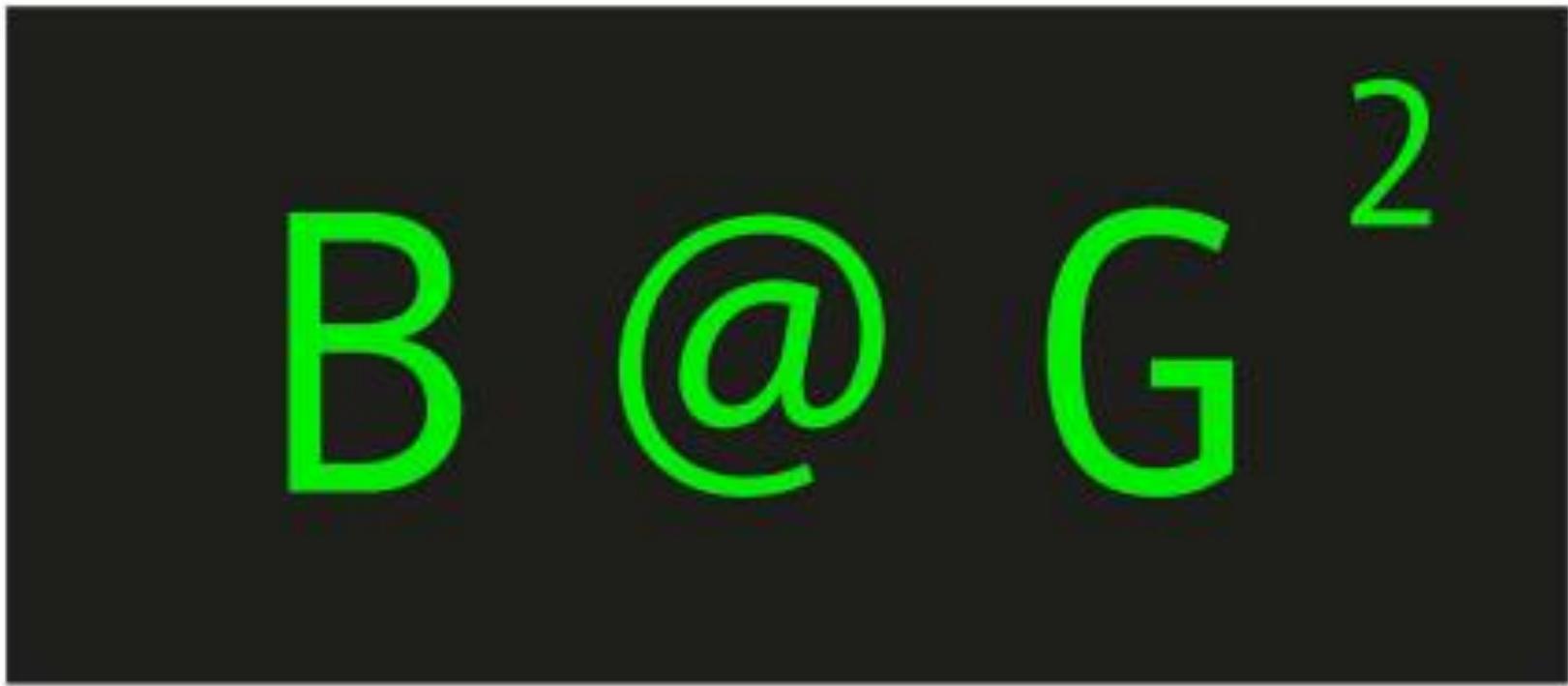


# Winter School 2016 • Weggis



## GWAS

Arthur Korte, Fri 4<sup>th</sup> March

# About me (Disclaimer)

**PhD in Molecular Biology (TU Munich)**

**Postdoc in Population Genetics  
(Gregor Mendel Institute of Molecular Plant Biology, Vienna)**

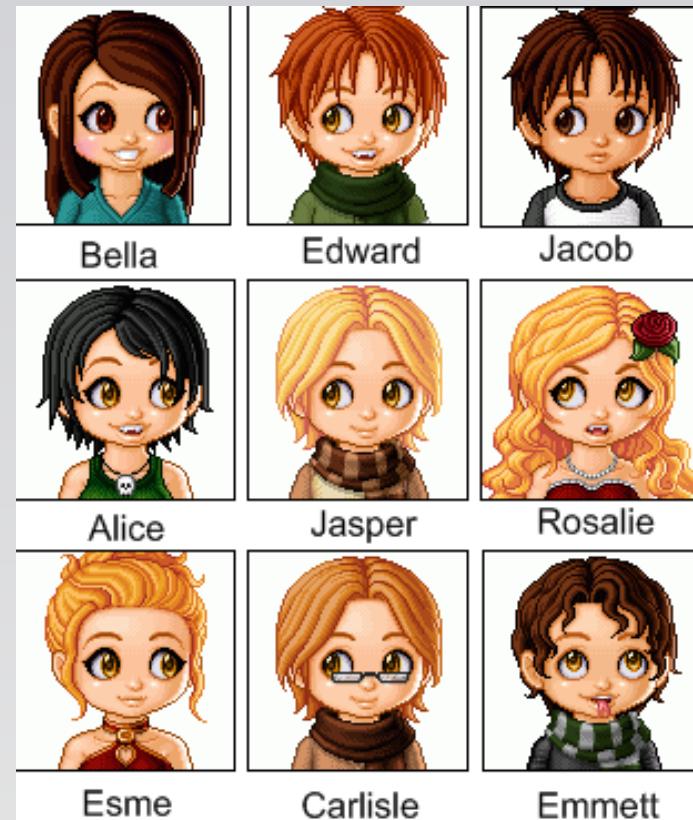
**Junior professor  
Center for Computational Theoretical Biology, University Würzburg**



- 1. Lecture I : Introduction to GWAS and population structure**
- 2. Practical I: GWAS with custom R scripts**
  
- 3. Lecture II: Different GWAS models / software**
- 4. GWAS with webtools**
- 5. Practical II**

# Introduction to GWAS and population structure

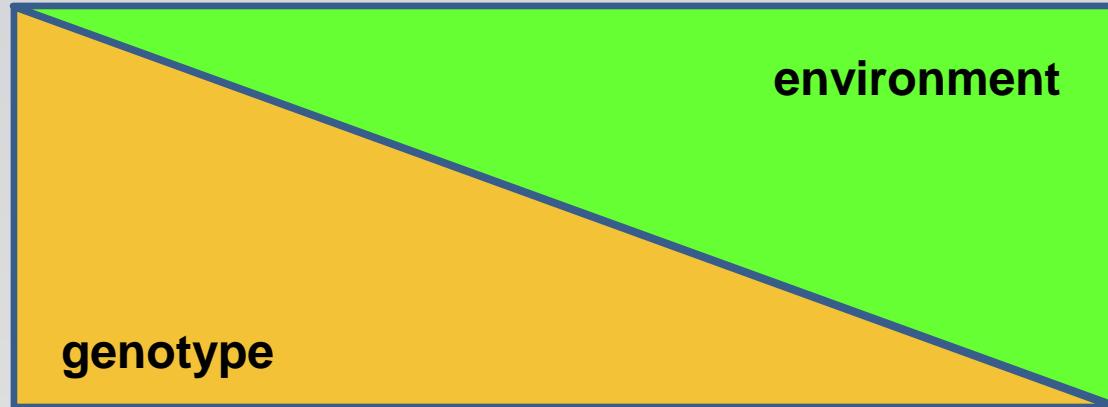
# GWAS



Associations between phenotype and genotype

$$Y = G + E + GxE$$

phenotype



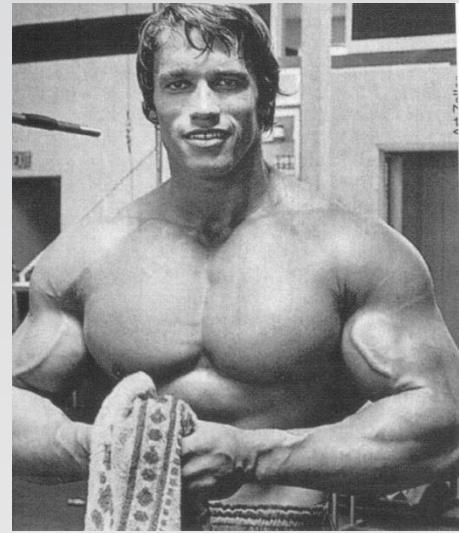
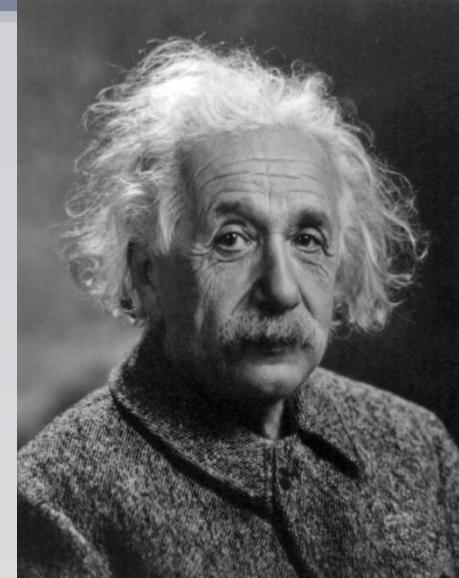
Cystic fibrosis

Breast cancer  
(BRCA1)

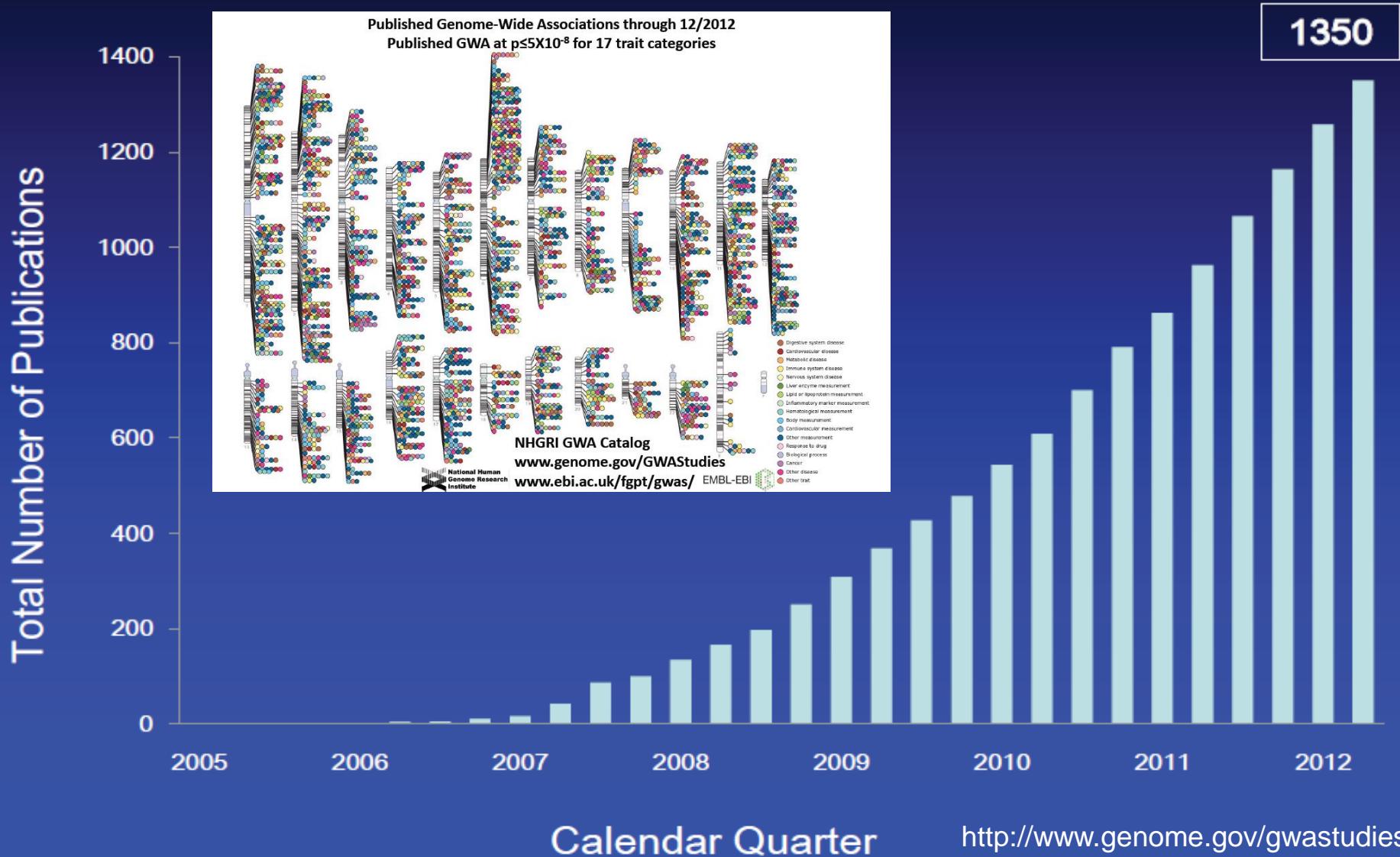
T2D

Lung cancer  
(smoking)

Broken arm



# Published GWA Reports, 2005 – 6/2012



gwas - PubMed - NCBI x

www.ncbi.nlm.nih.gov/pubmed/?term=gwas

NCBI Resources How To Sign in to NCBI

PubMed.gov US National Library of Medicine National Institutes of Health

PubMed gwas Create RSS Create alert Advanced Help

Article types Summary ▾ 20 per page ▾ Sort by Most Recent ▾ Send to: ▾ Filters: Manage Filters

Clinical Trial

Review

Customize ...

Text availability Items: 1 to 20 of 21846

Abstract

Free full text

Full text

PubMed Commons

Reader comments

Trending articles

Publication dates << First < Prev Page 1 of 1093 Next > Last >>

5 years

10 years

Custom range...

Species

Humans

Other Animals

**Search results**

Evaluating the Calibration and Power of Three Gene-Based Association Tests of Rare Variants for the X Chromosome.

Ma C, Boehnke M, Lee S; GoT2D Investigators. Genet Epidemiol. 2015 Oct 10. doi: 10.1002/gepi.21935. [Epub ahead of print] PMID: 26454253

A nonsense mutation in B3GALNT2 is concordant with hydrocephalus in Friesian horses.

Ducro BJ, Schurink A, Bastiaansen JW, Boegheim IJ, van Steenbeek FG, Vos-Lohuis M, Nijman IJ, Monroe GR, Hellinga I, Dibbits BW, Back W, Leegwater PA. BMC Genomics. 2015 Oct 9;16(1):761. PMID: 26452345

Parallelizing Epistasis Detection in GWAS on FPGA and GPU-Accelerated Computing Systems.

New feature Try the new Display Settings option - Sort by Relevance

Results by year Download CSV

Related searches

gwas review

gwas nature

gwas schizophrenia

gwas - PubMed - NCBI x

www.ncbi.nlm.nih.gov/pubmed/?term=gwas

NCBI Resources How To Sign in to NCBI

PubMed.gov US National Library of Medicine National Institutes of Health

PubMed Search Help

Article types Summary 20 per page Sort by Most Recent Send to: Filters: Manage Filters

Clinical Trial

Review

Customize ...

Text availability Items: 1 to 20 of 21846

Abstract << First < Prev Page 1 of 1093 Next > Last >>

Free full text

Full text

PubMed Commons

Reader comments

Trending articles

Publication dates

5 years

10 years

Custom range...

Species

Humans

Other Animals

Evaluating the Calibration and Power of Three Gene-Based Association Tests of Rare Variants for the X Chromosome.

Ma C, Boehnke M, Lee S; GoT2D Investigators. Genet Epidemiol. 2015 Oct 10. doi: 10.1002/gepi.21935. [Epub ahead of print] PMID: 26454253

A nonsense mutation in B3GALNT2 is concordant with hydrocephalus in Friesian horses.

Ducro BJ, Schurink A, Bastiaansen JW, Boegheim IJ, van Steenbeek FG, Vos-Lohuis M, Nijman IJ, Monroe GR, Hellinga I, Dibbits BW, Back W, Leegwater PA. BMC Genomics. 2015 Oct 9;16(1):761. PMID: 26452345

Parallelizing Epistasis Detection in GWAS on FPGA and GPU-Accelerated Computing Systems.

New feature Try the new Display Settings option - Sort by Relevance

Results by year Download CSV

Related searches gwas review gwas nature gwas schizophrenia

# Why GWAS ?



Center for Computational  
and Theoretical Biology



**How do the genetic differences in a population  
translate to distinct phenotypic differences ?**

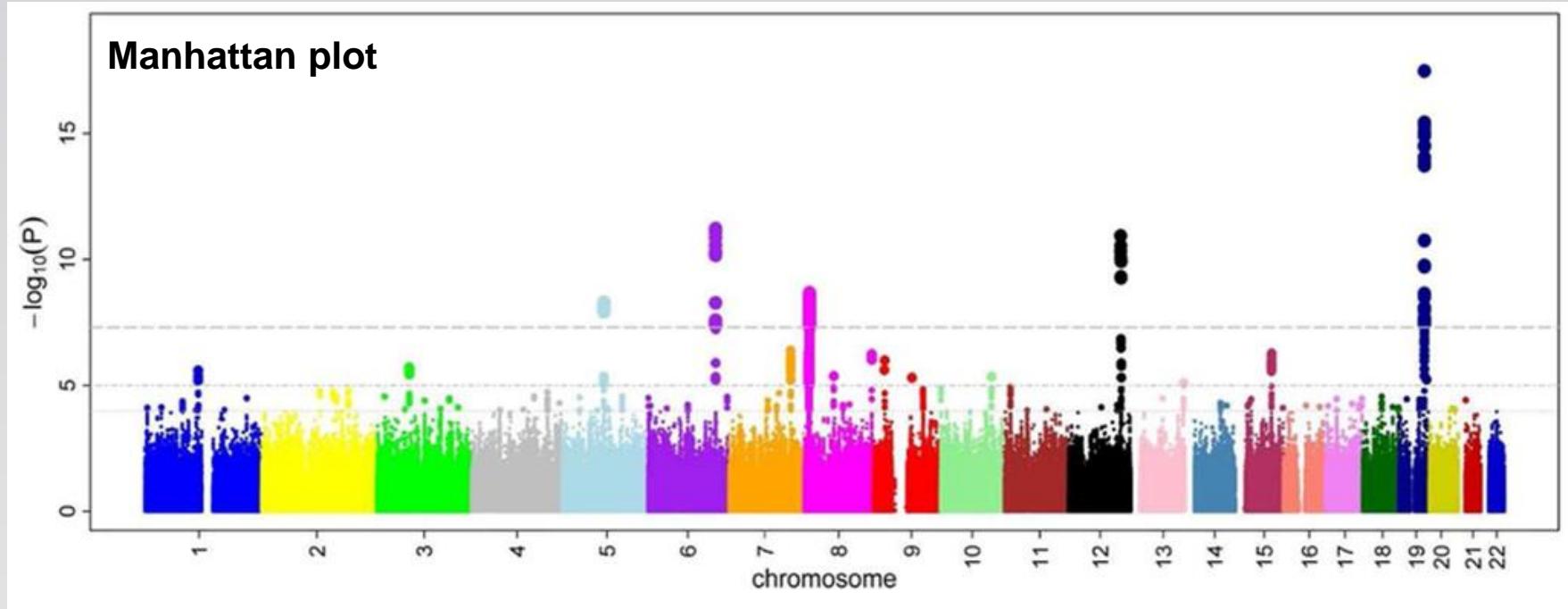
**Just do it !**



Center for Computational  
and Theoretical Biology



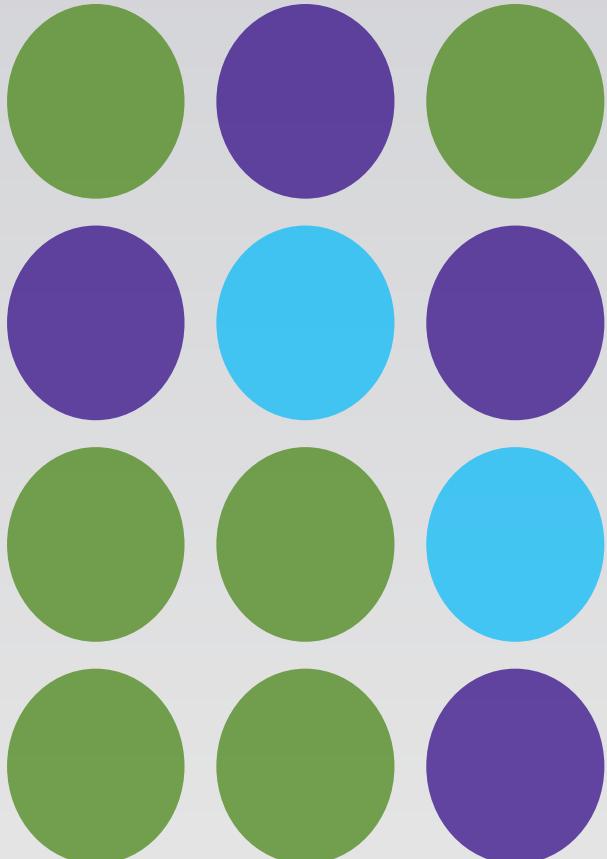
# The outcome



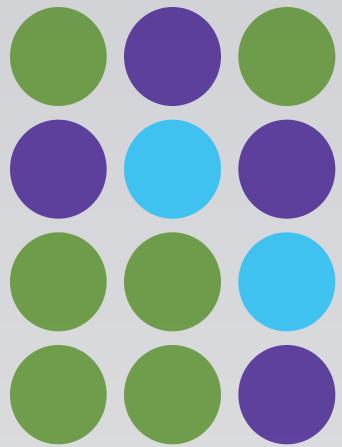
# Introduction: GWAS



Center for Computational  
and Theoretical Biology



# Introduction: GWAS



•••AGCCTG - - - TG**C**ACTAAGA**C**T•••  
•••AGCCTG - - - TG**C**ACTAAGA**C**T•••  
•••AGCCTG - - - TG**C**ACTAAGA**G**T•••  
•••AGCCTG - - - TG**C**ACTAAGA**C**T•••  
•••AGCCTG**A****G****T****T**G**C**ACTAAGA**G**T•••  
•••AGCCTG**A****G****T****T**G**C**ACTAAGA**G**T•••  
•••AGCCTG**A****G****T****T**G**T**ACTAAGA**C**T•••  
•••AGCCTG**A****G****T****T**G**T**ACTAAGA**G**T•••  
•••AGCCTG**A****G****T****T**G**T**ACTAAGA**C**T•••  
•••AGCCTG**A****G****T****T**G**T**ACTAAGA**C**T•••

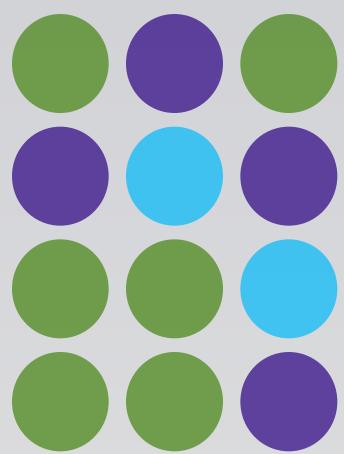
SNPs

Indels

CNVs

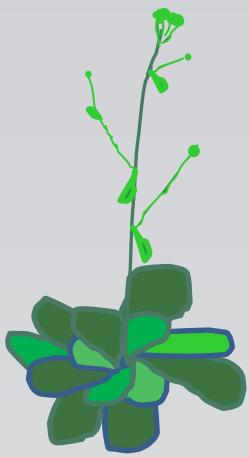
Epigenetic  
markers

# Introduction: GWAS



•••AGCCTG - - - TG**C**ACTAAGA**C**••• **SNPs**  
•••AGCCTG - - - TG**C**ACTAAGA**C**••• **Indels**  
•••AGCCTG - - - TG**C**ACTAAGA**G**••• **CNVs**  
•••AGCCTG - - - TG**C**ACTAAGA**C**•••  
•••AGCCTG**A****G****T****G**T**G**CACTAAGA**G**•••  
•••AGCCTG**A****G****T****G**T**G**CACTAAGA**G**•••  
•••AGCCTG**A****G****T****G**T**G**TACTAAGA**C**•••  
•••AGCCTG**A****G****T****G**T**G**TACTAAGA**G**•••  
•••AGCCTG**A****G****T****G**T**G**TACTAAGA**C**•••  
•••AGCCTG**A****G****T****G**T**G**TACTAAGA**G**•••  
•••AGCCTG**A****G****T****G**T**G**TACTAAGA**C**•••  
•••AGCCTG**A****G****T****G**T**G**TACTAAGA**C**••• **Epigenetic  
markers**

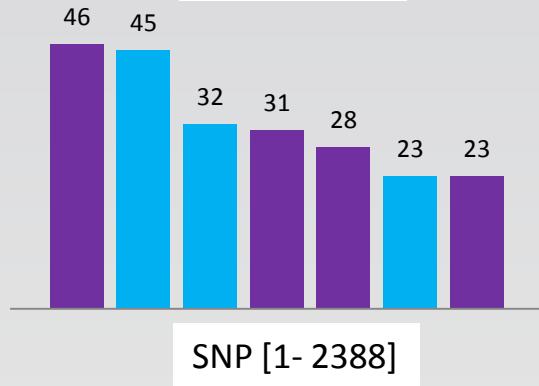
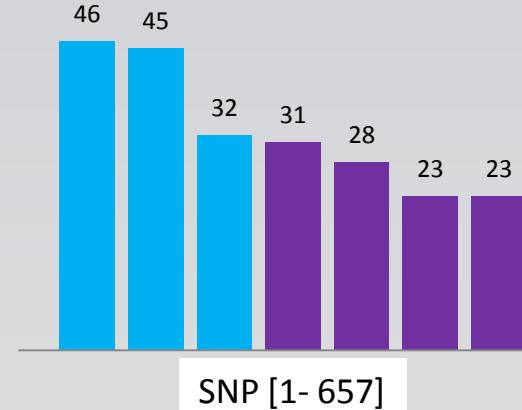
# Introduction: GWAS (linear model)



ecotype	phenotype
6909	23
6911	32
7344	28
7436	23
8233	45
8234	46
8673	31

SNP [1- 657]	SNP [1- 2388]
0	0
1	1
0	0
0	1
1	1
1	0
0	0

■ 0 Allele ■ 1 Allele



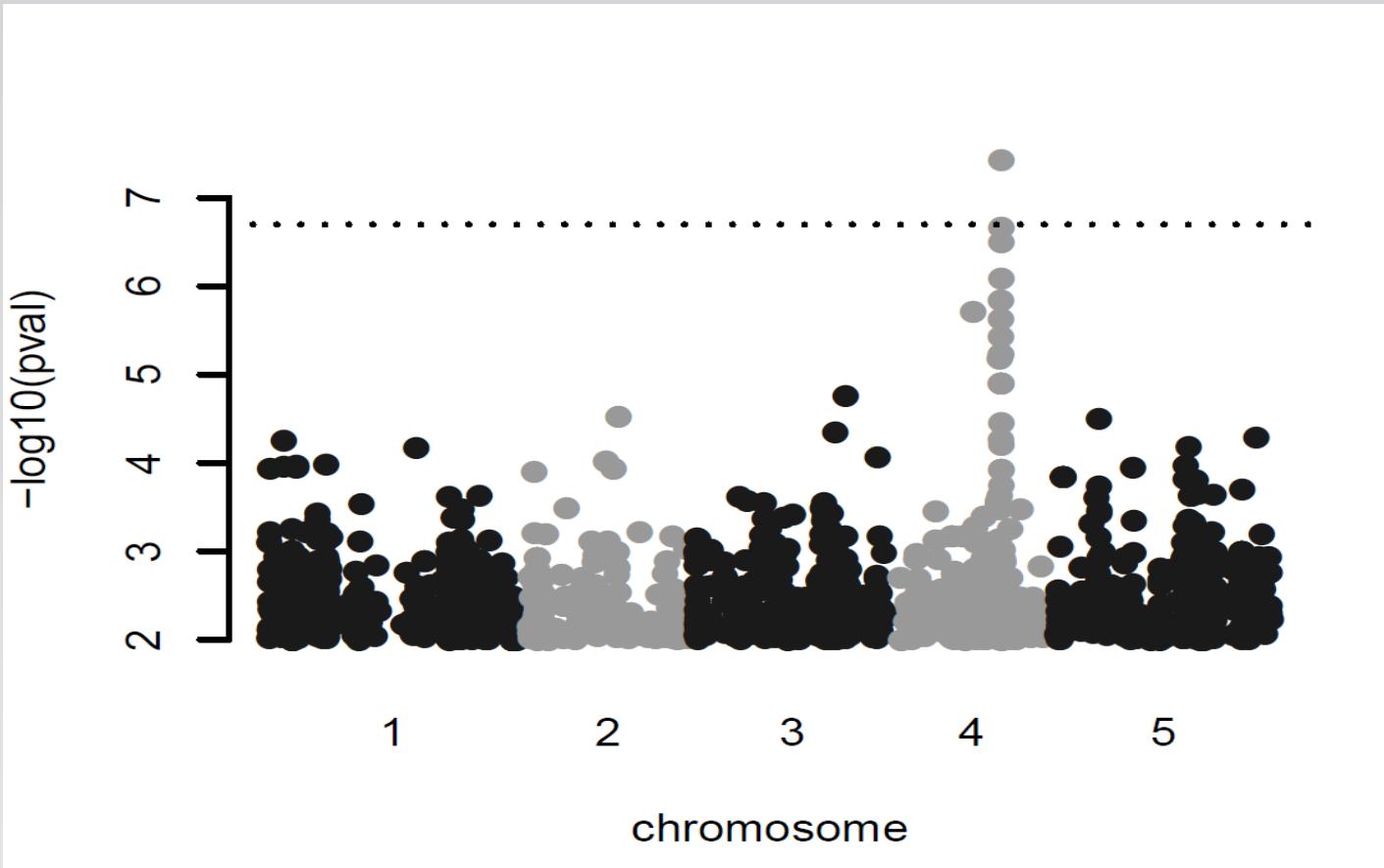
$$\begin{aligned} H_0 : Y &= \beta_0 + \varepsilon \\ H_1 : Y &= \beta_0 + \beta_1 G + \varepsilon \end{aligned}$$

Test for each SNP, if including the SNP in the model can explain the phenotype

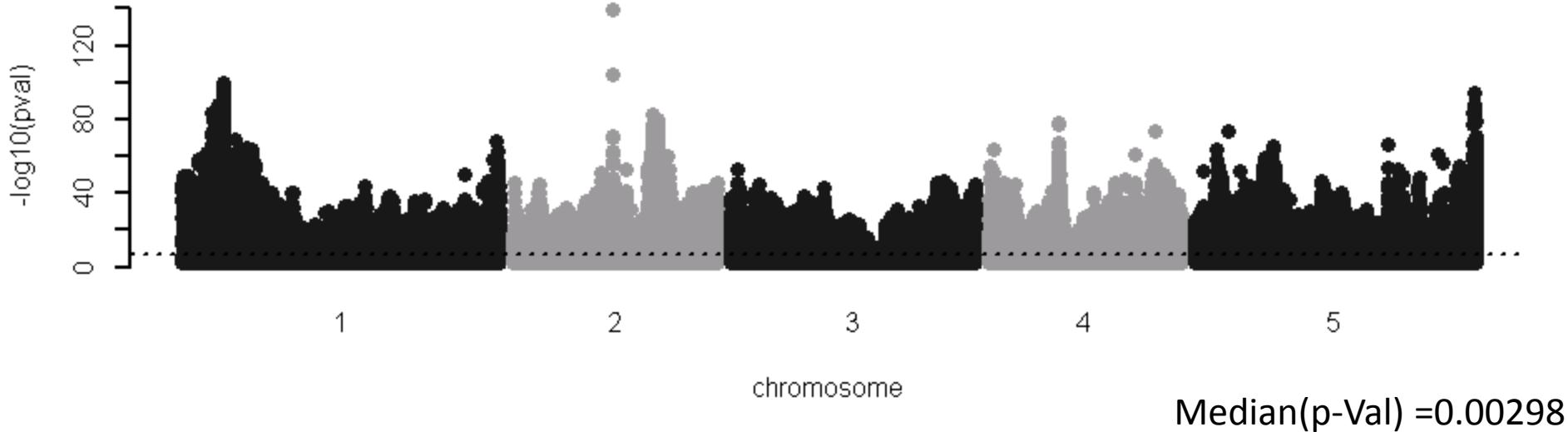
# Introduction: GWAS (linear model)



Center for Computational  
and Theoretical Biology

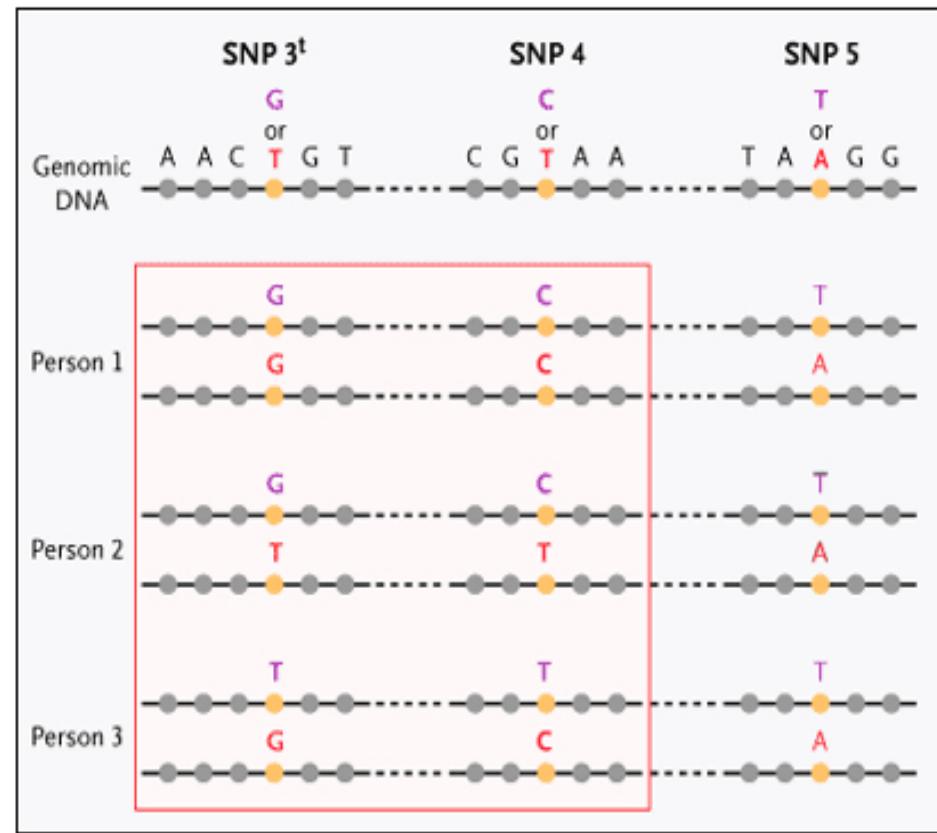
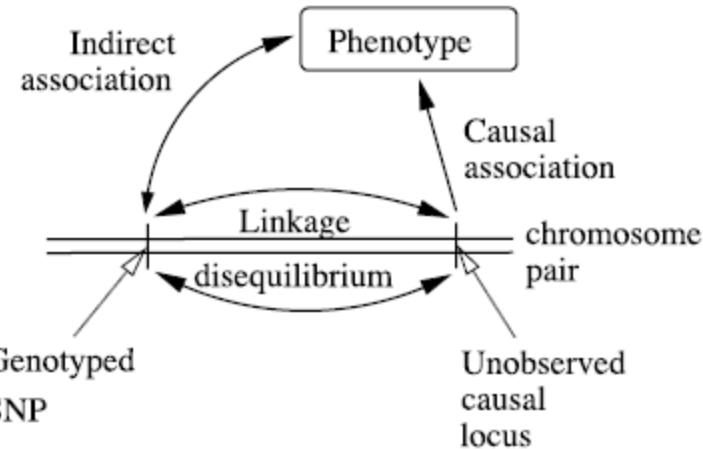


# GWAS (linear model) for 925 accessions on flowering time



**~30 % of the genome is causative for flowering time ?!?**

# Introduction: LD

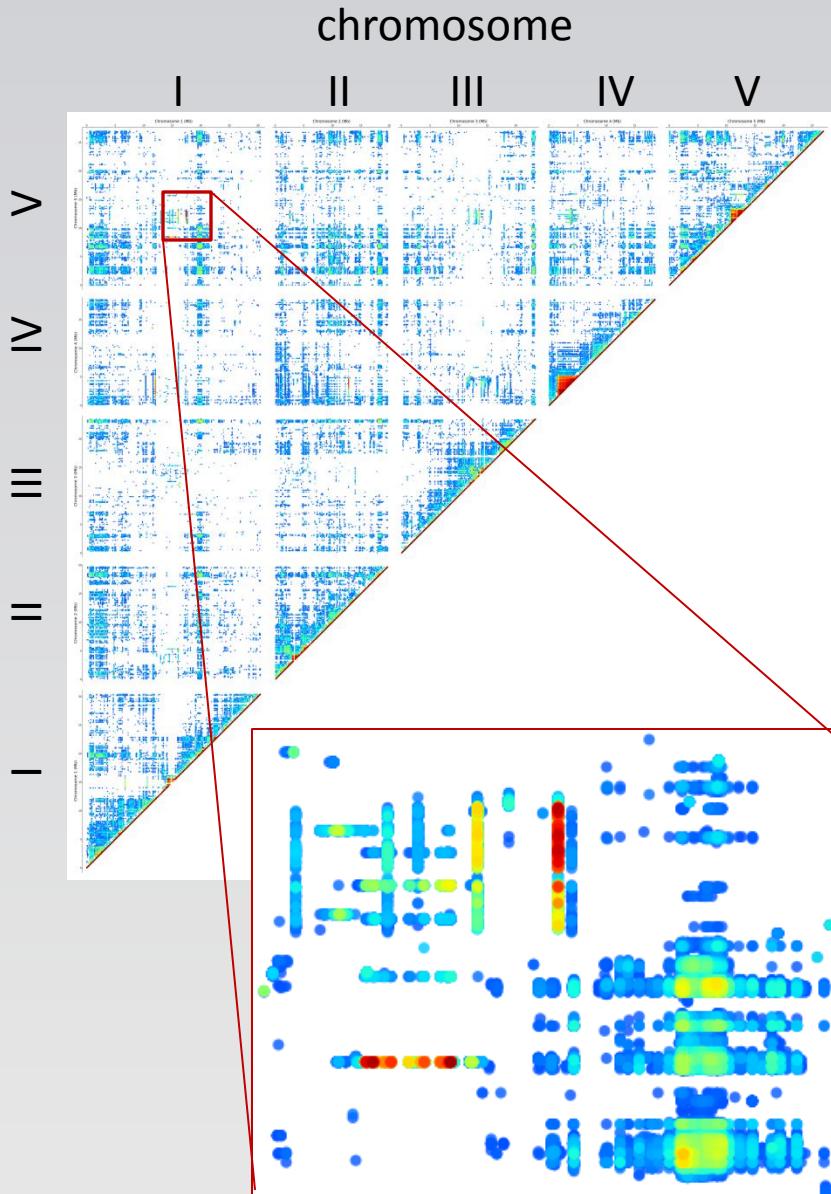


Due to Linkage disequilibrium one SNP can serve as a proxy for others, but also 'rent' power from them

# Pair wise correlation of SNPs across t



Center for Computational  
and Theoretical Biology



# Linkage Disequilibrium

$$p_1(\text{●}) = 0.4$$

$$q_1(\text{○}) = 0.2$$

$$\chi_{11}(\text{●○}) = 0.2$$

$$D = \chi_{11} - p_1 q_1$$

$$D = 0.12$$

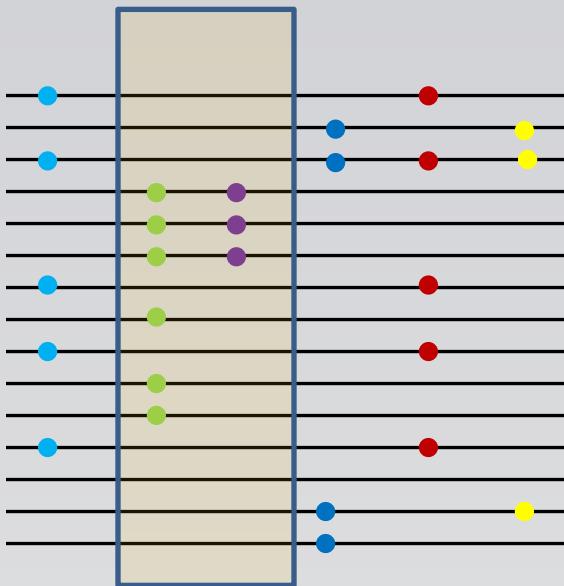
$$D' = D/D_{\max}$$

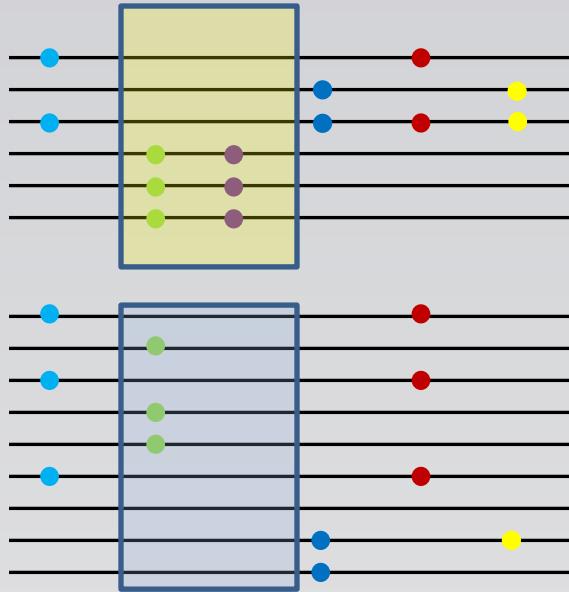
$$D_{\max} = \begin{cases} \min(p_1 q_1, p_2 q_2) & \text{when } D < 0 \\ \min(p_1 q_2, p_2 q_1) & \text{when } D > 0 \end{cases}$$

$$D' = 1$$

$$r^2 = D^2 / p_1 q_1 p_2 q_2$$

$$r^2 = 0.375$$





Subpopulation 1

Subpopulation 2

**Creates LD**

**Selection  
Admixture  
Drift**

**Breaks down LD**

**Recombination  
gene conversion**

**LD between two SNPs can differ in different Subpopulations**

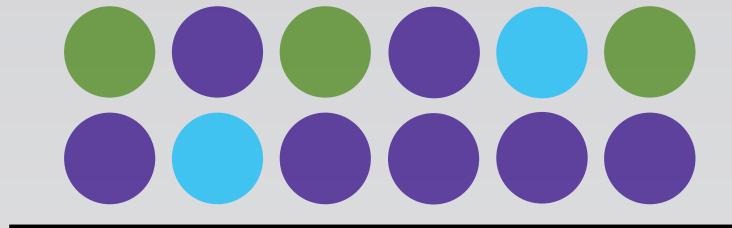
# Population structure



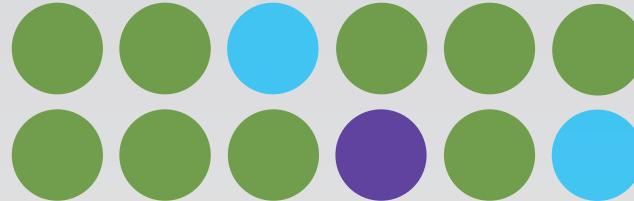
Center for Computational  
and Theoretical Biology

Confounding due to population structure may arise if it correlates with the trait in question

Sub-population 1



Sub-population 2

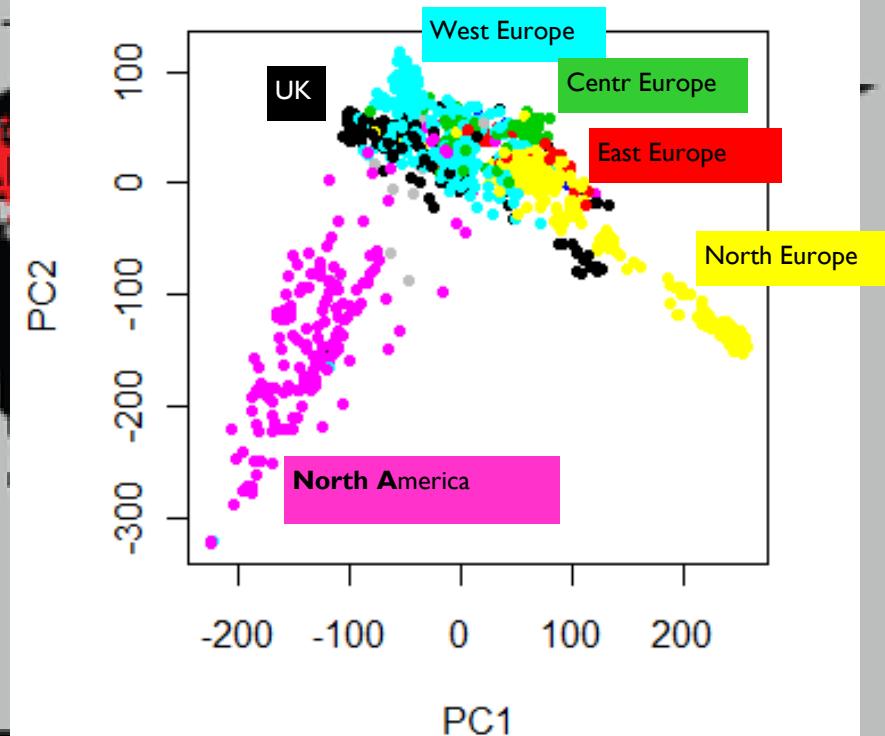


The penetrance of Cystic fibrosis is ~ 1:2.000 for Europeans, for Asians the penetrance is ~ 1:90.000.

# Population structure in *Arabidopsis thaliana*



# Population structure in *Arabidopsis thaliana*

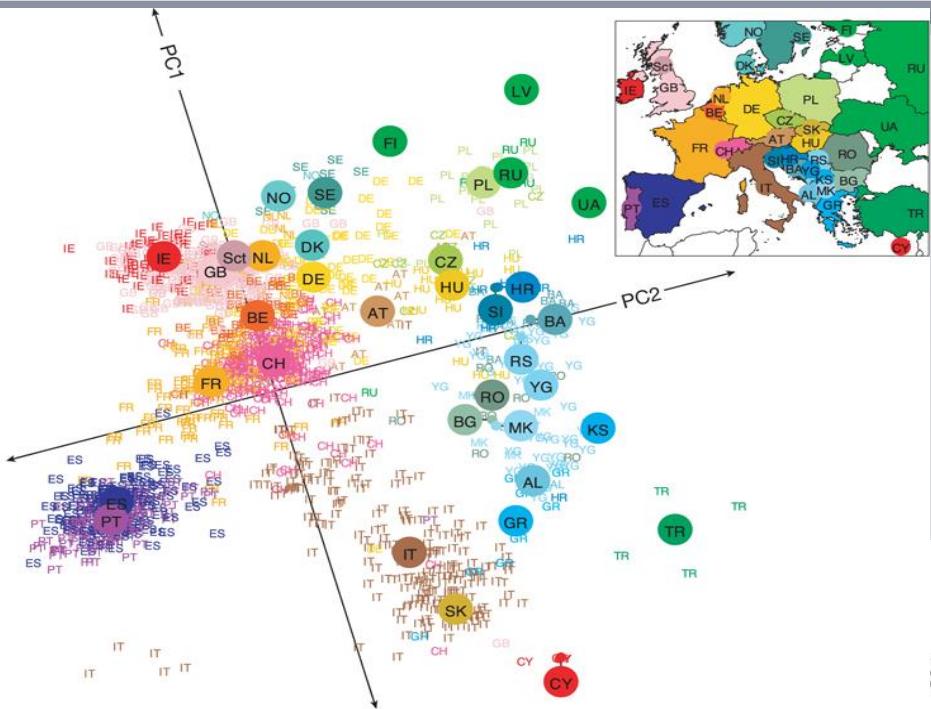


Closely related individuals share not only the causative variants, but have a lot of non-causative SNPs in LD

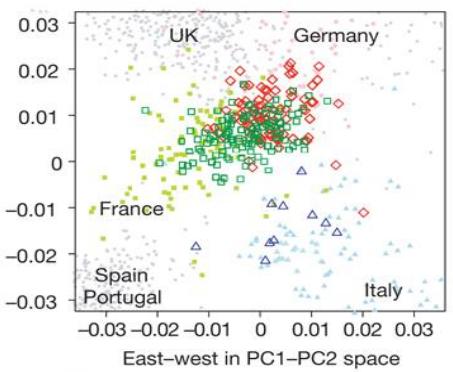
# Structure in human data



Center for Computational  
and Theoretical Biology

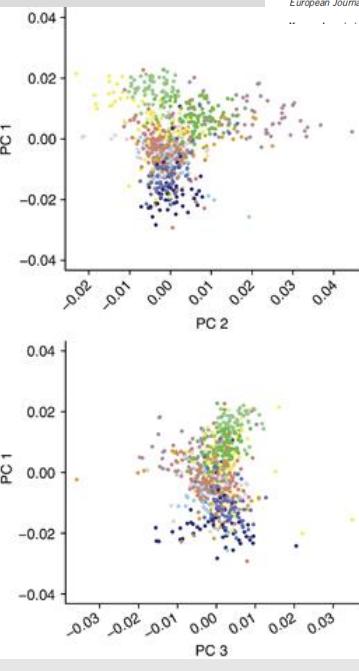
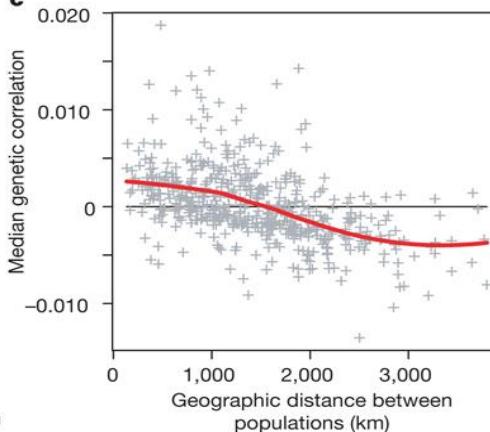


b



French-speaking Swiss      French  
German-speaking Swiss      German  
Italian-speaking Swiss      Italian

c



EJHG Open

European Journal of Human Genetics (2013), 1–7  
© 2013 Macmillan Publishers Limited. All rights reserved 1018-4813/13  
[www.nature.com/ejhg](http://www.nature.com/ejhg)

## ARTICLE

### The Genome of the Netherlands: design, and project goals

Dorret I Boomsma<sup>a,1,2\*</sup>, Cisca Wijmenga<sup>a,2,22</sup>, Eline P Slagboom<sup>3,22</sup>, Morris A Svertz<sup>2,22</sup>, Lennart C Karssen<sup>4</sup>, Abdel Abdellaoui<sup>1</sup>, Kai Ye<sup>5</sup>, Victor Gurfey<sup>3,6</sup>, Martijn Vermaat<sup>7,8,9</sup>, Freek van Dijk<sup>2</sup>, Laurent C Franciolli<sup>10</sup>, Jouke Jan Hottenga<sup>11</sup>, Jeroen FJ Laro<sup>7,8,9</sup>, Qibin Li<sup>11</sup>, Hongzhi Cao<sup>11</sup>, Ruoyan Chen<sup>11</sup>, Yuanding Du<sup>11</sup>, Ning Li<sup>12</sup>, Sujei Gao<sup>12</sup>, Jessica van Setten<sup>10</sup>, Androniki Melouka<sup>10</sup>, Sara L Pula<sup>10</sup>, Jayne Y Hehir-Kwa<sup>15</sup>, Marian Beekman<sup>16</sup>, Clara C Elbers<sup>16</sup>, Heerlijck Byelaas<sup>16</sup>, Anton JM de Craen<sup>16</sup>, Patrick Deelen<sup>16</sup>, Martin Dijksstra<sup>2</sup>, Johan T den Dunnen<sup>8,9</sup>, Peter de Knijff<sup>8,9</sup>, Jeanine Houwing-Duistermaat<sup>17</sup>, Vyacheslav Koval<sup>18</sup>, Karol Estrada<sup>18</sup>, Albert Hofman<sup>4</sup>, Alexandros Kanterakis<sup>2</sup>, David van Enckevort<sup>2</sup>, Halligang Mai<sup>7</sup>, Mathijs Kattenberg<sup>1</sup>, Elisabeth M van Leeuwen<sup>4</sup>, Pieter JT Neerinckx<sup>2</sup>, Ben Oostera<sup>19</sup>, Fernanndo Rivadeneira<sup>18</sup>, Eka HD Suchiman<sup>3</sup>, Andre G Uitdeelinden<sup>18</sup>, Gonnieke Willemsen<sup>1</sup>, Bruce H Wolfenbuttel<sup>20</sup>, Jun Wang<sup>11,13,14,22</sup>, Paul IW de Bakker<sup>10,22</sup>, Gert-Jan van Ommen<sup>21,22</sup> and Cornelia M van Duijn<sup>\*,4,22</sup>

Within the Netherlands a national network of biobanks has been established (Biobanking and Biomolecular Research Infrastructure-Netherlands (BBMRI-NL)) as a national node of the European BBMRI. One of the aims of BBMRI-NL is to enrich biobanks with different types of molecular and phenotype data. Here, we describe the Genome of the Netherlands (GoNL), one of the projects within BBMRI-NL. GoNL is a whole-genome-sequencing project in a representative sample consisting of 250 trios from the Dutch population, including 750 individuals from the Dutch Biobank and 1750 individuals from the Dutch population. The parent-offspring trios include adult individuals ranging in age from 19 to 87 years (mean = 53 years; SD = 16 years) from birth cohorts 1910–1994. Sequencing was done on blood-derived DNA from uncultured cells and accomplished coverage was 14–15x. The family-based design represents a unique resource to assess the frequency of regional variants, accurately reconstruct haplotypes by family-based phasing, characterize short indels and complex structural variants, and establish the rate of *de novo* mutational events. GoNL will also serve as a reference panel for imputation in the available genome-wide association studies in Dutch and other cohorts and for population-specific variants. GoNL will create a catalog of common genomic variants that are uniquely characterized with respect to their geographic location and a wide range of phenotypes. The resource will be made available to the research and medical community to guide the interpretation of sequencing projects. The present paper summarizes the global characteristics of the project.

European Journal of Human Genetics advance online publication, 29 May 2013; doi:10.1038/ejhg.2013.118



Novembre, Nature 2008

# HOW TO CONTROL FOR STRUCTURE ?



Center for Computational  
and Theoretical Biology

- Genomic control (Devlin & Roeder 1999, Biometrics)
- Structured association (Pritchard et al. 2000, Am.J.Hum.Genet.)
- Principal-components approach (Price et al. 2006, Nature Genet.)
- Mixed-model approach (Yu et al. 2006, Nature Genet.; Kang et al. 2008, Genetics)
- FaST-LMM -Select (Listgarten et al. 2013, Nature Genet.)
- ....

# The mixed model is the only one that works

- Mixed model handles different levels of relatedness generally:

$$Y = X\beta + u + \epsilon, \quad u \sim N(0, \sigma_g K), \quad \epsilon \sim N(0, \sigma_e I)$$

$$\text{LM: } Y = X\beta + \epsilon$$



Introduced by R.A. Fisher (1918) and further developed by Henderson in the 1950s

$$\text{LMM: } \mathbf{y} = \mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{u} + \boldsymbol{\varepsilon}$$

- $\mathbf{y}$  is a  $n \times 1$  vector of the observed phenotype
- $\mu$  is a fixed effect representing the phenotypic mean
- $\mathbf{X}$  is an  $n \times q$  matrix of fixed effects
- $\boldsymbol{\beta}$  is a  $q \times 1$  vector of the coefficients of the fixed effects
- $\mathbf{u}$  is the random effect of the mixed model with  $\text{var}(\mathbf{u}) = \sigma_g \mathbf{K}$
- $\mathbf{K}$  is the  $n \times n$  kinship matrix inferred from genotypes
- $\boldsymbol{\varepsilon}$  is a  $n \times n$  matrix of residual effects with  $\text{var}(\boldsymbol{\varepsilon}) = \sigma_e \mathbf{I}$

The overall phenotypic variance-covariance matrix can be represented as

$$\mathbf{V} = \sigma_g \mathbf{K} + \sigma_e \mathbf{I}$$

Pseudo heritability

$$\hat{h}^2 = \hat{\sigma}_g^2 / (\hat{\sigma}_g^2 + \hat{\sigma}_e^2)$$

$$\text{LMM: } \mathbf{y} = \mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{u} + \boldsymbol{\varepsilon}$$

Yu et al. (2006) : solving the model with best linear unbiased predictor (BLUP) via the Henderson equation (Henderson 1984)

Kang et al. (2008) : direct estimation of the variance components ( $\sigma_g$  and  $\sigma_e$ ) by maximizing the restricted likelihood (EMMA)

$$l_F(\mathbf{y}; \boldsymbol{\beta}, \sigma, \delta) = \frac{1}{2} \left[ -n \log(2\pi\sigma^2) - \log |H| - \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' H^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \quad (2)$$

$$l_R(\mathbf{y}; \sigma, \delta) = l_F(\mathbf{y}; \hat{\boldsymbol{\beta}}, \sigma^2, \delta) + \frac{1}{2} [q \log(2\pi\sigma^2) + \log |X'X| - \log |X'H^{-1}X|] \quad (3)$$

Kang et al. (2010) : fast heuristic approximation that assumes  $\sigma_g/\sigma_e$  is constant (EMMAX)

Zhang et al. (2010): P3D

Lippert et al.(2011): FaST-LM

There are four principal assumptions which justify the use of linear regression :

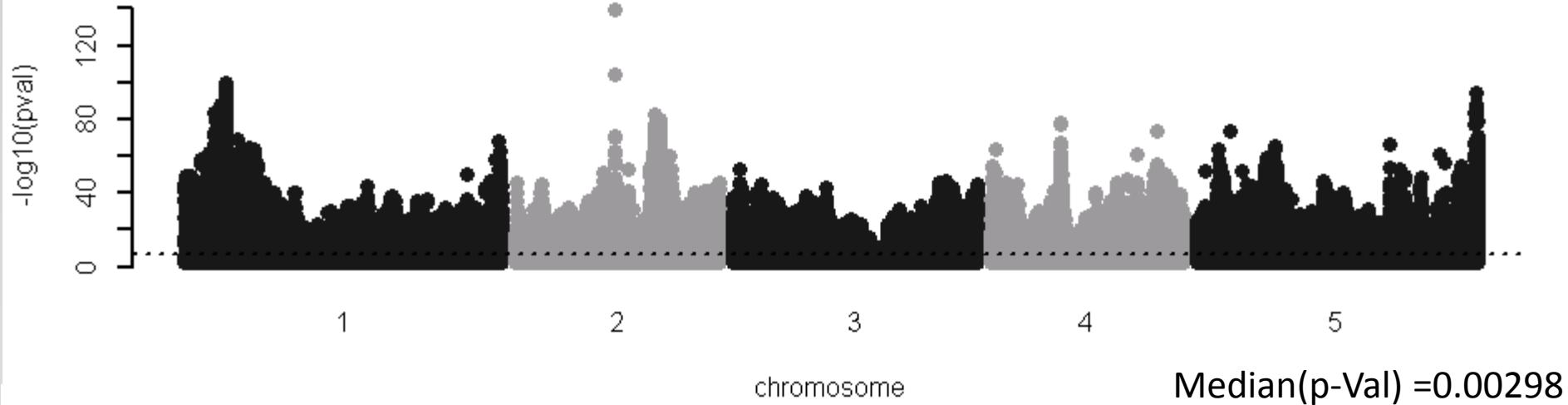
- (i)** linearity of the relationship between dependent and independent variables
- (ii)** independence of the errors
- (iii)** constant variance of the errors
- (iv)** normality of the error distribution

If any of these assumptions is violated model predictions can be wrong, which makes the interpretation of the results inefficient or even misleading.

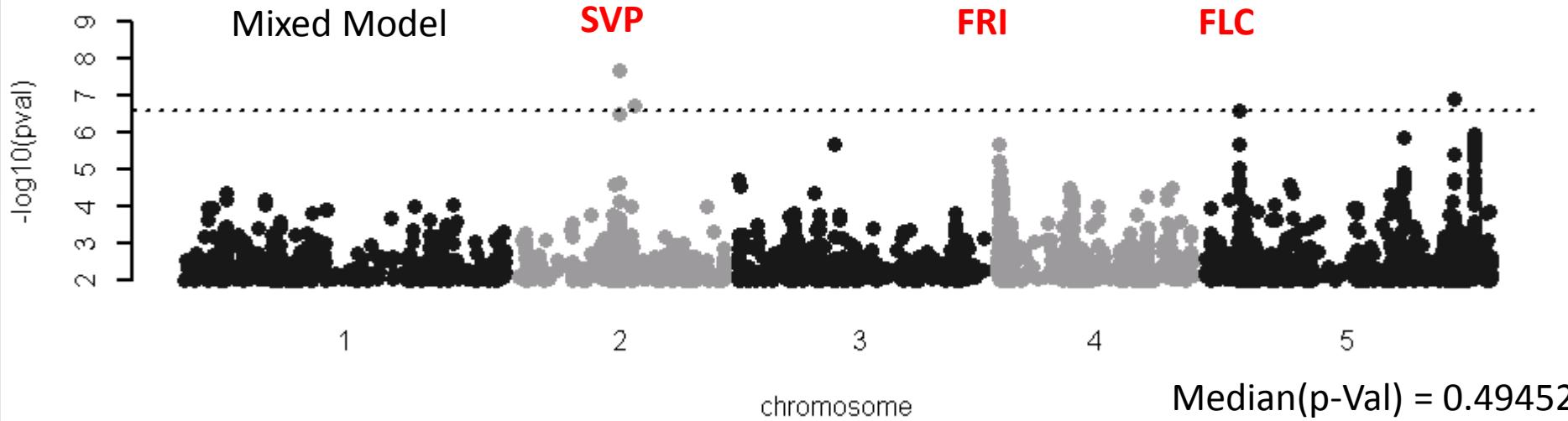
One alternative could be non-parametric models (e.g. Kruskal-Wallis or Anderson-Darling test) or generalized linear models

# It works !

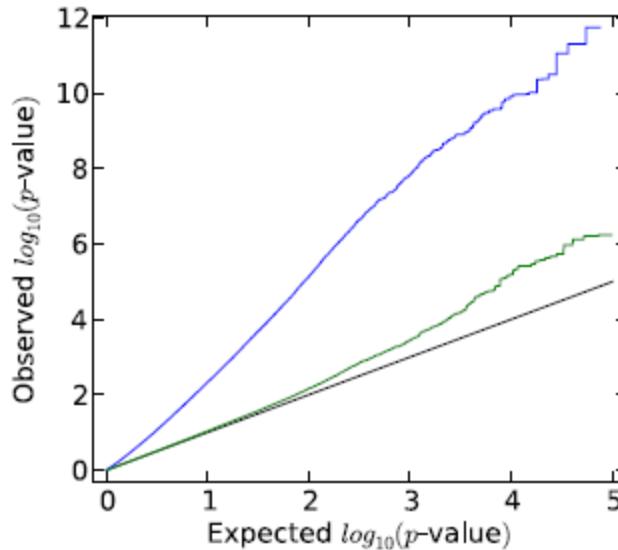
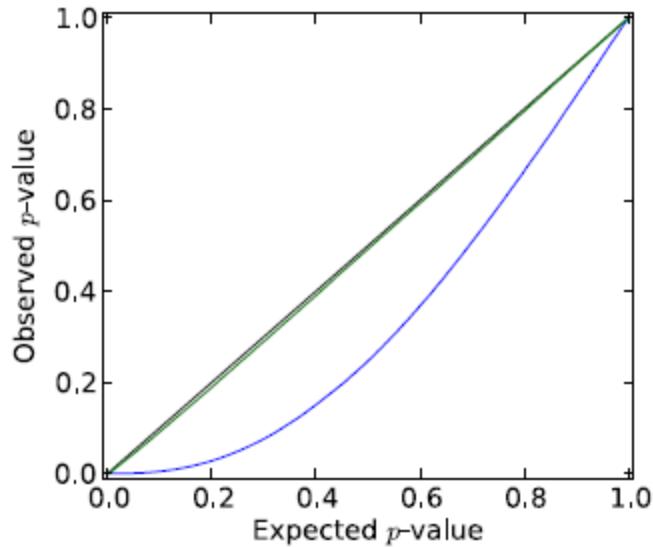
Linear model



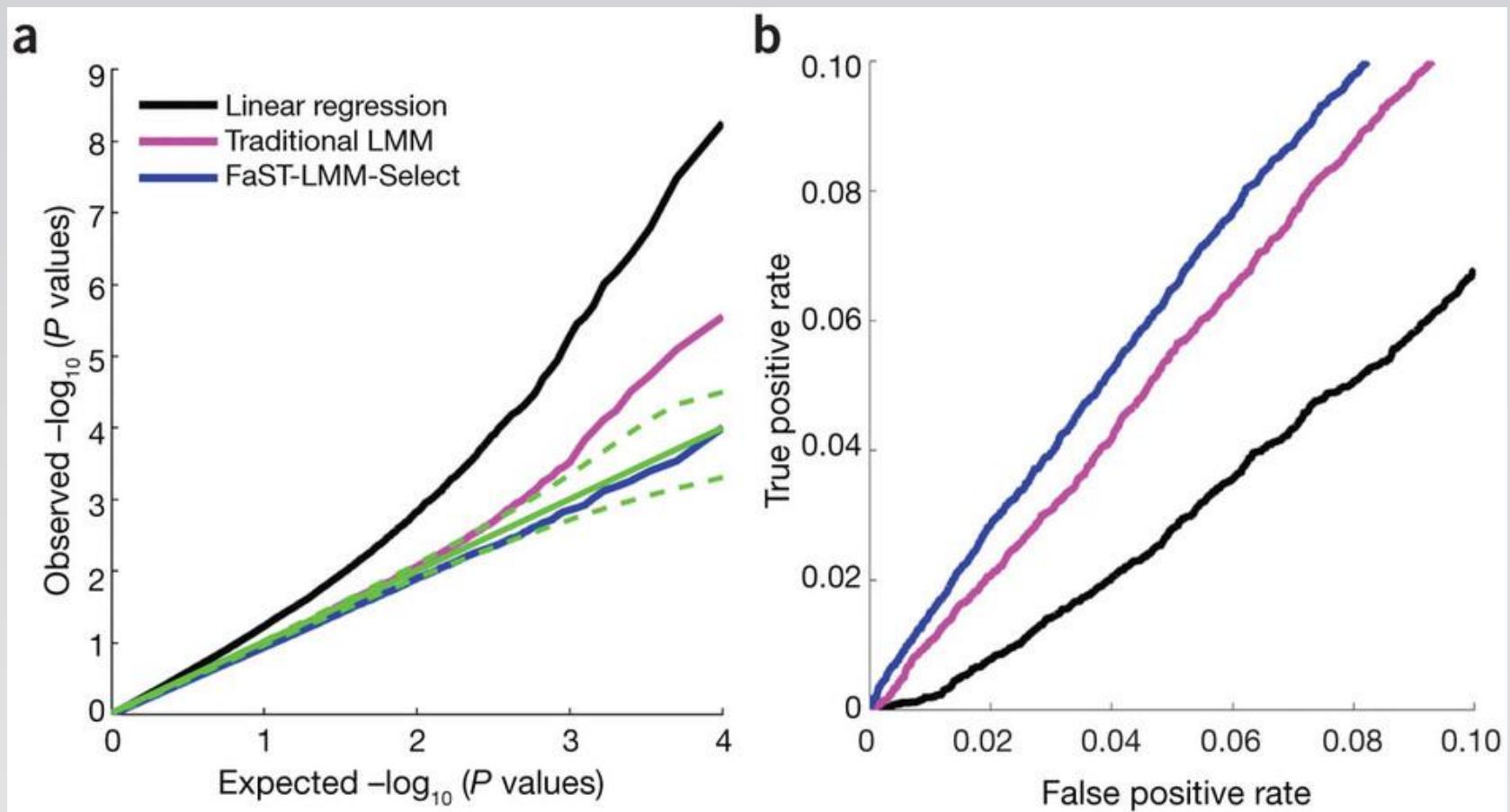
Mixed Model



## qq-plots



- We still have false positives — and have probably introduced false negatives...



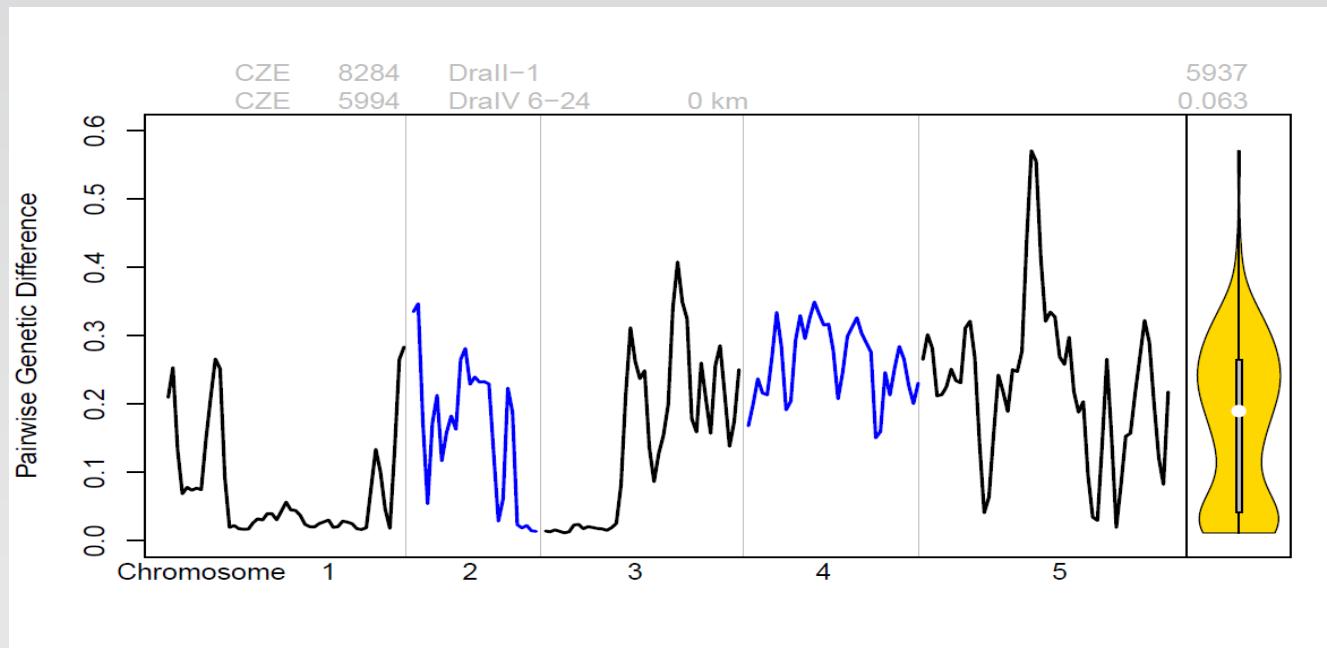
# Kinship matrix



**Kinship matrix could be calculated as IBD – matrix, but this assumes the contribution of each loci to the trait is independent of its allele frequency.**

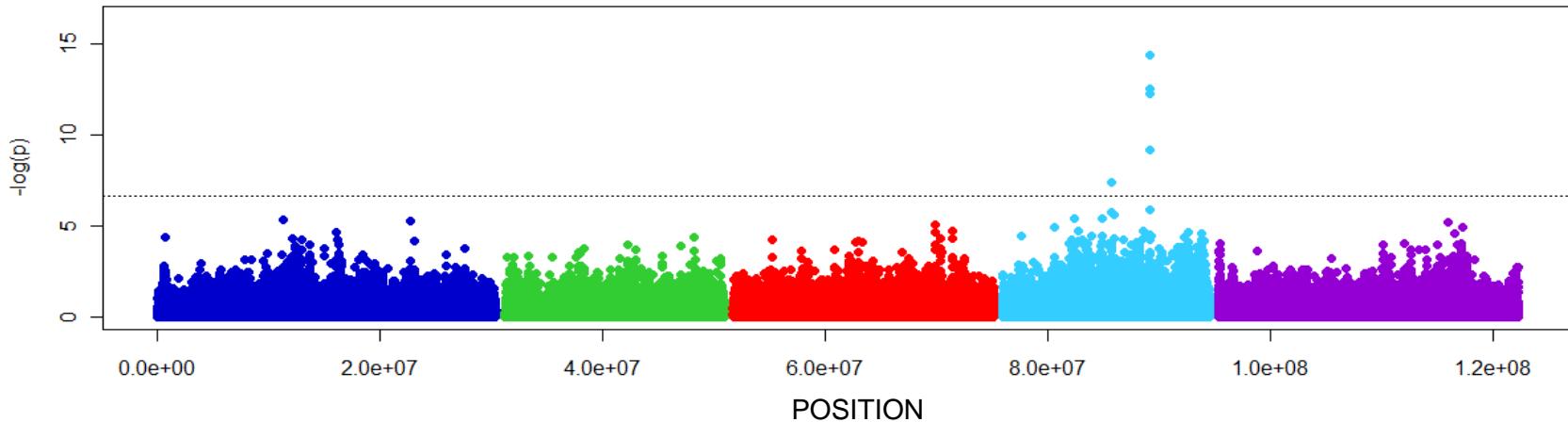
**A nother possibility is the calculation of an IBS matrix, estimated using all SNPs (infinitesimal model)**

**Additionally methods are proposed to select only subsets of SNPs to calculate the kinship matrix (see e.g. Listgar ten et al. 2013)**

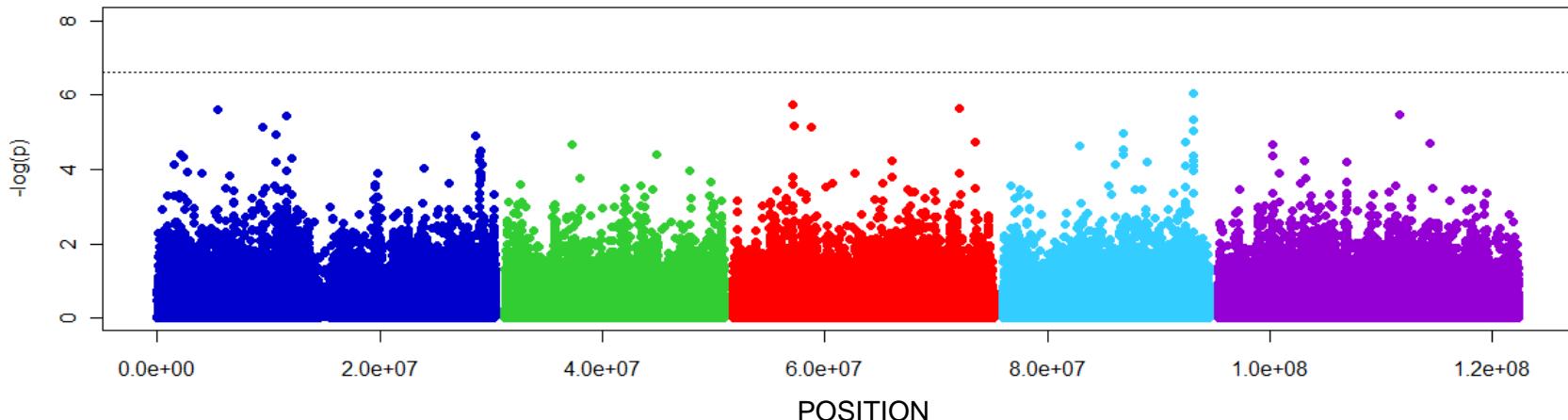


# Examples of real traits

A: GWA analysis of hypersensitive response to bacterial elicitor



B: GWA analysis of germination on MS medium

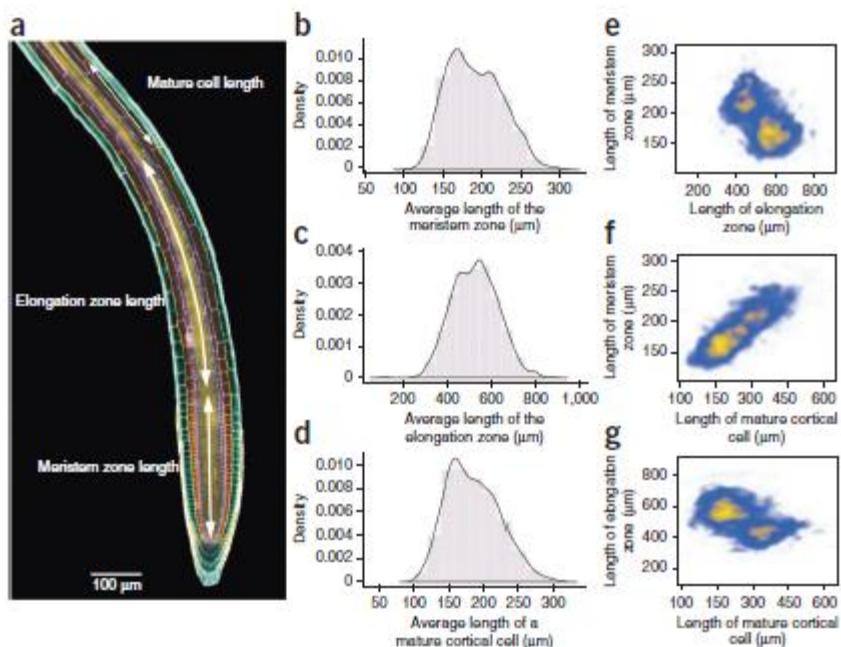


- Trait architecture (simple vs. complex traits, how many causative Alleles, allelic heterogeneity, epistasis etc....)
- Phenotypic distribution of the trait (normal ?, outliers ?)
- Correlation between the phenotype and non-genetic parameters

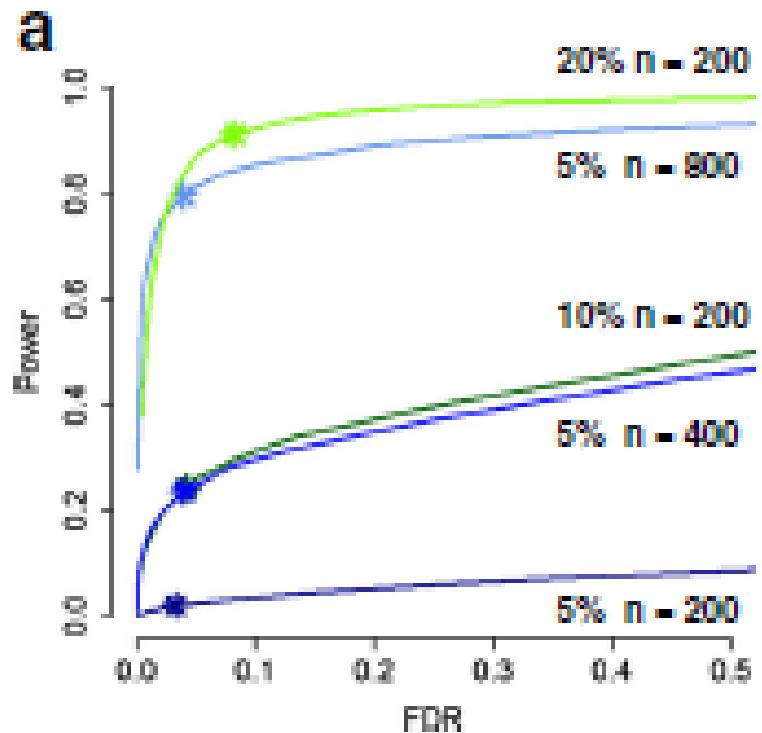
## LETTERS

nature  
genetics

## Genome-wide association study using cellular traits identifies a new regulator of root development in *Arabidopsis*

Mónica Meijón<sup>1,2</sup>, Santosh B Satbhai<sup>1</sup>, Takashi Tsuchimatsu<sup>1</sup> & Wolfgang Busch<sup>1</sup>With the increased availability of high-resolution sequence data, the identification of cellular features<sup>1–4</sup> and a large number of high-density genotypes<sup>5–7</sup>

# Power / FDR



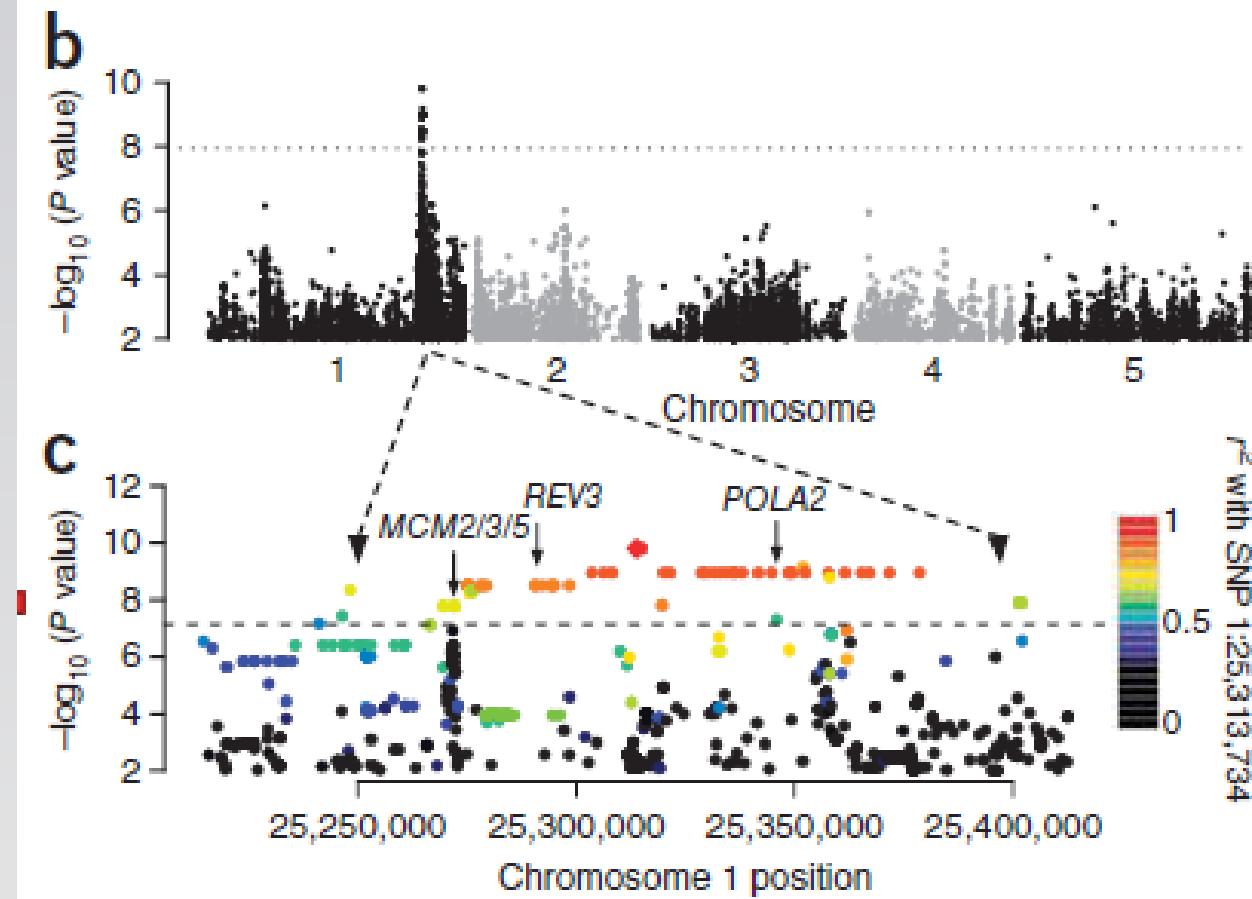
# Another problem

- We have seen that it is usually difficult to decide which peaks are significant
- In addition, peaks are often complex, making it difficult to pinpoint causal sites

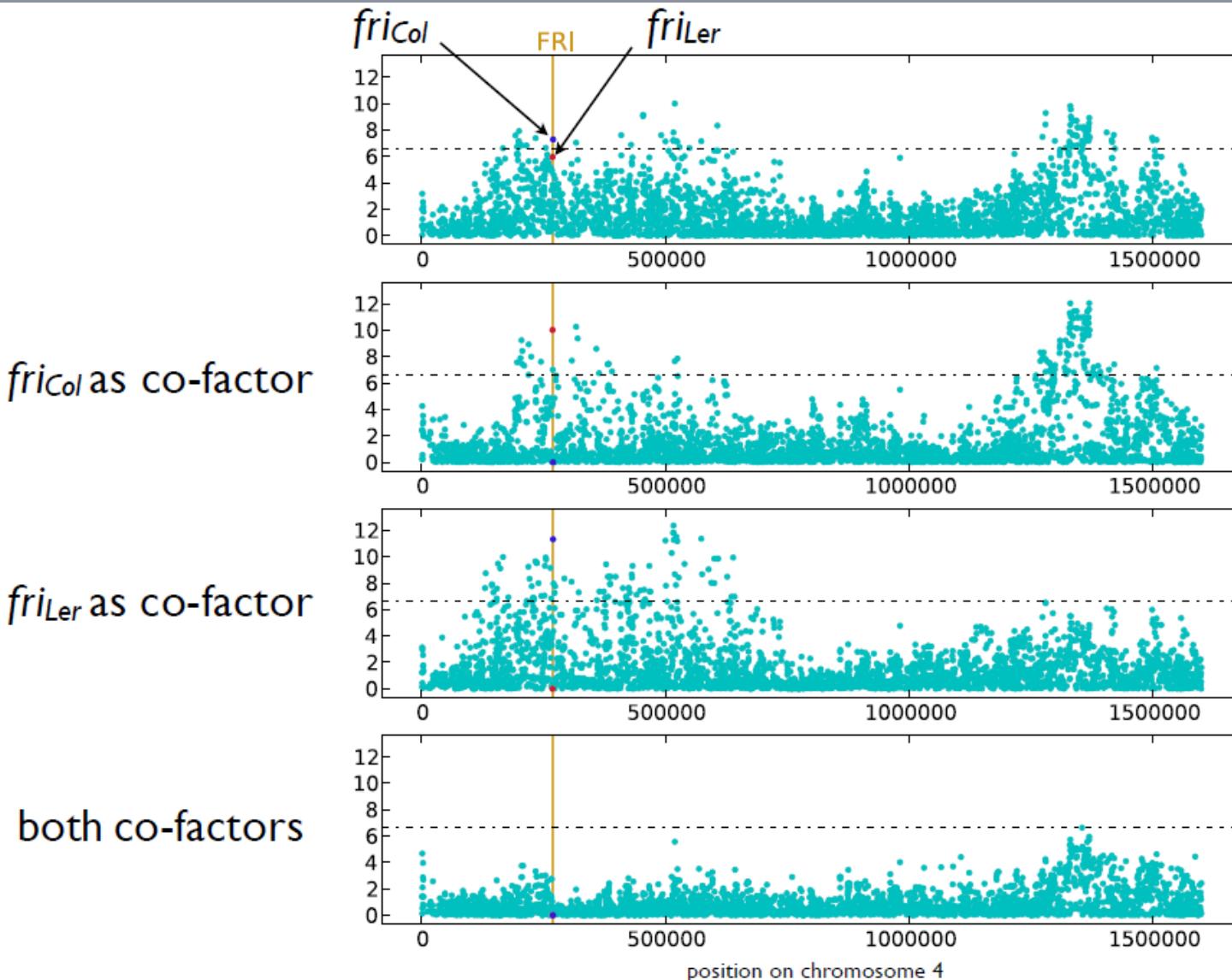
# GWAS peaks can be complex



Center for Computational  
and Theoretical Biology



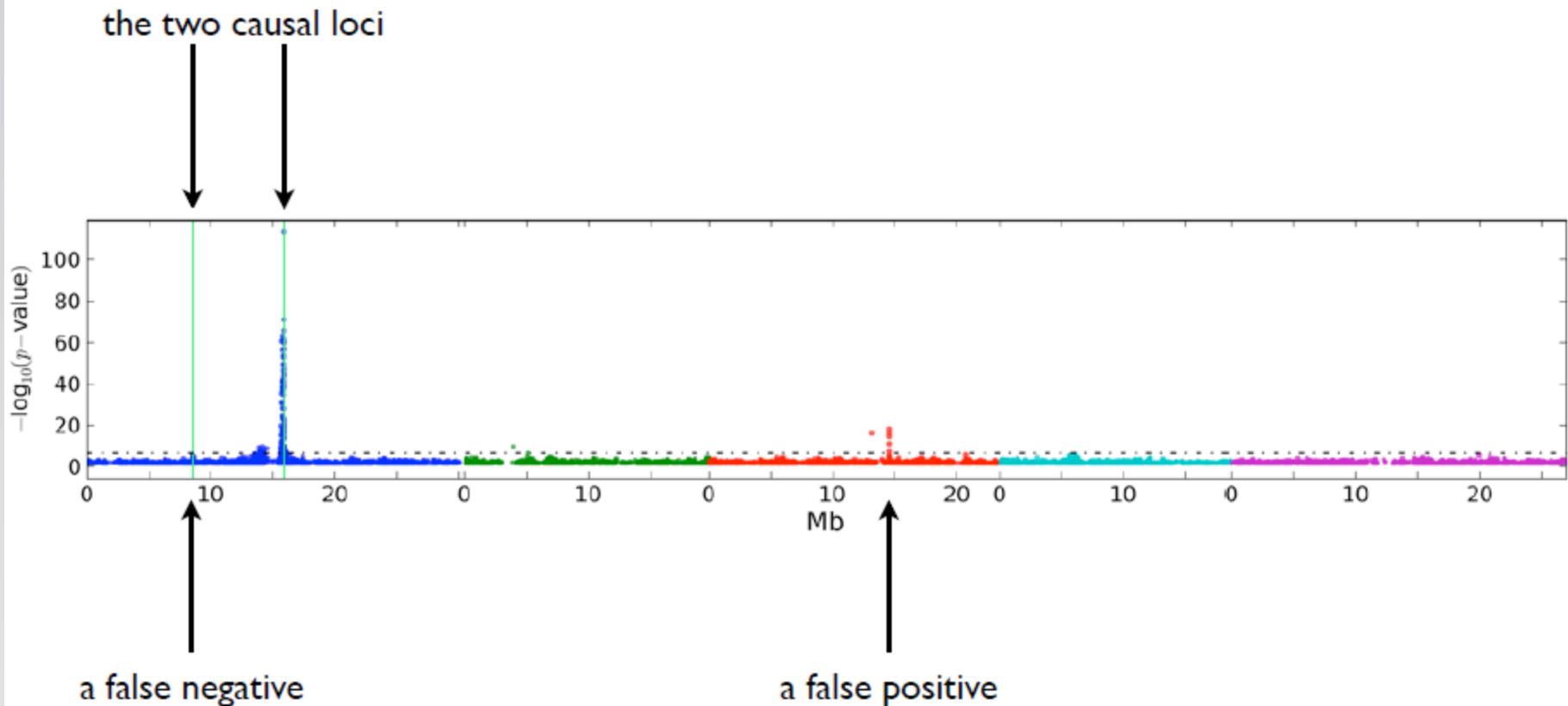
# Example of a complex GWAS peak

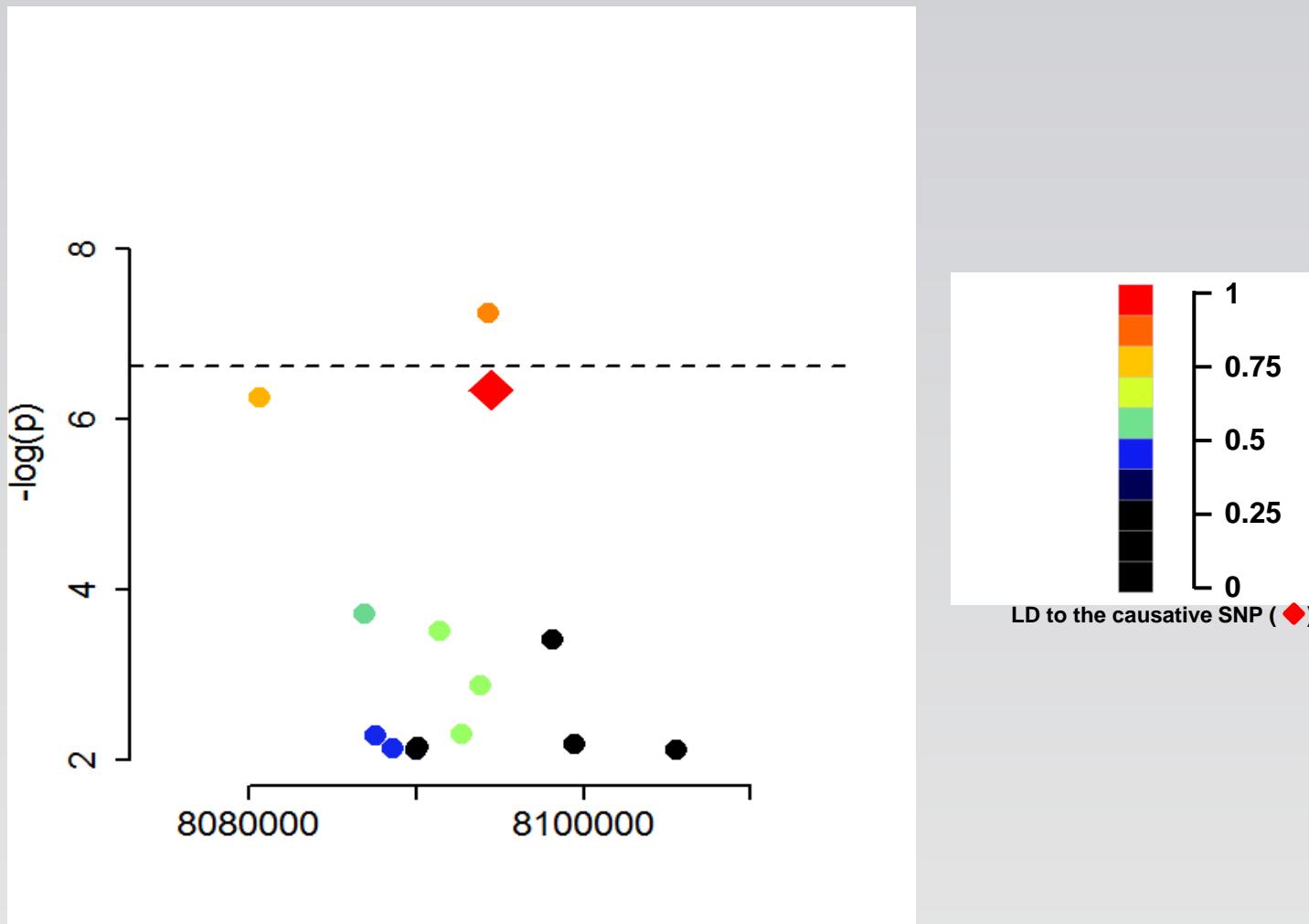


# Conditions under which GWAS will be positively misleading

- Correlation between causal factors and unlinked non-causal markers
- More than one causal factor
- Epistasis

(Platt et al., *Genetics* 2010)





- A non-causative SNP can show stronger association with a given phenotype than the actual causative SNP.

- Population structure *per se* is neither necessary or sufficient for GWAS to go wrong
- The fundamental problem is *not* population structure, it is model mis-specification — modeling something multifactorial using a single locus
- Population structure greatly aggravates the problem, however

# COMMENT

## The nature of confounding in genome-wide association studies

Bjarni J. Vilhjálmsson<sup>1,2</sup> and Magnus Nordborg<sup>3,4</sup>

The authors argue that population structure per se is not a problem in genome-wide association studies—the true sources are the environment and the genetic background, and the latter is greatly underappreciated. They conclude that mixed models effectively address this issue.

“ population structure is not the fundamental source of the problem, and removing it is not the solution ”

Thanks to dramatically decreasing genotyping sequencing costs, genome-wide association studies (GWASs) are becoming the default method for studying the genetics of natural variation. The increasing number and diversity of GWASs will require appropriate statistical analysis methods. The most basic problem is assessing the significance of an association in the presence of confounding effects that may cause spurious associations.

The aspect of this problem that has received most attention is the danger of false positives in structured populations. If the study population is a mixture of populations that differ with respect to allele frequencies as well as the trait of interest, spurious correlations may arise.

Korte and Farlow *Plant Methods* 2013, 9:29  
<http://www.plantmethods.com/content/9/1/29>

### REVIEW

### Open Access



## The advantages and limitations of trait analysis with GWAS: a review

Arthur Korte<sup>\*†</sup> and Ashley Farlow<sup>†</sup>

### Abstract

Over the last 10 years, high-density SNP arrays and DNA re-sequencing have illuminated the majority of the genotypic space for a number of organisms, including humans, maize, rice and *Arabidopsis*. For any researcher willing to define and score a phenotype across many individuals, Genome Wide Association Studies (GWAS) present a powerful tool to reconnect this trait back to its underlying genetics. In this review we discuss the biological and statistical considerations that underpin a successful analysis or otherwise. The relevance of biological factors including effect size, sample size, genetic heterogeneity, genomic confounding, linkage disequilibrium and spurious association, and statistical tools to account for these are presented. GWAS can offer a valuable first insight into trait architecture or candidate loci for subsequent validation.

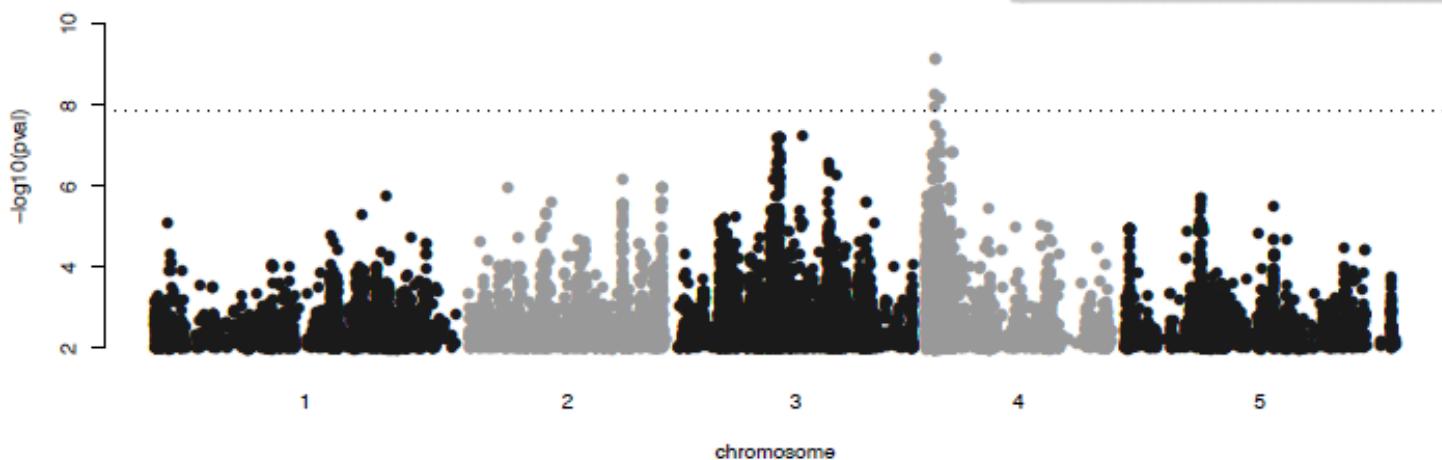
**Keywords:** GWAS, *Arabidopsis*, Mixed model, Effect size, Genetic heterogeneity



## The case of the missing heritability

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. **Brendan Maher** shines a light on six places where the missing loot could be stashed away.

## Slug damage GWAS identifies peak consistent with AOP2/3



# GWAS in *Arabidopsis thaliana*

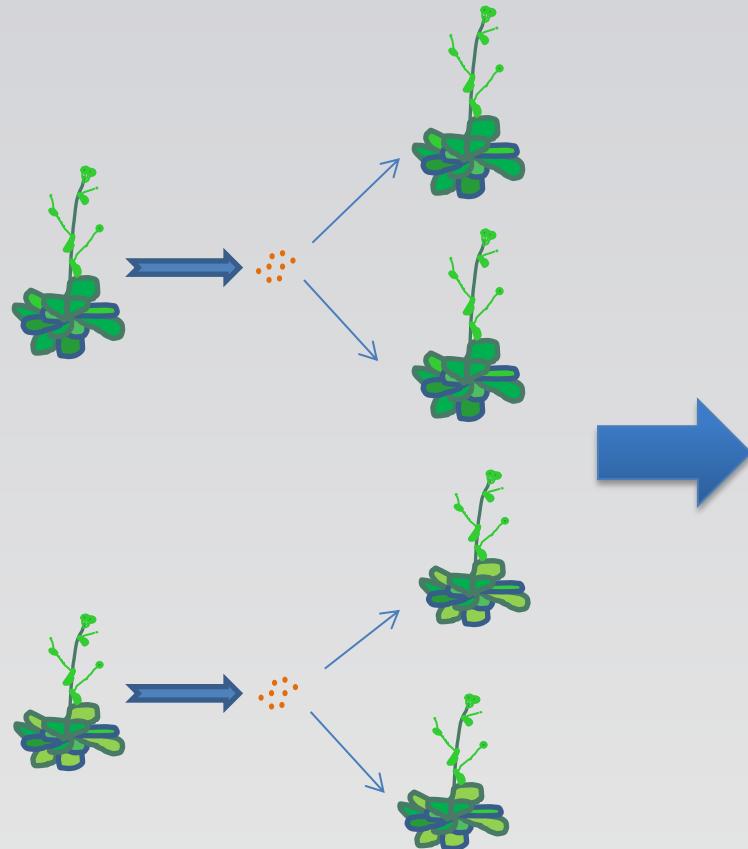


Center for Computational  
and Theoretical Biology

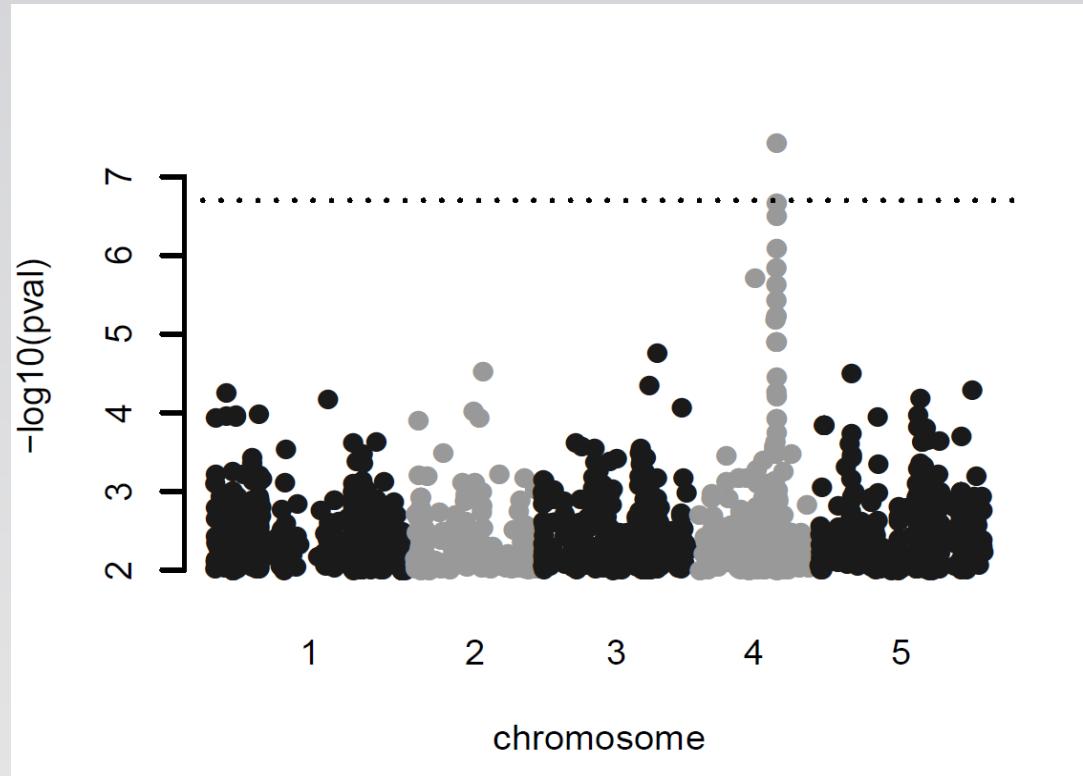


# GWAS in *Arabidopsis thaliana*

## Phenotype



## GWAS



- Robust (replicated) phenotypes in controlled environments
- Follow-up studies are easy

## 1001 Genomes

A Catalog of *Arabidopsis thaliana* Genetic Variation

[Home](#)[Data Providers](#)[Accessions](#)[Tools](#)[Software](#)[Data Center](#)[About](#)

### Welcome to the 1001 Genomes Project

---

## AUTHOR CONTRIBUTIONS

### Project coordination

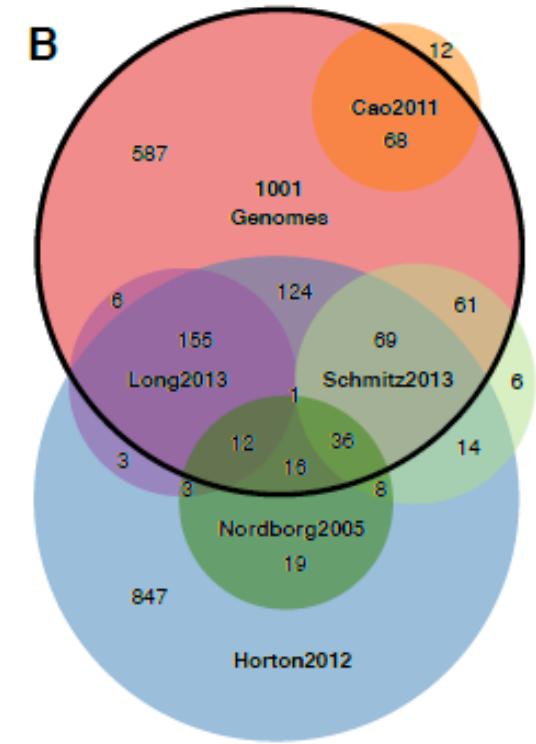
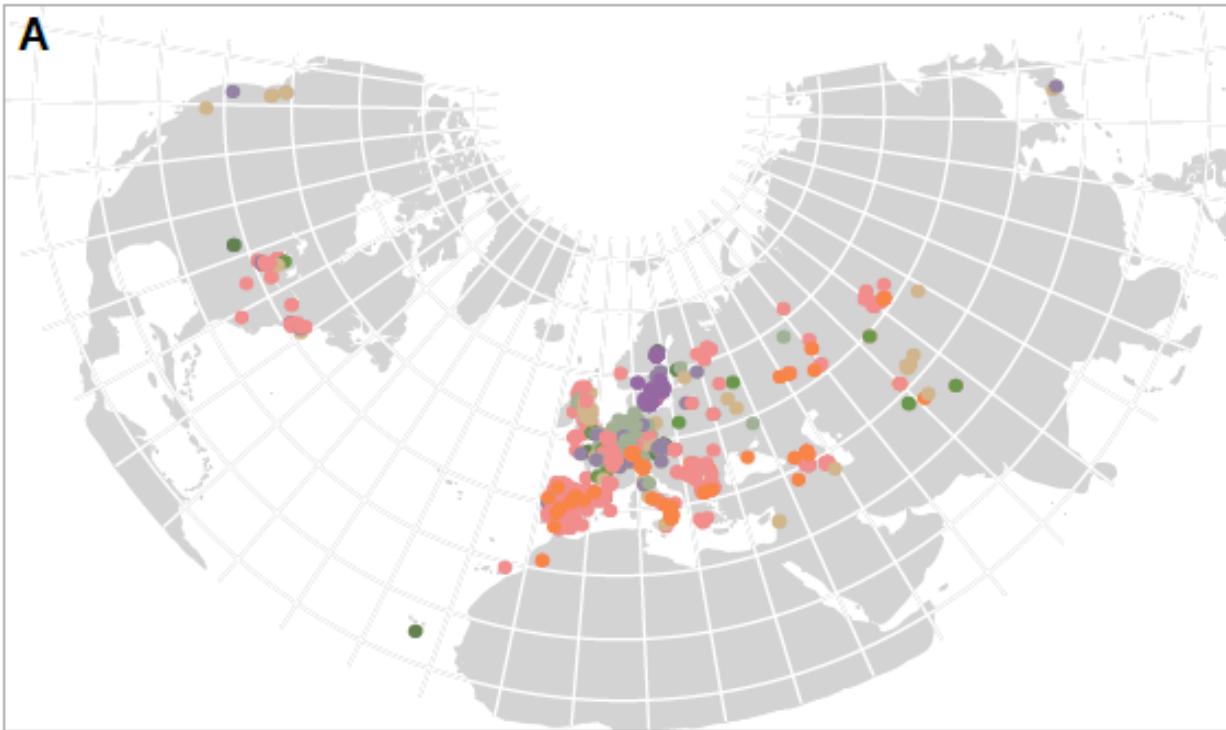
Joy Bergelson	UCHICAGO	jbergels@uchicago.edu
Joe R. Ecker	SALK	ecker@salk.edu
Magnus Nordborg	GMI	magnus.nordborg@gmi.oeaw.ac.at
Mitchell Sudkamp	MONSANTO	mitchell.sudkamp@monsanto.com
Detlef Weigel	MPI	weigel@weigelworld.org

[Go »](#)

80 strains (D. Weigel lab, MPI)  
195 strains (J. Ecker lab, Salk)  
180 strains (M. Nordborg Lab, GMI)

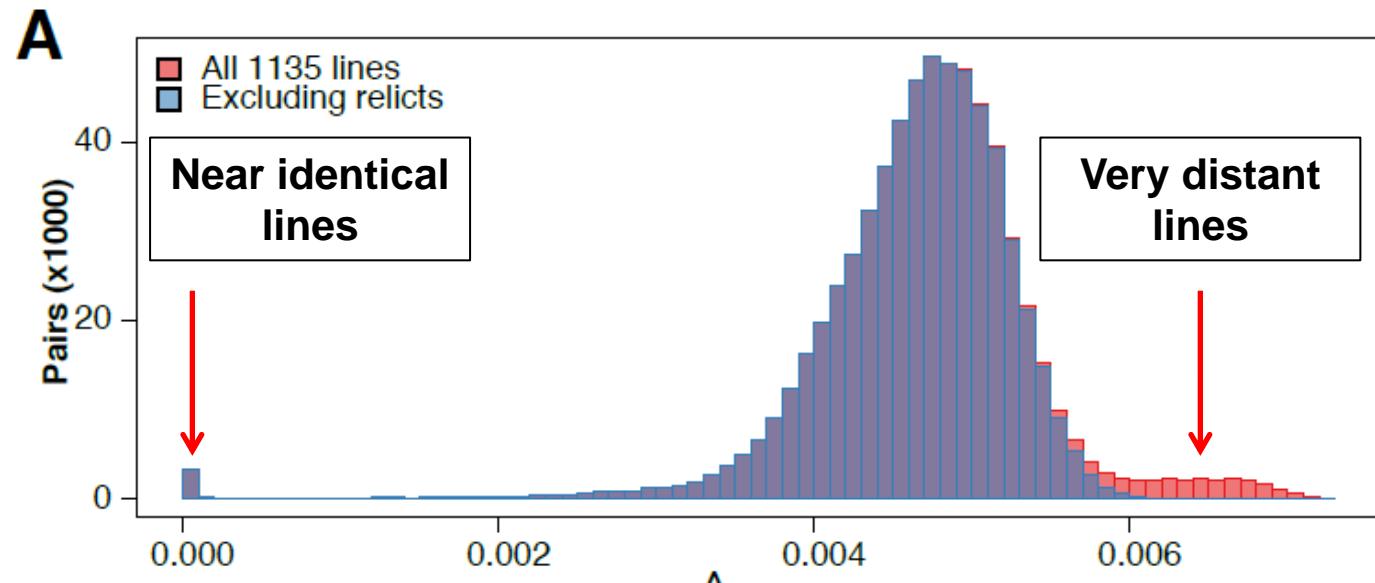
The 1001 Genomes Vision

# Origins of the 1,001 Genomes accessions



**Full genomes for a world wide collection of  
1135 different natural inbred lines**

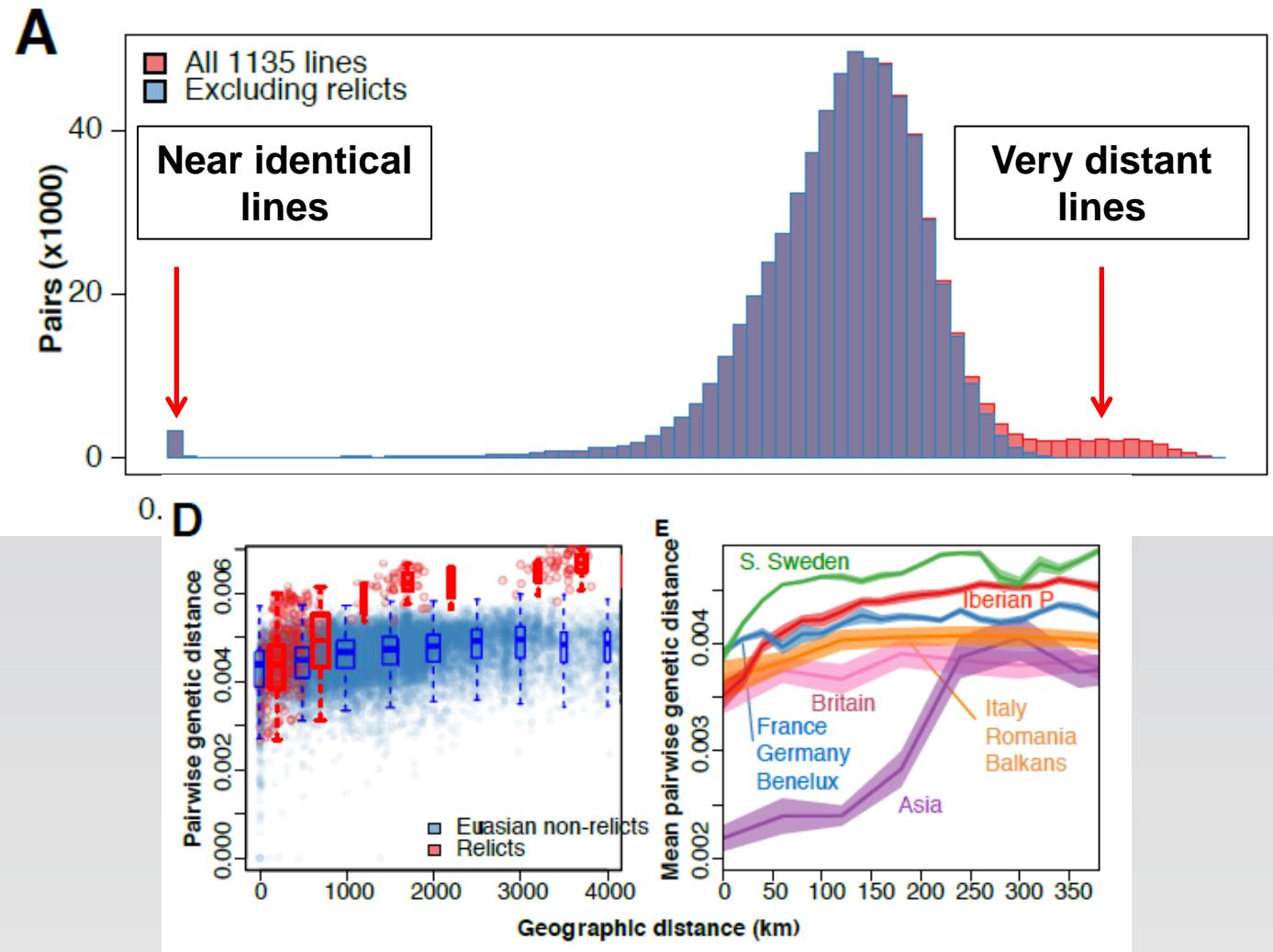
# Genetic and geographic distance



# Genetic and geographic distance



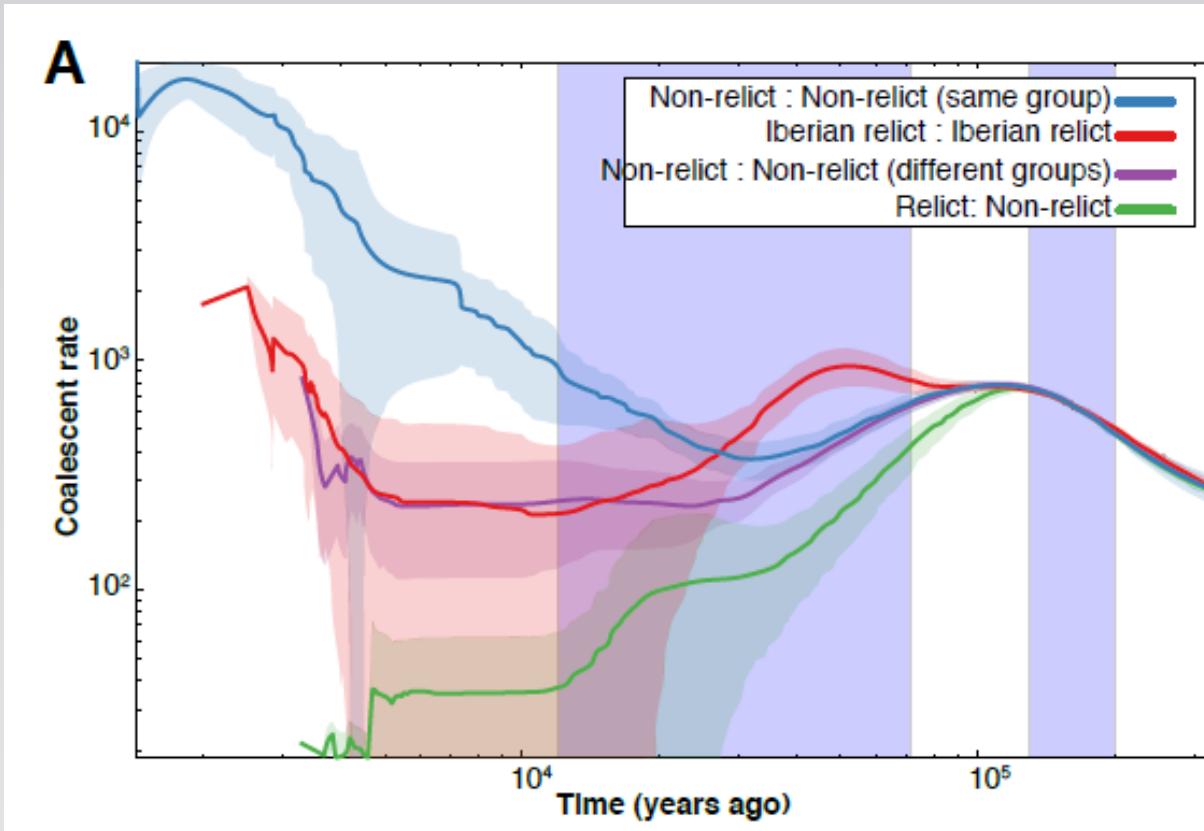
Center for Computational  
and Theoretical Biology



Isolation by distance seems to explain the observed diversity,  
but differs greatly between different subsets

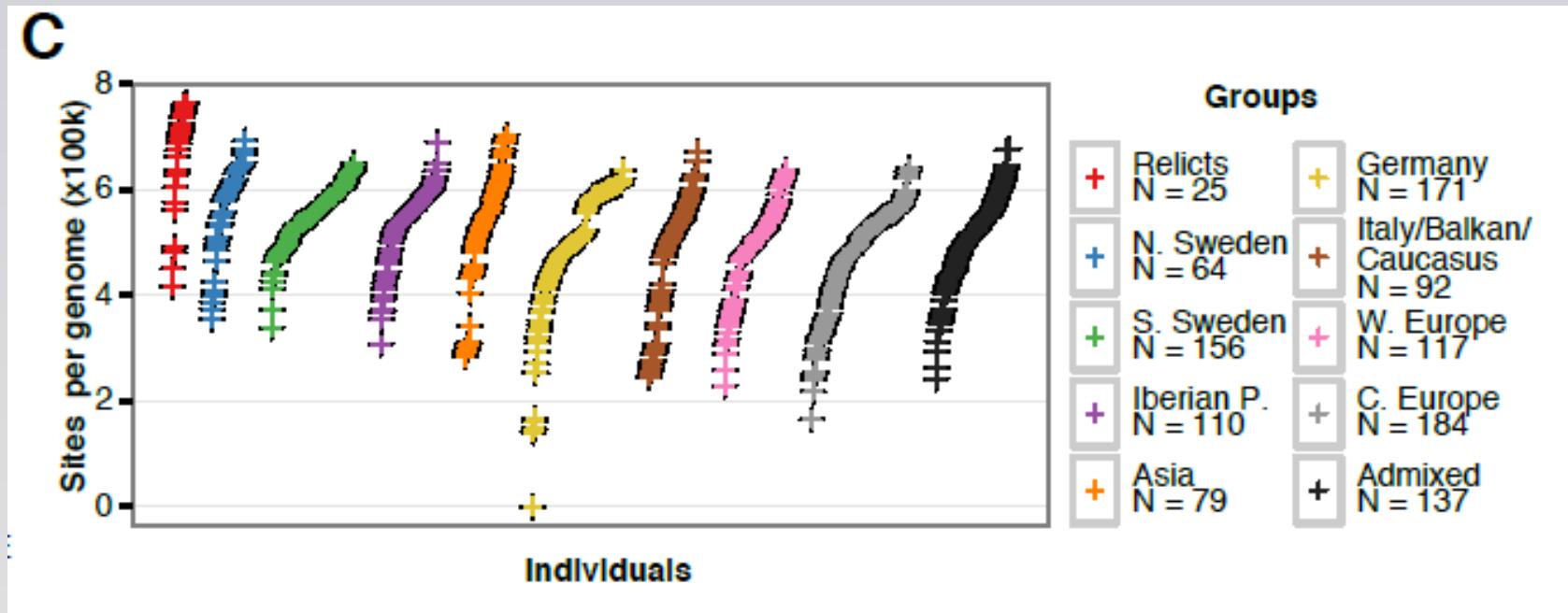
# Effect of the last glacial maximum

Inferred coalescence rates using MSMC (Schiffels and Durbin 2014)



Clear separation of groups during the last glacial maximum

# Genetic diversity in different groups



The total (sequenced) population contains 25 relicts and 8 different groups of *modern* accessions

## 1135 sequenced natural inbred lines reveal the global pattern of polymorphism in *Arabidopsis thaliana*

The 1001 Genomes Consortium\*

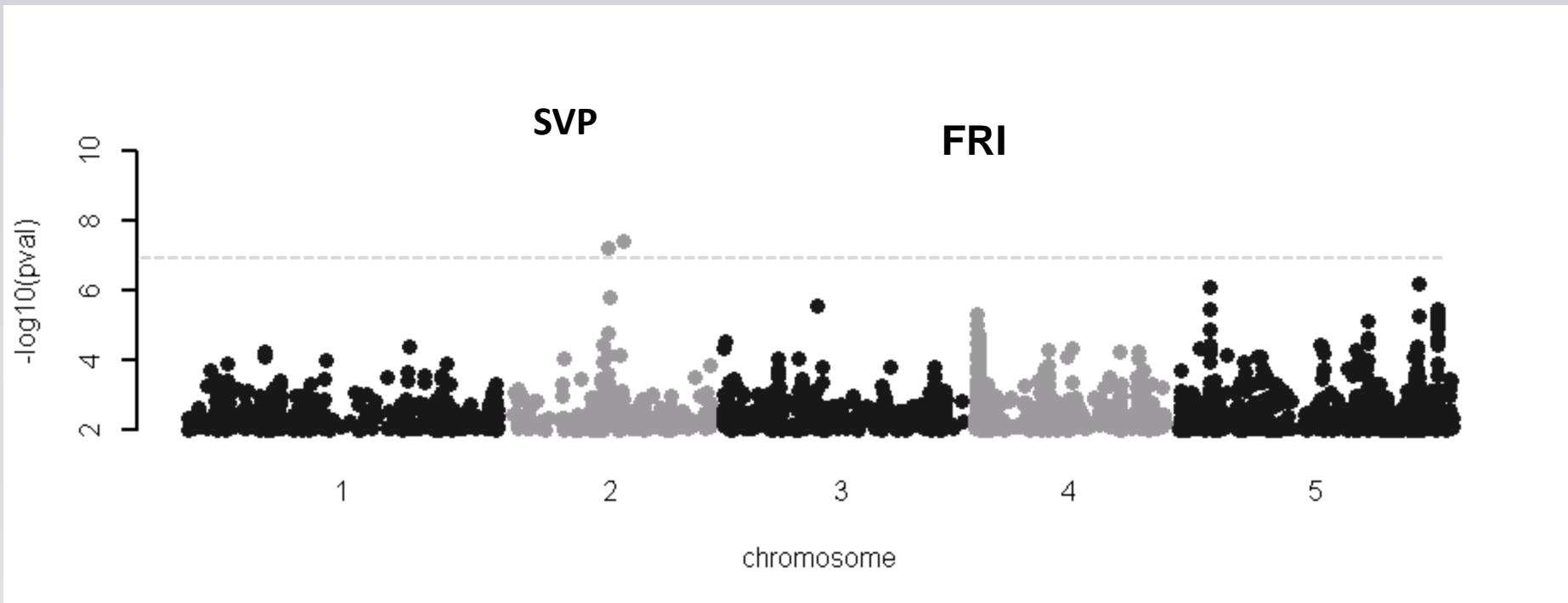
\*Lists of participants and their affiliations appear at the end of the paper.

**Data : 1,135 high quality genomes with more than 10 M SNP and 500k structural variants**

# GWAS: 250k data



Center for Computational  
and Theoretical Biology

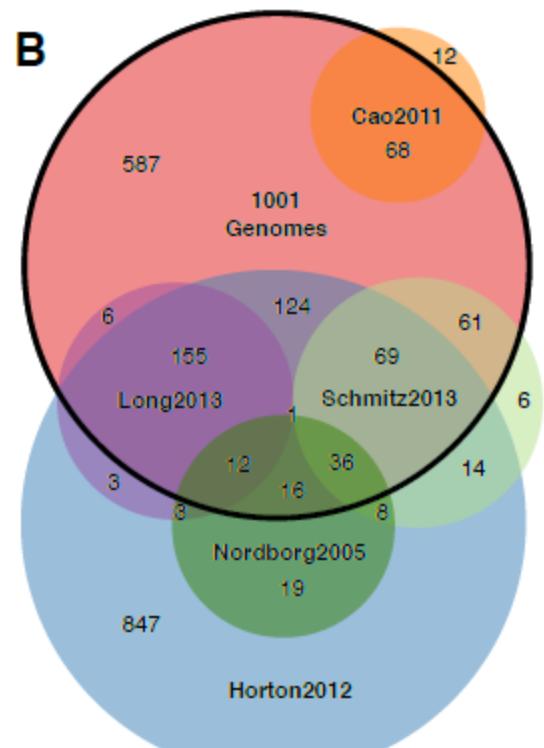


**250 k SNPs  
(1SNP/500 bp)**

# Inputted dataset for GWAS



Center for Computational  
and Theoretical Biology



**Full genome data for 1,135 accessions  
(10M SNPs)**

+

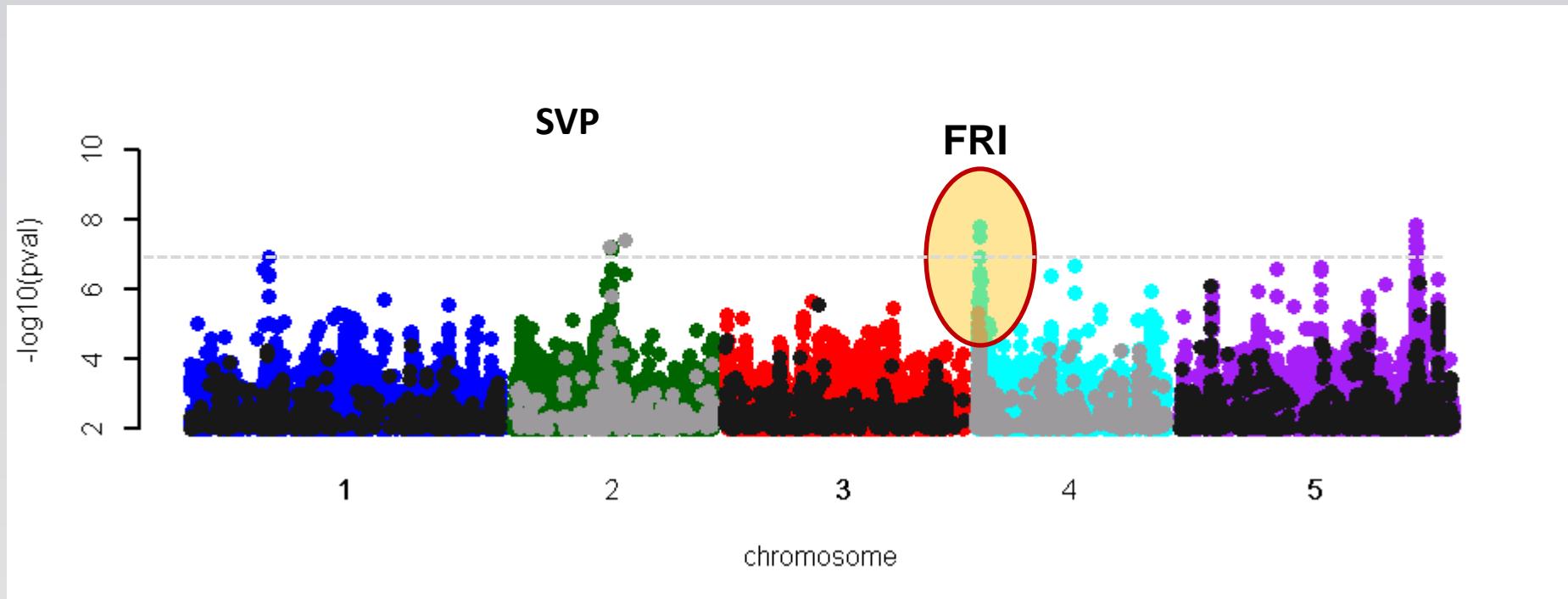
**250k SNPs for 1,300 accessions (Horton  
*et al.* 2012)**

**Generation of an imputed data set using  
BEAGLE  
for 2,029 accessions on 10 M SNPs**

# GWAS: Full sequence data



Center for Computational  
and Theoretical Biology

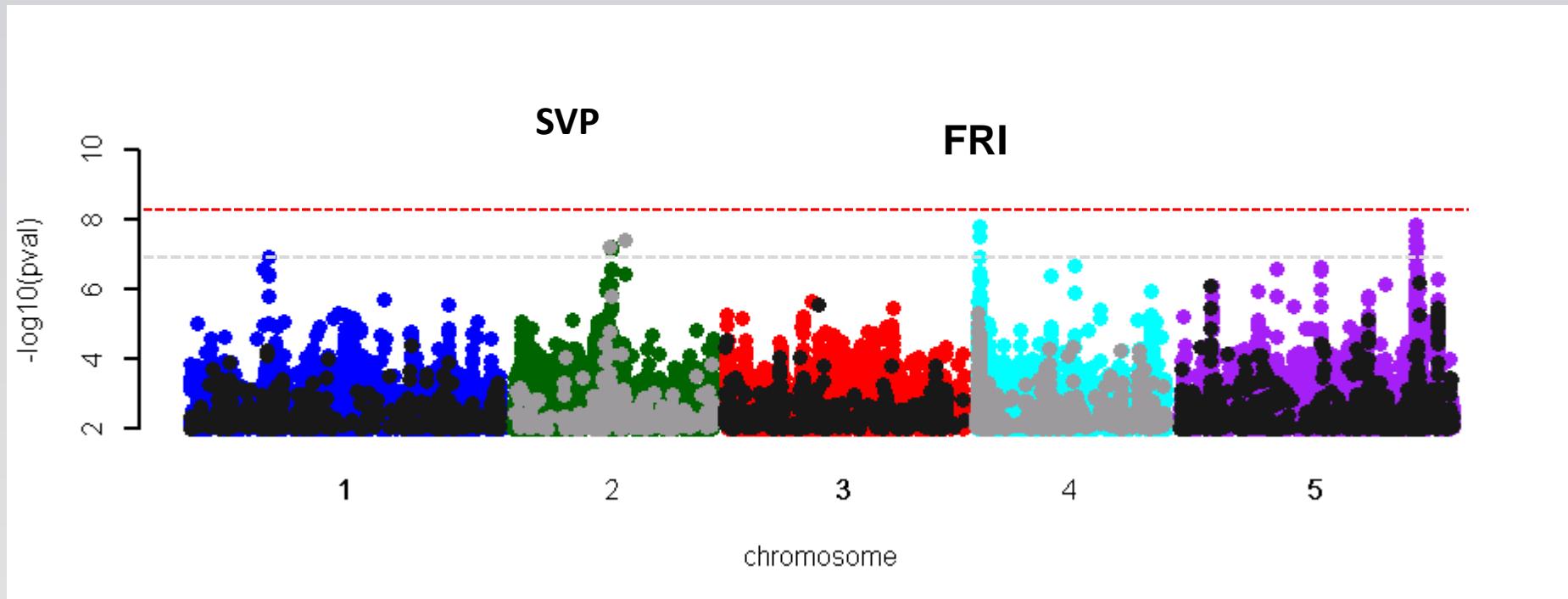


**10 M SNPs**

# GWAS: Full sequence data

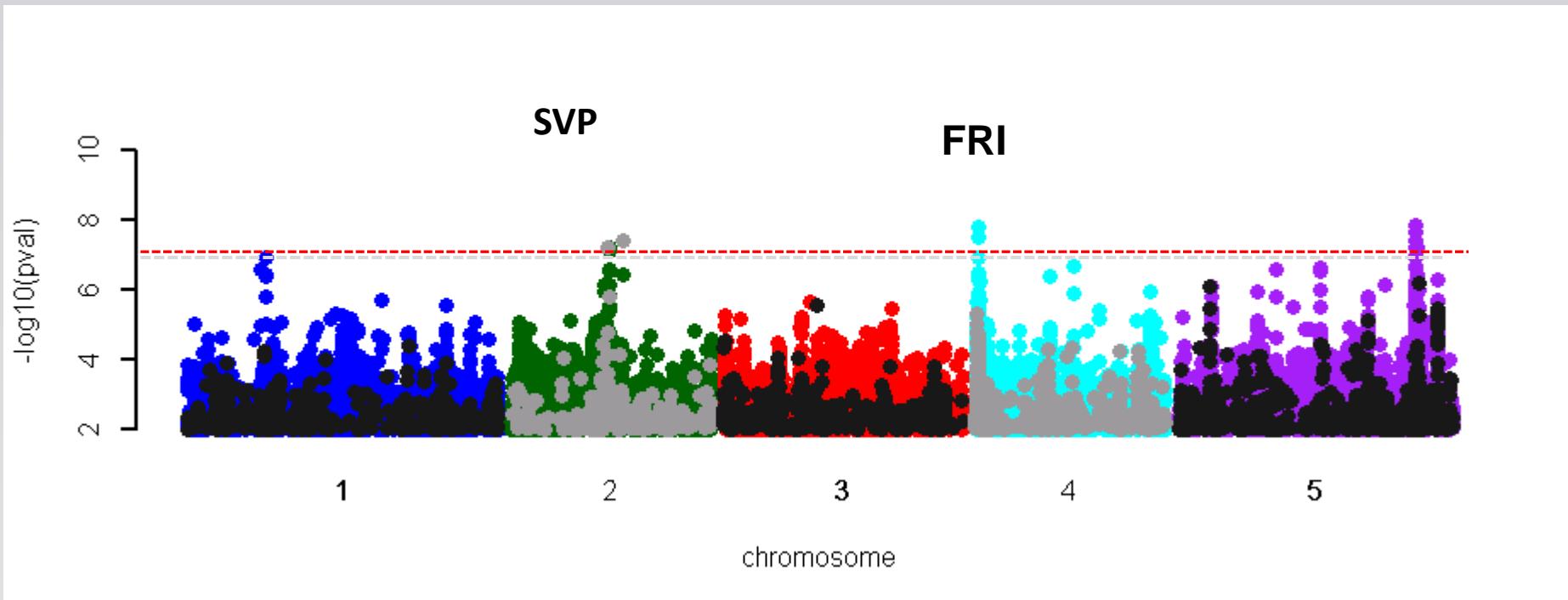


Center for Computational  
and Theoretical Biology



Threshold for multiple testing : Bonferoni

# GWAS: Full sequence data

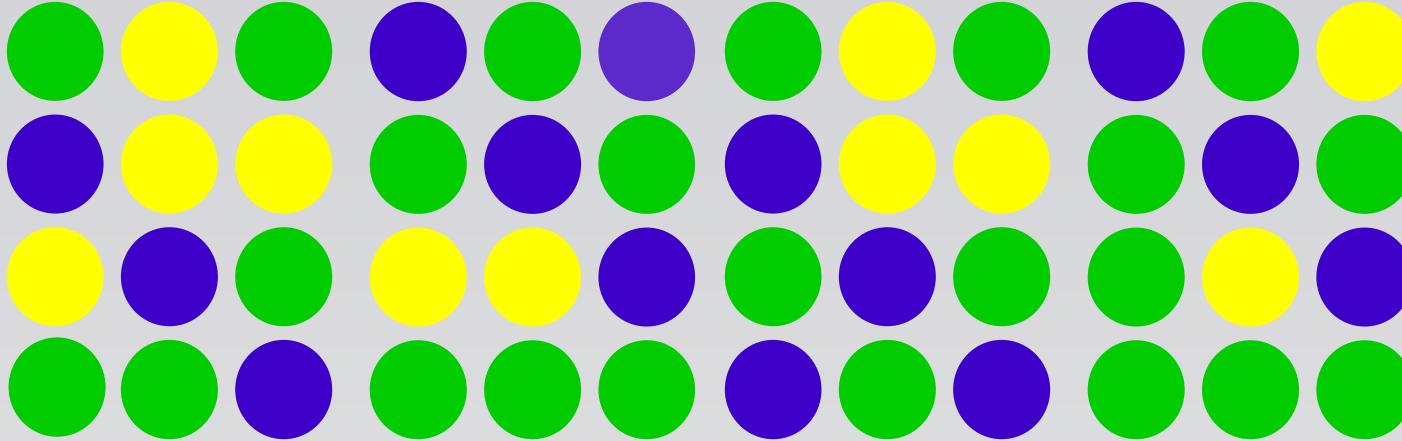


Permutation based threshold



Center for Computational  
and Theoretical Biology

# Practical part Ia



**Your project :**  
**Identify genetic factors that are  
important for making dots blue**

**Dataset : 48 individuals with phenotypes, and chip-data for 1.000  
homozygous SNPs**

# The data



```
> rm(list=ls())
> load('Y_dots.rda')
> load('X_dots.rda')
> load('K_dots.rda')
> source('emma.r')
> source('gwas.r')
> ls()
 [1] "amm_gwas"                  "emma.delta.ML.dLL.w.Z"    "emma.delta.ML.dLL.wo.Z"
 [4] "emma.delta.ML.LL.w.Z"      "emma.delta.ML.LL.wo.Z"    "emma.delta.REML.dLL.w.Z"
 [7] "emma.delta.REML.dLL.wo.Z"  "emma.delta.REML.LL.w.Z"   "emma.delta.REML.LL.wo.Z"
[10] "emma.eigen.L"              "emma.eigen.L.w.Z"         "emma.eigen.L.wo.Z"
[13] "emma.eigen.R"              "emma.eigen.R.w.Z"         "emma.eigen.R.wo.Z"
[16] "emma.kinship"              "emma.ML.LRT"               "emma.MLE"
[19] "emma.MLE.noX"              "emma.REML.t"               "emma.REMLE"
[22] "emma.test"                 "K_dots"                   "X_dots"
[25] "Y_dots"
> |
```

# The gwas.r script

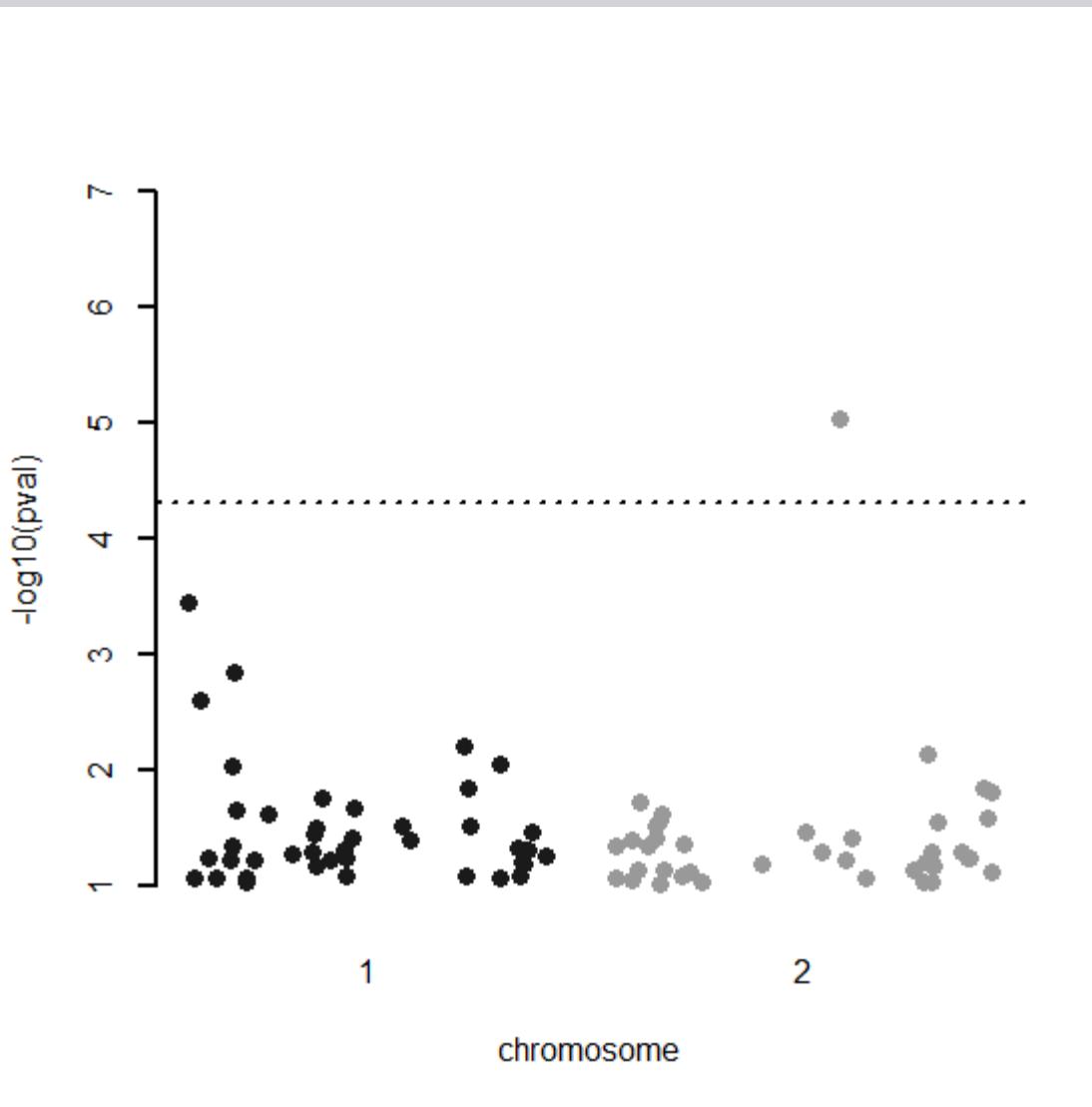
```
amm_gwas<-function(Y,X,K,p=0.001,n=2,run=T,calculate.effect.size=FALSE,  
include.lm=FALSE,use.SNP_INFO=FALSE,update.top_snps=FALSE,gen.data='binary') {
```

```
## amm_gwas.R -- a GWAS script for R 1.1.1
```

```
> GWAS_dots<-amm_gwas(Y_dots,X_dots,K_dots,n=3)  
GWAS performed on 48 ecotypes, 0 values excluded  
SNP_INFO file created  
pseudo-heritability estimate is 0.4381761  
> |
```

```
> source('plots_gwas.r')  
> plot_gwas(GWAS_dots,lower.limit=0.1)  
> |
```

# Manhattan plot



# The gwas.r script



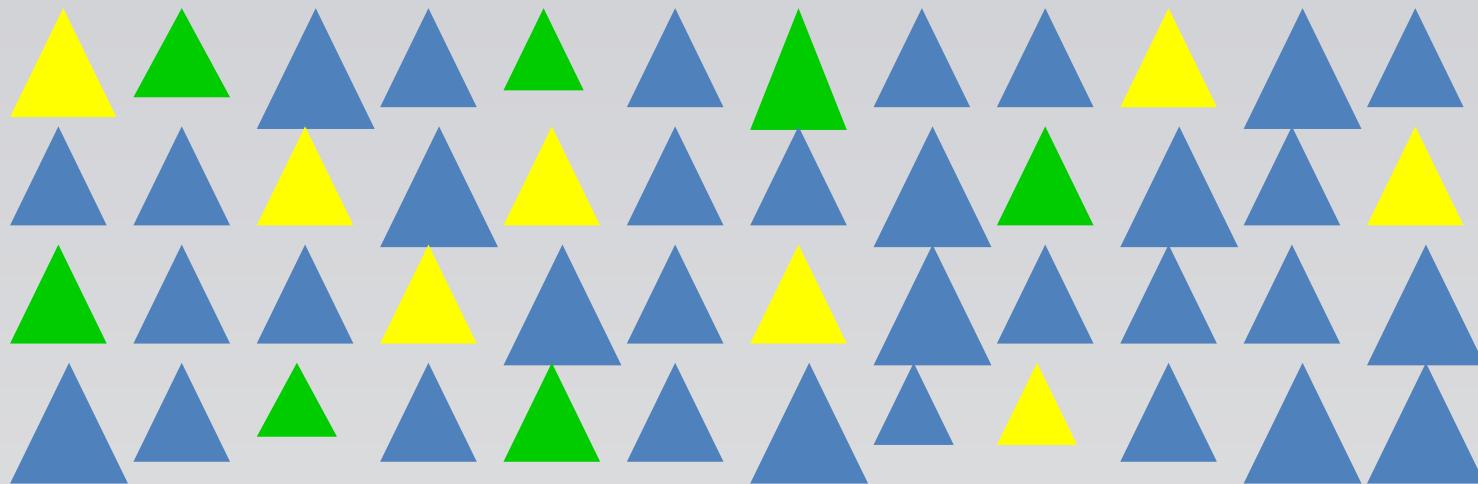
A screenshot of the RStudio IDE interface. The top menu bar includes File, Edit, Search, View, Encoding, Language, Settings, Macro, Run, Plugins, Window, and Help. Below the menu is a toolbar with various icons for file operations like Open, Save, Print, and Find. The bottom navigation bar shows several tabs: ideas.txt, grc\_2015.txt, tals\_Tübingen.txt, find\_clust.r, mtmm\_test.r, gwas.r (which is the active tab), AT4G04740\_t.txt, and plots\_. The main code editor area displays the 'gwas.r' script:

```
1 #####  
2 ### AMM -- R script for GWAS corecting for population structure (similar to EMMA)  
3 ###  
4 #####  
5 #  
6 ##  
7 ##  
8 ##  
9 #  
10  
11  
12 ##REQUIRED DATA & FORMAT  
13  
14 #requires functions from the original emma function (Kang et al. 2008, Genetics)  
15 #source('emma.r')  
16 #PHENOTYPE - Y: a n by 1 matrix, where n=number of individuals and the row  
17  
18 #GENOTYPE - X: a n by m matrix, where n=number of individuals, m=number of SNPs  
19 #KINSHIP - K: a n by n matrix, with rownames(K)=colnames(K)=individual names  
20 #each of these data being sorted in the same way, according to the individual  
21 #  
22 #  
23 #SNP INFORMATION - SNP_INFO: a data frame having at least 3 columns:  
24 # - 1 named 'SNP', with SNP names (same as colnames(X)),  
25 # - 1 named 'Chr', with the chromosome number to which belong each SNP  
26 # - 1 named 'Pos' with the position of the SNP onto the chromosome it belongs to
```

# Practical part Ib



Center for Computational  
and Theoretical Biology



**Your project :**  
**Identify genetic factors that are**  
**important for the size of**  
**triangles.**

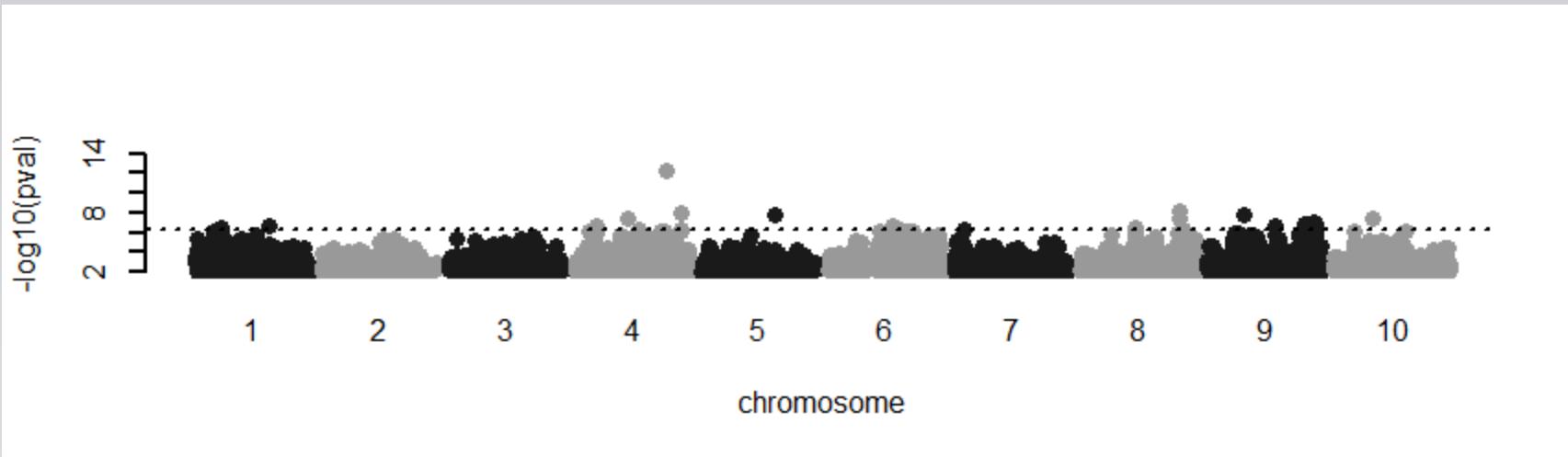
**Dataset : 480 individuals with 2 phenotypes, and full sequencing data**

**Y\_tri.rda, X\_tri.rda, X\_imp.rda, K\_tri.rda, SNP\_tri.rda**

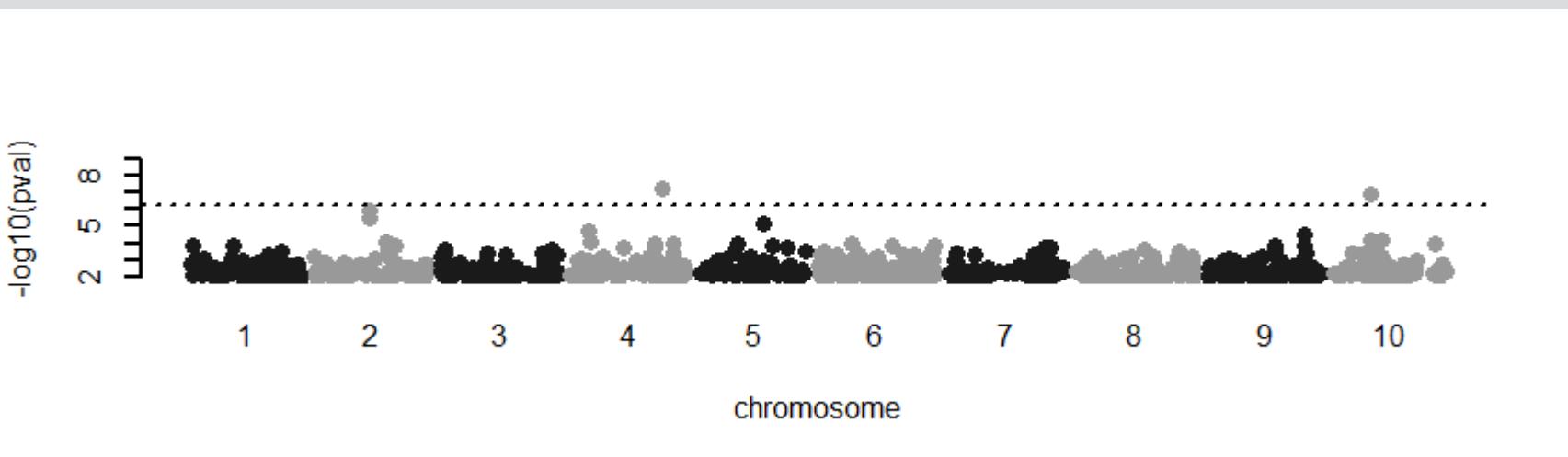
# Manhattan plots triangles



Center for Computational  
and Theoretical Biology



**Linear model**

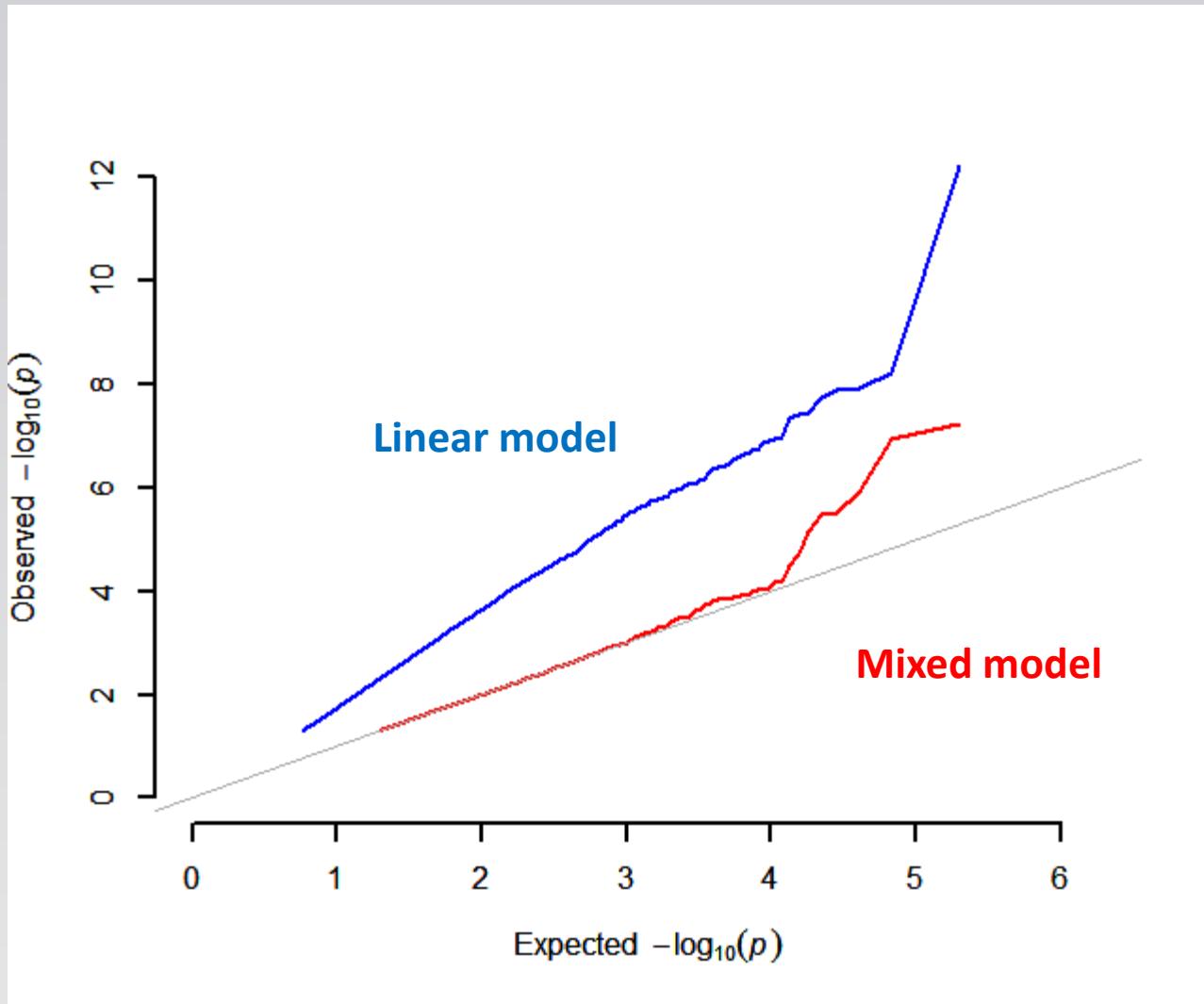


**Mixed model**

# qq\_plots triangles



Center for Computational  
and Theoretical Biology



# Batch effects



```
o      12 15.01777    o  45.2007    FRA  blue   1
> summary(lm(size~batch,data=Y_tri))

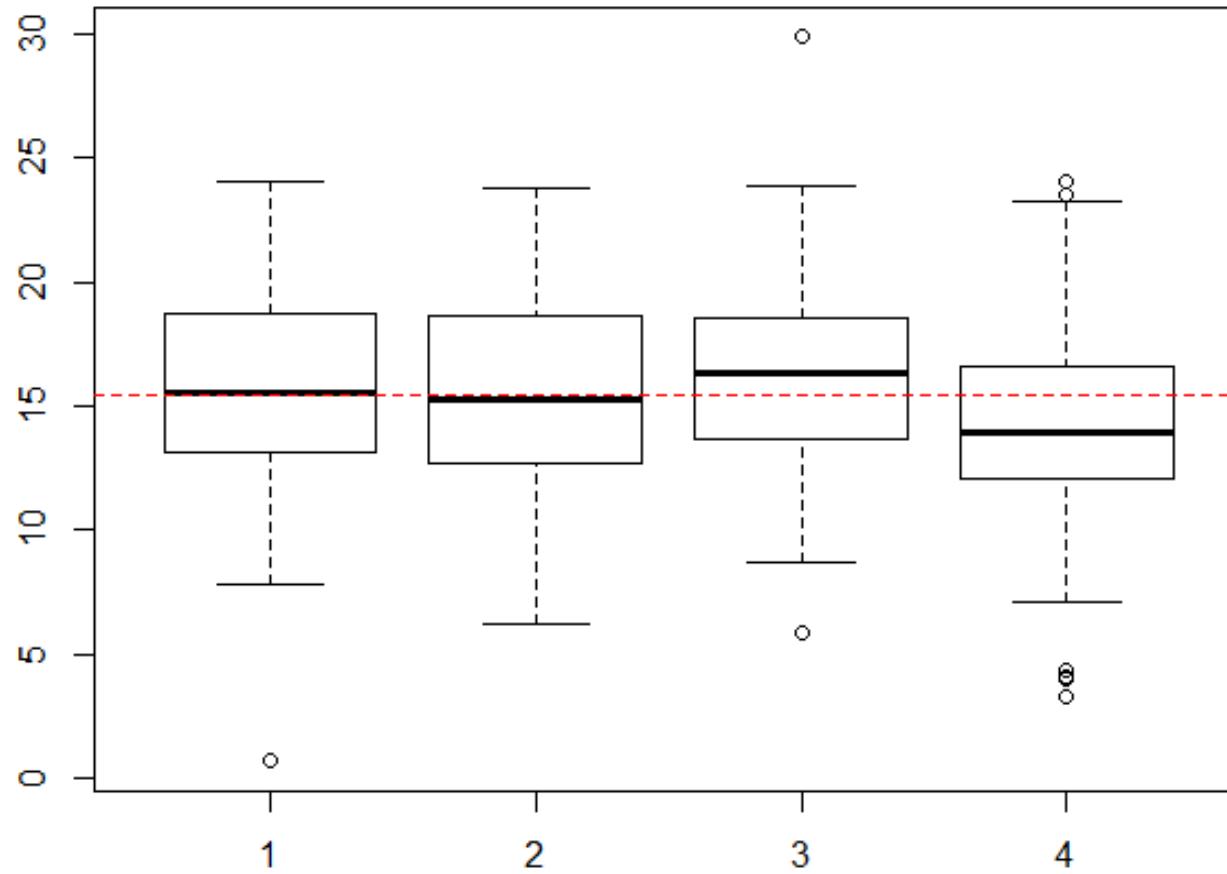
Call:
lm(formula = size ~ batch, data = Y_tri)

Residuals:
    Min      1Q  Median      3Q     Max 
-15.2259 -2.6554 -0.1001  2.9074 14.5458 

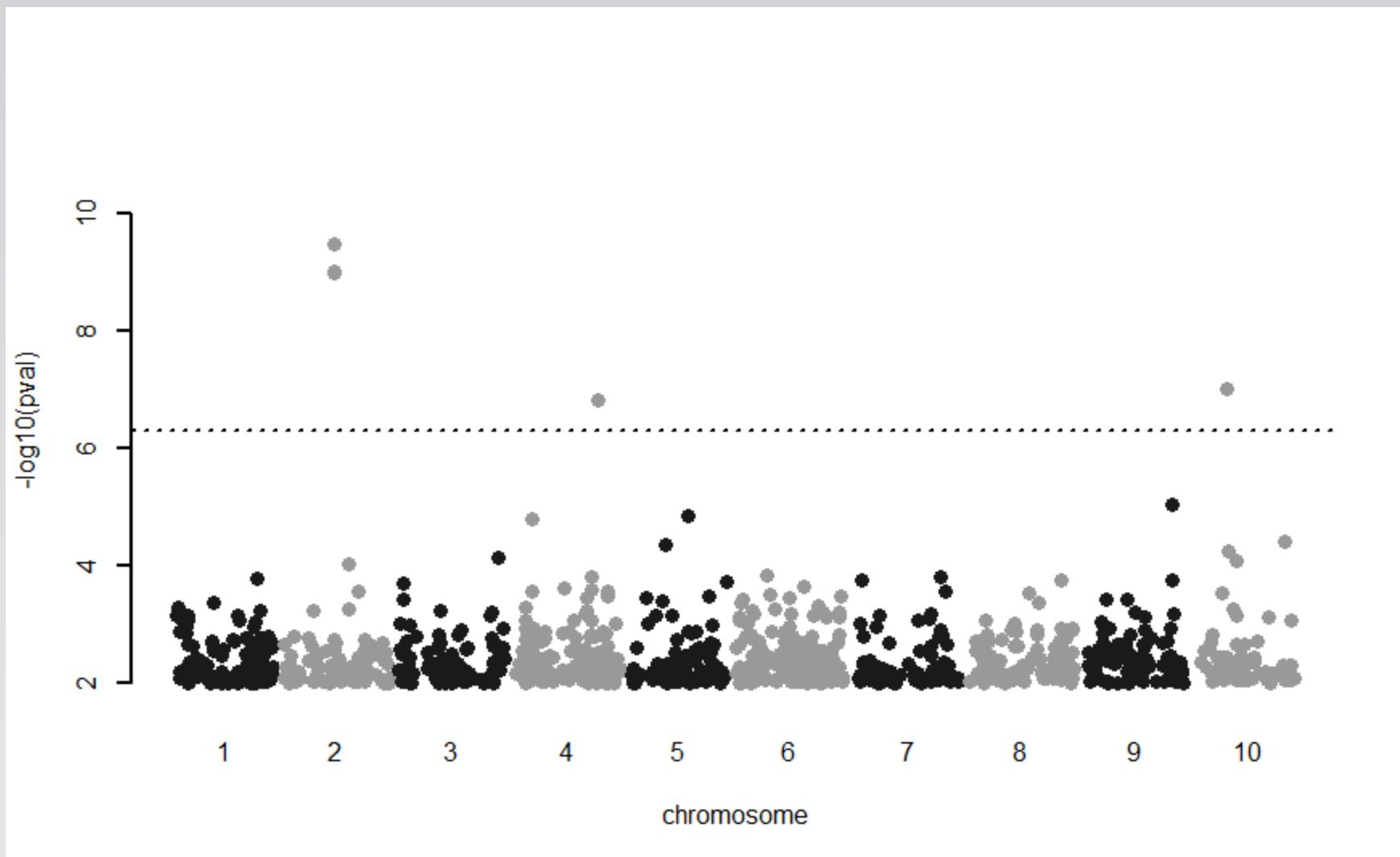
Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 16.1951     0.4466  36.264 <2e-16 ***
batch       -0.2909     0.1631  -1.784   0.0751 .  
---
Signif. codes:
0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.994 on 478 degrees of freedom
Multiple R-squared:  0.006613, Adjusted R-squared:  0.004535 
F-statistic: 3.182 on 1 and 478 DF,  p-value: 0.07508
```

# Batch effects



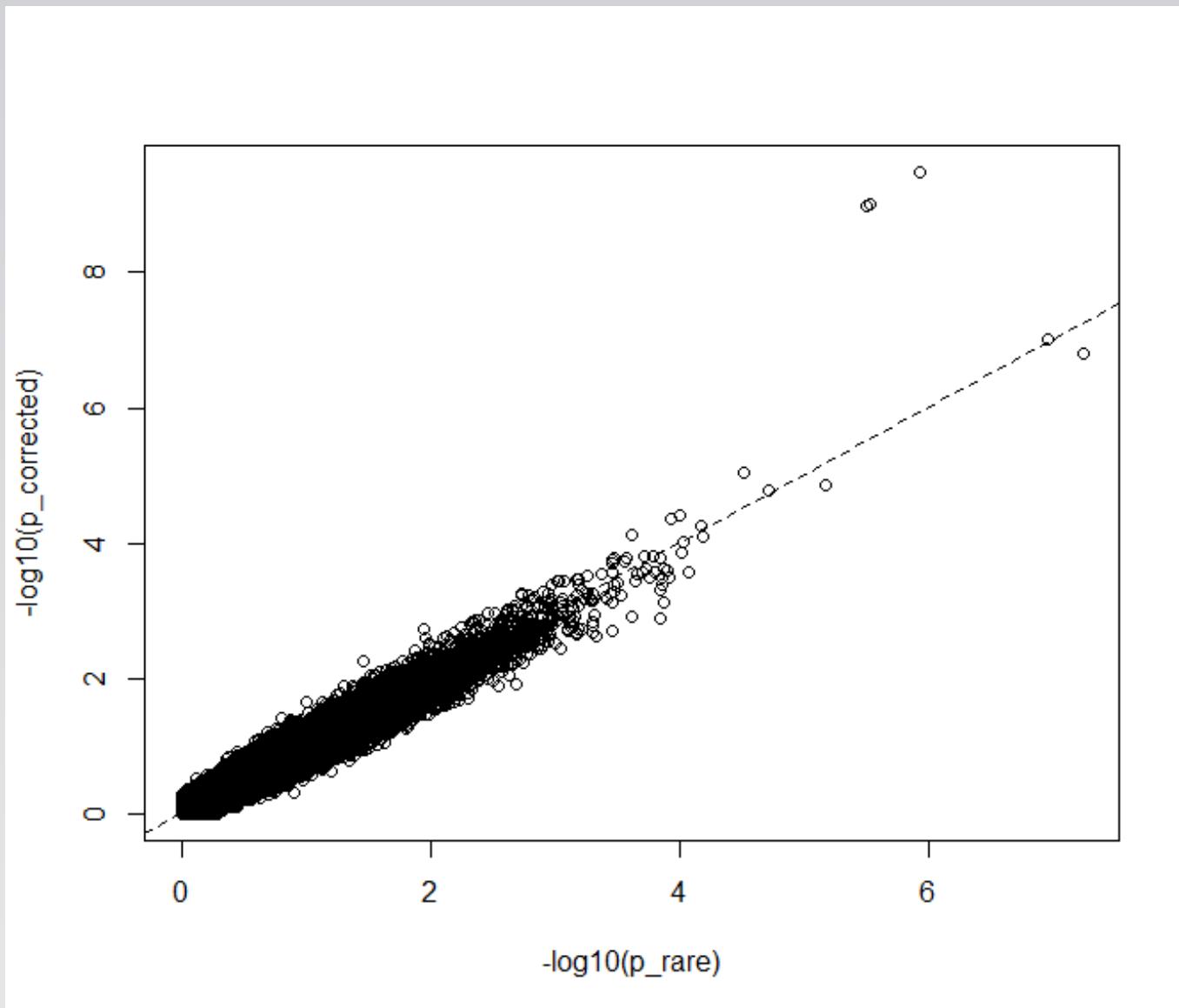
# Batch effects



# Batch effects



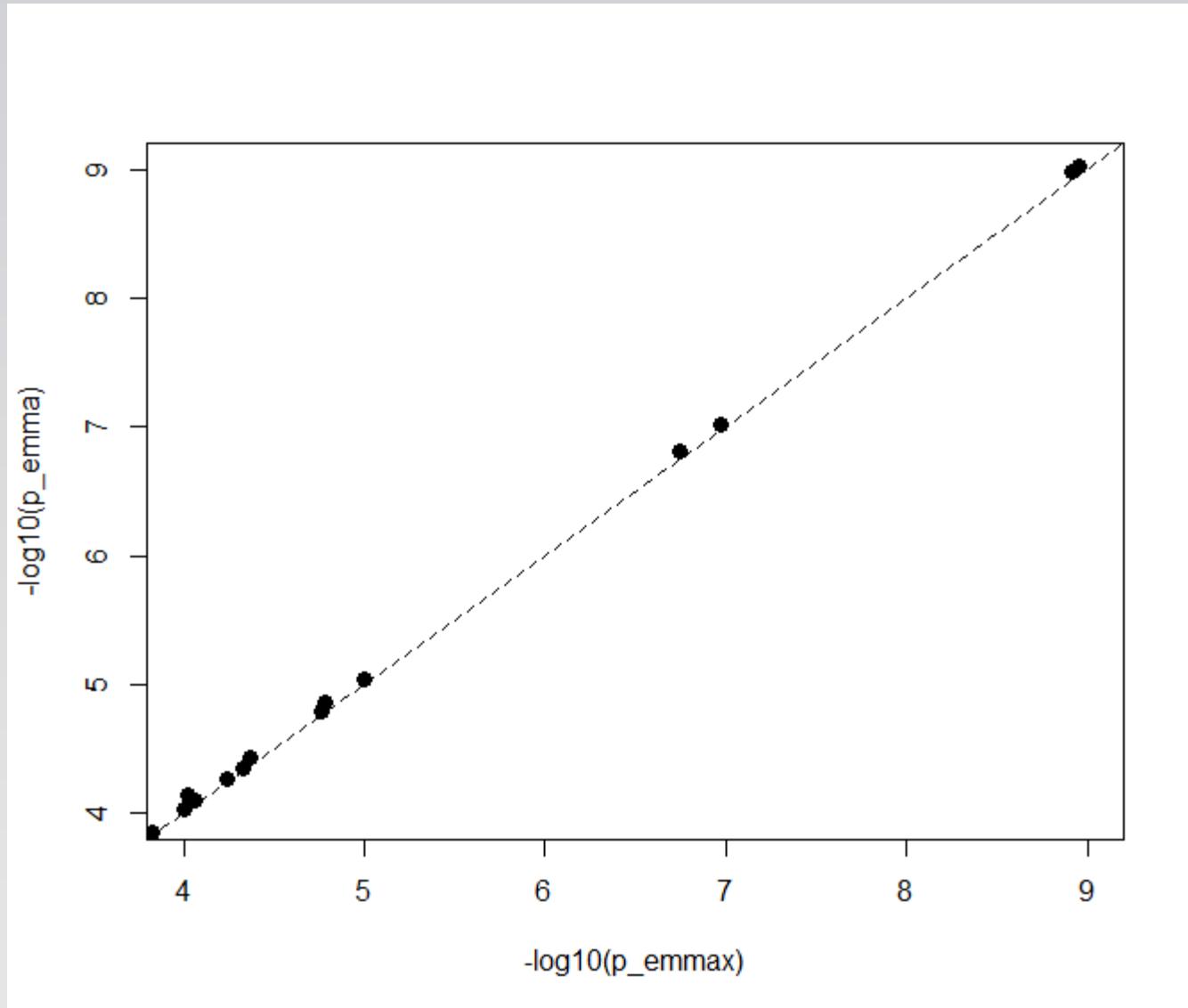
Center for Computational  
and Theoretical Biology



# Emma vs. EmmaX



Center for Computational  
and Theoretical Biology



# **lunch**



Center for Computational  
and Theoretical Biology

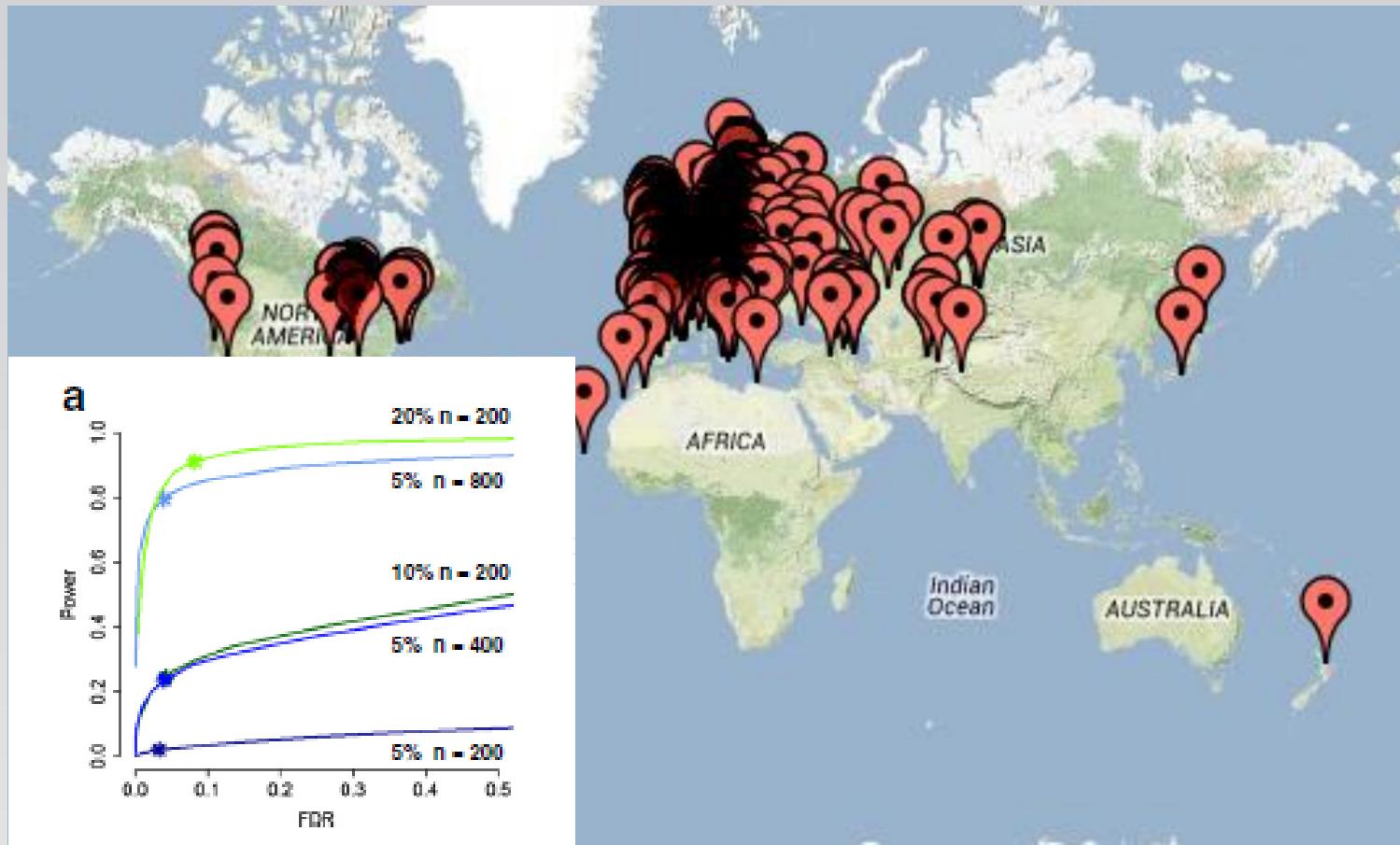


# Which lines should I phenotype?



- > 1,000 different lines completely sequenced
- more than 10M SNPs (imputed for over 2,000 lines) and 500K structural variants

# Which lines should I phenotype?

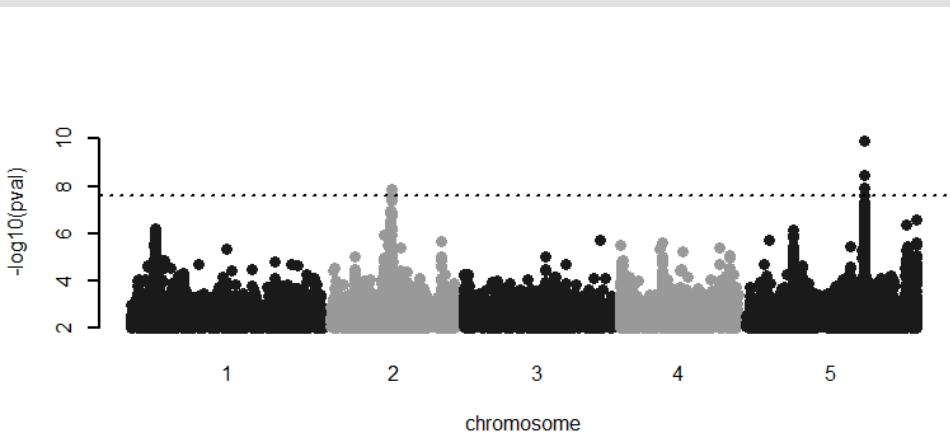
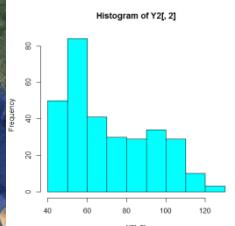
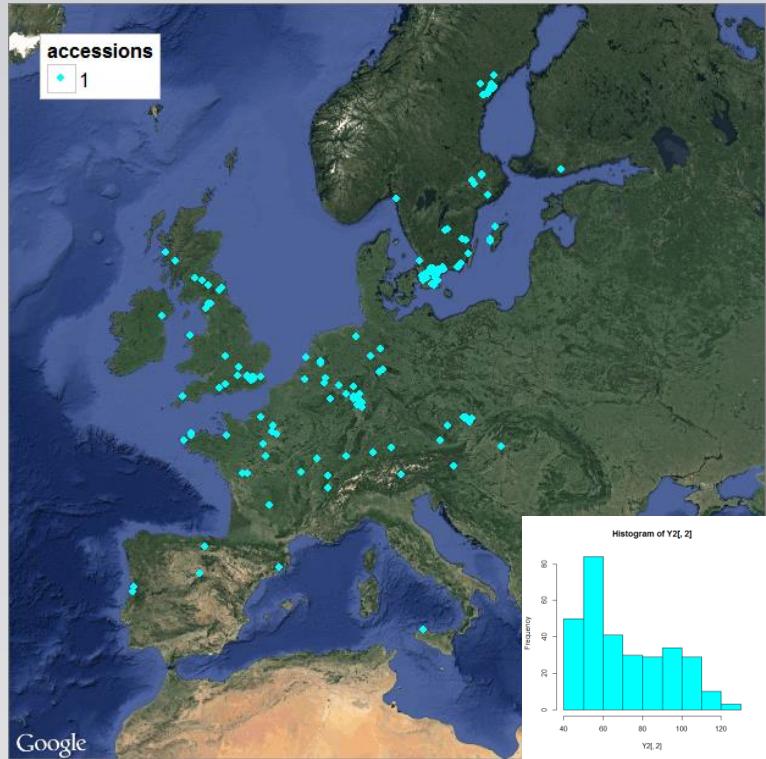


- > 1,000 different lines completely sequenced
- more than 10M SNPs (imputed for over 2,000 lines) and 500K structural variants

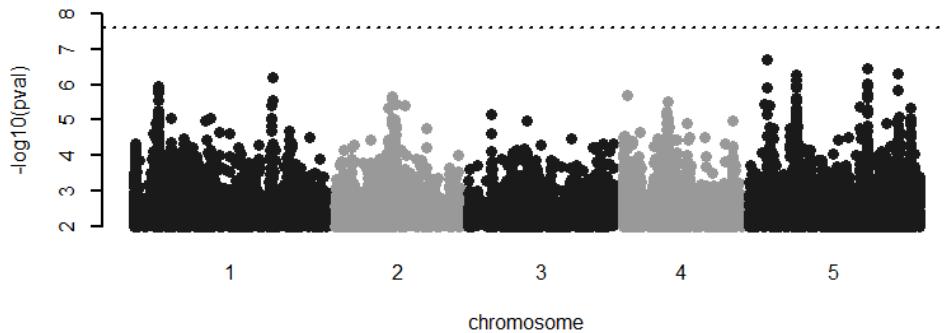
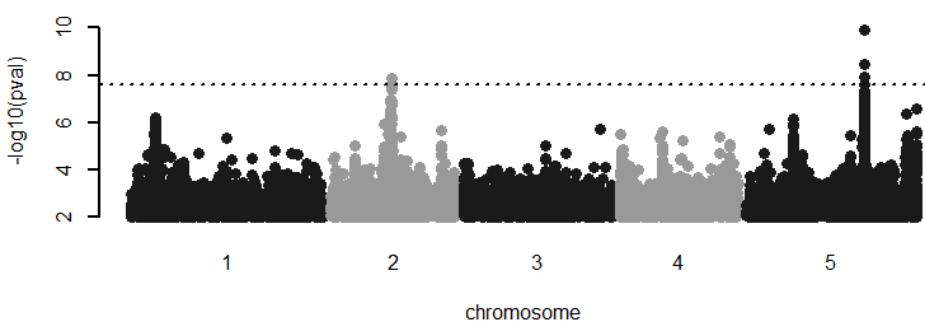
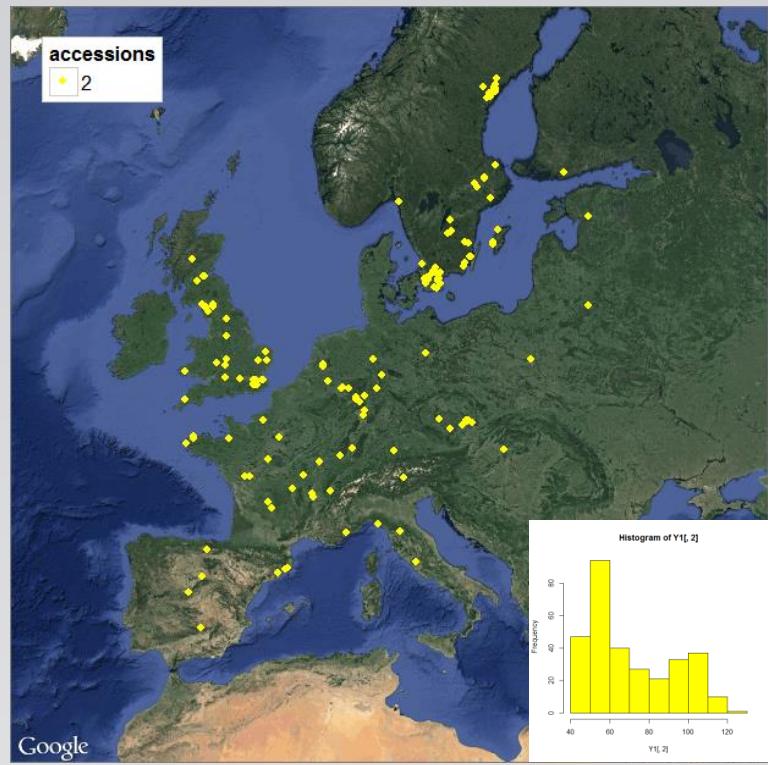
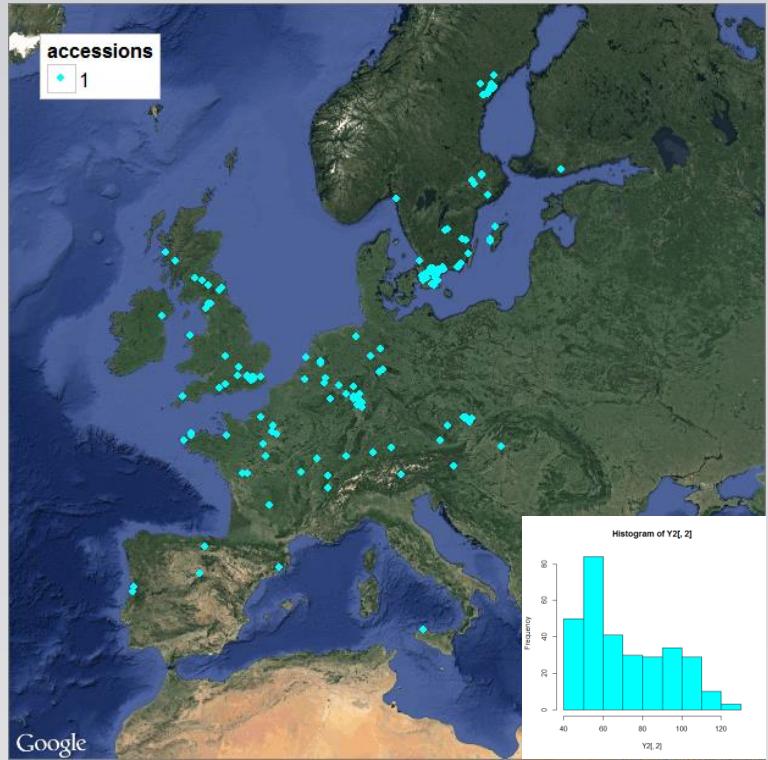
# GWAS on flowering time (200 accessions)



Center for Computational  
and Theoretical Biology



# GWAS on flowering time (200 accessions)

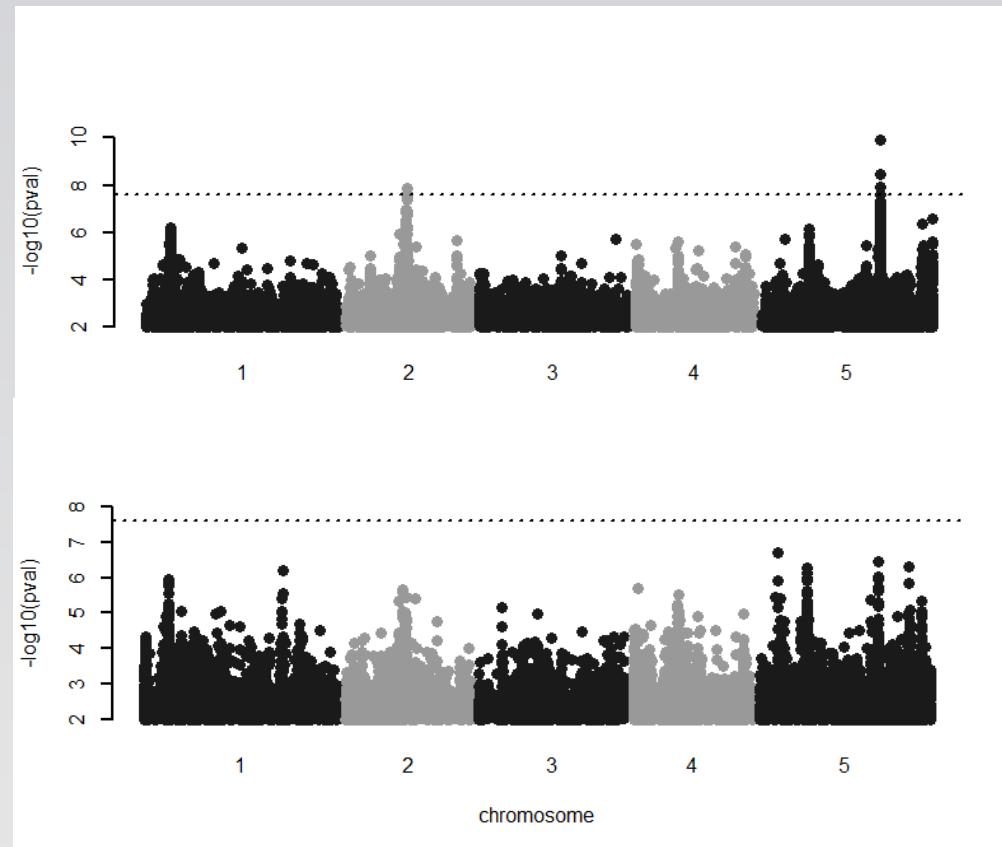


# GWAS on flowering time (200 accessions)



Center for Computational  
and Theoretical Biology

1000 different random subsets of 200 accessions

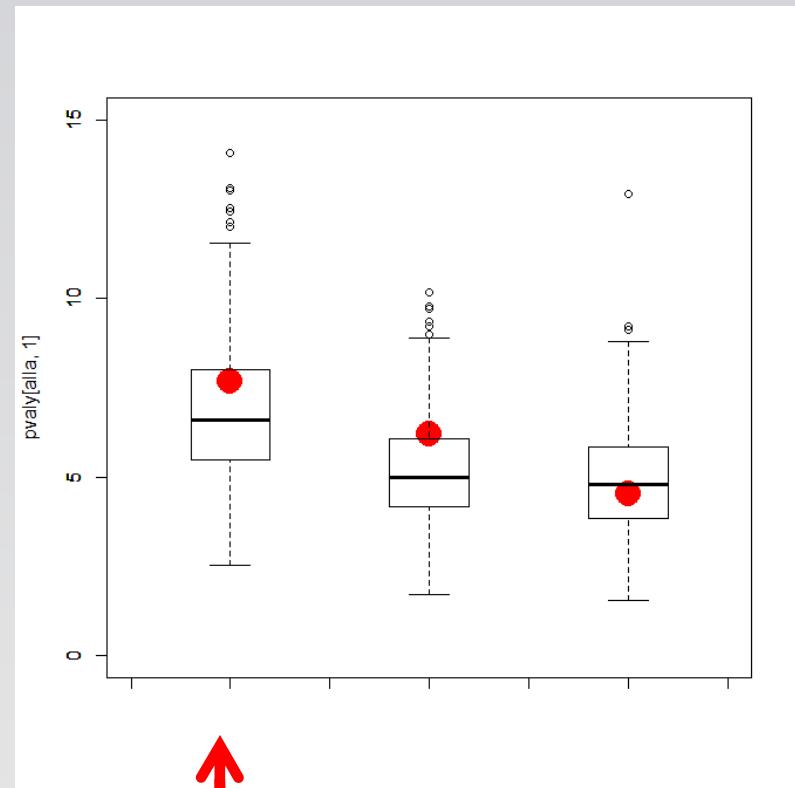
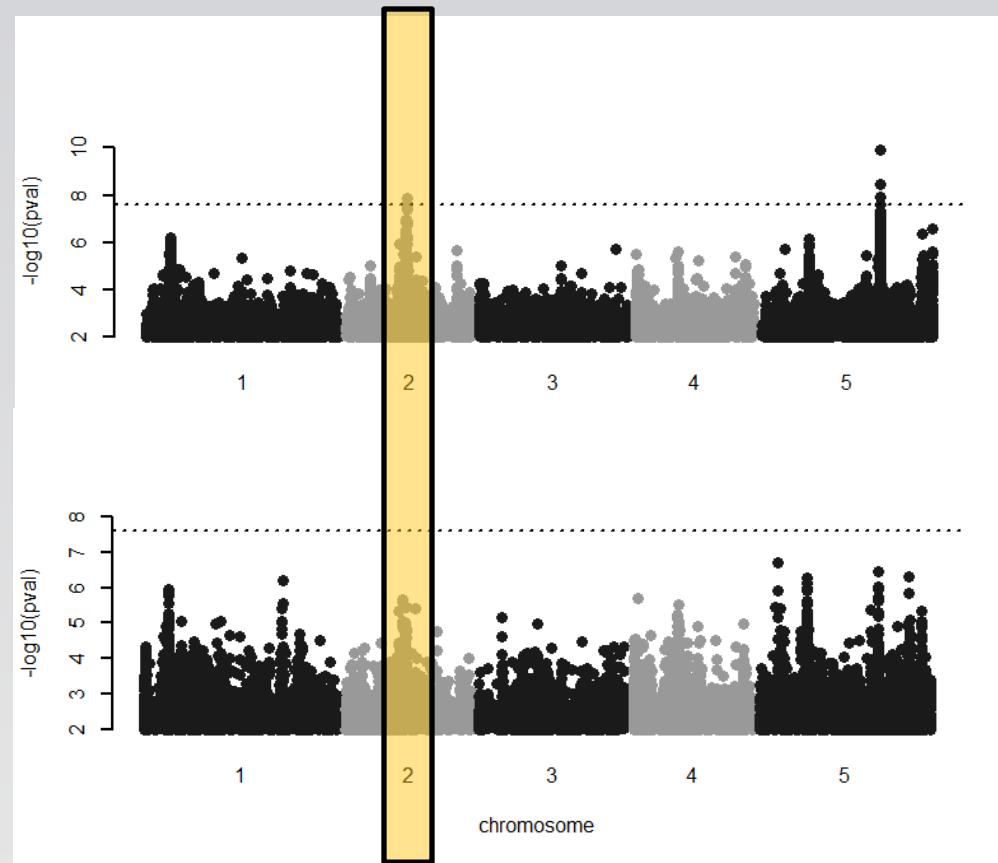


# GWAS on flowering time (200 accessions)

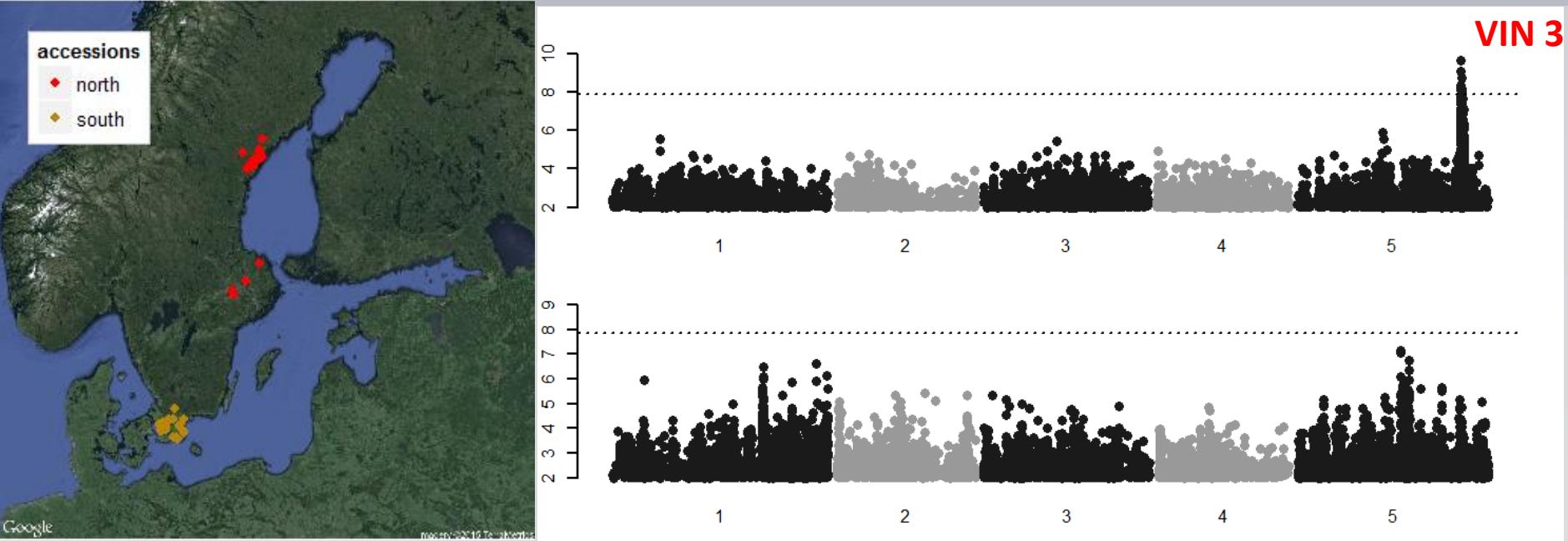


Center for Computational  
and Theoretical Biology

1000 different random subsets of 200 accessions

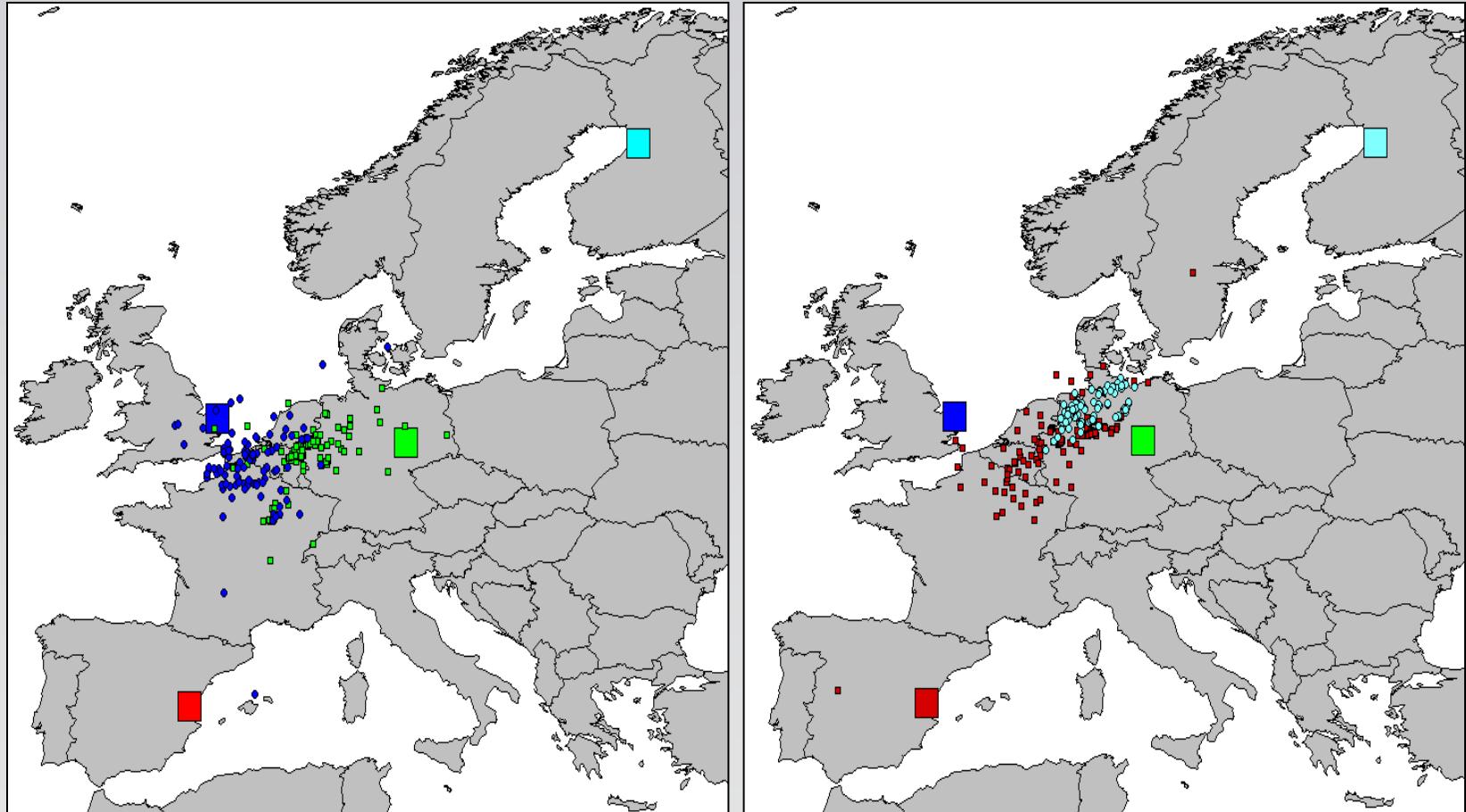


# GWAS in local subsets



Different results in different subsets

# Local Adaptation

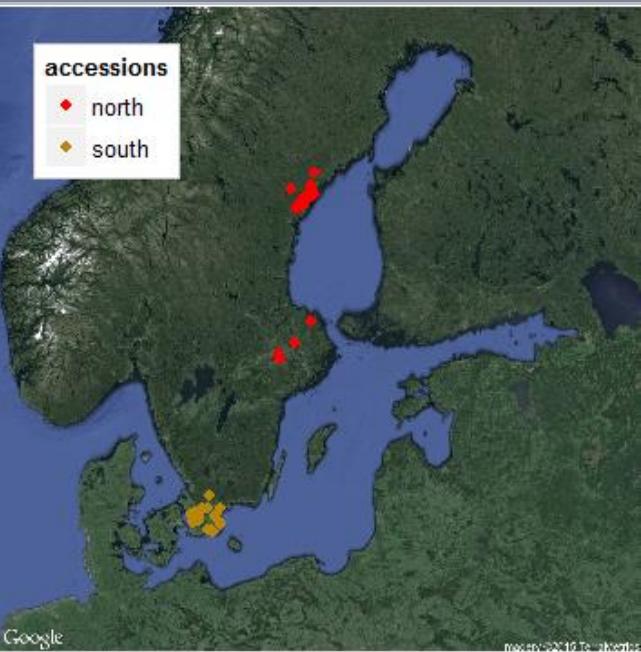


**Genes providing fitness under natural conditions are mostly local**

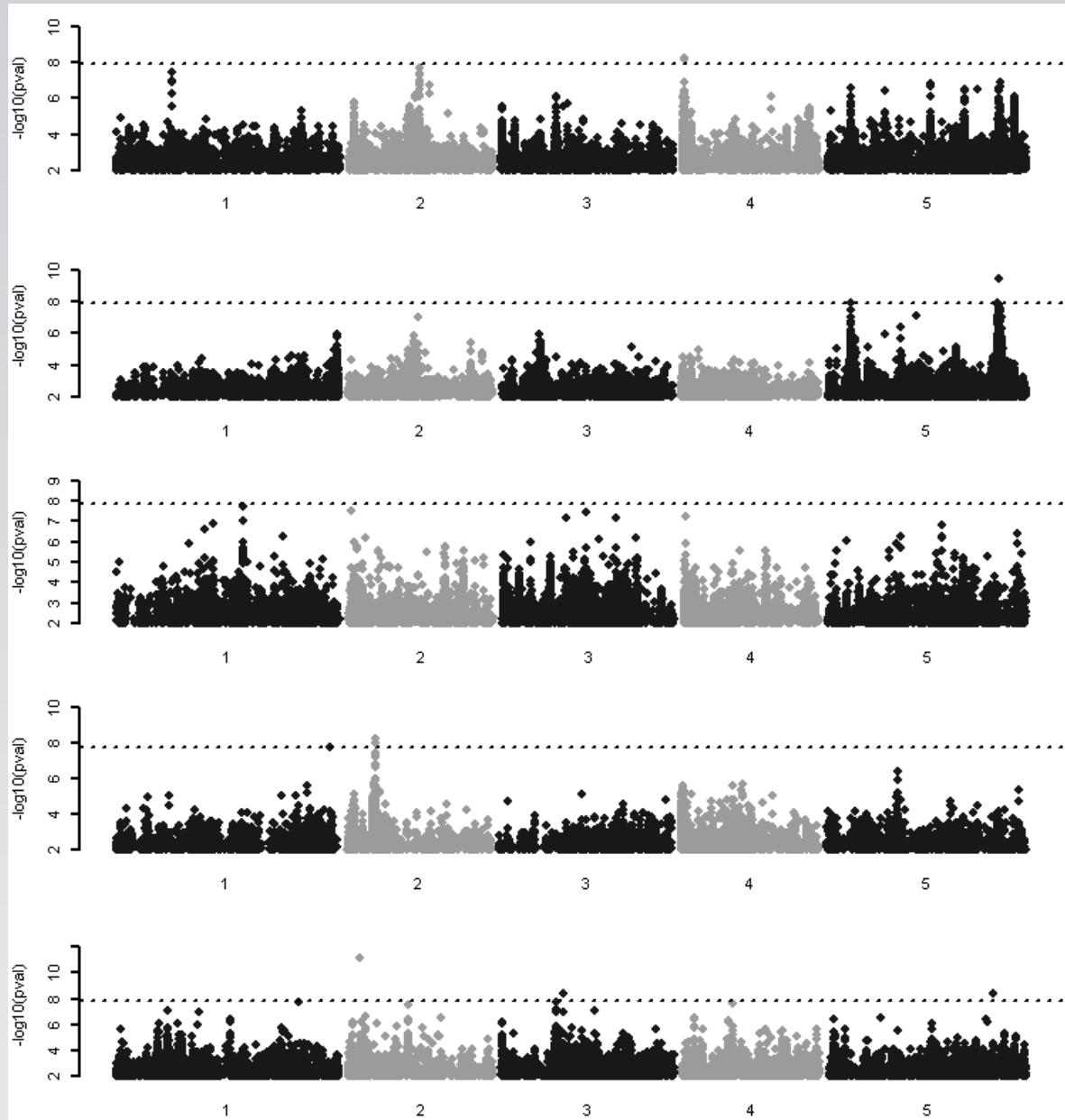
# GWAS in local subsets



Center for Computational  
and Theoretical Biology



# GWAS results in different subsets



All

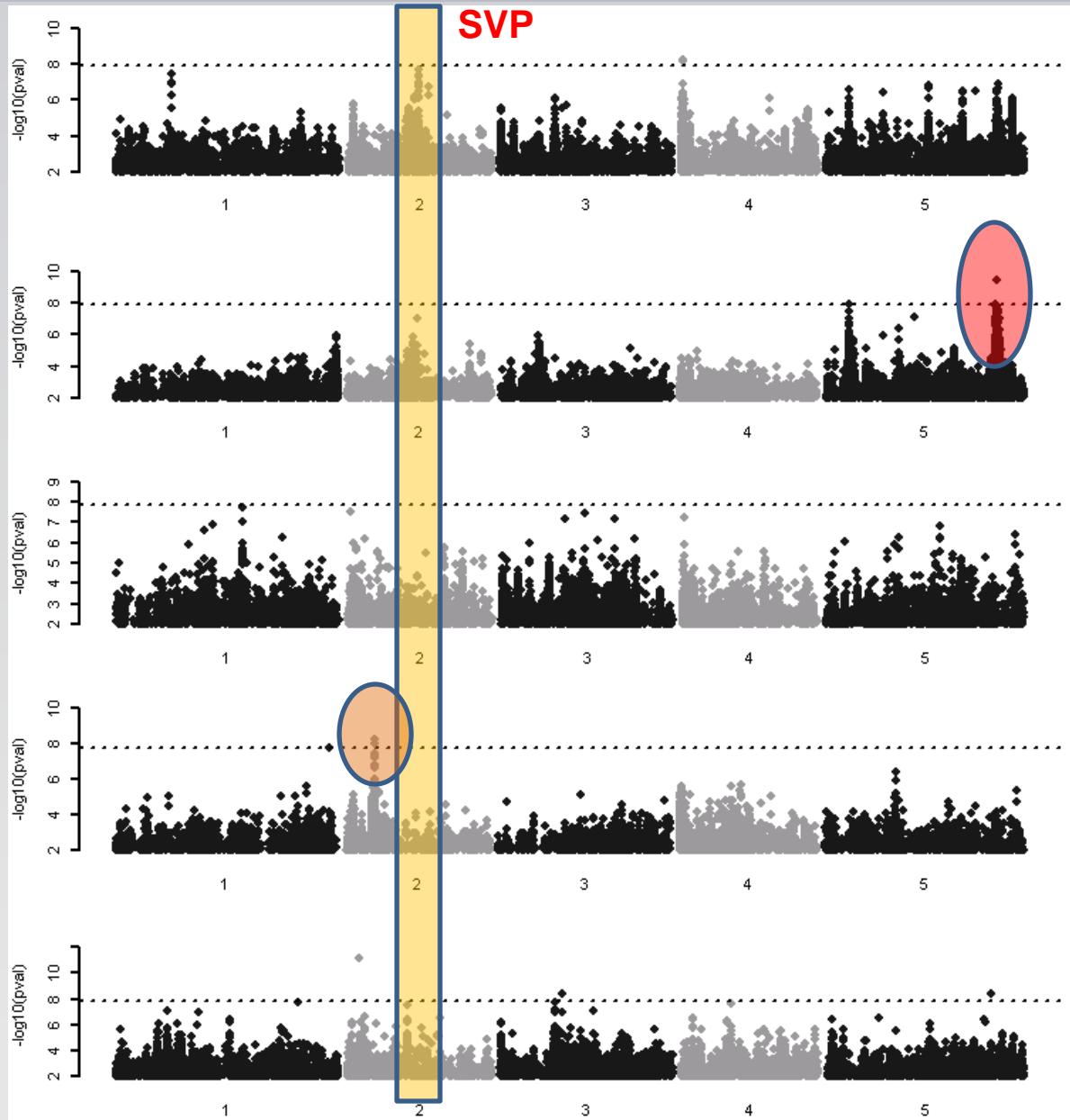
SWE

FRA

CZE

UK

# GWAS results in different subsets



All

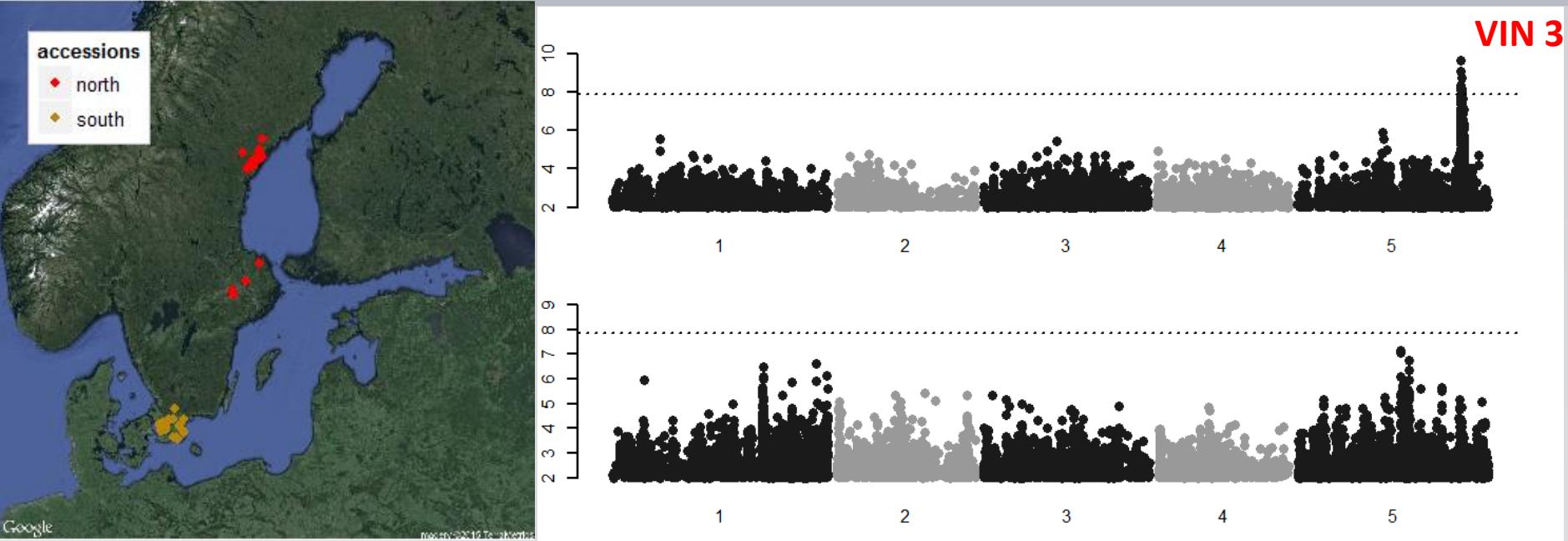
SWE

FRA

CZE

UK

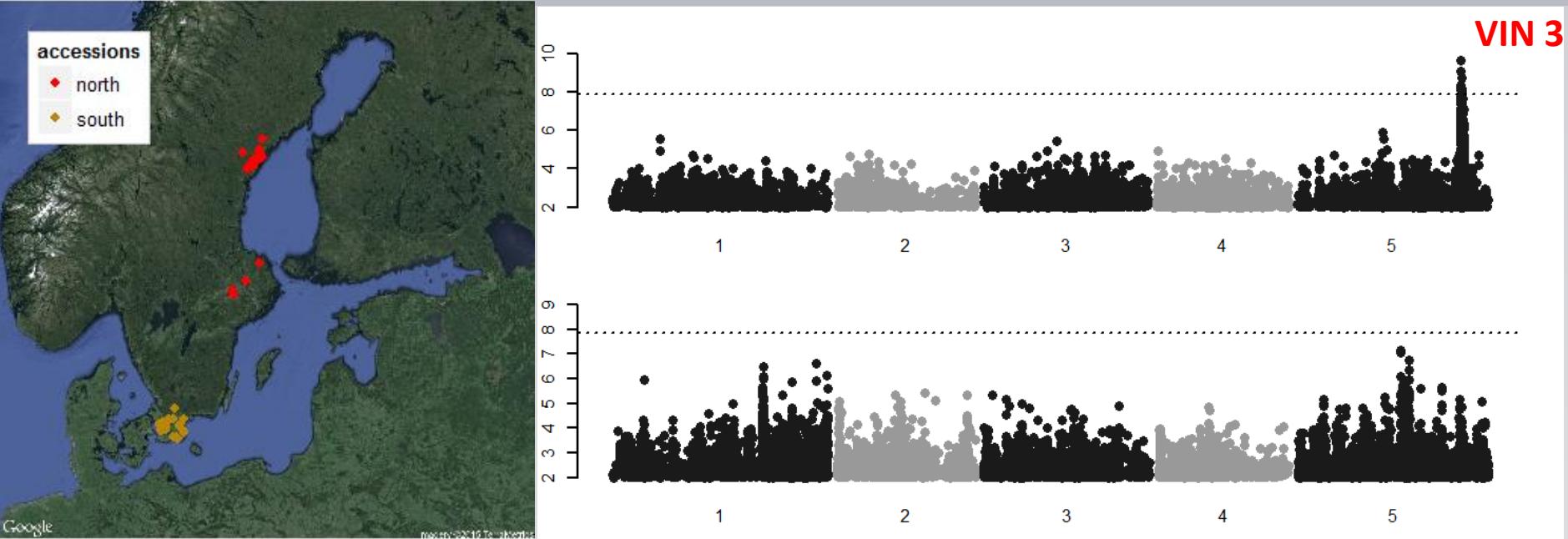
# GWAS in local subsets



Why ?

- Differences in Allele frequency of the causative marker

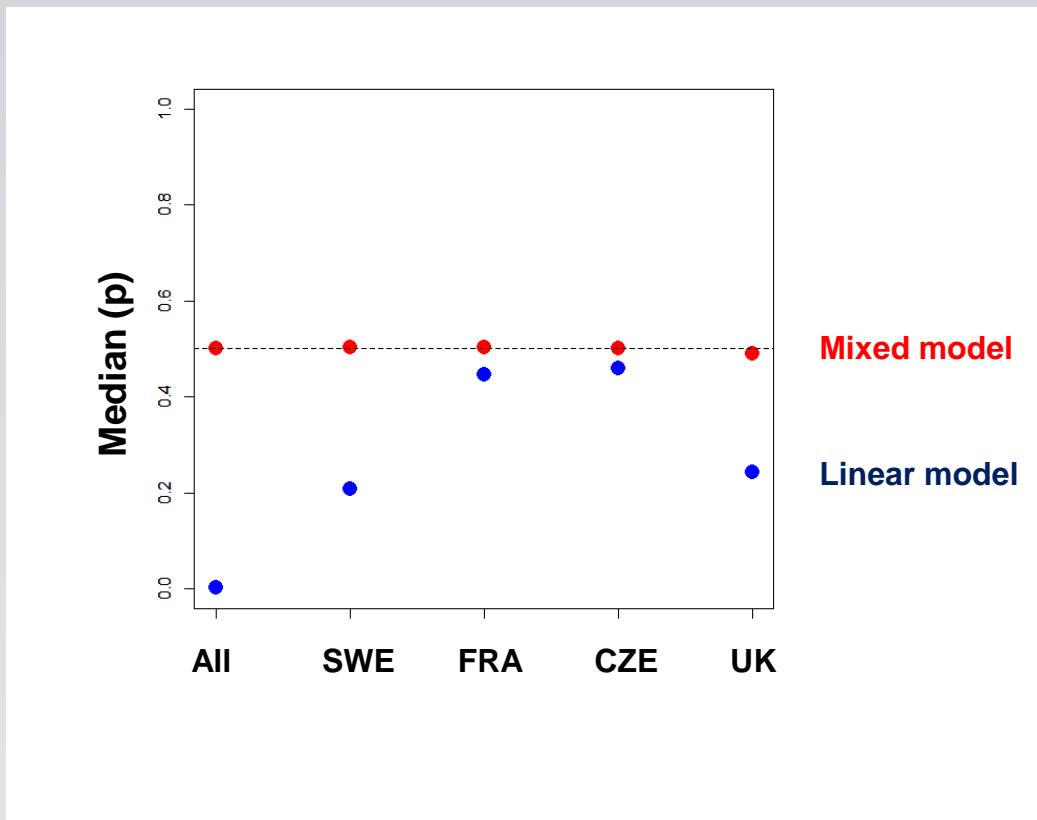
# GWAS in local subsets



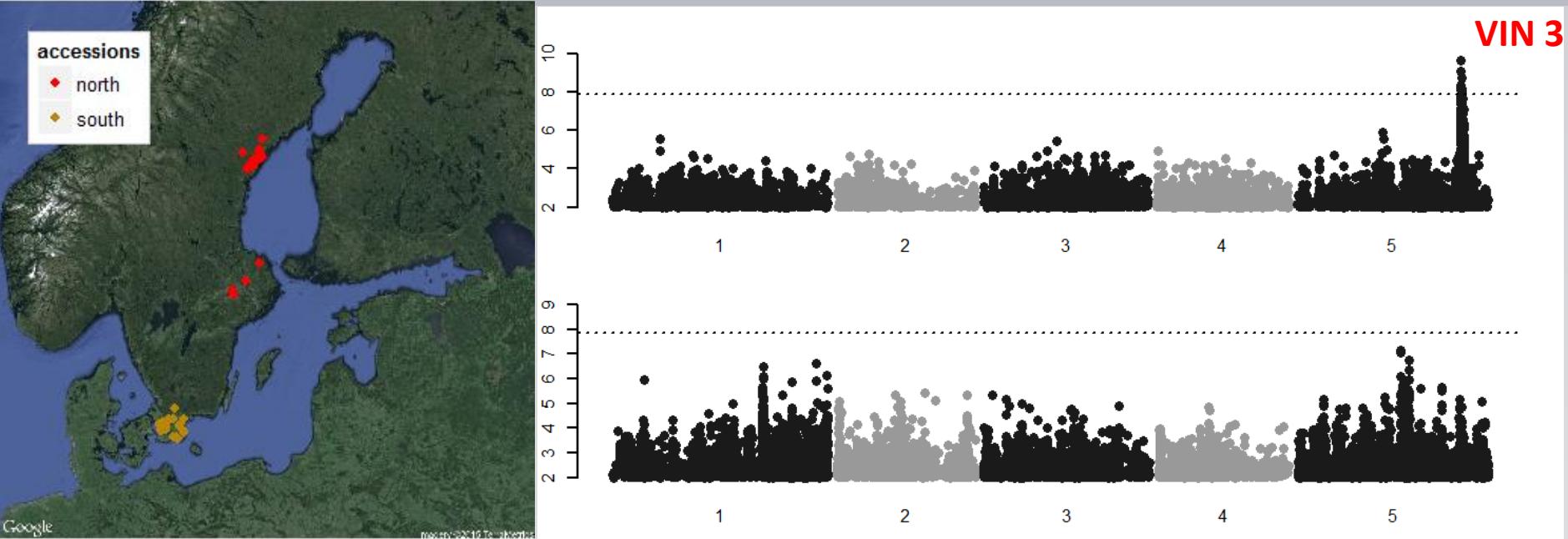
Why ?

- Differences in Allele frequency of the causative marker
- Artefact of the mixed model

# Population structure control in subsets



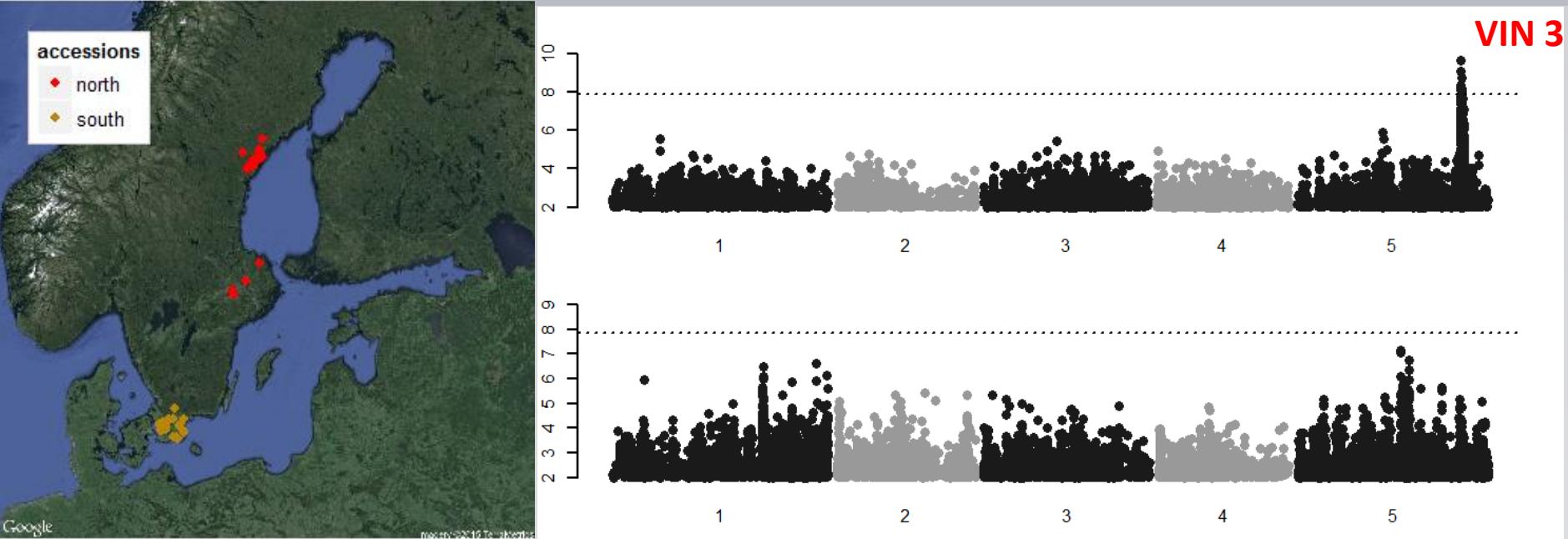
# GWAS in local subsets



Why ?

- Differences in Allele frequency of the causative marker
- Artefact of the mixed model
- The associated SNP is tagging different haplotypes

# GWAS in local subsets



Why ?

- Differences in Allele frequency of the causative marker
- Artefact of the mixed model
- The associated SNP is tagging different haplotypes
- Epistasis

## We cannot test genome-wide SNP by SNP Epistasis !

### Evidence for Network Evolution in an *Arabidopsis* Interactome Map

*Arabidopsis* Interactome Mapping Consortium\*†

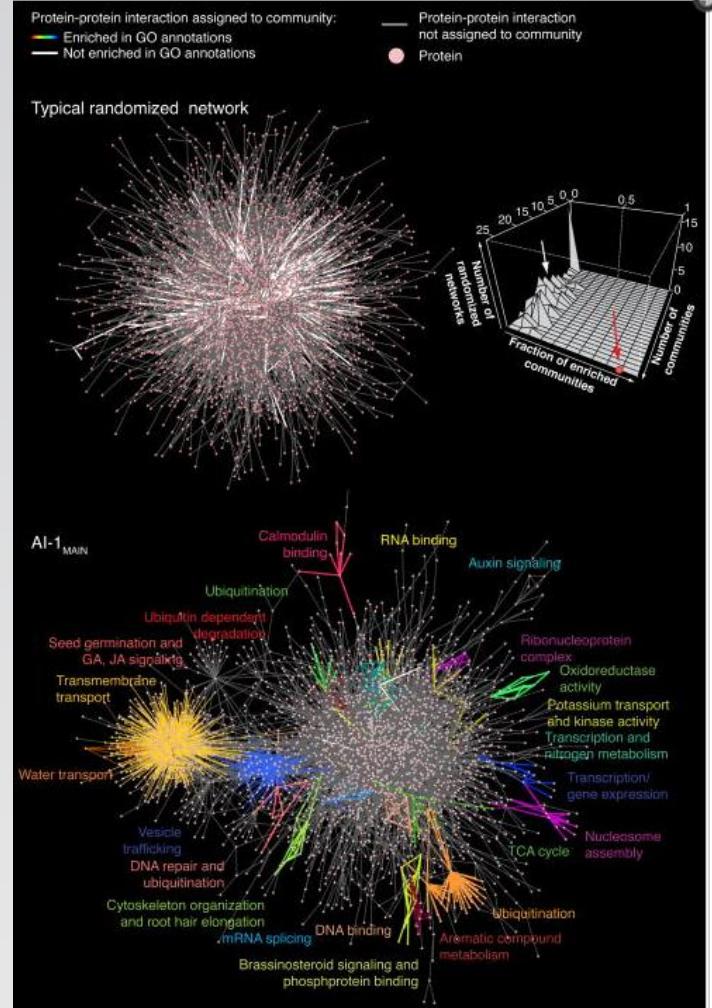
Plants have unique features that evolved in response to their environments and ecosystems. A full account of the complex cellular networks that underlie plant-specific functions is still missing. We describe a proteome-wide binary protein-protein interaction map for the interactome network of the plant *Arabidopsis thaliana* containing about 6200 highly reliable interactions between about 2700 proteins. A global organization of plant biological processes emerges from community analyses of the resulting network, together with large numbers of novel hypothetical functional links between proteins and pathways. We observe a dynamic rewiring of interactions following gene duplication events, providing evidence for a model of evolution acting upon interactome networks. This and future plant interactome maps should facilitate systems approaches to better understand plant biology and improve crops.

Classical genetic and molecular approaches have provided fundamental understanding of processes such as growth control or development and molecular descriptions of genotype-to-phenotype relationships for a varie-

ty of plant systems. Yet, more than 60% of the protein-coding genes of the model plant *Arabidopsis thaliana* (hereafter *Arabidopsis*) remain functionally uncharacterized. Knowledge about the biological organization of macromolecules in complex and dynamic “interactome” networks is lacking for *Arabidopsis* (fig. S1 and tables S1 and S2), depriving us of an understanding of how genotype-to-phenotype relationships are mediated at the systems level (1).

\*All authors with their affiliations and contributions are listed at the end of the paper.

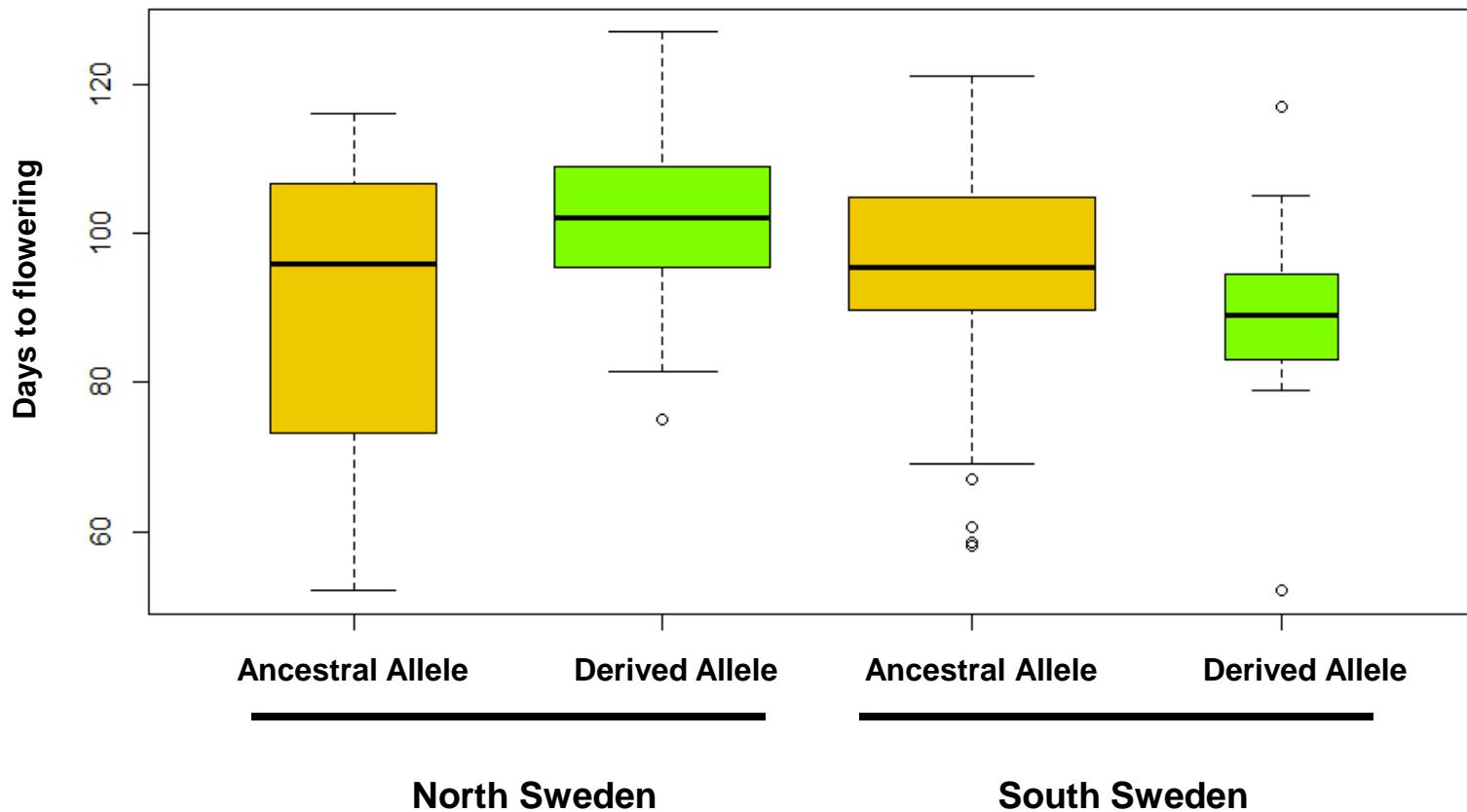
†To whom correspondence should be addressed. E-mail: marc\_vidal@dfci.harvard.edu; ecker@salk.edu; pascal\_braun@dfci.harvard.edu; david\_hill@dfci.harvard.edu



# Effect of the *VIN3* Allele on Flowering time



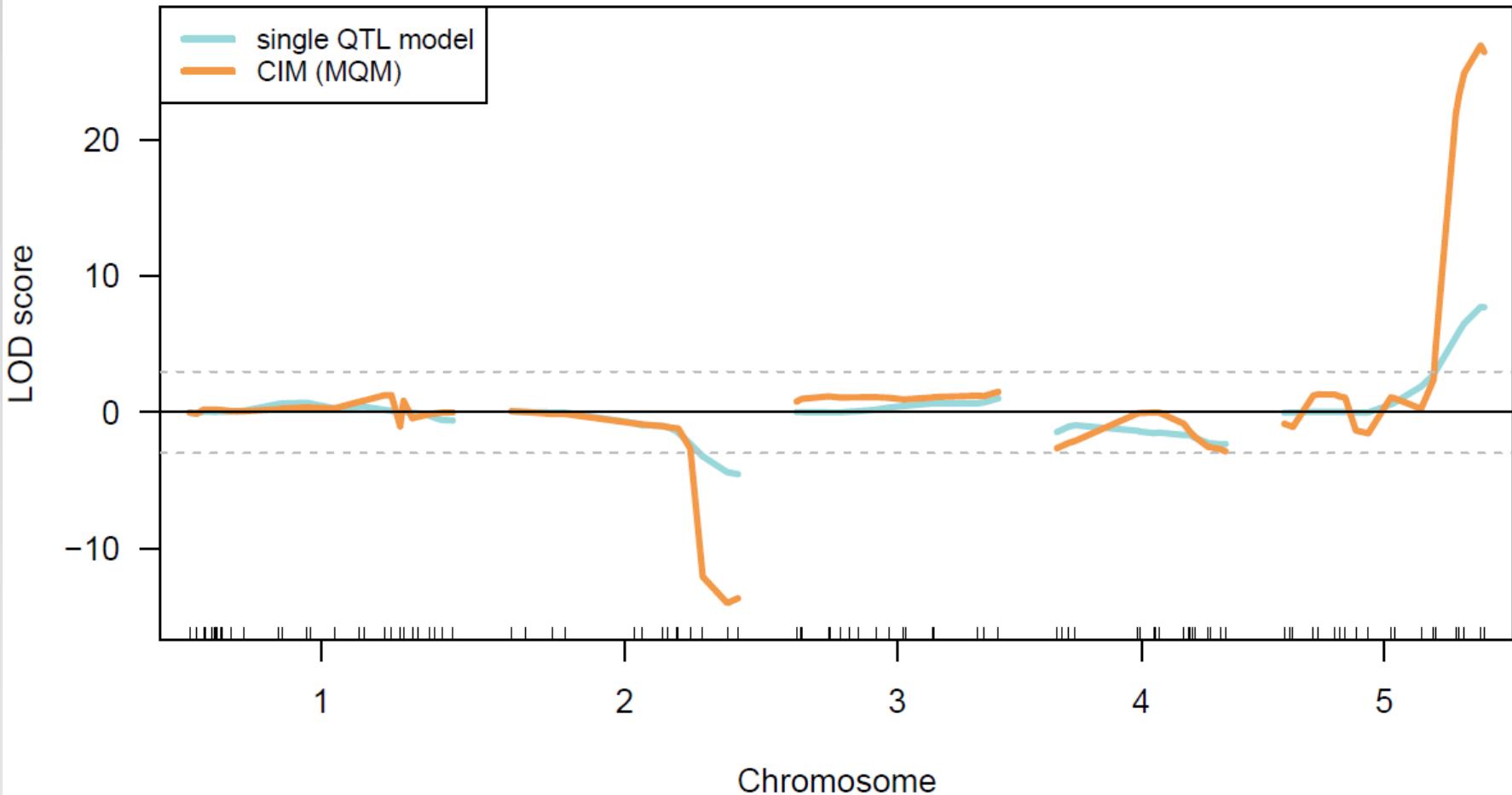
Center for Computational  
and Theoretical Biology



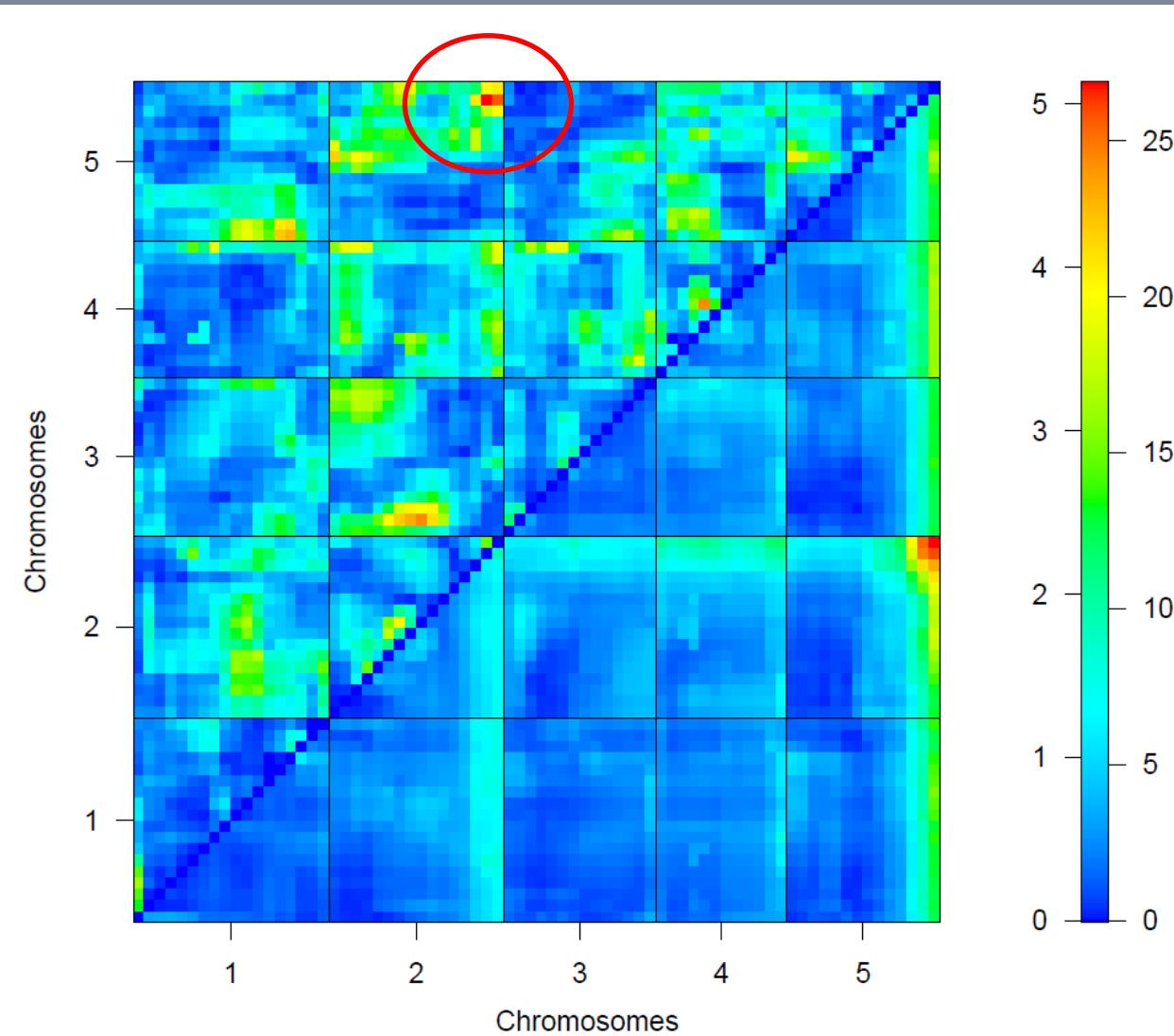
The effect of the respective allele depends on the genetic background

# QTL mapping in crosses

## QTL mapping – 9332x6974 – FT



# QTL mapping in crosses (2d plot)



Epistatic interaction in the QTL analysis

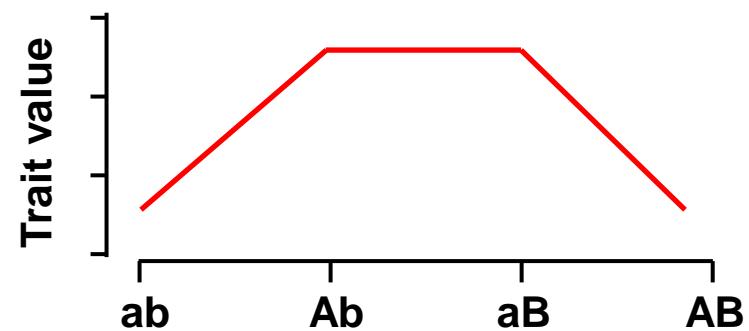
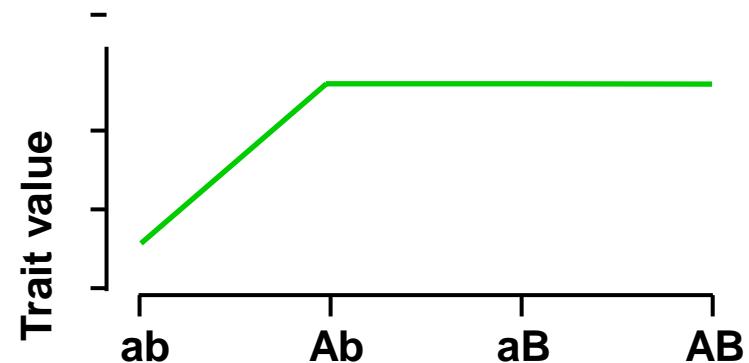
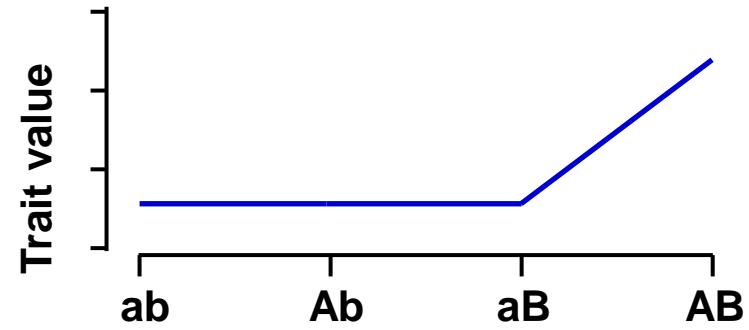
# GWAS models to test for epistasis

## SNP by SNP epistasis

Genome-wide epistasis with the lead SNP in *VIN3*

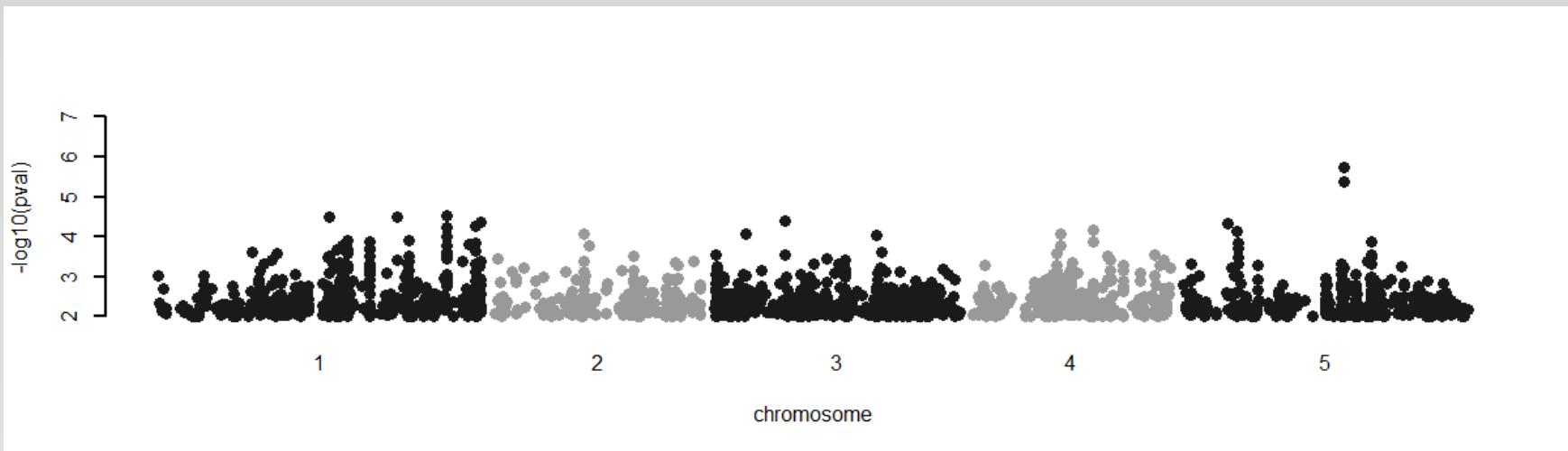
3 different models with two homozygous SNPs

AND, OR and XOR

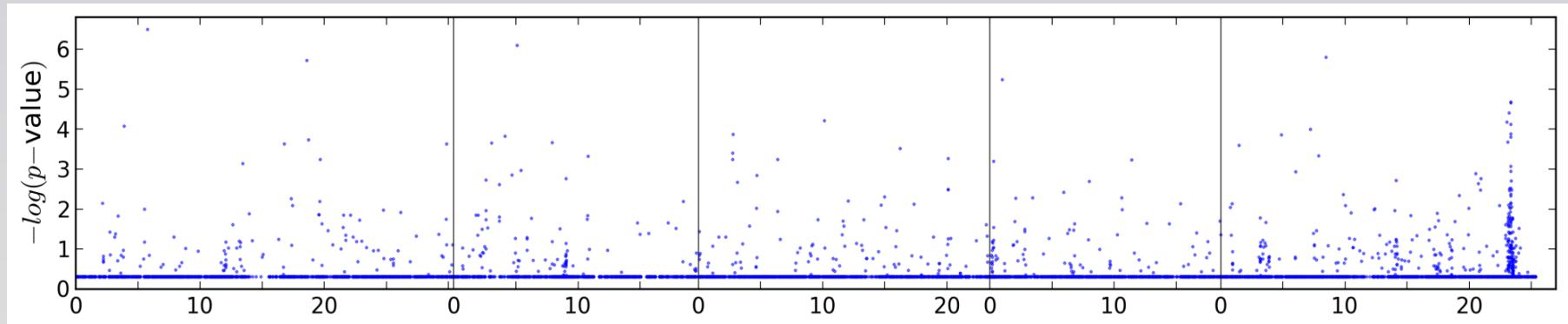


# GWAS models to test for epistasis

## SNP by SNP epistasis



## SNP-by-genomic background epistasis



Different SNPs in *VIN3* seems to interact with the genomic background

**There are different reason why associations don't replicate across different subsets of accessions :**

- **False Positives**
- **No local variation of the causative Allele**
- **Effect depends on the genetic background (Epistasis)**
- **...**

*Invited Paper*

# Genotype–Phenotype Correlation in Cystic Fibrosis: The Role of Modifier Genes

Francesco Salvatore,<sup>1,\*</sup> Olga Scudiero,<sup>1</sup> and Giuseppe Castaldo<sup>1,2</sup>

<sup>1</sup>CEINGE Scarl and Dipartimento di Biochimica e Biotecnologie Mediche Università di Napoli "Federico II," Naples, Italy

<sup>2</sup>Facoltà di Scienze Matematiche, Fisiche e Naturali, Università del Molise, Isernia, Italy

More than 1,000 mutations have been identified in the cystic fibrosis (CF) transmembrane regulator (*CFTR*) disease gene. The impact of these mutations on the protein and the wide spectrum of CF phenotypes prompted a series of Genotype–Phenotype correlation studies. The *CFTR* genotype is invariably correlated with pancreatic status—in about 85% of cases with pancreatic insufficiency and in about 15% of cases with pancreatic sufficiency. The correlations between the *CFTR* genotype and pulmonary, liver, and gastrointestinal expression are debatable. The heterogeneous phenotype in

mutations in  $\alpha$ -1 antitrypsin (*A1AT*) and mannose binding lectin genes were found to be independent risk factors for liver disease in CF patients. The body of evidence available suggests that the variegated CF phenotype results from complex interactions between numerous gene products.

© 2002 Wiley-Liss, Inc.

**KEY WORDS:** cystic fibrosis; modifier genes; Genotype–Phenotype correlation

# TECHNICAL REPORTS



## A mixed-model approach for genome-wide association studies of correlated traits in structured populations

Arthur Korte<sup>1,4</sup>, Bjarni J Vilhjálmsdóttir<sup>1,2,4</sup>, Vincent S  
Magnus Nordborg<sup>1,2</sup>

Genome-wide association studies (GWAS) are a standard approach for studying the genetics of natural variation. A major concern in GWAS is the need to account for the complicated dependence structure of the data, both between loci as well between individuals. Mixed models have emerged as a general and flexible approach for correcting for population structure.

TECHNICAL REPORT



## An efficient multi-locus mixed-model approach for genome-wide association studies in structured populations

Vincent Segura<sup>1,2,4</sup>, Bjarni J Vilhjálmsdóttir<sup>1,3,4</sup>, Alexander Platt<sup>1,3</sup>, Arthur Korte<sup>1</sup>, Ümit Seren<sup>1</sup>, Quan Long<sup>1</sup> & Magnus Nordborg<sup>1,3</sup>

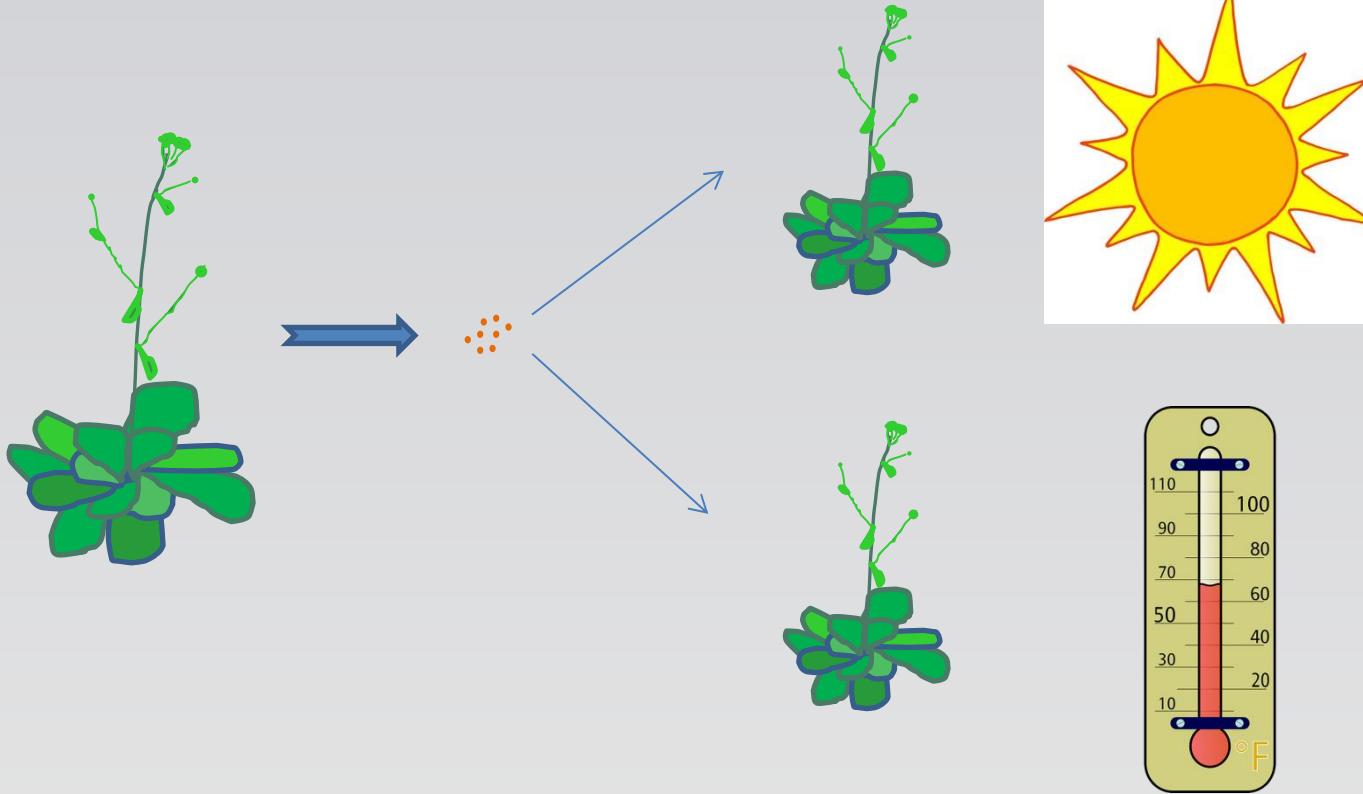
Population structure causes genome-wide linkage disequilibrium between unlinked loci, leading to statistical confounding in genome-wide association studies. Mixed models have been shown to handle the confounding effects of

the cluster memberships and principal-component loadings of individuals, respectively. Whereas these approaches are expected to perform well when the population structure is simple, they may perform poorly when the structure is more complex: for example, when

# Combining traits for a joint GW



Center for Computational  
and Theoretical Biology

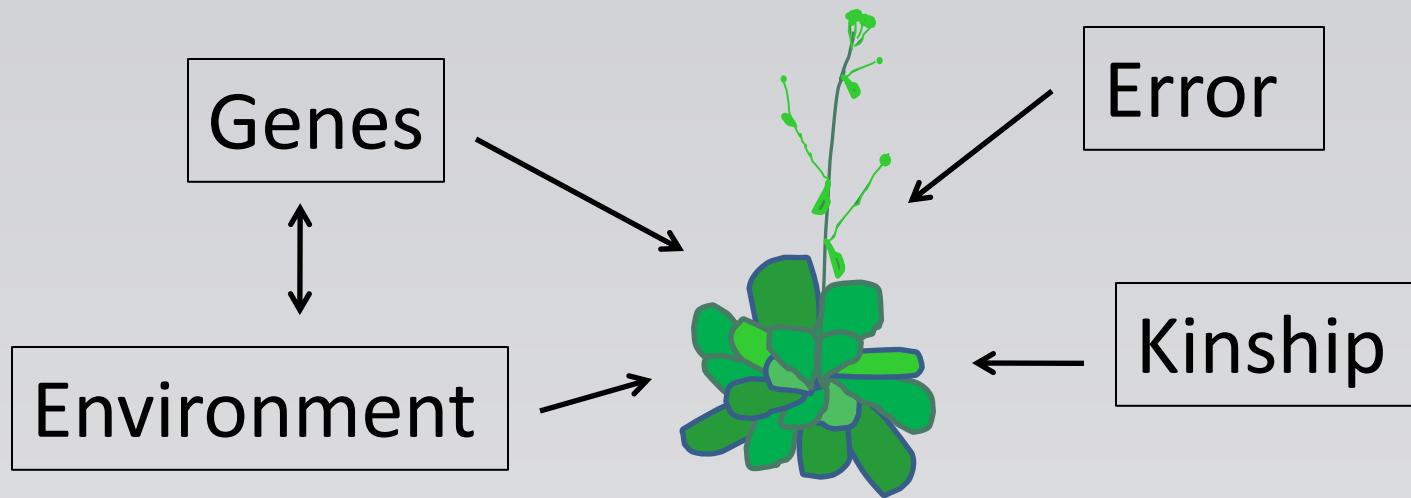


Perform a combined analysis on the data

# Many phenotypes are correlated

	X5_FT10	X6_FT16	X7_FT22	X43_FLC	X44_FRI	X80_LN10	X81_LN16	X82_LN22	X158_Siliq	X159_Siliq	X168_Chlc	X169_Chlc	X281_Stor	X282_Stora
X5_FT10	1	0.821332	0.780659	0.317835	0.280113	0.860099	0.697789	0.710058	-0.11299	-0.25084	0.066188	0.227002	0.313736	0.32493
X6_FT16	0.821332	1	0.882998	0.407995	0.354112	0.785828	0.829433	0.756154	-0.15999	-0.3158	0.094969	0.116744	0.352554	0.336236
X7_FT22	0.780659	0.882998	1	0.432569	0.332482	0.752075	0.767188	0.897305	-0.16028	-0.36557	0.113572	0.108052	0.358229	0.407532
X43_FLC	0.317835	0.407995	0.432569	1	0.399599	0.361573	0.431313	0.431833	0.128897	0.177295	0.060705	0.041051	-0.07754	-0.08906
X44_FRI	0.280113	0.354112	0.332482	0.399599	1	0.248296	0.33493	0.249646	-0.16378	-0.079	0.13947	0.031799	0.190994	0.294031
X80_LN10	0.860099	0.785828	0.752075	0.361573	0.248296	1	0.667059	0.657481	-0.17255	-0.31484	0.061653	0.207983	0.329016	0.3094
X81_LN16	0.697789	0.829433	0.767188	0.431313	0.33493	0.667059	1	0.6967	-0.10834	-0.17035	0.102676	0.141131	0.34098	0.296626
X82_LN22	0.710058	0.756154	0.897305	0.431833	0.249646	0.657481	0.6967	1	-0.19009	-0.32481	0.15016	0.09012	0.286067	0.363039
X158_Silique.16	-0.11299	-0.15999	-0.16028	0.128897	-0.16378	-0.17255	-0.10834	-0.19009	1	0.623108	0.036657	0.099576	-0.09551	-0.13473
X159_Silique.22	-0.25084	-0.3158	-0.36557	0.177295	-0.079	-0.31484	-0.17035	-0.32481	0.623108	1	0.070875	-0.11316	-0.3201	-0.30419
X168_Chlorosis.16	0.066188	0.094969	0.113572	0.060705	0.13947	0.061653	0.102676	0.15016	0.036657	0.070875	1	0.357333	-0.07354	0.010447
X169_Chlorosis.22	0.227002	0.116744	0.108052	0.041051	0.031799	0.207983	0.141131	0.09012	0.099576	-0.11316	0.357333	1	0.006146	0.077456
X281_Storage.7.days	0.313736	0.352554	0.358229	-0.07754	0.190994	0.329016	0.34098	0.286067	-0.09551	-0.3201	-0.07354	0.006146	1	0.844788
X282_Storage.28.days	0.32493	0.336236	0.407532	-0.08906	0.294031	0.3094	0.296626	0.363039	-0.13473	-0.30419	0.010447	0.077456	0.844788	1

Correlation of two traits can be due to  
shared genetics or a shared environment



$$Y = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \beta_0 + \beta_1 T + \beta_2 G + \beta_3 G^* T + u + \varepsilon$$

$$H_0 : Y = T$$

$$H_1 : Y = T + G$$

$$H_2 : Y = T + G + GT$$

Full p-value :

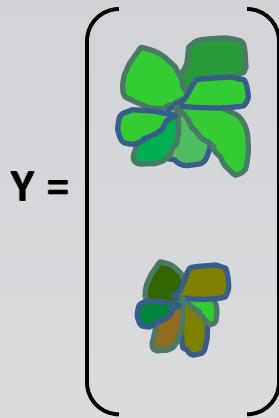
Trait specific p-value :

Trait common p-value :

GLS between  $H_2$  and  $H_0$

GLS between  $H_2$  and  $H_1$

GLS between  $H_1$  and  $H_0$

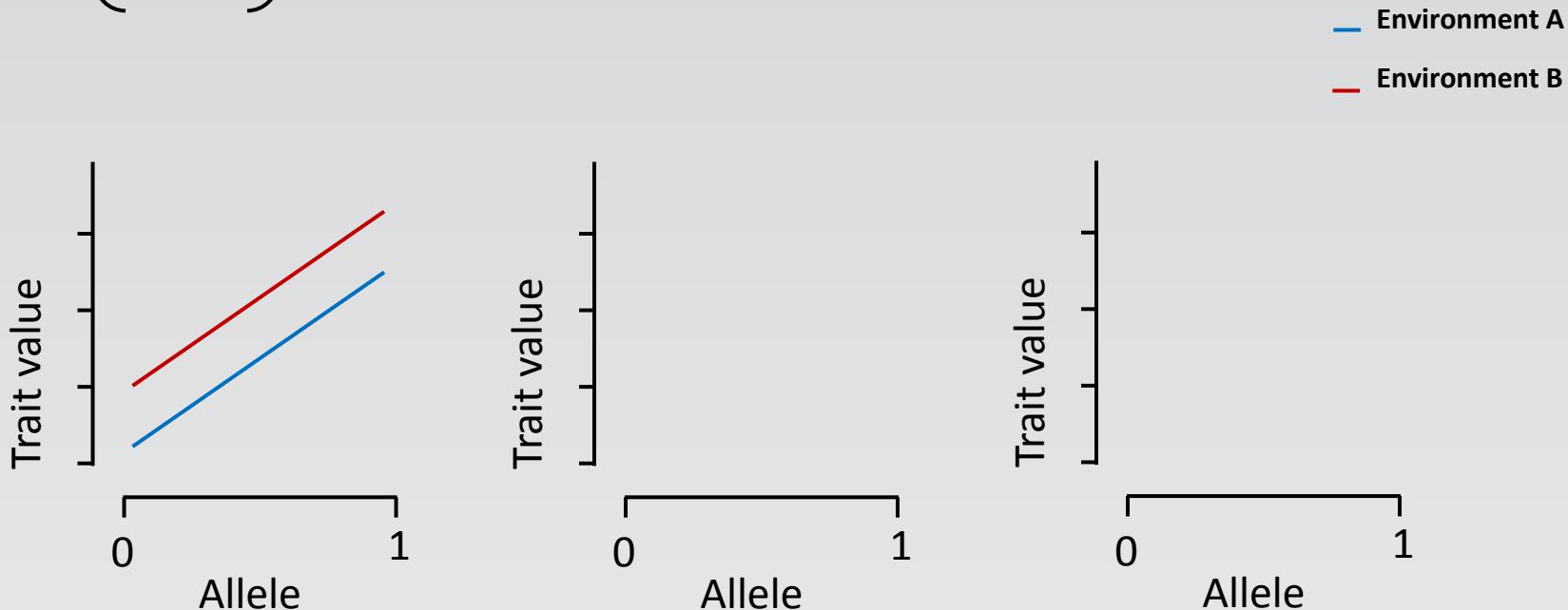


3 different tests :

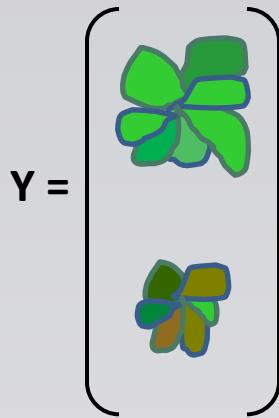
**Trait common** : *same effect on the trait in both environments*

**Trait specific** : *effect only in one environment /  
opposite effect in the two environments*

**Full test** : *combination of common and specific effect*



(same trait in two environments | longitudinal data)



$\gamma =$

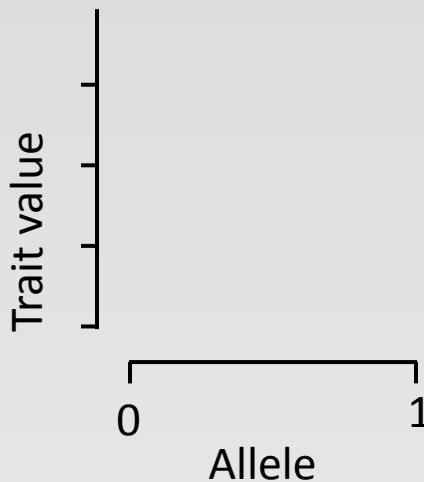
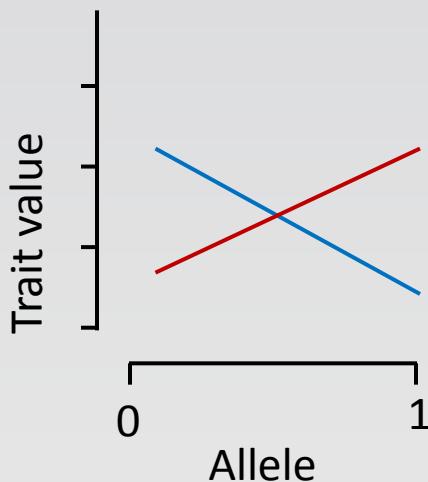
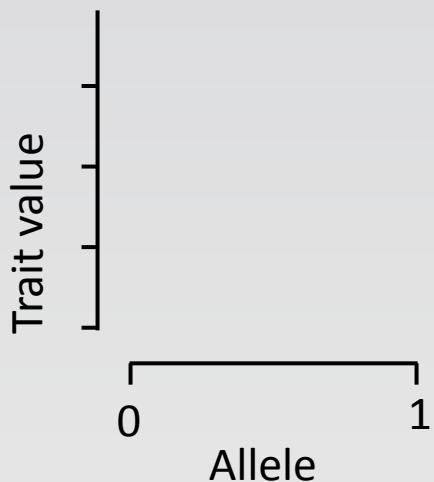
3 different tests :

**Trait common** : *same effect on the trait in both environments*

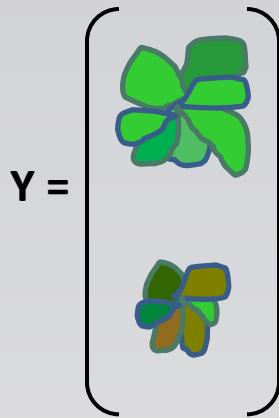
**Trait specific** : *effect only in one environment /  
opposite effect in the two environments*

**Full test** : combination of common and specific effect

Environment A  
Environment B



(same trait in two environments | longitudinal data)



$\mathbf{Y} =$

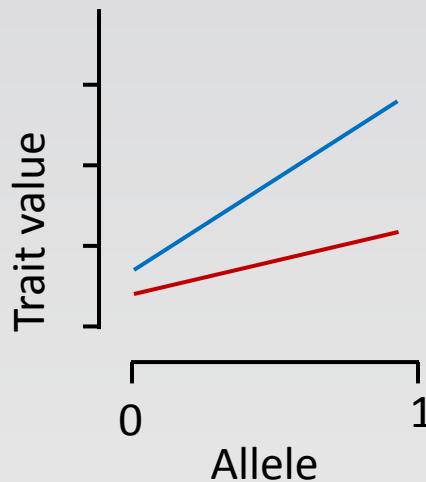
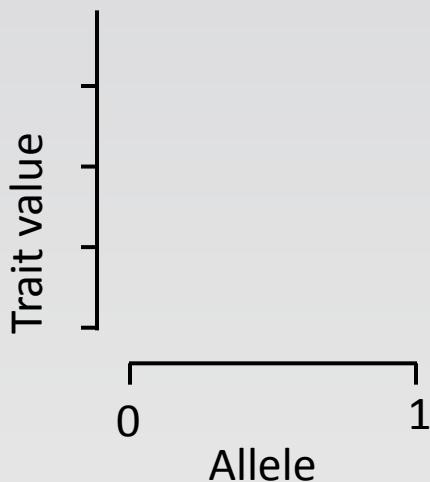
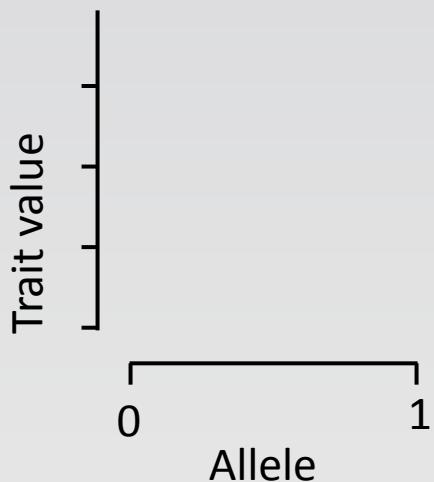
3 different tests :

**Trait common** : *same effect on the trait in both environments*

**Trait specific** : *effect only in one environment /  
opposite effect in the two environments*

**Full test** : combination of common and specific effect

— Environment A  
— Environment B



(same trait in two environments | longitudinal data)

### Mixed model for a pair of correlated traits

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \beta_0 + \mathbf{T}\beta_1 + \mathbf{G}\beta_2 + \mathbf{G}\mathbf{T}\beta_3 + \mathbf{Z}\mathbf{u} + \epsilon,$$

$$\text{var}(\mathbf{u}) = \begin{bmatrix} \sigma_{g11}^2 & \sigma_{g12} \\ \sigma_{g21} & \sigma_{g22} \end{bmatrix} \otimes \mathbf{K}, \text{var}(\epsilon) = \begin{bmatrix} \sigma_{\epsilon11}^2 & \sigma_{\epsilon12} \\ \sigma_{\epsilon21} & \sigma_{\epsilon22} \end{bmatrix} \otimes \mathbf{I}$$

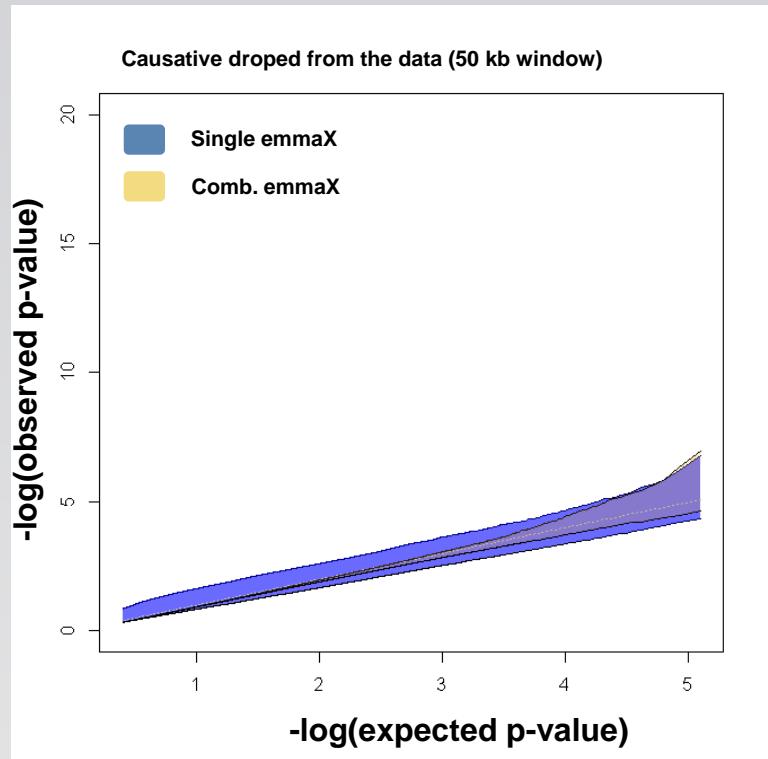
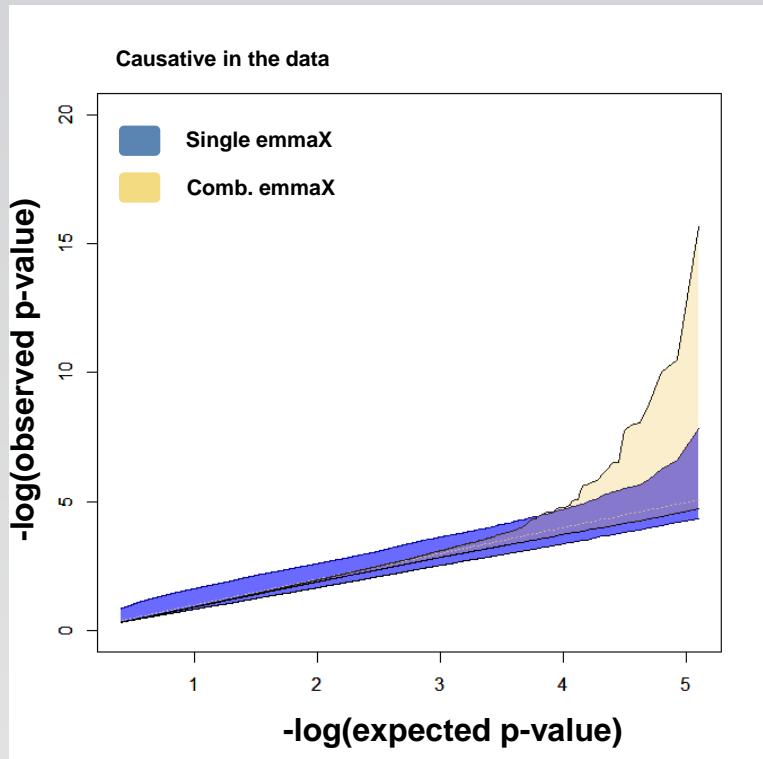
Under the phenotypic model assumptions described in Kang et al. (2008) which yield K to be the IBS kinship matrix, and assuming the causal effects have a genetic correlation  $\rho$ , but are otherwise independent (no environmental correlation), the genetic correlation turns out to be:

$$\rho_g = \frac{\rho_{\mathbf{y}_1, \mathbf{y}_2} \sqrt{(\sigma_{g1}^2 + \sigma_{e1}^2)(\sigma_{g2}^2 + \sigma_{e2}^2)}}{\sigma_{g1}\sigma_{g2}} = \frac{\rho_{\mathbf{y}_1, \mathbf{y}_2}}{h_1 h_2}$$

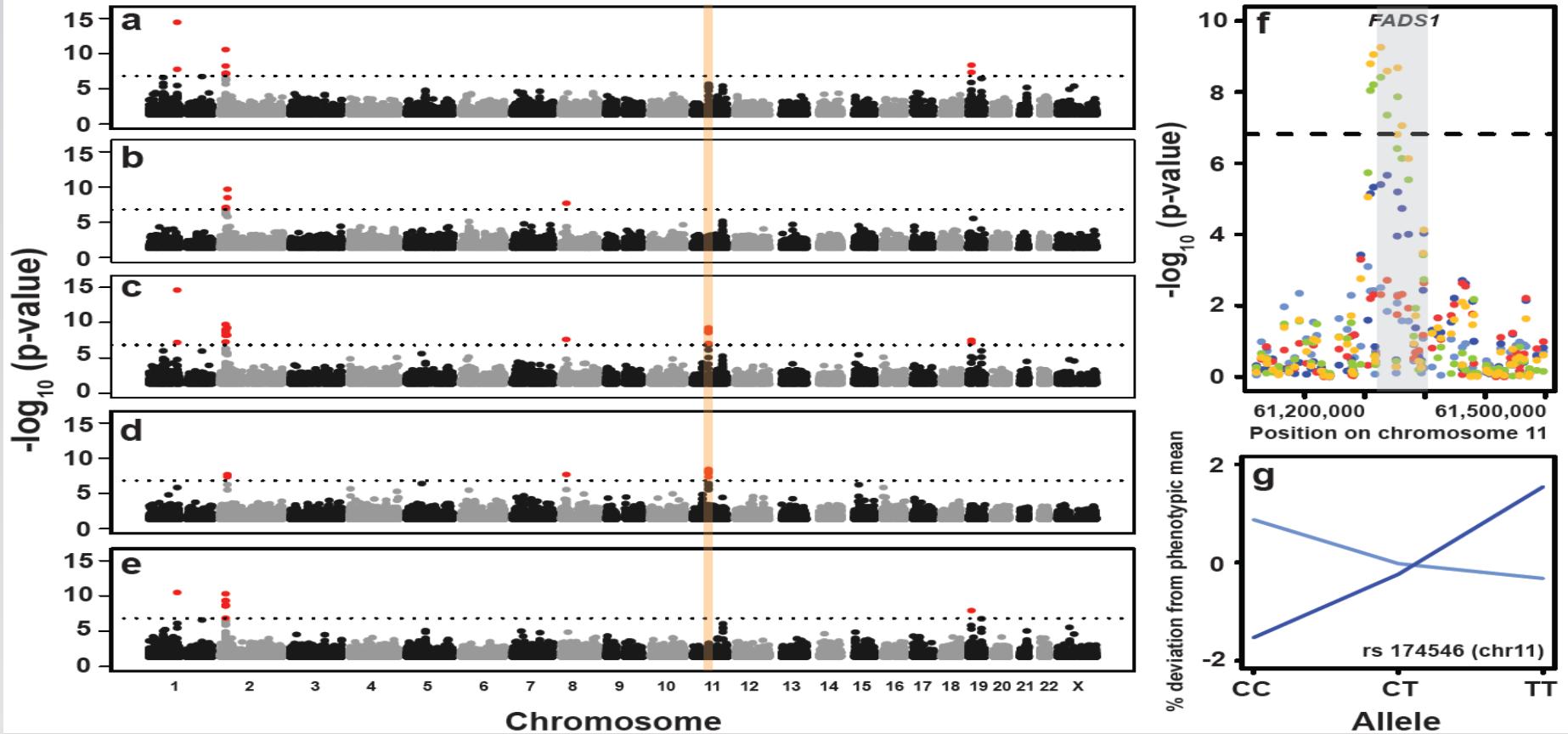
# Simulations

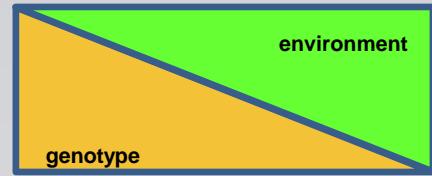


Center for Computational  
and Theoretical Biology



# MTMM (Multi-trait mixed model)





# LETTER

doi:10.1038/nature10781

## Genetic contributions to stability and change in intelligence from childhood to old age

Ian J. Deary<sup>1,2\*</sup>, Jian Yang<sup>3\*</sup>, Gail Davies<sup>1,2</sup>, Sarah E. Harris<sup>2,4</sup>, Albert Tenesa<sup>4,5</sup>, David Liewald<sup>1,2</sup>, Michelle Luciano<sup>1,2</sup>, Lorna M. Lopez<sup>1,2</sup>, Alan J. Gow<sup>1,2</sup>, Janie Corley<sup>1</sup>, Paul Redmond<sup>1</sup>, Helen C. Fox<sup>6</sup>, Suzanne J. Rowe<sup>5</sup>, Paul Haggarty<sup>7</sup>, Geraldine McNeill<sup>6</sup>, Michael E. Goddard<sup>8</sup>, David J. Porteous<sup>2,4</sup>, Lawrence J. Whalley<sup>6</sup>, John M. Starr<sup>2,9</sup> & Peter M. Visscher<sup>2,3,10,11\*</sup>

Understanding the determinants of healthy mental ageing is a priority for society today<sup>1,2</sup>. So far, we know that intelligence differences show high stability from childhood to old age<sup>3,4</sup> and there are estimates of the genetic contribution to intelligence at different ages<sup>5,6</sup>. However, attempts to discover whether genetic causes contribute to differences in cognitive ageing have been relatively uninformative<sup>7–10</sup>. Here we provide an estimate of the genetic and environmental contributions to stability and change in intelligence across most of the human lifetime. We used genome-wide single nucleotide polymorphism (SNP) data from 1,940 unrelated individuals whose intelligence was measured in childhood (age 11 years) and again in old age (age 65, 70 or 79 years)<sup>11,12</sup>. We use a statistical method that allows genetic (co)variance to be estimated from SNP data on unrelated individuals<sup>13–17</sup>. We estimate that causal genetic variants in linkage disequilibrium with common SNPs account for 0.24 of the variation in cognitive ability change from childhood to old age. Using bivariate analysis, we estimate a genetic correlation between intelligence at age 11 years and in old age of 0.62. These estimates, derived from rarely available data on lifetime cognitive measures, warrant the search for genetic causes of cognitive stability and change.

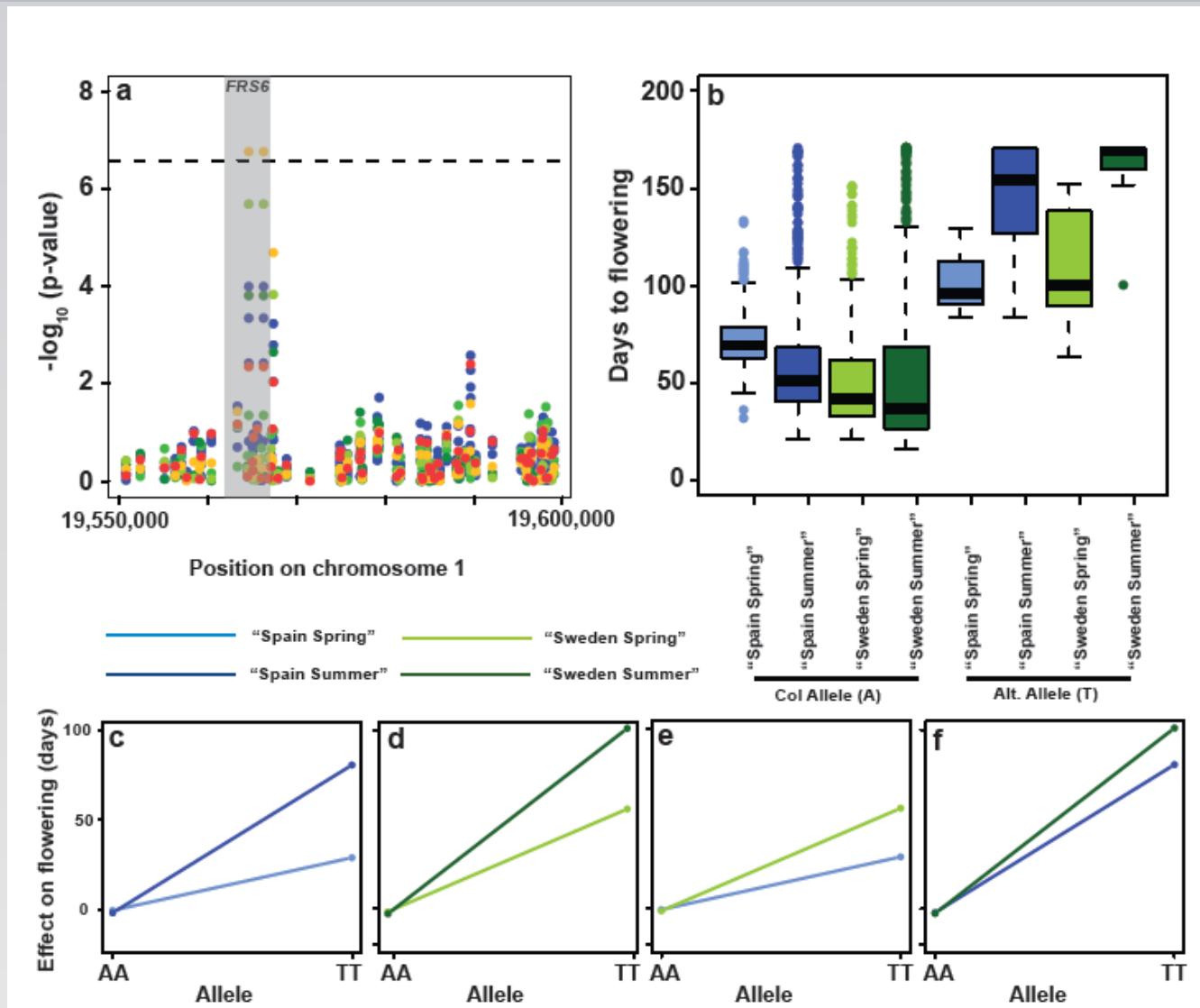
life course are largely unreplicated<sup>22</sup>. Therefore, an important novel contribution would be to partition the covariance between intelligence scores at either end of the human life course into genetic and environmental causes. To address this, the present study applies a new analytical method<sup>13–17</sup> to genome-wide association data from human participants with general cognitive ability test scores in childhood and again in old age.

Participants were members of the Aberdeen Birth Cohort 1936 (ABC1936) and the Lothian Birth Cohorts of 1921 and 1936 (LBC1921, LBC1936)<sup>11,12,17</sup>. They are community-dwelling, surviving

**We estimate a genetic correlation between intelligence at age 11 years and in old age of 0.62.**

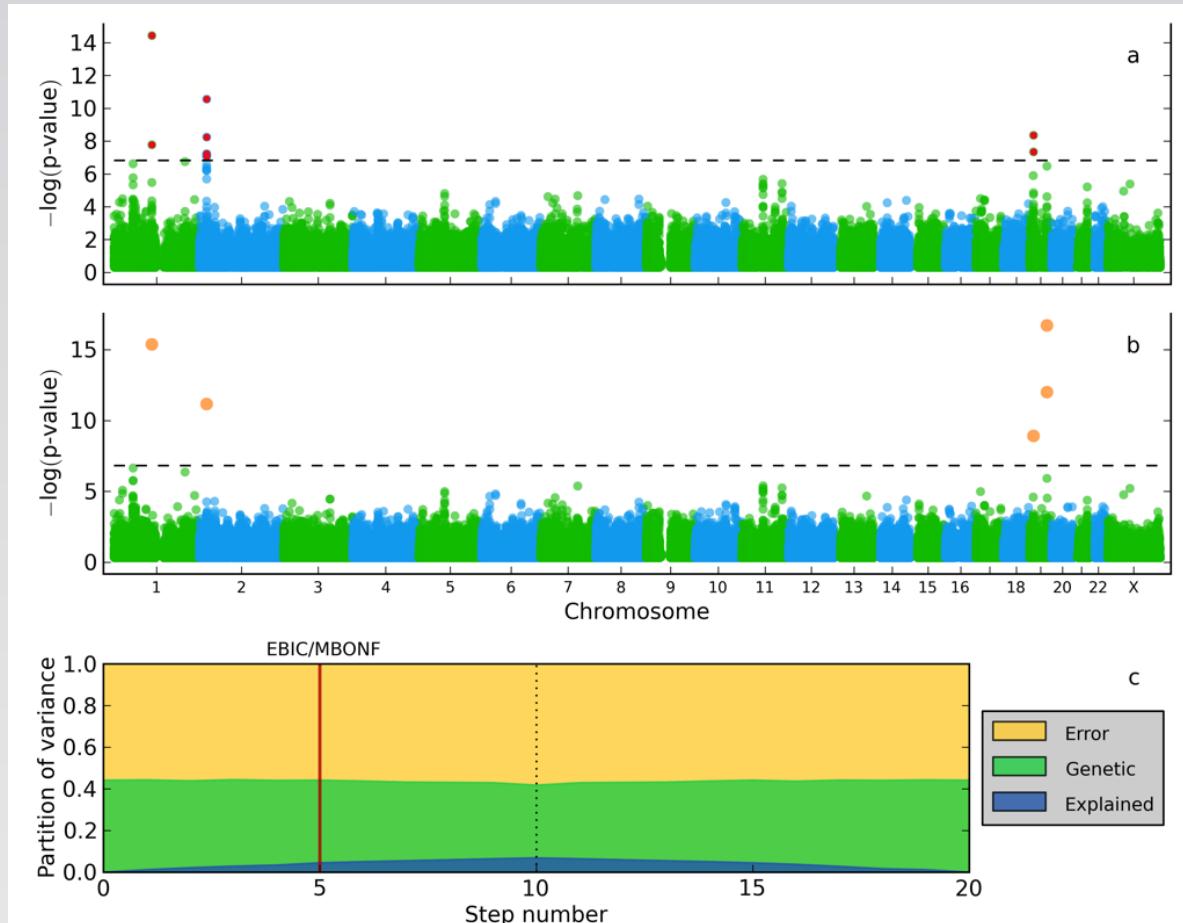
number of diverse cognitive tests. Additionally, the LBC1921 and LBC1936 cohorts re-took the Moray House Test in old age. Thus,

# MTMM for more traits



# Multi Locus Mixed Model

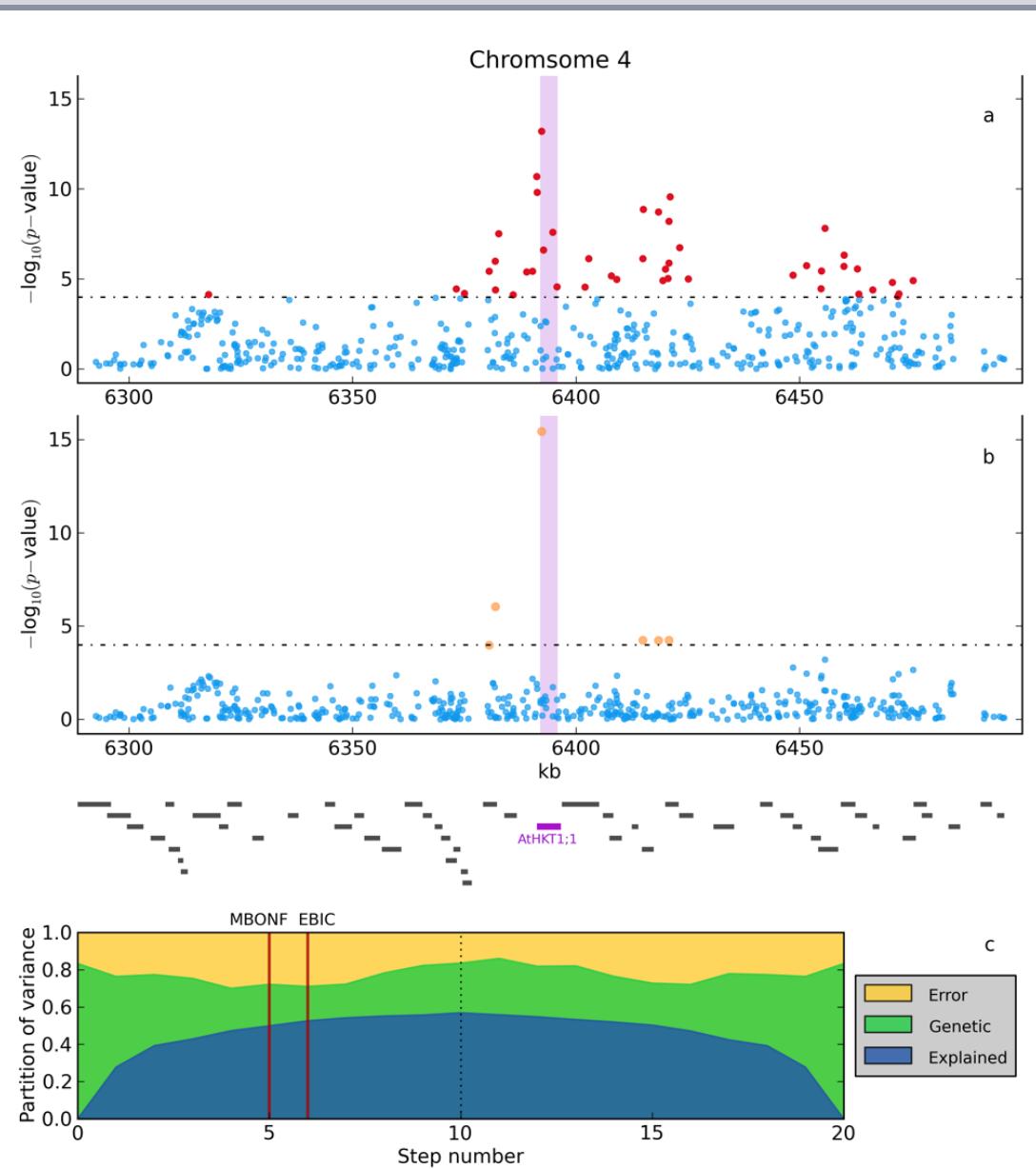
Adding the most significant SNP as Co-factor to the model



# *Arabidopsis* trait : Sodium levels in leaves



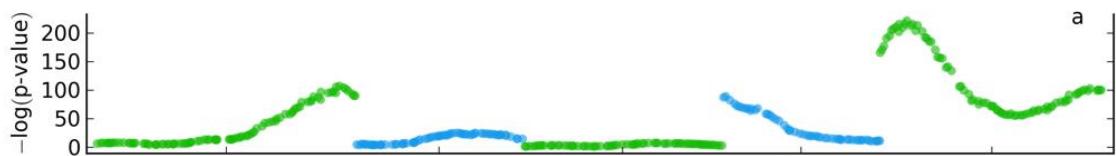
Center for Computational  
and Theoretical Biology



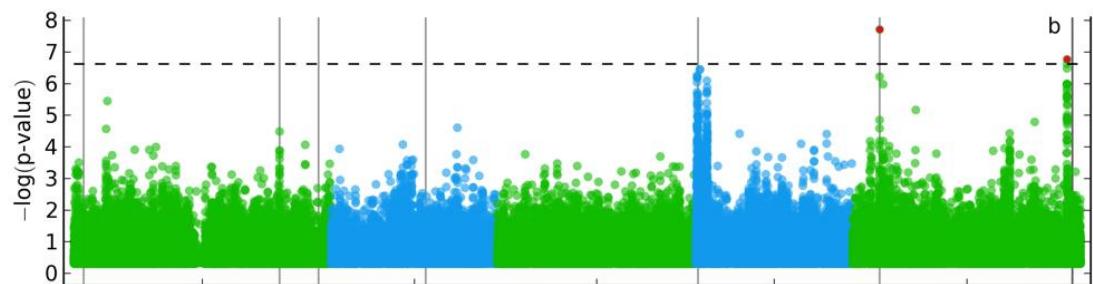
- Bayes factor:  $\text{BF} = \frac{\text{P}(D|M_1)}{\text{P}(D|M_0)} = \frac{\int_{\theta_1} \text{P}(D|\theta_1, M_1)d\theta_1}{\int_{\theta_0} \text{P}(D|\theta_0, M_0)d\theta_0}$ ,
  - If  $\pi$  is the prior probability that a SNP is causal then the posterior odds and posterior probability of associations are  $\text{PO} = \text{BF} \cdot \pi / (1 - \pi)$        $\text{PPA} = \text{PO} / (1 + \text{PO})$
- Following Kass and Raftery (1995), we can (roughly) approximate the Bayes factor which only depends on the likelihood and the degrees of freedoms as:

$$\log(\text{ABF}) = \log(\text{P}(D|\beta, M_1)) - \log(\text{P}(D|\beta, M_0)) - \frac{1}{2}(d_1 - d_0)\log(n)$$

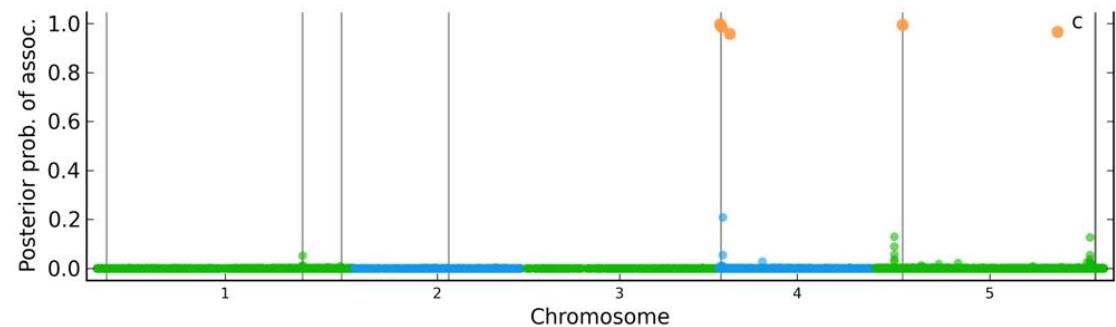
# Bayesian analysis combining the Linkage and GWAS study



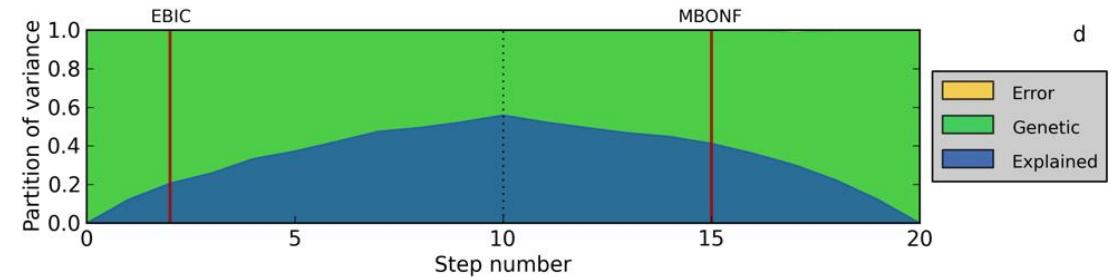
Linkage mapping score



EMMAX scan



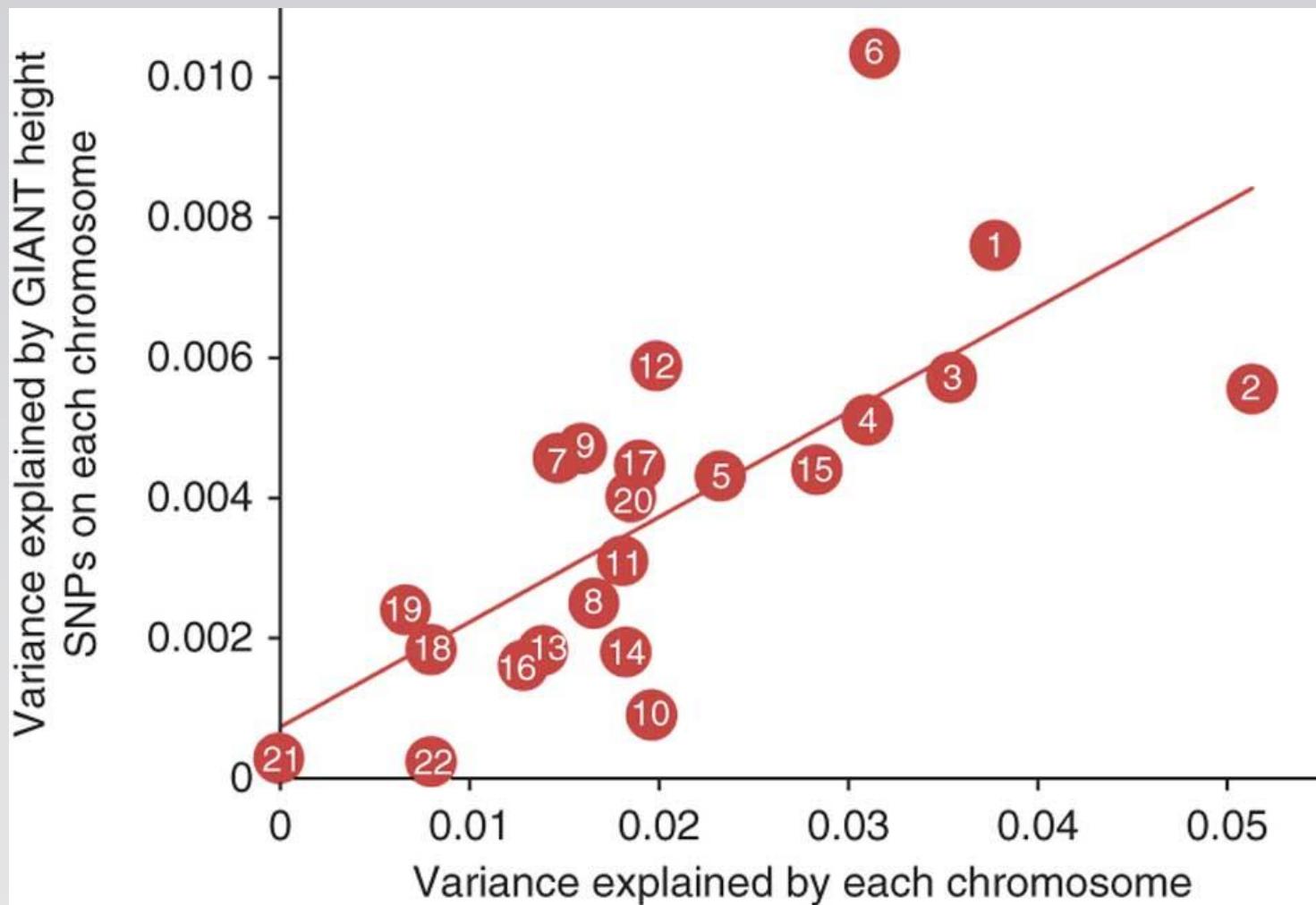
Bayesian analysis combining the  
Linkage and the GWAS study



# Partitioning of the phenotypic variance



Center for Computational  
and Theoretical Biology



# Analysis of RNAseq data



Center for Computational  
and Theoretical Biology

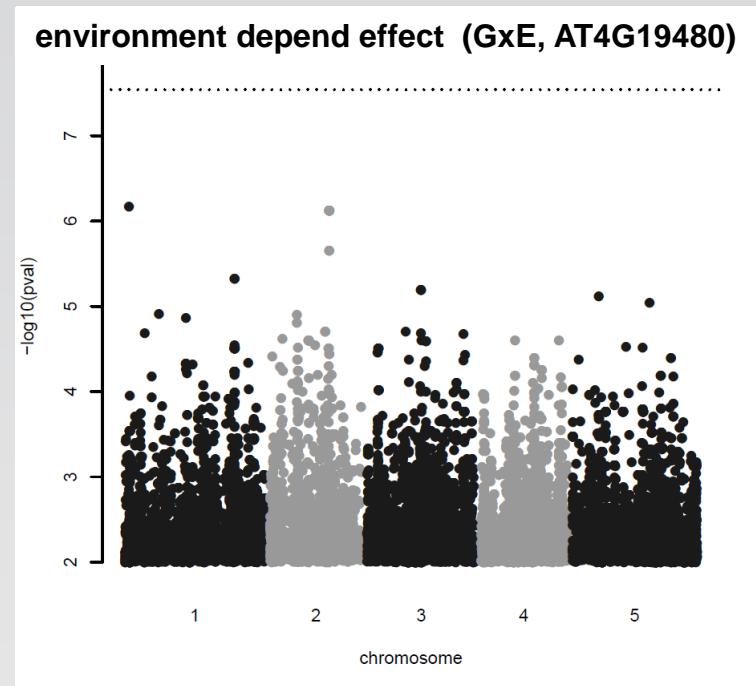
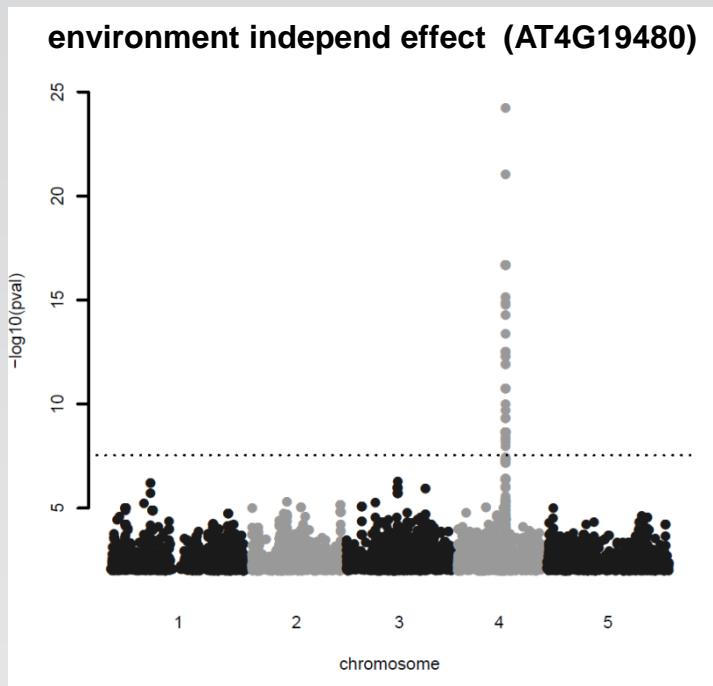
RNAseq data for 80 accessions grown in two different environments



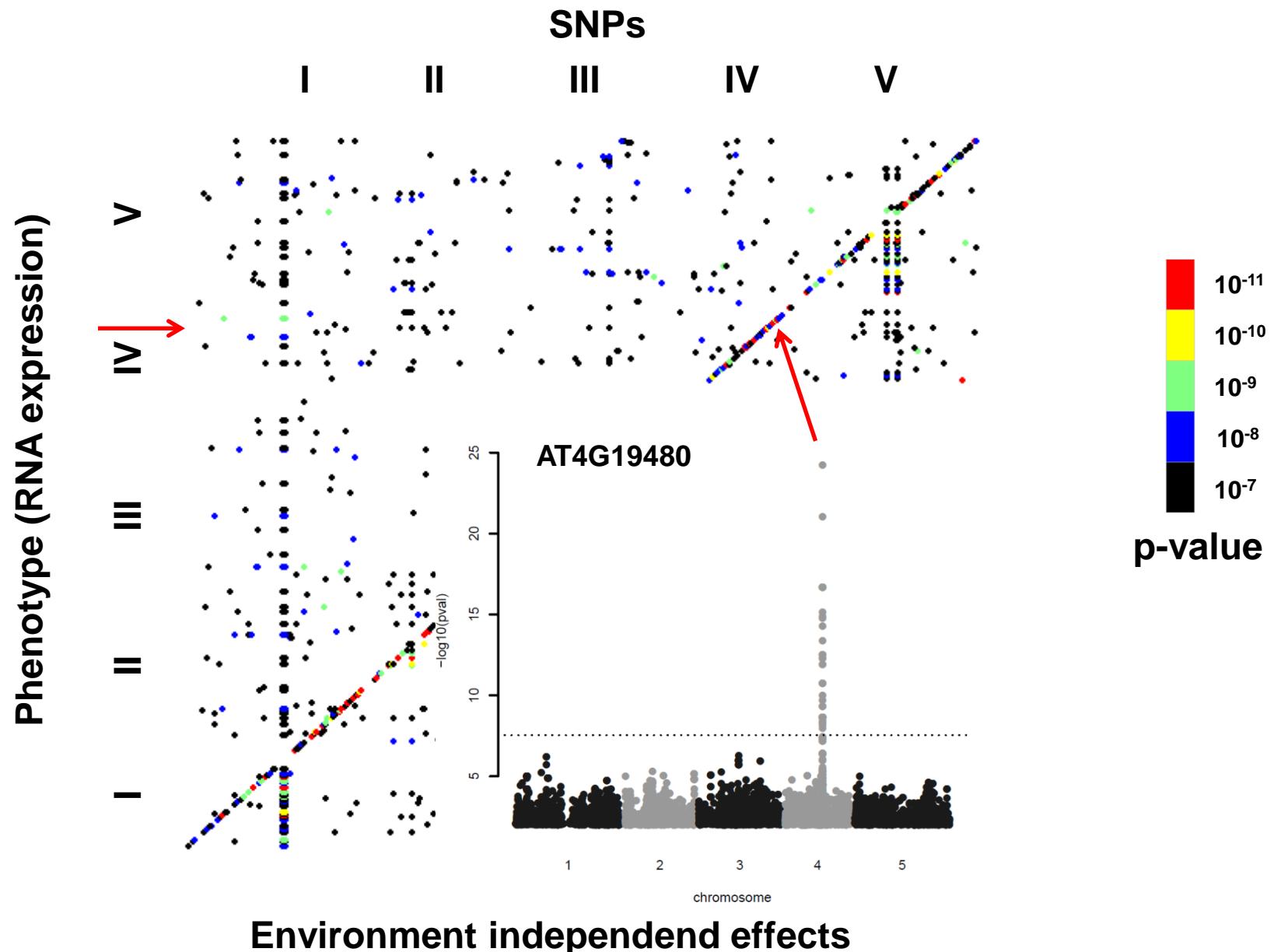
~ 56,000 GWAS !



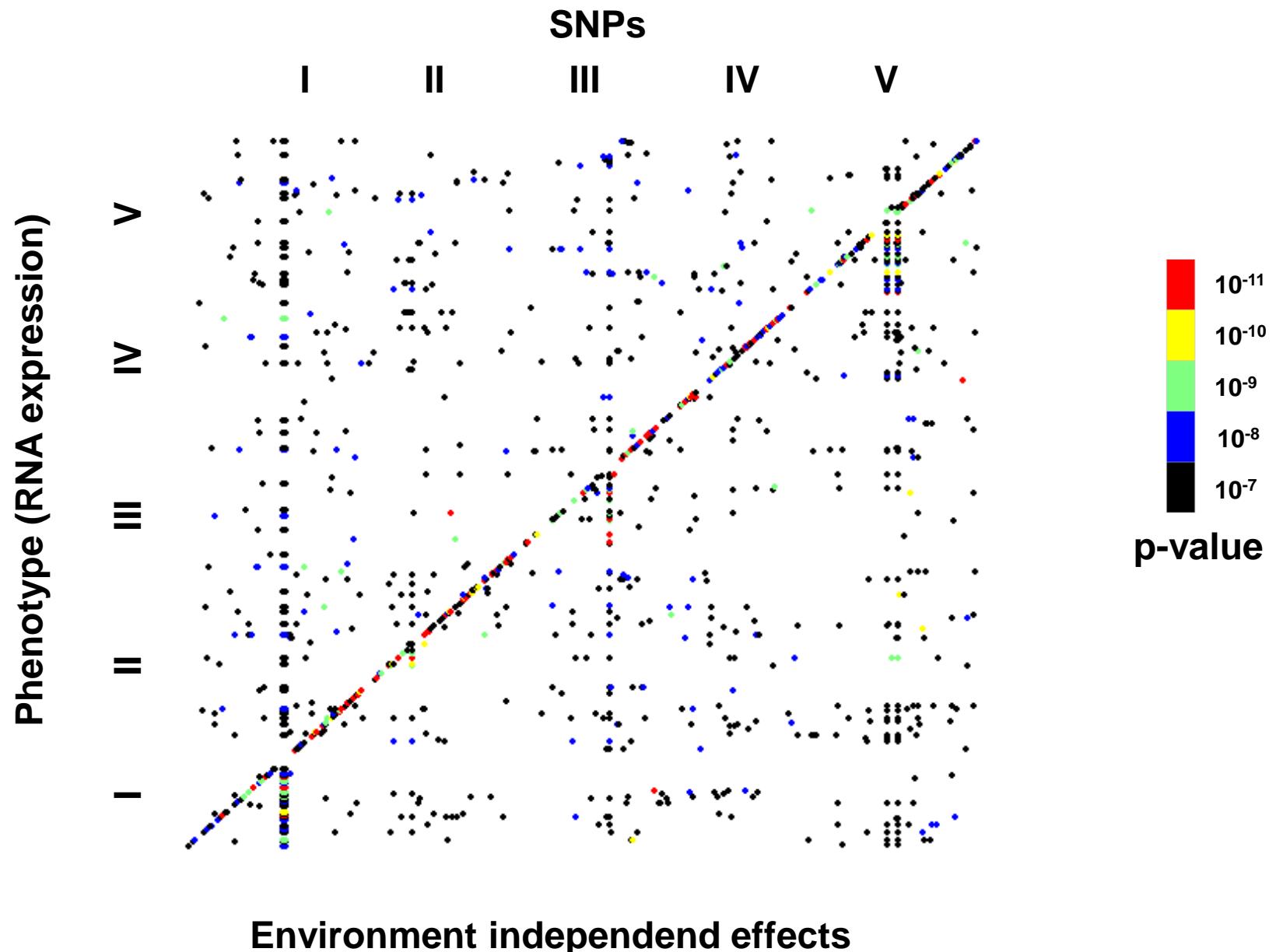
~ 4700 h computational time  
~ 2 d on a cluster with 100 nodes



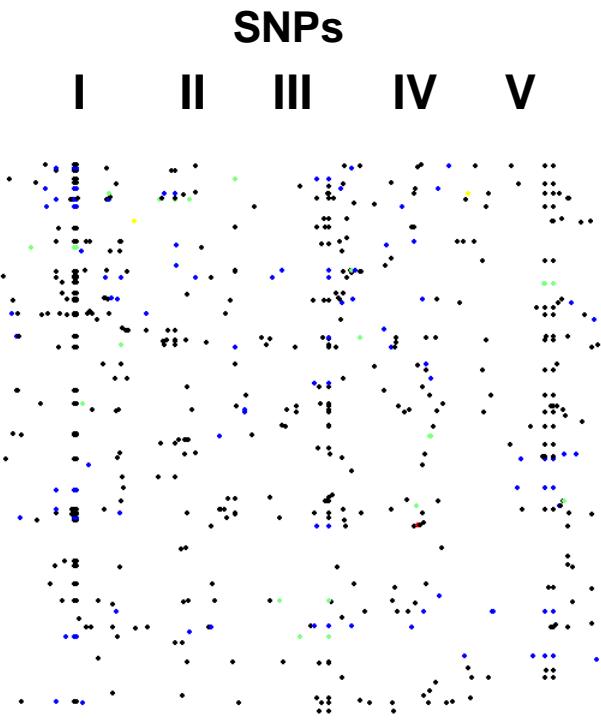
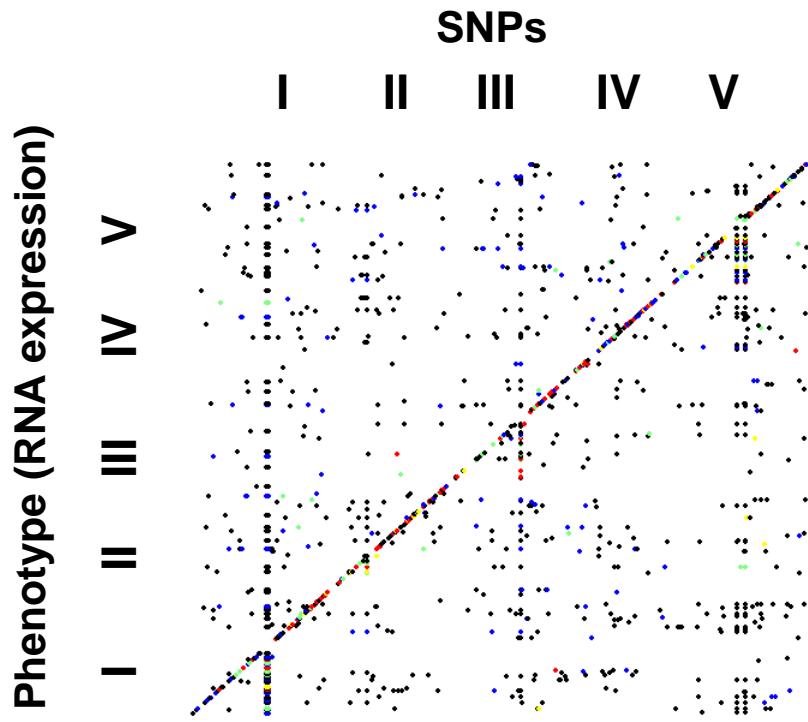
# Summary of MTMM on RNAseq data



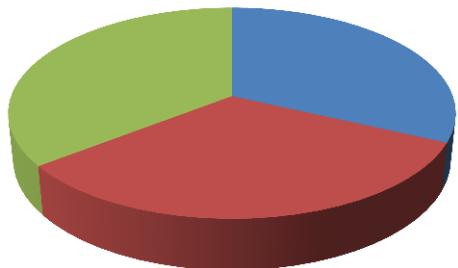
# Summary of MTMM on RNAseq data



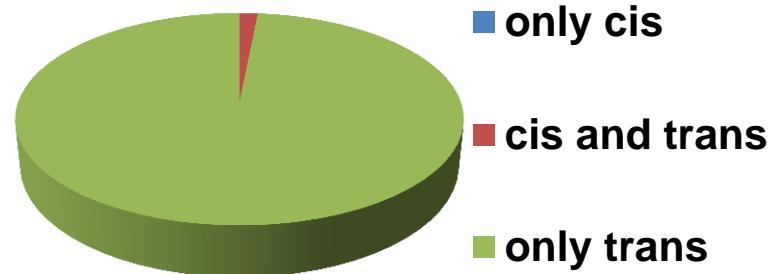
# Summary of MTMM on RNAseq data



**Environment independent**

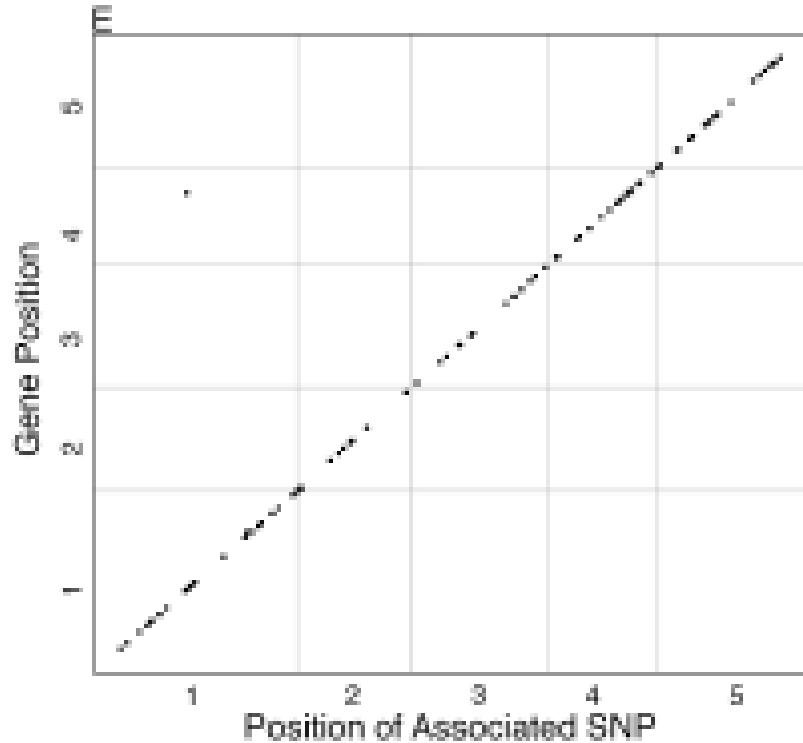


**G x E**

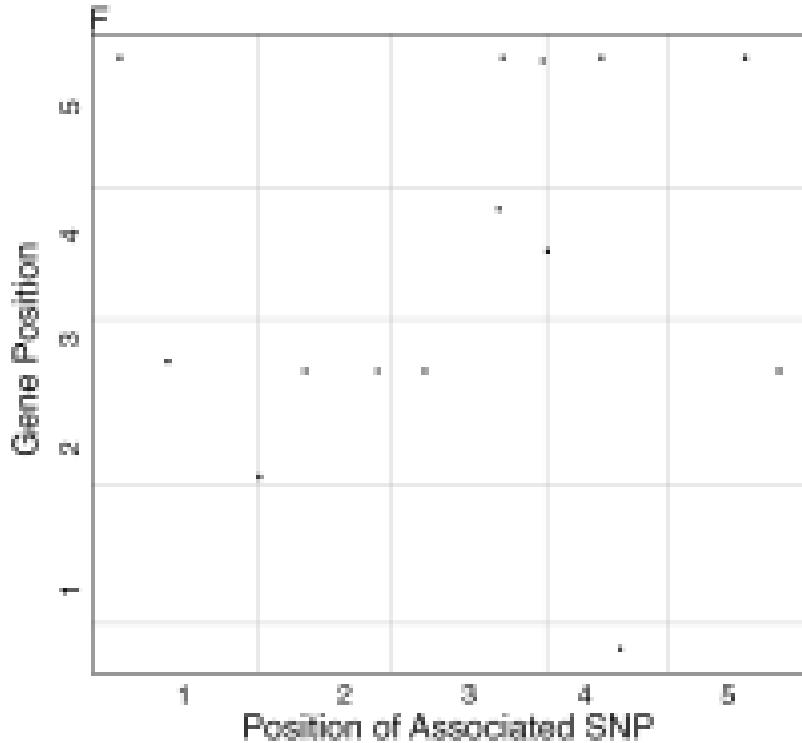


# Summary of MTMM on RNAseq data

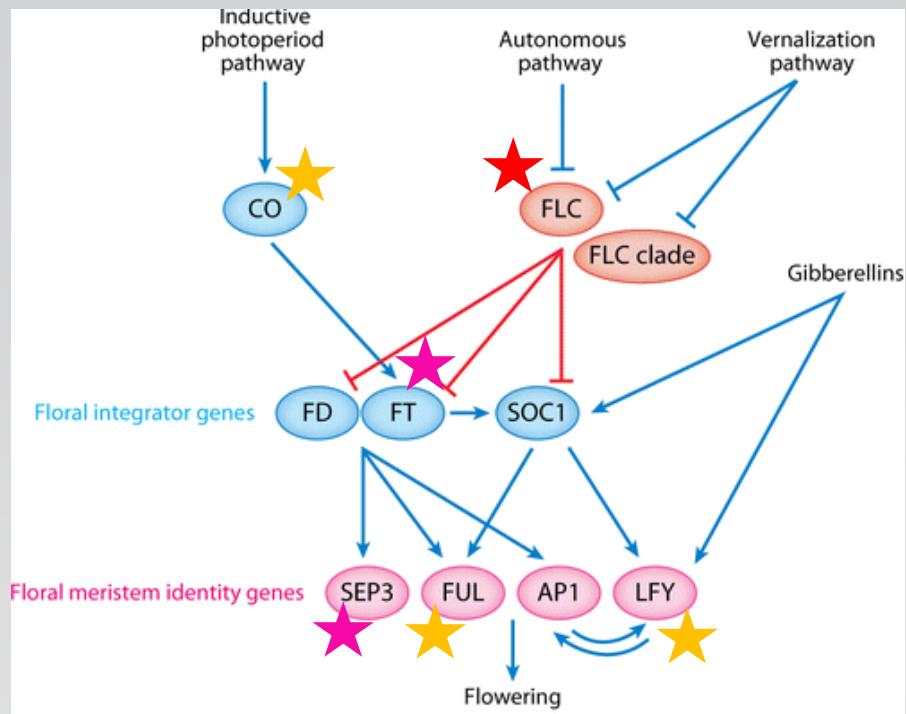
**Environment independent**



**G x E**



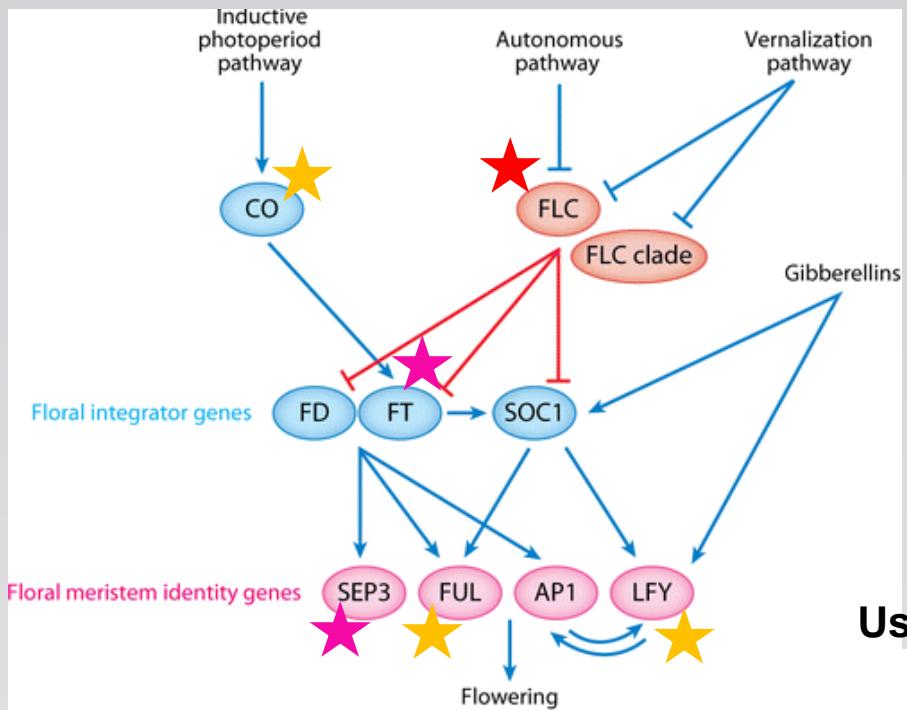
# Beyond GWAS: Understanding the function of pathways



## Significance in GWAS at 10°C

- |             |             |
|-------------|-------------|
| Yellow star | $< 10^{-3}$ |
| Pink star   | $< 10^{-4}$ |
| Red star    | $< 10^{-5}$ |

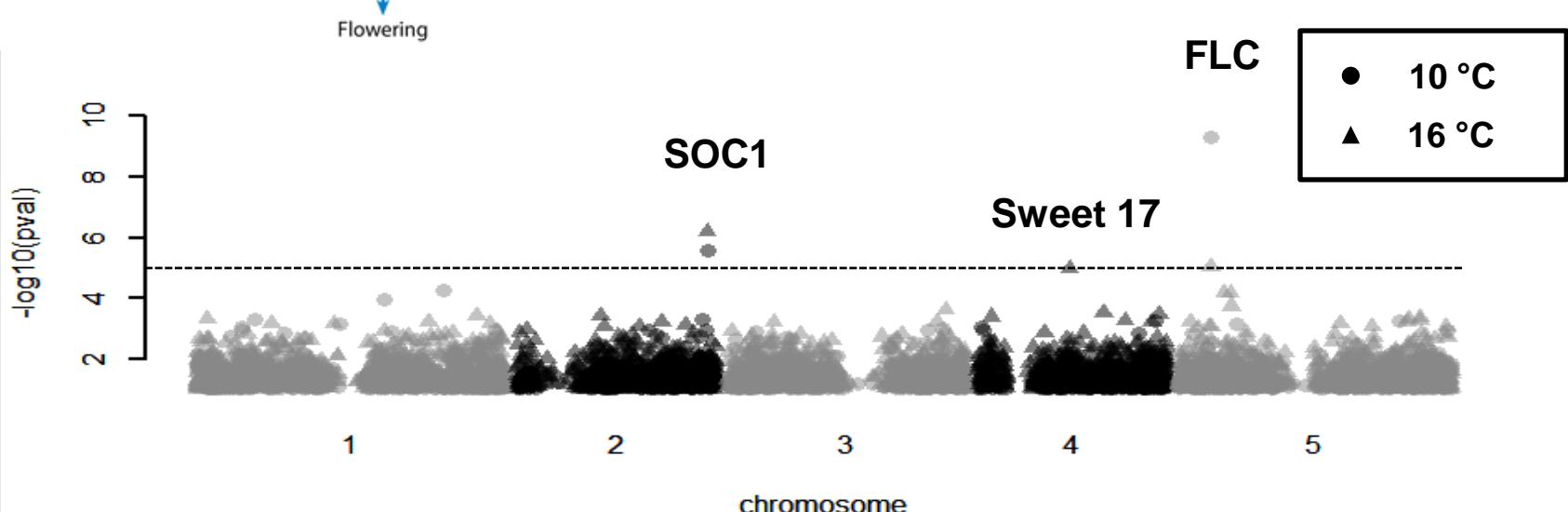
# Beyond GWAS: Understanding the function of pathways



## Significance in GWAS at 10°C

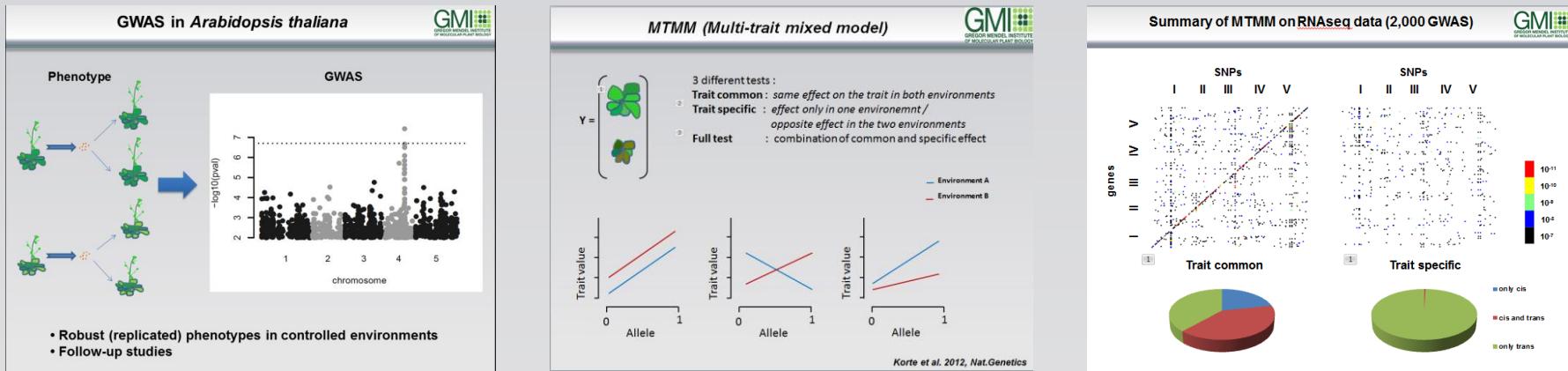
- Yellow star:  $< 10^{-3}$
- Pink star:  $< 10^{-4}$
- Red star:  $< 10^{-5}$

## Using RNA expression to explain the phenotype

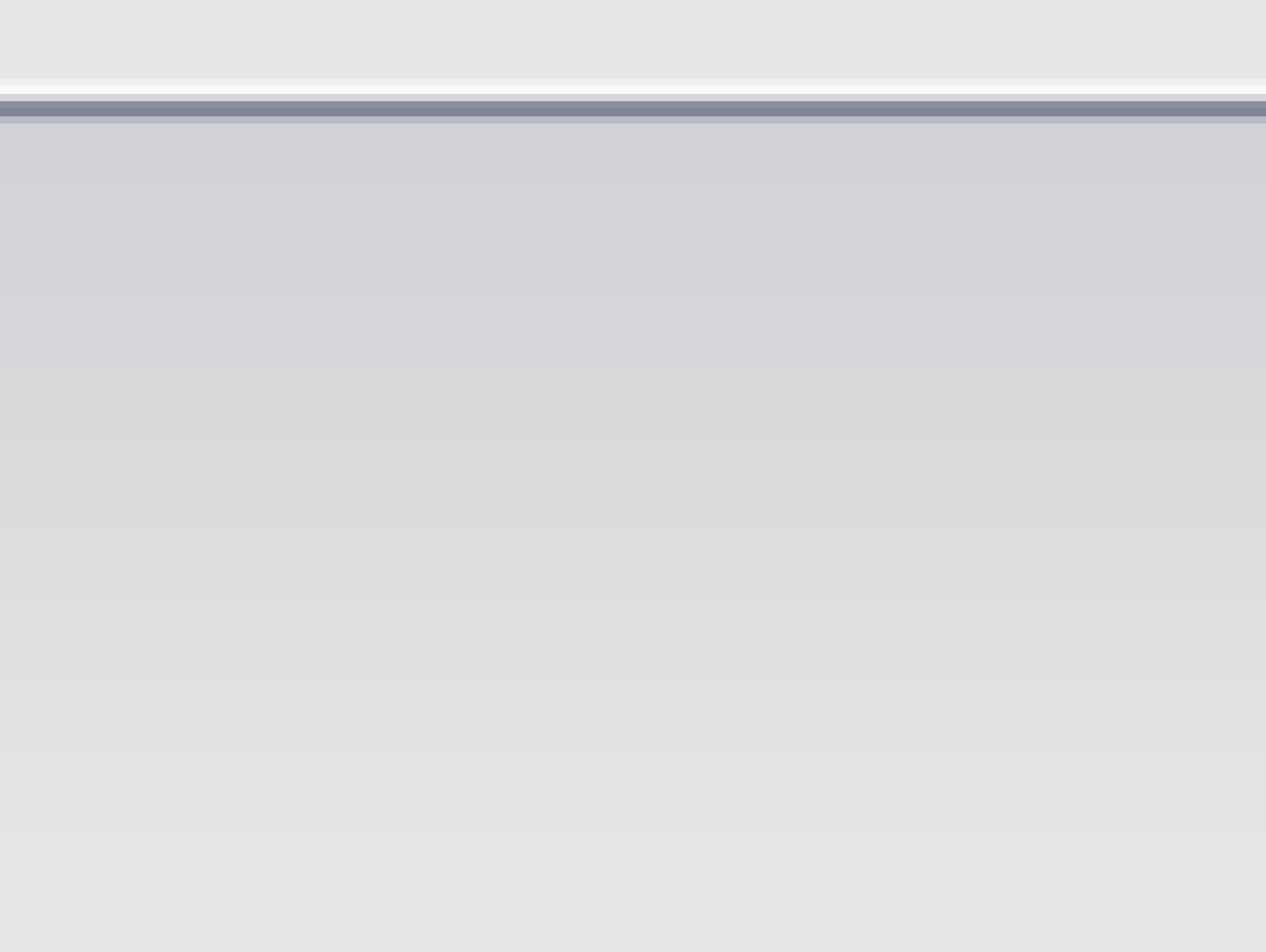


# Summary

## Understanding genotype – phenotype relationships



Integration of different data layers into joint models



# GWAS with webtools

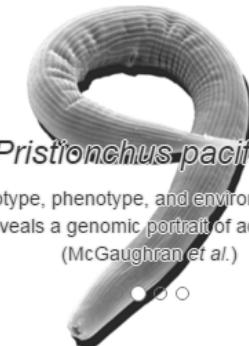


Center for Computational  
and Theoretical Biology

The browser or screen resolution is too low. The web-application is designed for a screen resolution larger than 1600x1000. Lower screen resolutions may cause display errors. X

We released easyGWAS version 2.5! An overview of all new features and visualisations can be found here ➔

### Publicly Available Studies or Data



*Pristionchus pacificus*  
Linking genotype, phenotype, and environment in an island nematode reveals a genomic portrait of adaptive divergence  
(McGaughran et al.)

### easyGWAS Statistics

👤 Registered Users	452
🌐 Public Projects	3
💻 Public GWAS	109
🐾 Public Species	3
📊 Public Datasets	5
🐦 Public Phenotypes	116
🔬 Performed GWAS	4834

## [Video](#)

gwas.gmi.oeaw.ac.at/#!home

Home Phenotypes Germplasm Genotype Log In

### Welcome to GWA-Portal

Resource for phenotypes and GWAS studies

Users can interactively browse and view public phenotypes.  
Logged-in users can create studies, upload phenotypes, run GWAS analysis using different methods on different genotype datasets and share the data with other users.

[Take a tour](#)



#### Phenotypes

Display studies, phenotypes and GWAS results.

[Browse](#)



#### New GWAS analysis

Logged-in users can create and run GWAS studies.

[Create](#)



#### GWAPP

Run GWAS on the fly (250k SNP datasets).

[Open](#)

#### Recent News

 **Filter for Meta-analysis added**  
by Umit Seren on Wed Oct 16 16:32:20 GMT+200 2013  
A new feature was added that allows users to filter the [Top-results](#) (Meta-analysis) by different filters. Aside from filtering by method (AMM, KW, etc) or Genotype (250k) it is also possible to filter by the recently introduced [candidate gene lists](#). This filter will only show those top SNPs that are in the vicinity (20kb up/down) of the genes in the candidate gene list.

 **Candidate gene lists added**  
by Umit Seren on Thu Oct 03 12:07:53 GMT+200 2013  
In the phenotype section users can create and share candidate gene lists. In the future it should be possible to run candidate enrichment analysis and filter the meta-analysis results by the candidate gene list

 **Beta version launched**  
by Umit Seren on Tue Jun 25 13:34:04 GMT+200 2013  
The beta version of GWA-Portal was launched. Basic functionality is in place

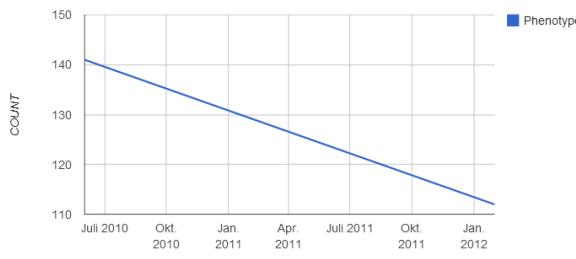
#### Quick Stats

 0	Users
 3	Studies
 253	Phenotypes
 572	Analysis
 2	Ontologies
 3	Publications

#### Public Phenotypes

Studies Phenotypes Analyses

Weeks Months Years



Month	Count
Jul. 2010	140
Okt. 2010	135
Jan. 2011	130
Apr. 2011	125
Juli 2011	120
Okt. 2011	115
Jan. 2012	110

@ 2012 GMI  transPLANT   

**Username:**  
**Passwd:**

[bag@workshop.org](mailto:bag@workshop.org)  
gwas

DATA on the server :

**X\_2029\_1.rda**  
**Annotation\_1001genomes.rda**  
**MTMM\_SAMPLE\_DATA.Rdata**  
**tair9.Rdata**

Scripts for MLMM are on the Gmi github site:

**<https://github.com/Gregor-Mendel-Institute>**

GEMMA is installed on the server : [akorte@gdcsrv1 ~]\$ cd /usr/local/gemma/

## More ideas

**Transform the phenotype (sqrt, log, boxcox) and rerun the analysis  
-> (compare results !)**

**Analyze subsets , remove outliers, permute the phenotype**

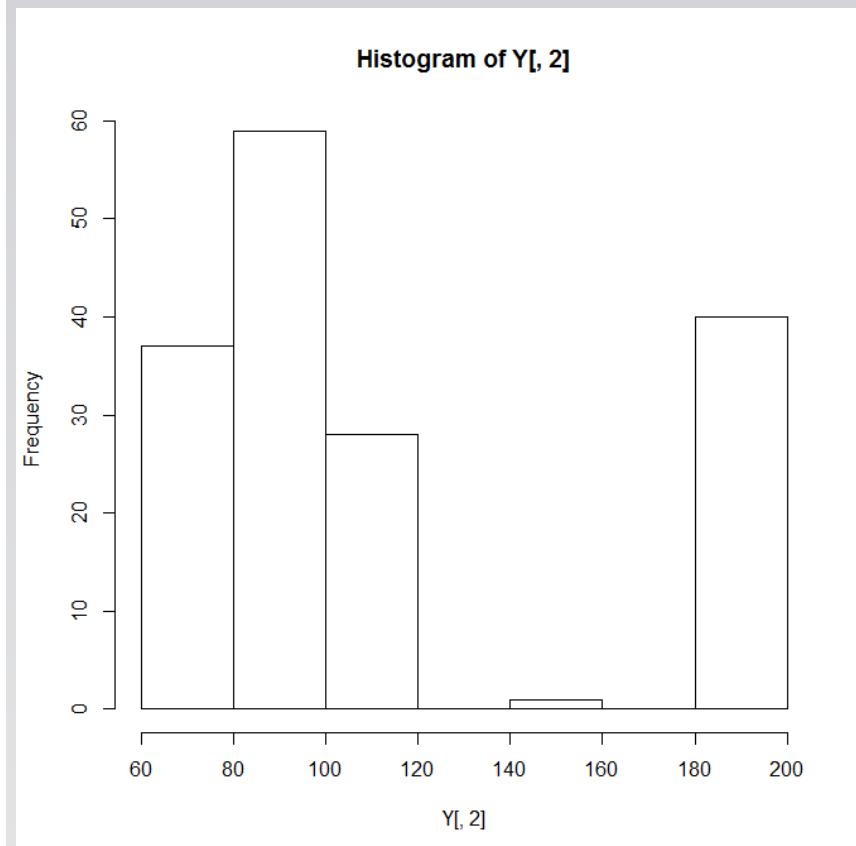
**Merge SNPs with genes; compare genic / intergenic SNPs ..**

R is a collaborative project with many contributors.  
 Type 'contributors()' for more information and  
 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
 'help.start()' for an HTML browser interface to help.

Type 'q()' to quit R.

```
> setwd('Z:/arthur/work/R')
> source('emma.r')
>
> source('gwas.r')
> source('plots_gwas.r')
>
> load('MTMM_SAMPLE_DATA.Rdata')
> head(Y)
  ecotype_id      SD      SDV
298      5837  75.6667     NA
308      6008  69.6667 58.0000
314      6016 200.0000 55.0000
321      6024 200.0000 72.3333
328      6039  89.8083 70.5000
329      6040 200.0000     NA
> dim(Y)
[1] 165   3
> dim(X)
[1] 1307 214051
> X[1:5,1:5]
  1- 657 1- 3102 1- 4648 1- 4880 1- 5975
1    0      1      0      1      0
2    0      0      0      0      0
4    1      1      1      0      0
5    0      0      0      0      0
6    1      1      1      0      0
> dim(K)
[1] 1307 1307
> K[1:5,1:5]
      1       2       4       5       6
1 1.0000000 0.7521432 0.7535681 0.7707930 0.7552406
2 0.7521432 1.0000000 0.7709938 0.7489477 0.7394032
4 0.7535681 0.7709938 1.0000000 0.7701436 0.7377167
5 0.7707930 0.7489477 0.7701436 1.0000000 0.7397116
6 0.7552406 0.7394032 0.7377167 0.7397116 1.0000000
> |
```



```
> hist(Y[,2])
> a<-Sys.time()
> amm_gwas(Y,X,K,include.lm=T,calculate.effect.size=FALSE)
GWAS performed on 165 ecotypes, 0 values excluded
SNP_INFO file created
pseudo-heritability estimate is 0.9999546
> Sys.time()-a
Time difference of 1.359486 mins
> head(output)
      SNP Chr      Pos AC_1 AC_0   MAC       MAF      Pval    Pval_lm
1 1- 10000173    1 10000173    35   130    35 0.21212121 0.05266552 0.004369247
2 1- 10000715    1 10000715    18   147    18 0.10909091 0.338888896 0.389214696
3 1- 10001699    1 10001699   108    57    57 0.34545455 0.78554188 0.862913234
4 1- 10002088    1 10002088   102    63    63 0.38181818 0.57424989 0.222504617
5 1- 10002264    1 10002264     6   159     6 0.03636364 0.65037890 0.616469715
6 1- 10002246    1 10002246    82    79    79 0.44242424 0.05552055 0.018774129
```

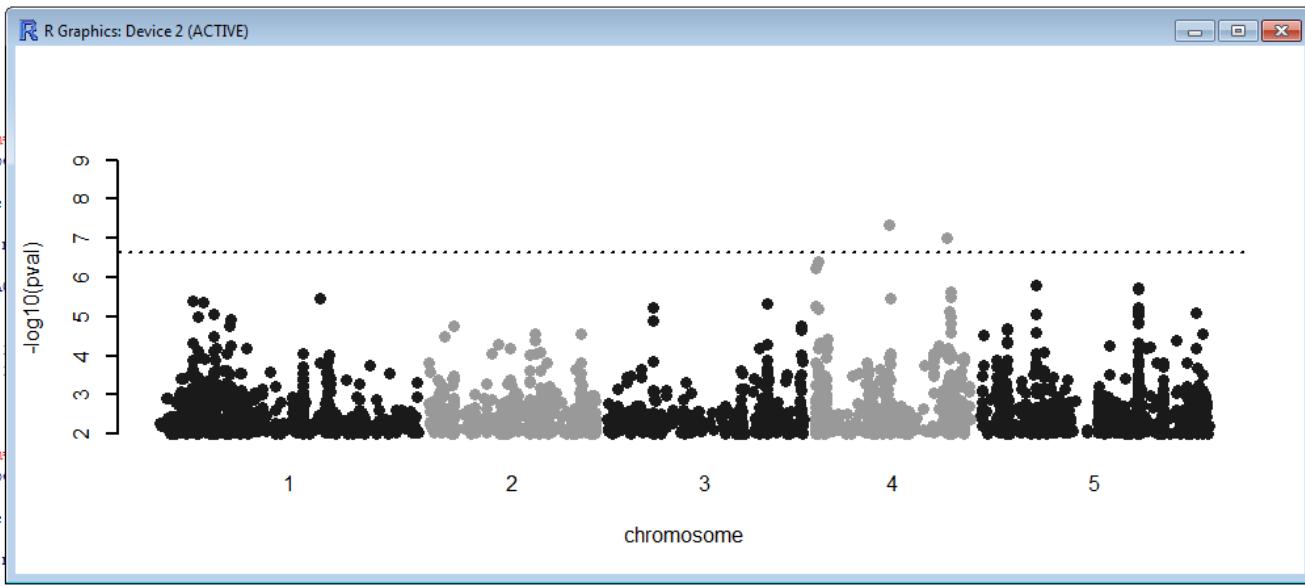
```
> ######, gwas -- ,  
> amm_gwas(Y,X,K,n=3)  
GWAS performed on 150 ecotypes, 0 values excluded  
SNP_INFO file created  
pseudo-heritability estimate is 0.9560748  
> head(output)  
          SNP Chr      Pos AC_1 AC_0 MAC        MAF      Pval  
1 1- 10000173    1 10000173    33   117   33 0.22000000 0.6862535  
2 1- 10000715    1 10000715    17   133   17 0.11333333 0.6342901  
3 1- 10001699    1 10001699    99    51   51 0.34000000 0.5829813  
4 1- 10002088    1 10002088    94    56   56 0.37333333 0.5040329  
5 1- 10002264    1 10002264     5   145    5 0.03333333 0.3683449  
6 1- 10003346    1 10003346    85    65   65 0.43333333 0.7158708  
> |
```

# Estimating the effect size

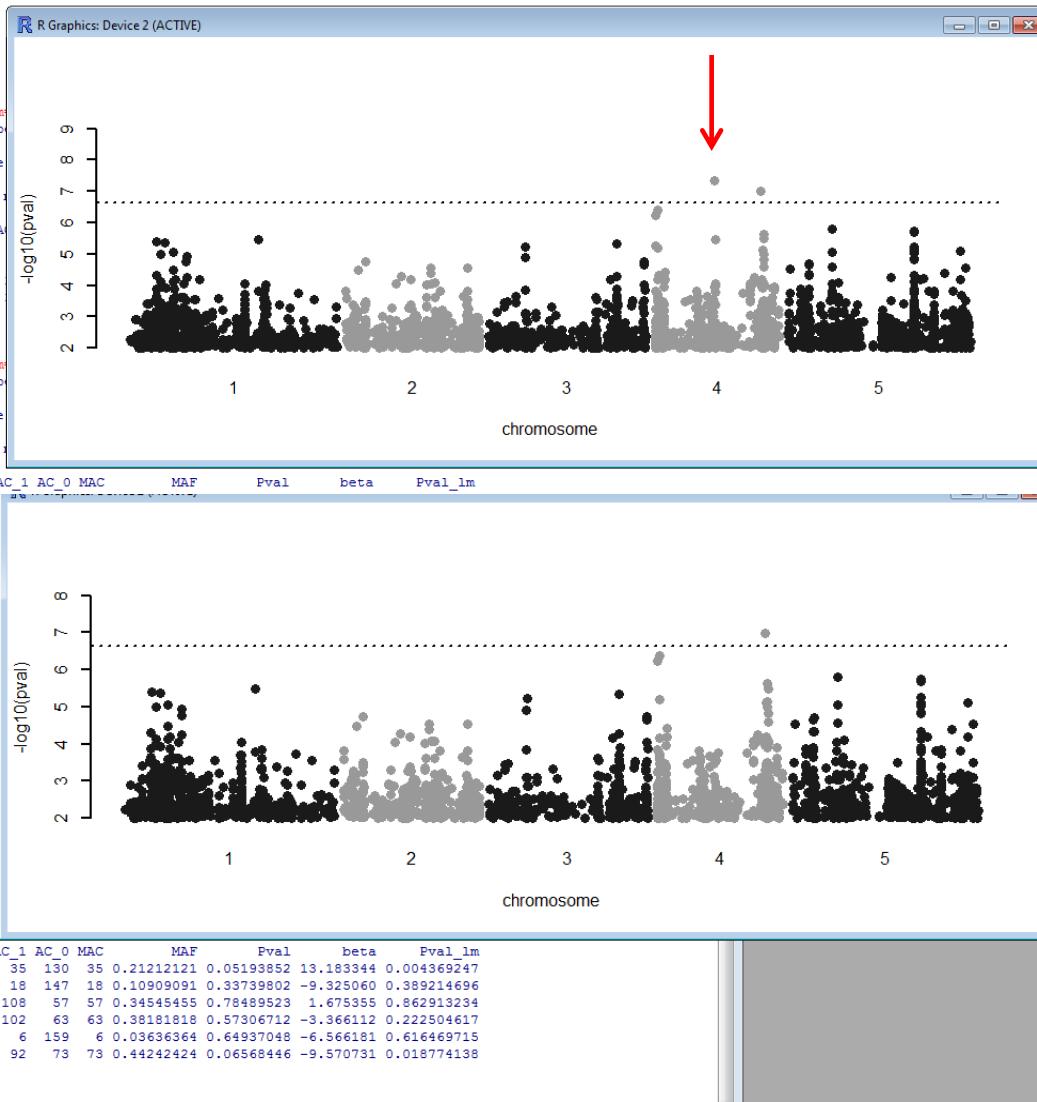


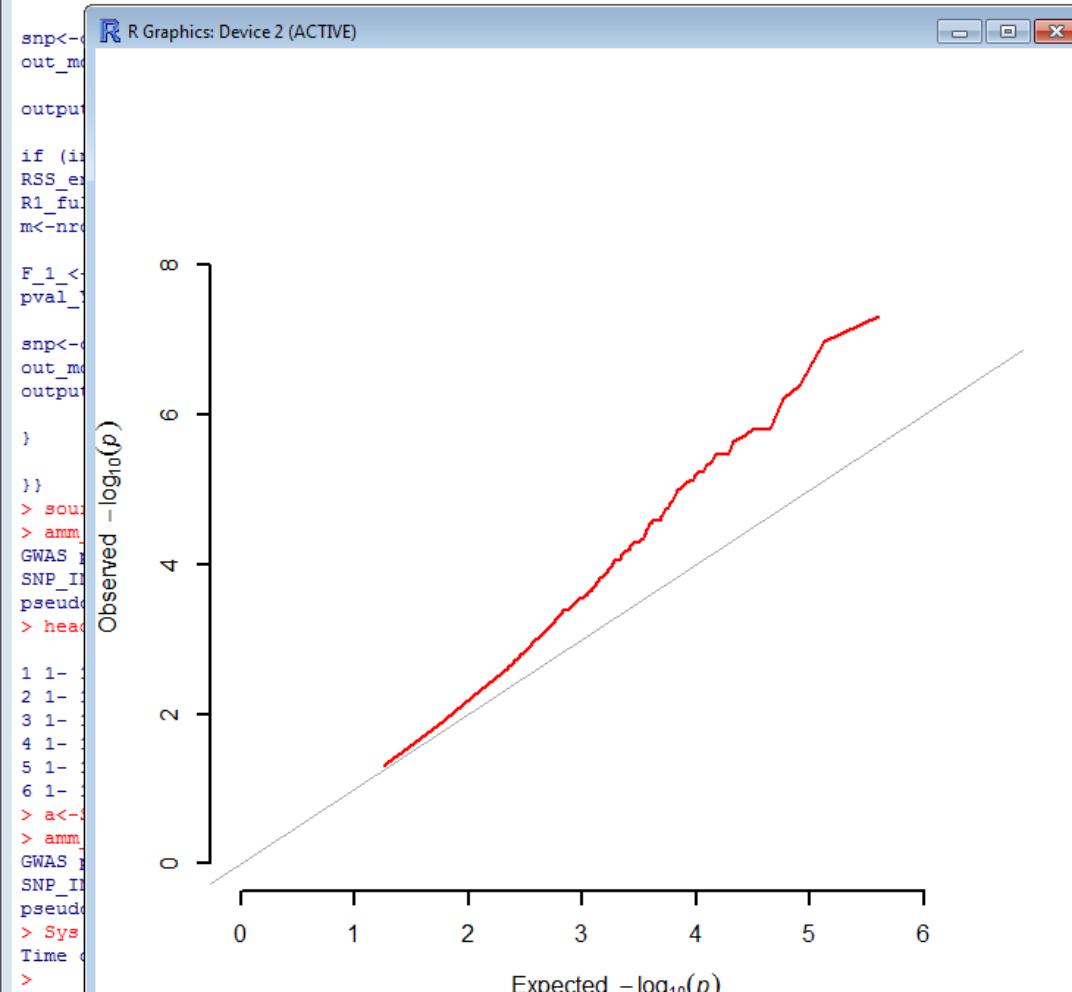
```
> a<-Sys.time()
> amm_gwas(Y_,X,K,include.lm=T)
GWAS performed on 165 ecotypes, 0 values excluded
SNP_INFO file created
pseudo-heritability estimate is 0.9999546
> Sys.time()-a
Time difference of 1.376254 mins
> head(output)
      SNP Chr      Pos AC_1 AC_0 MAC        MAF       Pval     Pval_lm
1 1- 10000173    1 10000173   35   35 0.21212121 0.05266552 0.004369247
2 1- 10000715    1 10000715   18   18 0.10909091 0.33888896 0.389214696
3 1- 10001699    1 10001699   108   57 0.34545455 0.78554188 0.862913234
4 1- 10002088    1 10002088   102   63 0.38181818 0.57424989 0.222504617
5 1- 10002264    1 10002264    6   159   6 0.03636364 0.65037890 0.616469715
6 1- 10003346    1 10003346   92   73 0.44242424 0.06652066 0.018774138
> a<-Sys.time()
> amm_gwas(Y_,X,K,include.lm=T,calculate.effect.size=T)
GWAS performed on 165 ecotypes, 0 values excluded
SNP_INFO file created
pseudo-heritability estimate is 0.9999546
> Sys.time()-a
Time difference of 7.383855 mins
> head(output)
      SNP Chr      Pos AC_1 AC_0 MAC        MAF       Pval      beta     Pval_lm
1 1- 10000173    1 10000173   35   35 0.21212121 0.05193852 13.183344 0.004369247
2 1- 10000715    1 10000715   18   18 0.10909091 0.33739802 -9.325060 0.389214696
3 1- 10001699    1 10001699   108   57 0.34545455 0.78489523  1.675355 0.862913234
4 1- 10002088    1 10002088   102   63 0.38181818 0.57306712 -3.366112 0.222504617
5 1- 10002264    1 10002264    6   159   6 0.03636364 0.64937048 -6.566181 0.616469715
6 1- 10003346    1 10003346   92   73 0.44242424 0.06568446 -9.570731 0.018774138
> |
```

```
6008 69.6667  
6016 200.0000  
6024 200.0000  
6039 89.8083  
6040 200.0000  
> hist(Y_)  
>  
> a<-Sys.time()  
> amm_gwas(Y_,X,K,include.lm=TRUE)  
GWAS performed on 165 ecotypes  
SNP_INFO file created  
pseudo-heritability estimate  
> Sys.time()-a  
Time difference of 1.376254  
> head(output)  
  SNP Chr    Pos A  
1 1- 10000173 1 10000173  
2 1- 10000715 1 10000715  
3 1- 10001699 1 10001699  
4 1- 10002088 1 10002088  
5 1- 10002264 1 10002264  
6 1- 10003346 1 10003346  
> a<-Sys.time()  
> amm_gwas(Y_,X,K,include.lm=TRUE)  
GWAS performed on 165 ecotypes  
SNP_INFO file created  
pseudo-heritability estimate  
> Sys.time()-a  
Time difference of 7.383855  
> head(output)  
  SNP Chr    Pos AC_1 AC_0 MAF      Pval     beta   Pval_lm  
1 1- 10000173 1 10000173 35 130 35 0.21212121 0.05193852 13.183344 0.004369247  
2 1- 10000715 1 10000715 18 147 18 0.10909091 0.33739802 -9.325060 0.389214696  
3 1- 10001699 1 10001699 108 57 57 0.34545455 0.78489523 1.675355 0.862913234  
4 1- 10002088 1 10002088 102 63 63 0.38181818 0.57306712 -3.366112 0.222504617  
5 1- 10002264 1 10002264 6 159 6 0.03636364 0.64937048 -6.566181 0.616469715  
6 1- 10003346 1 10003346 92 73 73 0.44242424 0.06568446 -9.570731 0.018774138  
> plot_gwas(output)
```



```
6008 69.666/  
6016 200.0000  
6024 200.0000  
6039 89.8083  
6040 200.0000  
> hist(Y_)  
>  
> a<-system.time()  
> amm_gwas(Y_,X,K,include.lm  
GWAS performed on 165 ecotyp  
SNP_INFO file created  
pseudo-heritability estimate  
> system.time()-a  
Time difference of 1.376254  
> head(output)  
          SNP Chr   Pos A  
1 1- 10000173 1 10000173  
2 1- 10000715 1 10000715  
3 1- 10001699 1 10001699  
4 1- 10002088 1 10002088  
5 1- 10002264 1 10002264  
6 1- 10003346 1 10003346  
> a<-system.time()  
> amm_gwas(Y_,X,K,include.lm  
GWAS performed on 165 ecotyp  
SNP_INFO file created  
pseudo-heritability estimate  
> system.time()-a  
Time difference of 7.383855  
> head(output)  
          SNP Chr   Pos A  
1 1- 10000173 1 10000173  
2 1- 10000715 1 10000715  
3 1- 10001699 1 10001699  
4 1- 10002088 1 10002088  
5 1- 10002264 1 10002264  
6 1- 10003346 1 10003346  
> plot_gwas(output)  
> plot_gwas(output,maf=0.01)  
> plot_gwas(output,maf=0.001  
> plot_gwas(output,maf=0.1)  
> head(output)  
          SNP Chr   Pos A  
1 1- 10000173 1 10000173 35 130 35 0.21212121 0.05193852 13.183344 0.004369247  
2 1- 10000715 1 10000715 18 147 18 0.10909091 0.33739802 -9.325060 0.389214696  
3 1- 10001699 1 10001699 108 57 57 0.34545455 0.78489523 1.675355 0.862913234  
4 1- 10002088 1 10002088 102 63 63 0.38181818 0.57306712 -3.366112 0.222504617  
5 1- 10002264 1 10002264 6 159 6 0.03636364 0.64937048 -6.566181 0.616469715  
6 1- 10003346 1 10003346 92 73 73 0.44242424 0.06568446 -9.570731 0.018774138  
> plot_gwas(output)  
> plot_gwas(output,maf=0.1)
```





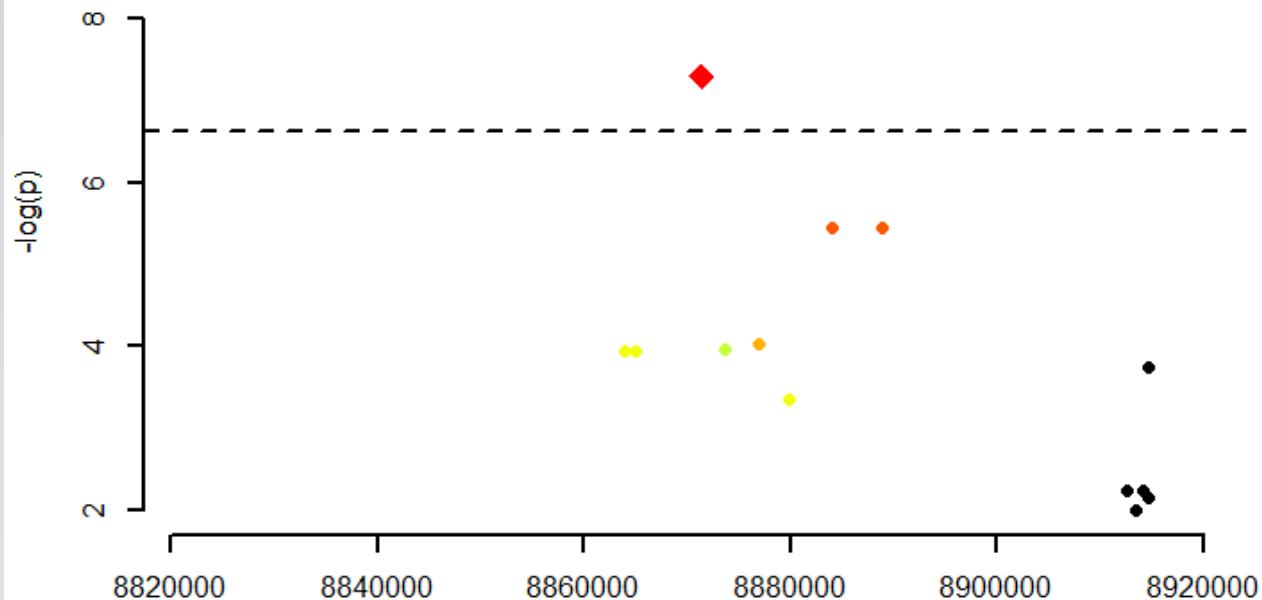
```
snp<-c()
out_m<-
output<-
if (i<
RSS_e<-
R1_ful<
m<-nr
F_1<-
pval_<-
snp<-
out_m<-
output<-
}
})
> sou<-
> amm<-
GWAS<-
SNP_ID<-
pseudo<-
> head<-
1 1-
2 1-
3 1-
4 1-
5 1-
6 1-
> a<-S<-
> amm<-
GWAS<-
SNP_ID<-
pseudo<-
> Sys<-
Time<-
> Sys<-
> head<-
1 1- 10000173 1 10000173 35 130 35 0.21212121 0.05193852 13.183344 0.004369247
2 1- 10000715 1 10000715 18 147 18 0.10909091 0.33739802 -9.325060 0.389214696
3 1- 10001699 1 10001699 108 57 57 0.34545455 0.78489523 1.675355 0.862913234
4 1- 10002088 1 10002088 102 63 63 0.38181818 0.57306712 -3.366112 0.222504617
5 1- 10002264 1 10002264 6 159 6 0.03636364 0.64937048 -6.566181 0.616469715
6 1- 10003346 1 10003346 92 73 73 0.44242424 0.06568446 -9.570731 0.018774138
> plot_gwas(output)
> qq_plot(output)
> local_plot(output)
> qq_plot(output)
> l
```

## Local plots

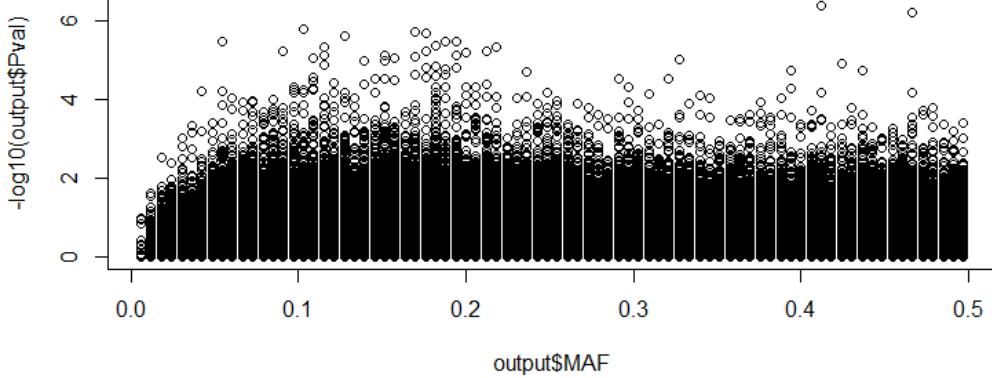


Center for Computational  
and Theoretical Biology

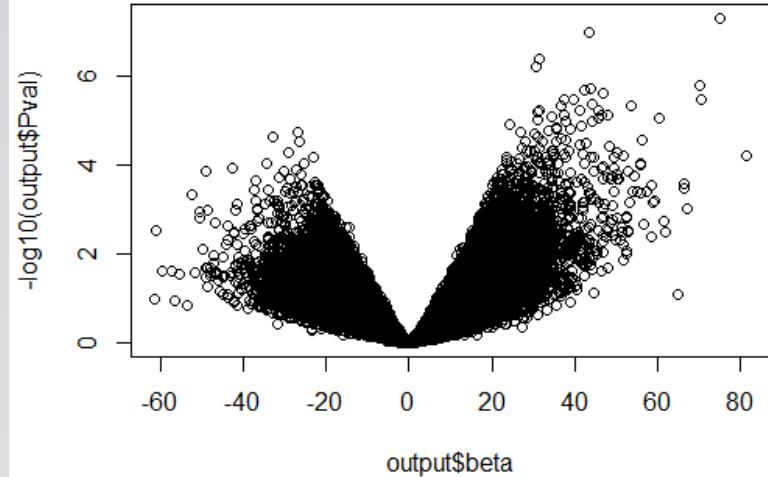
4- 8871394



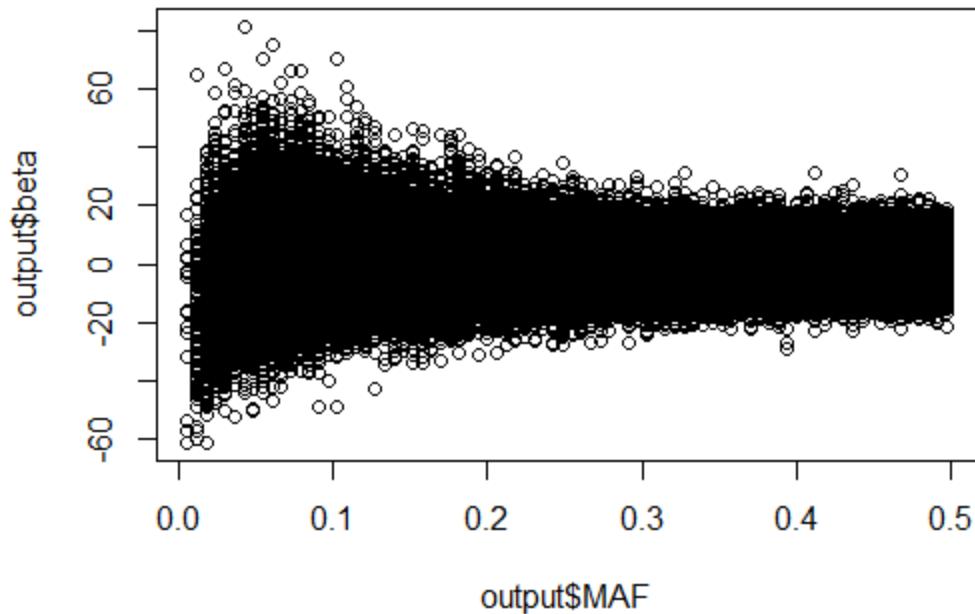
```
> local_plot(output, rank=10)
> local_plot(output, rank=1)
> color scale
```



```
> plot(output$MAF,-log10(output$Pval))  
> |
```

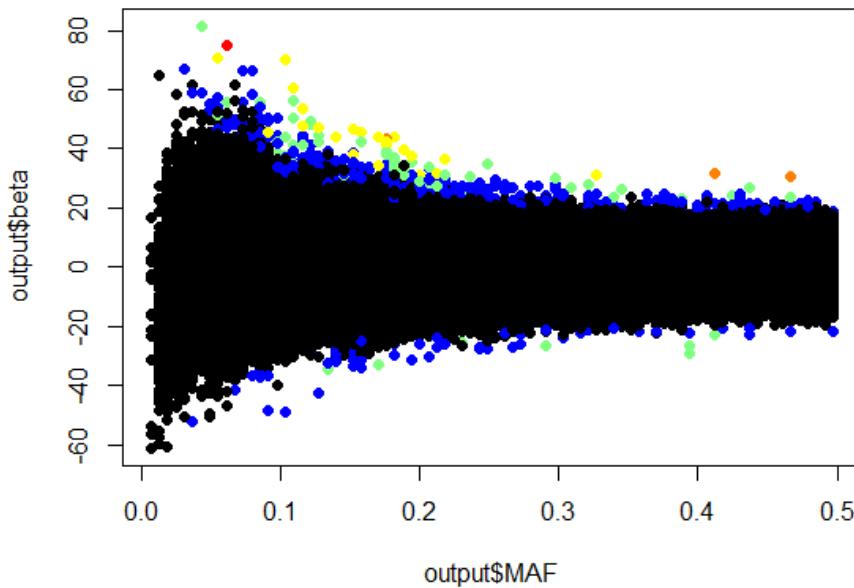


```
> plot(output$beta,-log10(output$Pval))  
> |
```

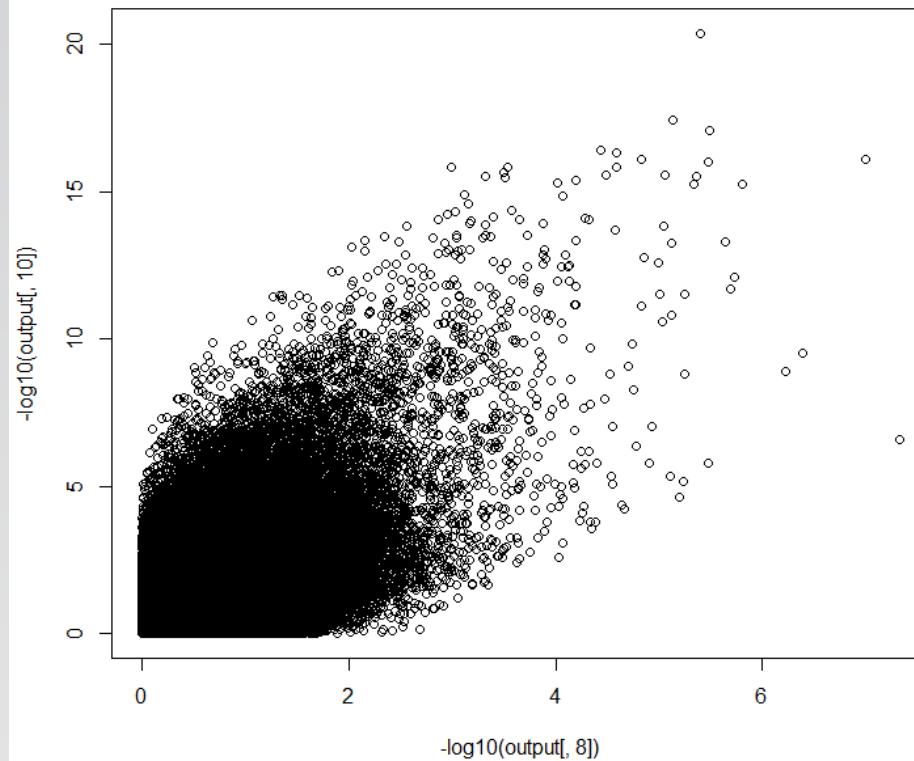


```
> plot(output$MAF, output$beta)
> |
```

# Color the SNPs accordingly to their p-value



```
> output[,11]<-ceiling(-log10(output[,8]))
> head(output)
  SNP Chr      Pos AC_1 AC_0 MAC      MAF      Pval      beta      Pval_lm V11
1 1- 10000173 1 10000173 35 130 35 0.21212121 0.05193852 13.183344 0.004369247 2
2 1- 10000715 1 10000715 18 147 18 0.10909091 0.33739802 -9.325060 0.389214696 1
3 1- 10001699 1 10001699 108 57 57 0.34545455 0.78489523 1.675355 0.862913234 1
4 1- 10002088 1 10002088 102 63 63 0.38181818 0.57306712 -3.366112 0.222504617 1
5 1- 10002264 1 10002264 6 159 6 0.03636364 0.64937048 -6.566181 0.616469715 1
6 1- 10003346 1 10003346 92 73 73 0.44242424 0.06568446 -9.570731 0.018774138 2
> jet.colors = colorRampPalette(c("black", "blue", "#7FFF7F", "yellow", "#FF7F00", "red"))
> plot(output$MAF, output$beta, col=jet.colors(8)[output[,11]], pch=16)
> |
```

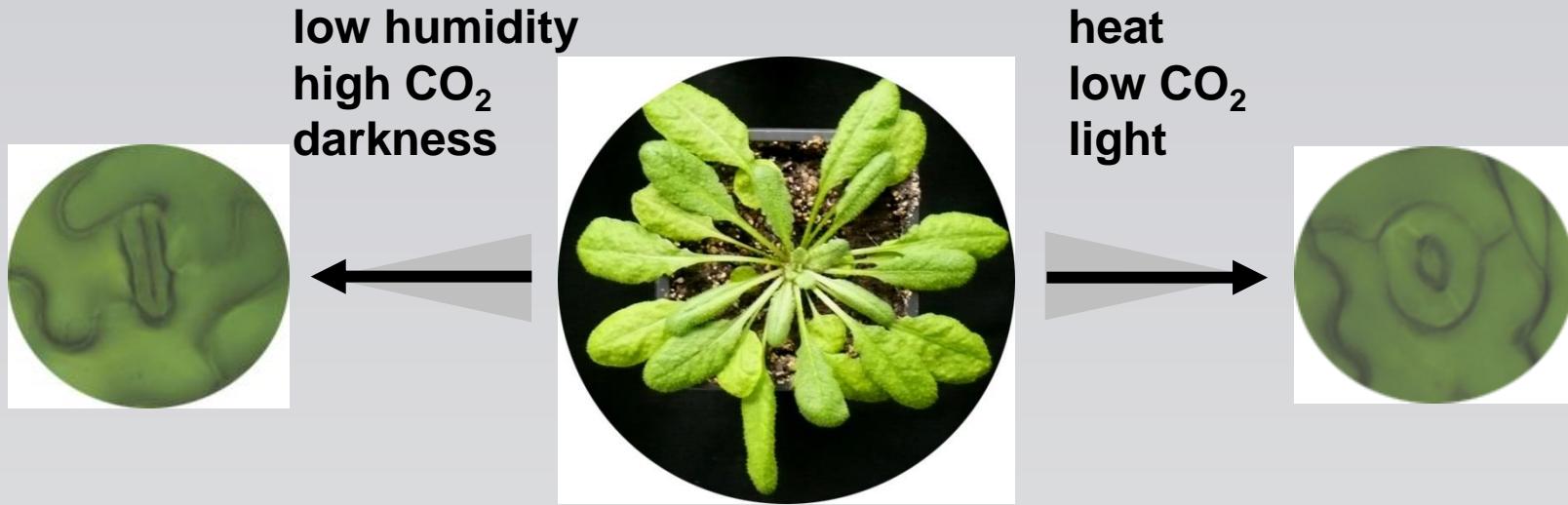


Kinship correction is not linear (unlike genomic control !)

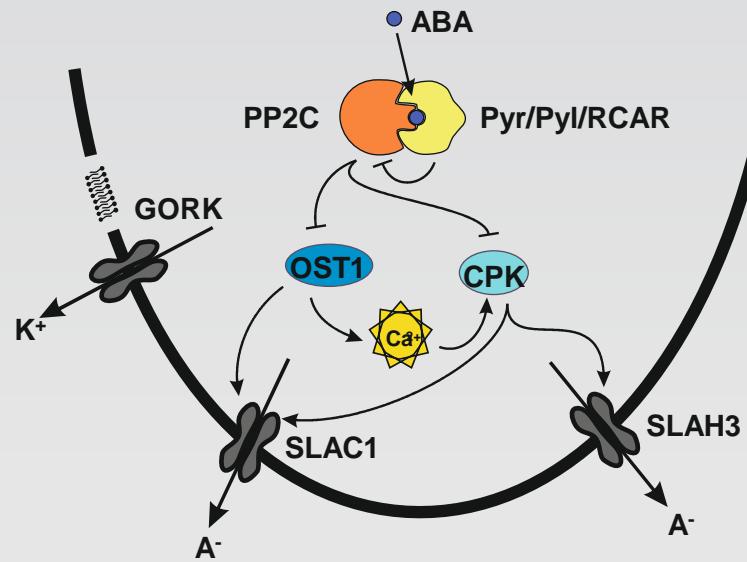
# Natural variation in guard cell signalling



Center for Computational  
and Theoretical Biology



## ABA-induced stomatal closing



## stomatal opening

