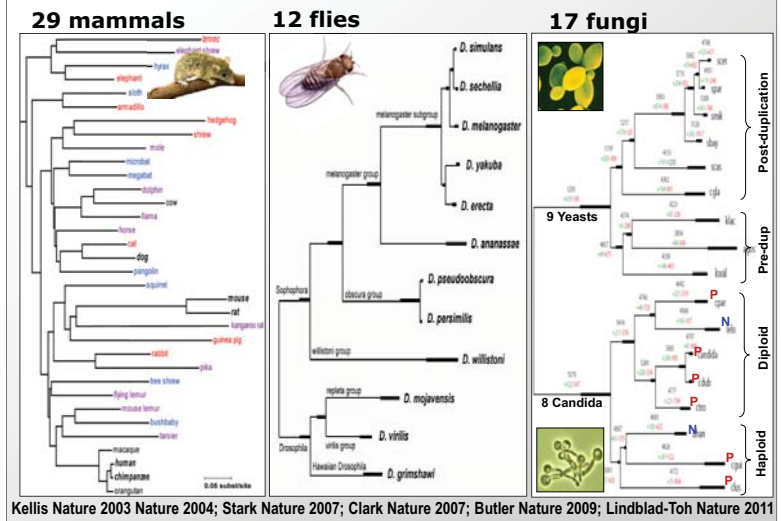


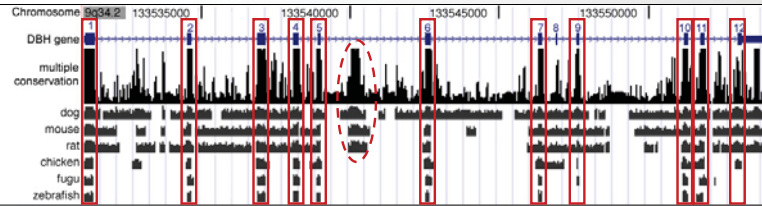
Systems-level views of disease epigenomics

- Evolutionary signatures → gene/genome annotation**
 - High-resolution annotation: genes, RNAs, motif instances
 - Measuring selection within the human population
- Chromatin states for interpreting disease association**
 - Annotate dynamic regulatory elements in multiple cell types
 - Activity-based linking of regulators → enhancers → targets
- Interpreting disease-associated sequence variants**
 - Mechanistic predictions for individual top-scoring SNPs
 - Functional roles of 1000s of disease-associated SNPs
- Systematic manipulation of 2000+ human enhancers**
 - Test effect of single-motif and single-nucleotide disruptions
 - Role of activator/repressor motifs, disease-associated SNPs
- Personal genomes/epigenomes in health and disease**
 - Allele-specific activity. Alzheimer's ↔ brain methylation ↔ SNP
 - Global repression of distal enhancers. NRSF, ELK1, CTCF

Large-scale comparative genomics datasets

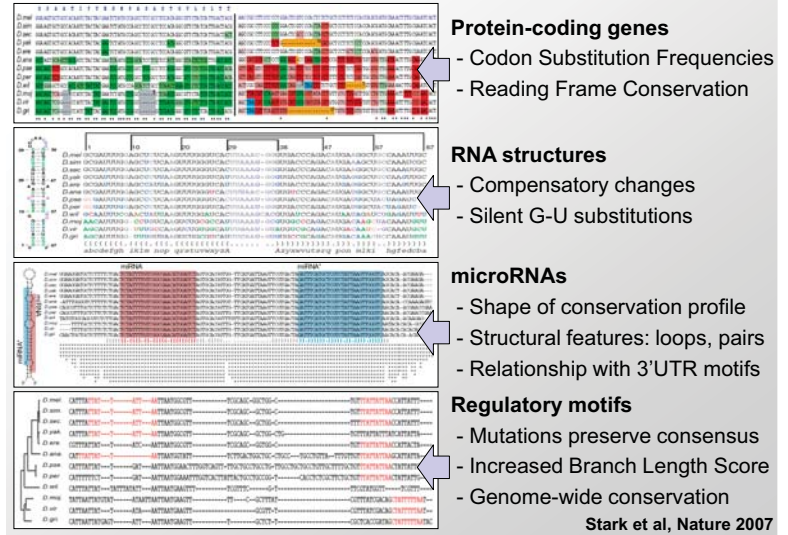


Comparative genomics and evolutionary signatures

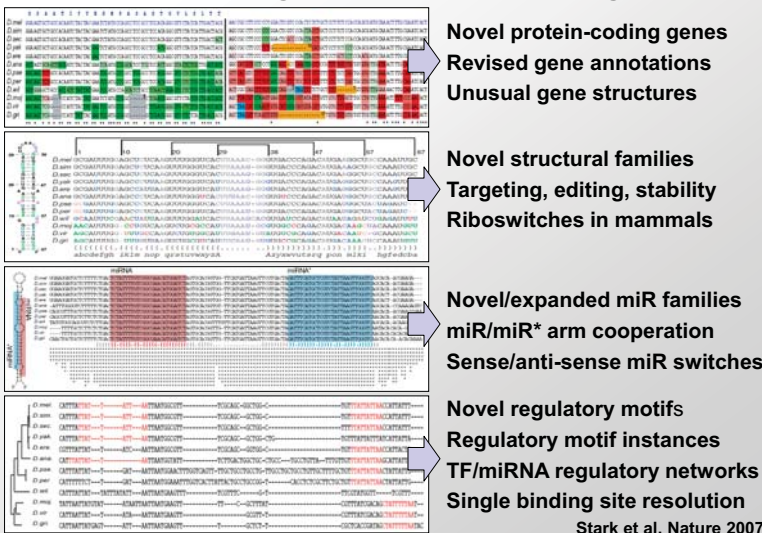


- Comparative genomics can reveal functional elements**
 - For example: exons are deeply conserved to mouse, chicken, fish
 - Many other elements are also strongly conserved: exons / regulatory?
- Can we also pinpoint specific functions of each region? Yes!**
 - Patterns of change distinguish different types of functional elements
 - Specific function ↔ Selective pressures ↔ Patterns of mutation/inse/del
- Develop evolutionary signatures characteristic of each function**
 Kellis Nature 2003 Nature 2004; Stark Nature 2007; Clark Nature 2007; Butler Nature 2009; Lindblad-Toh Nature 2011

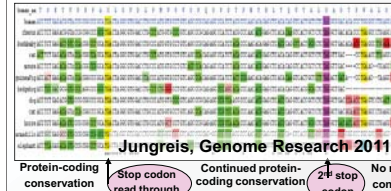
Evolutionary signatures for diverse functions



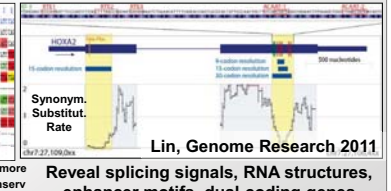
Implications for genome annotation / regulation



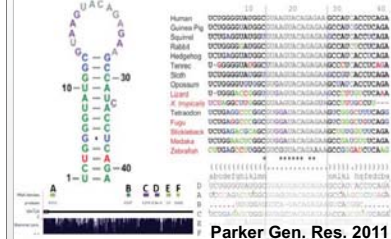
Translational read-through in human & fly



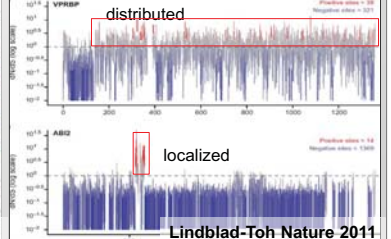
Overlapping selection in human exons



RNA structure families: ortholog/paralog cons

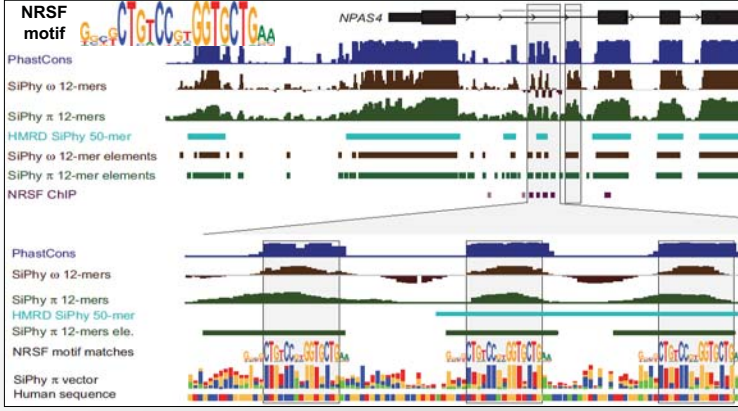


Regions of codon-level positive selection



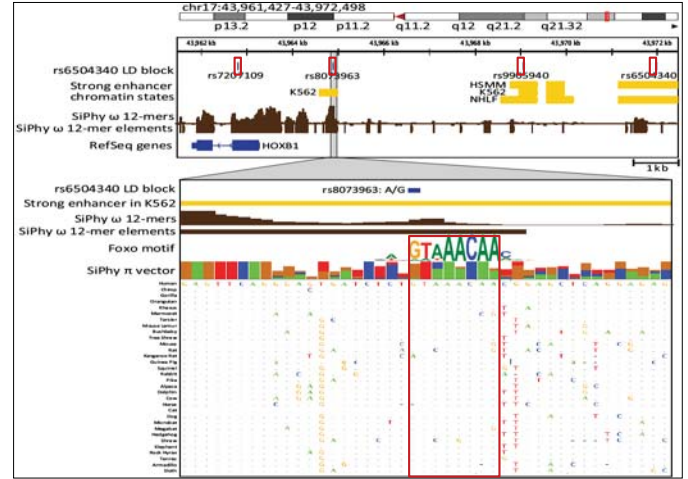
Distributed vs. localized positive selection
 Immunity/taste vs. retinal/bone/secretion

Measuring constraint at individual nucleotides



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

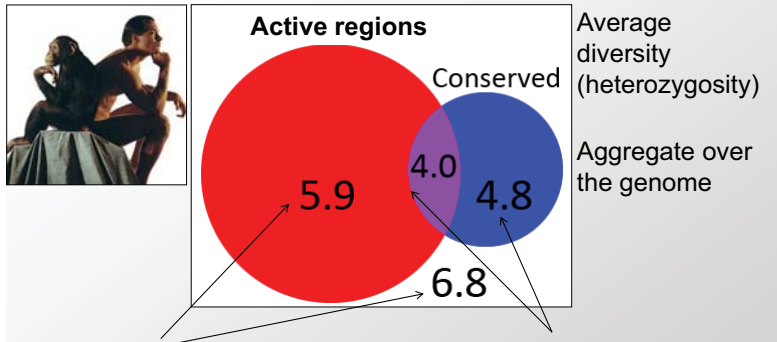
Detect SNPs that disrupt conserved regulatory motifs



- Functionally-associated SNPs enriched in states, constraint
- Prioritize candidates, increase resolution, disrupted motifs

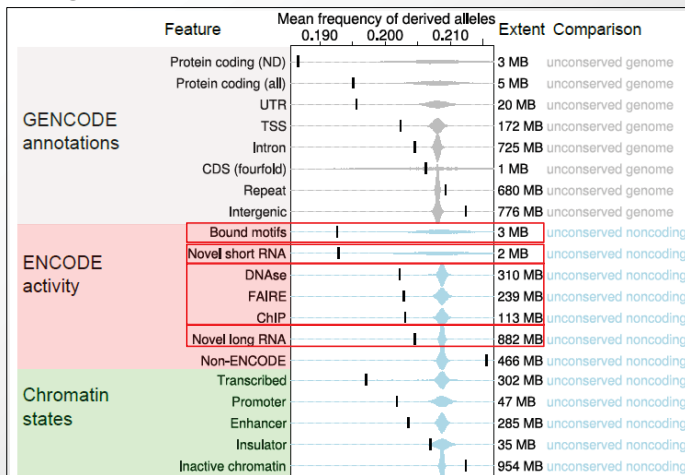
Measuring selection within the human lineage

Human constraint outside conserved regions



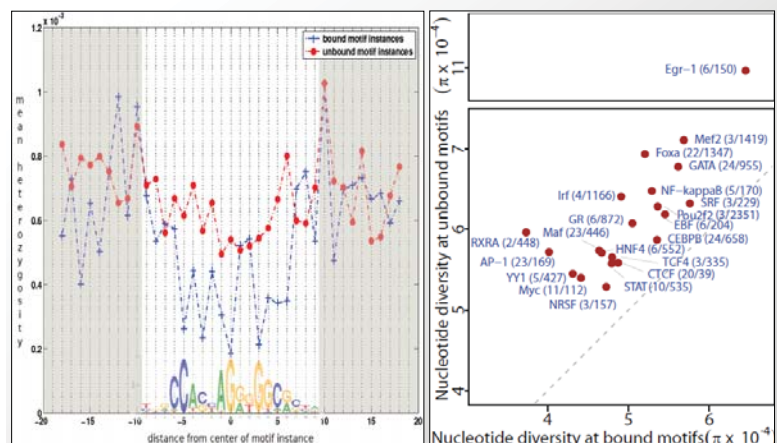
- **Non-conserved regions:**
 - ENCODE-active regions show reduced diversity
 - Lineage-specific constraint in biochemically-active regions
- **Conserved regions:**
 - Non-ENCODE regions show increased diversity
 - Loss of constraint in human when biochemically-inactive

Strongest: motifs, short RNA, Dnase, ChIP, IncRNA



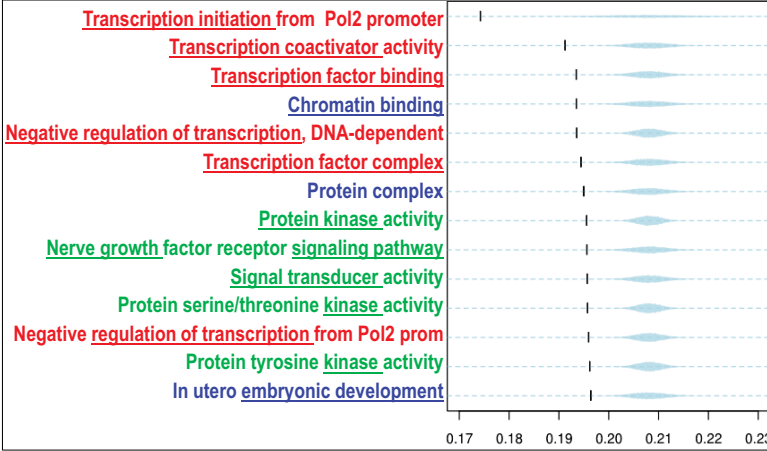
- Significant derived allele depletion in active features

Bound motifs show increased human constraint



- Position-specific reduction in bound motif heterozygosity
- Aggregate across thousands of CTCF motif instances

Most constrained human-specific enhancer functions

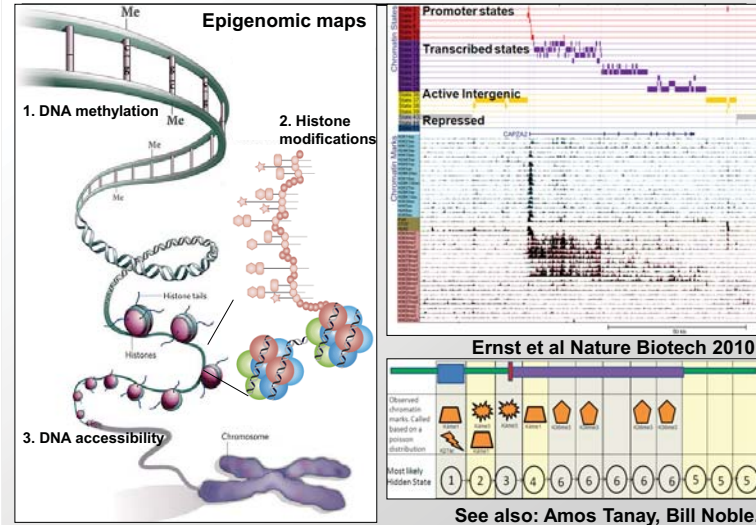


Regulatory genes: Transcription, Chromatin, Signaling.
Developmental enhancers: embryo, nerve growth

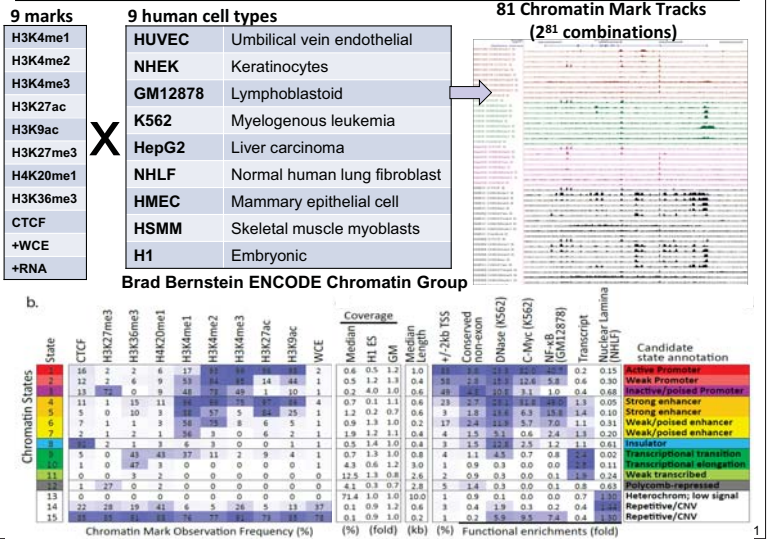
Systems-level views of disease epigenomics

- **Evolutionary signatures** → gene/genome annotation
 - High-resolution annotation: genes, RNAs, motif instances
 - Measuring selection within the human population
- **Chromatin states for interpreting disease association**
 - Annotate dynamic regulatory elements in multiple cell types
 - Activity-based linking of regulators → enhancers → targets
- **Interpreting disease-associated sequence variants**
 - Mechanistic predictions for individual top-scoring SNPs
 - Functional roles of 1000s of disease-associated SNPs
- **Systematic manipulation of 2000+ human enhancers**
 - Test effect of single-motif and single-nucleotide disruptions
 - Role of activator/repressor motifs, disease-associated SNPs
- **Personal genomes/epigenomes in health and disease**
 - Allele-specific activity. Alzheimer's ↔ brain methylation ↔ SNP
 - Global repression of distal enhancers. NRSF, ELK1, CTCF

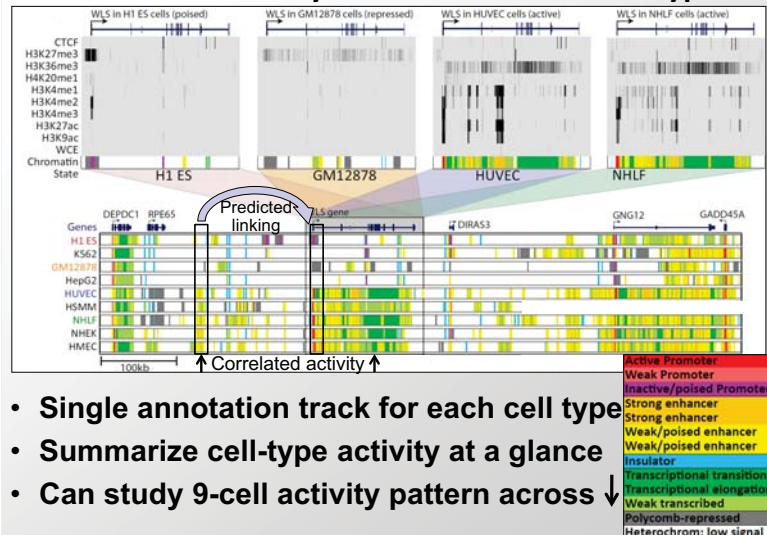
Chromatin signatures for genome annotation



ENCODE: Study nine marks in nine human cell lines

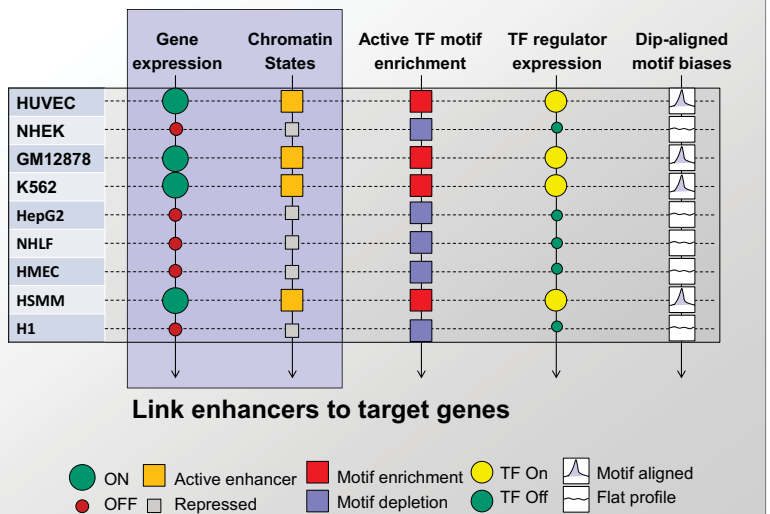


Chromatin states dynamics across nine cell types



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

Introducing multi-cell activity profiles

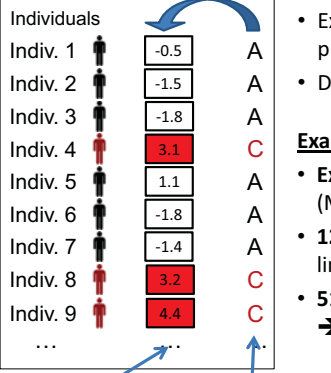


Link enhancers to target genes

- ON
- Active enhancer
- Motif enrichment
- TF On
- Motif aligned
- OFF
- Repressed
- Motif depletion
- TF Off
- Flat profile

Enhancer-gene links supported by eQTL-gene links

eQTL study



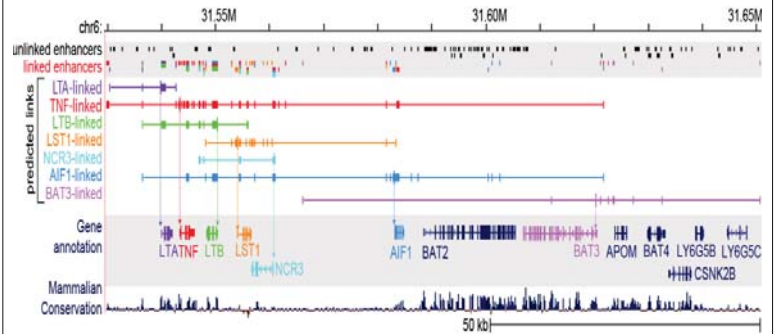
Validation rationale:

- 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)

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

Expression level of gene ↔ Sequence variant at distal position

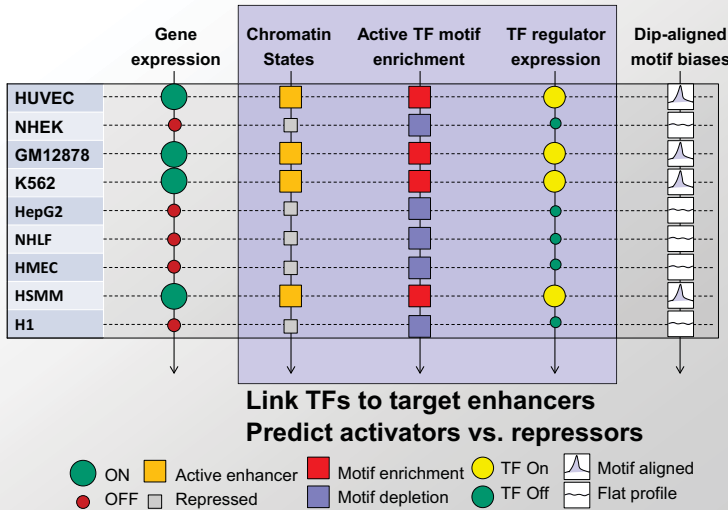
Visualizing 10,000s predicted enhancer-gene links



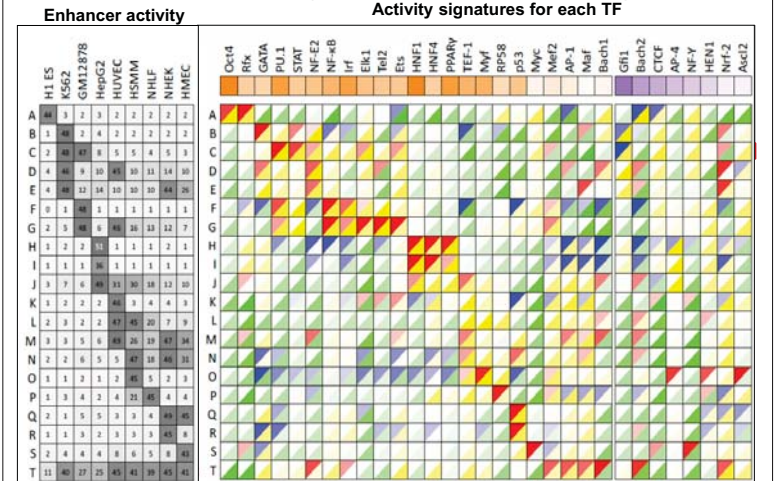
- Overlapping regulatory units, both few and many
- Both upstream and downstream elements linked
- Enhancers correlate with sequence constraint

26

Introducing multi-cell activity profiles

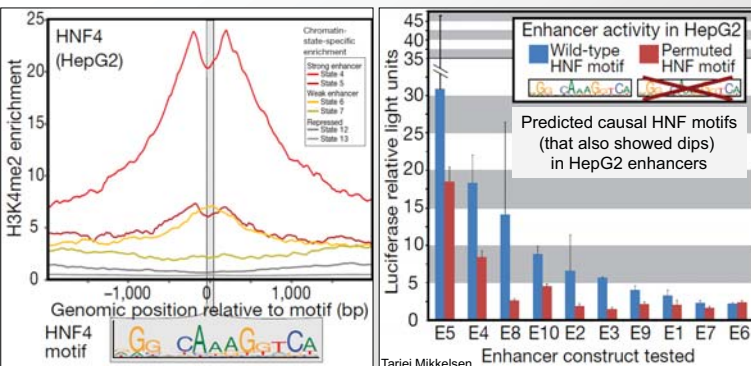


Coordinated activity reveals activators/repressors



- Enhancer networks: Regulator → enhancer → target gene

Causal motifs supported by dips & enhancer assays



→ Motifs bound by TF, contribute to enhancers

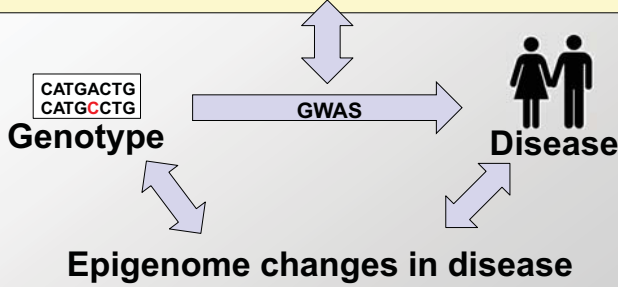
Systems-level views of disease epigenomics

- **Evolutionary signatures** → gene/genome annotation
 - High-resolution annotation: genes, RNAs, motif instances
 - Measuring selection within the human population
- **Chromatin states for interpreting disease association**
 - Annotate dynamic regulatory elements in multiple cell types
 - Activity-based linking of regulators → enhancers → targets
- **Interpreting disease-associated sequence variants**
 - Mechanistic predictions for individual top-scoring SNPs
 - Functional roles of 1000s of disease-associated SNPs
- **Systematic manipulation of 2000+ human enhancers**
 - Test effect of single-motif and single-nucleotide disruptions
 - Role of activator/repressor motifs, disease-associated SNPs
- **Personal genomes/epigenomes in health and disease**
 - Allele-specific activity. Alzheimer's ↔ brain methylation ↔ SNP
 - Global repression of distal enhancers. NRSF, ELK1, CTCF

Interpreting disease-association signals

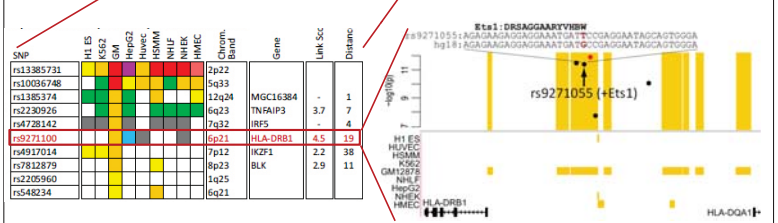
Interpret variants using Epigenomics

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



Revisiting disease-associated variants

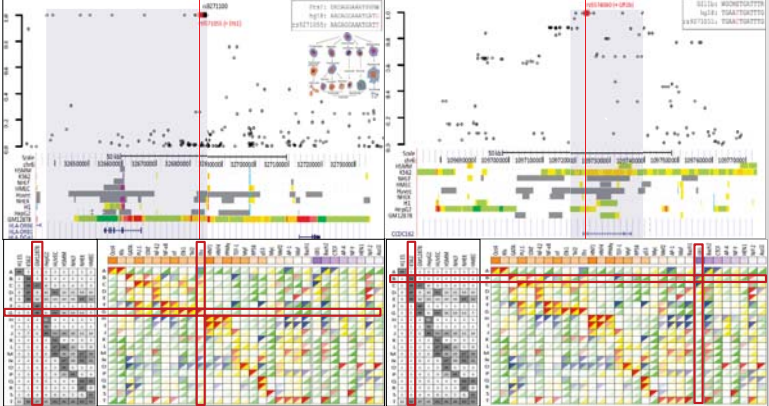
Phenotype	Top Cell Type	Total #SNPs from Study	#SNPs in enh States 4 and 5	p-value	FDR	HL ES	K562	GM12878	HepG2	HUVEC	HSMM	NHEK	NHEK	HMEC
Erythrocyte phenotypes (Ref. 38)	K562	35	9	<10 ⁻⁷	0.02	9	17	4	0	0	1	2	1	1
Blood lipids (Ref. 39)	HepG2	101	13	<10 ⁻⁷	0.02	3	5	0	11	2	3	3	4	3
Rheumatoid arthritis (Ref. 40)	GM12878	29	7	2.0 x 10 ⁻⁷	0.03	0	0	18	0	0	0	0	3	3
Primary biliary cirrhosis (Ref. 41)	GM12878	6	4	6.0 x 10 ⁻⁷	0.03	0	11	41	0	0	0	0	8	8
Systemic lupus erythematosus (Ref. 42)	GM12878	18	6	9.0 x 10 ⁻⁷	0.03	0	4	21	0	1	8	8	1	1
Lipoprotein cholesterol/triglycerides (Ref. 43)	HepG2	18	5	1.2 x 10 ⁻⁶	0.03	17	8	0	34	1	6	4	1	1
Hematological traits (Ref. 44)	K562	39	7	1.7 x 10 ⁻⁶	0.03	0	12	10	2	1	1	0	1	0
Hematological parameters (Ref. 45)	K562	28	6	2.2 x 10 ⁻⁶	0.03	0	15	7	0	5	7	7	3	2
Colorectal cancer (Ref. 46)	HepG2	3	3	3.8 x 10 ⁻⁶	0.03	0	0	0	46	0	13	0	13	13
Blood pressure (Ref. 47)	K562	9	4	5.0 x 10 ⁻⁶	0.04	0	0	10	4	10	6	7	1	11



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

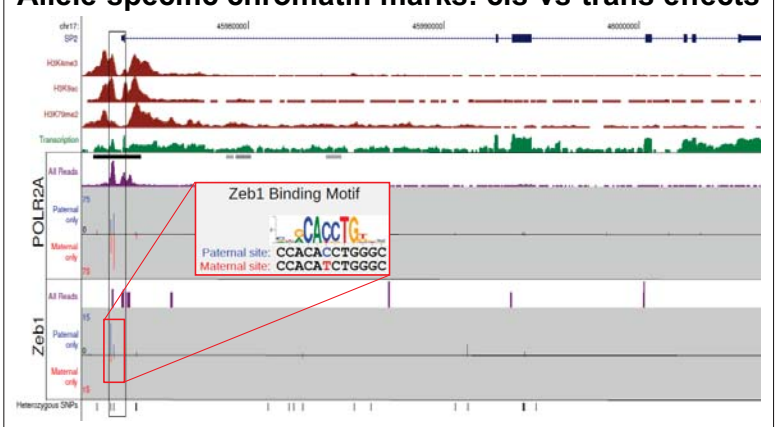
Mechanistic predictions for top disease-associated SNPs

Lupus erythematosus in GM lymphoblastoid Erythrocyte phenotypes in K562 leukemia cells



- Disrupt activator Ets-1 motif
- Loss of GM-specific activation
- Loss of enhancer function
- Loss of HLA-DRB1 expression
- Creation of repressor Gfi1 motif
- Gain K562-specific repression
- Loss of enhancer function
- Loss of CCDC162 expression

Allele-specific chromatin marks: cis-vs-trans effects



- Maternal and paternal GM12878 genomes sequenced
- Map reads to phased genome, handle SNPs indels
- Correlate activity changes with sequence differences

HaploReg: systematic ENCODE mining of variants (compbio.mit.edu/HaploReg)

Query SNP: rs17145713 and variants with r² >= 0.95

chr	pos (hg19)	LD variant	Ref	Alt	ASN	CEU	YRI	GERP	SfPhy	Promoter	Enhancer	DNAse	Proteins bound	Motifs changed	GENCODE genes	RefSeq genes	dbSNP func annot	
7	72842724	1	7.72480960	GAC	G	0	0.15	0					RANX-1		5.4kb 5' of FZD9	5.4kb 5' of FZD9		
7	72856430	1	rs1178979	T	C	0.13	0.18	0.3					DLL	GATA	BAZ1B	BAZ1B	intrinsic	
7	72857049	1	rs1178979	A	G	0.14	0.18	0.3						AREBB, DEC	BAZ1B	BAZ1B	intrinsic	
7	72857713	1	rs13064203	CA	G	0.13	0.11	0.2					8 cell types	Sox	BAZ1B	BAZ1B	intrinsic	
7	72866522	1	rs13064076	A	G	0.12	0.10	0.26							BAZ1B	BAZ1B	intrinsic	
7	72863105	1	rs13064144	T	C	0.14	0.18	0.29						HMEC, IPS, AS49, HL-HESC, HT-HESC, HA-h, HNPC, HPC, CHRE	BAZ1B	BAZ1B	intrinsic	
7	72865810	1	rs6879500	G	A	0.14	0.18	0.39							BAZ1B	BAZ1B	intrinsic	
7	72904810	1	rs17145713	C	T	0.14	0.18	0.3						ATF3	BAZ1B	BAZ1B	intrinsic	
7	72938244	1	rs11983997	G	C	0.13	0.18	0.26					GM12878, K562		2.8kb 5' of BAZ1B	2.8kb 5' 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 5' of BCL7B	
7	72980969	1	rs3559126	C	T	0.13	0.18	0.08							TBL2	TBL2	TBL2	intrinsic
7	72989141	1	rs34550818	C	CA	0.11	0.12	0							TBL2	TBL2	TBL2	intrinsic
7	72989390	1	rs11974409	A	G	0.13	0.18	0.14							TBL2	TBL2	TBL2	intrinsic
7	72989852	1	rs9638180	A	G	0.12	0.18	0.08							Zfp33	5.8kb 5' of TBL2	5.8kb 5' of TBL2	
7	72999105	1	rs9638182	T	G	0.12	0.18	0.14								8kb 5' of TBL2	6.1kb 5' of TBL2	
7	73007943	1	rs1051921	G	A	0.12	0.18	0.08					4 cell types	POL2	MLXPL	MLXPL	3' UTR	

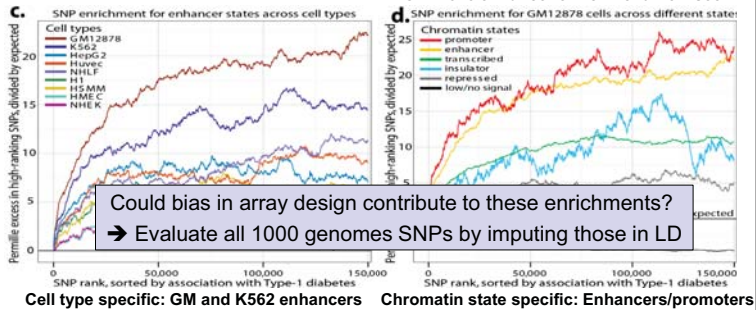
- Start with any list of SNPs or select a GWA study
 - Mine publically available ENCODE data for significant hits
 - Hundreds of assays, dozens of cells, conservation, motifs
 - Report significant overlaps and link to info/browser

Functional enrichment for 1000s of SNPs

Full T1D association spectrum → 1000s of causal SNPs

- Rank all SNPs by P-value
- Find chromatin states with enrichment in high ranks
- Signal spans 1000s of SNPs

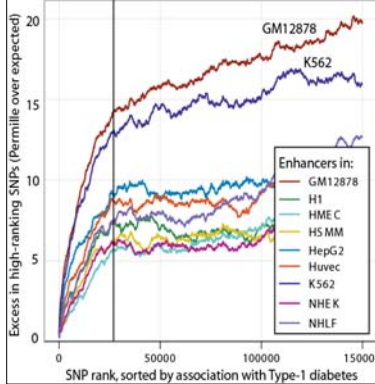
GM12878	Lymphoblastoid
K562	Myelogenous leukemia



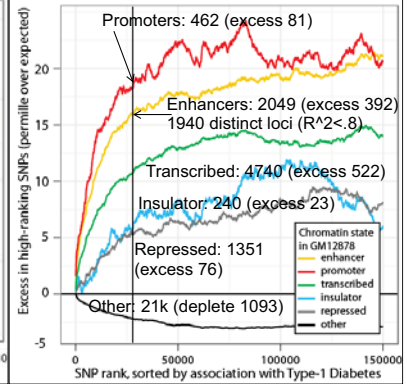
Could bias in array design contribute to these enrichments?
→ Evaluate all 1000 genomes SNPs by imputing those in LD

Imputing SNPs in LD → stronger cell/state separation

Enhancers across cell types



Chromatin states in GM12878

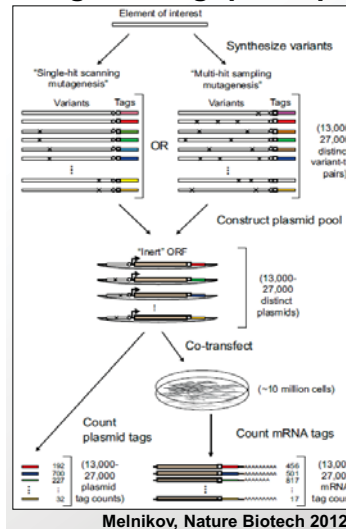


- Excess of 30,000 SNPs ⇔ 2049 enhancers (excess 392)
- Mostly found in independent loci (1730 with $R^2 < 0.2$)
- Systematically measure their regulatory contributions

Systems-level views of disease epigenomics

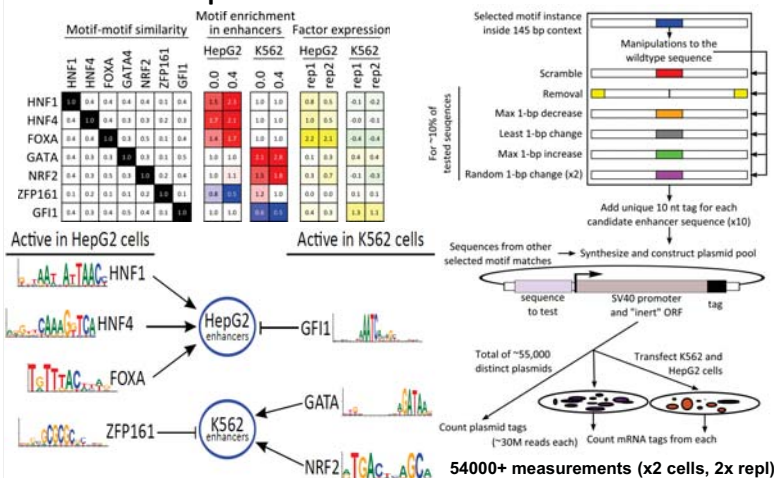
- Evolutionary signatures → gene/genome annotation**
 - High-resolution annotation: genes, RNAs, motif instances
 - Measuring selection within the human population
- Chromatin states for interpreting disease association**
 - Annotate dynamic regulatory elements in multiple cell types
 - Activity-based linking of regulators → enhancers → targets
- Interpreting disease-associated sequence variants**
 - Mechanistic predictions for individual top-scoring SNPs
 - Functional roles of 1000s of disease-associated SNPs
- Systematic manipulation of 2000+ human enhancers**
 - Test effect of single-motif and single-nucleotide disruptions
 - Role of activator/repressor motifs, disease-associated SNPs
- Personal genomes/epigenomes in health and disease**
 - Allele-specific activity. Alzheimer's ⇔ brain methylation ⇔ SNP
 - Global repression of distal enhancers. NRSF, ELK1, CTCF

High-throughput experiments: 10,000s enhancers

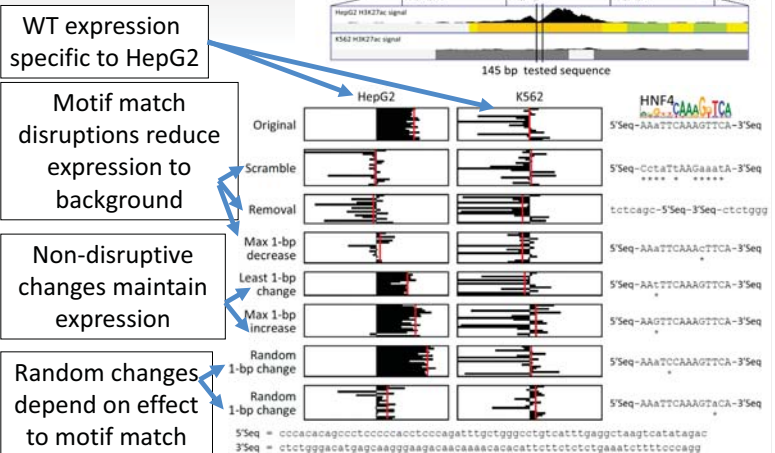


- Experiment features:**
 - Multiplexed enhancer assays
 - 10,000s of elements
 - Each w/ unique barcode
 - Multiple human cell types
 - Repeat experiments on same array / diff barcodes
- Applied to:**
 - Test enhancer offsets
 - Test causal motifs
- With: Tarjei Mikkelsen**
 - Broad Institute, ARRA funds
 - See also: Barak Cohen, Jay Shendure, Eran Segal

Systematic motif disruption for 5 activators and 2 repressors in 2 human cell lines

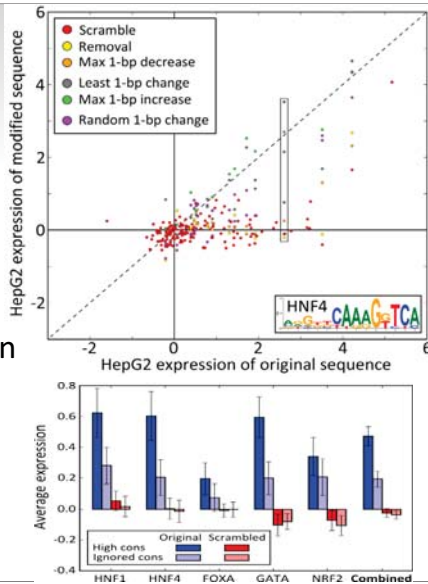


Example activator: conserved HNF4 motif match

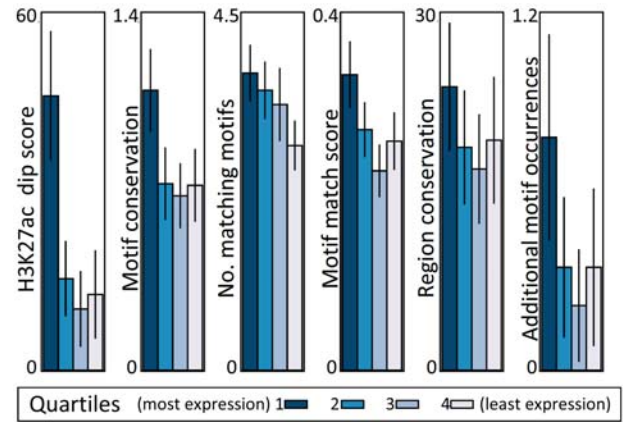


Results hold across 2000+ enhancers

- Scramble abolishes reporter expression
- Neutral mutations show no change
- Increasing mutations show more expression
- However, only 40% show wild-type expression: context?

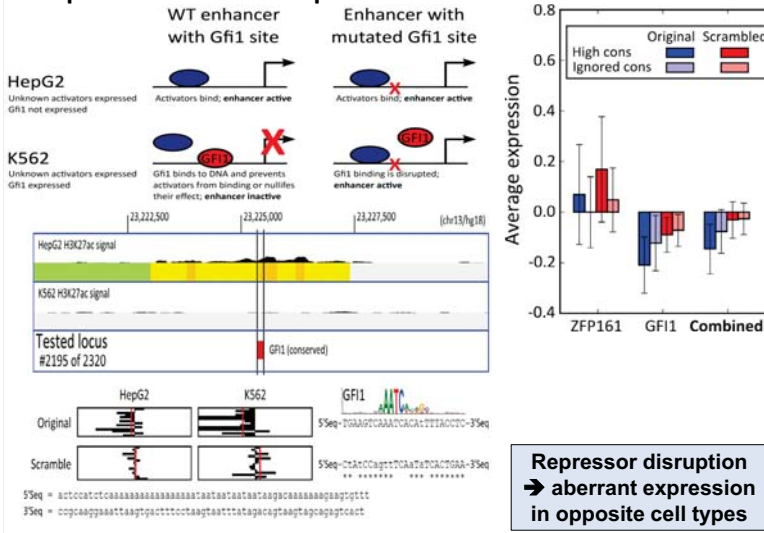


Features of functional wildtype enhancers



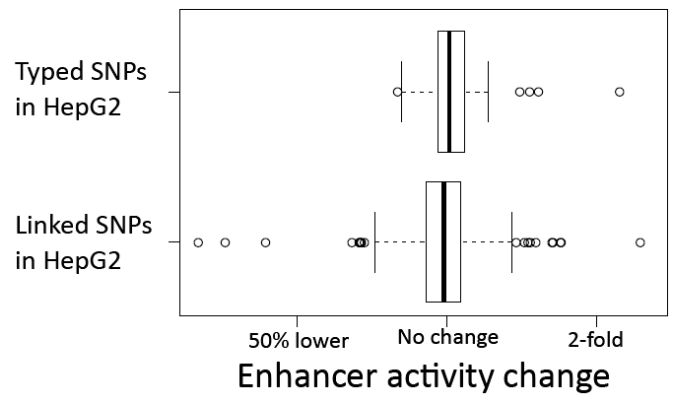
- Nucleosome exclusion, motif conservation, other TFs
- Each of these features is encoded in primary sequence

Repressors of HepG2 enhancer act in K562



Repressor disruption → aberrant expression in opposite cell types

Testing effect of SNP change in enhancer constructs



- SNPs in enhancer regions can lead to expression changes in downstream reporter genes
- Currently testing all T1D-associated enhancer SNPs

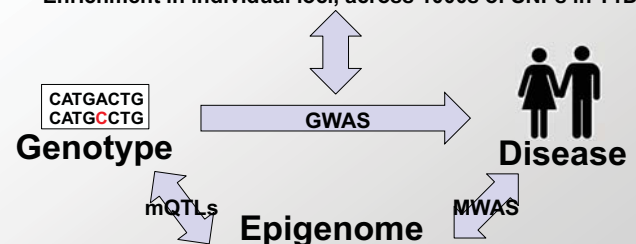
Systems-level views of disease epigenomics

- **Evolutionary signatures → gene/genome annotation**
 - High-resolution annotation: genes, RNAs, motif instances
 - Measuring selection within the human population
- **Chromatin states for interpreting disease association**
 - Annotate dynamic regulatory elements in multiple cell types
 - Activity-based linking of regulators → enhancers → targets
- **Interpreting disease-associated sequence variants**
 - Mechanistic predictions for individual top-scoring SNPs
 - Functional roles of 1000s of disease-associated SNPs
- **Systematic manipulation of 2000+ human enhancers**
 - Test effect of single-motif and single-nucleotide disruptions
 - Role of activator/repressor motifs, disease-associated SNPs
- **Personal genomes/epigenomes in health and disease**
 - Allele-specific activity. Alzheimer's ↔ brain methylation ↔ SNP
 - Global repression of distal enhancers. NRSF, ELK1, CTCF

Interpreting disease-association signals

(1) Interpret variants using Epigenomics

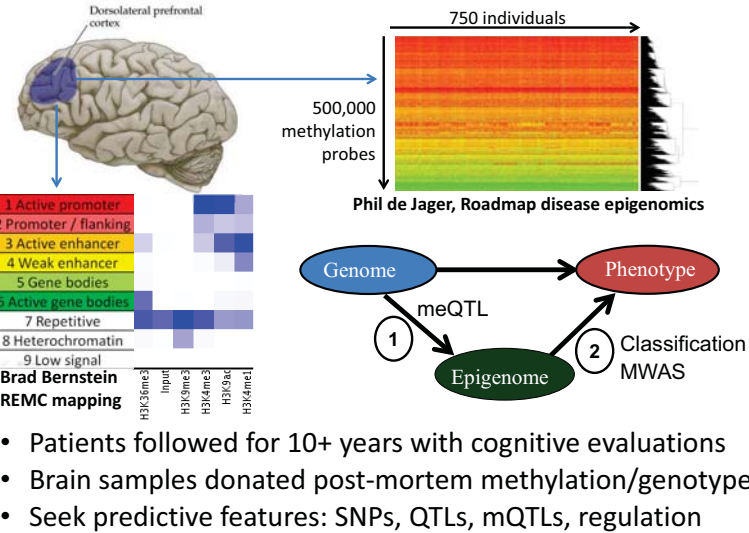
- Chromatin states: Enhancers, promoters, motifs
- Enrichment in individual loci, across 1000s of SNPs in T1D



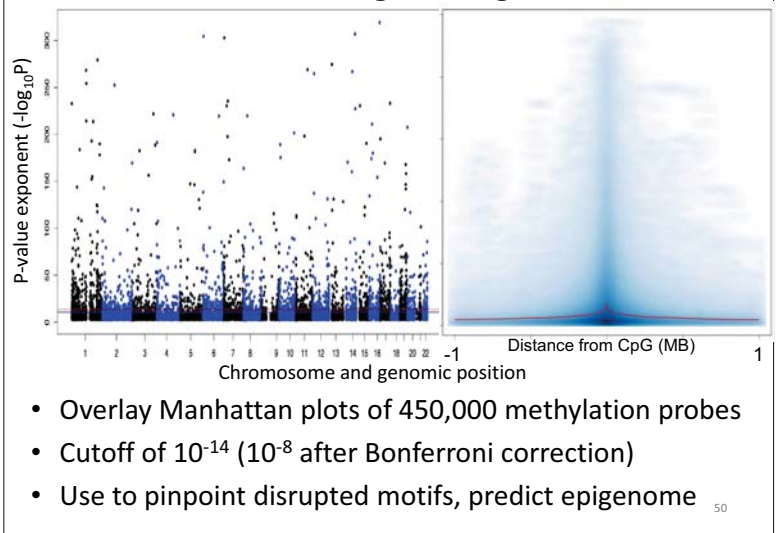
(2) Epigenome changes in disease

- Intermediate molecular phenotypes associated with disease
- Variation in brain methylomes of Alzheimer's patients

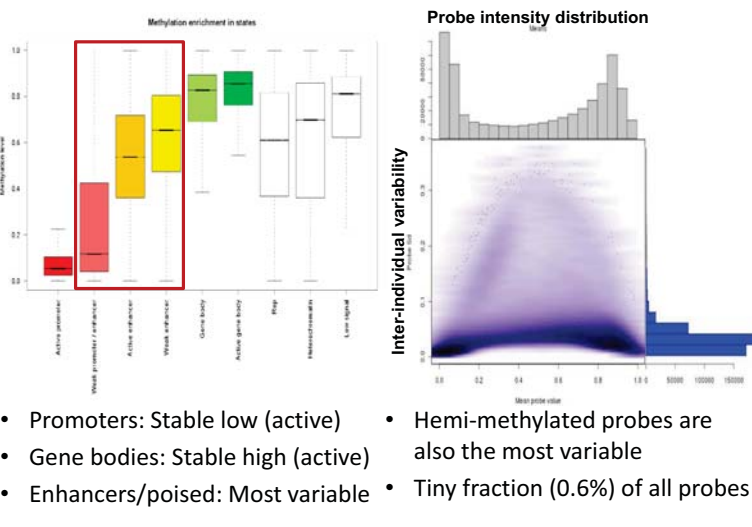
Phil de Jager: Methylation in 750 Alzheimer patients



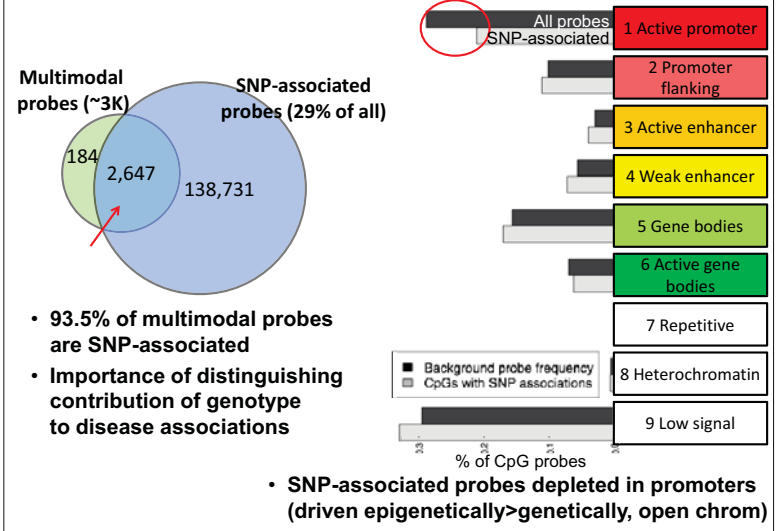
2,500 mQTLs for neighboring SNPs at 10^{-14}



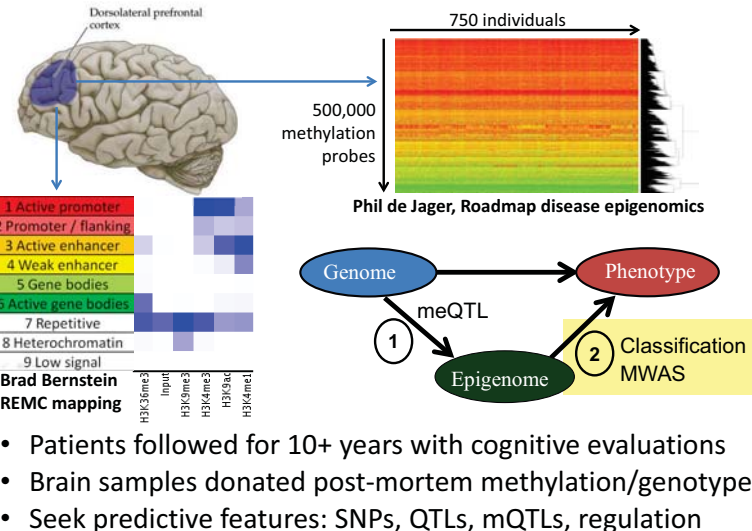
Focusing on 2831 most variable probes



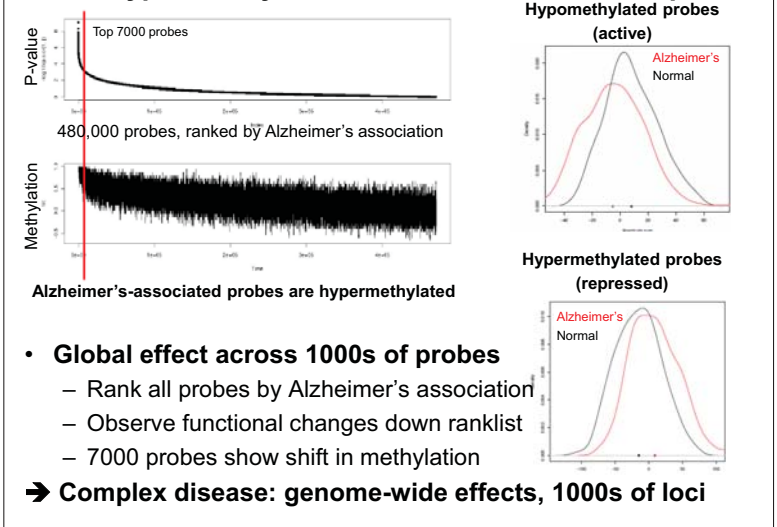
Multimodal ↔ SNP-associated ↔ Promoter-depleted



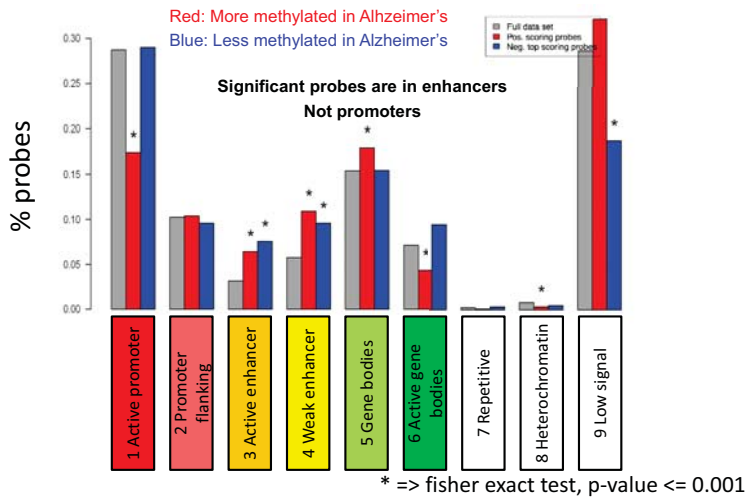
Phil de Jager: Methylation in 750 Alzheimer patients



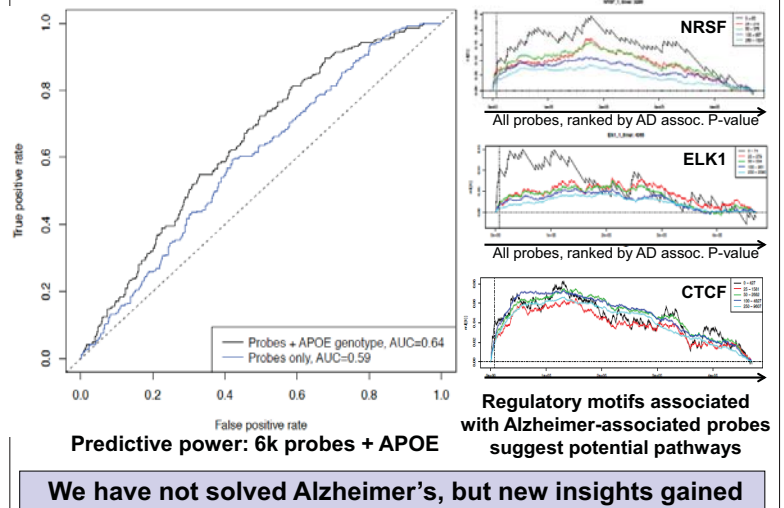
Global hyper-methylation trend in AD-associated probes



Chromatin state breakdown reveals ↓ activity



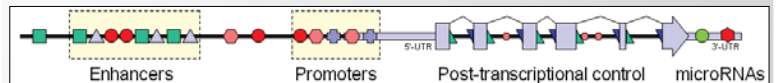
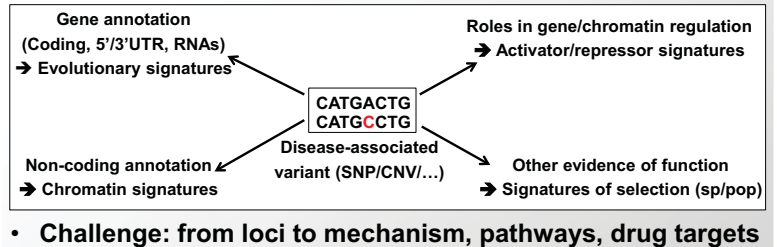
Alzheimer's prediction vs. likely biological pathways



Systems-level views of disease epigenomics

- Evolutionary signatures → gene/genome annotation**
 - High-resolution annotation: genes, RNAs, motif instances
 - Measuring selection within the human population
- Chromatin states for interpreting disease association**
 - Annotate dynamic regulatory elements in multiple cell types
 - Activity-based linking of regulators → enhancers → targets
- Interpreting disease-associated sequence variants**
 - Mechanistic predictions for individual top-scoring SNPs
 - Functional roles of 1000s of disease-associated SNPs
- Systematic manipulation of 2000+ human enhancers**
 - Test effect of single-motif and single-nucleotide disruptions
 - Role of activator/repressor motifs, disease-associated SNPs
- Personal genomes/epigenomes in health and disease**
 - Allele-specific activity. Alzheimer's ↔ brain methylation ↔ SNP
 - Global repression of distal enhancers. NRSF, ELK1, CTCF

Understanding human variation and human disease



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

Collaborators and Acknowledgements

- ENCODE**
 - Brad Bernstein, Tarjei Mikkelsen, Noam Shores, David Epstein
- Massively parallel enhancer reporter assays**
 - Tarjei Mikkelsen, Broad Institute
- Epigenome Roadmap**
 - Bing Ren, Brad Bernstein, John Stam, Joe Costello
- 2X mammals**
 - Kerstin Lindblad-Toh, Eric Lander, Manuel Garber, Or Zuk
- Funding**
 - NHGRI, NIH, NSF, Sloan Foundation

MIT Computational Biology group

Compbio.mit.edu

