

expression quantitative trait loci (eQTL)

mapping & interpretation

Ankeeta Shah

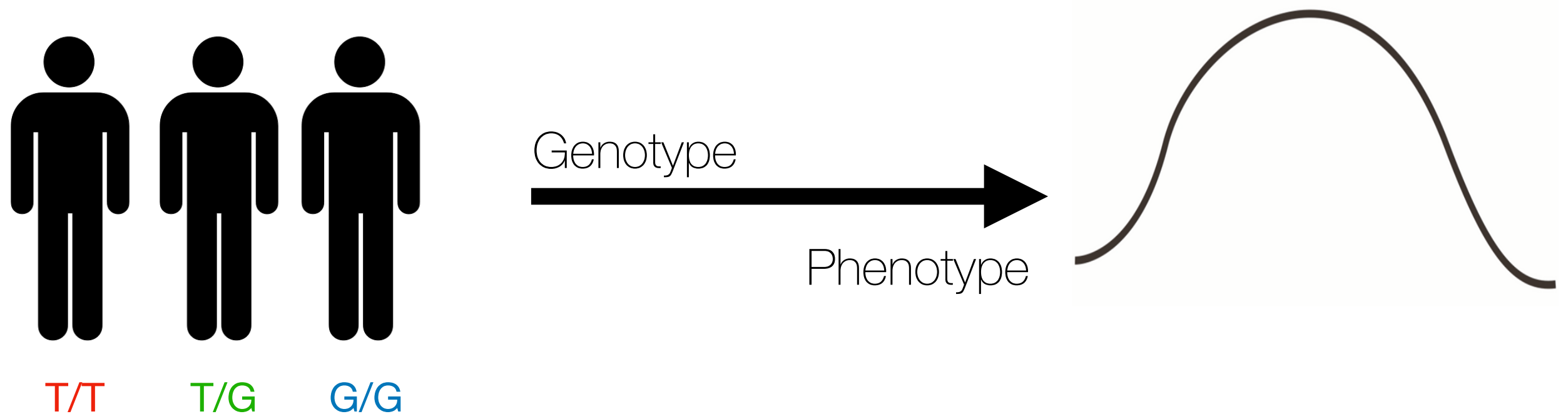
Li Lab, GGSB

HGEN 47000 Human Genetics I

Computational Workshop II

November 20, 2019

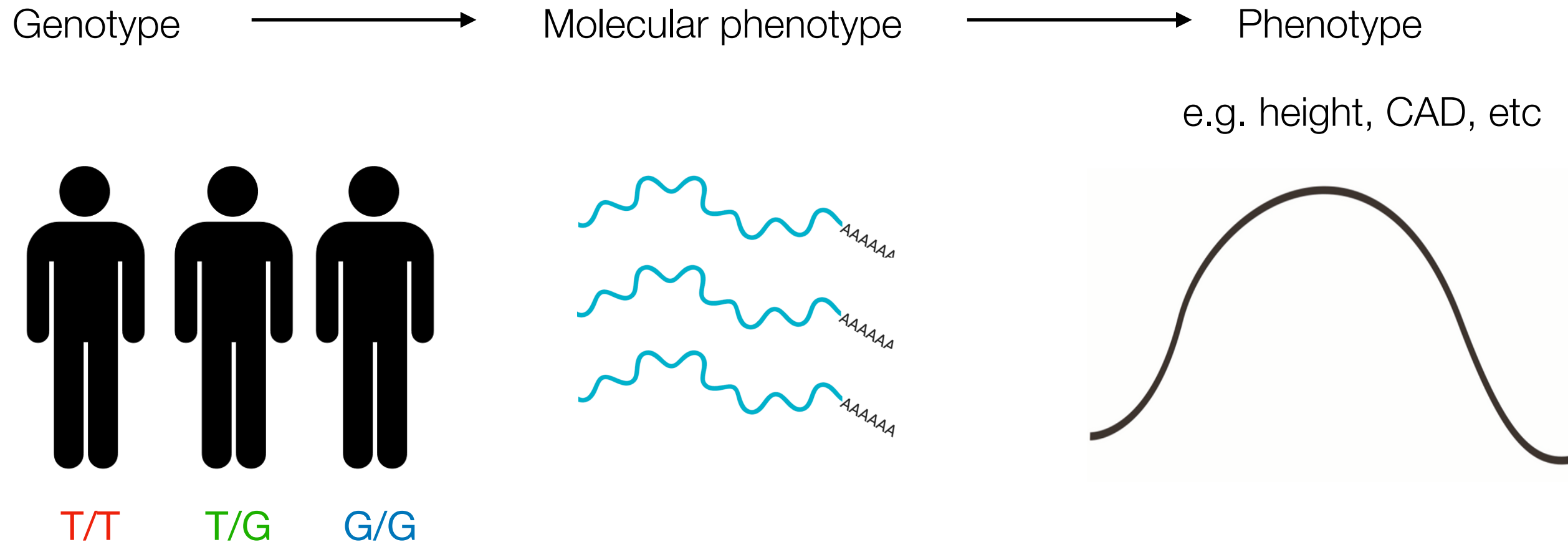
Understanding how genetic variation contributes to phenotypic variation in humans



Why should we map quantitative trait loci (QTL)?

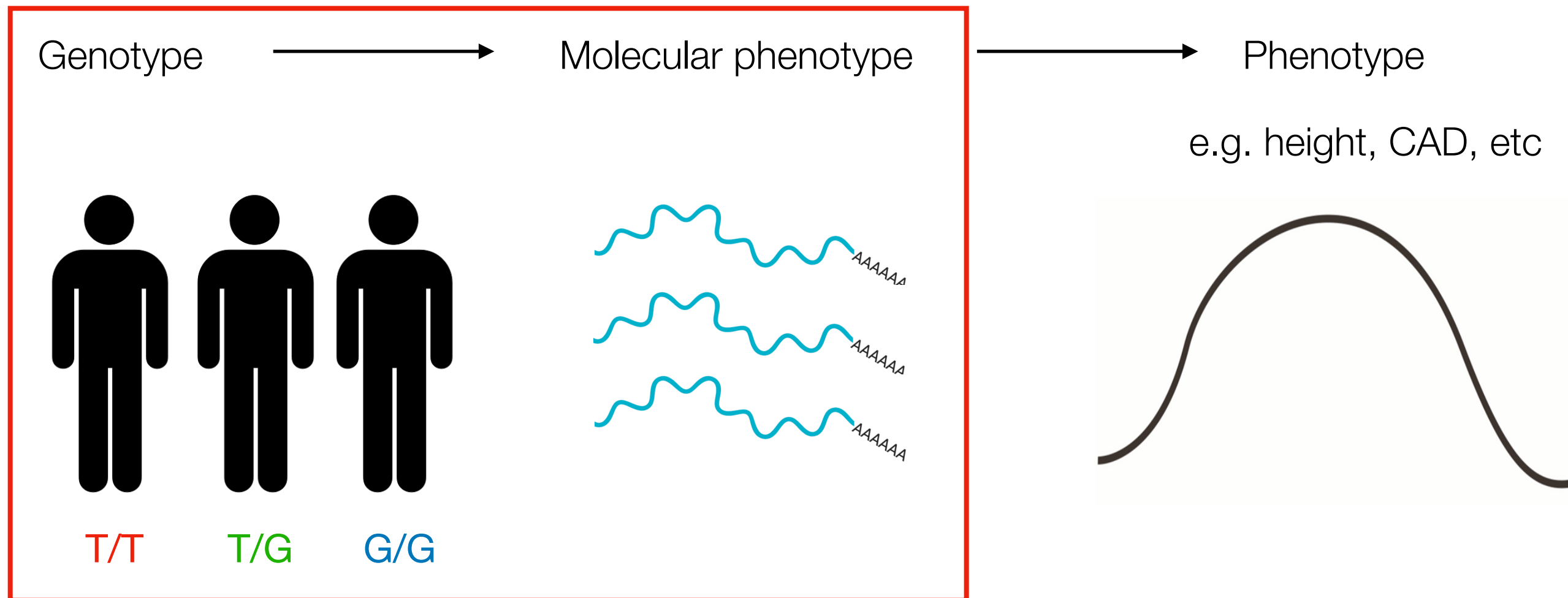
Why should we map QTL?

quantitative trait loci (QTL): regions of the genome that affect the levels of a heritable quantitative trait



Why should we map QTL?

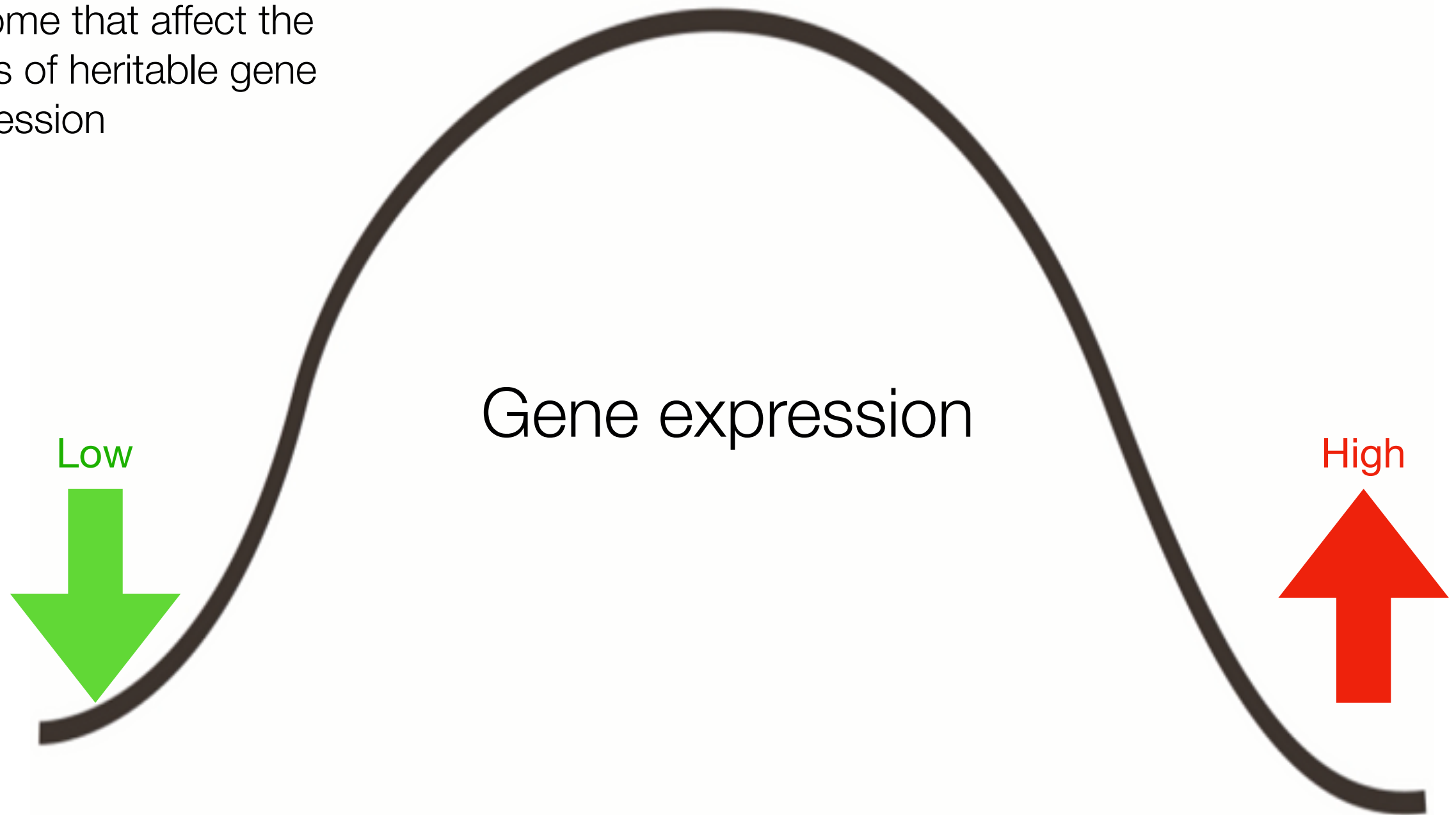
quantitative trait loci (QTL): regions of the genome that affect the levels of a heritable quantitative trait



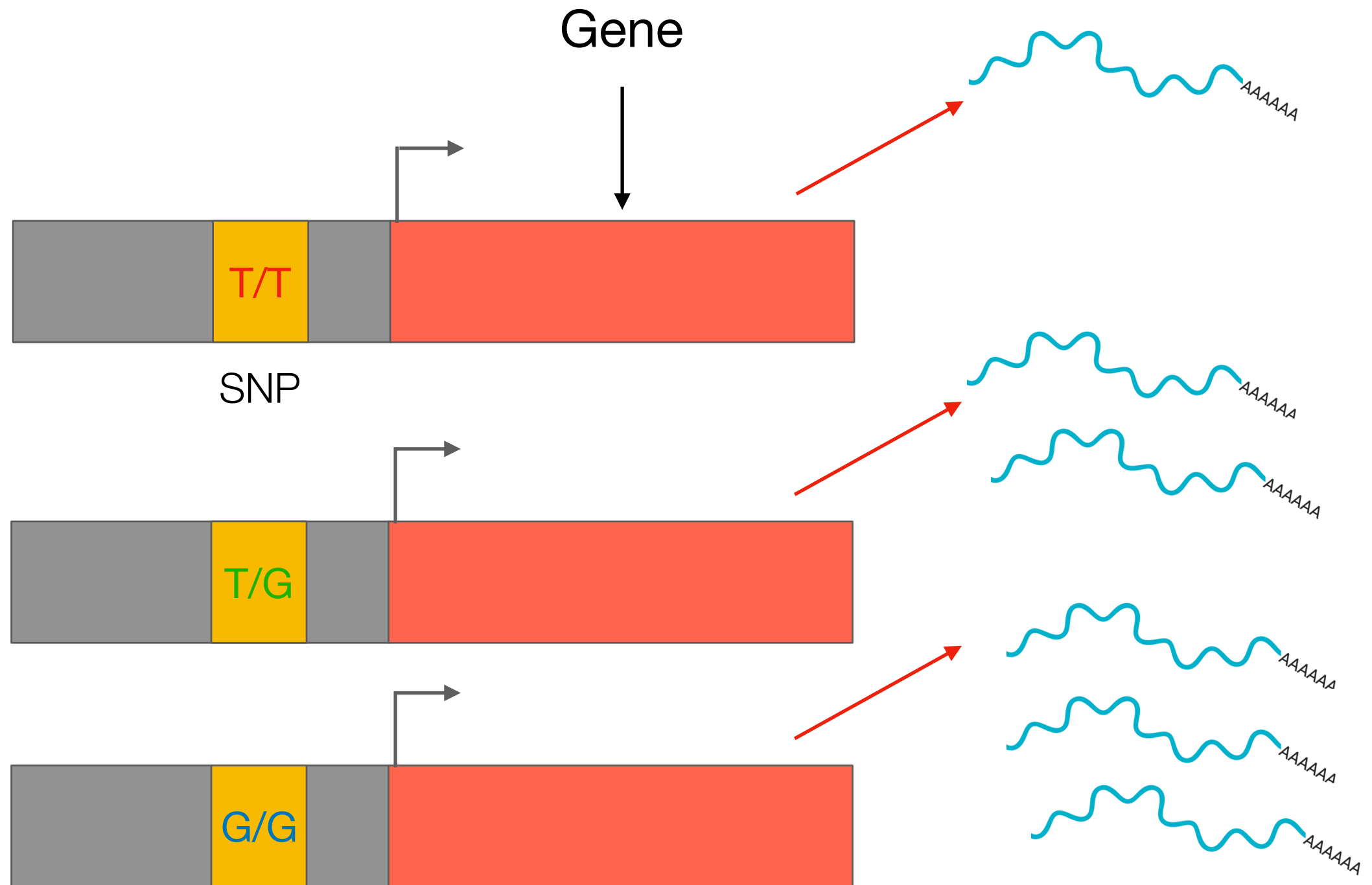
Why should we map expression QTL (eQTL)?

Gene expression is continuous (i.e. “quantitative trait”)

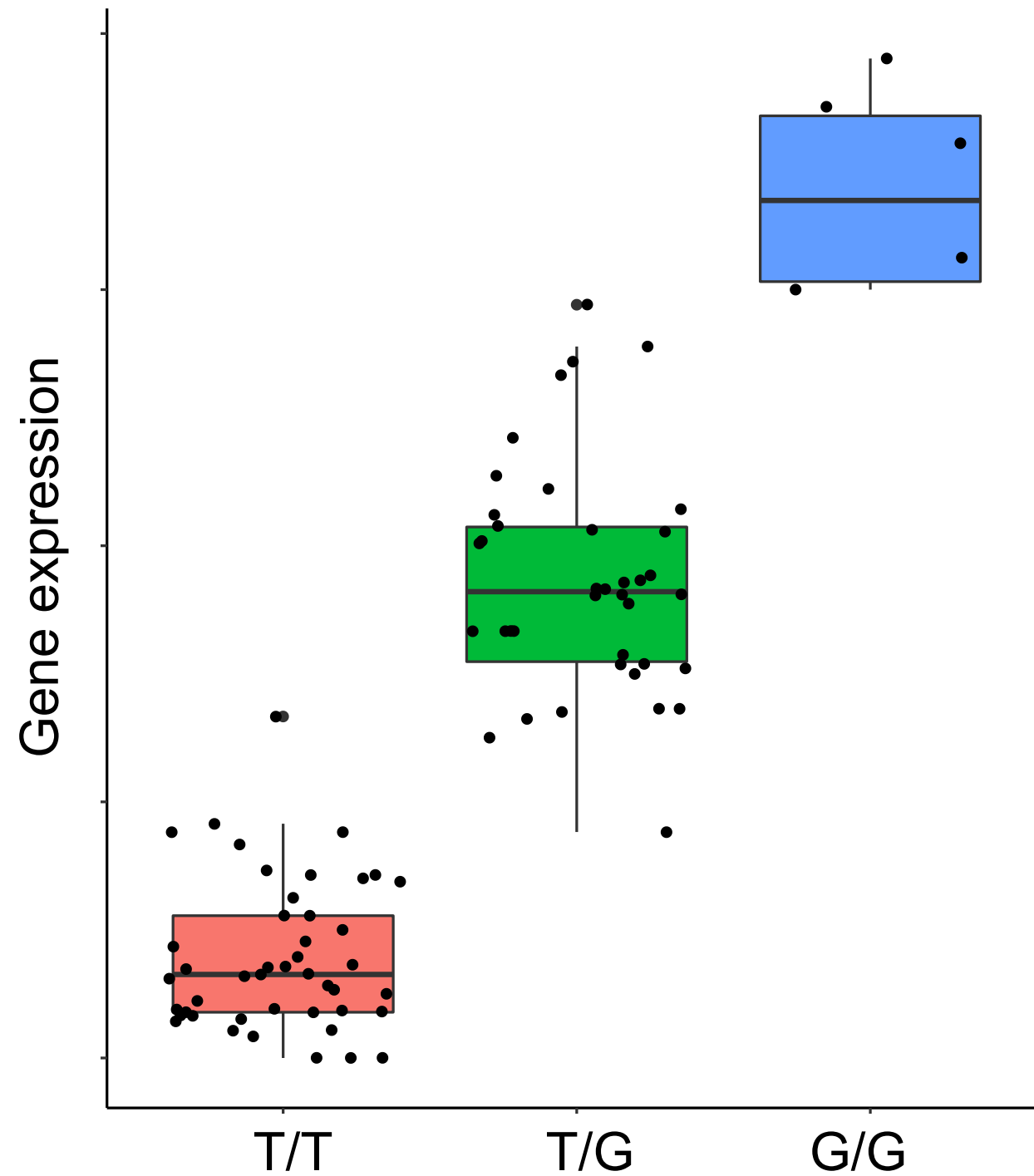
eQTLs: regions of the genome that affect the levels of heritable gene expression



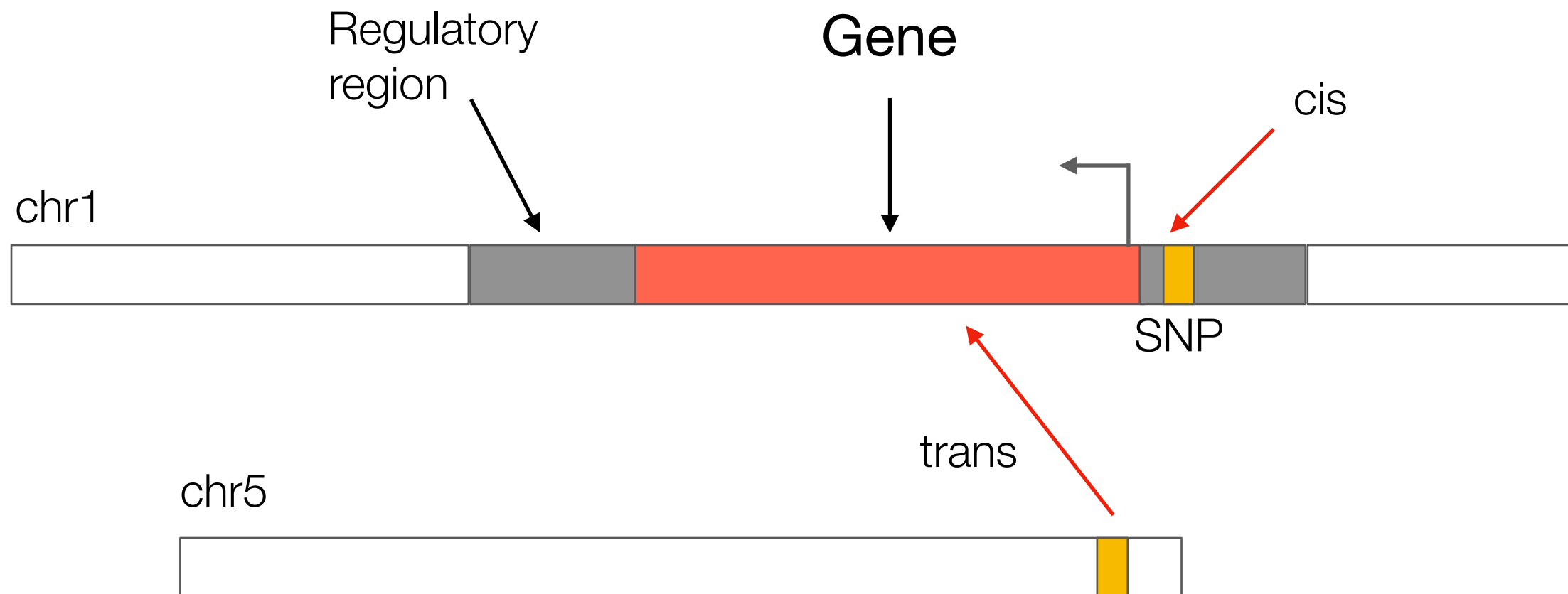
eQTL mapping



eQTL mapping



cis- vs trans-eQTL mapping



cis: test all SNPS within 1MB of the gene

trans: test all other SNPs in the genome

GEUVADIS* Consortium data

Entire dataset: Lymphoblastoid cell lines (LCLs)

Individuals from the 1000 Genomes Project

Across 5 populations:

CEPH (CEU),

Finns (FIN),

British (GBR),

Toscani (TSI),

and Yoruba (YRI))

Total of ~500 individuals with genotypes and RNA-seq

We are going to be working with the Yoruba (YRI) individuals (N = 89)

*Lappalainen et al., *Nature*, 2013

You can find the raw data here: [https://
www.internationalgenome.org/data-portal/
data-collection/geuvadis](https://www.internationalgenome.org/data-portal/data-collection/geuvadis)

Data processing steps

1. **Pre-processing:**

1. Filtering allele-specific biases
2. Normalizing gene counts (transcripts per million, TPM)

2. **Normalization:**

1. Quantile normalize the gene expression data (i.e. make an adjustment such that the expression data is normally distributed).

3. **Association:**

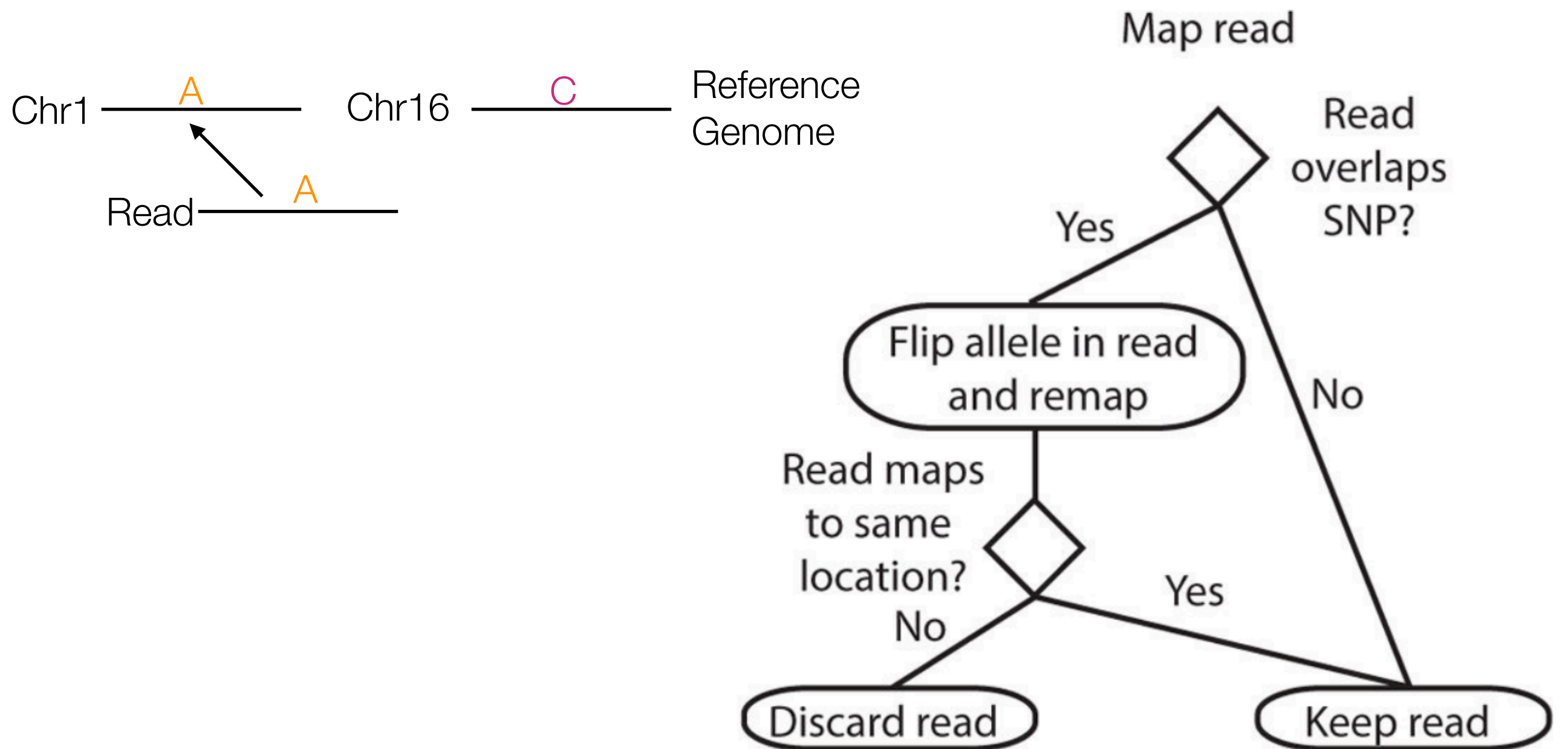
1. Test for association between SNPs and gene expression, while adjusting for covariates

4. **Significance testing:**

1. Permutations
2. Multiple testing correction

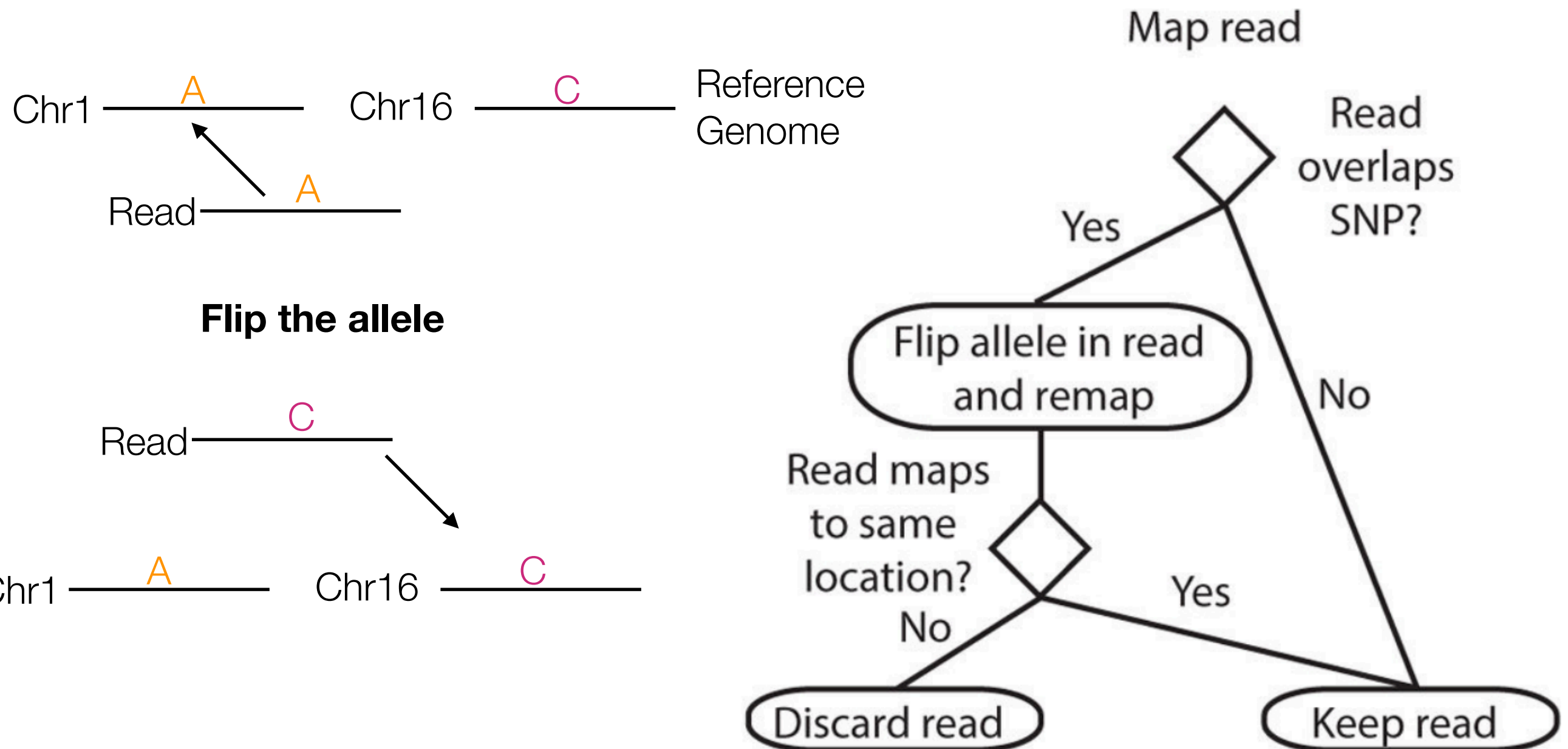
Filtering allele-specific biases

Read alignment with STAR (Dobin et al., *Bioinformatics*, 2013). Note: we have to be careful to remove allele-specific mapping biases (to avoid false positives).



Filtering allele-specific biases

Read alignment with STAR (Dobin et al., *Bioinformatics*, 2013). Note: we have to be careful to remove allele-specific mapping biases (to avoid false positives).



Transcripts per kilobase million (TPM)

1. Count reads overlapping genes (e.g. featureCounts, Liao et al., *Bioinformatics*, 2014)

Transcripts per kilobase million (TPM)

1. Count reads overlapping genes (e.g. featureCounts, Liao et al., *Bioinformatics*, 2014)
2. Normalize counts for gene length (divide rows by length in kb)

		Individuals		
Genes	5 kb	2	20	1
	7 kb	4	35	3
	6 kb	11	32	2

Transcripts per kilobase million (TPM)

1. Count reads overlapping genes (e.g. featureCounts, Liao et al., *Bioinformatics*, 2014)
2. Normalize counts for gene length (divide rows by length in kb)

		Individuals					Individuals		
Genes	5 kb	2	20	1	→	Genes	2/5	20/5	1/5
	7 kb	4	35	3			4/7	35/7	3/7
	6 kb	11	32	2			11/6	32/6	2/6

Transcripts per kilobase million (TPM)

1. Count reads overlapping genes (e.g. featureCounts, Liao et al., *Bioinformatics*, 2014)
2. Normalize counts for gene length (divide rows by length in kb)

		Individuals				Individuals				Individuals		
Genes	5 kb	2	20	1		$2/5$	$20/5$	$1/5$		0.4	4	0.2
	7 kb	4	35	3	Genes	$4/7$	$35/7$	$3/7$		0.57	5	0.43
	6 kb	11	32	2		$11/6$	$32/6$	$2/6$		1.8	5.3	0.33

3. Normalize counts for sequencing depth (divide columns by sum columns)

		Individuals		
Genes	0.4	4	0.2	
	0.57	5	0.43	
	1.8	5.3	0.33	
	2.8	14.3	0.96	

Transcripts per kilobase million (TPM)

1. Count reads overlapping genes (e.g. featureCounts, Liao et al., *Bioinformatics*, 2014)
2. Normalize counts for gene length (divide rows by length in kb)

		Individuals				Individuals				Individuals		
Genes	5 kb	2	20	1		$2/5$	$20/5$	$1/5$		0.4	4	0.2
	7 kb	4	35	3	Genes	$4/7$	$35/7$	$3/7$		0.57	5	0.43
	6 kb	11	32	2		$11/6$	$32/6$	$2/6$		1.8	5.3	0.33

3. Normalize counts for sequencing depth (divide columns by (sum columns $\times 10^{-6}$))

Individuals

0.4	4	0.2
0.57	5	0.43
1.8	5.3	0.33
2.8	14.3	0.96

Genes

→

Genes

0.4/2.8 x 10^-6	4/14.3 x 10^-6	0.2/0.96 x 10^-6
0.57/2.8 x 10^-6	5/14.3 x 10^-6	0.43/0.96 x 10^-6
1.8/2.8 x 10^-6	5.3/14.3 x 10^-6	0.33/0.96 x 10^-6

TPM CODE DEMO

Quantile normalization

There may be **systematic differences** between individuals that we can observe in gene expression data that have nothing to do with genotype.

Normalizing the data (e.g. TPM) will allow us to obtain identical distributions of gene expression across individuals.

	Individual 1	Individual 2	Individual 3	Individual 4
Gene 1	2	20	1	7
Gene 2	4	35	3	10
Gene 3	4	32	3	9
Gene 4	3	30	3	8
Gene 5	3	30	3	8

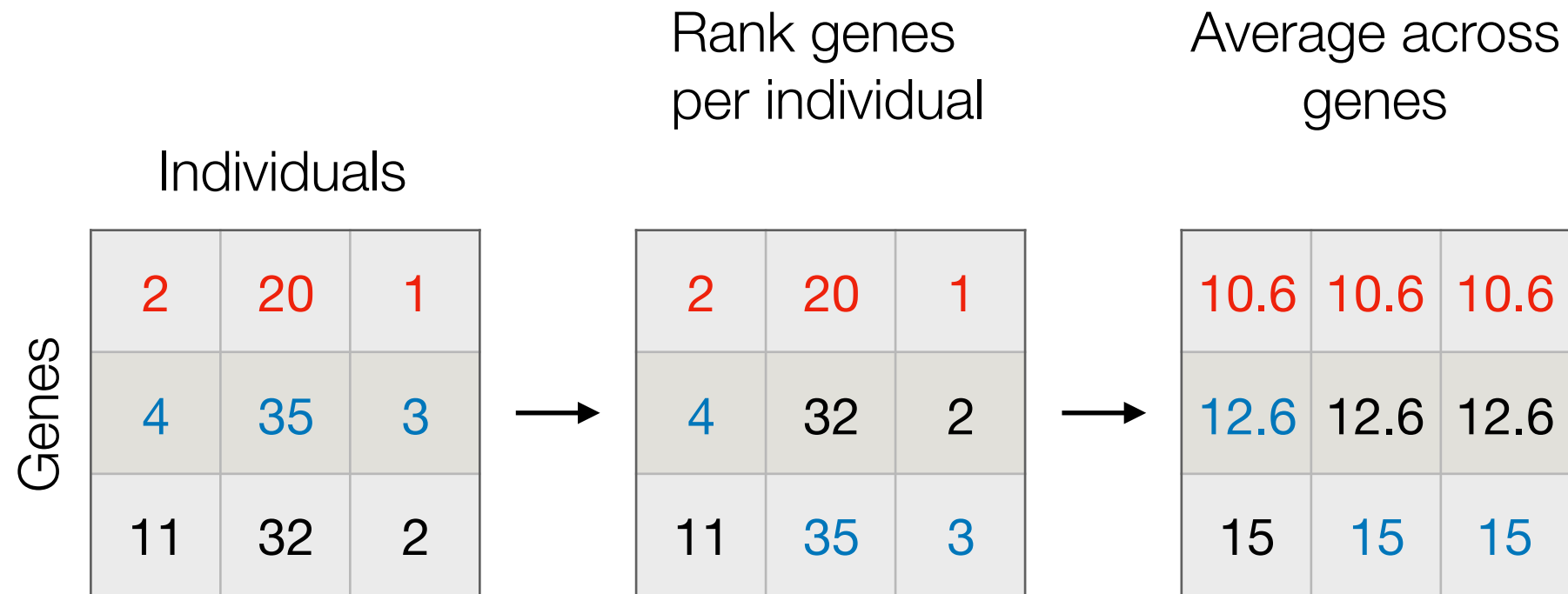
Quantile normalization

		Individuals		
Genes		2	20	1
		4	35	3
		11	32	2

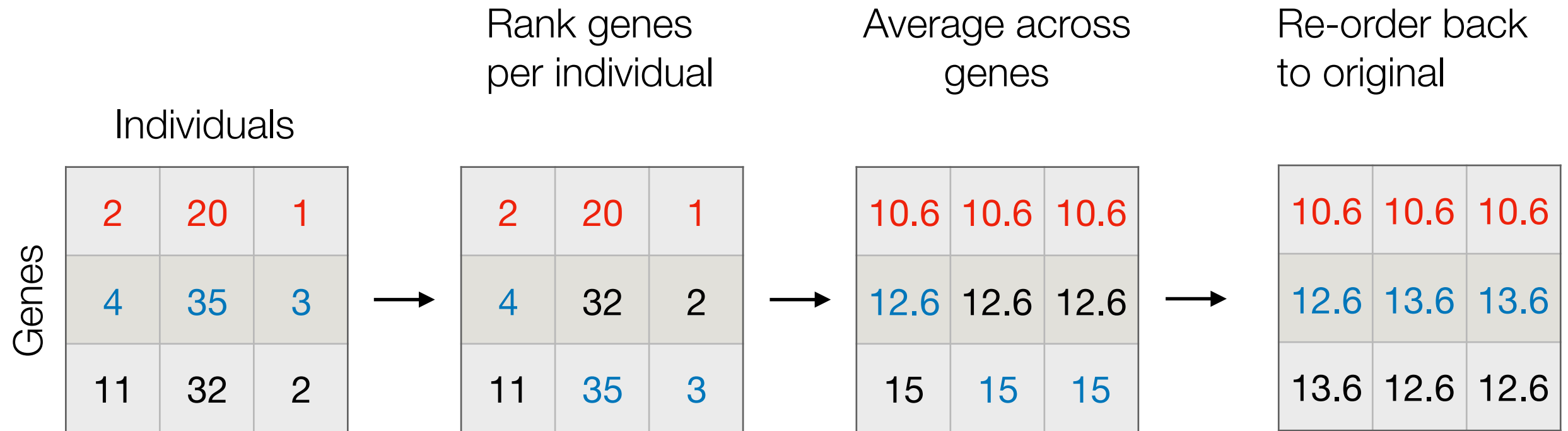
Quantile normalization

		Individuals					Rank genes per individual		
Genes	1	2	20	1	→	1	2	20	1
	2	4	35	3		2	4	32	2
	3	11	32	2		3	11	35	3

Quantile normalization



Quantile normalization



Linear mixed models (LMMs)

Goal: Use a linear model to test for association between gene expression (Y) and SNPs (X), while adjusting for covariates (Z)

The diagram shows the equation $Y = X\beta + Z + \epsilon$ with arrows pointing to each term from descriptive labels:

- An arrow from "SNPs" points to X .
- An arrow from "Effect sizes" points to β .
- An arrow from "Covariates" points to Z .
- An arrow from "Error" points to ϵ .
- An arrow from "Molecular phenotype (e.g. gene expression)" points to Y .

Covariates:

Known factors: sex, age, batch effects, etc

Unmeasured factors: cell type composition, cell growth rate, etc

Regressing out covariates

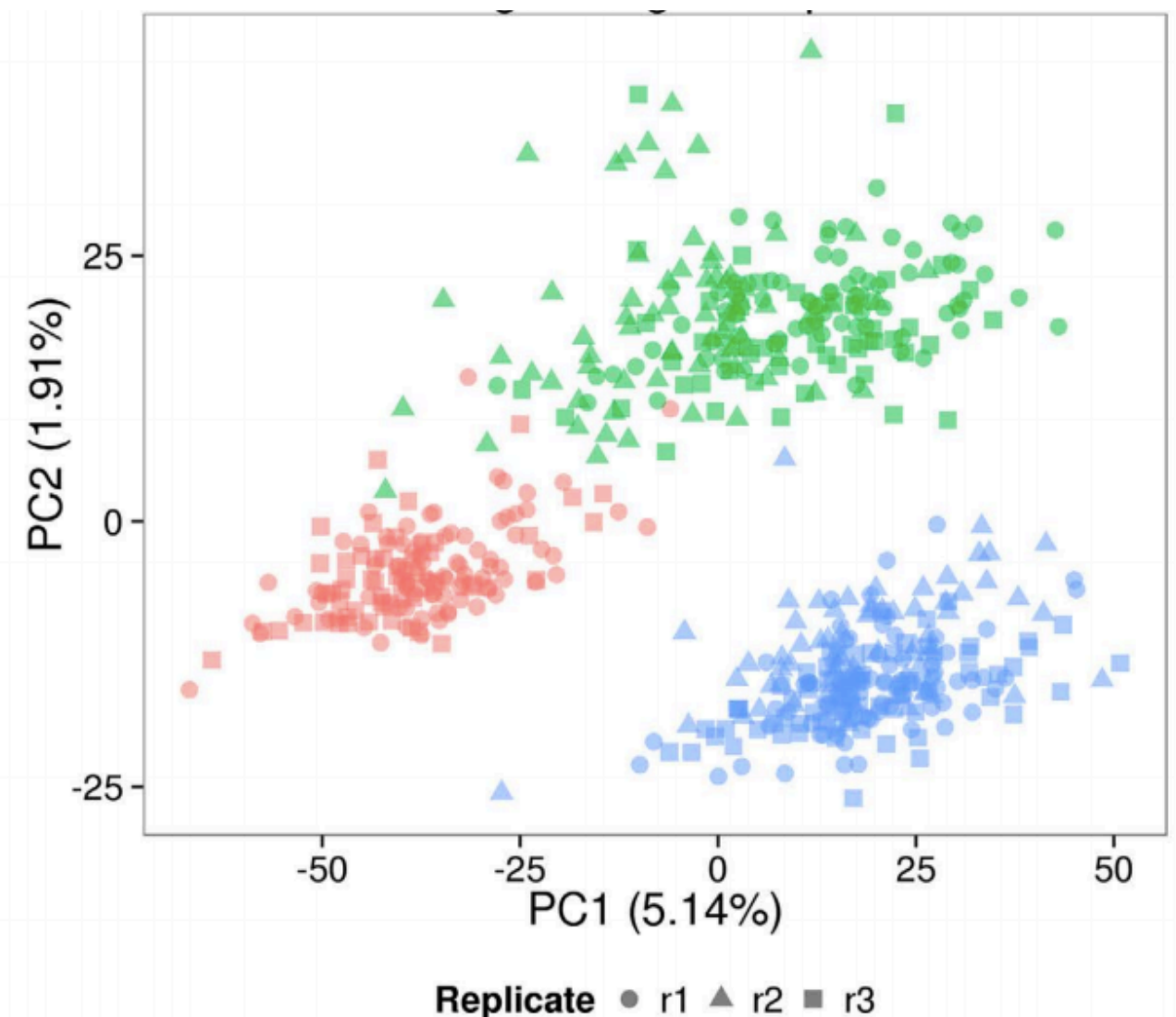
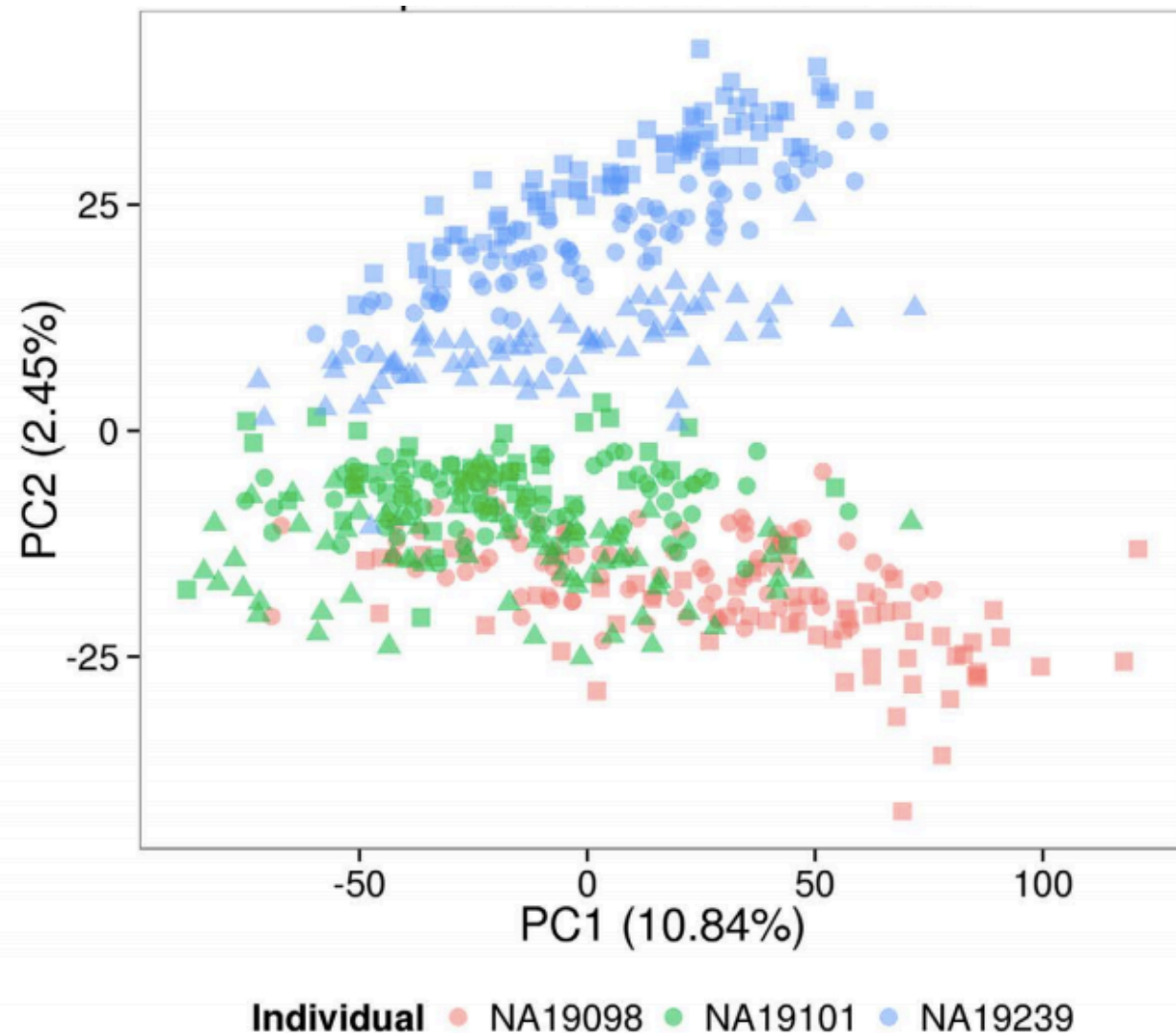
Covariates:

Known factors: sex, age, batch effects, etc

Unmeasured factors: cell type composition, cell growth rate, etc

Top PCs should be regressed out because they do not capture cis-genotype effects (the signal / association we want to measure)

We calculate principal components to be used as covariates in linear regression.



Tung et al., *Scientific Reports*, 2017

QUANTILE NORMALIZATION COVARIATES DEMO

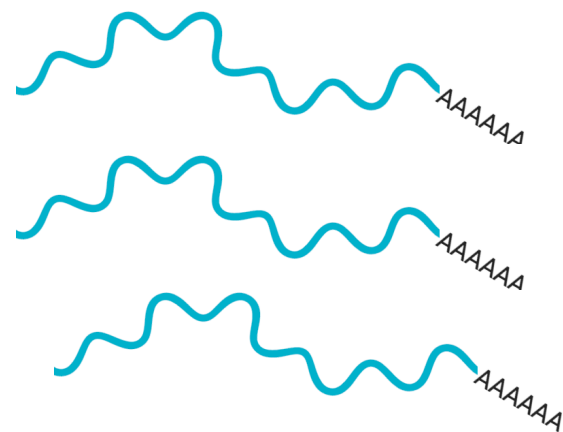
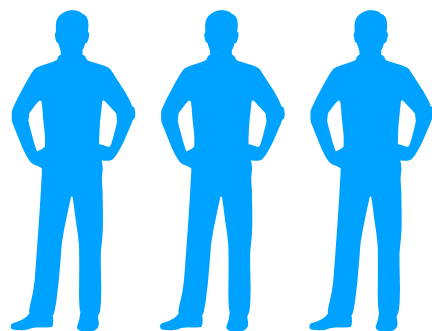
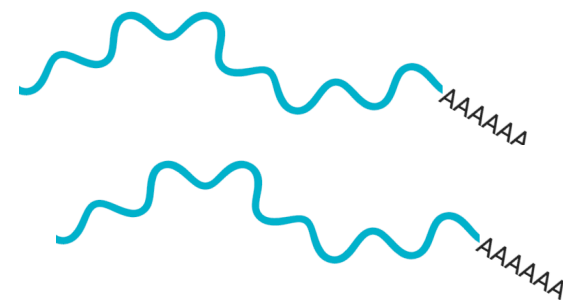
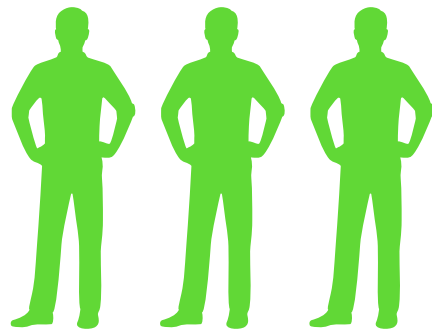
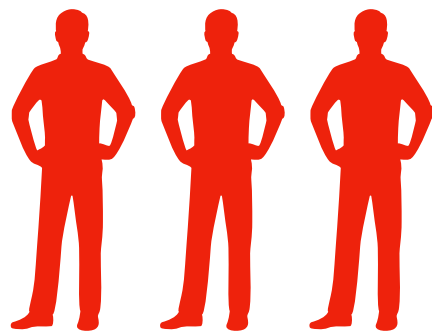
Permutation tests

Nominal p -values: assigned to every SNP-gene pair that is tested (independently)

} α

Genotypes

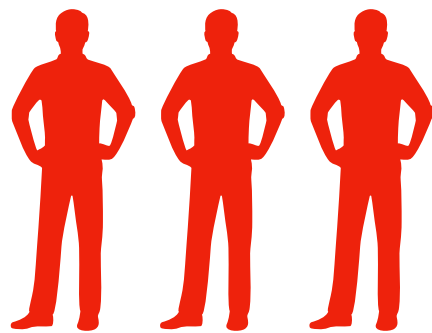
Gene Expression



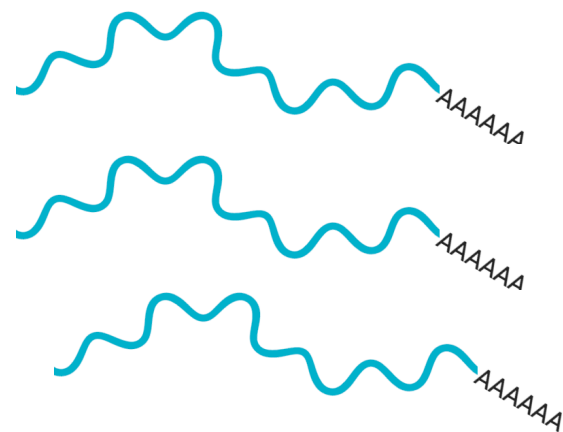
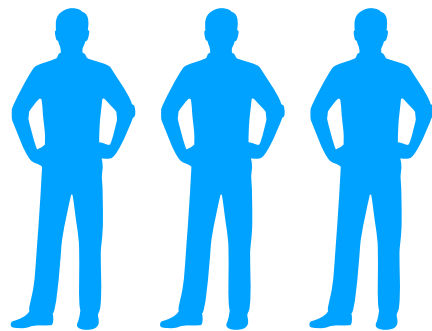
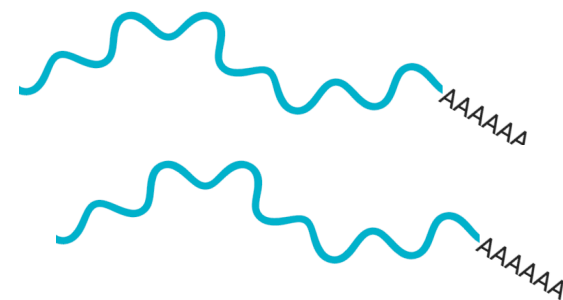
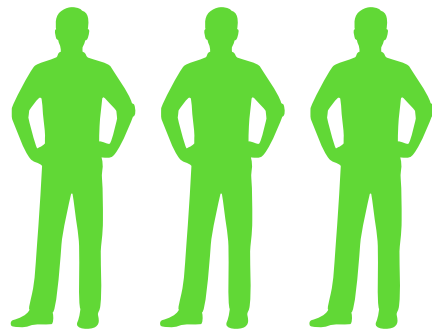
Permutation tests

Permutation tests allow us to assess statistical significance in our dataset.

Genotypes



Gene Expression

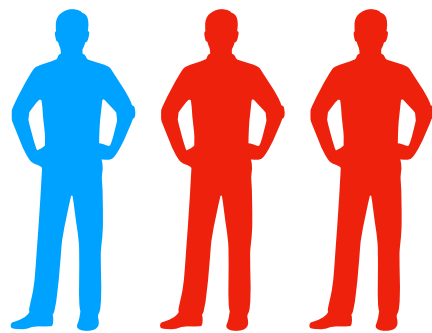


**Scramble the genotypes (1,000 times)
+ re-run associations**

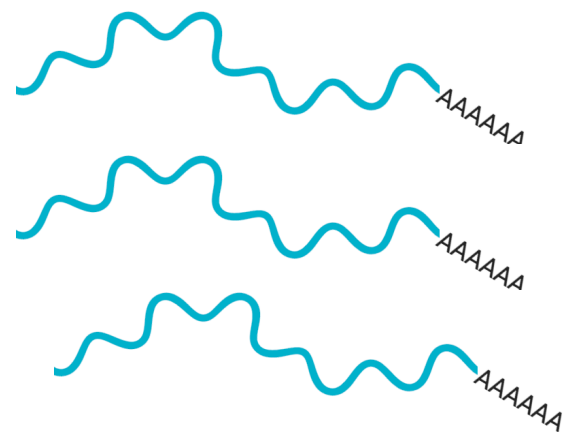
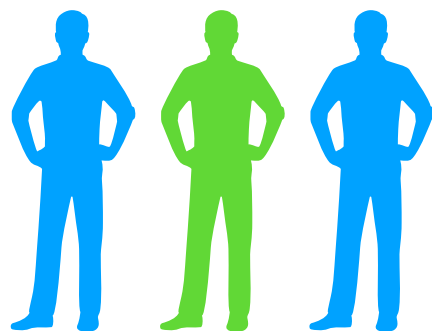
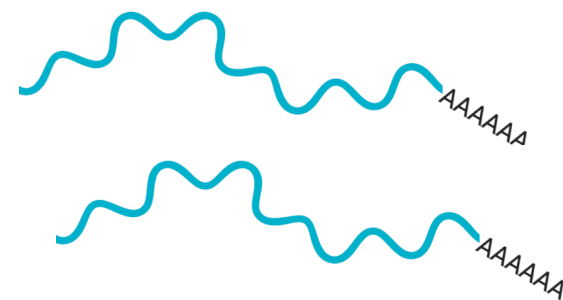
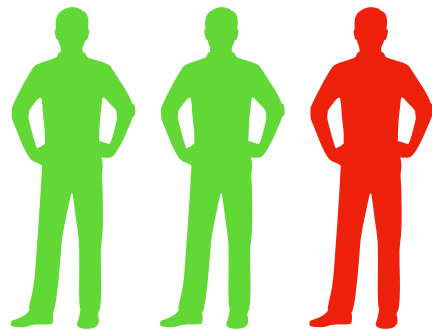
Permutation tests

Permutation tests allow us to assess statistical significance in our dataset.

Genotypes



Gene Expression

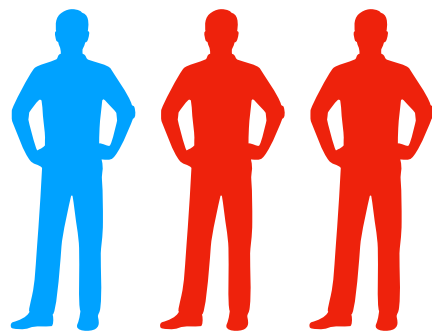


The permutation p -value is the proportion of permutations in which at least one SNP-gene pair p -value is $< \alpha$ (nominal association).

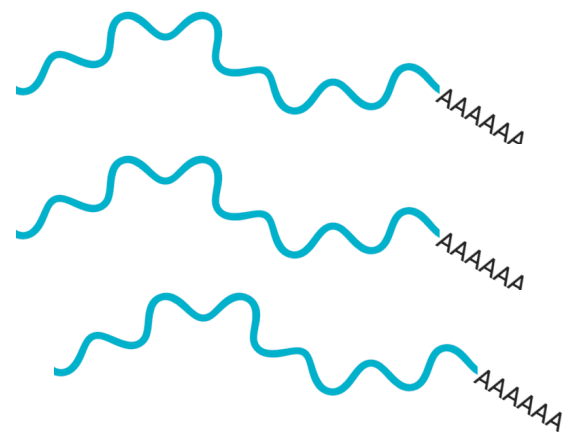
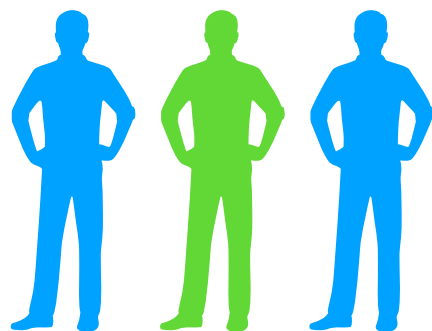
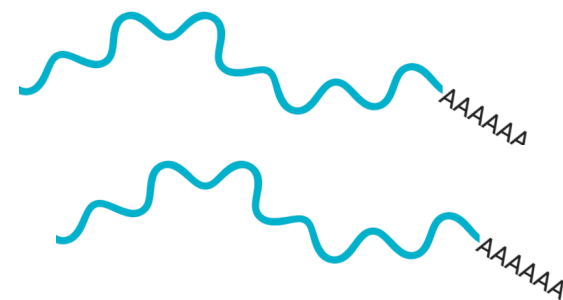
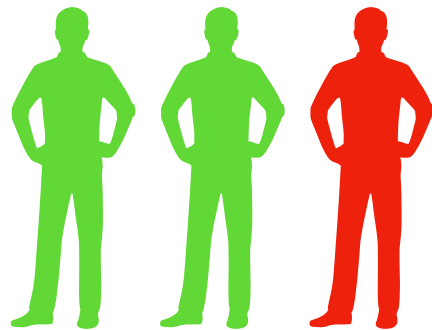
Permutation tests

Permutation tests allow us to assess statistical significance in our dataset.

Genotypes



Gene Expression



The permutation p -value is the proportion of permutations in which at least one SNP-gene pair p -value is $< \alpha$ (nominal association).

NOTE: In a real dataset, we do an *approximation* as we cannot permute all the data.

ASSOCIATION DEMO

Multiple testing correction

As the number of SNP-gene pairs we test increases, the probability of getting a significant association (by random chance) increases.

Multiple testing correction

As the number of SNP-gene pairs we test increases, the probability of getting a significant association (by random chance) increases.

Approach: false discovery rate (FDR) cutoffs can be applied. FDR is defined as the proportion of false positives among all significant results. In QTL mapping, we tend to assume an $FDR < 5\%$ or $FDR < 10\%$.

Multiple testing correction

As the number of SNP-gene pairs we test increases, the probability of getting a significant association (by random chance) increases.

Approach: false discovery rate (FDR) cutoffs can be applied. FDR is defined as the proportion of false positives among all significant results. In QTL mapping, we tend to assume an $FDR < 5\%$ or $FDR < 10\%$.

Benjamini-Hochberg (BH) procedure: controls the FDR at the value you specify (e.g. 10%)

	P-value	Rank	(Rank / Tests)* FDR cutoff	
	0.001	1	0.004	
	0.005	2	0.008	
	0.007	3	0.012	
	
*	0.09	24	0.096	Highest p-val smaller than critical value
	0.12	25	0.1	

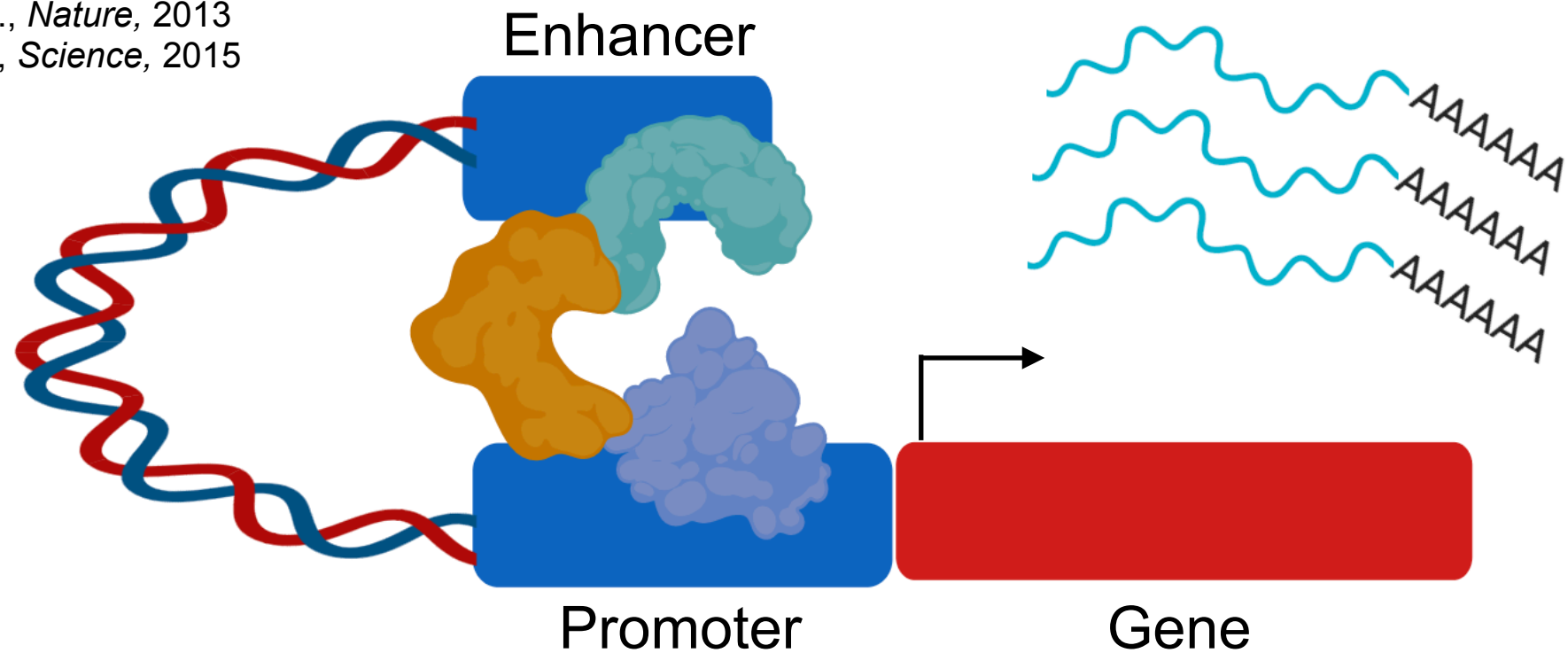
*All above (including this red row) are considered significant

MTC DEMO

What we have learned through extensive eQTL analysis...

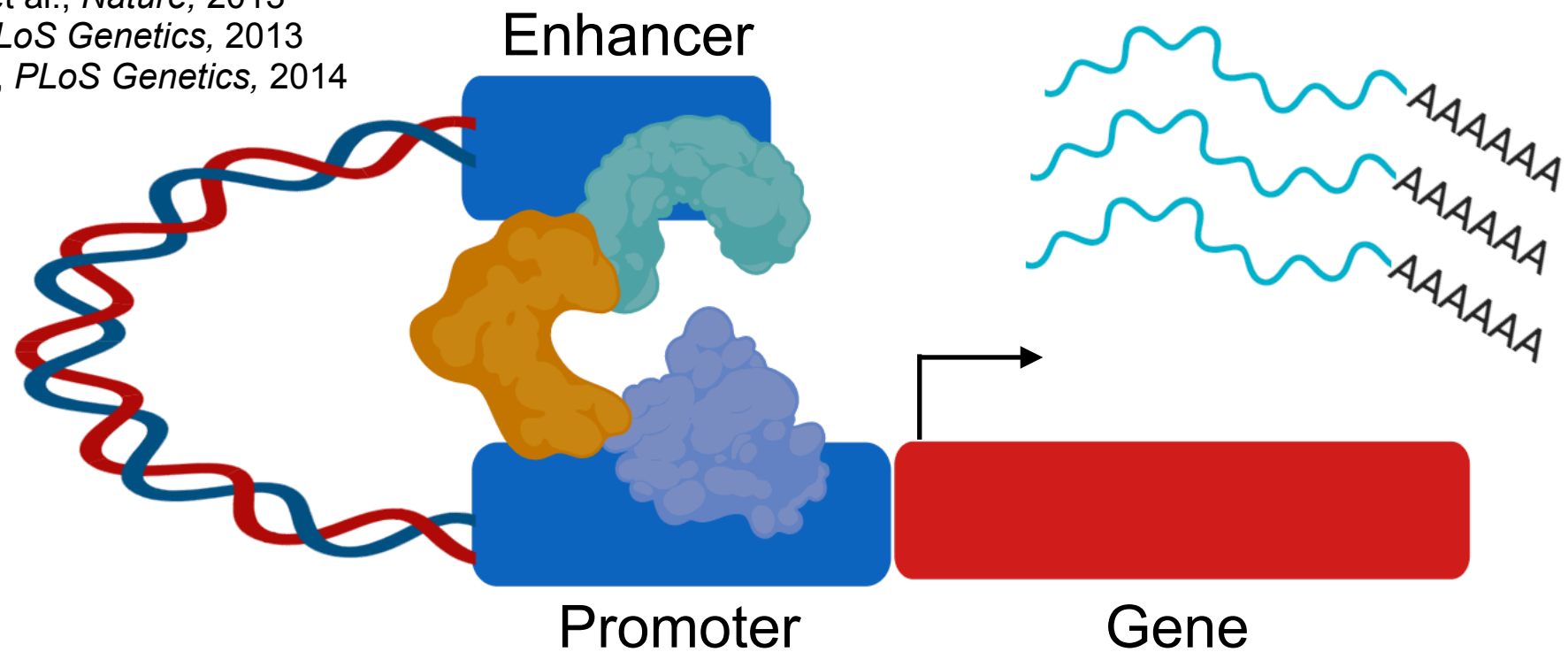
eQTL are enriched in enhancers and promoters

Lappalainen et al., *Nature*, 2013
GTEx Consortium, *Science*, 2015

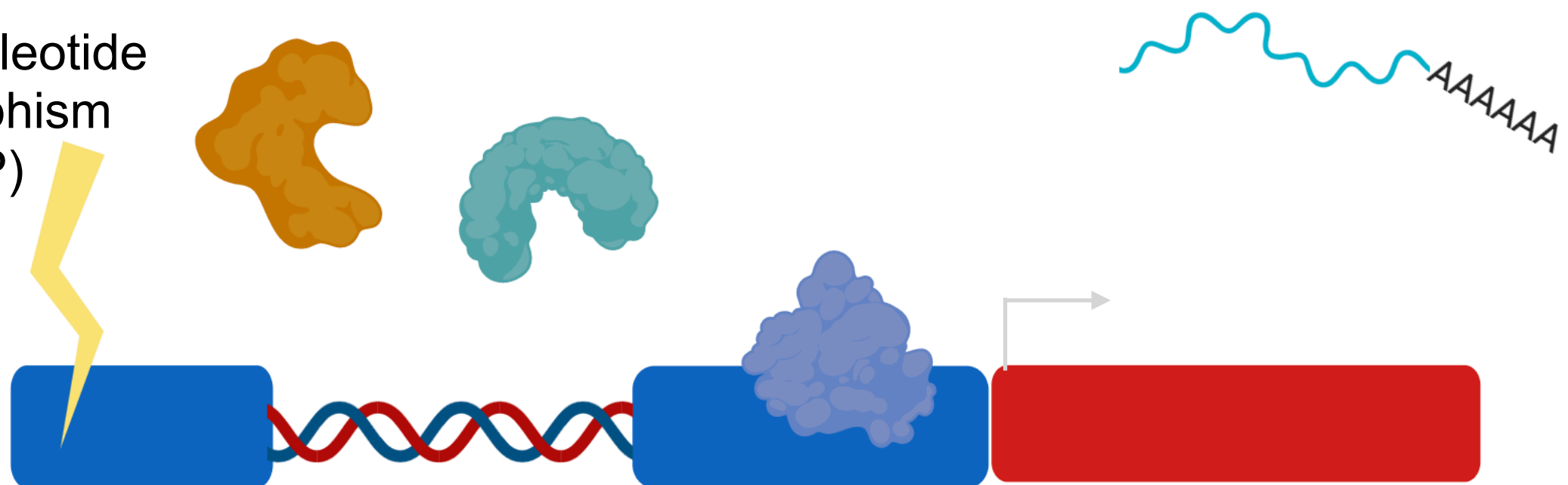


eQTL often affect transcription factor binding sites

Lappalainen et al., *Nature*, 2013
Brown et al., *PLoS Genetics*, 2013
Cusanovich et al., *PLoS Genetics*, 2014

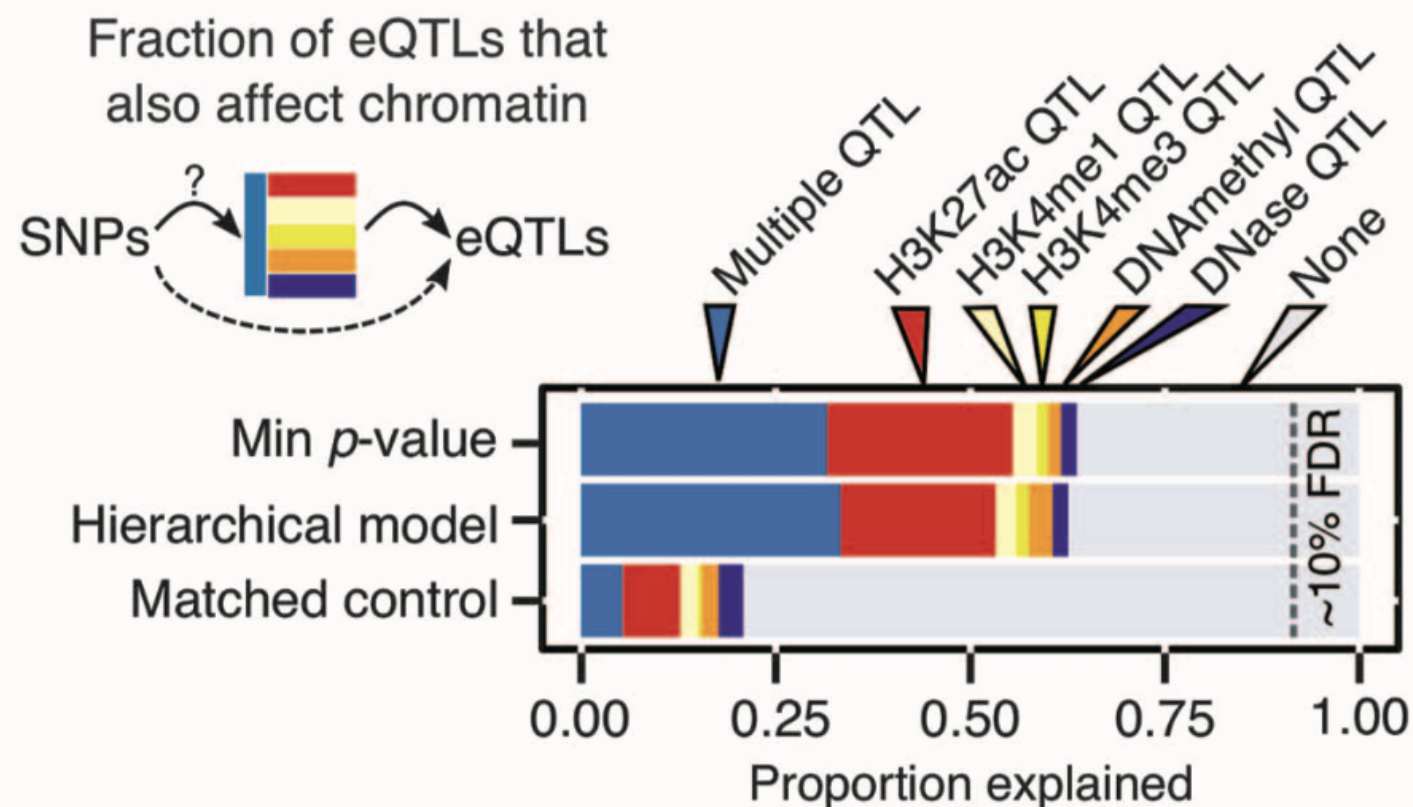
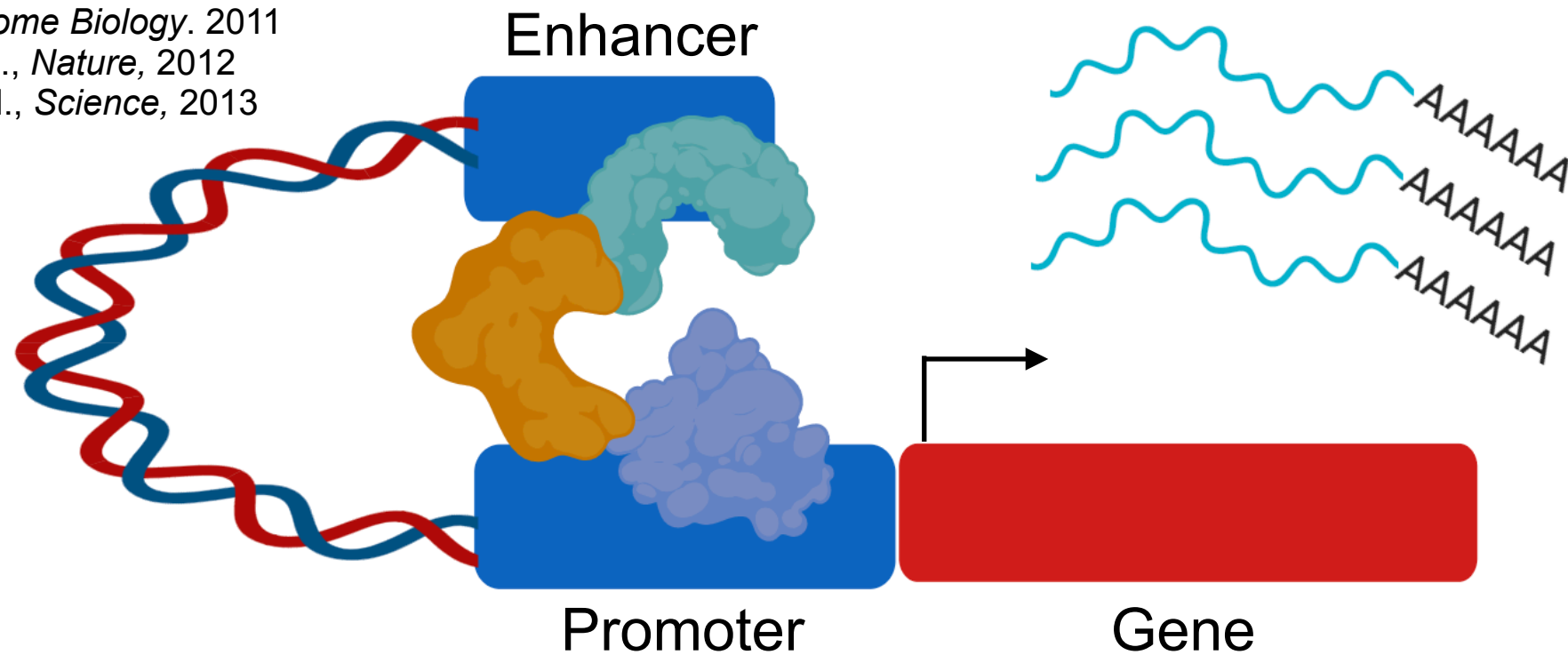


Single nucleotide
polymorphism
(SNP)

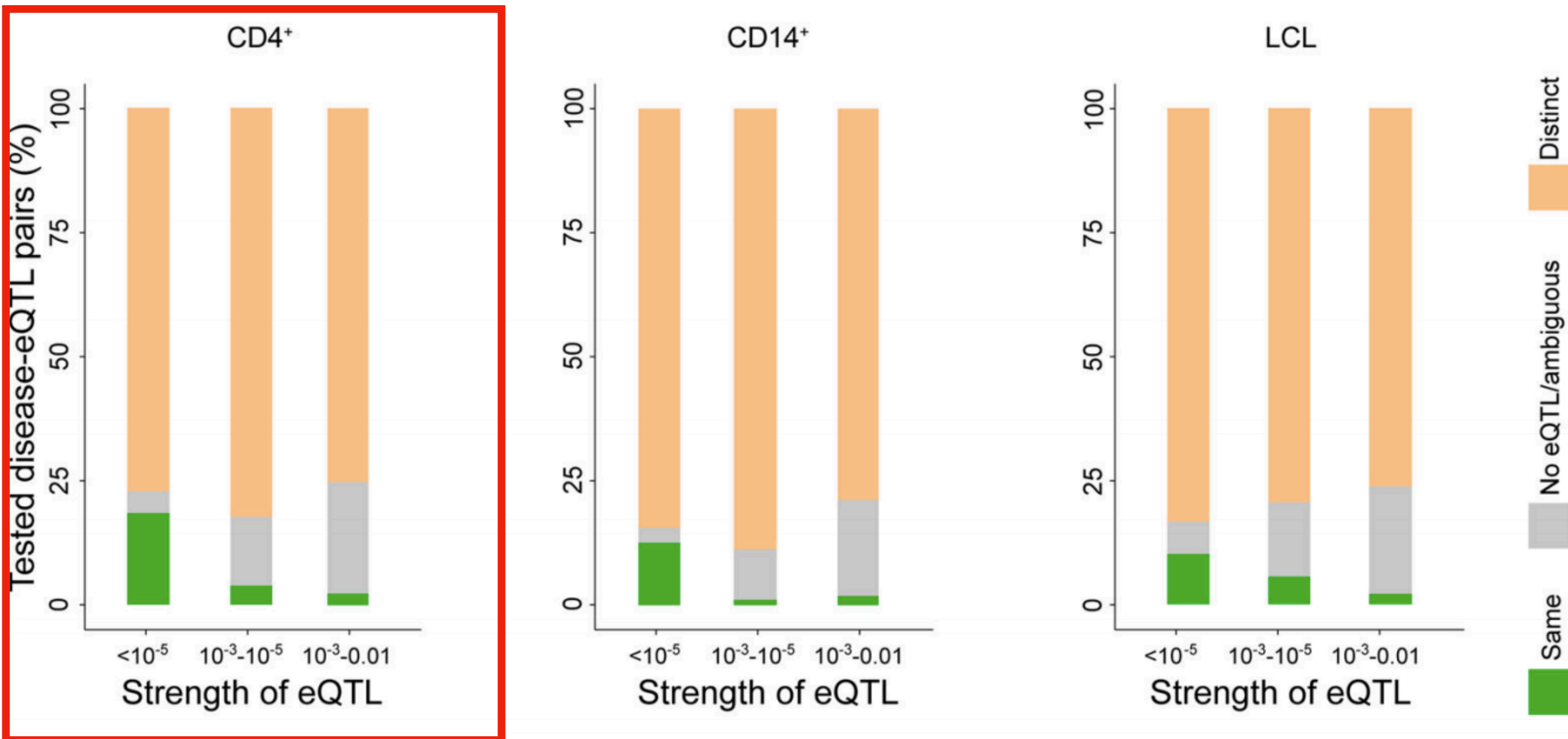


Many eQTL can be explained by chromatin-level phenotypes

Bell et al., *Genome Biology*, 2011
Degner et al., *Nature*, 2012
McVicker et al., *Science*, 2013



75% of *disease-associated* variants disrupt gene regulation in a manner *independent* of total mRNA levels / gene expression



Expanding the repertoire of molecular phenotypes

