# Genomic enrichment strategies: ChIP-seq and ATAC-seq Analysis of Next-Generation Sequencing Data

## Friederike Dündar

Applied Bioinformatics Core

Slides at https://bit.ly/2Kn0QHt<sup>1</sup>

April 9, 2019

### Weill Cornell Medicine

<sup>1</sup>https://physiology.med.cornell.edu/faculty/skrabanek/lab/angsd/schedule\_2018/

F. Dündar (ABC, WCM)

Genomic enrichment strategies: ChIP-seq and

April 9, 2019 1 / 76



- 2 Studying Chromatin
- 3 ATAC-seq principles
- 4 Processing ATAC-seq data
- 5 ChIP-seq principles
- 6 Processing of ChIP-seq data

## References

## From DNA to phenotype





## Epigenetics

### Waddington's definition of epigenetics

Epigenetics encompasses the molecular mechanisms by which the genes of the genotype bring about phenotypic changes [Waddington, 1942].



# Epigenetics: understanding how the genetic code is interpreted



F. Dündar (ABC, WCM)

Genomic enrichment strategies: ChIP-seq and

## DNA does not occur naked in eukaryotic cells



## Chromatin = DNA + proteins + ncRNA



### The most obvious function of chromatin is **DNA compaction**.

F. Dündar (ABC, WCM)

Genomic enrichment strategies: ChIP-seq and

## **DNA** compaction



## DNA compaction

## Example for relatively trivial compaction: 375 m (~1230 ft) of yarn packed into a ball of about 10 cm × 4 cm (4"×1.6") using **simple coils**





Studying Chromatin

## Studying Chromatin

## From DNA to phenotype: epigenetics

The current assumption is that the **chromatin structure** is an essential part of defining an individual cell's fate, i.e. by interacting tightly with DNA and regulating access to it, chromatin has a key role in how transcription is achieved in a highly time- and tissue-dependent manner.



"Understanding the chromatin structure can give a perspective of how a certain mRNA expression state was reached and how a cell might advance." [Winter et al., 2015]

# Chromatin signatures of regulatory regions

- ${\ensuremath{\, \bullet }}$  Trans-regulatory elements = DNA encoding transcription factors
  - the actual effectors are proteins
- **Cis-regulatory elements (CRE)** = non-protein-coding DNA that regulates transcription of neighboring genes
  - ▶ the effectors are thought to be (at least partially) the DNA sequences



# NGS-based features of regulatory regions

Using NGS, we've catalogued distinct features of different CRE types.



## • transcription start site (**TSS**)

- lots of Pol II & associated machinery
- H3K4me3, H3K27ac and more

## enhancers

- 100 bp to 1,000 bp
- enriched for H3K4me1 & p300
- bound by TF
- weak Pol II activity

### insulators

- enhancer blockers or barriers preventing chromatin condensation
- 300 bp to 2000 bp
- characterized by CTCF binding and intra- and inter-chromosomal interactions
- repressed chromatin
  - H3K27me3 & DNA

Studying Chromatin

## Chromatin states are cell-type-specific





Ernst & Kellis (2017) doi:10.1038/nprot.2017.124

Cell-type-specific chromatin states (including active cisreg. elements) are often defined by integrating genome-wide profiles of histone marks and transcription factor binding. Studying Chromatin

## Chromatin states are cell-type-specific



Different chromatin states are also characterized by different nucleosome occupancies.

## 2 basic chromatin states based on nucleosome occupancy

For transcription to occur, the RNA Pol II machinery needs to access the **naked** DNA strand, i.e. the chromatin needs to be made **accessible locally**.



#### Studying Chromatin

# NGS techniques for studying chromatin and DNA modifications



The majority of epigenomics data entails profiles of **nucleosome occupancy**, specific **histone marks** and **transcription factor** binding.

These information are all inferred based on which DNA sequences we find over-represented in our data set.

### Basic concept



Basic concept



### Basic concept





# NGS approaches for epigenomics



- DNA = more or less immutable code
- RNA = the code's local read-out
- "epigenome" = additional molecules or chemical DNA modifications that govern the process of DNA-to-RNA transcription
- technically, epigenetics only refers to *heritable* marks that influence transcription [Ptashne, 2013]
- in practice, epigenomics is often used to describe all kinds of aspects of transcription regulation, including highly dynamic ones!

ATAC-seq principles

## ATAC-seq principles

## Identifying accessible chromatin regions

Active CRE (promoters, gene bodies, enhancers, TFBS) are expected to be accessible.



# Assay for transposase-accessible chromatin (ATAC)



## ATAC-seq profiles



29 / 76

# Interpretation of ATAC-seq data



# ATAC-seq profiles are typically population snapshots



# ATAC-seq profiles are typically population snapshots, but scATAC-seq is possible



Processing ATAC-seq data

## Processing ATAC-seq data

Processing ATAC-seq data



# Established ATAC-seq pipelines

## • ENCODE

- lots of QC scores and guidelines for identfying samples that worked/failed
- somewhat cumbersome implementation
- Tom Carroll's R-based workflow
  - mostly follows ENCODE's guidelines
  - every command is shown including some explanations about important parameters
  - R is not the best-suited environment for some of the steps (e.g. bigWig generation)

### Harvard FAS

- some steps of the ENCODE pipeline are re-worked/re-thought
- alternative peak caller (not yet peer-reviewed, but more versatile/ATAC-seq-oriented than MACS2)

## Raw data processing: FASTQ to BAM

- **FastQC** the usual suspects: sequencing quality, duplications, contaminations
- adapter removal may be warranted
  - PE sequencing will often lead to frequent adapter sequences for ATAC-seq data because many *fragments* are shorter than 2x50bp

DNA fragment > 2x read length



DNA fragment < 2x read length



• genome aligners for short reads, e.g. Bowtie2 or BWA
# Raw data QC: filtering the BAM files

The following reads are removed:

- mitochondrial reads
- discordantly "paired" reads
- non-uniquely aligned reads
- PCR duplicates
- reads corresponding to fragments < 40 bp (see slides about fragment size distributions)
- reads overlapping with blacklisted regions





Processing ATAC-seq data

# PCR duplicates are frequent – more so for low cell numbers!



# The dominant fragment size distribution signal in ATAC-seq should reflect the nucleosome pattern



#### ATAC-seq data contains multiple levels of information



Processing ATAC-seq data

### But MNase usually beats ATAC-seq in terms of resolution



#### Examples of ATAC-seq frag. size distributions



Ou et al (2018). doi: 10.1186/s12864-018-4559-3

- typical problems seen here:
  - overdigestion/too much Tn5
  - too little Tn5/incomplete digestion
  - flawed size selection



#### Blacklisted regions: regions with spurious signals

- typically appear uniquely mappable
- often found at specific types of repeats such as centromeres, telomeres and satellite repeats
- especially important to remove these regions *before* computing measures of similarity

**Blacklists** were generated empirically by the (mod)ENCODE consortium: http://mitra.stanford.edu/kundaje/akundaje/release/blacklists/



bedtools intersect -abam reads.bam -b blacklisted.bed > filtered\_reads.bam

### Checking the signal enrichment for ATAC-seq

- fraction of reads in peaks (FRiP)
- enrichments around active TSS
- visual inspection (genome browser!)



# Checking the signal enrichment: generating coverage files



- deepTools [Ramírez et al., 2016] offers the bamCoverage function that is fairly versatile and flexible
  - check out the documentation!
  - ▶ 2 types of normalization to account for sequencing depth differences
    - **RPGC** (reads per genomic content) will divide the reads per bin by the *coverage* (calculated based on effective genome size); this will make different samples comparable to each other recommended
    - **RPKM**: division by total number of reads

```
bamCoverage --bam a.bam -o a.SeqDepthNorm.bw --binSize 10 \
    --normalizeUsing RPGC --effectiveGenomeSize 2150570000 \
    --ignoreForNormalization chrX -minFragmentLength 40
```

#### Checking the signal enrichment: TSS focus



deepTools offers functions for visualizations of the bigWig files

```
plotHeatmap -m ATAC_TSS.tab.gz \
  -out hm_ATAC.png \
  --heatmapHeight 15 \
  --refPointLabel center
```

### Checking the signal enrichment: peak calling

= identifying regions with higher read coverage than expected based on the background



# Checking the signal enrichment: peak calling

Starting from the BAM file:

- 1 generate a signal of *fragment* counts along the genome
- 2 identify regions of enrichment



We usually use MACS [Zhang et al., 2008]; mostly because it's part of most pipelines, not because it's such a great tool (but it has proven itself to be fairly robust and useful).

F. Dündar (ABC, WCM)	Genomic enrichment strategies: O	ChIP-seq and	April 9, 2019	48 / 76
----------------------	----------------------------------	--------------	---------------	---------

# Peak calling

Identifying and assessing regions of enrichment with MACS

- Sliding a window of length 2 x bandwidth (= half of estimated sonication size) across genome and determine read counts
- ② Retain windows with counts > MFOLD (fold-enrichment of treatment/back-ground)
- ③ PEAKS: probability of an enrichment being stronger than expected
  - H0: reads are randomly distributed throughout the genome following a Poisson distribution
  - Determine the background distribution (λ) by sliding a window of size 2 x fragment size across the background to estimate the local coverage

```
MACS2 callpeak -t pairedEnd.bam -f BAMPE --outdir path/to/output/ \
--name pairedEndPeakName -g 2.7e9
```

See Tom Carroll's pipeline for detailed MACS2 commands.

# The result of MACS is a BED file of regions with sign. enrichments, i.e. peaks.

### Checking signal enrichments: FRiP

 $FRiP = \frac{reads in peaks}{total reads}$ 



FRiP > 0.3 is optimal; FRiP > 0.2 acceptable by ENCODE standards.

# QC checklist ATAC-seq

- fragments of 40 100 bp size should be over-represented
- 1/3 of the reads should fall into peaks (FRiP)
- very sharp and not too broad enrichments around TSS of active genes
- IGV snapshots: the signal should look sharp and high



ChIP-seq principles

#### ChIP-seq principles

# NGS techniques for studying chromatin and DNA modifications

Depending on the type of insights you're interested in, there are different ways of *enrichment*.

How to enrich for the NA	Biological insights	Example technique
Nuclease susceptibility	nucleosome packaging	DNase-seq, MNase-seq
	regulatory regions	ATAC-seq
Affinity-based enrichments	protein-DNA interactions	ChIP-seq
	histone modifications	
	protein-RNA interactions	CLIP-seq
	chromatin-chromatin interactions	ChIA-PET
	RNA modifications	m6A-seq, MeRIP-Seq,
Proximity ligation	chromatin-chromatin interactions	3C, Hi-C, ChIA-PET,
Chemical susceptibility	DNA modifications	WGBS, RRBS

Table based on Friedman and Rando [2015]

F. Dündar (ABC, WCM)

### Identifying transcription factor binding sites with ChIP



The vast majority of TFBS has been found in regions of open chromatin.

F. Dündar (ABC, WCM)

ChIP-seq principles

#### Extracting DNA sites bound by a TF of interest



Genomic enrichment strategies: ChIP-seq and

# In contrast to ATAC-seq, nobody would say ChIP-seq was "easy"



#### • depends on antibodies

- expensive! (typically 1 vial per experiment)
- cross-reactivity
- lack of affinity/binding needs incredibly optimized conditions
- signal-to-noise ratio will depend on how abundantly the protein of interest binds to DNA
- sonication can be fickle and inherently favors open chromatin regions
- cross-linking is a frequent source of bias
- takes 3-4 days to complete
- requires lots of cells (1-10 mio)

See, for example, Jordán-Pla and Visa [2018] for how to optimize ChIP experiments.

# ChIP enrichments are often marginal and variable across experiments



# Different types of ChIP'ed factors will yield different types of signals



# ChIP experiment absolutely require an "input" control

= basically, the ChIP experiment without the antibody addition



Ideally, input samples should be done in parallel with the ChIP experiments; they should also be sequenced at least as deeply or **more deeply sequenced** than the ChIP samples.

Processing of ChIP-seq data

#### Processing of ChIP-seq data





many basic processing steps are the same for ATAC- and ChIP-seq data, but some QC scores differ

F. Dündar (ABC, WCM)

# Peak calling: different ChIP'ed factors require different peak callers

Identifying peaks for sharp, narrow, high enrichments is easy (-> MACS). Assigning stats to broad enrichment is still an unsolved issue.



F. Dündar (ABC, WCM)

Genomic enrichment strategies: ChIP-seq and

# Peak calling: different ChIP'ed factors require different peak callers

Identifying peaks for sharp, narrow, high enrichments is easy (-> MACS). Assigning stats to broad enrichment is still an unsolved issue.

★	Comprehensive list is at: https://omictools.com/peak-calling-category
---	---

MACS2 (MACS1.4)	Most widely used peak caller. Can detect narrow and broad peaks.
Epic (SICER)	Specialised for broad peaks
BayesPeak	R/Bioconductor
Jmosaics	Detects enriched regions jointly from replicates
Т-РІС	Shape based
EDD	Detects megabase domain enrichment
GEM	Peak calling and motif discovery for ChIP-seq and ChIP-exo
SPP	Fragment length computation and saturation analysis to determine if read depth is adequate.

F. Dündar (ABC, WCM)	Genomic enrichment strategies: ChIP-se	eq and April 9, 2019	63 / 76
----------------------	--	----------------------	---------

### Peak calling: take input samples into consideration!





Consider the bioconductor package GreyListChIP to define cell-type-specific regions of input biases.

### Signal check: fingerprints instead of FRiP

#### How well can signal & background be separated?

A very specific and strong ChIP enrichment will be indicated by a prominent and steep rise of the cumulative sum towards the highest rank. This means that a big chunk of reads from the ChIP sample is located in few bins which corresponds to high, narrow enrichments typically seen for transcription factors.



when counting the reads contained in 97% of all genomic bins, only 55% of the maximum number of reads are reached, i.e. 3% of the genome contain a very large fraction of reads!

this indicates very **localized**, very **strong** enrichments! (as every biologist hopes for in a ChIP for H3K4me3)

#### ## another deepTools function

\$ plotFingerprint -b testFiles/\*bam --labels H3K4me3 H3K4me1 H3K27me3 \
 --plotFile fingerprints.png --outRawCounts fingerprints.tab

### Signal check: fingerprints instead of FRiP

all (i.e. bins containing zero reads)



 more difficult to distinguish input and ChIP, it does not mean, however, that this particular ChIP experiment failed

#### Comparing different ChIP-seq experiments

- comparing the levels of ChIP (and ATAC)-seq enrichments across different conditions is more difficult than one would have hoped for [Guertin et al., 2018]
  - Steinhauser et al. [2016] did a comparison of differential ChIP-seq tools
  - the winner tends to be the bioconductor package DiffBind, which is basically a sophisticated wrapper around DESeq
- relatively few efforts have been made towards understanding ChIP-seq/ATAC-seq-specific data properties, but the general consensus is that particularly ChIP-seq is awfully noisy and dependent on too many experimental parameters

"Although we would ideally want to study the absolute levels of binding, we have to accept the limitations of ChIP-seq [and ATAC-seq] and adapt by designing experiments in such a way that meaningful conclusions can be drawn from relative levels." [Meyer and Liu, 2014]

References

### References

[Buenrostro et al., 2013, Chen et al., 2014, Qin et al., 2016, Cowie et al., 2013, Ramírez et al., 2016, Ernst and Kellis, 2017, Friedman and Rando, 2015, Gaffney et al., 2012, Klemm et al., 2019, Zhang et al., 2008, Meyer and Liu, 2014, Montefiori et al., 2017, Mundade et al., 2014, Nakato and Shirahige, 2017, Ou et al., 2018, Park, 2009, Ptashne, 2013, Splinter and De Laat, 2011, Stricker et al., 2016, Thomas et al., 2017, Waddington, 1942, Wei et al., 2018, Wilbanks and Facciotti, 2010, Winter et al., 2015]

#### References

- Jason D. Buenrostro, Paul G. Giresi, Lisa C. Zaba, Howard Y. Chang, and William J. Greenleaf. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nature Methods*, 2013. doi: 10.1038/nmeth.2688.
- R. Chen, R. Kang, X. G. Fan, and D. Tang. Release and activity of histone in diseases. *Cell Death and Disease*, 2014. doi: 10.1038/cddis.2014.337.
  Philip Cowie, Ruth Ross, and Alasdair MacKenzie. Understanding the Dynamics of Gene Regulatory Systems; Characterisation and Clinical Relevance of cis-Regulatory Polymorphisms. *Biology*, 2013. doi: 10.3390/biology2010064.
- Timothy Daley and Andrew D. Smith. Predicting the molecular complexity of sequencing libraries. *Nature Methods*, 2013. ISSN 15487091. doi: 10.1038/nmeth.2375.
- Jason Ernst and Manolis Kellis. Chromatin-state discovery and genome annotation with ChromHMM. *Nature Protocols*, 2017. doi: 10.1038/nprot.2017.124.

- Nir Friedman and Oliver J. Rando. Epigenomics and the structure of the living genome. *Genome Research*, 2015. doi: 10.1101/gr.190165.115.
  Daniel J. Gaffney, Graham McVicker, Athma A. Pai, Yvonne N. Fondufe-Mittendorf, Noah Lewellen, Katelyn Michelini, Jonathan Widom, Yoav Gilad, and Jonathan K. Pritchard. Controls of Nucleosome Positioning in the Human Genome. *PLoS Genetics*, 2012. doi: 10.1371/journal.pgen.1003036.
- Michael J. Guertin, Amy E. Cullen, Florian Markowetz, and Andrew N. Holding. Parallel factor ChIP provides essential internal control for quantitative differential ChIP-seq. *Nucleic Acids Research*, 2018. doi: 10.1093/nar/gky252.
- Antonio Jordán-Pla and Neus Visa. Considerations on experimental design and data analysis of chromatin immunoprecipitation experiments. In *Methods in Molecular Biology*. 2018. doi: 10.1007/978-1-4939-7380-4\_2.
- Sandy L. Klemm, Zohar Shipony, and William J. Greenleaf. Chromatin accessibility and the regulatory epigenome. *Nature Reviews Genetics*, 2019. doi: 10.1038/s41576-018-0089-8.
- Clifford A Meyer and X Shirley Liu. Identifying and mitigating bias in next-generation sequencing methods for chromatin biology. *Nature Reviews Genetics*, 2014. doi: 10.1038/nrg3788.
- Lindsey Montefiori, Liana Hernandez, Zijie Zhang, Yoav Gilad, Carole Ober, Gregory Crawford, Marcelo Nobrega, and Noboru Jo Sakabe. Reducing mitochondrial reads in ATAC-seq using CRISPR/Cas9. *Scientific Reports*, 2017. doi: 10.1038/s41598-017-02547-w.
- Rasika Mundade, Hatice Gulcin Ozer, Han Wei, Lakshmi Prabhu, and Tao Lu. Role of ChIP-seq in the discovery of transcription factor binding sites, differential gene regulation mechanism, epigenetic marks and beyond. *Cell Cycle*, 2014. doi: 10.4161/15384101.2014.949201.
- Ryuichiro Nakato and Katsuhiko Shirahige. Recent advances in ChIP-seq analysis: From quality management to whole-genome annotation. *Briefings in Bioinformatics*, 2017. doi: 10.1093/bib/bbw023.

- Jianhong Ou, Haibo Liu, Jun Yu, Michelle A. Kelliher, Lucio H. Castilla, Nathan D. Lawson, and Lihua Julie Zhu. ATACseqQC: A Bioconductor package for post-alignment quality assessment of ATAC-seq data. BMC Genomics, 2018. doi: 10.1186/s12864-018-4559-3.
- Peter J Park. ChIP-seq: advantages and challenges of a maturing technology. *Nature Reviews Genetics*, 10(10):669–80, Oct 2009. doi: 10.1038/nrg2641.
- M. Ptashne. Epigenetics: Core misconcept. Proceedings of the National Academy of Sciences, 2013. doi: 10.1073/pnas.1305399110.
  Qian Qin, Shenglin Mei, Qiu Wu, Hanfei Sun, Lewyn Li, Len Taing, Sujun Chen, Fugen Li, Tao Liu, Chongzhi Zang, Han Xu, Yiwen Chen, Clifford A. Meyer, Yong Zhang, Myles Brown, Henry W. Long, and X. Shirley Liu. ChiLin: A comprehensive ChIP-seq and DNase-seq quality control and analysis pipeline. *BMC Bioinformatics*, 2016. doi: 10.1186/s12859-016-1274-4.

Fidel Ramírez, Devon P. Ryan, Björn Grüning, Vivek Bhardwaj, Fabian Kilpert, Andreas S. Richter, Steffen Heyne, Friederike Dündar, and Thomas Manke. deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic acids research*, 2016. ISSN 13624962. doi: 10.1093/nar/gkw257.

- Erik Splinter and Wouter De Laat. The complex transcription regulatory landscape of our genome: Control in three dimensions. *EMBO Journal*, 2011. doi: 10.1038/emboj.2011.344.
- Sebastian Steinhauser, Nils Kurzawa, Roland Eils, and Carl Herrmann. A comprehensive comparison of tools for differential ChIP-seq analysis. *Briefings in Bioinformatics*, 2016. doi: 10.1093/bib/bbv110.
- Stefan H. Stricker, Anna Köferle, and Stephan Beck. From profiles to function in epigenomics. *Nature Reviews Genetics*, 2016. doi: 10.1038/nrg.2016.138.
- Reuben Thomas, Sean Thomas, Alisha K. Holloway, and Katherine S.Pollard. Features that define the best ChIP-seq peak calling algorithms.*Briefings in Bioinformatics*, 2017. doi: 10.1093/bib/bbw035.

- C. H. Waddington. The epigenotype. . International Journal of Epidemiology (2012), 1942. doi: 10.1093/ije/dyr184.
  Zheng Wei, Wei Zhang, Huan Fang, Yanda Li, and Xiaowo Wang. esATAC: an easy-to-use systematic pipeline for ATAC-seq data analysis. Bioinformatics, 2018. doi: 10.1093/bioinformatics/bty141.
  Elizabeth G. Wilbanks and Marc T. Facciotti. Evaluation of algorithm performance in ChIP-seq peak detection. PLoS ONE, 2010. doi: 10.1371/journal.pone.0011471.
- Deborah R. Winter, Steffen Jung, and Ido Amit. Making the case for chromatin profiling: A new tool to investigate the immune-regulatory landscape. *Nature Reviews Immunology*, 2015. doi: 10.1038/nri3884.
  Yong Zhang, Tao Liu, Clifford A. Meyer, Jérôme Eeckhoute, David S. Johnson, Bradley E. Bernstein, Chad Nussbaum, Richard M. Myers, Myles Brown, Wei Li, and X. Shirley Shirley. Model-based analysis of ChIP-Seq (MACS). *Genome Biology*, 2008. doi: 10.1186/gb-2008-9-9-r137.