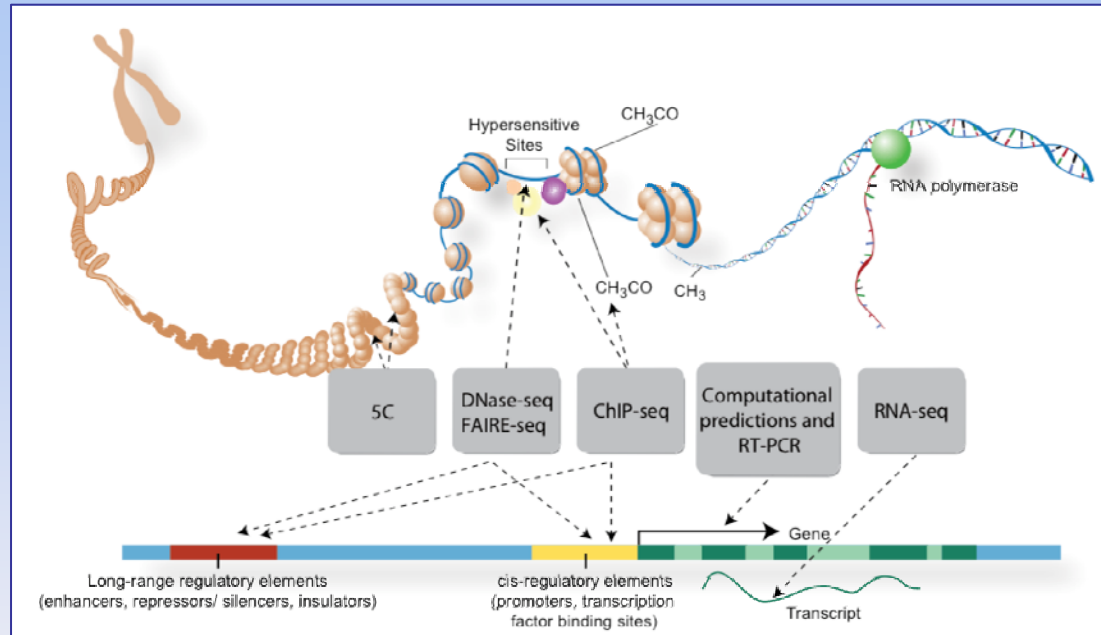


# ENCODE Project

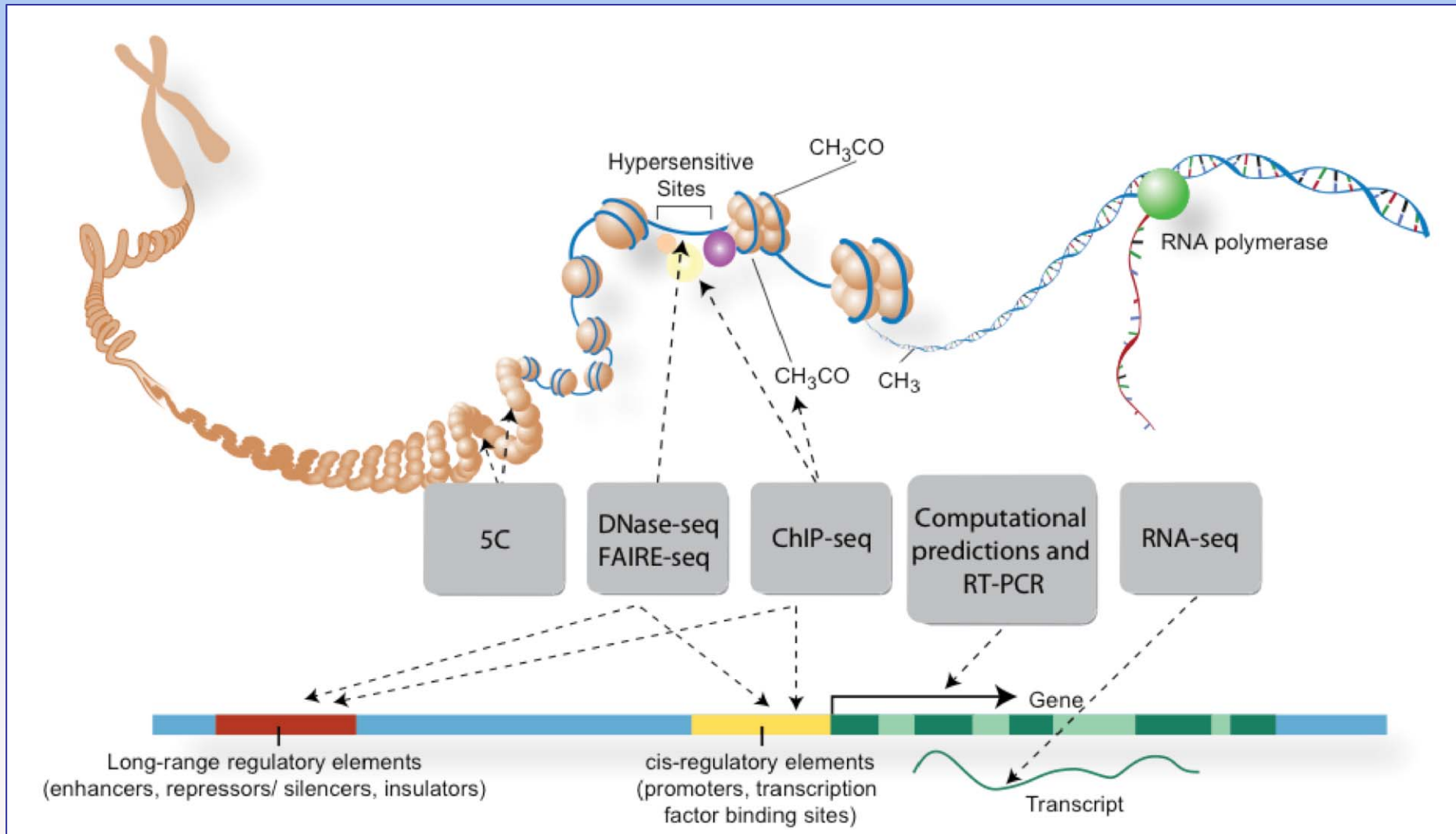


Robert Kuhn  
University of California Santa Cruz

Workshop: Working with ENCODE Data  
Korean Genome Organization  
February 5-7, 2014  
Yong Pyong, Korea

[http:// genome.ucsc.edu](http://genome.ucsc.edu)

[http:// encodeproject.org](http://encodeproject.org)



# ENCODE Project



- ENCODE pilot project covered only 1% of human genome. Phase II ENCODE is full-genome on human and mouse. Phase III is starting now. DCC at Stanford, UCSC hosts data.
- 32 biology labs organized into 19 grants + Analysis Working Group and Data Coordination Center (DCC)
- Goal: identify and characterize all **functional elements** of the genome.
- ENCODE DCC's job is to make data accessible and clear, to put it in UCSC Genome Browser, and to help other databases at NCBI, EBI, and elsewhere import ENCODE data as well.

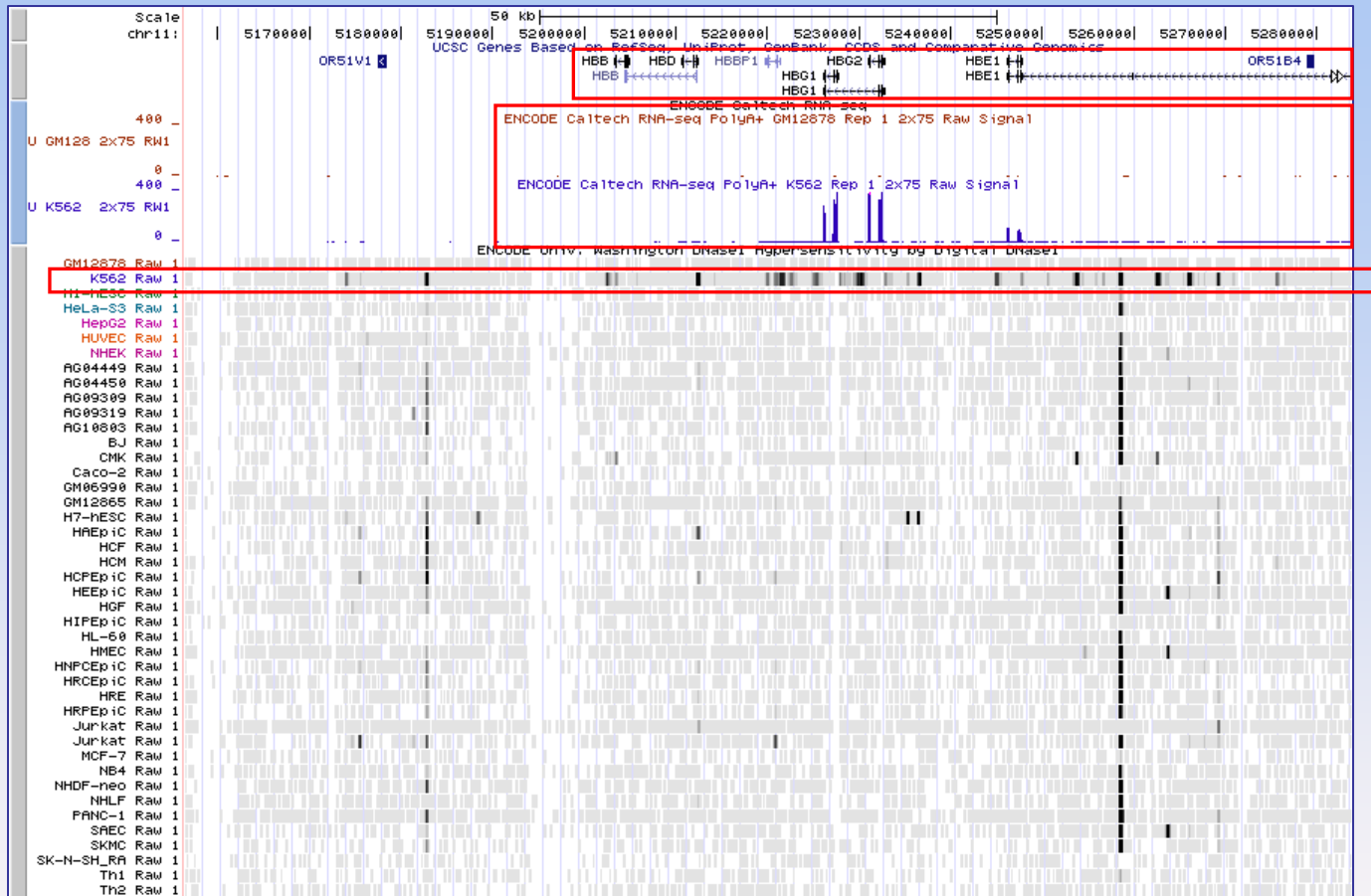
# ENCODE assays on regulation of transcription

- Opening/closing chromatin
  - DNase hypersensitivity, FAIRE-seq
  - Chromatin immunoprecipitation & sequencing (ChIP-seq) of histone marks
- Binding expressive/inhibitory transcription factors.
  - ChIP-seq of various transcription factors
- RNA transcription (or not)
  - mRNA sequencing of ENCODE cell lines
  - RNA seq fractionation: short/long polyA+/- localized to nucleus, cytoplasm, polysome, nucleoplasm, nuclear matrix, mitochondria, etc.

# ENCODE DNase Hypersensitivity

- Several genome-wide high-throughput methods used in ENCODE. All involve DNA-seq
- Data currently available for 388 cell lines and tissues
- Main alignment artifacts to watch for:
  - DNA present in cell in multiple copies:
    - Mitochondria, centromeric repeats, other repeats
    - Generally such regions ignored except in “raw” data.
  - Sequencing biases (highly GC-rich regions etc.)
  - In general, sequencing artifacts are easier to work around than those associated with DNA-chip based assays.

# UW DNaseI at Hemoglobin Beta



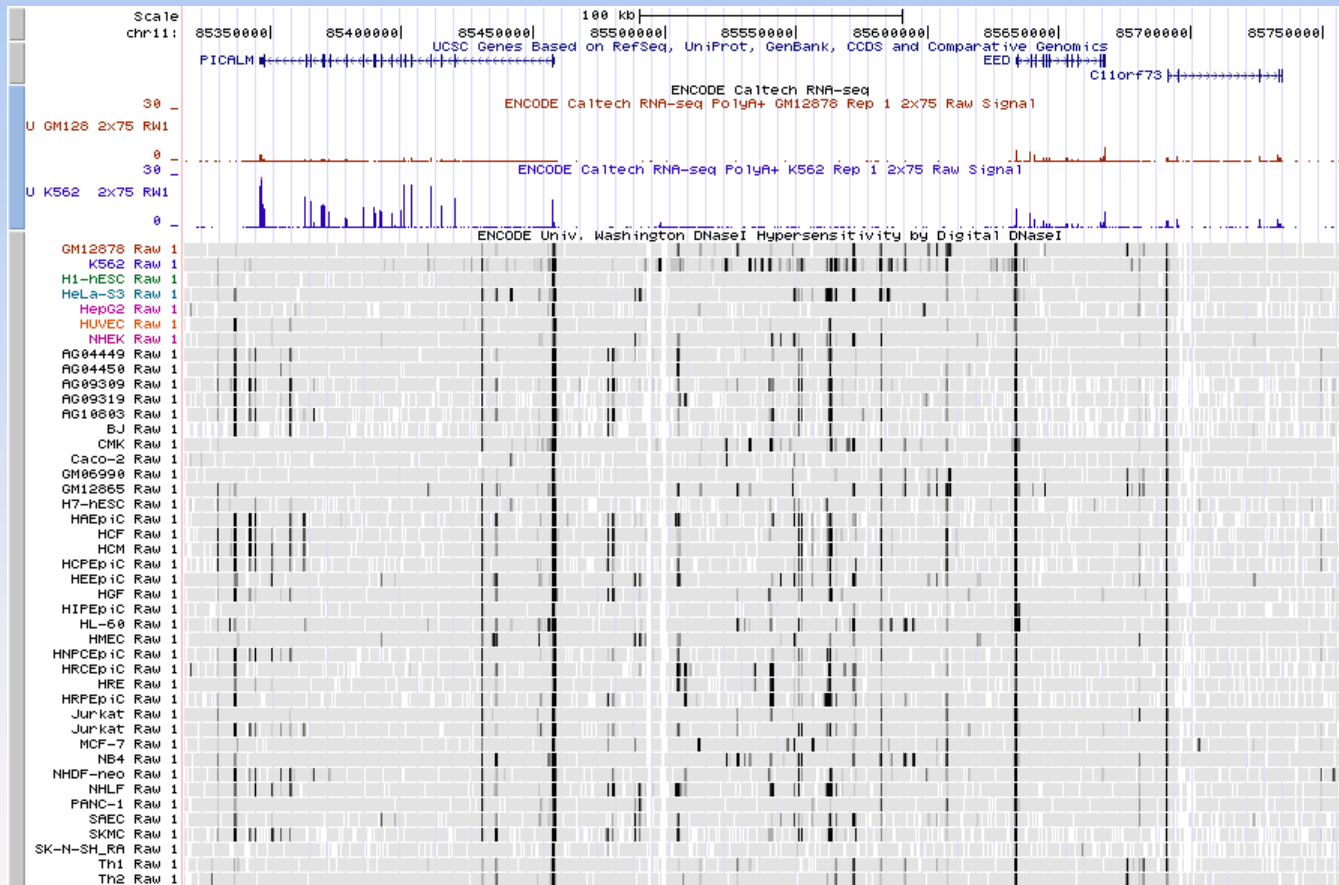
Top track shows **genes** in the Hemoglobin beta (HBB) locus.

Next track shows **RNA levels in GM12878 and K562** cell lines.

The last track is density plots of **DNase hypersensitivity** in many cell lines.

**K562**, a cell line similar to a red blood cell precursor, shows much RNA and DNase activity.

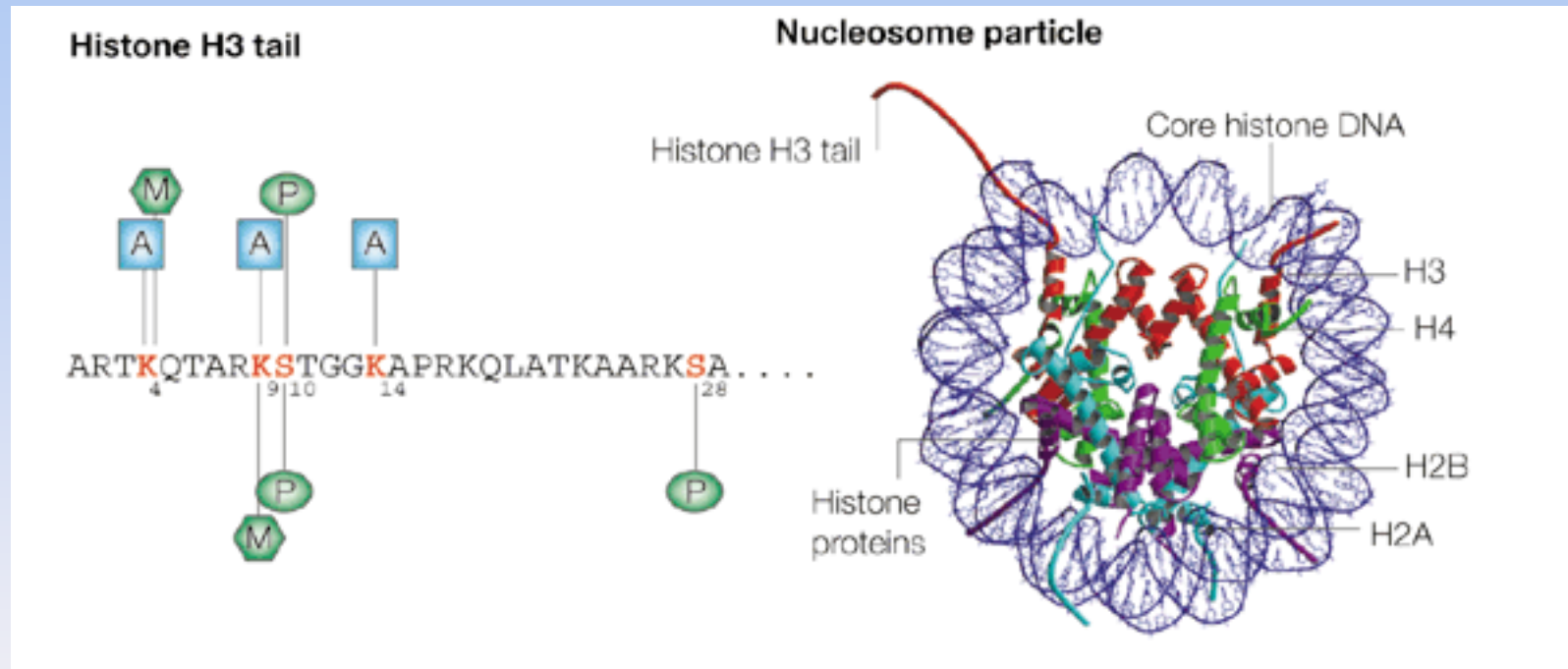
# A more typical locus - PICALM



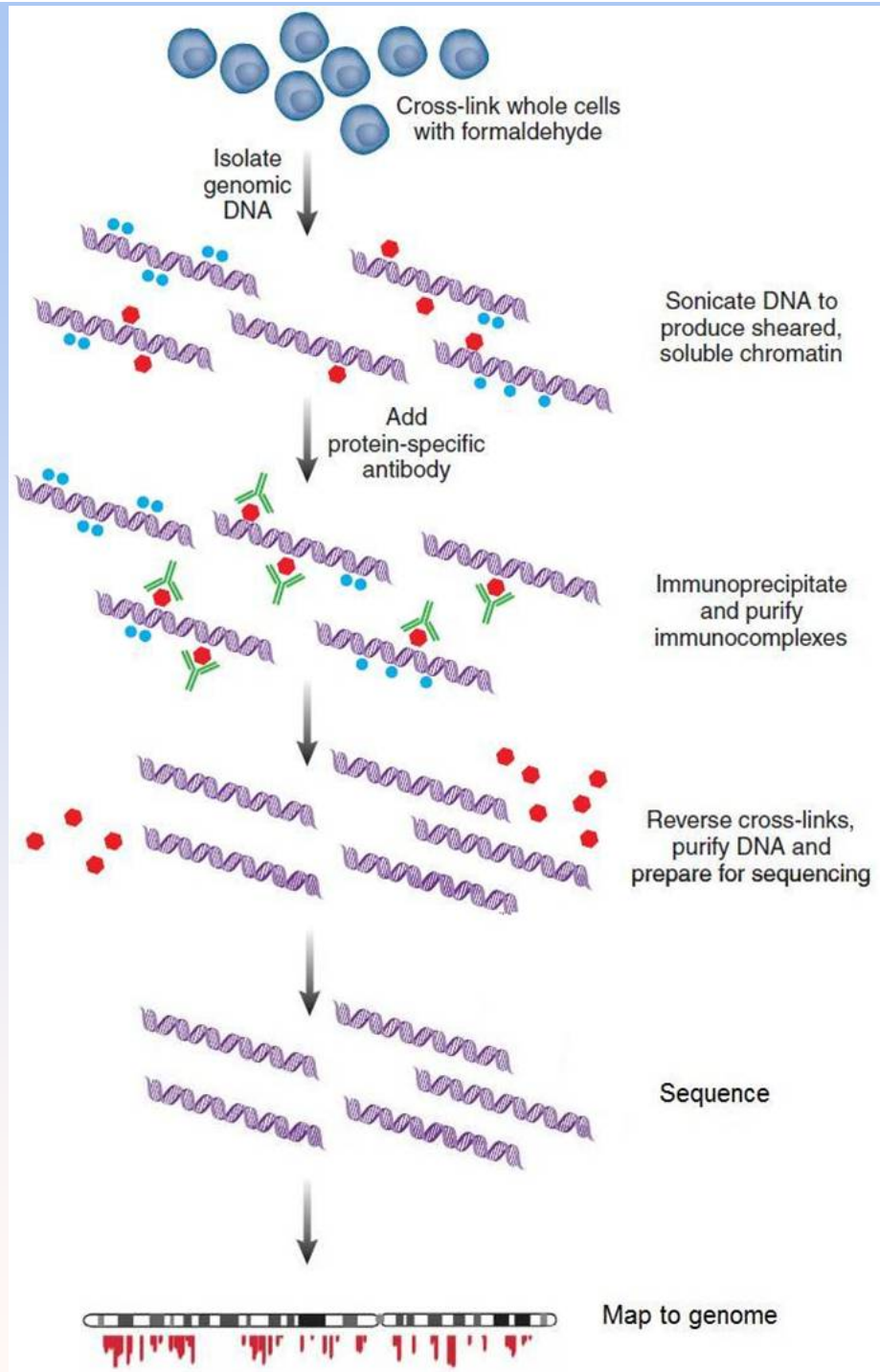
DNase patterns typically are less specific to a single cell type as seen here



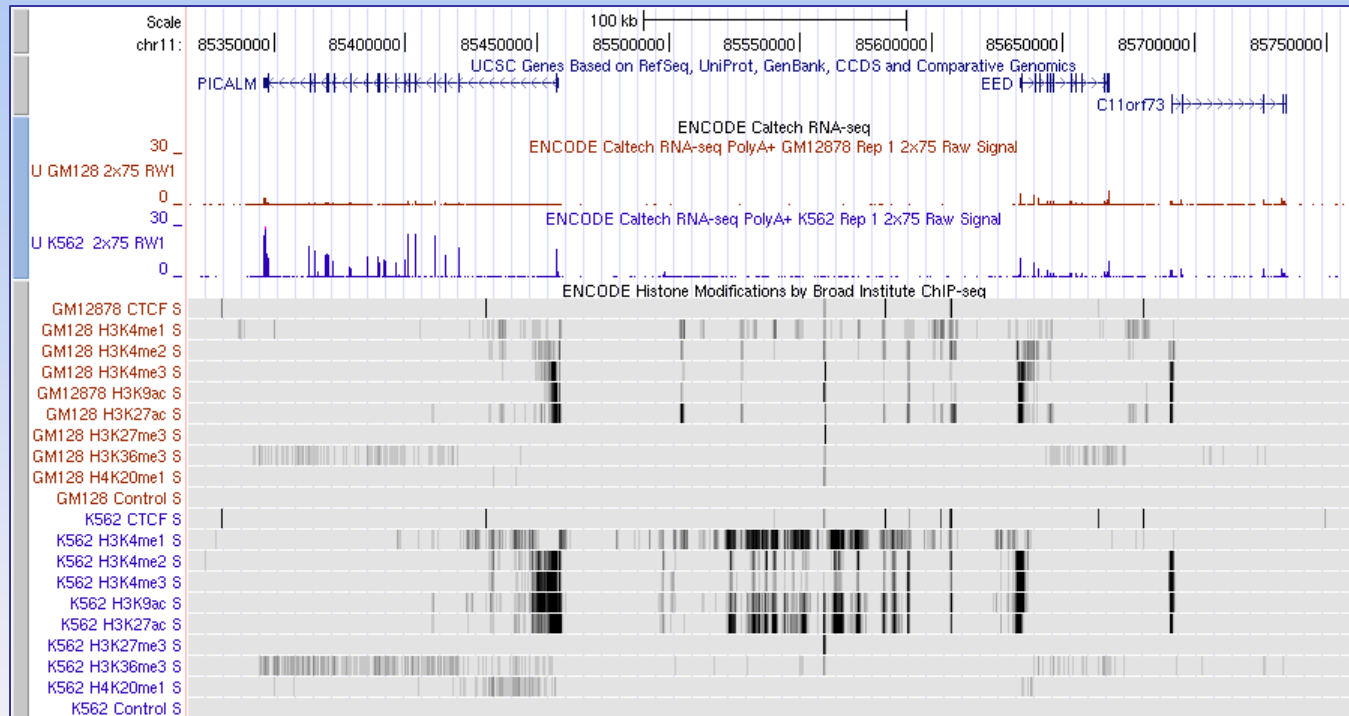
# Histone Mark and related ChIP-seq



- Various histone marks give a broad picture of promoters, enhancers, repressed regions, transcribed regions
- ENCODE data sets currently include 12 histone marks + CTCF (insulator mark) in 67 cell lines. ~12 cell lines have near complete histone mark coverage

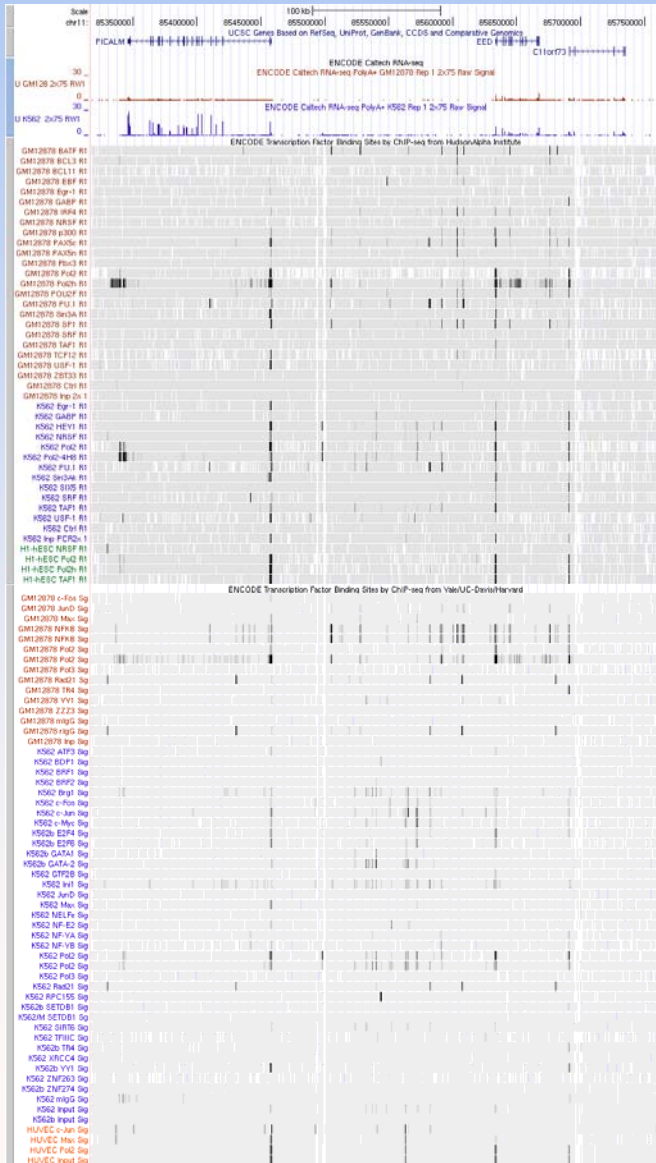


# Histone marks on 2 cell lines



Histone mark data at the same locus in two cell lines, GM12878 (red) and K562 (blue). Different marks are associated with promoters, transcribed regions, silencers, enhancers, etc. Most marks are darker in K562, which is more actively transcribing this region.

# Transcription Factor ChIP-Seq



ENCODE has data on 160 factors – in many cell lines where they are expressed.

# Making data fit on a single screen

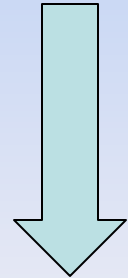
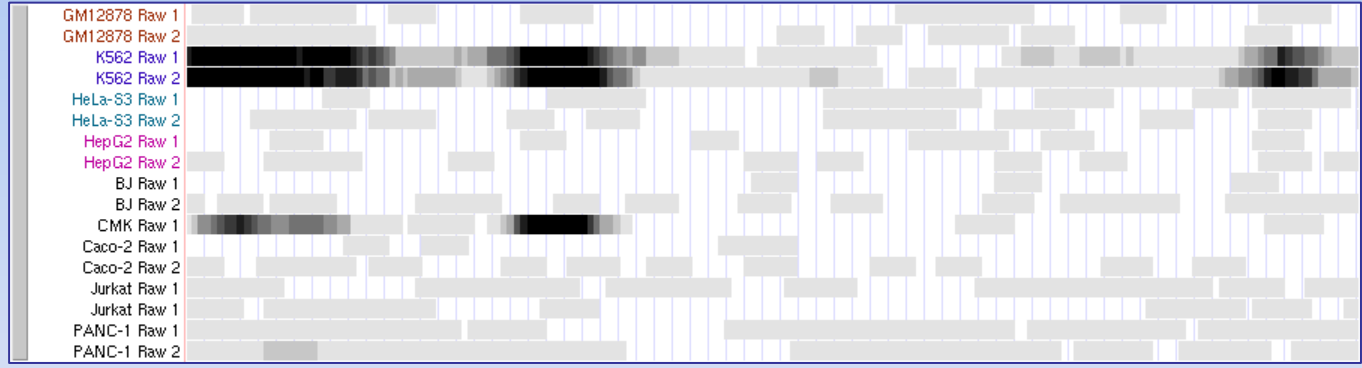
- All of the ENCODE data are excellent, but there is so much of it, it can be hard to know if you've seen everything relevant.
- Problem most acute in transcription factor ChIP-seq, but really a problem everywhere.
- UCSC has developed several ways of visually summarizing the data.

# Integrating DNase across cell lines

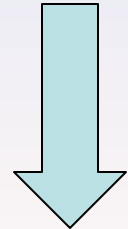
HBB Gene



DNaseI  
signal



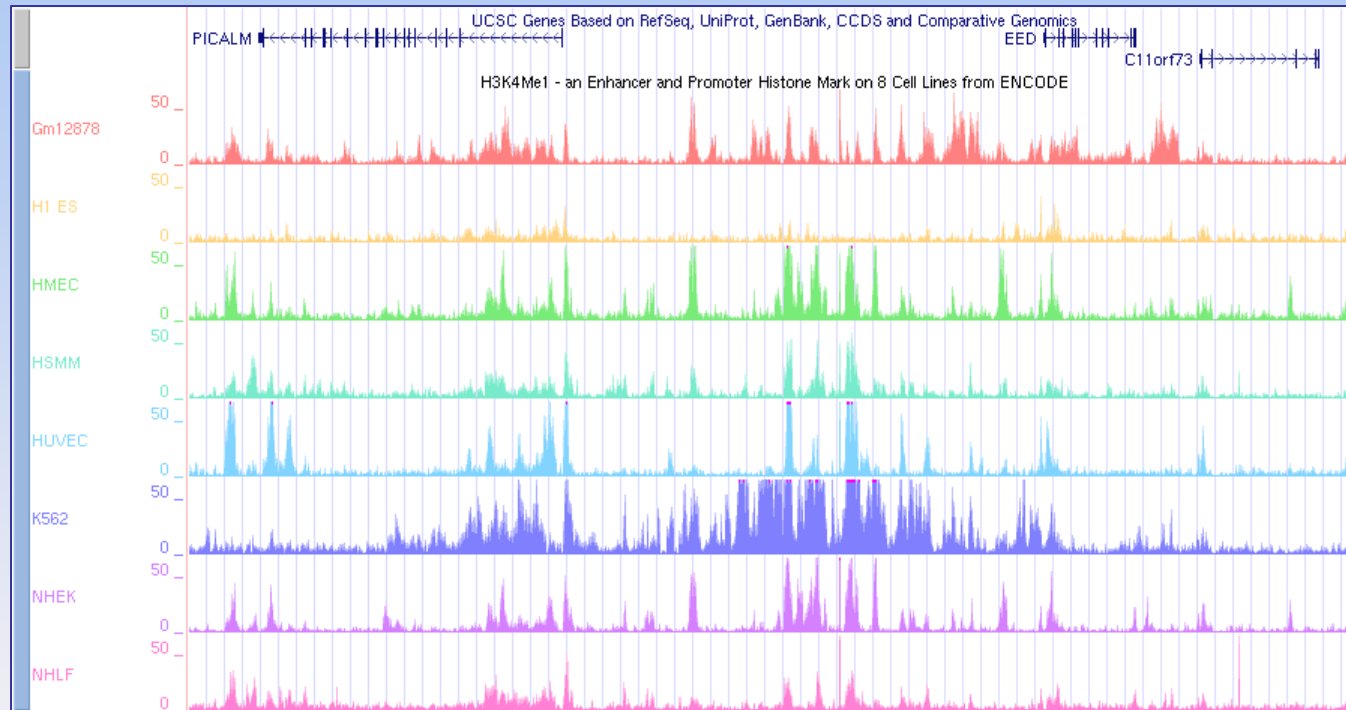
peaks



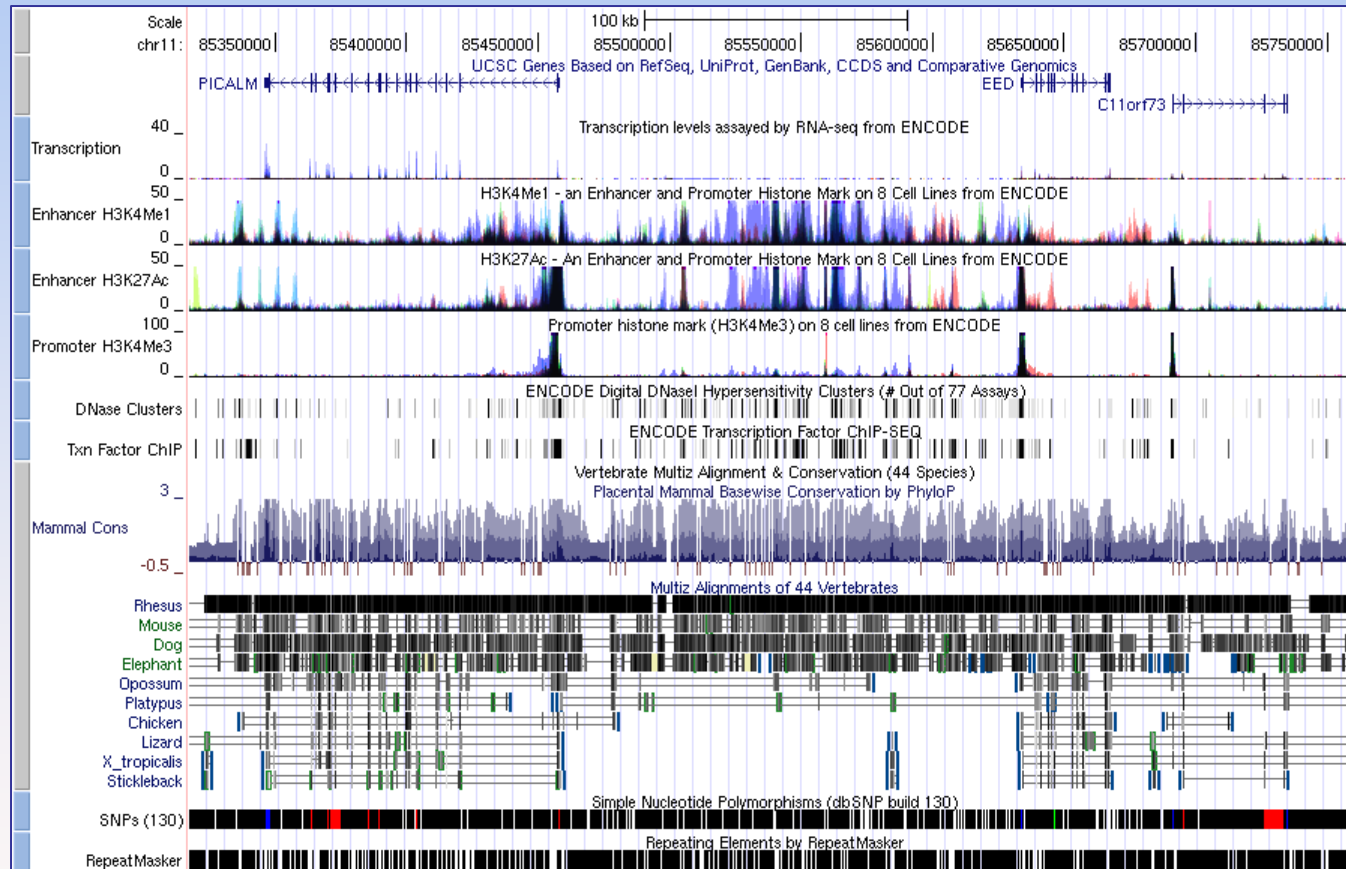
clustered  
peaks



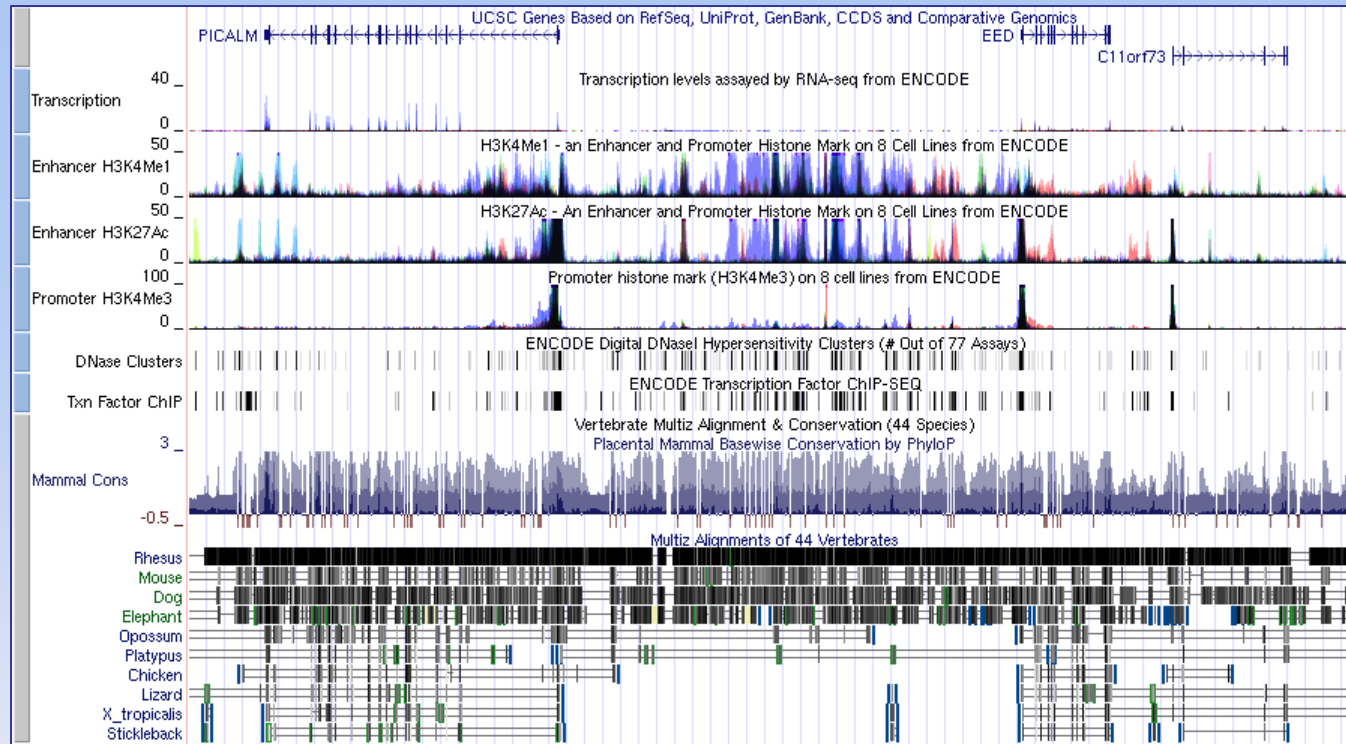
# Rainbow overlay for histone marks



# Integrated regulatory tracks in context with other genomics information at UCSC







- ENCODE regulatory data:
  - Histone marks –characterization of large regions into promoter/enhancer/repressed
  - DNase hypersensitivity - defines smaller regions as regulatory
  - Transcription factor chromatin immunoprecipitation – what regulatory factors bind in a smaller region.
  - Chromatin conformation capture – just starting to ramp up.
- Available at <http://genome.ucsc.edu>

# Accessing ENCODE Data at DCC

- <http://www.encodeproject.org>
  - ENCODE portal. Describes project overall, project news, tables and spreadsheets for all experiments
- <http://genome.ucsc.edu>
  - ENCODE data integrated into UCSC Genome Browser on hg19 and mm9 assemblies

Much of the data also is at NCBI (GEO) and Ensembl.



# Encyclopedia of DNA Elements

## Human

Experiment List

Search

Downloads

Genome Browser (hg19)

Preview Browser (hg19)

Session Gallery

Cell Types

## Mouse

Data Summary

Search

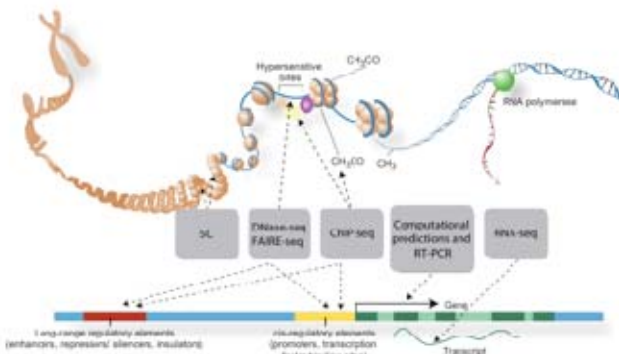
Downloads

Genome Browser (mm9)

Preview Browser (mm9)


## About ENCODE Data

The [Encyclopedia of DNA Elements \(ENCODE\) 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.



[Click to enlarge](#)

ENCODE data are now available for the entire human genome. **All ENCODE data are free and available for immediate use via :**

- [Search](#) for displayable tracks and downloadable files
- [Download](#) of data files
- [Visualization](#) in the UCSC Genome Browser (ENCODE data marked with the  NHGRI logo)
- [Data mining](#) with the UCSC Table Browser and other [UCSC Genome Bioinformatics tools](#)

To search for ENCODE data related to your area of interest and set up a browser view, use the UCSC [Track Search tool](#) (*Advanced features*). The [Data Summary](#) shows a comprehensive listing of ENCODE data that is released or in preparation. Early access to pre-release ENCODE data is provided at <http://genome-preview.ucsc.edu>. If you would like to receive notifications of ENCODE data releases and related news by email, subscribe to the [encode-announce mailing list](#). For more information about how to access this data, see the free online [OpenHelix ENCODE tutorial](#).

To complement the human ENCODE data, Mouse ENCODE experiments are currently underway. Early access to this data is available on the Mouse mm9/NCBI37 browser at the UCSC preview site. The [Mouse ENCODE Data Summary](#) lists experiments that are planned or in progress.

All ENCODE data is freely available for download and analysis. However, before publishing research that uses ENCODE data, please read the [ENCODE Data Release Policy](#), which places some restrictions on publication use of data for nine months following data release. [Read more](#) about ENCODE data at UCSC.

	DNA Methylation			Open Chromatin			RNA Binding Proteins				RNA Profiling					TFBS & Histones		Other		
	Methyl Array	Methyl RRBS	Methyl-seq	DNase-DGF	DNase-seq	FAIRE-seq	RIP Gene ST	RIP Tiling Array	RIP Validation	RIP-seq	CAGE	Exon Array	RNA-chip	RNA-PET	RNA-seq	ChIP-seq	view matrix	5C	ChIA-PET	Combined
•	2	1	1		2	1	7	4		4	6	2	6	2	14	112		2		2
•	2	1	1		2	1	3				4	1		1	13	63		1		2
•	2	1	1	1	3	3	6	4		4	9	3	9	6	24	178		2	2	2
•	1	1			2	1					3	2			17	48				
•											1									
•				1	1										1	2				
•	1	1	1		3	3	4				5	4		3	11	84		1	1	2
•	2	1	1	1	2	1	4				6	2	5	2	11	103		1		2
•	1			1	2	1					5	2		2	9	33				2
•	2	1									3				9					
•															2	7				
•	2	1			3	1					3	7			10	16		1	3	

Experiment matrix link off of ENCODE Portal, provides overview of number of experiments of various types on various cells. Clicking on a cell brings up list of individual tracks or files. It's a big matrix, note size on scrollbar.



# ChIP-seq Experiment Matrix hg19

## Antibody Targets

search for:  tracks  files

### Cell Types

Tier 1	
GM12878	•
H1-hESC	•
K562	•
Tier 2	
A549	•
CD20+_RO01778	•
CD20+_RO01794	•
HeLa-S3	•
HepG2	•
HUVEC	•

Histone Modification	Antibody Targets											Transcription Factor												
	H2AZ	H3K27ac	H3K27me3	H3K36me3	H3K4me1	H3K4me2	H3K4me3	H3K79me2	H3K9ac	H3K9me1	H3K9me3		H4K20me1	AP-2alpha	AP-2gamma	ATF2	ATF3	BAF155	BAF170	BATF	BCL11A	BCL3	BCLAF1	BDP1
GM12878	1	1	2	2	1	1	2	1	1		1	1			1	1			1	1	1	1		
H1-hESC	1	1	1	1	1	1	1	1	1		1	1			1	1				2				
K562	1	1	3	2	2	1	8	1	2	1	1	1				2					1	1	1	
A549							1									1					1			
CD20+_RO01778							1																	
CD20+_RO01794							1																	
HeLa-S3	1	1	2	2	1	1	2	1	1		1	1		1	1		1	1						1
HepG2	1	1	2	2	1	1	2	1	1			1				1								
HUVEC	1	1	2	2	1	1	2	1	1	1		1												

ChIP-seq experiments have their own submatrix. This is an even bigger matrix. Note size on both horizontal and vertical scrollbars.

# Track Search

- Can do a free-form (Google-style) search or search metadata field-by-field

The screenshot shows the UCSC Genome Browser Gateway search interface. At the top is a navigation bar with links: Home, Genomes, Blat, Tables, Gene Sorter, PCR, Session, FAQ, Help. Below this is the title "Human (*Homo sapiens*) Genome Browser Gateway". A paragraph of text states: "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." The search form contains several fields: "clade" (Mammal), "genome" (Human), "assembly" (Feb. 2009 (GRCh37/hg19)), "position or search term" (chr21:33,031,597-33,041,570), "gene" (empty), and "image width" (800). A "submit" button is to the right. Below the form, there is a link: "Click [here](#) to [reset](#) the browser user interface settings to their defaults." At the bottom, there are four buttons: "track search" (highlighted with a red box), "add custom tracks", "configure tracks and display", and "clear position".

# Free-text search

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

**Search for Tracks in the Human Mar. 2006 (NCBI36/hg18) Assembly**

Search **Advanced**

H3K4me K562 Chip-seq

search clear cancel

+ -	Visibility	Track Name
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me1 S</a> ENCODE Histone Mods, Broad ChIP-seq Signal (H3K4me1, K562) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me1 P</a> ENCODE Histone Mods, Broad ChIP-seq Peaks (H3K4me1, K562) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me3 S</a> ENCODE Histone Mods, Broad ChIP-seq Signal (H3K4me3, K562) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me3 P</a> ENCODE Histone Mods, Broad ChIP-seq Peaks (H3K4me3, K562) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me2 S</a> ENCODE Histone Mods, Broad ChIP-seq Signal (H3K4me2, K562) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me2 P</a> ENCODE Histone Mods, Broad ChIP-seq Peaks (H3K4me2, K562) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me3 H1</a> ENCODE UW Histone ChIP Hotspots - 1st (H3K4me3 in K562 cells) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me3 S1</a> ENCODE UW Histone ChIP Raw Signal - 1st (H3K4me3 in K562 cells) ...
<input type="checkbox"/>	hide	<a href="#">K562 H3K4me3 P1</a> ENCODE UW Histone ChIP Peaks (FDR 0.5%) - 1st (H3K4me3 in K562 cells) ...

Return to Browser (0 of 9 selected)

# Advanced field-by-field search

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

Search for Tracks in the Human Mar. 2006 (NCBI36/hg18) Assembly

**Search** **Advanced**

Track Name: contains

and Description: contains

and Group: is Any

and Data Format: is Signal (wig) - wiggle format

*ENCODE terms*

+ and Cell, tissue or DNA sample  is HUVEC  Cell, tissue or DNA sample

+ and Antibody or target protein  is CTCF  Antibody or target protein

+ -	Visibility	Track Name
<input type="checkbox"/>	hide	<a href="#">HUVEC CTCF S</a> ENCODE Histone Mods, Broad ChIP-seq Signal (CTCF, HUVEC) ...
<input type="checkbox"/>	hide	<a href="#">HUVEC CTCF FD</a> ENCODE Open Chromatin, UT ChIP-seq F-Seq Density Signal (CTCF in HUVEC cells) ...
<input type="checkbox"/>	hide	<a href="#">HUVEC CTCF BO</a> ENCODE Open Chromatin, UT ChIP-seq Base Overlap Signal (CTCF in HUVEC cells) ...
<input type="checkbox"/>	hide	<a href="#">HUVEC CTCF S1</a> ENCODE UW Histone ChIP Raw Signal - 1st (CTCF in HUVEC cells) ...
<input type="checkbox"/>	hide	<a href="#">HUVEC CTCF S2</a> ENCODE UW Histone ChIP Raw Signal - 2nd (CTCF in HUVEC cells) ...

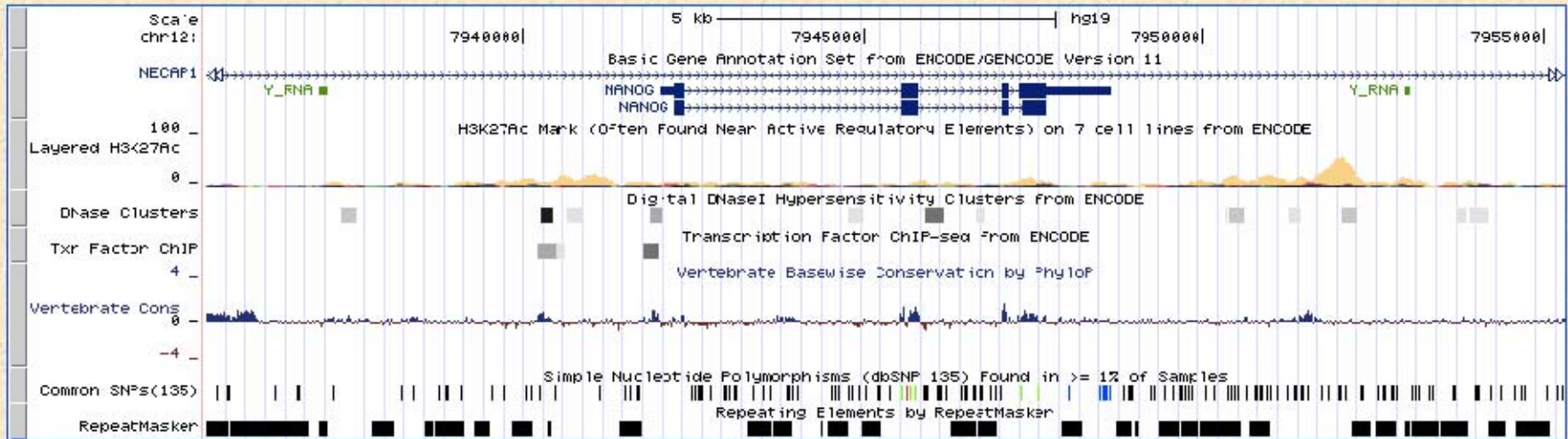
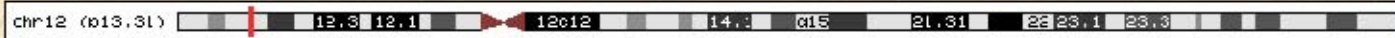
(0 of 5 selected)



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

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

position/search chr12:7,935,334-7,955,316  jump clear size 19,983 bp. configure



move start < 2.0 > Click on a feature for details. Click or drag in the base position track to zoom in. Click side bars for track options. Drag side bars or labels up or down to reorder tracks. Drag tracks left or right to new position. < 2.0 > move end

track search default tracks default order hide all add custom tracks track hubs configure reverse resize refresh

collapse all Use drop-down controls below and press refresh to alter tracks displayed. expand all Tracks with lots of items will automatically be displayed in more compact modes.

- + Mapping and Sequencing Tracks refresh
- + Phenotype and Disease Associations refresh
- Genes and Gene Prediction Tracks refresh

UCSC Genes Old UCSC Genes Alt Events GENCODE Genes V11 GENCODE Genes V10 GENCODE Genes V7

hide hide hide pack hide hide

# Acknowledgements

Jim Kent

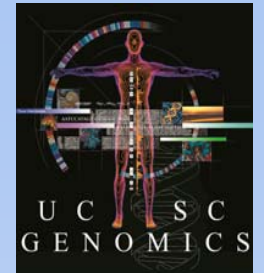
Kate Rosenbloom

Tim Dreszer

Katrina Learned

Brian Lee

# Browser Team



# Hands-on Exercises

p 31 – questions, steps (Module 3)

p 22 – solutions (screen grab) (Module2)

http:// genome.ucsc.edu

My data, Sessions:



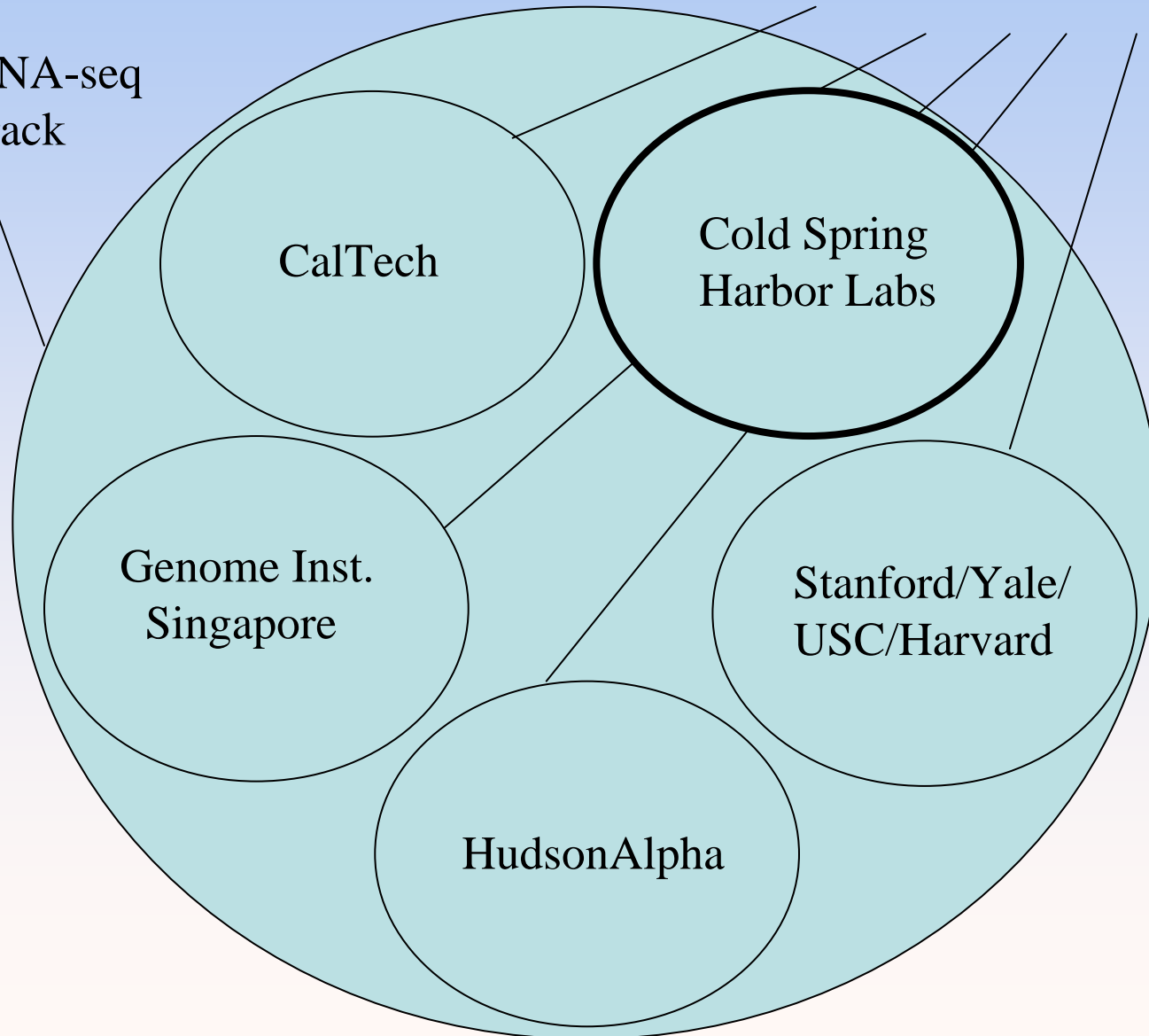
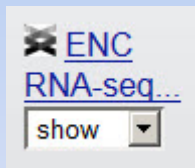
user: korea2014

pwd: \*\*\*\*\*

# ENCODE RNA-seq track

5 composite tracks,  
each with many subtracks

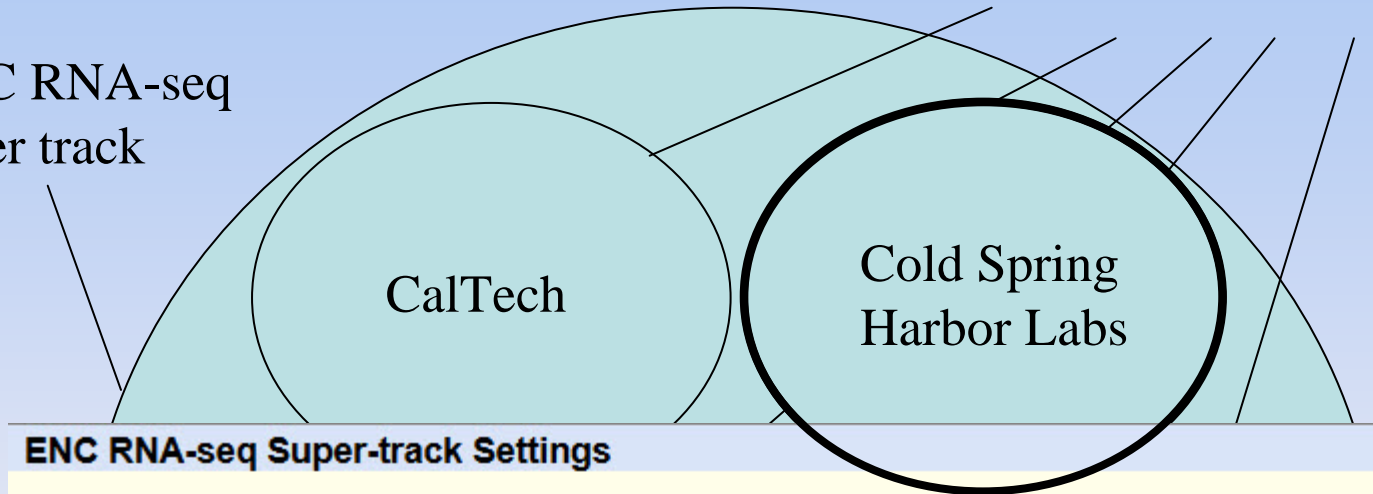
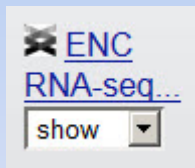
ENC RNA-seq  
super track



# ENCODE RNA-seq track

5 composite tracks,  
each with many subtracks

ENC RNA-seq  
super track



ENC RNA-seq Super-track Settings



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

Display mode:

**All**

- [Caltech RNA-seq](#) RNA-seq from ENCODE/Caltech
- [CSHL Long RNA-seq](#) Long RNA-seq from ENCODE/Cold Spring Harbor Lab
- [GIS RNA-seq](#) RNA-seq from ENCODE/Genome Institute of Singapore
- [HAIB RNA-seq](#) RNA-seq from ENCODE/HAIB
- [SYDH RNA-seq](#) RNA-seq from ENCODE/Stanford/Yale/USC/Harvard

# ENCODE RNA-seq track

Multi-dimensional data

Cold Spring  
Harbor Labs

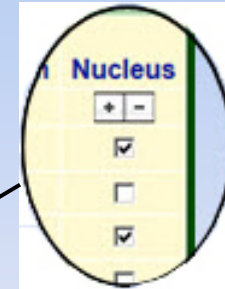
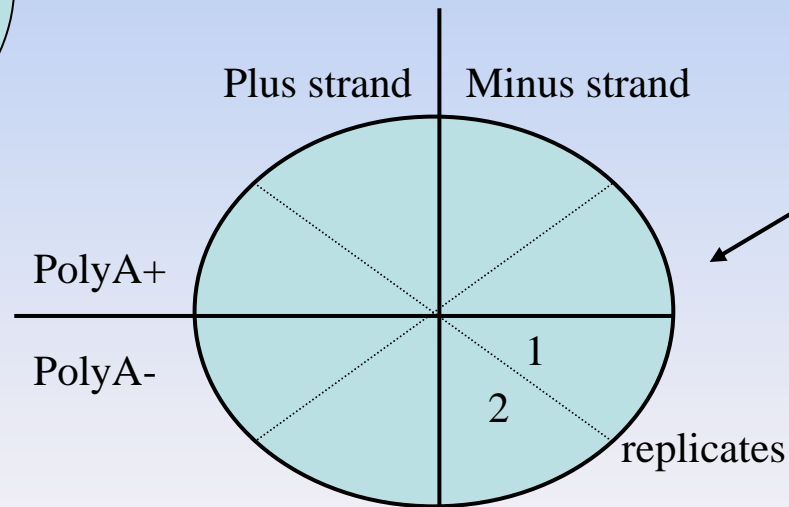
Cell Line	Localization	Whole	Chromatin	Cytosol	Nucleolus	Nucleoplasm	Nucleus
		Cell					
		+ -	+ -	+ -	+ -	+ -	+ -
GM12878 (Tier 1)		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>		<input checked="" type="checkbox"/>
H1-hESC (Tier 1)		<input type="checkbox"/>		<input type="checkbox"/>			<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"/>					
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"/>
HUVEC (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
IMR90 (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
MCF-7 (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
Monocytes CD14+ (Tier 2)		<input type="checkbox"/>					
SK-N-SH (Tier 2)		<input type="checkbox"/>		<input type="checkbox"/>			<input type="checkbox"/>
AG04450		<input type="checkbox"/>					
BJ		<input type="checkbox"/>					
CD34+ Mobilized		<input type="checkbox"/>					
HAoAF		<input type="checkbox"/>					
HAoEC		<input type="checkbox"/>					



# ENCODE RNA-seq track



Multi-dimensional data



Cell Line <sup>1</sup>	Localization <sup>2</sup>	RNA Extract <sup>3</sup>	Views <sup>4</sup>	Replicate
GM12878	Nucleus	PolyA-	Minus Signal	1st
GM12878	Nucleus	PolyA-	Minus Signal	2nd
GM12878	Nucleus	PolyA-	Plus Signal	1st
GM12878	Nucleus	PolyA-	Plus Signal	2nd
GM12878	Nucleus	PolyA+	Minus Signal	1st
GM12878	Nucleus	PolyA+	Minus Signal	2nd
GM12878	Nucleus	PolyA+	Plus Signal	1st
GM12878	Nucleus	PolyA+	Plus Signal	2nd