# Bioinformatics: Genomics Part II

## Applications of Sequencing Technology

Matt Simon
Dept. of Molecular Biophysics & Biochemistry
Chemical Biology Institute
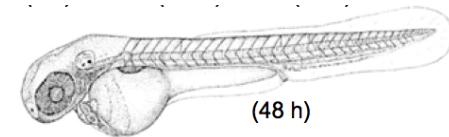January 16, 2015

# 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
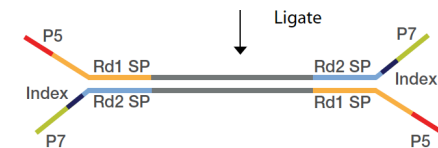
# Review of how a sequencing experiment works

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.

# Q. How many cycles of PCR are used in flow cell generation?
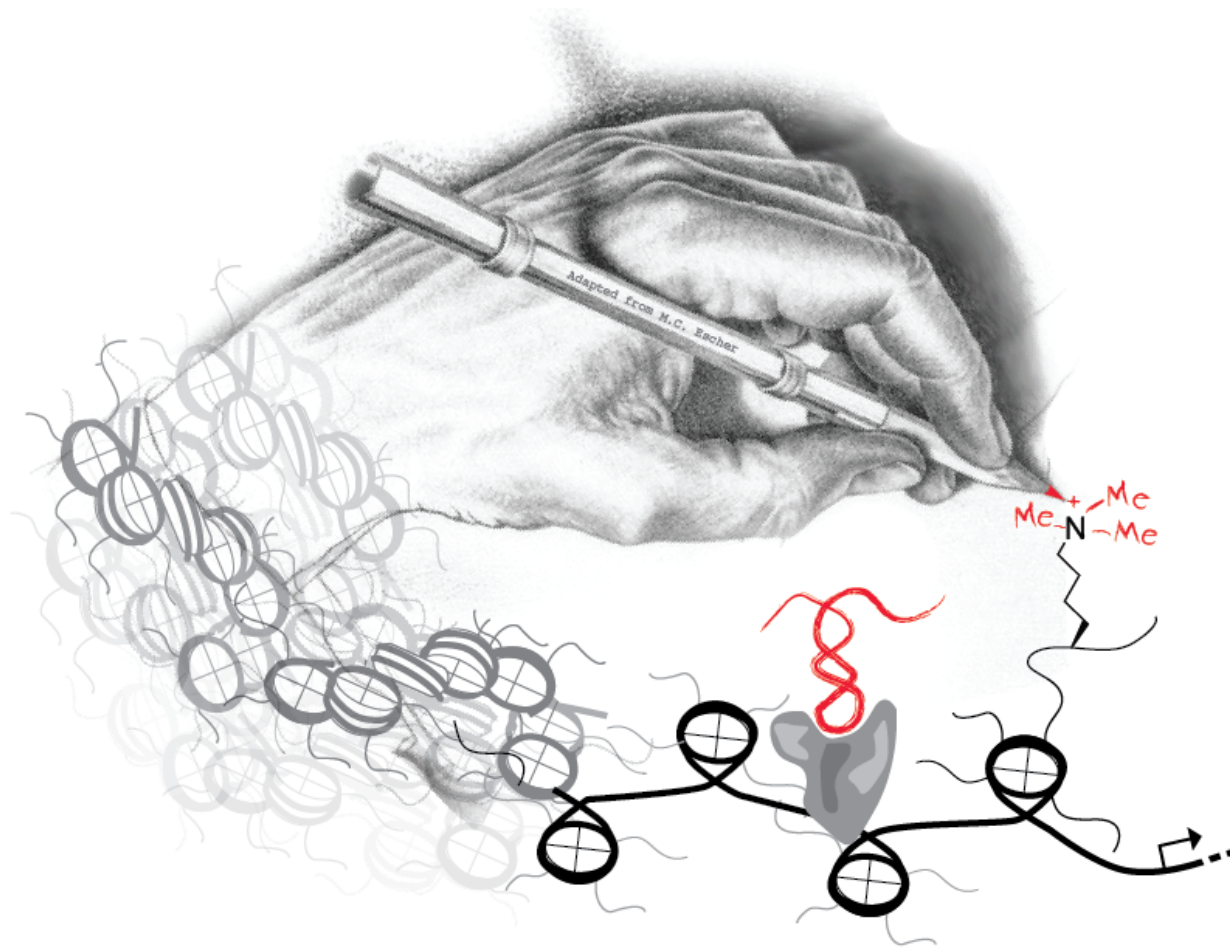
## Cluster Generation

Sequencing templates are immobilized on a proprietary flow cell surface (Figure 1) designed to present the DNA in a manner that facilitates access to enzymes while ensuring high stability of surface-bound template and low non-specific binding of fluorescently labeled nucleotides. Solid-phase amplification (Figures 2–7) creates up to 1,000 identical copies of each single template molecule in close proximity (diameter of one micron or less). Because this process does not involve photolithography, mechanical spotting, or positioning of beads into wells, densities on the order of ten million single-molecule clusters per square centimeter are achieved.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

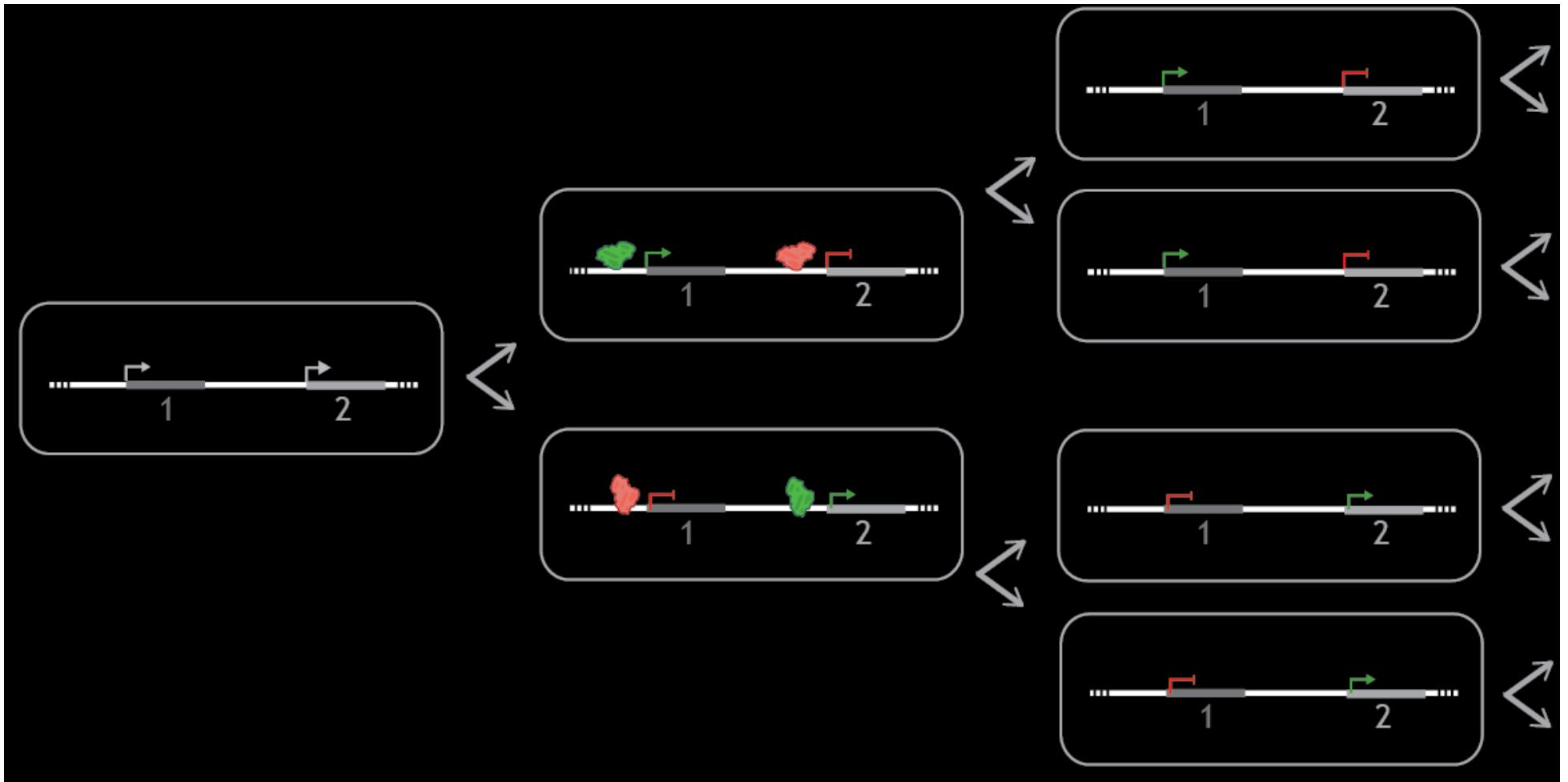1 molecule • $2^n$ = 1000 molecules

$n \approx 10$

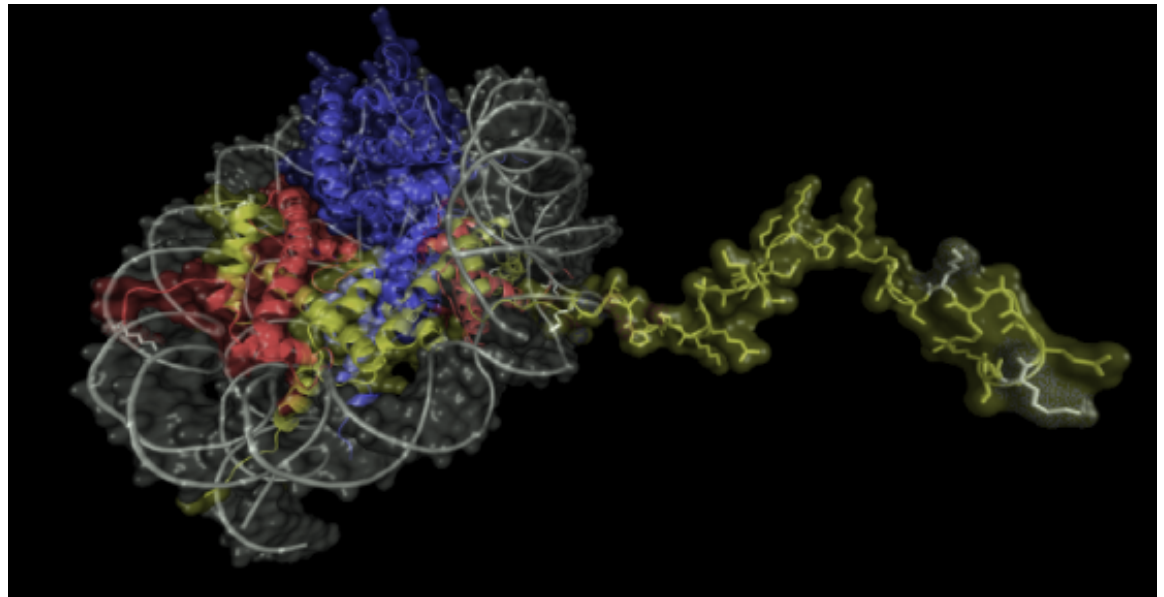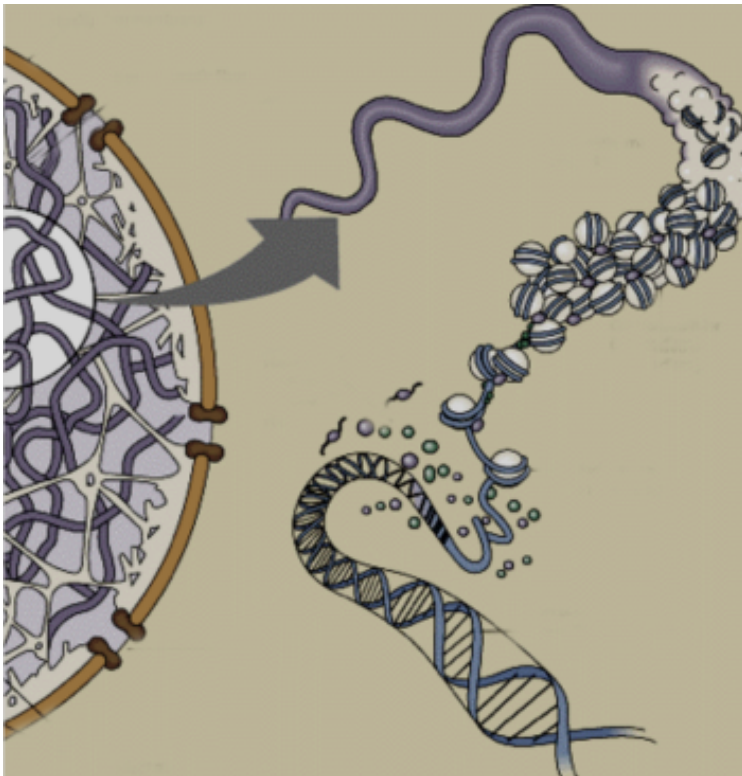# Part 1. How do cells annotate their genomes?

# 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.

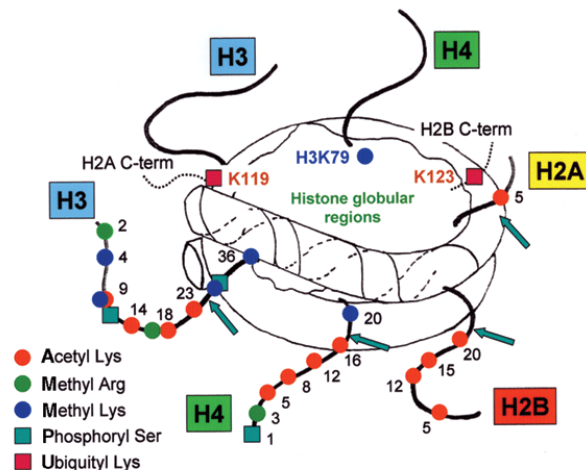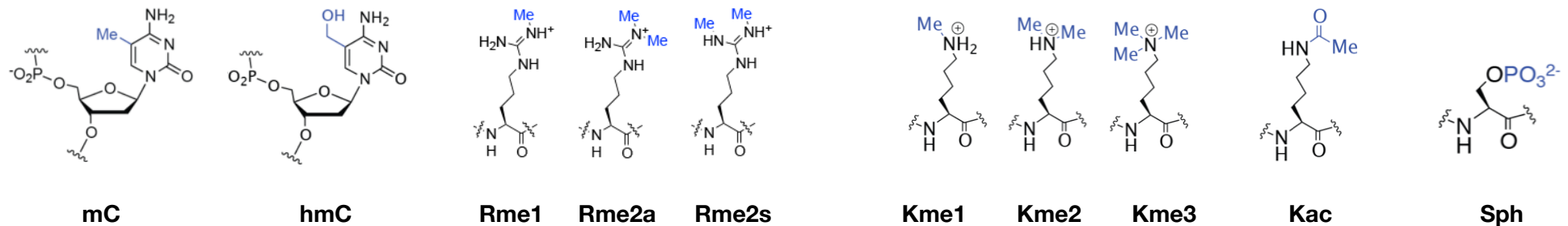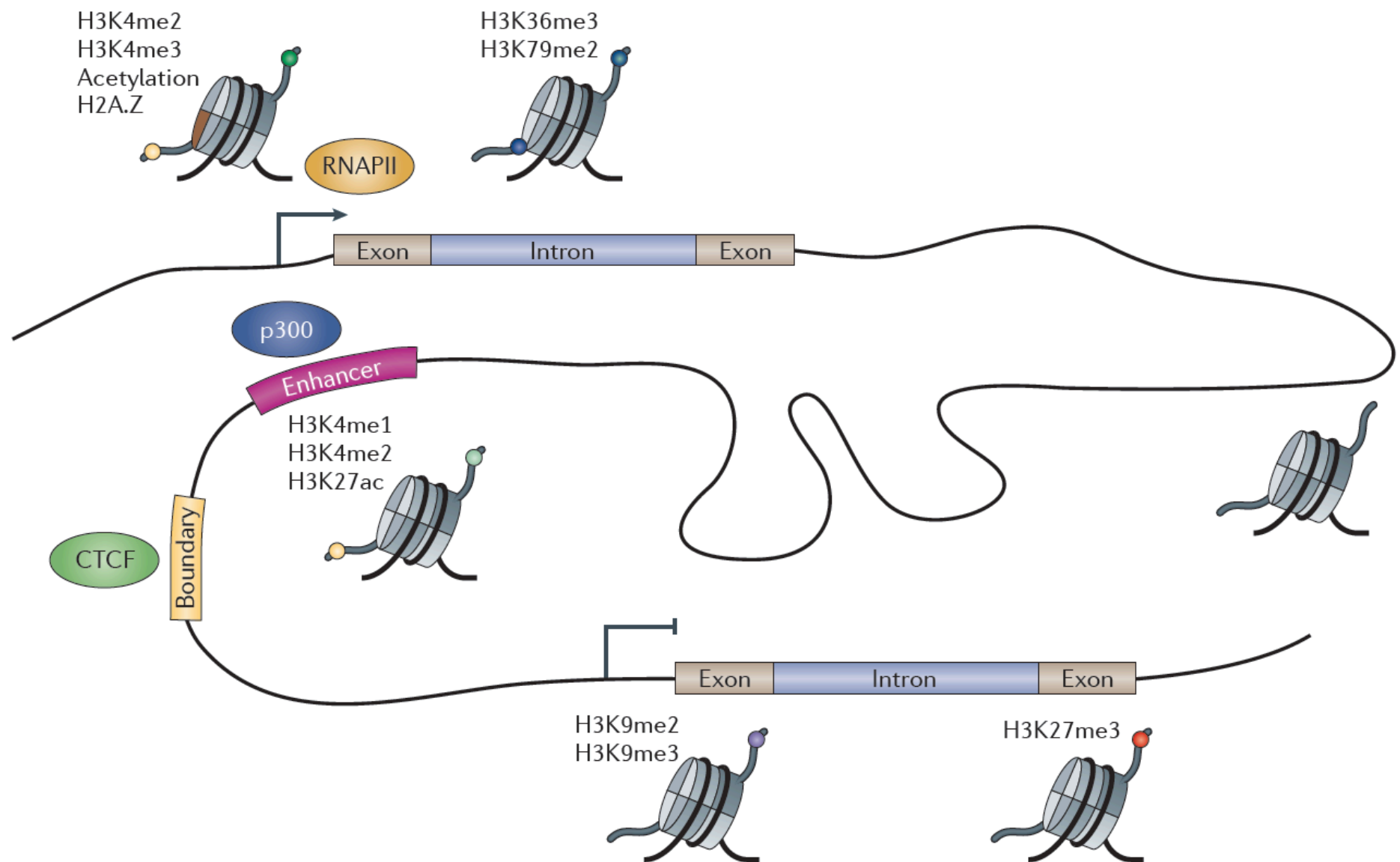# Summary and nomenclature of common covalent modifications.



mC     hmC     Rme1   Rme2a   Rme2s     Kme1   Kme2   Kme3     Kac     Sph



**H3**    **K27**    **ac**

Histone   Residue   Modification

**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 |
| Methyl- | Arginine | mono- | me1 | H3R17me1 |
| | Arginine | di-, symmetrical | me2s | H3R2me2s |
| | Arginine | di-, asymmetrical | me2a | H3R17me2a |
| | Lysine | mono- | me1 | H3K4me1 |
| | Lysine | di- | me2 | H3K4me2 |
| | Lysine | tri- | me3 | H3K4me3 |
| Phosphoryl- | Serine or threonine | mono- | ph | H3S10ph |
| Ubiquityl- | Lysine | mono-[c] | ub1 | H2BK123ub1 |
| SUMOyl- | Lysine | mono- | su | H4K5su[d] |
| ADP ribosyl- | Glutamate | mono- | ar1 | H2BE2ar1 |
| | Glutamate | poly- | arn | H2BE2arn[d] |

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

9

# Chromatin modifications correlate with different genomic functions.

# Installing, binding, and removing modifications



Tollervey and Lunyak (2012) Epigenetics 7:823
Ram *et al.*, *Cell* 147:1628 (2011)

# Regulation is temporally and specially controlled



Brain expression

Brain-expressed transcription factors

Limb expression

Limb-expressed transcription factors

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

# 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 and NET-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.

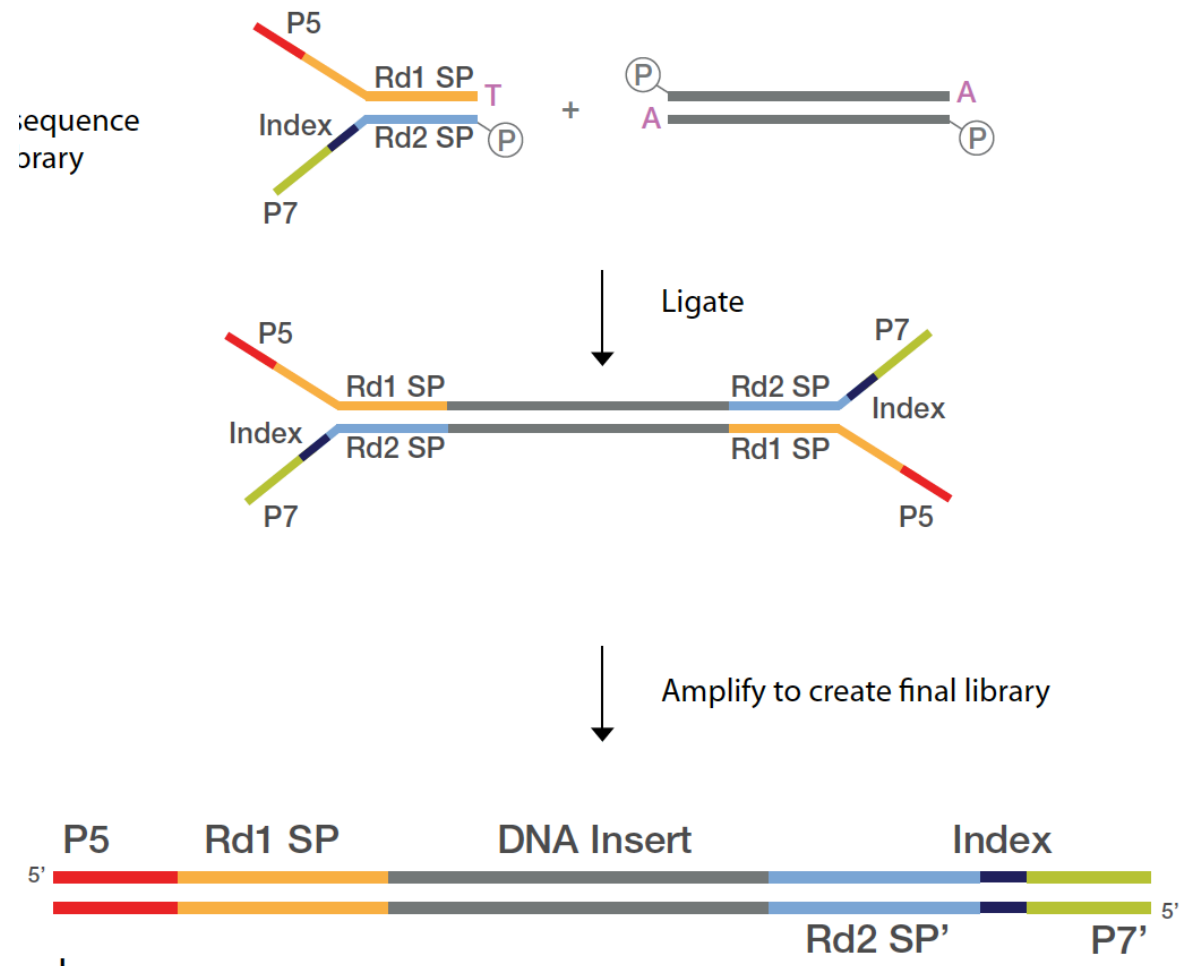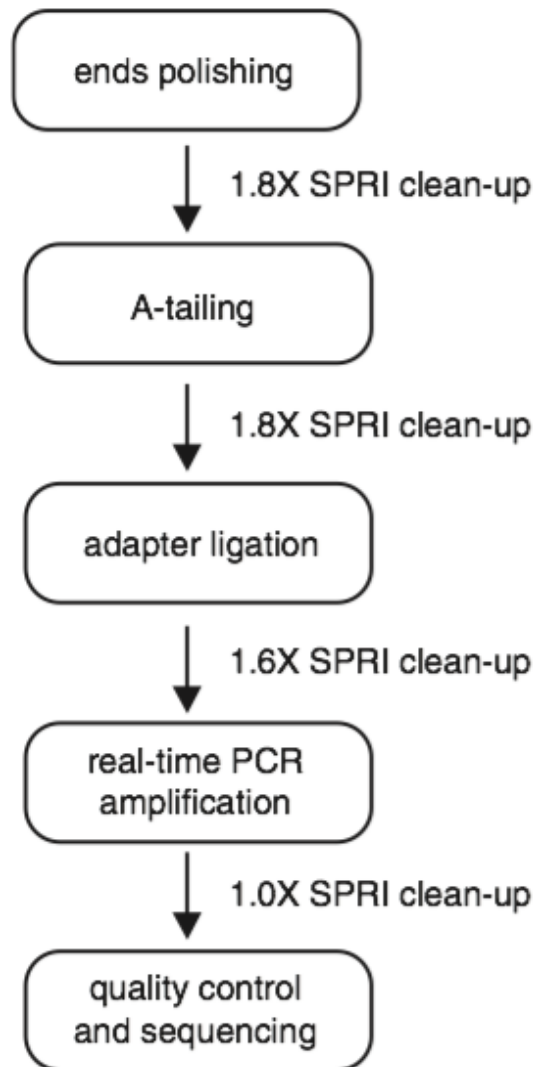2.  **Shear chromatin** to smaller pieces.

    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.

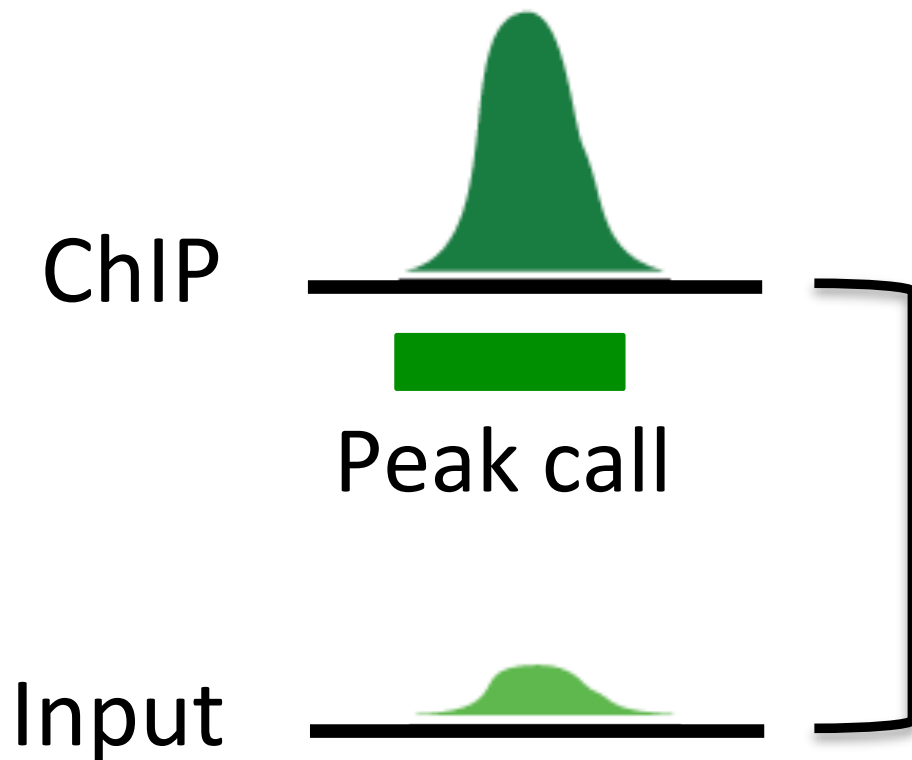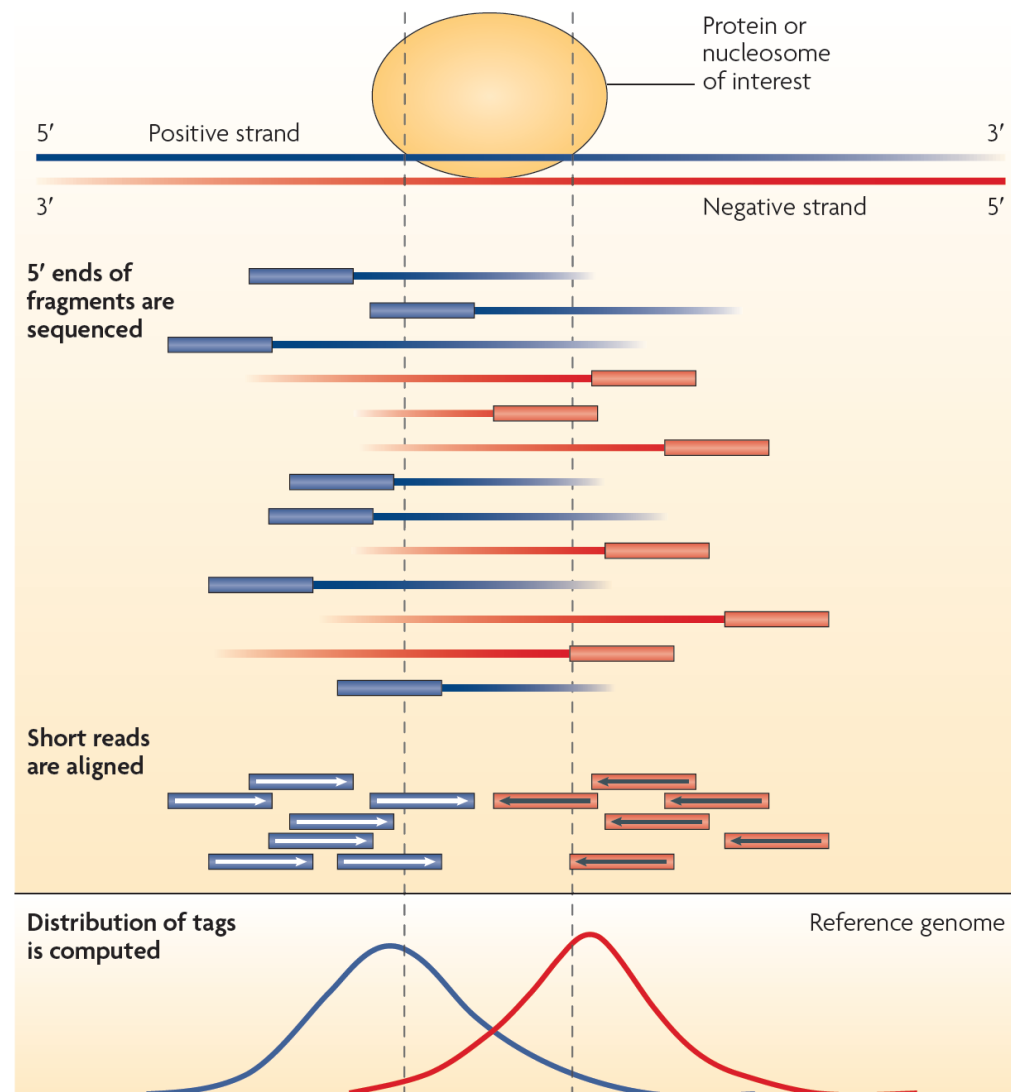# Determining sites of enrichment from ChIP-Seq

ChIP

Input

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

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

# Determining sites of enrichment from ChIP-Seq



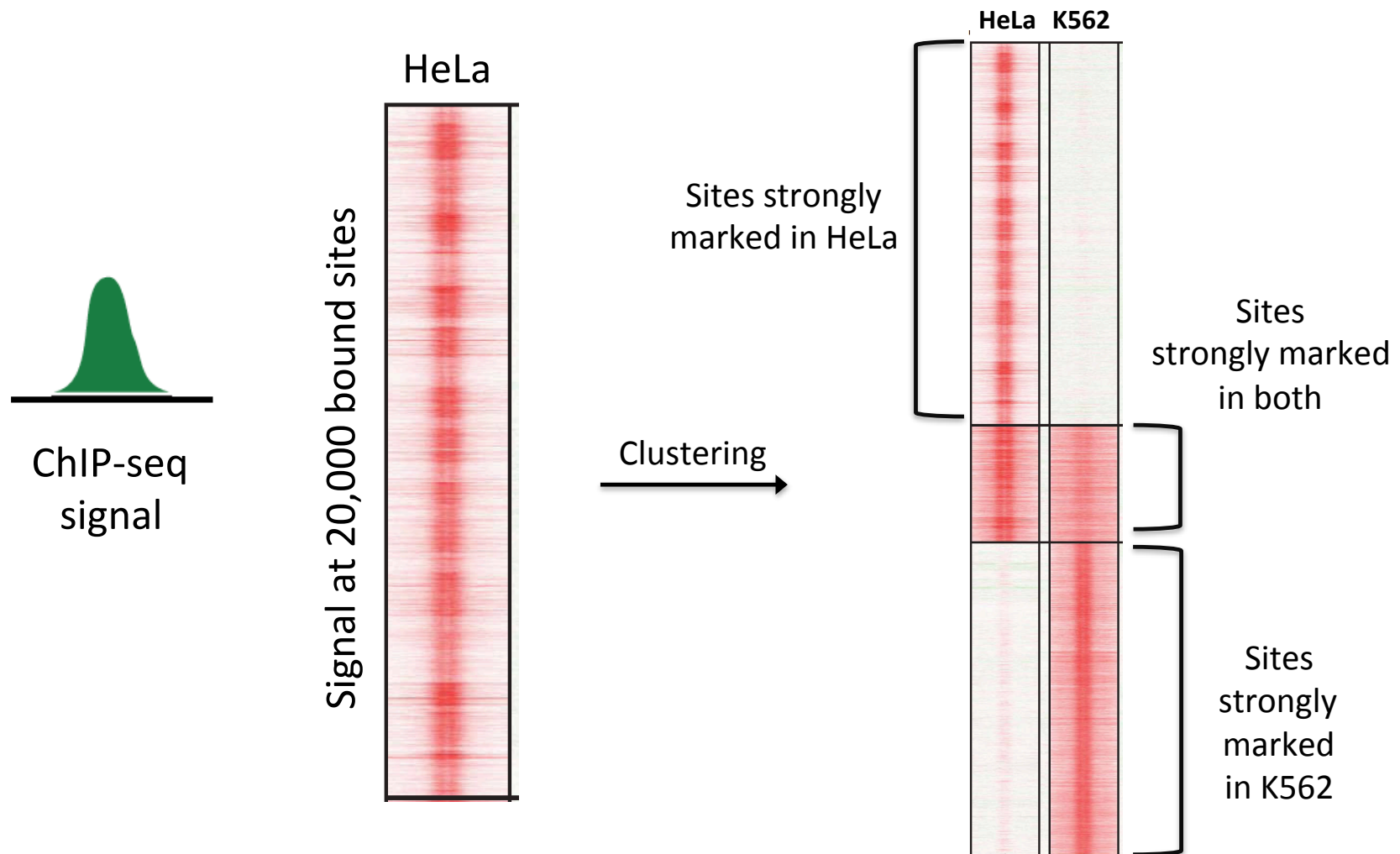ChIP

Peak call

Input

Signal

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.
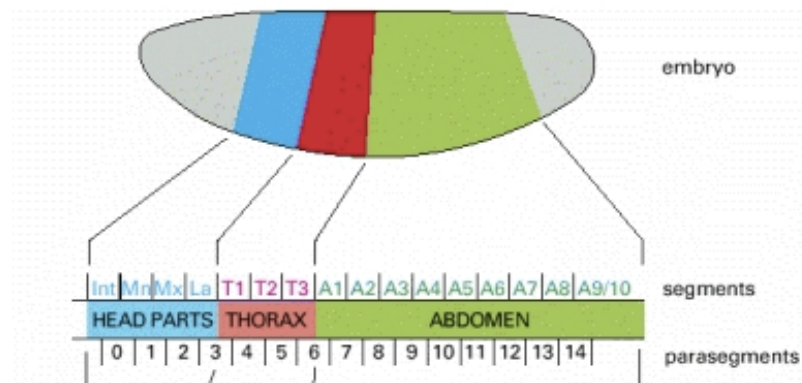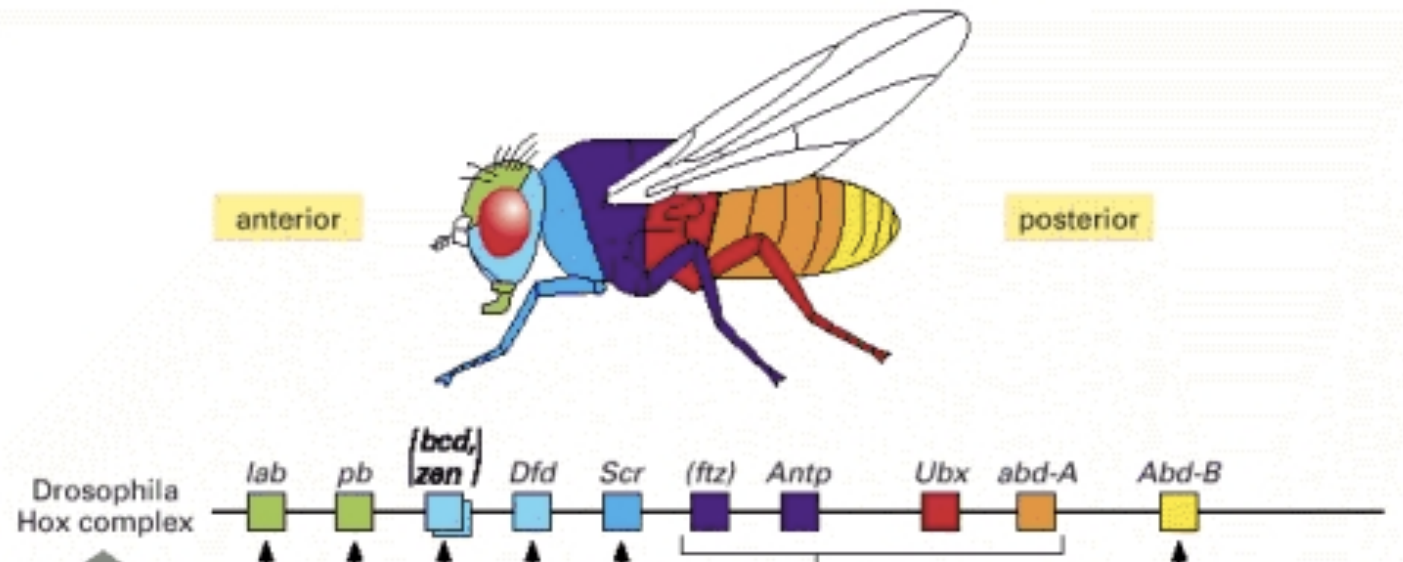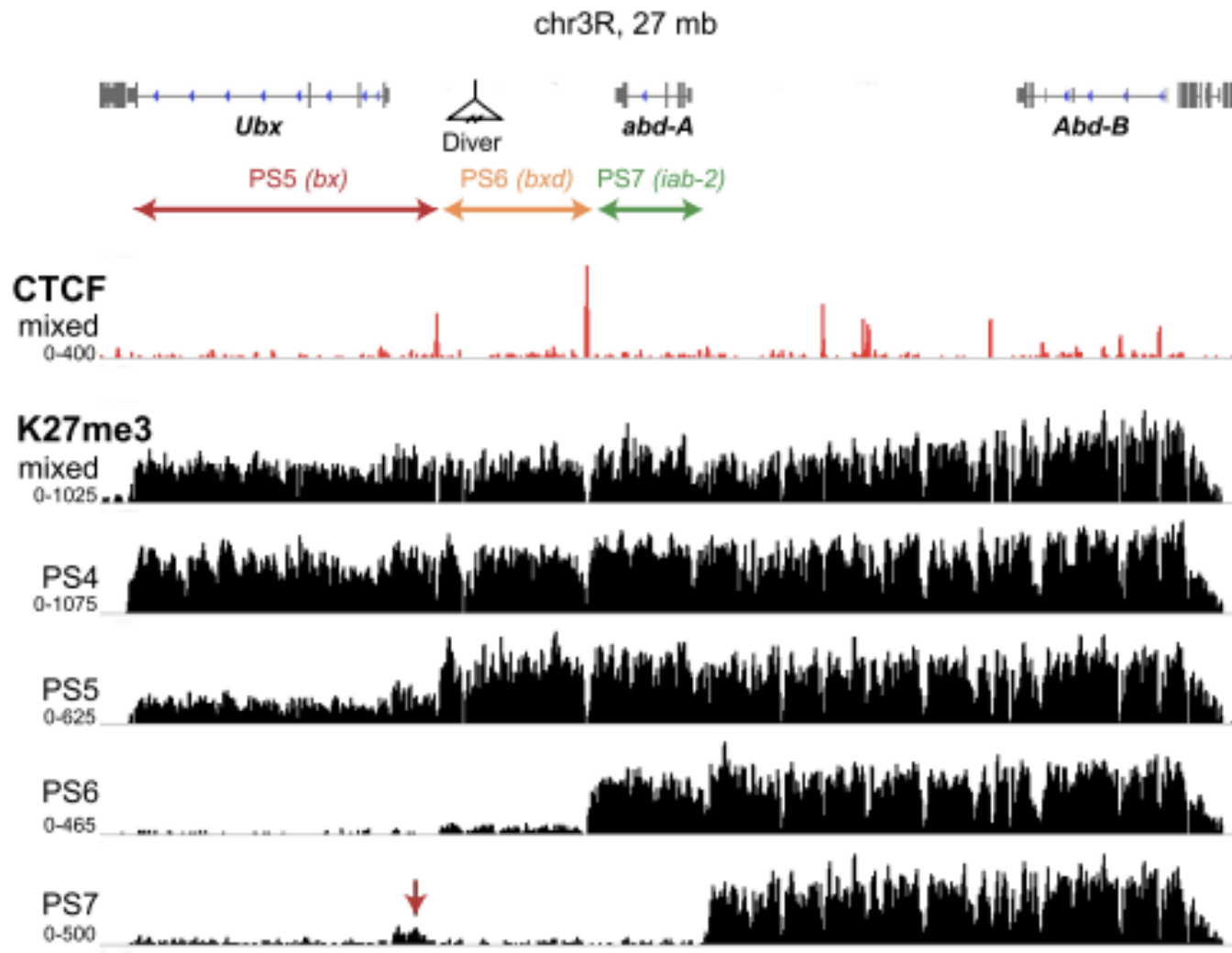
# Avoiding artifacts using features in Seq data



From Park (2009) Nat Rev Genet 10:669

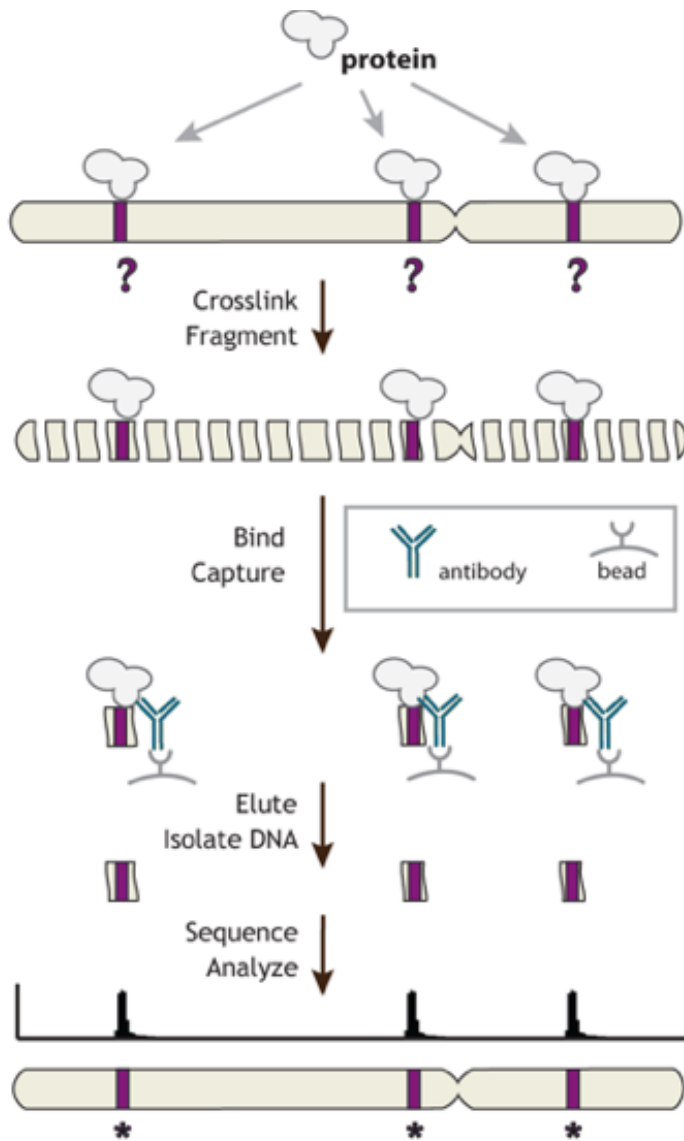# ChIP-Seq signals reveal difference between cells

# Example: Anterior-to-posterior body plan in flies

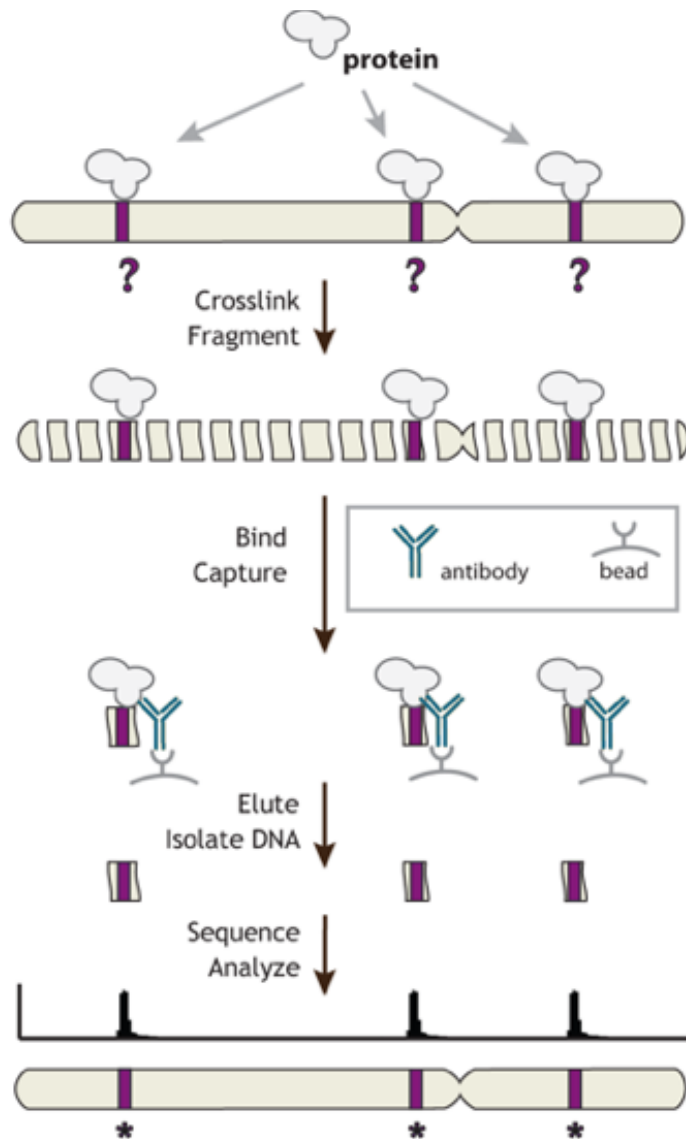# ChIP of CTCF and H3K27me3 in fly development

# Limitations of ChIP-Seq



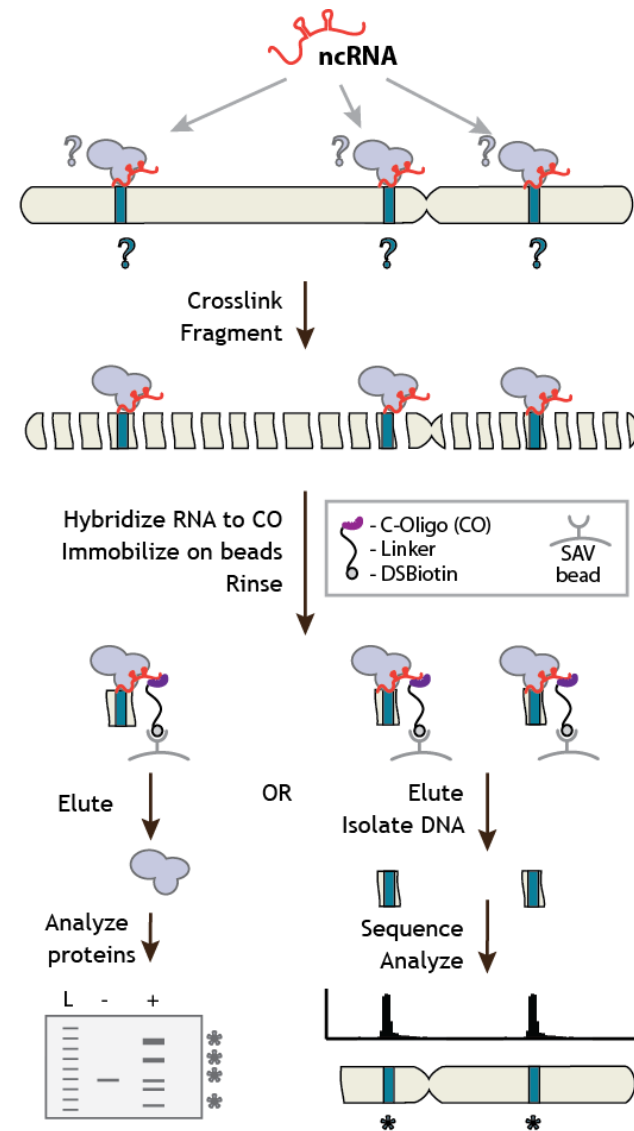1. **Cross linking** efficiency is not necessarily uniform.

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 get very **higher resolution** (ChIP-exo).
.

2. Analysis of **nucleosome turnover** and exchange.

3. Extension to **RNA factors**.

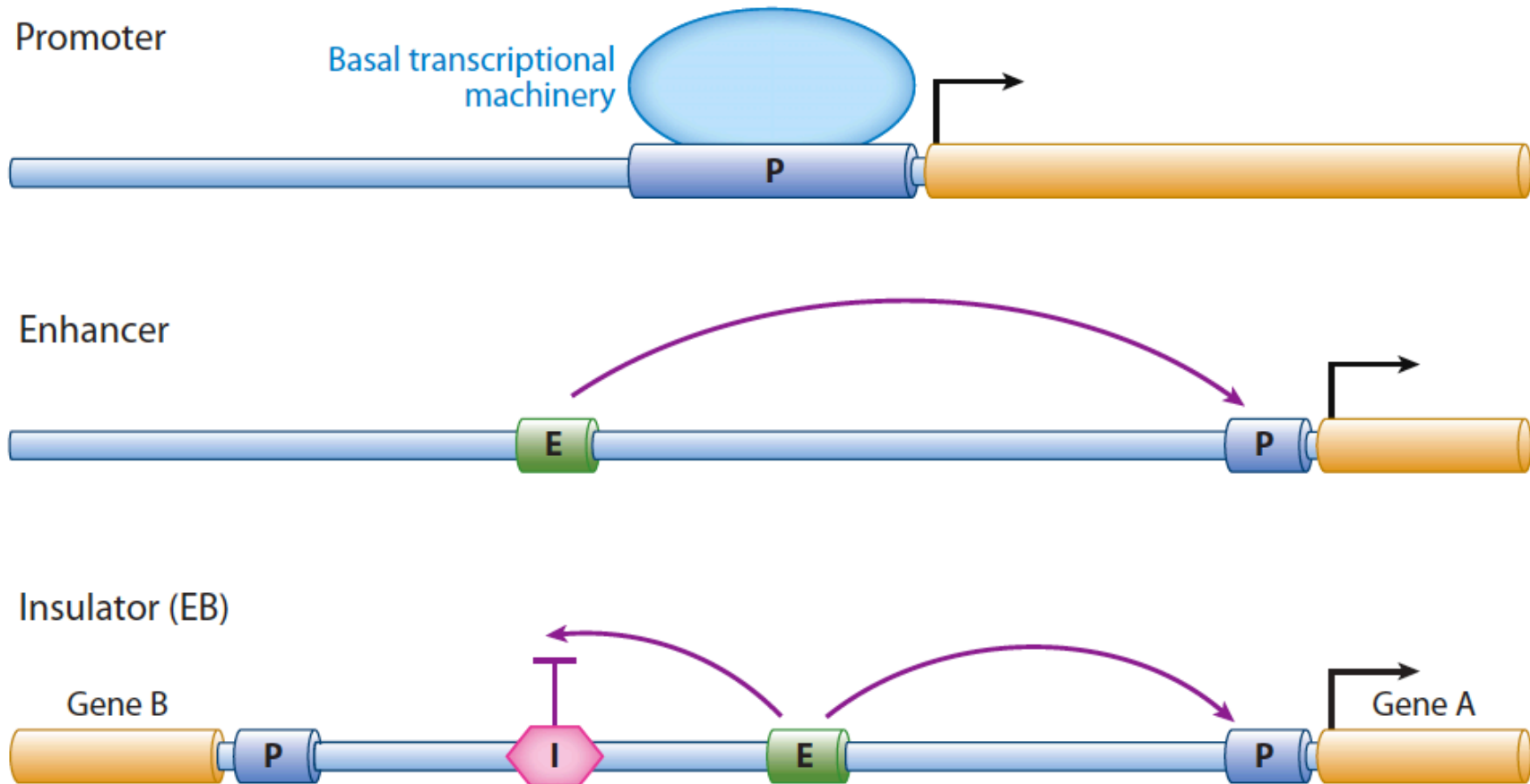# Extension to RNA factors: CHART, ChIRP and RAP
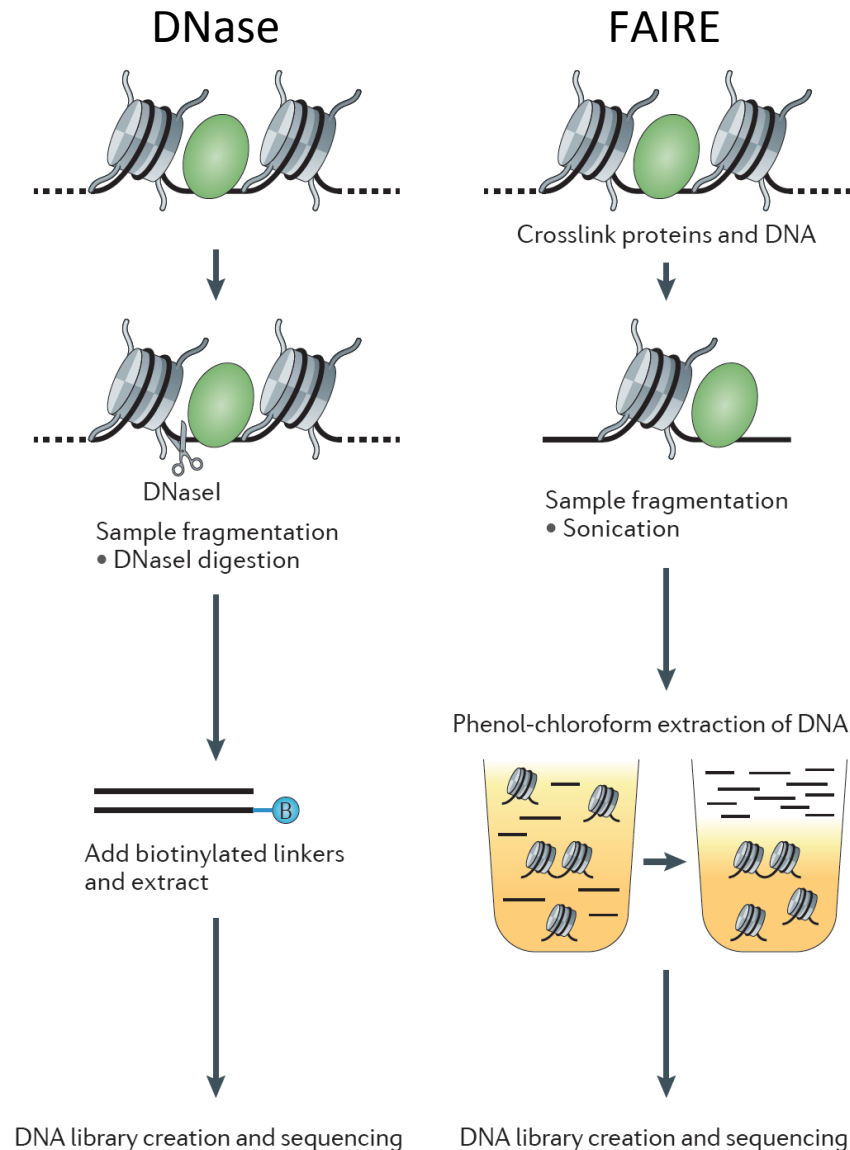
# Using sequencing to annotate the genome

1. **Where are the cis-acting regulatory elements in DNA?**
   A. DNase I hyper-sensitivity mapping (**DNase-Seq**).
   B. **FAIRE** to map regulatory elements.

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** and **NET-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.

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



**DNase**

Sample fragmentation
• DNaseI digestion

DNaseI

Add biotinylated linkers and extract

DNA library creation and sequencing

**FAIRE**

Crosslink proteins and DNA

Sample fragmentation
• Sonication

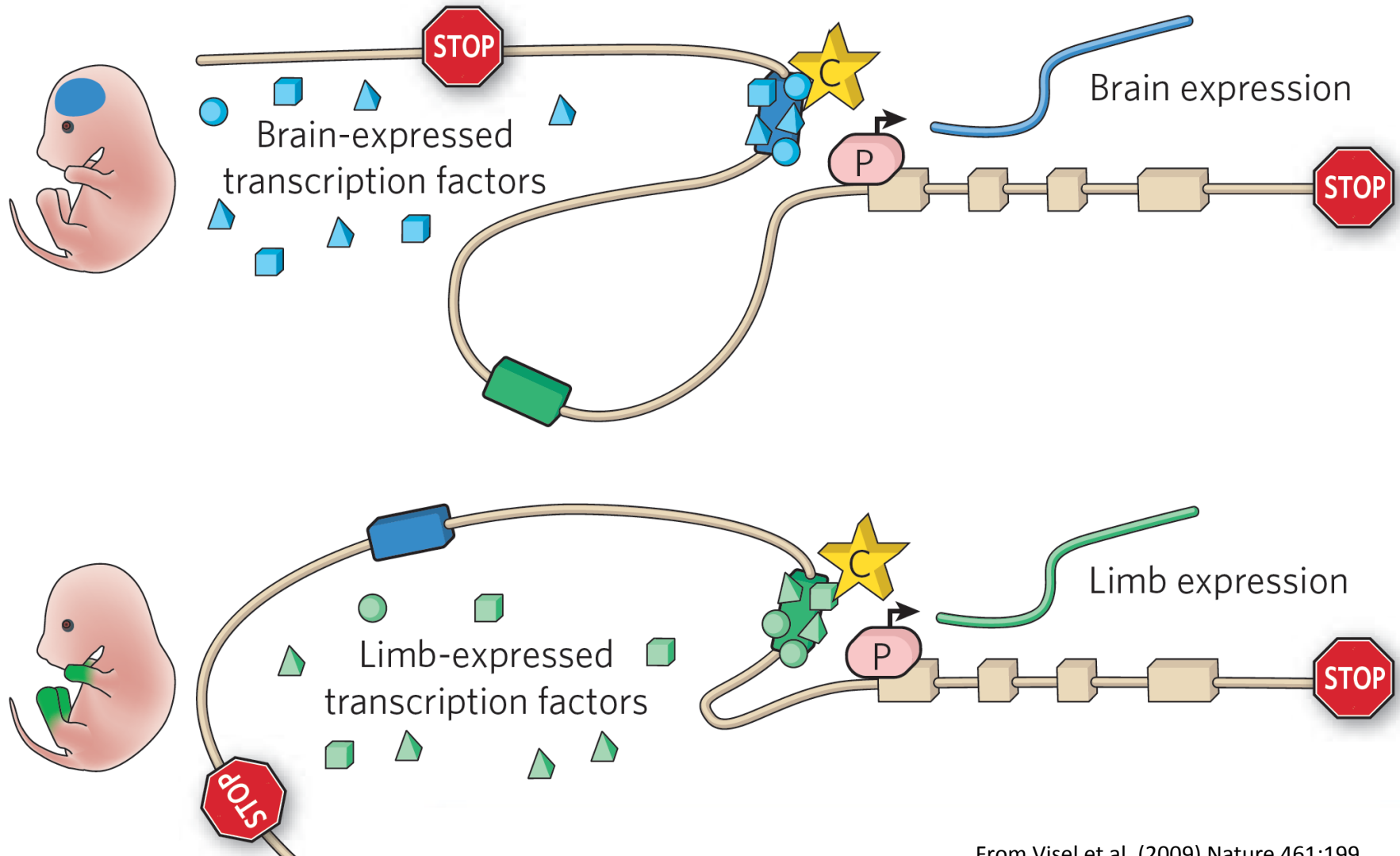Phenol-chloroform extraction of DNA

DNA library creation and sequencing

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

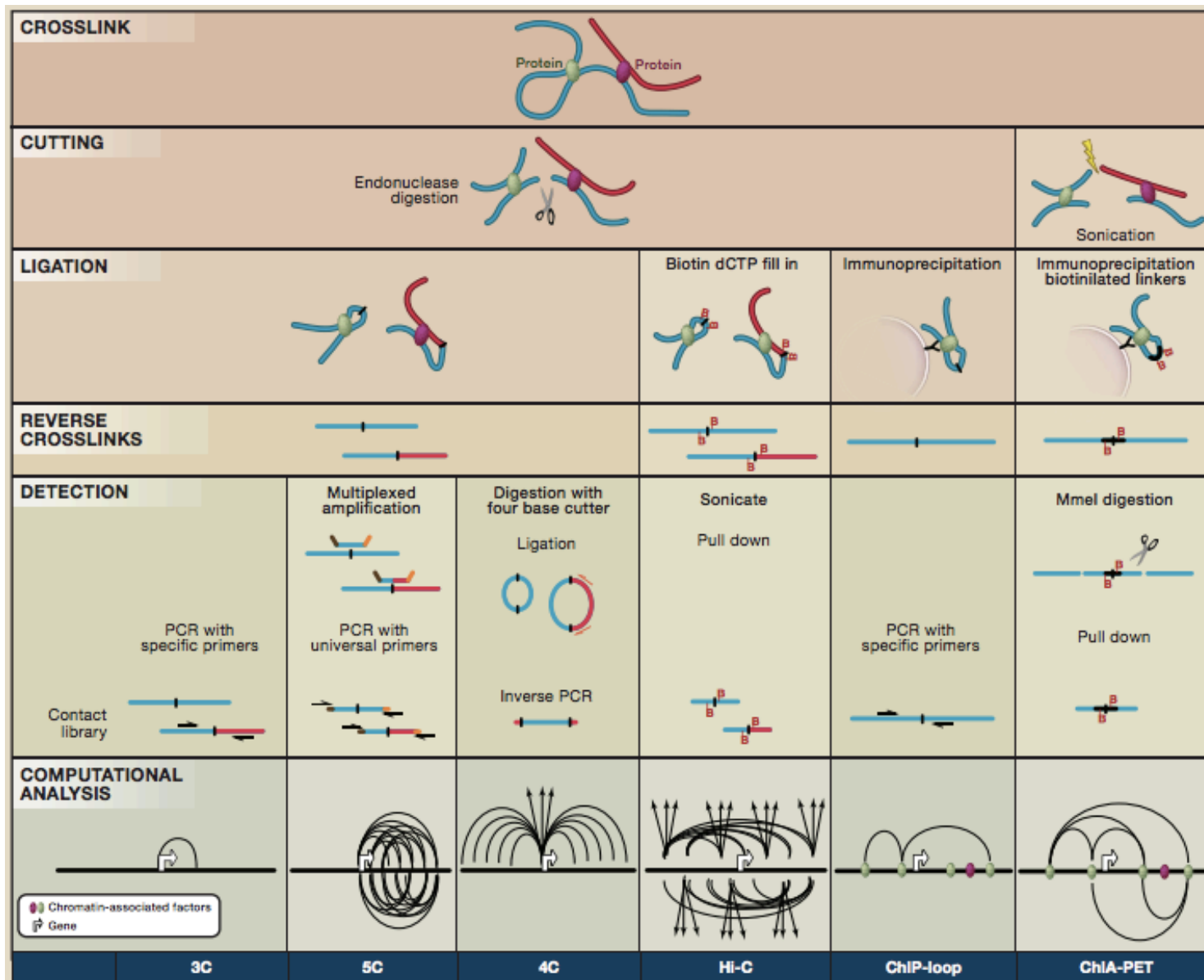2. Changes in **accessibility of chromatin** can provide information about regulation

   -FAIRE-seq (shown)

   -MNase-Seq (not shown).

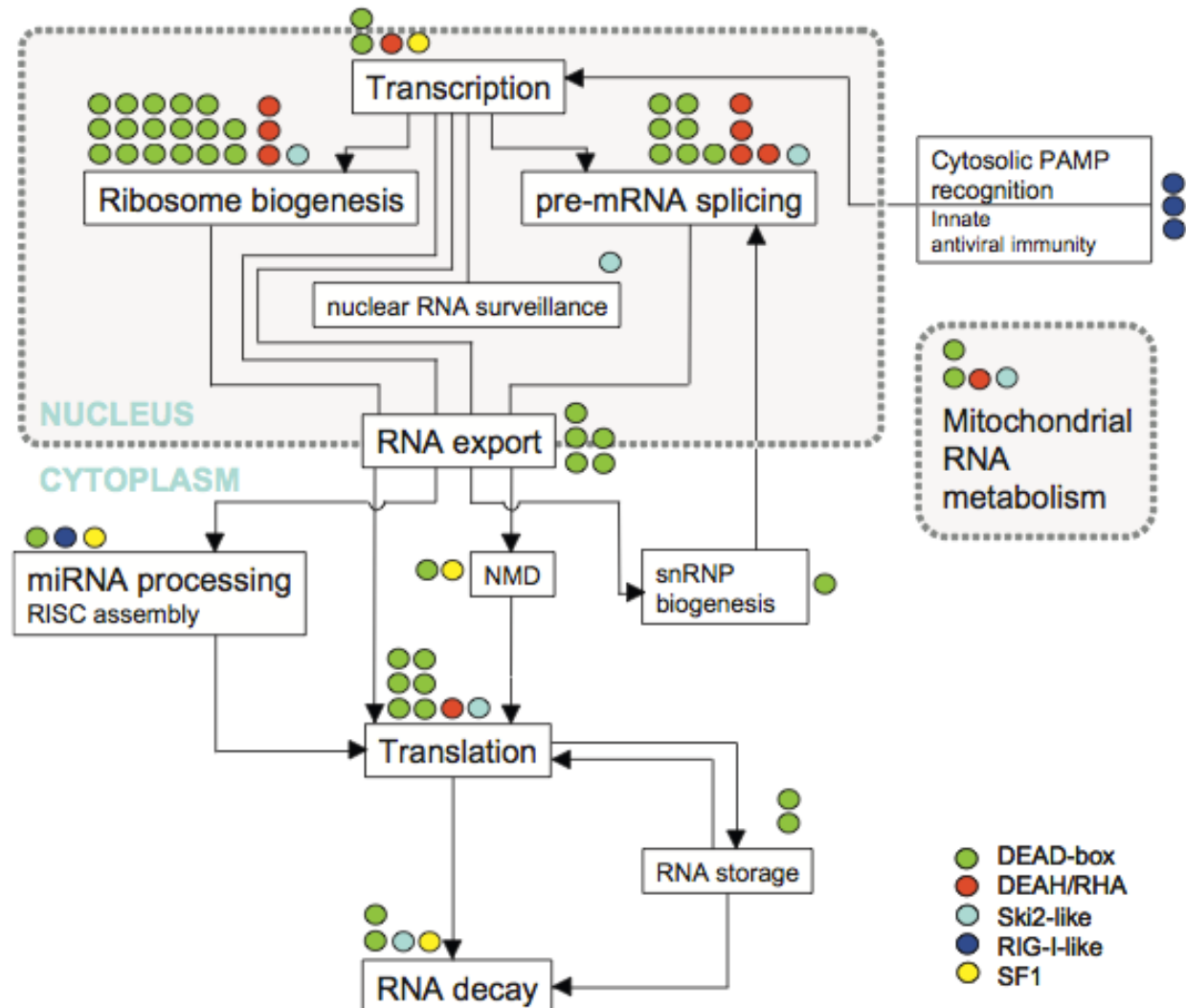# The 3D organization of the genome is important



STOP

Brain-expressed
transcription factors

Brain expression

STOP

STOP

Limb-expressed
transcription factors

Limb expression

STOP

STOP

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

# Techniques to analyze chromatin conformation



Hakim & Misteli, Cell (2012)

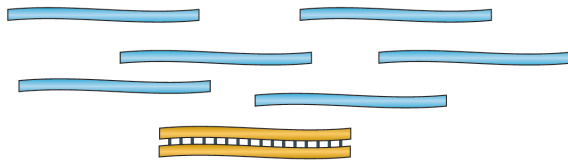# Gene expression is also controlled at the level of RNA

# Part 2: RNA-Seq and applications of RNA-Seq
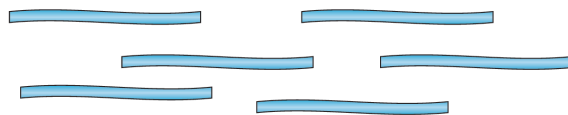
# Using RNA-Seq to examine RNA

- Technical methodology

- Read mapping and normalization

- Estimating isoform-level gene expression

- De novo transcript reconstruction

- Sensitivity and sequencing depth
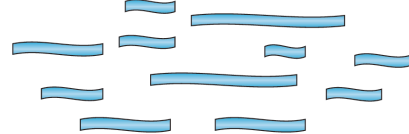
- Differential expression analysis

# RNA-Seq workflow

① mRNA or total RNA

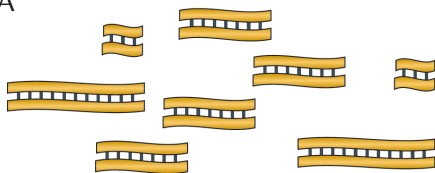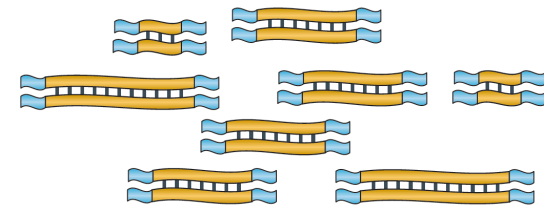② Remove contaminant DNA

Remove rRNA?
Select mRNA?

③ Fragment RNA

④ Reverse transcribe
into cDNA

Strand-specific RNA-seq?

Strand-specific RNA-seq?

⑤ Ligate sequence adaptors

PCR amplification?

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



RiboMinus

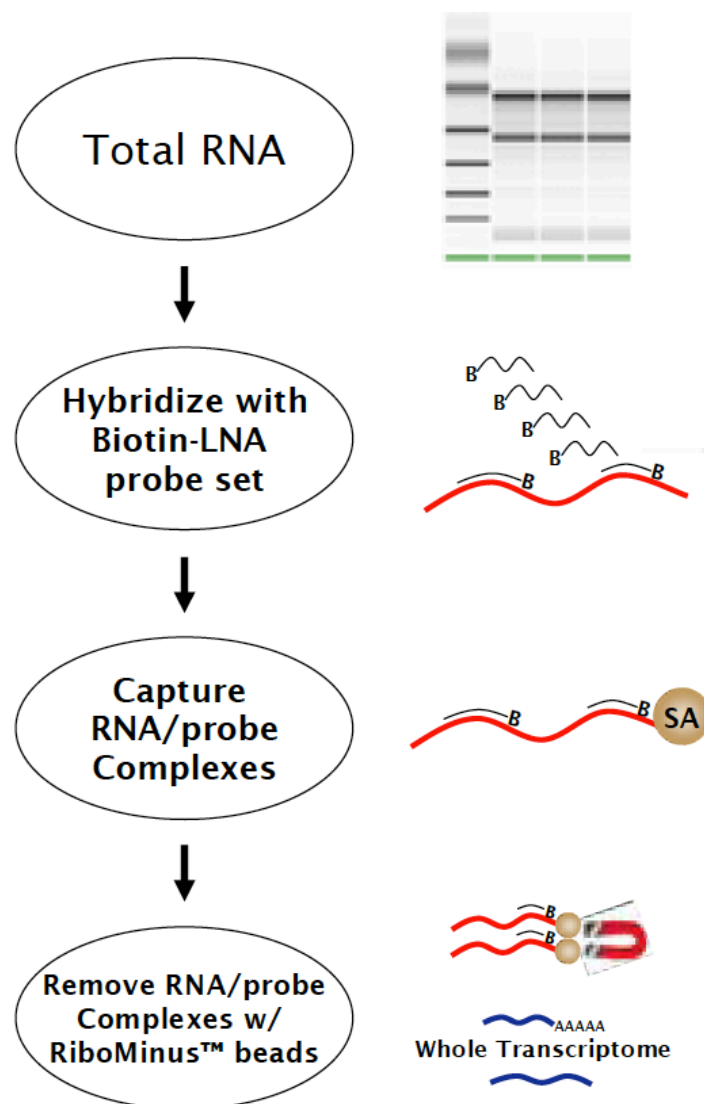# Illumina RNA-seq workflow

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

AAAAAAAA  mRNA

Fragment mRNA

RNA fragments

Strand-specific cDNA synthesis

Synthesize ds cDNA
Ligate adapters
Amplify

```
ATCACAGTGGGACTCCATAAATTTTTCT
CGAAGGACCAGCAGAAACGAGAGAAAAA
GGACAGAGTCCCCAGCGGGCTGAAGGGG
ATGAAACATTAAAGTCAAACAATATGAA
. . . . . .
```

Generate clusters and
sequence

# RNA-Seq reads map mostly to exons



ORF
sequence

**Exonic reads**

**Junction reads**

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

Scale
chr11:
70.4394

| 20 kb |
| 85690000 85695000 85700000 85705000 85710000 85715000 85720000 85725000 85730000 85735000 85740000 85745000 |

RNA_seq

1

UCSC Genes Based on RefSeq, UniProt, GenBank, CCDS and Comparative Genomics
Tbx4
Tbx4

# 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:**
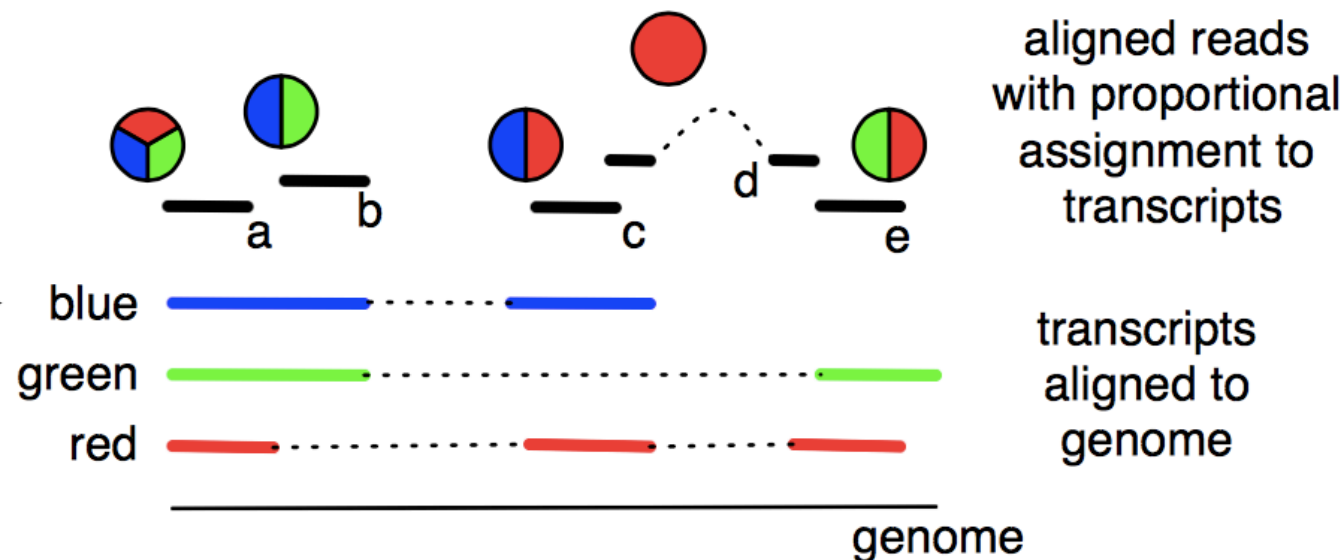Assemble transcripts directly from reads
Allows transcriptome analyses of species without reference genomes

# RNA-seq reads contain information about the abundance of different transcript isoforms
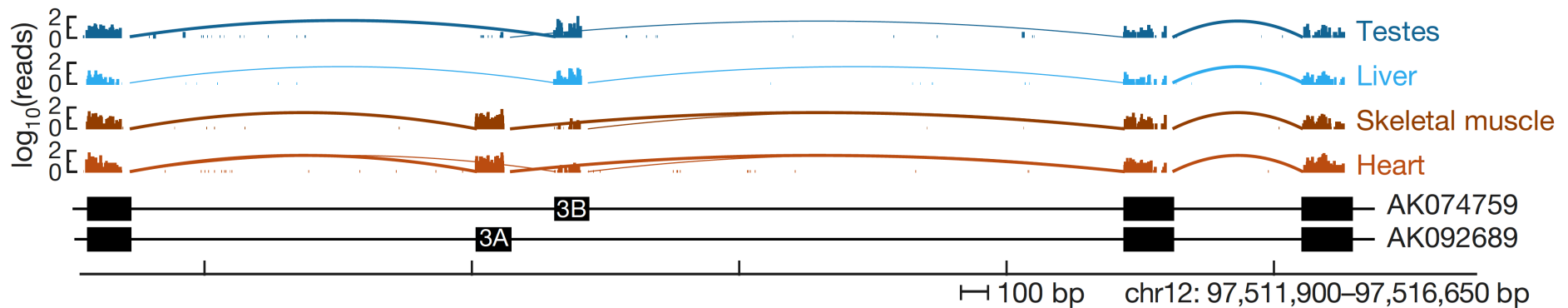
**Normalization** :
  **Internal**: *Reads or Fragments* per kilobase of feature length per million mapped reads (RPKM or FPKM)
  **External**: Reads relative to a standard "spike"

# There is a lot of functional diversity in transcript isoforms



| Alternative transcript events | | Total events (×10³) | Number detected (×10³) | 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
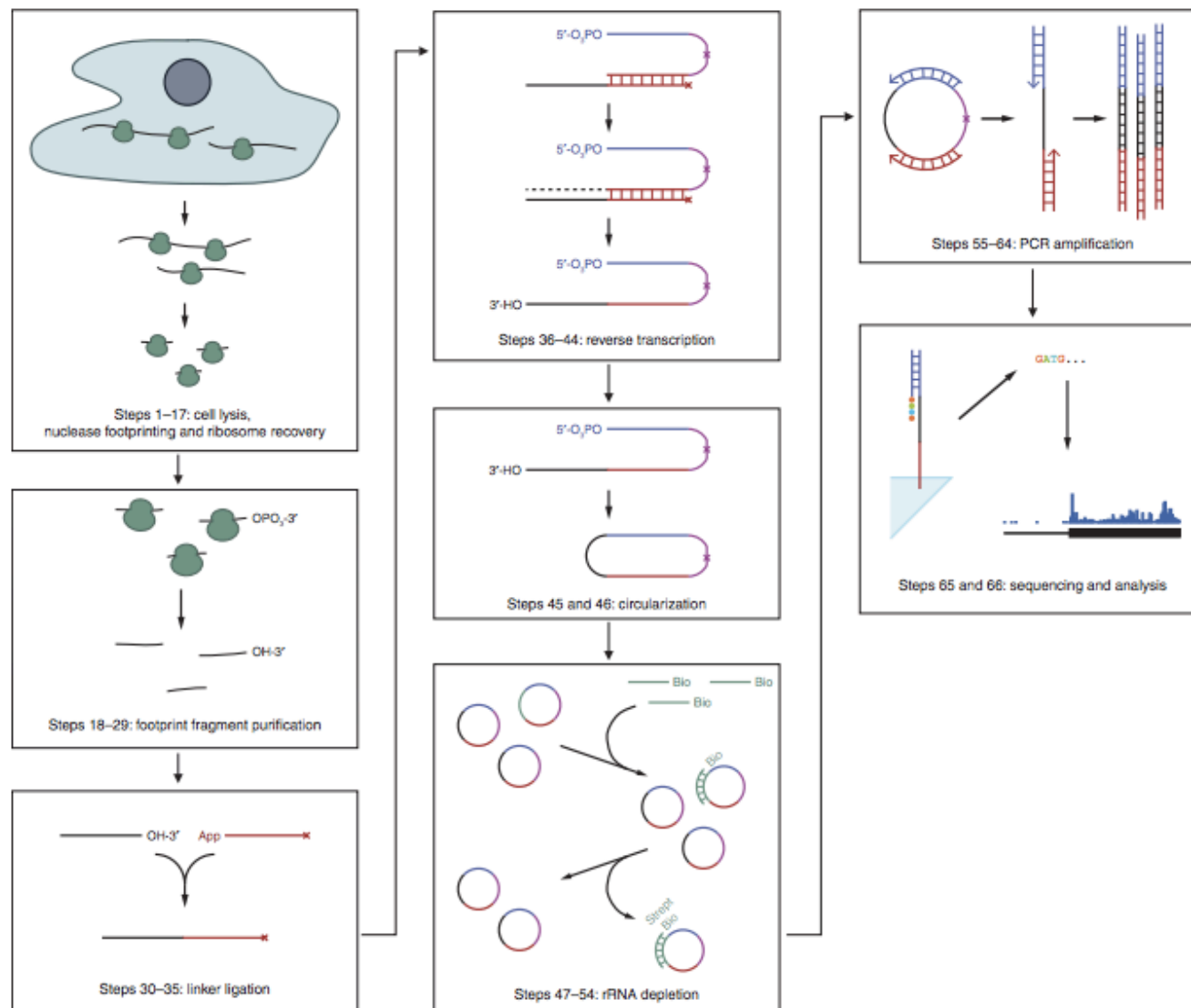
Examining protein-RNA interactions (CLIP, RIP, &c.)

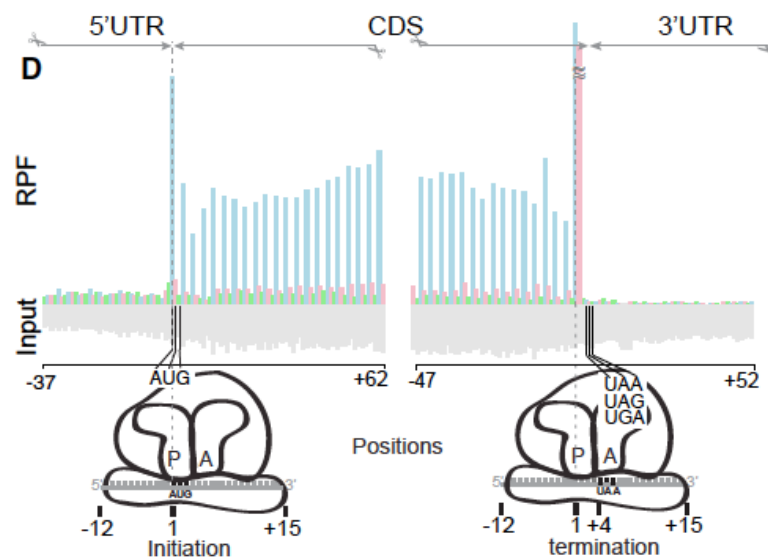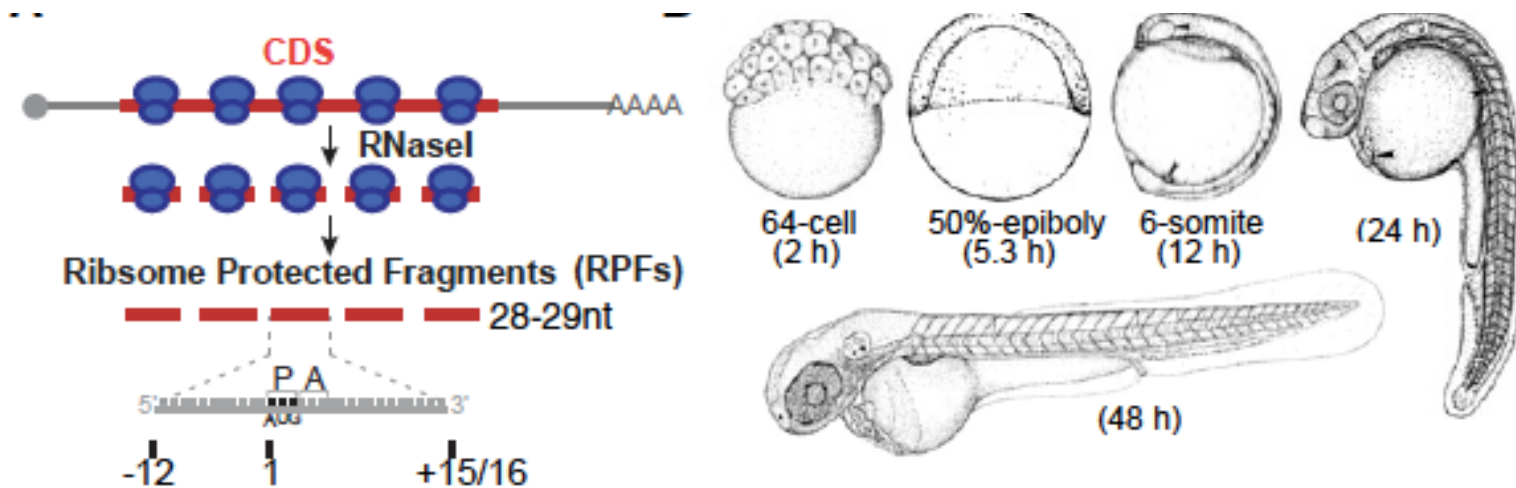Effect of genetic variation on gene expression
    Imprinting
    RNA editing
    Novel events

# Ribosome profiling to reveal translation



Steps 1–17: cell lysis, nuclease footprinting and ribosome recovery

Steps 18–29: footprint fragment purification

Steps 30–35: linker ligation

Steps 36–44: reverse transcription

Steps 45 and 46: circularization

Steps 47–54: rRNA depletion

Steps 55–64: PCR amplification

Steps 65 and 66: sequencing and analysis

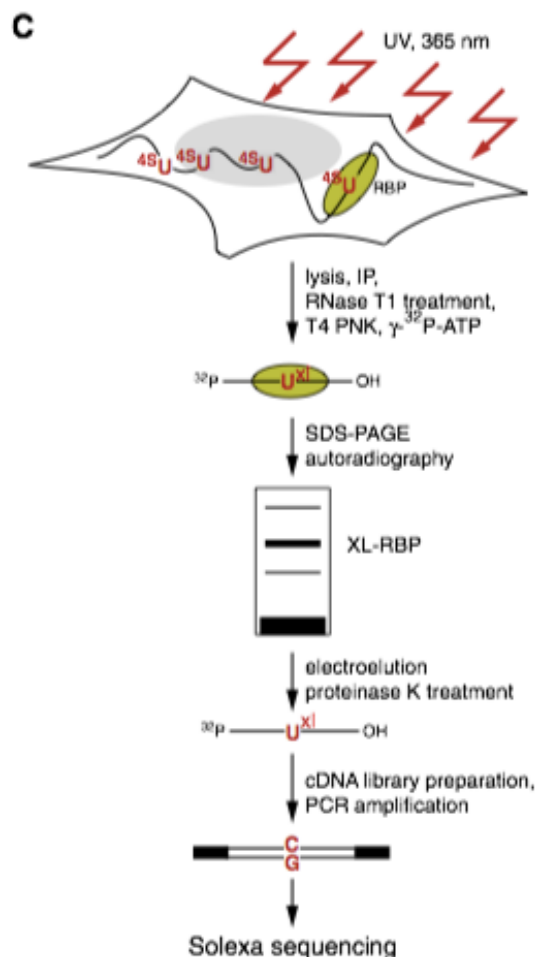# Ribosome foot printing can reveal which reading frame is translated.

1. Bazzini, A. A. *et al.* Identification of small ORFs in vertebrates using ribosome footprinting and evolutionary conservation. *EMBO J* **33,** 981–993 (2014).

# 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.

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

Using sequencing to study _____.
(noun)