Map Contigs Gap, BAC End Pairs, OC Percent, GRC Patch Release, Hg18 Diff, GRC Incident, Hi Seq Depth, Short Match, Repeat Elements, WAI Track, RIL ORGid, and Mapability.
- Phenotype and Disease Associations:** Includes GAD View, DECIPHER, OMM AV SNPs, OMM Genes, OMM Pheuo Loci, COSMIC, GWAS Catalog, ISCA, RGD Human OTL, RGD Rat OTL, MGI Mouse OTL, and GeneReviews.
- Genes and Gene Prediction Tracks:** This section is circled in red. It includes UCSC Genes, Old UCSC Genes, Alt Events, GENCODE Genes V11 (circled in red), GENCODE Genes V10, GENCODE Genes V7, RefSeq Genes, Other RefSeq, MGC Genes, ORFgene Clones, TransMap, Yaga Genes, Ensembl Genes, AceView Genes, SIB Genes, N-SCAN, SGP Genes, Genoid Genes, GenScan Genes, Exonify, Yale Pseudo60, tRNA Genes, and H-Iny 7.0. A red box with the number '7' is located to the right of this section.
- miRNA and EST Tracks:** Includes Human mRNAs, Spliced ESTs, Human ESTs, Other mRNAs, Other ESTs, H-Iny, Gene Boundaries, SIB Alt-Splicing, Poly(A), PolyA-Seq, CGAP SAGE, and Human RNA Editing.
- Expression:** Includes Affy Exon Array, Affy GNF11, Affy RNA Loc, Affy U133, Affy U133Plus2, Affy U95, Allen Brain, Burge RNA-seq, ENC Exon Array, ENC ProtGemo, ENC RNA-seq (circled in red), and GIS RNA PET. A red box with the number '8' is located to the right of this section.

This screenshot provides a detailed view of the TP53 gene region (75,737,720-75,900,812 bp). The tracks shown include:

- Gene Structure:** TP53 gene structure with exons and introns.
- Transcript Evidence:** Tracks for TP53 transcripts from various sources like RefSeq, Ensembl, and GENCODE.
- Expression Data:** The 'Expression' section is highlighted with a red circle. It includes tracks for Affy Exon Array, Affy GNF11, Affy RNA Loc, Affy U133, Affy U133Plus2, Affy U95, Allen Brain, Burge RNA-seq, ENC Exon Array, ENC ProtGemo, ENC RNA-seq (circled in red), and GIS RNA PET. A red box with the number '8' is located to the right of this section.
- Other Data:** Tracks for miRNAs, ESTs, and other genomic features.

### ENC RNA-seq Super-track Settings

## ENCODE RNA-seq Tracks [\(All Expression tracks\)](#)

Display mode: show

All

- dense [Caltech RNA-seq](#) RNA-seq from ENCODE/Caltech
- dense [CSHL Long RNA-seq](#) Long RNA-seq from ENCODE/Cold Spring Harbor Lab
- dense [CSHL Sm RNA-seq](#) Small RNA-seq from ENCODE/Cold Spring Harbor Lab
- dense [GIS RNA-seq](#) RNA-seq from ENCODE/Genome Institute of Singapore
- dense [HAIB RNA-seq](#) RNA-seq from ENCODE/HAIB
- dense [RIKEN CAGE Loc](#) RNA Subcellular CAGE Localization from ENCODE/RIKEN
- dense [SYDH RNA-seq](#) RNA-seq from ENCODE/Stanford/Yale/USC/Harvard

NOTE: Early access to additional track data may be available on the [Preview Browser](#).

### ENC RNA-seq Super-track Settings

## ENCODE RNA-seq Tracks [\(All Expression tracks\)](#)

Display mode: show

All

- hide [Caltech RNA-seq](#) RNA-seq from ENCODE/Caltech
- full [CSHL Long RNA-seq](#) Long RNA-seq from ENCODE/Cold Spring Harbor Lab
- hide [CSHL Sm RNA-seq](#) Small RNA-seq from ENCODE/Cold Spring Harbor Lab
- hide [GIS RNA-seq](#) RNA-seq from ENCODE/Genome Institute of Singapore
- hide [HAIB RNA-seq](#) RNA-seq from ENCODE/HAIB
- hide [RIKEN CAGE Loc](#) RNA Subcellular CAGE Localization from ENCODE/RIKEN
- hide [SYDH RNA-seq](#) RNA-seq from ENCODE/Stanford/Yale/USC/Harvard

NOTE: Early access to additional track data may be available on the [Preview Browser](#).

### ENC RNA-seq Super-track Settings

## Long RNA-seq from ENCODE/Cold Spring Harbor Lab [\(ENC RNA-seq\)](#)

Maximum display mode: full   [Reset to defaults](#)

Select views (help): Contigs hide Plus Raw Signal full Minus Raw Signal full Splice Junctions hide Alignments hide

Select subtracks by localization and cell line:

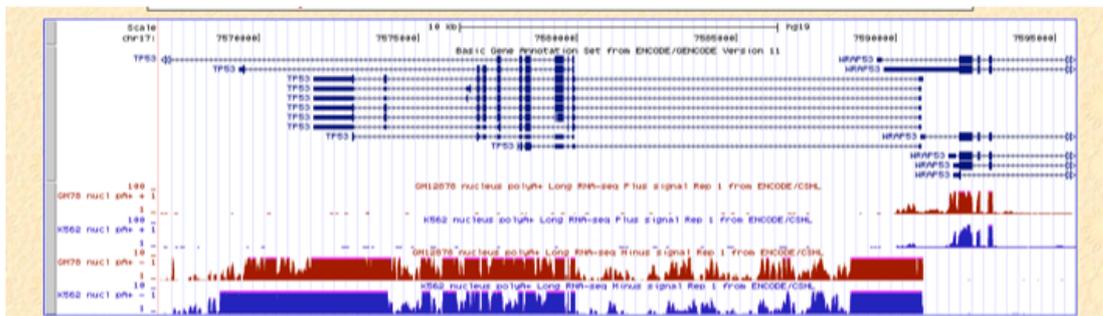
RNA Extract:  PolyA+  PolyA-  Total RNA

Rep:  1  2  Pooled

	Localization	Whole Cell	Cytosol	Nucleus	Nucleoplasm	Chromatin	Nucleolus
<span style="border: 1px solid red; border-radius: 50%; padding: 2px;">All</span>	Cell Line	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
GM12878 (Tier 1)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
H1hESC (Tier 1)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
K562 (Tier 1)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
HeLaS3 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
HepG2 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
HUVEC (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
A549	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
AG04450	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
BJ	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
HMEC	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
HSMM	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
MCF-7	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
NHEK	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
NHLF	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
SK-N-SH retinoic acid	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

List subtracks:  only selected/visible  all (4 of 473 selected) [Topf](#)

Views <sup>1</sup>	Cell Line <sup>2</sup>	Localization <sup>3</sup>	RNA Extract <sup>4</sup>	Rep <sup>5</sup>	Track Name	
<input checked="" type="checkbox"/> <span style="border: 1px solid gray; border-radius: 3px; padding: 2px;">full</span>	<input type="checkbox"/> Plus Raw Signal	GM12878	Nucleus	PolyA+	1	GM12878 nucleus polyA+ Long RNA-seq Plus signal Rep 1 from ENCODE/CSHL
<input checked="" type="checkbox"/> <span style="border: 1px solid gray; border-radius: 3px; padding: 2px;">full</span>	<input type="checkbox"/> Plus Raw Signal	K562	Nucleus	PolyA+	1	K562 nucleus polyA+ Long RNA-seq Plus signal Rep 1 from ENCODE/CSHL



12

## Worked Example 2: Exploring TFBS and Histone Marks in the TP53 region

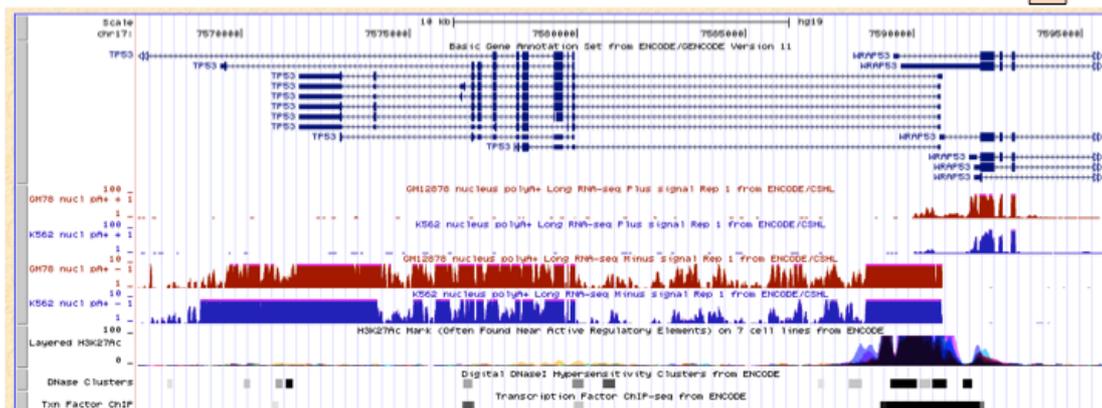
1

2

Regulation

<input checked="" type="checkbox"/> ENCODE Regulation... <input type="button" value="show"/>	<input checked="" type="checkbox"/> CD34 Dnase1 <input type="button" value="hide"/>	<input checked="" type="checkbox"/> CpG Islands <input type="button" value="hide"/>	<input checked="" type="checkbox"/> ENC Chromatin... <input type="button" value="hide"/>	<input checked="" type="checkbox"/> ENC DNA Methyl... <input type="button" value="hide"/>	<input checked="" type="checkbox"/> ENC DNase/FAIRE... <input type="button" value="hide"/>
<input checked="" type="checkbox"/> ENC Histone... <input type="button" value="hide"/>	<input checked="" type="checkbox"/> ENC RNA Binding... <input type="button" value="hide"/>	<input checked="" type="checkbox"/> ENC TF Binding... <input type="button" value="hide"/>	<input checked="" type="checkbox"/> FSU Repli-chip <input type="button" value="hide"/>	<input checked="" type="checkbox"/> ORegAnno <input type="button" value="hide"/>	<input checked="" type="checkbox"/> Stanf Nucleosome <input type="button" value="hide"/>
<input checked="" type="checkbox"/> SUNY SwitchGear <input type="button" value="hide"/>	<input checked="" type="checkbox"/> SwitchGear TSS <input type="button" value="hide"/>	<input checked="" type="checkbox"/> TFBS Conserved <input type="button" value="hide"/>	<input checked="" type="checkbox"/> TS miRNA sites <input type="button" value="hide"/>	<input checked="" type="checkbox"/> UMMS Brain Hist <input type="button" value="hide"/>	<input checked="" type="checkbox"/> UW Repli-seq <input type="button" value="hide"/>
<input checked="" type="checkbox"/> Vista Enhancers <input type="button" value="hide"/>	<input checked="" type="checkbox"/> NK1 Nuc Lamina... <input type="button" value="hide"/>	<input checked="" type="checkbox"/> UCSE Brain Methyl <input type="button" value="hide"/>	<input type="button" value="refresh"/>		

3



**ENCODE Integrated Regulation from ENCODE Tracks** (▲All Regulation tracks)

Display mode: show : Submit

**All**

- hide : [Transcription](#) Transcription Levels Assayed by RNA-seq on 7 Cell Lines from ENCODE
- hide : [Layered H3K4Me1](#) H3K4Me1 Mark (Often Found Near Regulatory Elements) on 7 cell lines from ENCODE
- hide : [Layered H3K4Me3](#) H3K4Me3 Mark (Often Found Near Promoters) on 7 cell lines from ENCODE
- full : [Layered H3K27Ac](#) H3K27Ac Mark (Often Found Near Active Regulatory Elements) on 7 cell lines from ENCODE
- dense : [DNase Clusters](#) Digital DNaseI Hypersensitivity Clusters from ENCODE
- full : [Txn Factor ChIP](#) Transcription Factor ChIP-seq from ENCODE

**Transcription Factor ChIP-seq from ENCODE (Po2-4H8)**

Factor: Po2-4H8  
 Cluster Score (out of 1000): 943  
 Position: chr17:7589024-7589555  
 Band: 17p13.1  
 Genomic Size: 532  
[View DNA for this feature \(hg19/Human\)](#)

#	signal	chr	cellType	factor	treatment	lab	more info
1	436.00	G	GM12878	Po2-4H8	None	HadsonAlpha	metadata *
2	185.00	g	GM12891	Po2-4H8	None	HadsonAlpha	metadata *
3	235.00	g	GM12892	Po2-4H8	None	HadsonAlpha	metadata *
4	943.00	l	H1-hESC	Po2-4H8	None	HadsonAlpha	metadata *
5	786.00	h	HCT-116	Po2-4H8	None	HadsonAlpha	metadata *
6	488.00	K	K562	Po2-4H8	None	HadsonAlpha	metadata *

Symbol	Cell Type
l	H1-hESC
a	A549+DEX_100nM
a	A549+EtOH_0.02pct
a	A549+DEX_50nM

### ENCODE H3K27Ac Mark (Often Found Near Active Regulatory Elements) on 7 cell lines from ENCODE

(• ENCODE Regulation)

Display mode:

Overlay method:

Type of graph:

Track height:  pixels (range: 11 to 100)

Vertical viewing range: min:  max:  (range: 0 to 3851)

Data view scaling:  Always include zero:

Transform function: Transform data points by:

Windowing function:  Smoothing window:  pixels

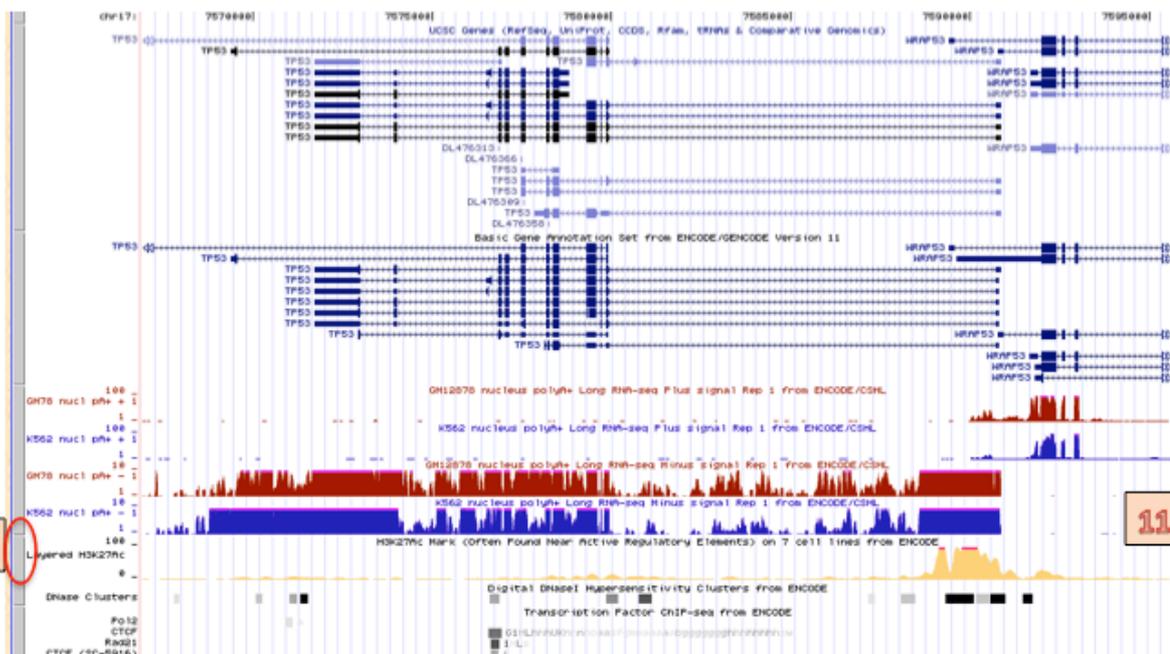
Draw y indicator lines: at y = 0.0:  at y =

[Graph configuration help](#)

List subtracks:  only selected/visible  all (1 of 7 selected) Restricted Until

<input checked="" type="checkbox"/>	GM12878	H3K27Ac Mark (Often Found Near Regulatory Elements) on GM12878 Cells from ENCODE	<a href="#">schema</a>	2009-10-05
<input checked="" type="checkbox"/>	H1-hESC	H3K27Ac Mark (Often Found Near Regulatory Elements) on H1-hESC Cells from ENCODE	<a href="#">schema</a>	2011-03-21
<input checked="" type="checkbox"/>	HMM	H3K27Ac Mark (Often Found Near Regulatory Elements) on HMM Cells from ENCODE	<a href="#">schema</a>	2010-09-16
<input checked="" type="checkbox"/>	HUVEC	H3K27Ac Mark (Often Found Near Regulatory Elements) on HUVEC Cells from ENCODE	<a href="#">schema</a>	2009-10-06
<input checked="" type="checkbox"/>	K562	H3K27Ac Mark (Often Found Near Regulatory Elements) on K562 Cells from ENCODE	<a href="#">schema</a>	2009-10-05
<input checked="" type="checkbox"/>	NHEK	H3K27Ac Mark (Often Found Near Regulatory Elements) on NHEK Cells from ENCODE	<a href="#">schema</a>	2009-10-07
<input checked="" type="checkbox"/>	NHLF	H3K27Ac Mark (Often Found Near Regulatory Elements) on NHLF Cells from ENCODE	<a href="#">schema</a>	2010-06-28

1 of 7 selected



## Worked Example 3:

Intersect NFKB binding sites with RNA-seq using the Table Browser.

1

UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

position/search chr17:7,566,934-7,595,649    size 28,716 bp.

chr17 (p10.1) 7,566,934 7,595,649

**Table Browser**

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade: Mammal genome: Human assembly: Feb. 2009 (GRCh37/hg19)

group: Regulation track: Tfn Factor ChIP

table: wgEncodeRegTfbsClustered

region:  genome  ENCODE Pilot regions  position chr21:1-48129895

identifiers (names/accessions):

filter:

intersection:

output format: all fields from selected table Send output to  Galaxy  GREAT

output file:  (leave blank to keep output in browser)

file type returned:  plain text  gzip compressed

2

**Table Browser**

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade: Mammal genome: Human assembly: Feb. 2009 (GRCh37/hg19)

group: Regulation track: Tfn Factor ChIP

table: wgEncodeRegTfbsClustered

region:  genome  ENCODE Pilot regions  position chr21:1-48129895

identifiers (names/accessions):

filter:

intersection:

output format: all fields from selected table Send output to  Galaxy  GREAT

output file:  (leave blank to keep output in browser)

file type returned:  plain text  gzip compressed

3

**Filter on Fields from hg19.wgEncodeRegTfbsClustered**

bin is ignored 0

chrom does match \*

chromStart is ignored 0 AND

chromEnd is ignored 0 AND

name does match nfkb

score is > 500 AND

strand does match \*

thickStart is ignored 0 AND

thickEnd is ignored 0 AND

reserved is ignored 0 AND

blockCount is ignored 0 AND

blockSizes does match \*

chromStarts does match \*

expCount is ignored 0 AND

explds does match \*

expScores does match \*

AND Free-form query:

4



**Table Browser**

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade:  genome:  assembly:

group:

table:

region:  genome  ENCODE Pilot

identifiers (names/accessions):

filter:

intersection:

output format:

output file:

file type returned:  plain text  gzip compressed

**Intersect with Txn Factor ChIP**

Select a group, track and table to intersect with:

group:  track:

table:

Note: Txn Factor ChIP has gene/alignment structure. Only the exons/blocks will be considered.

**Intersect Txn Factor ChIP items with bases covered by CSHL Long RNA-seq:**

These combinations will maintain the names and gene/alignment structure (if any) of Txn Factor ChIP:

All Txn Factor ChIP records that have any overlap with CSHL Long RNA-seq

All Txn Factor ChIP records that have no overlap with CSHL Long RNA-seq

All Txn Factor ChIP records that have at least  % overlap with CSHL Long RNA-seq

All Txn Factor ChIP records that have at most  % overlap with CSHL Long RNA-seq

**Intersect bases covered by Txn Factor ChIP and/or CSHL Long RNA-seq:**

These combinations will discard the names and gene/alignment structure (if any) of Txn Factor ChIP and produce a simple list of position ranges.

Base-pair-wise intersection (AND) of Txn Factor ChIP and CSHL Long RNA-seq

Base-pair-wise union (OR) of Txn Factor ChIP and CSHL Long RNA-seq

Check the following boxes to complement one or both tables. To complement a table means to include a base pair in the intersection/union if it is *not* included in the table.

Complement Txn Factor ChIP before base-pair-wise intersection/union

Complement CSHL Long RNA-seq before base-pair-wise intersection/union

**Table Browser**

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade:  genome:  assembly:

group:  track:

table:

region:  genome  ENCODE Pilot regions  position

identifiers (names/accessions):

filter:

intersection with

output format:    Galaxy  GREAT

output file:  (leave blank to keep output in browser)

file type returned:  plain text  gzip compressed

Note: The all fields and selected fields output formats are not available when an intersection has been specified.

Hyperlinks to Genome Browser

- [NFKB at chr21:15458907-15459240](#)
- [NFKB at chr21:19273936-19274302](#)
- [NFKB at chr21:26829134-26829459](#)
- [NFKB at chr21:26946025-26946474](#)
- [NFKB at chr21:26950582-26951164](#)
- [NFKB at chr21:27107050-27107546](#)
- [NFKB at chr21:30374917-30375340](#)
- [NFKB at chr21:30560714-30561033](#)

14

The screenshot shows the UCSC Genome Browser interface for human chromosome 21. The top navigation bar includes links for Home, Genomes, Blat, Tables, Gene Sorter, PCR, DNA, Convert, PDF, Session, Ensembl, and NCBI. The main title is "UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly". Below the title, there are navigation controls for moving and zooming, with a zoom level of 10x selected. The search bar shows the position "chr21:19,273,936-19,274,302" and the gene "NFKB". The main content area displays several tracks: Scale, chr21, UCSC Genes (including NFKB), Basic Gene Annotation, ENCODE data (such as DNase-seq, ChIP-seq, and Hi-C), and Digital Distal Interactions. The tracks are color-coded and show various genomic features and data points.

## Exercises for the ENCODE data in the UCSC Genome Browser

- 1) Using RNA-seq data, examine expression of RNA in the vicinity of the TP53 gene. From the CSHL Long RNA-seq track, determine which strand is transcribed into Poly-A+ RNA and then found in the nuclear fraction of K562 and GM12878 cells.

*Skills: Use RNA-seq data to evaluate RNA presence in a region; become aware of the cellular fraction data that is available.*

- 2) In the region we are exploring, let's add transcription factor binding data and histone marks that are often found near active regulatory elements. Let's also determine if these histone marks are indicated in human embryonic stem cells.

*Skills: Explore TFBS data; examine features associated with histone modifications; visualize cell type specific data.*

- 3) Use the Table Browser to locate NKFB transcription factor binding signals that are greater than 500 on chromosome 21. Let's intersect that with RNA-seq data indicating presence of RNA in epidermal keratinocyte cells (NHEK cells).

*Skills: Table Browser to query ENCODE data; use filters and intersections to generate a complex customized query of the data.*

**For additional guidance and ways to interact with the ENCODE data, access this open access publication in PLoS Biology: <http://bit.ly/plosENC>**

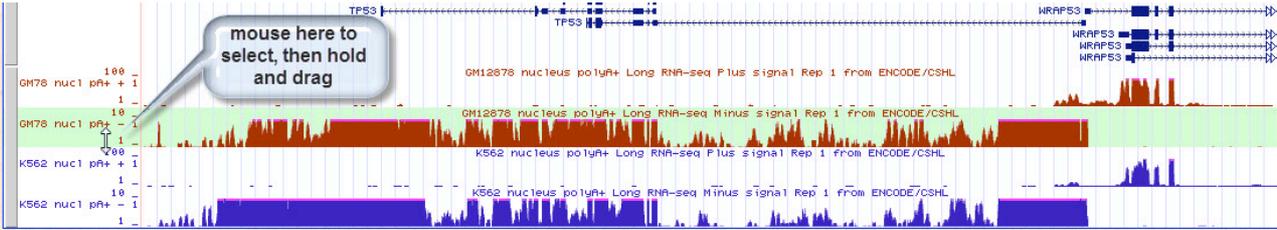
**Citation:** The ENCODE Project Consortium (2011) A User's Guide to the Encyclopedia of DNA Elements (ENCODE). PLoS Biol 9(4): e1001046. doi:10.1371/journal.pbio.1001046

**UCSC ENCODE 2 Exercises, version 1.  
Correspond to the data available in March 2012.**

**The materials and slides offered are for non-commercial use only. Reproduction, distribution and/or use for commercial purposes are strictly prohibited.  
Copyright 2012, OpenHelix, LLC.**

- 1) Using RNA-seq data, examine expression of RNA in the vicinity of the TP53 gene. From the CSHL Long RNA-seq track, determine which strand is transcribed into Poly-A+ RNA and then found in the nuclear fraction of K562 and GM12878 cells.

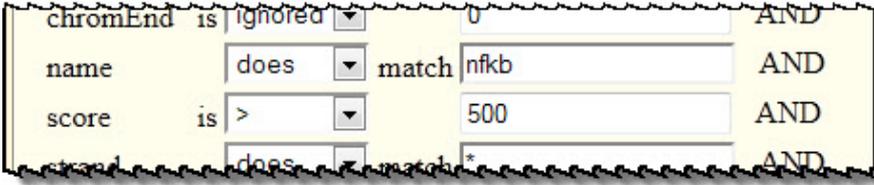
Step	Action	✓
1	Go to the UCSC Genome Browser homepage, <a href="http://genome.ucsc.edu">genome.ucsc.edu</a>	
2	From the blue navigation links on the left side of the page, <b>click the link for Genome Browser.</b>	
3	From the Gateway interface, <b>click the link that says “<u>Click here to reset the browser user interface settings to their defaults.</u>”</b> This will ensure that any prior activity on the Browser has been cleared out and that everyone is starting with default settings.	
4	Choose the Human <b>February 2009</b> assembly. Enter the text <b>tp53</b> in the gene box. Choose the TP53 item in the list. Click <b>submit.</b>	
5	In the TP53 region on the browser, examine the features briefly. Then <b>click the “zoom out” 1.5x button</b> near the top. Assess the features again.	
6	<b>Click the “hide all” button</b> in the middle of the resulting Genome viewer page. <i>(We want to reduce what’s in the display to reduce the burden on the servers, and to focus on our features of interest.)</i>	
7	<p><b>Add back 2 tracks</b> to the viewer:</p> <ul style="list-style-type: none"> <li>• <b>GENCODE Genes V11 in “pack”</b> visibility (from the Genes and Gene Predictions group)</li> <li>• <b>ENC RNA-seq...in “show”</b> (from the Expression group)</li> </ul> <p><b>Click a refresh button</b> to add these tracks back to the viewer. It may take a while for this to load, as there is a lot of data here.</p>	
8	<p>RNA seq data from multiple labs, cell lines, and experiment types are shown. Let’s focus on the Long RNA-seq data. You can see there is signal across this region indicating RNA transcription in this region of the genome in this mode. But we’d like to distinguish which RNA-seq data corresponds to which genes in this region.</p> <p><b>Return to the RNA-seq... hyperlink and click it to access individual tracks from this super-track.</b></p>	
9	Note all the RNA-seq... component tracks are in “dense” visibility at this time. <b>Turn all of them to “hide” except for the CSHL Long RNA-seq menu.</b>	
10	<b>Click the CSHL Long RNA-seq hyperlink. Examine the options you can set to explore this track’s data.</b> <i>(continued on next page)</i>	

<p>1 1</p>	<p><b>At the top, set the Maximum display mode to “full”.</b></p> <p><b>Make these changes</b> on the Long RNA-seq settings:</p> <ul style="list-style-type: none"> <li><b>*Set Contig view to “hide”.</b></li> <li><b>*Poly-A+ should be checked. Uncheck “Total RNA” in the extract row. Leave the other settings in that area unchanged.</b></li> <li><b>*Select the GM12878 cell line “nucleus” localization checkbox.</b></li> <li><b>*Unselect all other localization checkboxes except K562 “nucleus”.</b></li> </ul> <p><b>Click “Submit” when these changes have been made.</b></p>
<p>1 2</p>	<p>Back on the viewer, examine the data. <b>Use the select/drag feature of the left label area to move the GM12878 data sets together.</b></p>  <p>The screenshot displays a genomic track with four main signal tracks. From top to bottom: GM12878 nucleus polyA+ Long RNA-seq Plus signal (red), GM12878 nucleus polyA+ Long RNA-seq Minus signal (blue), K562 nucleus polyA+ Long RNA-seq Plus signal (red), and K562 nucleus polyA+ Long RNA-seq Minus signal (blue). Above the tracks, gene models for TP53 and WRAP53 are shown. A callout box on the left side of the GM12878 Plus signal track contains the text: "mouse here to select, then hold and drag".</p>
<p>1 3</p>	<p><b>Note the data which derives from the Plus strand and which from the Minus strand.</b> It appears that the WRAP53 RNA derives from the plus strand, and the TP53 RNA from the minus strand. This will help you to orient when looking for transcription factor binding sites or other genomic features.</p>
<p>This exercise was inspired by the Figure 3 illustration in the ENCODE User Guide paper. See that figure legend and the accompanying text for more assessments of the data and the features in this region: <a href="http://bit.ly/plosENC">http://bit.ly/plosENC</a></p>	

2) In the region we are exploring, let's add transcription factor binding data and histone marks that are often found near active regulatory elements. Let's also determine if these histone marks are indicated in human embryonic stem cells.

Step	Action	✓
1	On the browser view that we established in exercise 1, <b>scroll down to the Regulation Group.</b>	
2	<b>Locate the ENCODE Regulation...</b> track. Choose <b>“show”</b> in the pulldown menu. <b>Click a “refresh”</b> button.	
3	<b>Examine the display.</b> New data appears in the viewer beneath the RNA-seq data. *Note that the Transcription Factor ChIP-seq from ENCODE track shows data blocks, but not individual transcription factors. *Note that the H3K27Ac histone mark track appears to have multiple data sets of various colors.	
4	<b>Return to the ENCODE Regulation menu area. Click the hyperlink</b> to look at the component tracks of this super-track.	
5	By default Txn Factor ChIP is visible in “dense” mode. <b>Set that menu to “full”. Click the “Submit” button.</b>	
6	Examine the display again. Note that individual transcription factors can be identified by name using the labels on the left. Note that the letter codes near the blocks correspond to cell lines that have been used in experiments for this data. <b>Click some of the blocks</b> to note the cell lines and signal levels observed in them. Return to the viewer for the next steps.	
7	<b>Click the grey control button to the left of the Layered H3K27Ac</b> to go to the controls for that track.	
8	On this histone mark page, note that there are various cell line data sets, which have color codes. One of the lines is <b>H1-hESC</b> , which is a human embryonic stem cell line.	
9	<b>Uncheck all cell line boxes except H1-hESC.</b>	
10	<b>Click the “Submit” button at the top</b> to return to the genome viewer.	
11	<b>Note that we can now see</b> that there is signal associated with this histone mark in stem cells in this region. This was difficult to examine before because of the other color overlays.	
12	<b>Return to the histone mark page</b> by clicking the gray bar to the left of the browser track. <b>Turn on or off various cell lines</b> to view the data. <b>Return to the viewer</b> each time by clicking “Submit”.	
13	The various data types in this region should help you to understand possible features of regulation of the genes in this area.	
<p>This exercise was inspired by the Figure 5 illustration in the ENCODE User Guide paper. See that figure legend and the accompanying text for more assessments of the data and the features in this region: <a href="http://bit.ly/plosENC">http://bit.ly/plosENC</a></p>		

3) Use the Table Browser to locate NFKB transcription factor binding signals that are greater than 500 on chromosome 21. Let's intersect that with RNA-seq data indicating presence of RNA in epidermal keratinocyte cells (NHEK cells).

Step	Action	✓
1	From the genome browser, <b>click the navigation bar option called Tables.</b>	
2	<b>At the Table Browser, begin to establish the query with these choices:</b> *Mammal, human, February 2009 *Regulation group, Txn Factor ChIP track *table wgEncodeRegTfbsClustered *region: position chr21. <b>Click the lookup button to load the chr21 range.</b>	
3	Next we'll set a filter. <b>Click the "create" button.</b>  We want the factor NFKB, and signals to be over 500. *in the name area chose "does" match nfkb <i>[remove the asterisk]</i> *in the score choose "is >" and type 500 in the text box  <b>Your filter should look like this:</b>  <b>Click submit.</b>	
4	<b>Click the summary/statistics button to assess the results at this point.</b> This will provide a sense of how many results return with these settings. If there were too many or too few, you might want to adjust the filters accordingly. <b>Return to the table browser by clicking the Tables link.</b>	
5	Let's take a look at the output at this point. In the "output format" area <b>select "all fields from selected table".</b>	
6	<b>Click "get output"</b> to see the results in table form. Note that the Name field has our choice, and all the scores are over 500.	
	<i>Let's intersect this data with some other data. Let's require that this NFKB data also overlap with RNA-seq evidence in a particular cell type of our choice.</i>	
7	<b>Use the back button</b> to go back to the Table Browser interface. It should still have all of your previous choices and settings.	
8	<b>Find the "intersection" option, and click the "create" button.</b> <i>(continued on the next page)</i>	

9	<p>Make these choices in the “Intersect with” interface:</p> <p><b>*In the group menu, select Expression.</b></p> <p><b>*Track choice: select CSHL Long RNA-seq.</b> <i>(This is only because we are already familiar with this track and have some of it visible in the browser. You could choose any of the data sets later.)</i></p> <p><b>*Table selection: choose nuclear NHEK polyA+ first data set.</b> This looks like: <b>NHEK nucl pA+ + 1</b> <i>(This is the CSHL Long data set Plus track)</i></p>	
10	<p>Ensure that the Intersect <b>radio button is set to the first choice</b> for “any overlap”.</p>	
11	<p><b>Click submit.</b></p>	
12	<p>This time let’s <b>choose “output format” as “hyperlinks to Genome Browser”</b>. This will allow us to quickly inspect some of the results visually.</p>	
13	<p><b>Click the “get output” button.</b></p>	
14	<p><b>Click on some of the links to explore</b> the region of the browser that meets these criteria. Zoom out for larger scope.</p> <p>Do you see the NHEK data in the current view? If not, <b>go to the ENC RNA seq... super-track and access the CSHL Long RNA-seq track details</b> as we did before. <b>Select NHEK nuclear data</b> to add it to the viewer. <b>Submit.</b></p> <p>Show or hide various expression tracks, transcription factor tracks, or any other features you are interested in. You may need to turn on or off tracks in the browser because they were not on when we were using it before.</p>	

**Tasks:**

1.  
Search for the SLC25A29 gene in UCSC and view the GENCODE geneset.  
How many alternative splice variants are there and what are their biotypes?  
What regulatory information can you find by investigating the ENCODE data tracks not only for this gene but for the adjacent locus?
  
2.  
Search for RP4-550H1.6.1.  
What biotype is this and what can you observe about this locus?
  
3.  
Search for RPN2 and then zoom out so that you can also see C20orf132.  
How many alternative variants are there for each of these two loci and what are their biotypes? What can you observe about these two loci?

**Answers:**

1.

There are 21 alternative splice variants for this locus.

8 protein coding, 3 NMD, 7 processes transcript and 2 retained intron.

Observations: Histone marks H3K4Me3 shows the classic dip at the tss.

There is TF binding, transcriptional evidence and exonic conservation back to zebrafish. Contrast this with the adjacent locus SLC25A47. This just has 2 coding transcripts and a pretty much identical conservation pattern, but appears to be transcriptionally silent in ENCODE cell-types. You can see the difference in histone modifications (very low levels detected), open chromatin (weak signals), TF binding (lower levels) and general transcription (background).

2.

This is a single-exon lincRNA locus that has a biotype of lincRNA (protein coding LOC in UCSC genes).

There is no alternative splicing. The conservation is quite good for a lincRNA, right back to Tenrec. Looking at the Open Chromatin marks there is a string signal from DNaseq and FAIRE, and the histone marks H3K4Me1, Me3 and K27Ac are all strong. Looking at the TF binding there is a huge pileup of factors with a very tight distribution and transcription looks high across the full-length of the transcript. If you look at polyA+ CAGE TSS there are strong signals on the minus strand, with the highest on H1-hESC and K562 and some expression (but not much) in GM12878. The polyA site and signal annotation shows very tight correspondence with the transcription level.

3.

C20orf132 has 10 alternative variants: 9 protein-coding and 1 processed transcript.

RPN2 has 11 alternative variants: 9 protein-coding, 1 NMD and 1 processed transcript.

They have similar levels of conservation, with C20orf132 conserved down to platypus and RPN2 to zebrafish. The two loci have a head to head arrangement. The histone modifications have detected high levels here and the Open Chromatin marks give a strong signal from DNase1 and FAIRE. The TF binding has a big pileup of factors and a very tight distribution.

But, if you look at the transcription you can see that RPN2 has much higher transcription levels than C20orf132. This is shown by looking at the CAGE TSS data, as its much higher for the positive strand (RPN2) than the negative strand (C20orf132) and is true in all three cell types. There is also proteogenomic data for RPN2 e.g. peptides TGQEVVFVAEPDNK and FPEEEAPSTVLSQNLFTPQ are found in both mitochondrial and nuclear fraction.