

EN-TE_x / Personal Genomes subgroup

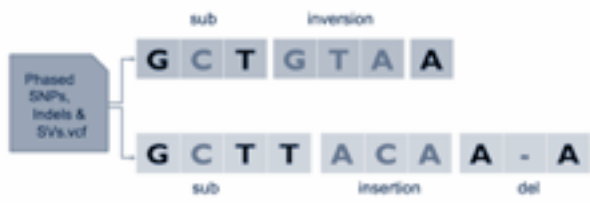
M Gerstein & E Aiden
'19 ENCODE meeting

Goals of EN-TEx

Collaboration between ENCODE & GTEx

Cataloguing genomic elements in personal genomes with matched ENCODE datasets & the full breadth of assays on many tissues

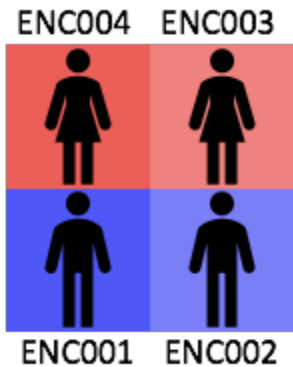
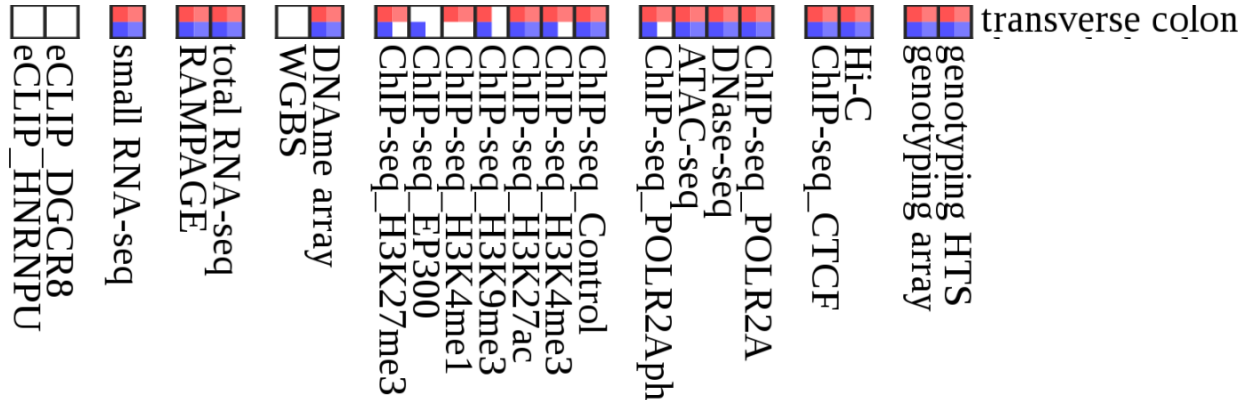
- 1. Collect **matched WGS & functional genomic** data across ~20 tissues for 4 individuals.
- 2. Catalogue genomic elements across individuals & tissues and study their **variation**
- 3. Study the utility of **phased personal diploid genomes** for analyzing functional genomics data



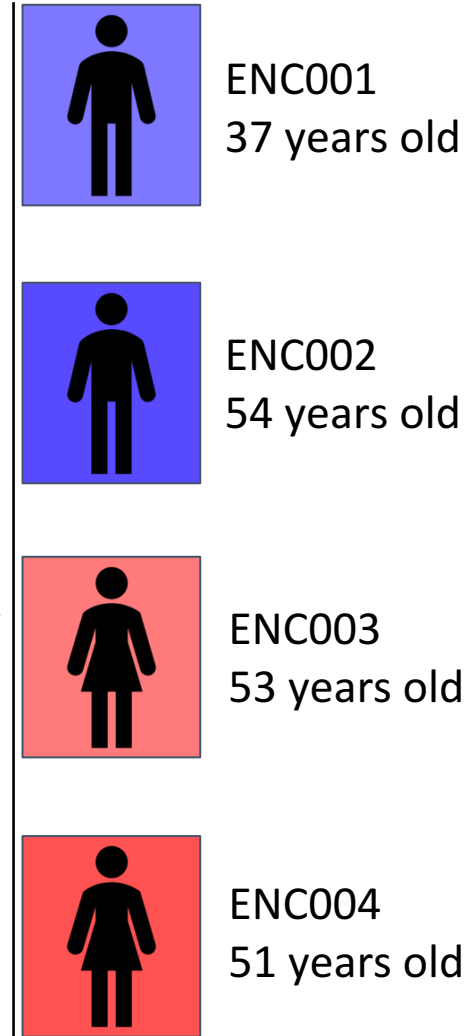
1. The Matched Data Matrix & Personal Genome

EN-TEX

Assays, Individuals & Tissue Types

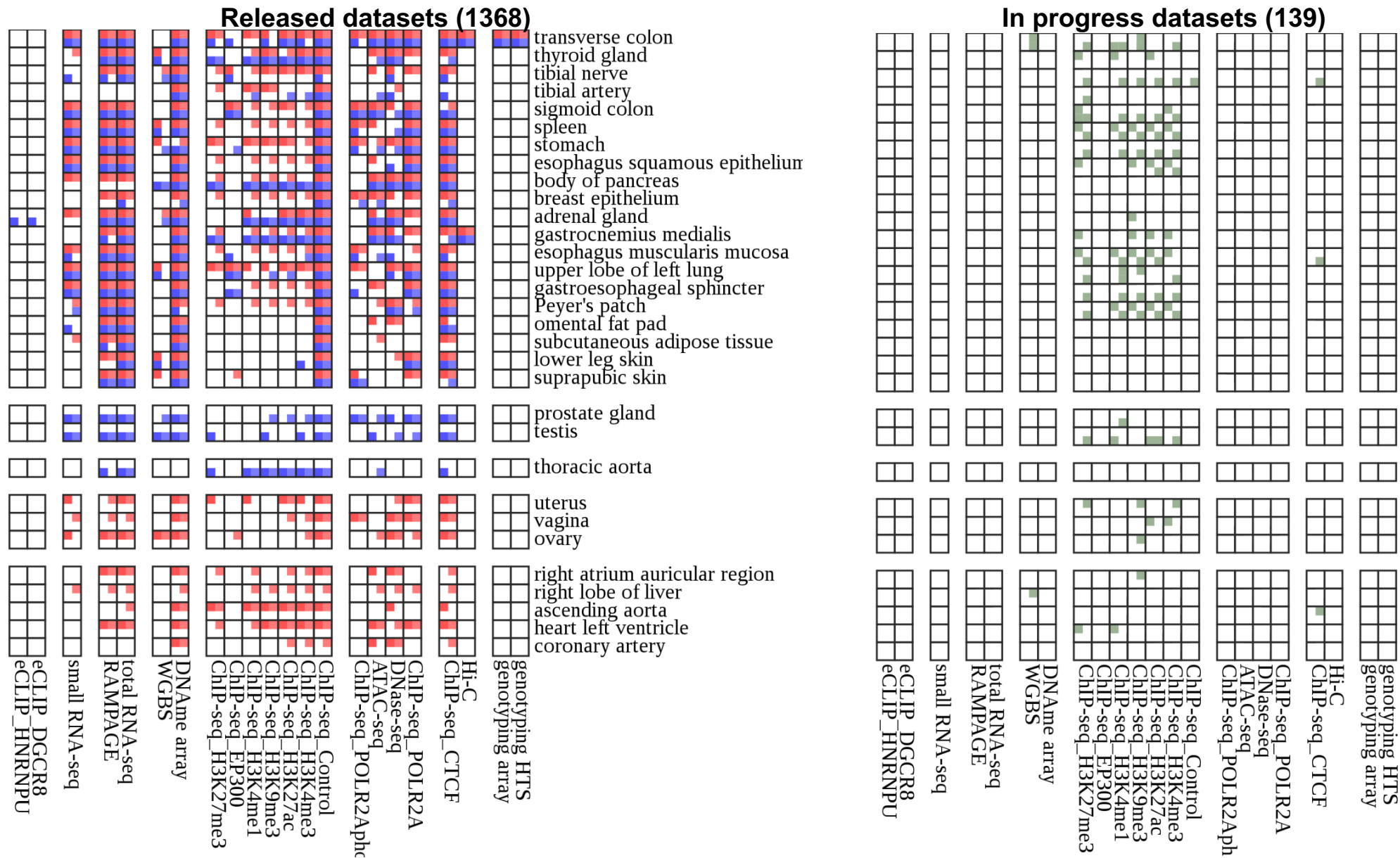


European-Americans



- Gastrocnemius medialis
- Stomach
- Transverse colon
- Peyer's patch
- Sigmoid colon
- Right lobe of liver
- Body of pancreas
- Esophagus squamous epithelium
- Gastroesophageal sphincter
- Esophagus muscularis mucosa
- Testis
- Prostate gland
- Uterus
- Vagina
- Ovary
- Thyroid gland
- Adrenal gland
- Spleen
- Right atrium auricular region
- Heart left ventricle
- Thoracic aorta
- Ascending aorta
- Omental fat pad
- Subcutaneous adipose tissue
- Tibial nerve
- Suprapubic skin
- Lower leg skin
- Breast epithelium
- Upper lobe of left lung

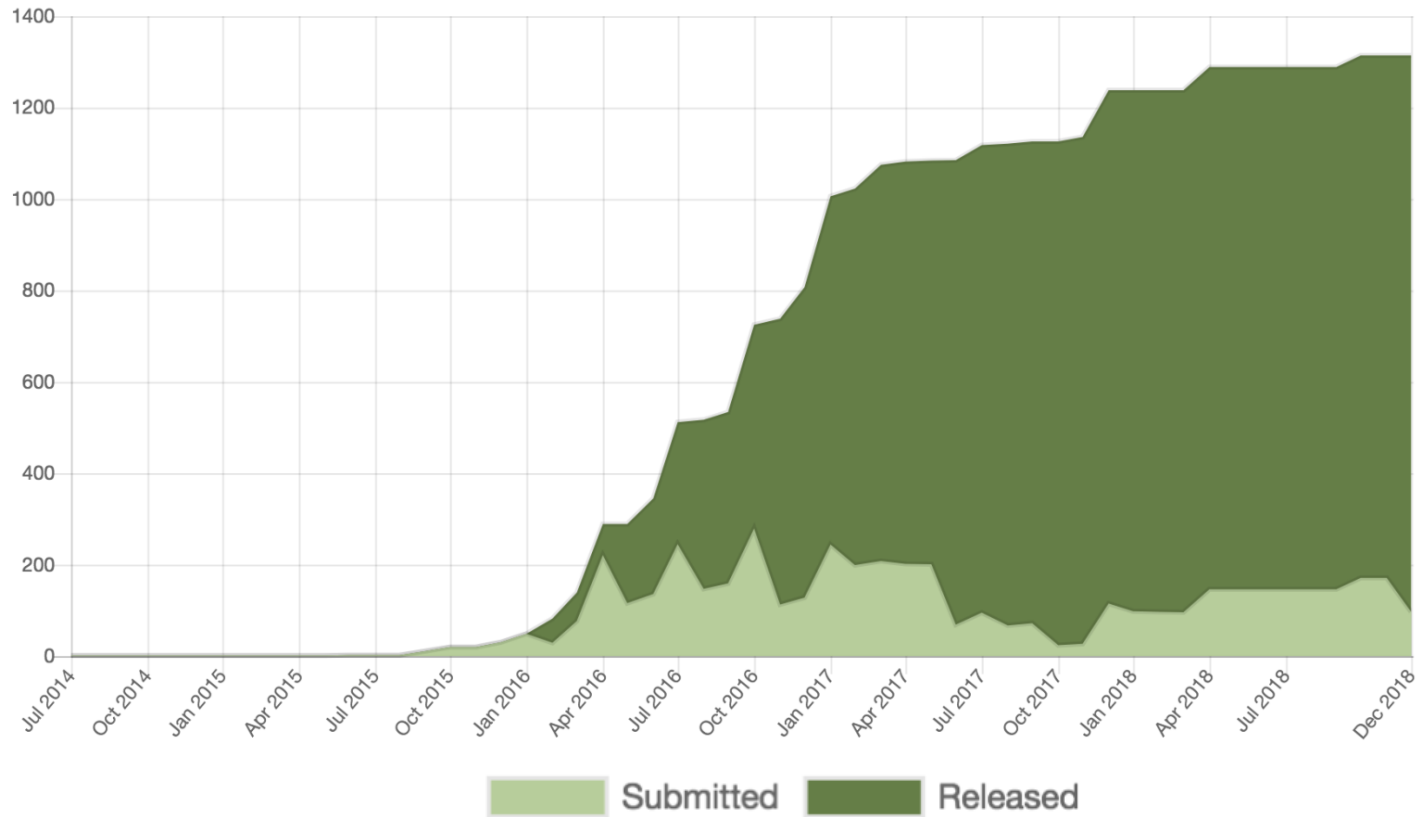
EN-TEEx Full Data Matrix as of Dec-2018



Additional data releases (beyond 1368+139) - expected in Jan & April

EN-TEch dataset (Number of assays)

Number of Assays by Status

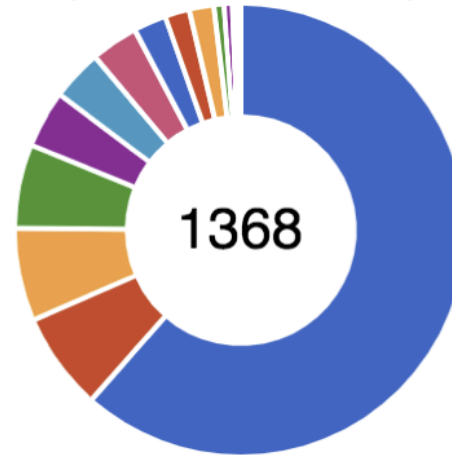


Lab (Gb)



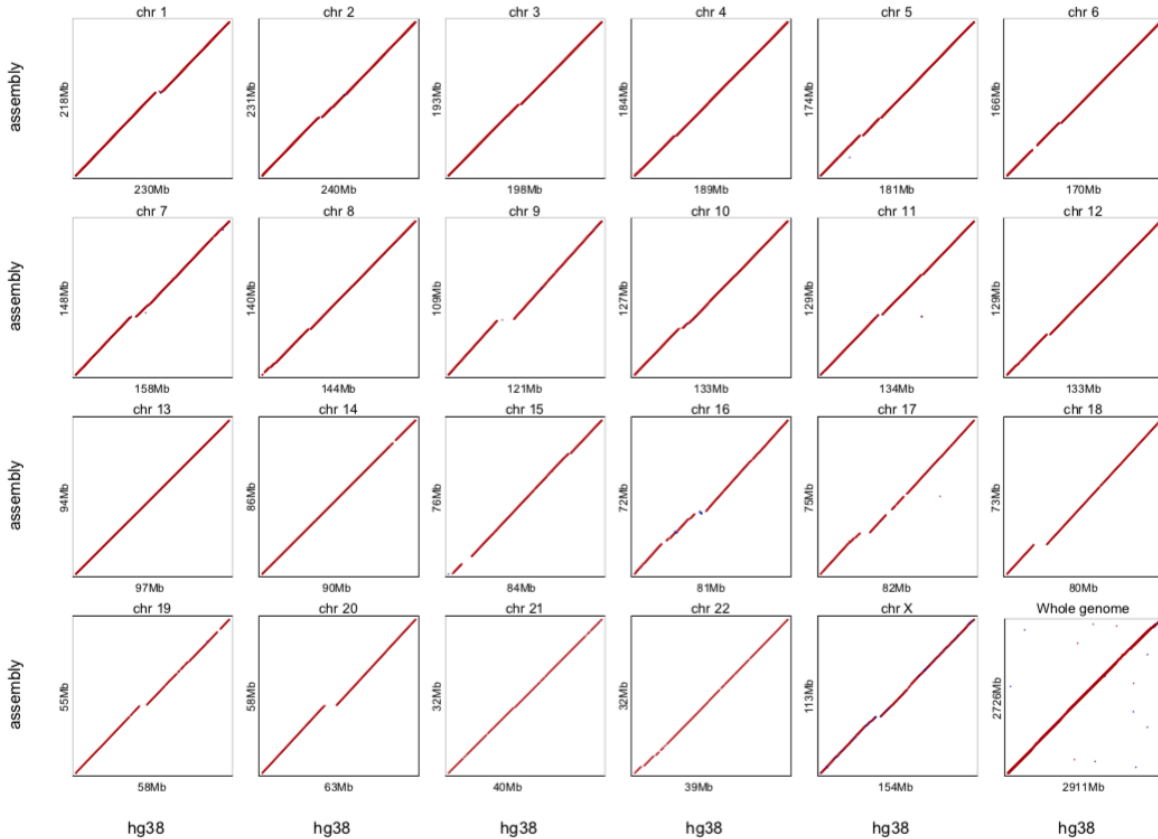
- richard-myers
- thomas-gingeras
- michael-snyder
- john-stamatoyannopoulos
- bradley-bernstein
- brenton-graveley
- bing-ren
- barbara-wold

Assay (# Experiments)



- ChIP-seq
- DNAME array
- total RNA-seq
- RAMPAGE
- DNase-seq
- small RNA-seq
- ATAC-seq
- WGBS
- microRNA counts
- microRNA-seq
- genotyping HTS
- Hi-C
- eCLIP
- genotyping array

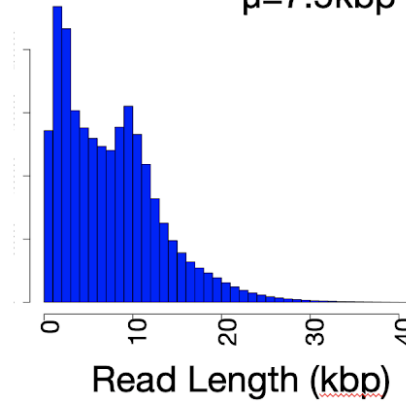
EN-TE_x Hi-C Assembly with chromosomal length scaffolds



Comparison with reference

PacBio
Long SVs (~50bp- > 25000kb)

$\mu=7.5\text{kbp}$



55x Long Reads
***Only for ENC-002 & ENC-003**

HiC Analysis

Illumina $m=350\text{bp}$;
> 120x Paired End
All 4 samples
Full chromosome phasing

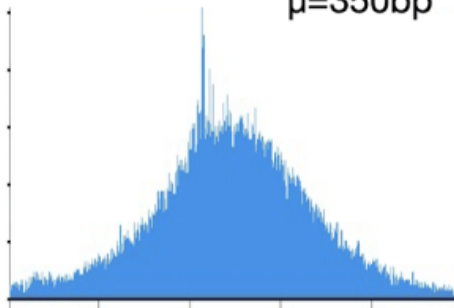
Ren Lab Hi-C data
Aiden Lab assembly
(more on tools later)

EN-TE_x WGS with whole chromosome phasing

ILLUMINA

SNPs

$\mu=350\text{bp}$



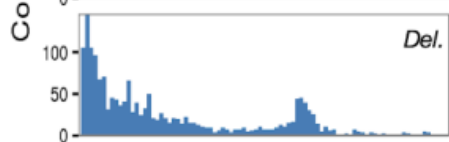
Fragment Length (kbp)

60x Paired End
All 4 samples

Variants 50 to 500 bp



Ins.



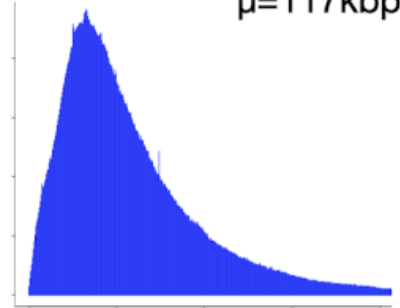
Del.

Variant size

10X Genomics

Phasing and long SV

$\mu=117\text{kbp}$



Molecule Length (kbp)

35x Linked Reads
All 4 samples

Variants 50 to 500 bp



Ins.



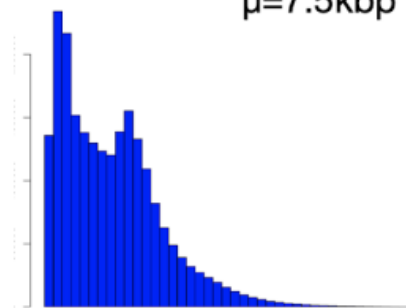
Del.

Variant size

PacBio

Long SVs (~50bp- > 25000kb)

$\mu=7.5\text{kbp}$



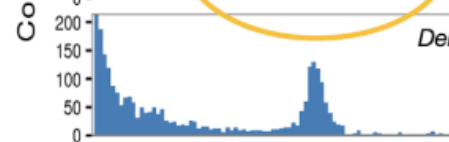
Read Length (kbp)

55x Long Reads
***Only for ENC-002 & ENC-003**

Variants 50 to 500 bp



Ins.



Del.

Variant size

HiC Analysis

ILLUMINA $m=350\text{bp}$;
> 120x Paired End
All 4 samples
Full chromosome phasing

Variant calls

	Sub	Ins (<50)	Del (<50)	SV (>50)
ENC-001	3,895,883	255,857	281,757	0
ENC-002	3,877,699	272,322	296,040	24,696
ENC-003	4,024,169	288,027	314,366	26,181
ENC-004	3,947,657	263,092	293,823	0

SV count by type

Sample	SV-Ins	SV-Del	Inv	Dup	Trans
ENC-002*	11,745	10,460	281	1,850	360
ENC-003*	13,077	10,813	248	1,746	297
GIAB son*	13,750	11,487	300	1,647	0
NA12878*	11,127	9,385	230	1,057	394
1000 Genomes‡	168	42,279	786	6,025	0

* PacBio Long Reads w/NGMLR & Sniffles

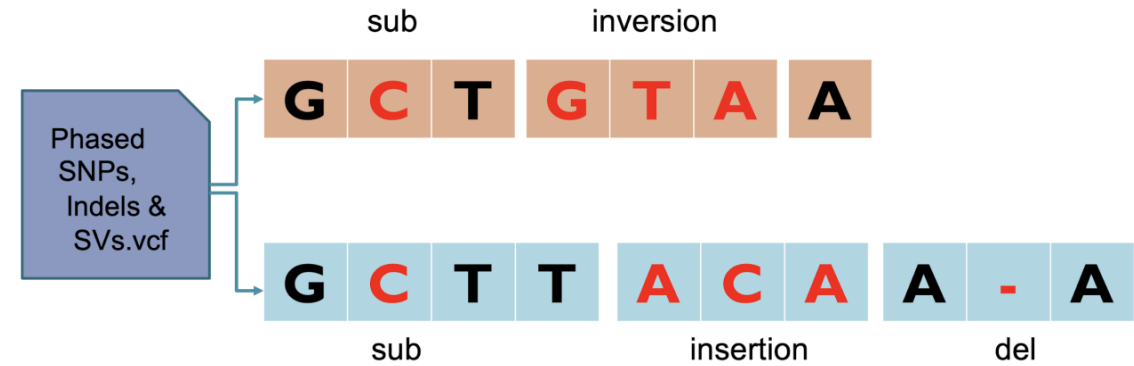
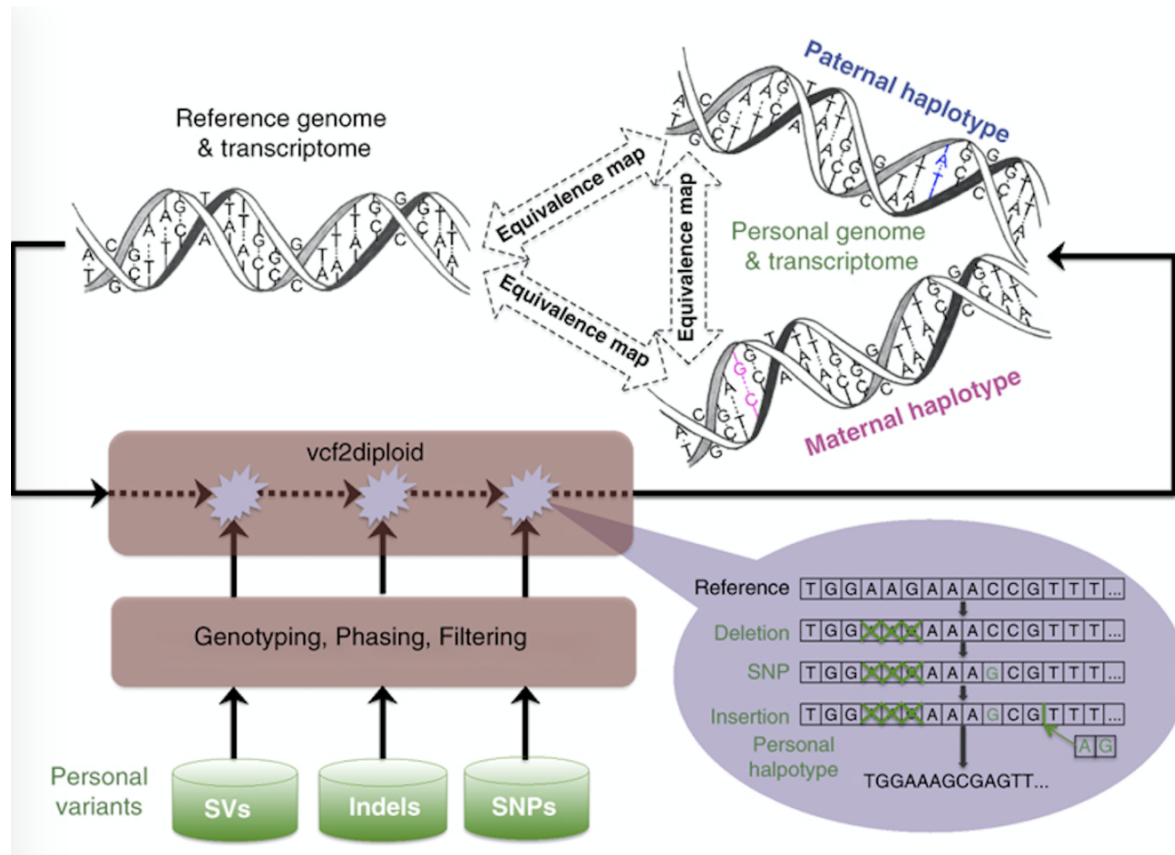
‡ Illumina short reads across 2,504 humans

EN-TEX SV calls:
very developed
& accurate

SV concordance among genomes

	ENC002	ENC003	NA12878	GIAB
ENC002	19725	12297	10746	12573
ENC003	12297	20693	10771	12551
NA12878	10746	10771	16966	11237
GIAB	12573	12551	11237	21036

Assembling a phased personal diploid genome from all the variant calls



Stitching based on AlleleSeq pipeline enhanced for SVs

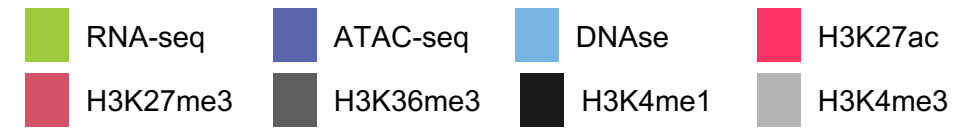
- Maintains a mapping from reference to personal genome coordinates for liftover

Using 10X + HiC + PacBio, assemble nearly perfect diploid human genomes

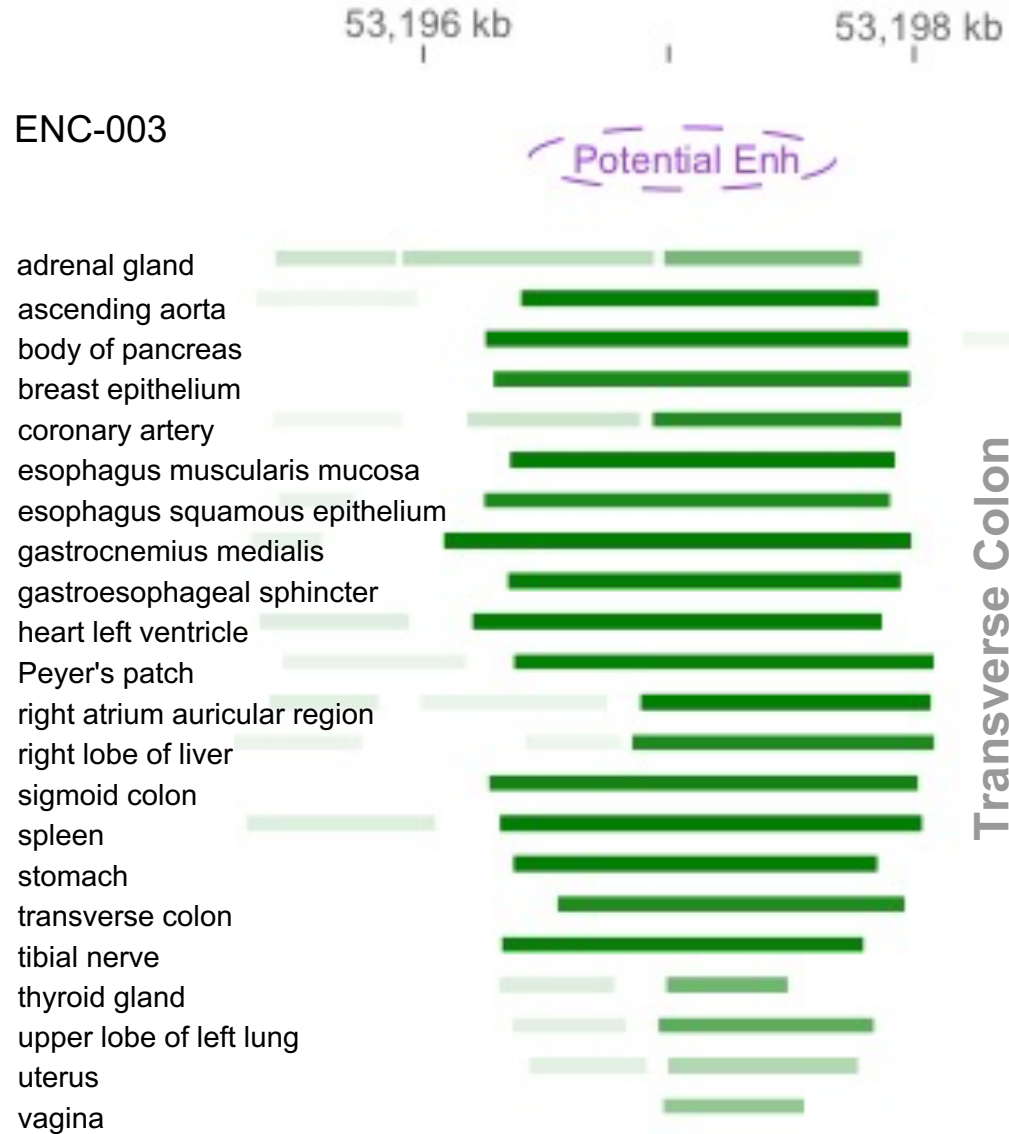
- Phased diploid genome can be aligned or aligned against

2. Variation in Transcriptome & Epigenome across Individuals & Tissues

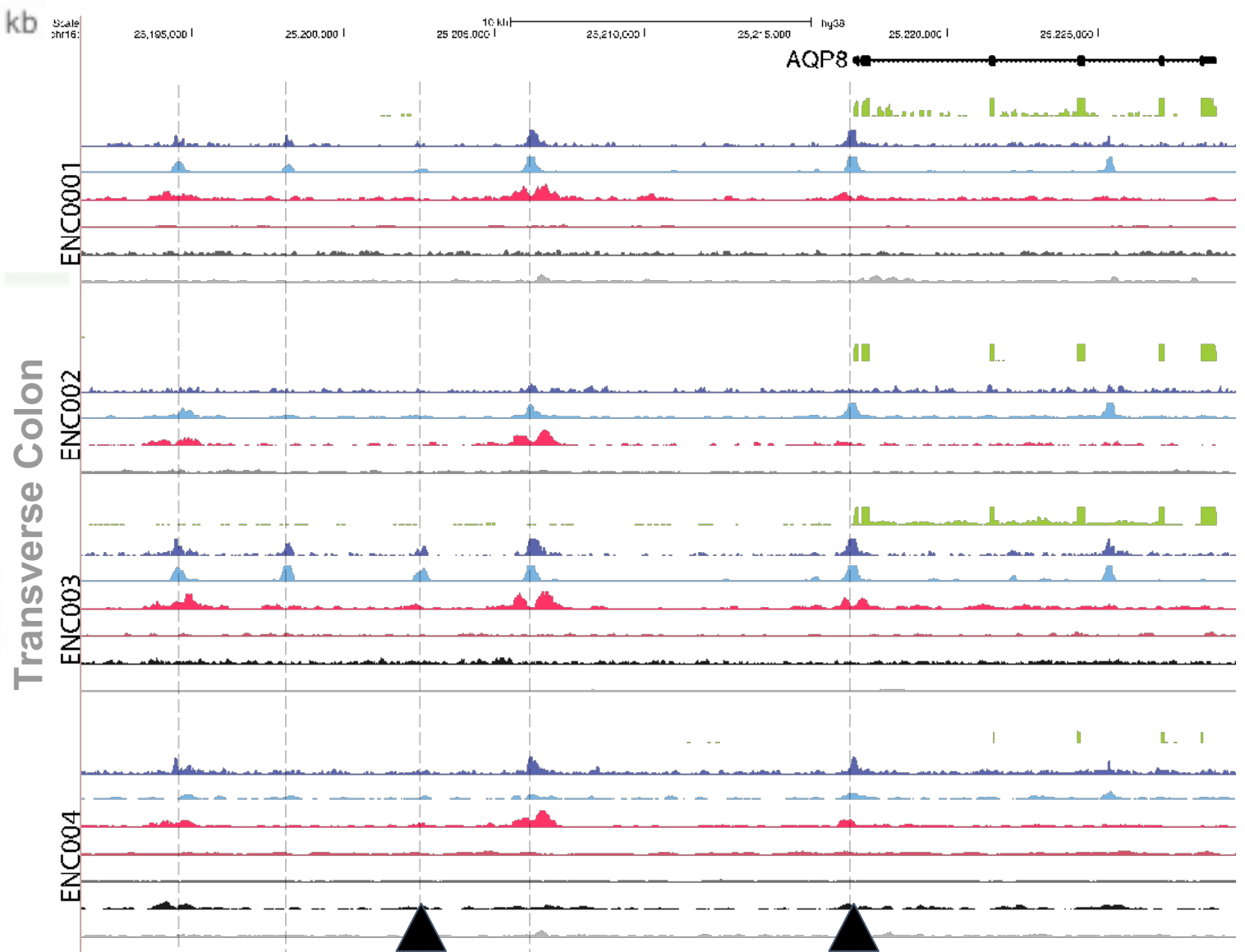
Var. across tissues & individuals (single locus)



Peak variation across tissues



Signal variation across individuals



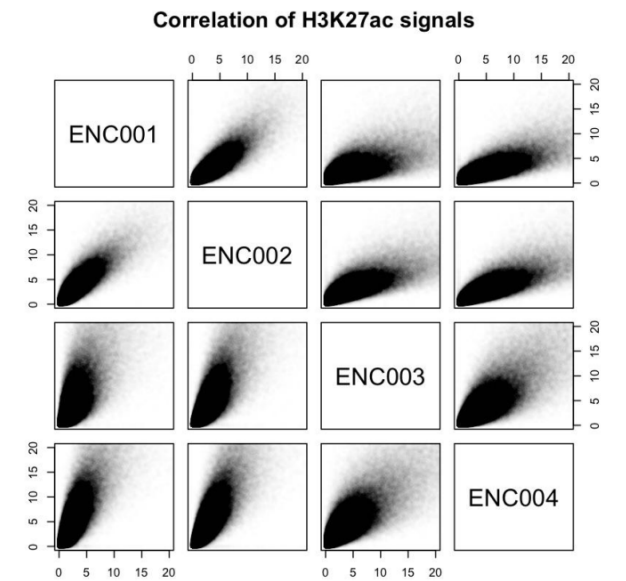
Chromatin variation across individuals (aggregated): Integration vs Single assay

Distal DHS

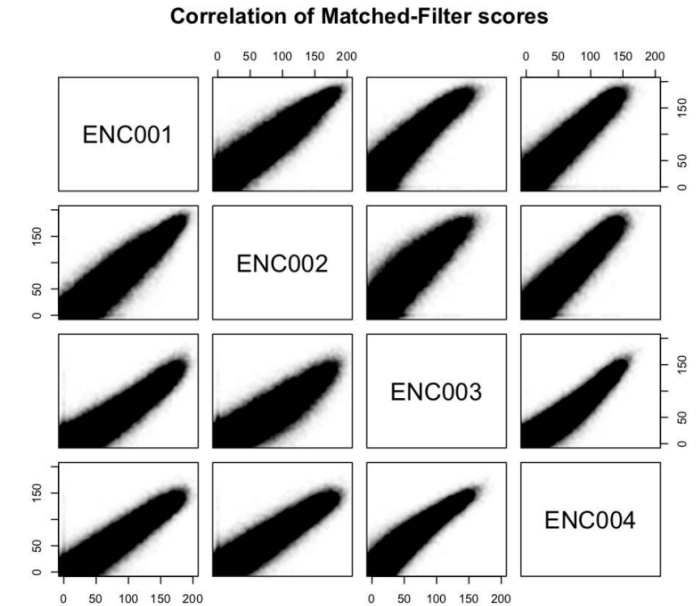
Transverse Colon

	ENC-001	ENC-002	ENC-003	ENC-004
ENC-001	161005	108157 (67%)	116780 (73%)	57606 (36%)
ENC-002	108493 (69%)	156610	104568 (67%)	62256 (40%)
ENC-003	116976 (75%)	104364 (66%)	156961	57176 (36%)
ENC-004	58353 (37%)	63135 (40%)	57818 (36%)	159323

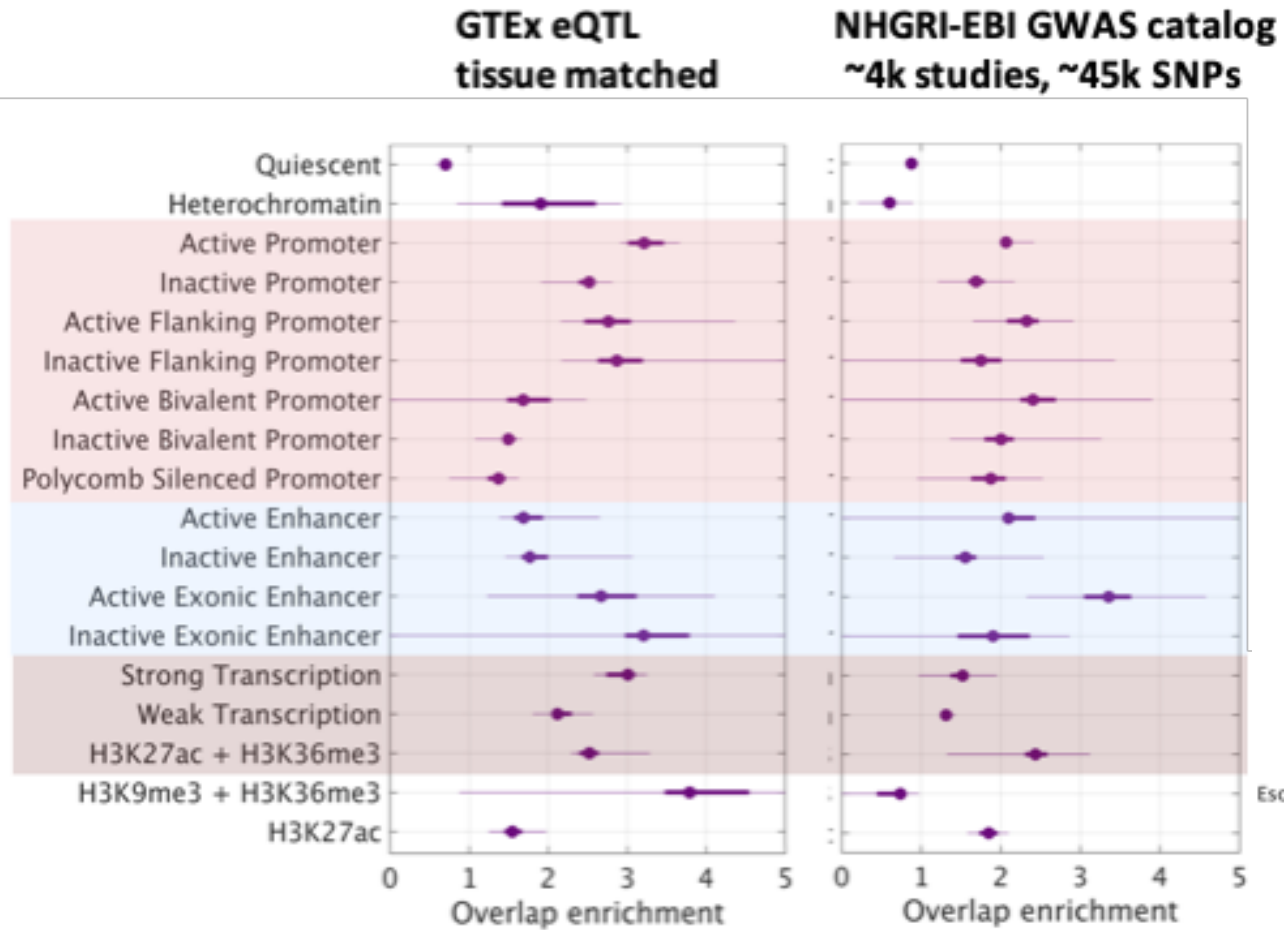
Correlation
of
K27
Signal



Corr. of
Integrated
Cross-
signal
Score
(MF Lin.
Combo)



Methods of integration (chromHMM)



EN-TEChromHMM segmentation: one donor, 10 tissues

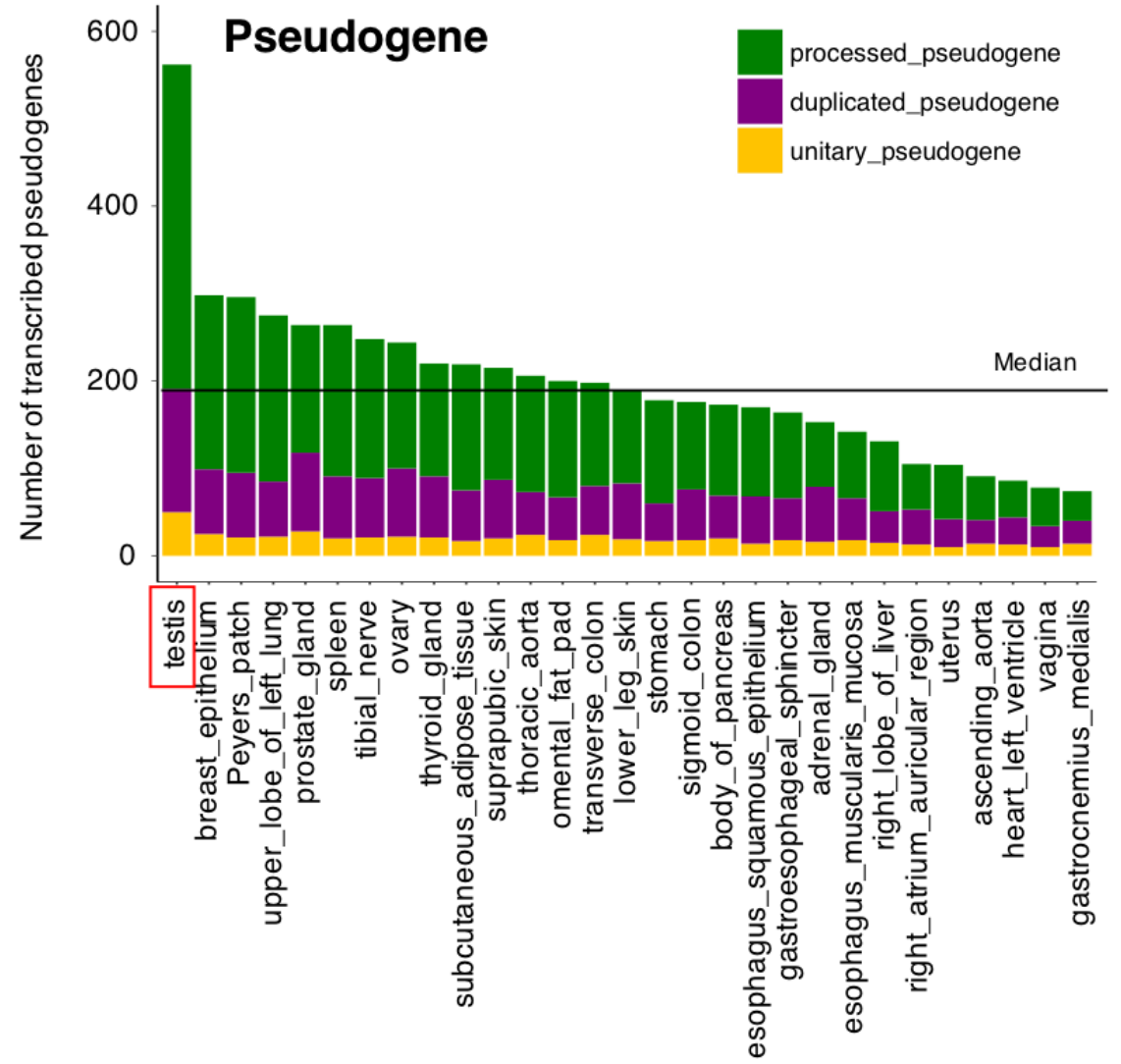
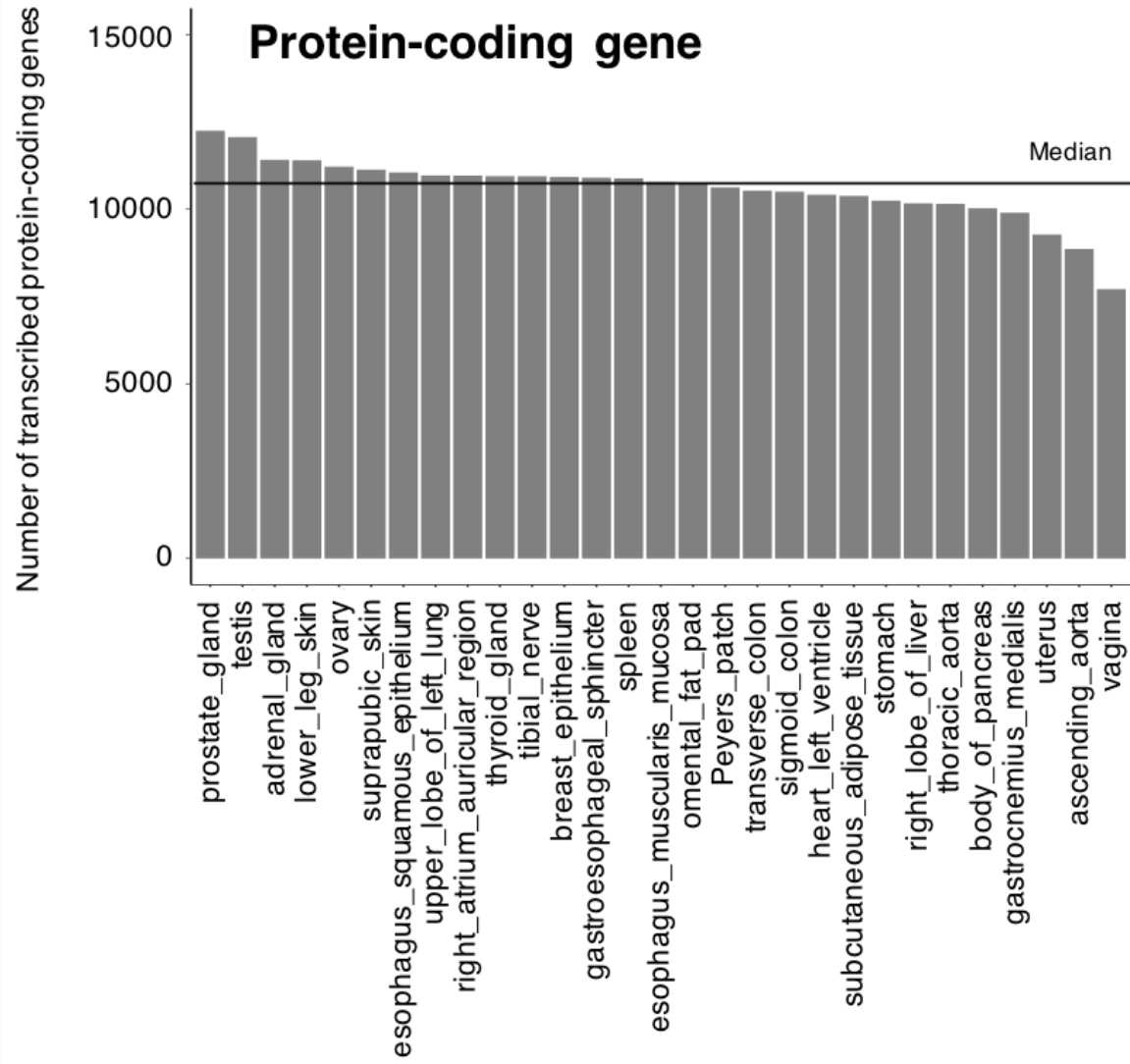
H3 K4 me3	H3 K4 me1	H3 K27 ac	H3 K36 me3	H3 K27 me3	H3 K9 me3	Genome coverage %	Exons	TSS +/- 2kb	Chromatin State
0	0	0	0	0	0	75.5	0.5	0.5	Quiescent
0	0	0	0	1	45	1.6	0.7	0.4	Heterochromatin
98	4	98	0	1	0	0.7	9.0	23.5	Active Promoter
79	7	9	0	2	0	0.3	5.1	16.5	Inactive Promoter
97	87	99	0	1	0	0.2	4.2	16.9	Active Flanking Promoter
90	82	7	0	2	0	0.1	5.0	19.6	Inactive Flanking Promoter
87	31	88	1	73	4	0.03	8.5	19.4	Active Bivalent Promoter
62	41	1	1	84	5	0.11	8.5	16.6	Inactive Bivalent Promoter
0	0	0	0	33	0	0.8	2.6	5.7	Polycomb Silenced Promoter
3	74	93	1	0	0	1.0	1.6	2.5	Active Enhancer
0	48	6	0	0	0	1.0	1.3	3.2	Inactive Enhancer
37	55	90	80	0	0	0.04	11.7	5.4	Active Exonic Enhancer
4	49	9	68	0	0	0.03	8.8	3.5	Inactive Exonic Enhancer
0	0	0	76	0	0	2.5	5.9	0.7	Strong Transcription
0	0	0	3	0	0	12.0	2.1	1.0	Weak Transcription
1	3	58	75	0	0	0.14	8.1	1.2	H3K27ac + H3K36me3
4	0	1	70	0	80	0.12	7.6	0.9	H3K9me3 + H3K36me3
1	6	65	0	0	0	1.9	1.3	2.4	H3K27ac

One of the most comprehensive datasets of consistent collection of histone modifications for multiple individuals and tissues allows better chromatin segmentation

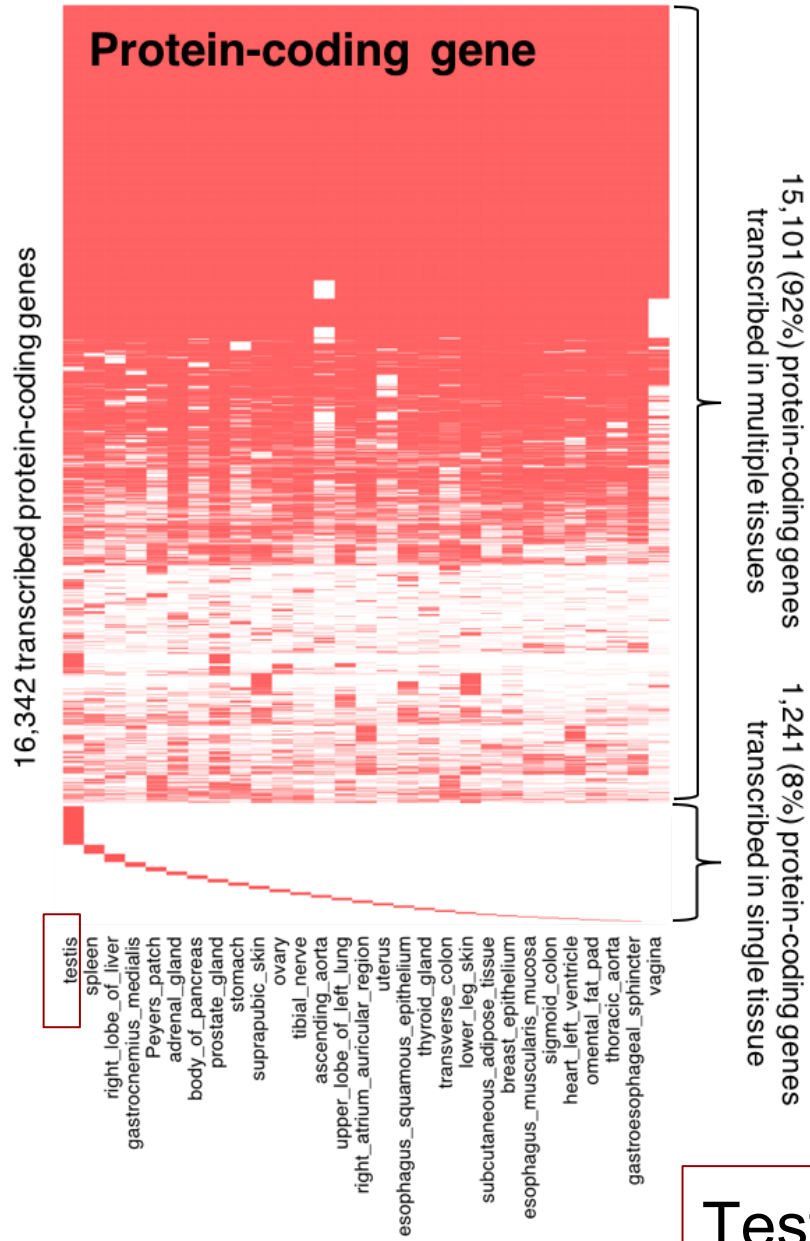
GWAS variants in "cholesterol" studies

83 studies; 1,183 SNPs; 680 unique SNPs

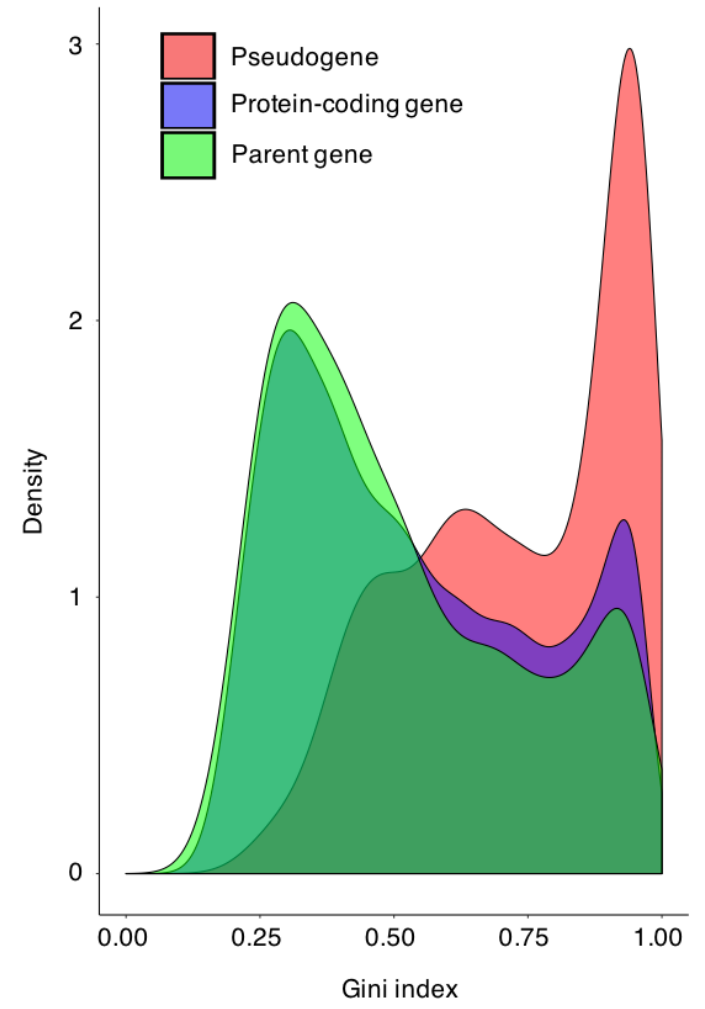
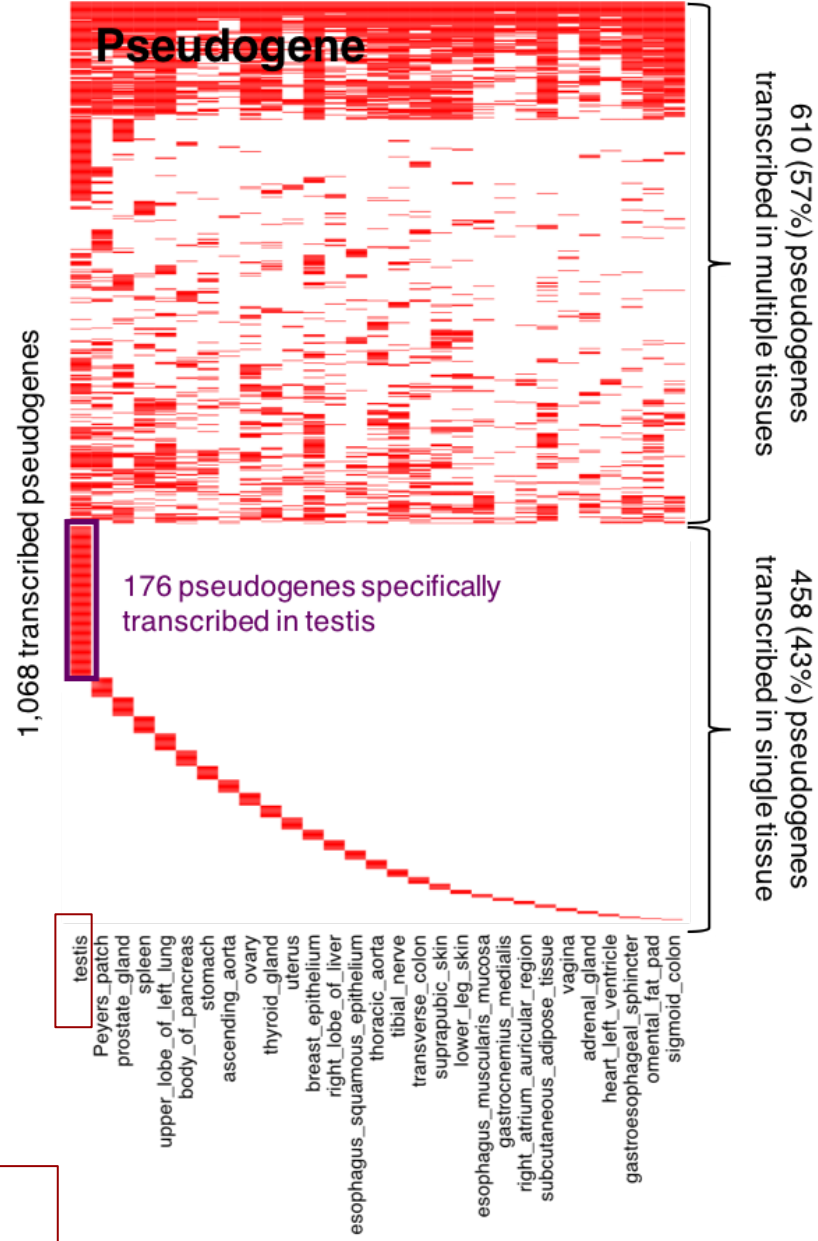
Landscape of EN-TE_x transcription



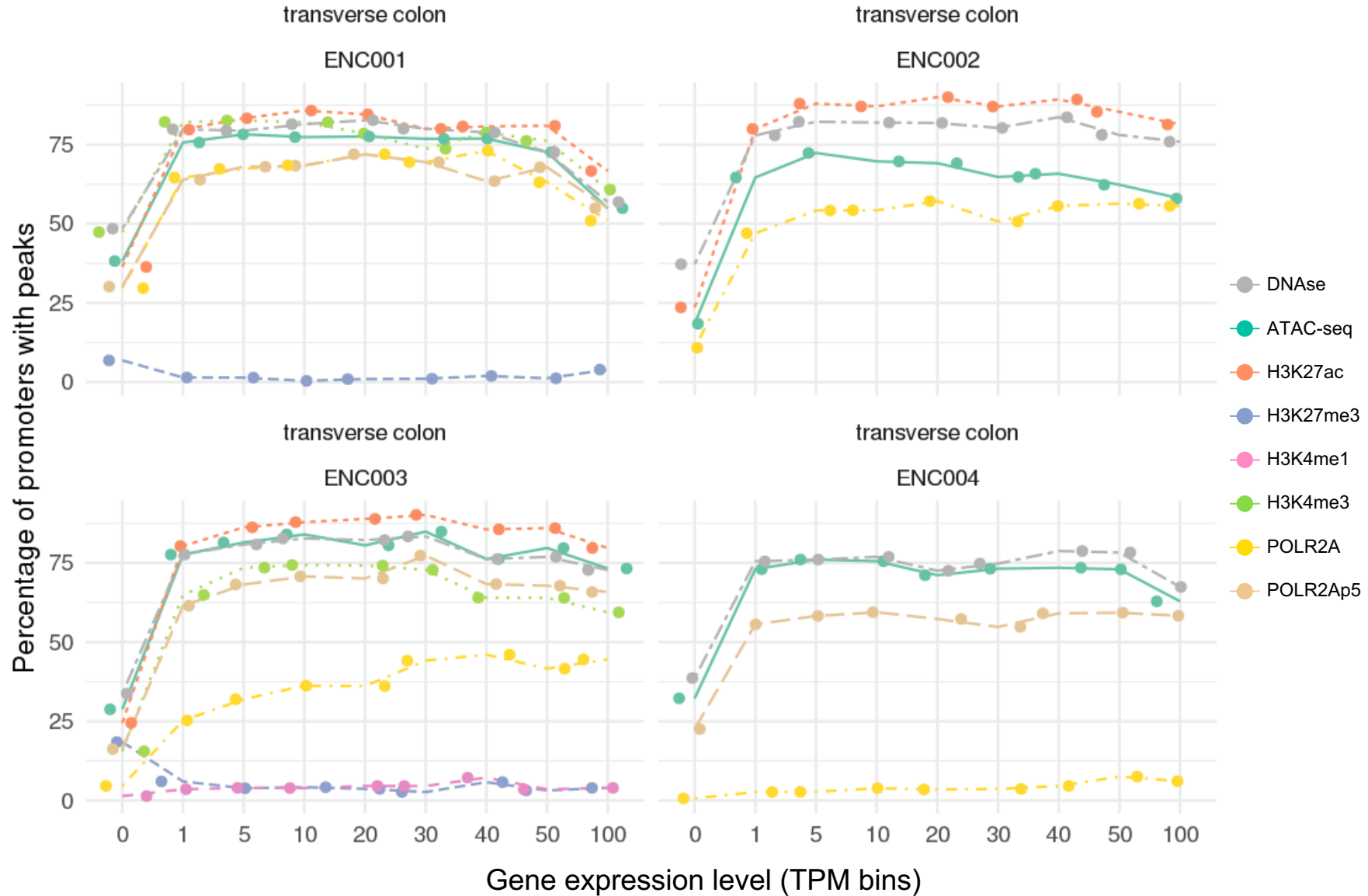
Tissue-specificity of gene vs pseudogene expression



Testis



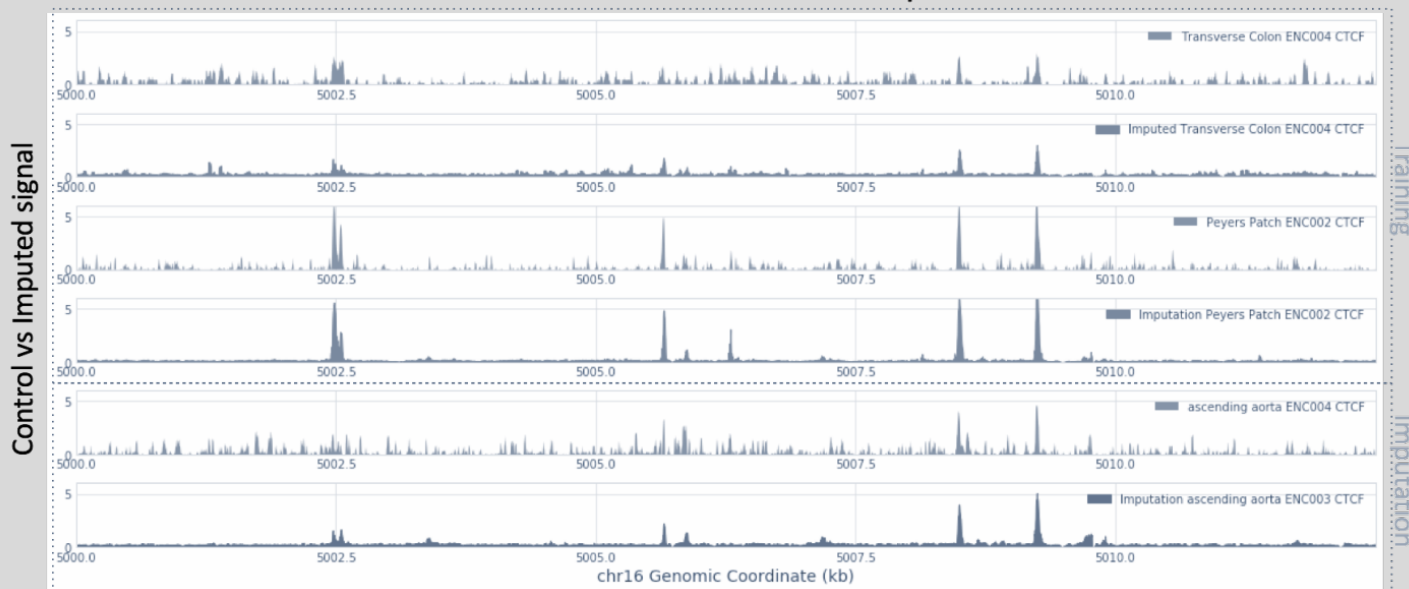
Comparing chromatin activity with gene expression



Imputation & RCA visualization

Application of Avocado to the EN-TEEx Dataset

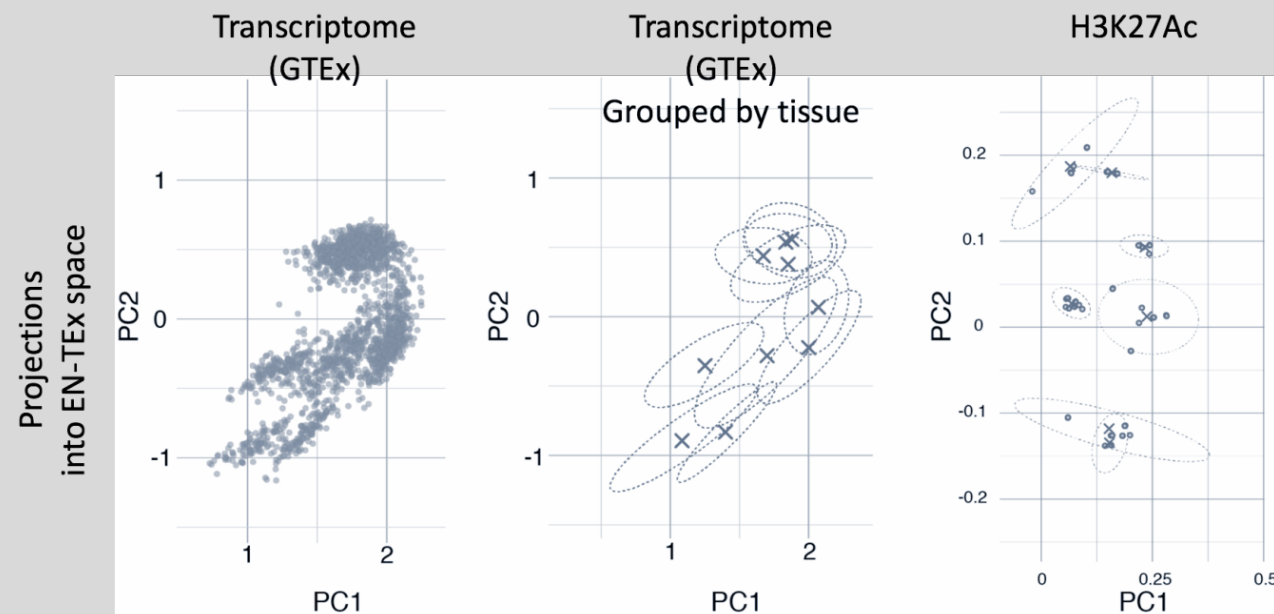
910 sample imputations from Avocado (Deep learning model)
Histone modifications + CTCF + RNAPol II



Jacob Schreiber, Bill Noble

Reference Component Analysis (RCA)

Applied to EN-TEEx data, giving consistent transcriptome v epigenome comparison



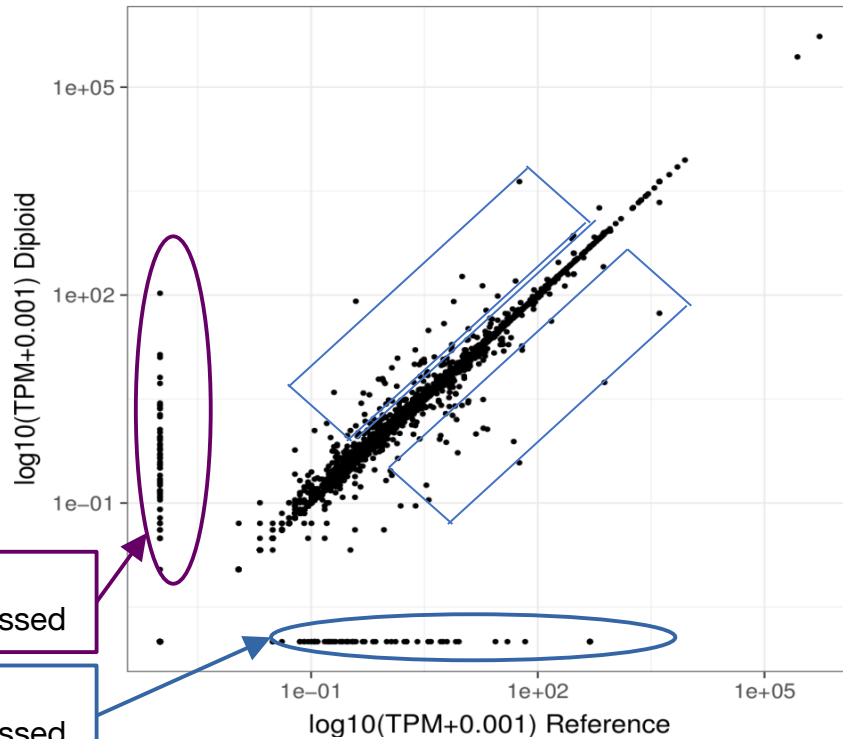
Already discussed in DAC presentation

3. Detailed Effect of Specific Genomic Variants from the Personal Genome

Diploid personal genomes: Effect of variants on Gene Expression

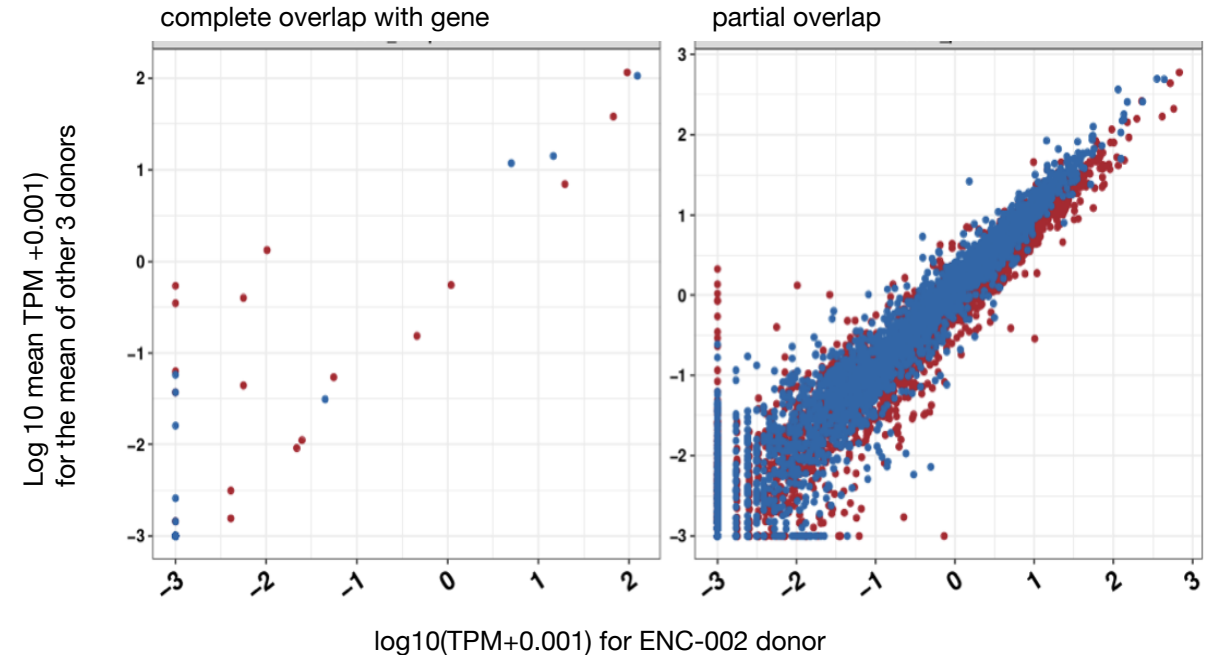
Genes quantified differently when mapped to diploid personal genome compared to the reference genome

ENC-002 stomach
($p_r=0.99$; $s_r=1$)



- 728 (252 protein-coding) genes quantified differently
- Examples of protein-coding genes:
 - with significantly higher expression across multiple tissues in personal genome: HLA-DRB1, HLA-DQA1, HLA-DQB1, IGHV4-31
 - in the reference genome: NBPF26, UGT2B15, HLA-DRB5, AC073333.1, RP11-514P8.6, FOXD4L3, NANOG, CTD-3126B10.5, RIMBP3

Effect of deletions on gene expression

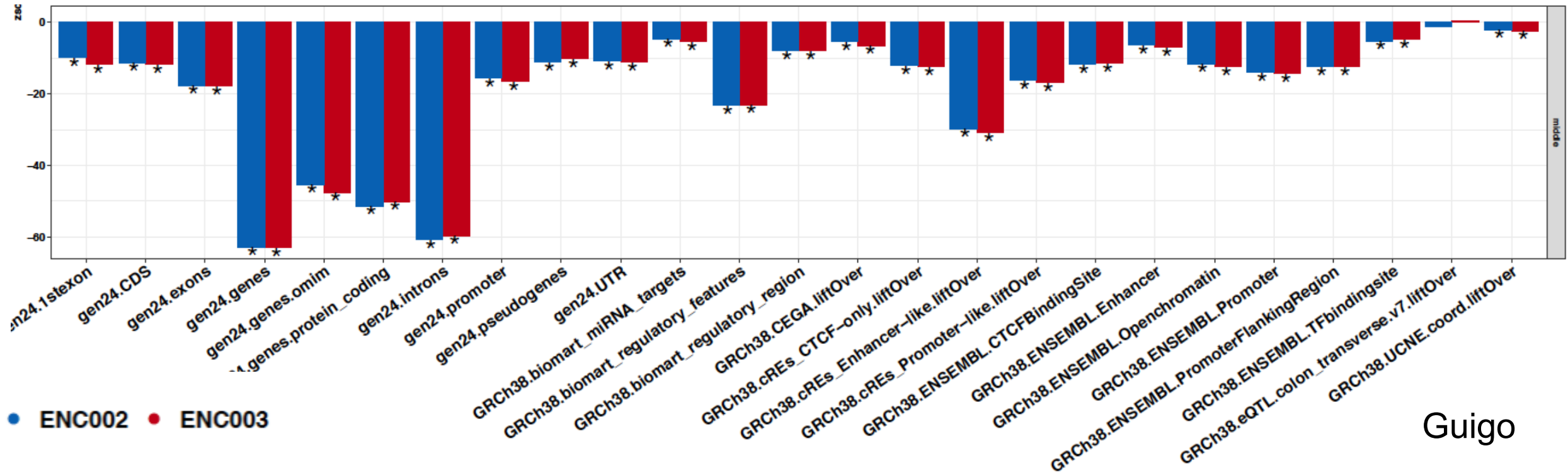


Functional elements are depleted in SVs

Calculated **partial** (1bp) and **complete** overlap between SVs and genomic elements

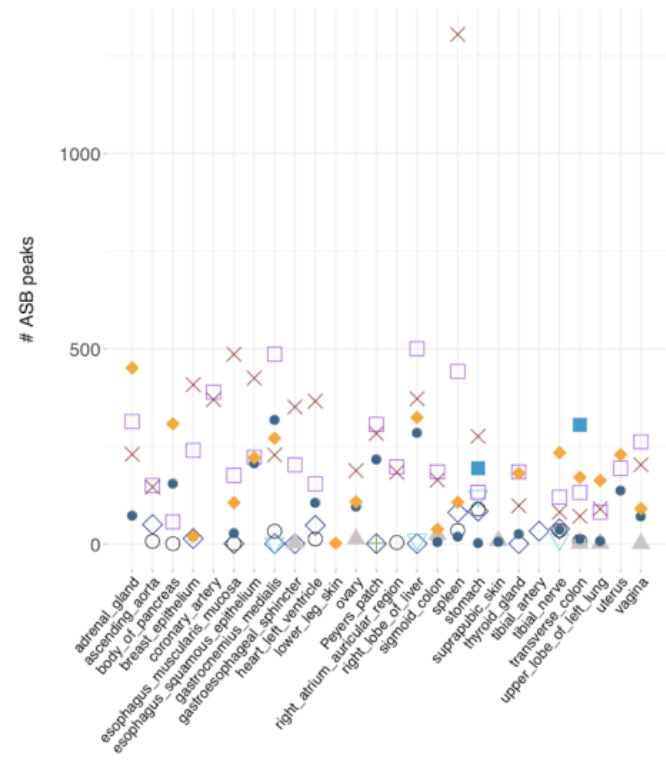
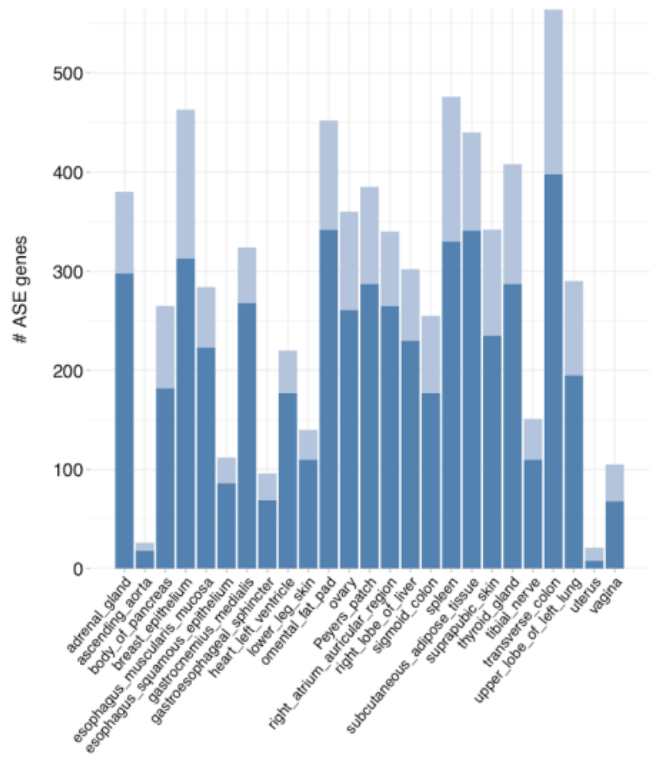
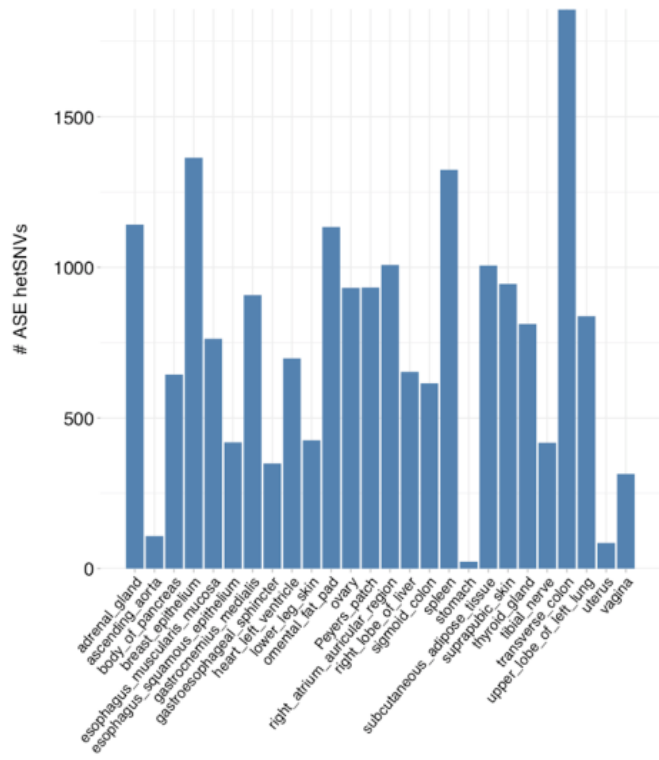


Complete overlap



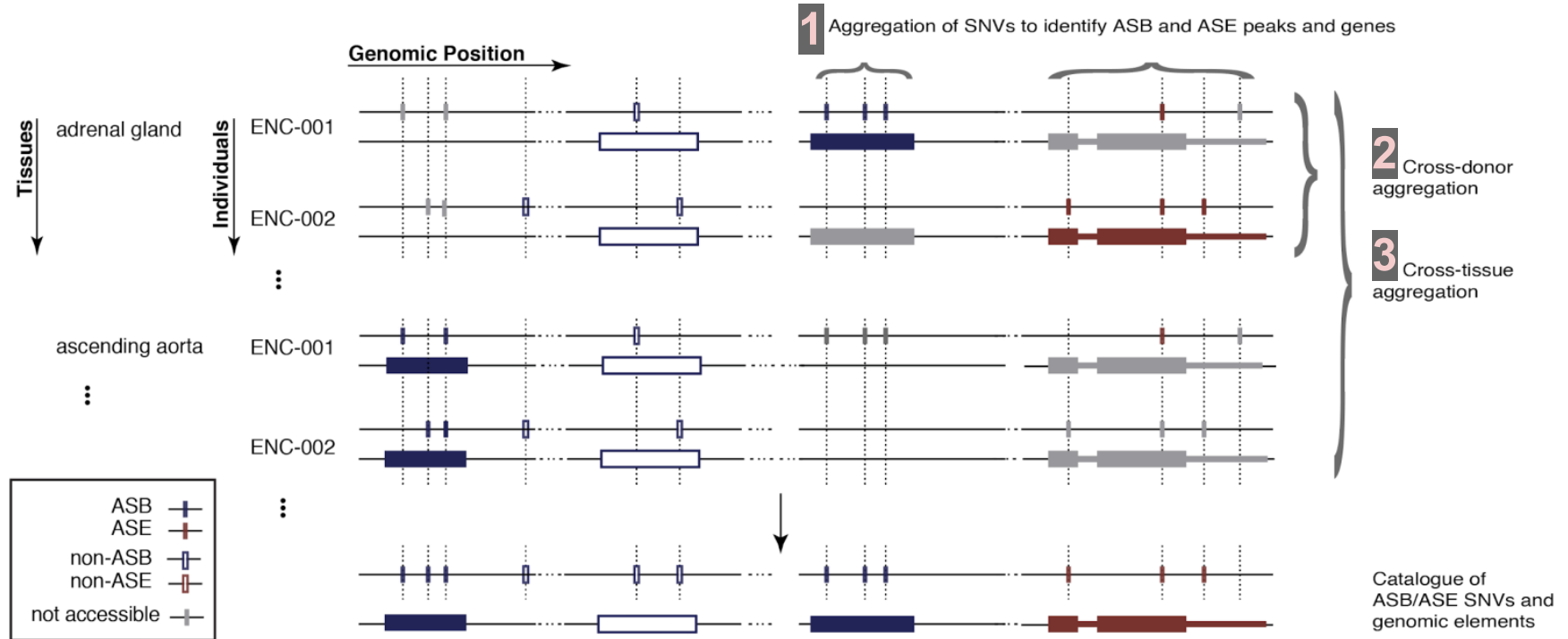
Generating a catalogue of SNVs, genes, & binding sites associated with allelic activity across different tissues

Counts of ASE hetSNVs and genes and ASB peaks; all tissues ENC-003

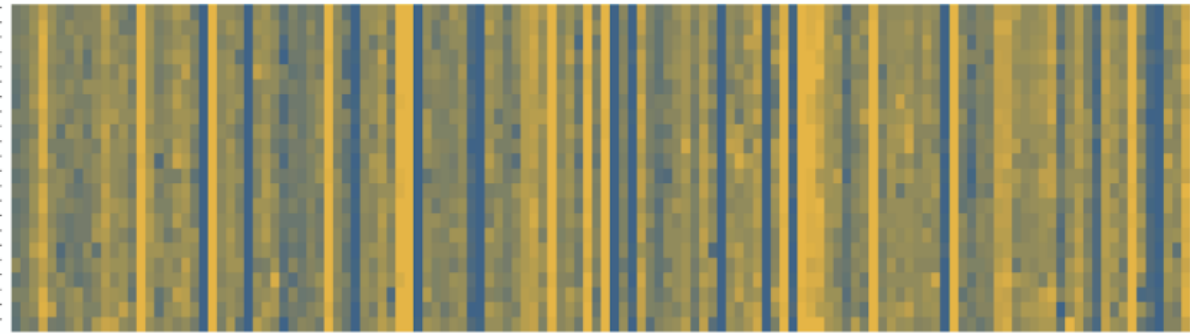
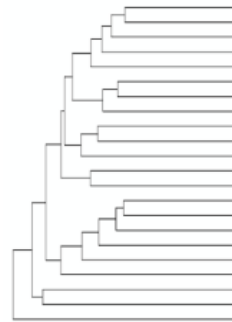
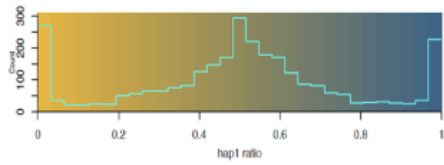


■ protein coding ■ other
◆ CTCF ■ POLR2AphosphoS5 □ H3K27ac ▽ H3K36me3 × H3K4me3
● POLR2A ▲ EP300 ○ H3K27me3 + H3K4me1 ◇ H3K9me3

Developing an integrated cross element (1), individual (2) & tissue (3) ASE and ASB annotation for genomic features (genes & ccREs)

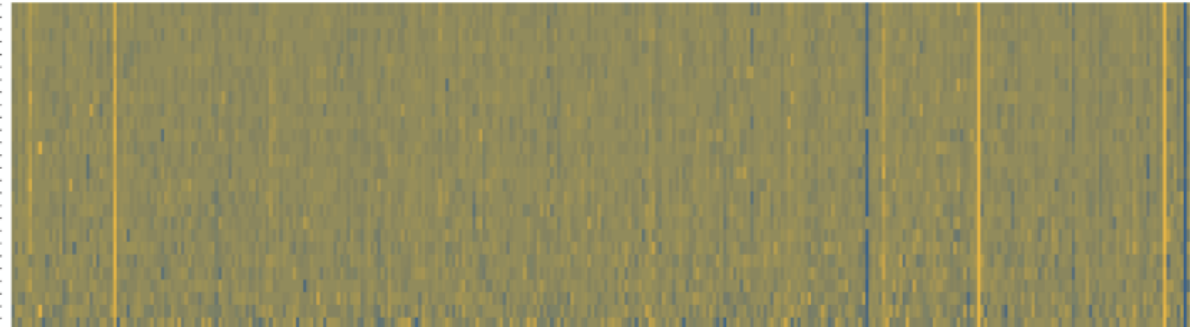
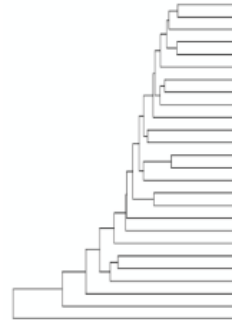
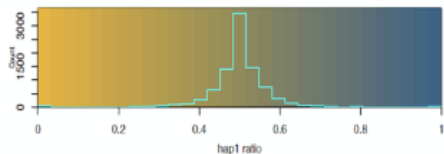


ASB (H3K27ac) in ccREs across ENC-003 tissues



ENC-003_vagina
 ENC-003_breast_epithelium
 ENC-003_coronary_artery
 ENC-003_gastroesophageal_splincter
 ENC-003_adrenal_gland
 ENC-003_heart_left_ventricle
 ENC-003_sigmoid_colon
 ENC-003_esophagus_muscularis_mucosa
 ENC-003_tibial_nerve
 ENC-003_stomach
 ENC-003_gastrocnemius_medialis
 ENC-003_transverse_colon
 ENC-003_right_atrium_auricular_region
 ENC-003_right_lobe_of_liver
 ENC-003_esophagus_squamous_epithelium
 ENC-003_Peyers_patch
 ENC-003_spleen
 ENC-003_upper_lobe_of_left_lung
 ENC-003_thyroid_gland
 ENC-003_uterus
 ENC-003_ascending_aorta
 ENC-003_body_of_pancreas

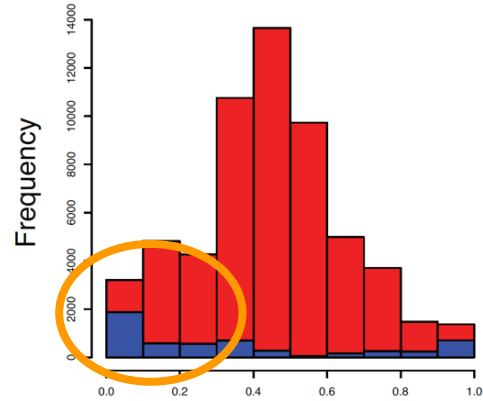
ASE in protein-coding genes across ENC-003 tissues



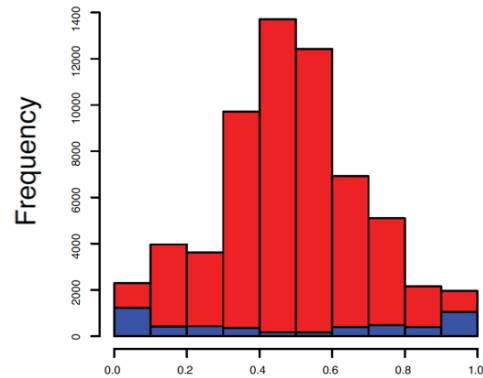
ENC-003_omental_fat_pad
 ENC-003_breast_epithelium
 ENC-003_tibial_nerve
 ENC-003_thyroid_gland
 ENC-003_upper_lobe_of_left_lung
 ENC-003_suprapubic_skin
 ENC-003_sigmoid_colon
 ENC-003_ovary
 ENC-003_spleen
 ENC-003_subcutaneous_adipose_tissue
 ENC-003_esophagus_muscularis_mucosa
 ENC-003_vagina
 ENC-003_transverse_colon
 ENC-003_Peyers_patch
 ENC-003_body_of_pancreas
 ENC-003_right_atrium_auricular_region
 ENC-003_heart_left_ventricle
 ENC-003_adrenal_gland
 ENC-003_right_lobe_of_liver
 ENC-003_gastrocnemius_medialis
 ENC-003_esophagus_squamous_epithelium
 ENC-003_gastroesophageal_splincter
 ENC-003_lower_leg_skin
 ENC-003_ascending_aorta
 ENC-003_uterus
 ENC-003_stomach

Using personal genomes alleviates reference mapping bias

Reference Genome



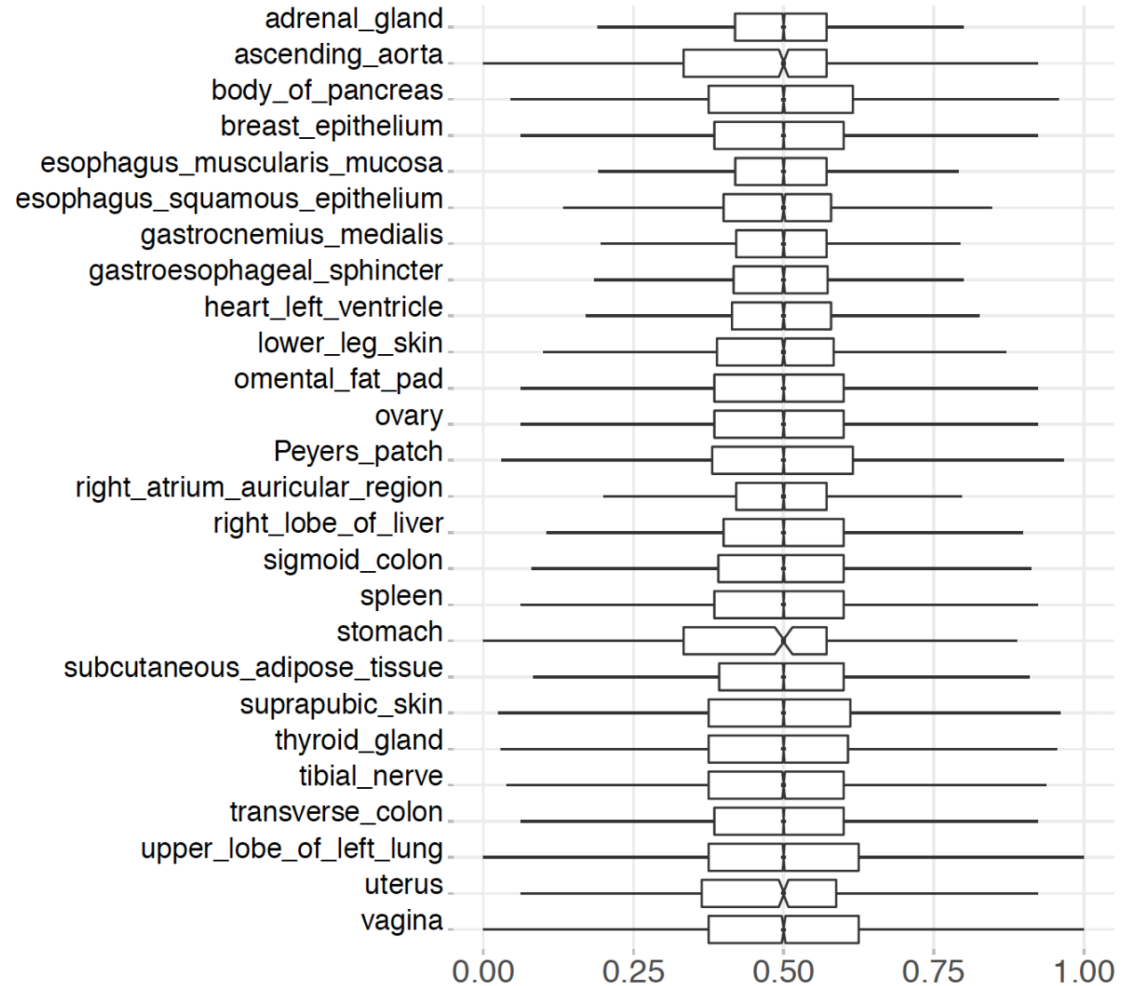
Personal Genome NA12878



alternative allele ratio

(fraction of RNA-Seq reads mapping to alternative allele per heterozygous SNP)

Personal Genome ENC-003; all RNA-seq samples



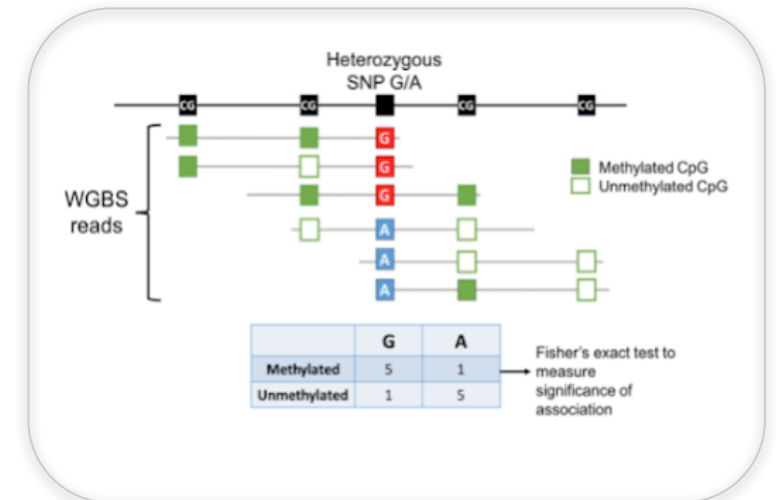
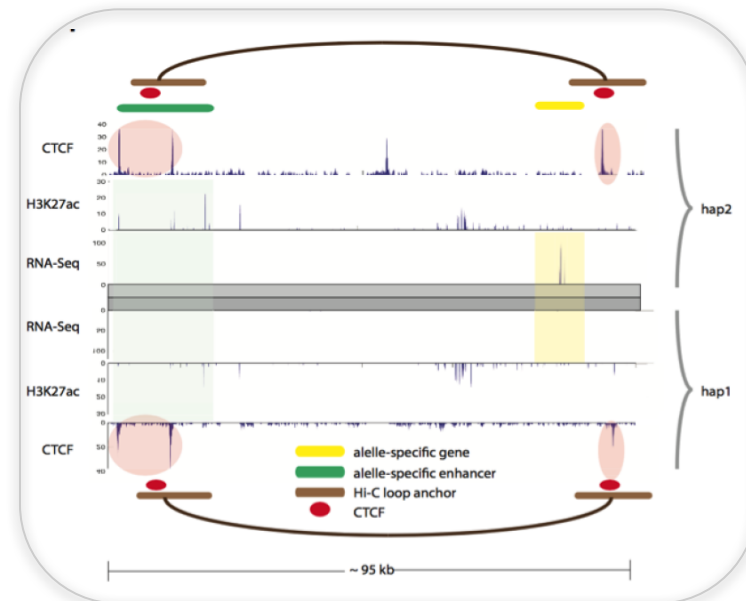
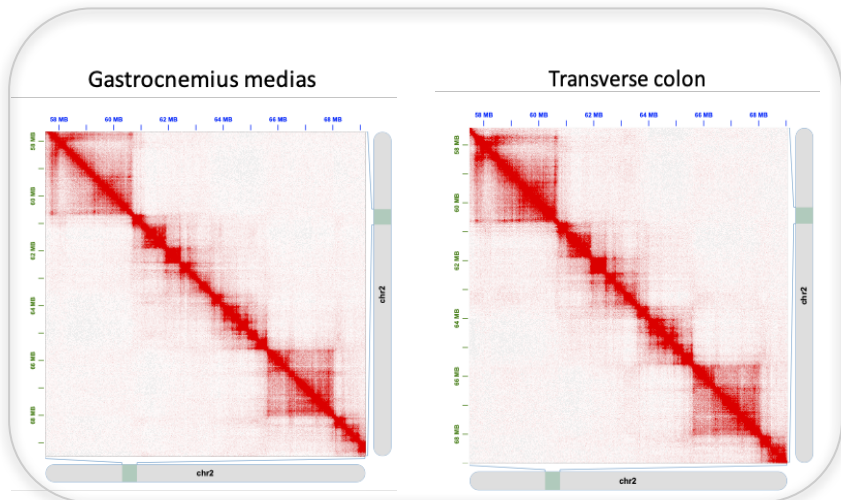
alternative allele ratio

EN-TE_x Summary

- 1) High-quality assembly of personal genomes from integration of long & short read sequencing
 - Determination of many types of SVs
 - Matched data matrix almost filled in
- 2) Great variation of raw transcriptome & epigenome data across individuals and tissues
 - Assay integration gives a more stable picture for epigenetics across individuals
 - Pseudogenes are more tissue specific than genes
- 3) Personal genomic variants have clear impact on transcription
 - Generating a catalogue of allelic elements and genes, by integrating across individuals and tissues

Future Directions for the EN-TEEx Analysis

- Variation in 3D organization of genome across tissues and individuals
 - loops
 - domains
 - compartments
- Allele-specific Hi-C
- Variation in methylation across tissues and individuals
- Allele-specific methylation



EN-TEX

HudsonAlpha Institute

Richard M. Myers Lab

Stanford University

J. M Cherry Lab

I. Gabdank DCC

M. Snyder Lab

Broad Institute

Bradley E. Bernstein

Charles B. Epstein

Baylor College of Medicine

Erez L. Aiden Lab

Salk Institute

Joseph R. Ecker Lab

Yale University

Mark B. Gerstein

Joel Rozowsky

Fábio C. P. Navarro

Gamze Gursoy

Timur R. Galeev

Mengting Gu

Yucheng T. Yang

Chengfei Yan

Berkeley National Laboratory

Len A Pennacchio Lab

California Institute of Technology

Barbara Wold Lab

Pennsylvania State University

Ross Hardison Lab

University of California San Diego

Bing Ren Lab

Cold Spring Harbor Laboratory

Thomas Gingeras

Alexander Dobin

Jorg Drenkow

Cassidy Danyko

Jesse Gillis

Carrie A. Davis

Chris Zaleski

Alex Scavelli

NHGRI (funding agency)

Mike Pazin

Elise Feingold

Dan Gilchrist

Michael Pagan

Eileen Cahill

Center for Genomic Regulation

Roderic Guigo

Julien Lagarde

Anna Vlasova

Alessandra Breschi

Sarah Djebali

Dimitri Pervouchine

University of Connecticut

Brenton Graveley Lab

University of Washington

John A. Stamatoyannopoulos

Rajinder Kaul

Jacob Schreiber

Bill Noble

University of Massachusetts

Zhiping Weng

Jill Moore

Henry Pratt

University of California, Irvine

Ali Mortazavi

Rabi Murad

Johns Hopkins University

Michael Schatz

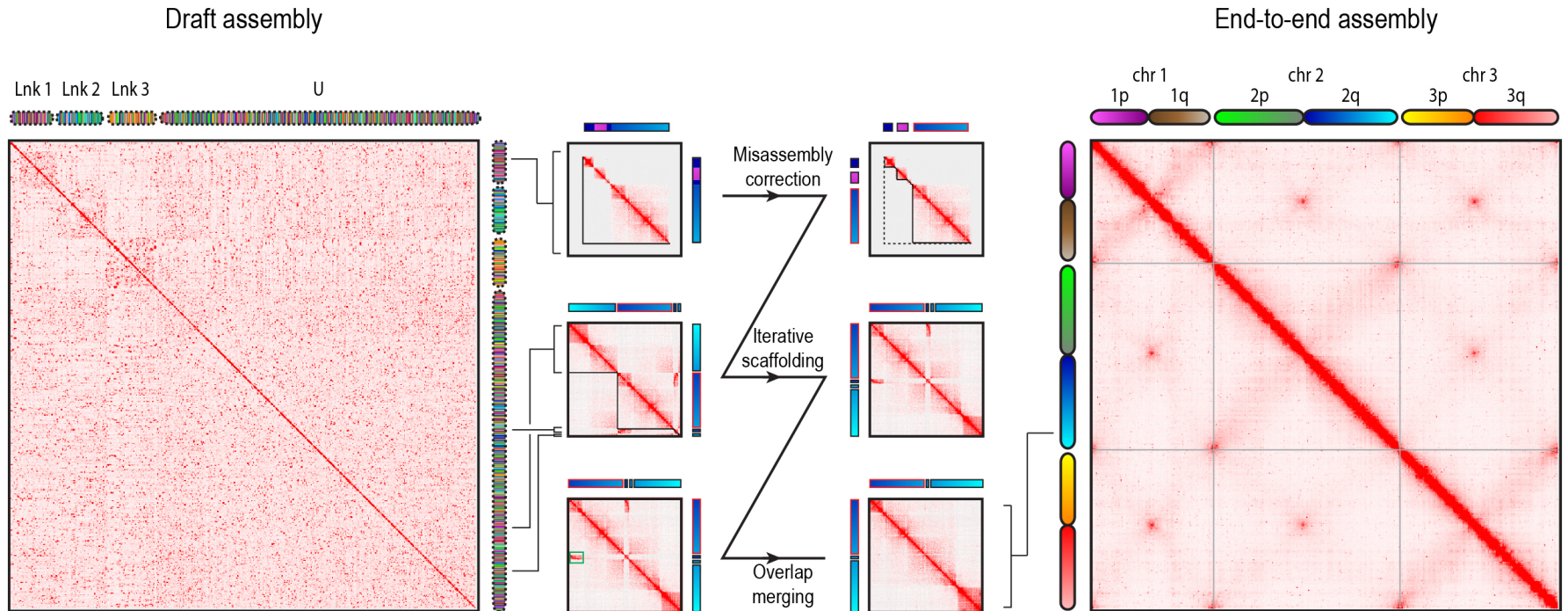
Fritz Sedlazeck

Han Fang

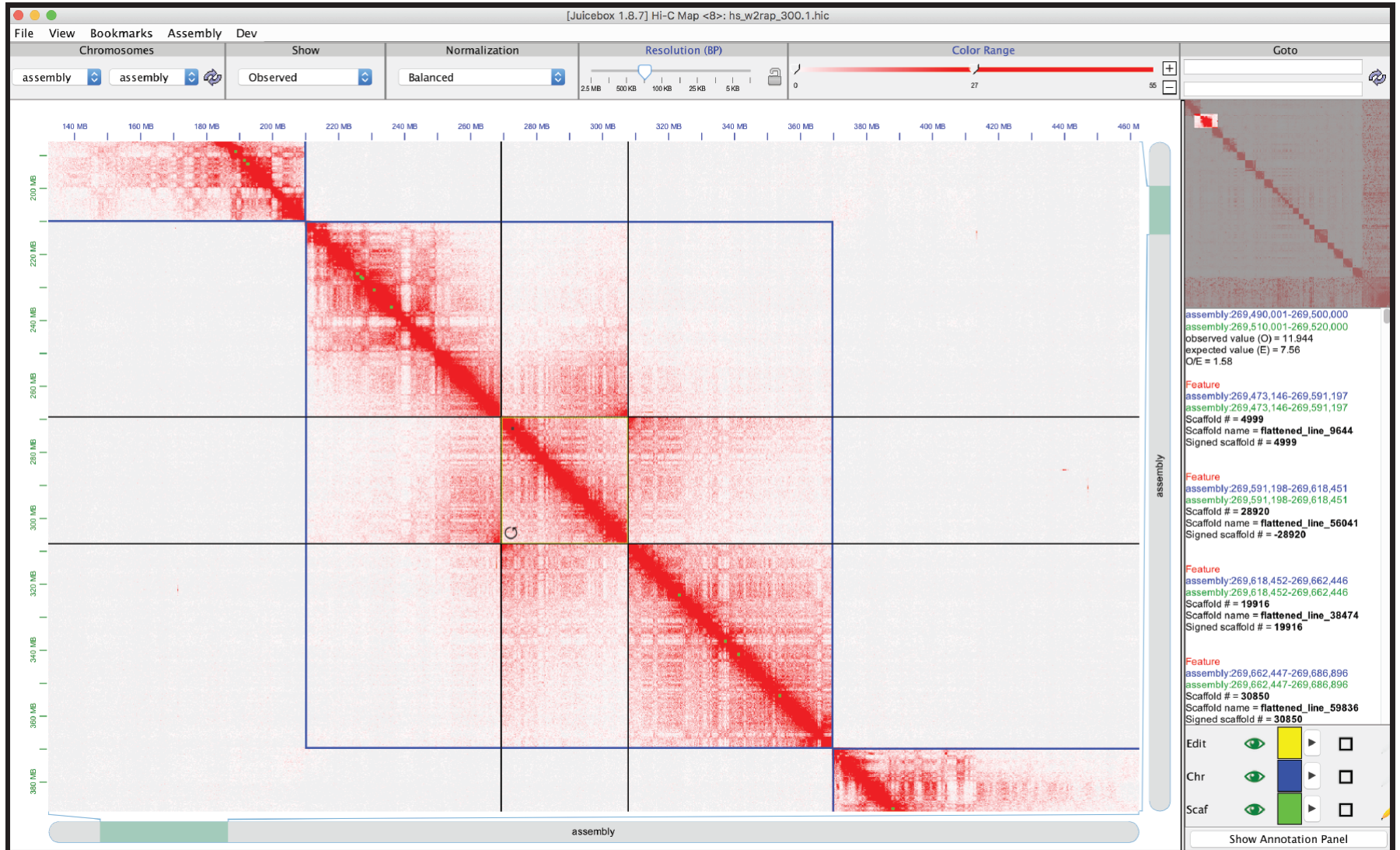
Maria Nattestad

PERSONALIZED GENOMES

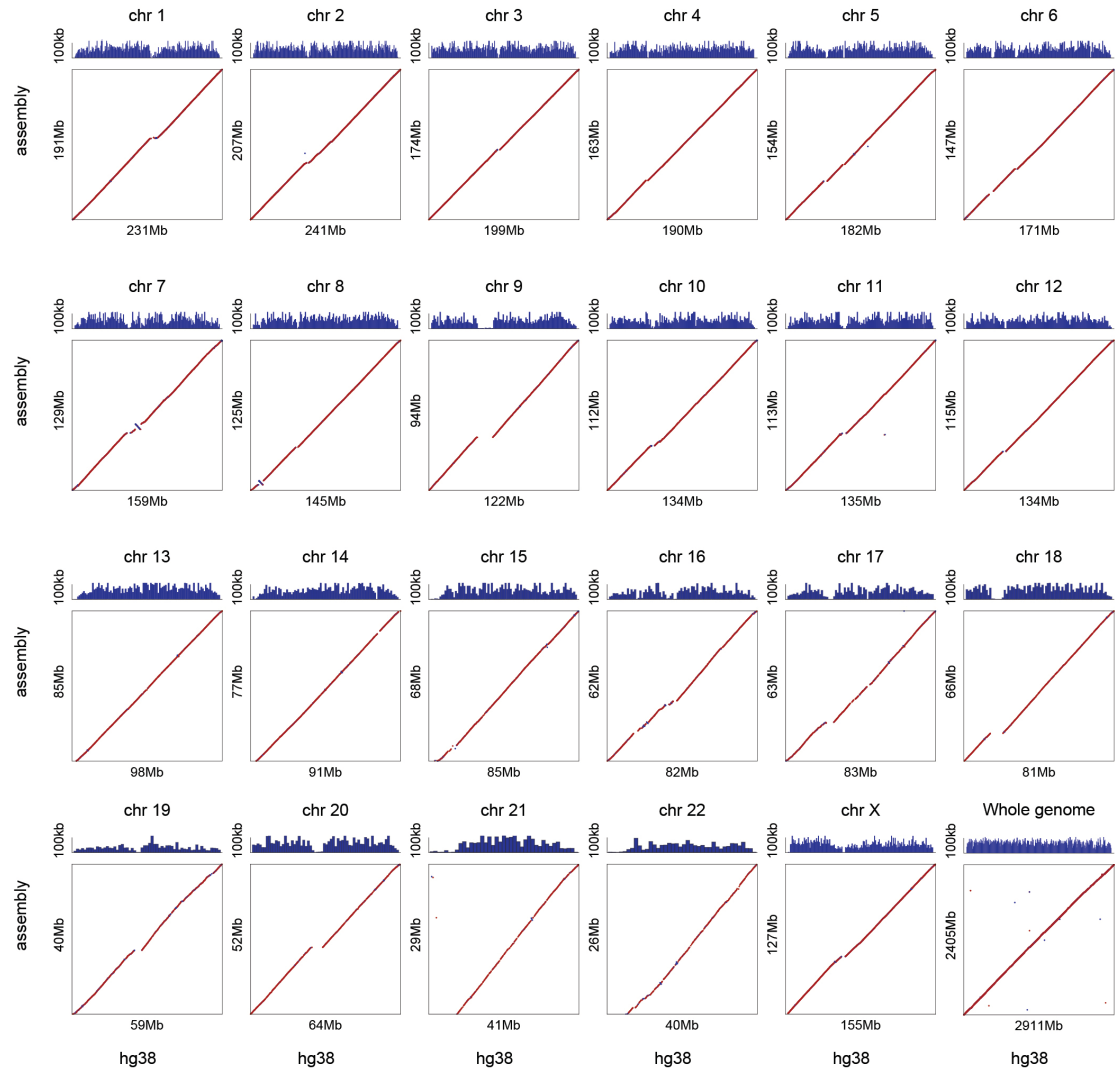
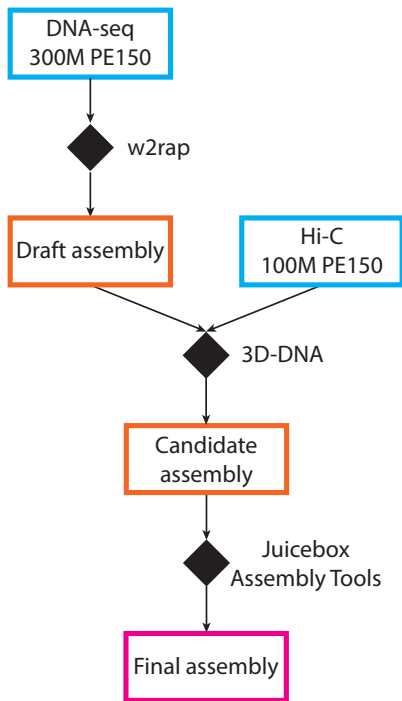
3D-DNA: AUTOMATIC 3D DE NOVO GENOME ASSEMBLY



JUICEBOX ASSEMBLY TOOLS



HERE IS A HUMAN GENOME WE MADE FOR \$1000



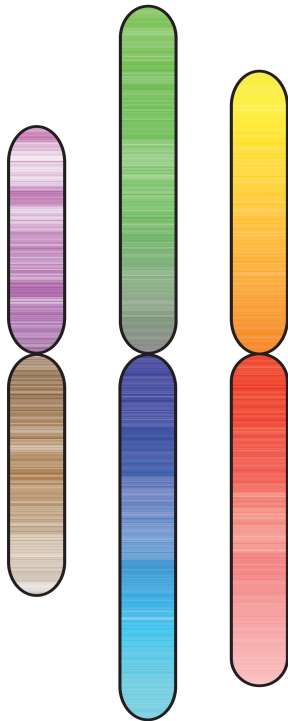
END-TO-END (GAPPY) ASSEMBLY: \$1000

PLATINUM ASSEMBLY: ~\$20K

AaegL4:

Sanger DNA-Seq + Illumina Hi-C

Contig NG50: 82 Kb



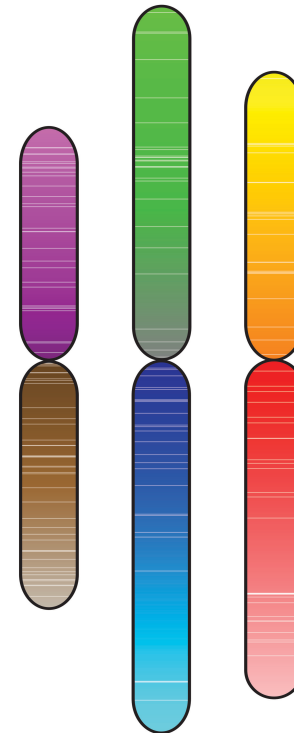
Number of gaps: 25,582

Dudchenko et al., *Science*, 2017

AaegL5:

PacBio DNA-Seq + Illumina Hi-C

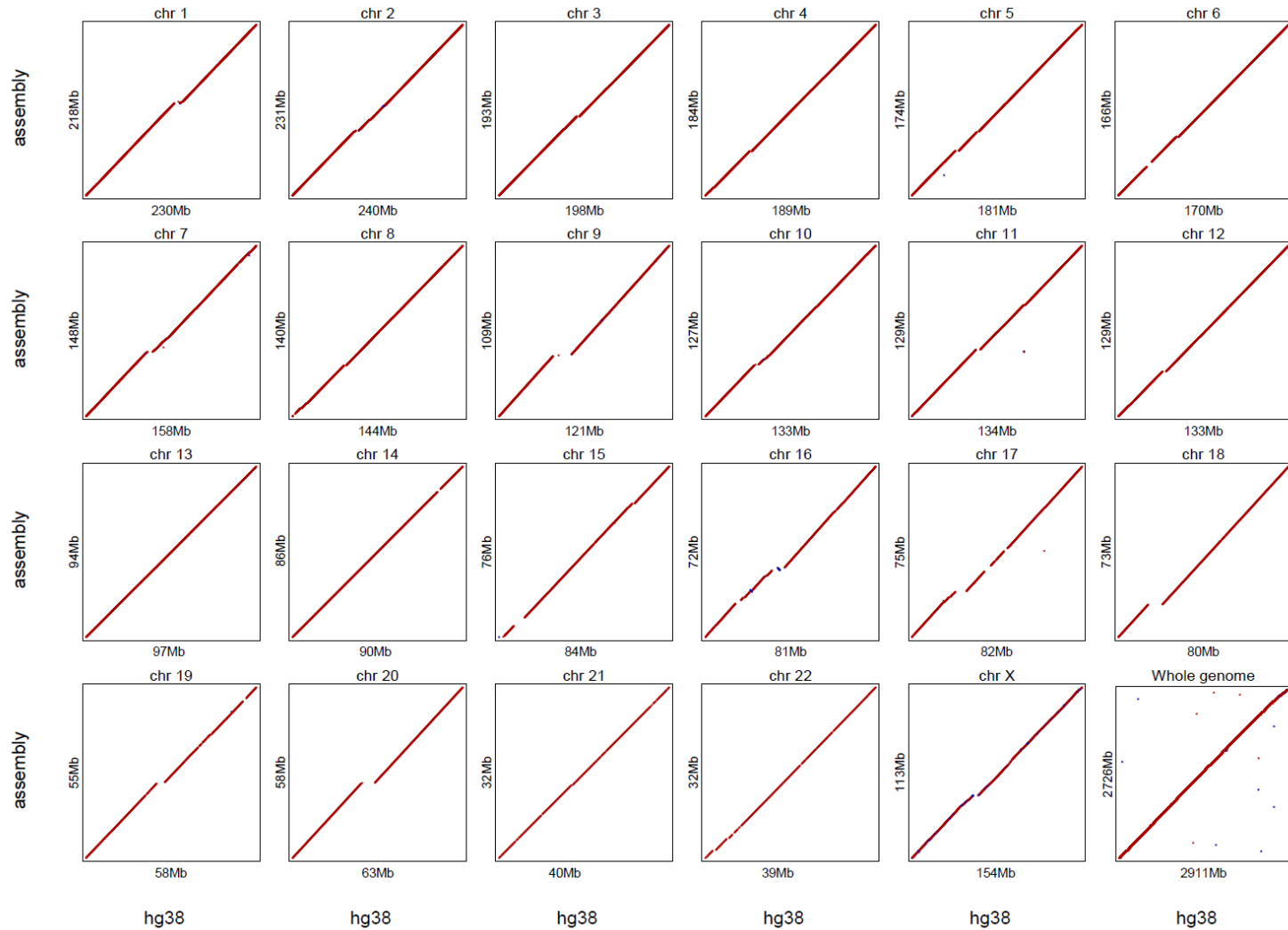
Contig NG50: 11.8 Mb



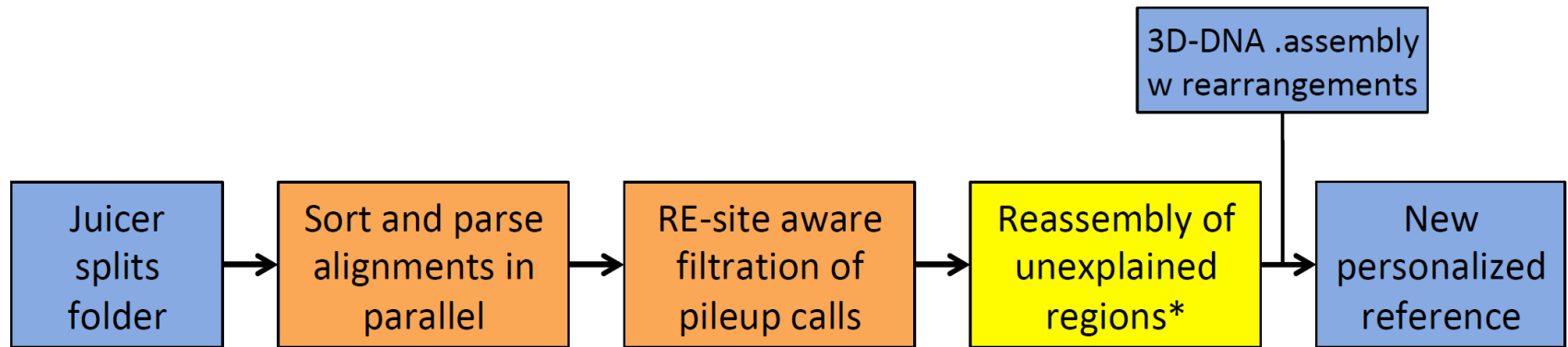
Number of gaps: 173

Matthews, Dudchenko, Kingan et al., *Nature*, 2018

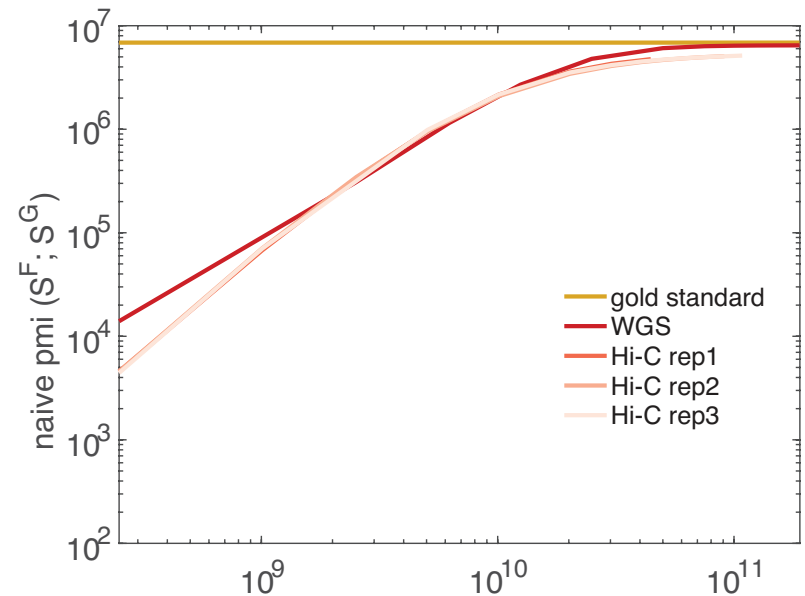
OPPORTUNISTIC PROJECT (\$0!): ENTEX 002 & ENTEX 003



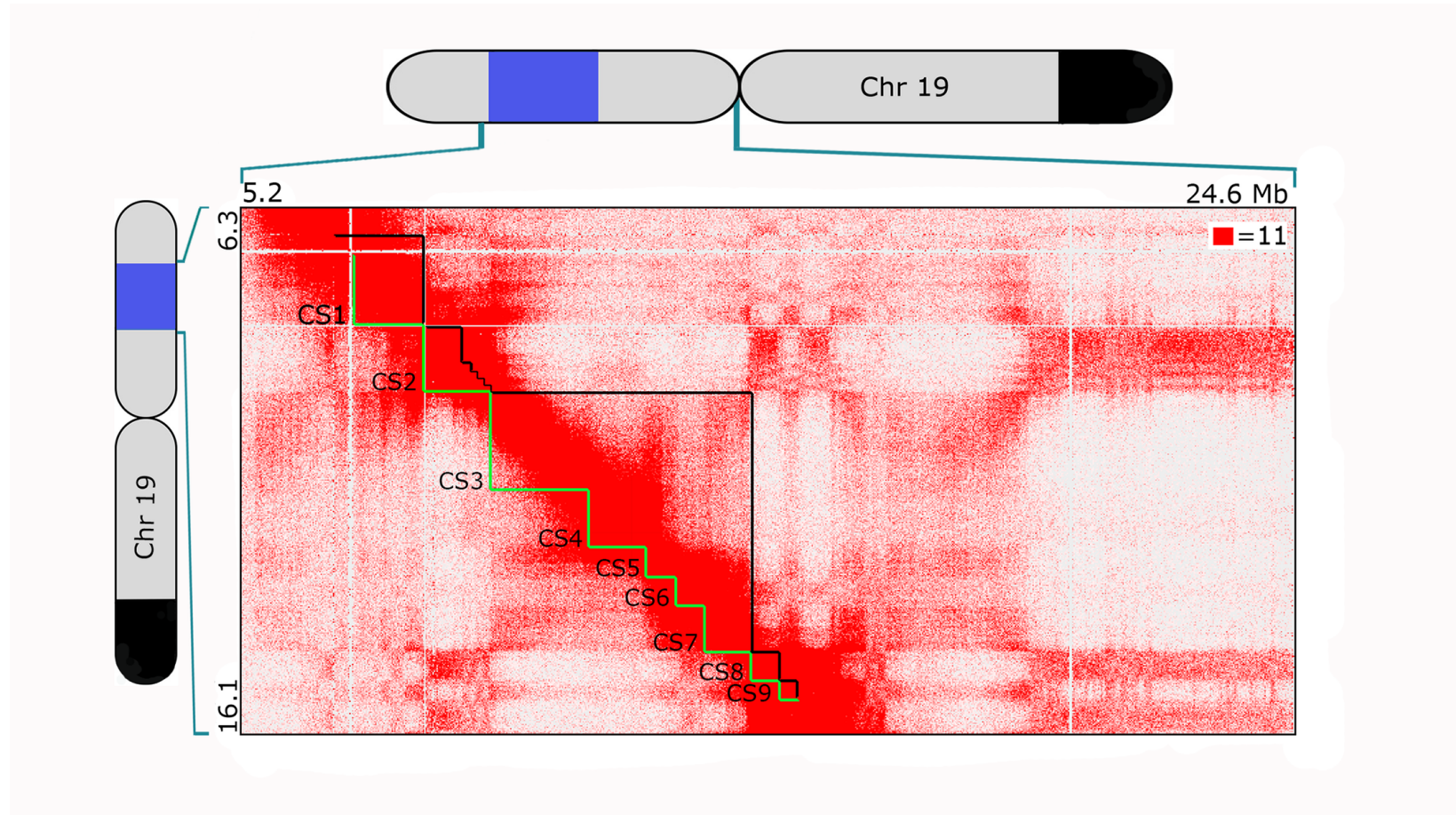
OPPORTUNISTIC PROJECT (\$0!): HI-C DATA YIELDS GENOME-WIDE SNP CALLING



- Available personalized genome modules for Hi-C data include:
 - variant calling
 - phasing
 - assembly



ENCODE4 PILOT PROJECT (\$1000): PGP1 & HEPG2



Data generation in progress

QUESTIONS FOR DISCUSSION

- What is the value of the personal genome & having functional genomic assays matched to it?
- How to think about genomic element calls relative to substantial epigenetic variation across individuals?
- How best to fill out the EN-TE_x data matrix in the coming months?
 - What types of integrated calculations are best done over the EN-TE_x matrix?
- How to best include opportunistic Hi-C personalized genome data into ENCODE?
- Should we prioritize additional reference cell lines & samples for *de novo* genome assembly efforts?
 - Backfilling? (K562, IMR90, HCT-116, HMEC, NHEK etc.)
 - EN-TE_x 1&4?