# Epigenomics data resources and the genome browser

Oana Ursu

# Motivation

- 2001: First draft of the human genome

- ~20 000 genes covering 1-3% of the genome



First printout of the human genome

Image from
http://en.wikipedia.org/wiki/Human_Genome_Project#mediaviewer/File:Wellcome_
genome_bookcase.png
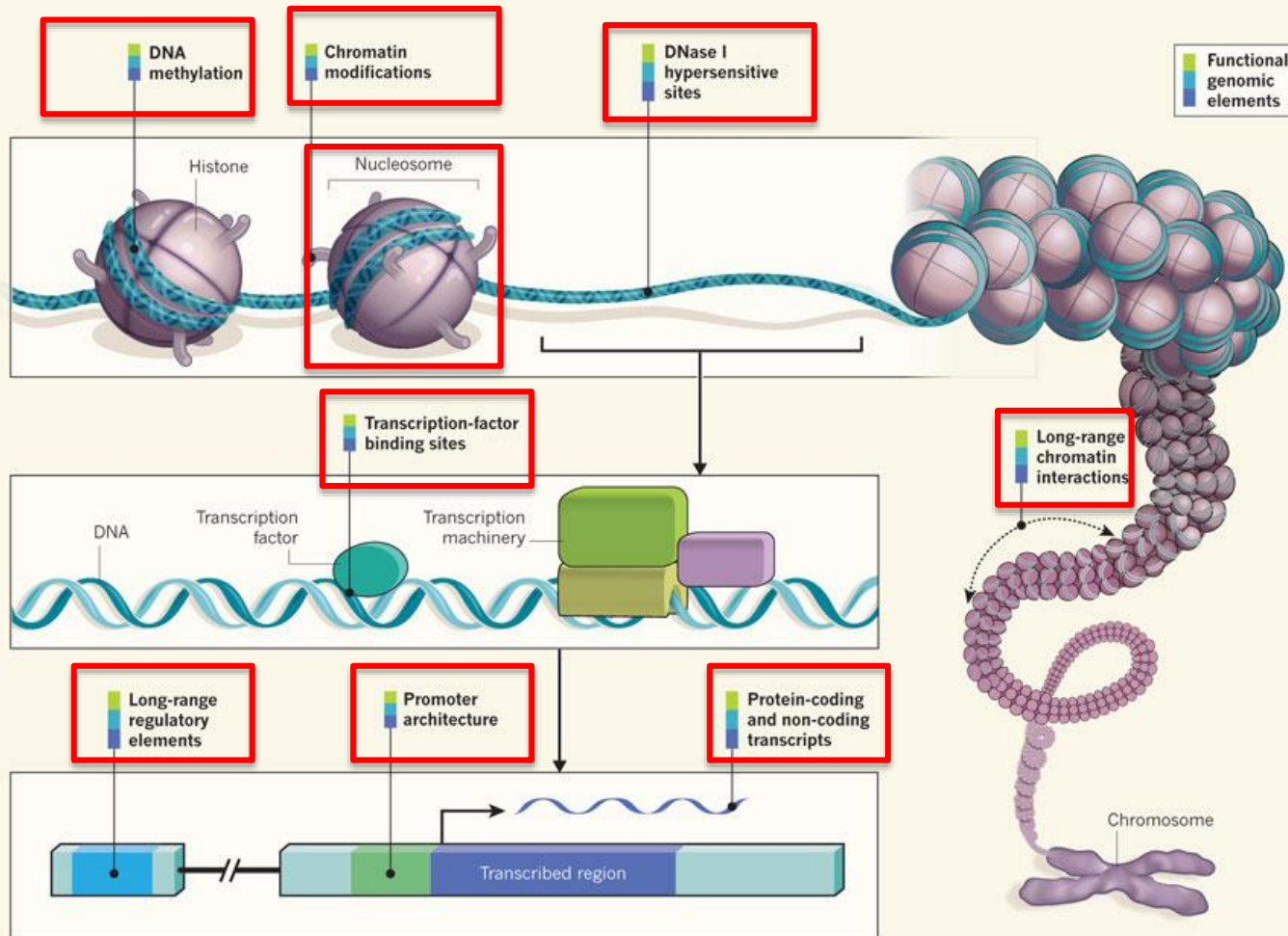
**What about the remaining 97-99% of the genome?**

**ENCODE**
(Encyclopedia Of DNA Elements)
set out to answer this question

# ENCODE: An <u>En</u>cyclopedia <u>Of</u> <u>D</u>NA <u>E</u>lements

- Goal: Complete catalog of all functional elements
  - Protein-coding genes: heavily studied, sequence-based, transcription
  - Non-coding DNA: much less studied, some motifs, diverse assays
- Dimensions for catalog completion
  - Genome-wide: systems-level view
  - Cell types: hundreds of human tissues and cell types
  - Dynamics: time, conditions, stimulation, environment, response


- Pilot phase 2003-: Small-scale targeted experiments in 1% of genome (30Mb)
  - Single gene, single pathway, few TFs, few cell-types, **tiling array-based**
- Scale-up 2007-: 100-fold increase in scale (3Gb), more assays, tech dev
  - Big change: RNA-seq, ChIP-seq, DNase-seq, **next-gen seq technologies**
  - Game changer: complete view, integration possible, networks and circuits
- Build-up 2012-: Further increases in all dimensions
  - Deeper sequencing, more assays, more conditions, more TFs.
  - More validation

# Diversity of assayed biochemical events



RNA-seq
CAGE-seq
Exon Arrays

TF ChIP-seq

Chromatin ChIP-seq

DNase-seq
FAIRE-seq

Methyl RRBS
Methyl Arrays

3C, 4C, 5C
ChIA-PET
HiC

# ENCODE data at a glance



http://genome.ucsc.edu/ENCODE/cellTypes.html

http://genome.ucsc.edu/ENCODE/dataMatrix/encodeDataSummaryHuman.html

# Outline

- **The ENCODE project: experiments, data, findings**
  - Genes and transcripts: RNAseq
  - Open chromatin:          DNaseI-seq
  - DNA-binding proteins: ChIP-seq
  - Chromatin state:         Histone ChIP-seq
  - Genome 3D:               3C
- **The genome browser**

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - Open chromatin:         DNaseI-seq
  - DNA-binding proteins: ChIP-seq
  - Chromatin state:         Histone ChIP-seq
  - Genome 3D:               3C
- **The genome browser**

# Studying genes and transcripts: RNAseq

## RNA-seq workflow

Remove rRNA
(>97% of your sample is rRNA!)



AAAAAAAAA    1) PolyA+ RNA captured
TTTTTTTTTTT B

2) RNA fragmented and primed

3) First strand cDNA synthesized

4) Second strand cDNA synthesized

5) 3' ends adenylated and 5' ends repaired

6) DNA sequencing adapters ligated

Barcode

Rd1

7) Ligated fragments PCR amplified

Rd2    Index

8) Sequencing
9) Map reads to transcriptome

Figure adapted from Corney, 2014

# Studying genes and transcripts: RNAseq

## Example



## Goals

- Transcriptome assembly
- Gene expression quantification
- Splicing

# Studying genes and transcripts: RNAseq

**Main findings**

- Pervasive transcription
  - "62% of genomic bases are reproducibly represented in sequenced long (>200 nucleotides) RNA molecules or GENCODE exons" (ENCODE, 2012)

Proportion of genomic bases included in a primary transcript, by number of technologies supporting the transcribed base



Figure from ENCODE 1 paper (1% of the genome)

ENCODE, 2007, Nature
ENCODE, 2012, Nature

# Studying genes and transcripts: RNAseq

**Main findings**

- Pervasive transcription
  - "62% of genomic bases are reproducibly represented in sequenced long (>200 nucleotides) RNA molecules or GENCODE exons" (ENCODE, 2012)

- Many flavors of RNAs
  - ~8000 small RNAs, ~9000 lncRNAs
  - lncRNAs more cell-type restricted, lower expression levels (compared to protein-coding genes)

**Outstanding questions**

- lncRNA functions
- Gene regulation through the act of transcription, not the transcript

ENCODE, 2007, Nature
ENCODE, 2012, Nature

# Studying genes and transcripts: RNAseq

**RNAseq experiment variants**

- RNA selection: polyA selection
- Location: cytoplasmic, nuclear
- Read length: short reads, long-read RNAseq (PacBio, Moleculo)
- More: CAGE-seq (TSS), Ribo-seq (translated transcripts)

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - Open chromatin:          DNaseI-seq
  - DNA-binding proteins: ChIP-seq
  - Chromatin state:          Histone ChIP-seq
  - Genome 3D:                3C
- **The genome browser**

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - **Open chromatin:        DNaseI-seq**
  - DNA-binding proteins: ChIP-seq
  - Chromatin state:         Histone ChIP-seq
  - Genome 3D:              3C
- **The genome browser**

# Studying open chromatin: DNaseI-seq

## Chromosomal Subunits in Active Genes Have an Altered Conformation

Globin genes are digested by deoxyribonuclease I in red blood cell nuclei but not in fibroblast nuclei.

Harold Weintraub and Mark Groudine

Knowledge of the structure of DNA has provided many insights into its biological function (1). In higher cells, a detailed understanding of the structure of chromatin will probably provide analogous insights into how genes are regulated. Already, there are a number of important observations demonstrating a relation between the structure of chromatin and its biological activity (2, 3).

The packaging of most of the nuclear DNA is now thought to be based on repeating units of about 180 to 200 base pairs of DNA associated with specific complexes of histones (4, 5), possibly tw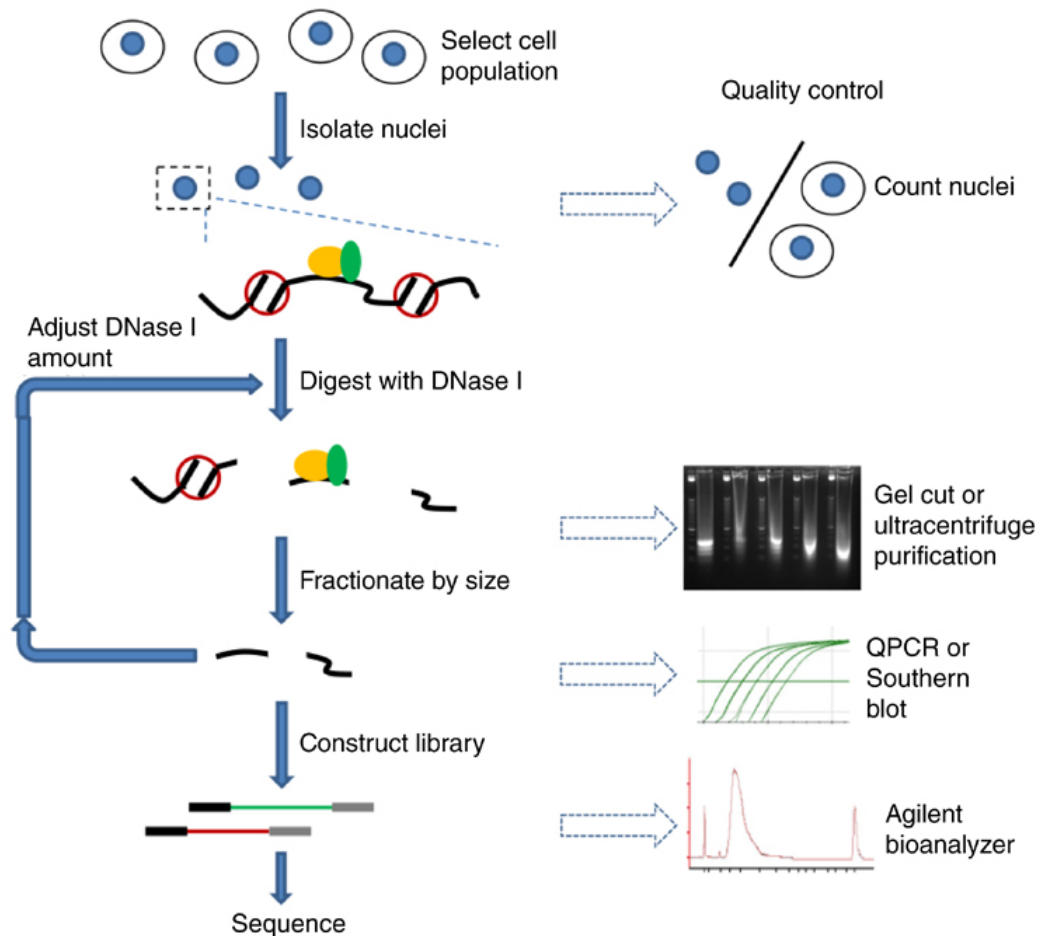o self-complementary tetramers each containing one of the four major histones (6). These two tetramers could define the twofold axis of symmetry within the nucleosome. These complexes interact through 70 to 90 amino acid residues at their carboxyl terminal end to produce a tight, trypsin-resistant core (7). The positively charged histone amino terminal residues extend outward from this core and define what may prove to be a "kinked" or "coiled" pathway for the DNA (5, 8) about the histone complexes. These so-called "particles-on-a-string" or "nu" bodies constitute the primary level of folding for the bulk of the chromosome. Through their mutual interactions higher levels of DNA packaging can be achieved, although details of this organization are not known. At present there is no proof that nu bodies are homo-

Dr. Weintraub is an assistant professor in the Department of Biochemical Sciences, Frick Laboratories, Princeton University, Princeton, New Jersey 08540. Dr. Groudine was a visiting fellow in the same department and is now at the Department of Radiation Oncology, University of Washington Hospital, Seattle 98105.

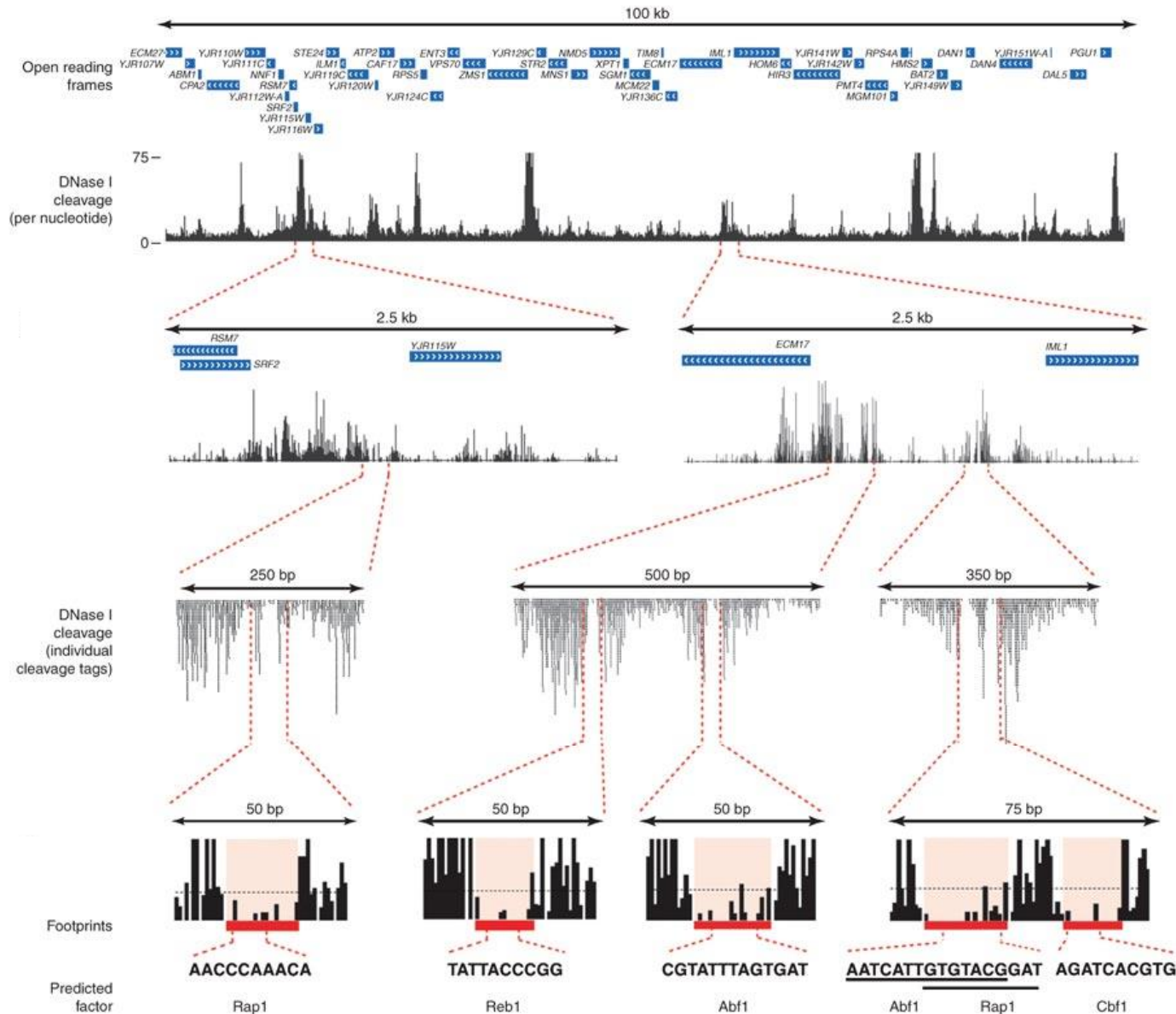# Studying open chromatin: DNaseI-seq

## DNaseI-seq workflow



This is hard!

Figure from Zeng et al., Nature, 2012

# Studying open chromatin: DNaseI-seq
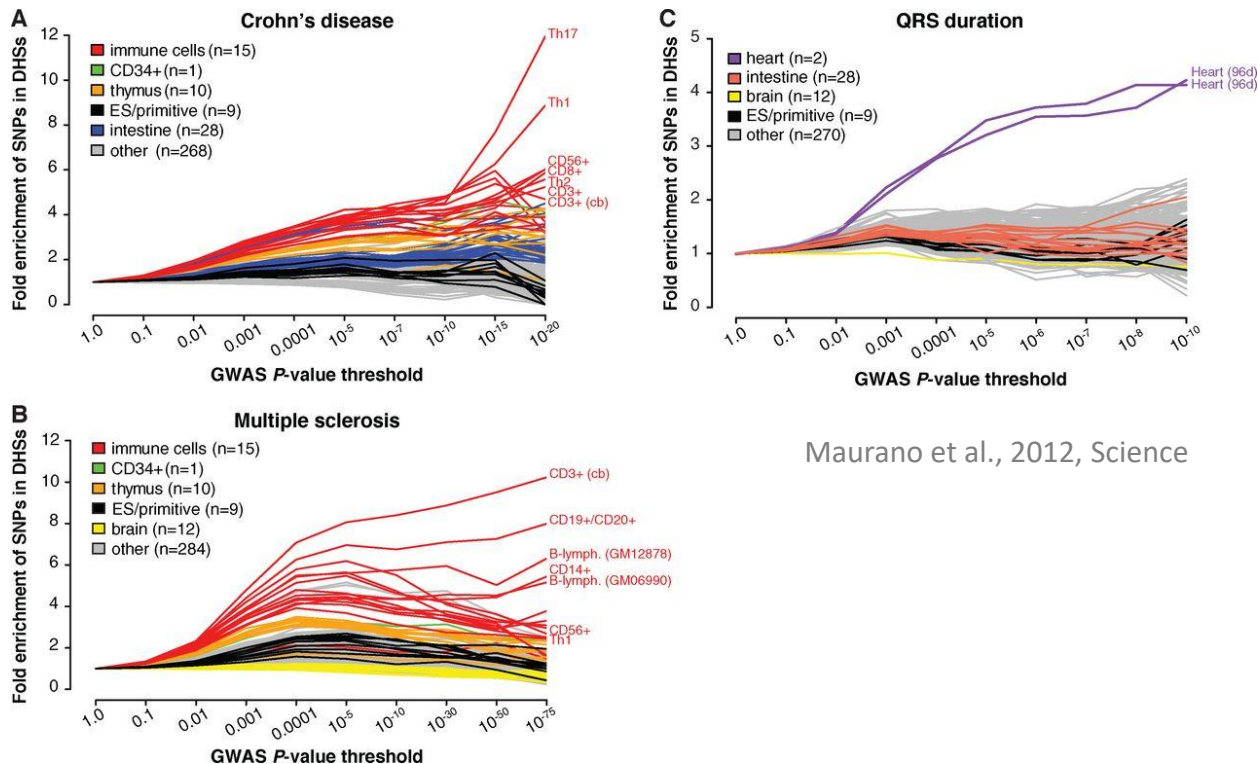


**Goals**
- Open chromatin

- TF footprinting

Figure from Hesselberth et al., Nature, 2009

# Studying open chromatin: DNaseI-seq

## Main findings

- DNaseI-seq sites at TSS, at enhancers, at protein-bound regions
- DNaseI-seq very cell-type specific (modules)
- Using DNaseI sites, can match disease with most likely affected cell type (because disease mutations fall in cell-type specific DNaseI sites)



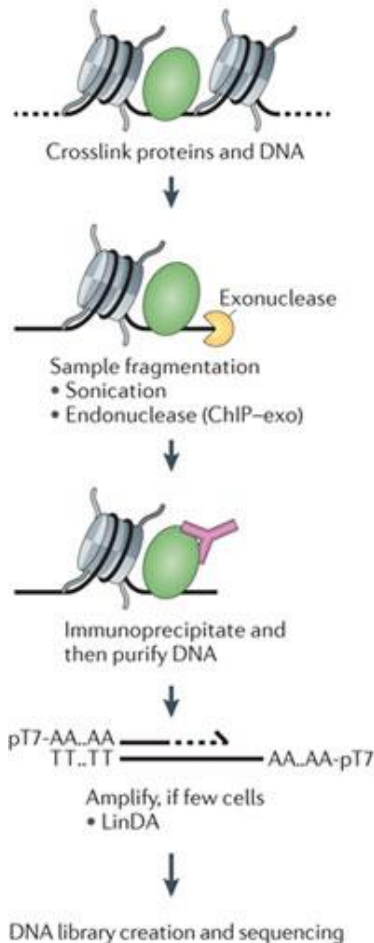Maurano et al., 2012, Science

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - **Open chromatin:        DNaseI-seq**
  - DNA-binding proteins: ChIP-seq
  - Chromatin state:        Histone ChIP-seq
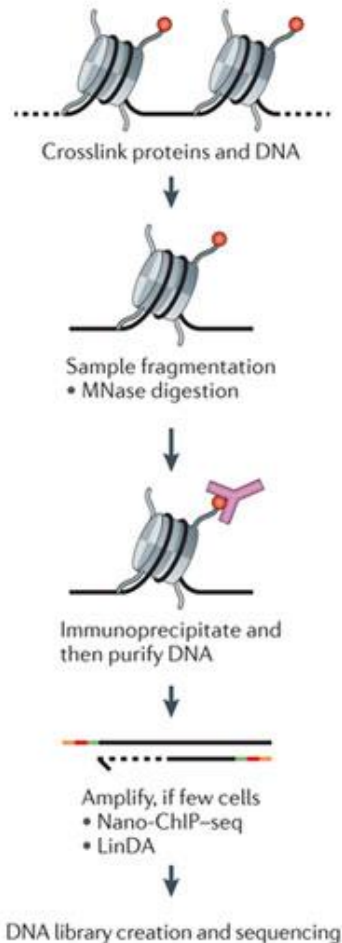  - Genome 3D:            3C
- **The genome browser**

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - **Open chromatin:          DNaseI-seq**
  - **DNA-binding proteins: ChIP-seq**
  - Chromatin state:          Histone ChIP-seq
  - Genome 3D:                3C
- **The genome browser**

# Studying DNA-binding proteins: ChIP-seq
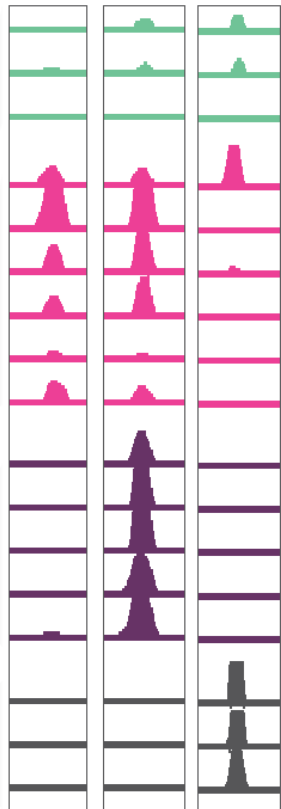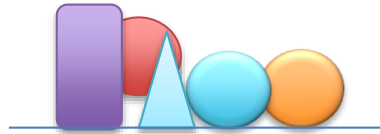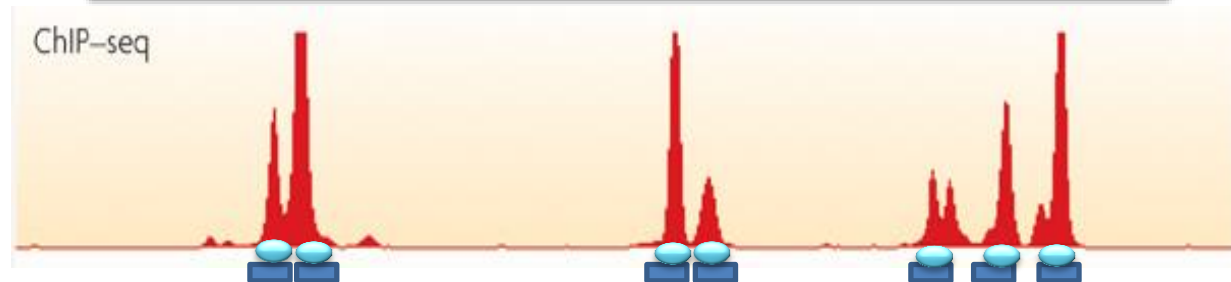
**ChIP-seq workflow**

# Studying DNA-binding proteins: ChIP-seq



**100s of binding maps** of different regulatory proteins!

Genome-wide protein-DNA binding maps
**CONTINUOUS**

ChIP-seq

Signal Peaks (potential binding sites)
**DISCRETE**

# Studying DNA-binding proteins: ChIP-seq

## Main findings

- High quality TF binding motifs



Kheradpour et al., 2012, Genome Research

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - **Open chromatin:          DNaseI-seq**
  - **DNA-binding proteins: ChIP-seq**
  - Chromatin state:          Histone ChIP-seq
  - Genome 3D:                3C
- **The genome browser**

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - **Open chromatin:           DNaseI-seq**
  - **DNA-binding proteins: ChIP-seq**
  - **Chromatin state:         Histone ChIP-seq**
  - Genome 3D:                3C
- **The genome browser**

# Studying chromatin state: ChromHMM

**ChromHMM workflow**

- Find functional elements from histone marks using Hidden Markov Models

# Studying chromatin state: ChromHMM

**ChromHMM workflow**
- Find functional elements from histone marks using Hidden Markov Models
- Chromatin states annotated using known genomic features

# The dynamic chromatin state

# Outline

- **The ENCODE project: experiments, data, findings**
  - **Genes and transcripts: RNAseq**
  - **Open chromatin:        DNaseI-seq**
  - **DNA-binding proteins: ChIP-seq**
  - **Chromatin state:        Histone ChIP-seq**
  - Genome 3D:                3C
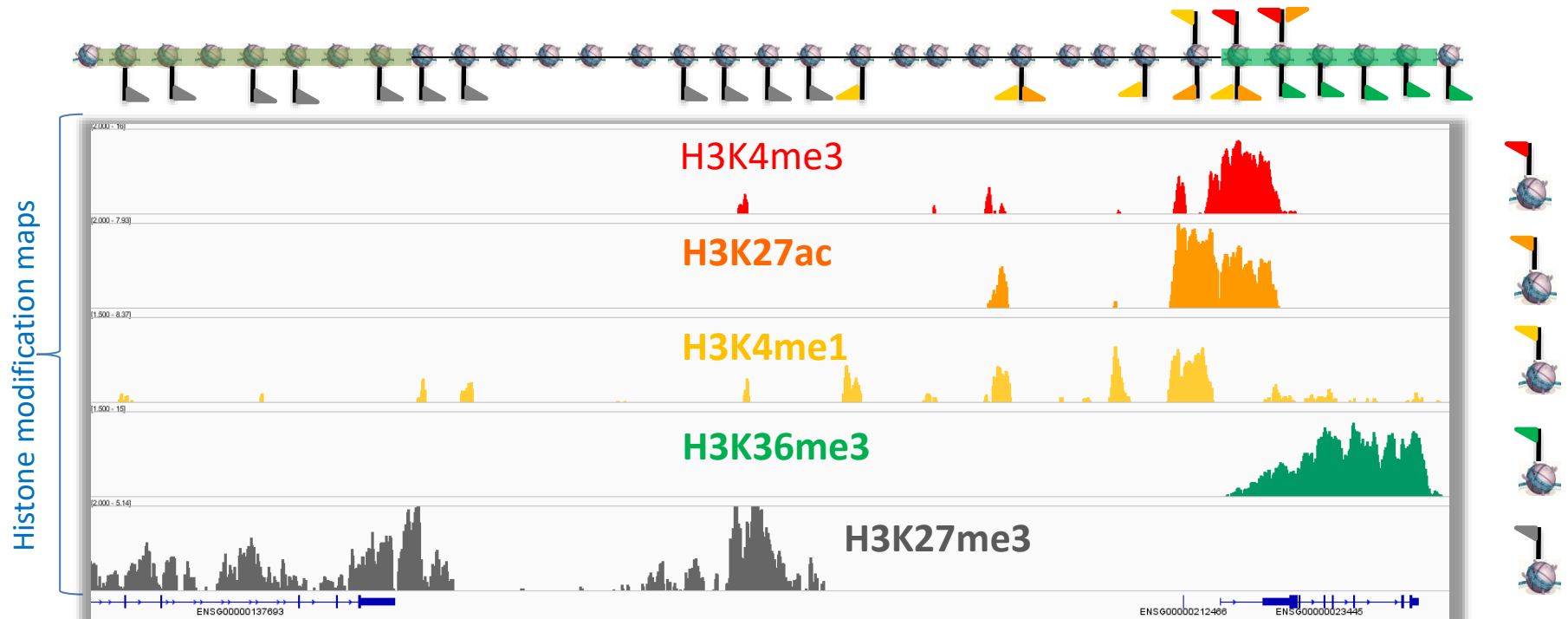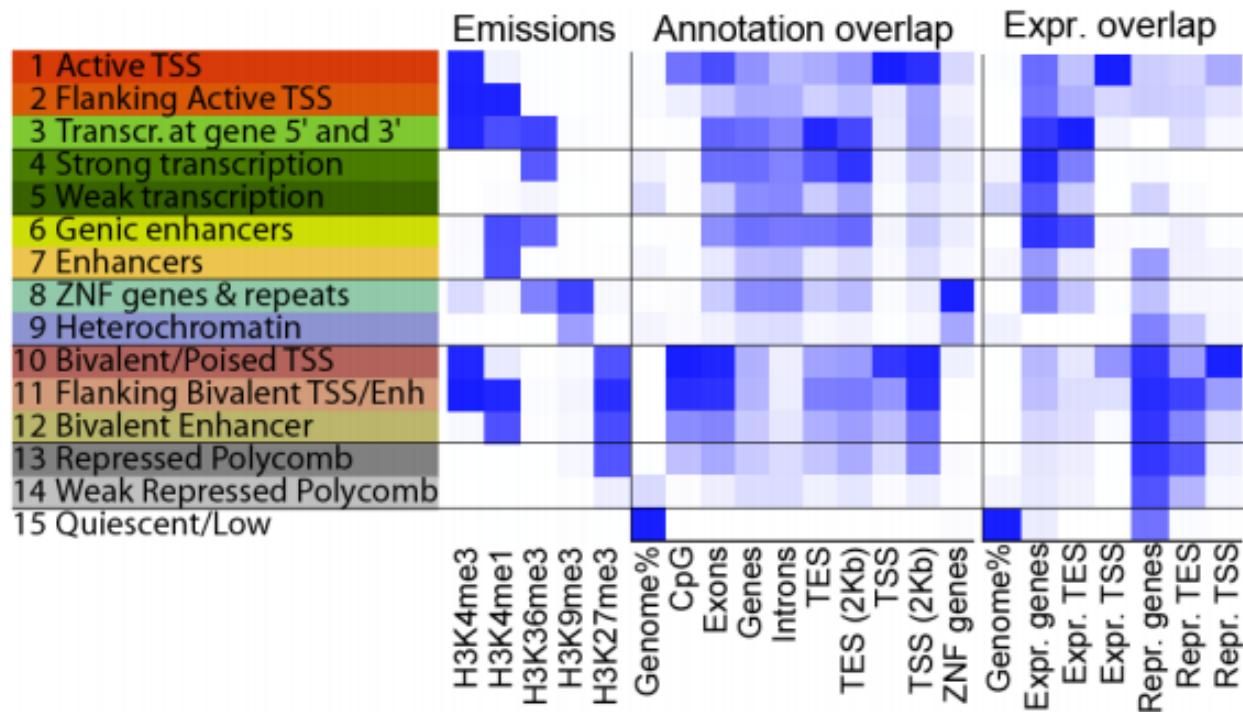- **The genome browser**

# Outline

- **The ENCODE project: experiments, data, findings**
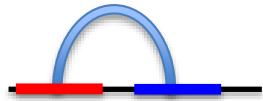  - **Genes and transcripts: RNAseq**
  - **Open chromatin:           DNaseI-seq**
  - **DNA-binding proteins: ChIP-seq**
  - **Chromatin state:          Histone ChIP-seq**
  - **Genome 3D:                3C**
- **The genome browser**

# Studying the genome 3D: 3C



New!
DNaseHiC, CaptureC



Crosslinking of interacting loci → Fragmentation → Ligation → DNA purification

Dekker et al. 2013

# Studying the genome 3D: 3C

**Main findings**

# Beyond ENCODE

## > 150 Cell-Types/Tissues

- 6 histone marks (Histone ChIP-seq)
- Open chromatin (DNase-seq)
- DNA methylation (WGBS, RRBS)
- Gene expression (RNA-seq)

### Roadmap Epigenomics Project
- Primary tissues

# Beyond ENCODE

**Roadmap Epigenomics Project**  http://roadmapepigenomics.org
- Primary tissues: chromatin state, open chromatin, expression

**BLUEPRINT**

**GTEx (Genotype-Tissue Expression) https://www.gtexportal.org/home/**
- Genetic variation