

# Functional Genomics

Nov 29, 2017  
Sarah Gagliano  
Sarah.Gagliano@umich.edu

Biostat666: Statistical methods in human genetics

# Goals for this lecture

- Provide an overview of functional genomics
- Understand its importance/uses in the context of GWAS

# Outline

- Functional Genomics 101
- Using functional genomics to prioritize risk variants

# Outline

- Functional Genomics 101
- Using functional genomics to prioritize risk variants

DNA → mRNA → amino acids chain (protein)

Transcription

Translation

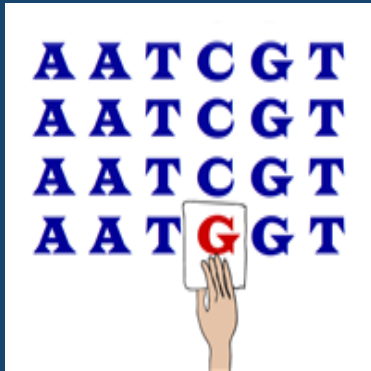
The “Central Dogma”  
of Biology

Coding  
variation



phenotype

Genetic variation → amino acid → protein

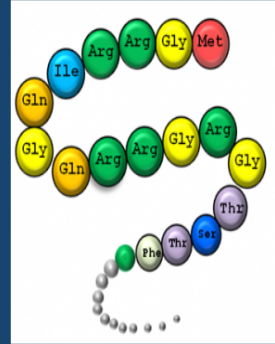
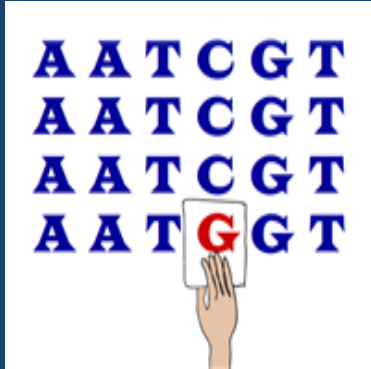


Coding  
variation



phenotype

Genetic variation → amino acid → protein

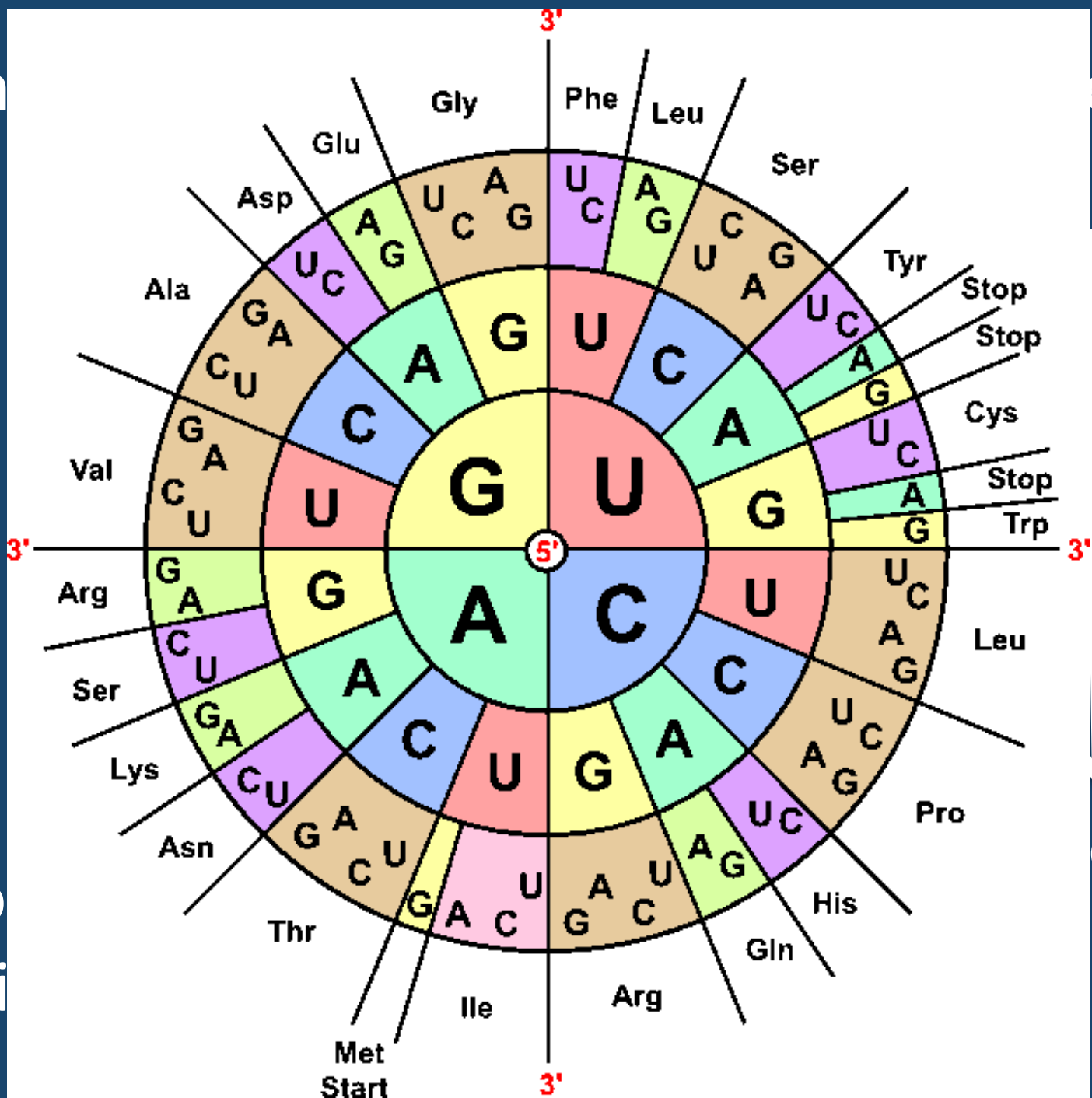


phenotype

Coding  
variation

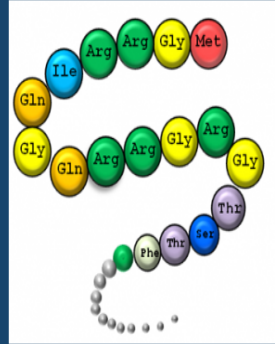
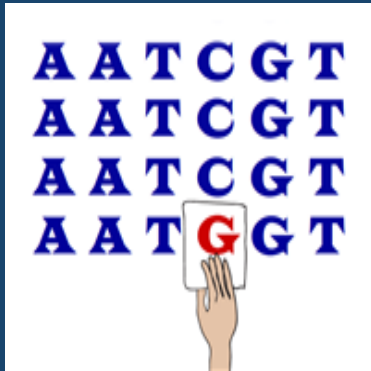
Co  
vari

otype





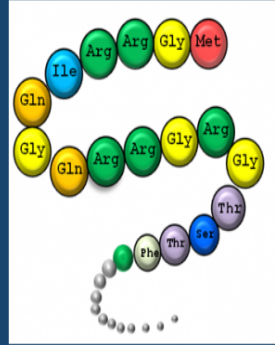
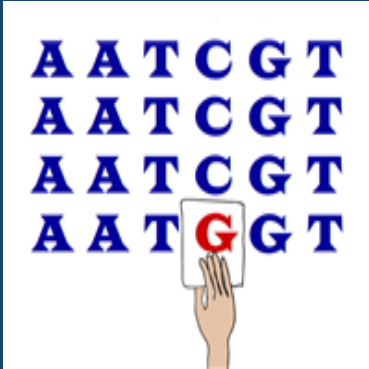
Genetic variation → amino acid → protein



phenotype

Coding  
variation

Genetic variation → amino acid → protein

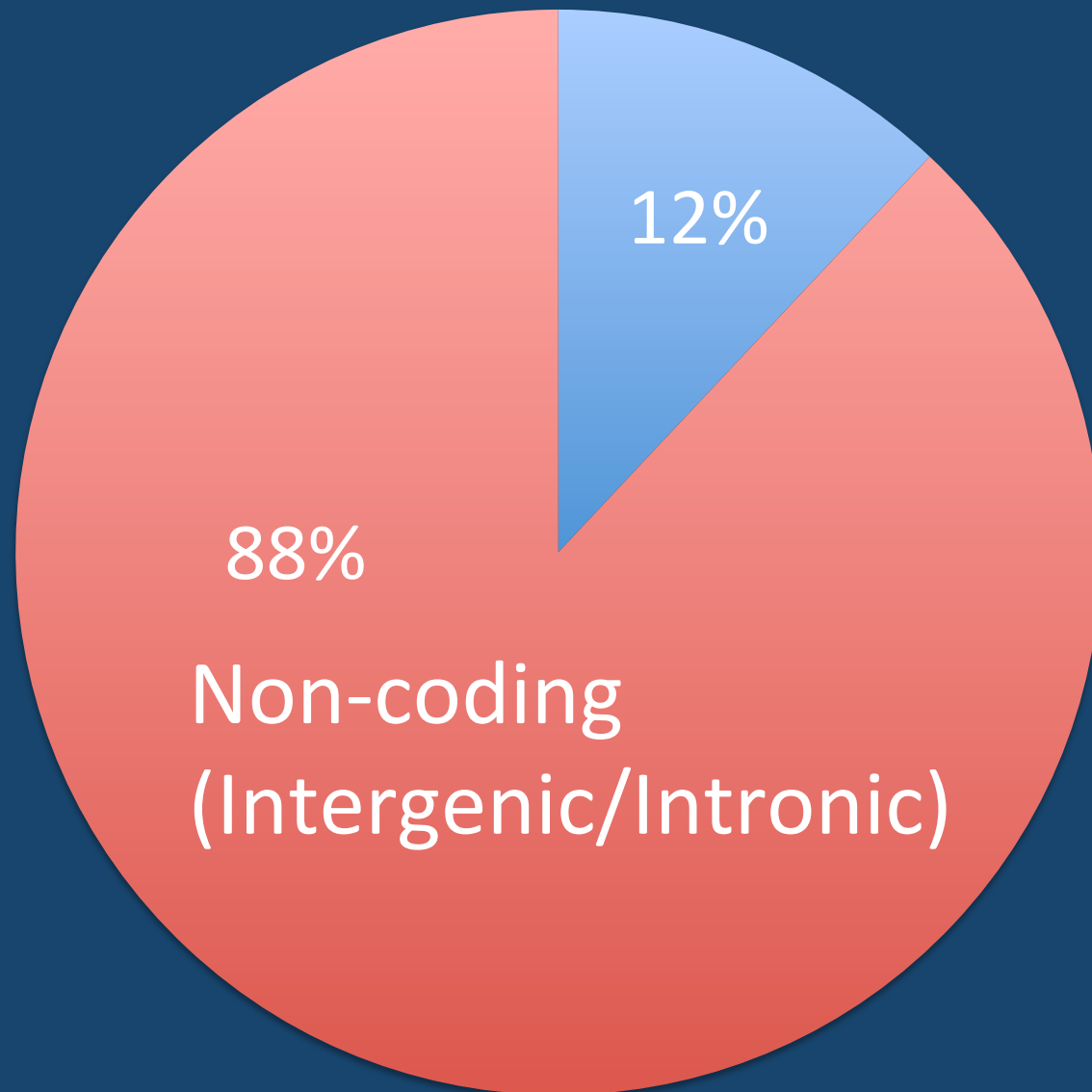


phenotype



Coding  
variation

# Most GWAS hits are non-coding



## Scientists Find 'Junk DNA' Useful After All



*Come on! Help me find it, I think this one controls  
hair loss!*

ENCODE Project creates a  
map of functional regions

Non-coding  
regions are  
not “junk”!

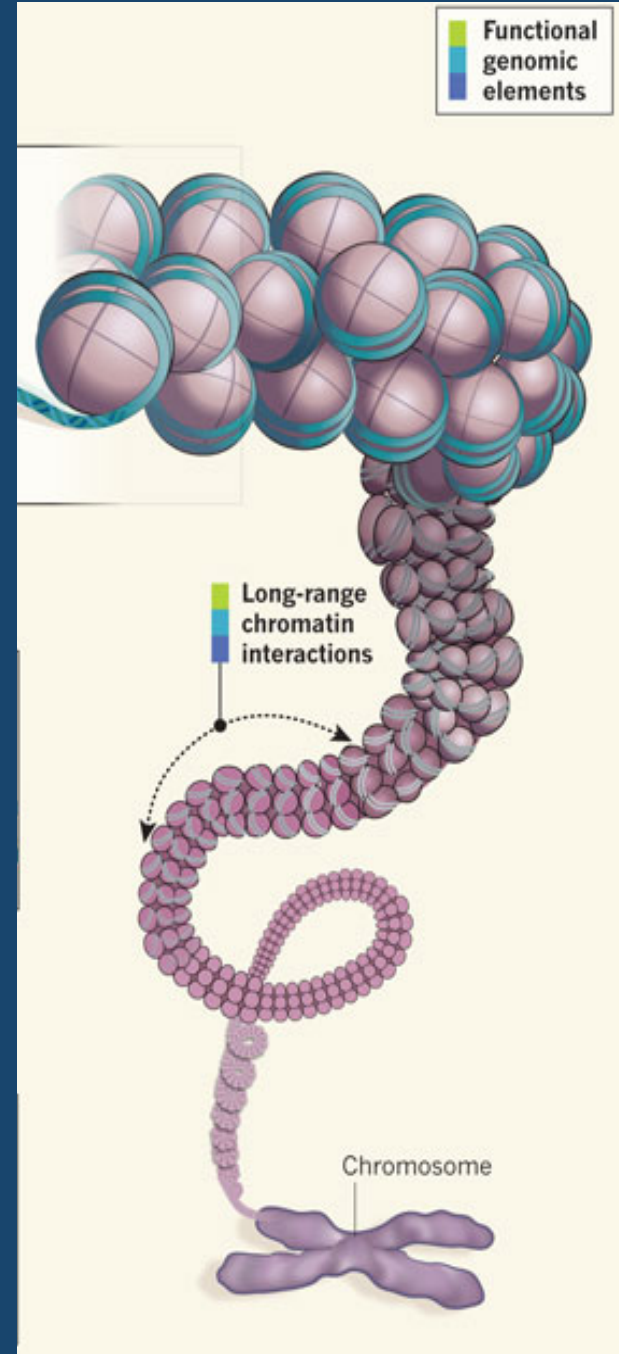


Figure from Ecker et al. (2012) *Nature*, 489: 52-55



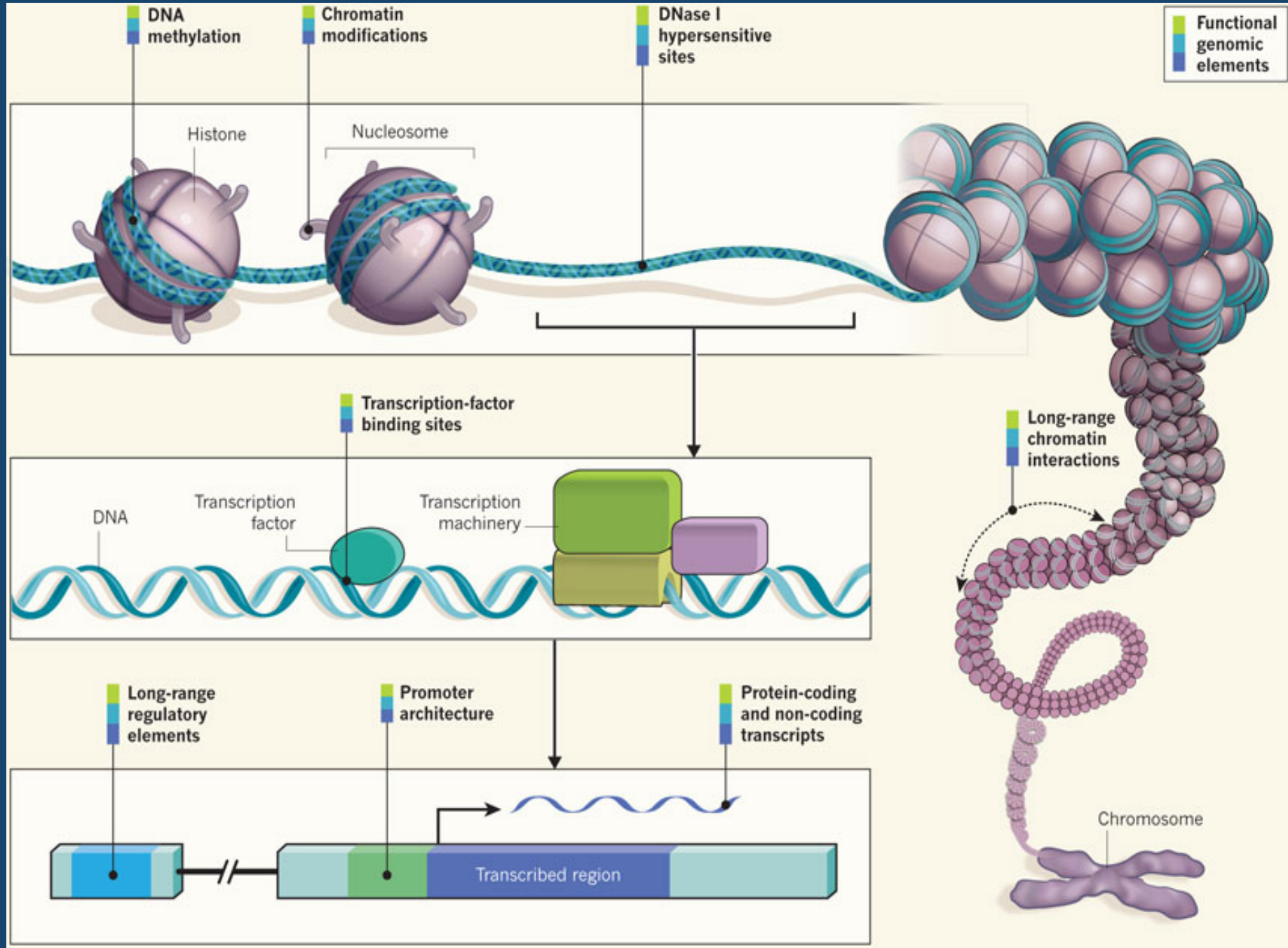


Figure from Ecker et al. (2012) *Nature*, 489: 52-55

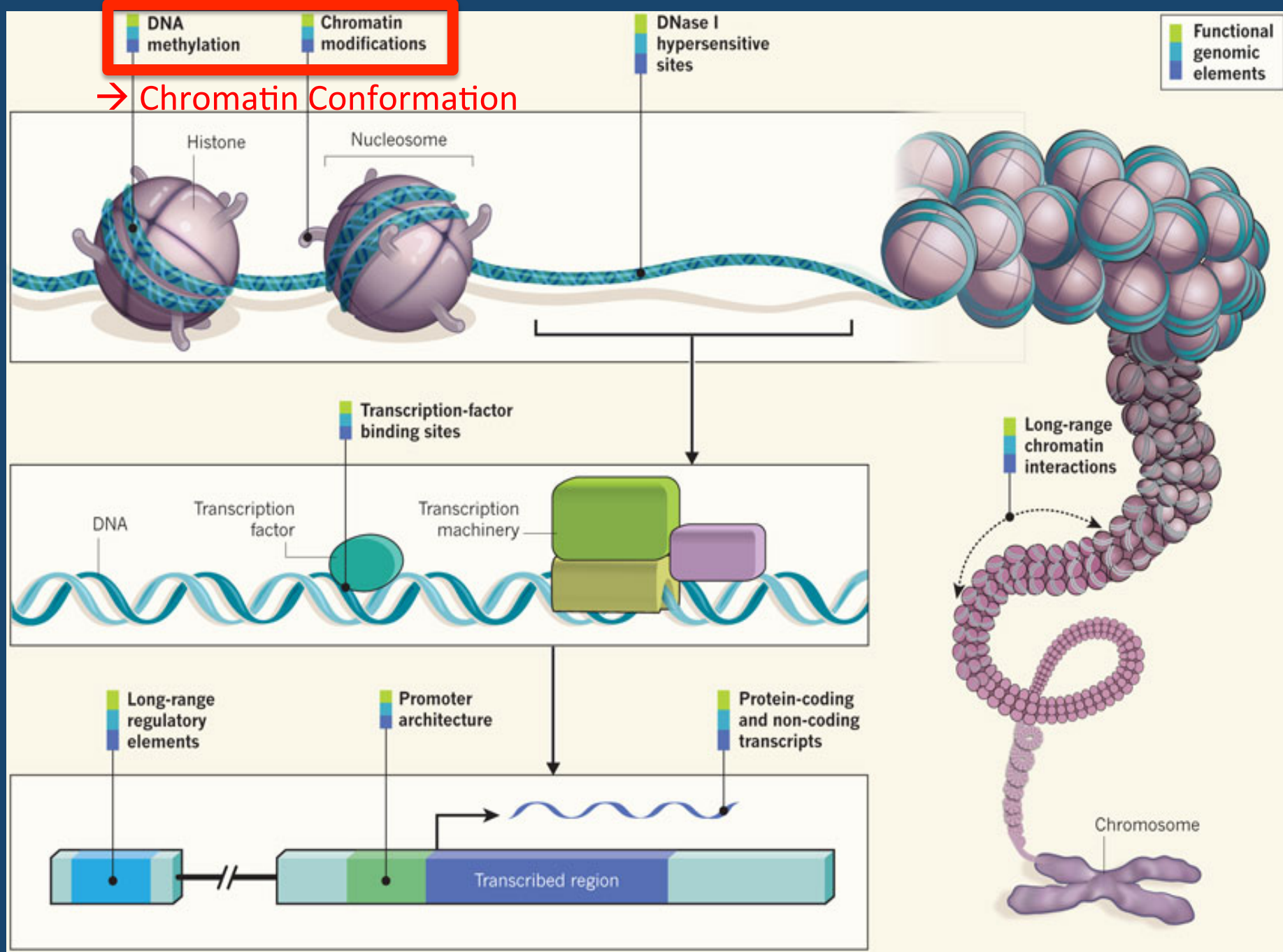
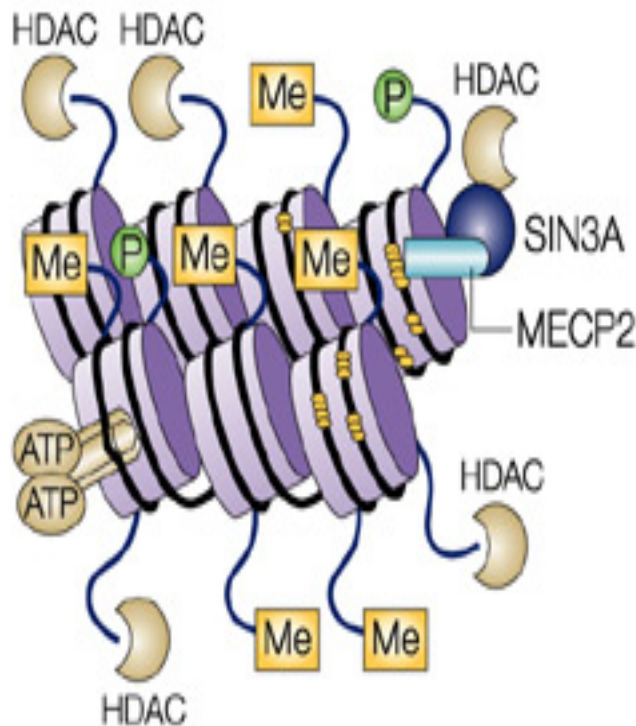


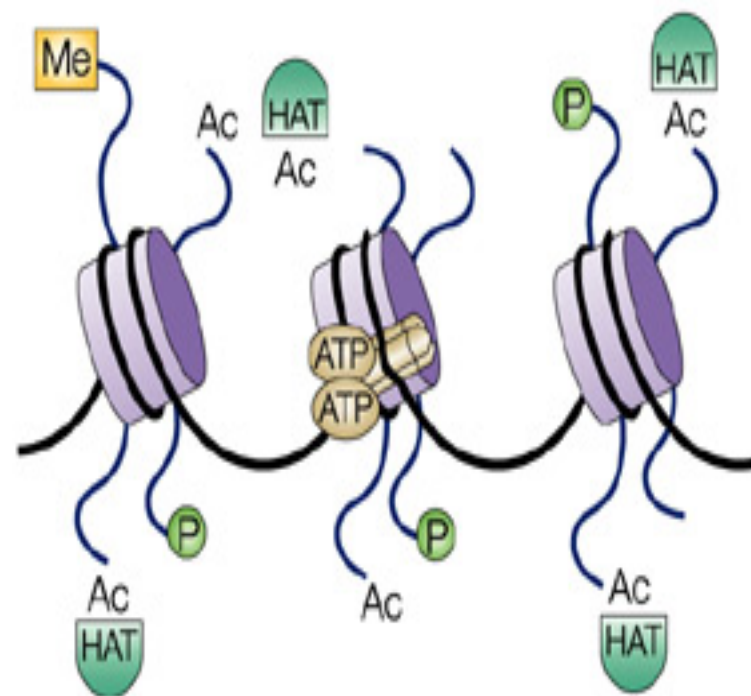
Figure from Ecker et al. (2012) *Nature*, 489: 52-55

# Marks associated with closed vs. open chromatin

**a** Closed chromatin: transcriptional repression



**b** Open chromatin: transcriptional activation



- Developmental time-point
- Tissue
- Sex



# Histone Modifications

# Histone marks define functional regions



chr16 (q12.2) 16p13.3 12.3 12.1 p11.2 16q11.2 q12.1 16q21 22.1 q23.1

Mark of an active enhancer

UCSC Genome Browser

# Enhancer-Promoter Looping

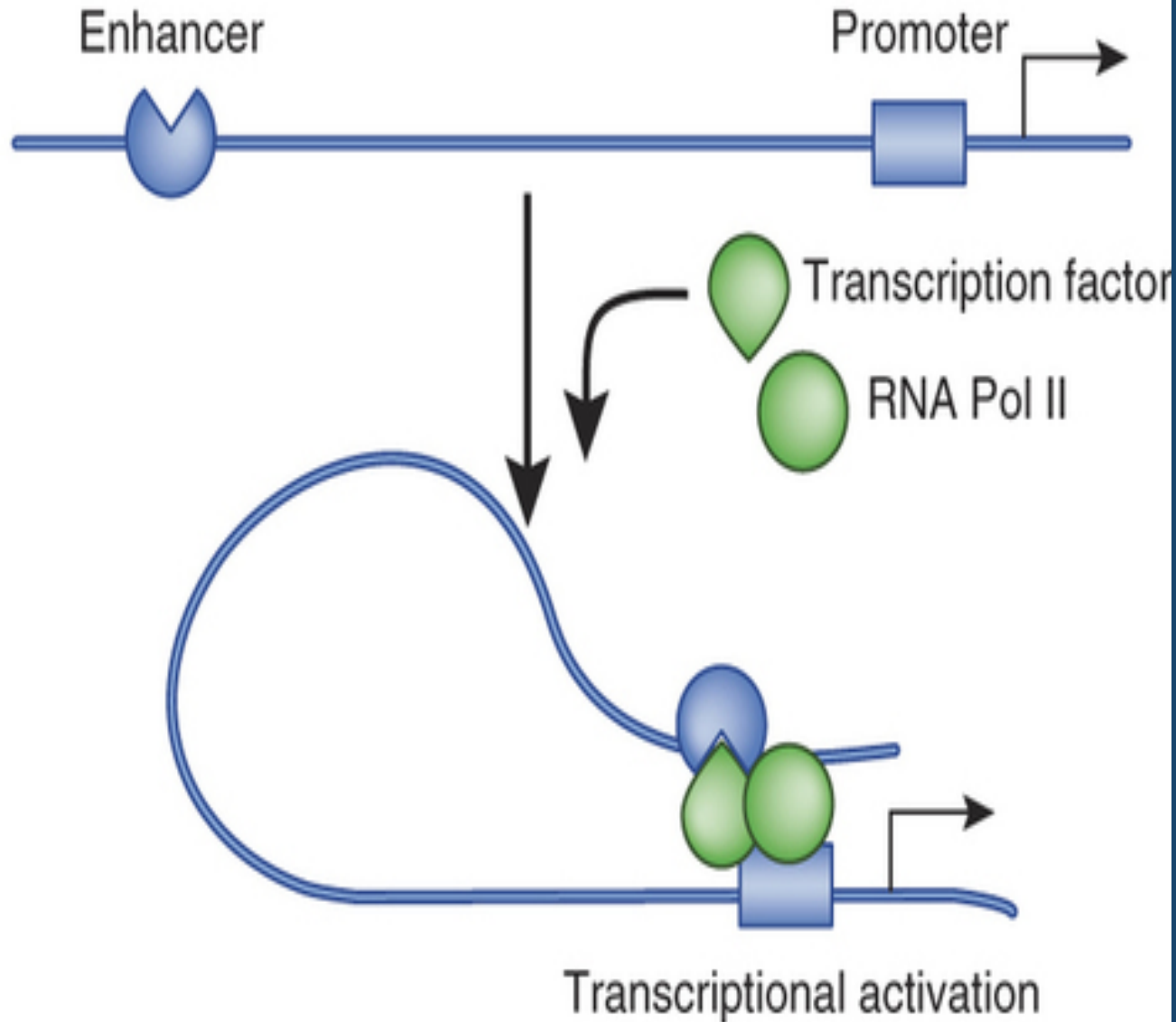


Fig 2b from Cavalli & Misteli (2012) *Nature Structural & Molec Bio*

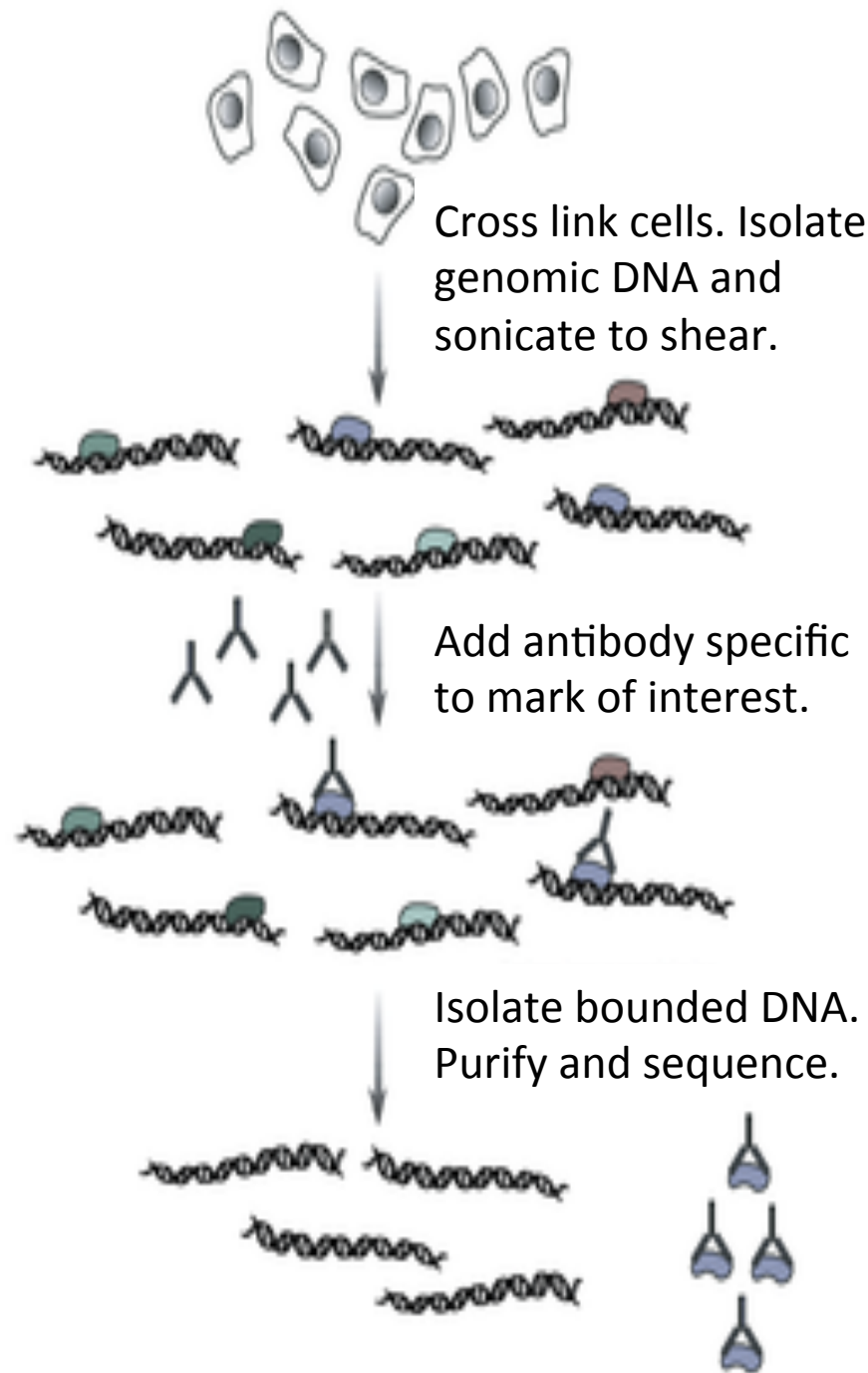
# Detecting Histone Modifications

1. Wet lab

2. Impute

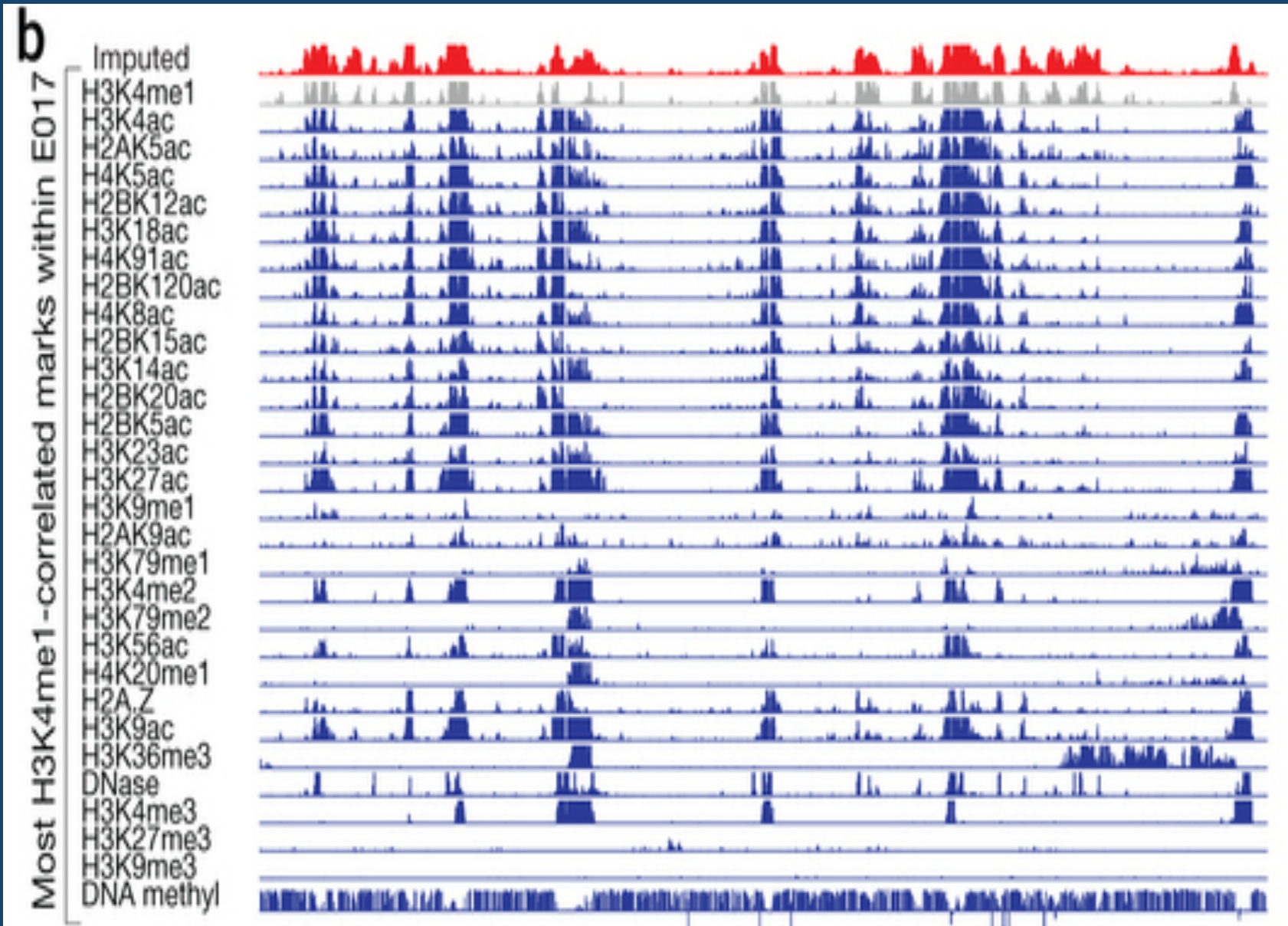
# ChIP-seq to map out histone marks

*Newer:*  
Chromatin  
conformation  
capture  
techniques (e.g.  
Hi-C)



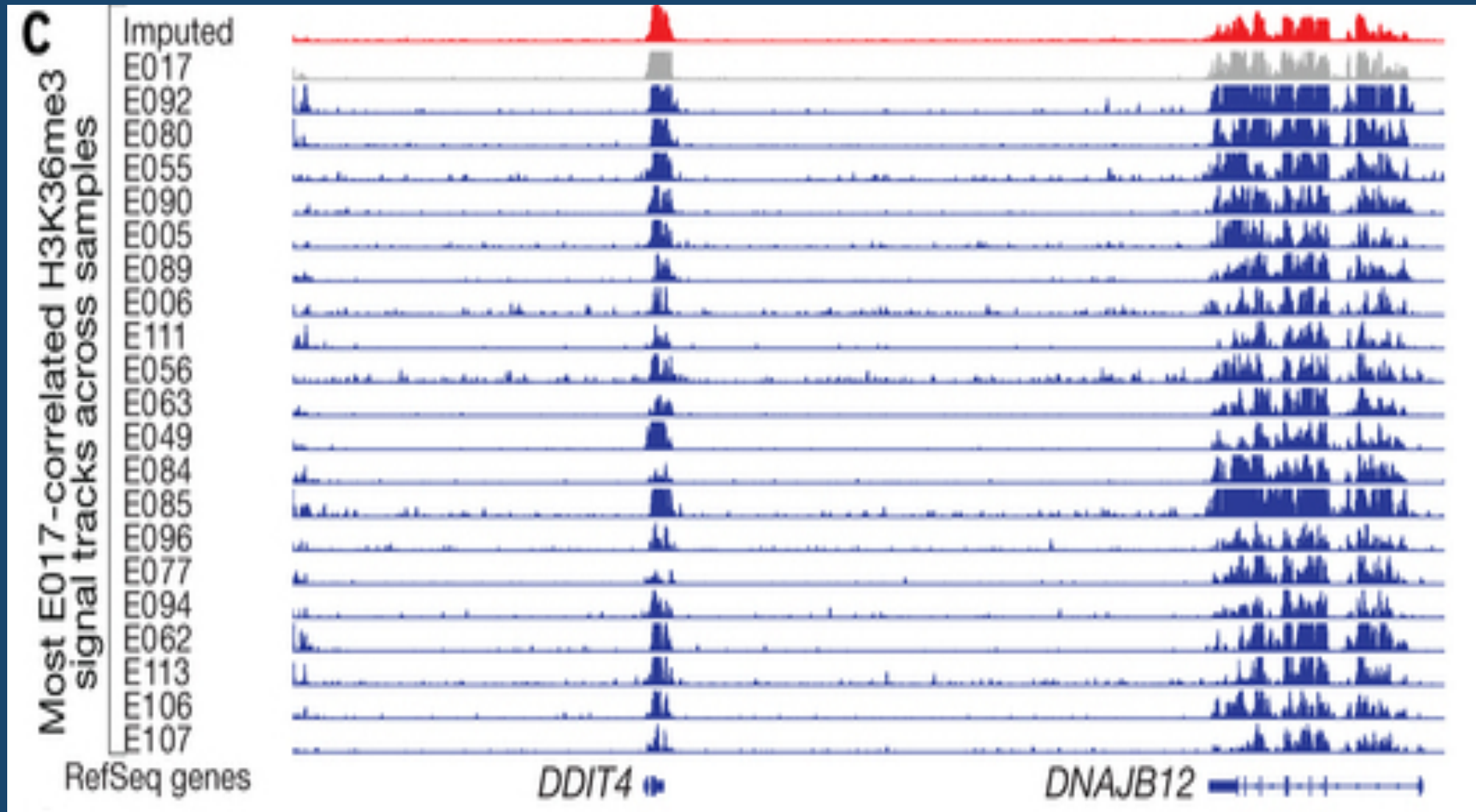
[http://mmg-233-2013-genetics-genomics.wikia.com/wiki/Chromatin\\_Immunoprecipitation\\_\(ChIP\)](http://mmg-233-2013-genetics-genomics.wikia.com/wiki/Chromatin_Immunoprecipitation_(ChIP))

# Correlations across marks per sample

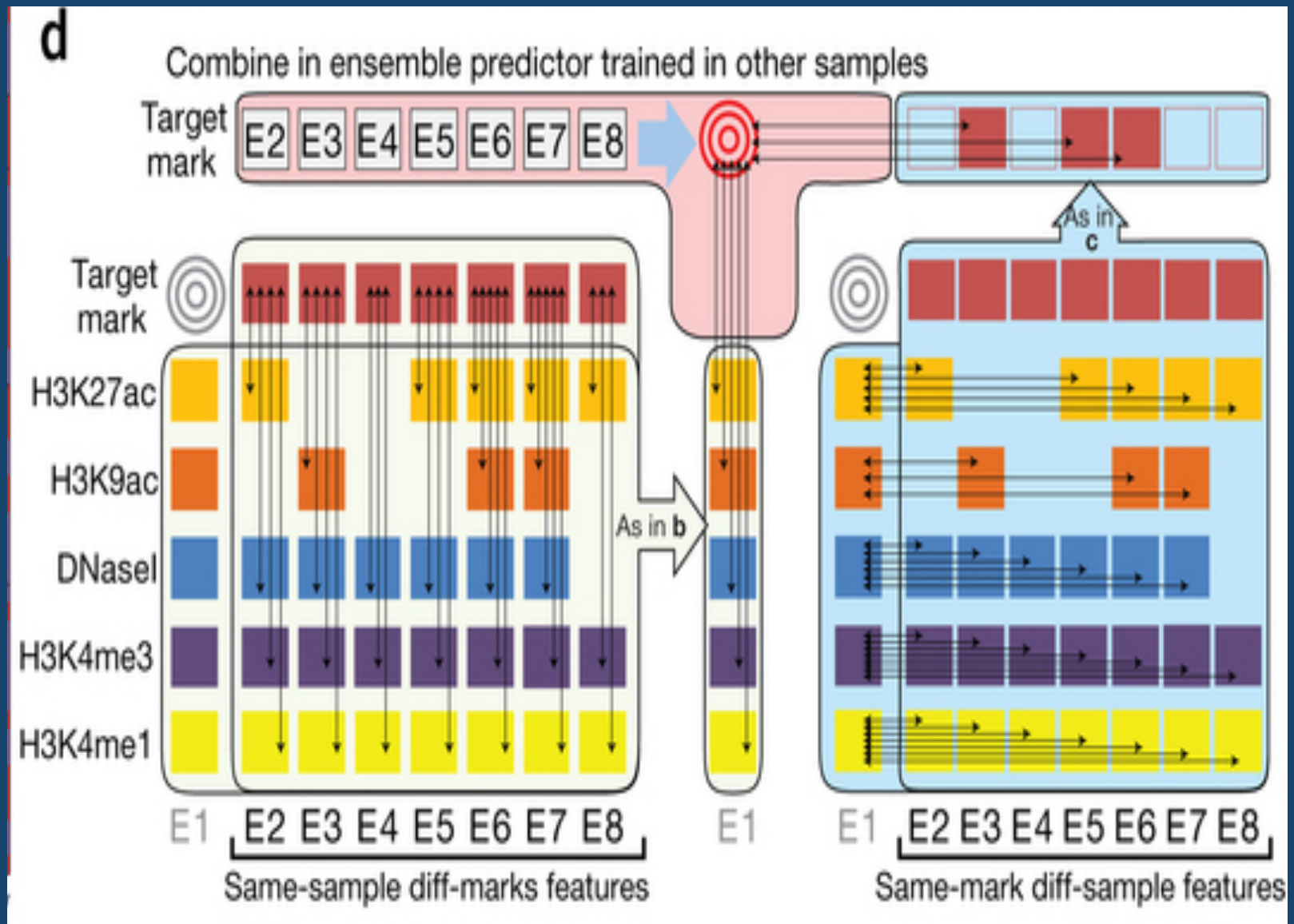




# Correlations across samples per mark



# Use correlations btwn marks & btwn samples





# Access to Histone Modifications

Roadmap Epigenomics Project (lots of tissues)

ENCODE (lots of cell lines)

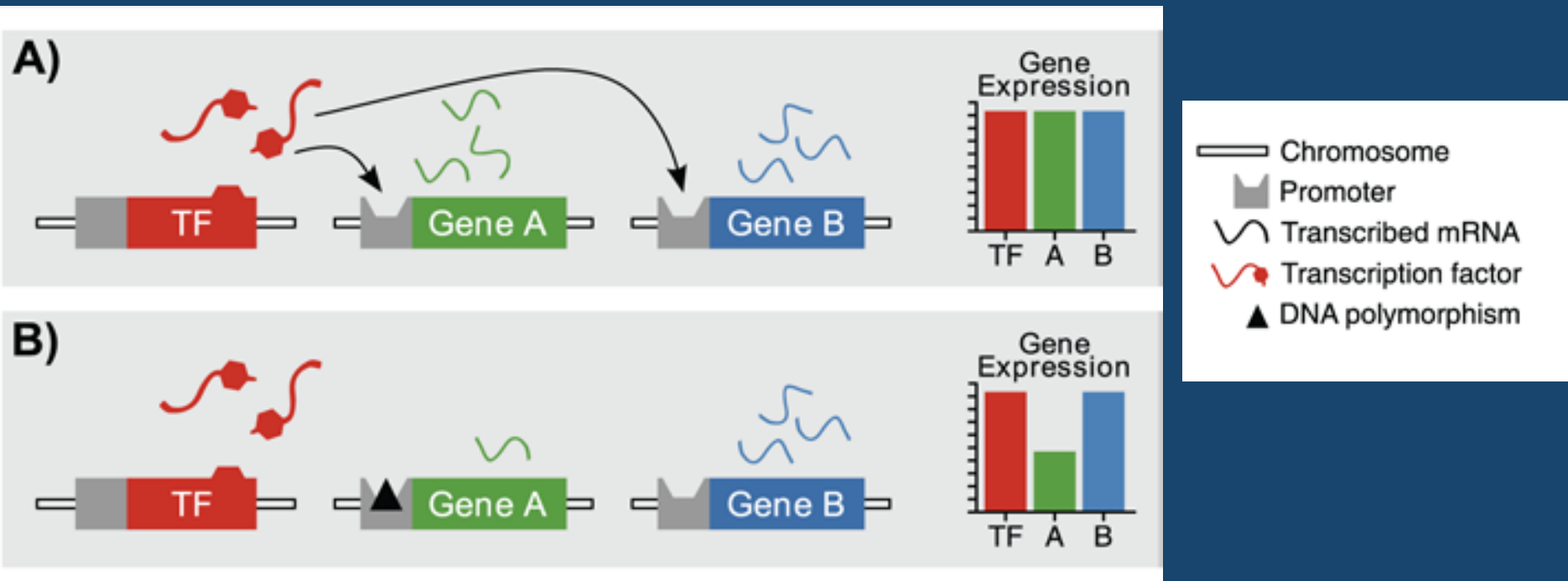
PsychENCODE (brain from control and/or psychiatric samples)

UCSC Genome Browser to visualize

eQTLs

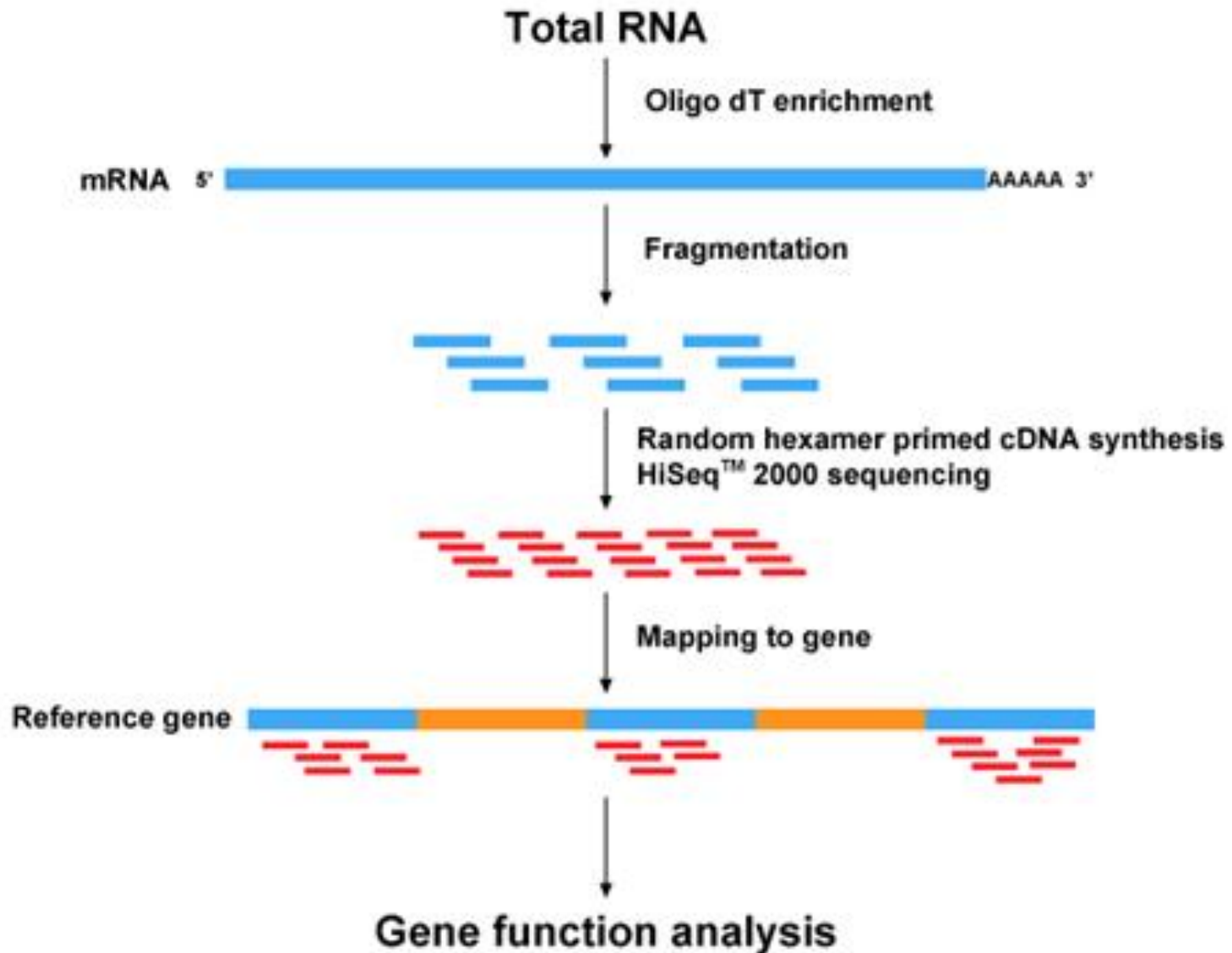
# What is an expression quantitative trait locus?

**Expression quantitative trait loci (eQTLs)** = DNA loci that regulate expression levels of RNAs



Adapted from: Wolen AR, Miles MF. Identifying gene networks underlying the neurobiology of ethanol and alcoholism. *Alcohol Res.* 2012;34(3):306-17.

# RNA-seq to identify eQTLs



# Access to eQTLs

GTEx (44 tissues with eQTLs)

<http://www.gtexportal.org/>

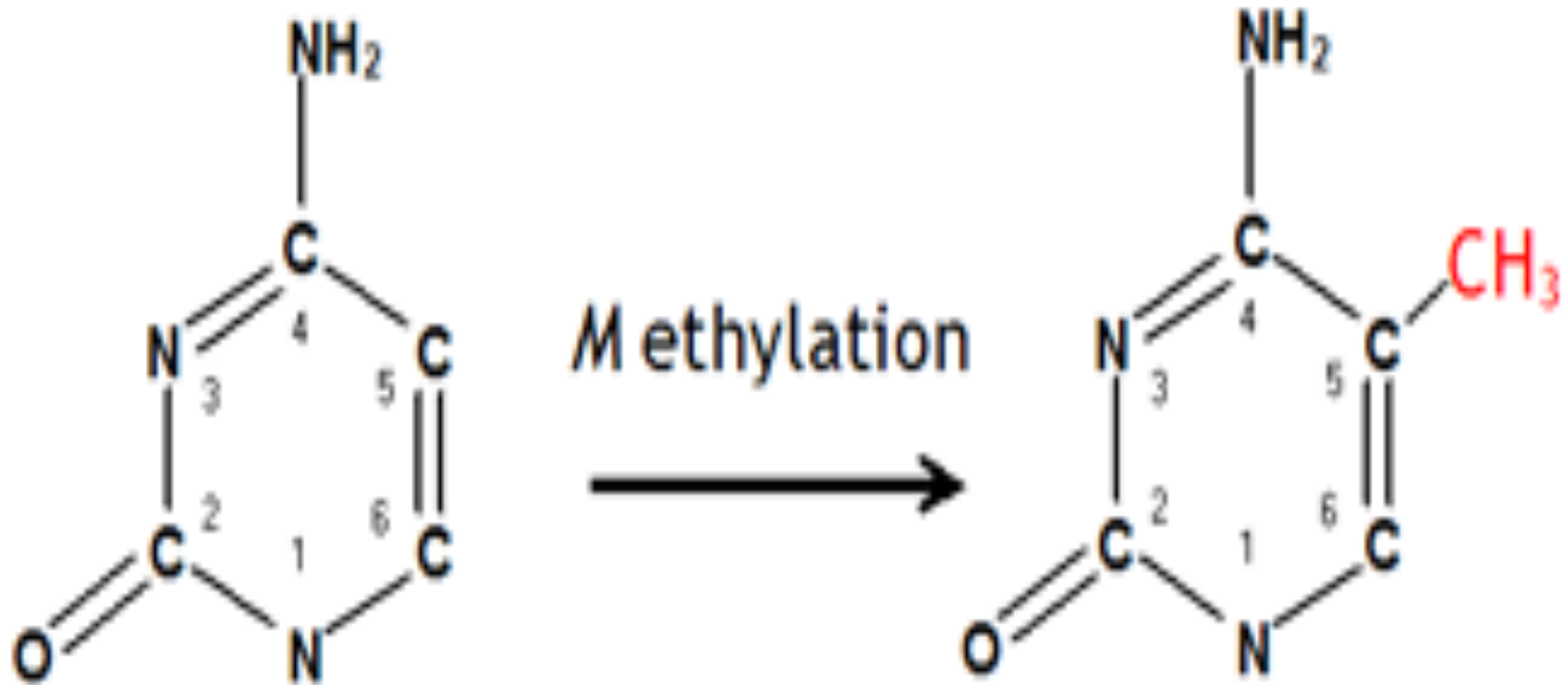
UKBEC (10 brain tissues with eQTLs)

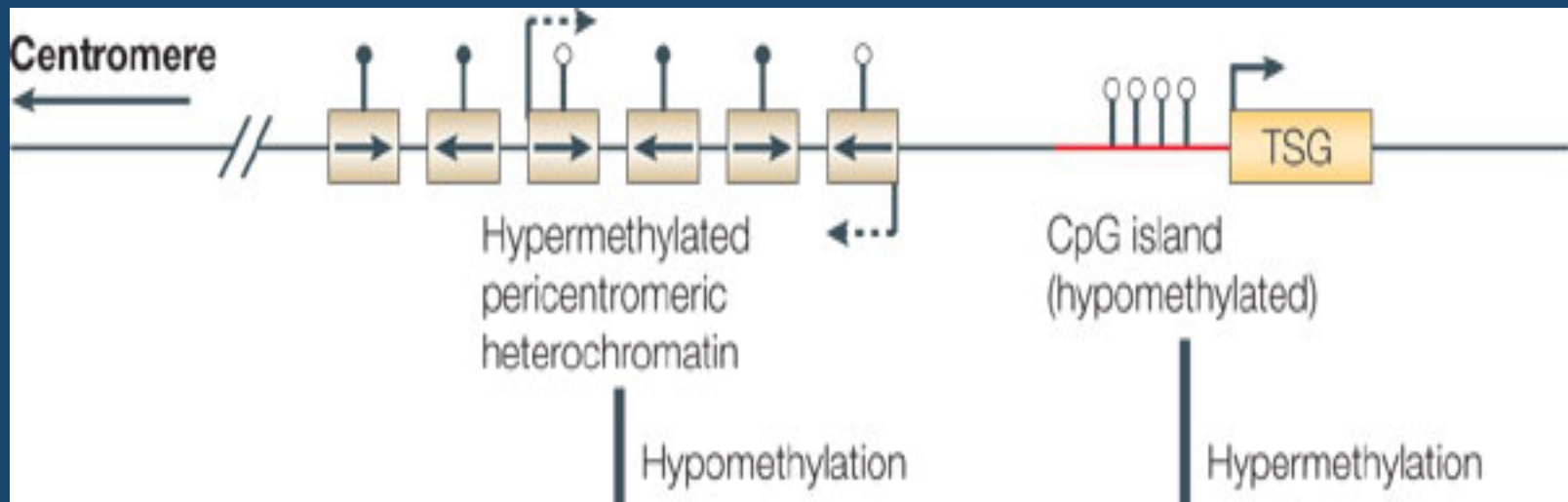
<http://braineac.org/>

Individual studies

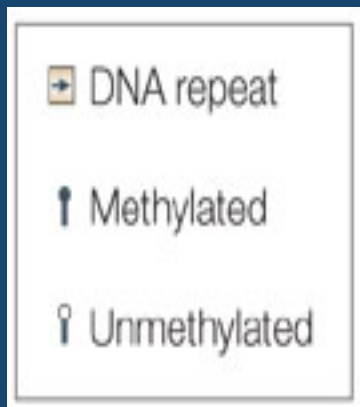
# DNA methylation

# 5-methyl-Cytosine

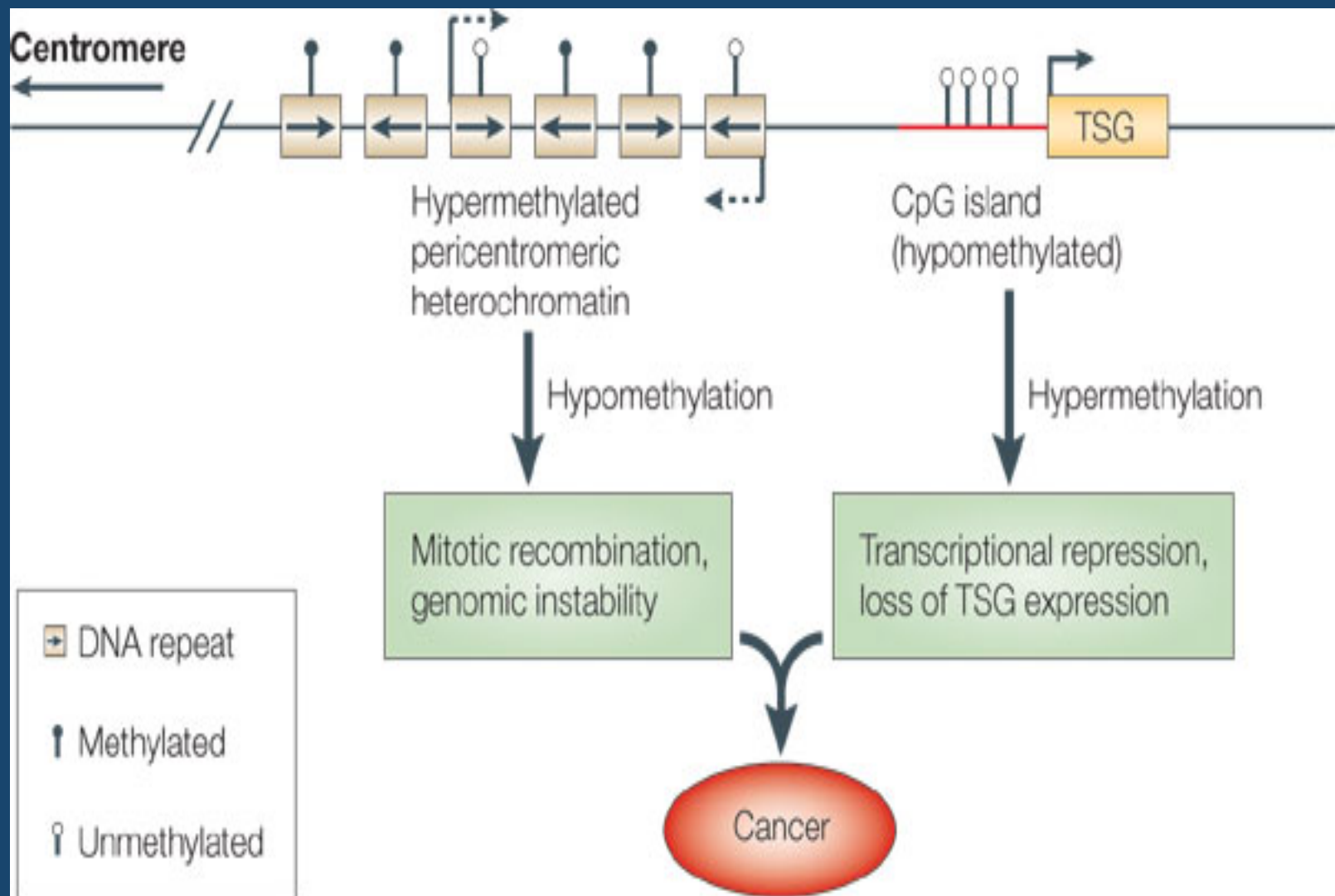




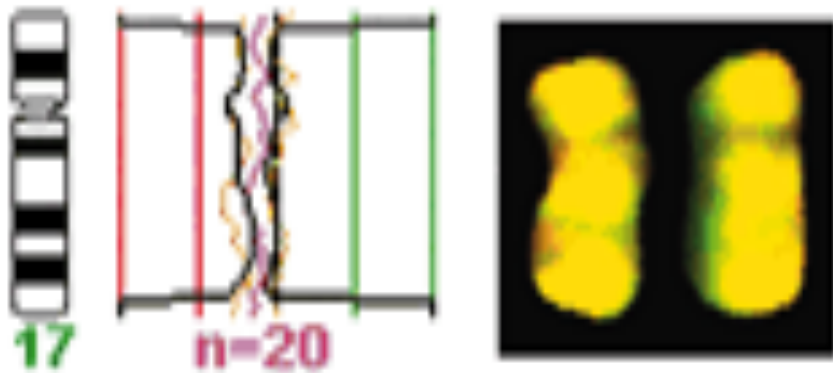
TSG= tumor suppressor gene





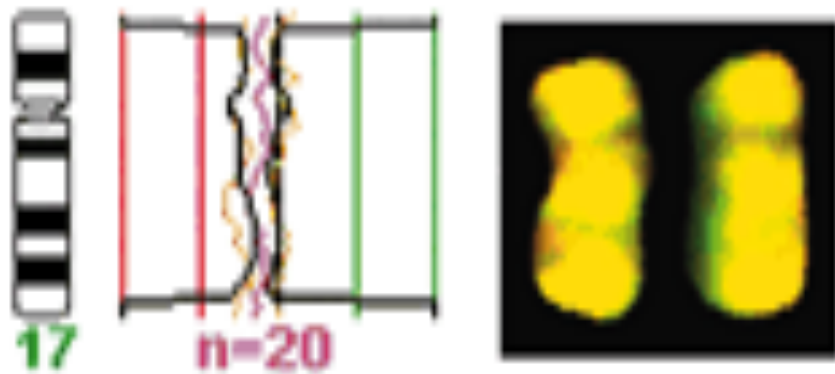


# DNA Methylation changes over time

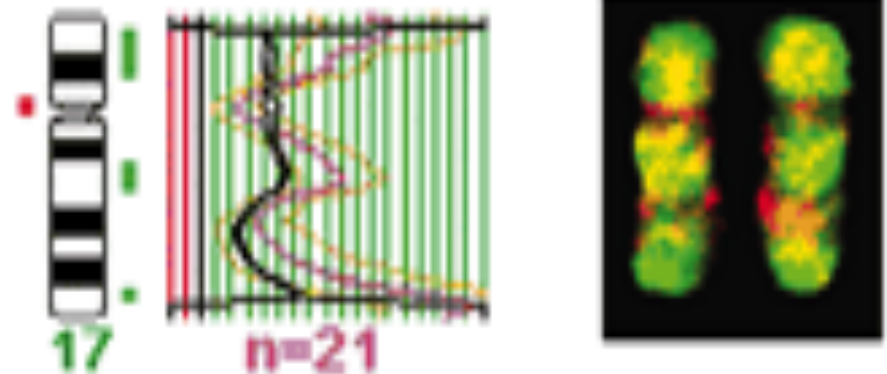


3-year-old twins

# DNA Methylation changes over time



3-year-old twins



50-year-old twins



# Arrays to assess DNA CH<sub>3</sub>

Illumina HumanMethylation450K

Illumina EPIC (>850K CH<sub>3</sub> sites)

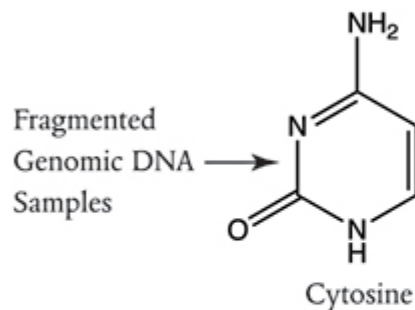
- >90% content on 450K
- New: CH<sub>3</sub> sites in ENCODE open chromatin & enhancers

# (Gold Standard) Bisulfite Conversion of gDNA

## Step 1

### Denaturation

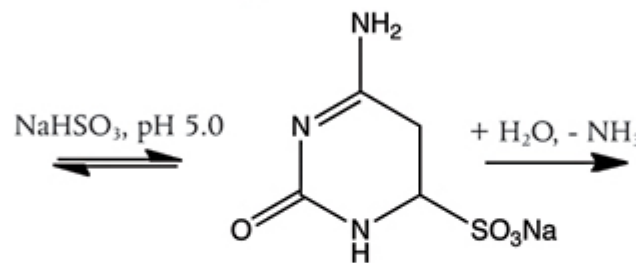
Incubation at 95°C  
fragments genomic DNA



## Step 2

### Conversion

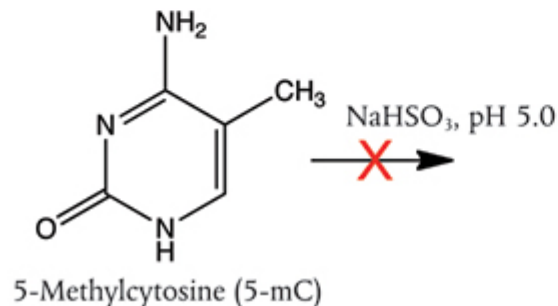
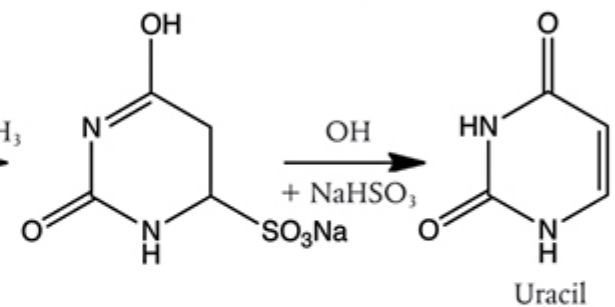
Incubation with sodium bisulfite  
at 65°C and low pH (5-6)  
deaminates cytosine residues  
in fragmented DNA



## Step 3

### Desulphonation

Incubation at high pH  
at room temperature for 15 min  
removes the sulfite moiety,  
generating uracil



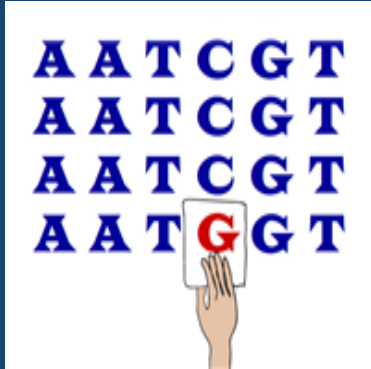
5-mC and 5-hmC (not shown) are not susceptible  
to bisulfite conversion and remain intact



Genetic variation



gene regulation

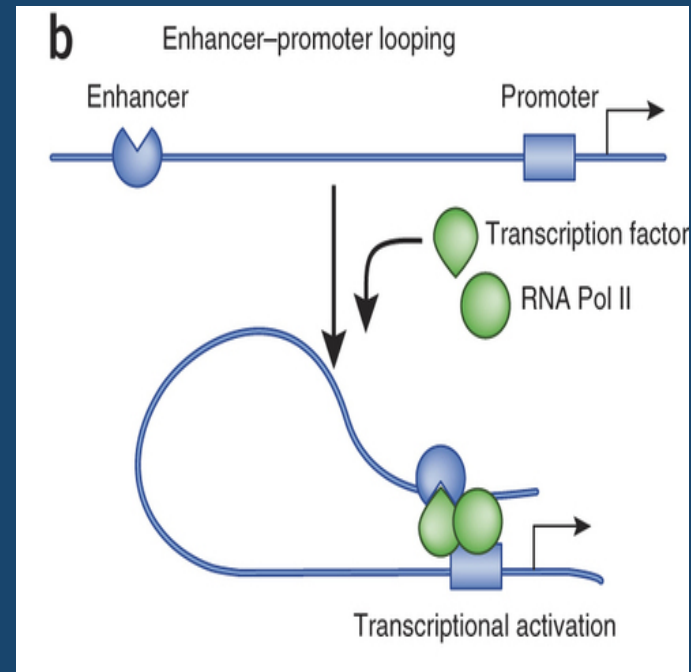
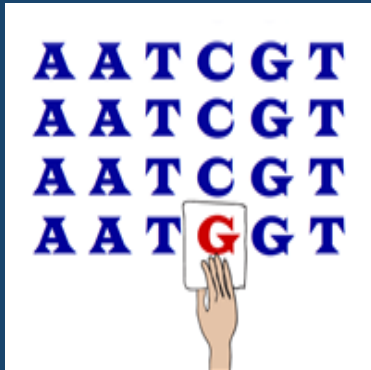


Non-coding  
variation



phenotype

Genetic variation → gene regulation

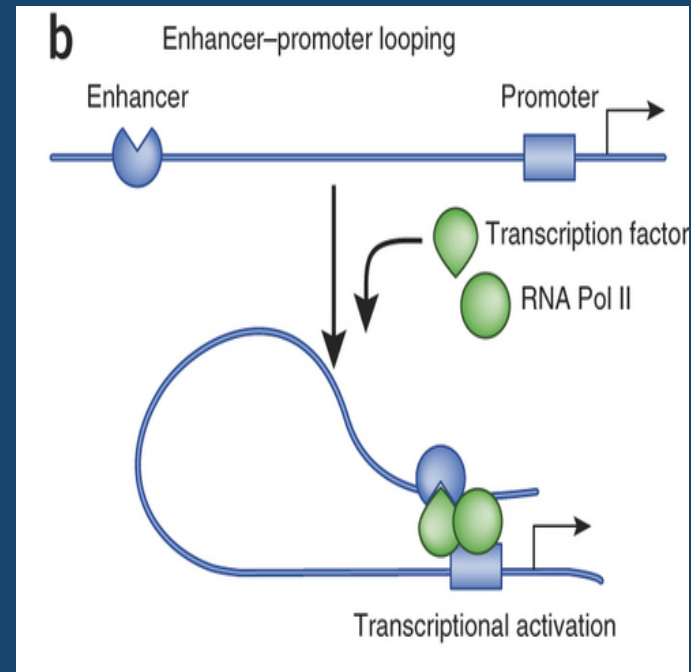
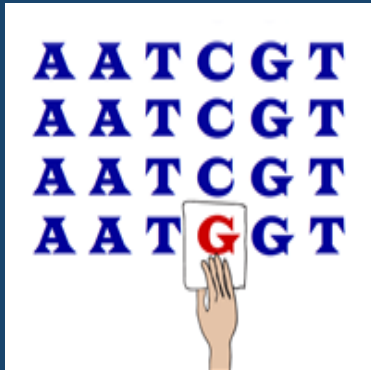


Non-coding  
variation

phenotype



Genetic variation → gene regulation



Non-coding  
variation



phenotype

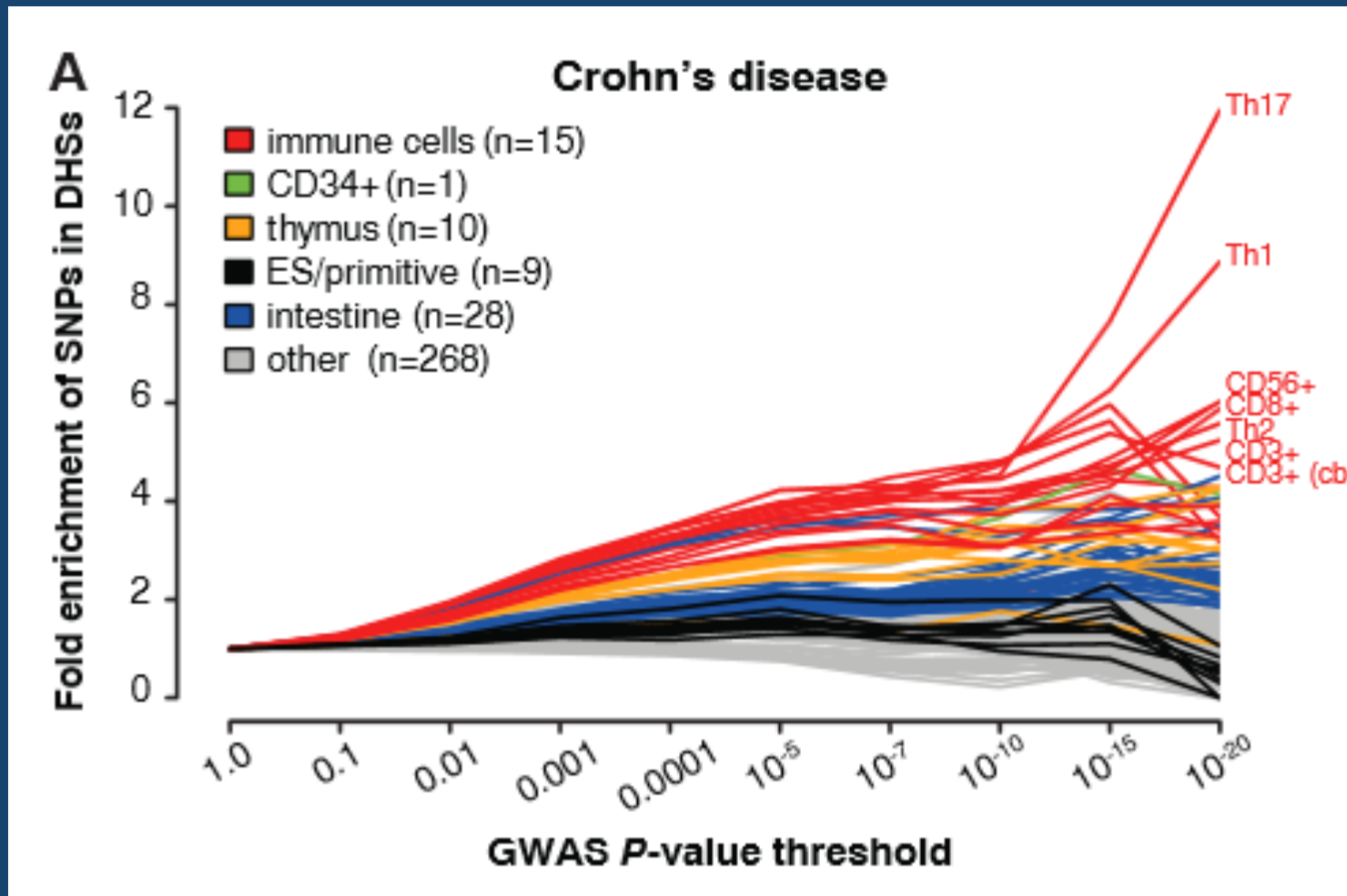
# Key Points

- Most variants identified through GWAS are noncoding
- Non-coding DNA help regulate transcription (e.g. histone marks, eQTLs, DNA methylation)
- Functional annotations are dynamic

# Outline

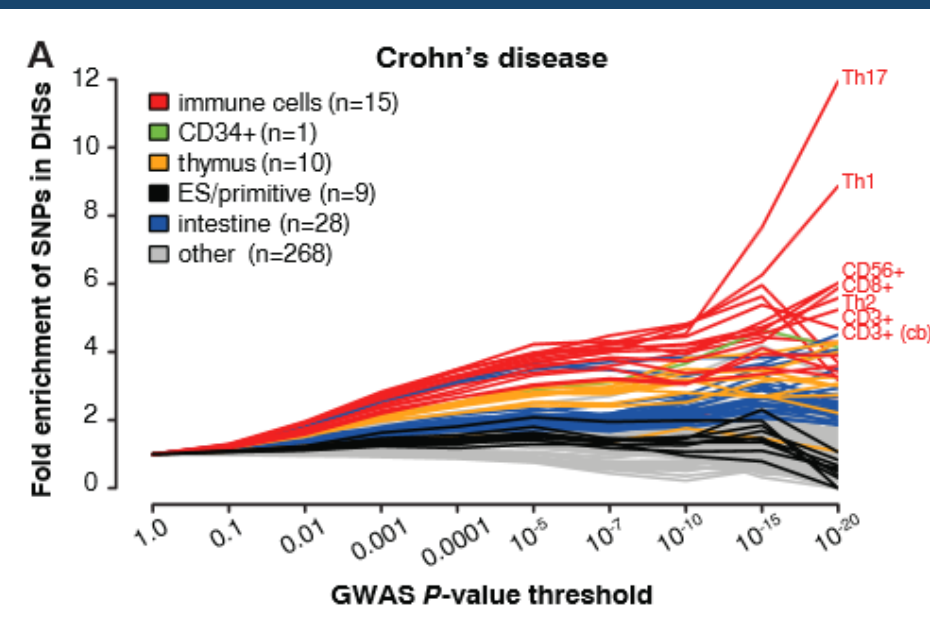
- Functional Genomics 101
- Using functional genomics to prioritize risk variants

# GWAS hits overlap with functional regions

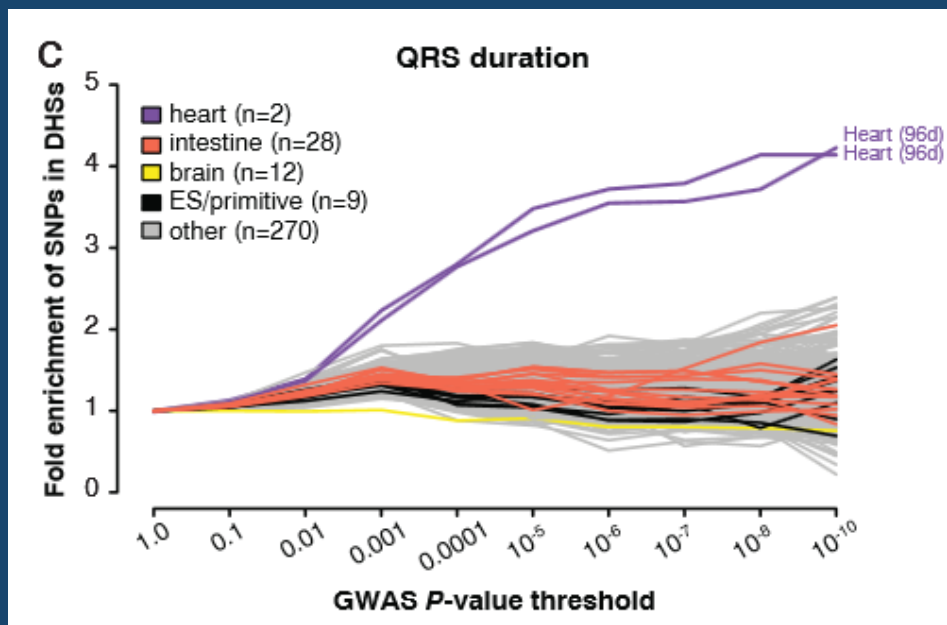


(Franke et al. 2010)  
N SNPs= 938,703 N  
inds= 6,333 cases &  
15,056 controls

# GWAS hits overlap with functional regions

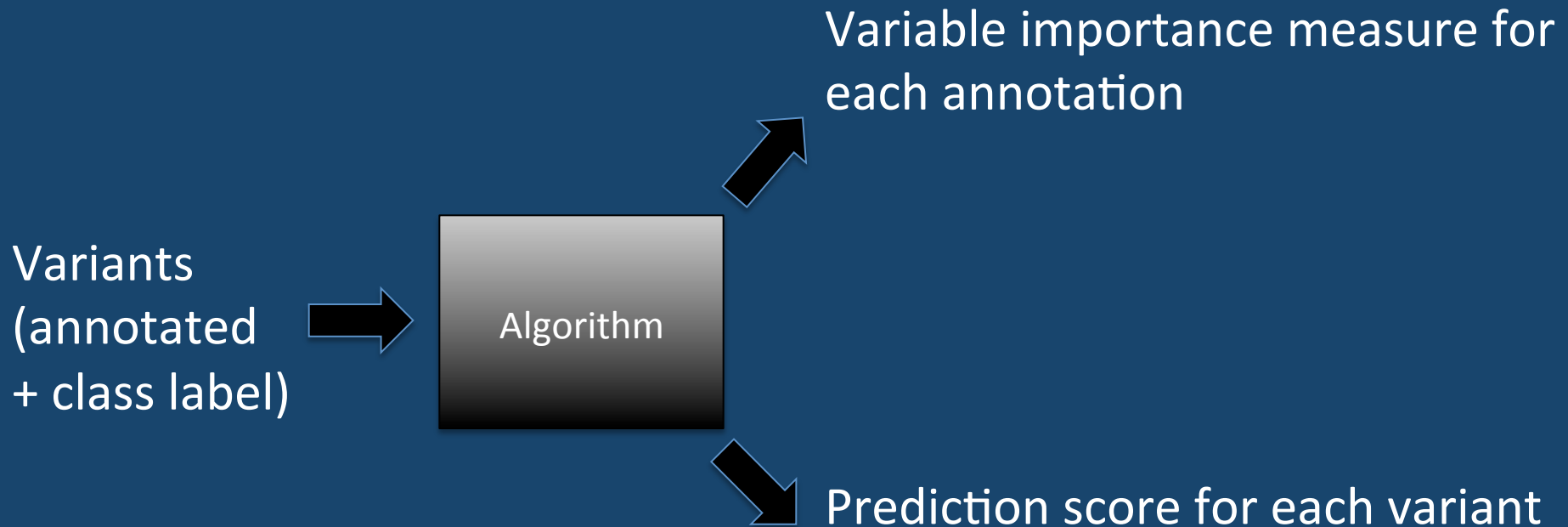


(Franke et al. 2010)  
N SNPs= 938,703  
N inds= 6,333 cases & 15,056 controls



(Sotoodehnia et al. 2010)  
N SNPs ~2.5 M  
N inds~ 40K

# Use functional info to identify hits



# Machine Learning Steps



# Supervised vs. Unsupervised



# Supervised vs. Unsupervised

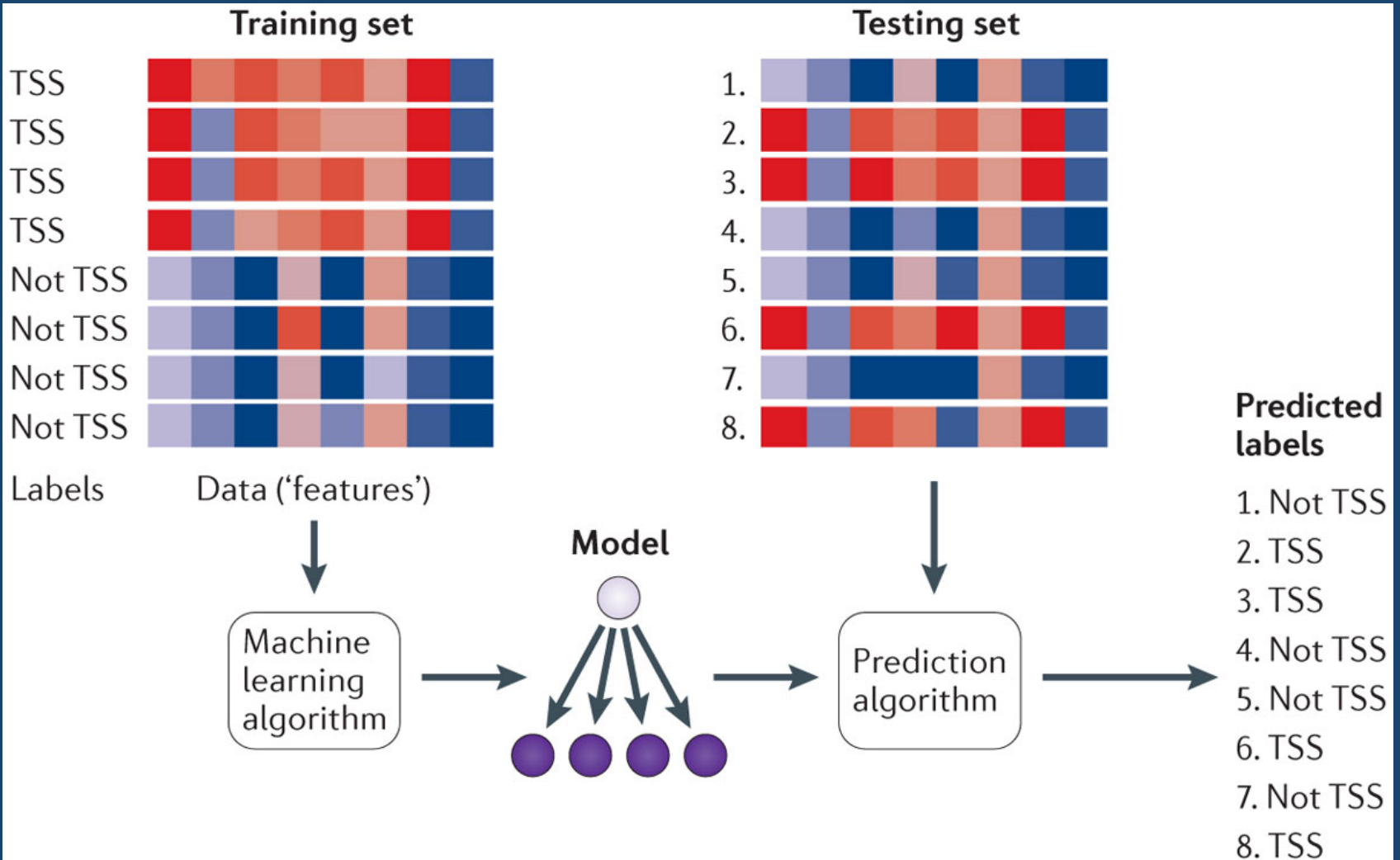


Class assignment  
(data are labelled)



Pattern discovery  
(data aren't labelled)

# A simple ML application



# Quiz! Supervised or unsupervised?

# Quiz! Supervised or unsupervised?

Find novel TFBS based on annotated TFBS in the JASPAR database vs. non-TFBS sequences

|

# Quiz! Supervised or unsupervised?

Find novel TFBS based on annotated TFBS in the JASPAR database vs. non-TFBS sequences

Identify clusters of similar tumors in cancer patients

# Quiz! Supervised or unsupervised?

Find novel TFBS based on annotated TFBS in the JASPAR database vs. non-TFBS sequences

Identify clusters of similar tumors in cancer patients

Identify risk variants from other variants trained on genomic characteristics for *Human Gene Mutation Database* SNPs vs. 1KG (GWAVA, Ritchie et al. 2014)

# Quiz! Supervised or unsupervised?

Find novel TFBS based on annotated TFBS in the JASPAR database vs. non-TFBS sequences

Identify clusters of similar tumors in cancer patients

Identify risk variants from other variants trained on genomic characteristics for *Human Gene Mutation Database* SNPs vs. 1KG (GWAVA, Ritchie et al. 2014)

Segment the genome into X chromatin states using chromatin mark patterns (Segway, Hoffman et al. 2012)

# A Bayesian Method to Incorporate Hundreds of Functional Characteristics with Association Evidence to Improve Variant Prioritization

Sarah A. Gagliano<sup>1,2</sup>, Michael R. Barnes<sup>3</sup>, Michael E. Weale<sup>4,5</sup>, Jo Knight<sup>1,2,5,\*</sup>

## TECHNICAL REPORTS

nature  
genetics

A general framework for estimating the relative pathogenicity of human genetic variants

Martin Kircher<sup>1,5</sup>, Daniela M Witten<sup>2,5</sup>, Preti Jain<sup>3,4</sup>, Brian J O’Roak<sup>1,4</sup>, Gregory M Cooper<sup>3</sup> & Jay Shendure<sup>1</sup>

## BRIEF COMMUNICATIONS

# Functional annotation of noncoding sequence variants

Graham R S Ritchie<sup>1,2</sup>, Ian Dunham<sup>1</sup>, Eleftheria Zeggini<sup>2</sup> & Paul Flicek<sup>1,2</sup>

## TECHNICAL REPORTS

nature  
genetics

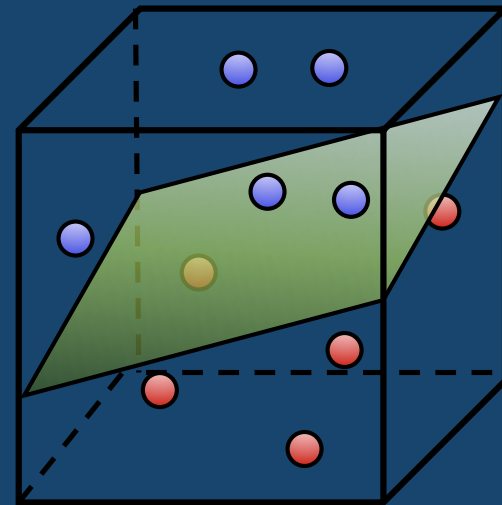
A method to predict the impact of regulatory variants from DNA sequence

Dongwon Lee<sup>1,4</sup>, David U Gorkin<sup>1,3,4</sup>, Maggie Baker<sup>1</sup>, Benjamin J Strober<sup>2</sup>, Alessandro L Asoni<sup>2</sup>, Andrew S McCallion<sup>1</sup> & Michael A Beer<sup>1,2</sup>

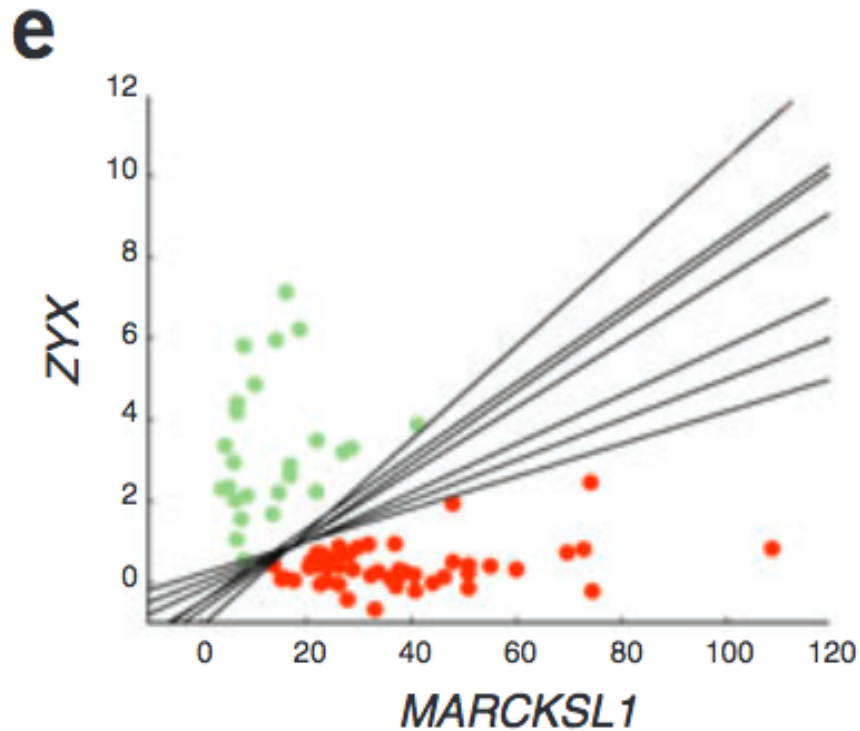


# Support Vector Machines

Maximum-margin separating  
hyperplane in multi-dimensional  
space

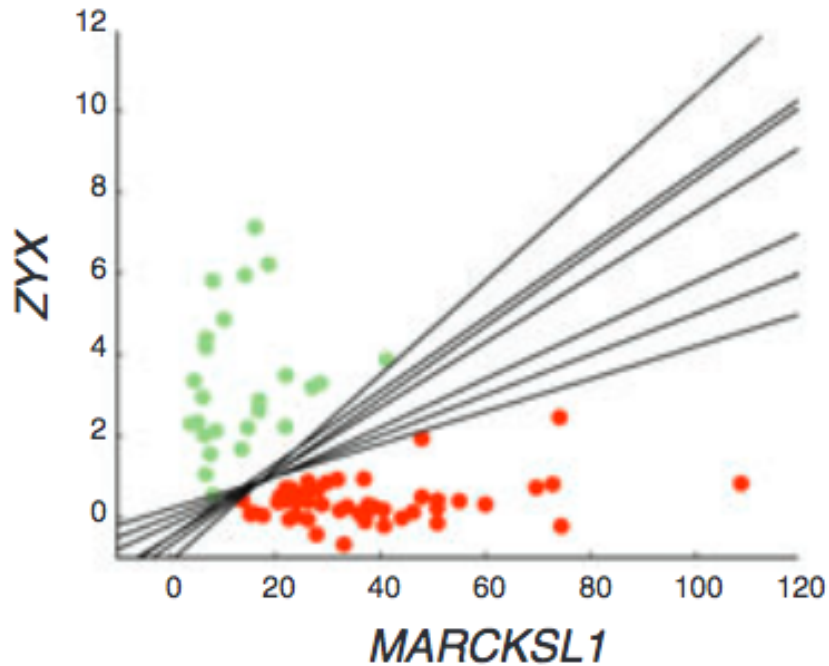


# Which line provides the best classifier?

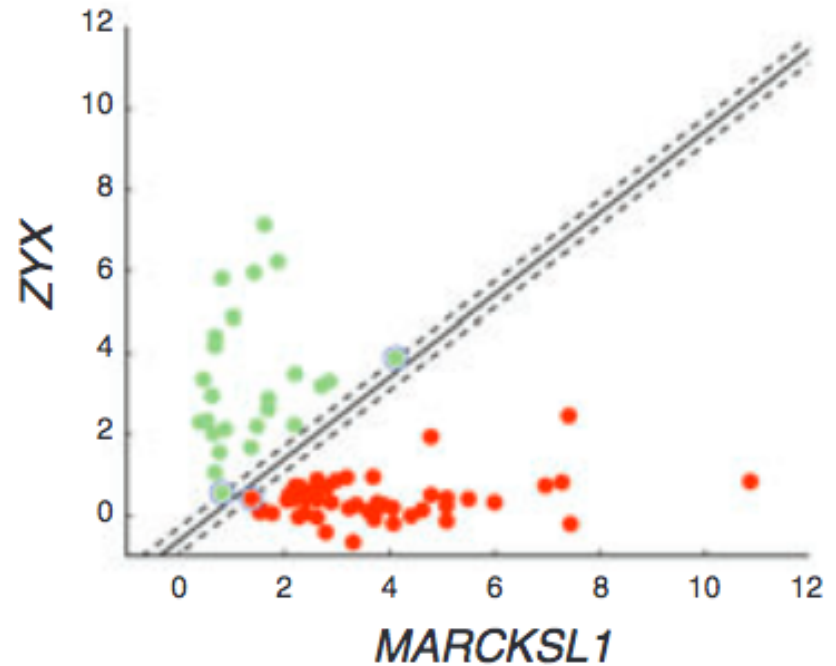


# Which line provides the best classifier?

**e**



**f**



# The “Soft Margin”

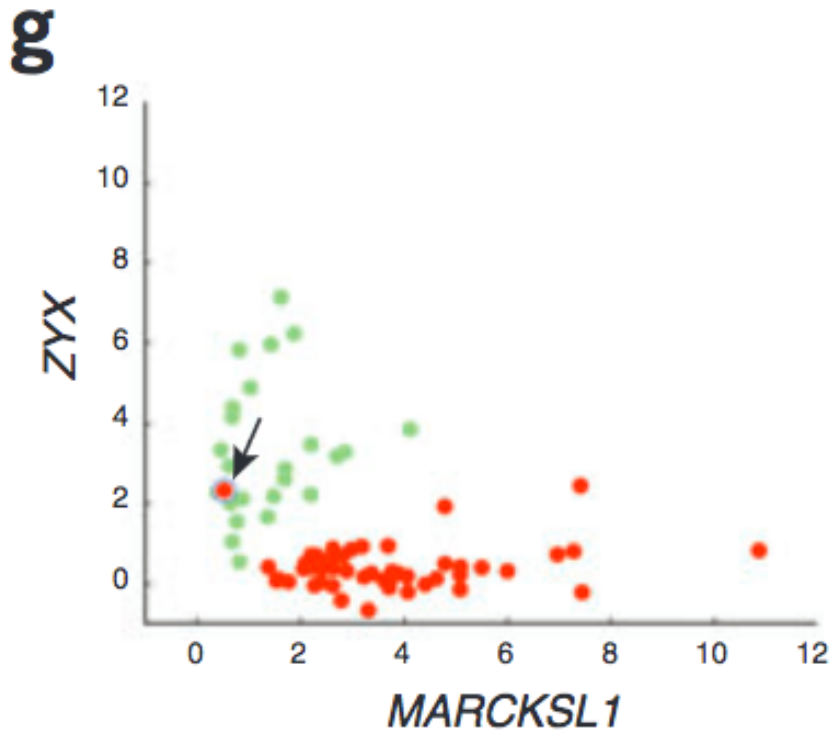


Figure 1 Noble 2016 *Nat.Biotech*

# The “Soft Margin”

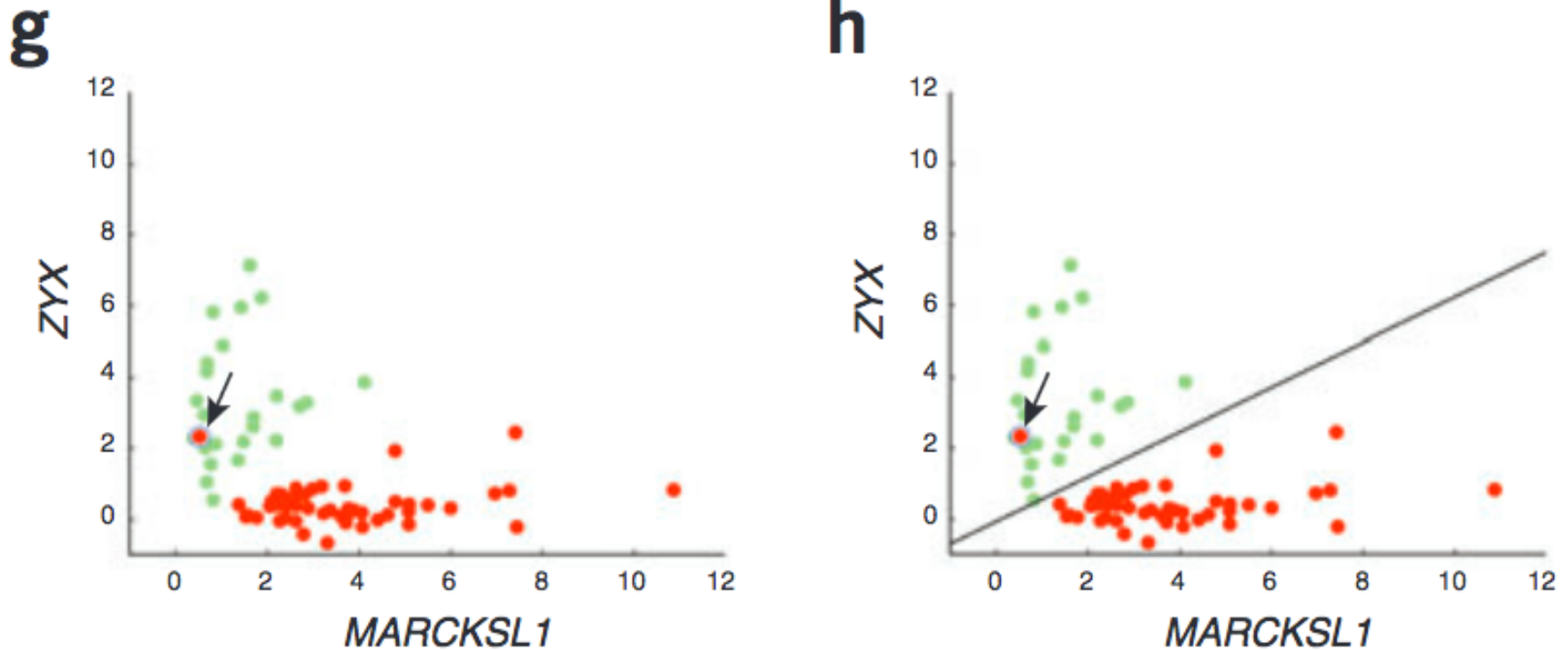
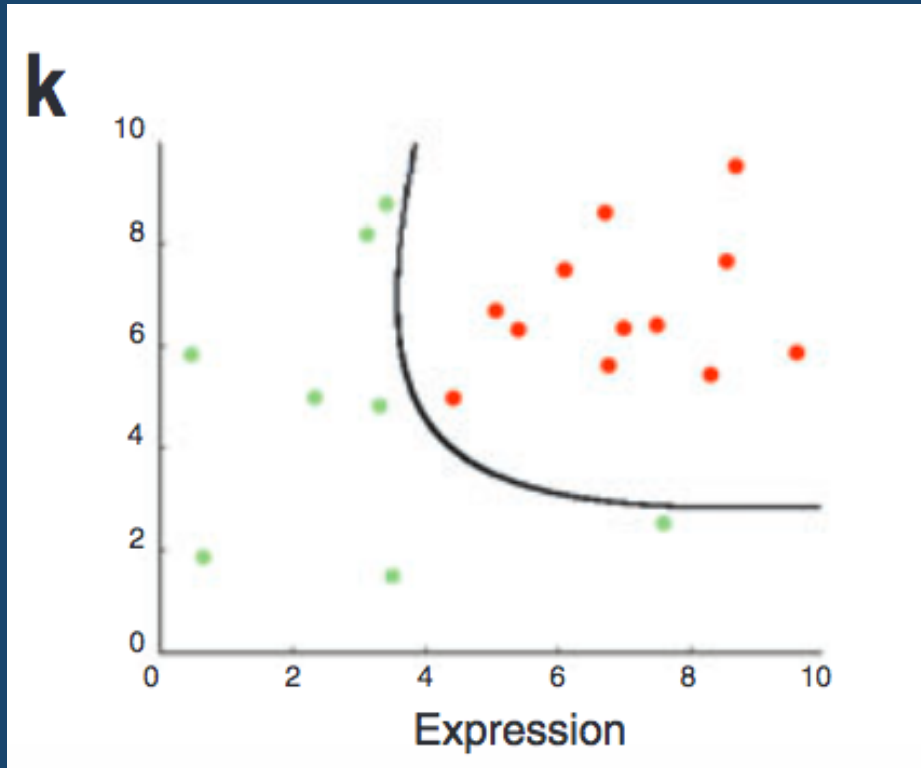


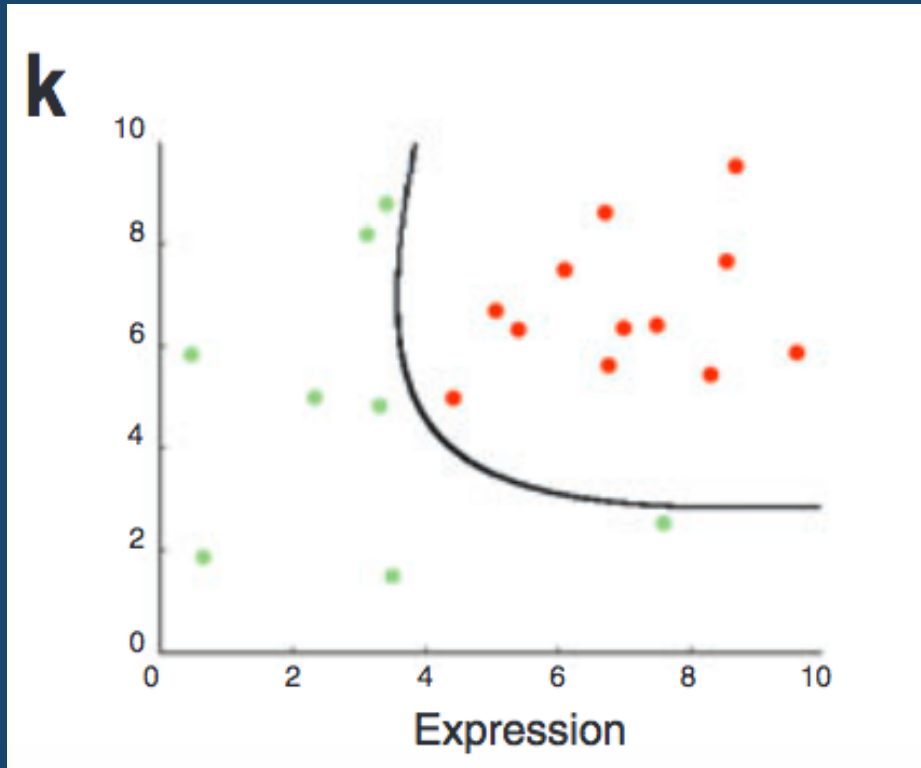
Figure 1 Noble 2016 *Nat.Biotech*

# Linearly nonseparable? Kernel Trick



1. Kernel function to project the data in 2-D to 4-D space

# Linearly nonseparable? Kernel Trick



1. Kernel function to project the data in 2-D to 4-D space
2. Project SVM hyperplane in 4-D to original 2-D space

1. Classifier
2. SVM protocol
3. Some results

## TECHNICAL REPORTS

nature  
genetics

A general framework for estimating the relative pathogenicity of human genetic variants

Martin Kircher<sup>1,5</sup>, Daniela M Witten<sup>2,5</sup>, Preti Jain<sup>3,4</sup>, Brian J O’Roak<sup>1,4</sup>, Gregory M Cooper<sup>3</sup> & Jay Shendure<sup>1</sup>



# 1. Classifier

# Simulated vs. Observed variants

Simulated (14.7 million variants):

- Empirical model of sequence evolution with local adjustment of mutation rates

# Simulated vs. Observed variants

Simulated (14.7 million variants):

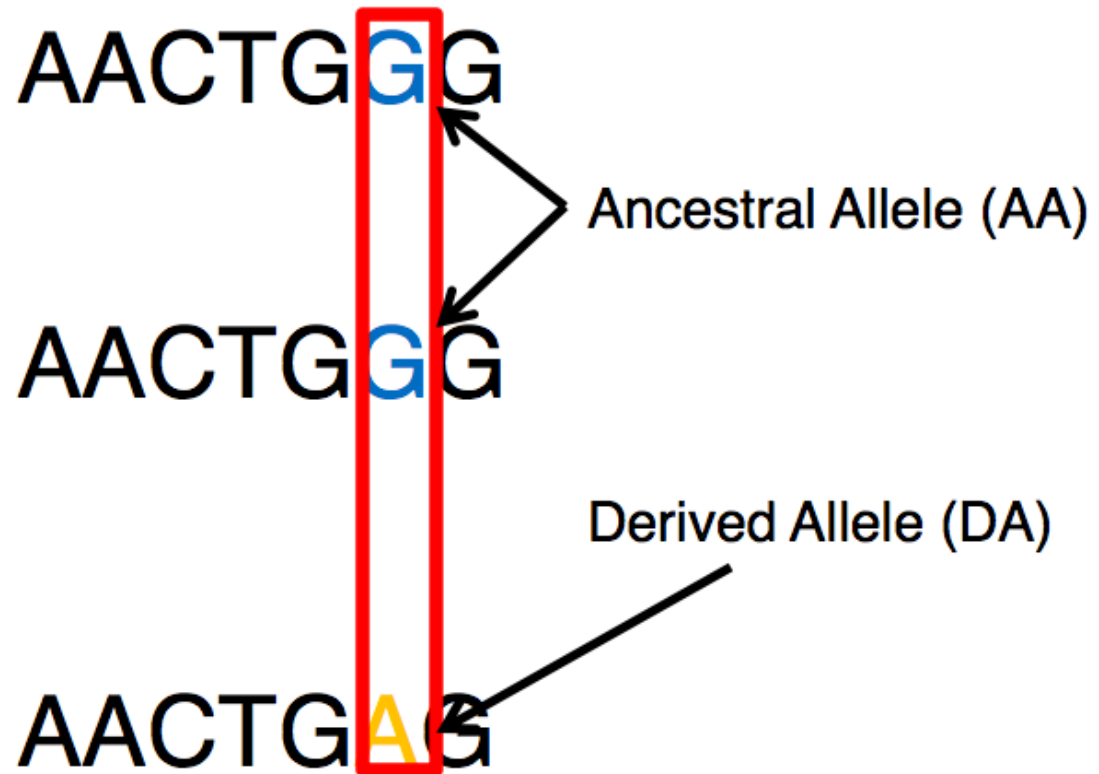
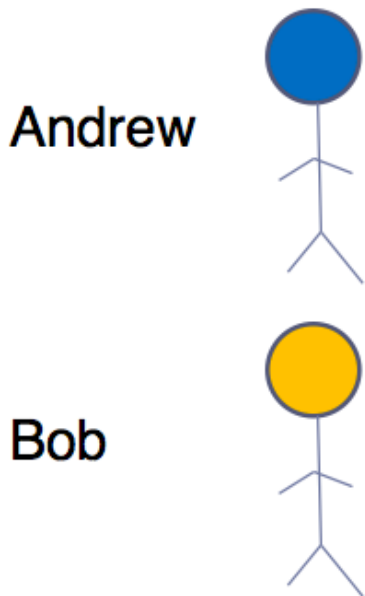
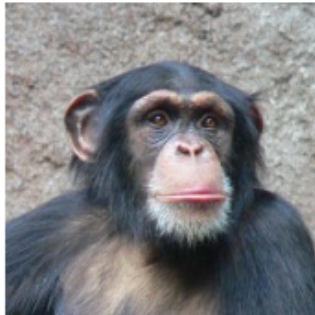
- Empirical model of sequence evolution with local adjustment of mutation rates

Observed (14.7 million variants):

- Human derived allele >95% (exposed to many generations of natural selection)

1% reserved for testing

# Ancestral Allele (AA) and Derived Allele (DA)



# Peppered Moth Evolution

Ancestral



Derived

No pollution

# Peppered Moth Evolution

Ancestral



Derived



No pollution

Pollution

Photo from: <http://www.truthinscience.org.uk/tis2/index.php/evidence-for-evolution-mainmenu-65/127-the-peppered-moth.html>

## 2. SVM Protocol

# Kircher's use of SVM

- Linear kernel
- Prior feature selection (univariate analysis)
- Imputed missing values
- Boolean variables for categorical variables
- Interaction terms (include the few that improve the two-feature linear regression models)



## 3. Some results

# Performance in GWAS

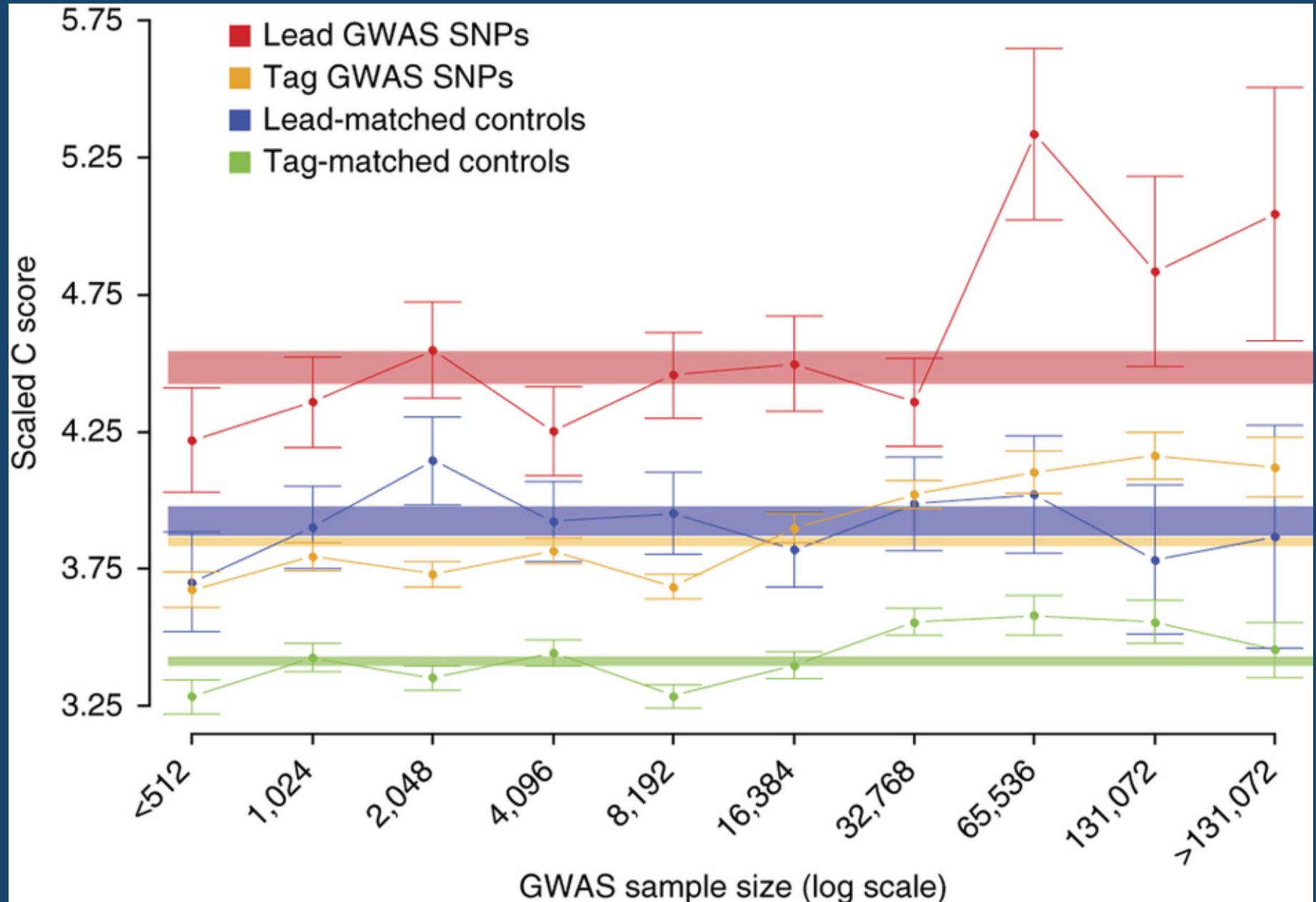


Figure 5 from Kircher *et al.* 2014 *Nat.Genetics*

# deltaSVM

- Takes into account tissue-specificity

1. Classifier
2. SVM protocol
3. Some results

## TECHNICAL REPORTS

nature  
genetics

A method to predict the impact of regulatory variants  
from DNA sequence

Dongwon Lee<sup>1,4</sup>, David U Gorkin<sup>1,3,4</sup>, Maggie Baker<sup>1</sup>, Benjamin J Strober<sup>2</sup>, Alessandro L Asoni<sup>2</sup>,  
Andrew S McCallion<sup>1</sup> & Michael A Beer<sup>1,2</sup>

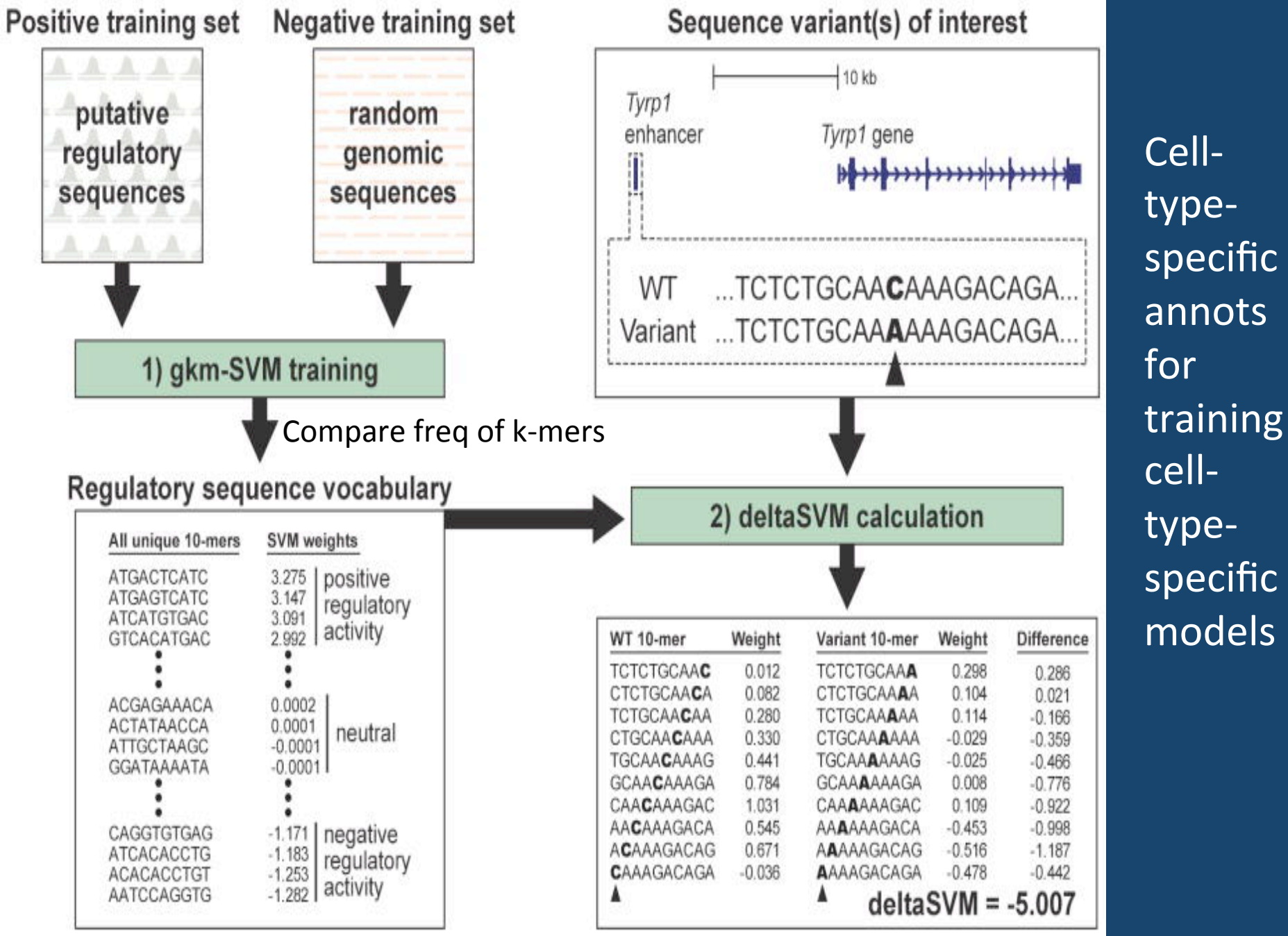


Figure 1 from Lee *et al.* 2015 *Nat.Genetics*

# deltaSVM

- Takes into account tissue-specificity

1. Classifier
2. SVM protocol
3. Some results

## TECHNICAL REPORTS

nature  
genetics

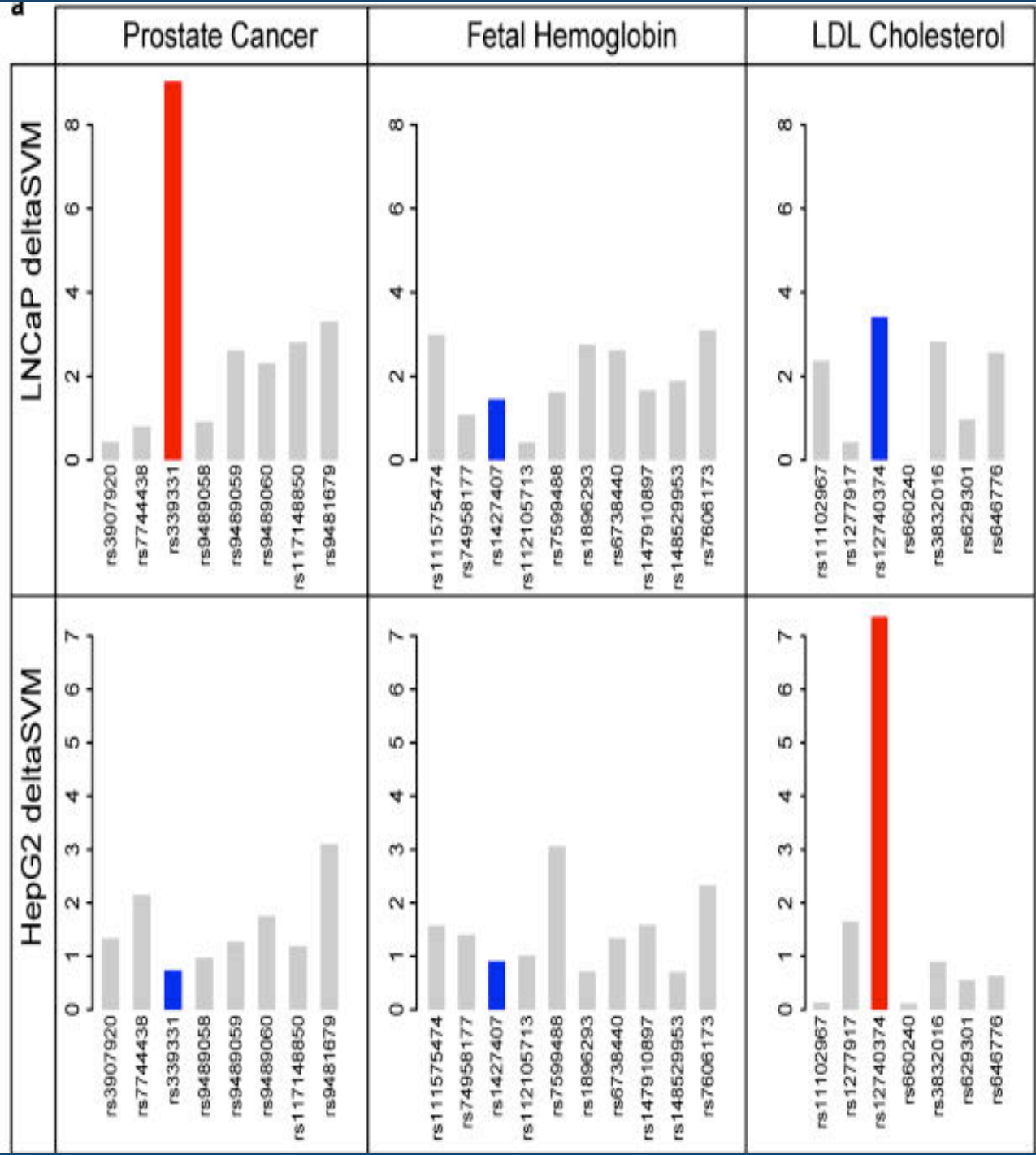
A method to predict the impact of regulatory variants  
from DNA sequence

Dongwon Lee<sup>1,4</sup>, David U Gorkin<sup>1,3,4</sup>, Maggie Baker<sup>1</sup>, Benjamin J Strober<sup>2</sup>, Alessandro L Asoni<sup>2</sup>,  
Andrew S McCallion<sup>1</sup> & Michael A Beer<sup>1,2</sup>

delta SVM  
values for 3  
experimentally  
validated SNPs

human prostate  
adenocarcinoma

Hepatocytes  
(main cell in liver)



# Key Points

- GWAS variants are enriched for functional information in a tissue-specific manner
- Machine learning can be used to find patterns in functional data to identify novel variants associated with disease

