

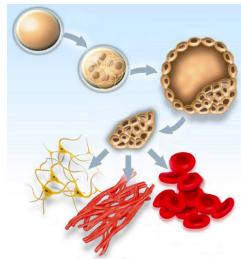
UNDERSTANDING GENE REGULATION – PART I

Myrto
Kostadima

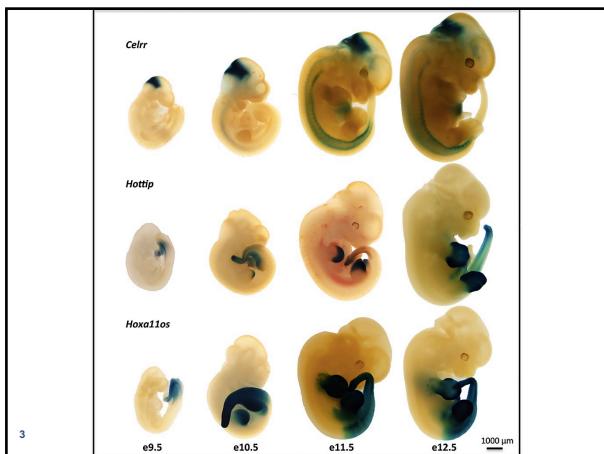
Ensembl Regulation Project Leader, EMBL-EBI

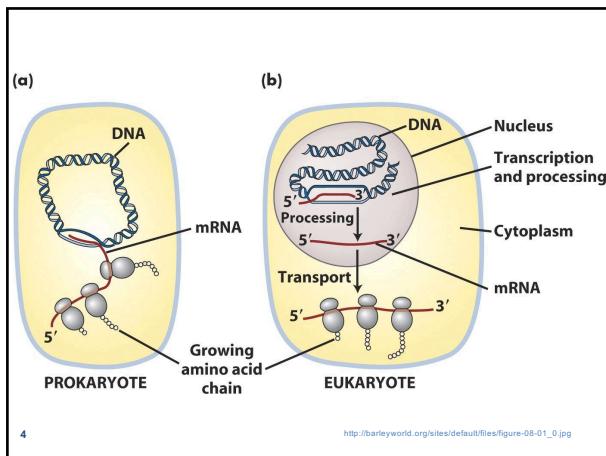
One Genome – Many cell types

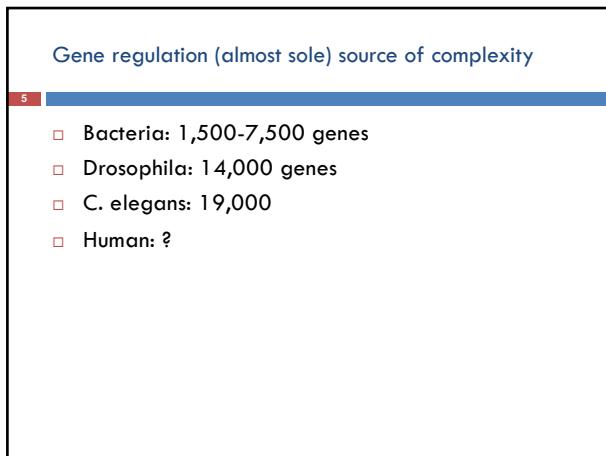
2

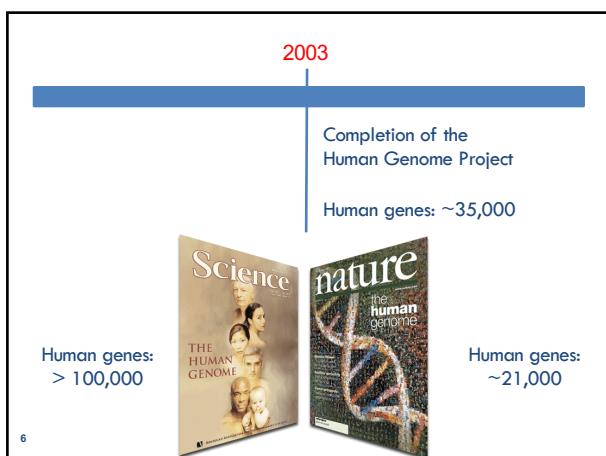


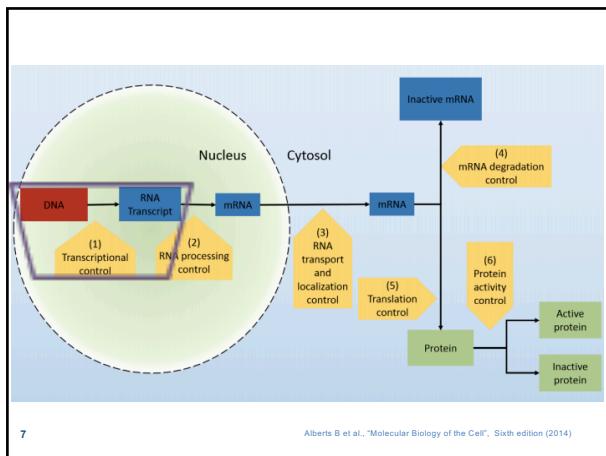
- Essentially all cells of an organism share
 - the same genome
 - the same genes
- Hundreds of different cell types with a clearly distinct phenotype
- Difference?
 - Gene expression profiles



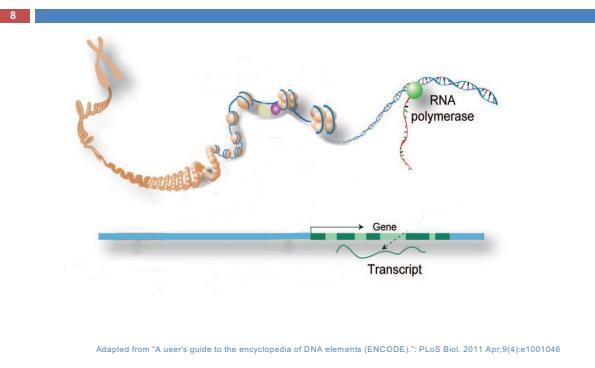




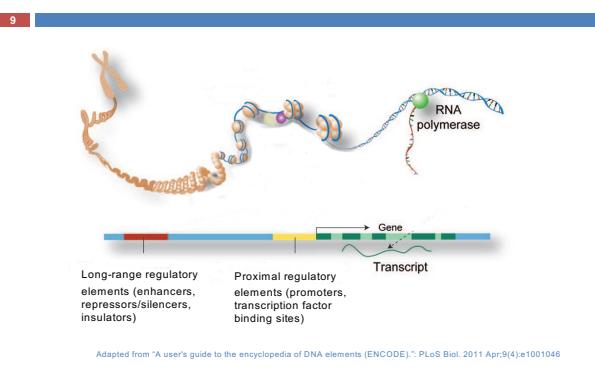




From protein-coding DNA...



... to regulatory DNA



Elements of gene regulation

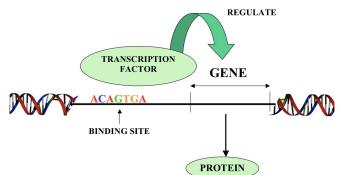
10

- Transcription factors
 - Promoters, enhancers, insulators..
- RNA Pol II machinery
- Epigenetics
 - DNA methylation
 - Nucleosome positioning/Open chromatin
 - Histone modification
 - Euchromatin/heterochromatin
- 3D Chromatin organisation

Transcription factors (TFs)

11

- TFs are a subset of genes that encode proteins that bind to sequences to regulate transcription



- ~1,000 human genes are TFs (depending on definition)

Vaquerizas JM et al., Nat Rev Genet. 2009 Apr;10(4):252-63.

Prokaryotes vs Eukaryotes

12

- | | |
|------------------------------|-------------------------------------|
| □ simple process | □ complex process |
| □ genes organised in operons | □ many TFs per gene |
| □ one TF per operon | □ many regulatory elements per gene |

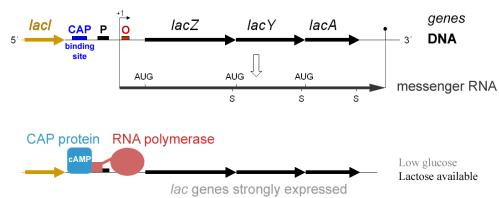
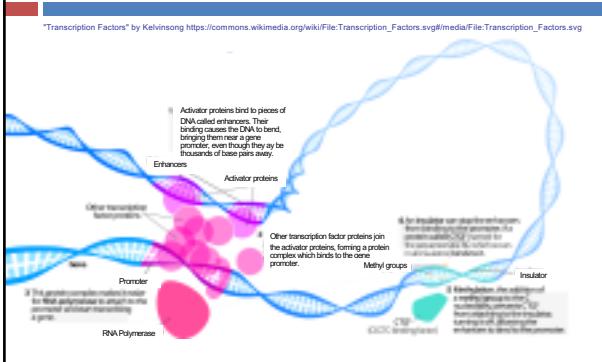
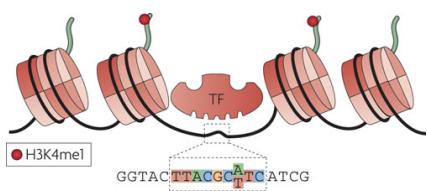


Image adapted from: https://commons.wikimedia.org/wiki/File:Lac_operon-2010-21-01.png

TFs and regulatory regions



TF binding site and motif



Nature Reviews | Genetics

Position Weight Matrices

Consensus sequence	
GAGGTAAC	□ A PWM is a commonly used representation of motifs in biological sequences.
TCCGTAAGT	□ PWMs are often derived from a set of aligned sequences that are thought to be functionally related.
CAGGTTGGA	□ They have become an important part of many software tools for computational motif discovery.
ACAGTCAGT	
TAGGTCATT	
TAGGTACTG	
ATGGTAACT	
CAGGTATAAC	
TGTGTGAGT	
AAGGTAAGT	
TAGGTAAGT	

Position Weight Matrices

16

POS 1|2|3|4|5|6|7|8|9
 1. G|A|G|G|T|A|A|A|C
 2. T|C|C|G|T|A|A|G|T
 3. C|A|G|G|T|T|G|G|A
 4. A|C|A|G|T|C|A|G|T
 5. T|A|G|G|T|C|A|T|T
 6. T|A|G|G|T|A|C|T|G
 7. A|T|G|G|T|A|A|C|T
 8. C|A|G|G|T|A|T|A|C
 9. T|G|T|G|T|G|A|G|T
 10. A|A|G|G|T|A|A|G|T

POS 1 2 3 4 5 6 7 8 9
 A [- - - - - - - - -]
 C [- - - - - - - - -]
 G [- - - - - - - - -]
 T [- - - - - - - - -]

Position Weight Matrices

17

POS 1|2|3|4|5|6|7|8|9
 1. G|A|G|G|T|A|A|A|C
 2. T|C|C|G|T|A|A|G|T
 3. C|A|G|G|T|T|G|G|A
 4. A|C|A|G|T|C|A|G|T
 5. T|A|G|G|T|C|A|T|T
 6. T|A|G|G|T|A|C|T|G
 7. A|T|G|G|T|A|A|C|T
 8. C|A|G|G|T|A|T|A|C
 9. T|G|T|G|T|G|A|G|T
 10. A|A|G|G|T|A|A|G|T

POS 1 2 3 4 5 6 7 8 9
 A [3 6 1 0 0 6 7 2 1]
 C [2 2 1 0 0 2 1 1 2]
 G [1 1 7 10 0 1 1 5 1]
 T [4 1 1 0 10 1 1 2 6]

Position Weight Matrices

18

POS 1|2|3|4|5|6|7|8|9
 1. G|A|G|G|T|A|A|A|C
 2. T|C|C|G|T|A|A|G|T
 3. C|A|G|G|T|T|G|G|A
 4. A|C|A|G|T|C|A|G|T
 5. T|A|G|G|T|C|A|T|T
 6. T|A|G|G|T|A|C|T|G
 7. A|T|G|G|T|A|A|C|T
 8. C|A|G|G|T|A|T|A|C
 9. T|G|T|G|T|G|A|G|T
 10. A|A|G|G|T|A|A|G|T

POS 1 2 3 4 5 6 7 8 9
 A [.3 .6 .1 0 0 .6 .7 .2 .1]
 C [.2 .2 .1 0 0 .2 .1 .1 .1]
 G [.1 .1 .7 1 0 .1 .1 .5 .1]
 T [.4 .1 .1 0 1 .1 .1 .2 .6]

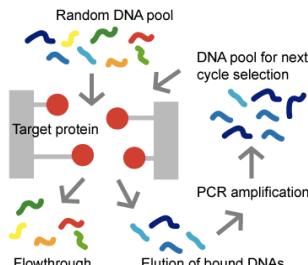
Phylogenetic Footprinting

19

YDR374C

Systematic Evolution of Ligands by Exponential Enrichment (SELEX)

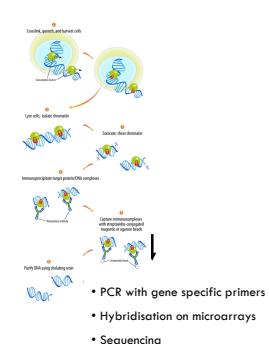
20



[Image from: <http://altair.sci.hokudai.ac.jp/g6/Projects/Sclex-e.html>]

Chromatin ImmunoPrecipitation (ChIP)

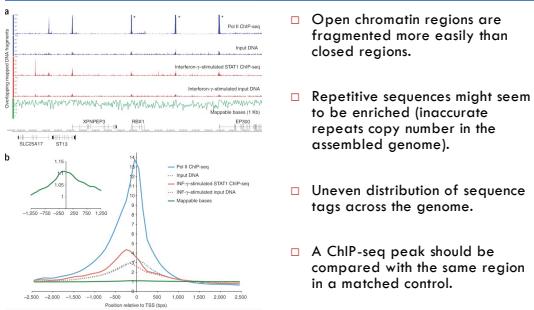
- Transcription factors
 - RNA Polymerase II



Chromatin immunoPrecipitation followed by sequencing (ChIP-seq)

- One of the early applications of high-throughput sequencing
- First studies published in 2007
 - Johnson et al (Science) - NRSF
 - Barski et al (Cell) - histone methylation
 - Robertson et al (Nature Methods) - STAT1
 - Mikkelsen et al (Nature) - histone modification
- Thousands of publications currently in PubMed

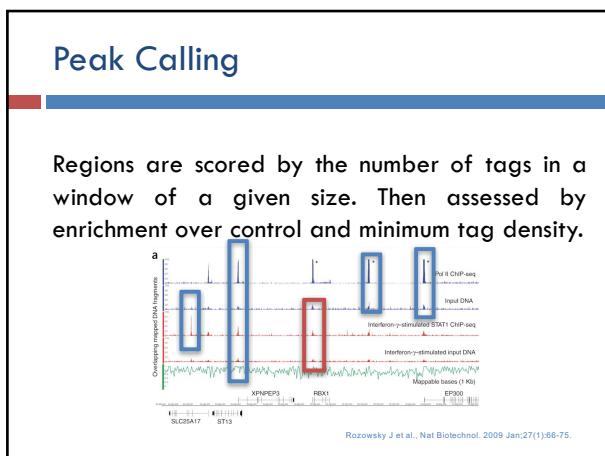
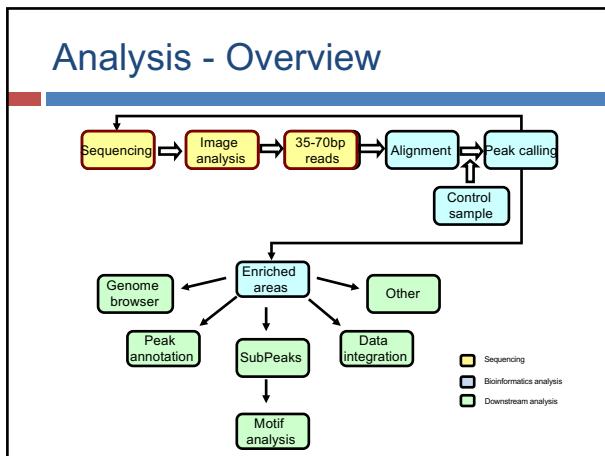
Why we need a control sample



- Open chromatin regions are fragmented more easily than closed regions.
- Repetitive sequences might seem to be enriched (inaccurate repeats copy number in the assembled genome).
- Uneven distribution of sequence tags across the genome.
- A ChIP-seq peak should be compared with the same region in a matched control.

Control types in ChIP-seq

- Input DNA
- Mock IP - DNA obtained from IP without antibody
 - Very little material can be pulled down leading to inconsistent results of multiple mock IPs.
- Nonspecific IP - using an antibody against a protein that is not known to be involved in DNA binding
- Sequencing a control can be avoided when looking at:
 - time points
 - differential binding pattern between conditions



- ## Peak Calling - Challenges
- Adjust for sequence mappability - regions that contain repetitive elements have different expected tag count
 - Different ChIP-seq applications produce different type of peaks. Most current tools have been designed to detect sharp peaks (TF binding)
 - Alternative tools exist for broader peaks (histone modifications that mark domains - transcribed or repressed), e.g. SICER

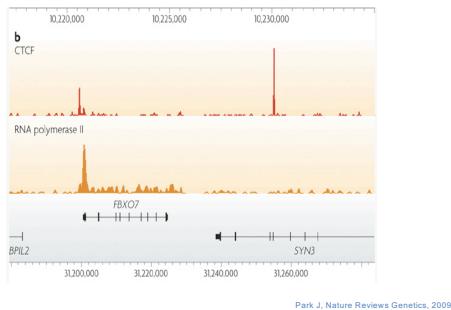
MACS - Peak detection

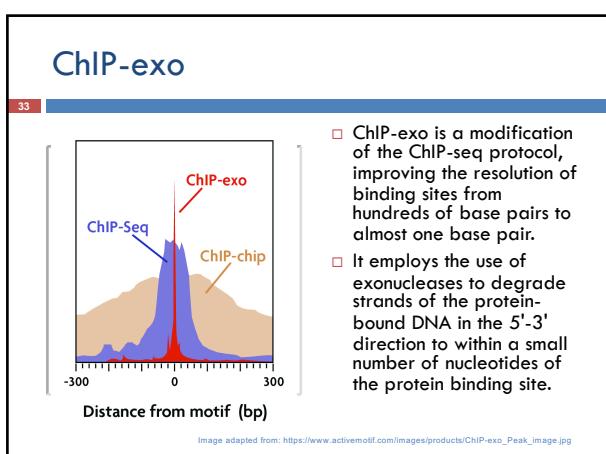
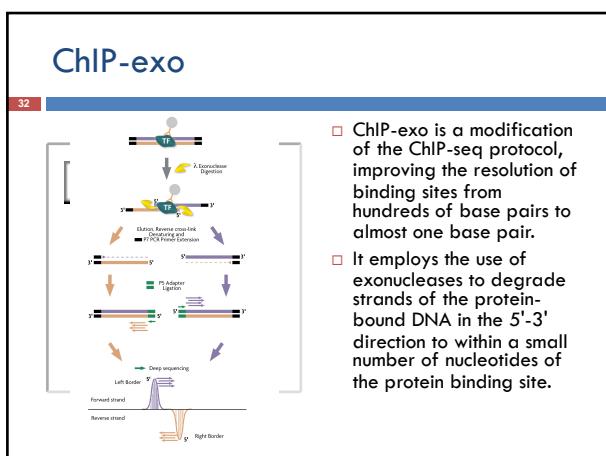
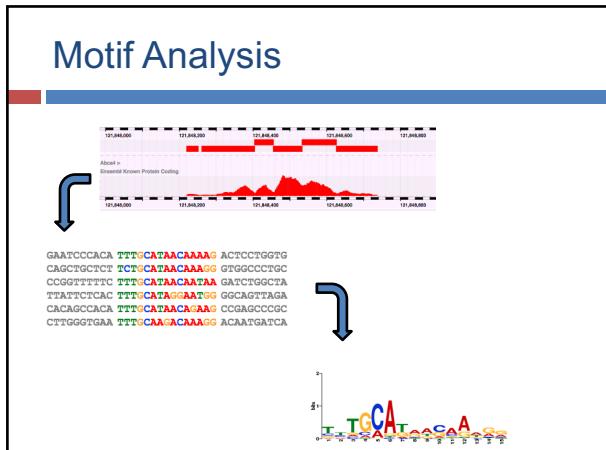
- Remove duplicate tags (in excess of what can be expected by chance)
- Slide window across the genome to find candidate peaks with a significant tag enrichment (Poisson distribution, global background, default p-value 10e-5)
- Merge overlapping peaks, and extend each tag d bases from its center
- Also looks at local background levels and eliminates peaks that are not significant with respect to local background
- Uses the control sample to eliminates peaks that are also present there

Analysis downstream to peak calling

- Visualisation - genome browser: Ensembl, UCSC, IGV
- Peak Annotation - finding interesting features surrounding peak regions
- Correlation with expression data
- Discovery of binding sequence motifs
- Gene Ontology analysis on genes bound by the same transcription factor
- Correlation with SNP data to find allele-specific binding

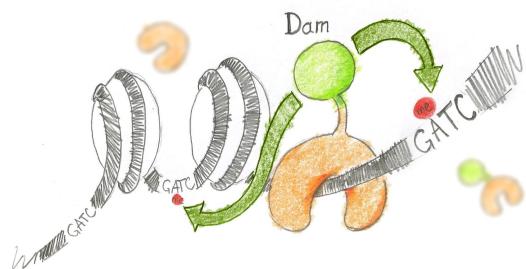
Visualisation in a genome browser





DNA adenine methyltransferase identification

34

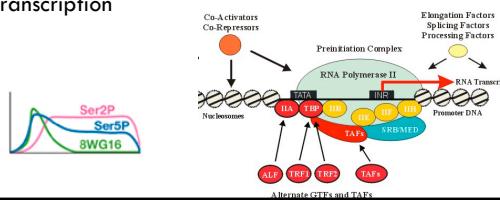


"DamID concept" by Guillaume Filion - Own work. Licensed under Public Domain via Commons
https://commons.wikimedia.org/wiki/File:DamID_concept.jpg#/media/File:DamID_concept.jpg

RNA Pol II and the transcriptional machinery

35

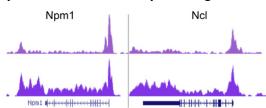
- RNA Pol II does not act alone
- Needs a Pre-initiation complex (TFs, mediator proteins)
- RNA Pol II is chemically modified in different phases of transcription



RNA Pol II pausing

36

- RNA Pol II modification & elongation factors needed for full transcripts
- RNA Pol II Pre-initiation complex can be present, but no (or low) transcription: RNA Pol II pausing



- RNA Pol II ChIP: Compute Pausing Index/Travelling ratio (signal at TSS vs gene body)

Further reading

37

REPORT

Suboptimization of developmental enhancers

Emma K. Farley^{1,2,*}, Katrina M. Olson^{1,2}, Wei Zhang³, Alexander J. Brandt⁴, Daniel S. Rokhsar¹, Michael S. Levitt^{1,2,*}

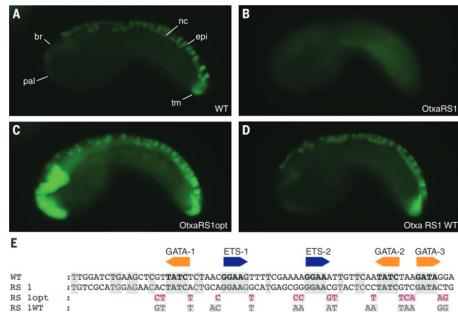
+ Author Affiliations

*Corresponding author. E-mail: msl2@princeton.edu (M.S.L.); ekfarley@princeton.edu (E.K.F.)

Science 16 Oct 2015;
Vol. 350, issue 6258, pp. 325-328
DOI: 10.1126/science.aac6948

Further reading

38



39

Gene Regulation Practical

Prerequisites for next week's practical:

1. Install R packages:

1. rtracklayer
2. ggplot2

2. Add the following to your UNIX \$PATH:

1. /local/data/genome_informatics/programs/bedtools2