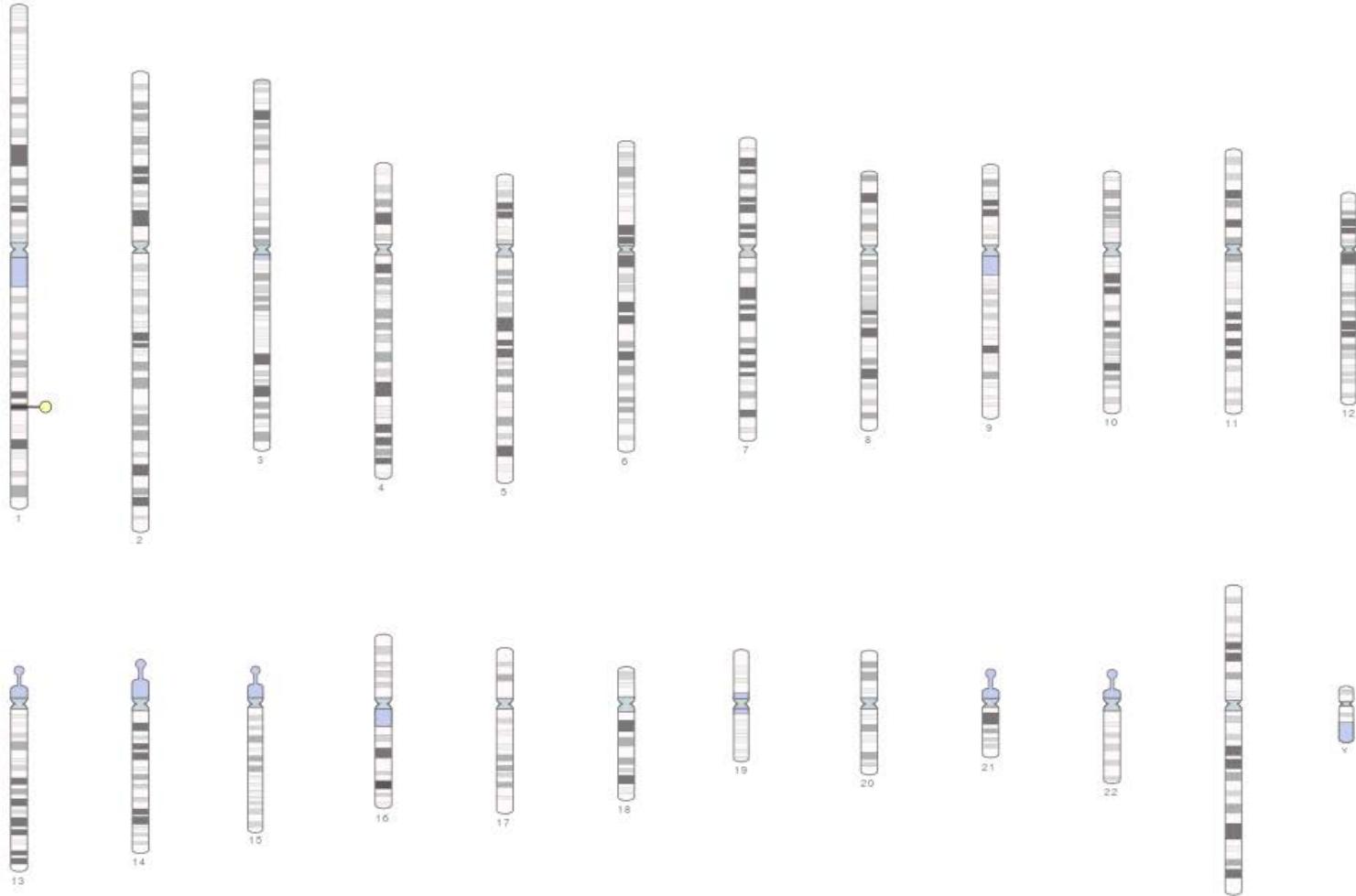


Effect of genotypes on gene expression

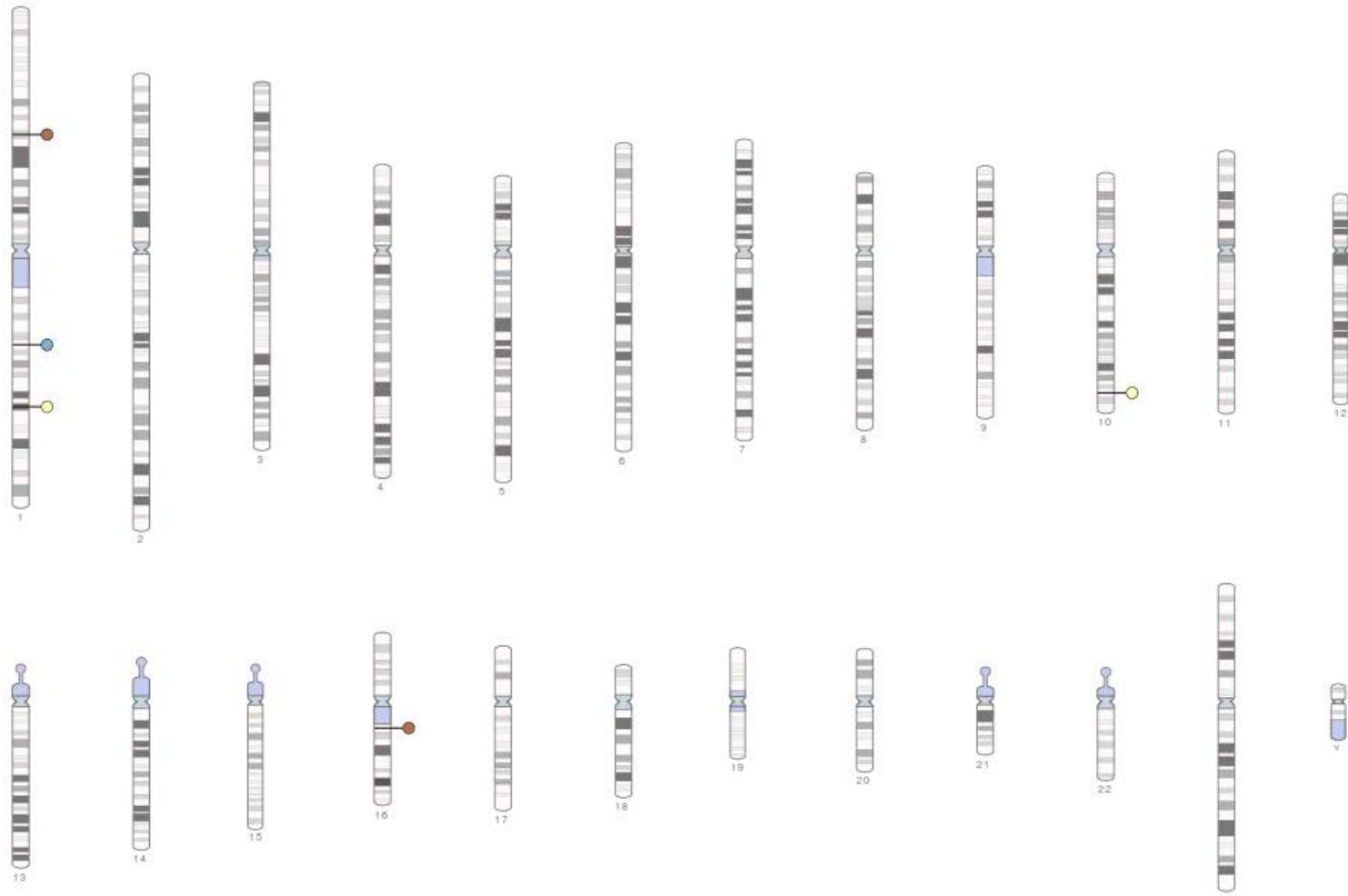
Patrick Deelen, 25-09-2018



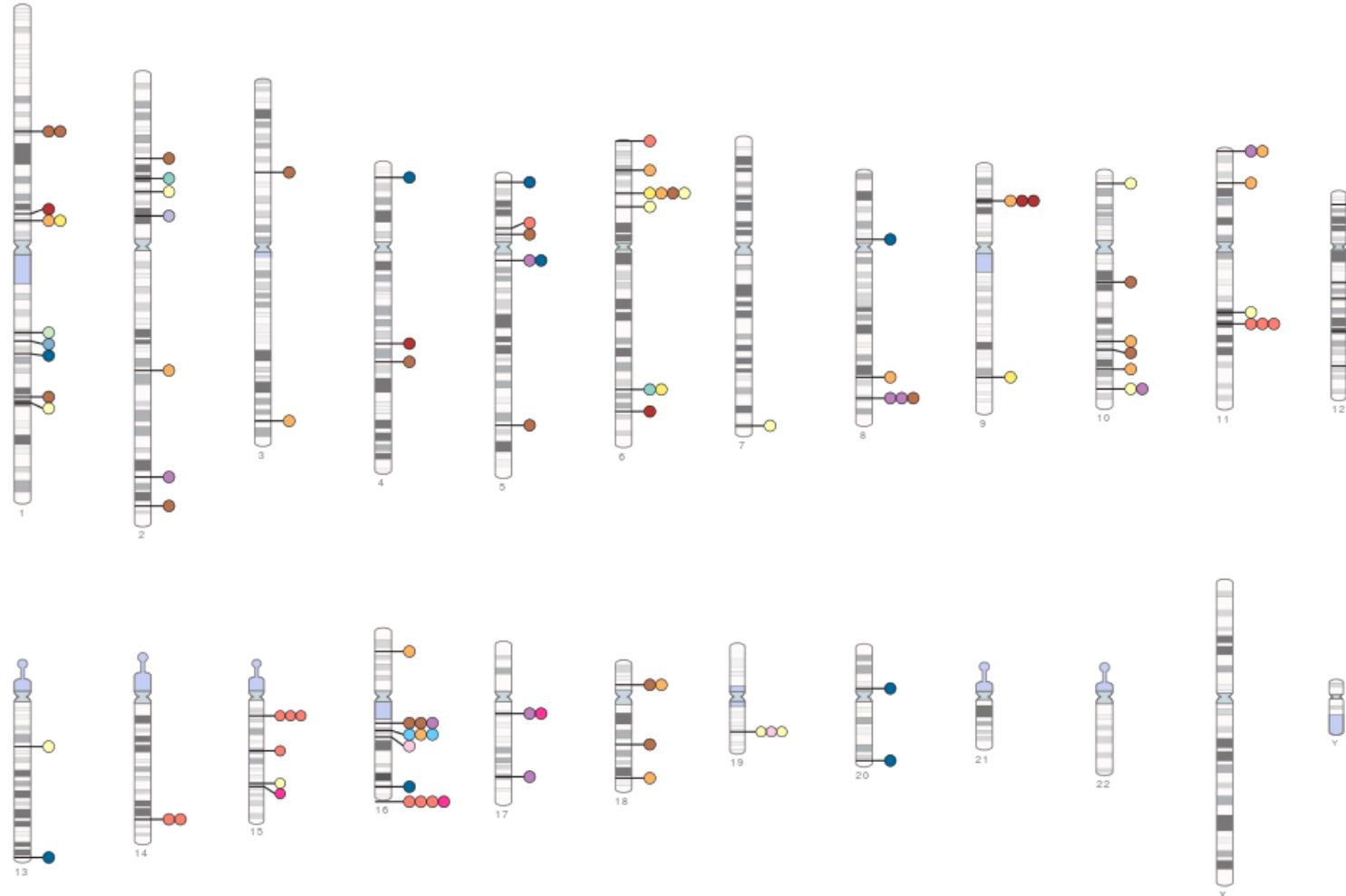
2005



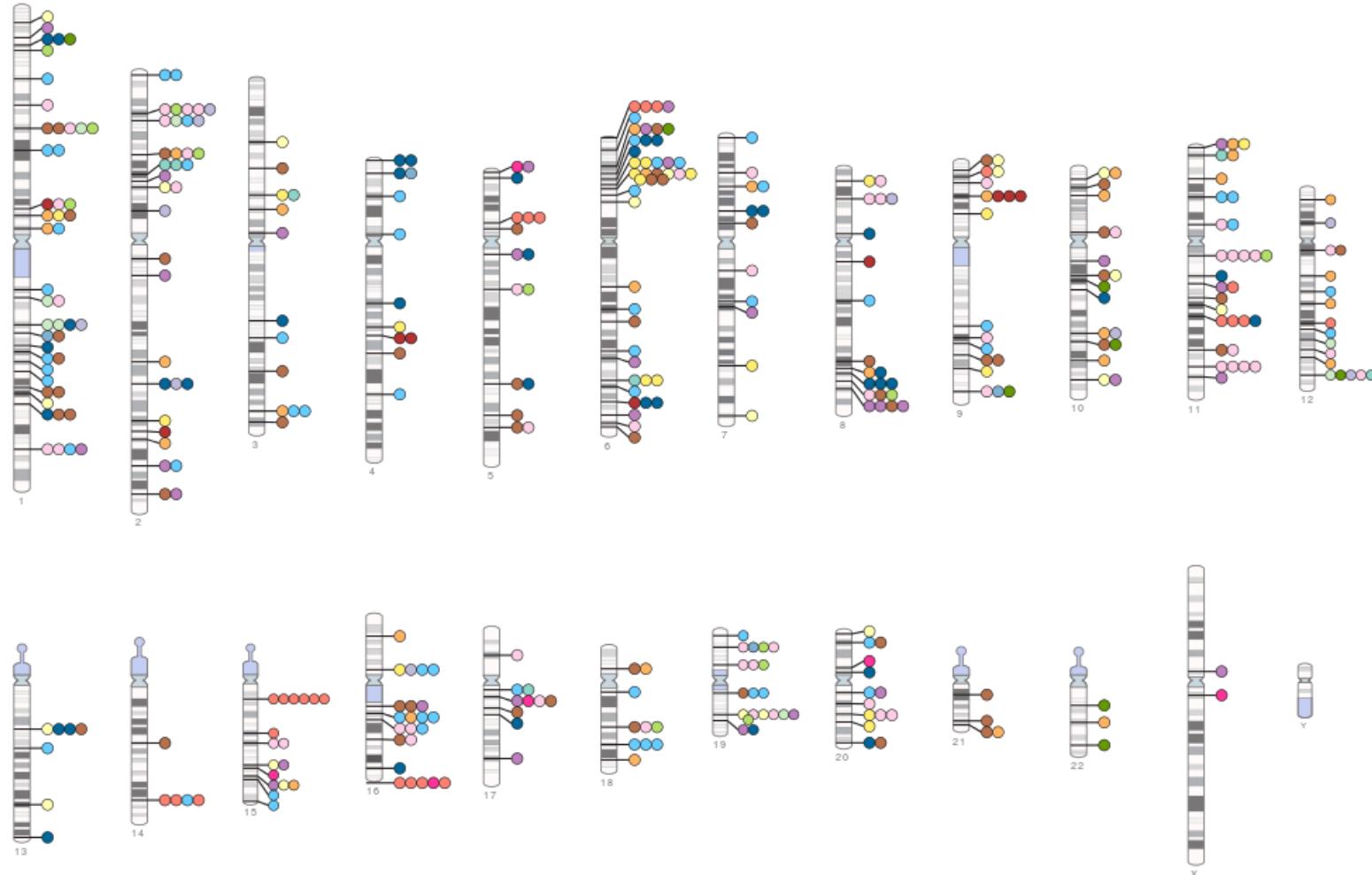
2006



2007



2008



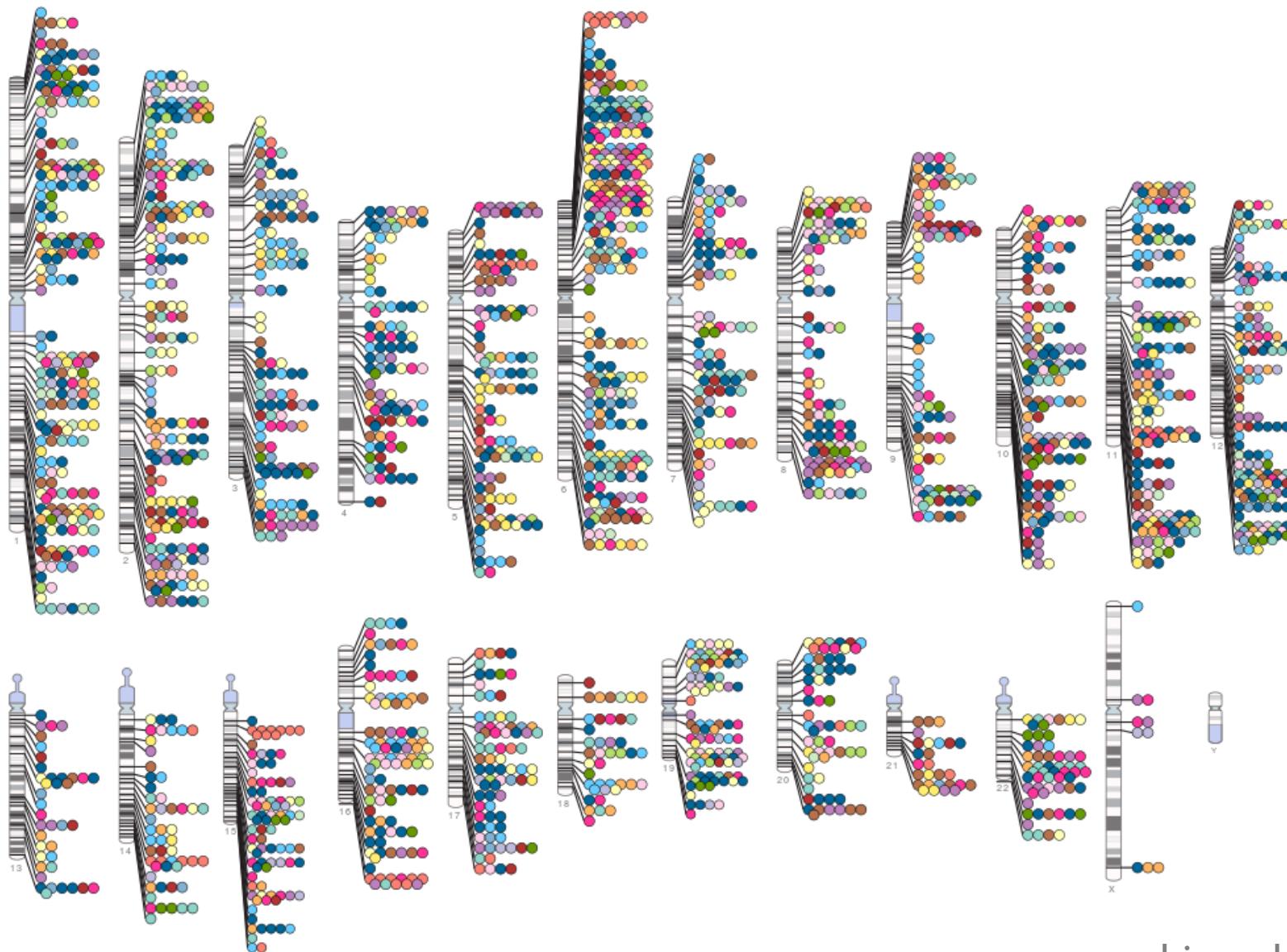
2009



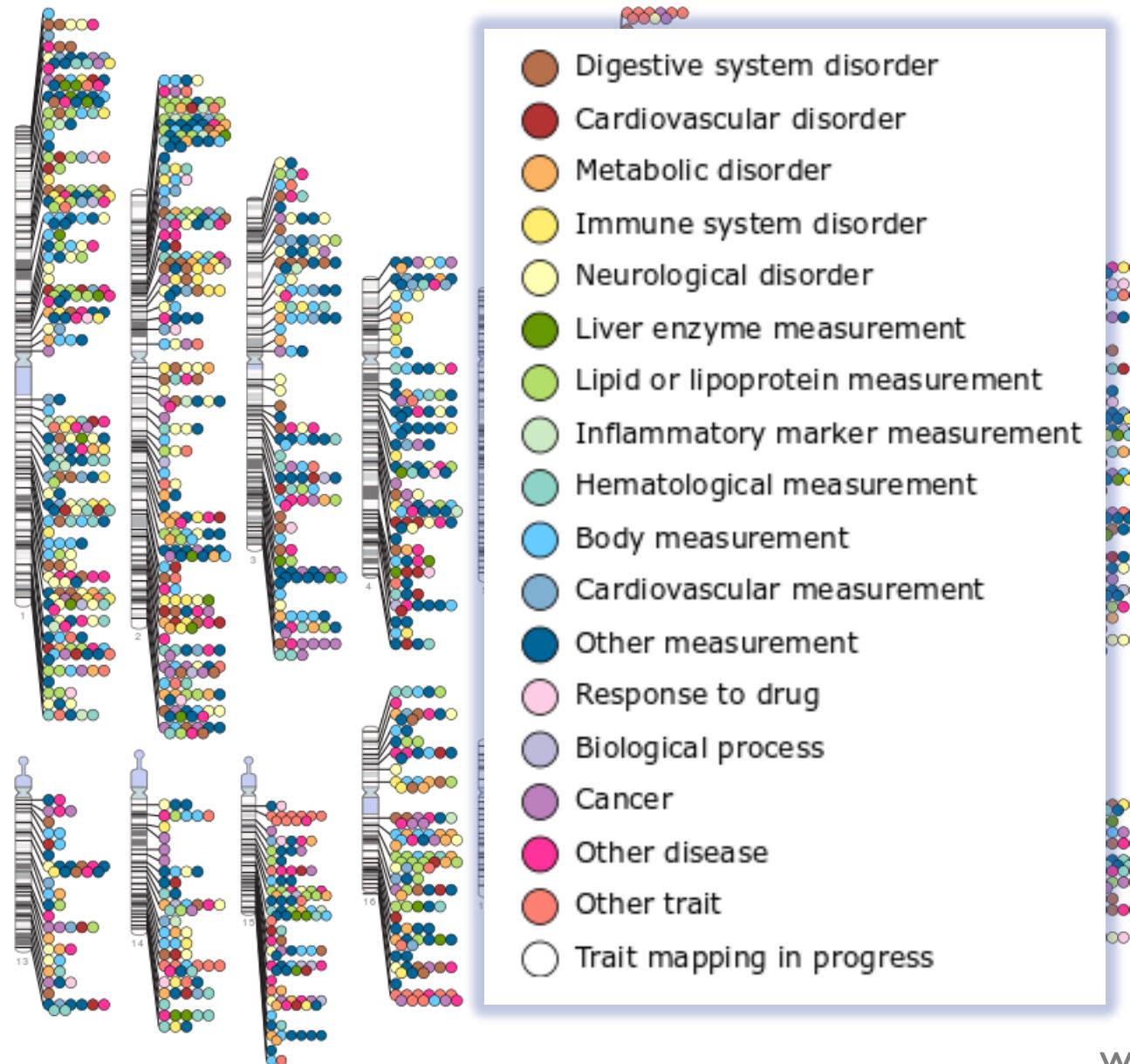
2010



2011



2012



GWAS variant



Causal variant

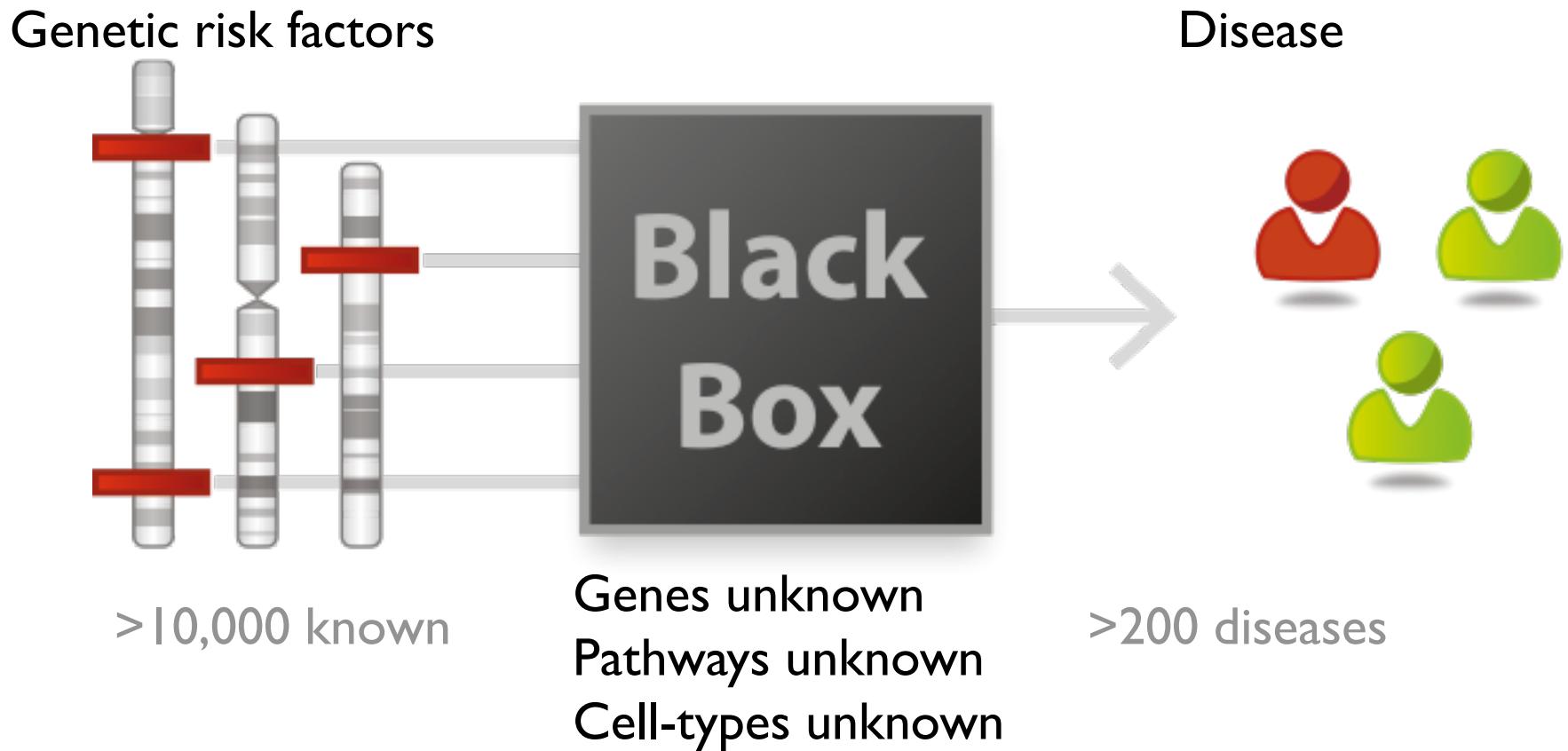


Affected gene



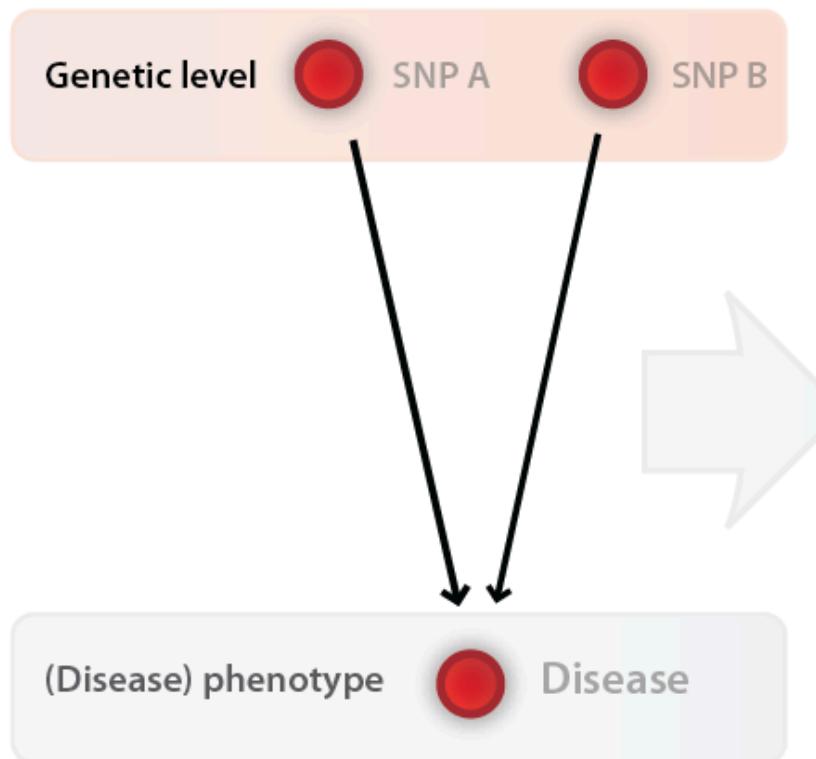
Biological
function

Problem of life science community

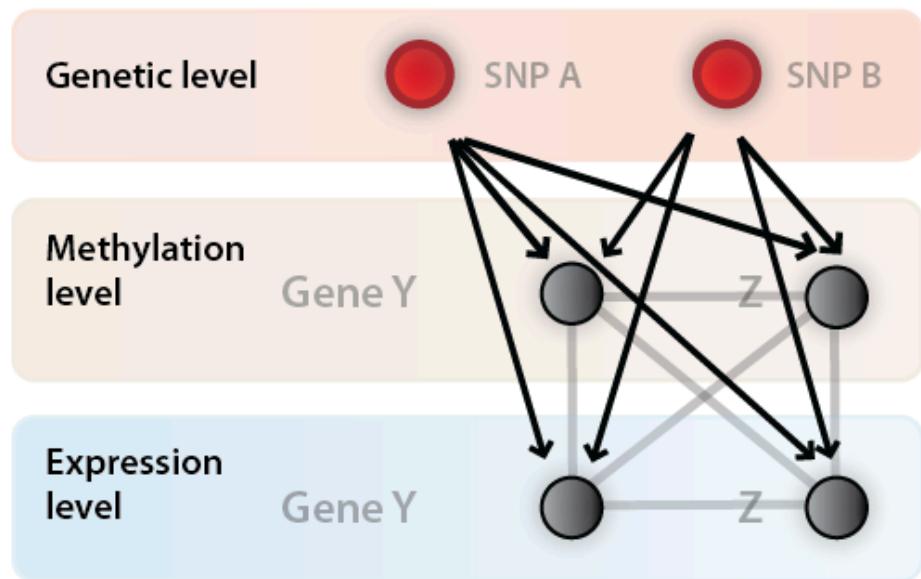


Why look at eQTL

Current situation

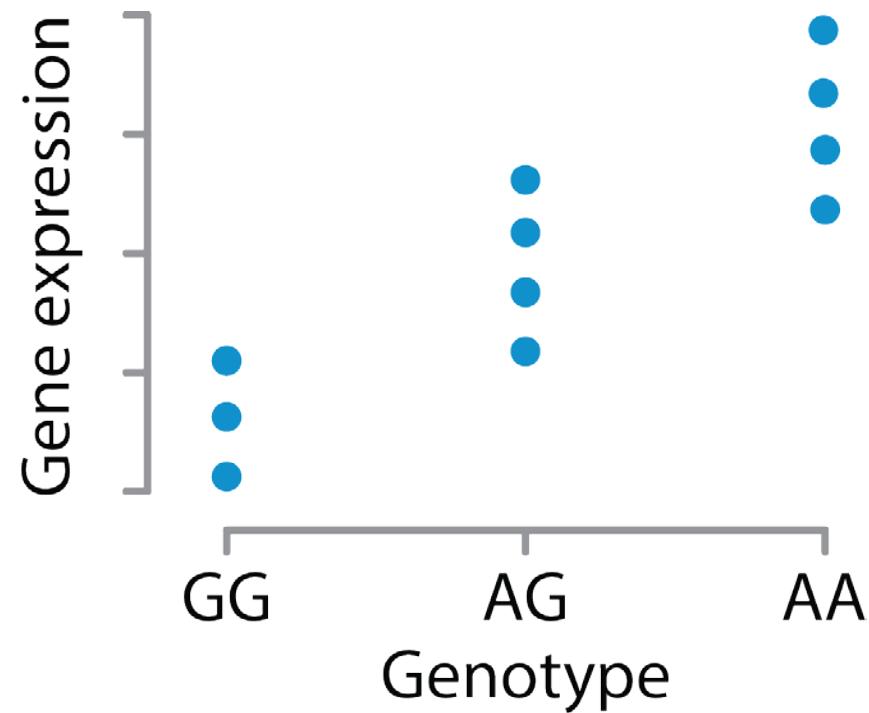


Identification of intermediate downstream affected genes



Moving from variants to genes using eQTLs

- ▶ Most GWAS hits do **not** change protein structure
- ▶ They often **do** affect expression levels
- ▶ Expression quantitative trait locus - eQTL
- ▶ Genetic variation correlated to expression
- ▶ Large dataset with for same individual
 - ▶ Genotypes
 - ▶ Expression levels



eQTL dataset

Genotypes

Variant	Sample1	Sample2	Sample3
SNP 1	A/A	A/T	T/T
SNP 2	T/T	T/T	T/C
SNP 3	G/C	G/C	G/G
SNP 4	G/G	G/G	G/G

Expression levels

Gene	Sample1	Sample2	Sample3
Gene 1	10	20	30
Gene 2	20	10	18
Gene 3	10	11	12
Gene 4	10	10	10

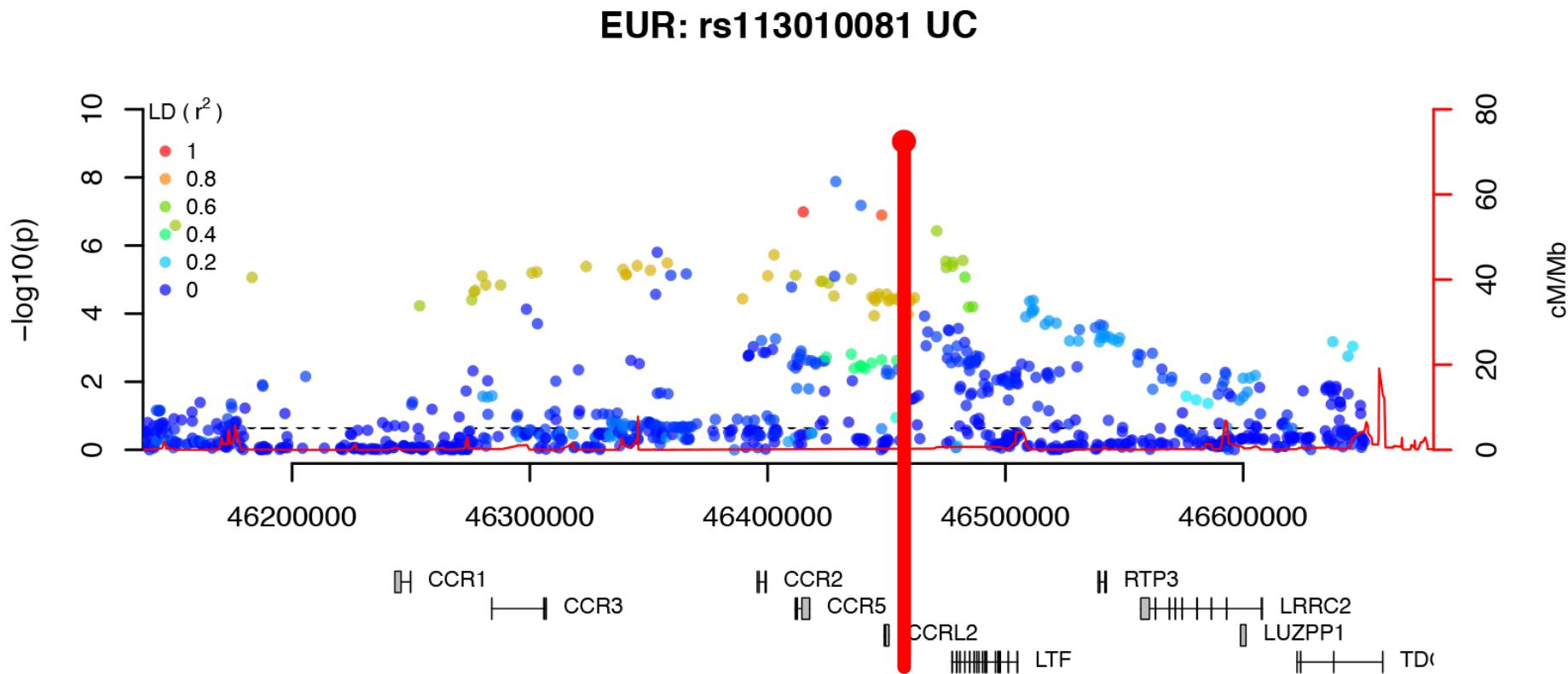
Multiple testing

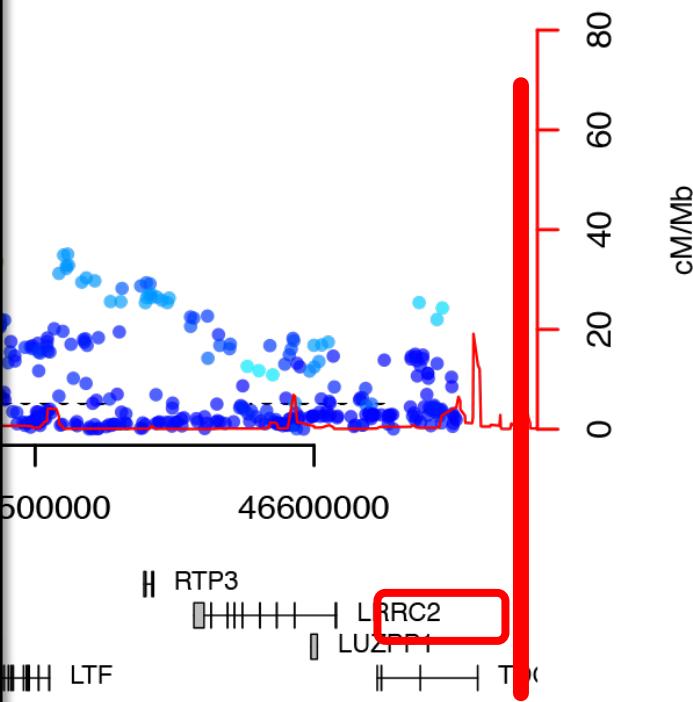
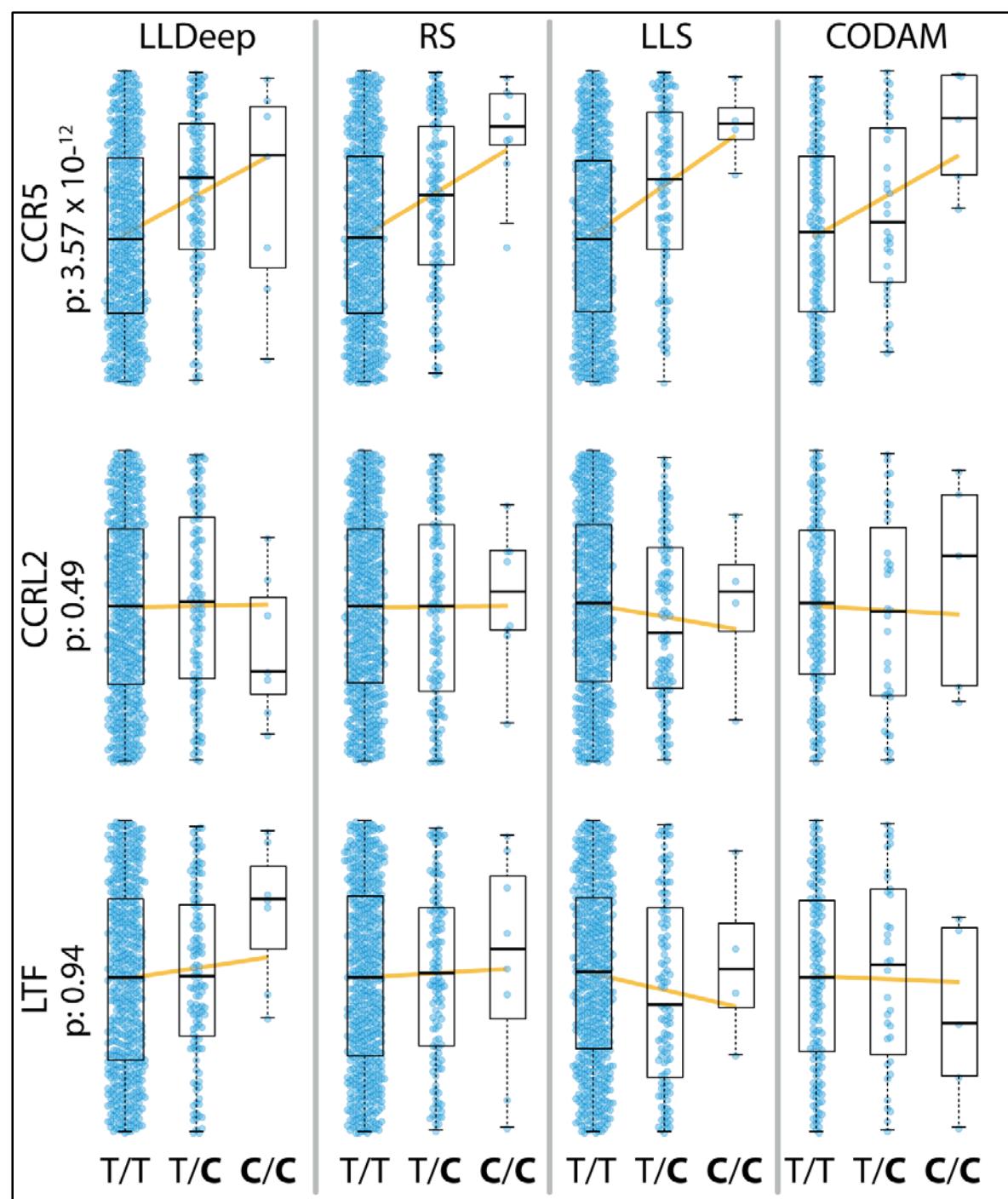
- ▶ Normally we say p-value 0.05 is significant
 - ▶ 5% ($1/20$) chance that a result is a false positive.
- ▶ Suppose we perform 20 test
 - ▶ When using p-value cut-off of 0.05, we expect 1 false positive
- ▶ Cis eQTL mapping: 20,000 genes \times 10 SNPs = 200,000 tests

False discovery rate (FDR)

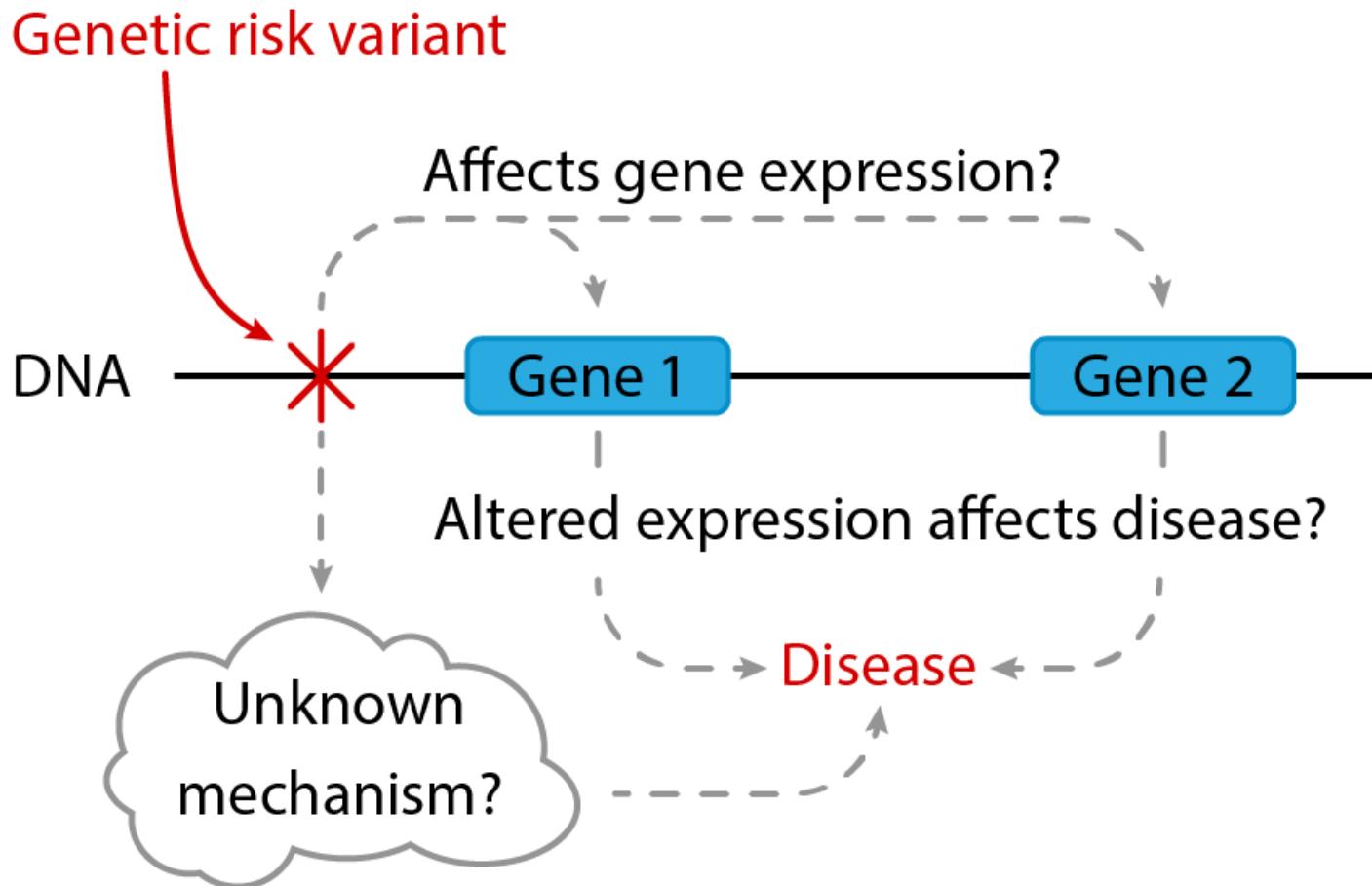
- ▶ One of the solutions when doing a lot of tests
- ▶ Expected fraction of false positives results
- ▶ Determine p-value cut-off that would result in a specific result

Example UC locus

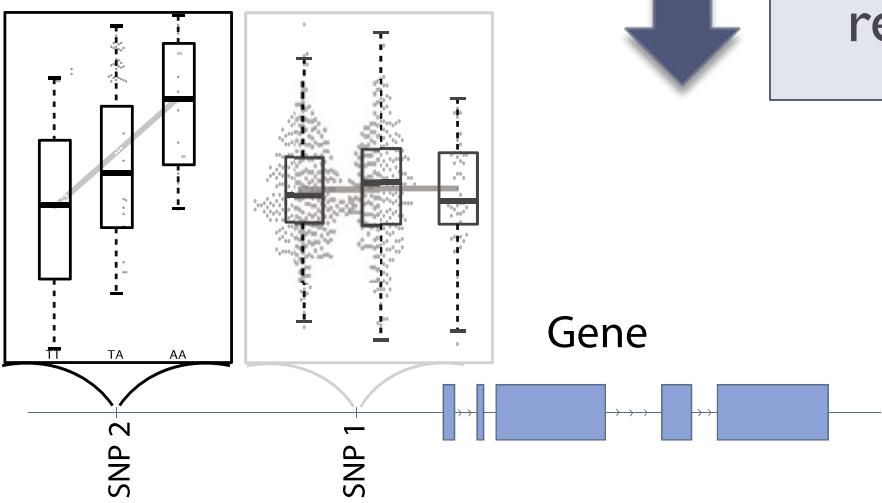
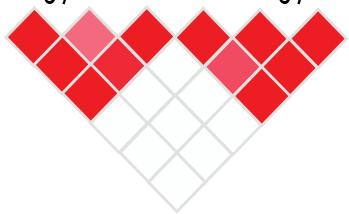
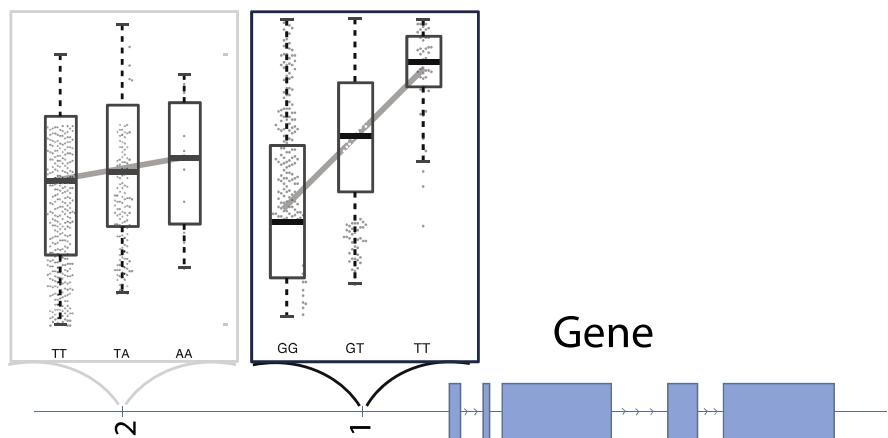




Pleiotropy

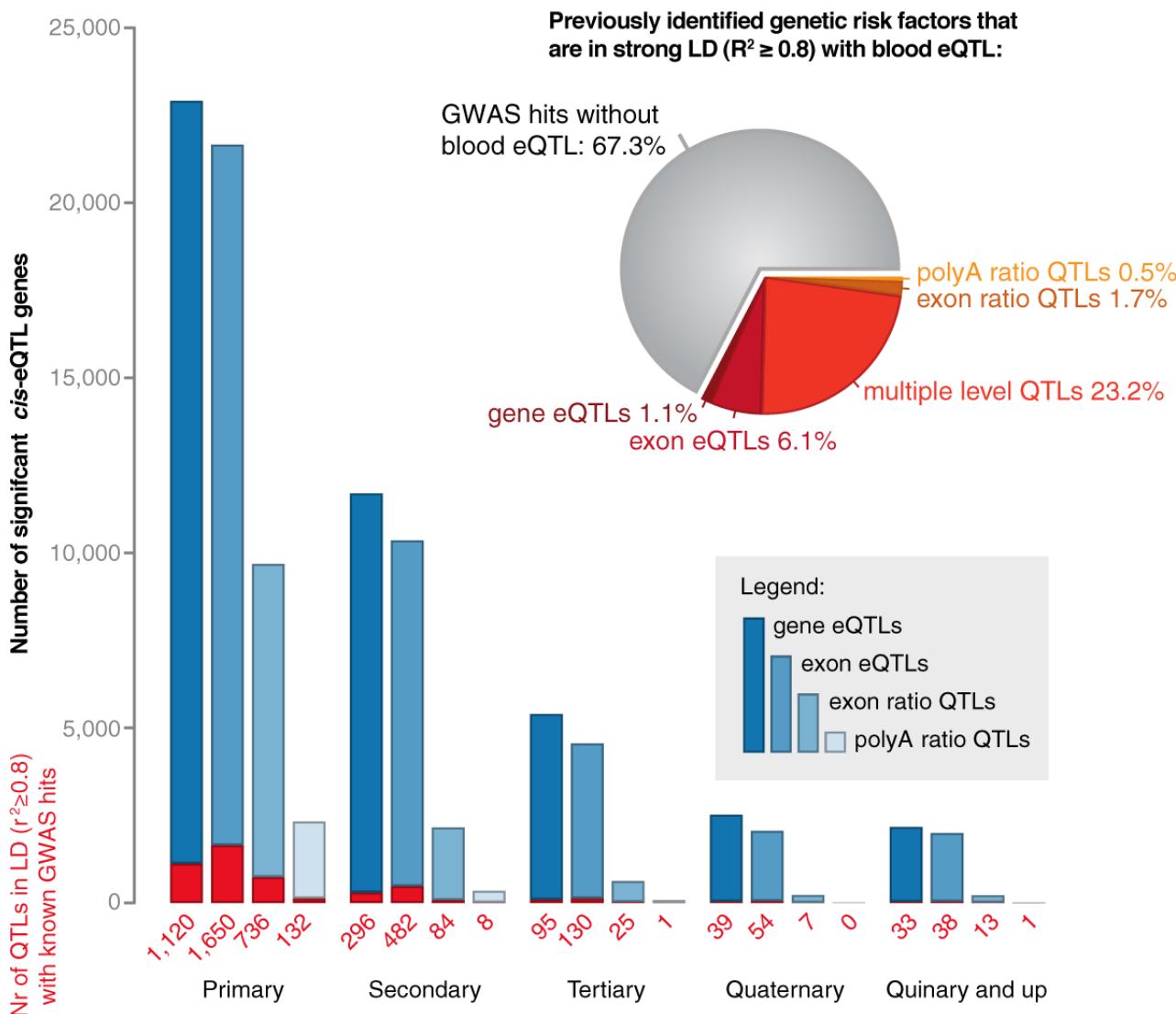


Detecting independent eQTLs

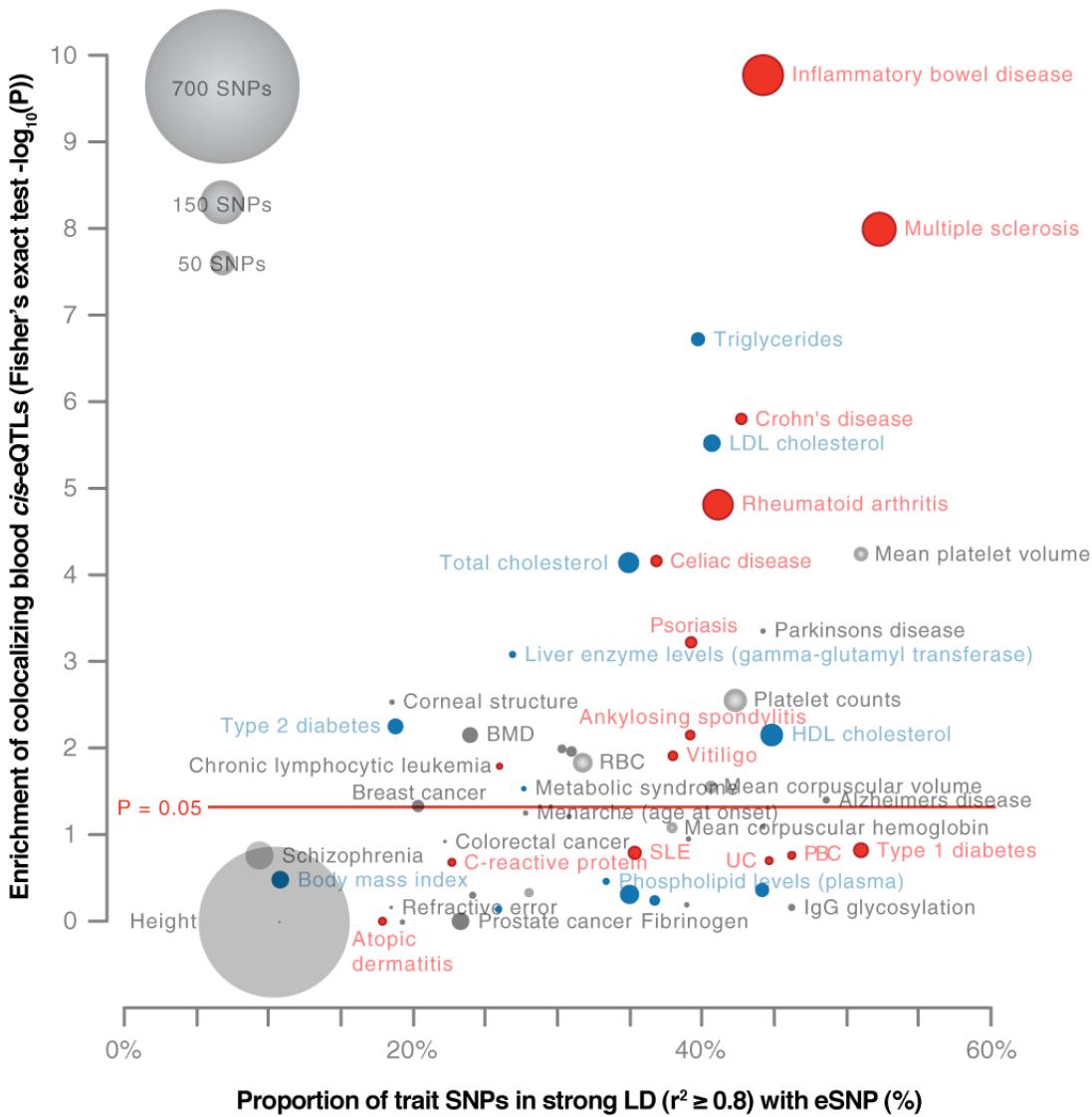


Regress out the effect of the top SNP
Run secondary eQTL mapping on residual expression

Many genetic risk factors affect gene expression



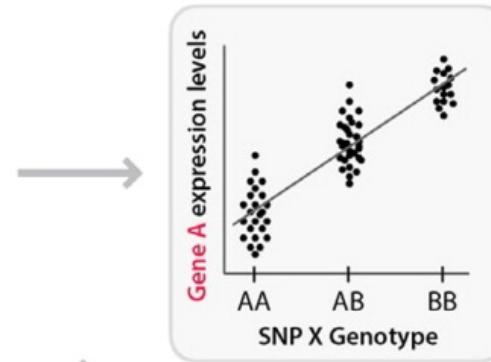
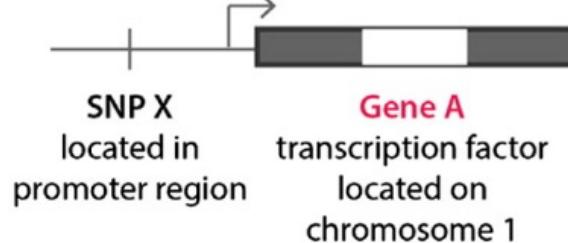
Enrichments of eQTL



Local vs distal effects

Cis-eQTL

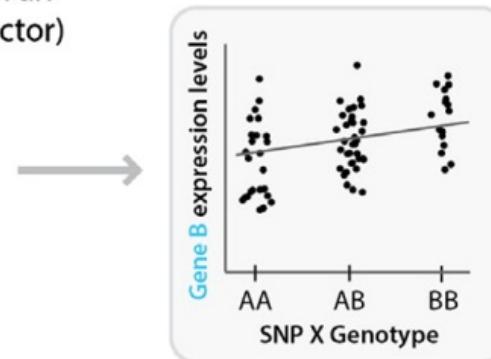
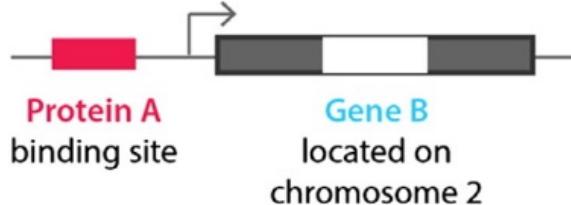
SNP X has an effect on local Gene A



Altered Protein A levels,
effect on the binding to
the transcription factor
binding sites of
downstream genes

Trans-eQTL

SNP X has an effect on distant Gene B through an intermediary factor (such as a transcription factor)



Goal

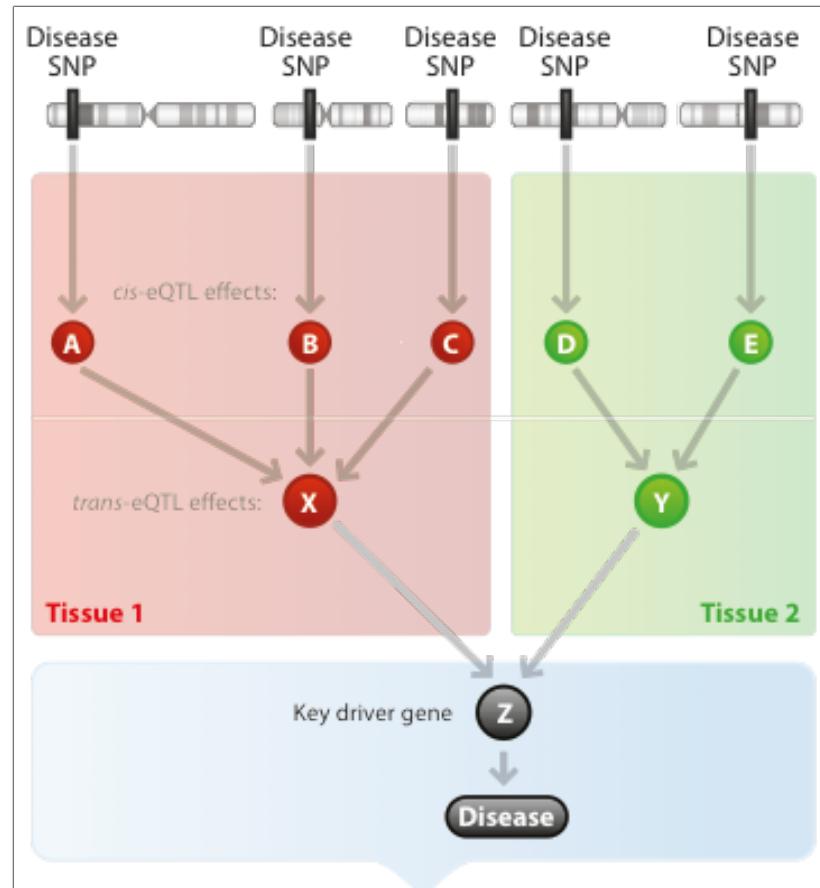
**Big-data
techniques**

Genome-wide
association studies

cis-eQTL mapping

trans-eQTL mapping

Key driver gene
identification



Systemic lupus erythematosis risk factor:



Chr. 7

Local expression effect:



Chr. 7

Type 1 interferon response:

(in Monocytes)

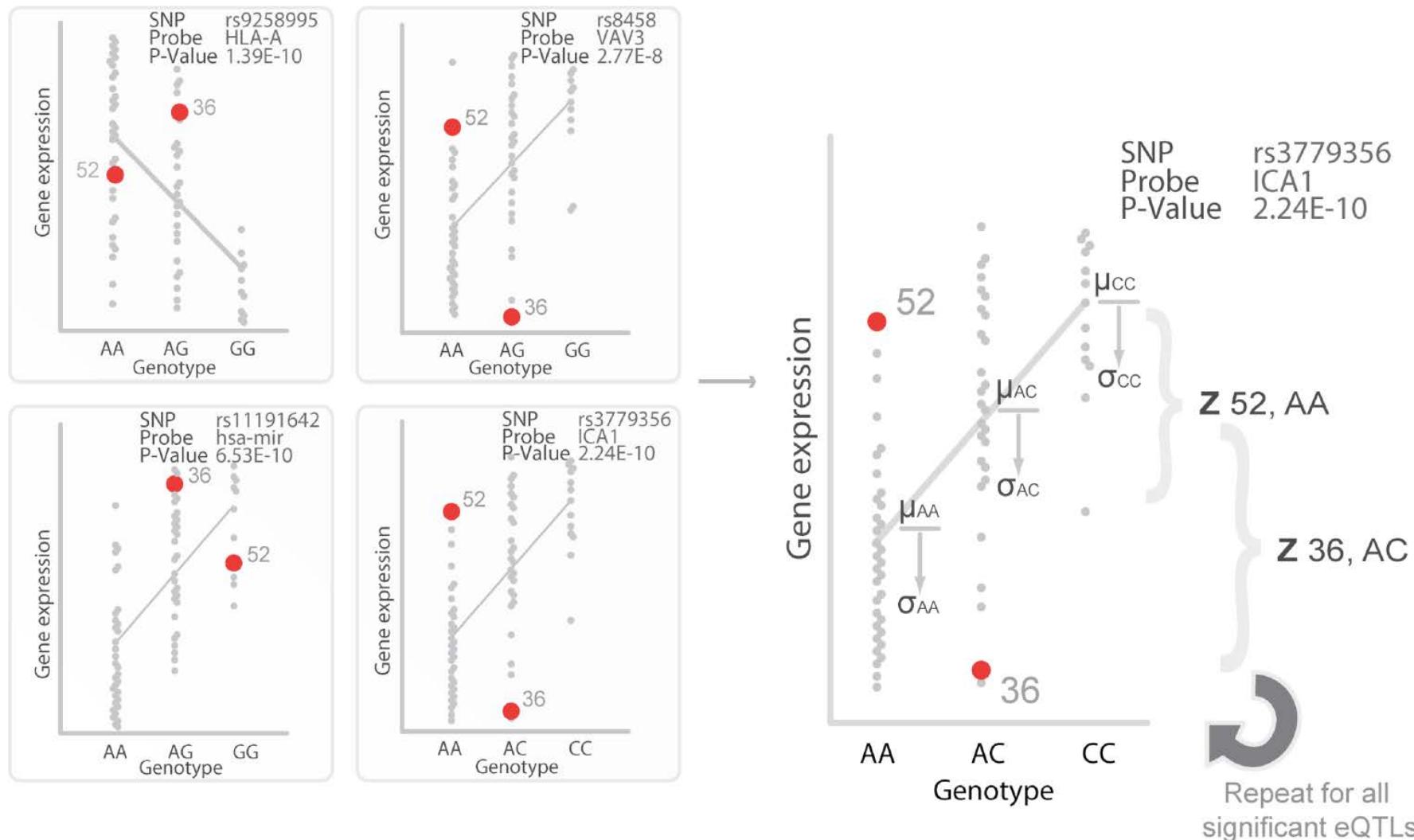


Downstream
trans-eQTL
effects

Downstream effects identified for
>200 genetic risk factors
New meta-analysis ongoing
in 25,000 blood samples ongoing

Sample mix-ups: how to identify them

Lude Franke



What happened to our data

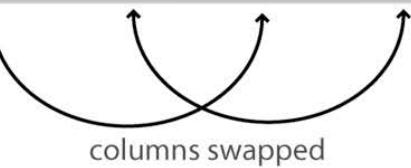
Lude Franke

Assumed plate layout

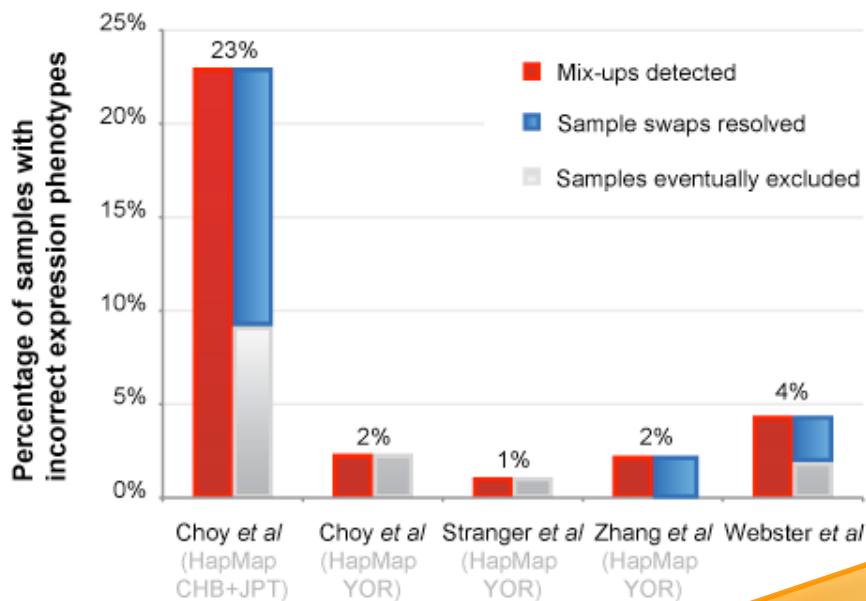
	1	2	3	4	5	6
A	65	101	70	106	68	103
B	54	108	63	112	58	110
C	42	115	52	41	47	37
D	113	45	40	53	36	48
E	107	55	111	64	109	62
F	100	66	104	71	102	69

Actual plate layout

	1	2	3	4	5	6
A	100	101	102	103	104	106
B	107	108	109	110	111	112
C	113	115	115	110	111	112
D	42	45	45	48	40	41
E	54	55	55	48	52	53
F	65	66	66	69	70	71

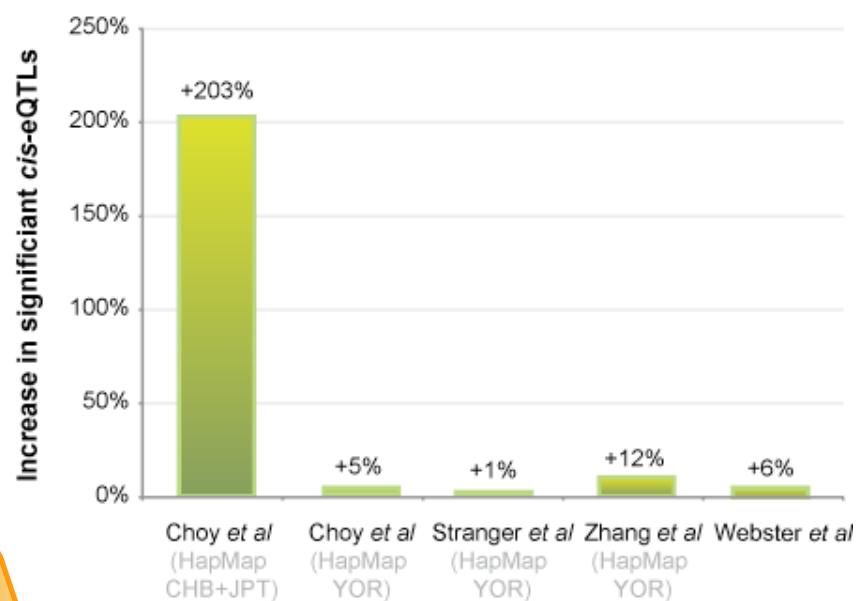


eQTL datasets with mix-ups

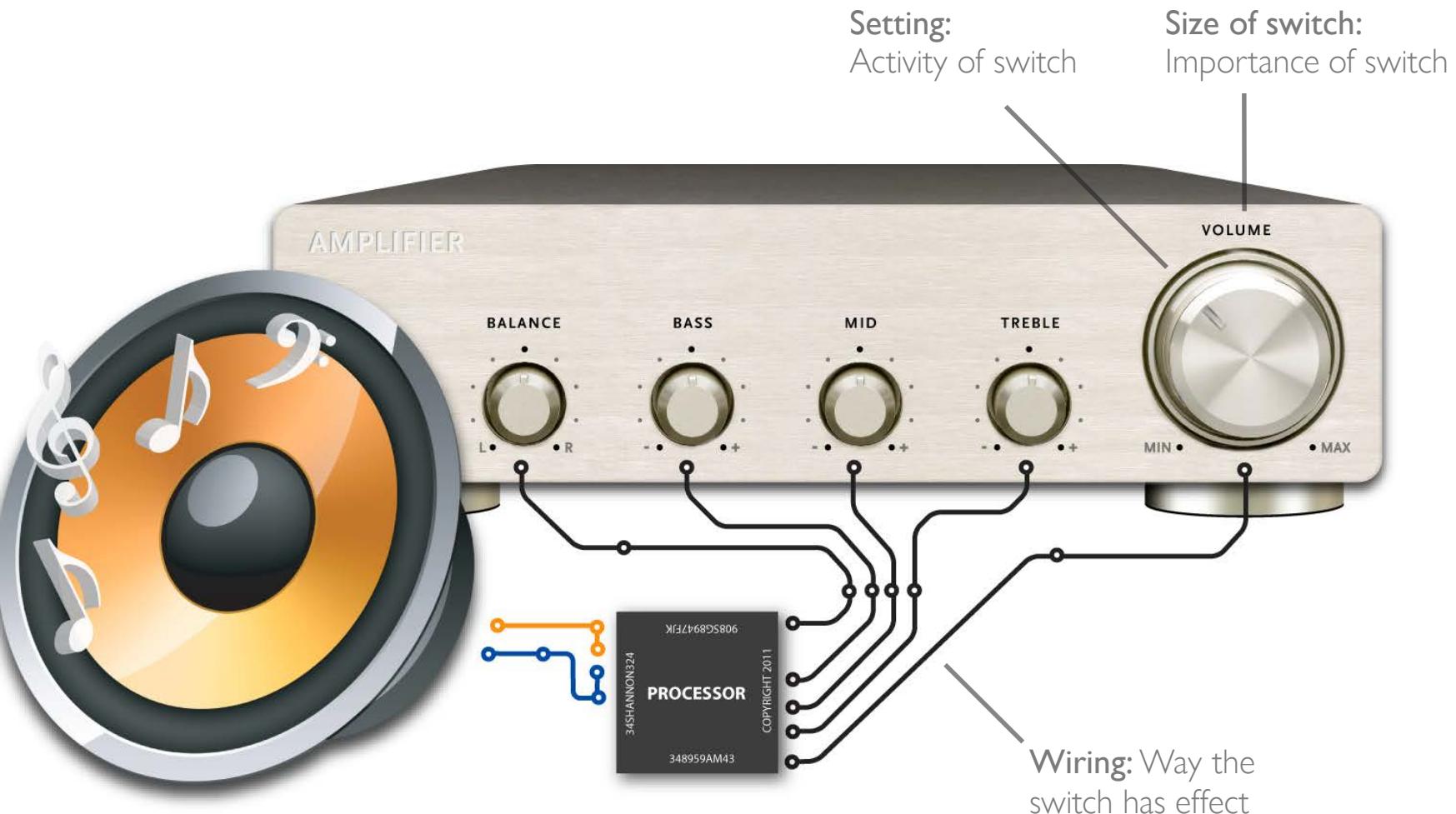


On average 3% of eQTL samples are mixed-up!

Effect of correcting for these mix-ups



Amplifier can change many aspects of music



A control panel that determines gene expression?

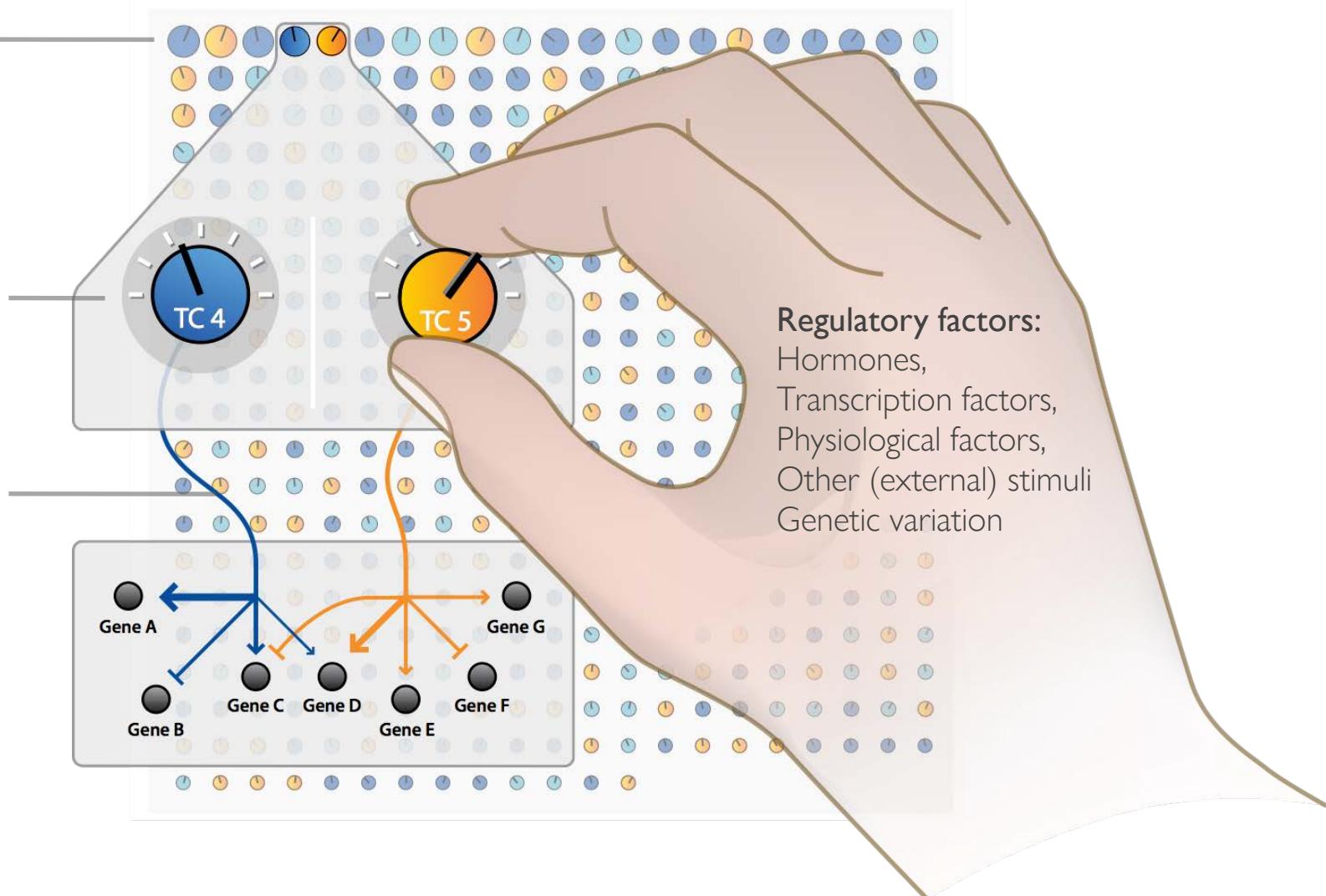
Size of switch:

Importance

Setting: State of
a certain sample

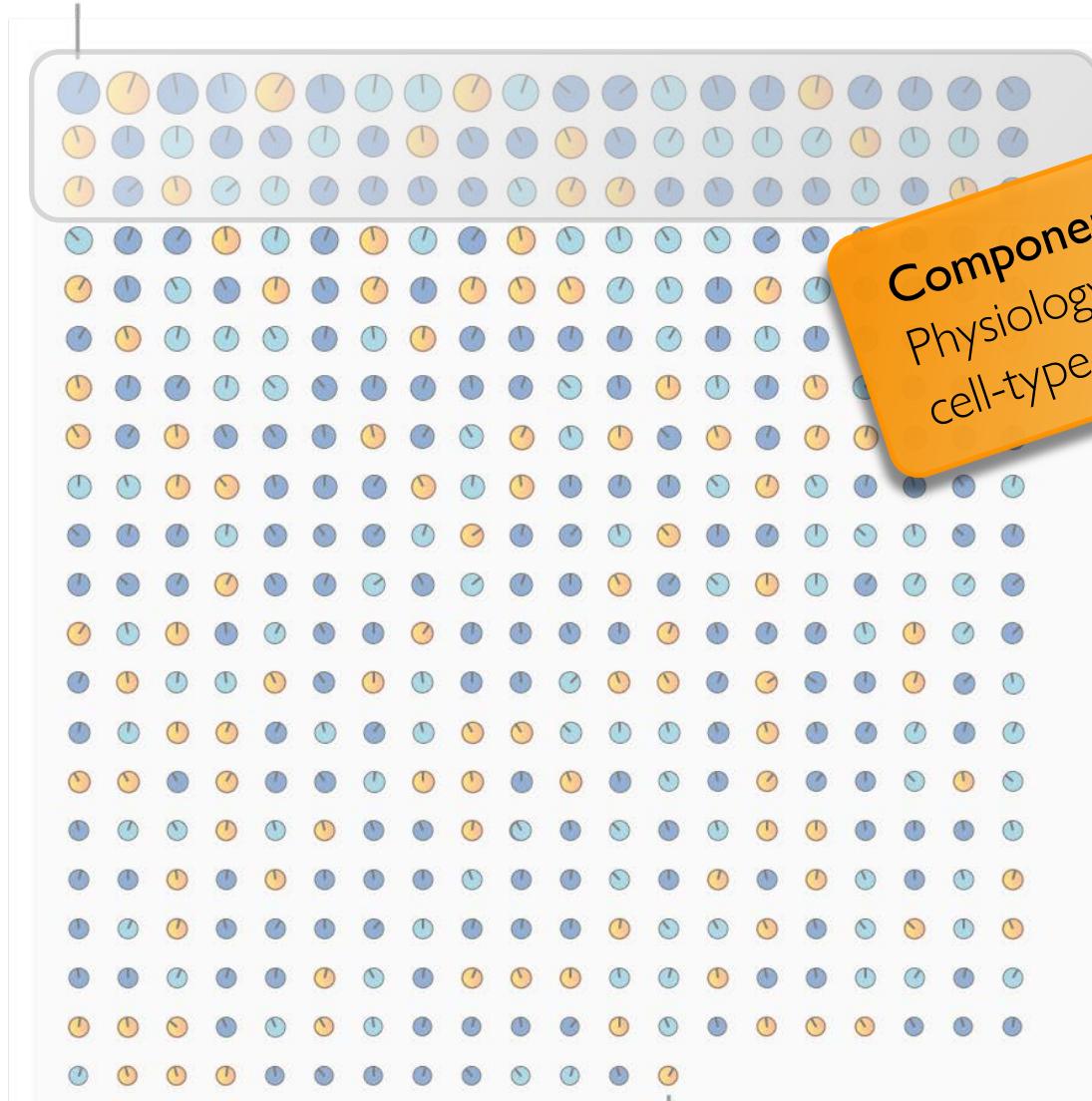
Wiring: Effect on
individual genes

Regulatory factors:
Hormones,
Transcription factors,
Physiological factors,
Other (external) stimuli
Genetic variation



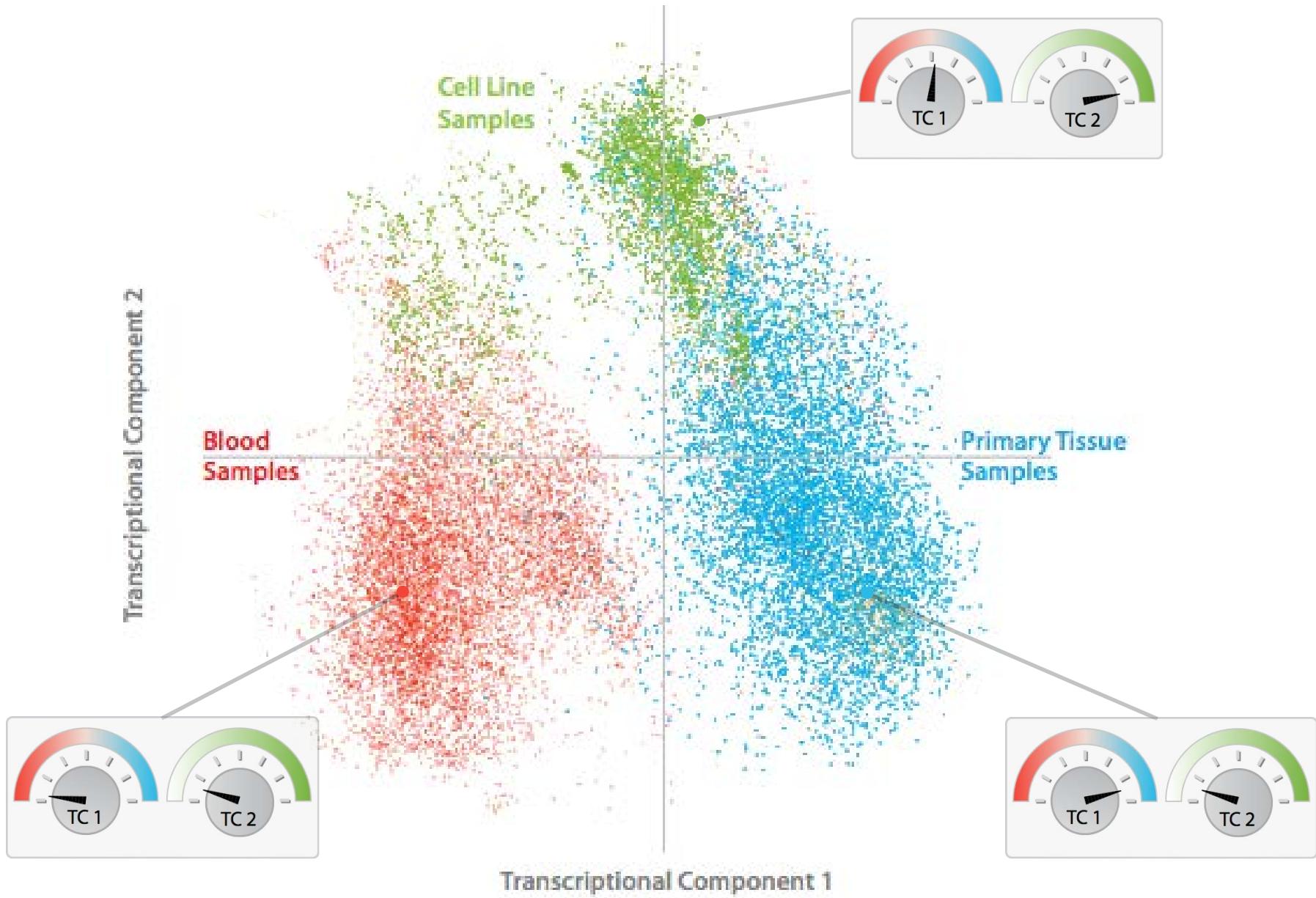
800 ‘transcriptional components’: Component 1 - 50

Component 1

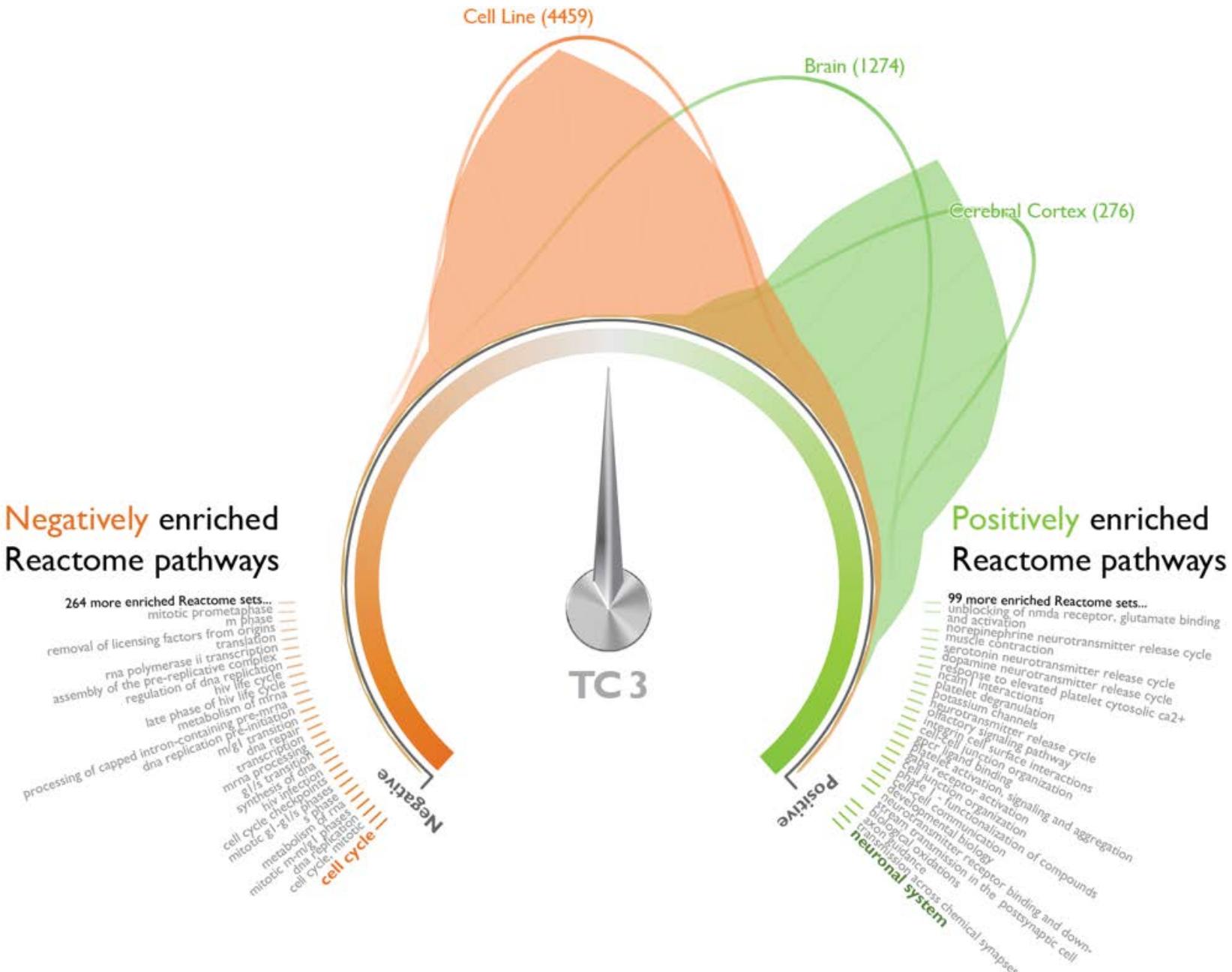


Components 1 - 50:
Physiology, metabolism,
cell-type differences

Component 1 and 2



Transcriptional component 3



Detection cytogenetic aberration in expression data

Chromosome

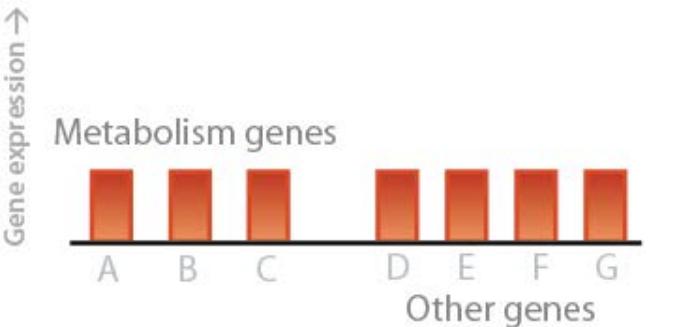
Down Syndrome patient: dup 21



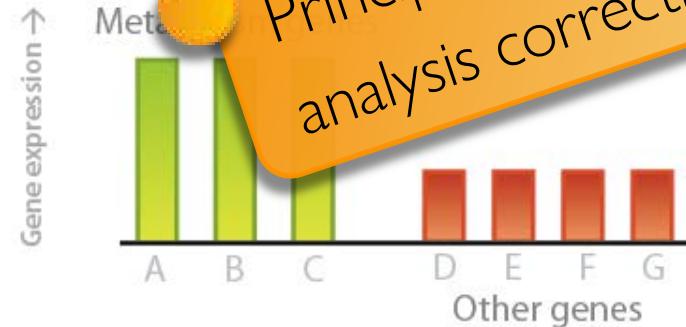
Large proportion of expression variation is determined by genetic variation but due to e.g.:

- Physiological state of samples
- Environmental state of samples (e.g. fasting vs. non-fasting)

RNA blood expression
when you wake up

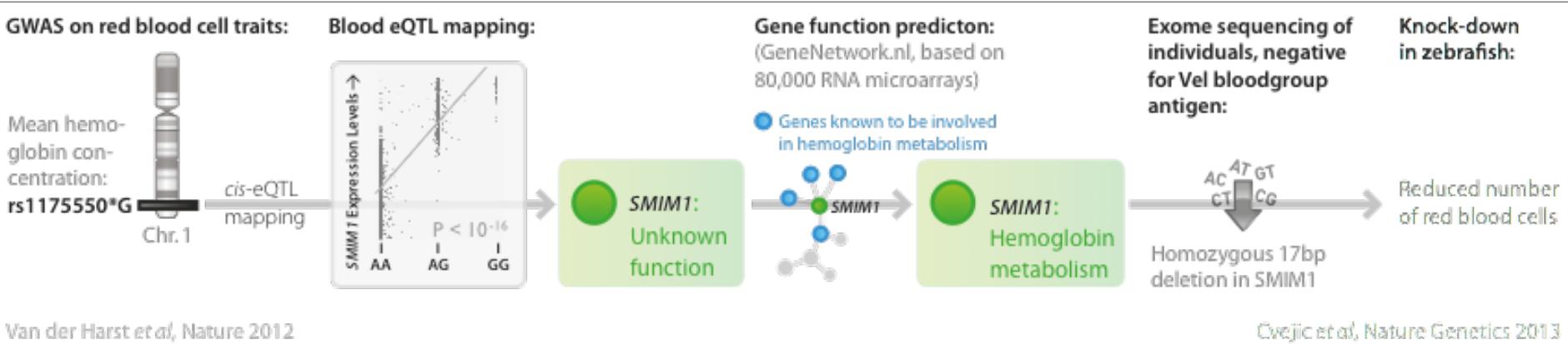


RNA blood expression after dinner



Get rid of this ‘noise’?
Principal component
analysis correction

GeneNetwork gene function predictions



Amounts of data integrated:

GWAS in 135,000 samples

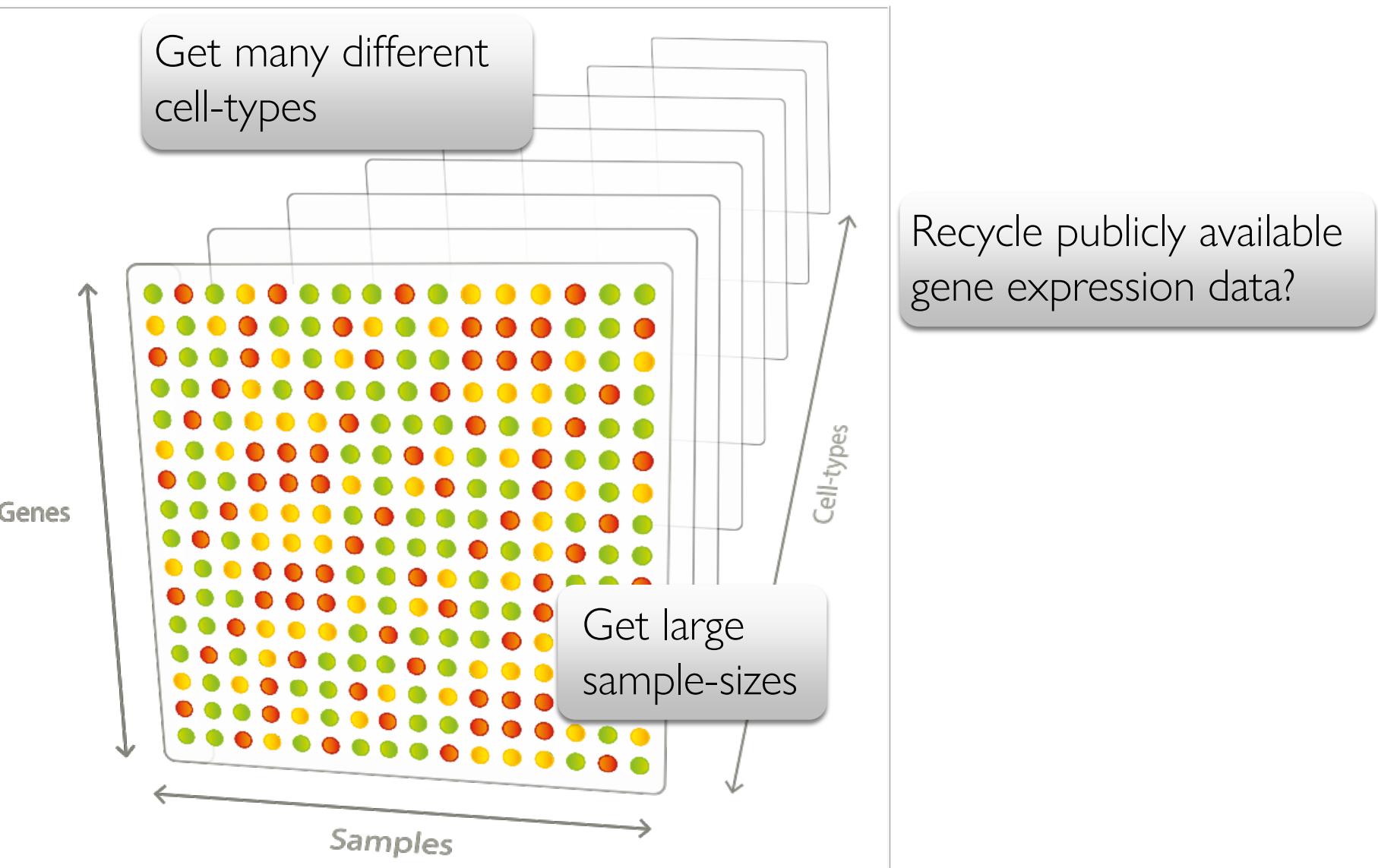
eQTL mapping in 1,500 samples

Transcriptomics in 80,000 samples

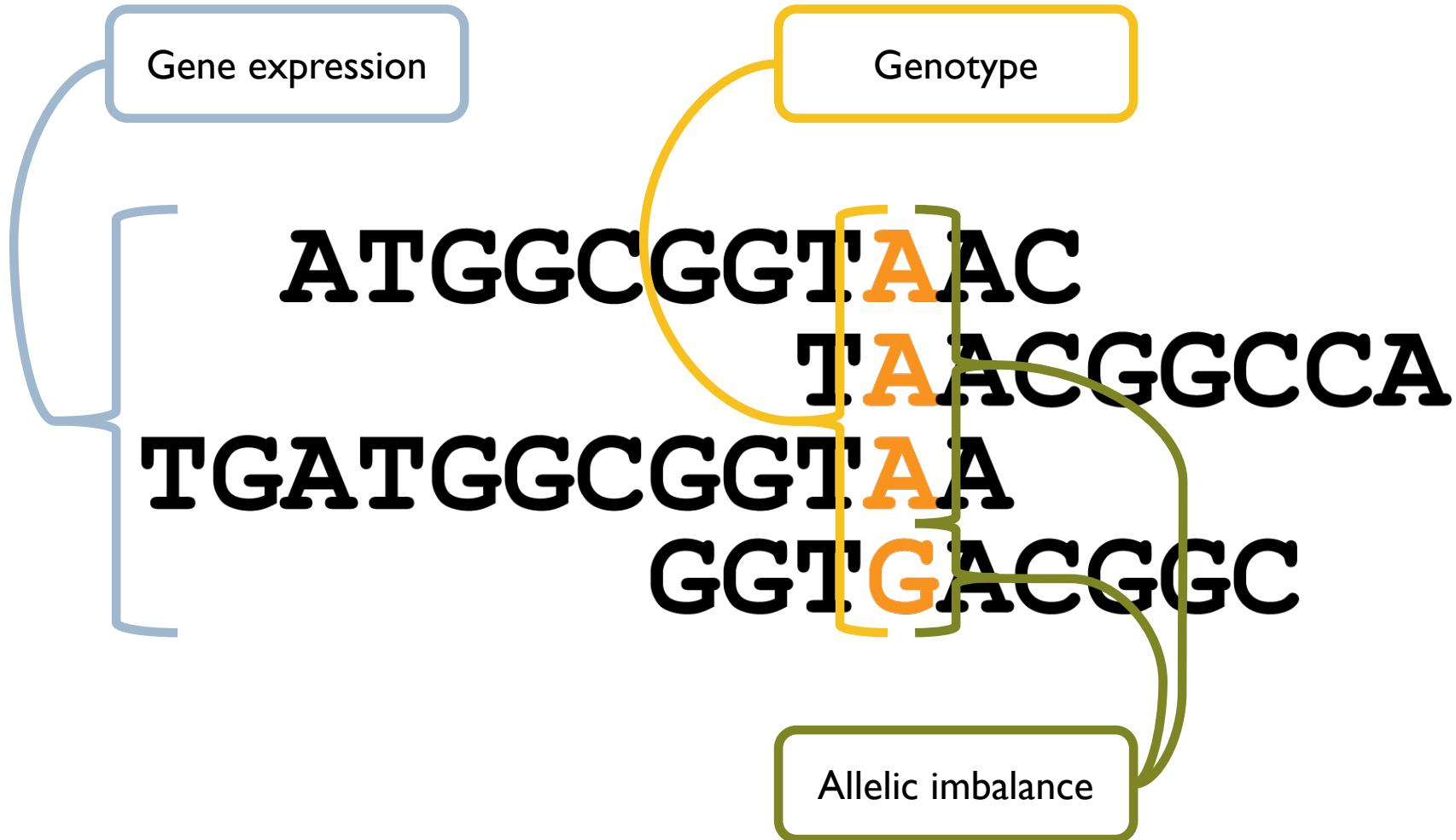
Exome sequencing

Wet lab proof

Using public data

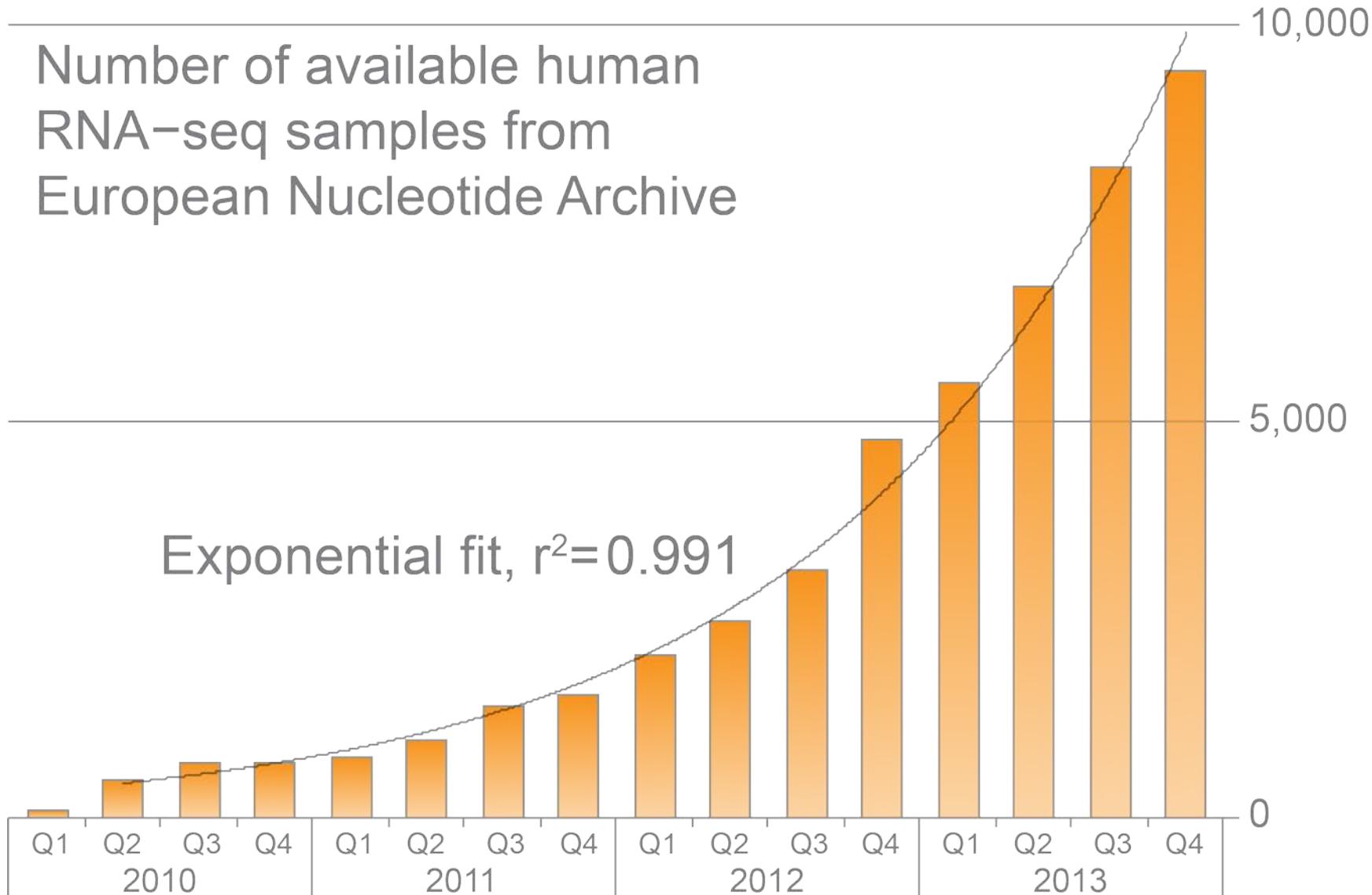


RNA sequencing

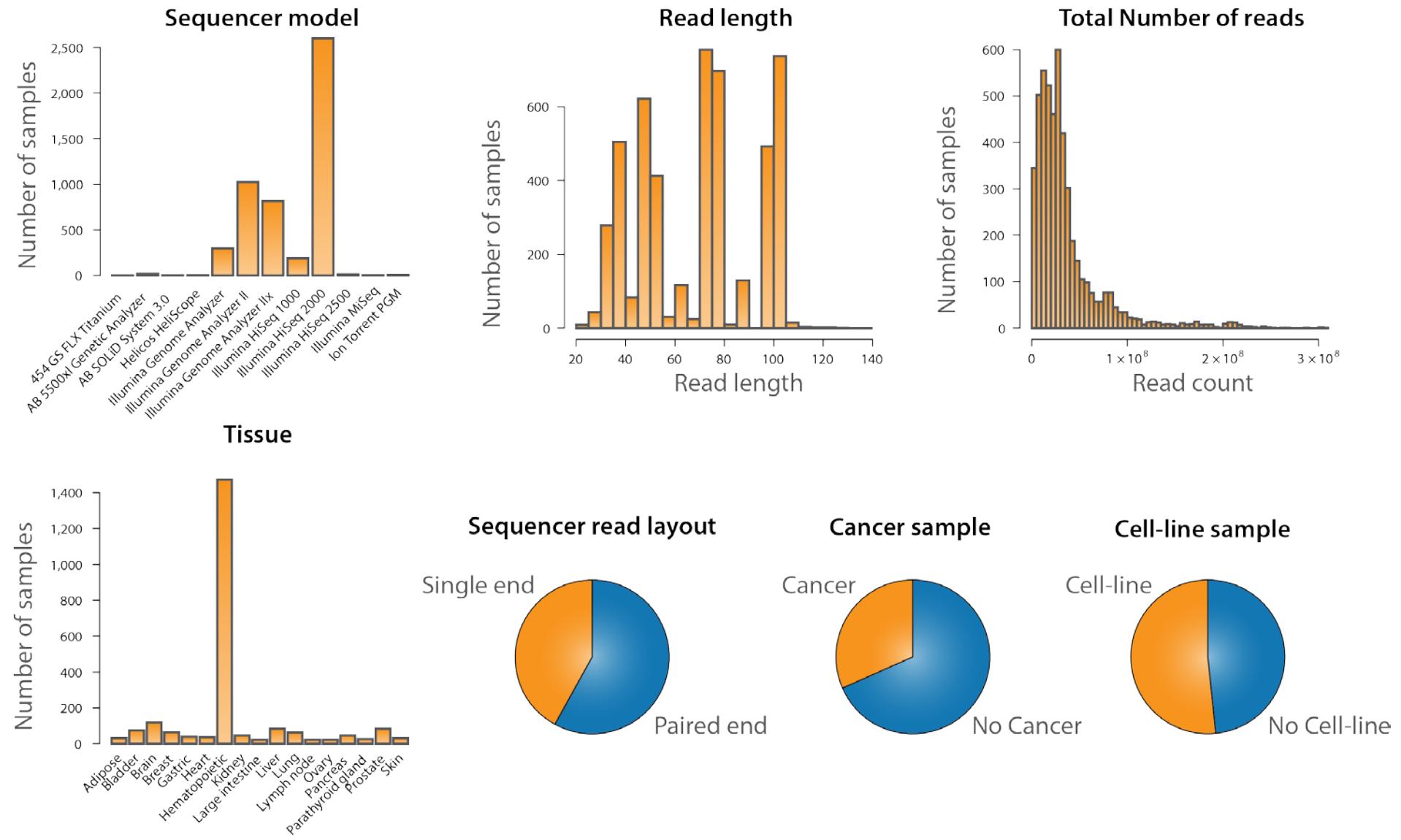


Available data is growing rapidly

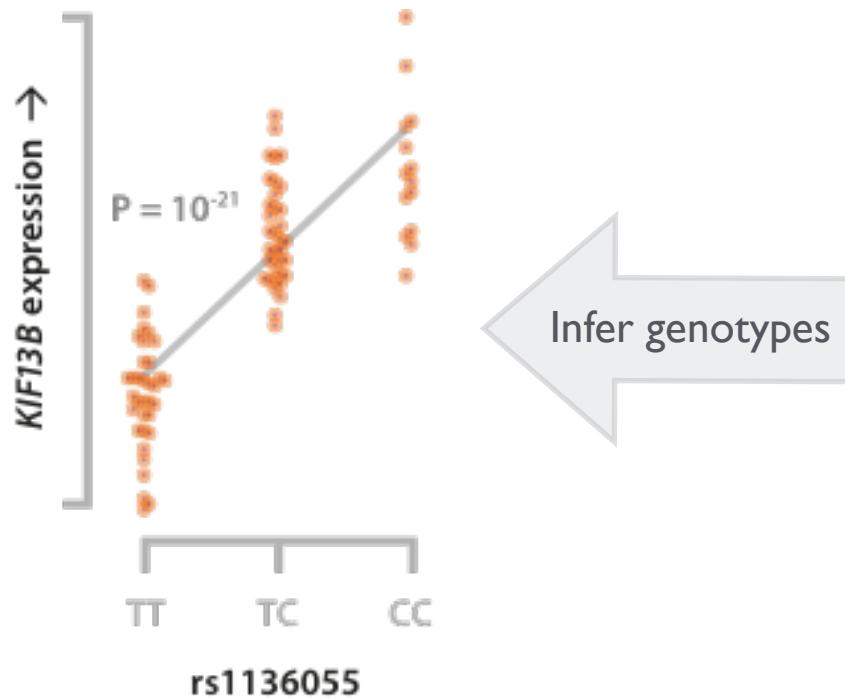
Number of available human
RNA-seq samples from
European Nucleotide Archive



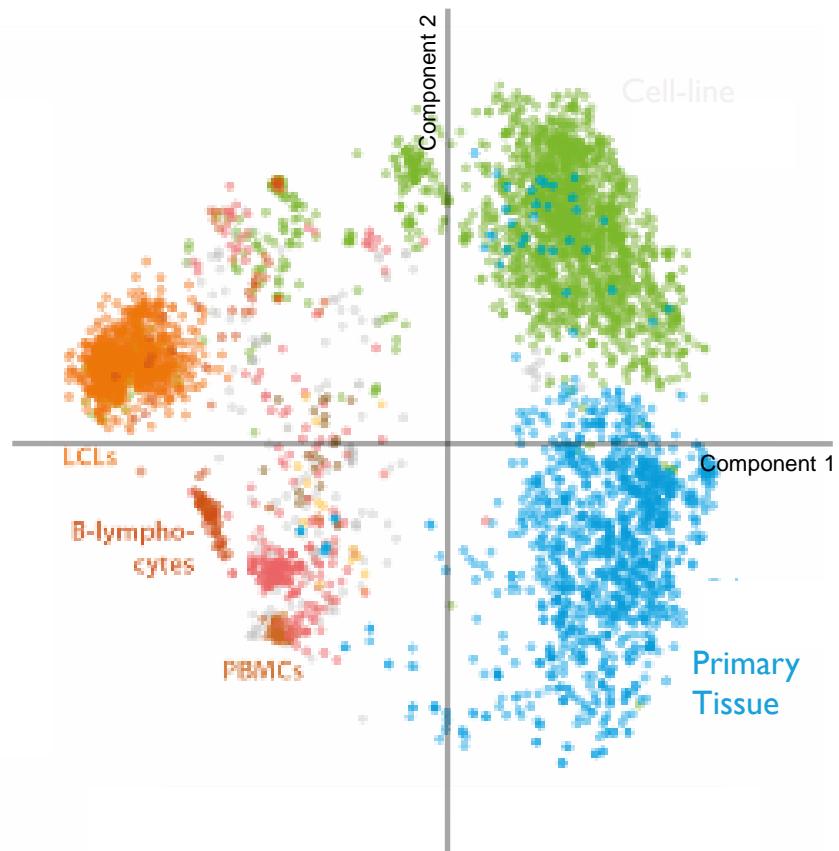
Heterogeneous dataset



Derive SNP genotypes from RNA-seq data

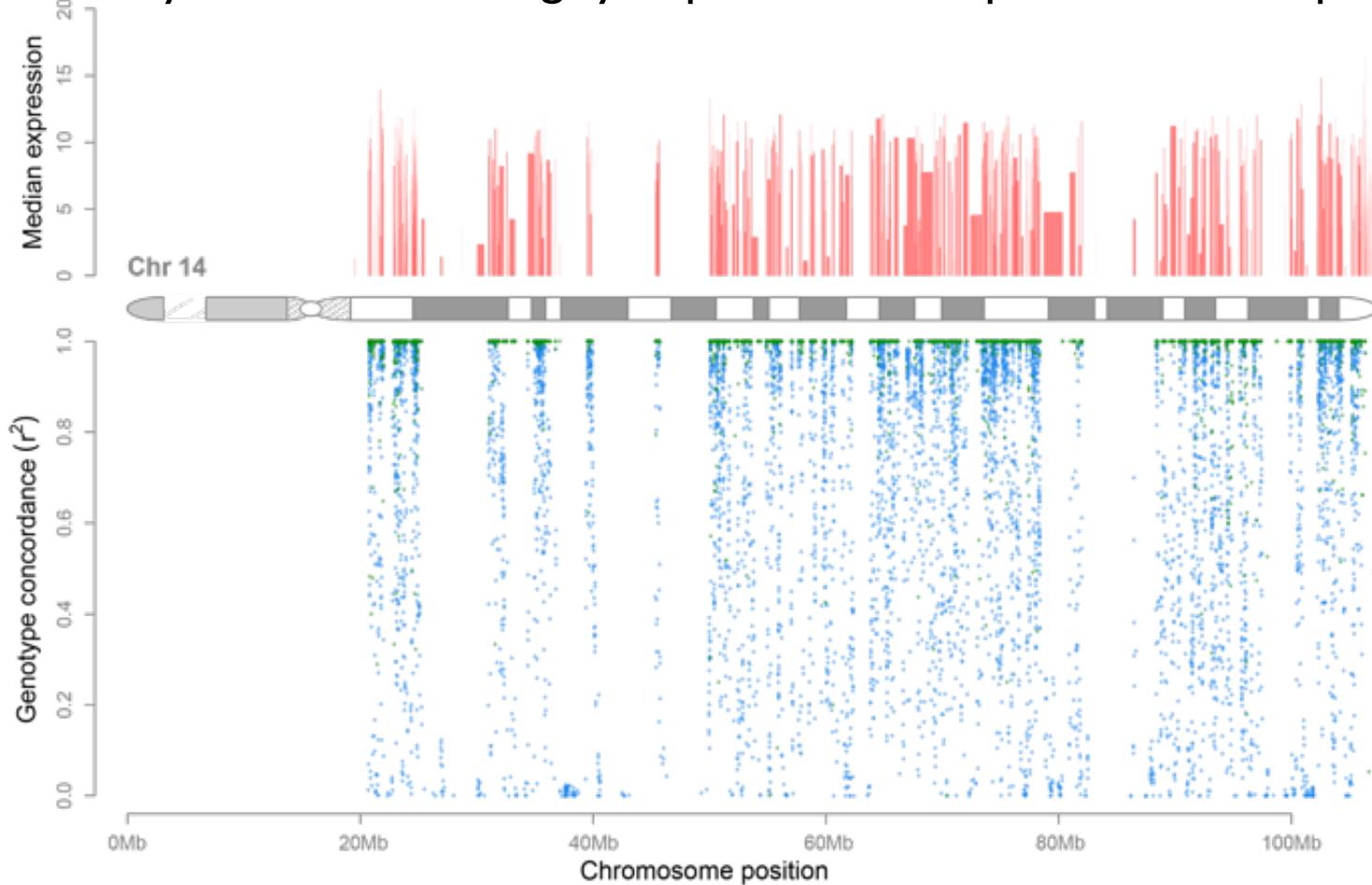


Public RNA-seq data (5,000 samples)

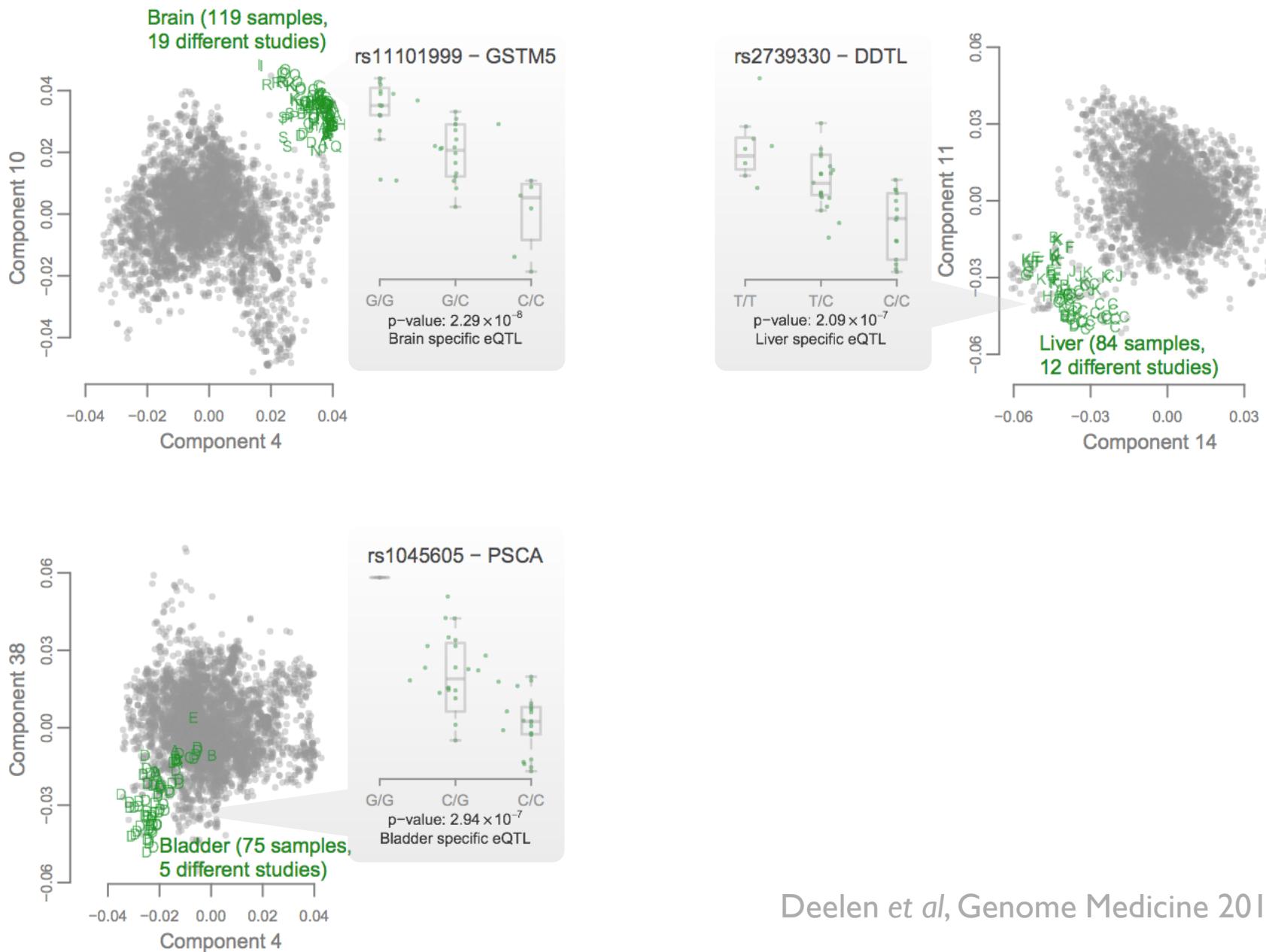


Calling genotypes in RNA-seq data

Ability to call SNP is largely dependent on expressed transcripts



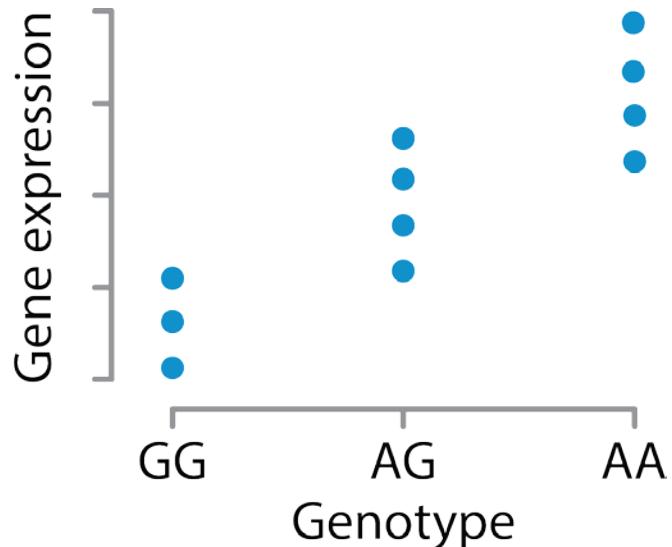
Tissue-specific eQTL mapping for free



Investigation of rare variants

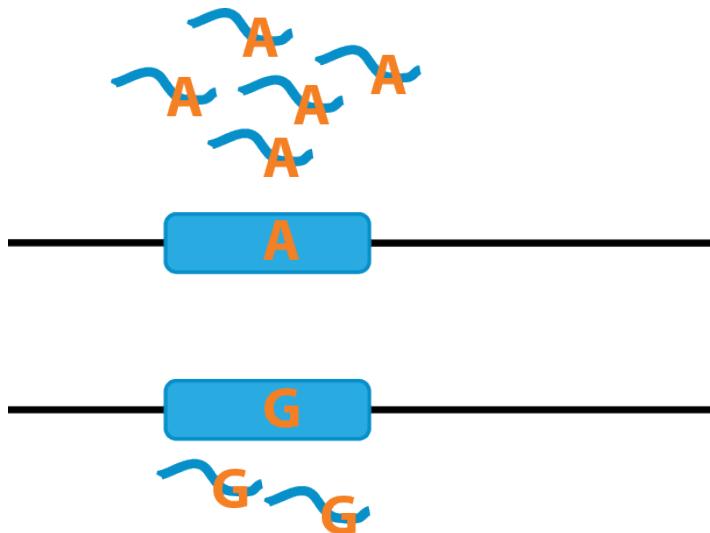
eQTL

- ▶ Correlation between genotype and expression
- ▶ Common variants



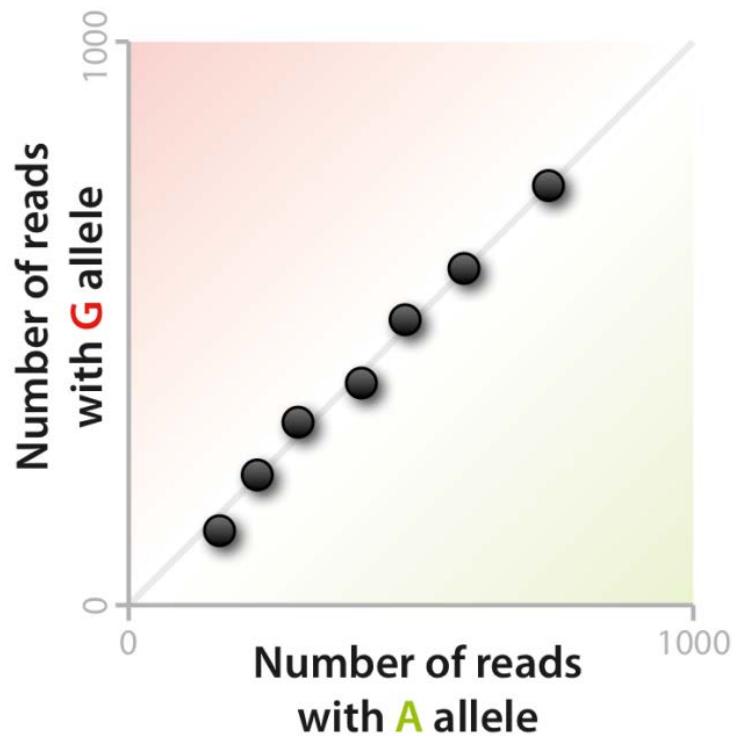
Allele Specific Expression

- ▶ Imbalance in expression
- ▶ Rare variants

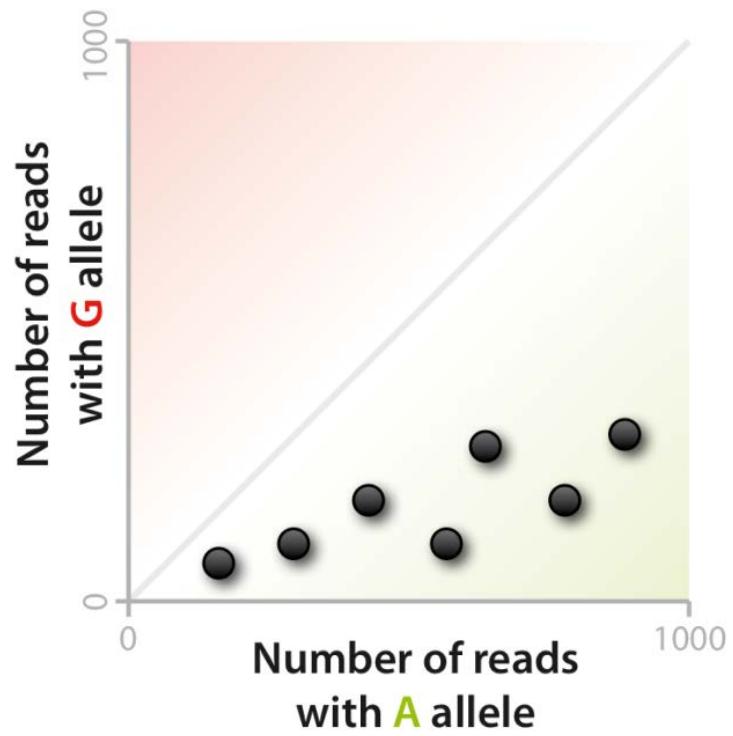


Allele Specific Expression (ASE)

Alleles equally expressed



Bias towards A allele



Why focus on rare variants?

Common
complex
diseases

Rare
monogenetic
diseases

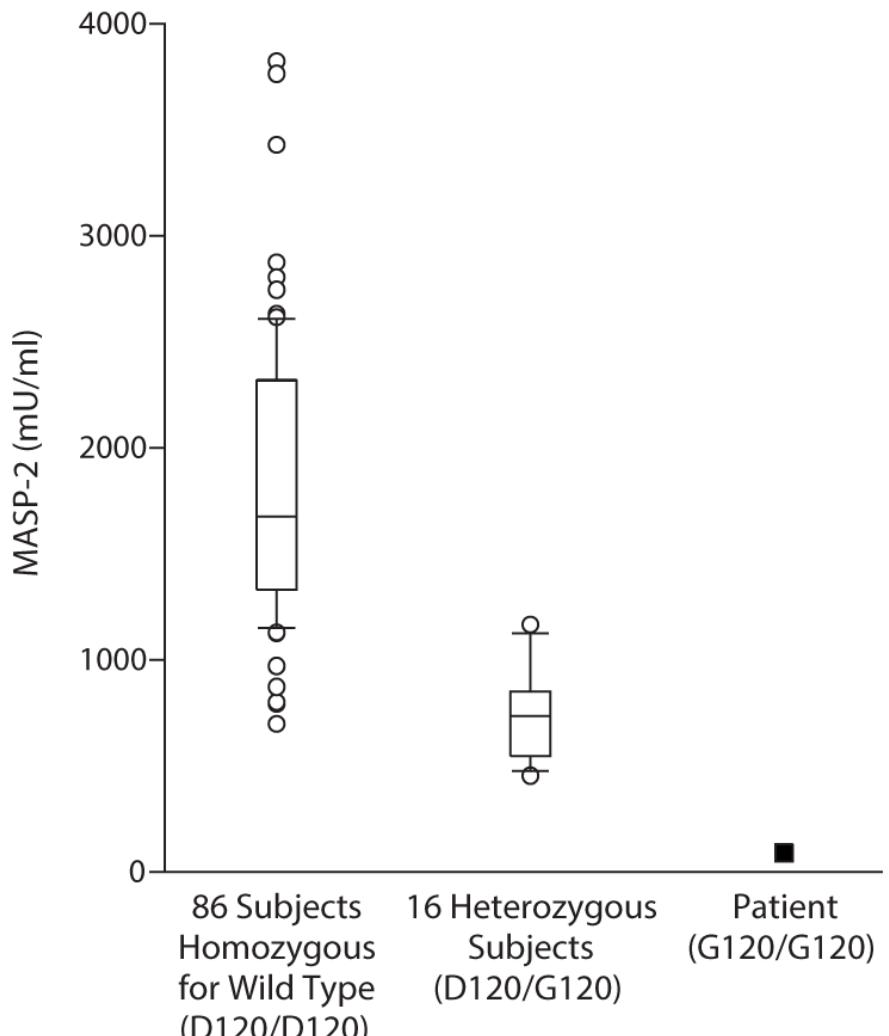
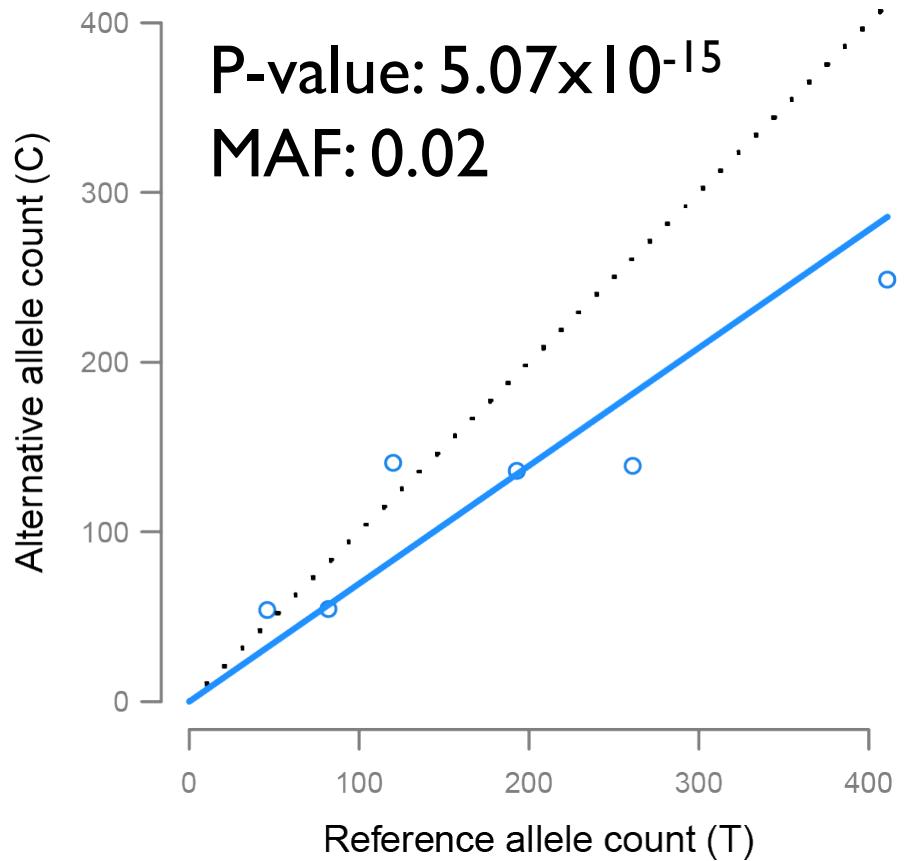
Many common
variants

Single rare
genetic variant

>50% of variants
affect gene
expression

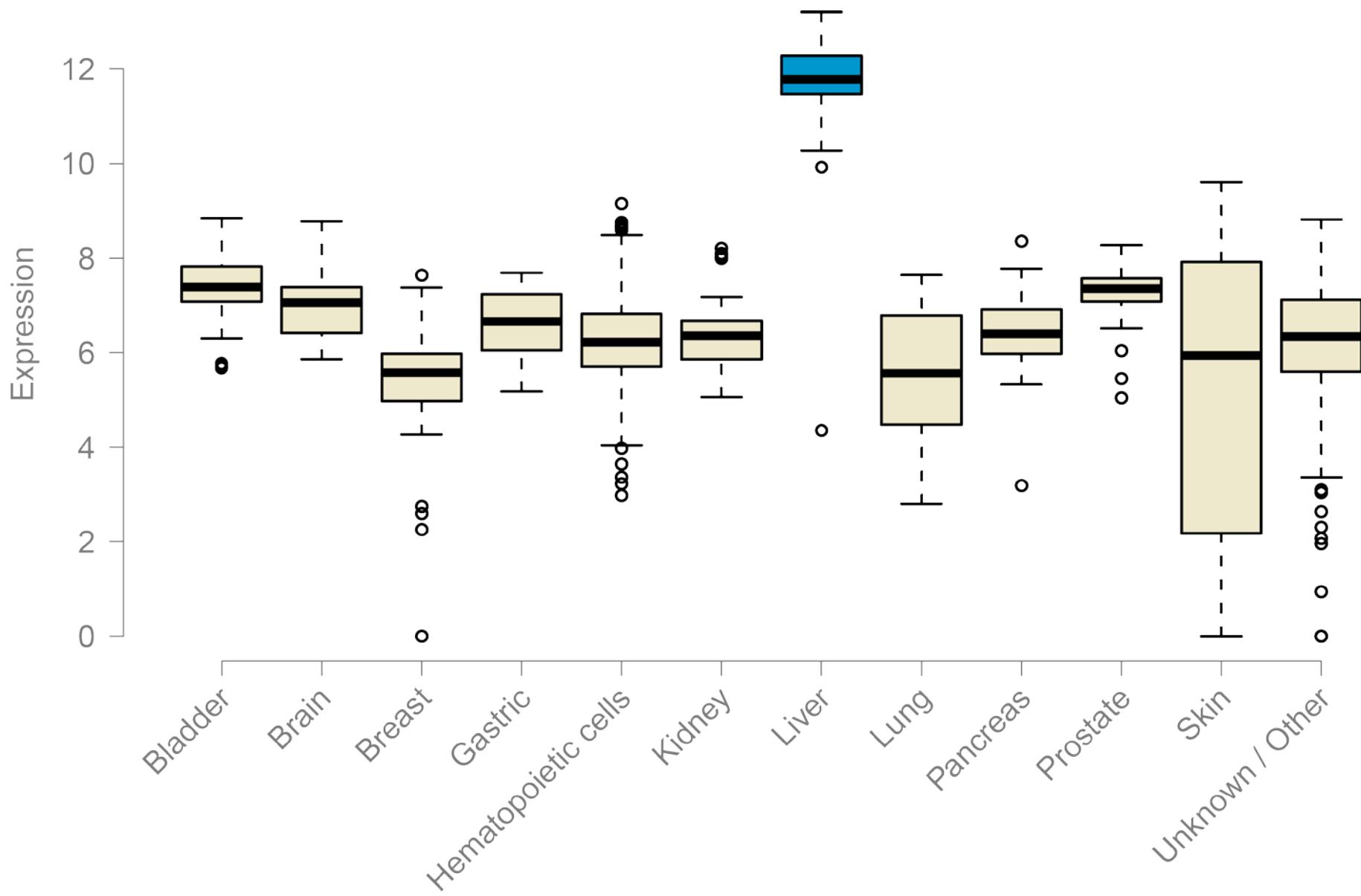
Focus on
protein coding
variants

Example: MASP2 deficiency



Stengaard-Pedersen et.al NEJM 2003

MASP2 expression levels



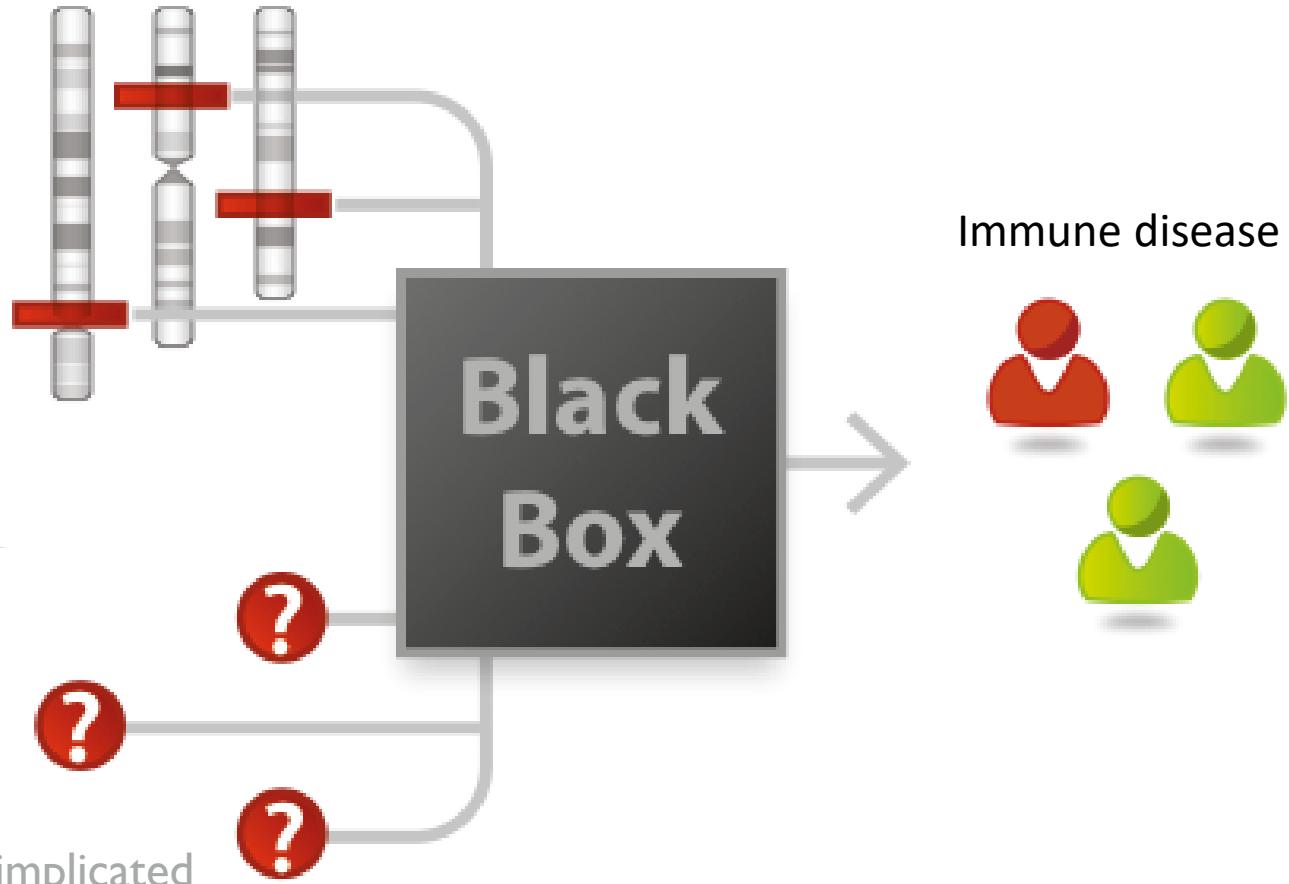
Interplay of genetic & environmental risk factors unknown

Genetic
risk factors

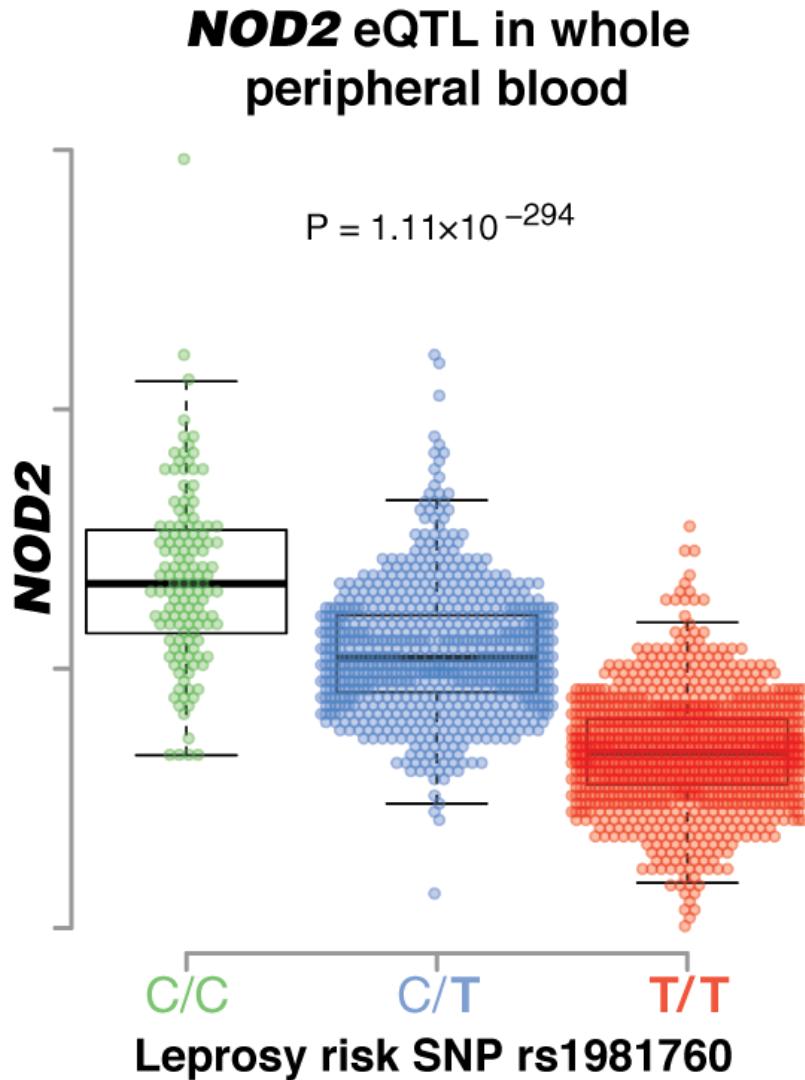
>1,000 known,
But small effects

Environmental
risk factors

Mostly unknown,
Viruses / bacteria implicated

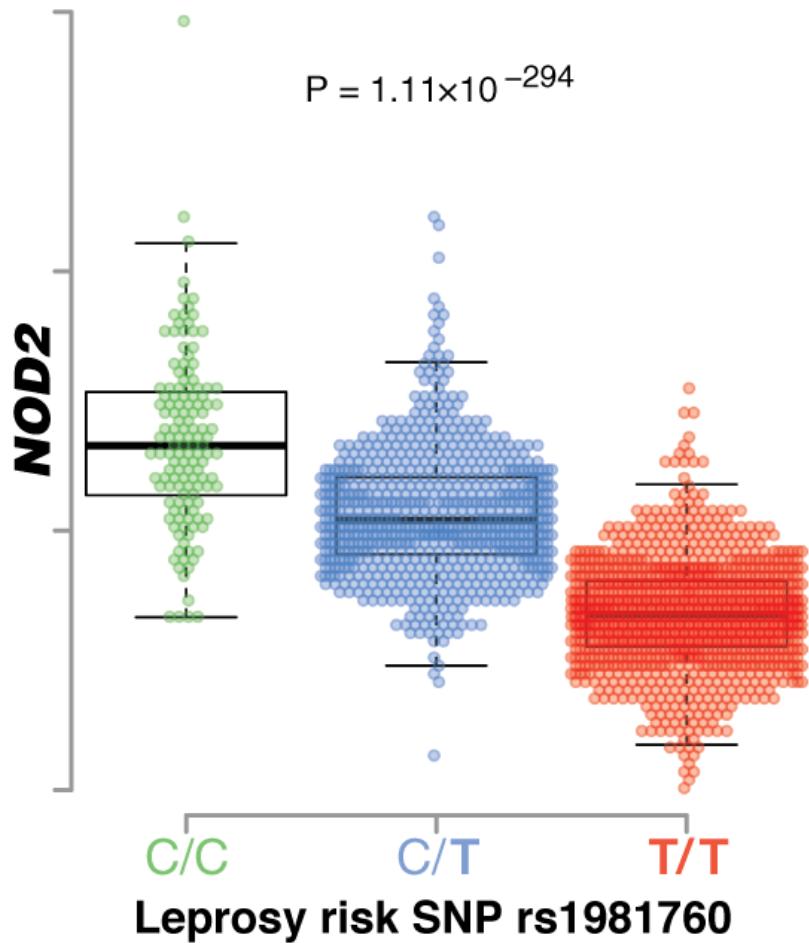


Detecting cell-type dependent eQTLs in whole blood

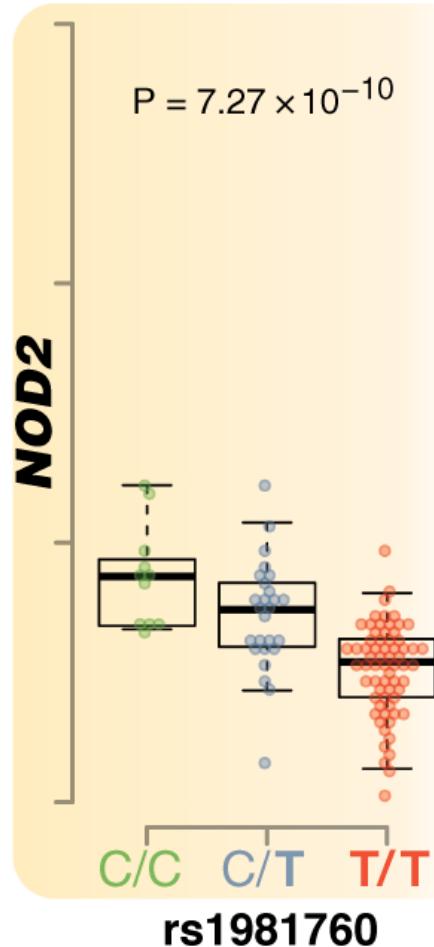


Detecting cell-type dependent eQTLs in whole blood

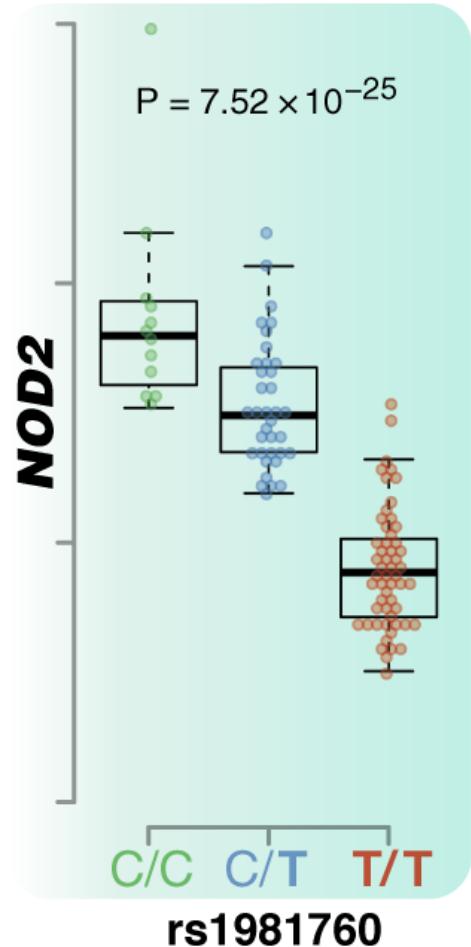
***NOD2* eQTL in whole peripheral blood**



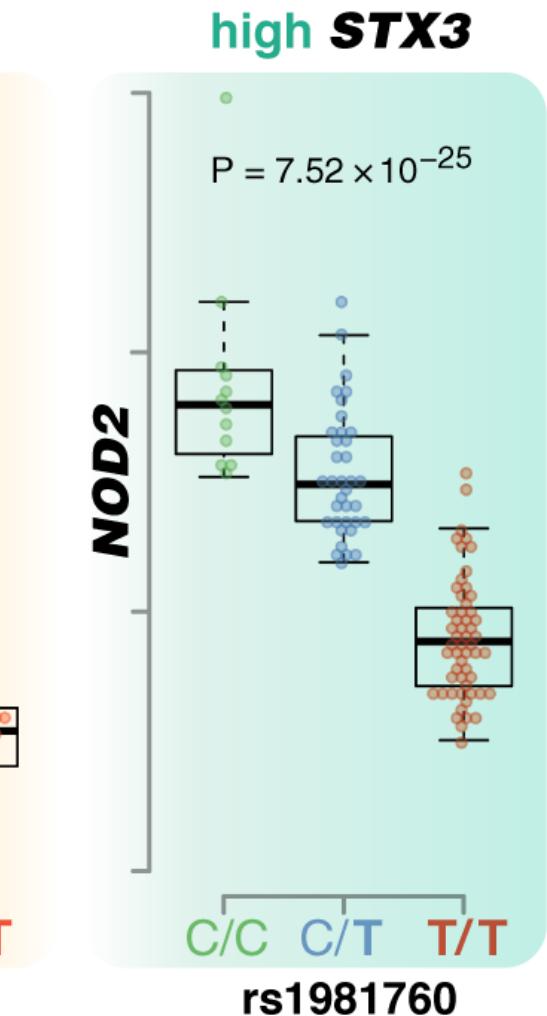
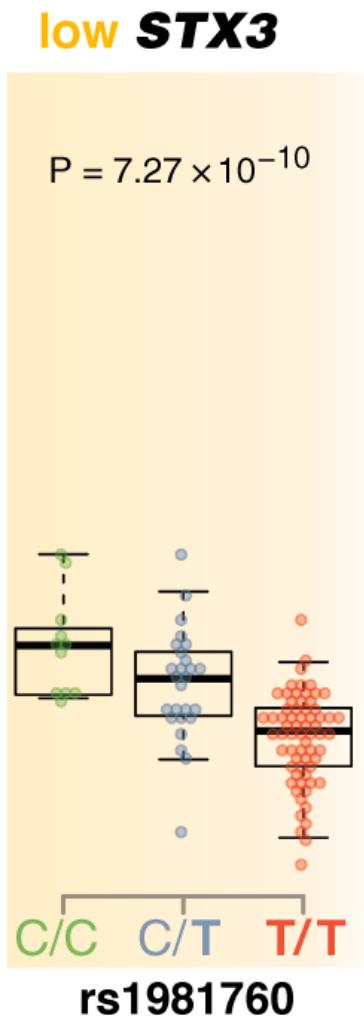
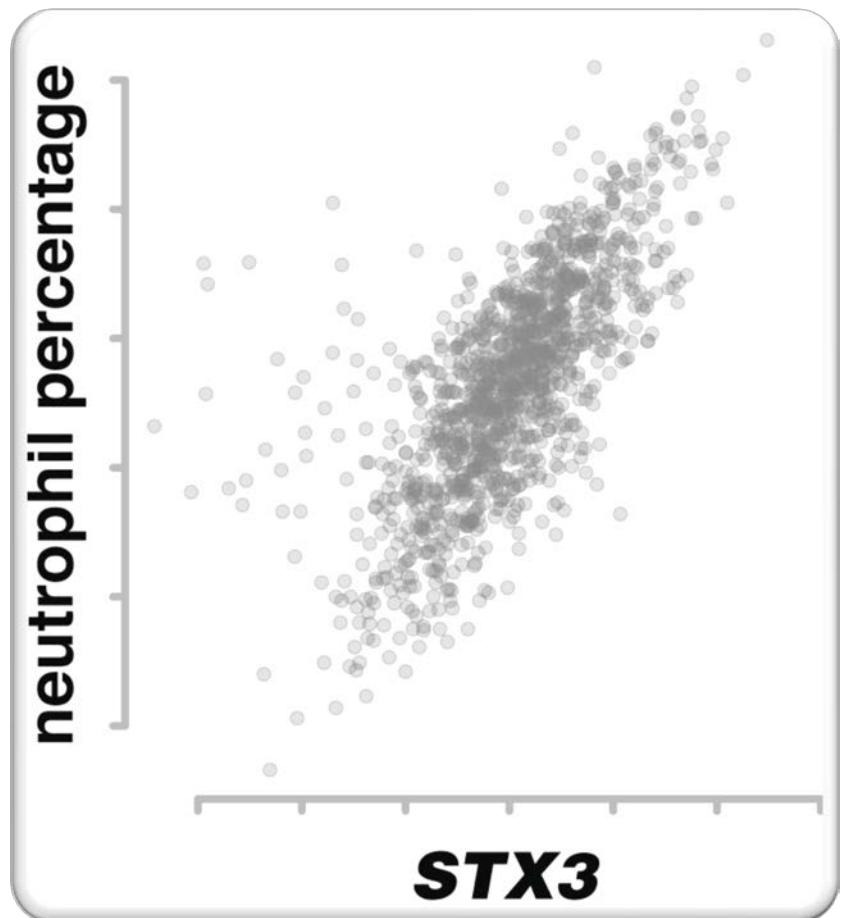
low ***STX3***



high ***STX3***



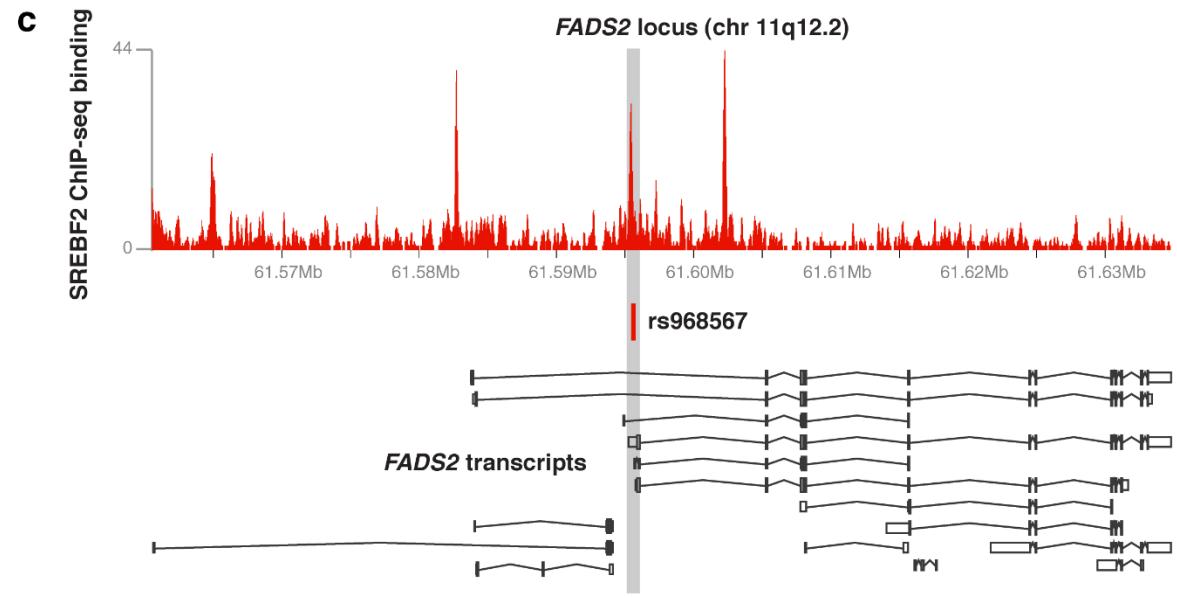
Detecting cell-type dependent eQTLs in whole blood



Many different cell-types show specific eQTL effects

Module number	Module description	Number of affected eQTLs # eQTLs in strong LD with known GWAS hits	GO biological process top enriched pathway
1	Neutrophils 1	917	Detection of bacterium
2	CD4+ T-cells	337	T cell selection
3	NK cells / CD8+ T-cells	226	Cellular defense response
4	Erythrocytes	188	Hemoglobin metabolic process
5	Monocytes / Macrophages	181	Defense response to virus
6	Growth factor	156	Nerve growth factor receptor signaling pathway
7	Type 1 interferon	145	Regulation of defense response
8	Neutrophils 2	121	Detection of bacterium
9	B-cells	123	B cell receptor signaling pathway
10	Eosinophil	120	Regulation of myeloid leukocyte mediated immunity

Detection of regulatory networks



Conclusions

- ▶ GWAS allows identification of genomic regions that harbor a genetic risk factor for complex traits and diseases
 - ▶ They do not directly point to function
- ▶ eQTLs can be used to link genes to disease
- ▶ Combing eQTL with environment allows deeper insight into systems