

**Ciências ULisboa**

Abordagens Ómicas em Biomedicina e Biotecnologia

**ChIP-seq 101**

Paulo Matos  
2020/2021

**Regulation of Gene Expression...**

• One genome → multiple cell types!

*Precise regulation of gene expression is determinant for the stimulus-cued and temporally and spatially coordinated differentiation and function of different cell-types...*

https://www.cryo-cells.com/CryoCells/media/CryoCells/Class%20imagery/stem-cell-differentiation.png

**Regulation of Gene Expression...**

Epithelial cell monolayer

**Regulation of Gene Expression...**

• Different stimuli are temporally and spatially integrated through complex signal transduction pathways and relayed to the nucleus...

**Transcriptional regulation in eukaryotes**

**Levels of transcriptional regulation:**

- 3D structure of chromatin in the nucleus - Chromatin compaction by histone proteins and their post-translational modifications;
- Reversible modifications of DNA (e.g., methylation of promoters);
- Transcription factor activation and binding to DNA...

Balaban 2003, ESHN, Physiol 4: 238

**Transcription factors in eukaryotes**

**Transcription factors (TFs) are proteins that bind specific DNA sequences to control the transcription of genes.**

- General TFs bind to core promoter regions.
- Site-specific TFs bind at the proximal promoter or distal enhancer or repressor regions.
- TFs often work in protein complexes, integrating the outputs of several signaling pathways...

Frederick, 2005, Nature Reviews Genetics, 6(10): 519-36

**Identification of TF binding sites (TFBS)**

To understand how different conditions or signals influence gene expression and how particular regulatory molecules or events are involved in the observed responses, we need to know:

- All of the sites in the genome to which the transcription regulators of interest bind, under those conditions.
- Possibly correlate this information with chromatin status, and transcript expression levels.

Wang et al., 2005, Nature 435: 826-30

**ChIP-seq**

Chromatin Immunoprecipitation-next generation Sequencing

**ChIP-seq overview**

**ChIP-Seq principle**

Find all regions in the genome bound by a specific transcription factor (or bearing a specific histone modification) in a certain experimental setting (a particular cell type, a specific stimulus, a drug treatment, etc...).

Fu, 2008, Nature Reviews Genetics, 9(12): 845-56



### Mapping reads to the genome

Bowtie2 is a commonly used mapping tool for CHIP-Seq data

Langmead B, Salzberg S. Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 2012, 9:375-380.

**Bowtie output files are in the SAM (Sequence Alignment/Map) format**

- A generic format for storing large nucleotide sequence alignments.
- The related **BAM format** is a compressed binary representation of a SAM files.

Field	Description
CONNAME	Query template (read/NAME)
FLAG	A bit value summarizing general alignment control parameters
RNAME	Reference sequence NAME
POS	1-based leftmost mapping POSITION
MAPQ	Mapping quality
CIGAR	Match/mismatch/insertion/deletion
REFLEN	Ref. name of the match/mismatch
PRECISE	Position of the match/mismatch
TLEN	observed template length
SEQ	sequence
QUAL	ASCII of Phred-scaled base QUALity+33

### Mapping reads to the genome

Mapped CHIP-seq data can be visualized as tracks in genome browser tools (e.g., the Integrated Genome Browser (IGV) or the UCSC Genome Browser) together with other annotations (e.g., genes in the vicinity)...

Reads aligned to the + strand

Reads aligned to the - strand

Mapped reads cluster at enriched regions - "peaks"

### Peak Calling

Strand asymmetry in CHIP-seq read densities is actually expected (peak recognition pattern)...

TF "protected" site

~200 bp fragments flanking the TF binding site

~40 bp reads from sequence from one end of the fragments (+ or - strand)

read densities on +/- strand

TF

CHIP-seq on DNA binding TF

+/- strand read "peaks" do not represent the true location of the binding site

### Peak Calling

- ✓ The input signal for peak calling is the number of mapped reads at each genomic position.
- ✓ Reads mapping to the forward and reverse strand have to be combined to increase the signal and peak resolution.

- **Tag shifting:** reads are shifted by  $d/2$  (i.e. towards the middle of the IP fragment), where  $d$  represent the estimated mean fragment length.
- **Tag elongation:** Each tag is computationally extended in 3' to a total length of  $d$ .

### Peak Calling

- ✓ Strand asymmetry is blurred by the binding of several TFs as complexes...

CHIP-seq on DNA binding TF

read densities on +/- strand

CHIP-seq on histone modifications

read densities on +/- strand

...and is virtually lost in datasets of CHIP for histone modifications.

### For example...

TF peaks

Histone peaks

TF reads

Histone reads

### CHIP = target enrichment...

Potential sources of "noise" in CHIP-seq experiments:

- DNA-shearing is not uniform across genome, which results in more reads in open chromatin regions.
- PCR amplification bias (GC content) in library preparation.
- Repetitive regions might appear enriched due to underestimated repeat copies in the reference genome.

### How do we estimate the noise?

www.VADL0.com

"Data don't make any sense, we will have to resort to statistics."

### Poisson model for read counts

If we assume that the "noise" reads within a given genomic window have a random distribution then we can use a theoretical distribution to estimate the probability of observed number of reads being above what is expected (i.e., the "noise").

The Poisson distribution is considered adequate since it models the number of random events (reads) in a fixed interval (genomic window) (Note - some authors prefer the binomial or the negative binomial):

$$P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

$P$  = the probability that  $k$  reads map to a given genomic region.

$\lambda$  = mean number of reads in a genomic window of a certain size.

Certain algorithms use a uniform  $\lambda$  value estimated from annotated whole genome data. But...

### Modelling noise levels

Noise is **not uniform** (chromatin conformation, technical biases, repetitive regions...etc)

**Solution:** sequence **Input DNA** to use as reference (control)

"Input DNA" is isolated from cells that have been cross-linked and fragmented under the same conditions as the immunoprecipitated DNA.

### Defining "significant peaks"

Slide a window across the genome. At each location:

- Evaluate the enrichment of the ChIP signal
- Calculate  $\lambda$  from input and  $p$  based on the poisson distribution
- retain regions with significant  $p$ -values

### MACS (Model-based Analysis of ChIP-Seq)

Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

<https://macs3-project.github.io/MACS/>

MACS is the pipeline recommended by the **ENCODE** project for ChIP-seq data analysis

ENCODE: The Encyclopedia of DNA Elements is a project launched by the NIH's Human Genome Research Institute in 2003 that aims to identify all functional elements in the human genome sequence.

### MACS (Model-based Analysis of ChIP-Seq)

Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

#### 1: Modeling the shift size of ChIP-Seq tags (estimating $d$ )

- Slides a window of  $2x$  the user defined sonication size (BANDWIDTH)
- Detects regions with high read enrichment in ChIP vs. input (MFOLD - default 10-30 fold)
- Plots average +/- strand read densities  $\rightarrow$  estimates  $d'$  (i.e., mean fragment size)
- Shifts read tags  $d'/2$  to increase ChIP signal resolution...

### MACS (Model-based Analysis of ChIP-Seq)

Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

#### 2: "Significant Peak" detection

##### I. Normalize datasets:

- Datasets generally do not have the same sequencing depth
- ChIP (=signal + noise) and input (=noise)

Input =  $N$  reads  
ChIP-seq dataset =  $M > N$  reads  
ChIP scaled by library size:  
 $M' = M * (N1/N2)$

### MACS (Model-based Analysis of ChIP-Seq)

Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

#### 2: "Significant Peak" detection

##### II. Identify local noise parameter

- slide a window of size  $2*d'$  across treatment and input
- estimate parameter  $\lambda_{local}$  for Poisson distribution

$P\text{-value} = 1E-30$

### MACS (Model-based Analysis of ChIP-Seq)

Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

#### 3: "Significant Peak" detection

##### III. Calling peaks and scoring

- Candidate peaks with  $p$ -values below a user-defined threshold (PVALUE - default 1E-05) are called (Peak regions).
- The ratio between the ChIP-Seq tag count and  $\lambda_{local}$  is reported as **fold\_enrichment**.
- The location with the **highest fragment pileup (summit)** is predicted as the TF binding location.

### MACS (Model-based Analysis of ChIP-Seq)

Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

#### 3: Estimation of False Discovery Rate (FDR)

- The FDR is an estimation of the proportion of incorrectly identified peaks among those that are found to be of a certain significance ( $p$ -value).

#### How MACS does it?

- Swaps "ChIP" and "input" datasets and calculates  $p$ -values for "negative peaks".
- Ranks  $p$ -values from both **Positive** (ChIP vs input) and **Negative** (input vs. ChIP) peaks.
- Calculates FDR for a peak with a given  $p$  as:

$FDR(p) = \frac{\# \text{negative peaks with } Pval < p}{\# \text{positive peaks with } Pval < p}$

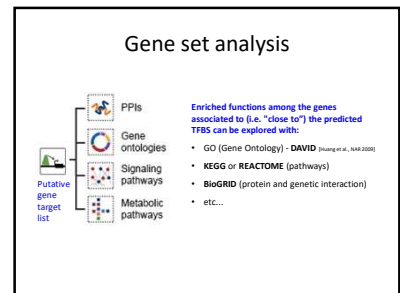
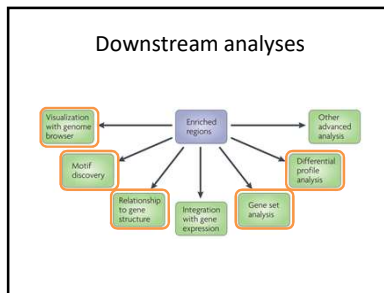
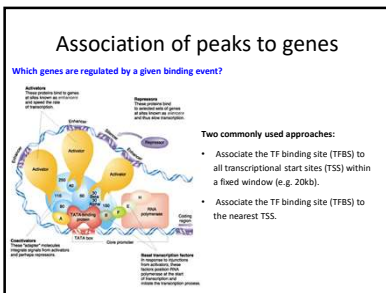
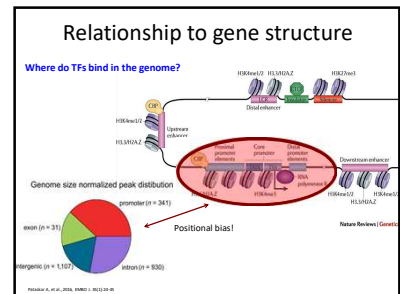
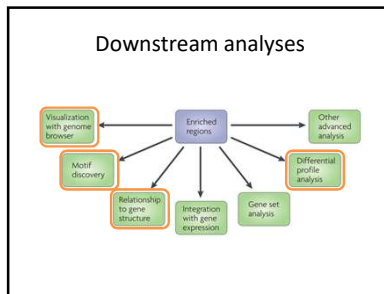
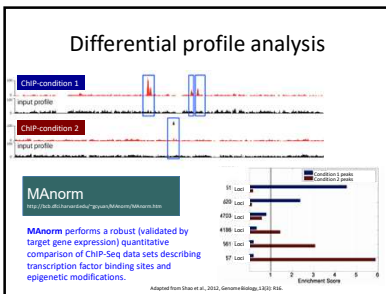
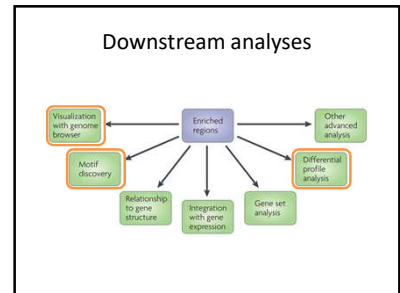
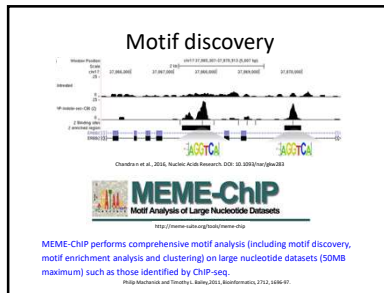
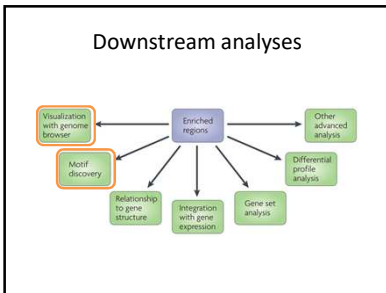
### MACS (Model-based Analysis of ChIP-Seq)

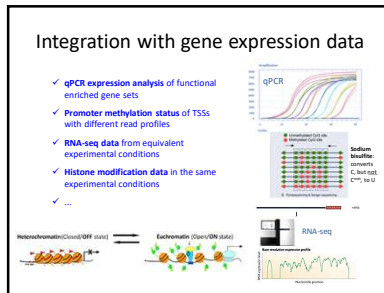
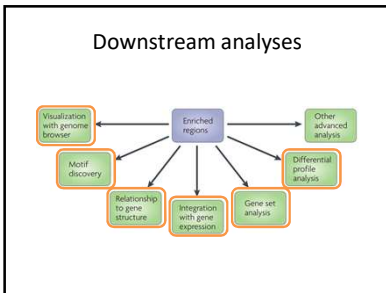
Zhang et al. Model-based Analysis of ChIP-Seq (MACS) Genome Biol (2008) vol. 9 (3) pp. R137

#### 4: Output a list of peaks coordinates in BED format

- BED** (Browser Extensible Data) is a UCSC Genome Browser format that provides a flexible way to define the data lines that are displayed in an annotation track.

chr	start	end	length	summit	tags	$-\log_{10}(p\text{-value})$	fold_enrichment	FDR (%)
chr1	8595	8611	422	286	58	146.98	11.88	0.14
chr1	94382	36984	605	405	99	315.23	27.05	0.24
chr1	70070	70070	309	272	14	85.40	5.74	1.19
chr1	289684	339029	362	118	13	62.50	7.17	1.19
chr1	61083	61648	506	279	61	234.60	20.85	0.19
chr1	68009	68068	366	66	90	348.81	20.10	0.19
chr1	60227	60268	405	245	77	349.87	18.81	0.20
chr1	51246	51329	720	330	66	479.24	36.90	0.17
chr1	51204	51418	846	130	60	389.72	27.98	0.19





### Advantages and disadvantages of ChIP-seq

**Advantages:**

- Genome-wide binding sites.
- High resolution (~200 bp)
- Relatively cheap and scalable (ENCORE project performed hundreds of ChIP-seq experiments for many TFs and histone marks in many tissues, conditions and organisms).

**Disadvantages:**

- Requires large cell populations (millions) (although single-cell ChIP-seq approaches are being developed)
- Depends on highly specific antibodies
- Can not distinguish between direct and indirect TF binding
- Finding enriched binding regions (peaks) relies heavily on statistical methods - bias towards genomic amplifications (cancer cells), repetitive regions, and **high affinity binding sites...**

### Suggested reading

- ✓ Nakatani, R. and Sakatani, T. (2020) Methods for ChIP-seq analysis: A practical workflow and advanced applications. *Methods*, in press; DOI: 10.1016/j.ymeth.2020.03.005
- ✓ Pavani, G. (2017) ChIP-Seq Data Analysis to Define Transcriptional Regulatory Networks. *Adv Biochem Eng Biotechnol.*, 160:1-14.
- ✓ Furey, T. S. (2012). ChIP-seq and beyond: new and improved methodologies to detect and characterize protein-DNA interactions. *Nature reviews. Genetics*, 13:840-52.