

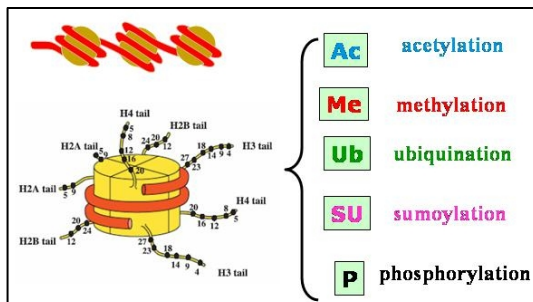
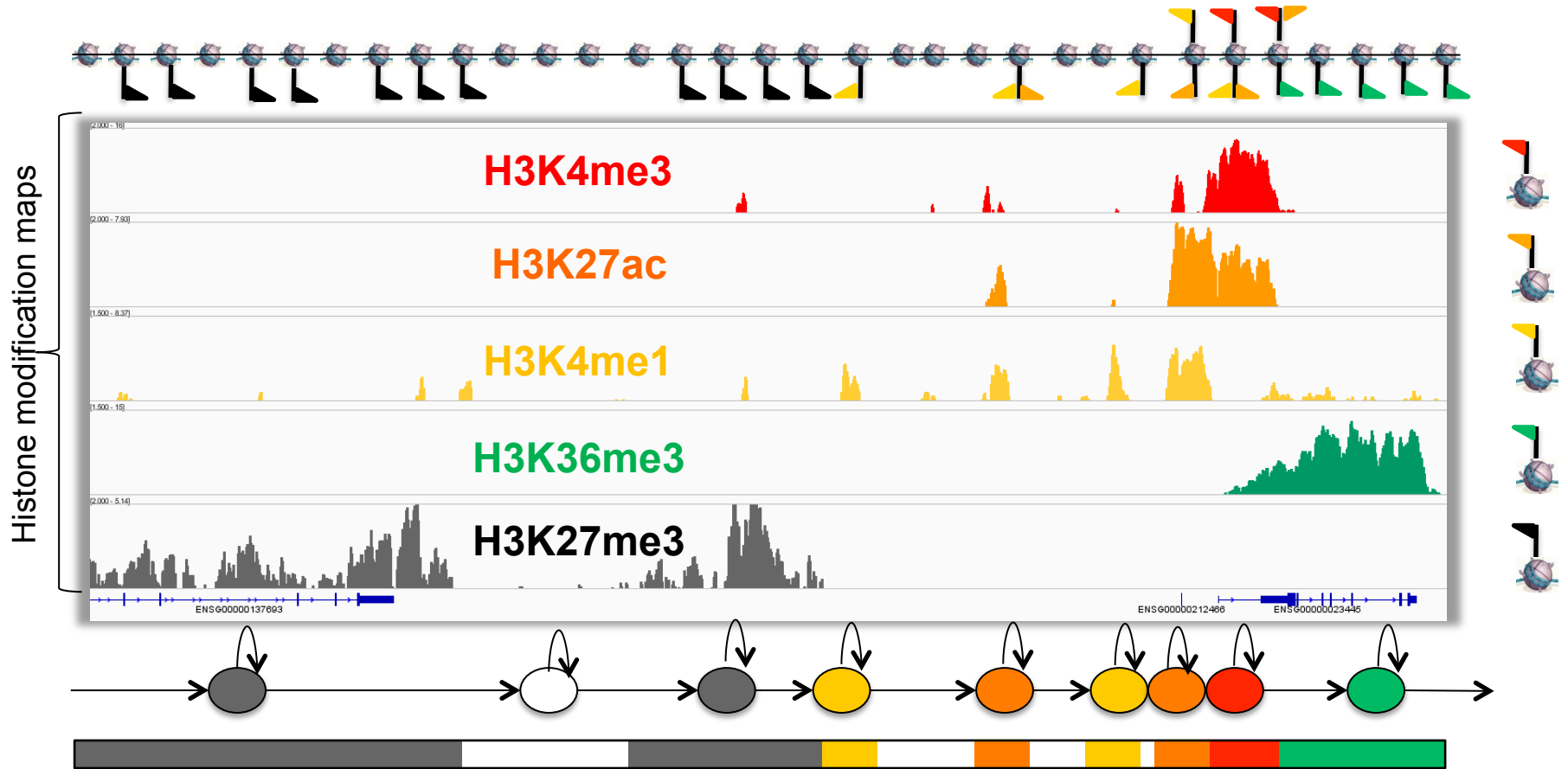
ENCODE3 chromatin state calls based on
Roadmap Epigenomics ChromHMM models

WM20151112

Goal: chromatin state calls for ENCODE3 data

- The latest ENCODE3 ChIP-seq data freeze contained a limited amount of new data. Instead of building new chromatin state models on these data, we may be able to use existing models.
- The Roadmap Epigenomics Project has released chromatin state models across 127 epigenomes, spanning 100+ cell types.
- These models can be applied to incoming ENCODE3 data, provided that these data are processed identically to what was used for Roadmap.
- Benefits: validate ENCODE3 cell types against large body of reference cell types, fairly rapid general QC, identify novel regions.

Learning genome-wide chromatin state maps



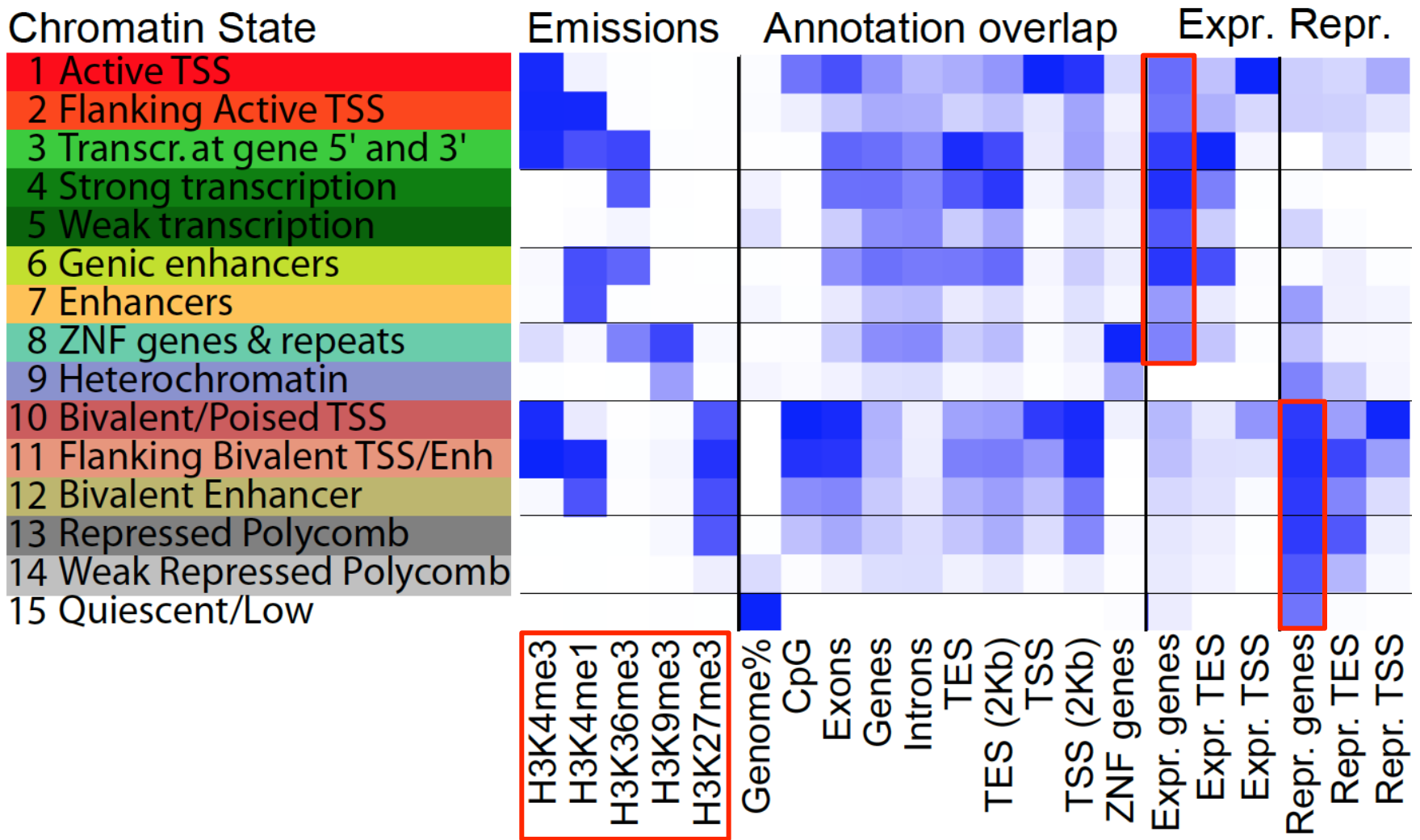
Chromatin States: Hidden Markov model

ChromHMM: automating chromatin-state discovery and characterization

Jason Ernst & Manolis Kellis

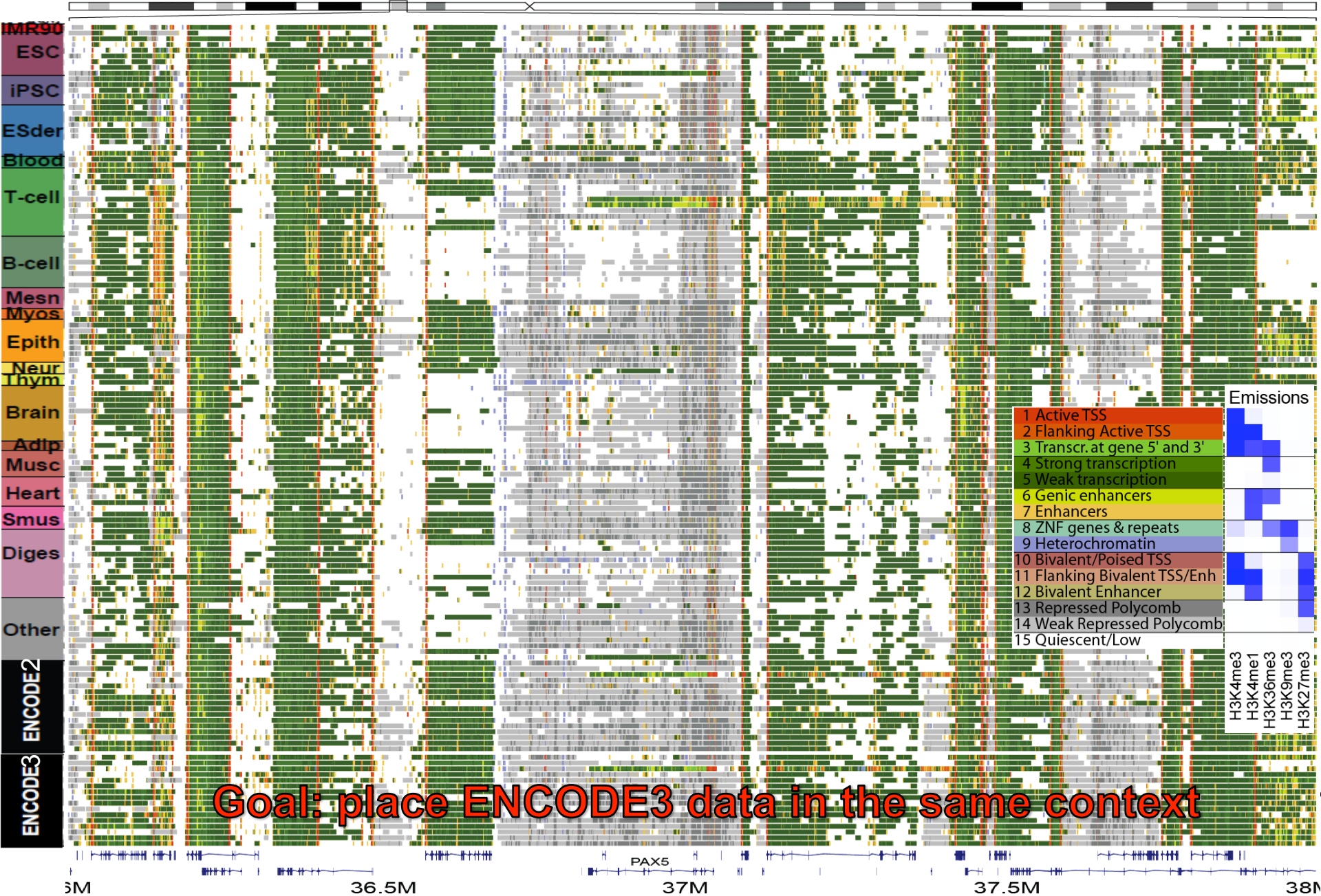
Nature Methods **9**, 215–216 (2012) | doi:10.1038/nmeth.1906

The 5 core histone modification marks have been used to build a ChromHMM model with 15 chromatin states



States make a clear distinction based on transcriptional activity

Chromatin state annotations across 127 epigenomes



ENCODE3 human ChIP-seq data on the DCC portal

Summary for project file set ENCSR301PBV

Status: released

Assay(s): ChIP-seq
Accession: ENCSR301PBV
Description: Uniformly processed ENCODE3 human histone ChIP-seq October 2015
Project type: Evaluation
Biosample term name: Loucy, DOHH2, OCI-LY7 ← 3 cell types only?
Biosample type: immortalized cell line
Organism: human
Lab: Zhiping Weng, UMass
External resources: None submitted

Files in project file set ENCSR301PBV

Visualize Data

Raw data

Accession	Originating dataset	File type	Biological replicate	Technical replicate	Read length	Run type	Paired end	Mapping assembly	Lab	Date added
ENCF002ASS Download 1.14 GB	ENCSR623GZY	fastq	2	1	36 nt	single-ended			Bradley Bernstein, Broad	2014-03-19
ENCF002ATI Download 1.1 GB	ENCSR005SXO	fastq	1	1	36 nt	single-ended			Bradley Bernstein, Broad	2014-03-19
ENCF002ATK Download 1 GB	ENCSR447ZGY	fastq	2	1	36 nt	single-ended			Bradley Bernstein, Broad	2014-03-19

etcetera

<https://www.encodeproject.org/projects/ENCSR301PBV/>

Navigated to via:

<https://www.encodeproject.org/search/?type=Project&lab.title=Zhiping+Weng%2C+UMass&status=released>

ENCODE3 chromatin state calls

- Current ENCODE3 data I'm working with:
 - 19 cell types (next slide)
 - 6 epitopes (K4me1/3, K27me3/ac, K9me3, K36me3) + WCE
- Strategy:
 - process all mapped sequencing reads of ENCODE3 data in a similar manner as done for Roadmap Epigenomics.
 - Remove duplicates & pool data across replicates.
 - Aim for at least 30M reads per cell type / epitope combination. If more, subsample to 30M.
 - Call chromatin states for processed ENCODE3 data, using Roadmap Epigenomics models.

QC: number of (non-duplicate) reads in millions

Broad pipeline

Roadmap/ENCODE3 pipeline

56	116	63	61	57	66	80	Fibroblast-of-arm
102	56	95	98	93	101	72	PC9
76	69	65	61	66	98	252	H1-derived-neurons
96	68	81	89	83	79	135	MM.1S-Myeloma-cell-line
34	42	55	55	69	135	97	Neuroectoderm
87	80	79	92	85	77	104	PC3-prostate
29	40	44	61	59	50	82	Epiblast
55	83	97	66	41	38	76	Radial-Glia-VZ
47	62	58	55	84	35	86	MCF-7
61	48	56	57	42	58	116	DOHH2
14	62	75	77	45	71	68	Neuroepithelial
135	39	83	40	43	39	43	SuDHL-6
49	55	49	54	44	111	43	Oci-Ly-3
83	40	50	47	53	56	61	Karpas-422
48	42	38	36	16	42	145	Oci-Ly-7
39	45	40	41	69	50	62	LOUCY
33	73	32	33	31	30	48	Oci-Ly-1
28	29	37	37	32	21	103	KOPT-K1-KOPzero
8	15	8	6	16	11	77	A673
H3K27me3	H3K9me3	H3K36me3	H3K4me1	H3K27ac	H3K4me3	WCE	

93	177	99	104	101	112	126	Fibroblast-of-arm
116	100	107	112	113	116	56	PC9
60	79	50	50	54	87	232	H1-derived-neurons
75	68	78	70	66	62	124	MM.1S-Myeloma-cell-line
25	30	57	44	120	164	76	Neuroectoderm
69	61	62	77	72	63	84	PC3-prostate
39	44	56	69	63	61	104	Epiblast
58	61	75	78	47	29	58	Radial-Glia-VZ
37	46	44	45	81	28	90	MCF-7
49	35	44	46	34	48	89	DOHH2
9	42	56	62	58	57	52	Neuroepithelial
105	28	64	31	33	29	33	SuDHL-6
38	40	37	44	36	91	33	Oci-Ly-3
63	28	37	38	43	43	48	Karpas-422
35	31	29	30	14	37	112	Oci-Ly-7
28	32	31	34	55	40	49	LOUCY
26	53	24	28	26	27	37	Oci-Ly-1
22	16	27	29	26	18	79	KOPT-K1-KOPzero
7	9	7	7	15	15	63	A673
H3K27me3	H3K9me3	H3K36me3	H3K4me1	H3K27ac	H3K4me3	WCE	

Processing pipeline

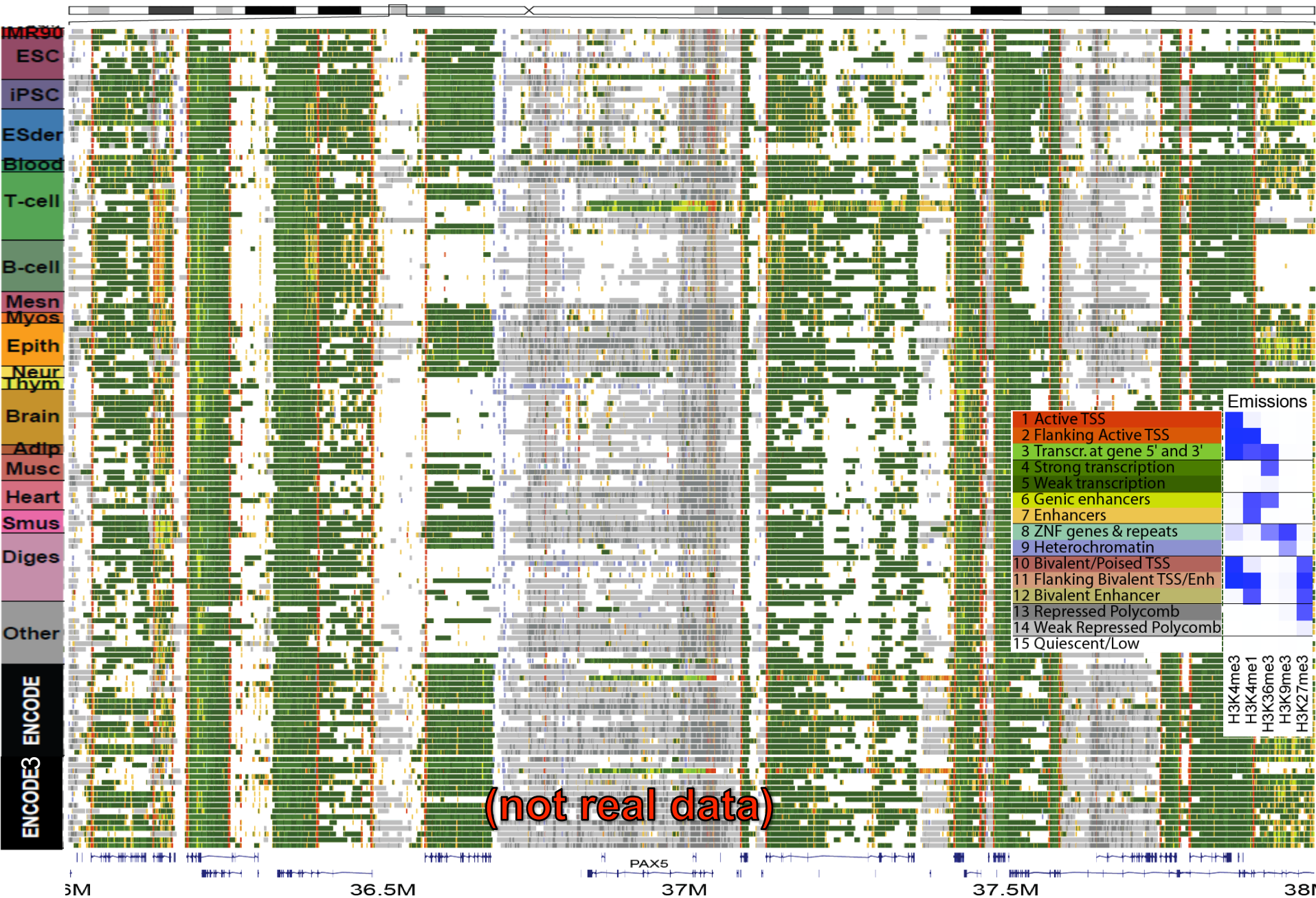
Largely based on code from Anshul and the (draft?) ENCODE3 processing pipeline*

Processing steps:

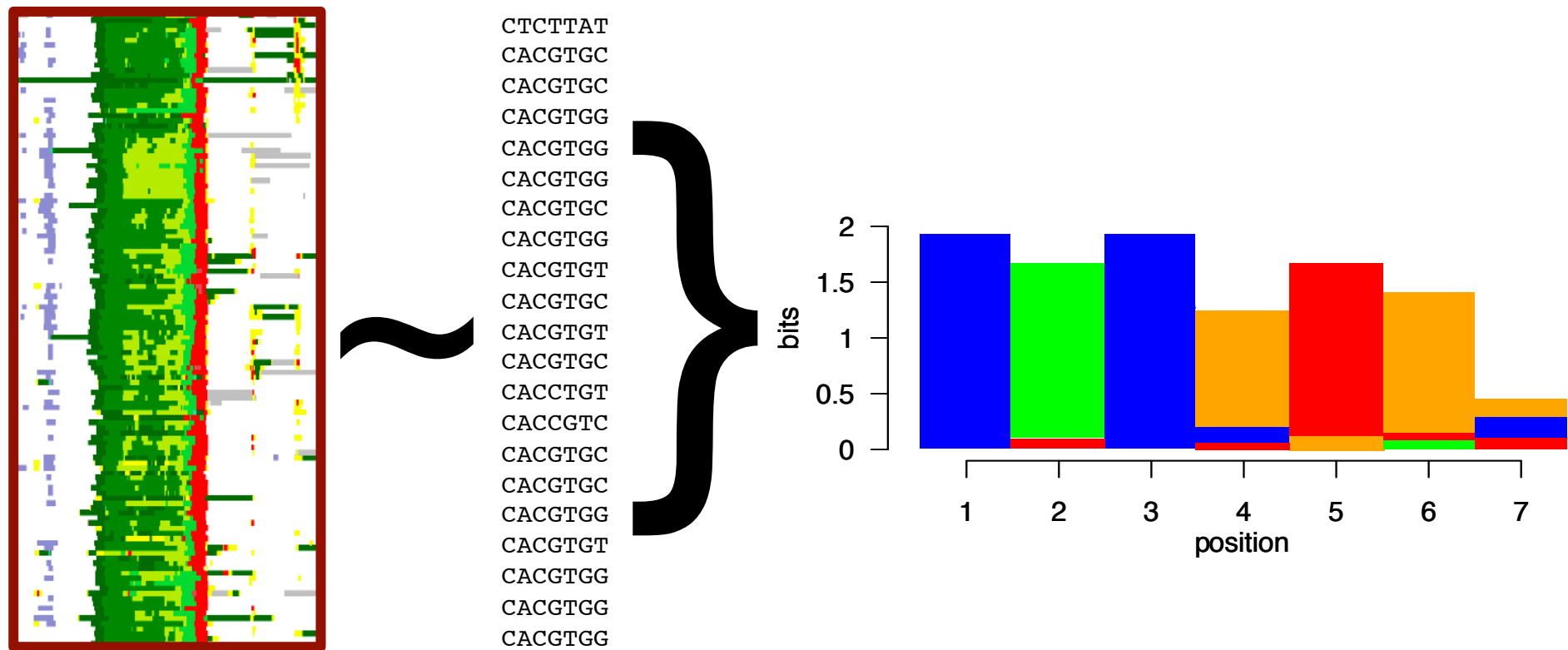
1. Remove un/mis-mapped, low quality reads, mark & remove duplicates, compute library complexity QC, filter for mappability
2. Pool reads across replicates and (if needed) sub-sample to 30M reads, use phantomPeakQualTools to obtain QC info
3. Binarize data using ChromHMM BinarizeBed, create segmentation based on existing Roadmap models and output states per 200bp bins

*: https://docs.google.com/document/d/1IG_Rd7fnYgRpSlqrlfuVIAz2dW1VaSQThzk836Db99c/edit#heading=h.9ecc41kilcvq

Bonus: **epi**logos for ENCODE3 data

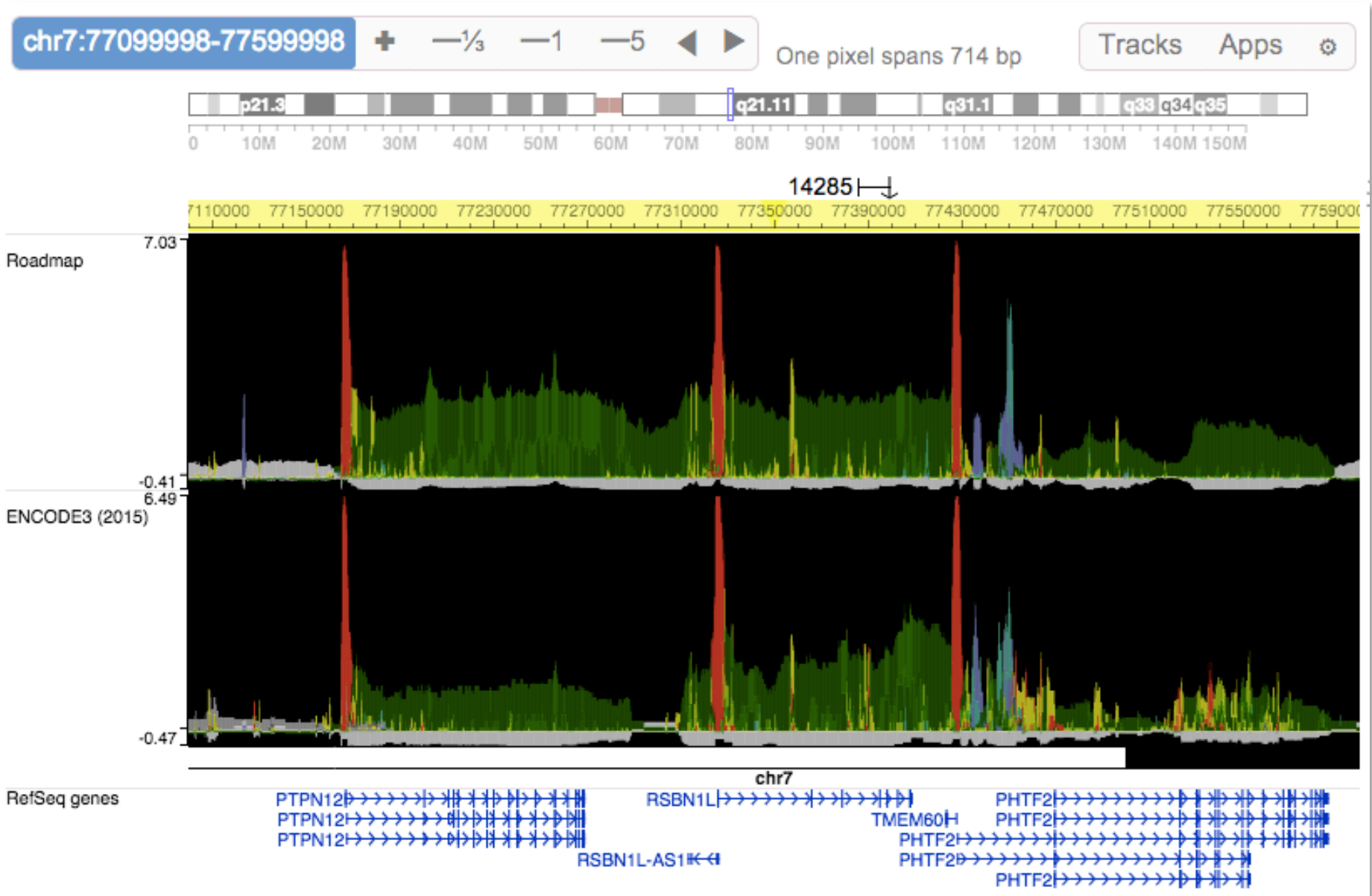


Chromatin state calls across many epigenomes can be viewed as an alignment of sequences with a finite alphabet



There are good ways of modeling such alignments: logos!
Information content of a region, considering background

ENCODE3 data epilogos



<http://tinyurl.com/encode-vs-roadmap>