

6.047/6.878/HST.507

Computational Biology: Genomes, Networks, Evolution

Lecture 11 - Epigenomics

read mapping – peak calling – multivariate HMMs

Module III: Epigenomics and gene regulation

- **Computational Foundations**
 - L10: Gibbs Sampling: between EM and Viterbi training
 - L11: Rapid linear-time sub-string matching
 - L11: Multivariate HMMs
 - L12: Post-transcriptional regulation
- **Biological frontiers:**
 - L10: Regulatory motif discovery, TF binding
 - L11: Epigenomics, chromatin states, differentiation
 - L12: Post-transcriptional regulation

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

One Genome – Many Cell Types

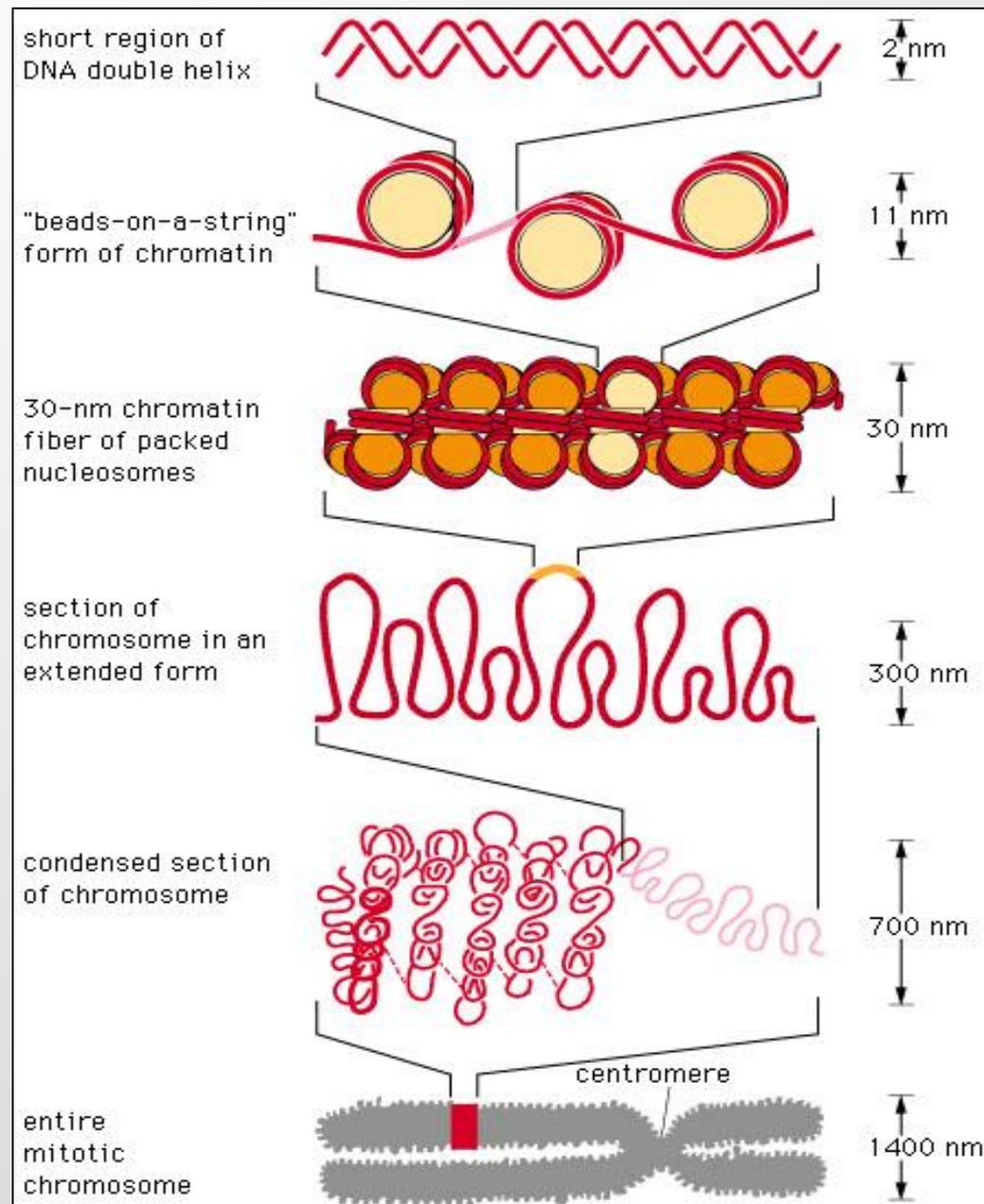
ACCAGTTACGACGGTCA
GGGTACTGATACCCCAA
ACCGTTGACCGCATTTA
CAGACGGGGTTTGGGTT
TTGCCCCACACAGGTAC
GTTAGCTACTGGTTTAG
CAATTTACCGTTACAAC
GTTTACAGGGTTACGGT
TGGGATTTGAAAAAAG
TTTGAGTTGGTTTTTTC
ACGGTAGAACGTACCGT
TACCAGTA



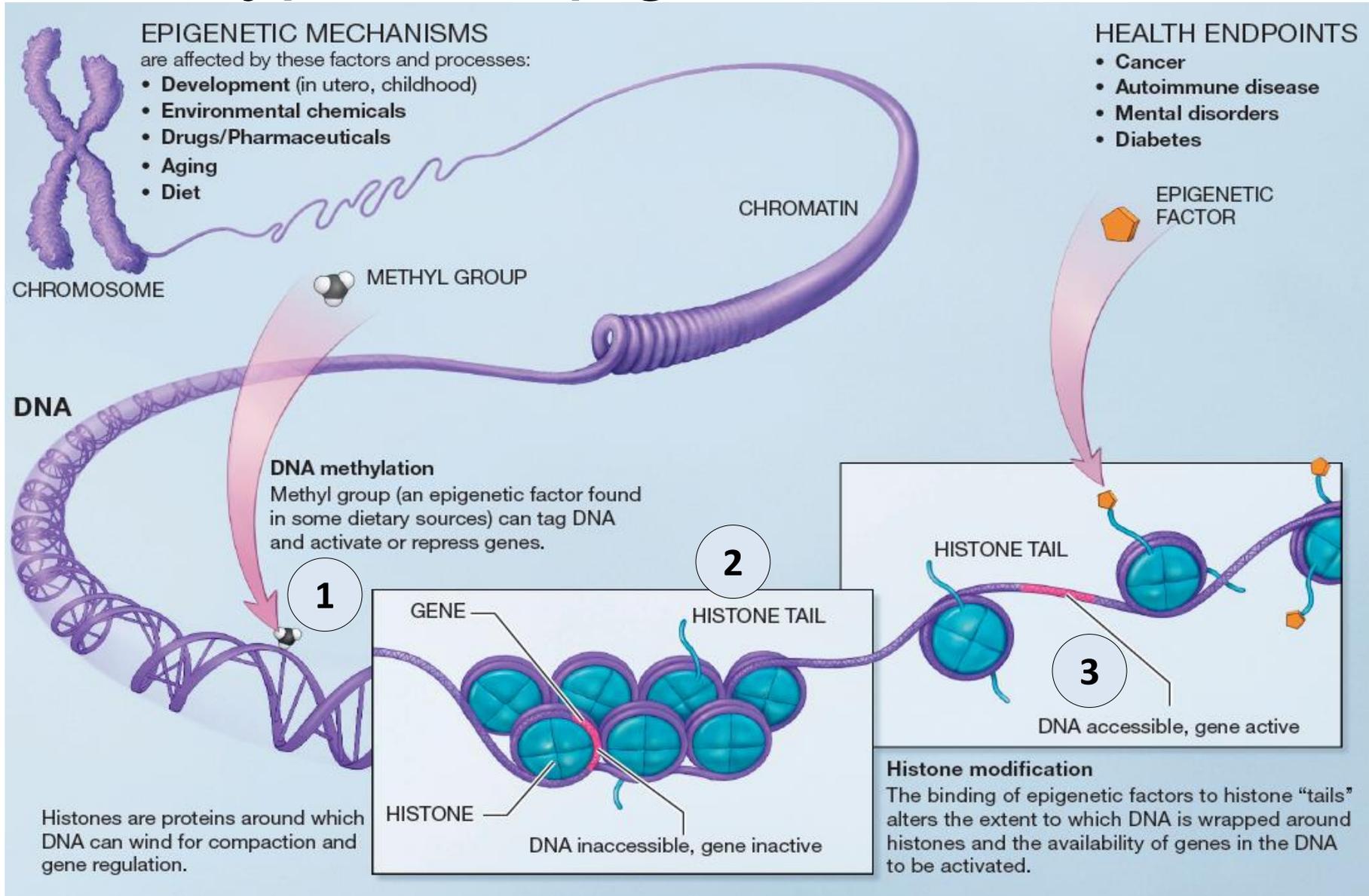
Images of skin, heart, a red blood cell, and a human brain removed due to copyright restrictions.

DNA packaging

- Why packaging
 - DNA is very long
 - Cell is very small
- Compression
 - Chromosome is 50,000 times shorter than extended DNA
- Using the DNA
 - Before a piece of DNA is used for anything, this compact structure must open locally
- Now emerging:
 - Role of accessibility
 - State in chromatin itself
 - Role of 3D interactions



Three types of epigenetic modifications

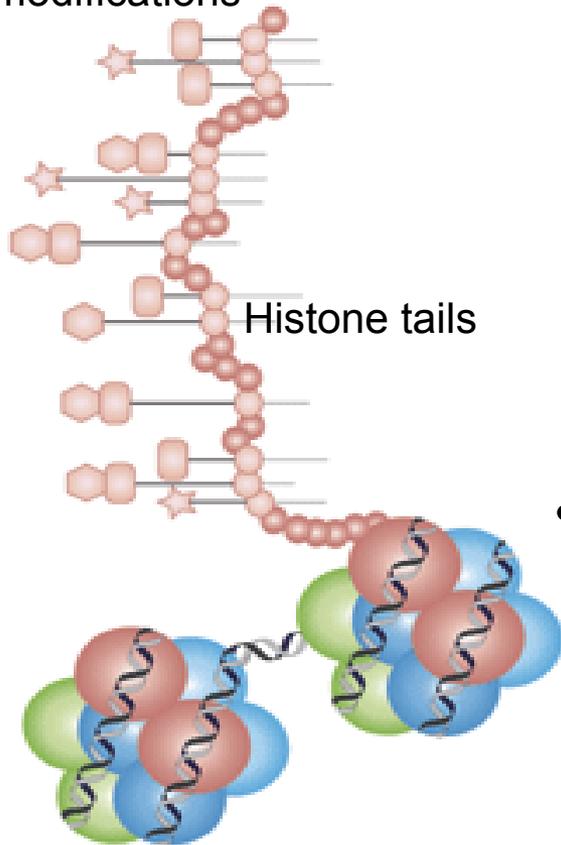


Courtesy of [National Institutes of Health](http://www.nih.gov). Image in the public domain.

Image source: <http://nihroadmap.nih.gov/epigenomics/>

100s of histone tail modifications

modifications



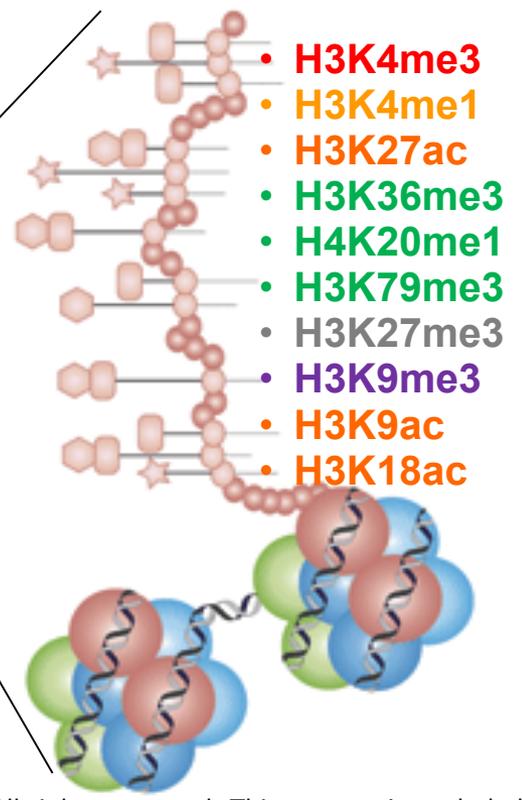
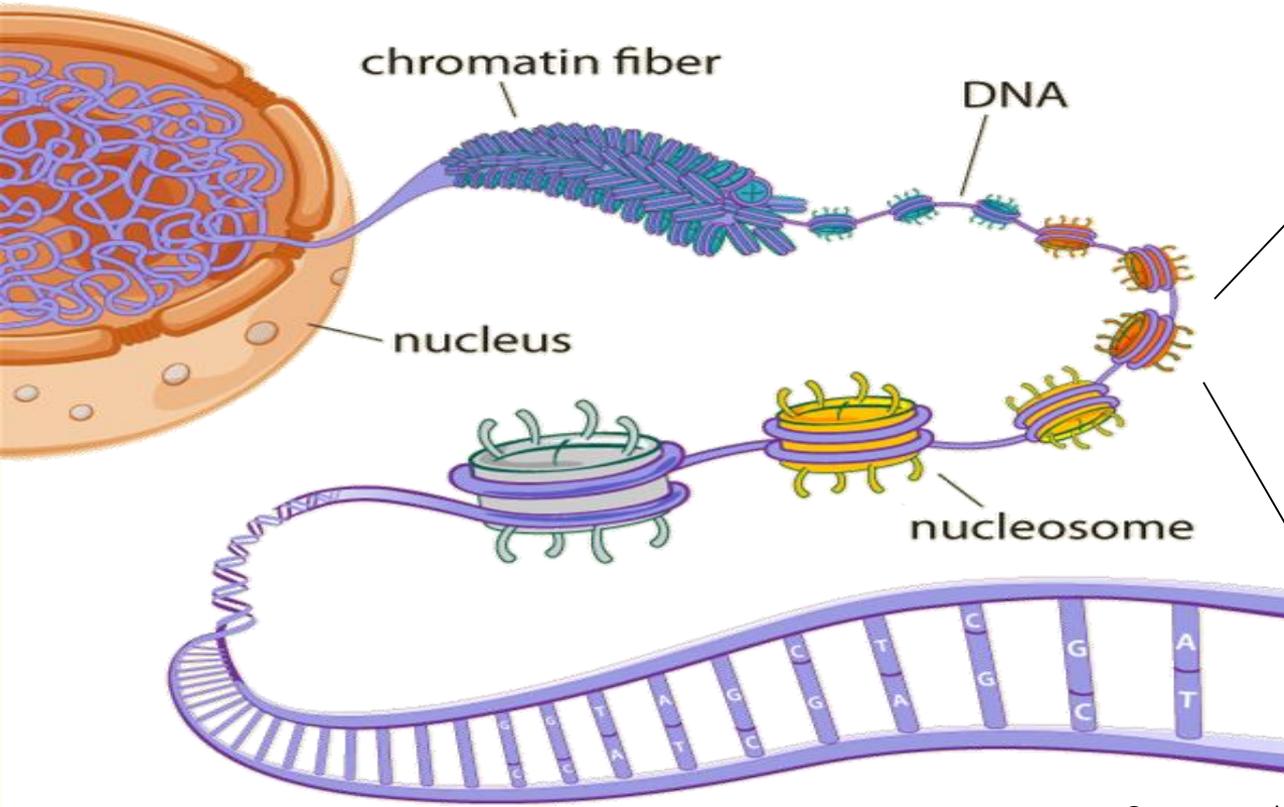
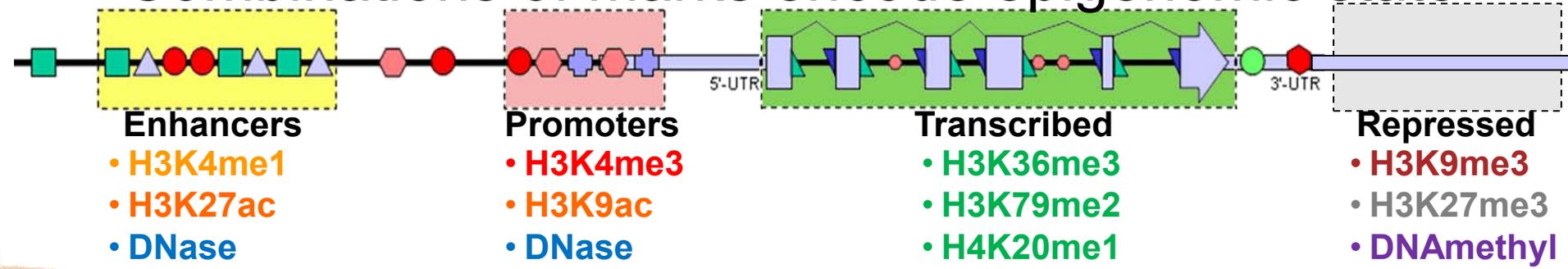
Histone tails

- 100+ different histone modifications
 - Histone protein → H3/H4/H2A/H2B
 - AA residue → Lysine4(K4)/K36...
 - Chemical modification → Met/Pho/Ubi
 - Number → Me-Me-Me(me3)
 - Shorthand: H3K4me3, H2BK5ac
- In addition:
 - DNA modifications
 - Methyl-C in CpG / Methyl-Adenosine
 - Nucleosome positioning
 - DNA accessibility
- The constant struggle of gene regulation
 - TF/histone/nucleo/GFs/Chrom compete

DNA wrapped around
histone proteins

© source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Combinations of marks encode epigenomic state

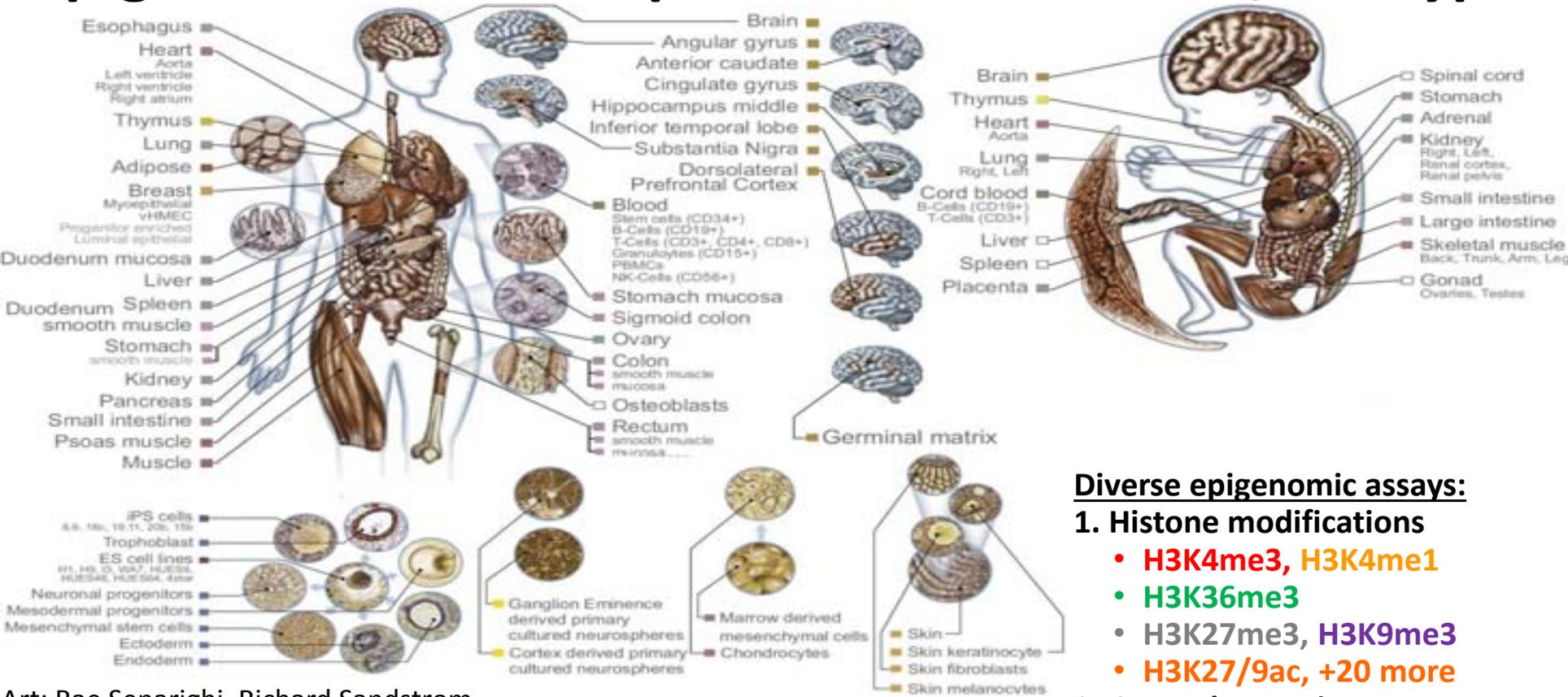


Courtesy of Broad Communications. Used with permission.

© source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- 100s of known modifications, many new still emerging
- Systematic mapping using ChIP-, Bisulfite-, DNase-Seq

Epigenomics Roadmap across 100+ tissues/cell types



Diverse epigenomic assays:

1. Histone modifications

- H3K4me3, H3K4me1
- H3K36me3
- H3K27me3, H3K9me3
- H3K27/9ac, +20 more

2. Open chromatin:

- DNase

3. DNA methylation:

- WGBS, RRBS, MRE/MeDIP

4. Gene expression

- RNA-seq, Exon Arrays

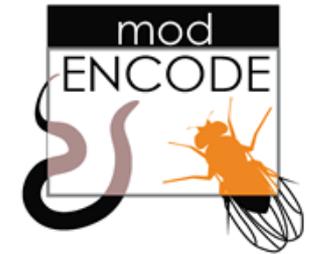
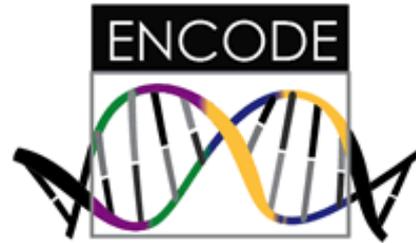
Art: Rae Senarighi, Richard Sandstrom

Courtesy of [NIH Roadmap Epigenomics Mapping Consortium](#). Used with permission.

Diverse tissues and cells:

1. Adult tissues and cells (brain, muscle, heart, digestive, skin, adipose, lung, blood...)
2. Fetal tissues (brain, skeletal muscle, heart, digestive, lung, cord blood...)
3. ES cells, iPS, differentiated cells (meso/endo/ectoderm, neural, mesench, trophobl) ⁹

Ongoing epigenomic mapping projects



- Mapping multiple modifications
 - In multiple cell types
 - In multiple individuals
 - In multiple species
 - In multiple conditions
 - With multiple antibodies
 - Across the whole genome
- First wave published
 - Lots more in pipeline
 - Time for analysis!

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

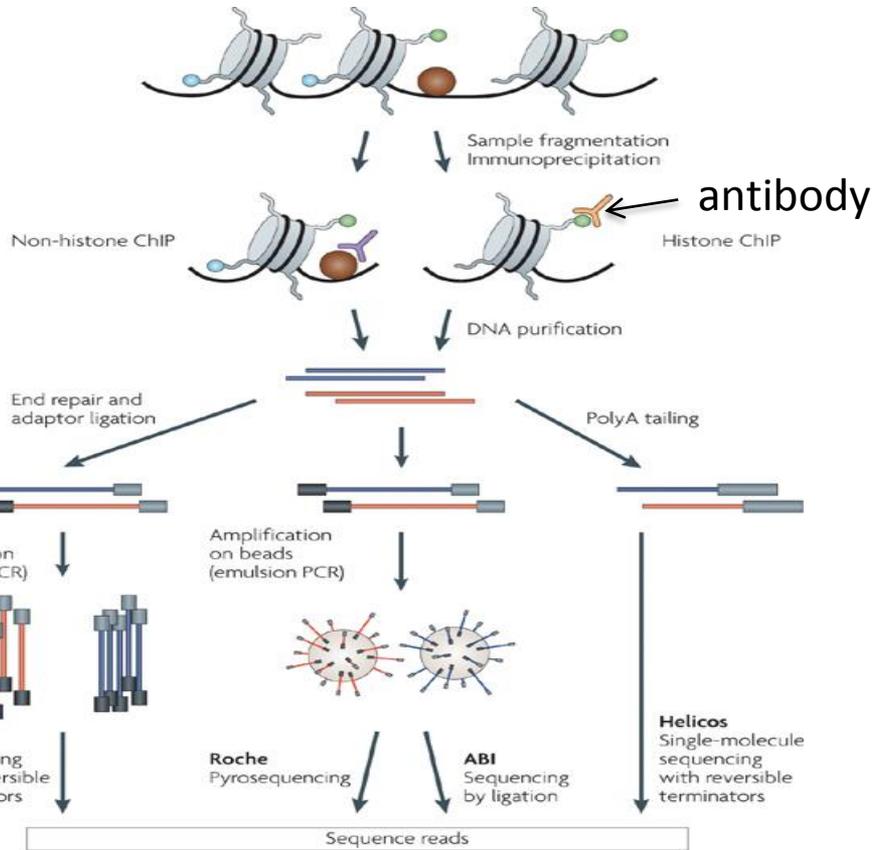
5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

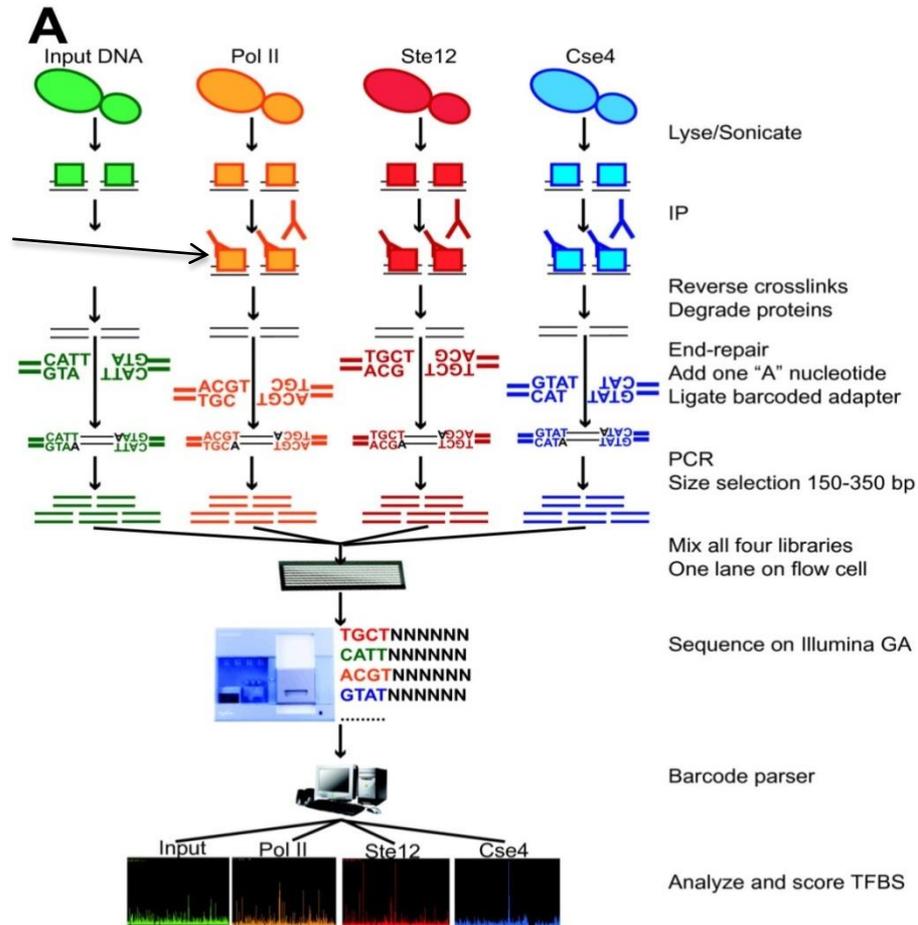
(Future: Chromatin states to interpret disease-associated variants)

ChIP-seq review

(Chromatin immunoprecipitation followed by sequencing)



Nature Reviews | Genetics

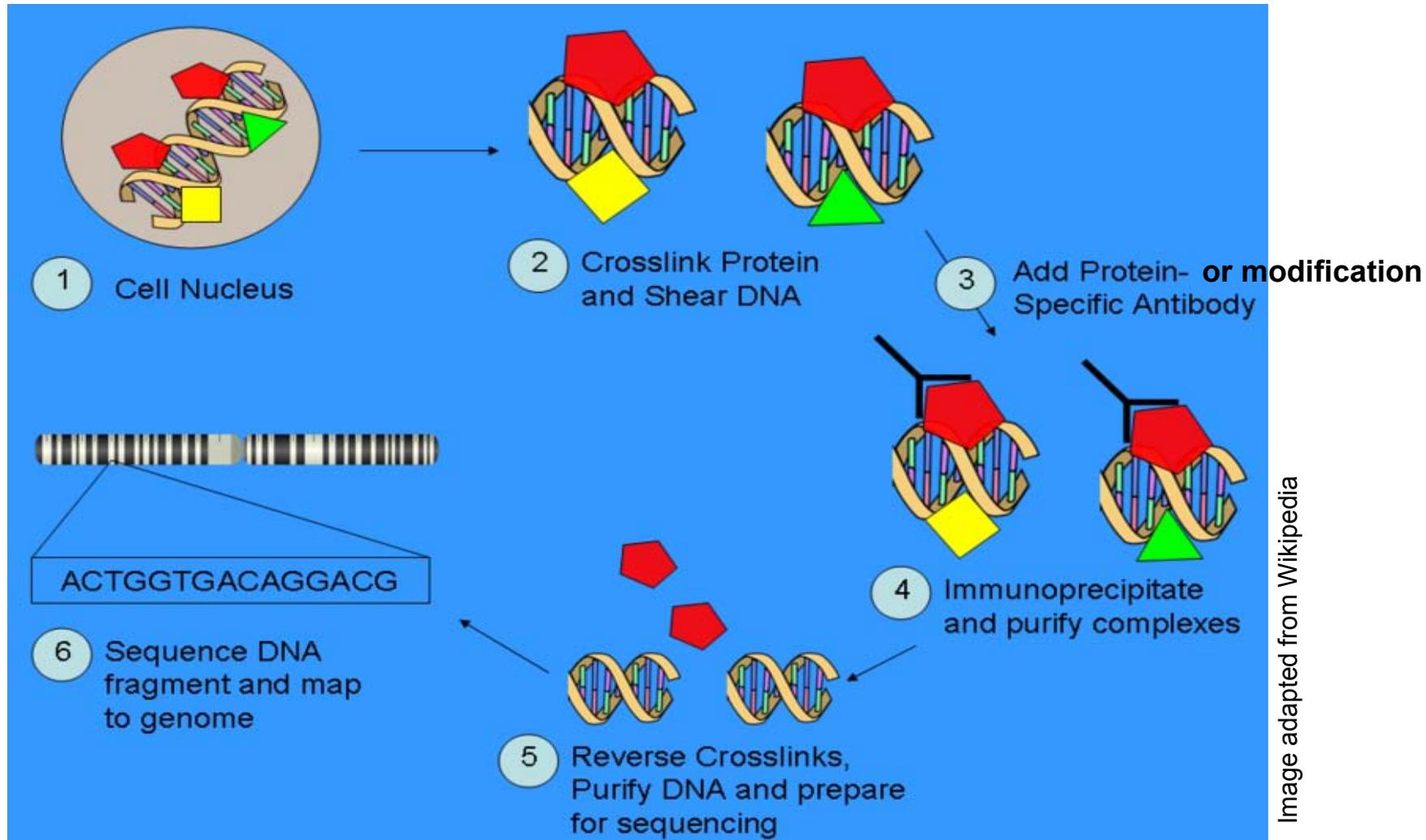


© Illumina, Inc. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Source: Lefrançois, Philippe et al. "Efficient yeast ChIP-Seq using multiplex short-read DNA sequencing." BMC genomics 10, no. 1 (2009): 1.

Bar-coded multiplexed sequencing

ChIP-chip and ChIP-Seq technology overview

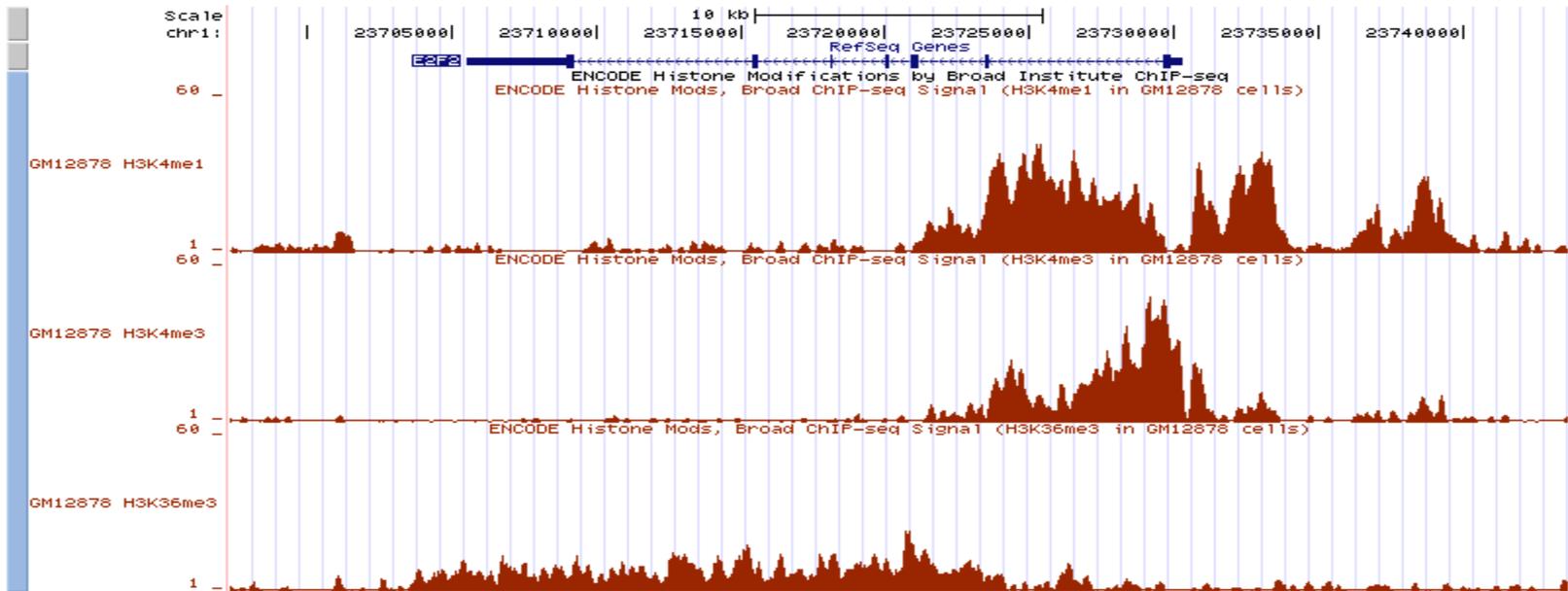


Modification-specific antibodies → Chromatin Immuno-Precipitation

followed by: ChIP-chip: array hybridization

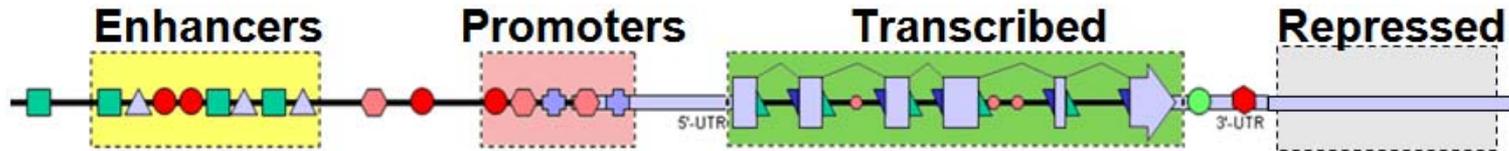
ChIP-Seq: Massively Parallel Next-gen Sequencing¹³

ChIP-Seq Histone Modifications: What the raw data looks like

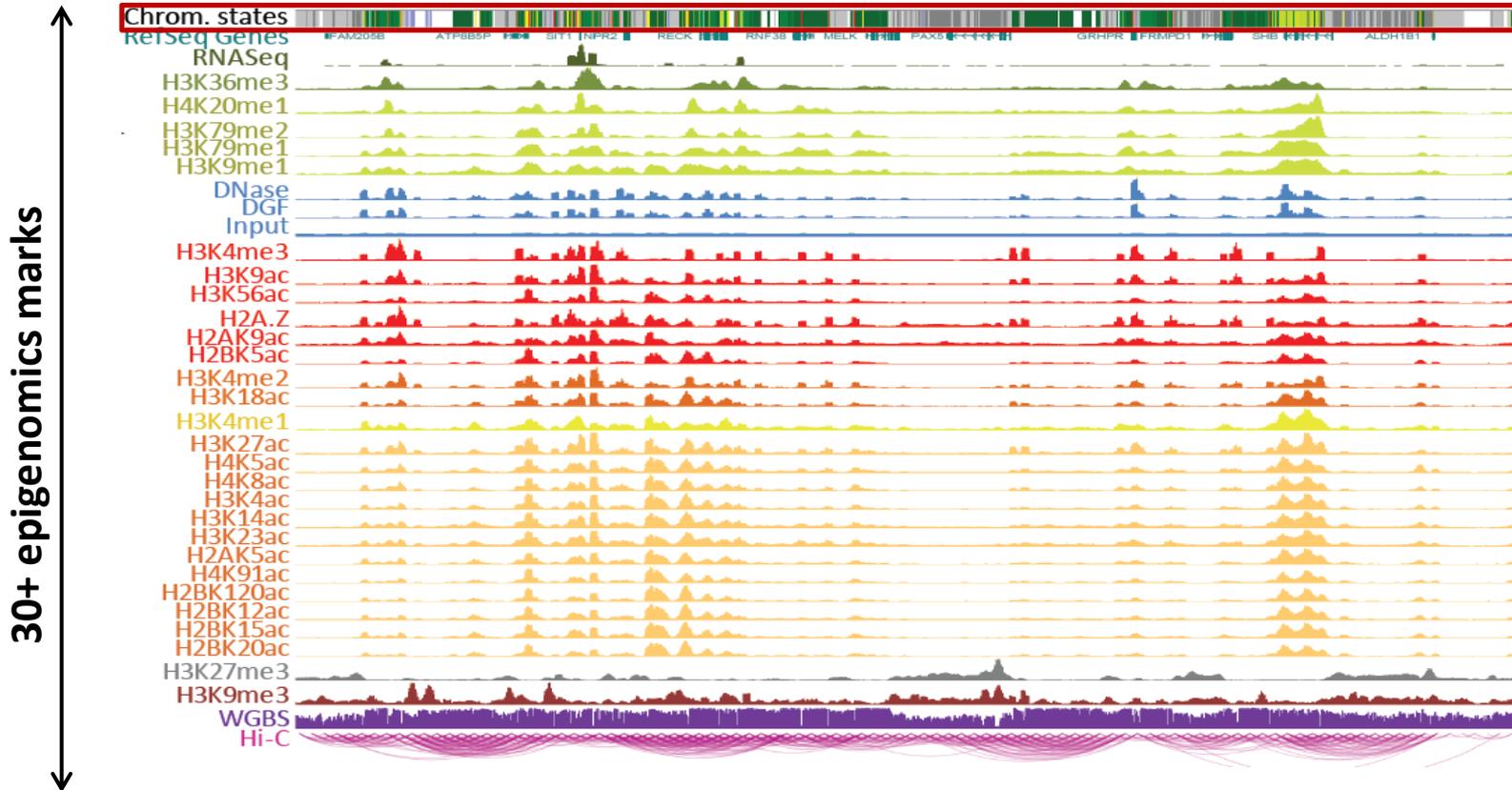


- Each sequence tag is 30 base pairs long
- Tags are mapped to unique positions in the ~3 billion base reference genome
- Number of reads depends on sequencing depth. Typically on the order of 10 million mapped reads.

Summarize multiple marks into chromatin states



Chromatin state track summary



© WashU Epigenome Browser. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

WashU Epigenome Browser

ChromHMM: multi-variate hidden Markov model

Mapping millions of short reads to the genome

Traditional Hashing Schemes

Burrows-Wheeler Transform (BWT)

Mapping Reads to the Genome

- Assign reads to best matching location in reference genome
- 10,000,000s of reads, ~30 bases long
- Example: **CAGGGCTGATTGAGGACATTCATCACG**
- Allow mismatches: sequencing errors, or SNPs
- Algorithmic and memory efficiency is critical

```
...ATAGTCTTCCCTGCATAGTCCTTTCTGCCAGACGGTAATTACAACCTTTTGTATAAAAAATAGAGAAGACTTAAAATTCTGCAGTAGGAGTGTCTGTATTCCTCCG
CAATCACTTCAATGTGTCTATTTTTGTGATCTAAAAATAACGGCTCCTGCAGATAAACTCGGATATGAGAGTTTCATAATGACAACCTAGCATATATTTGTCCAGAG
TTATTTAAAACGGTCTAGACGAGACTATCATTTTCCCTAAAATACCAAAGATTAAGTACACCGGAAGACTCAGAAAAACACCTACAGAGACCTCACAGAAGTTTCTAG
TTTAAAGTATGTGAGTGTGCACACTTTCATCTTAGTCTAAGCATCAGGGGGAACGTTGGGTAAACATTACTAAAGCTGAAACAGTGCCACGATGCCAGATATTAGG
TCATAAATATGAACTTTTTTTTTTTTGGAGATGGAGTCTTGCTCTGTTGCCAGGCTGCAGTGCAGTGGCACAATCTCAGCTCACTGCAGCCTCCGCCTCCAGGCTC
AAGCAATTCTCCTGCCTCAGCCTCCTGAGTAGCTAGGATTACAGATACCCACCACCATGCCCGGCTAATTTTTGTATTTTTTAGTAGAGACAGGGTTTTACCATGTT
GGCCAGGCTGGTCTCGAACTCCTGGCCTTAAGTGATCTGCCACCTCTGCCTTCCAAAGTGCTGGGATTACAGGCCTGAGCCATCGCGCCTGGCTATAAGTATGAA
CTTTTAAAGAATCTAGAAATGAGGCCCTCCAAAAAGAGATGAGCTGGTAACAGAGCCGAACACACAGAAAAATAGTTTCAGGAAGGGCCTGGGCAGAGGAAGGCCTAA
TAAGCAAGGAAGCCACAAACATGTAGCCCAGCAATACACACACACAAACAATTCCTACATGCAGAGCCCTTTAGGAATGGCAGACCTTTGTTTCTACAACAGATGA
AGCTGTGAATAGCCTAAAGAACAATGCTCCTGGGGGTGGCCTGTAGAGTGTGATAAAAGTCTGAATAAAACGGGCTGGGTGGAGCTGGATGATCACGTGTGTGGT
TCCACAGGGTGAAGACAGCATCCGGTTCACAGTCACAGGTTCTGTGTGTAAGGCGTGCATGTGGAGAAACGCCTTTGAGGAAAAGGCGTGTGAAAGGGTCTTTGGGG
GGGACGGGCTAGACACAGGCTCAGAGAAGTGGATGGTTCTCAGGATGCAGATGAGTGTGGTAACTGGAGTCTAAATCCAGTGGTAAGACTGTGCTGTCAAGAGACA
CTGGGGTGACACAGGGCAAATGGAGGCAGAAGAGCAGGTCCCACCTGAAGAAGGGCTCAGGGGCTGGAATCTAGGGCAGGAAGTCCAGGCTGAGAGCCTGCCACAGGC
TGGTATGGTGCCATCTTAAGCAGGAAGAACTCGCACAAGCCCCACCCAGGGGTGGAGTGTGTGGTGACTGTGGGCACCCAGAGACACCCAGGGAGGATTGGCT
GAGGGGGAAAGGAGGAGATTCACTGGACCTGATACCCCTCCGCCTAAGATGGGGGGCTCTACTGGATGGACTCTGAAGCTAGGATGGGATCCTAAAGTGGCTCTGT
TTGCCCCGTGCCACCCTGTCTAACATGGGACCTACAAGCGGGCCCTGCCCTGCCAGGGCCCAGGAAGCTCTCCCCGCTCCTATGTCTGTTTTCCCTCCAGGTCC
ACTCACCCCATGAGACTCAAAGGCCCTTTCAGGACAAAGACAATCGCTTCACCATTTCTTCTTCAACTCCTGGCACAGAGTCTGGCCACTGGGAGACACCCAGCC
AATAAGGCAAGGAGAGAGGACTGAGGAGGGAAGGGGGCAGATCAAGTGATGAGAAGATCCCTCTTTAGAATCAGGTGGGGCCCTCGCACAGAAAGGGCGGCCTCC
CCCACAGGAACCCAGGGCAGGTCCAGAGCAGCAGGAAGGAGGAGGCGGCCAATGGGAAGGCAACCGAGCCCCAGGGACACACTGCGTCCATCGTGGCTCCTGAGG
GATGGGCCACCCACTTCCGACCCGGCCACTAGAACCTGCTTTTCAGTTTTGTTTATGCTCTGAGCACTGGGGTCTCAGCCCCCTCTCTTCCCTCAAGGAGGCTGT
TGTCTCTTGGTTCCTGCTGTGGGGCAGCTATGAATTTACGATGCTCAGGGCTGATTGAGGACATTCATCAGGATATCGGGGAAAAGAATGGAGAATCAAAACAGTAA
GAAAAAAGTCTGAAATACCTTCCAAGTCTATTTTCATAGCCTTGGAAAACATAACAATAAATTTACTTTTATGTCTACCTTTGAAAATTATCTTAACATAGATGCCAA
TTTCAAACCCTCCCAGTACTGGGAGACAAATGGCATACTGGTTTCTACAAGCCTCCTTCATTCATCTGCTAACTGTGAAGGCCTCATCTCTGAACGCCCAGGGCC
GGGACCCGTGCCCTGGATCAGGCAGGATGCTCAATACGCGGTTGTGAGATGAGTAACAGGCAGACACCGTAGAACAGCACTTGTGAGGCCTGCTGATT...
```

How would you do it:

- L2: Sequence alignment: $O(m*n)$
- L3: Hashing / BLAST: $O(m+n)$
 - Solution until 2008 (e.g. MAQ, Li et al, GR 2008)
- Other advanced algorithms:
 - Linear-time string matching: $O(m+n)$. L3 addendum
 - Suffix trees and suffix arrays: $O(m)$. L13 addendum
- Challenge: memory requirements
 - Hash table, suffix tree/array require $O(m*n)$ space
- Today: Burrows-Wheeler transformation $O(m)$
 - Ultrafast/memory efficient. New norm since 2009.
 - Introduced in: Bowtie (Langmead GB 2009).

Second Generation Mappers have Leveraged the Burrows Wheeler Transformation

Software

Ultrafast and memory-efficient alignment of short DNA sequences to the human genome

Ben Langmead, Cole Trapnell, Mihai Pop and Steven L Salzberg

Open Access

Address: Center for Bioinformatics and Computational Biology, Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA.

Correspondence: Ben Langmead. Email: langmead@es.umd.edu

Published: 4 March 2009

Received: 21 October 2008

Genome Biology 2009, 10:R25 (doi:10.1186/gb-2009-10-3-r25)

Revised: 19 December 2008

Accepted: 4 March 2009

The electronic version of this article is the complete one and can be found online at <http://genomebiology.com/2009/10/3/R25>

© 2009 Langmead et al.; licensee BioMed Central Ltd.

This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract



Bowtie is an ultrafast, memory-efficient alignment program for aligning short DNA sequence reads to large genomes. For the human genome, Burrows-Wheeler indexing allows Bowtie to align more than 25 million reads per CPU hour with a memory footprint of approximately 1.3 gigabytes. Bowtie extends previous Burrows-Wheeler techniques with a novel quality-aware backtracking algorithm that permits mismatches. Multiple processor cores can be used simultaneously to achieve even greater alignment speeds. Bowtie is open source <http://bowtie.cbcb.umd.edu>.

BIOINFORMATICS APPLICATIONS NOTE

Vol. 24 no. 5 2008, pages 713-714
doi:10.1093/bioinformatics/btn025

Sequence analysis

SOAP: short oligonucleotide alignment program

Ruiqiang Li^{1,2}, Yingrui Li¹, Karsten Kristiansen² and Jun Wang^{1,2,*}

¹Beijing Genomics Institute at Shenzhen, Shenzhen 518083, China and ²Department of Biochemistry and Molecular Biology, University of Southern Denmark, Odense M, DK-5230, Denmark

Received on November 10, 2007; revised on December 20, 2007; accepted on January 14, 2008

Advance Access publication January 28, 2008

Associate Editor: Keith Cantall

ABSTRACT

Summary: We have developed a program SOAP for efficient gapped and ungapped alignment of short oligonucleotides onto reference sequences. The program is designed to handle the huge amounts of short reads generated by parallel sequencing using the new generation Illumina-Solexa sequencing technology. SOAP is compatible with numerous applications, including single-read or pair-end sequencing, small RNA discovery and mRNA tag sequence mapping. SOAP is a command-driven program, which supports multi-threaded parallel computing, and has a batch module for multiple query sets.

Availability: <http://soap.genomics.org.cn>

Contact: soap@genomics.org.cn

SOAP will allow either a certain number of mismatches or one continuous gap for aligning a read onto the reference sequence. The best hit of each read which has minimal number of mismatches or smaller gap will be reported. For multiple equal-best hits, the user can instruct the program to report all, or randomly report one, or disregard all of them. Since the typical read length is 25-50bp, hits with too many mismatches are unreliable which are hard to distinguish with random matches. By default, the program will allow at most two mismatches. Between two haplotype genome sequences, occurrence of single nucleotide polymorphism is much higher than that of small insertions or deletions, so ungapped hits have precedence over gapped hits. For gapped alignment only one continuous gap with a size ranging from 1 to 3bp is accepted, while no

BIOINFORMATICS APPLICATIONS NOTE

Vol. 25 no. 15 2008, pages 1966-1967
doi:10.1093/bioinformatics/btp036

Sequence analysis

SOAP2: an improved ultrafast tool for short read alignment

Ruiqiang Li^{1,2,*}, Chang Yu¹, Yingrui Li¹, Tak-Wah Lam³, Siu-Ming Yiu², Karsten Kristiansen² and Jun Wang^{1,2,*}

¹Beijing Genomics Institute at Shenzhen, Shenzhen, 518083, China, ²Department of Biochemistry and Molecular Biology, University of Southern Denmark, Odense M, DK-5230, Denmark and ³Department of Computer Science, University of Hong Kong, Hong Kong, China

Received on January 23, 2008; revised on April 27, 2008; accepted on May 24, 2008

Advance Access publication June 3, 2008

Associate Editor: Joaquin Dopazo

ABSTRACT

Summary: SOAP2 is a significantly improved version of the short oligonucleotide alignment program that both reduces computer memory usage and increases alignment speed at an unprecedented rate. We used Burrows-Wheeler Transformation (BWT) compression index to substitute the seed strategy for indexing the reference sequence in the main memory. We tested it on the whole human genome and found that this new algorithm reduced memory usage from 147 to 5.4GB and improved alignment speed by 20-30 times. SOAP2 is compatible with both single- and paired-end reads. Additionally, this tool now supports multiple text and compressed file formats. A consensus builder has also been developed for consensus assembly and SNP detection from alignment of short reads on a reference genome.

Availability: <http://soap.genomics.org.cn>

Contact: soap@genomics.org.cn

variation map, will generate about 15Tb of sequence using next-generation sequencing technologies. With even the fastest programs currently available, one would need ~1000 CPU months to align these short reads onto the human reference genome. Additionally, new methods are now needed to support longer reads as the existing methods were primarily designed for very short reads with typical lengths shorter than 50bp. With improvements in sequencing chemistry and data processing algorithms, the Illumina Genome Analyzer can now generate up to 75-100bp high-quality reads, and longer reads are expected in the near future.

Here, we have developed an improved version of SOAP, called SOAP2. The new program uses the Burrows-Wheeler Transformation (BWT) compressed index instead of the seed algorithm that was used in the previous version for indexing the reference sequence in the main memory. Use of BWT substantially improved alignment speed; additionally, it significantly reduced memory usage.

BIOINFORMATICS ORIGINAL PAPER

Vol. 25 no. 14 2008, pages 1754-1760
doi:10.1093/bioinformatics/btn024

Sequence analysis

Fast and accurate short read alignment with Burrows-Wheeler transform

Heng Li and Richard Durbin*

Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK

Received on February 20, 2008; revised on May 12, 2008; accepted on May 12, 2008

Advance Access publication May 15, 2008

Associate Editor: John Quackenbush

ABSTRACT

Motivation: The enormous amount of short reads generated by the new DNA sequencing technologies call for the development of fast and accurate read alignment programs. A first generation of hash-table-based methods has been developed, including MAQ, which is accurate, feature rich and fast enough to align short reads from a single individual. However, MAQ does not support gapped alignment for single-end reads, which makes it unsuitable for alignment of longer reads where indels may occur frequently. The speed of MAQ is also a concern when the alignment is scaled up to the resequencing of hundreds of individuals.

Results: We implemented Burrows-Wheeler Alignment tool (BWA), a new read alignment package that is based on backward search with Burrows-Wheeler Transform (BWT), to efficiently align short sequencing reads against a large reference sequence such as the human genome, allowing mismatches and gaps. BWA supports both large space reads, e.g. from Illumina sequencing machines, and color space reads from ABI SOLiD machines. Evaluations on both simulated and real data suggest that BWA is ~10-20x faster than MAQ, while achieving similar accuracy. In addition, BWA outputs alignments in the new standard SAM (Sequence Alignment/Map) format. Source code and binary files are available at <http://code.google.com/p/bwa/>.

Availability: <http://maq.sourceforge.net>

Contact: rd@sanger.ac.uk

of scanning the whole genome when few reads are aligned. The second category of software, including SOAP1 (Li et al., 2008b), PASS (Carpignani et al., 2009), MOM (Eaves and Cai, 2009), ProbeMatch (Jiang, Kim et al., 2009), NovoAlign (<http://www.novocell.com>), RESEQ (<http://bioinformatics.hawaii.edu/novocell/RESEQ/>), Mosaik (<http://bioinformatics.hawaii.edu/novocell/Mosaik/>) and BFAST (<http://genome.sci.uci.edu/bfast/>), handle the genome. These programs can be easily parallelized with multi-threading, but they usually require large memory to build an index for the human genome. In addition, the iterative strategy frequently introduced by these software may make their speed sensitive to the sequencing error rate. The third category includes slider (Mullis et al., 2009) which does alignment by merge-sorting the reference subsequences and read sequences.

Recently, the theory on string matching using Burrows-Wheeler Transform (BWT) (Burrows and Wheeler, 1994) has drawn the attention of several groups, which has led to the development of SOAP2 (<http://soap.genomics.org.cn/>), Bowtie (Langmead et al., 2009) and BWA, our new aligner described in this article. Essentially, using backward search (Ferragina and Manzoni, 2000; Lipson, 2005) with BWT, we are able to effectively mimic the top-down traversal on the prefix tree of the genome with relatively small memory footprint (Lam et al., 2008) and to count the number of exact hits of a string of length l in $O(l)$ time independent of the size of the genome. For exact search, BWA samples from the implicit prefix tree the distinct substrings that are less than k edit distance

Mapping short DNA sequencing reads and calling variants using mapping quality scores

Heng Li,¹ Jue Ruan,² and Richard Durbin^{1,3}

¹The Wellcome Trust Sanger Institute, Hinxton CB10 1SA, United Kingdom; ²Beijing Genomics Institute, Chinese Academy of Science, Beijing 100029, China

New sequencing technologies promise a new era in the use of DNA sequence. However, some of these technologies produce very short reads, typically of a few tens of base pairs, and to use these reads efficiently requires new algorithms and software. In particular, there is a major issue in efficiently aligning short reads to a reference genome and handling ambiguity or lack of accuracy in this alignment. Here we introduce the concept of mapping quality, a measure of the confidence that a read actually comes from the position it is aligned to by the mapping algorithm. We describe the software MAQ that can build assemblies by mapping shotgun short reads to a reference genome, using quality scores to derive genotype calls of the consensus sequence of a diploid genome, e.g., from a human sample. MAQ makes full use of mate-pair information and estimates the error probability of each read alignment. Error probabilities are also derived for the final genotype calls, using a Bayesian statistical model that incorporates the mapping qualities, error probabilities from the raw sequence quality scores, sampling of the two haplotypes, and an empirical model for correlated errors at a site. Both read mapping and genotype calling are evaluated on simulated data and real data. MAQ is accurate, efficient, versatile, and user-friendly. It is freely available at <http://maq.sourceforge.net>.

(Supplemental material is available online at www.genome.org. Short-read sequences have been deposited in the European Read Archive (ERA) under accession no. ERA000002 (<http://ftp.era.ac.uk/ERA000002/>).

18185-1828 (2008) by Cold Spring Harbor Laboratory Press. ISSN 1088-9651/08 www.genome.org

Genome Research

www.genome.org

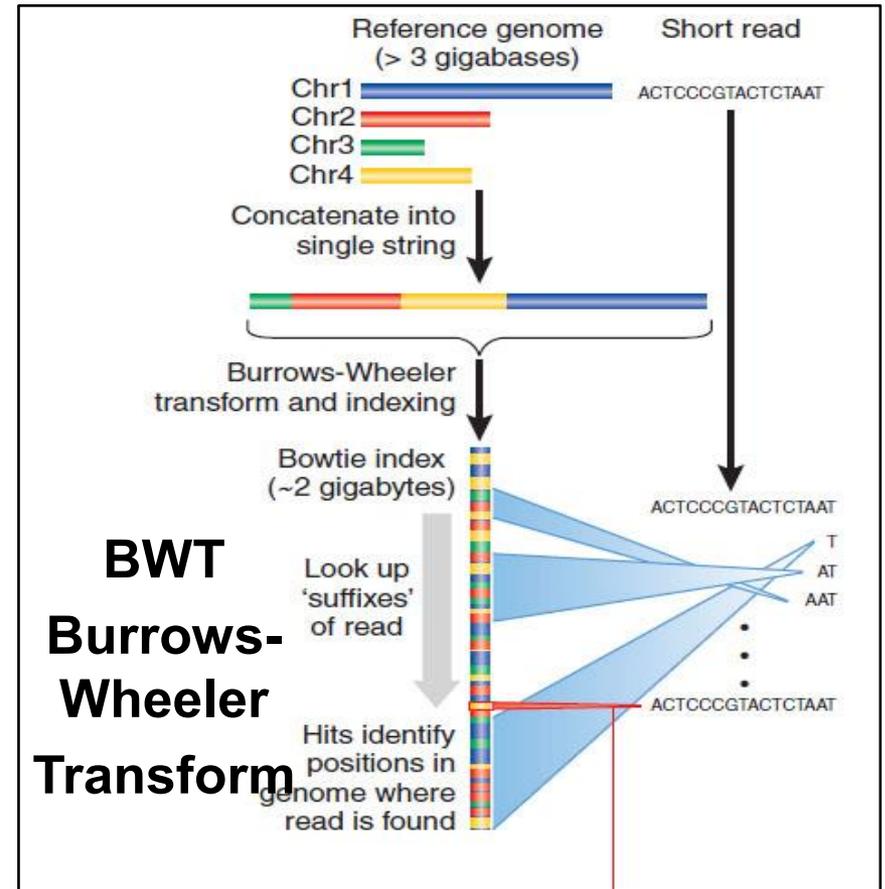
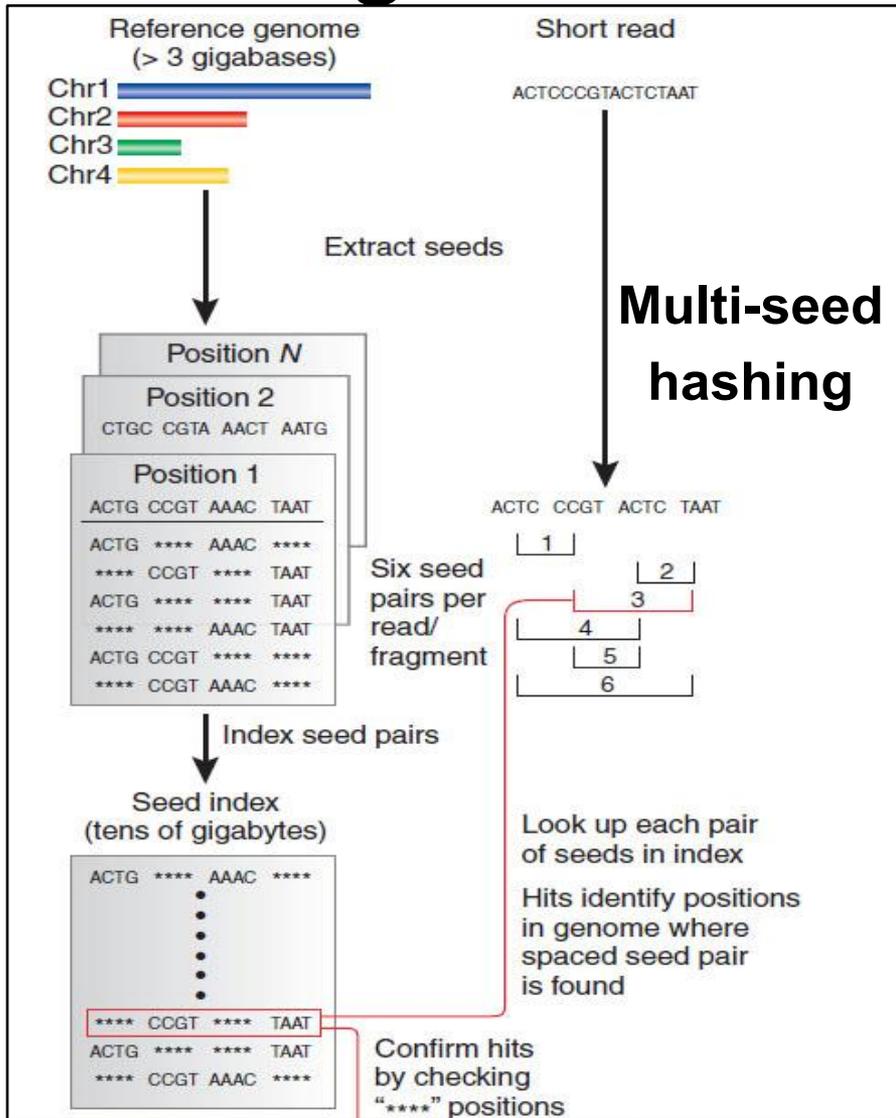
18185-1828 (2008) by Cold Spring Harbor Laboratory Press. ISSN 1088-9651/08 www.genome.org

Genome Research

www.genome.org

“...35 times faster than Maq and 300 times faster than SOAP under the same conditions”

Hashing vs. Burrows Wheeler Transform



Today: How does the BW transform actually work?

Convert each hit back to genome location



Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Trapnell, Cole and Steven L. Salzberg. "How to map billions of short reads onto genomes." Nature Biotechnology 27, no. 5 (2009): 455.

Burrows-Wheeler Transform (BWT)

<http://www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.pdf>

- Transform: \wedge BANANA@ INTO: BNN \wedge AA@A

function BWT (*string s*)

create a table, rows are all possible rotations of s
sort rows alphabetically

return (last column of the table)

All Rotations	Sorted List of Rotations	Output Last Column
\wedge BANANA@	ANANA@ \wedge B	BNN \wedge AA@A
@ \wedge BANANA	ANA@ \wedge BAN	
A@ \wedge BANAN	A@ \wedge BANAN	
NA@ \wedge BANA	BANANA@ \wedge	
ANA@ \wedge BAN	NANA@ \wedge BA	
NANA@ \wedge BA	NA@ \wedge BANA	
ANANA@ \wedge B	\wedge BANANA@	
BANANA@ \wedge	@ \wedge BANANA	

- Reversible

function inverseBWT (*string s*)

create empty table

repeat length(s) times

insert s as a column of table before first column of the table // first insert creates first column
sort rows of the table alphabetically

return (row that ends with the 'EOF' character)

Last column only suffices to reconstruct entire matrix, and thus recover original string

Add 1	Sort 1	Add 2	Sort 2	Add 3	Sort 3	Add 4	Sort 4	Add 5	Sort 5	Add 6	Sort 6	Add 7	Sort 7	Add 8	Sort 8
B	A	BA	AN	BAN	ANA	BANA	ANAN	BANAN	ANANA	BANANA	ANANA@	BANANA@	ANANA@ \wedge	BANANA@ \wedge	ANANA@ \wedge B
N	A	NA	AN	NAN	ANA	NANA	ANA@	NANA@	ANA@ \wedge	NANA@ \wedge	ANA@ \wedge B	NANA@ \wedge B	ANA@ \wedge BA	NANA@ \wedge BA	ANA@ \wedge BAN
N	A	NA	A@	NA@	A@ \wedge	NA@ \wedge	A@ \wedge B	NA@ \wedge B	A@ \wedge BA	NA@ \wedge BA	A@ \wedge BAN	NA@ \wedge BAN	A@ \wedge BANA	NA@ \wedge BANA	A@ \wedge BANAN
\wedge	B	\wedge B	BA	\wedge BA	BAN	\wedge BAN	BANA	\wedge BANA	BANAN	\wedge BANAN	BANANA	\wedge BANANA	BANANA@ \wedge	\wedge BANANA@	BANANA@ \wedge
A	N	AN	NA	ANA	NAN	ANAN	NANA	ANANA	NANA@	ANANA@	ANANA@ \wedge	ANANA@ \wedge	NANA@ \wedge B	ANANA@ \wedge B	NANA@ \wedge BA
A	N	AN	NA	ANA	NA@	ANA@	NA@ \wedge	ANA@ \wedge	NA@ \wedge B	ANA@ \wedge B	NA@ \wedge BA	ANA@ \wedge BA	NA@ \wedge BAN	ANA@ \wedge BAN	NA@ \wedge BANA
@	\wedge	@ \wedge	\wedge B	@ \wedge B	\wedge BA	@ \wedge BA	\wedge BAN	@ \wedge BAN	\wedge BANA	@ \wedge BANA	\wedge BANAN	@ \wedge BANAN	\wedge BANANA	@ \wedge BANANA	\wedge BANANA@
A	@	A@	@ \wedge	A@ \wedge	@ \wedge B	A@ \wedge B	@ \wedge BA	A@ \wedge BA	@ \wedge BAN	A@ \wedge BAN	@ \wedge BANA	A@ \wedge BANA	@ \wedge BANAN	A@ \wedge BANAN	@ \wedge BANANA
last	1st col	pairs	2nd col	triples	3rd col	4mers	4th col	5mers	5th col	6-mers	6th col	7-mers	7th col	8-mers	Full matrix

Searching for an Exact Match

e.g. Searching for **OLIS**

In MANOLISKELLIS

For simplicity (here):

- only exact matches
- Show entire matrix

In practice: only pointers

Algorithm 3 EXACTMATCH($P[1, p]$)

```

1:  $c \leftarrow P[p]$ 
2:  $sp \leftarrow C[c] + 1$ 
3:  $ep \leftarrow C[c + 1] + 1$ 
4:  $i \leftarrow p - 1$ 
5: while  $sp < ep$  and  $i \geq 1$  do
6:    $c \leftarrow P[i]$ 
7:    $sp \leftarrow C[c] + \text{Occ}(c, sp) + 1$ 
8:    $ep \leftarrow C[c] + \text{Occ}(c, ep) + 1$ 
9:    $i \leftarrow i - 1$ 
10: end while
11: return  $sp, ep$ 

```

P is the input substring

$C[c]$ – is how many characters occur before c lexographically in the genome

$\text{Occ}(c, k)$ is the number of occurrence of the character c before index k in the far right column

OLIS

```

1. $MANOLISKELLIS
2. ANOLISKELLIS$M
3. ELLIS$MANOLISK
4. IS$MANOLISKELL
5. ISKELLIS$MANOL
6. LIS$MANOLISKEL
7. LISKELLIS$MANO
8. LLIS$MANOLISKE
9. KELLIS$MANOLIS
10. MANOLISKELLIS$
11. NOLISKELLIS$MA
12. OLISKELLIS$MAN
13. S$MANOLISKELLI
14. SKELLIS$MANOLI

```

OLIS

```

1. $MANOLISKELLIS
2. ANOLISKELLIS$M
3. ELLIS$MANOLISK
4. IS$MANOLISKELL
5. ISKELLIS$MANOL
6. LIS$MANOLISKEL
7. LISKELLIS$MANO
8. LLIS$MANOLISKE
9. KELLIS$MANOLIS
10. MANOLISKELLIS$
11. NOLISKELLIS$MA
12. OLISKELLIS$MAN
13. S$MANOLISKELLI
14. SKELLIS$MANOLI

```

OLIS

```

1. $MANOLISKELLIS
2. ANOLISKELLIS$M
3. ELLIS$MANOLISK
4. IS$MANOLISKELL
5. ISKELLIS$MANOL
6. LIS$MANOLISKEL
7. LISKELLIS$MANO
8. LLIS$MANOLISKE
9. KELLIS$MANOLIS
10. MANOLISKELLIS$
11. NOLISKELLIS$MA
12. OLISKELLIS$MAN
13. S$MANOLISKELLI
14. SKELLIS$MANOLI

```

OLIS

```

1. $MANOLISKELLIS
2. ANOLISKELLIS$M
3. ELLIS$MANOLISK
4. IS$MANOLISKELL
5. ISKELLIS$MANOL
6. LIS$MANOLISKEL
7. LISKELLIS$MANO
8. LLIS$MANOLISKE
9. KELLIS$MANOLIS
10. MANOLISKELLIS$
11. NOLISKELLIS$MA
12. OLISKELLIS$MAN
13. S$MANOLISKELLI
14. SKELLIS$MANOLI

```



Key properties of Burrows-Wheeler Transform

- **Very little memory usage. Same as input (or less)**
 - Don't represent matrix, or strings, just pointers
 - Encode: Simply sort pointers. Decode: follow pointers
- **Original application: string compression (bZip2)**
 - Runs of letters compressed into (letter, runlength) pairs
- **Bioinformatics applications: substring searching**
 - Achieve similar run time as hash tables, suffix trees
 - But: very memory efficient → practical speed gains
- **Mapping 100,000s of reads: only transform once**
 - Pre-process once; read counts in transformed space.
 - Reverse transform once, map counts to genome coords

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

Quality control metrics

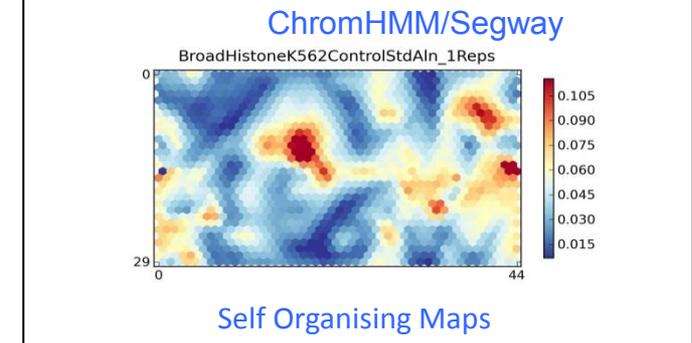
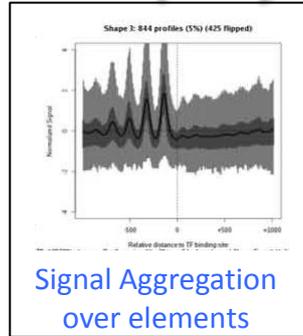
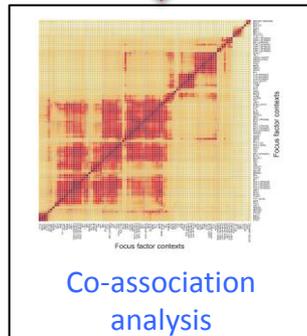
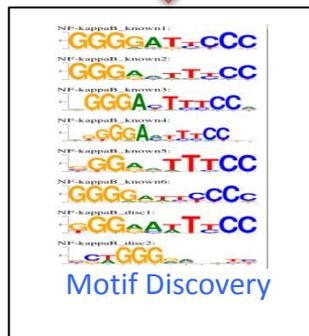
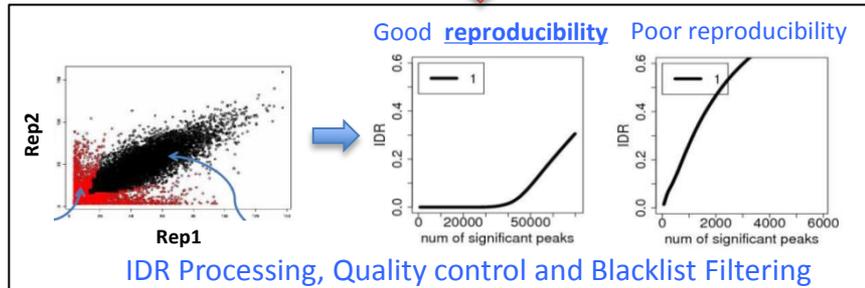
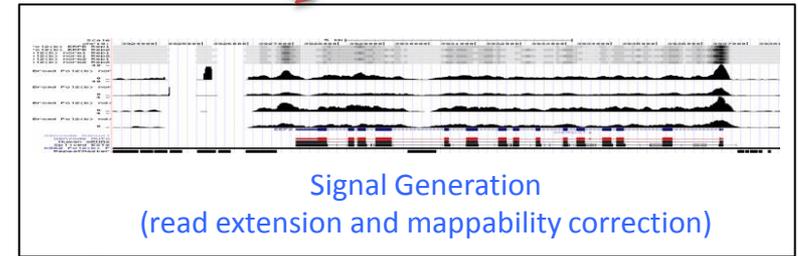
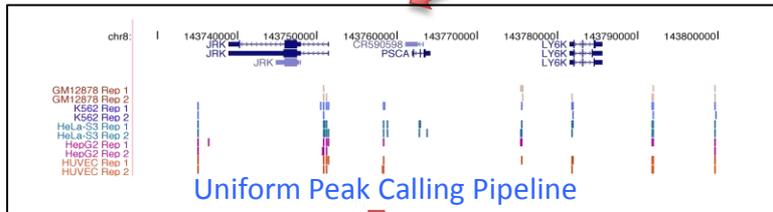
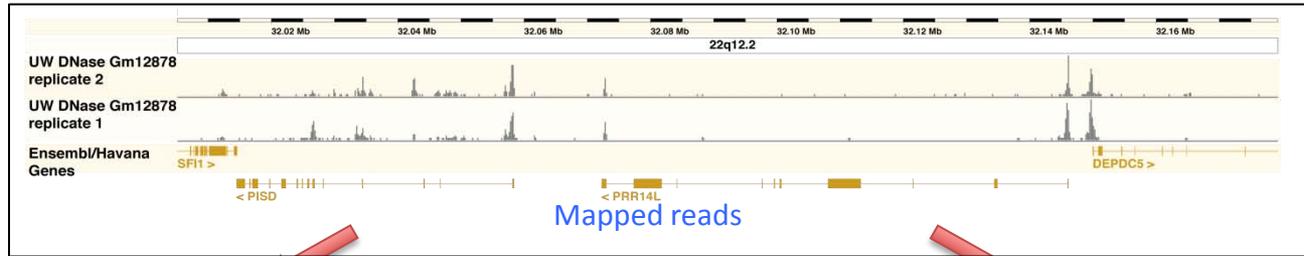
ChIP vs. Input DNA

Read quality

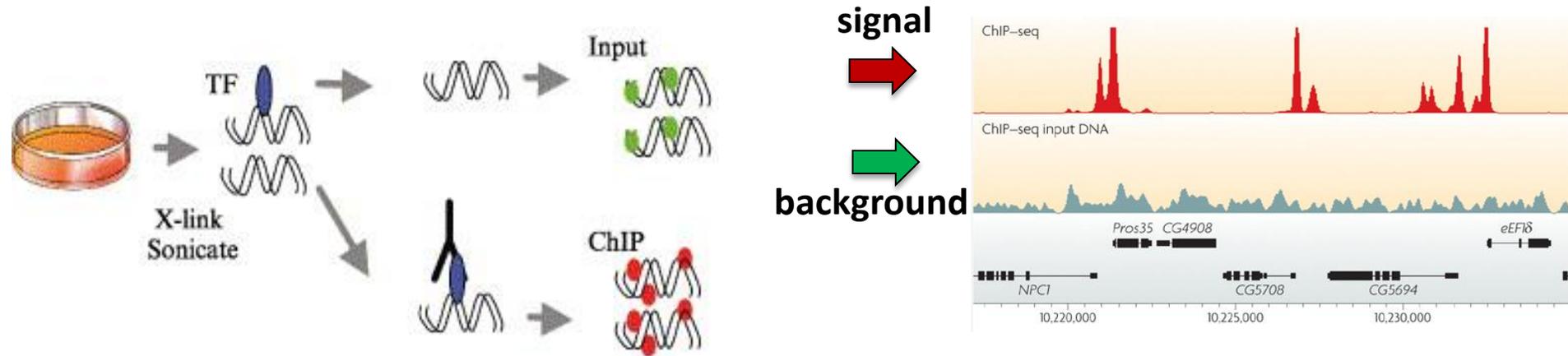
Mappability

Library complexity

ENCODE uniform processing pipeline



QC1: Use of input DNA as control dataset



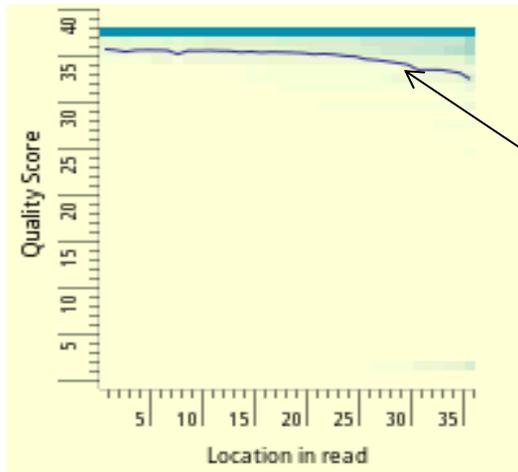
© sources unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- **Challenge:**
 - Even without antibody: Reads are **not** uniformly scattered
- **Sources of bias in input dataset scatter:**
 - Non-uniform fragmentation of the genome
 - Open chromatin fragmented more easily than closed regions
 - Repetitive sequences over-collapsed in the assembled genome.
- **How to control for these biases:**
 - Remove portion of DNA sample before ChIP step
 - Carry out control experiment without an antibody (input DNA)
 - Fragment input DNA, sequence reads, map, use as background

QC2: Read-level sequencing quality score $Q > 10$

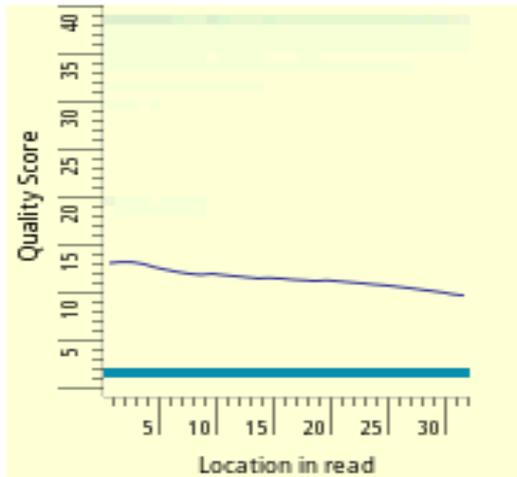
Read quality histograms

High quality reads



average base score
per position

Low quality reads

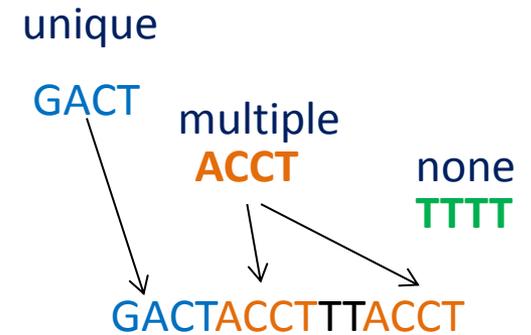


- Each column is a color-coded histogram
- Encodes fraction of all mapped reads that have base score Q (y-axis) at each position (x-axis)
- **Darker blue = higher density**
- Read quality tends to drop towards the ends of reads
- Low average per base score implies greater probability of mismappings.
- Typically, reject reads whose average score $Q < 10$

QC3: Fraction of short reads mapped >50%

Reads can map to:

- exactly one location (uniquely mapping)
- multiple locations (repetitive or multi-mapping)
- no locations (unmappable)



Dealing with multiply-mapping reads:

- Conservative approach: do not assign to any location
- Probabilistic approach: assign fractionally to all locations
- Sampling approach: pick one location at random, averages across many reads
- EM approach: map according to density, estimated from unambiguous reads
- Pair-end approach: use paired end read to resolve ambiguities in repeat reads

Absence of reads in a region could be due to:

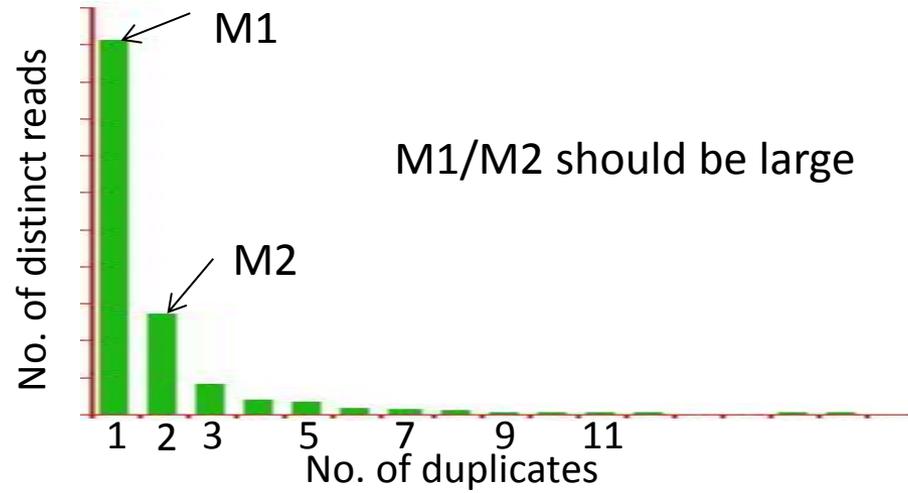
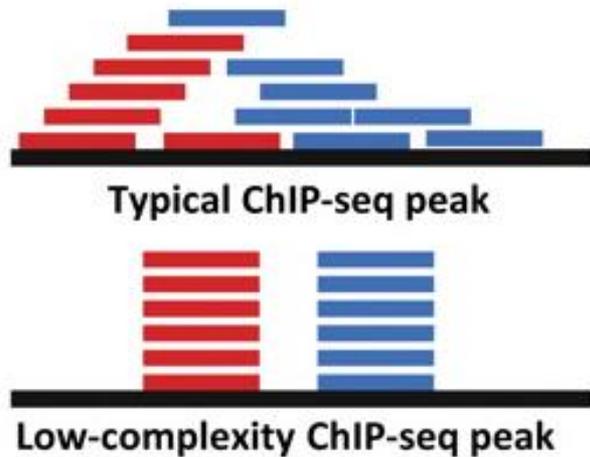
- No assembly coverage in that region (e.g. peri-centromeric region)
- Too many reads mapping to this location (e.g. repetitive element)
- No activity observed in this location (e.g. inactive / quiescent / dead regions)

Dealing with mappability biases:

- 'Black-listed' regions, promiscuous across many datasets
- 'White-listed' regions, for which at least some dataset has unique reads
- Treat unmappable regions as missing data, distinguish from 'empty' regions

QC4: Library complexity: non-redundant fraction

Library complexity



© sources unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

How many distinct uniquely mapping read? How many duplicates?

If your sample does not contain sufficient DNA and/or you over-sequence, you will simply be repeatedly sequencing PCR duplicates of a restricted pool of distinct DNA fragments. This is known a **low-complexity library** and is not desirable.

- **Histogram of no. of duplicates**
- Non-redundant fraction (NRF) =
$$\frac{\text{No. of 'distinct' unique-mapping reads}}{\text{No. of unique-mapping reads}}$$
- NRF should be > 0.8 when $10\text{M} < \# \text{reads} < 80\text{M}$ unique-mapping reads

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

Cross-correlation analysis

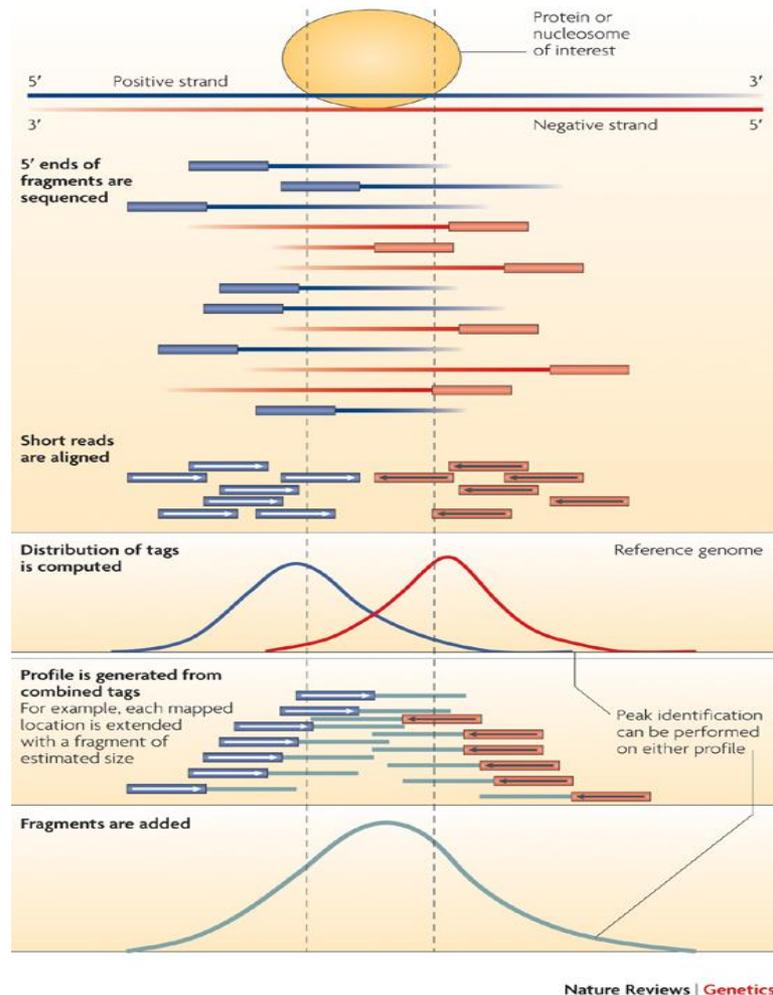
Exploiting forward and reverse reads

Fragment-length peak

Phantom read-length peak

ChIP-seq: exploiting forward and reverse reads

(Chromatin immunoprecipitation followed by sequencing)



Multiple IP fragments are obtained corresponding to each binding event

Ends of the fragments are sequenced i.e. "Short-reads/tags"

- Typically ~36 bp, 50 bp, 76 bp or 101 bp

Single-end (SE) sequencing

- Randomly sequence one of the ends of each fragment

Paired-end (PE) sequencing

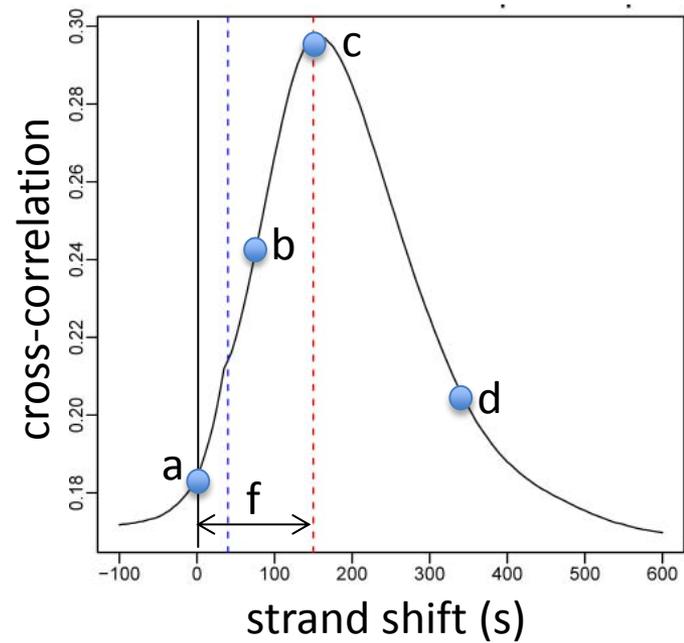
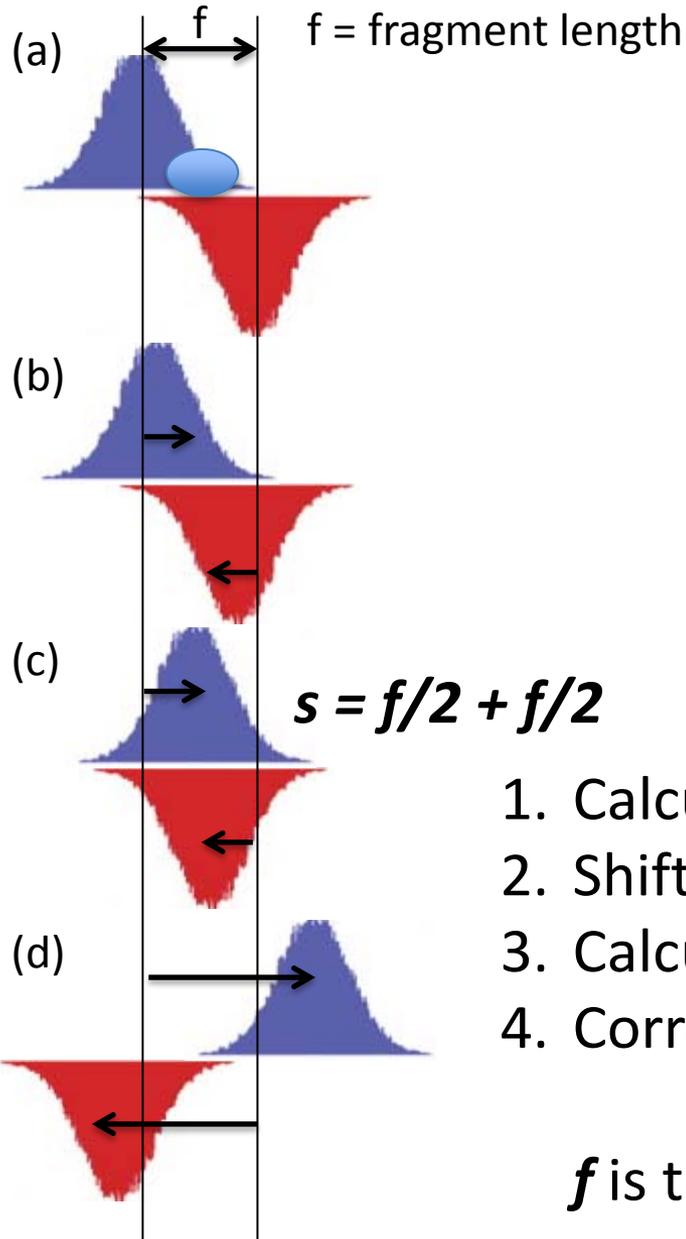
- sequence both ends of each fragment

Canonical "stranded mirror distribution of short-reads" after mapping reads to genome

- Heaps of reads on the **+ strand** and **- strand** separated by a distance \approx fragment length

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Park, Peter J. "ChIP-seq: advantages and challenges of a maturing technology." Nature Reviews Genetics 10, no. 10 (2009): 669-680.

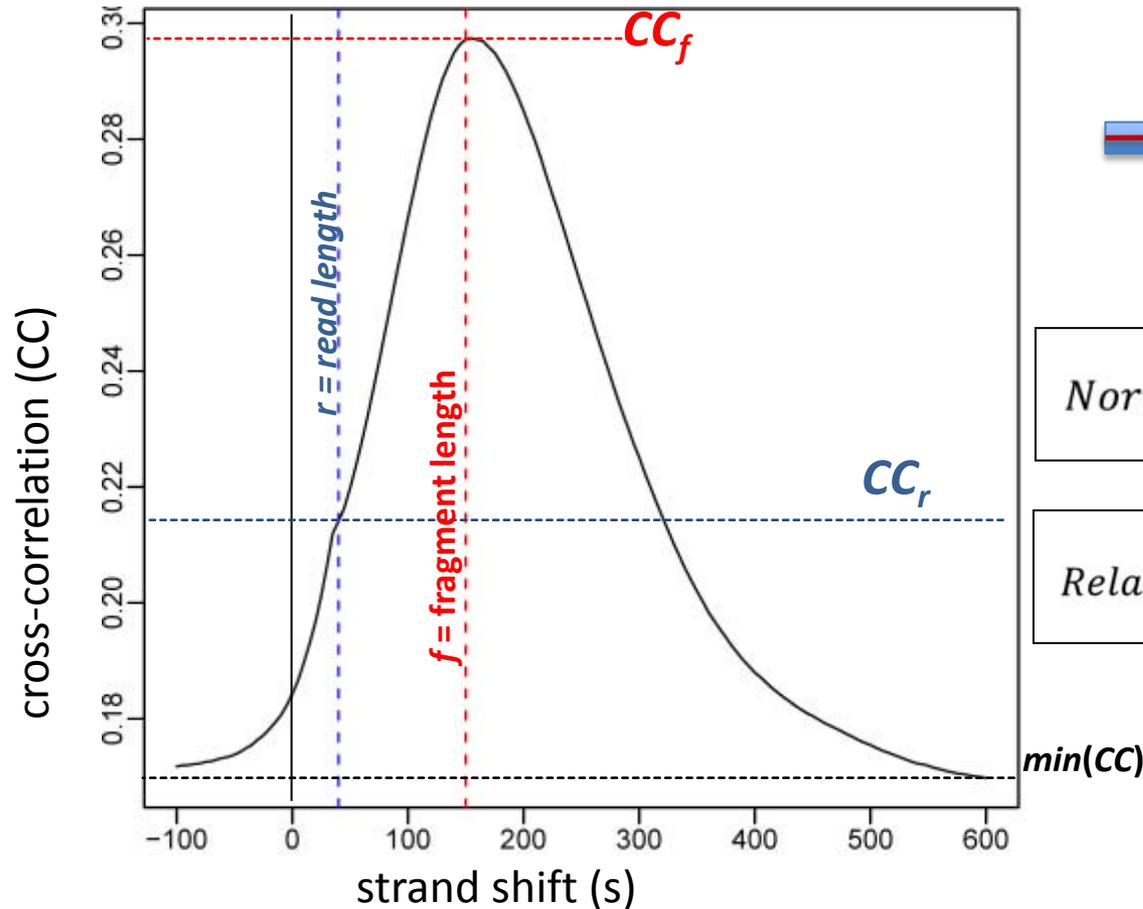
Strand cross-correlation (CC) analysis



1. Calculate forward and reverse strand signals
2. Shift both by specified offset towards each other
3. Calculate correlation of two signals at that shift
4. Correlation peaks at **fragment length** offset f

f is the length at which ChIP DNA is fragmented

Cross-correlation at *read* vs. *fragment* length

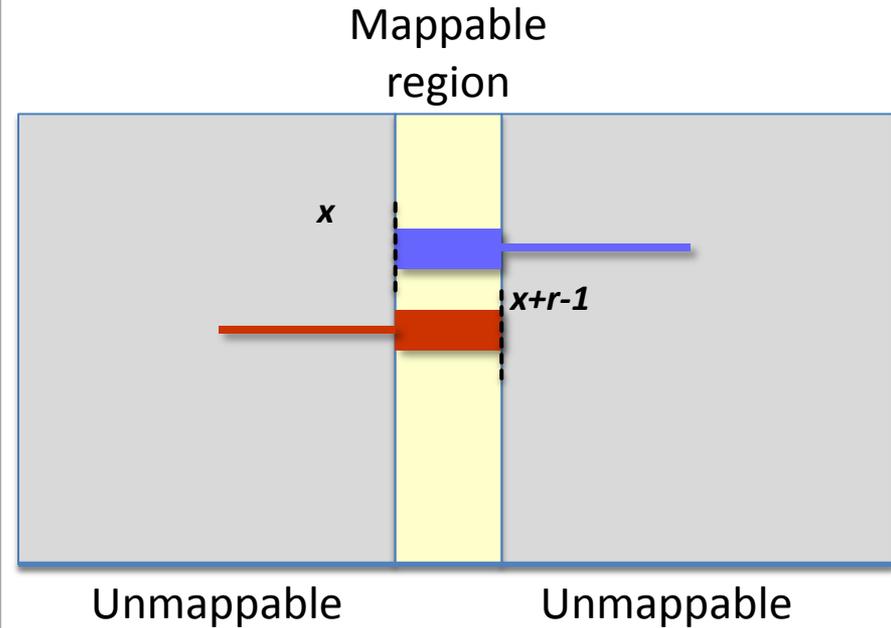
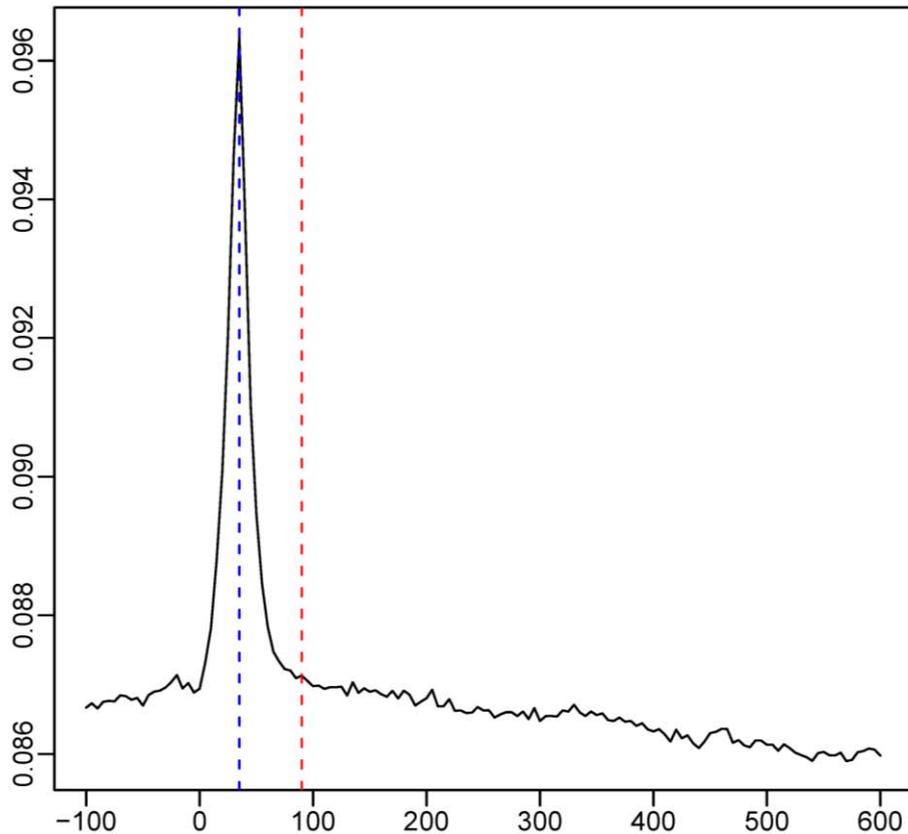


$$\text{Normalized strand CC (NSC)} = \frac{CC_f}{\min(CC)}$$

$$\text{Relative strand CC (RSC)} = \frac{(CC_f - \min(CC))}{(CC_r - \min(CC))}$$

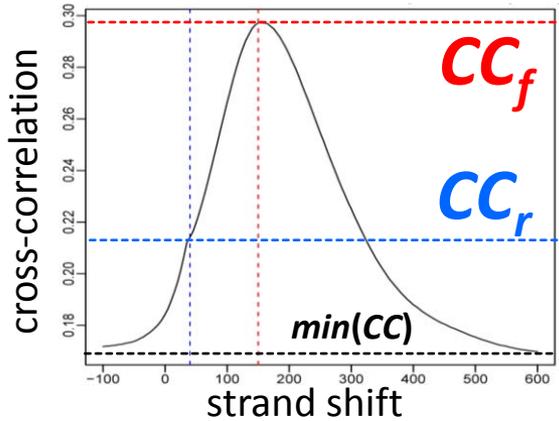
- Sign of a good dataset:
 - High absolute cross-correlation at *fragment* length (NSC)
 - High *fragment* length CC relative to *read* length CC (RSC)

Where does *read* cross-correlation come from?

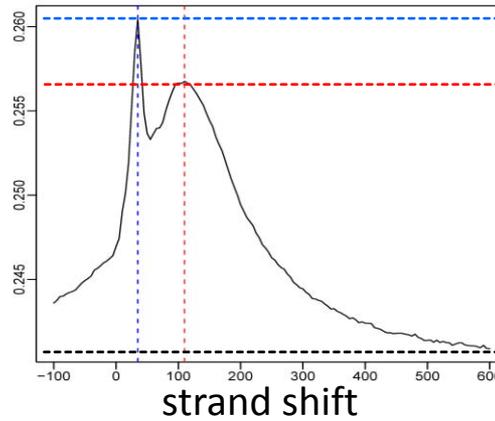


- Input dataset (no ChIP) shows 'phantom' peak at *read* length only
- Due to read mappability:
 - If position 'x' is uniquely mappable on + strand
 - Then position 'x+r-1' is uniquely mappable on - strand
- *Fragment*-length peak should always dominate the read-length peak

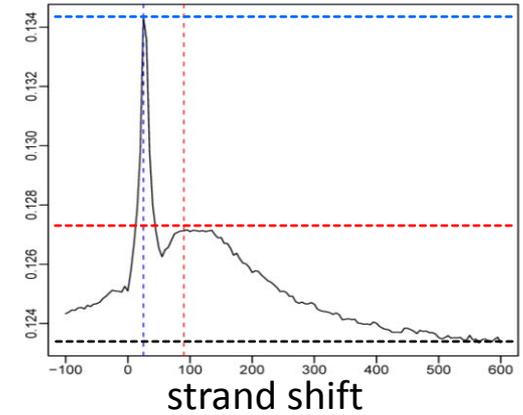
Example of good, medium, bad CC datasets



Highly quality



Medium quality



Low quality

$$\text{Normalized strand CC (NSC)} = \frac{CC_f}{\min(CC)}$$

$$\text{Relative strand CC (RSC)} = \frac{(CC_f - \min(CC))}{(CC_r - \min(CC))}$$

For highly enriched datasets, fragment length cross-correlation peak should be able to beat read-length phantom peak

RSC should be > 1

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

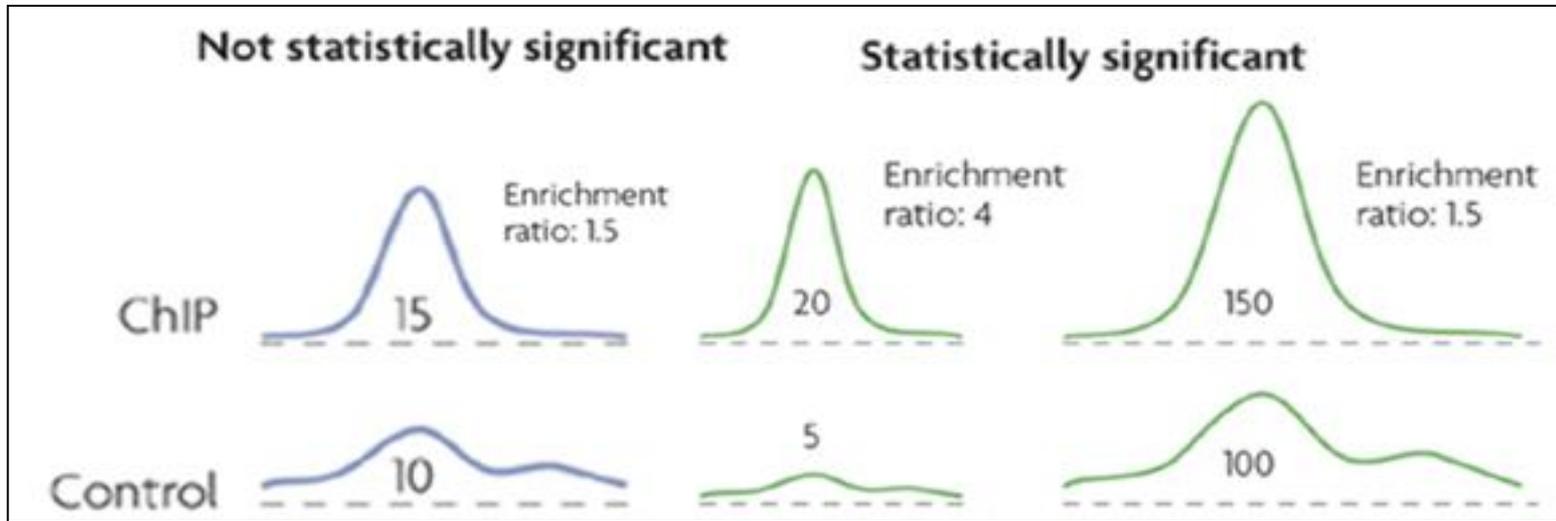
- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

Peak Calling

Continuous signal → Intervals

Peak calling: detect regions of enrichment



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Goal: Transform read counts into **normalized intensity signal**

Steps:

1. Estimate fragment-length f using strand cross-correlation analysis
2. Extend each read from 5' to 3' direction to fragment length f
3. Sum intensity for each base in 'extended reads' from both strands
4. Perform same operation on input-DNA control data (correct for sequencing depth differences)
5. Calculate enrichment ratio value for every position in the genome

Result: Enrichment fold difference for ChIP / control signal

Peak calling: identify discrete intervals

Program	Reference	Version	Graphical user interface?	Window-based scan	Tag clustering	Gaussian kernel density estimator	Strand-specific density	Peak height or fold enrichment (FE)	Background subtraction	Compensates for genomic duplications or deletions	False Discovery Rate	Compare to normalized control data (FE)	Compare to statistical model fitted with control data	Statistical model or test
CisGenome	28	1.1	X*	X			X	X		X		X		conditional binomial model
Minimal ChipSeq Peak Finder	16	2.0.1		X			X				X			
E-RANGE	27	3.1		X			X				X	X		chromosome scale Poisson dist.
MACS	13	1.3.5		X			X			X		X		local Poisson dist.
QuEST	14	2.3			X		X			X**		X		chromosome scale Poisson dist.
HPeak	29	1.1		X			X					X		Hidden Markov Model
Sole-Search	23	1	X	X			X		X			X		One sample t-test
PeakSeq	21	1.01		X			X					X		conditional binomial model
SISSRS	32	1.4		X		X					X			
spp package (wtd & mtc)	31	1.7		X		X		X	X'	X				
				Generating density profiles			Peak assignment		Adjustments w. control data		Significance relative to control data			

X* = Windows-only GUI or cross-platform command line interface

X** = optional if sufficient data is available to split control data

X' = method excludes putative duplicated regions, no treatment of deletions

Courtesy of the authors. License: CC BY.

Source: Wilbanks, Elizabeth G. and Marc T. Facciotti. "Evaluation of algorithm performance in ChIP-seq peak detection." PLOS ONE 5, no. 7 (2010): e11471.

Peak calling thresholds

Poisson p-value thresholds

- Read count model: Locally-adjusted' Poisson distribution

$$P(\text{count} = x) = \frac{\lambda_{local}^x \exp(-\lambda_{local})}{x!}$$

- $\lambda_{local} = \max(\lambda_{BG}, [\lambda_{1k},] \lambda_{5k}, \lambda_{10k})$ estimated from control data
 - Poisson p -value = $P(\text{count} \geq x)$
 - q -value : Multiple hypothesis correction

Peaks: Genomic locations that pass a user-defined p -value (e.g. $1e-5$) or q -value (e.g. 0.01) threshold

Empirical False discovery rates

- Swap ChIP and input-DNA tracks
 - Recompute p -values
- At each p -value, eFDR = Number of control peaks / Number of ChIP peaks
 - Use an FDR threshold to call peaks

Issues with peak calling thresholds

Cannot set a universal threshold for empirical FDRs and p-values

- Depends on ChIP and input sequencing depth
- Depends on binding ubiquity of factor
- Stronger antibodies get an advantage

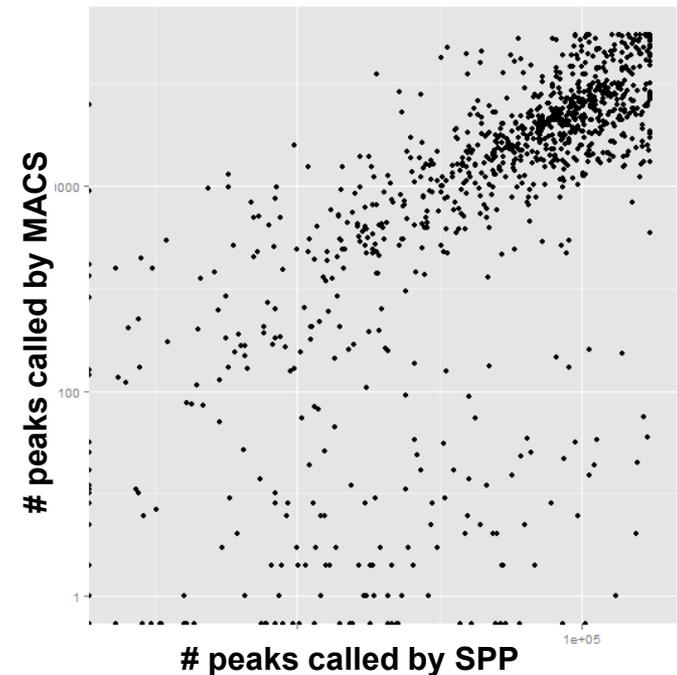
FDRs quite unstable

- Small changes in threshold => massive changes in peak numbers

Difficult to compare results across peak callers with a fixed threshold

- Different methods to compute eFDR or q-values

(at FDR = 1% cutoff)



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

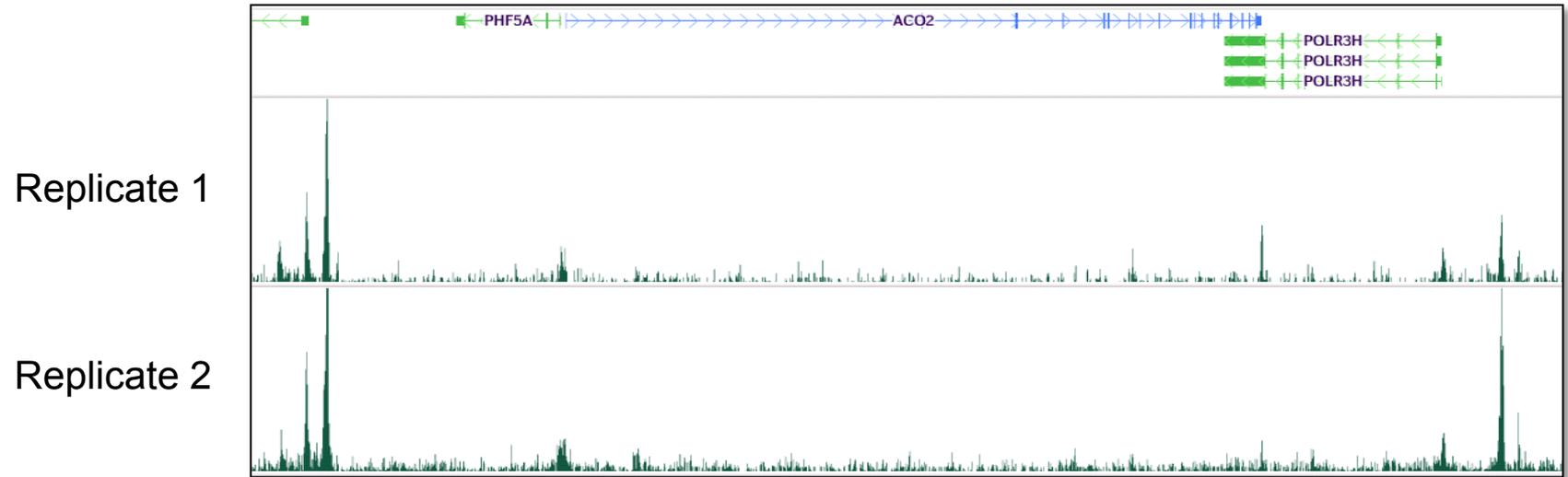
Selecting meaningful peaks using reproducibility

Use peak ranks in replicate experiments

IDR: Irreproducible Discovery Rate

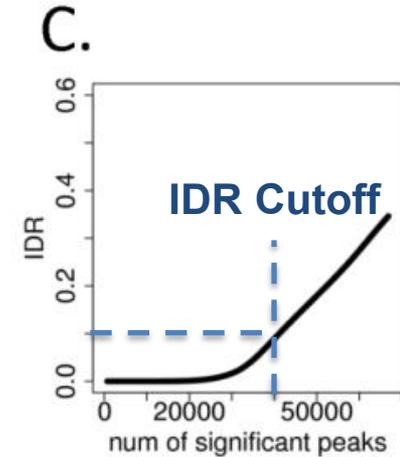
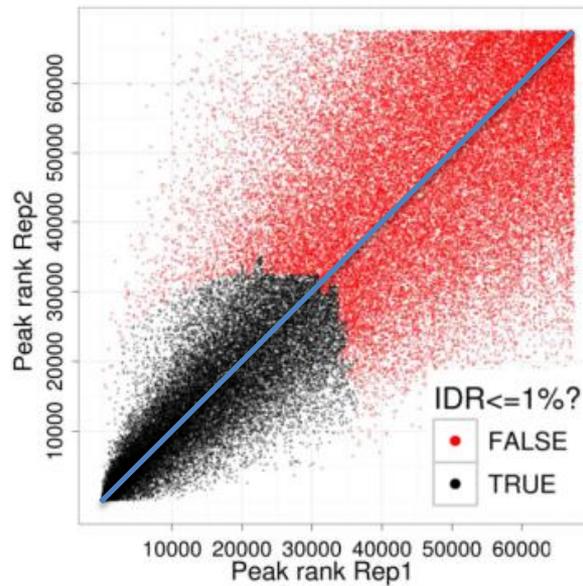
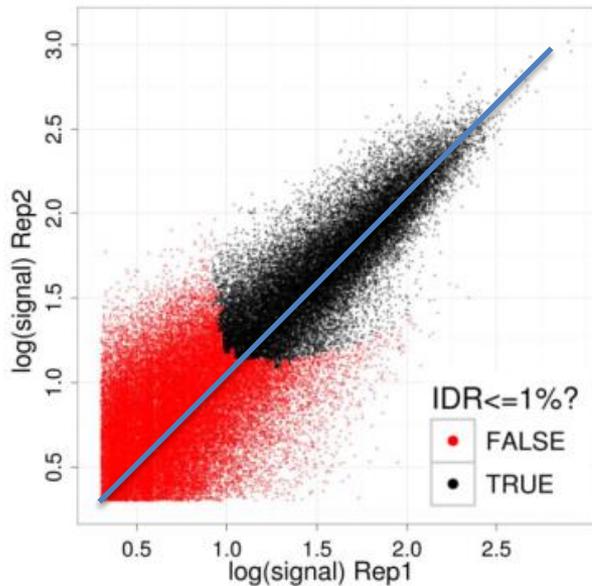
<http://anshul.kundaje.net/projects/idr>

How to combine two replicates



- Challenge:
 - Replicates show small differences in peak heights
 - Many peaks in common, but many are unique
- Problem with simple solutions:
 - Union: too lenient, keeps garbage from both
 - Intersection: too stringent, throws away good peaks
 - Sum: does not exploit independence of two datasets

IDR idea: Exploit peak rank similarity in replicates



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- Key idea: True peaks will be highly ranked in both replicates
 - Keep going down rank list, until ranks are no longer correlated
 - This cutoff could be different for the two replicates
 - The actual peaks included may differ between replicates
 - Adaptively learn optimal peak calling threshold
 - FDR threshold of 10% → 10% of peaks are false (widely used)
 - IDR threshold of 10% → 10% of peaks are not reproducible

The IDR model: A two component mixture model

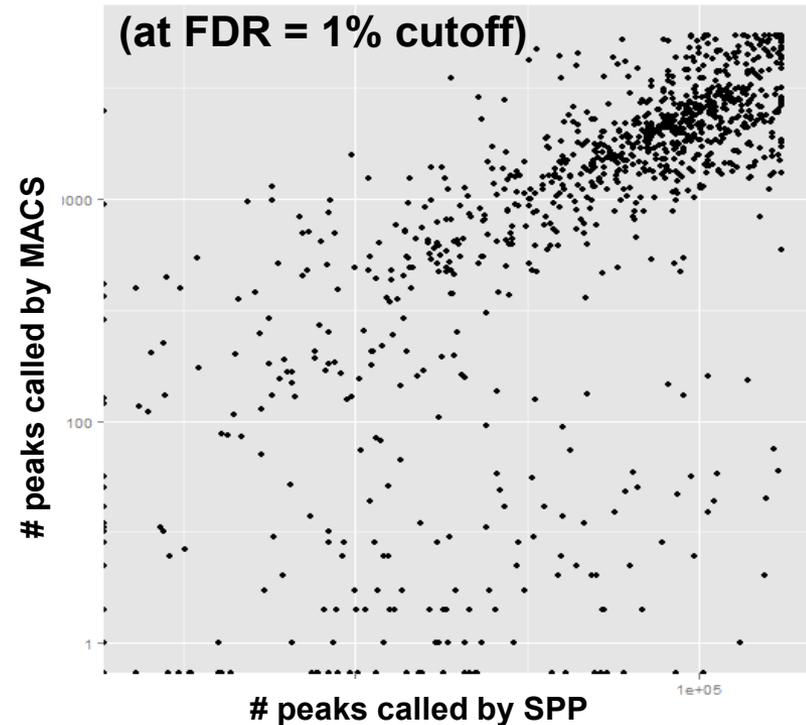
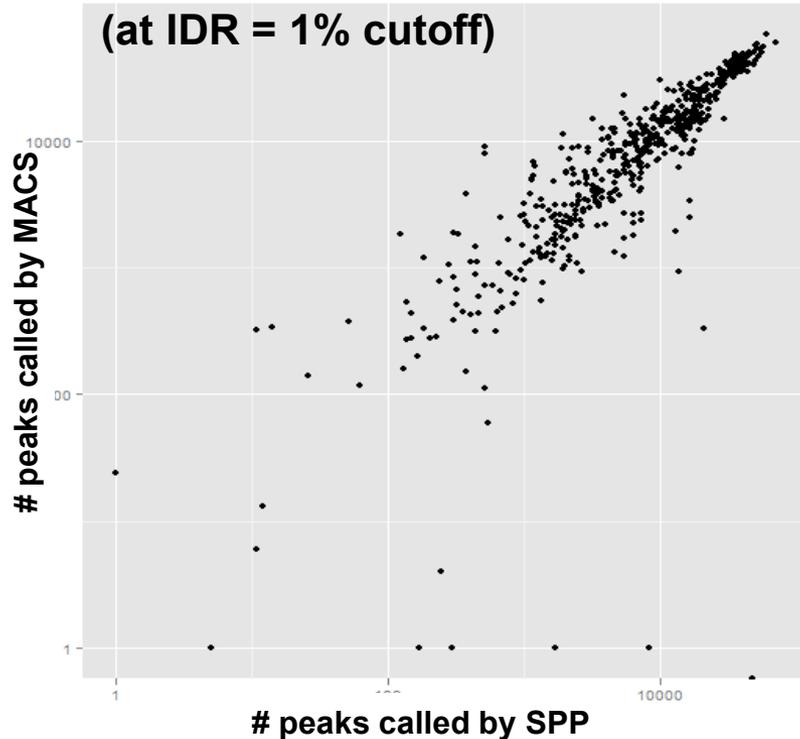
- Looking only at ranks means that the marginals are uniform, so all the information is encoded in the joint distribution.
- Model the joint distribution of ranks as though it came from a two component Gaussian mixture model:

$$(x, y) \sim pN(\mu, \mu, \sigma, \sigma, \rho) + (1 - p)N(0, 0, 1, 1, 0)$$

- This can be fit via an EM-like algorithm.

IDR leads to higher consistence between peak callers

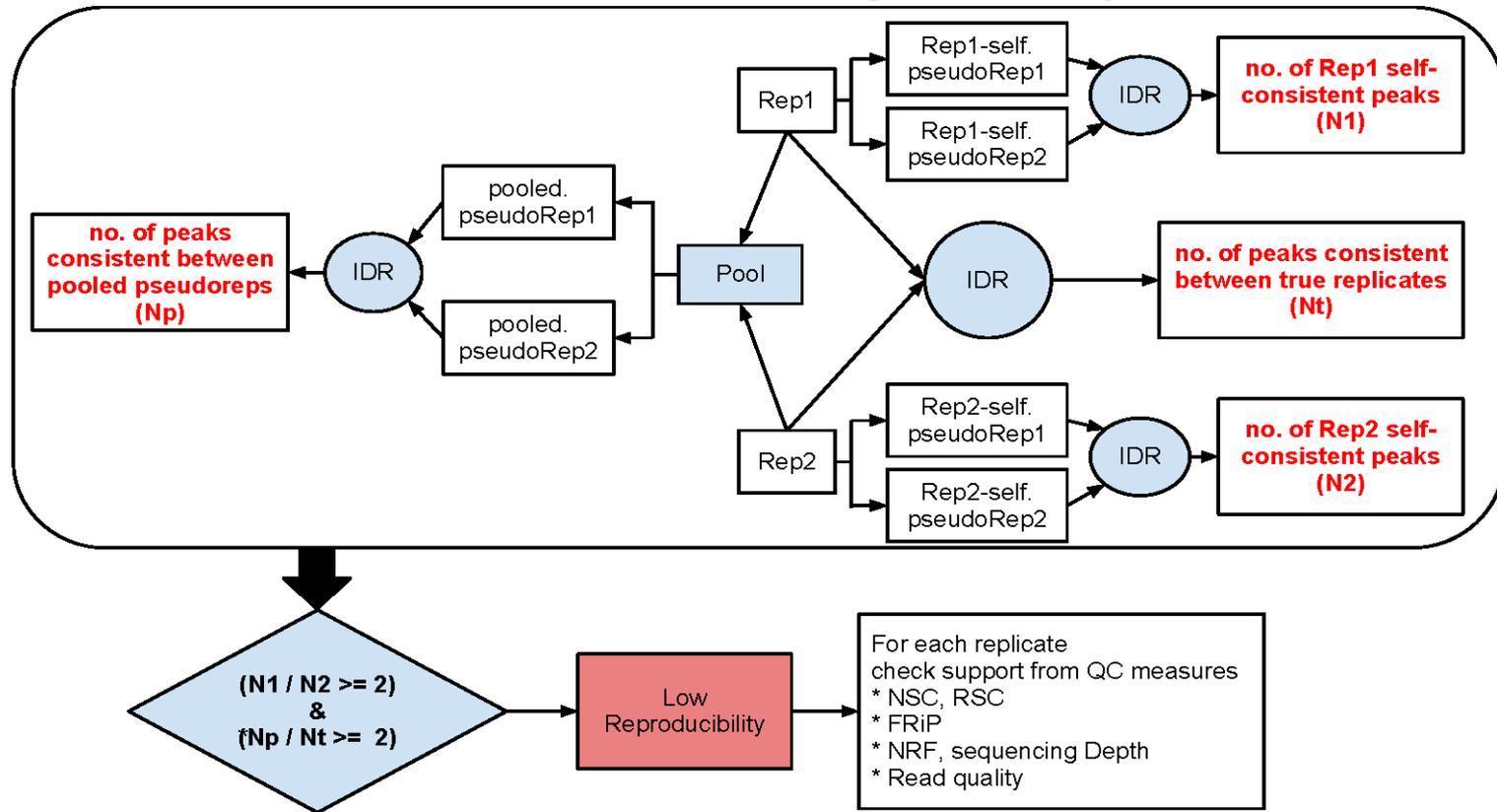
IDR = Irreproducible Discovery Rate FDR = False Discovery Rate



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- Compare number of peaks found by two different peak callers
- IDR thresholds are far more robust and comparable than FDR
- FDR only relies on enrichment over input, IDR exploits replicates

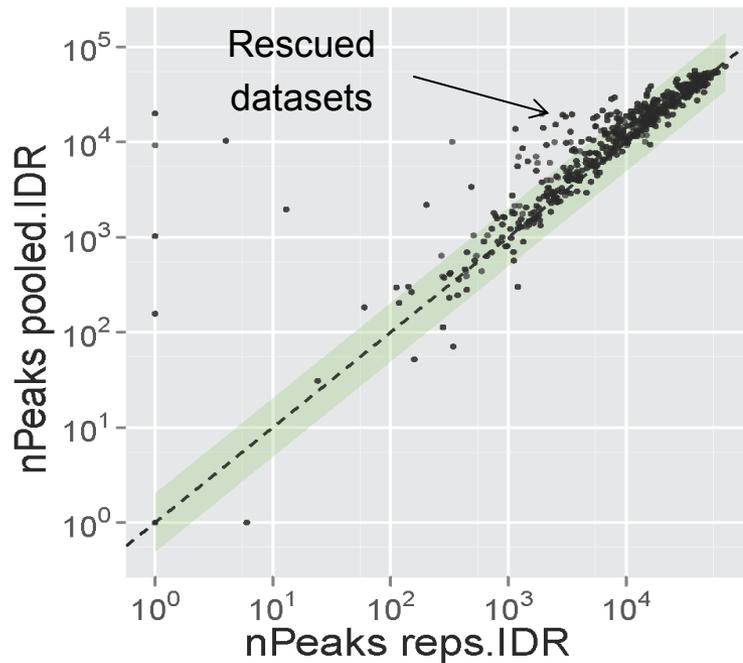
What if we don't have good replicates?



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- IDR pipeline uses replicates when they are available
 - IDR pipeline also evaluates each replicate individually
 - Pooling strategy to generate pseudo-replicates
- ➔ Can pin-point 'bad' replicates that may lead to low reproducibility
- ➔ Can estimate IDR thresholds when replicates are not available

Only one good replicate: Pseudo-replicates



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- IDR pipeline can be used to rescue datasets with only one good replicate (using pseudo-replicates)
- IDR pipeline can also be used to call optimal thresholds on a dataset with a single replicate (e.g. when there isn't enough material to perform multiple reps)

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

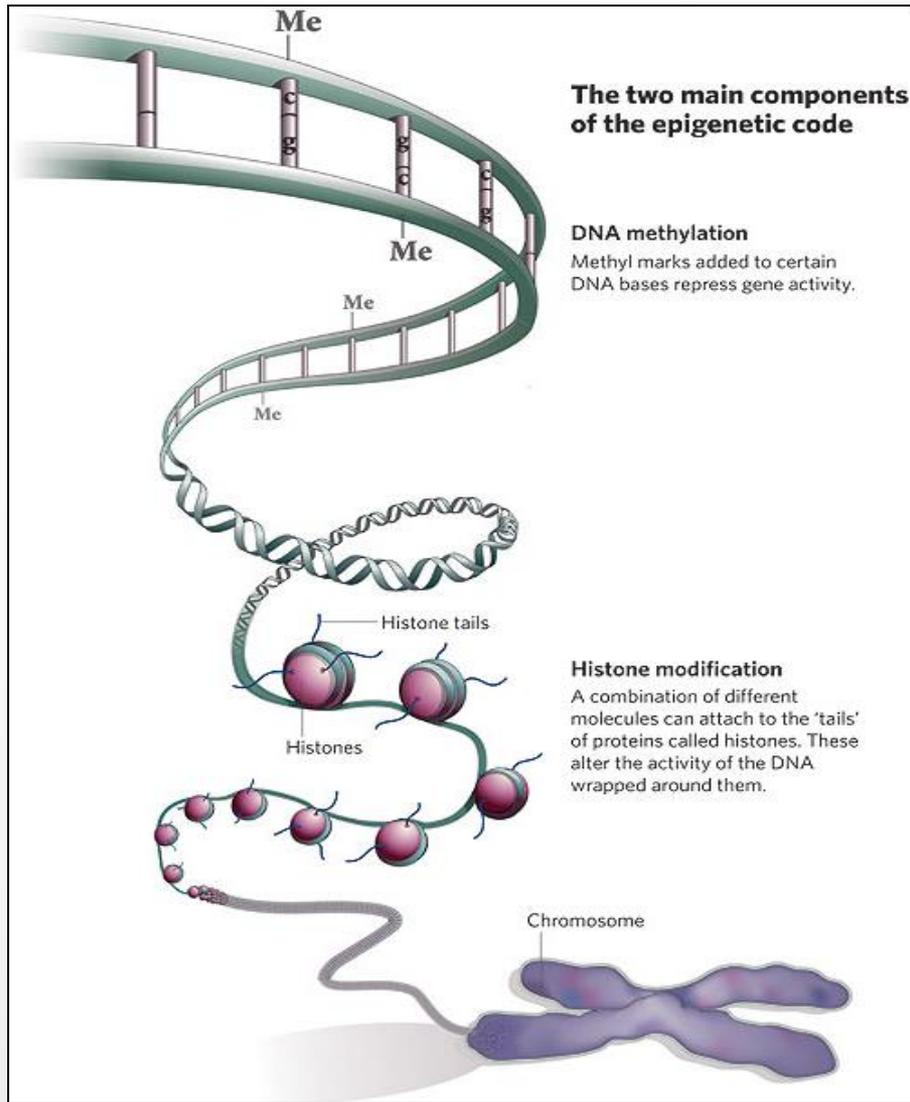
- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

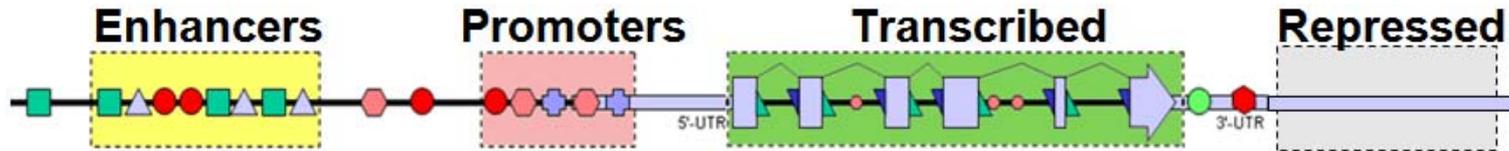
Chromatin signatures for genome annotation



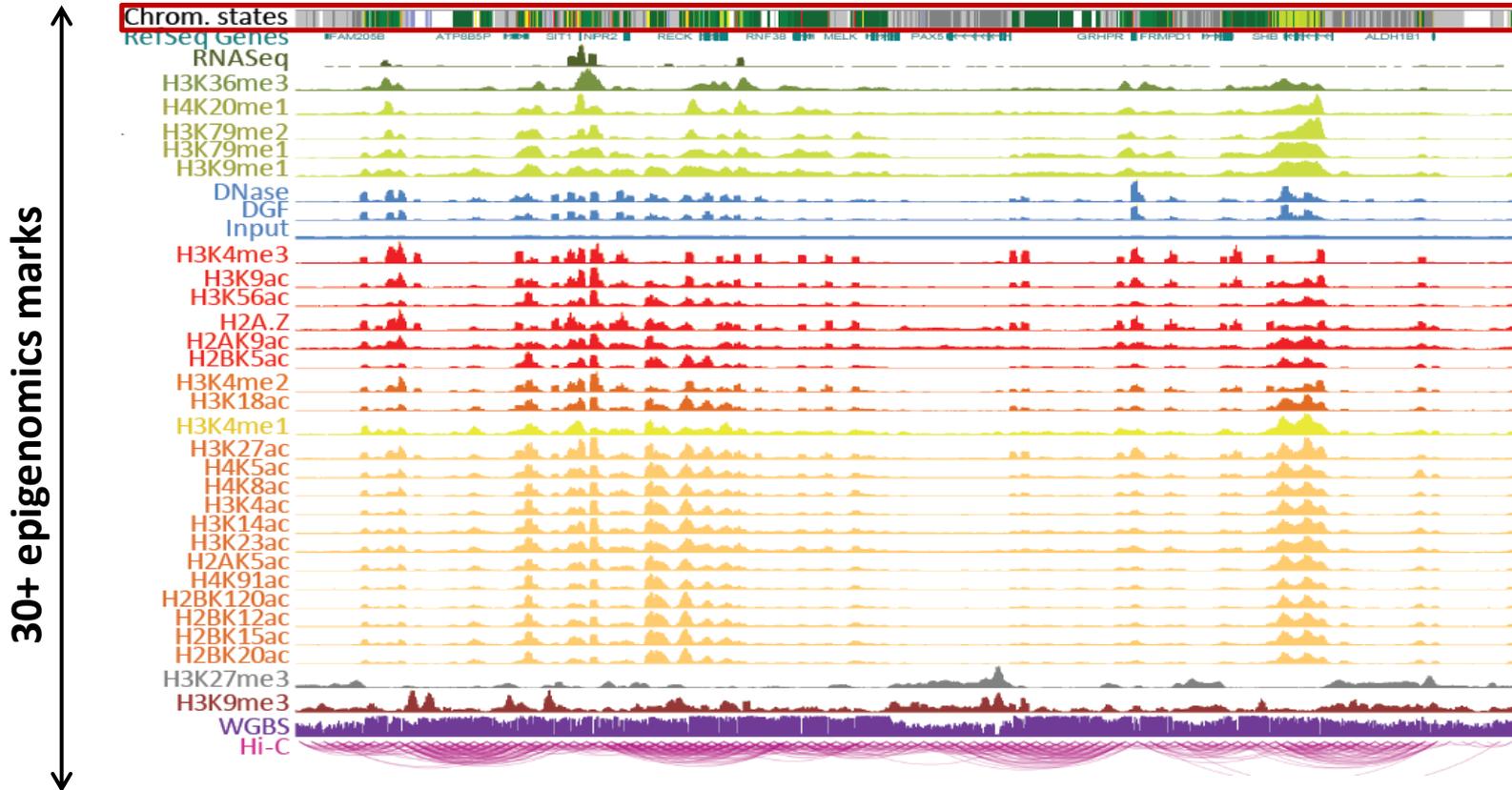
- **Challenges**
 - Dozens of marks
 - Complex combinatorics
 - Diversity and dynamics
- **Histone code hypothesis**
 - Distinct function for distinct combinations of marks?
 - Both additive and combinatorial effects
- **How do we find biologically relevant ones?**
 - Unsupervised approach
 - Probabilistic model
 - Explicit combinatorics

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Qiu, Jane. "Epigenetics: Unfinished Symphony." *Nature* 441, no. 7090 (2006): 143-145.

Summarize multiple marks into chromatin states



Chromatin state track summary

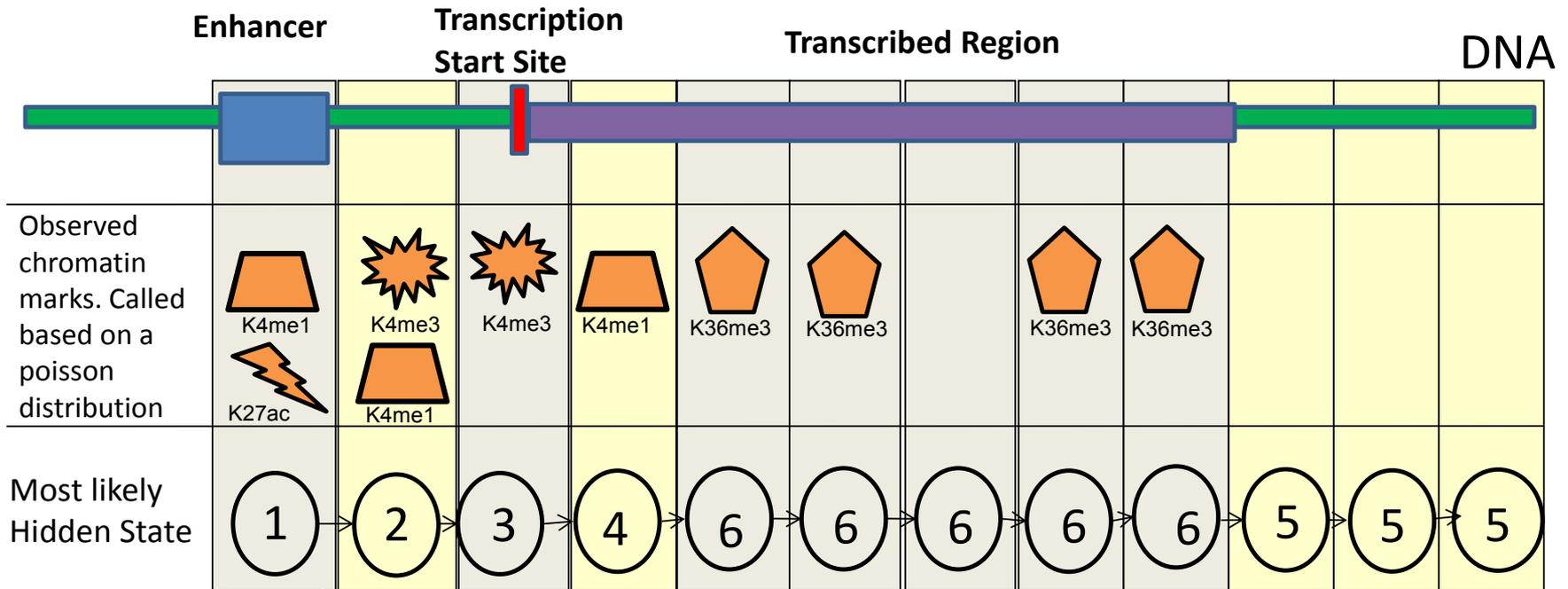


© WashU Epigenome Browser. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

WashU Epigenome Browser

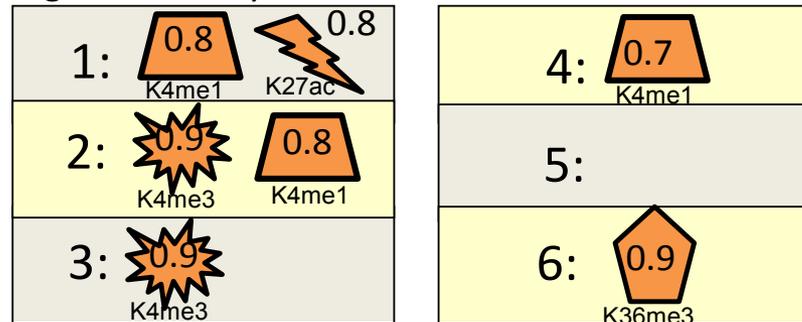
ChromHMM: multi-variate hidden Markov model

Multivariate HMM for Chromatin States



200bp intervals

High Probability Chromatin Marks in State



All probabilities are learned from the data

state	H3K14ac	H3K23ac	H4K12ac	H2AK9ac	H4K16ac	H2AK5ac	H4K91ac	H3K4ac	H2BK20ac	H3K18ac	H2BK120ac	H3K27ac	H2BK5ac	H2BK12ac	H3K36ac	H4K5ac	H4K8ac	H3K9ac	PoII	CTCF	H3AZ	H3K4me3	H3K4me2	H3K4me1	H3K9me1	H3K79me3	H3K79me2	H3K79me1	H3K27me1	H2BK5me1	H4K20me1	H3K36me3	H3K36me1	H3R2me1	H3R2me2	H3K27me2	H3K27me3	H4R3me2	H3K9me2	H3K9me3	H4K20me3
1	3.8	23.6	24.2	18.0	37.7	25.5	95.2	94.8	94.3	99.2	99.6	99.7	98.9	79.1	88.6	93.6	96.9	83.6	51.6	15.7	87.5	94.2	93.8	64.2	87.0	3.8	3.3	12.0	19.4	11.6	3.8	0.5	2.6	1.9	2.1	0.2	0.1	0.2	0.5	0.1	1.8

Ernst and Kellis
Nature Biotech 2010

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Ernst, Jason and Manolis Kellis. "Discovery and characterization of chromatin states for systematic annotation of the human genome." Nature Biotechnology 28, no. 8 (2010): 817-825.

Design Choice

- How to model the emission distribution
 - Model the signal directly
 - Locally binarize the data
- For M input marks each state k has a vector of (p_{k1}, \dots, p_{kM}) of parameters for independent Bernoulli random variables which determine the emission probability for an observed combination of marks

Data Binarization

- Leads to biologically interpretable models that can be robustly learned
- Let c_{ij} be the number of reads for mark i . mapping to bin j . λ_i be the average number of reads mapping to a bin for modification i . The input for feature i becomes '1' if

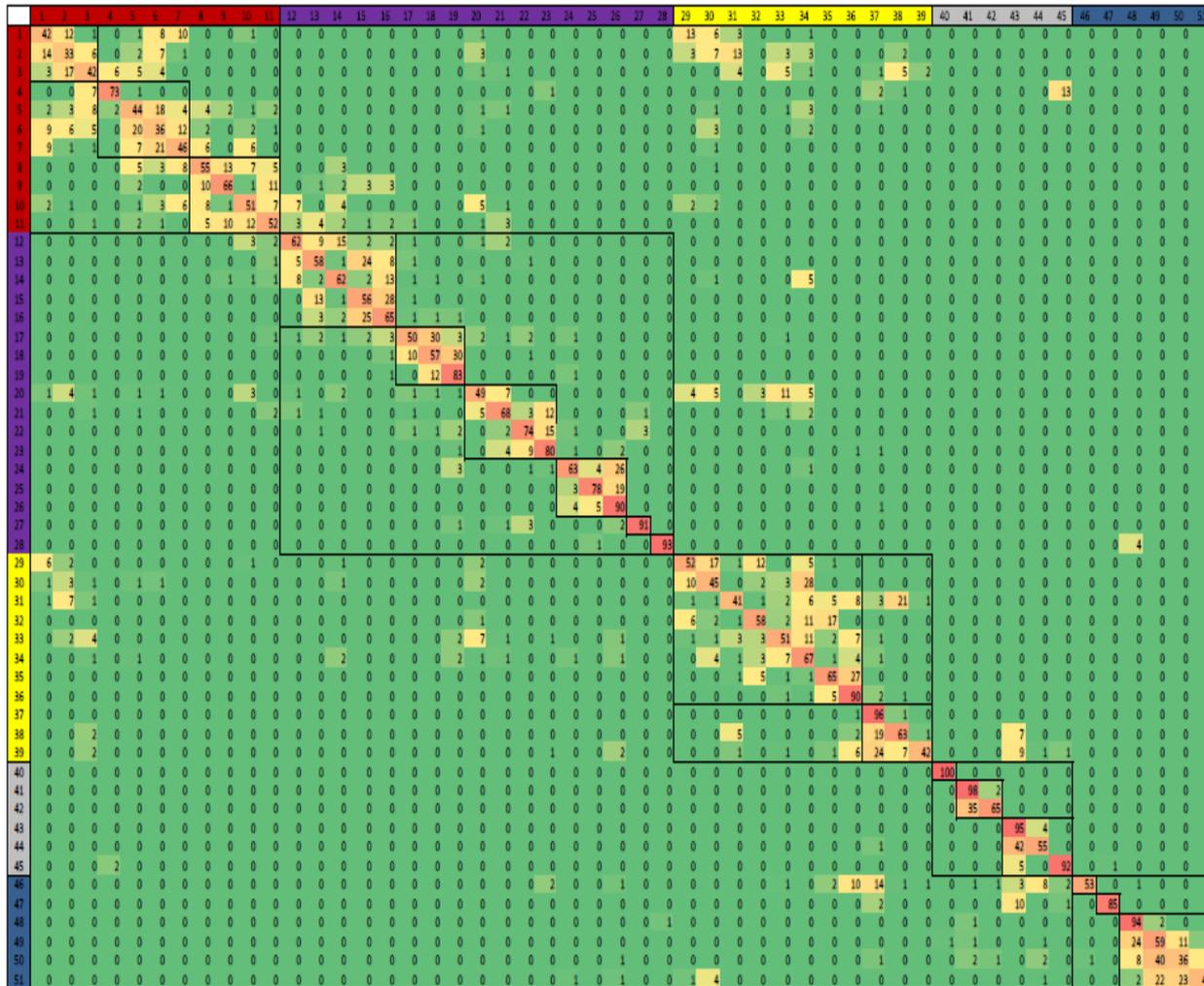
$$P(X > c_{ij}) < 10^{-4}$$

where X is a Poisson random variable with mean λ_i

Emission Parameter Matrix $e_k(\vec{x}_i)$

state	H3K1ac	H3K23ac	H3K27ac	H3K36ac	H4K1ac	H4K20ac	H4K31ac	H4K36ac	H4K91ac	H4K93ac	H2BK20ac	H3K18ac	H2BK120ac	H3K27ac	H2BK3ac	H2BK12ac	H3K36ac	H4K53ac	H4K83ac	H3K93ac	PoII	CTCF	H2AZ	H3K4me3	H3K4me2	H3K4me1	H3K9me1	H3K9me3	H3K79me2	H3K79me1	H3K27me1	H2BK5me1	H4K20me1	H3K36me3	H3K36me1	H3K27me1	H3K27me2	H3K27me3	H4K3me2	H3K9me2	H3K9me3	H4K20me3
1	3.8	23.6	24.2	18.0	37.7	25.5	95.2	94.8	94.3	99.2	99.6	99.7	98.9	79.1	88.6	93.6	86.9	83.6	51.6	15.7	87.5	94.2	93.8	64.2	87.0	3.8	3.3	12.0	19.4	11.6	3.8	0.5	2.6	1.9	2.1	0.2	0.1	0.2	0.5	0.1	1.8	
2	2.5	17.5	9.2	3.2	5.9	6.3	44.6	44.4	47.0	73.2	74.1	85.9	71.2	22.1	33.5	61.9	63.3	35.4	18.1	10.9	91.2	86.7	90.4	65.9	78.3	2.4	2.2	7.9	17.6	8.7	2.3	0.6	2.2	1.7	1.5	0.4	0.5	0.2	0.4	0.1	1.4	
3	0.5	5.8	1.8	1.0	0.9	1.2	12.3	9.5	8.8	22.6	21.3	22.8	12.1	2.2	4.2	8.4	12.8	7.1	11.2	16.3	77.1	93.9	80.3	45.6	74.2	1.5	1.3	4.6	4.3	7.0	8.8	0.2	1.4	2.1	1.5	0.2	2.1	0.1	0.1	0.1	1.2	
4	0.1	0.8	0.1	0.4	0.3	0.2	1.5	0.9	0.7	2.1	0.5	0.2	0.1	0.1	0.1	0.3	0.3	1.9	6.2	19.0	77.9	20.8	21.1	26.4	0.3	0.0	0.1	0.0	1.3	14.9	0.1	0.2	1.4	0.9	0.0	10.2	0.1	0.1	0.4	0.1	1.3	
5	0.0	0.2	0.8	1.3	2.0	0.4	26.6	12.6	6.7	15.8	23.1	26.8	24.3	1.7	5.3	20.6	21.8	87.0	11.2	2.7	15.7	6.3	3.8	2.5	0.0	5.5	14.2	0.0	0.3	0.2	0.6	0.0	5.0	0.2	0.2	0.6	0.0	0.0	0.1	0.1	0.0	1.5
6	0.1	1.8	3.6	6.9	6.1	1.9	74.5	63.5	53.0	75.7	84.3	89.4	86.7	20.8	41.5	20.5	21.6	62.7	69.2	25.5	61.2	98.3	37.4	7.1	40.3	5.3	2.7	5.6	0.6	6.0	11.6	0.0	0.5	0.4	0.9	0.0	0.0	0.0	0.0	0.0	0.1	1.6
7	1.7	8.7	20.5	43.0	53.7	9.8	98.7	98.6	95.7	99.4	98.9	99.9	99.9	76.5	93.3	81.8	76.6	99.2	88.0	26.9	77.1	99.7	38.3	2.8	37.9	32.1	24.9	14.0	2.6	6.5	16.2	0.1	1.0	0.5	1.2	0.0	0.0	0.0	0.1	0.0	1.9	
8	1.1	12.7	5.2	11.9	5.6	6.8	56.9	56.1	37.5	52.4	69.8	89.1	85.8	21.3	24.8	16.7	10.3	69.8	52.1	12.0	31.4	96.7	51.7	14.3	45.3	86.3	80.1	23.8	1.7	6.7	42.2	4.5	0.4	1.1	0.9	0.0	0.0	0.0	0.0	0.0	0.5	
9	0.5	7.2	3.0	1.1	0.5	2.1	4.0	7.4	2.4	2.5	11.6	35.3	28.0	2.7	2.8	2.0	1.8	8.6	34.7	4.2	4.5	79.2	41.5	23.8	36.1	86.0	82.6	12.0	1.9	6.7	43.5	7.4	0.2	0.6	0.7	0.0	0.0	0.0	0.0	0.0	0.2	0.4
10	4.1	24.8	13.1	17.5	24.4	37.0	90.4	88.6	82.0	89.8	95.5	97.0	95.1	54.0	56.4	67.2	45.7	55.7	46.6	10.2	40.8	84.6	92.3	91.4	92.8	67.1	67.2	63.4	29.2	53.8	65.2	4.5	6.8	5.7	3.4	0.1	0.0	0.3	0.1	0.0	1.0	
11	1.6	21.0	3.9	3.8	3.2	8.4	28.0	26.1	14.5	22.6	37.6	56.8	47.4	6.0	6.4	13.5	8.8	18.4	30.9	6.9	20.1	92.4	93.5	94.3	94.5	73.8	74.2	55.1	24.0	57.2	79.9	8.9	6.6	4.4	2.5	0.1	0.0	0.2	0.1	0.1	0.7	
12	3.6	17.0	8.9	2.2	14.1	34.9	60.3	51.0	38.8	35.6	56.6	53.3	55.3	11.9	11.5	30.0	15.4	17.0	28.7	0.7	1.4	5.8	26.0	77.8	87.4	87.0	76.3	41.6	79.8	82.2	13.4	3.1	6.4	3.6	0.3	0.0	0.7	0.2	0.1	0.4		
13	1.2	10.8	3.6	0.7	2.5	7.4	9.1	6.5	2.6	2.5	6.5	7.7	5.5	0.5	1.0	5.5	3.3	0.3	10.2	1.5	2.0	5.6	83.7	82.9	92.7	92.6	64.4	38.2	80.0	89.8	12.0	2.4	3.3	2.1	0.3	0.4	0.2	0.1	0.3	0.1		
14	0.7	5.3	7.9	1.0	2.4	18.0	19.8	20.7	14.6	7.9	24.5	20.1	21.8	6.7	6.6	11.3	8.6	0.3	6.9	1.7	2.1	1.3	8.0	33.0	16.3	61.8	62.1	37.9	9.7	14.3	18.3	9.7	0.2	1.7	1.2	0.2	0.0	0.1	0.3	0.4	1.0	
15	0.2	1.9	2.9	0.3	0.3	1.5	0.8	1.3	0.3	0.2	1.2	1.5	1.2	0.2	0.4	0.7	1.2	0.1	5.0	0.7	0.0	0.2	11.0	17.3	29.3	84.7	82.2	33.1	8.0	26.4	56.2	5.2	0.2	0.7	0.9	0.1	0.0	0.1	0.1	0.6	0.2	
16	0.0	0.3	0.4	0.1	0.1	0.5	0.4	0.6	0.2	0.1	0.5	0.4	0.3	0.1	0.2	0.2	0.3	0.0	0.6	0.3	0.1	1.2	2.8	3.9	29.0	25.2	8.5	0.7	1.3	7.9	1.2	0.0	0.1	0.1	0.0	0.0	0.0	0.0	0.4	0.1		
17	1.8	9.8	2.8	0.9	2.4	7.8	6.8	6.1	2.3	4.0	3.5	8.4	3.5	0.3	1.8	9.3	6.6	0.5	3.5	1.1	0.4	1.0	52.3	68.8	83.7	22.7	23.5	61.2	48.3	64.7	57.3	21.8	2.8	4.9	2.0	1.0	0.1	0.5	0.5	0.1	0.5	
18	0.3	2.6	2.1	0.4	0.5	2.4	1.4	1.6	0.5	0.6	0.9	1.6	0.7	0.2	0.5	1.9	1.9	0.1	1.6	0.7	0.1	10.4	9.7	29.1	11.3	34.0	14.9	19.9	21.4	10.6	0.5	1.2	0.9	0.5	0.1	0.1	0.3	0.2	0.2			
19	0.1	0.3	0.5	0.2	0.1	0.5	0.1	0.3	0.0	0.0	0.1	0.2	0.1	0.1	0.1	0.1	0.3	0.0	0.4	0.3	0.0	0.5	0.2	2.0	4.4	2.9	7.3	1.0	1.0	2.3	1.4	0.0	0.2	0.3	0.1	0.0	0.0	0.1	0.4	0.1		
20	2.5	10.7	5.4	3.1	9.9	26.2	58.2	48.8	41.7	49.3	54.8	57.1	51.5	13.0	14.1	31.6	21.7	4.0	14.5	6.7	15.6	20.9	56.8	97.1	70.5	5.4	5.6	33.8	31.1	52.6	38.4	7.4	3.4	6.9	3.8	0.5	0.1	0.7	0.3	0.0	1.0	
21	0.2	0.8	2.0	1.3	7.2	11.3	32.3	15.4	11.5	5.7	18.6	8.4	12.4	2.1	1.2	3.2	2.3	0.5	1.5	6.5	0.6	4.7	17.0	68.7	33.3	8.7	6.4	37.2	9.6	65.4	87.7	9.2	1.3	7.1	5.3	0.1	0.2	1.4	0.0	0.0	0.8	
22	0.1	0.1	1.1	0.6	6.2	2.4	7.8	1.8	0.6	0.1	1.5	0.8	1.5	0.1	0.1	0.7	0.7	0.1	8.5	1.0	0.0	5.3	8.4	14.6	15.5	9.1	50.0	9.6	77.5	94.1	22.9	0.5	5.6	4.6	0.1	0.0	1.4	0.0	0.1	0.7		
23	0.0	0.1	0.1	0.4	2.0	1.6	4.9	1.2	0.5	0.2	0.9	0.2	0.1	0.0	0.1	0.1	0.1	0.0	1.4	1.1	0.0	0.5	2.6	1.4	0.5	0.1	5.4	1.3	19.4	36.8	2.5	0.1	1.4	1.5	0.1	0.2	0.3	0.0	0.0	0.2		
24	0.3	1.8	2.1	0.9	3.2	3.8	4.0	2.3	0.9	0.6	1.1	3.6	1.9	0.1	0.4	3.7	3.9	0.3	2.2	1.0	0.1	6.0	4.5	17.2	1.3	0.2	15.6	29.8	29.3	7.1	49.5	1.3	4.7	2.2	1.0	0.0	0.6	0.3	0.1	0.3		
25	0.1	0.3	0.8	0.5	0.5	0.6	0.3	0.3	0.1	0.0	0.1	0.5	0.3	0.1	0.1	0.4	0.7	0.1	0.8	0.4	0.0	0.1	0.4	0.1	1.4	0.8	0.1	2.0	8.8	2.8	0.4	6.0	1.4	1.1	0.9	0.6	0.1	0.2	0.3	1.7		
26	0.1	0.2	0.6	0.2	0.2	0.2	0.1	0.2	0.0	0.1	0.3	0.2	0.0	0.1	0.2	0.4	0.0	0.6	0.3	0.0	0.0	0.3	0.1	0.8	0.2	0.0	0.8	2.3	0.9	0.2	4.2	0.1	0.2	0.4	0.1	0.0	0.0	0.1	0.0	0.5		
27	0.0	0.5	4.4	0.4	1.3	1.2	1.3	0.7	0.3	0.1	0.7	2.1	2.4	0.1	0.1	1.6	2.7	0.1	21.7	1.4	0.0	1.1	1.1	3.5	4.6	1.2	9.9	3.8	7.1	31.7	34.0	0.2	0.7	1.1	0.0	0.0	0.1	0.0	0.3			
28	0.0	0.0	0.2	0.3	0.0	0.4	0.1	0.2	0.0	0.0	0.3	0.1	0.1	0.0	0.0	0.1	0.1	0.0	1.3	0.3	0.0	0.2	0.5	0.5	1.3	3.4	5.7	1.3	3.0	1.5	6.8	0.1	2.7	2.7	0.3	0.2	0.8	0.4	4.3	74.9		
29	4.6	8.4	11.1	6.6	20.4	54.7	88.5	88.1	89.6	86.6	95.3	96.3	86.9	68.1	60.2	67.6	42.6	4.0	13.6	3.8	24.2	7.6	24.7	84.4	25.8	4.9	5.6	14.9	17.6	21.7	5.0	2.9	0.9	4.8	3.1	0.4	0.1	0.4	0.8	0.2	4.4	
30	1.2	3.6	8.4	2.4	2.6	13.5	24.9	34.5	34.4	24.6	52.1	60.1	64.6	27.4	23.8	20.7	16.7	3.2	9.1	2.9	12.0	6.9	8.8	35.6	6.8	2.9	4.4	3.7	3.0	1.0	2.8	0.1	0.9	1.0	0.0	0.1	0.1	0.4	0.6	3.4		
31	1.7	7.6																																								

Transition matrix a_{kl}

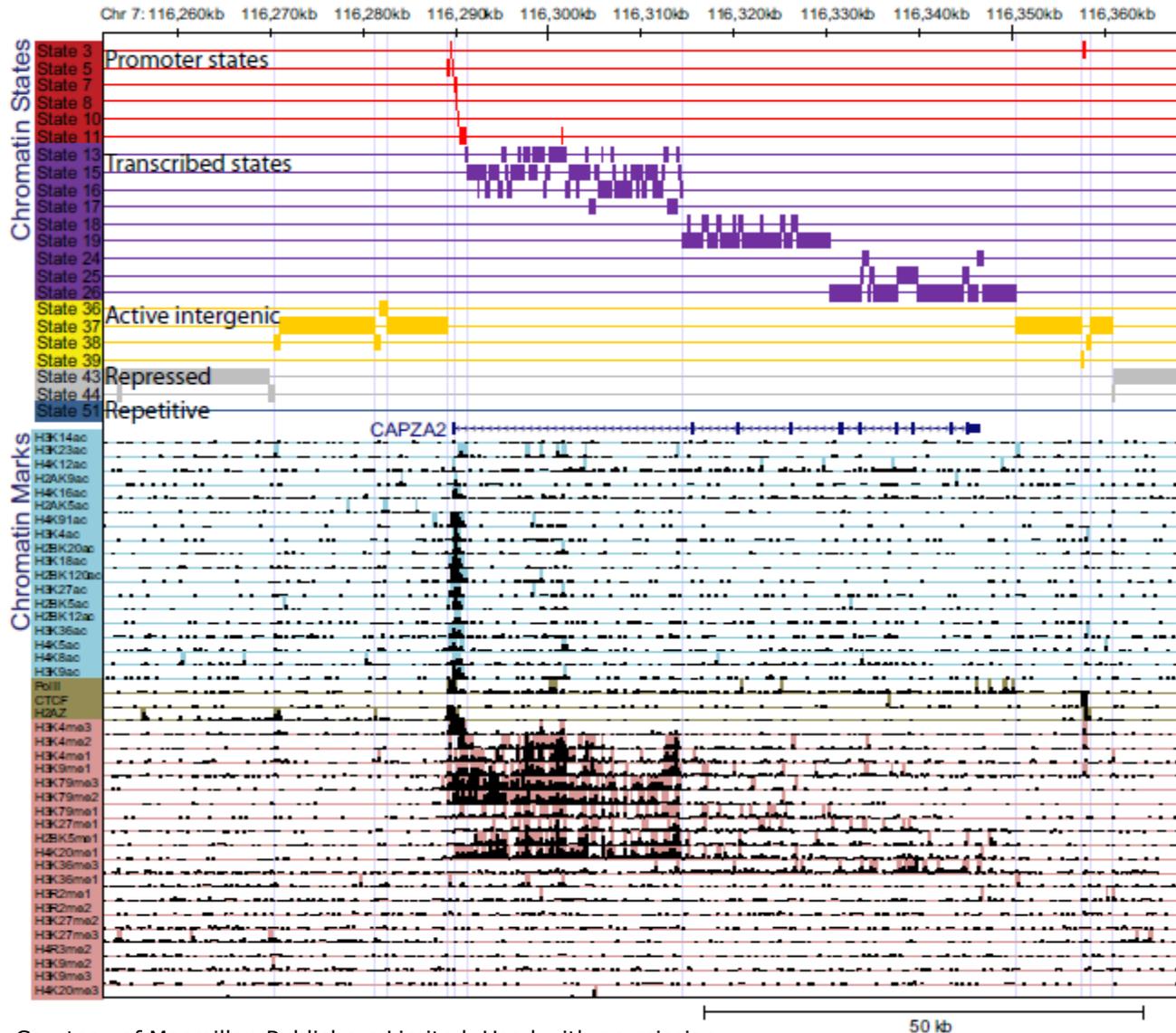


- Learns spatial relationships between neighboring states
- Reveals distinct sub-groups of states
- Reveals transitions between different groups

© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Source: Ernst, Jason and Manolis Kellis. "Discovery and characterization of chromatin states for systematic annotation of the human genome." Nature Biotechnology 28, no. 8 (2010): 817-825.

Example Chromatin State Annotation

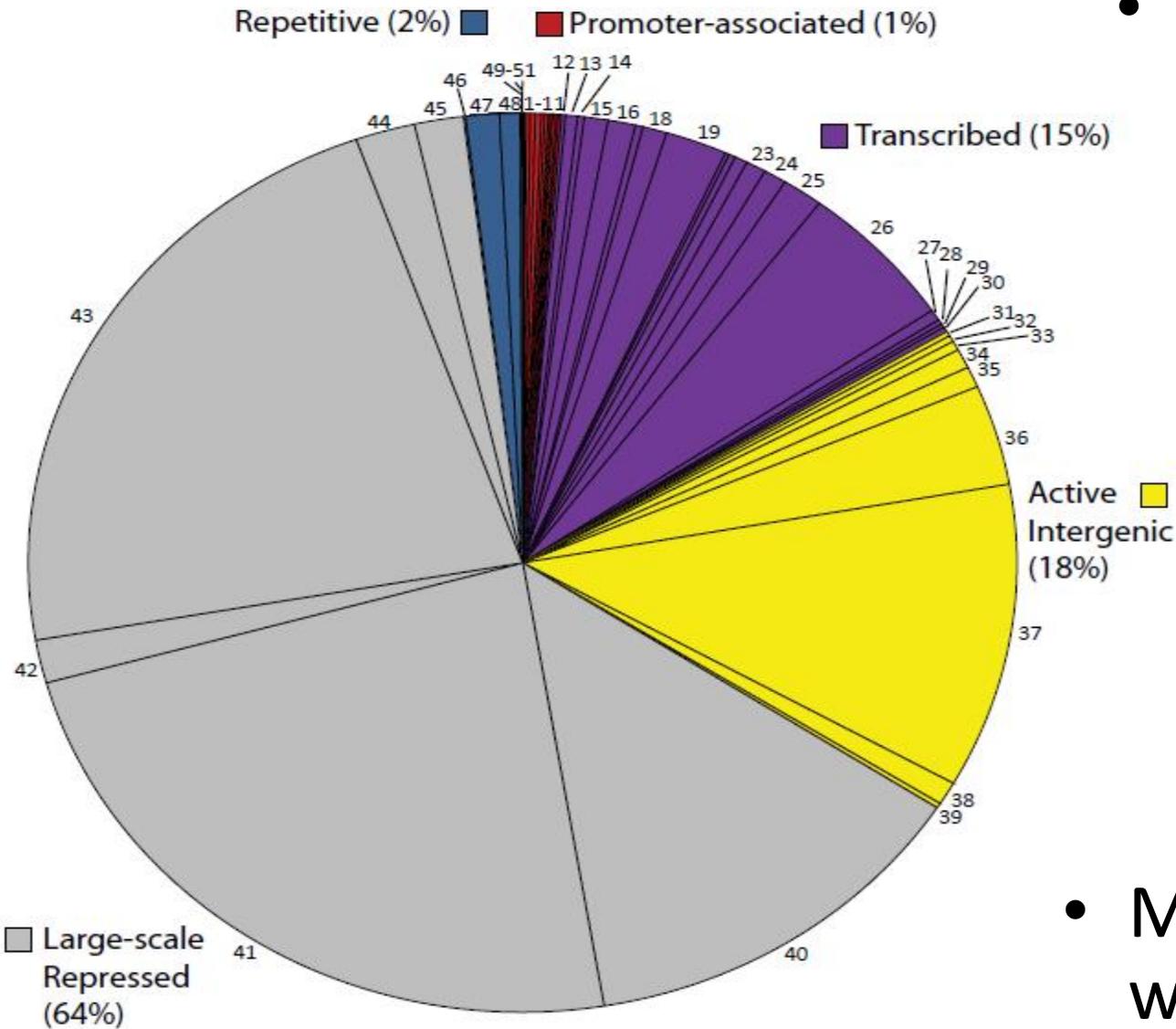


- Use Baum Welch to learn hidden states and their annotations
- Learned states correspond to known functional elements
- *De novo* discovery of major types of chromatin

Courtesy of Macmillan Publishers Limited. Used with permission.

Source: Ernst, Jason and Manolis Kellis. "Discovery and characterization of chromatin states for systematic annotation of the human genome." *Nature Biotechnology* 28, no. 8 (2010): 817-825.

Model complexity matches that of genome

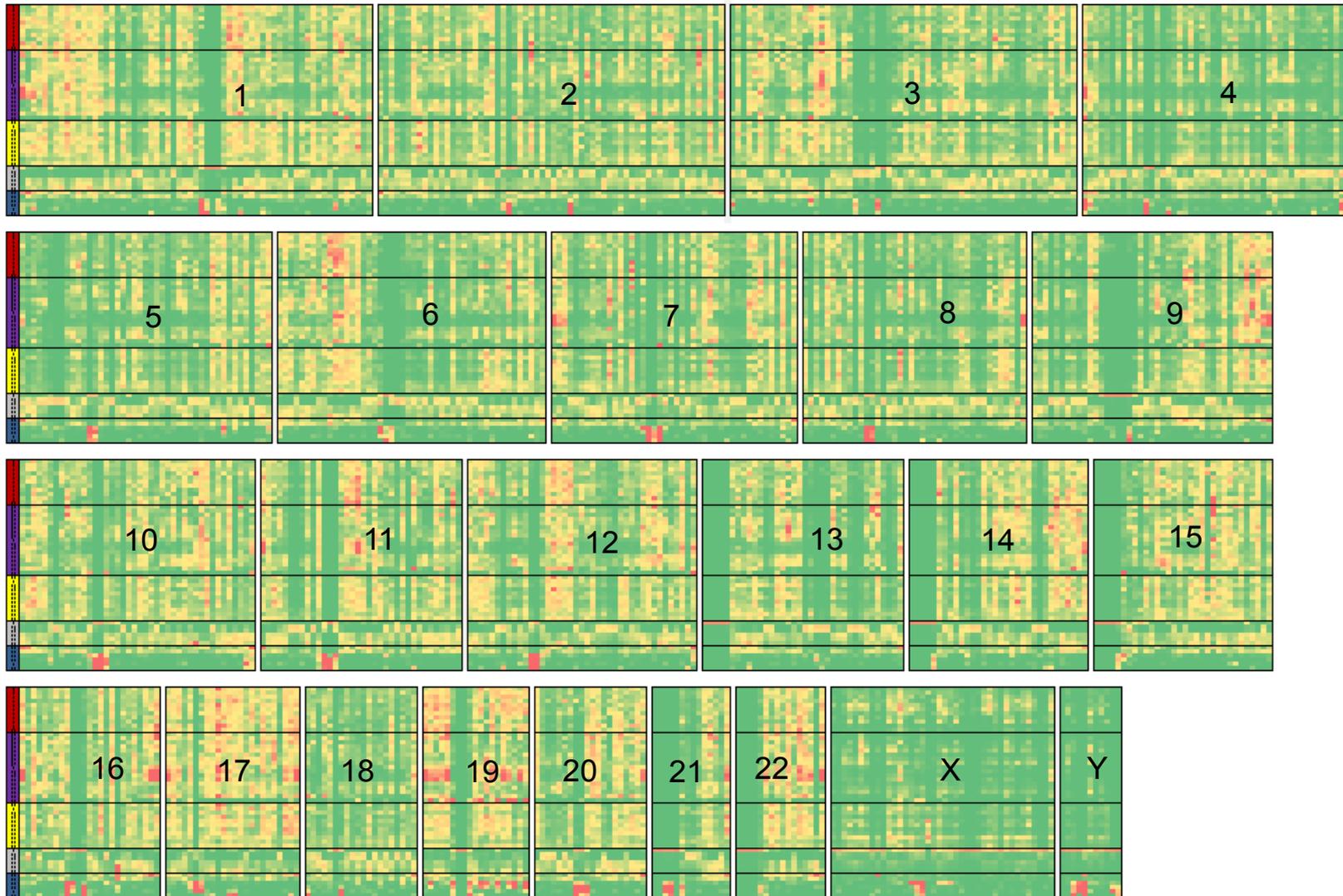


- Handful of repressed states capture vast majority of genome

- Only 1% of genome split in 14 promoter states

- Modeling power well distributed where needed

Apply genome wide to classify chromatin states *de novo*



© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
Source: Ernst, Jason and Manolis Kellis. "Discovery and characterization of chromatin states for systematic annotation of the human genome." Nature Biotechnology 28, no. 8 (2010): 817-825.

Now what? Interpret these states biologically



Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Chromatin state characterization: Functional/positional enrichment

4. Model complexity: selecting the number of states/marks

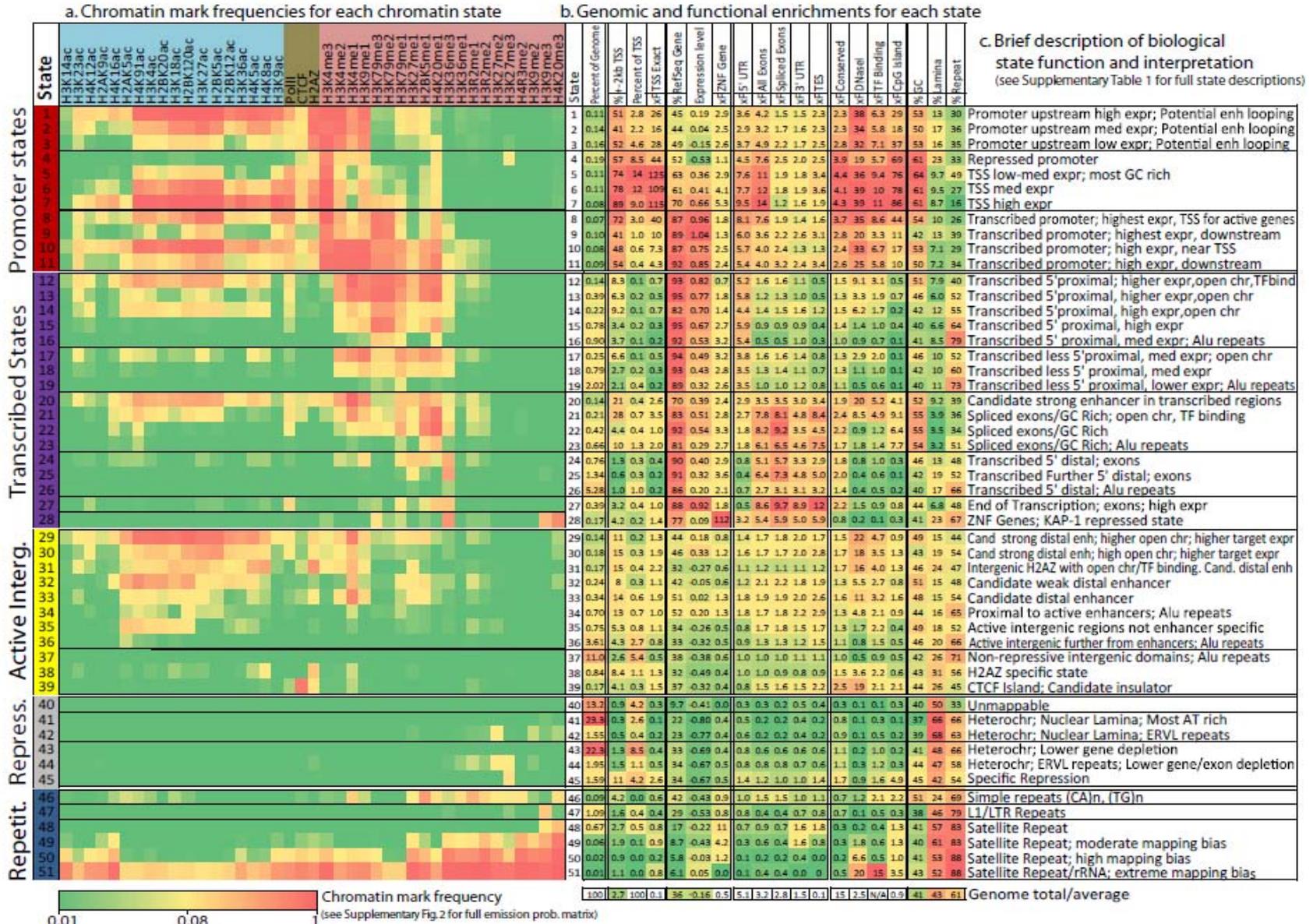
- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

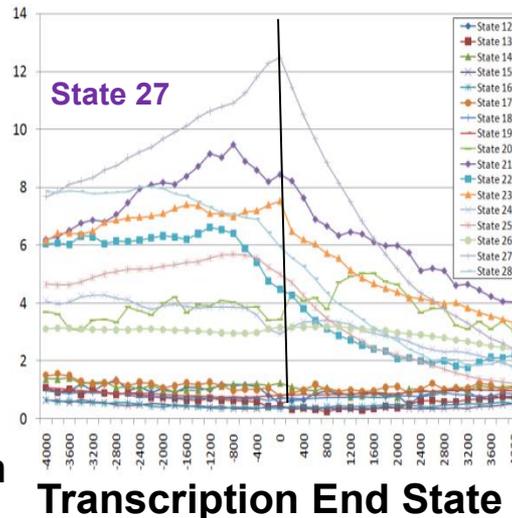
Functional enrichments enable annotation of 51 distinct states



Functional properties of discovered chromatin states

GO Category	State 3	State 4	State 5	State 6	State 7	State 8
Cell Cycle Phase	2.10 (2x10 ⁻⁷)	0.57 (1)	1.61 (0.001)	1.45 (1)	1.15 (1)	1.51 (1)
Embryonic Development	1.24 (1)	2.82 (9x10 ⁻²³)	1.07 (1)	0.85 (1)	0.54 (1)	1.00 (1)
Chromatin	1.20 (1)	0.48 (1)	2.2 (1.4x10 ⁻⁷)	1.64 (1)	0.85 (1)	0.85 (1)
Response to DNA Damage Stimulus	1.20 (1)	0.35 (1)	1.55 (0.074)	2.13 (6.5x10 ⁻¹¹)	1.97 (1.0x10 ⁻⁴)	0.84 (1)
RNA Processing	0.49 (1)	0.26 (1)	1.31 (1)	1.91 (4.2x10 ⁻¹¹)	2.64 (8.7x10 ⁻²⁴)	2.45 (3.0x10 ⁻⁴)
T cell Activation	0.77 (1)	0.88 (1)	1.27 (1)	0.70 (1)	0.79 (1)	4.72 (2x10 ⁻⁷)

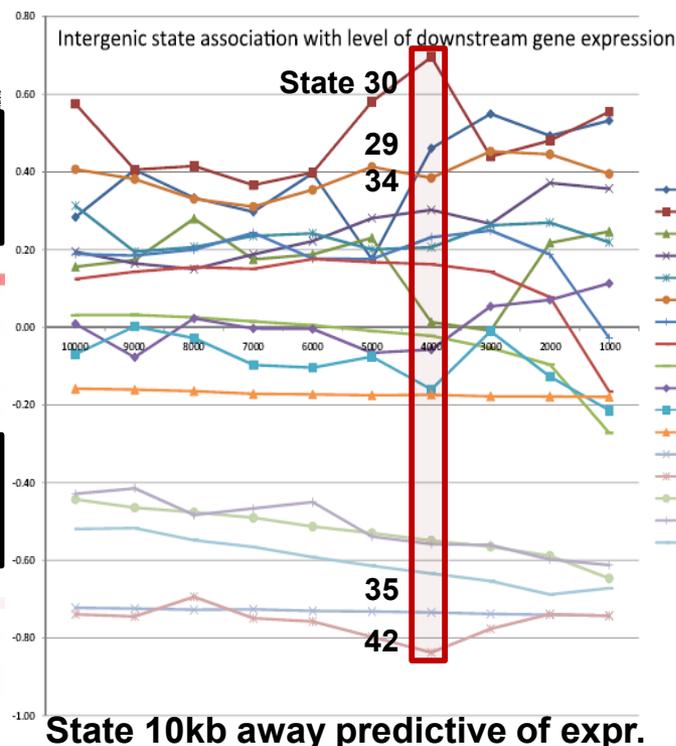
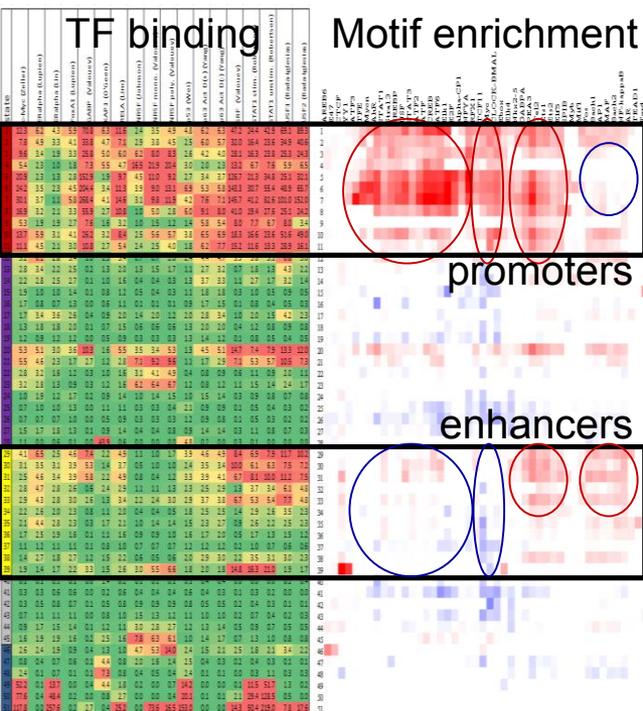
Promoter state → gene GO function



State 28: 112-fold ZNF enrich

“The achievement of the repressed state by wild-type KAP1 involves decreased recruitment of RNA polymerase II, reduced levels of histone **H3 K9 acetylation** and **H3K4 methylation**, an increase in histone occupancy, enrichment of **trimethyl histone H3K9, H3K36**, and **histone H4K20** ...” MCB 2006.

ZNF repressed state recovery



state	stalk	variable heterochromatic	acrocentric	gneg	gpos25	gpos50	gpos75	gpos100
40	7.6	6.6	6.1	0.6	0.4	0.5	0.3	0.5
41	0.0	0.2	0.4	0.5	0.4	0.9	1.7	2.5
42	0.0	0.2	0.3	0.5	0.5	1.1	1.7	2.2
43	0.0	0.1	0.1	1.1	1.2	1.3	1.1	0.7
44	0.0	0.1	0.1	1.2	1.3	1.3	1.0	0.6
45	0.0	0.1	0.1	1.3	1.6	1.3	0.8	0.4
46	0.0	0.1	0.2	1.7	1.7	0.7	0.4	0.2
47	0.0	0.2	0.1	1.2	1.3	1.3	0.9	0.6
48	0.0	3.2	6.2	0.8	2.2	0.2	0.5	0.4
49	0.0	3.6	11.2	0.5	1.8	0.2	0.4	0.1
50	0.0	4.7	12.0	0.6	0.6	0.2	0.2	0.1
51	0.0	4.4	12.7	0.5	1.4	0.2	0.2	0.1
% Overall	0.6	3.9	3.6	42.1	6.8	13.6	13.1	16.2

Distinct types of repression

- Chrom bands / HDAC resp
- Repeat family / composition

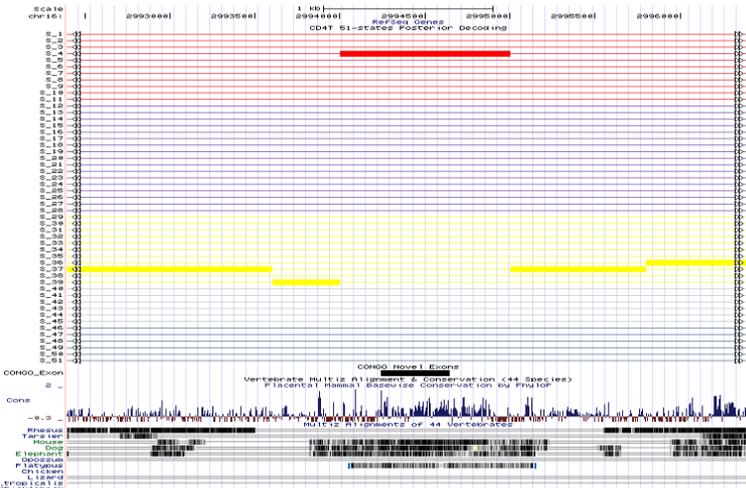
Promoter vs. enhancer regulation

State 10kb away predictive of expr.

© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

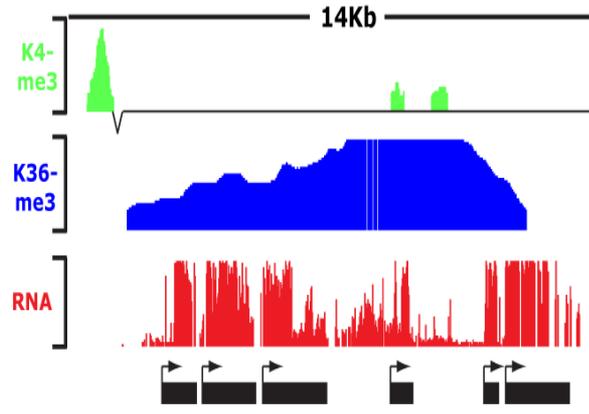
Applications to genome annotation

New protein-coding genes

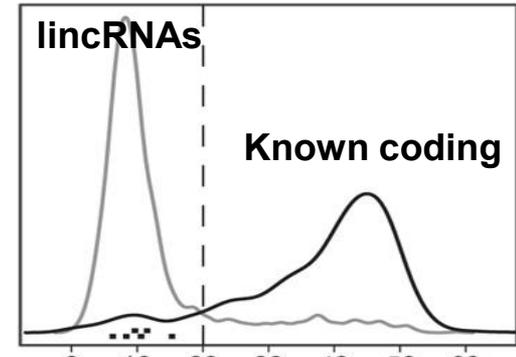


In promoter(short)/low-expr states

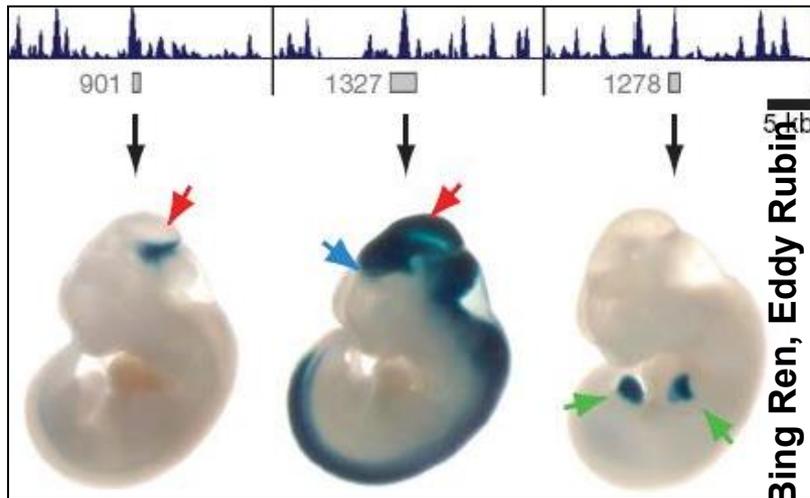
Long intergenic non-coding RNAs/lincRNAs



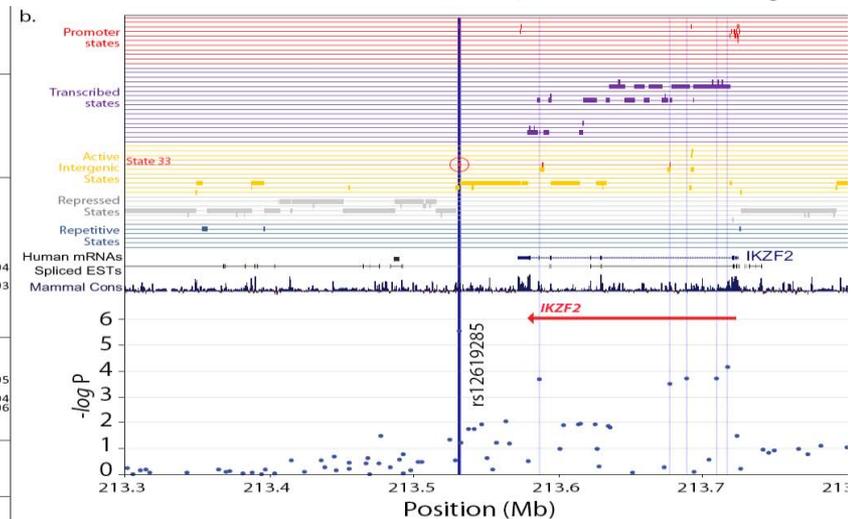
Chromatin signature:
promoter / transcribed



Evolutionary signature:
not protein-coding



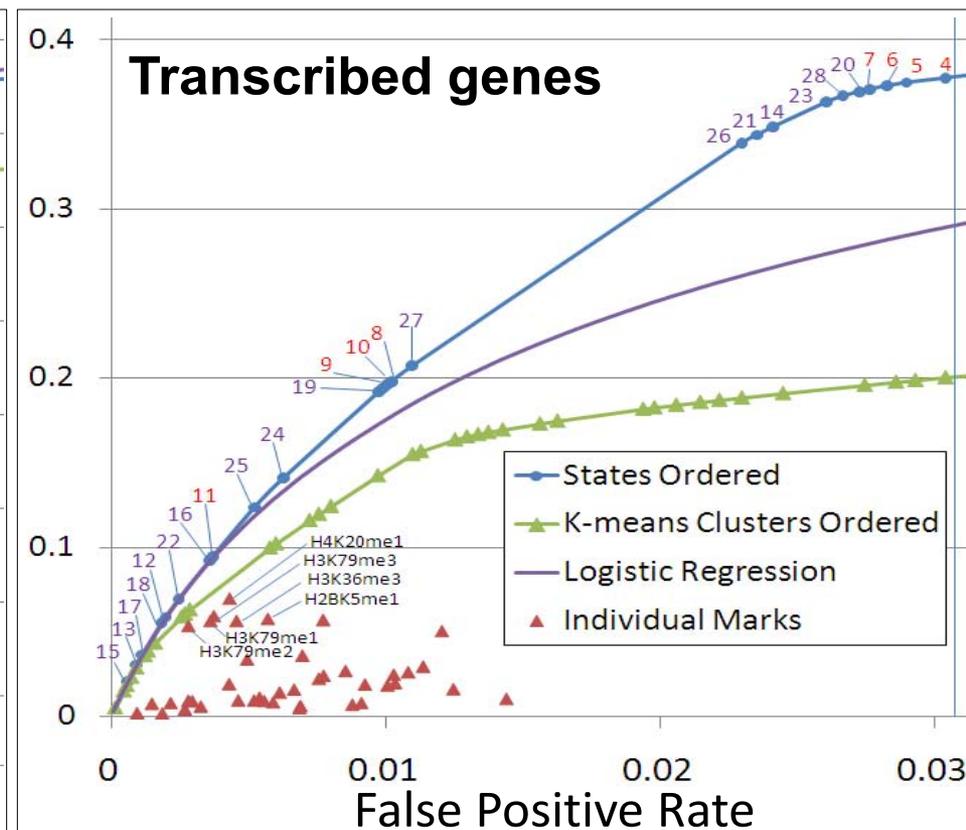
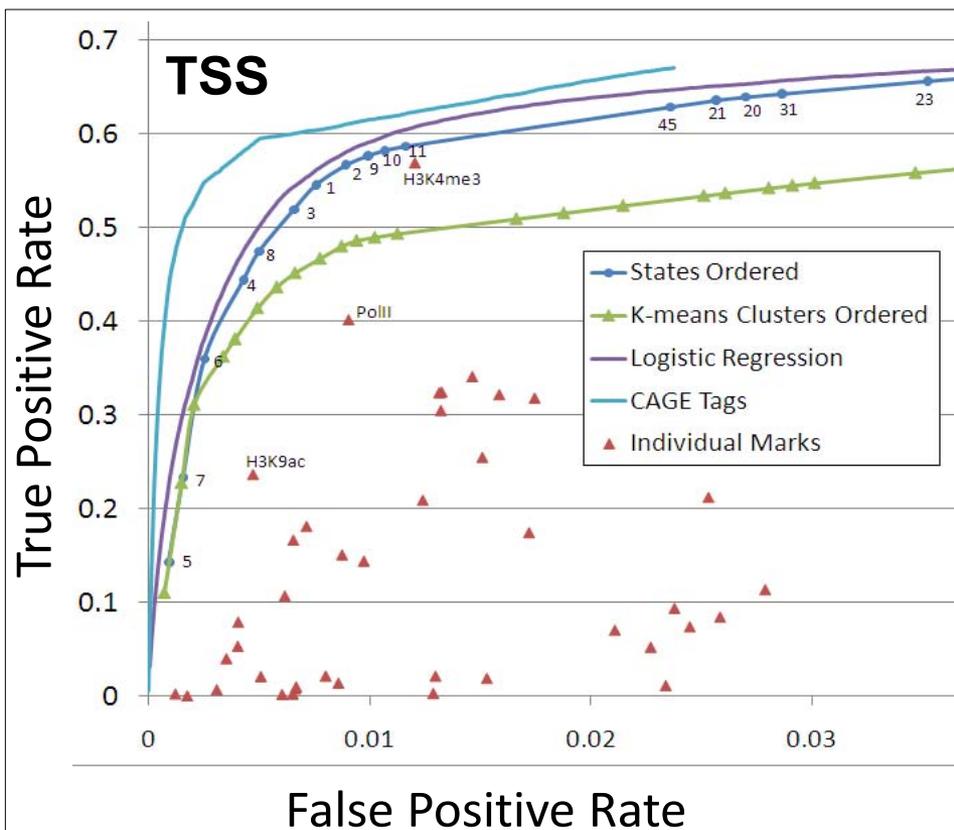
state	% Genome	RegMap	CEU	SNP	GWAS	RegMap	CEU	SNP	GWAS	p-value
1	0.1	1.2	1.7	1.3						
2	0.1	1.2	1.2	1.0						
3	0.2	0.9	2.5	2.6						
4	0.2	0.6	1.4	2.4						
5	0.1	0.8	2.1	2.6						
6	0.1	1.0	1.3	1.2						
7	0.1	0.9	2.4	2.6						
8	0.1	0.9	2.0	2.1						
9	0.1	1.1	1.4	1.3						
10	0.1	1.0	1.2	1.2						
11	0.1	1.0	1.3	1.2						
12	0.1	1.1	1.6	1.5						
13	0.1	1.0	1.3	1.2						
14	0.2	1.3	3.4	2.6						
15	0.8	0.9	0.8	0.8						
16	0.9	0.7	0.9	1.4						
17	0.2	1.2	1.7	1.0						
18	0.8	1.1	1.6	1.5						
19	2.0	0.8	0.9	1.1						
20	0.1	1.2	2.1	1.7						
21	0.2	1.0	5.9	4.1	4.6E-04					
22	0.4	1.1	2.9	1.8						
23	0.7	0.8	2.0	2.5	3.2E-03					
24	0.8	1.2	2.1	1.7						
25	1.3	1.2	1.6	1.4						
26	5.3	0.9	1.0	1.1						
27	0.4	1.1	1.9	1.7						
28	0.2	1.0	2.6	2.5						
29	0.2	1.3	3.8	2.1						
30	0.2	1.2	0.9	0.8						
31	0.2	1.3	3.0	2.3						
32	0.2	1.3	1.4	1.0						
33	0.3	1.2	3.8	3.3	5.2E-05					
34	0.7	1.0	1.5	1.6						
35	0.7	1.4	3.0	2.2	5.8E-04					
36	3.6	1.1	1.9	1.7	3.6E-06					
37	11.0	0.9	1.1	1.2						
38	0.8	1.3	1.8	1.4						
39	0.2	1.3	1.6	1.1						
40	13.2	0.2	0.1	0.5						
41	23.8	1.2	0.7	0.5						
42	1.5	1.4	0.8	0.6						
43	22.3	1.2	1.2	1.0						
44	2.0	1.4	1.4	1.0						
45	1.6	1.2	1.5	1.2						
46	0.1	1.0	0.9	0.9						
47	1.1	1.1	1.5	1.3						
48	0.7	0.7	0.3	0.5						
49	0.1	0.7	0.3	0.4						
50	0.0	0.5	0.0	0.0						
51	0.0	0.4	0.0	0.0						



Assign candidate functions to intergenic SNPs
from genome-wide association studies

New developmental enhancer regions

Discovery power for promoters, transcripts



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- Significantly outperforms single-marks
- Similar power to supervised learning approach
- CAGE experiments give possible upper bound

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Capturing dependencies. State-conditional mark independence
- Selecting the number of states, selecting number of marks

5. Learning chromatin states jointly across multiple cell types

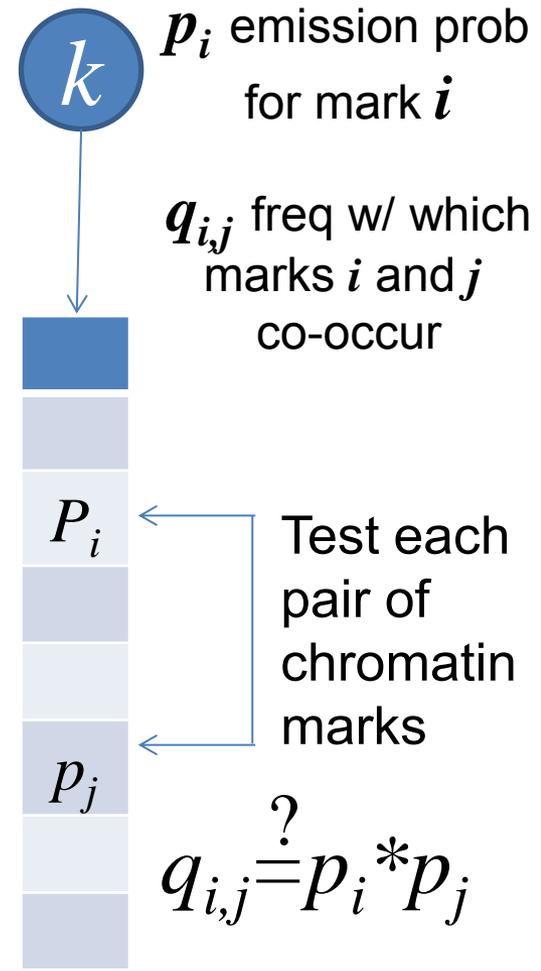
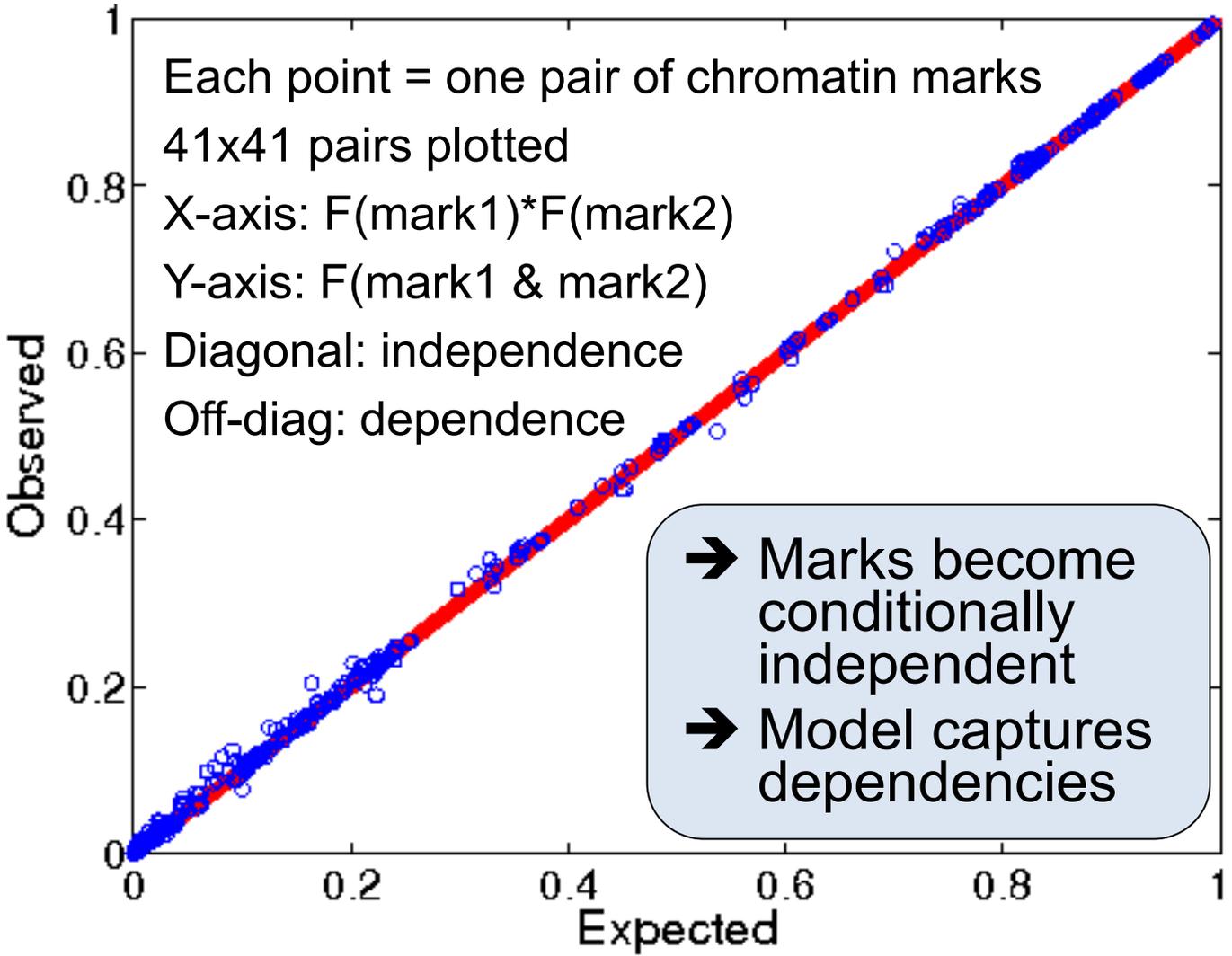
- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

State-conditional mark independence

Do hidden states actually capture dependencies between marks?

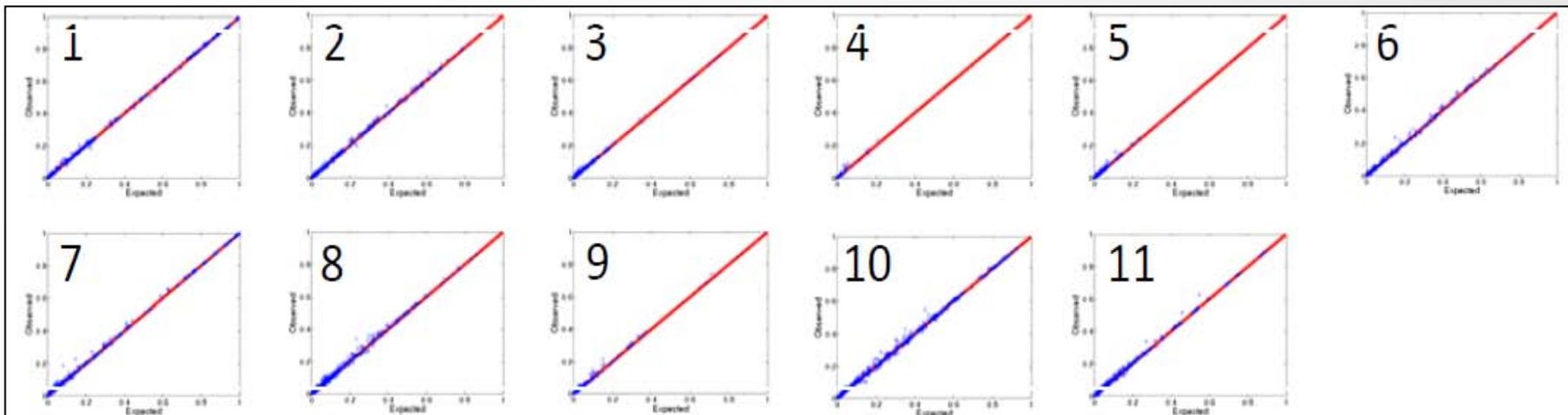
Pairwise Expected vs. Observed Mark Co-Occurrence



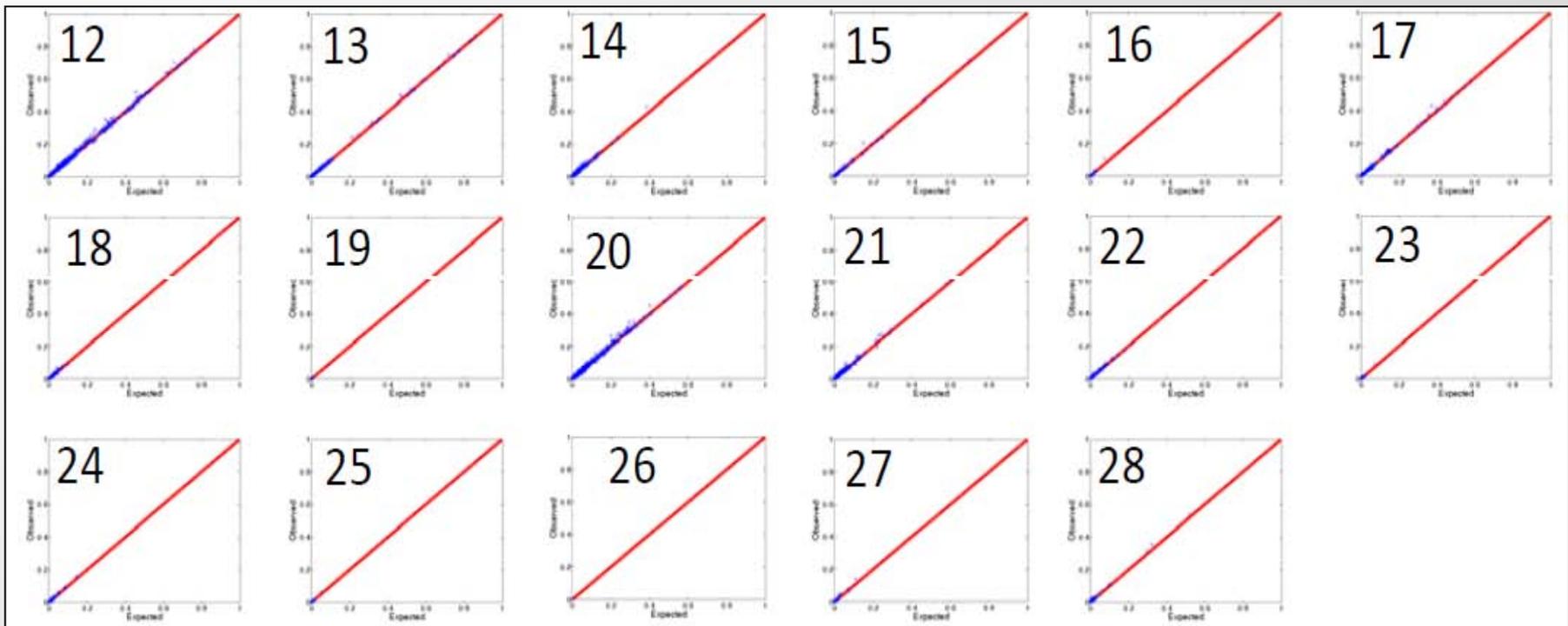
Multi-variate HMM emits entire vector of marks at a time
 Model assumes mark independence *conditional* upon state
 In fact, it specifically seeks to *capture* these dependencies

Test conditional independence for each state

Promoter states

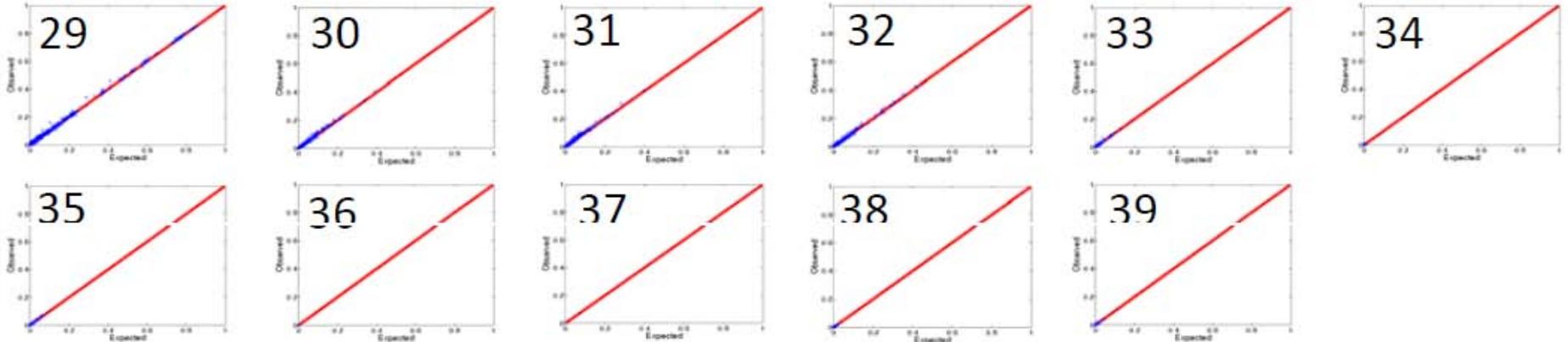


Transcribed states

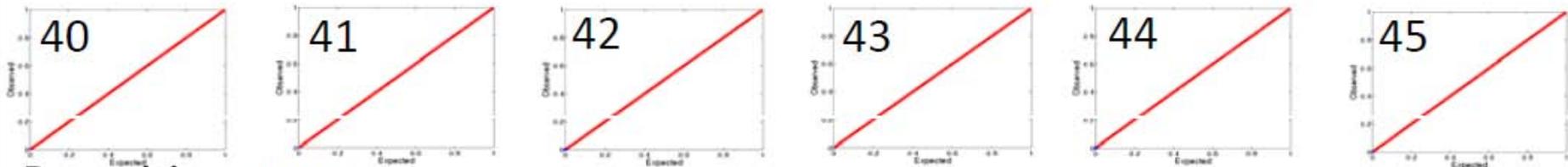


Non-independence reveals cases of model violation

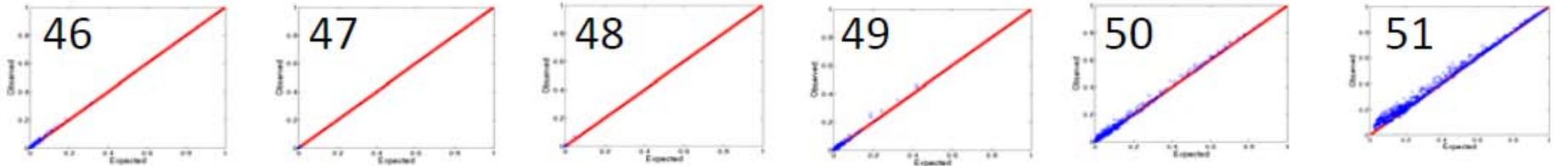
Active Intergenic states



Repressed states

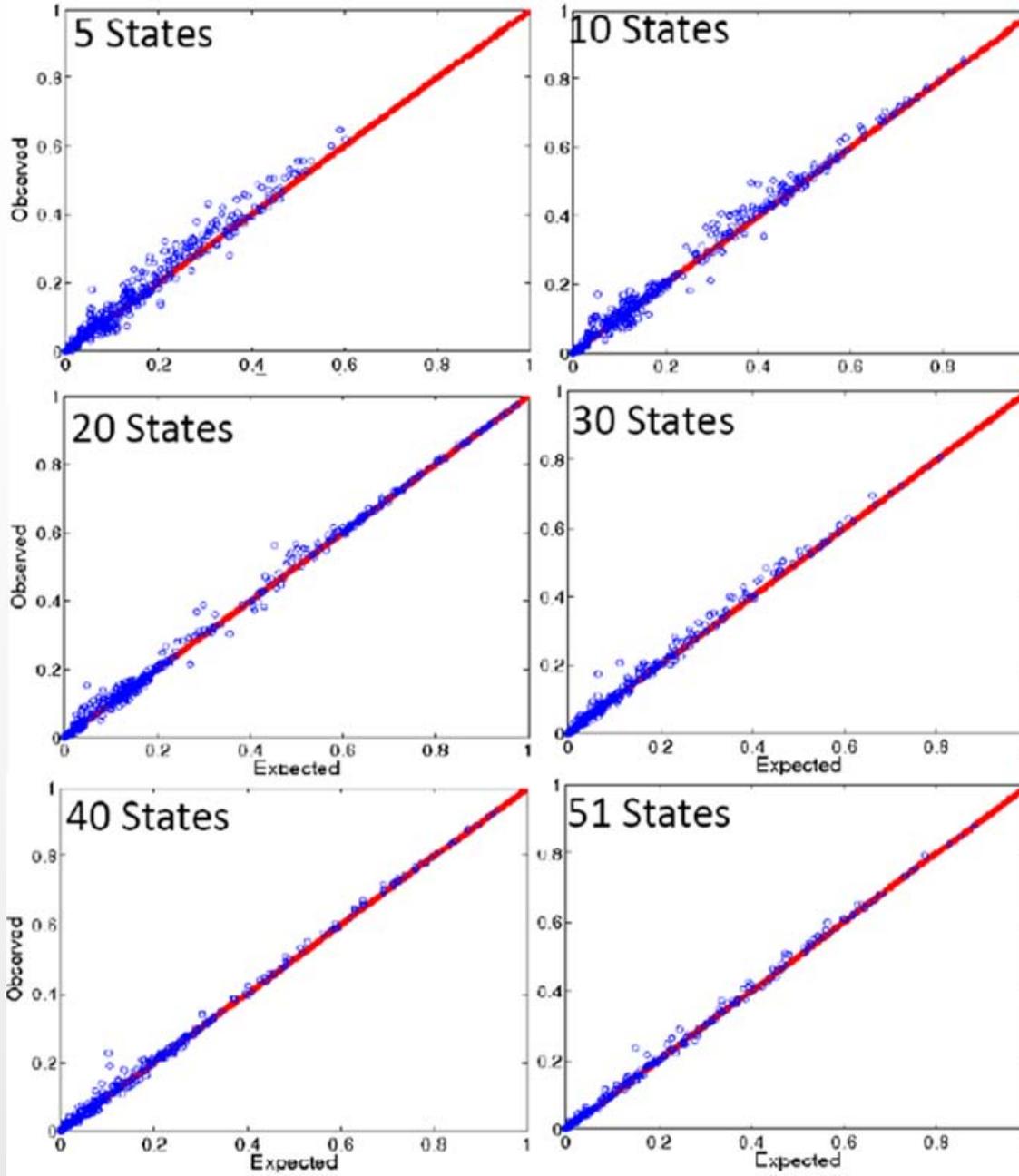


Repetitive states



- Repetitive states show more dependencies
- Conditional independence does not hold

As more states are added, dependencies captured



- With only 5 states in HMM, not enough power to distinguish different properties
- Dependencies remain
- As model complexity increases, states learned become more precise
- Dependencies captured

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

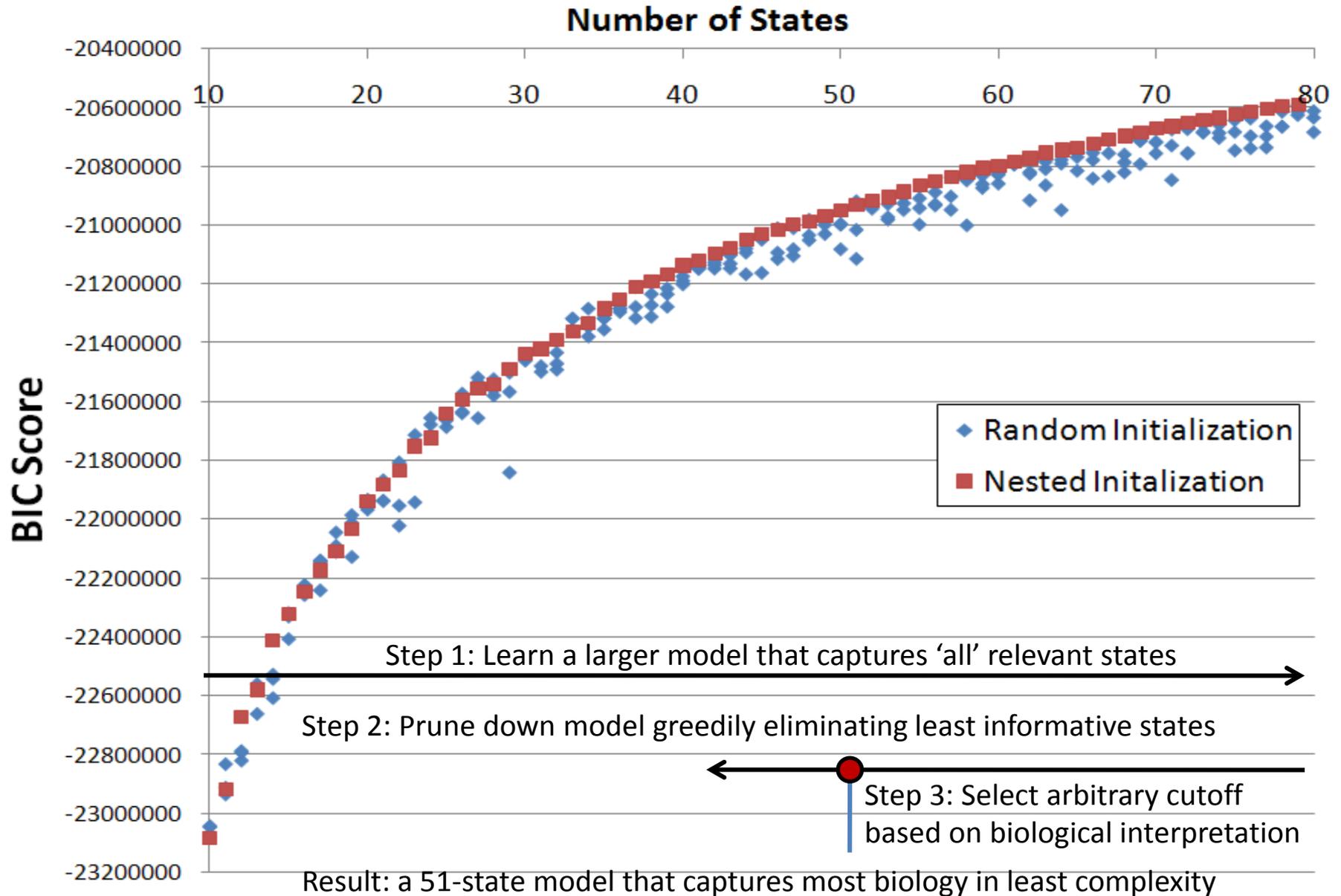
- Capturing dependencies. State-conditional mark independence
- Selecting the number of states, selecting number of marks

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

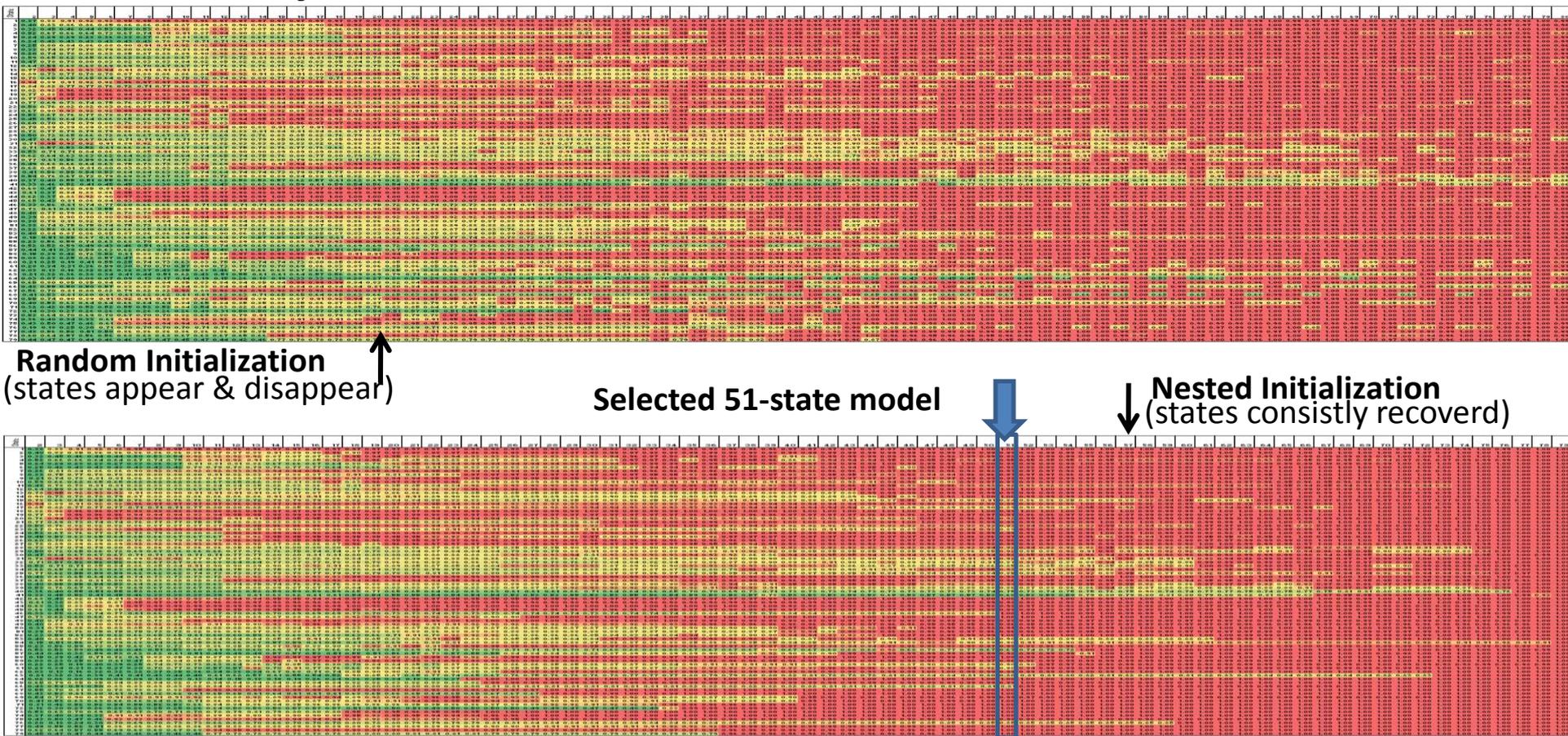
(Future: Chromatin states to interpret disease-associated variants)

Comparison of BIC Score vs. Number of States for Random and Nested Initialization



- Standard model selection criteria fail due to genome complexity: more states always preferred
- Instead: Start w/complex model, keep informative states, prune redundant states. Pick cutoff

Recovery of 79-state model in random vs. nested initialization



Random Initialization
(states appear & disappear)

Selected 51-state model

Nested Initialization
(states consistently recovered)

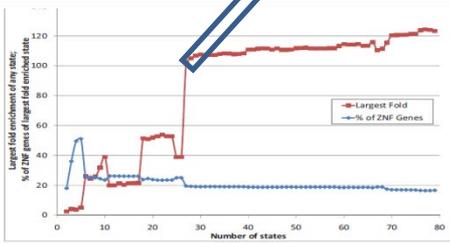
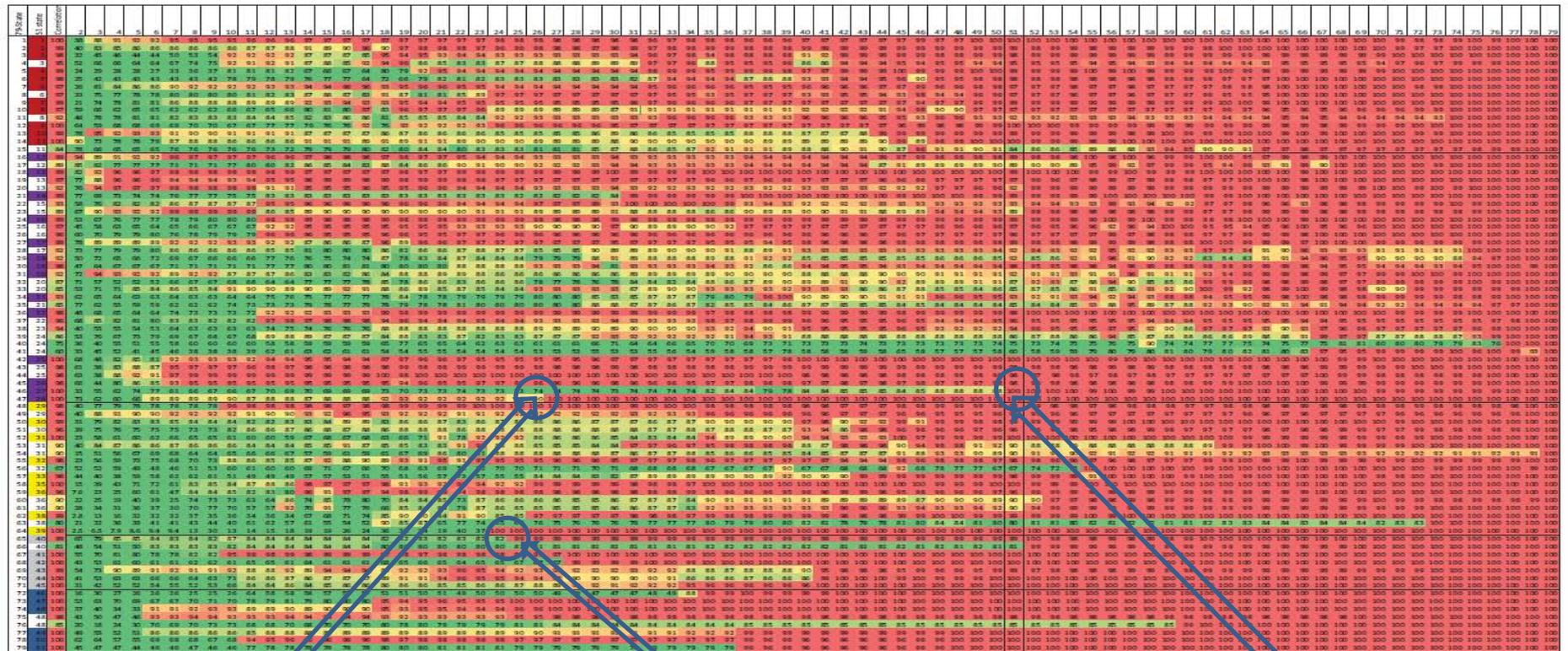
Nested initialization approach:

- **First pass:** learn models of increasing complexity
- **Second pass:** form nested set of emission parameter initializations by greedily removing states from best BIC model found

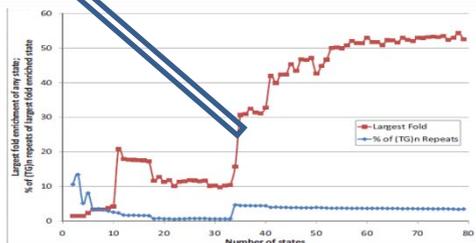
Nested models criteria:

- Maximize sum of correlation of emission vectors with nested model
- Models learned in parallel

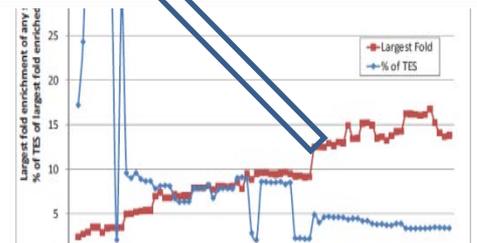
Functional recovery with increasing numbers of states



Zinc Finger state



Simple Repeat state



Transcription End State

- Red: Maximum fold enrichment for corresponding biological category
- Blue: Percent of that functional category that overlaps regions annotated to this state
- Top plot: Correlation of emission parameter vector for that state to closest state

Chromatin state recovery with increasing numbers of marks

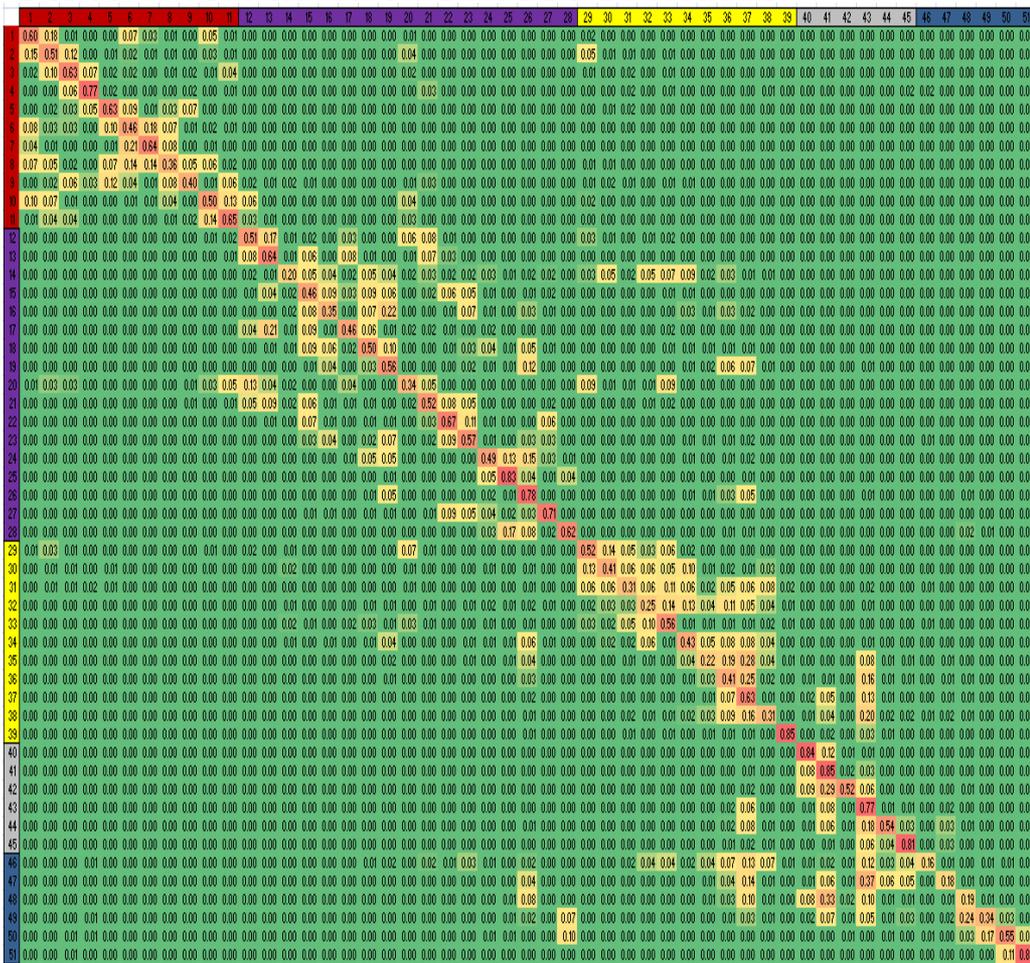
Which states are well-recovered?

Precisely what mistakes are made?

Increasing numbers of marks (greedy)

(for a given subset of 11 ENCODE marks)

State Inferred with subset of marks



State Inferred with all 41 marks

State Inferred with all 41 marks

State confusion matrix with 11 ENCODE marks

Recovery of states with increasing number of marks

© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Ernst, Jason and Manolis Kellis. "Discovery and characterization of chromatin states for systematic annotation of the human genome." Nature Biotechnology 28, no. 8 (2010): 817-825.

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

ENCODE: Study nine marks in nine human cell types

9 marks

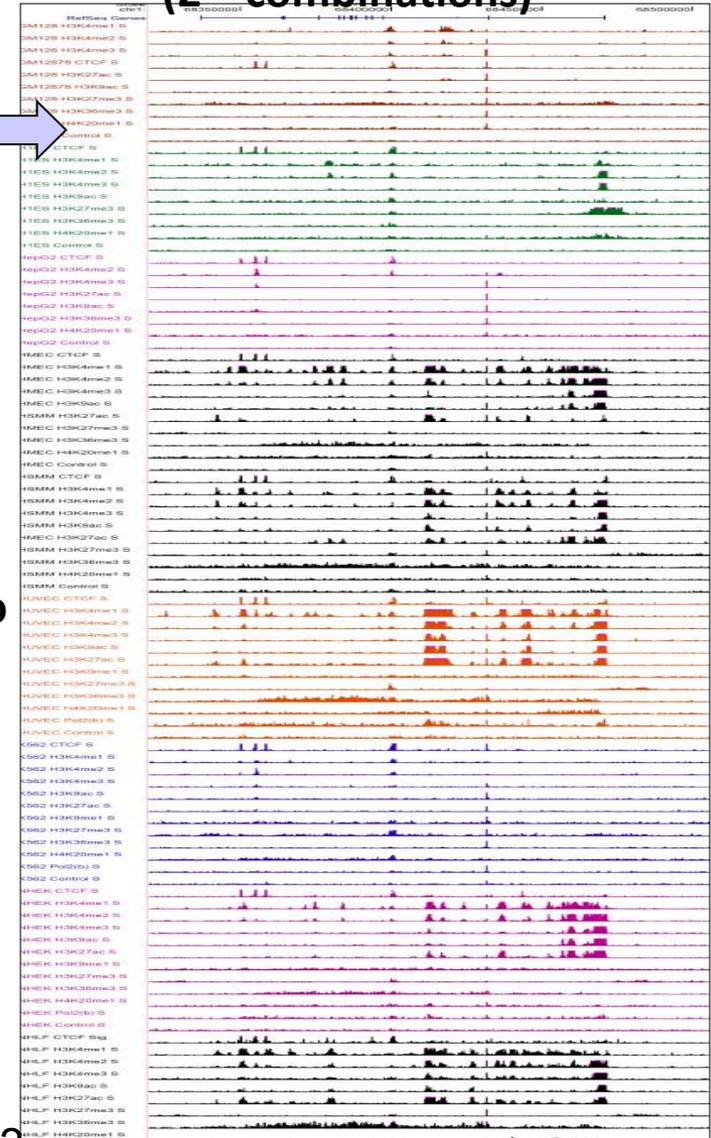
H3K4me1
H3K4me2
H3K4me3
H3K27ac
H3K9ac
H3K27me3
H4K20me1
H3K36me3
CTCF
+WCE
+RNA

X

9 human cell types

HUVEC	Umbilical vein endothelial
NHEK	Keratinocytes
GM12878	Lymphoblastoid
K562	Myelogenous leukemia
HepG2	Liver carcinoma
NHLF	Normal human lung fibroblast
HMEC	Mammary epithelial cell
HSMM	Skeletal muscle myoblasts
H1	Embryonic

81 Chromatin Mark Tracks
(2^8 combinations)

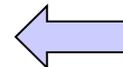


Brad Bernstein ENCODE Chromatin Group

b.

Chromatin States	State	Chromatin Mark Observation Frequency (%)								
		CTCF	H3K27me3	H3K36me3	H4K20me1	H3K4me1	H3K4me2	H3K4me3	H3K27ac	H3K9ac
1	16	2	2	6	17	93	99	96	98	2
2	12	2	6	9	53	94	95	14	44	1
3	13	72	0	9	48	78	49	1	10	1
4	11	1	15	11	96	99	75	97	86	4
5	5	0	10	3	88	57	5	84	25	1
6	7	1	1	3	58	75	8	6	5	1
7	2	1	2	1	56	3	0	6	2	1
8	92	2	1	3	6	3	0	0	1	1
9	5	0	43	43	37	11	2	9	4	1
10	1	0	47	3	0	0	0	0	0	1
11	0	0	3	2	0	0	0	0	0	0
12	1	27	0	2	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0
14	22	28	19	41	6	5	26	5	13	37
15	85	85	91	88	76	77	91	73	85	78

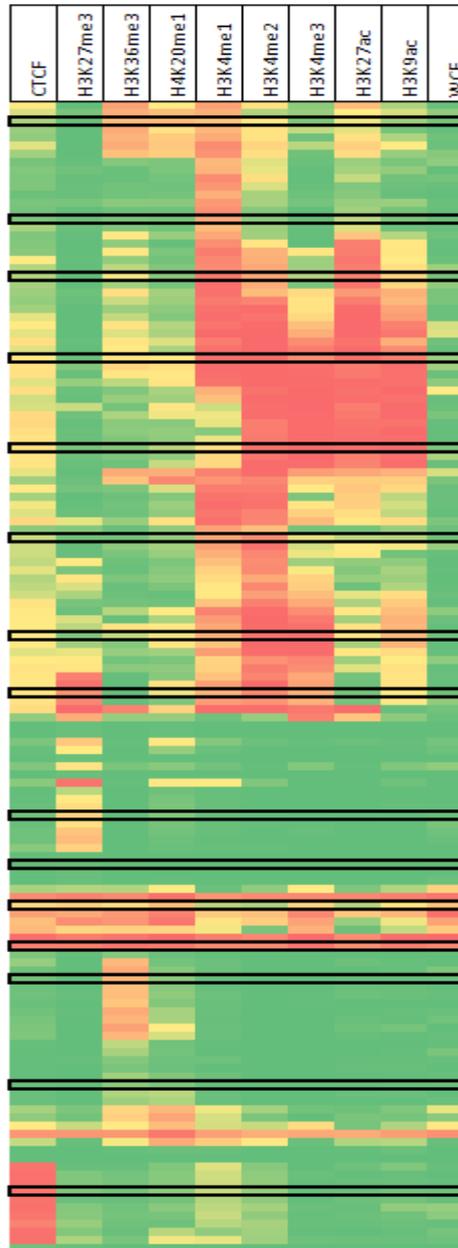
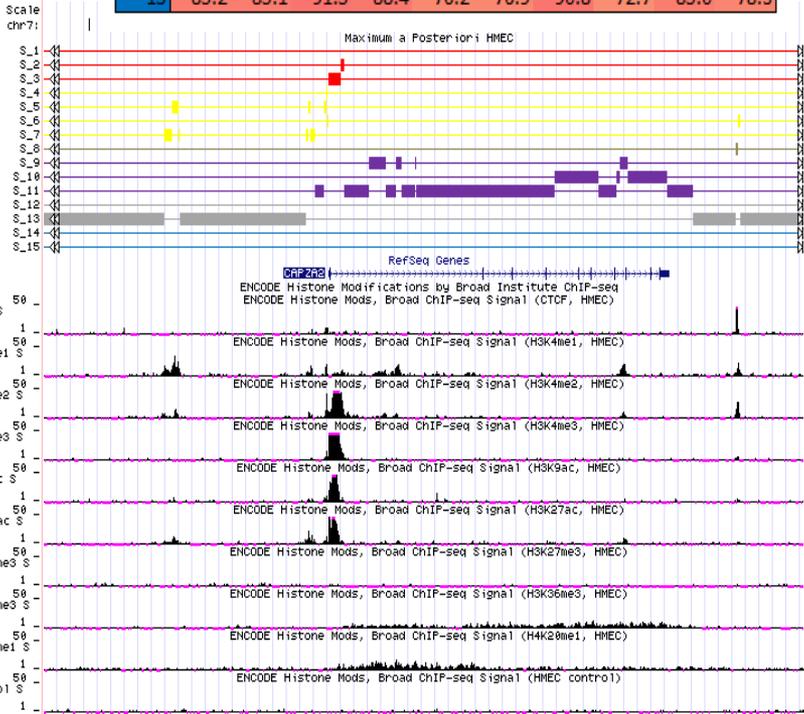
Chromatin Mark Observation Frequency (%)



How to learn single set of chromatin states?

Solution 1: Learn independent models and cluster

	state	CTCF	H3K27me3	H3K36me3	H4K20me1	H3K4me1	H3K4me2	H3K4me3	H3K27ac	H3K9ac	WCE
Promoter	1	13.2	72.0	0.2	9.1	47.9	77.8	49.5	1.3	10.2	0.7
	2	11.9	1.9	6.1	9.0	52.7	93.7	95.0	14.1	44.1	0.9
	3	16.4	1.5	2.4	5.5	17.0	92.6	99.0	95.7	98.1	1.9
Candidate enhancer	4	11.4	0.6	14.5	11.3	96.3	99.3	75.1	97.2	85.7	3.7
	5	5.3	0.2	9.5	2.6	88.1	56.8	5.3	84.4	24.9	1.5
	6	6.7	0.9	1.0	3.2	58.3	74.7	8.4	5.8	5.4	0.8
	7	1.6	0.6	1.6	1.3	56.5	2.7	0.4	5.9	1.6	0.6
Insulator	8	91.5	1.8	0.9	2.8	6.3	3.3	0.4	0.5	1.0	0.8
	9	4.6	0.3	43.2	43.1	36.5	11.5	1.9	9.1	3.9	1.3
Transcribed	10	1.2	0.1	47.2	2.7	0.4	0.0	0.1	0.3	0.3	0.5
	11	0.4	0.1	2.7	1.7	0.2	0.1	0.1	0.2	0.3	0.4
	12	0.9	26.8	0.0	2.1	0.4	0.1	0.1	0.1	0.1	0.4
Repressive	13	0.2	0.4	0.0	0.1	0.1	0.0	0.0	0.0	0.1	0.1
Repetitive	14	21.9	27.9	19.1	41.0	5.7	4.8	25.9	5.3	13.1	37.5
	15	85.2	85.1	91.5	88.4	76.2	76.9	90.8	72.7	85.0	78.3



Basic approach:

- Train a k-state model in each cell type independently
- Cluster models learned independently
- Merge clusters and re-apply to each cell type

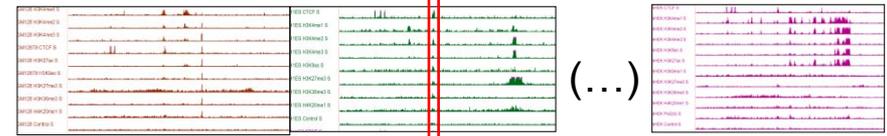
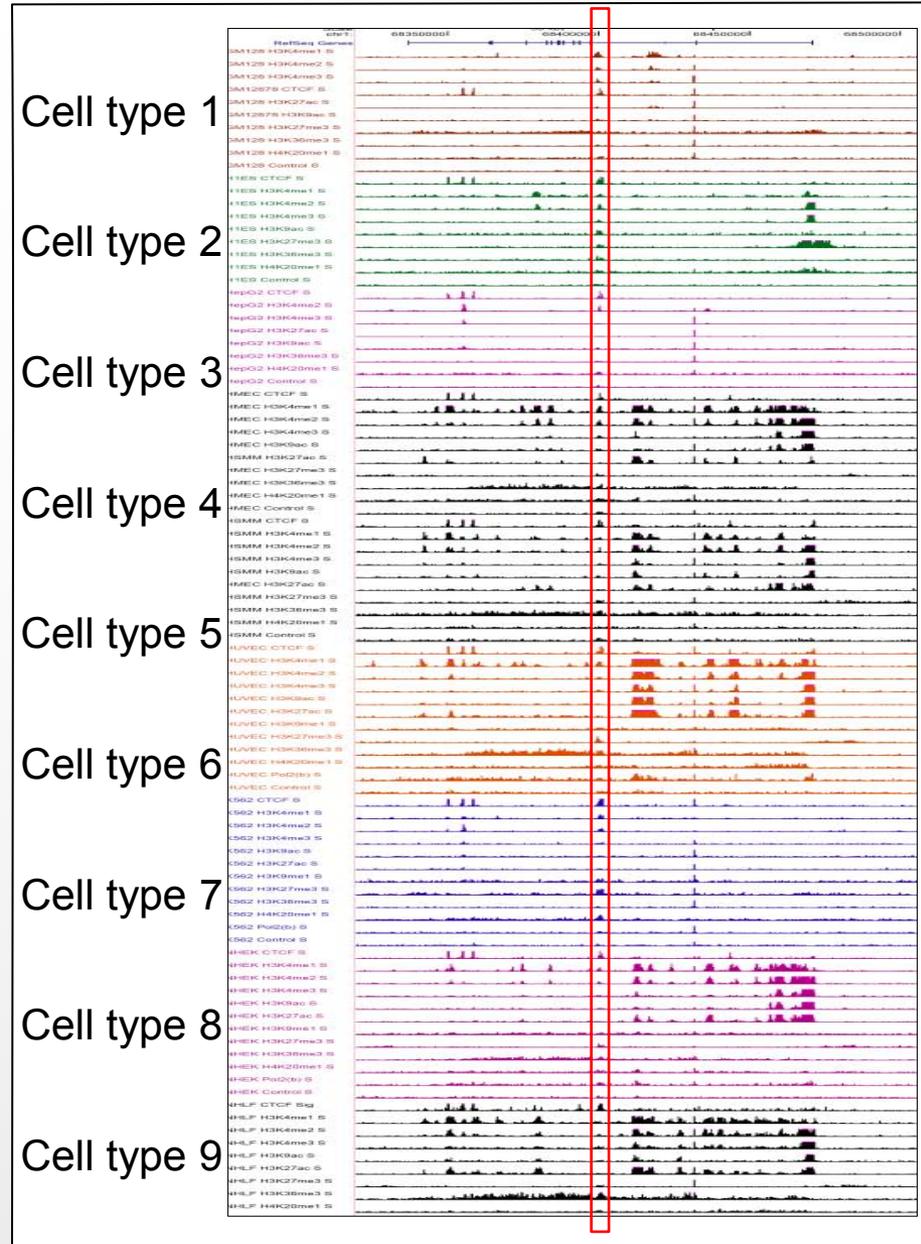
How to cluster

- Using emission probability matrix: most similar definitions
- Using genome annotation: posterior probability decoding

Joint learning of states across multiple cell types

Solution 2: Stacking

- Learns each combination of activity as a separate state
- Ex: ES-specific enhancers: enhancer marks in ES, no marks in other cell types

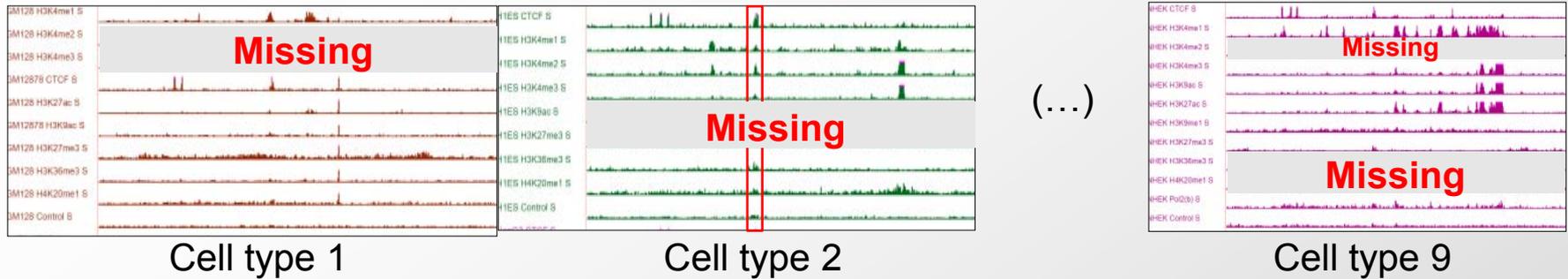


Cell type 1 Cell type 2 Cell type 9

Solution 3: Concatenation

- Requires that profiled marks are the same (or treat as missing data)
- Ensures common state definitions across cell types

Joint learning with different subsets of marks (Solution 3)



Option (a) Treat missing tracks as missing data

- EM framework allows for unspecified data points
- As long as pairwise relationship observed in some cell type

Option (b) Chromatin mark imputation

- Explicitly predict max-likelihood chromatin track for missing data
- Less powerful if ultimate goal is chromatin state learning

ENCODE: Study nine marks in nine human cell lines

9 marks

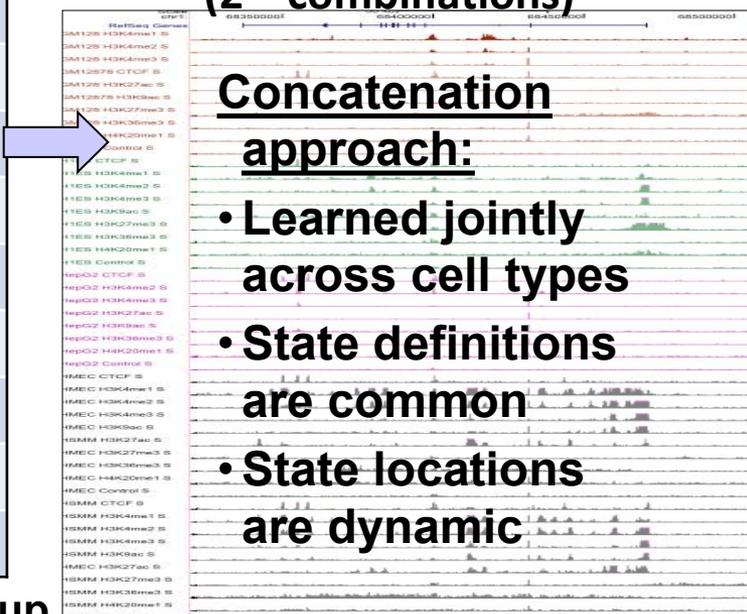
9 human cell types

81 Chromatin Mark Tracks
(2^{81} combinations)

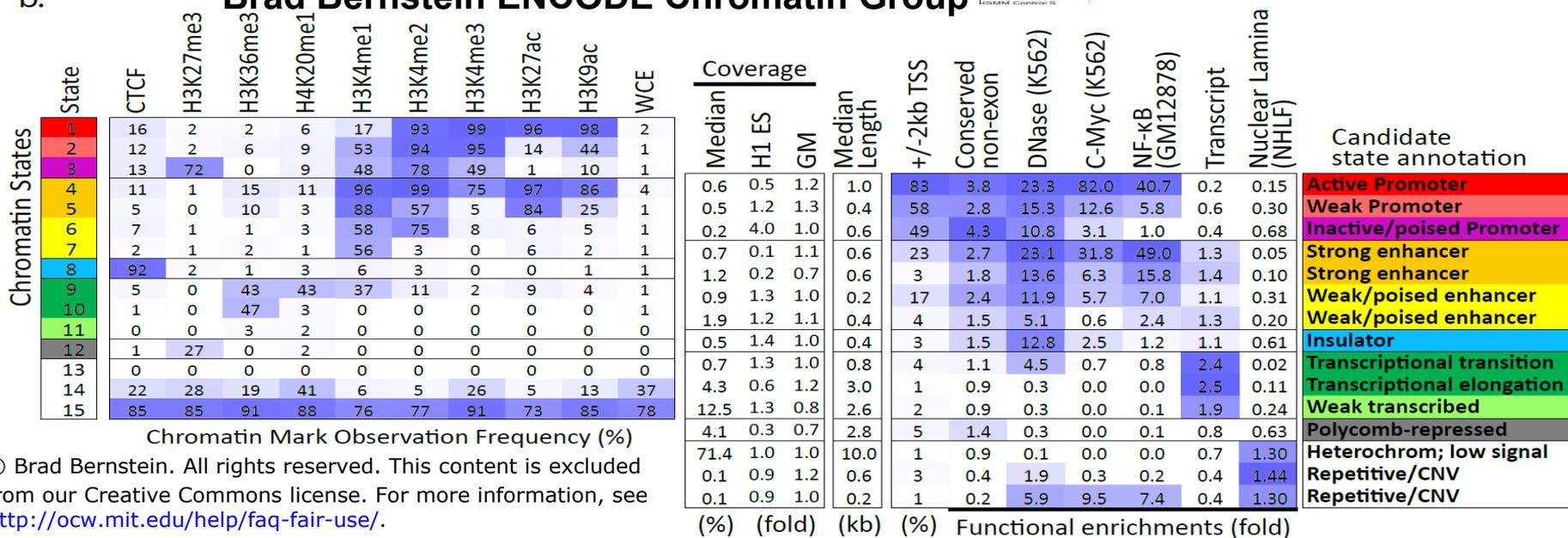
H3K4me1
H3K4me2
H3K4me3
H3K27ac
H3K9ac
H3K27me3
H4K20me1
H3K36me3
CTCF
+WCE
RNA

X

HUVEC	Umbilical vein endothelial
NHEK	Keratinocytes
GM12878	Lymphoblastoid
K562	Myelogenous leukemia
HepG2	Liver carcinoma
NHLF	Normal human lung fibroblast
HMEC	Mammary epithelial cell
HSMM	Skeletal muscle myoblasts
H1	Embryonic

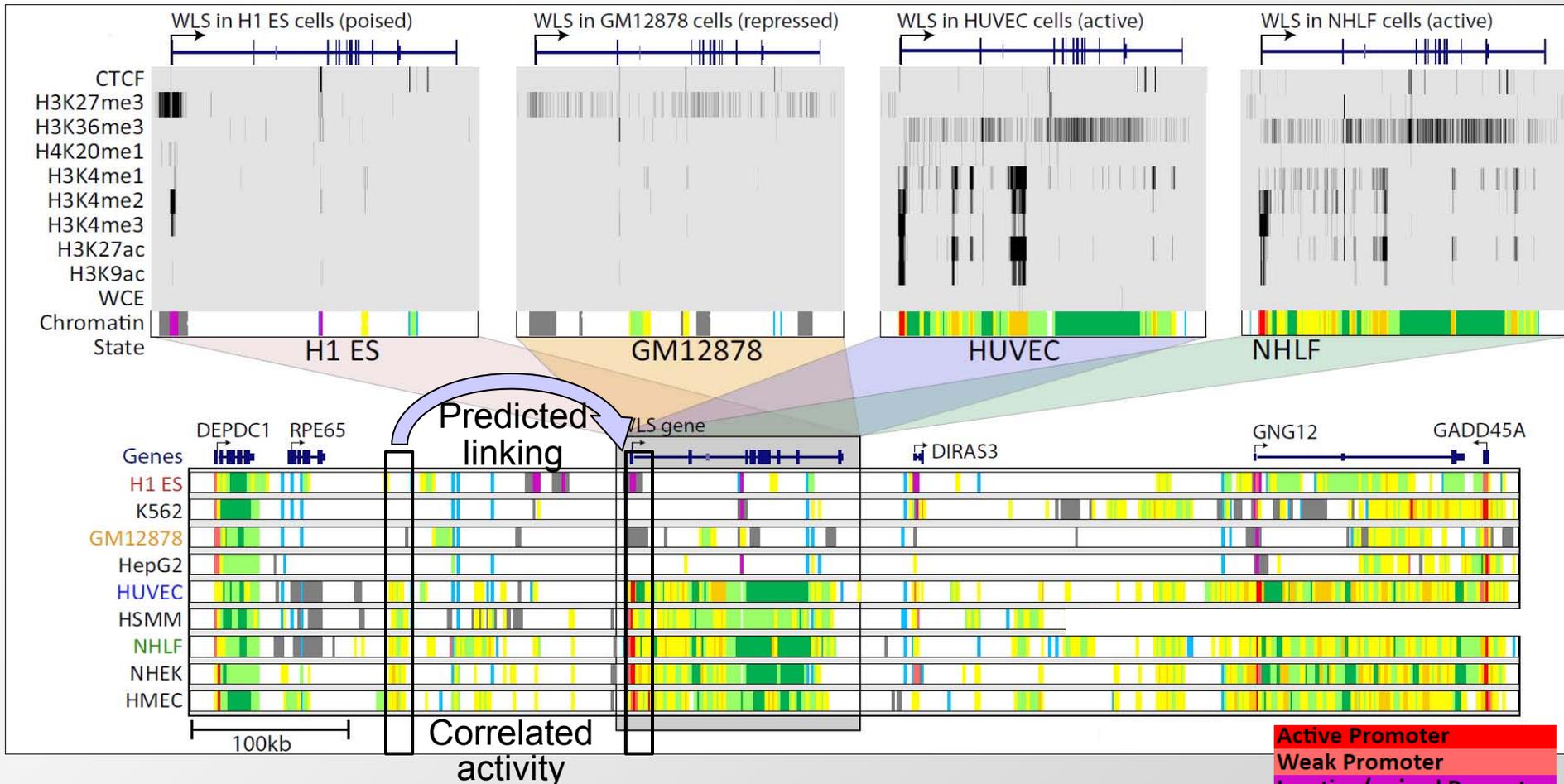


b. Brad Bernstein ENCODE Chromatin Group



© Brad Bernstein. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Chromatin states dynamics across nine cell types

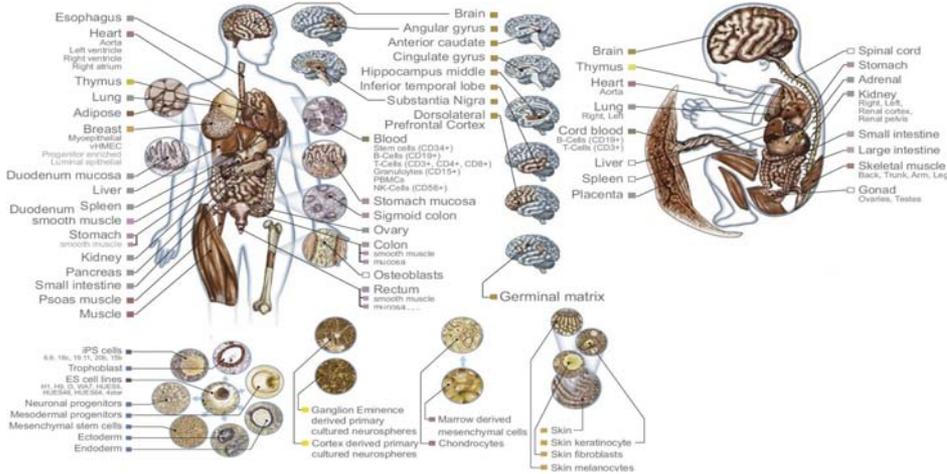


© Brad Bernstein. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- Single annotation track for each cell type
- Summarize cell-type activity at a glance ↓
- Can study 9-cell activity pattern across

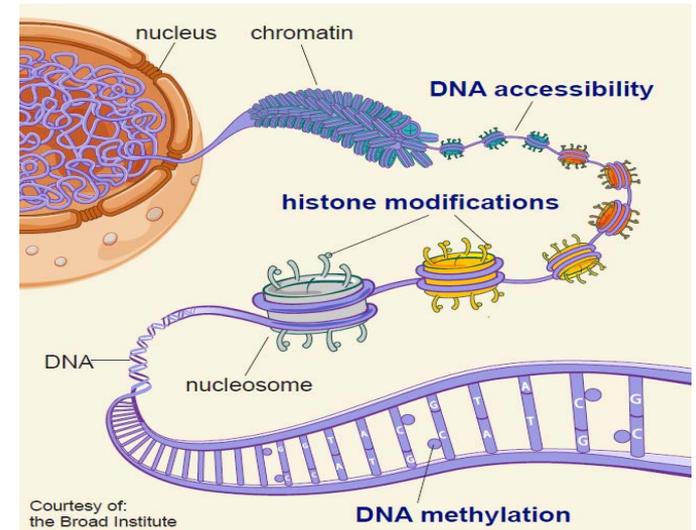
Epigenomic mapping across 100+ tissues/cell types

Diverse tissues and cells



X

Diverse epigenomic assays



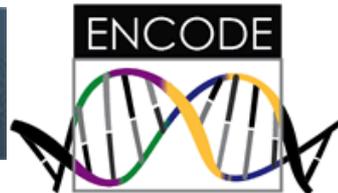
Courtesy of Broad Communications. Used with permission.

Adult tissues and cells (brain, muscle, heart, digestive, skin, adipose, lung, blood...)

Fetal tissues (brain, skeletal muscle, heart, digestive, lung, cord blood...)

ES cells, iPS, differentiated cells

(meso/endo/ectoderm, neural, mesench...)



Histone modifications

- H3K4me3, H3K4me1, H3K36me3
- H3K27me3, H3K9me3, H3K27/9ac
- +20 more

Open chromatin:

- DNA accessibility

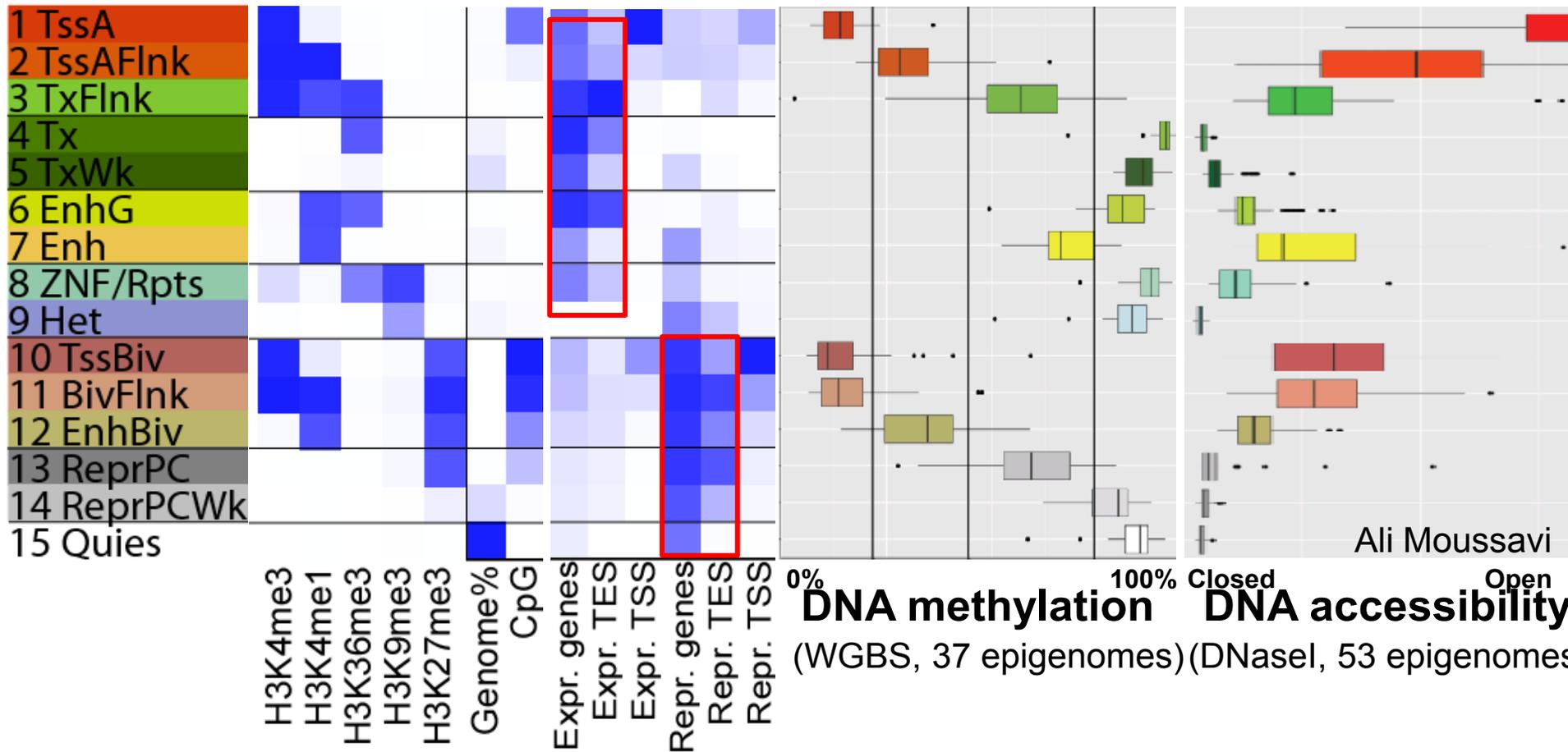
DNA methylation:

- WGBS, RRBS, MRE/MeDIP

Gene expression

- RNA-seq, Exon Arrays

States show distinct mCpG, DNase, Tx, Ac profiles



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

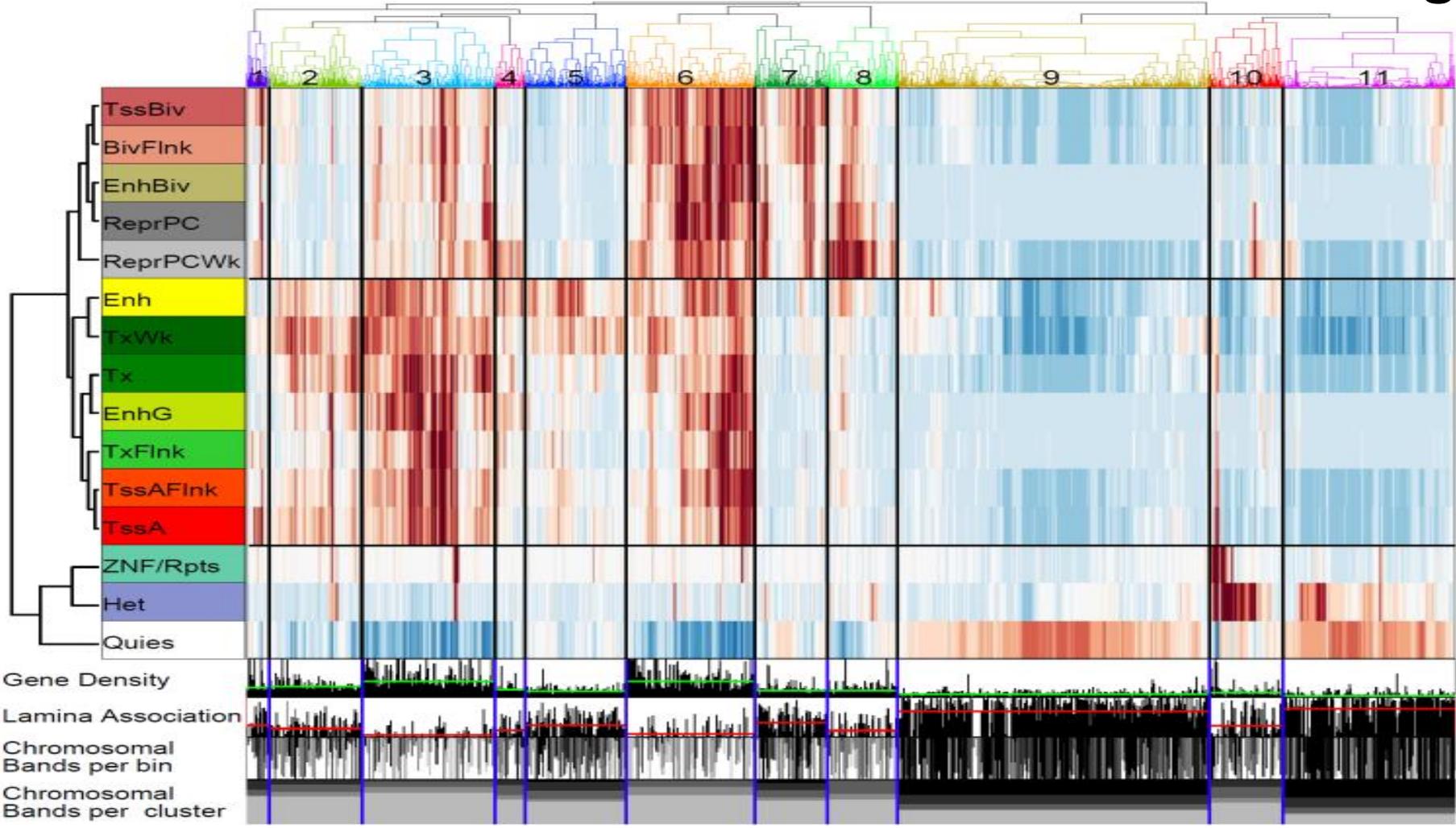
TssA vs. **TssBiv**: diff. activity, both open, both unmethylated!

Enh vs. **ReprPC**: diff. activity, both intermediate DNase/Methyl

Tx: Methylated, closed, actively transcribed

→ Distinct modes of repression: **H3K27me3** vs. **DNase** vs. **Het**

Chromosomal 'domains' from chromatin state usage

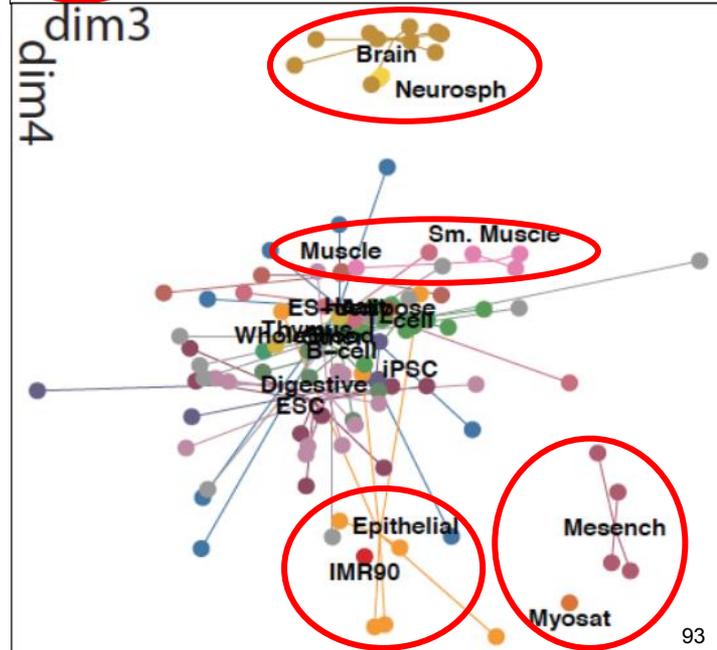
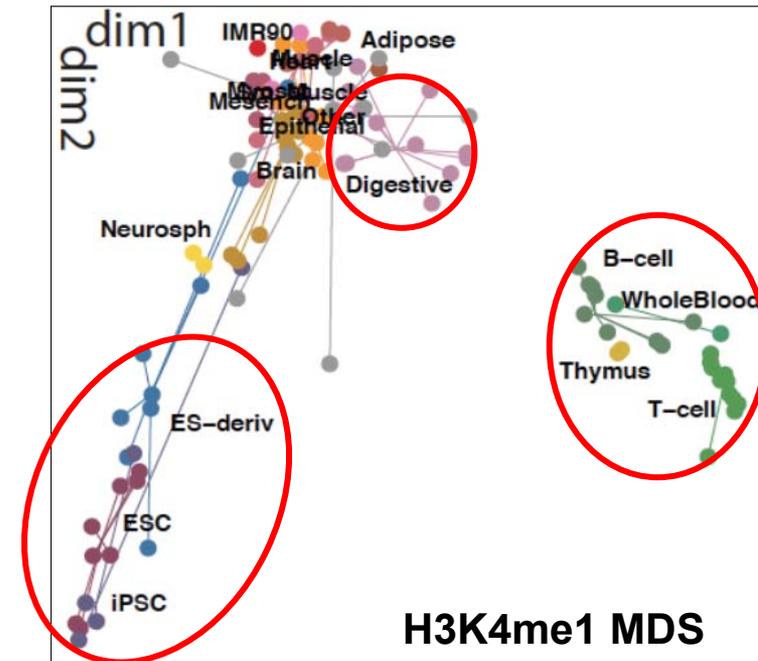
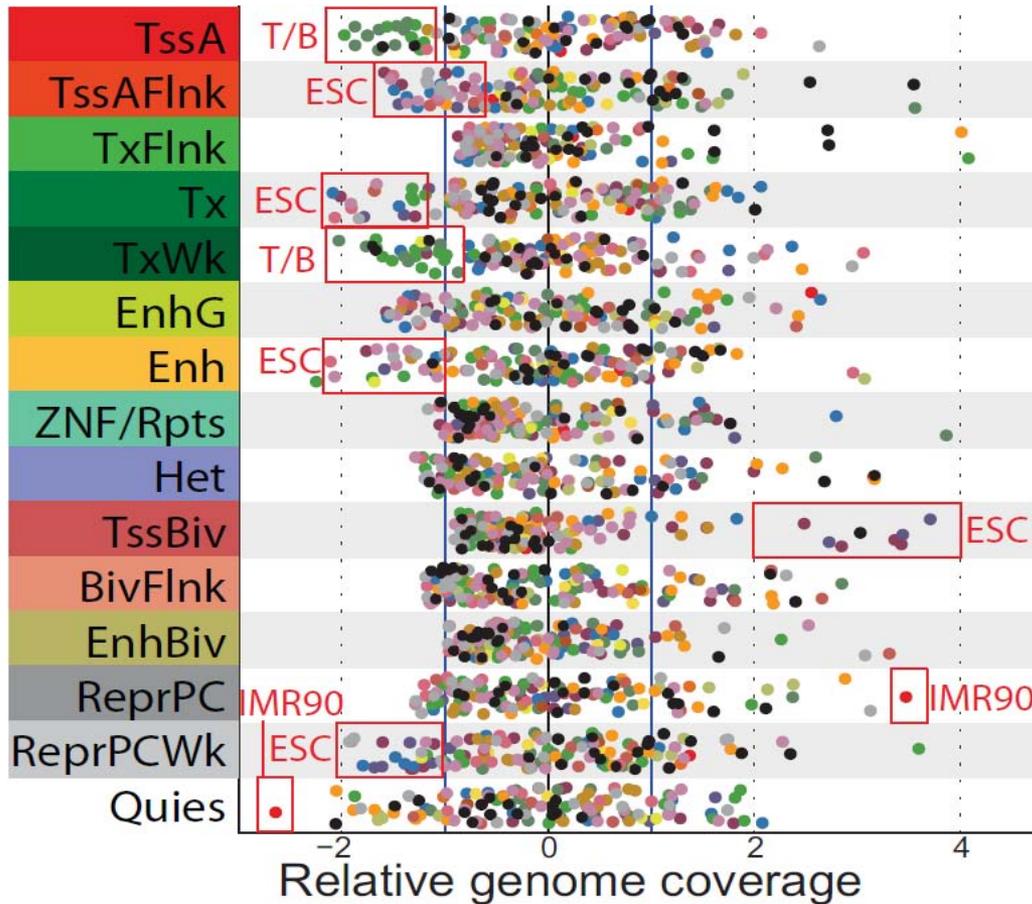


© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

Misha Bilenky

- State usage → gene density, lamina, cytogenetic bands
- Quies/ZNF/het | gene rich/poor, each active/repressed

Cells/Tissues at extremes of epigenomic variation



© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

- **ES/Immune/IMR90 most extreme**
- **ES: ↑Biv, ↓Enh/Tx/TssFlnk/PCwk**
- **Immune: ↓TssA, ↓TxWk**
- **IMR90: ↑ReprPC, ↓Quies**

Misha Bilenky, Wouter Meuleman

Chromatin state annotations across 127 epigenomes

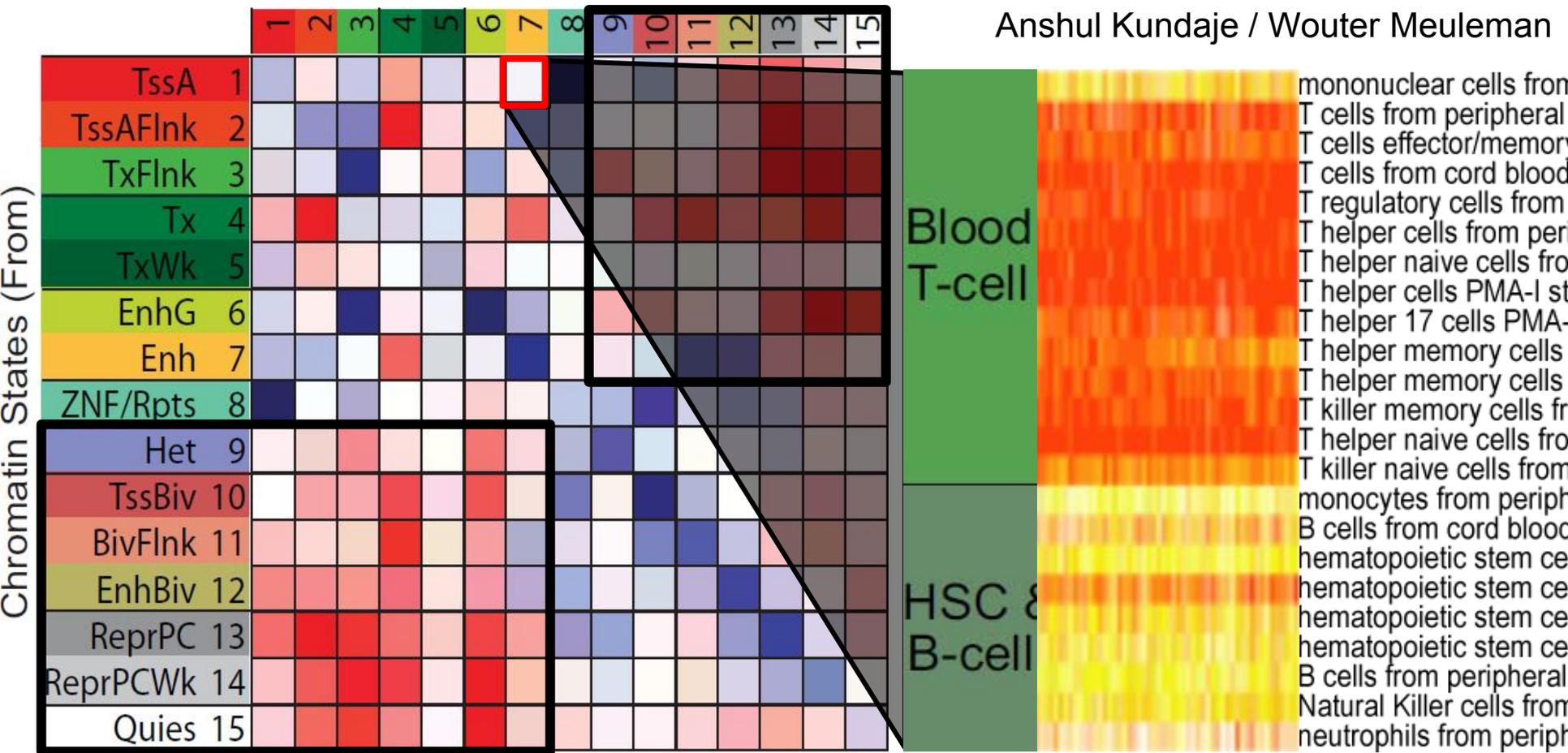


© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

Reveal epigenomic variability: enh/prom/tx/repr/het
 Anshul Kundaje 94

State switching: active/inactive, mostly keep identity

Anshul Kundaje / Wouter Meuleman

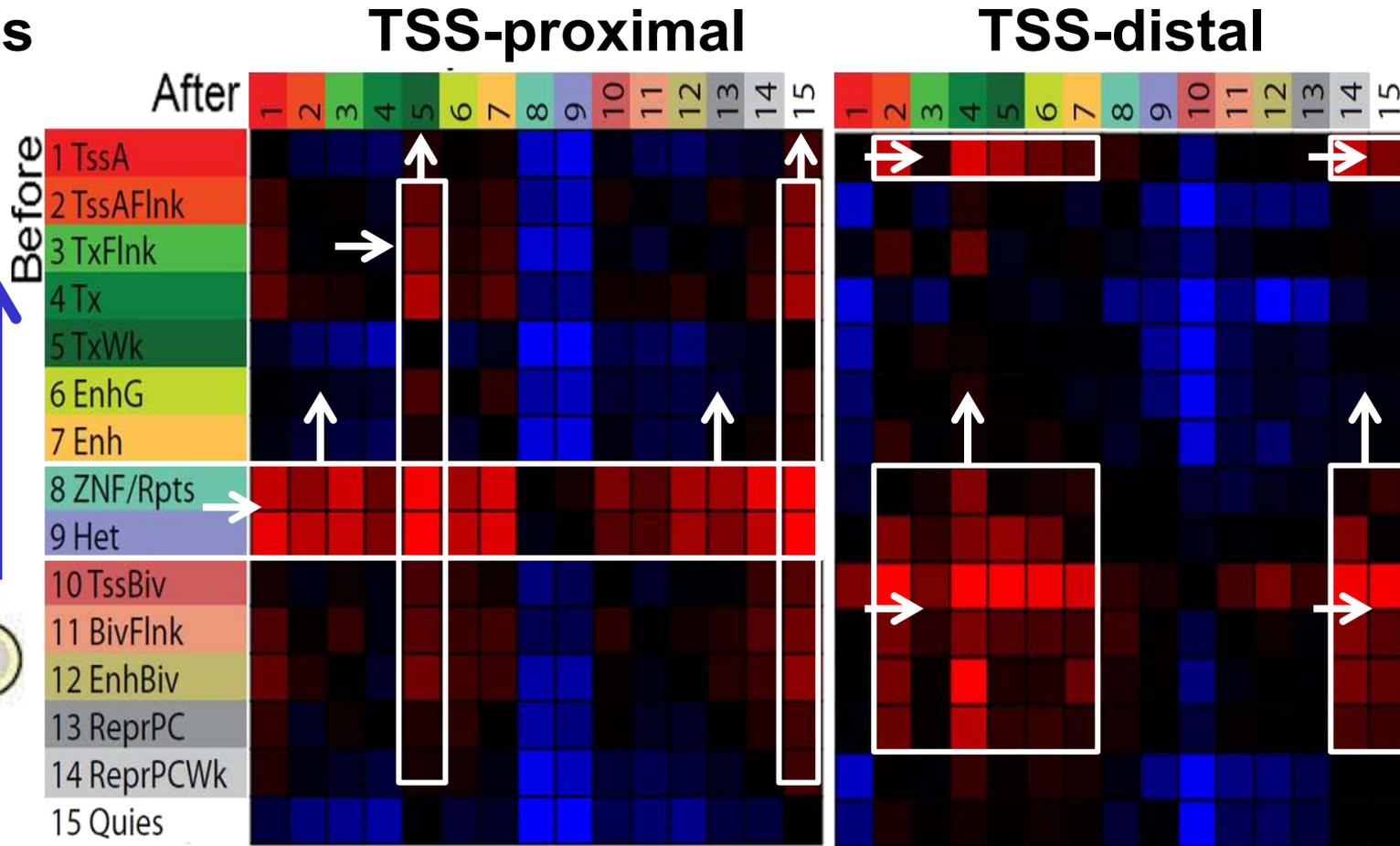
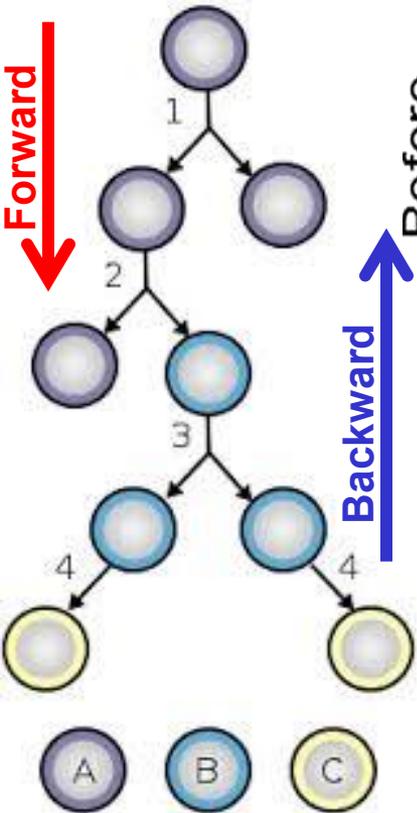


© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

- **Most variable: Enhancers. Least: TssA/Tx/Quies**
- **State switching: Active (1-7) ⇔ Inactive (10-15)**
- **Exception: Dyadic regions: enhancer ⇔ promoter**

Chromatin state changes during differentiation

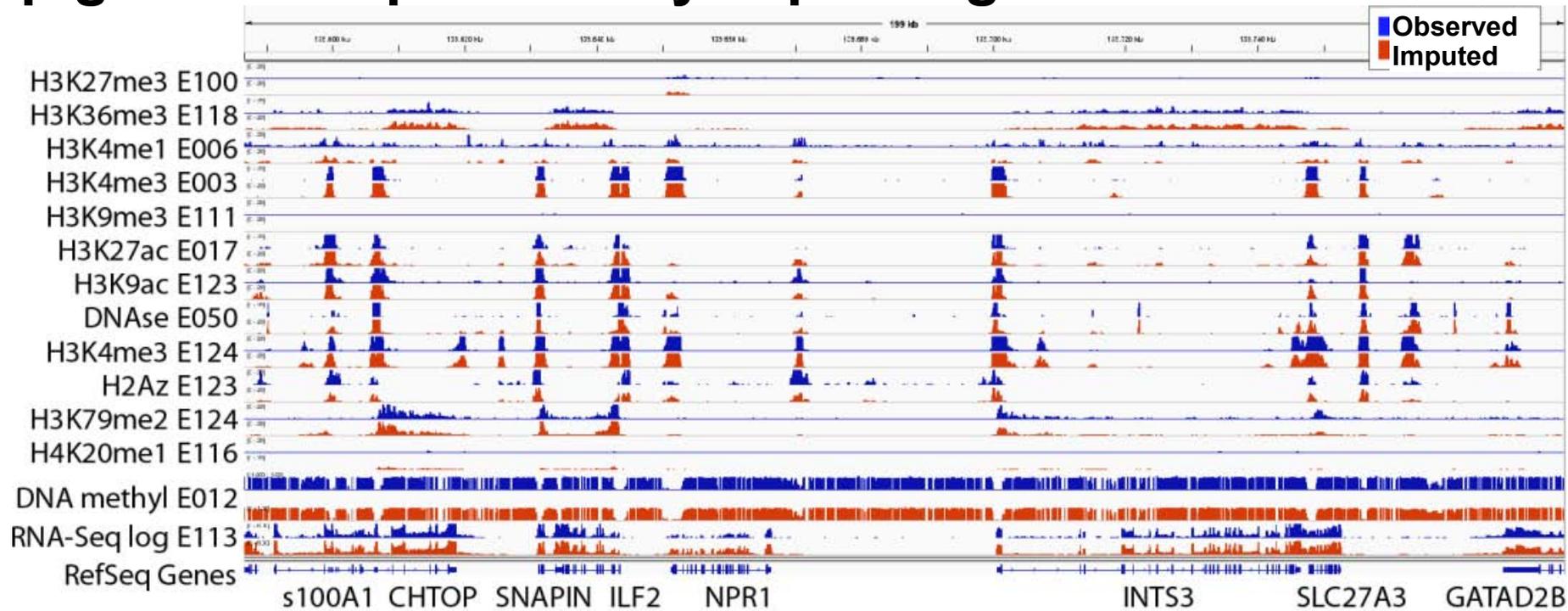
Classify cells



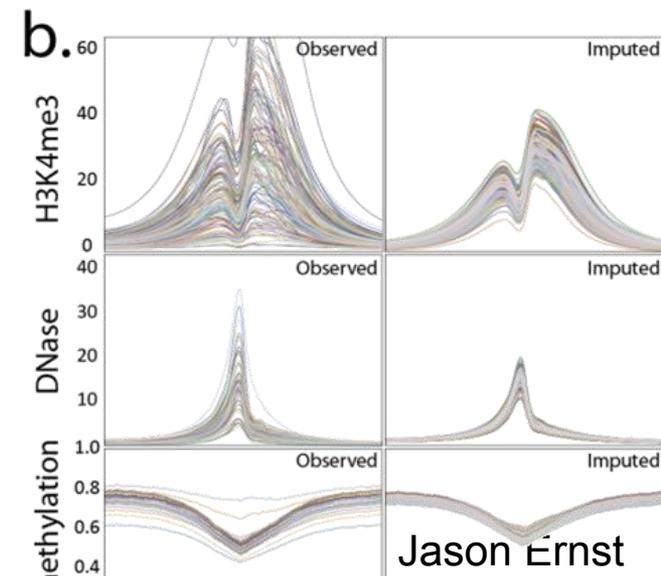
© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- **Epigenomic features can predict directionality: AUC 78%**
- **TSS-proximal:** (1) Loss of Het/ZNF. (2) Gain of TxWk, Quies
- **TSS-distal:** Bivalent, PCrepressed → Enhancer, Tx, TssFLnk

Epigenome imputation by exploiting mark correlations



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



- **Two types of features**
 - Other marks + context in same tissue
 - Same mark in 'closest' tissues
- **Impute missing datasets**
 - Predict DNase, marks @ 25bp res
 - Predict RNA-Seq @ 25 bp res
 - Predict DNA methylation @ 1bp res

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

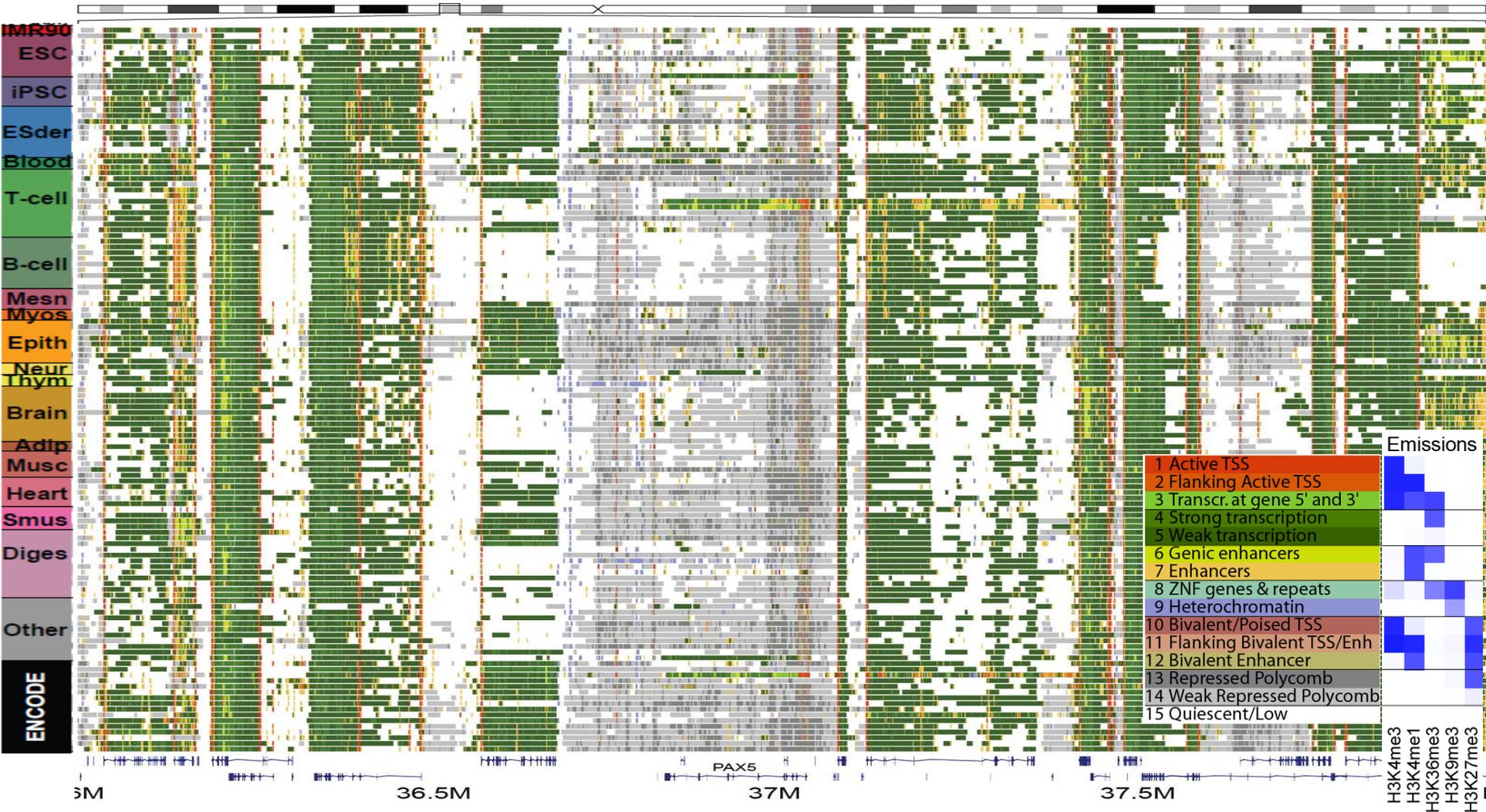
- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

5. Correlation-based links of enhancer networks

Regulators → Enhancers → Target genes

Chromatin state annotations across 127 epigenomes

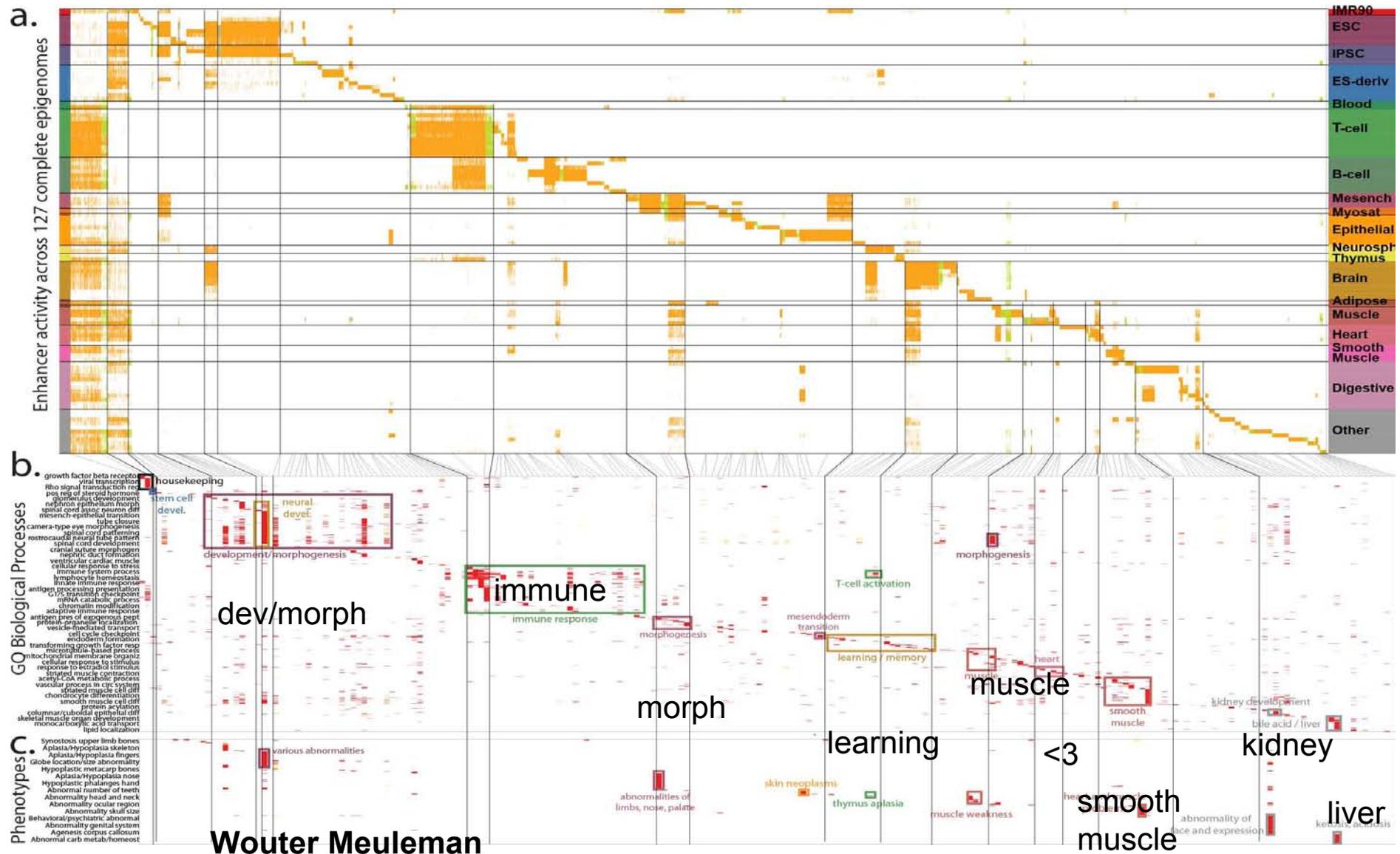


Courtesy of Macmillan Publishers Limited. Used with permission.

Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

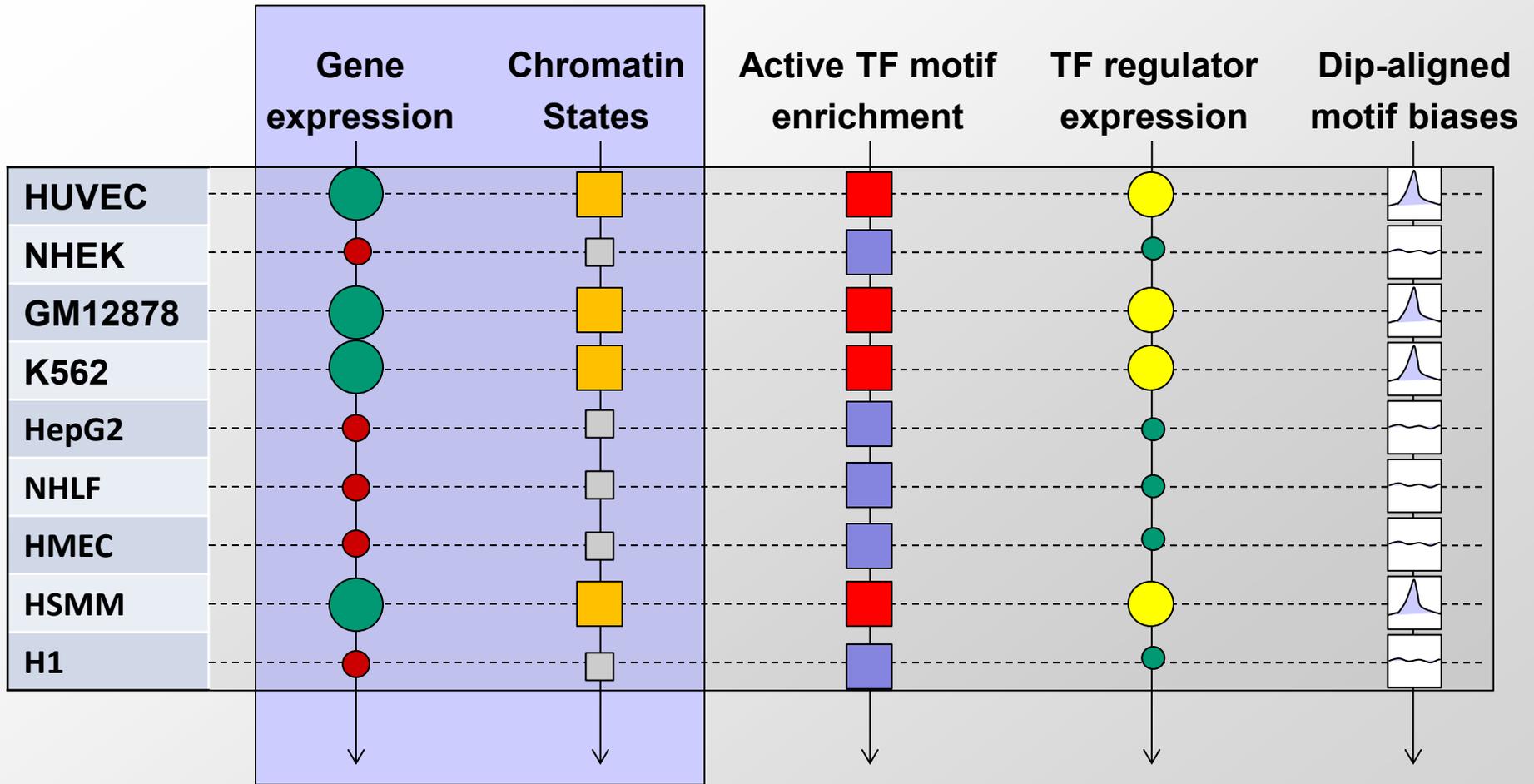
Reveal epigenomic variability: enh/prom/tx/repr/het

2.3M enhancer regions ⇔ only ~200 activity patterns

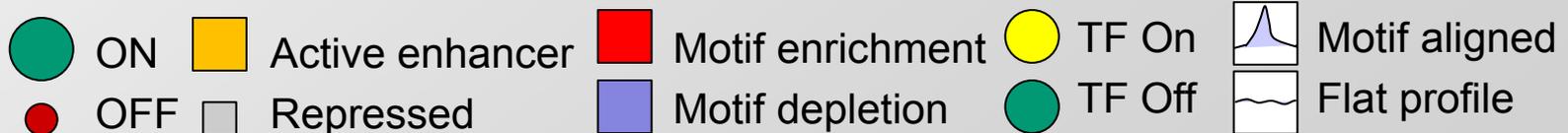


© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

Introducing multi-cell activity profiles

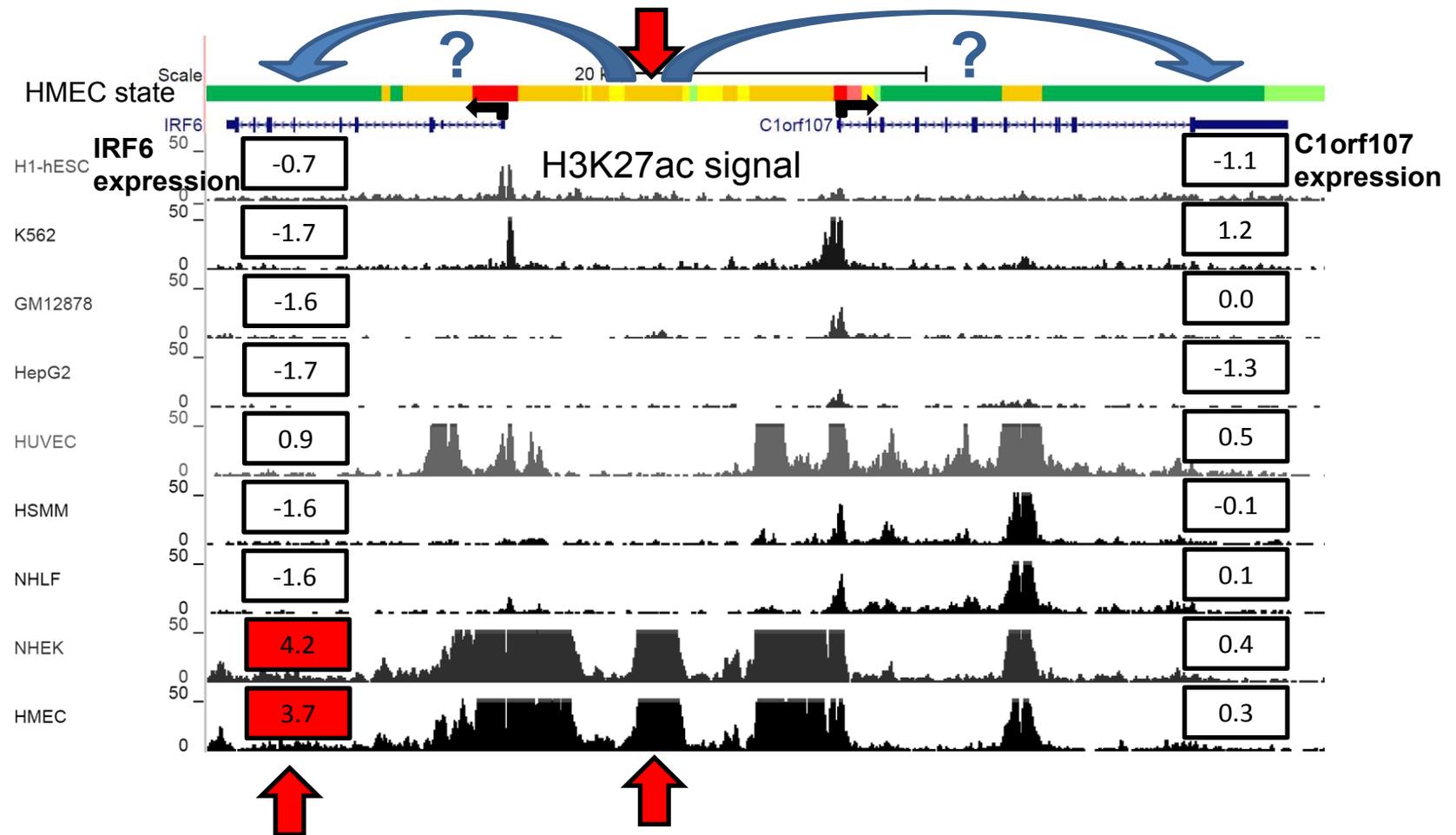


Link enhancers to target genes



Activity-based linking of enhancers to target genes

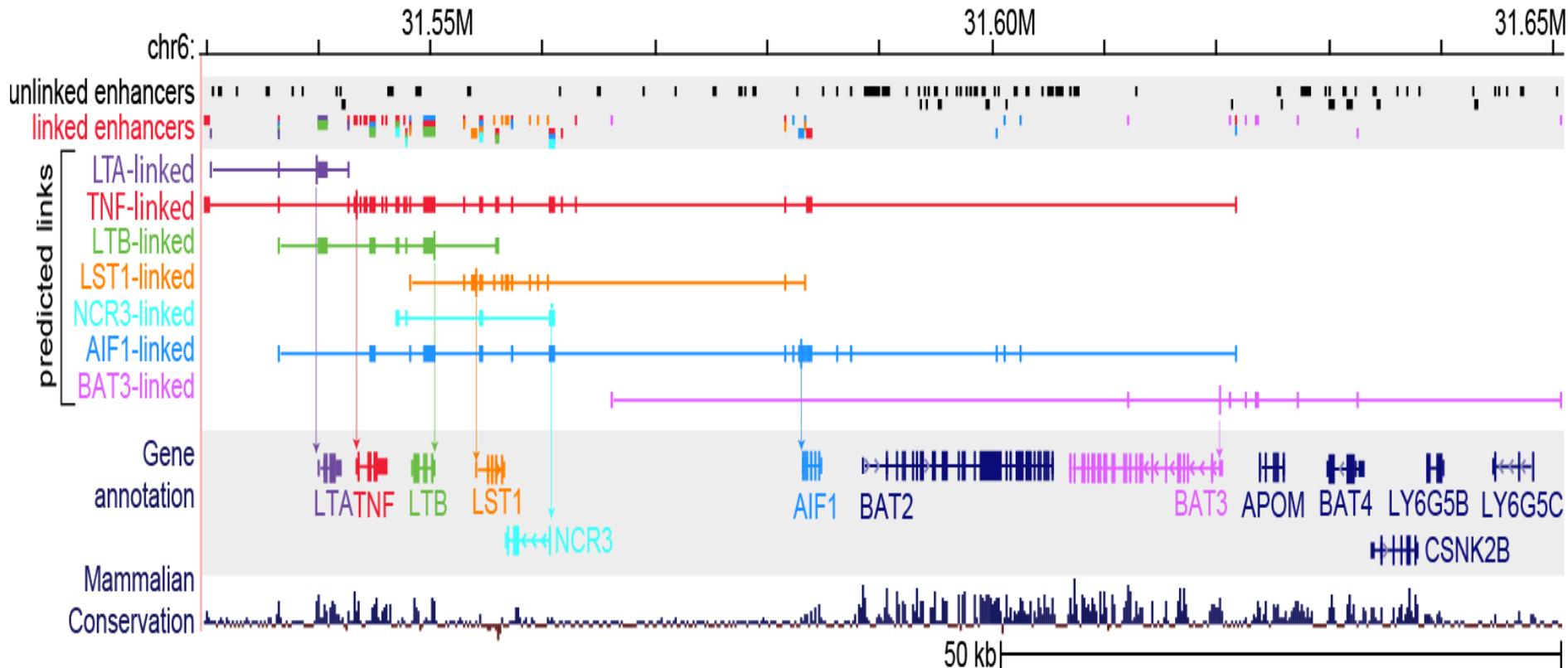
Finding correct target of enhancer in divergently transcribed genes



© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Compute correlations between gene expression levels and enhancer associated histone modification signals

Visualizing 10,000s predicted enhancer-gene links



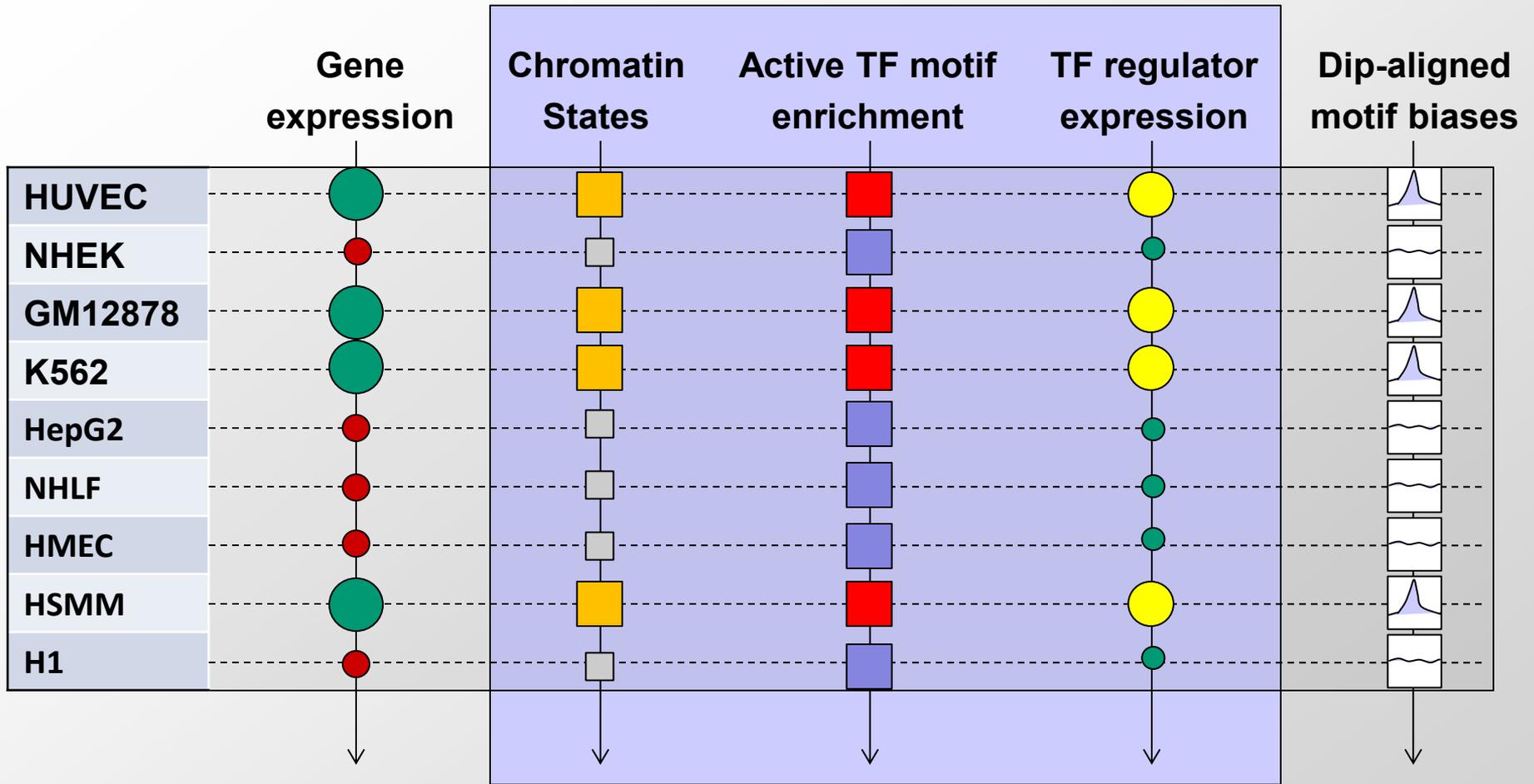
© Source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

- Overlapping regulatory units, both few and many
- Both upstream and downstream elements linked
- Enhancers correlate with sequence constraint

Chromatin dynamics: linking enhancer networks

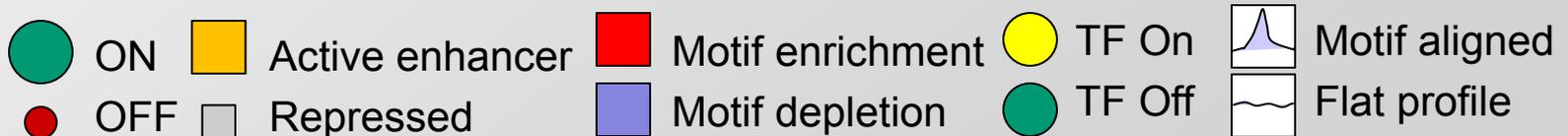
TFs → enhancers → target genes

Introducing multi-cell activity profiles



Link TFs to target enhancers

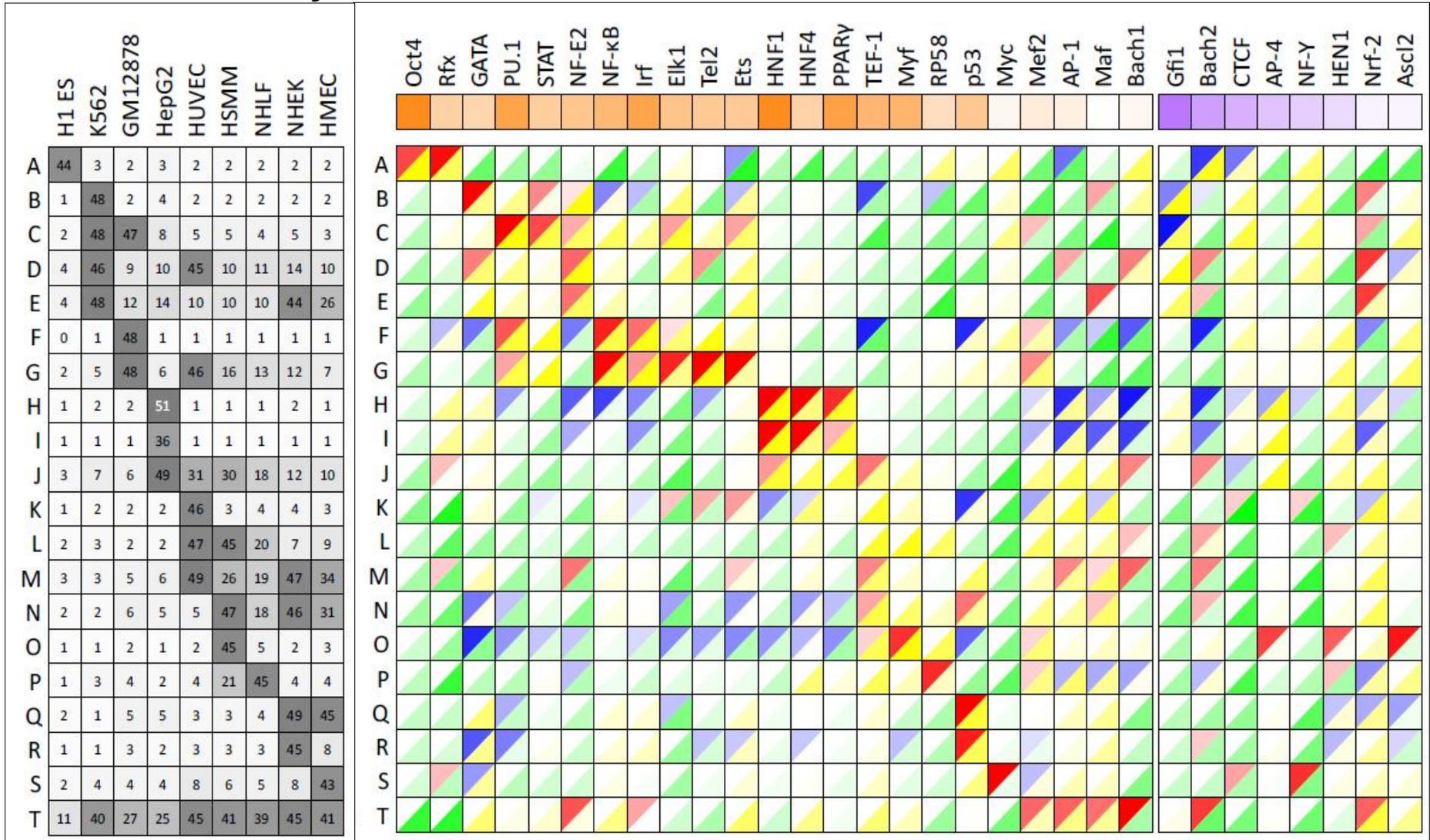
Predict activators vs. repressors



Coordinated activity reveals activators/repressors

Enhancer activity

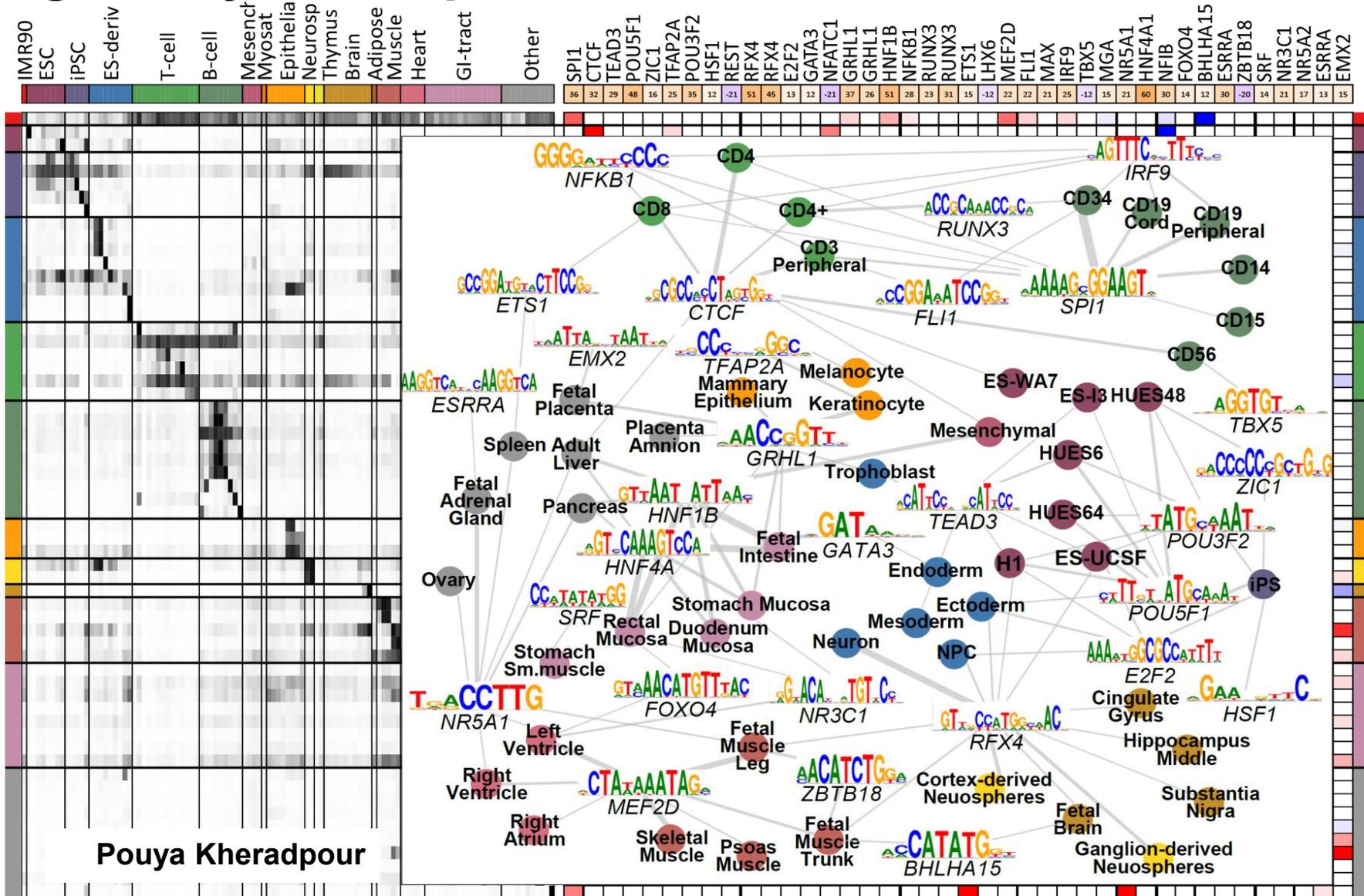
Activity signatures for each TF



Courtesy of Macmillan Publishers Limited. Used with permission.
 Source: Ernst, Jason et al. "Mapping and analysis of chromatin state dynamics in nine human cell types." Nature 473, no.7345 (2011): 43-49.

- Enhancer networks: Regulator → enhancer → target gene ¹⁰⁷

Regulatory motifs predicted to drive enhancer modules



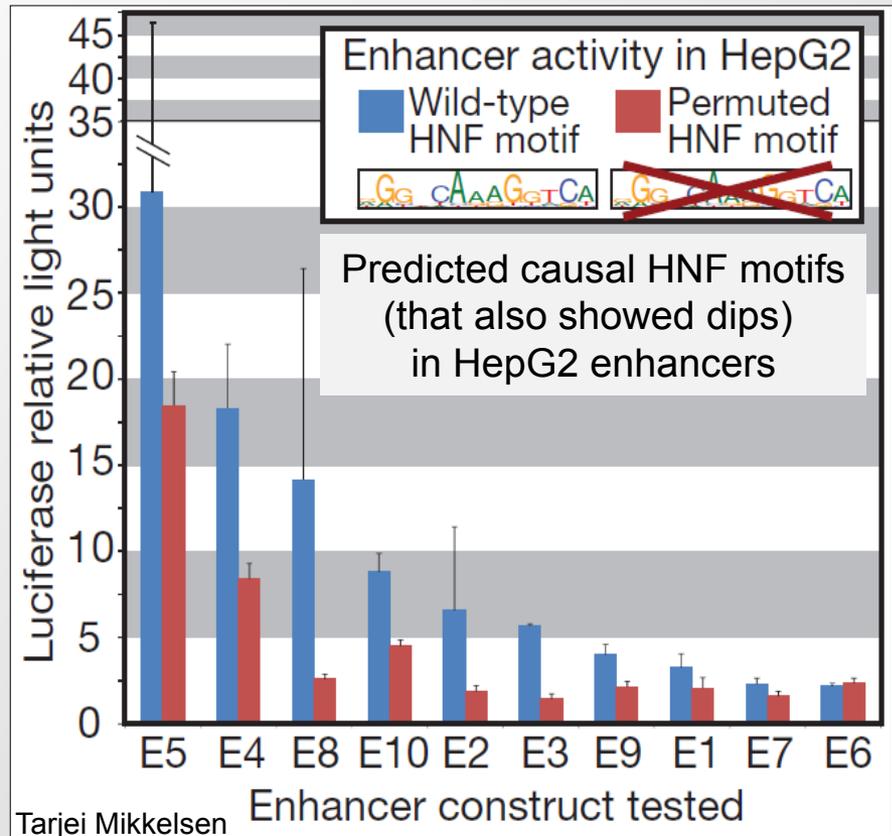
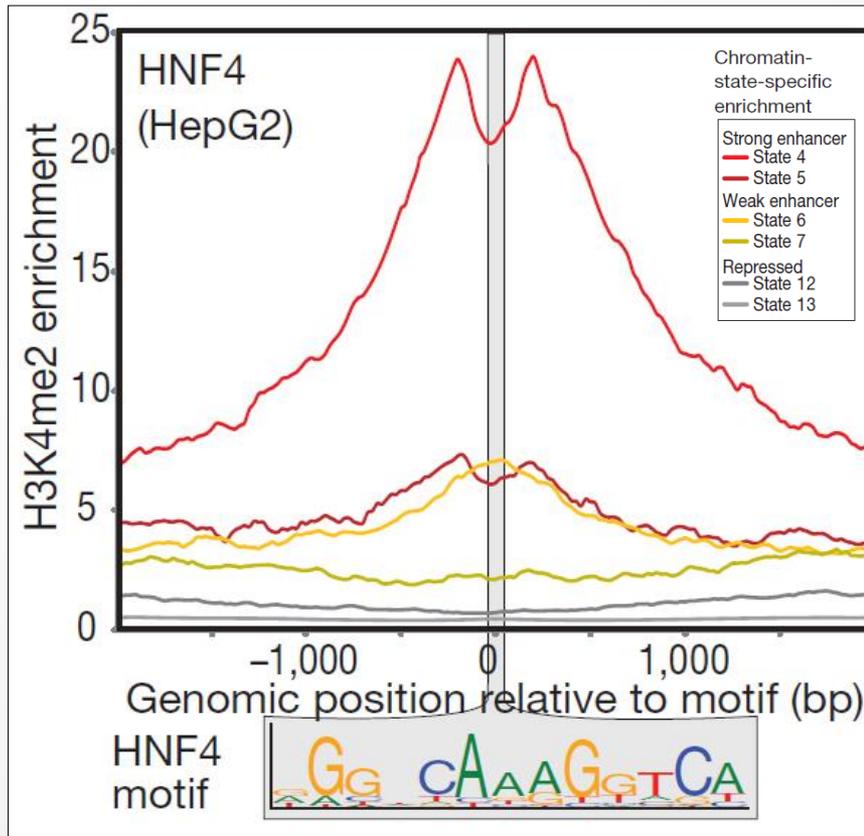
Pouya Kheradpour

© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

• Activator and repressor motifs consistent with tissues 108

Causal motifs supported by dips & enhancer assays



Courtesy of Macmillan Publishers Limited. Used with permission.
 Source: Ernst, Jason et al. "Mapping and analysis of chromatin state dynamics in nine human cell types." Nature 473, no.7345 (2011): 43-49.

Dip evidence of TF binding (nucleosome displacement)

Enhancer activity halved by single-motif disruption

➔ Motifs bound by TF, contribute to enhancers

Goals for today: Computational Epigenomics

1. Introduction to Epigenomics

- Overview of epigenomics, Diversity of Chromatin modifications
- Antibodies, ChIP-Seq, data generation projects, raw data

2. Primary data processing: Read mapping, Peak calling

- Read mapping: Hashing, Suffix Trees, Burrows-Wheeler Transform
- Quality Control, Cross-correlation, Peak calling, IDR (similar to FDR)

3. Discovery and characterization of chromatin states

- A multi-variate HMM for chromatin combinatorics
- Promoter, transcribed, intergenic, repressed, repetitive states

4. Model complexity: selecting the number of states/marks

- Selecting the number of states, selecting number of marks
- Capturing dependencies and state-conditional mark independence

5. Learning chromatin states jointly across multiple cell types

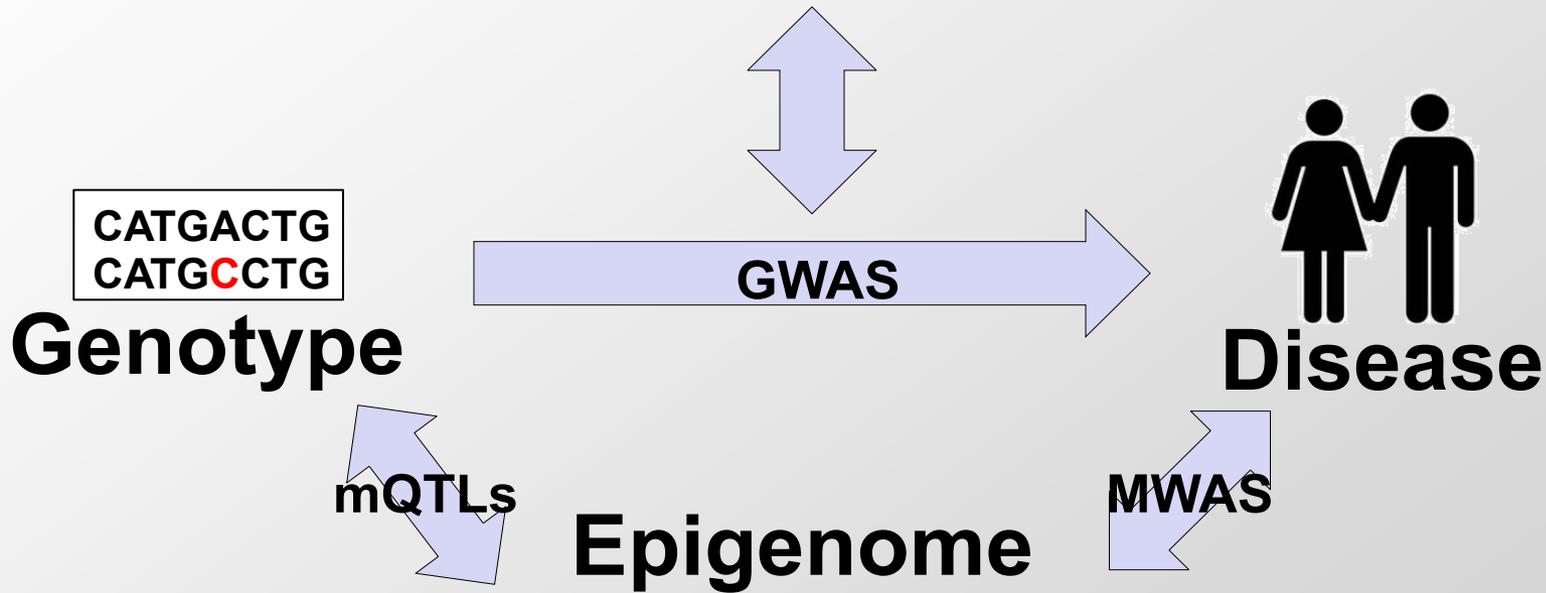
- Stacking vs. concatenation approach for joint multi-cell type learning
- Defining activity profiles for linking enhancer regulatory networks

(Future: Chromatin states to interpret disease-associated variants)

Interpreting disease-association signals

Interpret variants using reference states

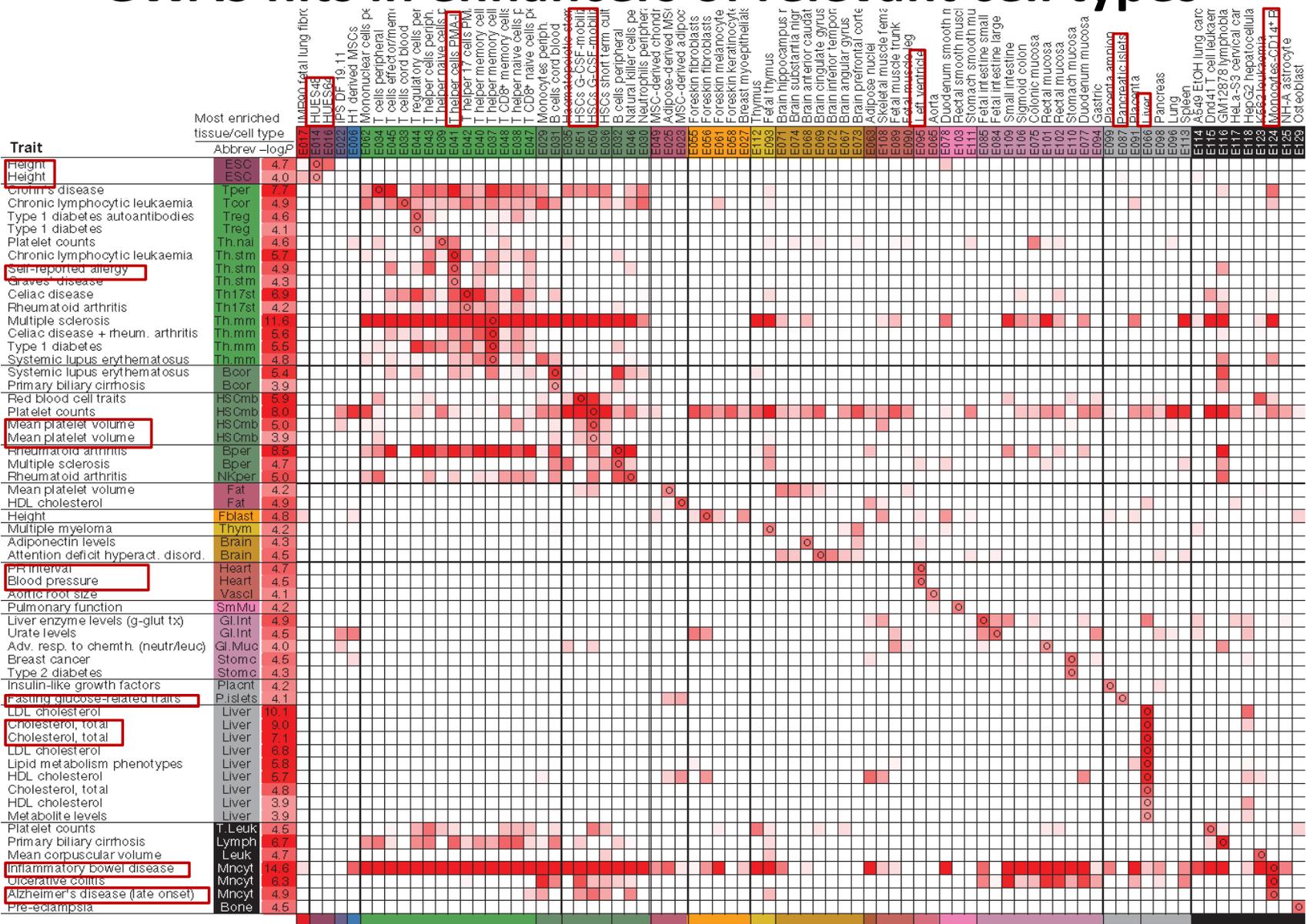
- Chromatin states: Enhancers, promoters, motifs
- Enrichment in individual loci, across 1000s of SNPs in T1D



Epigenome changes in disease

- Molecular phenotypic changes in patients vs. controls
- Small variation in brain methylomes, mostly genotype-driven
- 1000s of brain-specific enhancers increase methylation in Alzheimer's

GWAS hits in enhancers of relevant cell types



© Macmillan Publishers Limited. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.
 Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

HaploReg: systematic mining of GWAS variants

Query SNP: rs4684847 and variants with $r^2 \geq 0.8$

pos (hg19)	pos (hg38)	LD (r ²)	LD (D')	variant	Ref	Alt	AFR freq	AMR freq	ASH freq	EUR freq	SiPhy cons	Promoter histone marks	Enhancer histone marks	DNAse	Proteins bound	eQTL tissues	Motifs changed	Drivers disrupted	GENCODE genes	dbSNP func annot
chr3:12329783	chr3:12288284	0.95	0.97	rs17036160	C	T	0.01	0.06	0.04	0.12		24 organs	7 organs	4 organs			4 altered motifs		PPARG	intronic
chr3:12336507	chr3:12296008	0.95	0.97	rs11709077	G	A	0.01	0.07	0.04	0.12		LNG	9 organs	15 organs			4 altered motifs		PPARG	intronic
chr3:12344730	chr3:12303231	0.94	0.97	rs11712037	C	G	0.01	0.08	0.04	0.12			8 organs	BLD			AP-1,TCF11::MafG		PPARG	intronic
chr3:12351521	chr3:12310022	0.95	0.97	rs36000407	T	G	0.01	0.07	0.04	0.12		LNG	5 organs					Smad	PPARG	intronic
chr3:12360884	chr3:12319385	0.95	0.97	rs150732434	TG	T	0.01	0.07	0.04	0.12		FAT	7 organs	MUS,VAS	CFOS		Hdx,Sox,TATA		PPARG	intronic
chr3:12365308	chr3:12323809	0.95	0.97	rs13083375	G	T	0.01	0.07	0.04	0.12		BLD	BLD, FAT				Homez,Sox,YY1		PPARG	intronic
chr3:12369401	chr3:12327902	0.95	0.97	rs13064780	C	T	0.01	0.07	0.04	0.12			7 organs						PPARG	intronic
chr3:12375956	chr3:12334457	0.95	0.97	rs2012444	C	T	0.01	0.07	0.04	0.12			SKIN, FAT, BLD						PPARG	intronic
chr3:12383285	chr3:12341786	0.96	0.99	rs13085211	G	A	0.18	0.10	0.04	0.12			FAT, SKIN						PPARG	intronic
chr3:12383714	chr3:12342215	0.96	0.99	rs7638903	G	A	0.18	0.10	0.04	0.12			8 organs	CRVX					PPARG	intronic
chr3:12385828	chr3:12344329	0.95	1	rs11128603	A	G	0.18	0.10	0.04	0.12			CRVX						PPARG	intronic
chr3:12386337	chr3:12344838	1	1	rs4684847	C	T	0.01	0.07	0.04	0.12			8 organs						PPARG	intronic
chr3:12388409	chr3:12346910	0.99	1	rs7610055	G	A	0.17	0.09	0.04	0.12			BLD						PPARG	intronic
chr3:12389313	chr3:12347814	0.99	1	rs17036326	A	G	0.17	0.09	0.04	0.12			FAT, BL						PPARG	intronic
chr3:12390484	chr3:12348985	0.99	1	rs17036328	T	C	0.17	0.09	0.04	0.12			FAT, CR						PPARG	intronic
chr3:12391207	chr3:12349708	0.99	1	rs6802898	C	T	0.61	0.15	0.04	0.12			FAT, BL						PPARG	intronic
chr3:12391583	chr3:12350084	0.99	1	rs2197423	G	A	0.17	0.09	0.04	0.12			FAT, LIV						PPARG	intronic
chr3:12391813	chr3:12350314	0.99	1	rs7647481	G	A	0.17	0.09	0.04	0.12			4 organs						PPARG	intronic
chr3:12392272	chr3:12350773	0.99	1	rs7649970	C	T	0.17	0.09	0.04	0.12			5 organs						PPARG	intronic
chr3:12393125	chr3:12351626	1	1	rs1801282	C	G	0.01	0.07	0.04	0.12			FAT, LIV						PPARG	missense
chr3:12393682	chr3:12352183	0.99	1	rs17036342	A	G	0.17	0.09	0.04	0.12			FAT						PPARG	intronic
chr3:12394840	chr3:12353341	0.99	1	rs1899951	C	T	0.61	0.15	0.04	0.12			FAT						PPARG	intronic
chr3:12395645	chr3:12354146	0.99	1	rs4684848	G	A	0.61	0.15	0.04	0.12			FAT, BLD						PPARG	intronic
chr3:12396845	chr3:12355346	0.93	1	rs4135250	A	G	0.17	0.09	0.04	0.13			4 organs	ADRL,GI,CRVX	5 bound proteins				PPARG	intronic
chr3:12396913	chr3:12355414	0.98	1	rs71304101	G	A	0.01	0.07	0.04	0.12			4 organs	PLCNT					PPARG	intronic
chr3:12396955	chr3:12355456	0.96	1	rs2881654	G	A	0.61	0.15	0.04	0.12			4 organs	PLCNT					PPARG	intronic

Courtesy of the authors. License: CC BY-NC.

Source: Ward, Lucas D. and Manolis Kellis. "HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants." Nucleic Acids Research 40, no. D1 (2012): D930-D934.

- **Start with any list of SNPs or select a GWA study**
 - Mine ENCODE and Roadmap epigenomics data for hits
 - Hundreds of assays, dozens of cells, conservation, motifs
 - Report significant overlaps and link to info/browser
- **Try it out: <http://compbio.mit.edu/HaploReg>** Ward, Kellis NAR 2011₁₁₄

MIT OpenCourseWare
<http://ocw.mit.edu>

6.047 / 6.878 / HST.507 Computational Biology
Fall 2015

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.