

Genomics Part II

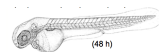
Applications of Sequencing Technology

Biomedical Data Science: Mining and Modeling
CB&B 752 • MB&B 452
Matt Simon
Feb 8, 2021

Workflow

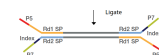
1. Isolation of sample.

e.g., Isolate DNA and shear.



2. Library preparation

e.g., Clean up and ligate Y-adaptors.



3. Sequencing

e.g., Illumina HiSeq



4. Analysis

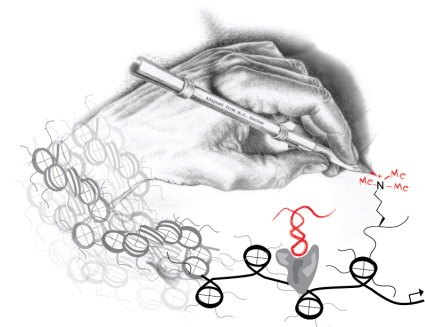
e.g., Map to genome and interpret.



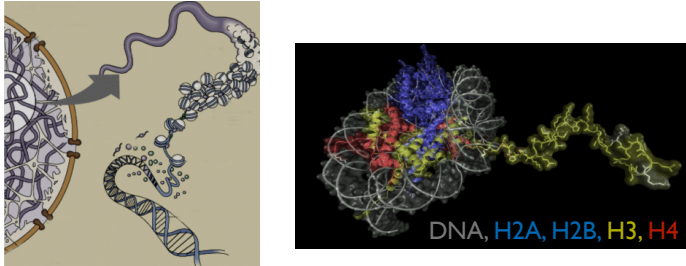
Overview

- Genomics I (Wednesday's lecture): Focus on sequencing technology and genomes.
- Genomics II: (Today's lecture): Focus on applications of sequencing technology.
 1. Annotation of the genome in chromatin
 2. Regulation of gene expression at the level of RNA

Part 1. How do cells annotate their genomes?



DNA in the cell is packaged into chromatin



Modeled nucleosome based on Luger et al., *Nature* **1997** 389, 251.

Summary and nomenclature of common covalent modifications.

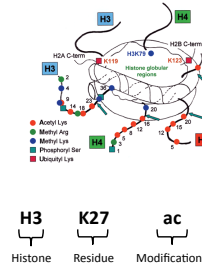
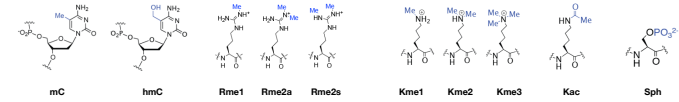
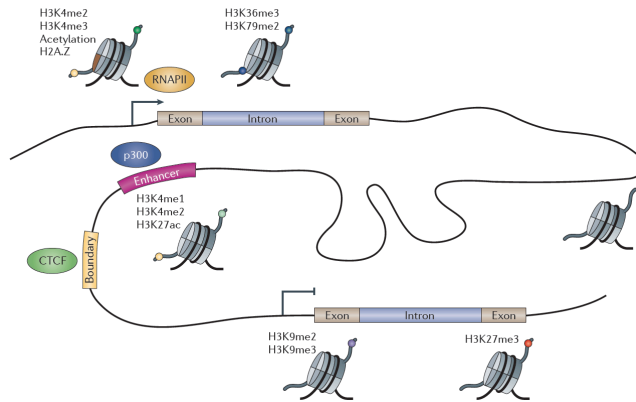


Table 1 The Brno nomenclature for histone modifications

Modifying group	Amino acid(s) modified	Level of modification	Abbreviation for modification ^a	Examples of modified residues ^b
Acetyl-	Lysine	mono-	ac	H3K9ac
	Arginine	mono-	me1	H3R17me1
		di-, symmetrical	me2s	H3R2me2s
	Arginine	di-, asymmetrical	me2a	H3R17me2a
Methyl-	Lysine	mono-	me1	H3K4me1
	Lysine	di-	me2	H3K4me2
	Lysine	tri-	me3	H3K4me3
	Phosphoryl-	Serine or threonine	mono-	ph
Ubiquityl-	Lysine	mono- ^c	ub1	H2BK122ub1
SUMOyl-	Lysine	mono-	su	H4K5su ^d
ADP-ribosyl-	Glutamate	mono-	ar1	H2BE2ar1
	Glutamate	poly-	am	H2BE2ar ^e

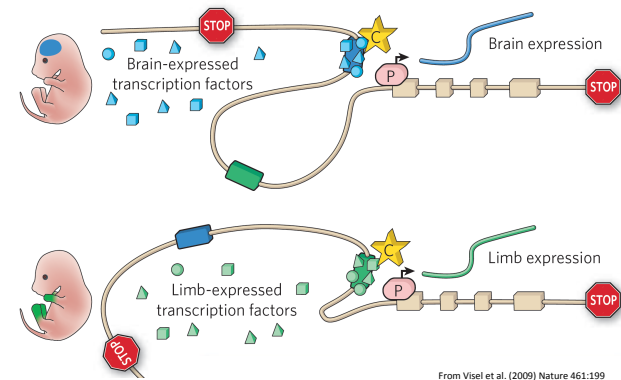
Turner, B. M. Reading signals on the nucleosome with a new nomenclature for modified histones. *Nat Struct Mol Biol* **12**, 110–112 (2005).

Chromatin modifications correlate with different genomic functions.



Zhou et al. *Nat Rev Genet* **12:7** (2011)

Regulation is temporally and specially controlled



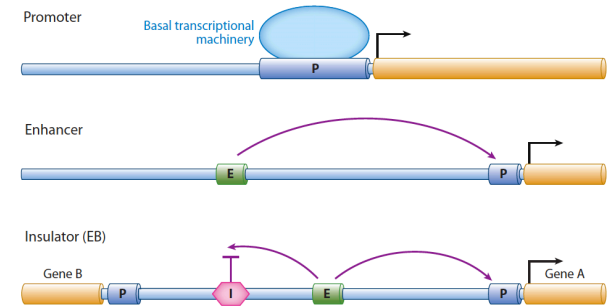
From Visel et al. (2009) *Nature* **461:199**

Using sequencing to annotate the genome

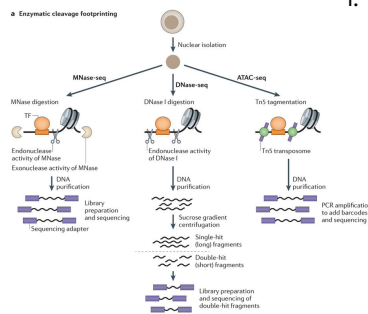
- Where are the cis-acting regulatory elements in DNA?
 - DNase I hyper-sensitivity mapping (**DNase-Seq**).
 - FAIRE** to map regulatory elements.
 - ATAC-Seq** to map regulatory elements.
- How does the chromatin composition vary across the genome?
 - ChIP-seq** of transcription factors (or in high res, ChIP-exo)
 - CUT&RUN** and **CUT&Tag** for small scale/single cell analysis.
- Where is RNA polymerase transcribing?
 - ChIP-Seq** of polymerase.
 - GRO-Seq**, **PRO-Seq** and **NET-Seq** to measure RNA polymerase activity.

Targeted approaches v Global approaches

How do we identify regulatory elements in the genome?



Using differences in biochemical properties of regulatory elements to identify them by Seq

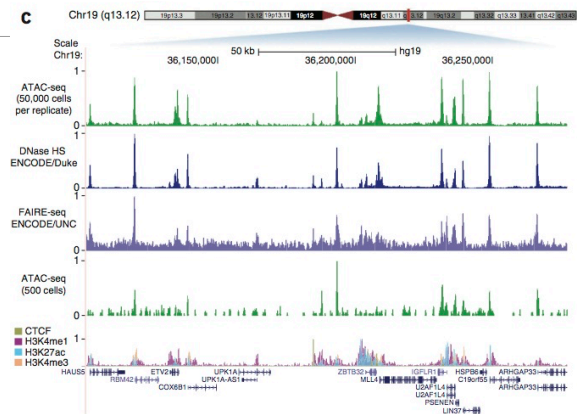


- Transcription factor binding** frequently deforms the B-form DNA, making it hypersensitive to DNase I and transposases.

Changes in **accessibility of chromatin** can provide information about regulation

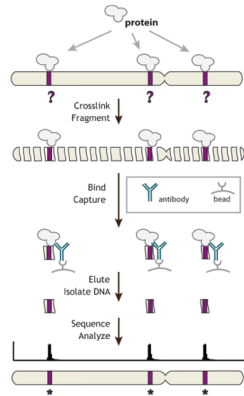
- ATAC-seq (shown)
- MNase-Seq (shown).
- DNase-Seq (shown).
- FAIRE-Seq (not shown).

Zentgraf GE, Henikoff S. High-resolution digital profiling of the epigenome. *Nat Rev Genet.* 2014;15: 814-827. doi:10.1038/nrg3798



Buenrostro JD, Giresi PG, Zaba LC, Chang HY, and Greenleaf WJ. (2013) "Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position." *Nature Methods*

Localization of *specific proteins* in the genome with chromatin immunoprecipitation (ChIP-Seq)



1. **Crosslink** the cells with formaldehyde to “fix” factors in place.

Exception: Native ChIP with histone antibodies.

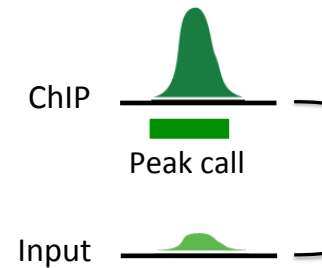
2. **Shear chromatin** to smaller pieces.

Shear size determines resolution. Note: ChIP-exo uses an exonuclease at a later step to increase resolution.

3. **Enrich** target using an antibody.

Enrichment is only as good as the antibody.

Determining sites of enrichment from ChIP-Seq

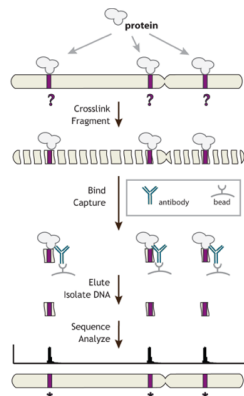


1. **Align** reads to the genome.

2. **Compare to input** to look for enrichment. Input coverage is not even.

3. **Call peaks** to determine statistically significant sites of enrichment.

Limitations of ChIP-Seq



1. **Cross linking** efficiency is not necessarily uniform.

2. Enrichment is dependent on the **quality of antibody**. e.g., Site and degree of histone modifications.

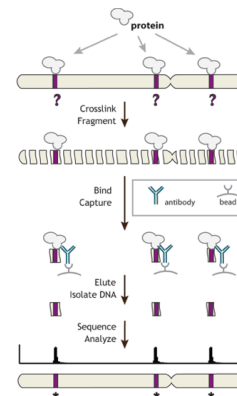
3. Enrichment is dependent on the **accessibility of the epitope**.

Comparing different sites to each other in the genome can be problematic.

4. Output is **descriptive**.

Hard to infer function without more experimentation.

Extensions of ChIP

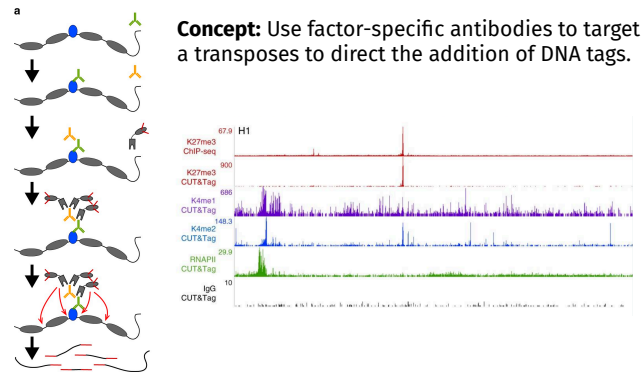


1. Using a nuclease to achieve **higher resolution** (ChIP-exo).

2. Analysis of **small samples or single cells** (CUT&RUN or CUT&Tag).

3. Extension to **RNA factors**.

Extensions of ChIP: CUT&Tag



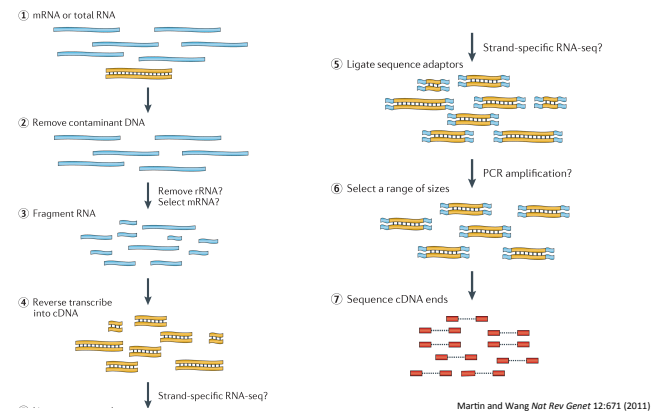
Kaya-Okur... & Henikoff (2019) CUT&Tag for efficient epigenomic profiling of small samples and single cells. *Nat Commun*

Part 2: RNA-Seq and applications of RNA-Seq

Using RNA-Seq to examine RNA

- Technical methodology
- Read mapping and normalization
- Estimating isoform-level gene expression
- De novo transcript reconstruction
- Sensitivity and sequencing depth
- Differential expression analysis

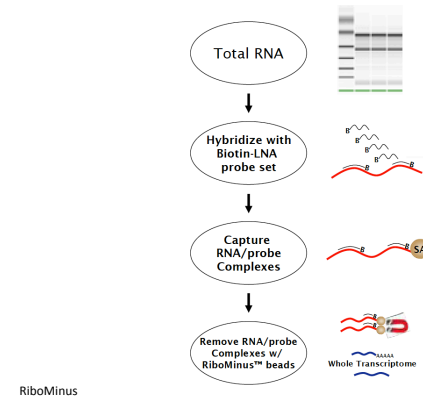
RNA-Seq workflow



Some technical details specific to RNA-Seq

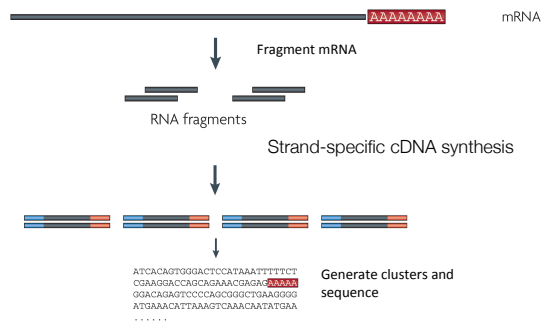
- Wide dynamic range of RNA concentrations.
- RNA is strand specific (unlike dsDNA)
- RNA degrades easily (RNase and spontaneous)
- RNA is processed (e.g., spliced)
- RNA has secondary structure (possible blocks to reverse transcriptase).

Ribosomal RNA will dominate the sequenced reads unless removed

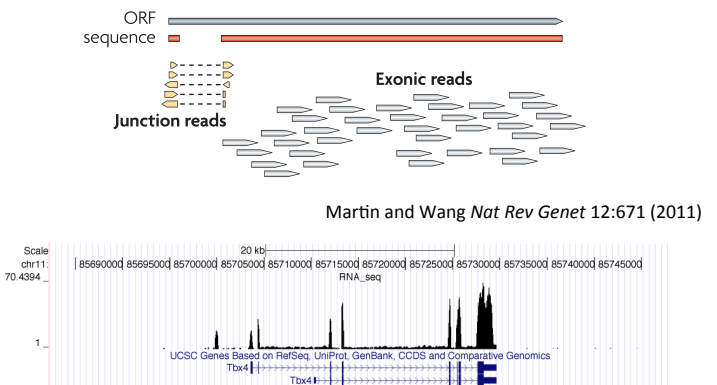


Illumina RNA-seq workflow

- Capture poly-A RNA with poly-T oligo attached beads (100 ng total) (2x)
- RNA quality must be high – degradation produces 3' bias
 - Non-poly-A RNAs are not recovered



RNA-Seq reads map mostly to exons



How does one analyze RNA levels from RNA-Seq?

Use existing gene annotation:

- Align to genome plus annotated splices
- Depends on high-quality gene annotation
- Which annotation to use: RefSeq, GENCODE, UCSC?
- Isoform quantification?
- Identifying novel transcripts?

Reference-guided alignments:

- Align to genome sequence
- Infer splice events from reads
- Allows transcriptome analyses of genomes with poor gene annotation

De novo transcript assembly:

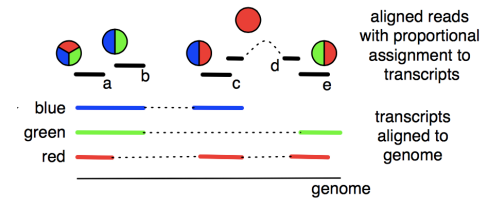
- Assemble transcripts directly from reads
- Allows transcriptome analyses of species without reference genomes

RNA-seq reads contain information about the abundance of different transcript isoforms

Normalization :

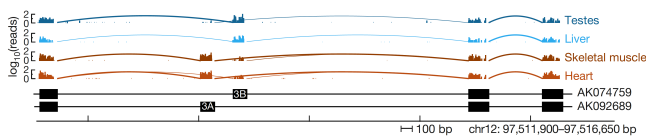
Internal: Reads or Fragments per kilobase of feature length per million mapped reads (RPKM or FPKM)

External: Reads relative to a standard "spike"



<http://arxiv.org/pdf/1104.3889v2.pdf>

Functional diversity in transcript isoforms



Examples of applications of RNA-seq

Characterizing transcriptome complexity

Alternative splicing

Differential expression analysis

Gene- and isoform-level expression comparisons

Novel RNA species

lncRNAs and eRNAs

Pervasive transcription

Translation

Ribosome profiling

Allele-specific expression

Measuring RNA half-lives and decay

Examining protein-RNA interactions (CLIP, RIP, &c.)

Effect of genetic variation on gene expression

Imprinting

RNA editing

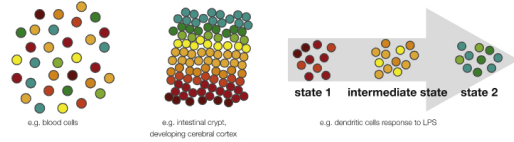
Novel events

Examining cell heterogeneity with scRNA-seq

Bulk RNA-seq averages over the RNA content of many cells masking differences.

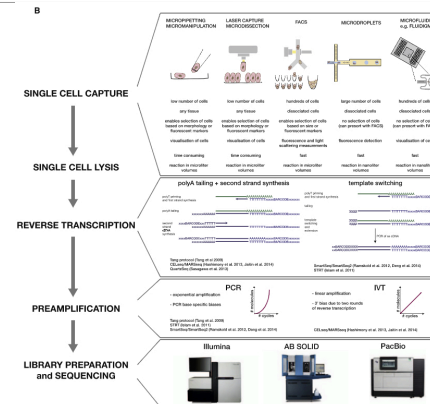
These differences can be revealed by sequencing the RNA from individual cells using single cell RNA-seq (scRNA-seq)

Analysis of RNA transcripts in individual cells can reveal rare cell populations and lineage trajectories.



Kolodziejczyk ... & Teichmann (2015). The technology and biology of single-cell RNA sequencing. Mol Cell

Examining cell heterogeneity with scRNA-seq



Kolodziejczyk ... & Teichmann (2015). The technology and biology of single-cell RNA sequencing. Mol Cell

Summary

- Genomics I: Deep sequencing gives us access to information on a genomic level.
- Genomics II: These approaches provide a diverse set of tools to study life at a genomic scale.
- * Sophisticated use of data from genomics requires an integrated understanding of the biological experiment, sample preparation and down stream computational analyses of the data.