

Regulatory genomics to interpret complex disease genetics

Luke Ward

Manolis Kellis lab, MIT

July 18, 2013

Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

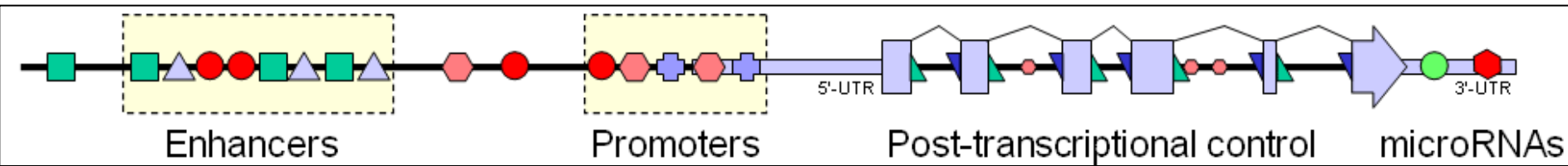
Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

There are lots of noncoding regulatory mechanisms for disease genetics

Noncoding element disrupted	Molecular function and effect of mutations.	Disease association
Splice-junction and splicing-enhancer	Splicing is constitutive for some transcripts and highly tissue-specific for others, relying on both canonical sequences at the exon-intron junction as well as weakly-specified sequence motifs distributed throughout the transcript. Mutations affecting constitutive splice sites can have an effect similar to nonsense or missense mutations, resulting in aberrantly included introns or skipped exons, sometimes resulting in nonsense-mediated decay (NMD).	Splicing regulatory variants are implicated in several diseases. A recent analysis suggests that the majority of disease-causing point mutations in OMIM may exert their effects through splicing. Alternative splice site variants in the <i>WT1</i> gene are involved in Frasier Syndrome (FS) Skipping of exon 7 of the <i>SMN</i> gene is involved in spinal muscular atrophy (SMA)
Sequences regulating translation, stability, and localization	Sequences in the 5'-untranslated regions (UTRs) of mRNAs can influence translation regulation, such as upstream ORFs, premature AUG or AUC codons, and palindromic sequences that form inhibitory stem loops. Sequence motifs in the 3'-UTR are recognized by microRNAs and RNA-binding proteins (RBPs).	Loss-of-function mutations in the 5'-UTR of <i>CDKN2A</i> predispose individuals to melanoma. A rare mutation that creates a binding site for the miRNA hs-miR-189 in the transcript of the gene <i>SLITRK1</i> is associated with Tourette's syndrome.
Genes encoding trans-regulatory RNA	Non-coding RNAs participate in a panoply of regulatory functions, ranging from the well-understood transfer and ribosomal RNA to the recently-discovered long non-coding RNAs.	Both rare and common mutations in the gene <i>RMRP</i> encoding an RNA component of the mitochondrial RNA processing ribonuclease have been associated with cartilage-hair hypoplasia Non-coding RNA mutations can cause many other diseases.
Promoter	Promoter regions are an essential component of transcription initiation and the assembly of RNA polymerase and associated regulators. Mutations can affect binding of activators or repressors, chromatin state, nucleosome positioning, and also looping contacts of promoters with distal regulatory elements. Genes with coding disease mutations can also harbor independently-associated regulatory variants that correlate with expression, are bound by proteins in an allele-specific manner, and disrupt or create regulatory motifs	Mutations in the promoter of the HIV1-progression associated gene <i>CCR5</i> , are correlated with expression of the receptor it encodes and bind differentially to at least three transcription factors <i>APOE</i> promoter mutations are associated with Alzheimer's disease Heme oxygenase-1 (<i>HO-1</i>) promoter mutations lead to expression changes and are associated with many diseases
Enhancer	Enhancers are distal regulatory elements that often lie 10,000 to 100,000 nucleotides from the start of their target gene. Mutations within them can disrupt sequence motifs for sequence-specific transcription factors, chromatin regulators, and nucleosome positioning signals. Structural variants including inversions and translocations can disrupt their regulatory activity by moving them away from their targets, disrupting local chromatin conformation, or creating interactions with insulators or repressors that can hinder their action. While it is thought that looping interactions with promoter regions play a role, the rules of enhancer-gene targeting are still poorly understood.	The role of distal enhancers in disease was suggested even before GWAS by many Mendelian disorders for which some patients had translocations or other structural variants far from the promoter In one early study, point mutations were mapped in an unlinked locus in the intron of a neighboring gene, a million nucleotides away from the developmental gene <i>Shh</i> ; this distal locus acted as an enhancer of <i>Shh</i> and recapitulated the polydactyly phenotype in mouse. A number of GWAS hits have been validated as functional enhancers; for example, common variants associated with cancer susceptibility map to a gene desert on chromosome 8, with one SNP demonstrated to disrupt a TCF7L2 binding site and to inhibit long-range activation of the oncogene <i>MYC</i> .
Synonymous mutations within protein-coding sequences	All of the aforementioned regulatory elements can also be encoded within the protein-coding exons themselves. Thus, synonymous mutations within protein-coding regions may be associated with non-coding functions, acting pre-transcriptionally at the DNA level, or post-transcriptionally at the RNA level.	A synonymous variant in the dopamine receptor gene <i>DRD2</i> associated with schizophrenia and alcoholism has been shown to modulate receptor production through differences in mRNA folding and stability.

What we can model with regulatory annotations



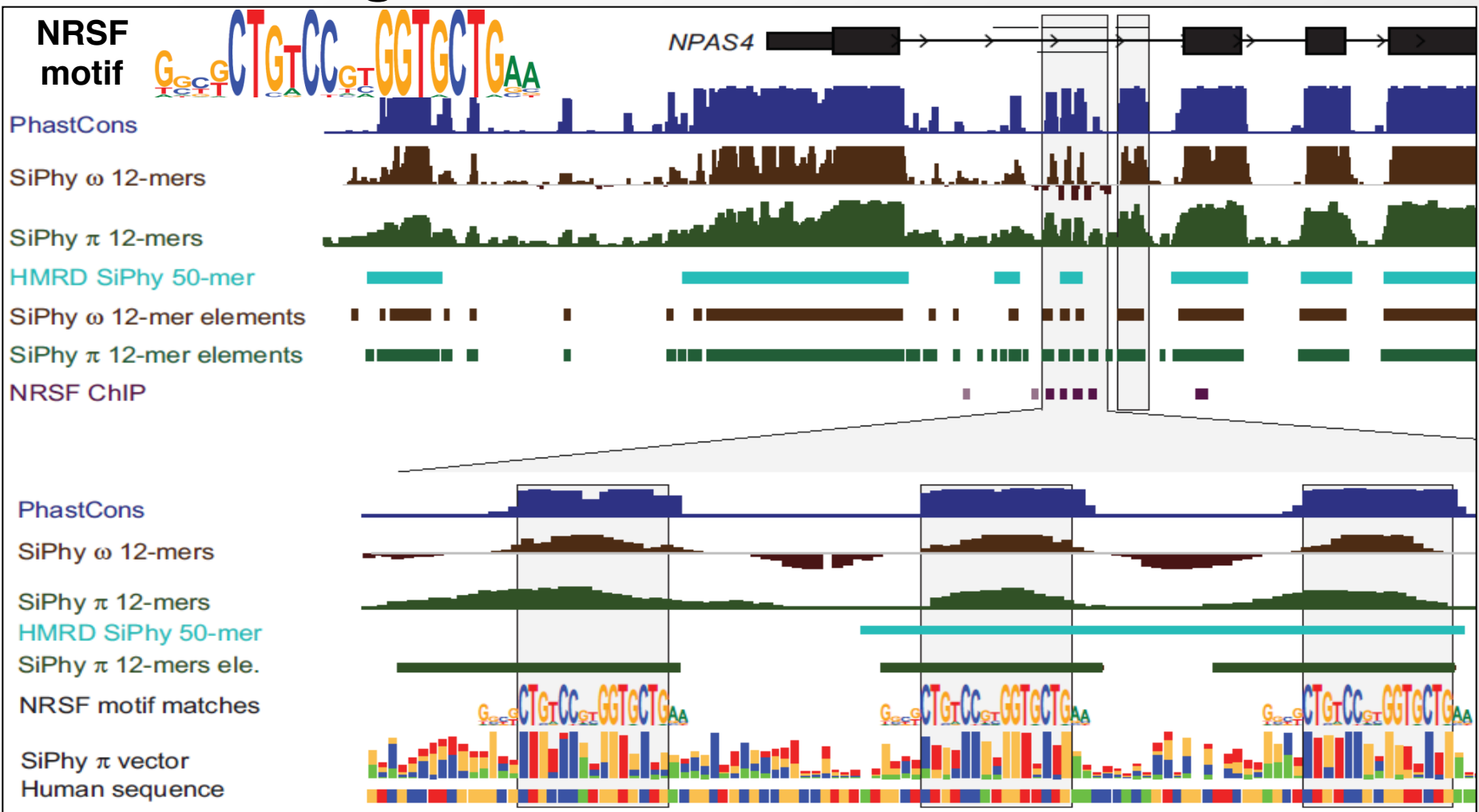
Goal: A systems-level understanding of genomes and gene regulation:

- The regulators: Transcription factors, microRNAs, sequence specificities
 - The regions: enhancers, promoters, and their tissue-specificity
 - The targets: TFs → targets, regulators → enhancers, enhancers → genes
 - The grammars: Interplay of multiple TFs → prediction of gene expression
- ➔ The parts list = Building blocks of gene regulatory networks

Our tools: Comparative genomics & large-scale experimental datasets.

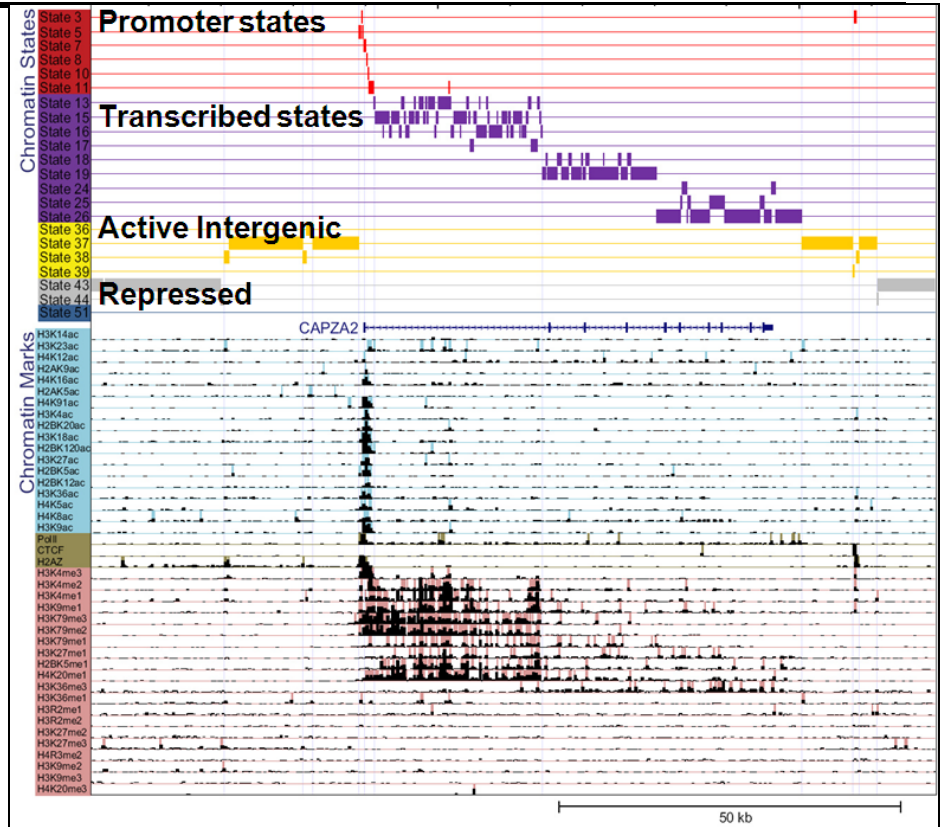
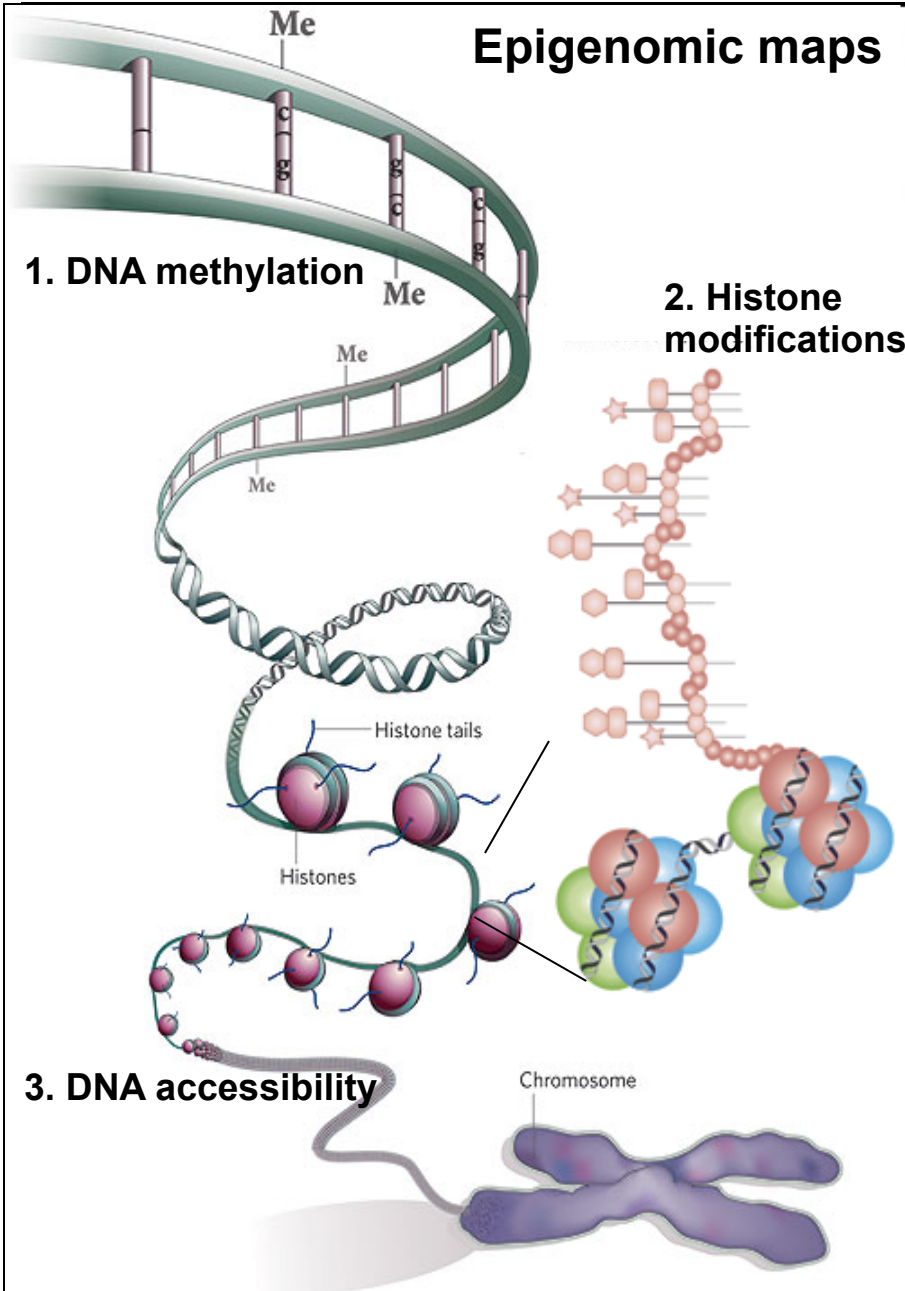
- Evolutionary signatures for coding/non-coding genes, microRNAs, motifs
 - Chromatin signatures for regulatory regions and their tissue specificity
 - Activity signatures for linking regulators → enhancers → target genes
 - Predictive models for gene function, gene expression, chromatin state
- ➔ Integrative models = Define roles in development, health, disease

Measuring constraint at individual nucleotides

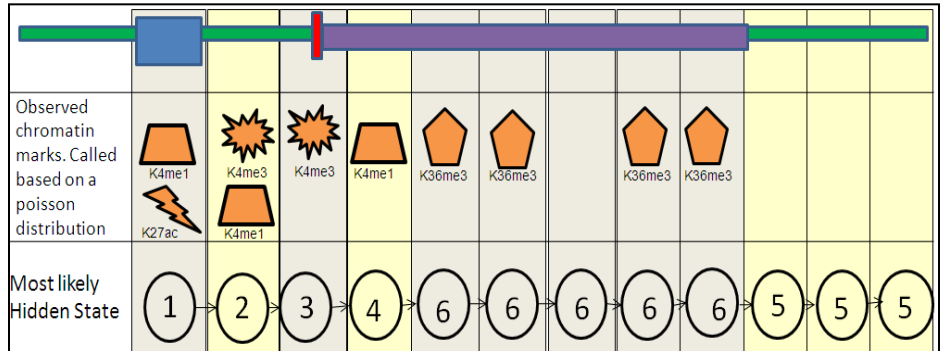


- Reveal individual transcription factor binding sites
- Within motif instances reveal position-specific bias
- More species: motif consensus directly revealed

Chromatin signatures for genome annotation



Ernst et al Nature Biotech 2010



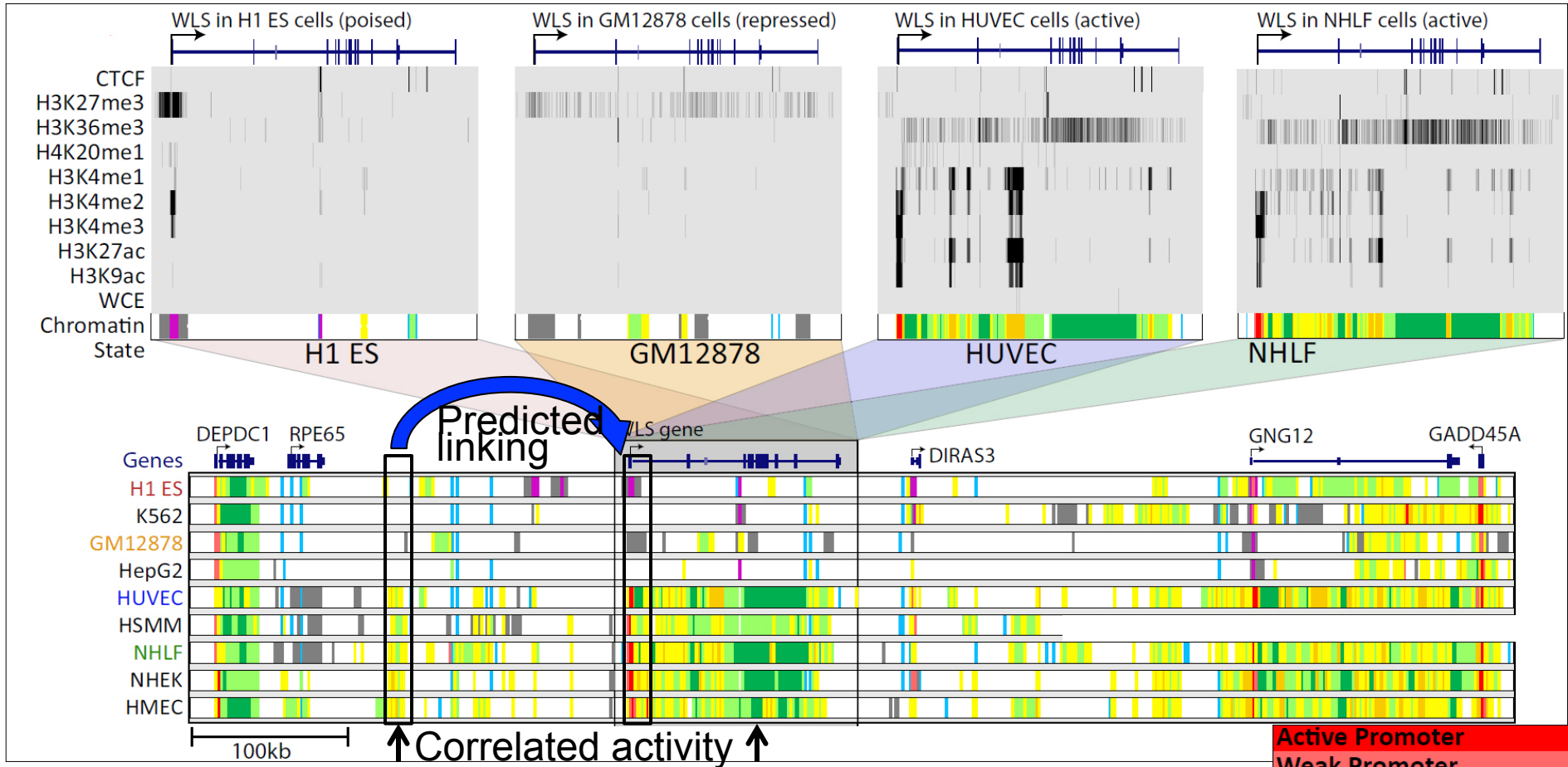
Multi-variate HMM captures combinations

Epigenomics Roadmap: 90 reference epigenomes



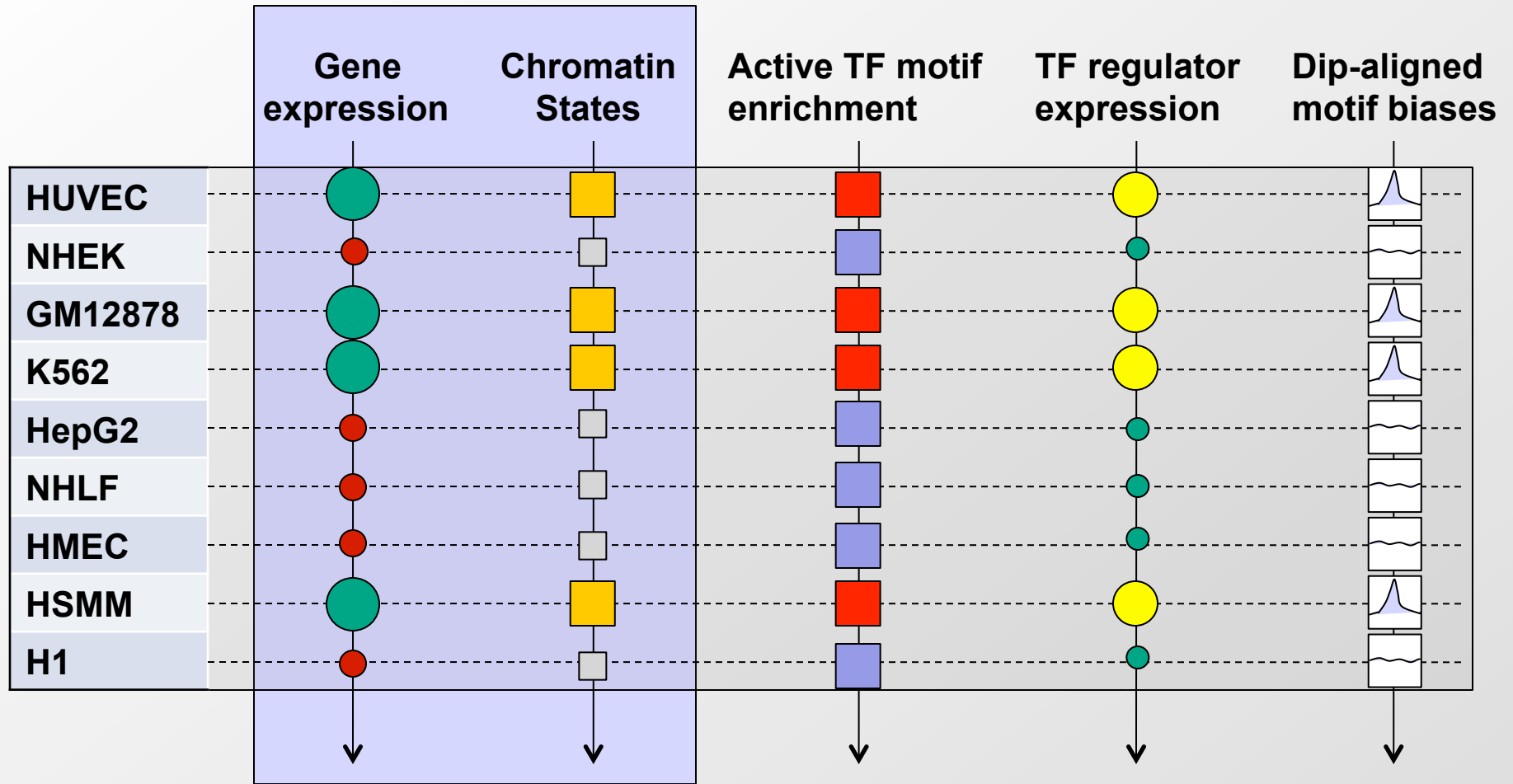
Interpret GWAS, global effects, reveal relevant cell types

Chromatin state dynamics reveal linking/regulators

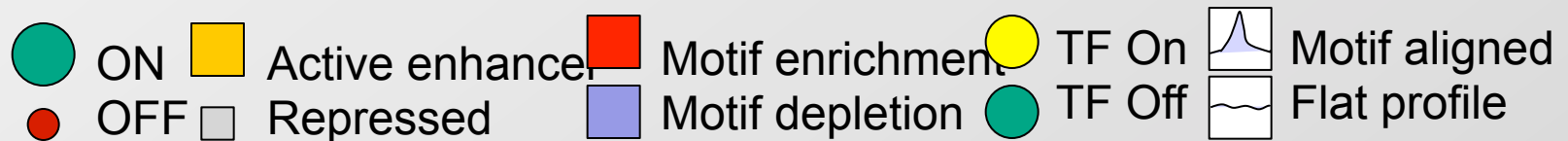


- Single annotation track for each cell type
- Summarize cell-type activity at a glance
- Study activity pattern across tissues

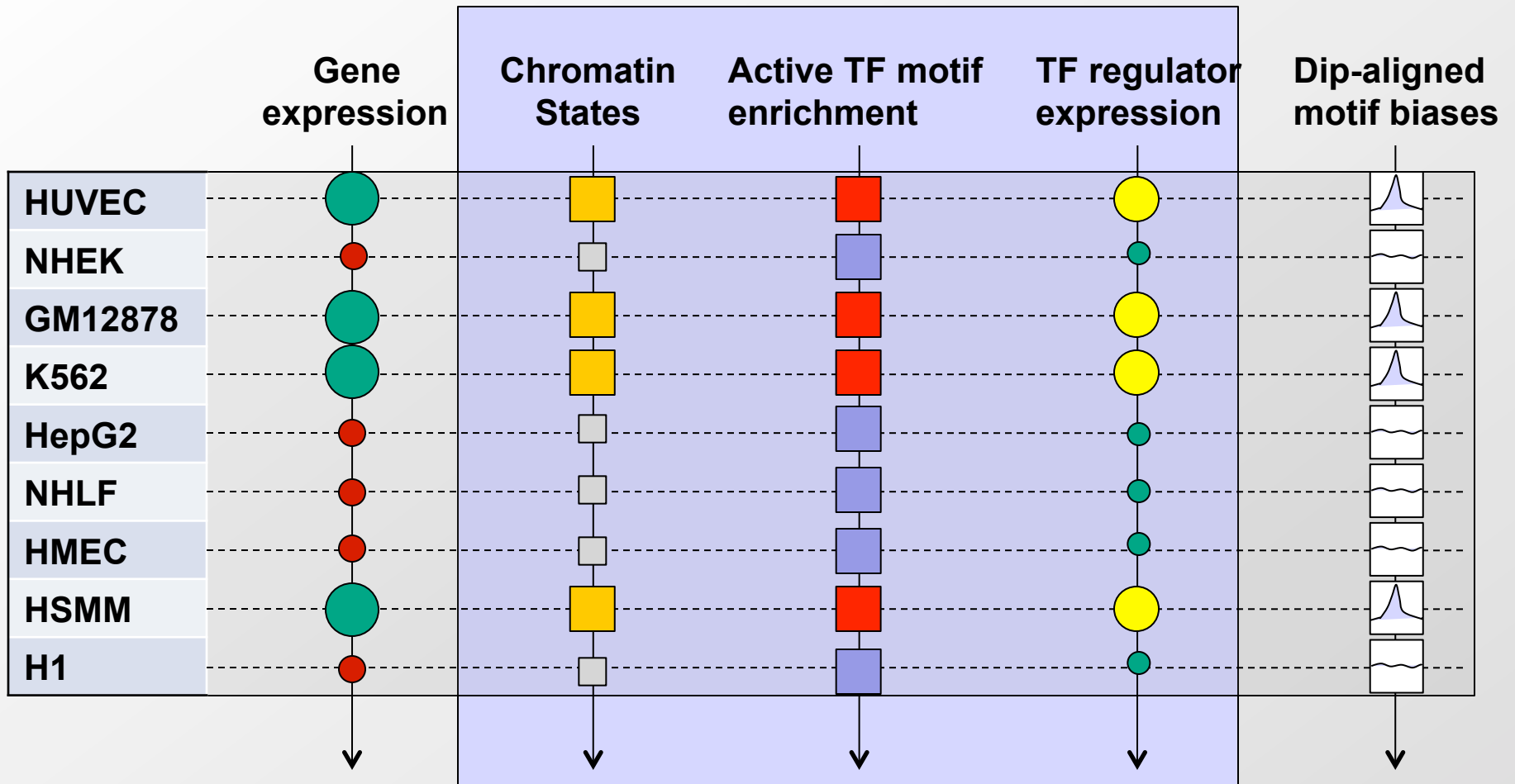
Multi-cell activity profiles connect enhancers



Link enhancers to target genes

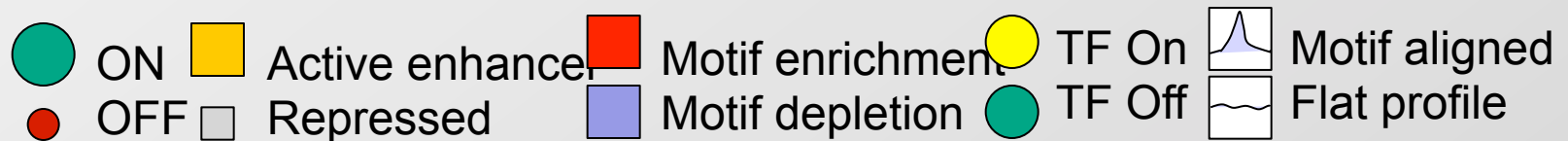


Multi-cell activity profiles connect enhancers



Link TFs to target enhancers

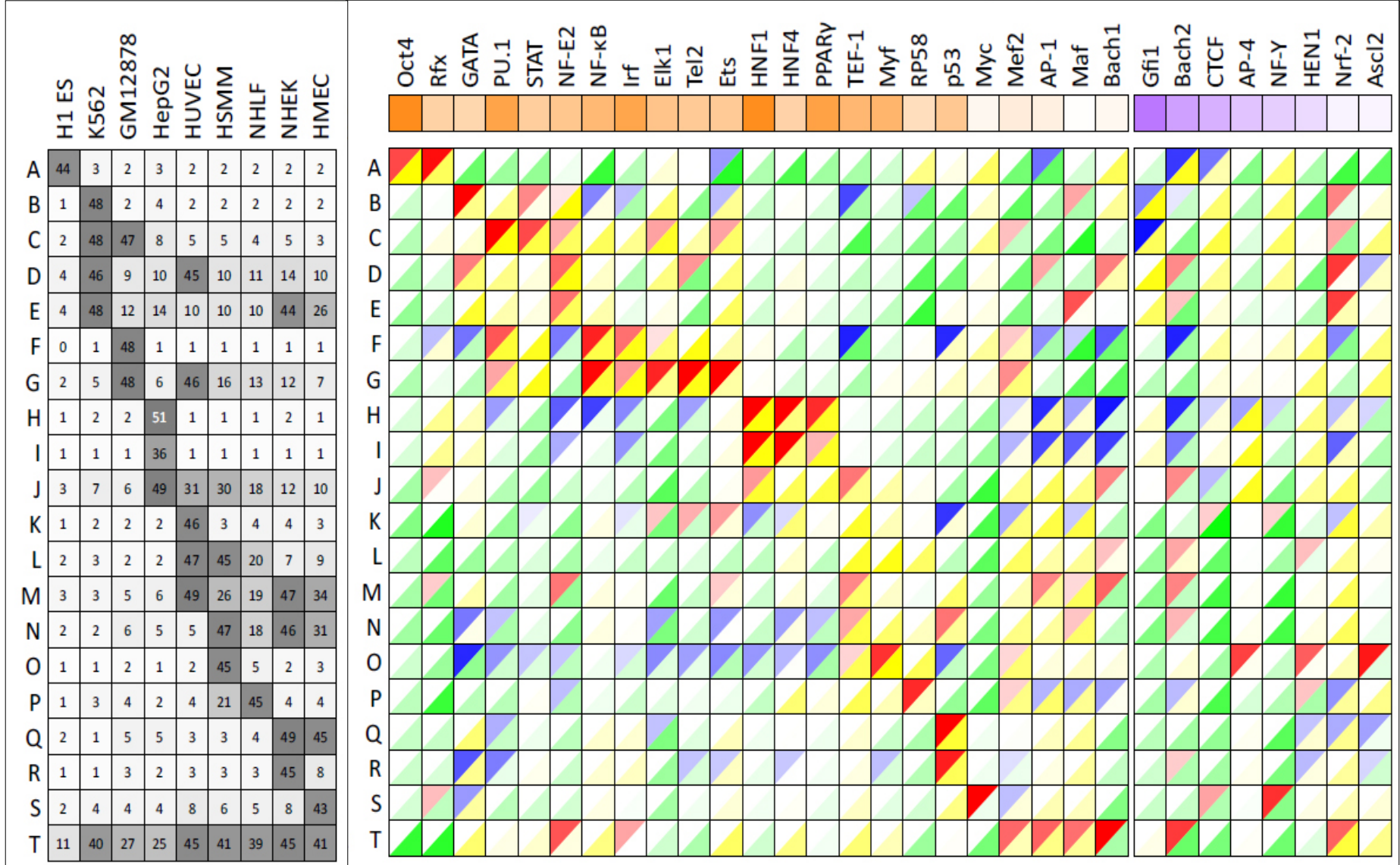
Predict activators vs. repressors



Coordinated activity reveals activators/repressors

Enhancer activity

Activity signatures for each TF

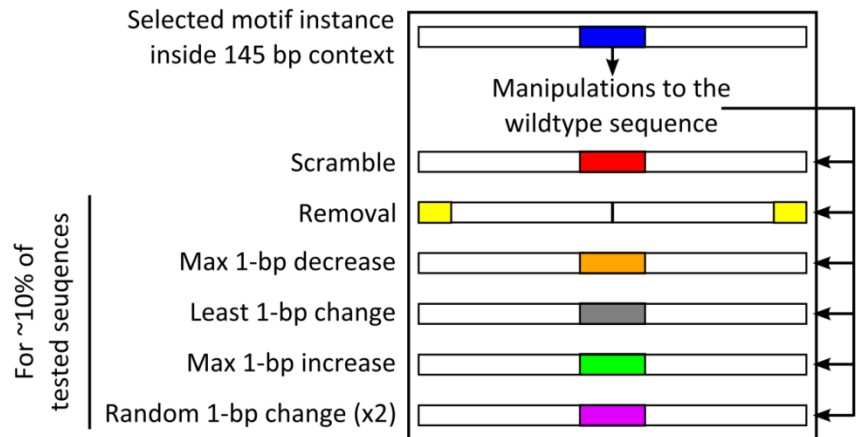


• Enhancer networks: Regulator → enhancer → target gene

Experimental validation of motif activity in tissue-specific enhancers

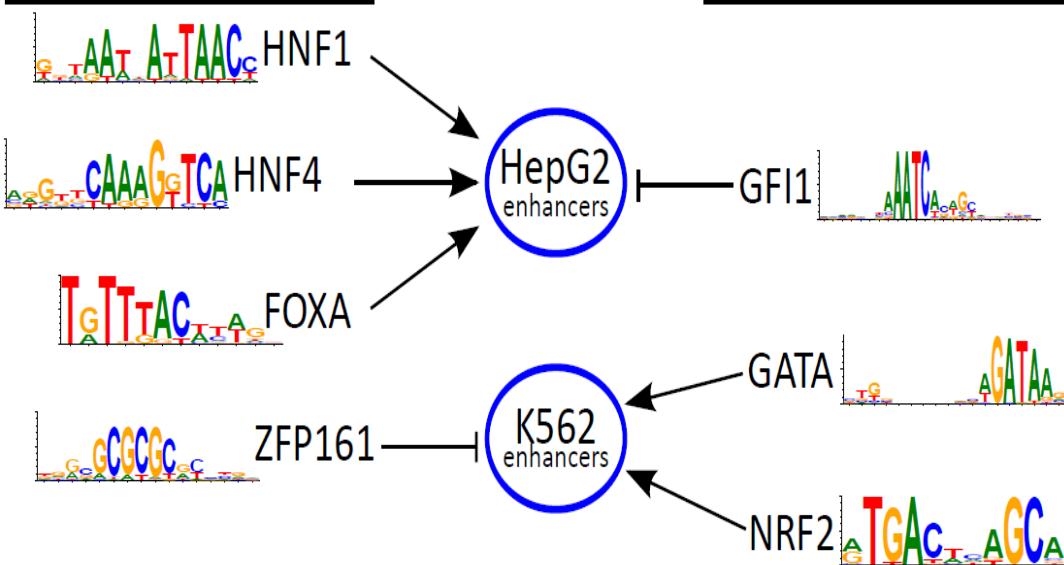
Motif enrichment in enhancers

Motif-motif similarity	Motif-motif similarity							Motif enrichment in enhancers				Factor expression			
	HNF1	HNF4	FOXA	GATA4	NRF2	ZFP161	GFI1	HepG2		K562		HepG2		K562	
	0.0	0.4	0.0	0.4	0.0	0.4	0.0	0.4	0.0	0.4	rep1	rep2	rep1	rep2	
HNF1	1.0	0.4	0.4	0.4	0.4	0.1	0.4	1.5	2.3	1.0	1.0	0.8	0.5	-0.1	-0.2
HNF4	0.4	1.0	0.4	0.3	0.3	0.2	0.3	1.7	2.1	1.0	1.0	1.0	0.5	-0.0	-0.1
FOXA	0.4	0.4	1.0	0.3	0.5	0.1	0.4	1.4	1.7	1.0	1.0	2.2	2.1	-0.4	-0.4
GATA	0.4	0.3	0.3	1.0	0.3	0.1	0.5	1.0	1.0	2.1	2.8	0.1	0.3	0.4	0.4
NRF2	0.4	0.3	0.5	0.3	1.0	0.2	0.4	1.0	1.1	1.5	1.8	0.3	0.7	-0.1	-0.3
ZFP161	0.1	0.2	0.1	0.1	0.2	1.0	0.1	0.8	0.5	1.2	1.0	0.0	0.0	0.1	0.1
GFI1	0.4	0.3	0.4	0.5	0.4	0.1	1.0	1.0	1.0	0.6	0.5	0.4	0.3	1.3	1.1

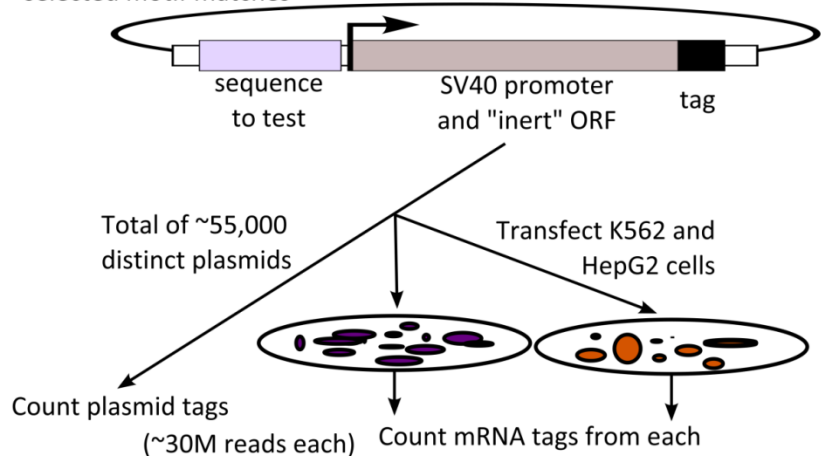


Active in HepG2 cells

Active in K562 cells

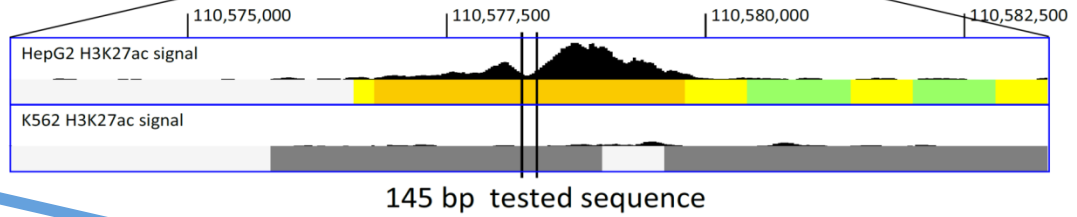
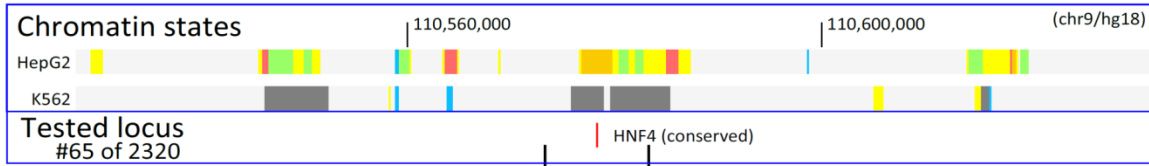


Sequences from other selected motif matches → Synthesize and construct plasmid pool



54000+ measurements (x2 cells, 2x repl)

Example activator: conserved HNF4 motif match

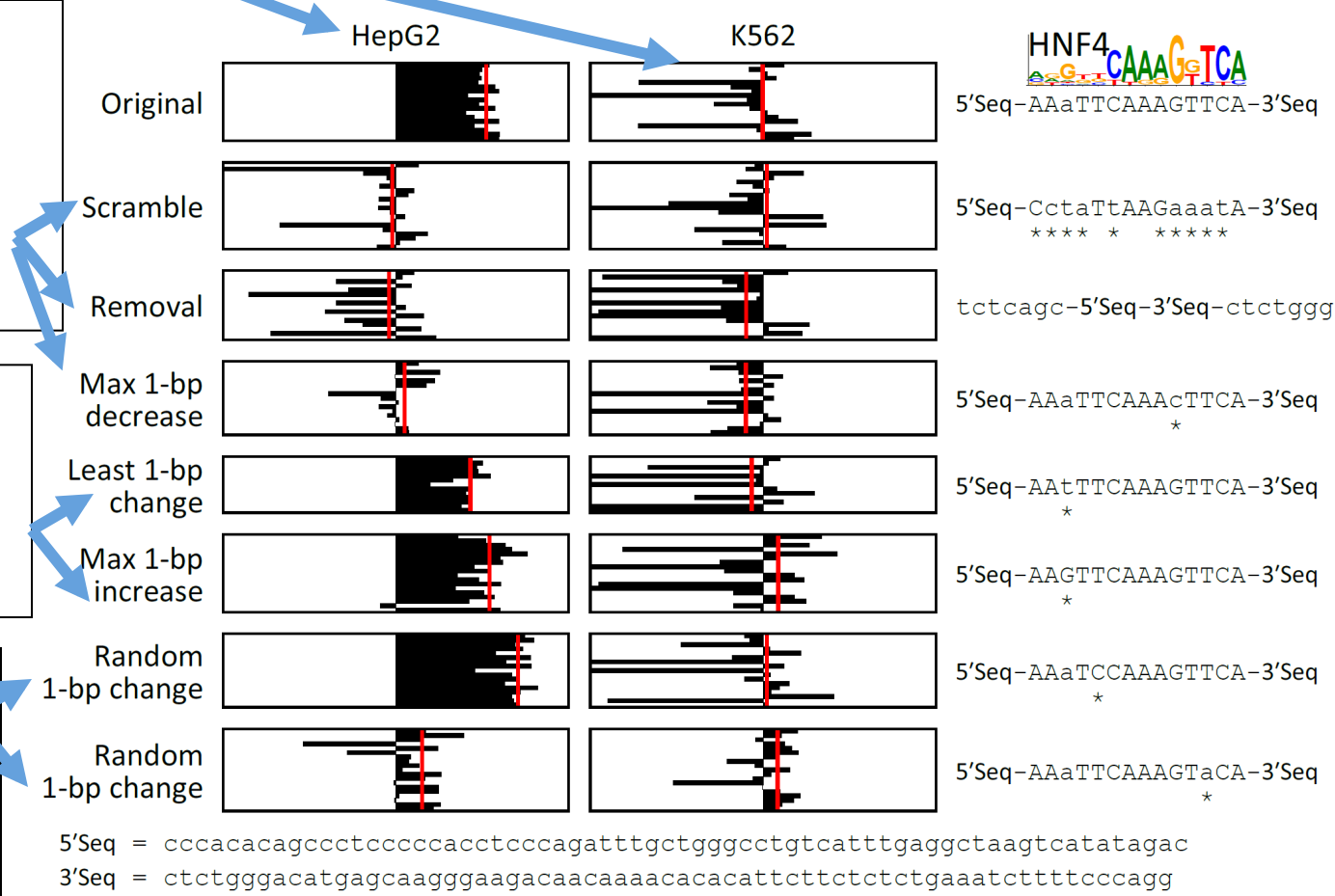


WT expression
specific to HepG2

Motif match
disruptions reduce
expression to
background

Non-disruptive
changes maintain
expression

Random changes
depend on effect
to motif match



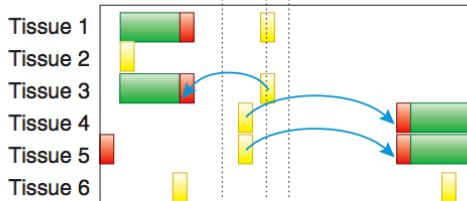
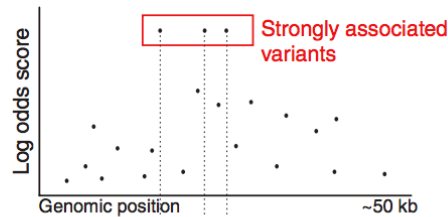
Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

Faced with resolution-limiting LD, annotations can help

Interpreting GWAS signals using functional and comparative genomics datasets

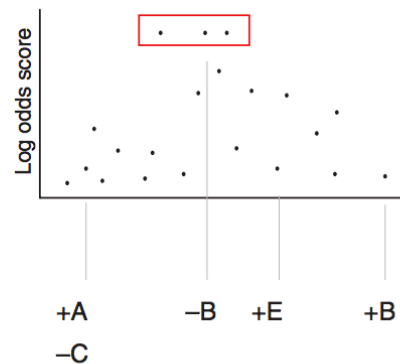
a Dissect associated haplotype using functional genomics



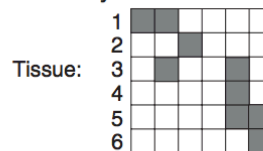
- Enhancer histone marks
- Promoter histone marks
- Transcribed region marks
- ↪ Enhancer-gene links

Chromatin state annotations

b Dissect associated haplotype using regulatory genomics

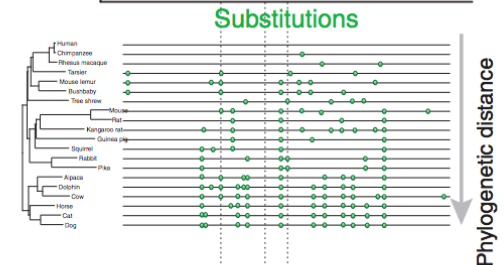
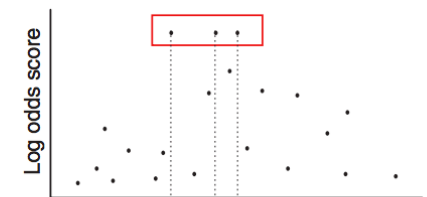


Enhancer motif enrichment analysis
Motif library: A B C D E F



Motifs altered by variants

c Dissect associated haplotype using comparative genomics



Constrained elements

Mammalian constraint

HaploReg: systematic regulatory mining of variants

(compbio.mit.edu/HaploReg)

Query SNP: **rs17145713** and variants with $r^2 \geq 0.95$

chr	pos (hg19)	LD	variant	Ref	Alt	ASN freq	CEU freq	YRI freq	GERP cons	SiPhy cons	Promoter histone marks	Enhancer histone marks	DNAse	Proteins bound	Motifs changed	GENCODE genes	RefSeq genes	dbSNP func annot
7	72842724	1	7:72480660	GAC	G	0	0.15	0					PANC-1			5.4kb 5' of FZD9	5.4kb 5' of FZD9	
7	72856430	1	rs1178979	T	C	0.13	0.18	0.3					CLL		GATA	BAZ1B	BAZ1B	intronic
7	72857049	1	rs1178977	A	G	0.14	0.18	0.3							AREB6,DEC	BAZ1B	BAZ1B	intronic
7	72857713	1	rs34604283	CA	C	0.13	0.1	0.2					8 cell types		Sox	BAZ1B	BAZ1B	intronic
7	72868522	1	rs1306476	A	G	0.12	0.18	0.36								BAZ1B	BAZ1B	intronic
7	72883106	1	rs62465144	T	C	0.14	0.18	0.29								BAZ1B	BAZ1B	intronic
7	72885810	1	rs6976930	G	A	0.14	0.18	0.39								BAZ1B	BAZ1B	intronic
7	72904810	1	rs17145713	C	T	0.14	0.18	0.3							ATF3	BAZ1B	BAZ1B	intronic
7	72939244	1	rs11983997	G	C	0.13	0.18	0.26				GM12878, K562	GM12864, GM12878, K562			2.6kb 5' of BAZ1B	2.6kb 3' of BAZ1B	
7	72977249	1	rs34594435	C	T	0.12	0.18	0.03				K562	CMK	KAP1, SETDB1		4.9kb 5' of BCL7B	5.2kb 3' of BCL7B	
7	72988069	1	rs35659126	C	T	0.13	0.18	0.08								TBL2	TBL2	intronic
7	72989141	1	rs34550818	C	CA	0.11	0.12	0						POL2		TBL2	TBL2	intronic
7	72989390	1	rs11974409	A	G	0.13	0.18	0.14								TBL2	TBL2	intronic
7	72998952	1	rs9638180	A	G	0.12	0.18	0.08							Zbtb3	5.8kb 5' of TBL2	5.9kb 3' of TBL2	
7	72999105	1	rs9638182	T	G	0.12	0.18	0.14								6kb 5' of TBL2	6.1kb 3' of TBL2	
7	73007943	1	rs1051921	G	A	0.12	0.18	0.08					4 cell types	POL2		MLXIPL	MLXIPL	3'-UTR

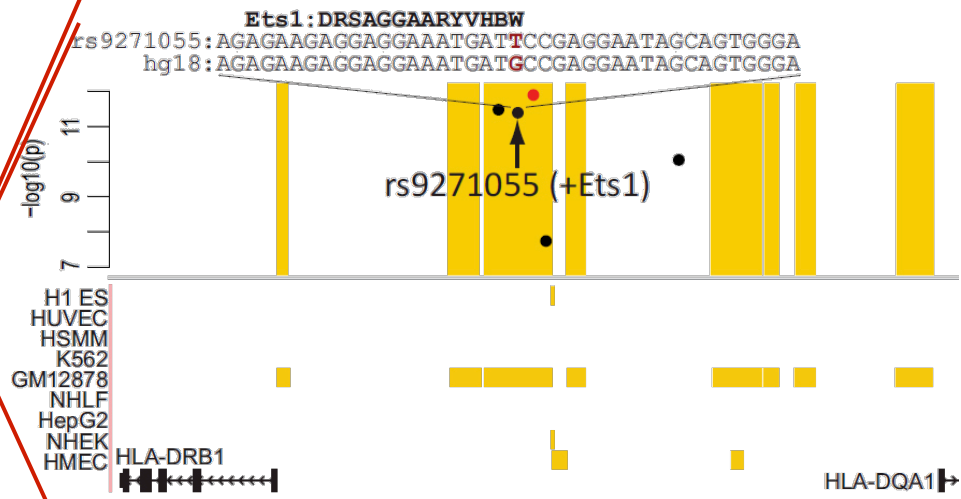
- Start with any list of SNPs or select a GWA study
 - Mine publically available ENCODE+Roadmap data for hit LD blocks
 - Hundreds of assays, dozens of cells, conservation, motifs
 - Integration with published eQTL from GTEx portal (future: ASE)
 - Report systems-level enrichments (future: R tools)

Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

Phenotype	Top Cell Type	Total #SNPs from Study	#SNPs in enh. States 4 and 5	p-value	FDR	H1 ES	K562	GM12878	HepG2	HUVEC	HSMM	NHLF	NHEK	HMEC
Erythrocyte phenotypes (Ref. 38)	K562	35	9	$<10^{-7}$	0.02	9	17	4	0	0	1	2	1	1
Blood lipids (Ref. 39)	HepG2	101	13	$<10^{-7}$	0.02	3	5	0	11	2	3	3	4	3
Rheumatoid arthritis (Ref. 40)	GM12878	29	7	2.0×10^{-7}	0.03	0	0	15	0	2	0	0	2	3
Primary biliary cirrhosis (Ref. 41)	GM12878	6	4	6.0×10^{-7}	0.03	0	11	41	0	0	0	0	8	8
Systemic lupus erythromatosus (Ref. 42)	GM12878	18	6	9.0×10^{-7}	0.03	0	4	21	0	5	8	0	3	5
Lipoprotein cholesterol/triglycerides (Ref. 43)	HepG2	18	5	1.2×10^{-6}	0.03	17	8	0	24	3	6	4	3	3
Hematological traits (Ref. 44)	K562	39	7	1.7×10^{-6}	0.03	0	12	10	2	1	0	0	1	0
Hematological parameters (Ref. 45)	K562	28	6	2.2×10^{-6}	0.03	0	15	7	0	5	7	7	3	2
Colorectal cancer (Ref. 46)	HepG2	4	3	3.8×10^{-6}	0.03	0	0	0	66	0	12	0	12	12
Blood pressure (Ref. 47)	K562	9	4	5.0×10^{-6}	0.04	0	30	14	0	10	6	7	5	11

SNP	H1 ES	K562	GM	HepG2	Huvec	HSMM	NHLF	NHEK	HMEC	Chrom. Band	Gene	Link Sc	Distanc
rs13385731	■	■	■	■	■	■	■	■	■	2p22			
rs10036748	■	■	■	■	■	■	■	■	■	5q33			
rs1385374	■	■	■	■	■	■	■	■	■	12q24	MGC16384	-	1
rs2230926	■	■	■	■	■	■	■	■	■	6q23	TNFAIP3	3.7	7
rs4728142	■	■	■	■	■	■	■	■	■	7q32	IRF5	-	4
rs9271100	■	■	■	■	■	■	■	■	■	6p21	HLA-DRB1	4.5	19
rs4917014	■	■	■	■	■	■	■	■	■	7p12	IKZF1	2.2	38
rs7812879	■	■	■	■	■	■	■	■	■	8p23	BLK	2.9	11
rs2205960	■	■	■	■	■	■	■	■	■	1q25			
rs548234	■	■	■	■	■	■	■	■	■	6q21			



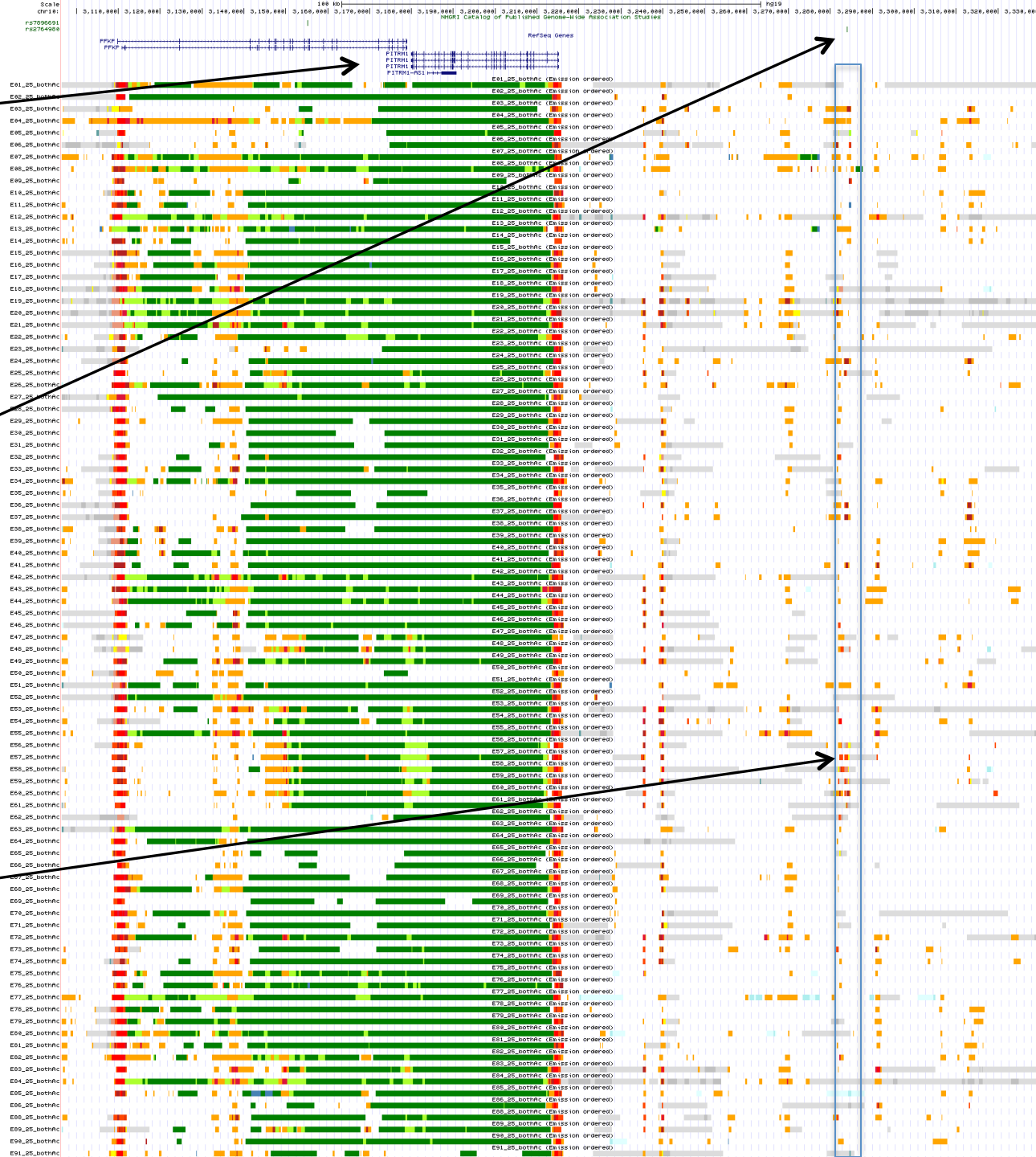
- Disease-associated SNPs enriched for enhancers in relevant cell types
- E.g. lupus SNP in GM enhancer disrupts Ets1 predicted activator

Heatmap showing gene expression levels (log2) across various cell types and tissues. The color scale ranges from 0 (blue) to 14 (red). Genes are grouped into clusters such as 'ADHD / brain: cingulate gyrus', 'Adiponectin and CAD / adipocytes', and 'CD4+ CD25+ IL17+ PMA-ionomycin-stimulated_Th17_Primary_Cells'. Tissues include Brain, Bone Marrow, Chondrocytes, and Stomach.

PITRM1
gene

ADHD SNP
rs2764980

Intergenic
brain-specific
promoter



HaploReg view of ADHD SNP rs17658378

Query SNP: **rs17658378** and variants with $r^2 \geq 0.8$

chr	pos (hg19)	LD (r ²)	LD (D')	variant	Ref	Alt	AFR freq	AMR freq	ASN freq	EUR freq	SiPhy cons	Promoter histone marks	Enhancer histone marks	DNAse	Proteins bound	eQTL tissues	Motifs changed	GENCODE genes	dbSNP func annot
8	116394075	1	1	rs17658378	A	G	0.02	0.13	0.00	0.09		BN SN, BN ITL, BN AC	9 cell types				4 altered motifs	27kb 3' of TRPS1	

Regulatory chromatin states (Roadmap)

Cell ID	Cell description	State (25-state HMM)
BN.SN	Brain Substantia Nigra	6_TssD2
BN.ITL	Brain Anterior Temporal Lobe	6_TssD2
BN.AC	Brain Anterior Caudate	6_TssD2
BR.H35	Breast vHMEC.Donor RM035	11_EnhWk1
IPS.18	IPS-18 Cell Line	12_EnhWk2
H1.BMP4DM	H1 BMP4 Derived Mesendoderm Cultured Cells	12_EnhWk2
BN.CC	Brain Cingulate Gyrus	12_EnhWk2
PFK.2	Penis Foreskin Keratinocyte Primary Cells.Donor skin02	12_EnhWk2
HUES6	HUES6 Cell Line	12_EnhWk2
BN.HM150	Brain Hippocampus Middle.Donor 150	12_EnhWk2
ESO	Esophagus	12_EnhWk2
BN.AG	Brain Angular Gyrus	12_EnhWk2

Regulatory motifs altered

PWM	Strand	Ref	Alt	Match on:
				Ref: CTGTTCCTCTTCCAGGCCATAGCGGCTATCAGGAACITGTAGCCATCTGGGGGTCAG Alt: CTGTTCCTCTTCCAGGCCATAGCGGCTGTCAGGAACITGTAGCCATCTGGGGGTCAG
PU_1_disc1	-	-33.3	-21.6	AWGRGGAAGT
PU_1_known3	-	11.7	9.8	NHASTTCBYHWHN
Pbx3_disc3	-	6.6	11.2	TGGGYVNNBNSCYGYCMVT
SETDB1_disc1	+	2.3	9.9	CRNDGMHYBMYGGRARWKGTAGTYY
Znf143_disc3	+	9.9	11.2	GSVBBSBGGGVVNBGBRGB

rs17658378 / TRPS1
 From Lasky-Su et al (PMID 2008), who do not consider it one of the interesting loci
 TRPS1 is a TF; locus has been associated by GWAS with weight fluctuation (PMID 22911880) and major depressive disorder (PMID 22472876)



mesh_term Intracranial Aneurysm
 mesh_term Immune System Diseases
 mesh_term Head Diseases
 mesh_term Digestive System Diseases
 mesh_term Crohn Disease
 mesh_term Multiple Sclerosis
 mesh_term Diabetes Mellitus, Type 1
 analysis_name New gene functions in megakaryopoiesis and platelet formation.
 mesh_term Platelet Count
 mesh_term Celiac Disease
 analysis_name Biological, clinical and population relevance of 95 loci for blood lipids.
 analysis_name Primary role for cell-mediated immune mechanisms in multiple sclerosis.
 analysis_name Multiple common variants for celiac disease influencing immune gene expression.
 analysis_name Genome-wide meta-analysis increases to 71 the number of confirmed Crohn's disease
 mesh_term Crohn Disease
 analysis_name Genome-wide association defines more than 30 distinct susceptibility loci for Crohn's disease
 mesh_term Liver Cirrhosis, Biliary
 mesh_term Diabetes Mellitus
 mesh_term Lupus Erythematosus, Systemic
 analysis_name Genome-Wide Association Study in Type 1 Diabetes
 mesh_term Cellitis, Ulcerative
 mesh_term Behcet Syndrome
 analysis_name Genome-wide Association Study of Behçet's disease
 analysis_name Genome-wide association study and meta-analysis find that over 40 loci affect risk of type
 mesh_term Lipoproteins, HDL
 mesh_term Body Height
 analysis_name Genome-wide association analysis of autoantibody positivity in type 1 diabetes cases.
 analysis_name Hundreds of variants clustered in genomic loci and biological pathways affect human height.
 analysis_name Meta-analysis of genome scans and replication identify CD6, IRF8 and TNFRSF1A as new
 analysis_name Meta-analysis identifies 29 additional ulcerative colitis risk loci, increasing the number of con
 analysis_name Many sequence variants affecting diversity of adult human height.
 analysis_name Genome-wide association study identifies 12 new susceptibility loci for primary biliary cirrho
 analysis_name Genome-wide association analysis identifies 20 loci that influence adult height.
 analysis_name Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk
 analysis_name Meta-analysis of genome-wide association studies in 8,678,000 subjects identifies multiple
 mesh_term Laboratory Techniques and Procedures
 mesh_term Autoimmune Diseases
 analysis_name Meta-analysis identifies nine new susceptibility
 analysis_name Meta-analysis identifies 13 new loci associated with waist-hip ratio and reveals sexual dimor
 analysis_name Genome-wide Association scan in women with systemic lupus erythematosus
 analysis_name Genome-wide association study identifies loci influencing concentration of liver enzymes in
 mesh_term Arthritis, Rheumatoid
 mesh_term Leprosy
 mesh_term Glutamate Transaminase
 analysis_name Genome-wide association study of ulcerative colitis identifies three new susceptibility loci, i
 mesh_term Coronary Artery Disease
 mesh_term Autoimmune Diseases
 analysis_name Meta-analysis of genome-wide association studies in celiac disease and rheumatoid arthritis
 analysis_name Meta-analysis of genome-wide scans for human adult stature identifies novel Loci and asso
 analysis_name Genome-wide association study of leprosy in Chinese population
 analysis_name Newly identified loci that influence lipid concentrations and risk of coronary artery disease.
 mesh_term Colorectal Neoplasms
 analysis_name Common variants at 30 loci contribute to polygenic dyslipidemia.
 analysis_name Interaction between ERAP1 and HLA-B27 in ankylosing spondylitis implicates peptide handling
 analysis_name Common variants in 22 loci are associated with QRS duration and cardiac ventricular condu
 analysis_name A genome-wide meta-analysis identifies 22 loci associated with eight hemological paramet
 analysis_name Genetic variants in novel pathways influence blood pressure and cardiovascular disease risk
 mesh_term Waist-Hip Ratio
 mesh_term Cholesterol
 mesh_term Apolipoproteins E
 mesh_term Erythrocyte Indices
 analysis_name Six new loci associated with blood low-density lipoprotein cholesterol, high-density lipoprotein
 mesh_term Tumor Necrosis Factor-alpha
 analysis_name Large-scale genome-wide association study of Asian populations uncovers genetic factors
 analysis_name Large-scale association analysis identifies 13 new susceptibility loci for coronary artery dise
 analysis_name Genome-wide association study of hematological and biochemical traits in a Japanese popu
 mesh_term Macular Degeneration
 analysis_name Genome-wide association study identifies six new loci influencing pulse pressure and mean
 analysis_name A genome-wide association study in Europeans and South Asians identifies five new loci for
 analysis_name Loci influencing lipid levels and coronary heart disease risk in 16 European population coho
 analysis_name Genome-wide association yields new sequence variants at seven loci that associate with m
 analysis_name Population-based genome-wide association studies reveal six loci influencing plasma levels
 mesh_term Cornea
 mesh_term Spondylitis, Ankylosing
 mesh_term Kidney Diseases
 mesh_term Cholesterol, HDL
 mesh_term Stomatognathic Diseases
 mesh_term Bacterial Infections and Mycoses
 mesh_term Heart Function Tests
 analysis_name Several common variants modulate heart rate, P-R interval and QRS duration.
 analysis_name Whole Genome Association Study of Systemic Lupus Erythematosus
 mesh_term Optic Disk
 mesh_term Life Expectancy
 analysis_name A genome-wide association study identifies two new risk loci for Graves' disease.
 analysis_name A genome-wide association study of optic disc parameters.
 mesh_term Coronary Disease
 mesh_term Scleroderma, Systemic
 analysis_name Identification, replication, and fine-mapping of Loci associated with adult height in individua
 analysis_name A meta-analysis and genome-wide association study of platelet count and mean platelet volu

Enhancers (by tissue)

DNase (by tissue)

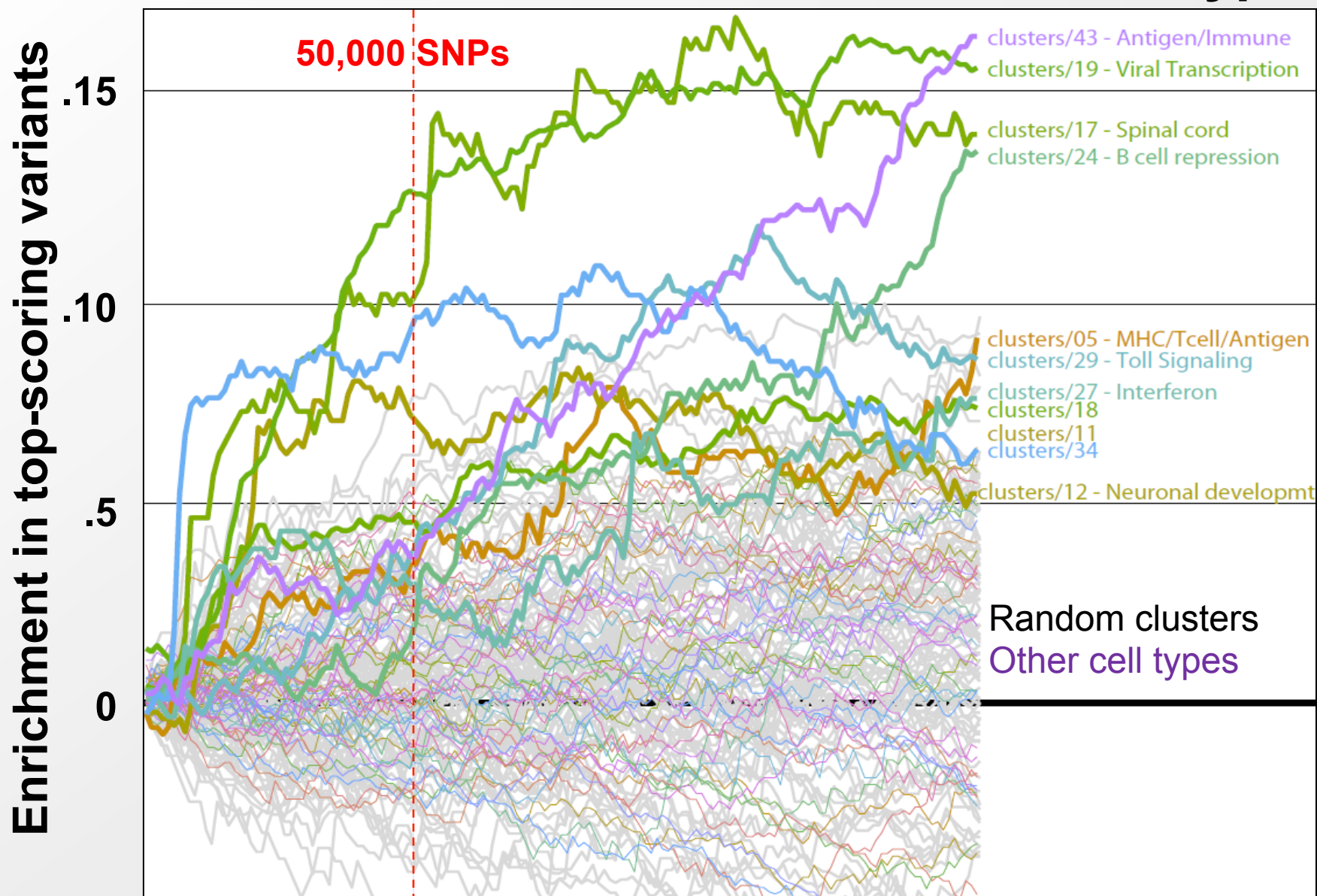
Enhancers (by cluster)

DNase (by cluster)

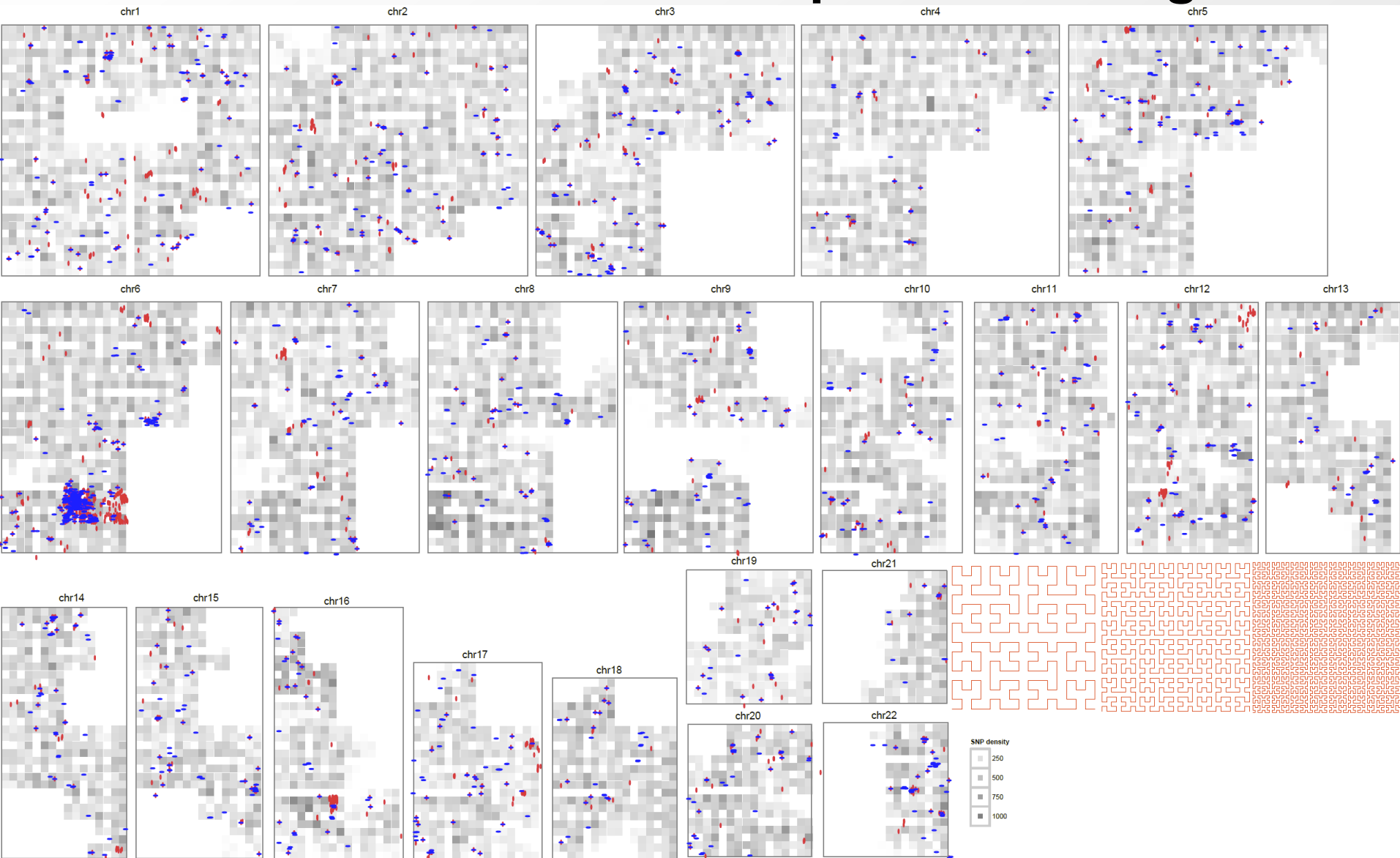
Examples where the clustered enhancers give a previously undetectable enrichment

(Cardiac traits <-> enhancer cluster with fluid shear GO)

T1D enrichment relative to randomized cell types



T1D/RA-enriched enhancers spread across genome

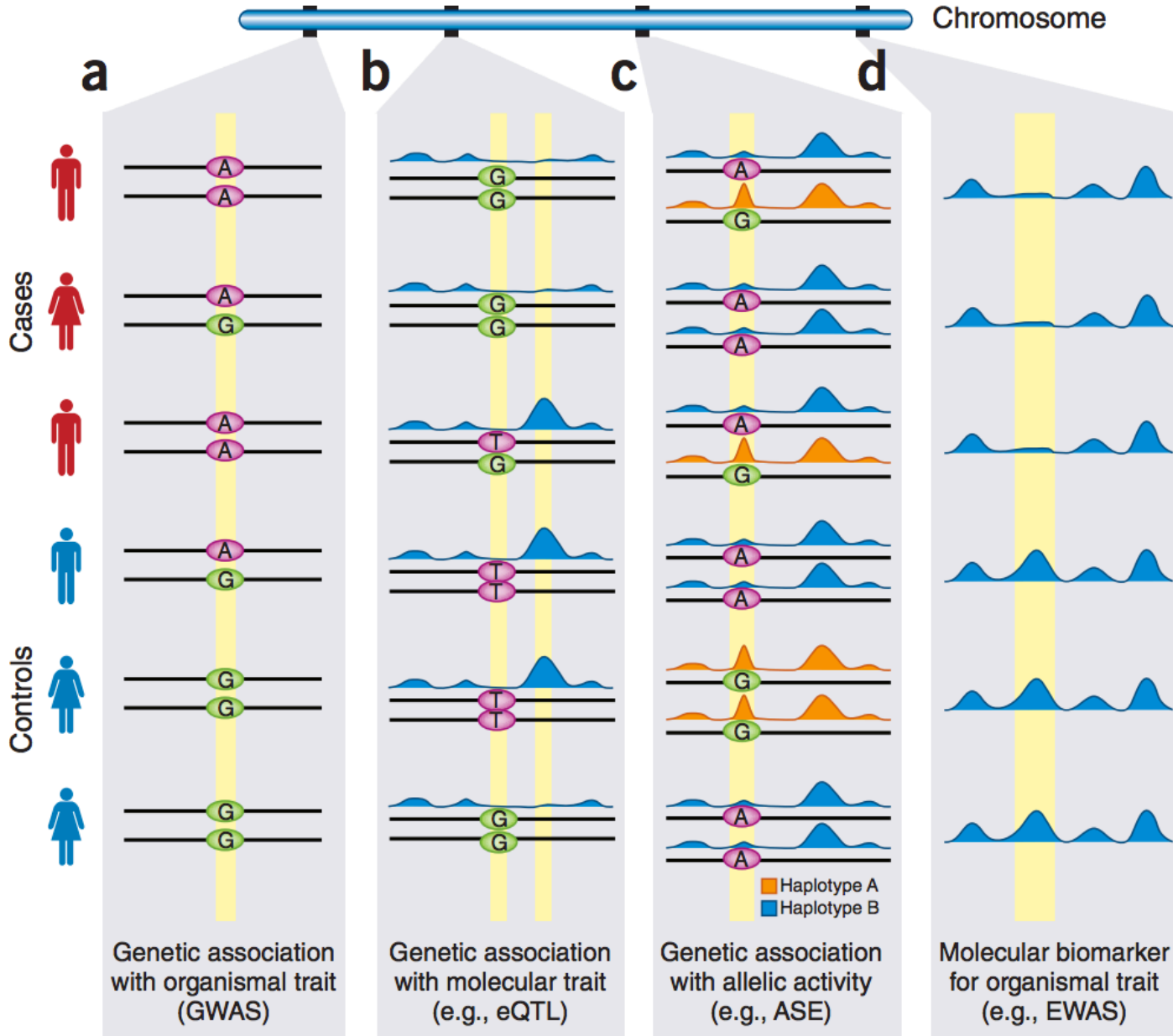


- High concentration of loci in MHC, high overlap
- Yet: many distinct regions, 1000s of distinct loci

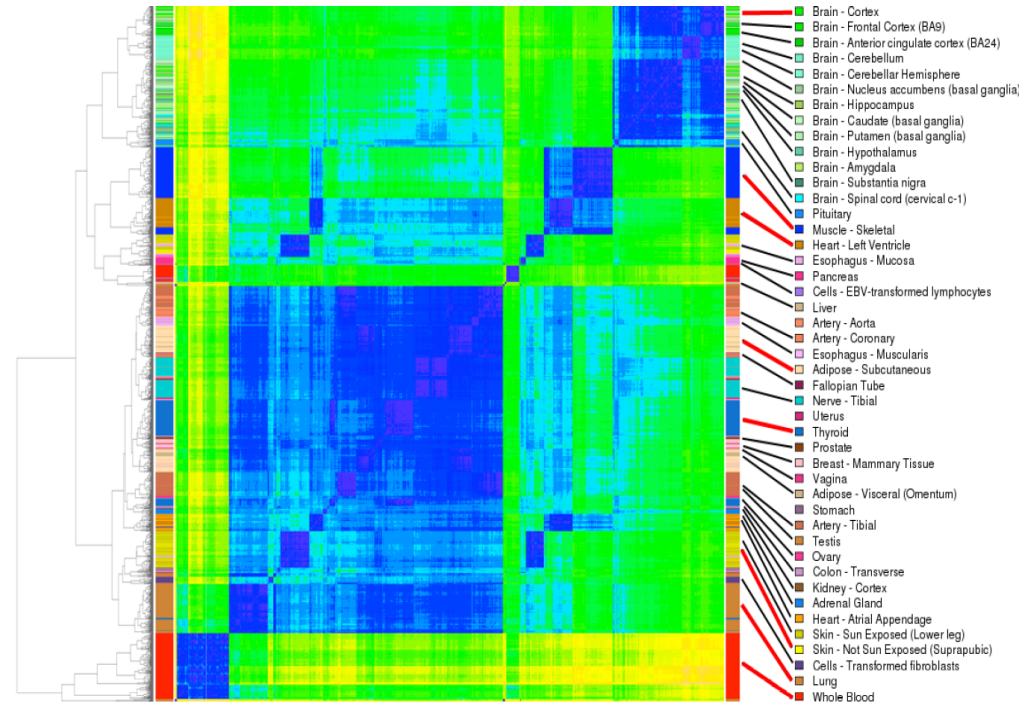
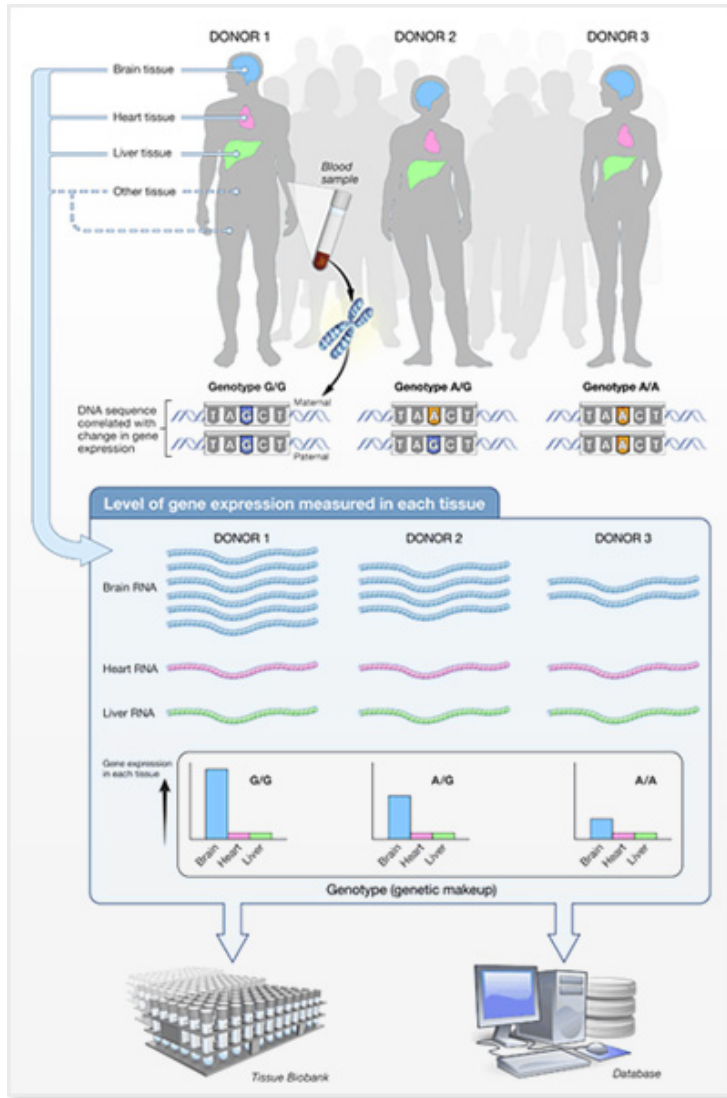
Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

Beyond GWAS: Molecular variability



eQTLs: The GTEx Project



17 GTEx tissues w/ close Roadmap match

GTEx tissues

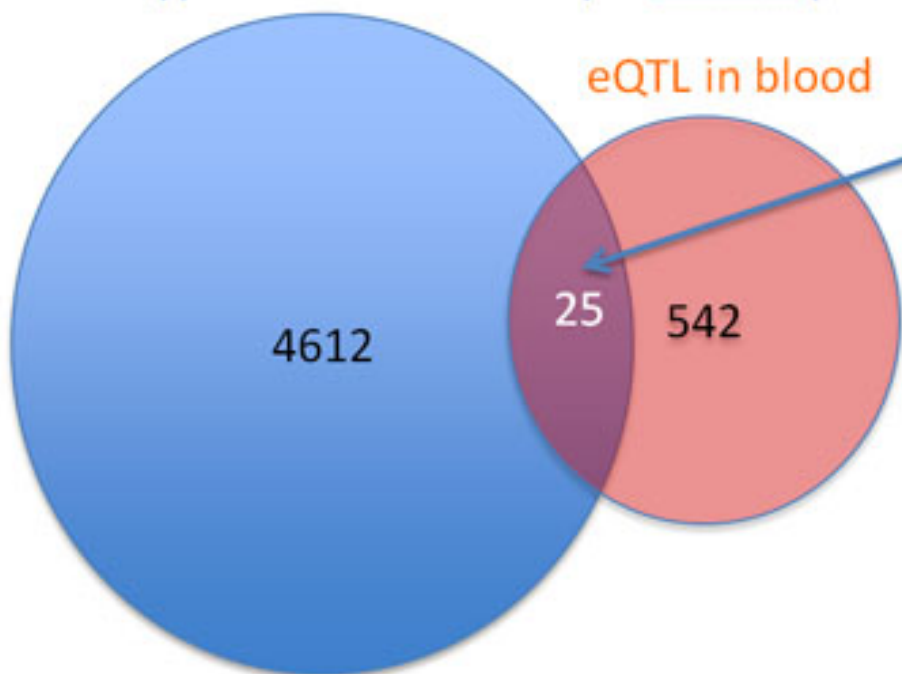
Current Roadmap tissues

1	Adipose	Adipose nuclei (E54)
2	Blood	Peripheral blood mononuclear primary cells (E87)
3	Brain - Midbrain - Substantia nigra	Brain - midbrain - substantia nigra (E61)
4	Brain - Cortex - Frontal Cortex (BA9)	Brain - cerebral cortex - mid frontal lobe (E60)
5	Brain - Cerebellar Hemisphere	Brain - cerebral cortex - inferior temporal lobe (E59)
	Brain - Cerebellum	Brain - cerebral cortex - angular gyrus (E56)
	Brain - Cortex - Anterior cingulate cortex (BA24)	Brain - cerebral cortex - cingulate gyrus (E58)
	Brain - Cortex	
6	Brain - Cerebrum - Subcortical - Hippocampus	Brain - cerebrum - basal ganglia - anterior caudate (E57)
	Brain - Cerebrum - Subcortical - Caudate (basal ganglia)	Brain - cerebrum - hippocampus (E25)
	Brain - Cerebrum - Subcortical - Basal ganglia - Putamen	
	Brain - Cerebrum - Subcortical - Basal ganglia - Nucleus accumbens	
7	Heart	Heart - left ventricle (E82)
		Heart - fetal heart (E04)
8	Lung	Lung - fetal lung (E06)
9	Muscle - Skeletal	Muscle - skeletal muscle (E47, E48, E89)
10	Pancreas	Pancreas (E85)
		Pancreatic islets (E86)
11	Skin	Penis foreskin fibroblast (E19, E20)
		Penis foreskin keratinocyte (E21, E42)
		Penis foreskin melanocyte (E22, E33, E44)

- Expected to increase with additional coverage
- Expression correlation metric for unbiased matching

Example of functional overlap

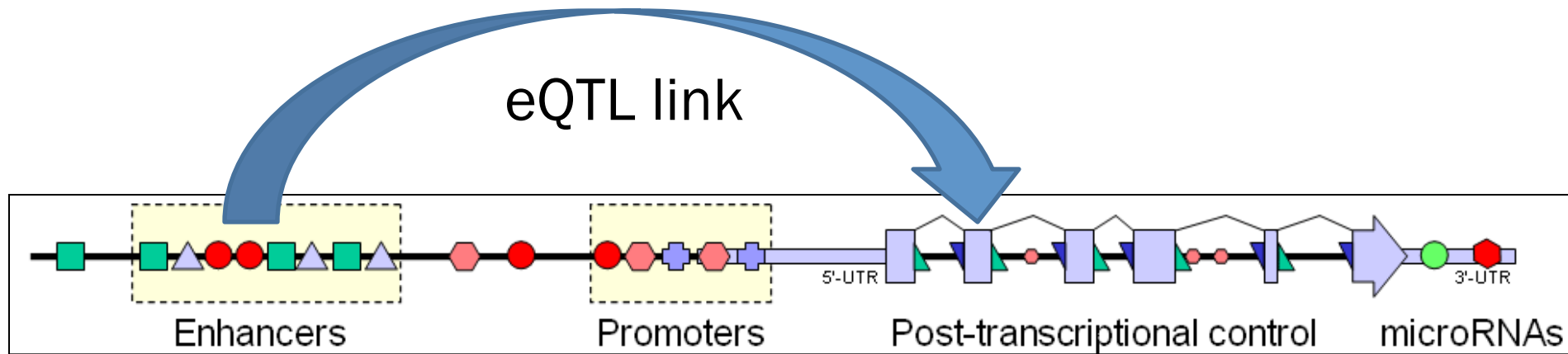
Genotyped SNPs in monocyte DNase peaks



Background model	Expected	P-value
All genes	5	4.7×10^{-12}
Genes expressed in tissue	7.3	4.4×10^{-8}
Genes expressed & eQTL TSS distance	9.4	5.2×10^{-6}

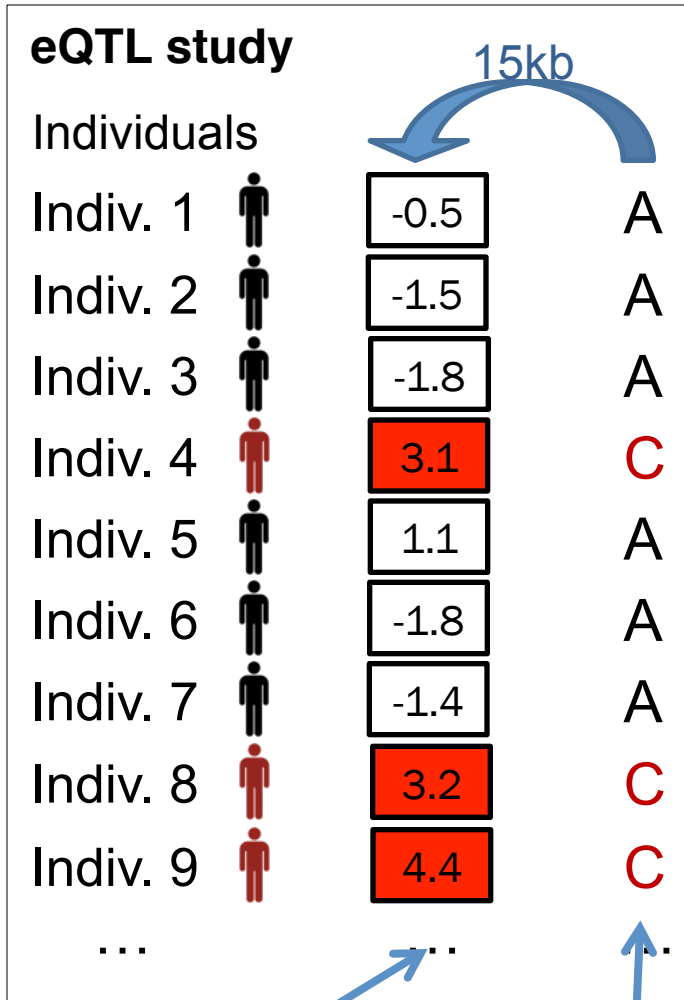
- 25 eQTLs in blood from GTEx lie within monocyte DNase peaks from ENCODE
- Significantly more than expected by chance
- Even after correcting for TSS distance and tissue³¹
→ eQTLs identify likely regulatory elements

Functional characterization of eQTLs



1. **Functional roles:** specific region, motif, link, ASE
2. Exploit GTEx matrix for **systems biology** studies
3. **Disease roles:** modules, tissues, genome-wide

Enhancer-gene links supported by eQTL-gene links



Validation rationale:

- Expression Quantitative Trait Loci (eQTLs) provide **independent SNP-to-gene** links
- Do they agree with activity-based links?

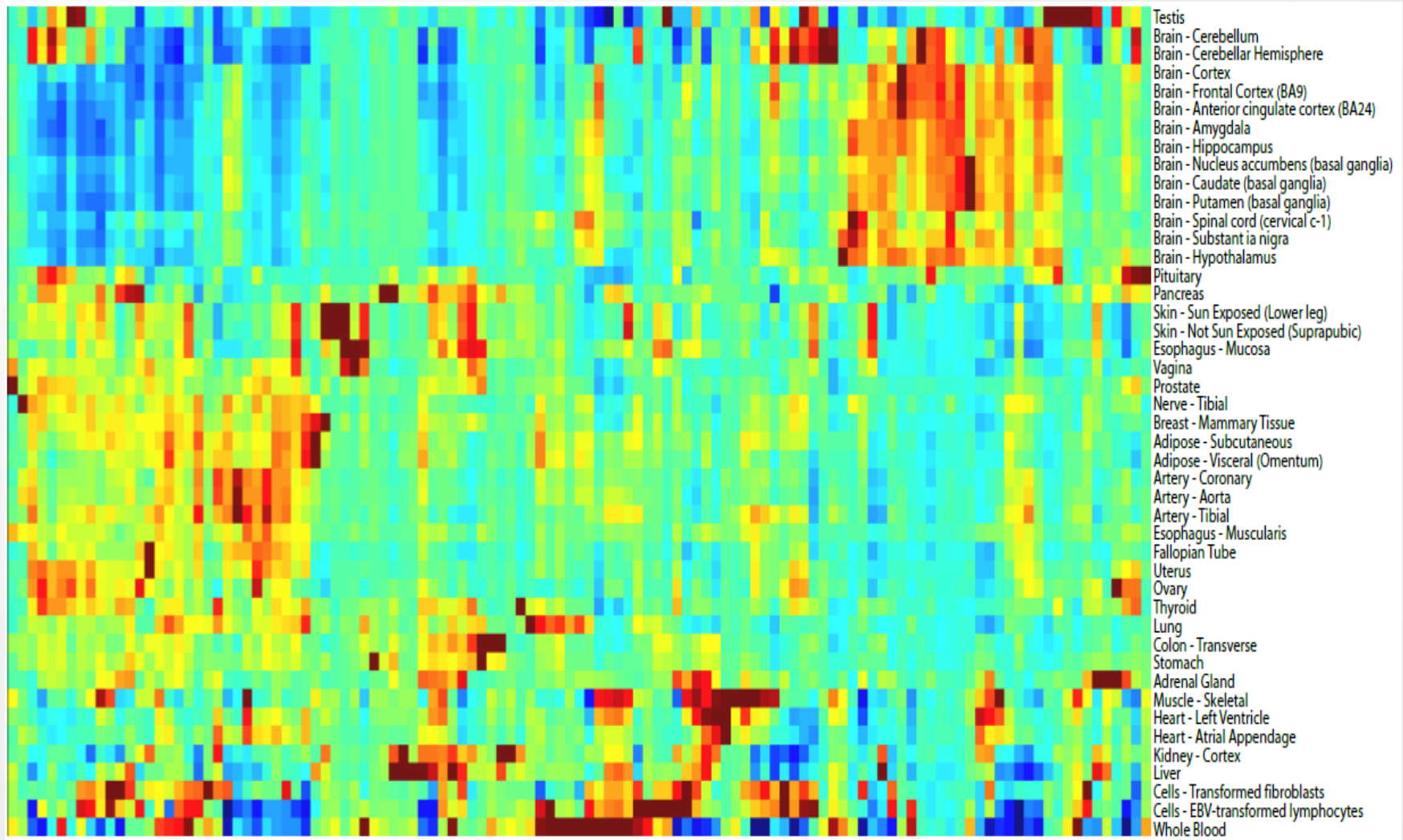
Example: Lymphoblastoid (GM) cells study

- **Expression/genotype** across 60 individuals (Montgomery *et al*, Nature 2010)
- **120** eQTLs are eligible for enhancer-gene linking based on our datasets
- **51** actually linked (43%) using predictions
→ **4-fold enrichment** (10% exp. by chance)

Expression level of gene ↔ **Sequence variant at distal position**

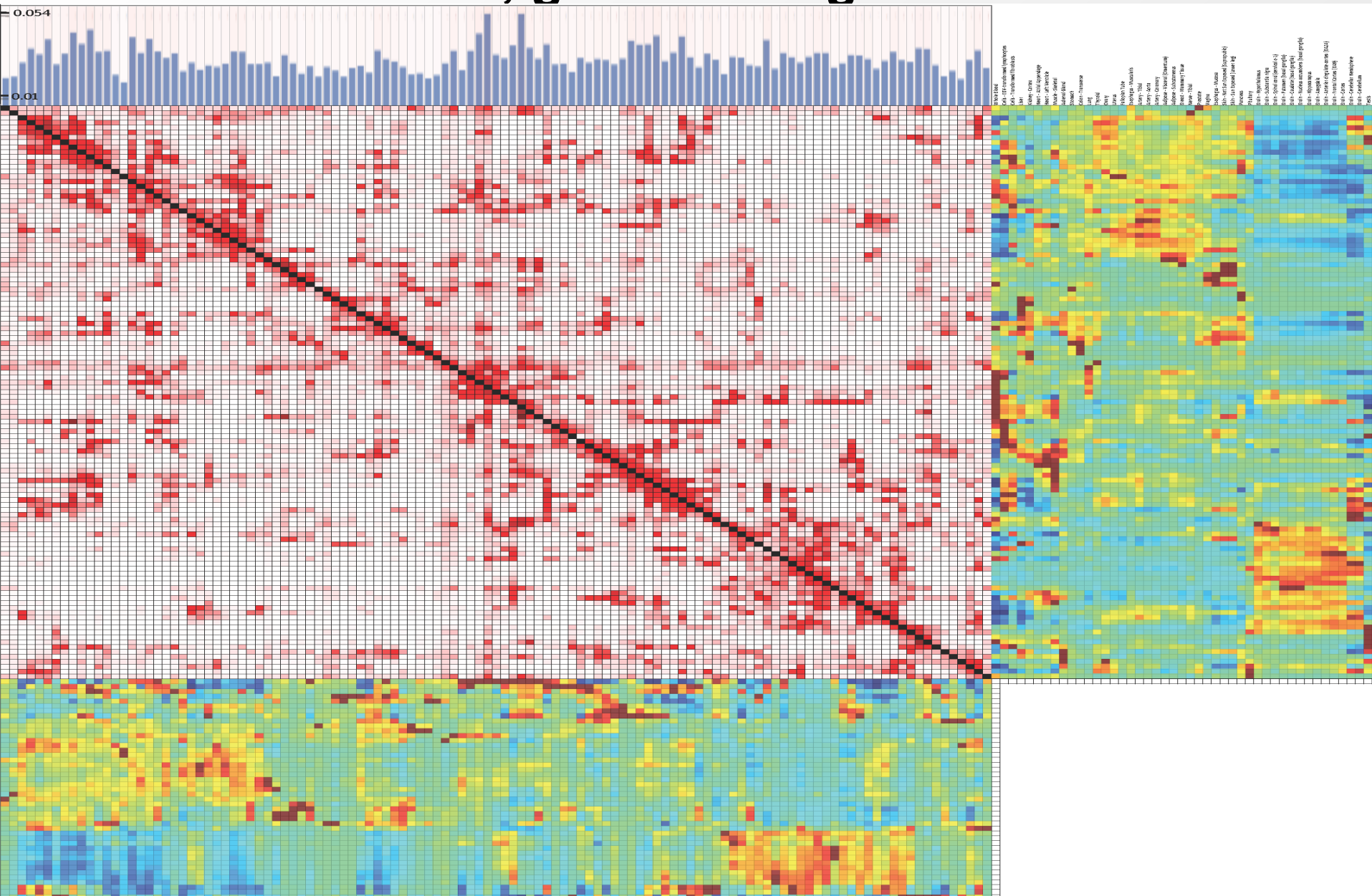
- Independent validation of links.
- Relevance to disease datasets.

GTEX eQTLs enable us to define multi-tissue expression modules in individuals



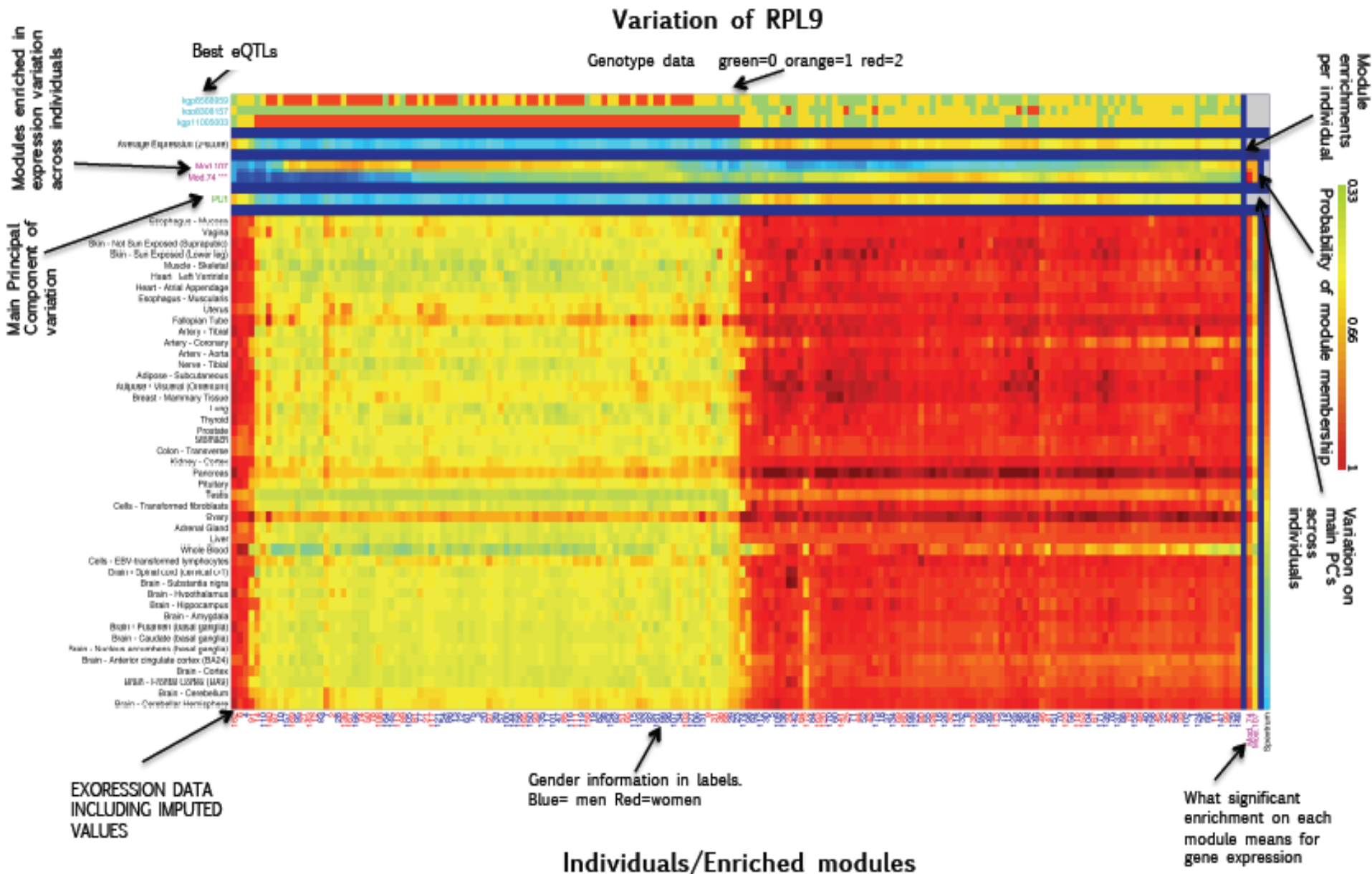
- Regulators and motifs affect gene expression patterns
- Exploit multi-tissue multi-individual nature of dataset

Across individuals, genes change between modules



→ Use module membership prob. as quantitative traits

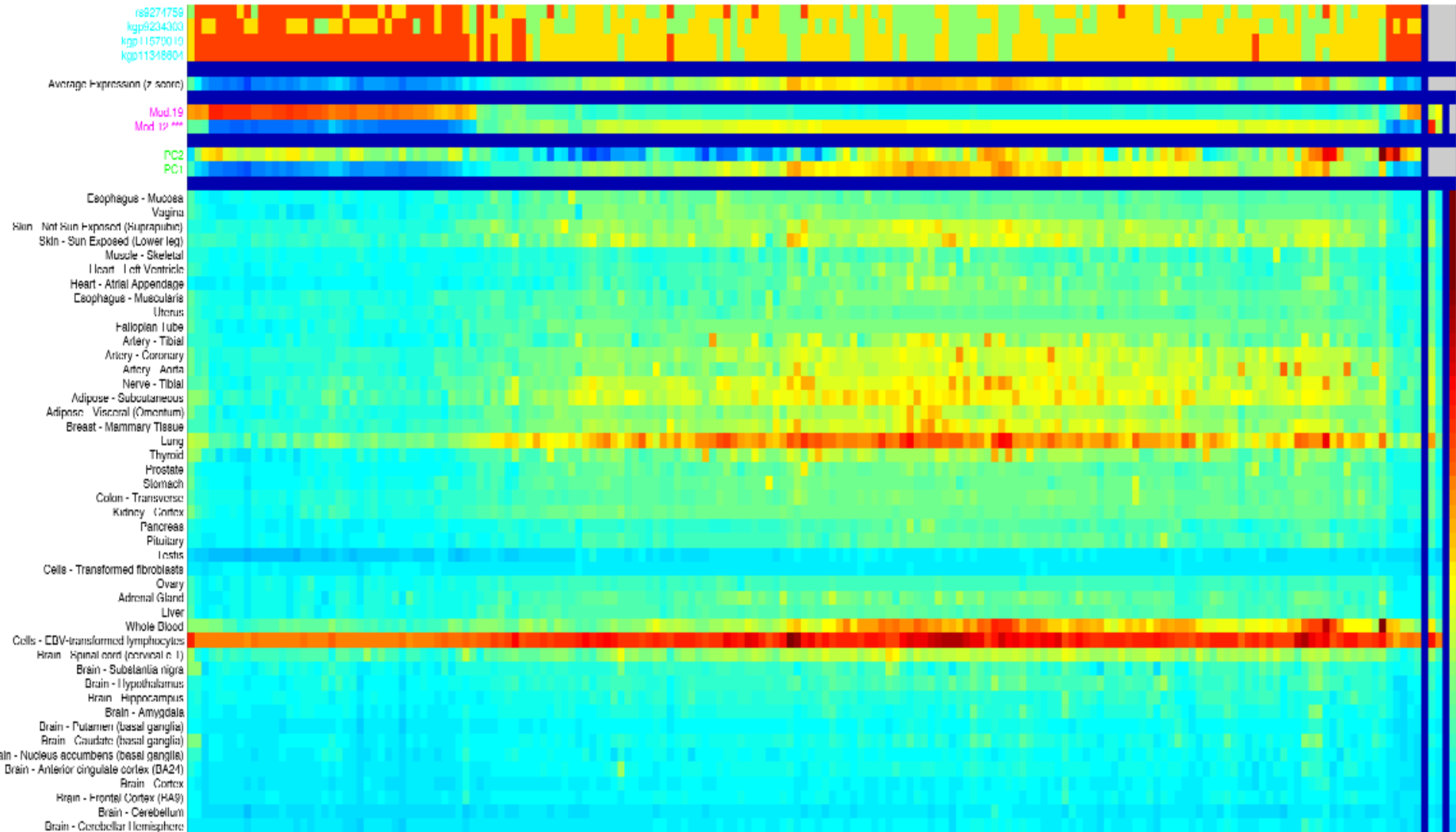
Identify SNPs underlying module changes: netQTLs



- Single-tissue expression patterns only a partial picture

Ex2: Mod19(lymph only) ⇔ Mod12(mesdoerm-wide)

Variation of
HLA-DQA1

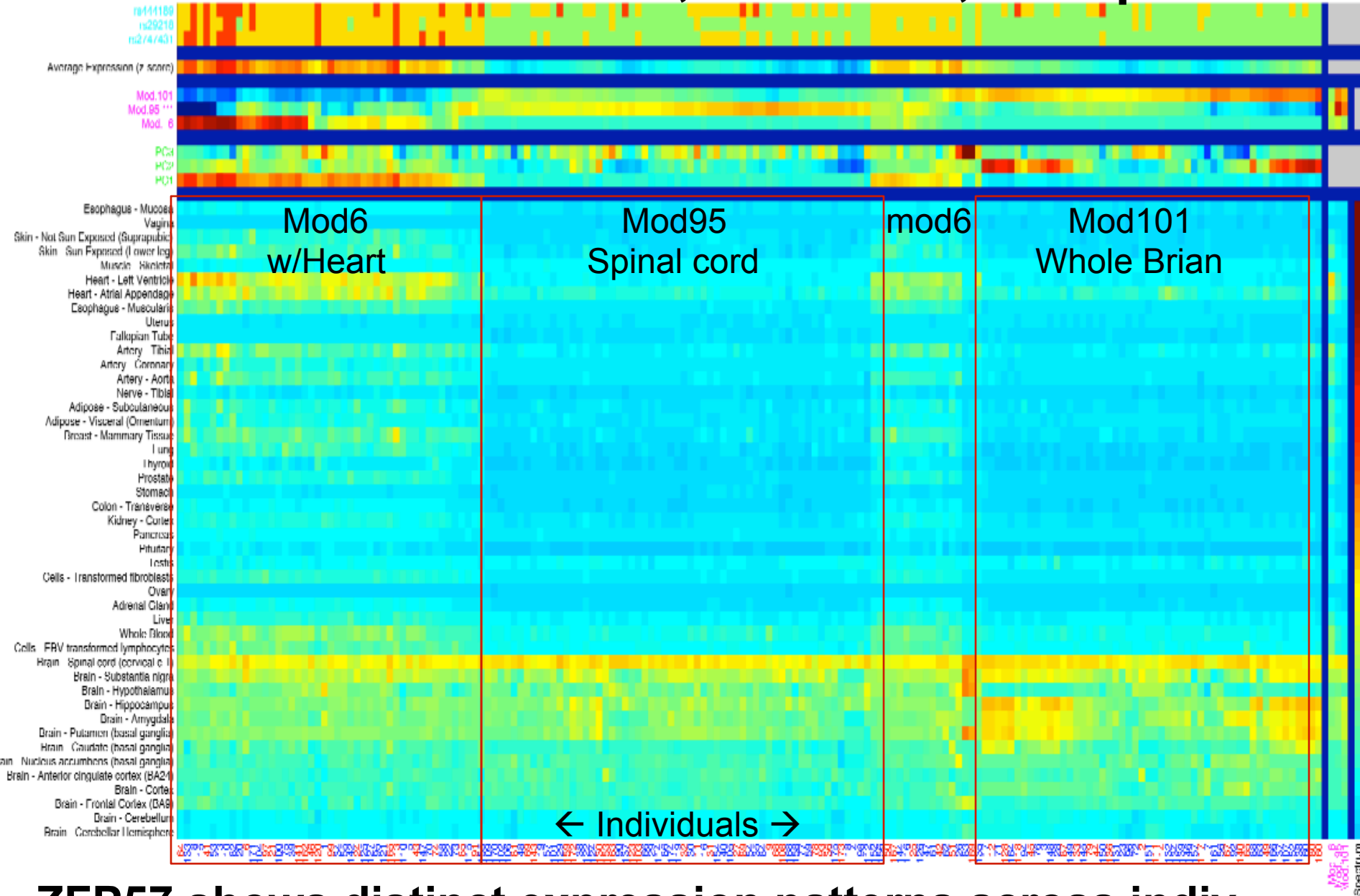


← Individuals →

Mod 19
Spectrum

- Multiple SNPs associated with different PCs of variation

Ex3: ZFP57 three modules, three PCs, multiple SNPs

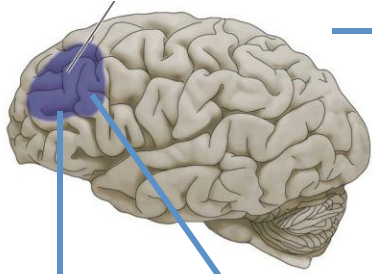


- ZFP57 shows distinct expression patterns across indiv.
- ZNF-KRAB regulator, transient neonatal diabetes mellitus

EWAS: Global association of brain enhancers with AD

MAP Memory and Aging Project
+ ROS Religious Order Study

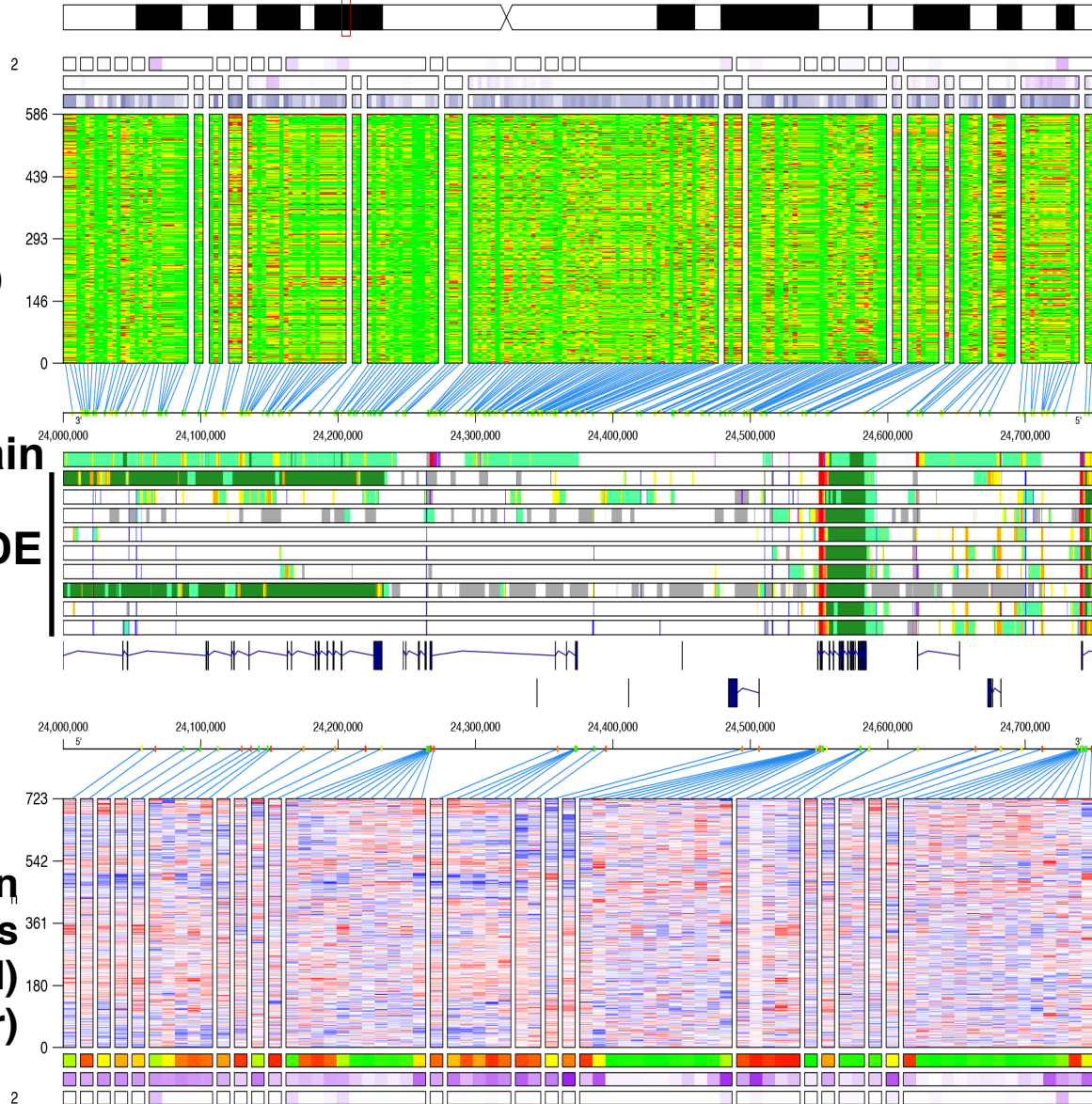
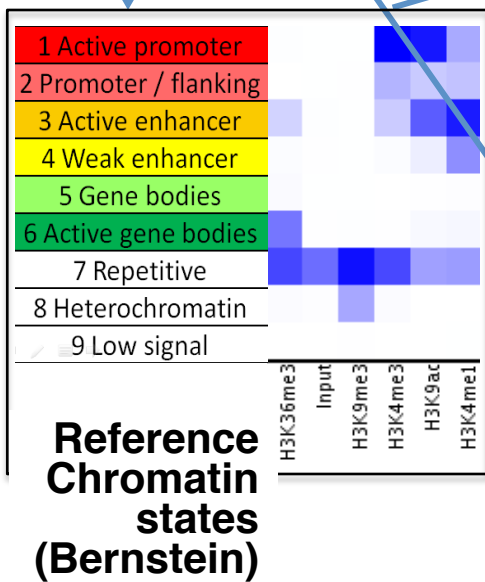
Dorsolateral PFC



Genotype
(1M SNPs
x700 ind.)
(De Jager)

**Brain
ENCODE**

Methylation
(450k probes
x 700 ind)
(De Jager)



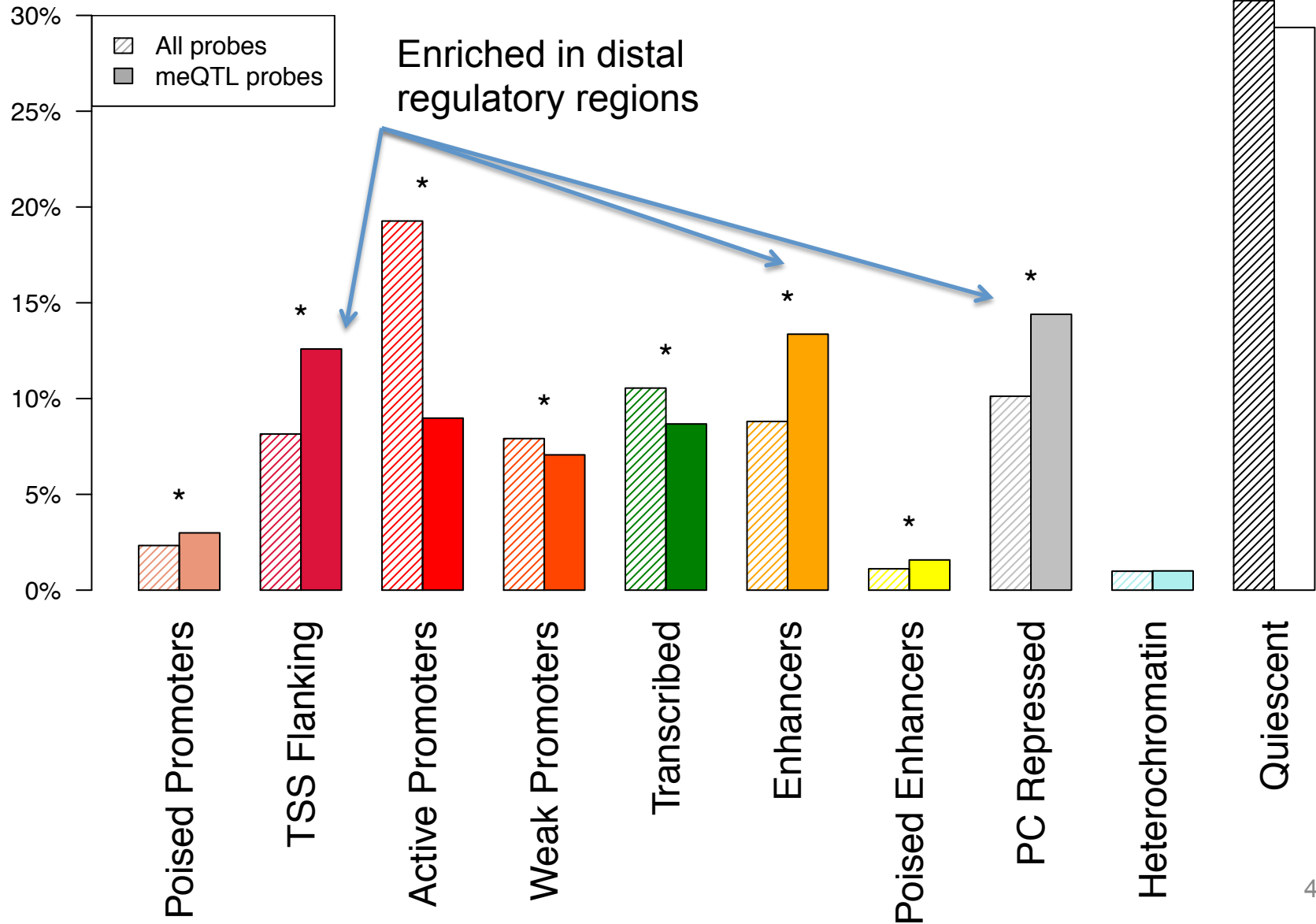
750 subjects, initially cognitively normal, Alzheimer's diagnosed by pathology. (Bennett)

Majority of AD-associated GWAS SNPs are meQTLs

Rank	ad.rsid	Gene	Description	meQTL P-value	meQTL SNP	meQTL Gene	SNP state
1	rs11767557	EPHA1	Ephrene A receptor 1	1.12E-13	rs12703526	cg18997129	24_Quies3
2	rs1532278	CLU	Clusterin	6.68E-125	rs17057441	cg18814083	22_Quies1
3	rs3865444	CD33	Myeloid transmembrane receptor	1.13E-10	rs12971624	cg11581627	22_Quies1
4	rs561655	PICALM	Phosphatidylinositol binding clathrin assembly protein	7.11E-77	rs17817919	cg24166175	22_Quies1
5	rs610932	MS4A2	Immunoglobulin receptor subunit	1.37E-33	rs562028	cg16954525	22_Quies1
6	rs6701713	CR1	Complement Receptor 1	1.39E-21	rs3849266	cg19373649	22_Quies1
7	rs7561528	BIN1	Bridging Integrator Nucleocytoplasmic adaptor protein	3.73E-176	rs4663104	cg02887598	2_TssF
8	rs9349407	CD2AP	Actin Cytoskeleton Regulating Scaffold	7.12E-63	rs2275446	cg16361253	22_Quies1

- Importance of mapping intermediate phenotypes
- Genetic ↔ Molecular ↔ Cellular ↔ Neural ↔ Disease

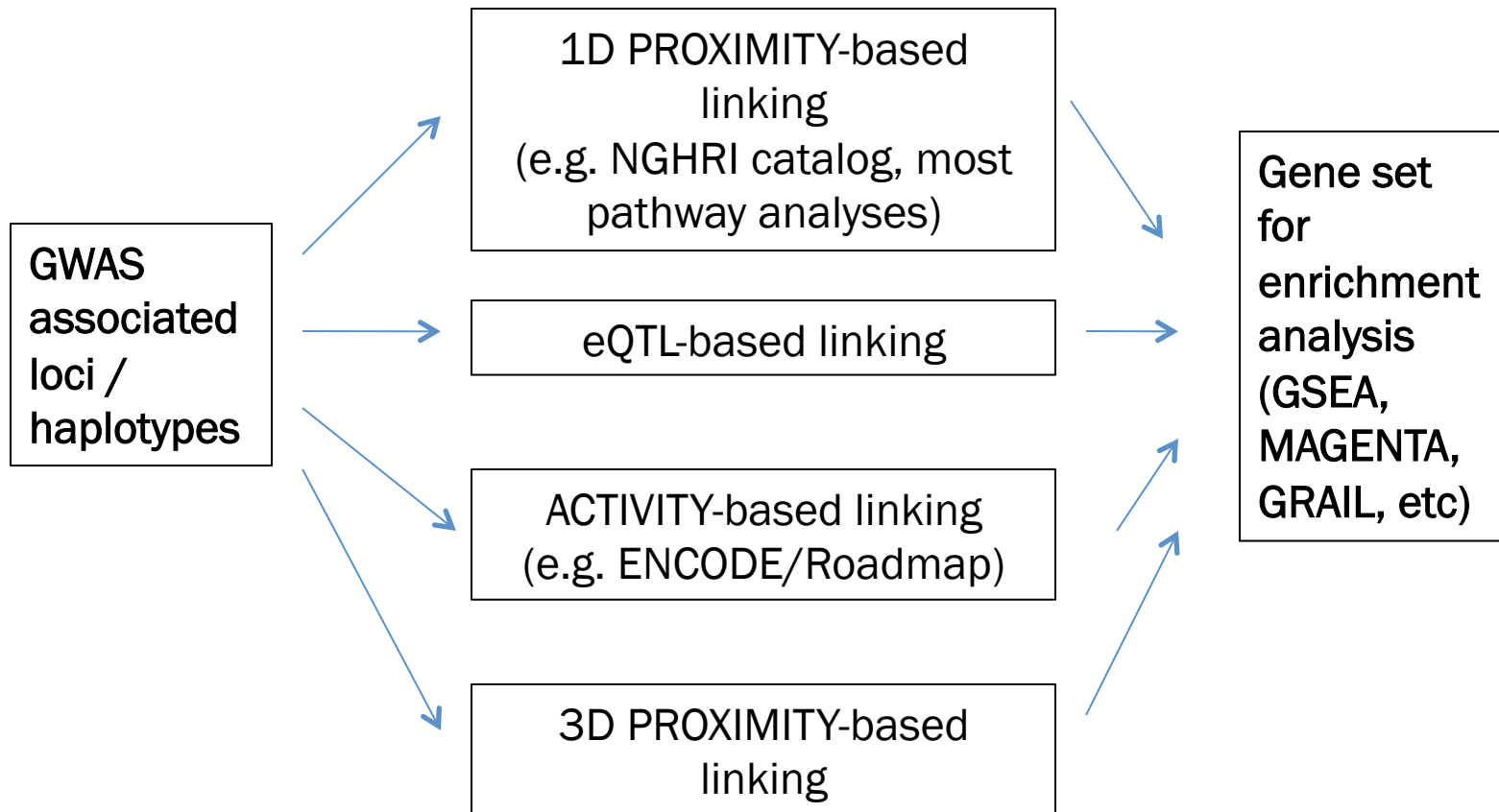
SNP-associated CpGs depleted in promoters, enriched in enhancers



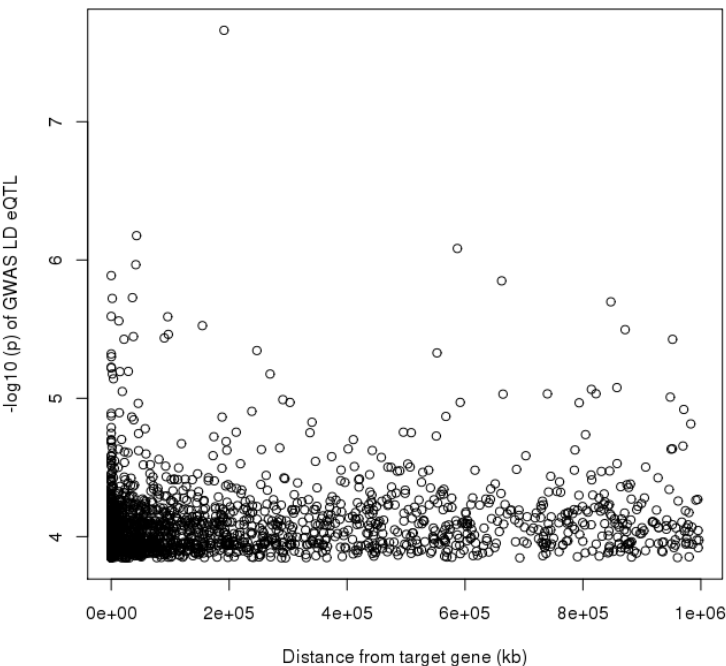
Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

Pathway analysis of GWAS relies on accurately linking cis-regulatory regions to their targets

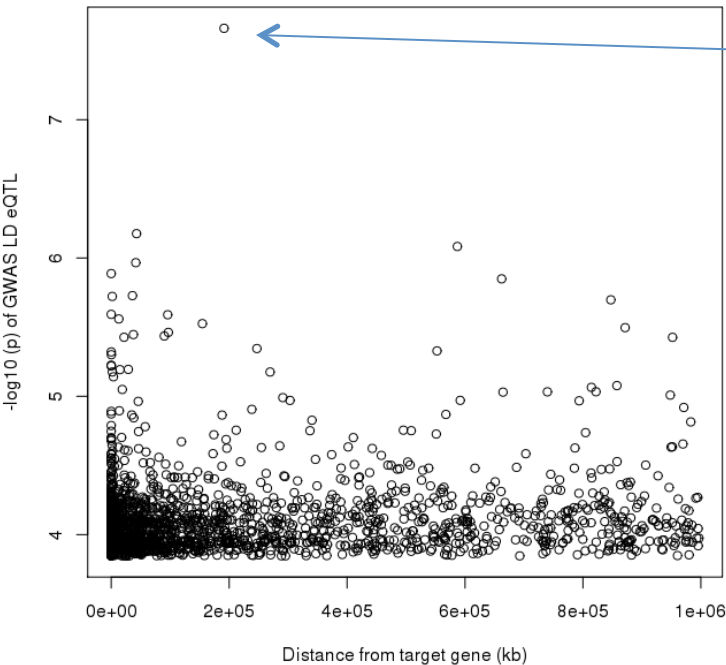


eQTLs reveal that closest isn't always best



- 1877 GWAS SNPs are within $r^2 \geq 0.8$ of the best GTEx eQTL for a gene/tissue combination
- Only 345/1810 (19%) of phenotype-associated GTEx eQTLs show agreement between their strongest eQTL-linked gene(s) and their physically closest or overlapping gene(s)

eQTLs for improved target gene prediction of regulatory GWAS SNPs



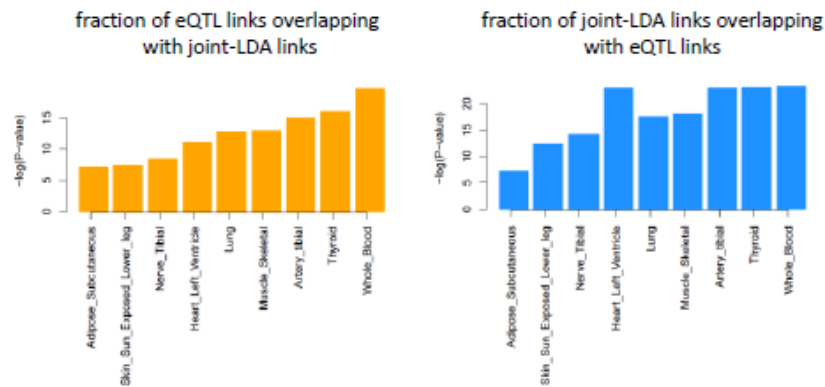
- rs919129 (myocardial infarction) is 25 kb from proximity-based target FBN. Strongest eQTL target is a novel lincRNA 191 kb away, RP11-506F22.2
- Alters a TBX5 motif instance (role in heart development)

Regulatory motifs altered

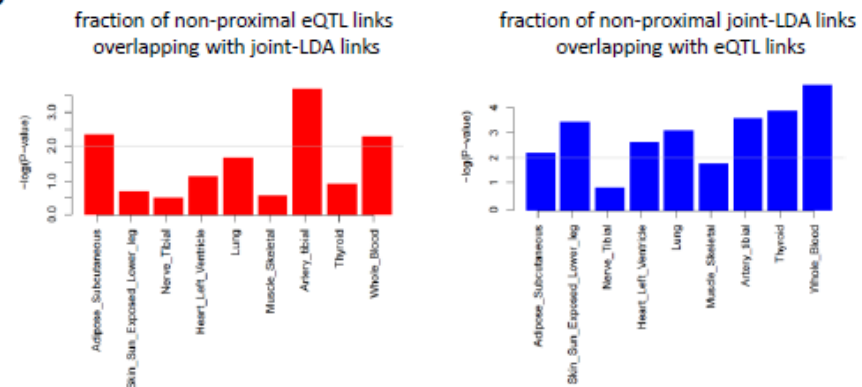
PWM	Strand	Ref	Alt	Match on:
				Ref: AAGTGCAAAAAGGAGGTATGGGAGAATGTGTGAGT Alt: AAGTGCAAAAAGGAGGTATGGGAGAATGTATGAGT
TBX5_3	+	6.7	-5.3	MAGGTGTGAR

Activity, eQTL, and conformation methods agree on links

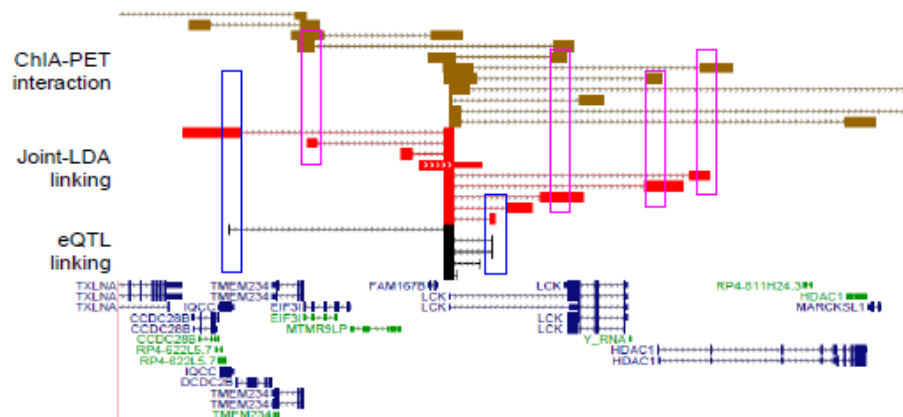
A



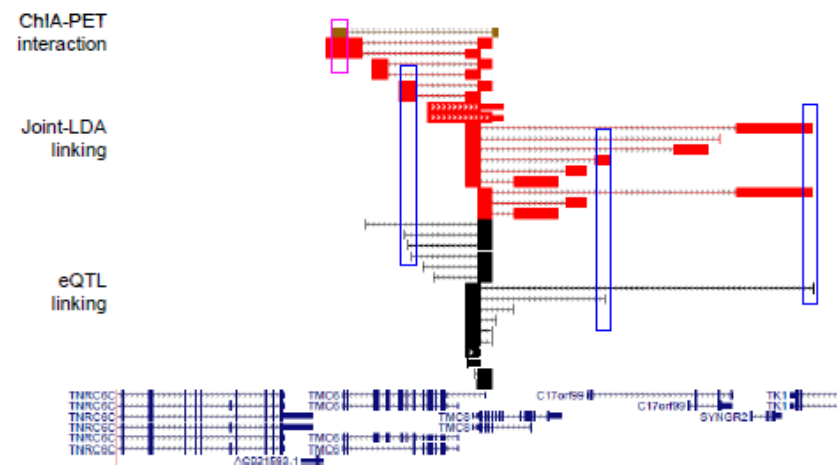
B



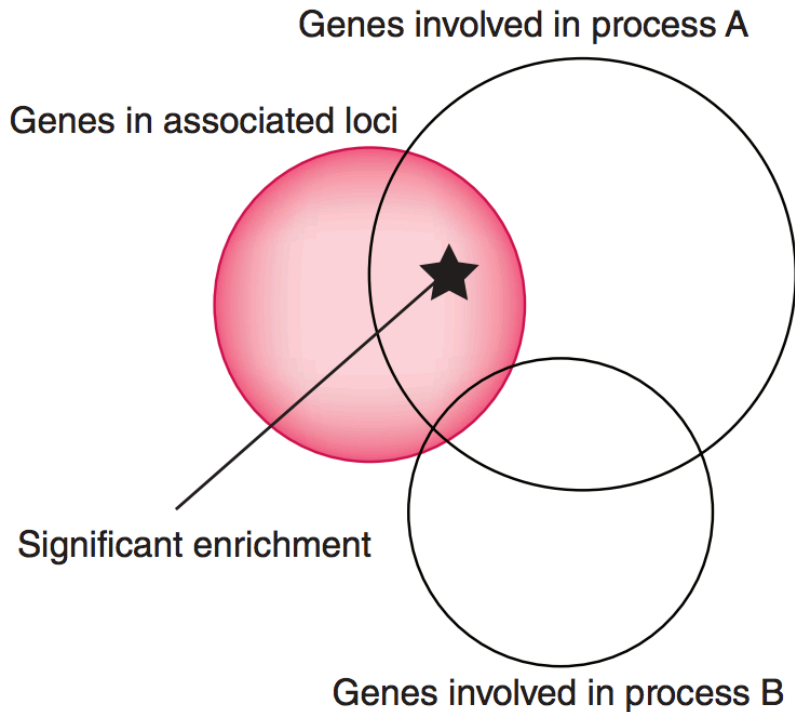
C



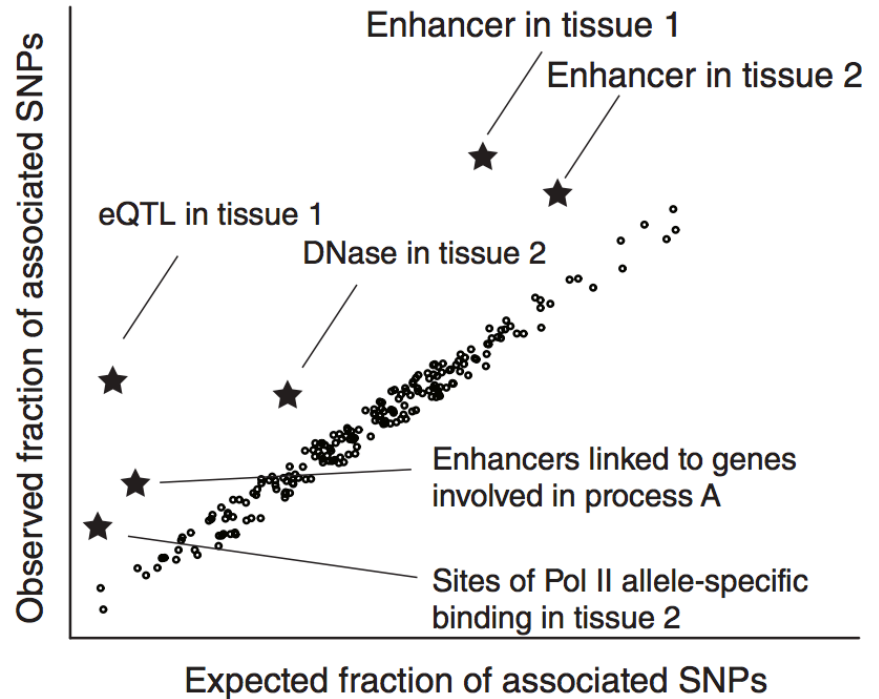
D



Linking as another way to improve GWAS analysis

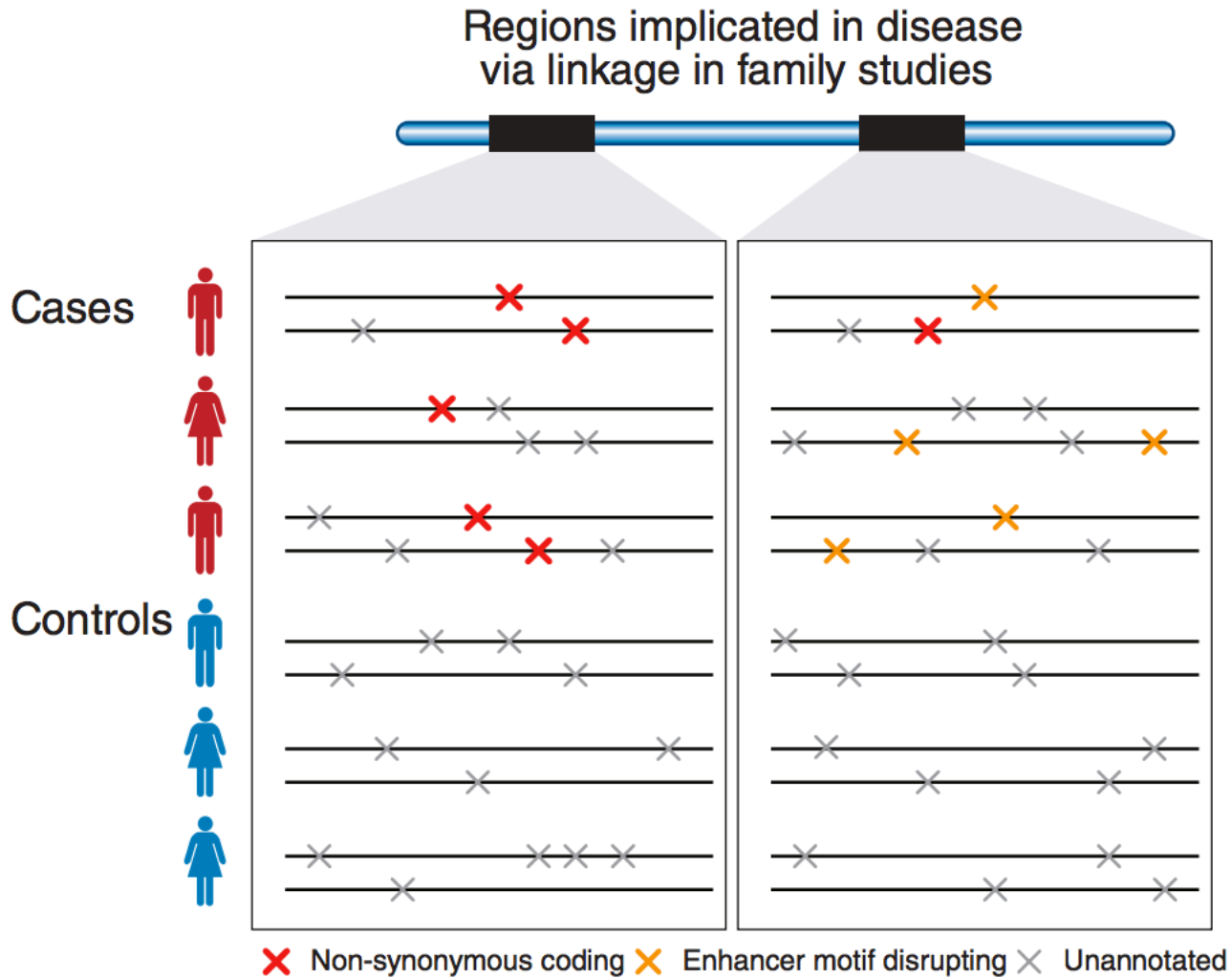


Current methods

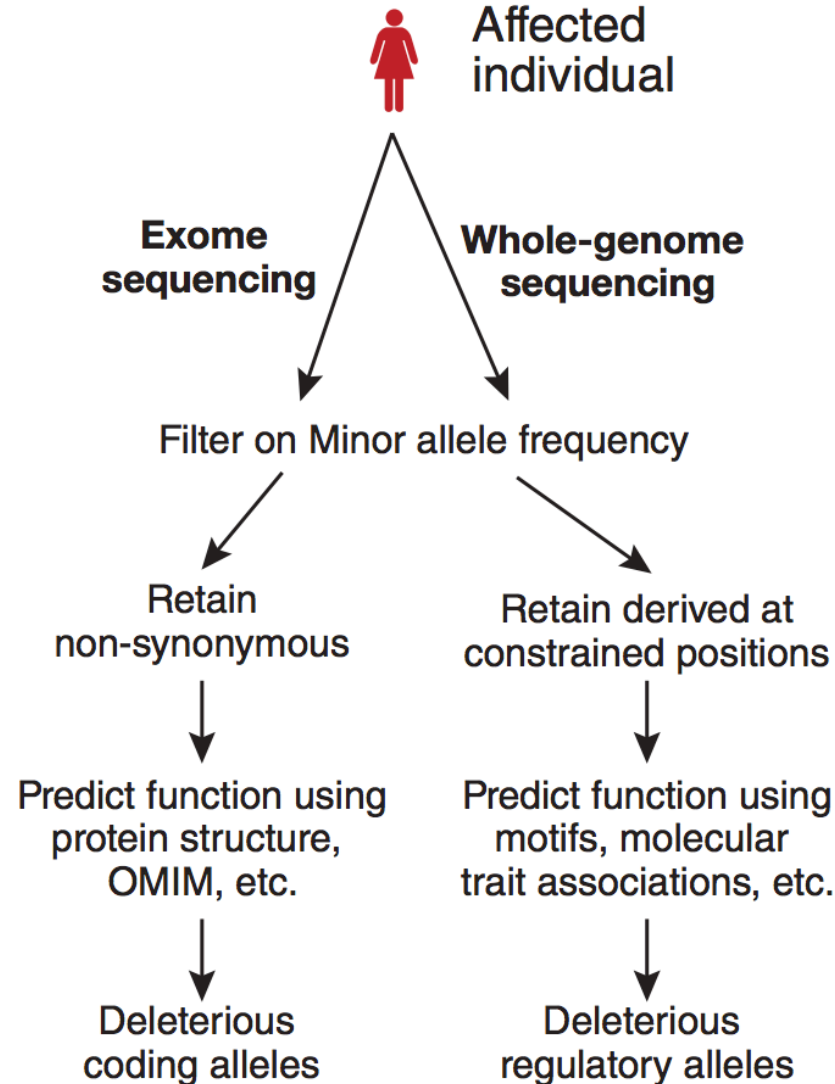


Regulatory genomics informed methods

Annotating rare variants: for burden tests on WGS data



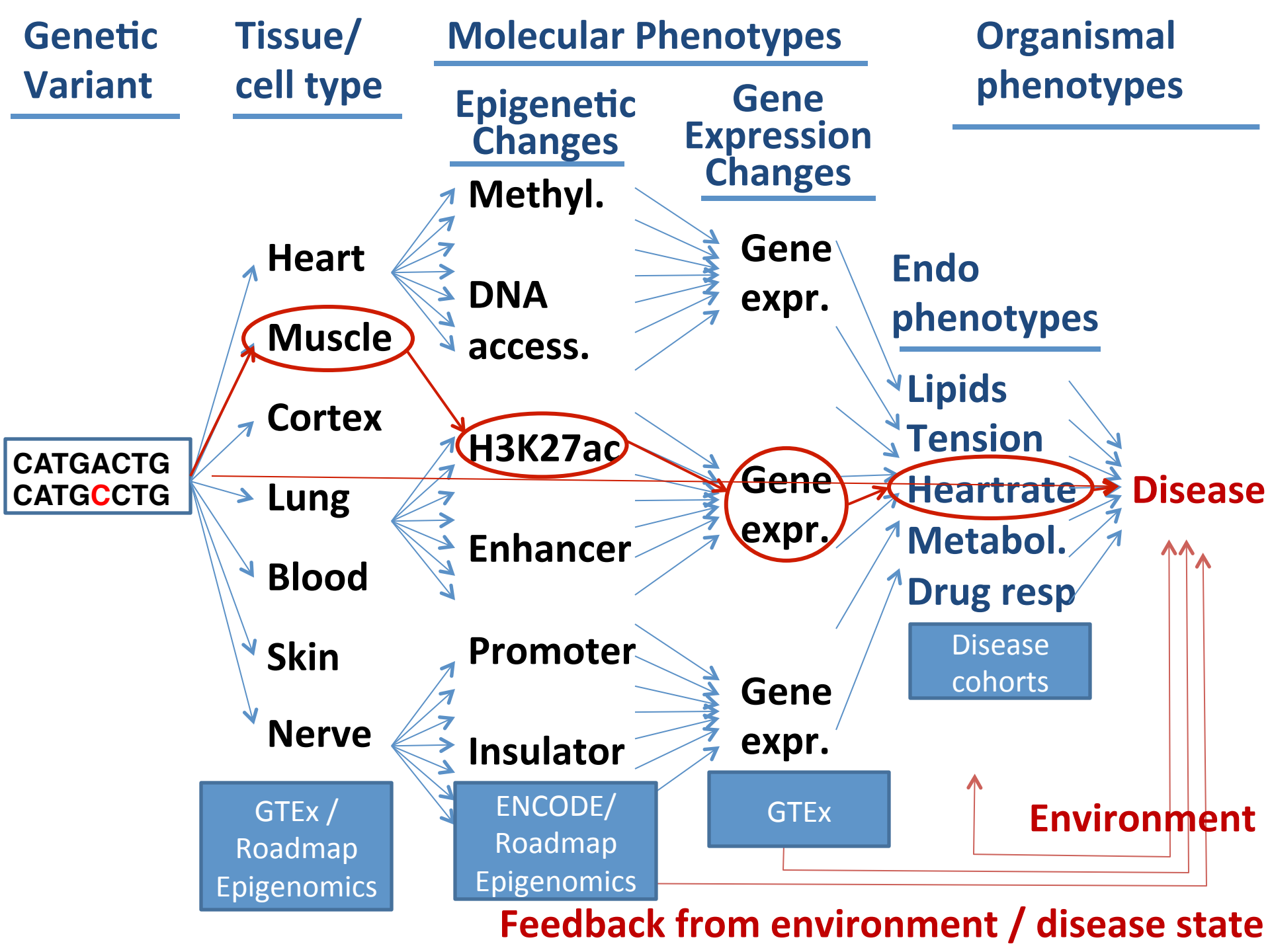
Annotating rare variants: for clinical interpretation pipelines



Regulatory genomics to interpret complex disease genetics

1. Regulatory annotations of the human genome: an overview
2. Using regulatory annotations to interpret GWAS
 - a. Locus level
 - b. Systems level
3. Beyond GWAS
 - a. Molecular variability
 - b. Empowering rare-variant and pathway analysis

The (very) big picture



Molecular Variation in reference individuals

GTEx

- eQTLs, meQTLs, enhQTLs
- Pinpoint regulatory regions
- Identify cell types where common variants act
- Link SNPs to target genes/regions

Functional Genomics and Epigenomics in reference cell types

**ENCODE
Roadmap
Epigenomics**

- Fine-map top-scoring loci
- Identify relevant cell types
- Identify relevant pathways
- Detect additional loci

Genetic Variation

CATGACTG
CATG**C**CTG

GWAS

- Top-scoring loci
- P-values, effect sizes
- Agnostic to mechanism

Disease

Reference variation maps

- Define LD blocks
- Common/rare variants
- Imputation / LD pruning

**1000 Genomes
HapMap**

Molecular Variation in cases/controls

**Disease
epigenomics**

- Intermediate molecular phenotypes
- Measured in disease-relevant tissues
- Capture environmental effects
- Capture downstream disease effects

**Environment
Covariates**

Acknowledgements

Kellis lab:

- Manolis Kellis
- Pouya Kheradpour
- Jason Ernst
- Benjamin Iriarte
- Abhishek Sarkar
- Matthew Eaton
- Wouter Meuleman
- Anshul Kundaje
- Jianrong Wang

ENCODE Project Consortium

Roadmap Epigenome Mapping Consortium

GTEx Project Consortium



Regulatory genomics to interpret complex disease genetics

Luke Ward

Manolis Kellis lab, MIT

July 18, 2013