

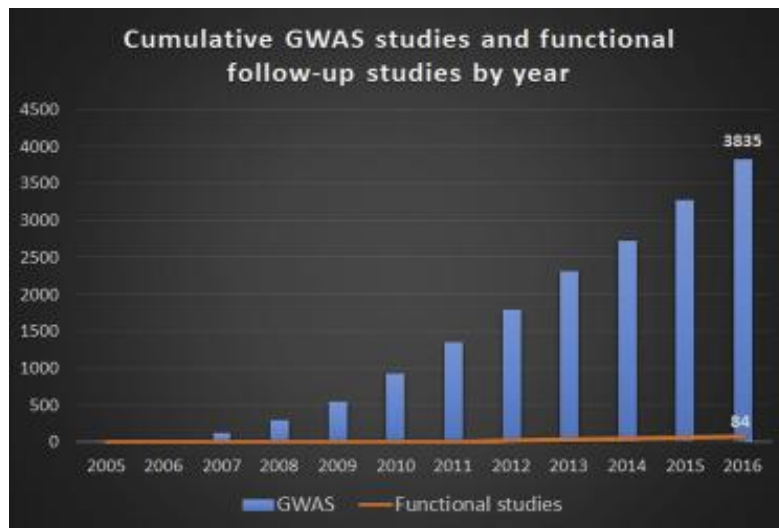
Linking variants to genes and variant effect prediction

Gerard Bouland

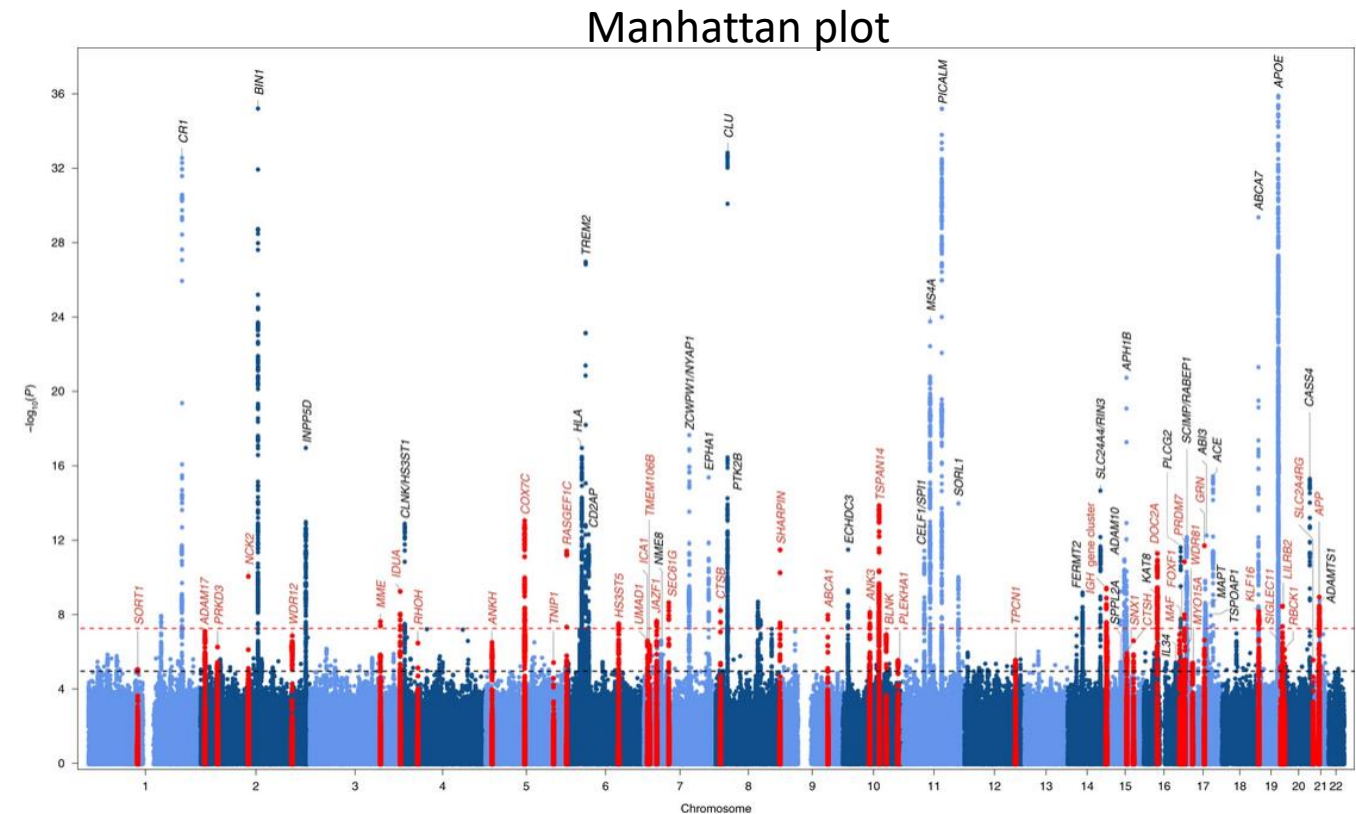
The Delft Bioinformatics Lab, TU Delft

Human Genetics, LUMC

A lot of risk variants but not a lot of understanding

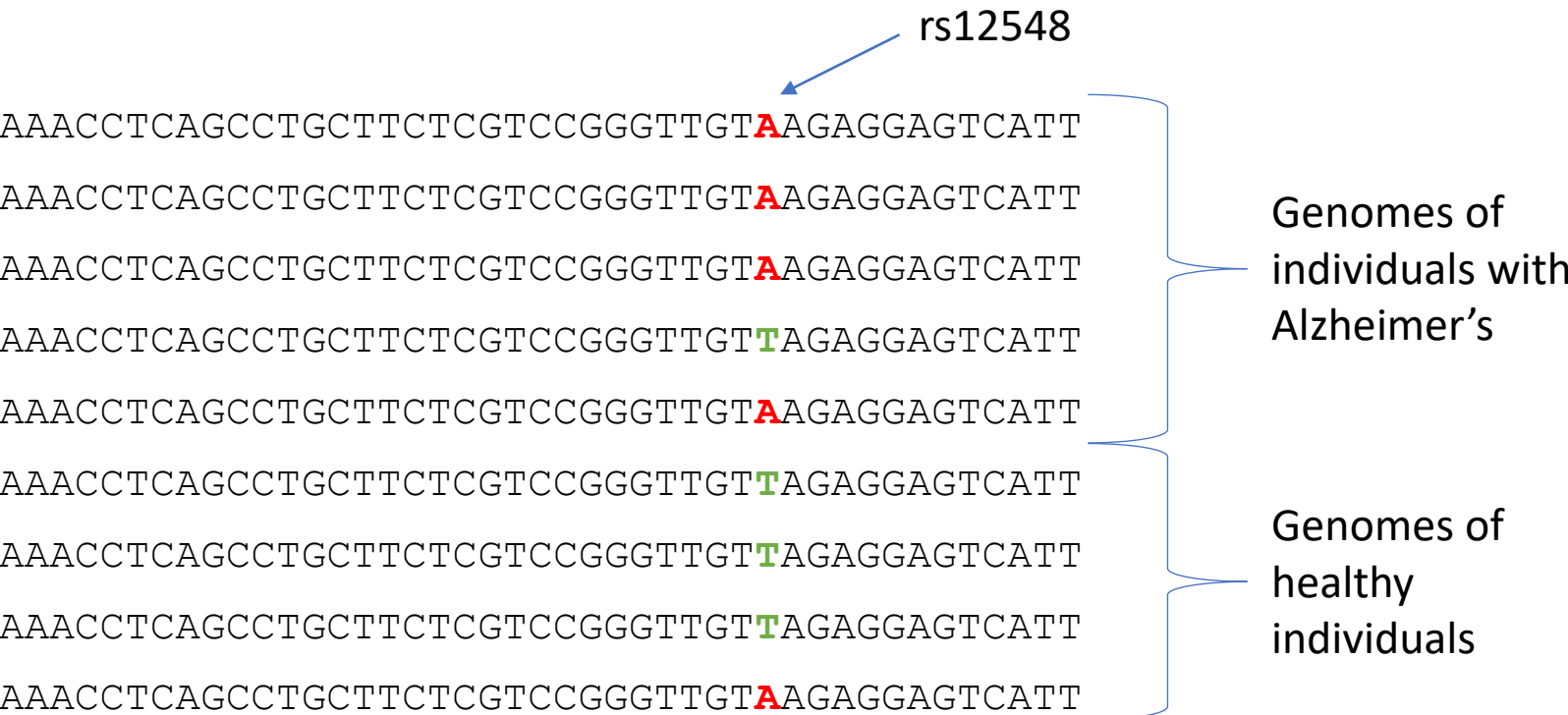


ref: The Post-GWAS Era: From Association to Function
(<https://doi.org/10.1016/j.ajhg.2018.04.002>)



ref: New insights into the genetic etiology of Alzheimer's disease and related dementias (<https://www.nature.com/articles/s41588-022-01024-z>)

How does a genome-wide association study find variants?



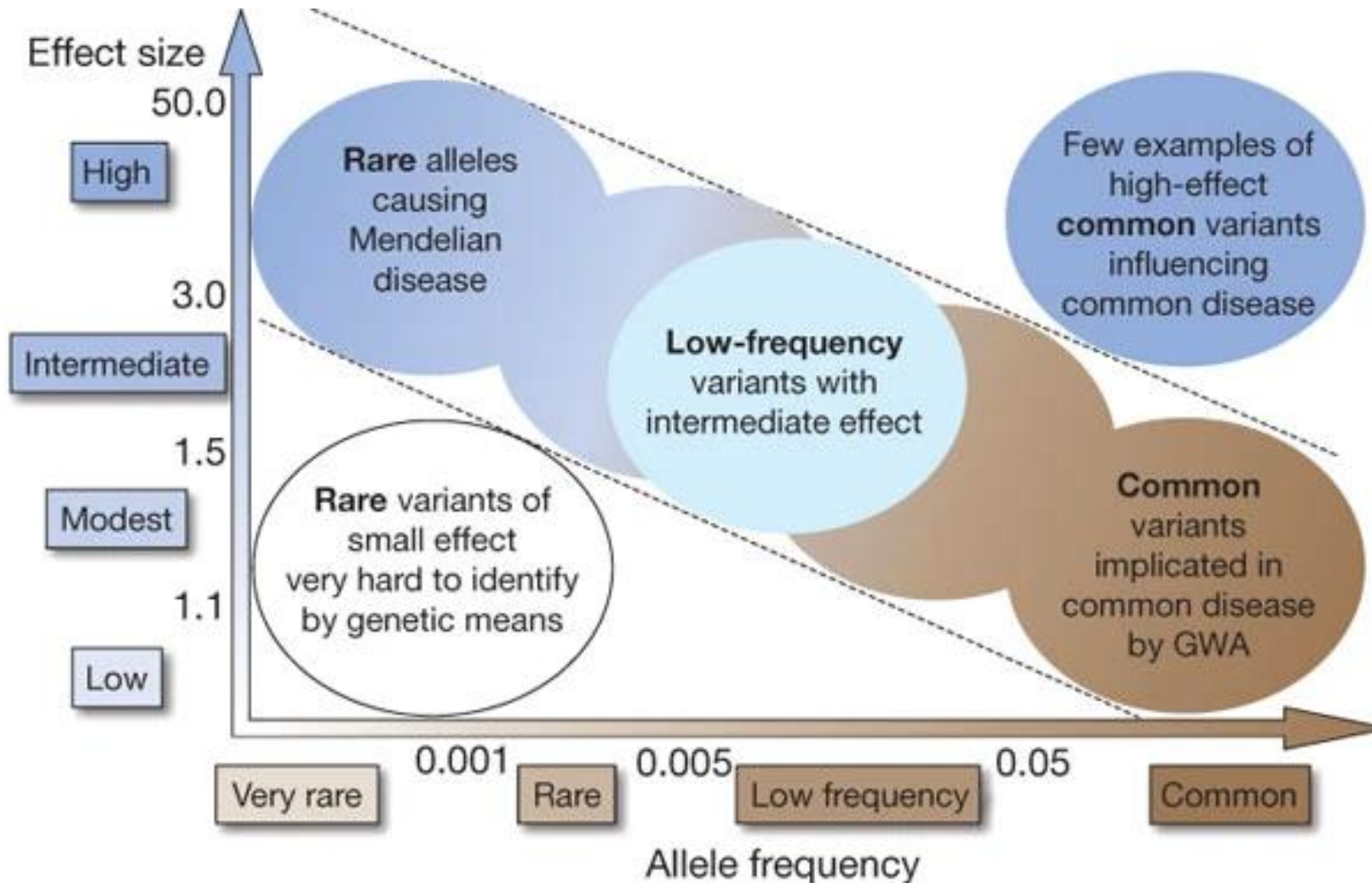
Chi-sq = 25.043
P = 5.61× 10⁻⁷

Observed	Allele A (N = 50)	Allele T (N = 40)
Alzheimer's (n =50)	40	10
Healthy (n = 40)	10	30

Expected	Allele A (N = 50)	Allele T (N = 40)
Alzheimer's (n =50)	27.8	22.2
Healthy (n = 40)	22.2	17.8

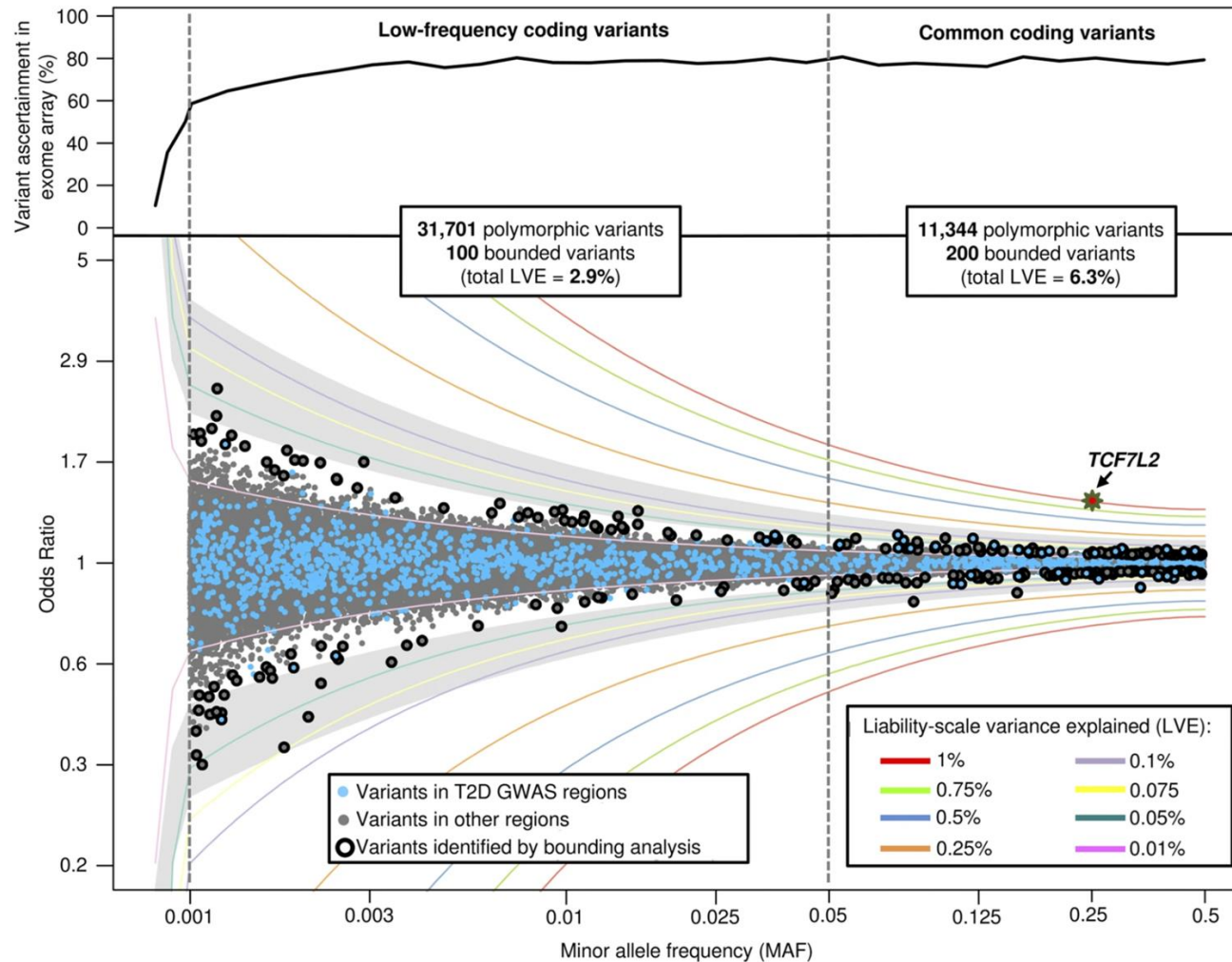
$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Most common variants have small effects



Most common variants have small effects

Type 2 diabetes



Fuchsberger, et al. 2016

Multiple hypothesis testing

- Usually, millions of variants are tested in GWAS studies.
- If we assume significance at $P = 0.05$, then you allow that there is a 5% chance that the association you found is random.
- So, if you find a 100 times a $P = 0.05$, then at least 5 of them are likely to be based on chance.
- Each study has hundreds of $P < 0.001$ purely by statistical chance (no real relationship to disease)
- “Genome-wide significance” often set at $P = 5 \times 10^{-8}$ ($= .05 / 1 \text{ million tests}$)

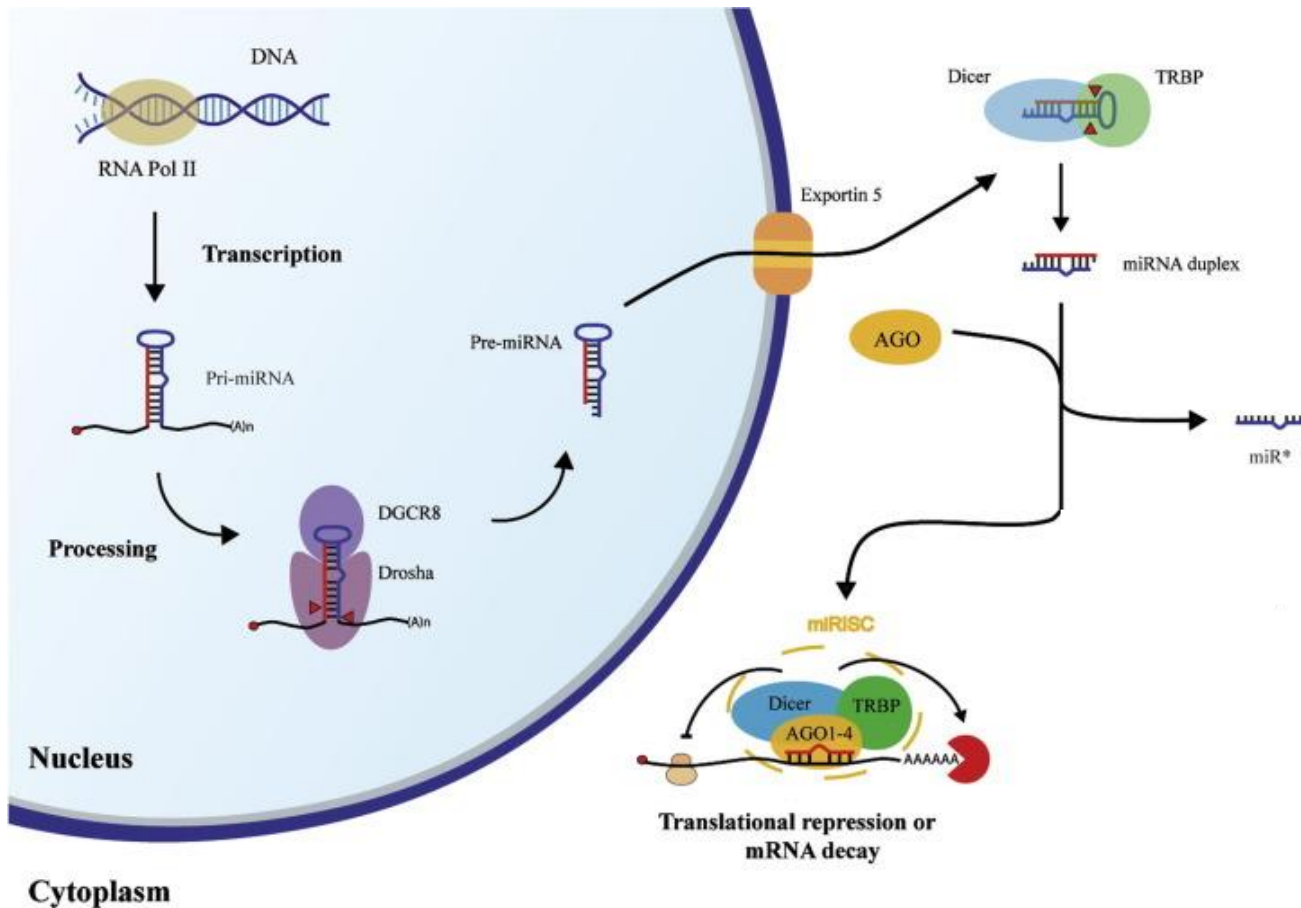
Disease associated variants

- Alzheimer's (~ 80 variants)
- Schizophrenia (~290 variants)
- Type 2 Diabetes (~240 variants)
- Chronic Kidney Disease (~40 variants)
- Bipolar disorder (~ 65 variants)
- Only 5% of disease-associated SNPs are in gene coding sequences
- So, what is happening with the other 95%???

Biological mechanisms of non-coding variants

1. Variant in 3'-UTR region, altering microRNA binding site
2. Variants altering splice sites
3. Variants altering transcription factor binding site
4. Variants altering DNA methylation

Variant in 3'-UTR region, altering microRNA binding site



Saraiva et al, 2017

	Predicted consequential pairing of target region (top) and miRNA (bottom)
Position 21-28 of HMGA2 3' UTR	5' ...GCCAACGUUCGAUUUCUACCUCA...
hsa-let-7f-5p	3' UUGAUUAGUUAGAUGAUGGAGU

<https://www.targetscan.org>

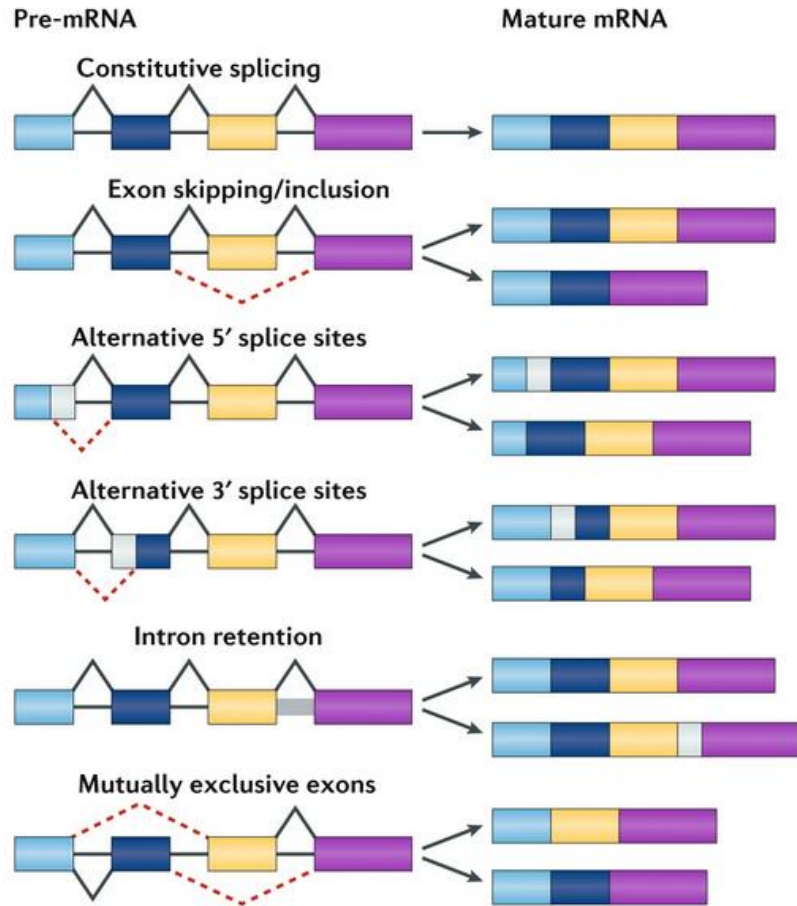
Total N	Significance threshold	N in 3'UTR (% to the UTR; % to the total)
57,671	p value<9×10 ⁻⁶	1,652 (2.9%)
	p value<5×10 ⁻⁸	1,041 (1.8%)

Steri et al, 2019

Variant in 3'-UTR region, altering microRNA binding site

- Alzheimer's Disease
 - 10 known variants (Ghanbari et al, 2016)
- Autism
 - 21 known variants (Vaishnavi et al, 2014)
- Metastasis in osteosarcoma
 - 1 known variant (Zhang S et al, 2016)

Variants altering splice sites



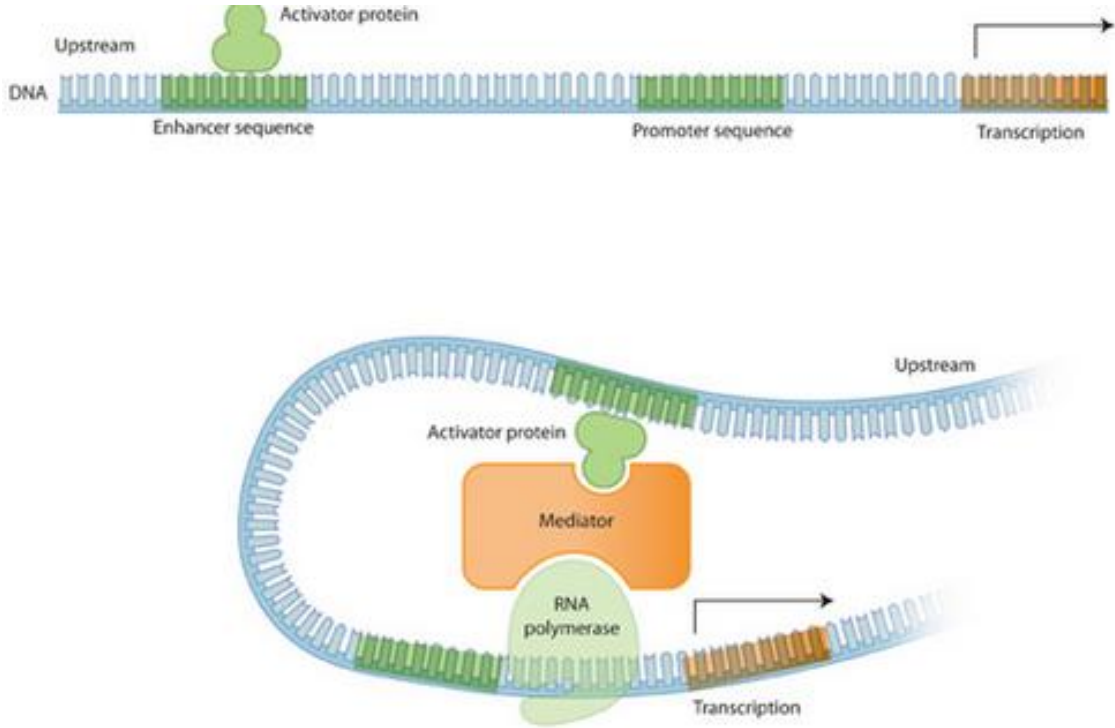
Frankiw et al, 2019

- >99% splice sites GT and AG in the donor and acceptor sites, respectively. (Kurmangaliyev et al, 2013)
- Most disease-causing mutations of splice sites occur at these dinucleotides (Kurmangaliyev et al, 2013)

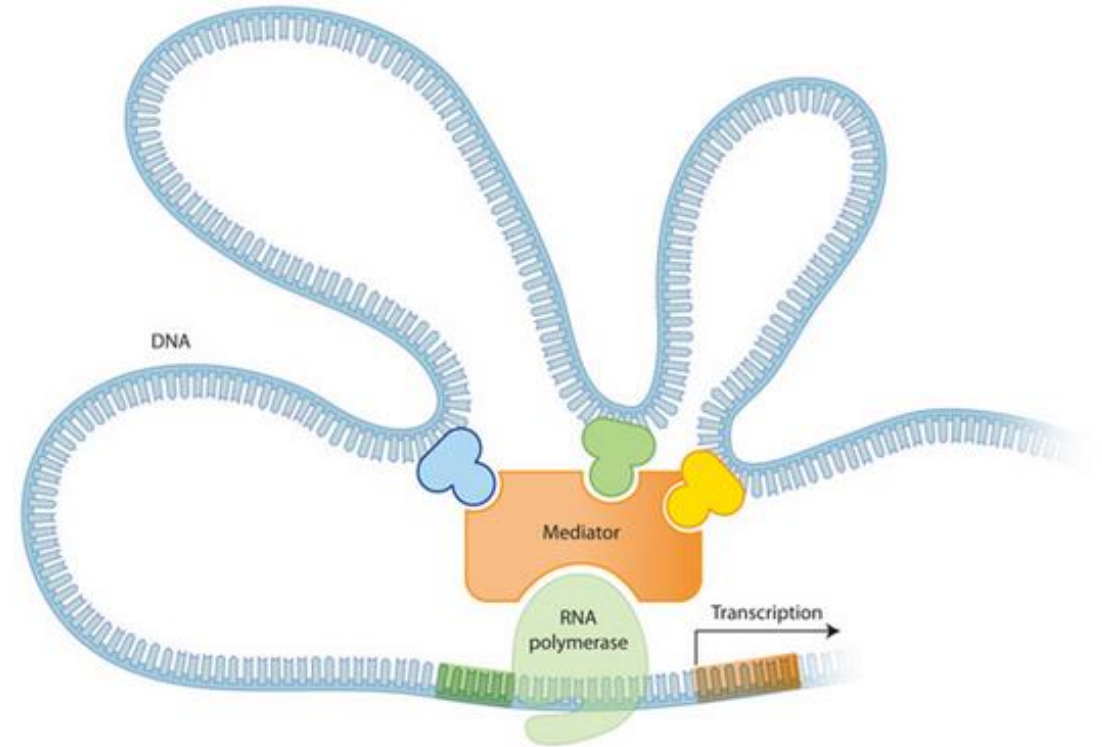
Variants altering splice sites

- Adipose-related traits
 - Seven exon-exon junctions in four genes (*AKTIP*, *DTNBP1*, *FTO* and *UBE2E1*) (**Ma et al, 2015**)
- Alzheimer's Disease
 - 21 genes (**Raj, et al, 2018**)
- Schizophrenia
 - Four genes (*NEK4*, *FXR1*, *SNAP91* or *APOPT1*)(**Takata et al, 2017**)
- Breast Cancer
 - Six genes (*BABAM1*, *DCLRE1B/PHTF1*, *PEX14*, *RAD51L1*, *SRGAP2D* and *STXBP4*) (**Caswell et al, 2015**)

Variants altering transcription factor binding site



www.nature.com/scitable



www.nature.com/scitable

Variants altering transcription factor binding site

- 8% of all polymorphisms reside in TFBS, yet it is estimated that of all disease associated polymorphisms 31% reside in TFBS
- BMI
 - *ARID5B, IRX3, IRX5* (Claussnitzer et al, 2015)
- Chronic kidney disease
 - *TCF7L2* (Köttgen et al, 2009)
- Diabetes
 - *NEUROD1, MTNR1B* (Lyssenko et al, 2009)
- Hypertension
 - *PHOX2, SCG2* (Wen et al, 2007)
- Many many more

Variants altering DNA methylation

- Methylation most often happens at CpG sites
- A methylated cytosine represses gene expression and can even lead to a more condensed chromatin structure.
- 23% of all polymorphisms are CpG site related SNPs (**Zhou et al, 2015**)
- Cardiovascular disease (**Jiantao Ma et al 2022**)
 - *APOB*
 - *SREBF1*
 - *PM20D1*

Summary of biological mechanisms

1. Variant in 3'-UTR region, altering microRNA binding site
2. Variants altering splice sites
3. Variants altering transcription factor binding site
4. Variants altering DNA methylation

Finding expression quantitative trait loci (eQTLs)

- Expression quantitative trait loci (eQTL) is a SNP that is associated with the expression of a gene.
- An eQTL always forms a pair with an eGene.
- eQTLs are usually identified using linear models.

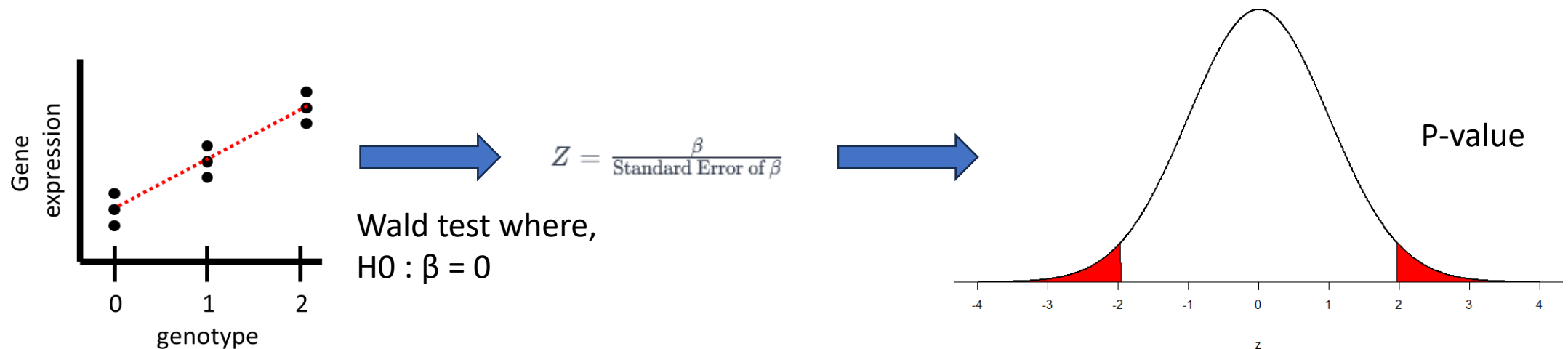
Finding expression quantitative trait loci (eQTLs)

1. Encode genotypes as 0,1,2

- A/A = 0
- A/T = 1
- T/T = 2

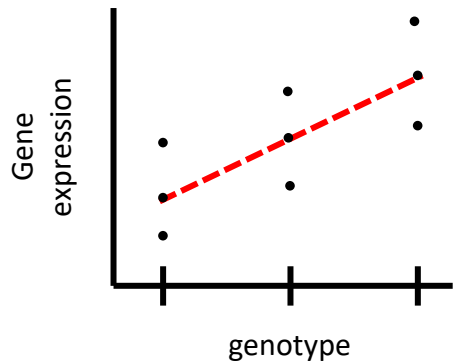
2. Associate genotype with gene expression (often linear model)

- expression = $\beta G + e$



Finding expression quantitative trait loci (eQTLs)

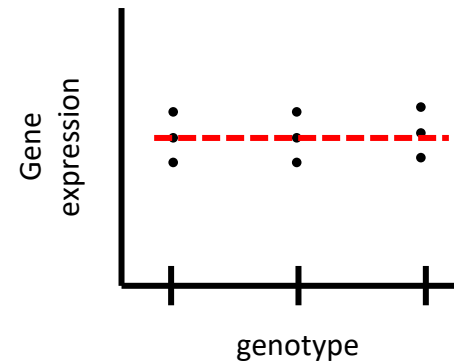
$$Z = \frac{\beta}{\text{Standard Error of } \beta}$$



$$\beta = 1$$
$$\text{SE } \beta = 0.33$$

$$Z: 1 / 0.33 = \mathbf{3}$$

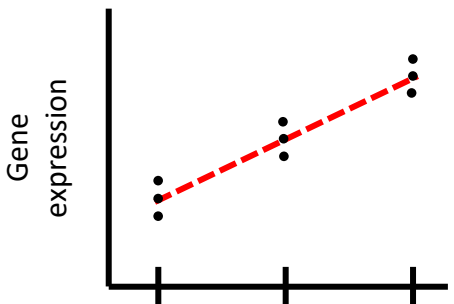
Genotype seems to be predictive of gene expression, but there is still some variability. Small p-value



$$\beta = 0$$
$$\text{SE } \beta = 0.33$$

$$Z: 0 / 0.33 = \mathbf{0}$$

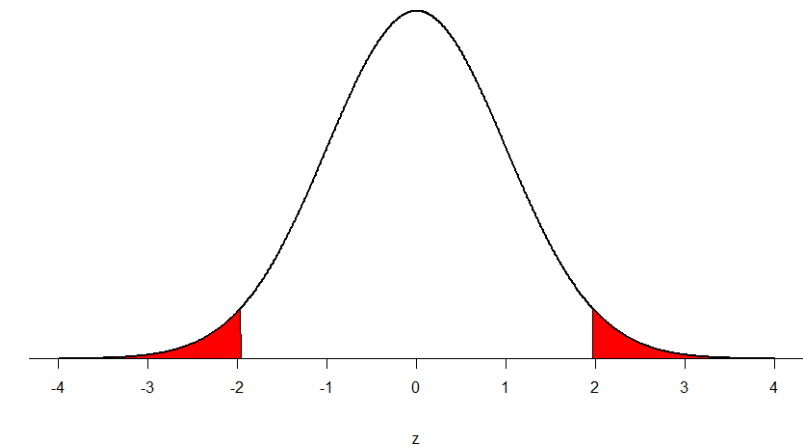
Genotype says nothing about gene expression.
P-value = 1



$$\beta = 1$$
$$\text{SE } \beta = 0.2$$

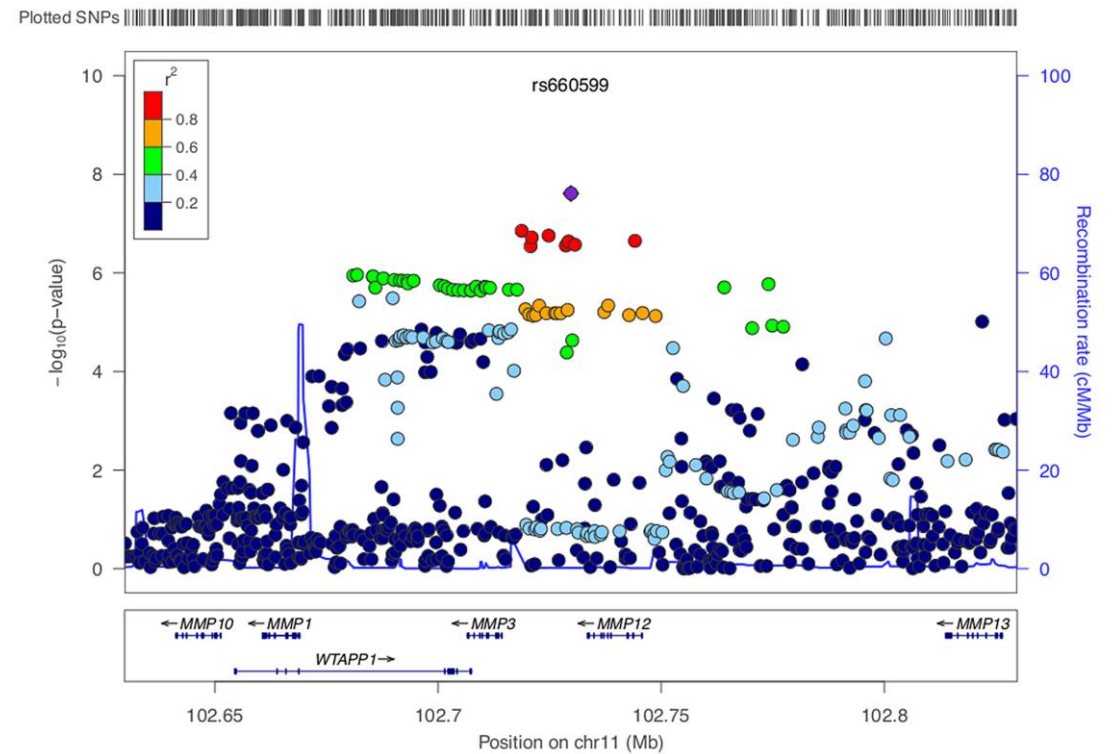
$$Z: 1 / 0.20 = \mathbf{5}$$

Genotype seems to be predictive of gene expression, very little variability. Very small p-value

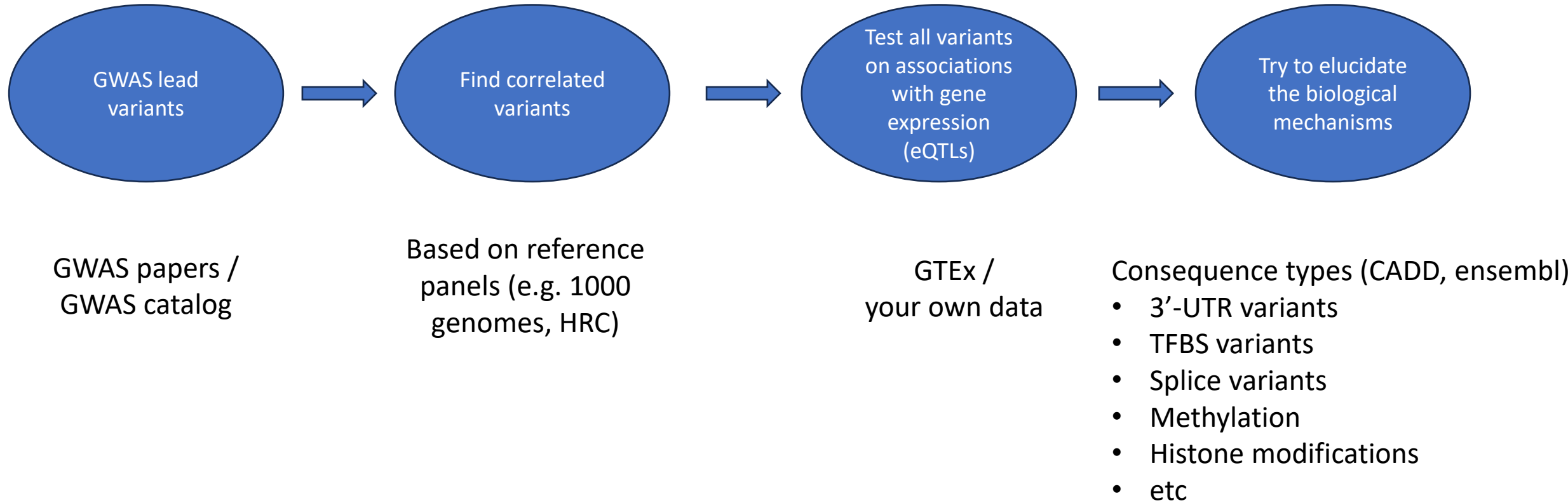


Finding expression quantitative trait loci (eQTLs)

- Usually, for every gene SNPs within 1MB of the TSS of the respective gene are tested. (Why? See biological mechanism.)
- Because many SNPs are correlated with each other, a gene is often associated with a whole LD-block.
- Before multiple testing, “independent signals” are identified using clumping.
- So instead of correcting for the total number of SNPs, you correct for the total number of “independent signals”



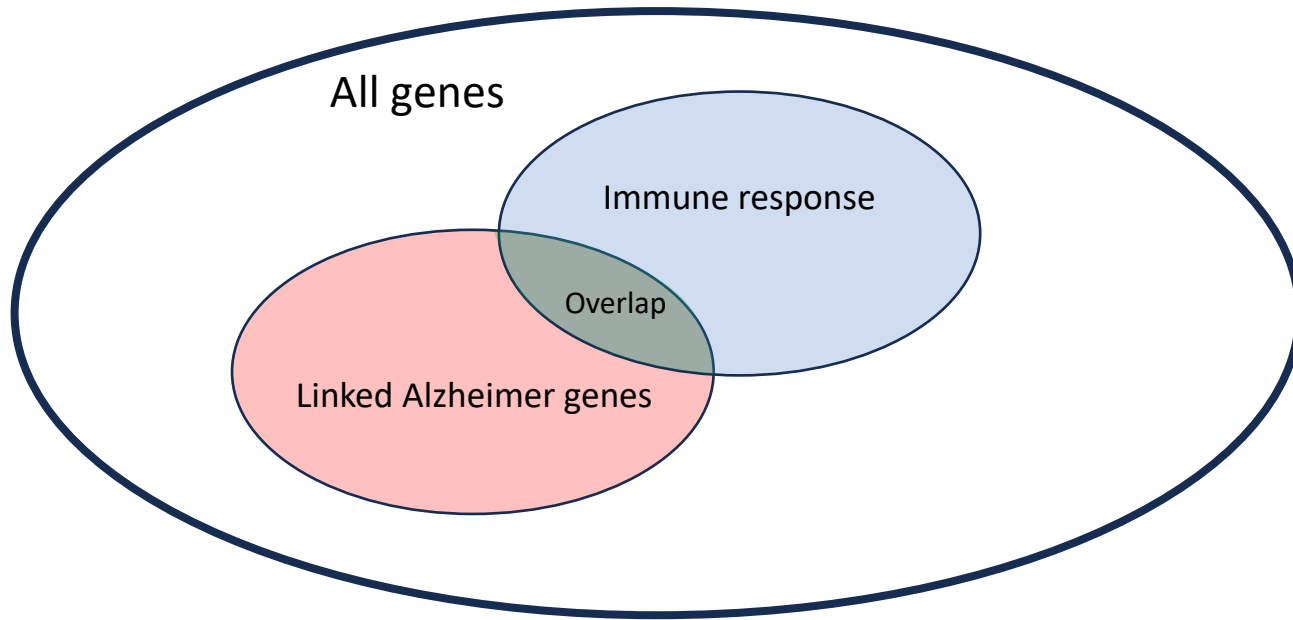
Combining GWASs with eQTLs



Gene set enrichment analysis

- Assume now that we linked all Alzheimer's variants to genes.
- What could be next thing that we would want to investigate?
- We can test whether the genes are involved in the same pathway or active in the same biological process, molecular function, cellular component or are they interacting with the same metabolite, and many more!

Gene set enrichment analysis



GO:0006955   

immune response

Biological Process

Definition ([GO:0006955 GONUTS page](#))

Any immune system process that functions in the calibrated response of an organism to a potential internal or invasive threat.

408,913 annotations

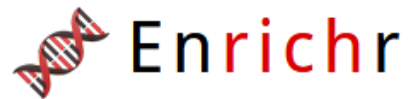
<https://www.ebi.ac.uk/QuickGO/term/GO:0006955>

Observed	Not Alzheimer Genes (N = 50)	Alzheimer Genes (N = 40)
Not in immune (n = 50)	40	10
In Immune (n = 40)	10	30

Expected	Not Alzheimer Genes (N = 50)	Alzheimer Genes (N = 40)
Not in immune (n = 50)	27.8	22.2
In Immune (n = 40)	22.2	17.8

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Gene set enrichment analysis

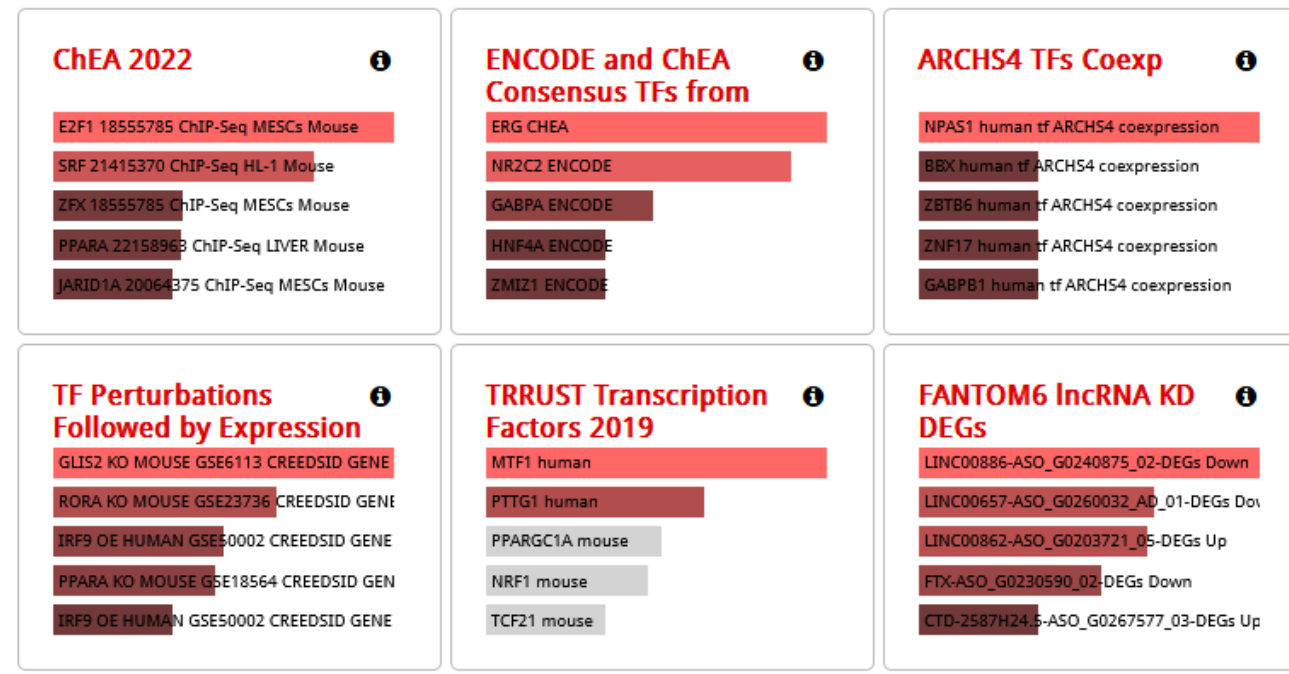


[Login](#) | [Register](#)

Transcription Pathways Ontologies Diseases/Drugs Cell Types Misc Legacy Crowd

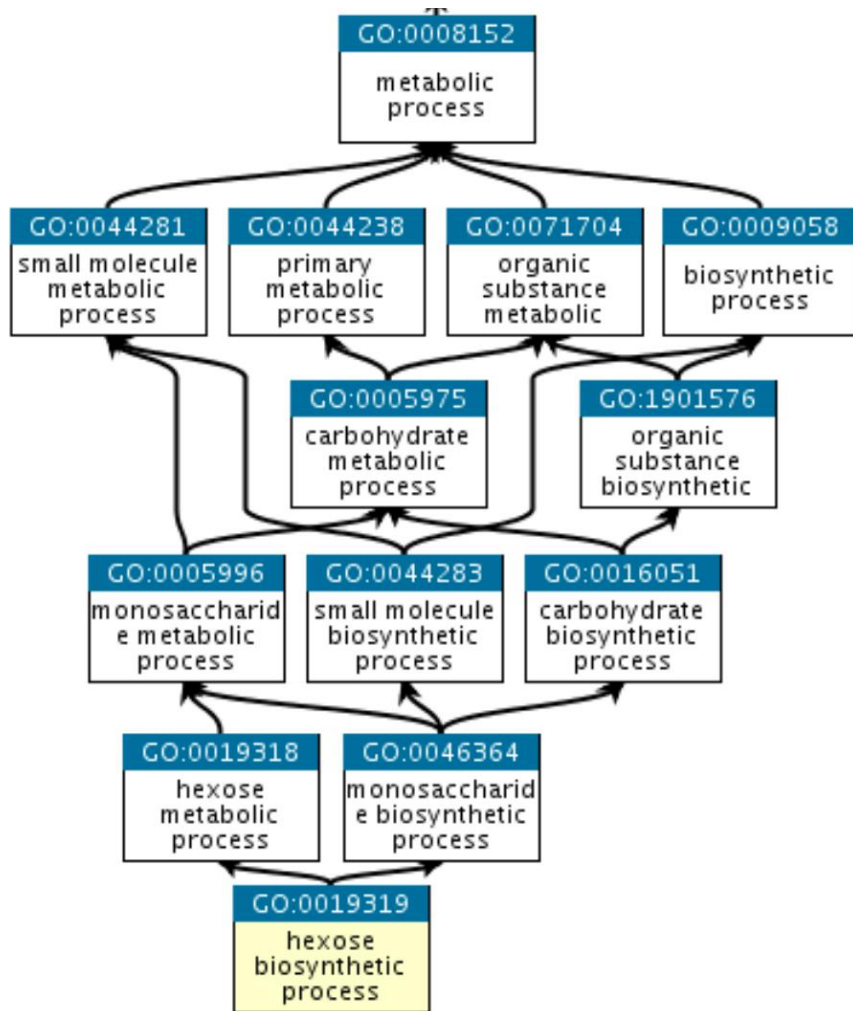
225 different libraries

Description Sample gene set (375 genes)  



<https://maayanlab.cloud/Enrichr/>

Gene ontology



Hierarchy of terms

- Biological processes
- Cellular Component
- Molecular function

Very extensive (~7.028 terms)

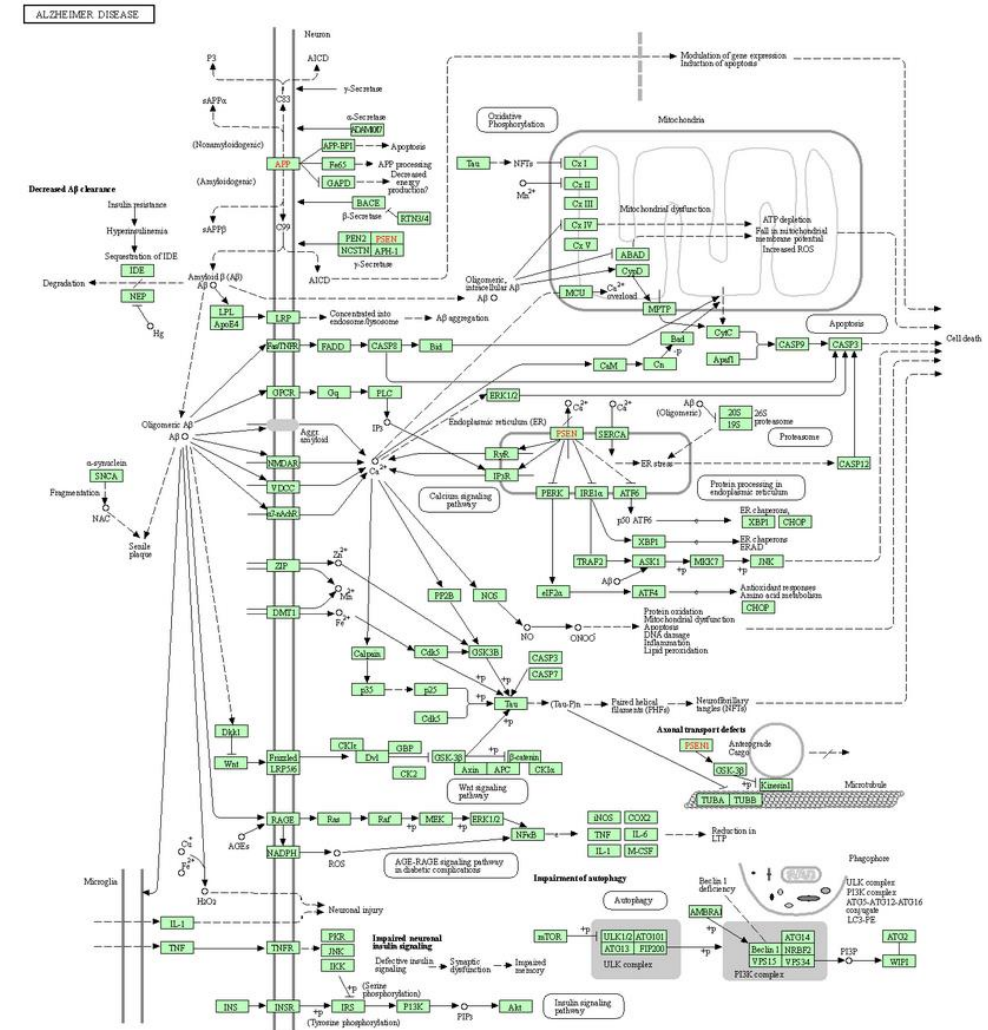
But a lot of overlap between terms.

KEGG PATHWAY Database

- “KEGG PATHWAY is a collection of manually drawn pathway maps representing our knowledge of the molecular interaction, reaction and relation networks for:”

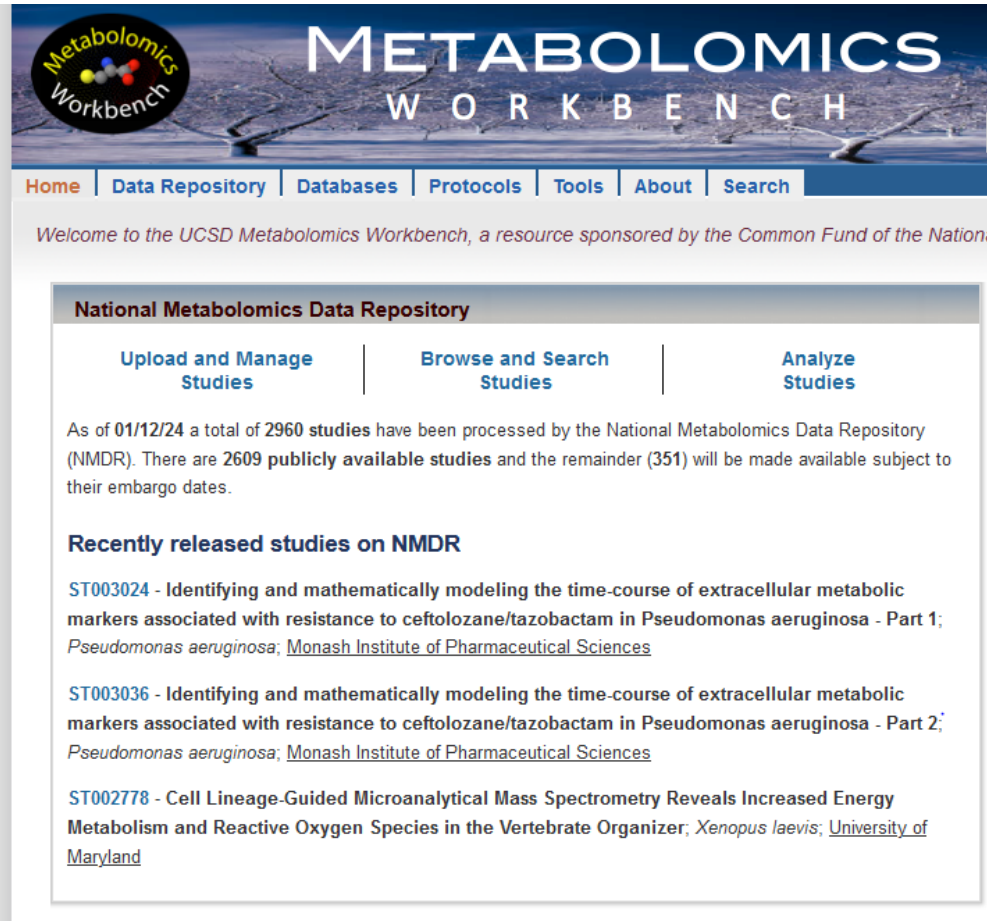
~320 terms

1. Metabolism
2. Genetic Information Processing
3. Environmental Information Processing
4. Cellular Processes
5. Organismal Systems
6. Human Diseases
7. Drug Development



Metabolomics Workbench Metabolite Database

- Database of protein metabolite interactions



The screenshot shows the homepage of the Metabolomics Workbench. At the top, there is a logo for 'Metabolomics Workbench' and the text 'METABOLOMICS WORKBENCH'. Below this is a navigation bar with links: Home, Data Repository, Databases, Protocols, Tools, About, and Search. A welcome message follows: 'Welcome to the UCSD Metabolomics Workbench, a resource sponsored by the Common Fund of the National Institutes of Health'. The main content area is titled 'National Metabolomics Data Repository' and features three columns: 'Upload and Manage Studies', 'Browse and Search Studies', and 'Analyze Studies'. A text block states: 'As of 01/12/24 a total of 2960 studies have been processed by the National Metabolomics Data Repository (NMDR). There are 2609 publicly available studies and the remainder (351) will be made available subject to their embargo dates.' Below this is a section titled 'Recently released studies on NMDR' with three entries: ST003024, ST003036, and ST002778, each with a brief description and a link to the study page.

National Metabolomics Data Repository

[Upload and Manage Studies](#) | [Browse and Search Studies](#) | [Analyze Studies](#)

As of 01/12/24 a total of 2960 studies have been processed by the National Metabolomics Data Repository (NMDR). There are 2609 publicly available studies and the remainder (351) will be made available subject to their embargo dates.

Recently released studies on NMDR

[ST003024](#) - Identifying and mathematically modeling the time-course of extracellular metabolic markers associated with resistance to ceftolozane/tazobactam in *Pseudomonas aeruginosa* - Part 1; *Pseudomonas aeruginosa*; [Monash Institute of Pharmaceutical Sciences](#)

[ST003036](#) - Identifying and mathematically modeling the time-course of extracellular metabolic markers associated with resistance to ceftolozane/tazobactam in *Pseudomonas aeruginosa* - Part 2; *Pseudomonas aeruginosa*; [Monash Institute of Pharmaceutical Sciences](#)

[ST002778](#) - Cell Lineage-Guided Microanalytical Mass Spectrometry Reveals Increased Energy Metabolism and Reactive Oxygen Species in the Vertebrate Organizer; *Xenopus laevis*; [University of Maryland](#)

233 terms


Co-expression QTLs

Alzheimer's Disease associated with emerging and disrupted protein correlation patterns in Gyrus Temporalis Medialis through genetic variants and pathology

Gerard A. Bouland^{1,2}, Niccolò Tesi^{1,3,4}, Meng Zhang¹, Andrea B. Ganz⁴, Marc Hulsman, Sven van der Lee, Marieke Graat, Annemieke Rozemuller, Martijn Huisman, Natasja van Schoor, Wiesje van der Flier, Jeroen Hoozemans, August B Smit⁴, Marcel J.T. Reinders^{1,2,6*}, Henne Holstege^{1,3,4*}

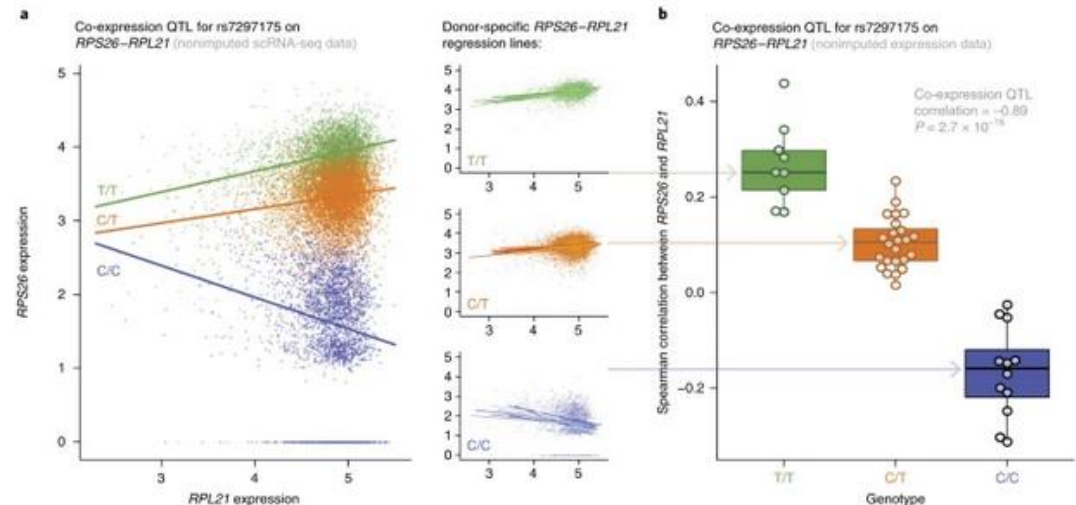
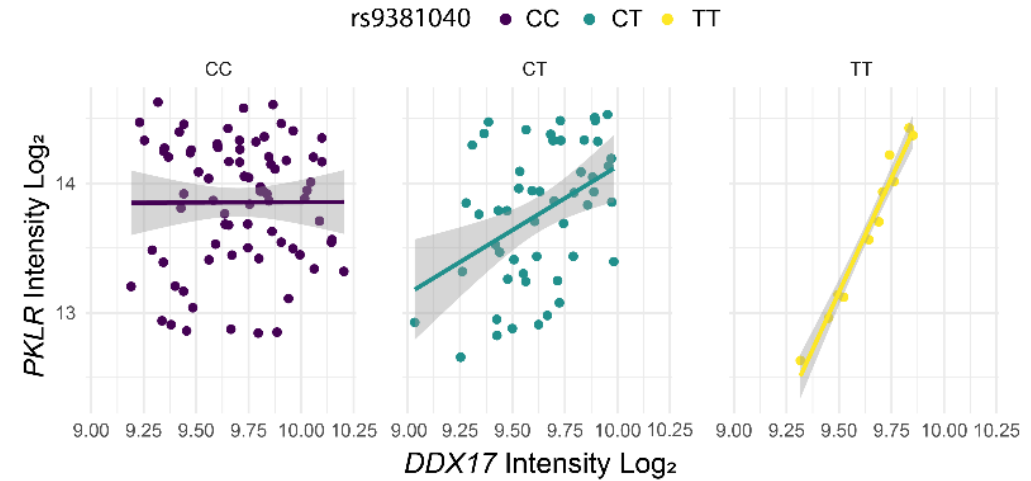
Letter | [Published: 02 April 2018](#)

Single-cell RNA sequencing identifies celltype-specific cis-eQTLs and co-expression QTLs

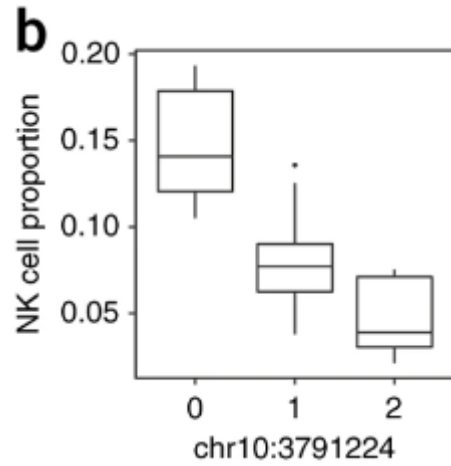
[Monique G. P. van der Wijst](#), [Harm Brugge](#), [Dylan H. de Vries](#), [Patrick Deelen](#), [Morris A. Swertz](#), [LifeLines Cohort Study](#), [BIOS Consortium](#) & [Lude Franke](#) 

[Nature Genetics](#) **50**, 493–497 (2018) | [Cite this article](#)

24k Accesses | **170** Citations | **91** Altmetric | [Metrics](#)



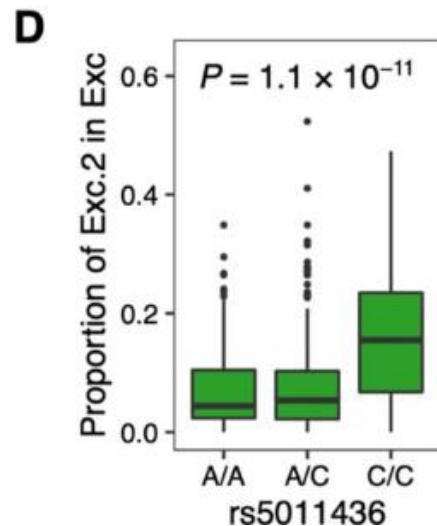
Cell state abundance quantitative trait loci (csaQTLs)






Published: 11 December 2017

Multiplexed droplet single-cell RNA-sequencing using natural genetic variation

Hyun Min Kang , Meena Subramaniam, Sasha Targ, Michelle Nguyen, Lenka Maliskova, Elizabeth McCarthy, Eunice Wan, Simon Wong, Lauren Byrnes, Cristina M Lanata, Rachel E Gate, Sara Mostafavi, Alexander Marson, Noah Zaitlen, Lindsey A Criswell & Chun Jimmie Ye 

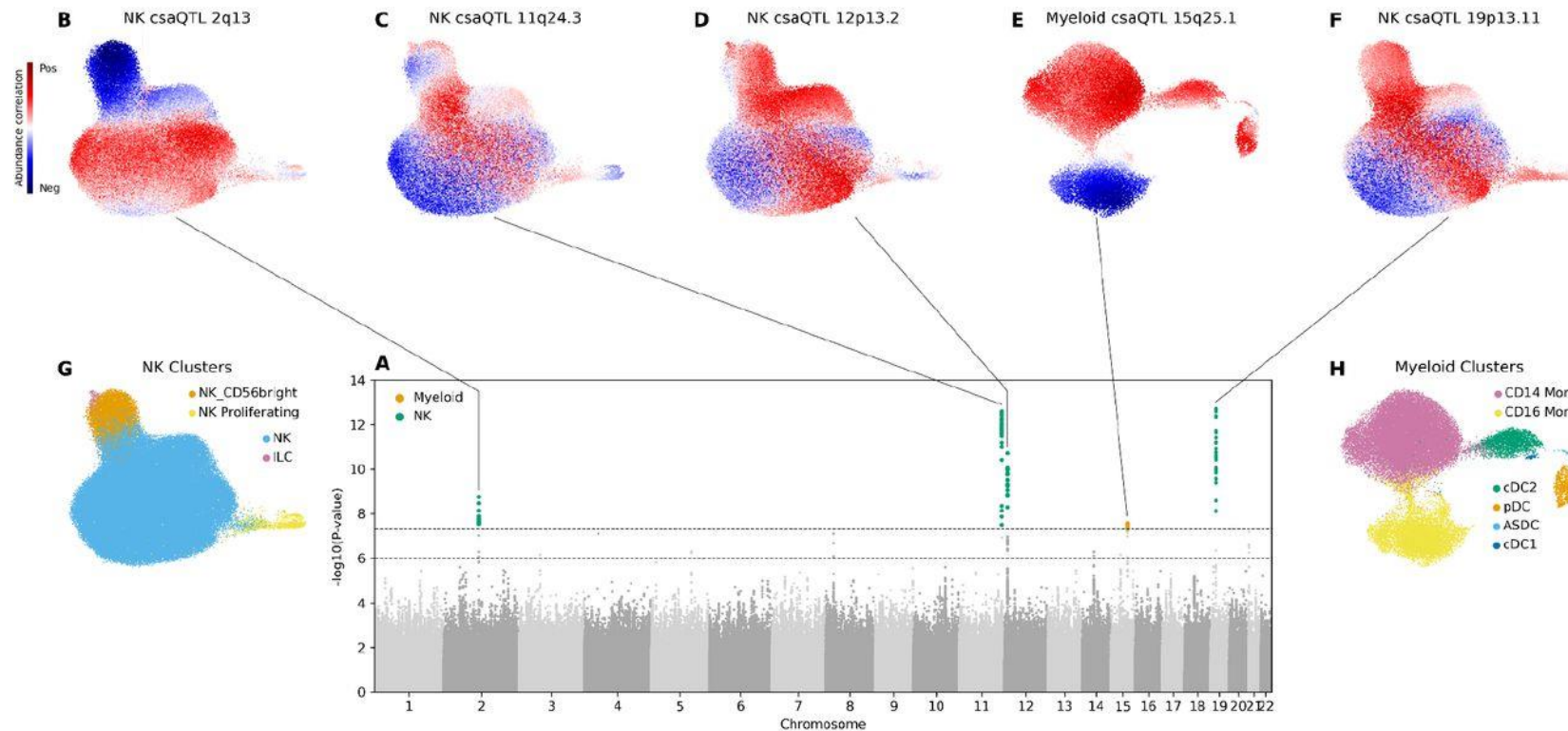


Cell-subtype specific effects of genetic variation in the aging and Alzheimer cortex


 Masashi Fujita, Zongmei Gao, Lu Zeng, Cristin McCabe, Charles C. White, Bernard Ng, Gilad Sahar Green, Orit Rozenblatt-Rosen, Devan Phillips, Liat Amir-Zilberstein, Hyo Lee, Richard V. Pearse II, Atlas Khan, Badri N. Vardarajan, Krzysztof Kiryluk,  Chun Jimmie Ye, Hans-Ulrich Klein, Gao Wang, Aviv Regev, Naomi Habib, Julie A. Schneider, Yanling Wang, Tracy Young-Pearse, Sara Mostafavi, David A. Bennett, Vilas Menon,  Philip L. De Jager

doi: <https://doi.org/10.1101/2022.11.07.515446>

Cell state abundance quantitative trait loci (csaQTLs)



Identifying genetic variants that influence the abundance of cell states in single-cell data

Laurie Rumker, Saori Sakaue, Yakir Reshef, Joyce B. Kang, Seyhan Yazar, Jose Alquicira-Hernandez, Cristian Valencia, Kaitlyn A Lagattuta, Annelise Mah-Som, Aparna Nathan, Joseph E. Powell, Po-Ru Loh,  Soumya Raychaudhuri

doi: <https://doi.org/10.1101/2023.11.13.566919>