

# Identify Disease-Causal Genes from GWAS Loci by 3D Genome Structure, Regulatory Landscapes & Deep Learning

Yi-Hsiang Hsu, MD, ScD

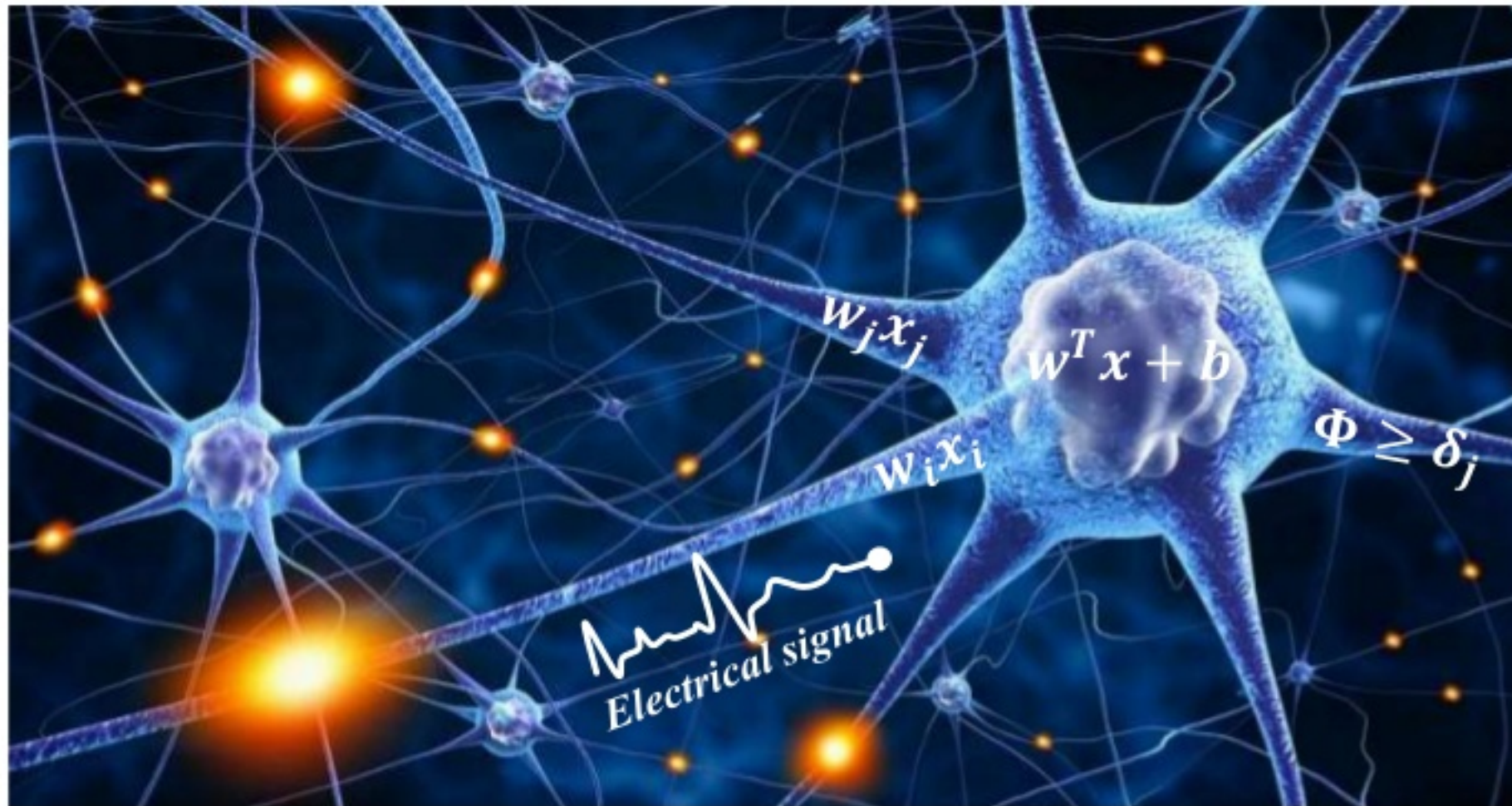


YongSheng Huang, Ph.D



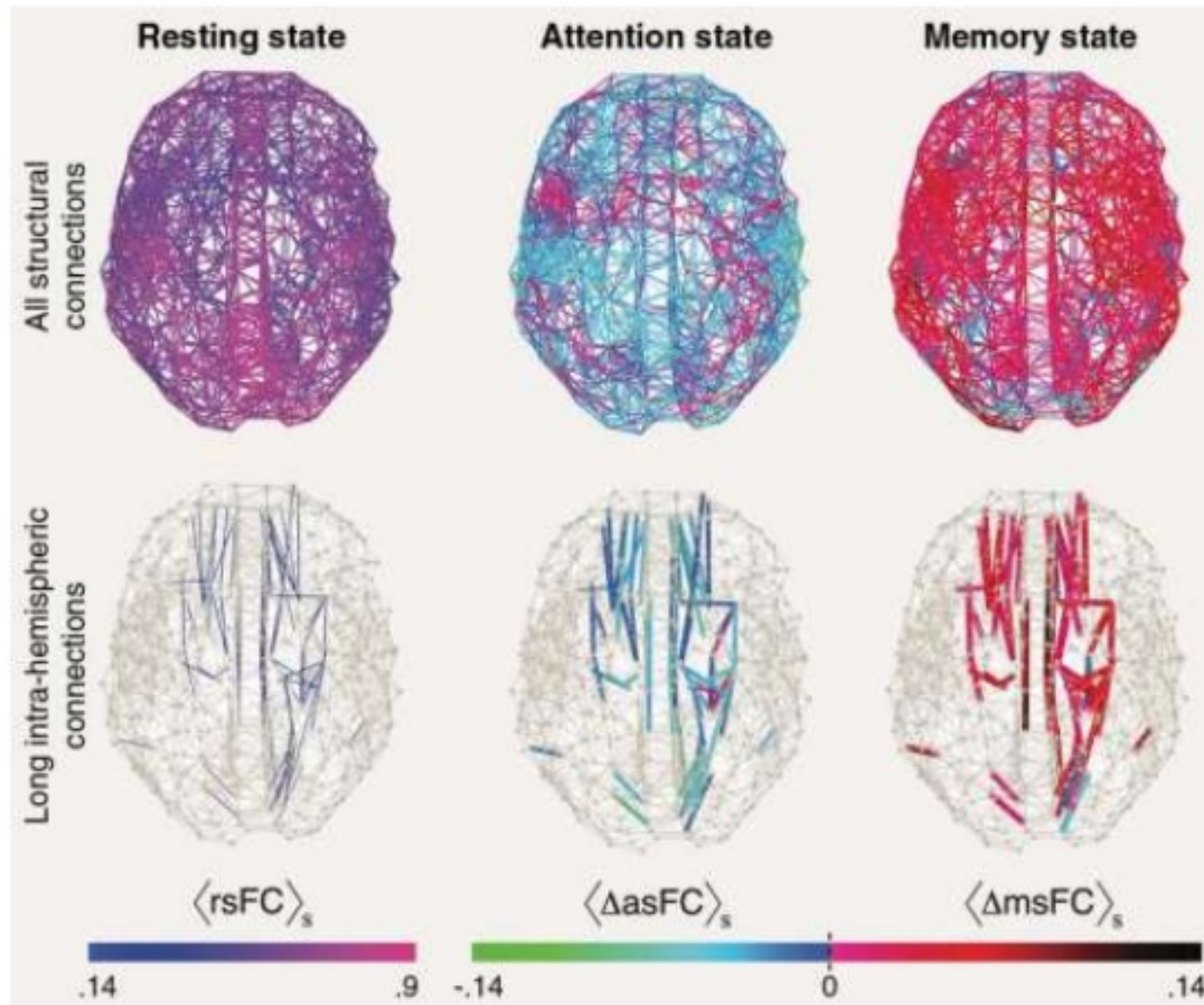
# Deep Learning: The Inspiration

*“the deepest concepts in mathematics are those which link one world of ideas with another”*  
---- Freeman Dyson





# Deep Learning: The Natural Form



# Deep-Learning: The Renaissance

The New York Times

## *Scientists See Promise in Deep-Learning Programs*

by JOHN MARKOFF

Nov. 23, 2012

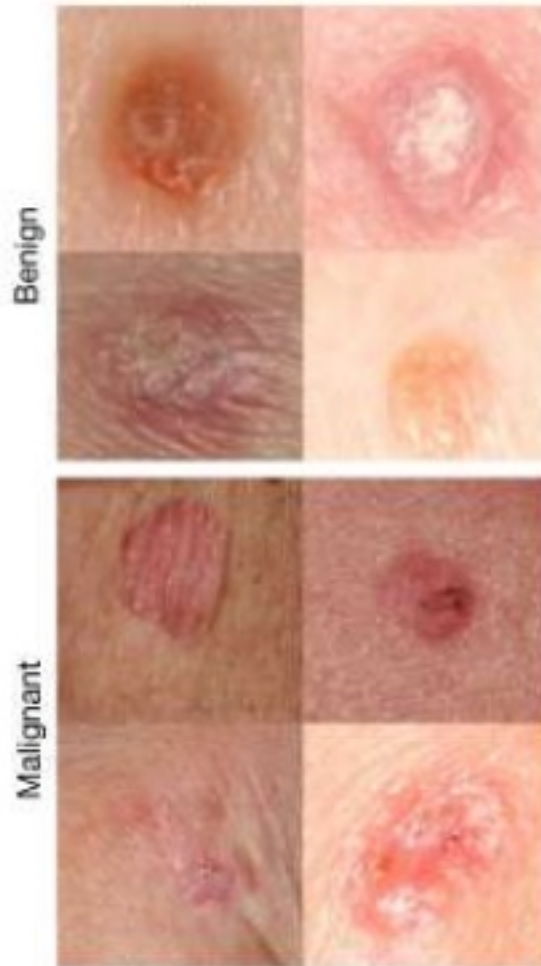
In the 1960s, ..... believed that a workable artificial intelligence system was just 10 years away. In the 1980s, a wave of commercial start-ups collapsed, leading to what some people called the “**A.I. winter.**”

But recent achievements have impressed ..... In October, for example, a team of graduate students studying with the University of Toronto computer scientist [Geoffrey E. Hinton](#) won the top prize in a contest sponsored by Merck to design software to **help find molecules that might lead to new drugs.**

# Deep Learning: Impact on Medicine

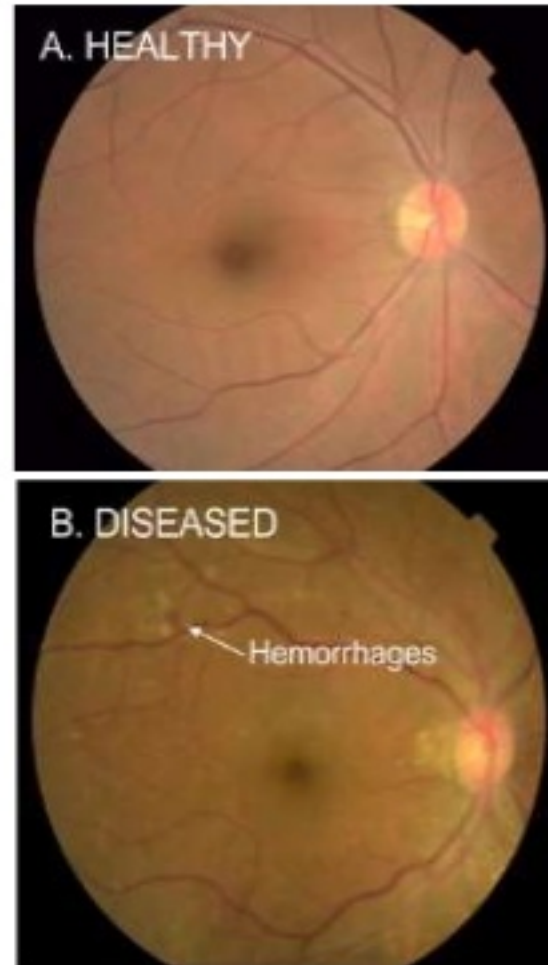
*On par performance as 21  
board-certified pathologists*

Epidermal lesions

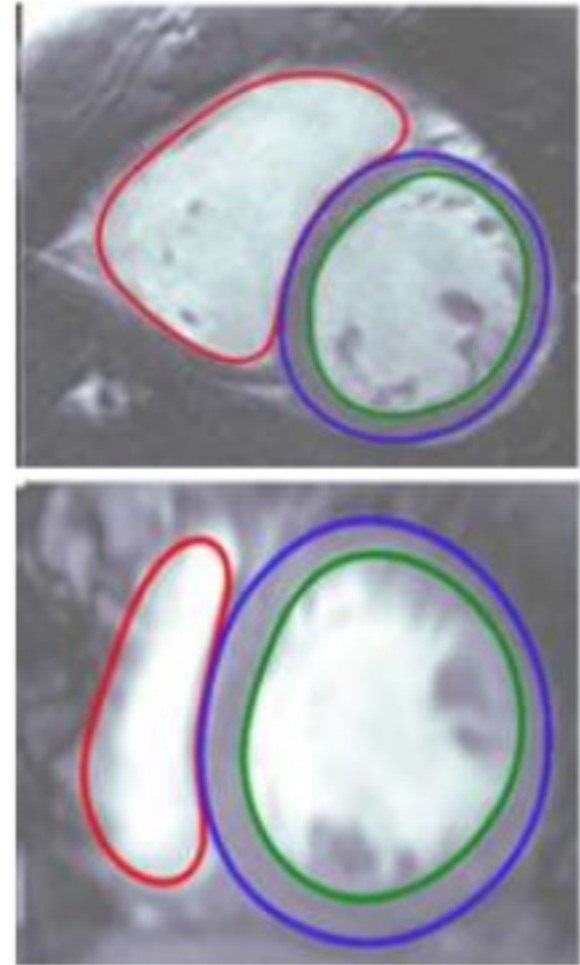


Nature 2017 Feb

*>90% specificity and  
sensitivity as board-certified  
ophthalmologists*



*Artery's Cardio DL wins  
FDA approval for clinical  
diagnosis (10-sec vs. 1hr)*





# Deep Learning: The New Disruption

*Can we leverage DL to identify genetic variants that are disease causal, so that we can treat diseases at its root level per individual patient ?*

***Yi-Hsiang Hsu, MD, ScD***

[yihsianghsu@hsl.harvard.edu](mailto:yihsianghsu@hsl.harvard.edu)  
[yihsiang@broadinstitute.org](mailto:yihsiang@broadinstitute.org)

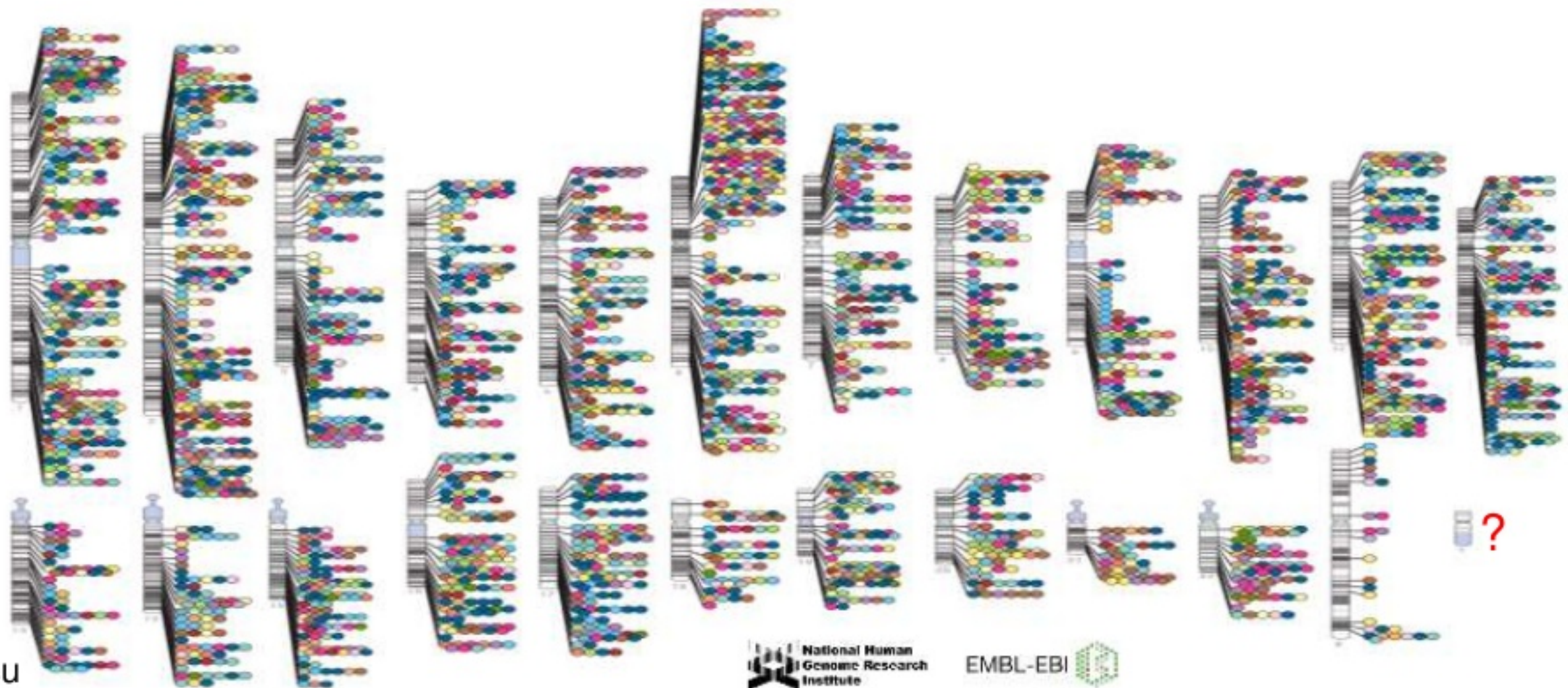
***Director & Associate Professor, HSL GeriOmics Center, Harvard Medical Sch  
Program for Quantitative Genomics, Harvard School of Public Health  
Associate Member, BROAD Institute of MIT and Harvard  
NHLBI Framingham Heart Study Investigator***





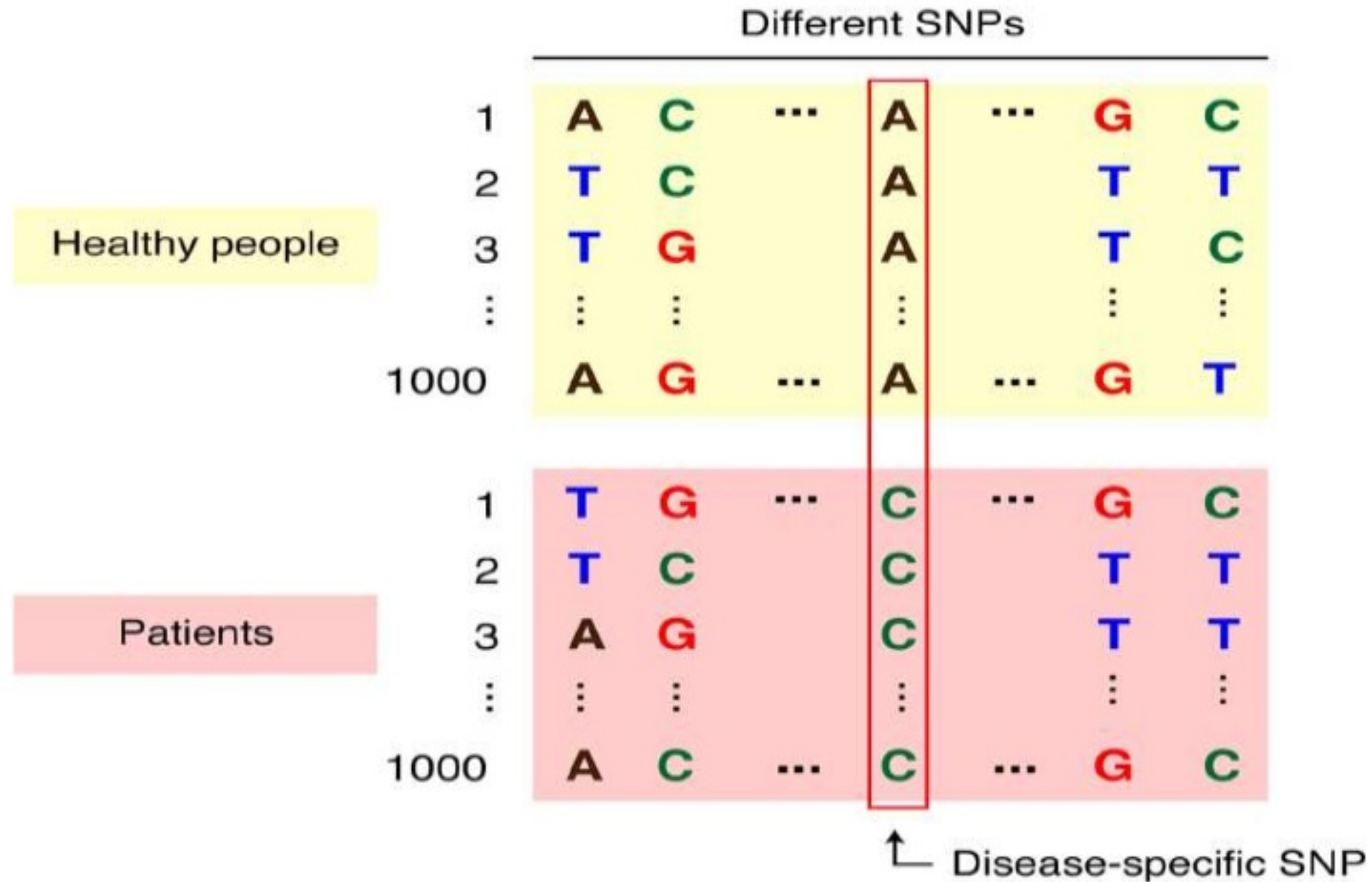
# Genome-Wide Association Studies (GWAS) Catalog

- Identified ~13,000 genetic variants (single nucleotide mutations/polymorphisms) to be associated with ~2,000 diseases/phenotypes

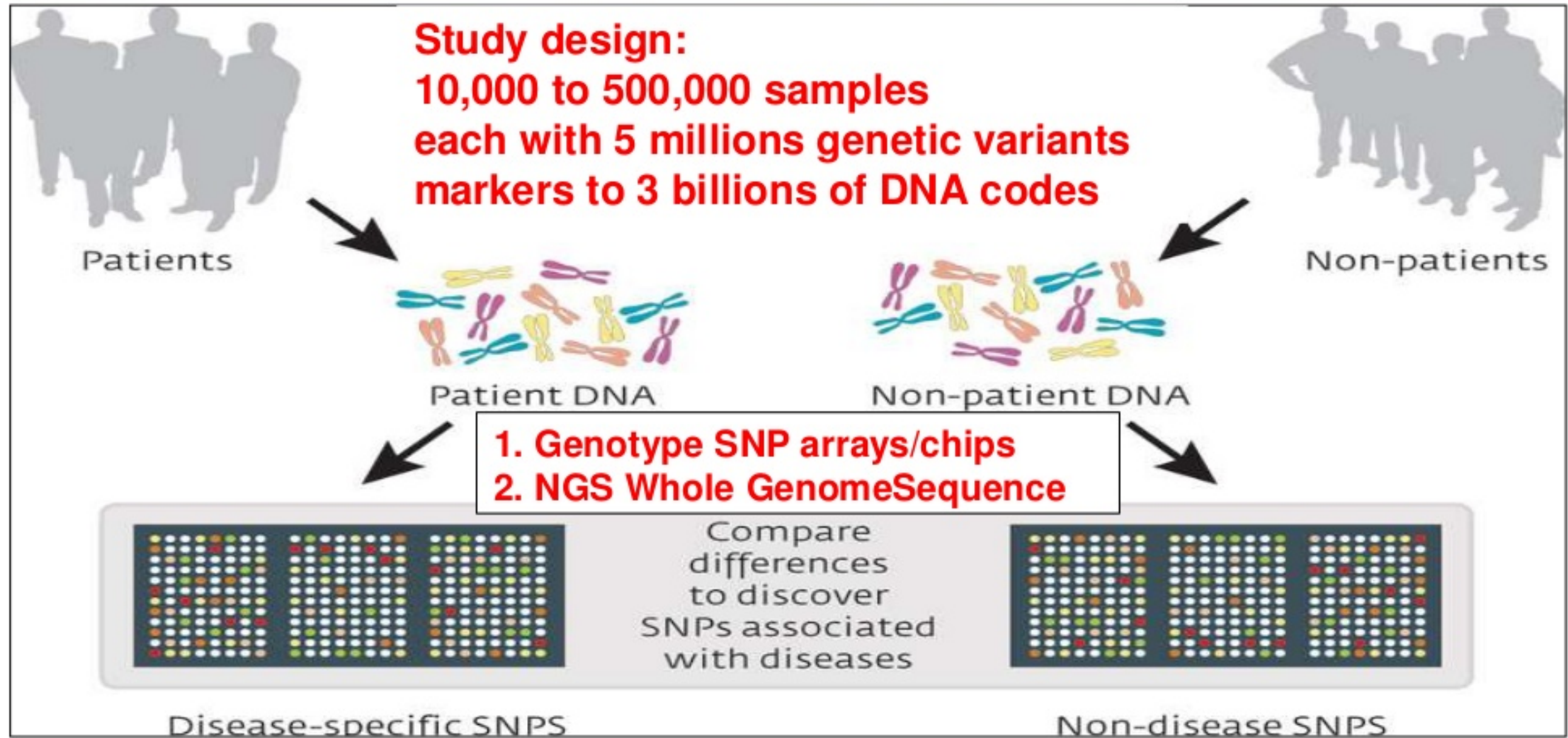




# Genome-Wide Association Scans



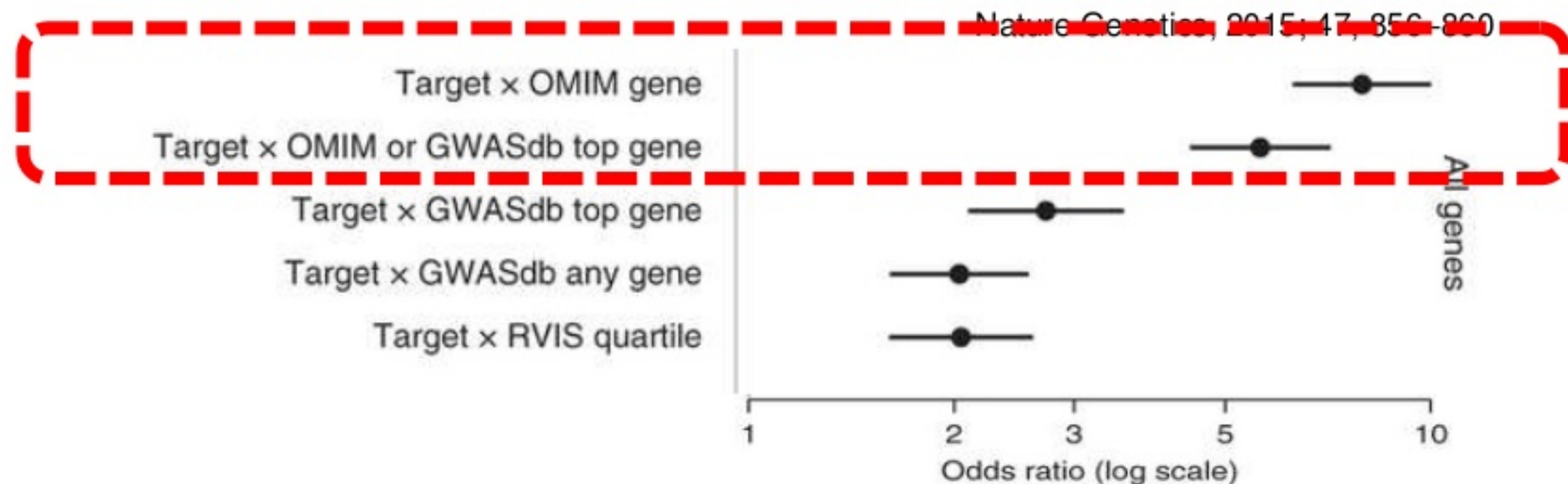
# GWAS (Whole Genome Association) Scans





# % Successfully Approved Drugs & Human Genetics

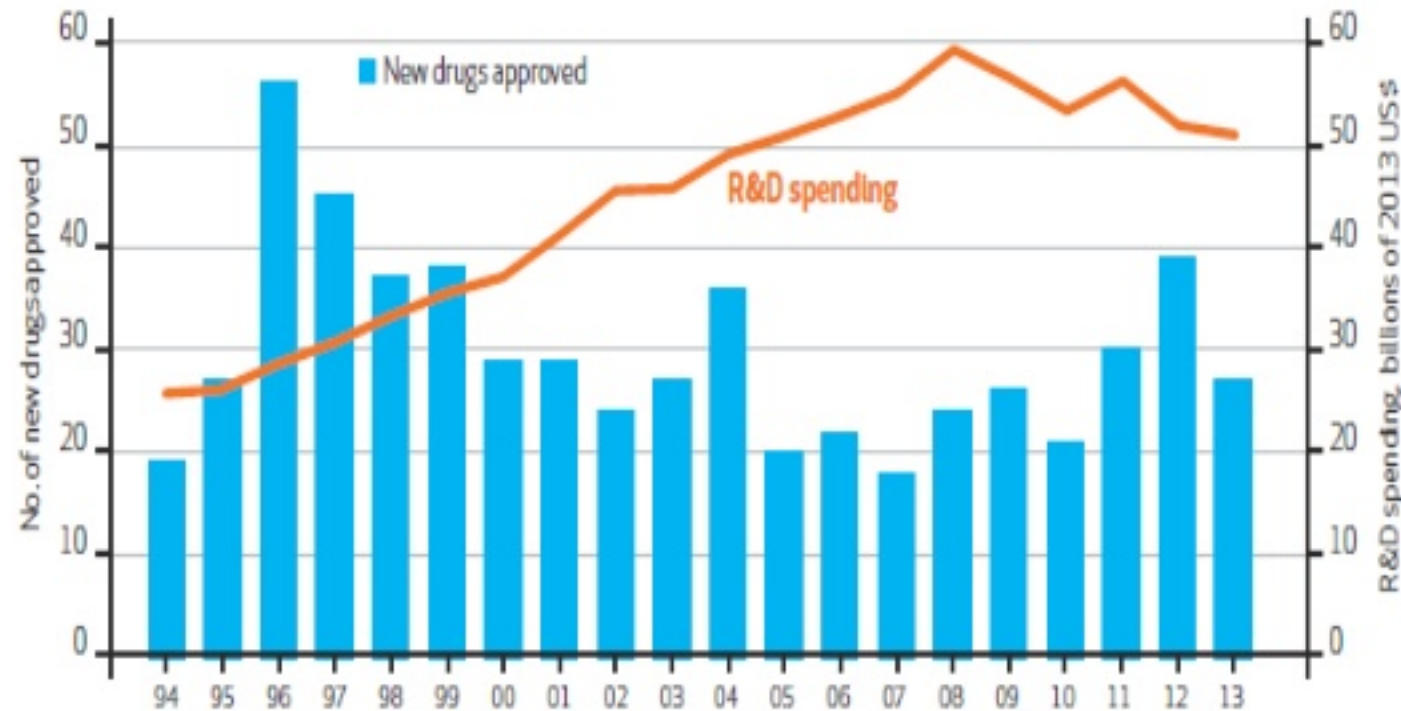
- The impact on medical care from GWAS could potentially be substantial



- FDA approved drugs with human genetic information are 5~10X more likely to be successful
- Failure targets at each drug development stage (pre-clinical, phase I, II, III) are more likely to be those targets without genetic validation

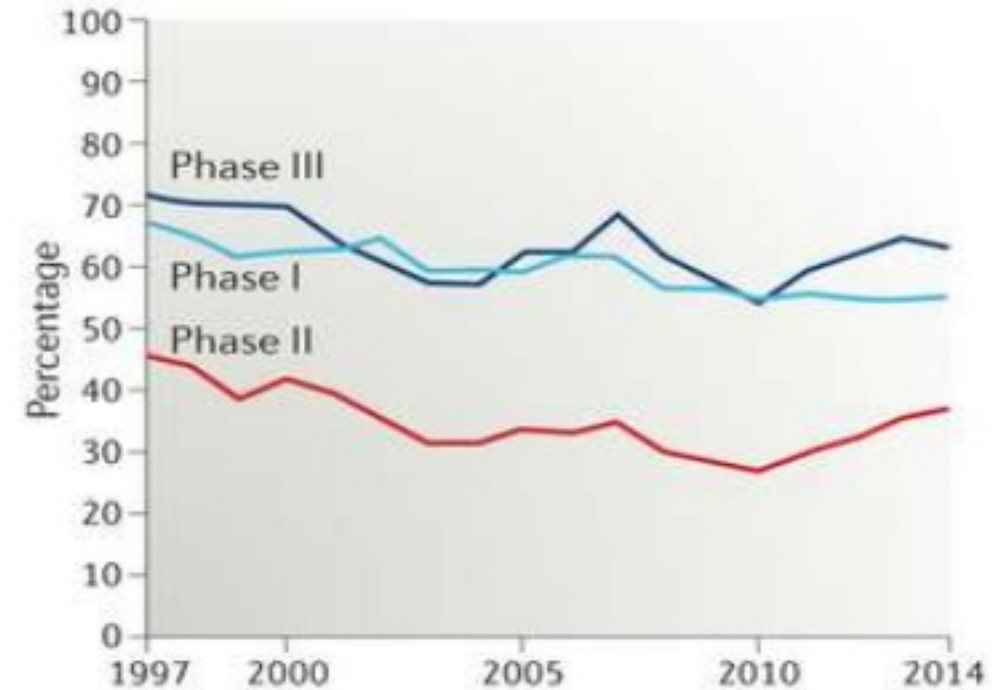
# R&D Spending on New Drugs $\neq$ Drug Approvals

Annual New Drug Approvals By The Food And Drug Administration (FDA) And Industry Spending On Research And Development (R&D), 1994-2013



**a** Success rates by phase

Percentage likelihood of moving to next phase, 3-year rolling average\*

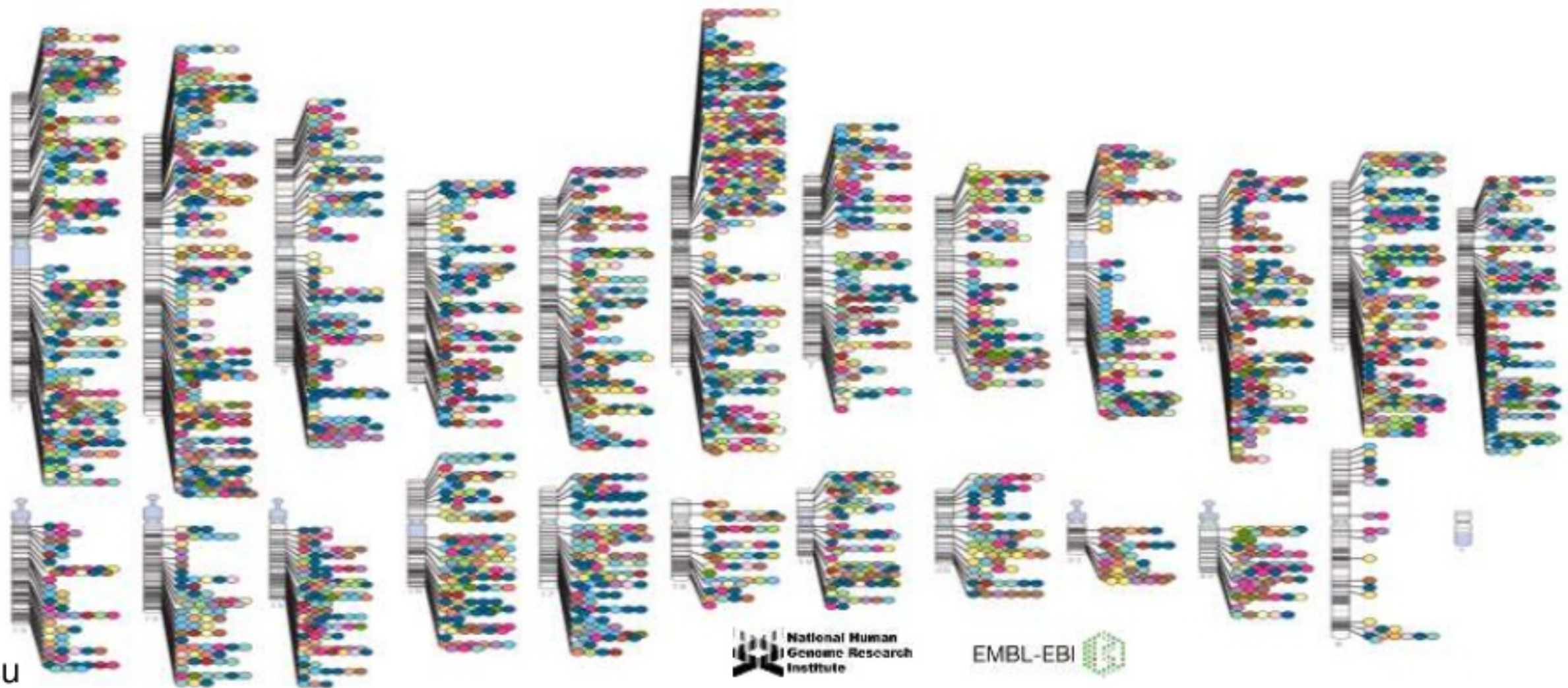


- New a better drug development pipeline
- Utilizing human genetic information/validation is the key



# Genome-Wide Association Studies (GWAS) Catalog

- Identified ~13,000 genetic variants (single nucleotide mutations/polymorphisms) to be associated with ~2,000 diseases/phenotypes

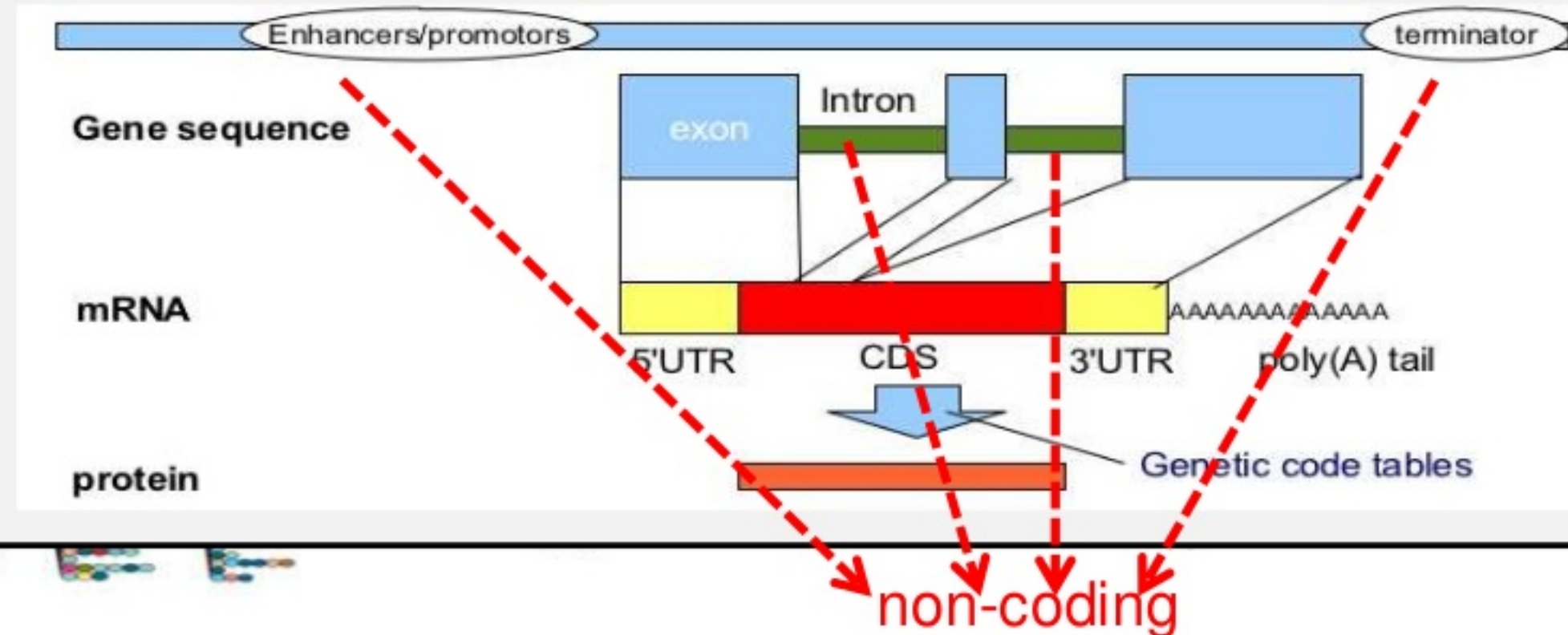




# Genome-Wide Association Studies (GWAS) Catalog

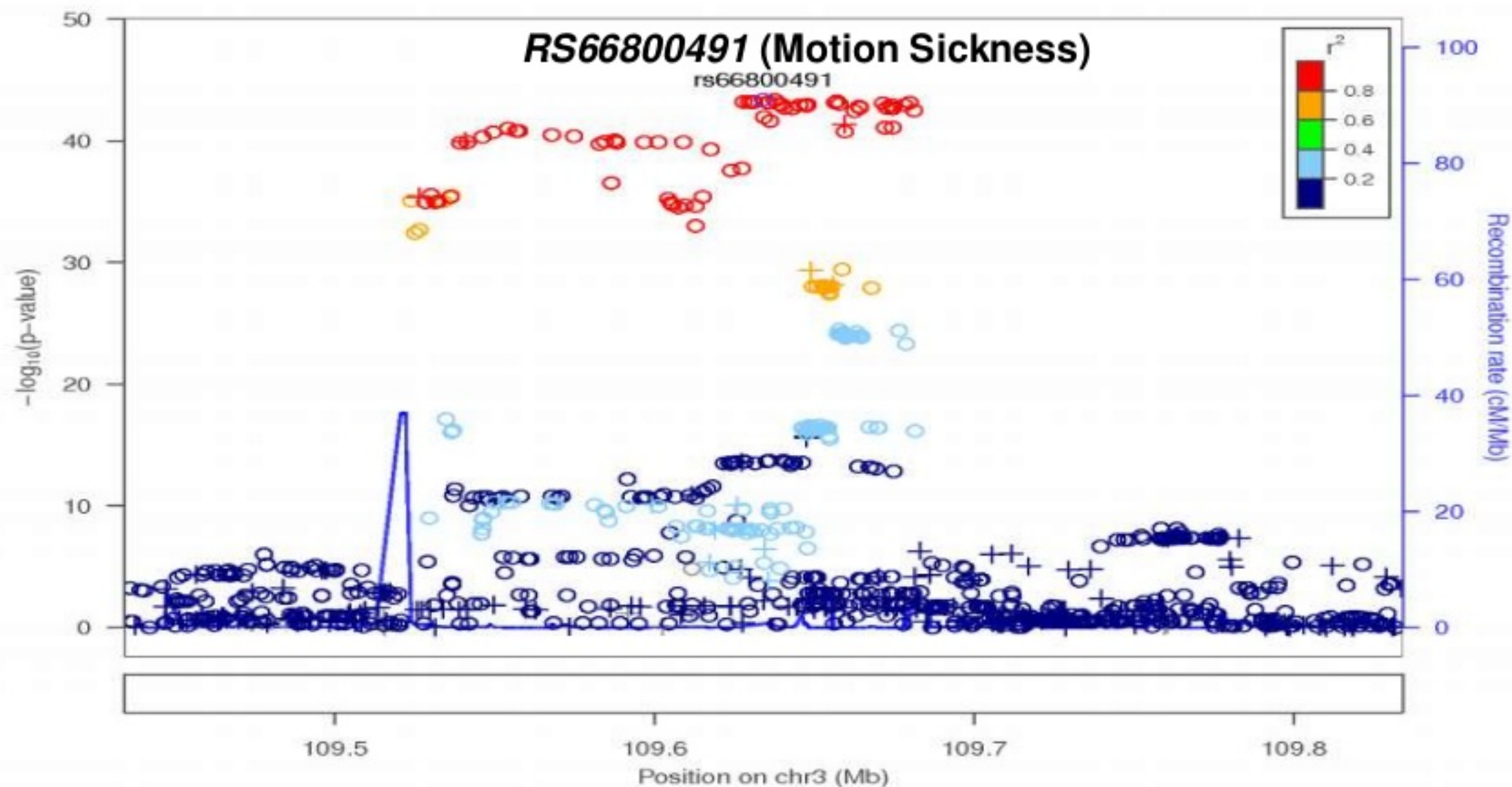
- Identified ~13,000 genetic variants (single nucleotide mutations/polymorphisms) to be associated with ~2,000 diseases/phenotypes

- **91% of disease-associated genetic variants are located in non-protein-coding regions; used to call “junk DNA”**
- **Unknown function, difficult to translate findings into clinical use**



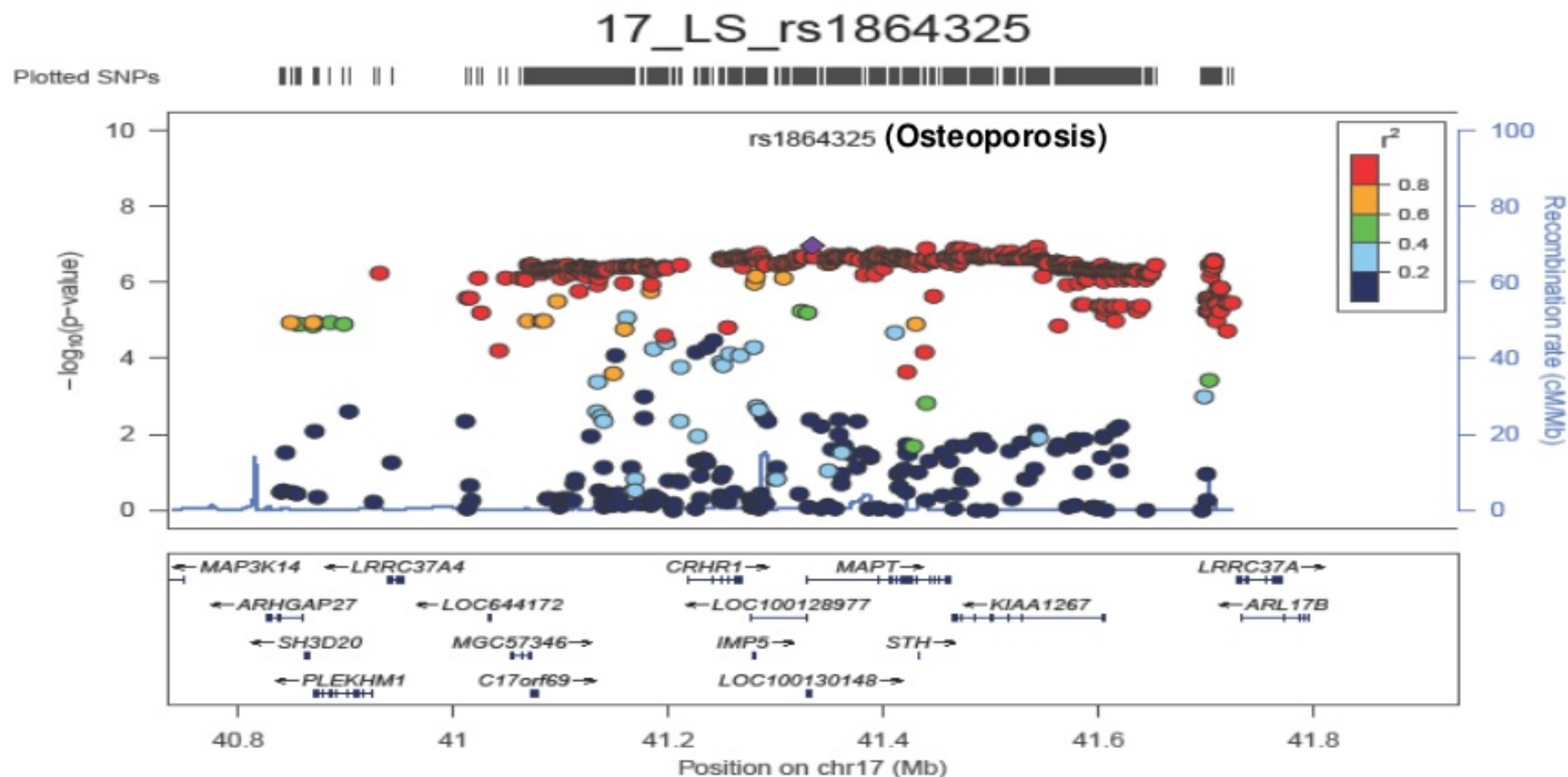


# Associated Variants Located in Gene Desert



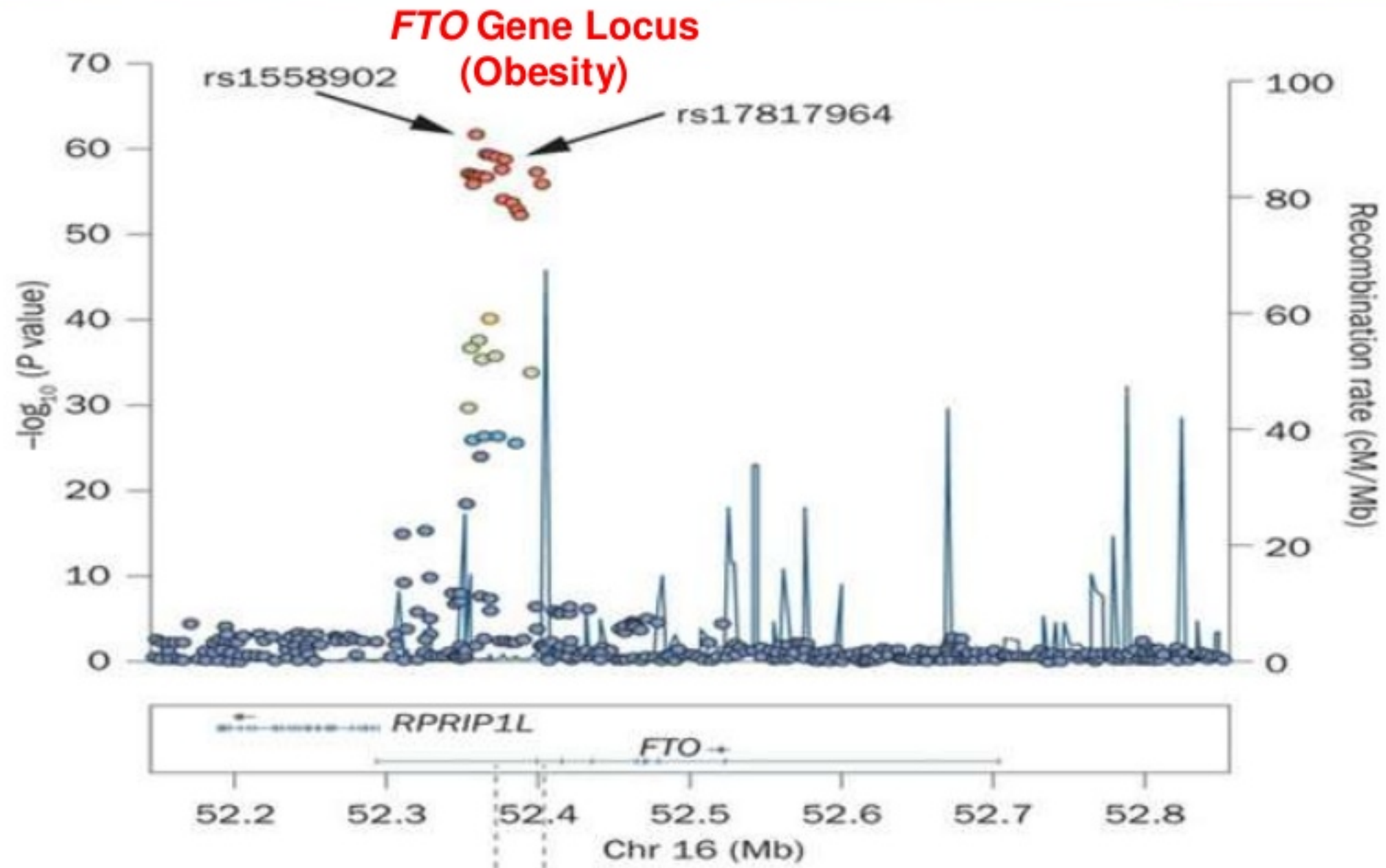
Genetic Coordination: 1D Physical Location on Linear DNA Sequences

# Too Many Genes: Which Gene(s)?

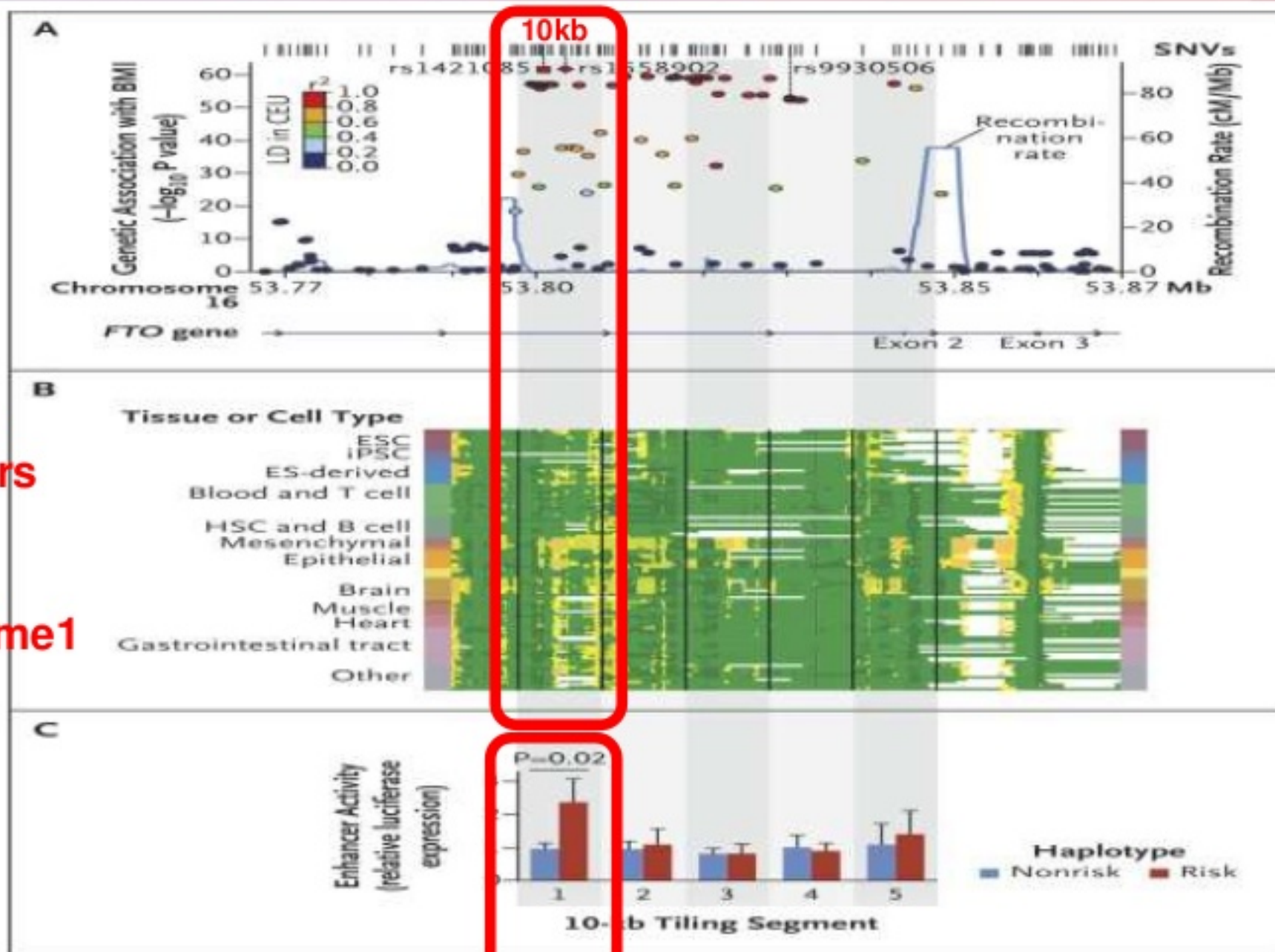




# Associated Variants Located in Introns: Looks Promising?



# Functional Genomics Approaches

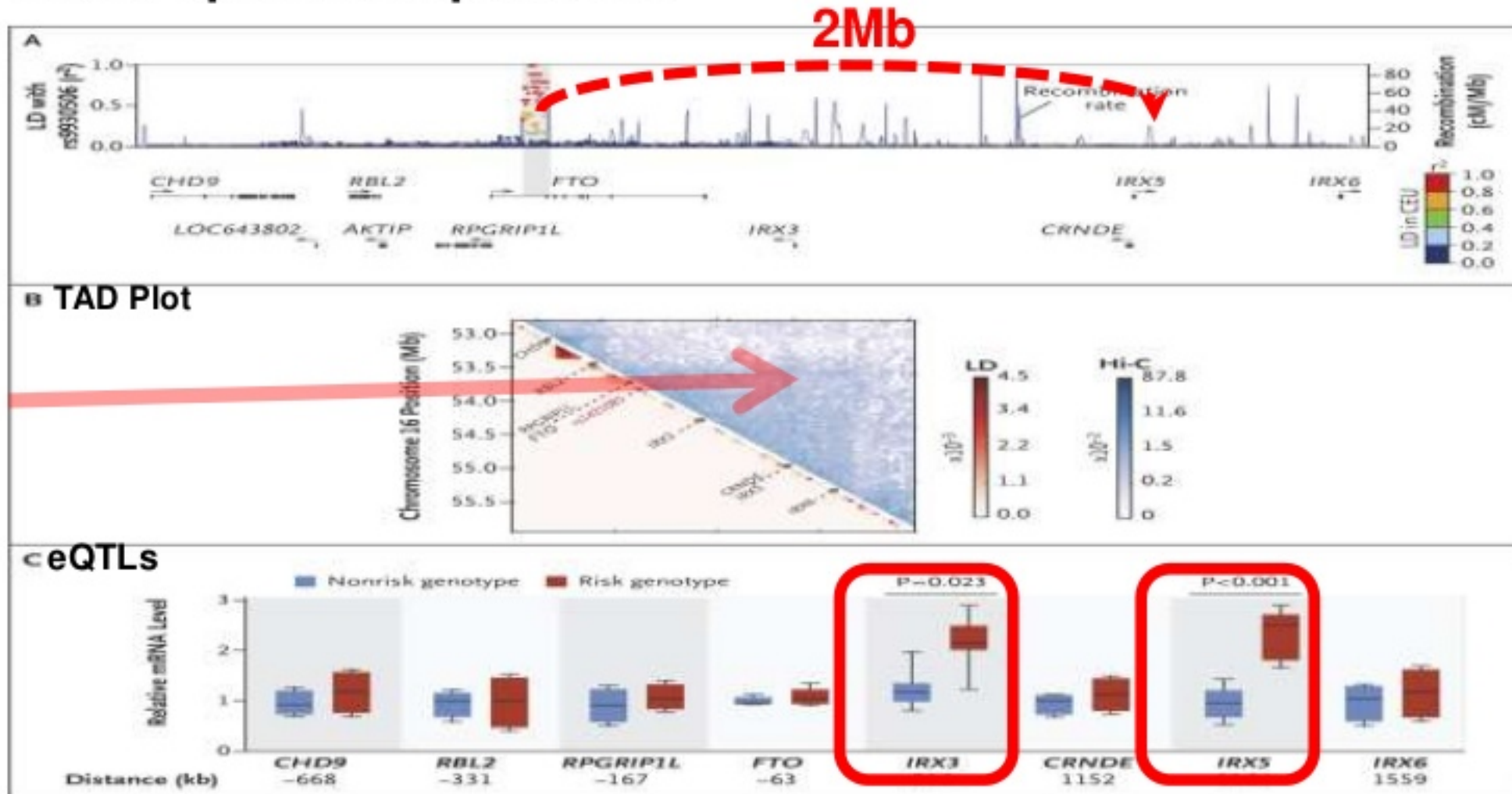


Tissue-Specific  
Active Enhancers  
predicted by  
Histone Marks  
H3K27ac, H3K4me1  
P300



# 3D Genome Interaction Structure with *IRX5* Gene

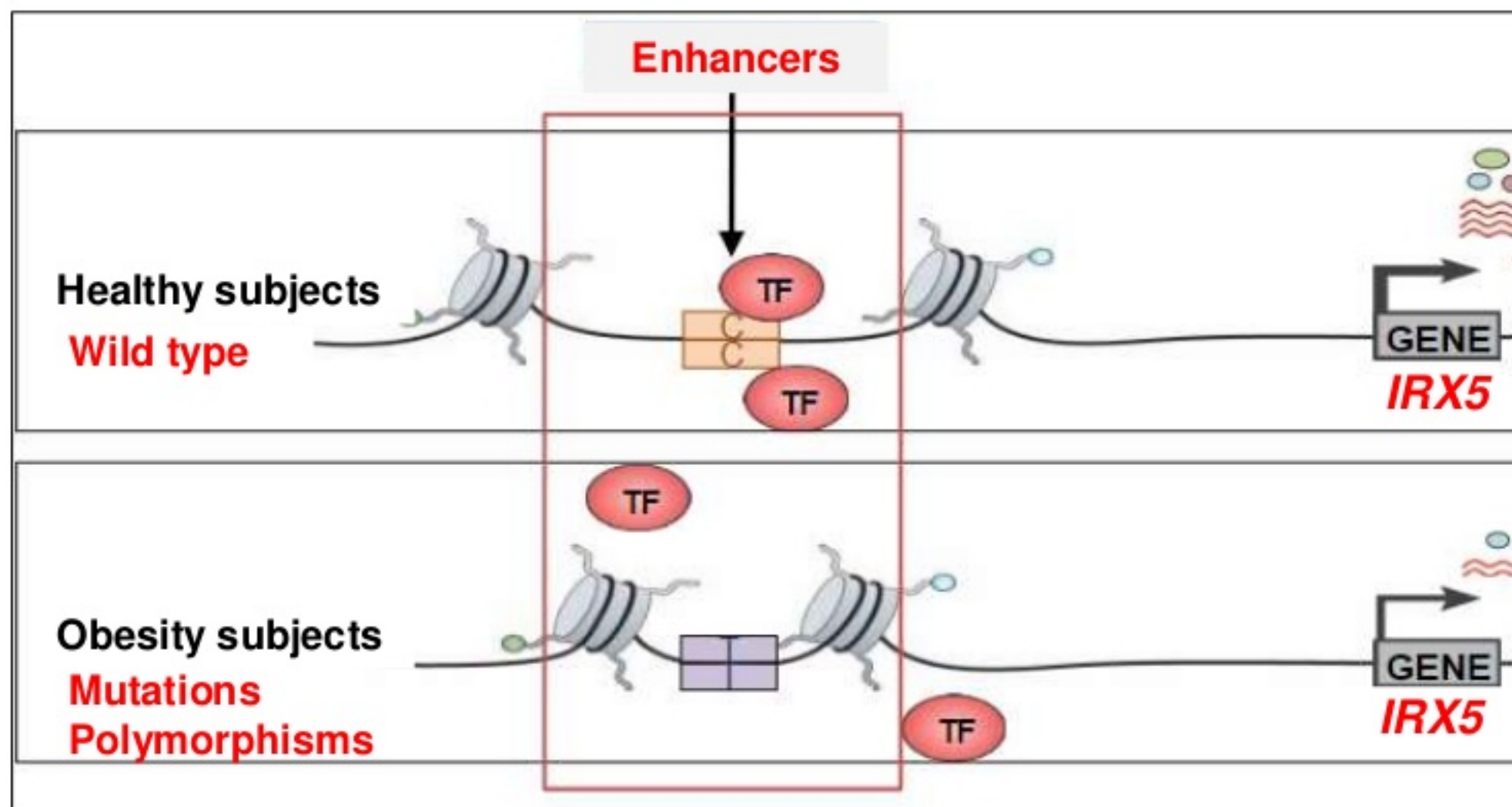
- Tissue-Specific Chromatin Confirmation Capture (3C Tech)
- eQTLs (associations between variants and gene expression)
- Allele-specific expression



Intensity of  
3D Physical  
Interaction  
by Hi-C seq

# *FTO* Genetic Variants and *IRX5* Gene Regulation

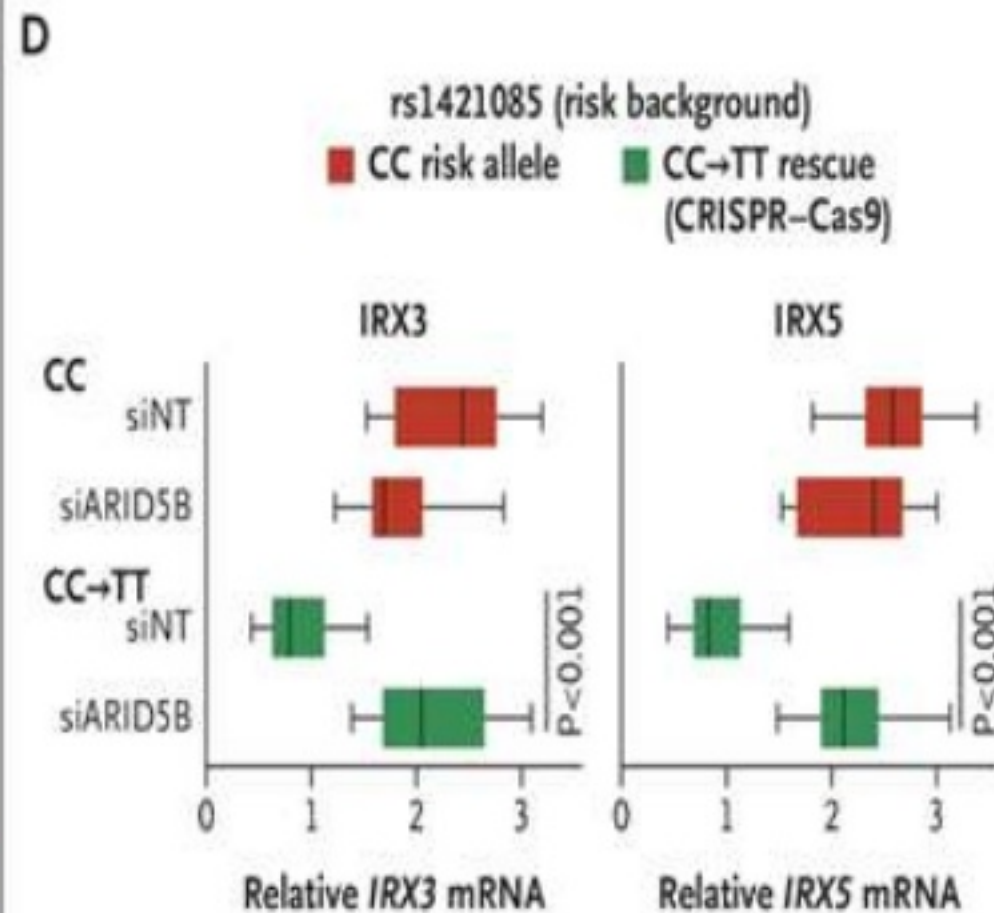
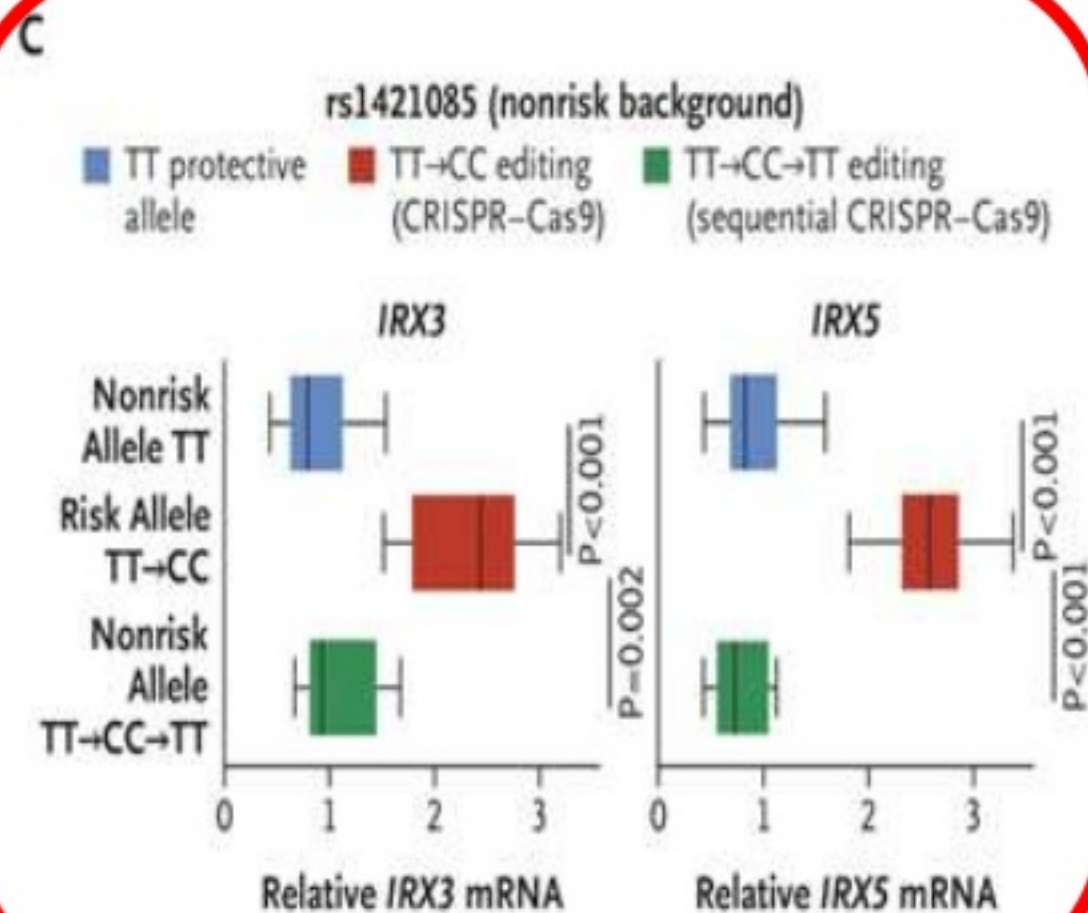
- Obesity associated genetic variants disrupt TF binding and then reduce *IRX5* gene expression



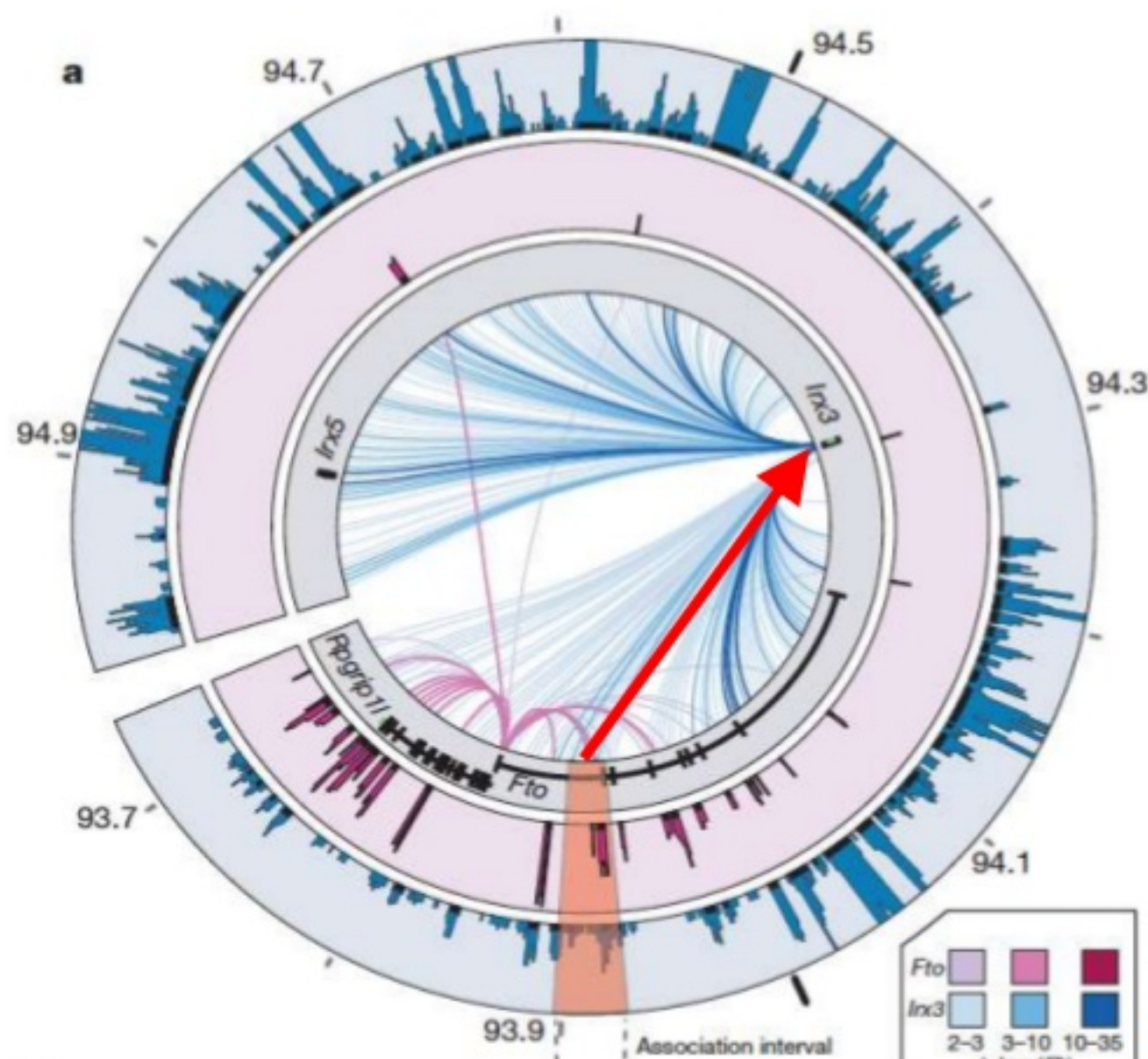


# Gene-Editing: Functional Validation

- Gene Editing by CRISPR/Cas9 in Human adipocytes from subjects carried “risk allele” and subjects carried “protective allele”
- The Risk Allele C: Gain-of-function



# FTO Variants Link to *Irx3* Gene in Brain

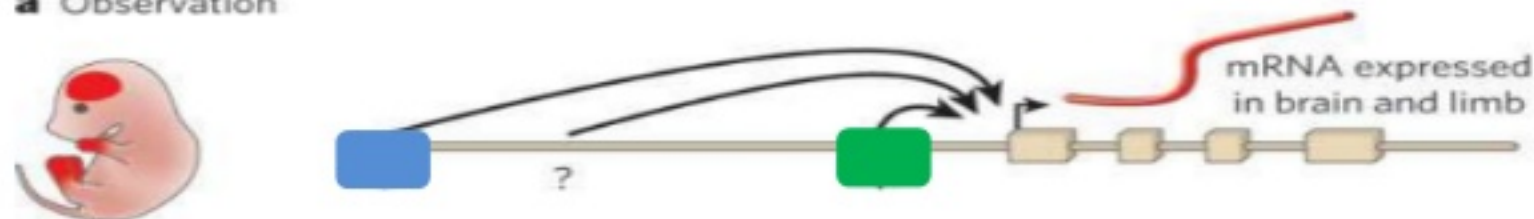


- The obesity associated variants physically interacts with promoter of *Irx3* gene, but **not *Fto***, **not *Irx5*** in mouse brain by 4C-seq
- 4C-seq: Regional Chromatin Confirmation Capture (3C Tech)

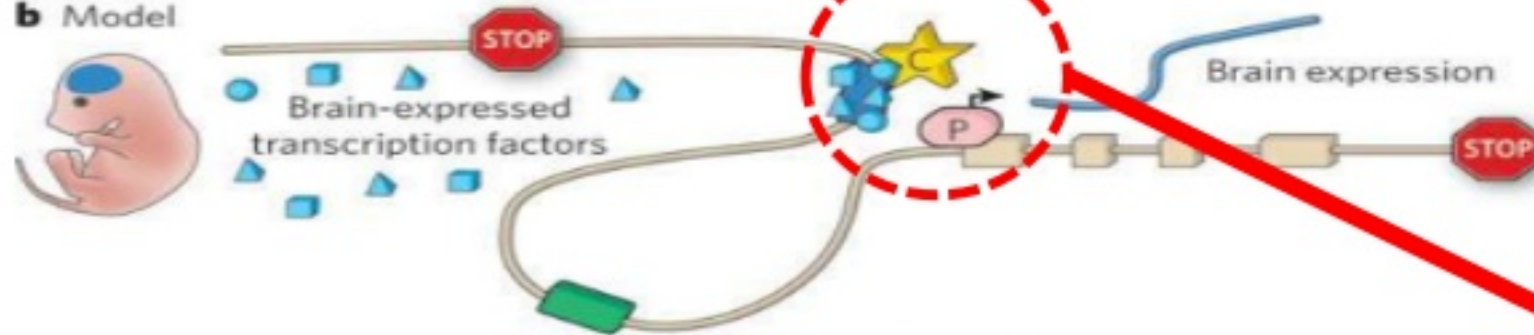


# Gene Regulatory Models

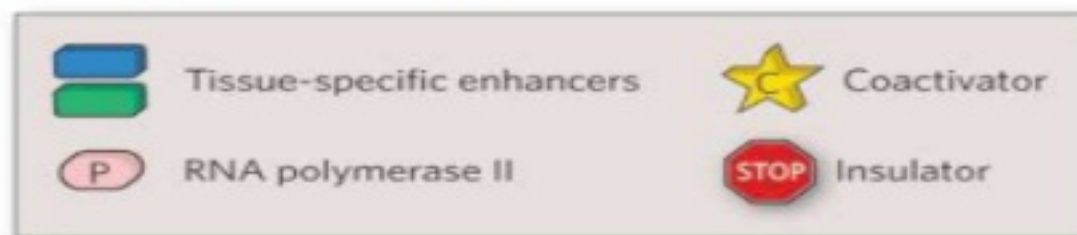
**a** Observation



**b** Model



**c** Model

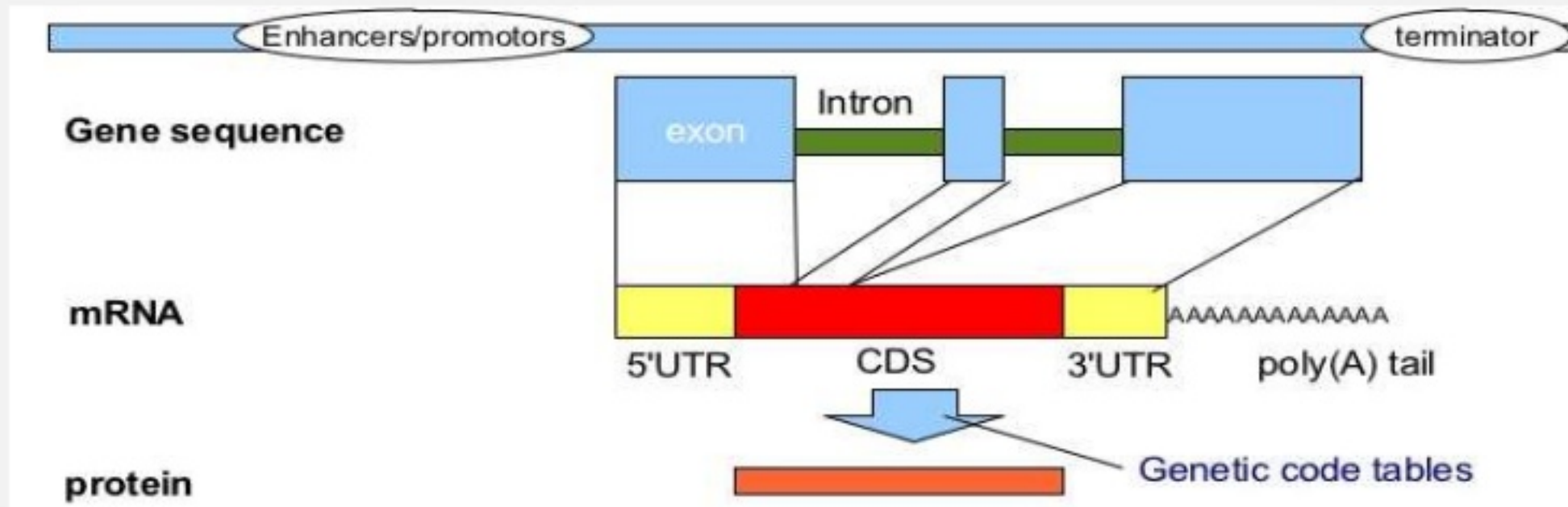


Gene regulatory elements in physical proximity (3D space) with the gene promoters via looping mechanisms

**Tissue (Cell)-Specific DNA Loops: Enhancer-Promoter Interactions**

# Genome-Widely Identify/Predict Targeted Genes?

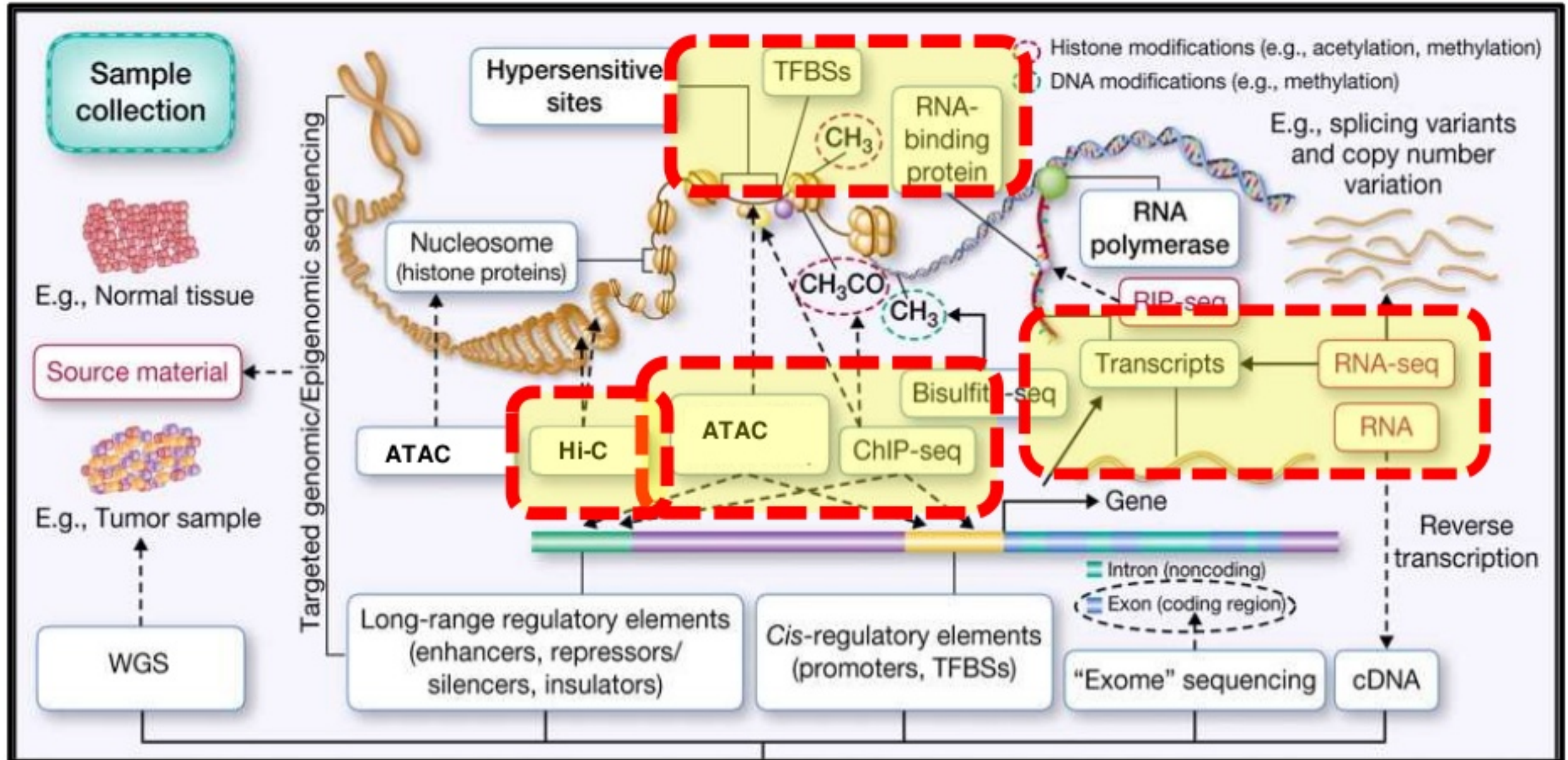
- Identified ~13,000 genetic variants (single nucleotide mutations/polymorphisms) to be associated with ~2,000 diseases/phenotypes
- **91% of disease-associated genetic variants are located in non-protein-coding regions; used to call “junk DNA”**
- **Unknown function, difficult to translate findings into clinical use**
- **May involve in tissue/cell type-specific gene regulation**







# Building Tissue-Specific Gene Regulatory Circuits





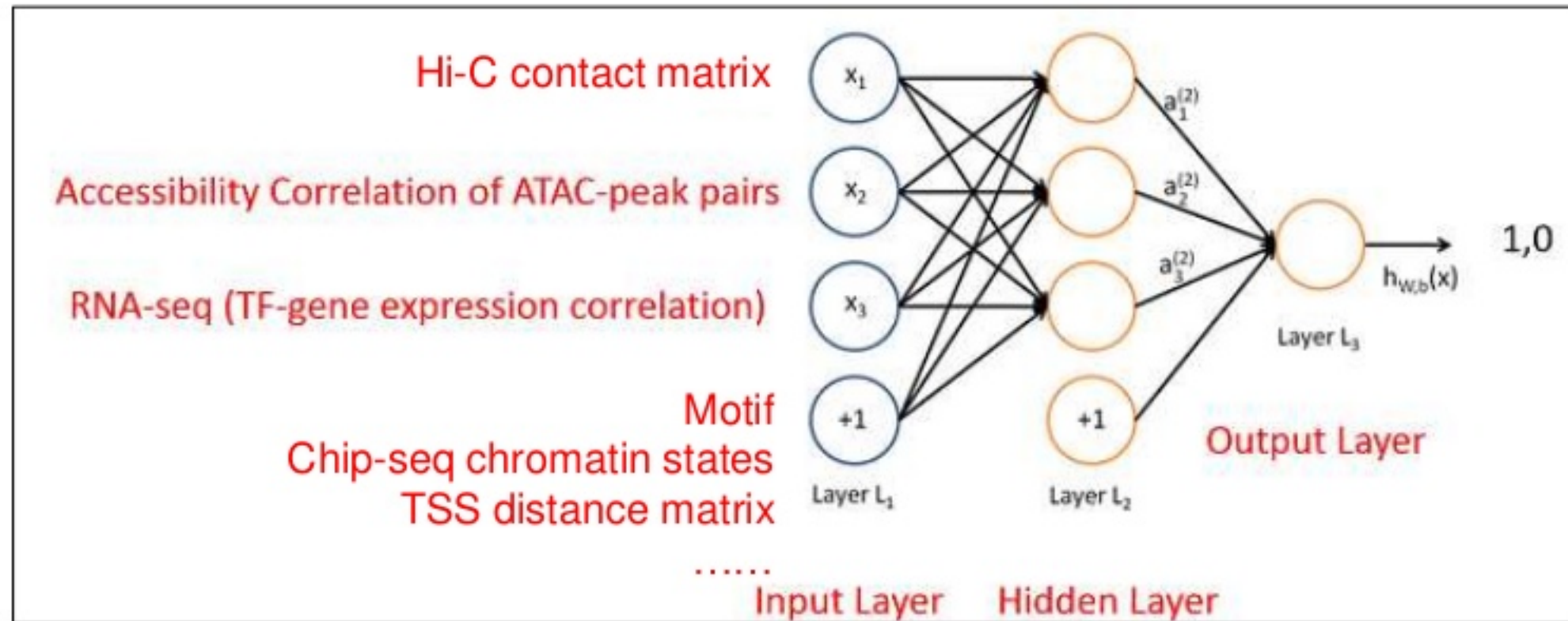
# Building Gene Regulatory Circuits On Human Heart

- Omics experiments on normal human primary cardiac fibroblasts and myocytes from atrium and ventricle; HMSC; skeletal muscle cells
- Publicly available (low resolution): Left ventricle, right ventricle and aorta tissues

Experiments	Functions	Notes
ATAC-seq	Active cis-regulatory region	Active TF binding
Hi-C	Chromatin confirmation capture	1.5 to 2 kb resolution (DpnII, 2 Billions Reads, 2Tb)
ChIP-seq	H3K4me3: Aactive promoter	Predicted chromatin states by HMM
	H3K27ac: Active enhancer/promoter	
	CTCF: Insulator	
	Cohesin: Insulator-RAD21	
	Cohesin: Insulator-SMC3	
	H3K27me3: Polycomb repressed/bivalent promoter/enhancer	
	H3K9me3: Heterochromatin	
	H3K36me3: Transcribed region	
RNA-seq	mRNA (active and )	Isoforms; coexpression with TF
	microRNA and small RNA	Enhancer RNA

# Model Gene Regulation with Deep Neural Network (DNN)

- DNN implemented in the TensorFlow to predict enhancer-promoter gene pairs



- **Training sets** (VISTA: enhancer elements are in 100kb of genes):
- 1,564 Enhancer-promoter gene pairs (the positive set) functionally validated to have regulatory relationships in mouse models
  - 1,207 EP pairs without regulatory relationships (the negative set)



# Acknowledgements

