

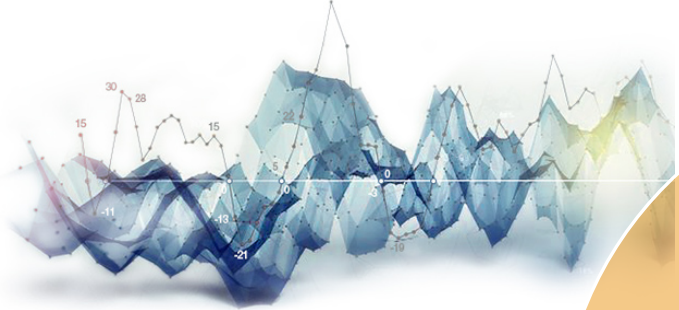
# Multivariate Multiscale Impacts of Genetic Variants on Gene Expression Variability in Humans

---

JAMES CAI

1/20/2017





**Computational  
Statistics**

**Data Science**

**Medical  
Genetics**



# Outline

Additive, epistatic, and environmental effects through the lens of **evQTLs**

Exploiting **aberrant gene expression** in autism for gene discovery and diagnosis



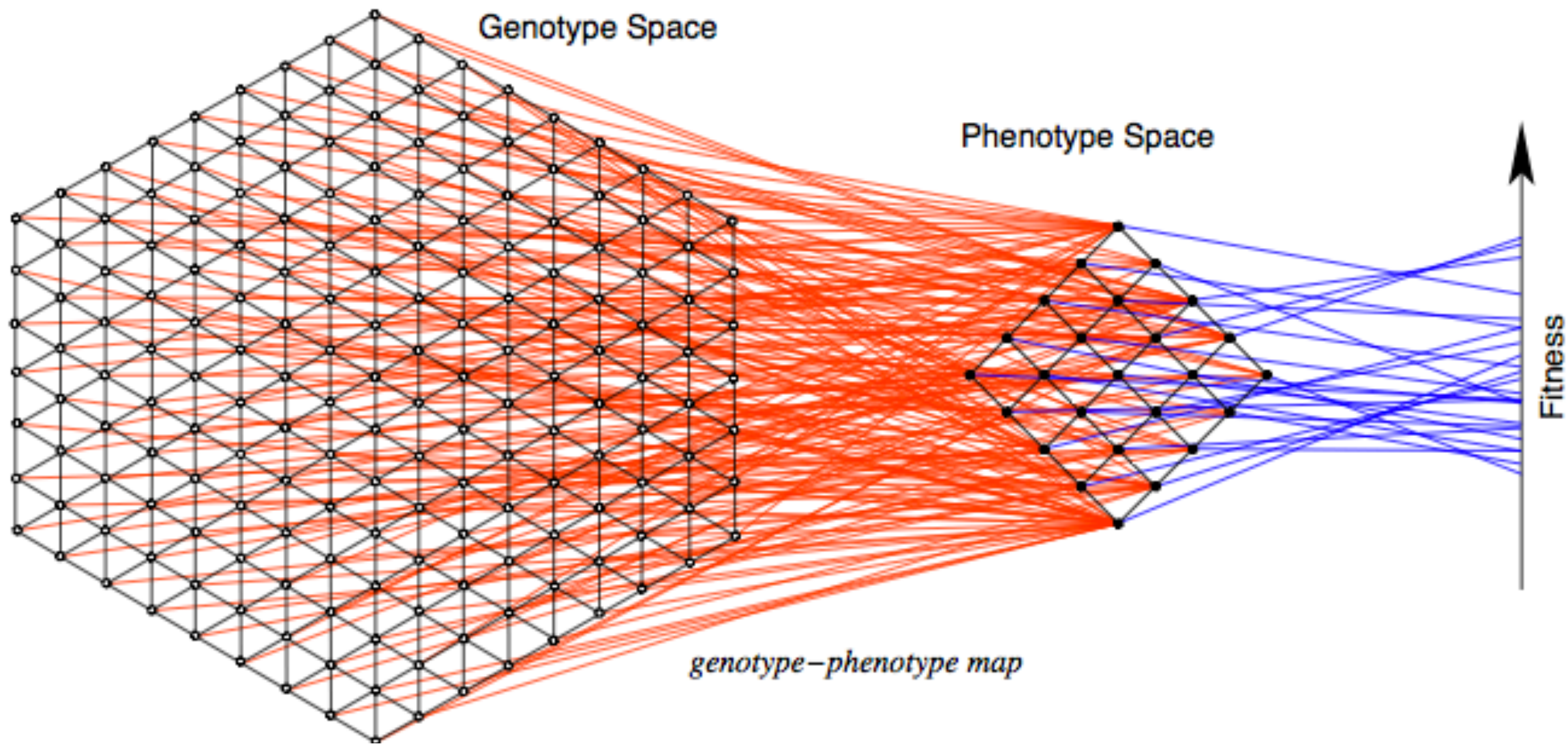
# Additive, epistatic, and environmental effects through the lens of **evQTLs**

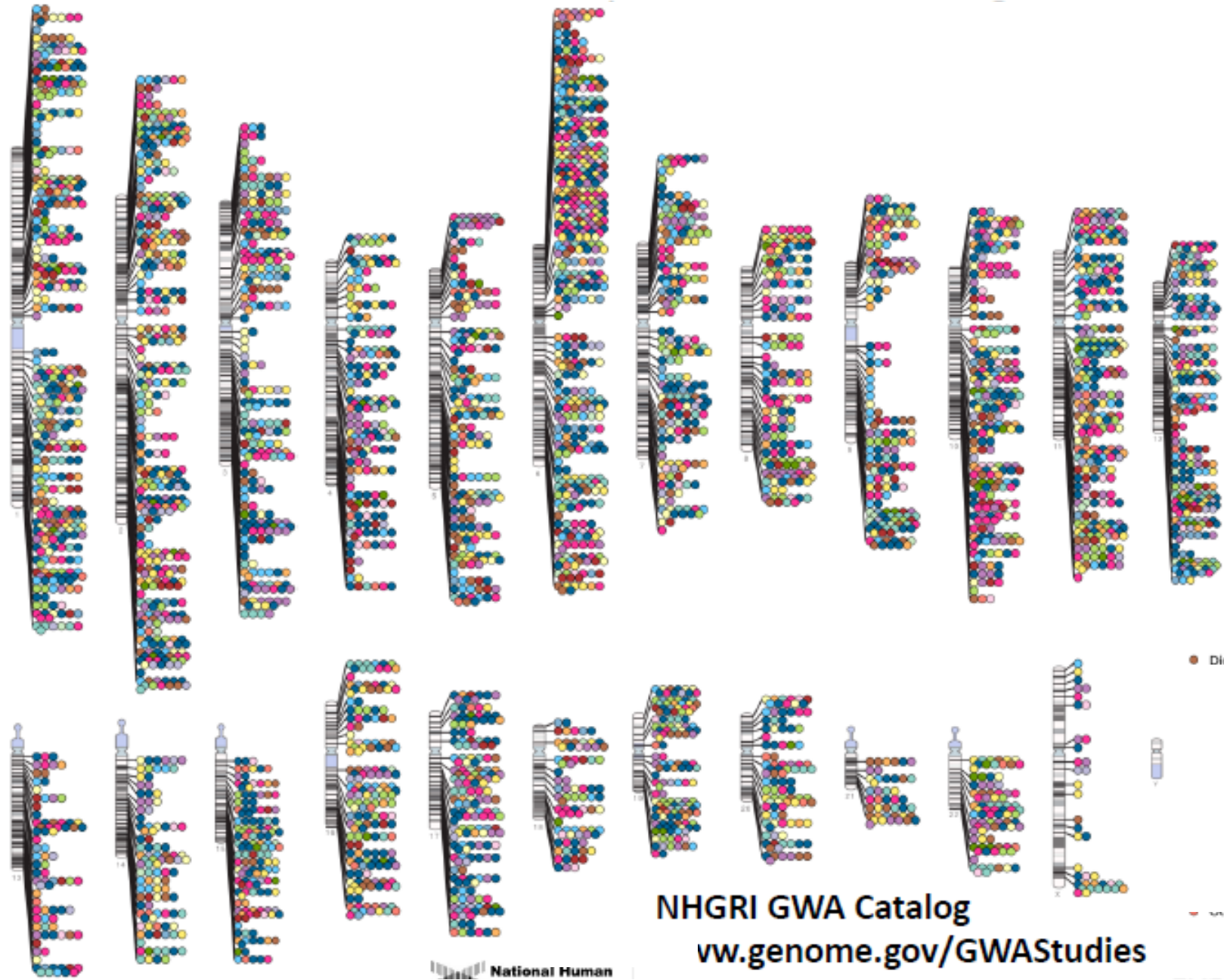
---

Effect of **common** genetic variants on gene expression variability







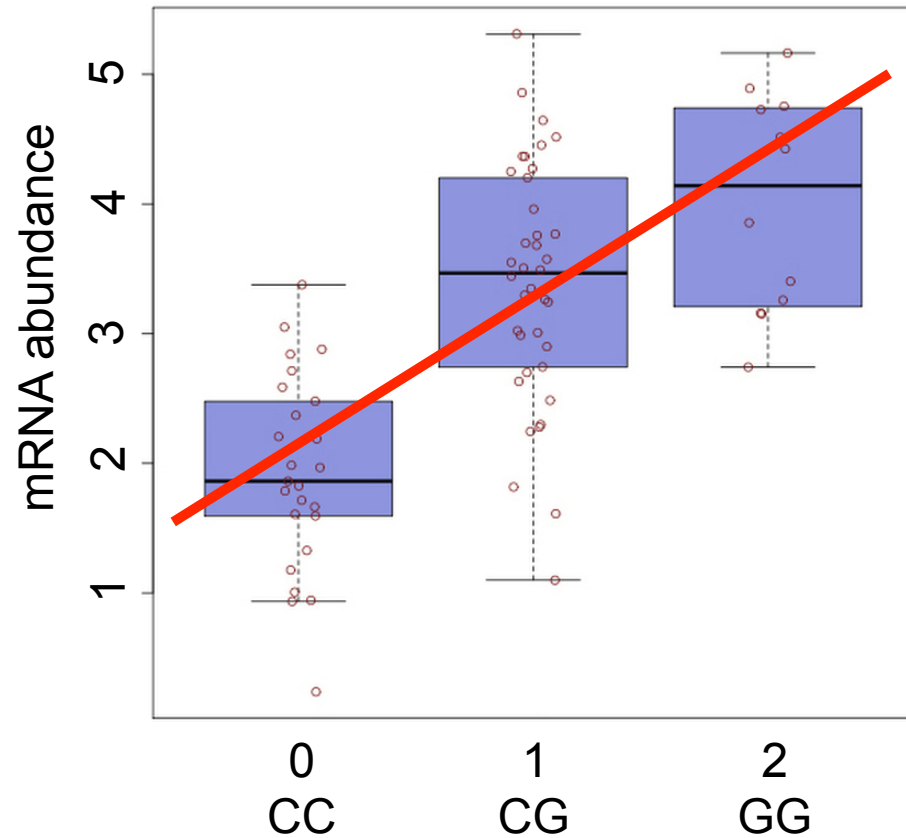


- Digestive system disease
- Cardiovascular disease
- Metabolic disease
- Immune system disease
- Nervous system disease
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurement
- Other measurement
- Response to drug
- Biological process
- Cancer
- Other disease
- Other trait

NHGRI GWA Catalog  
[www.genome.gov/GWAStudies](http://www.genome.gov/GWAStudies)

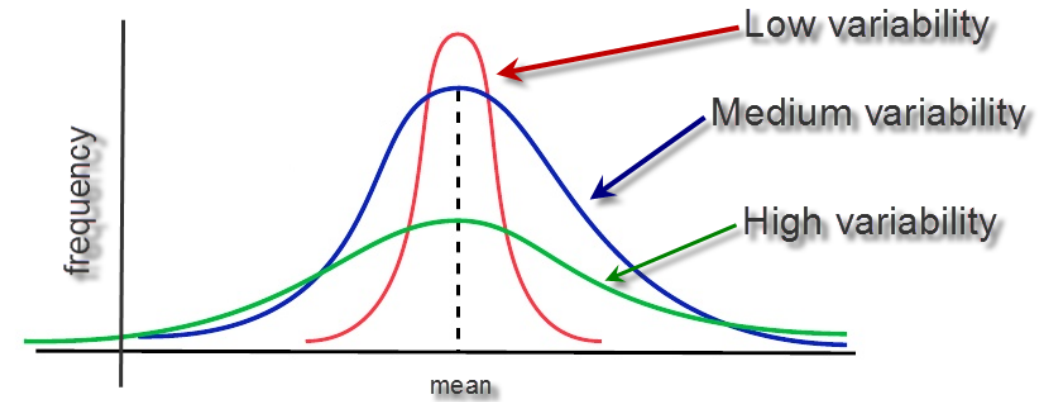
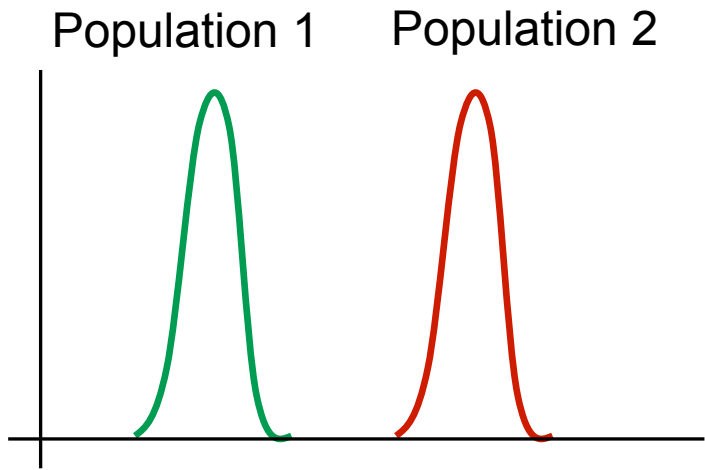
# Expression QTLs (eQTLs)

---



Gene expression level as an “intermediate phenotype”

$$y_i = \mu + x_i \beta + g_i \alpha + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$



# Variation vs. Variability

---



# New evidence: **phenotypic variability** (**variance**) is genetically controlled

---

*FTO* genotype is associated with phenotypic variability of body mass index (Yang *et al. Nature* 2012)

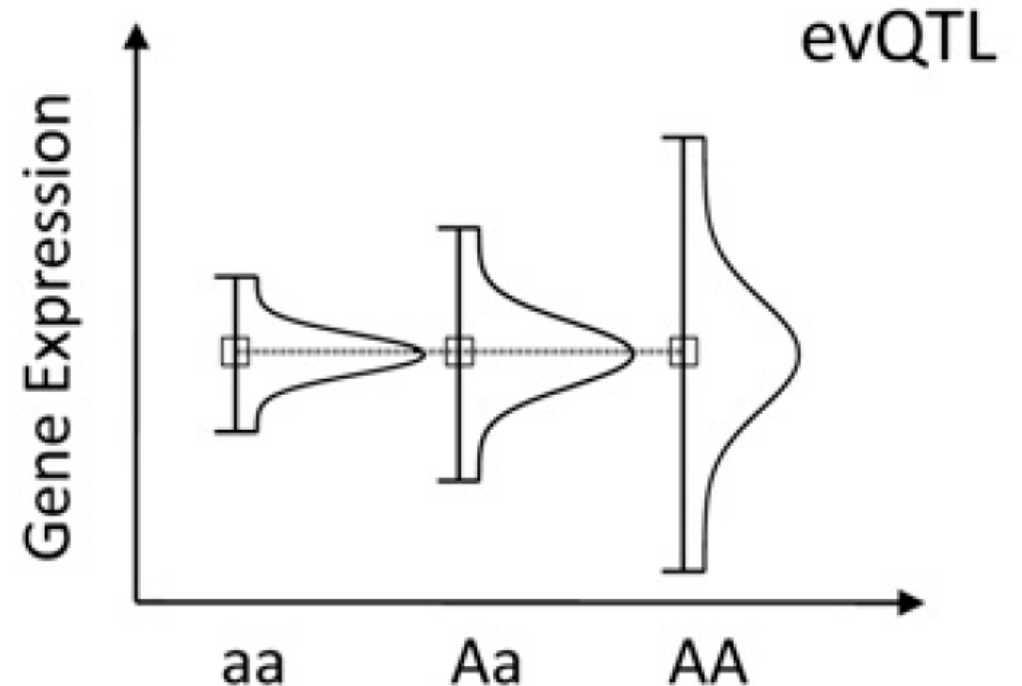
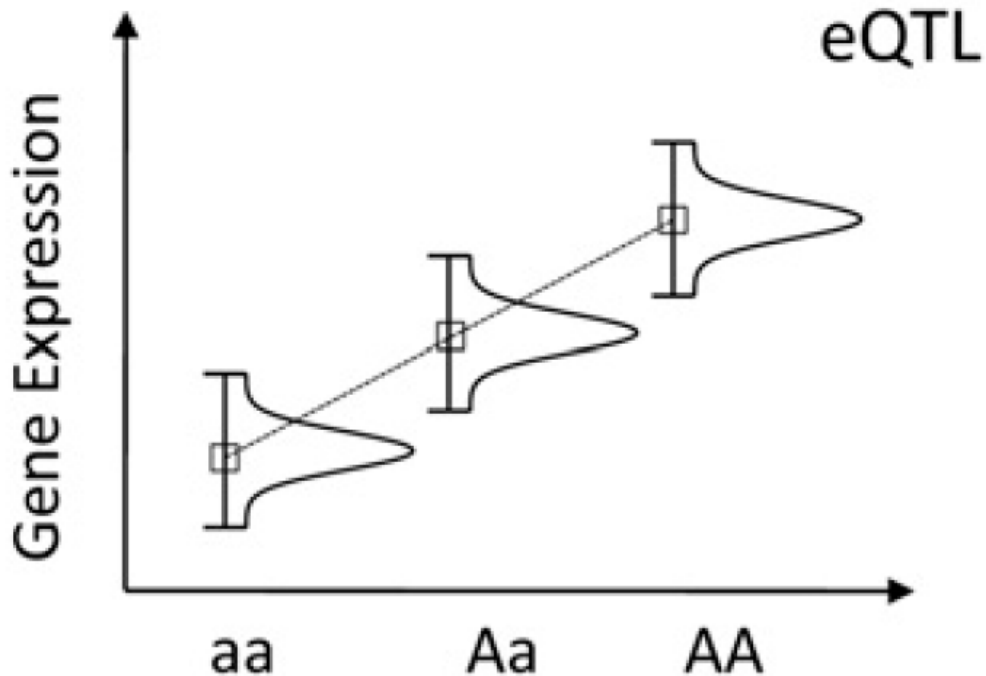
Inheritance beyond plain heritability: variance-controlling genes in *Arabidopsis thaliana* (Shen *et al. PLoS Genet* 2012)

Behavioral idiosyncrasy reveals genetic control of phenotypic variability (Julien *et al. PNAS* 2015)

Selection on noise constrains variation in a eukaryotic promoter (Metzger *et al. Nature* 2015)

# Expression variability QTL – **evQTL**

i.e., genetic loci linked to or associated with expression variance





# Detection of evQTLs

---

## Linear regression model

$$y_i = \mu + x_i \beta + g_i \alpha + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$$

## Double generalized linear model (DGLM)

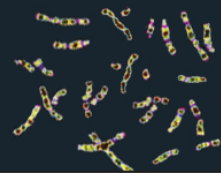
$$y_i = \mu + x_i \beta + g_i \alpha + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2 \exp(g_i \theta))$$

# Genome scan for evQTLs

## Data Sets:

1. Genotype data from the **1000G** project

**1000 Genomes**  
A Deep Catalog of Human Genetic Variation



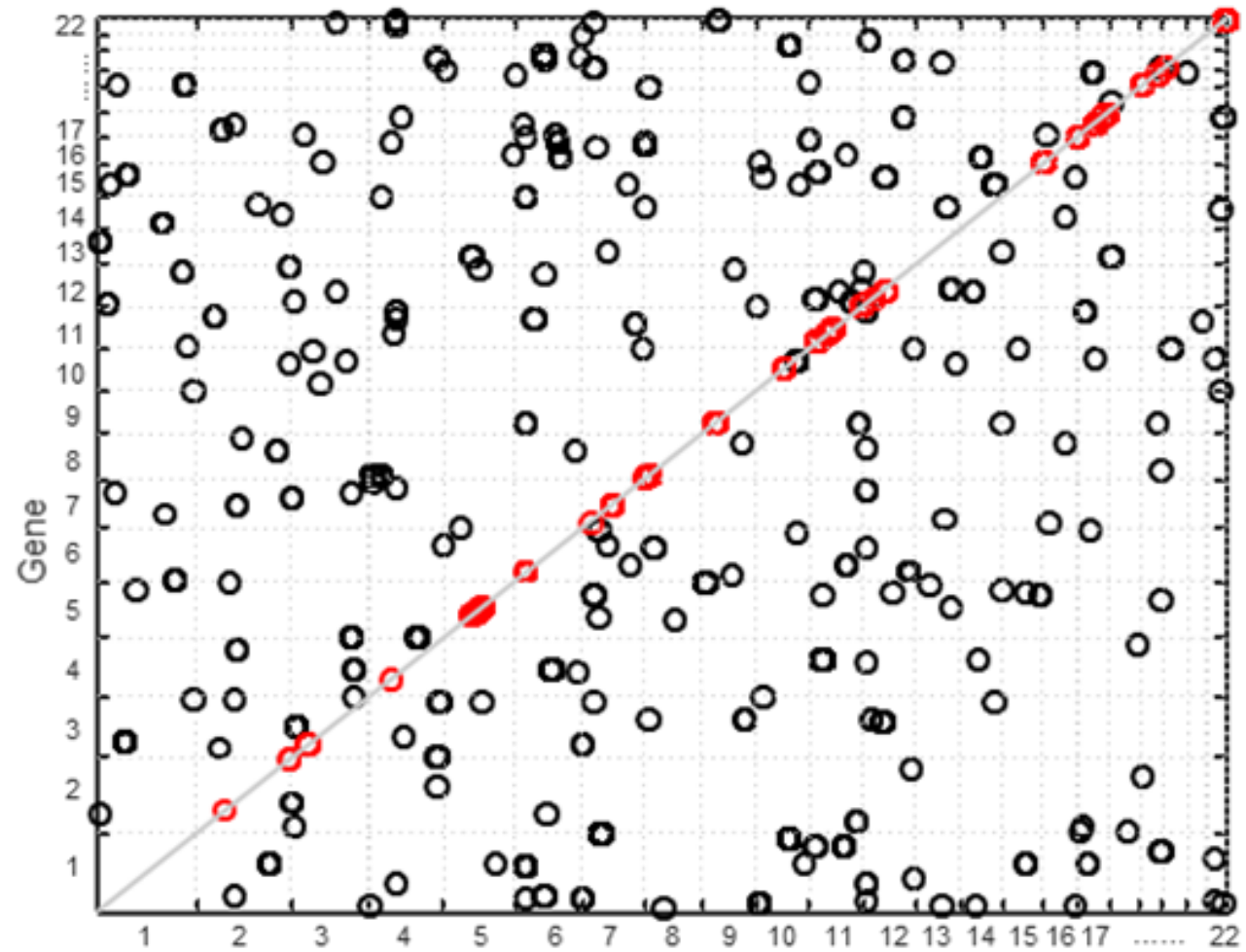
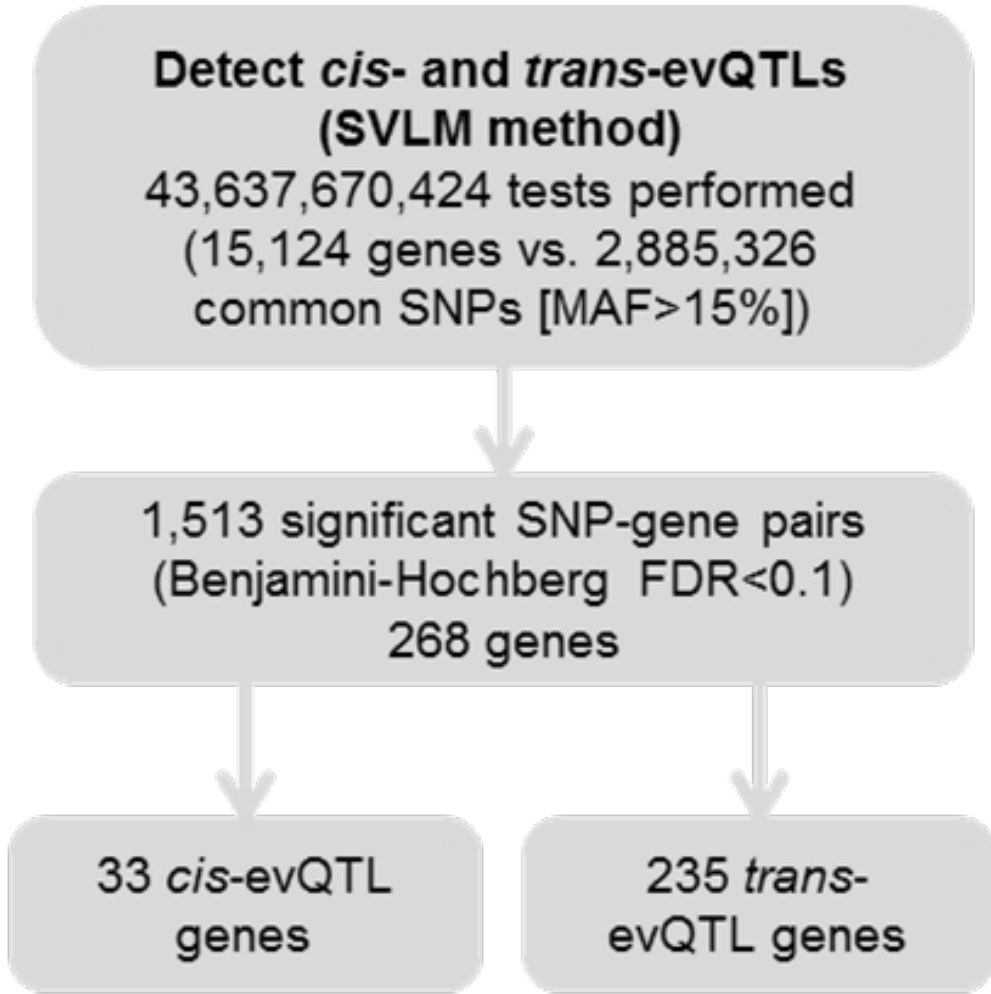
2. RNA-seq data from the **Geuvadis** project



**Detect *cis*-evQTLs  
(DGLM method)**

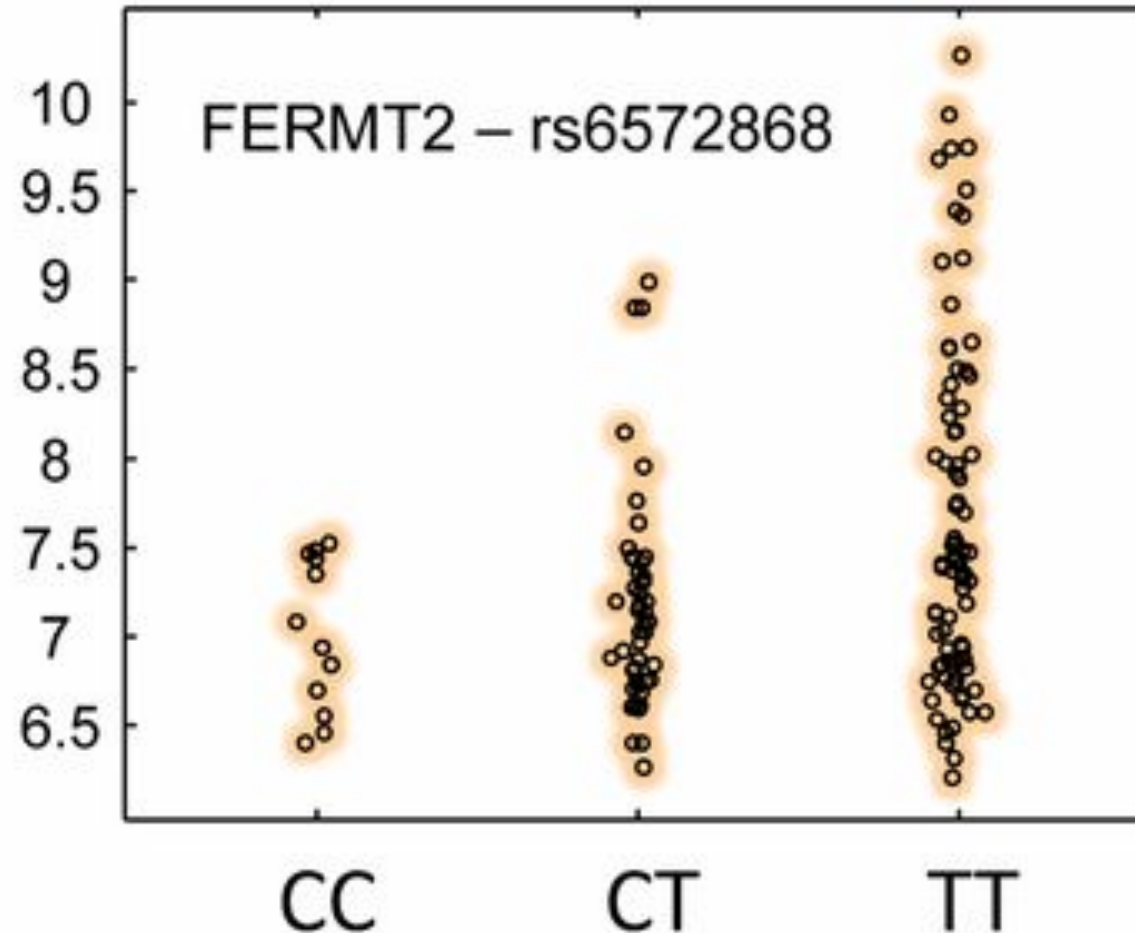
28,494,473 tests performed  
(15,124 genes vs. common SNPs  
[MAF>15%])

17,949 significant SNP-gene pairs  
(Bonferroni corrected  $P < 0.05$ )  
1,304 genes



# Expression variability QTL – **evQTL**

i.e., genetic loci linked to or associated with expression variance

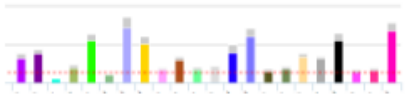




## Current Release

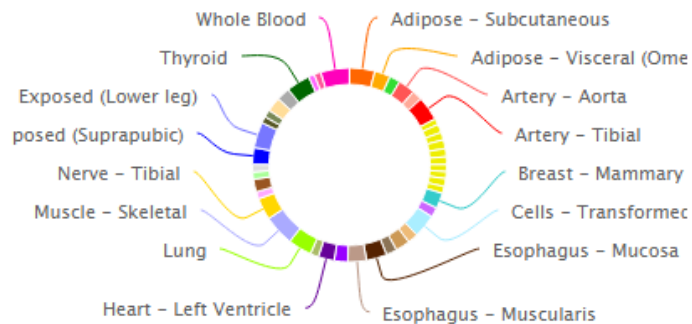
Latest Version: V6p

[Dataset Summary Statistics Report](#)



Browse eQTL Tissues

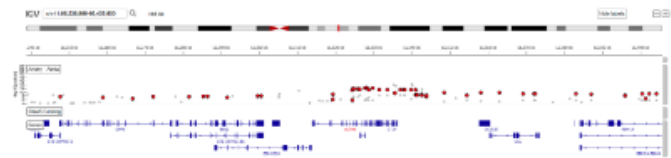
Total samples in all eQTL tissues: 7051



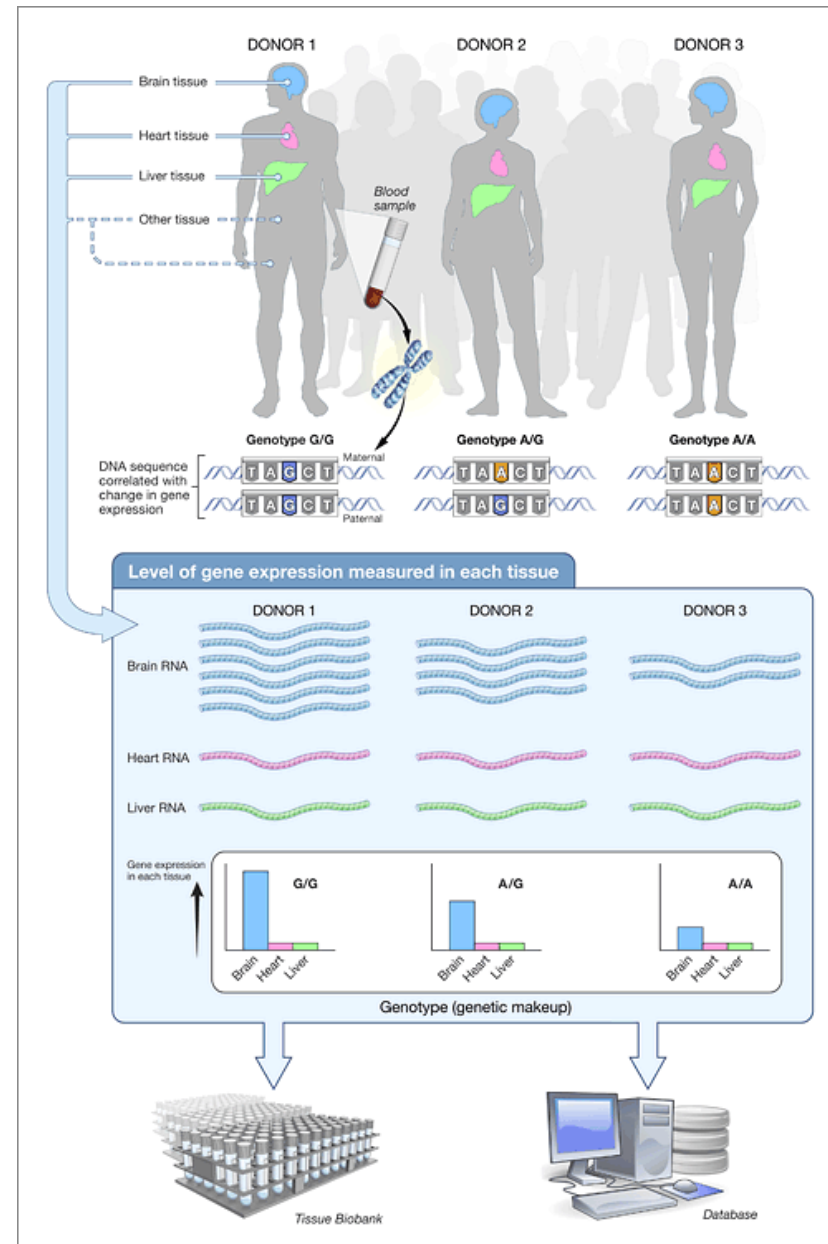
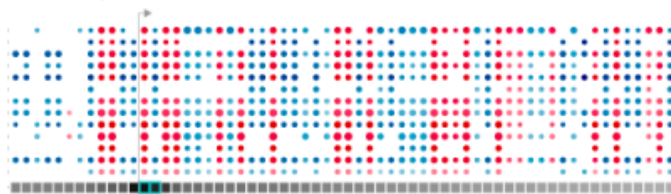
## Genetic Association

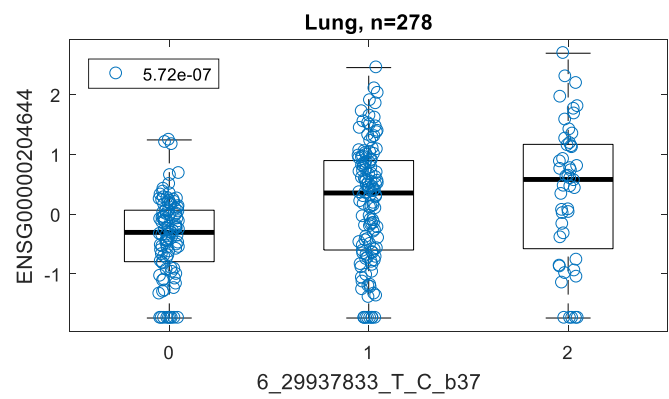
Single Tissue eQTLs

eQTL IGV Browser

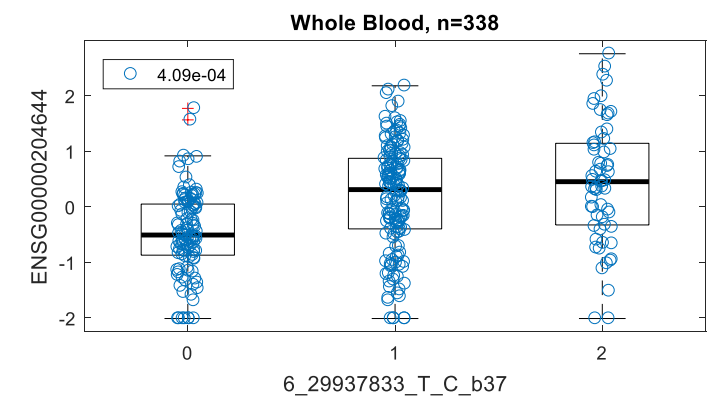
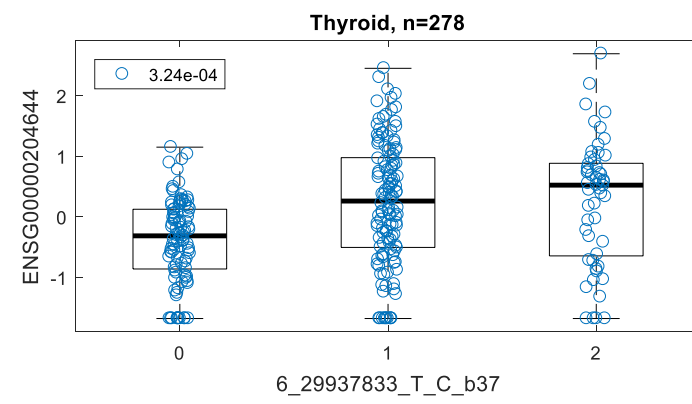
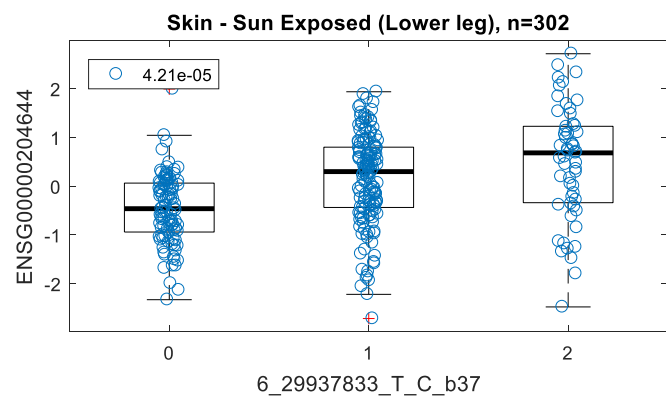
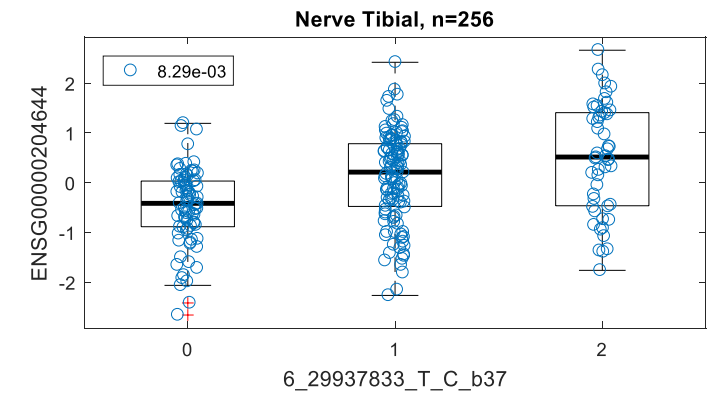
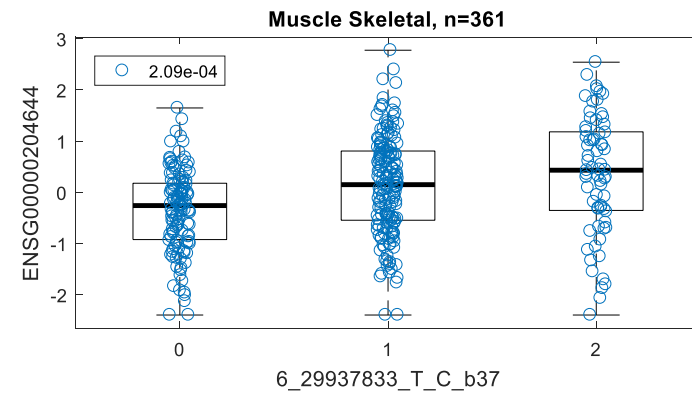
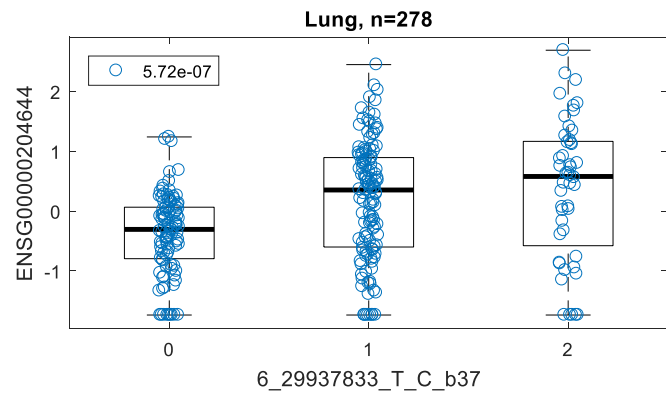
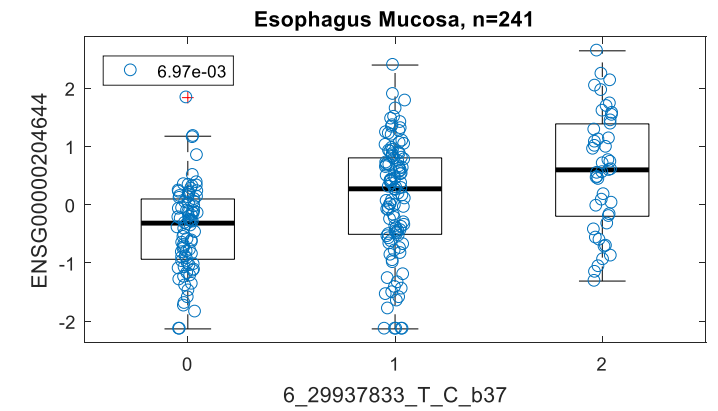
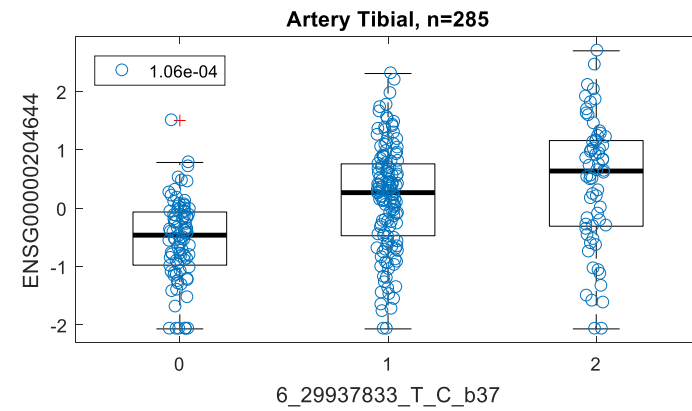
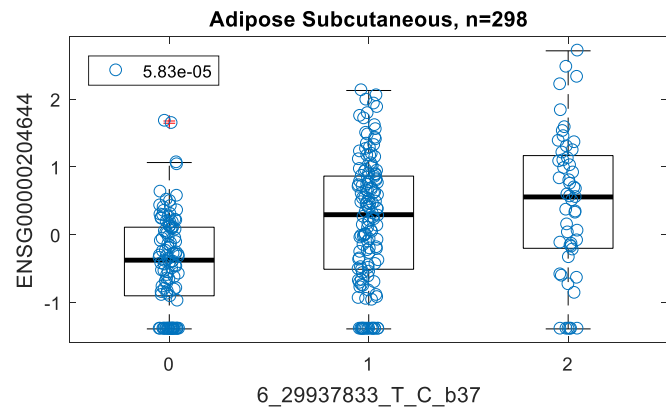


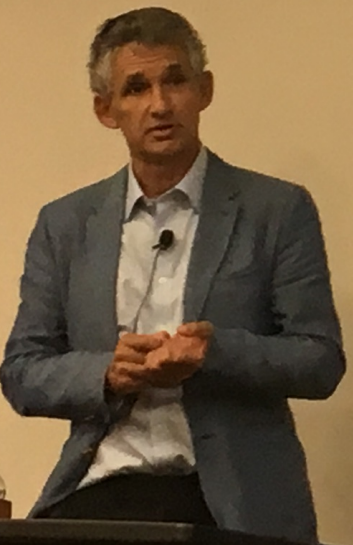
Gene eQTL Visualizer









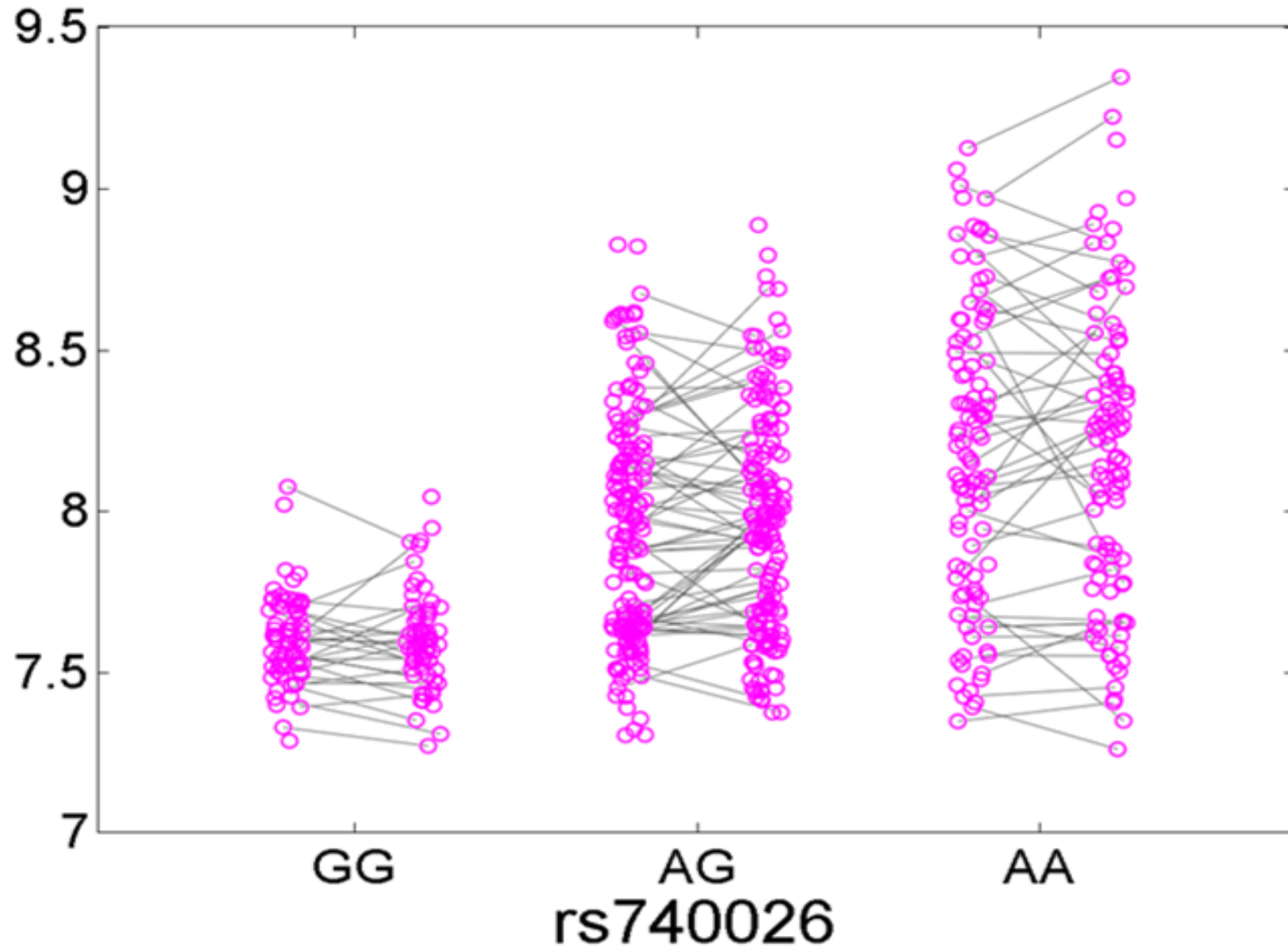


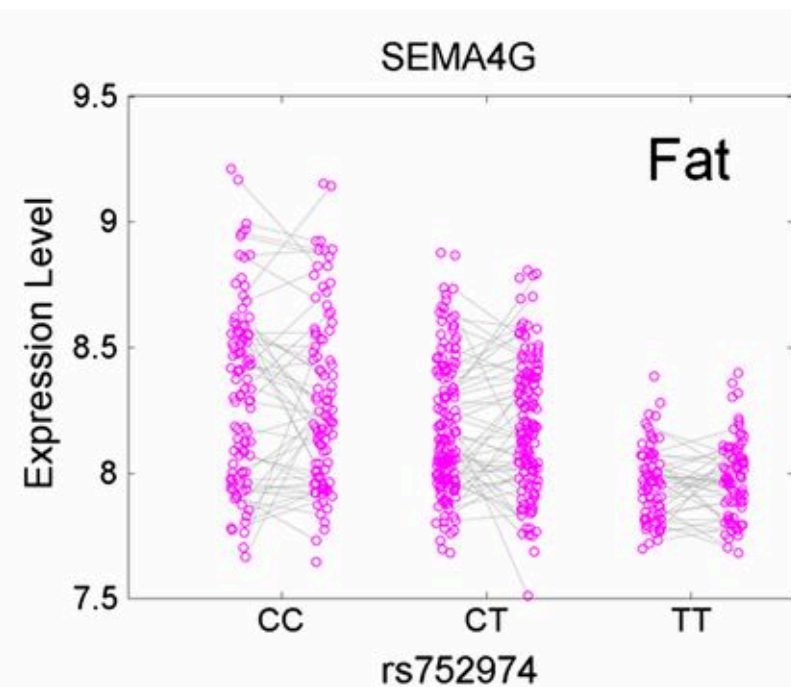
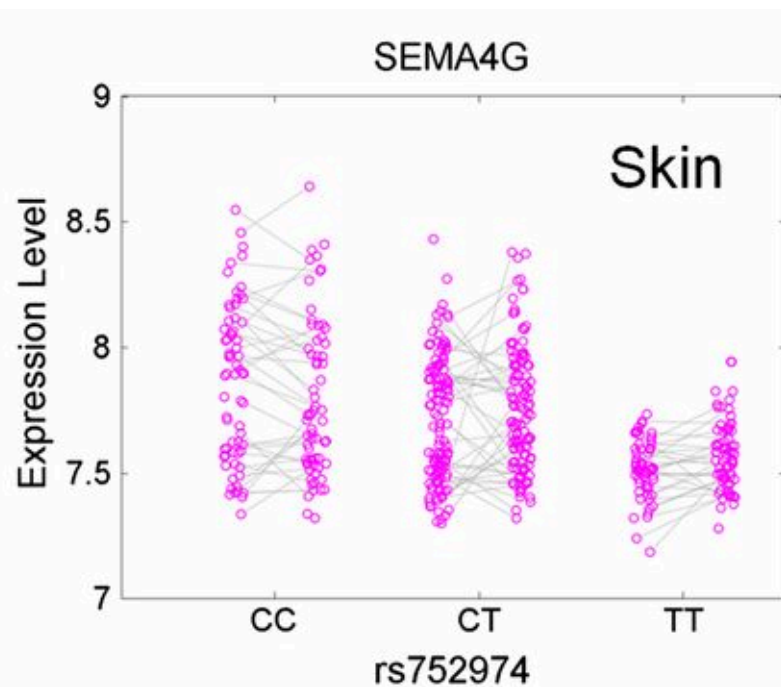
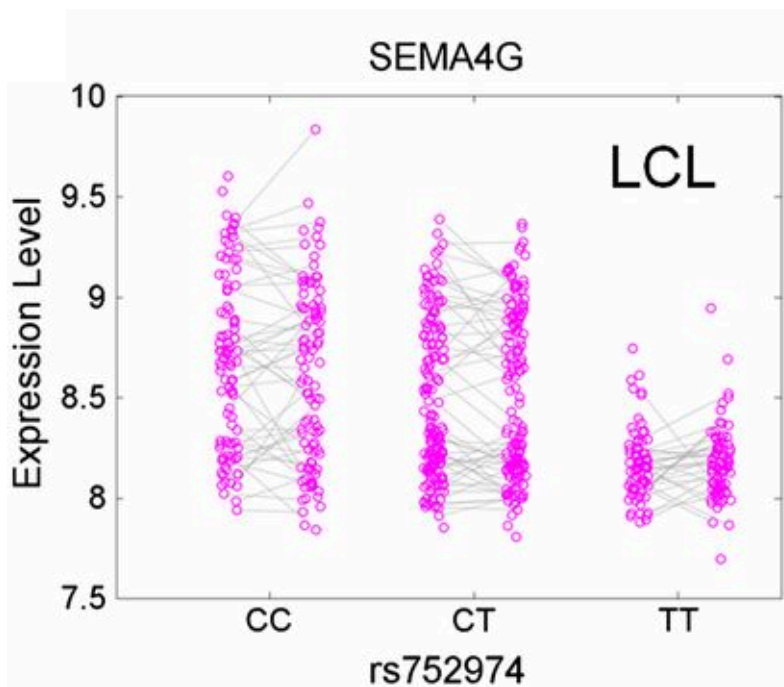
A large grid of diverse human faces is projected on a screen. The faces are arranged in many rows and columns, showing a wide variety of ethnicities, ages, and genders. Overlaid on the center of the grid is a semi-transparent purple banner with the text "TwinsUK" in large white letters, and below it, "The best studied people on the planet" in smaller white letters. To the left of the grid, there is a vertical toolbar with various icons for navigation and interaction, including arrows, a magnifying glass, and a search icon.

Tim Spector



# AXIN2





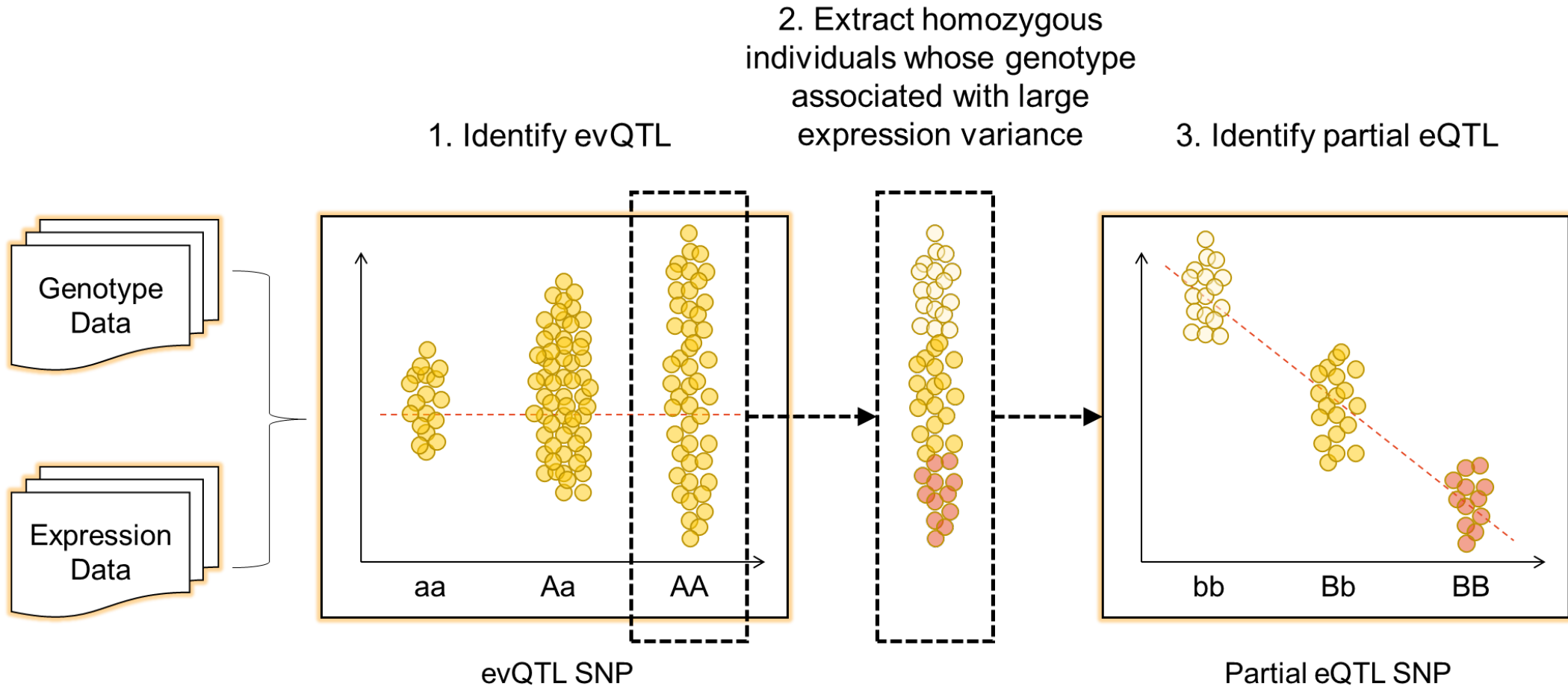
# Two distinct models explaining the creation of evQTLs

---

**GxG (epistasis):** the interaction between genotypes

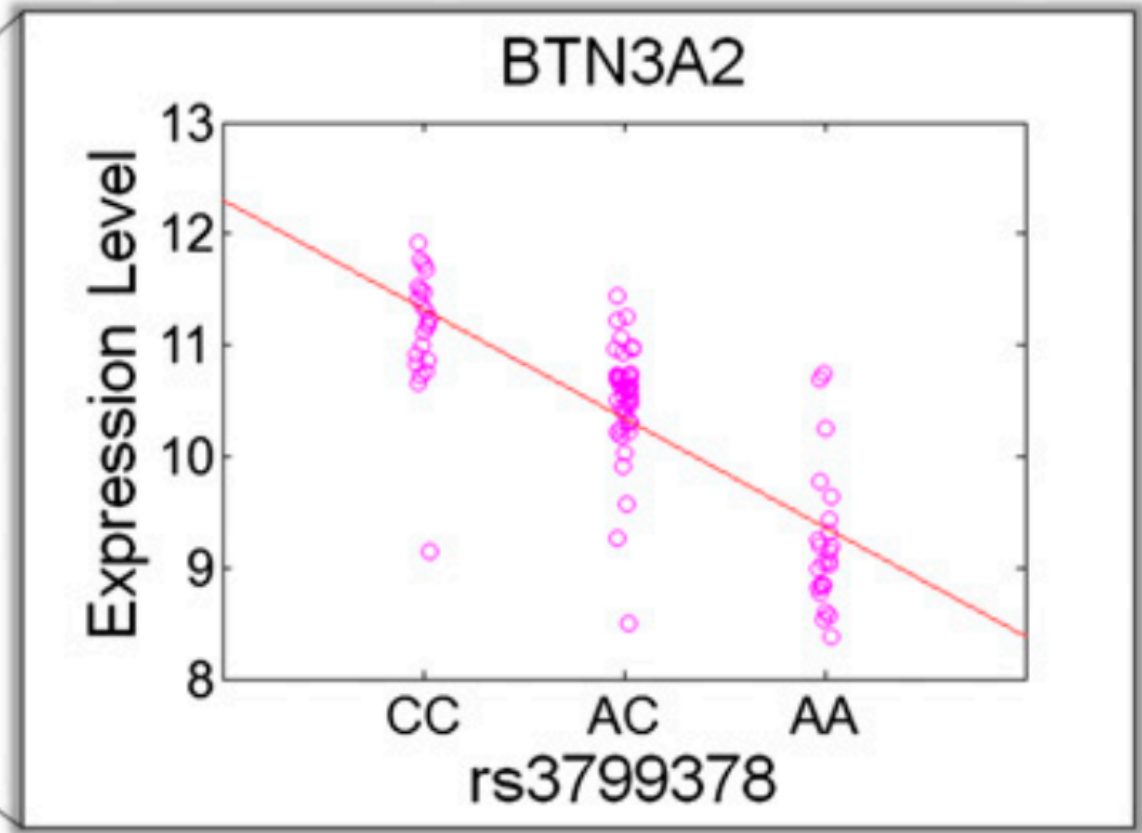
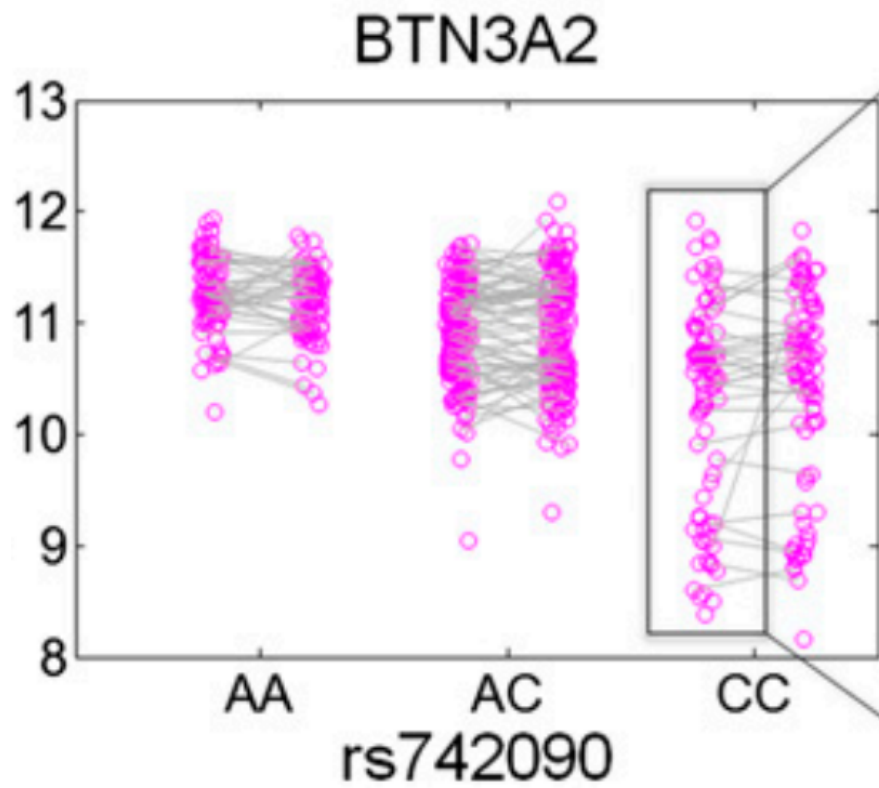
**GxE (destabilization):** the interaction between genotype and environment

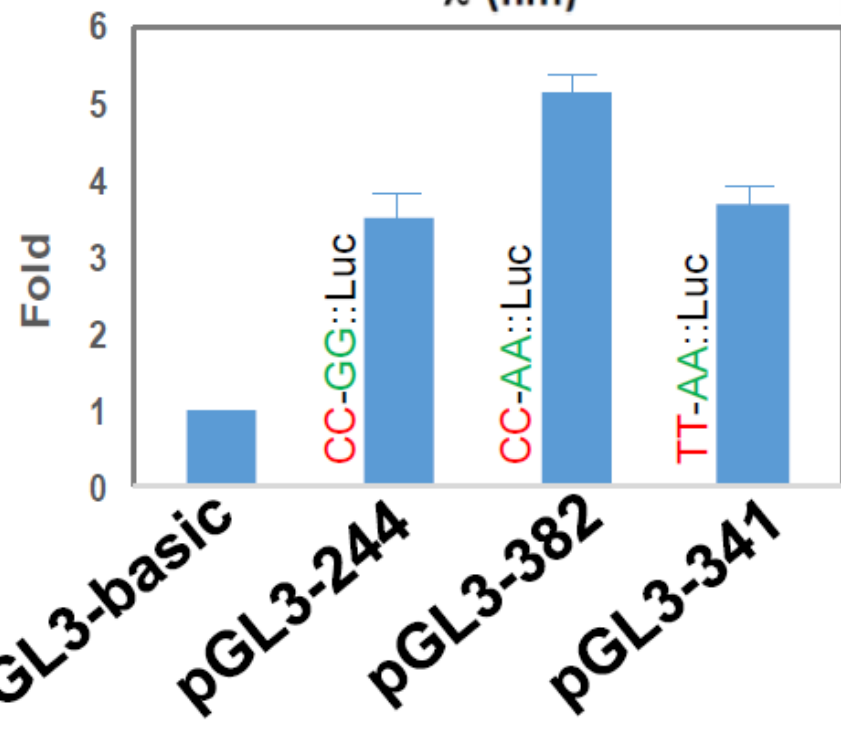
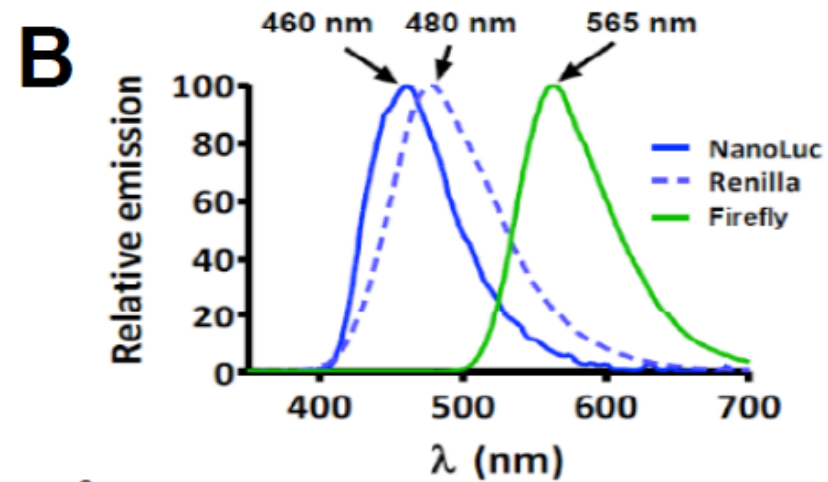
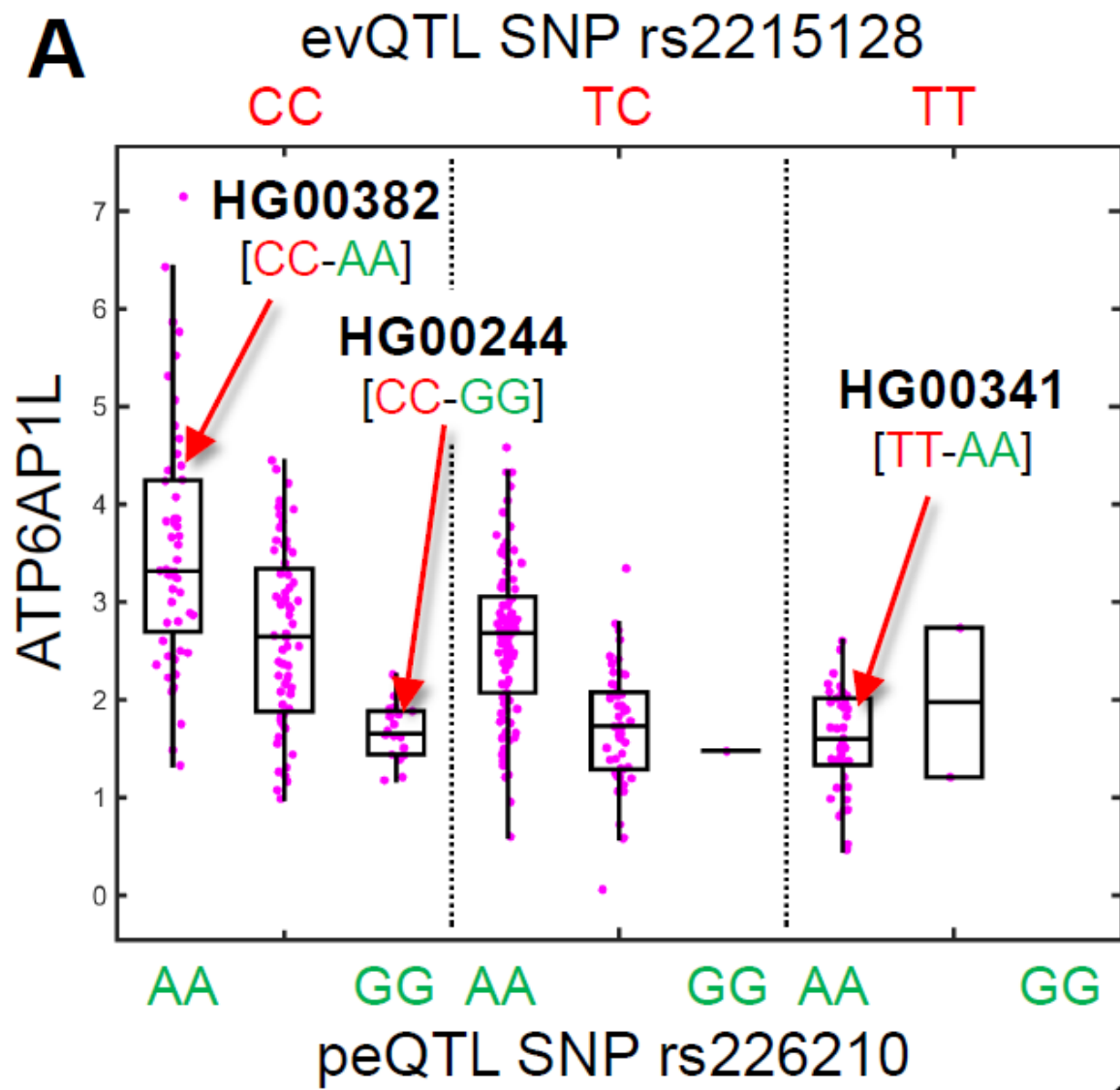
# GxG (epistasis) model





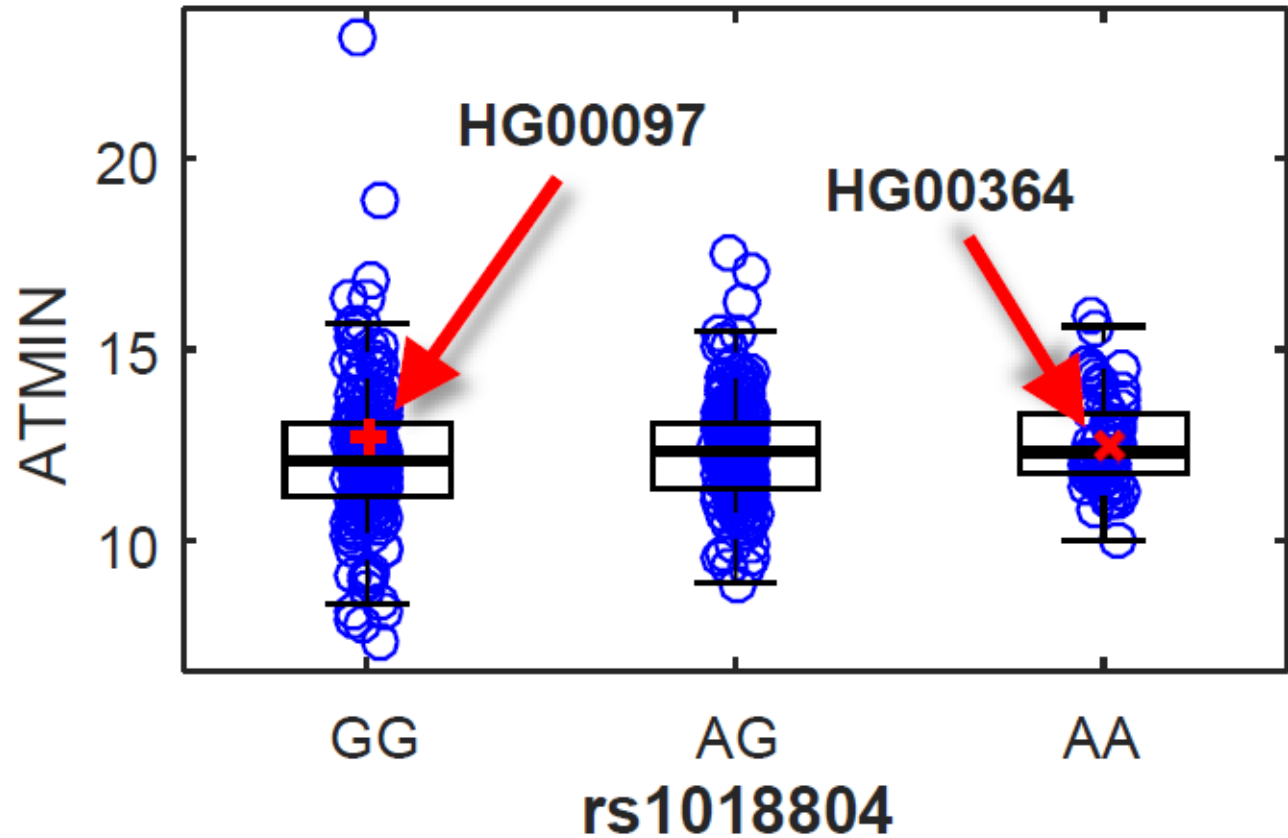
# GxG (epistasis) model

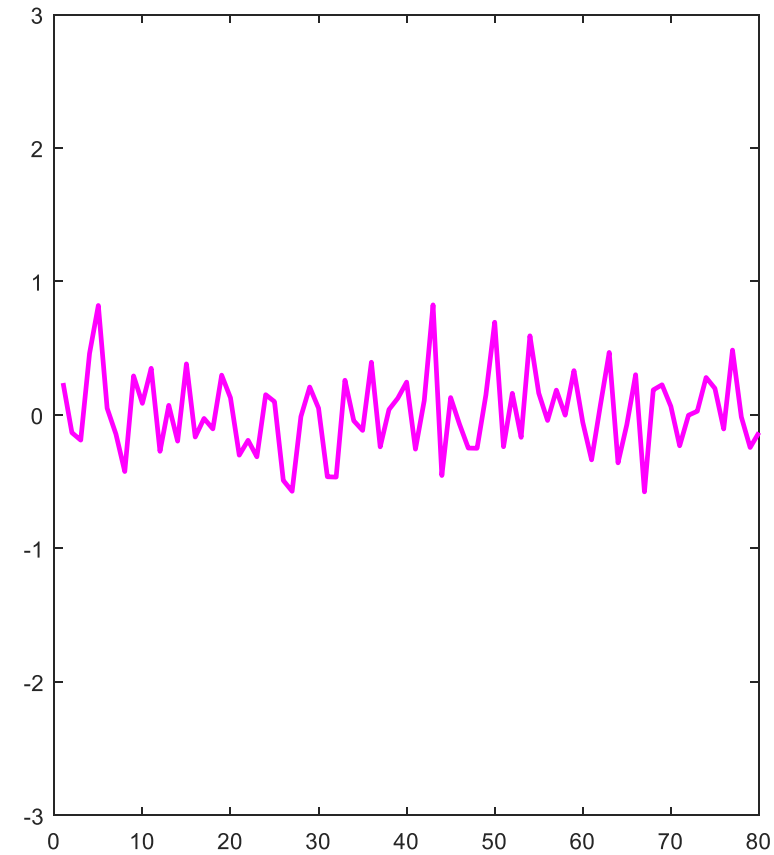
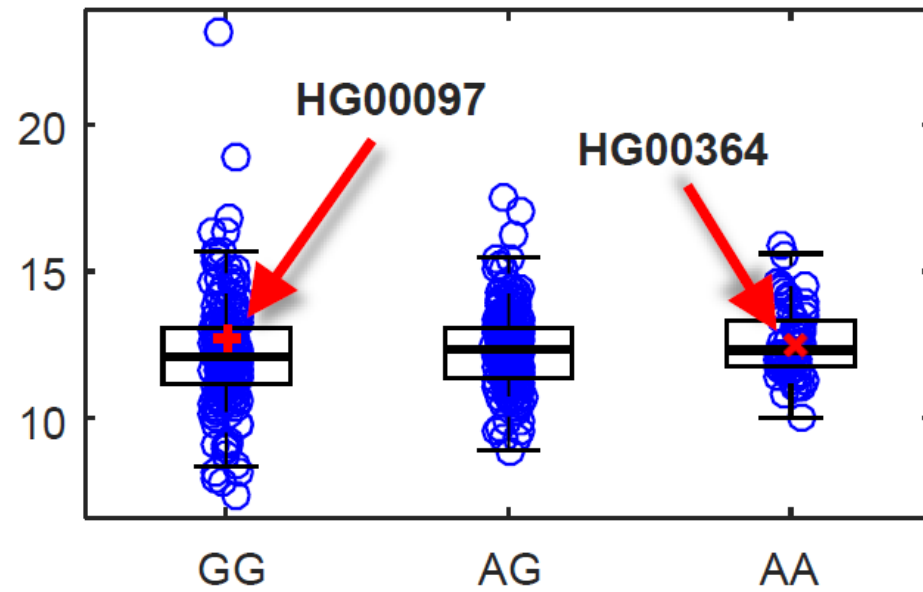
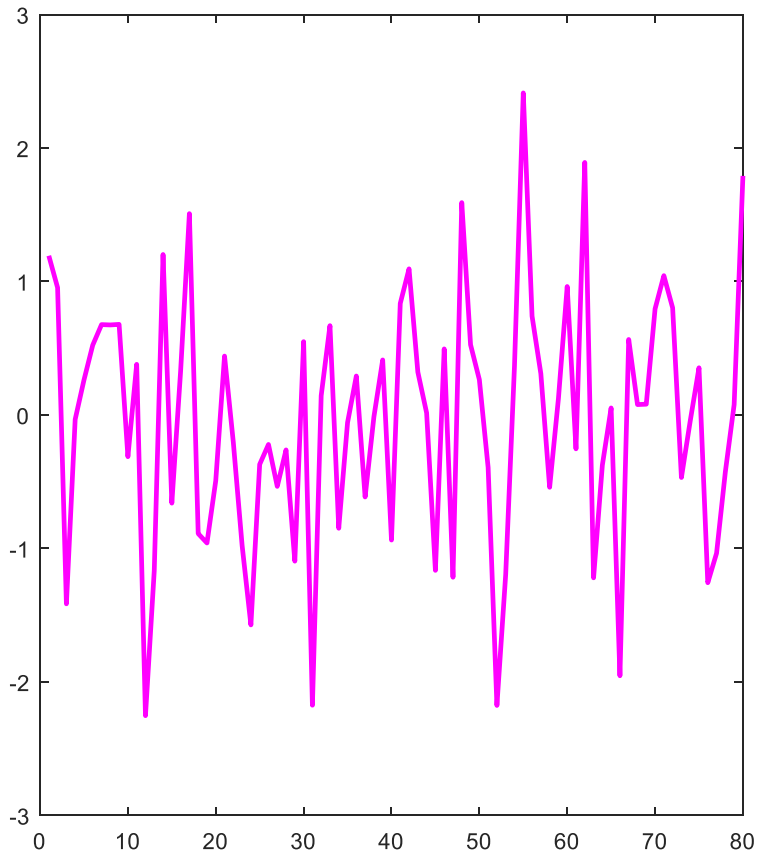




# GxE (destabilization) model – repetitive qPCR

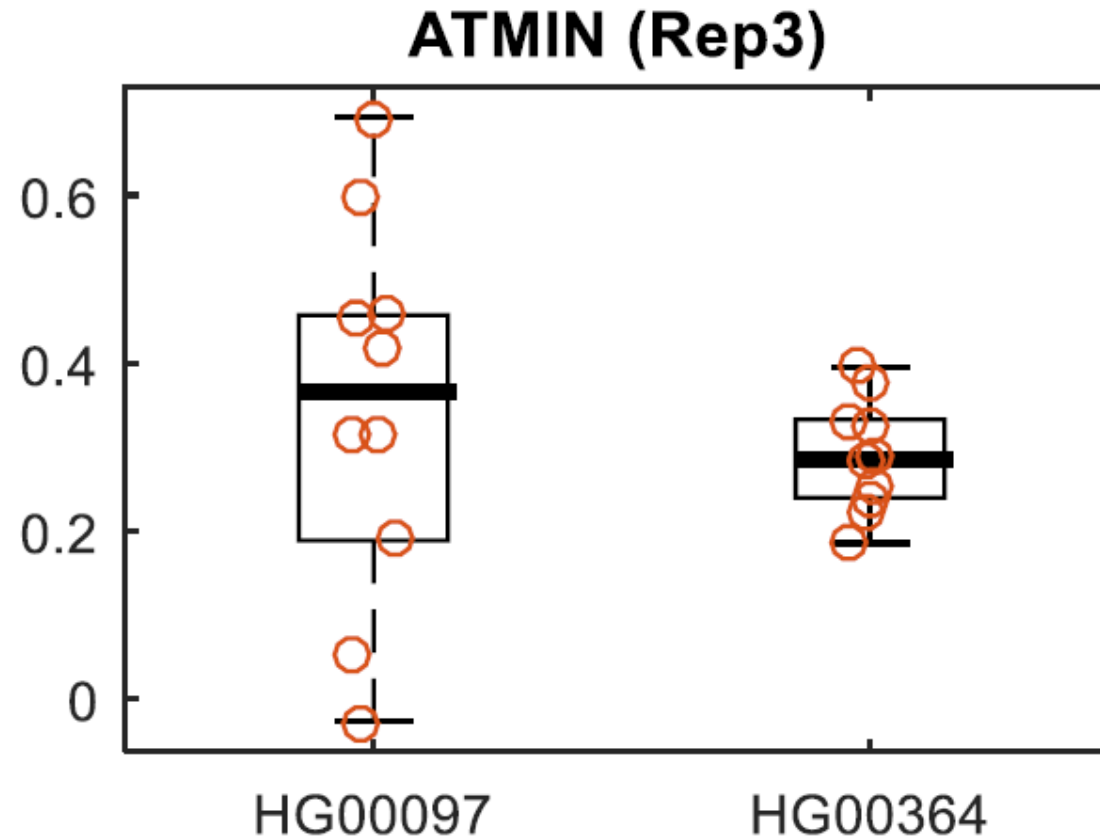
Select **two cell lines** from groups with **large** and **small** expression variability.



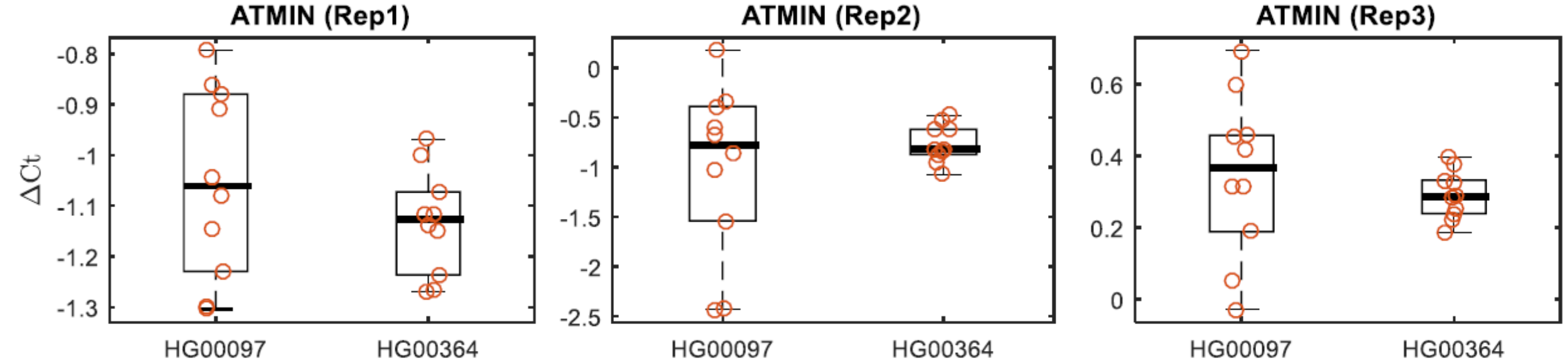


# GxE (destabilization) model – repetitive qPCR

qRT-PCR assay was repeated 10 times for each sample.

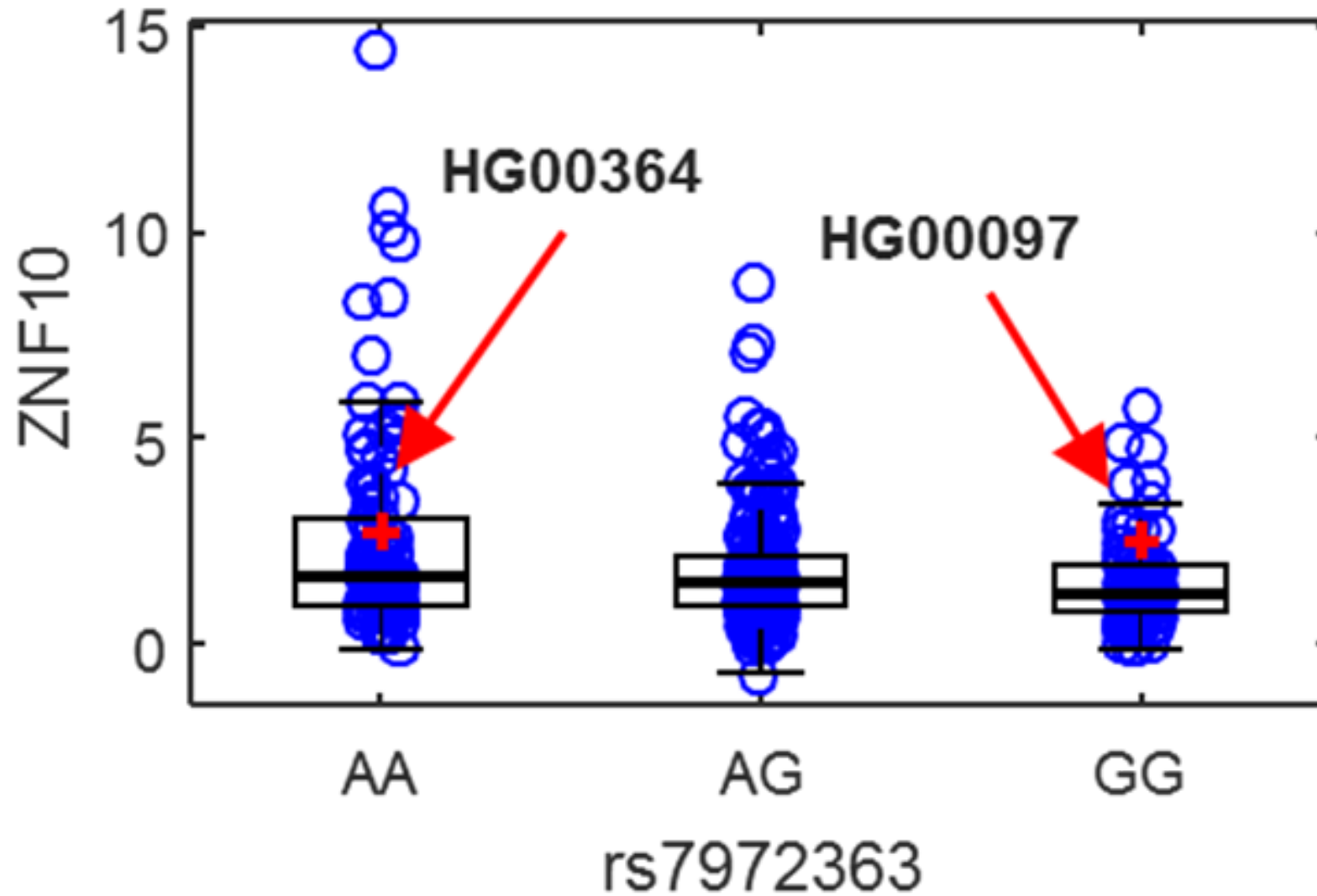


# GxE (destabilization) model – repetitive qPCR

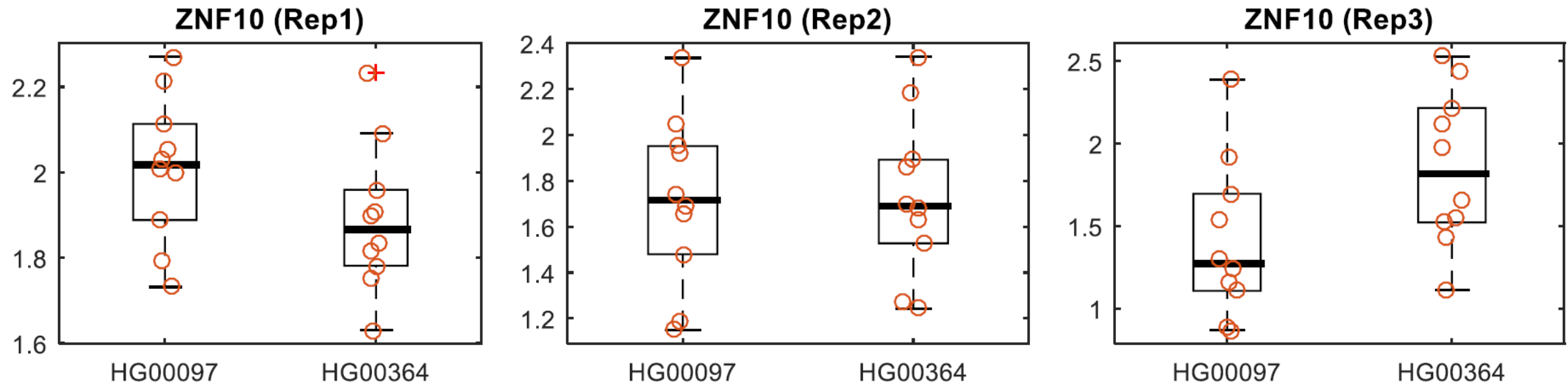




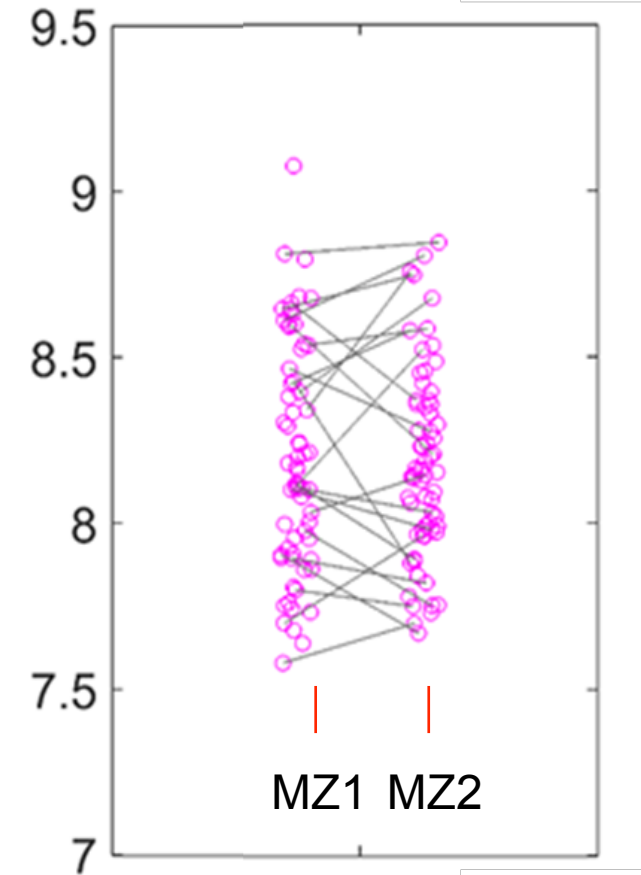
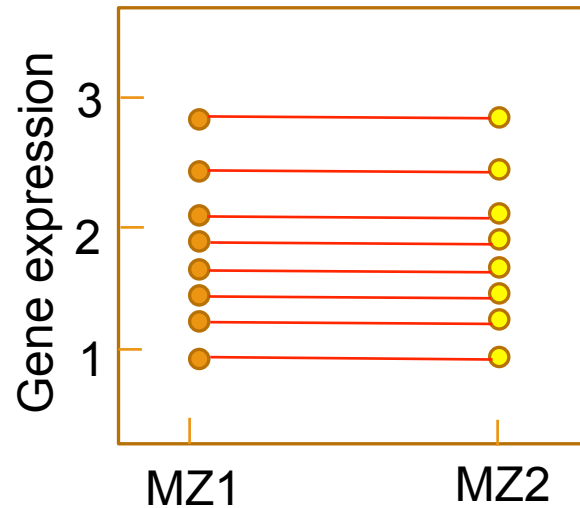
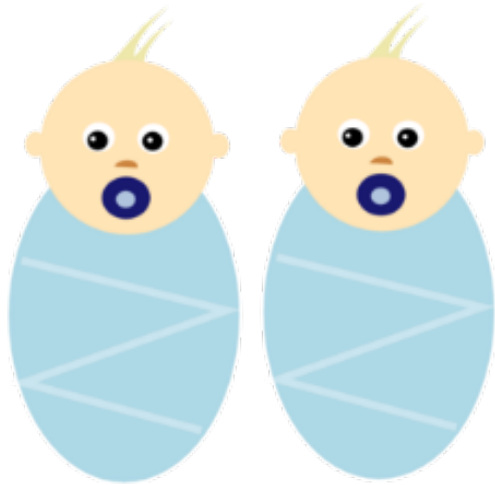
An evQTL explained by the **GxG** (epistasis) model



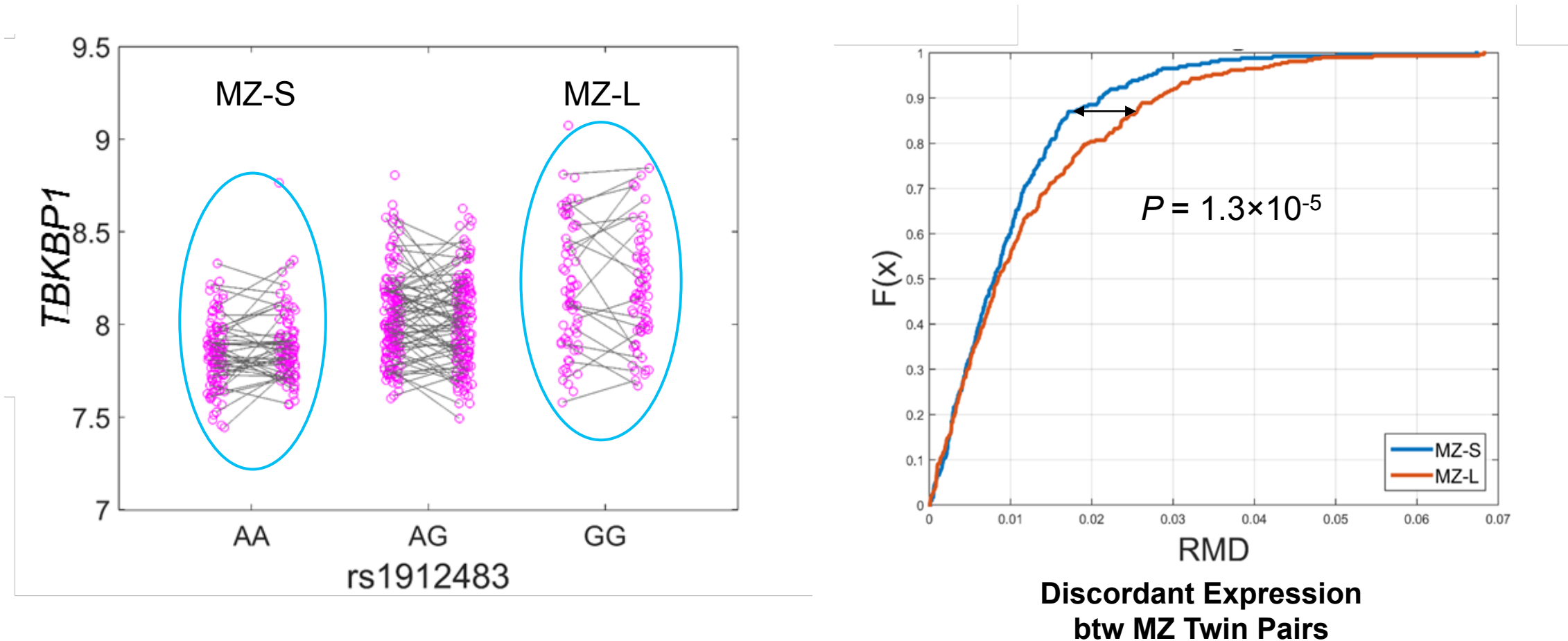
# An evQTL explained by the **GxG** (epistasis) model



# GxE (destabilization) model – discordant expression between monozygotic (MZ) twins



# GxE (destabilization) model – discordant expression between monozygotic (MZ) twins



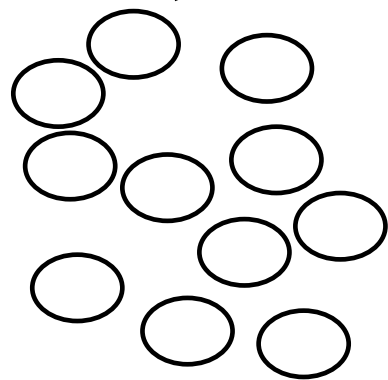
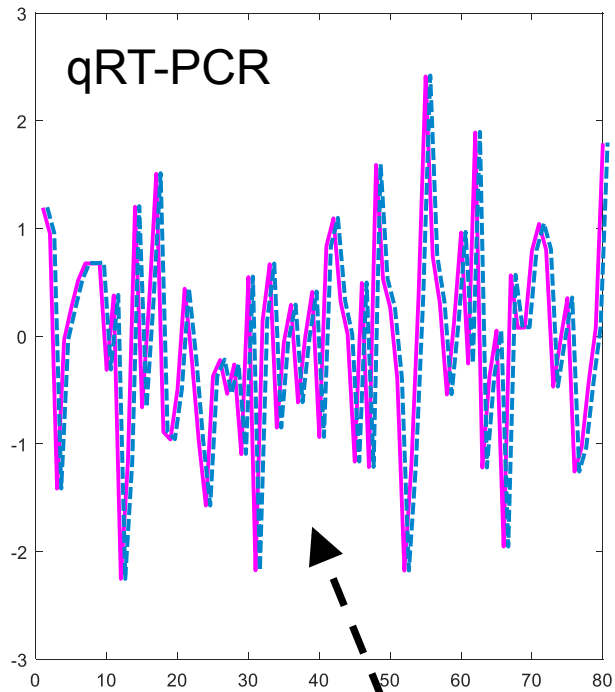
## Future plans

---

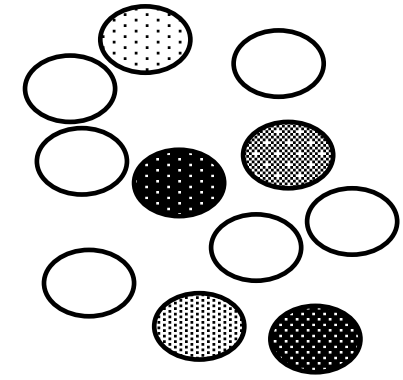
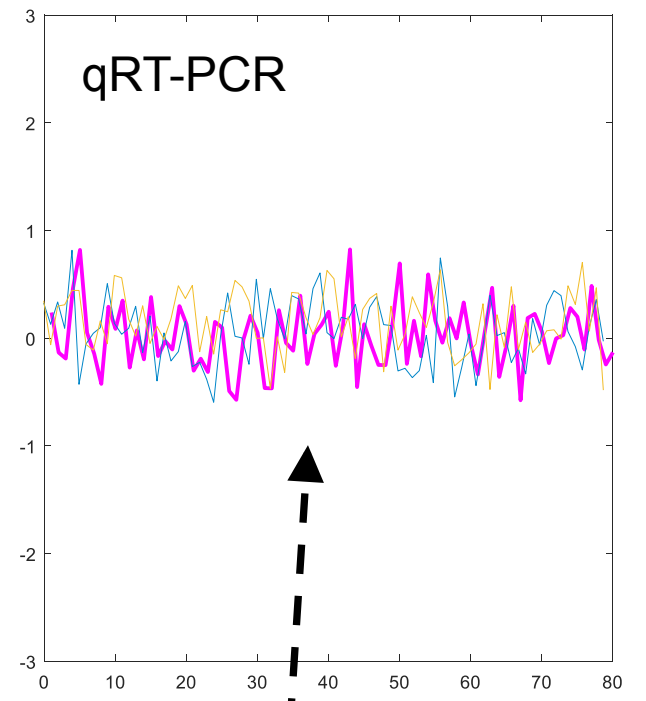
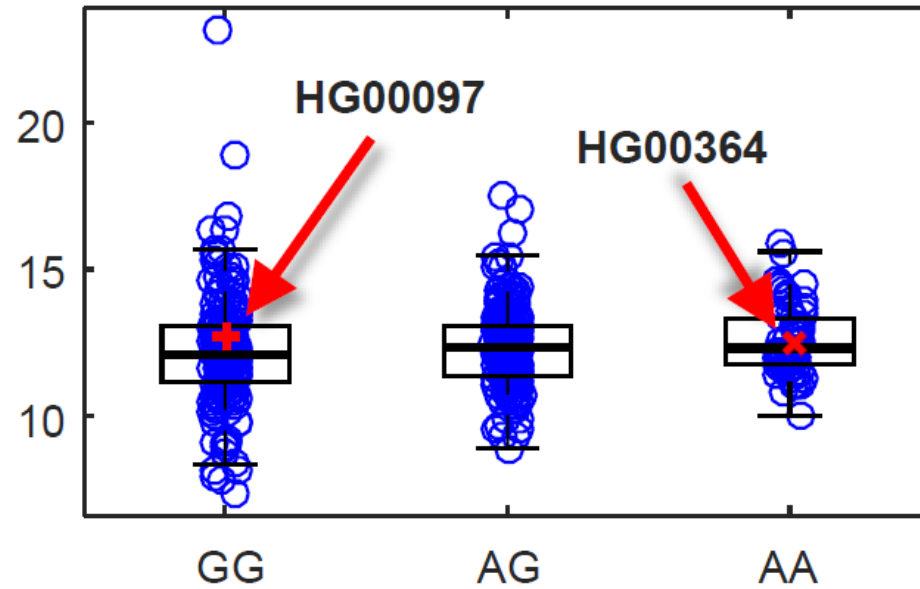
**Circadian rhythm** gene expression analysis (D. Earnest)

**Single-cell** gene expression analysis (A. Raj)

**CRISPR/Cas9**-based gene editing (D. Segal)



Single cells



Single cells

# Summary

---

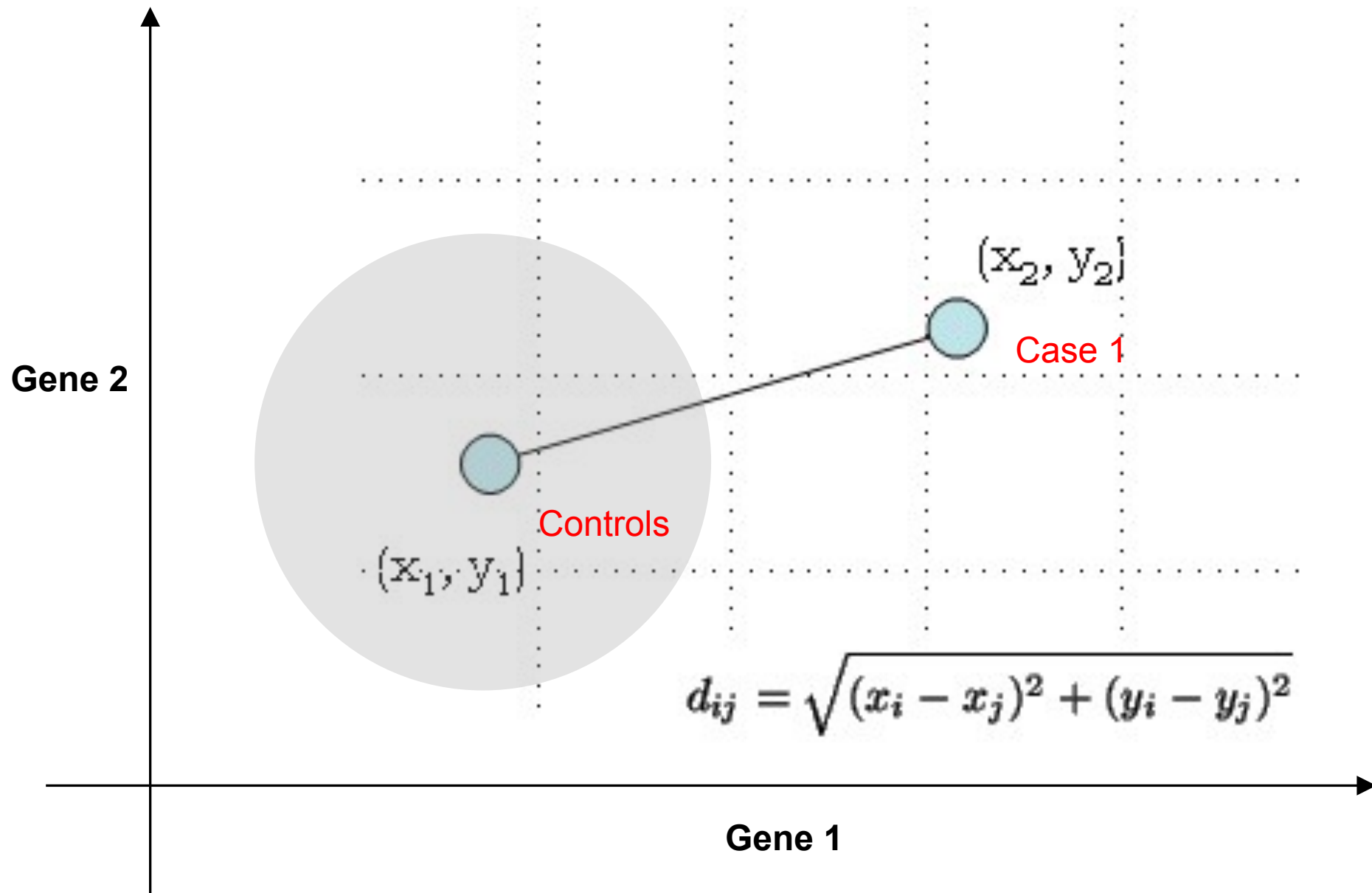
- Two distinct modes of action — **epistasis** and **destabilization**.
  - Genetic variants work either interactively (GxG) or independently (GxE) to influence gene expression variance.

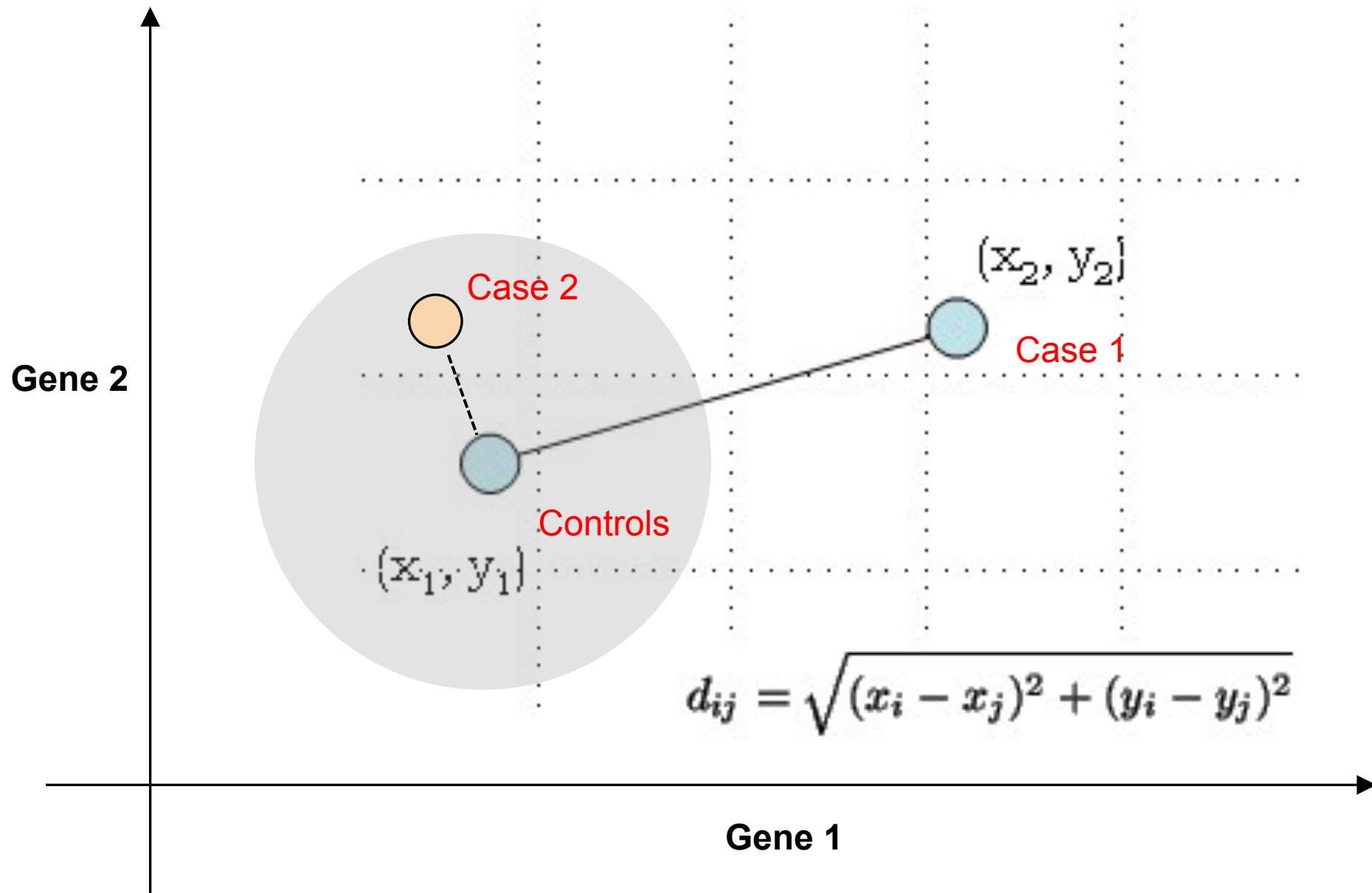
# Exploiting aberrant gene expression in autism for discovery and diagnosis

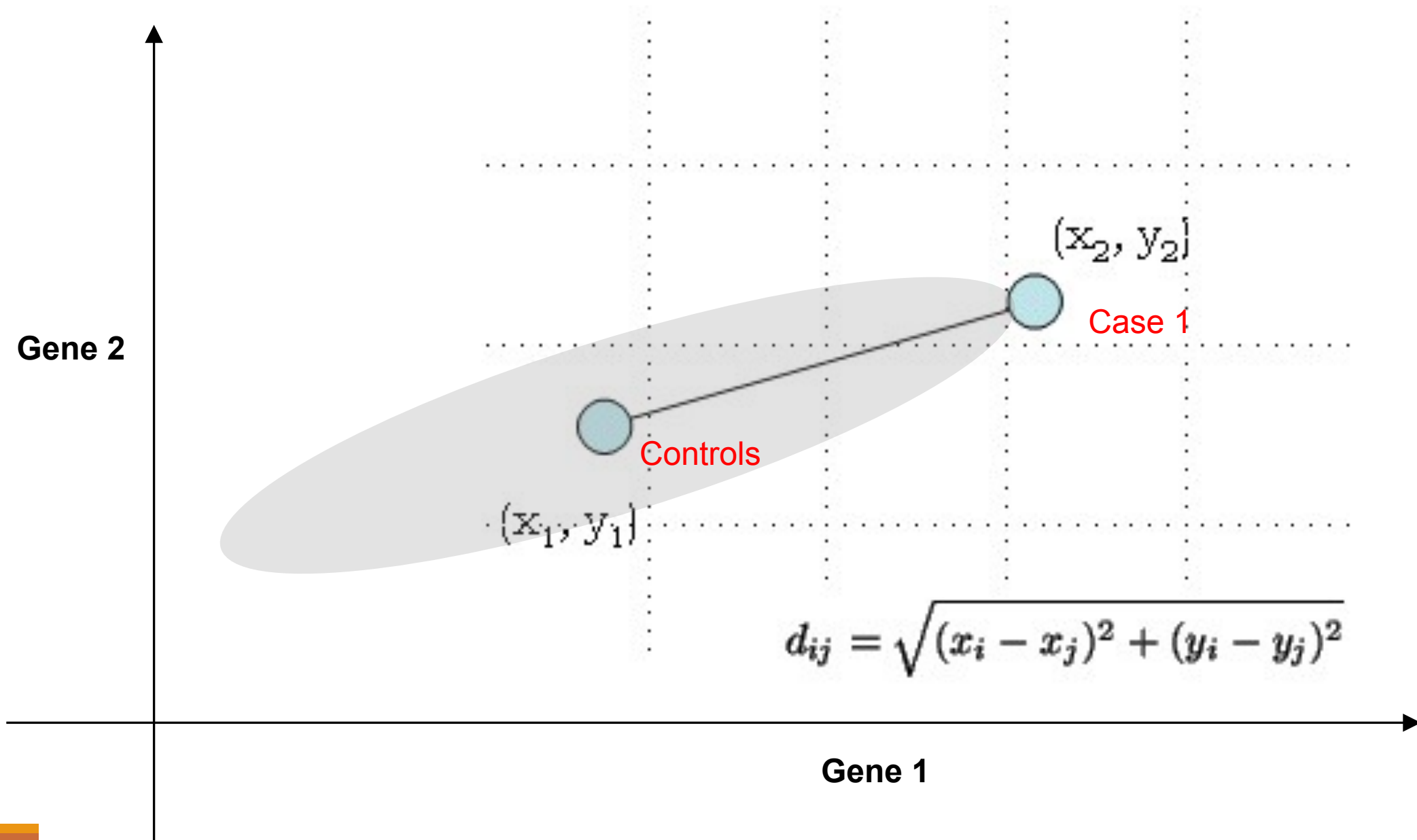
---

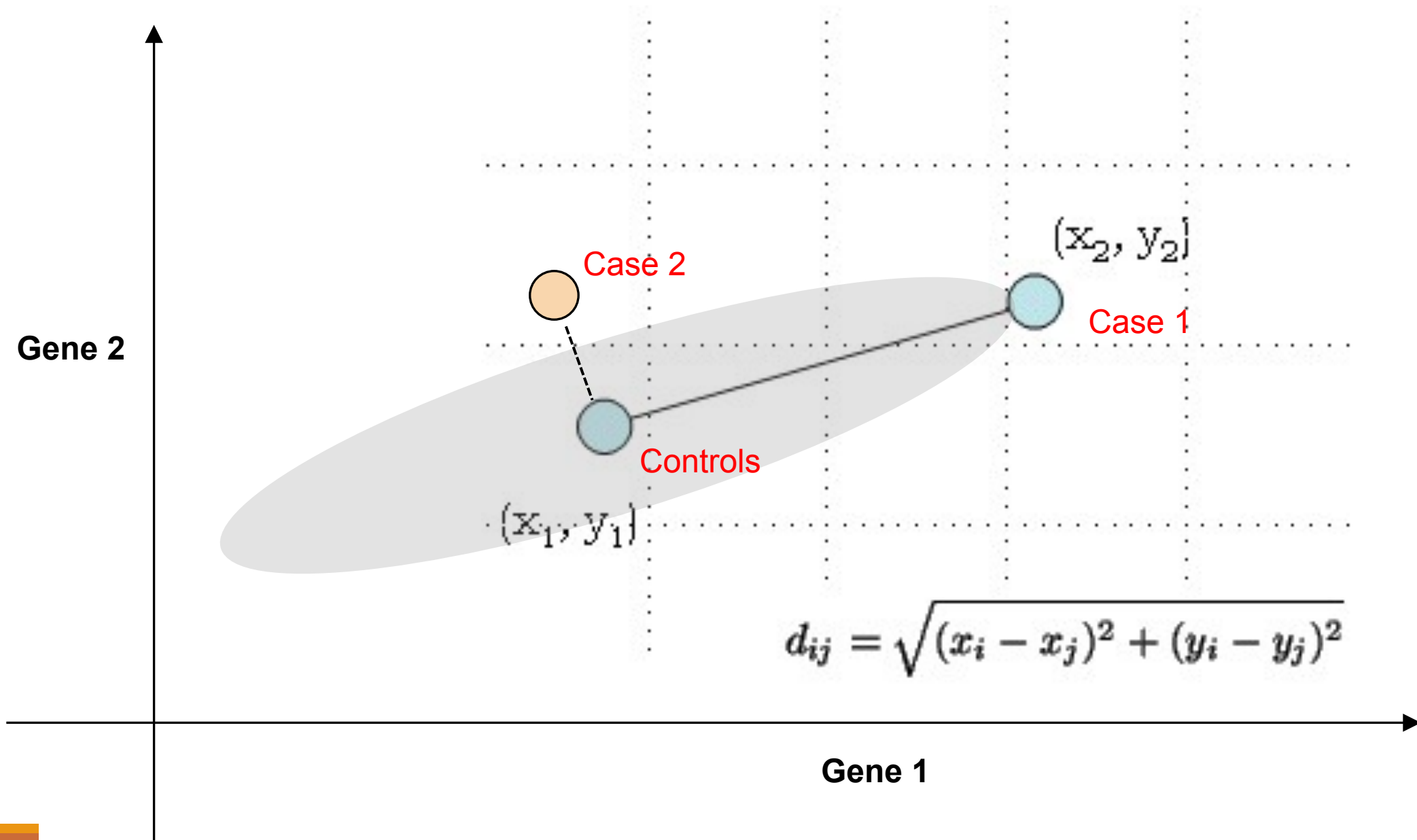
Effect of **rare** genetic variants on gene expression variability









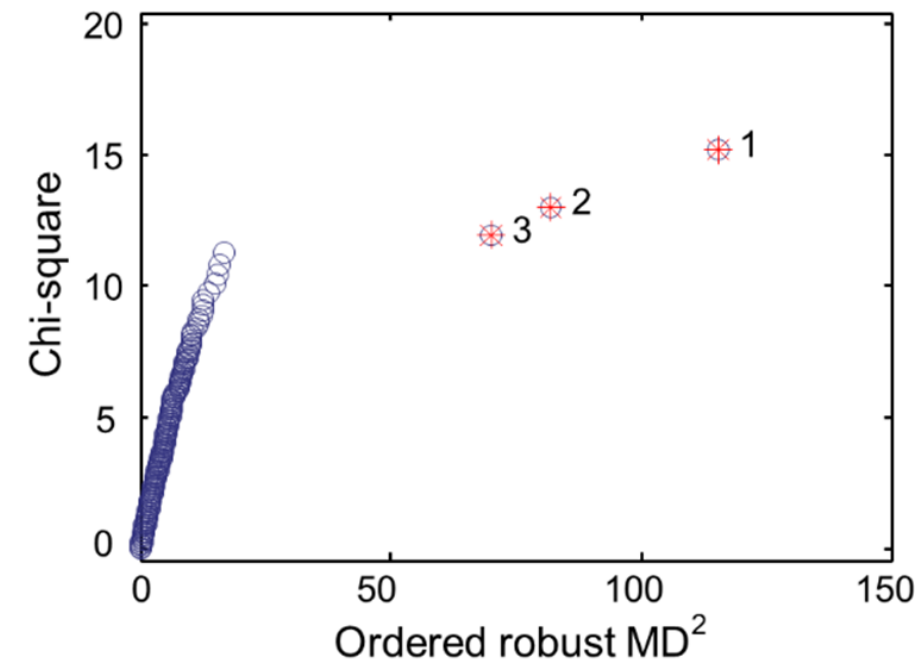
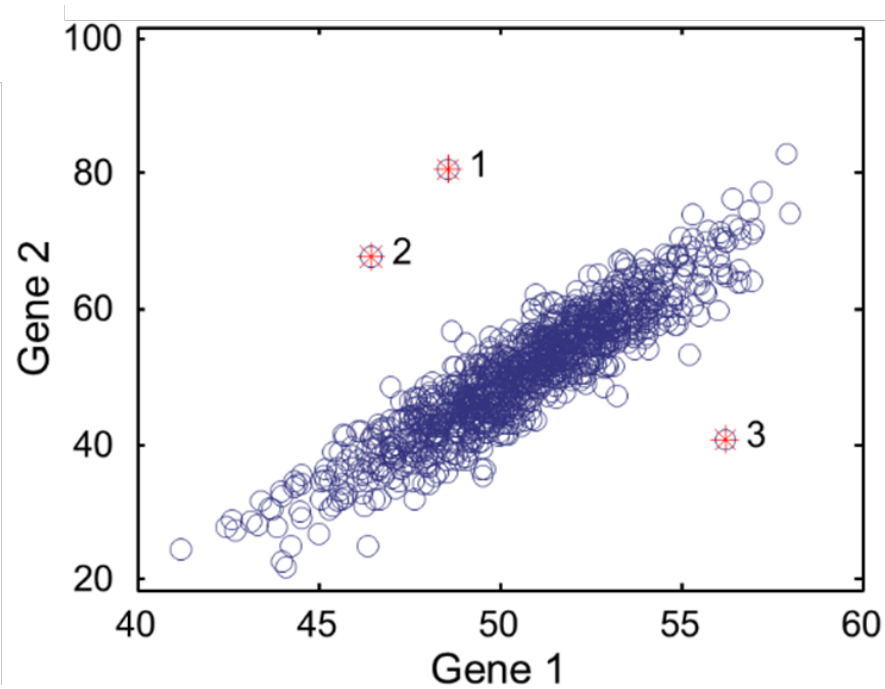


# Mahalanobis distance (MD) is used to detect outliers



P.C. Mahalanobis

1893 – 1972

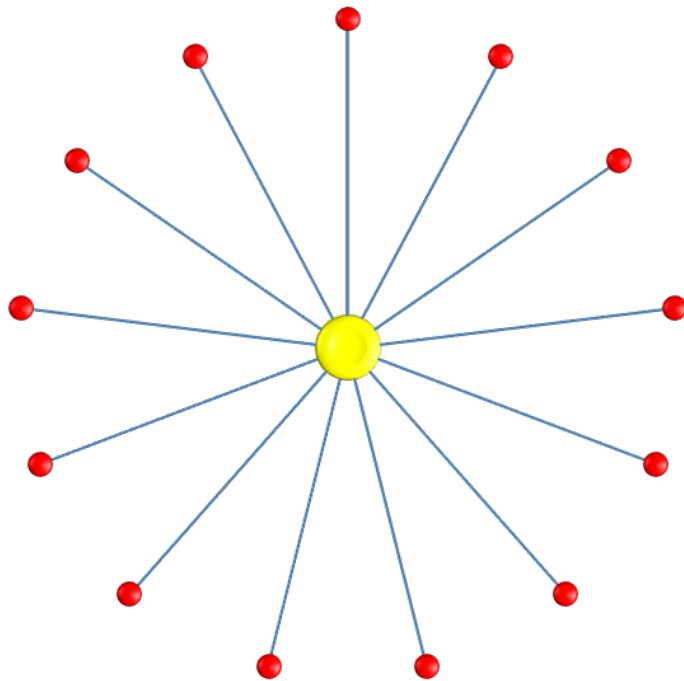


$$MD_i = \sqrt{(E_i - \bar{\mu})^T Cov^{-1} (E_i - \bar{\mu})}$$

# MD measures the level **gene expression dispersion** for a population

---

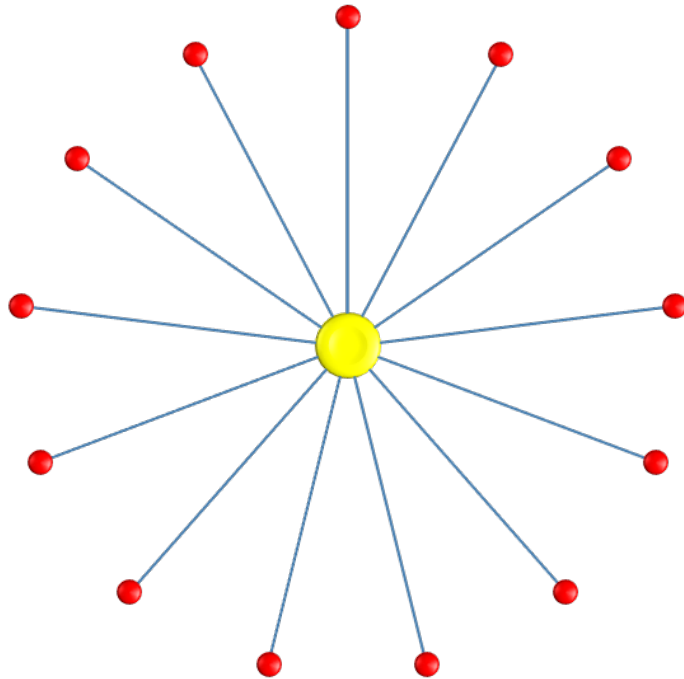
GENE SET 1



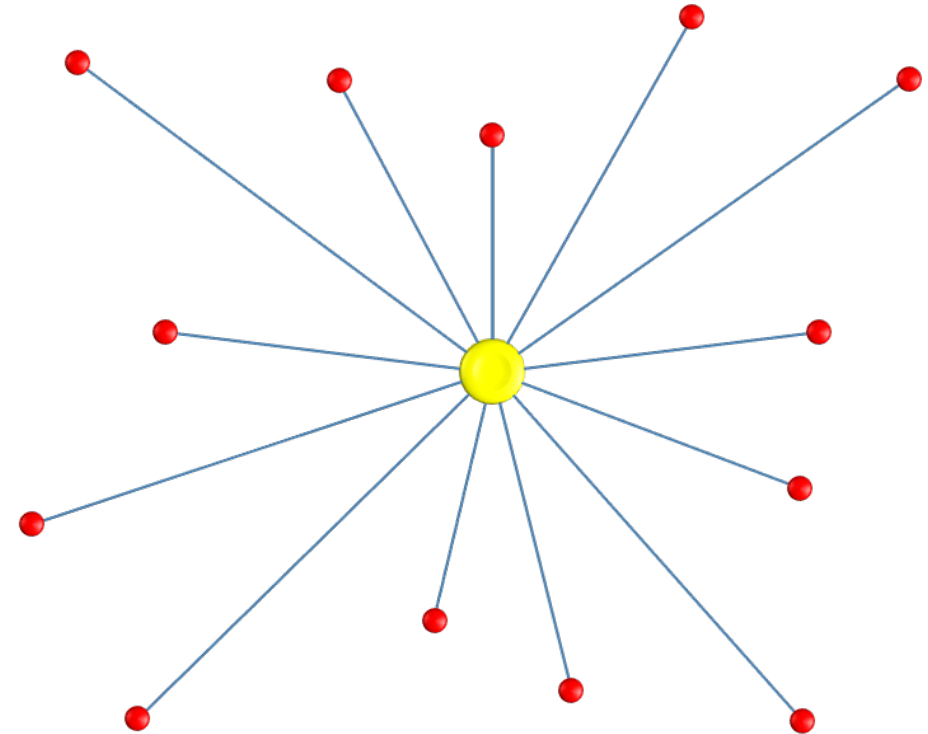
# MD measures the level gene **expression dispersion** for a population

---

GENE SET 1



GENE SET 2





# Sum of squared MD (**SSMD**) – **Overall dispersion** level of a gene set

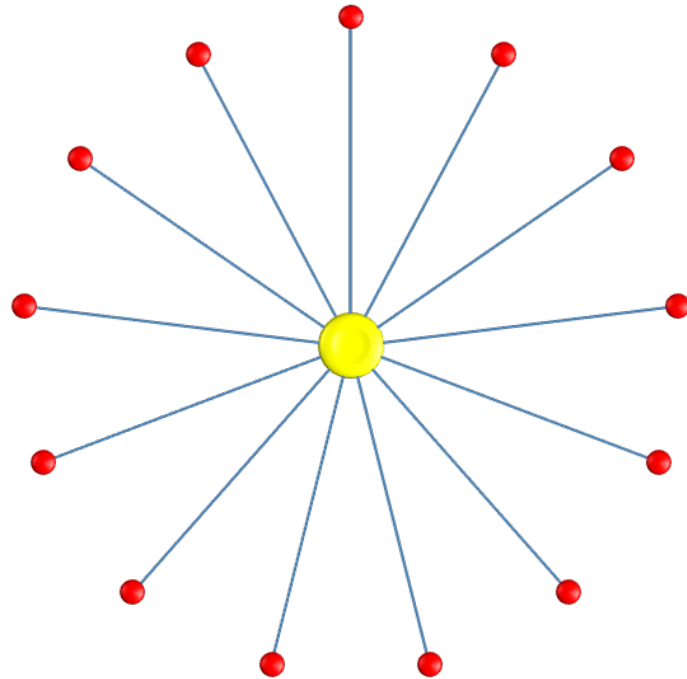
---

$$SSMD = \sum_{i=1}^M MD_i^2$$

# SSMD – overall dispersion level of a gene set

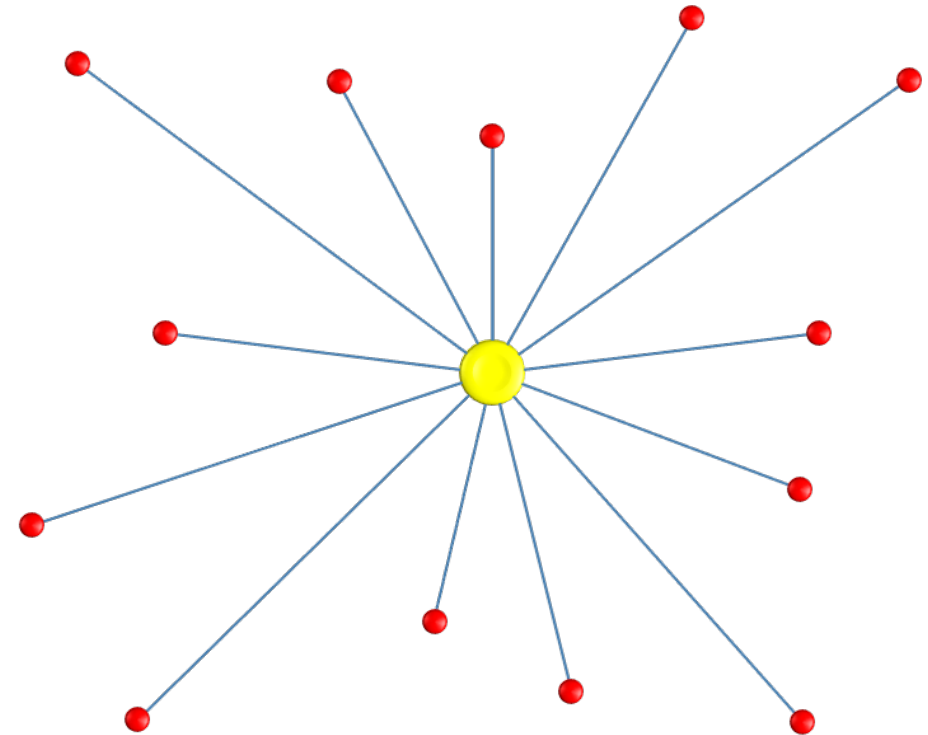
---

GENE SET 1



SSMD ↓↓

GENE SET 2




SSMD ↑↑

# Gene sets (L-SSMD) that tend to be aberrantly expressed

---

MSigDB: molecular signatures database from the Broad Institute  
31 gene sets

- **G-protein coupled receptor activity**
- **Transmission of nerve impulse**
- **Ligand-gated ion channel transportation**
- **Cyclic guanosine monophosphate (cGMP) effects**



**Regulation of cellular processes and modulation of signal transduction**

# Gene sets (S-SSMD) that tend **not** to be aberrantly expressed

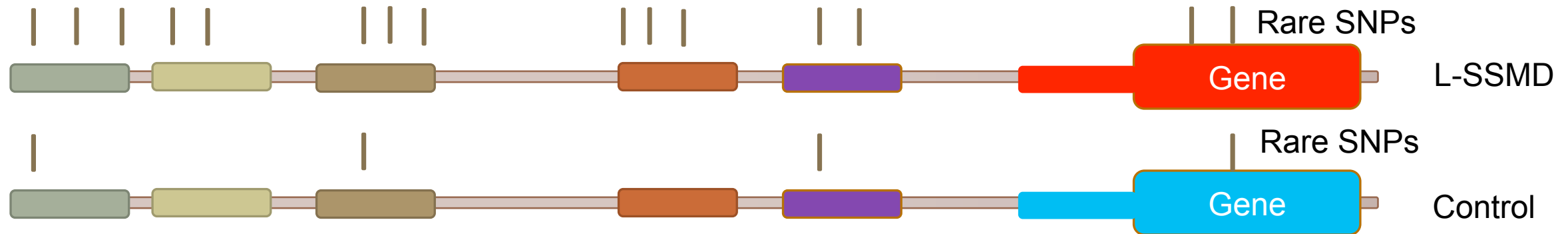
---

MSigDB: molecular signatures database from the Broad Institute  
13 gene sets

- **Homologous recombination repair of replication-independent double-strand breaks**
- **Transfer of a phosphate group to a carbohydrate substrate**
- **Cell cycle control**

**Fundamental  
molecular functions  
and metabolic  
pathways**

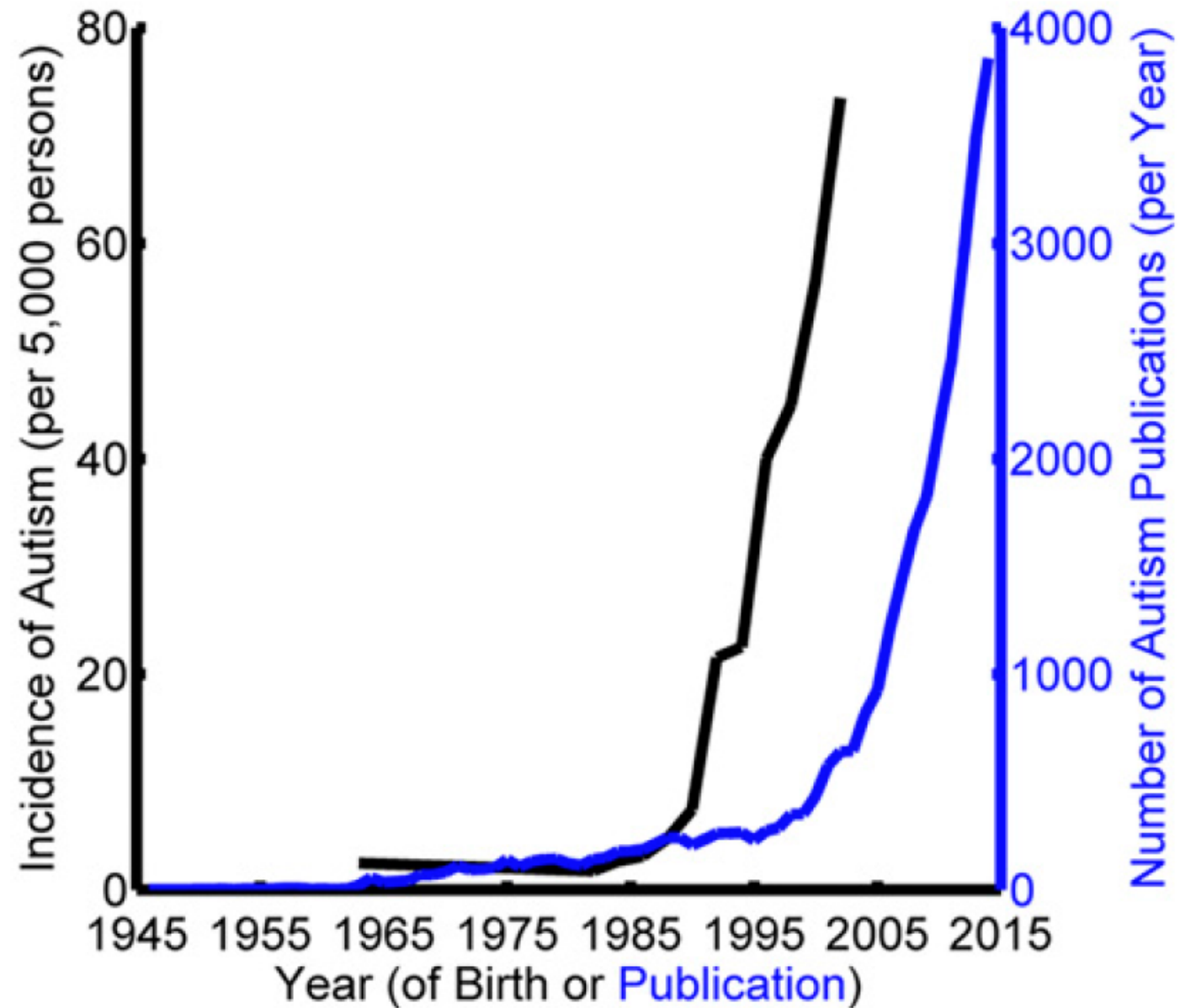
# SNP density in regulatory regions of L-SSMD genes in outlier individuals

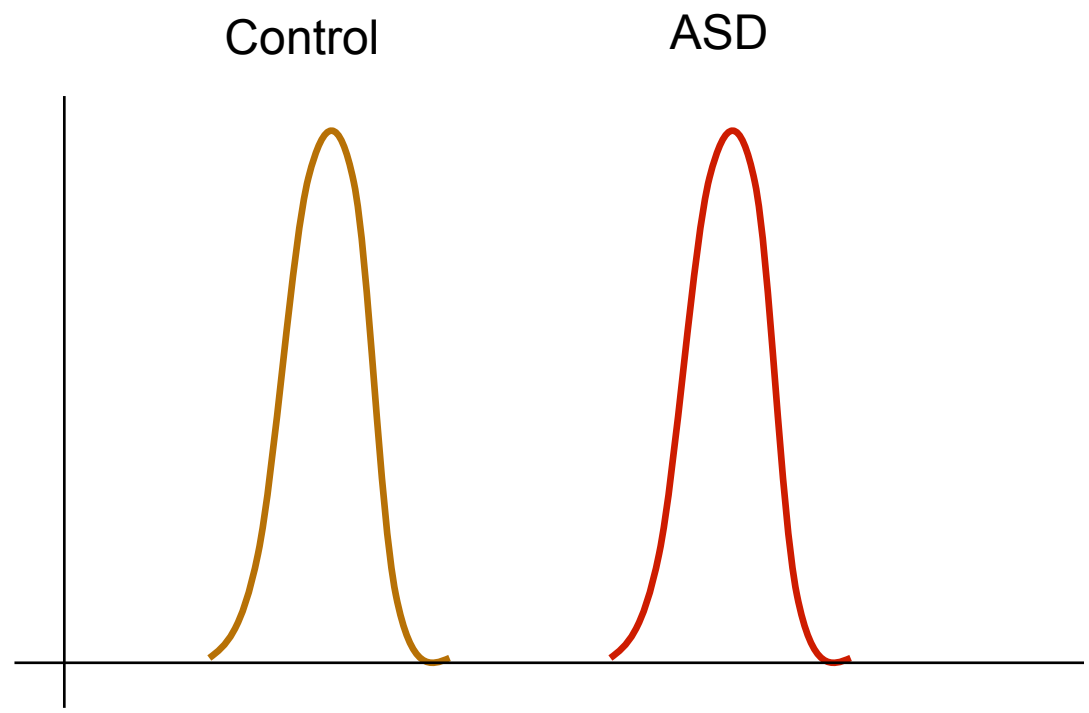


## ENCODE regulatory regions

- E: enhancer
- TSS: transcription start site
- T: transcribed region
- PF: predicted promoter flanking region
- CTCT: CTCF-enriched element
- R: repressed or low-activity region
- WE: weak enhancer or open chromatin cis-regulatory element

# Autism Spectrum Disorder (ASD)

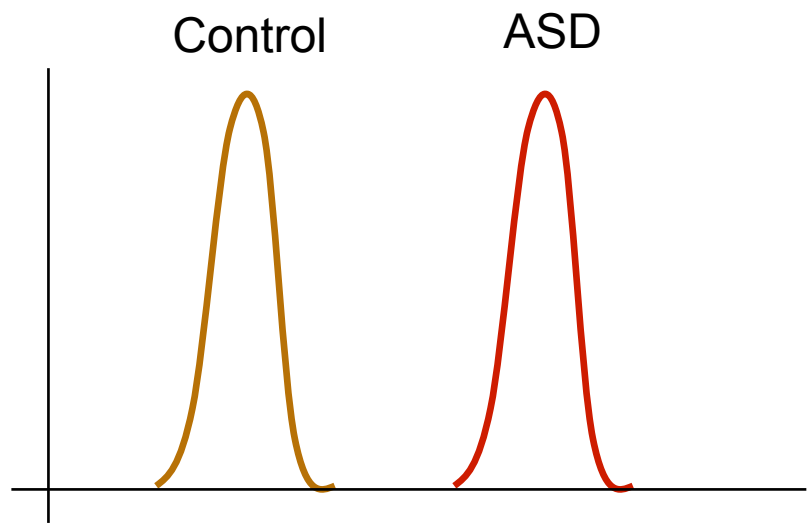




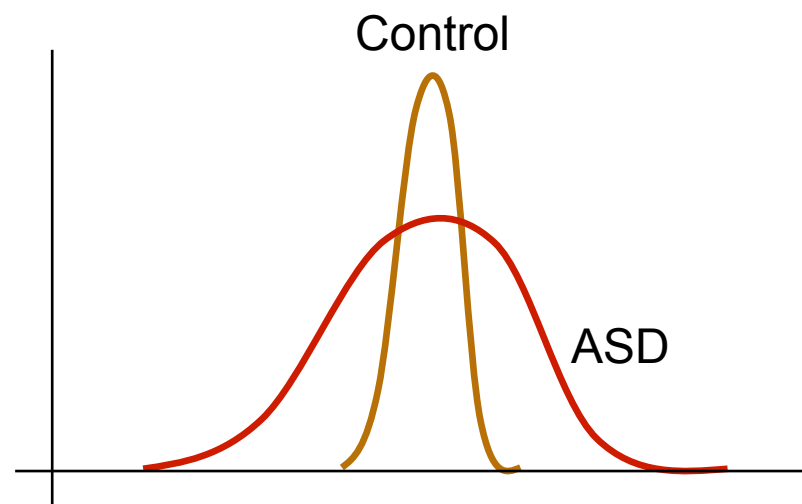
**DE**







**DE**



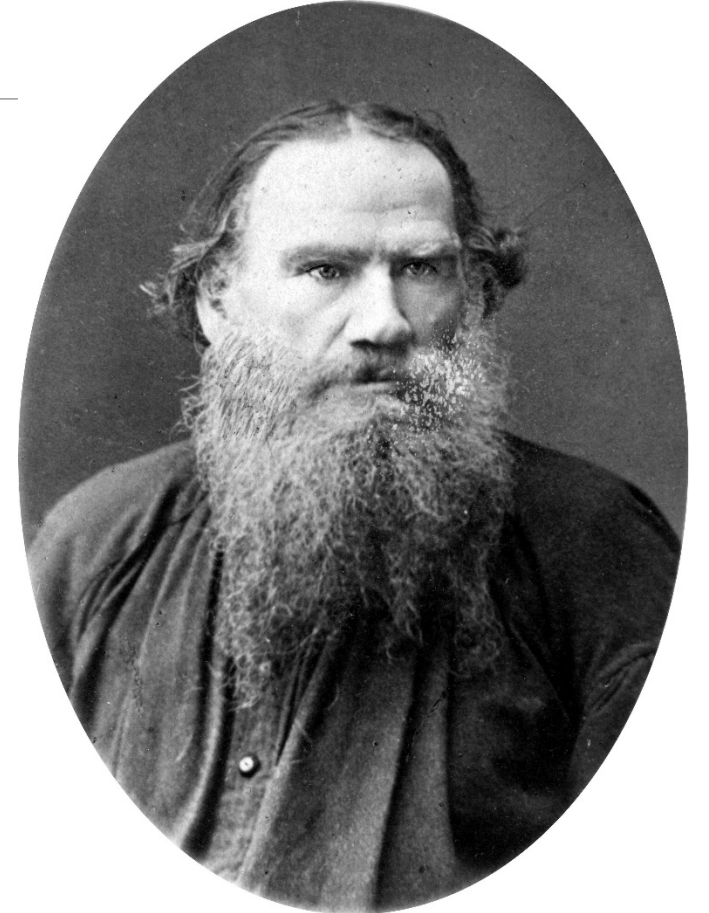
**DV**

# Anna Karenina Principle

---

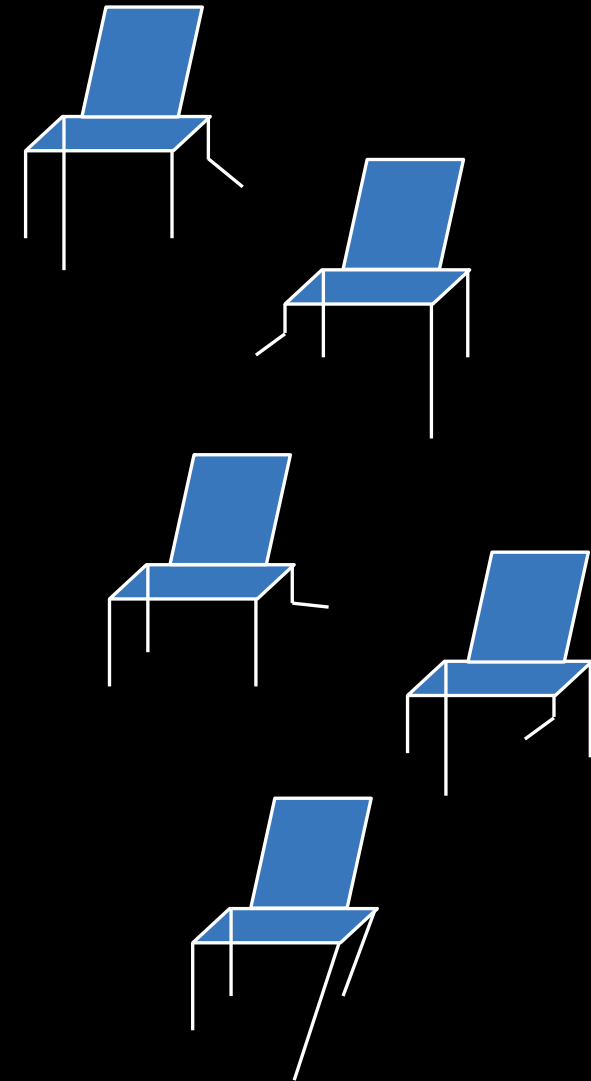
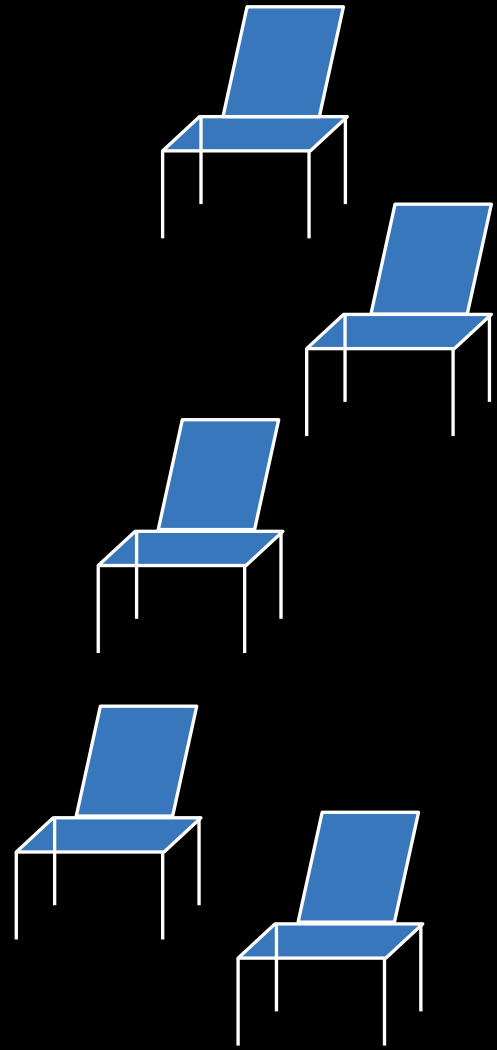
*“Happy families are all alike; every unhappy family is unhappy in its own way.”*

*All **healthy people** are alike; each **sick person** is sick in his or her own way.*



Leo Tolstoy  
1828 – 1910

# Chair Model



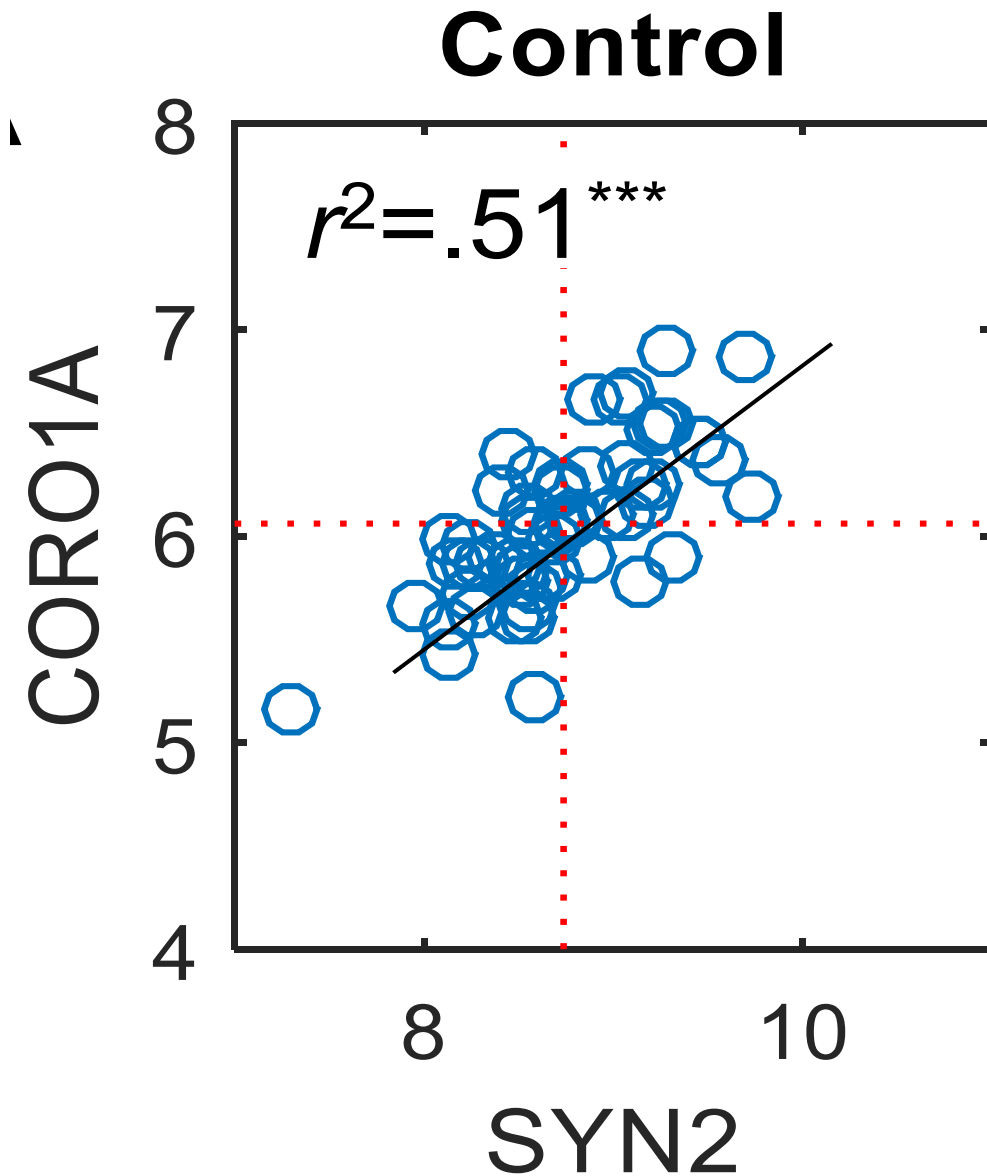
Brain RNA-seq:

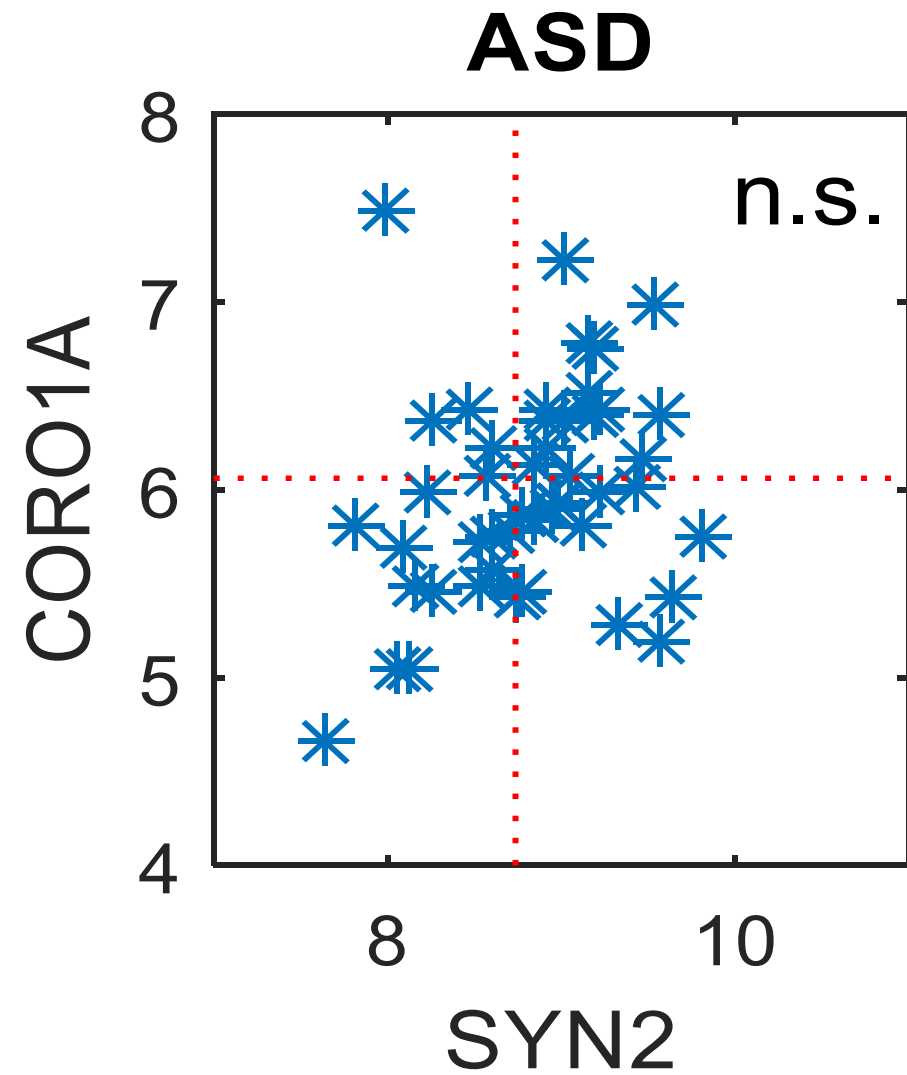
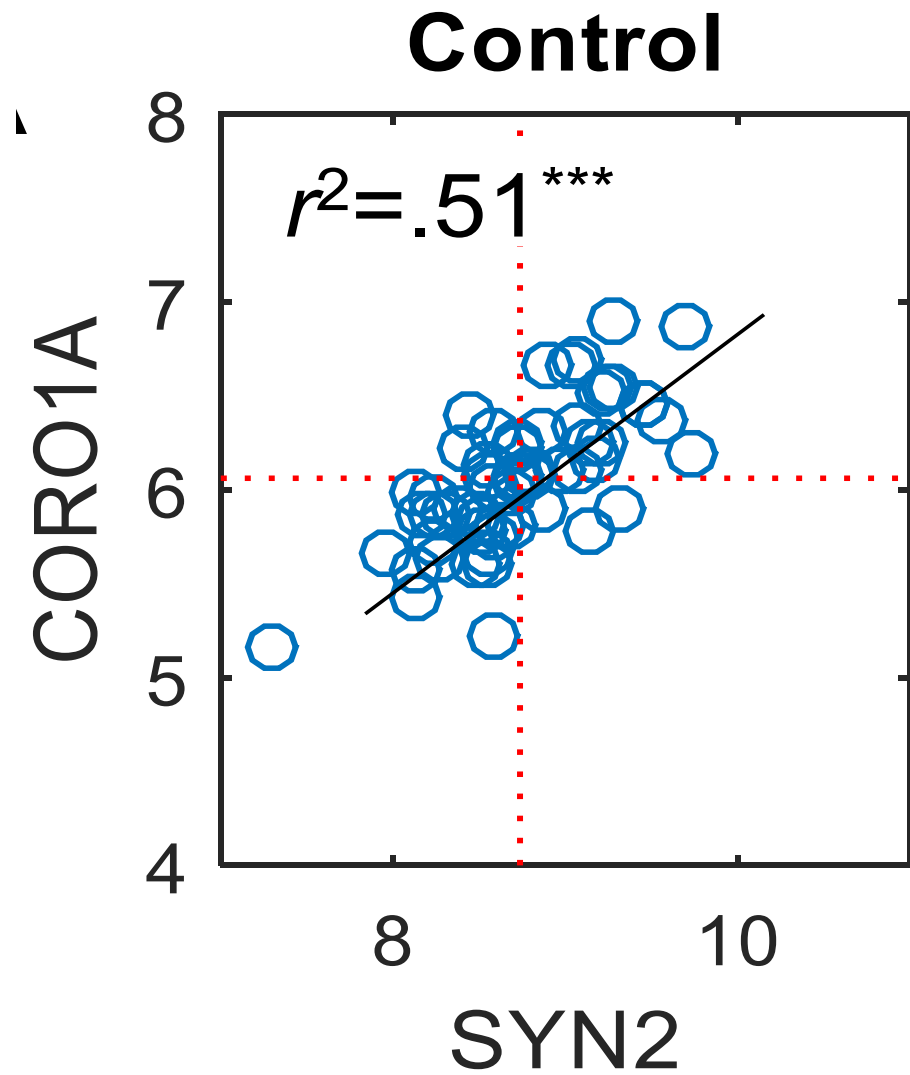
- 47 ASD
- 57 controls

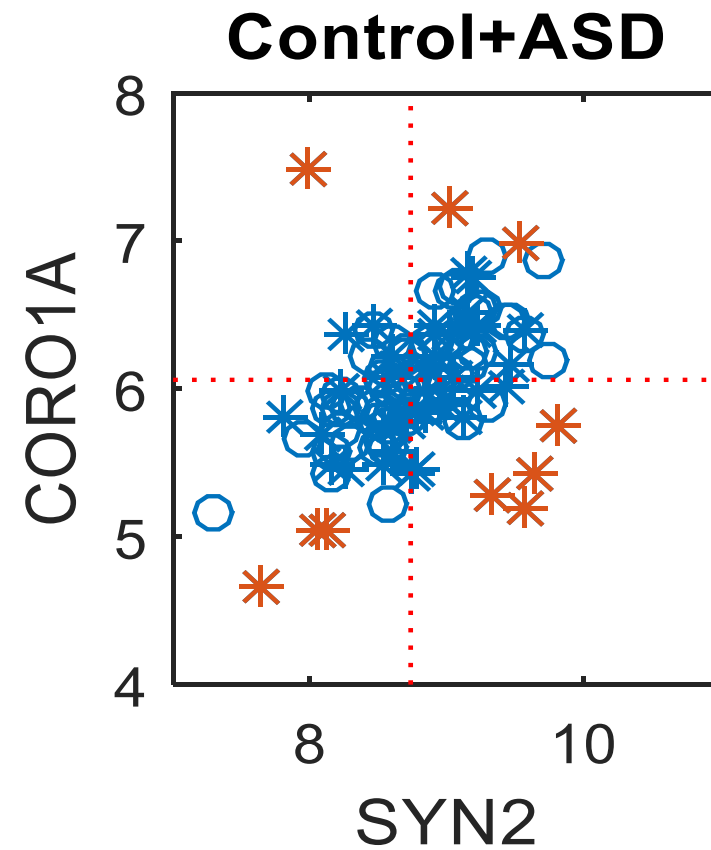
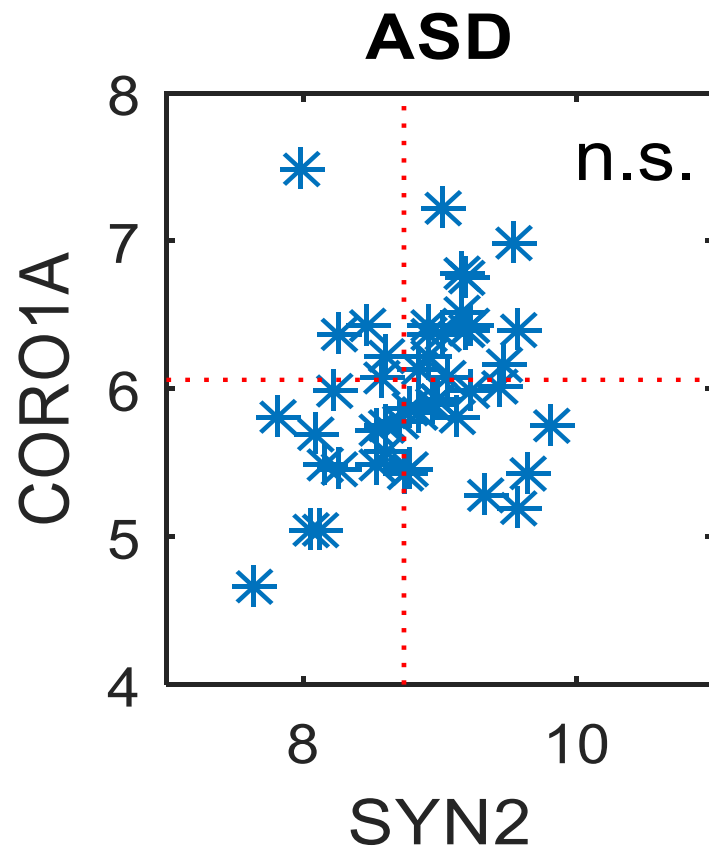
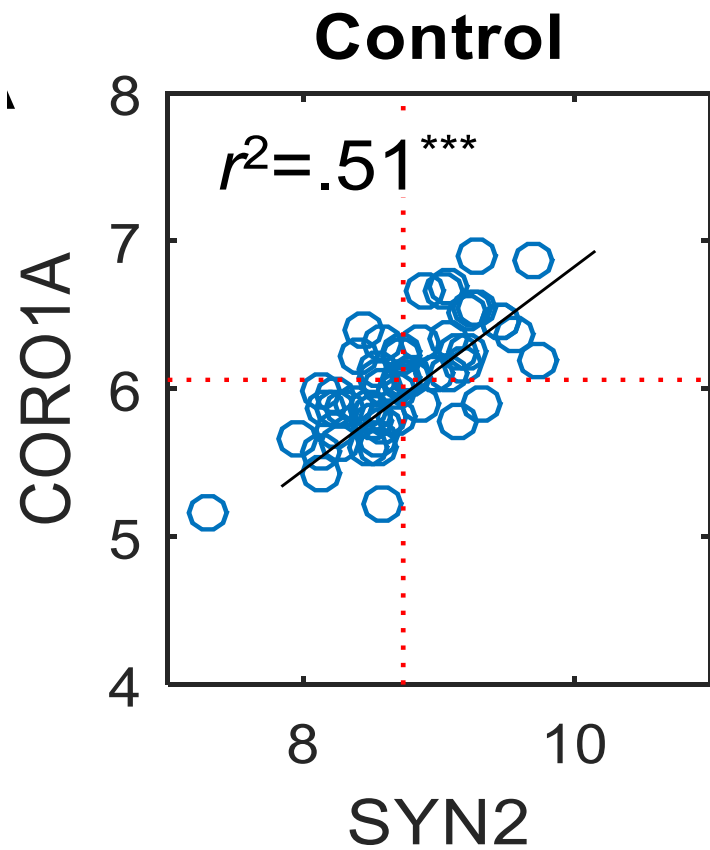
Gupta et al. (2014) *Nat Commun* 5:5748.

Coronin 1A facilitates formation of heterotrimeric or multiprotein complexes.

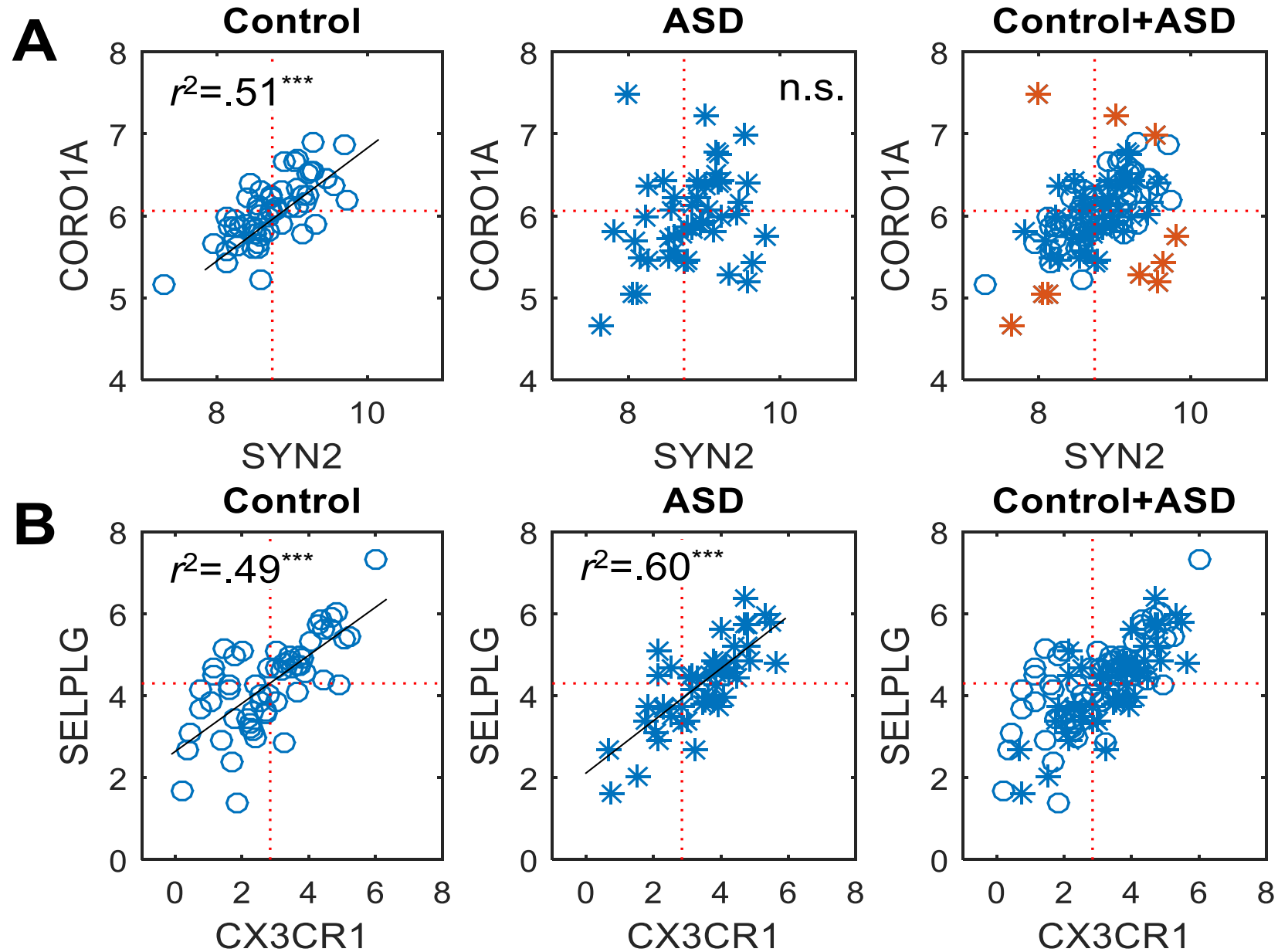
Synapsin II encodes neuronal phosphoprotein associated with the cytoplasmic surface of synaptic vesicles.

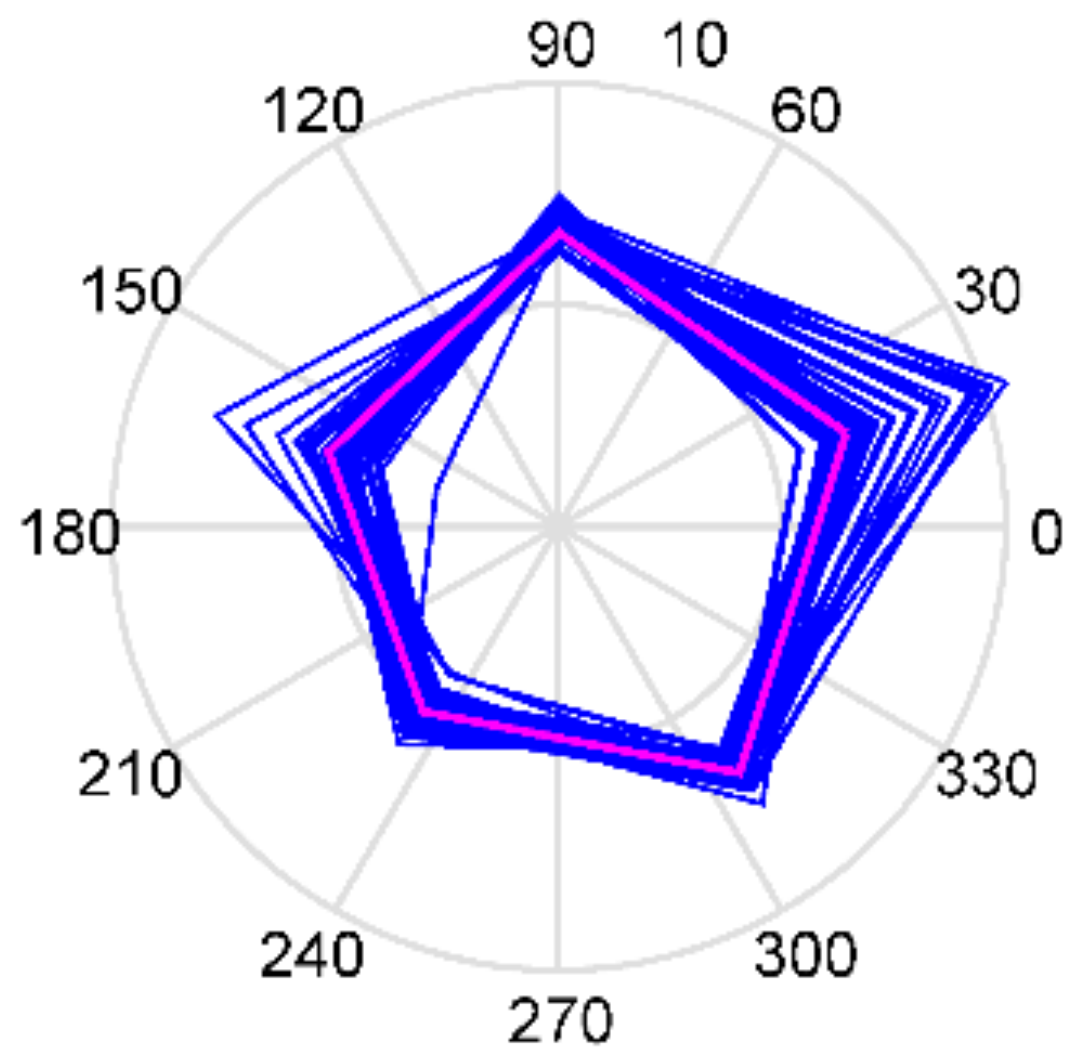
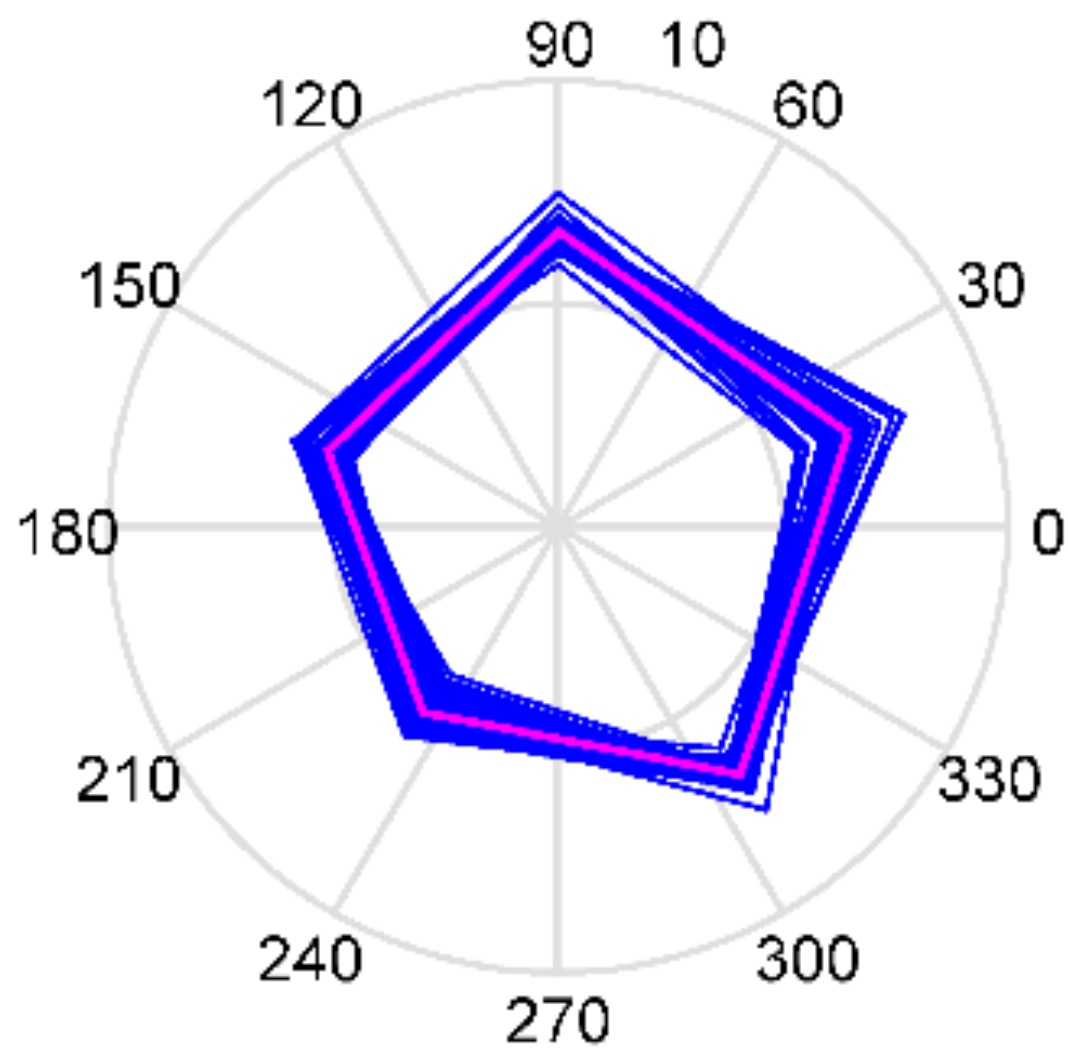






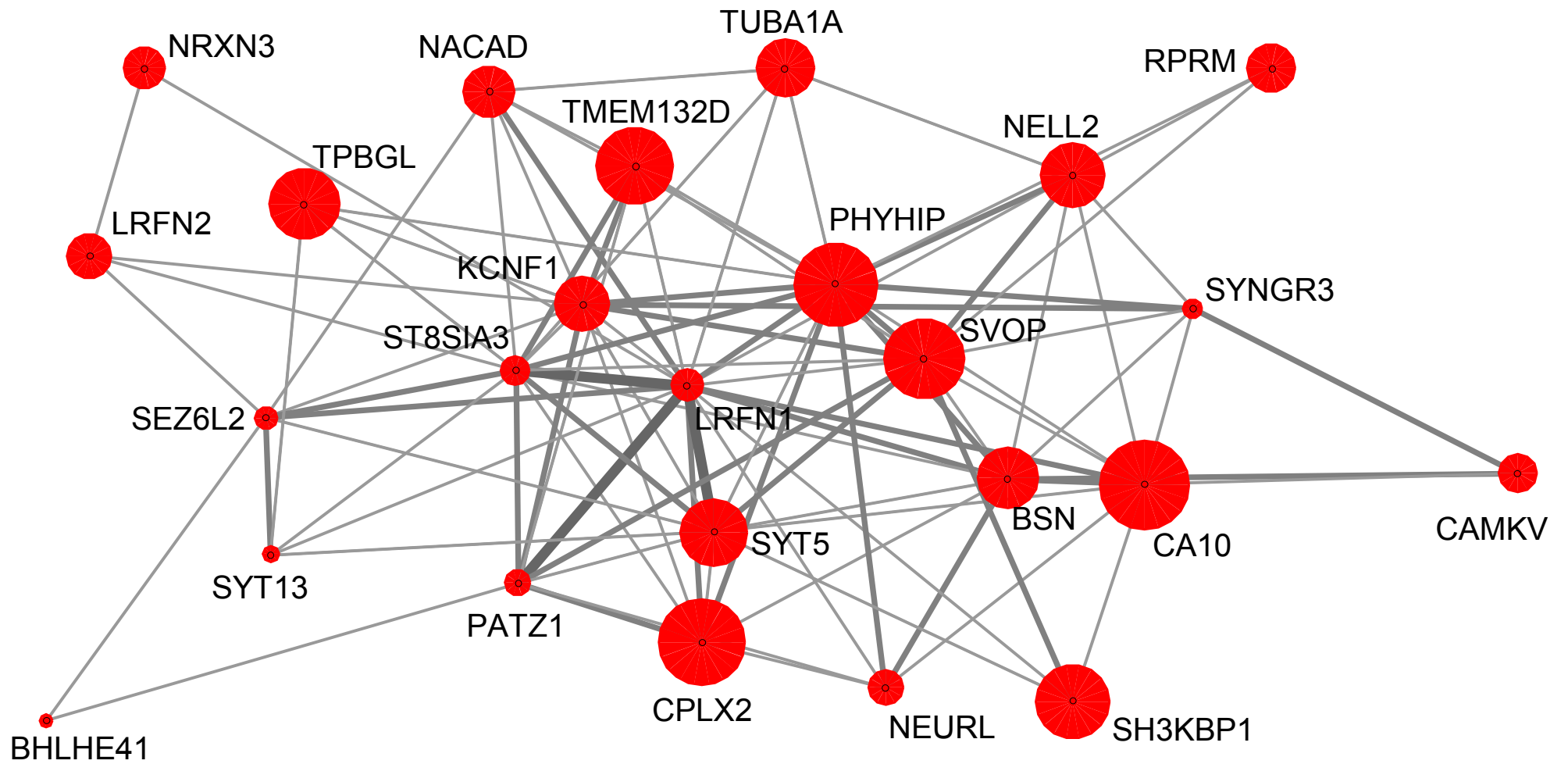


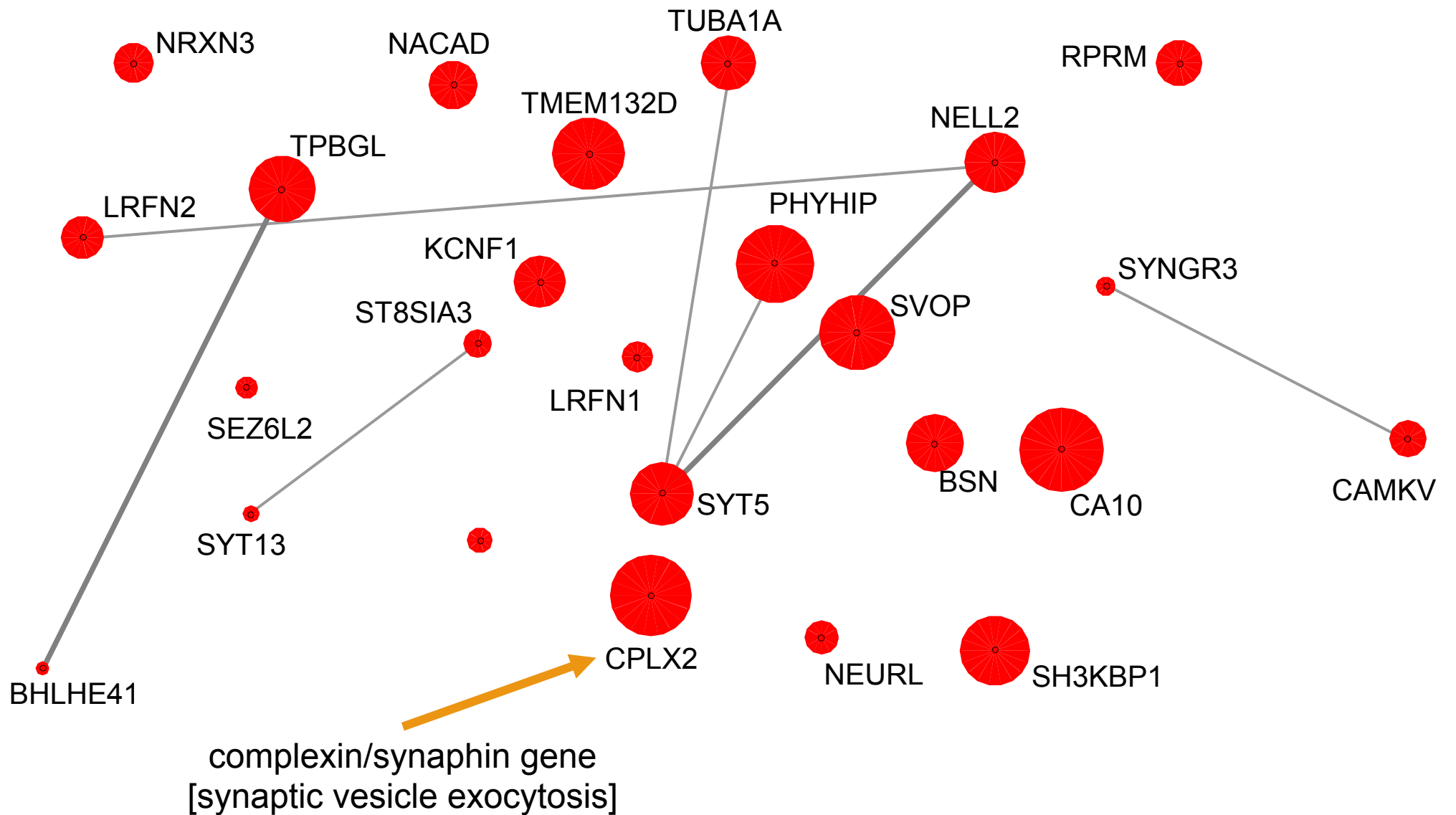


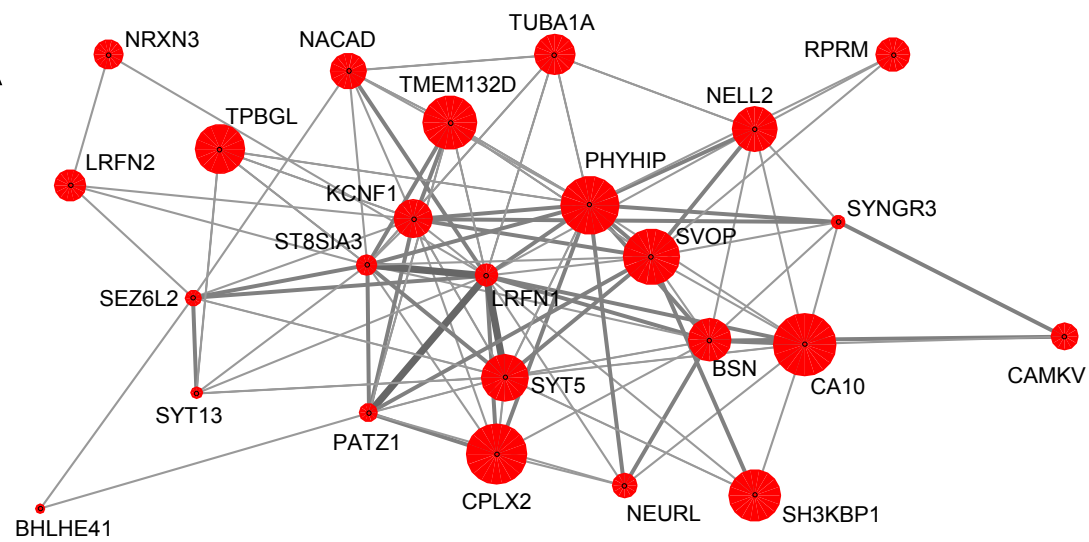
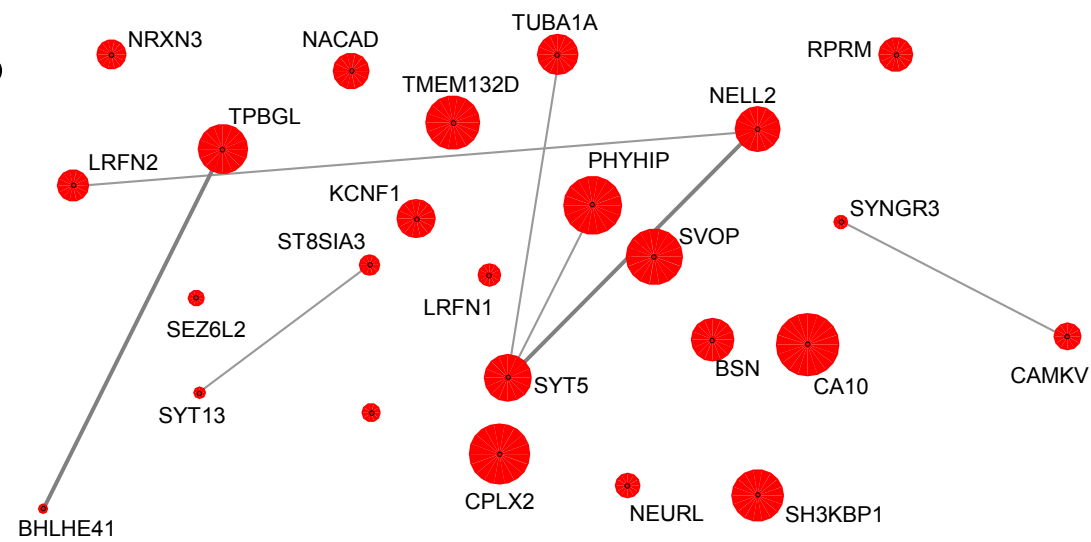
**ASD****Control**

|                                     | GSEA gene set   | # of genes* | Top $\Delta$ SSMD gene |
|-------------------------------------|---|-------------|------------------------|
| <b>Metabolism and biosynthesis</b>  |   |             |                        |
|                                     | KEGG_PENTOSE_PHOSPHATE_PATHWAY                            | 19/27       | H6PD, PRPS2, PFKP      |
|                                     | KEGG_STEROID_BIOSYNTHESIS                                 | 14/17       | SC5DL, NSDHL, DHCR7    |
|                                     | REACTOME_CHOLESTEROL_BIOSYNTHESIS                         | 20/24       | SQLE, HSD17B7, HMGCR   |
|                                     | REACTOME_BRANCHED_CHAIN_AMINO_ACID_CATABOLISM             | 16/17       | DLD, HIBADH, MCCC2     |
| <b>Immune/Inflammatory response</b> |   |             |                        |
|                                     | BIOCARTA_LAIR_PATHWAY                                     | 4/17        | SELPLG, C3, ITGB1      |
|                                     | BIOCARTA_41BB_PATHWAY                                     | 12/17       | MAPK8, ATF2, MAPK14    |
|                                     | REACTOME_IL1_SIGNALING                                    | 25/39       | CHUK, RBX1, BTRC       |
|                                     | REACTOME_REGULATION_OF_IFNA_SIGNALING                     | 6/24        | STAT1, PTPN1, JAK1     |
| <b>Signaling pathway</b>            |   |             |                        |
|                                     | BIOCARTA_IGF1_PATHWAY                                     | 20/21       | JUN, CSNK2A1, ELK1     |
|                                     | PID_S1P_S1P2_PATHWAY                                      | 21/24       | MAPK8, MAPK14, JUN     |
|                                     | PID_HNF3APATHWAY (FOXA1/HNF3A TF network)                 | 22/44       | NDUFV3, PISD, FOS      |
|                                     | REACTOME_ENERGY_DEPENDENT_REGULATION_OF_MTOR_BY_LKB1_AMPK | 15/18       | PRKAA1, CAB39, TSC1    |
| <b>Vitamins and supplements</b>     |   |             |                        |
|                                     | BIOCARTA_VITCB_PATHWAY                                    | 6/11        | SLC2A3, COL4A2, SLC2A1 |
|                                     | REACTOME_TETRAHYDROBIOPTERIN_BH4_SYNTHESIS_               | 9/13        | GCHFR, PTS, AKT1       |

|                                 |   |       |                        |
|---------------------------------|---|-------|------------------------|
|                                 | OF_MTOR_BY_LKB1_AMPK  |       |                        |
| <b>Vitamins and supplements</b> |   |       |                        |
|                                 | BIOCARTA_VITCB_PATHWAY  | 6/11  | SLC2A3, COL4A2, SLC2A1 |
|                                 | REACTOME_TETRAHYDROBIOPTERIN_BH4_SYNTHESIS_<br>RECYCLING_SALVAGE_AND_REGULATION | 9/13  | GCHFR, PTS, AKT1       |
| <b>Miscellaneous</b>            |   |       |                        |
|                                 | REACTOME_ACTIVATED_POINT_MUTANTS_OF_FGFR2                                       | 4/16  | FGF9, FGFR2, FGF1      |
|                                 | REACTOME_ACTIVATION_OF_THE_AP1_FAMILY_OF_<br>TRANSCRIPTION_FACTORS              | 10/10 | MAPK14, MAPK3, ATF2    |
|                                 | REACTOME_INWARDLY_RECTIFYING_K_CHANNELS   | 20/31 | KCNJ10, KCNJ4, GNG4    |
|                                 | REACTOME_G2_M_CHECKPOINTS   | 22/45 | MCM2, RFC5, RPA2       |



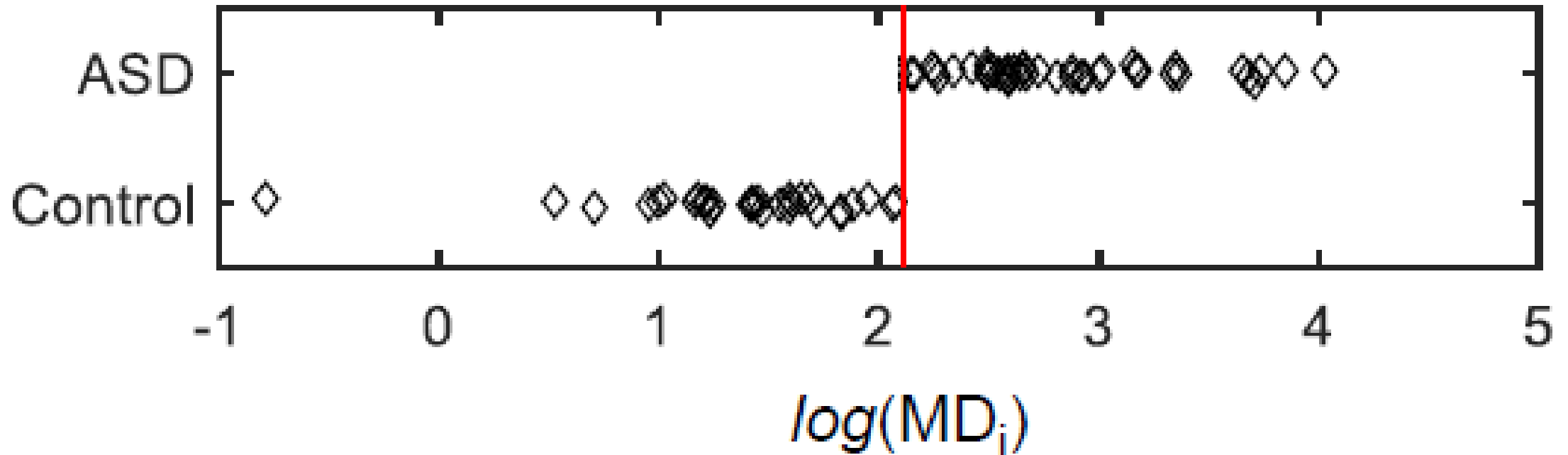


**A****B**



# Search for **gene expression markers** for **early diagnosis**

---



# Search for **gene expression markers** for **early** diagnosis

---

(■N@1)=1.000494500

(■N@3)=2.2177e+11

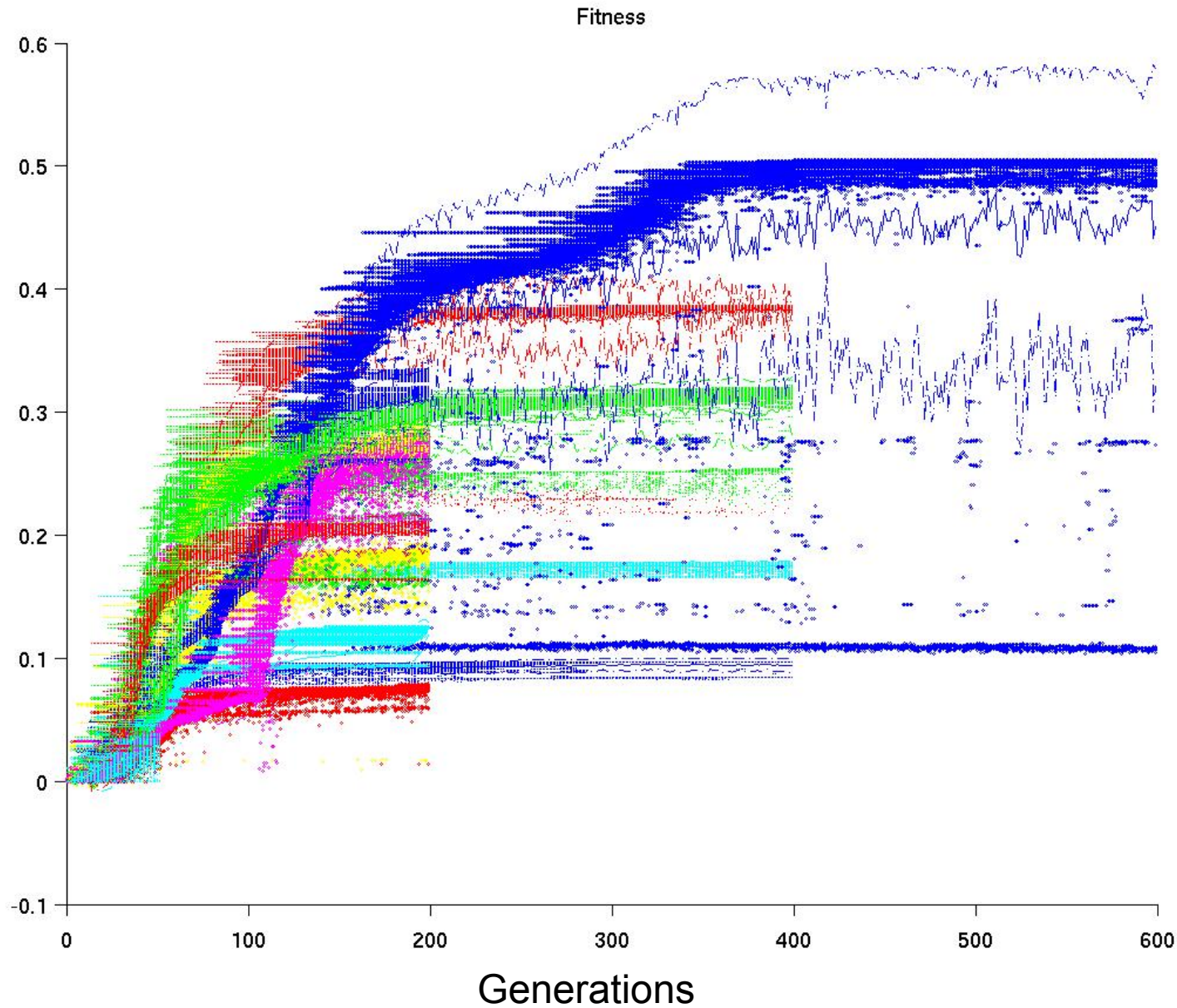
(■N@4)=6.0971e+14

(■N@5)=1.3409e+18

# Genetic Algorithm Optimization

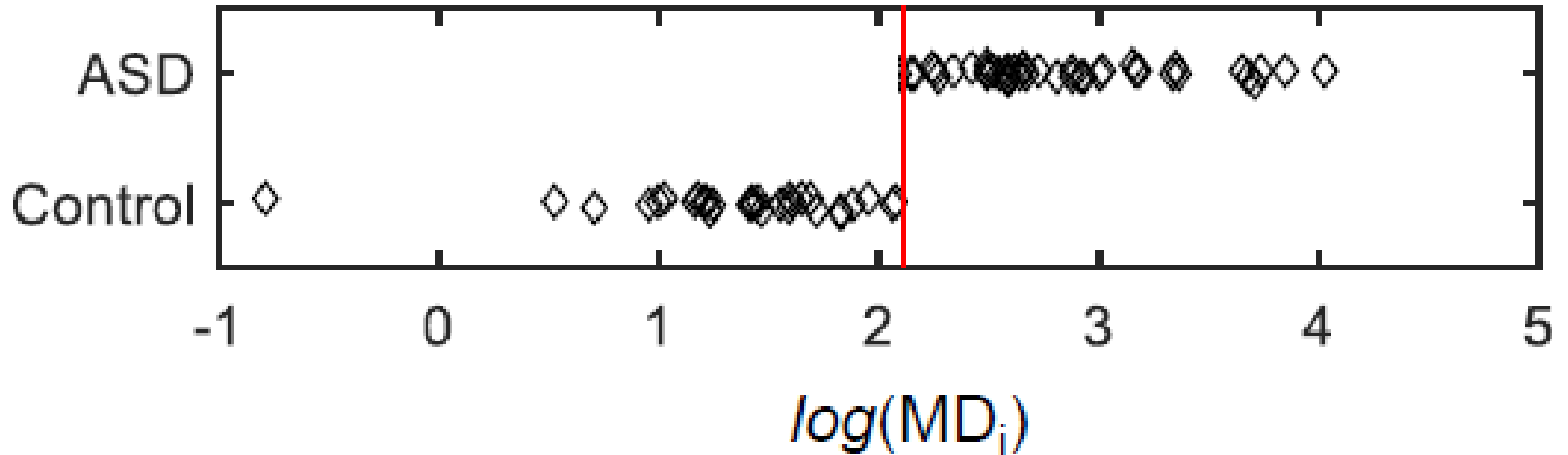
Fitness  
(deltaSSMD)

| Fitness | Initial Population |   |                            |
|---------|--------------------|---|----------------------------|
| 22      | 101010100111110101 | ← | Selected parent string one |
| 9       | 110011010101011100 | ← | 110011010101011100         |
| 8       | 111110101111010101 |   |                            |
| 70      | 111001111100001001 |   |                            |
| 19      | 110011010101011100 |   |                            |
| 48      | 101110101111001001 | ← | Selected parent string two |
| 23      | 110011010101011100 | ← | 111001111100001001         |
| 38      | 111001111100001001 | ← |                            |



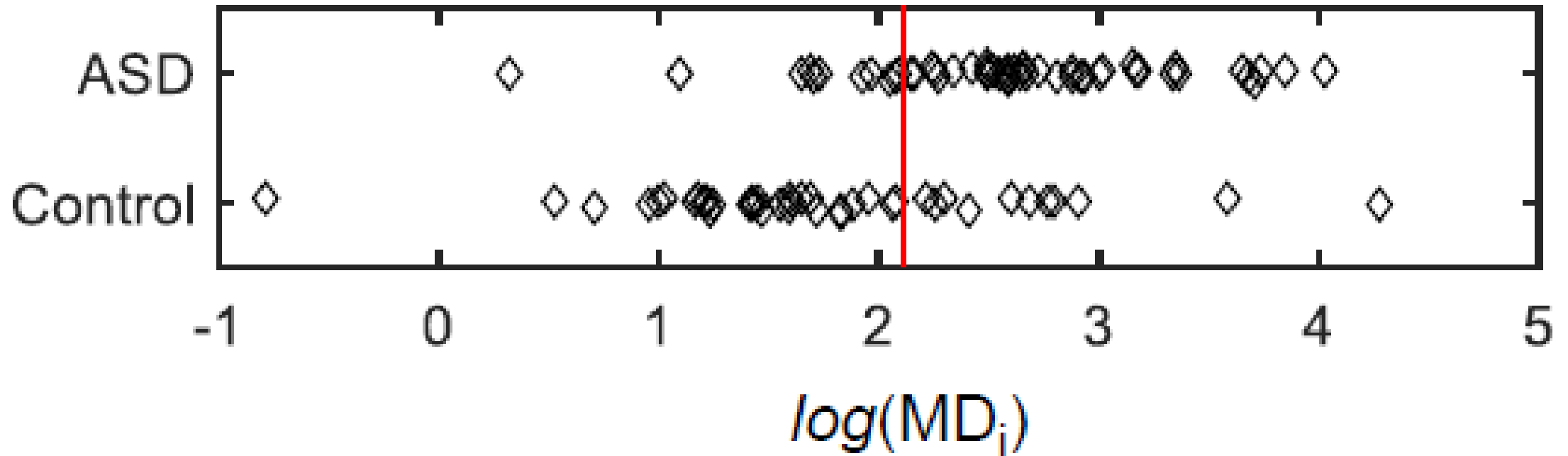
# Search for **gene expression markers** for **early diagnosis**

---



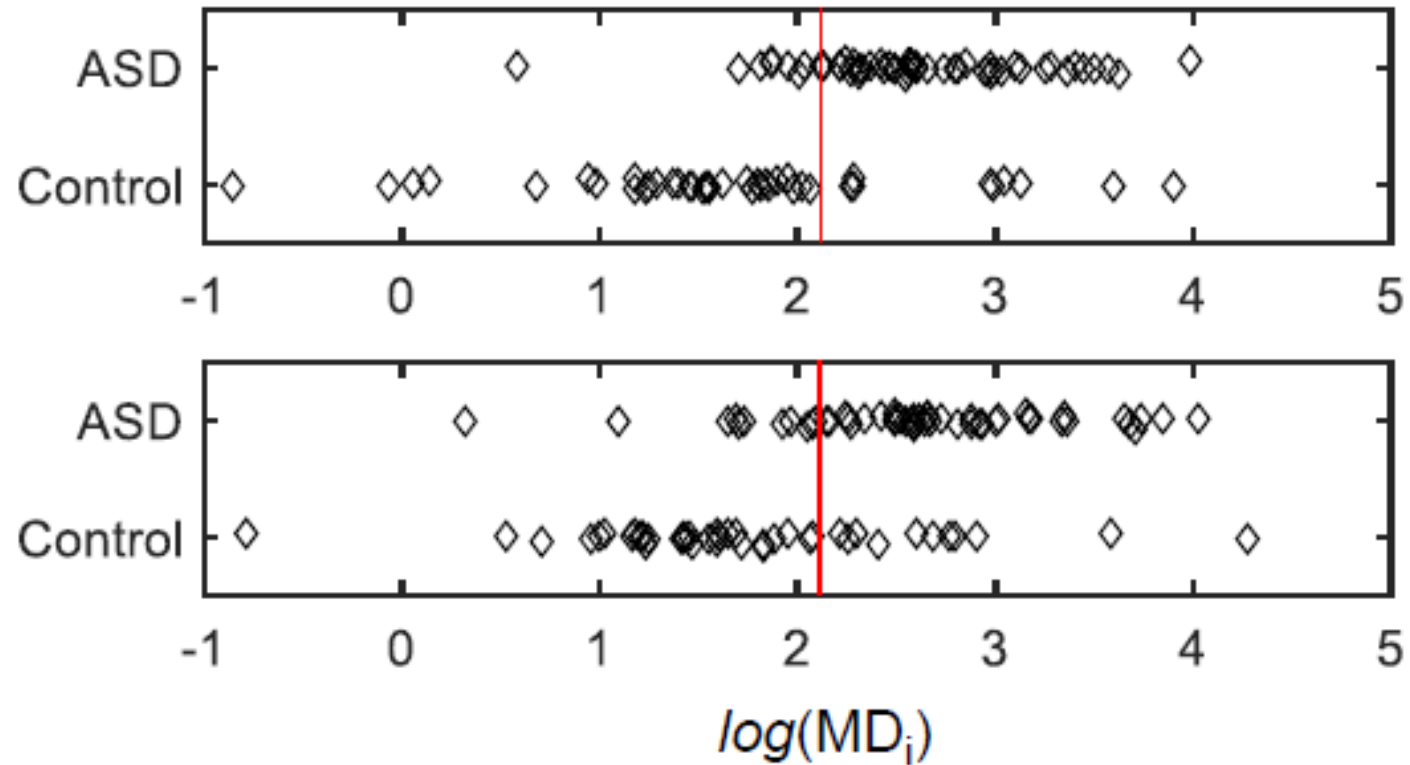
***{EVI2B, MYLIP, OR11G2, TSPAN16, ZNF594}***

---

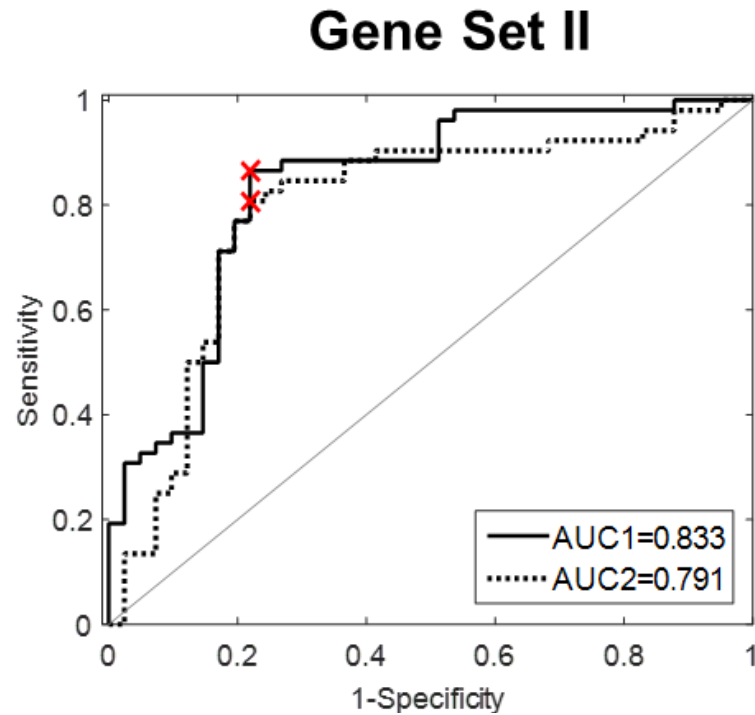


***{EVI2B, MYLIP, OR11G2, TSPAN16, ZNF594}***

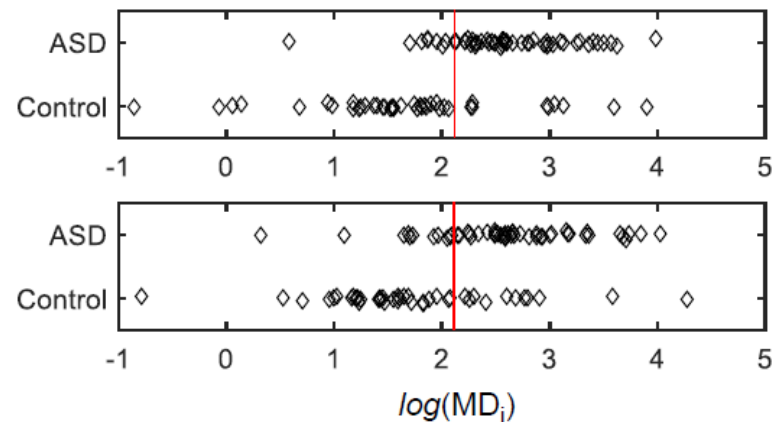
---



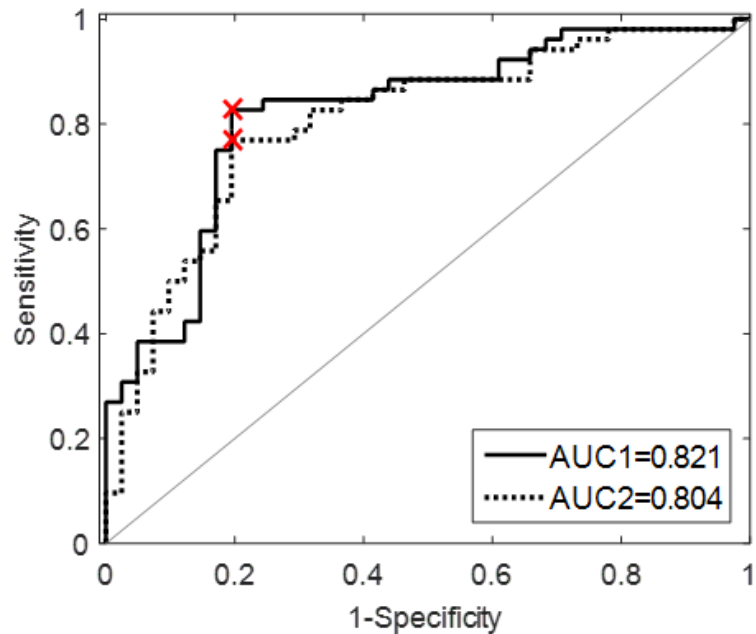
# Receiver Operating Characteristic (ROC) Curve



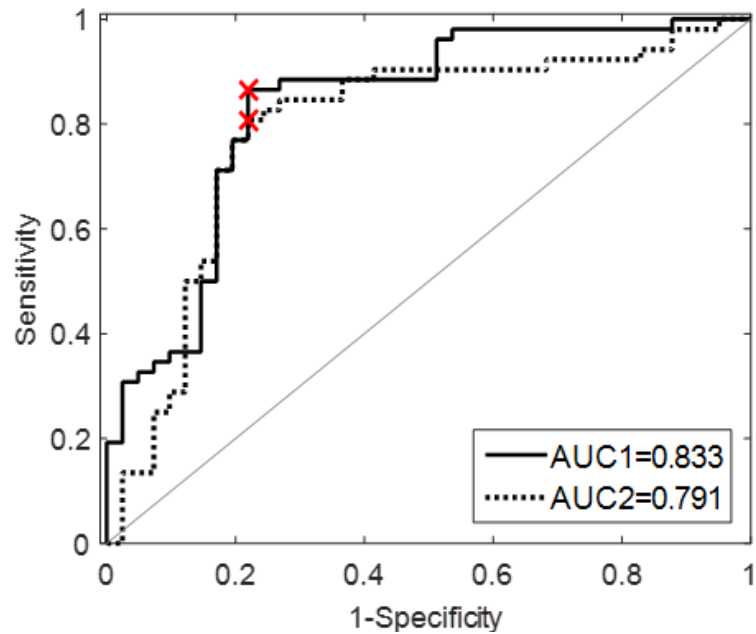
{**EVI2B**, **MYLIP**, OR11G2, TSPAN16, **ZNF594**}



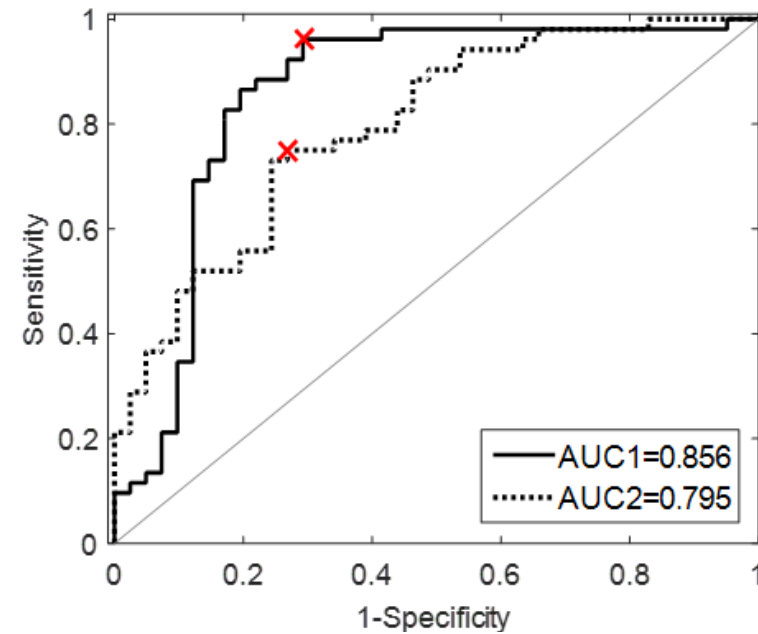


**Gene Set I**

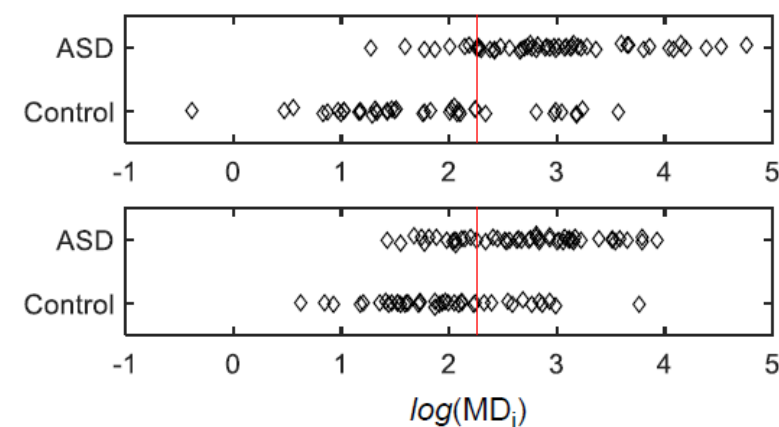
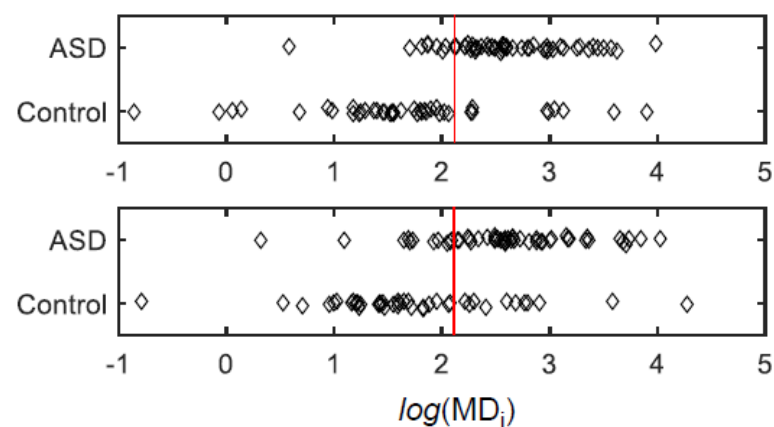
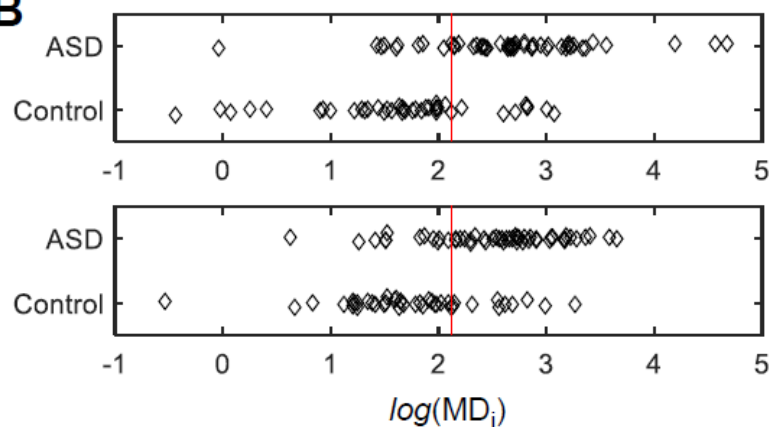
{FAM120A, HDC, **OR13C8**, PSAP, **RFX8**}

**Gene Set II**

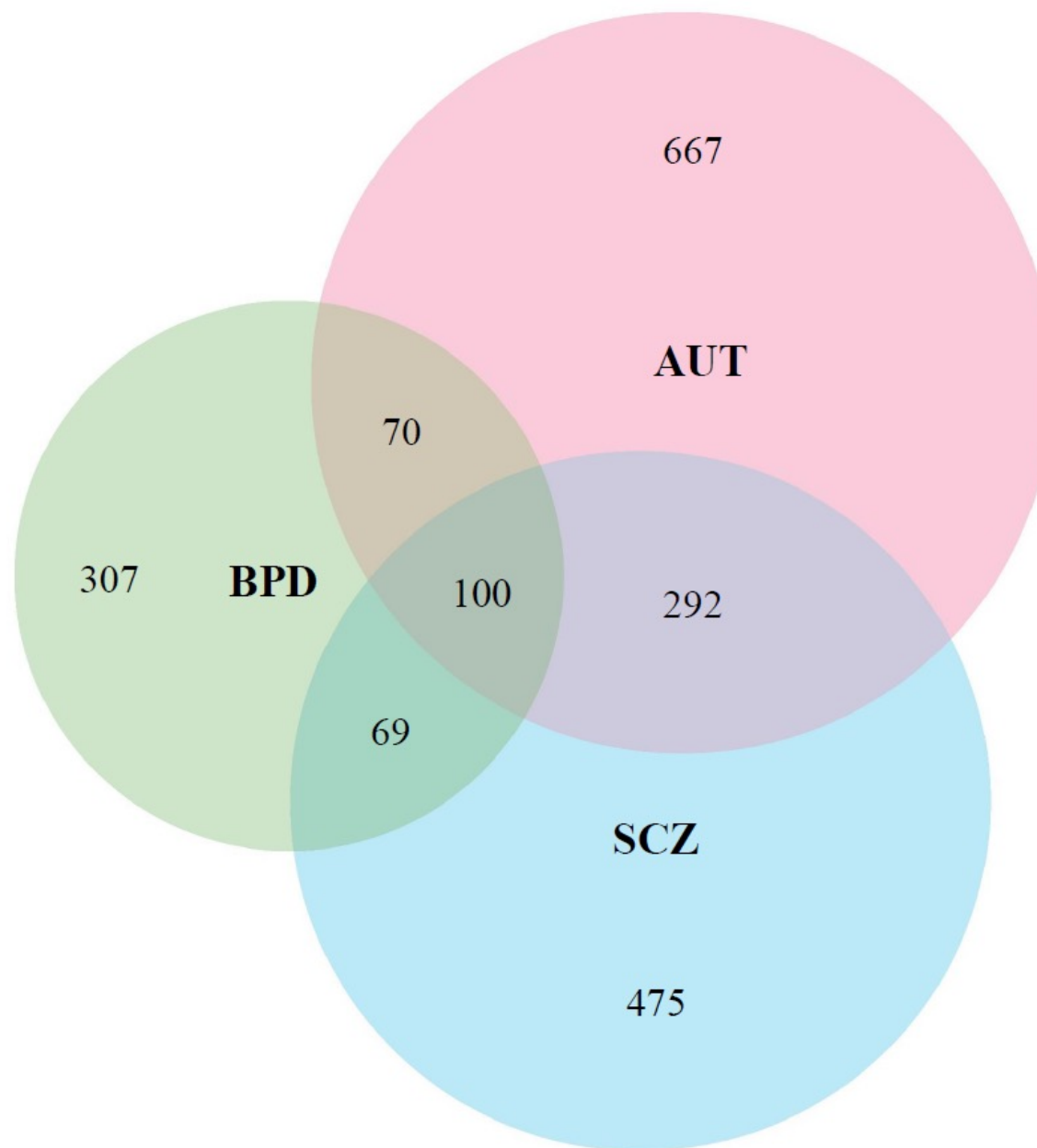
{**EVI2B**, **MYLIP**, OR11G2, TSPAN16, **ZNF594**}

**Gene Set III**

{**BCL11A**, **DST**, ORM2, RBM14, **SERAC1**}

**B**

Common  
dysregulated  
gene sets  
between **AUT**,  
**SCZ**, and **BPD**



# Summary

---

- Detecting **aberrant gene expression** and identifying underlying genes and mutations represent a **new discovery and diagnostic strategy** for *genetically heterogeneous* disorders such as autism.

