

1529: Machine Learning in Bioinformatics (Spring 2017)

Applications of HMMs in Epigenomics

Yuzhen Ye
School of Informatics and Computing
Indiana University, Bloomington
Spring 2017

Contents

- Background: chromatin structure & DNA methylation
- Epigenomic projects
 - ENCODE
 - modENCODE
 - Human epigenome atlas
- Techniques
 - ChIP-Seq
 - MeDIP (MeDIP-chip & MeDIP-seq)
- Two applications
 - MeDIP-HMM (second-order HMM)
 - ChromHMM (multivariate HMM)

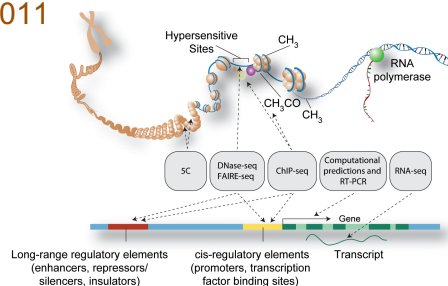
Background

- Genomic DNA is packaged into a complex molecular structure known as chromatin. This structure mediates the interaction between the genome and all types of regulatory and transcriptional molecules.
- In vertebrate genomes, methylation at position 5 of the cytosine in CpG dinucleotides is a heritable “epigenetic” mark that has been connected with both **transcriptional silencing and imprinting**
 - Ref: **DNA methylation patterns and epigenetic memory** (*Genes & Dev.* 2002. 16: 6-21)

ENCODE

- Encyclopedia of DNA Elements
 - “The ENCODE Consortium is integrating multiple technologies and approaches in a collective effort to discover and define the functional elements encoded in the human genome, including **genes, transcripts, and transcriptional regulatory regions, together with their attendant chromatin states and DNA methylation patterns.**”
 - Ref: **A User’s Guide to the Encyclopedia of DNA Elements (ENCODE)** (PLoS Biology, 2011)
- Initial phase launched in 2003—1% of the human genome
 - **Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project** (Nature, June 13, 2007)

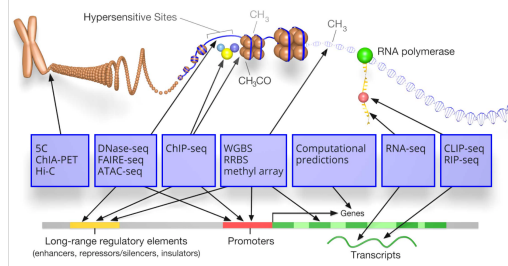
2011



The ENCODE Project Consortium (2011) A User's Guide to the Encyclopedia of DNA Elements (ENCODE). *PLoS Biol* 9(4): e1001046. doi:10.1371/journal.pbio.1001046
<http://www.plosbiology.org/article/info:doi/10.1371/journal.pbio.1001046>

PLOS BIOLOGY

2017



ENCODE

Based on an image by Darryl Lajo (NHGRI), Ian Dunham (EBI), Michael Piatn (NHGRI)

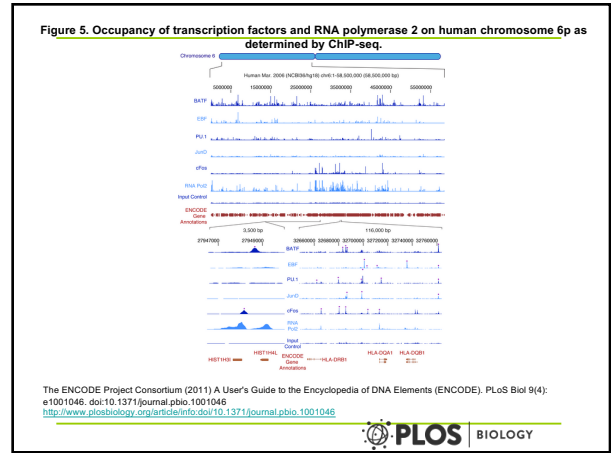
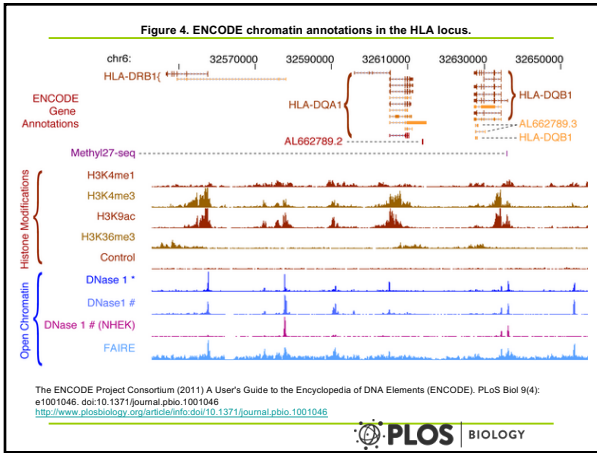
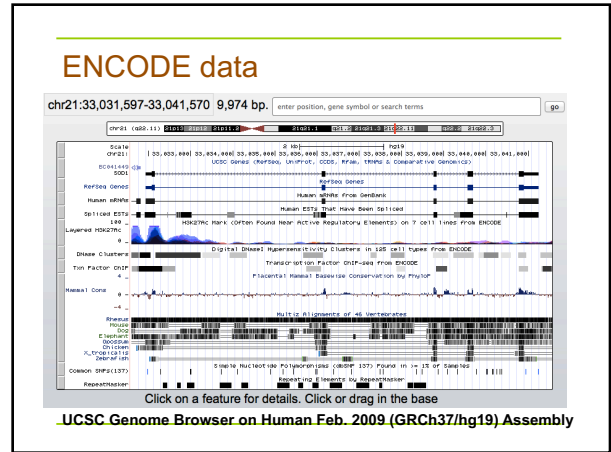
Ref: <https://www.encodeproject.org>

Table 1. Experimental assays used by the ENCODE Consortium.

Gene/Transcript Analysis		
Region/Feature	Method	Group
Gene annotation	GENCODE	Wellcome Trust
PolyA coding regions	RNA-seq 3' tag DNA microarray; PET	CSHL, Stanford/Harvard/Catoh
Total RNA coding regions	RNA-seq 3' tag DNA microarray; PET	CSHL
Coding regions in subcellular RNA fractions (e.g. nuclear, cytoplasmic)	PET	CSHL
Small RNA	short RNA-seq	CSHL
Transcription initiation (5'-end and termination (3-end) sites)	CAGE; 5TAGS	RIKEN, GS
Full-length RNA	RACE	University of Geneva; University of Louisiana
Protein-bound RNA coding regions	RP; CLIP	SUNY-Albany, CSHL
Transcription Factors/Chromatin		
Elements/Regions	Method(s)	Group(s)
Transcription Factor Binding Sites (TFBS)	ChIP-seq	Stanford/Harvard/UCSD/Novartis; HudsonAlpha/Catoh; Duke/Cornell; UW, U, Chicago/Stanford
Chromatin structure (accessibility, etc.)	DNAse1 hypersensitivity; FAIRE	UW, Duke, UNC
Chromatin modifications (H3K27ac, H3K9me3, H3K36me3, etc.)	ChIP-seq	Broad, UW
DNAse1 footprints	Digital genomic footprinting	UW
Other Elements/Features		
Feature	Method(s)	Group(s)
DNA methylation	HiSeq; Illumina Methy27; Methy-seq	HudsonAlpha
Chromatin interactions	3C; ChIA-PET	UMass; UW; GS
Genotyping	Illumina 1M Duo	HudsonAlpha

doi:10.1371/journal.pbio.1001046

The ENCODE Project Consortium (2011) A User's Guide to the Encyclopedia of DNA Elements (ENCODE). PLoS Biol 9(4): e1001046. doi:10.1371/journal.pbio.1001046
<http://www.plosbiology.org/article/info:doi/10.1371/journal.pbio.1001046>



modENCODE

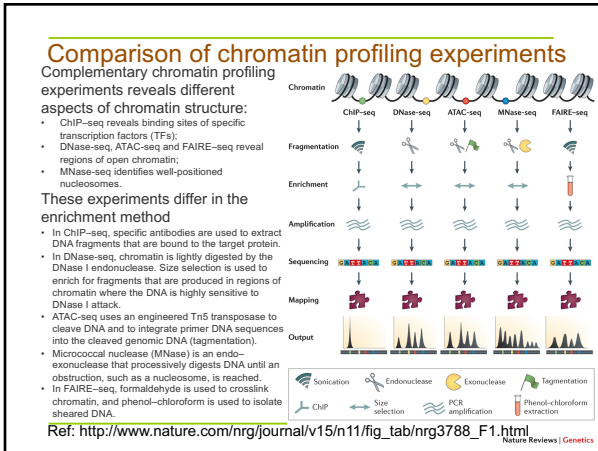
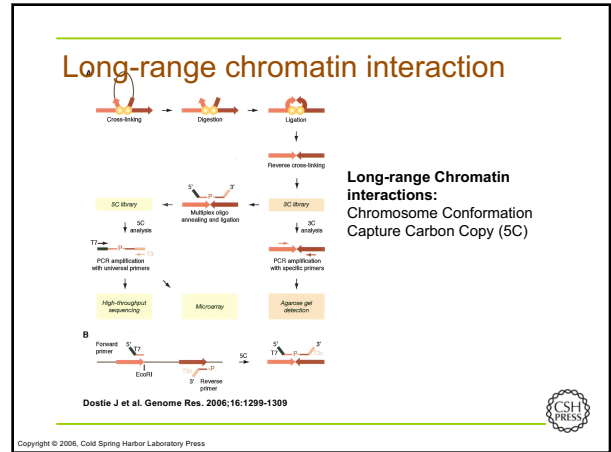
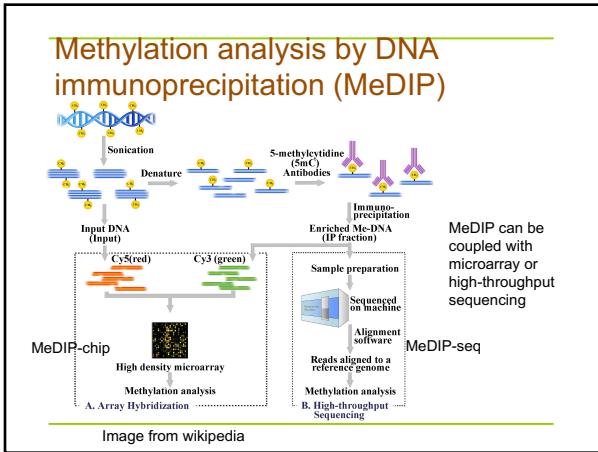
<http://www.modencode.org/>

"The modENCODE Project will try to identify all of the sequence-based functional elements in the *Caenorhabditis elegans* and *Drosophila melanogaster* genomes."

- Chromatin structure
- Copy Number Variation
- Gene Structure
- Genome Sequence
- Histone modification and replacement
- Metadata only
- Other chromatin binding sites
- RNA expression profiling
- Replication
- TF binding sites

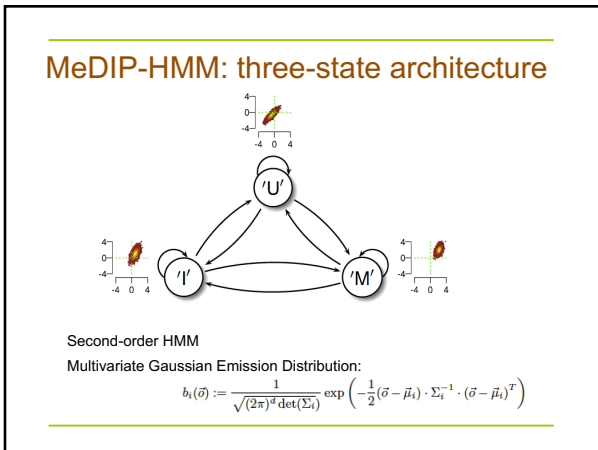
Human epigenome atlas

- Successive releases of the Atlas will provide progressively more detailed insights into locus-specific epigenomic states, including histone marks and DNA methylation marks across specific tissues and cell types, developmental stages, physiological conditions, genotypes, and disease states.



A HMM application for the inference of DNA methylation

- **MeDIP-HMM: genome-wide identification of distinct DNA methylation states from high-density tiling arrays**
- MeDIP-HMM utilizes a higher-order state-transition process improving modeling of spatial dependencies between chromosomal regions
- Enables a differentiation between **unmethylated, methylated and highly methylated** genomic regions.
- Training algorithm: a Bayesian Baum-Welch algorithm integrating prior knowledge on methylation levels.
- Application of MeDIP-HMM to the analysis of the Arabidopsis root methylome and systematically investigate the benefit of using higher-order HMMs.
- *Bioinformatics (2012) doi: 10.1093*



Chromatin-state decoding

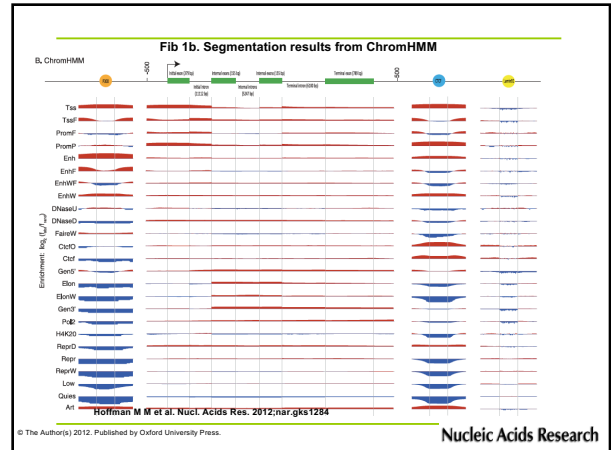
- **Automated mapping of large-scale chromatin structure in ENCODE**
– *Bioinformatics (2008) 24 (17): 1911-1916.*
- **ChromHMM: automating chromatin-state discovery and characterization**
– *Nature Methods 9, 215–216 (2012)*

Integrative annotation of chromatin elements from ENCODE data

Table 1. Major differences between ChromHMM and Segway as applied to the ENCODE data

	ChromHMM	Segway
Modeling framework	Hidden Markov model	Dynamic Bayesian network
Genomic resolution	200 bp	1 bp
Data resolution	Boolean	Real value
Handling missing data	Interpolation	Marginalization
Emission modeling	Bernoulli distribution	Gaussian distribution
Length modeling	Geometric distribution	Geometric plus hard and soft constraints
Training set	Entire genome	ENCODE regions (1%)
Decoding algorithm	Posterior decoding	Viterbi
Learning across six cell types	Single model for all cell types	One model per cell type

Ref: *Nucl. Acids Res.* (2013) 41 (2): 827-841.



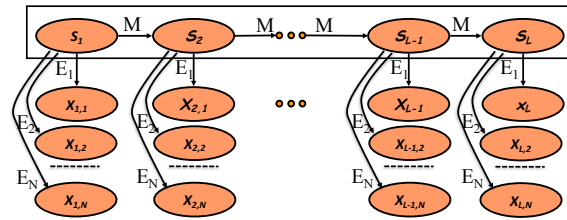
ChromHMM is a *multivariate* HMM

- ChromHMM uses a multivariate HMM that explicitly models the combination of marks

Cell	chr1	Mark1	Mark2	Mark3
0	0	0		
0	1	0		
0	0		1	

absent present

Multivariate HMM



Multivariate HMM (formal definition)

- A multivariate HMM θ has
 - N sets of observation symbols, each for one given observation sequence n ($n=1, 2, \dots, N$)
 - A set of hidden states
 - Transition probabilities a_{ij} , for any pair of hidden states i and j
 - Initial probabilities $B_j = a_{0j}$ for any hidden states j
 - N sets of emission probabilities $e_s^n(x_n)$ for the observation symbol being emitted in the n th observation sequence from the hidden state s .

Multivariate HMM

- Given N observation sequences of the same length L , $X = \{(x_{1,1} \dots x_{1,L}), \dots, (x_{N,1} \dots x_{N,L})\}$ and the hidden state sequence $S = (s_1 \dots s_L)$, the full probability from the multivariate HMM is,

$$P(S, X | \theta) = \prod_{j=1}^L \left[a_{s_{j-1}, s_j} \prod_{n=1}^N e_{s_j}^n(x_{n,j}) \right]$$

Let $e_{s_j}(x_{n,1}, \dots, x_{n,j}) = \prod_{n=1}^N e_{s_j}^n(x_{n,j})$, the multivariate HMM can be reduced to conventional HMM, except the observation symbol becomes a vector $(x_{n,1} \dots x_{n,j})$ at position j . The same algorithms for model inference (Viterbi and forward/backward) and learning can be used.