



Centre for Molecular Medicine  
and Therapeutics



## xQTL Analysis

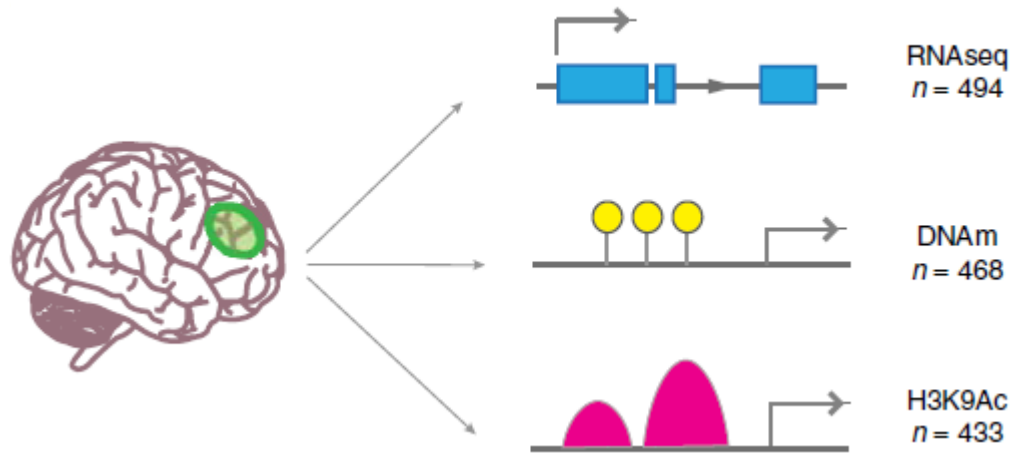
**Bernard Ng**



**Department of Statistics, UBC**

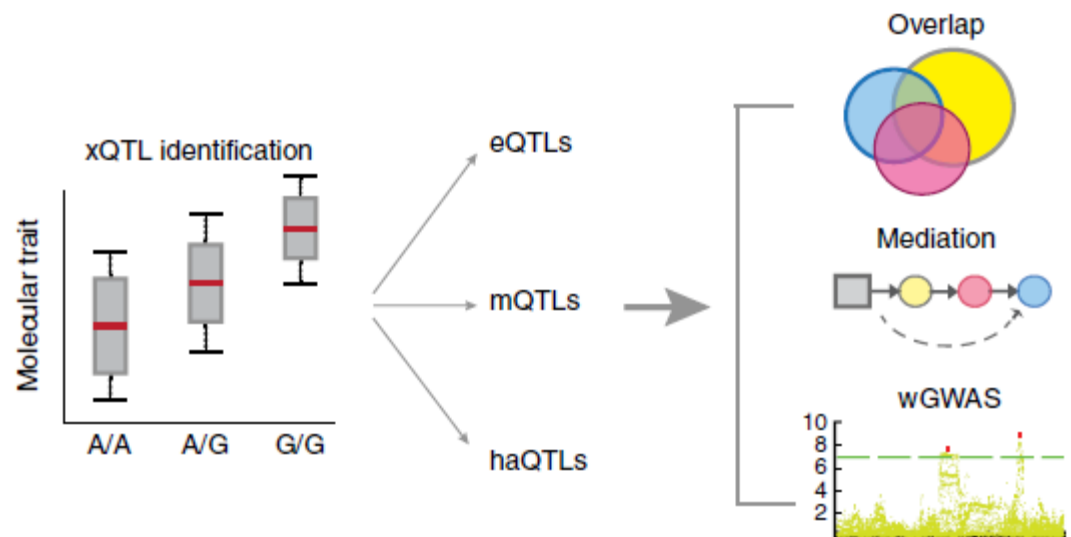
[www.cmmmt.ubc.ca](http://www.cmmmt.ubc.ca)

# Outline

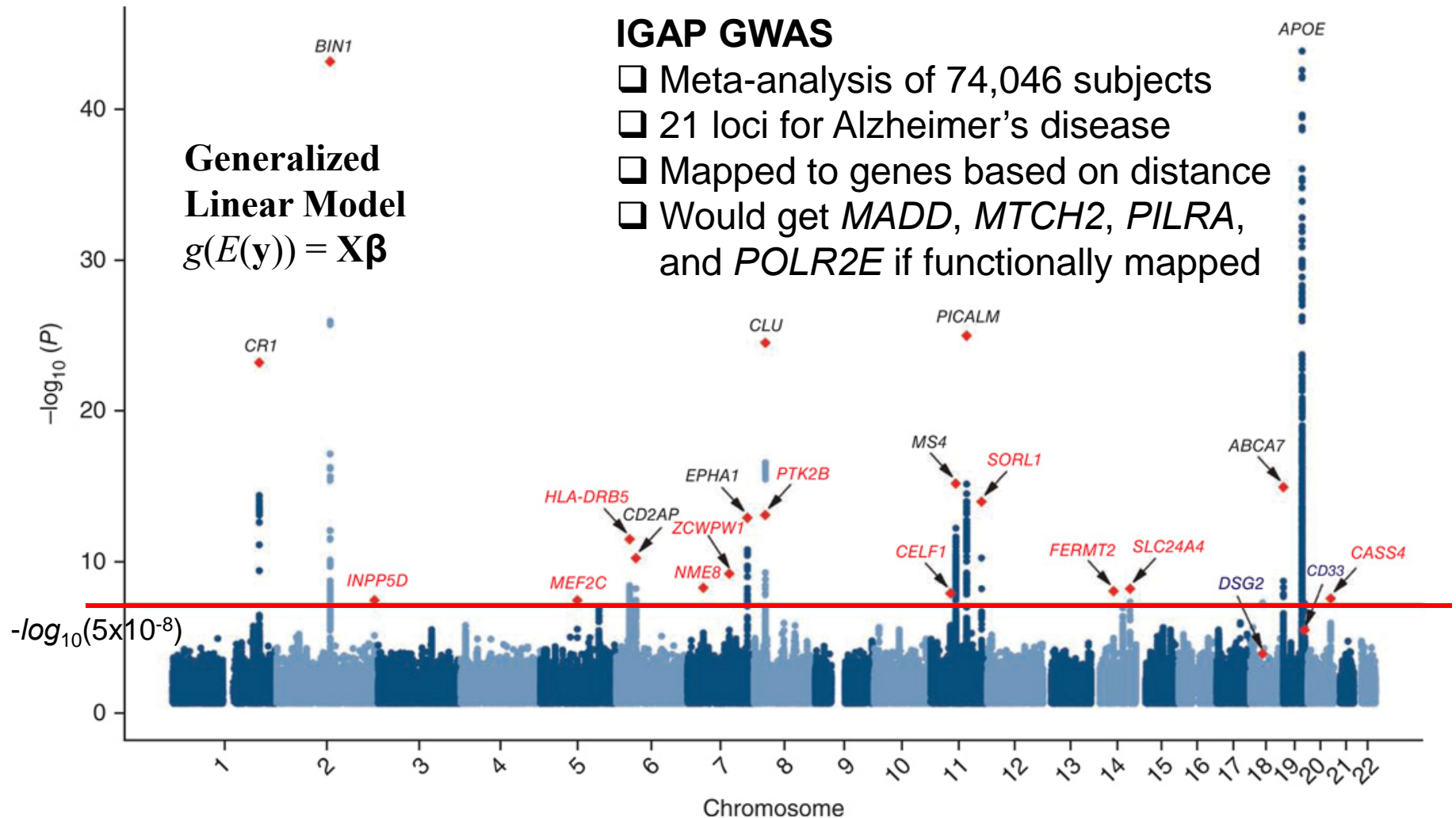


- Recap of GWAS
- xQTL Analysis
- Challenges

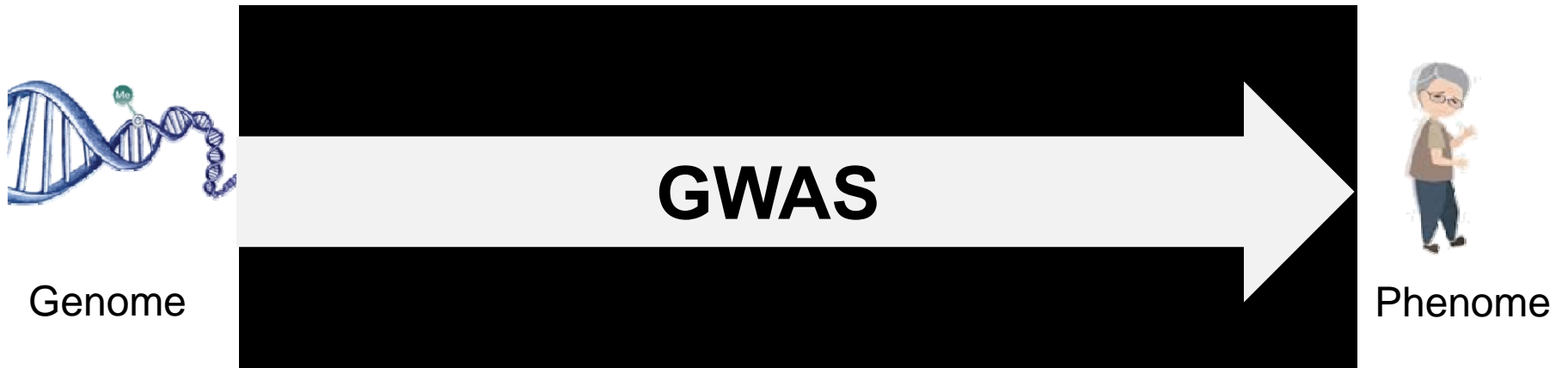
- Replication
- xQTL Sharing
- Mediation Analysis
- Cell Specific Analysis
- Weighted GWAS



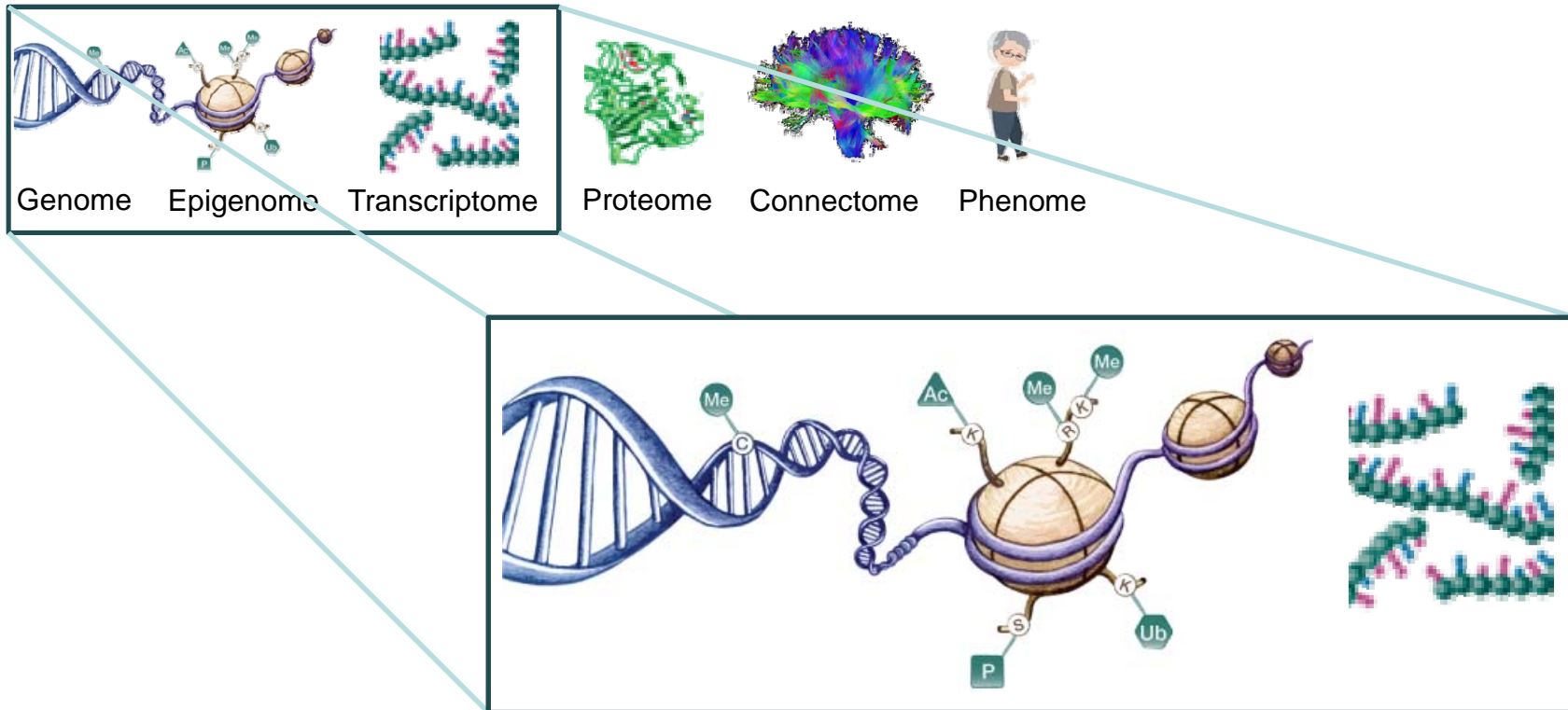
# Motivation



# xQTL Analysis



# xQTL Analysis



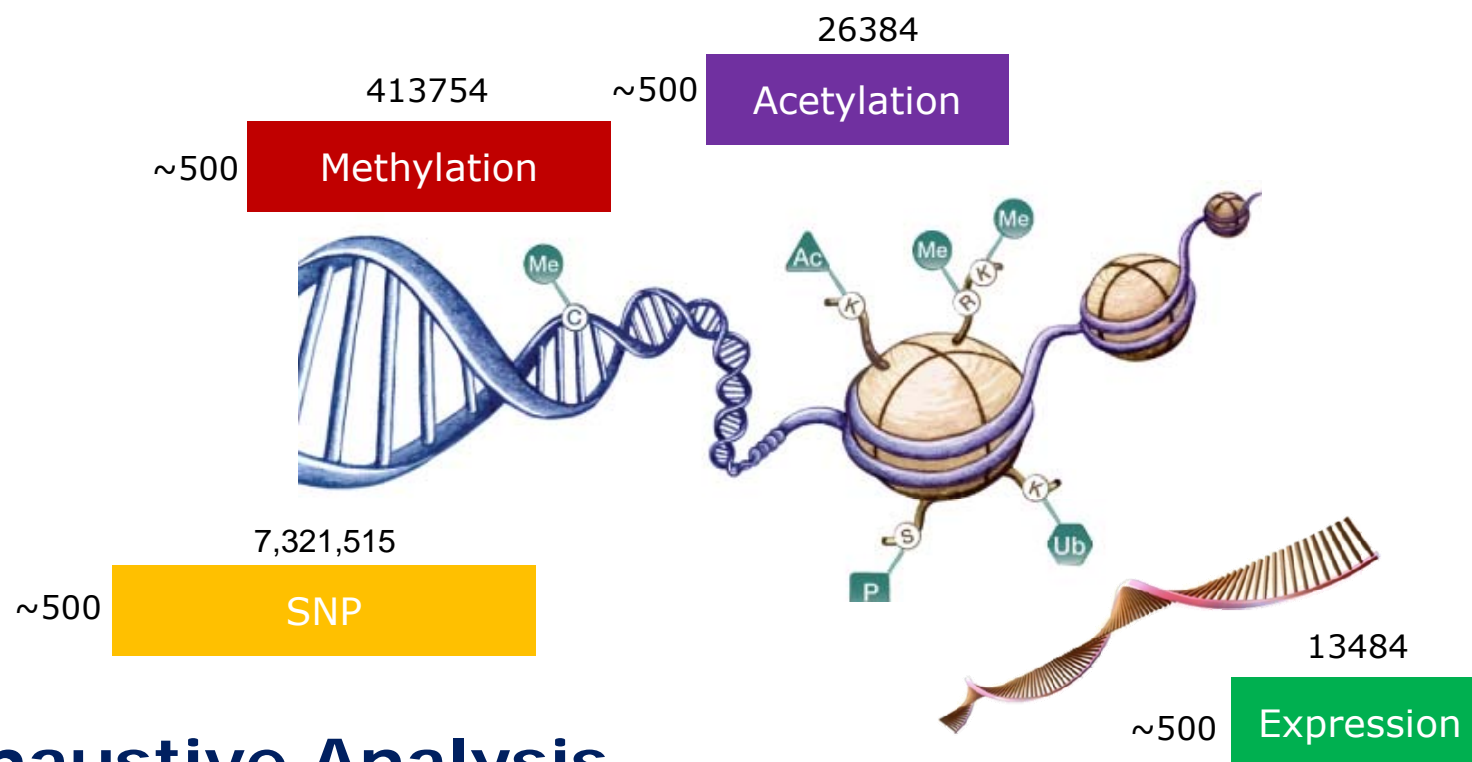
Linear Regression Model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Spearman's correlation

$$\text{corr}(\text{rank}(\mathbf{y}^c), \text{rank}(\mathbf{X}_i))$$

# Challenges: High Dimensionality

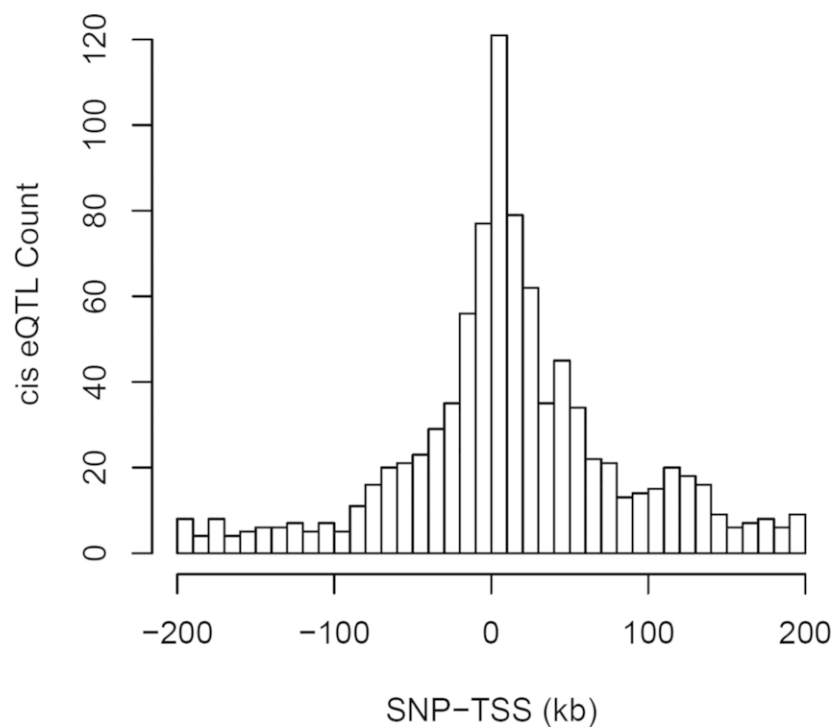


## Exhaustive Analysis

- eQTL: 98,723,308,260
- mQTL: 3,029,306,117,310
- haQTL: 193,170,851,760

Trillion of tests =>  
no statistical power

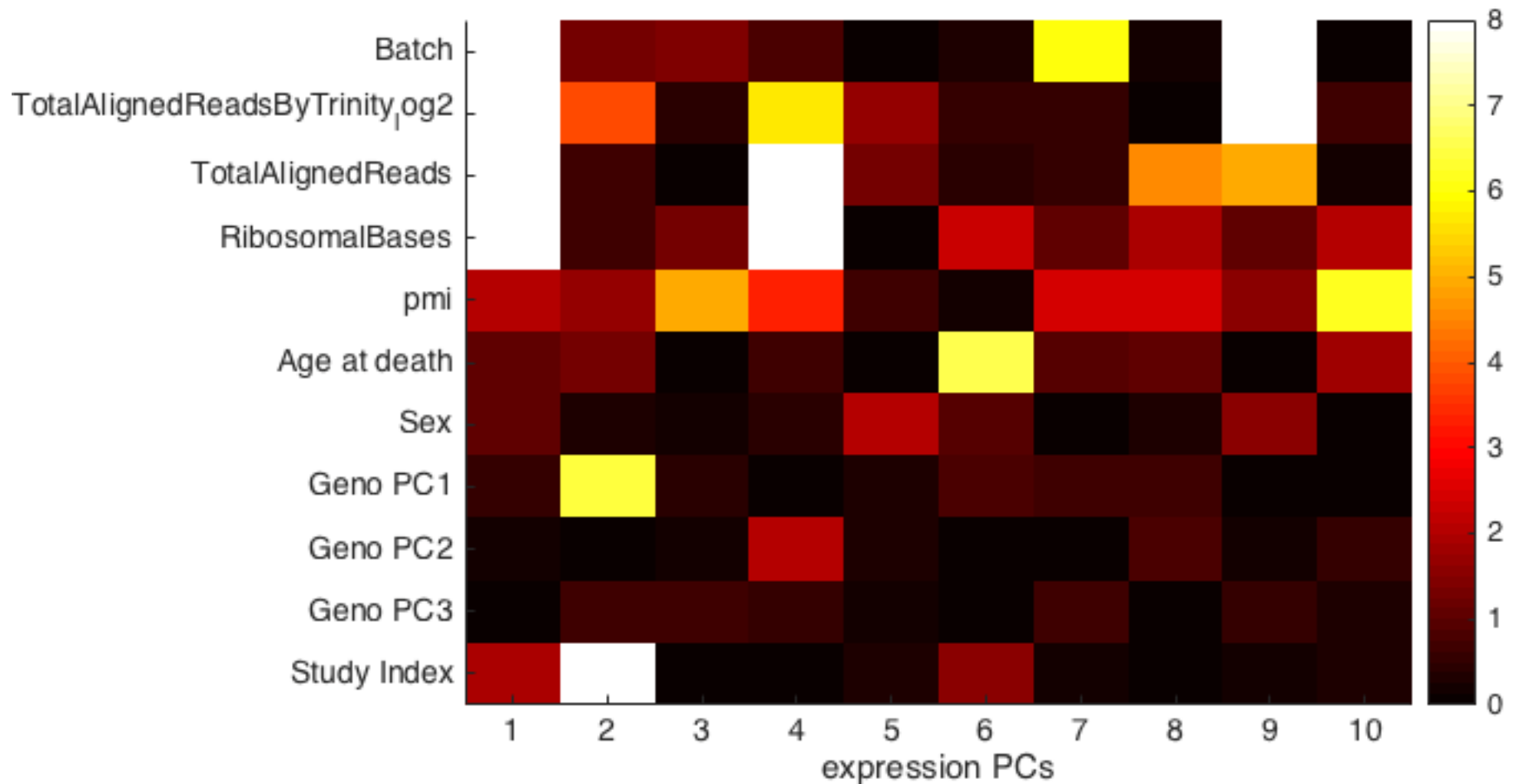
# Challenges: High Dimensionality



Windowing to  
reduce the  
number of tests

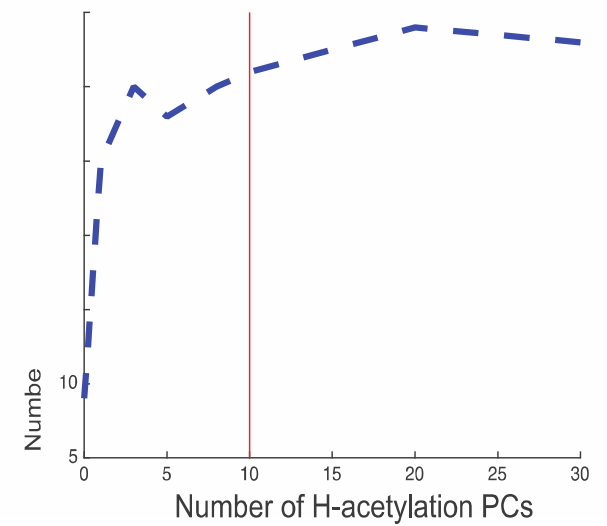
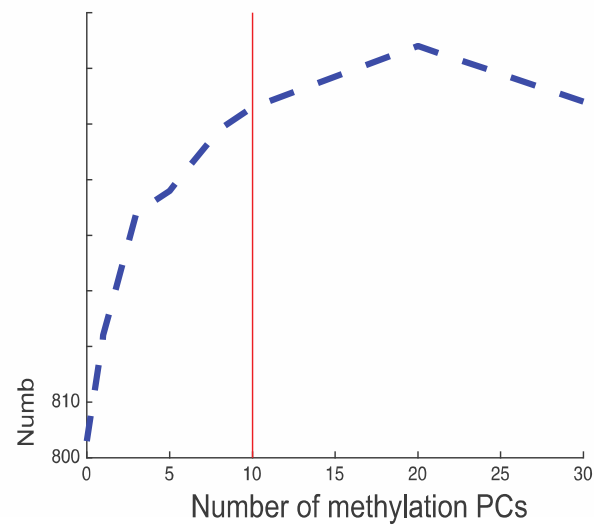
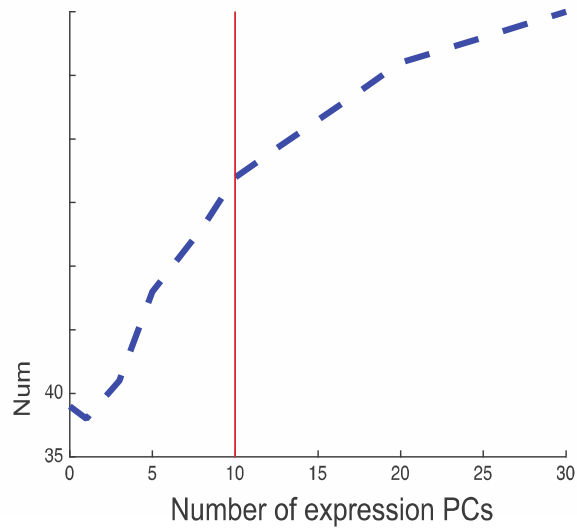
	#SNP-feature pairs	#Features	#SNPs
eQTL (1Mb)	60,456,556	12,979	6,442,864
mQTL (5Kb)	9,939,236	412,152	2,358,873
haQTL (1Mb)	125,100,450	25,720	6,756,597

# Challenges: Hidden Confounds



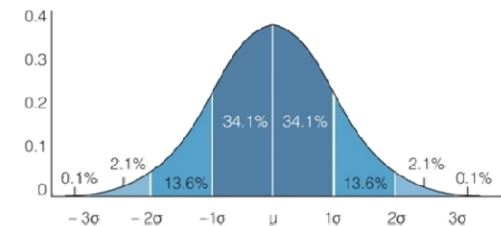


# Challenges: Hidden Confounds



# Challenges: Multiple Testing

- $P(\text{not rejecting 1 hypothesis}) = 1 - \alpha$
- $P(\text{not rejecting all } n \text{ hypotheses}) = (1 - \alpha)^n$
- $\alpha_{\text{FWER}} = 1 - (1 - \alpha)^n$ 
  - $\alpha = 0.05, n = 10: \alpha_{\text{FWER}} = 0.4013$
  - $\alpha = 0.05, n = 10^2: \alpha_{\text{FWER}} \approx 1$
- So if e.g. run 100 experiments, then  $\alpha_{\text{FWER}} \cdot 100$  of them would have  $\geq 1$  hypothesis falsely rejected.
- Intuition is that the more we sample the variable space, the more “likely” we will get some “extreme” samples.



# Bonferroni Correction

## Procedures

- Recall  $\alpha_{\text{FWER}} = 1 - (1 - \alpha)^n$
- Set  $\alpha = 1 - (1 - \alpha_{\text{FWER}})^{1/n} \approx 1 - (1 - \alpha_{\text{FWER}}/n) = \alpha_{\text{FWER}}/n$

## Examples

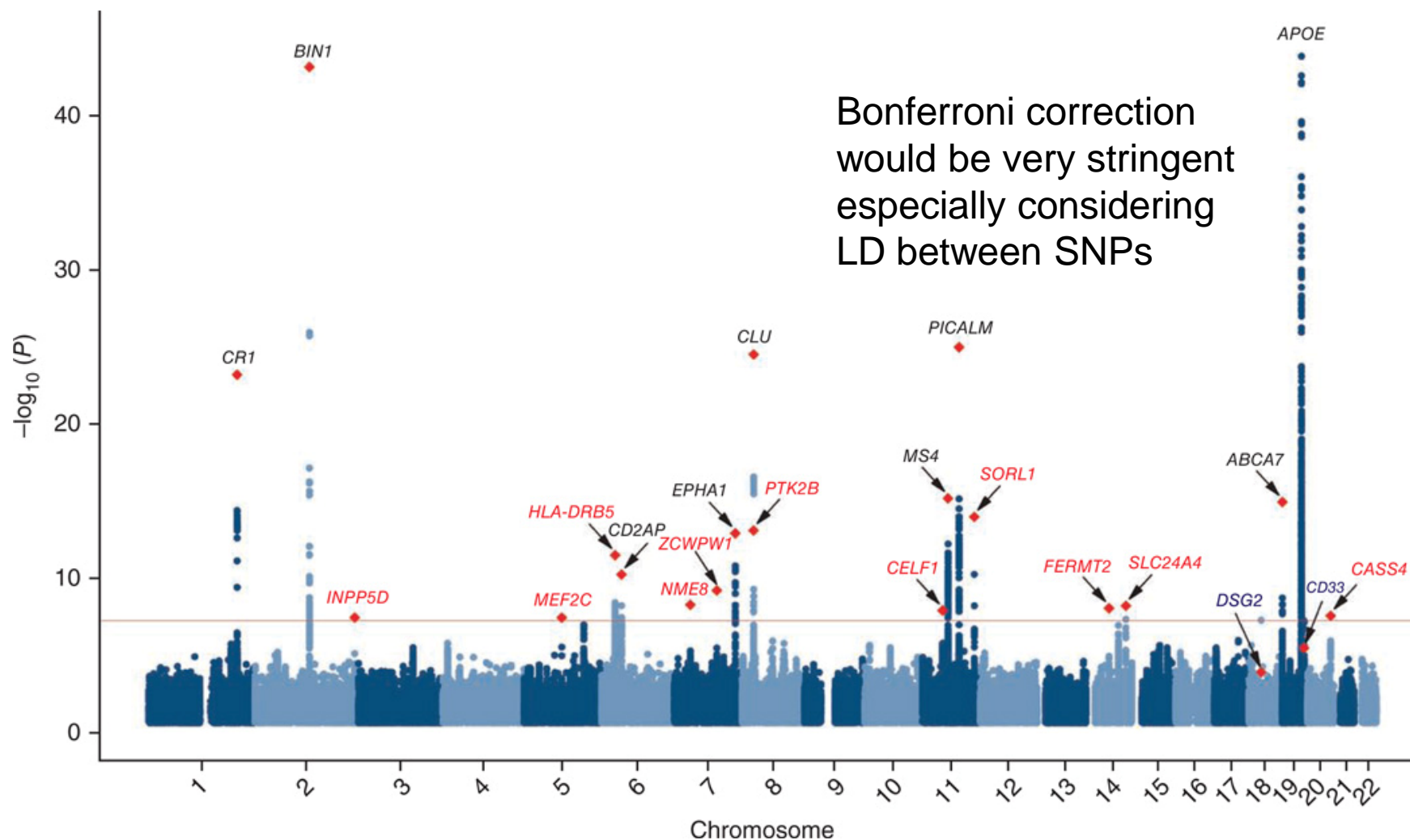
- $\alpha_{\text{FWER}} = 0.05$  and  $n = 10$ , needs  $\alpha = 0.05/10 = 0.005$
- $\alpha_{\text{FWER}} = 0.05$  and  $n = 10^6$ , needs  $\alpha = 0.05/10^6 = 5 \times 10^{-8}$

## Properties

- Controls FWER =  $P(V \geq 1)$  in *strong* sense.
- Can handle correlated hypotheses.
- Very stringent

		Predicted		
		True	False	
Ground Truth	True	U	V	$n_0$
	False	T	S	$n - n_0$
		$n - R$	R	$n$

# Challenges: Multiple Testing



# False Discovery Rate

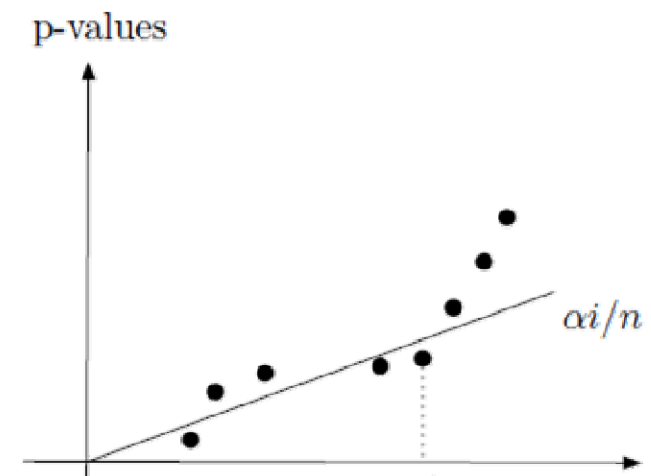
## Idea

- Benjamini & Hochberg, 1995
- Recall  $\text{FWER} = P(V \geq 1)$
- $\text{Fdp} = V/\max(R, 1)$
- But  $V$  unobserved, so:  $\text{FDR} = E(\text{Fdp})$

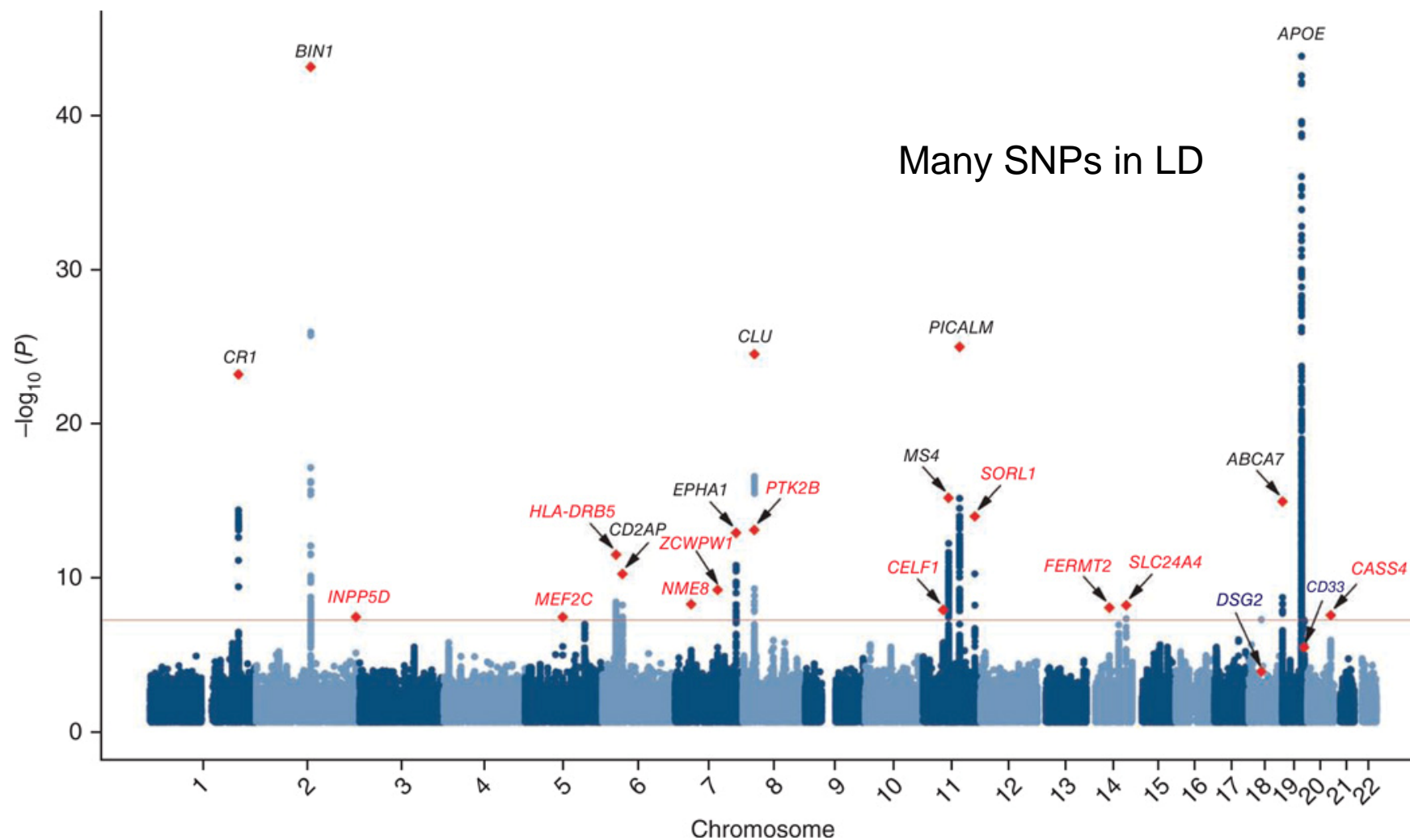
## Procedures

- Sort  $p$  in ascending order.
- Find  $i_0 = \max i \text{ s.t. } p(i) \leq i \cdot q/n$

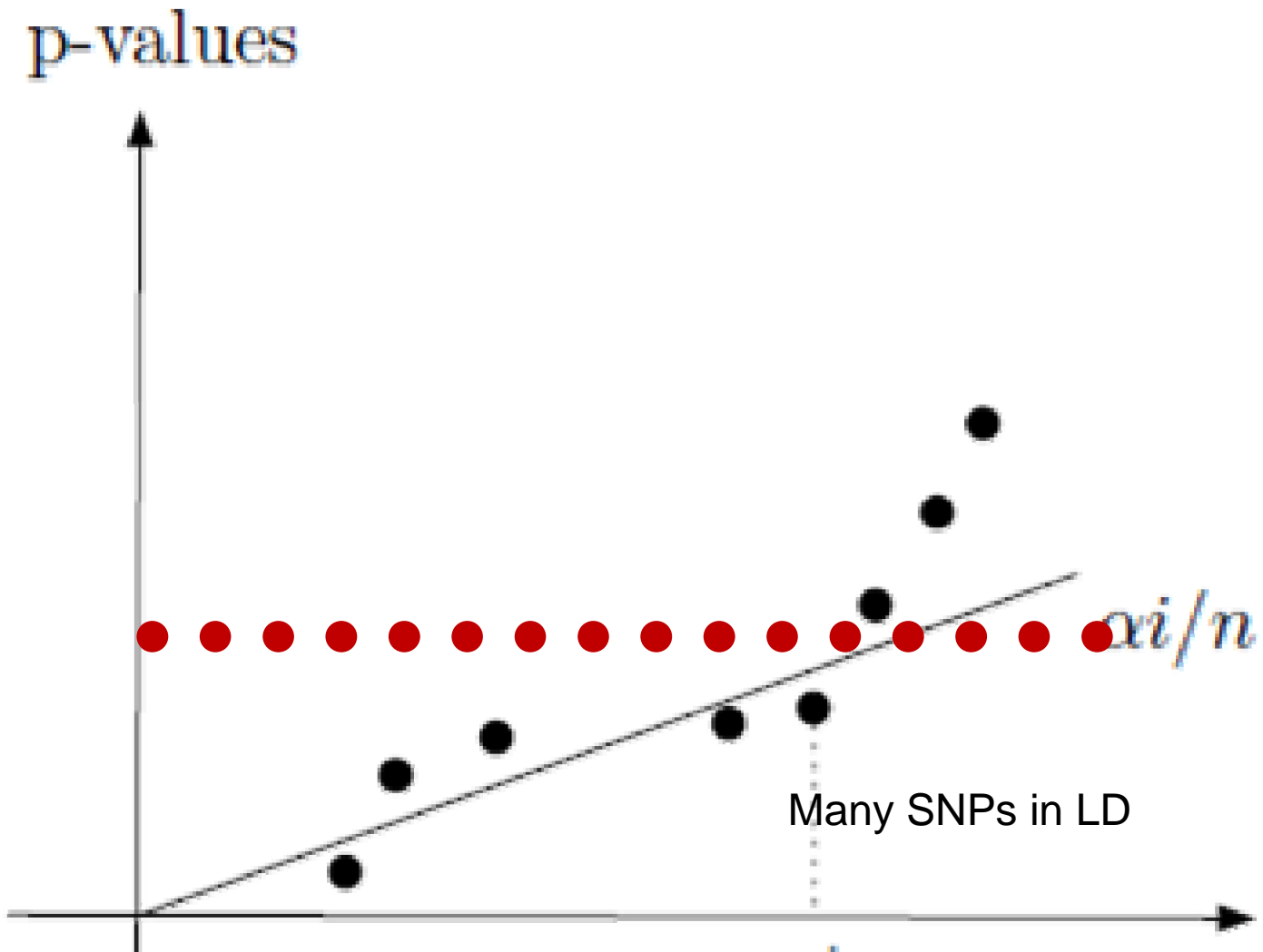
		Predicted		
		True	False	
Ground Truth	True	U	V	$n_0$
	False	T	S	$n - n_0$
		$n - R$	R	n



# Challenges: Multiple Testing



# False Discovery Rate



# False Discovery Rate

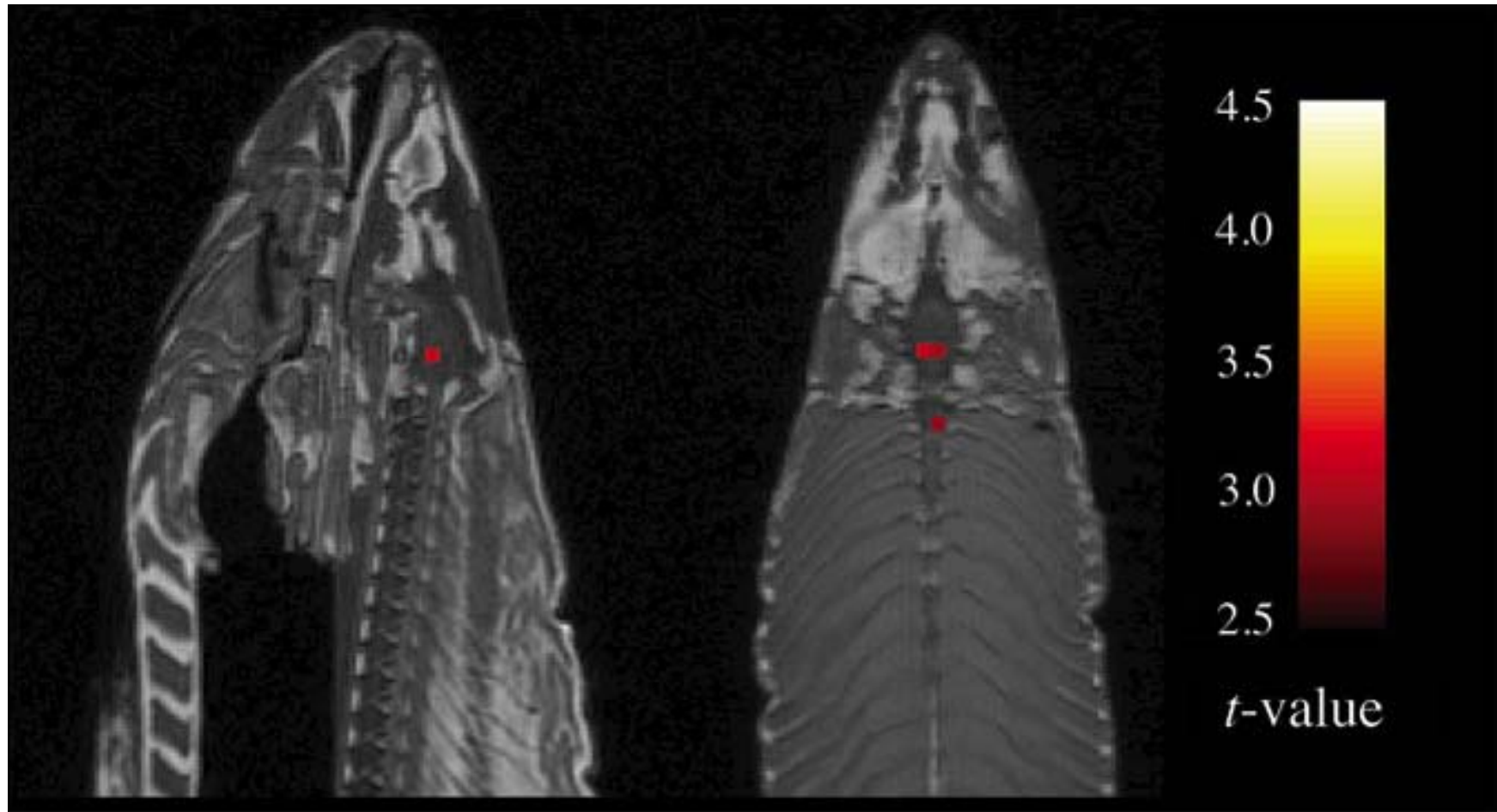
## Properties

- If hypotheses are independent, then  $FDR < q$  for *all* configurations of hypotheses.
- If hypotheses are correlated,  
 $FDR < q \cdot (\log(n) + 0.577)$   
 $\Rightarrow p(i) < i \cdot q/n / (\log(n) + 0.577)$   
BUT  $i = 1$ ,  $p(i) < q/n / (\log(n) + 0.577) < q/n$

		Predicted		
		True	False	
Ground Truth	True	U	V	$n_0$
	False	T	S	$n - n_0$
		$n - R$	R	n



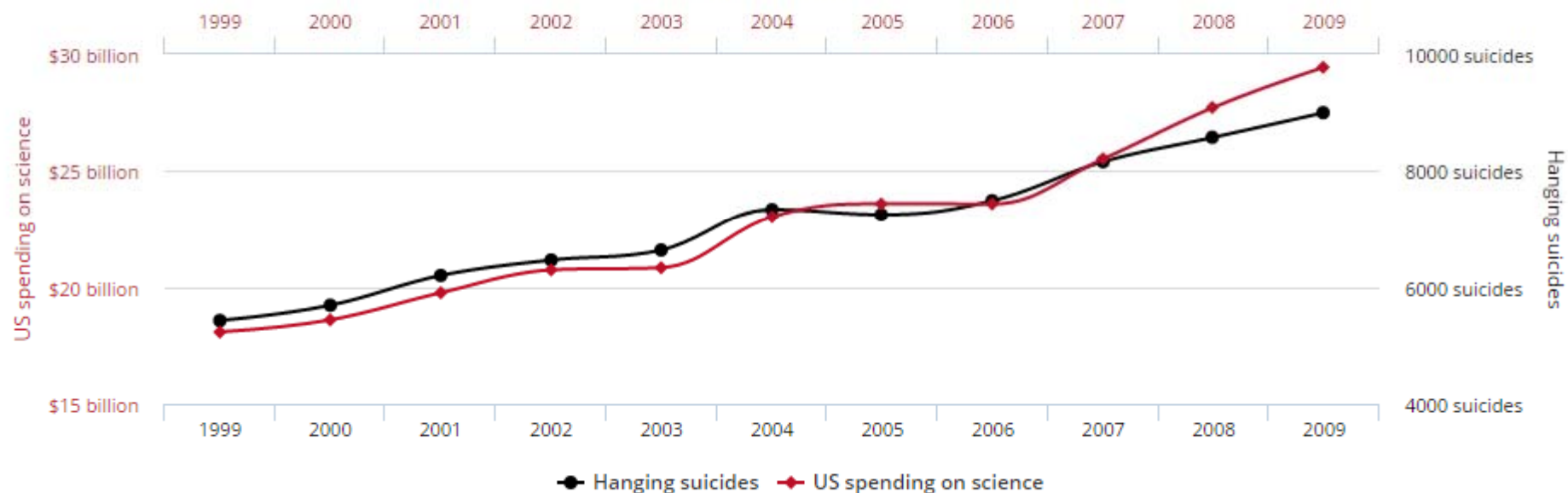
## Why Important?



# Why Important?

US spending on science, space, and technology  
correlates with  
Suicides by hanging, strangulation and suffocation

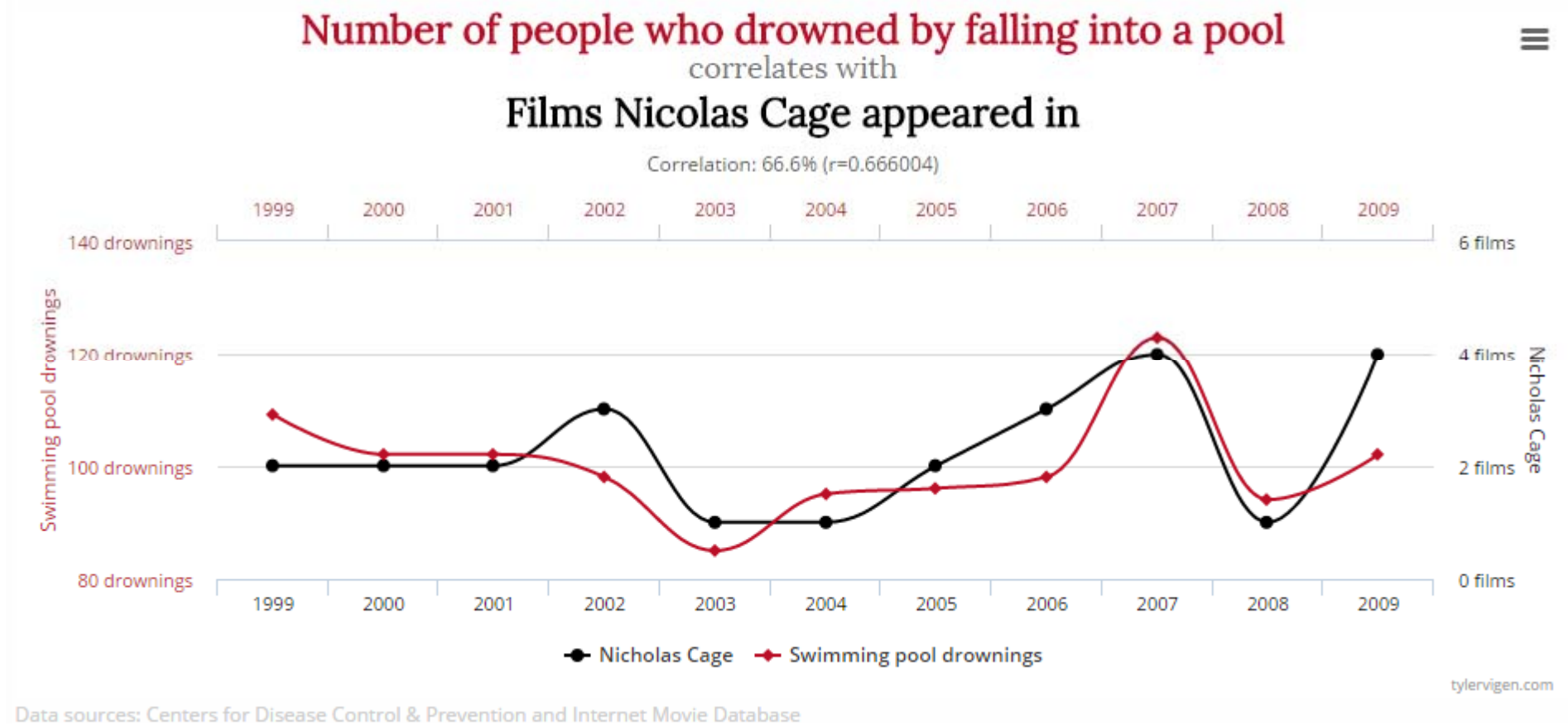
Correlation: 99.79% ( $r=0.99789126$ )



tylervigen.com

Data sources: U.S. Office of Management and Budget and Centers for Disease Control & Prevention

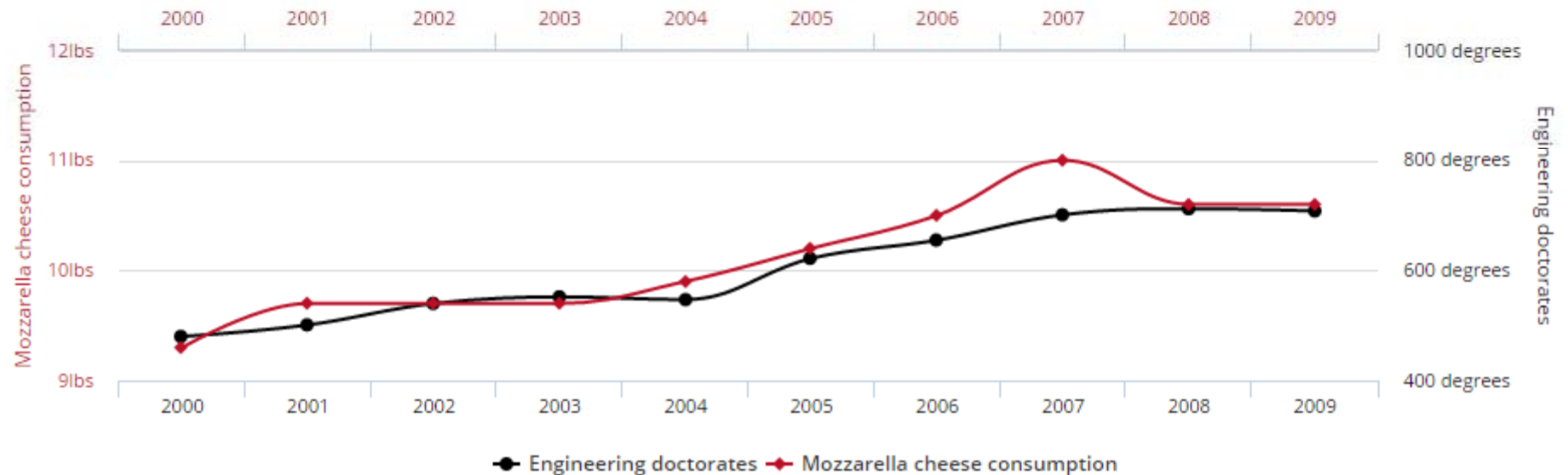
# Why Important?



# Why Important?

Per capita consumption of mozzarella cheese  
correlates with  
Civil engineering doctorates awarded

Correlation: 95.86% ( $r=0.958648$ )



Data sources: U.S. Department of Agriculture and National Science Foundation

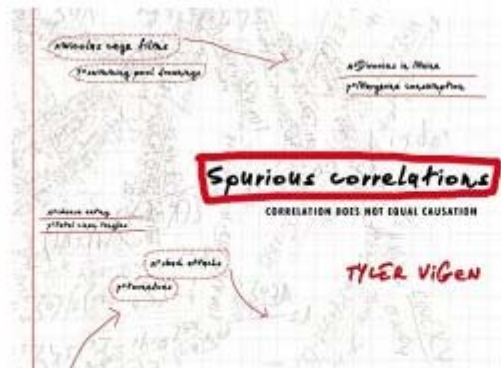
tylervigen.com

# Why Important?

tylervigen.com

about | twitter | email | subscri

## Spurious correlations



Now a ridiculous book!

- Spurious charts
- Fascinating factoids
- Commentary in the footnotes

Amazon | Barnes & Noble | Indie Bound

<http://tylervigen.com/spurious-correlations>

# Why Important?

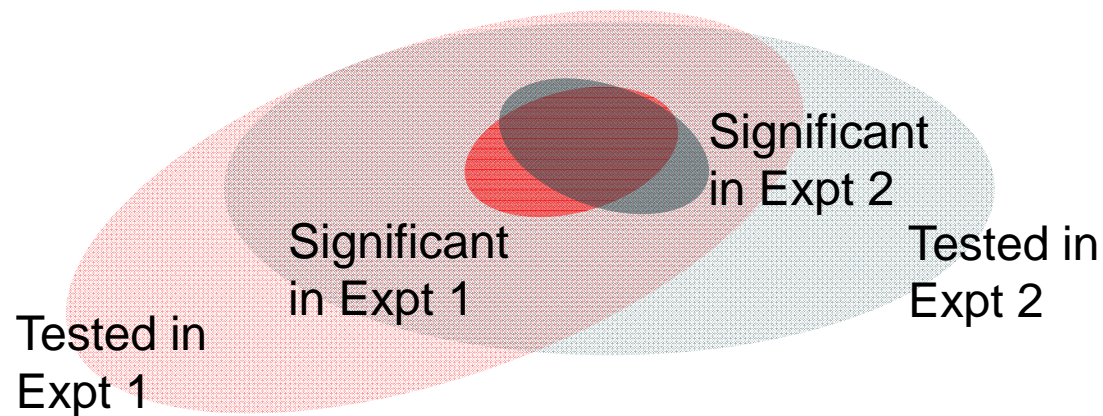
Problems with scientific research

## How science goes wrong

Scientific research has changed the world. Now it needs to change itself

Last year researchers at one biotech firm, Amgen, found they could reproduce just **six of 53** “landmark” studies in cancer research. Earlier, a group at Bayer, a drug company, managed to repeat just a **quarter of 67** similarly important papers. A leading computer scientist frets that **three-quarters** of papers in his subfield are bunk. In 2000-10 roughly **80,000 patients** took part in clinical trials based on research that was **later retracted** because of mistakes or improprieties.

# Replication

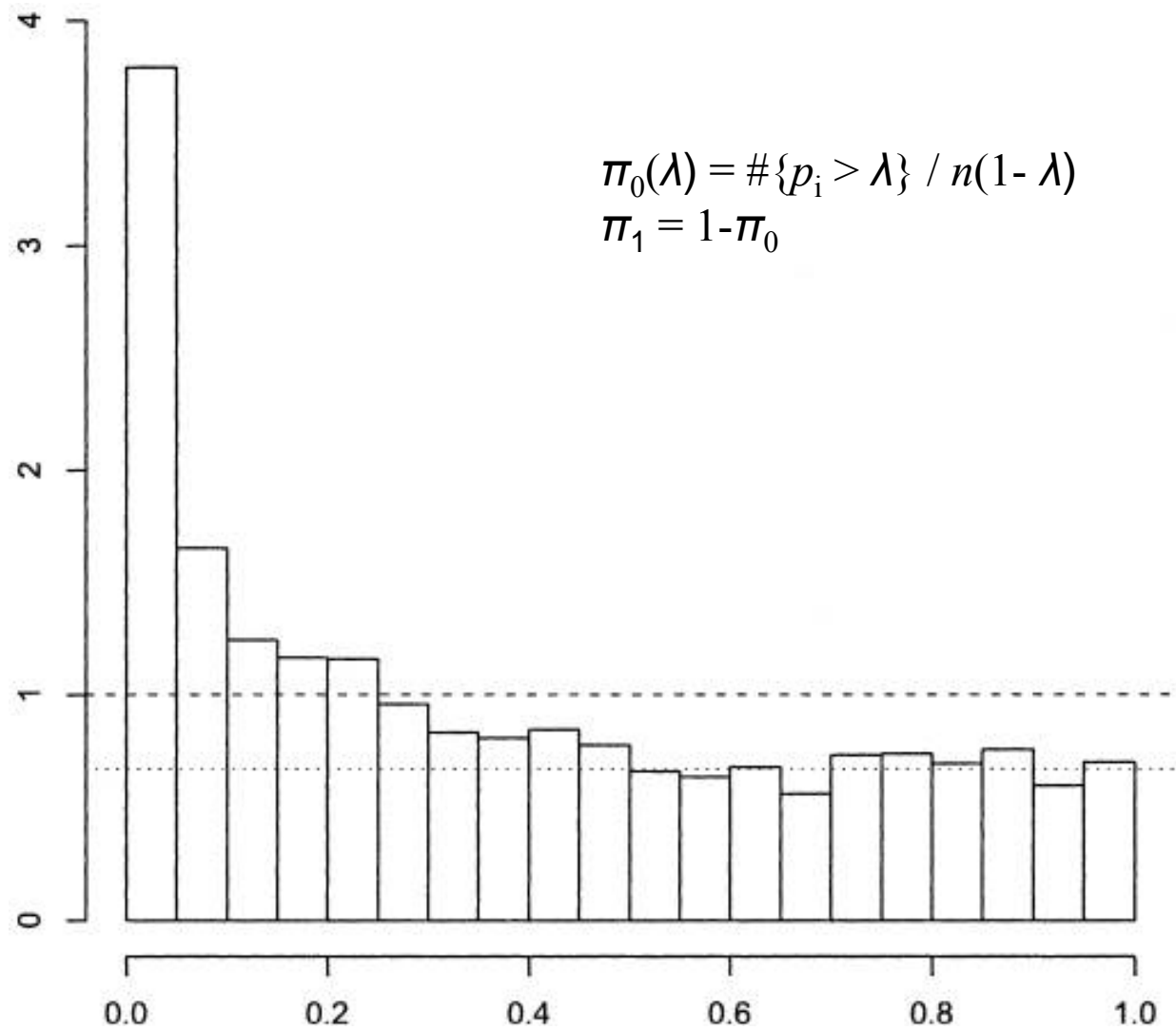


## Strategies

- Lenient: Significant at  $\alpha/n$  in Expt 1 and at  $\alpha$  in Expt 2
- Stringent: Significant at  $\alpha/n$  in Expt 1 and Expt 2
- In between:  $\pi_1$  statistics = proportion of non-nulls

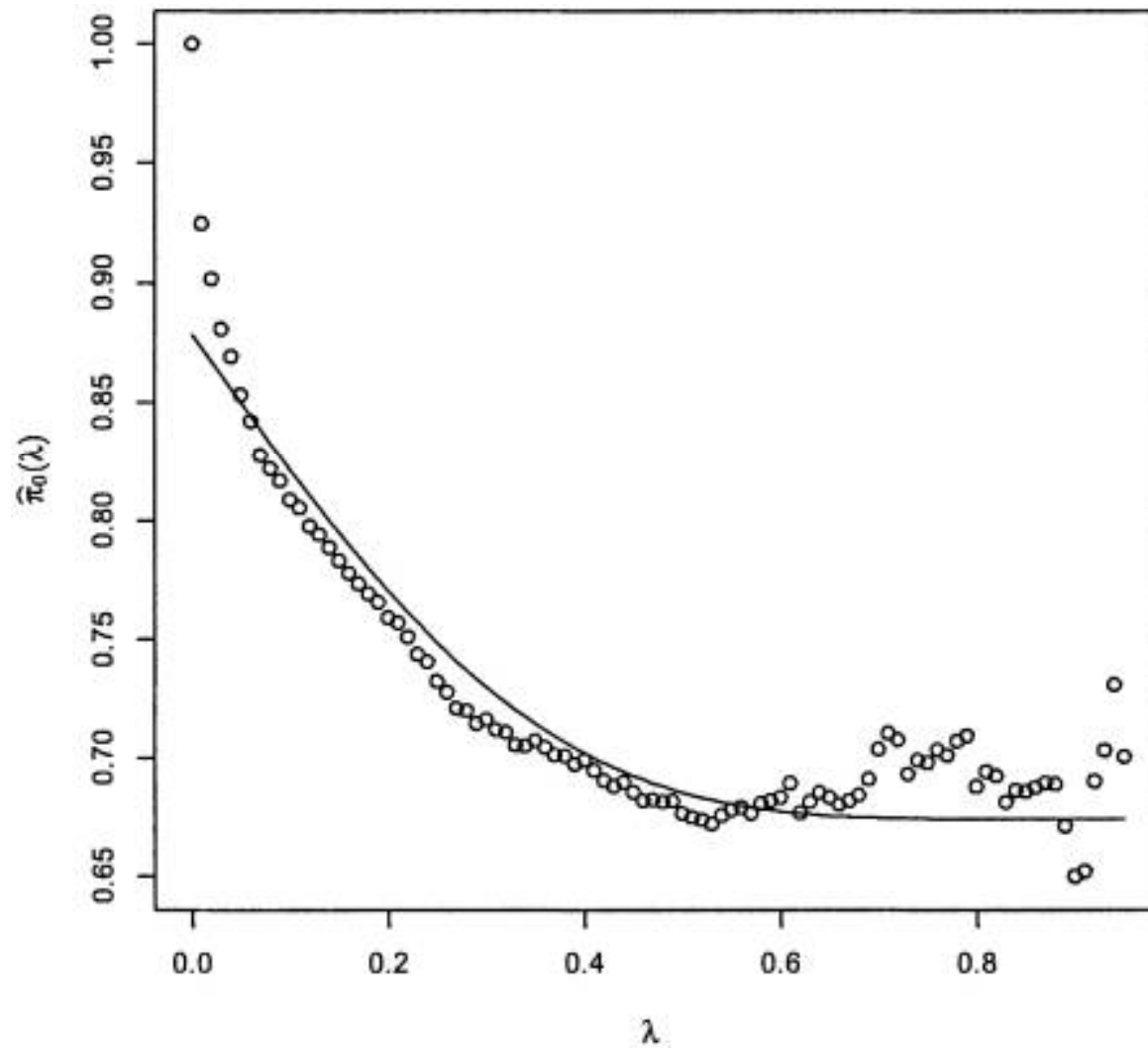


## $\pi_1$ Statistics

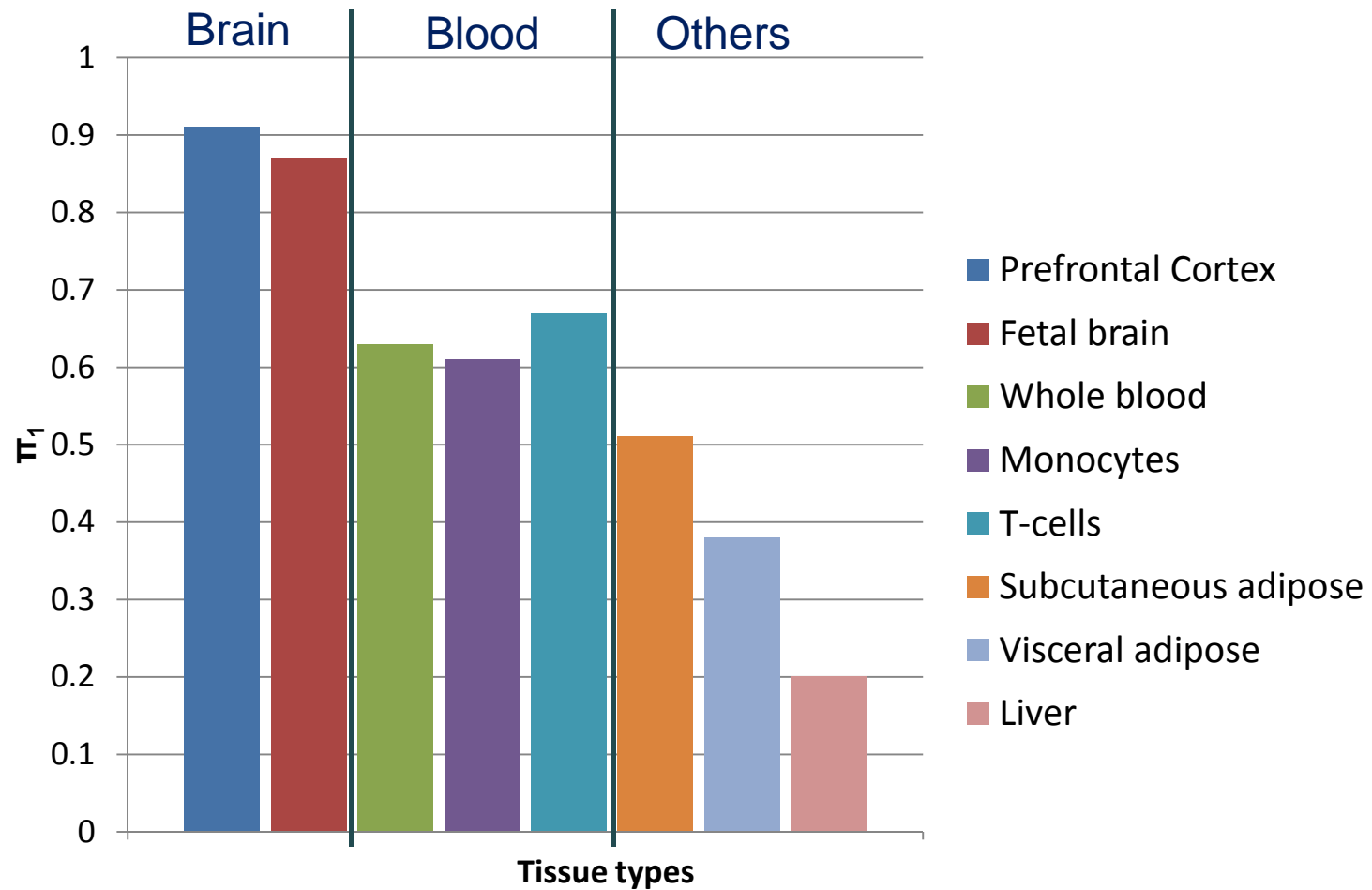




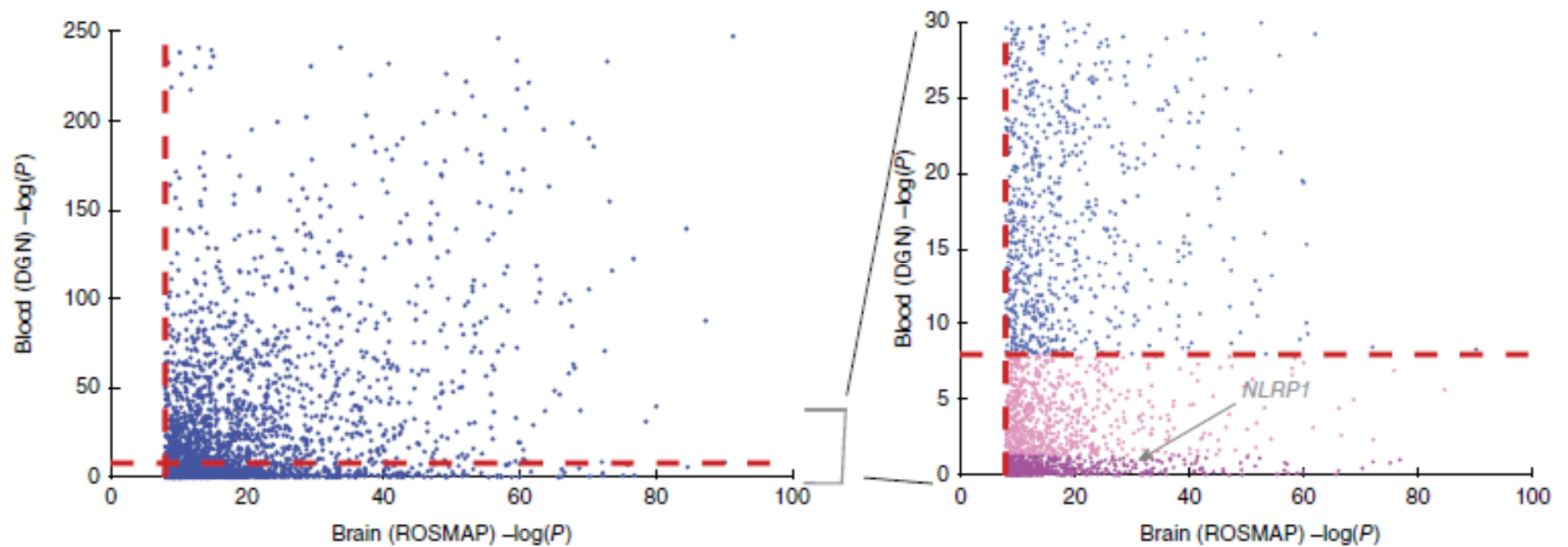
# $\Pi_1$ Statistics



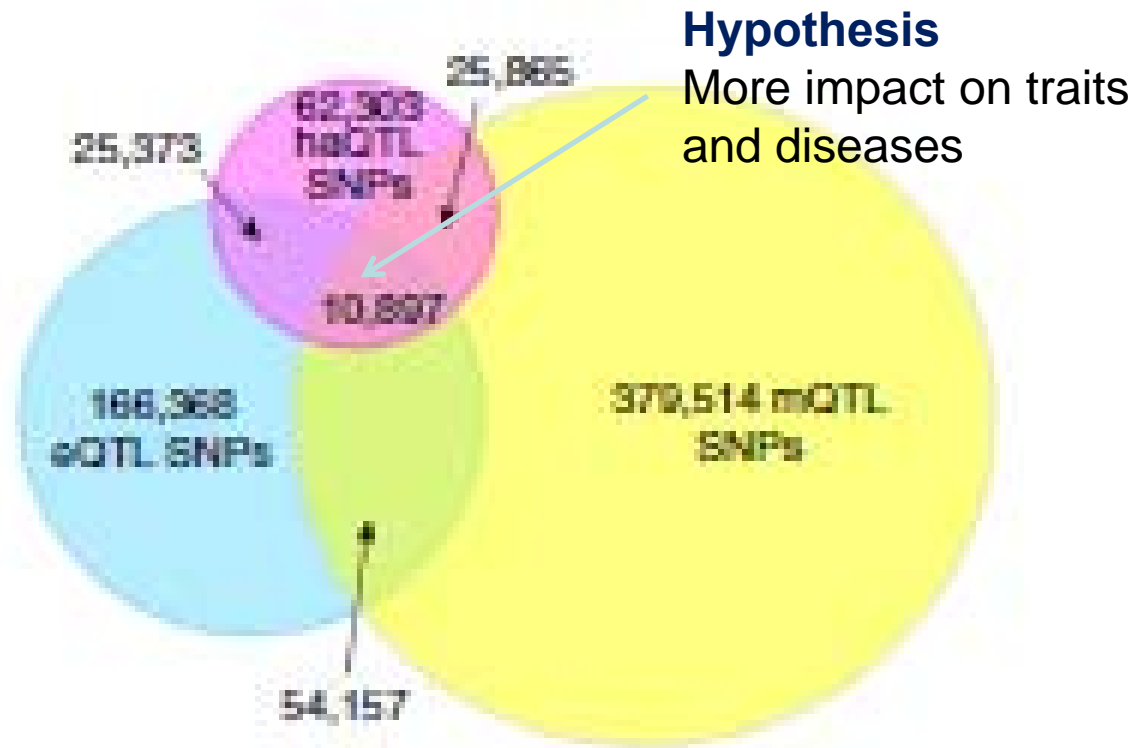
# Replication



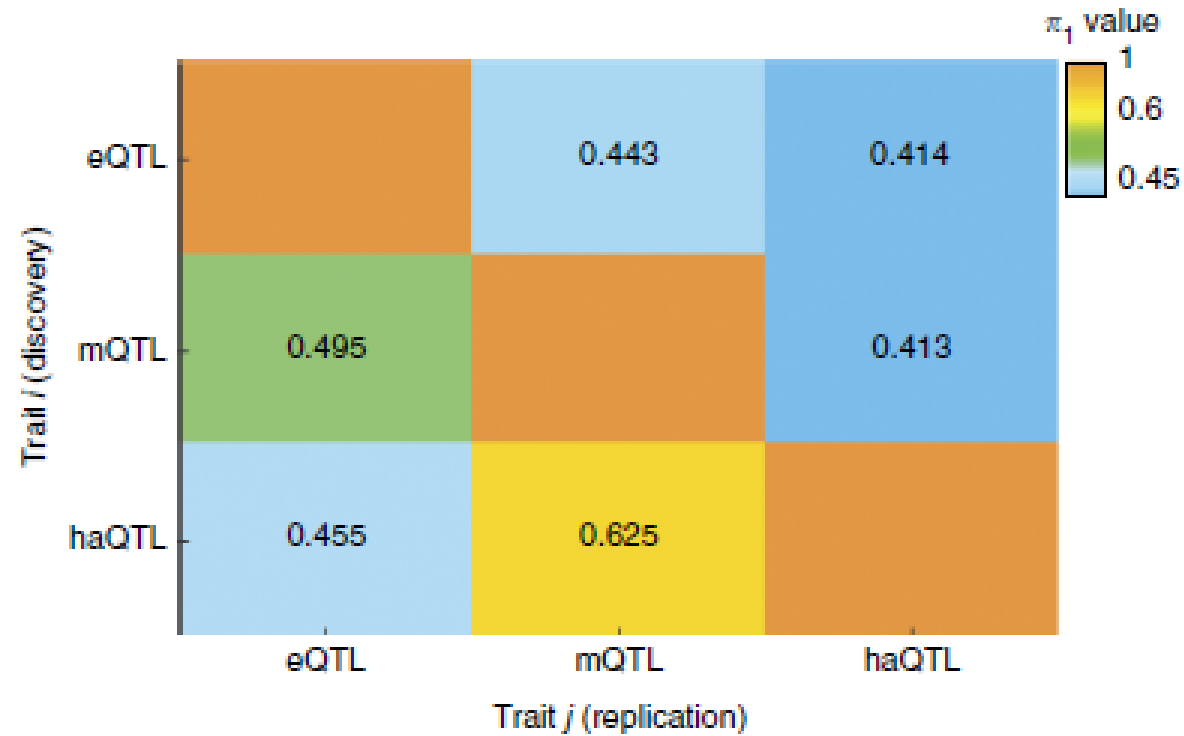
# Replication



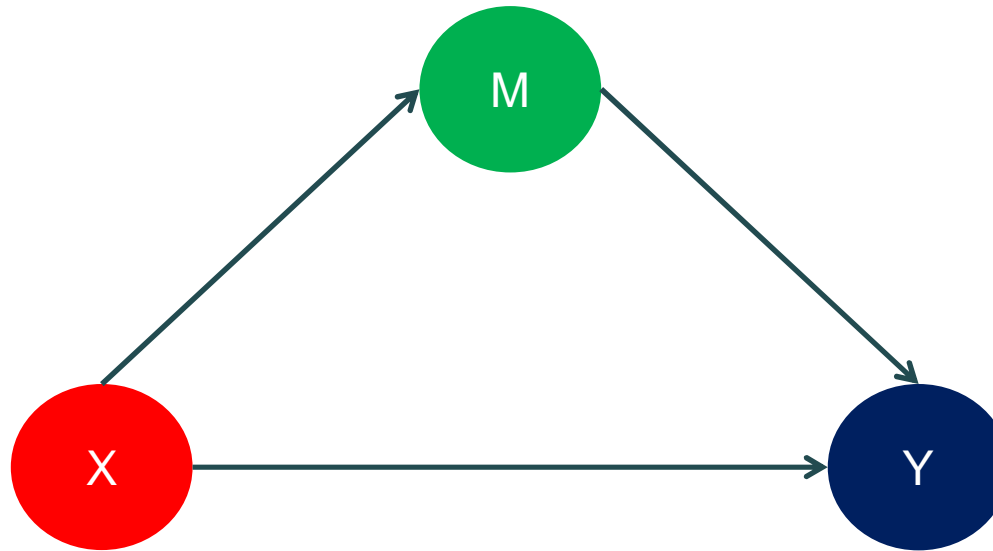
# xQTL Sharing



# xQTL Sharing



# Mediation Analysis

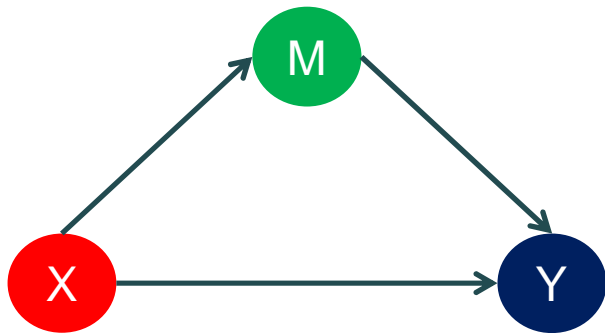


## Casual Inference Test

- X & Y associated:  $Y = X\beta + \varepsilon$  with  $H_1: \beta \neq 0$
- X & M|Y associated:  $M = X\beta + Y\gamma + \varepsilon$  with  $H_1: \beta \neq 0$
- M & Y|X associated:  $Y = X\beta + M\alpha + \varepsilon$  with  $H_1: \alpha \neq 0$
- $X \perp Y|M$ :  $Y = X\beta + M\alpha + \varepsilon$  with  $H_1: \beta = 0$

Requires  
equivalence test

# Mediation Analysis



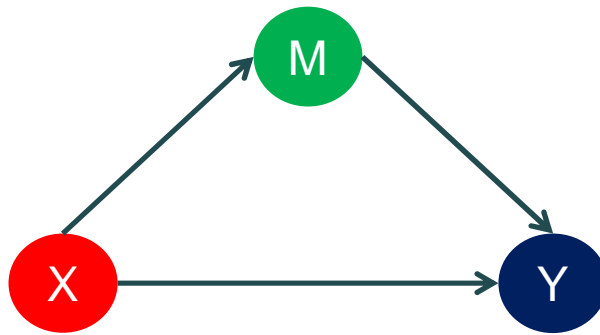
## Casual Inference Test

- X & Y associated:  $Y = X\beta + \varepsilon$  with  $H_1: \beta \neq 0$
- X & M|Y associated:  $M = X\beta + Y\gamma + \varepsilon$  with  $H_1: \beta \neq 0$
- M & Y|X associated:  $Y = X\beta + M\alpha + \varepsilon$  with  $H_1: \alpha \neq 0$
- $X \perp Y|M$ :  $Y = X\beta + M\alpha + \varepsilon$  with  $H_1: \beta = 0$

## Equivalent Test for $H_1: \beta = 0$ under independence model

- $M = X\beta + \varepsilon$
- $M^p = X\beta + \varepsilon^p$
- $Y = X\beta^p + M^p\alpha^p + \varepsilon$
- $\beta^p \rightarrow F^p$
- Generate  $m$   $F^p$ , then p-value =  $\#\{F > F^p\}/m$

# Mediation Analysis



## Casual Inference Test

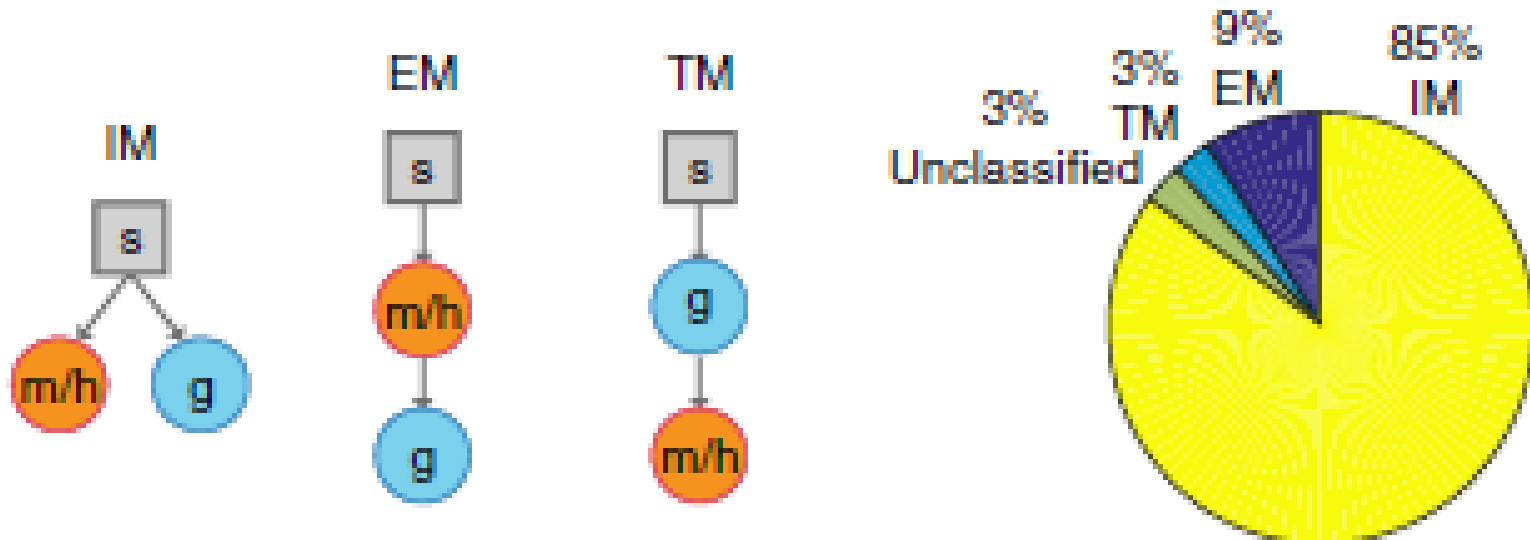
- $X$  &  $Y$  associated:  $p_1$
- $X$  &  $M|Y$  associated:  $p_2$
- $M$  &  $Y|X$  associated:  $p_3$
- $X \perp Y|M$ :  $p_4$

## Multiple Testing

- $p = \max(p_1, p_2, p_3, p_4)$  based on intersection-union test, i.e. weakest link
- Correct for number of mediation tested

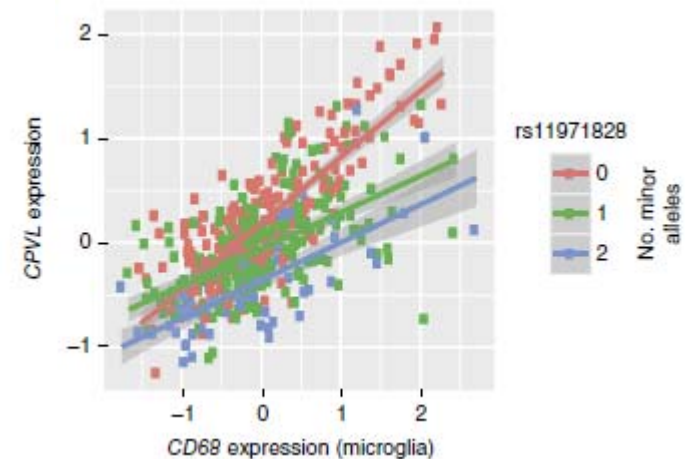


# Mediation Analysis

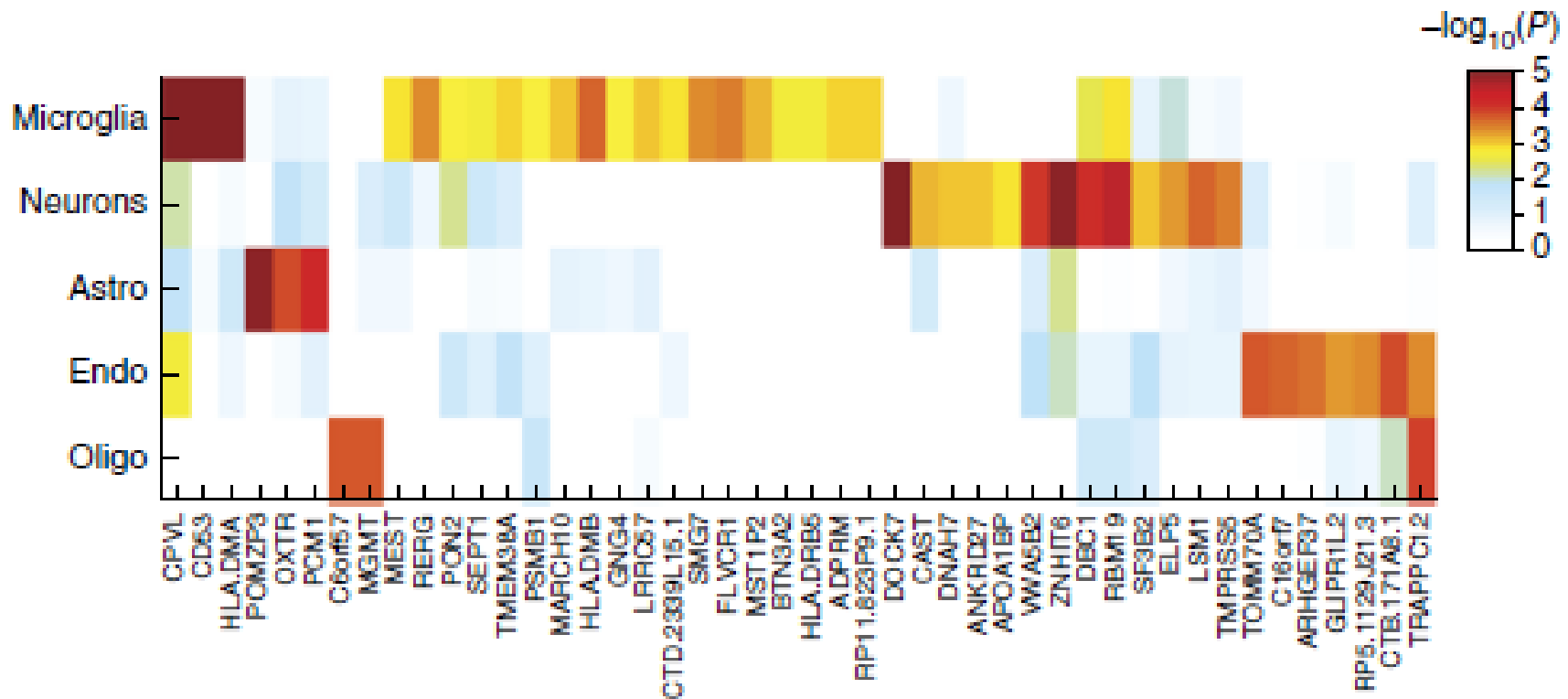


# Cell Type Specific Analysis

- $E_i = S_j\alpha + C\beta + S \cdot C\gamma + \varepsilon$ 
  - $E_i$  = Expression of gene  $i$
  - $S_j$  = Genotype values of SNP  $j$
  - $C$  = cell type proportion
  - $S \cdot C$  = element-wise product
- $C$  estimated based on expression markers
  - ENO2 for Neuron
  - OLIG2 for Oligodendrocyte
  - GFAP for Astrocyte
  - CD68 for Microglia
  - CD34 for Endothelial



# Cell Type Specific Analysis



- 46 genes at liberal FDR of 0.2
- This type of analysis falls under the general area of GxE interaction, which is nontrivial to detect

# Weighted GWAS

## Theory

- Given  $\{p_i\}$ ,  $i=1,\dots,n$
- If  $\{w_i p_i\}$ ,  $w_i \geq 0$  and  $\sum_i w_i = 1$ , then FWER controlled

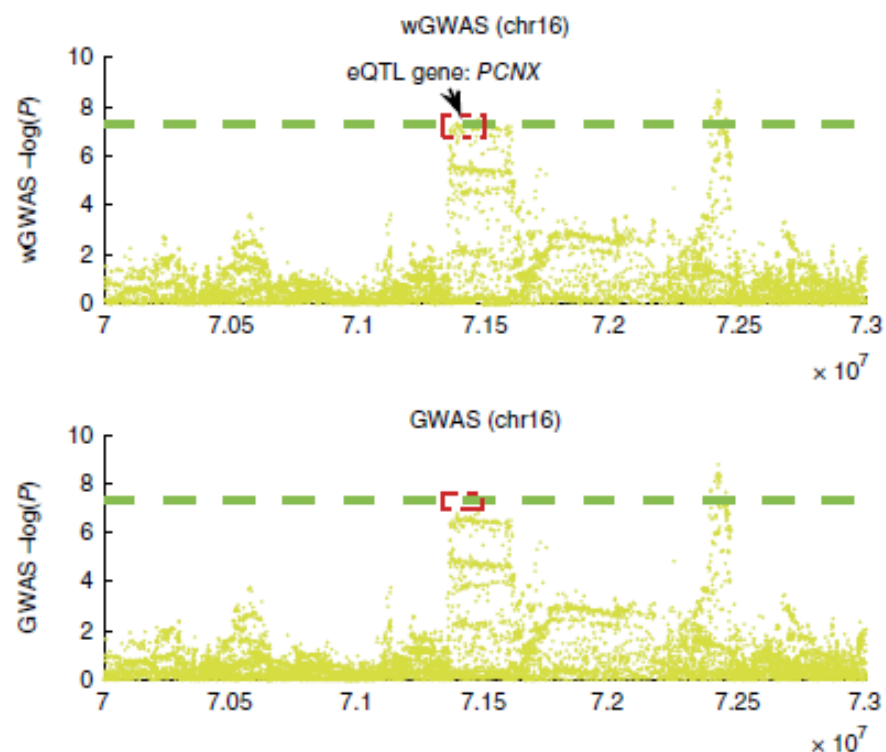
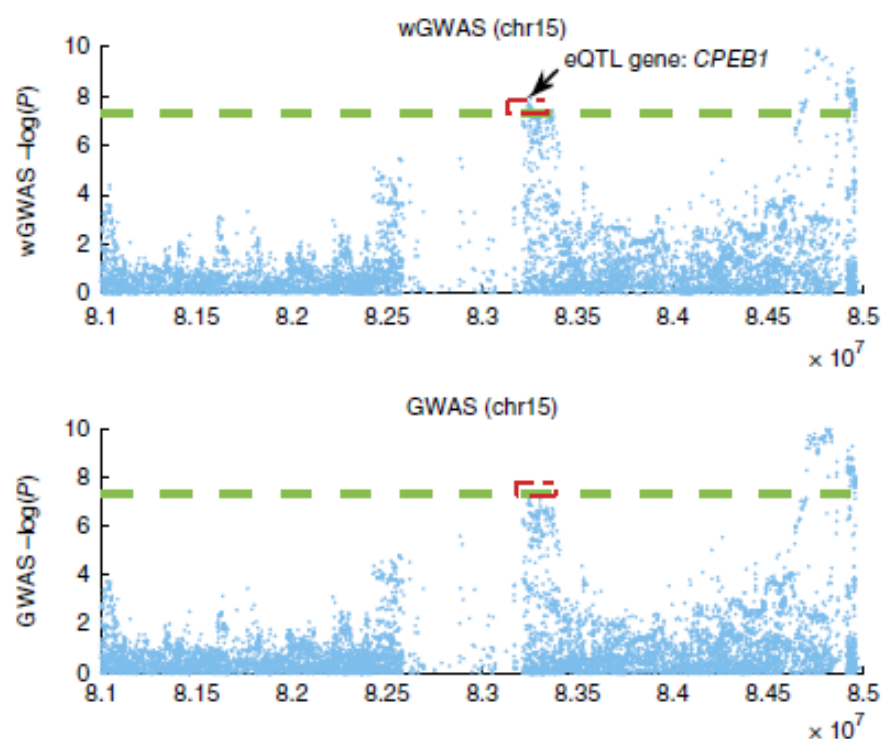
## Binary Strategy

- $w_1 = s/[1 + (s - 1)n_1/n]$
- $w_0 = 1/[1 + (s - 1)n_1/n]$
- $s = w_1/w_0$  ranging from 1 to 100
- $n_1 = \#xQTL\ SNPs$ ,  $n = \#SNPs$

## s Selection

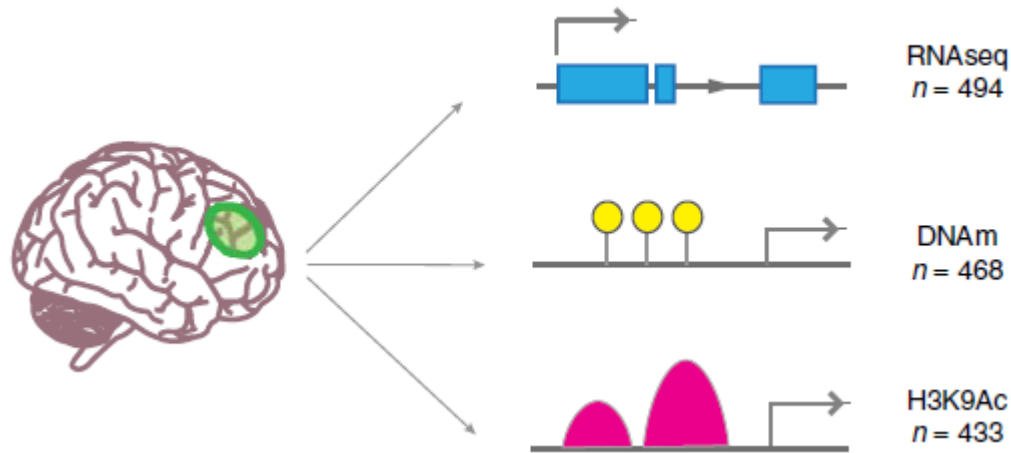
- Randomly split the SNPs into 2 halves
- $J(s) = (D^1(s)/\pi_1^1 + D^2(s)/\pi_1^2) / |D^1(s)/\pi_1^1 - D^2(s)/\pi_1^2|$

# Weighted GWAS



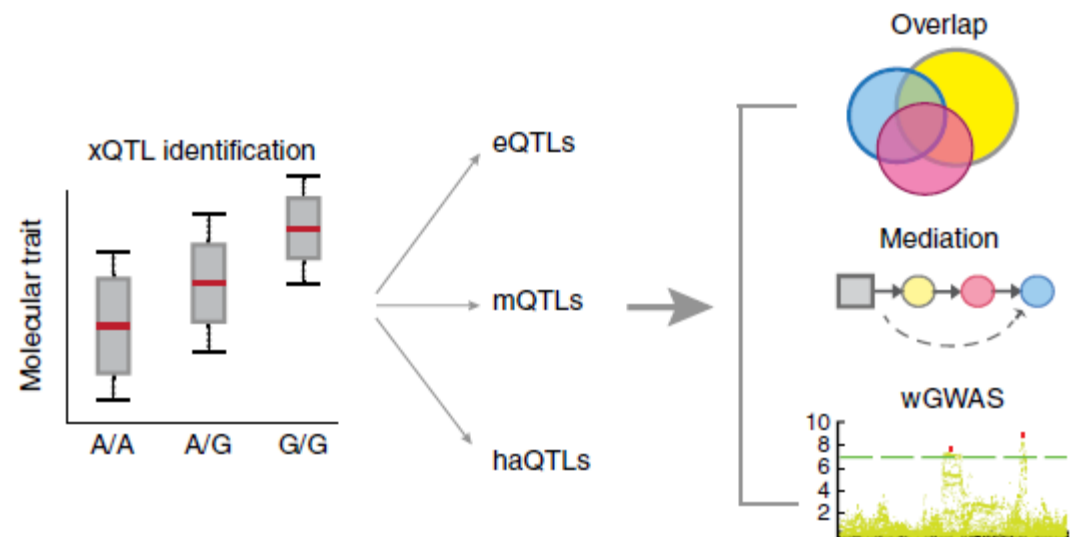
Applied on the largest Schizophrenia data set that found 108 loci, we found **18** additional loci that met genome-wide significance.

# Summary



- Recap of GWAS
- xQTL Analysis
- Challenges

- Replication
- xQTL Sharing
- Mediation Analysis
- Cell Specific Analysis
- Weighted GWAS






# xQTL Serve

## RESOURCE

nature  
neuroscience

**An xQTL map integrates the genetic architecture of the human brain's transcriptome and epigenome**

Bernard Ng<sup>1,2</sup>, Charles C White<sup>3</sup>, Hans-Ulrich Klein<sup>3,4</sup>, Solveig K Sieberts<sup>5</sup>, Cristin McCabe<sup>3</sup>, Ellis Patrick<sup>3</sup>, Jishu Xu<sup>3</sup> , Lei Yu<sup>6</sup>, Chris Gaiteri<sup>6</sup>, David A Bennett<sup>6</sup>, Sara Mostafavi<sup>1,2,7,8</sup>  & Philip L De Jager<sup>3,4,8</sup> 

**<http://mostafavilab.stat.ubc.ca/xQTLServe/>**