

Mapping and navigating the human regulatory genome

Wouter Meuleman

Principal Investigator

Altius Institute for Biomedical Sciences
Seattle, WA USA



@nameluem

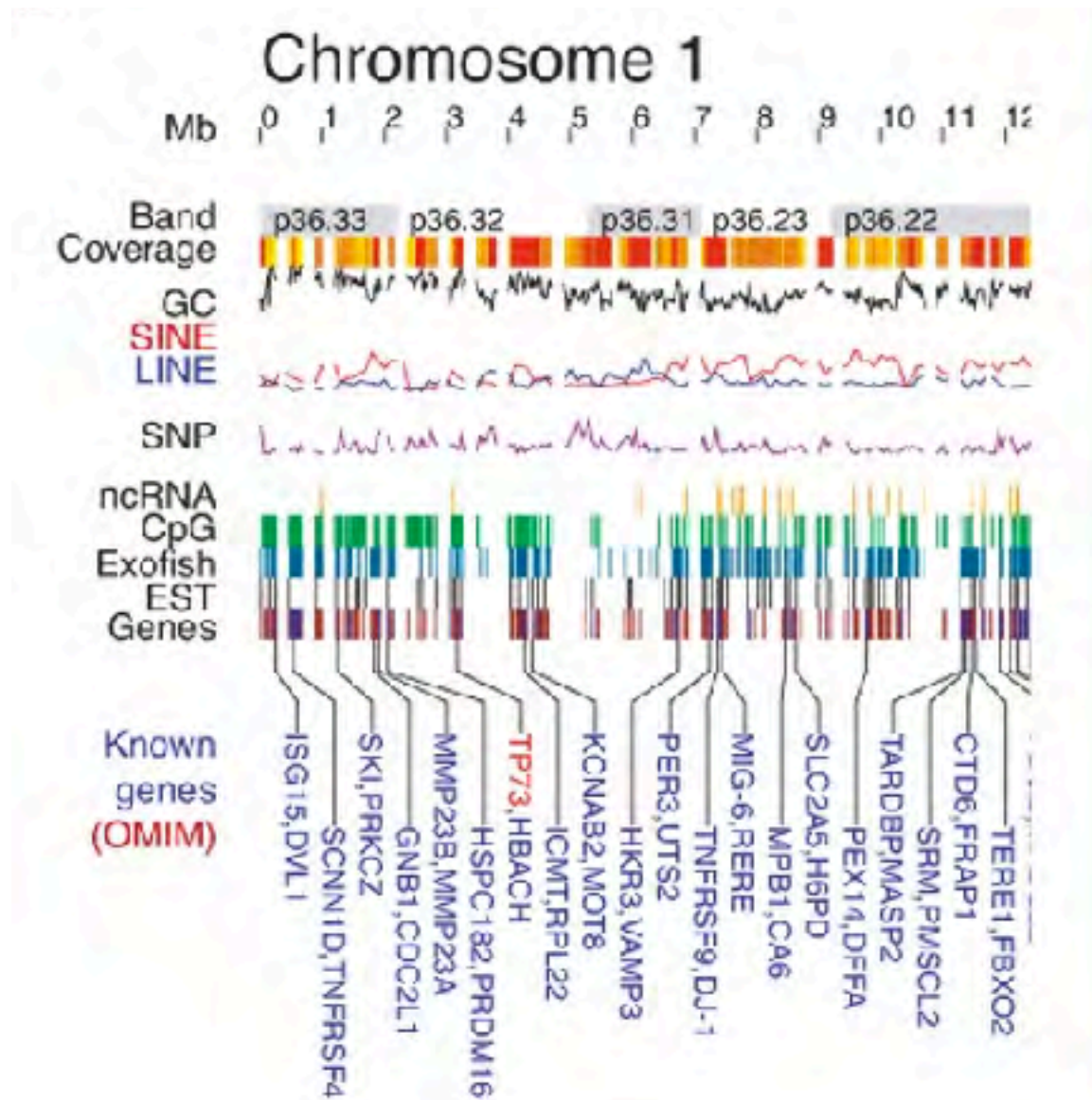
CSE 590C seminar

University of Washington

April 26, 2021

The Human Genome

The Map (2001)



Lander *et al.* Nature (2001)

20 years
↔

The Navigation System (1981)

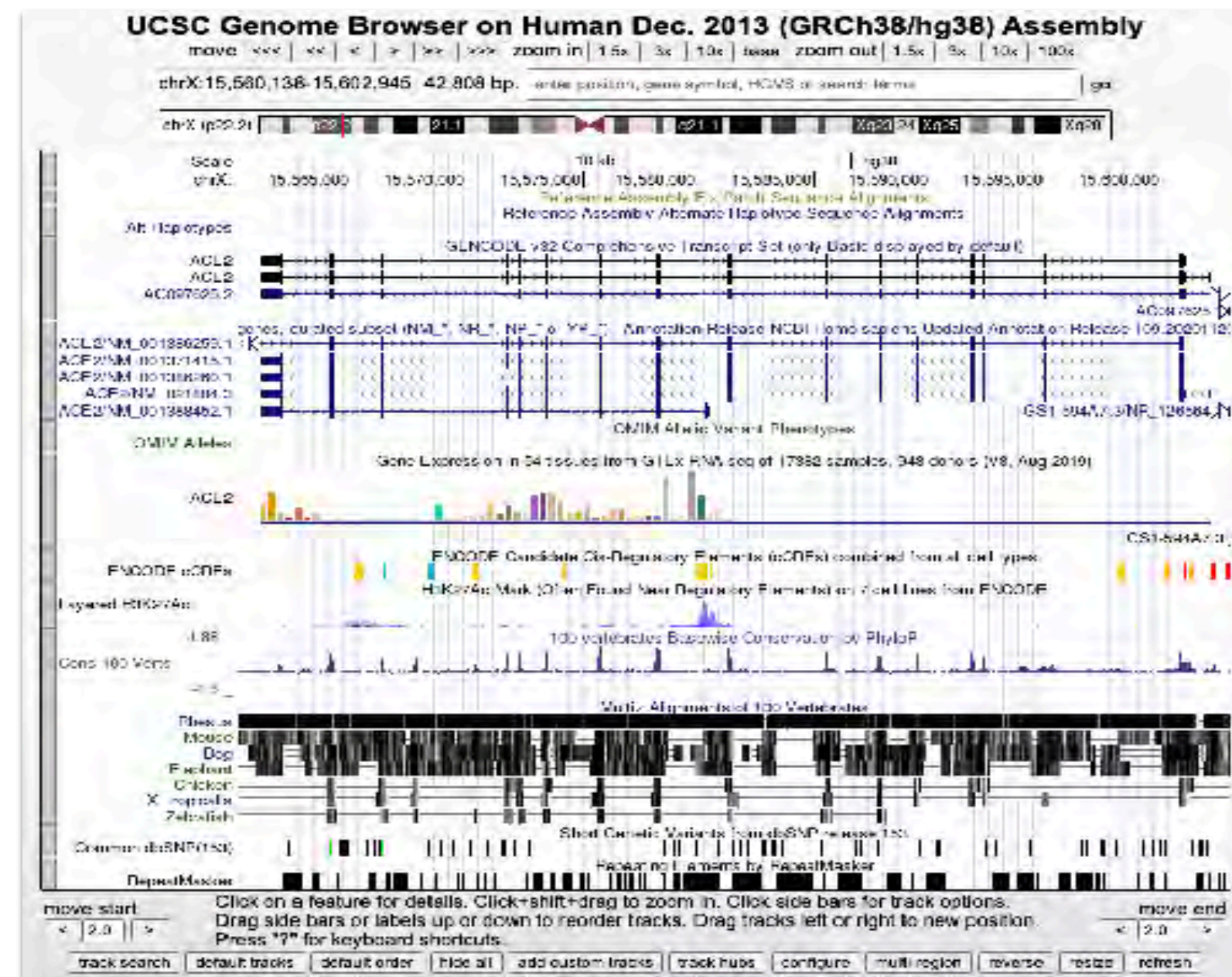


m66roepers @ Flickr

Fast-forward 20 years!

The Human Genome

The Map (2021)



UCSC Genome Browser (2021)

The Navigation System (2001)



TomTom Navigator (2001)

20 years
↔

From paper to screen, but still hard to access and interpret at scale

Maps (should) encourage exploration



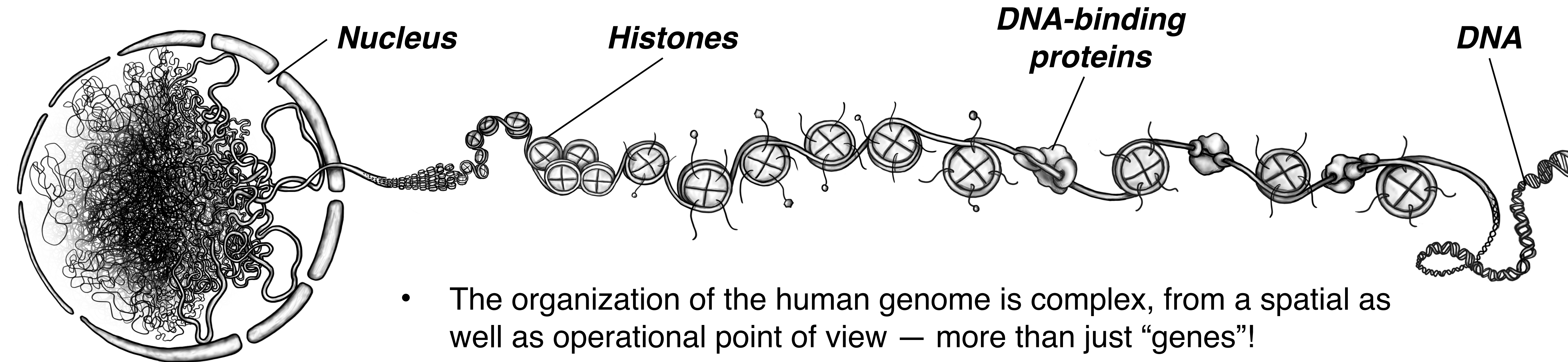
Walt Disney World® Resort

Maps are useful because they are wrong*

- Maps provide summarized representations of reality, highlighting only the most relevant information
- What is considered '**relevant**' depends on e.g. the mapped subject matter, map resolution, and context
- Towards a Disney Map of Genomics, we need to **annotate** the most 'exciting attractions' of the genome

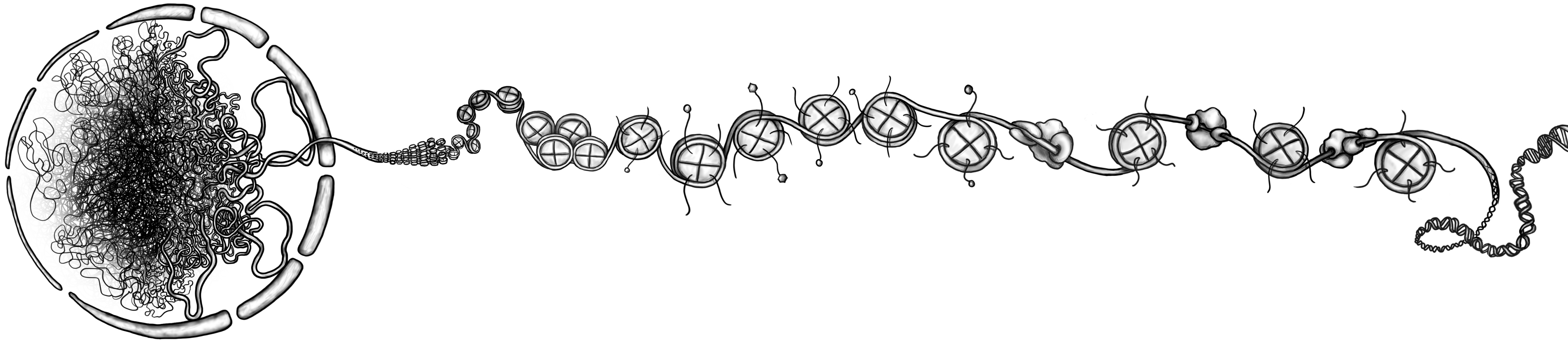


6 In search of 'relevance': annotating the (regulatory) human genome



- The organization of the human genome is complex, from a spatial as well as operational point of view — more than just “genes”!
- Protein-coding regions make up at most a few % of the genome, with regulatory elements hidden in its vast non-coding portion.
- To interrogate the (non-coding) genome, many experimental methods are available, most utilizing high-throughput sequencing.
- The resulting genome-wide datasets can be hard to interpret on their own, but offer lots of opportunities for creating useful annotations.

7 In search of 'relevance': annotating the (regulatory) human genome



Chromatin domains

Scale: 10kbp-Mbp

Guelen *et al.*, Nature (2008)
Peric-Hupkes, Meuleman *et al.* Mol. Cell (2010)
Meuleman *et al.*, Genome Res. (2013)

Chromatin states

Scale: 200bp-1kbp

Kundaje, Meuleman *et al.*, Nature (2015)
Claussnitzer *et al.*, NEJM (2015)
Marco, Meuleman *et al.*, Nature Comm. (2017)

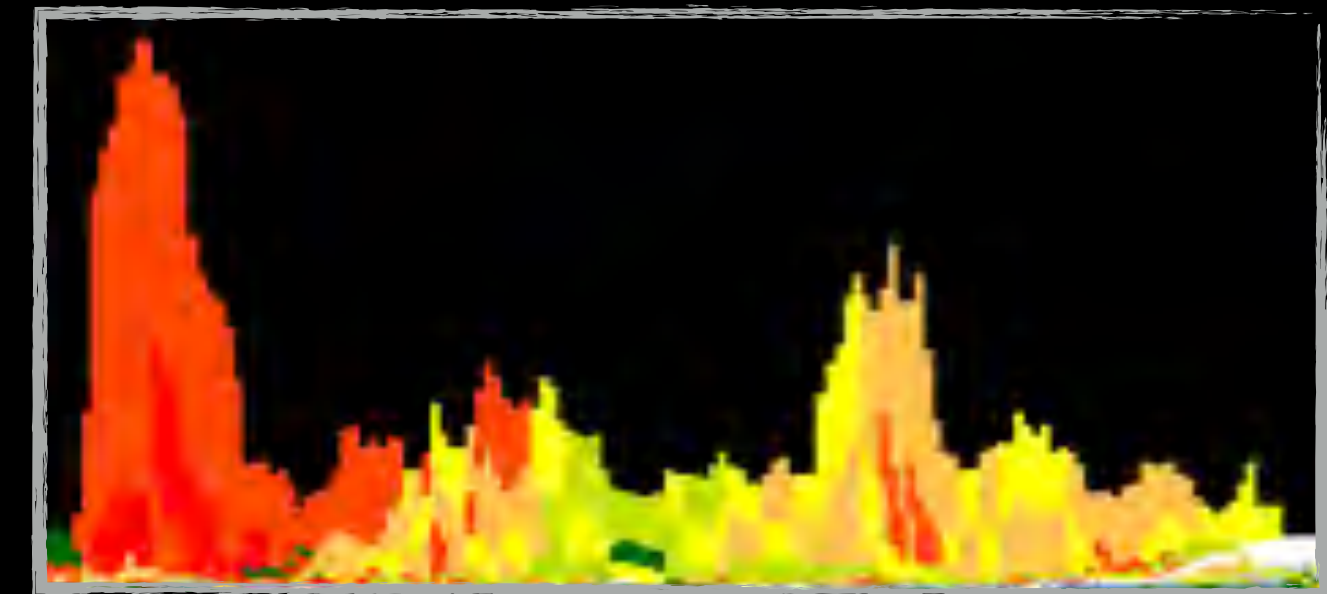
Chromatin accessibility

Scale: <200bp

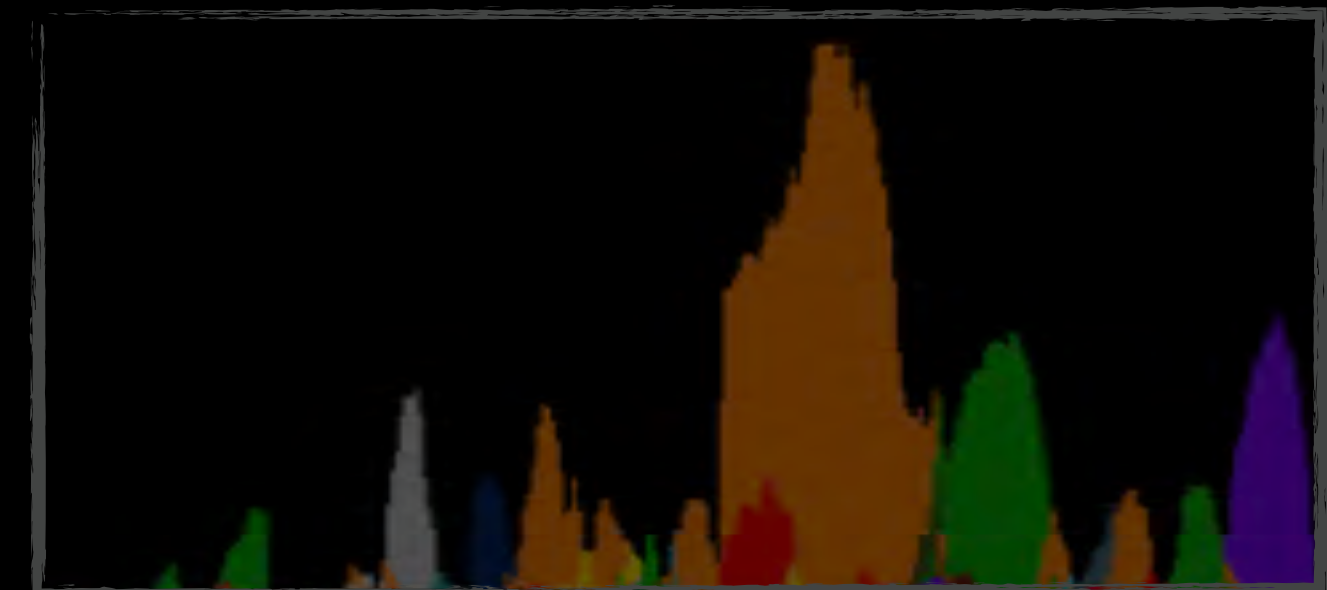
Meuleman *et al.*, Nature (2020)
Vierstra *et al.*, Nature (2020)
Boix *et al.*, Nature (2021)

In search of 'relevance': two types of genomic annotations

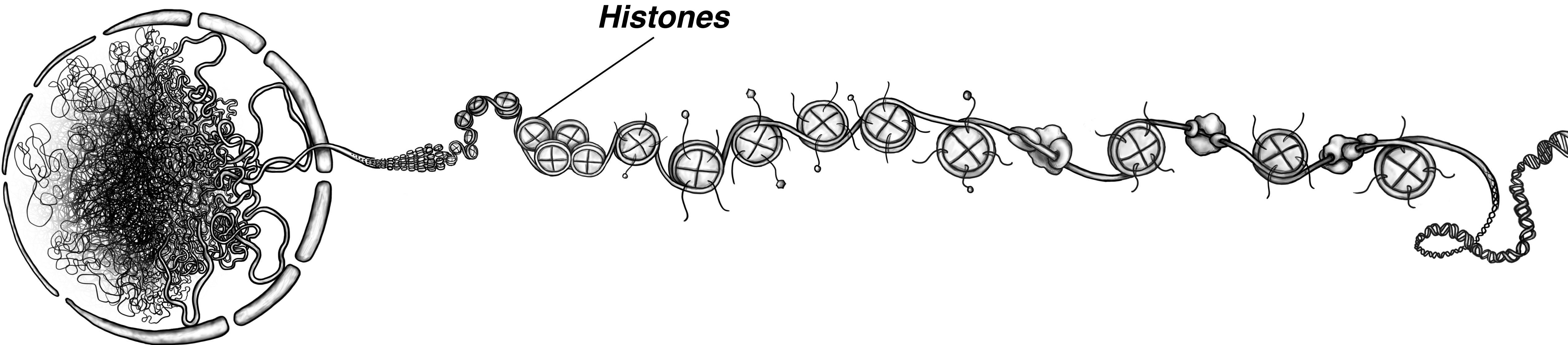
- **Chromatin states (epilogos):**
“What type of functionality does a genomic region encode?”
(e.g. **promoter**, **enhancer**, repressor)
- **Chromatin accessibility (DHS Index):**
“In which cellular contexts are regulatory regions utilized?”
(e.g. **cardiac**, **lymphoid**, **neural**)



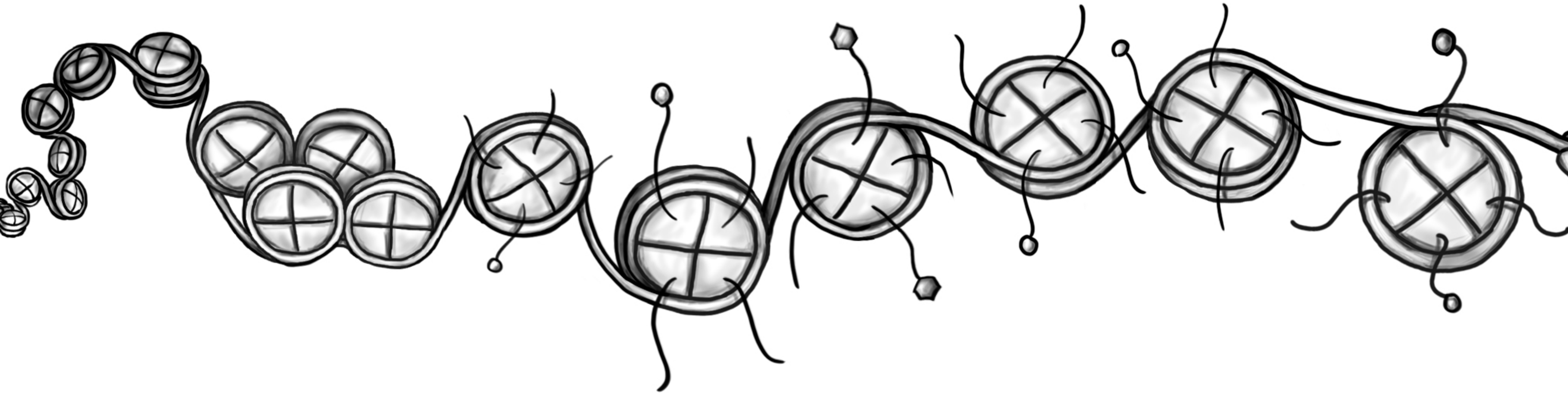
<https://epilogos.net>



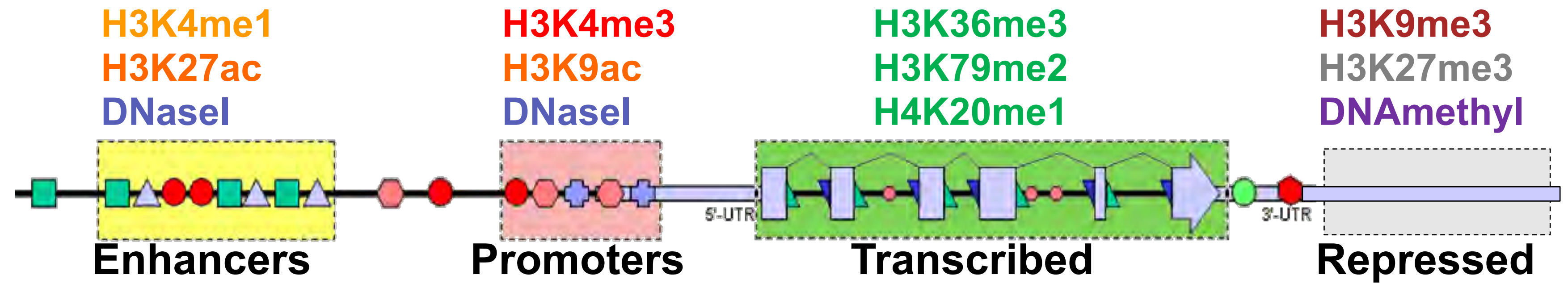
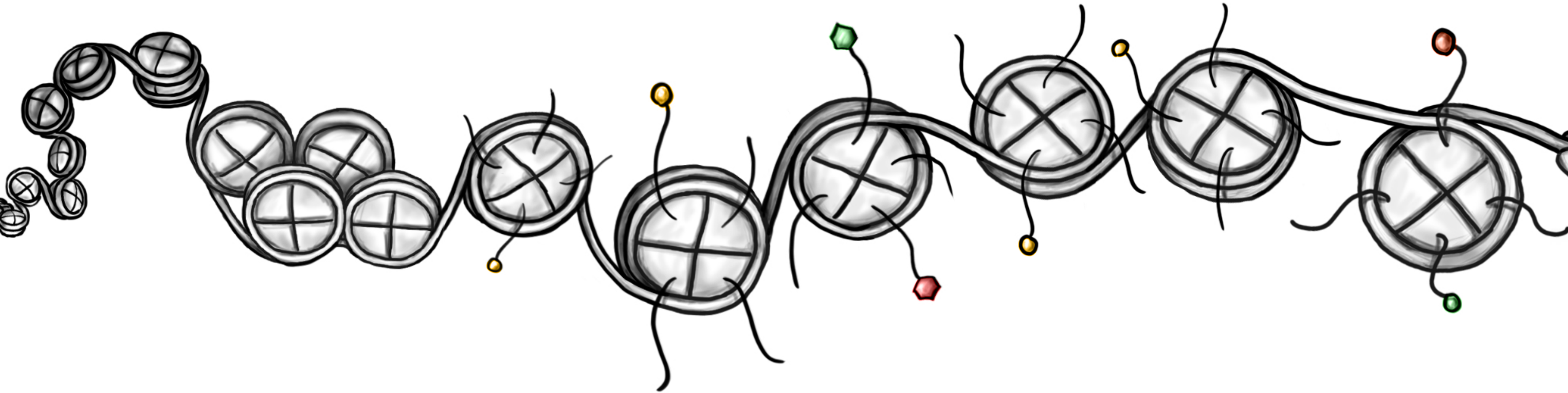
Meuleman *et al.*, 2020 & ongoing



10 Histone tails can be chemically tagged with *epigenomic marks*



11 These epigenomic marks are associated with *functional elements*



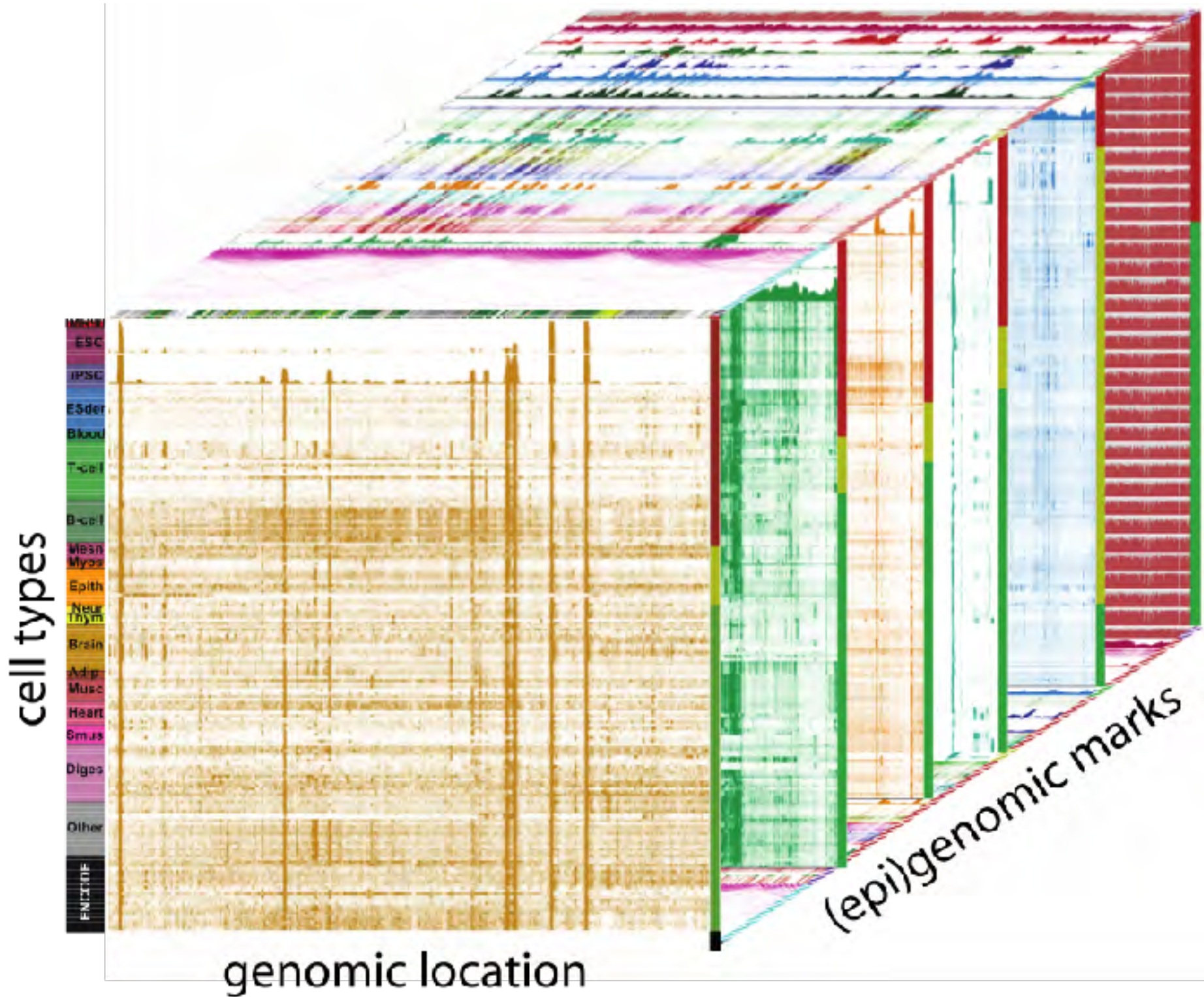
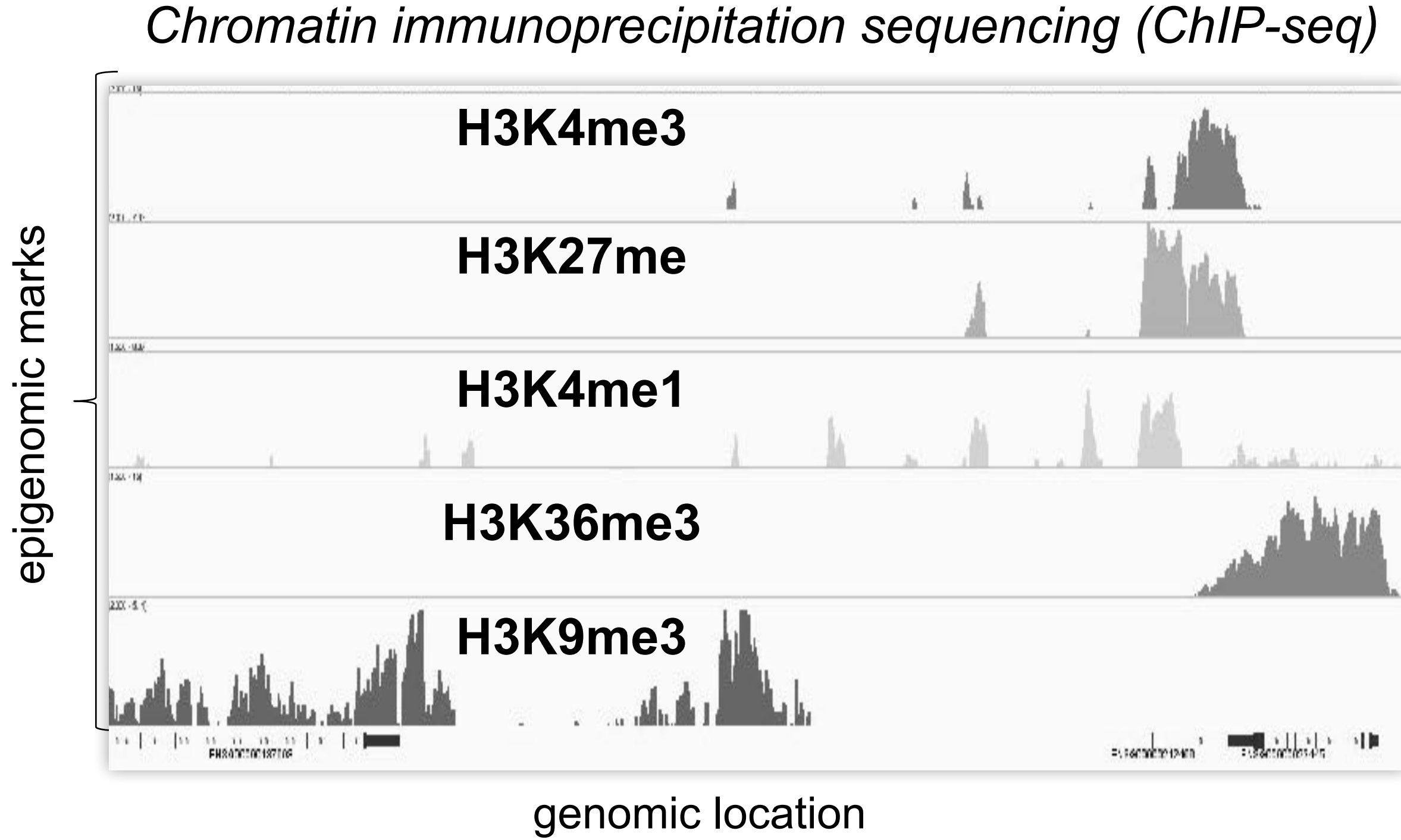
12 The NIH Roadmap Epigenomics Project (2008-2017)

Goal: create reference maps of a wide variety of epigenomic marks across many cell types from healthy individuals.

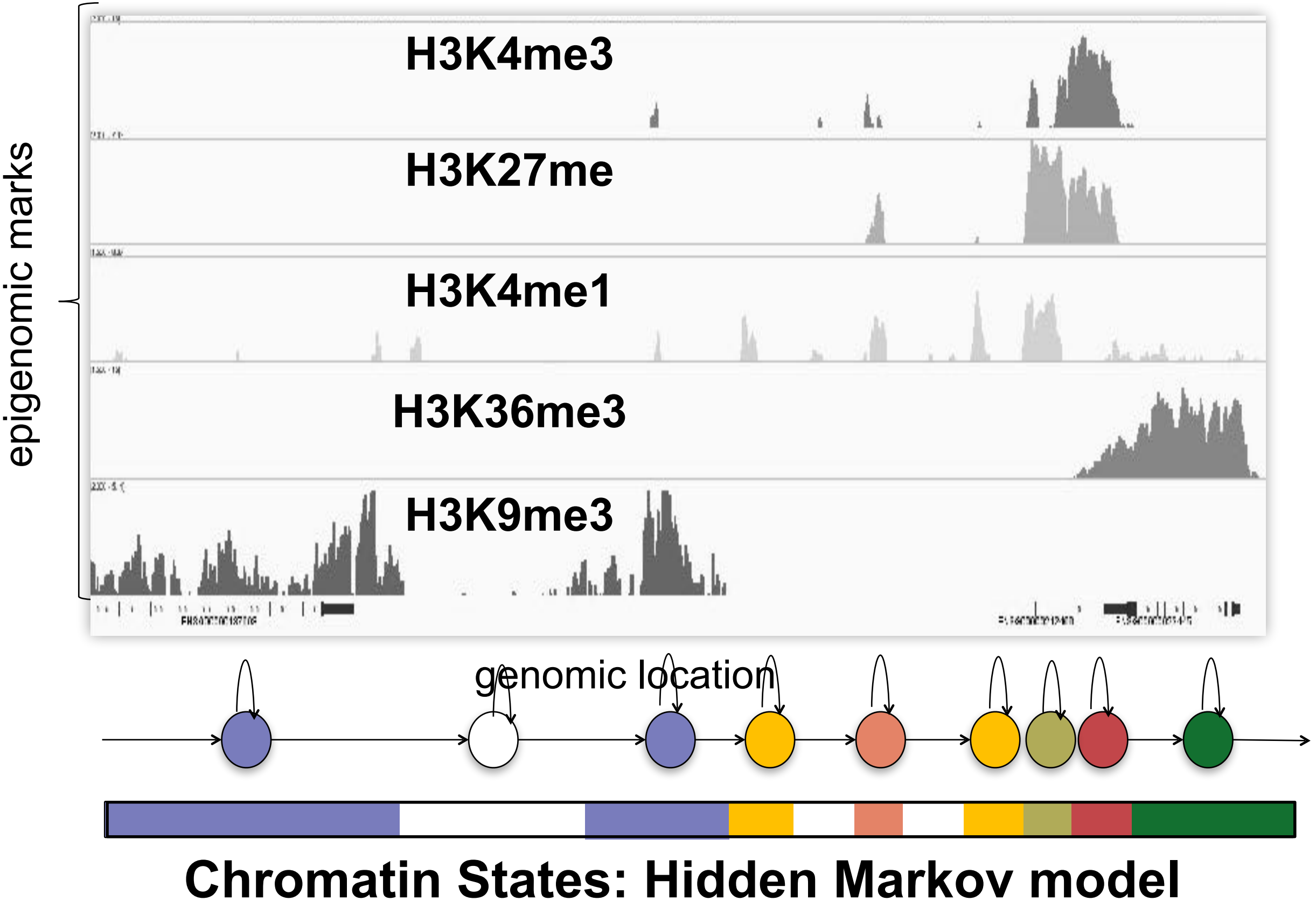


- **4 Reference Epigenome Mapping Centers (REMCs)**
- **Central data repository and read mapping at Baylor**
- **Uniform processing and integrative analysis at MIT**

13 Genome-wide profiling of epigenomic marks has resulted in giant data cubes



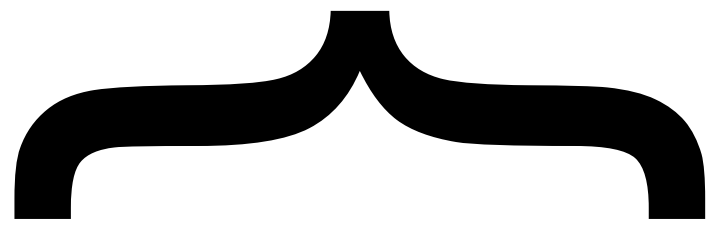
14 Data can be summarized by learning a limited number of ***chromatin states***



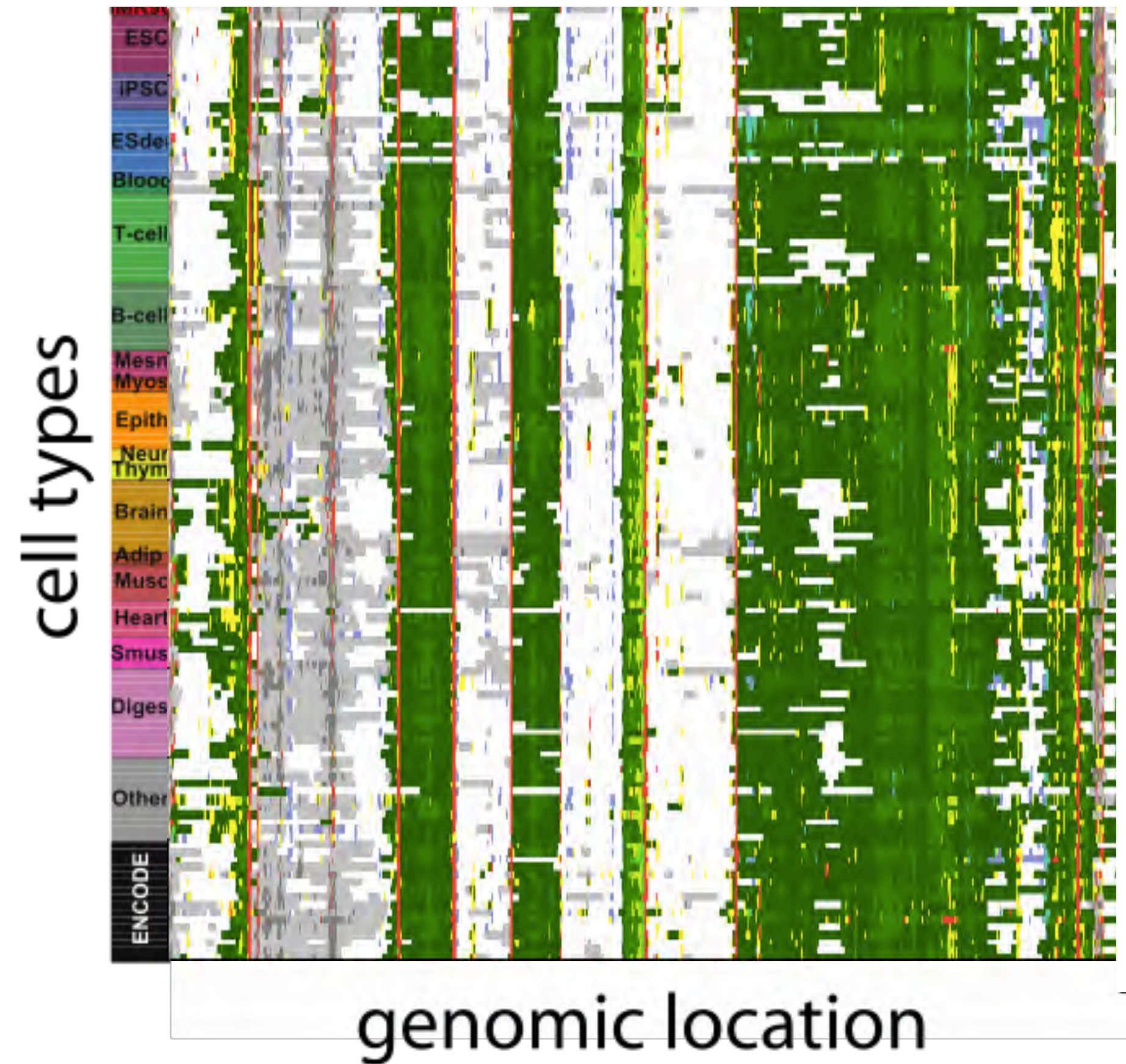
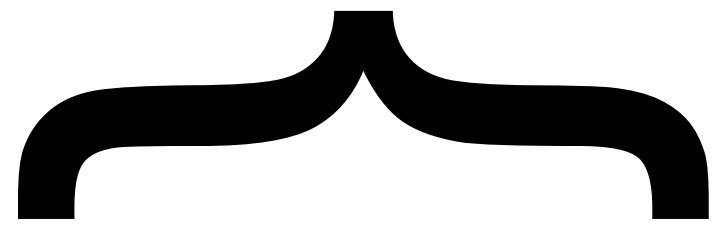
Abbreviation	emissions	Cov.	
TssA	[H3K4me3]	0.7%	<i>Promoters</i>
TssAFlnk	[H3K4me3, H3K27me]	0.5%	
TxFlnk	[H3K4me3, H3K27me, H3K4me1]	0.1%	
Tx	[H3K4me3, H3K27me]	3.6%	<i>Genes (Tx)</i>
TxWk	[H3K4me3, H3K27me, H3K4me1]	11.6%	
EnhG	[H3K4me1, H3K27me]	0.4%	<i>Enhancers</i>
Enh	[H3K4me1, H3K27me, H3K4me3]	2.8%	
ZNF/Rpts	[H3K4me1, H3K27me, H3K36me3]	0.2%	<i>Repeats</i>
Het	[H3K4me1, H3K27me, H3K36me3, H3K9me3]	2.6%	
TssBiv	[H3K4me3, H3K27me]	0.1%	<i>Bivalent regions</i>
BivFlnk	[H3K4me3, H3K27me, H3K4me1]	0.1%	
EnhBiv	[H3K4me3, H3K27me, H3K4me1, H3K27me]	0.1%	
ReprPC	[H3K9me3, H3K27me]	1.2%	<i>Repressed regions</i>
ReprPCWk	[H3K9me3, H3K27me, H3K4me1]	8.3%	
Quies	[H3K9me3, H3K27me]	67.8%	

15 This allows us to transform the 3D cube into a 2D matrix:

A color here...



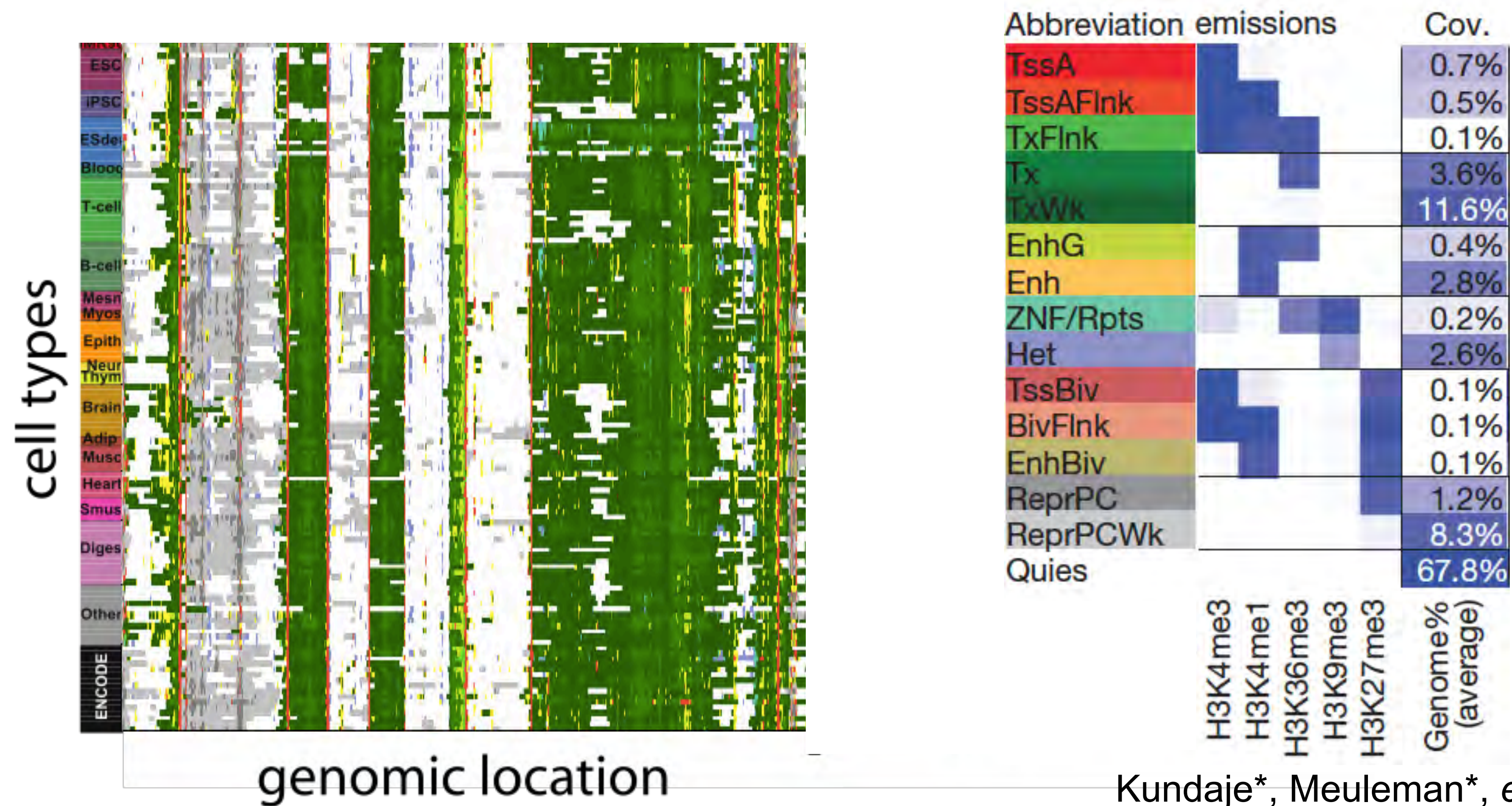
...corresponds to a state here:



Abbreviation	emissions	Cov.
TssA		0.7%
TssAFlnk		0.5%
TxFlnk		0.1%
Tx		3.6%
TxWk		11.6%
EnhG		0.4%
Enh		2.8%
ZNF/Rpts		0.2%
Het		2.6%
TssBiv		0.1%
BivFlnk		0.1%
EnhBiv		0.1%
ReprPC		1.2%
ReprPCWk		8.3%
Quies		67.8%

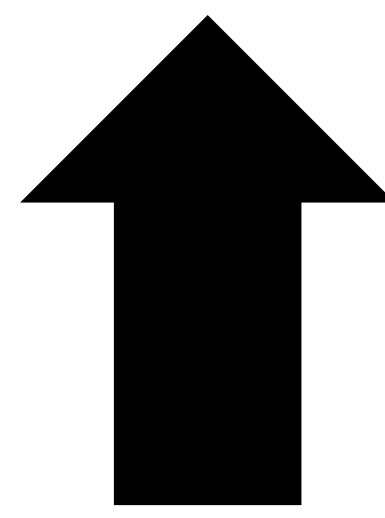
16 A reference map of chromatin states across 127 epigenomes

Documents the dynamics of chromatin states between cell types, e.g. during cell differentiation

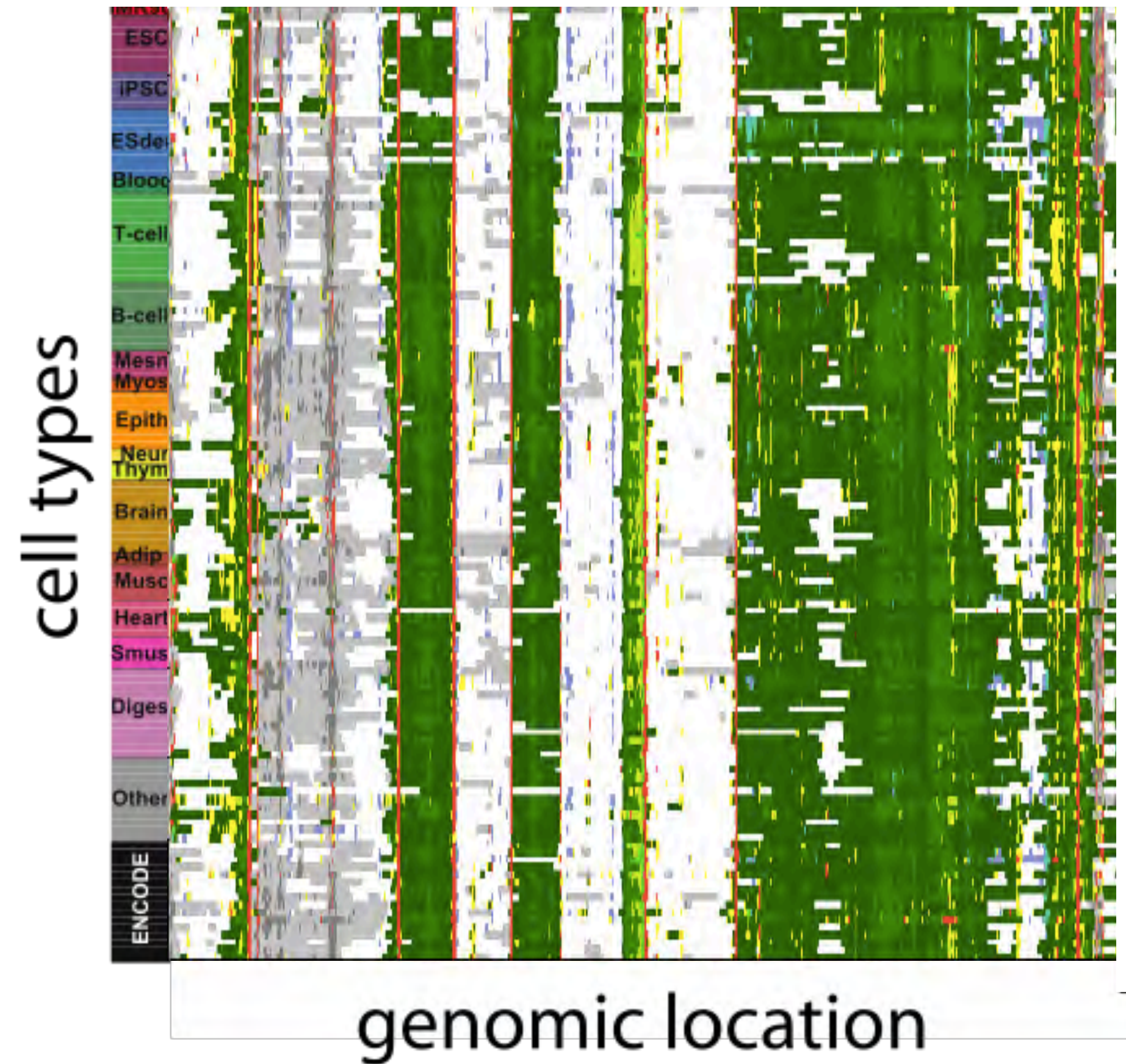


Kundaje*, Meuleman*, et al., Nature (2015)

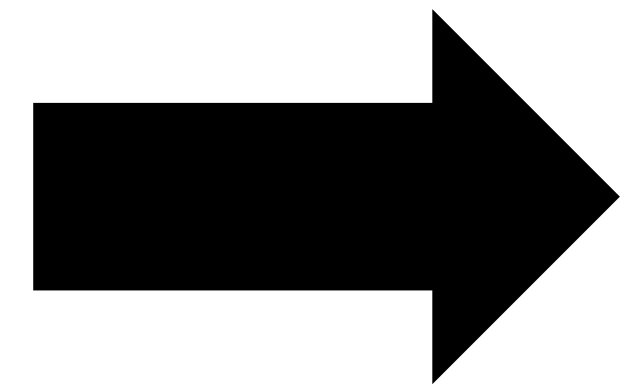
17 New methods are needed to navigate these maps



Many more epigenomes
are being profiled
(cell types, disease states,
personal epigenomics, etc)

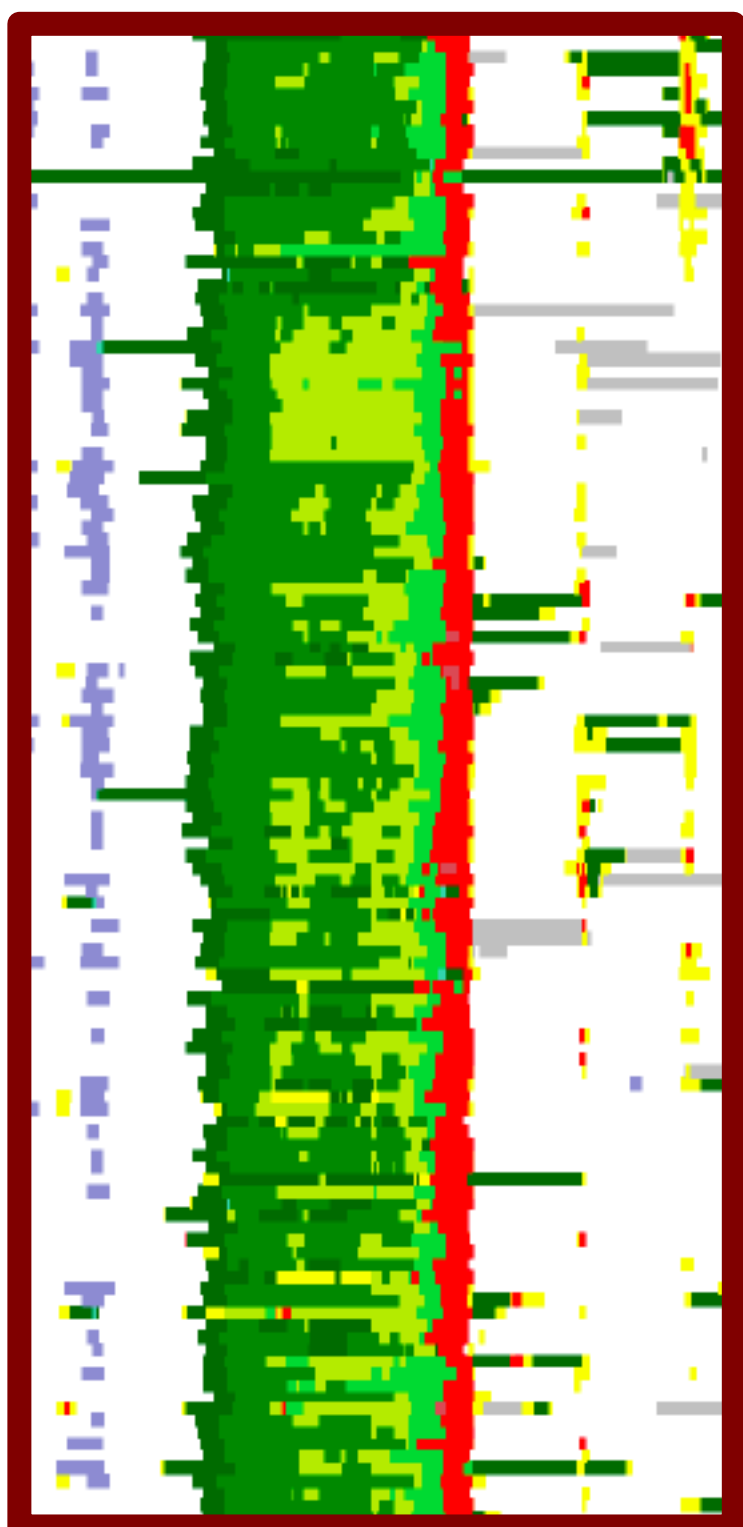


The genome is large
(shown here is only
0.0267% of the genome)

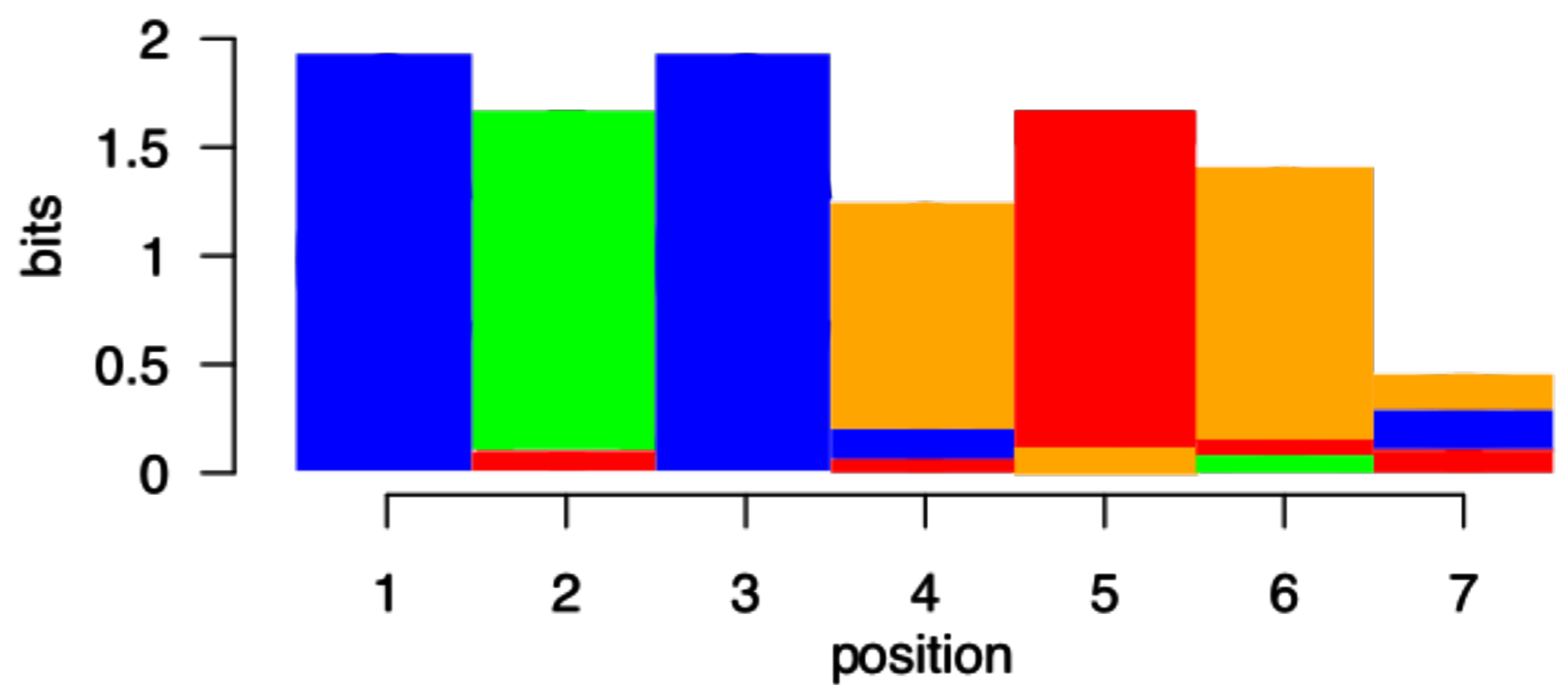
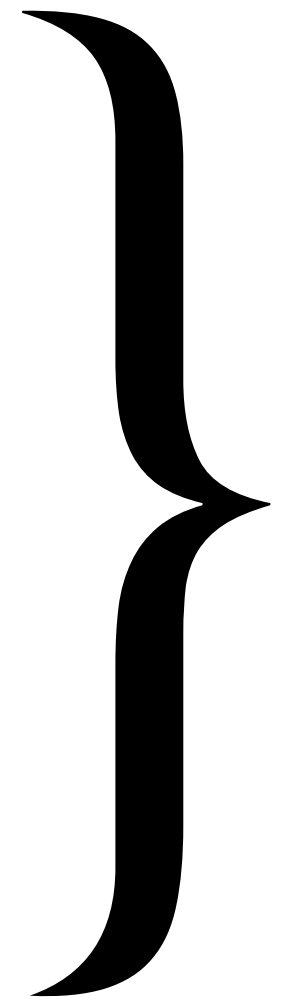


18 Chromatin states across many epigenomes: analogy with sequence motifs

Alignments of multiple sequences

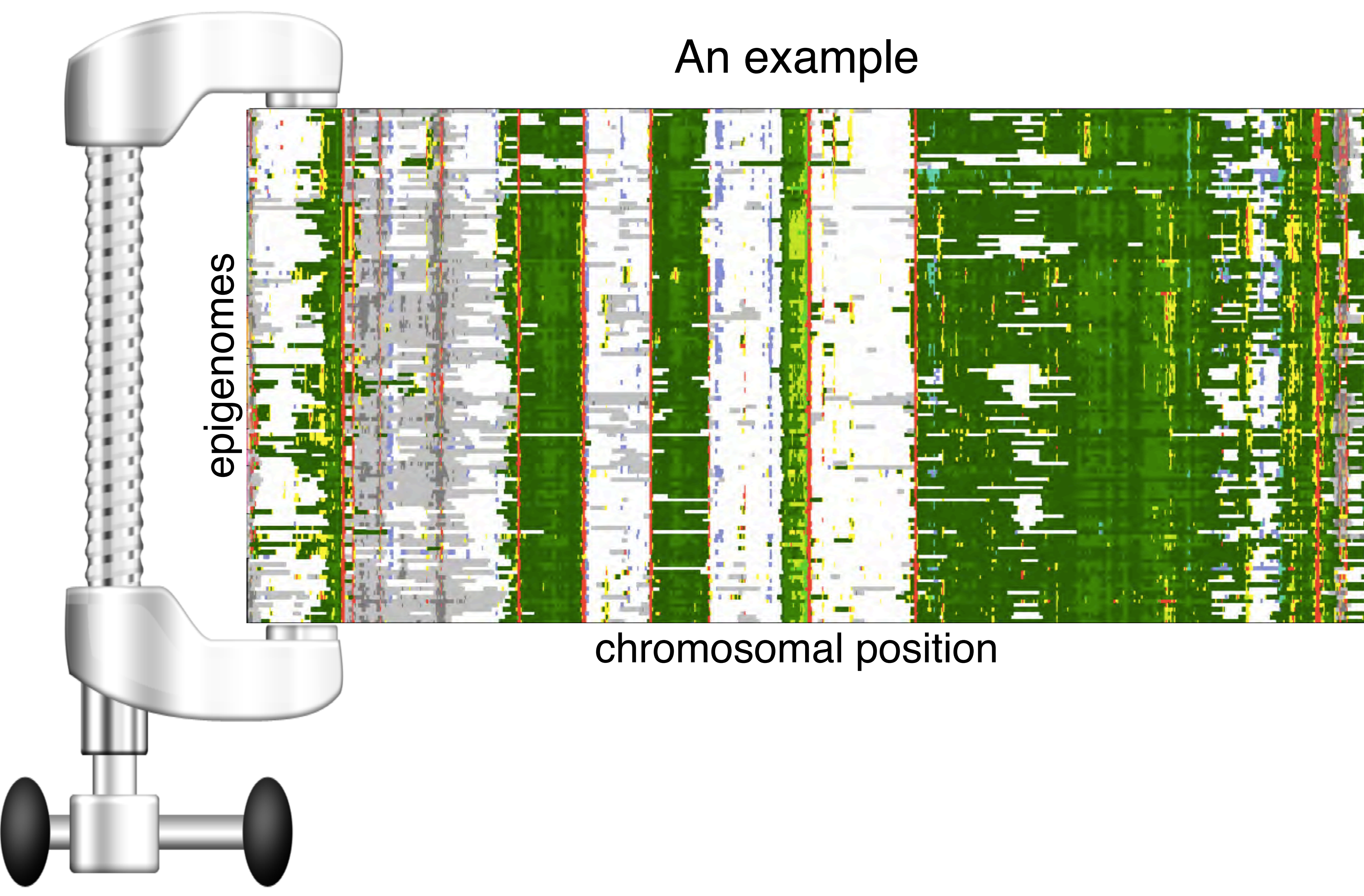


CTCTTAT
CACGTGC
CACGTGC
CACGTGG
CACGTGG
CACGTGG
CACGTGC
CACGTGG
CACGTGT
CACGTGC
CACGTGT
CACGTGC
CACCTGT
CACCGTC
CACGTGC
CACGTGC
CACGTGG
CACGTGT
CACGTGG
CACGTGG
CACGTGG

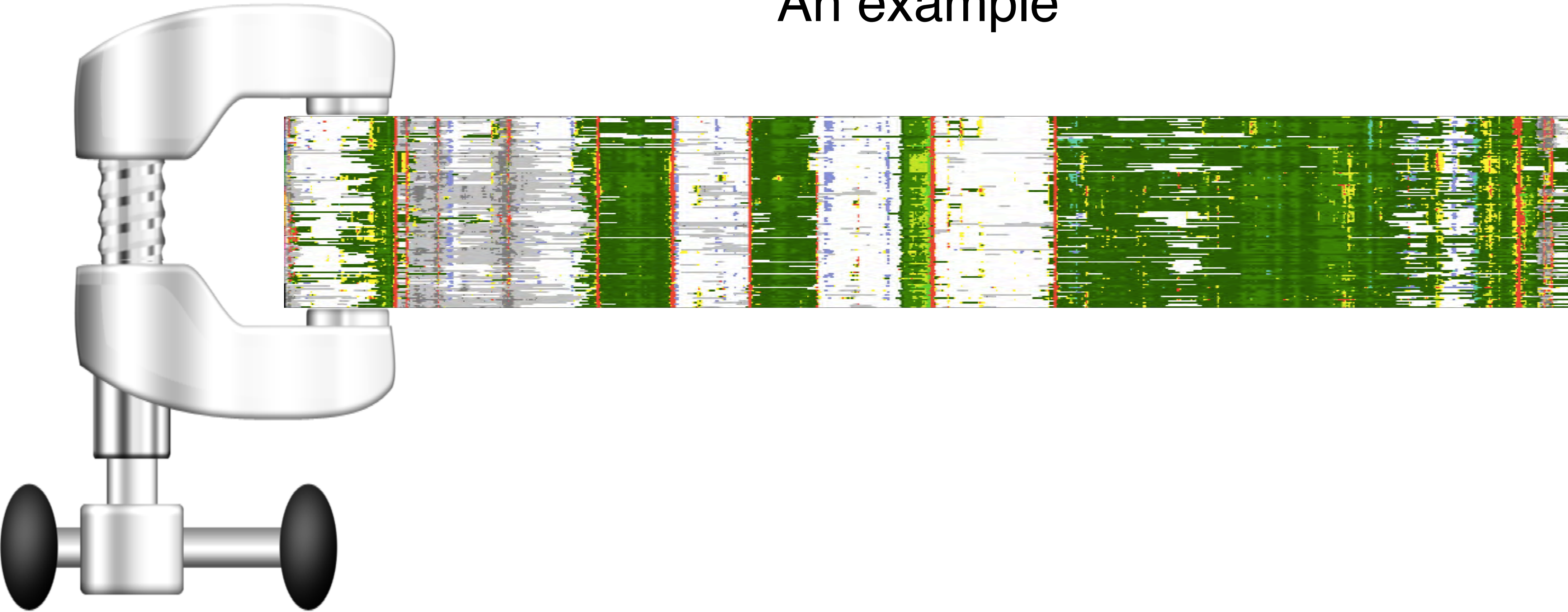


There are good ways of modeling such alignments: logos!
Information content of a region, considering background

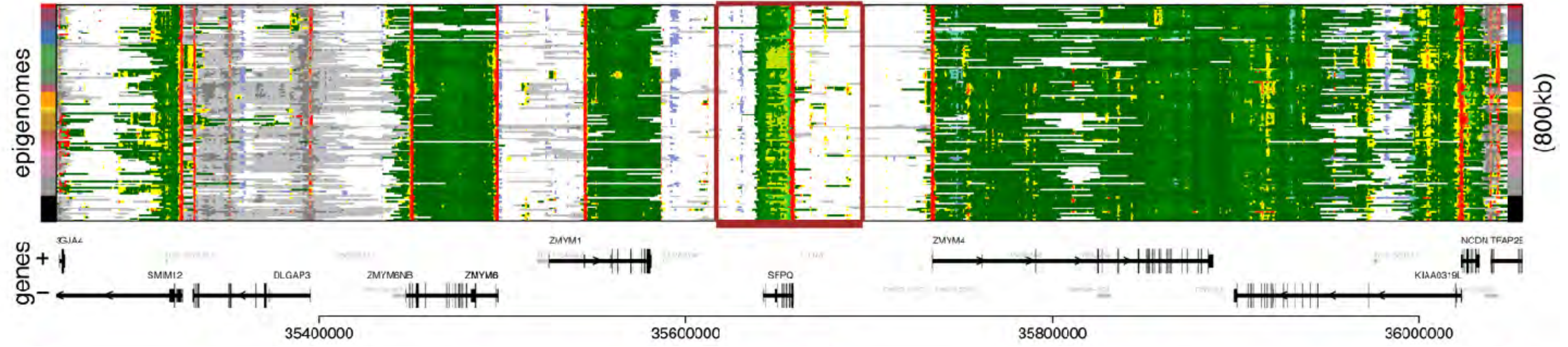
An example



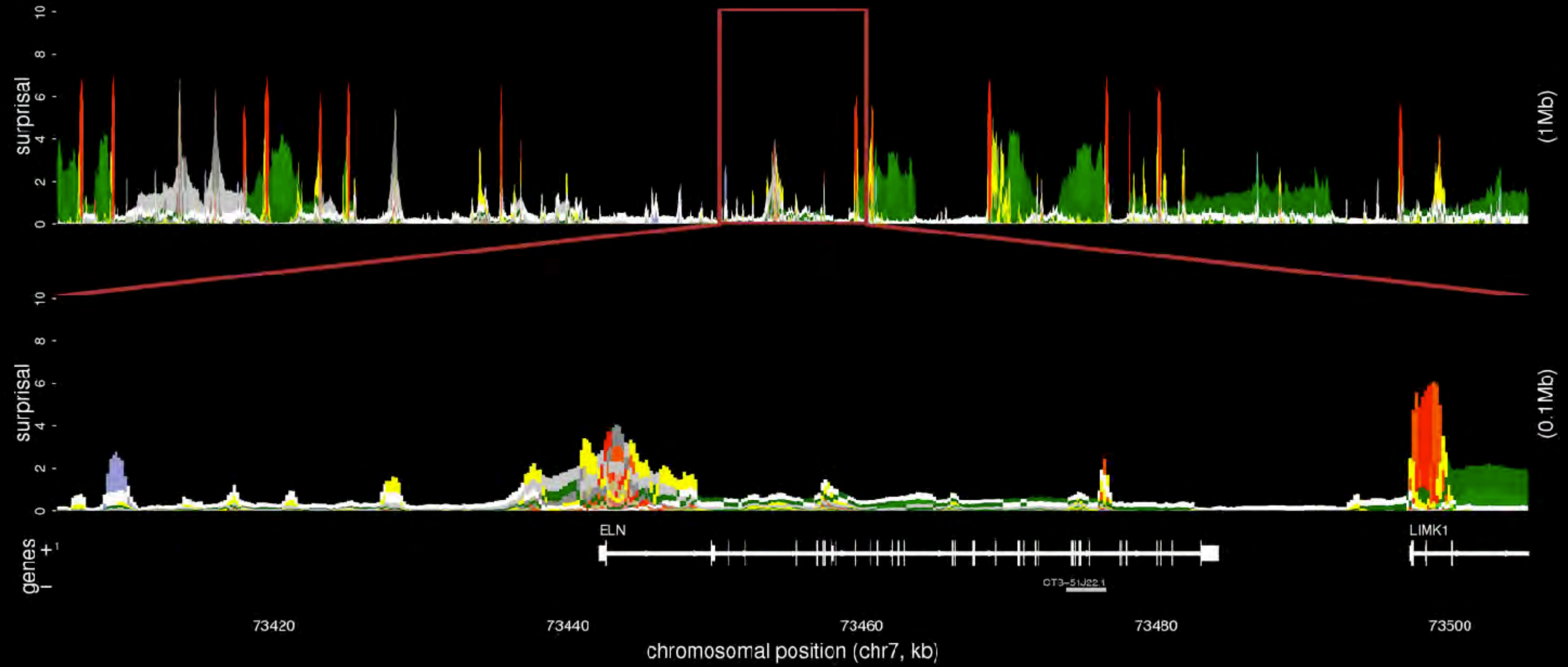
An example



epilogos

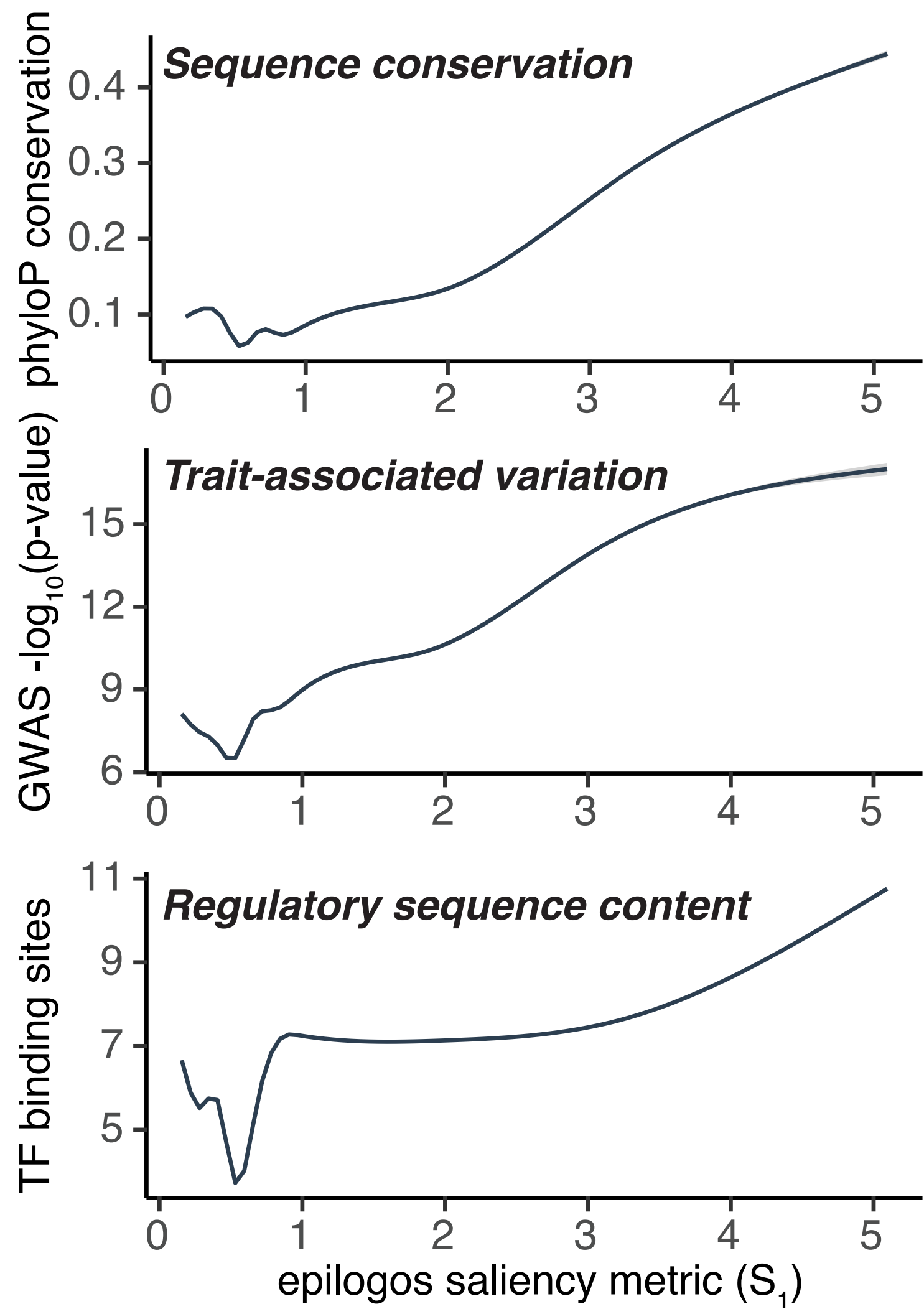
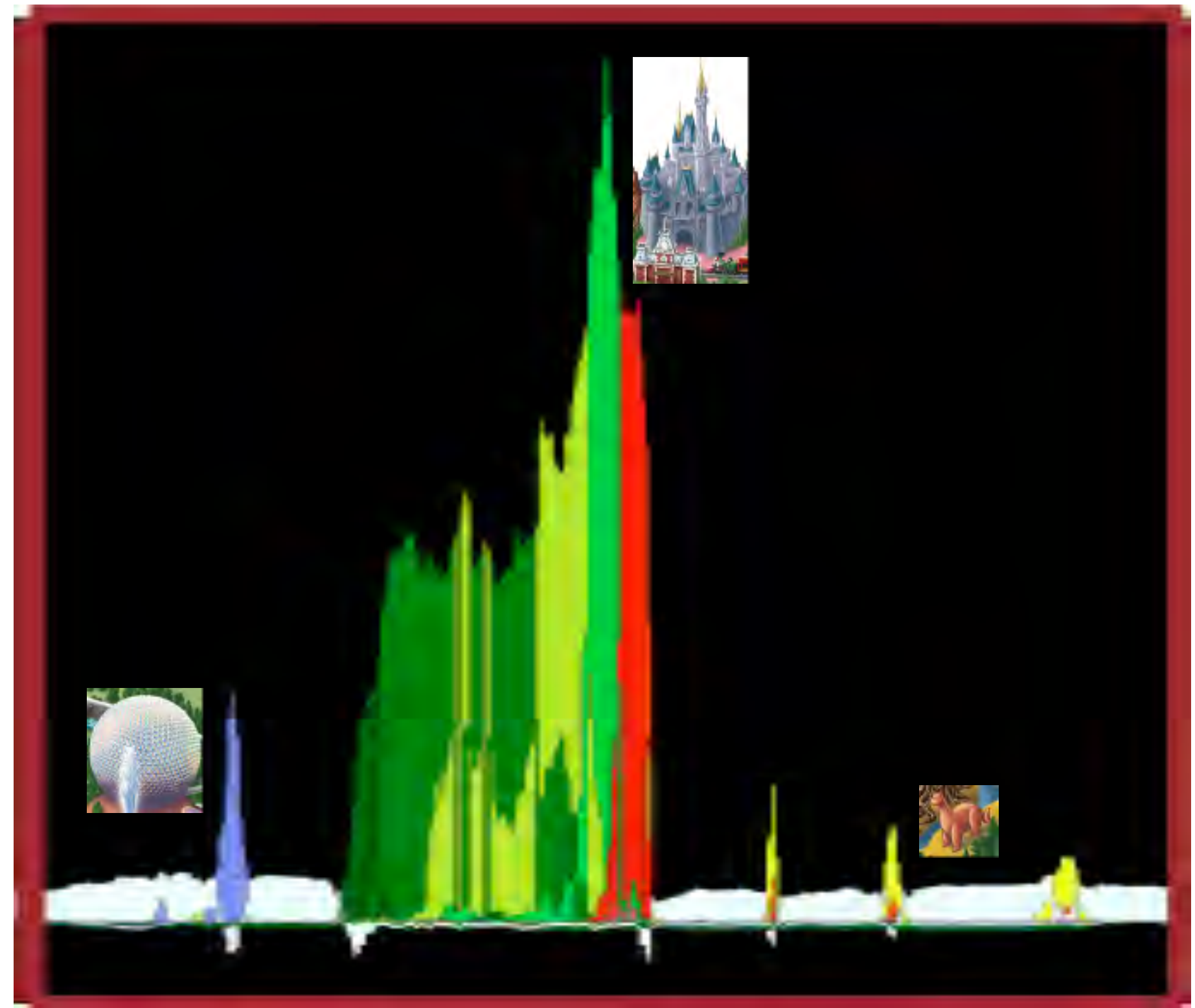


epilogos



23 Epilogos saliency metrics enrich for functionally relevant regions

epilogos saliency metric ↑



You're looking at a summarization of ~5,000 genome-wide datasets

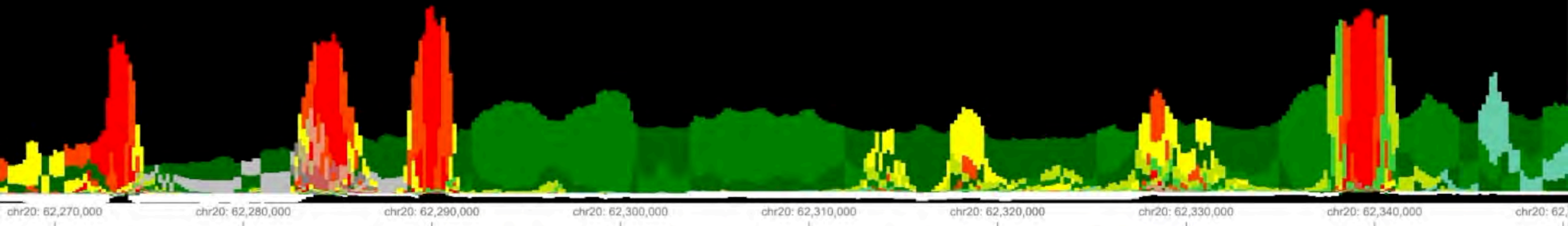
epilogos.net

Specify an interval or gene

Go

I'm Feeling Lucky

e.g., use query terms like HGNC symbols
(HOXA1, NFKB1, etc.) or genomic regions
(chr17:41155790-41317987, etc.)



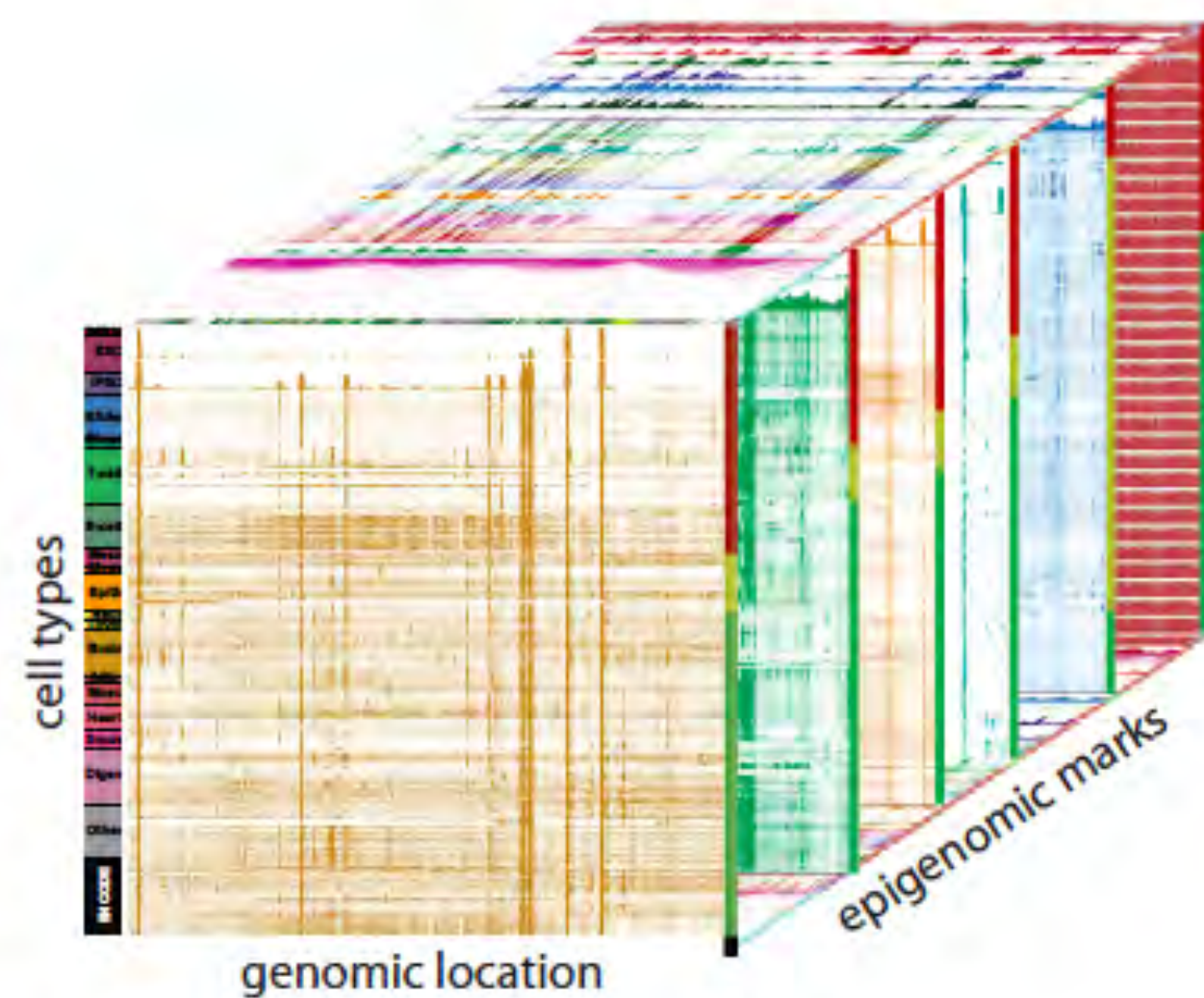
see how  it works

epilogos.net

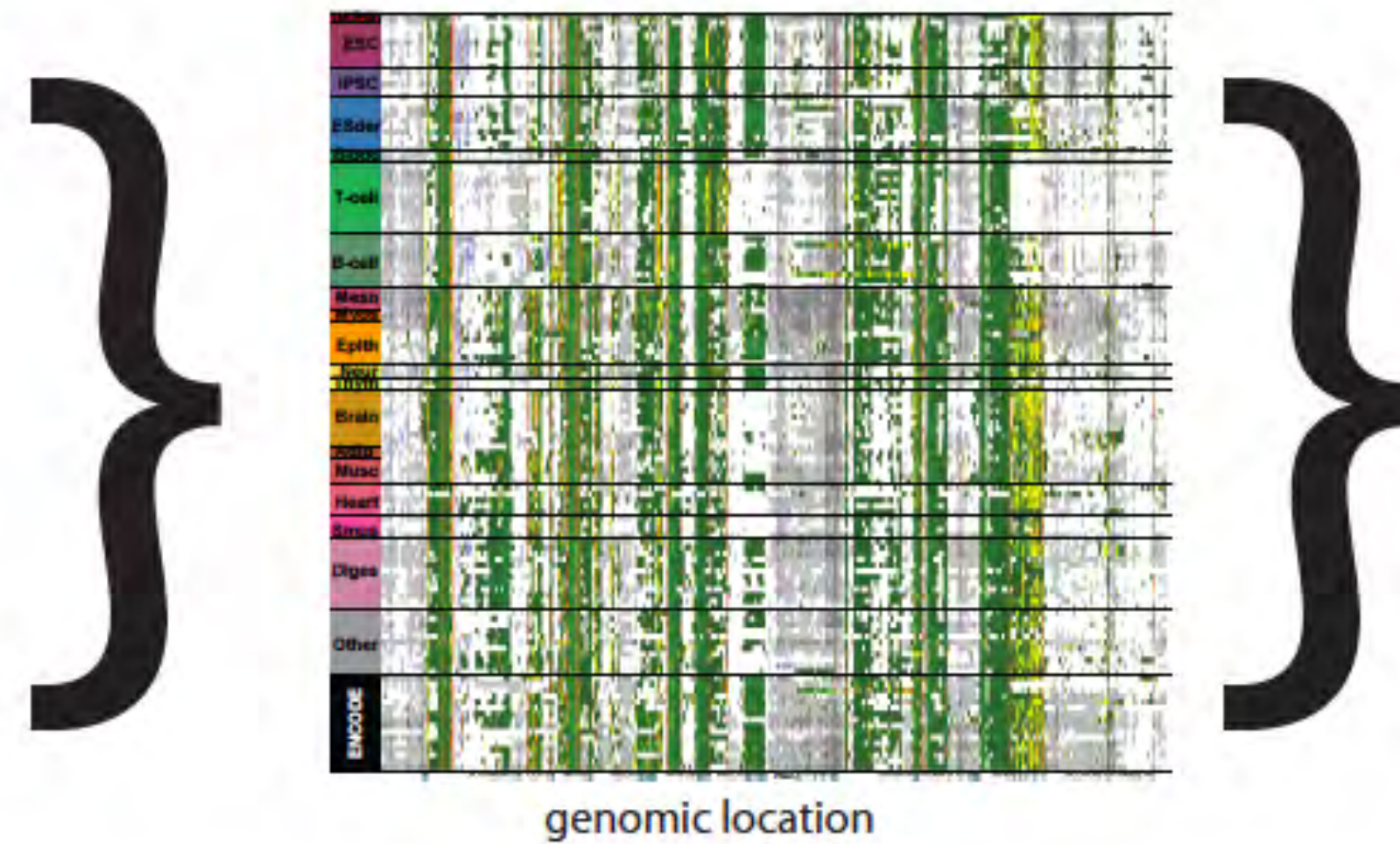
epilogos

Interpretation of large-scale (epi)genomic datasets through information-based dimensionality reduction

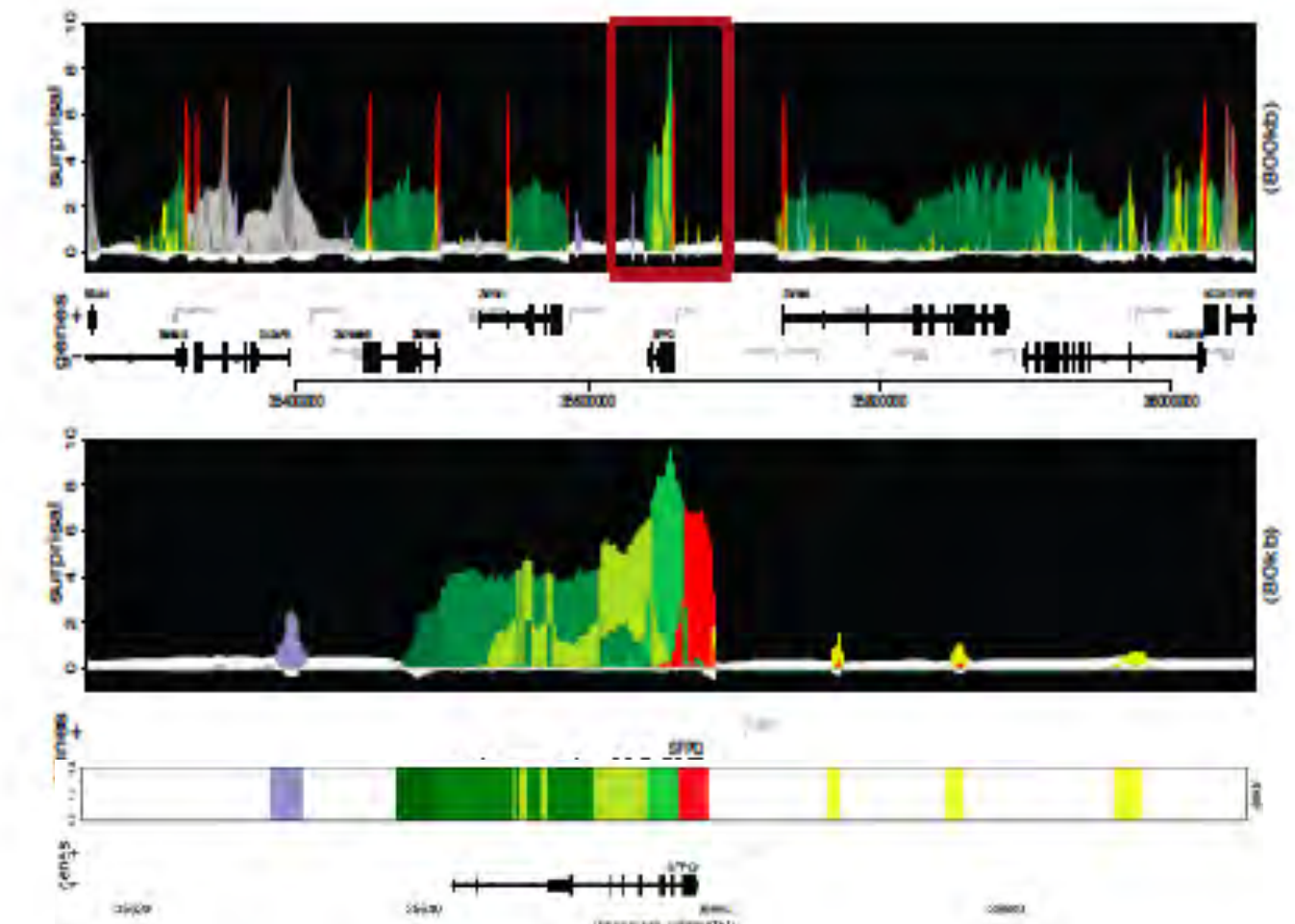
3D



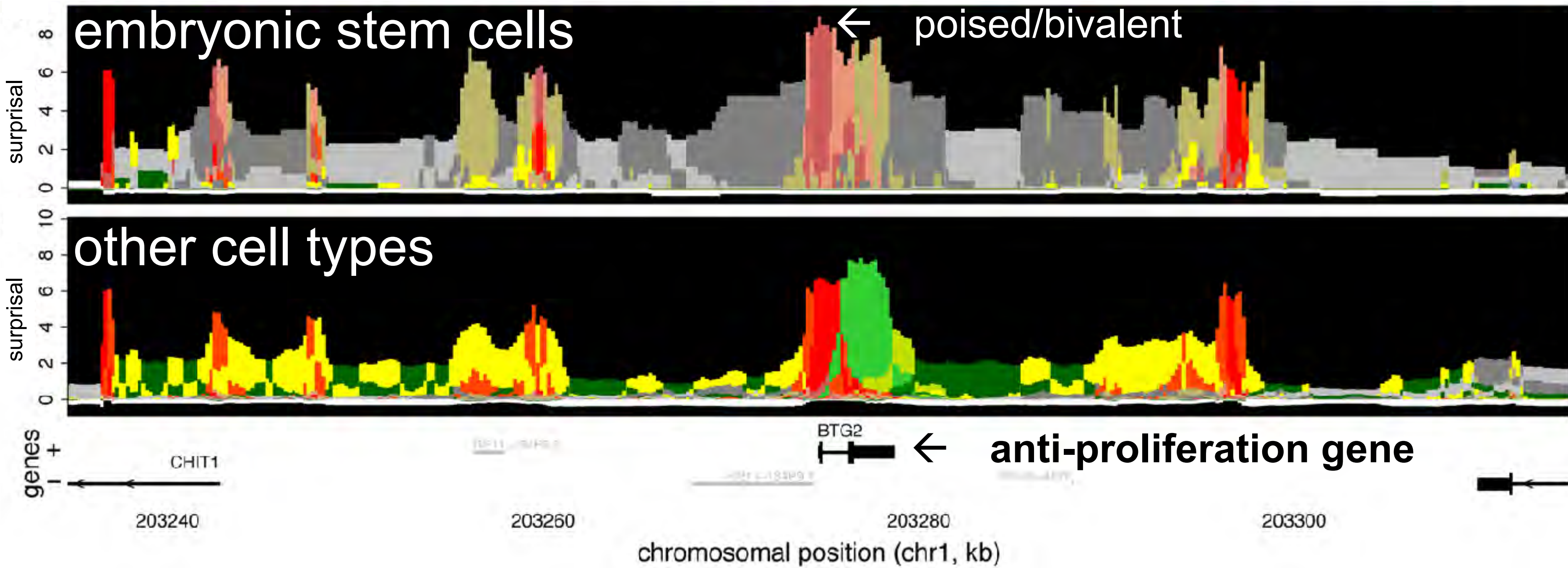
2D



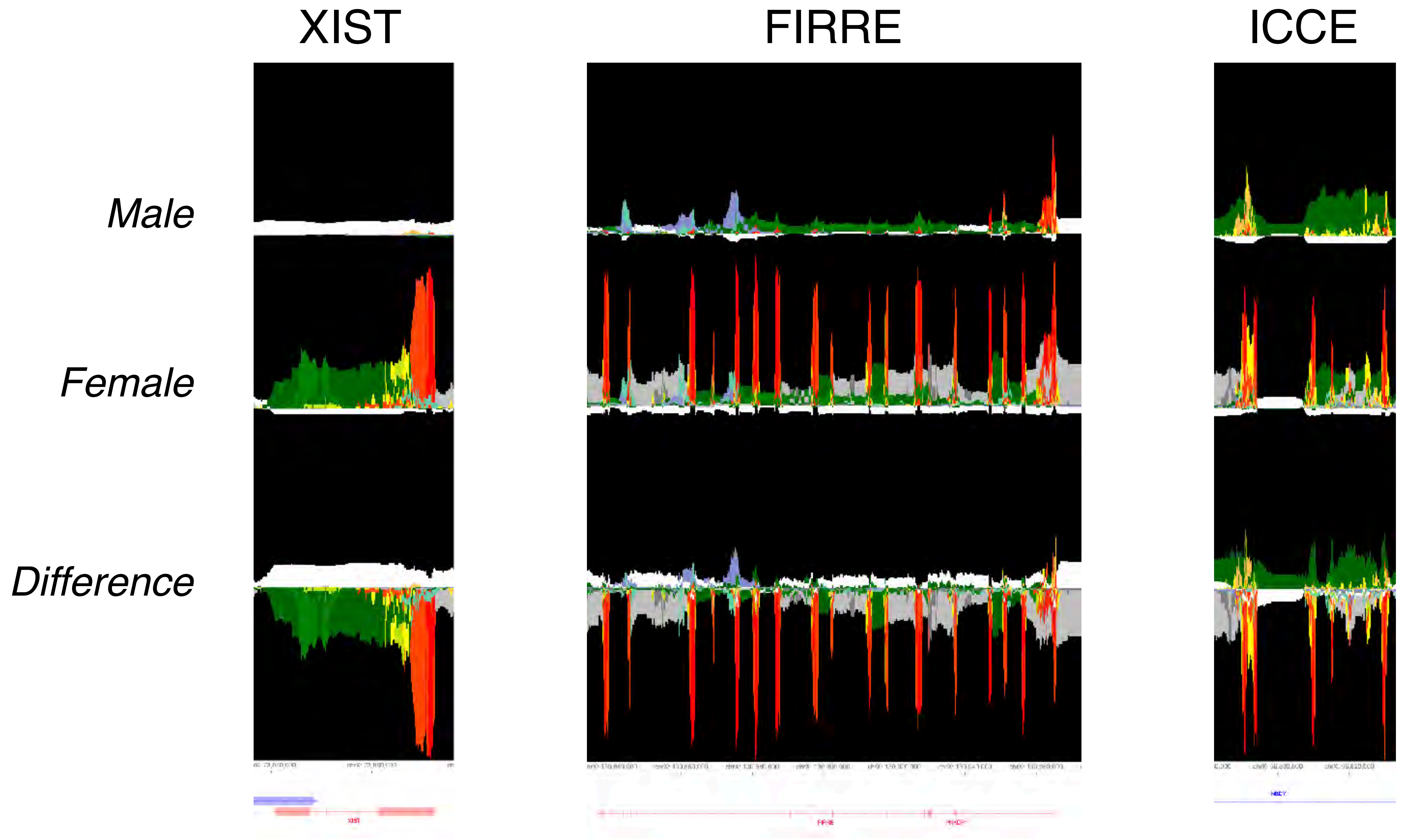
1D



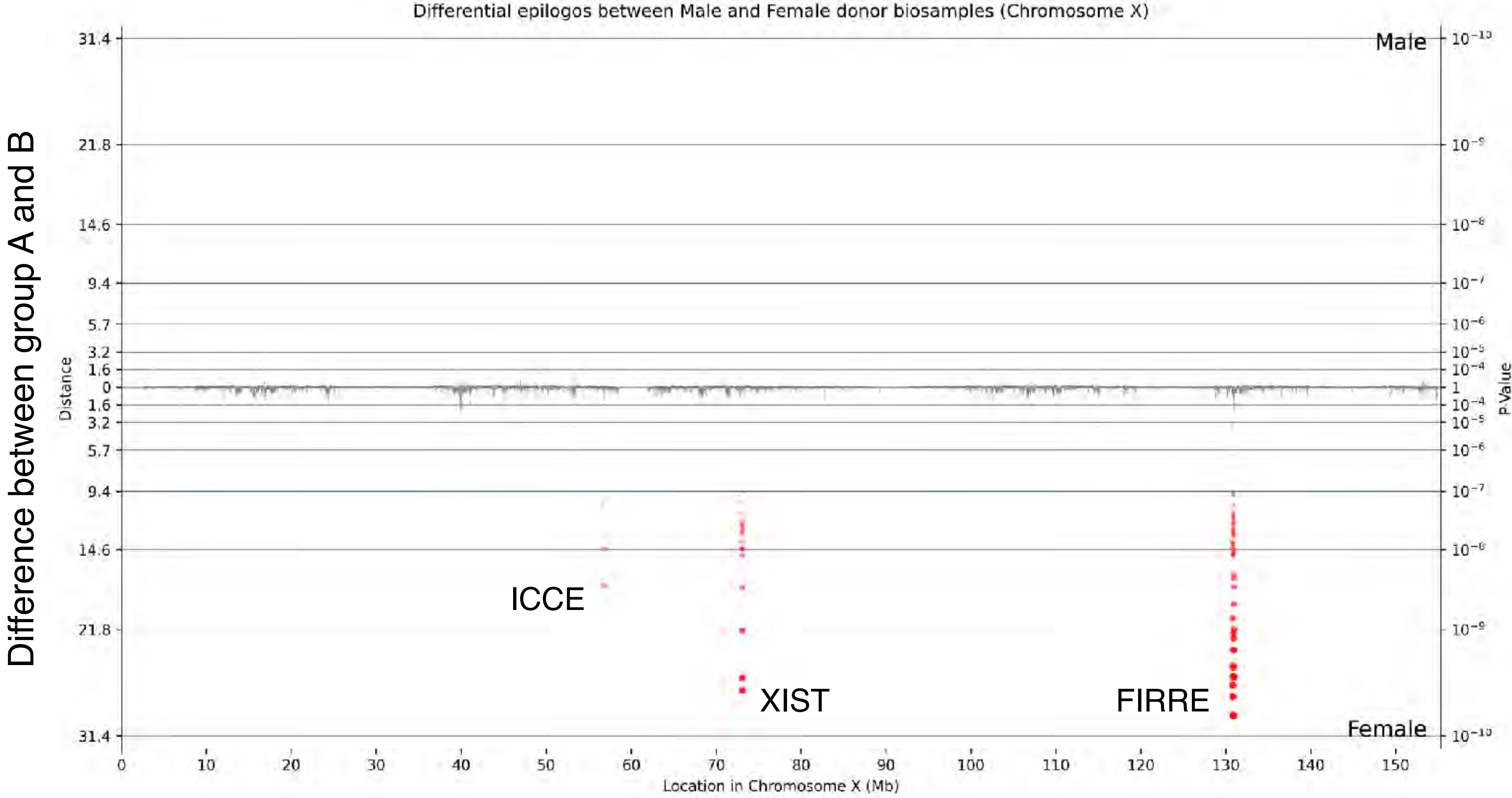
26 Pairwise comparison of groups of biosamples or interest



27 A comparison of male vs. female donors

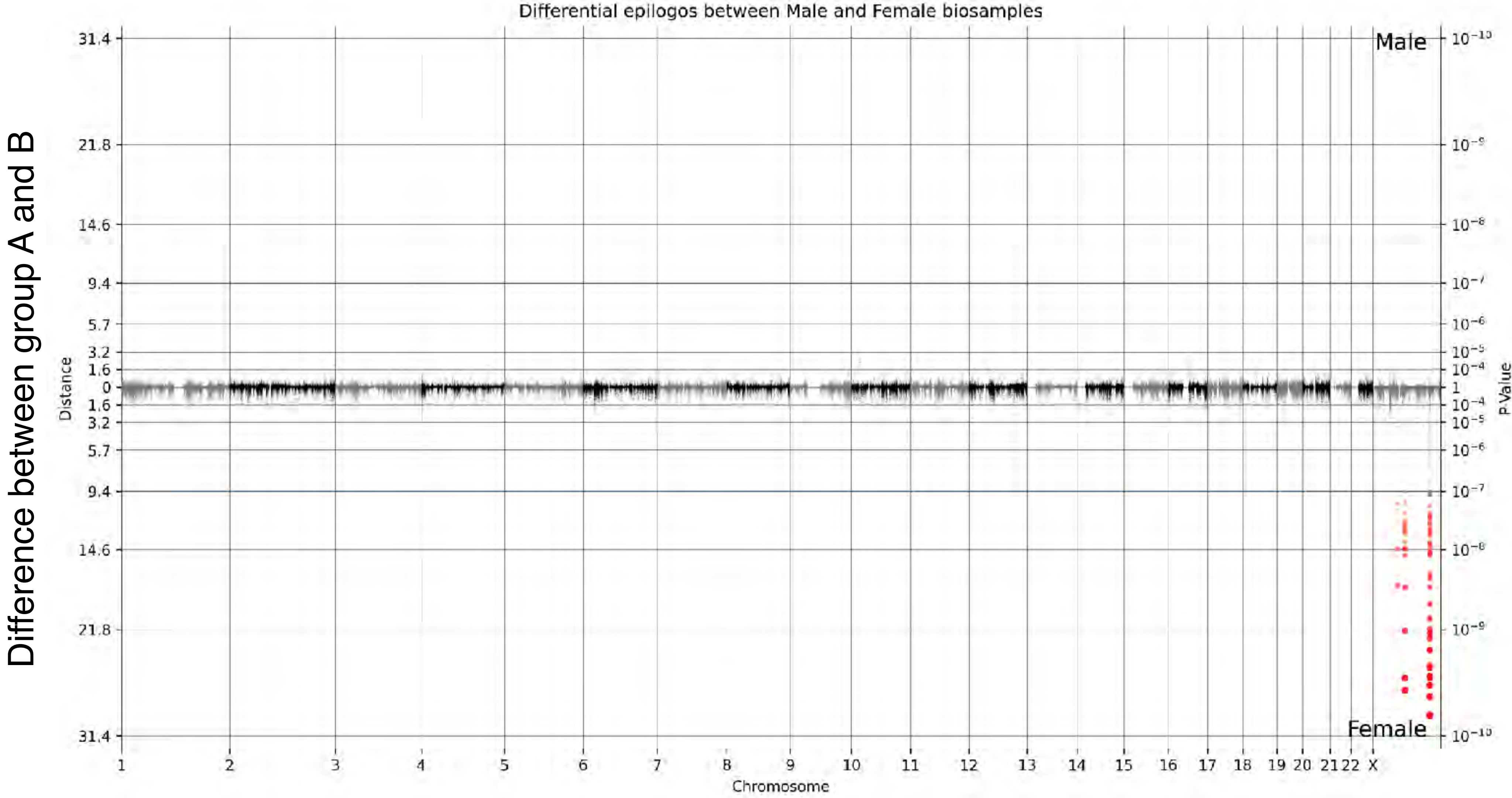


28 Male vs. female donors (total of 684 biosamples)



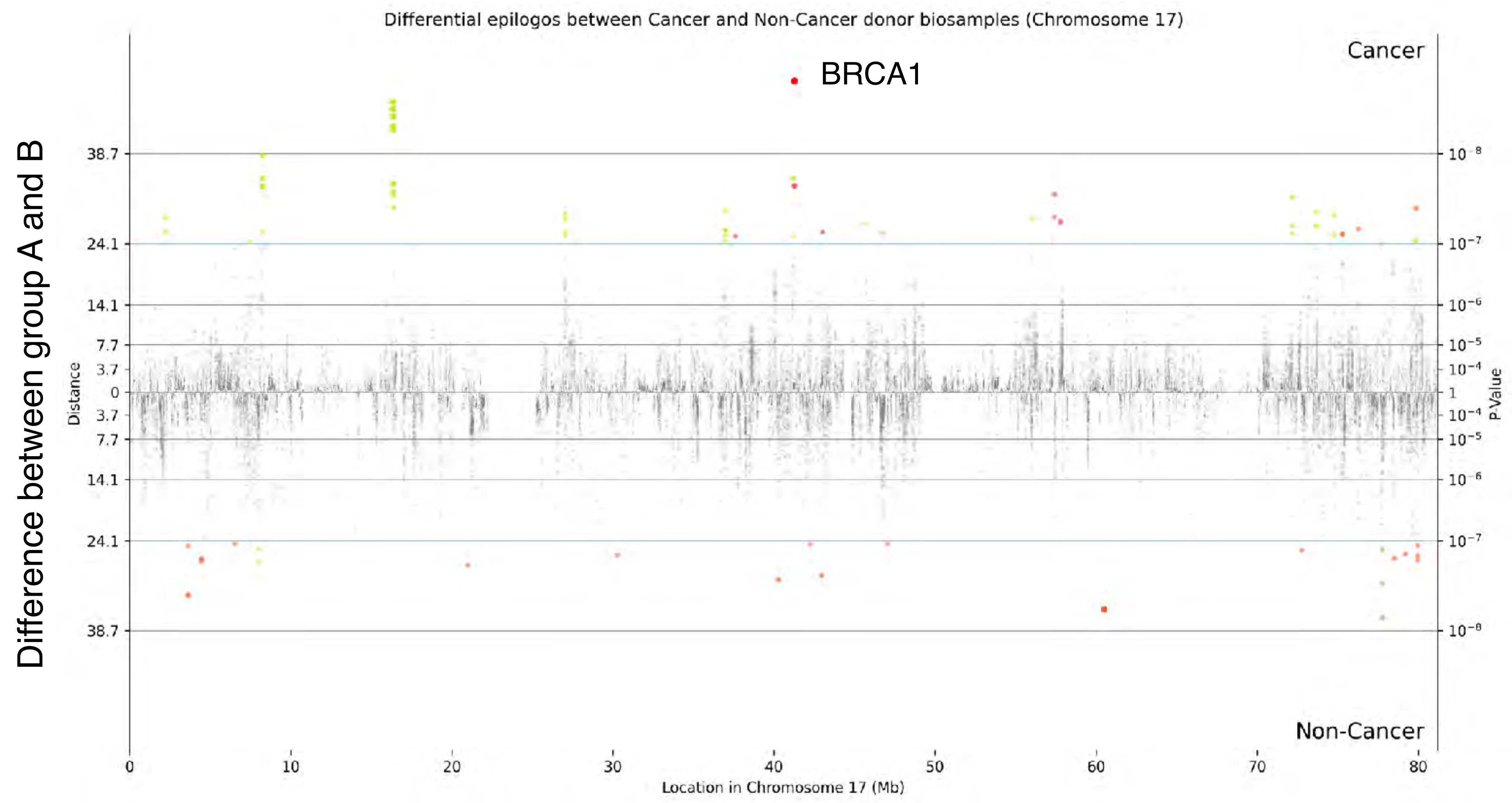
1% FWER

29 Male vs. female donors (total of 684 biosamples)



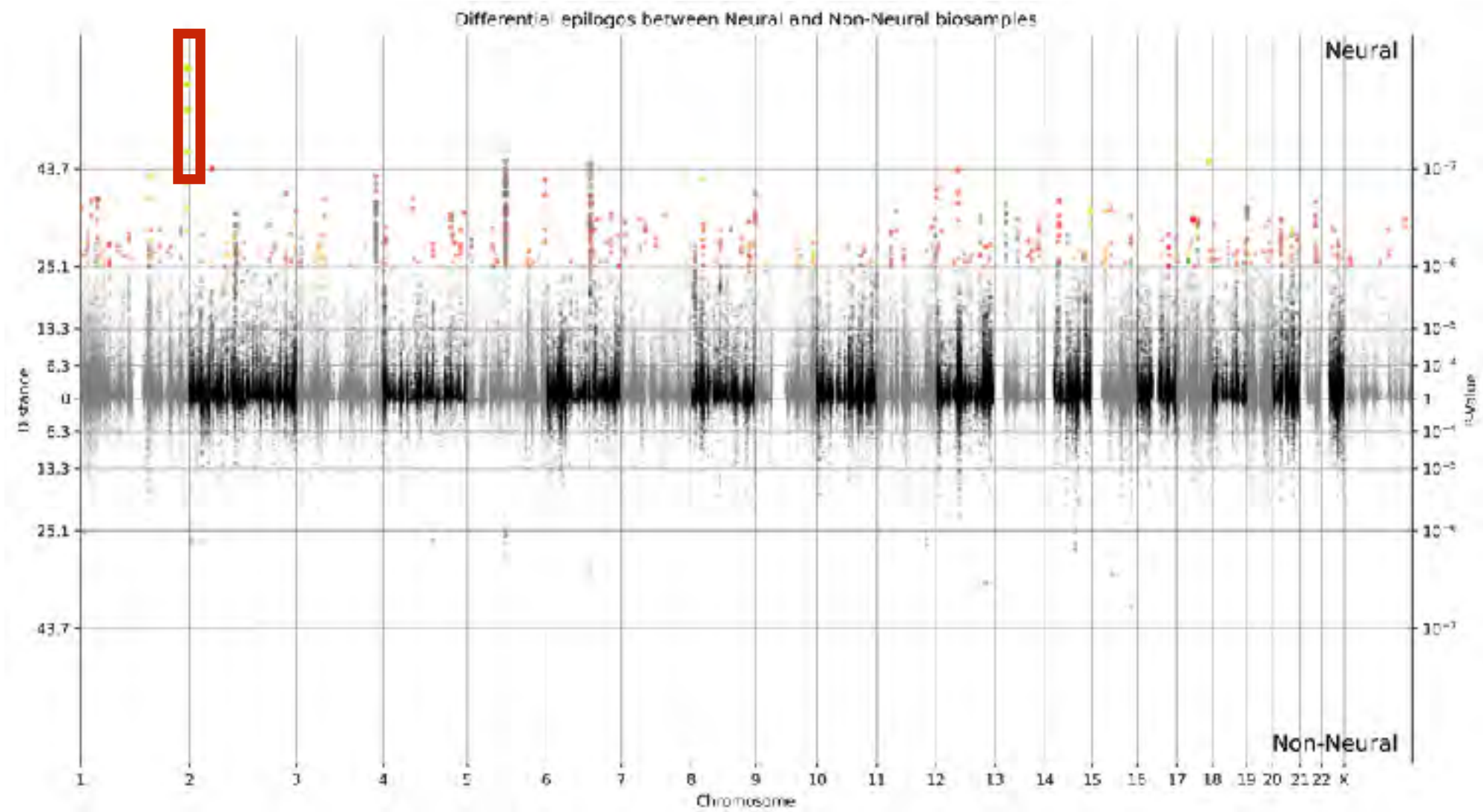
1% FWER

31 Cancer vs. non-cancer biosamples (chr17)



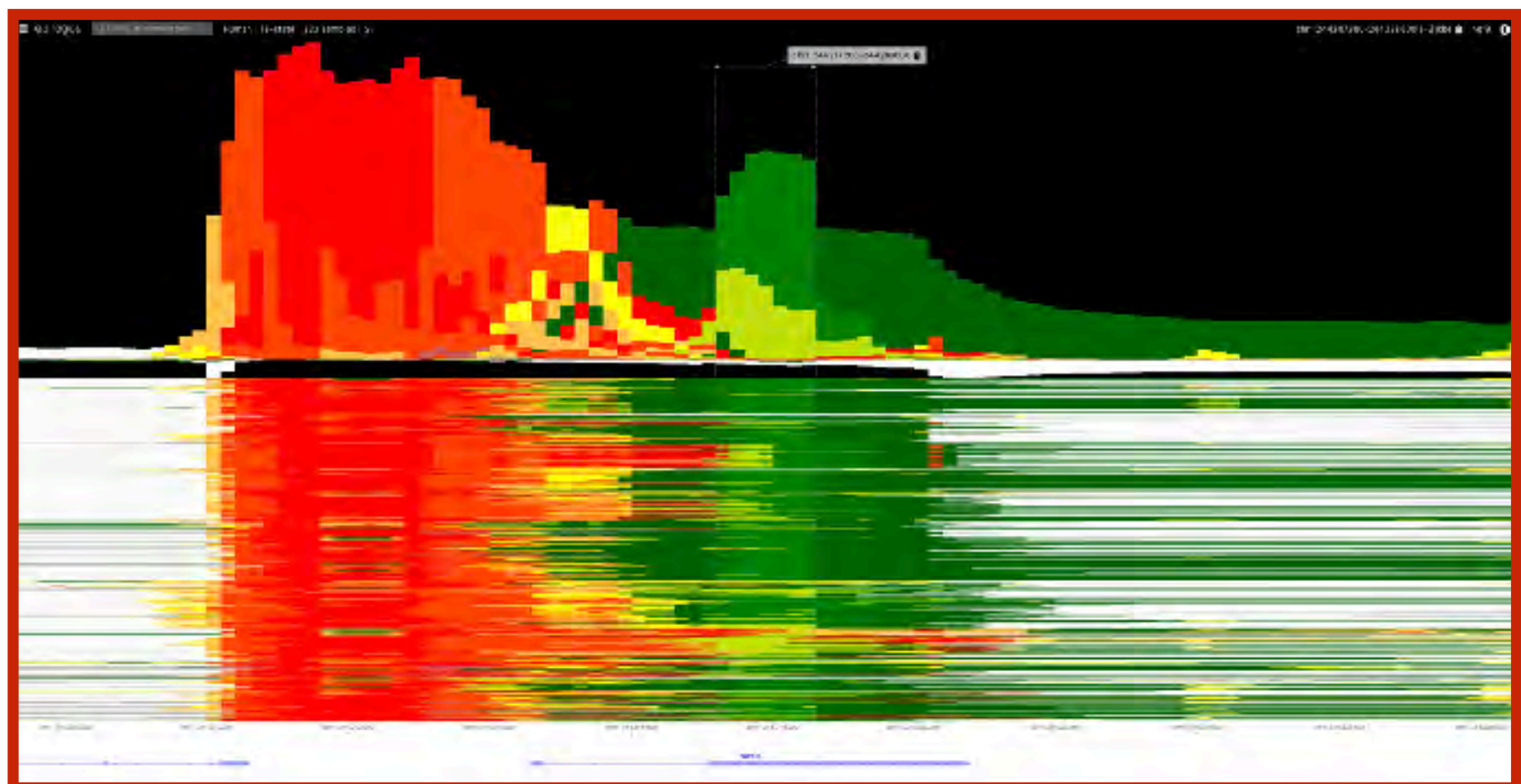
1% FWER

neural vs. non-neural

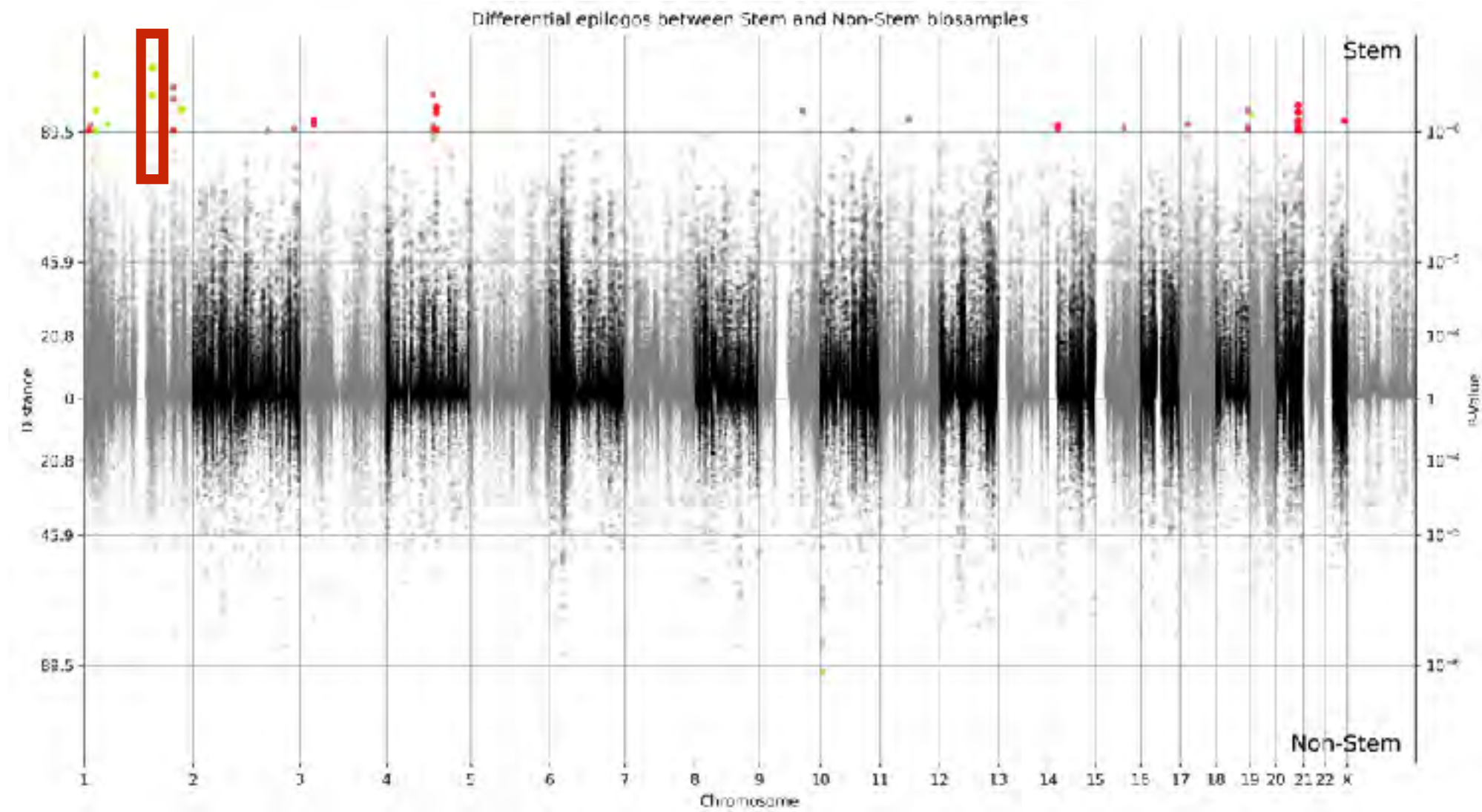


ZBTB18

transcriptional repressor involved in neuronal development

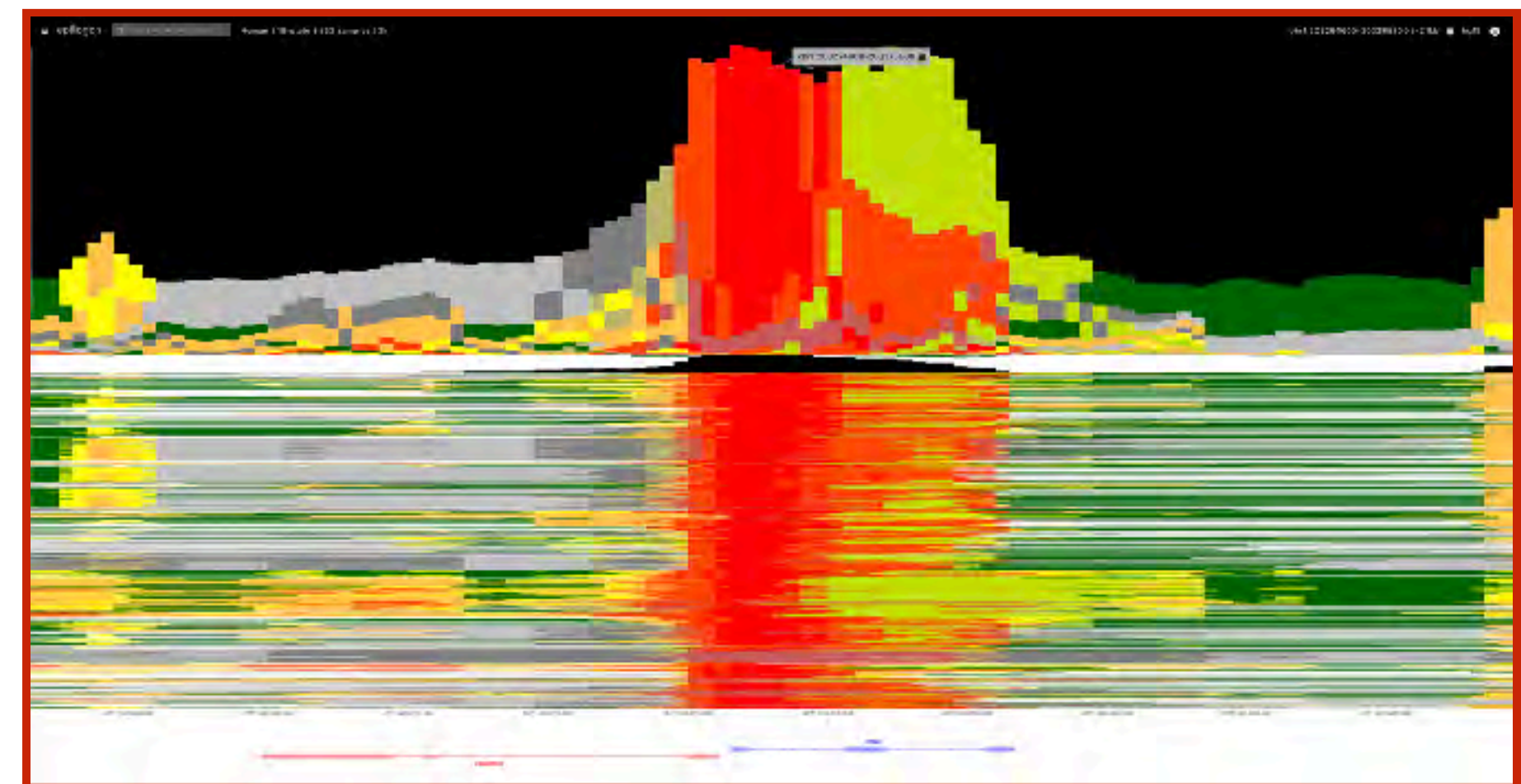


stem cell-like vs. other



BTG2

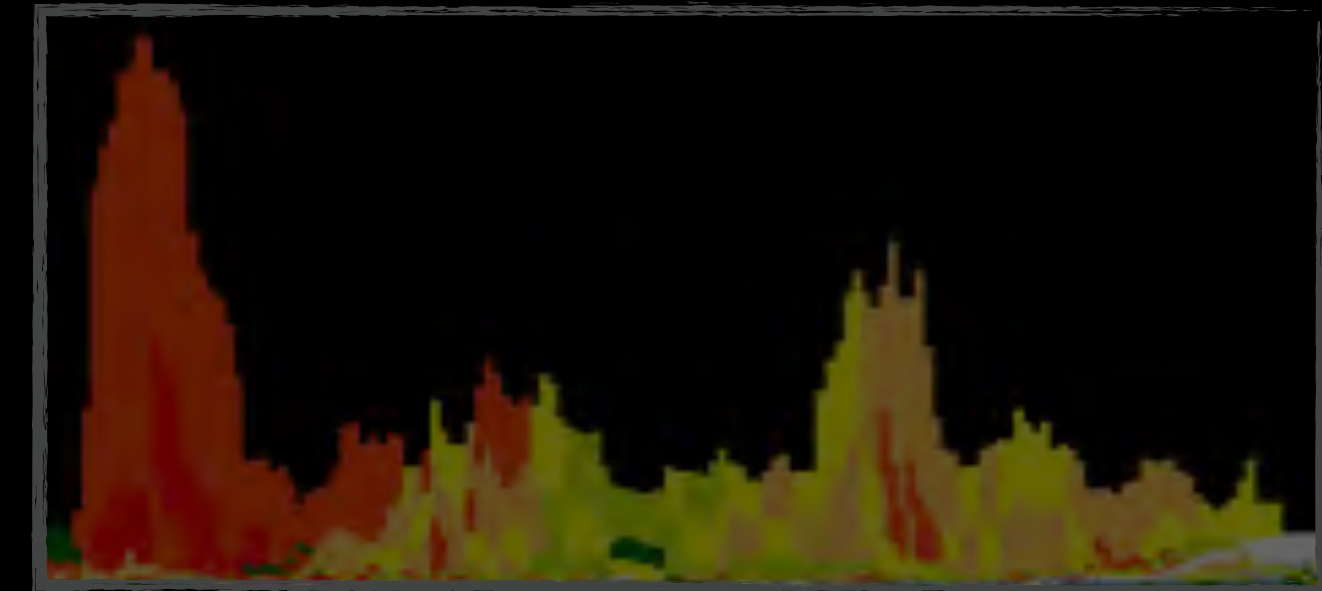
regulation of the G1/S transition of the cell cycle



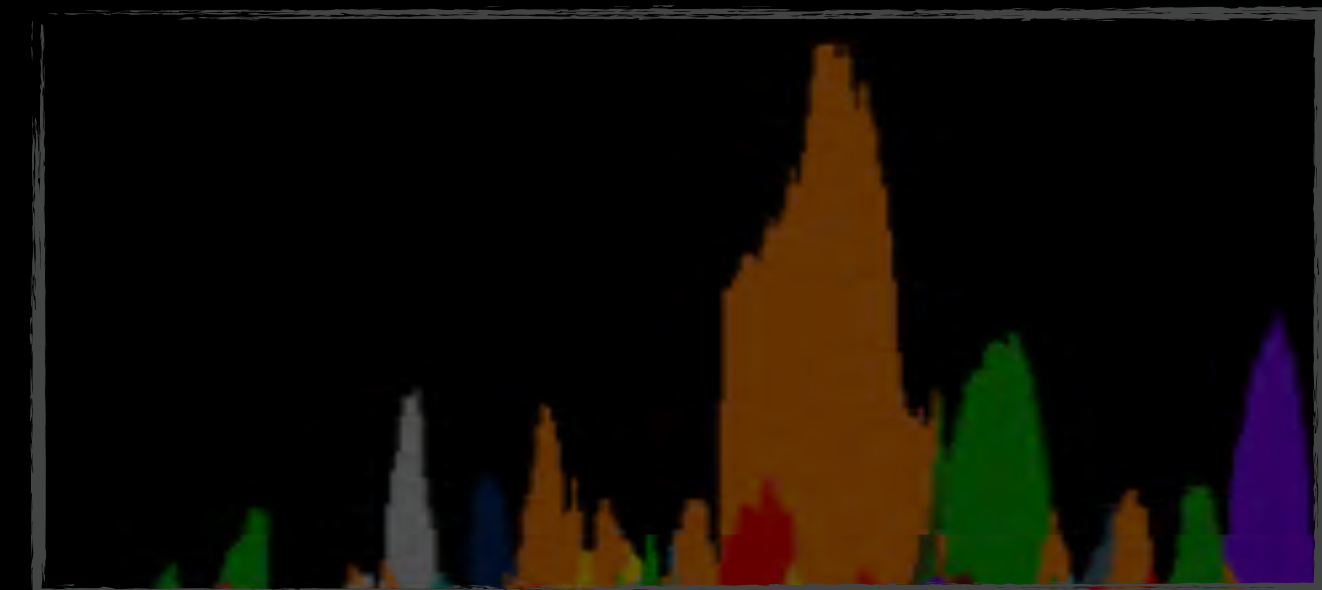
10% FWER

In search of 'relevance': two types of genomic annotations

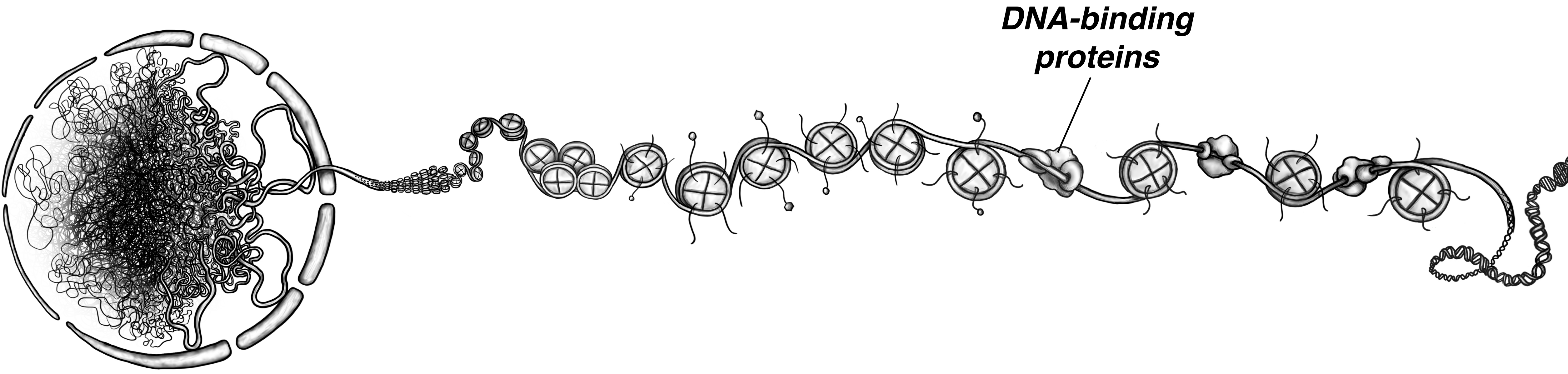
- **Chromatin states (epilogos):**
“What type of functionality does a genomic region encode?”
(e.g. **promoter**, **enhancer**, repressor)
- **Chromatin accessibility (DHS Index):**
“In which cellular contexts are regulatory regions utilized?”
(e.g. **cardiac**, **lymphoid**, **neural**)



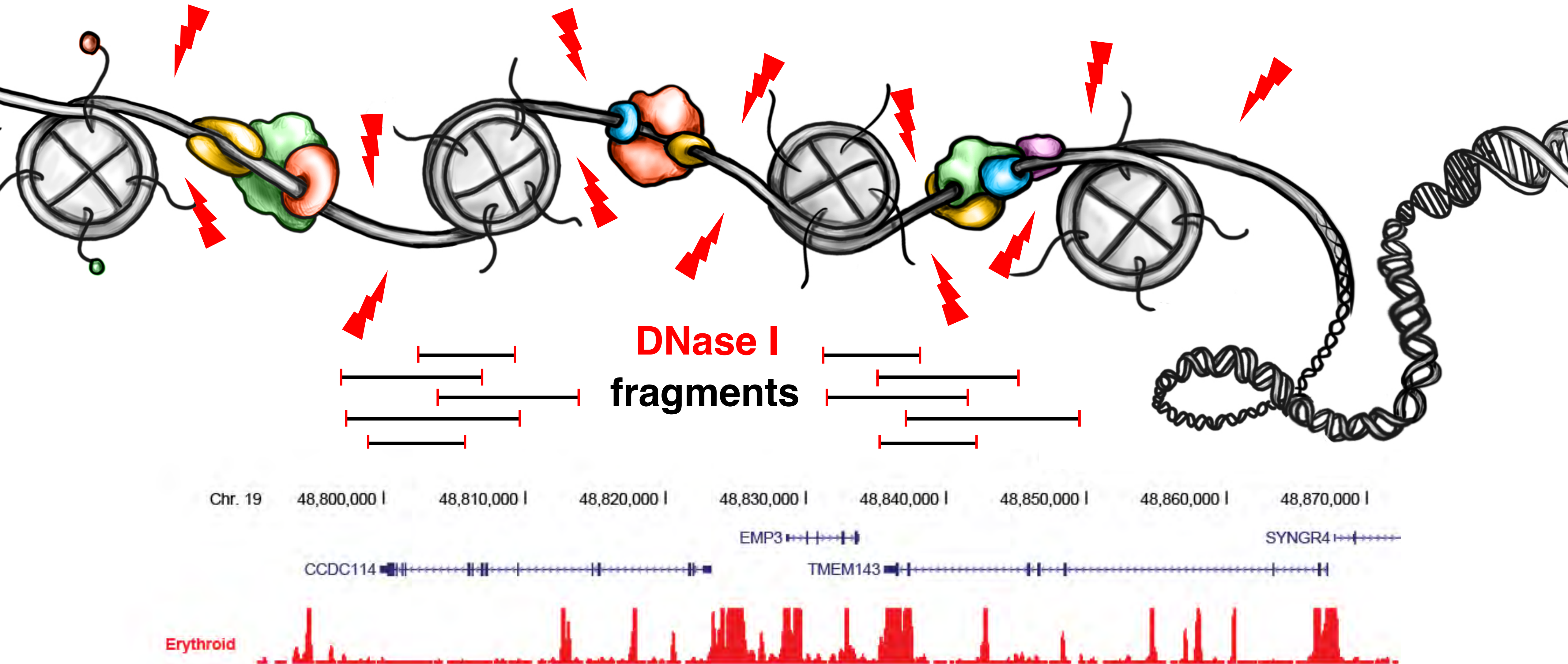
<https://epilogos.net>



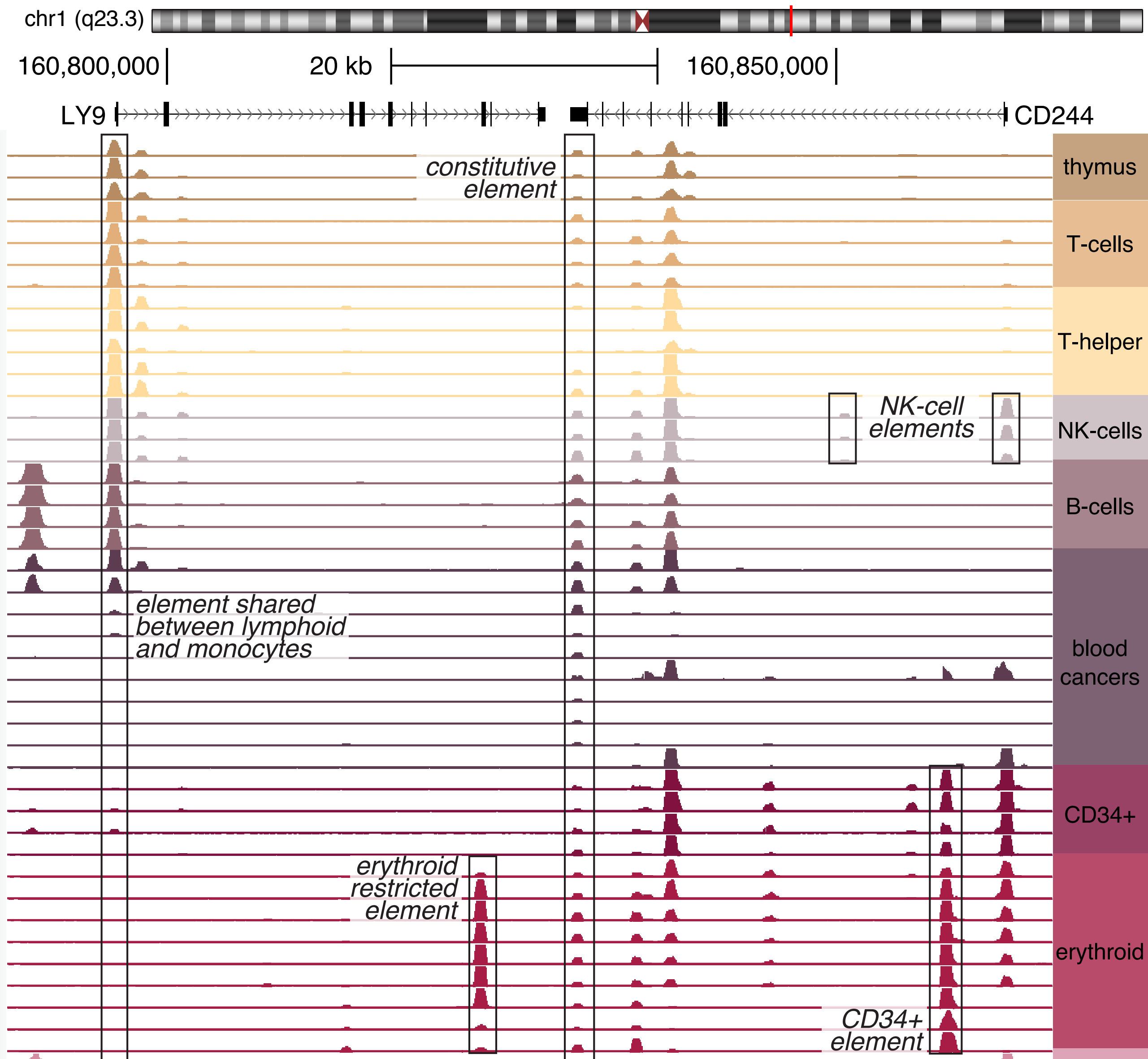
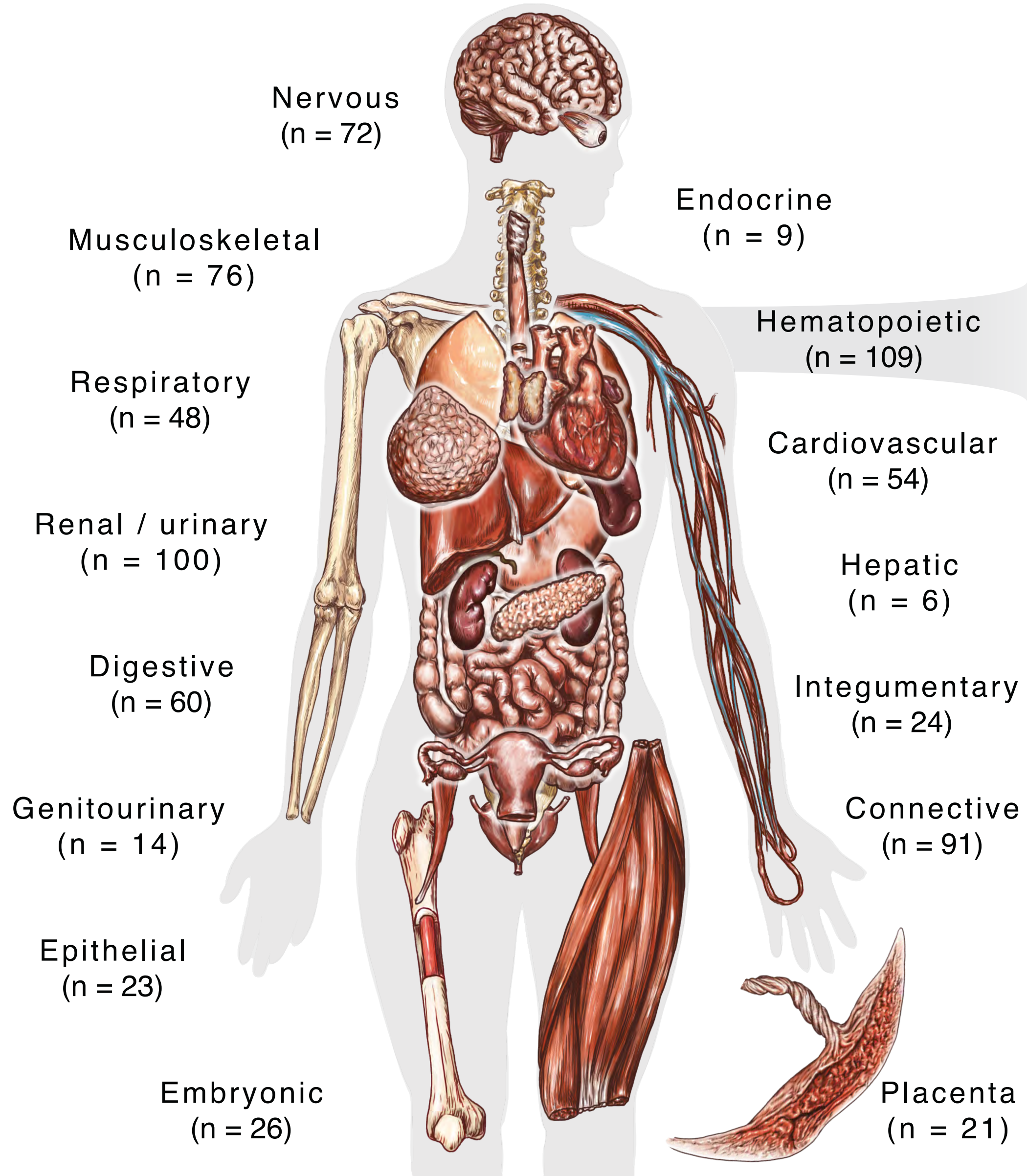
Meuleman *et al.*, 2020 & ongoing



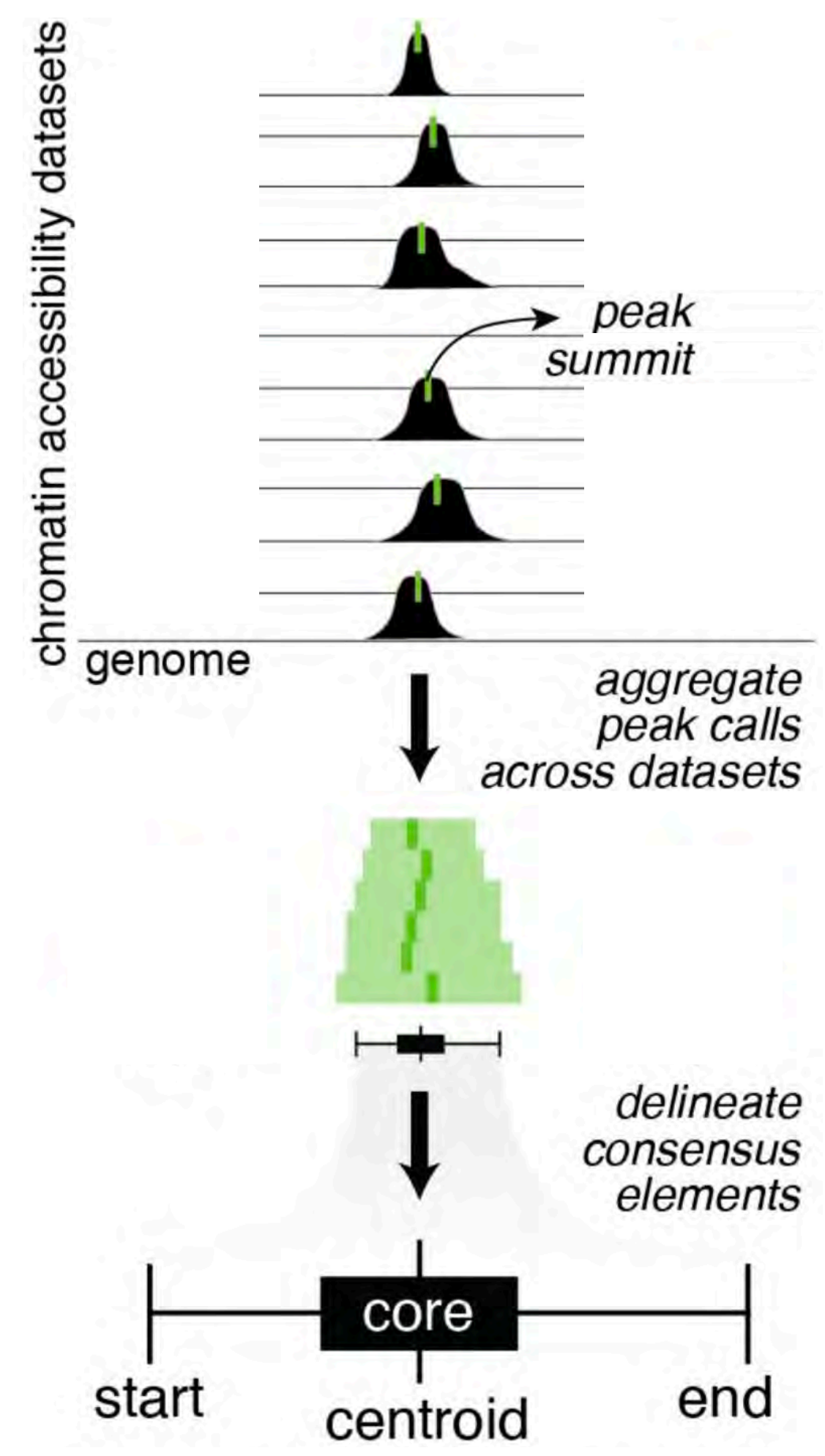
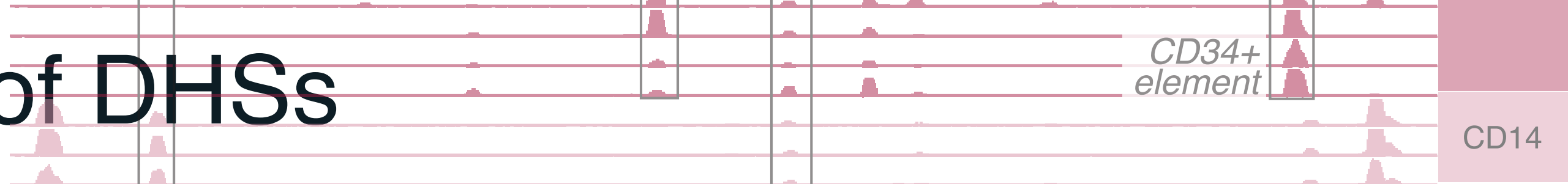
The regulatory genome can be mapped using DNase I digestion



A survey of chromatin accessibility across 400+ cell types and states

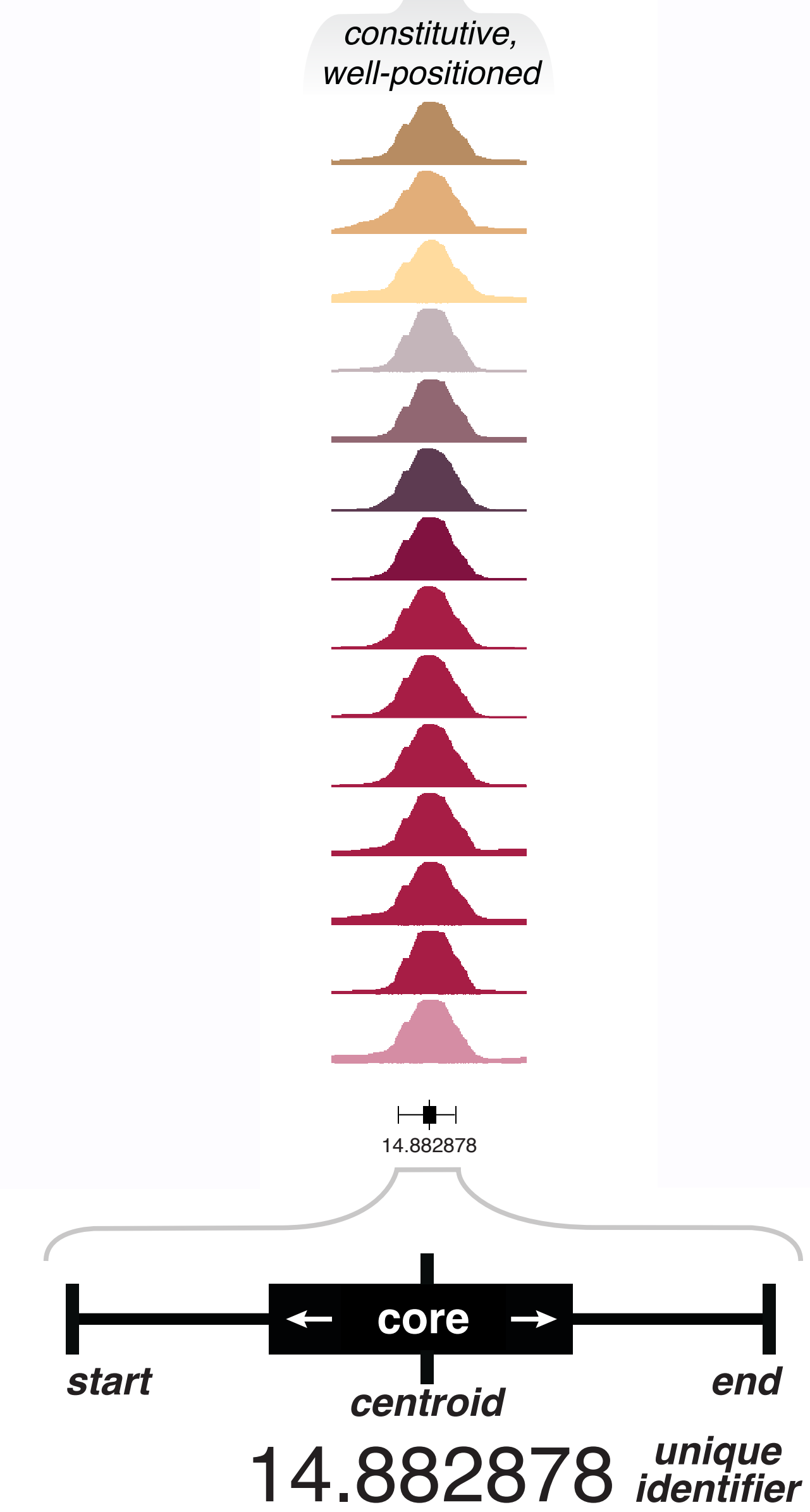


38 A common coordinate system of DHSs

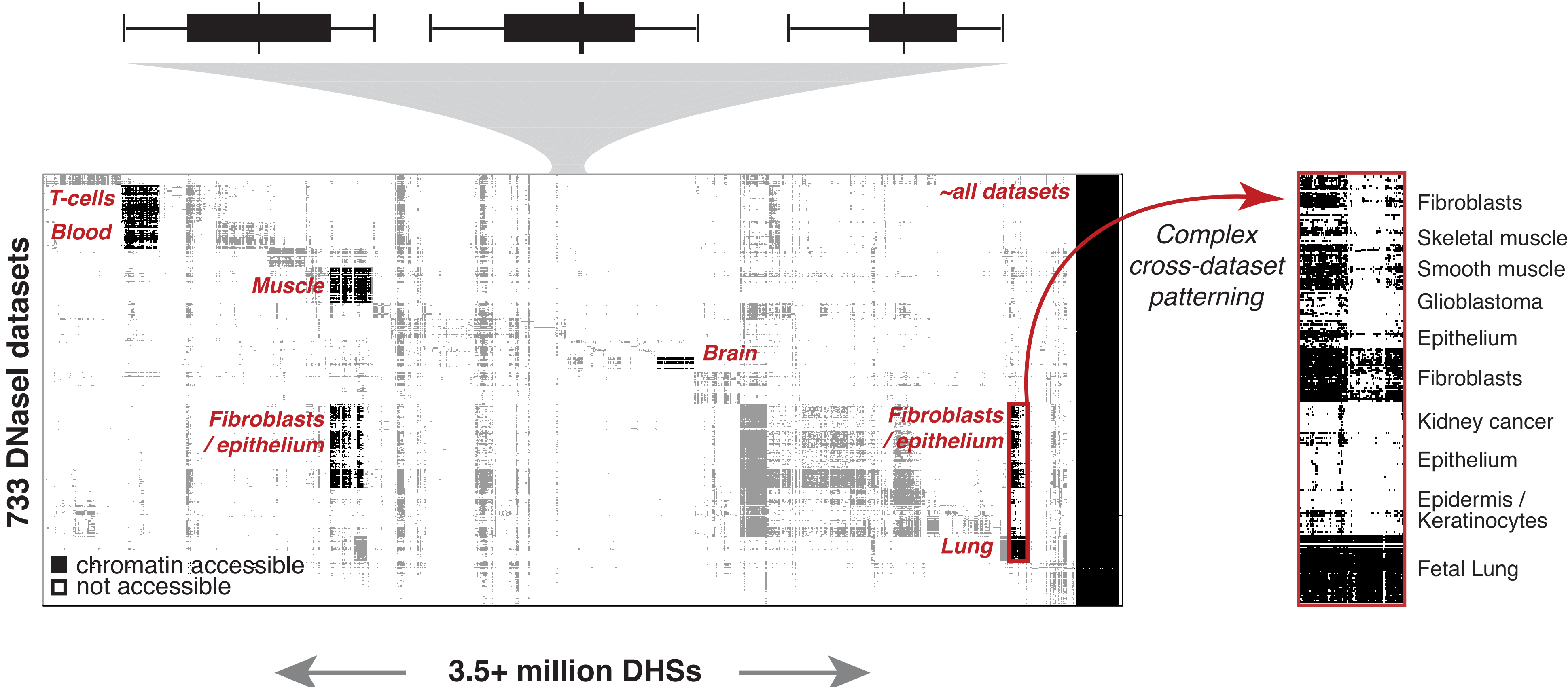


3.5M+ consensus elements

- Thymus
- CD3+
- T2 helper cells, induced
- CD56+
- CD20+
- OCI Ly7
- CD34+
- CD34 (day 4)
- CD34 (day 6)
- CD34 (day 8)
- CD34 (day 15)
- CD34 (day 17)
- CD34 (day 18)
- CD14+

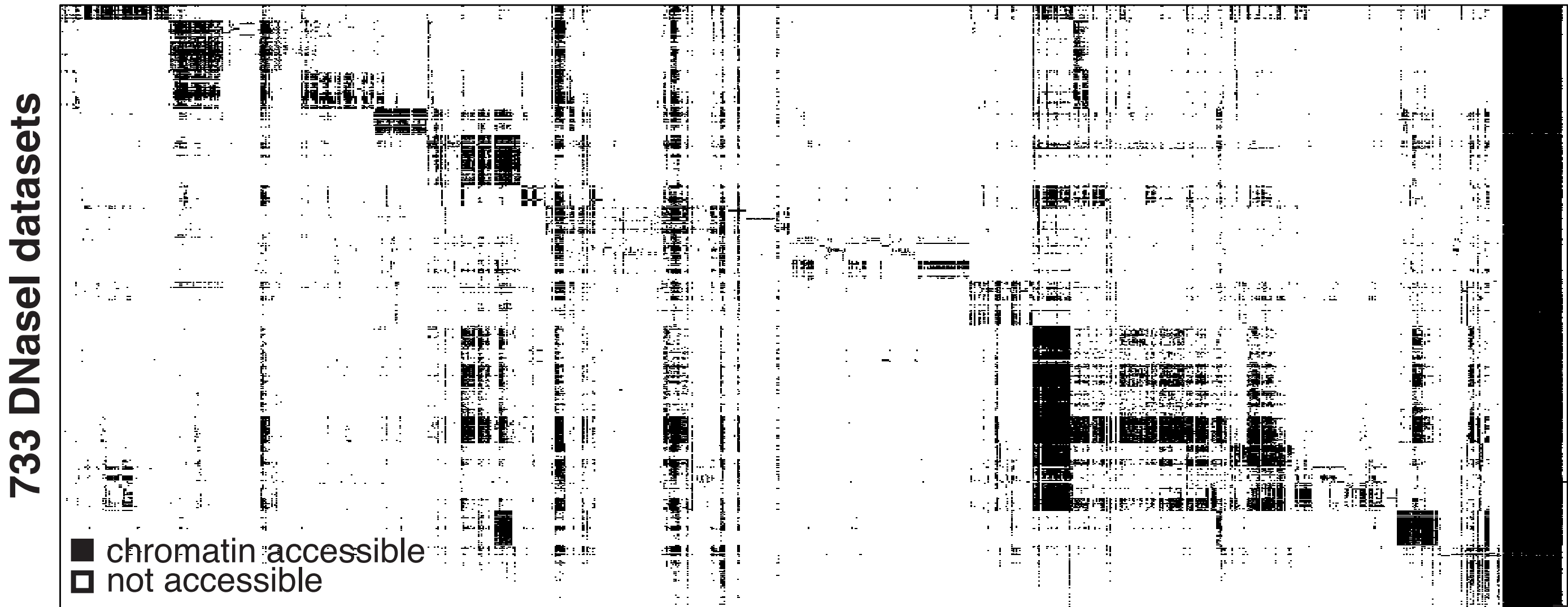


39 Complex DHS patterning across cell types and states



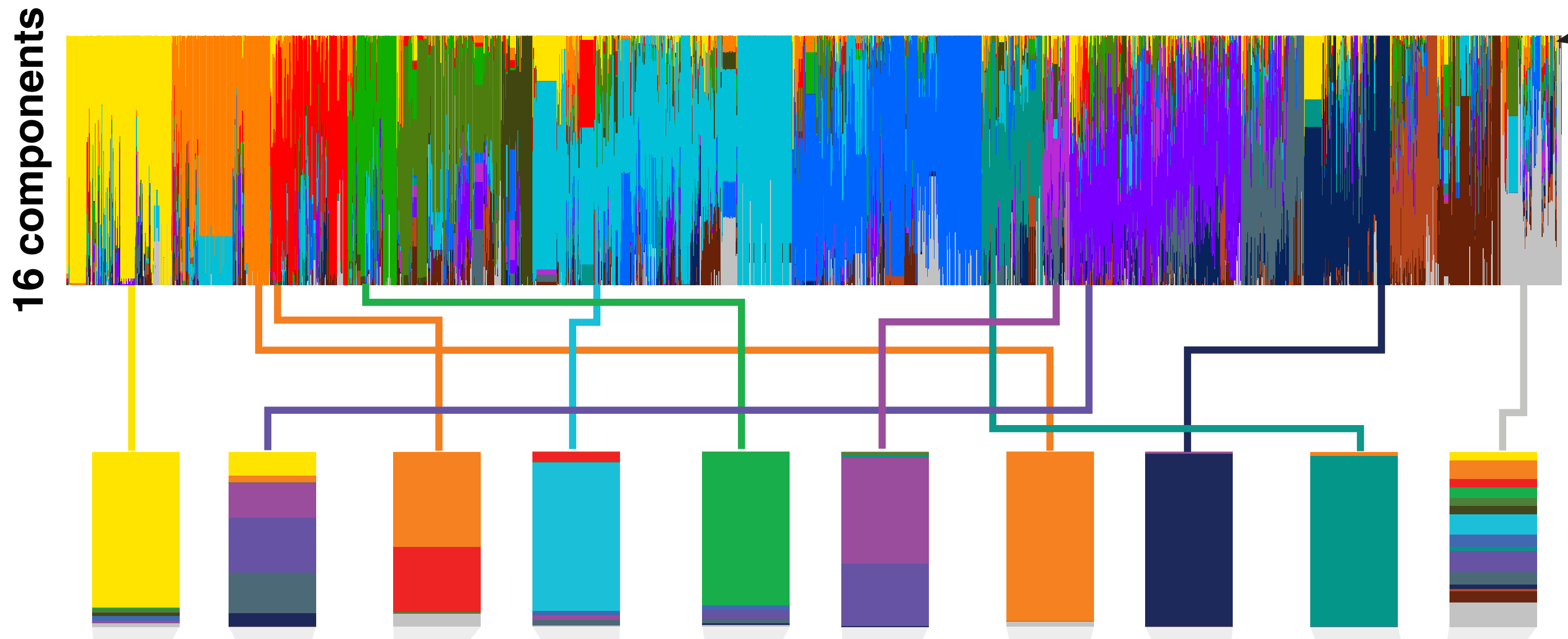
Although generally quite cell type/state selective, DHSs are often shared between broader cellular contexts

40 DHS patterns can be decomposed into *components*



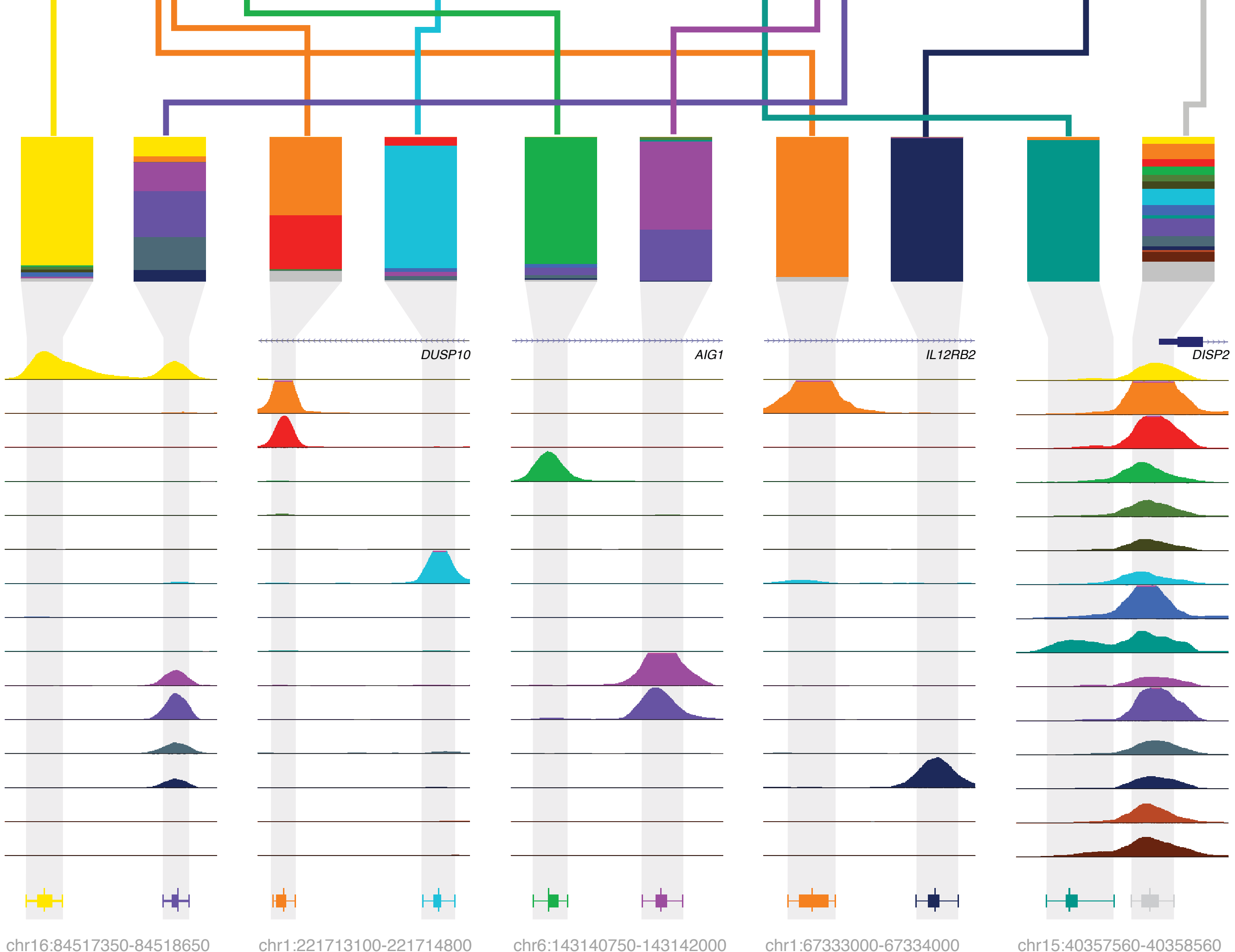
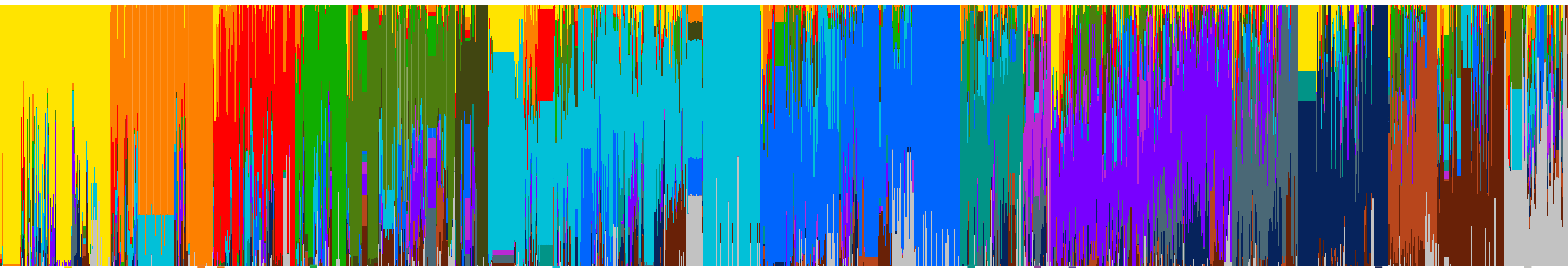
← 3.5+ million DHSs →

Non-Negative Matrix Factorization



Each DHS is described by a mixture of components

Components reflect distinct biological contexts and regulatory signal



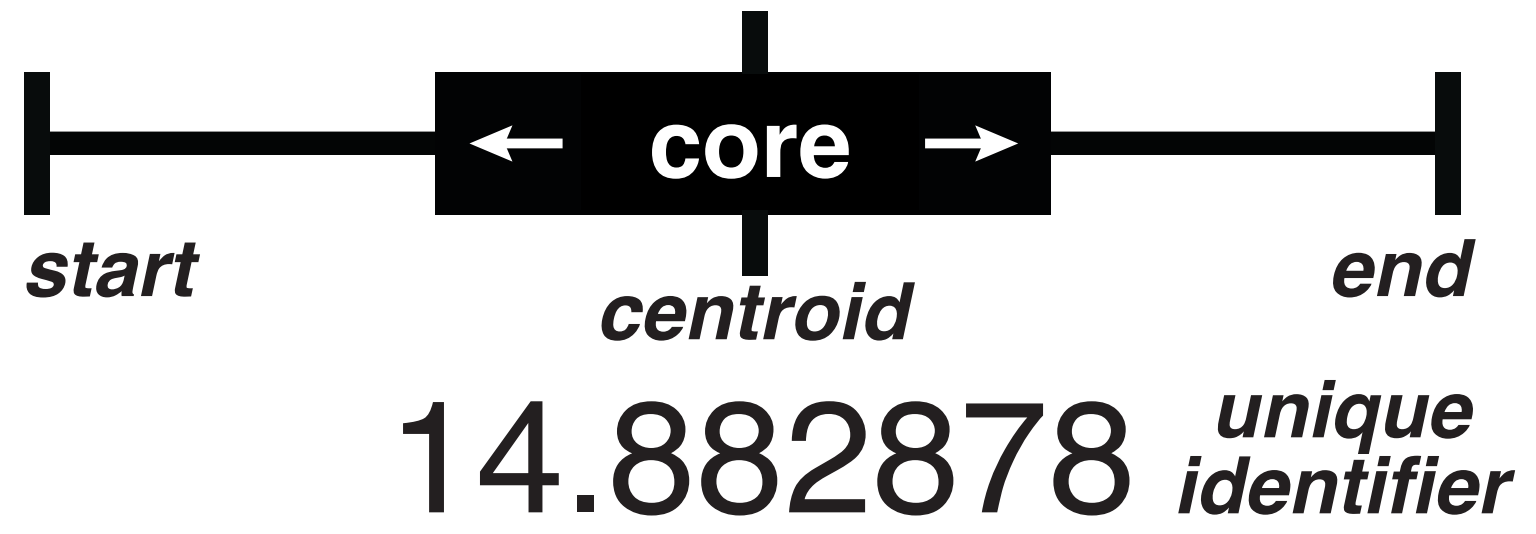
- Datasets**
- Yellow: Trophoblasts
 - Orange: CD8 cells
 - Red: CD34 cells
 - Green: Fetal heart
 - Olive: Fetal muscle (leg)
 - Dark Green: HMVEC
 - Cyan: hESC
 - Blue: Fetal brain
 - Teal: Large intestine
 - Purple: Skin fibroblasts (leg)
 - Dark Purple: PGP1 fibroblasts
 - Dark Blue-Gray: Renal cell carcinoma
 - Dark Blue: Esophageal epithelial
 - Brown: Fetal lung
 - Dark Brown: Fetal kidney

- Motifs**
- GCM1
 - IRF
 - ETS/SPI1
 - MEF2
 - E-box
 - ETS/ERG
 - Oct-4
 - NeuroD
 - HNF4A
 - AP-1
 - AP-1
 - HNF1
 - TP63
 - FOX
 - PAX2

- DHS Vocabulary**
- Yellow: Placental
 - Orange: Lymphoid
 - Red: Myeloid / erythroid
 - Green: Cardiac
 - Olive: Musculoskeletal
 - Dark Green: Vascular / endothelial
 - Cyan: Embryonic / primitive
 - Blue: Neural
 - Teal: Digestive
 - Purple: Stromal A
 - Dark Purple: Stromal B
 - Dark Blue-Gray: Renal / cancer
 - Dark Blue: Cancer / epithelial
 - Brown: Pulmonary devel.
 - Dark Brown: Organ devel. / renal
 - Gray: Tissue invariant

42 DHS Index and Vocabulary: a novel annotation of regulatory DNA

DHS Index



3.5M+ DHSs
Richly annotated and
indexed across cell types

DHS Vocabulary

- Placental
- Lymphoid
- Myeloid / erythroid
- Cardiac
- Musculoskeletal
- Vascular / endothelial
- Embryonic / primitive
- Neural
- Digestive
- Stromal A
- Stromal B
- Renal / cancer
- Cancer / epithelial
- Pulmonary devel.
- Organ devel. / renal
- Tissue invariant

index.altius.org

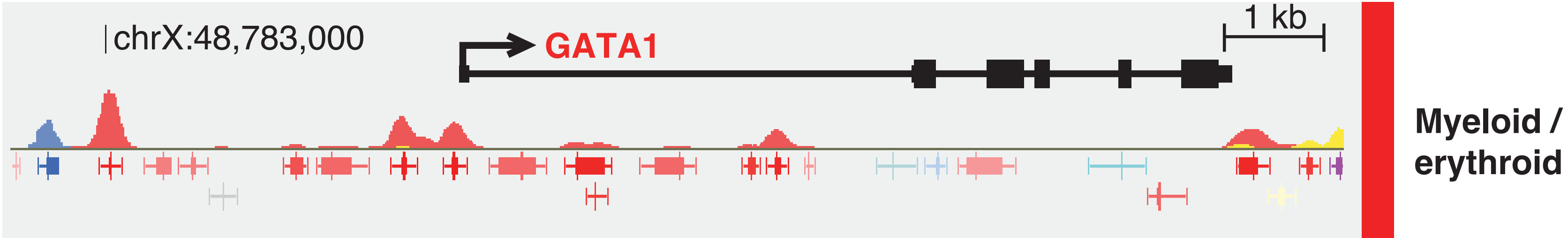
Go

I'm Feeling Lucky

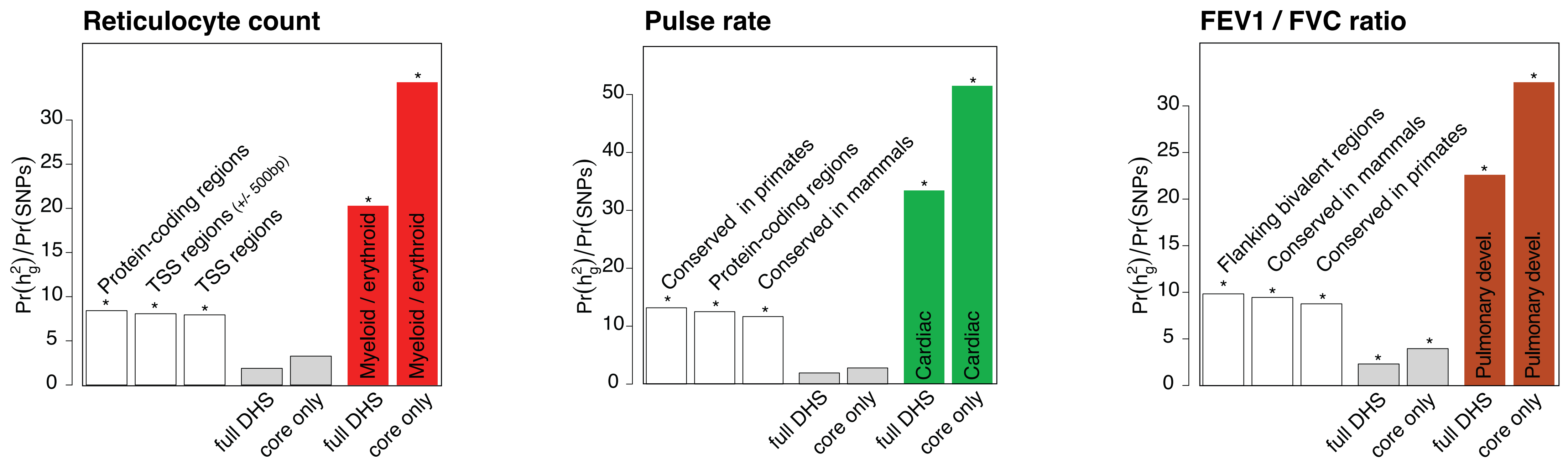
e.g., use query terms like HGNC symbols (HOXA1, NFKB1, etc.) or genomic regions (chr17:41165790-41317987, etc.)

How can we use these data to further identify regions of “relevance”?

44 Regions around genes show component-specific patterning

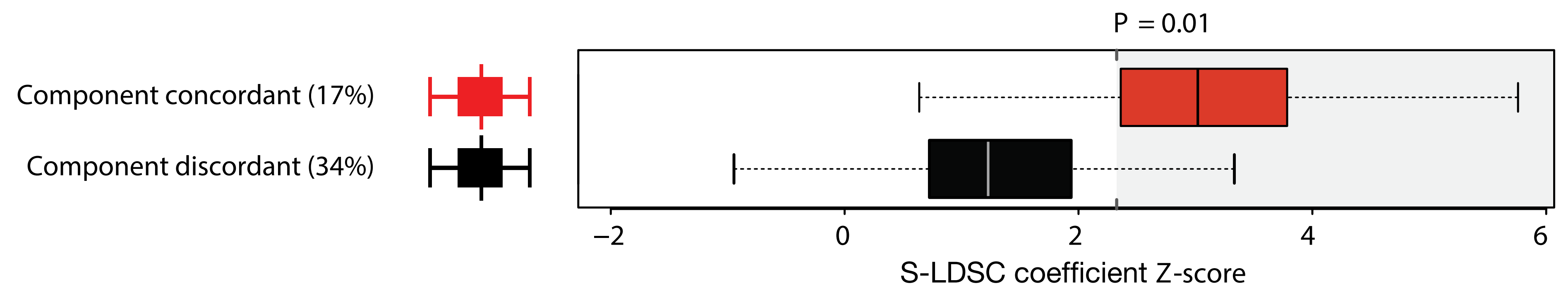
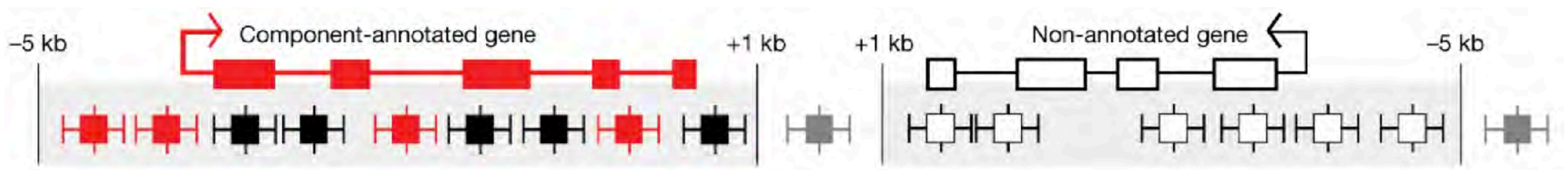


45 GWAS signal is strongly enriched in relevant component-associated DHSs

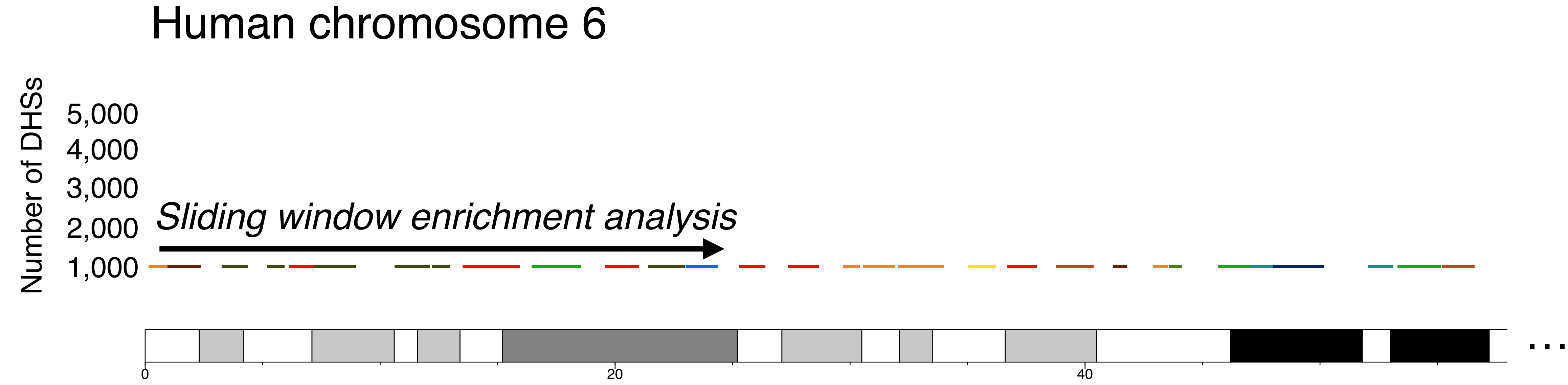


Strong heritability enrichment in relevant DHS components, relative to all DHSs or 85 other genome-wide annotations

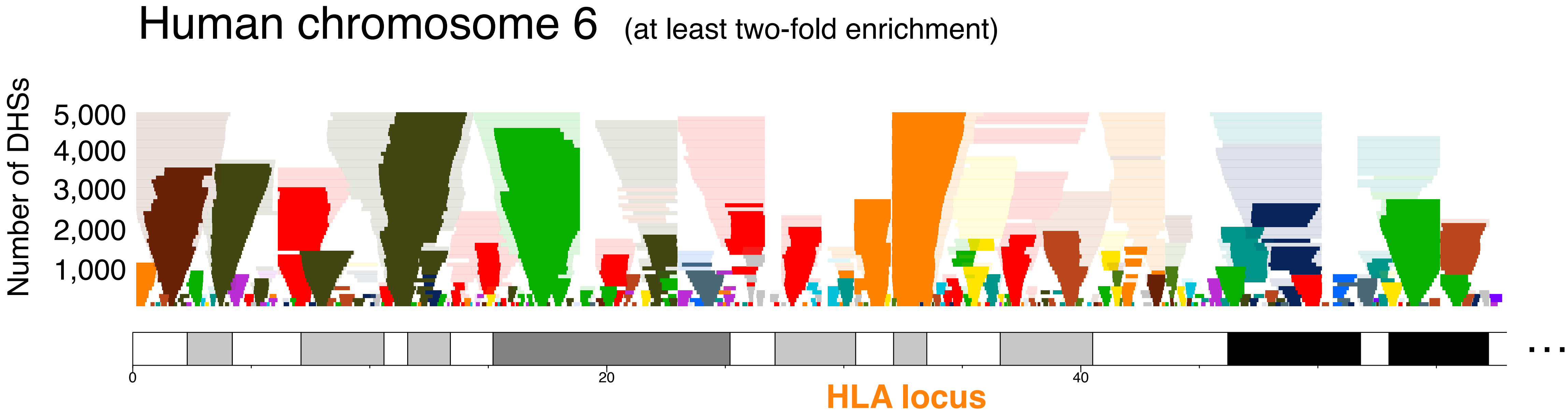
46 GWAS signal is spread across congruently annotated genic DHSs



47 Domain-level organization of component-associated DHSs



48 Domain-level organization of component-associated DHSs



The human leukocyte antigen (HLA) super-locus is a genomic region on chromosome 6 that encodes genes with important roles in the regulation of the immune system

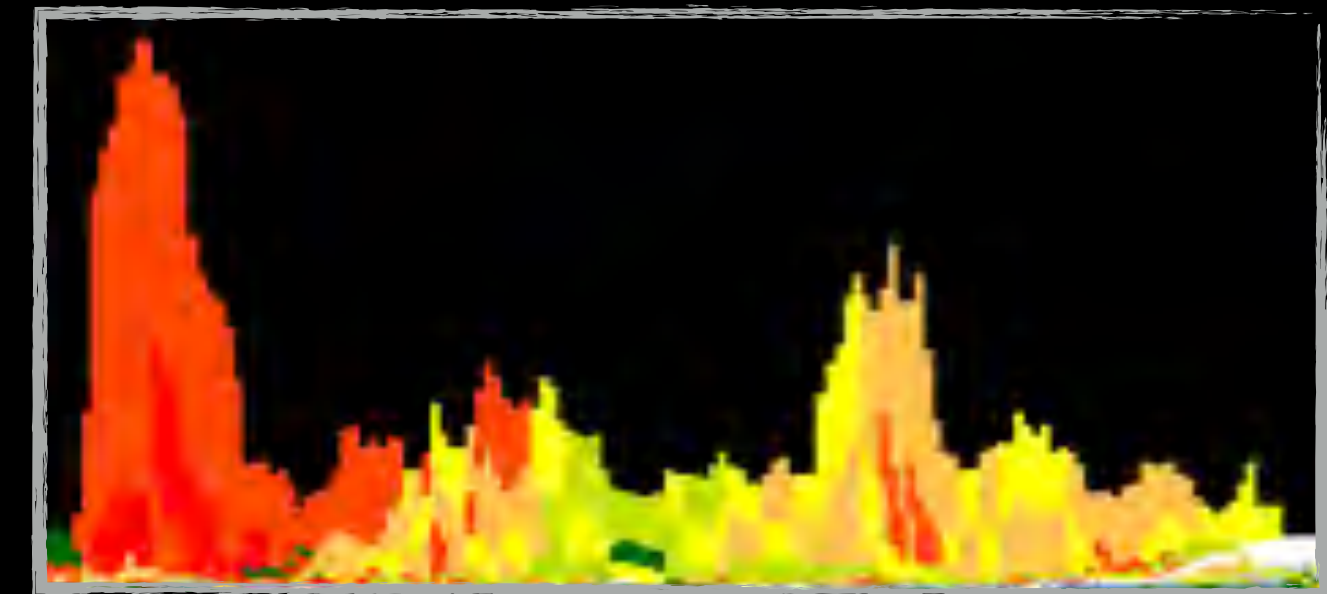
Domain-level organization of component-associated DHSs



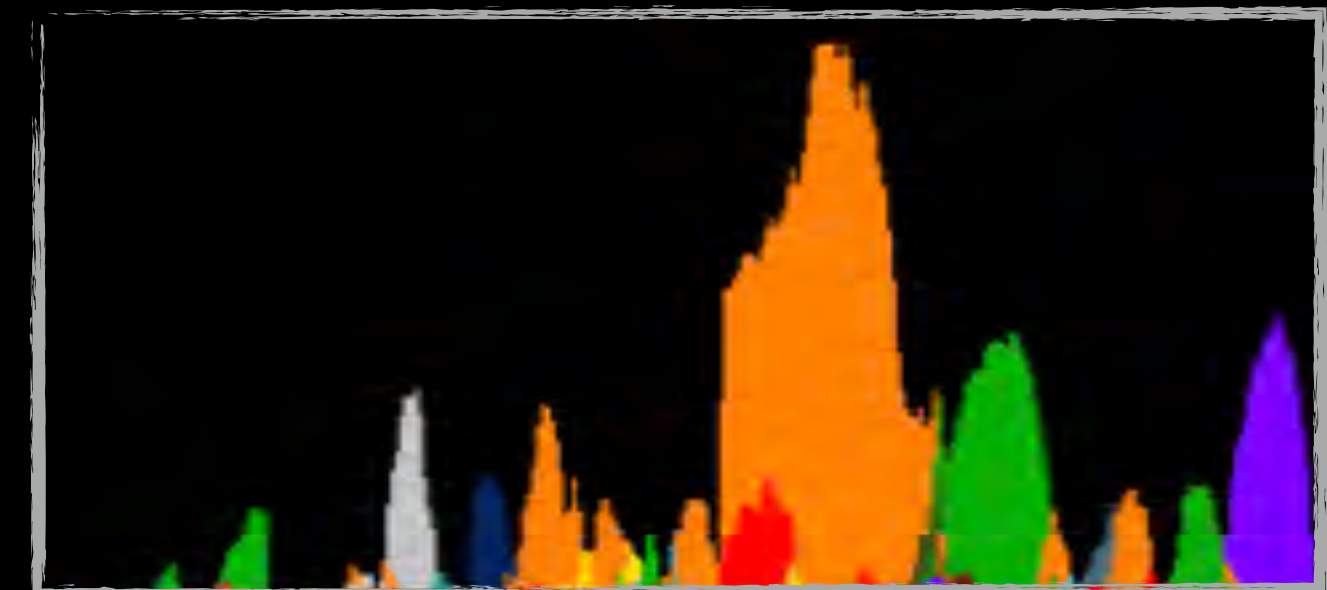
100s of multi-DHS “tornado” domains genome-wide

In search of 'relevance': two types of genomic annotations

- **Chromatin states (epilogos):**
“What type of functionality does a genomic region encode?”
(e.g. **promoter**, **enhancer**, repressor)
- **Chromatin accessibility (DHS Index):**
“In which cellular contexts are regulatory regions utilized?”
(e.g. **cardiac**, **lymphoid**, **neural**)

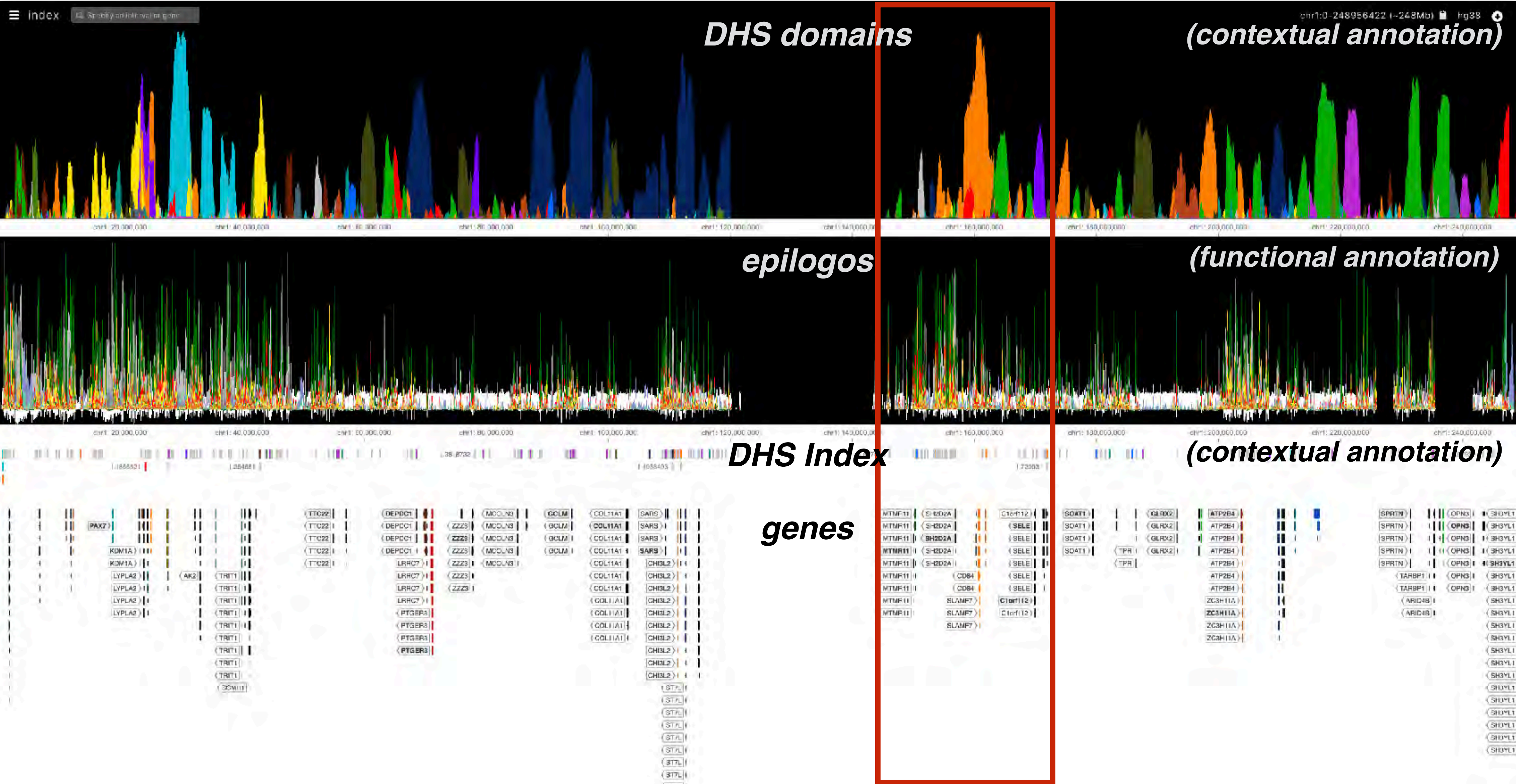


<https://epilogos.net>



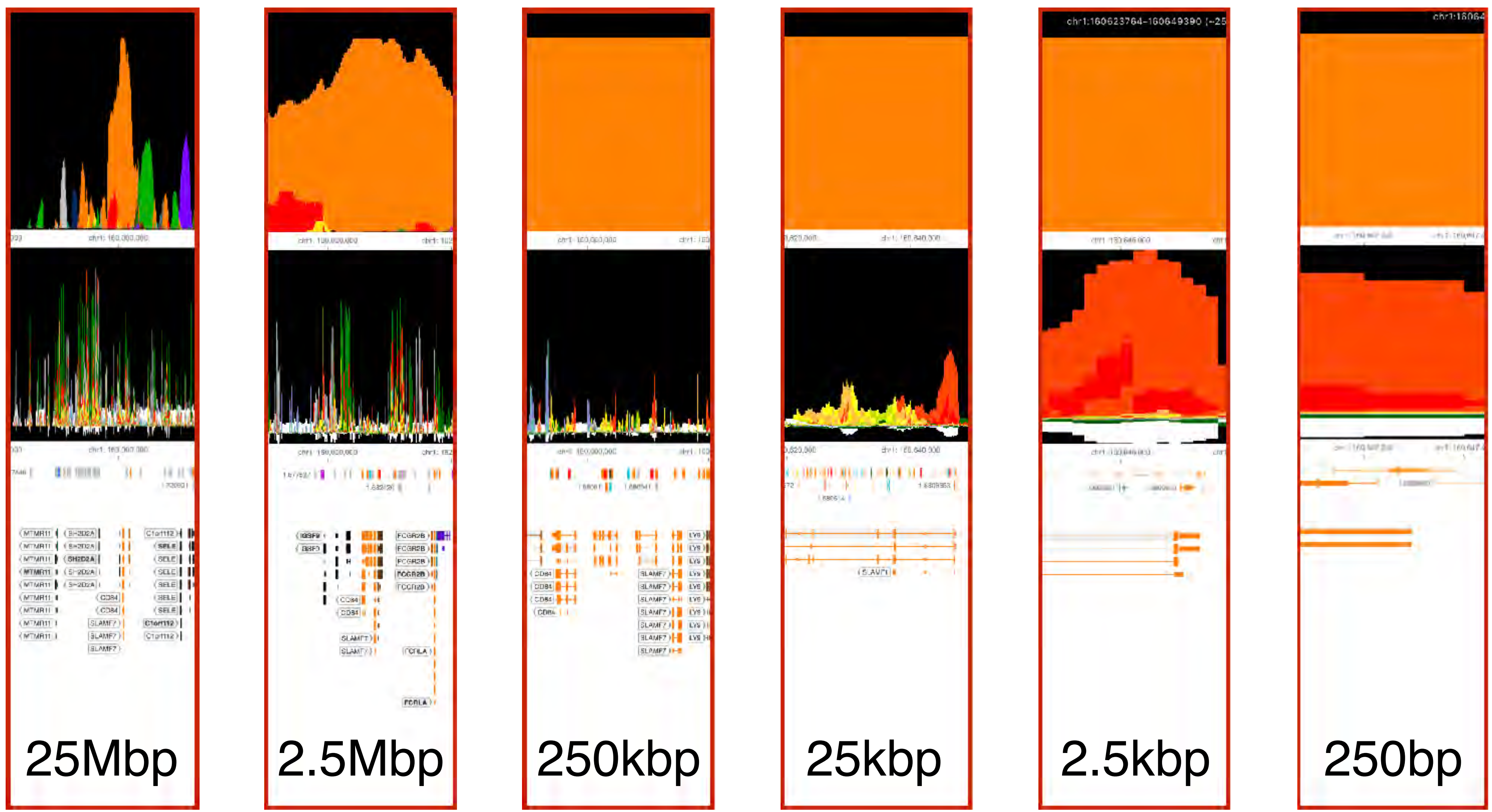
Meuleman *et al.*, 2020 & ongoing

51 Towards multi-scale regulatory reference annotations



Entirety of human chromosome 1 (~250Mbp)

59 Different annotations for different scales



25Mbp

2.5Mbp

250kbp

25kbp

2.5kbp

250bp

DHS domains

epilogos

DHS Index

...and beyond
(footprints, motifs,
genetic variation)

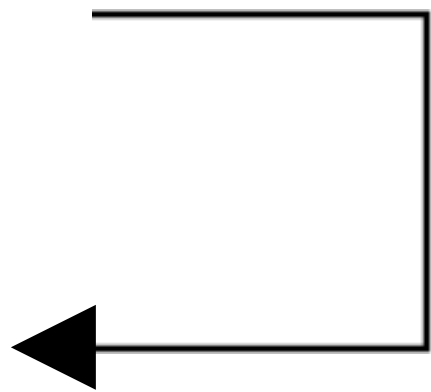
“But where’s my Disney map?”



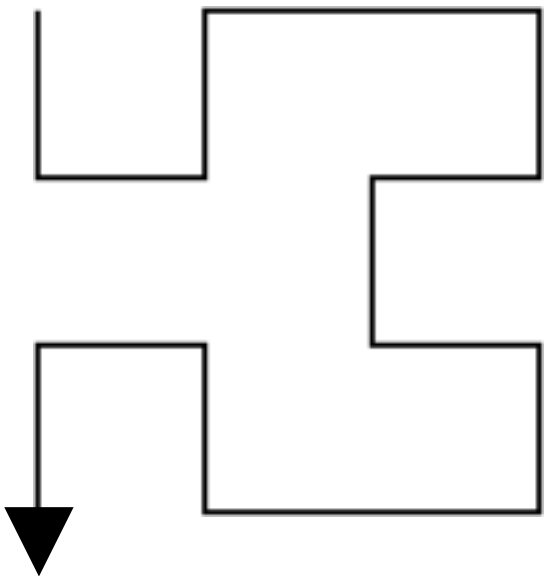
Mapping the linear 1D genome to 2D using Hilbert curves

the genome

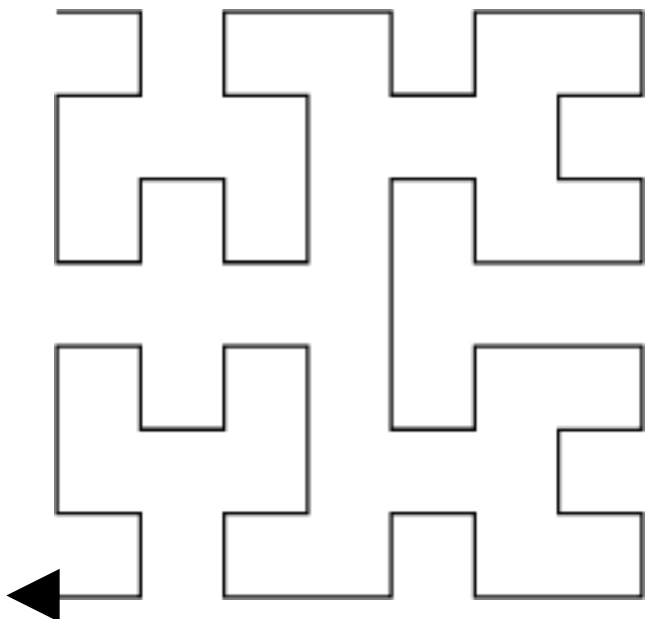
0th order



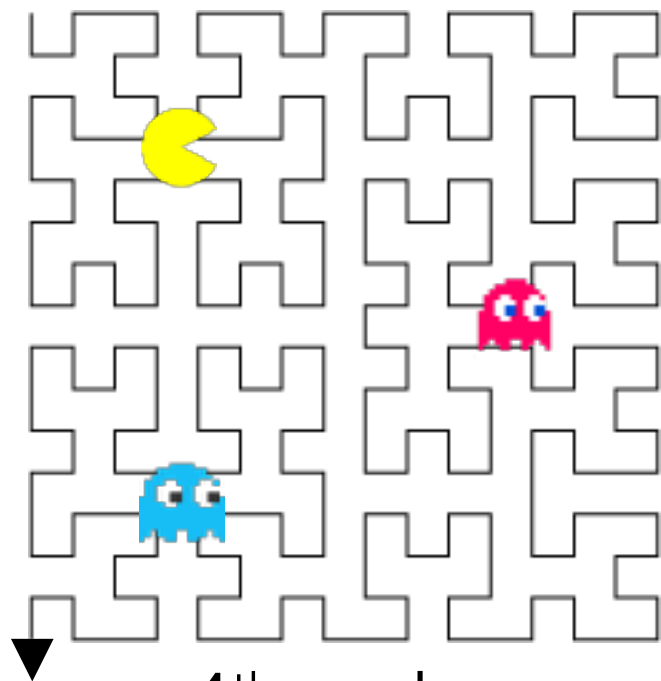
1st order



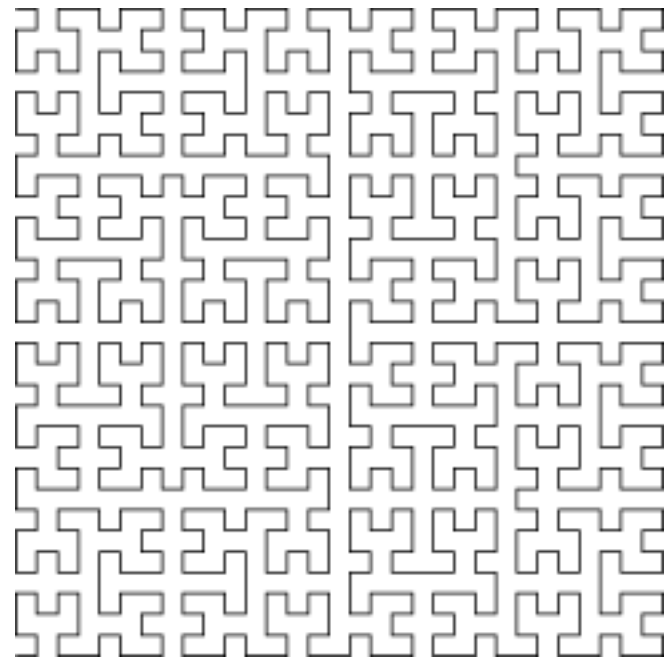
2nd order



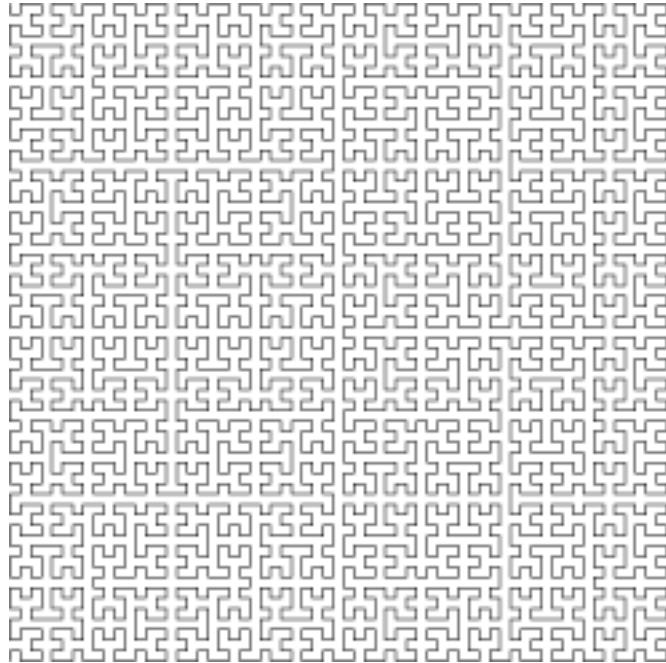
3rd order



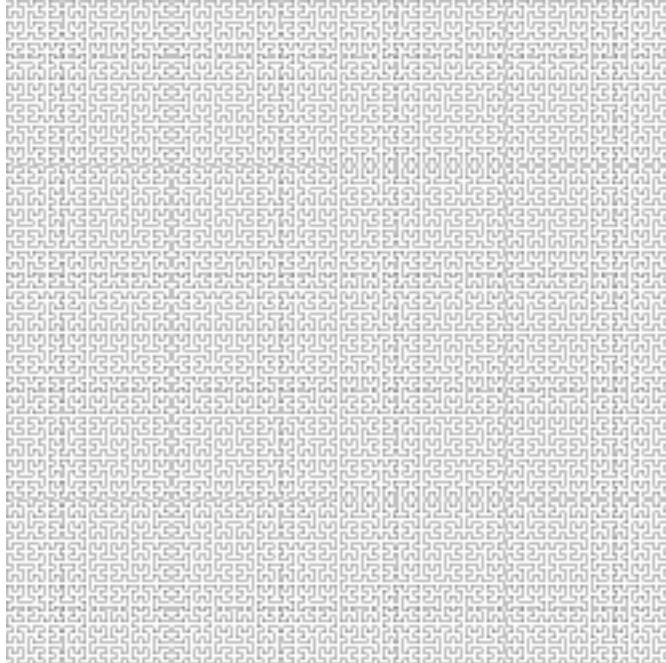
4th order



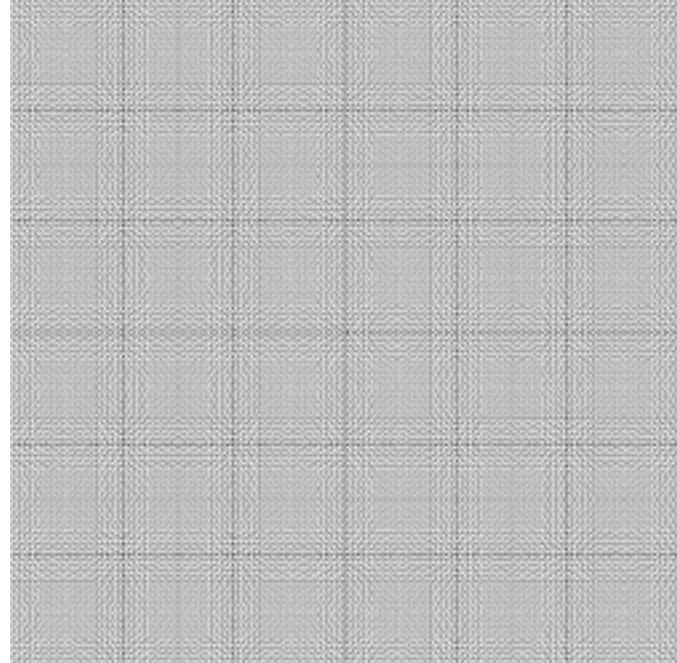
5th order



6th order

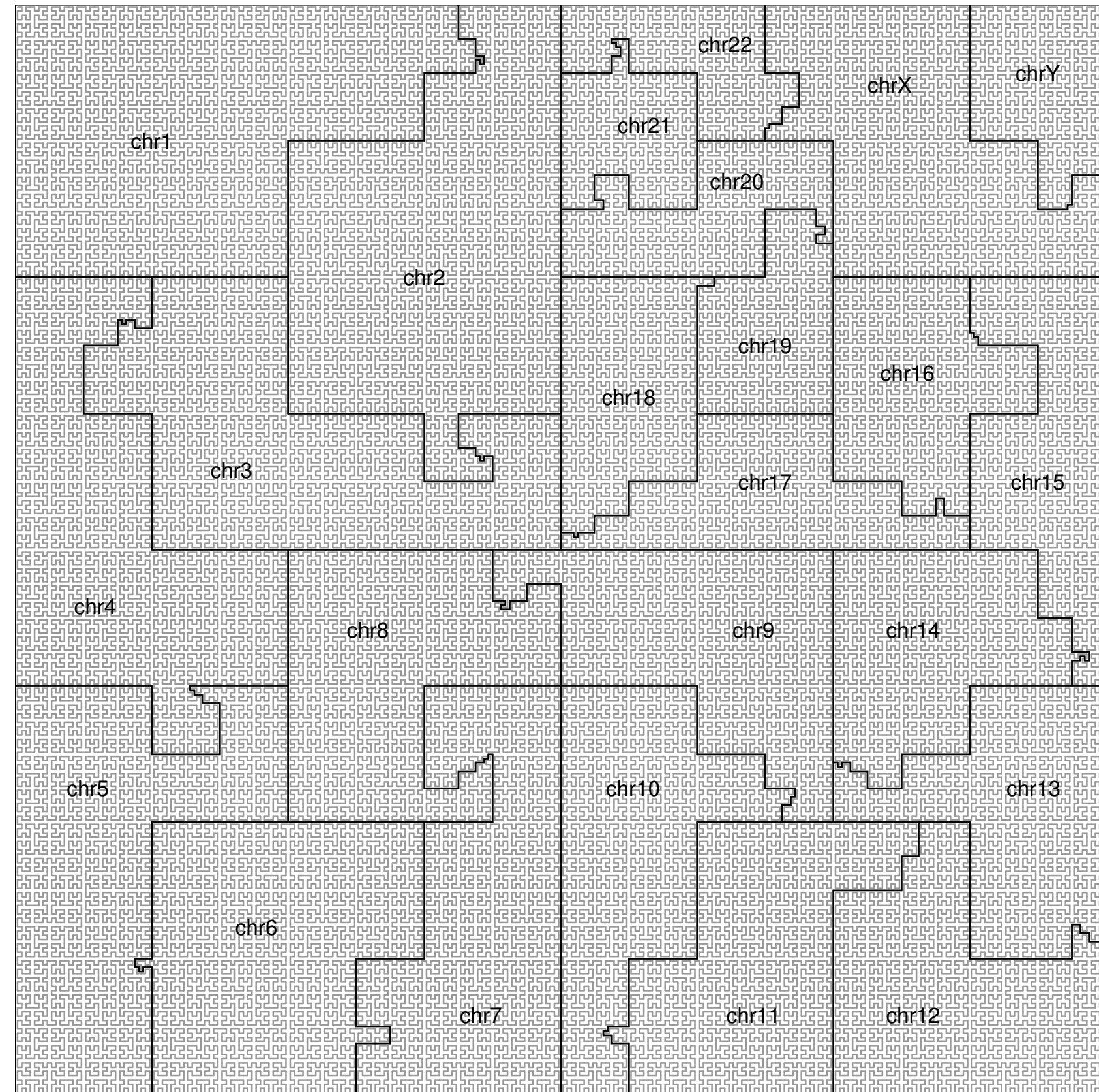


7th order



8th order

62 Hilbert curve of the human genome

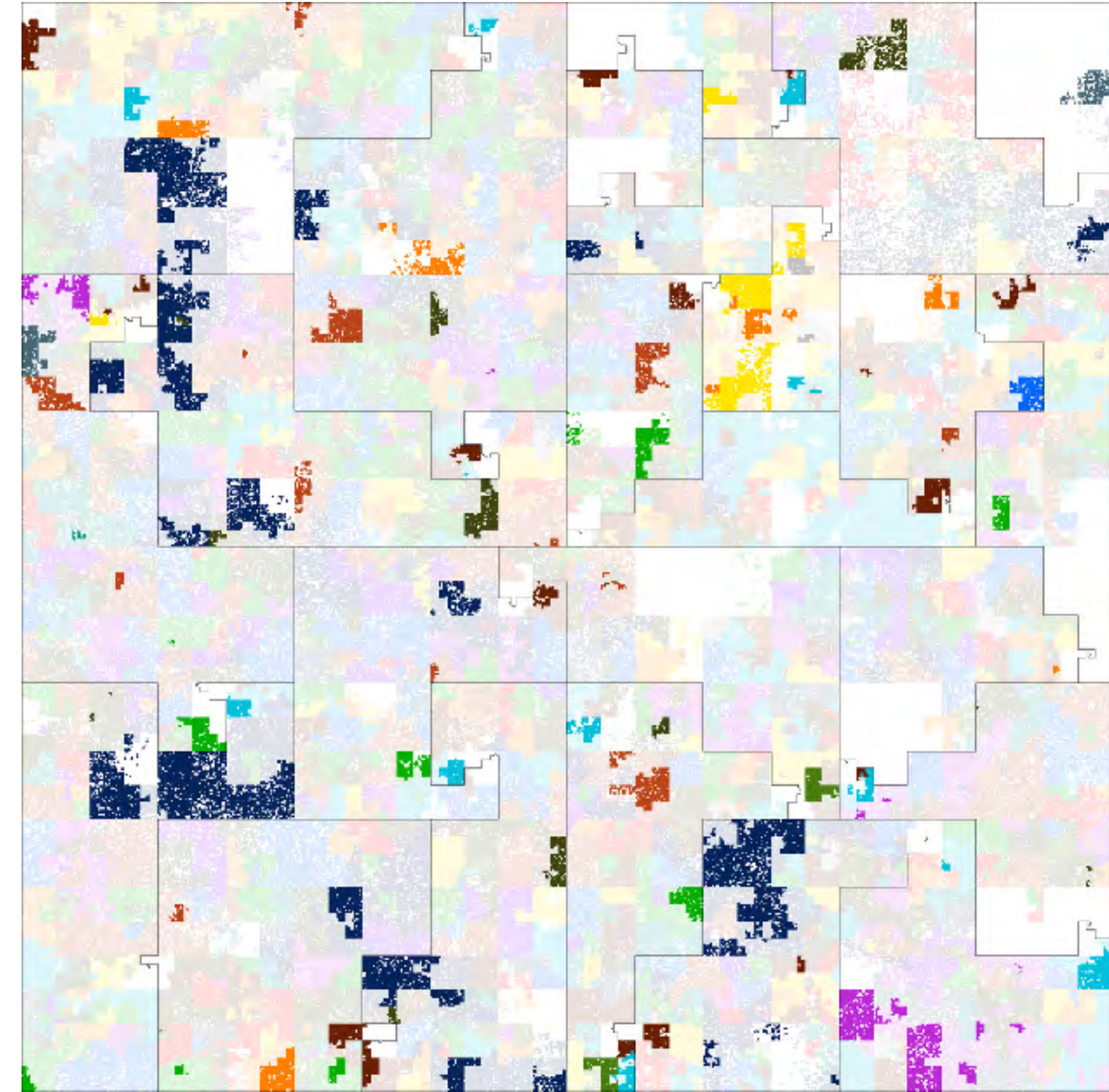


- Regions close in 1D are close in 2D
- The genome provides a (fixed) scaffold to project annotations on
- Full coverage across the genome and not limited by resolution

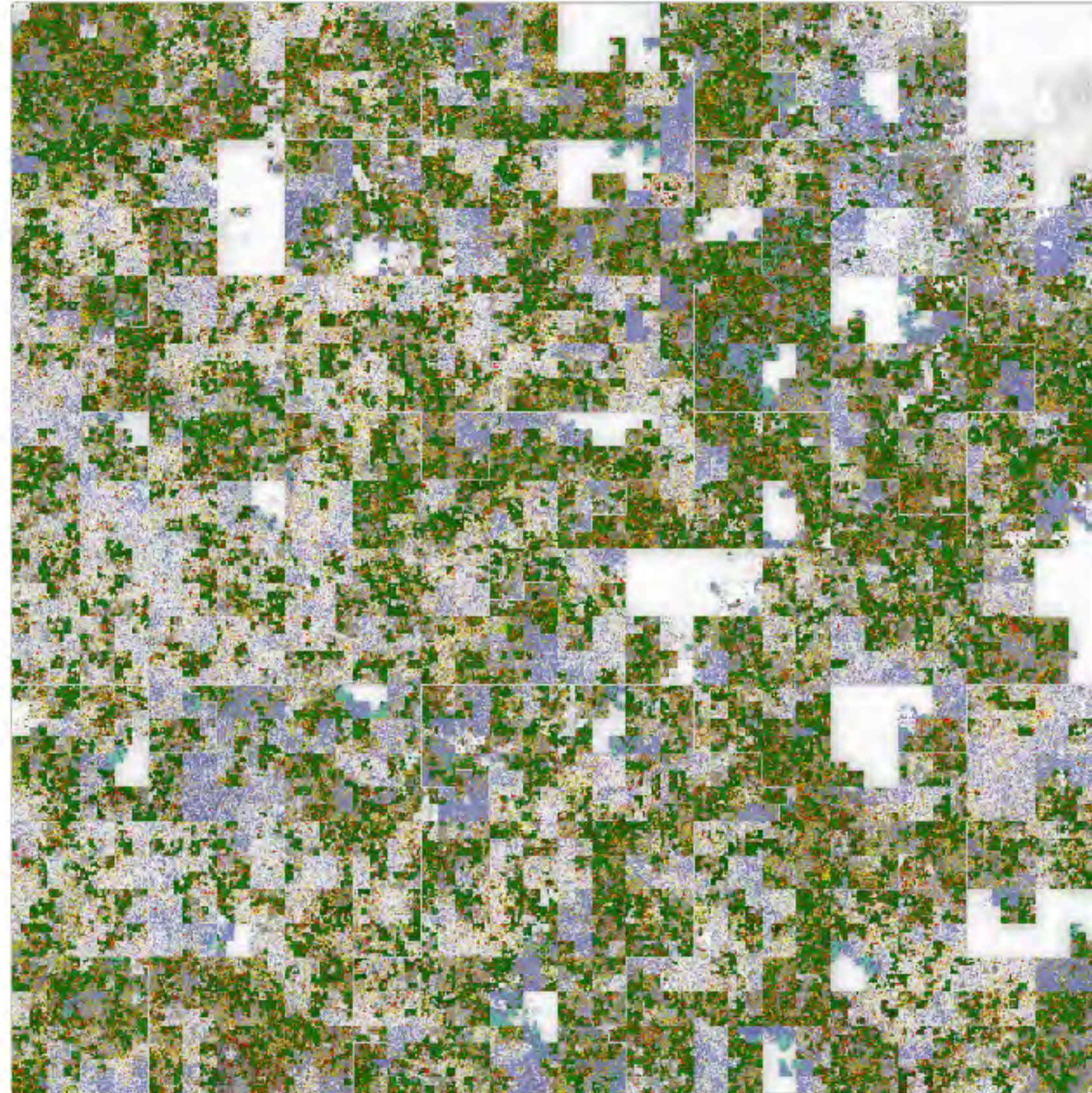
63 Hilbert curves of **functional** and **contextual** annotations



Functional annotation
(Chromatin states and epilogos)



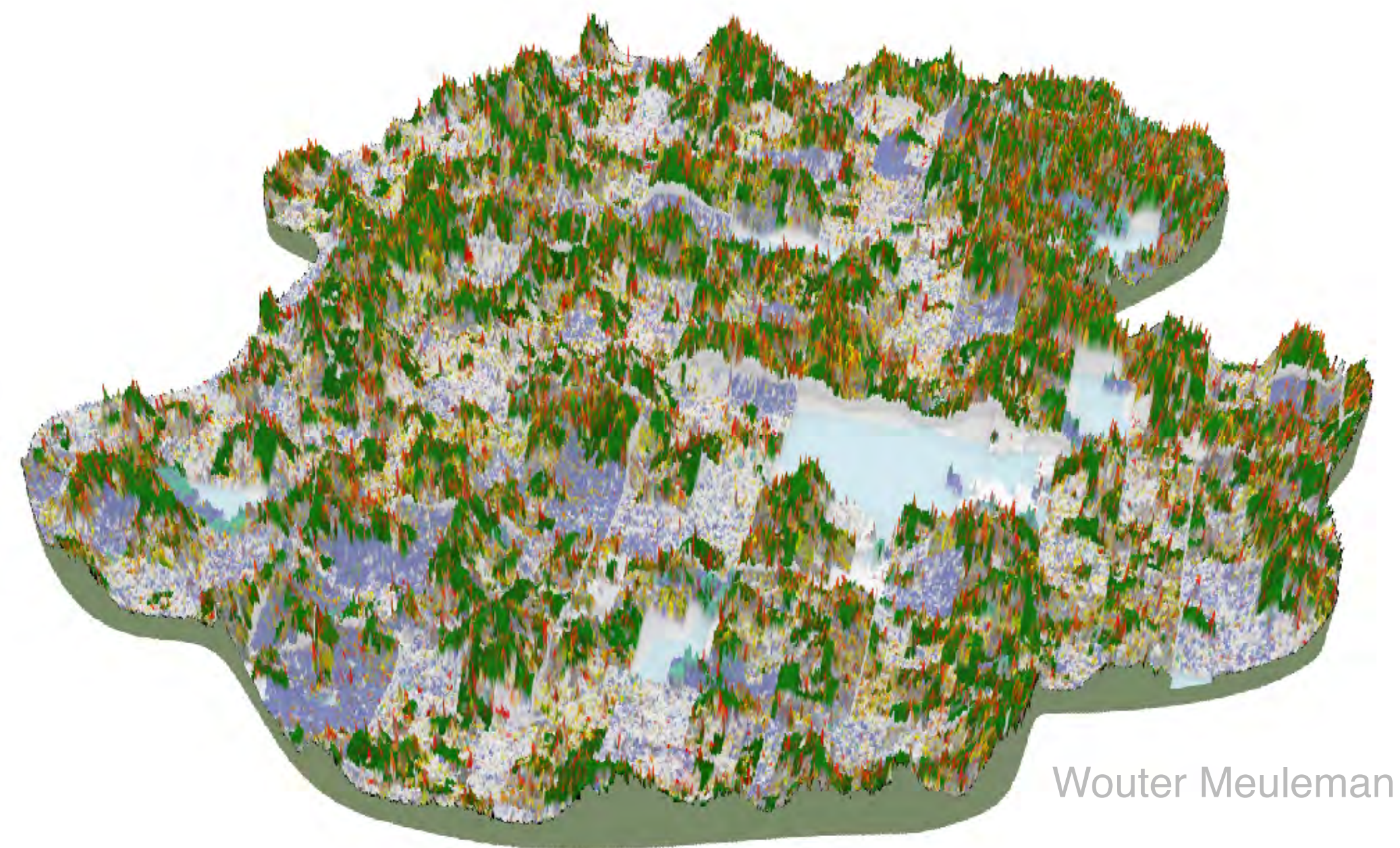
Contextual annotation
(DHS Index and domains)



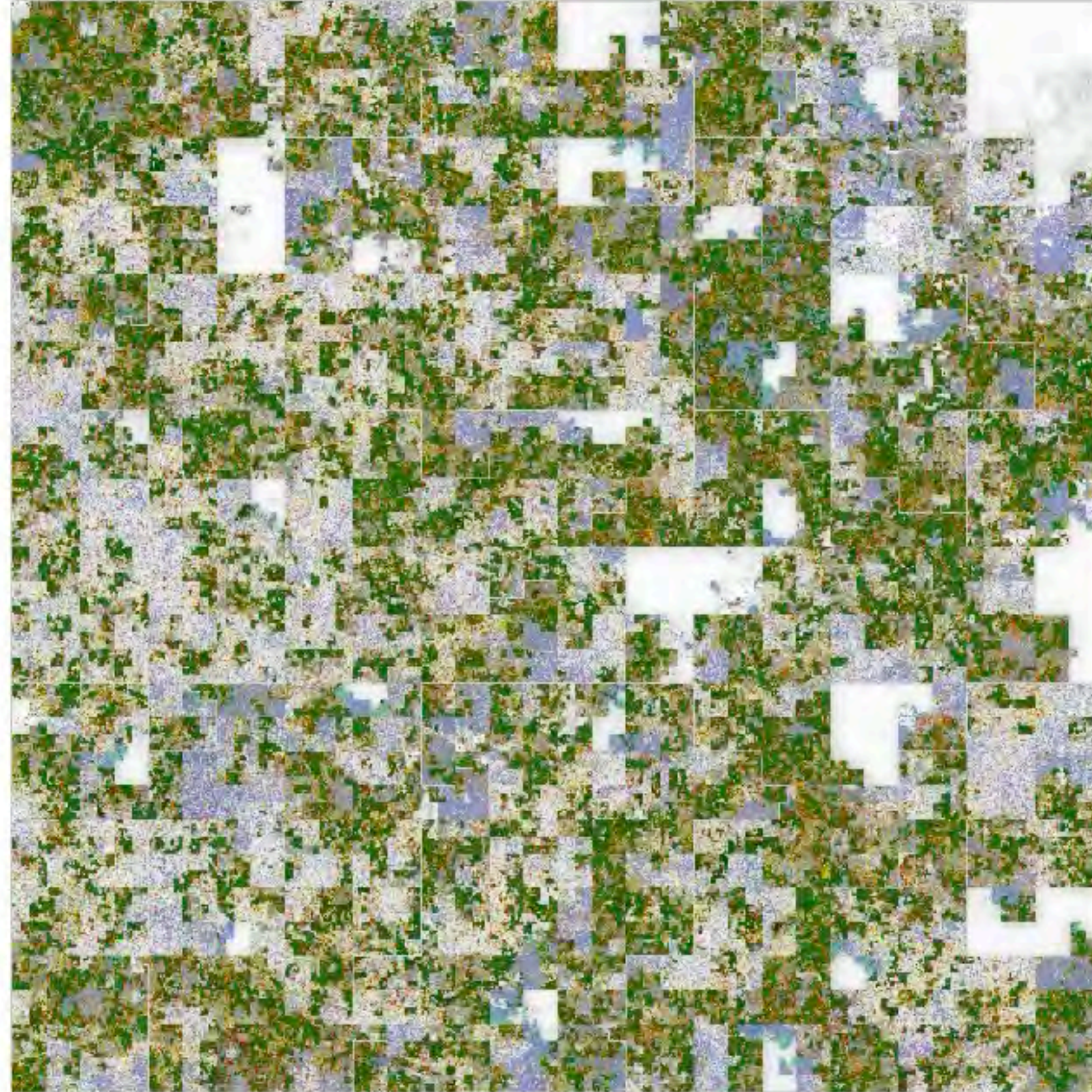
Integrative annotation
(functional+contextual+more)



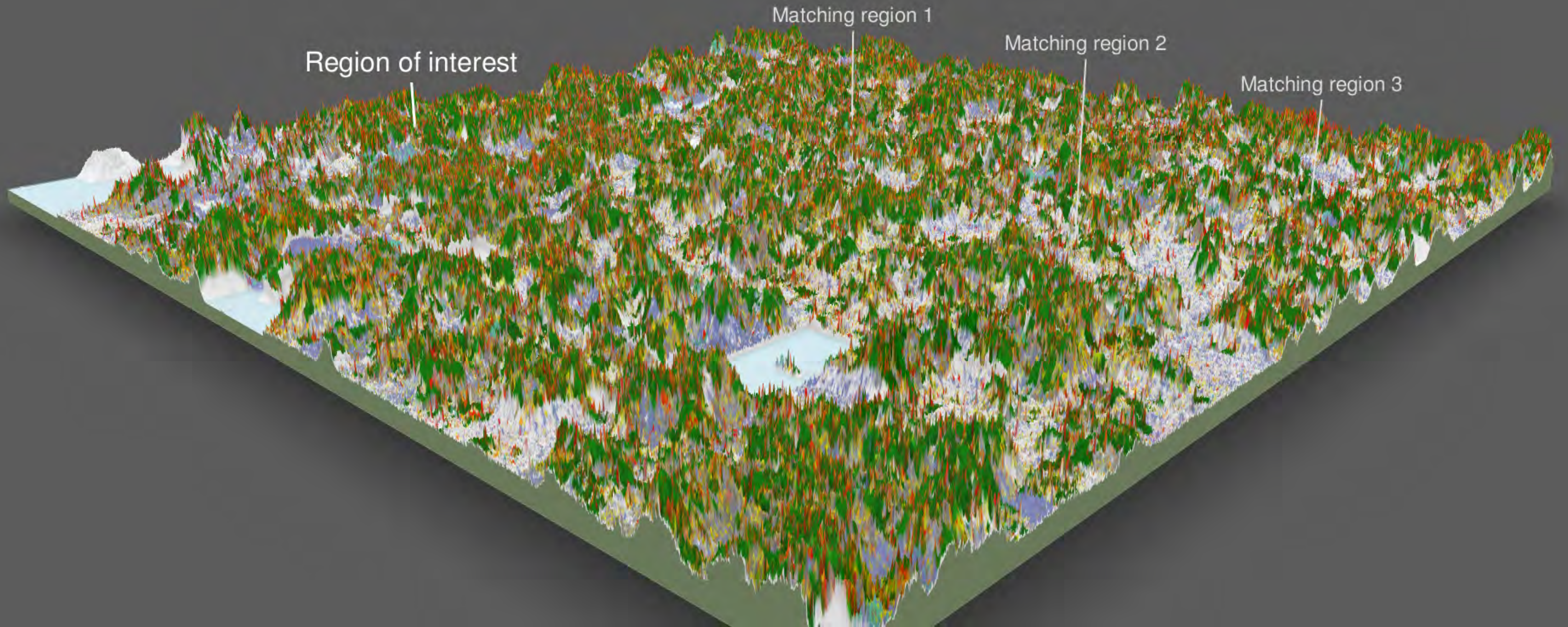
“A Disney map for Genomics”



Maps encourage exploration



Many opportunities for data driven exploration of these maps



with “~~Man versus~~ Machine”



- Show only “relevant” information to humans...
- ...while machines provide full data-driven guidance
- Human decisions get *augmented* by machines

We need better Navigation Systems

I consider these efforts part of a new field: "**augmented genomics**",
in which the work of genome scientists is supplemented — not
replaced! — by data-driven machine intelligence

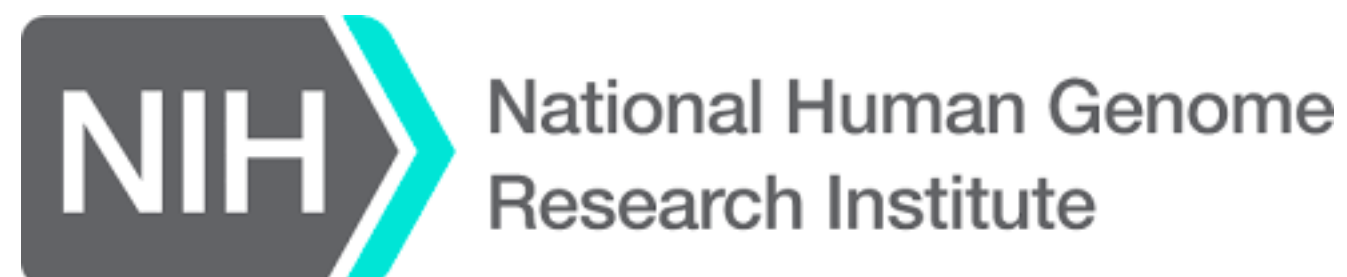
76 Now hiring!

Looking for students, postdocs and alternatively experienced folks

Are you curious about the regulatory genome and how it is organized in a cell nucleus? Do you have affinity with squeezing information out of large datasets? Want to have an impact in next generation regulatory annotations?



Positions available immediately, until filled



@nameluem 
www.meuleman.org/hiring

Acknowledgements



Data analysis

Sasha Muratov
Eric Rynes
Alex Reynolds
Jacob Quon
Nalu Tripician
Nasi Teodosiadis
Eric Haugen

Miscellaneous

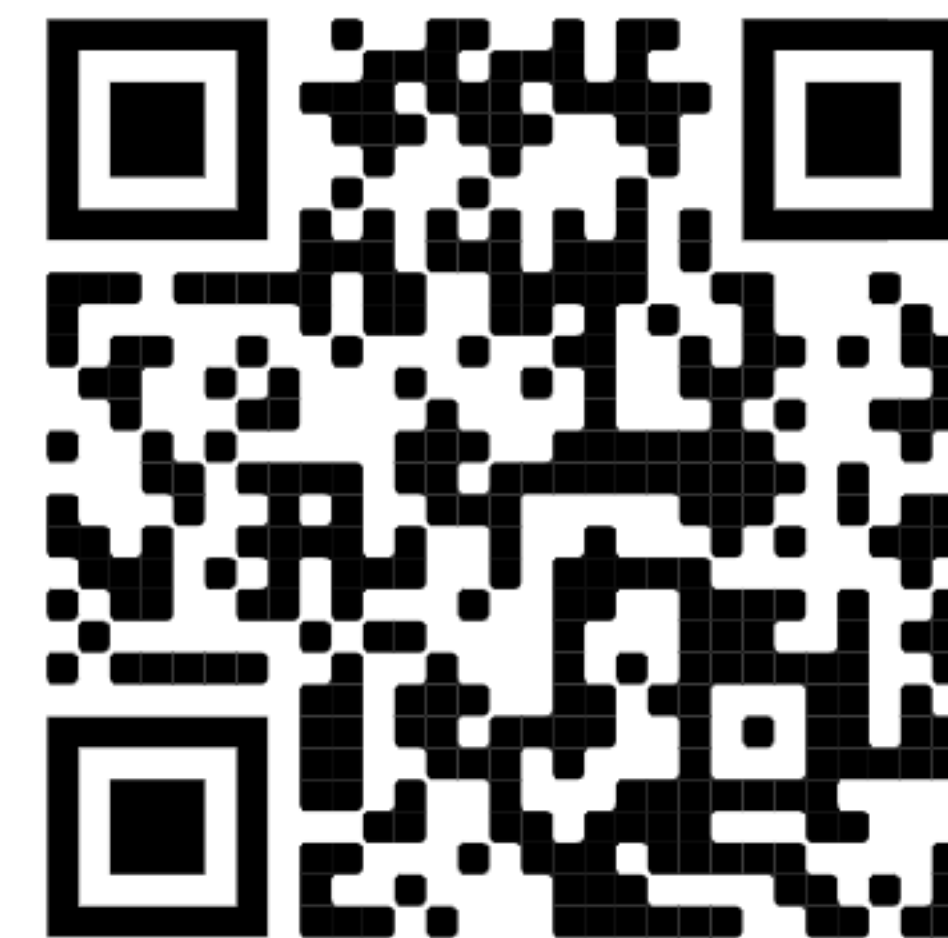
Chad Lundberg
Rae Senarighi
Tim Mercer
Jeff Vierstra
John Stam.

Data generation

Jessica Halow
Kristen Lee
Daniel Bates
Morgan Diegel
Douglass Dunn
Fidencio Neri

Data coordination

Richard Sandstrom
Audra Johnson
Jemma Nelson
Mark Frerker
Michael Buckley
Rajinder Kaul



National Human Genome
Research Institute



@nameluem 
www.meuleman.org/hiring