

Characterizing the effect of genetic variants within promoters with distal enhancer functions

Alejandra Medina Rivera, PhD

Laboratorio Internacional de Investigación en Genoma Humano

Universidad Nacional Autónoma de México

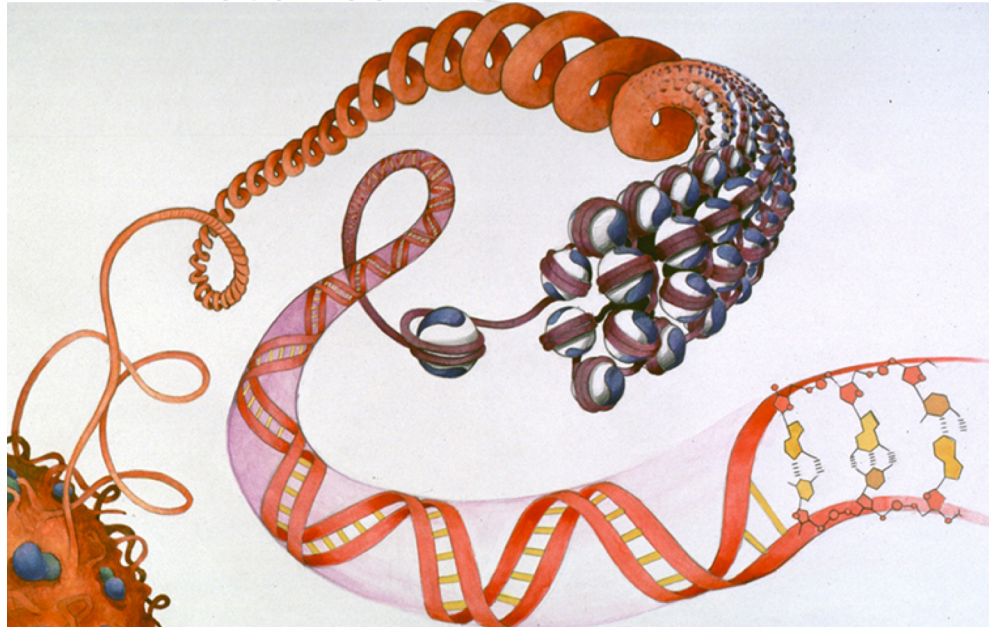
Oaxaca, November 2018



Regulatory Genome

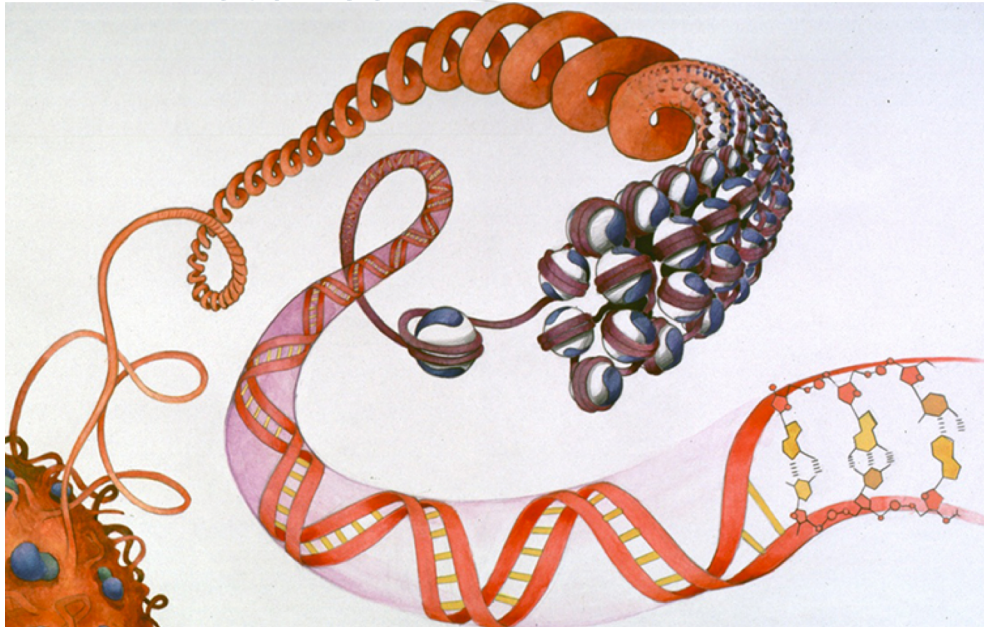
Regulatory Genome

Histones

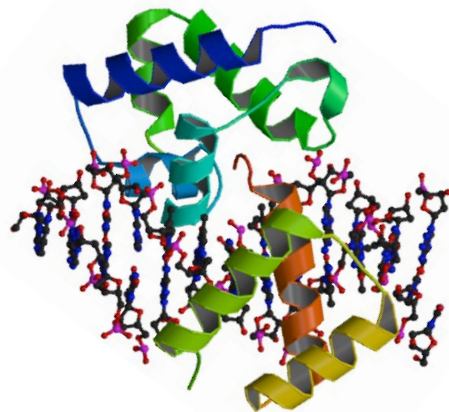


Regulatory Genome

Histones

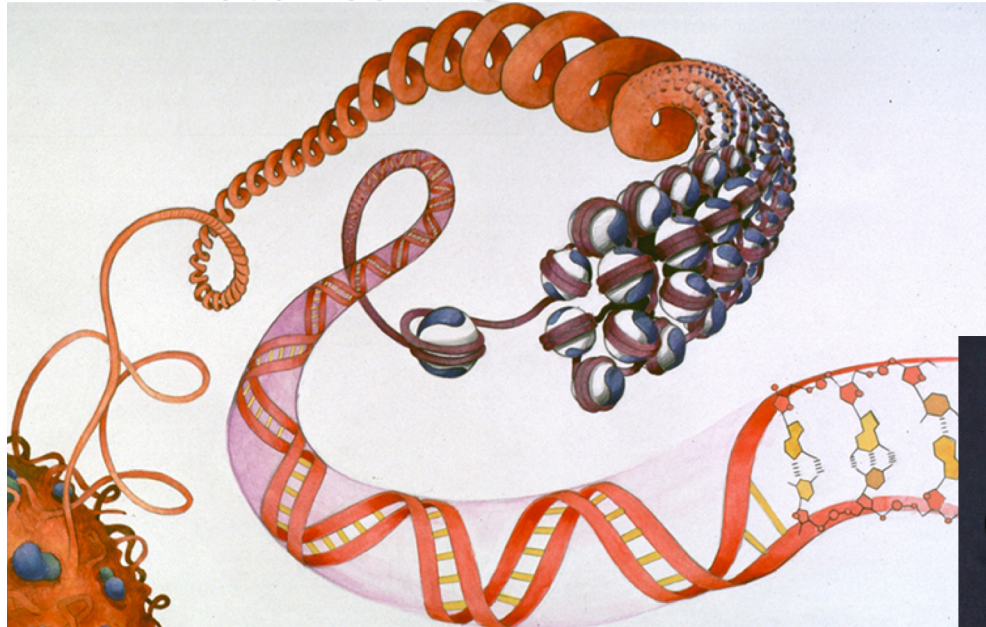


Transcriptional Factors

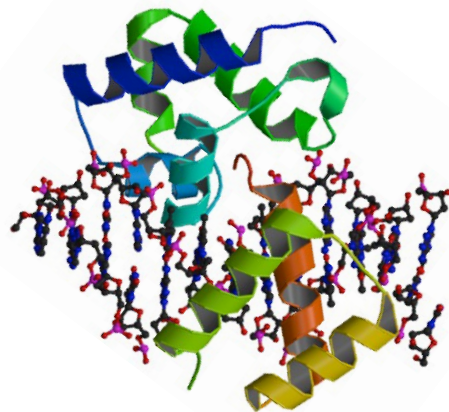


Regulatory Genome

Histones

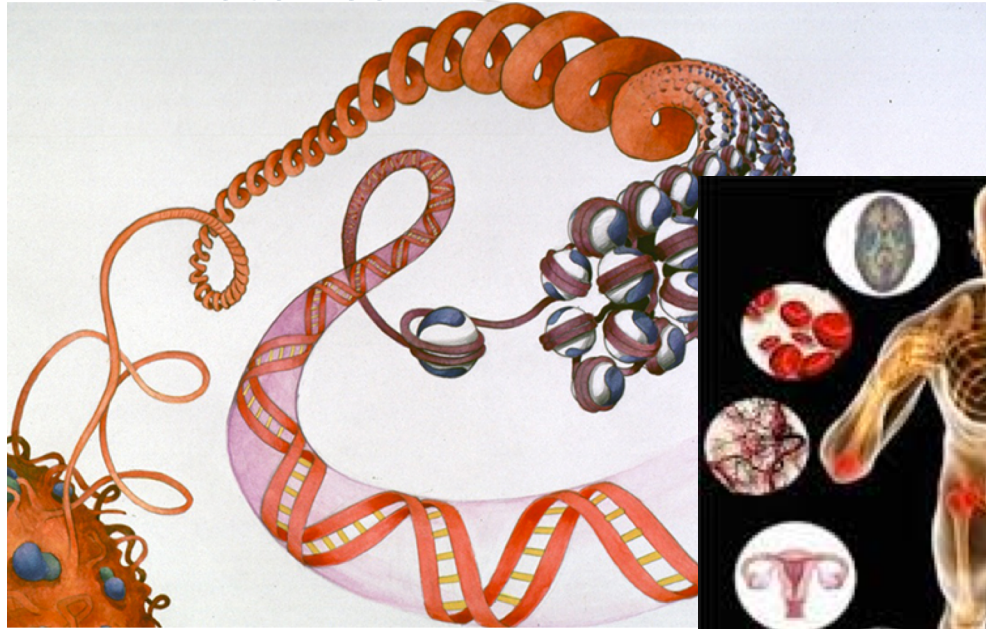


Transcriptional Factors

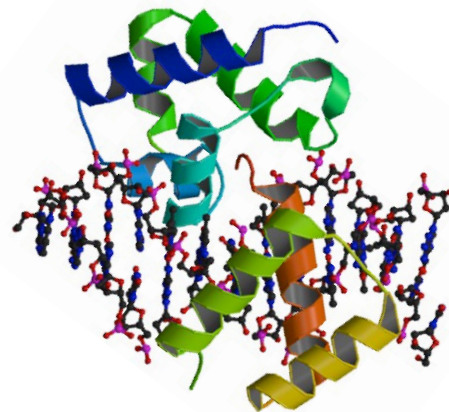


Regulatory Genome

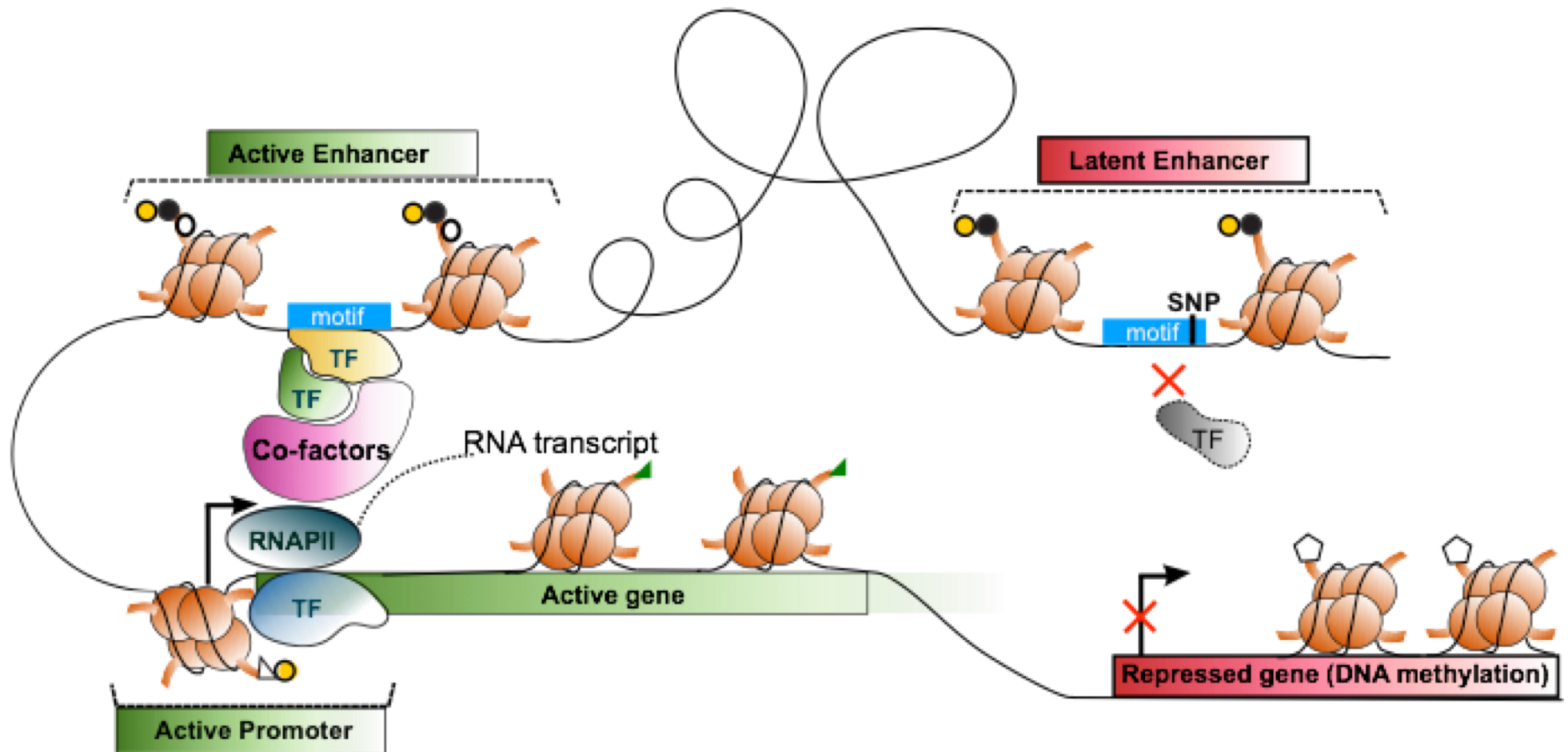
Histones



Transcriptional Factors



Annotating the Regulatory Genome



● H3K4me1

△ H3K4me3

▲ H3K36me3

● H3K4me2

○ H3K27Ac

◻ H3K9me3

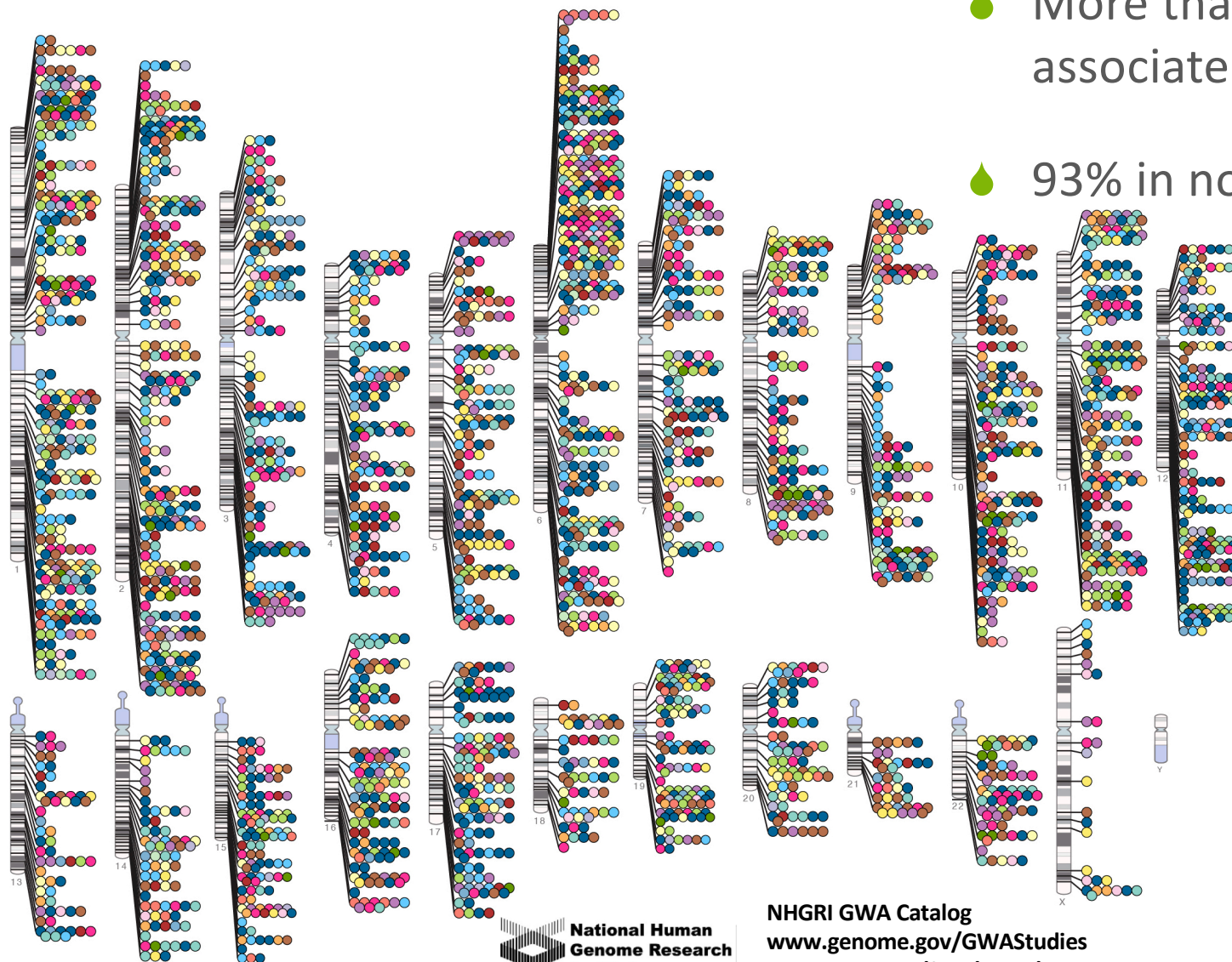


Published Genome-Wide Associations since 12/2012

Published GWA at $p \leq 5 \times 10^{-8}$ for 17 trait categories

More than 17,000 reported associated SNPs

93% in non-coding sequence



- Digestive system disease
- Cardiovascular disease
- Metabolic disease
- Immune system disease
- Nervous system disease
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurement
- Other measurement
- Response to drug
- Biological process
- Cancer
- Other disease
- Other trait

First characterized Enhancers

Gene	Origin	Size (bp)	Distance from TSS	
<i>Early gene</i>	SV40	196	~200	Banerji, J. et al. (1981); Benoist, C. and Chambon, P. (1981)
<i>Early gene</i>	Cytomegalovirus	406	-524 to -118	Boshart, M. et al. (1985)
<i>Hsp70</i>	<i>Xenopus</i>	160	-260 to -100	Bienz, M. and Pelham, H.R. (1986)
<i>Fos</i>	Human	340	-404 to -64	Deschamps, J. et al. (1985)
<i>hMT-IIA</i>	Human	327	-366 to -39	Serfling, E. et al. (1985)
<i>Mmt-IA</i>	Mouse	114	-187 to -73	Serfling, E. et al. (1985)
		155	-194 to -39	
<i>H2A</i>	Urchin	28	-139 to -111	Grosschedl, R. and Birnstiel, M.L. (1982)
<i>IFNβ</i>	Human	40	-77 to -37	Goodbourn, S. et al. (1985)

Regulatory Regions: Enhancers vs Promoters

Gene	Origin	Size (bp)	Distance from TSS	
<i>Early gene</i>	SV40	196	~200	Banerji, J. et al. (1981); Benoist, C. and Chambon, P. (1981)
<i>Early gene</i>	Cytomegalovirus	406	-524 to -118	Boshart, M. et al. (1985)
<i>Hsp70</i>	<i>Xenopus</i>	160	-260 to -100	Bienz, M. and Pelham, H.R. (1986)
<i>Fos</i>	Human	340	-404 to -64	Deschamps, J. et al. (1985)
<i>hMT-IIA</i>	Human	327	-366 to -39	Serfling, E. et al. (1985)
<i>Mmt-IA</i>	Mouse	114 155	-187 to -73 -194 to -39	Serfling, E. et al. (1985)
<i>H2A</i>	Urchin	28	-139 to -111	Grosschedl, R. and Birnstiel, M.L. (1982)
<i>IFNβ</i>	Human	40	-77 to -37	Goodbourn, S. et al. (1985)

How to classify Regulatory Regions?

Features (active elements)	Promoter	Enhancer
Intrinsic property	Induce transcription of a heterologous reporter gene	Activate a distal (heterologous) promoter
Transcription initiation	Unidirectional or divergent	Mainly divergent
Ratio between sense and antisense transcripts	Biased toward sense transcription	Equilibrated
Transcription elongation	Produce long polyadenylated transcripts	Some enhancers can produce low levels of polyadenylated transcripts
Histone modifications	H3K27ac (H3K4me1 < H3K4me3)	H3K27ac (H3K4me1 > H3K4me3)
RNAPII and GTF	Present	Present
GpG islands	Majority	Very rare

How to assess the enhancer function of promoters?

Major challenge: How to assess the enhancer function of promoters?



Salvatore Spicuglia



Detection of long-range regulatory regions

Enhancers

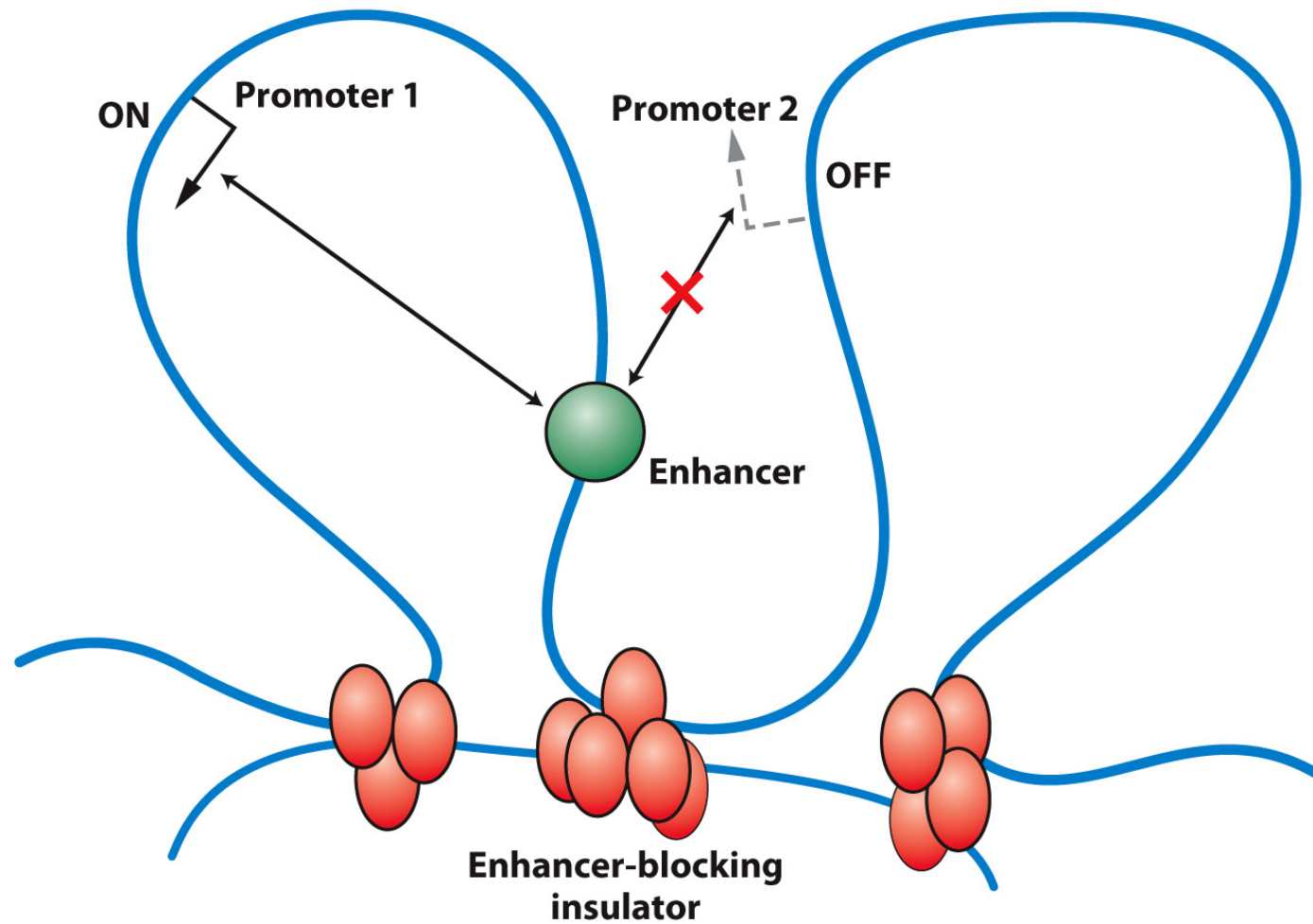
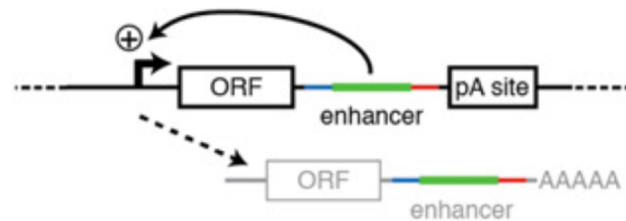


Figure 12-21
Introduction to Genetic Analysis, Eleventh Edition
© 2015 W. H. Freeman and Company

Detection of long-range regulatory regions Enhancers

- Starr-seq: High-throughput assessment of sequence enhancer potential

Arnold, C.D. et al. (2013). *Science* (80-.),,

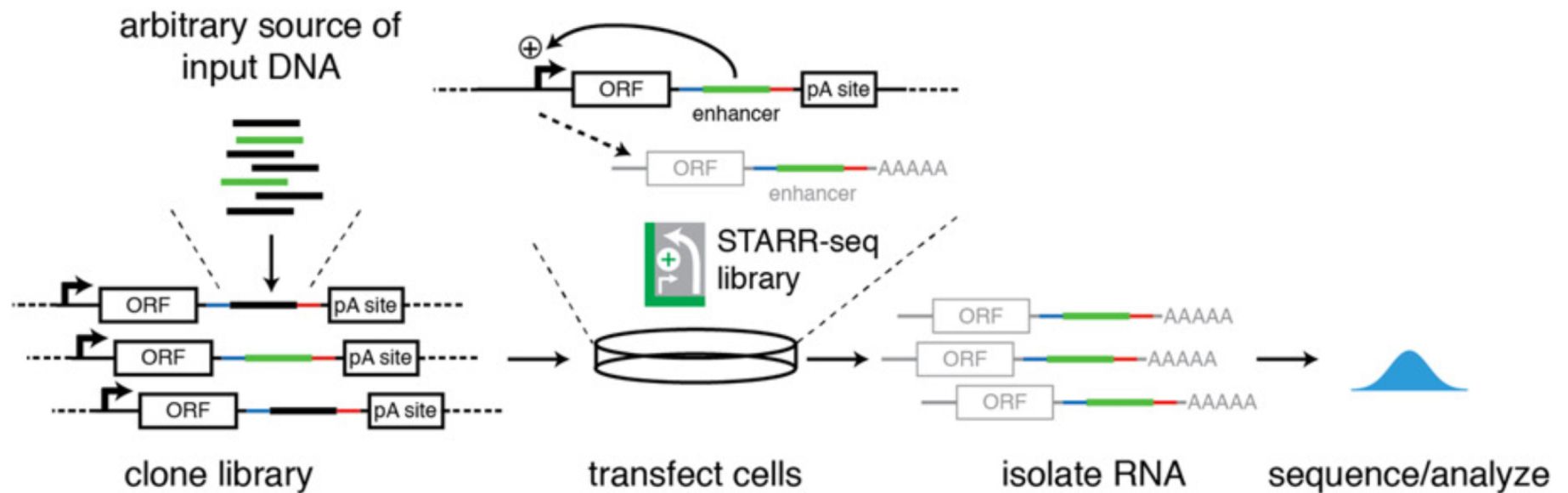


Muerdter, F. et al. (2015). *106*, 145–150

Detection of long-range regulatory regions Enhancers

- Starr-seq: High-throughput assessment of sequence enhancer potential

Arnold, C.D. et al. (2013). Science (80-.),

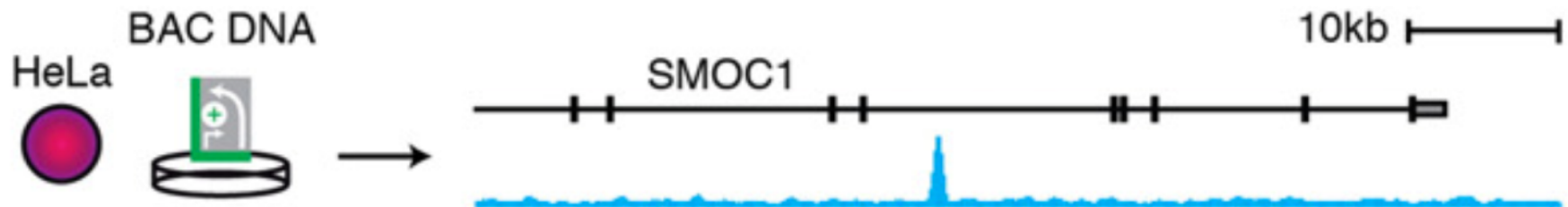


Muerdter, F. et al. (2015).106, 145–150

Detection of long-range regulatory regions Enhancers

- Starr-seq: High-throughput assessment of sequence enhancer potential

Arnold, C.D. et al. (2013). *Science* (80-.),

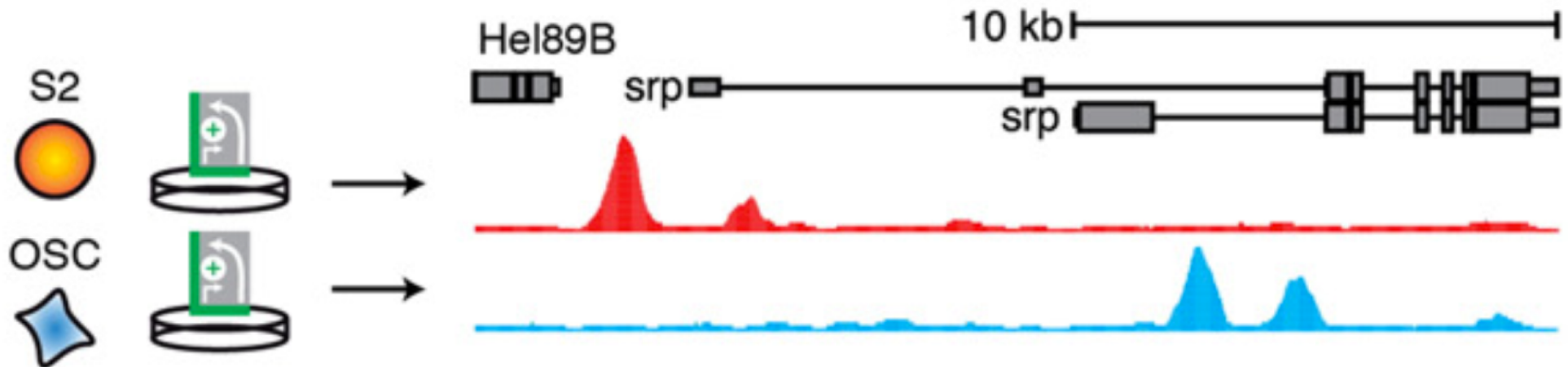


Muerdter, F. et al. (2015). *106*, 145–150

Detection of long-range regulatory regions Enhancers

- Starr-seq: High-throughput assessment of sequence enhancer potential

Arnold, C.D. et al. (2013). *Science* (80-.),



Muerdter, F. et al. (2015). *106*, 145–150

Major challenge: How to assess the enhancer function of promoters?

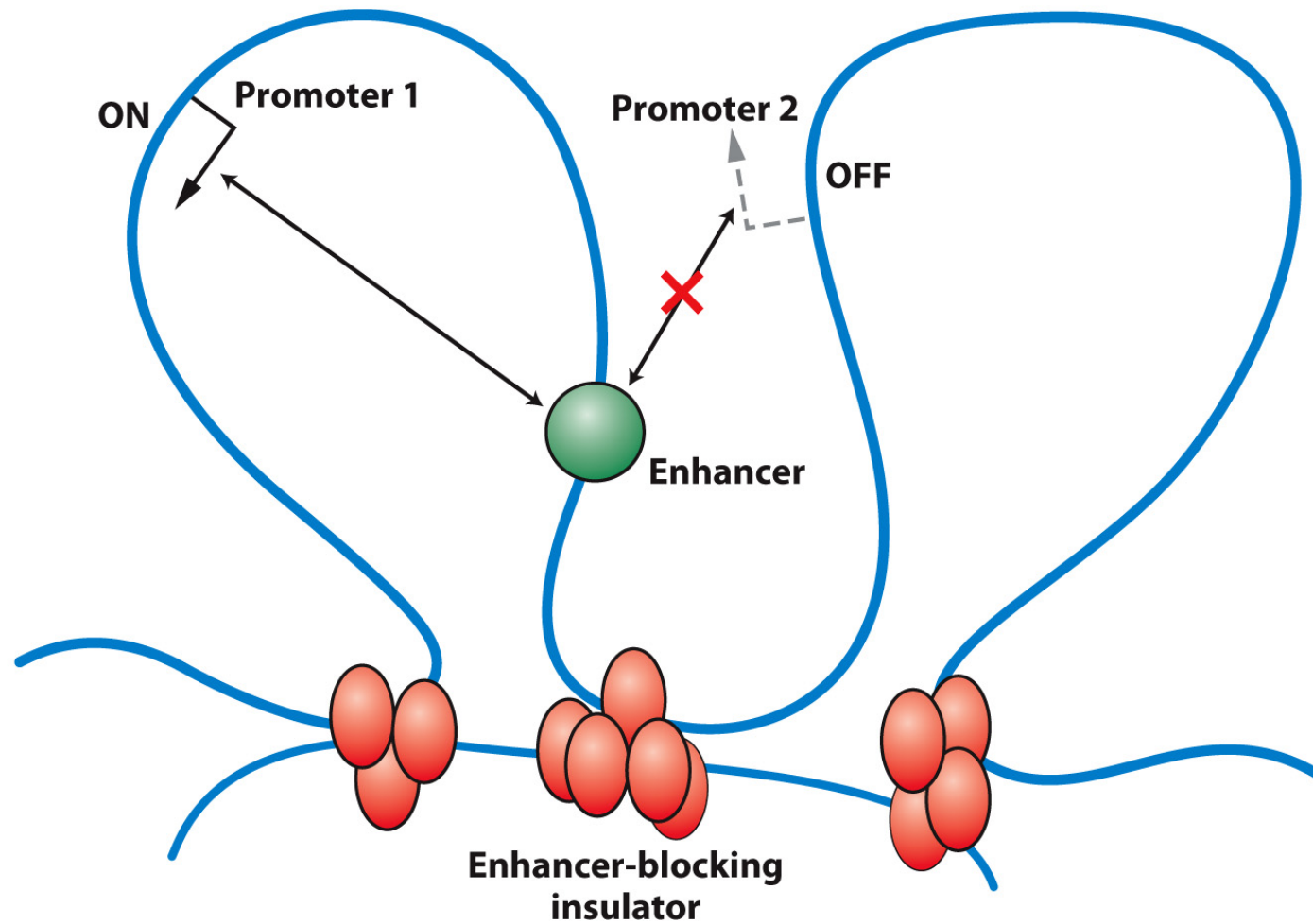


Figure 12-21
Introduction to Genetic Analysis, Eleventh Edition
© 2015 W. H. Freeman and Company

Major challenge: How to assess the enhancer function of promoters?

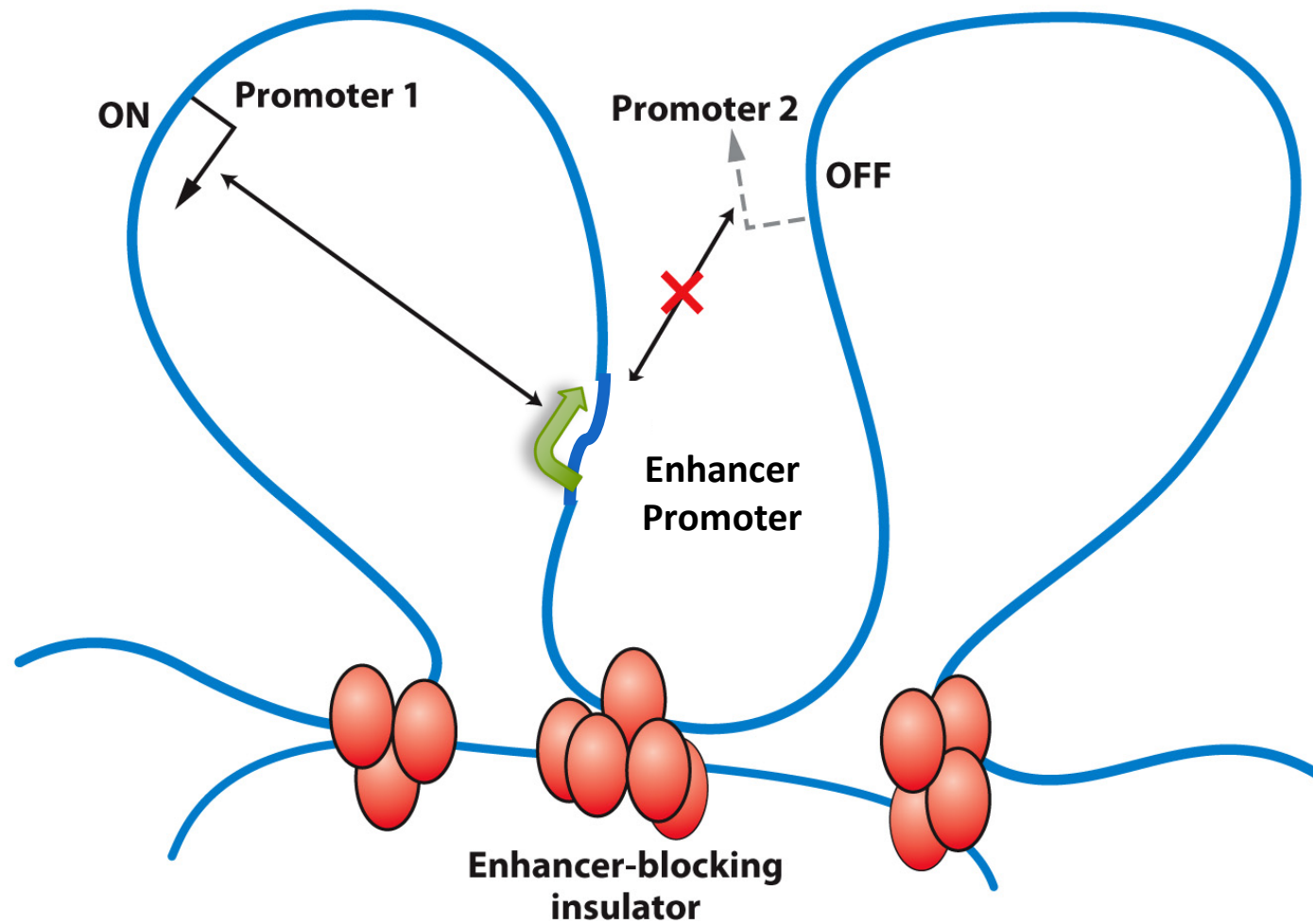
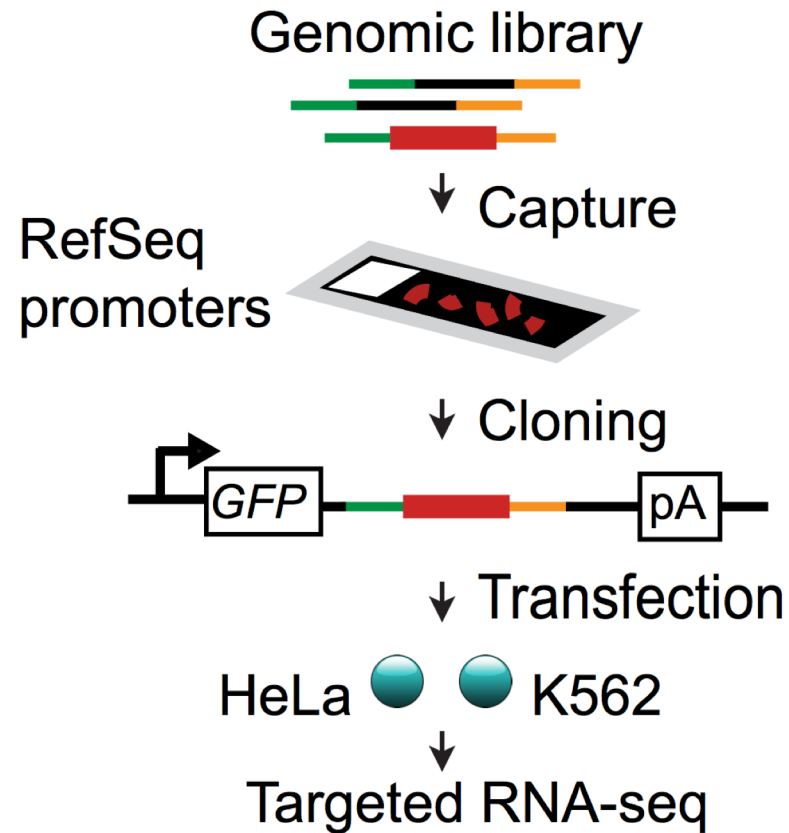


Figure 12-21
Introduction to Genetic Analysis, Eleventh Edition
© 2015 W. H. Freeman and Company

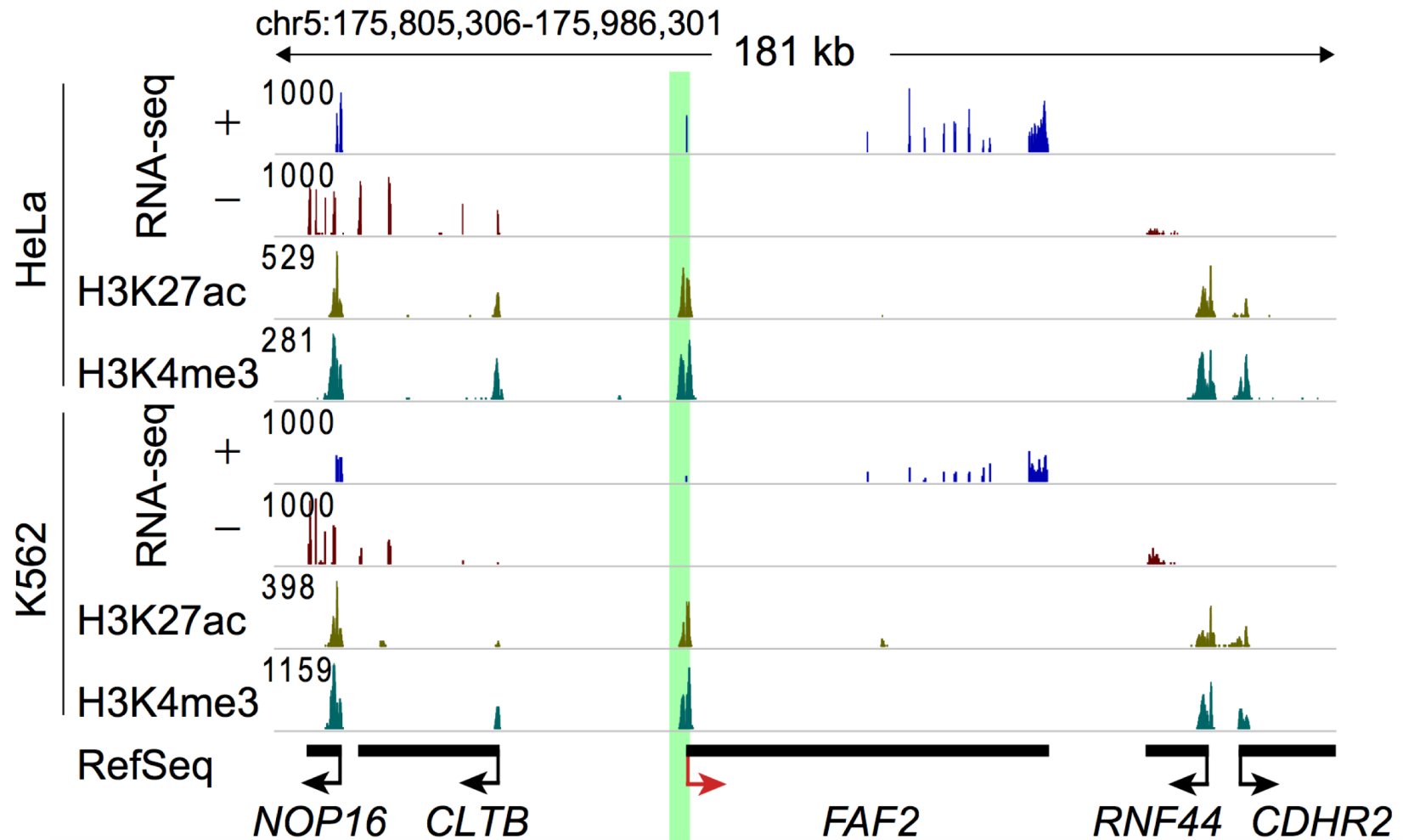
Can promoters work as enhancers?

CapStarr-seq

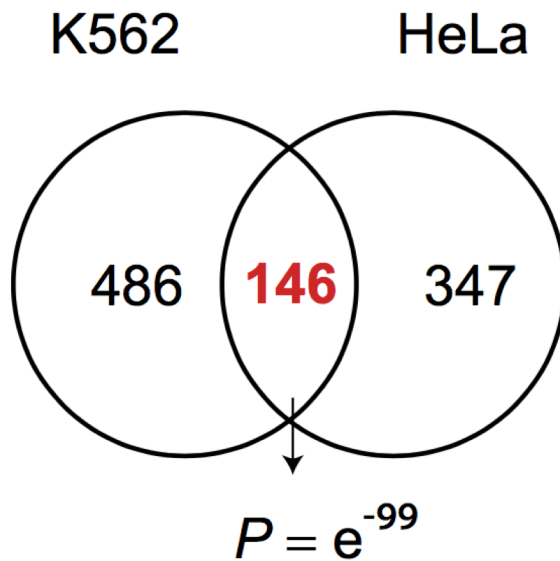


Salvatore Spicuglia

Can promoters work as enhancers?



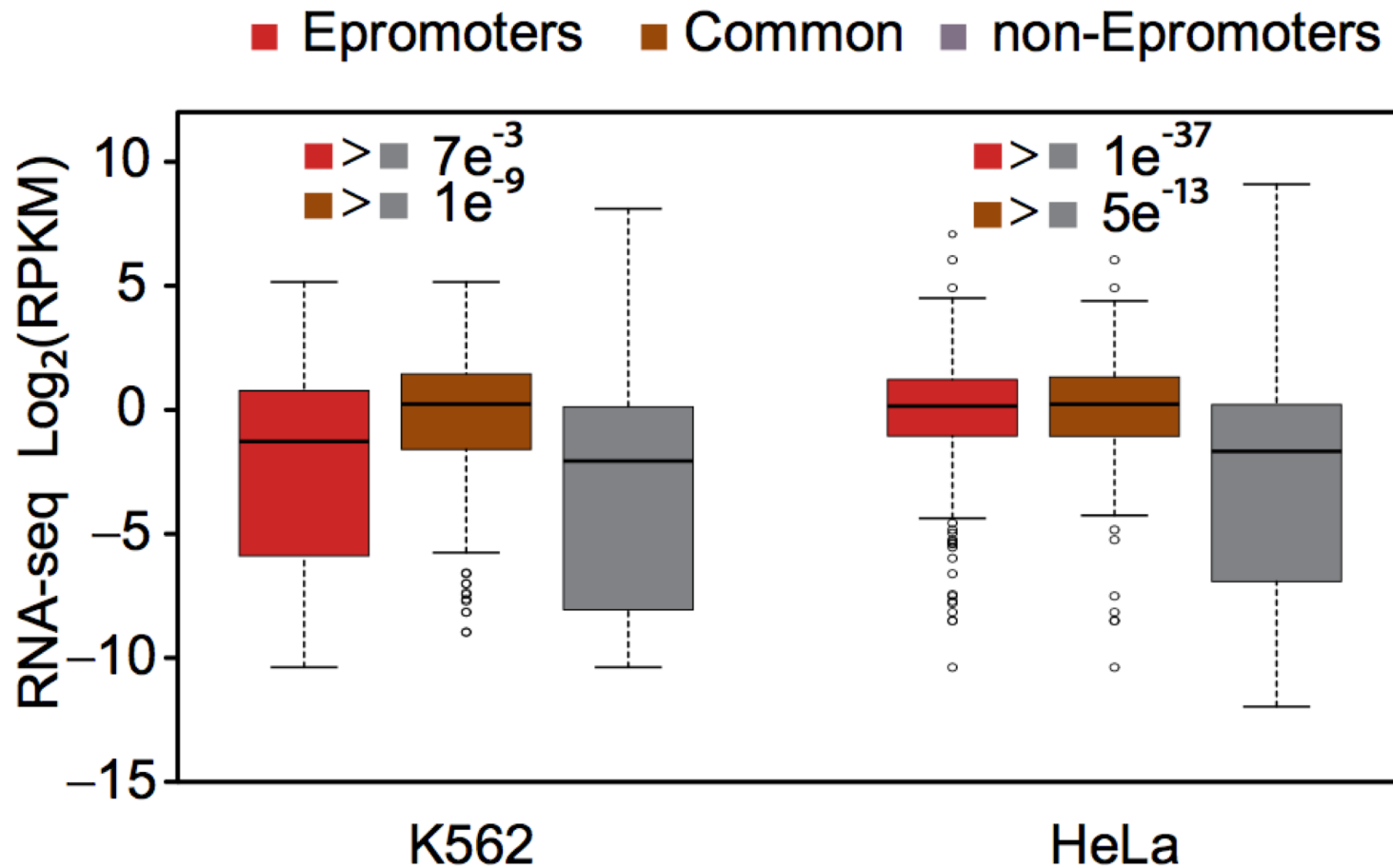
How many promoters work also as enhancers?



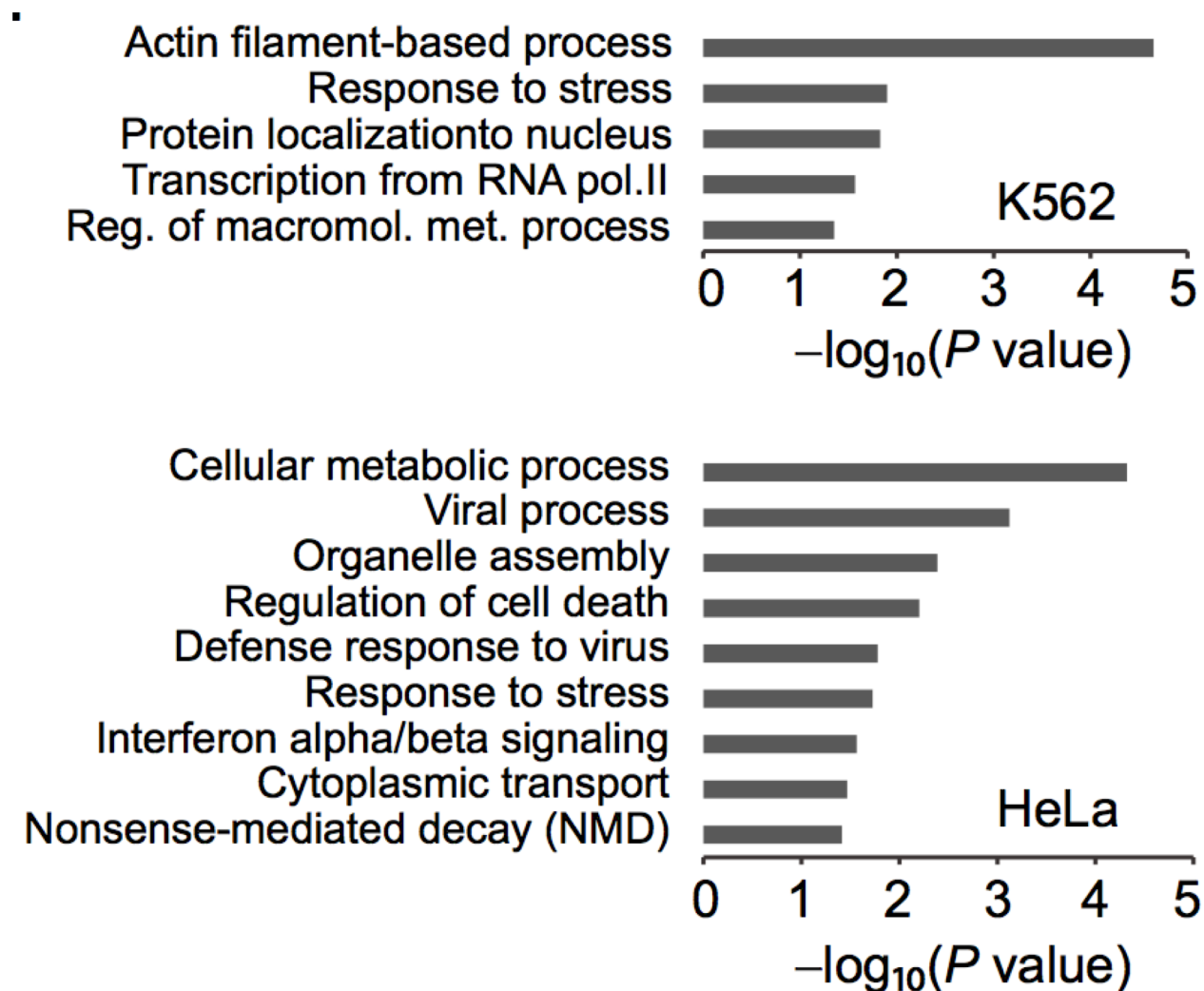
Total of ePromoters:

- K562= 632 (3%)
- HELA= 493 (2.37%)
- Total analyzed promoters= 20,719 (100%)

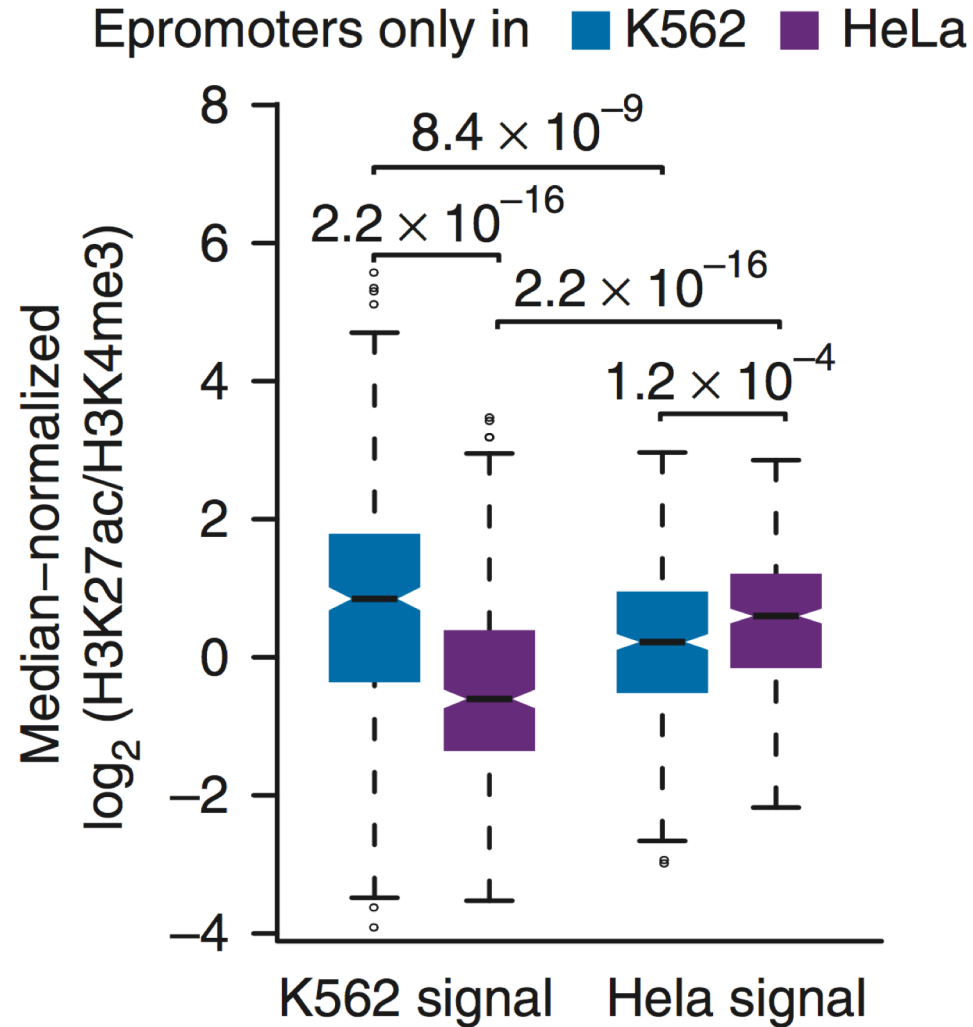
RNA expression of genes downstream promoters



Functional annotation of genes downstream Epromoters



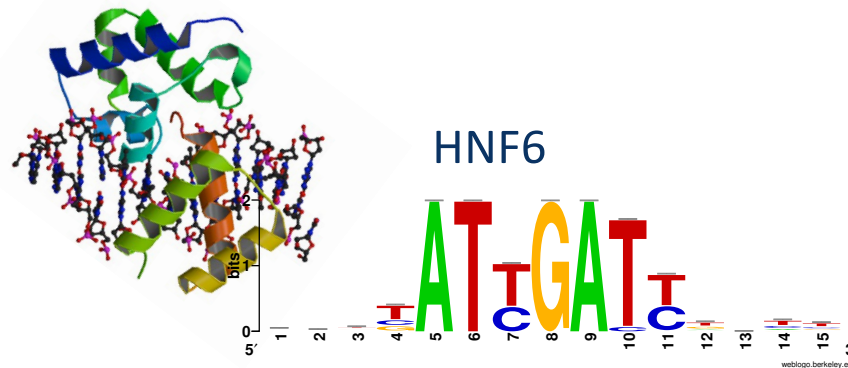
Epromoters had a higher H3K27ac/H3K4me3 ratio in the cell type where they were found to be active



Which are the regulatory mechanisms affecting Epromoters?

Transcription factors control gene regulation by binding to specific DNA sequences

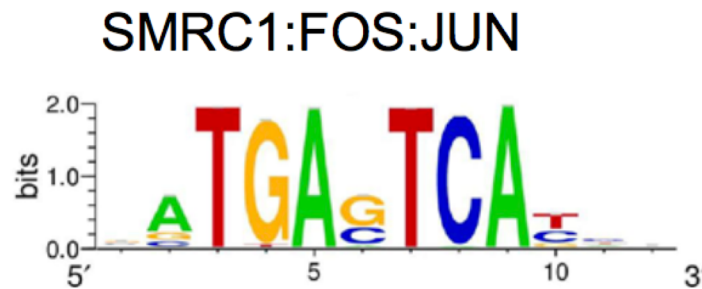
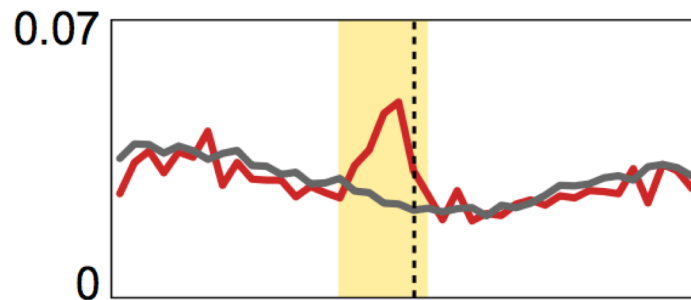
- Transcription Factors interact with DNA binding to sequence specific sites.



Transcription Factors enriched motifs in Epromoters

K562

— Epromoters
— non-Epromoters



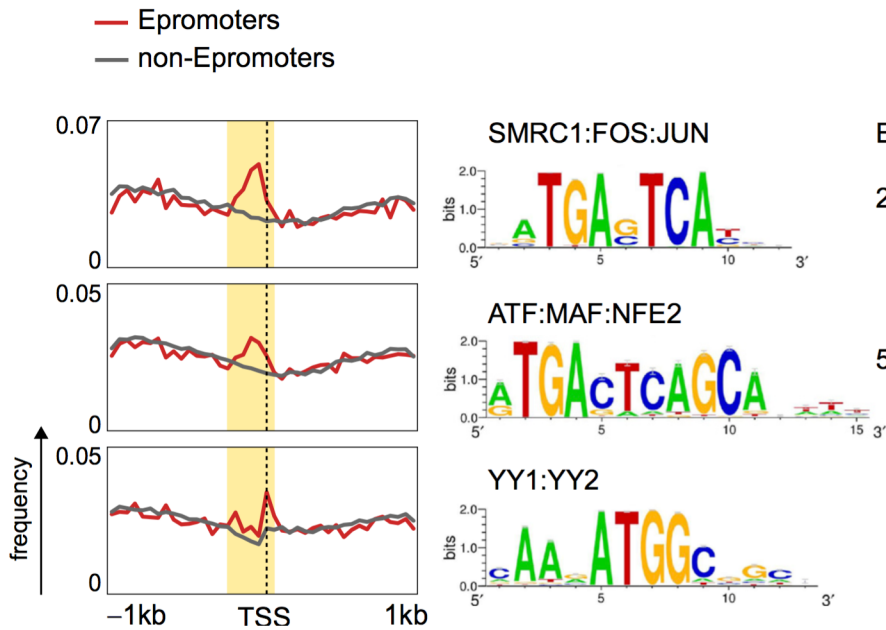
E-value
2.33E-17



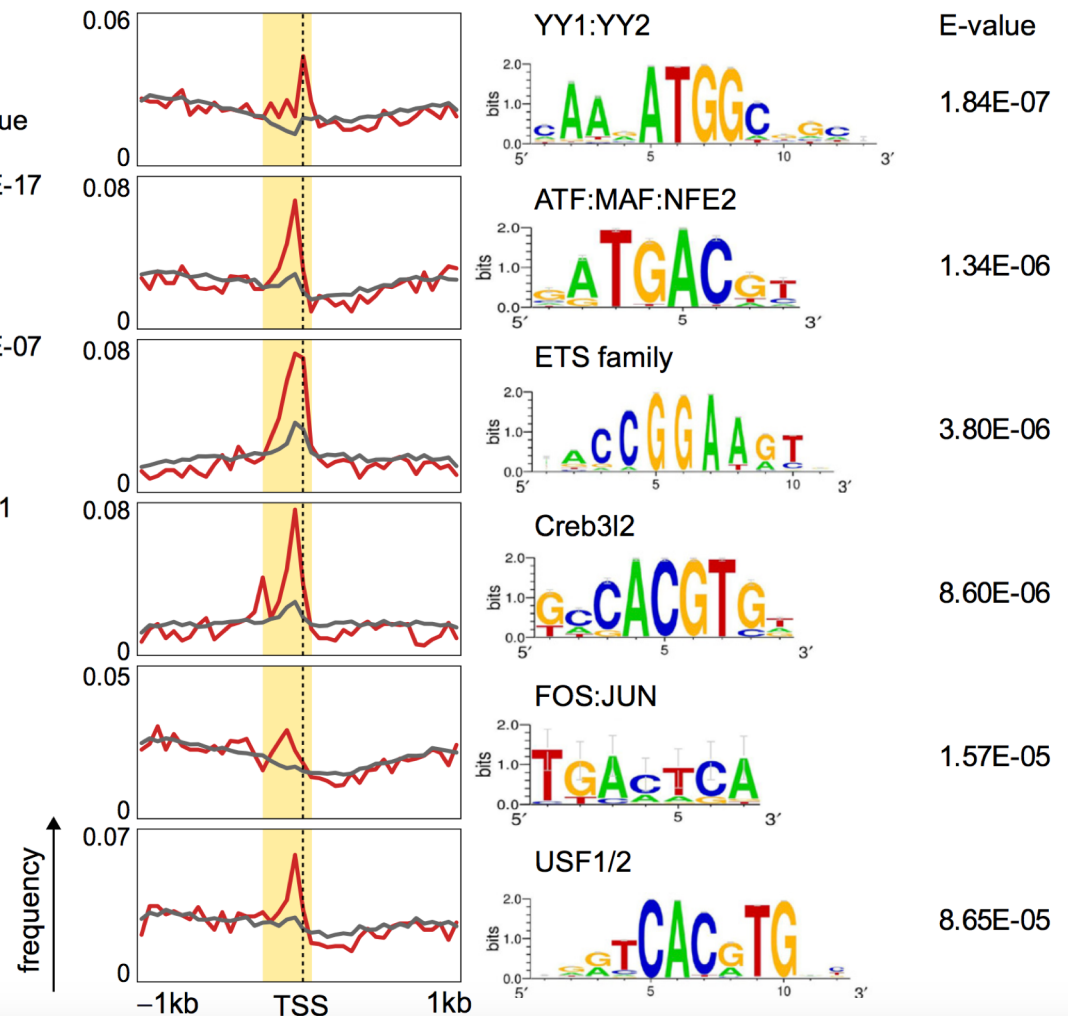
- *position-scan*: Motif enrichment in a given position of a set of sequences

Transcription Factors enriched motifs in Epromoters

K562

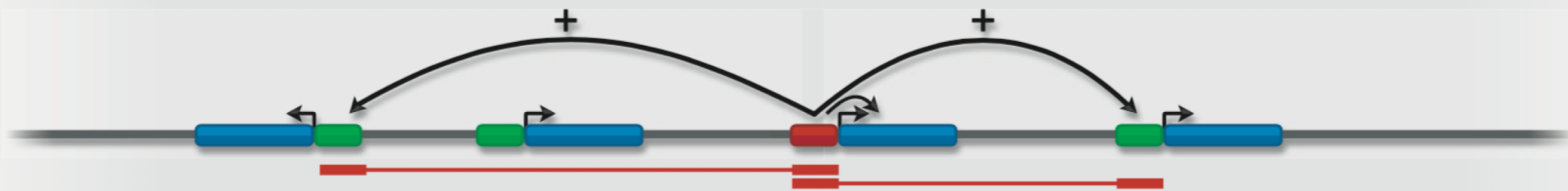


HeLA

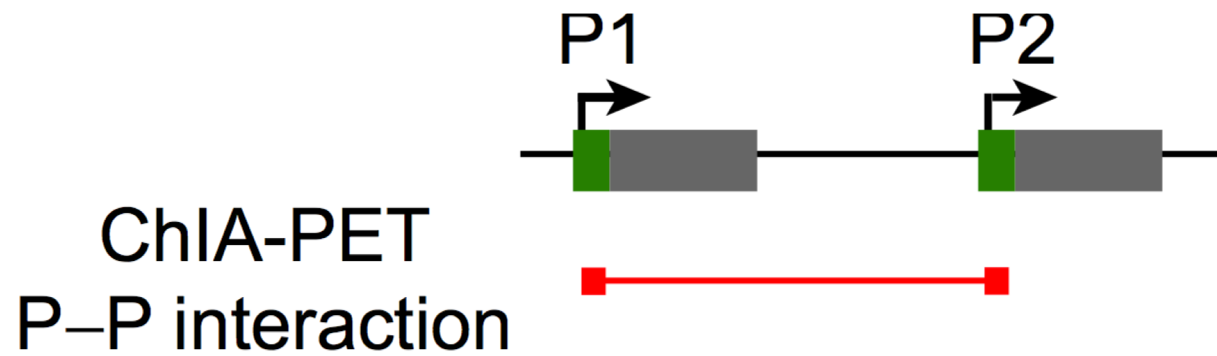


- *position-scan*: Motif enrichment in a given position of a set of sequences

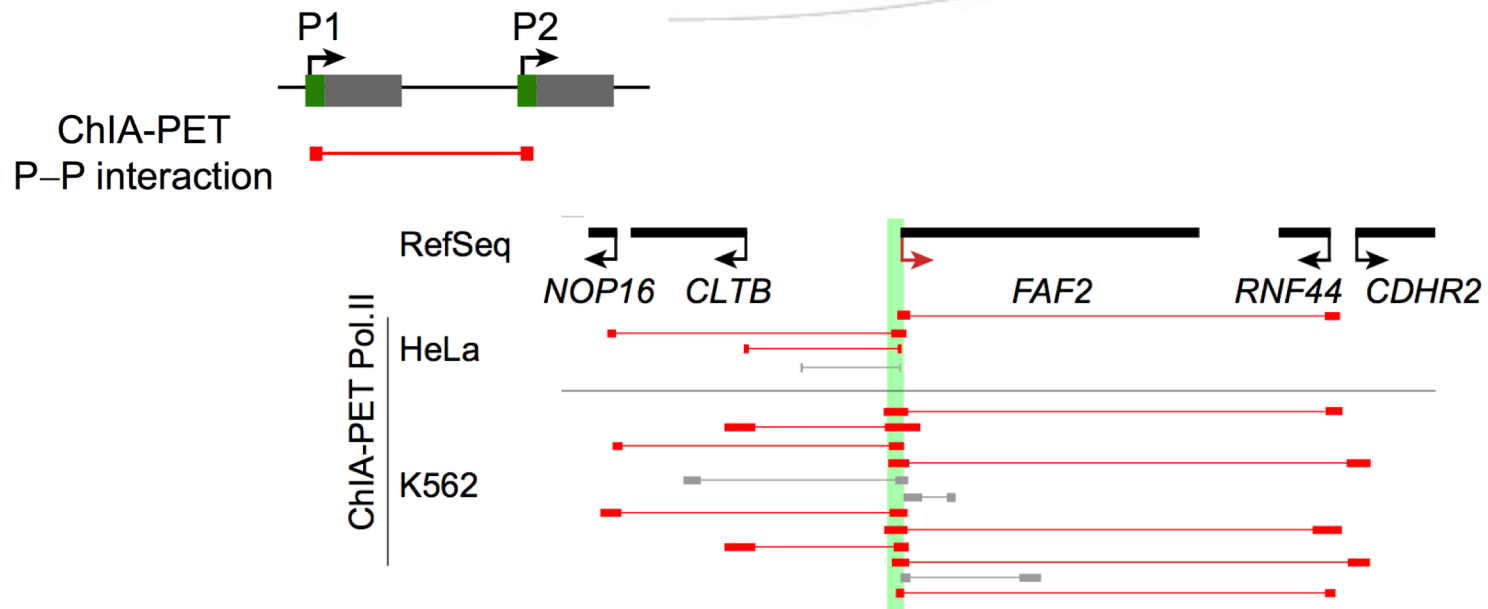
Which are the genes regulated by Epromoters?



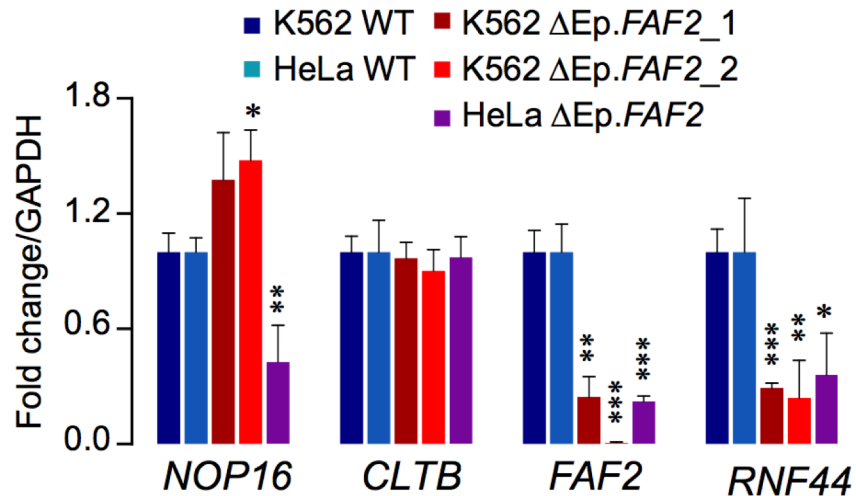
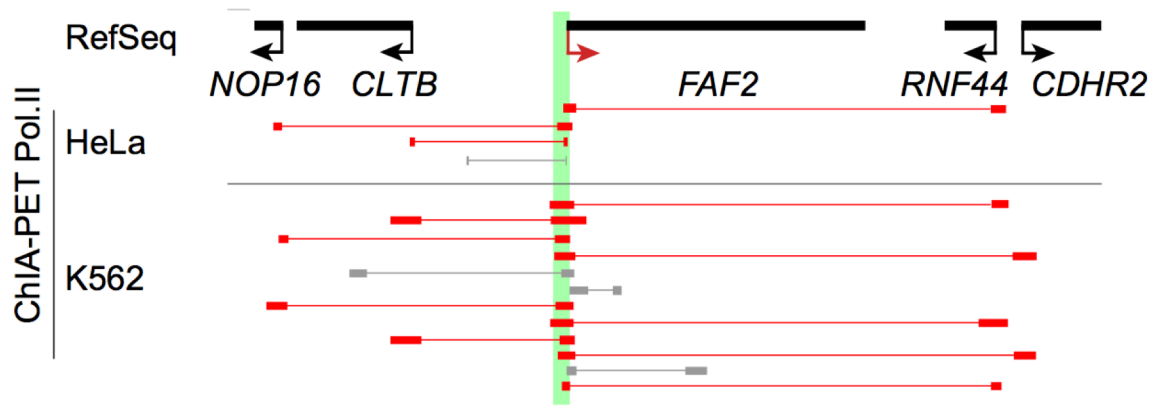
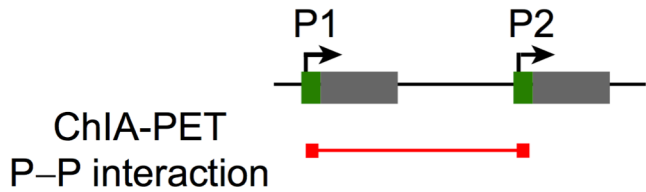
Which are the genes regulated by Epromoters?



Which are the genes regulated by Epromoters?



Which are the genes regulated by Epromoters?

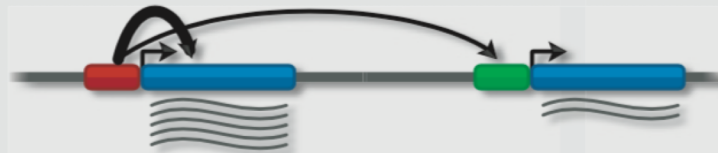


What is the effect of Epromoters on long range regulated genes?



Cell type A

Inverse correlation ?



Positive correlation ?

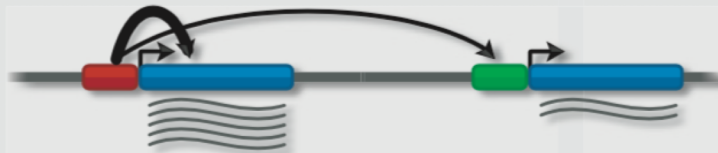


What is the effect of Epromoters on long range regulated genes?



Cell type A

Inverse correlation ?



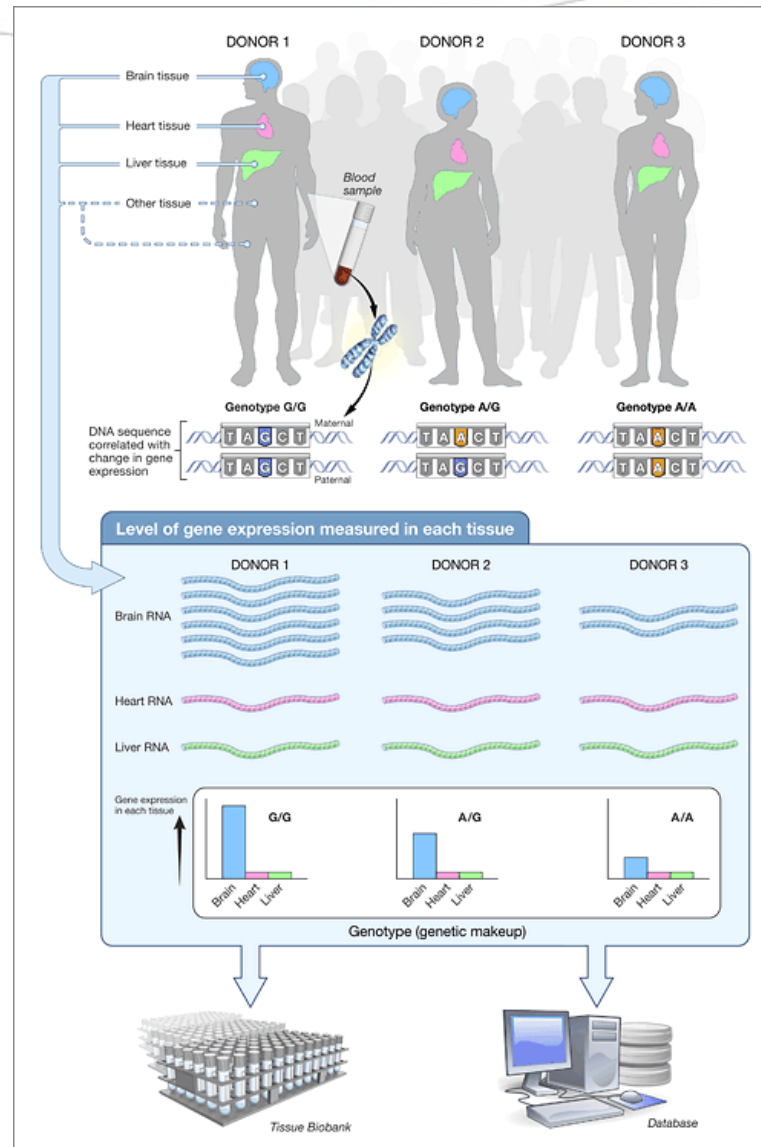
Positive correlation ?



Cell type B

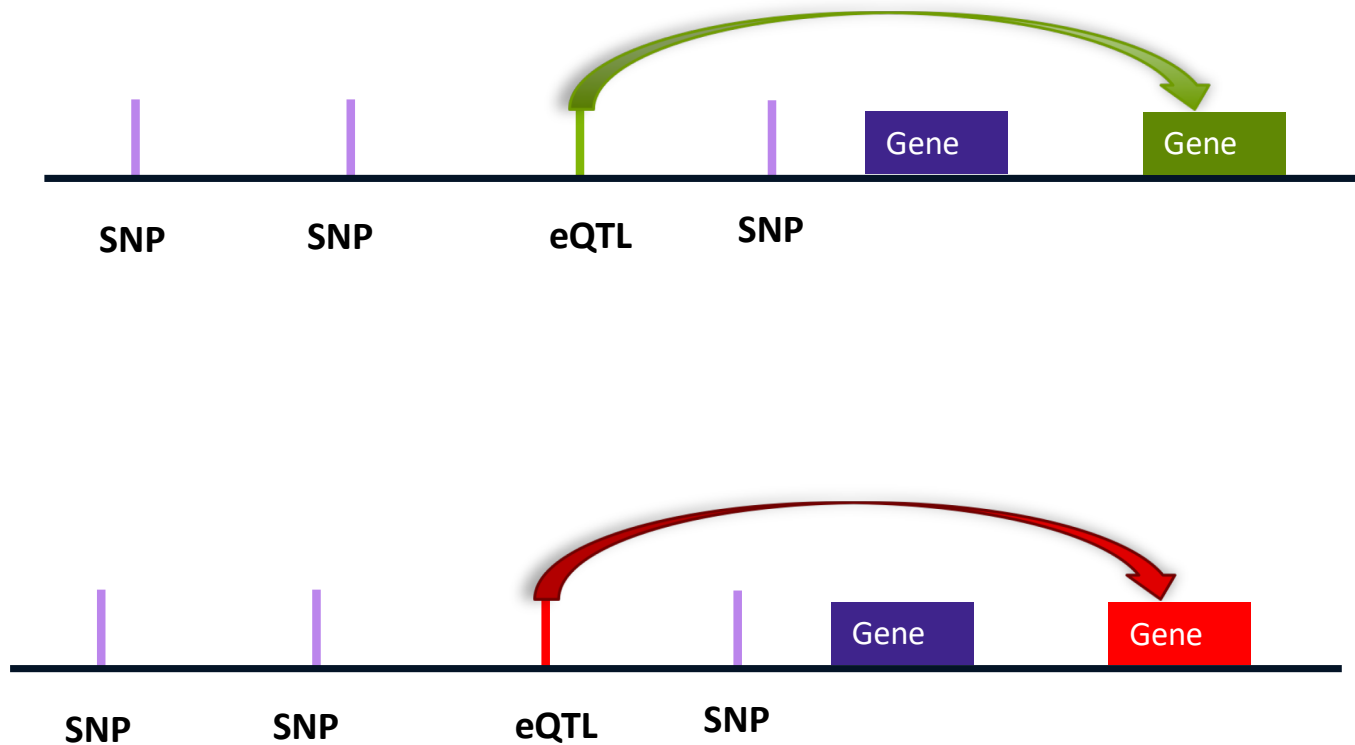


Genotype Tissue Expression Project

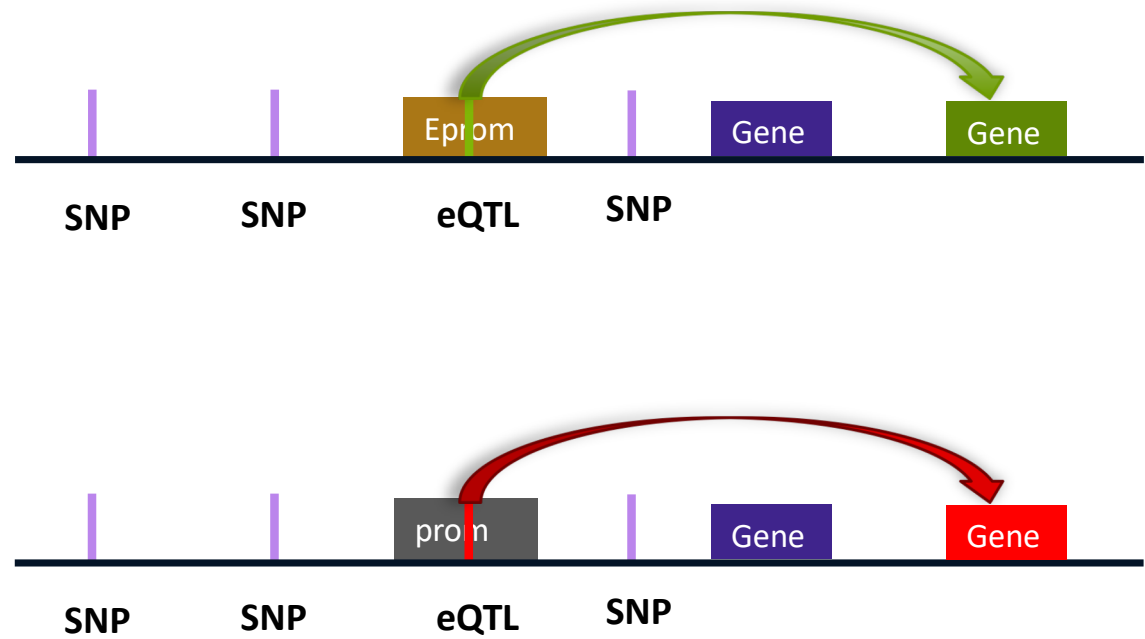


eQTLs can have negative or positive correlation to gene expression

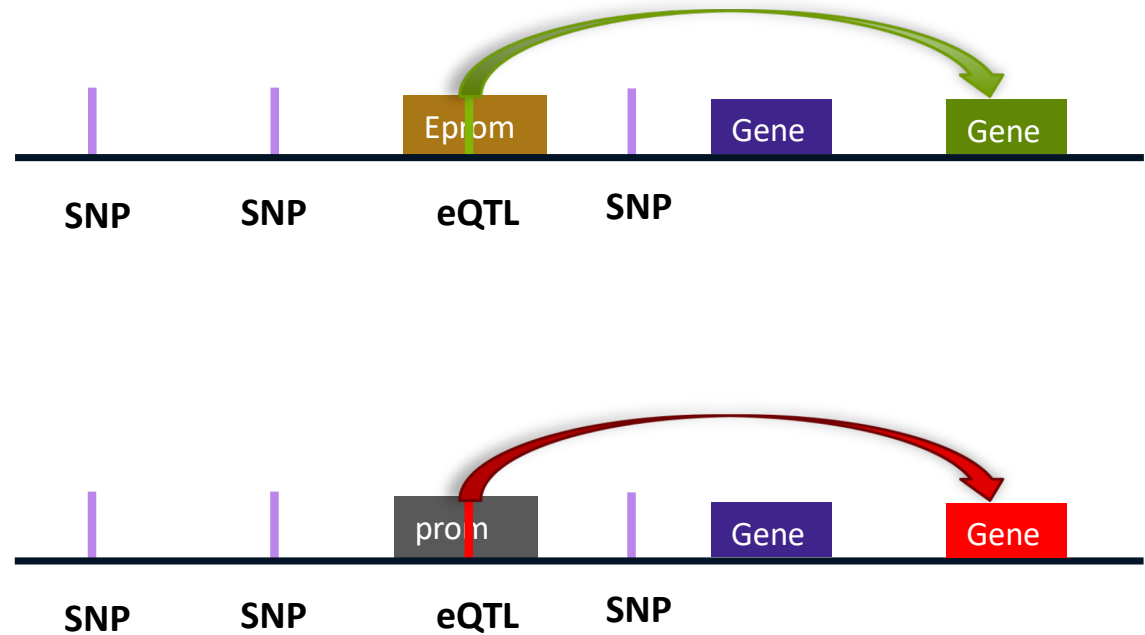
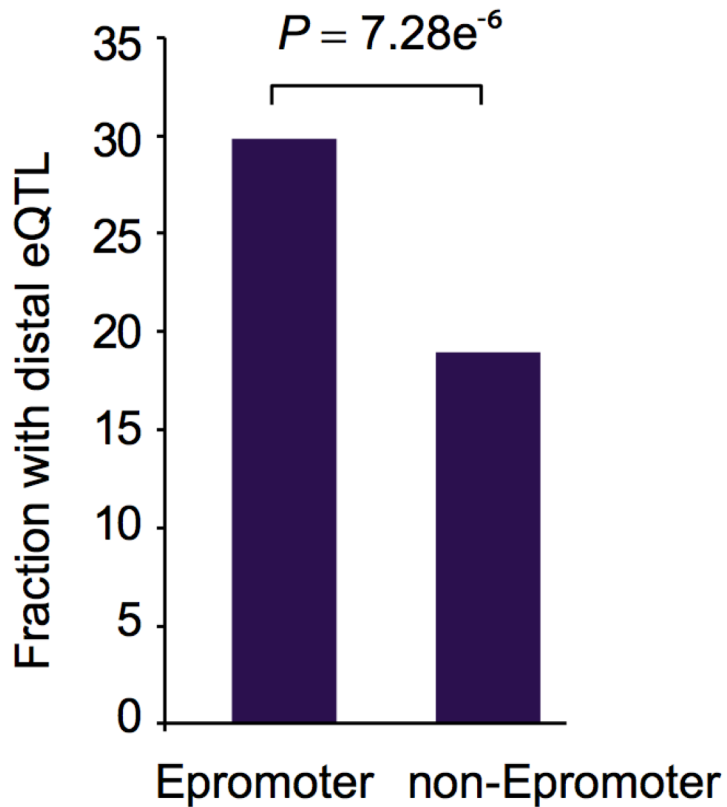
- Beta value is a measure of the effect the variant has on gene expression.
- Beta values is bound between -2 and 2.



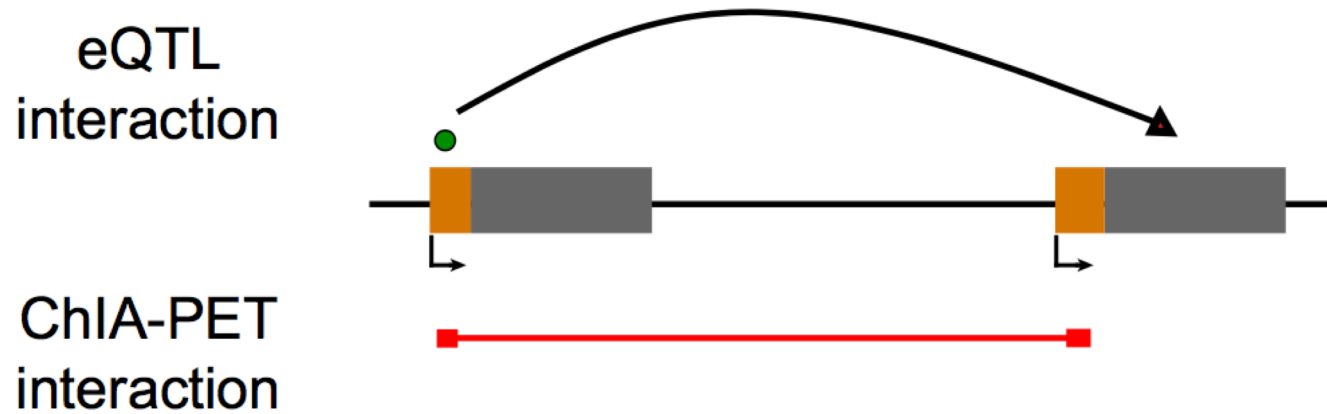
Are eQTLs overlapping Epromoters?



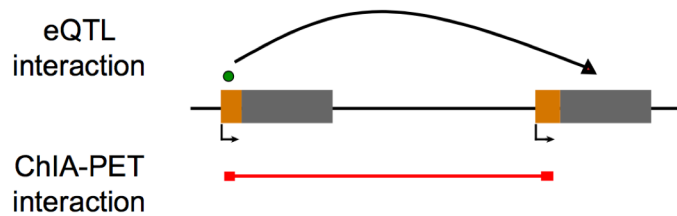
It is more likely to find an eQTL in an Epromoter than in a non-Epromoter



What is the effect of Epromoters on long range regulated genes?

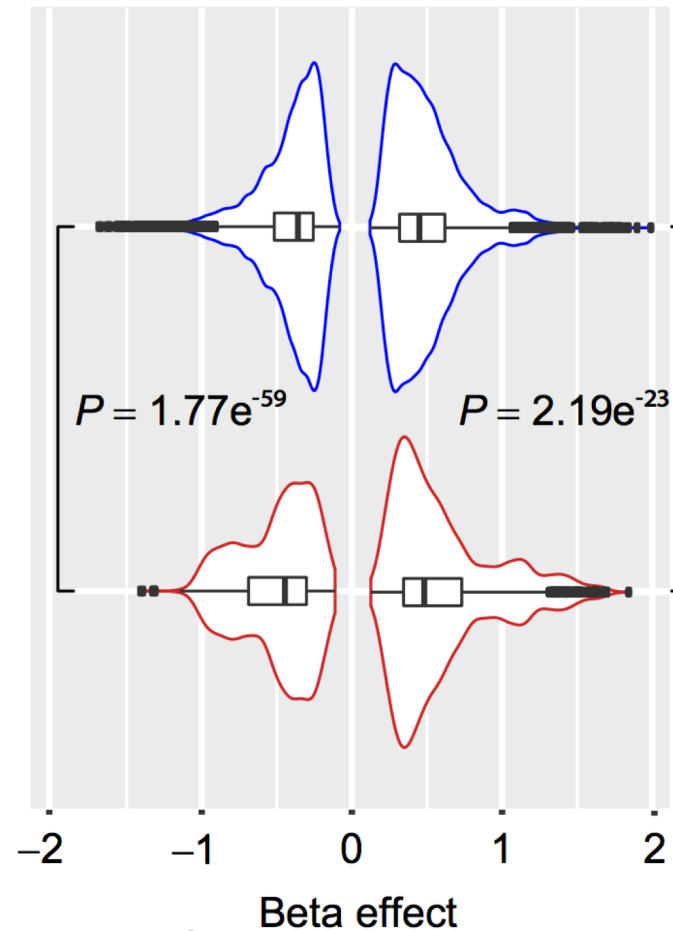


Epromoter eQTLs have a significantly stronger effect on distal gene expression



non-Epromoter

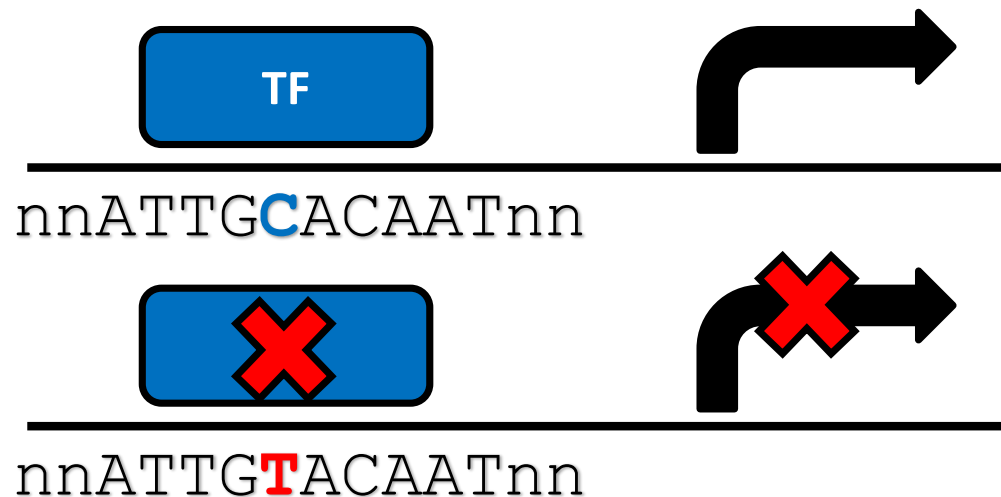
Epromoter



What is the effect of eQTLs that disrupt TF binding?

What is the effect of eQTLs that disrupt TF binding?

- Genetic variants in gene regulatory regions are known to cause many diseases
- Can be informative of transcription factors and regulatory mechanisms

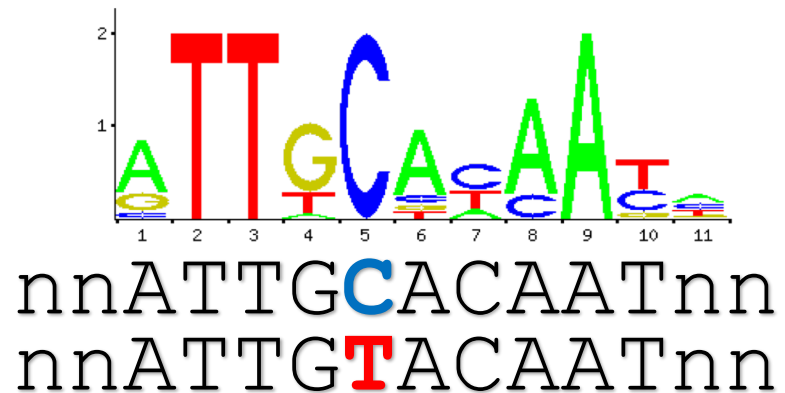


Detecting eQTL effect on TF binding

- *variation-scan*: PSSM based algorithm as part of the RSAT suite
- Flexible tool that can be used to scan any variable with any motif.
- Multiple organisms are available.



Walter Santana



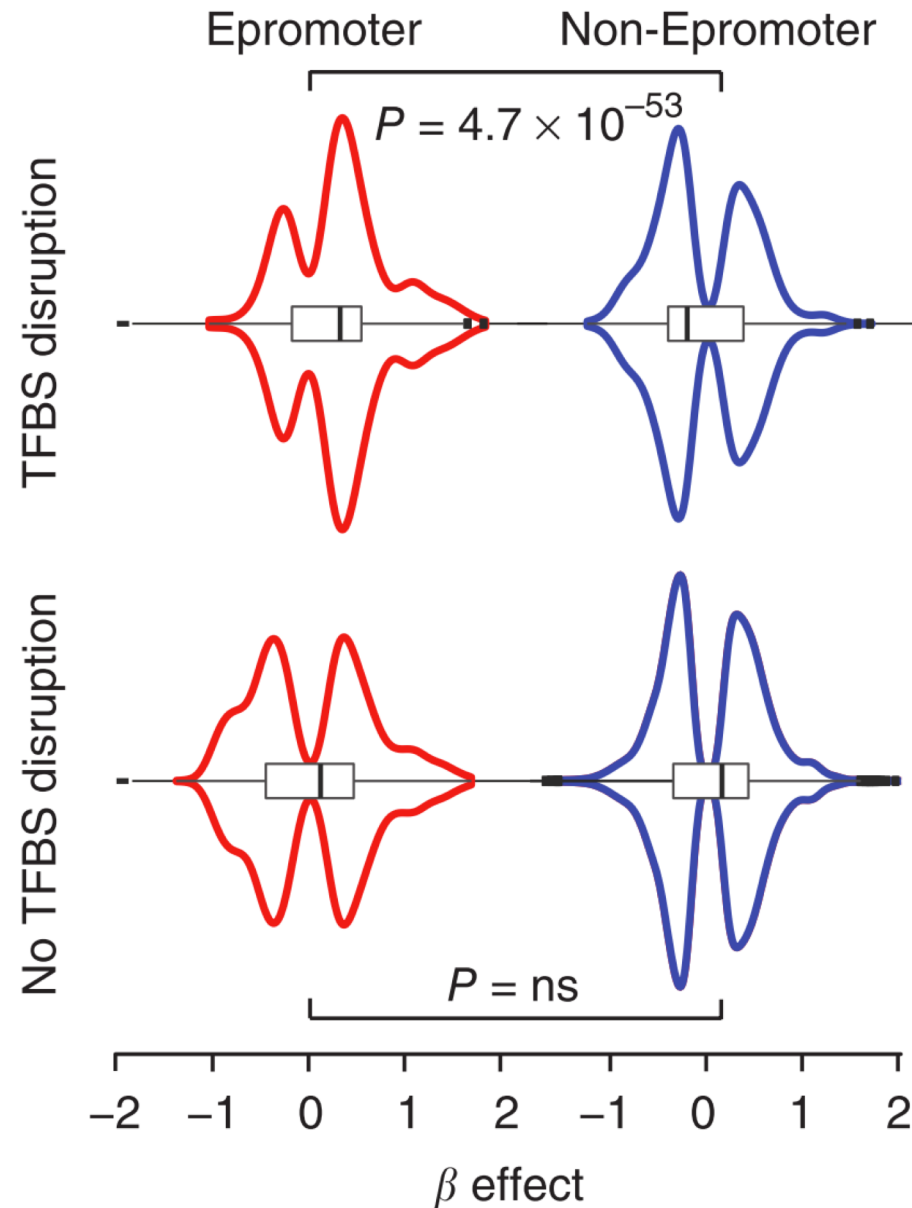
Weight Score

REF: 8.1

MUT: 2.1

DIF: 6.0

eQTL disrupting TF binding tend to have a positive beta-value



Summary

- High-throughput techniques can be used to identify long range regulatory sequences.
- Regulatory sequences can have several functions (Epromoters).
- Variation lying within Epromoters might impact distal gene expression beyond their cognate target gene.
- eQTLs that might disrupt TF binding tend to have a positive beta value in relation to their linked genes.

Future Steps

- Identify ePromoters for additional cell lines: Jukart, CEM, RPMI, GM12878
- Identify ePromoters under inflammatory conditions (IFN- alpha) in K562.
- Leverage public STARR-seq data in HELA with IFN repressors and hESC to recover additional ePromoters.
- Further develop position-scan to become user friendly.
- Establish a link between ePromoter regulated genes and disease (collaboration with Nancy J Cox's laboratory).
- Identify regulatory variants laying in ePromoters related to diseases that might be disrupting TF binding sites.

Acknowledgments:

Regulatory Genomics Lab

Karen Nuñez

Ana Villaseñor

Walter Santana

Monica Padilla

Lucia Ramirez

Diego Terol

Marisol Alvarez



INSERM

Salvatore Spicuglia

Lan Dao

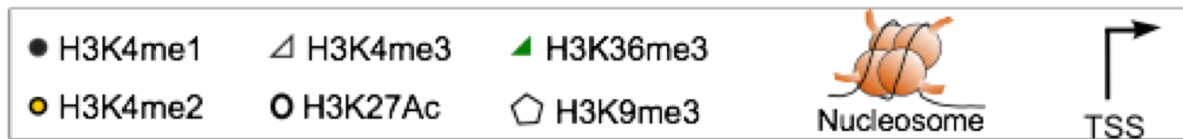
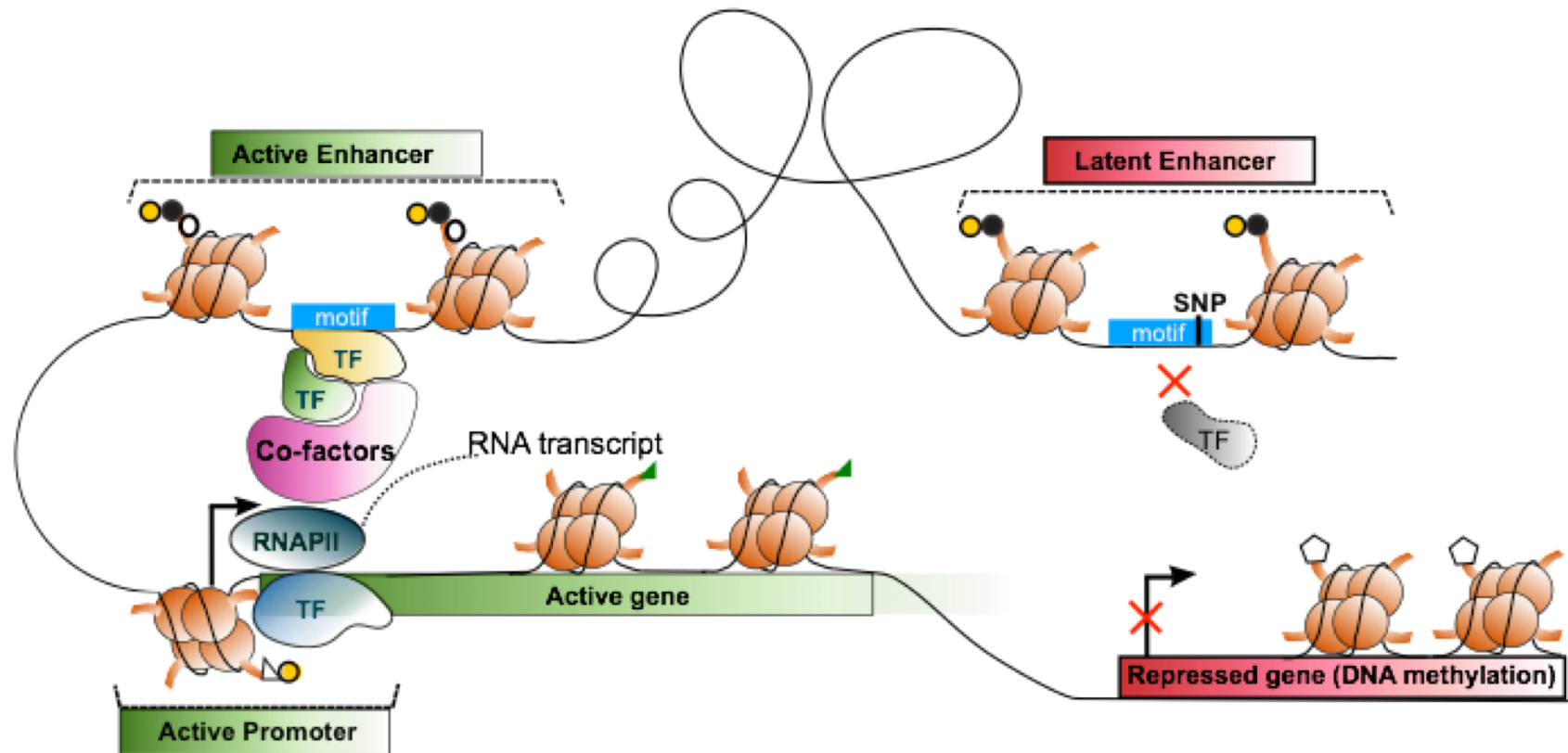
Ariel Galindo

Jacques van Helden

Jaime Castro

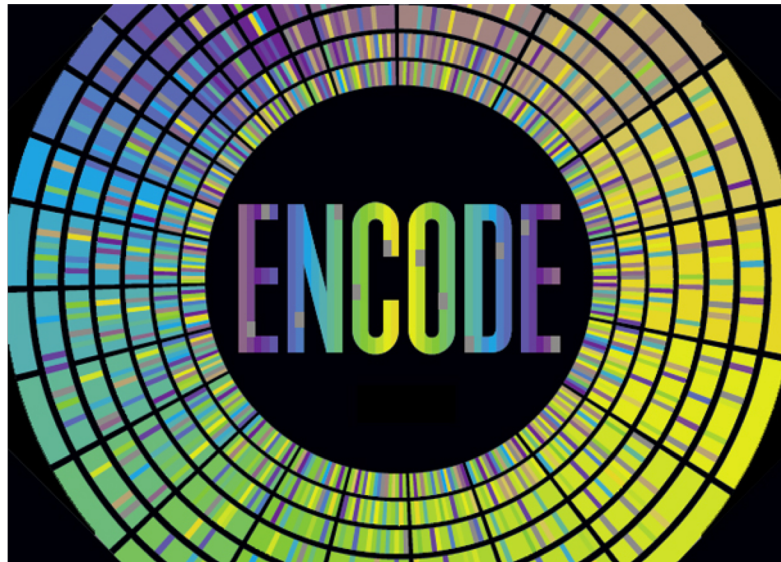


Challenging biological concepts



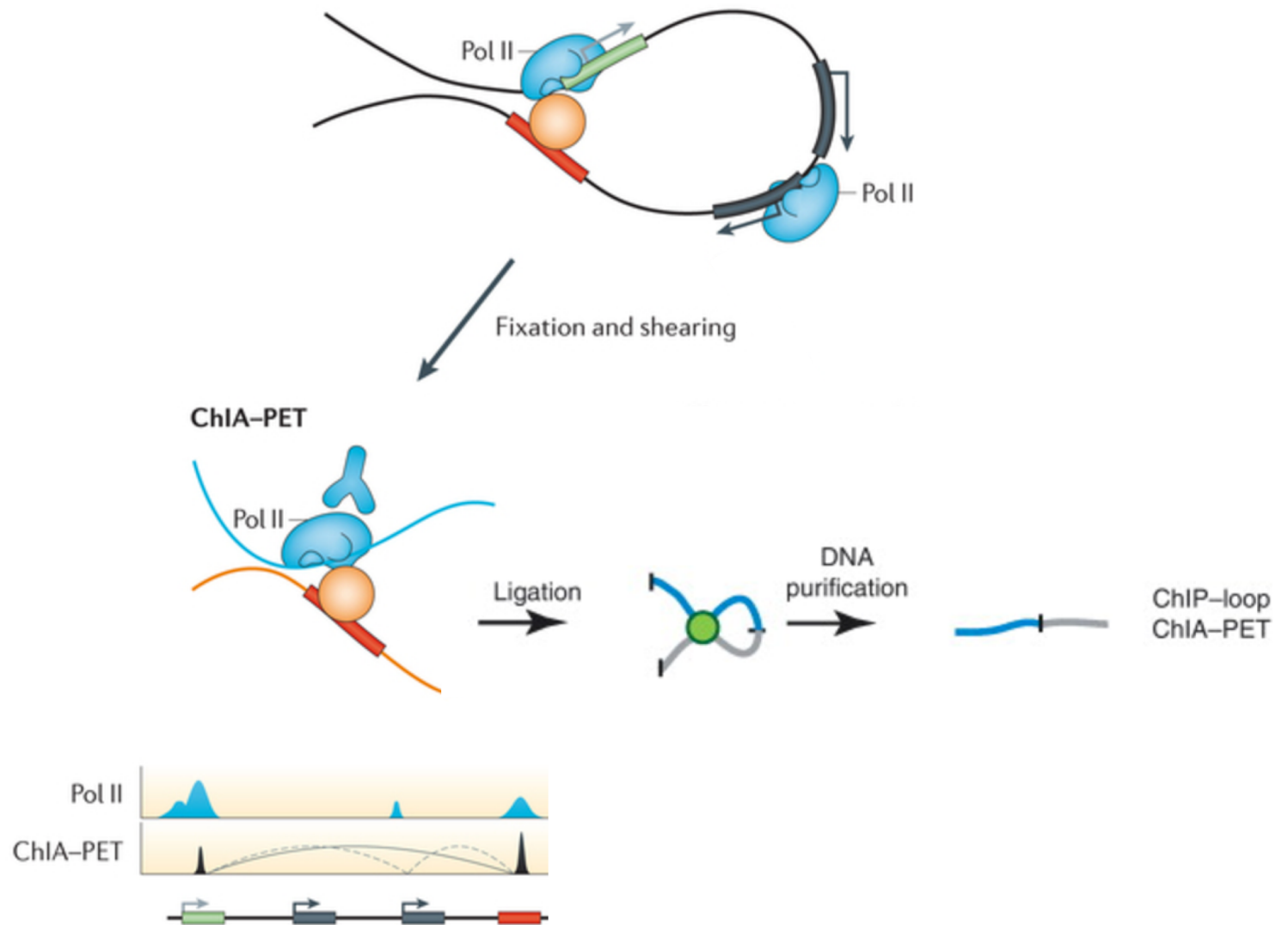
The Non-Coding Genome

- Around 80% of the genome has shown to have some biochemical interaction.
- Most non-coding region have a regulatory function.



How is regulation being affected by genetic variants?

Promoter interactions detected with Pol II ChIA-PET

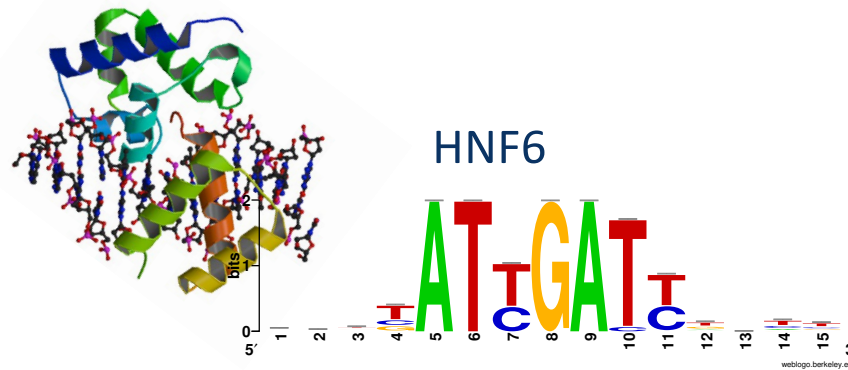


Major challenges to connect the regulatory genome to human pathologies

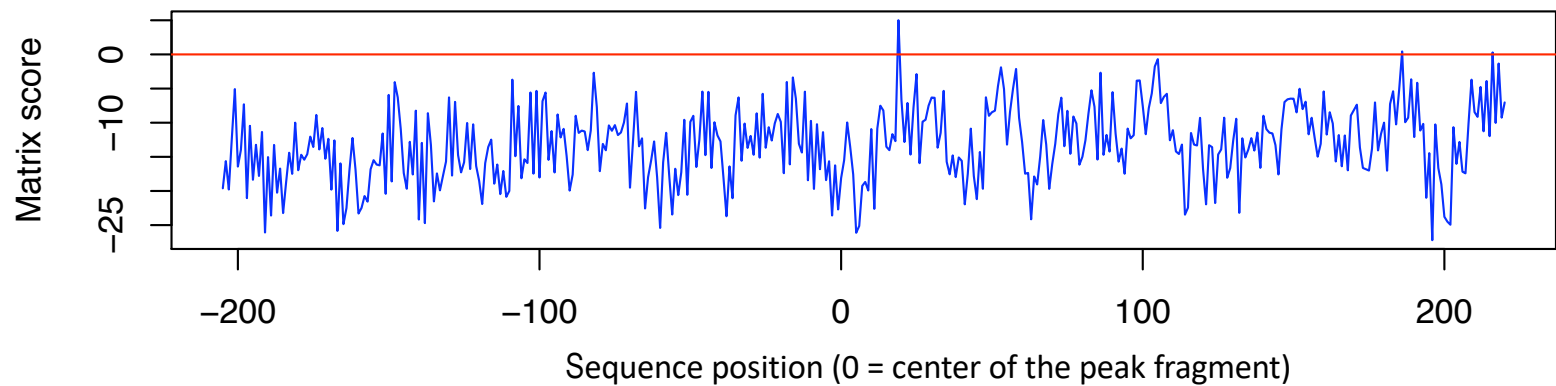
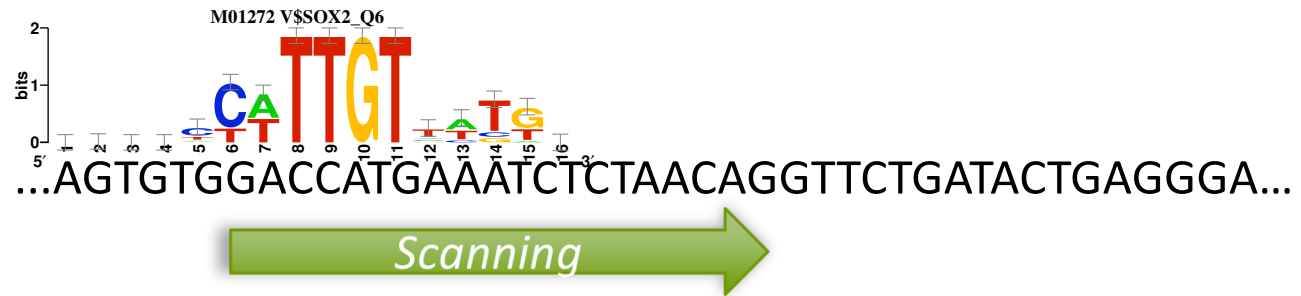
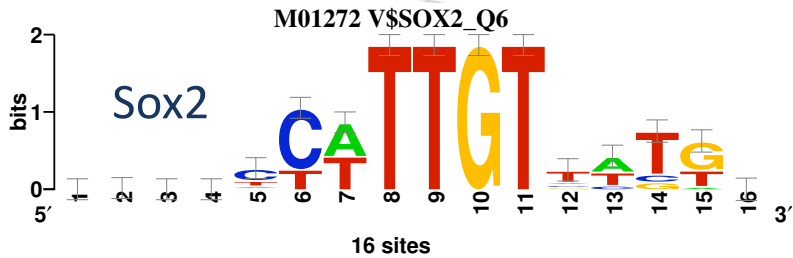
- Detect long-range regulatory regions (enhancers) and connect them to the regulated genes.
- Evaluate the effect genetic variants within regulatory sequences have on gene expression.

Transcription factors control gene regulation by binding to specific DNA sequences

- Transcription Factors interact with DNA binding to sequence specific sites.



Scanning a peak sequence with a TF binding motif

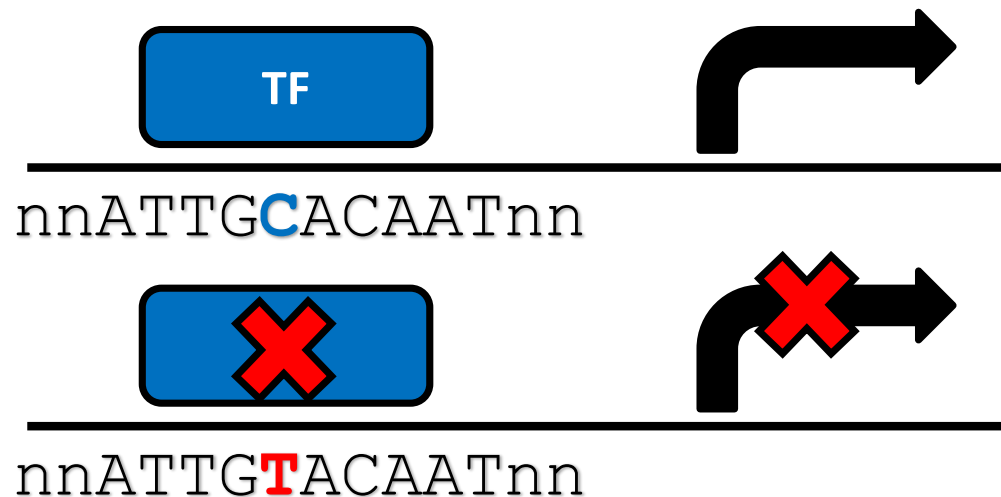


Detecting protein disruption of binding at base pair resolution

- Genetic variants in gene regulatory regions are known to cause many diseases
- Can be informative of transcription factors and regulatory mechanisms

Detecting protein disruption of binding at base pair resolution

- Genetic variants in gene regulatory regions are known to cause many diseases
- Can be informative of transcription factors and regulatory mechanisms

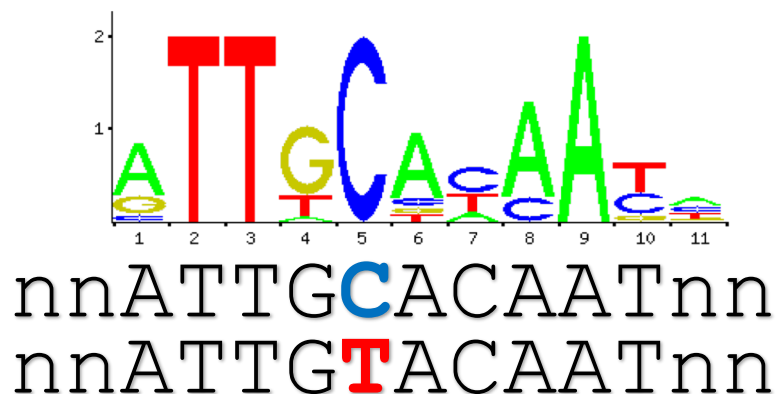


Detecting protein disruption of binding at base pair resolution: *variation-scan*

- *variation-scan*: PSSM based algorithm as part of the RSAT suite
- Flexibility with species, SNPs, motifs
- High-throughput



Jacques van Helden Walter Santana



Weight Score

REF: 8.1

MUT: 2.1

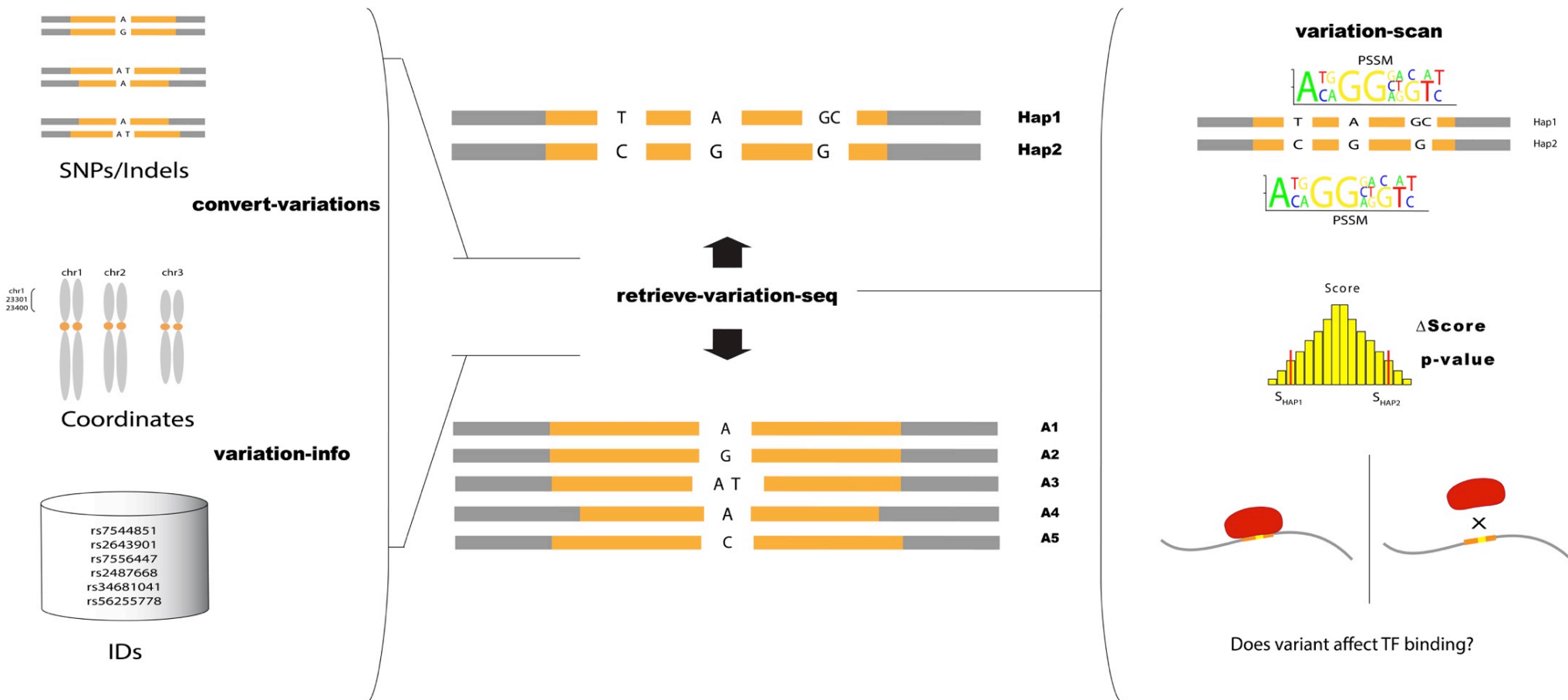
DIF: 6.0

Detecting protein disruption of binding at base pair resolution

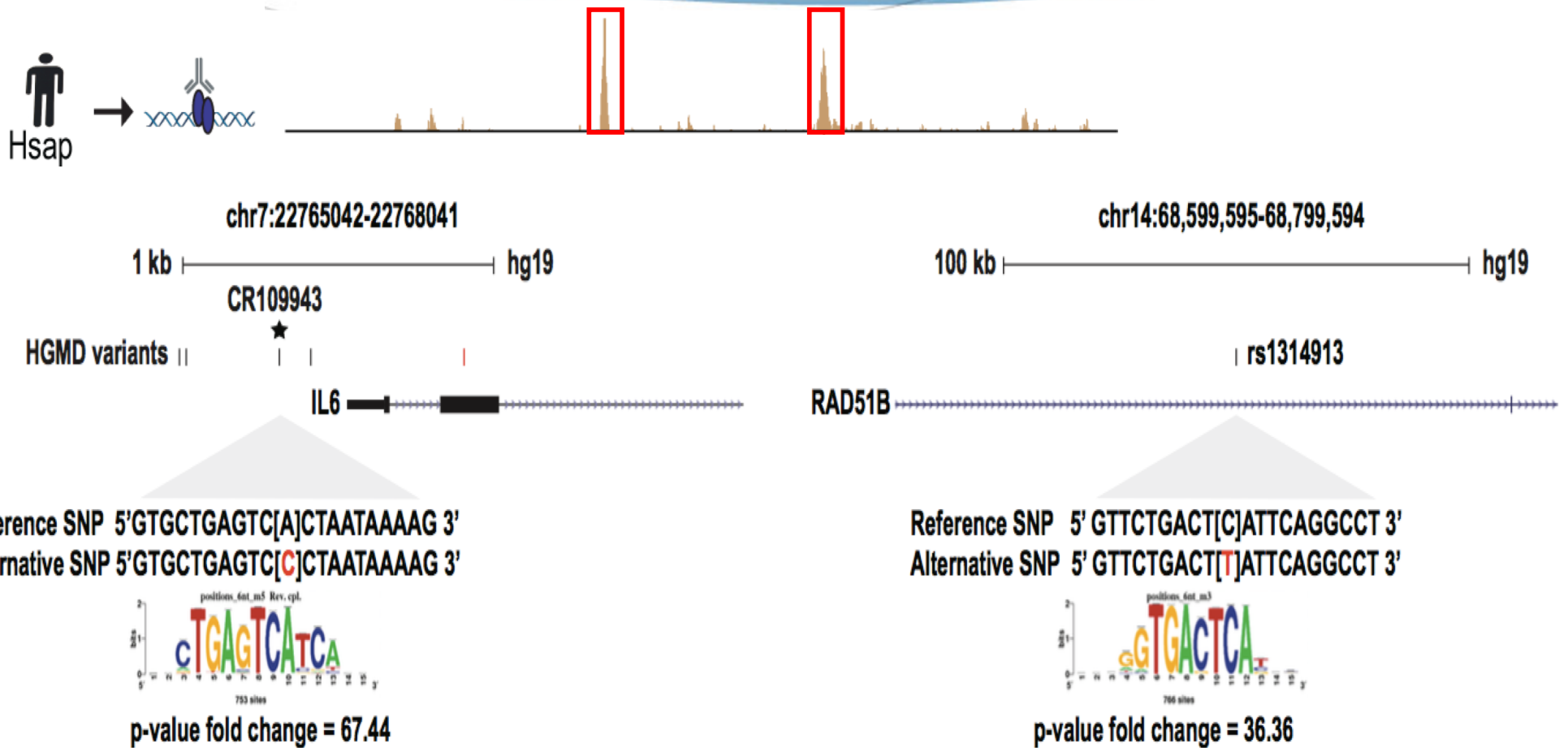
1. Genome-wide variations

2. Variation sequence reconstruction

3. Assess impact



Regulatory mutations alter JUN binding motifs in promoter region of *IL6* and intronic region of *RAD51B*.



Michael Wilson

DNA mutations in protein coding genes can explain many human diseases

Mutation Type	HGMD entries
Missense/nonsense	82176
Small deletions	22610
Splicing	13641
Gross deletions	10968
Small insertions	9423
Regulatory	2884
Gross insertions/duplications	2600
Small indels	2173
Complex rearrangements	1504
Repeat variations	434
Mutation total	148413
Gene total	6137

DNA mutations in protein coding genes can explain many human diseases

Mutation Type	HGMD entries
Missense/nonsense	82176
Small deletions	22610
Splicing	13641
Gross deletions	10968
Small insertions	9423
Regulatory	2884
Gross insertions/duplications	2600
Small indels	2173
Complex rearrangements	1504
Repeat variations	434
Mutation total	148413
Gene total	6137



DNA mutations in protein coding genes can explain many human diseases

Mutation Type	HGMD entries
Missense/nonsense	82176
Small deletions	22610
Splicing	13641
Gross deletions	10968
Small insertions	9423
Regulatory	2884
Gross insertions/duplications	2600
Small indels	2173
Complex rearrangements	1504
Repeat variations	434
Mutation total	148413
Gene total	6137



DNA mutations in protein coding genes can explain many human diseases

Mutation Type	HGMD entries
Missense/nonsense	82176
Small deletions	22610
Splicing	13641
Gross deletions	10968
Small insertions	9423
Regulatory	2884
Gross insertions/duplications	2600
Small indels	2173
Complex rearrangements	1504
Repeat variations	434
Mutation total	148413
Gene total	6137

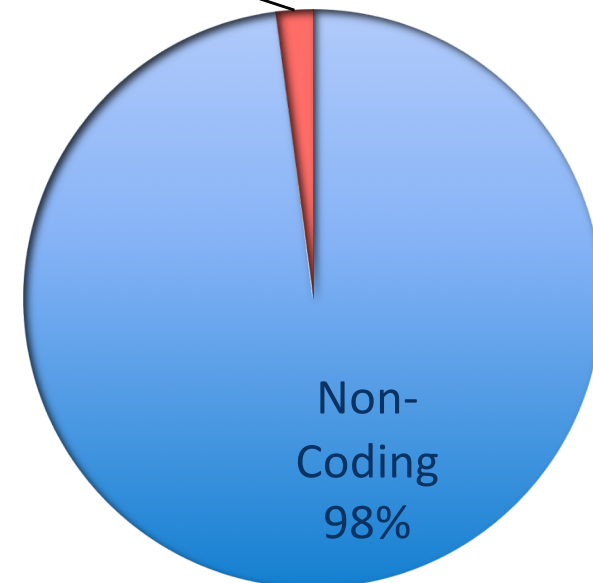


DNA mutations in protein coding genes can explain many human diseases

Mutation Type	HGMD entries
Missense/nonsense	82176
Small deletions	22610
Splicing	13641
Gross deletions	10968
Small insertions	9423
Regulatory	2884
Gross insertions/duplications	2600
Small indels	2173
Complex rearrangements	1504
Repeat variations	434
Mutation total	148413
Gene total	6137

Coding
2%

Human Genome

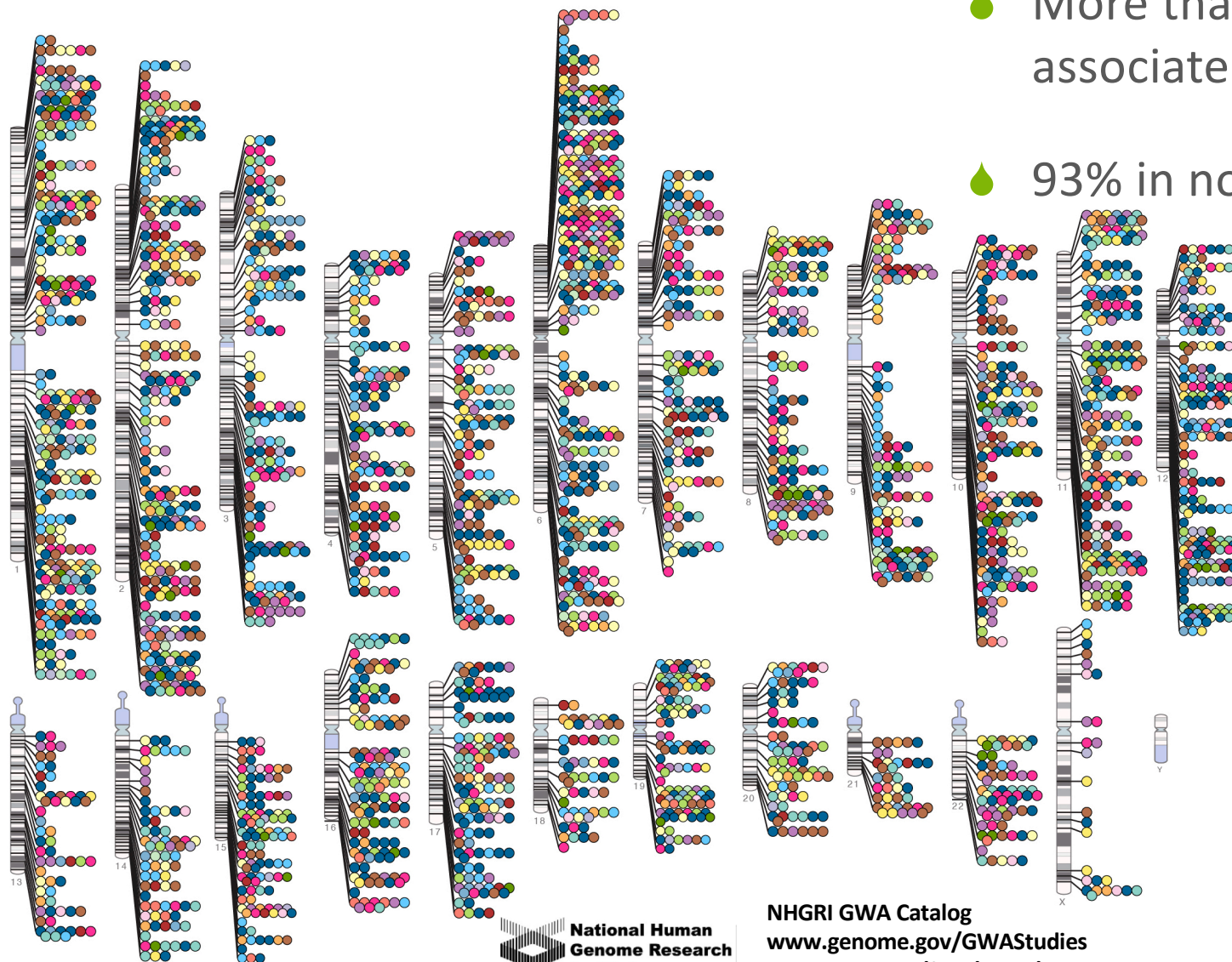


Published Genome-Wide Associations since 12/2012

Published GWA at $p \leq 5 \times 10^{-8}$ for 17 trait categories

More than 17,000 reported associated SNPs

93% in non-coding sequence



- Digestive system disease
- Cardiovascular disease
- Metabolic disease
- Immune system disease
- Nervous system disease
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurement
- Other measurement
- Response to drug
- Biological process
- Cancer
- Other disease
- Other trait

Weight score difference has proven not to be enough to reliably detect variations affecting TF binding in a genome wide scale

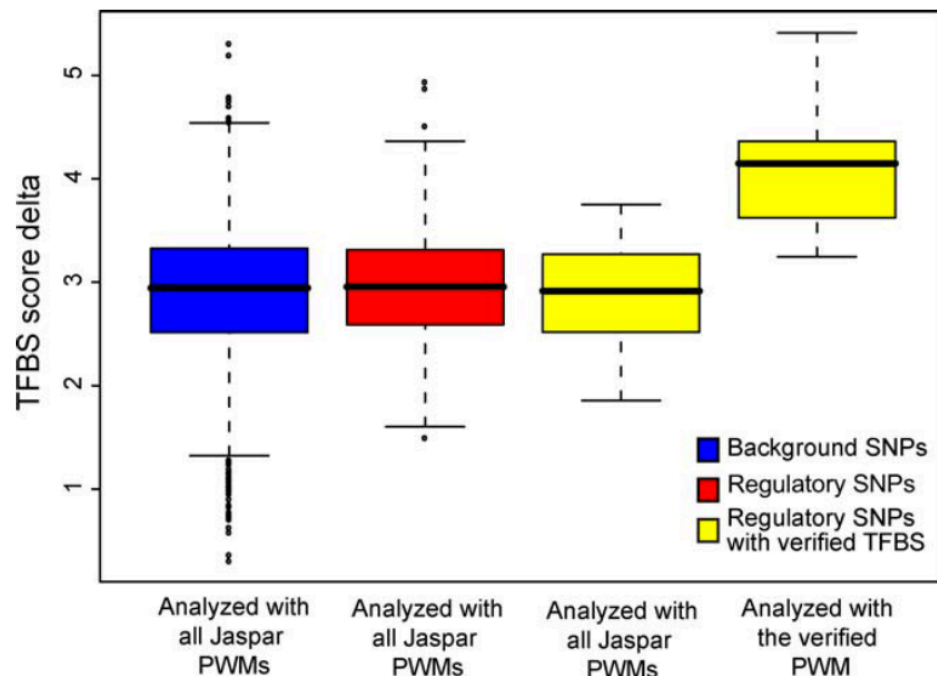
Weight score difference has proven not to be enough to reliably detect variations affecting TF binding in a genome wide scale

OPEN ACCESS Freely available online

PLoS COMPUTATIONAL BIOLOGY

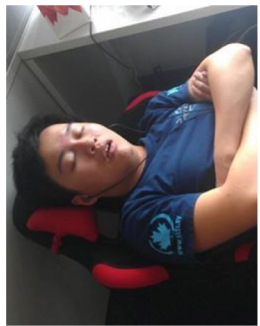
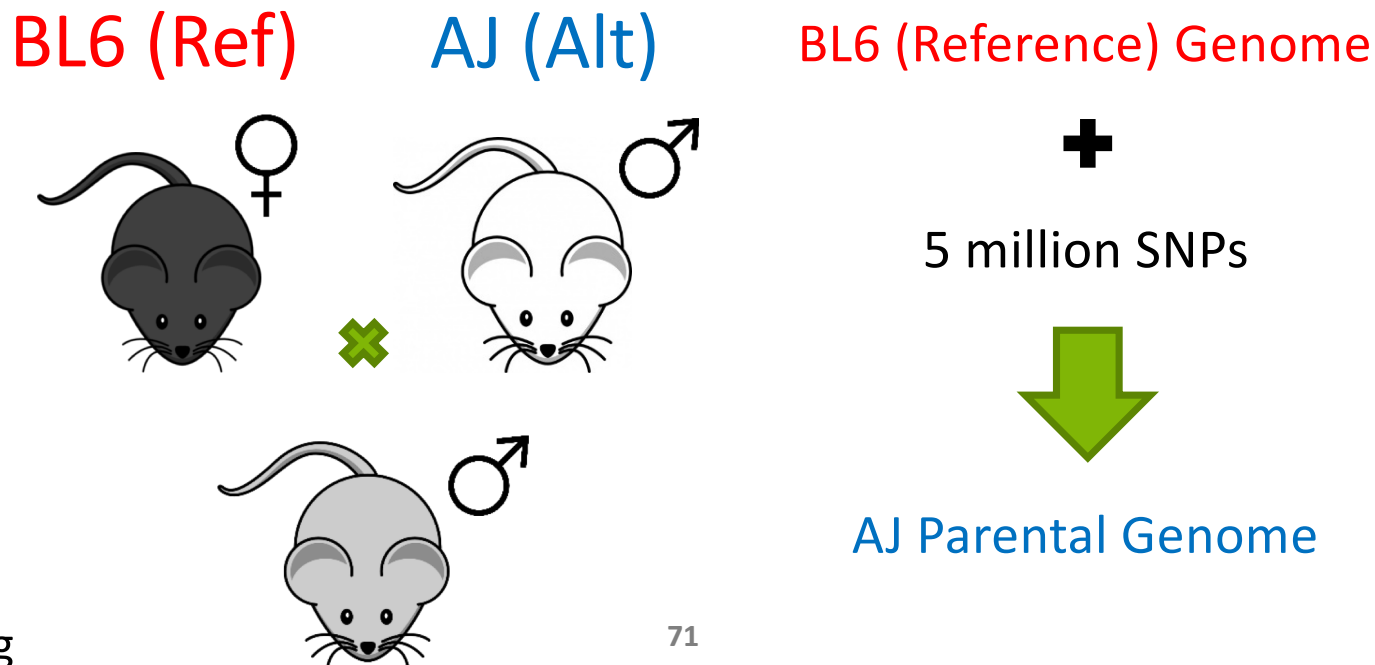
In Silico Detection of Sequence Variations Modifying Transcriptional Regulation

Malin C. Andersen^{1,2}, Pär G. Engström^{3,4}, Stuart Lithwick^{5,6}, David Arenillas⁶, Per Eriksson², Boris Lenhard³, Wyeth W. Wasserman^{6*}, Jacob Odeberg^{1,2,7*}



Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*

- F1 mice with inbred parents
- All SNP positions known – reduce false positives and search space
- Each allele compared within same organism and genomic context



Minggao Liang

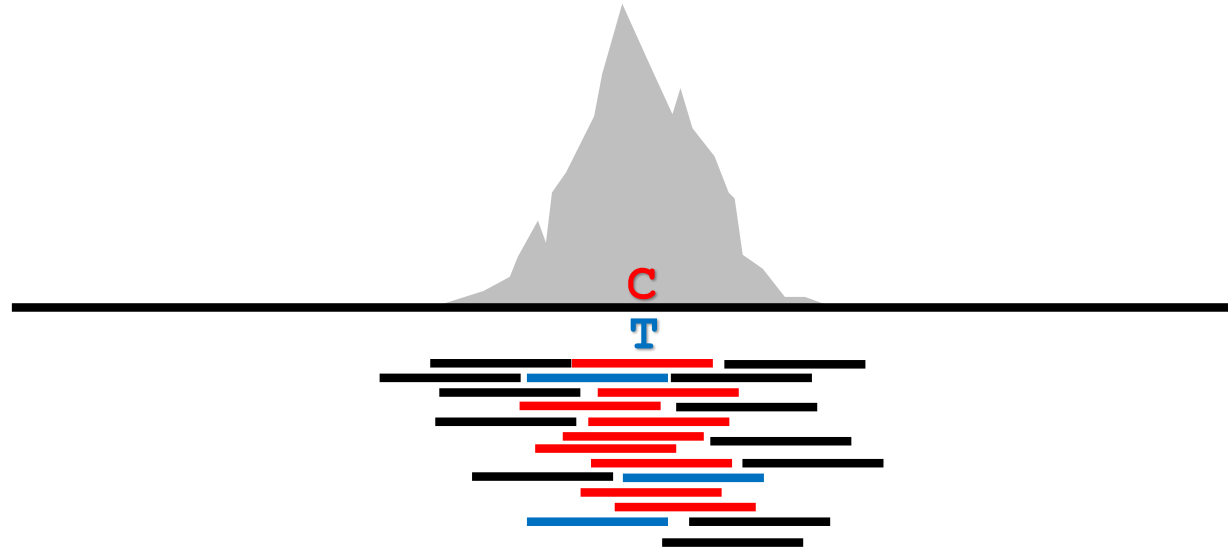
Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*



HNF6
CEBPA

Processing	Purpose	Remaining mapped reads
mm9 alignment		33.0 million
Filter against Encode Blacklist regions	<ul style="list-style-type: none">• Remove known sequencing artifact regions	31.9 million
MACS	<ul style="list-style-type: none">• Peak calling	
WASP	<ul style="list-style-type: none">• Extract reads overlapping known alleles• Remove reads with alignment bias• Remove duplicates	1.68 million 1.50 million 0.97 million
ALEA	<ul style="list-style-type: none">• Alignment to BL6 Parental genome (REF)• Alignment to A/J Parental genome (ALT)	0.43 million 0.47 million

Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*



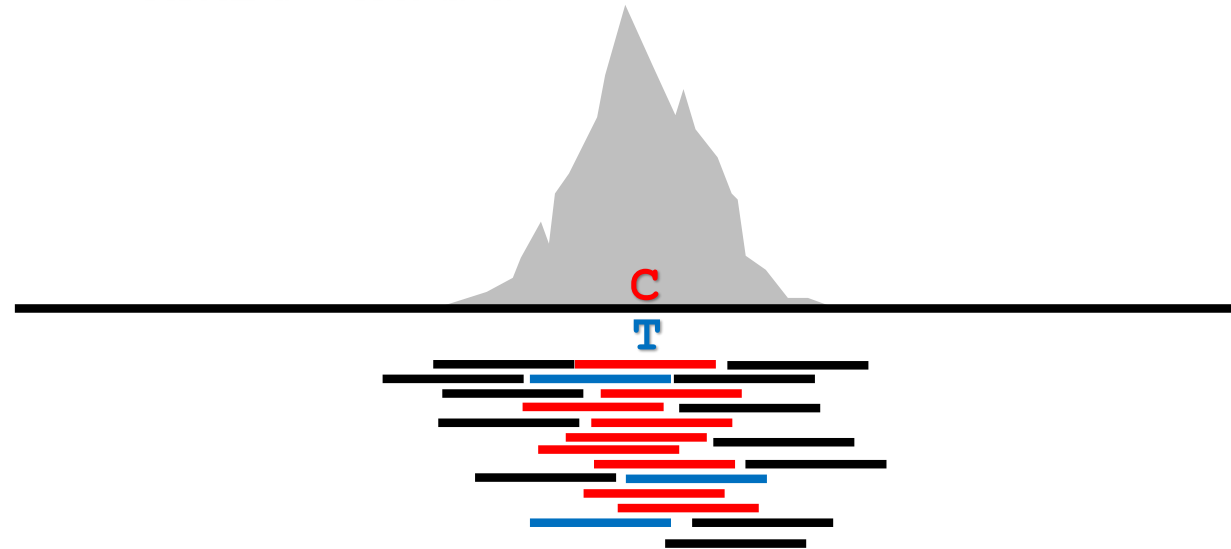
Count Reads :

REF: 9

ALT: 3

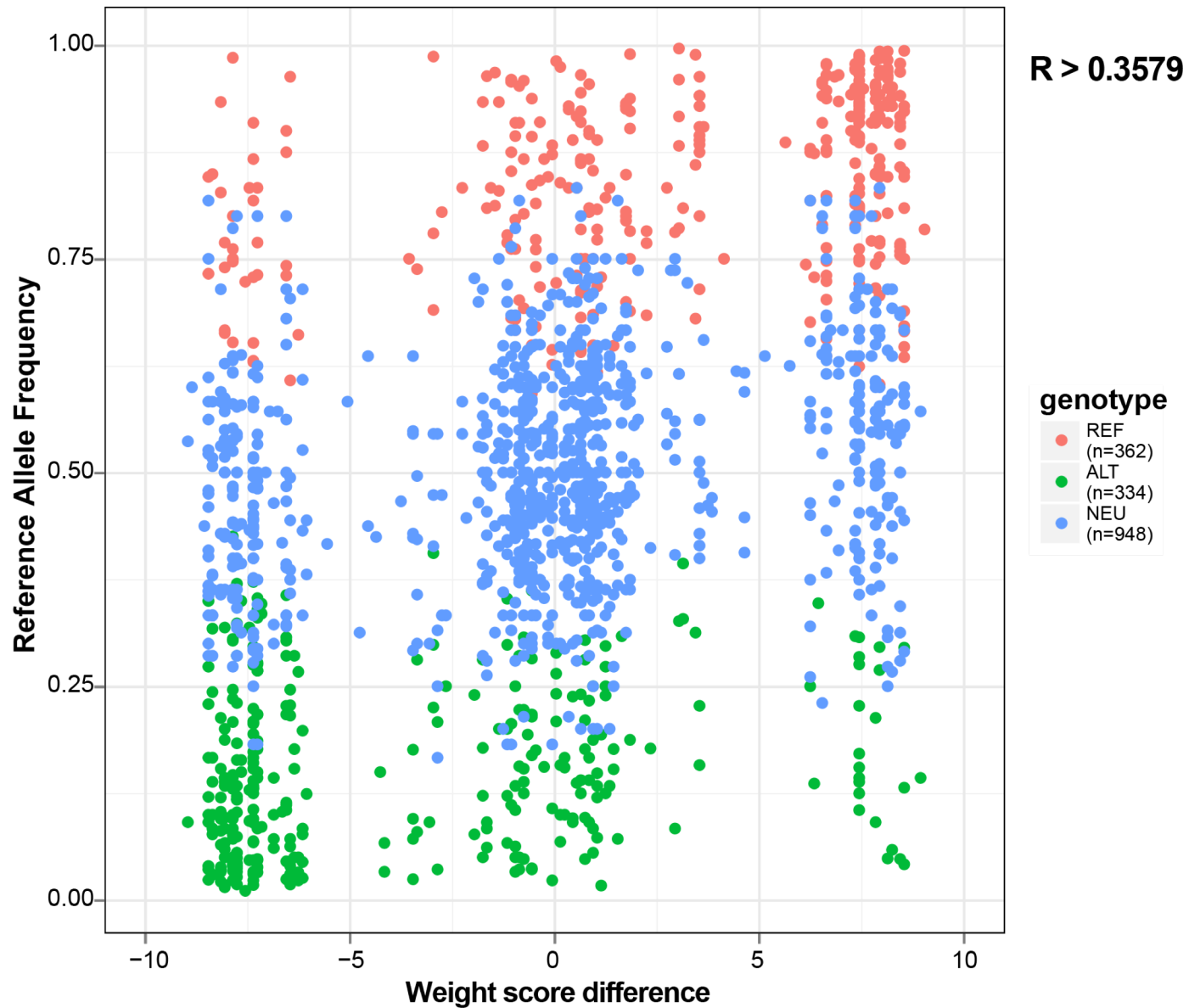
RAF: 0.75

Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*



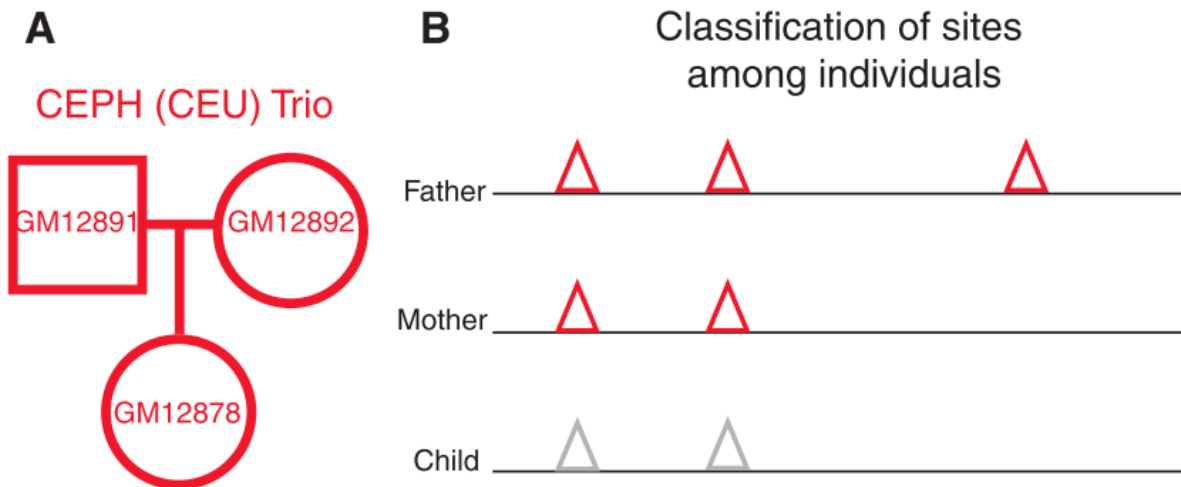
Count	Reads:	Weight	score:
REF:	9	REF:	7.0
ALT:	3	ALT:	2.0
RAF:	0.75	DIF:	5.0

Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*



Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*

- Human trio families with genotypes

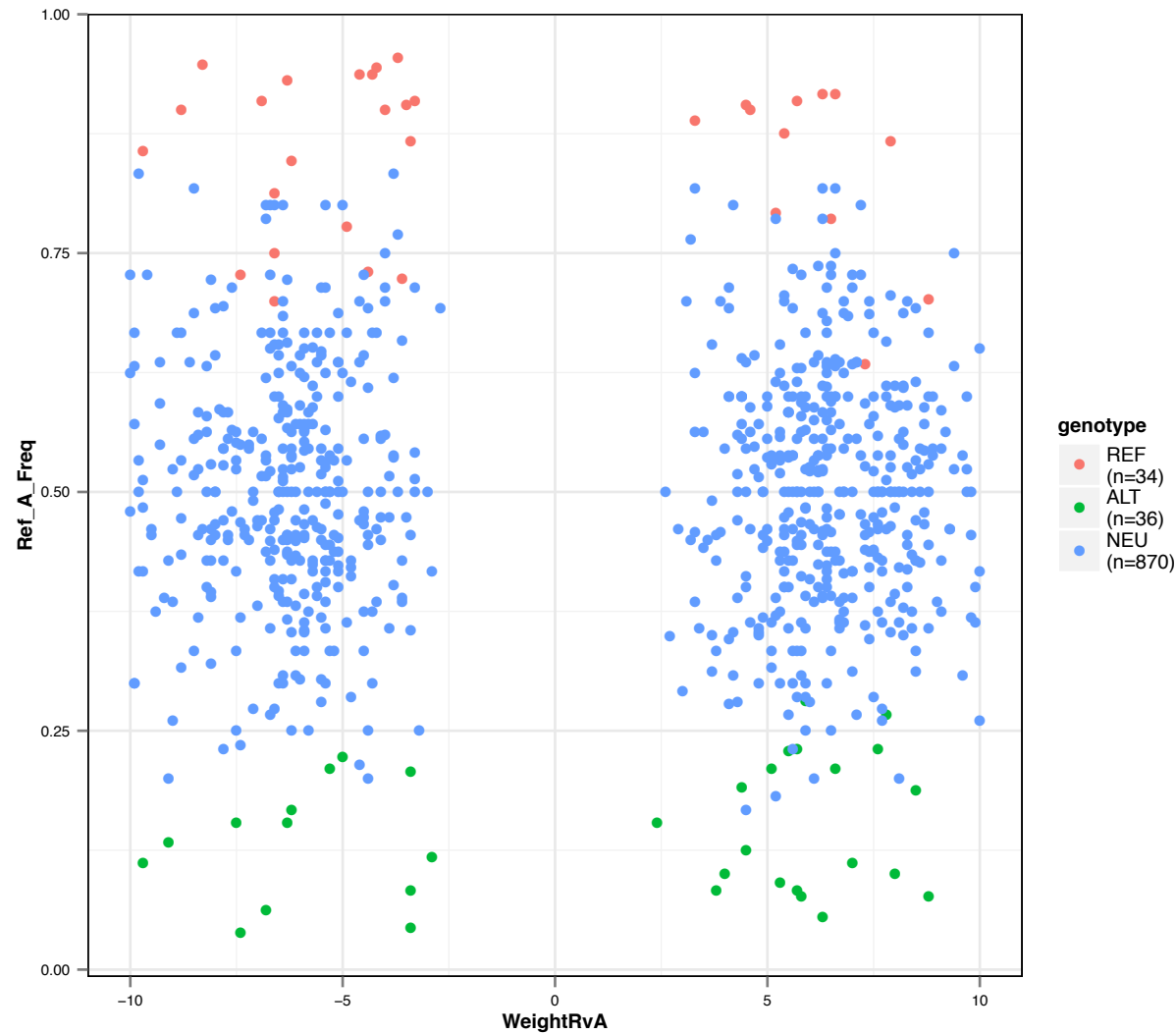


Extract *in vivo* motif
from ChIP data



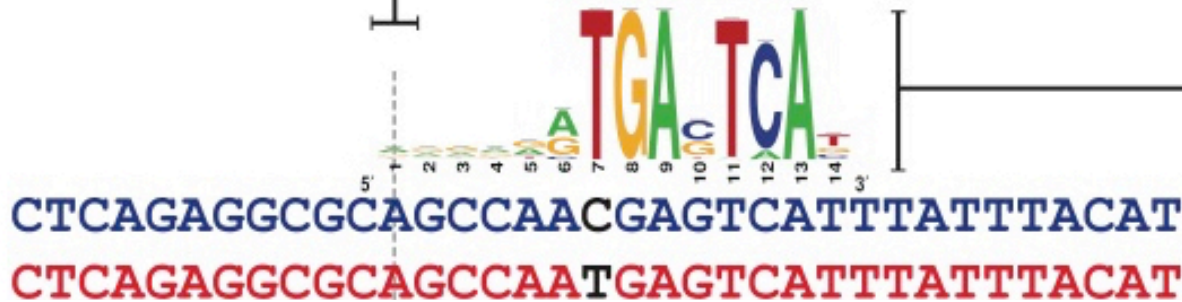
Detecting protein disruption of binding at base pair resolution: *Heterozygous ChIP-seq*

- Human trio families with genotypes



Detecting protein disruption of binding at base pair resolution: *variation-scan*

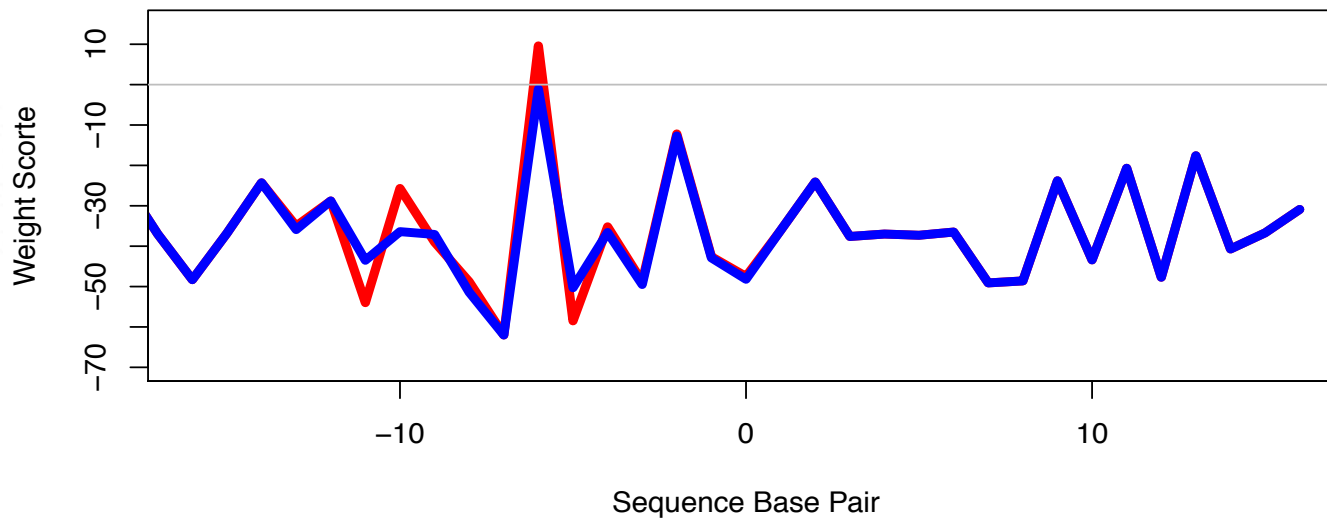
Scanning Variant rs148805532
JUN motif



3. Scan all windows overlapping SNP & choose best Pvalue for motif match

1. Motif of interest

2. SNP of interest & flanking bases



4. Report Pvalue fold change between best and worst allele

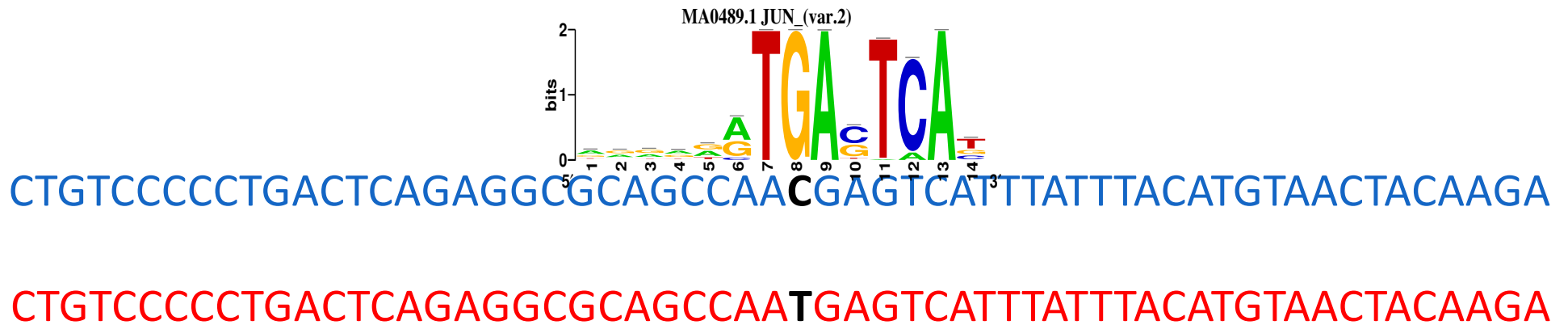
Detecting protein disruption of binding at base pair resolution: Variation Scan

CTGTCCCCCTGACTCAGAGGCGCAGCCAAC**C**GAGTCATTTATTTACATGTA ACTACAAGA

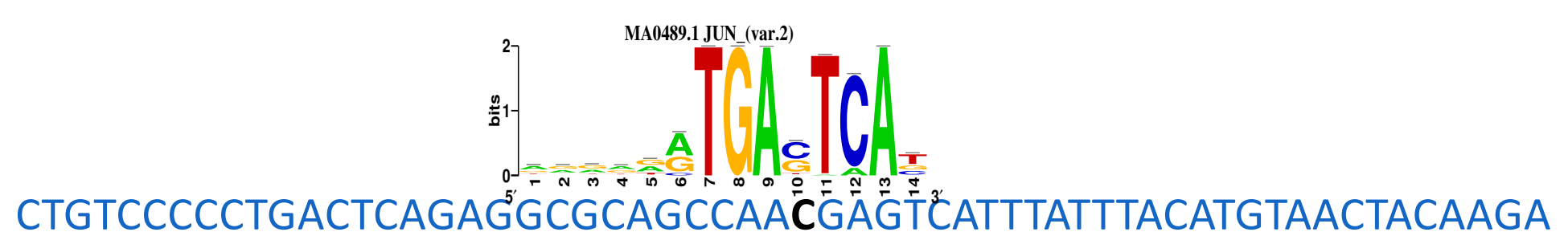
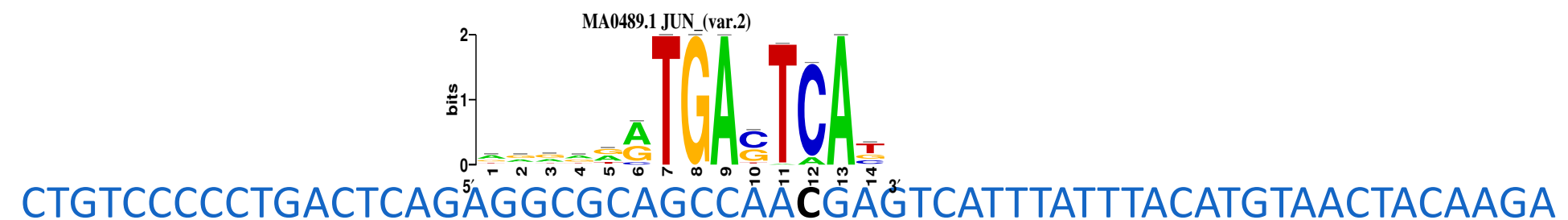
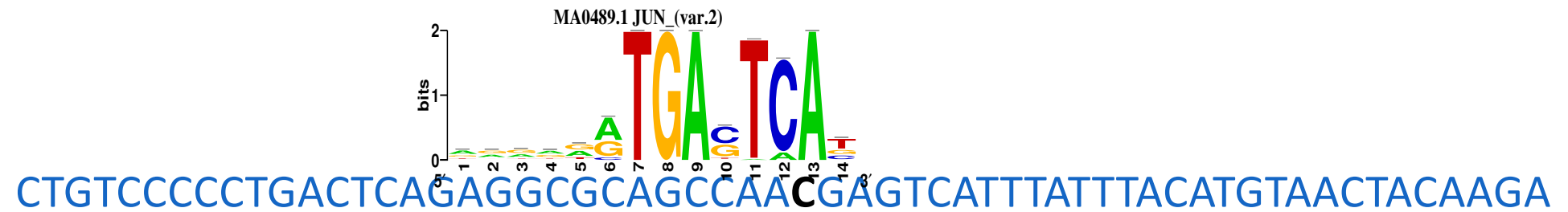
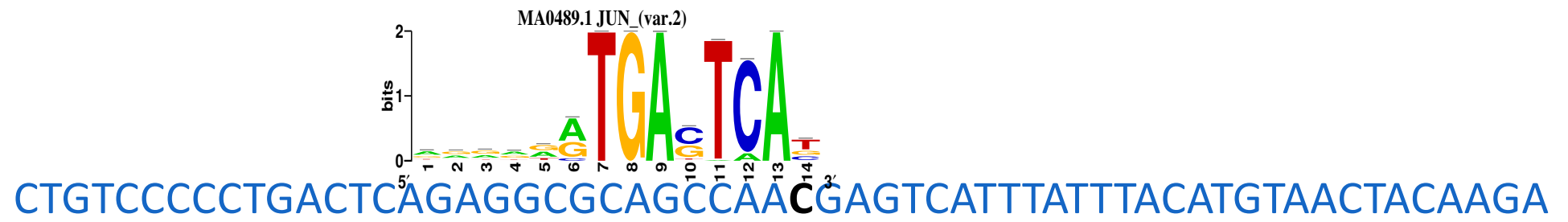
CTGTCCCCCTGACTCAGAGGCGCAGCCAAT**T**GAGTCATTTATTTACATGTA ACTACAAGA

Collaboration with Dr. van Helden, Université d'Aix-Marseille (AMU).

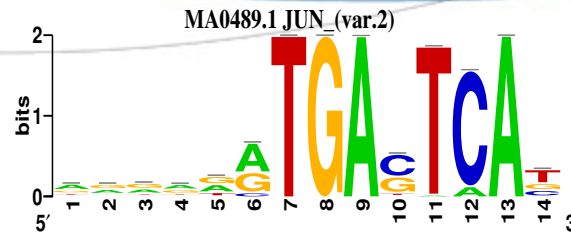
Detecting protein disruption of binding at base pair resolution: Variation Scan



Detecting protein disruption of binding at base pair resolution: Variation Scan

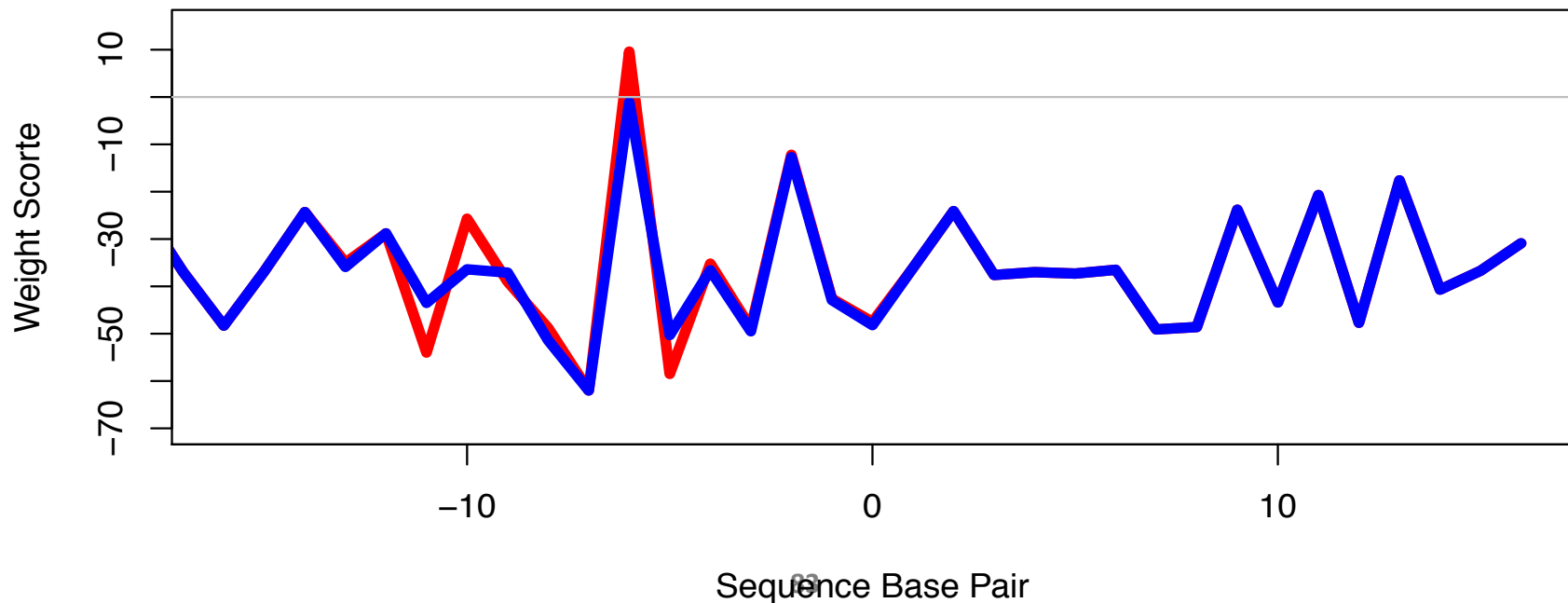


Detecting protein disruption of binding at base pair resolution: Variation Scan

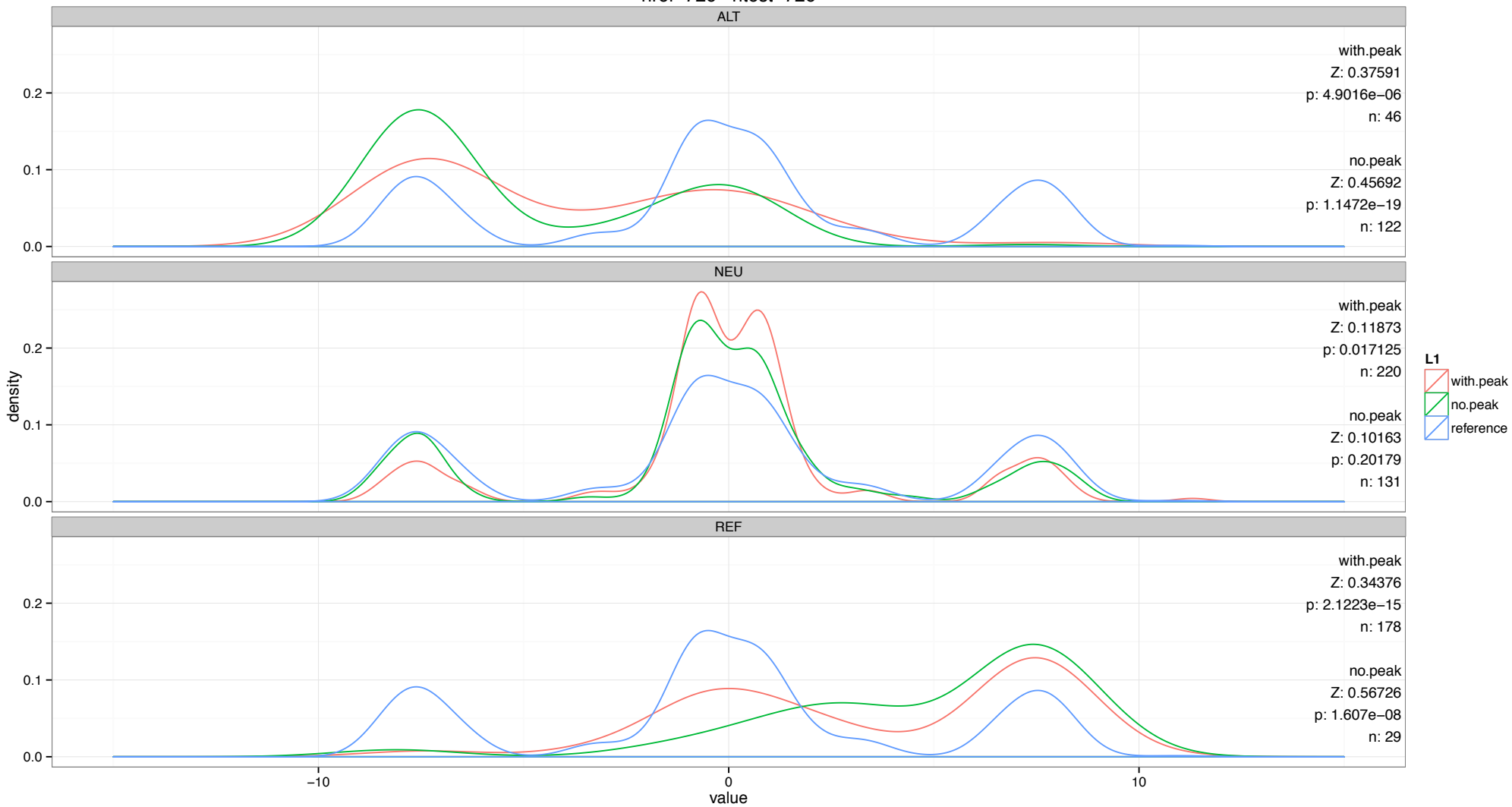


CTGTCCCCCTGACTCAGAGGCGCAGCCAAC**C**GAGTCATTTATTTACATGTA ACTACAAGA
CTGTCCCCCTGACTCAGAGGCGCAGCCAAT**T**GAGTCATTTATTTACATGTA ACTACAAGA

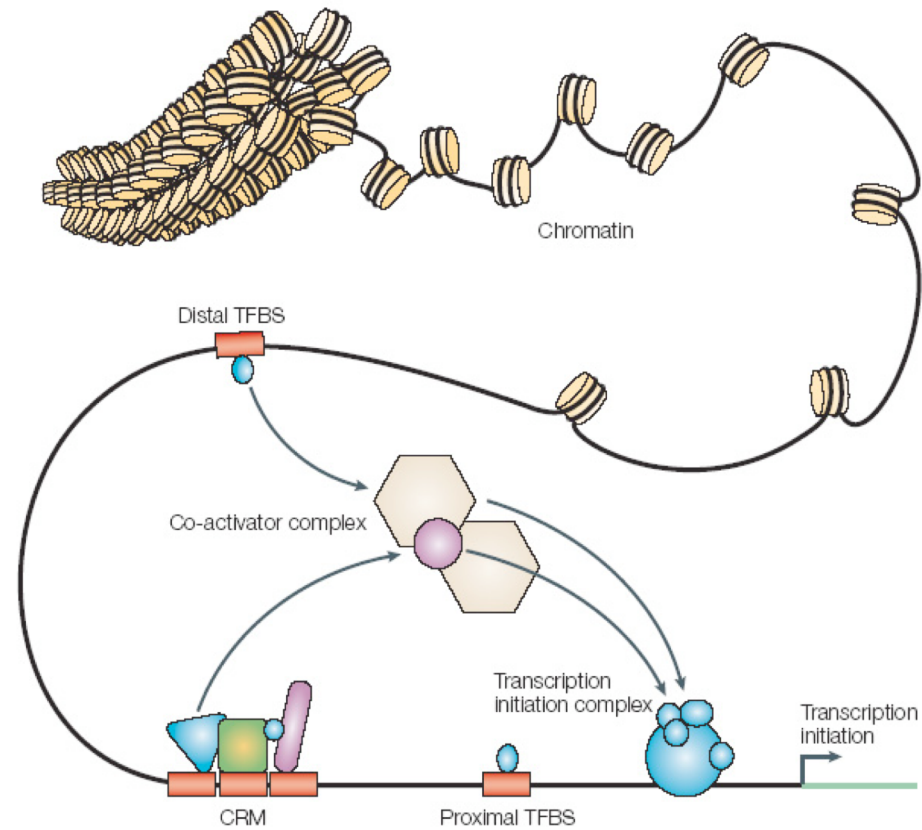
Scanning Variant rs148805532
JUN motif



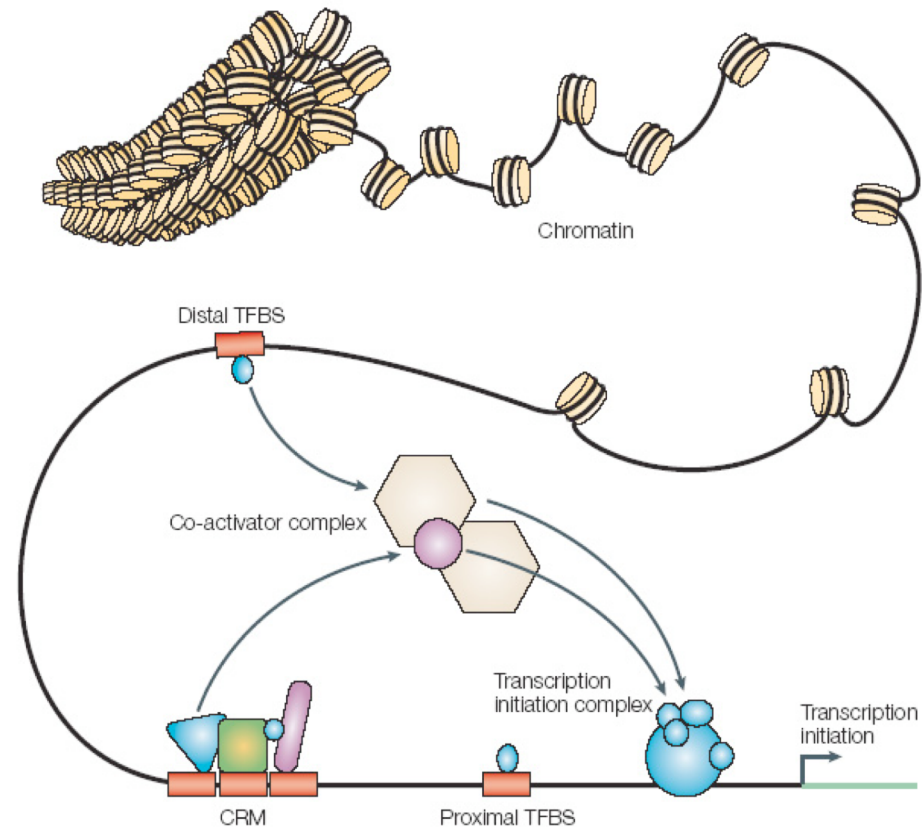
HNF6
nref=726 ntest=726



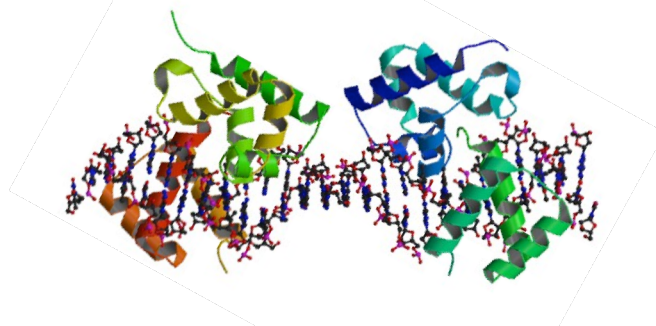
Multiple transcription factors can cooperatively regulate gene regulation



Multiple transcription factors can cooperatively regulate gene regulation

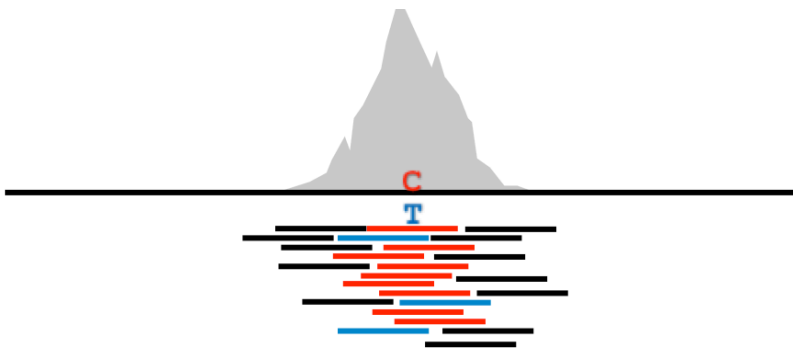


- ◆ Transcription Factors can bind in cis – Known as Cluster of Regulatory Modules (CRMs)

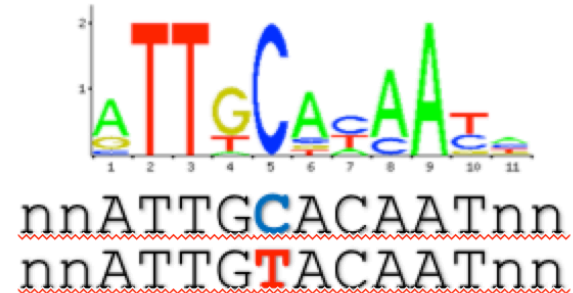


Can we detect when alleles disrupt cofactors?

- Why do we expect co-factors?
 - TFs in Eukaryotes tend to bind in clusters of regulatory modules.



Count Reads: Weight score:
REF: 9 **REF: 7.0**
ALT: 3 **ALT: 2.0**
RAF: 0.75 **DIF: 5.0**



Weight Score

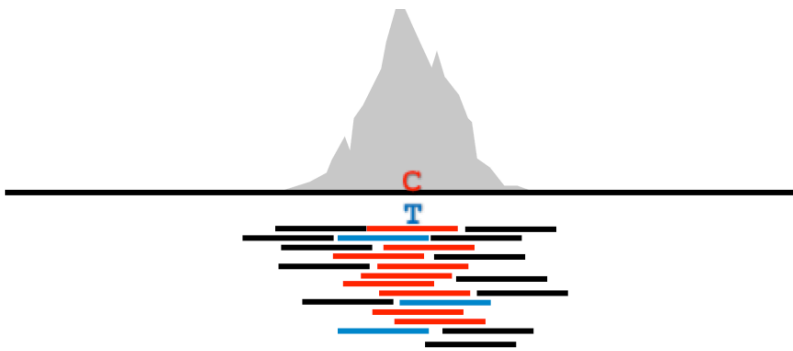
REF: 8.1

MUT: 2.1

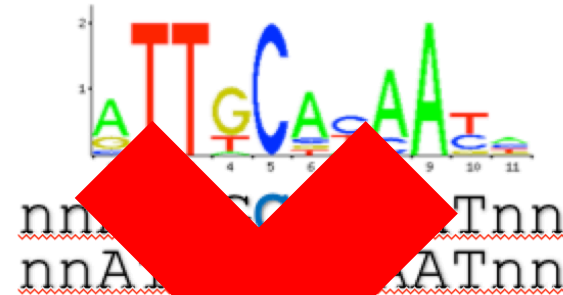
DIF: 6.0

Can we detect when alleles disrupt cofactors?

- Why do we expect co-factors?
 - TFs in Eukaryotes tend to bind in clusters of regulatory modules.

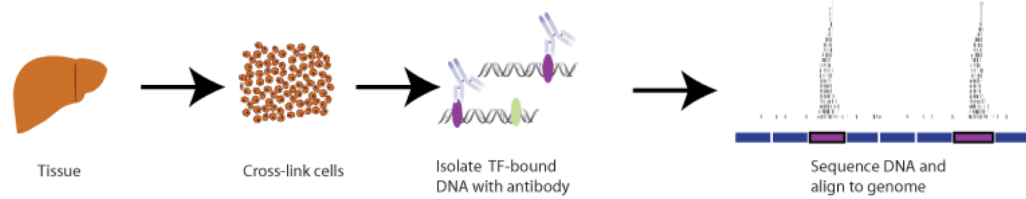


Count	Reads:	Weight	score:
REF: 9		REF: 7.0	
ALT: 3		ALT: 2.0	
RAF: 0.75		DIF: 5.0	



REF: 8.1
MUT: 2.1
DIF: 6.0

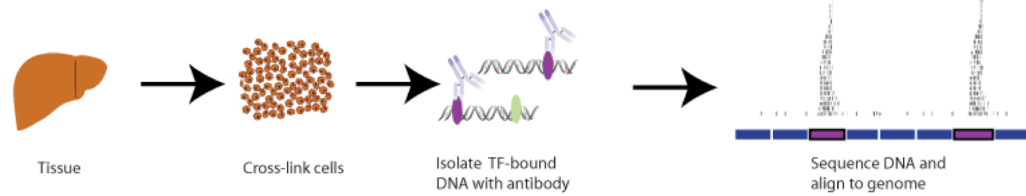
Experimentally determine genome wide binding of 4 TFs relevant for liver from five mammalian species.



Collaboration:
Benoit Ballester (INSERM)
Duncan Odom (Cambridge U)
Paul Flicek (EBI)



Experimentally determine genome wide binding of 4 TFs relevant for liver from five mammalian species.



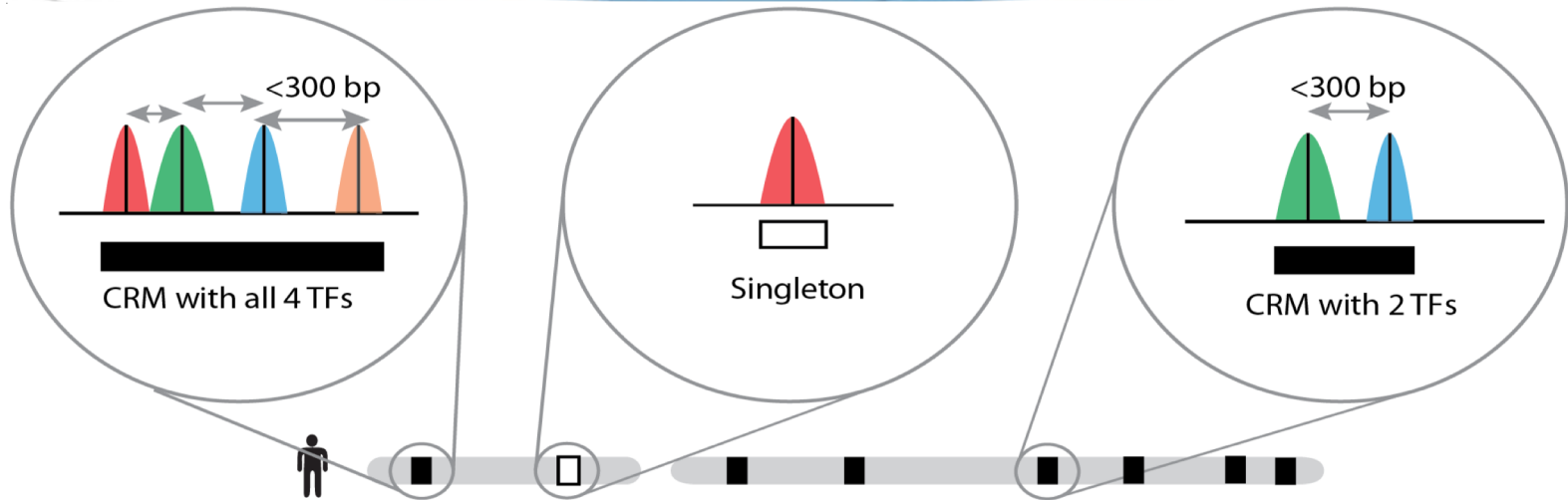
Collaboration:
 Benoit Ballester (INSERM)
 Duncan Odom (Cambridge U)
 Paul Flicek (EBI)



Regulator	Mouse knockout phenotype	Liver phenotype	TF PFAM Class
FoxA1	No gastrulation	no liver formation	Forkhead
HNF4a	No gastrulation	disrupted liver gene expr.	Nuclear Receptor
ONECUT1	Poor hepatocyte, b-cell proliferation	no gall bladder	Onecut domain
CEBPA	Disturbed liver architecture	No glycogen storage	bZIP

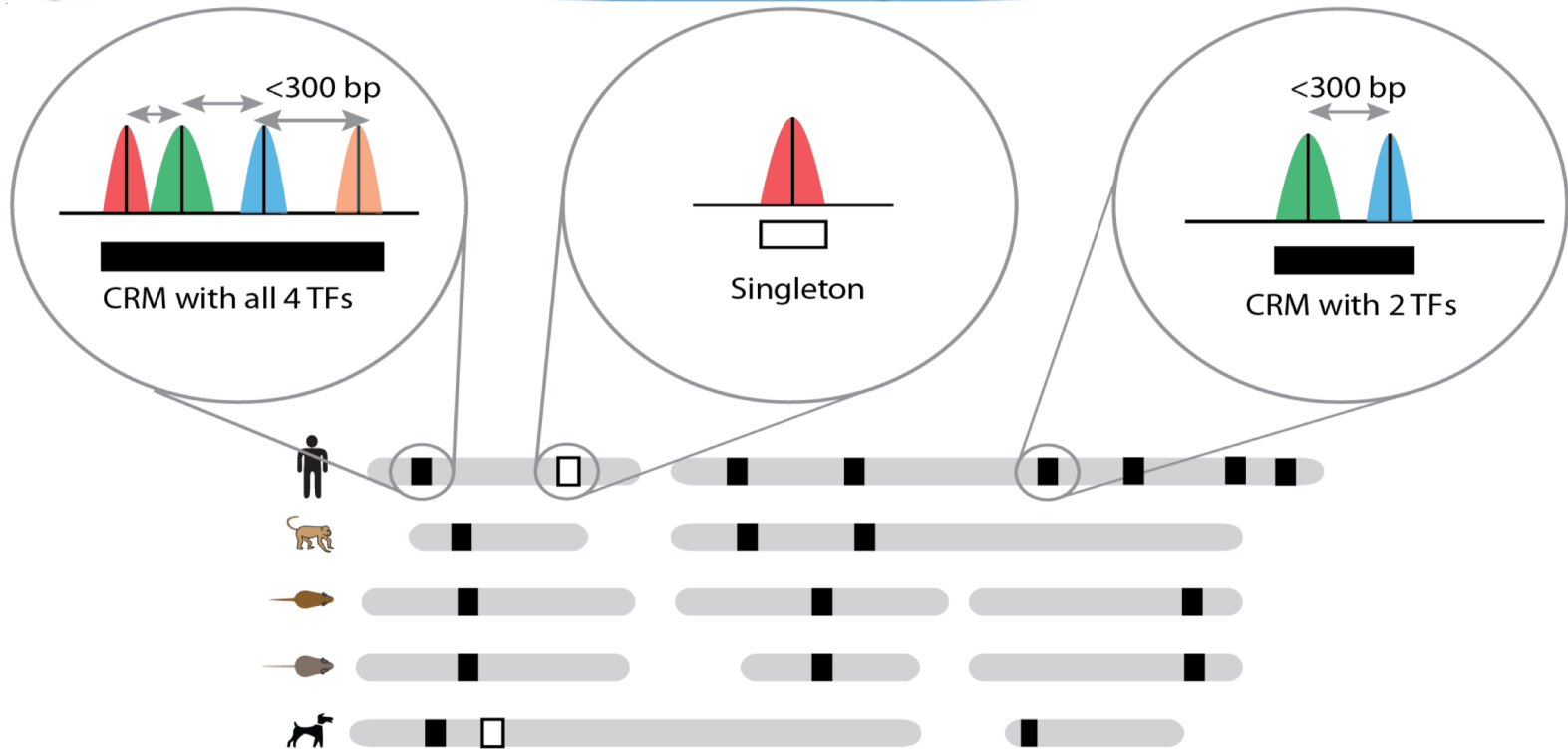
Transcription factors cluster together to regulate genes

CRMs



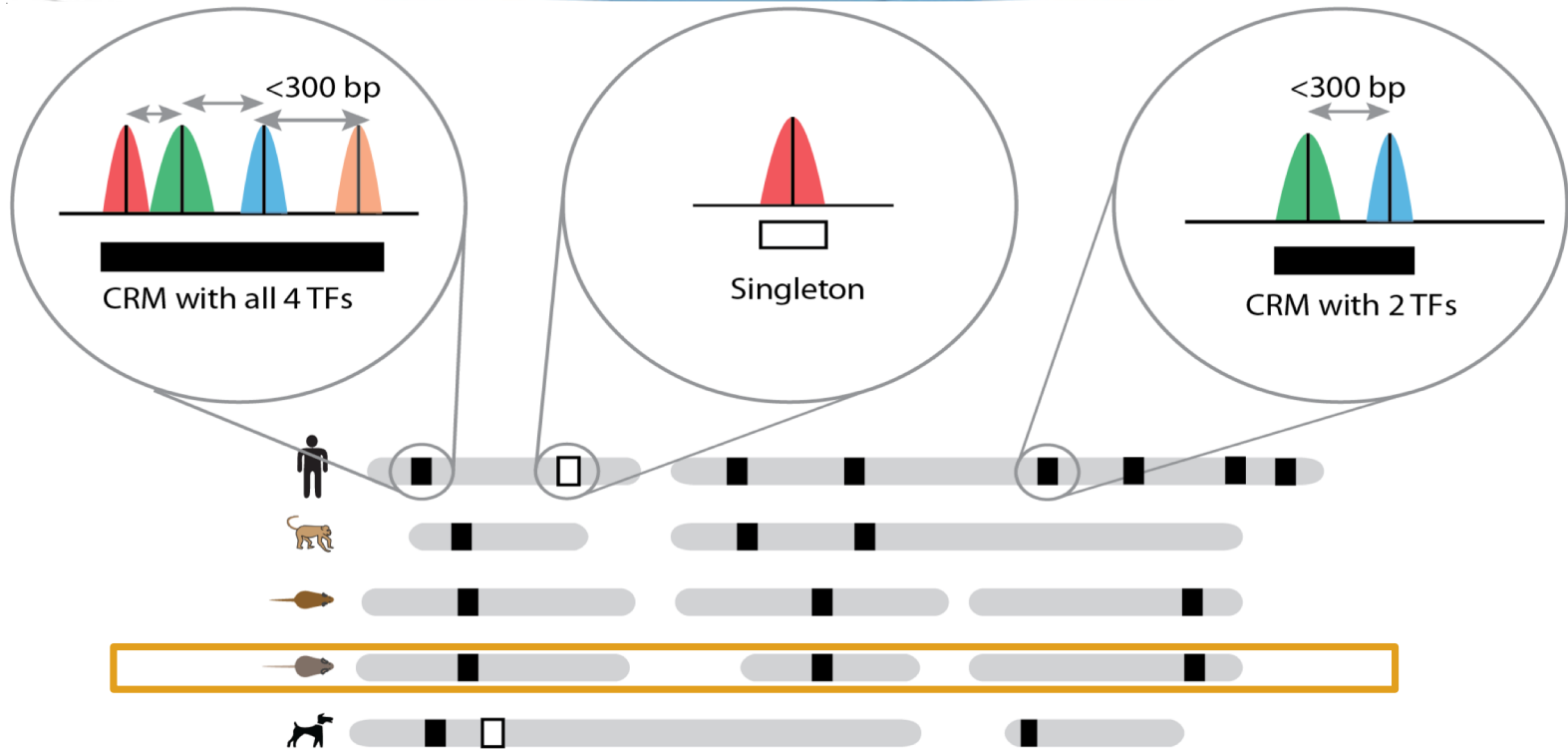
Transcription factors cluster together to regulate genes

CRMs



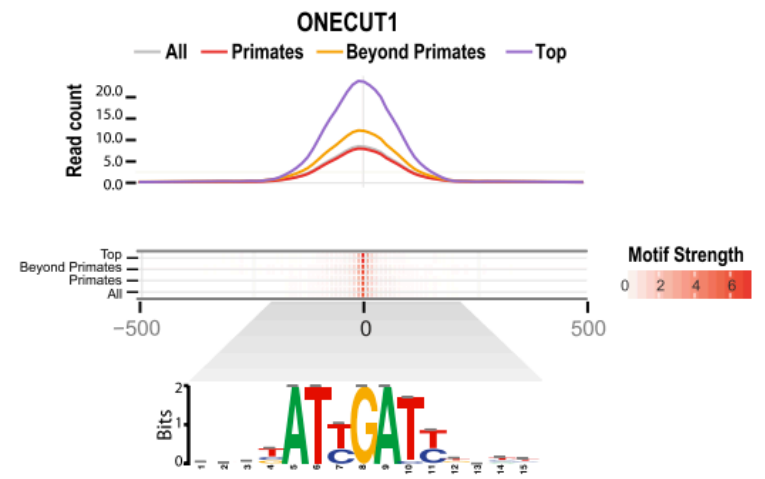
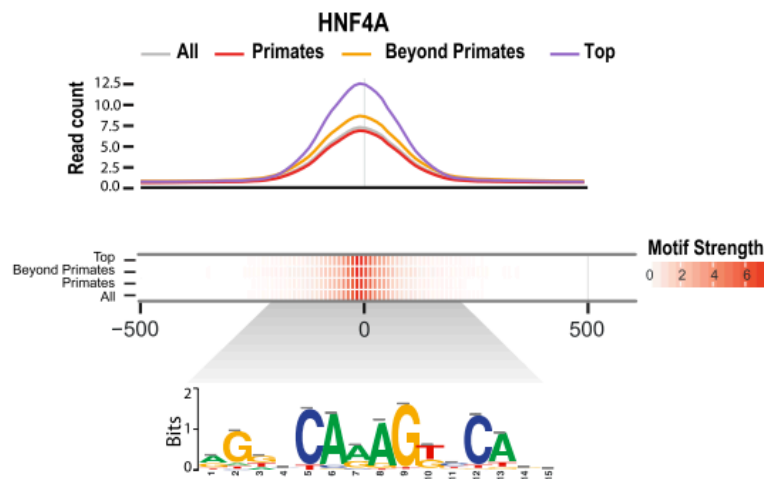
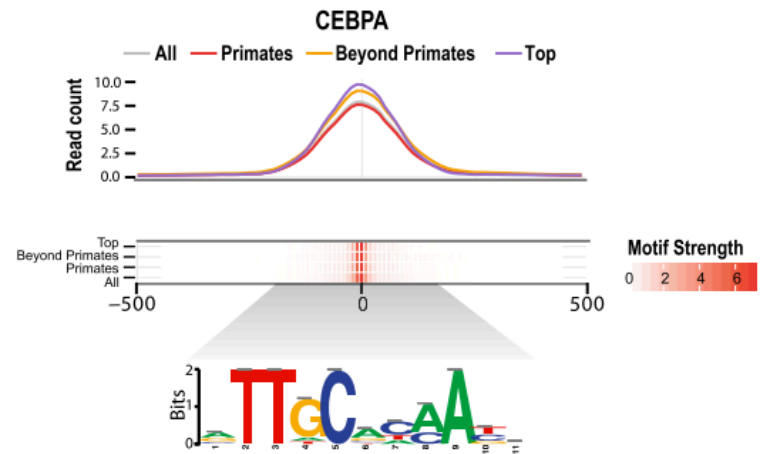
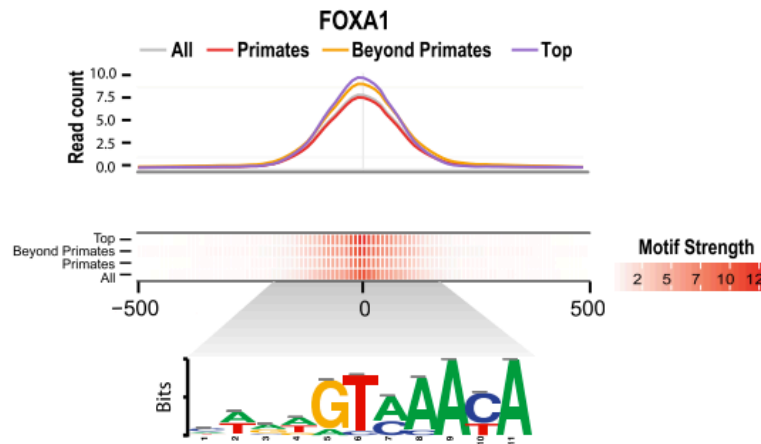
Transcription factors cluster together to regulate genes

CRMs



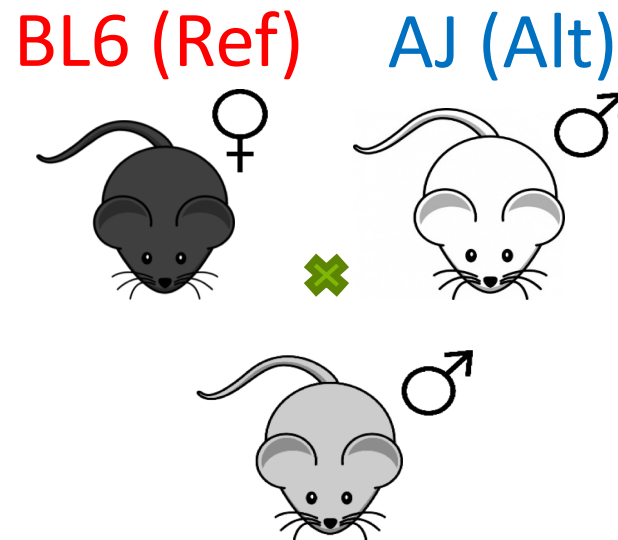
Transcription factors cluster together to regulate genes

CRMs



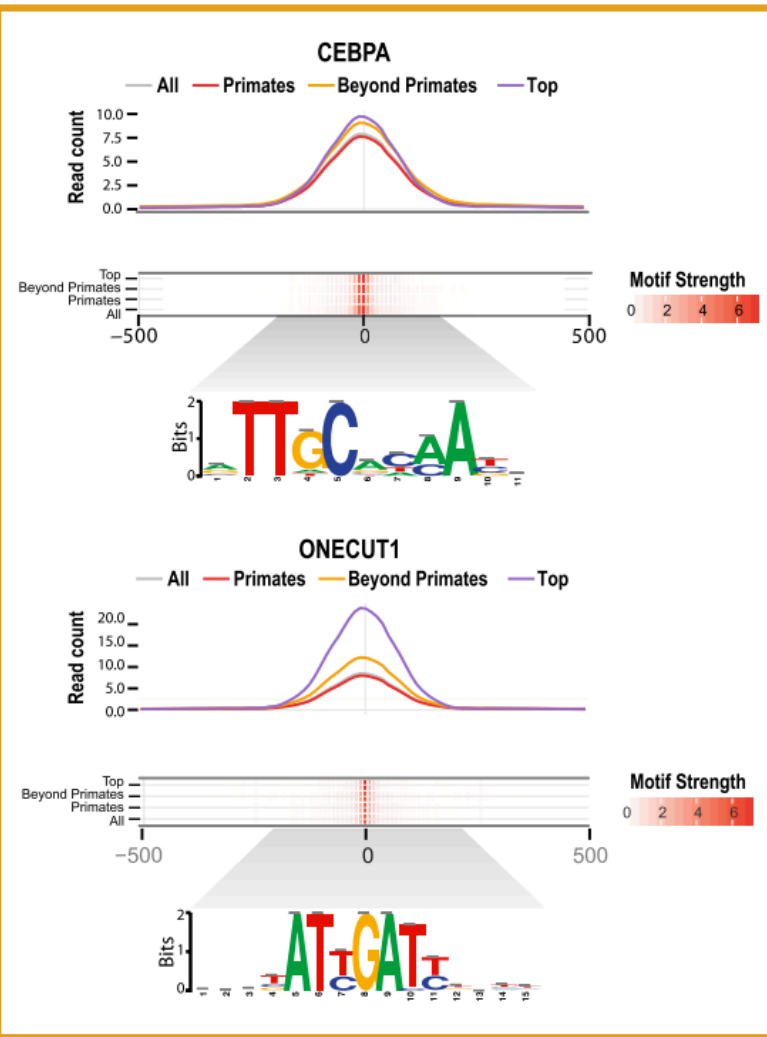
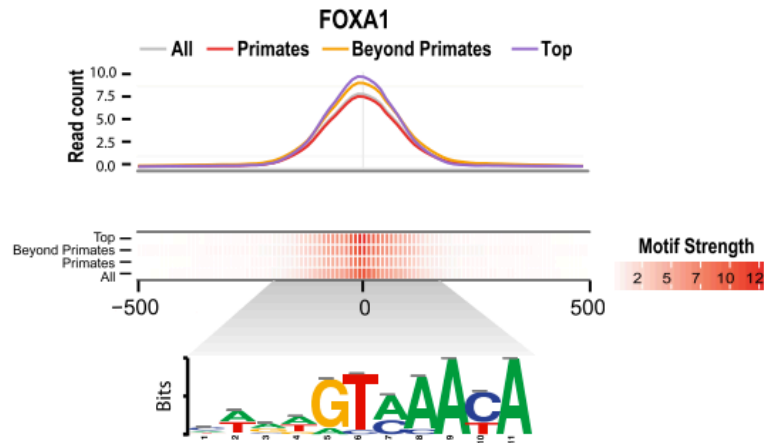
Heterozygous ChIP-seq: Can we detect co-factor effects?

- F1 mice with inbred parents



HNF6
CEBPA

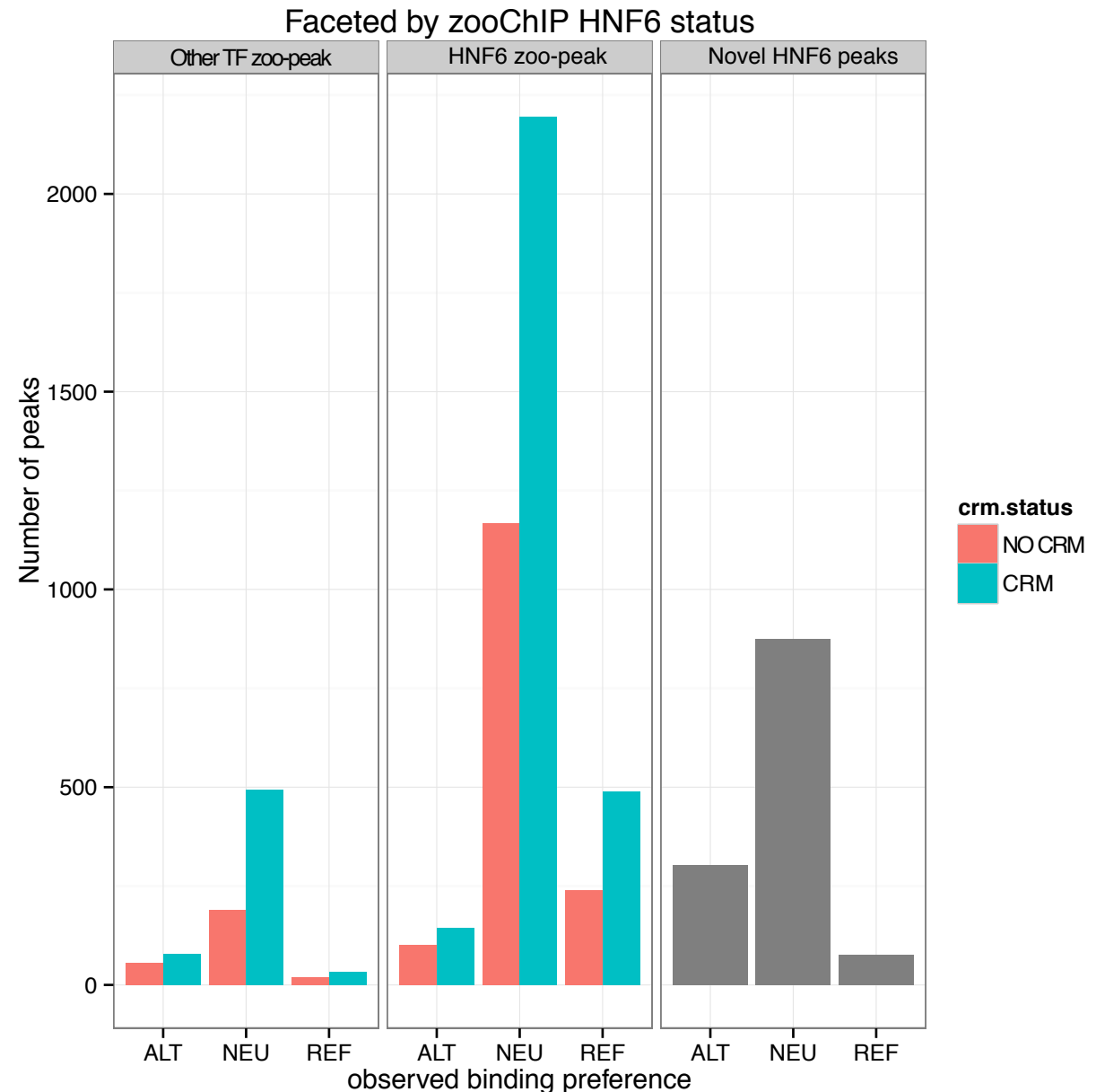
Transcription factors cluster together to regulate genes CRMs



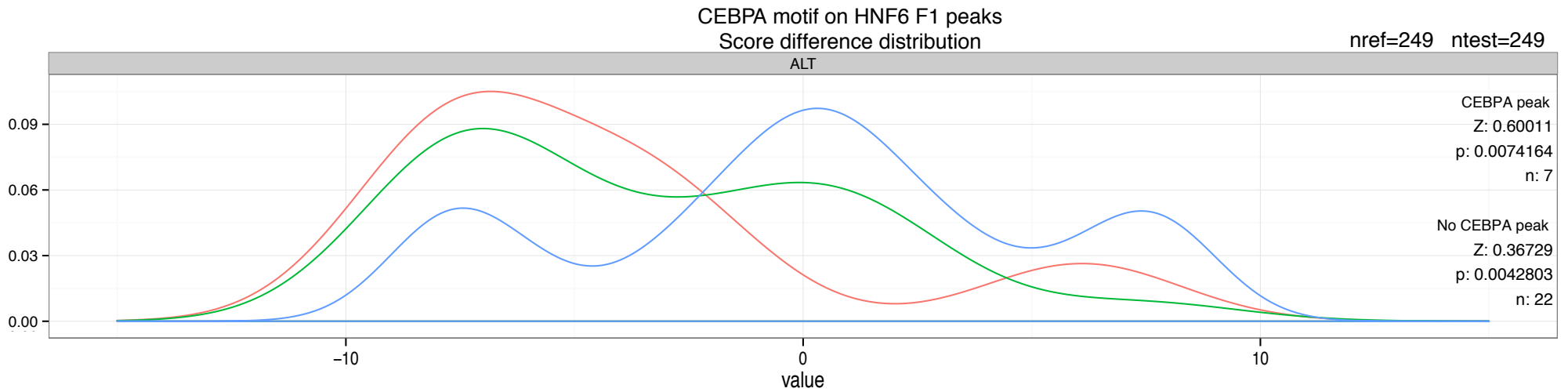
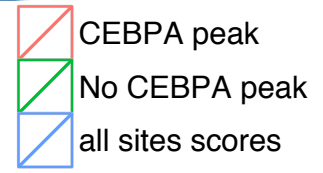
Heterozygous ChIP-seq: Can we detect co-factor effects?

HNF6 F1 data set

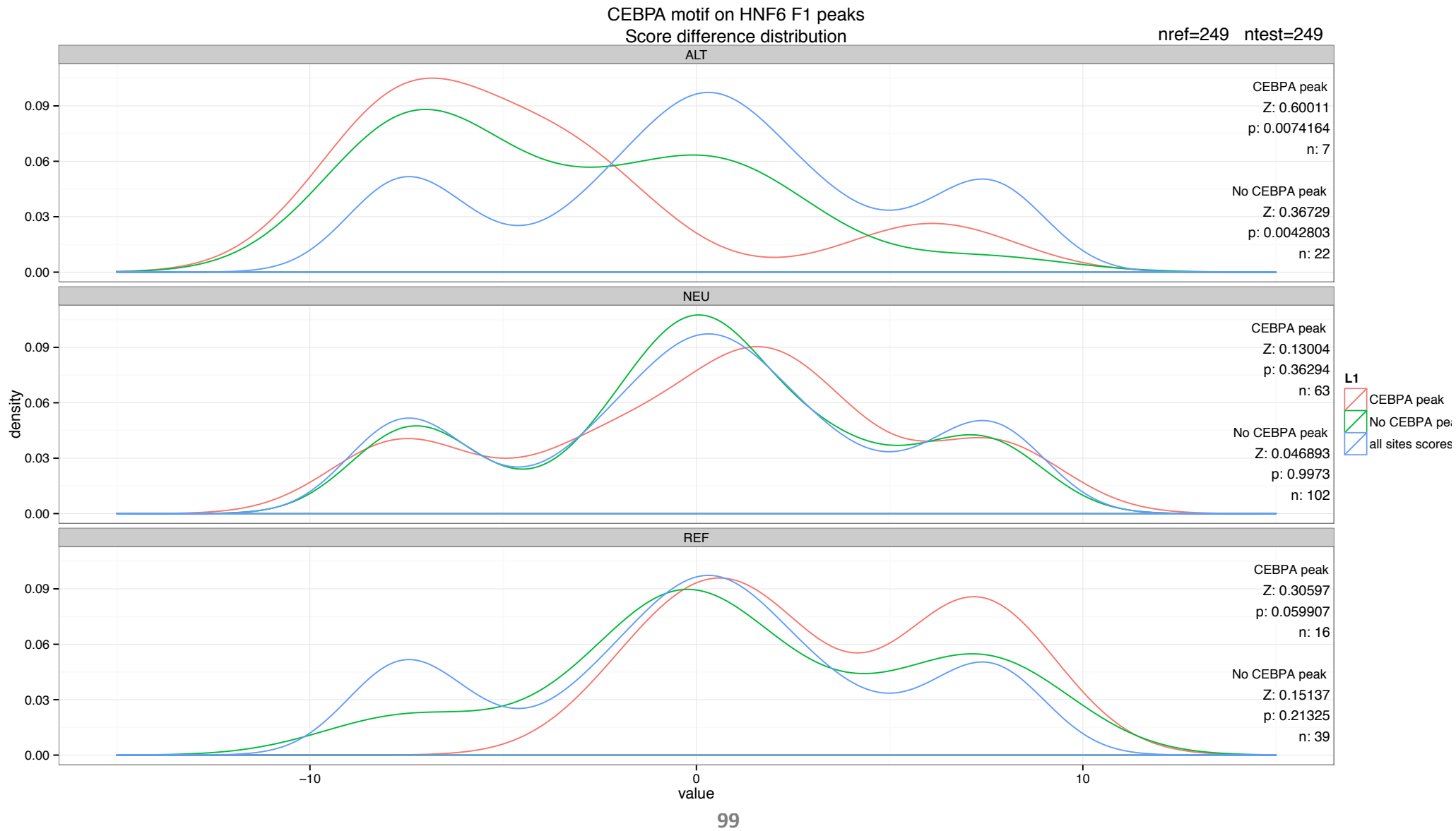
- 6458 Peaks with SNPs
- 2123 sites are novel HNF6 binding.
- 3433 overlap an zoo-Chip HNF6.
- 868 are gains at existing binding sites for other TFs.
- 1253 represent novel binding sites that don't overlap any peak from the Zoo-Chip mouse data.



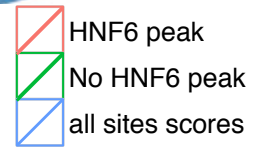
Heterozygous ChIP-seq: Can we detect co-factor effects?



Heterozygous ChIP-seq: Can we detect co-factor effects?

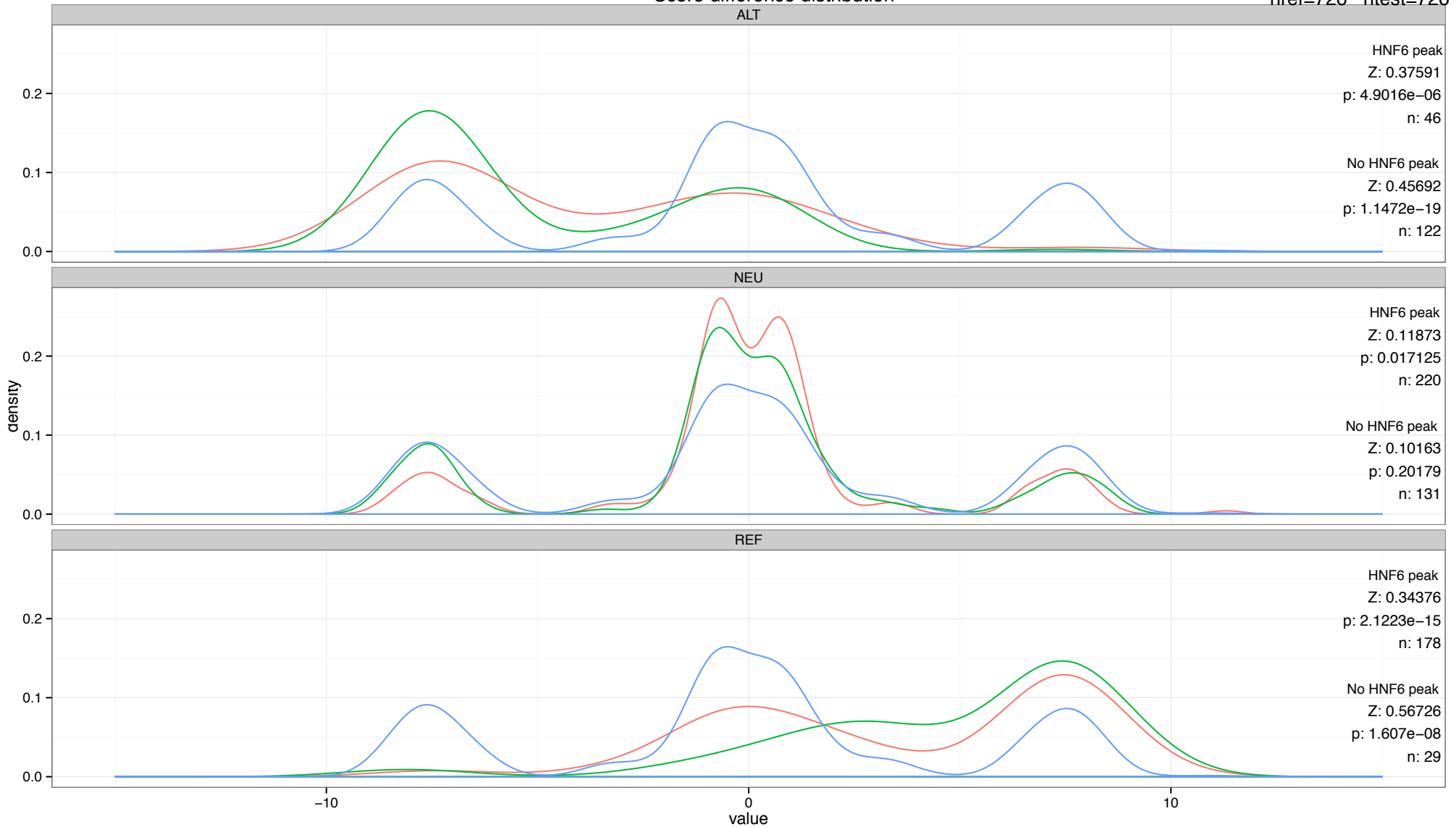


Heterozygous ChIP-seq: Can we detect co-factor effects?

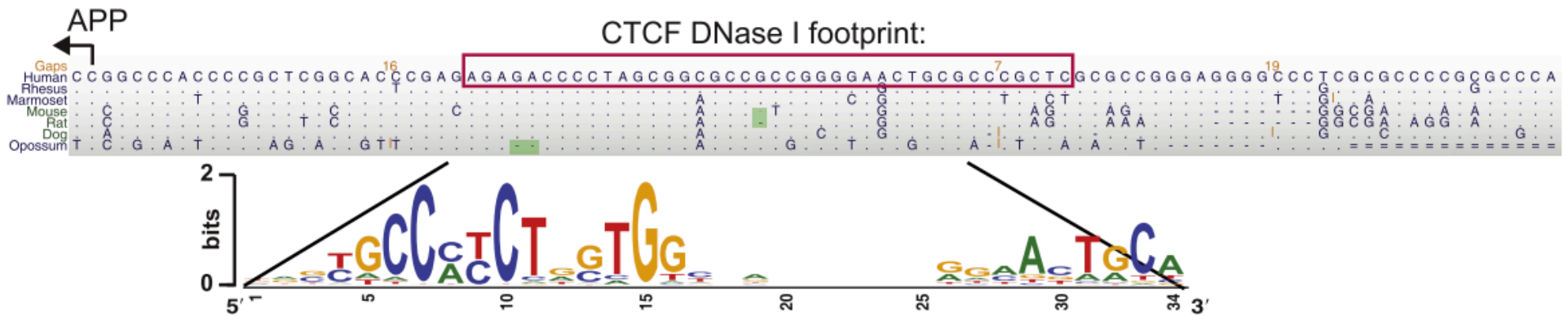


HNF6 motif on HNF6 F1 peaks
Score difference distribution

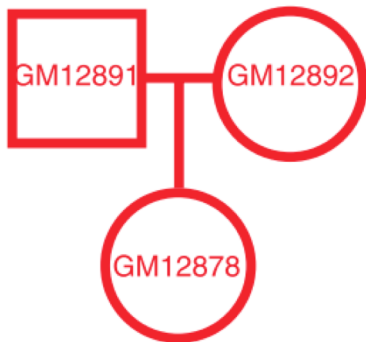
nref=726 ntest=726



Heterozygous ChIP-seq: CTCF trio



CEPH (CEU) Trio



Can we detect allele effects on secondary motifs?

Schmidt, D., Schwalie, P. C., Wilson, M., Ballester, B., Gonçalves, A. et al (2012). Waves of retrotransposon expansion remodel genome organization and CTCF binding in multiple mammalian lineages. CELL