

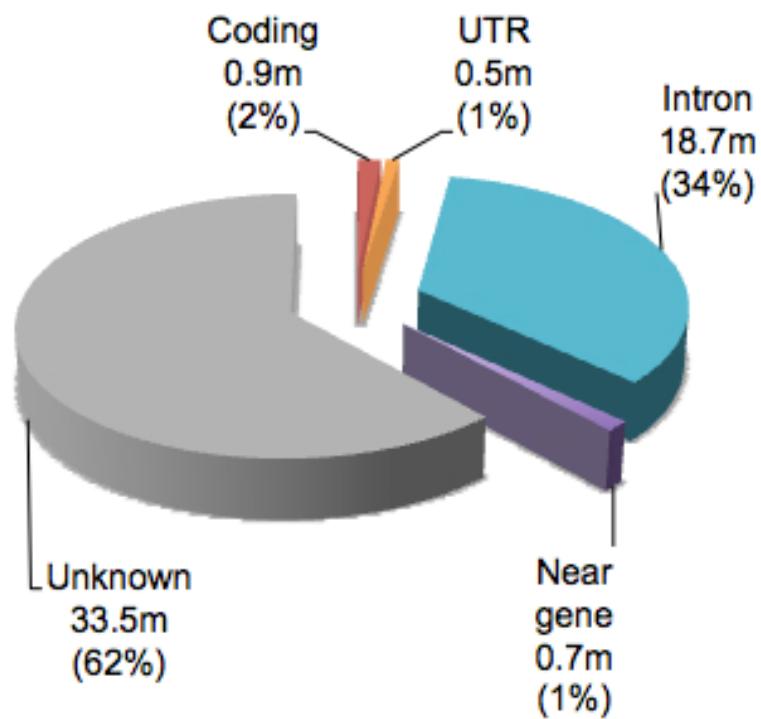
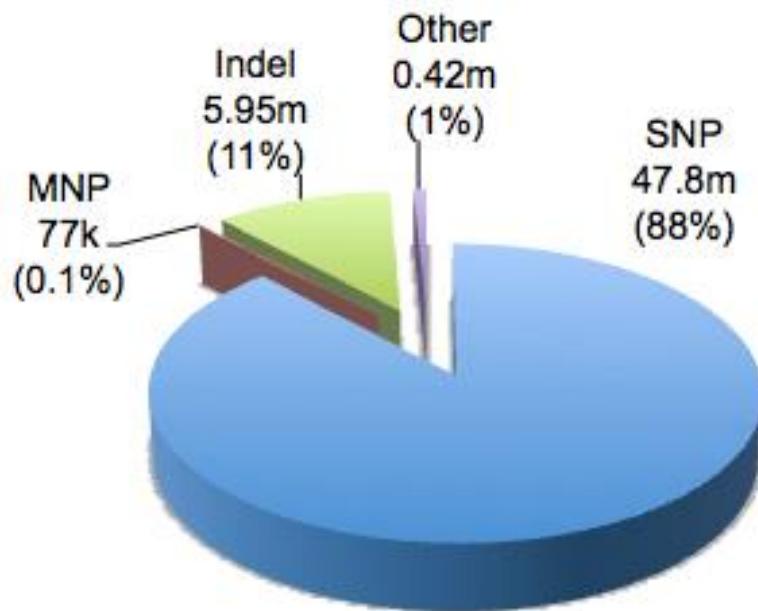
# Transcriptional Regulation by Protein Factors and Epigenetic Mechanisms

Introduction to epigenetics, transcription factors, and related next-generation sequencing methods

Alexander Predeus  
Bioinformatics Institute, Saint Petersburg

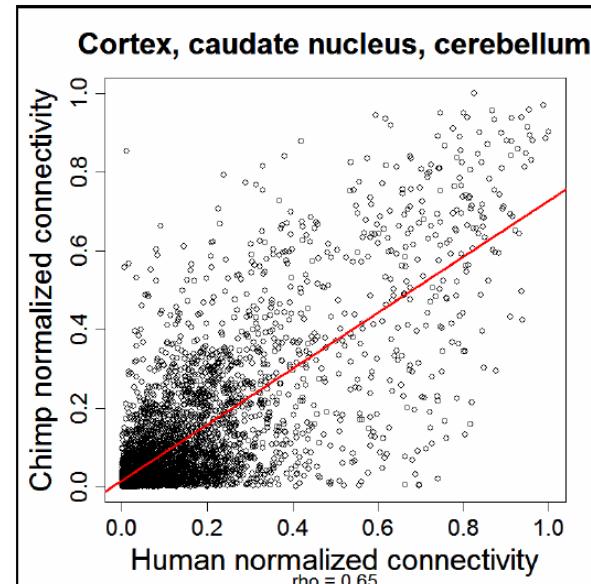
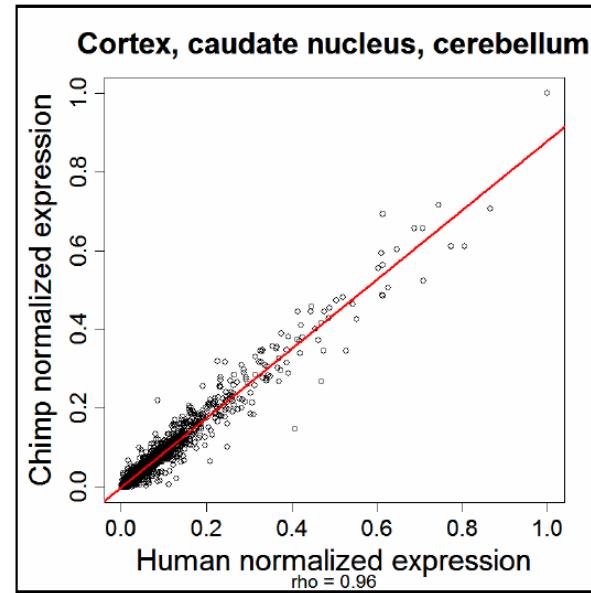
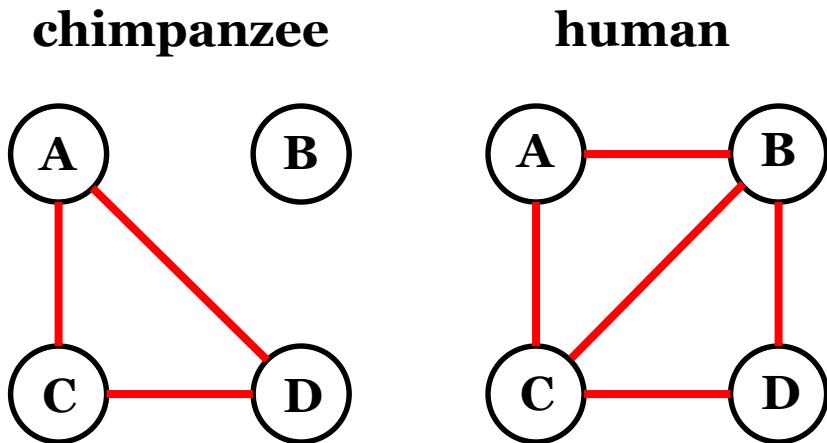
# Why study regulation?

- Evolution mainly works on regulatory, not protein-coding, DNA
- If we know regulation, we can find key spots in large pathways
- Next-generation sequencing! (i.e. because we can)



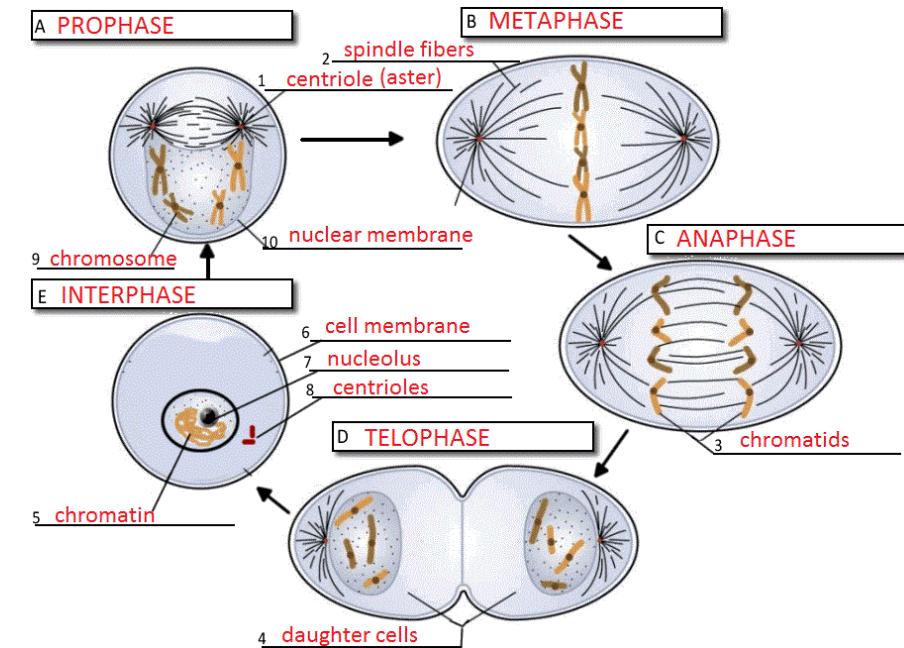
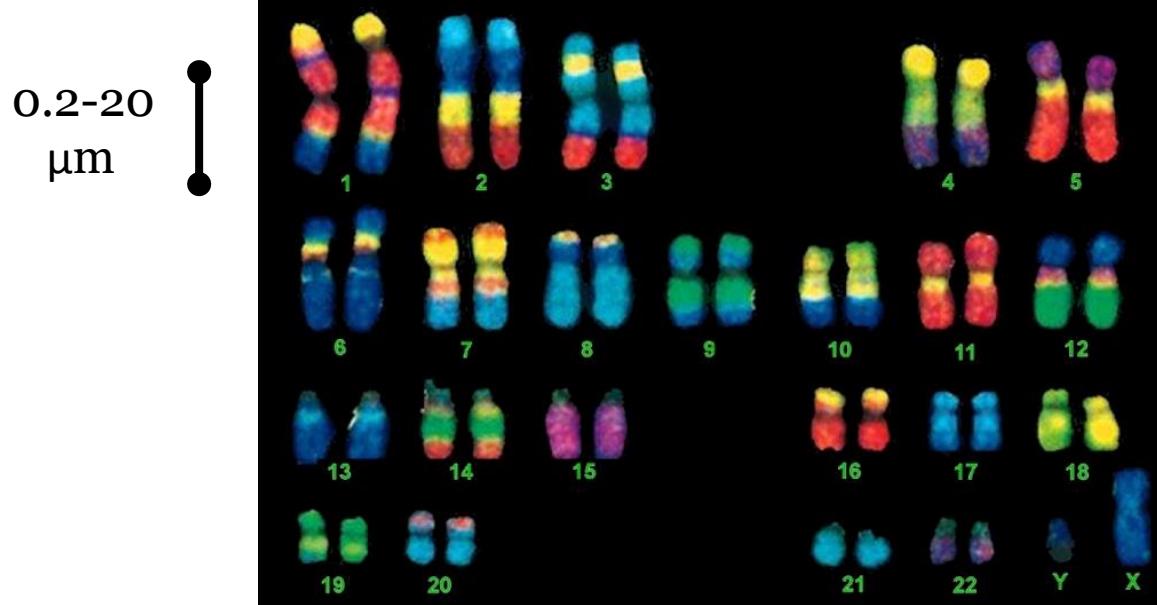
# Primate to human: it's in the regulation

- gene structure and expression are well conserved
- however, **gene coexpression is not**
- the difference lies in gene regulation



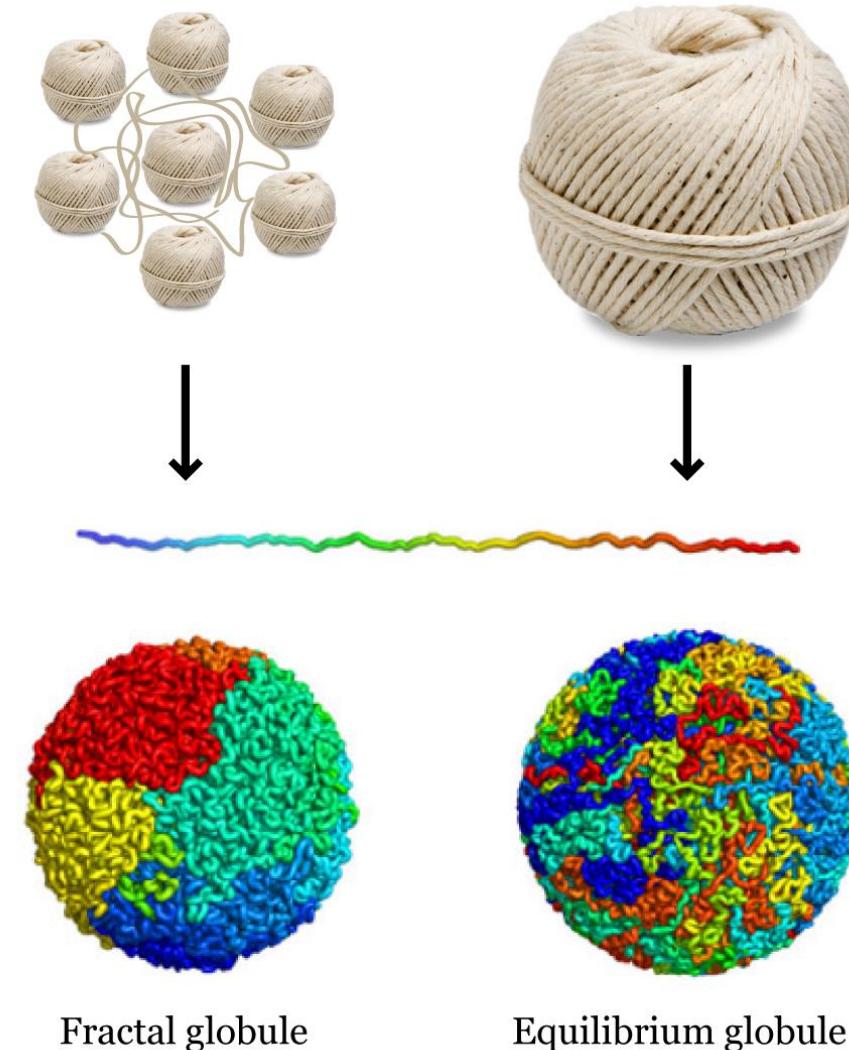
# Chromosomes and chromatin

- chromosomes are dense complexes of **DNA** and **proteins**
- each human chromosome contains on average **5 cm** of DNA
- this is about **2 m** of DNA overall – too much!
- chromatin = euchromatin + heterochromatin



# Chromatin is a fractal globule

- equilibrium ~ roll of twine
- fractal globule ~ many rolls of twine loosely put together
- **easier access, relative independence** of regions

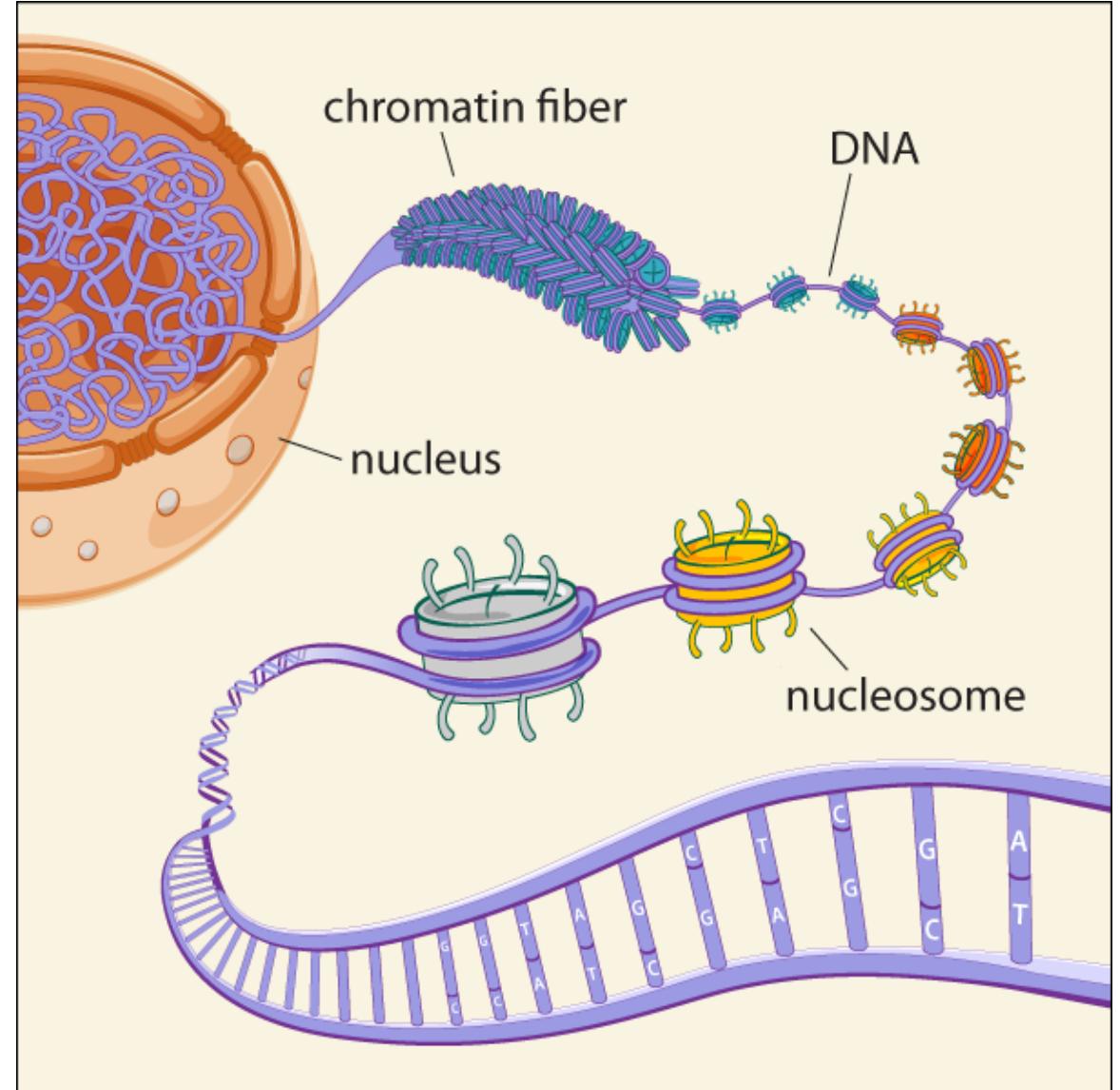


Fractal globule

Equilibrium globule

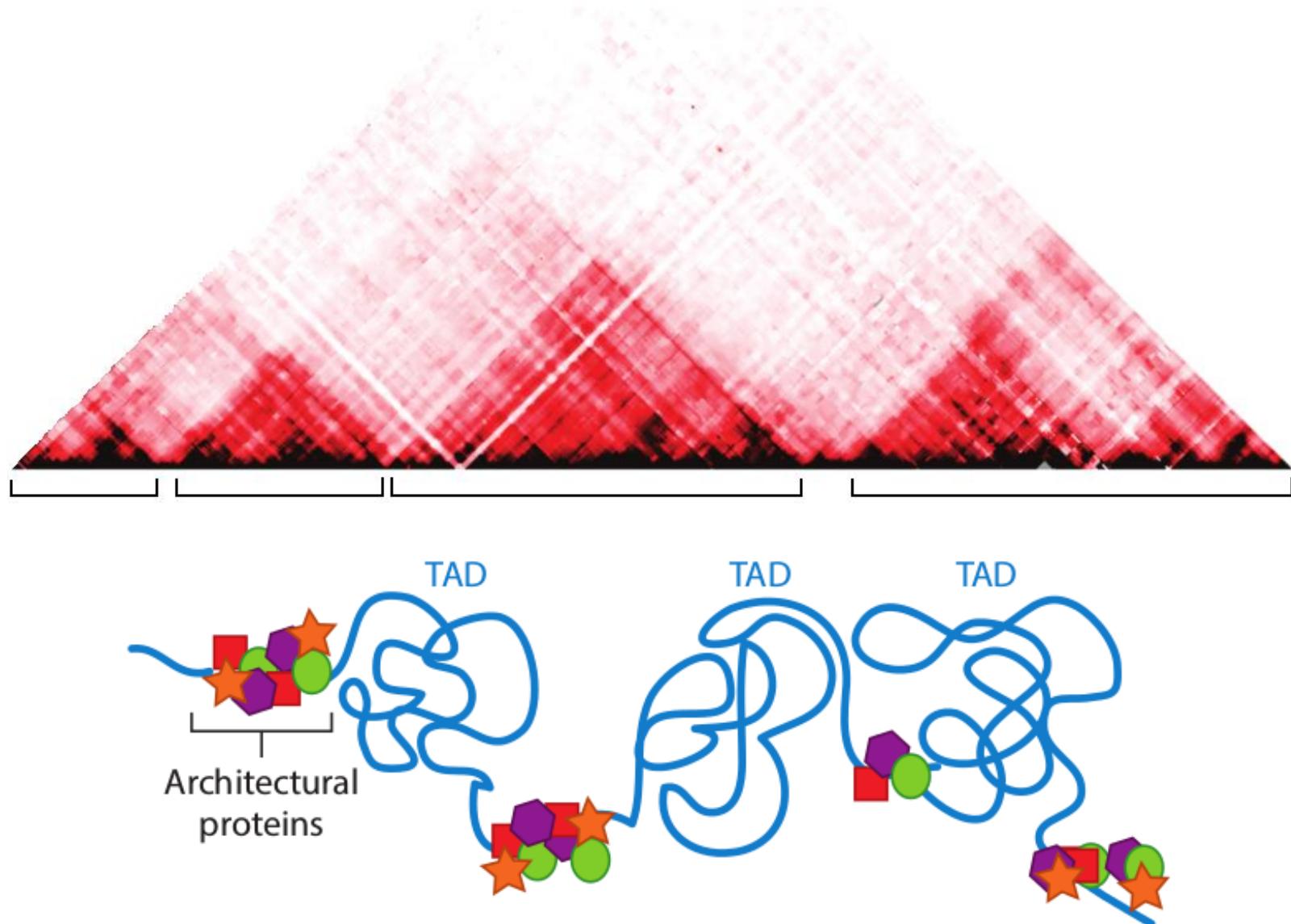
# Chromatin organization

- there are three major levels of DNA packaging:
  - “beads on a string” nucleosomes
  - 30 nm fibre (still euchromatin!)
  - densely packed structures (heterochromatin)
- nucleosome beads include histone proteins



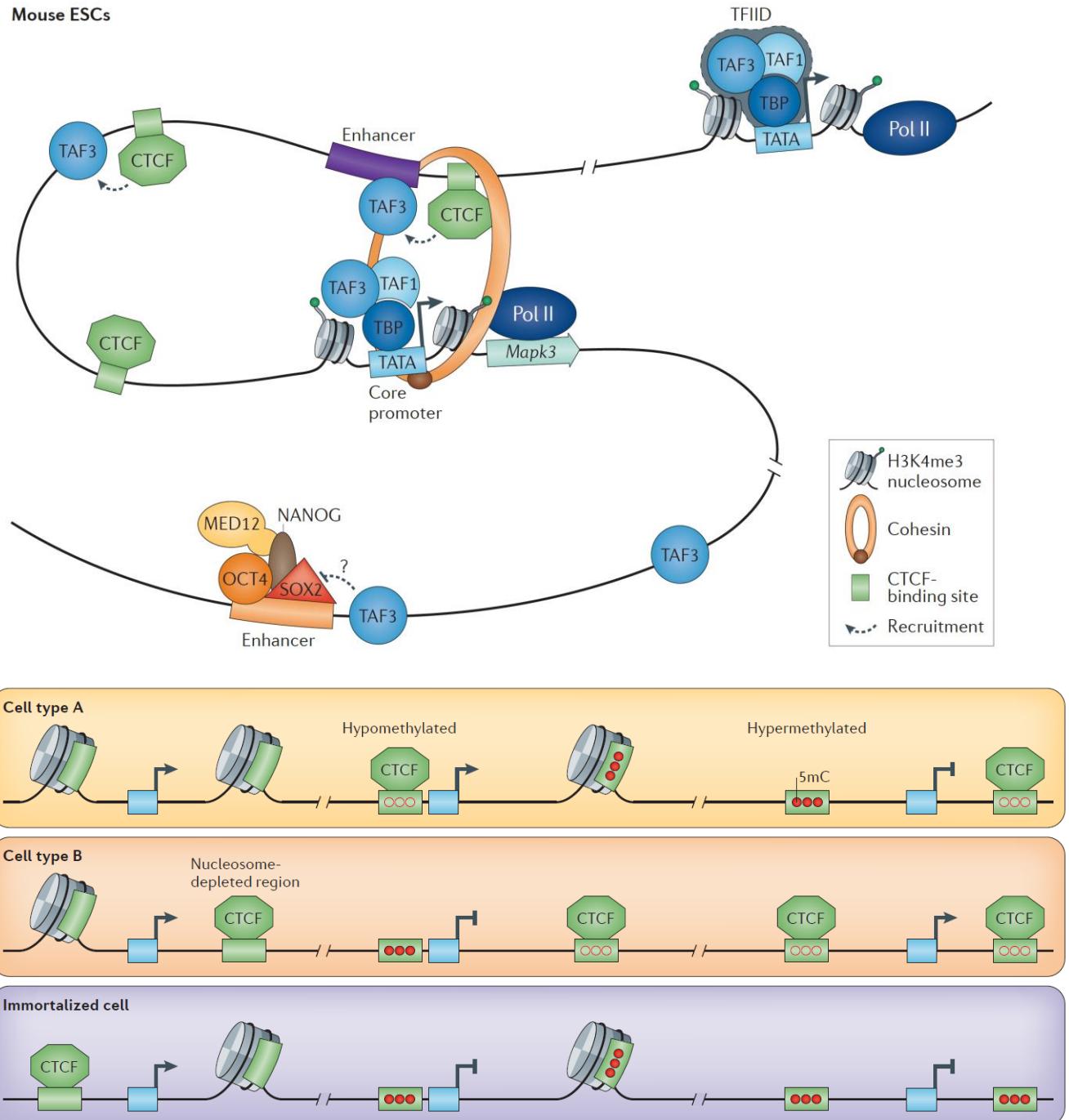
# Topologically Associated Domains

- TADs are regions that keep most physical interactions within
- separated by architectural proteins
- relatively conserved
- experimentally determined by Hi-C method



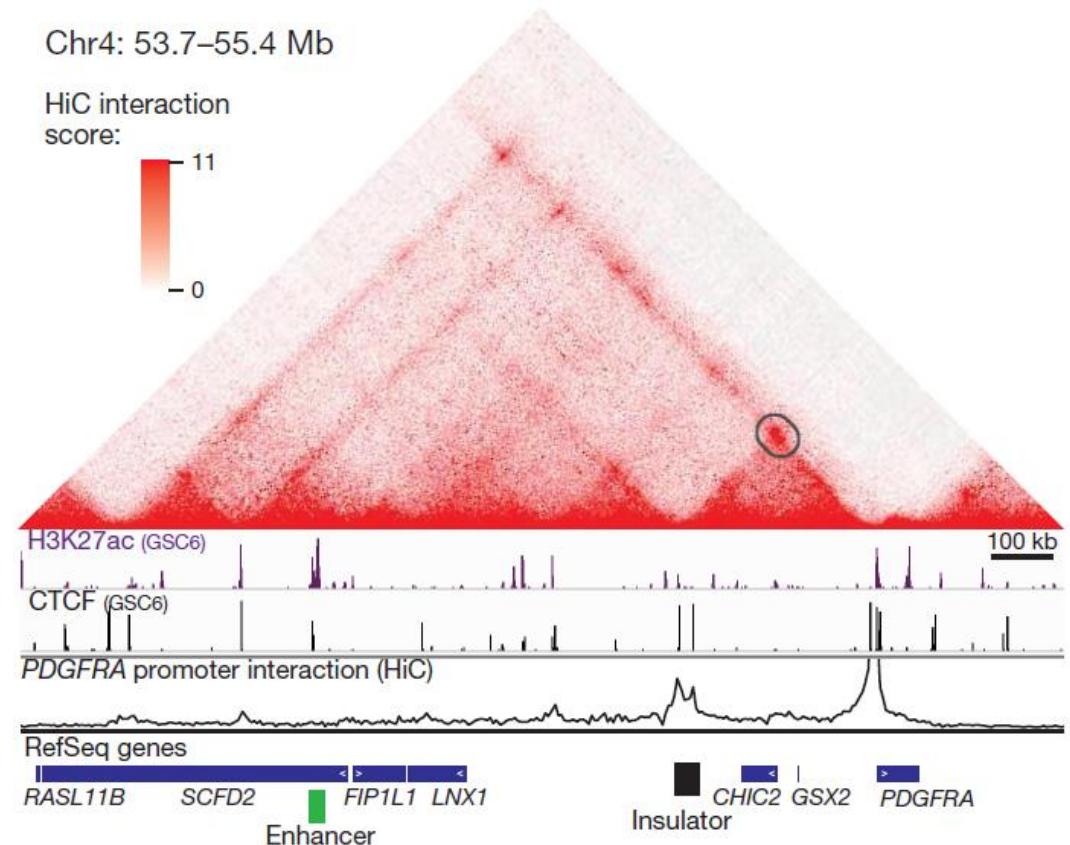
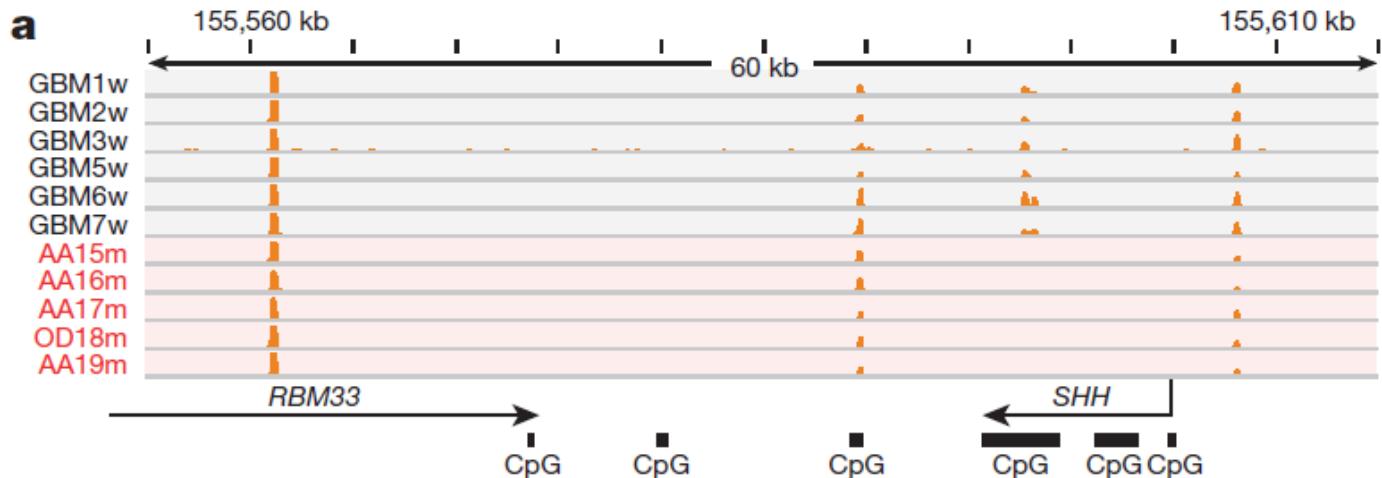
# CTCF and cohesin

- CTCF and cohesin are major architectural proteins
- enhancer-promoter contacts usually involves both
- CTCF binding is often controlled by DNA methylation



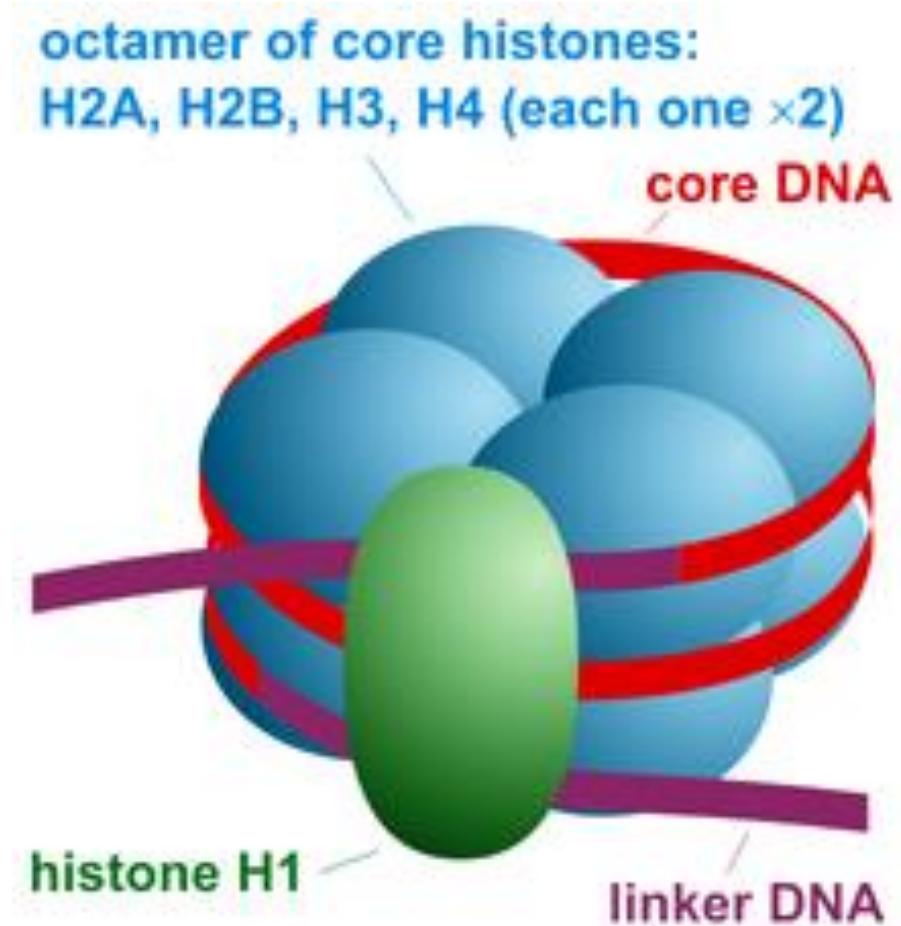
# TADs and glioma

- IDH gain-of-function mutations are known to be driver in gliomas
- Mutant IDH makes new metabolite (2-hydroxyglutarate) that inhibits TET enzymes that demethylate DNA
- methylation disrupts CTCF binding
- enhancers drive up PDGFRA oncogene



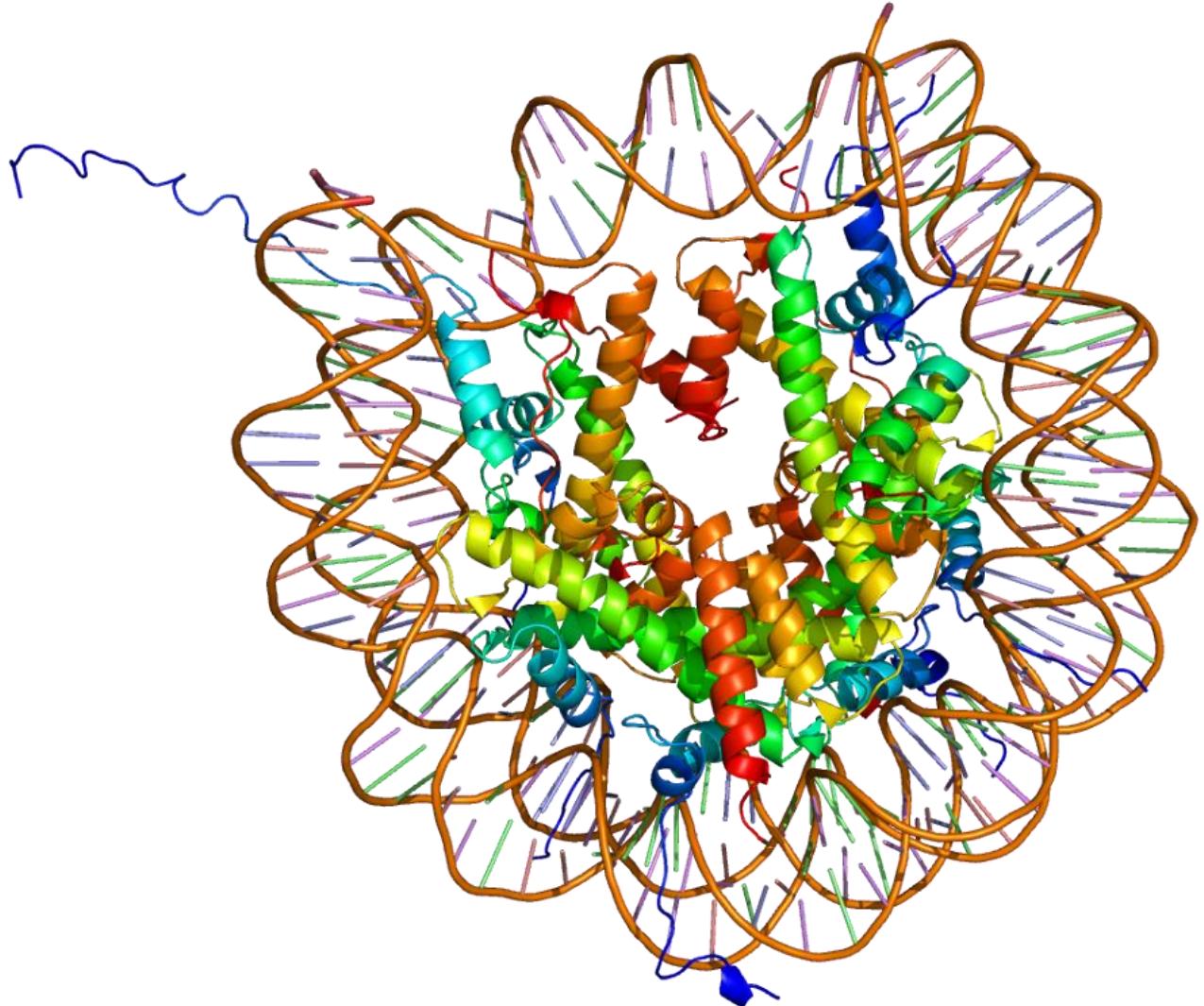
# Nucleosome

- nucleosomes are made of 8 histone proteins (2x H2A, H2B, H3, and H4) and include 147 nt of DNA
- H1 is a linker histone
- DNA is almost never histone-free, but can be **nucleosome-poor** (easy access, DNase I digestion)



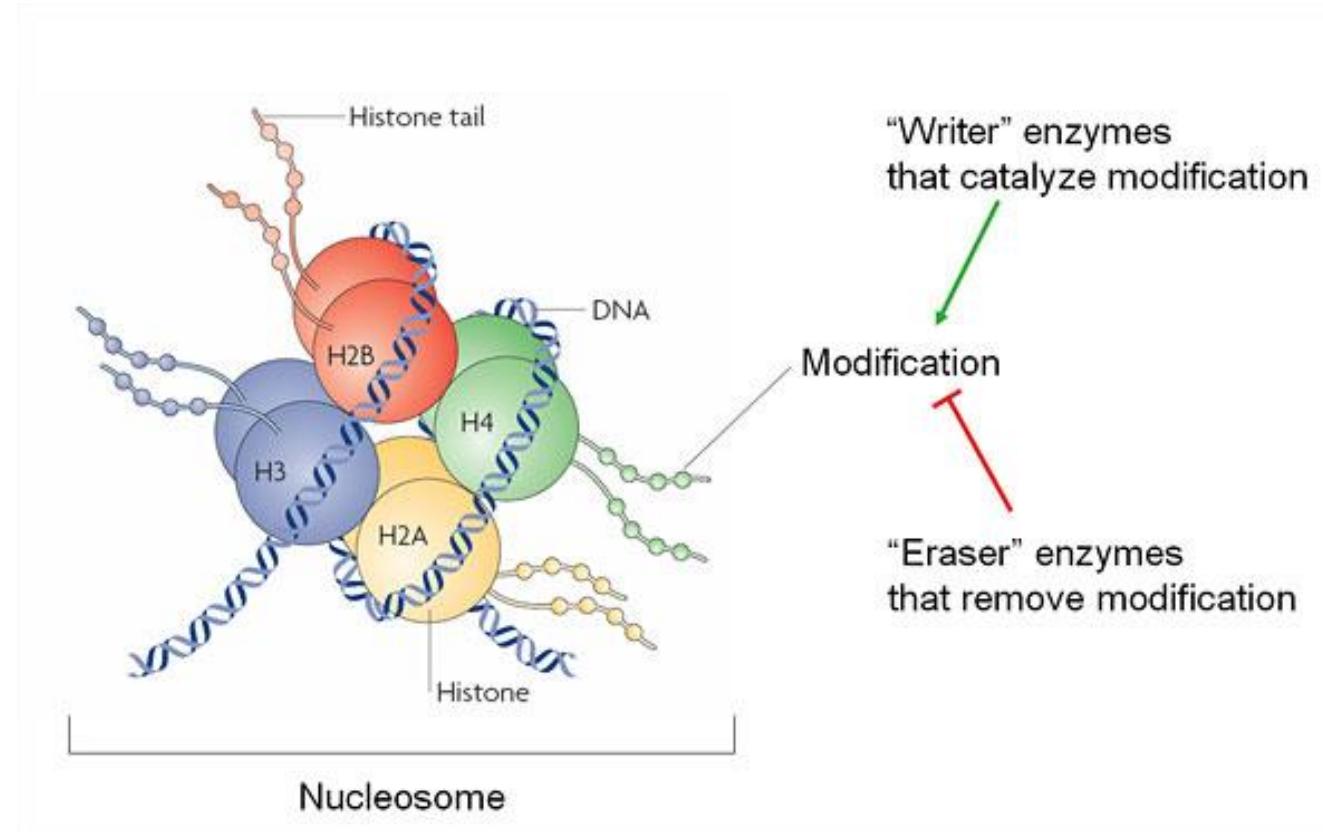
# Nucleosome

- DNA is negatively charged, so histones have positive charge to achieve **neutrality**  
- important for packing!



# Nucleosome

- histones tails stick outside and can be recognized
- chemical modifications of histones influence DNA accessibility
- histone modifications can be **read**, **erased**, and **recognized**



# Major protein components of transcription machinery

- transcription of regular protein-coding genes is done by **RNA polymerase II**
- proteins that attach to DNA are called **transcription factors** (TFs)



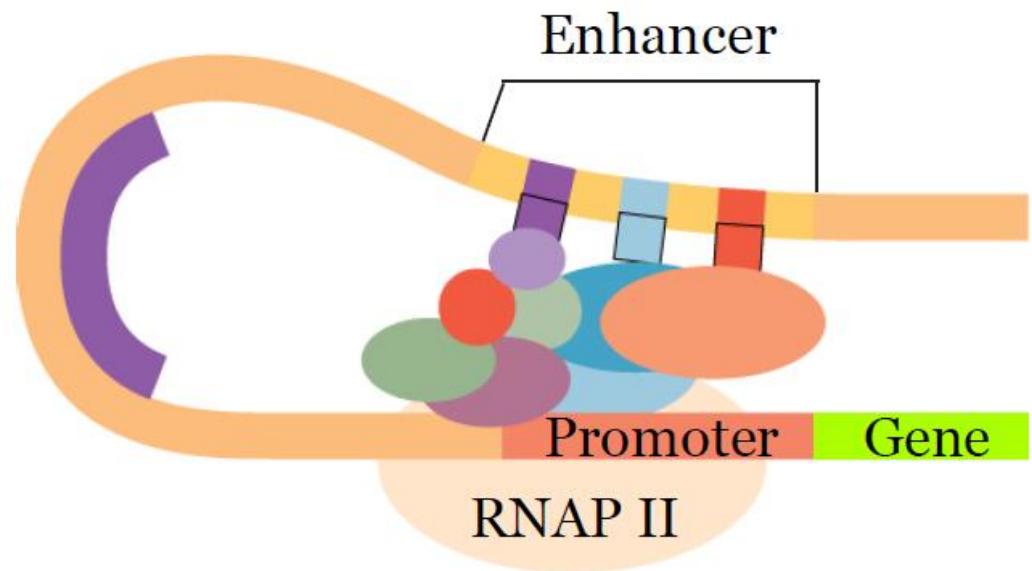
# Major protein components of transcription machinery

- transcription of regular protein-coding genes is done by **RNA polymerase II**
- proteins that attach to DNA are called **transcription factors** (TFs)
- there are common **general transcription factors** (GTFs - TFIIA-H) necessary for pre-initiation complex formation
- mediator complex, GTFs, and RNAPII combine into **holoenzyme** (100+ protein subunits!)



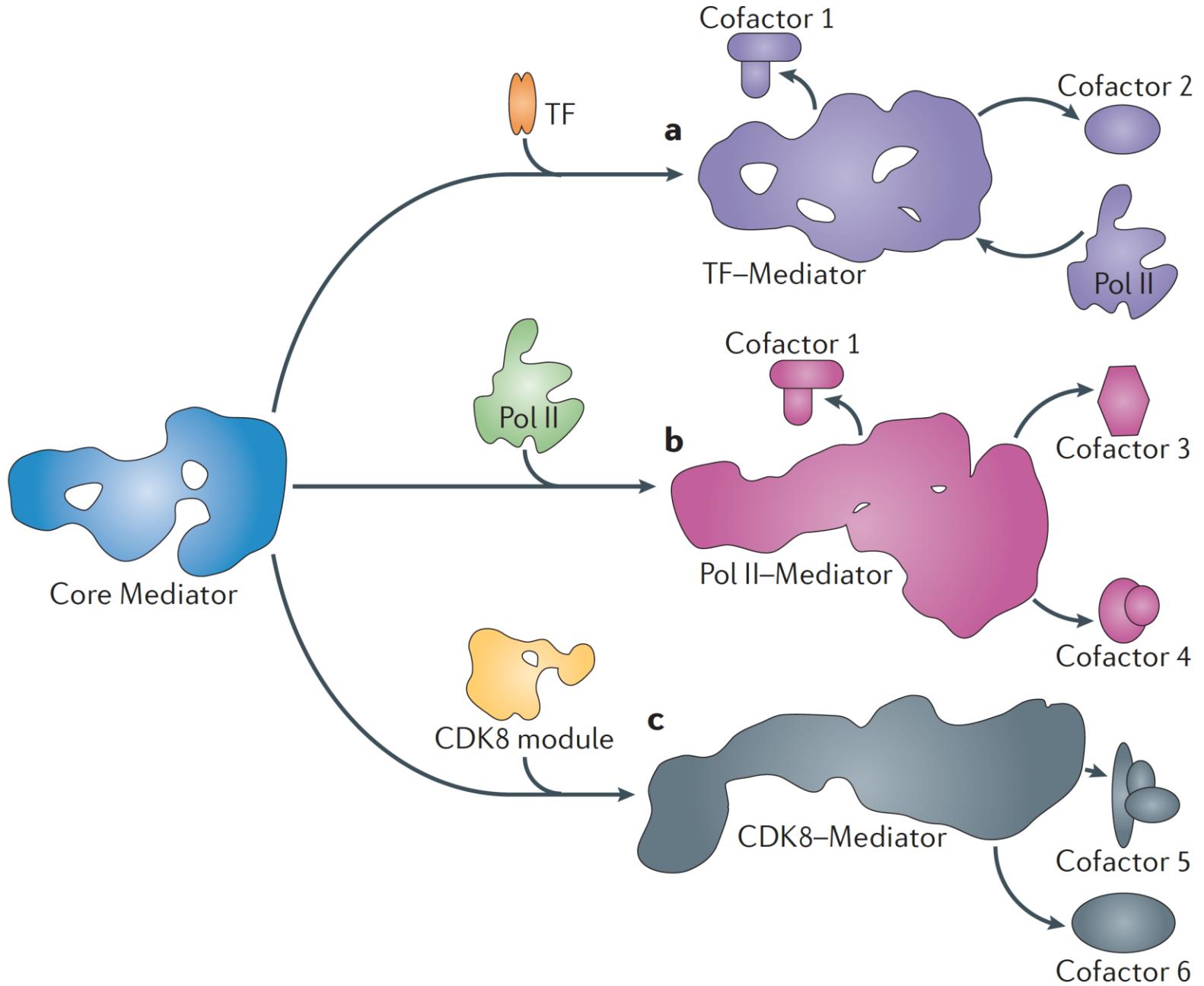
# Major genomic regulators of transcription

- basal transcription: general transcription factors bind the promoter and RNA polymerase II
- activator proteins bind DNA spots named enhancers
- enhancers are often located far and have to loop



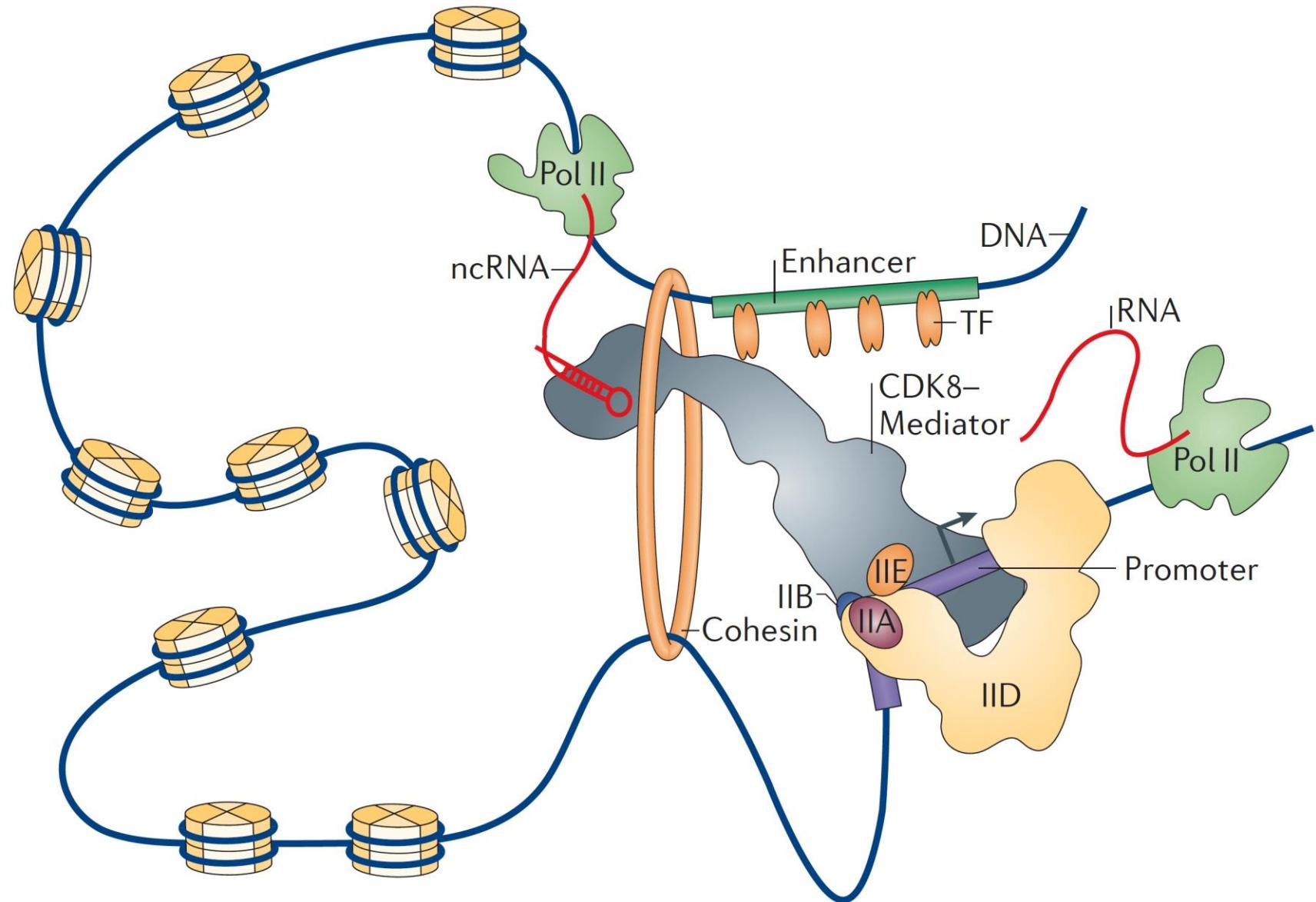
# Mediator complex

- many intrinsically disordered regions
- complexity is thought to be an evolutionary advantage
- mediator is involved in most active transcription



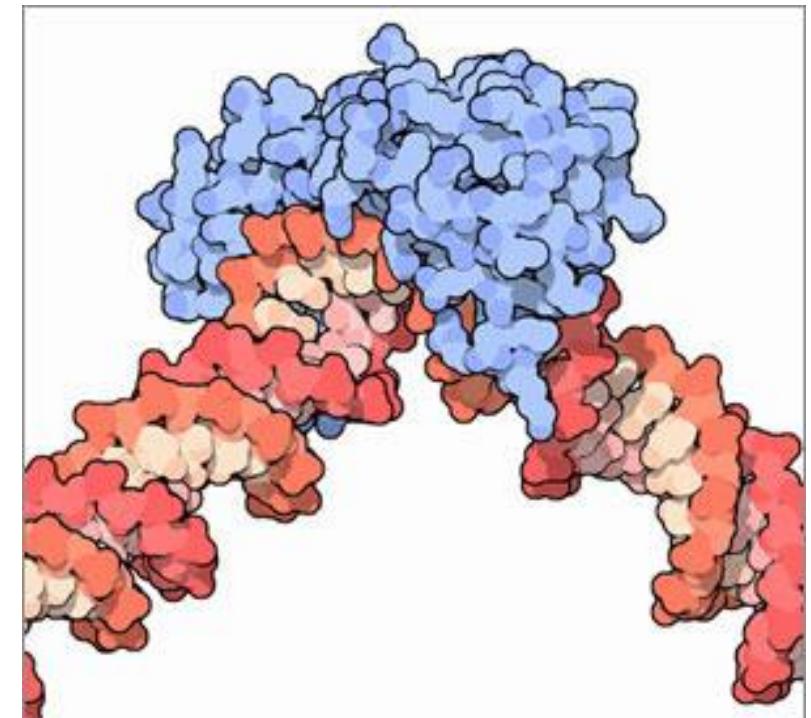
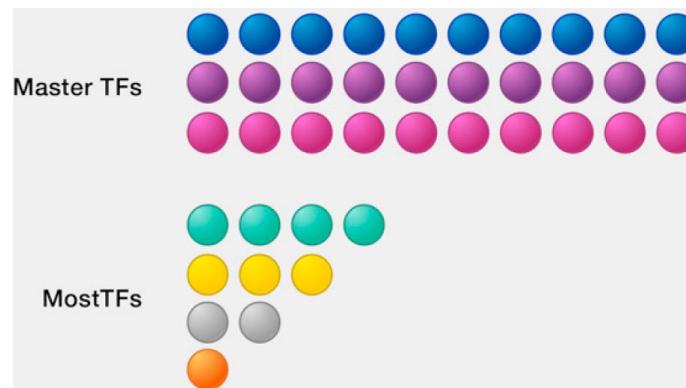
# Mediator and loops

- mediator does not facilitate long range enhancer-promoter loops in yeast



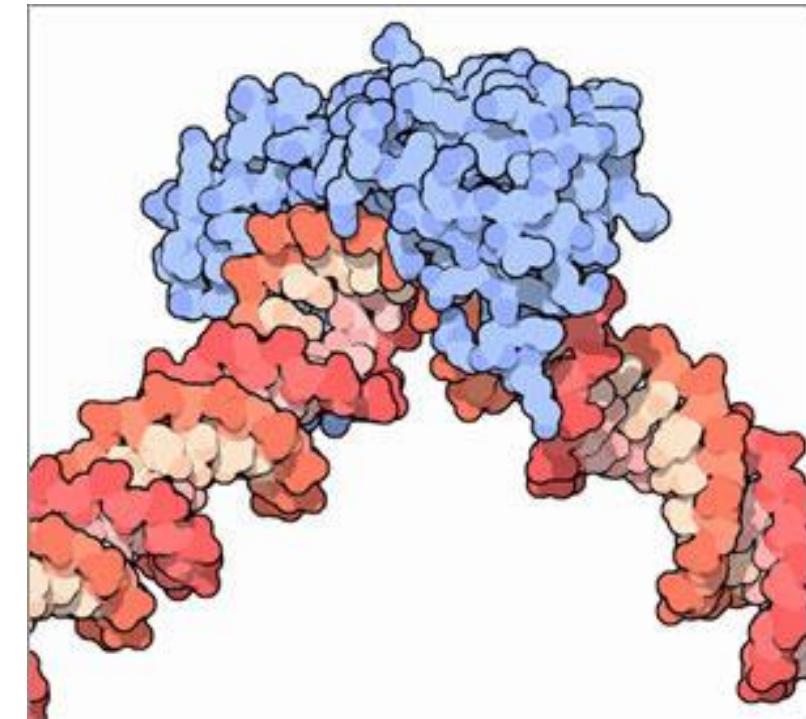
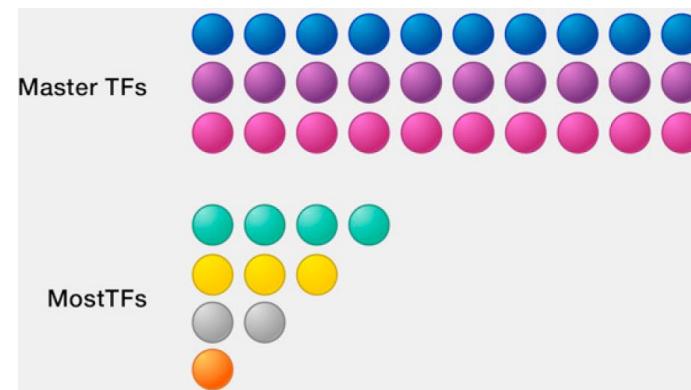
# Transcription factors

- there are estimated 1,500 transcription factors in humans
- **binding motif** represented by consensus sequence



# Transcription factors

- there are estimated 1,500 transcription factors in humans
- binding motif represented by consensus sequence
- **master regulators** exist but are not always known
- mechanisms of expression modulation involve:
  - stabilize/block RNAP II binding to DNA
  - catalyze histone acetylation or deacetylation
  - recruit coactivator or corepressor



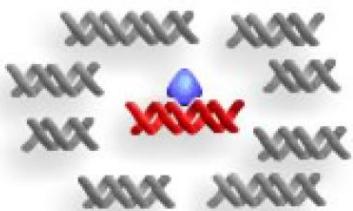
# Chromatin Immunoprecipitation (ChIP)

- simple but effective:  
reversible chemical cross-linking of proteins to DNA

Add formaldehyde to **chemically bind** proteins and DNA



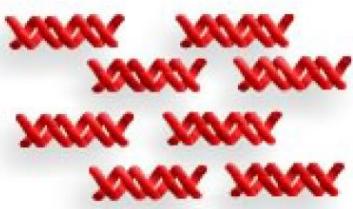
**Break** DNA into small fragments  
(usually ultrasound)



Precipitate fragments with  
the target protein  
using **selective antibody**



Amplify the DNA  
using PCR  
and **sequence it**



# Who came first?



Science, 2007 Jun 8;316(5830):1497-502. Epub 2007 May 31.

## Genome-wide mapping of *in vivo* protein-DNA interactions.

Johnson DS<sup>1</sup>, Mortazavi A, Myers RM, Wold B.

### Author information

#### Abstract

In vivo protein-DNA interactions connect each transcription factor with its direct targets to form a gene network scaffold. To map these protein-DNA interactions comprehensively across entire mammalian genomes, we developed a large-scale chromatin immunoprecipitation assay (ChIPSeq) based on direct ultrahigh-throughput DNA sequencing. This sequence census method was then used to map *in vivo* binding of the neuron-restrictive silencer factor (NRSF; also known as REST, for repressor element-1 silencing transcription factor) to 1946 locations in the human genome. The data display sharp resolution of binding position [~10 base pairs (bp)], which facilitated our finding motifs and allowed us to identify noncanonical NRSF-binding motifs. These ChIPSeq data also have high sensitivity and specificity [ROC (receiver operator characteristic) area >= 0.96] and statistical confidence ( $P < 10^{-4}$ ), properties that were important for inferring new candidate interactions. These include key transcription factors in the gene network that regulates pancreatic islet cell development.

Cell

Volume 129, Issue 4, 18 May 2007, Pages 823–837

Cell  
PRESS

#### Resource

## High-Resolution Profiling of Histone Methylation in the Human Genome

Artem Barski<sup>1,3</sup>, Suresh Cuddapah<sup>1,3</sup>, Kairong Cui<sup>1,3</sup>, Tae-Young Roh<sup>1,3</sup>, Dustin E. Schones<sup>1,3</sup>, Zhibin Wang<sup>1,3</sup>, Gang Wei<sup>1,3</sup>, Iouri Chepelev<sup>2</sup>, Keji Zhao<sup>1</sup>,

nature

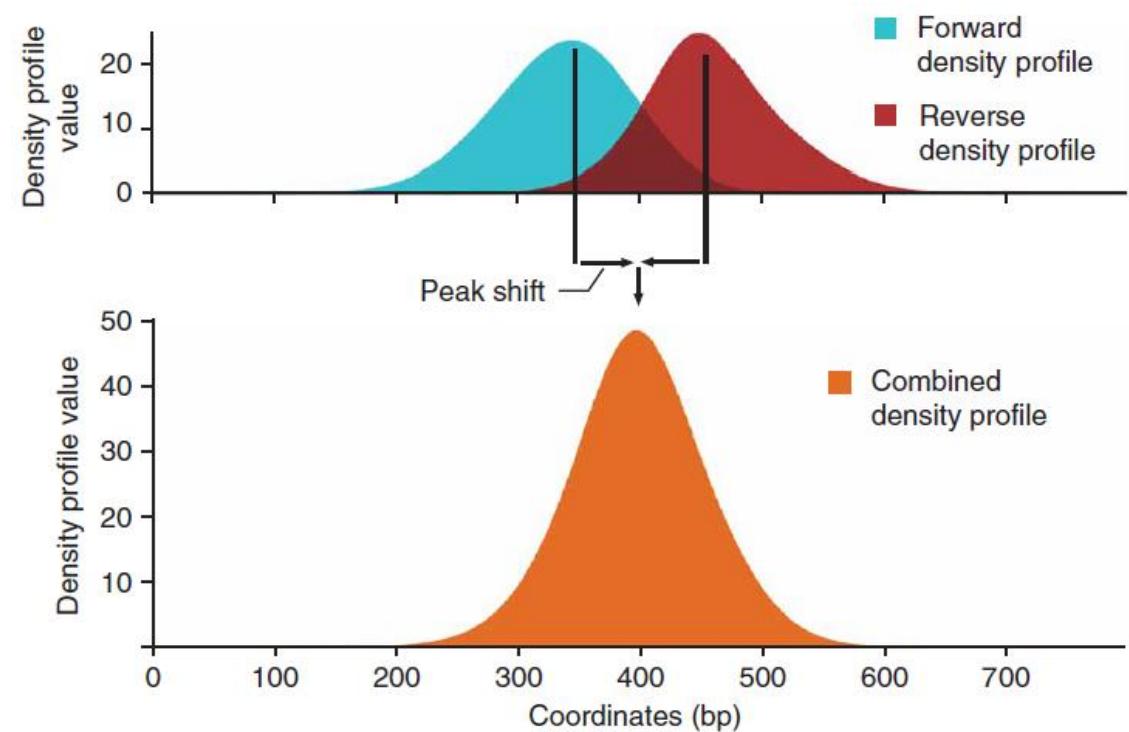
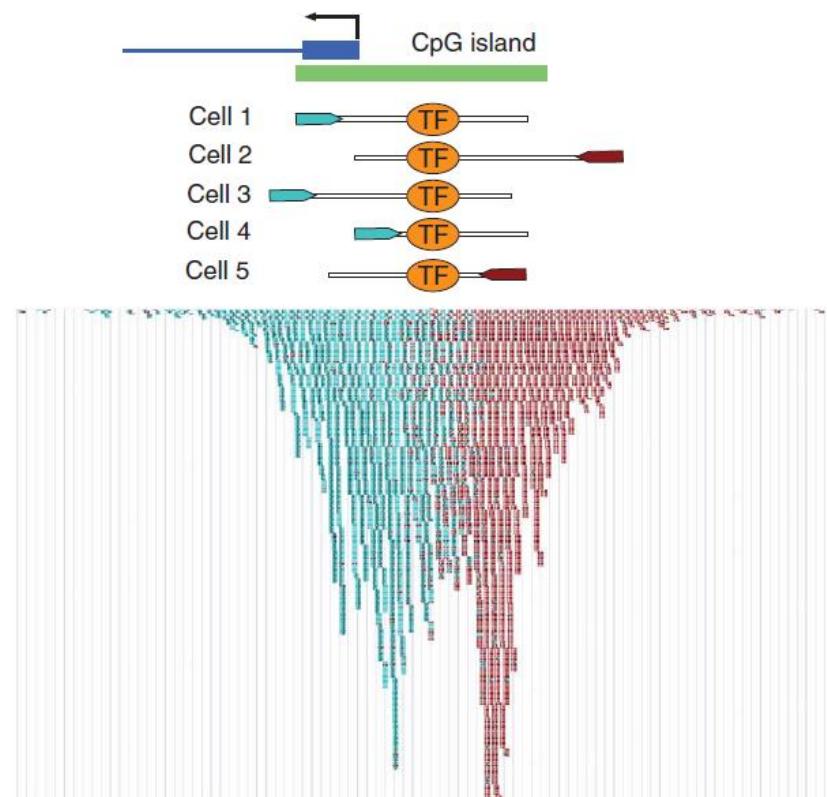
ARTICLES

## Genome-wide maps of chromatin state in pluripotent and lineage-committed cells

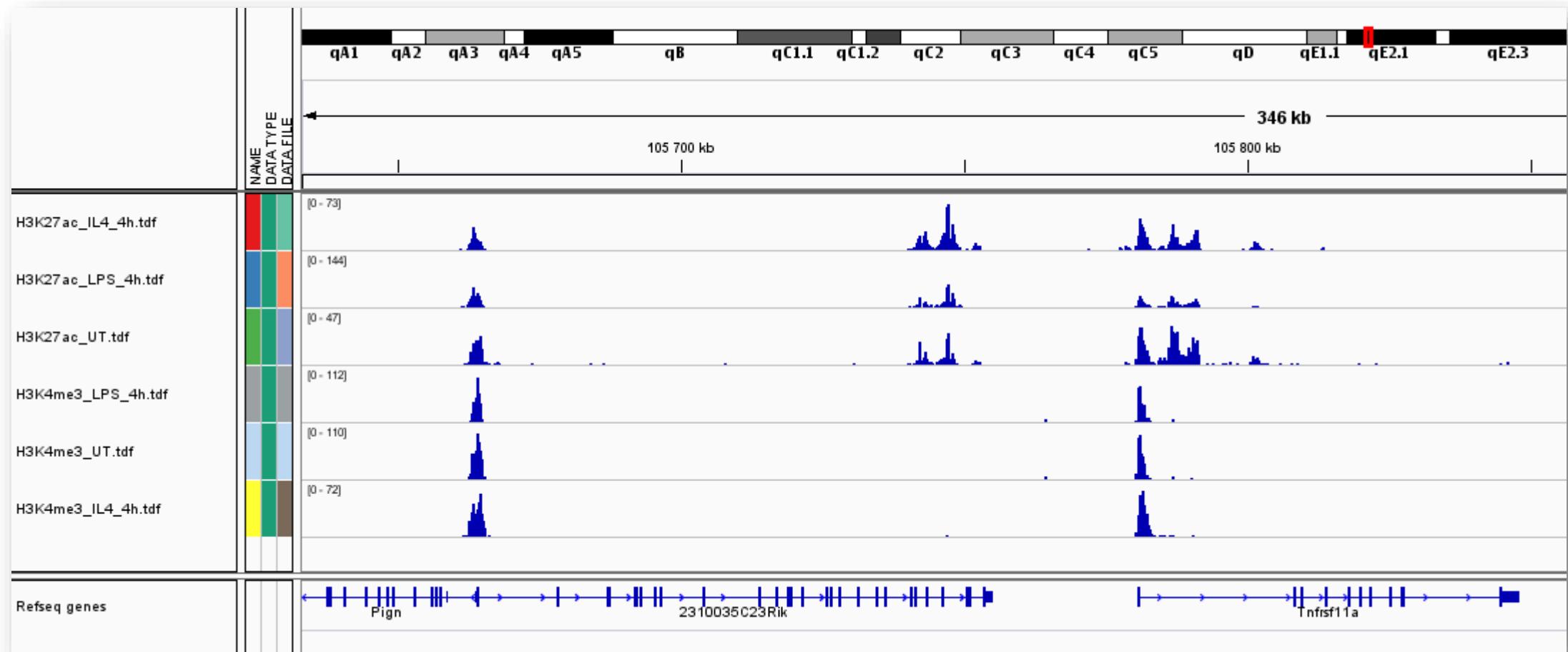
Tarjei S. Mikkelsen<sup>1,2</sup>, Manching Ku<sup>1,4</sup>, David B. Jaffe<sup>1</sup>, Biju Issac<sup>1,4</sup>, Erez Lieberman<sup>1,2</sup>, Georgia Giannoukos<sup>1</sup>, Pablo Alvarez<sup>1</sup>, William Brockman<sup>1</sup>, Tae-Kyung Kim<sup>5</sup>, Richard P. Koche<sup>1,2,4</sup>, William Lee<sup>1</sup>, Eric Mendenhall<sup>1,4</sup>, Aisling O'Donovan<sup>4</sup>, Aviva Presser<sup>1</sup>, Carsten Russ<sup>1</sup>, Xiaohui Xie<sup>1</sup>, Alexander Meissner<sup>3</sup>, Marius Wernig<sup>3</sup>, Rudolf Jaenisch<sup>3</sup>, Chad Nusbaum<sup>1</sup>, Eric S. Lander<sup>1,3\*</sup> & Bradley E. Bernstein<sup>1,4,6\*</sup>

# ChIP-Seq

- DNA library obtained after ChIP can be amplified and **sequenced**
- ChIP-Seq can be used for both **transcription factors** and **histone modifications**, like H3K4me3



# A typical ChIP-seq experiment

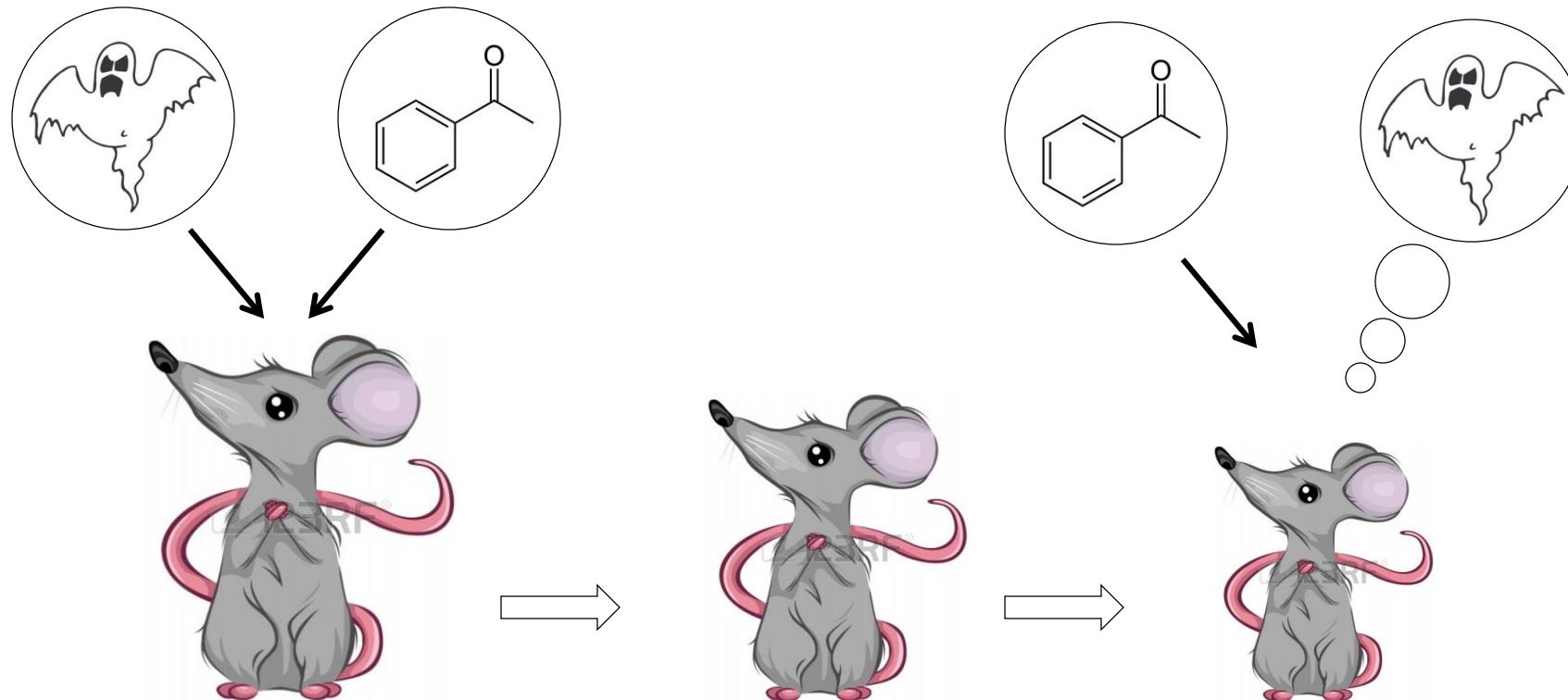


# Epigenetics

- Two major definitions that mean very different things:
  - Heritable changes that do not involve DNA modification (still heavily debated)
  - DNA modification and packaging changes that influence expression (well established)
- Very different mechanisms and meaning in
  - differentiating cells
  - pluripotent cells
  - differentiated cells

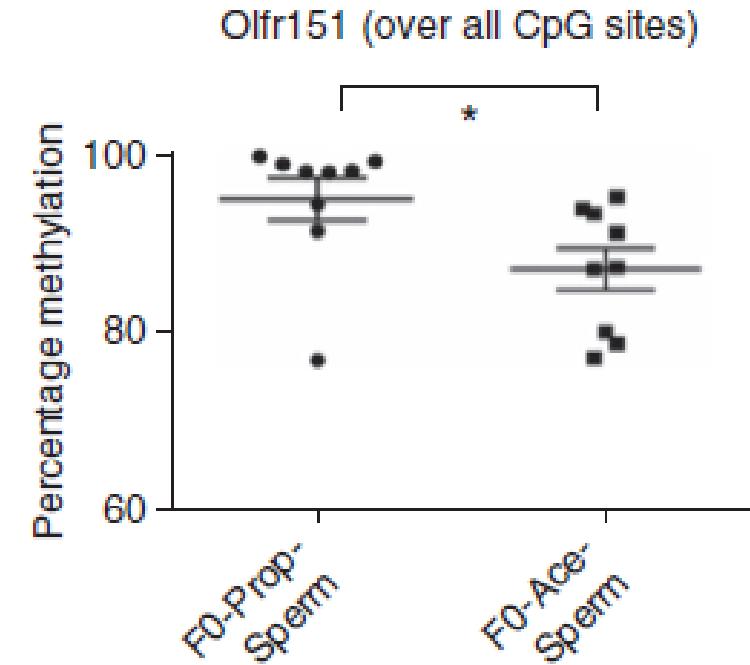
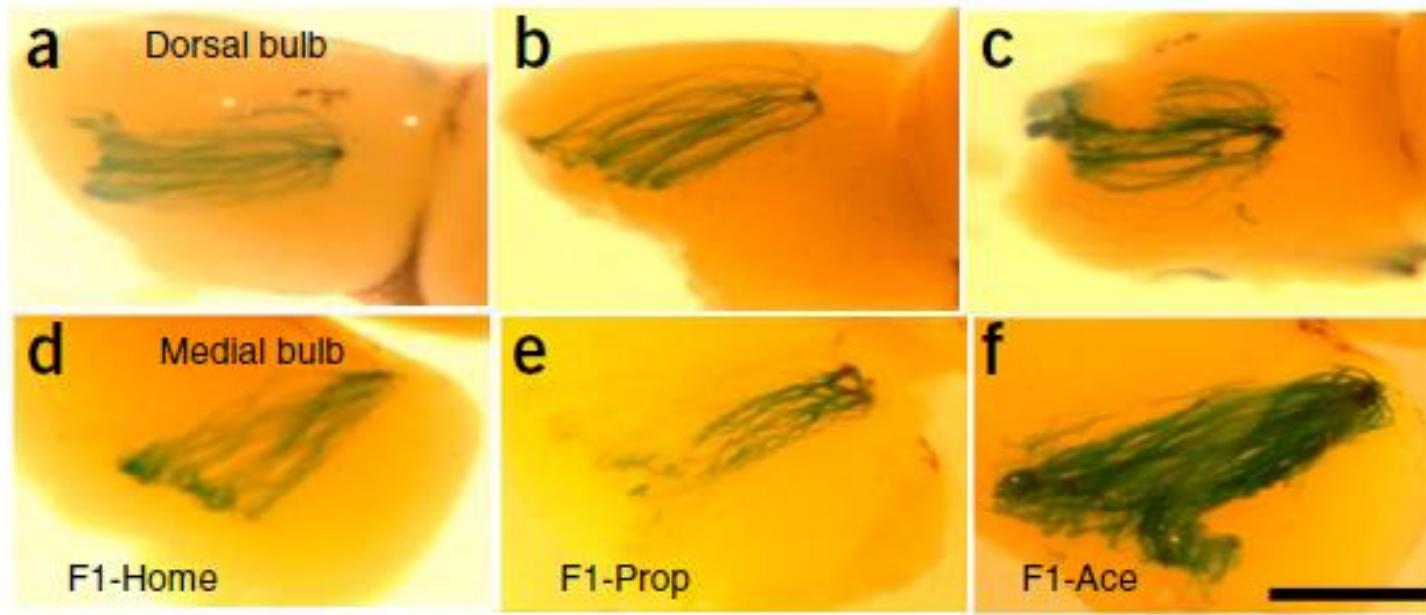
# Fear and loathing in Mouse-Vegas

- odor fear conditioning (acetophenone)
- the effect lasted 2 generations



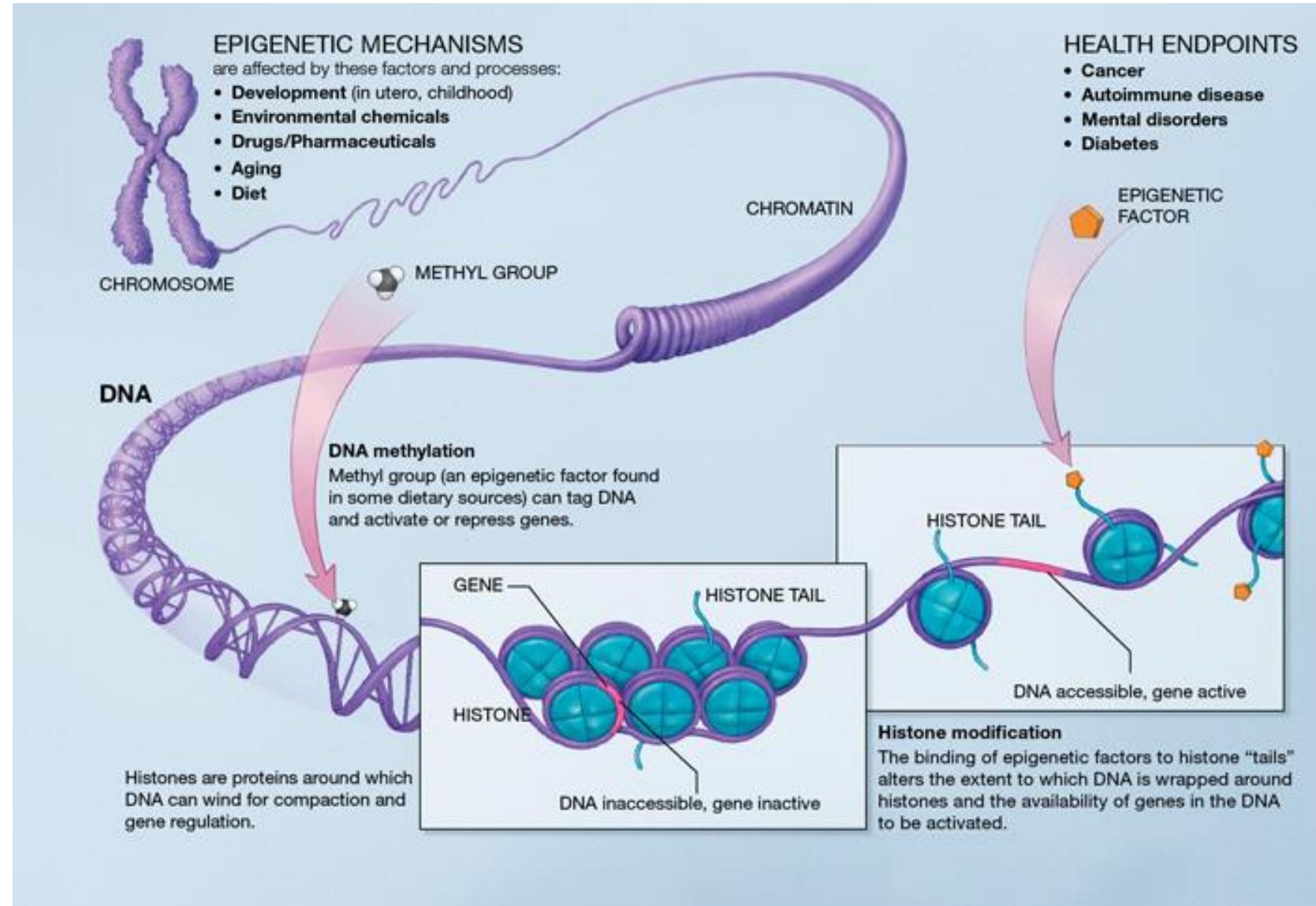
# Mechanism: *Olf151* promoter hypomethylation

- activates *Olf151* that expresses olfactory glomerulus M71
- bisulfite sequencing revealed that Olf151 was hypomethylated at CpG promoter



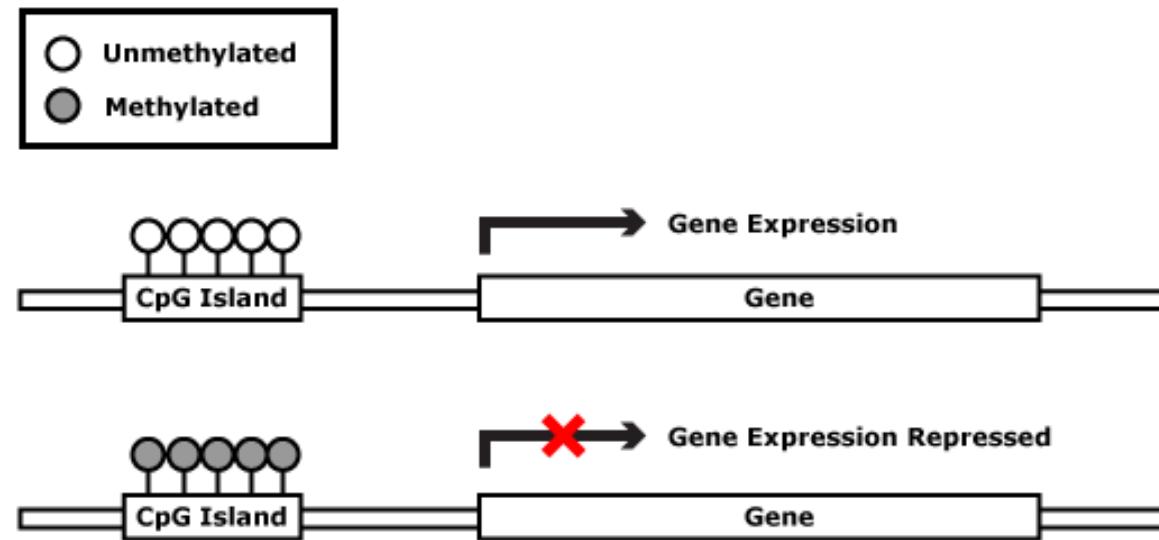
# Epigenetic regulation

- histone modifications
- DNA methylation
- noncoding RNA

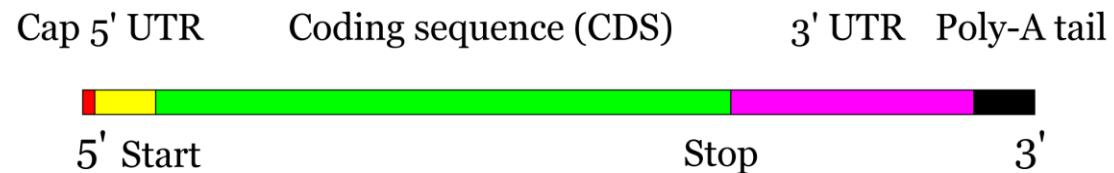


# Promoters and CpG islands

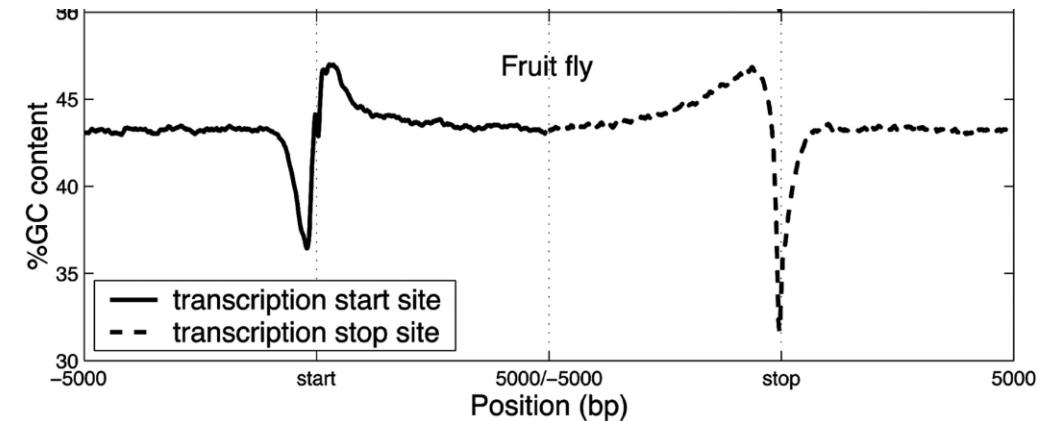
- many human promoters are **CpG rich** (two types of promoters?)
- in the rest of the genome GC mutates to GT
- methylation of the promoter CpG silences the gene



# Evolution: Human vs. Fruit Fly



- Drosophila is a "**Dnmt2 only**" organism
- no Dnmt1 and Dnmt3
- no functional homologs of known 5-methylcytosine reader proteins

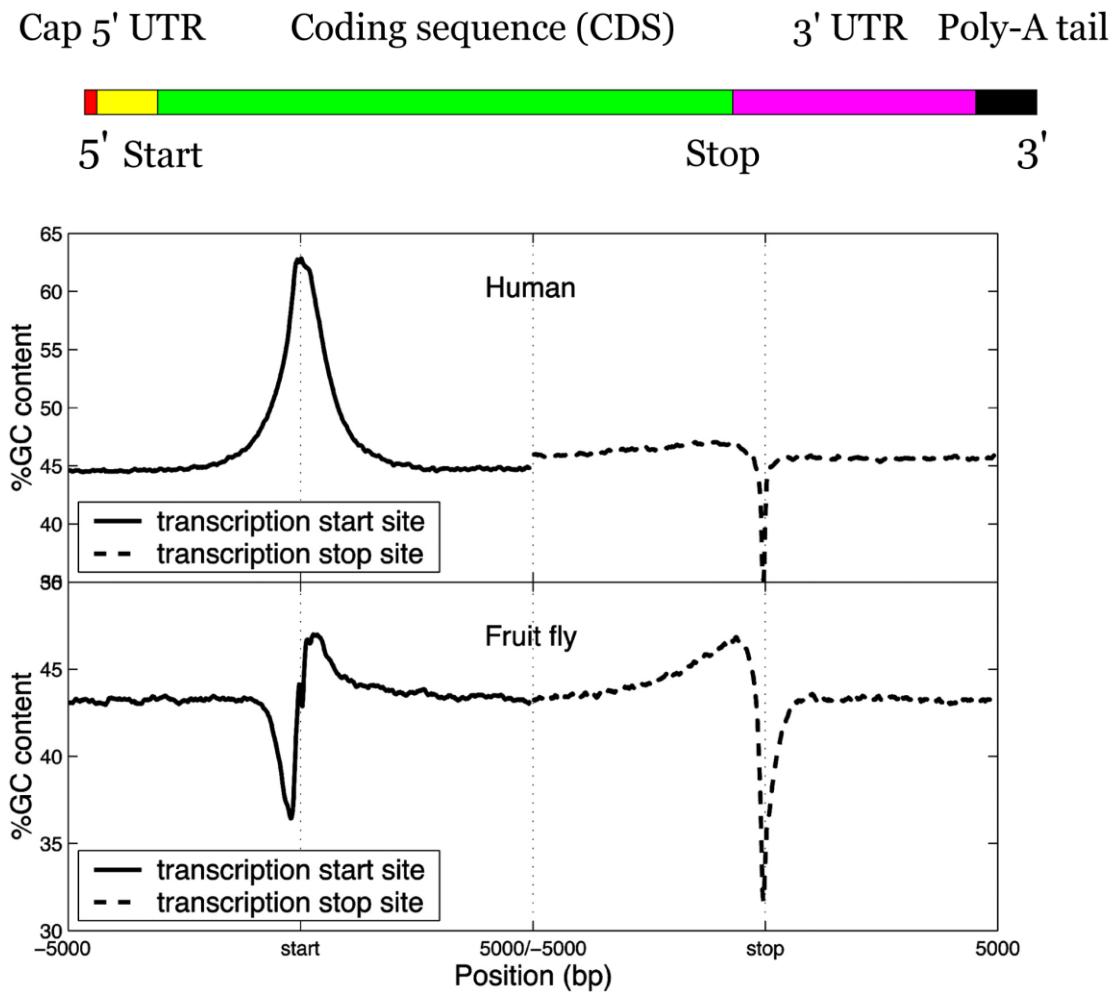


[http://en.wikipedia.org/wiki/Messenger\\_RNA](http://en.wikipedia.org/wiki/Messenger_RNA)

Zhang L. et al, Proc Nat Acad Sci USA, 101, pp 16855-16860 (2004)

# Evolution: Human vs. Fruit Fly

- Drosophila is a "**Dnmt2 only**" organism
- no Dnmt1 and Dnmt3
- no functional homologs of known 5-methylcytosine reader proteins
- human promoter GC% indicates pronounced conservation

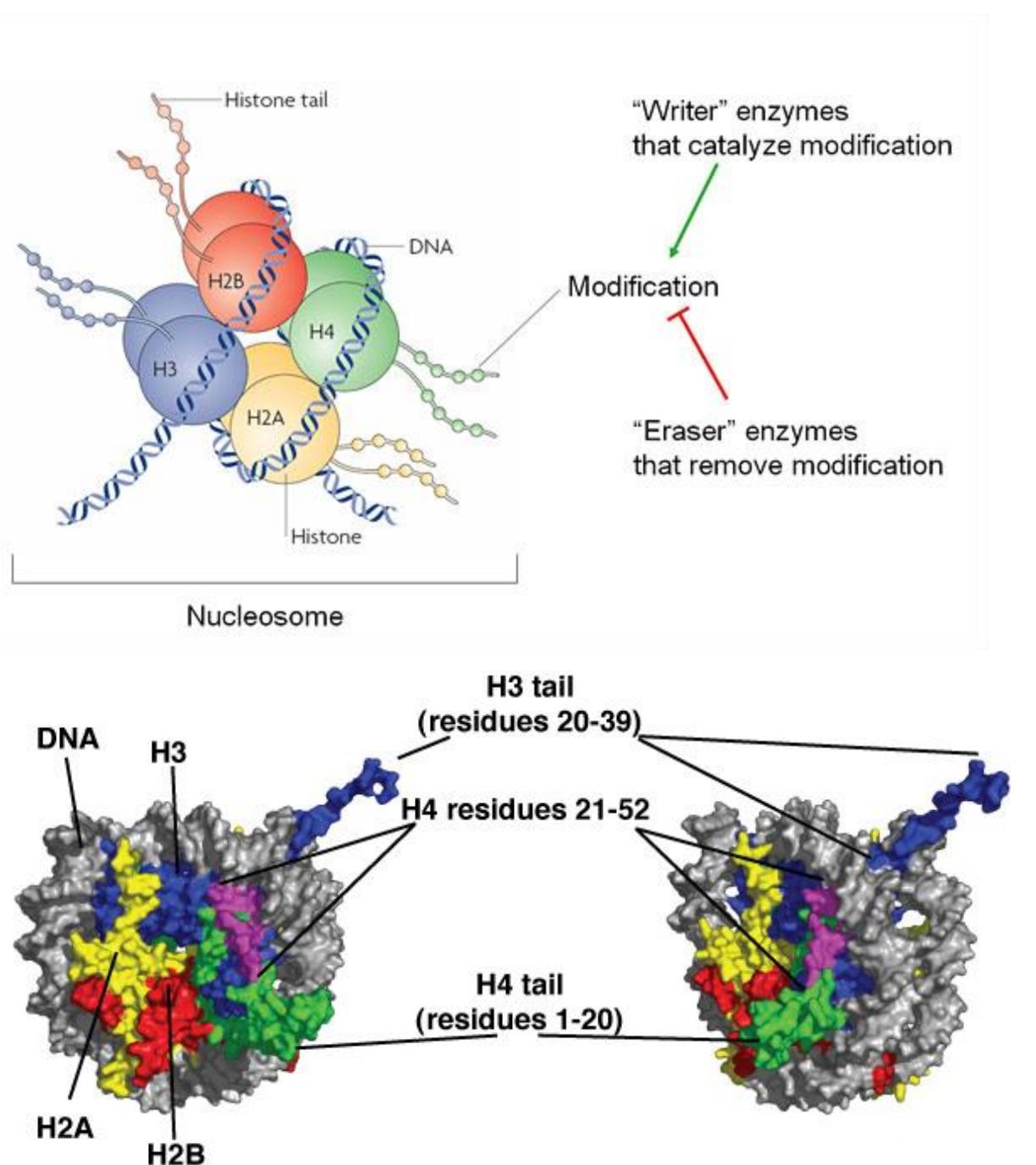


[http://en.wikipedia.org/wiki/Messenger\\_RNA](http://en.wikipedia.org/wiki/Messenger_RNA)

Zhang L. et al, Proc Nat Acad Sci USA, 101, pp 16855-16860 (2004)

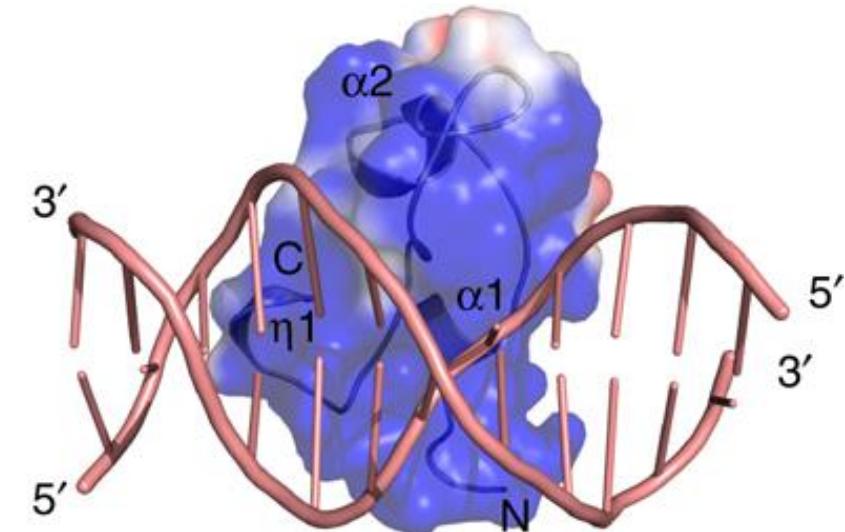
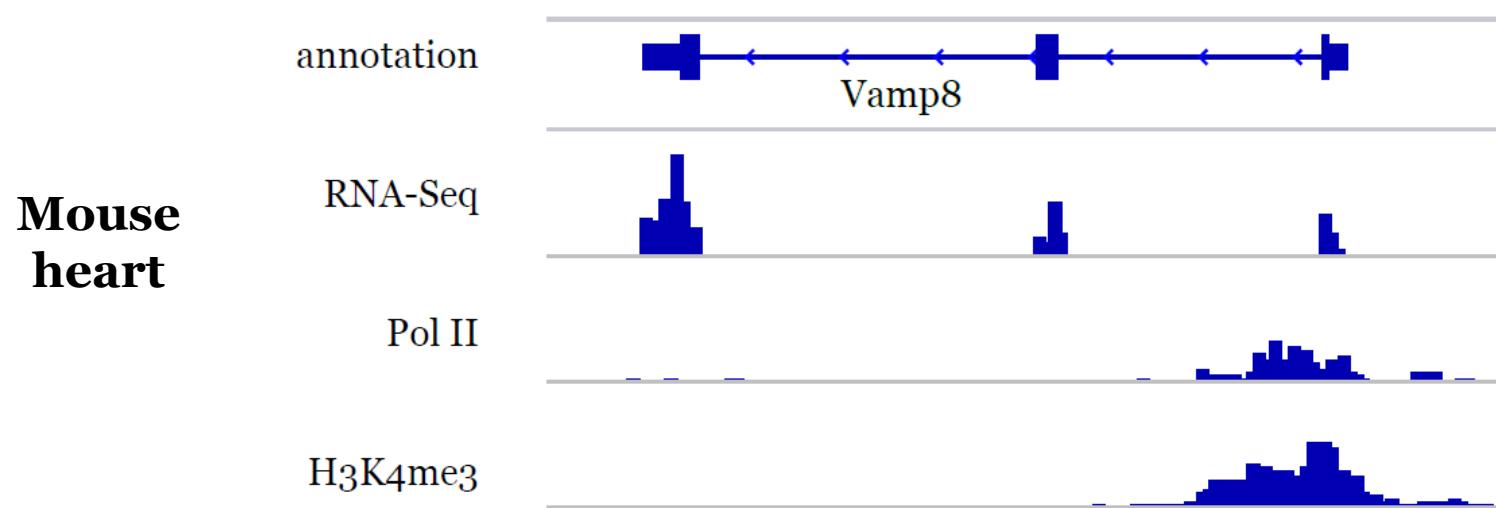
# Histone tail modifications

- histone tail can be chemically modified
- are in fact normal post-translational modifications
- histone modifications change transcription



# Promoter histone marks

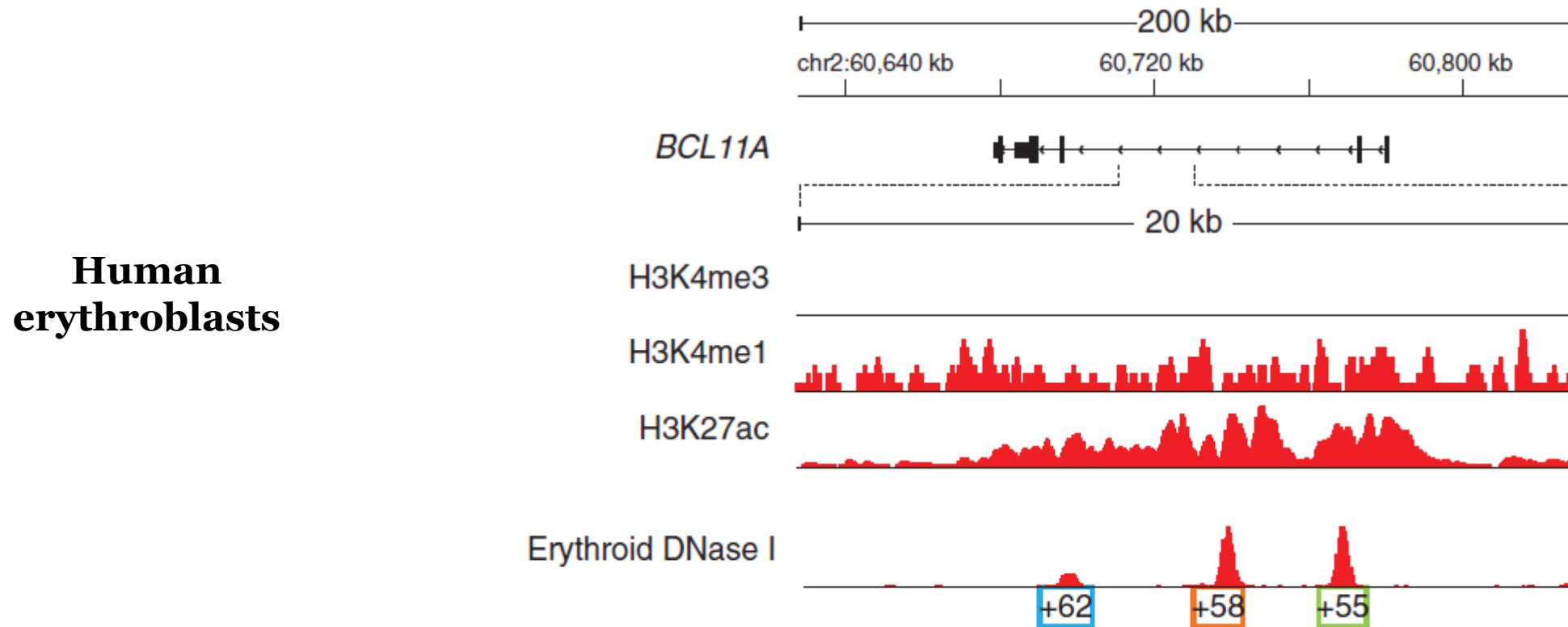
- narrow peaks of H3K4me3 mark promoters.
- enzyme that methylates K4 binds only to non-CpG-methylated promoters!



The EMBO Journal, 31, pp 3130–3146 (2012)  
C. Xu et al, Nature Communications , 2(227),pp 1-8 (2011)

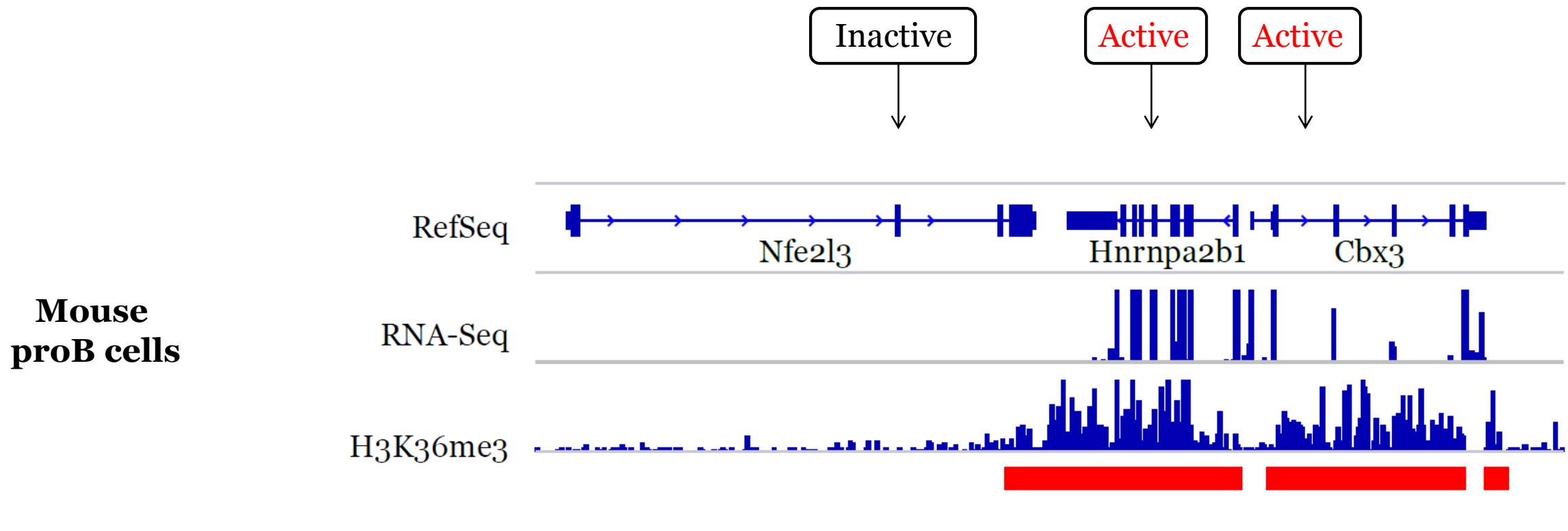
# Enhancer histone marks

- enhancers are associated with H3K4me1 and H3K27ac
- H3K27ac is thought to distinguish active enhancers from poised



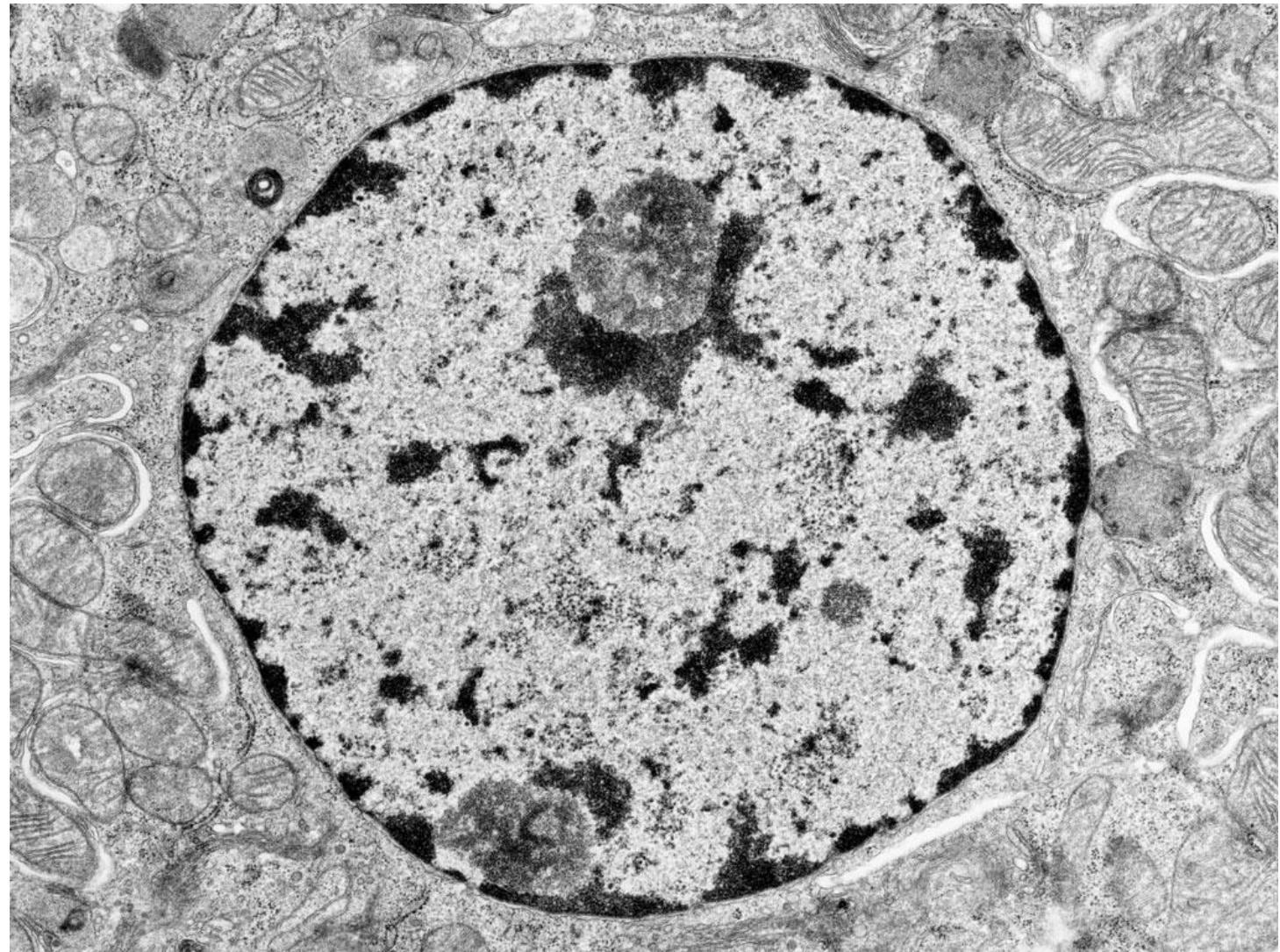
# Marks of transcriptional elongation

- Elongation is marked by H3K36me3 and H3K79me2



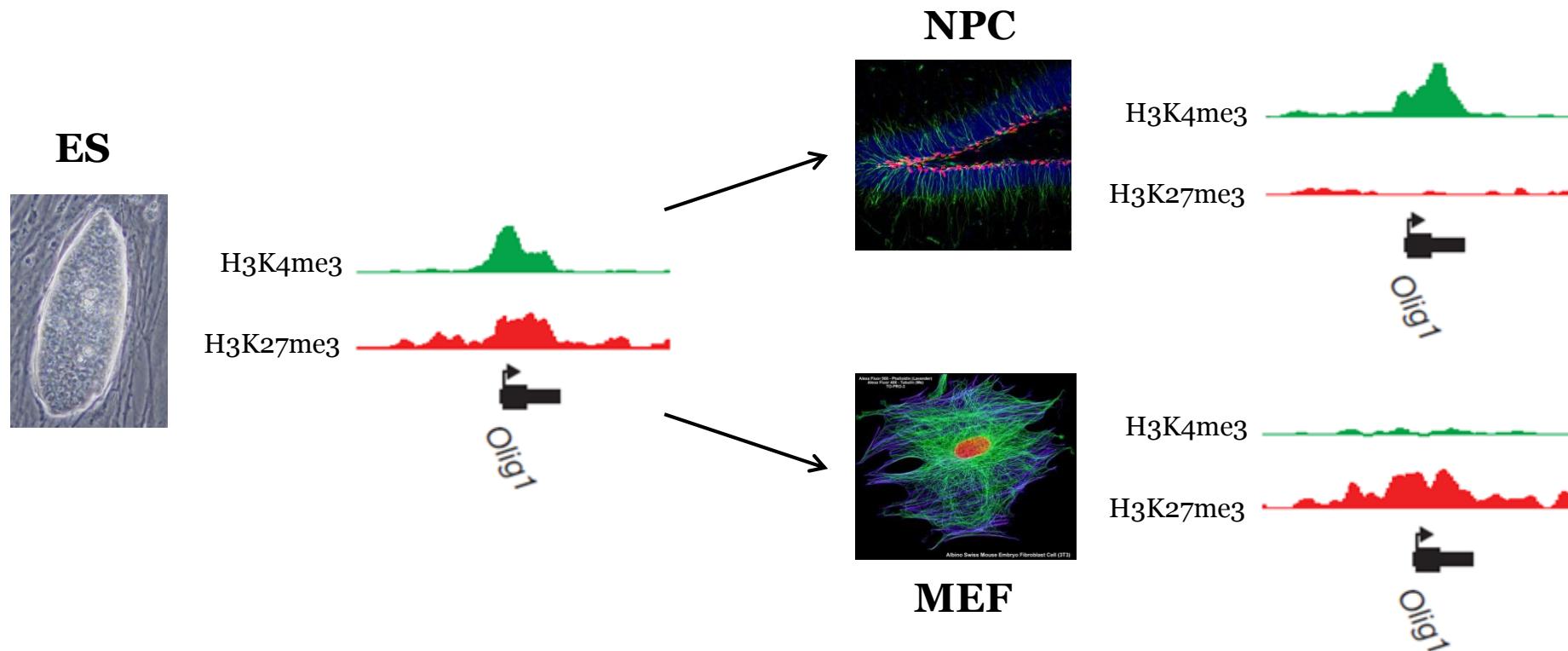
# Gene repression by chromatin

- Gene silencing can be **reversible** and (relatively) **irreversible**
- irreversible = heterochromatin
- reversible can be removed in the process of cell differentiation



# Reversible repression

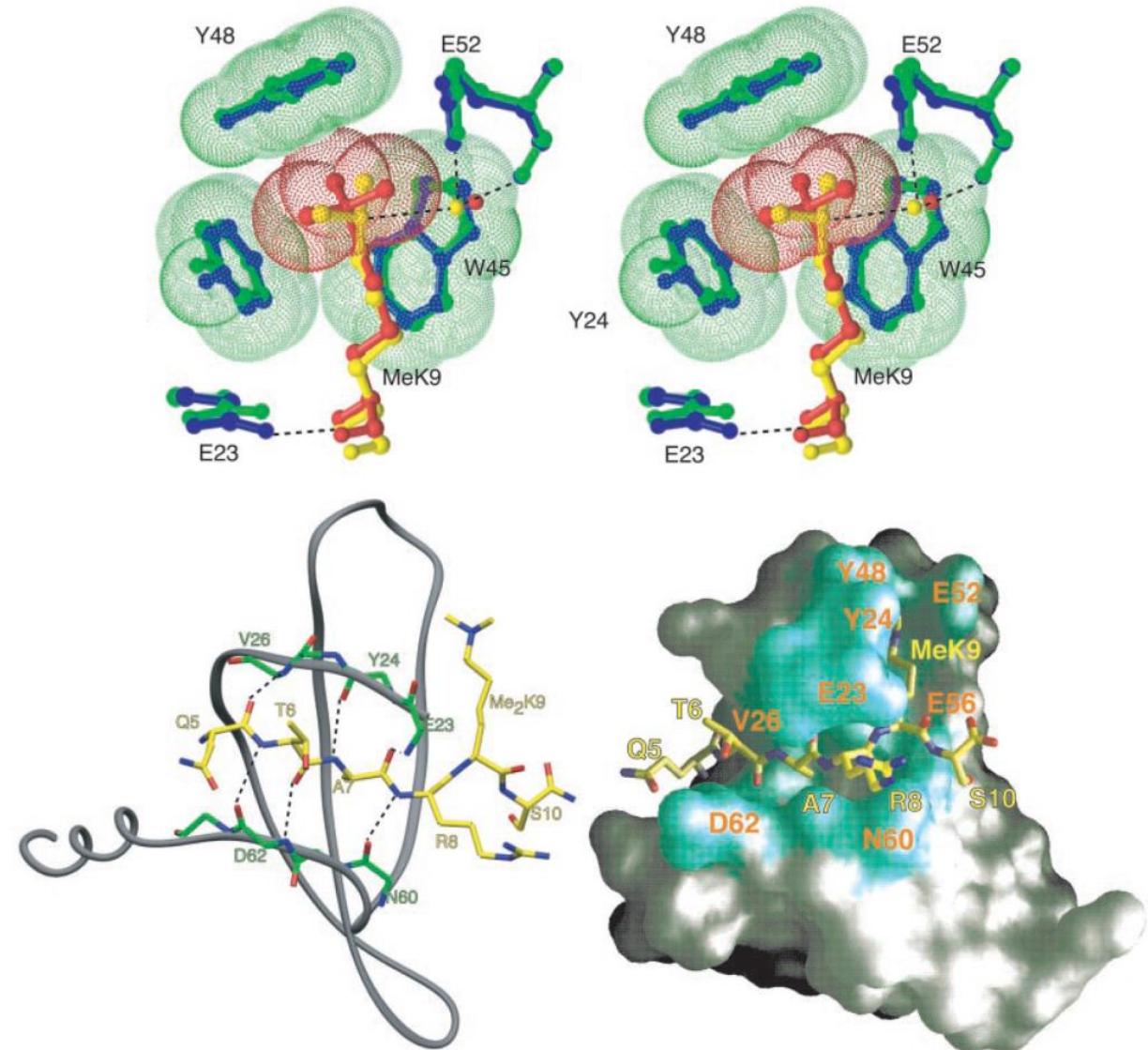
- H3K27me3 marks suppressed genes poised to be activated
- Stem cells can have both H3K4me3 and H3K27me3 – unique!



Mikkelsen, T.S.; Ku, M. et al, Nature 448, pp 553-560 (2007)  
Vastenhouw N.L.; Zhang Y. et al, Nature 464, pp 922-6 (2010)

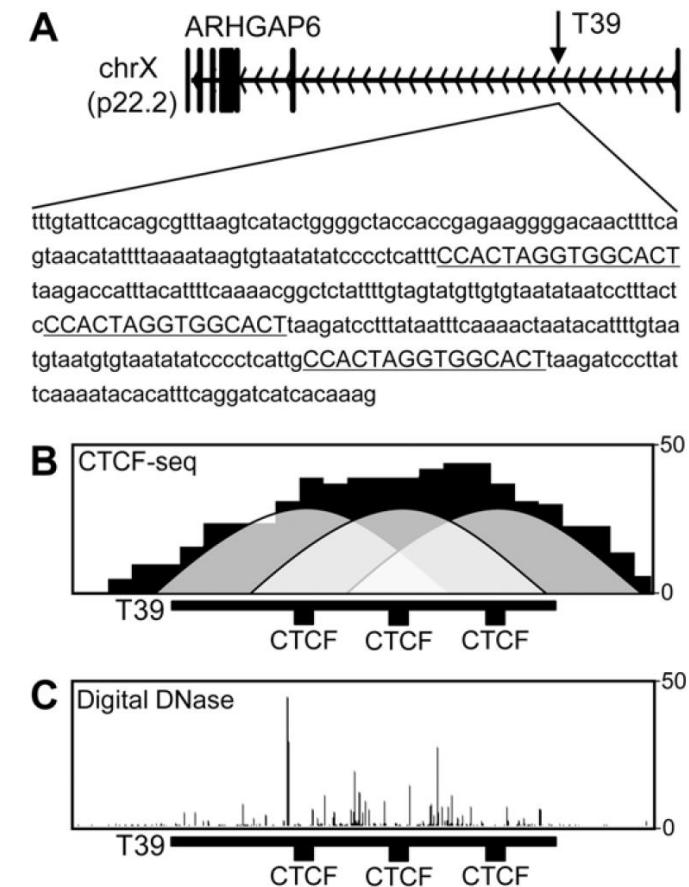
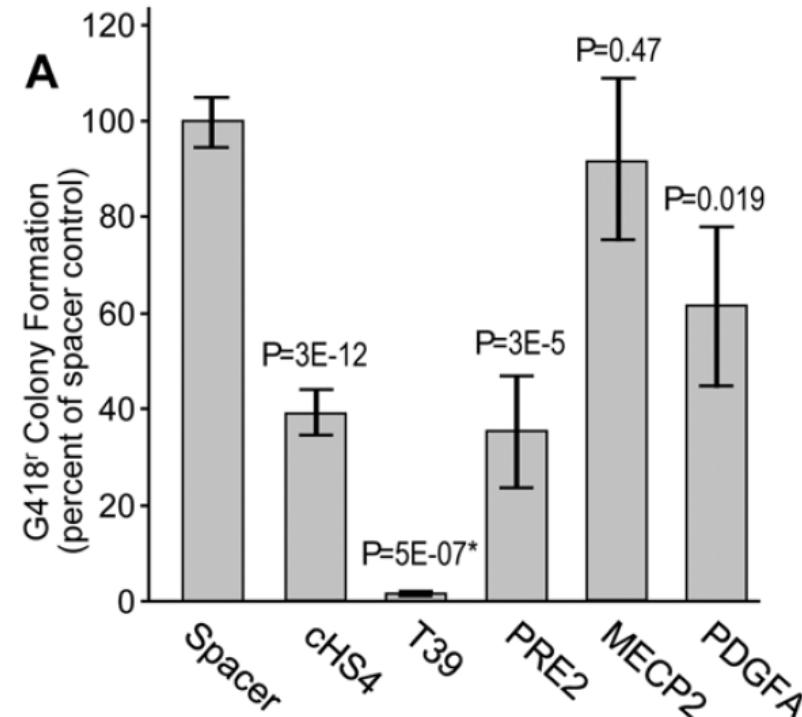
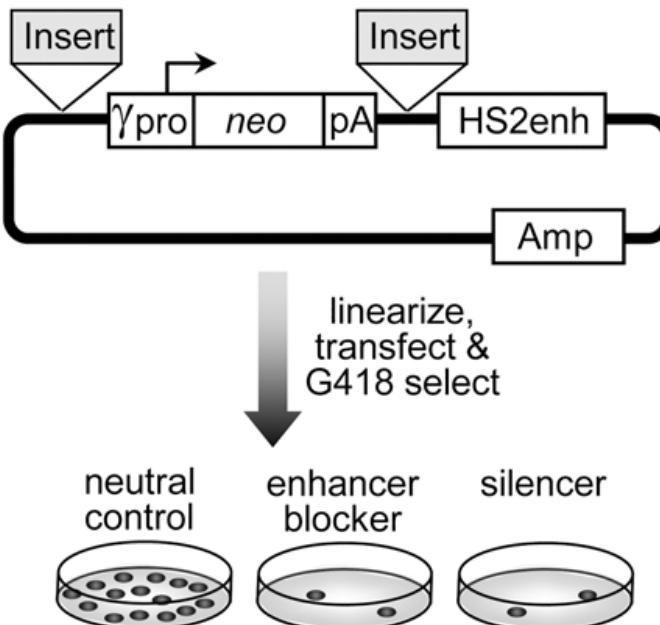
# Heterochromatin marks

- Marks **H3K9me2** and **H3K9me3** are strongly associated with heterochromatin
- Binding with protein HP1



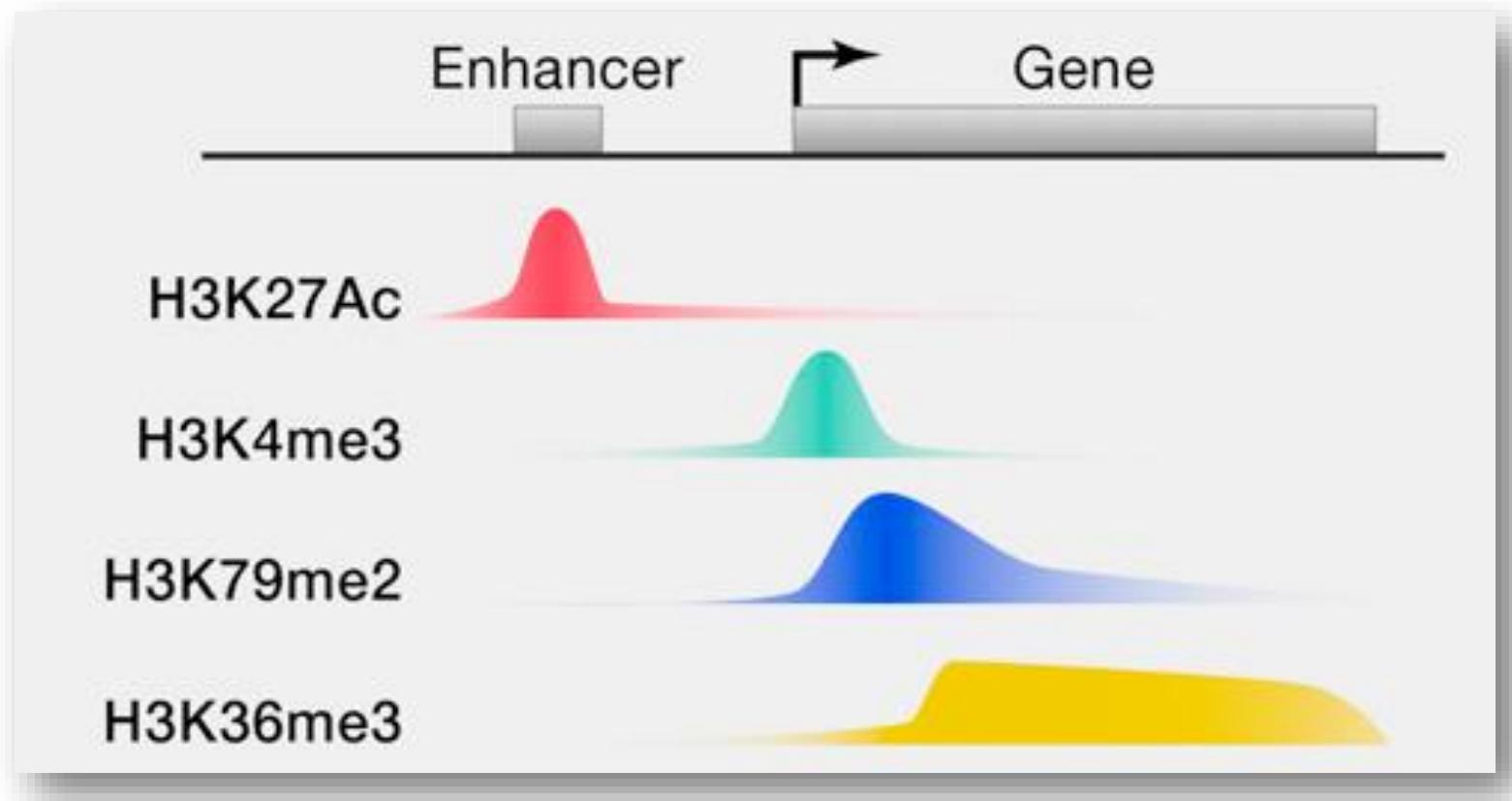
# Silencers in eukaryotes

- well-characterized silencers in eukaryotes are extremely rare
- most repression is thought to be done via chromatin

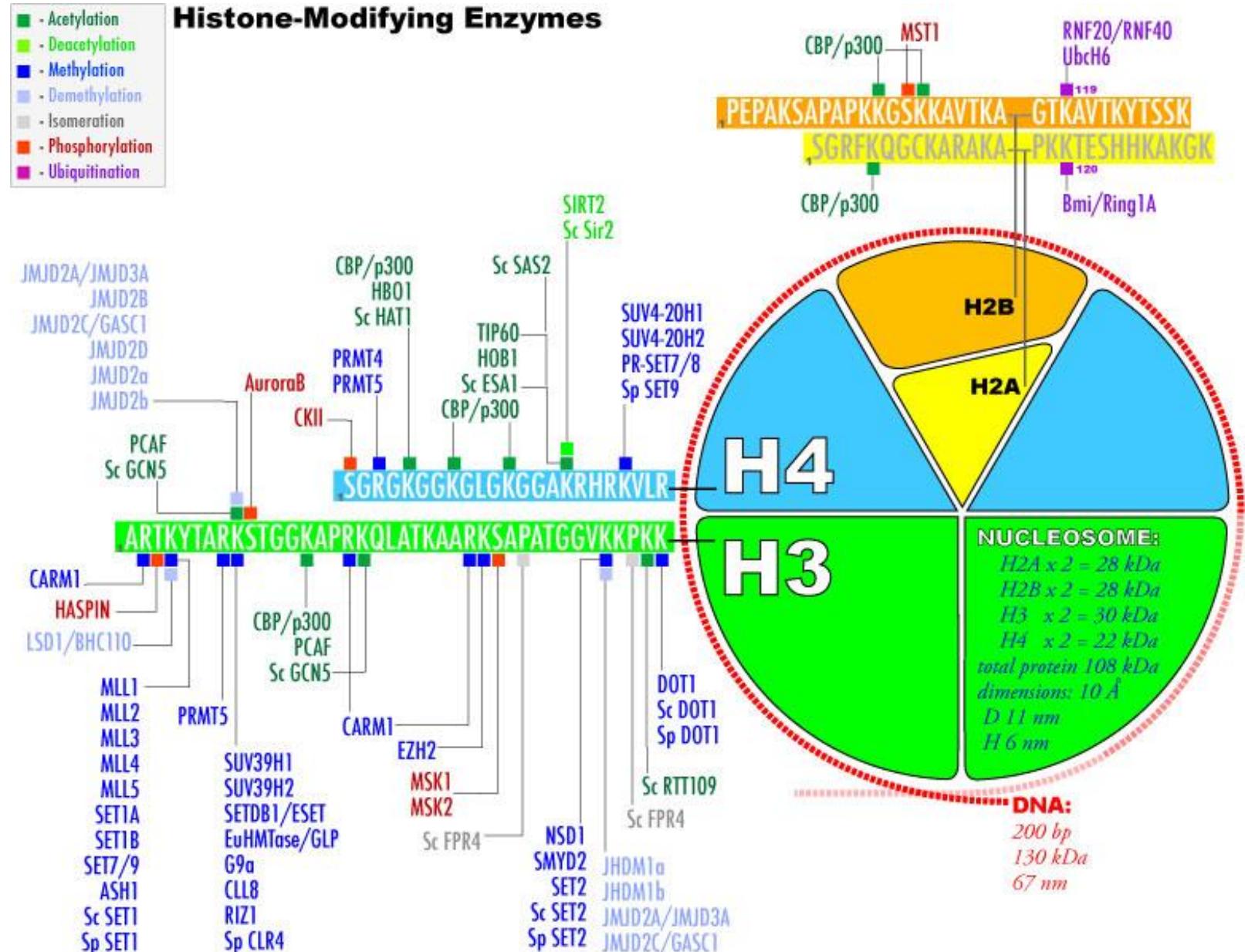


# Easy, right?

- Histone modifications can be our **guide** in the genome!

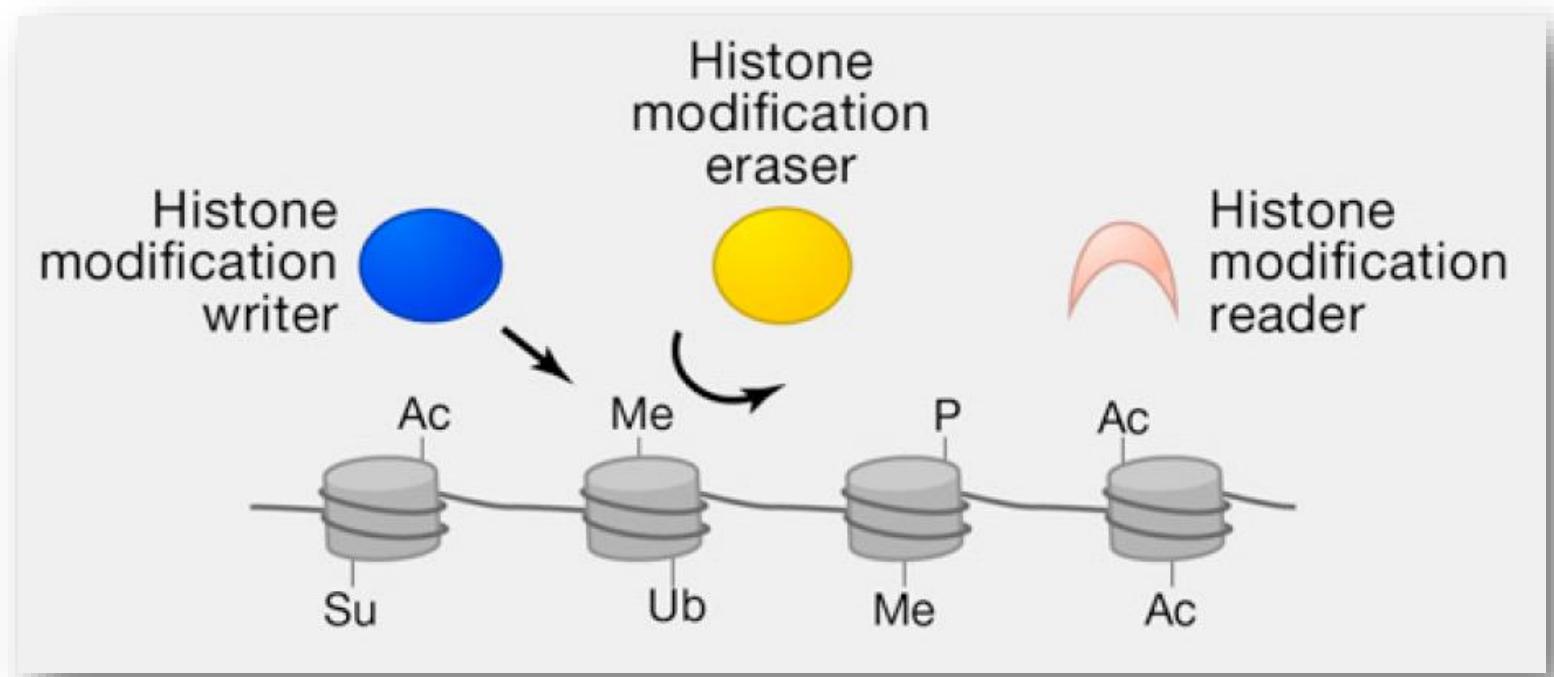


# Wrong...



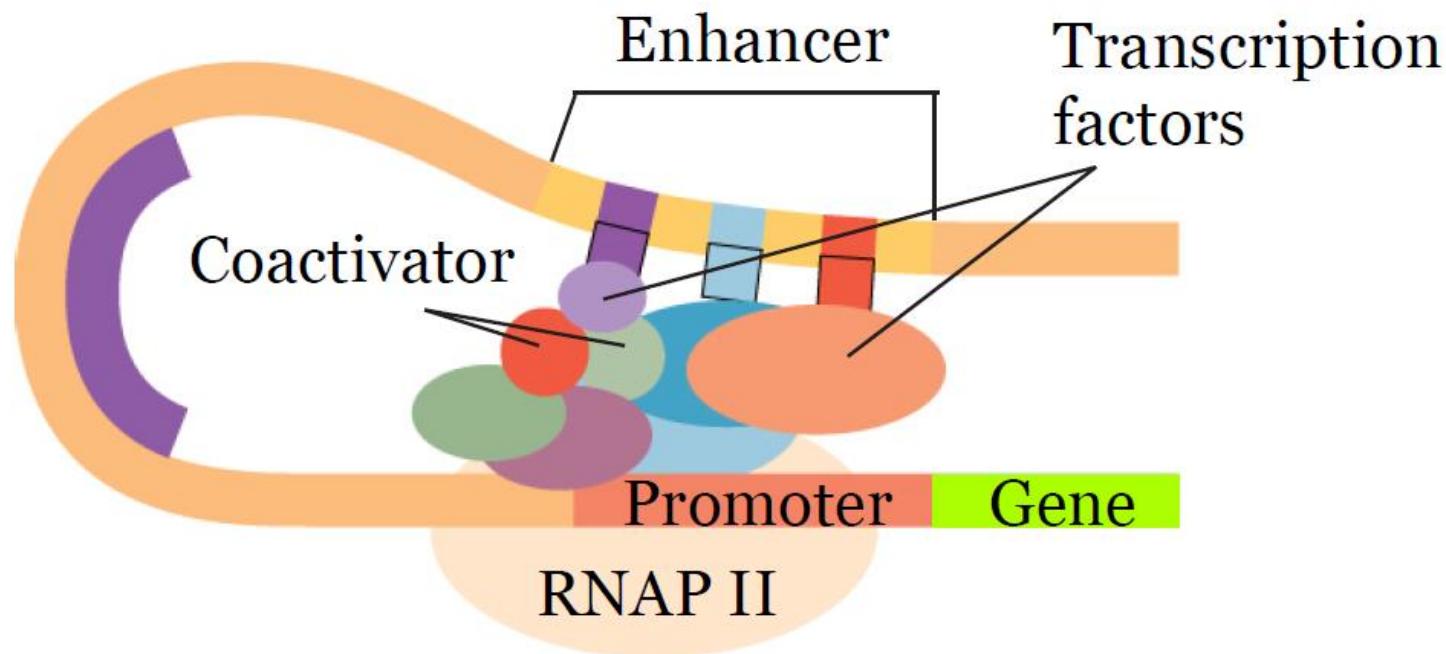
# Histone modification summary

- there are protein domains that add, remove, and recognize various histone modifications
- histone additions are basically post-translational modifications
- main modifications are methylation and acetylation, main histone is H3



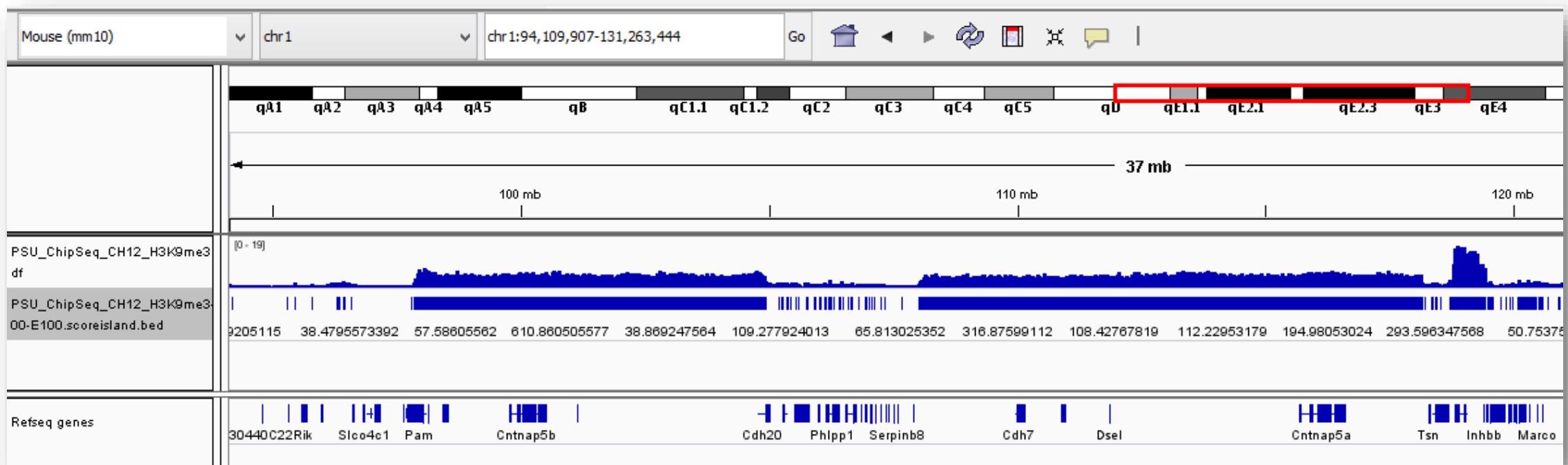
# Extended model of transcription

- transcriptional regulation is **very complex** and is influenced by (at least):
  - DNA methylation
  - histone modifications
  - transcription factors (10-20 per gene!)
- transcription factors recruit proteins that often also modify histones
- regulation can happen at transcriptional pre-initiation, initiation, or elongation



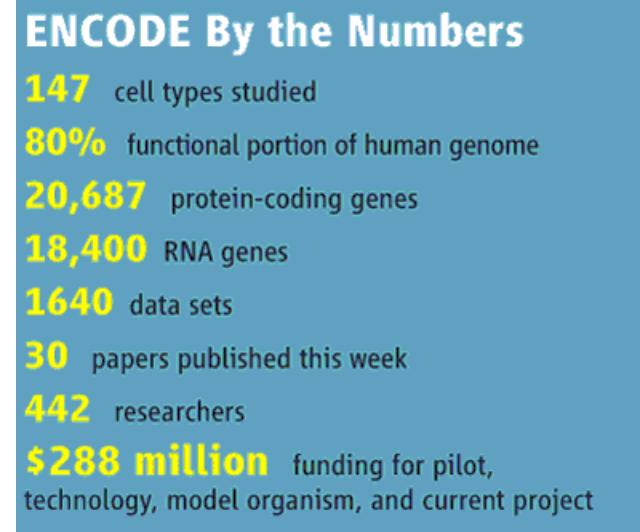
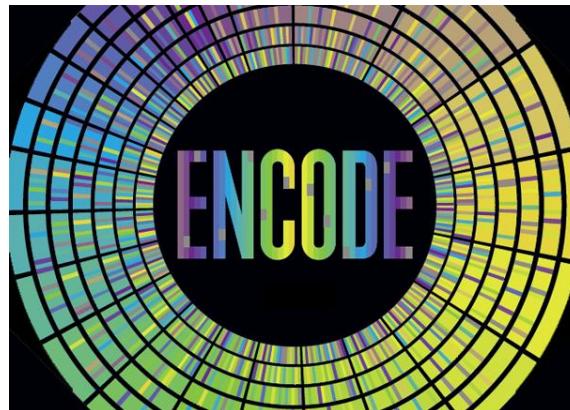
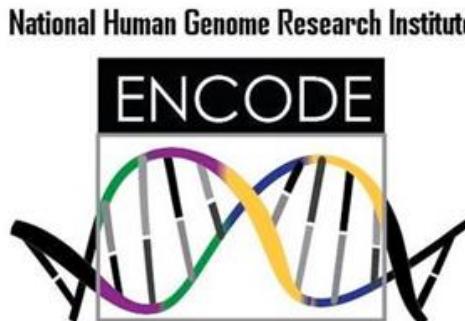
# TFs vs. Histone Modifications in ChIP-Seq

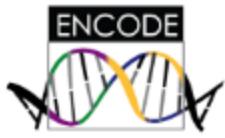
- TFs and histone modifications are quite different:
  - TFs have narrow binding motifs (5-20 bps) and are very localized
  - histone marks are often limited only by insulators (e.g. CTCF) and show broad signal patterns
- broad signal marks require higher library complexity (20M reads vs 3-10M for TFs)



# ENCODE project

- ENCODE = ENCyclopedia Of DNA Elements
- pilot cost (2007): \$55M, up to date: ~\$300M
- do RNA-Seq, ChIP-seq of major TFs and histone modifications, DNA methylation
- series of publications in the Fall of 2012 (6 Nature papers, 30 papers overall)





ENCODE ChIP-seq Experiment Matrix hg19

# Unit Matrix

## Antibody Targets

*search for:*  tracks  files

## Cell Types

Tier 1

GM12878

H1-hESC

K562

10

CD20+

CD20+\_RO01778

CD20+ RO01794

## H1-neurons

Help-S3

HepG2

HUVEG

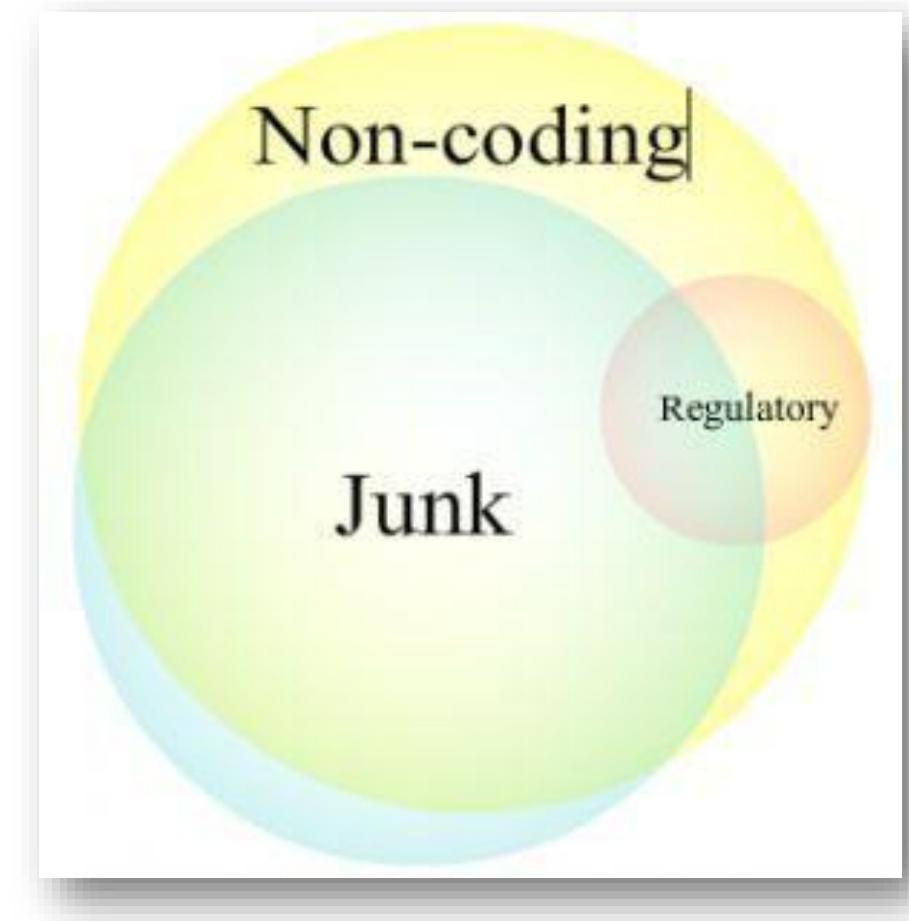
# ENCODE project discoveries

- 400,000 enhancers and 70,000 promoters
- more than 90% of genomic variation are in noncoding areas
- DNase I footprint is not that big
- mRNAs are more abundant in cytosol, other RNAs – in the nucleus
- “*More than 80% of human genome is functionally active*”



# ENCODE analysis criticism

- 80% of DNA cannot be truly functional, since only about 10% (5-15%) is conserved
- This means ~70% of genome is either
  - impervious to deleterious mutations, or
  - does not mutate, or
  - does not have deleterious mutations



# ENCODE criticism



**Mike White** @genologos

@homolog\_us @dangraur One of our Nature reviewers said this paper should not be published in any journal, ever

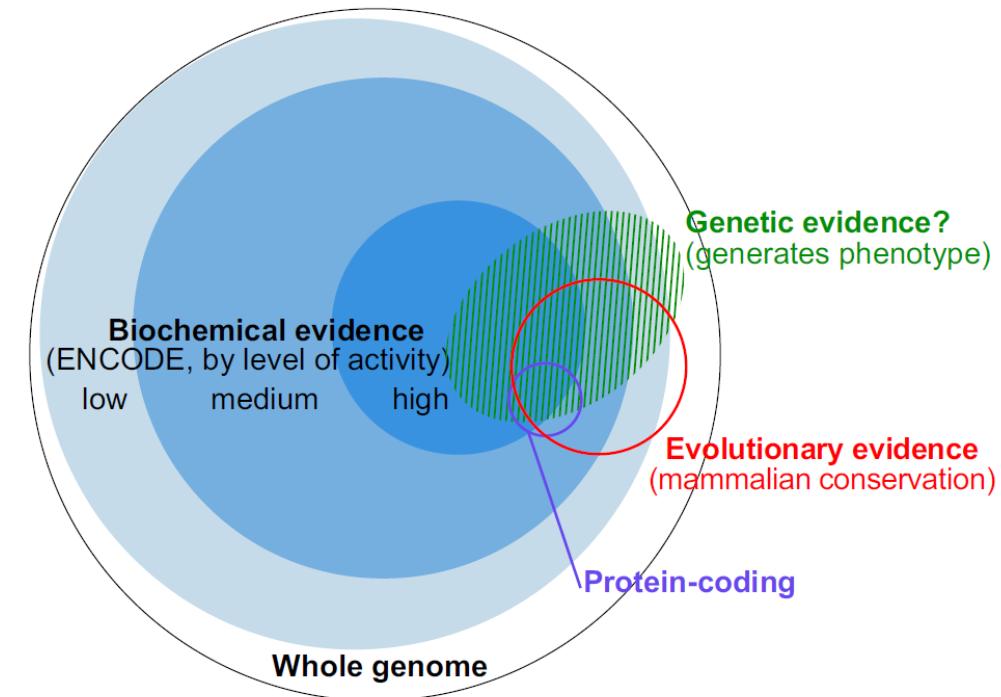
[View conversation](#)

10m

PERSPECTIVE

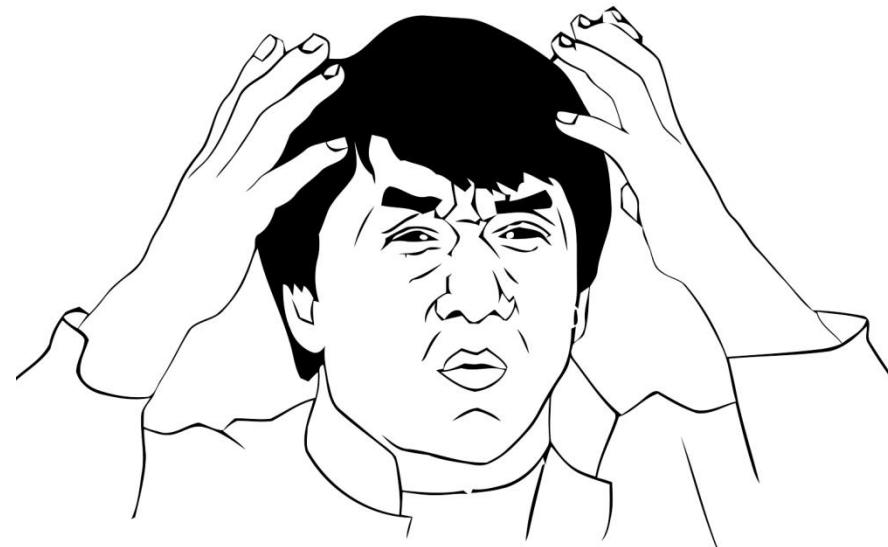
## Defining functional DNA elements in human genome

Manolis Kellis<sup>a,b,1,2</sup>, Barbara Wold<sup>c,2</sup>, Michael P. Snyder<sup>d,2</sup>, Bradley E. Bernstein<sup>b,e,f,2</sup>, Anshul Kundaje<sup>a,b,3</sup>, Georgi K. Marinov<sup>c,3</sup>, Lucas D. Ward<sup>a,b,3</sup>, Ewan Birney<sup>g</sup>, Gregory E. Crawford<sup>h</sup>, Job Dekker<sup>i</sup>, Ian Dunham<sup>g</sup>, Laura L. Elnitski<sup>j</sup>, Peggy J. Farnham<sup>k</sup>, Elise A. Feingold<sup>j</sup>, Mark Gerstein<sup>l</sup>, Morgan C. Giddings<sup>m</sup>, David M. Gilbert<sup>n</sup>, Thomas R. Gingeras<sup>o</sup>, Eric D. Green<sup>j</sup>, Roderic Guigo<sup>p</sup>, Tim Hubbard<sup>q</sup>, Jim Kent<sup>r</sup>, Jason D. Lieb<sup>s</sup>, Richard M. Myers<sup>t</sup>, Michael J. Pazin<sup>j</sup>, Bing Ren<sup>u</sup>, John A. Stamatoyannopoulos<sup>v</sup>, Zhiping Weng<sup>j</sup>, Kevin P. White<sup>w</sup>, and Ross C. Hardison<sup>x,1,2</sup>



# Flavors of DNA sequencing

- Sono-Seq, identical to ChIP-Seq but skipping the IP step to see open regions
- HITS-CLIP, PAR-CLIP, etc – for RNA-interacting proteins
- ChiRP-Seq, to measure RNA-bound DNA and proteins
- ChIP-exo uses exonuclease treatment to achieve up to single base-pair resolution
- **Hundreds** more!



# Example: ChIP-Exo

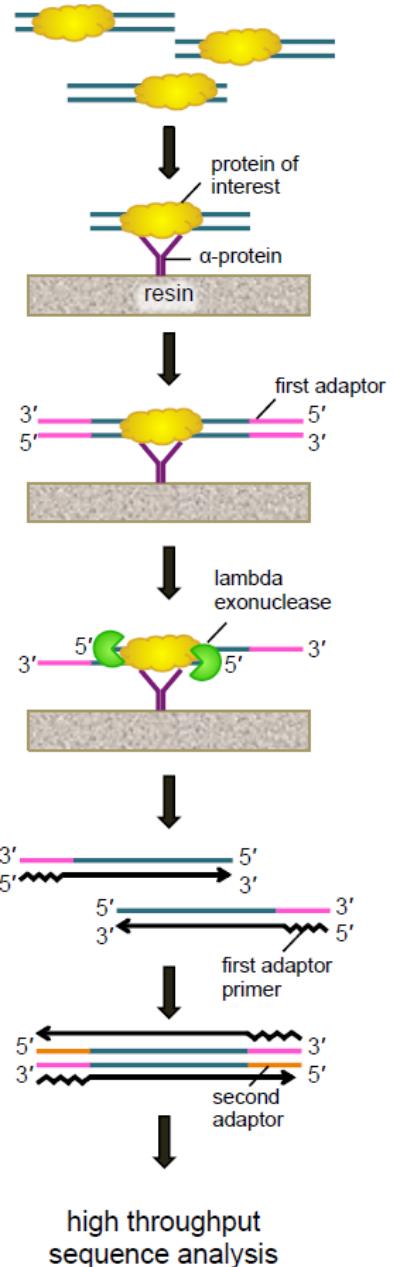
- ChIP-Exo allows to greatly increase **precision** of TF positioning on the DNA
- It uses endonuclease and ligation-mediated PCR
- Drawbacks – requires a lot of cells, generates libraries of poor complexity

## Frank Pugh

Evan Pugh University Professor  
Willaman Chair in Molecular Biology and  
Professor of Biochemistry and Molecular Biology

456A North Frear Laboratory  
University Park, PA 16802  
Email: [bfp2@psu.edu](mailto:bfp2@psu.edu)  
Phone: (814) 863-8252  
[Download as vCard](#)

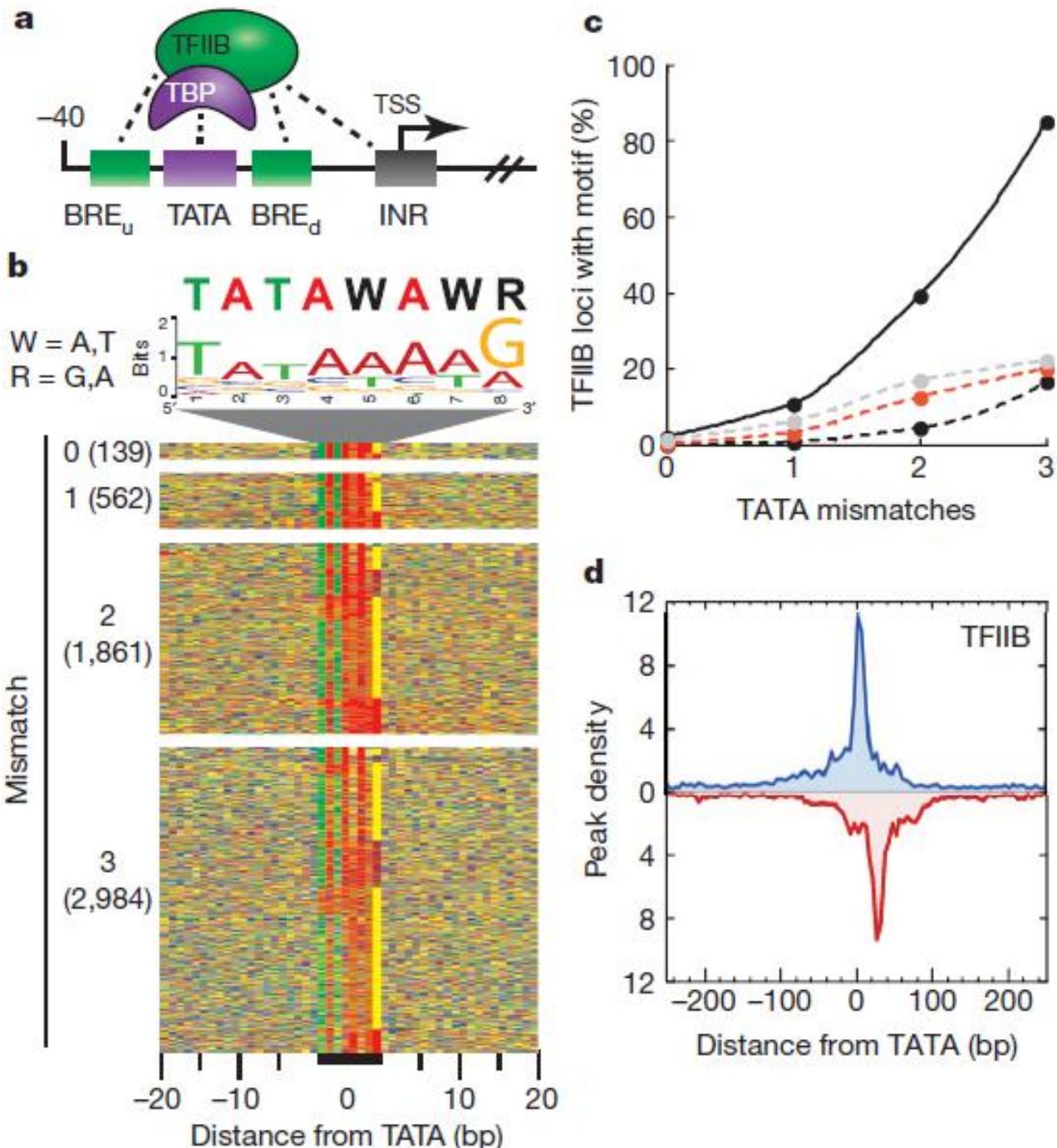
Websites



Venters, B.J.; Pugh, B.F., Nature., 502, pp 53-8 (2013)

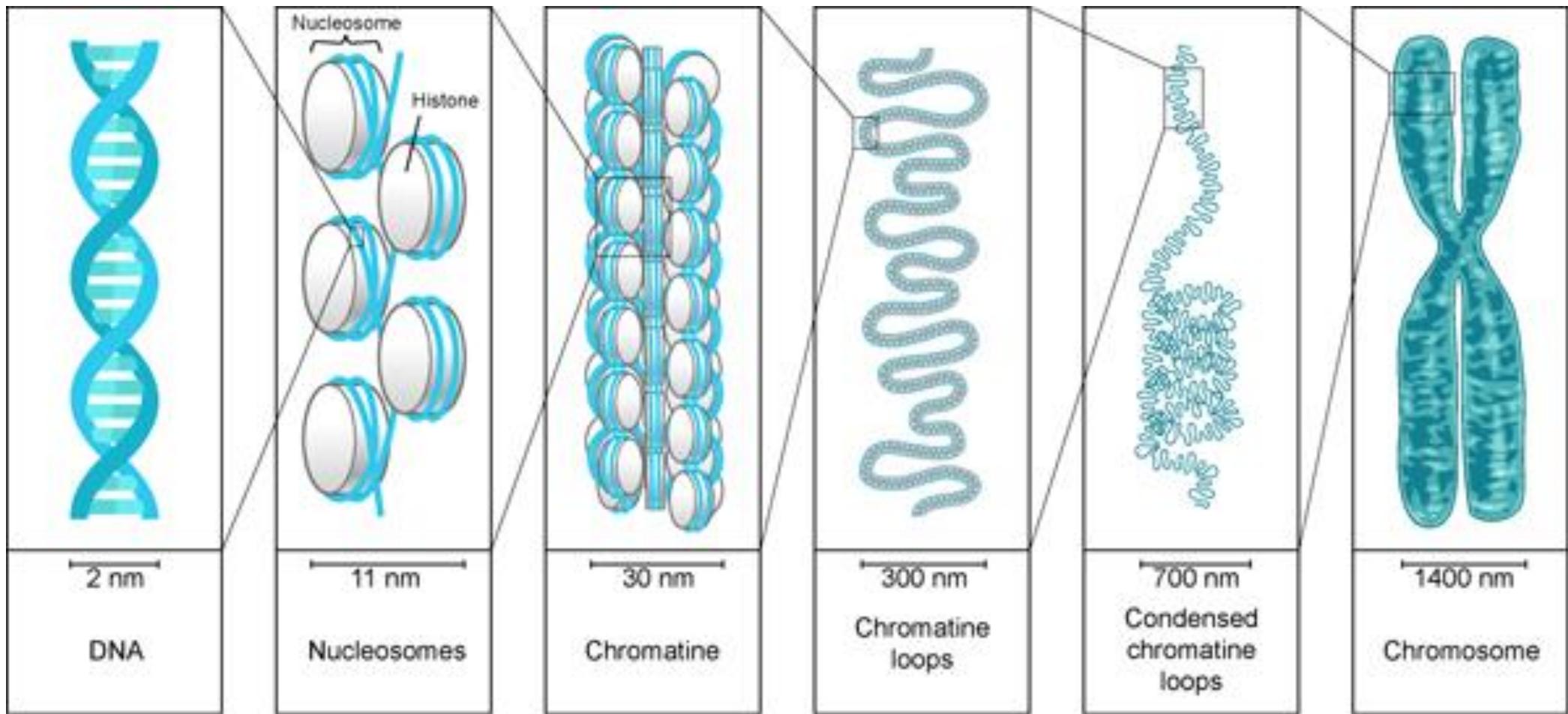
# ChIP-Exo

- virtually all gene promoters have the **same general structure**, including (often degenerate) TATA box!

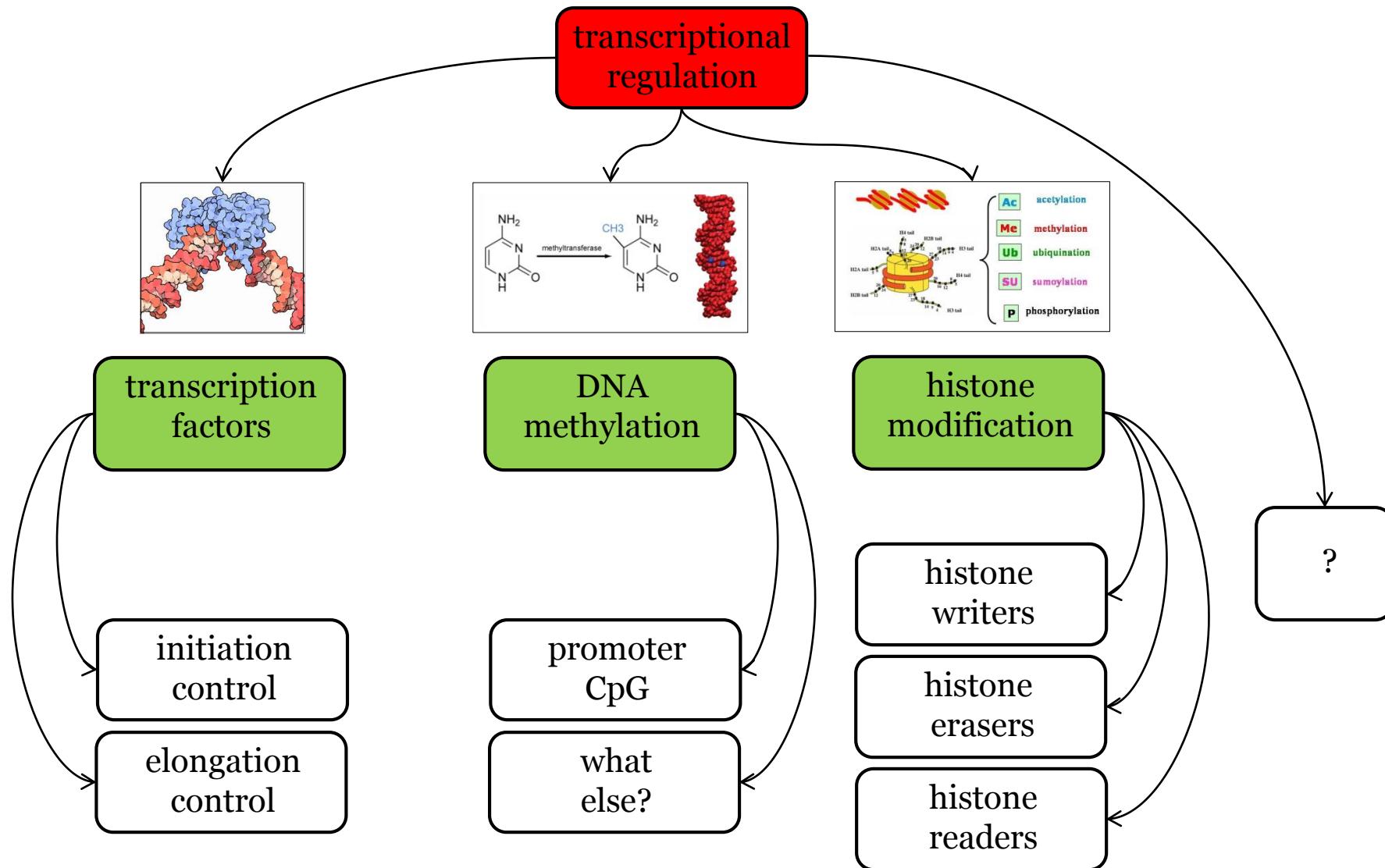


# Summary

- DNA packaging in eukaryotes is complicated and participates in regulation



# Summary



Tired yet?

