Epigenetics and ChIP-seq



CSAMA 2015, Brixen 16. 06. 2015. Aleksandra Pekowska aleksandra.pekowska@embl.de



Outline of the lecture

Purpose: introduce basic steps and key considerations in ChIP-seq analysis

- 1. Epigenetics fundamental concepts
- 2. The ChIP-seq method
- 3. What kind of information can we obtain from ChIP-seq?
- 4. Study design
- 5. ChIP-seq analysis workflow:
 - a. Preprocessing
 - b. Quality controls
 - c. Isolation of enriched regions
 - d. Analysis of enriched regions
 - e. Visualization
 - f. Average profiles
 - g. Comparative analysis of enriched regions

Epigenetics - inheritance, but not as we know it

Non-genic memory of function transmitted from generation to generation (A. Bird)



Adapted from Conrad Hal Waddington (1942)

Factors which are analysed:

- DNA methylation
- nucleosome occupancy
- histone modifications
- transcription factors
- RNA-polymerases
- chromatin modifying enzymes

Chromatin Immunoprecipitation



What kind of information can we obtain from the ChIP-seq experiments ?

Resource Е Nucleosomes Pol II H3K4me1 Promoter H3K4me2 H3K4me3 **High-Resolution Profiling of Histone** H2A.Z H3K27me3 H3K36me3 **Methylations in the Human Genome** Enhancer H2BK5me1 or Insulator H3K9me1 H4K20me1 Artem Barski,^{1,3} Suresh Cuddapah,^{1,3} Kairong Cui,^{1,3} Tae-Young Roh,^{1,3} Dustin E. Schones,^{1,3} Zhibin Wang,¹ Gang Wei,^{1,3} louri Chepelev,² and Keji Zhao^{1,7} ¹Laboratory of Molecular Immunology, National Heart, Lung, and Blood Institute, NIH, Bethesda, MD 20892, USA ² Department of Human Genetics, Gonda Neuroscience and Genetics Research Center, University of California, Los Angeles, Los Angeles, CA 90095, USA ³These authors contributed equally to this work and are listed alphabetically. *Correspondence: zhaok@nhlbi.nih.gov Active Genes Pol II H3K4me3 2.5e-07 2e-07 Normalized Counts high B 1.8e-07 medium 2e-07 1.6e-07 low silent 1.4e-07 1.2e-07 1.5e-07 1e-07 1e-07 8e-08 6e-08 5e-08 4e-08 2e-08 n 0 -2000 2000 4000 -2000 -1000 2000 0 -4000 0 1000

What kind of information can we obtain from the **ChIP-seq experiments ?**



Distance from Smc1a occupied region (kb)

Kagey 2010

What kind of information can we obtain from the ChIP-seq experiments ?





Chen 2008

To summarize - the most frequent tasks are:

- 1. Visualization along the genome
- 2. Peak finding and analysis (localization, cooccurrences, motifs)
- 3. Heatmaps of signal and average profiles at various genomic *loci*

But before we start the analysis... ChIP-seq: considerations for study design

- Distribution of modification number of sequenced reads
- Paired vs. single end sequencing fragment length estimation
- IgG control (pros and cons)
- Input control
- Biological replication!



ChIP-seq: sequencing depth matters



Landt 2012

ENCODE consortium guidelines

For mammalian genomes such as human and mouse:

20M aligned reads for broad marks
2. > 10M aligned reads for TFs

Paired vs. single end sequencing

- paired end sequencing is always useful (nucleosome positioning) however not absolutely necessary



The estimation of the length of the ChIP fragments



- •Binning visualization and signal distribution analysis
- •Quality control check
- Peak finding

Fragment length estimation - quality controls



Figure 4. (Legend on next page)

ChIP-seq: considerations for study design

- IgG control (pros and cons)
- Input control
- Biological replication

Finding enriched regions



150 -

Enriched regions ('peaks') regions with signal which is significantly higher than the background - input or IgG

<u>Input reads</u> - background reads' distribution exhibits a degree of clustering that is significantly greater than expected from a homogenous Poisson process (*P*-value< 10⁻⁶, Kharchenko et al., 2008)



Model-based analysis of ChIP-seq (MACS)

Method

Model-based Analysis of ChIP-Seq (MACS)

Yong Zhang^{**}, Tao Liu^{**}, Clifford A Meyer^{*}, Jérôme Eeckhoute[†], David S Johnson[‡], Bradley E Bernstein^{§¶}, Chad Nusbaum[¶], Richard M Myers[¥], Myles Brown[†], Wei Li[#] and X Shirley Liu^{*}

- removes PCR duplicates

d is estimated by picking highly enriched regions and looking at the distance between modes of positive and negative strand read pileups. Reads are extended towards this midpoint (building peak model)

- Sliding window of 2*d* to find significantly enriched bins using λ_{local} . We obtain enrichment P-value

- eFDR by swapping control and treatment



Several examples of peak callers

SICER - designed to deal with histone type data PeakSeq, chromHMM ...



MOSAiCS - suitable for TF and histone modification data

BayesPeak - suitable for TFs and histone modifications displaying peaklike signal

<u>ChIPseqR</u> - suitable for nucleosome positioning analysis

PICS CSAR NarrowPeaks CSSP

Peak processing - quality controls

- how do we decide whether samples and peaks are OK?





Visualization - seeing is believing



Visualization - other tools

IGB - Integrated Genome Browser http://bioviz.org/igb/index.html



IGV - Integrative Genomics Viewer https://www.broadinstitute.org/igv/

Visualization - file formats





Peak analysis

Frequently asked questions include:

- Localization of peaks with respect to functional elements in the genome (promoters, gene body, introns, transcription termination sites, intergenic regions etc.)

- Co-ocurrence between enriched regions
- The distribution of signal at the peaks

ChIPpeakAnno - provides functions performing peak annotation to promoters etc.

biomaRt - easy access to data bases including gene annotation, sequence conservation, sequence retrieval etc.

<u>GenomicRanges</u> - fast comparison between genomic intervals:

findOverlaps()

countOverlaps()

nearest()

Easy peak annotation to pre-established or new genomic features, cross-comparisons between peak locations and any kind of imaginable analysis

VennDiagram - visualization of two or multi-sample overlaps

Rcade - integrates ChIP-seq analysis with differential expression

Peak analysis - GREAT tool

GREAT Input: Genomic Re × +										
,GRE ≜ T,	Overview	News	Use GREAT	Demo	Video	How to Cite	Help	Forum		
GREAT version 3.0.0 current (02/15/2015 to now)						•				

GREAT predicts functions of cis-regulatory regions.

Many coding genes are well annotated with their biological functions. Non-coding regions typically lack such annotation. GREAT assigns biological meaning to a set of non-coding genomic regions by analyzing the annotations of the nearby genes. Thus, it is particularly useful in studying cis functions of sets of non-coding genomic regions. Cis-regulatory regions can be identified via both experimental methods (e.g. ChIP-seq) and by computational methods (e.g. comparative genomics). For more see our Nature Biotech Paper.

News

- Feb 15, 2015: GREAT version 3.0 switches to Ensembl genes, adds the mouse mm10 assembly, and adds new ontologies.
- Apr 3, 2012: GREAT version 2.0 adds new annotations to human and mouse ontologies and visualization tools for data exploration.
- · Feb 18, 2012: The GREAT forums are released, allowing increased user-to-user interaction

More news items...

Species Assembly

- Human: GRCh37 (UCSC hg19, Feb/2009)
- Mouse: NCBI build 37 (UCSC mm9, Jul/2007)
- Mouse: NCBI build 38 (UCSC mm10, Dec/2011)
- Zebrafish: Wellcome Trust Zv9 (danRer7, Jul/2010) Zebrafish CNE set

Can I use a different species or assembly?

Test regions

BED file: Browse...
No file selected.

BED data:



Peak analysis - motifs

<u>MEME</u> - provides functions performing motif discovery

RSAT - complete suite for motif finding

Position Weight Matrix (PWM) - describes the probability of each nucleotide at each position of a motif

JASPAR/TRANSFAC - data bases of PWM

R: MotifDb, FIMO and others





Co-enrichment and signal distribution analysis





Region divided in to tiles

Count how many fragments fall into each tile



0

Visualization



- Heatmaps of signal enrichment at - promoters
- loci enriched with factors of interest

We will see an example of such an analysis using R package *GenomicRanges*

A nice alternative: *HT-Seq* (python)

Comparative peak analysis



Threshold issues affecting all qualitative analyses

Comparative peak analysis

<u>DiffBind</u>

- 1. Count reads in peaks in all the replicates and conditions
- 2. Perform edgeR or DESeq2 analysis dba.analyze()
- 3. Provides various plotting functions

<u>MMDiff</u>

- 1. Count reads in peaks in all the replicates and conditions
- 2. Performs *DESeq* normalisation
- 3. Compares peak shapes using kernel based statistical tests

<u>ChIPQC</u> package for quality control checks and quantitative analysis of peak strengths

- 1. Plotting coverage histograms for peaks
- 2. Cross-coverage analysis in the function of shift sizes
- 3. Plotting peak profiles
- 4. Sample clustering



References (I)

- Anders S, Pyl PT, Huber W. 2015. HTSeq—a Python framework to work with highthroughput sequencing data. Bioinformatics 31 (2): 166-169
- Bailey T, Krajewski P, Ladunga I, Lefebvre C, Li Q, Liu T, Madrigal P, Taslim C, Zhang J. 2013. Practical Guidelines for the Comprehensive Analysis of ChIP-seq Data. PLoS Comput Biol 9: 5–12.
- Barski A, Cuddapah S, Cui K, Roh T, Schones DE, Wang Z, Wei G, Chepelev I, Zhao K. 2007. High-Resolution Profiling of Histone Methylations in the Human Genome. Cell 129: 823–837.
- Chen X, Xu H, Yuan P, Fang F, Huss M, Vega VB, Wong E, Orlov YL, Zhang W, Jiang J, et al. 2008. Integration of External Signaling Pathways with the Core Transcriptional Network in Embryonic Stem Cells. Cell 133: 1106–1117.
- Creyghton MP, Cheng AW, Welstead GG, Kooistra T, Carey BW, Steine EJ, Hanna J, Lodato MA, Frampton GM, Sharp PA, et al. 2010. Histone H3K27ac separates active from poised enhancers and predicts developmental state. PNAS 107.
- Feng J, Liu T, Qin B, Zhang Y, Liu XS. 2012. Identifying ChIP-seq enrichment using MACS. Nat Protoc 7: 1728–1740.
- Jothi R, Cuddapah S, Barski A, Cui K, Zhao K. 2008. Genome-wide identification of in vivo protein-DNA binding sites from ChIP-Seq data. Nucleic Acids Res 36: 5221–5231.

References (II)

- Kagey MH, Newman JJ, Bilodeau S, Zhan Y, Orlando DA, van Berkum NL, Ebmeier CC, Goossens J, Rahl PB, Levine SS, et al. 2010. Mediator and cohesin connect gene expression and chromatin architecture. Nature 467: 430–435.
- Kharchenko P V, Tolstorukov MY, Park PJ. 2008. Design and analysis of ChIP-seq experiments for DNA-binding proteins. Nat Biotechnol 26: 1351–1359.
- Landt S, Marinov G. 2012. ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. Genome Research 1813–1831.
- Li Q, Brown B, Huang H, Bickel P. 2010. IDR analysis 101 Measuring consistency between replicates in high-throughput experiments. 1–7.
- Lun ATL, Smyth GK. 2014. De novo detection of differentially bound regions for ChIPseq data using peaks and windows: Controlling error rates correctly. Nucleic Acids Res 42: 1–11.
- Thomas-Chollier M, Herrmann C, Defrance M, Sand O, Thieffry D, Van Helden J. 2012. RSAT peak-motifs: Motif analysis in full-size ChIP-seq datasets. Nucleic Acids Res 40.
- Zhang Y, Liu T, Meyer C a, Eeckhoute J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM, Brown M, Li W, et al. 2008. Model-based analysis of ChIP-Seq (MACS). Genome Biol 9: R137.