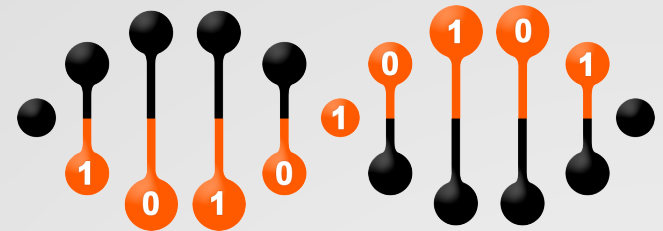


CBC Data Therapy

ChIP-Seq Discussion



Computational Biology Core



ChIP-seq big picture

Combine “Next-Generation” sequencing with Chromatin Immunoprecipitation to identify genomewide chromatin binding sites.

Select (and identify) fragments of DNA that interact with specific proteins such as:

- Transcription factors

- Modified histones

- RNA Polymerase (survey actively transcribe portions of the genome)

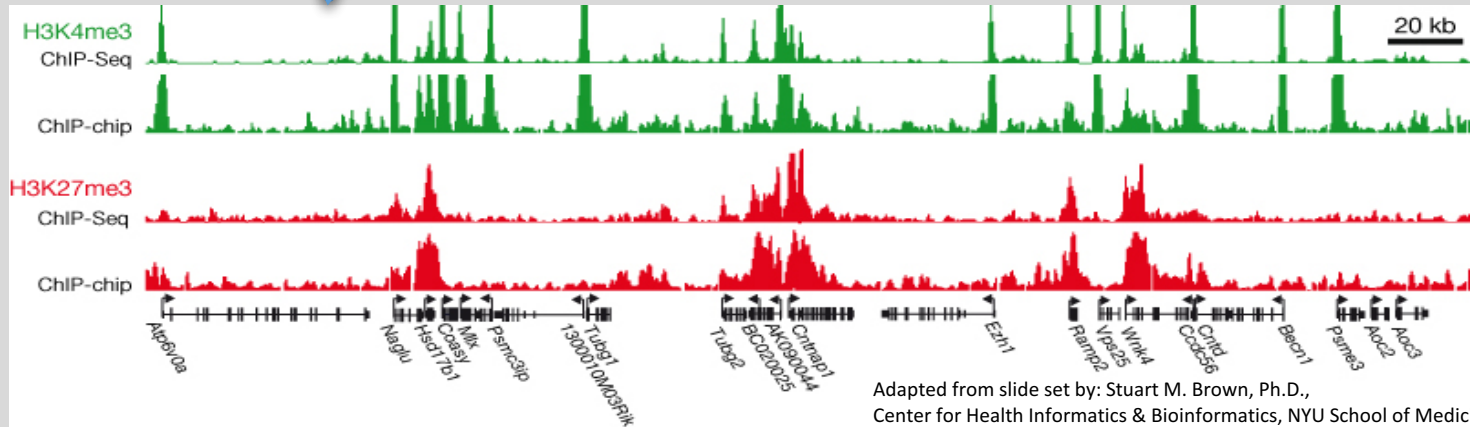
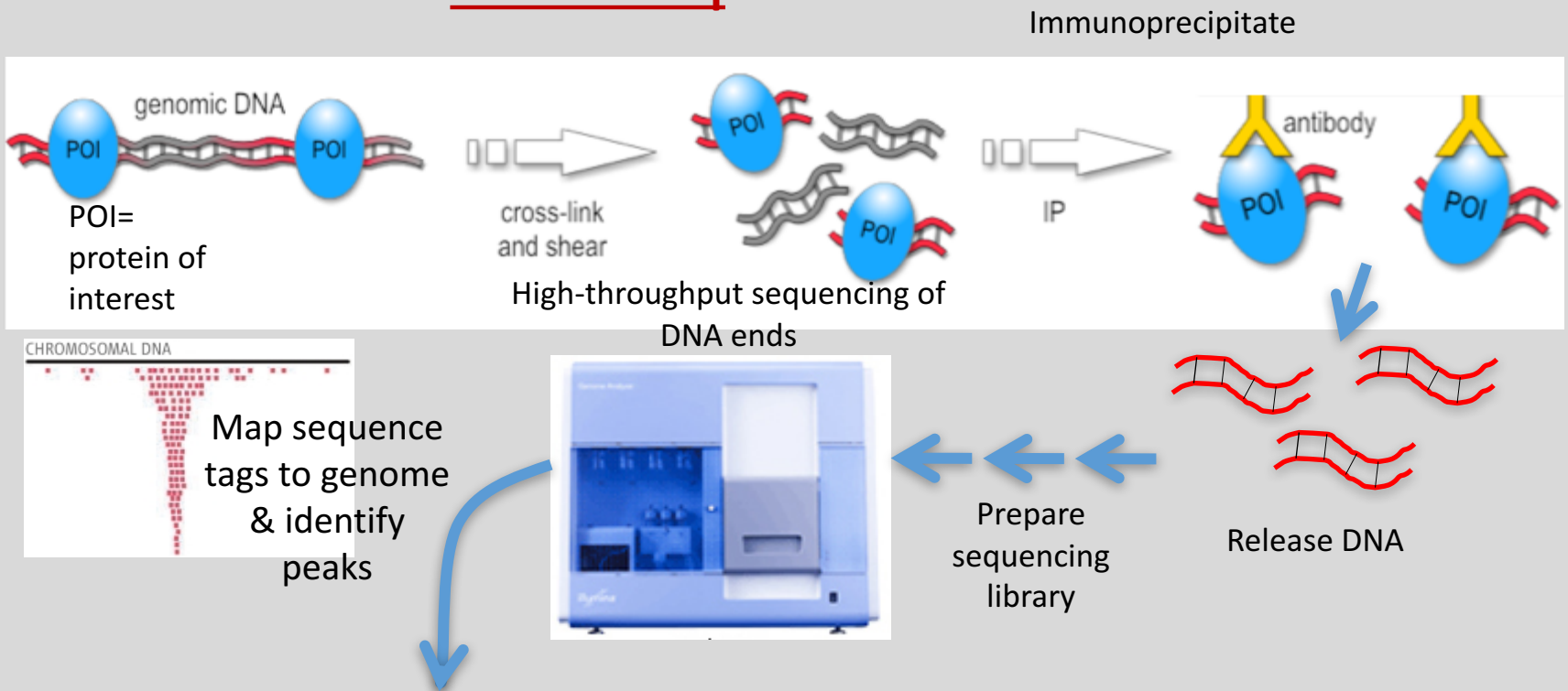
- DNA polymerase (investigate DNA replication)

- DNA repair enzymes

- Or fragments of DNA that are modified: e.g. CpG methylation



ChIP-seq



Adapted from slide set by: Stuart M. Brown, Ph.D.,
Center for Health Informatics & Bioinformatics, NYU School of Medicine



How many reads do I need?

- (a) Number of binding sites
- (b) Genome size
- (c) Specificity of antibody/method

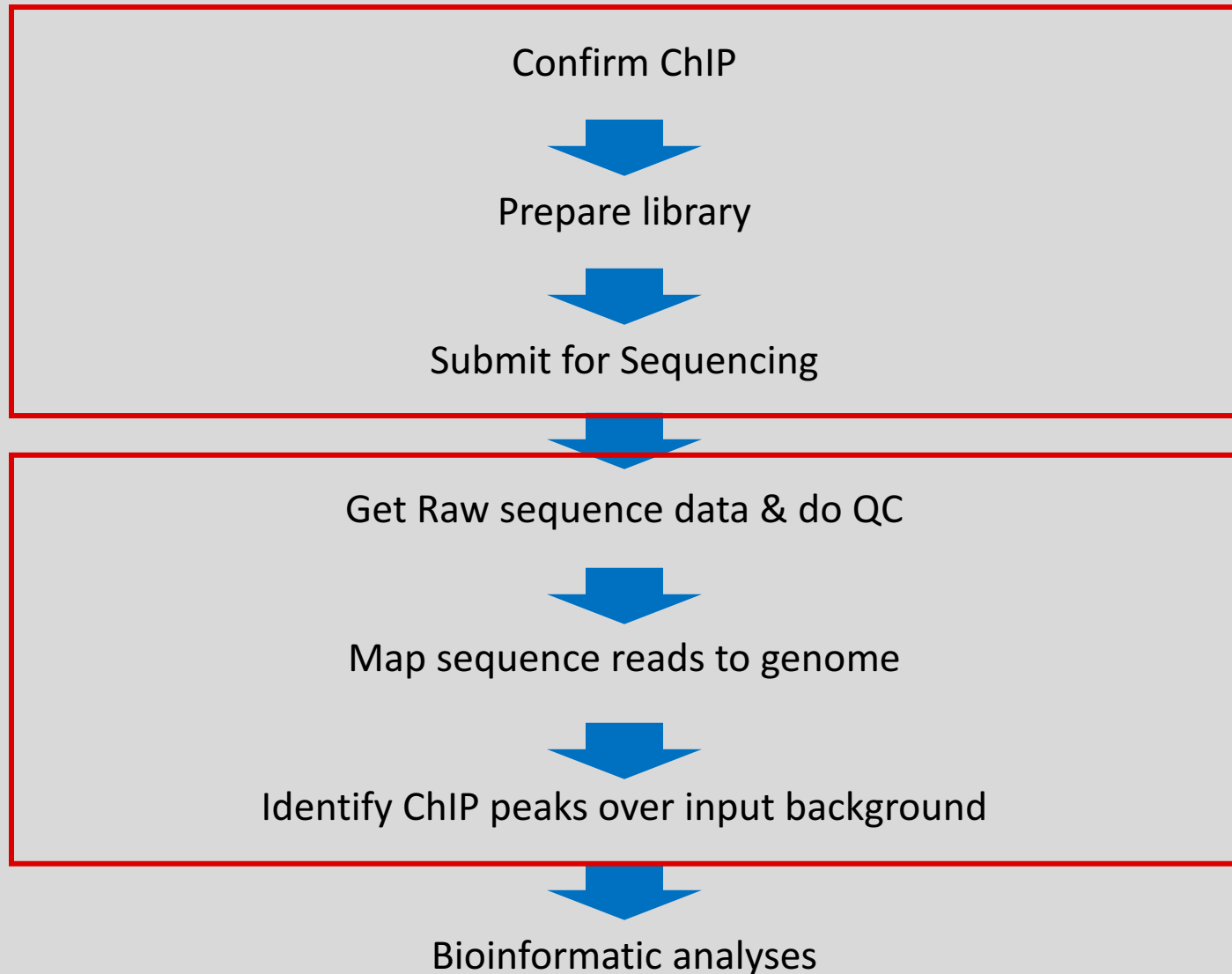
Minimum for ChIP-seq of a transcription factor with $< \sim 30,000$ binding sites in a mammalian genome:

- **3 replicates per condition**
- **20+ million reads per sample** ($>10X$ coverage of binding genome size, proportionately less for smaller genomes & fewer binding peaks)
- **Paired end shorter reads almost always good**

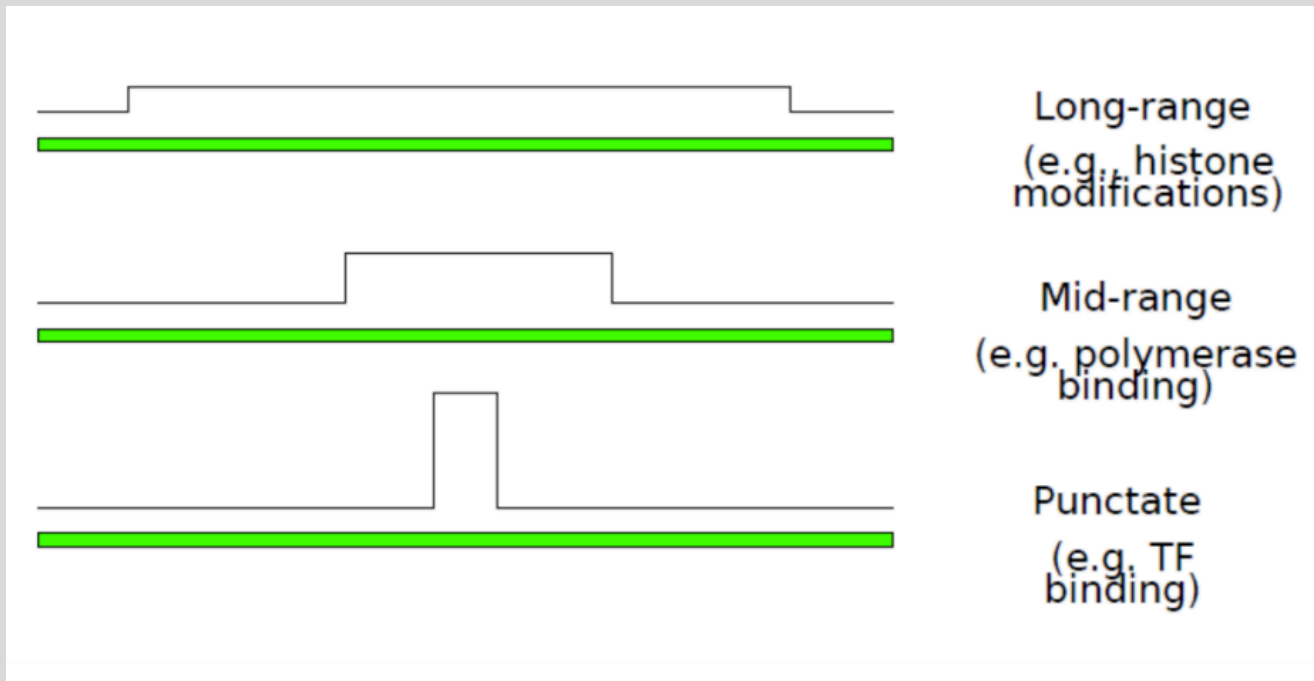
For some applications need many more reads (e.g. mapping nucleosome positions need >400 M). Make your best estimate. If you have too few you can re-sequence the same samples or add additional samples. Reads from all runs can be pooled in the end.



ChIP-seq Workflow



Types of events of interest



Peak calling algorithms/Softwares

1. The most popular peak caller by Tao Liu: [MACS2](#). Now --broad flag supports broad peaks calling as well.
2. [TF ChIP-seq peak calling using the Irreproducibility Discovery Rate \(IDR\) framework](#) and many [Software Tools Used to Create the ENCODE Resource](#)
3. [SICER](#) for broad histone modification ChIP-seq
4. [HOMER](#) can also used to call Transcription factor ChIP-seq peaks and histone modification ChIP-seq peaks.
5. [MUSIC](#)
6. [permseq](#) R package for mapping protein-DNA interactions in highly repetitive regions of the genomes with prior-enhanced read mapping. [Paper](#) on PLoS Comp.
7. [Ritornello](#): High fidelity control-free chip-seq peak calling. No input is required!
8. Tumor samples are heterogeneous containing different cell types. [MixChIP: a probabilistic method for cell type specific protein-DNA binding analysis](#)
9. [Detecting broad domains and narrow peaks in ChIP-seq data with hiddenDomains tool](#)
10. [BroadPeak: a novel algorithm for identifying broad peaks in diffuse ChIP-seq datasets](#)
11. [epic: diffuse domain ChIP-Seq caller based on SICER](#). It is a re-written of SICER for faster processing using more CPUs. (Will try it for broad peak for sure).
12. [Cistrome](#): The best place for wet lab scientist to check the binding sites. Developed by Shierly Liu lab in Harvard.
13. [Accounting for GC-content bias reduces systematic errors and batch effects in ChIP-Seq peak callers](#) tool in [github](#)



MACS2 : Examine the output

Start with your .macsinfo bsub -oo file.

vi LiE_ERaIPvINPUT_chr19.macsinfo

Use the arrow keys to go to the top, where you'll see all of the parameters you put in to run MACs. After some runtime info (including possible warnings, that you can ignore if there are not millions of them), you'll see:

INFO @ Sun, 10 Feb 2013 21:27:51: #1 total tags in treatment: 370513

INFO @ Sun, 10 Feb 2013 21:27:51: #1 user defined the maximum tags...

INFO @ Sun, 10 Feb 2013 21:27:51: #1 filter out redundant tags at the same location and the same strand by allowing at most 1 tag(s)

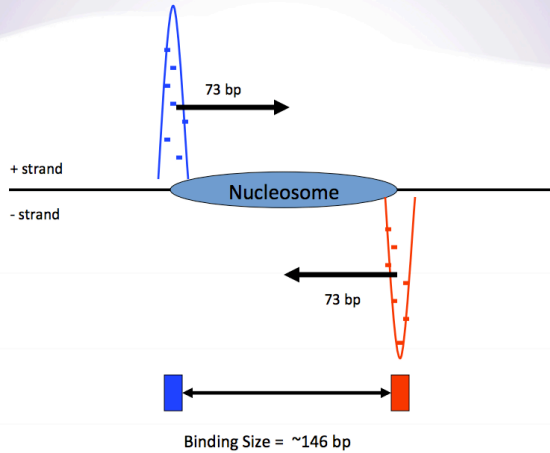
INFO @ Sun, 10 Feb 2013 21:27:51: #1 tags after filtering in treatment: 275955 (High is good)

INFO @ Sun, 10 Feb 2013 21:27:51: #1 Redundant rate of treatment: 0.26 (low is good)

This is useful information. It tells you how many different reads you had (out of all of the reads which mapped to only one place in the mouse genome- from Bowtie). You want this number to be high and the “redundant rate” to be low.



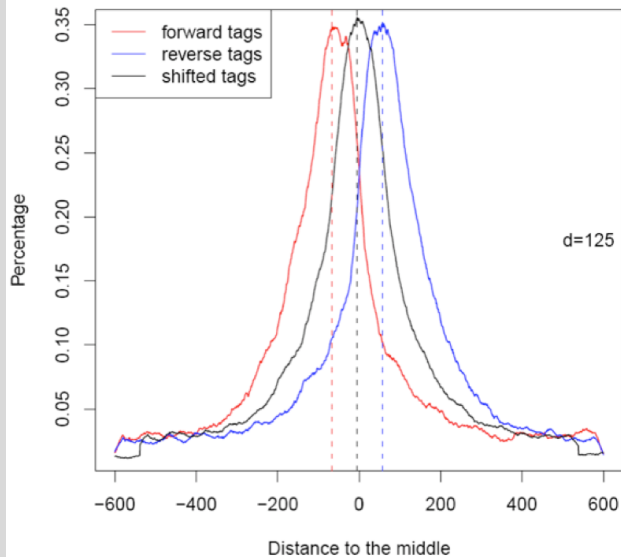
Finding the binding sizes



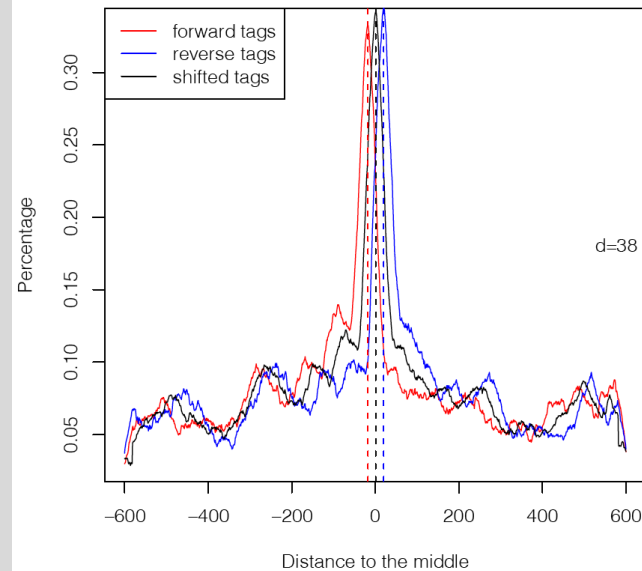
Be on the lookout for MACS building a model from short-separation noise peaks (that may arise from sonication sensitive breakpoints or other things unrelated to your protein binding).

To avoid this, you can decrease the maximum “mfold” so that these strong irrelevant peaks are ignored when the model is built.

Peak Model



Peak Model



Peaks & negative peaks

Keep scrolling down your .macsinfo file until you find...

...

INFO @ Sun, 10 Feb 2013 21:36:47: #3 **Finally, 364 peaks are called!**

INFO @ Sun, 10 Feb 2013 21:36:47: #3 find negative peaks by swapping treat and control

INFO @ Sun, 10 Feb 2013 21:36:52: #3 **Finally, 36 peaks are called!**

INFO @ Sun, 10 Feb 2013 21:36:52: #4 Write output...

This is the pay-off, where MACS identifies your ER alpha peak locations!

364 peaks on chromosome 19 (which is $\sim 1/50^{\text{th}}$ of the genome) suggests $\sim 20,000$ peaks for the whole genome, which is not bad!

Equally critical, MACS now swaps treat & control (pretending your INPUT data is your IP & your ChIP data is your input) and looks again for peaks.

The number of “negative” peaks found in this way should be far less than the positive peaks, and the 10:1 ratio here is fine.

#####

Peaks/Negative Peaks ratio is poor or too few peaks are detected:

- Adjust model settings to see if you can improve both. Otherwise, you may have to conclude that 1) your library was no good or 2) the factor just doesn't bind to many places in the genome.



MACS2 output : .bdg files

Broad IGV, an alternative to UCSC browser

IGV layout

The image shows a screenshot of the Integrative Genomics Viewer (IGV) interface. The interface is divided into several sections: a top menu bar with 'File', 'View', 'Tracks', and 'Help'; a search bar containing 'Human hg18', 'chr1', and 'hotair'; a cytoband track showing chromosome bands; a genomic coordinate track with a scale from 38 mb to 46 mb; and a track list on the left with entries like 'GM12878 H3K4me1' and 'GM12878 H3K27me3'. Red arrows point from yellow labels to these specific features: 'Genome' points to the search bar, 'Chromosome' points to 'chr1', 'Cytoband' points to the cytoband track, 'Track Names' points to the track list, and 'Genomic Coordinates' points to the coordinate scale. The IGV logo and 'Integrative Genomics Viewer' text are in the top right corner.

<http://www.broadinstitute.org/igv/>

You will need to register, but they don't send you spam.



GTRD: Gene Transcription Regulation Database

<http://gtrd.biouml.org/>

GTRD

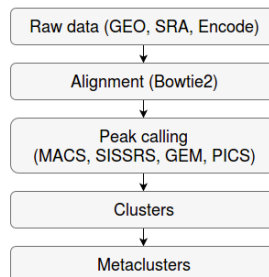
Gene Transcription Regulation Database

The most complete collection of uniformly processed ChIP-seq data to identify transcription factor binding sites for human and mouse. Convenient web interface with advanced search, browsing and genome browser based on the BioUML platform. For support or any questions contact ivan@dote.ru

[Start »](#) [Documentation »](#) [Download »](#)

Workflow

How it was constructed?



ChIP-seq experiment information and raw data were collected from publically available sources. Sequenced reads were aligned using Bowtie2 and ChIP-seq peaks were called using 4 different methods. Peaks were merged into clusters and metaclusters to produce non-redundant set of transcription factor binding sites.

Statistics

version 16.07

ChIP-seq experiments	5078
Transcription factors	542
ChIP-seq reads	183 770 420 298
Reads aligned	146 855 318 517
ChIP-seq peaks	296 112 068
Clusters	176 007 357
Metaclusters	30 178 267

[Learn more »](#)



Table 1. Comparison of databases that are based on ChIP-seq data

Database, URL	Source of human and mouse data	Number of samples (TF-related)*	Number of TFs	Number of ChIP-seq peak callers used	Metaclus-ter approach	Uniform data processing	Genome browser
ChIPBase (http://rna.sysu.edu.cn/chipbase)	GEO, ENCODE	total 3549 human 2498 mouse 1036 rat 15	252 TFs and non-TFs for 10 species	>10 in total, but no uniform pipeline, each ChIP-seq is processed by different peak caller	No	No	Self-developed: deepView genomeView
Cistrome DB (http://dc2.cistrome.org/#/)	GEO, SRA, ENA, ENCODE	total 10 276 (TF+non-TF) human 5774 mouse 4502 rat 0	260 TFs and non-TFs	1 (MACS2)	No	Yes	UCSC genome browser
ENCODE (https://www.encodeproject.org)	ENCODE	total 1448 human 1254 mouse 194 rat 0	295 TFs and non-TFs for human, 52 TFs and non-TFs for mouse	5 (SPP, GEM, PeakSeq, MACS, Hotspot/Hotspot2)	No	Yes	Self-developed: UCSC genome browser and WashU epigenome browser
Factorbook (http://www.factorbook.org)	ENCODE	total 1007 human 837 mouse 170 rat 0	167 TFs, co-factors and chromatin remodeling factors for human, 51—for mouse	None	No	No	No
GTRD (http://gtrd.biouml.org)	GEO, SRA, ENCODE	total 5078 human 2955 mouse 2107 rat 16	476 human and 257 mouse sequence specific TFs, corresponding to 542 TFClass classes.	4 (MACS, SISSRs, GEM, PICS)	Yes	Yes	Self-developed
ChIP-Atlas (http://chip-atlas.org)	SRA	total 10 774 human 5914 mouse 4860 rat 0	699 human and 502 mouse TFs and others.	1(MACS2)	No	Yes	IGV
GeneProf (http://www.geneprof.org)	SRA, ENCODE, literature	total 1692 human 693 mouse 999 rat 0	133 human and 131 mouse TFs	1(MACS)	No	Yes	Self-developed: based on GenomeGraphs
NGS-QC (http://www.ngs-qc.org)	GEO	total 6672 human 4234 mouse 2438 rat 0	unknown	None	No	Yes	No



GALAXY/CISTROME

<http://cistrome.dfci.harvard.edu/ap/root>

Galaxy tools specially designed for ChIP-seq analysis.

Cistrome allows easy access to many analyses that give you some quick insights into your data.

CISTROME TOOLBOX: ASSOCIATION STUDY

GCA: Gene centered annotation Find the nearest interval in the given intervals set for every annotated coding gene (e.g. where's the nearest ER binding site for each gene in the genome)

peak2gene: Peak Center Annotation Input a peak file, and It will search each peak on UCSC GeneTable to get the refGenes near the peak center (e.g. where's the nearest gene to each ERa binding site in the genome).

