

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 (2007-2012).

In Phase III, beginning in 2012, the DCC is managed at Stanford University and ENCODE production data continues to be routed to UCSC for validation, quality review, database storage, visualization, and dissemination to other public databases. At this time more than 2700 distinct ENCODE experiments have been processed by the DCC and made publicly available.

Other organisms represented at the UCSC Genome Browser site include 11 non-human primates, 34 other mammals including marsupials and a monotreme, 18 non-mammalian vertebrates, 23 invertebrates and yeast. 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. For human and mouse assemblies, we also offer a locally generated UCSC Genes track based upon RefSeq, GenBank and CCDS 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 offer a data-hub mechanism allowing visualization of user data hosted remotely.

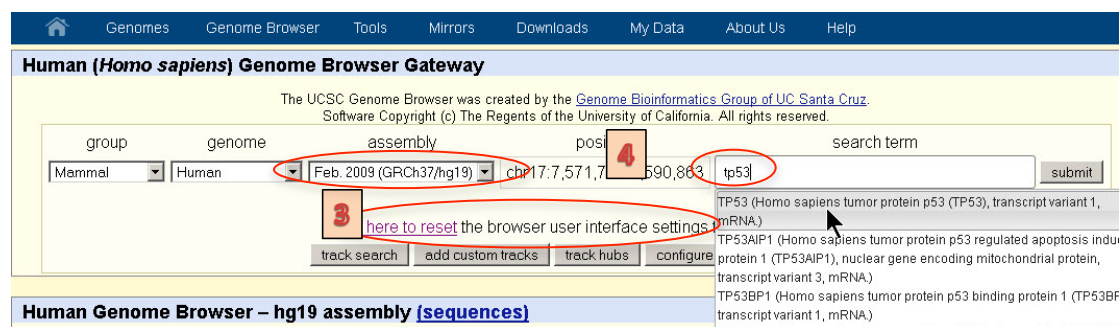
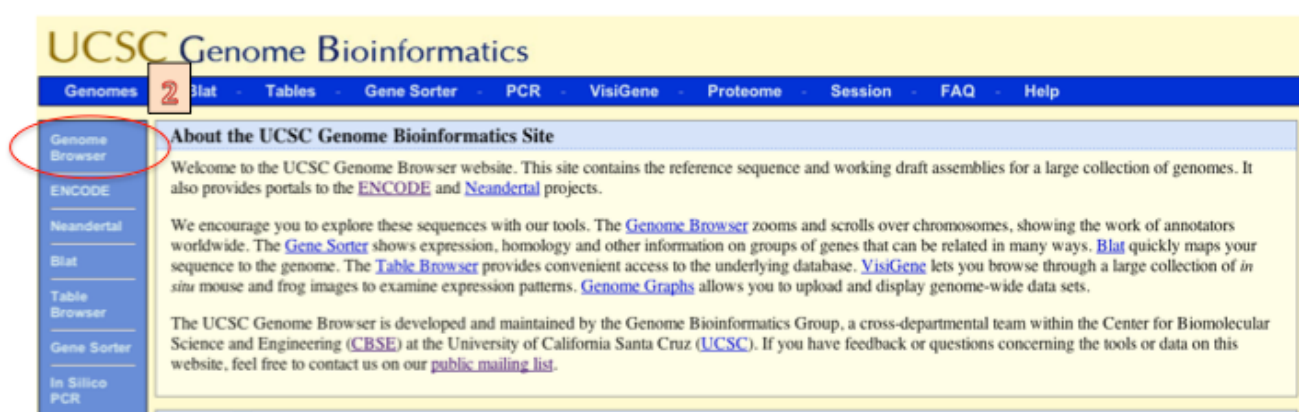
Step-by-step screenshots for the Exercises.

Worked Example 1:

Examining RNA expression in the vicinity of the TP53 gene

1

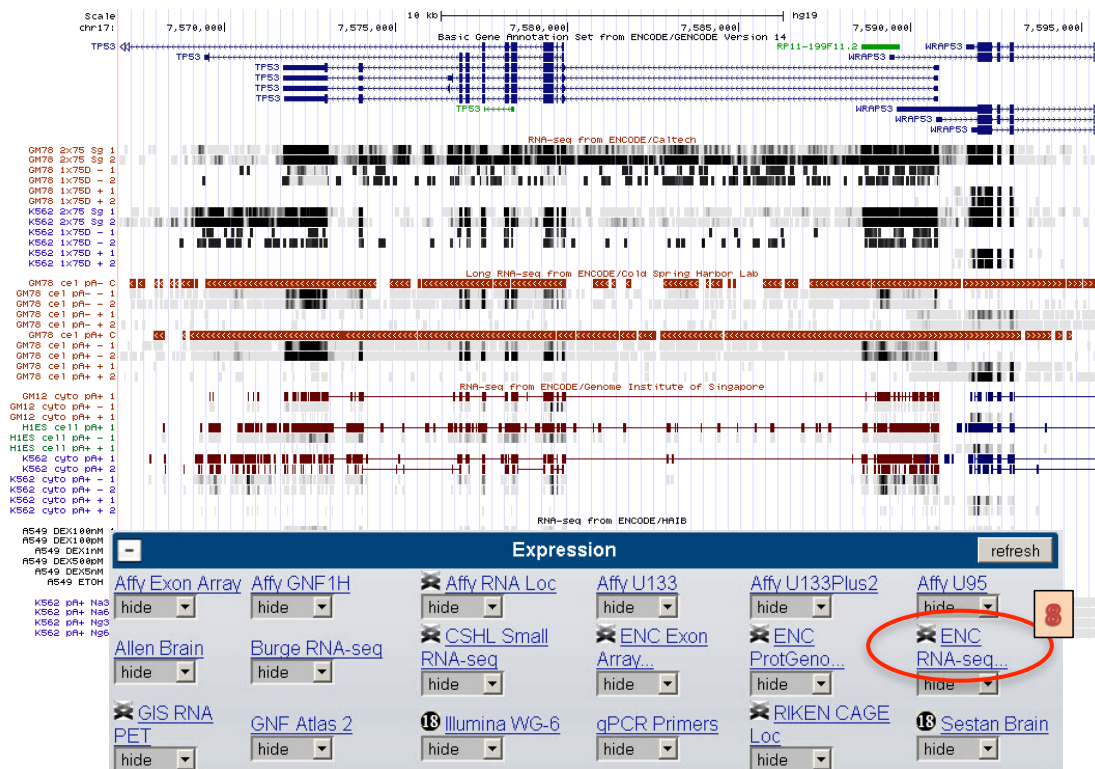
Browse to genome.ucsc.edu.



The screenshot displays the UCSC Genome Browser interface for the TP53 gene region on chromosome 17 (hg19 assembly). The main header shows the browser title and navigation options. Below the header, there are zoom controls (1.5x, 3x, 10x) and a search bar. The main content area contains several tracks: UCSC Genes (RefSeq, UniProt, CCDS, Rfam, tRNAs & Comparative Genomics), RefSeq Genes, Sequences (SNPs), Human mRNAs from GenBank, Spliced ESTs, Layered H3K27ac, DNase Clusters, Transcription Factor ChIP-seq, Placental Mammal Basewise Conservation by PhyloP, Multiz Alignments of 46 Vertebrates, Common SNPs (137), and RepeatMasker. A red circle highlights the 'base' zoom option in the zoom controls. A red box highlights the 'hide all' button in the track controls at the bottom. The interface also includes a 'move start' and 'move end' section with a zoom level of 2.0.

The screenshot displays the UCSC Genome Browser interface for Human Feb. 2009 (GRCh37/hg19) Assembly. The main view shows chromosome 17 with a scale from 7,575,000 to 7,595,000 bp. The interface is organized into several track categories:

- Mapping and Sequencing Tracks:** Includes Base Position, Chromosome Band, STS Markers, FISH Clones, Recomb Rate, deCODE Recomb, ENCODE Pilot, Map Contigs, Assembly, GRC Map Contigs, Gap, BAC End Pairs, Fosmid End Pairs, GC Percent, GRC Patch Release, Hg18 Diff, GRC Incident, Hi Seq Depth, Wiki Track, BU ORChID, Mapability, Short Match, and Restr Enzymes.
- Phenotype and Disease Associations:** Includes GAD View, DECIPHER, OMIM AV SNPs, OMIM Genes, OMIM Pheno Loci, COSMIC, GWAS Catalog, ISCA, Coriell CNVs, RGD Human QTL, RGD Rat QTL, MGI Mouse QTL, and GeneReviews. A red box with the number '7' highlights the 'refresh' button.
- Genes and Gene Prediction Tracks:** Includes UCSC Genes, GENCODE (circled in red), Old UCSC Genes, Alt Events, CCDS, RefSeq Genes, Other RefSeq, MGC Genes, ORFeome Clones, TransMap, Vega Genes, Pfam in UCSC Gene, Ensembl Genes, AceView Genes, SIB Genes, N-SCAN, SGP Genes, Geneid Genes, GenScan Genes, Exoniphy, Yale Pseudo60, tRNA Genes, H-Inv 7.0, and EvoFold.
- Literature:** Includes Publications.
- mRNA and EST Tracks:** Includes Human mRNAs, Spliced ESTs, Human ESTs, Other mRNAs, Other ESTs, H-Inv, UniGene, Gene Bounds, SIB Alt-Splicing, Poly(A), PolyA-Seq, CGAP SAGE, Human RNA Editing.
- Expression:** Includes Affy Exon Array, Affy GNF1H, Affy RNA Loc, Affy U133, Affy U133Plus2, Affy U95, Allen Brain, Burge RNA-seq, CSHL Small RNA-seq, ENC Exon Array, ENC ProtGeno, ENC RNA-seq (circled in red), GIS RNA PET, GNF Atlas 2, Illumina WG-6, qPCR Primers, RIKEN CAGE, and Sestan Brain Loc. A red box with the number '7' highlights the 'refresh' button.
- Regulation:** Includes various regulatory tracks.



ENC RNA-seq Super-track Settings

ENCORE RNA-seq Tracks (▲All Expression tracks)

Display mode: show Submit

All 8

- dense Caltech RNA-seq RNA-seq from ENCODE/Caltech
- dense CSHL Long RNA-seq Long 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 SYDH RNA-seq RNA-seq from ENCODE/Stanford/Yale/USC/Harvard

ENC RNA-seq Super-track Settings

ENCORE RNA-seq Tracks (▲All Expression tracks)

Display mode: show Submit

RNA sequencing

All 10

- hide Caltech RNA-seq RNA-seq from ENCODE/Caltech
- full CSHL Long RNA-seq Long 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 SYDH RNA-seq RNA-seq from ENCODE/Stanford/Yale/USC/Harvard

Description

RNA sequencing, or RNA-seq, is a method for mapping and quantifying the total amount of RNA transcripts in a cell.

CSHL Long RNA-seq Track Settings [ENCODE](#) [Downloads](#) [Subtracks](#)

Long RNA-seq from ENCODE/Cold Spring Harbor Lab

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

Select views [\(help\)](#):
[Contigs](#) [Plus Signal](#) [Minus Signal](#) [Splice Junctions](#) [Alignments](#)

Select subtracks by localization and cell line:

	Localization	Whole Cell	Chromatin	Cytosol	Nucleolus	Nucleoplasm	Nucleus
<input type="text" value="+ -"/> All							
<i>Cell Line</i>		<input type="text" value="+ -"/>	<input type="text" value="+ -"/>	<input type="text" value="+ -"/>	<input type="text" value="+ -"/>	<input type="text" value="+ -"/>	<input type="text" value="+ -"/>
GM12878 (Tier 1)		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>		<input checked="" type="checkbox"/>
H1-hESC (Tier 1)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
K562 (Tier 1)		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
A549 (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
B cells CD20+ (Tier 2)		<input type="checkbox"/>					<input type="checkbox"/>
HeLa-S3 (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
HepG2 (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>

11

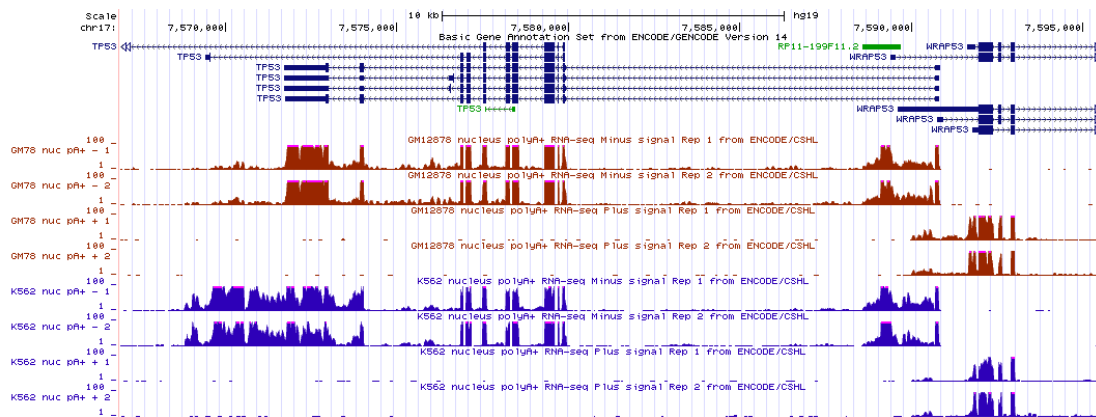
Cell Line

Localization **Whole Cell** **Chromatin** **Cytosol** **Nucleolus** **Nucleoplasm** **Nucleus**

All

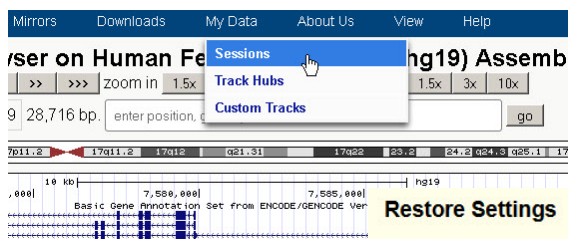
Select subtracks further by: (select multiple categories and items - [help](#))

RNA Extract: **Replicate rank:**



12

This view can be obtained directly using the session tool. Under "My Data, Sessions" ... "Restore Settings":
 user: example
 session name: hg19_korea2014



13

Restore Settings

Use settings from another user's saved session:
 user: session name:

Use settings from a local file: No file selected.

Use settings from a URL (http://..., ftp://...):

Select subtracks further by: (select)

RNA Extract:

Replicate rank:

- All
- 1st
- 2nd
- Pooled

List subtracks: full

14

Worked Example 2:

Exploring TFBS and Histone Marks in the TP53 region

1

2

3

4

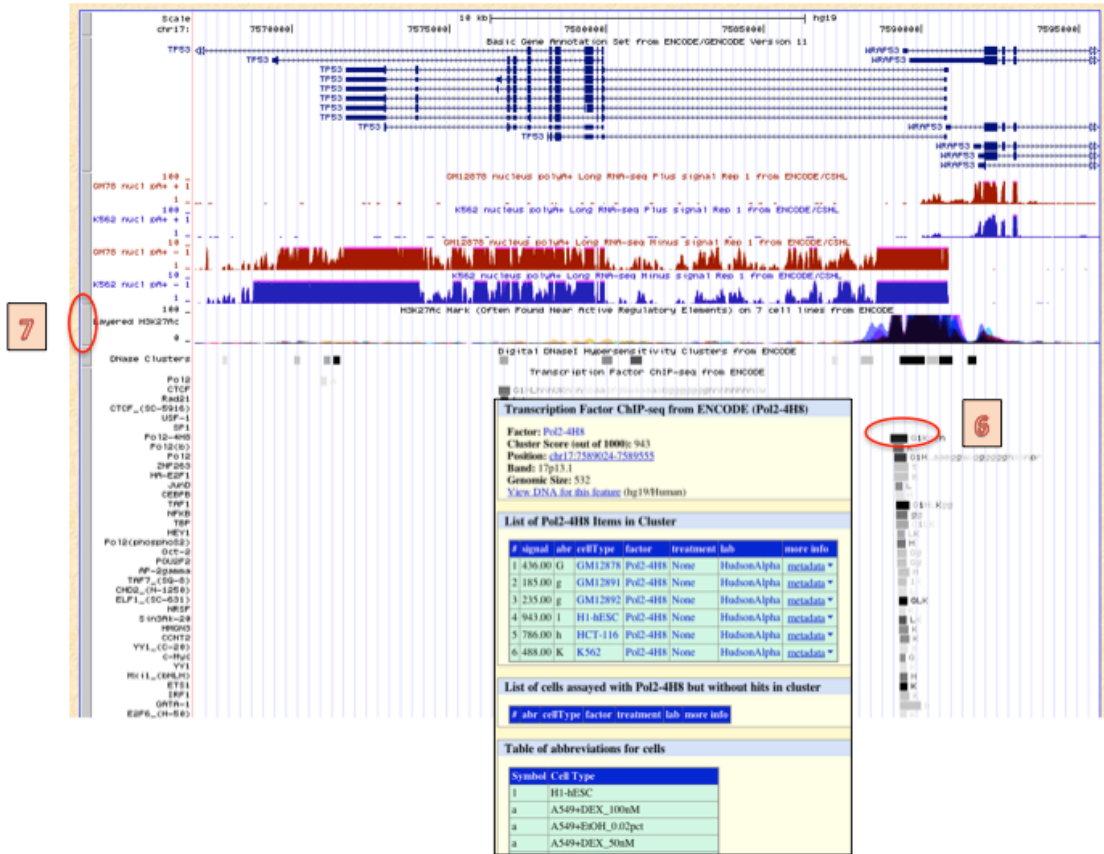
ENCODE Regulation Super-track Settings

Integrated Regulation from ENCODE Tracks (▲ [All Regulation tracks](#))

Display mode: show **Submit**

5

- All
- hide [Transcription](#) Transcription Levels Assayed by RNA-seq on 9 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 in 125 cell types from ENCODE
- hide [DNase Clusters V1](#) Digital DNaseI Hypersensitivity Clusters in 74 cell types (2 reps) from ENCODE
- full [Tsn Factor ChIP](#) Transcription Factor ChIP-seq from ENCODE



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

(★ ENCODE Regulation)

Display mode: 10

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 = 0:

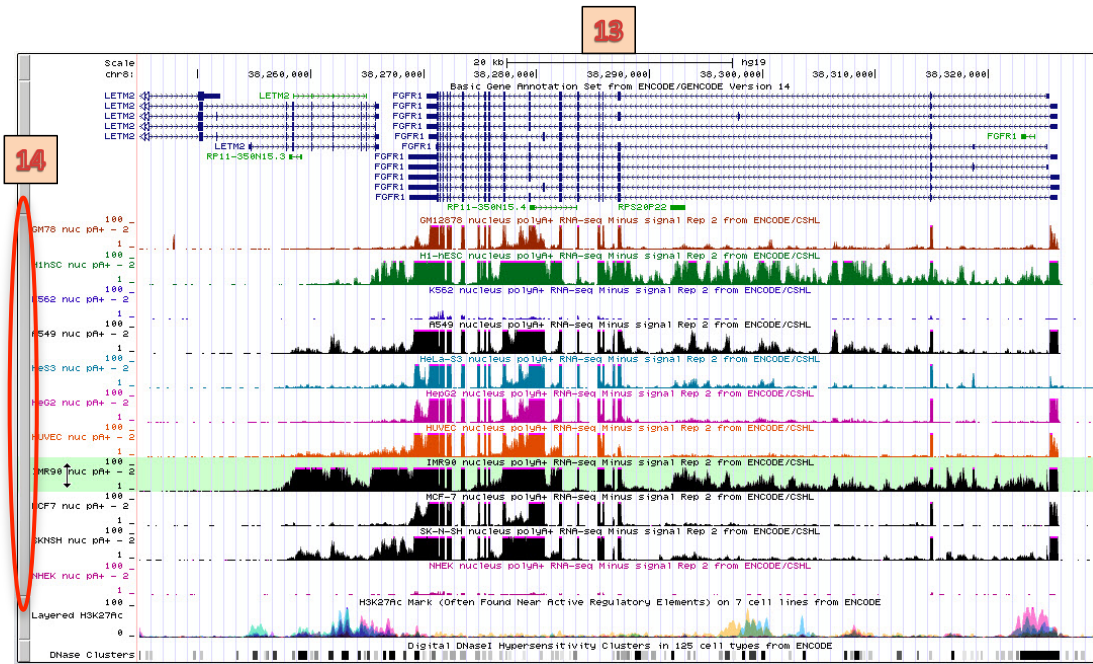
[Graph configuration help](#)

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

- GM12878 H3K27Ac Mark (Often Found Near Regulatory Elements) on GM12878 Cells from ENCODE [schema](#) 2009-10-05
- H1-hESC H3K27Ac Mark (Often Found Near Regulatory Elements) on H1-hESC Cells from ENCODE [schema](#) 2011-03-21
- HSMM H3K27Ac Mark (Often Found Near Regulatory Elements) on HSMM Cells from ENCODE [schema](#) 2010-09-16
- HUVEC H3K27Ac Mark (Often Found Near Regulatory Elements) on HUVEC Cells from ENCODE [schema](#) 2009-10-06
- K562 H3K27Ac Mark (Often Found Near Regulatory Elements) on K562 Cells from ENCODE [schema](#) 2009-10-05
- NHEK H3K27Ac Mark (Often Found Near Regulatory Elements) on NHEK Cells from ENCODE [schema](#) 2009-10-07
- NHLF H3K27Ac Mark (Often Found Near Regulatory Elements) on NHLF Cells from ENCODE [schema](#) 2010-06-28

1 of 7 selected

12



All		Whole
Cell Line	Localization	Cell Chromatin
	+ -	+ -
GM12878 (Tier 1)	<input type="checkbox"/>	<input type="checkbox"/>
H1-hESC (Tier 1)	<input type="checkbox"/>	<input type="checkbox"/>
K562 (Tier 1)	<input type="checkbox"/>	<input type="checkbox"/>
A549 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
B cells CD20+ (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
HeLa-S3 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
HepG2 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
HUVEC (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
IMR90 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
MCF-7 (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
Monocytes CD14+ (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
SK-N-SH (Tier 2)	<input type="checkbox"/>	<input type="checkbox"/>
AG04450	<input type="checkbox"/>	<input type="checkbox"/>

14

Cell, tissue or DNA sample: Cell line or tissue used as the source of experimental material.

cell	Tier	Description	Lineage	Tissue	Karyotype	Sex	Documents	Vendor ID	Term ID	Label
IMR90	2	fetal lung fibroblasts, newly promoted to tier 2. not in 2011 analysis	endo-derm	lung	normal	F	Stam	ATCC CCL-186	BTO:0001229	IMR90