

Bioinformatics: Genomics Part II

Applications of Sequencing Technology

Matt Simon
Dept. of Molecular Biophysics & Biochemistry
Chemical Biology Institute
January 22, 2018

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

Genomics lecture 1 summary

1. Isolation of sample.

e.g., Isolate DNA and shear.

2. Library preparation

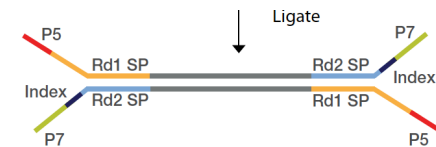
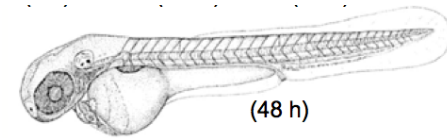
e.g., Clean up and ligate Y-adaptors.

3. Sequencing

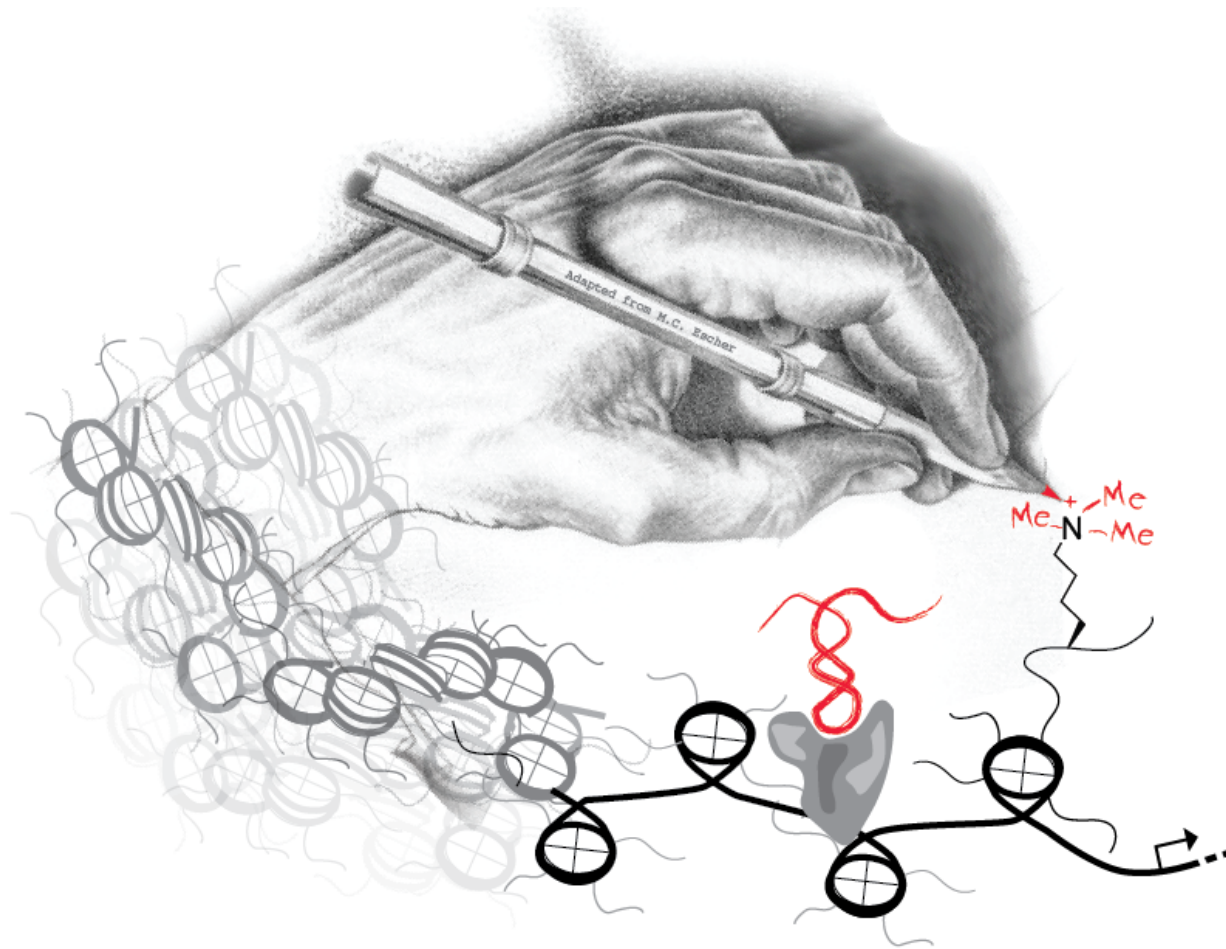
e.g., Illumina HiSeq

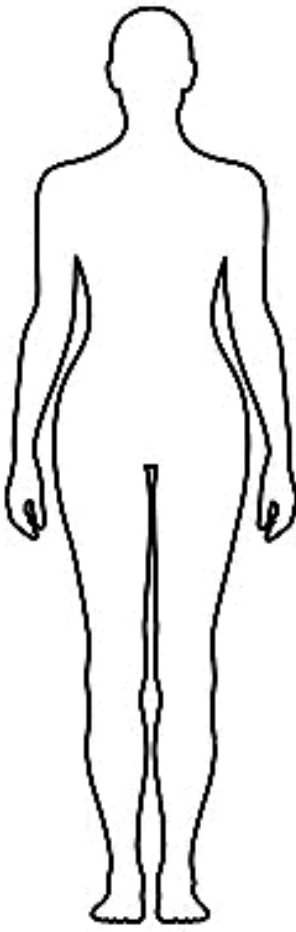
4. Analysis

e.g., Map to genome and interpret.

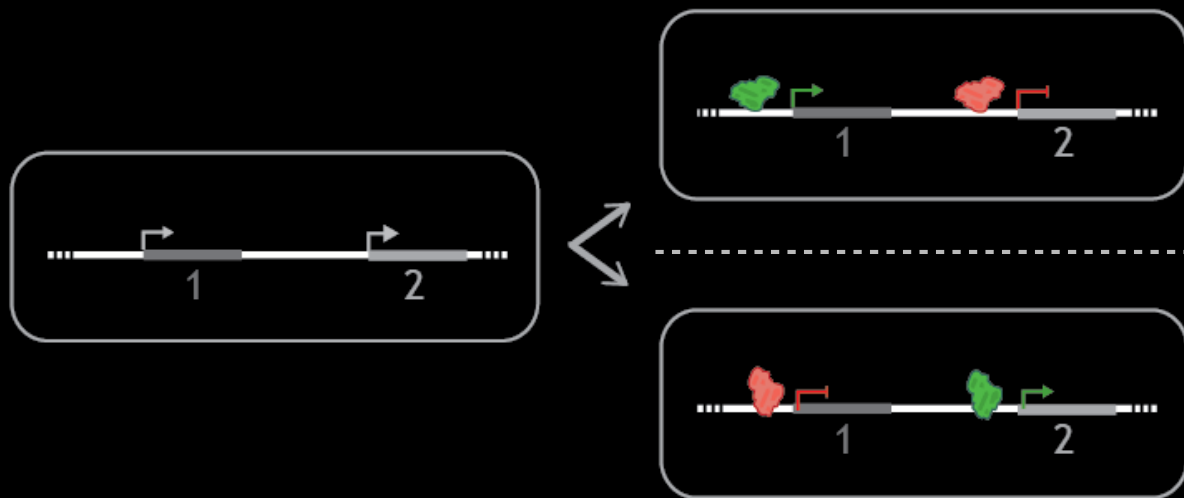


Part 1. How do cells annotate their genomes?

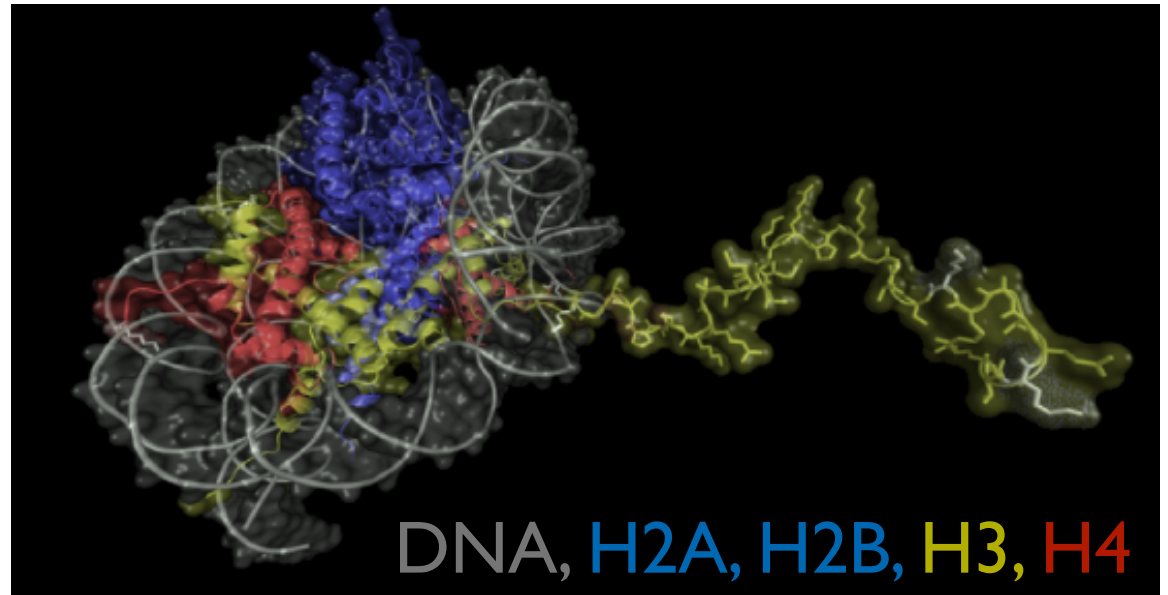
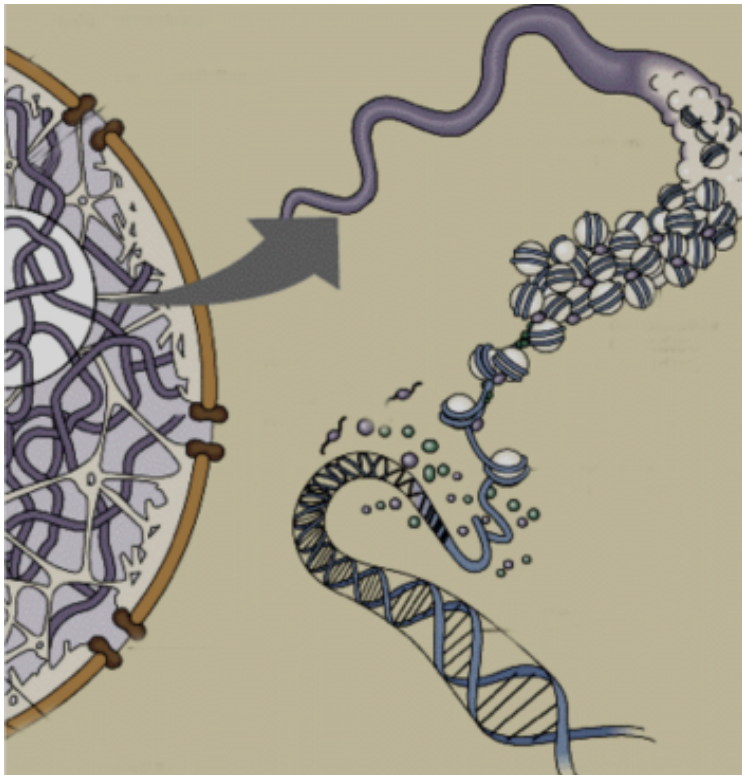




How is gene expression regulated and faithfully inherited?

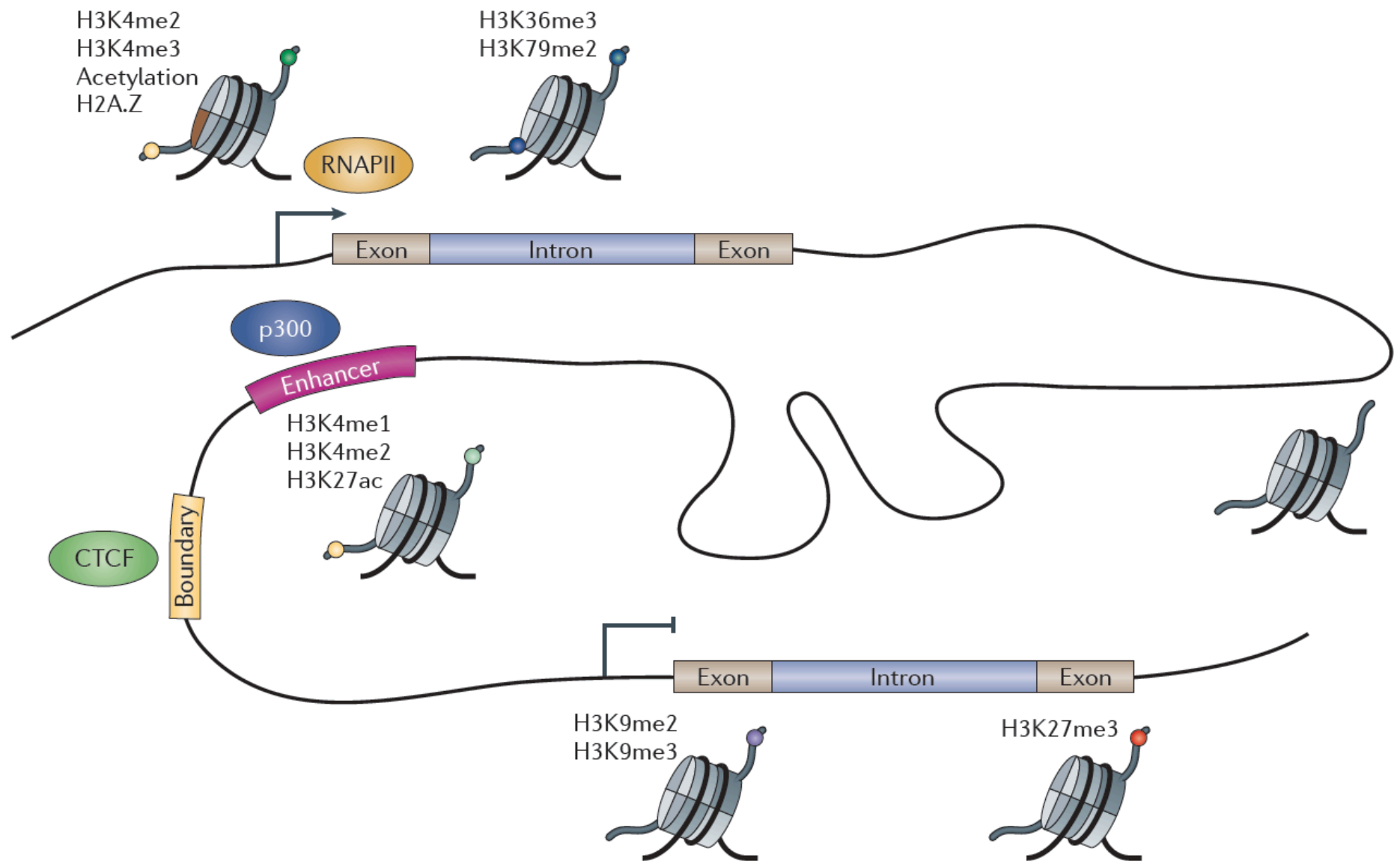


DNA in the cell is packaged into chromatin

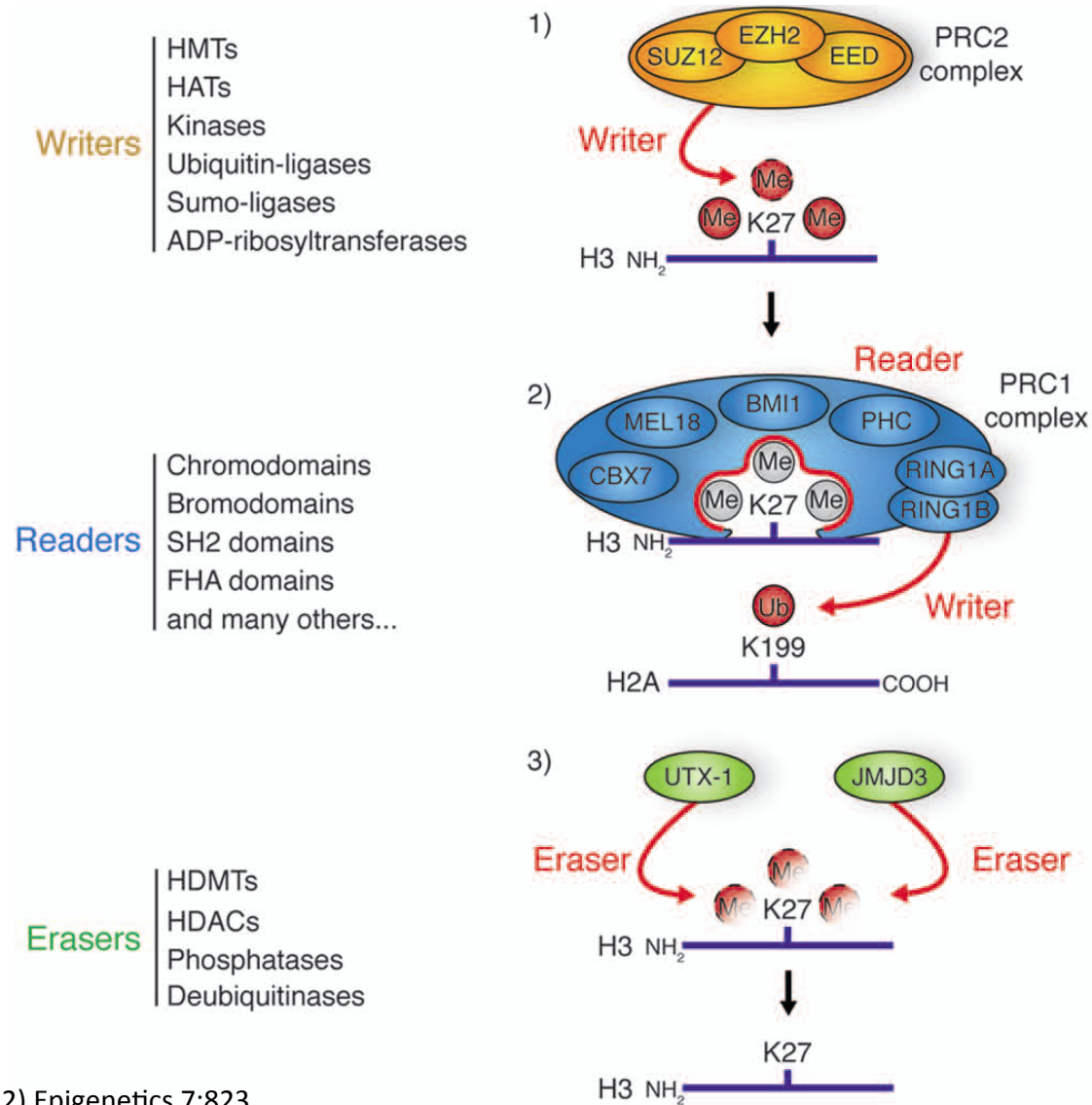


Modeled nucleosome based on Luger et al., *Nature* **1997** 389, 251.

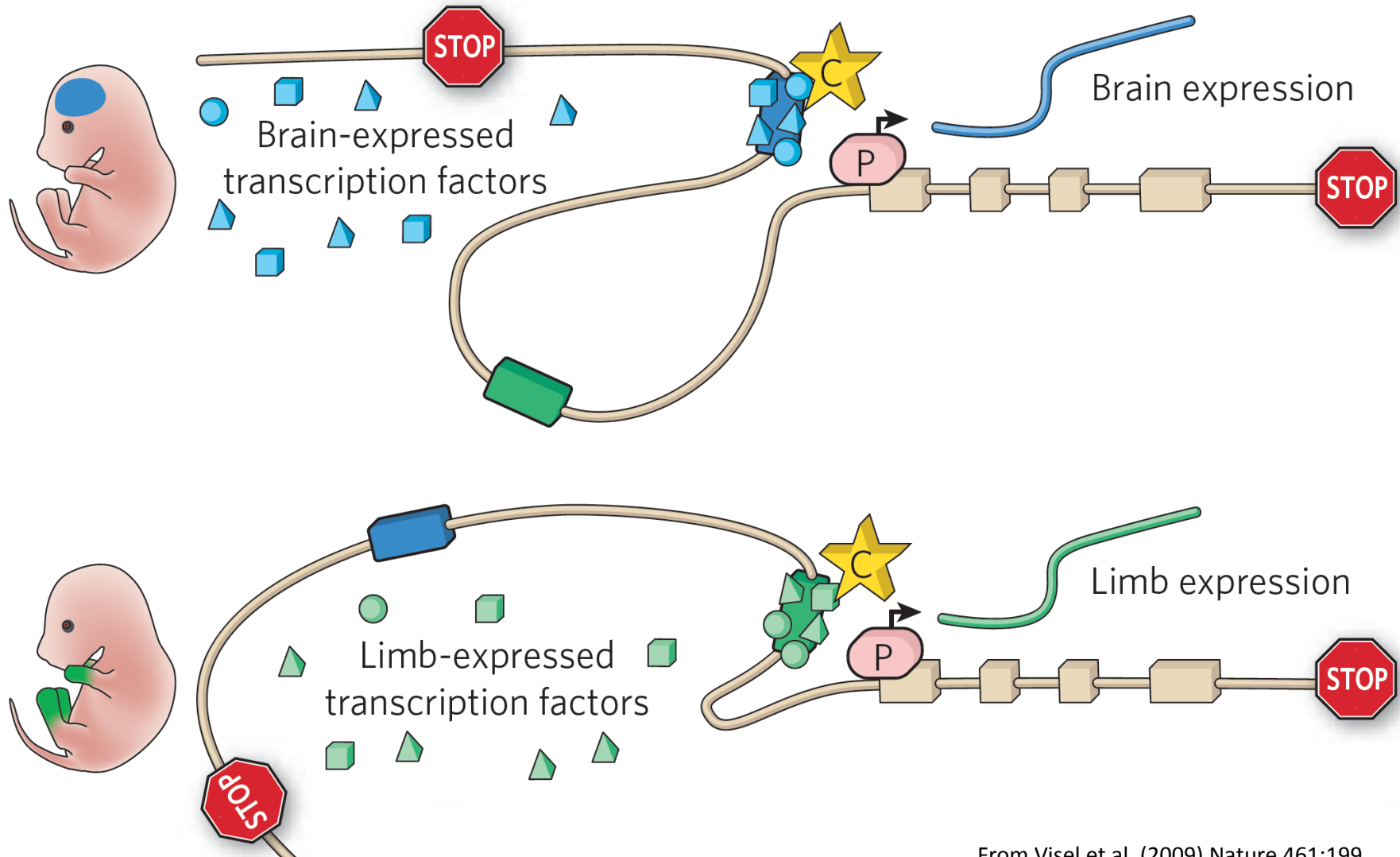
Chromatin modifications correlate with different genomic functions.



Installing, binding, and removing modifications



Regulation is temporally and specially controlled



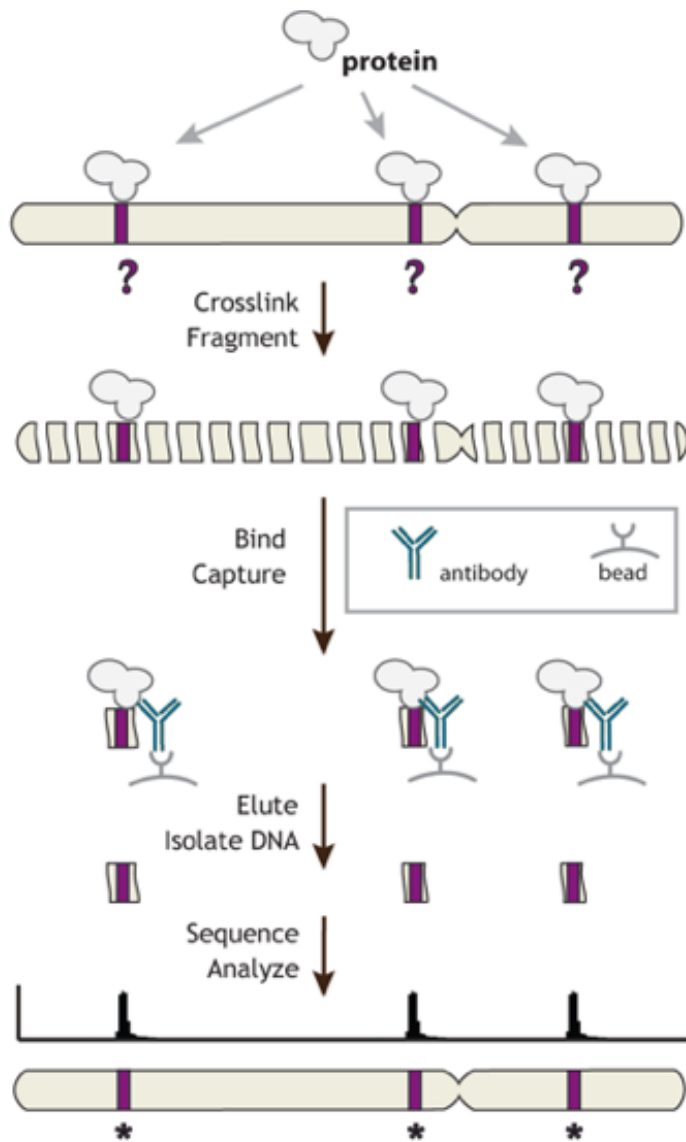
Using sequencing to annotate the genome

1. Where are the cis-acting regulatory elements in DNA?
 - A. DNase I hyper-sensitivity mapping (DNase-Seq).
 - B. FAIRE to map regulatory elements.
2. Where do transcription factors bind?
 - C. ChIP-seq of transcription factors (or in high res, ChIP-exo)
 - D. Nucleosome mapping (MNase-Seq).
3. Where are different histone modifications found?
 - E. ChIP-Seq of histone modifications.
 - F. ChIP-Seq of chromatin writers, readers and erasers.
4. Where is RNA polymerase transcribing?
 - G. ChIP-Seq of polymerase.
 - H. GRO-Seq, NET-Seq and TT-Seq to measure RNA in the polymerase active site..
5. How is the genome organized in 3D?
 - I. 4C/5C/Hi-C to measure chromatin conformation.

Using sequencing to annotate the genome

1. Where are the cis-acting regulatory elements in DNA?
 - A. DNase I hyper-sensitivity mapping (DNase-Seq).
 - B. FAIRE to map regulatory elements.
2. Where do transcription factors bind?
 - C. ChIP-seq of transcription factors (or in high res, ChIP-exo)
 - D. Nucleosome mapping (MNase-Seq).
3. Where are different histone modifications found?
 - E. ChIP-Seq of histone modifications.
 - F. ChIP-Seq of chromatin writers, readers and erasers.
4. Where is RNA polymerase transcribing?
 - G. ChIP-Seq of polymerase.
 - H. GRO-Seq, NET-Seq and TT-Seq to measure RNA in the polymerase active site..
5. How is the genome organized in 3D?
 - I. 4C/5C/Hi-C to measure chromatin conformation.

Localization of 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.

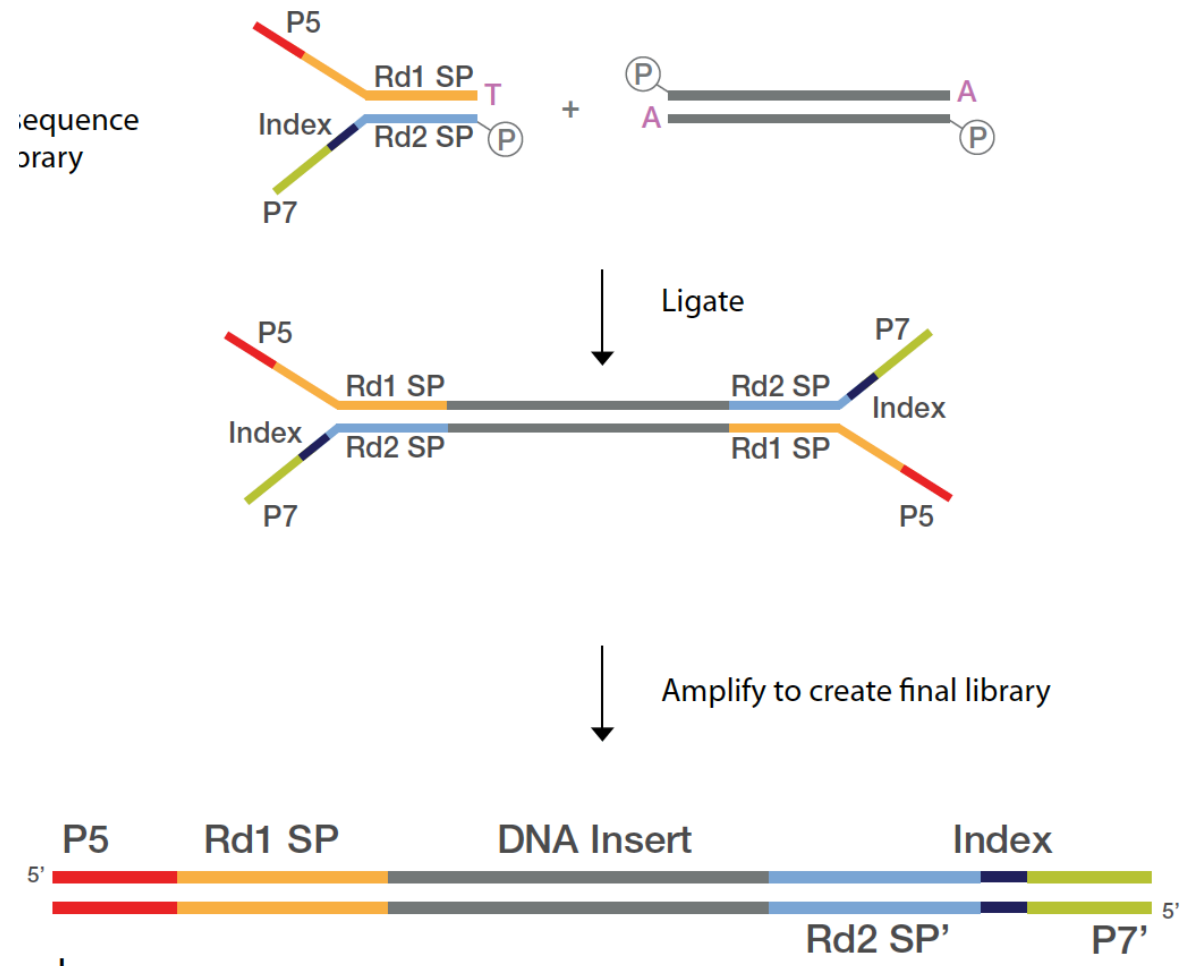
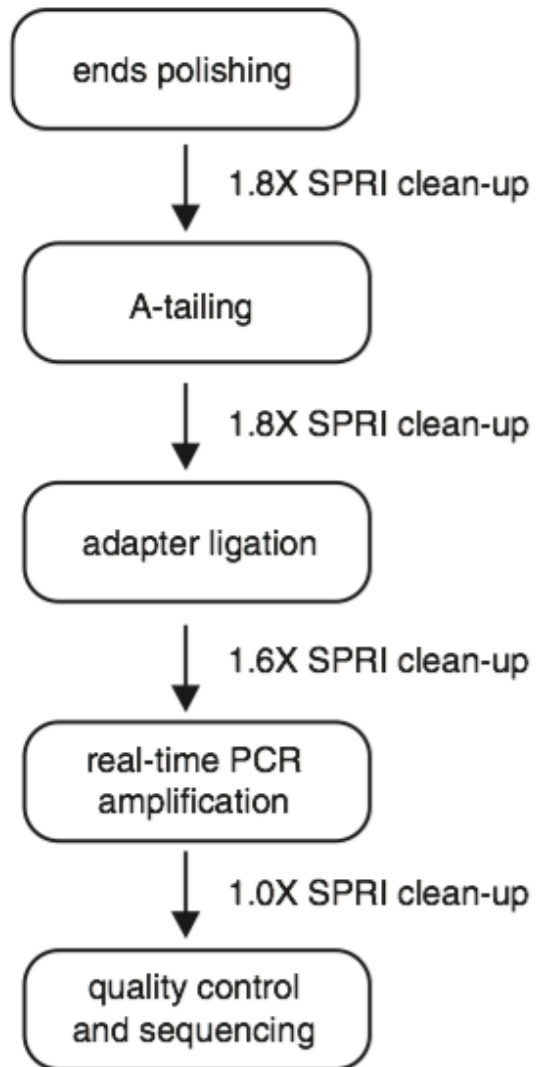
Shear size determines resolution.

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

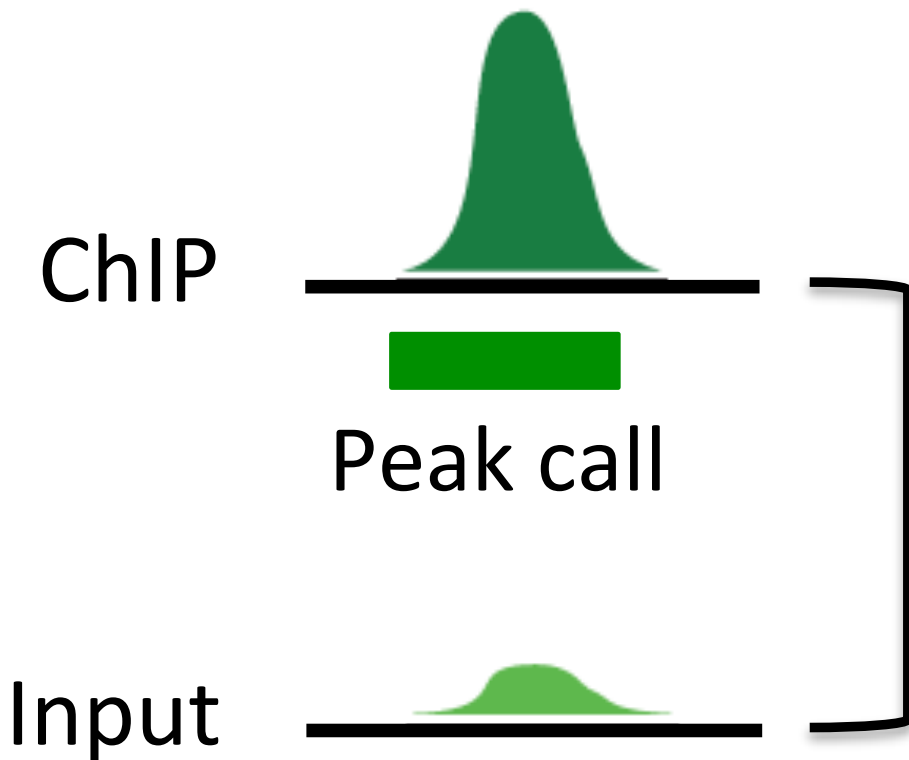
3. **Enrich** target using an antibody.

Enrichment is only as good as the antibody.

Preparing a Seq library using ChIP-enriched DNA.

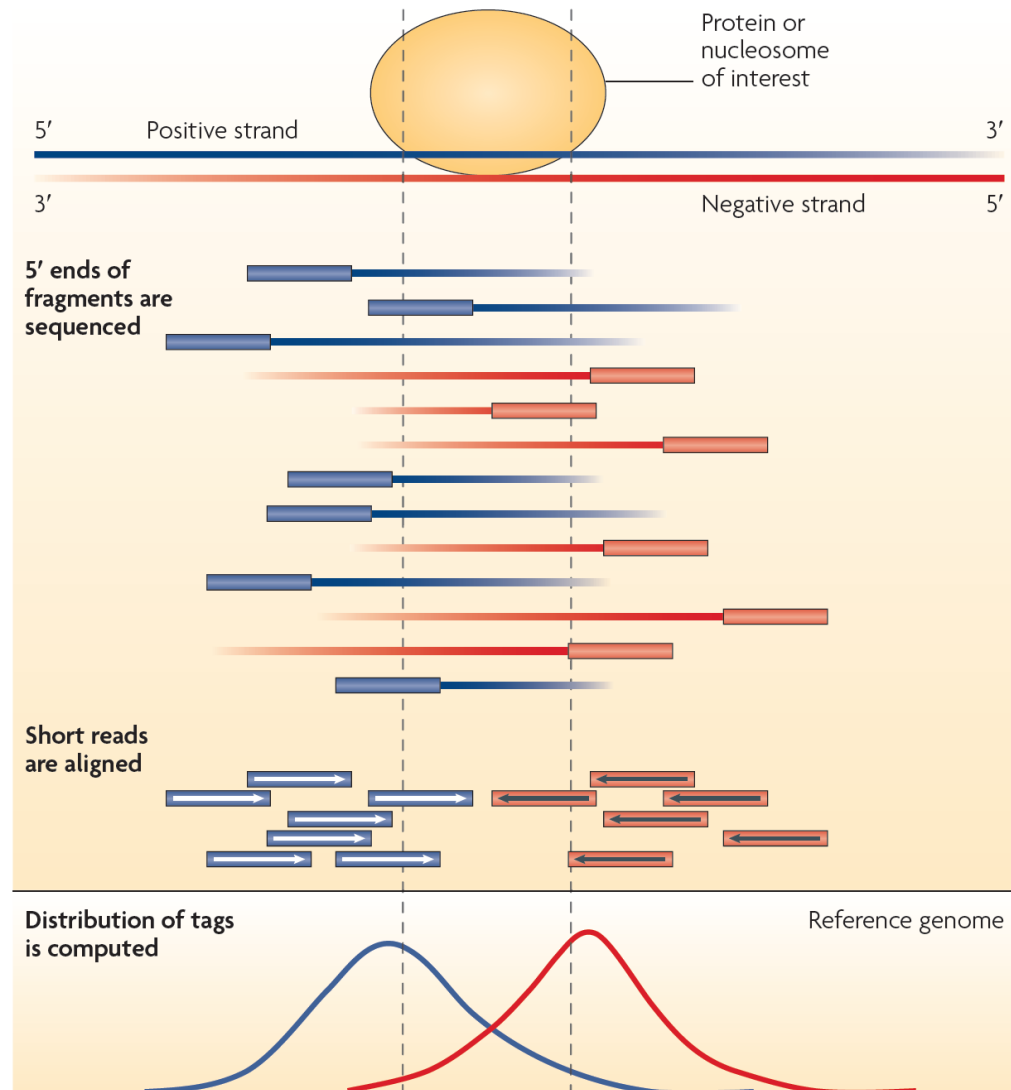


Determining sites of enrichment from ChIP-Seq

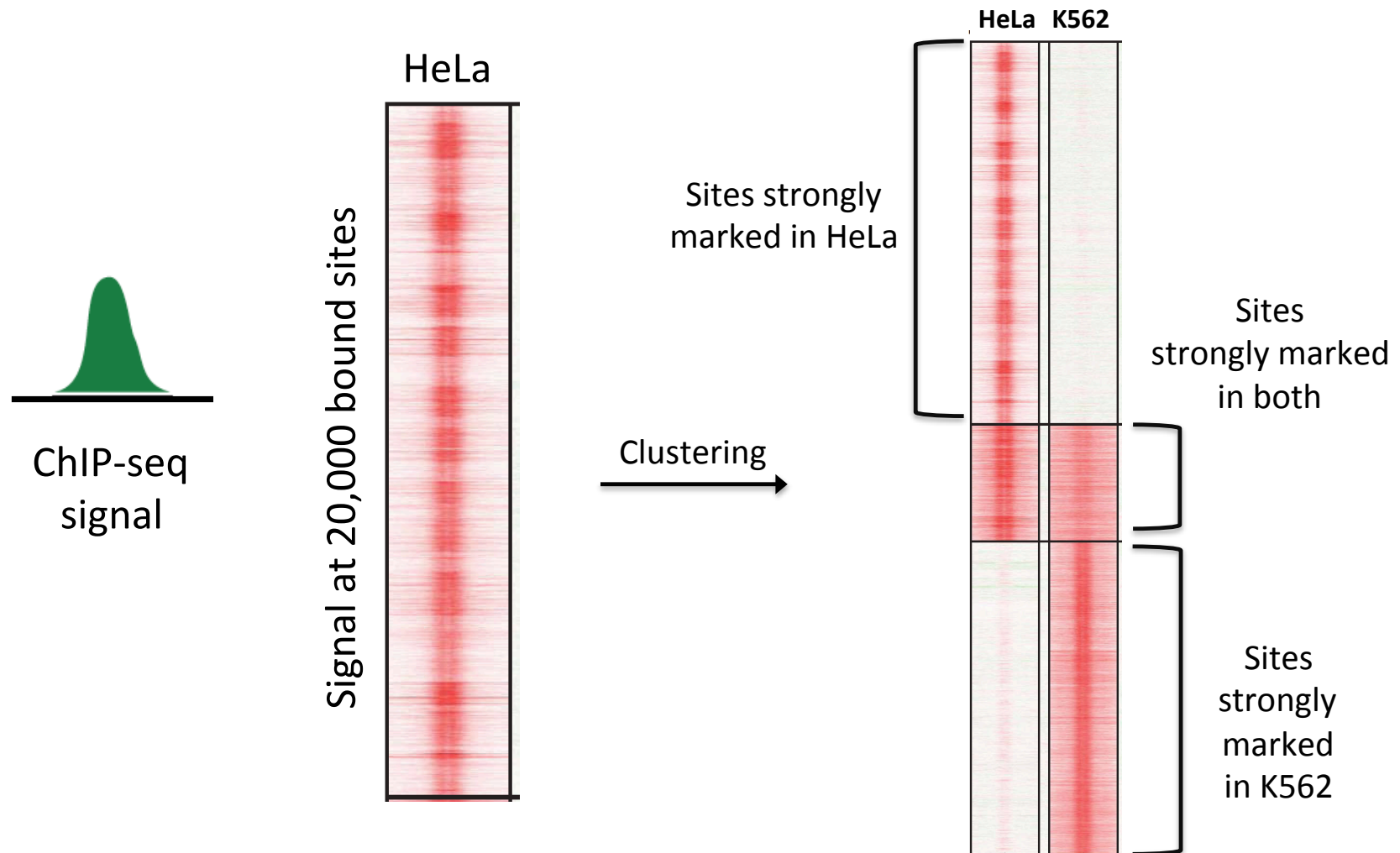


1. **Align** reads to the genome.
2. **Compare to input** to look for enrichment.
Input coverage is not even.
3. **Call peaks** to determine statistically significant sites of enrichment.

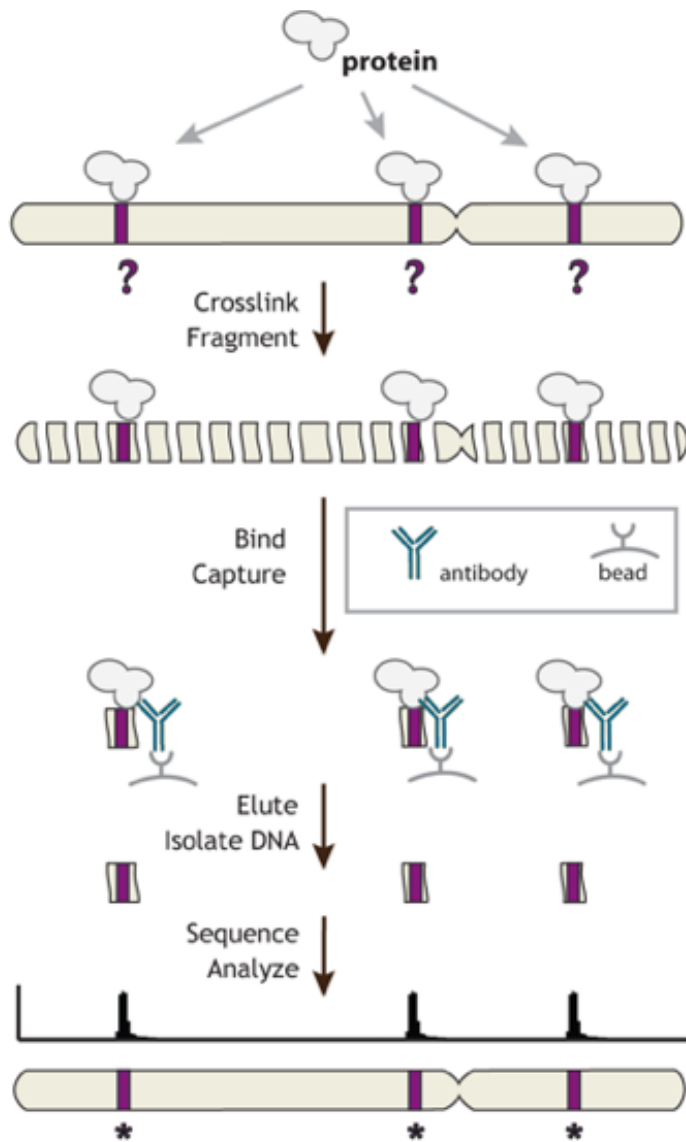
Avoiding artifacts using features in Seq data



ChIP-Seq signals reveal difference between cells

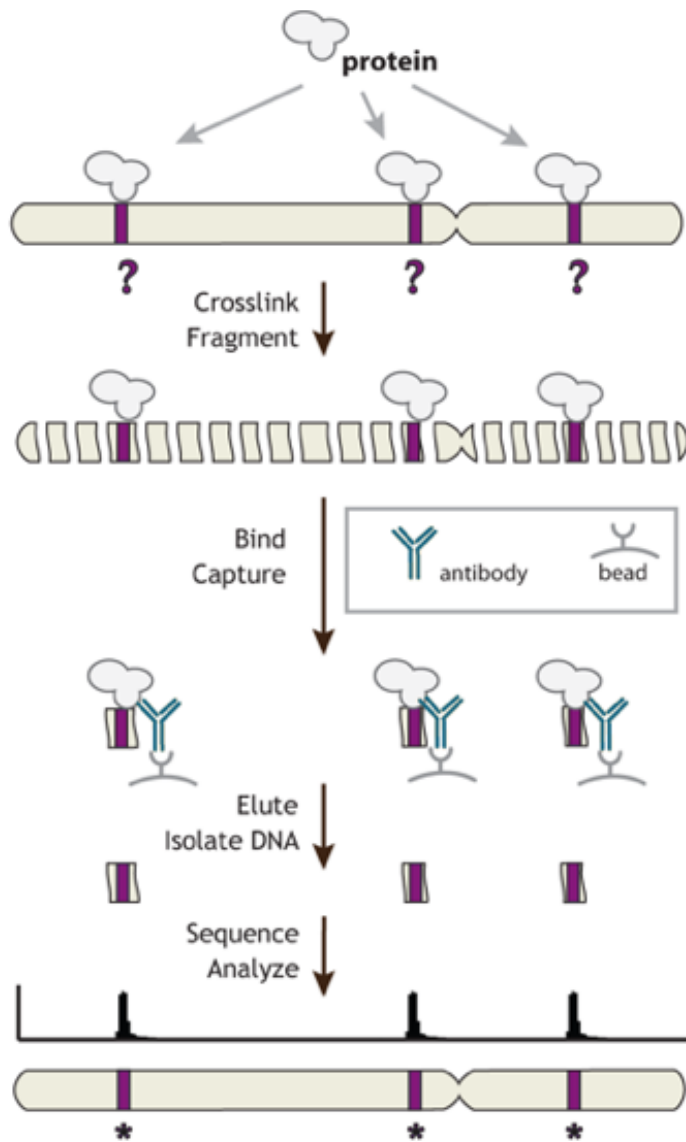


Limitations of ChIP-Seq



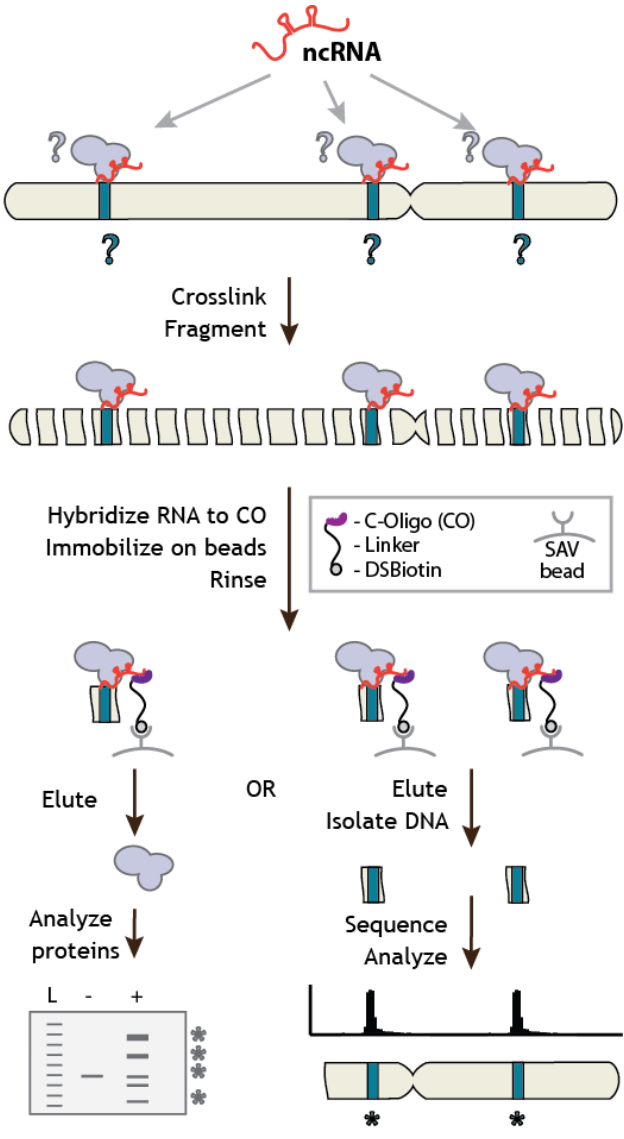
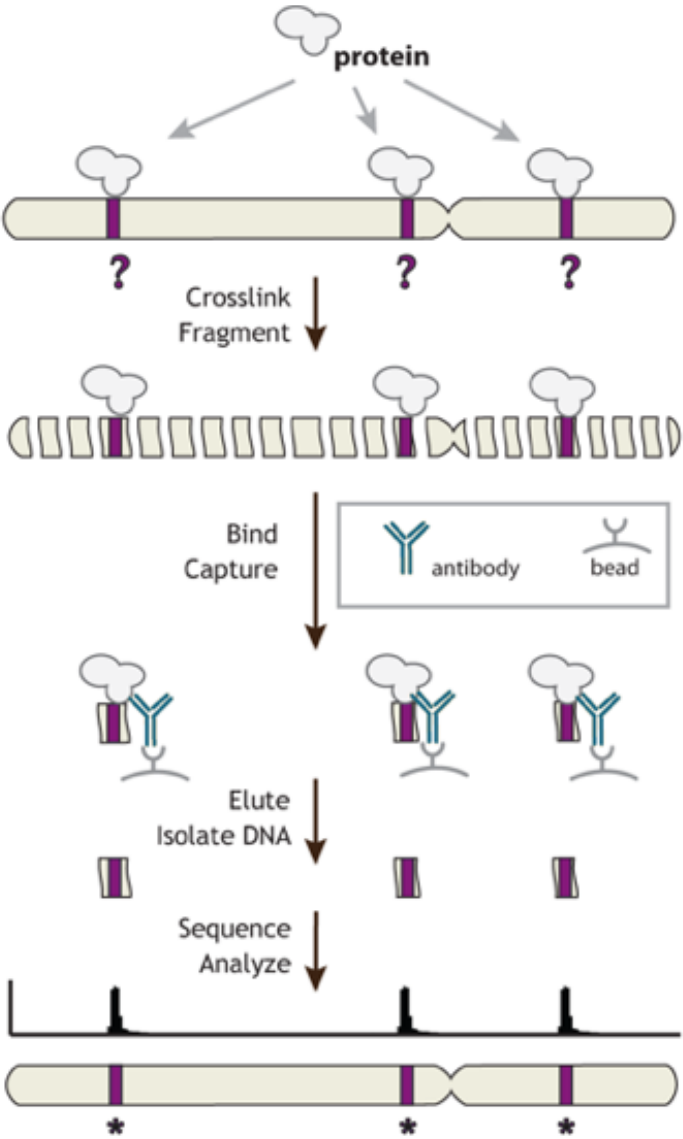
1. **Cross linking** efficiency is not necessarily uniform.
2. Enrichment is dependent on the **quality of antibody**.
e.g., Site and degree of histone modifications.
3. Enrichment is dependent on the **accessibility of the epitope**.
Comparing different sites to each other in the genome can be problematic.
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 **nucleosome turnover** and exchange.
3. Extension to **RNA factors**.

Extension to RNA factors: CHART, ChIRP and RAP

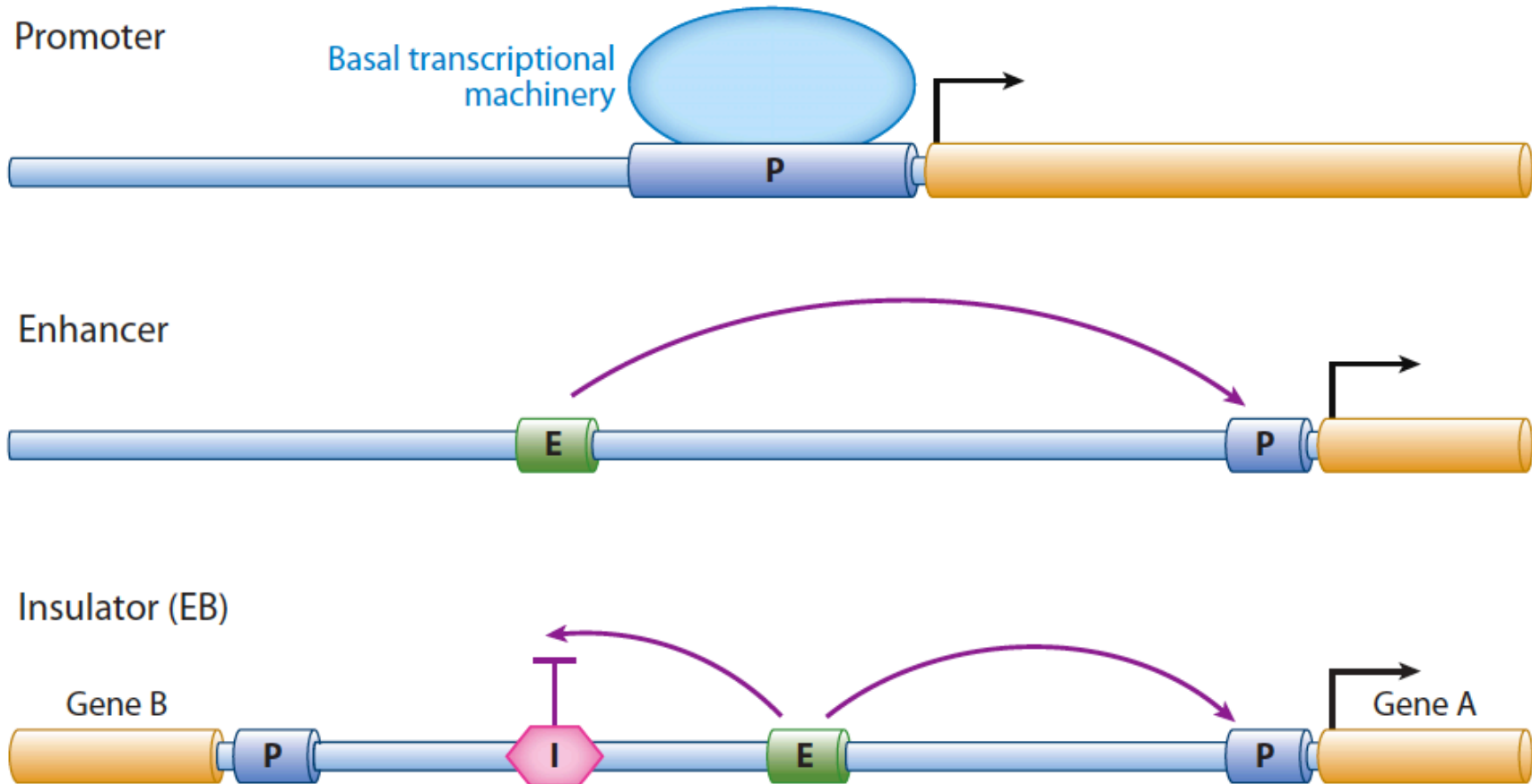


Using sequencing to annotate the genome

1. Where are the cis-acting regulatory elements in DNA?
 - A. DNase I hyper-sensitivity mapping (**DNase-Seq**).
 - B. **FAIRE** to map regulatory elements.
 - C. **ATAC-Seq** to map regulatory elements.
2. Where do transcription factors bind?
 - D. **ChIP-seq** of transcription factors (or in high res, ChIP-exo)
 - E. Nucleosome mapping (**MNase-Seq**).
3. Where are different histone modifications found?
 - F. **ChIP-Seq** of histone modifications.
 - G. **ChIP-Seq** of chromatin writers, readers and erasers.
4. Where is RNA polymerase transcribing?
 - H. **ChIP-Seq** of polymerase.
 - I. **GRO-Seq** and **NET-Seq** to measure RNA in the polymerase active site..
5. How is the genome organized in 3D?
 - J. **4C/5C/Hi-C** to measure chromatin conformation.

Targeted approaches v **Global** approaches

How do we identify regulatory elements in the genome?

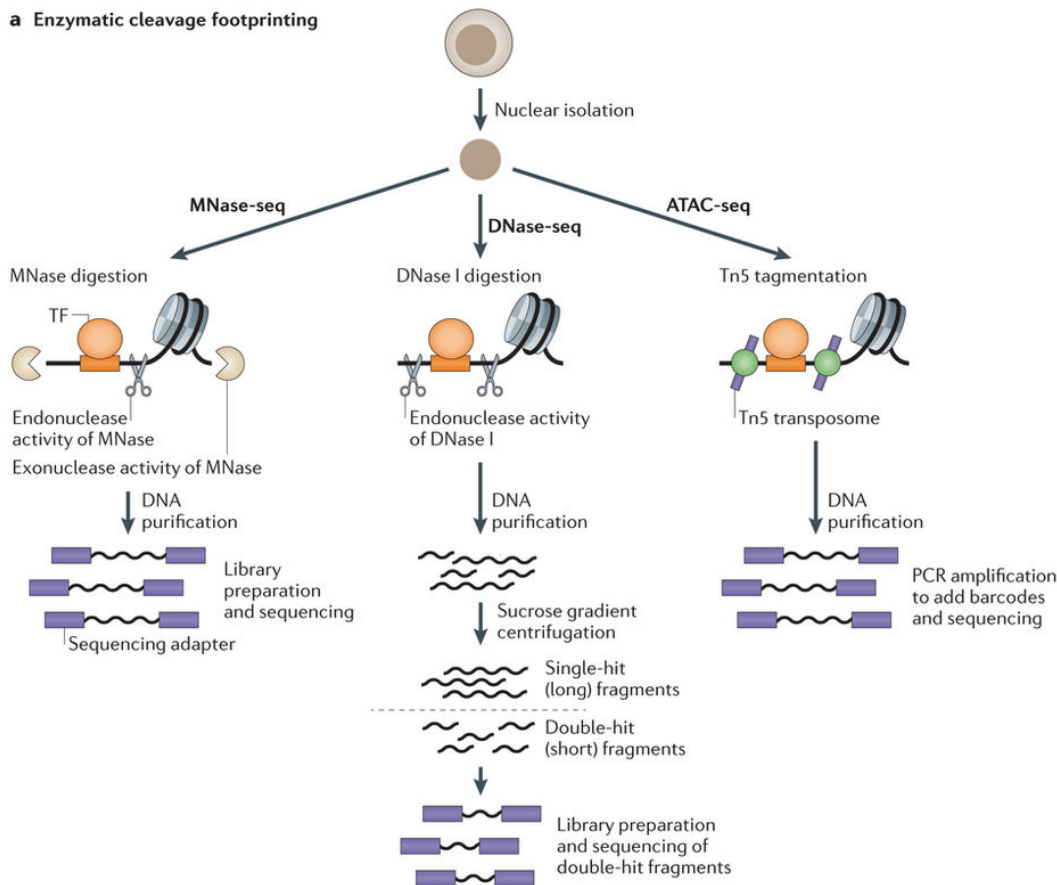


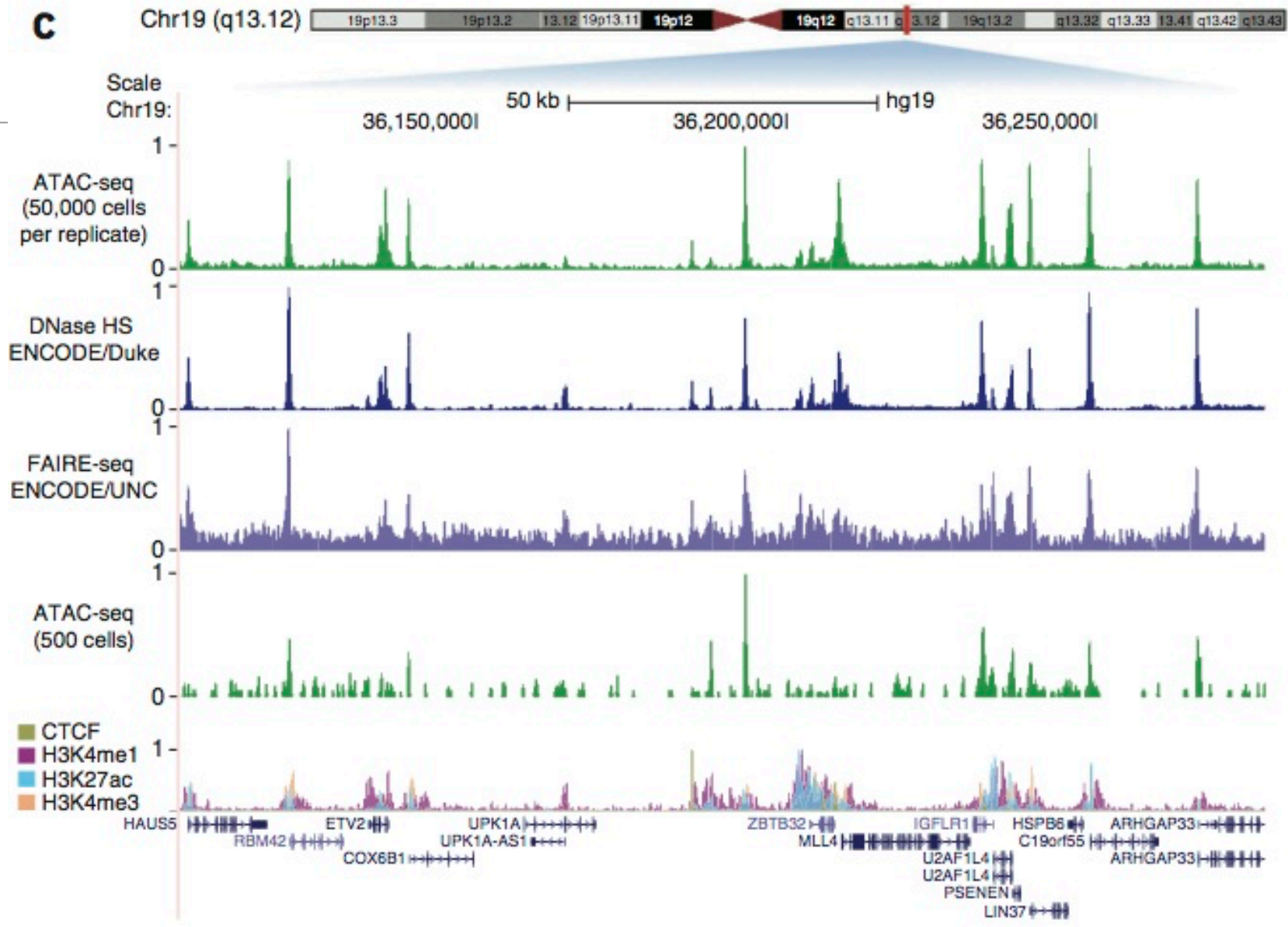
Using differences in biochemical properties of regulatory elements to identify them by Seq

1. **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

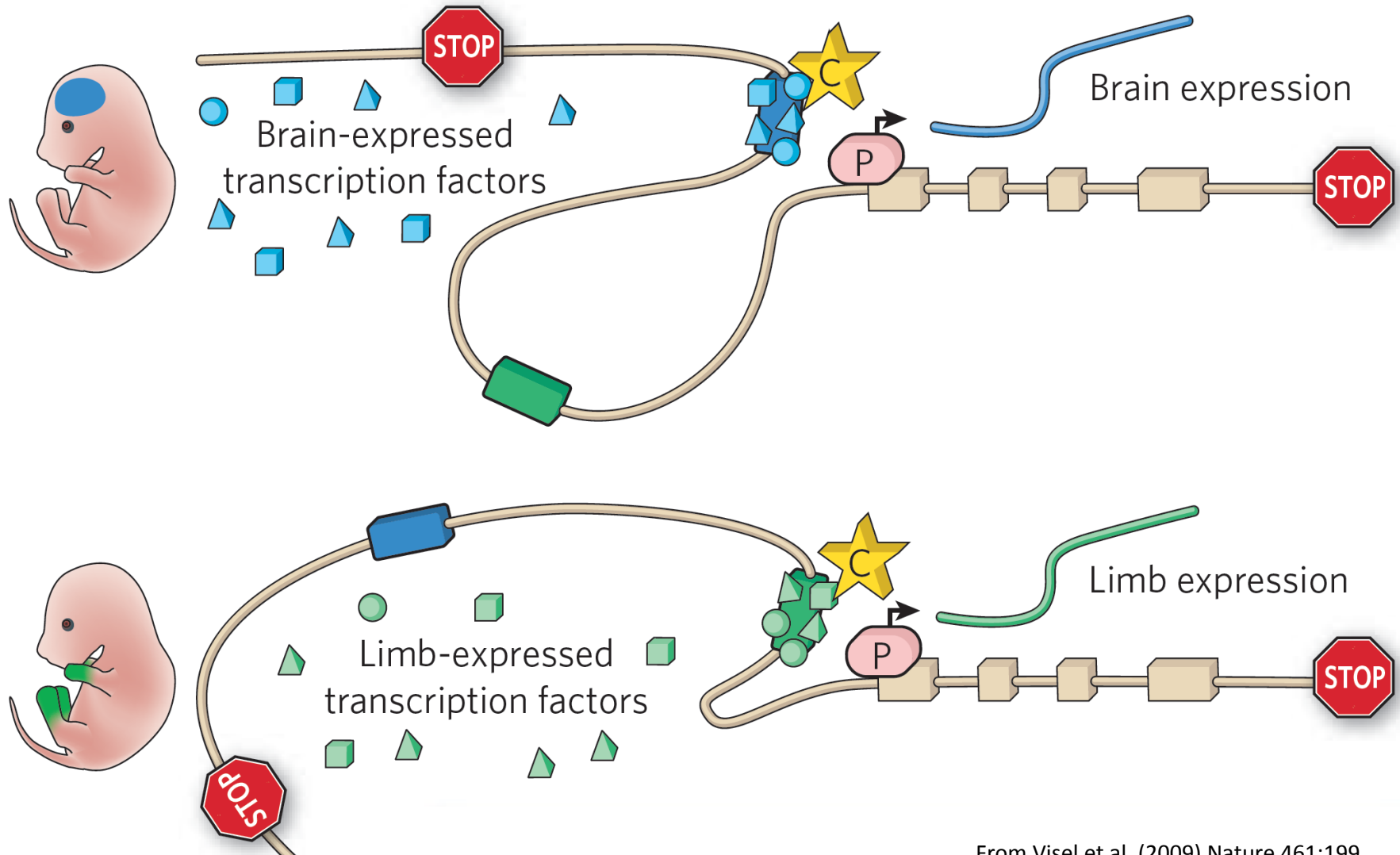
- ATAC-seq (shown)
- MNase-Seq (shown).
- DNase-Seq (shown).
- FAIRE-Seq (not shown).





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*

The 3D organization of the genome is important



Ligation-based methods to study 3D conformation



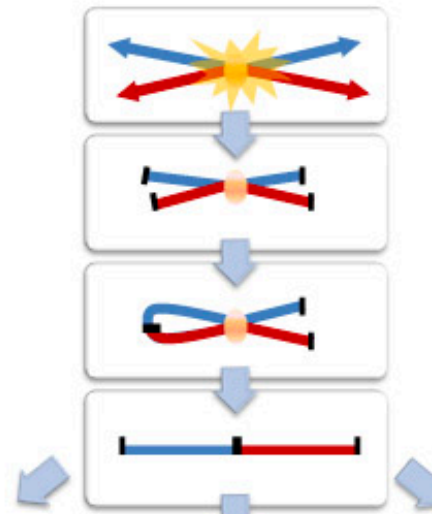
Image: David Goodsell

**Cross-link
DNA**

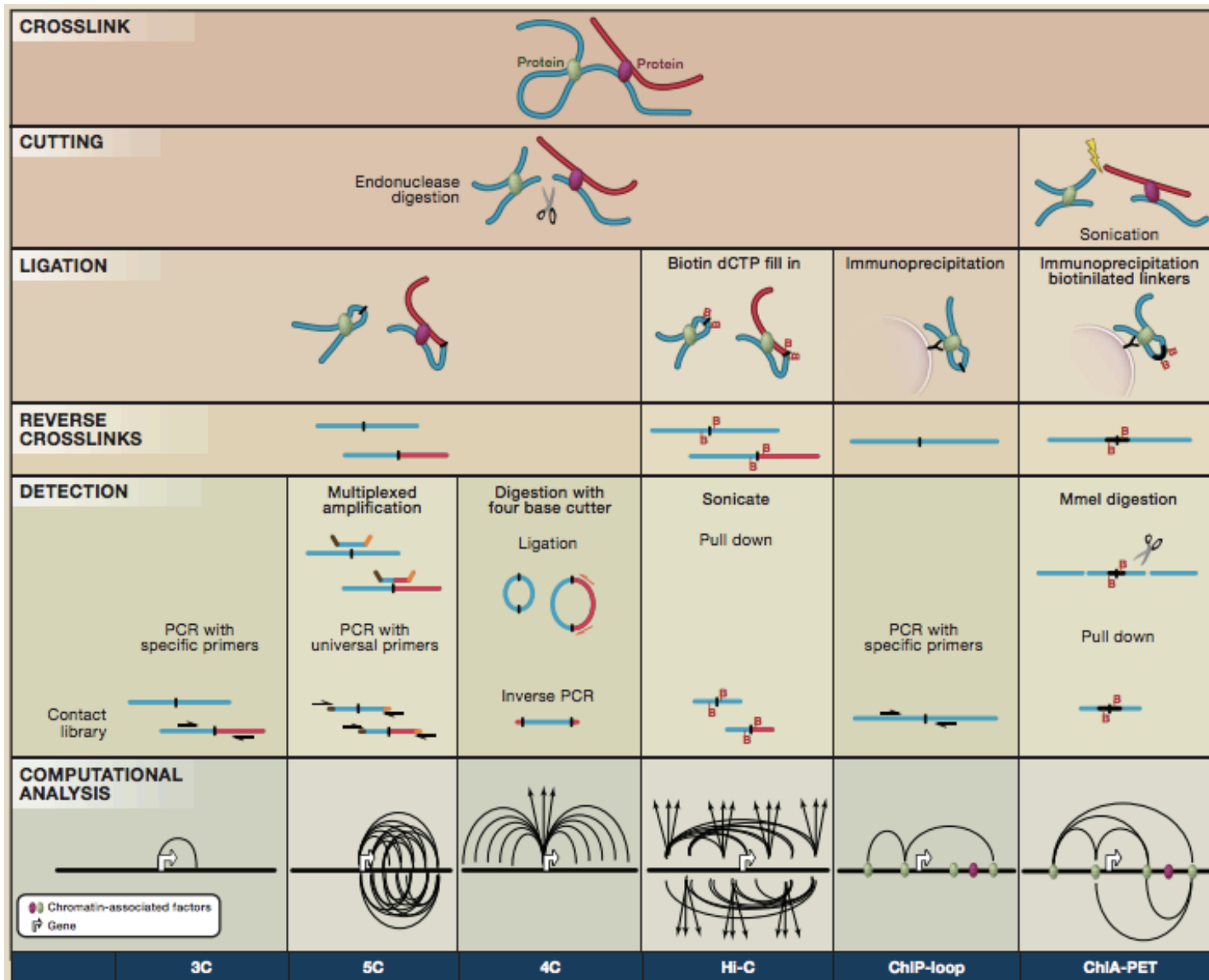
**Restriction
Digest**

**Intramolecular
Ligation**

**Reverse
Cross-link**

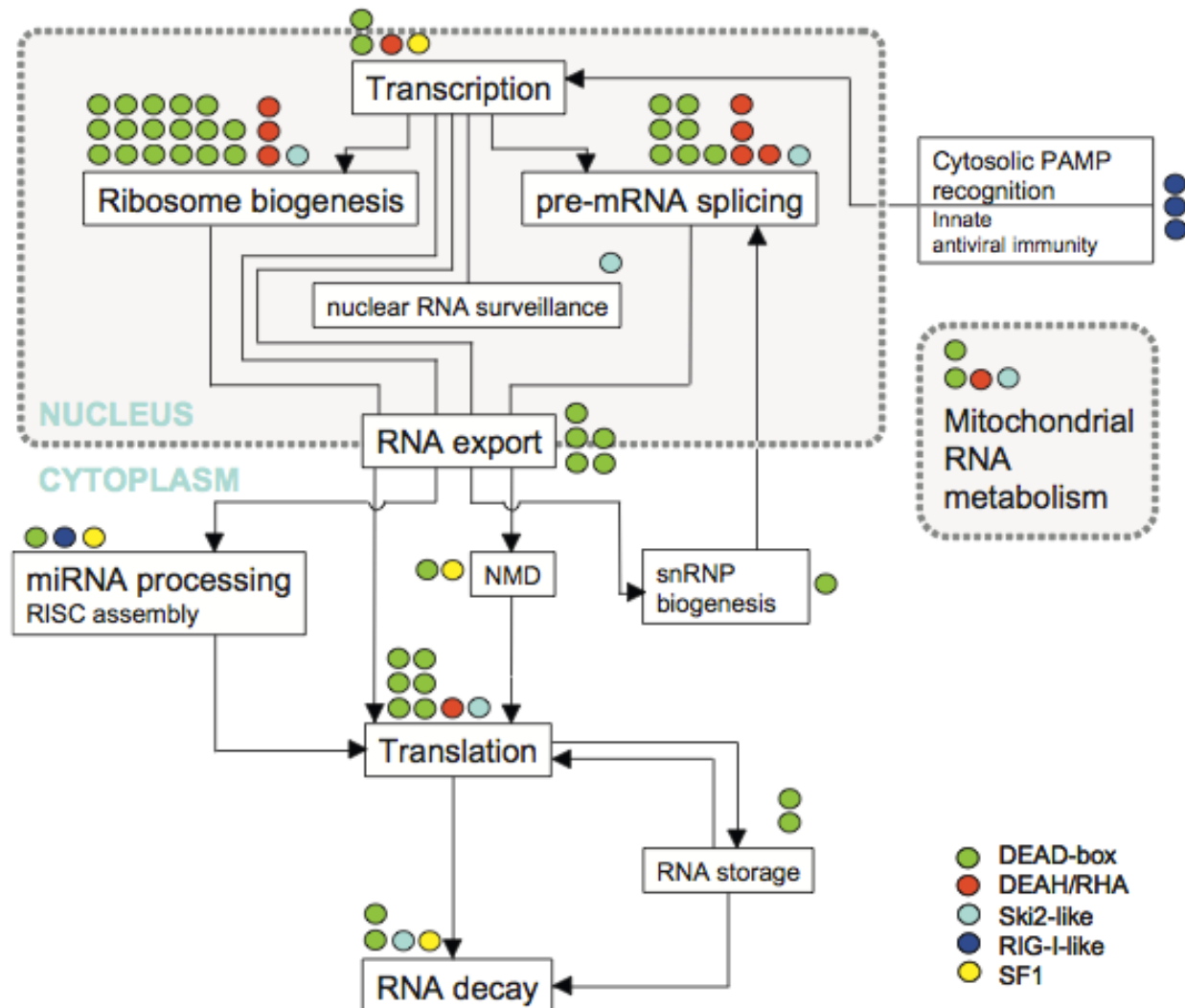


Many techniques to analyze chromatin conformation



Hakim & Misteli,
Cell (2012)

Gene expression is also controlled at the level of RNA



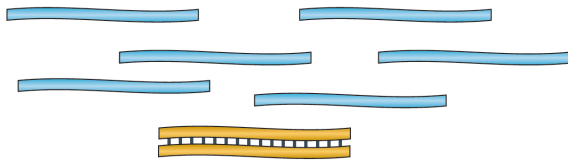
Part 2: RNA-Seq and applications of RNA-Seq

Using RNA-Seq to examine RNA

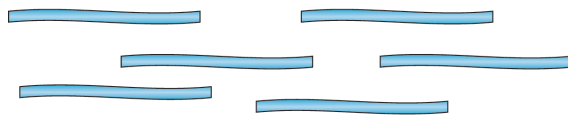
- Technical methodology
- Read mapping and normalization
- Estimating isoform-level gene expression
- De novo transcript reconstruction
- Sensitivity and sequencing depth
- Differential expression analysis

RNA-Seq workflow

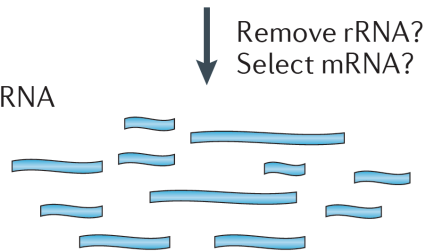
① mRNA or total RNA



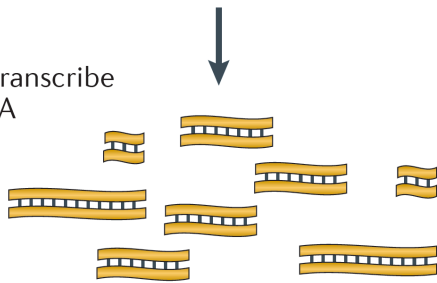
② Remove contaminant DNA



③ Fragment RNA

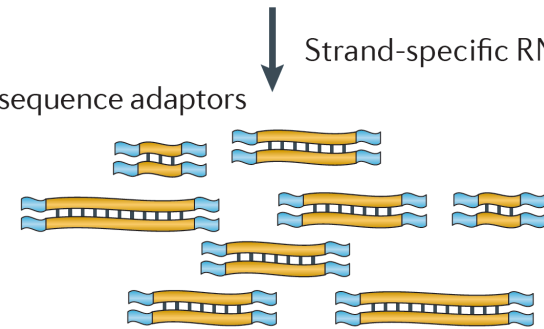


④ Reverse transcribe into cDNA

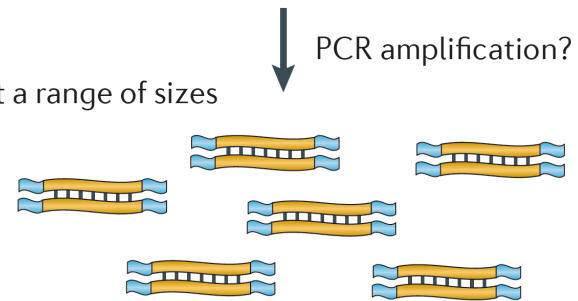


Strand-specific RNA-seq?

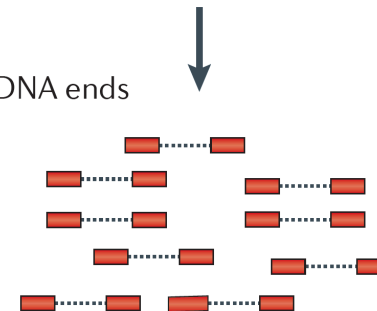
⑤ Ligate sequence adaptors



⑥ Select a range of sizes



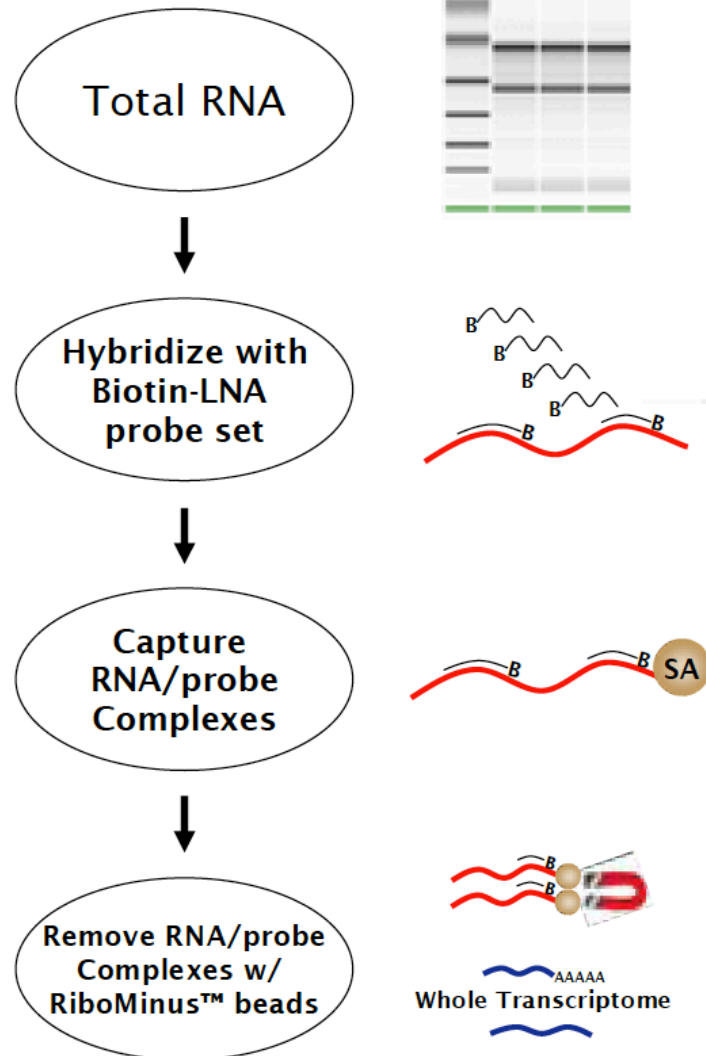
⑦ Sequence cDNA ends



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).

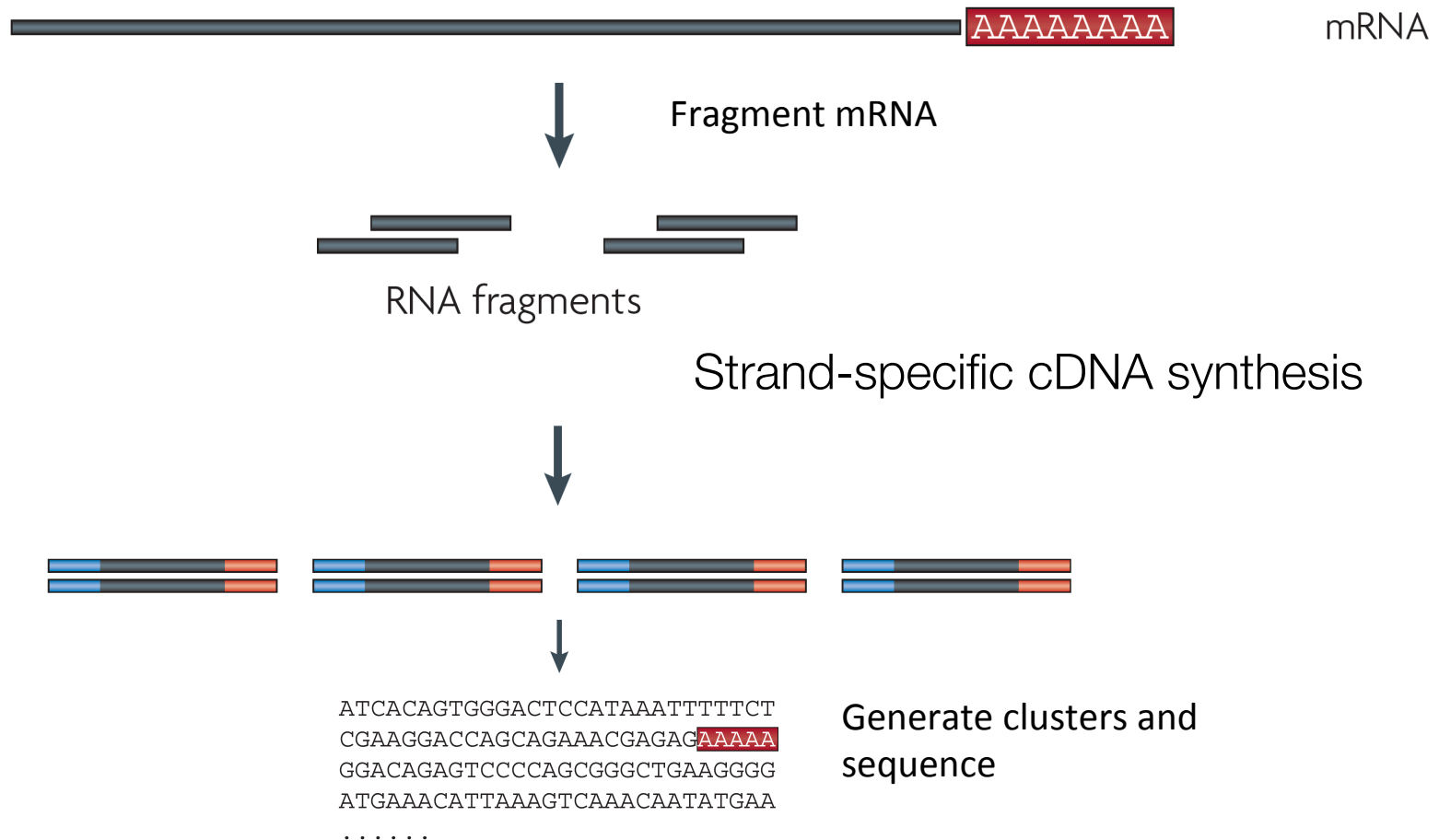
Ribosomal RNA will dominate the sequenced reads unless removed



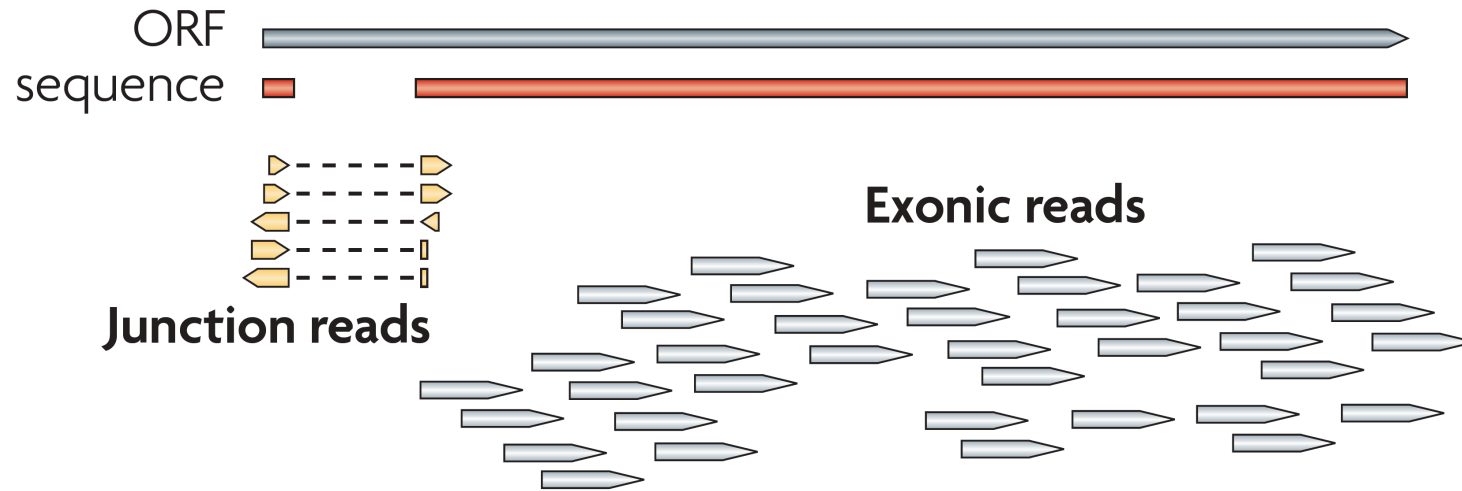
Illumina RNA-seq workflow

Capture poly-A RNA with poly-T oligo attached beads (100 ng total) (2x)

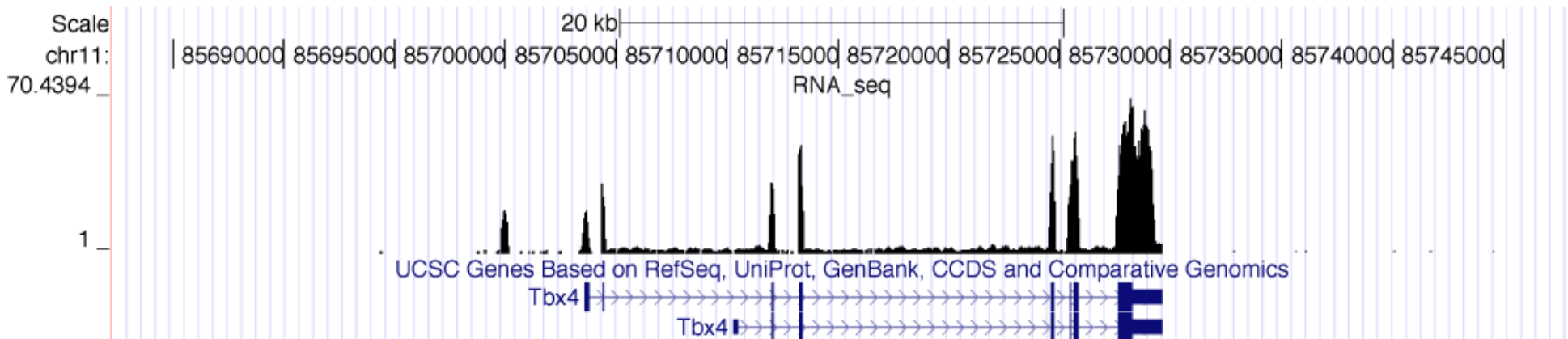
- RNA quality must be high – degradation produces 3' bias
- Non-poly-A RNAs are not recovered



RNA-Seq reads map mostly to exons



Martin and Wang *Nat Rev Genet* 12:671 (2011)



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

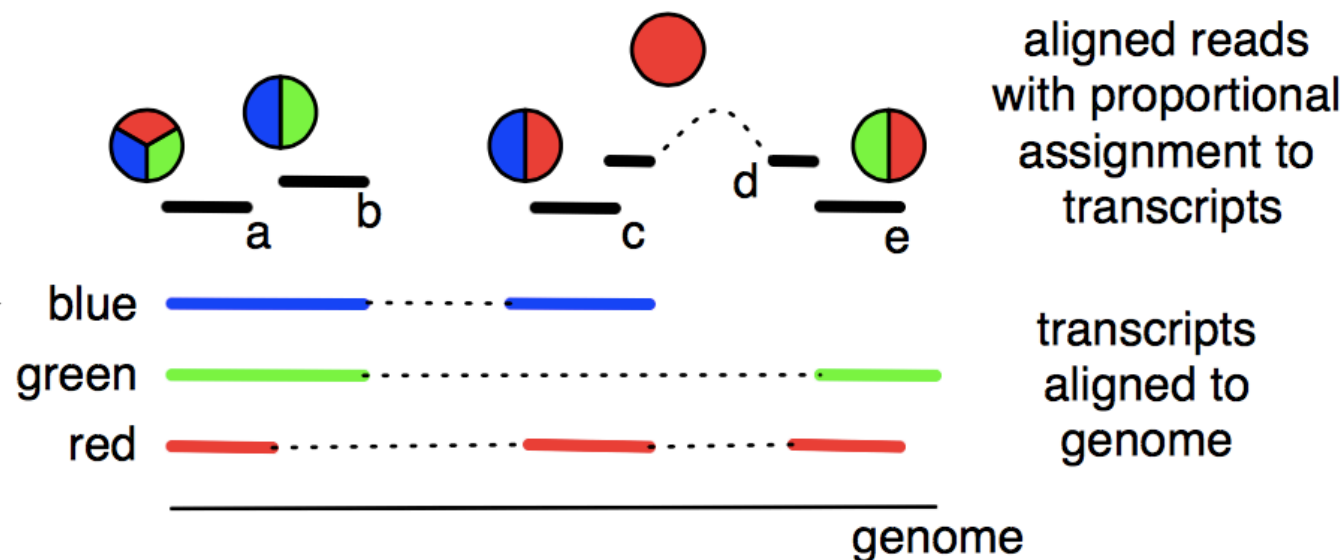
Allows transcriptome analyses of species without reference genomes

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

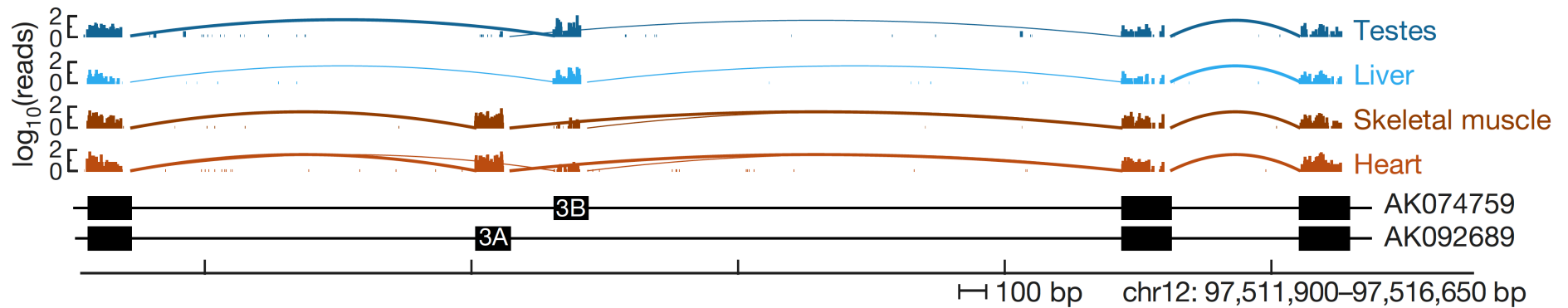
Normalization :

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

External: Reads relative to a standard “spike”



There is a lot of functional diversity in transcript isoforms



Alternative transcript events	Total events (×10 ³)	Number detected (×10 ³)	Both isoforms detected	Number tissue-regulated	% Tissue-regulated (observed)	% Tissue-regulated (estimated)
Skipped exon	37	35	10,436	6,822	65	72
Retained intron	1	1	167	96	57	71
Alternative 5' splice site (A5SS)	15	15	2,168	1,386	64	72
Alternative 3' splice site (A3SS)	17	16	4,181	2,655	64	74
Mutually exclusive exon (MXE)	4	4	167	95	57	66
Alternative first exon (AFE)	14	13	10,281	5,311	52	63
Alternative last exon (ALE)	9	8	5,246	2,491	47	52
Tandem 3' UTRs	7	7	5,136	3,801	74	80
Total	105	100	37,782	22,657	60	68

Constitutive exon or region
 Body read
 Junction read
 pA Polyadenylation site
 Alternative exon or extension
Inclusive/extended isoform
Exclusive isoform
Both 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

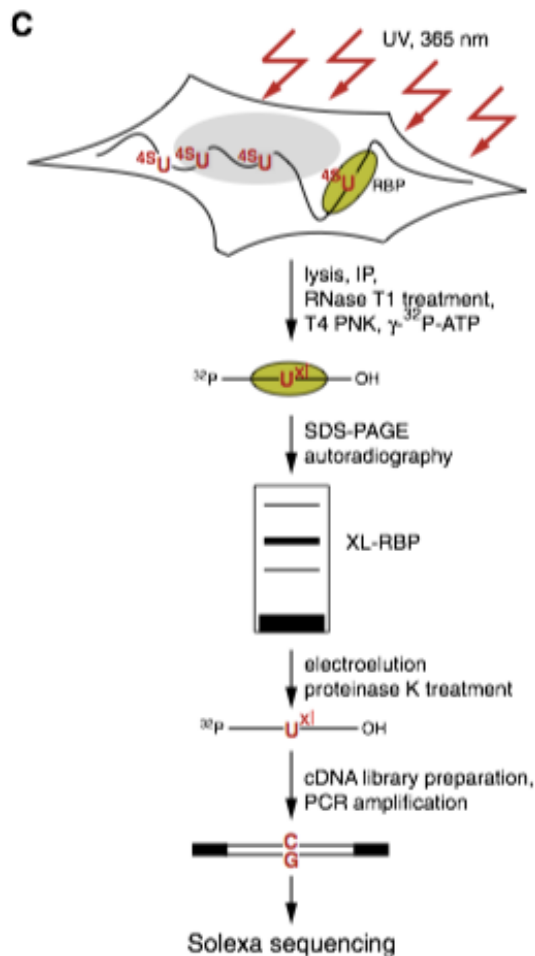
- Imprinting

- RNA editing

- Novel events

RNA-seq to examine protein-RNA interactions

PAR-CLIP Photoactivatable-Ribonucleoside- Enhanced Crosslinking and Immunoprecipitation



There are several methods to look at protein-RNA interactions using RNA-Seq such as RIP, CLIP and similar protocols.

3'UTR of ELF1

Sequence	# reads	error
AAATGTTTTAGATTACTTTTTCAACTGTAAATAATGTACATTTAATGTCACAAGAAAA	501	1
-----ATTACTTTTTCAACTGTAAACAATGTACATTT-----	239	1
-----ATTACTTTTTCAACTGTAAATAATGTACACTT-----	113	0
-----ACTTTTTCAACTGTAAACAATGTACATTTAAT-----	82	1
-----ATTACTTTTTCAACTGTAAATAATGTACATCT-----	67	1

3'UTR of HES1

Sequence	# reads	error
GTGACTGACCATGCACTATATTTGTATATATTTTATATGTTTCATATGCAATTCGCGCTT	527	1
-----CACTATATTTGTATACAATTTTATATG-----	130	1
-----CACTATATTTGTATACAATTTTATATG-----	48	1
-----ACTATATTTGTATACAATTTTATATG-----	40	1
-----CACTATATTTGTATATATTTTATATGTTTCA-----	22	1

1. Hafner, M. *et al.* Transcriptome-wide Identification of RNA-Binding Protein and MicroRNA Target Sites by PAR-CLIP. *CELL* **141**, 129–141 (2010).

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.