

# Integrative functional annotations of the human genome and their applications in GWAS analysis

Qiongshi Lu  
University of Wisconsin-Madison



Friday Harbor, 09/08/2017



## 1. Background

## 2. Introduction to (narrow-sense) functional annotations

- Functional annotation in protein-coding genes
- Functional annotation in non-coding regions
- Other useful tools

## 3. Applications of functional annotations

- Functional SNP fine-mapping
- Partitioning heritability and genetic covariance
- Gene-level analysis
- Effect size estimation and risk prediction



# Background

Despite the advancements, GWAS has its limitations

## It is difficult to identify all associations

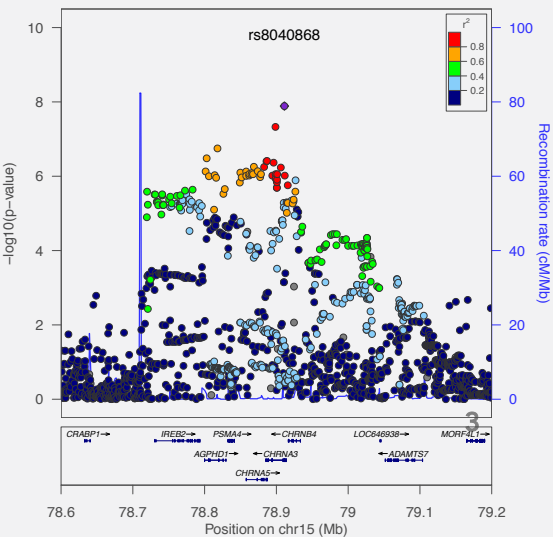
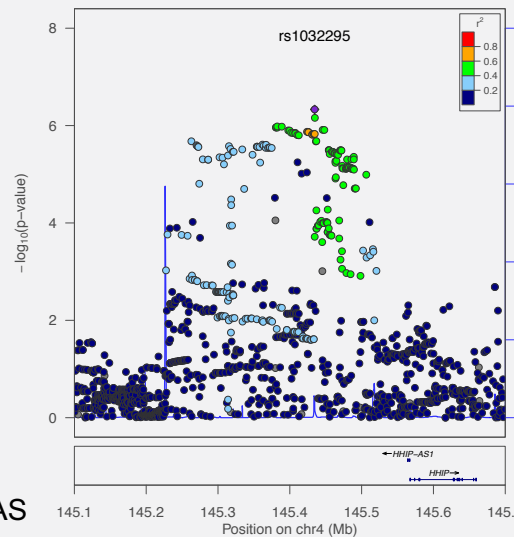
- Polygenicity and low effect size

## It remains challenging to interpret the findings

- 88% of significant associations are in the non-coding genome
- LD makes it challenging to identify biologically functional SNPs



Two examples from COPD GWAS

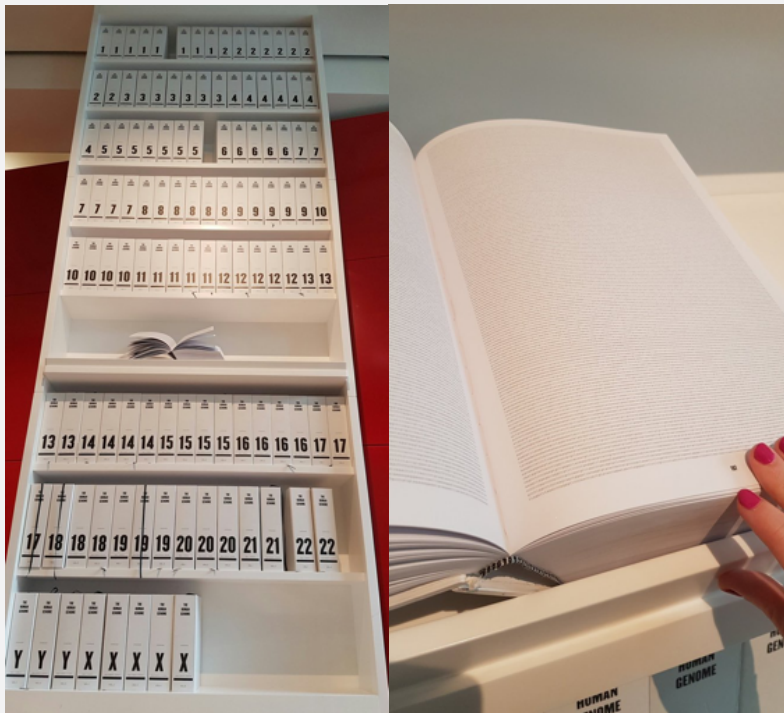


# Background



To solve these problems, we need better **functional annotations** for the (non-coding) human genome.

Only ~2% of the human genome encodes proteins. However, the rest 98% may be critically involved in a variety of regulatory machinery.



@susejohnston

gene?

```
ACGTTGCAAATTCAGTCCGGTACTTTAACCTACCTACCGTACTGGTATTGTCAGGTTGTTCAACT
CATGACACTCCAGATAGACAGATTGTCGGTGTATVVGACTTGGAACTGTAGGCCCTTGAATCT
TGGCAGTCCCTACGTACCGTCCGGTACTGGTACCGTGGTACGTTGTTCAACTCATCCAGGA
GAAATATCTCCGGATAATTAACAGATACACACCCTTAGACCATTTAATCCCTGGGAAAAGGCACTA
CGTACCAGTCTTTCCAGGCCACTGACAGATAGACAGATTGTCGGTATVVGACTTGGAACTGTA
GGCCCTTGAATCTTGGCAGTCCGTAACGTACGTACGGTACTGGTAACTGAGGTCAGGTTGTTT
AACTCATCGTGACTGATTACCAGGATCCTAGCCGGATCCTACTGACCTGACGTACGTAATGCAGT
GGTCAAGTTGTTCAACTCGATGACTAGAATATATCCAGGAAAAATCCCTGGGAAAAAATGGGCC
TACGTGTCGTAACGTACGTACGGTACTGGTAACTGAGCCAGGAAAAATCCCTGGGAAAAAATGG
GGCCCTATCGTGACTGATTACCAGGATCCTAGCCGGATCCTACTGACCTGACGTACGTAATGCAG
TGGTCAAGTTGTTCAACTCGATGACTAGAATATATCCAGGAAAAAATGGGCCCTACGTACC
GTAACGTTGCAAATTCAGTCCGGTACGTTTCCAGGCTACACATTGTCGGTGTATVVGACTTGGAA
CTGTAGCURLYHAIRGCCCTTGAATCTTGGCAGTCCGTAACGTACGTACGTACGTACGTACGTAC
AACTCATCCAGGAATGGGCCCTACGTACCGTAACTGCAAATTCAGTCCGGTACGTTTCCAGG
CTACACACACACTGACAGATAGACAGATTGTCGGTGTATVVGACTTGGAACTGTAGGCCCTTGA
ATCTTGGCAGTCCGTAACGTACGTACGGTACTGGTAACTGAGGTCAGGTTGTTTCAATACAGGA
TCTACTAGAAGAAAAATGGGCCCTACGTACCGTAACTGCAAATTCAGTCCGGTACGTTTCCAGG
GGCTACACACACACTGACAGATAGACAGATTGTCGGTGTATVVGACTTGGAACTGTAGGCCCTT
GAATCTTGGCAGTCCGTAACGTACGTACGGTACTGHEARTDISLTTTCAACTCATCCAGGAAAT
CCCTGGGAAAAAATTCAGGCTACGTAACCGTAACTGCAAATTCAGTCCGGTACGTTTCCAGGC
TACACACACACTGACAGATAGACAGATTGTCGGTGTATVVGACTTGGAACTGTAGGCCCTTCCAGG
ATGTAATGCACTGGTCAAGGTTGTTCAACTCGATGACTAGAATATATCCAGGAAAAATCCCTGGGA
```

related to neurological function?

enhancer?



## 1. Background

## 2. Introduction to (narrow-sense) functional annotations

- Functional annotation in protein-coding genes
- Functional annotation in non-coding regions
- Other useful tools

## 3. Applications of functional annotations

- Functional SNP fine-mapping
- Partitioning heritability and genetic covariance
- Gene-level analysis
- Effect size estimation and risk prediction

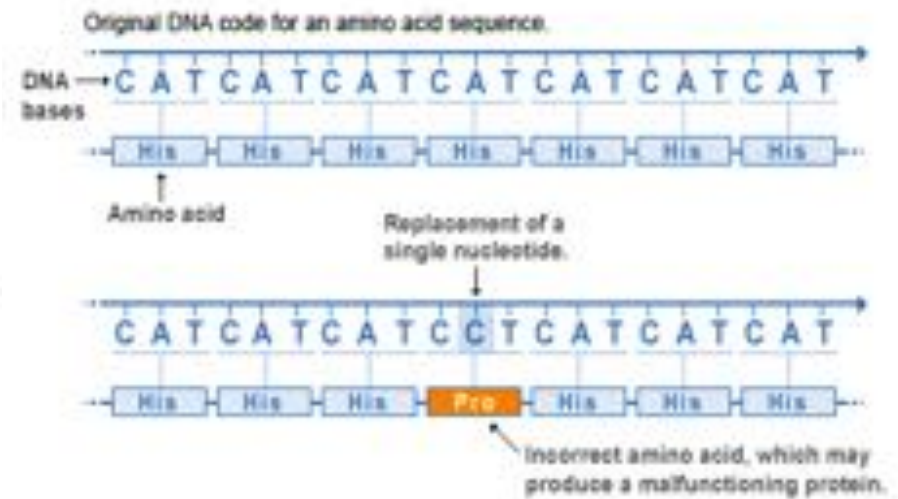


# Functional annotation in protein-coding genes

- Synonymous variant
- Missense variant

		Second Letter					
		U	C	A	G		
1st letter	U	UUU   Phe UUC UUA   Leu UUG	UCU   Ser UCC UCA UCG	UAU   Tyr UAC UAA Stop UAG Stop	UGU   Cys UGC UGA Stop UGG Trp	3rd letter	U C A G
	C	CUU   Leu CUC CUA CUG	CCU   Pro CCC CCA CCG	CAU   His CAC CAA   Gln CAG	CGU   Arg CGC CGA CGG		U C A G
	A	AUU   Ile AUC AUA AUG   Met	ACU   Thr ACC ACA ACG	AAU   Asn AAC AAA   Lys AAG	AGU   Ser AGC AGA   Arg AGG		U C A G
	G	GUU   Val GUC GUA GUG	GCU   Ala GCC GCA GCG	GAU   Asp GAC GAA   Glu GAG	GGU   Gly GGC GGA GGG		U C A G

## Missense mutation



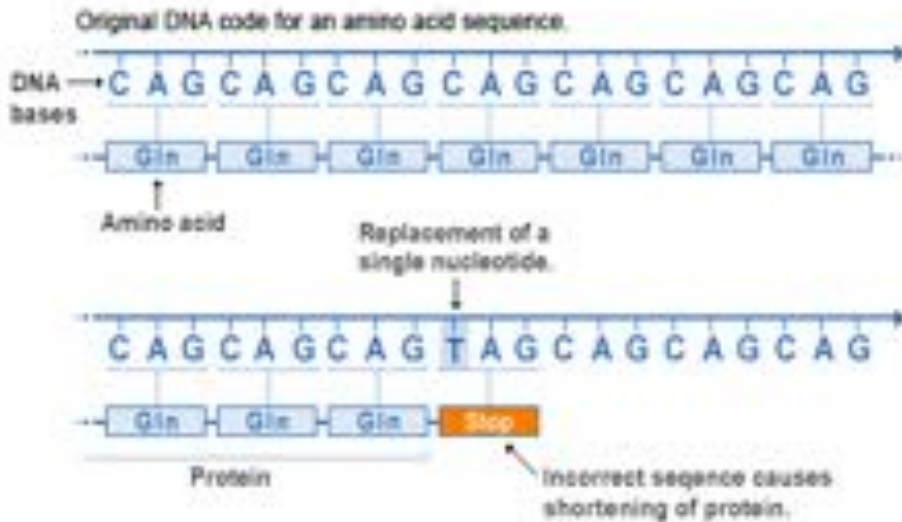
U.S. National Library of Medicine



# Functional annotation in protein-coding genes

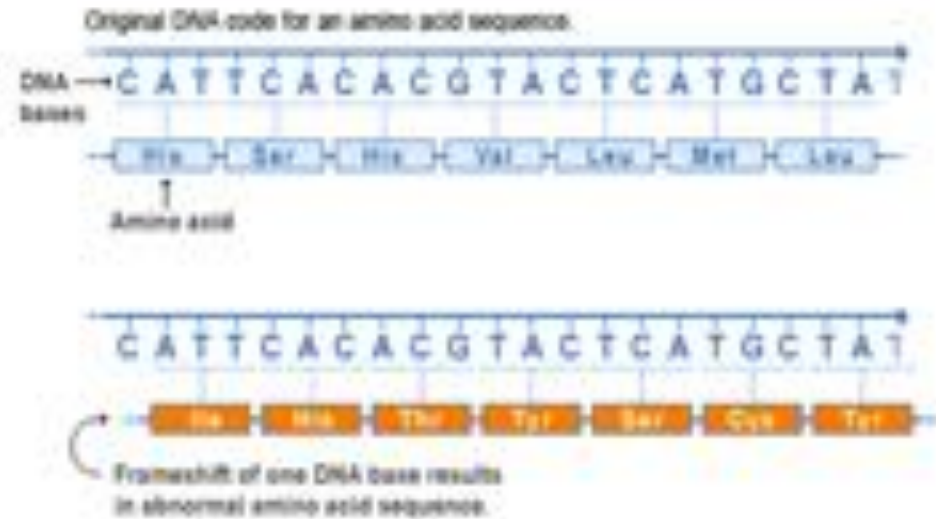
- Loss-of-function variant

## Nonsense mutation



U.S. National Library of Medicine

## Frameshift mutation



U.S. National Library of Medicine

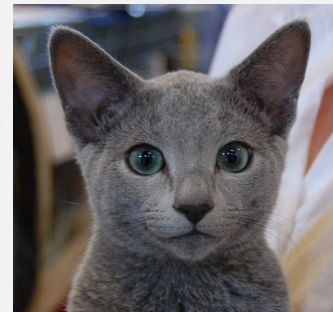
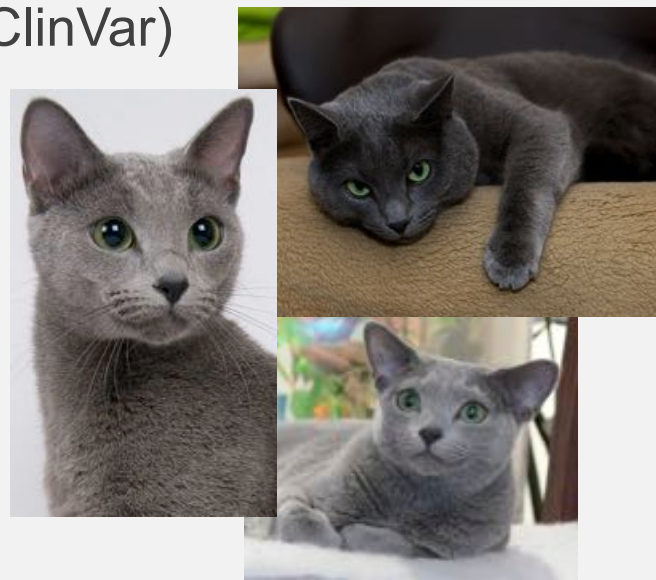


# Functional annotation in protein-coding genes

## Computation methods (supervised)

We understand the functional mechanism of genes. Training data are also available (OMIM, ClinVar)

- SIFT
- PolyPhen2
- MetaSVM



Maine Coon  
?  
Russian Blue





# Functional annotation in protein-coding genes

## Application – de novo mutation analysis

**Science** AAAS

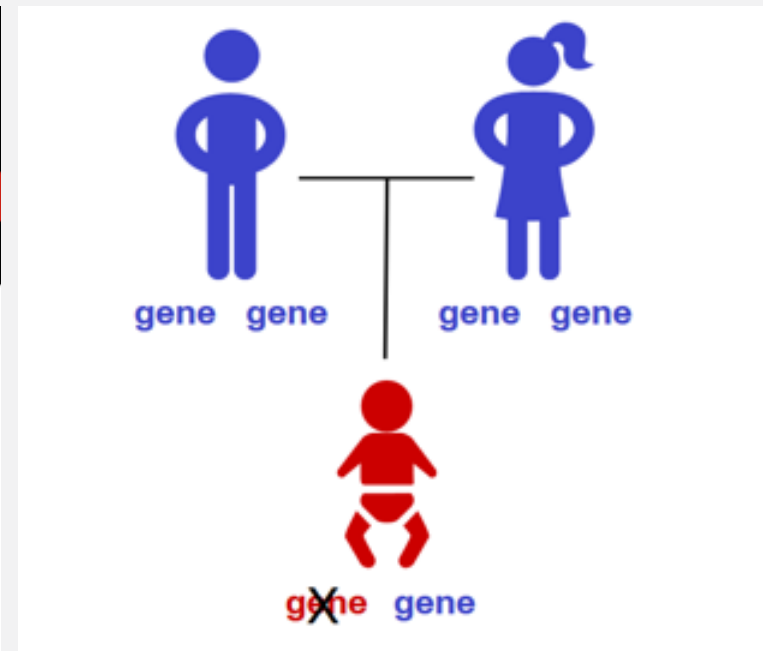
Home News Journals Topics Careers

Science Science Advances Science Immunology Science Robotics Science Signaling Science Translational Medicine

**SHARE** **REPORT**

**De novo mutations in congenital heart disease with neurodevelopmental and other congenital anomalies**

Jason Homsy<sup>1,2,\*</sup>, Samir Zaidi<sup>3,\*</sup>, Yufeng Shen<sup>4,\*</sup>, James S. Ware<sup>1,5,6,\*</sup>, Kaitlin E. Samocha<sup>1,7</sup>, Konrad J. Karczewski<sup>1,7</sup>, Steven R. DePalma<sup>1,8</sup>, David McKean<sup>1</sup>, Hiroko Wakimoto<sup>1</sup>, Josh Gorham<sup>1</sup>, Sheng Chih Jin<sup>3</sup>, John Deanfield<sup>9</sup>, Alessandro Giardini<sup>9</sup>, George A. Porter Jr.<sup>10</sup>, Richard Kim<sup>11</sup>, Kaya Bilguvar<sup>3,12</sup>, Francesc López-Giráldez<sup>13</sup>, Irina Tikhonova<sup>17</sup>, Shrikant Mane<sup>12</sup>, Angela Romano-Adesman<sup>13</sup>, Hongjian Qi<sup>14,15</sup>, Badri Vardarajan<sup>15</sup>, Lijiang Ma<sup>16</sup>, Mark Daly<sup>1,7</sup>, Amy E. Roberts<sup>17</sup>, Mark W. Russell<sup>18</sup>, Seema Mital<sup>19</sup>, Jane W. Newburger<sup>20</sup>, J. William Gaynor<sup>20</sup>, Roger E. Breitbart<sup>20</sup>, Ivan Iossifov<sup>22</sup>, Michael Ronemus<sup>22</sup>, Stephan J. Sanders<sup>23</sup>, Jonathan R. Kallman<sup>24</sup>, Jonathan G. Seidman<sup>1</sup>, Martina Brueckner<sup>3,4</sup>, Bruce D. Gelb<sup>25,1</sup>, Elizabeth Goldmuntz<sup>26,27,1</sup>, Richard P. Lifton<sup>3,28,1,1</sup>, Christine E. Seidman<sup>1,8,29,1,1</sup>, Wendy K. Chung<sup>30,1,1</sup>





# Functional annotation in protein-coding genes

## Application – de novo mutation analysis

	Cases, N = 1213						Controls, N = 900					
	Observed		Expected		Enrichment	P	Observed		Expected		Enrichment	P
	n	Rate	n	Rate			n	Rate	n	Rate		
<b>All genes</b>												
Total	1273	1.05	1312.7	1.08	1.0	0.87	925	1.03	979.7	1.09	0.9	0.96
Synonymous	277	0.23	371.4	0.31	0.7	1	229	0.25	277.4	0.31	0.8	1
Missense	846	0.70	824.9	0.68	1.0	0.24	614	0.68	615.6	0.68	1.0	0.53
D-Mis	212	0.17	133.1	0.11	1.6	$1.8 \times 10^{-10}$	119	0.13	99.3	0.11	1.2	0.03
LoF	150	0.12	116.5	0.10	1.3	<b>0.0016</b>	82	0.09	86.7	0.10	0.9	0.71
Damaging	362	0.30	249.5	0.21	1.4	$1.5 \times 10^{-11}$	201	0.22	186.0	0.21	1.1	0.14
<b>HHE genes</b>												
Total	448	0.37	372.4	0.31	1.2	$7.8 \times 10^{-05}$	271	0.30	277.7	0.31	1.0	0.66
Synonymous	81	0.07	103.5	0.09	0.8	0.99	80	0.09	77.3	0.09	1.0	0.39
Missense	288	0.24	234.3	0.19	1.2	<b>0.00038</b>	163	0.18	174.7	0.19	0.9	0.82
D-Mis	99	0.08	40.6	0.03	<b>2.4</b>	$7.7 \times 10^{-15}$	37	0.04	30.3	0.03	1.2	0.13
LoF	79	0.07	34.5	0.03	<b>2.3</b>	$6.2 \times 10^{-11}$	28	0.03	25.7	0.03	1.1	0.35
Damaging	178	0.15	75.1	0.06	<b>2.4</b>	$5.1 \times 10^{-24}$	65	0.07	55.9	0.06	1.2	0.13
<b>LHE genes</b>												
Total	825	0.68	940.3	0.78	0.9	1	654	0.73	702.1	0.78	0.9	0.97
Synonymous	196	0.16	267.8	0.22	0.7	1	149	0.17	200.1	0.22	0.7	1
Missense	558	0.46	590.5	0.49	0.9	0.91	451	0.50	440.9	0.49	1.0	0.32
D-Mis	113	0.09	92.4	0.08	1.2	0.021	82	0.09	69.0	0.08	1.2	0.069
LoF	71	0.06	82.0	0.07	0.9	0.9	54	0.06	61.1	0.07	0.9	0.83
Damaging	184	0.15	174.4	0.14	1.1	0.24	136	0.15	130.1	0.14	1.1	0.31



## What about non-coding regions?

### DNA conservation



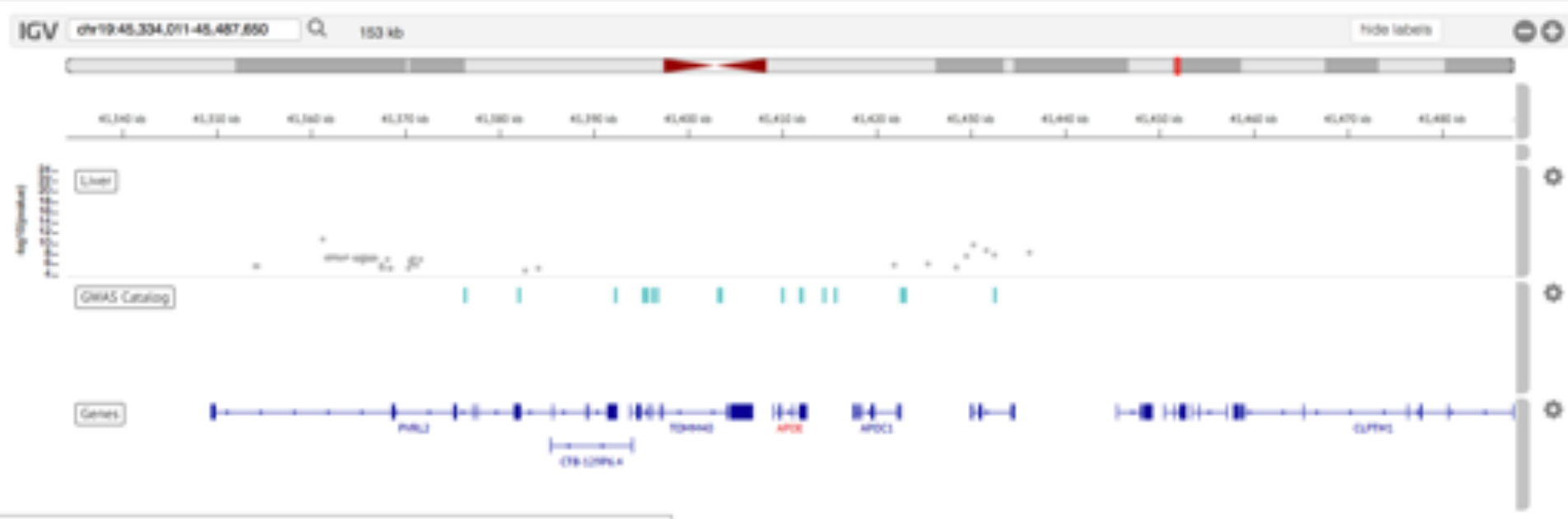
UCSC genome browser



# Functional annotation in non-coding regions

## Transcriptomic information

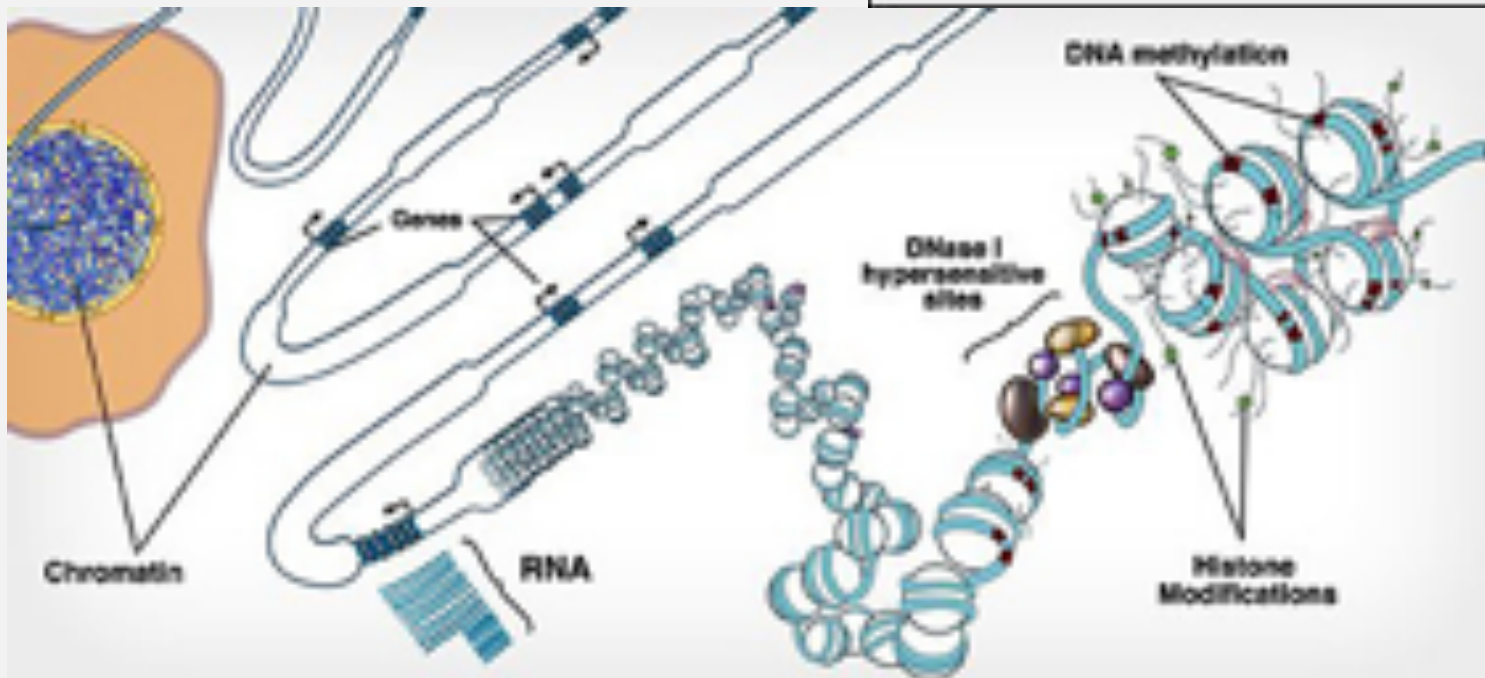
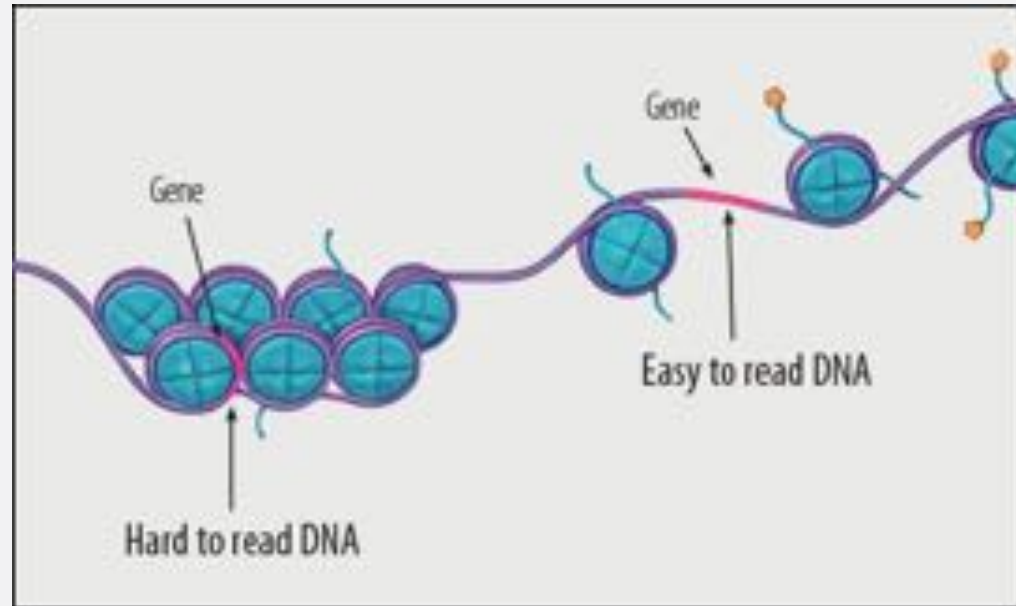
- ncRNA
- eQTL



# Functional annotation in non-coding regions



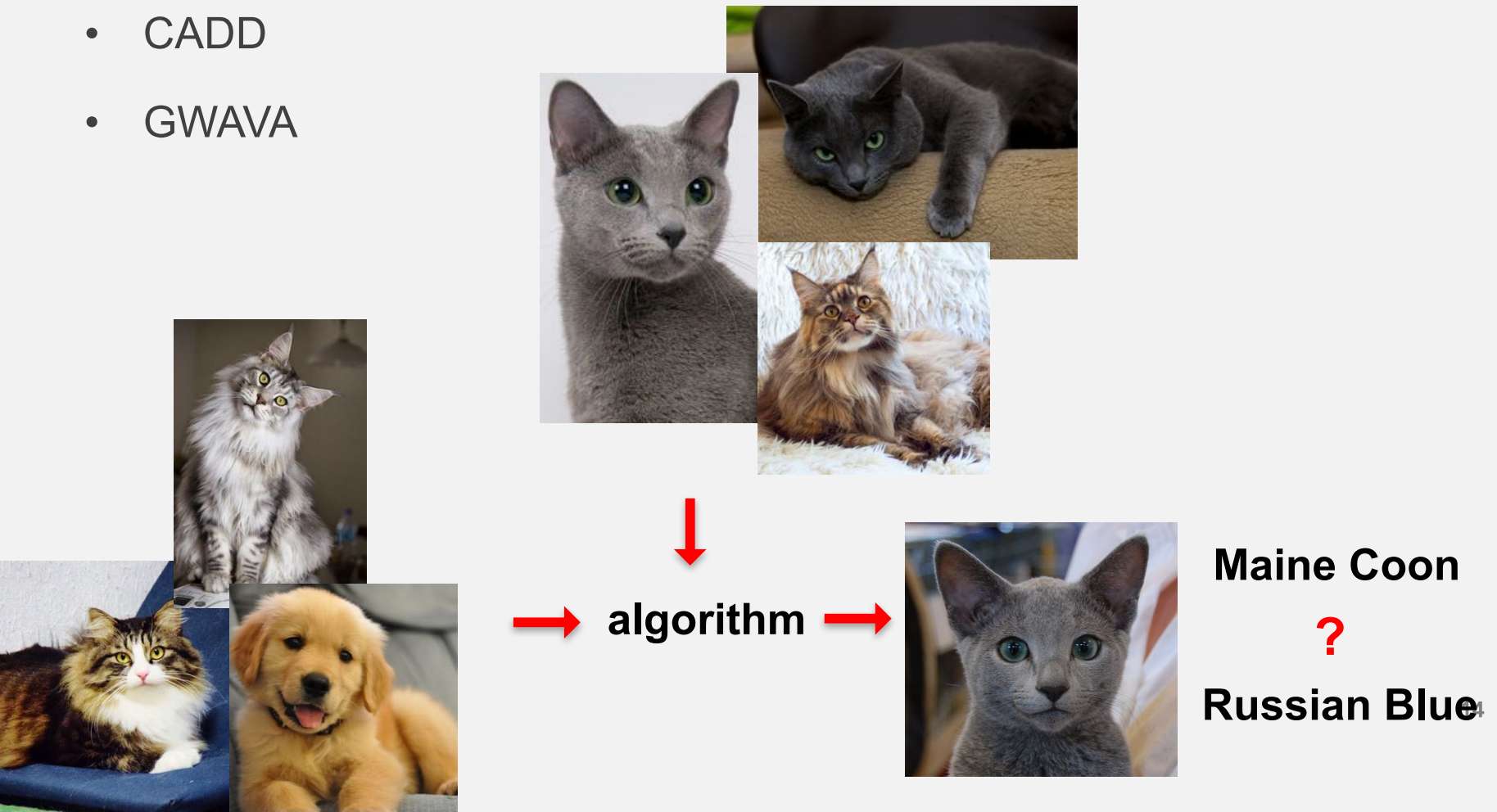
## Epigenetic information



# Functional annotation in non-coding regions

## Computational methods based on supervised learning

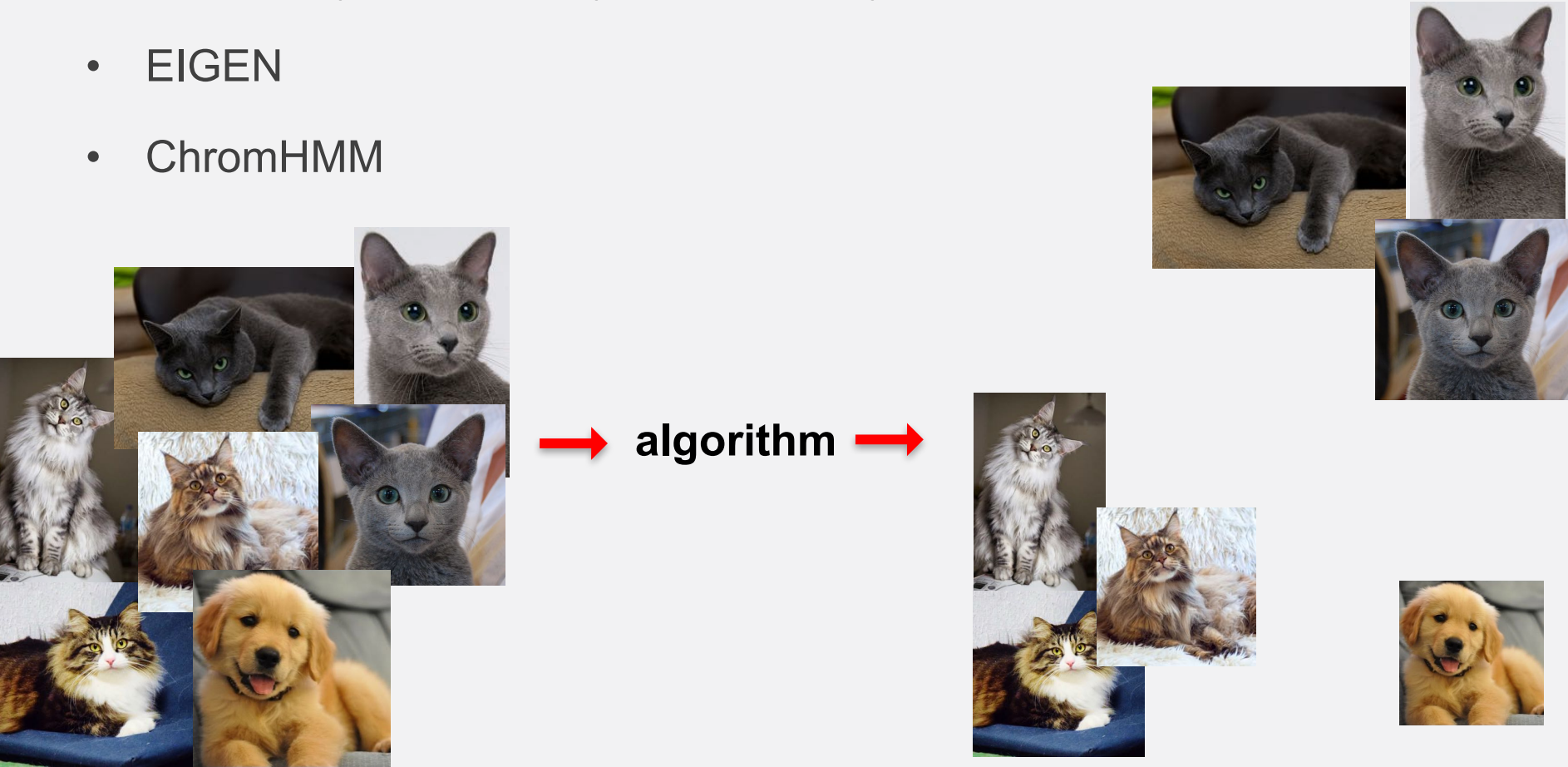
- **Labeled data** + Predictive features + Algorithm = Score
- CADD
- GWAVA



# Functional annotation in non-coding regions

## Computational methods based on unsupervised learning

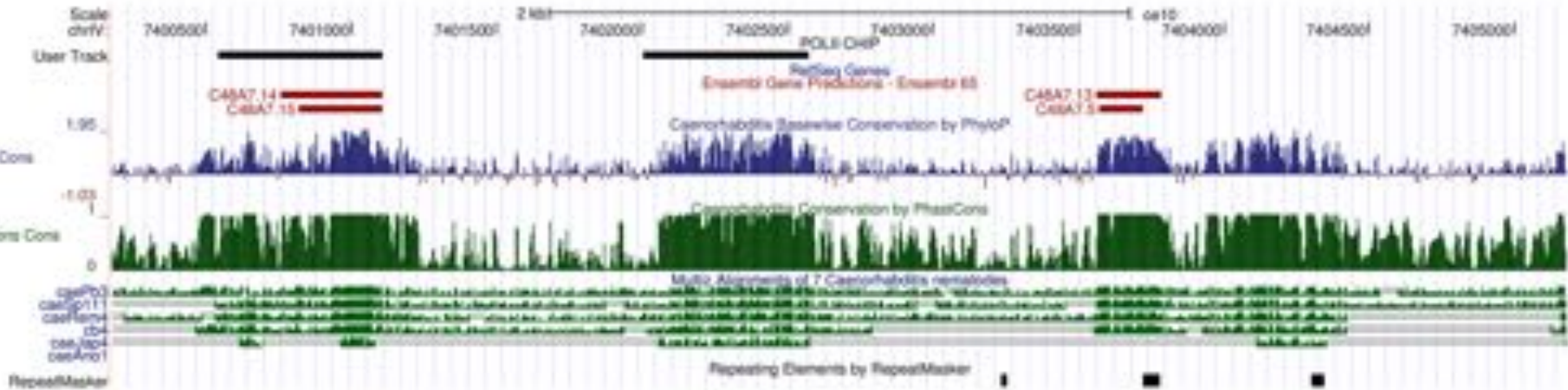
- **Unlabeled data** + Predictive features + Algorithm = Score
- GenoCanyon, GenoSkyline, GenoSkyline-Plus
- EIGEN
- ChromHMM





# Other useful tools

## UCSC genome browser



## Annotvar – annotate your variant list using many functional annotations

Gene_refG	GeneData	ExonicFun	AAChange	pLI.refGer	pRec.refG	pNull.refG	Gene_full	Function_description.refG	Disease_descriptio	Tissue_specificity(Uni	Expression(egz	Expression(GNF/Atla	P[de].refG	P[rec].ref	RVIS.refG
ISG15	NM_005111	.	.	0.009848	0.600249	0.389903	ISG15 ubiq	FUNCTION: Ubiquitin-like	(DISEASE: Immuned	TISSUE SPECIFICITY: De	.	.	0.1	0.22633	-0.11361
ATAD3C	NM_005001	.	.	4.91E-05	0.867306	0.132645	ATPase fa	.	.	.	.	.	0.16989	.	2.882598
NPHP4	NM_005121	.	.	1.29E-17	0.420065	0.579935	nephron	FUNCTION: Involved in the	DISEASE: Note=Cili	TISSUE SPECIFICITY: Ex	.	.	0.12343	0.16808	0.369383
DDR2	.	.	.	0.990992	0.009008	3.79E-09	discoidin	FUNCTION: Tyrosine kinase	(DISEASE: Spondylo	TISSUE SPECIFICITY: De	.	.	0.85011	0.1349	-0.77529
DNASE2B	.	.	.	3.79E-14	0.003092	0.996908	deoxyrib	FUNCTION: Hydrolyzes DN	.	TISSUE SPECIFICITY: Hi	.	.	0.20864	0.10705	0.88176
PRAMEF1	dist=1156	.	.	.	.	.	.	.	.	.	.	.	.	.	.
UBIAD1	dist=4396	.	.	.	.	.	.	.	.	.	.	.	.	.	.
LOC10012	dist=8725	.	.	.	.	.	.	.	.	.	.	.	.	.	.
IL23R	.	nonsynon	IL23R:NM_006495	0.98349	0.004014	interleuki	FUNCTION: Associates wif	DISEASE: Inflamma	TISSUE SPECIFICITY: Ex	.	.	.	0.11254	0.33307	0.795049
ATG16L1	.	nonsynon	ATG16L1:7	0.999737	0.000263	1.67E-11	autophag	FUNCTION: Plays an essen	(DISEASE: Inflamma	.	myocardium;sr dorsal root ganglion;	.	0.3463	0.10646	0.15076





## 1. Background

## 2. Introduction to (narrow-sense) functional annotations

- Functional annotation in protein-coding genes
- Functional annotation in non-coding regions
- Other useful tools

## 3. Applications of functional annotations

- Functional SNP fine-mapping
- Partitioning heritability and genetic covariance
- Gene-level analysis
- Effect size estimation and risk prediction



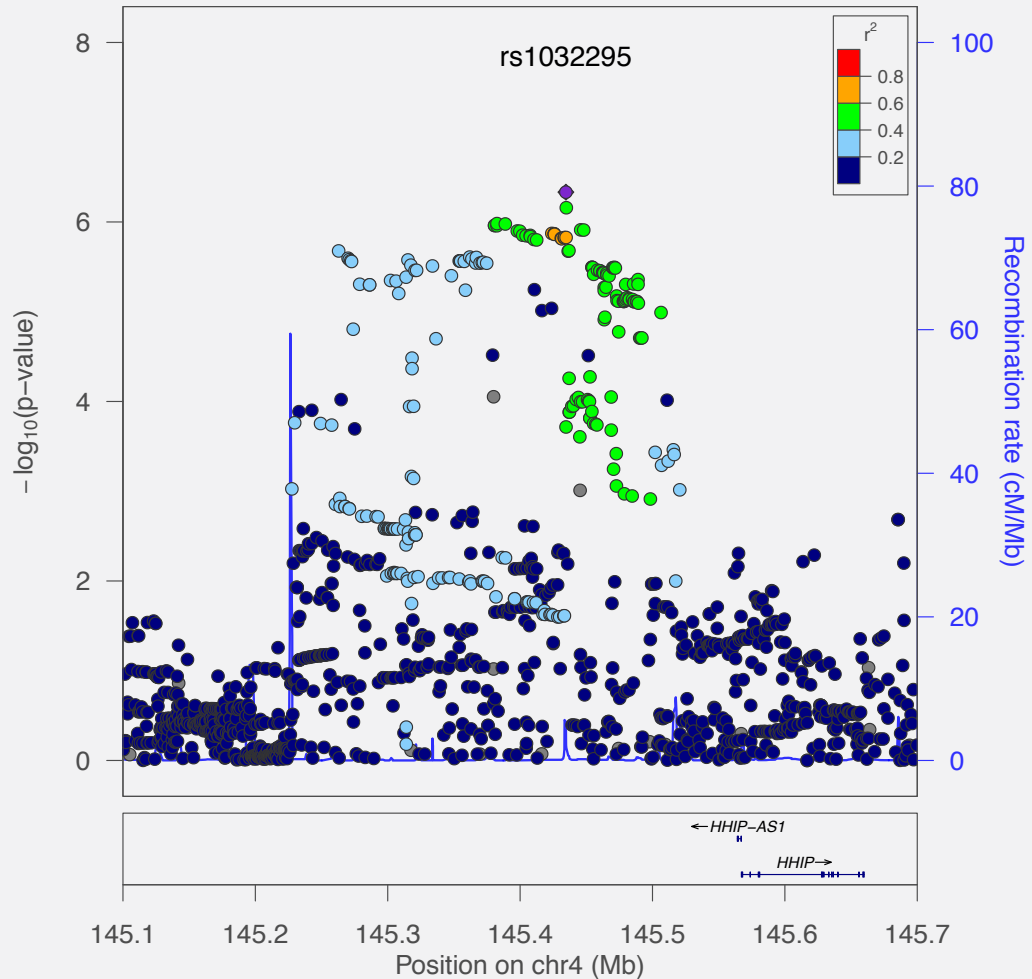
# Functional SNP fine-mapping

## How to identify functional SNP at a GWAS locus?

- PAINTOR
- FGWAS
- GenoWAP

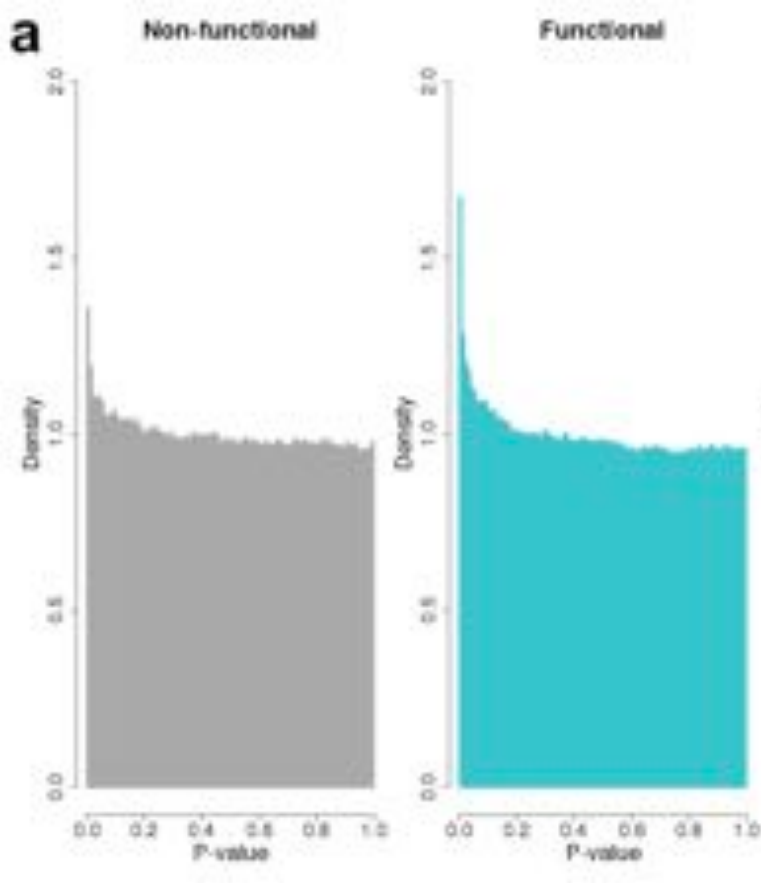
Goal:

$$\mathbb{P}(C \mid P, A)$$

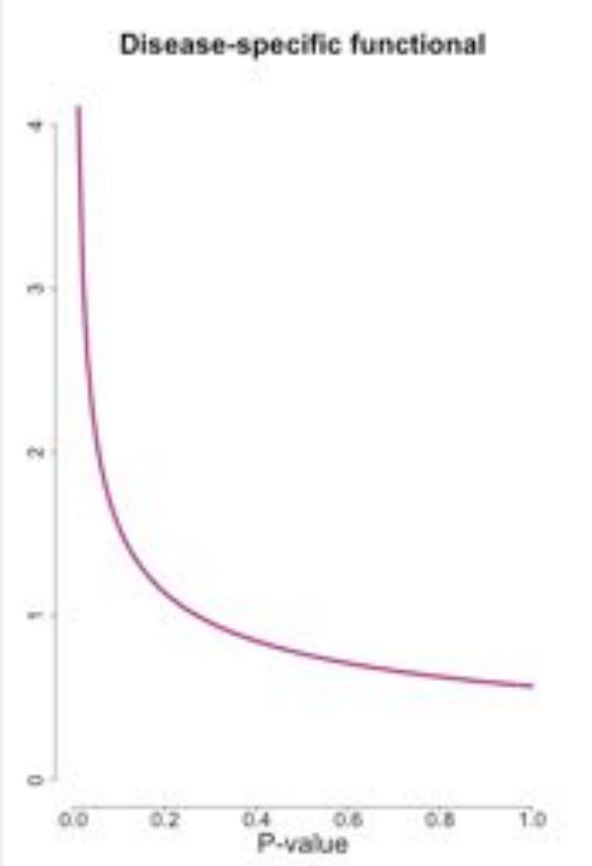


# Functional SNP fine-mapping

- GenoWAP



EM algorithm

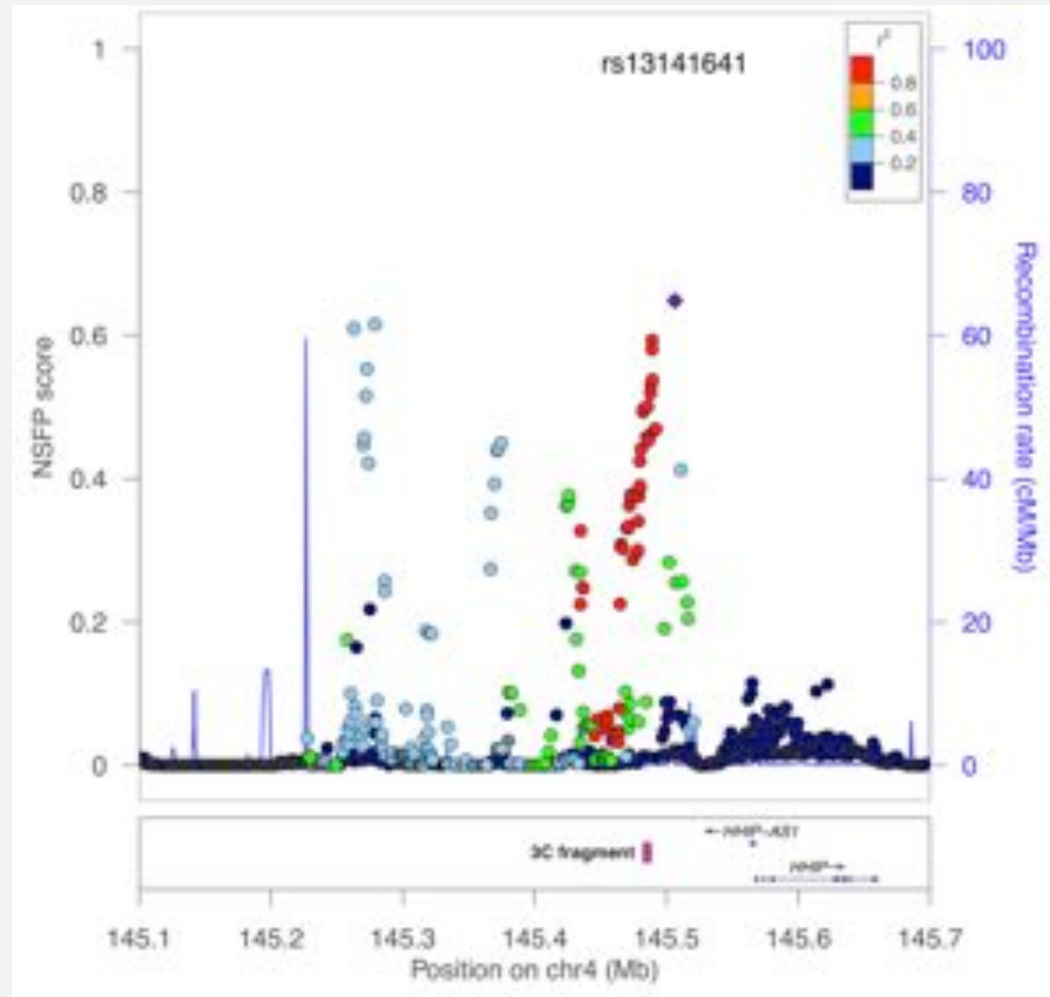
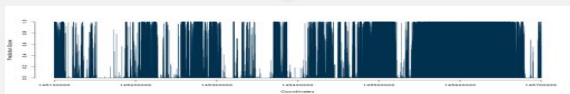
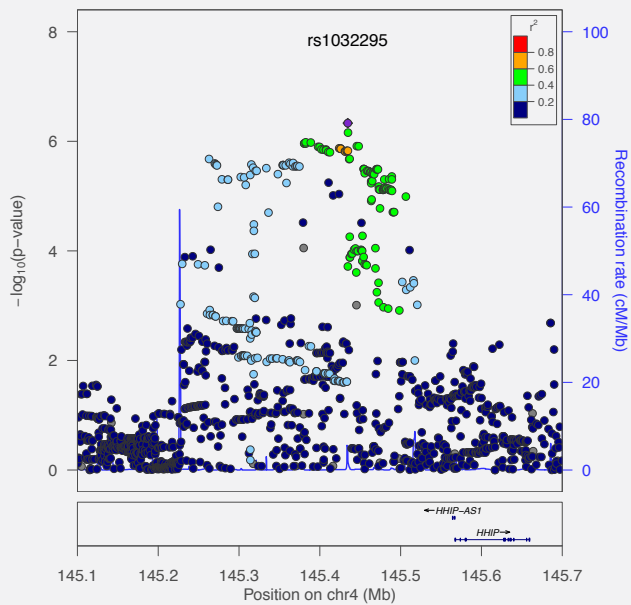


$$\mathbb{P}(C \mid P, A)$$

# Functional SNP fine-mapping



- GenoWAP





## 1. Background

## 2. Introduction to (narrow-sense) functional annotations

- Functional annotation in protein-coding genes
- Functional annotation in non-coding regions
- Other useful tools

## 3. Applications of functional annotations

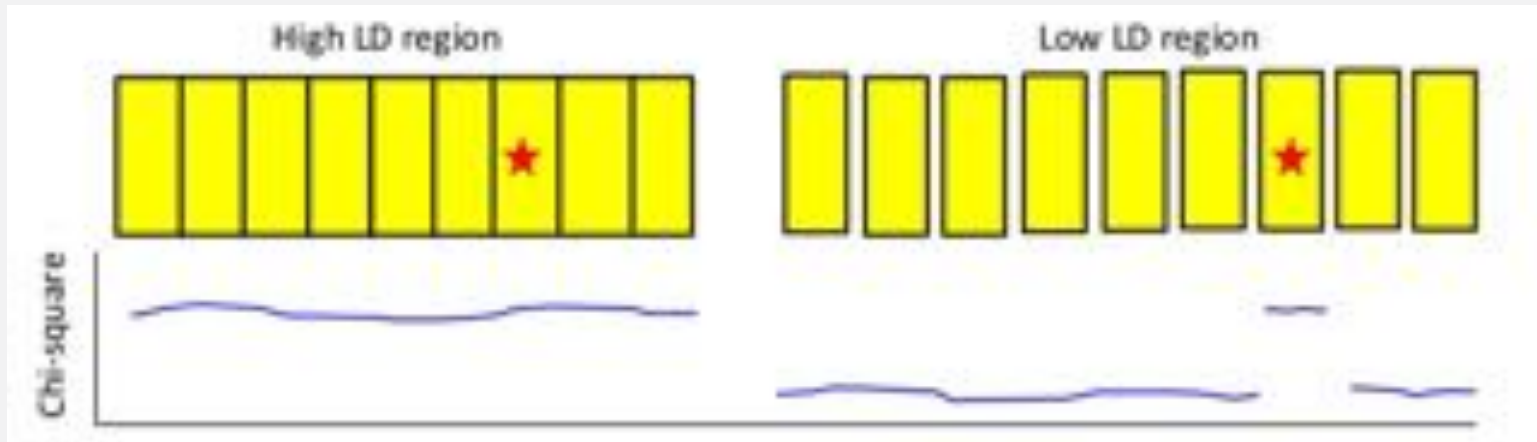
- Functional SNP fine-mapping
- Partitioning heritability and genetic covariance
- Gene-level analysis
- Effect size estimation and risk prediction



# Partitioning heritability and genetic covariance

## LD score regression

- Estimate (partition) heritability using GWAS summary statistics
- **We expect to see stronger associations in regions with high LD**



- It can be shown that this relationship is linear!

Finucane et al. 2015

$$\mathbb{E}(z_j^2) = \frac{Nh^2}{m} l_j + 1$$

sample size  $N$  points to the numerator  $Nh^2$

heritability  $h^2$  points to the numerator  $Nh^2$

number of SNPs  $m$  points to the denominator  $m$

LD score  $l_j$  points to the term  $l_j$

GWAS associations  $z_j^2$  points to the left side of the equation  $\mathbb{E}(z_j^2)$

Bulik-Sullivan et al. 2014

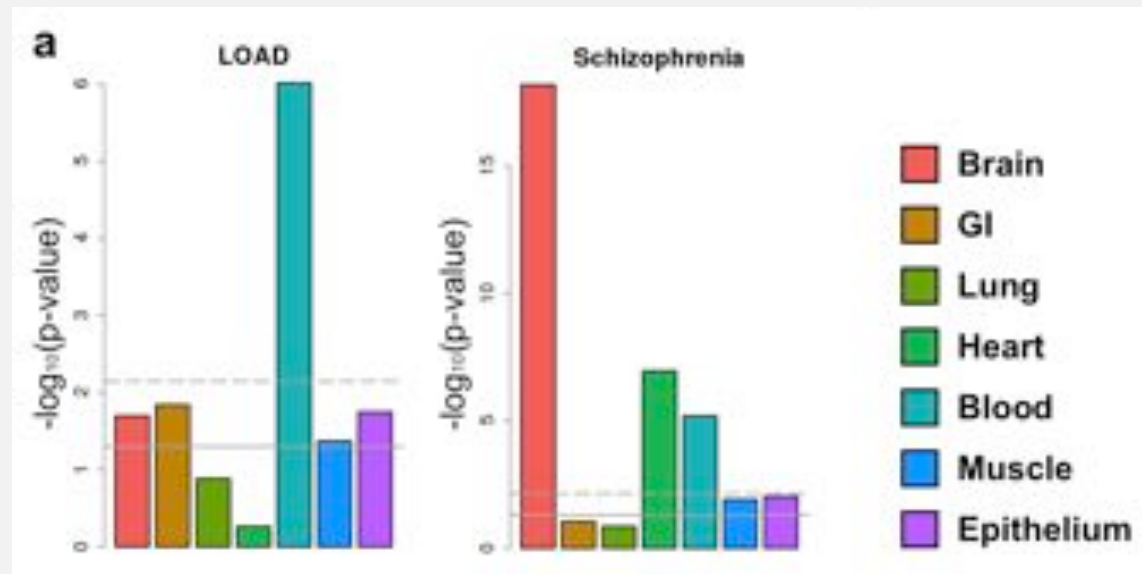


# Partitioning heritability and genetic covariance

## LD score regression

- The model can be extended to partition heritability by **functional annotation**
- This makes it possible to calculate enrichment

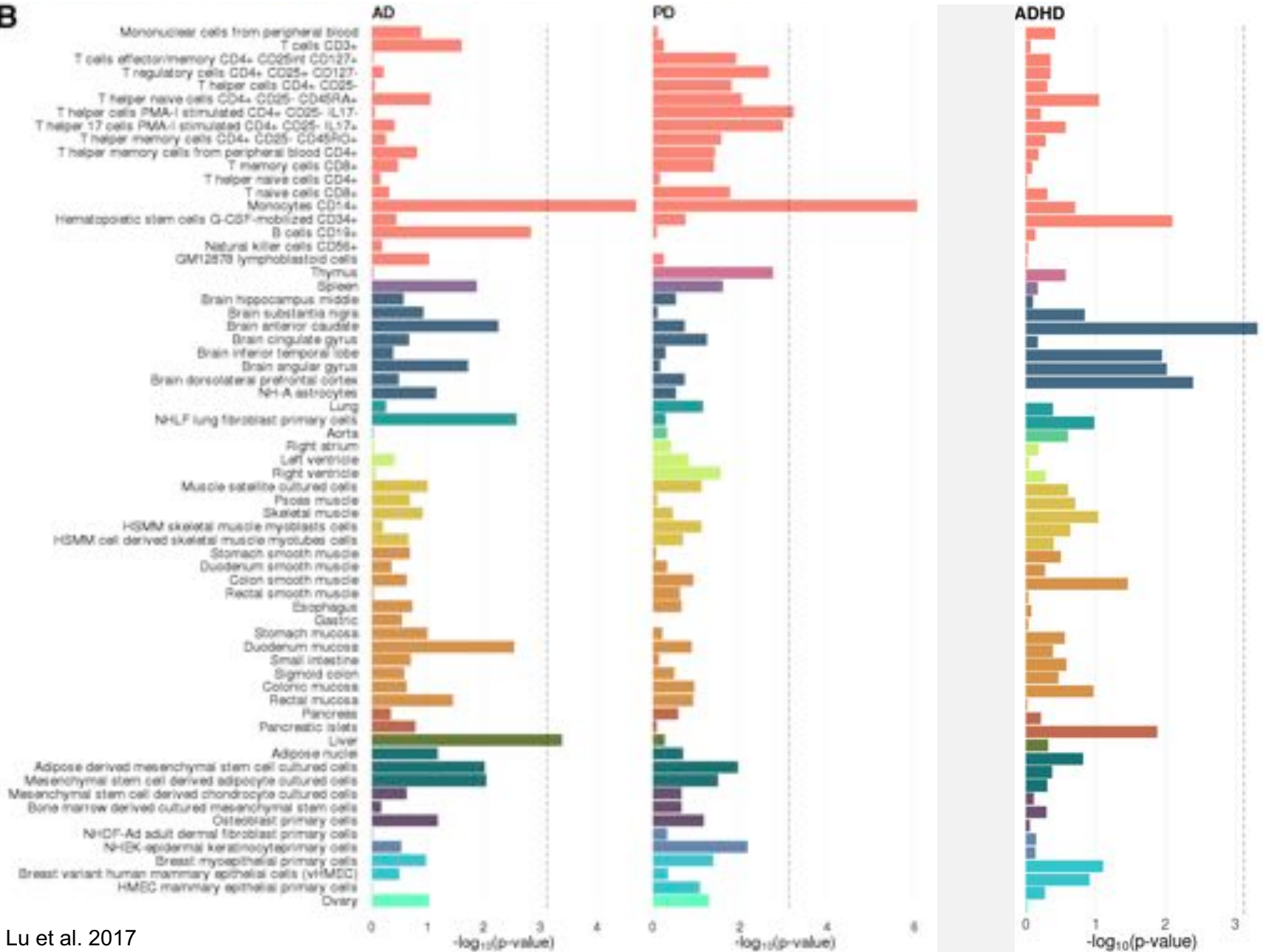
$$\text{Enrichment} = \frac{\% \text{ heritability explained}}{\% \text{ genome covered}}$$





# Partitioning heritability and genetic covariance

B





# Partitioning heritability and genetic covariance



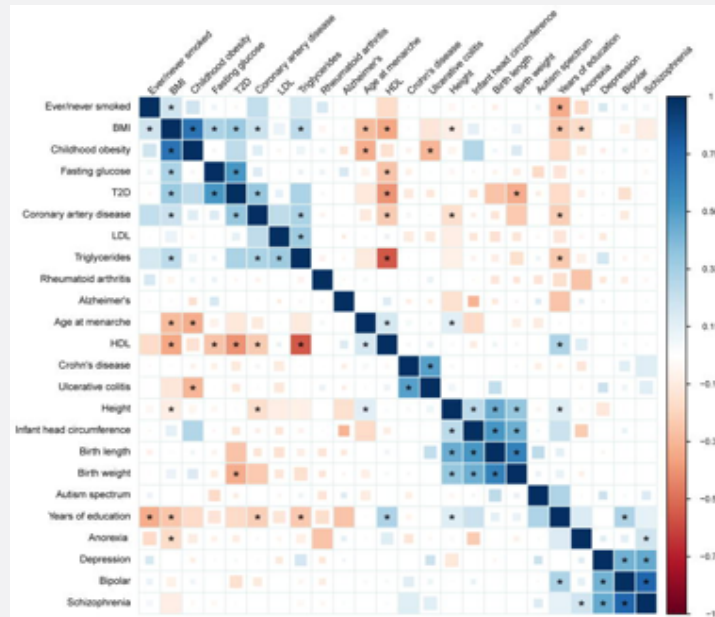
- **Genetic covariance** quantifies shared genetics among complex traits

$$y_1 = \sum_{i=1}^K X_i \beta_i + \epsilon$$

$$y_2 = \sum_{i=1}^K Z_i \gamma_i + \delta$$

genetic covariance

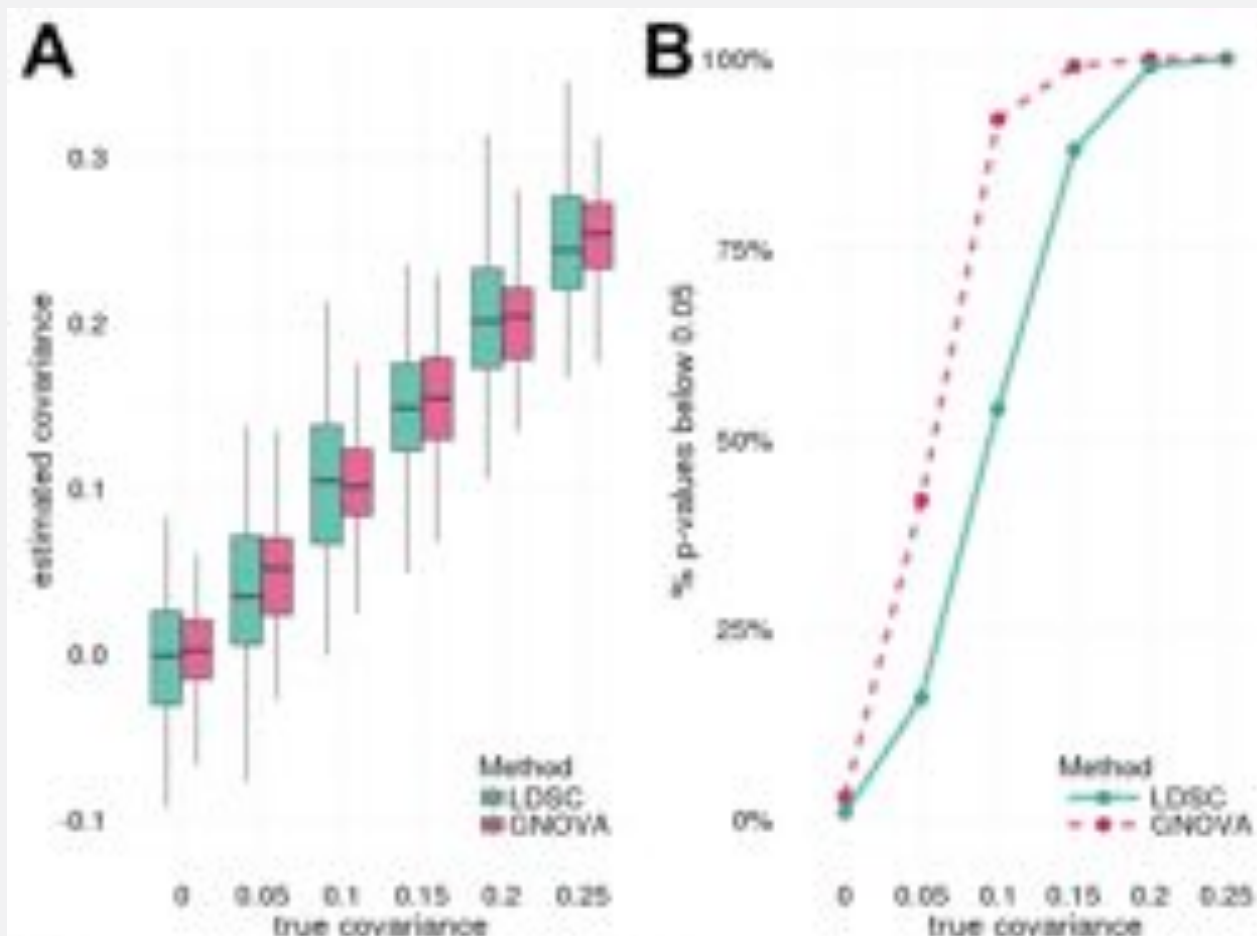
$$\mathbb{E}(\beta_i) = \mathbb{E}(\gamma_i) = 0 \text{ and } \mathbb{E}(\gamma_i \beta_i^T) = \frac{\rho_i}{m_i} I, \quad i = 1, \dots, K$$





# Partitioning heritability and genetic covariance

- **GNOVA**, a principled framework to perform annotation-stratified genetic covariance estimation



# Partitioning heritability and genetic covariance



We dissected the genetic covariance between late-onset Alzheimer's disease (LOAD) and amyotrophic lateral sclerosis (ALS)

**LOAD:** IGAP phase-I (17,008 cases, 37,154 controls)

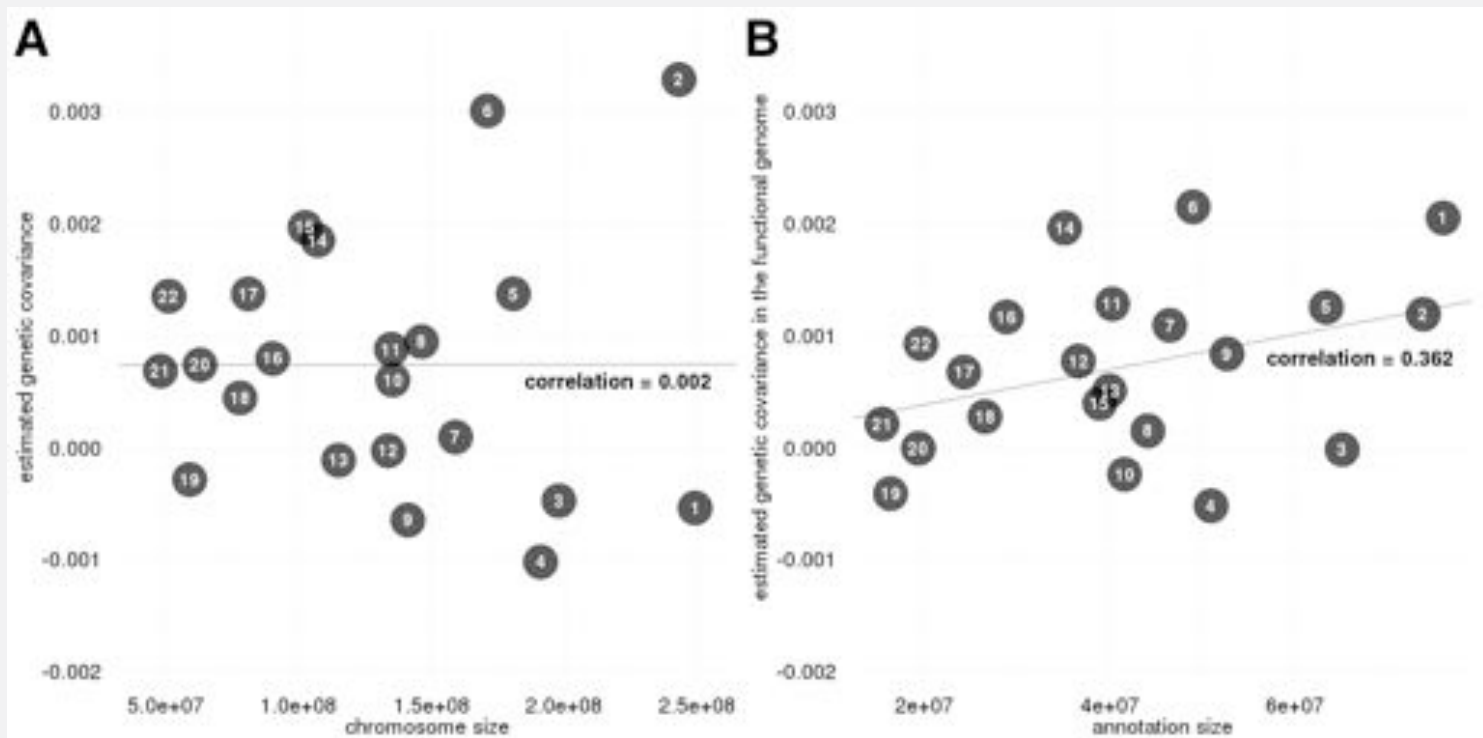
**ALS:** MinE project (12,577 cases, 23,475 controls)

Annotation	Category	Covariance	P-value
Non-stratified	GNOVA	0.016 (0.004)	<b><math>2.0 \times 10^{-4}</math></b>
	LDSC	0.012 (0.007)	0.075
GenoCanyon	functional	0.016 (0.004)	<b><math>8.2 \times 10^{-5}</math></b>
	non-functional	0.003 (0.004)	0.377
MAF	Q1	-0.001 (0.003)	0.842
	Q2	0.003 (0.004)	0.361
	Q3	0.004 (0.004)	0.327
	Q4	0.008 (0.003)	<b>0.005</b>



# Partitioning heritability and genetic covariance

Genetic covariance between LOAD and ALS is proportional to the size of the functional genome on each chromosome. This suggests a **polygenic** genetic covariance structure.





## 1. Background

## 2. Introduction to (narrow-sense) functional annotations

- Functional annotation in protein-coding genes
- Functional annotation in non-coding regions
- Other useful tools

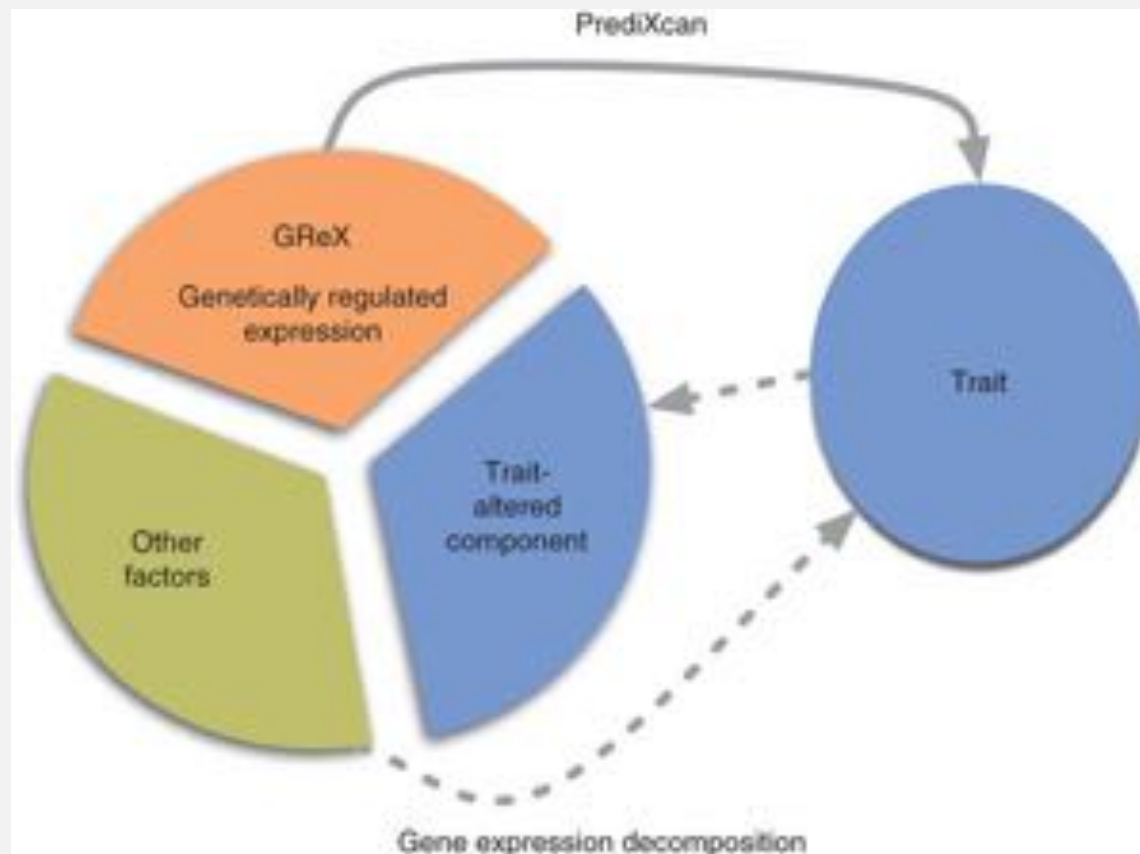
## 3. Applications of functional annotations

- Functional SNP fine-mapping
- Partitioning heritability and genetic covariance
- **Gene-level analysis**
- Effect size estimation and risk prediction

# Gene-level analysis

## PrediXcan and TWAS

- Impute gene expression using genotype information
- Perform association test using imputed expression

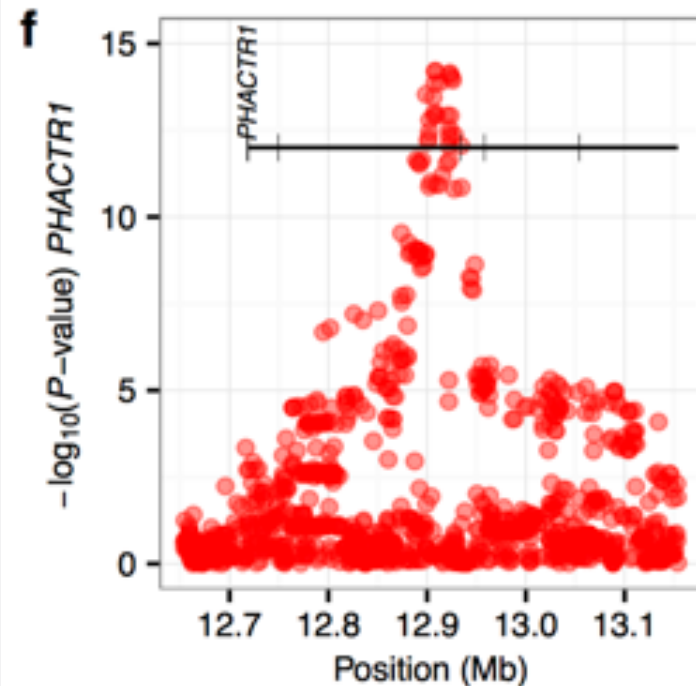
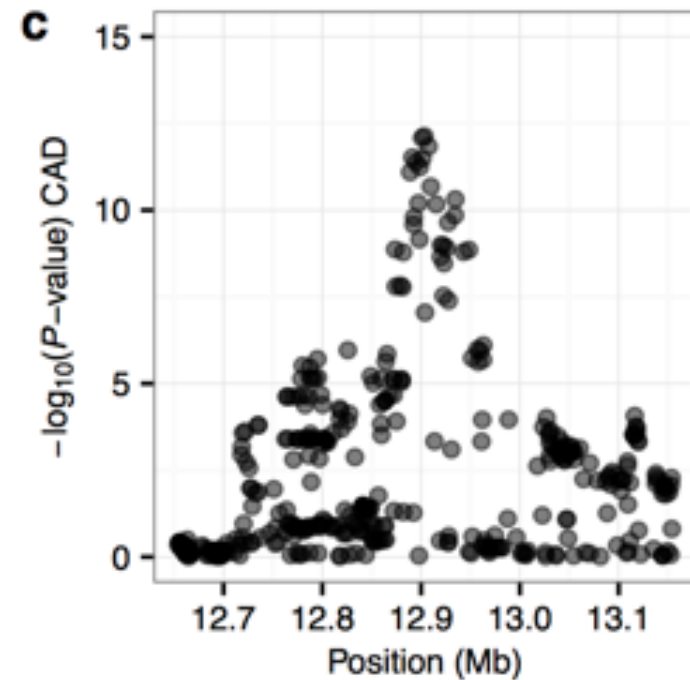
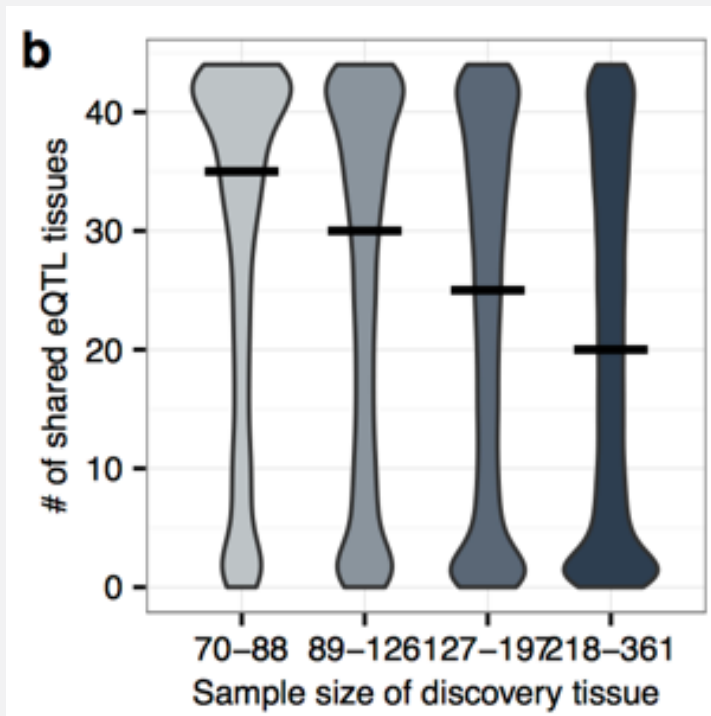


# Gene-level analysis

## PrediXcan and TWAS

Idea is similar to colocalization

Often identify signals in irrelevant tissues

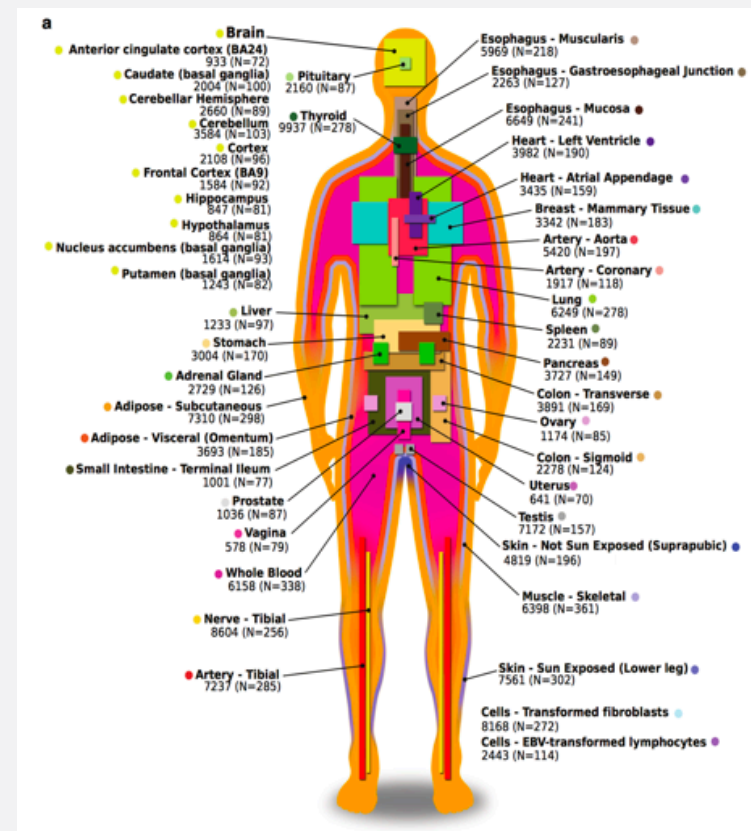
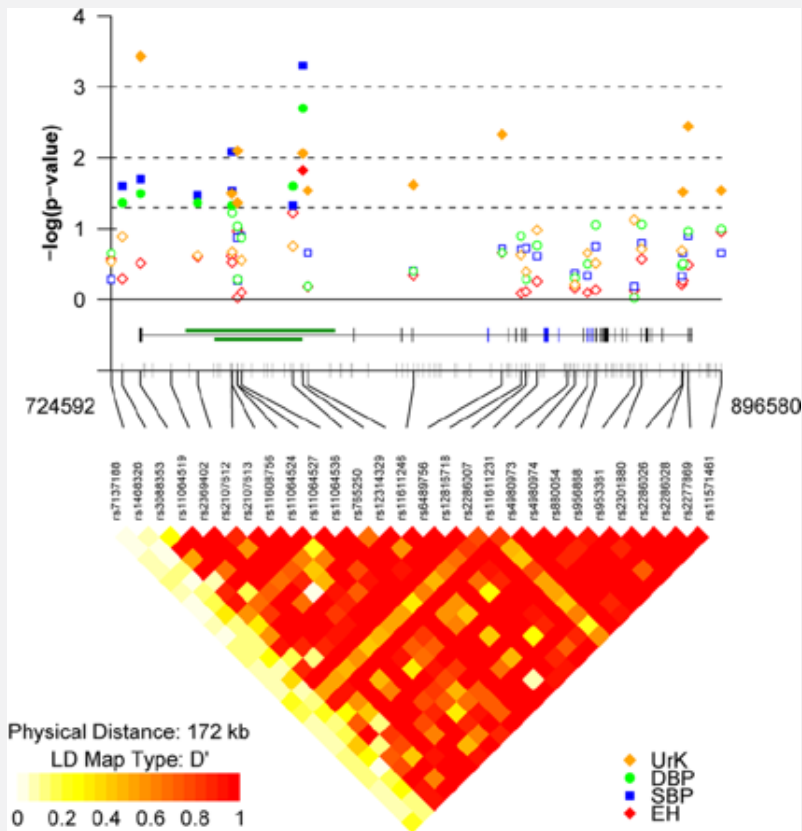


# Gene-level analysis

## PrediXcan and TWAS

Need a metric to summarize information across all tissues

This is very similar to “burden test for common SNPs”







## 1. Background

## 2. Introduction to (narrow-sense) functional annotations

- Functional annotation in protein-coding genes
- Functional annotation in non-coding regions
- Other useful tools

## 3. Applications of functional annotations

- Functional SNP fine-mapping
- Partitioning heritability and genetic covariance
- Gene-level analysis
- Effect size estimation and risk prediction

# Effect size estimation and risk prediction

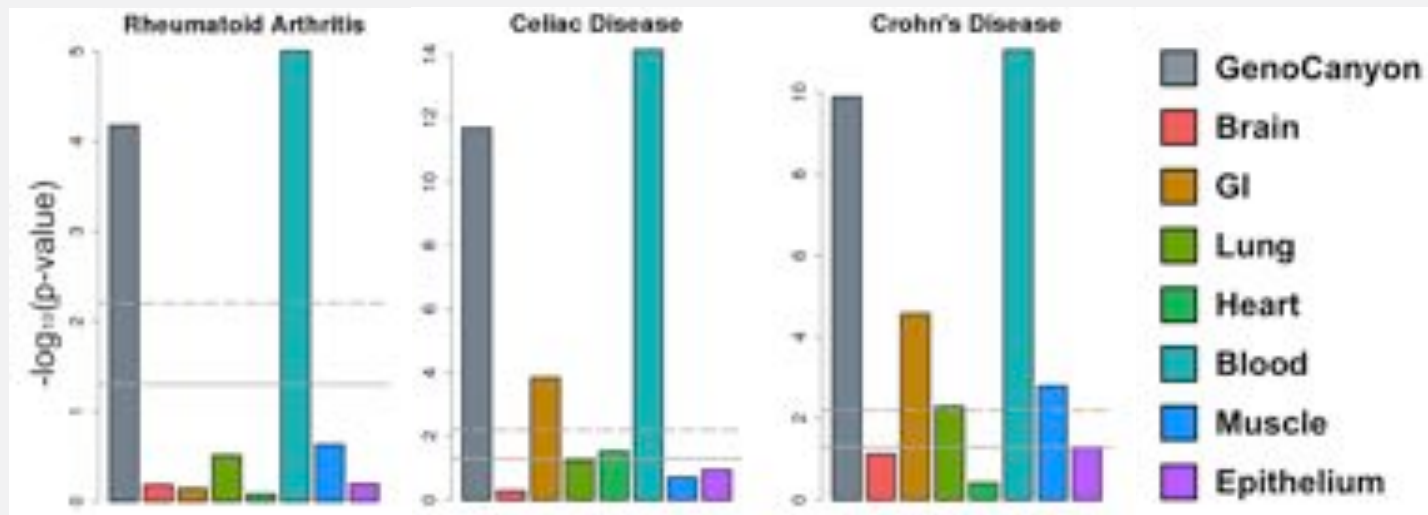
## AnnoPred

Remember we can connect disease with tissues?

$$Enrichment = \frac{\% \text{ heritability explained}}{\% \text{ genome covered}}$$

Such connections can help estimate SNPs' effect sizes

$$\mathbb{E}(\beta \mid \hat{\beta}, \hat{D})$$

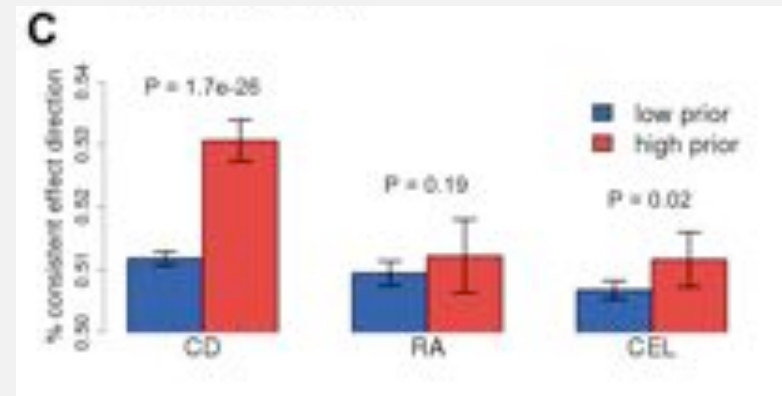
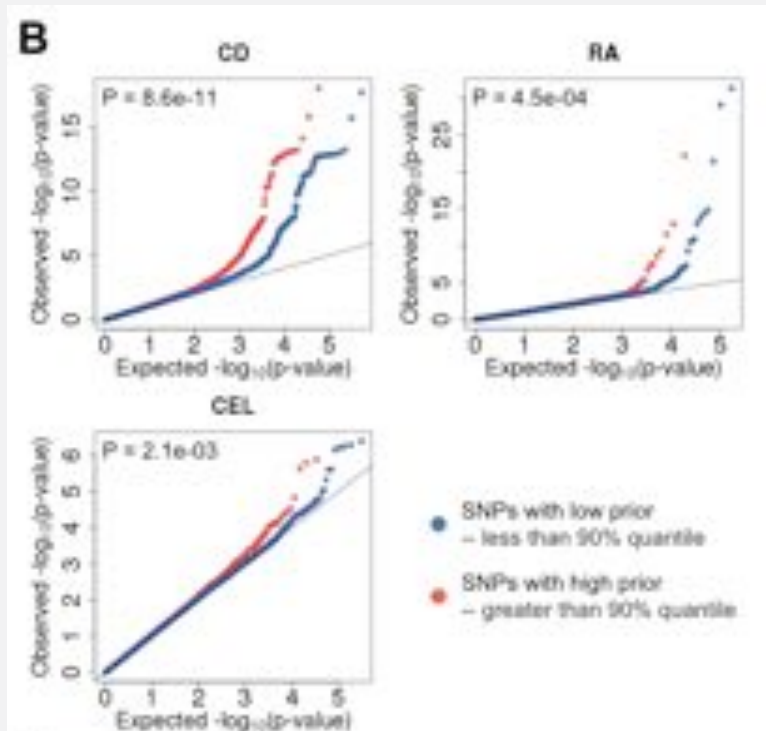




# Effect size estimation and risk prediction

## AnnoPred

SNPs with high prior show stronger associations and more consistent effect directions in validation cohorts





# Effect size estimation and risk prediction

## AnnoPred

We achieved higher risk prediction accuracy across five complex diseases

$$\hat{y} = X\mathbb{E}(\beta|\hat{\beta}, \hat{D})$$

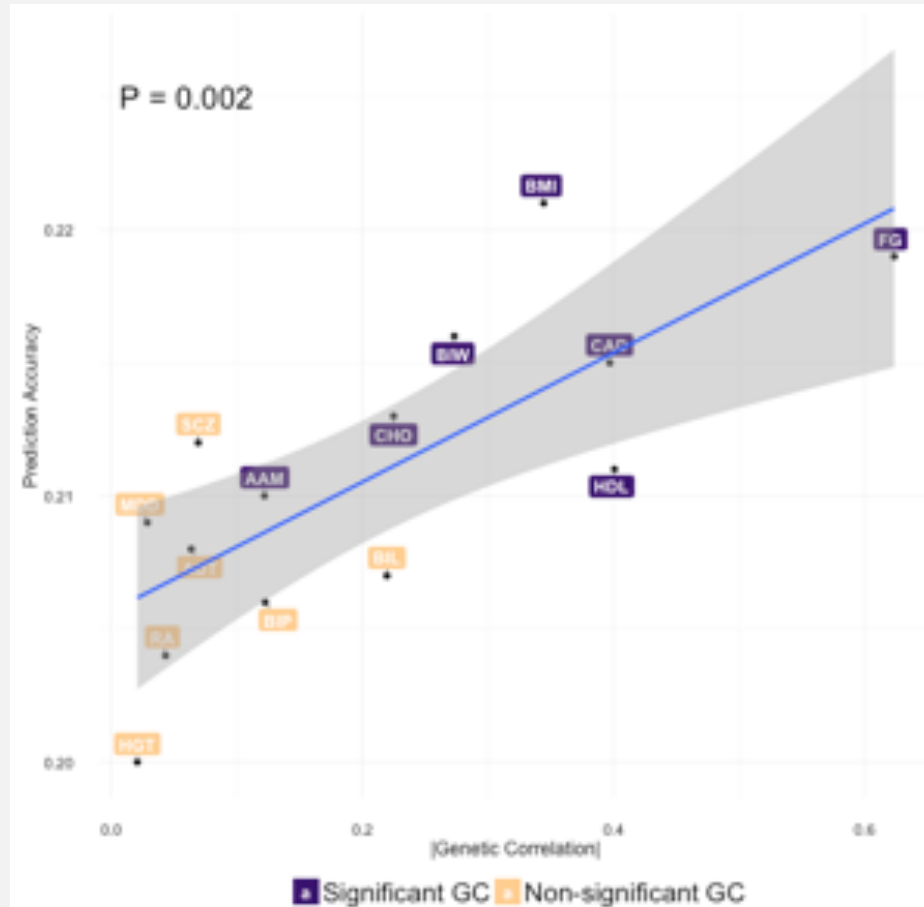
**Table 2. CORs of different methods.** The highest CORs are highlighted in boldface.

Disease/Trait	PRS <sub>sig</sub>	PRS <sub>all</sub>	PRS <sub>P+T</sub>	LDpred	AnnoPred
Crohn's Disease	0.27	0.229	0.32	0.325	<b>0.343</b>
Breast Cancer	0.084	0.055	0.12	0.122	<b>0.137</b>
Rheumatoid Arthritis	0.204	0.114	0.248	0.282	<b>0.287</b>
Type-II Diabetes	0.165	0.156	0.204	0.202	<b>0.22</b>
Celiac Disease	0.11	0.136	0.18	0.197	<b>0.213</b>



## PleioPred

We have further extended the model to incorporate multiple GWAS for genetically correlated diseases





## Annotation = External Information

- Conservation
- Epigenetic data
- Transcriptomic data
- Functional prediction scores (supervised / unsupervised)
- Quantitative trait loci (eQTL, sQTL, pQTL)
- Allele frequency (ExAC, gnomad)
- LD (1000 Genomes)
- Additional GWAS
- Chromatin interaction



## Applications

- Fine-mapping (GenoWAP)
- Partition heritability and infer relevant tissue (LDSC+GenoSkyline)
- Partition genetic covariance (GNOVA)
- Gene-level association test (new method coming soon)
- Risk prediction (AnnoPred, PleioPred)

**THANK YOU**

