

Personal Genomics

& Data Science:

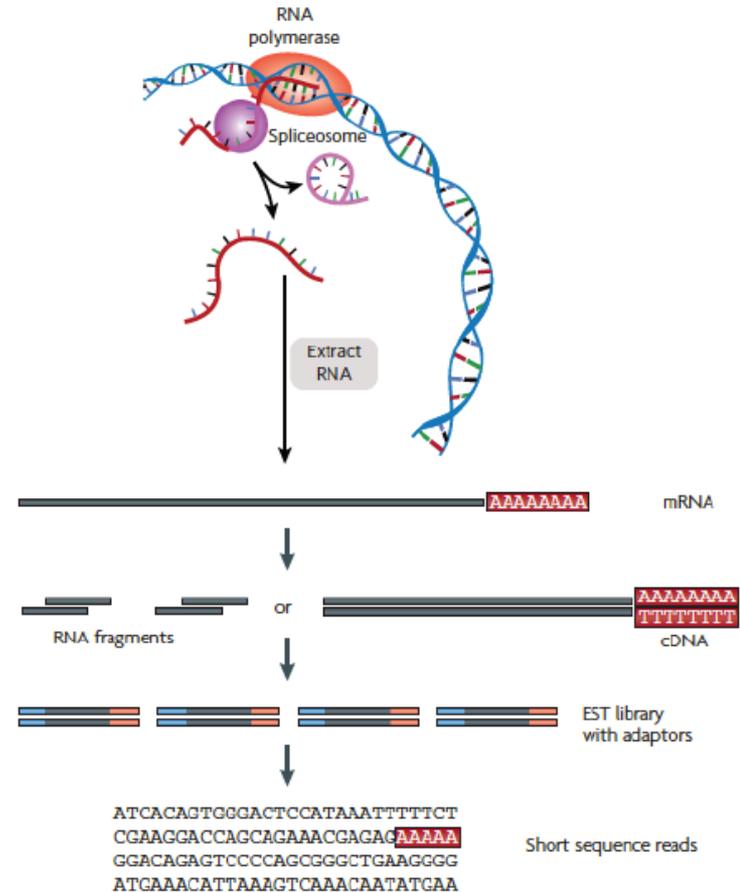
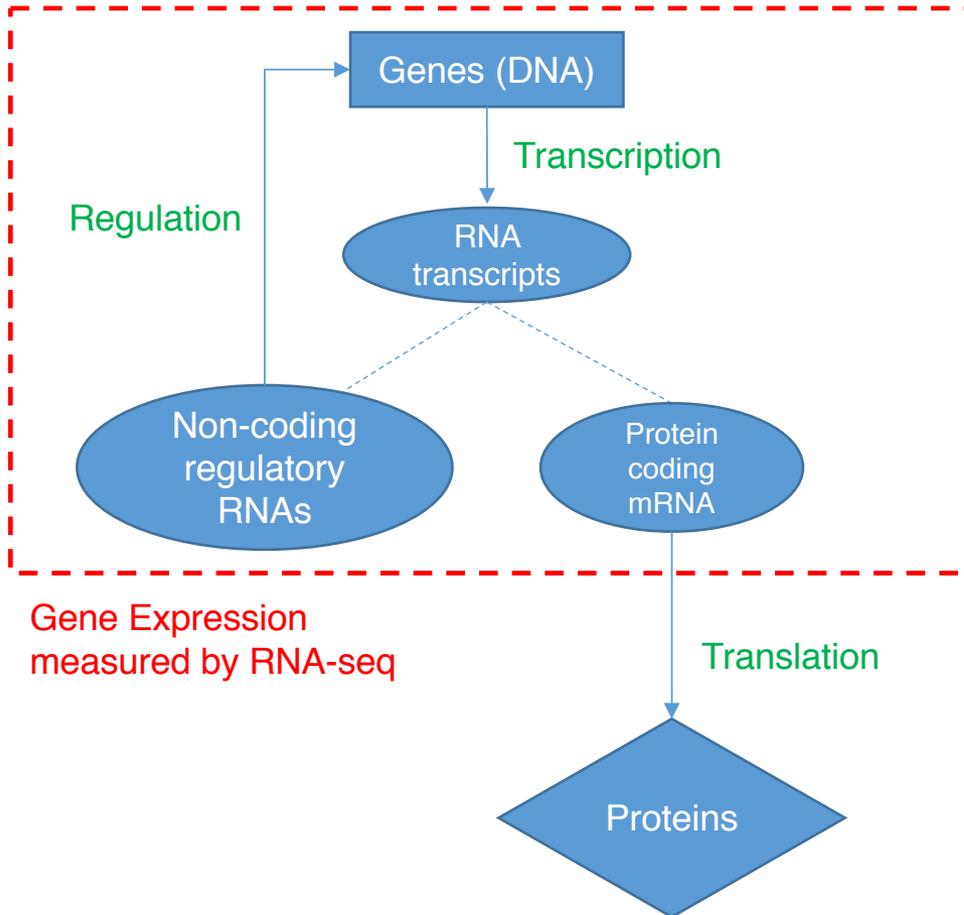
**Using population-scale  
functional genomics  
to understand  
neuropsychiatric  
disease  
& interpreting  
the data exhaust  
from this activity**

Mark Gerstein  
Yale

Slides freely downloadable from  
**Lectures.GersteinLab.org**  
& “tweetable” (via **@markgerstein**).

See last slide for more info.

# Transcriptome = Gene Activity of All Genes in the Genome, usually quantified by RNA-seq



Expression of genes is quantified by transcription:  
RNA-Seq measures mRNA transcript amounts

# RNA-Seq Overview

```
ATACAAGCAAGTATAAGTTCGTATGCCGTCTT
GGAGGCTGGAGTTGGGGACGTATGCGGCATAG
TACCGATCGAGTCGACTGTAAACGTAGGCATA
ATTCTGACTGGTGTCATGCTGATGTACTTAAA
```

Fastq sequence files  
~5-10 GB

Index-building + Alignment to reference genome

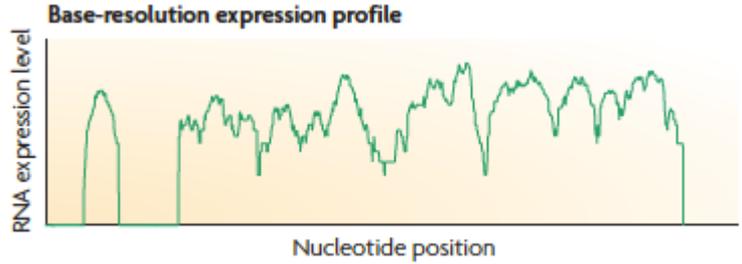
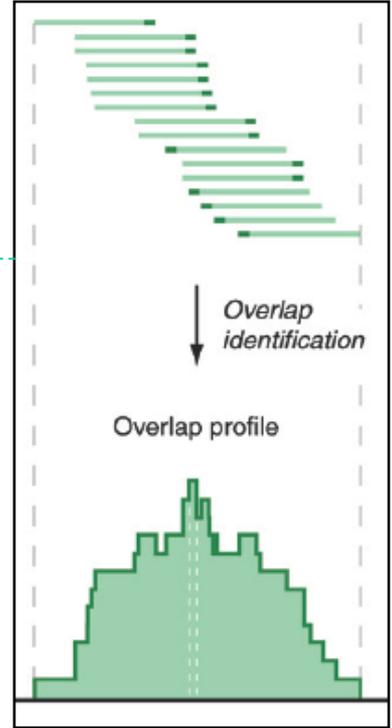
BAM files  
~1-2-fold reduction

Conversion to signal track by overlapping reads

BigWig files  
~25-fold reduction

Mapping to genes

Gene/Transcript expression matrix  
~20-fold reduction



Quantitative information from RNA-seq signal:  
average signals at exon level (RPKM)

Reads => Signal

Successive steps of Data Reduction

[NAT. REV. 10: 57; PLOS CB 4:e1000158; PNAS 4:107: 5254 ]



## Activity Patterns

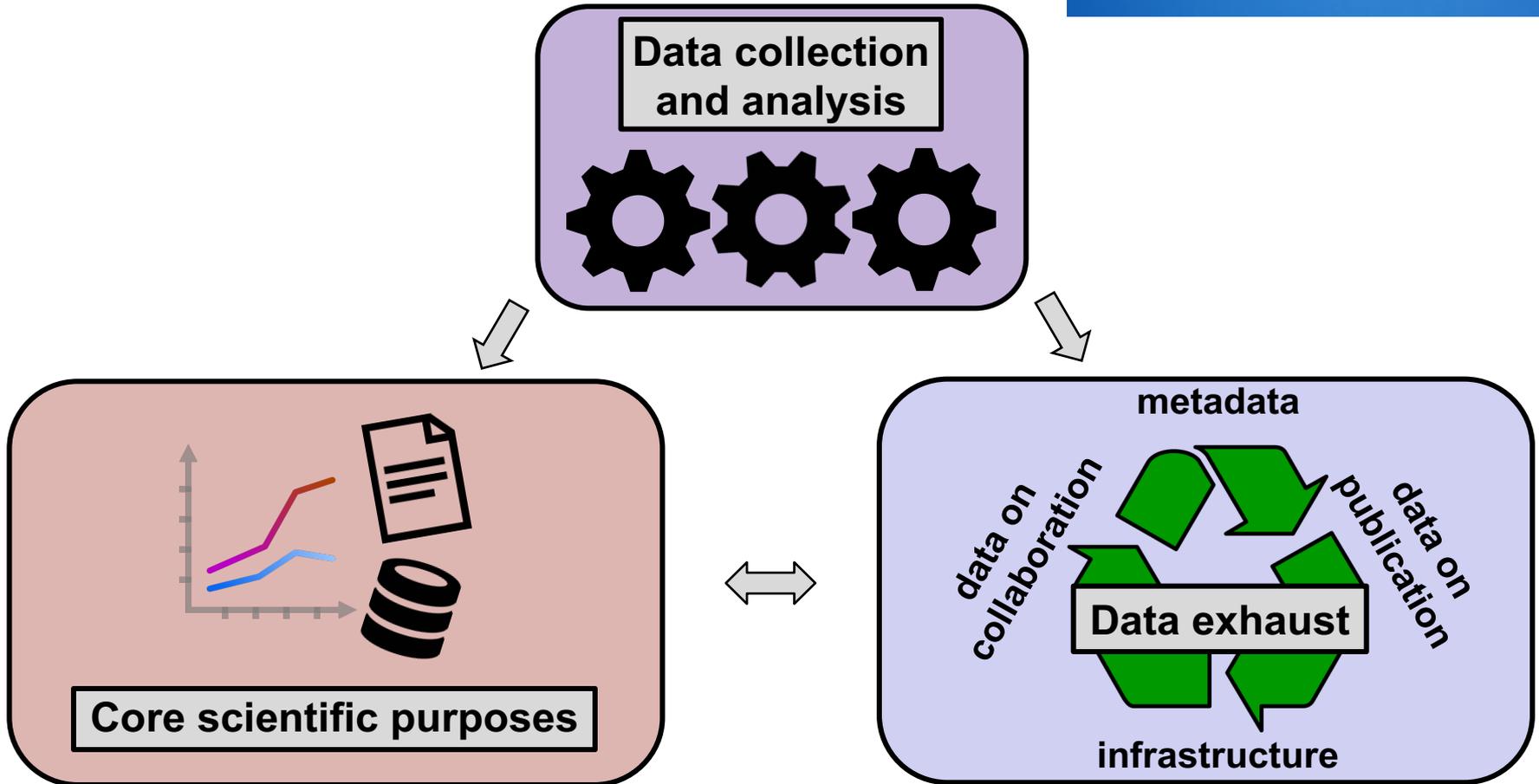
- RNA Seq. gives rise to activity patterns of genes & regions in the genome

# Some Core Science Qs Addressed by RNA-seq

- Gene activity as a function of:
  - **Developmental** stage: basic patterns of co-active genes across development
  - **Cell-type** & Tissue: relationship to specialized functions
  - **Evolutionary** relationships: behavior preserved across a wide range of organisms; patterns in model organisms in relation to those in humans
  - **Individual**, across the human population
  - **Disease** phenotypes: disruption of patterns in disease
- Some overarching Qs:
  - Are there core patterns of gene activity ?**
  - How do they vary across individual ?**
  - Are they disrupted by disease?**

# Data Exhaust

- Creative use of data is key to data science!
- Data exhaust = exploitable byproducts of big data collection and analysis



[photos: wikipedia/wikimedia]

## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

# Sample Sources: >2,500 brains

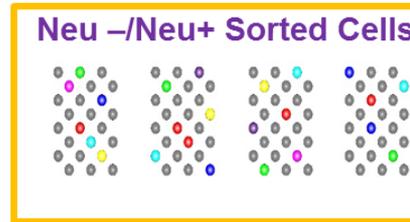
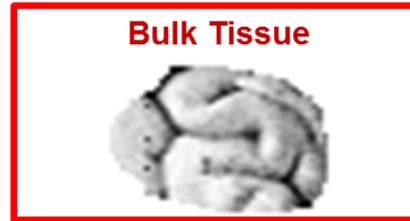
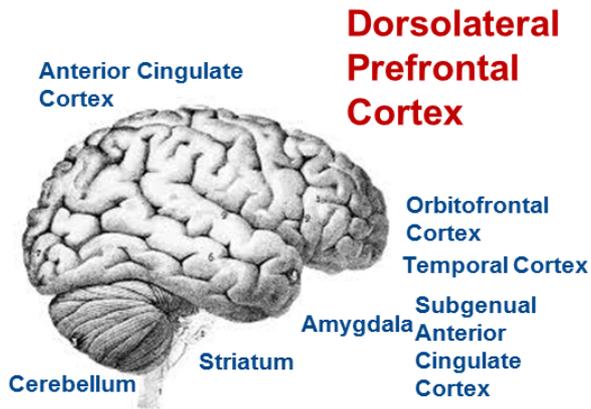
**Genome:**  
WGS, genotype

**Epigenome:**  
ChIP-seq, ATAC-seq, HiC, ERRBS, Array Methylation, NOMeSeq

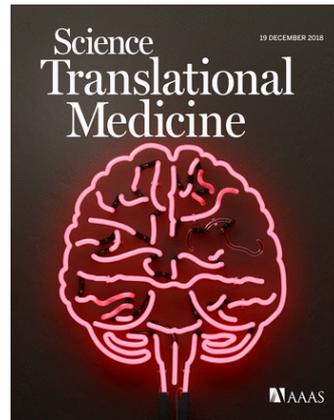
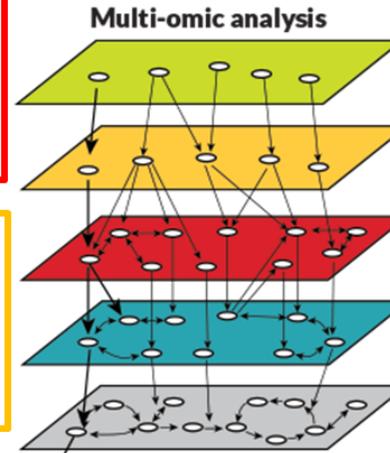
**Transcriptome:**  
RNA-seq, IncRNAseq,

**Proteome:**  
MWP, LC-MS/MS

Cross-disorder: ASD, SCZ, BP, Neurodevelopmental, Neurotypical



Single Cell



# PsychENCODE

## '18 rollout in Science

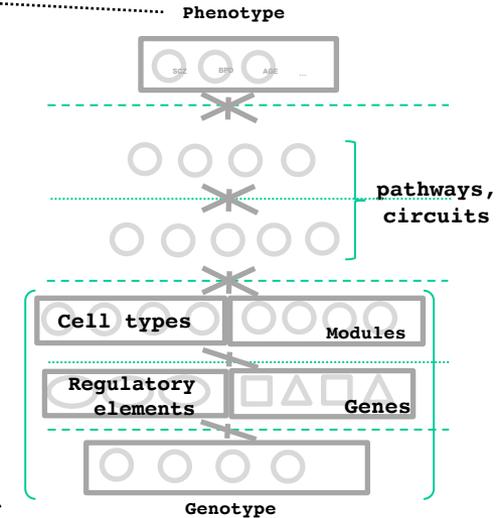
11 papers in total.

Major material in the 3 capstones:

Wang et al. ('18), Li et al. ('18), Gandal et al. ('18)

# A core issue addressed by PsychENCODE: Using functional genomics to reveal molecular mechanisms between genotype and phenotype in brain disorders

| Disease                    | Heritability* | Molecular <b>Mechanisms</b>             |
|----------------------------|---------------|---|
| <b>Schizophrenia</b>       | <b>81%</b>    | <b>(C4A)</b>                            |
| <b>Bipolar disorder</b>    | 70%           | -                                       |
| <b>Alzheimer's disease</b> | 58 - 79%      | Apolipoprotein E (APOE), Tau            |
| <b>Hypertension</b>        | 30%           | Renin–angiotensin–aldosterone           |
| <b>Heart disease</b>       | 34-53%        | Atherosclerosis, VCAM-1                 |
| <b>Stroke</b>              | 32%           | Reactive oxygen species (ROS), Ischemia |
| <b>Type-2 diabetes</b>     | 26%           | Insulin resistance                      |
| <b>Breast Cancer</b>       | 25-56%        | BRCA, PTEN                              |



Many psychiatric conditions are highly heritable

Schizophrenia: up to 80%

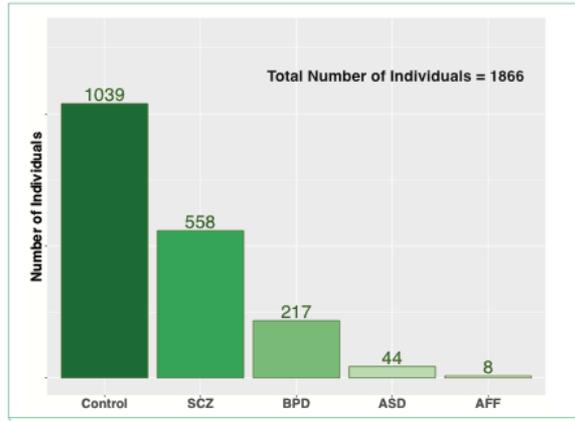
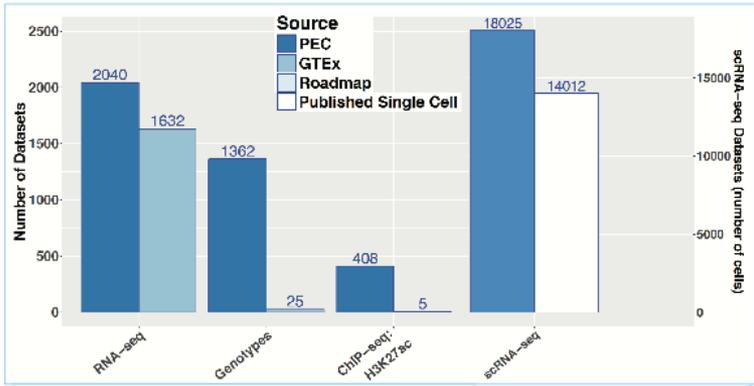
But we don't understand basic molecular mechanisms underpinning this association

(in contrast to many other diseases such as cancer & heart disease)

Thus, interested in developing predictive models of psychiatric traits which:

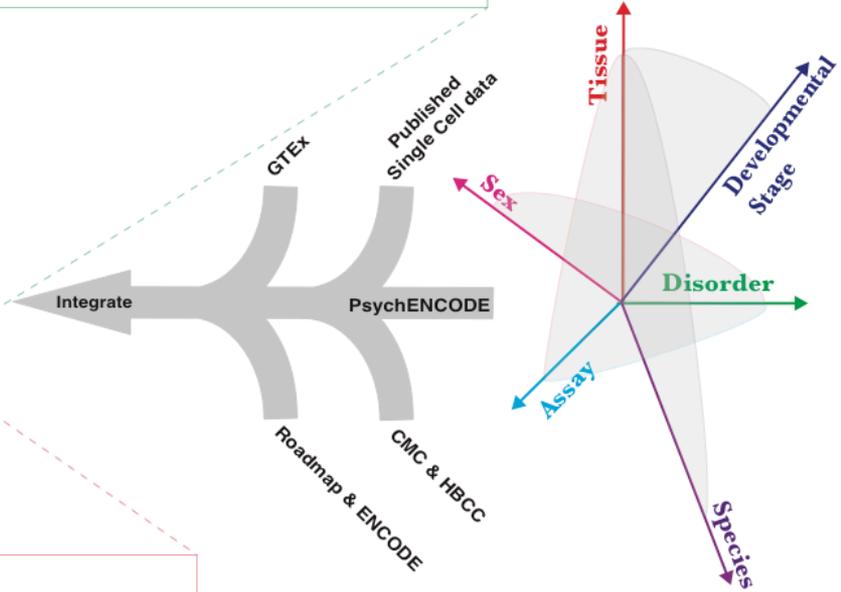
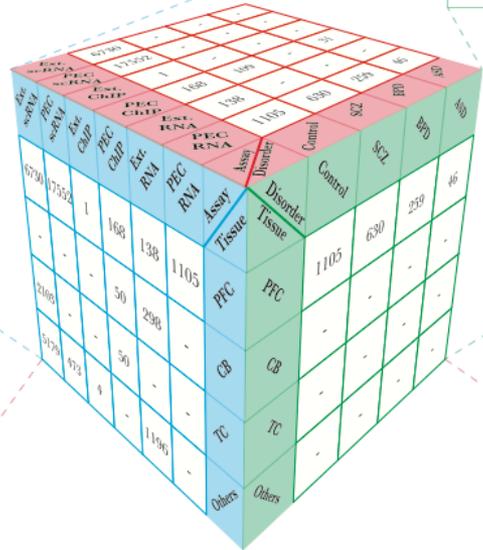
Use observations at intermediate (molecular levels) levels to inform latent structure

Use the predictive features of these “molecular endo phenotypes” to begin to suggest actors involved in mechanism



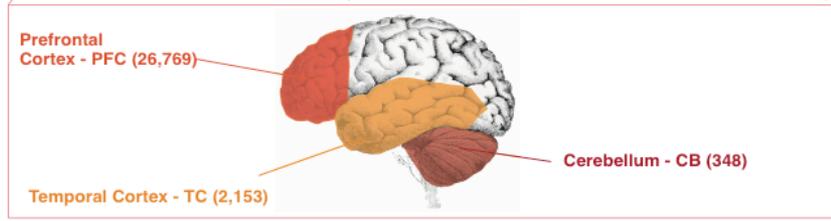
**1866**  
Individuals  
~3.7K bulk RNA-seq  
~32K single-cells

**Disorder**



**Collecting functional genomic datasets for the adult brain**

from PsychENCODE, other large consortia & single cell studies

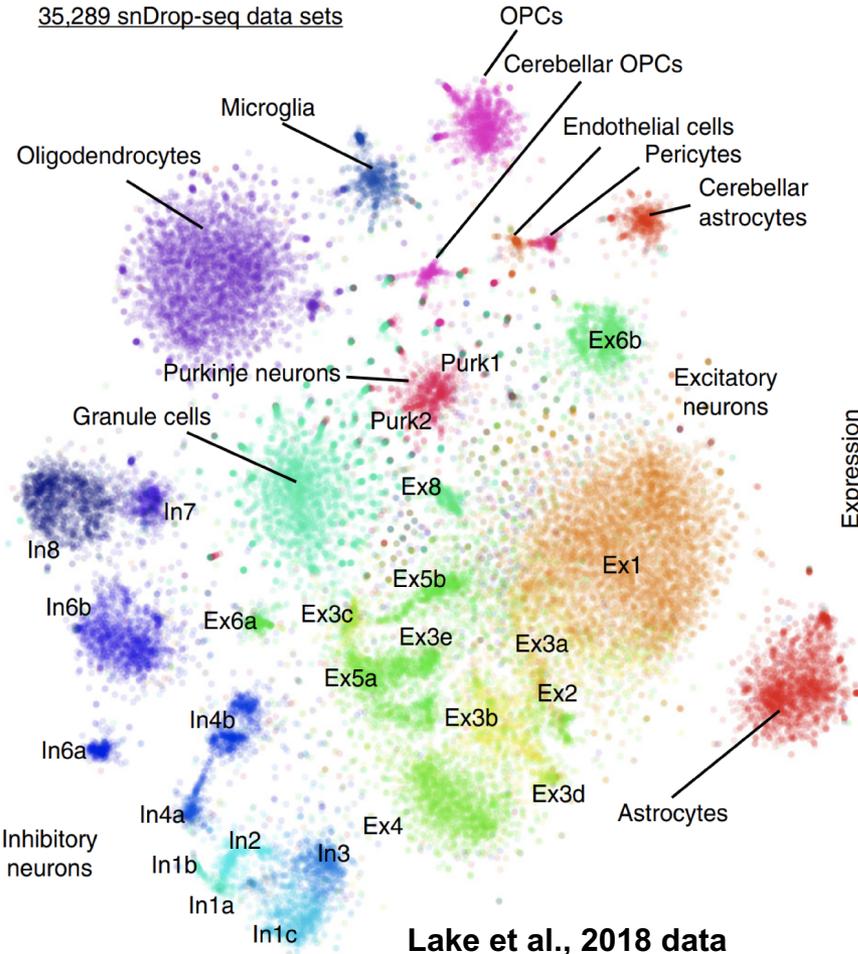


# Merging & Clustering Single Cell Data Sets

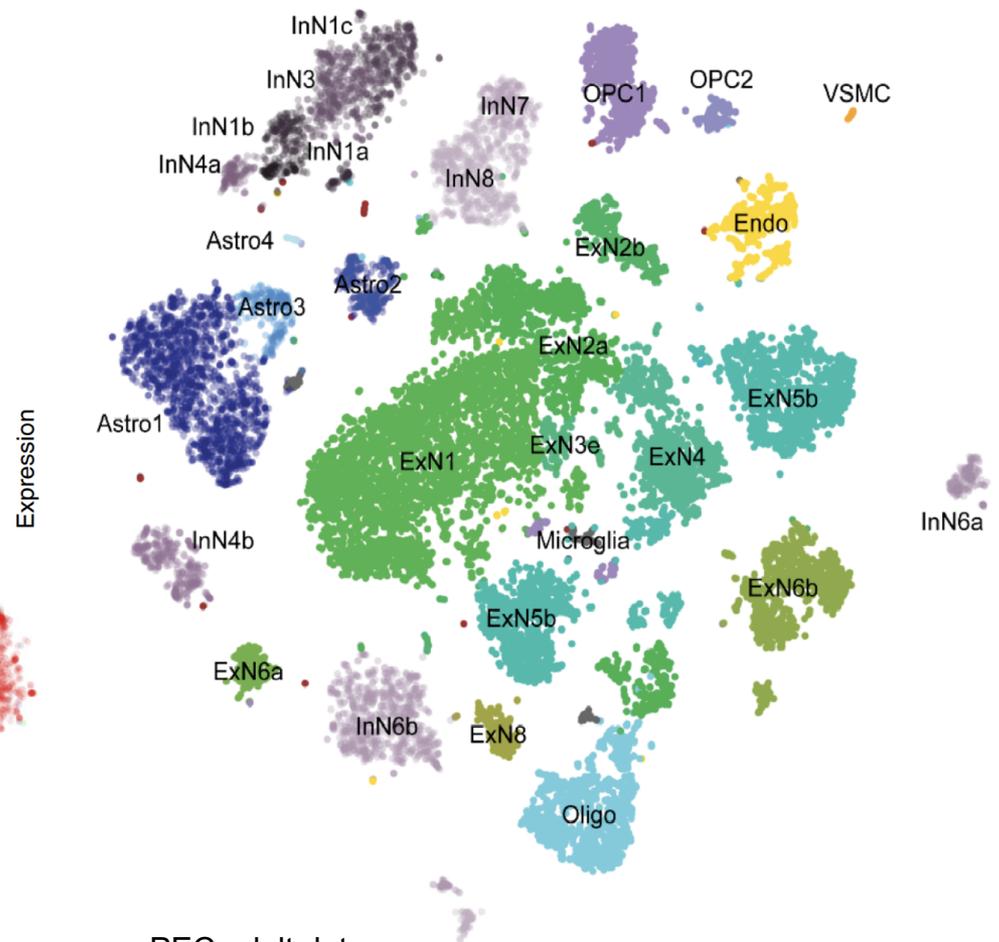
Single cell signatures, from:

- ~14K cells (Lake et al., '16 & '18)
- ~400 cells (Darmanis et al., PNAS, '15)
- ~18K cells (PsychENCODE)

35,289 snDrop-seq data sets



Lake et al., 2018 data

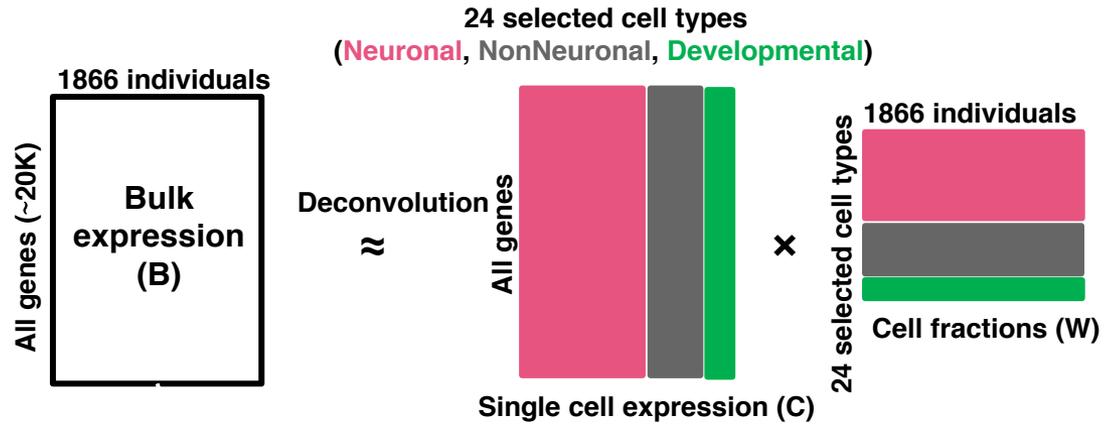


PEC adult data

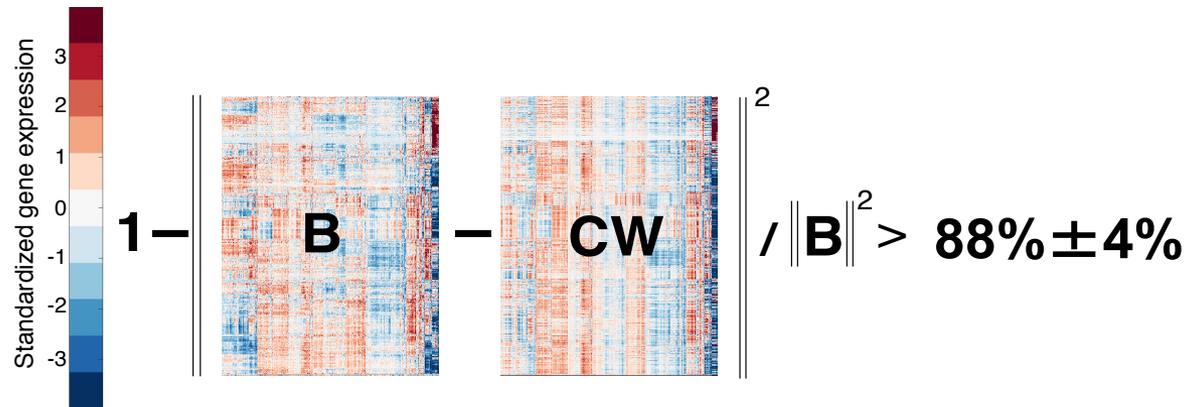
[Li et al. ('18), Science. Wang et al. ('18). Science]

# Single-cell deconvolution

## Step 1:

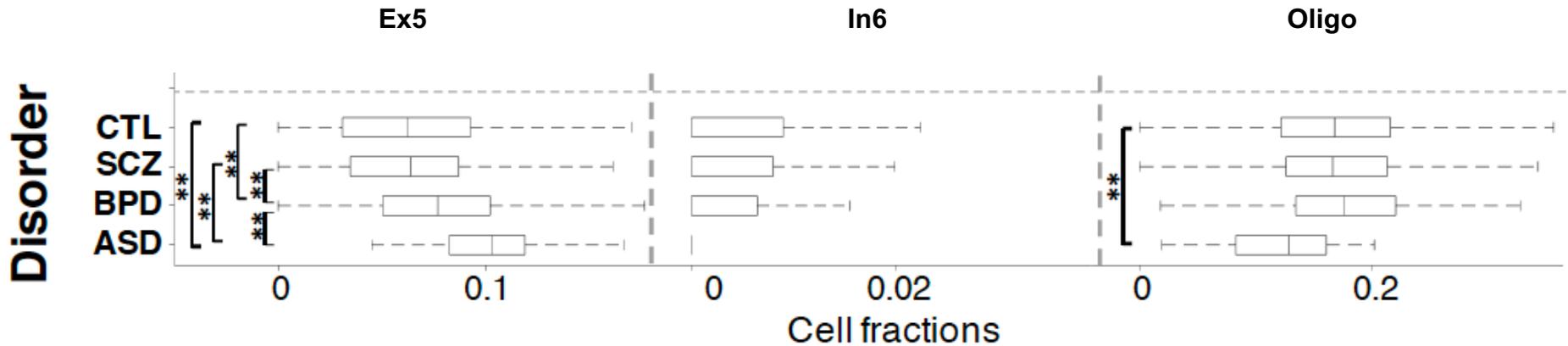


# Supervised learning to estimate cell fractions



Individual and cross-population reconstruction accuracy via deconvolution

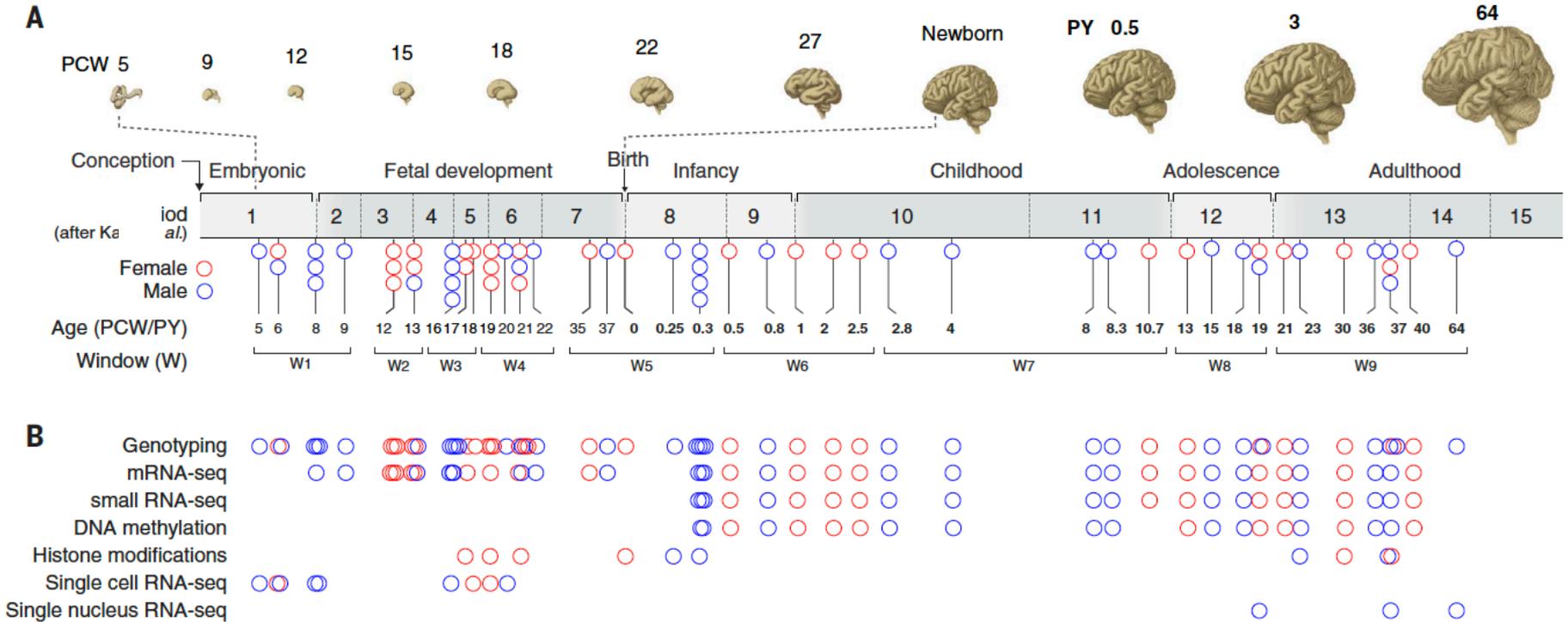
# Different neuronal & glial cell fractions across disorders



Excitatory to Inhibitory imbalance at neuronal subtype level for ASD\*

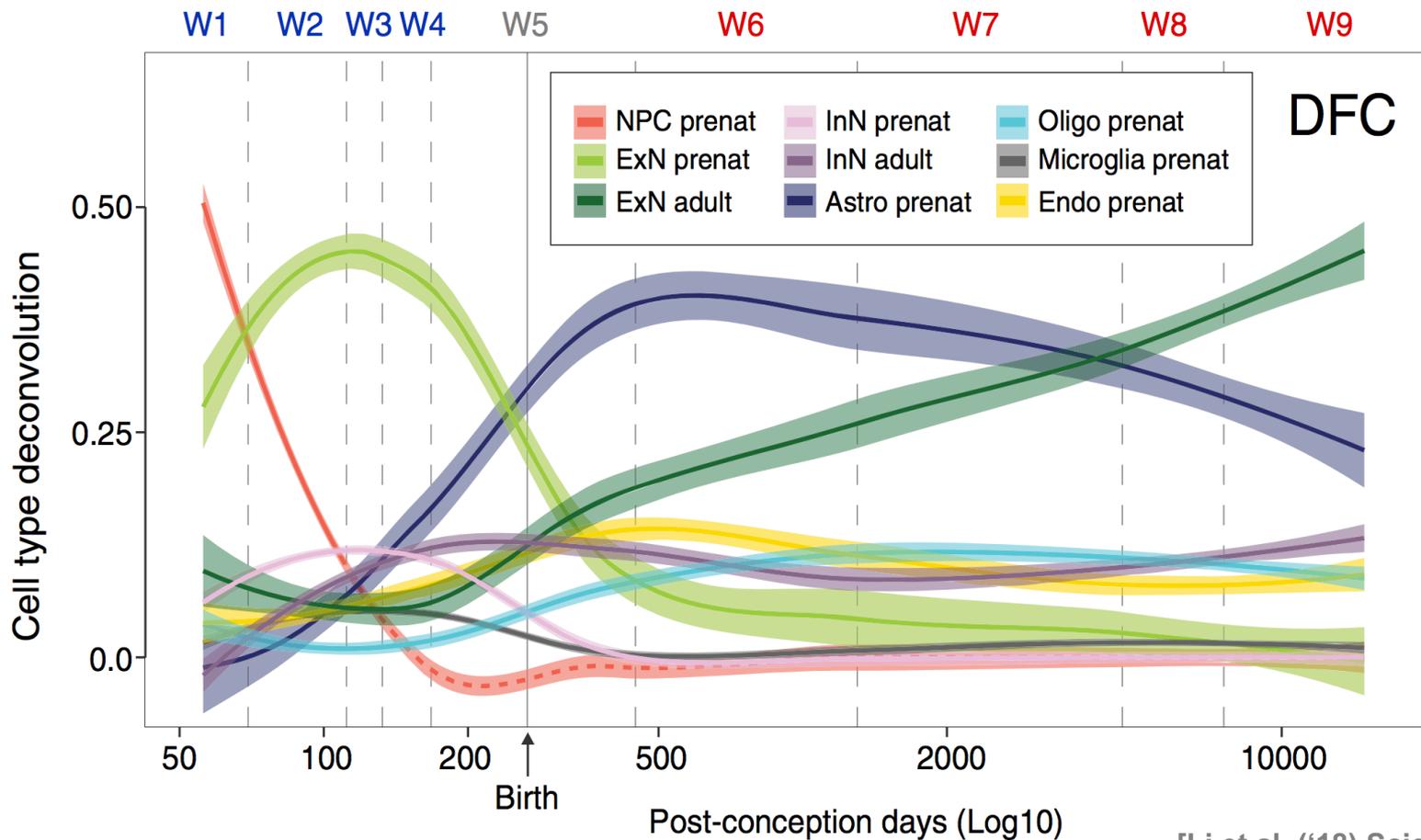
\* Rubenstein et al., Model of autism: increased ratio of excitation/inhibition in key neural systems, Genes Brain Behav. 2003

# Developmental Capstone Data Set



- 60 Individuals in total
- Ages from 5 PCW to 64 yrs.
- 16 brain regions for > 9 PCW

# Different neuronal & glial cell fractions across ages

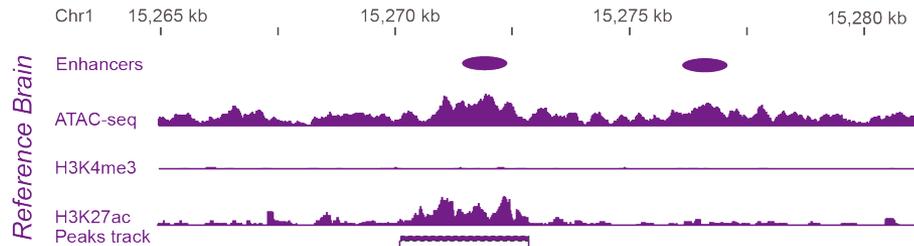


[Li et al. ('18) Science]

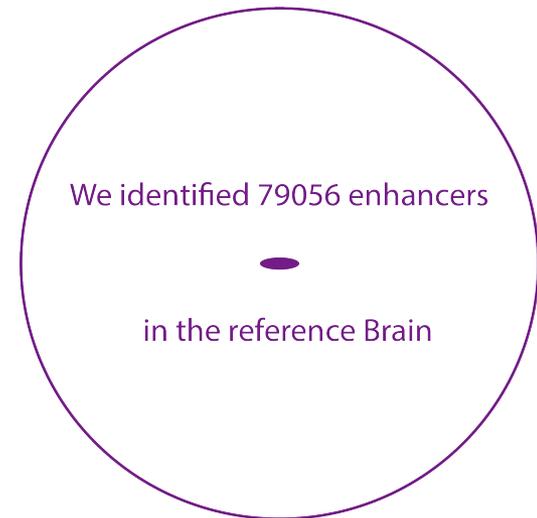
## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

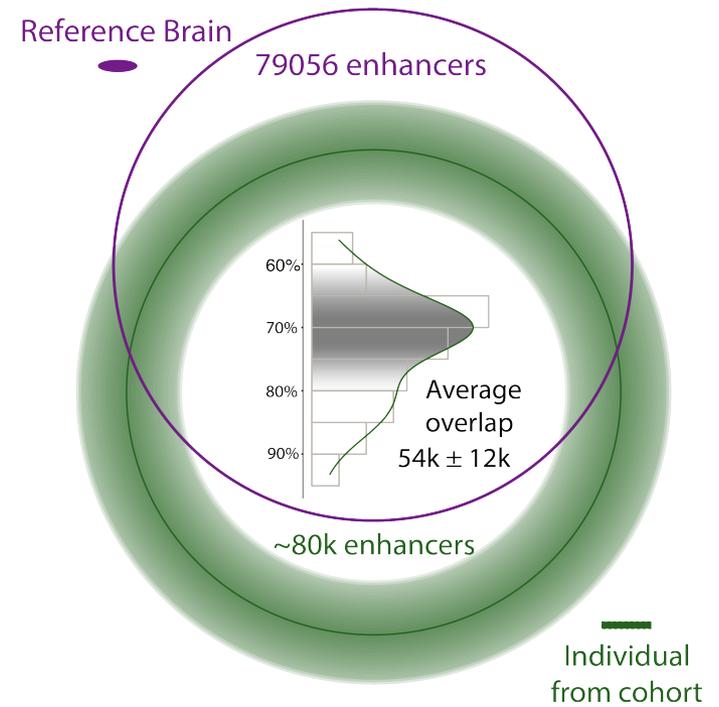
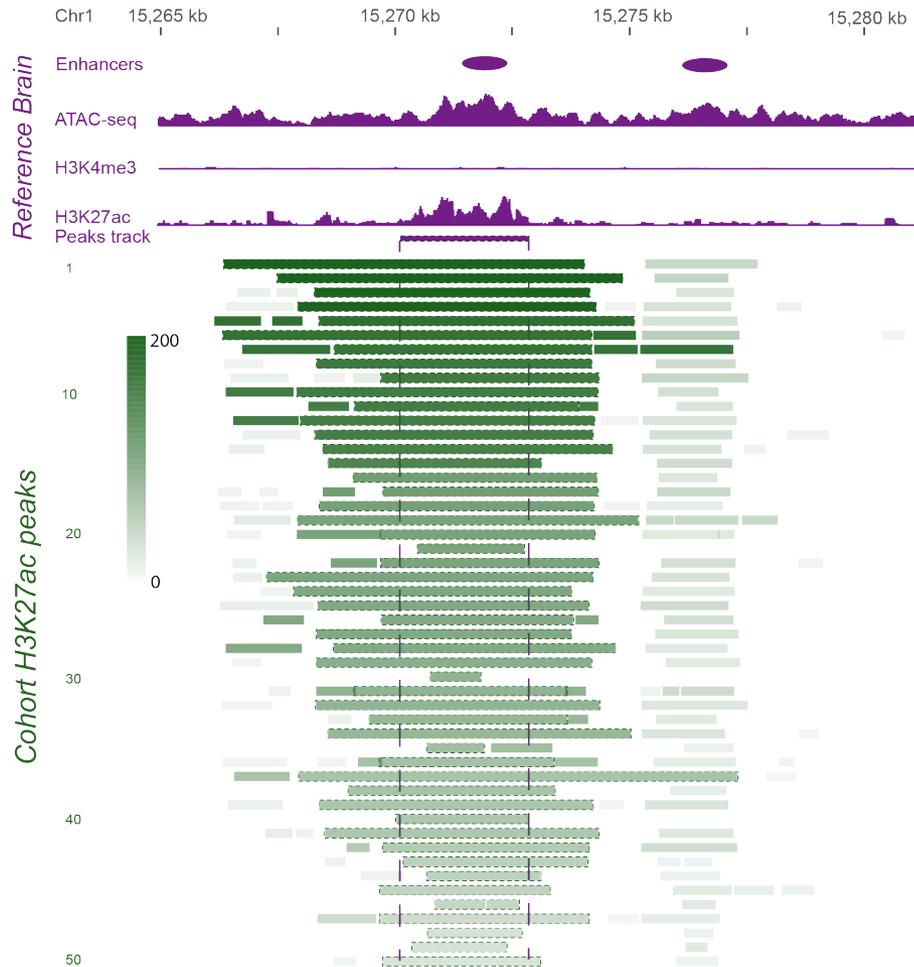
# Developing a Reference Set of ~79K PFC Enhancers & Studying Their Population Variation



**Consistent with ENCODE, active enhancers are identified as open chromatin regions enriched in H3K27ac and depleted in H3K4me3**



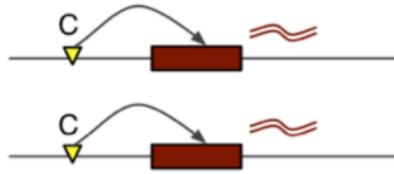
# Developing a Reference Set of ~79K PFC Enhancers & Studying Their Population Variation



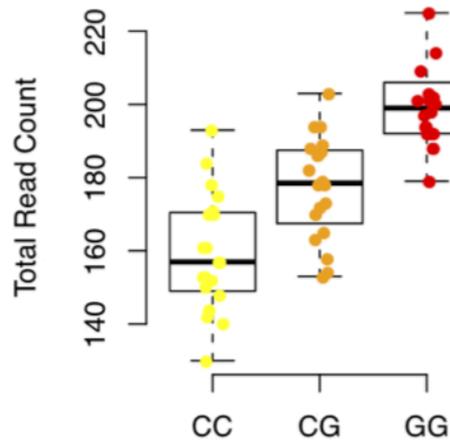
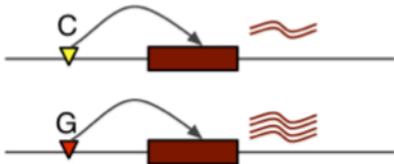
# Quantitative Trait Loci (QTLs) associated with variation

Gene expression (eQTL)

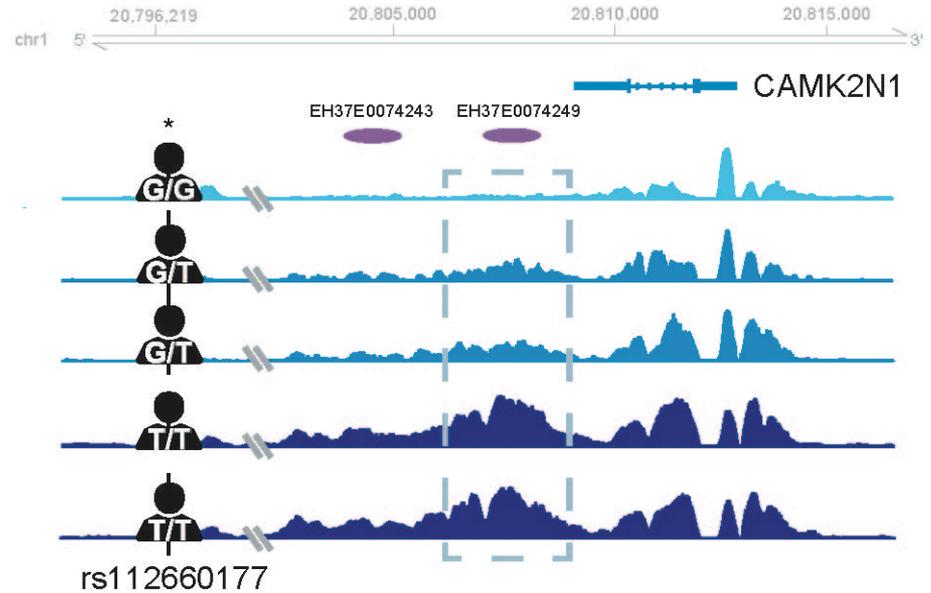
Sample 1: genotype CC



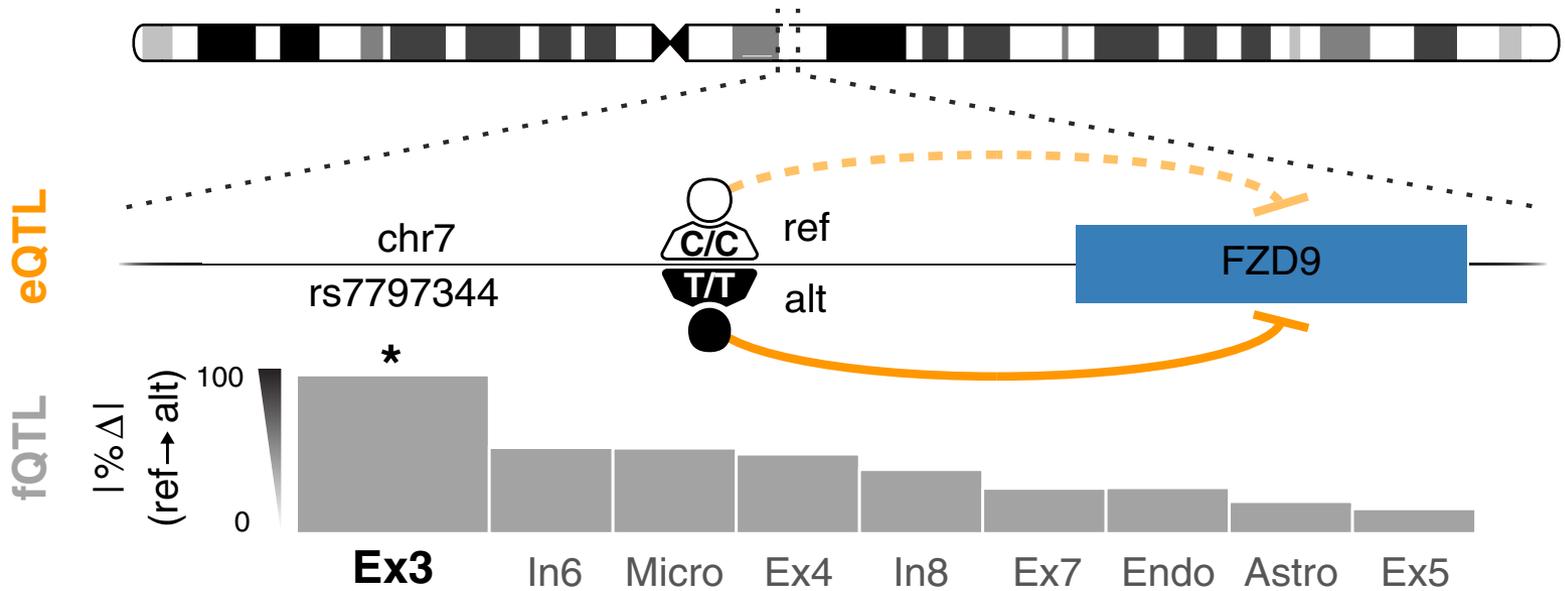
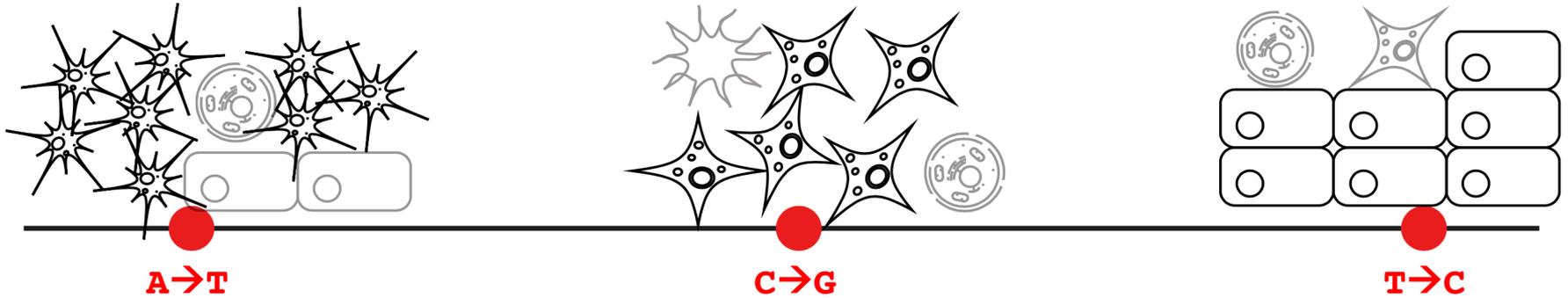
Sample 2: genotype CG



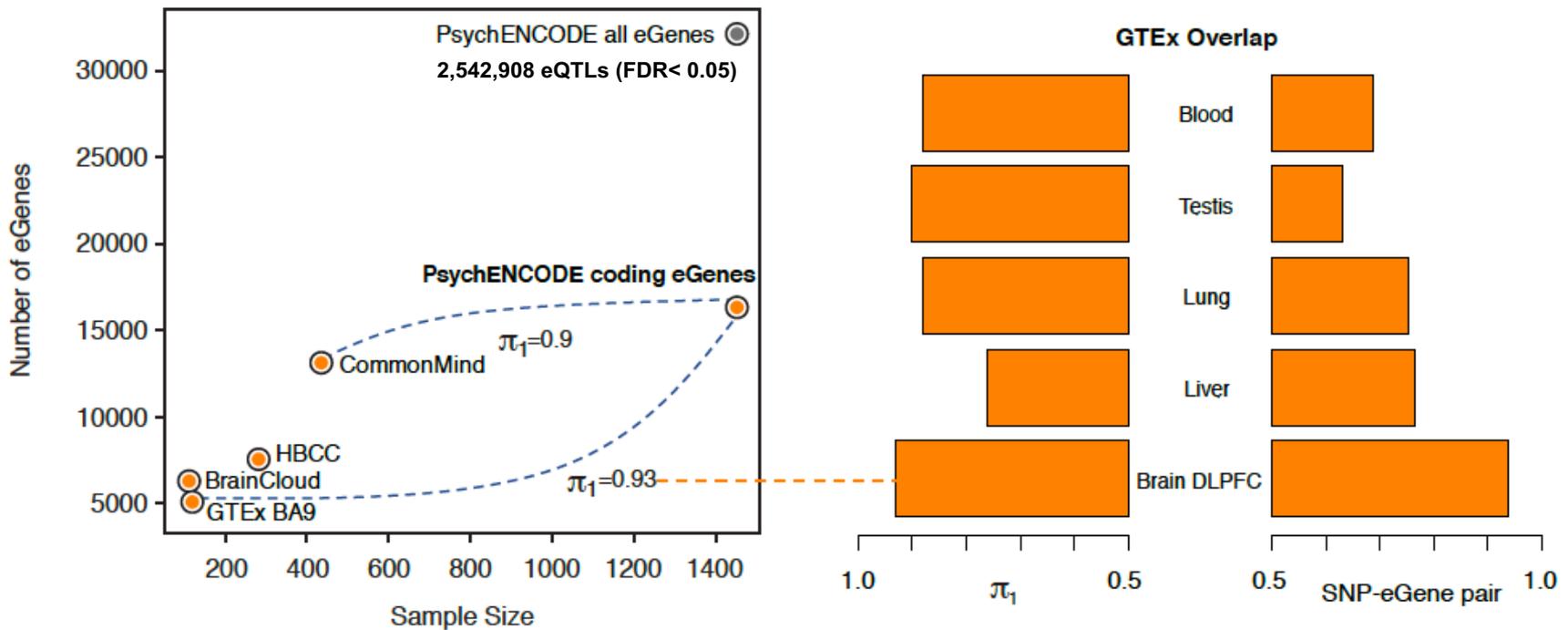
Chromatin (cQTL)



# Cell fraction QTLs (fQTLs)

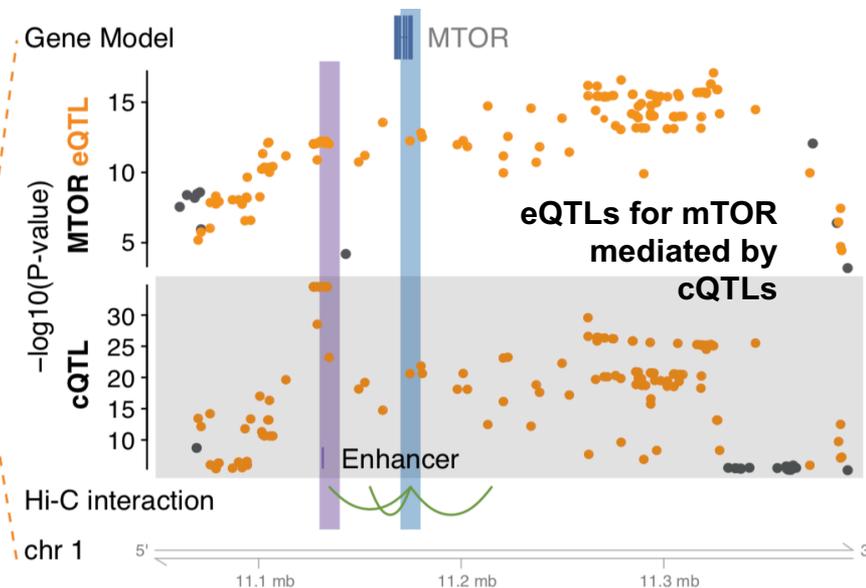
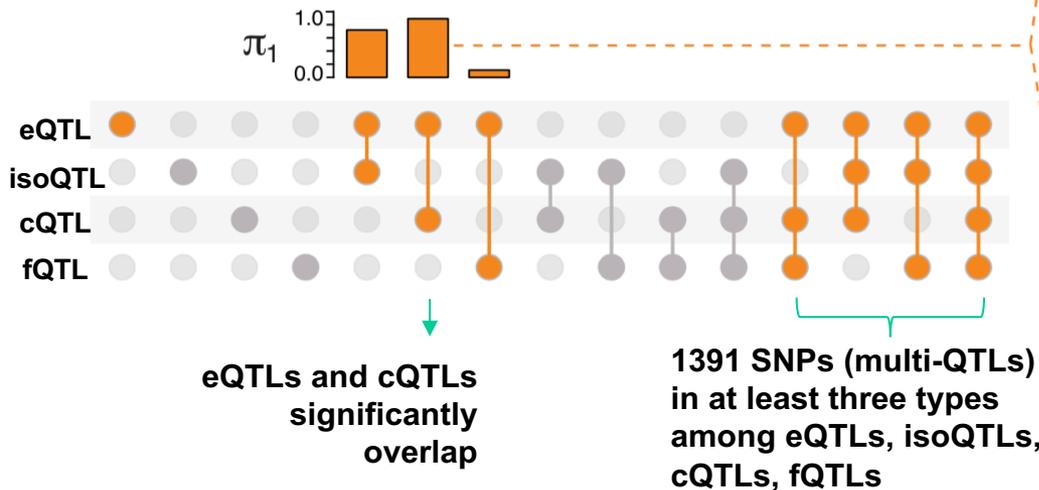


# Larger brain eQTL sets than previous studies, but strong overlap with them

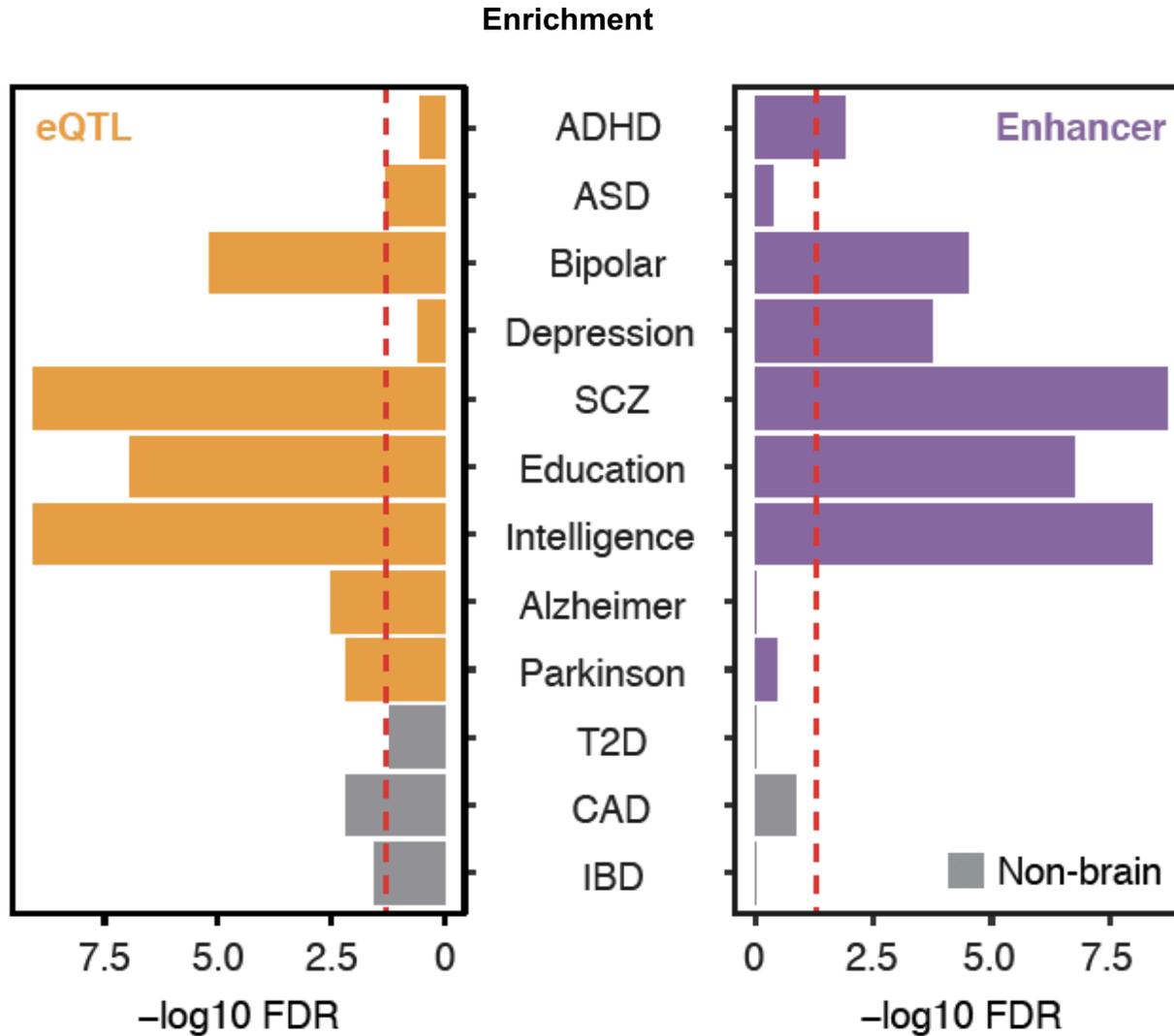


# multi-QTLs from overlapping different types of QTLs: cQTL, fQTL, eQTL & isoQTL

|        | Numbers of QTLs | eGenes Enhancers Cell types | SNPs      |
|--------|-----------------|-----------------------------|-----------|
| eQTL   | 2,542,908       | 32,944                      | 1,341,182 |
| isoQTL | 2,628,259       | 19,790                      | 1,052,939 |
| cQTL*  | 8,464           | 8,484                       | 7,983     |
| fQTL   | 4,199           | 9                           | 1,672     |

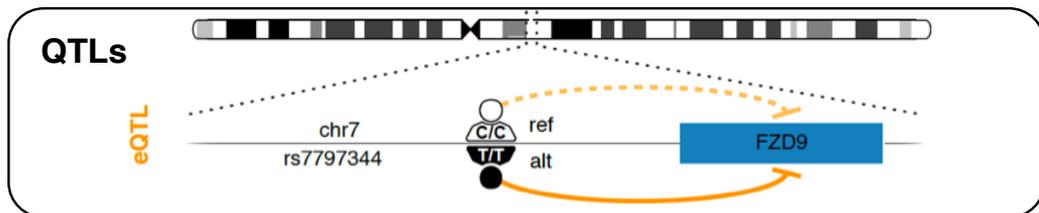
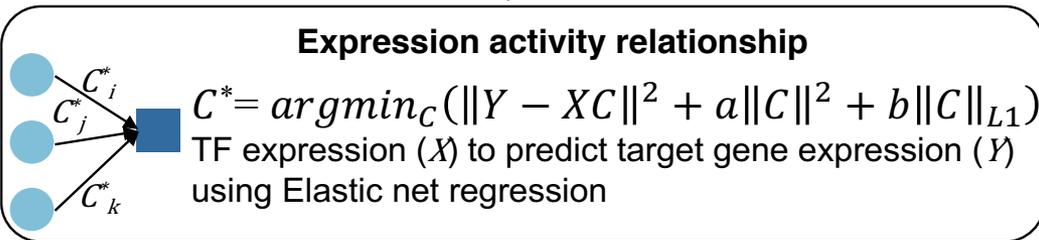
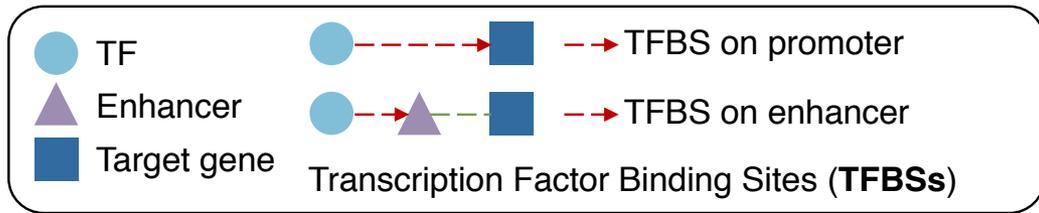
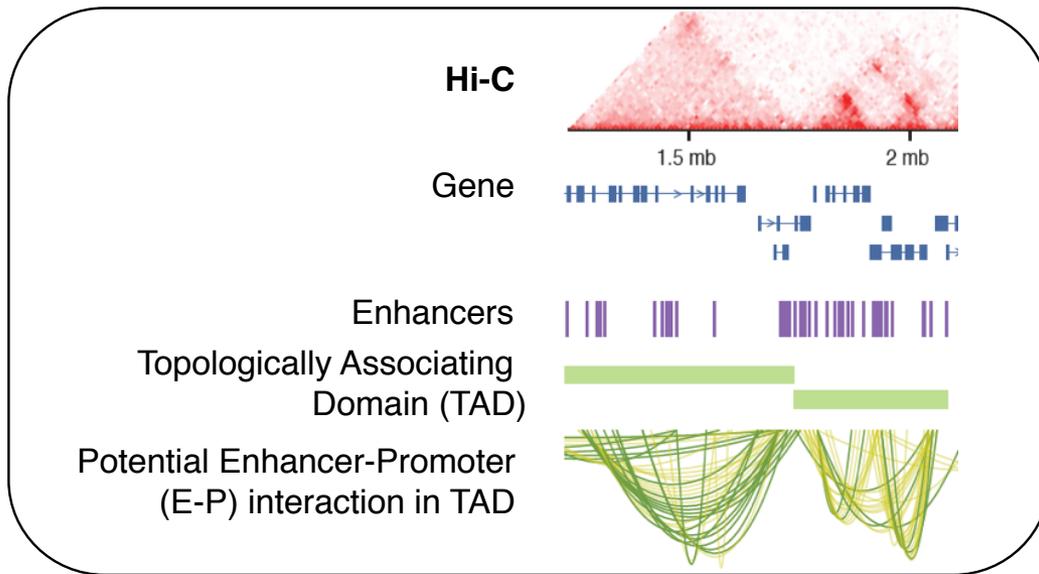


# Brain eQTLs and enhancers enriched with GWAS SNPs for brain disorders



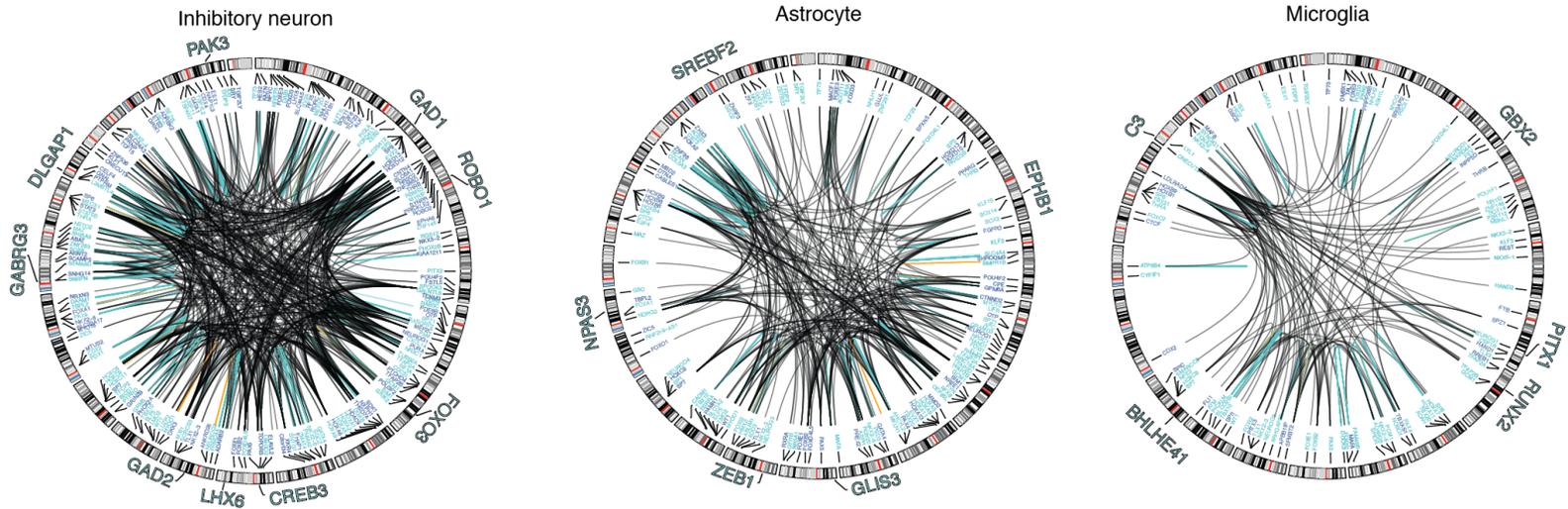
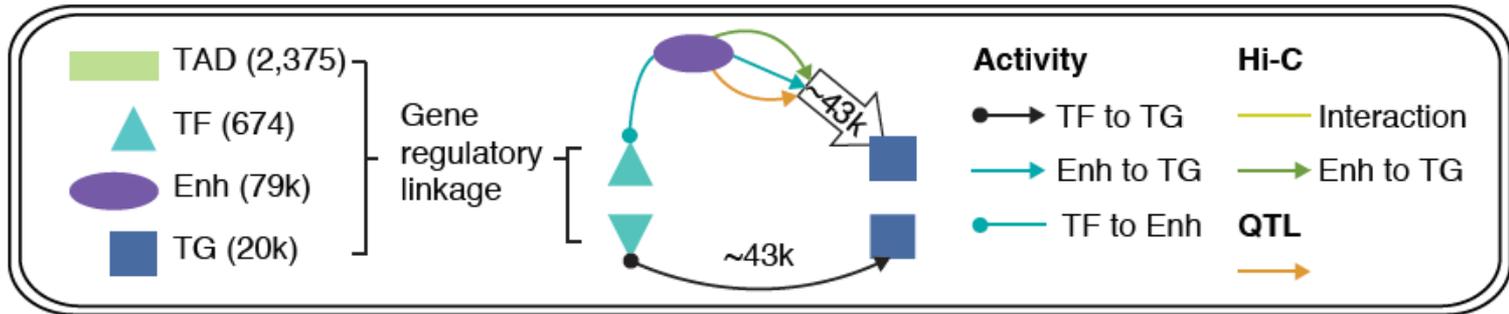
## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping



# Gene regulatory network inference from Hi-C, QTLs & Activity Correlations

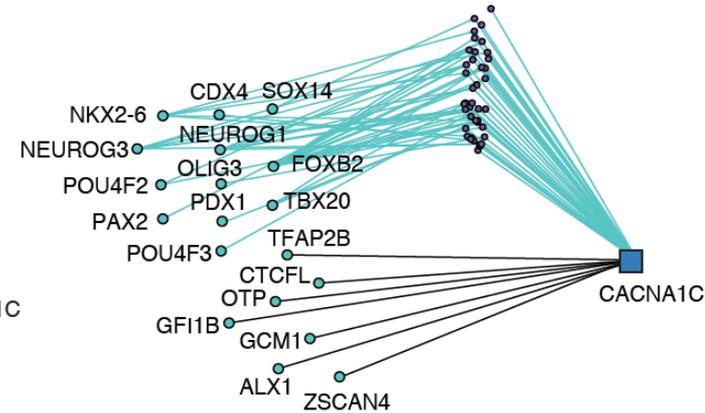
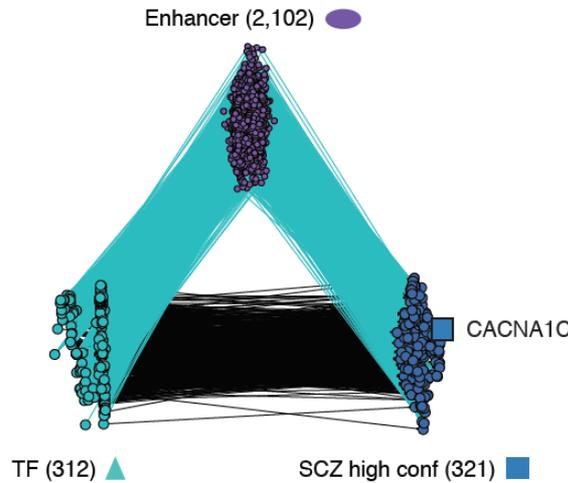
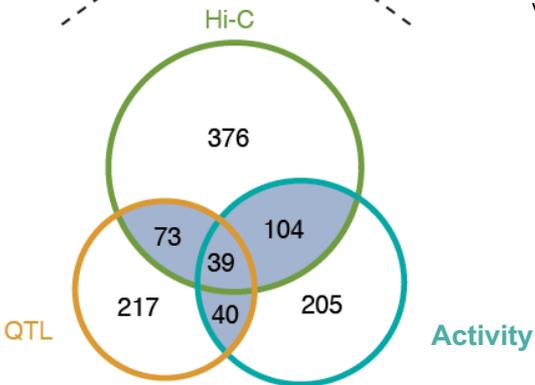
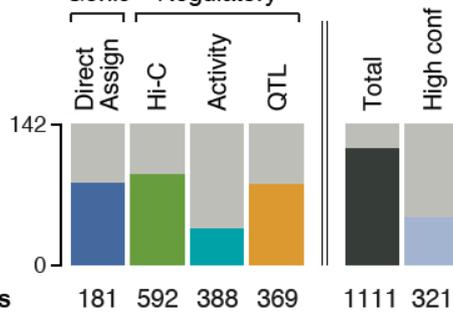
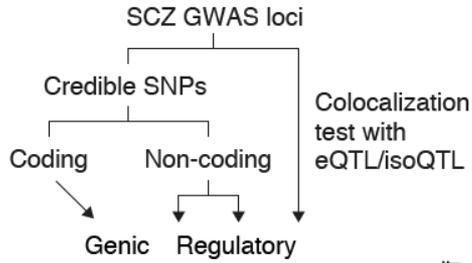
# Imputed gene regulatory network for the human brain



subnetworks targeting single cell marker genes

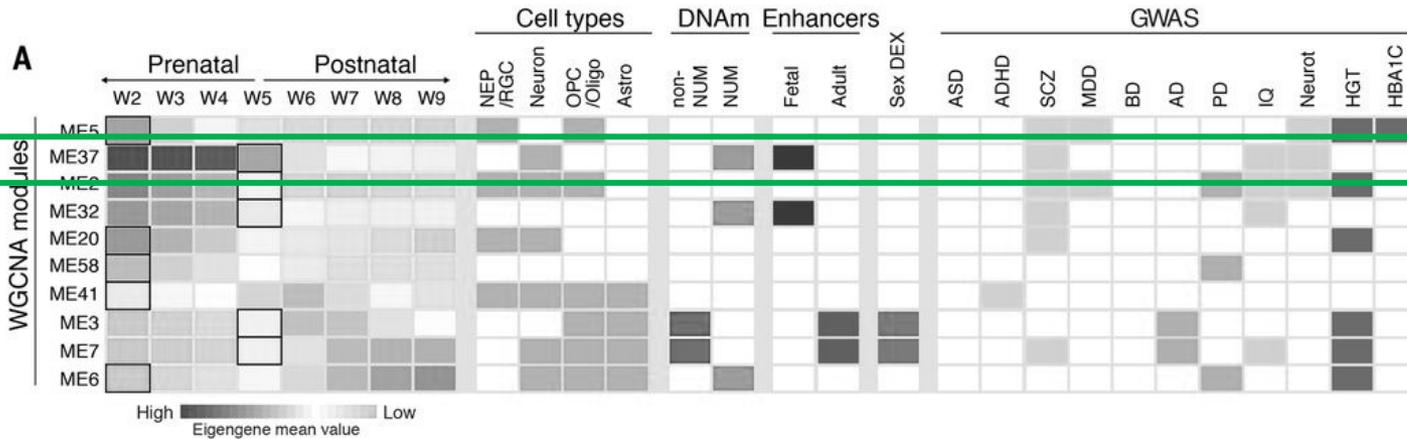
# 142

## Linking GWAS SNPs to disease genes using the regulatory network



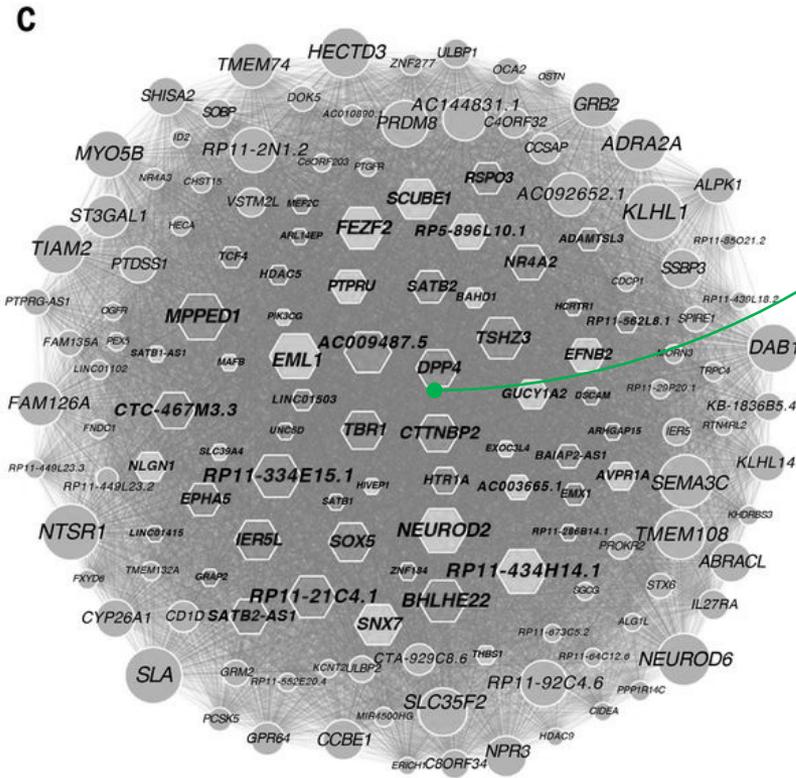
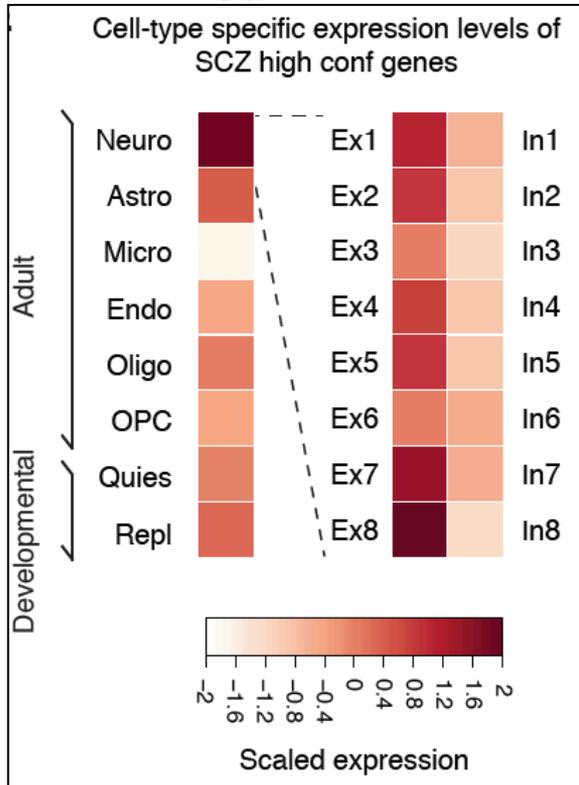
**321**  
 high-confident  
 SCZ genes

# Genes associated with SCZ enriched in specific neuronal cell types & co-expression modules, active prenatally



ME37

Gene co-expression network



Mingfeng Li et al. Science 2018

Wang, et al., Science, 2018

## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

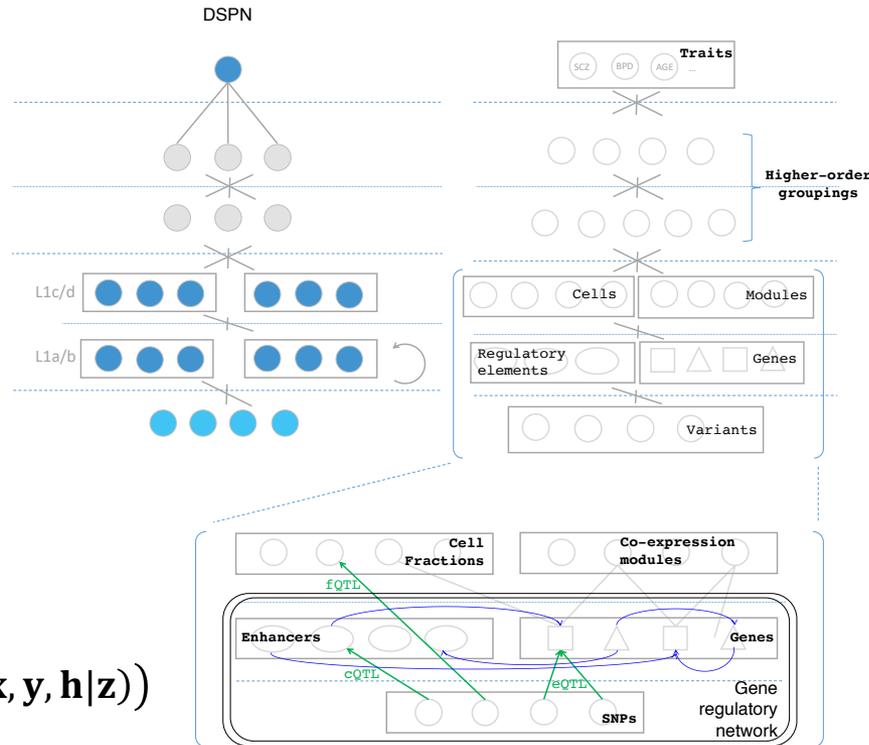
- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

# Deep Structured Phenotype Network (DSPN)

Gene regulatory network builds skeleton

Energy model:

$$p(\mathbf{x}, \mathbf{y}, \mathbf{h} | \mathbf{z}) \propto \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{h} | \mathbf{z}))$$



Boltzmann machine

**y**: phenotypes

**h**: hidden units (e.g., circuits)

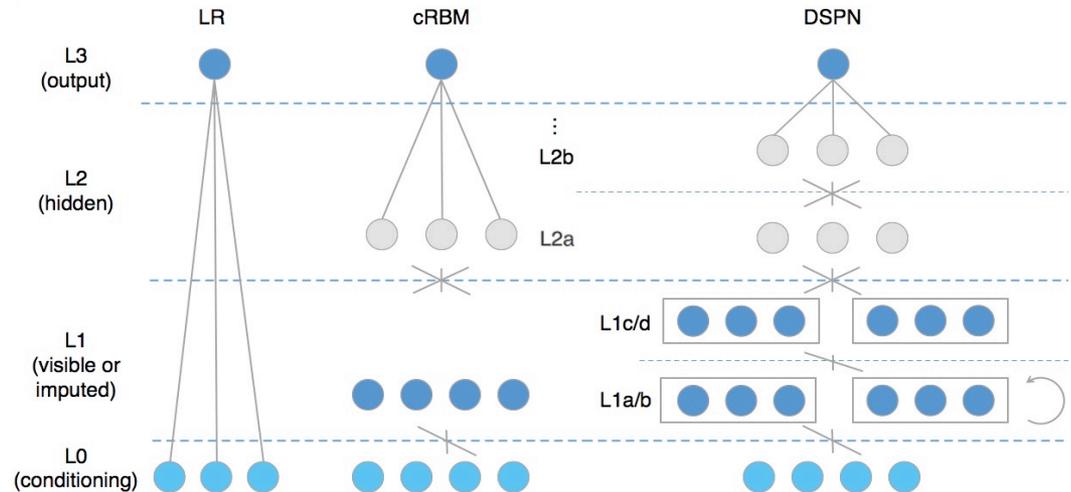
**x**: intermediate phenotypes (e.g., genes, enhancers)

**z**: genotypes (e.g., SNPs)

**W**: weights (e.g., regulatory network)

$$E(\mathbf{x}, \mathbf{y}, \mathbf{h} | \mathbf{z}) = -\mathbf{z}^T \mathbf{W}_1 \mathbf{x} - \mathbf{x}^T \mathbf{W}_2 \mathbf{x} - \mathbf{x}^T \mathbf{W}_3 \mathbf{h} - \mathbf{h}^T \mathbf{W}_4 \mathbf{h} - \mathbf{h}^T \mathbf{W}_5 \mathbf{y} - \text{Bias}$$

# DSPN improves brain disease prediction by adding deep layers



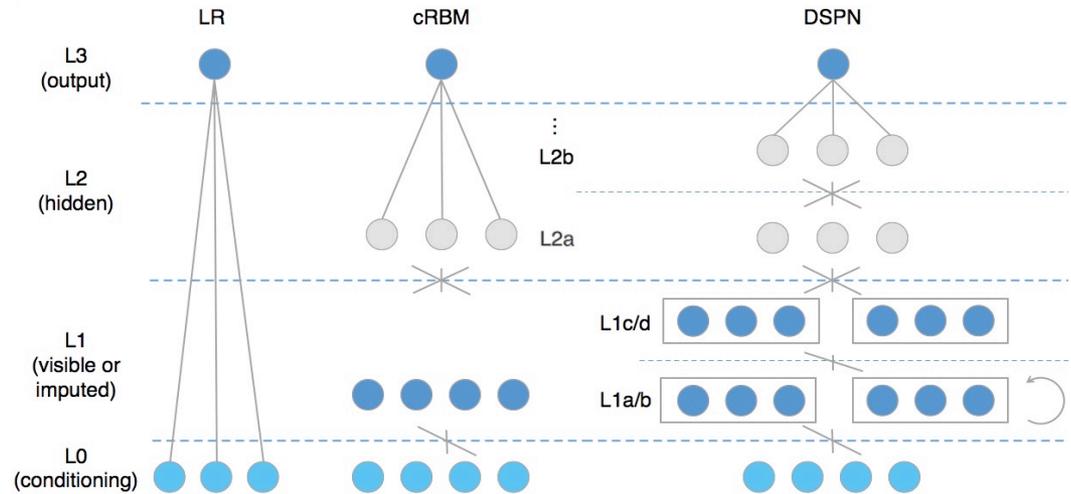
| Method                   | LR-genotype  | LR-transcriptome | cRBM  | DSPN-imputation | DSPN-full    |
|--------------------------|--------------|------------------|-------|-----------------|--------------|
| Schizophrenia            | <b>54.6%</b> | 63.0%            | 70.0% | 59.0%           | <b>73.6%</b> |
| Bipolar Disorder         | <b>56.7%</b> | 63.3%            | 71.1% | 67.2%           | <b>76.7%</b> |
| Autism Spectrum Disorder | <b>50.0%</b> | 51.7%            | 67.2% | 62.5%           | <b>68.3%</b> |



**X 6.0**

Accuracy = chance to correctly predict disease/health

# DSPN improves brain disease prediction by adding deep layers



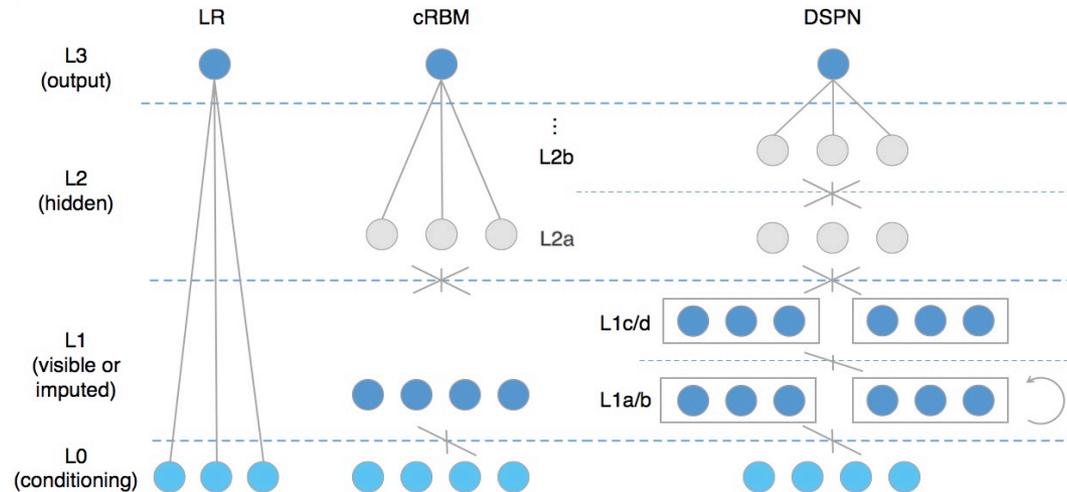
| Method                   | LR-genotype | LR-transcriptome | cRBM  | DSPN-imputation | DSPN-full |
|--------------------------|-------------|------------------|-------|-----------------|-----------|
| Schizophrenia            | 54.6%       | 63.0%            | 70.0% | 59.0%           | 73.6%     |
| Bipolar Disorder         | 56.7%       | 63.3%            | 71.1% | 67.2%           | 76.7%     |
| Autism Spectrum Disorder | 50.0%       | 51.7%            | 67.2% | 62.5%           | 68.3%     |



X 2.5

Accuracy = chance to correctly predict disease/health

# DSPN improves brain disease prediction by adding deep layers

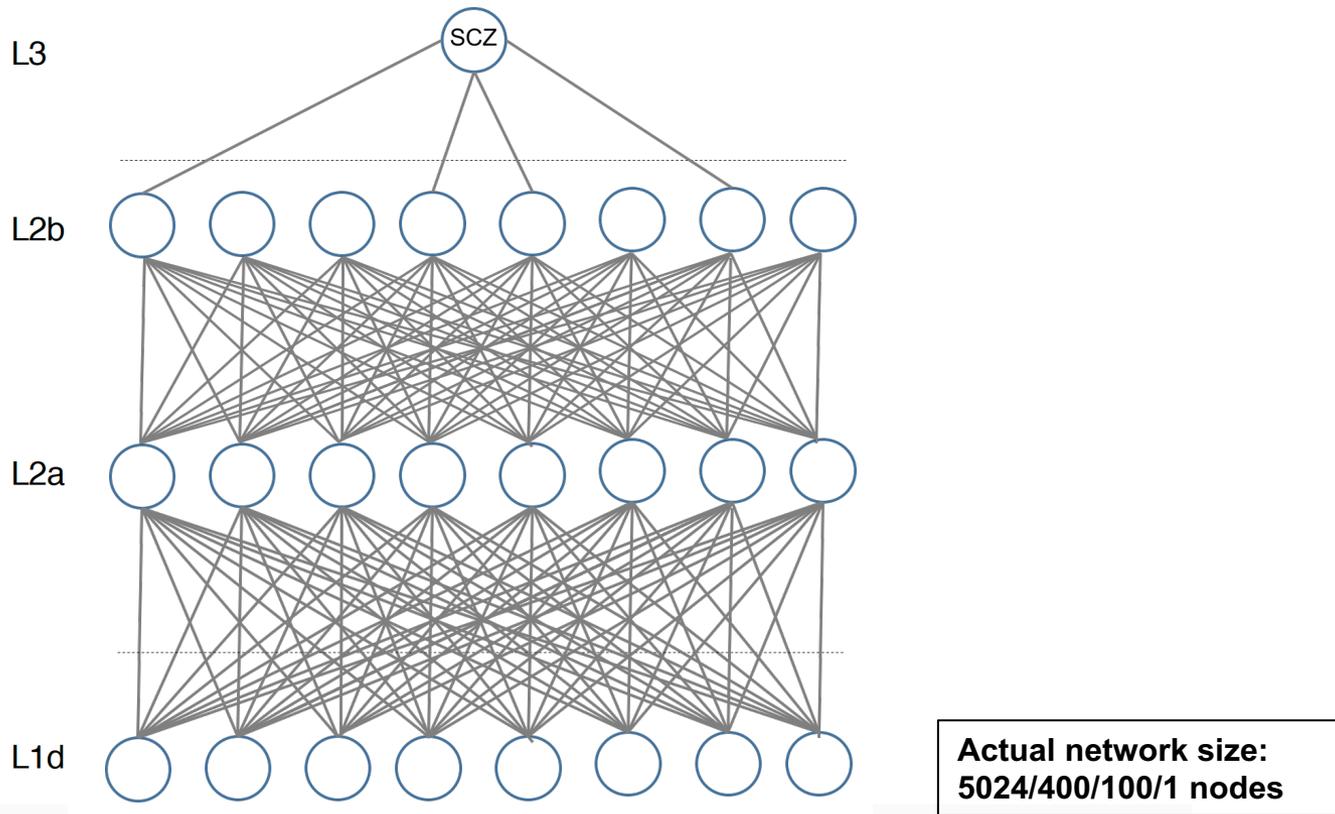


| Method                   | LR-genotype  | LR-transcriptome | cRBM  | DSPN-imputation | DSPN-full |
|--------------------------|--------------|------------------|-------|-----------------|-----------|
| Schizophrenia            | <b>54.6%</b> | 63.0%            | 70.0% | <b>59.0%</b>    | 73.6%     |
| Bipolar Disorder         | <b>56.7%</b> | 63.3%            | 71.1% | <b>67.2%</b>    | 76.7%     |
| Autism Spectrum Disorder | <b>50.0%</b> | 51.7%            | 67.2% | <b>62.5%</b>    | 68.3%     |

**X 3.1**

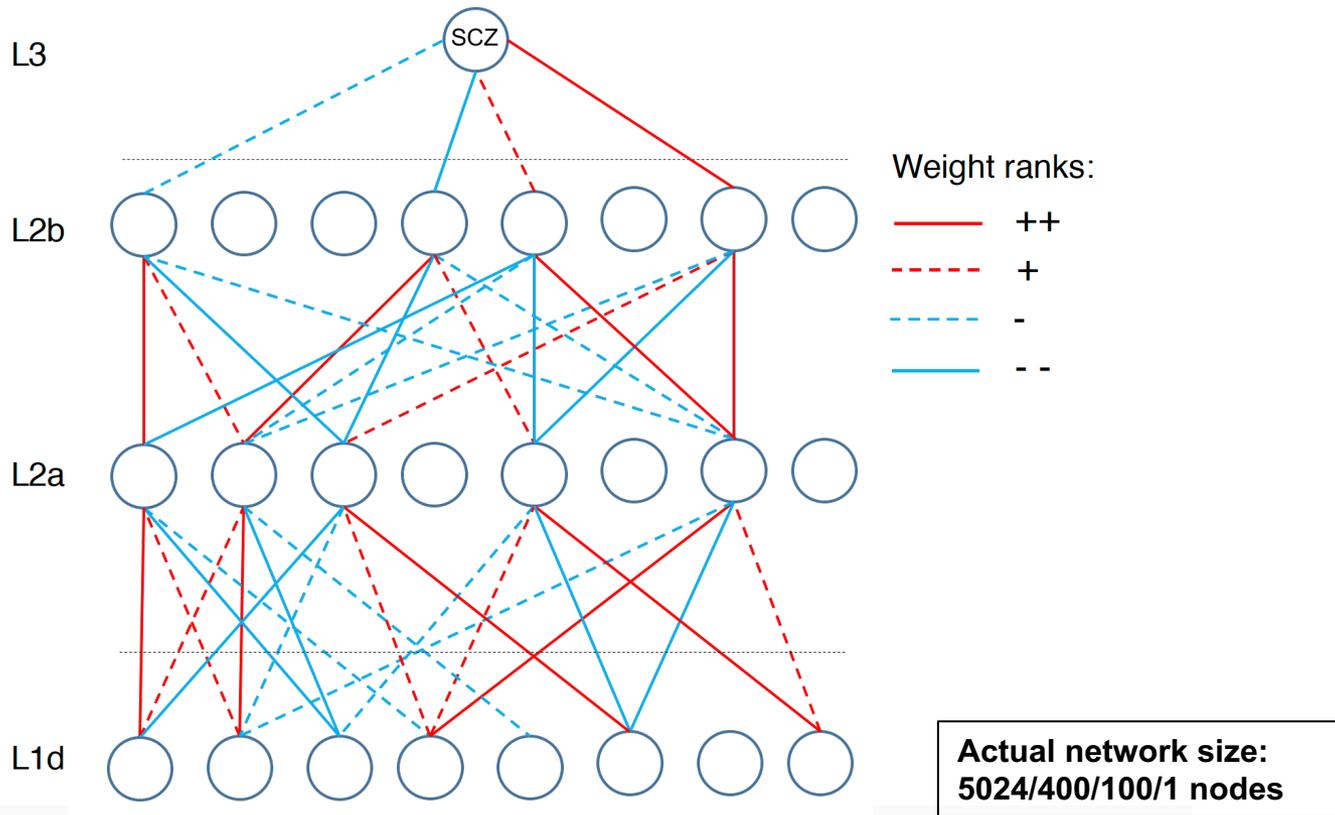
Accuracy = chance to correctly predict disease/health

# Multilevel Network Interpretation



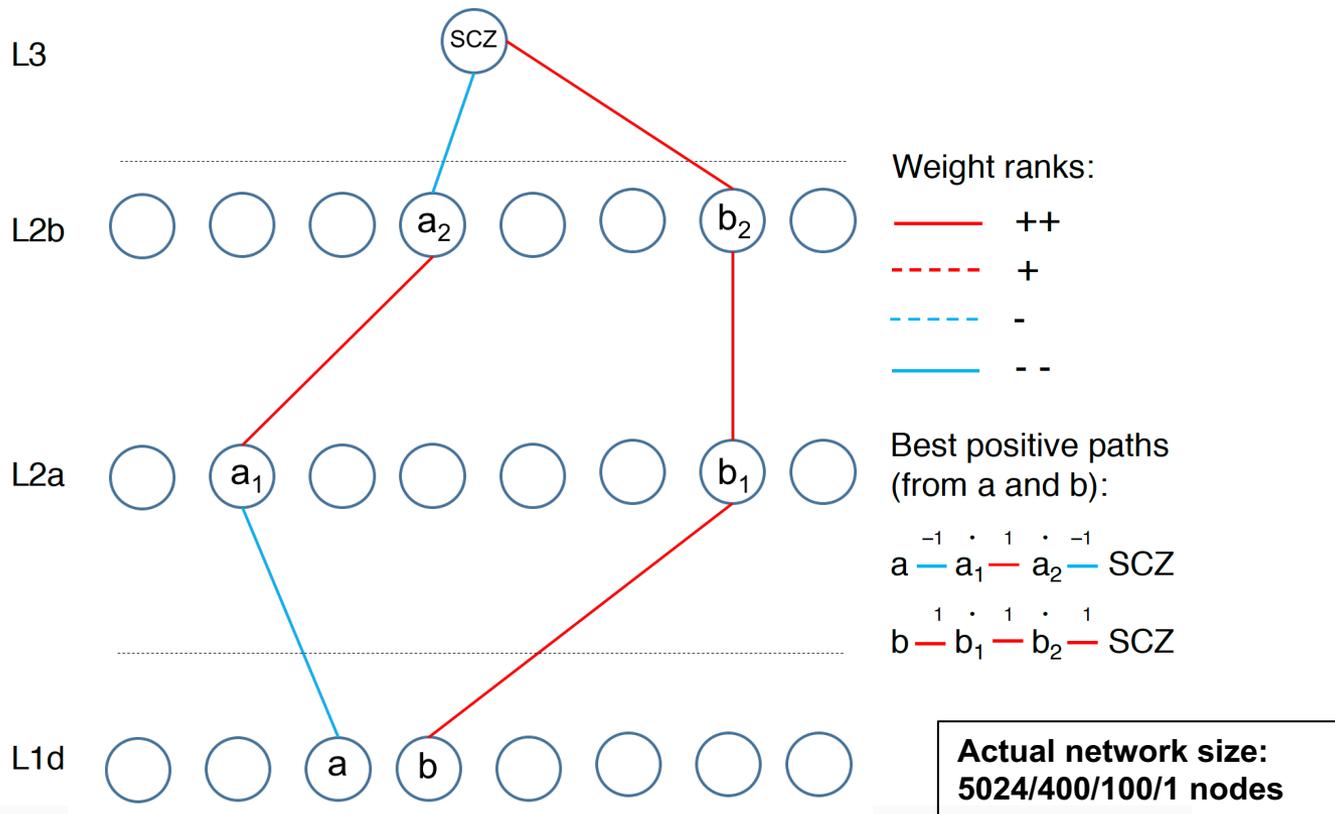
- Start with a fully connected trained network

# Multilevel Network Interpretation



- Start with a fully connected trained network
- Sparsify network using edges with largest absolute weights (+/-)

# Multilevel Network Interpretation



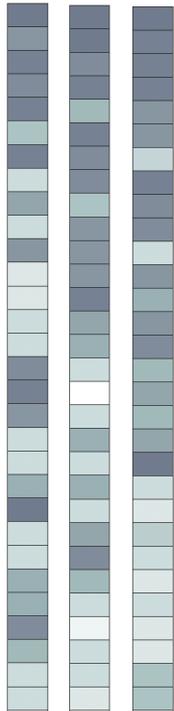
- Start with a fully connected trained network
- Sparsify network using edges with largest absolute weights (+/-)
- Extract 'best positive paths' to each prioritized module (e.g. a-a<sub>1</sub>-a<sub>2</sub>-SCZ) by summing weights and multiplying signs

# DSPN discovers enriched pathways and linkages to genetic variation

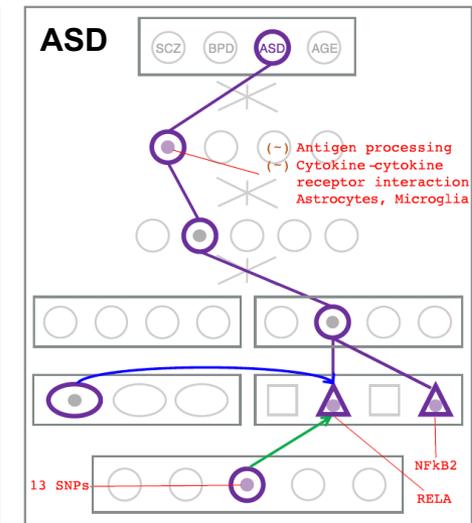
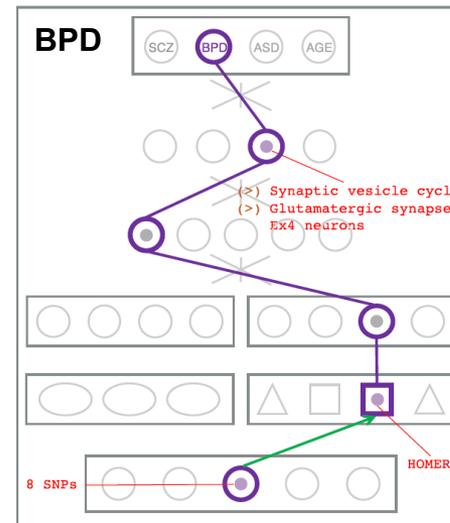
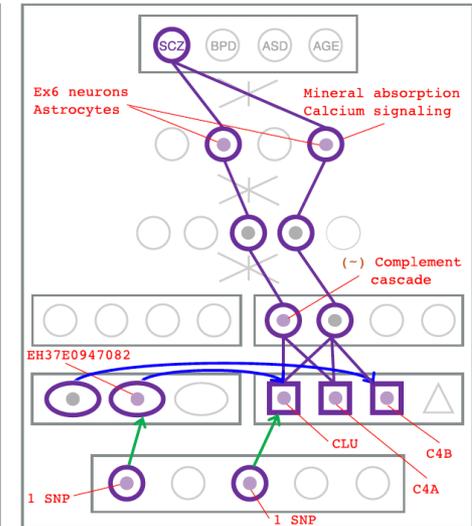
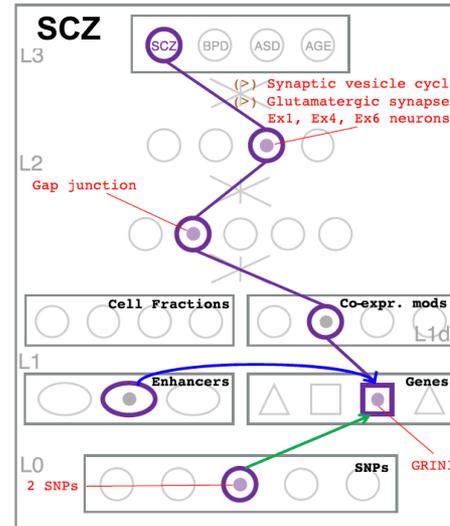
## Cross-disorder MOD/HOG enrichment ranking



SCZ BPD ASD



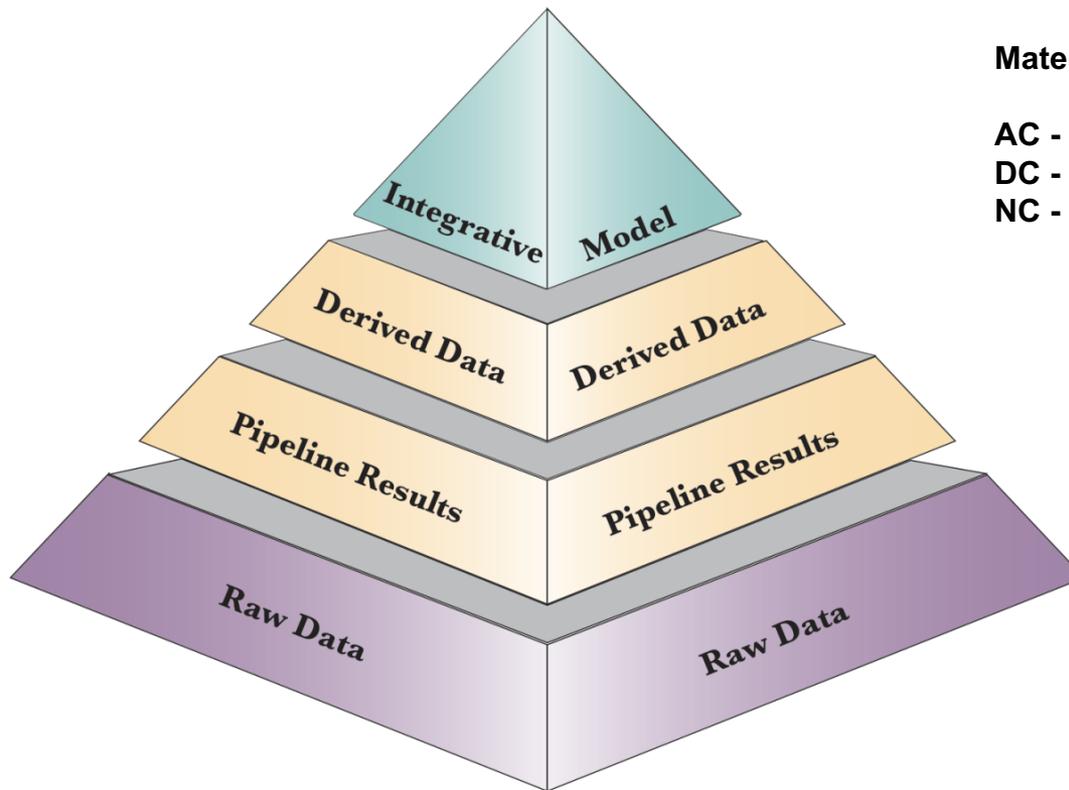
- (\*) Spliceosome / RNA splicing
- (>) Synaptic vesicle cycle
- (~) Antigen proc. and presentation
- Vesicle localization
- Proteasome
- (\*) mRNA processing
- Chromatin modification
- (#) Oxidative phosphorylation
- Retrograde endocannabinoid sig.
- (>) Chemical synaptic transmission
- Peptidyl-lysine modification
- Endocytosis
- Ubiquitin mediated proteolysis
- (>) Anterograde trans-synaptic sig.
- (\*) mRNA transport
- Phosphatidylinositol signaling
- Hippo signaling pathway
- (~) Staph./ Epstein-Barr virus inf.
- (>) Synaptic signaling
- Autophagy
- (>) Dop./GABA/Glutamatergic synapse
- (>) Calcium signaling
- (>) Endocrine calcium reabsorption
- (\*) RNA degradation / transport
- (#) Ribosome
- Neuron projection morphogenesis
- (~) Fc receptor signaling pathway
- cGMP-PKG signaling pathway
- (~) mTOR signaling pathway
- (~) Cytokine-cytokine receptor int.



## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

## Phase 1 PsychENCODE capstone resource: Layers of distributed information



**Material in the 3 capstones:**

**AC - Wang et al. ('18)**

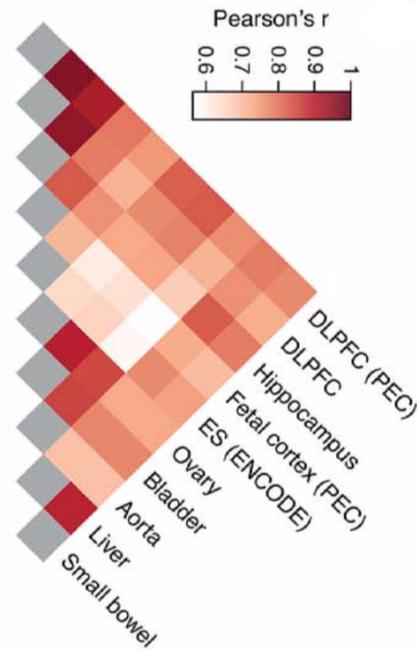
**DC - Li et al. ('18)**

**NC - Gandal et al. ('18)**

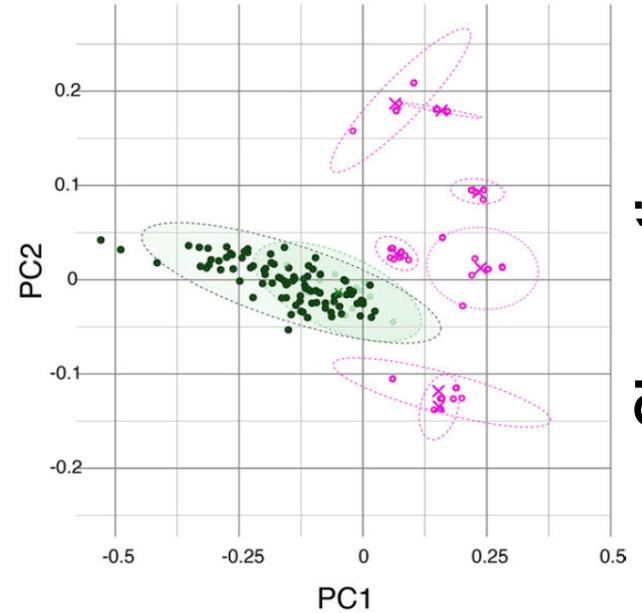
**Resource.psychencode.org**  
**Development.psychencode.org**

# Cross tissue variation in Chromatin & Expression

Placing the **Brain** in context of all other **Body Tissues**

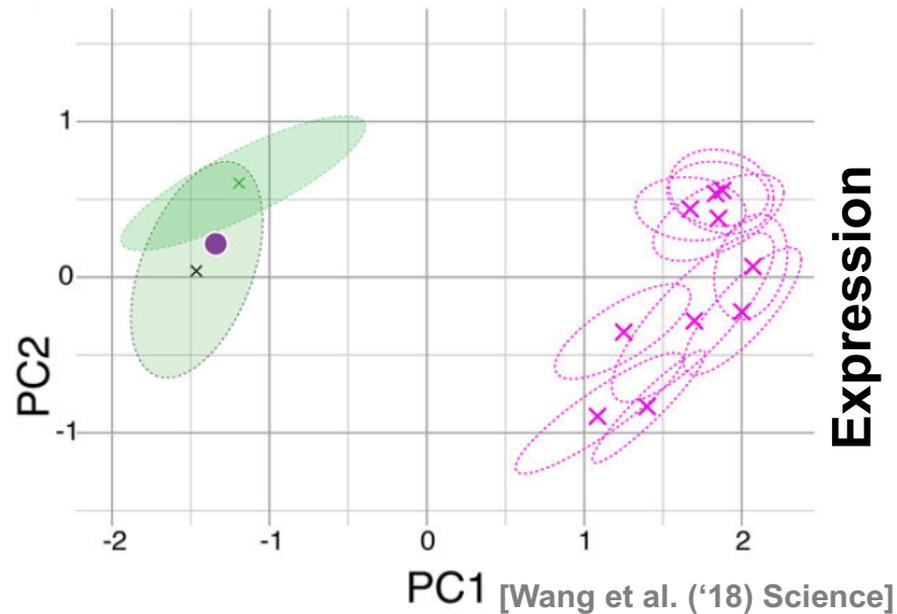
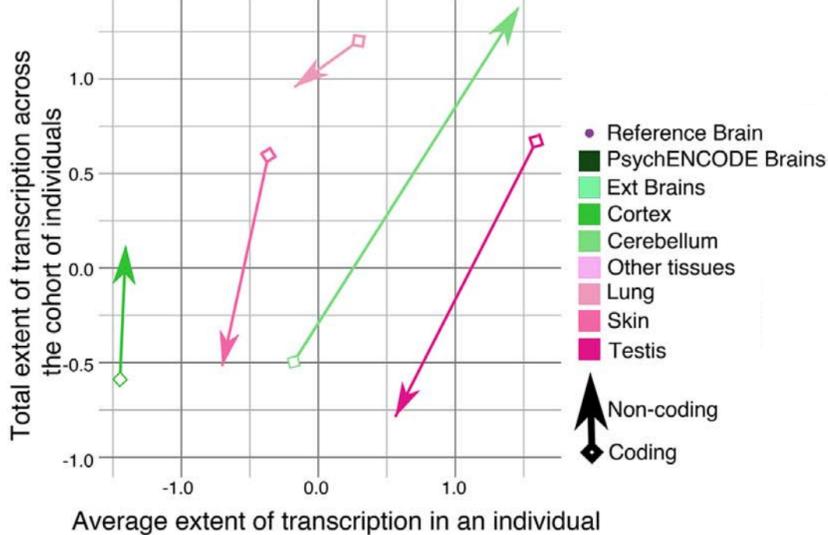


Hi-C



Chromatin

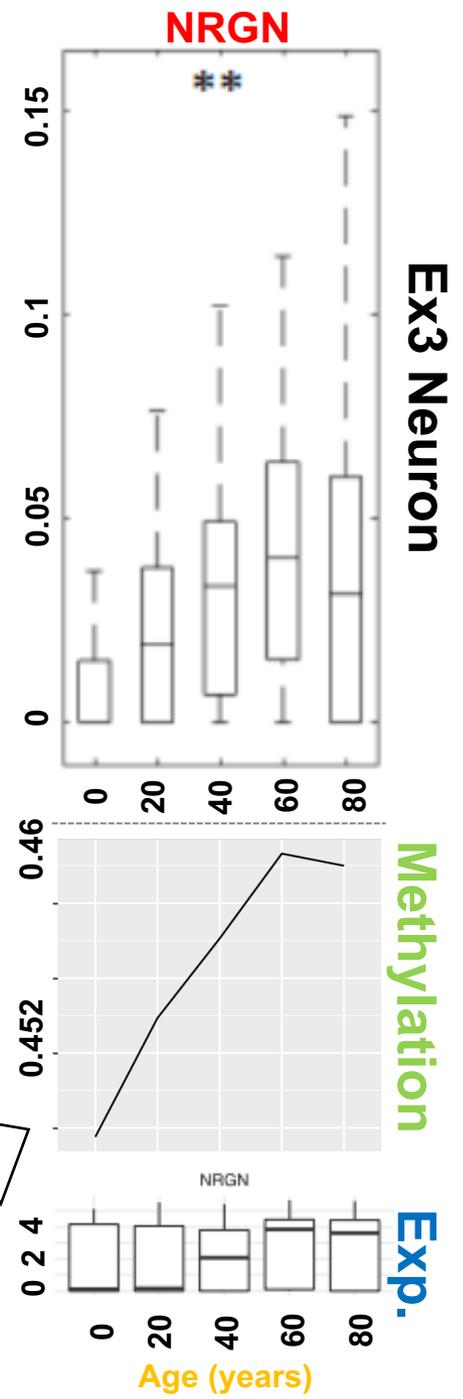
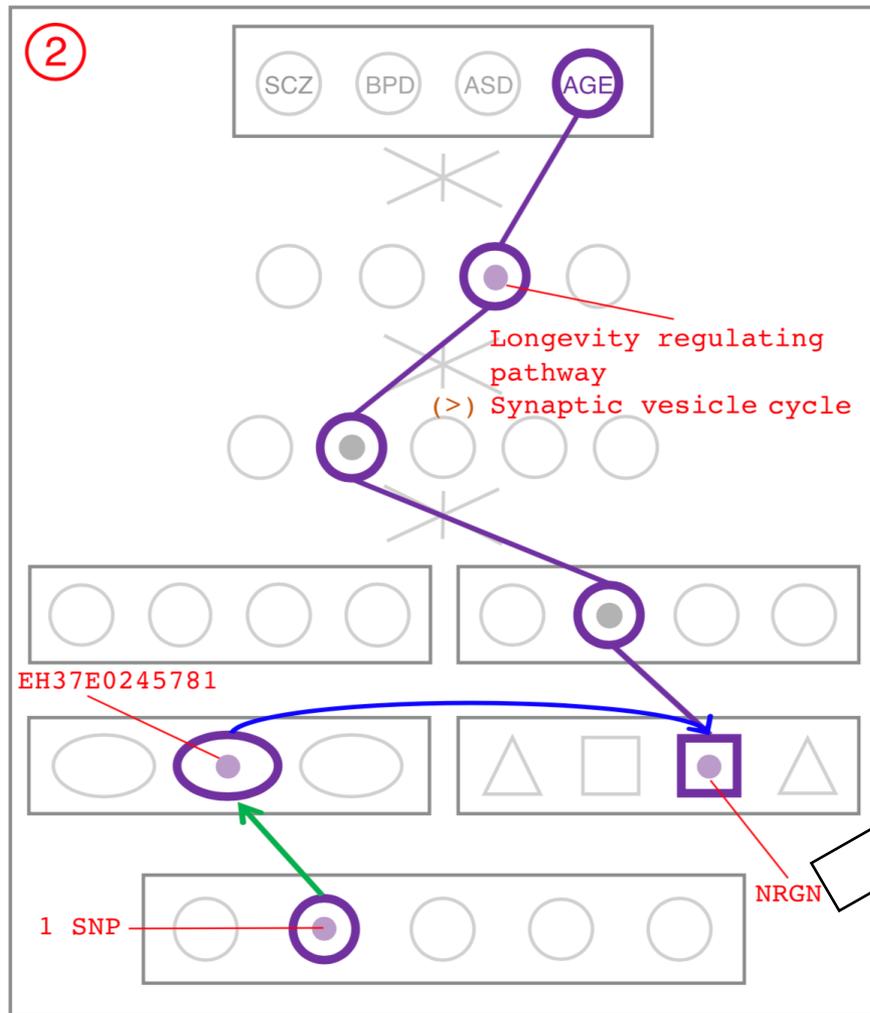
**Transcriptome diversity** increases in the non-coding portion of the **brain genome** while decreases in **other tissues**



Expression

# NRGN has variable expression over age and is in Synaptic vesicle cycle pathway is enriched in SCZ, BPD, ASD

**NRGN** is a gene associated with the **Synaptic vesicle pathway** and **NRGN expression** and **methylation** is correlated with **Age**



## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

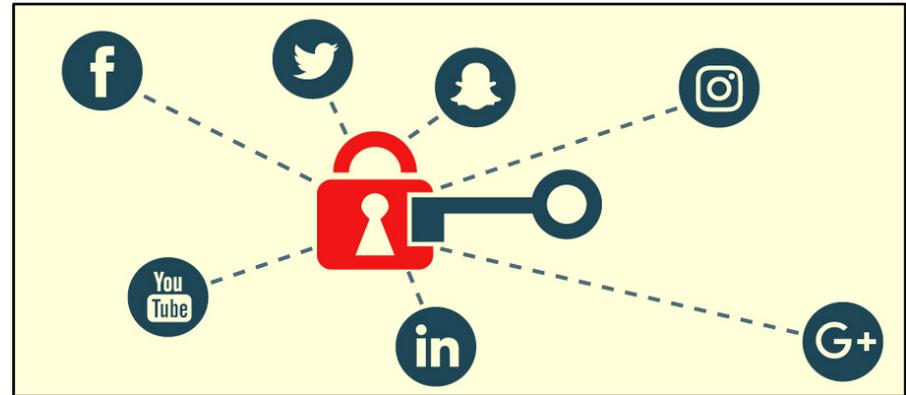
## 2-sided nature of functional genomics data: Analysis can be very **General/Public** or **Individual/Private**



- **General quantifications** related to overall aspects of a condition – ie gene activity as a function of:
  - Developmental stage, Evolutionary relationships, Cell-type, Disease
- **Above are not tied to an individual's genotype. However, data is derived from individuals & tagged with their genotypes**
- (Note, a few calculations aim to use explicitly genotype to derive general relations related to sequence variation & gene expression - eg allelic activity)

# Privacy: Does Genomics has similar "Big Data" Dilemma as in the Rest of Society?

- We confront privacy risks every day we access the internet (e.g., social media, e-commerce).
- Sharing & "peer-production" is central to success of many new ventures, with analogous risks to genomics
  - **EG web search**: Large-scale mining essential



## Genetic Exceptionalism :

The Genome is very fundamental data, potentially very revealing about one's identity & characteristics

**Personal Genomic info. essentially meaningless currently but will it be in 20 yrs? 50 yrs?**

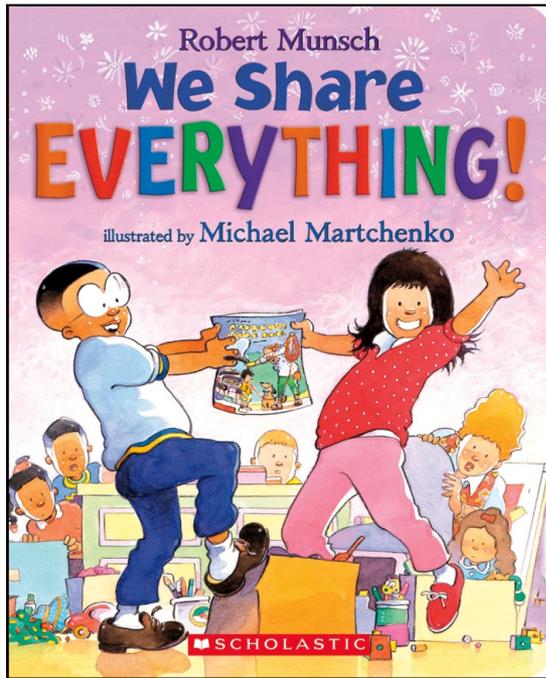
Genomic sequence very revealing about one's children. Is true consent possible?

Once put on the web it can't be taken back

**Ethically challenged** history of genetics

Ownership of the data & what consent means (Hela)

Could your genetic data give rise to a product line?

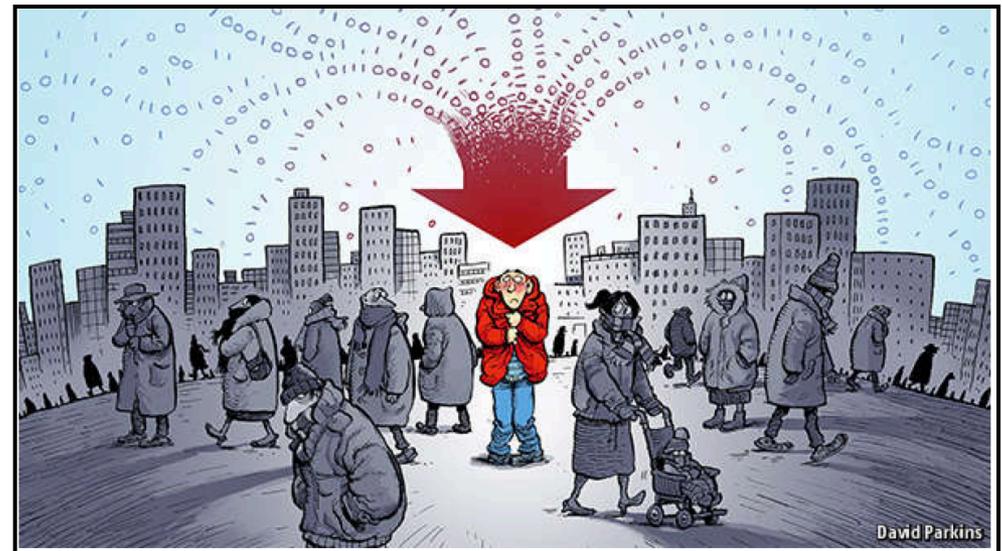


## The Other Side of the Coin for Genomics: Why we should share

- Sharing helps **speed research**
  - Large-scale mining of this information is important for medical research
  - Statistical power
  - Privacy is cumbersome, particularly for big data

### The Dilemma

- The individual (harmed?) v the collective (benefits)
  - But do sick patients care about their privacy?
- How to balance risks v rewards
  - Quantification



[Economist, 15 Aug '15]

[Yale Law Roundtable ('10). Comp. in Sci. & Eng. 12:8; D Greenbaum & M Gerstein ('09). Am. J. Bioethics; D Greenbaum & M Gerstein ('10). SF Chronicle, May 2, Page E-4; Greenbaum et al. PLOS CB ('11)]

# Current Social & Technical Solutions: The quandary where are now

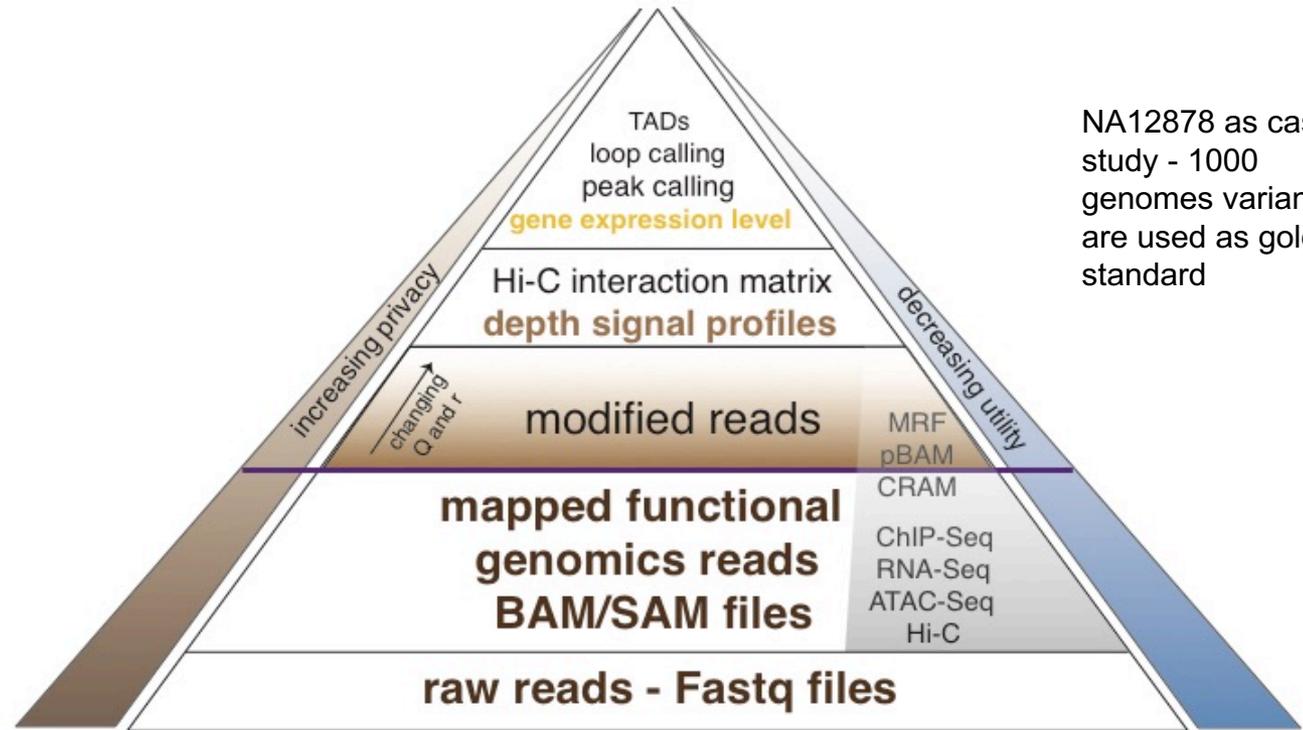
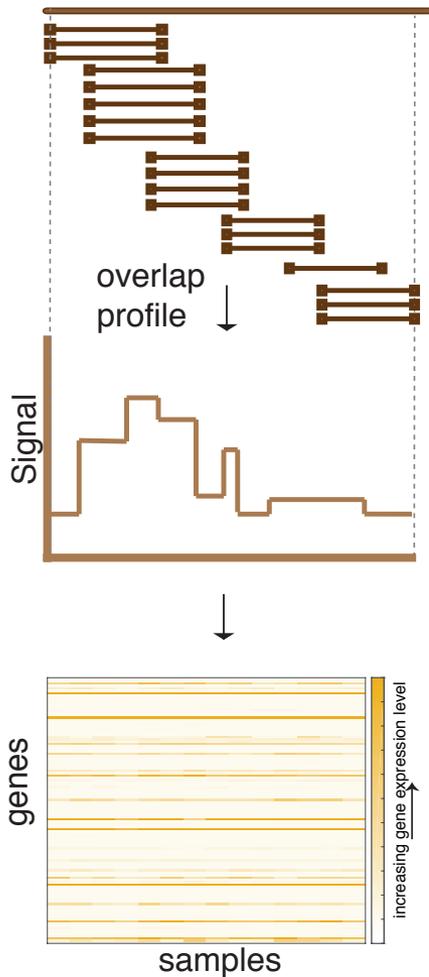
- **Closed Data** Approach
  - Consents
  - “Protected” distribution via dbGAP
  - Local computes on secure computer
- Issues with Closed Data
  - Non-uniformity of consents & paperwork
    - Different, confusing int’l norms
  - Computer security is burdensome
  - Many schemes get “hacked” .
  - **Tricky aspects of high-dimensional data** (leakage & ease of creating quasi-identifiers)
- **Open Data**
  - Genomic “test pilots” (ala PGP)?
    - Sports stars & celebrities?
  - Some public data & data donation is helpful but is this a realistic solution for an unbiased sample of ~1M



# Strawman Hybrid **Social** & **Tech** Proposed Solution?

- Fundamentally, researchers have to keep genetic secrets.
  - **Need for an (international) legal framework**
  - Genetic Licensure & training for individuals (similar to medical license, drivers license)
- Technology to make things easier
  - Cloud computing & enclaves (eg solution of Genomics England)
- Technological barriers shouldn't create a social incentive for “hacking”
- **Quantifying Leakage & allowing a small amounts of it**
- **Careful separation & coupling of private & public data**
  - **Lightweight, freely accessible secondary datasets coupled to underlying variants**
  - Selection of stub & "test pilot" datasets for benchmarking
  - Develop programs on public stubs on your laptop, then move the program to the cloud for private production run

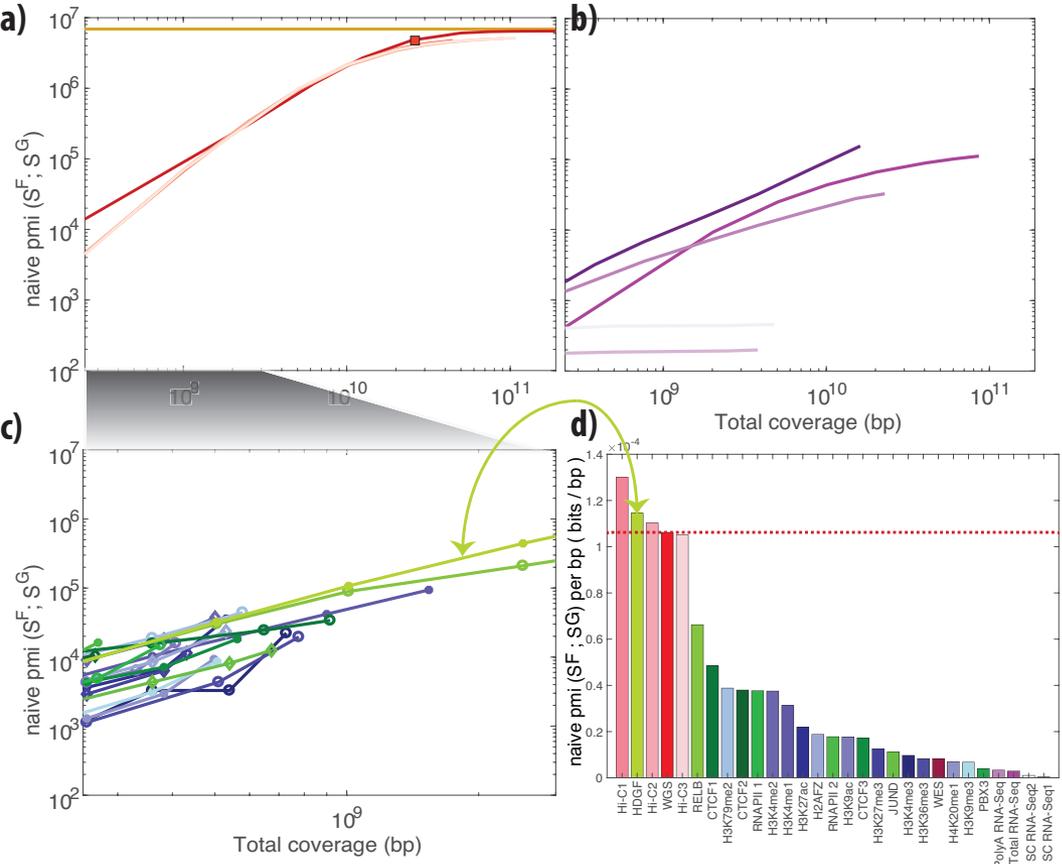
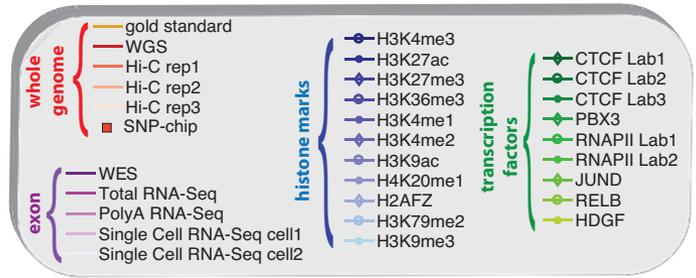
# Functional genomics data comes with a great deal of sequencing; We can quantify amount of leakage at every step of the data summarization process.



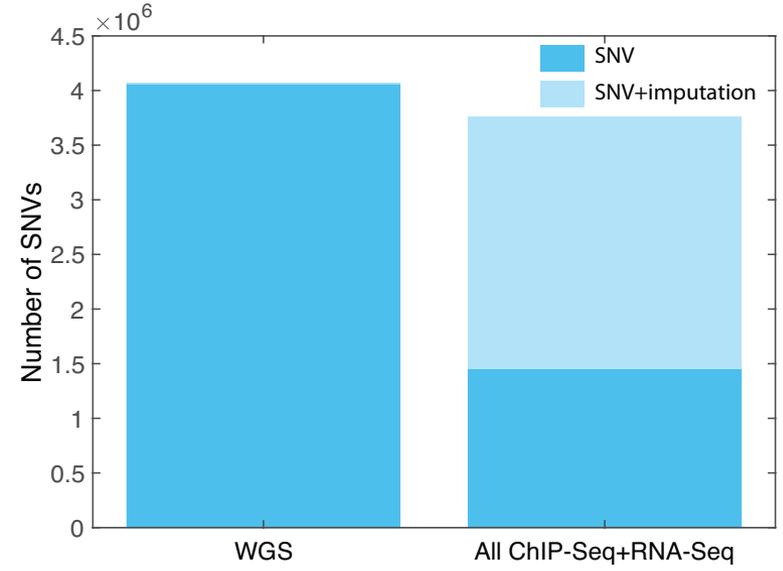
NA12878 as case study - 1000 genomes variants are used as gold standard

| Leakage Source                     | Leaking Variants | # of potential variants | Average leakage per variant (bits) | Maximum leakage per variant (bits) | # of accessible variants | Total leakage (bits) |
|------------------------------------|------------------|-------------------------|------------------------------------|------------------------------------|--------------------------|----------------------|
| Raw reads                          | Exonic variants  | 2,682,417               | 0.10 ± 0.28                        | 9.88 ± 2.12                        | 246,893                  | 24,689               |
| Modified reads<br>Q = {indels}     | Exonic SNVs      | 2,607,969               | 0.09 ± 0.27                        | 9.95 ± 2.02                        | 231,031                  | 207,92               |
| Modified reads<br>Q = {mismatches} | Exonic indels    | 51,408                  | 0.33 ± 0.47                        | 7.64 ± 2.42                        | 15,862                   | 5234                 |
| Signal profiles                    | Exonic deletions | 48,019                  | 0.29 ± 0.45                        | 7.97 ± 2.42                        | 1,067                    | 298                  |
| Gene expression quantification     | eQTLs            | 3,175                   | 1.19 ± 0.36                        | 4.00 ± 1.92                        | 158                      | 188                  |

- How much information, for example, do RNA-Seq reads (or ChIP-Seq) reads contain? Does that information enough to identify individuals?

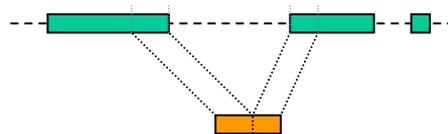
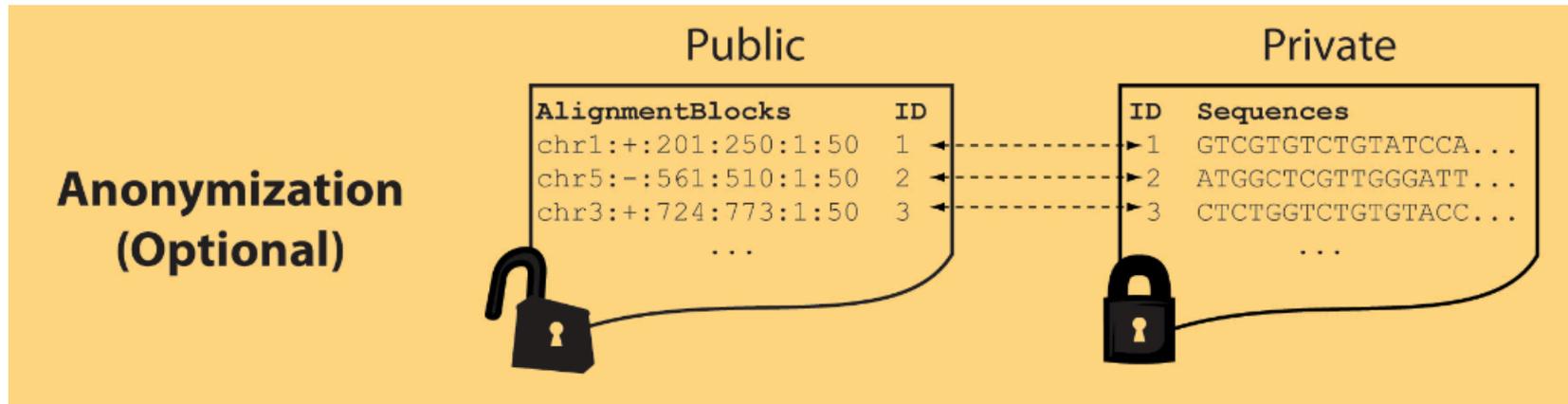


- It might seem like we don't infer much information from single ChIP-Seq and RNA-Seq experiments compared to WGS
  - However putting 10 different ChIP-Seq experiments and RNA-Seq together with imputation provides a great deal of information about the individual



# Light-weight formats to Hide Most of the Read Data (Signal Tracks)

- Some lightweight format clearly separate public & private info., aiding exchange
- Files become much smaller. Similar to CRAM
- Distinction between formats to compute on and those to archive with – become sharper with big data



**Mapping coordinates without variants (MRF)**

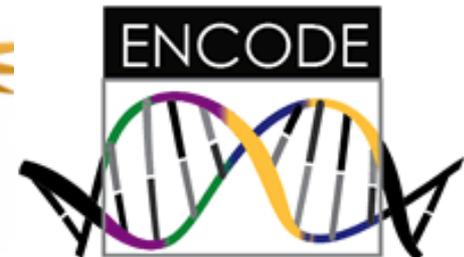
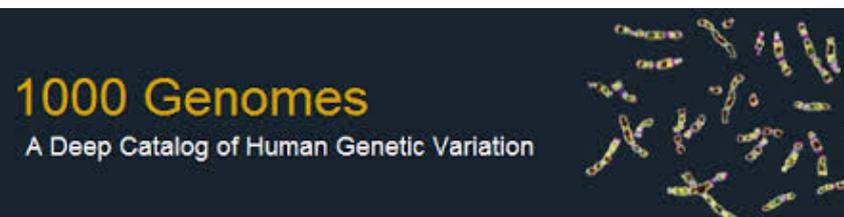
**Reads (linked via ID, 10X larger than mapping coord.)**

## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

# Representative Functional Genomics, Genotype, eQTL Datasets

- Genotypes are available from the 1000 Genomes Project
- mRNA sequencing for 462 individuals from gEUVADIS and ENCODE
  - Publicly available quantification for protein coding genes
- Functional genomics data (ChIP-Seq, RNA-Seq, Hi-C) available from ENCODE
- Approximately 3,000 cis-eQTL (FDR<0.05)



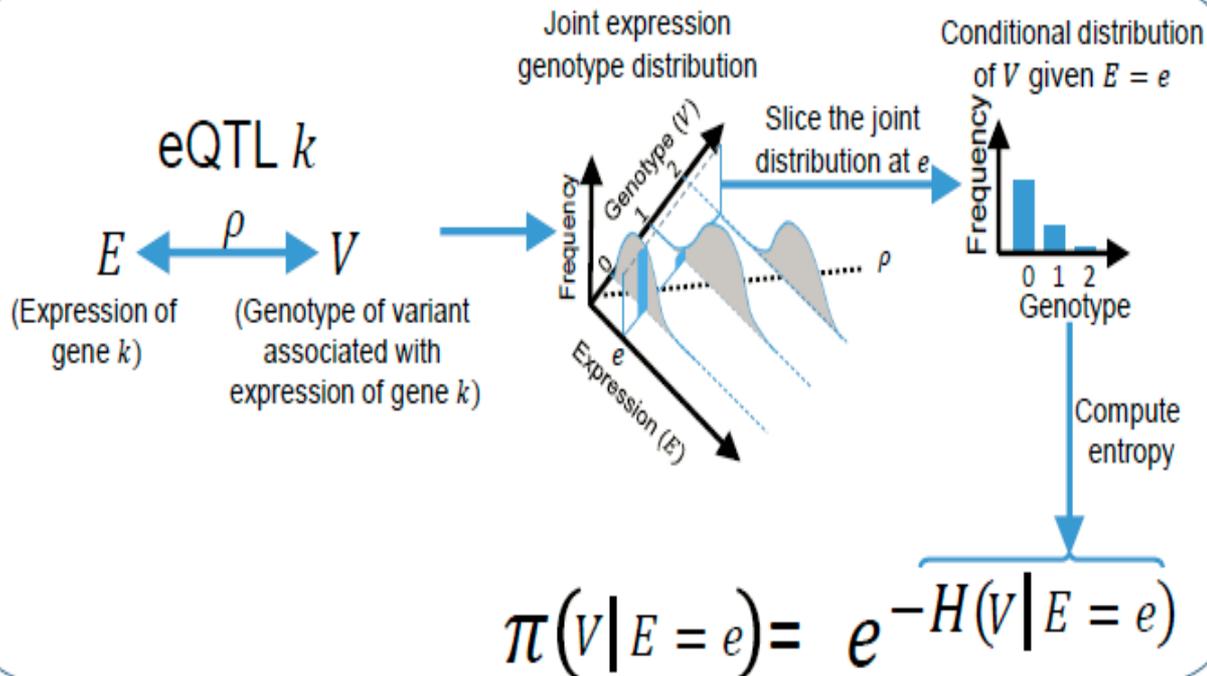
# Information Content and Predictability

$$ICI \left( \begin{array}{l} \text{Individual has variant} \\ \text{genotypes } g_1, g_2, \dots, g_n \\ \text{for variants } V_1, V_2, \dots, V_n \end{array} \right) = \log \left( \frac{1}{\text{Frequency of } V_1 \text{ genotype}} \right) + \log \left( \frac{1}{\text{Frequency of } V_2 \text{ genotype}} \right) + \dots + \log \left( \frac{1}{\text{Frequency of } V_n \text{ genotype}} \right)$$

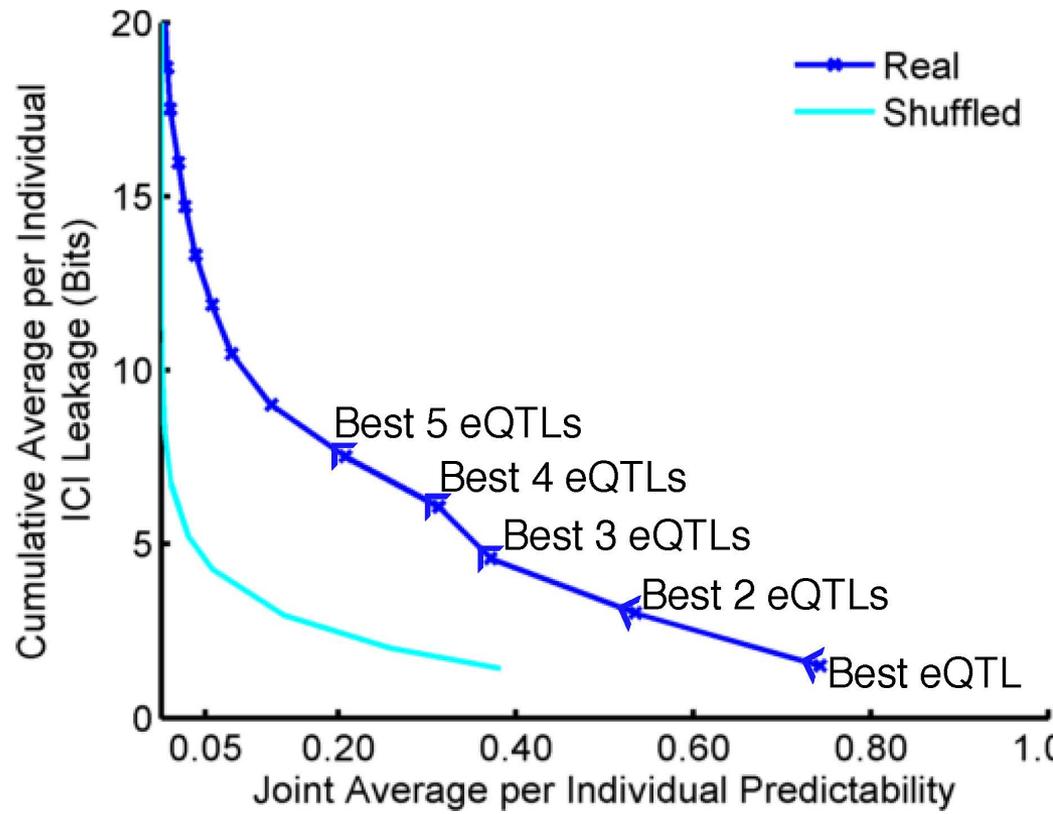
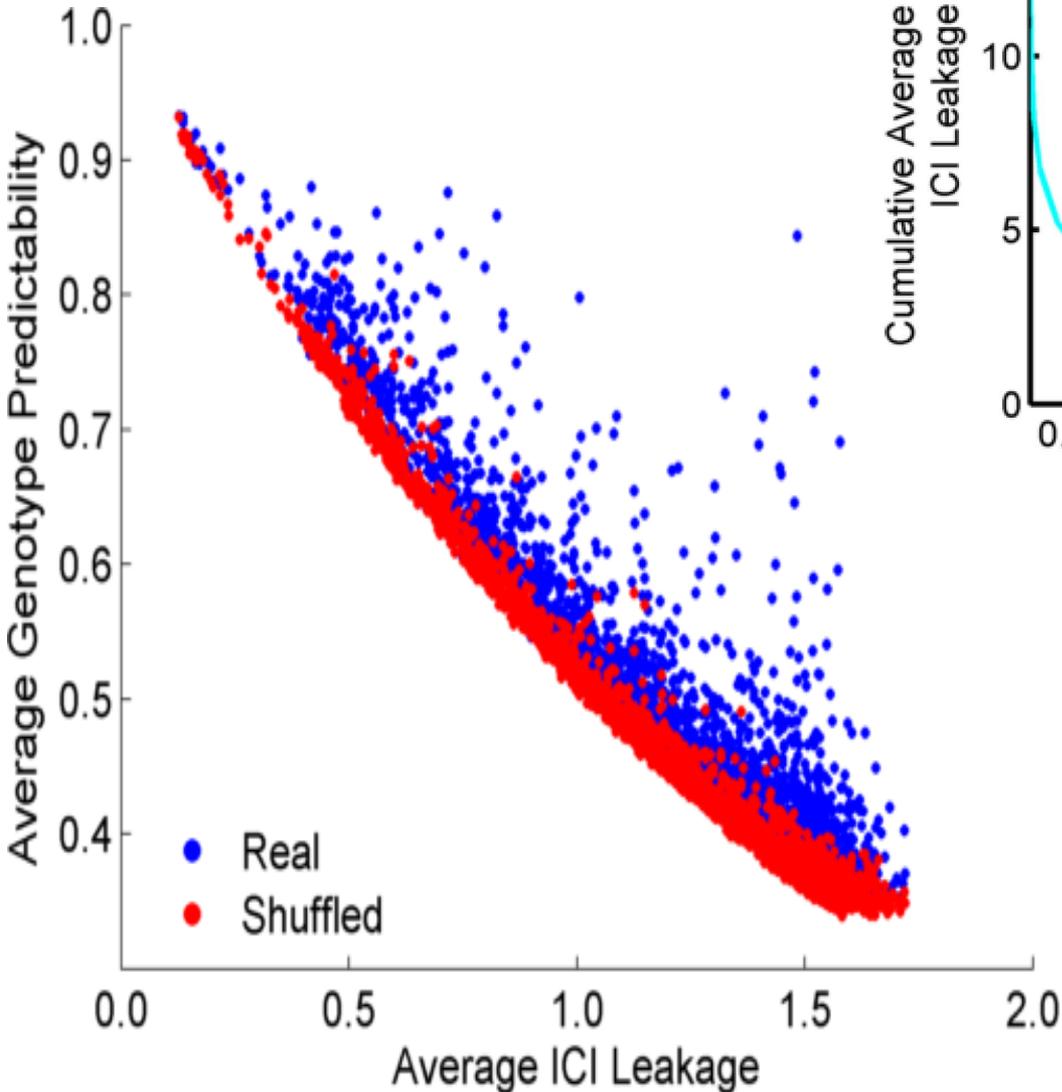
$g_1 = 2$                        $g_2 = 1$                        $g_n = 2$

$V_1$  genotype frequencies                       $V_2$  genotype frequencies                       $V_n$  genotype frequencies

- Naive measure of information (no LD, distant correlations, pop. struc., &c)
- Higher frequency: Lower ICI
- Additive for multiple variants

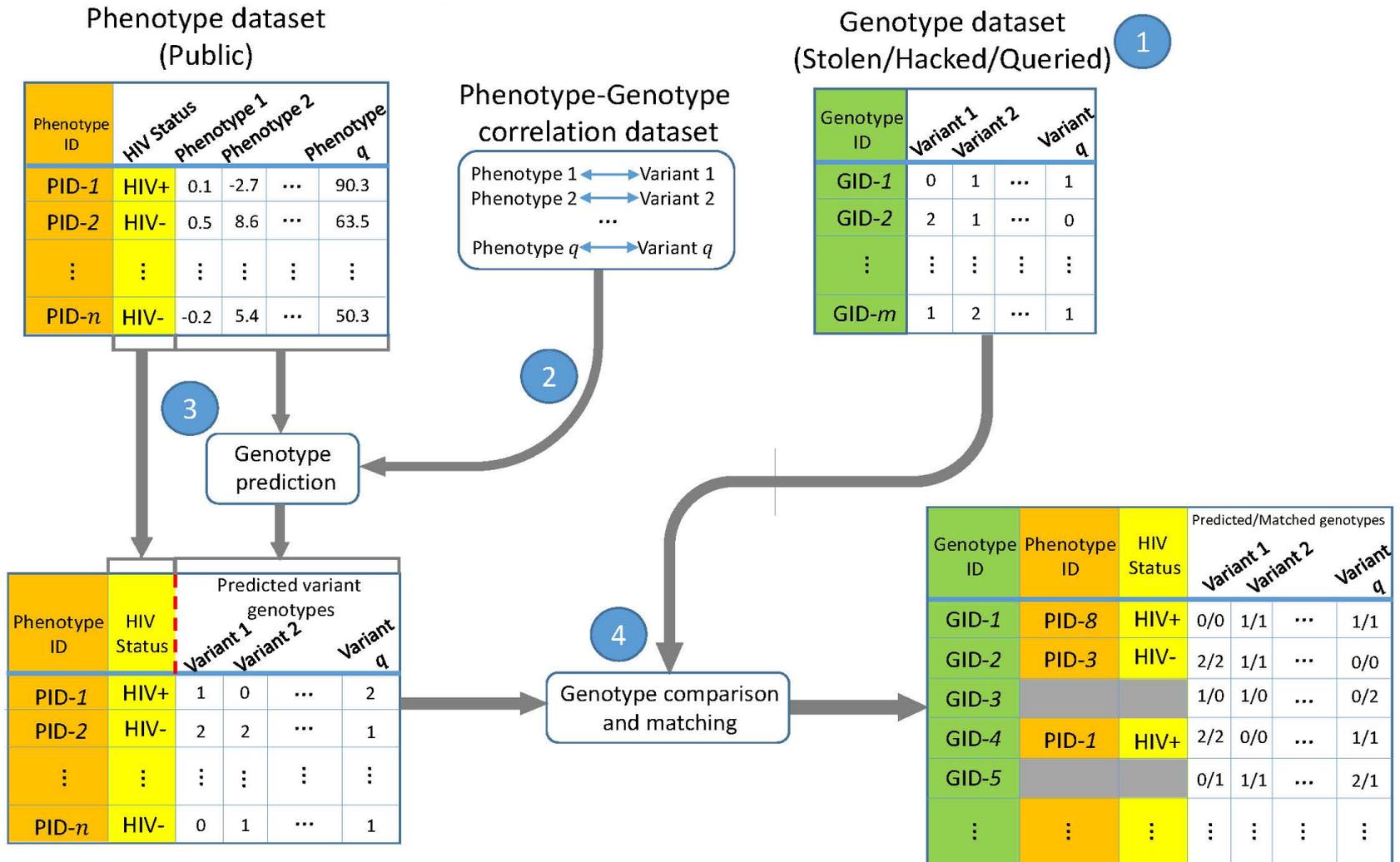


- Condition specific entropy
- Higher cond. entropy: Lower predictability
- Additive for multiple eQTLs

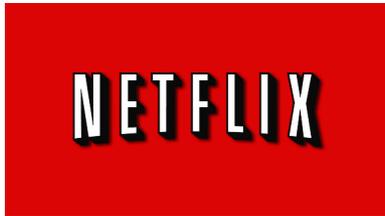


# ICI Leakage versus Genotype Predictability

# Linking Attack Scenario



# Linking Attacks: Case of Netflix Prize



Names available for many users!

| User (ID) | Movie (ID) | Date of Grade | Grade [1,2,3,4,5] |
|-----------|------------|---------------|-------------------|
| NTFLX-0   | NTFLX-19   | 10/12/2008    | 1                 |
| NTFLX-1   | NTFLX-116  | 4/23/2009     | 3                 |
| NTFLX-2   | NTFLX-92   | 5/27/2010     | 2                 |
| NTFLX-1   | NTFLX-666  | 6/6/2016      | 5                 |
| ...       | ...        | ...           | ...               |
| ...       | ...        | ...           | ...               |

| User (ID) | Movie (ID) | Date of Grade | Grade [0-10] |
|-----------|------------|---------------|--------------|
| IMDB-0    | IMDB-173   | 4/20/2009     | 5            |
| IMDB-1    | IMDB-18    | 10/18/2008    | 0            |
| IMDB-2    | IMDB-341   | 5/27/2010     | -            |
| ...       | ...        | ...           | ...          |
| ...       | ...        | ...           | ...          |
| ...       | ...        | ...           | ...          |

- Many users are shared
- The grades of same users are correlated
- A user grades one movie around the same date in two databases

Anonymized Netflix Prize Training Dataset  
made available to contestants

# Linking Attacks: Case of Netflix Prize



| User (ID) | Movie (ID) | Date of Grade | Grade [1,2,3,4,5] |
|-----------|------------|---------------|-------------------|
| NTFLX-0   | NTFLX-19   | 10/12/2008    | 1                 |
| NTFLX-1   | NTFLX-116  | 4/23/2009     | 3                 |
| NTFLX-2   | NTFLX-92   | 5/27/2010     | 2                 |
| NTFLX-1   | NTFLX-666  | 6/6/2016      | 5                 |
| ...       | ...        | ...           | ...               |
| ...       | ...        | ...           | ...               |

| User (ID) | Movie (ID) | Date of Grade | Grade [0-10] |
|-----------|------------|---------------|--------------|
| IMDB-0    | IMDB-173   | 4/20/2009     | 5            |
| IMDB-1    | IMDB-18    | 10/18/2008    | 0            |
| IMDB-2    | IMDB-341   | 5/27/2010     | -            |
| ...       | ...        | ...           | ...          |
| ...       | ...        | ...           | ...          |
| ...       | ...        | ...           | ...          |

- Many users are shared
- The grades of same users are correlated
- A user grades one movie around the same date in two databases
- IMDB users are public
- NetFLIX and IMdB moves are public

# Linking Attacks: Case of Netflix Prize

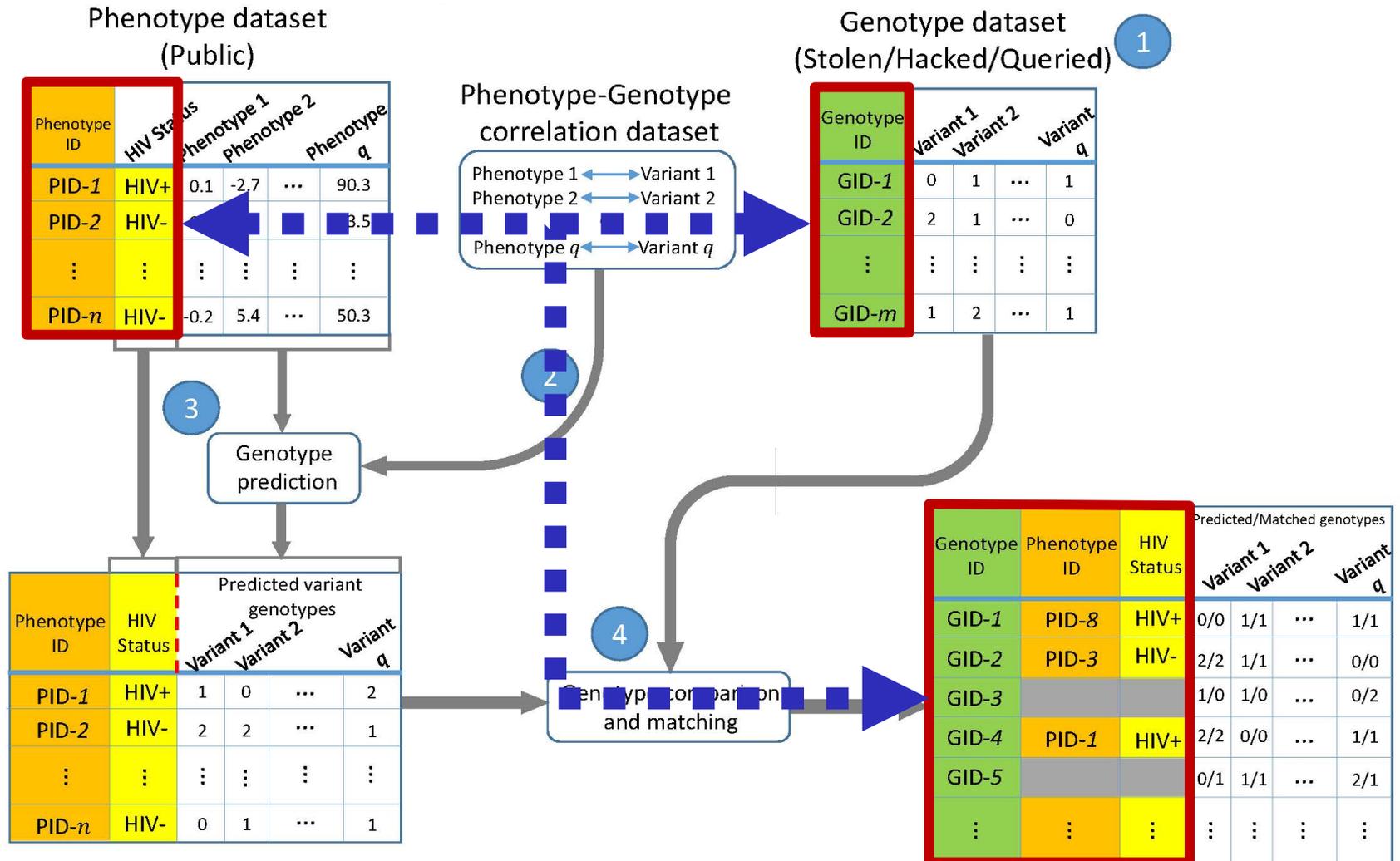


| User (ID) | Movie (ID)       | Date of Grade | Grade [1,2,3,4,5] |
|-----------|------------------|---------------|-------------------|
| NTFLX-0   | NTFLX-19         | 10/12/2008    | 1                 |
| NTFLX-1   | NTFLX-116        | 4/23/2009     | 3                 |
| NTFLX-2   | NTFLX-92         | 5/27/2010     | 2                 |
| NTFLX-1   | <b>NTFLX-666</b> | 6/6/2016      | 5                 |
| ...       | ...              | ...           | ...               |
| ...       | ...              | ...           | ...               |

| User (ID) | Movie (ID) | Date of Grade | Grade [0-10] |
|-----------|------------|---------------|--------------|
| IMDB-0    | IMDB-173   | 4/20/2009     | 5            |
| IMDB-1    | IMDB-18    | 10/18/2008    | 0            |
| IMDB-2    | IMDB-341   | 5/27/2010     | -            |
| ...       | ...        | ...           | ...          |
| ...       | ...        | ...           | ...          |
| ...       | ...        | ...           | ...          |

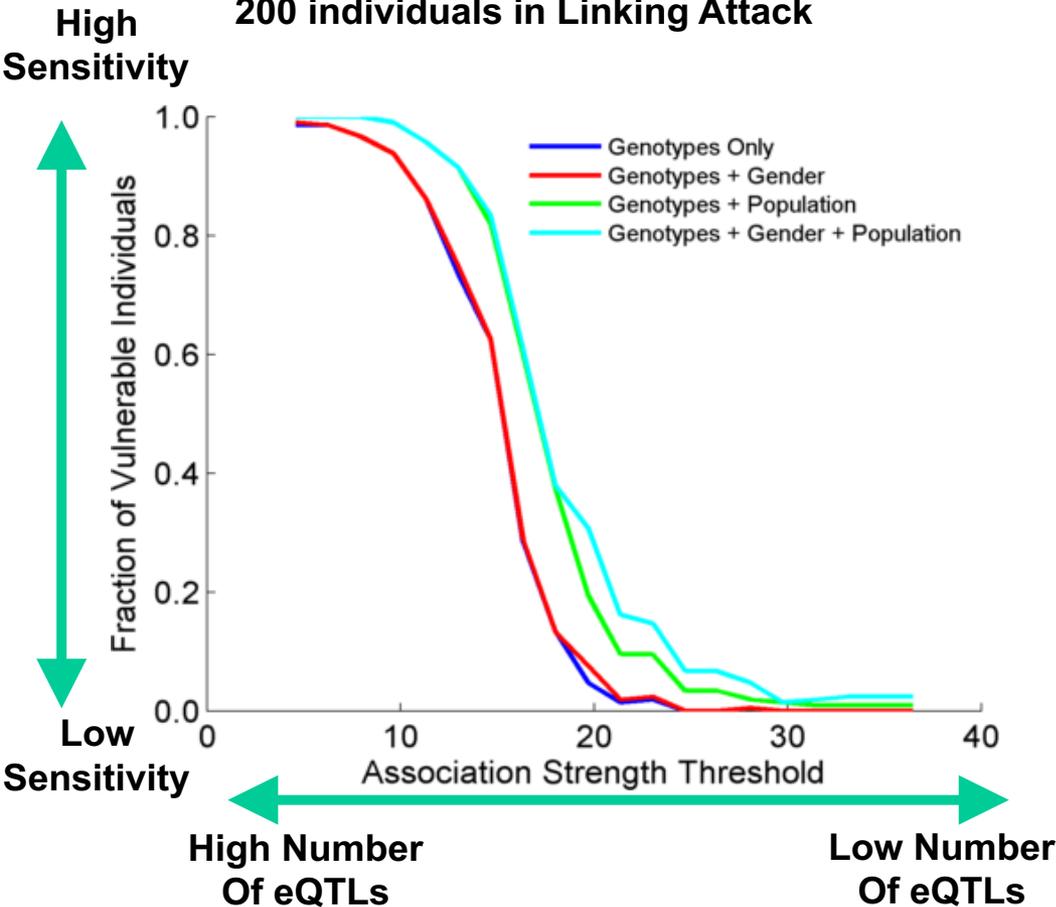
- Many users are shared
- The grades of same users are correlated
- A user grades one movie around the same date in two databases

# Linking Attack Scenario



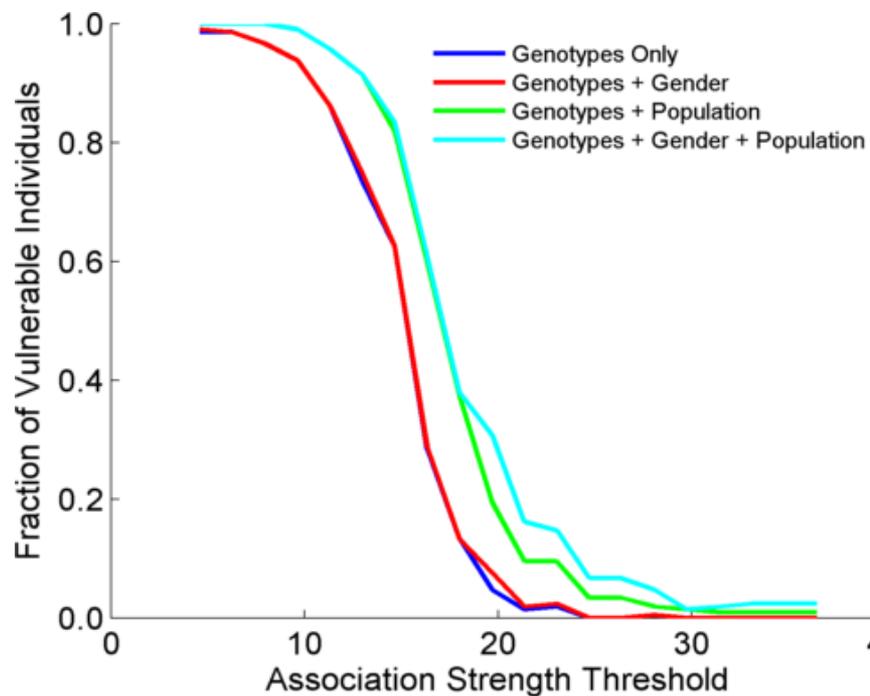
# Success in Linking Attack with Extremity based Genotype Prediction

200 individuals eQTL Discovery  
200 individuals in Linking Attack

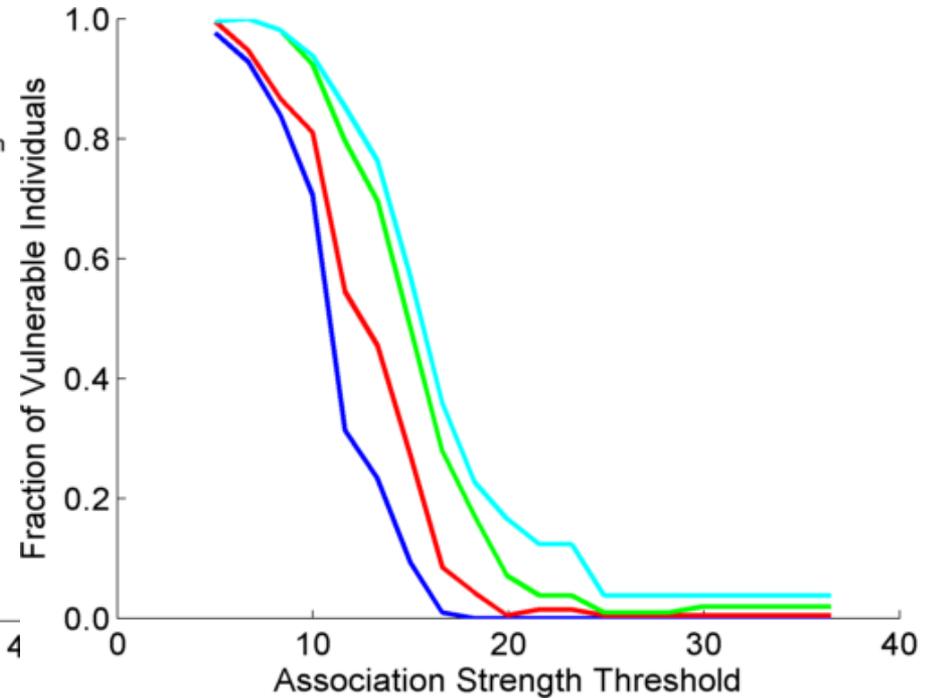


# Success in Linking Attack with Extremity based Genotype Prediction

200 individuals eQTL Discovery  
200 individuals in Linking Attack



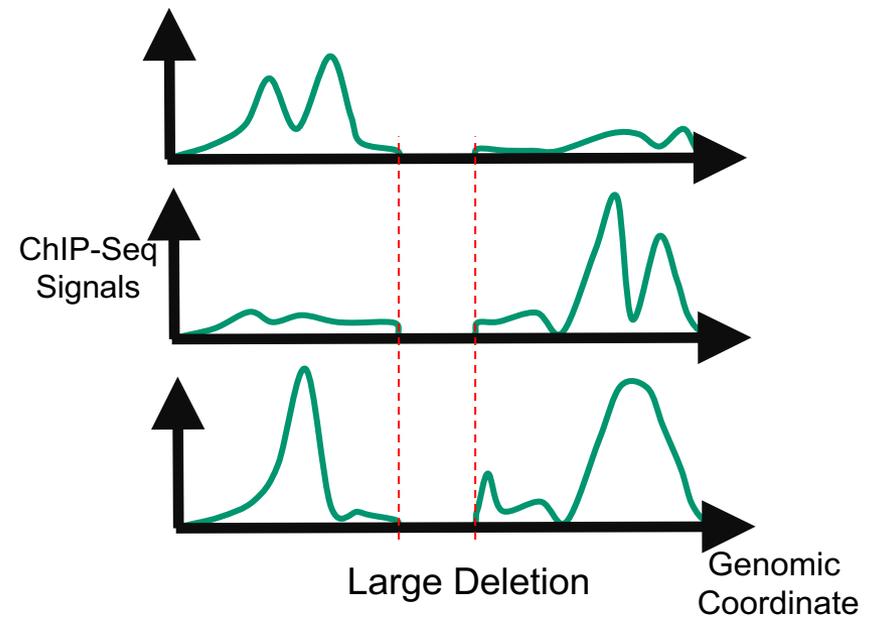
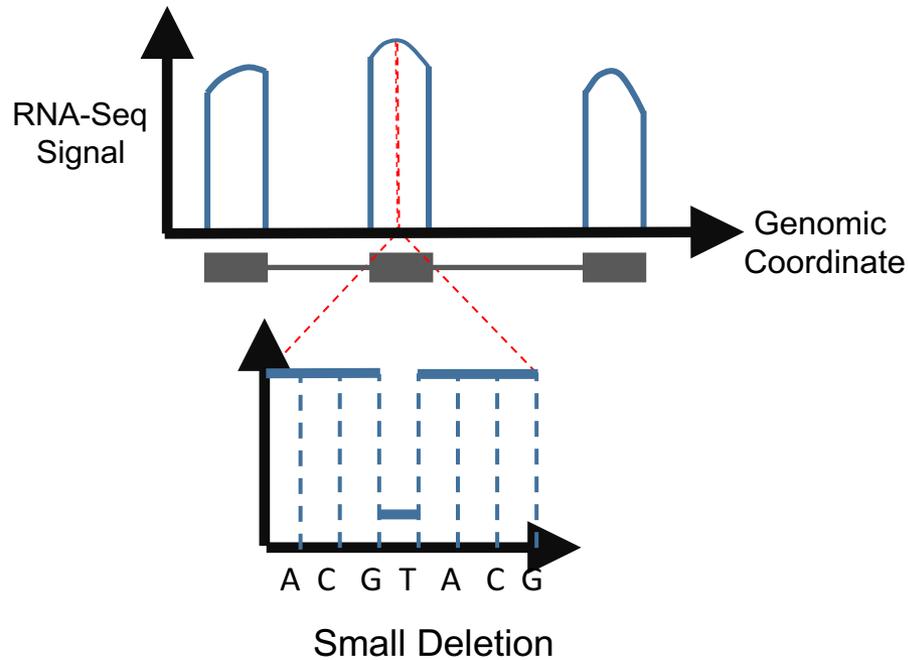
200 individuals eQTL Discovery  
100,200 individuals in Linking Attack



## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

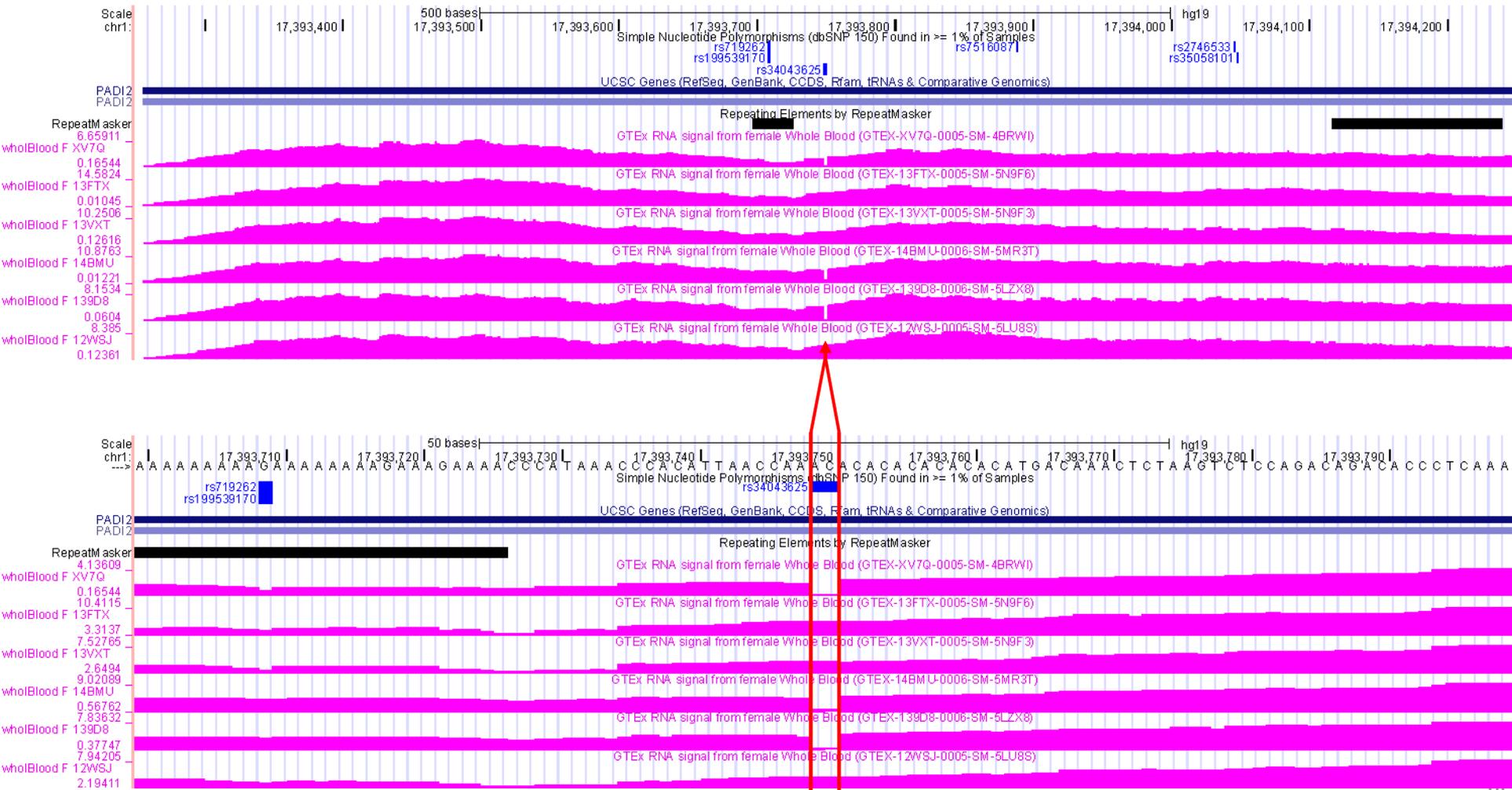
- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

# Detection & Genotyping of small & large SV deletions from signal profiles



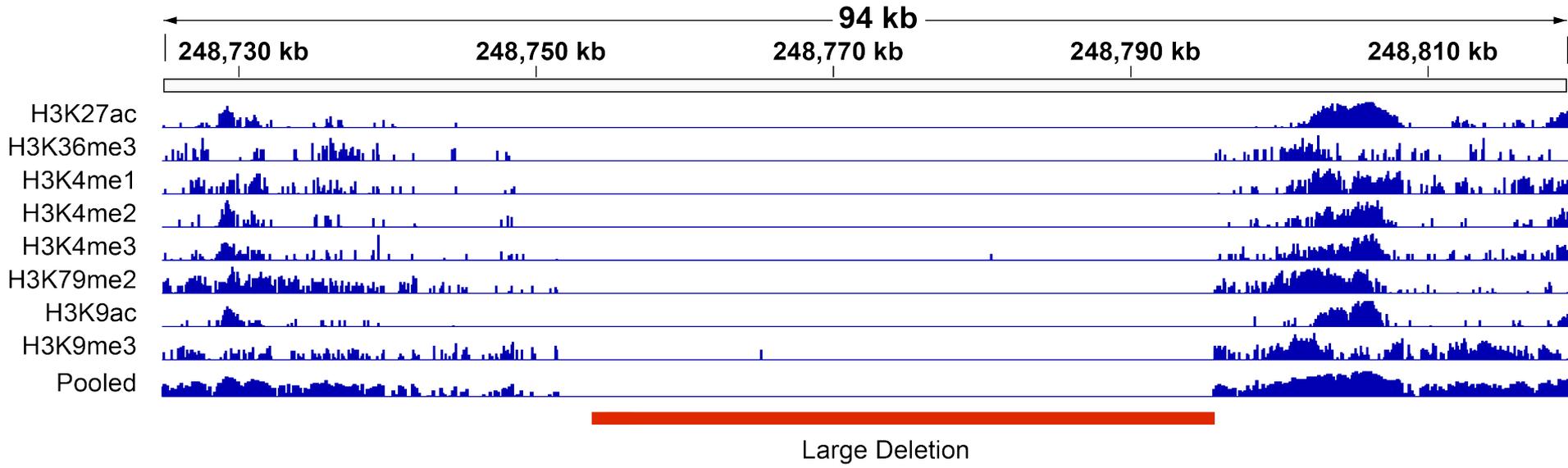
RNA-seq also shows large deletions

# Example of Small Deletion Evident in Signal Profile



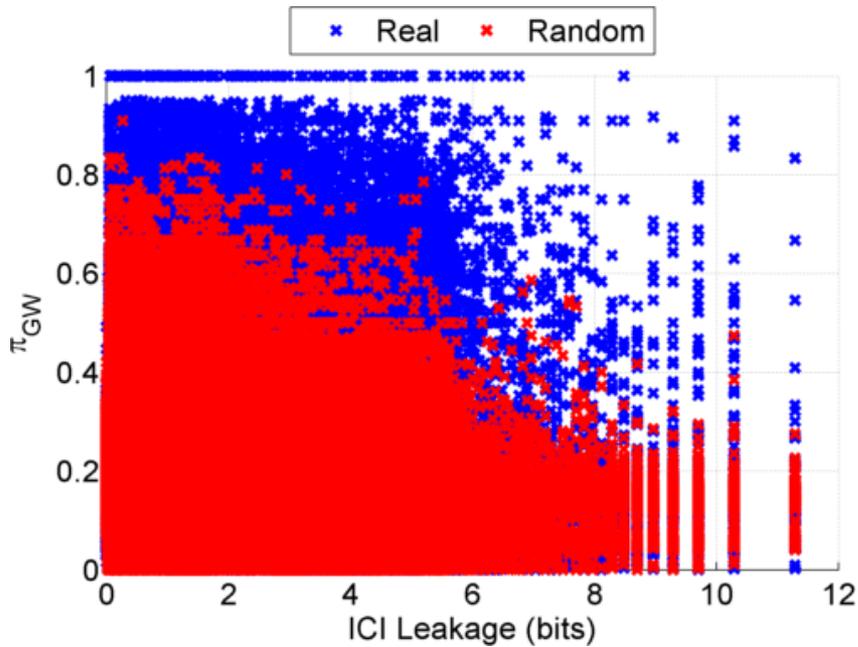
[Harmanci & Gerstein, *Nat. Comm.* ('18)]

# Example of Large Deletion Evident in Signal Profile

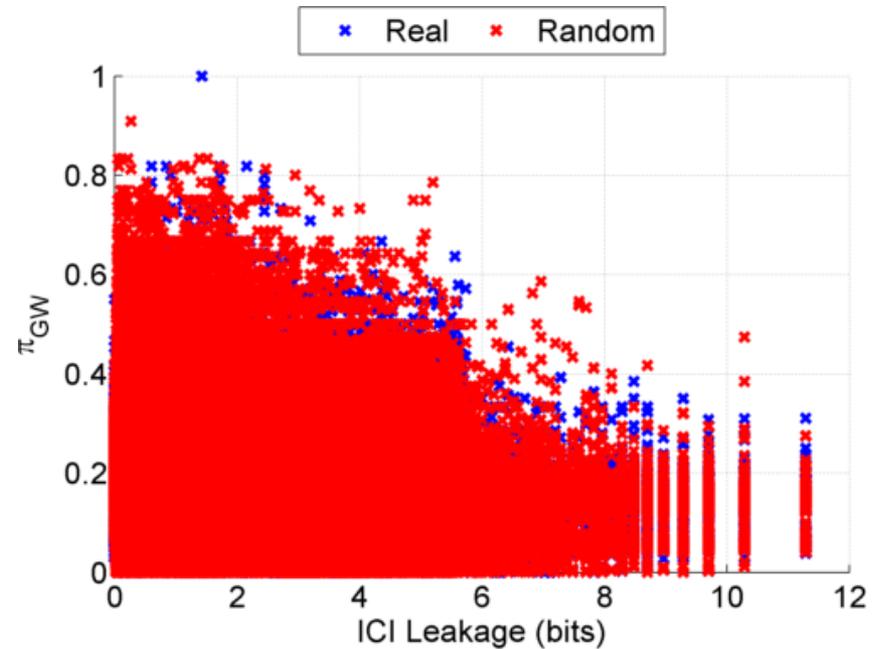


# Information Leakage from SV Deletions

a) Before Anonymization

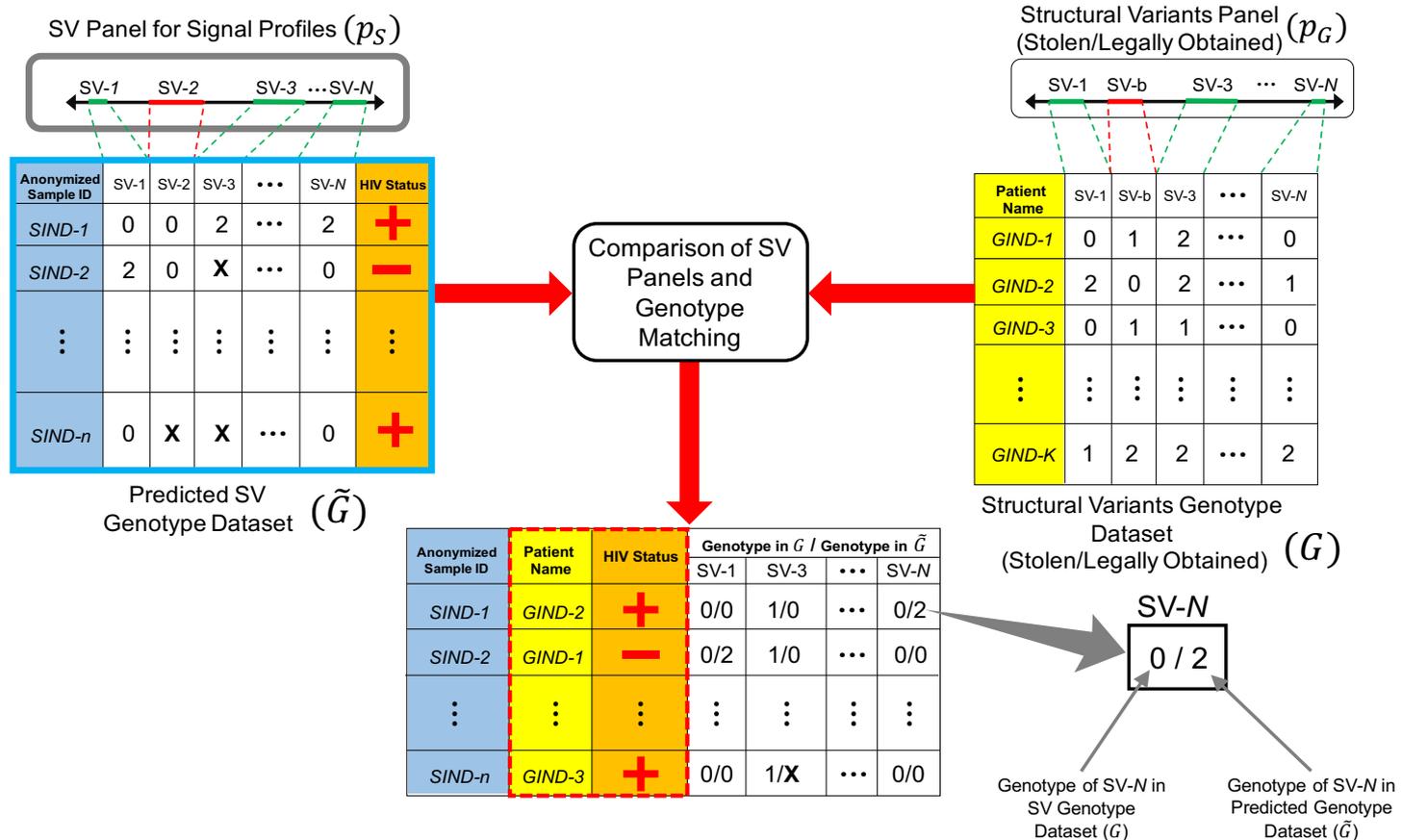


b) After Anonymization

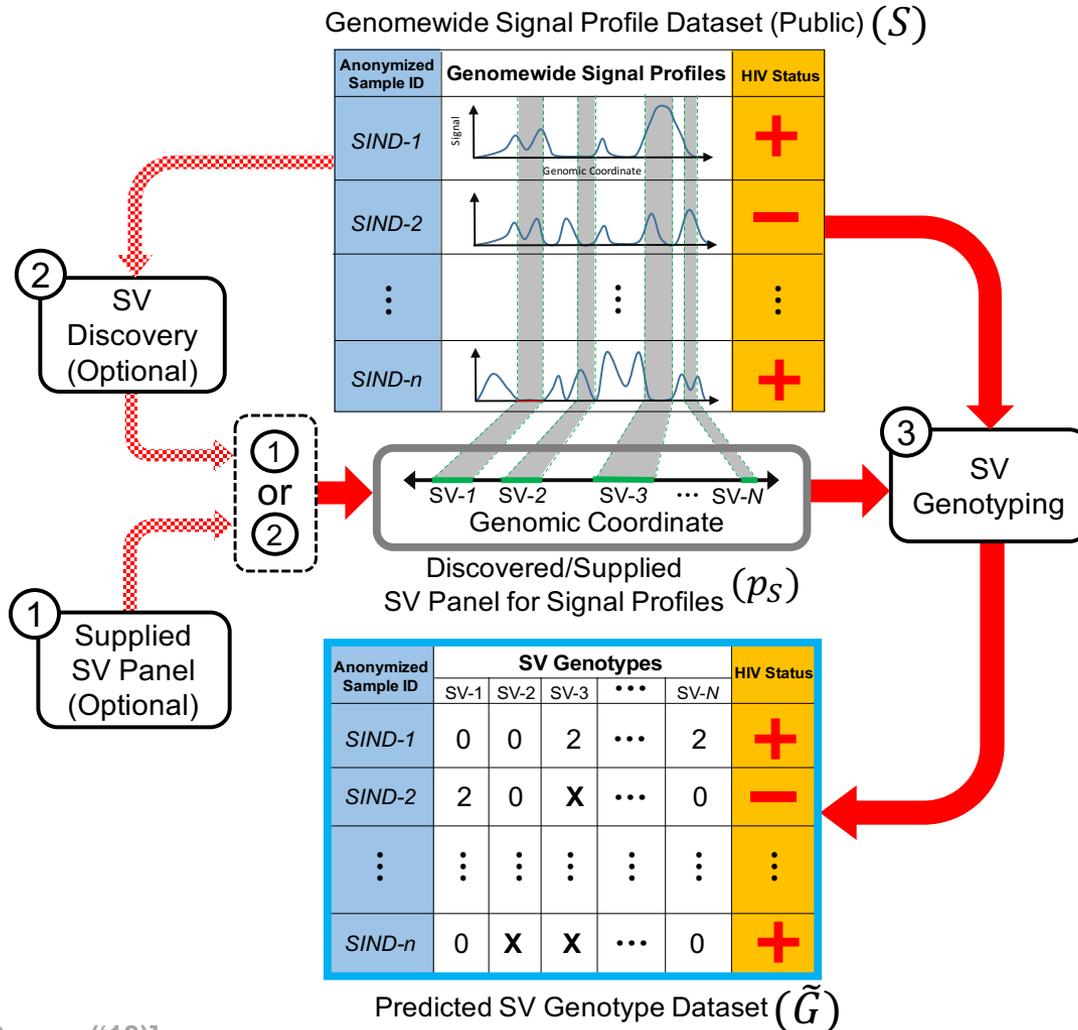


Simple anonymization procedure (filling in deletion by value at endpoints) has dramatic effect

# Another type of Linking Attack: Linking based on SV Genotyping

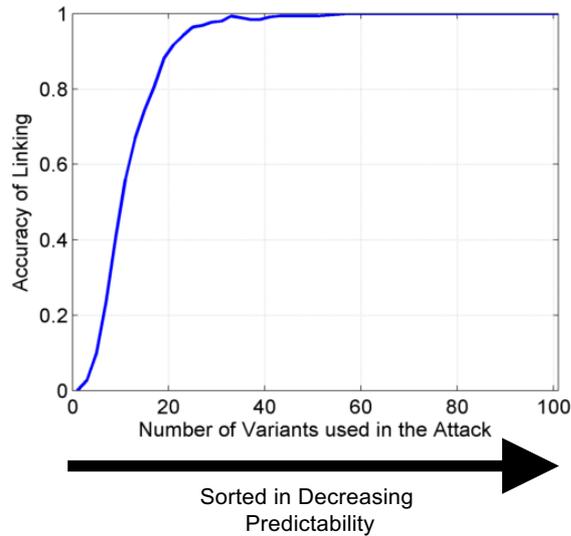


# Another type of Linking Attack: First Doing SV Genotyping

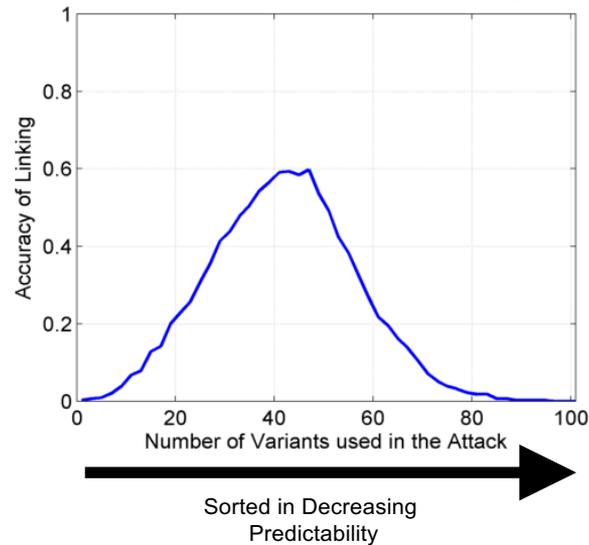


# Linking Attack Based on SV Deletions in gEUVADIS Dataset

c) Genotyping  
(1kG MAF>0.01)



d) Discovery + Genotyping



## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

## Using population-scale functional genomics to understand neuropsychiatric disease & interpreting the data exhaust from this activity

- *[Core]* **PsychENCODE**: Population-level analysis of functional genomics data related to neuropsychiatric disease
  - Construction of an adult brain resource with 1866 individuals + dev. time-course
  - Using the changing proportions of cell types (via **single-cell deconvolution**) to account for expression variation across a population, disorders & development
  - Large-scale processing defines ~79K PFC **enhancers & creates a comprehensive QTL** resource (~2.5M eQTLs + cQTLs & fQTLs)
  - Connecting QTLs, enhancer activity relationships & Hi-C into a **brain regulatory network** & using this to link SCZ GWAS SNPs to genes
  - Embedding the reg. network in a **deep-learning model** to predict disease from genotype & transcriptome. Using this to suggest specific pathways & genes, as targets.
- *[Exhaust]* **Other uses** for the resource
  - Highlighting aging related genes + consistently comparing the brain to other organs
- *[Exhaust]* **Genomic Privacy**
  - The **Dilemma**
    - The genome as fundamental, inherited info that's very private v. need for large-scale mining for med. research
    - 2-sided nature of RNA-seq presents tricky privacy issues
  - **eQTLs**: Quantifying & removing variant info from expression levels with ICI & predictability. Instantiating a practical linking attack with noisy quasi-identifiers
  - **Signal Profiles**: Manifest appreciable leakage from large & small deletions. Linking attacks possible but additional complication of SV discovery in addition to genotyping

## “Adult Capstone” Team – 1 of 3 capstones

# PsychENCODE Acknowledgment



- Geetha Senthil
- Lora Bingaman
- David Panchision
- Alexander Arguello
- Thomas Lehner

**Daifeng Wang, Shuang Liu, Jonathan Warrell, Hyejung Won, Xu Shi, Fabio Navarro, Declan Clarke, Mengting Gu, Prashant Emani, Yucheng T. Yang, Min Xu, Michael Gandal, Shaoke Lou, Jing Zhang, Jonathan J. Park, Chengfei Yan, Sunh Kyong Rhie, Kasidet Manakongtreecheep, Holly Zhou, Aparna Nathan, Mette Peters, Eugenio Mattei, Dominic Fitzgerald, Tonya Brunetti, Jill Moore, Yan Jiang, Kiran Girdhar, Gabriel Hoffman, Selim Kalayci, Zeynep Hulya Gumus, Greg Crawford,**

**PsychENCODE Consortium,**

**Panos Roussos, Schahram Akbarian, Andrew E. Jaffe, Kevin White, Zhiping Weng, Nenad Sestan,**

**Daniel H. Geschwind, James A. Knowles**

Dedicated to **Pamela Sklar**

**Resource.psychencode.org**

**The PsychENCODE Consortium:** Allison E Ashley-Koch, Duke University; Gregory E Crawford, Duke University; Melanie E Garrett, Duke University; Lingyun Song, Duke University; Alexias Safi, Duke University; Graham D Johnson, Duke University; Gregory A Wray, Duke University; Timothy E Reddy, Duke University; Fernando S Goes, Johns Hopkins University; Peter Zandi, Johns Hopkins University; Julien Bryois, Karolinska Institutet; Andrew E Jaffe, Lieber Institute for Brain Development; Amanda J Price, Lieber Institute for Brain Development; Nikolay A Ivanov, Lieber Institute for Brain Development; Leonardo Collado-Torres, Lieber Institute for Brain Development; Thomas M Hyde, Lieber Institute for Brain Development; Emily E Burke, Lieber Institute for Brain Development; Joel E Kleiman, Lieber Institute for Brain Development; Ran Tao, Lieber Institute for Brain Development; Joo Heon Shin, Lieber Institute for Brain Development; Schahram Akbarian, Icahn School of Medicine at Mount Sinai; Kiran Girdhar, Icahn School of Medicine at Mount Sinai; Yan Jiang, Icahn School of Medicine at Mount Sinai; Marija Kundakovic, Icahn School of Medicine at Mount Sinai; Leanne Brown, Icahn School of Medicine at Mount Sinai; Bibi S Kassim, Icahn School of Medicine at Mount Sinai; Royce B Park, Icahn School of Medicine at Mount Sinai; Jennifer R Wiseman, Icahn School of Medicine at Mount Sinai; Elizabeth Zharovsky, Icahn School of Medicine at Mount Sinai; Rivka Jacobov, Icahn School of Medicine at Mount Sinai; Olivia Devillers, Icahn School of Medicine at Mount Sinai; Elie Flatow, Icahn School of Medicine at Mount Sinai; Gabriel E Hoffman, Icahn School of Medicine at Mount Sinai; Barbara K Lipska, Human Brain Collection Core, National Institutes of Health, Bethesda, MD; David A Lewis, University of Pittsburgh; Vahram Haroutunian, Icahn School of Medicine at Mount Sinai and James J Peters VA Medical Center; Chang-Gyu Hahn, University of Pennsylvania; Alexander W Charney, Mount Sinai; Stella Dracheva, Mount Sinai; Alexey Kozlenkov, Mount Sinai; Judson Belmont, Icahn School of Medicine at Mount Sinai; Diane DelValle, Icahn School of Medicine at Mount Sinai; Nancy Francoeur, Icahn School of Medicine at Mount Sinai; Evi Hadjimichael, Icahn School of Medicine at Mount Sinai; Dalila Pinto, Icahn School of Medicine at Mount Sinai; Harm van Bakel, Icahn School of Medicine at Mount Sinai; Panos Roussos, Mount Sinai; John F Fullard, Mount Sinai; Jaroslav Bendl, Mount Sinai; Mads E Hauberg, Mount Sinai; Lara M Mangravite, Sage Bionetworks; Mette A Peters, Sage Bionetworks; Yooree Chae, Sage Bionetworks; Junmin Peng, St. Jude Children's Hospital; Mingming Niu, St. Jude Children's Hospital; Xusheng Wang, St. Jude Children's Hospital; Maree J Webster, Stanley Medical Research Institute; Thomas G Beach, Banner Sun Health Research Institute; Chao Chen, Central South University; Yi Jiang, Central South University; Rujia Dai, Central South University; Annie W Shieh, SUNY Upstate Medical University; Chunyu Liu, SUNY Upstate Medical University; Kay S. Grennan, SUNY Upstate Medical University; Yan Xia, SUNY Upstate Medical University/Central South University; Ramu Vadukapuram, SUNY Upstate Medical University; Yongjun Wang, Central South University; Dominic Fitzgerald, The University of Chicago; Lijun Cheng, The University of Chicago; Miguel Brown, The University of Chicago; Mimi Brown, The University of Chicago; Tonya Brunetti, The University of Chicago; Thomas Goodman, The University of Chicago; Majd Alsayed, The University of Chicago; Michael J Gandal, University of California, Los Angeles; Daniel H Geschwind, University of California, Los Angeles; Hyejung Won, University of California, Los Angeles; Damon Polioudakis, University of California, Los Angeles; Brie Wamsley, University of California, Los Angeles; Jiani Yin, University of California, Los Angeles; Tarik Hadzic, University of California, Los Angeles; Luis De La Torre Ubieta, UCLA; Vivek Swarup, University of California, Los Angeles; Stephan J Sanders, University of California, San Francisco; Matthew W State, University of California, San Francisco; Donna M Werling, University of California, San Francisco; Joon-Yong An, University of California, San Francisco; Brooke Sheppard, University of California, San Francisco; A Jeremy Willsey, University of California, San Francisco; Kevin P White, The University of Chicago; Mohana Ray, The University of Chicago; Gina Giase, SUNY Upstate Medical University; Amira Kefi, University of Illinois at Chicago; Eugenio Mattei, University of Massachusetts Medical School; Michael Purcaro, University of Massachusetts Medical School; Zhiping Weng, University of Massachusetts Medical School; Jill Moore, University of Massachusetts Medical School; Henry Pratt, University of Massachusetts Medical School; Jack Huey, University of Massachusetts Medical School; Tyler Borrmann, University of Massachusetts Medical School; Patrick F Sullivan, University of North Carolina - Chapel Hill; Paola Giusti-Rodriguez, University of North Carolina - Chapel Hill; Yunjung Kim, University of North Carolina - Chapel Hill; Patrick Sullivan, University of North Carolina - Chapel Hill; Jin Szatkiewicz, University of North Carolina - Chapel Hill; Sunh Kyong Rhie, University of Southern California; Christopher Armoskus, University of Southern California; Adrian Camarena, University of Southern California; Peggy J Farnham, University of Southern California; Valeria N Spitsyna, University of Southern California; Heather Witt, University of Southern California; Shannon Schreiner, University of Southern California; Oleg V Evgrafov, SUNY Downstate Medical Center; James A Knowles, SUNY Downstate Medical Center; Mark Gerstein, Yale University; Shuang Liu, Yale University; Daifeng Wang, Stony Brook University; Fabio C. P. Navarro, Yale University; Jonathan Warrell, Yale University; Declan Clarke, Yale University; Prashant S. Emani, Yale University; Mengting Gu, Yale University; Xu Shi, Yale University; Min Xu, Yale University; Yucheng T. Yang, Yale University; Robert R. Kitchen, Yale University; Gamze Guirsoy, Yale University; Jing Zhang, Yale University; Becky C Carlyle, Yale University; Angus C Nairn, Yale University; Mingfeng Li, Yale University; Sirisha Pochareddy, Yale University; Nenad Sestan, Yale University; Mario Skarica, Yale University; Zhen Li, Yale University; Andre M.M. Sousa, Yale University; Gabriel Santpere, Yale University; Jinmyung Choi, Yale University; Ying Zhu, Yale University; Tianliuyun Gao, Yale University; Daniel J Miller, Yale University; Adriana Cherskov, Yale University; Mo Yang, Yale University; Anahita Amiri, Yale University; Gianfilippo Coppola, Yale University; Jessica Mariani, Yale University; Soraya Scuderi, Yale University; Anna Szekely, Yale University; Flora M Vaccarino, Yale University; Feinan Wu, Yale University; Sherman Weissman, Yale University; Tanmoy Roychowdhury, Mayo Clinic Rochester; Alexej Abyzov, Mayo Clinic Rochester;

## Developmental Capstone

- **M Li, G Santpere, Y Imamura Kawasawa, OV Evgrafov, FO Gulden, S Pochareddy, SM Sunkin, Z Li, Y Shin,**

Y Zhu, AMM Sousa, DM Werling, RR Kitchen, HJ Kang, M Pletikos, J Choi, S Muchnik, X Xu, D Wang, B Lorente-Galdos, S Liu, P Giusti-Rodriguez, H Won, CA de Leeuw, AF Pardini,

**BrainSpan Consortium, PsychENCODE Consortium,  
PsychENCODE Developmental Subgroup,**

M Hu, F Jin, Y Li, MJ Owen, MC O'Donovan, JTR Walters, D Posthuma, MA Reimers, P Levitt, DR Weinberger, TM Hyde, JE Kleinman, DH Geschwind, MJ Hawrylycz, MW State, SJ

Sanders, PF Sullivan,

**ES Lein, JA Knowles, N Sestan**

**psychencode.org**

# Acknowledgements



[papers.gersteinlab.org/subject/privacy](https://papers.gersteinlab.org/subject/privacy)

**PrivaSig**.gersteinlab.org

**PrivaSeq**.gersteinlab.org

Hiring Postdocs. See

**JOB**S.gersteinlab.org !

A **Harmanci**,

D **Greenbaum**, G **Gürsoy**

**Extra**



# Info about content in this slide pack

- General PERMISSIONS
  - This Presentation is copyright Mark Gerstein, Yale University, 2017.
  - Please read permissions statement at [www.gersteinlab.org/misc/permissions.html](http://www.gersteinlab.org/misc/permissions.html) .
  - Feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).
  - Paper references in the talk were mostly from Papers.GersteinLab.org.
- PHOTOS & IMAGES. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see <http://streams.gerstein.info> .
  - In particular, many of the images have particular EXIF tags, such as kwpotppt , that can be easily queried from flickr, viz: <http://www.flickr.com/photos/mbgmbg/tags/kwpotppt>