

# ATAC-seq analysis

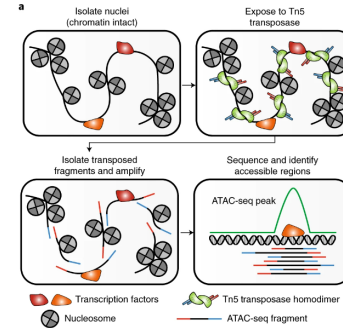
Bingbing Yuan

BaRC Hot Topics –April 4<sup>th</sup> 2024  
Bioinformatics and Research Computing  
Whitehead Institute

[http://barc.wi.mit.edu/hot\\_topics/](http://barc.wi.mit.edu/hot_topics/)



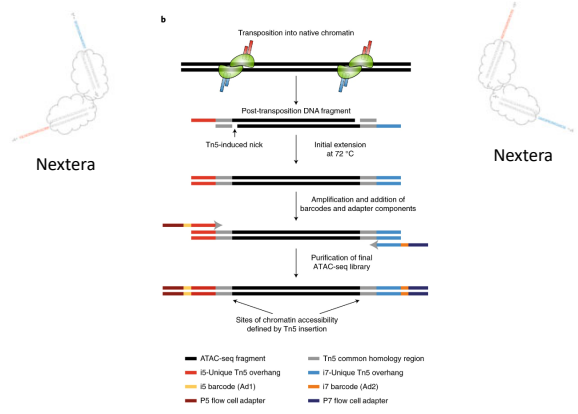
# Assay for Transposase-Accessible Chromatin with high-throughput sequencing (ATAC-seq)



*Nat Protoc* 17, 1518–1552 (2022)



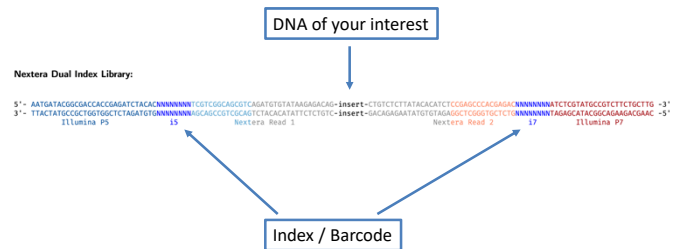
# ATAC-seq



*Nat Protoc* 17, 1518–1552 (2022)  
<https://www.ploeb.org/article/11443.html>



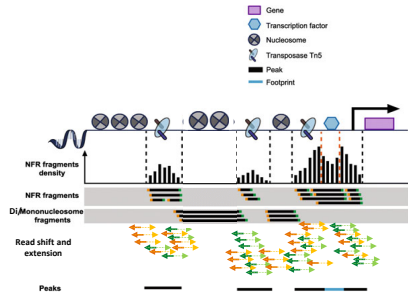
# Nextera library



Modified based on [https://teichlab.github.io/scg\\_lib\\_structs/methods\\_html/illumina.html](https://teichlab.github.io/scg_lib_structs/methods_html/illumina.html)



# Profile generation in ATAC-seq



Modified from Feng Yan et al. Genome Biology 21 (2020)

- > Open chromatin
- > Motif enrichment
- > TF footprinting
- > Nucleosome positions

5



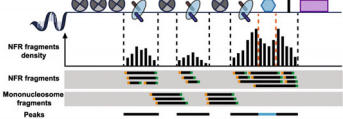
# ATAC-seq pipeline goal

- Using short-read sequencing, identify genome-wide regions of open chromatin

6



# ATAC-seq analysis workflow



NFR: nucleosome free regions  
Feng Yan et al.  
Genome Biology 21 (2020)

## Steps in ATAC-seq data analysis

Goal: Find the open chromatin regions  
Peak = open chromatin region

1. Quality control  
Remove adapters if necessary
2. Mapping  
Tailored to paired end reads
3. Post-alignment filtering
4. Peak calling and differential analysis
  - i) Read shift and extension and signal profile generation.
  - ii) Peak assignment
5. Peak interpretation
  - i) Find genes next the open chromatin regions
  - ii) Find motifs within peaks

7



# ENCODE ATAC-seq pipeline

- > Human and mouse with biological replicates
  - > QC report ([a sample report](#))
  - > Includes all the steps in a single run  
[Instructions to run pipeline using Whitehead server](#)
1. Pre-alignment quality control
  2. Aligning reads to genome
  3. Post-alignment filtering
  4. Post-alignment quality control
  5. Peak (accessible regions) calls
  6. Assessing Peak Calls with FRiP score (same as ChIP-seq)
  7. Blacklist filtering for peaks (same as ChIP-seq)

8



## What sequencing works best for ATAC-seq?

- Read depth recommended by ENCODE:
  - 50 million for paired-ended sequencing (25 million fragments)
- No input control sample
- Shorter read lengths (50x50 or 75x75) better than longer reads (100x100 or longer)
  - Ensure NFR
- Pair-end (PE) reads are recommended over single reads
  - Both ends of DNA fragments are considered
  - More detailed information about nucleosome packing and positioning in PE
  - Avoid losing reads due to mapping to repetitive elements



9



## Illumina data format

### • Fastq format:

[http://en.wikipedia.org/wiki/FASTQ\\_format](http://en.wikipedia.org/wiki/FASTQ_format)

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhghhhhhhhhehhhedhhhhfhhhhhh
```

- /1 or /2 paired-end
- @seq identifier
- seq
- +any description
- seq quality values

Input qualities	Illumina versions
--solexa-quals	<= 1.2
--phred64	1.3-1.7
--phred33	>= 1.8



10



## Check read quality with FastQC

(<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>)

1. Run FastQC to check read quality

fastqc sample.fastq

2. Open output file: "fastqc\_report.html"

### FastQC Report

#### Summary

- ✔ Basic Statistics
- ✔ Per base sequence quality
- ✔ Per tile sequence quality
- ✔ Per sequence quality scores
- ⚠ Per base sequence content
- ✔ Per sequence GC content
- ✔ Per base N content
- ✔ Sequence Length Distribution
- ✔ Sequence Duplication Levels
- ✔ Overrepresented sequences
- ✔ Adapter Content
- ⚠ Kmer Content

#### Basic Statistics

Measure	Value
Filename	Hepg2H3kme3_subset.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	1160004
Filtered Sequences	0
Sequence length	36
%GC	45



## Pre-alignment quality control

1. Check reads quality with fastqc

fastqc read1.fq read2.fq

2. Remove adapter when necessary:

trim\_galore --fastqc -nextera --paired --length 30 read1.fq read2.fq

- -nextera: ATAC-seq experiments use the Nextera DNA Library Prep Kit
- --paired: both reads need to pass or they are both removed.
- --length: discard trimmed reads shorter than this length.



12



## Local genomic files needed for mapping

tak: /nfs/genomes/

- Human, mouse, zebrafish, *C.elegans*, fly, yeast, etc.
- Different genome builds
  - mm10: mouse\_mm10\_dec\_11
  - mm39: mouse\_mm39\_jun20
- human\_hg38\_dec13 vs human\_hg38\_dec13\_no\_random
  - human\_hg38\_dec13 includes \*\_random.fa, \*hap\*.fa, etc.
- Sub directories:
  - bowtie
    - Bowtie1: \*.ebwt
    - Bowtie2: \*.bt2
  - fasta: one file per chromosome
  - fasta\_whole\_genome: all sequences in one file
  - gtf: gene models from Refseq, Ensembl, etc.



13



## Map reads to reference genome

Map reads with a non-spliced mapping tools: bowtie2 or BWA

```
bowtie2 --very-sensitive --no-discordant -p 2 -X 2000 -x hg38 -1  
read1.fq -2 read2.fq | samtools view -ub - | samtools sort - > |  
bowtie_out.bam
```

- --no-discordant:
  - Suppress discordant alignments for paired reads
- -X 2000:
  - Increase maximum fragment length to 2k to include nucleosome distribution
  - Used for plotting fragment size distribution



14



## Post alignment filtering

- Remove reads with low quality score: MAPQ < 30  
`alignmentSieve -b file.bam --minMappingQuality 30 --  
samFlagInclude 2 -o MAPQ30.bam"`
- Remove duplicates with Picard's 'MarkDuplicates'  
`java -jar picard.jar MarkDuplicates I=foo.bam O=noDups.bam  
M=foo.marked_dup_metrics.txt REMOVE_DUPLICATES=true`
- Remove reads mapped to mitochondria  
`samtools view -h file.bam | grep -v chrM | samtools view -b -  
-h -f 0x2 - | samtools sort - > file.sorted.bam`



`-f 0x2`  
`samFlagInclude 2` } keep only properly paired reads

15



## Post-alignment quality control

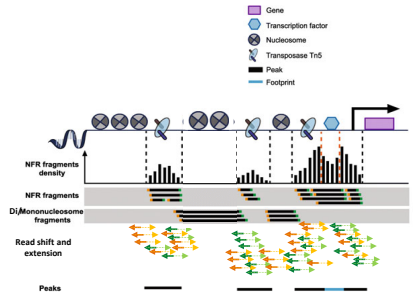
- Fragment size distribution
- TSS enrichment score



16



# Profile generation in ATAC-seq



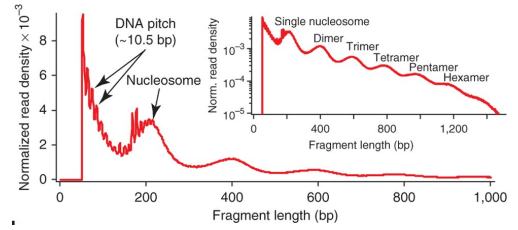
Modified from Feng Yan et al Genome Biology 21 (2020)

- > Open chromatin
- > Motif enrichment
- > TF footprinting
- > Nucleosome positions



# Fragment size distribution

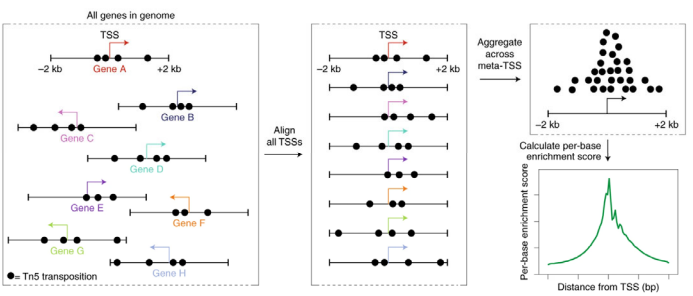
```
calculate_ATACseq_fragment_size_distribution.R
library("ATACseqQC")
pdf("sample.fragment_sizes.pdf", w=11, h=8.5)
fragSizeDist("sample.bam", "sample")
dev.off()
```



Jason D Buenrostro, et al. Nature Methods 10 (2013)



# (Transcription Start Site) TSS enrichment score



Nat Protoc 17, 1518–1552 (2022)



# (Transcription Start Site) TSS enrichment score

ENCODE standard with one TSS per gene:  
 • Score with TSS per transcript is smaller than with one TSS per gene

Annotation used	Value	Resulting Data Status
hg19 Refseq TSS annotation	< 6	Concerning
	6 - 10	Acceptable
	> 10	Ideal
GRCh38 Refseq TSS annotation	< 5	Concerning
	5 - 7	Acceptable
	> 7	Ideal
mm9 GENCODE TSS annotation	< 5	Concerning
	5 - 7	Acceptable
	> 7	Ideal
mm10 Refseq TSS annotation	< 10	Concerning
	10 - 15	Acceptable
	> 15	Ideal

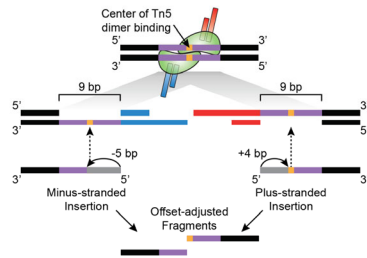
calculate\_TSS\_enrichment\_score.py





## Adjusted the read start sites to represent the center of the transposon binding event

- The Tn5 adapters are inserted in a staggered manner into the 5' ends of target sequence strands with a 9-bp gap between them
- The center of the Tn5 binding is 4 bp to the right of the edge on positive-strand reads, or 5 bp to the left on negative-strand reads.



Nat Protoc 17, 1518–1552 (2022)

25



## Shift reads

- Reads should be shifted + 4 bp and - 5 bp for positive and negative strand respectively, to account for the 9-bp duplication created by DNA repair of the nick by Tn5 transposase

```
cat foo.bed | awk -F '\t' 'BEGIN {OFS = FS}{ if ($6 == "+") {$2 = $2 + 4} else if ($6 == "-") {$3 = $3 - 5} print $0}' > |
foo_tn5.bed
```

Bed format:

```
Chr start end read_name score strand
chr1 3199164 3199215 K00168:88:HFF7YBBXX:1:1203:28371:20304/2 42 +
chr1 3199702 3199753 K00168:88:HFF7YBBXX:1:1203:28371:20304/1 42 -
```

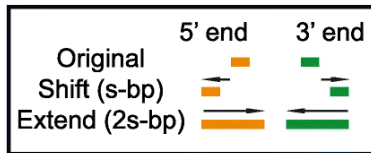
26



## Call peaks

```
macs2 callpeak -t foo_tn5.bed -n foo -f BED -g mm -q
0.01 --nomodel --shift -75 --extsize 150 --call-summits --
keep-dup all
```

Shift the reads to create a signal profile centered around the cutting sites



Feng Yan et al Genome Biology 21 (2020)

27



## MACS output

Output files:

- Excel peaks file (“\_peaks.xls”) contains the following columns  
Chr, start, end, length, abs\_summit, pileup, -LOG10(pvalue), fold\_enrichment, -LOG10(qvalue), name
- “\_summits.bed”: contains the peak summits locations for every peaks. The 5th column in this file is -log10qvalue
- “\_peaks.narrowPeak” is BED6+4 format file. Contains the peak locations together with peak summit, fold-change, pvalue and qvalue.

28



## Visualize peaks in IGV



29



## FRiP (Fraction of reads in peaks) score

- Fraction of all mapped reads that fall into the called peak regions
- The higher the score, the better
  - According to [ENCODE](#), score is preferably over 0.3, values greater than 0.2 are acceptable.

```
calculate_FRiP_score.py Sample.bam Sample_peaks.narrowPeak
```

30



## Blacklist filtering for peaks

- Anomalous, unstructured, or high signal in next-generation sequencing experiments independent of cell line or experiment
- Often in repeats (centromeres, telomeres, satellite repeats)
  - 0.5% of genome, but could account for >10% total signal
- BaRC\_datasets -> ENCODE\_blacklists
  - ENCODE: human, mouse, fly, or C. elegans

```
bedtools intersect -v -a foo_peaks.narrowPeak -b  
blacklist.bed > bfiltered_peaks.narrowPeak  
# -v option: only report those "A" peaks with no overlaps in "B"
```

31



## Compare open chromatin between different experimental conditions

- Without replicates, you can use bedtools to compare two samples:
  - intersectBed (finds the subset of peaks **common** in 2 samples **or unique** to one them)
- If you have replicates, you can use:
  - 1) bedtools merge to merge the peaks
  - 2) bedtools coverage to count the number of reads in peaks
  - 3) DESeq2/EdgeR on reads assigned the peaks to get differentially open peaks

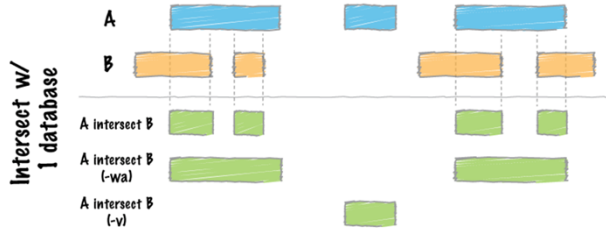
32



# bedtools

- intersectBed

**intersect**



<https://bedtools.readthedocs.io/>

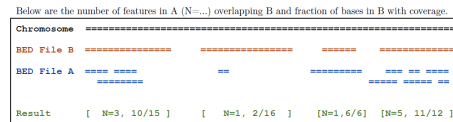


# bedtools

**merge**



**coverage**



## Linking peaks to genes: Bedtools

**slopBed**: extend a feature by a user-defined number of bases

```

$ cat A.bed
chr1 5 100
chr1 500 900
$ cat my.genome
chr1 3500
$ slopBed -i A.bed -g my.genome -b 3
chr1 0 105
chr1 795 985
$ slopBed -i A.bed -g my.genome -l 2 -r 3
chr1 5 102
chr1 798 983
    
```

**groupBy**

It groups rows based on the value of a given column/s and it summarizes the other columns

**closestBed**



## Linking peaks to nearby genes:

Get all the genes at a certain distance (*i.e.* 3Kb) of the peak.  
The distance we use depends on the area where we want to find regulatory interactions.

- Take all genes and add 3Kb up and down with slopBed
  - `slopBed -b 3000 -i GRCh37.p13.HumanENSEMBLgenes.bed -g /nfs/genomes/human_gp_feb_09_no_random/anno/chromInfo.txt > HumanGenesPlusMinus3kb.bed`
- Intersect the slopped genes with peaks and get the list of unique genes overlapping
 

```
intersectBed -wa -a HumanGenesPlusMinus3kb.bed -b peaks.bed | awk '{print $4}' | sort -u > Genesat3KborlessfromPeaks.txt
```

```

intersectBed -wa -a HumanGenesPlusMinus3kb.bed -b peaks.bed | head -3
chr1 45956538 45968751 ENSG00000236624_CCDC163P
chr1 45956538 45968751 ENSG00000236624_CCDC163P
chr1 51522509 51528577 ENSG00000265538_MIR4421
    
```



## Link peaks to closest gene

For each region find the closest gene and filter based on the distance to the gene

The commands below are an example where we are looking for interactions at 3Kb or less.

```
closestBed -d -a peaks.bed -b GRCh37.p13.HumanENSEMBLgenes.bed | head
chr1 20870 21204 H3k4me3_chr1_peak_1 5.77592 chr1 14363 29806 ENSG00000227232_WASH7P 0
chr1 28482 30214 H3k4me3_chr1_peak_2 374.48264 chr1 29554 31109 ENSG00000243485_MIR1302-10 0
chr1 28482 30214 H3k4me3_chr1_peak_2 374.48264 chr1 14363 29806 ENSG00000227232_WASH7P 0
#the next two steps can also be done on excel
closestBed -d -a peaks.bed -b GRCh37.p13.HumanENSEMBLgenes.bed | groupBy -g 9,10 -c
6,7,8, -o distinct,distinct,distinct | head -3
ENSG00000227232_WASH7P 0 chr1 14363 29806
ENSG00000243485_MIR1302-10 0 chr1 29554 31109
ENSG00000227232_WASH7P 0 chr1 14363 29806
```

```
closestBed -d -a peaks.bed -b GRCh37.p13.HumanENSEMBLgenes.bed | groupBy -g 9,10 -c 6,7,8, -o
distinct,distinct,distinct | awk 'BEGIN {OFS="\t"} {if ($2<3000) {print $3,$4,$5,$1,$2}}' | head -5
chr1 14363 29806 ENSG00000227232_WASH7P 0
chr1 29554 31109 ENSG00000243485_MIR1302-10 0
chr1 14363 29806 ENSG00000227232_WASH7P 0
chr1 134901 139379 ENSG00000237683_AL627309.1 0
chr1 135141 135895 ENSG00000268903_RP11-34P13.15 0
```



37



## Link peaks to closest gene (1 command)

For each region find the closest gene and filter based on the distance to the gene

The command below is an example where we are looking for interactions at 3Kb or less.

```
closestBed -d -a peaks.bed -b
GRCh37.p13.HumanENSEMBLgenes.bed | groupBy -g 9,10 -c 6,7,8,
-o distinct,distinct,distinct | awk 'BEGIN {OFS="\t"} {if ($2<3000)
{print $3,$4,$5,$1,$2}}' > closestGeneAt3KborLess.bed
```

### closestBed

-d print the distance to the feature in -b

### groupBy

-g columns to group on

-c columns to summarize

-o operation to use to summarize



38



## Identify footprints

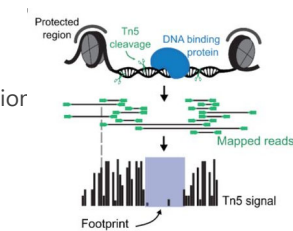
- DNA sequences directly bound by TFs
- High depth of coverage:  
At least ~200 million of reads
- Shift reads to account for the 9-bp duplication
- Tools:

**HINT:** (HmM-based IdeNticator footprints)

*Genome Biol* **20**, 45 (2019)

**TOBIAS:**

*Nat Commun* **11**, 4267 (2020)



Bentsen, M., Goymann, P., Schultheis, H. et al. *Nat Commun* **11** (2020)



39



## Single cell ATAC-seq (scATAC-seq)

- scATAC-seq is good for heterogeneous samples with multiple cell types or cell states, where bulk ATAC-seq considers average cell signals and could miss signals from a subset of cells
- Signal is sparse in scATAC-seq
- Remove cells with low quality or doublets
- Make cell-by-feature matrix
- Generate and annotate cell clusters

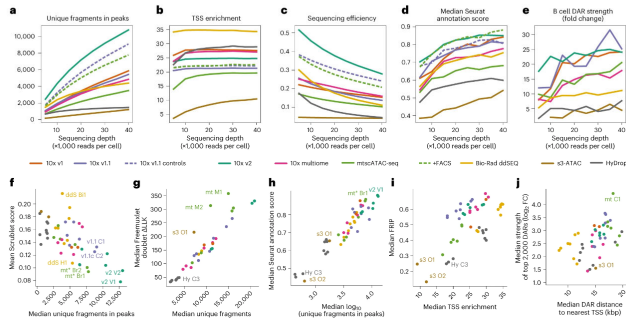


40



**Fig. 2: Differences in automated cell-type annotation accuracy and differential region calling between techniques.**

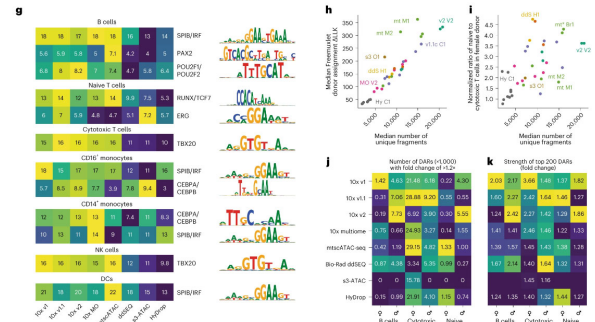
From: Systematic benchmarking of single-cell ATAC-sequencing protocols



De Rep, F.V., Hulselmans, G., Flerin, C. et al. *Nat Biotechnol* (2023)

41

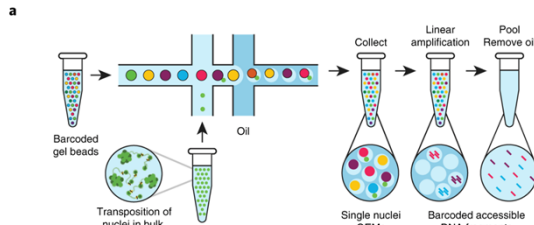
## Performance differences



De Rep, F.V., Hulselmans, G., Flerin, C. et al. *Nat Biotechnol* (2023)

42

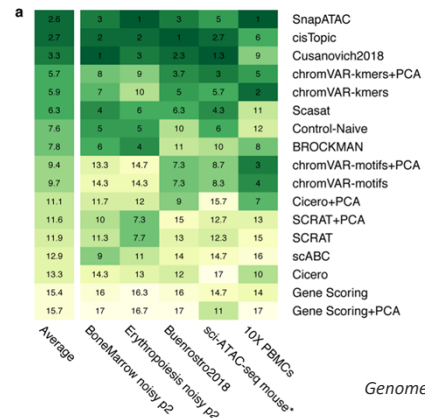
## Single cell ATAC-seq (scATAC-seq)



*Nat Biotechnol* 37, 925–936 (2019).

43

## scATAC-seq analysis tools



The darker shades of green indicates the better performance in clustering

44

# scATAC-seq analysis tools

**Table 1**  
Summary of scATAC-seq analysis software packages.

Tool	Platform	Feature Matrix	Preprocessing	Clustering	DAR	Motif/k-mer	Gene activity	Co-accessibility	Trajectory	Pathway	Enrichment analysis	scRNA integration	Reference
ChromVAR	R	TF motifs, k-mer	0	0	X	0	X	X	X	X	X	X	[17]
SCRAT	R/Web	Selectable feature	0	0	0	X	X	X	X	X	X	X	[18]
scABC	R	Peak	0	0	X	0 (ChromVAR)	X	X	X	X	X	X	[19]
Cicero	R	TSS	0	0	0	X	0	0	X	X	X	X	[20]
Scasat	Python/R	Peak	0	0	0	X	X	X	X	0 (GREAT)	X	X	[21]
cisTopic	R	Peak	0	0	X	X	0	X	X	0	0	X	[22]
snaiATAC	Python/R	Bin, peak	0	0	0	0 (ChromVAR, Homer)	0	X	X	0 (GREAT)	X	0 (Seurat)	[23]
epiScanpy	Python	Peak	0	0	X	X	X	X	X	X	X	X	[24]
Destin	R	Peak	0	0	0	X	X	X	X	0	X	X	[25]
SCALE	Python	Peak	0	0	0 (ChromVAR)	X	X	X	X	X	X	X	[26]
scATAC-pro	Python/R	Peak	0	0	0 (ChromVAR)	0	0 (Cicero)	X	X	0 (GREAT)	X	X	[27]
Signac	R	Peak	0	0	0 (ChromVAR)	0	X	X	X	X	0 (Seurat)	0 (Seurat)	[7]
ArchR	R	Bin, peak	0	0	0 (ChromVAR)	0	0	0	X	0	0 (Seurat)	0 (Seurat)	[28]

Tools used in junction are indicated in parentheses.

*Computational and Structural Biotechnology Journal 18 (2020)*

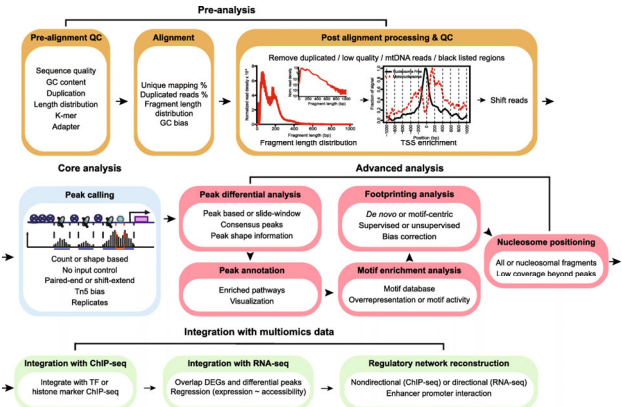
[Seurat scATAC-seq analysis pipeline:](#)

scRNA integration implements packages ( MACS2, monocle, Cicero, etc)

45



# ATAC-seq analysis



Feng Yan et al  
*Genome Biology 21 (2020)*

46



## Exercises:

- [Mouse lung tissue postnatal \(0 days\) from ENCODE](#)  
Randomly chosen 100,000 pair-end reads from the 1st replicate
- \\wi-files1\BaRC\_Public\Hot\_Topics\ ATAC-seq\ ATAC-seq\_2024\_commands.txt

47



## References

- o From reads to insight: a hitchhiker's guide to ATAC-seq data analysis. *Genome Biol* 21, 2020
- o Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* 10, 2013
- o ENCODE ATAC-seq guideline: <https://github.com/ENCODE-DCC/atac-seq-pipeline>
- o MACS:
  - o Model-based Analysis of ChIP-Seq (MACS). *Genome Biol* 9, 2008 <https://liulab-dfci.github.io/software/>
  - o Using MACS to identify peaks from ChIP-Seq data. *Curr Protoc Bioinformatics*. 2011
- o Picard Tools: <https://broadinstitute.github.io/picard/>
- o Bedtools: <https://code.google.com/p/bedtools/>
  - o BEDTools: a flexible suite of utilities for comparing genomic features, *Bioinformatics* 15, 2010
- o Massively parallel single-cell chromatin landscapes of human immune cell development and intratumoral T cell exhaustion. *Nat Biotechnol* 37, 925–936 (2019)
- o Single-cell ATAC sequencing analysis: From data preprocessing to hypothesis generation. *Computational and Structural Biotechnology Journal*, 18, 2020
- o Assessment of computational methods for the analysis of single-cell ATAC-seq data. *Genome Biol* 20, 241 (2019).
- o Identification of transcription factor binding sites using ATAC-seq. *Genome Biol* 20, 45 (2019)
- o ATAC-seq footprinting unravels kinetics of transcription factor binding during zygotic genome activation. *Nat Commun* 11, 4267 (2020)
- o Grandi, F.C., Modi, H., Kampman, L. et al. Chromatin accessibility profiling by ATAC-seq. *Nat Protoc* 17, 1518–1552 (2022).

48



## Other resources

- Previous Hot Topics  
Quality Control  
[http://barc.wi.mit.edu/education/hot\\_topics/NGS\\_QC\\_2017/slides4perPage.pdf](http://barc.wi.mit.edu/education/hot_topics/NGS_QC_2017/slides4perPage.pdf)
- Best Practices:  
[http://barcwiki.wi.mit.edu/wiki/SOPs/atac\\_Seq](http://barcwiki.wi.mit.edu/wiki/SOPs/atac_Seq)



49



## Upcoming Hot Topics sessions

- CUT & Tag analysis April 25<sup>th</sup>

Let us know if you're interested:

- Protein structure visualization
- CRISPR pooled screen analysis
- Enrichment analysis
- Other suggestions



50

