



# Evolutionary Genomics

-

## And why ML might have potential to analyze genomic data



@arthurkorte



arthurkorte

Arthur Korte

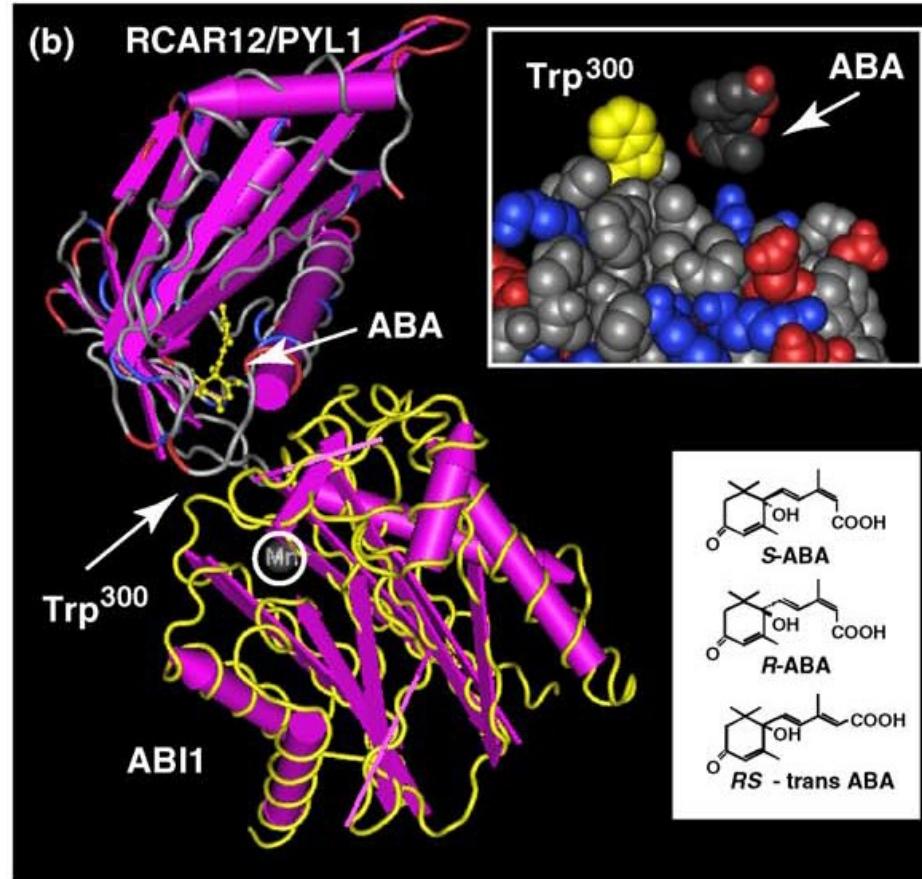
Sep 24<sup>th</sup> 2019



**“We are still in the Wild West times of Machine Learning for Genomics”**

**Andrew Kern, U Oregon**

# PhD in Molecular Biology

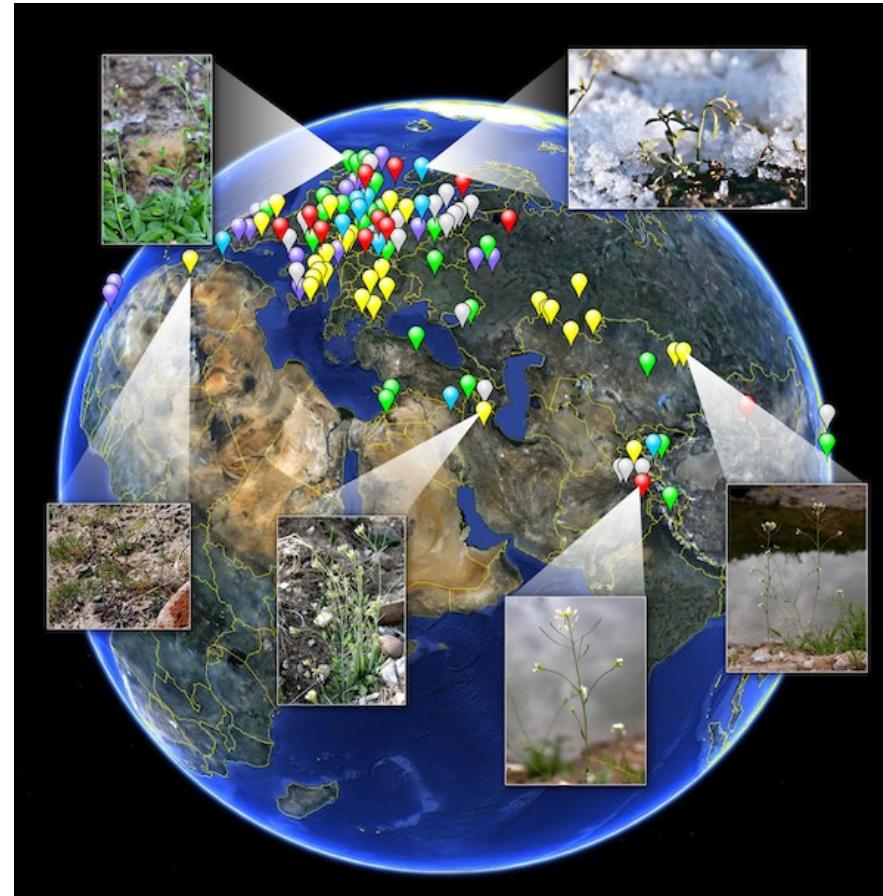


## How do plants react to stress ?

# Postdoc in Population Genetics



By Emmanuel Boutet - Own work, CC BY-SA 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=1437488>

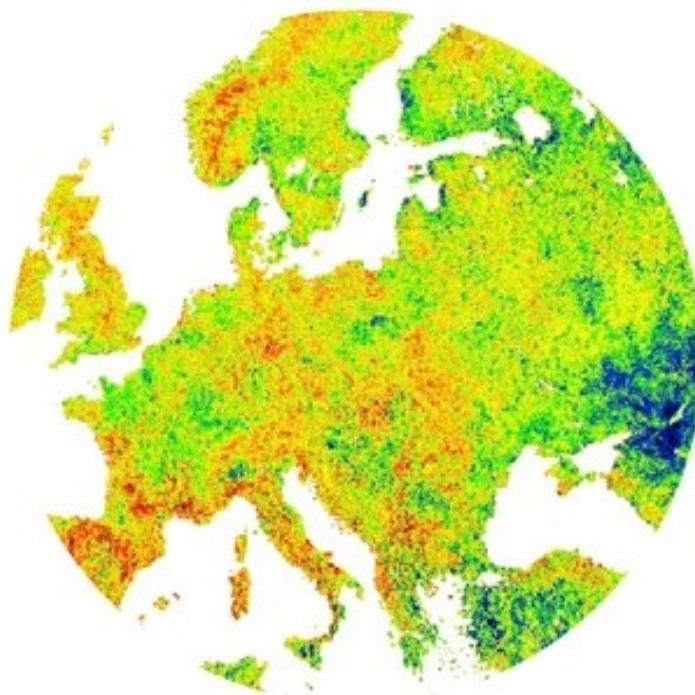


Credit:  
Photographs/images  
were contributed by  
Patrick Gooden,  
Kathleen Donohue and  
Google Earth. Graphic  
design: Jamie Simon,  
Salk Institute.

## How do populations adapt to the environment ?

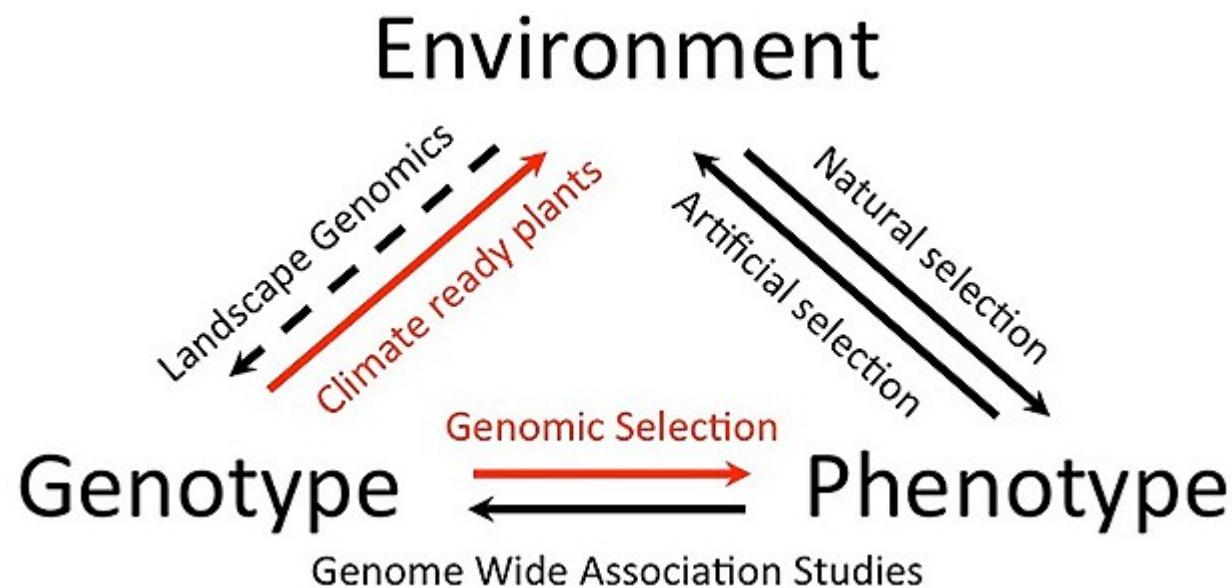
# **Quantitative Genetics & Computational Biology**

# How do plants adapt to drought

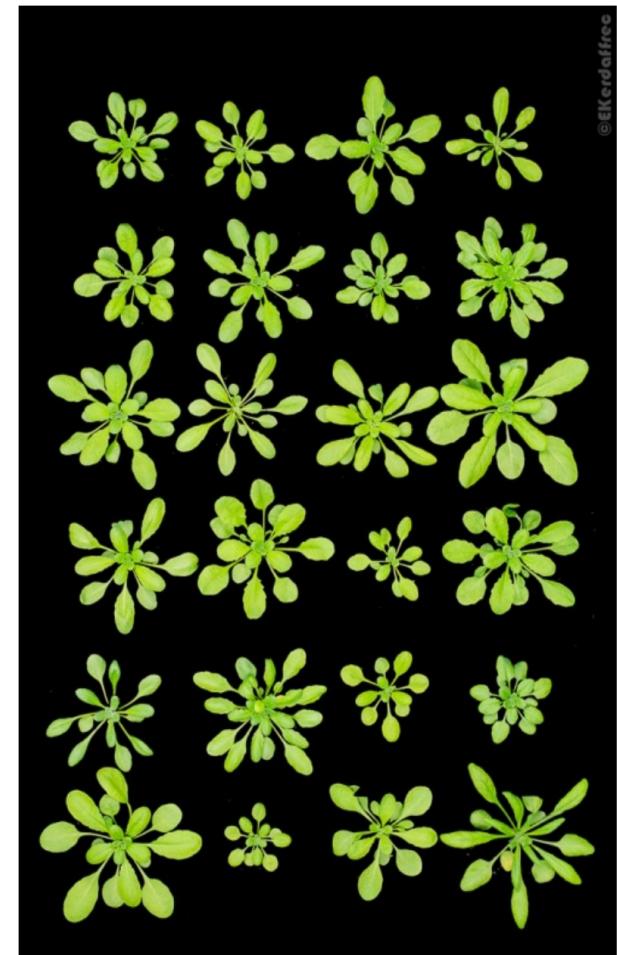
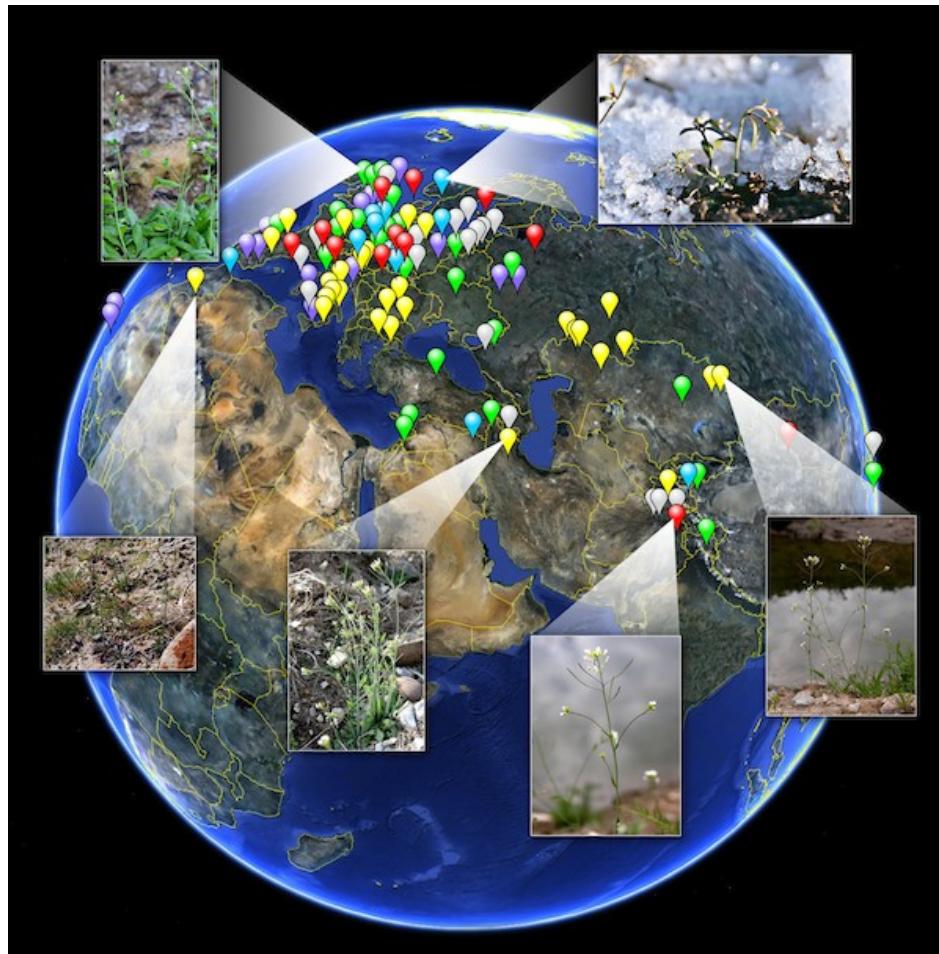


...ACATTAA  
GCGCAATTCCAC  
TCAGCATAAGGAAGA  
GTGAGGGTCACATGGC  
CTCAGTTAGCGTACATCCTA  
GTCAAAAGTGCGGGTTTCG  
TTCAGAACATCATGACCTGCAC  
GAGCTCGACATGGCAATAAC  
GAGAACATGTGTAAACGAAGC  
AAGGGGTCTGTGAATTGTTA  
GTGGGGGATAGTCACCG  
GAACAGGTTGGTTAA  
TATACTGTAGTAA  
^ATCTTTAC'

# Evolutionary Genomics



# *Arabidopsis thaliana*

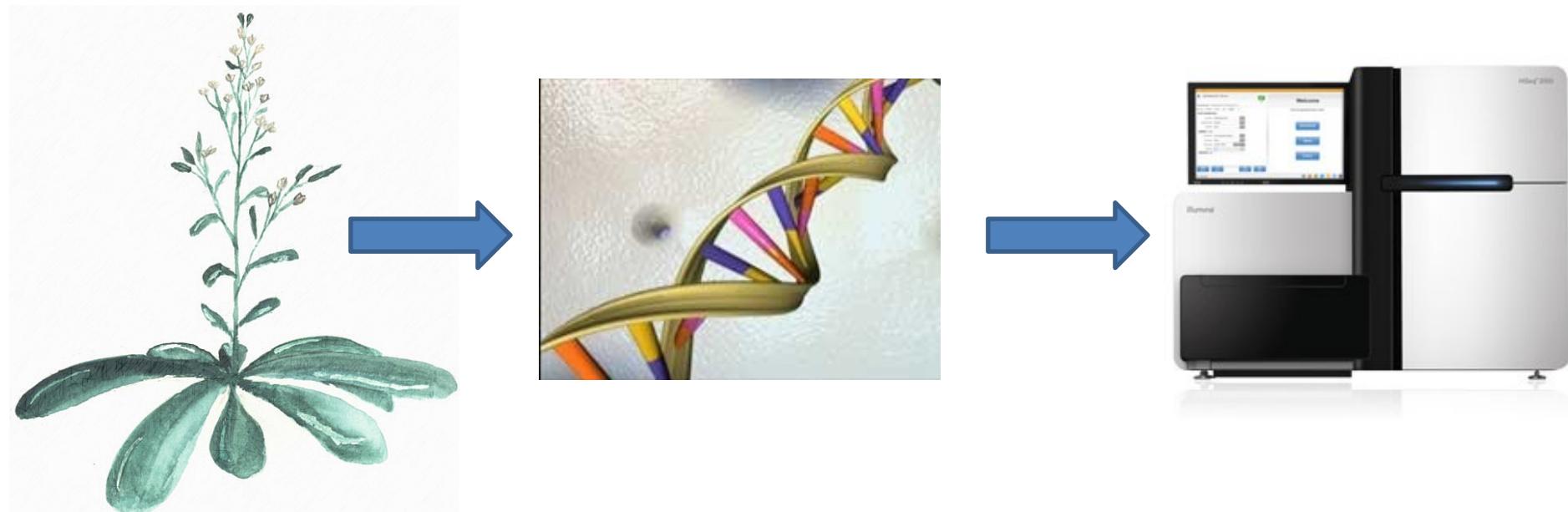


# Genome-wide association studies (GWAS)



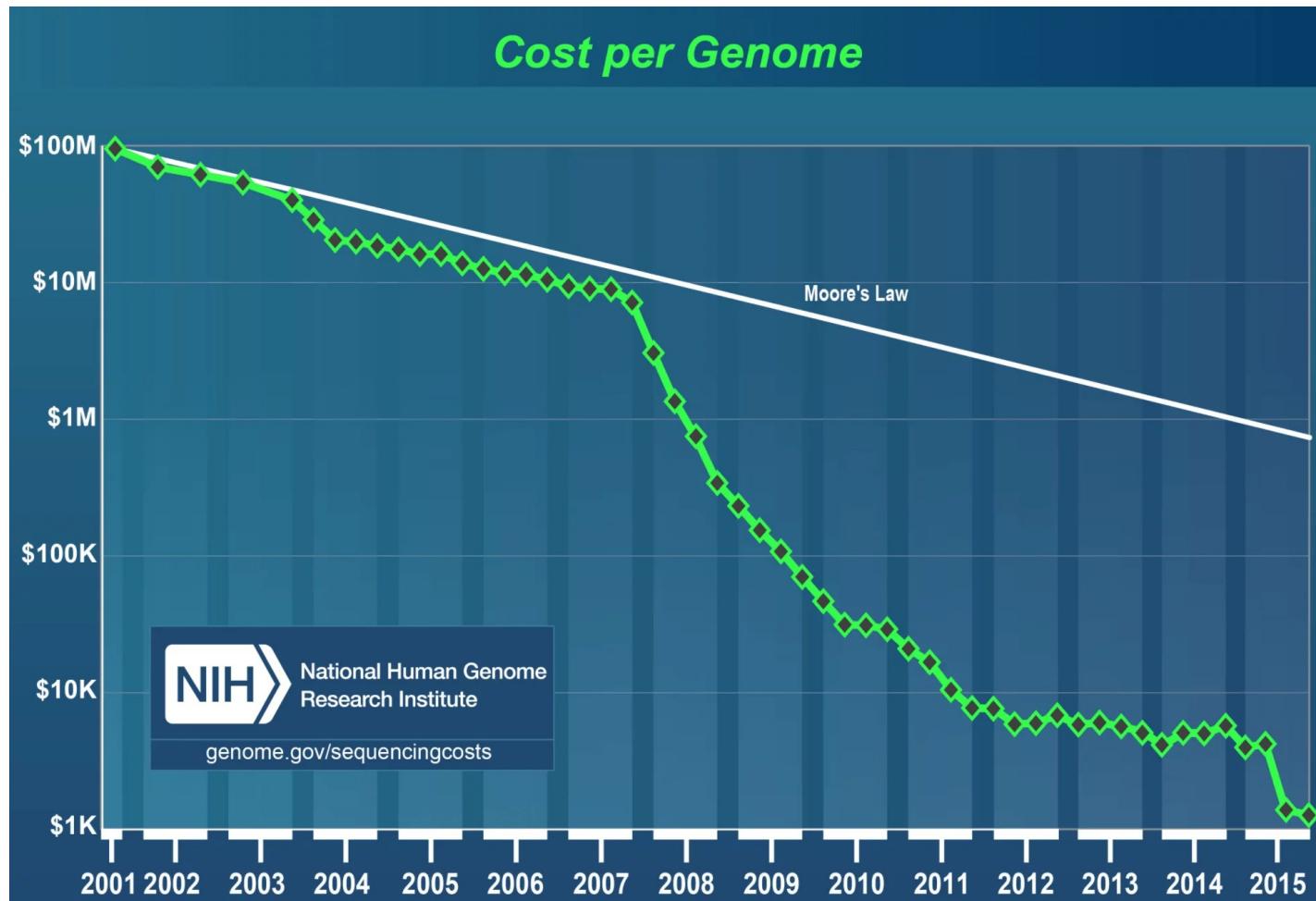
Associations between phenotype and genotype

# Genome sequencing



1. DNA isolation
2. DNA sequencing
3. Genome assembly

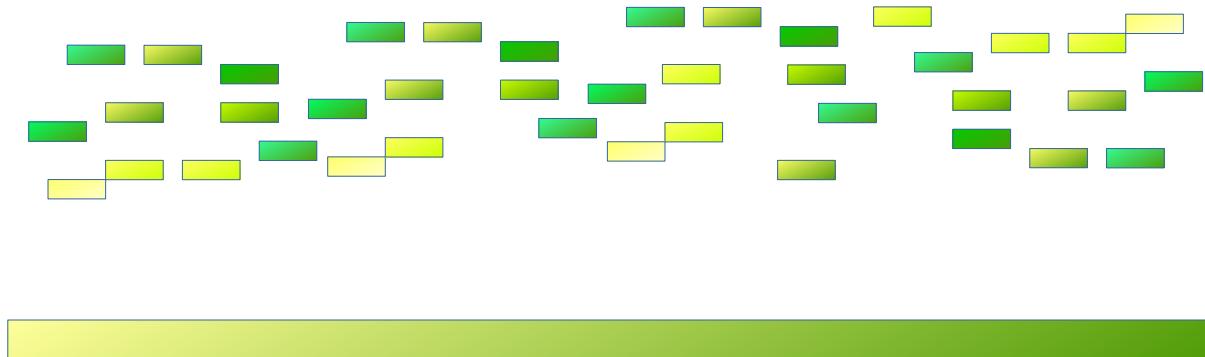
# Sequencing costs



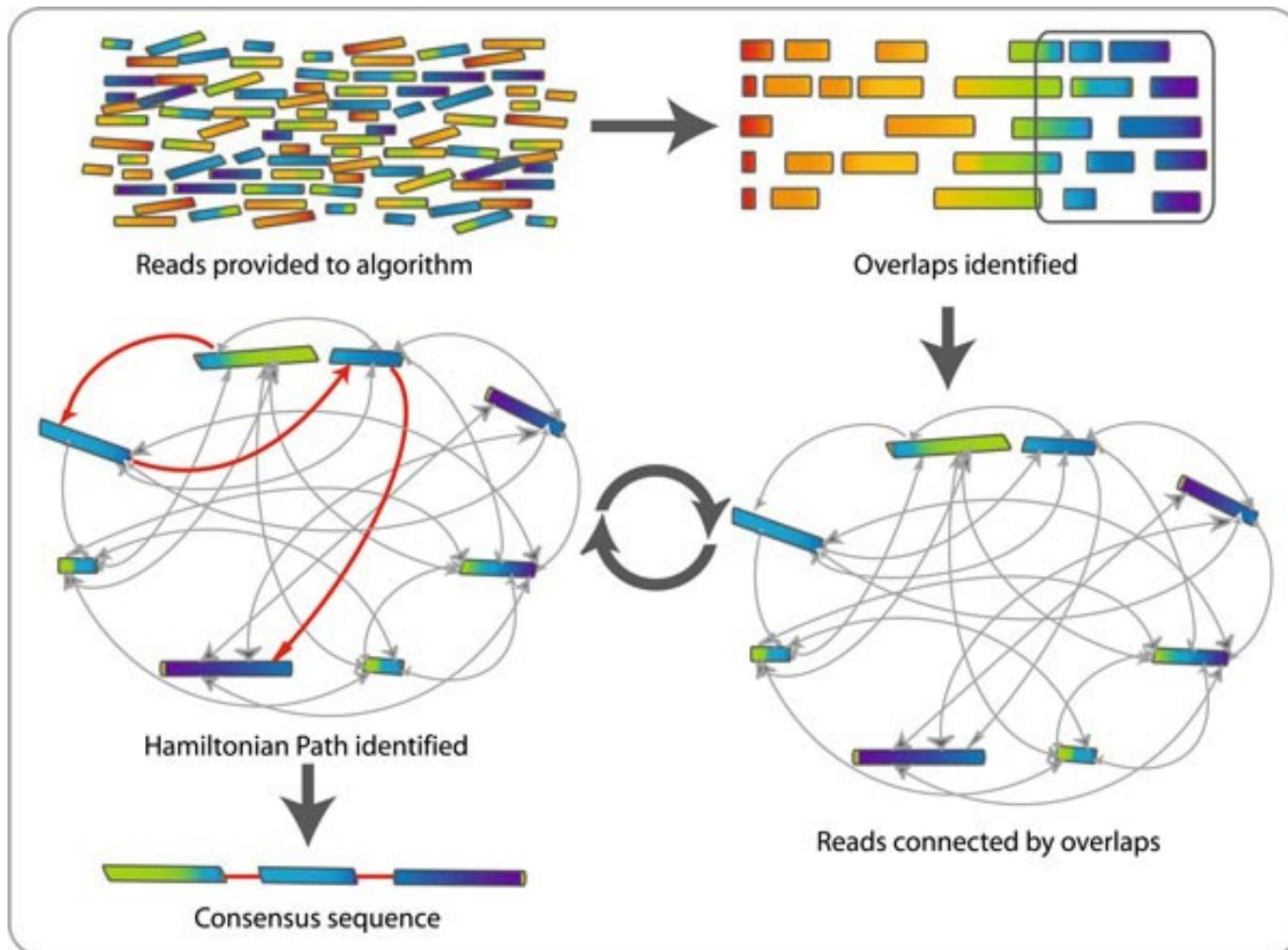
# Genome assembly : A giant jigsaw puzzle

We want a continuous strand of DNA (*aka* whole chromosomes)  
(e.g. Human 200 Mbp, Arabidopsis 20 Mbp)

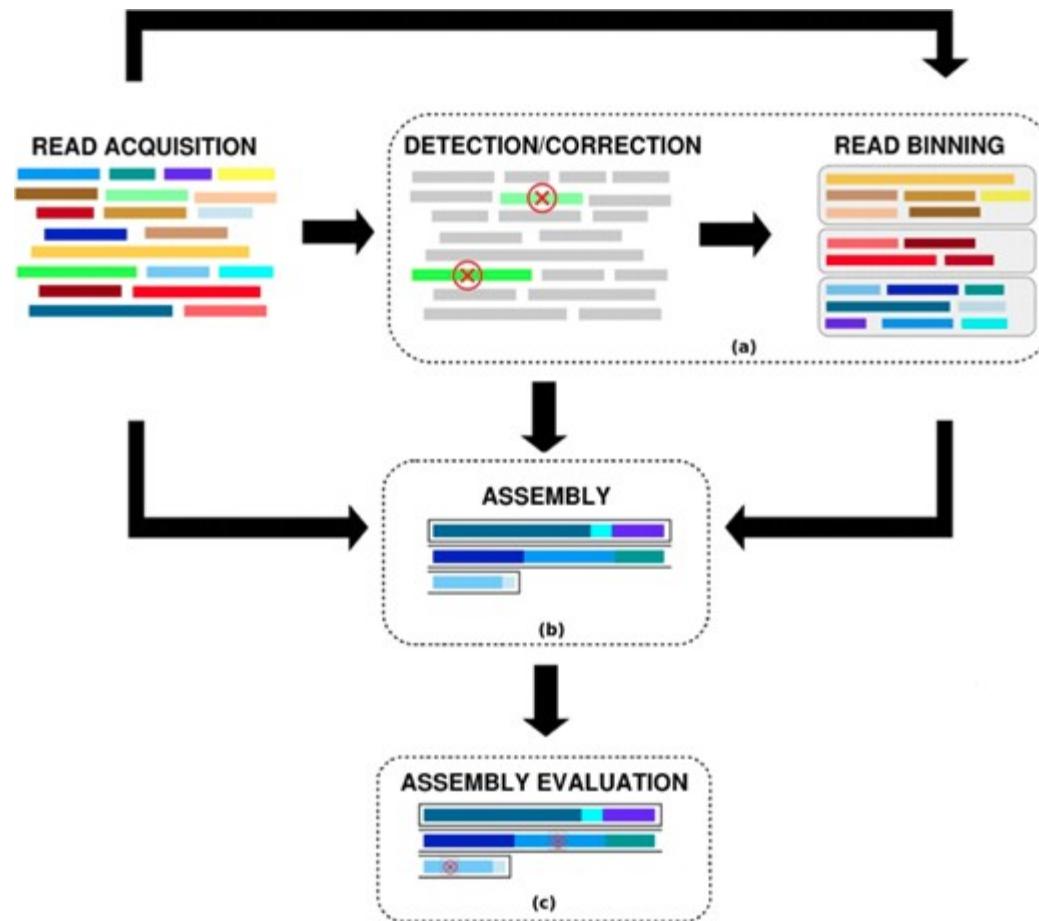
Sequencing generates 100 M reads with a length of 100 bp.



# Genome assembly : A giant jigsaw puzzle



# Genome assembly : A giant jigsaw puzzle



Step	Identification	Problem addressed	Main ML model	Data set	Simulated	Model features
<b>Real</b>						
Pre-assembly	Angeleri (1999)	Binning of reads	Artificial neural networks	✓	✗	Nucleotide sequences
	Constantinescu (2015)	Binning of reads	Artificial neural networks		✗	Nucleotide sequences
	Krachunov (2017)	Detection of sequencing errors	Artificial neural networks	✗		K-mer frequencies
Random forest						
Hoeffding Trees						
RIPPER						
Post-assembly	Choi (2008)	Assembly evaluation	Decision tree	✓	✗	Error identification metrics (e.g., RCN, RCX, CE, etc.)
			Random forest			
			Random trees			
			Bayesian networks			
			Naïve bayes			
	Lanc (2013)	Assembly evaluation	K-Means	✓	✗	Adjusted AMOS validate features
	Bodily (2014)	Contig evaluation for scaffolding step	Binary classifiers		✗	Abundance information (e.g., contig size and read coverage)
	Kuhring (2015)	Assembly evaluation	Random forest	✓		Contig size, read coverage, read length, read quality, etc.
Auto-assembly						
	Palmer (2010)	Read overlap evaluation	Random tree	✓		K-mer frequencies
			Random forest			
			Naïve Bayes			
	Bocicor (2011)	Read order definition	RL		✓	Overlapping score
	Zhu (2014)	Contig extension evaluation	SVMs	✓		Branch features (coverage level, path weight, etc.)
Read assembly						
	Wang (2015)	Read assembly into contigs	Hidden Markov Models	✓		Target-gene features
	Afiahayati (2015)	Read assembly into contigs and scaffolds	SVMs	✓	✗	Number, length and coverage of sequences and K-mer frequency
	Ji (2017)	Read assembly into contigs and scaffolds	Unsupervised learning SVMs	✓	✗	Codon frequency and TNF



# Machine learning models in error and variant detection in high-variation high-throughput sequencing datasets

Milko Krachunov <sup>1</sup>✉, Maria Nisheva <sup>1</sup>, Dimitar Vassilev <sup>1</sup>

 [Show more](#)

<https://doi.org/10.1016/j.procs.2017.05.242>

Under a Creative Commons license

[Get rights and content](#)

open access

## 4.2 Learning input

The ML-based models would classify an input example for a given evaluated base  $r_k$ —preselected by (1) or otherwise—producing a Boolean prediction (correct or incorrect) as its output. Each example in our models would be a vector of frequency statistics about the base  $r_k$ . In particular, each number would be its frequency in a different subset of the reads.

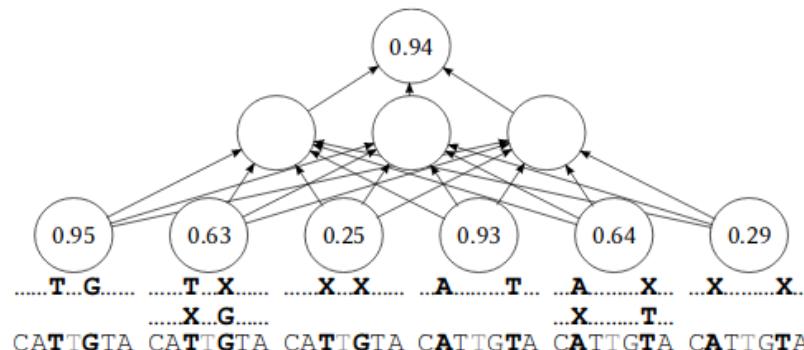
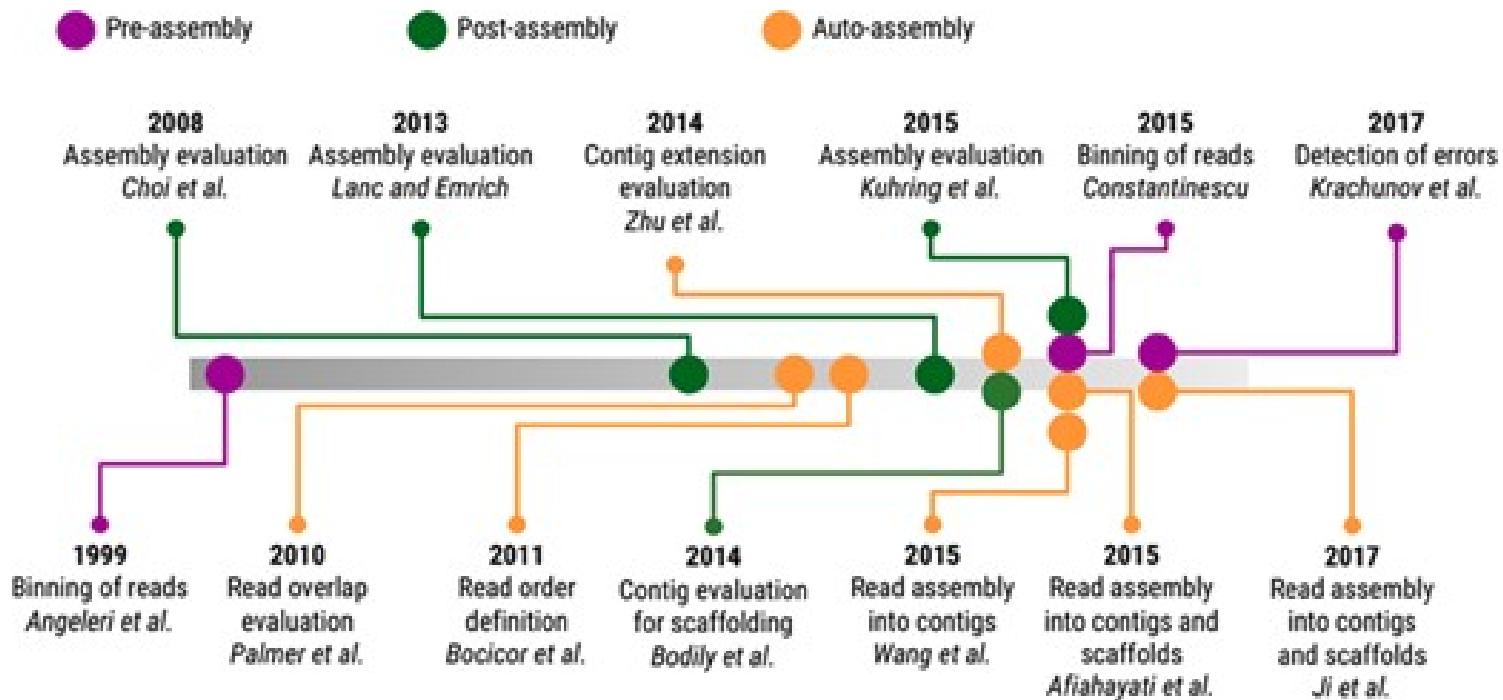
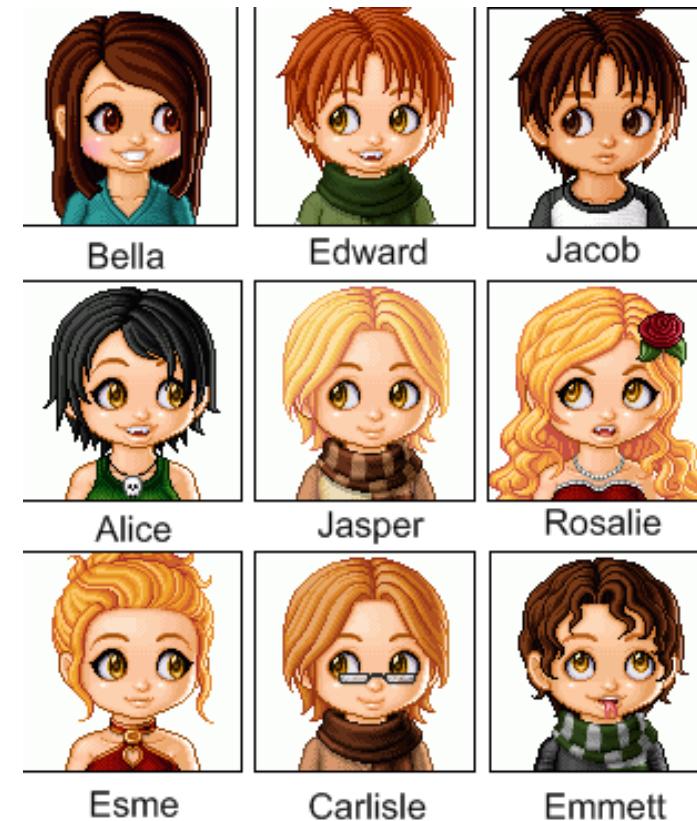


Figure 1: Example learning input

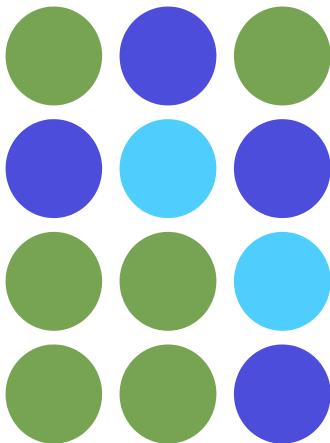
# ML in genome assembly



# Genome-wide association studies (GWAS)



Associations between phenotype and genotype



# Introduction: GWAS

•••AGCCTG-----TGCACTAAGA**Ct**•••

•••AGCCTG-----TGCACTAAGA**Ct**•••

•••AGCCTG-----TGCACTAAGA**Gt**•••

•••AGCCTG-----TGCACTAAGA**Ct**•••

•••AGCCTGAGTGTGCACTAAGA**Gt**•••

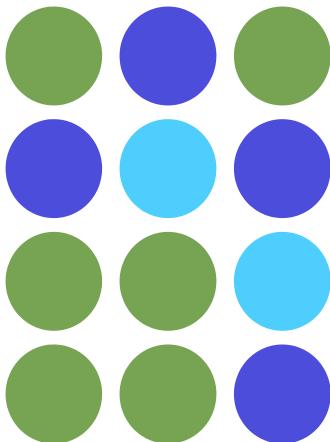
•••AGCCTGAGTGTGCACTAAGA**Gt**•••

•••AGCCTGAGTGTACTAAGA**Ct**•••

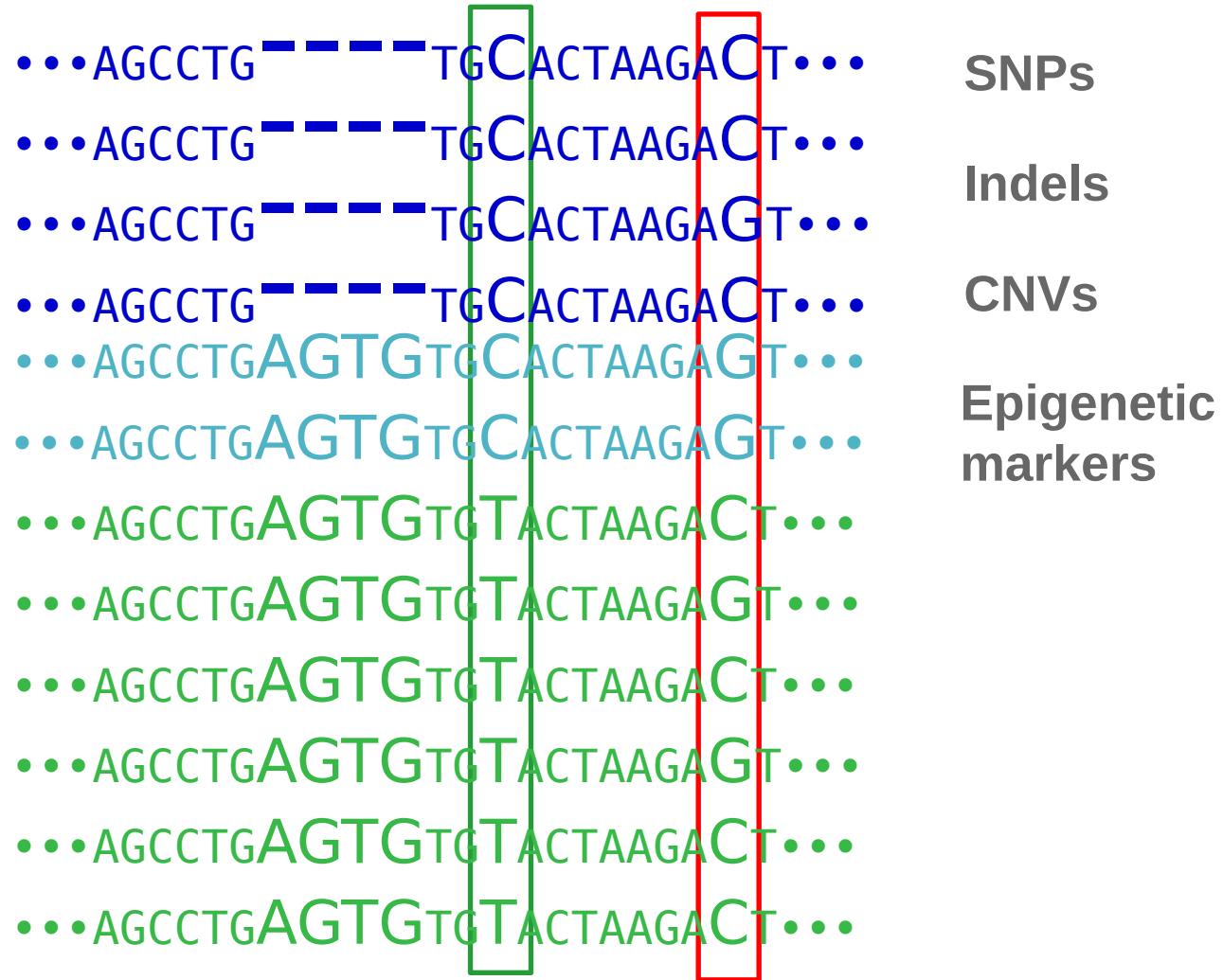
•••AGCCTGAGTGTACTAAGA**Gt**•••

•••AGCCTGAGTGTACTAAGA**Ct**•••

•••AGCCTGAGTGTACTAAGA**Ct**•••



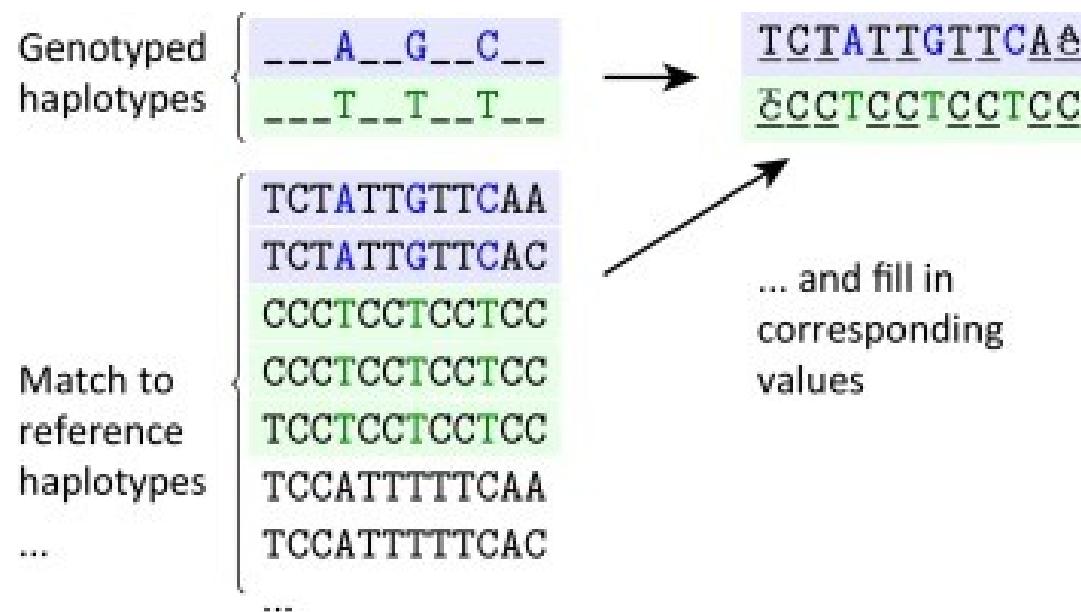
# Introduction: GWAS



# Real data are not perfect

...•AGCCTG-----TGCACTAAGACt...  
...•AGCCTG [REDACTED] TAAGACt...  
...•AGCCTG-----TGCACTAAGAGt...  
...•AGCCTG-----TGCACTAAGACt...  
...•AGCCTGAGTGTG [REDACTED] TAAGAGt...  
...•AGCCTGAGTGTGCACTAAGAGt...  
...•AGCCTGAGTGTG [REDACTED] TAAGACt...  
...•AGCCTGAGTGTGTACTAAGAGt...  
...•AGCCTGAGTGTG [REDACTED] TAAGACt...  
...•AGCCTGAGTGTGTACTAAGAGt...  
...•AGCCTGAGTGTGTACTAAGACt...  
...•AGCCTGAGTGTGTACTAAGACt...

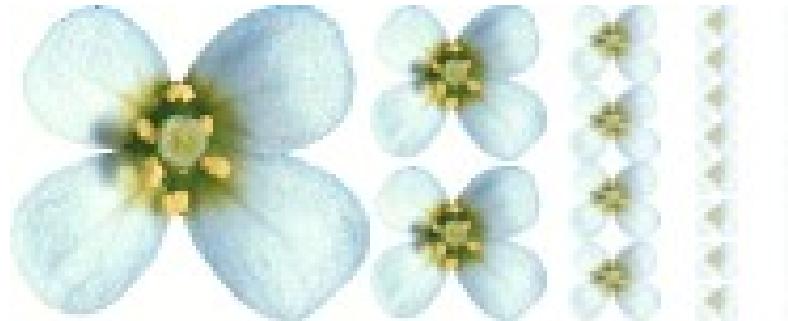
# Imputation of missing data



Trends in Genetics

**It is a simple classification problem**

# The 1001 Genomes Project



[www.1001genomes.org](http://www.1001genomes.org)

## 1001 Genomes

A Catalog of *Arabidopsis thaliana* Genetic Variation

[Home](#)

[Data Providers](#)

[Accessions](#)

[Tools](#)

[Software](#)

[Data Center](#)

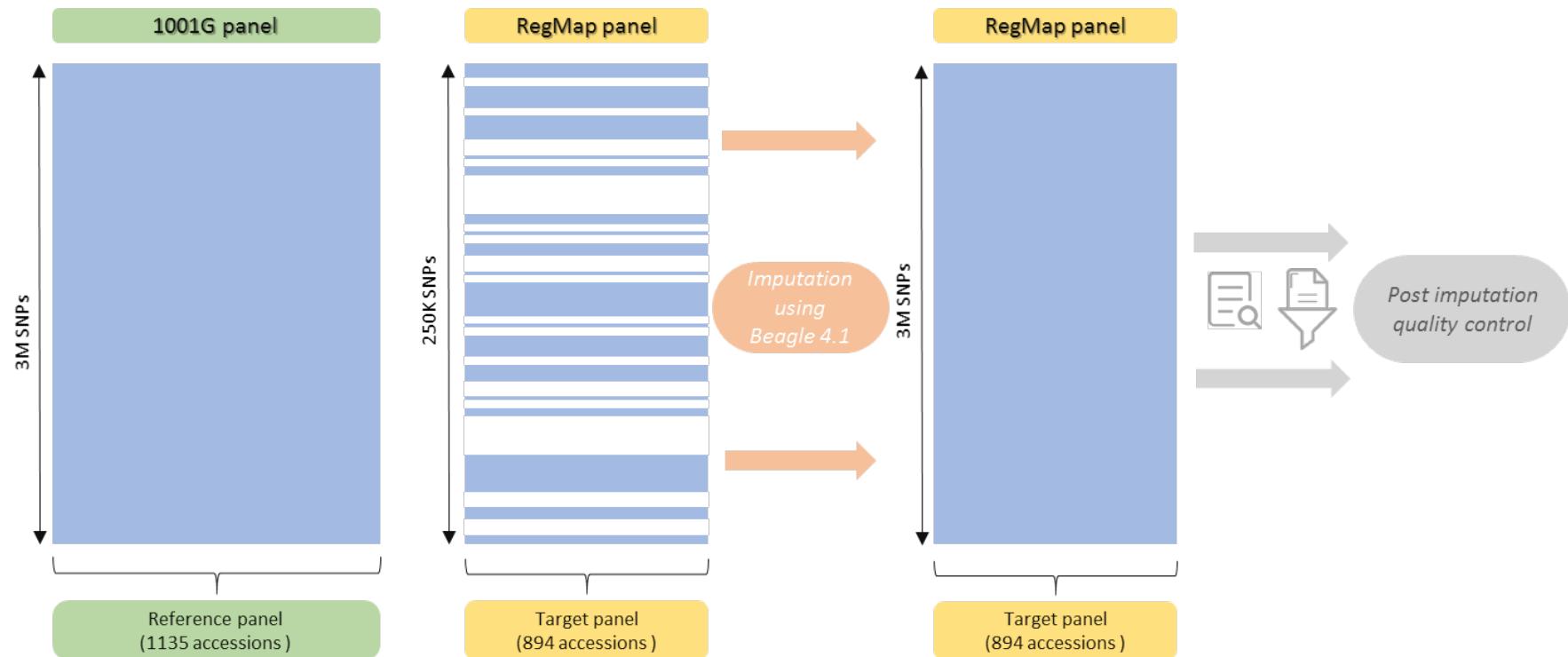
[About](#)

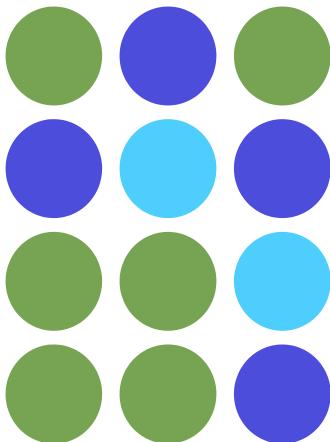
Welcome to the 1001 Genomes Project

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

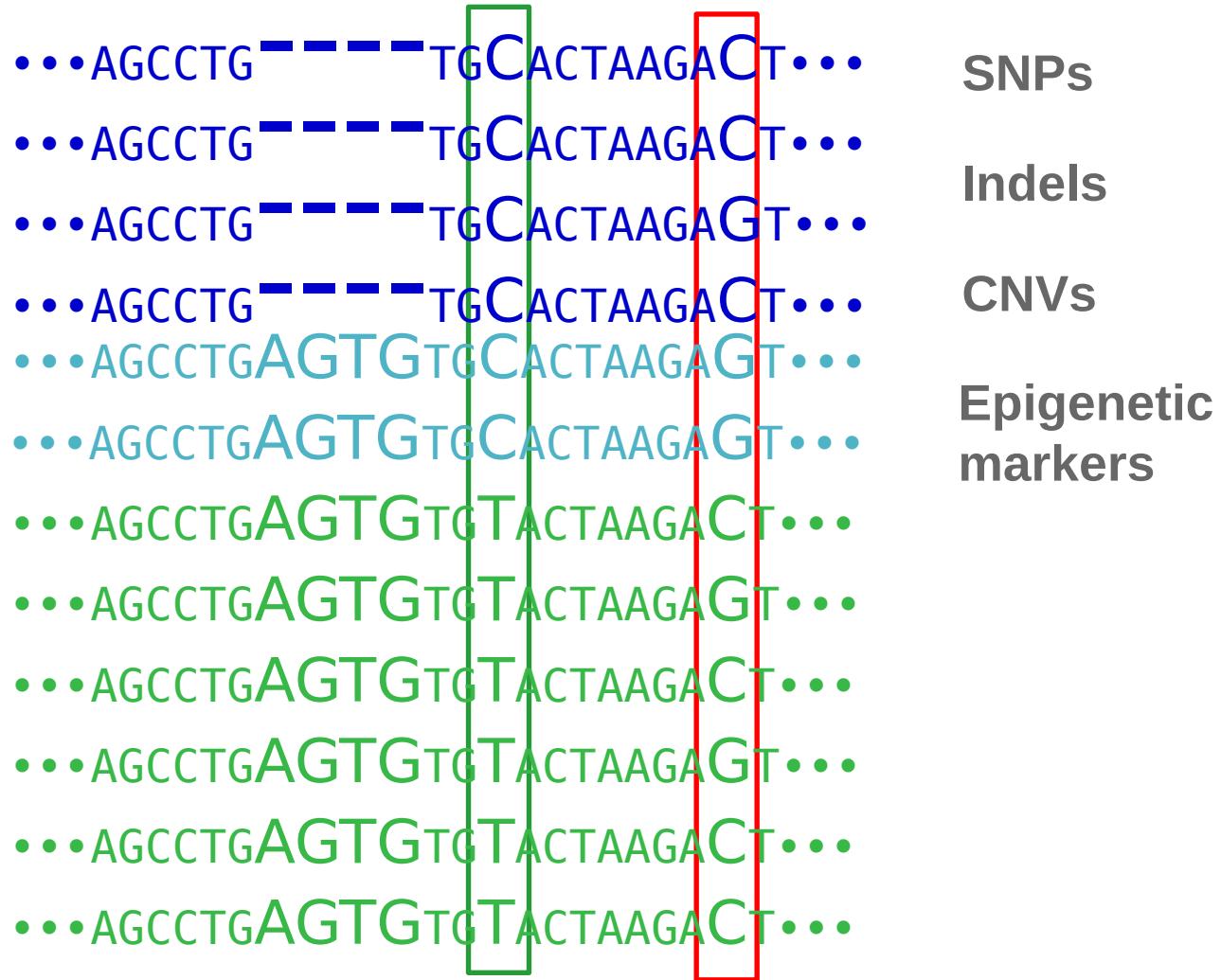
**Previous data : 1,307 SNP arrays with 250 k SNPs**

# Imputation of missing data



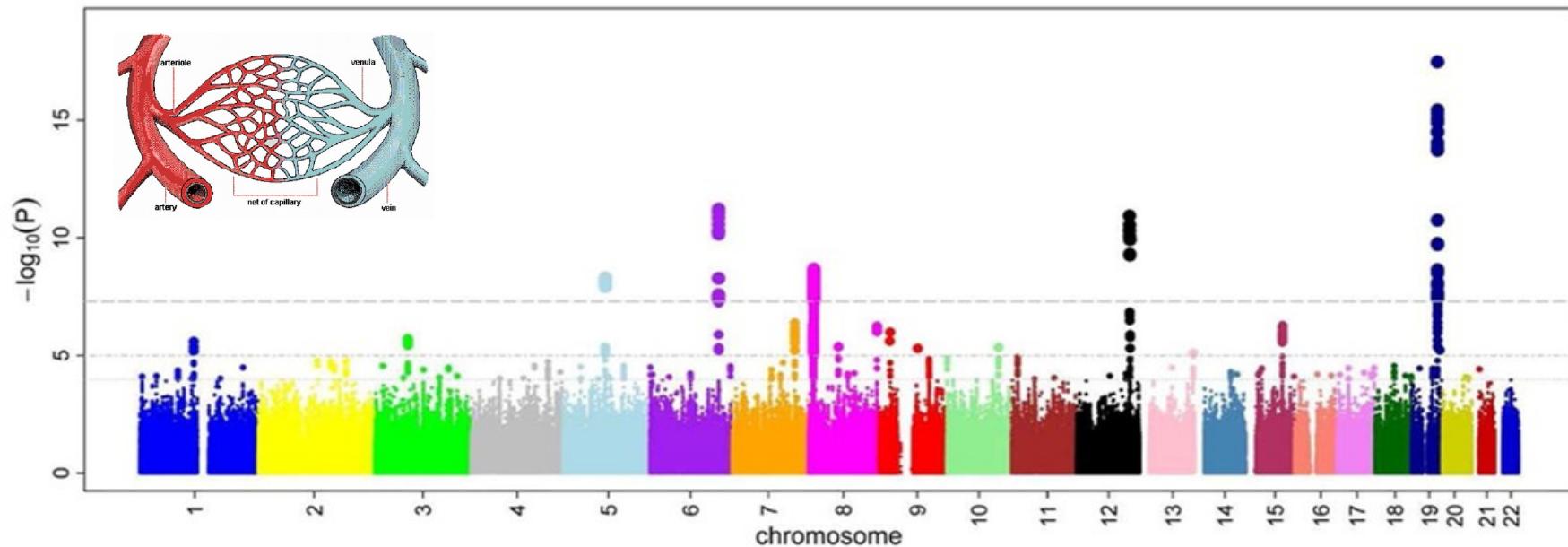


# Introduction: GWAS



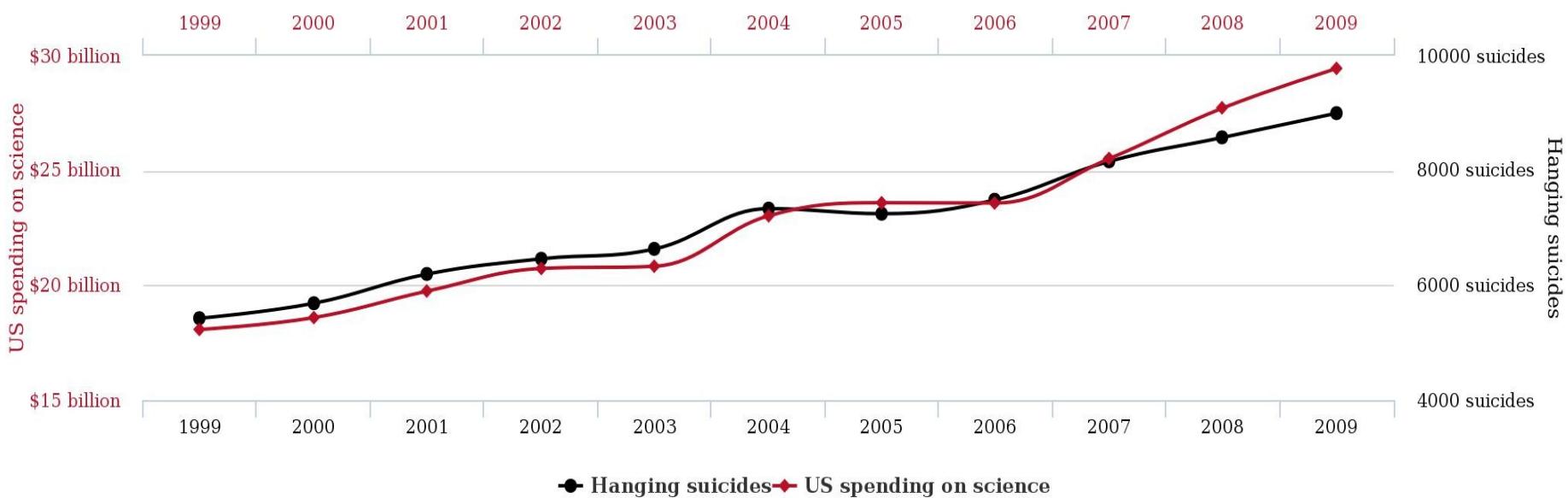
# GWAS on blood microcirculation

Manhattan plot



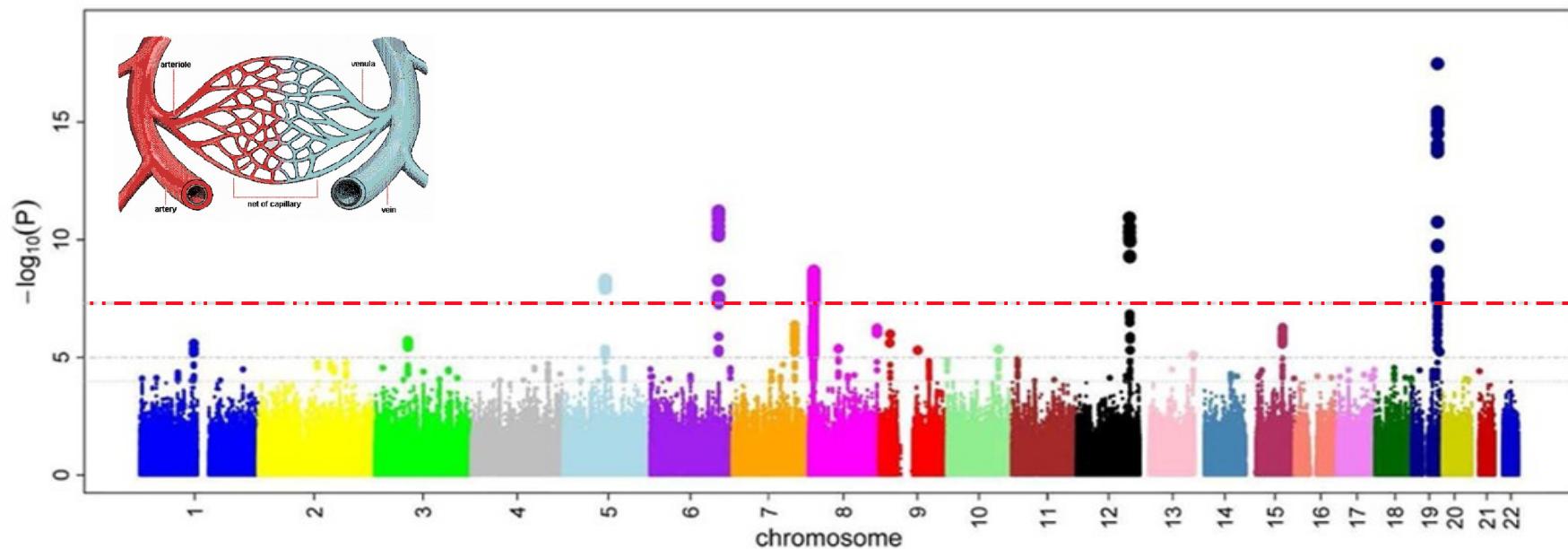
# Correlation and Causality

**US spending on science, space, and technology**  
correlates with  
**Suicides by hanging, strangulation and suffocation**



# GWAS on blood microcirculation

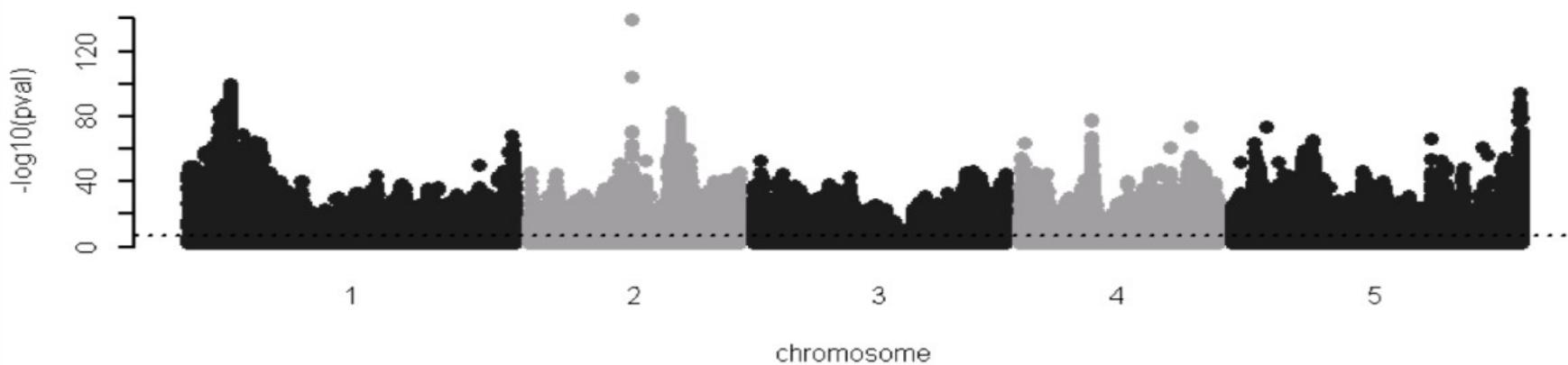
Manhattan plot



# GWAS on flowering time in *A.thaliana*



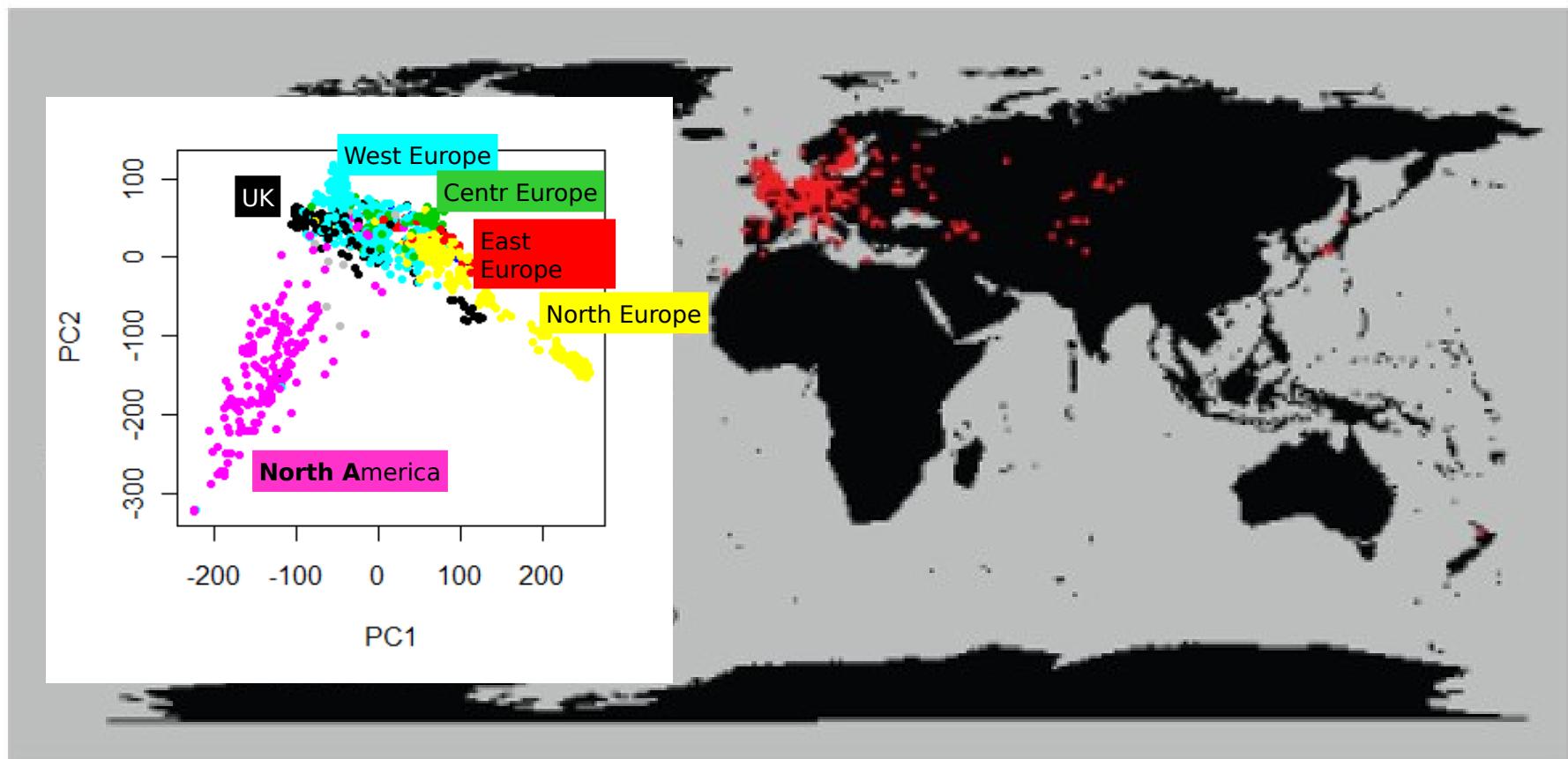
Experiment with 925 different ecotypes in a controlled environment



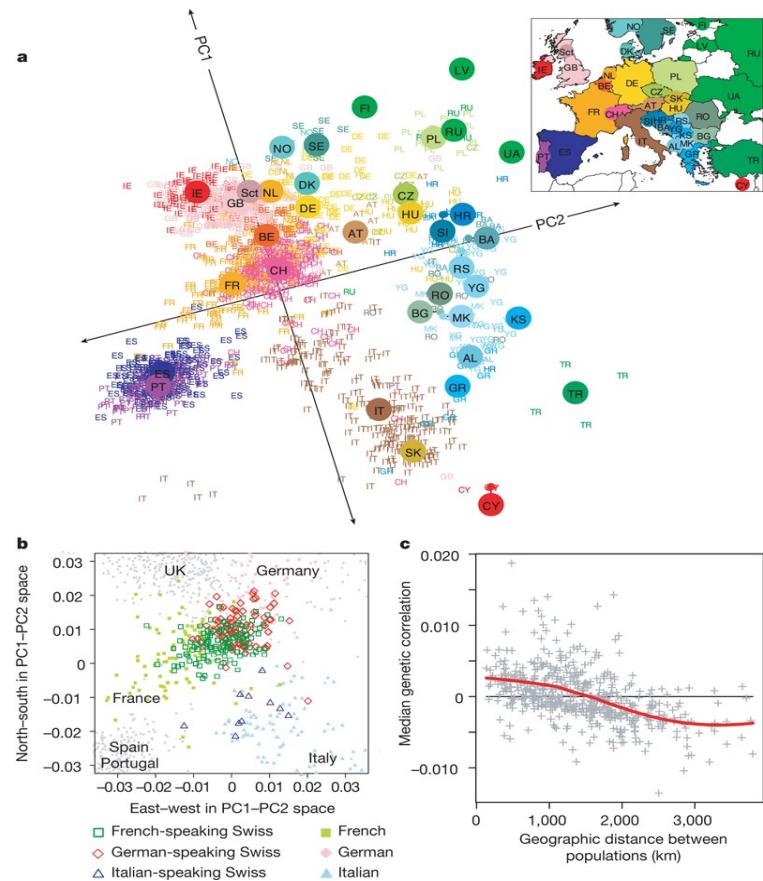
# Population structure



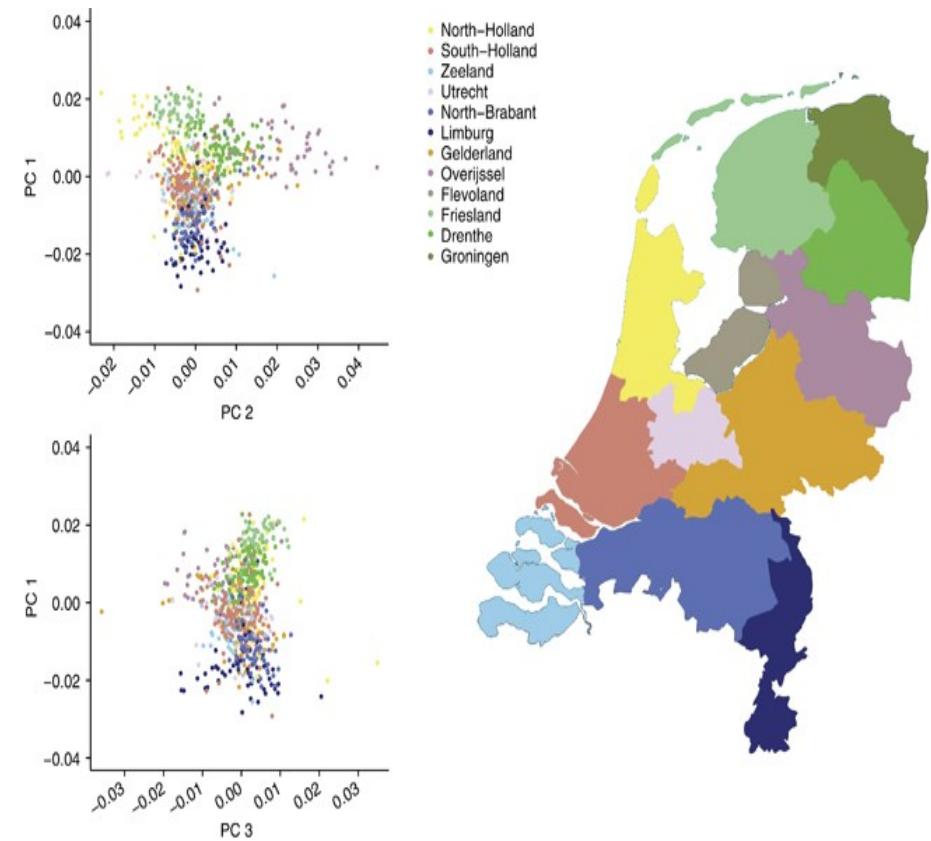
# Population structure



# Population structure in humans



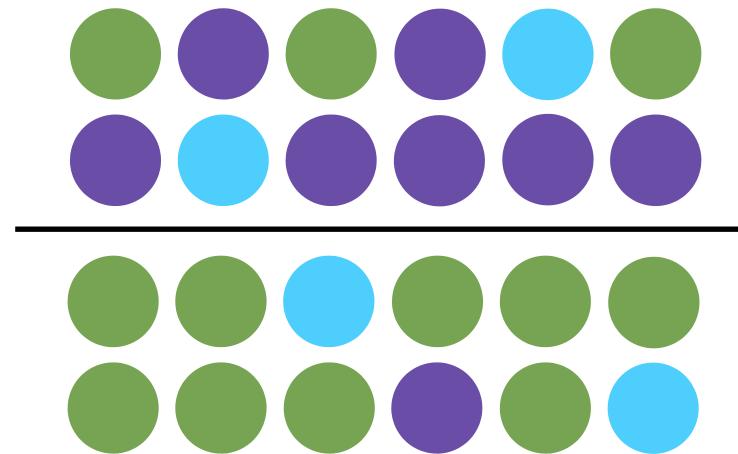
Novembre et al. (2008) Nature



Boomsma et al. (2013) EJHG

# Population structure leads to confounding

Sub-population 1



Sub-population 2

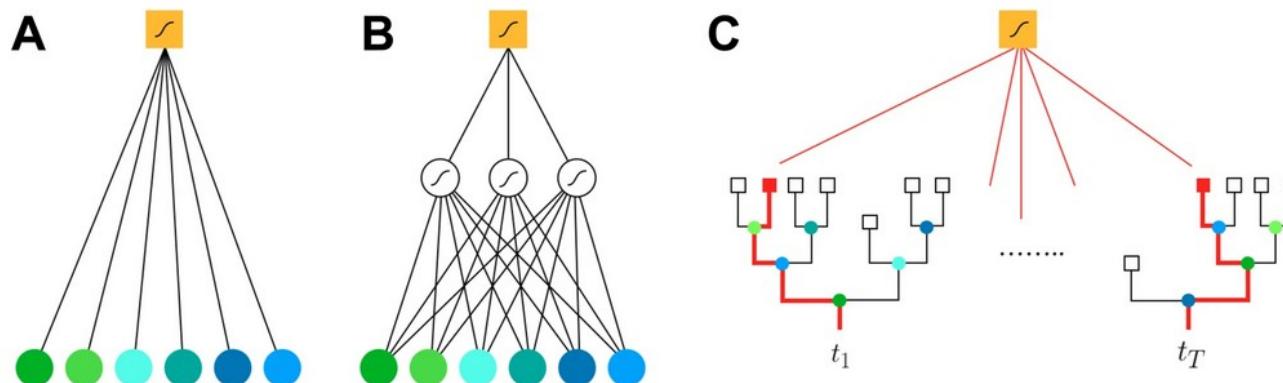
**Flowering time correlates with latitude**  
**Penetrance of many diseases differs in different**  
**ethnic groups**

# How to correct for population structure

- 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, Nat Genet.; Kang et al. 2008, Genetics)
- FaST-LMM -Select (Listgarten et al. 2013, Nature Genet.)
- ....
- ML (e.g. Rogmanoni et al. 2019, Scientific Reports)

# Comparative performances of machine learning methods for classifying Crohn Disease patients using genome-wide genotyping data

Alberto Romagnoni, Simon Jégou, Kristel Van Steen, Gilles Wainrib, Jean-Pierre Hugot  &  
International Inflammatory Bowel Disease Genetics Consortium (IIBDGC)



$$\text{Y} = \sigma\left(\sum_i w_i * x_i\right)$$

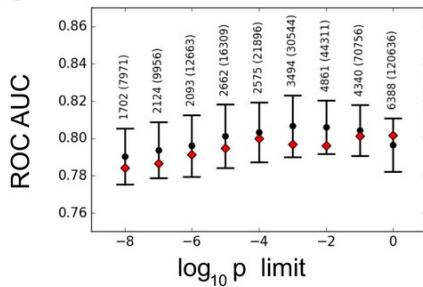
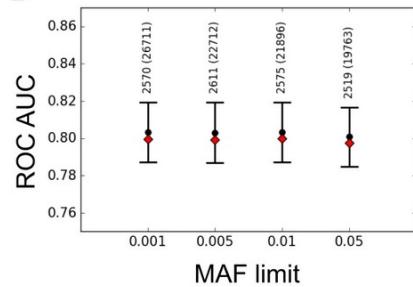
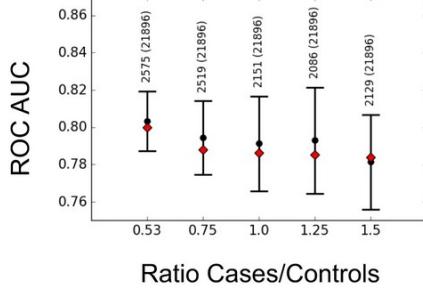
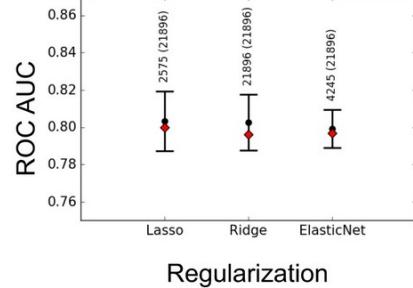
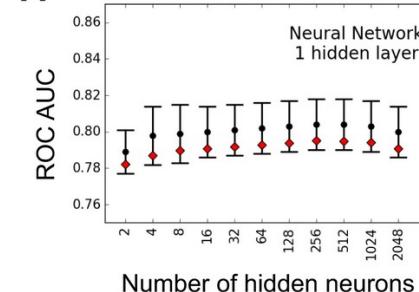
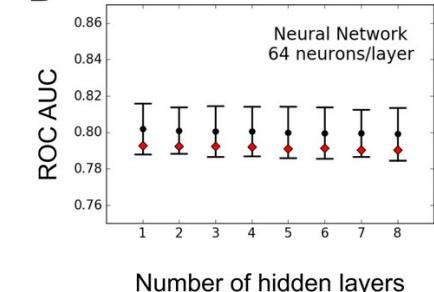
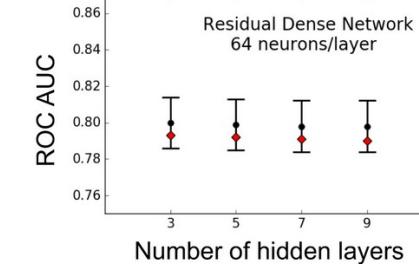
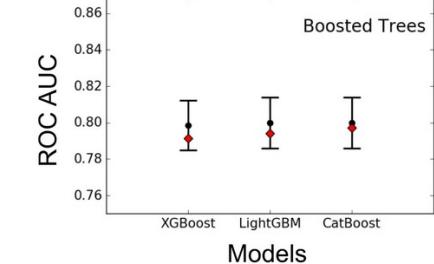
$$\odot h_j = \sigma\left(\sum_i w_i * x_i\right)$$

$$\text{Y} = \sigma\left(\sum_j w_j * h_j\right)$$

$$\blacksquare \quad \text{leaf}(t) = \text{Tree}_t(x_1, \dots, x_n)$$

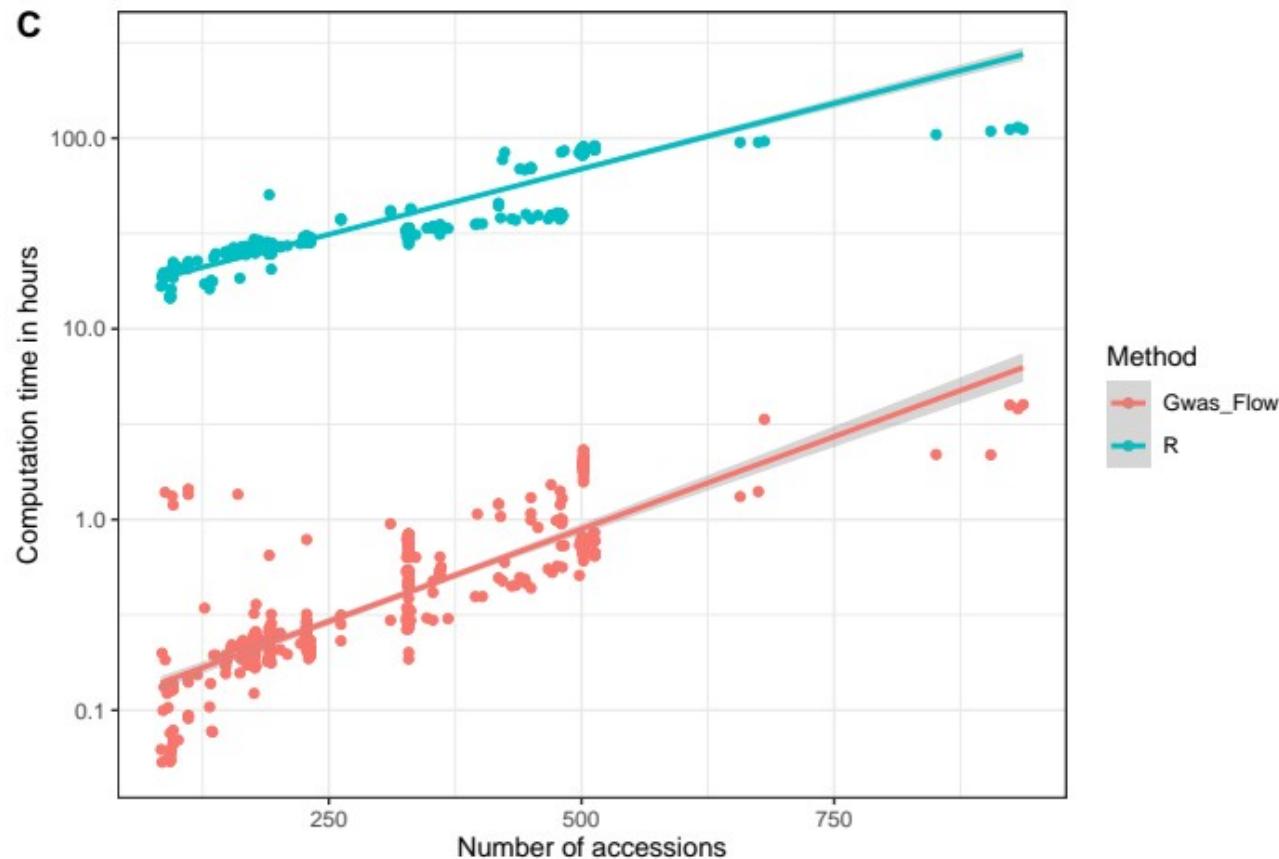
$$\text{Y} = \sigma\left(\sum_{t_i \in \text{tree}} \eta^i * \text{leaf}(t_i)\right)$$

# Network architecture and feature selection

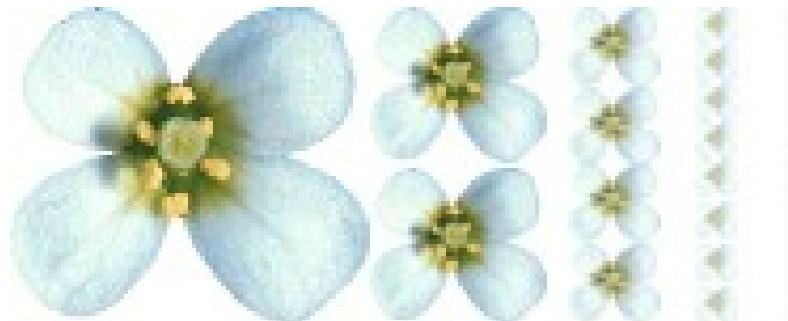
**A****B****C****D****A****B****C****D**

For GWAS ML (non-linear) and linear models  
perform on a comparable scale

# Using TensorFlow for GWAS



# The 1001 Genomes Project



[www.1001genomes.org](http://www.1001genomes.org)

## 1001 Genomes

A Catalog of *Arabidopsis thaliana* Genetic Variation

Home

Data Providers

Accessions

Tools

Software

Data Center

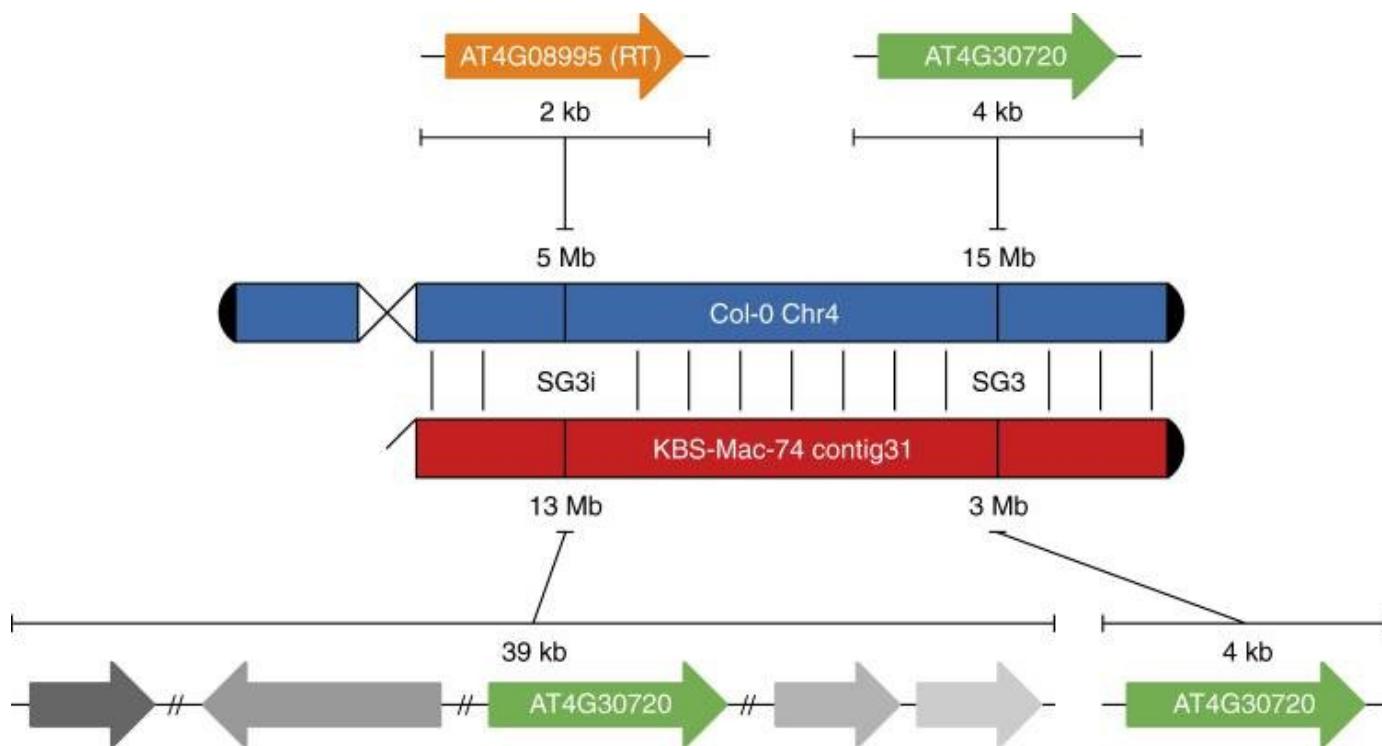
About

Welcome to the 1001 Genomes Project

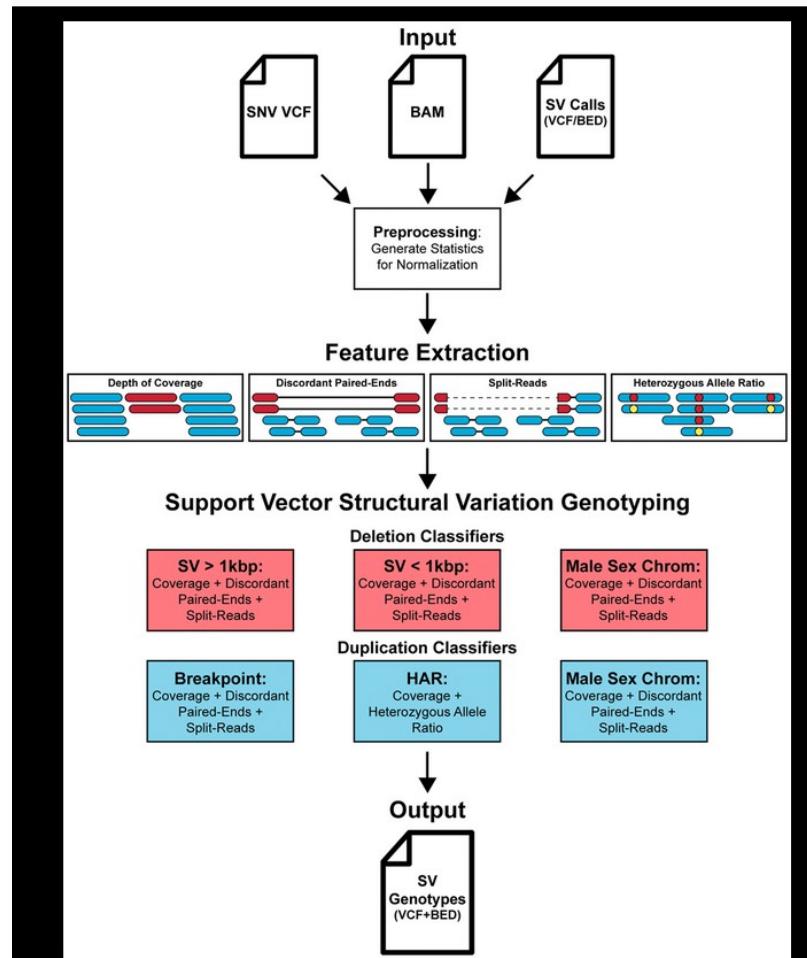
**Data : 1,135 high quality genomes with more than 10 M SNP and  
500k structural variants**

**Previous data : 1,307 SNP arrays with 250 k SNPs**

# Structural variation in the *Arabidopsis* genome



# ML for the detection of structural variants



## **Summary GWAS**

**GWAS is a powerful tool to detect genotype-phenotype association which are valuable for functional follow-up studies or trait prediction,**

**but commonly**

- only marginal effects of markers a tested
  - trait correlation is ignored
  - model assumptions might be violated

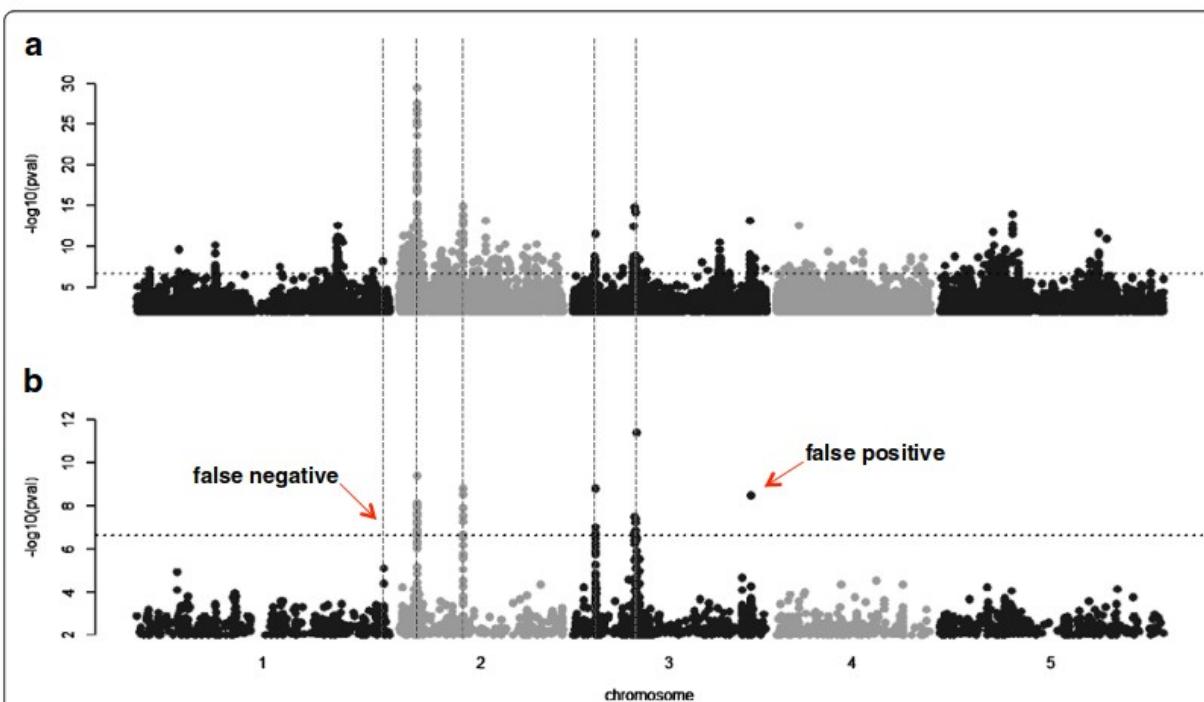
...



## REVIEW

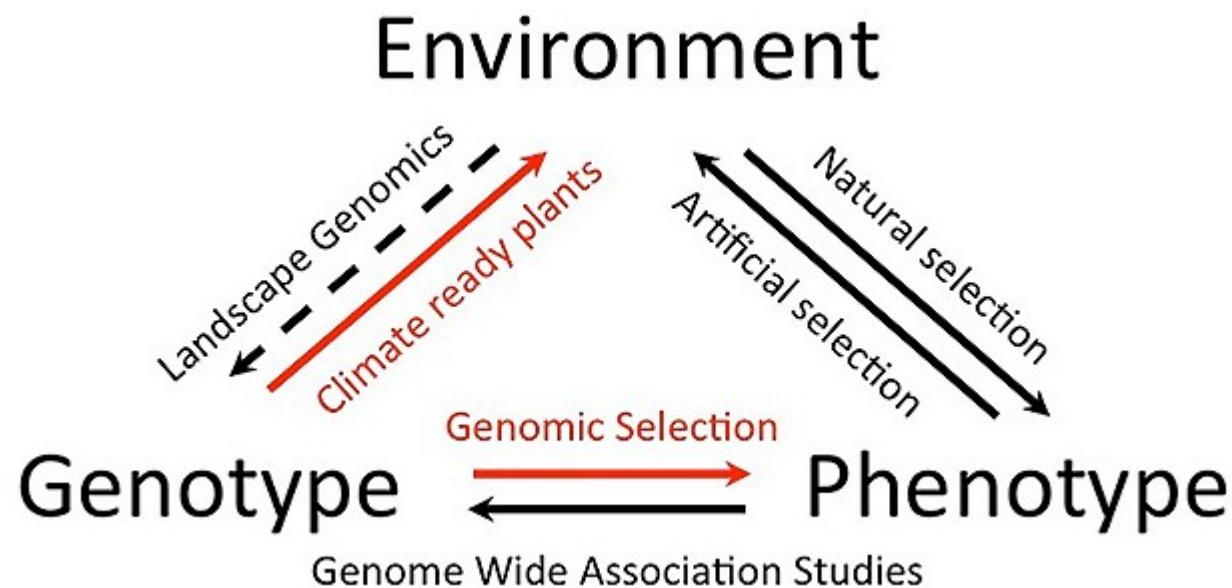
## Open Access

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

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

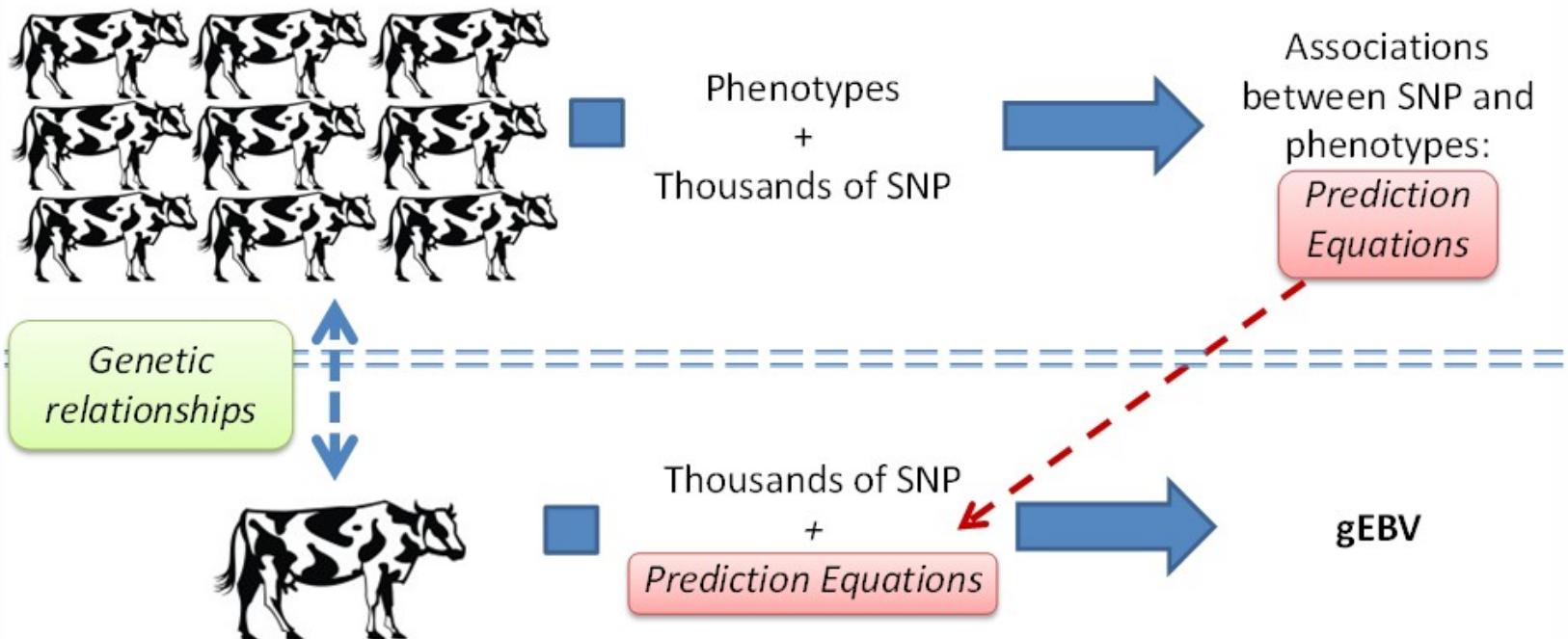
**Figure 3** Taking genetic background into account improves the performance of GWAS. Manhattan plots for a simulated trait, in which each data point represents a genotyped SNP, ordered across the five chromosomes of *Arabidopsis*. Five SNPs (indicated by vertical dashed lines) were randomly chosen to be 'causative' and account for up to 10% of the phenotypic variance each. GWAS using **a)** a linear model, and **b)** a mixed model that accounts for population structure and other background genomic factors. The simple linear model leads to heavily inflated p-values and the five causative markers are not the strongest associations. The mixed model is superior, but still leads to one false negative and one false positive. A dashed horizontal line denotes the 5% Bonferroni threshold.

# Evolutionary Genomics



# Genomic prediction

Reference population: Development of prediction equations



Main population: Application of prediction equations

# Genomic Prediction

1	1	1	1	0	0	1	1	1	1	1	1
0	1	1	1	1	1	1	1	0	0	1	1
1	1	1	1	1	1	1	1	1	1	0	
1	1	0	1	0	1	1	1	1	1	1	
1	1	1	1	1	1	1	1	1	1	1	
1	1	1	1	1	1	1	1	1	1	1	
1	1	1	1	1	1	1	1	1	1	1	
1	1	1	1	1	1	1	1	1	1	0	
1	1	1	1	1	1	1	1	1	1	1	
1	1	1	1	1	1	1	1	1	1	1	
1	1	1	1	1	1	1	1	1	1	1	

→ phenotype value (e.g. height, flowering time ...)

SNP Matrix

$$\text{GWAS : } \mathbf{Y} = \beta_0 + \mathbf{X}\beta_1 + \mathbf{u} + \boldsymbol{\varepsilon}, \mathbf{u} \sim N(\mathbf{0}, \sigma_g K), \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma_e I)$$

$$\text{GP : } \beta_0 + \mathbf{X}\beta_1 + \mathbf{uZ} + \boldsymbol{\varepsilon} = \mathbf{Y}$$

Might be possible to learn the betas instead of using a Gaussian distribution

## Decomposition of the phenotypic variance

$$\sigma_P = \sigma_G + \sigma_E + \sigma_{GxE}$$

$$\sigma_G = \sigma_A + \sigma_D + \sigma_I$$

## Decomposition of the phenotypic variance

$$\text{GP : } \beta_0 + X\beta_1 + uZ + \varepsilon = Y$$

$$\sigma_g = \sigma_A + \sigma_I$$

# Decomposition of the phenotypic variance

$$\text{GP : } \beta_0 + X\beta_1 + uZ + \varepsilon = Y$$

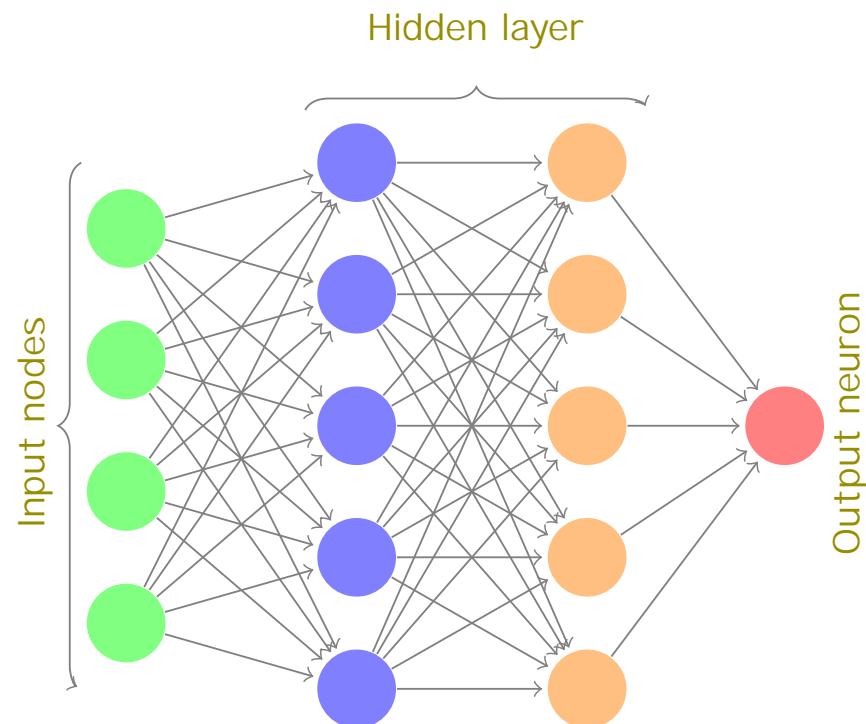
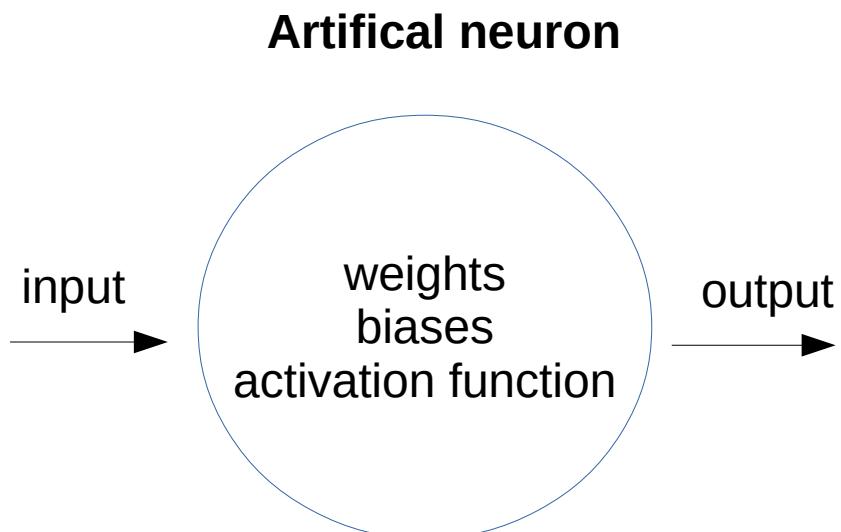
$$\sigma_g = \sigma_A + \sigma_I$$



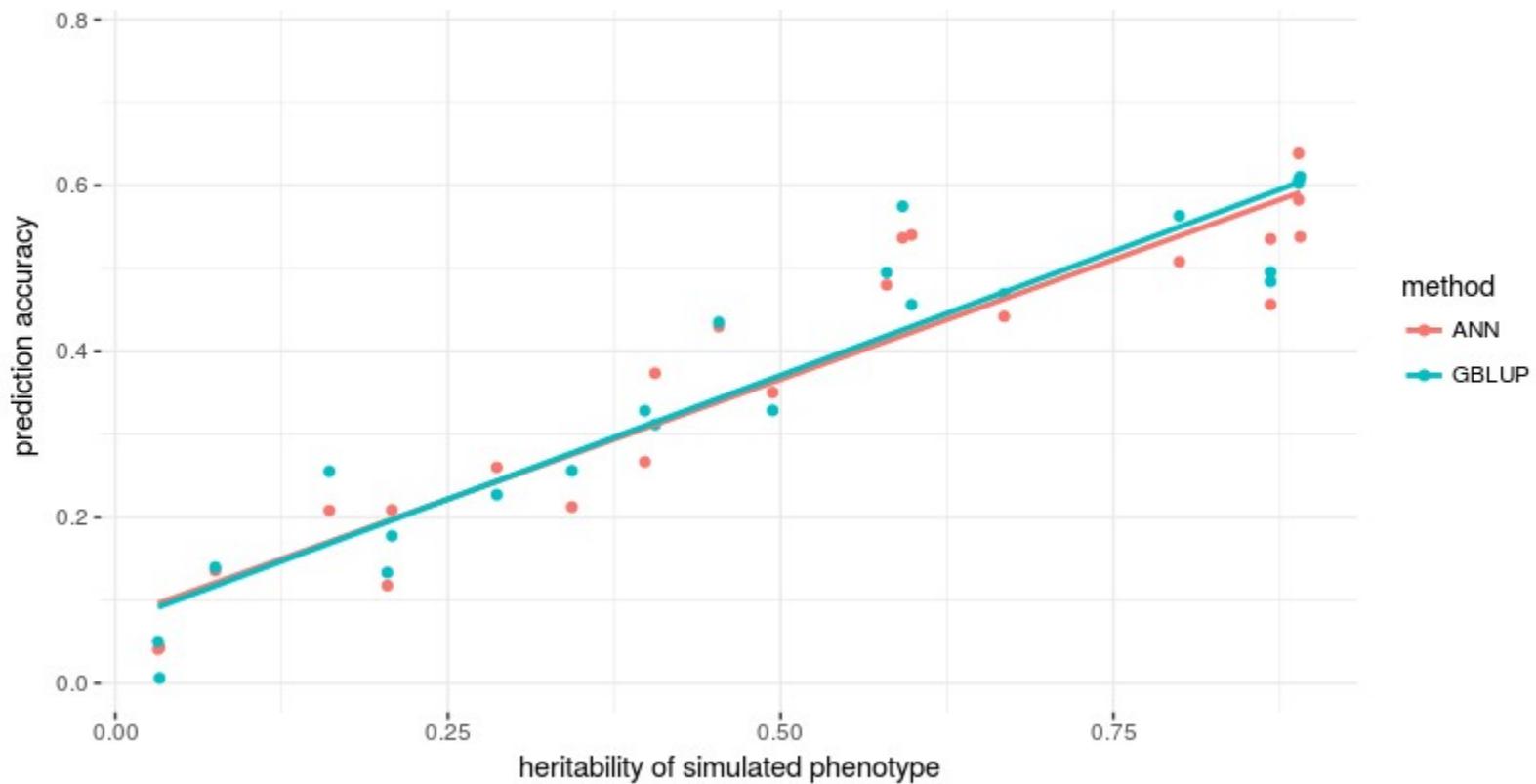
**perfect for linear models**

**Non-linear interaction,  
perfect for ML**

# ML for Genomic Prediction



# ML for Genomic Prediction



# ML does not outperform classical linear models

$$\text{GP : } \beta_0 + X\beta_1 + uZ + \varepsilon = Y$$

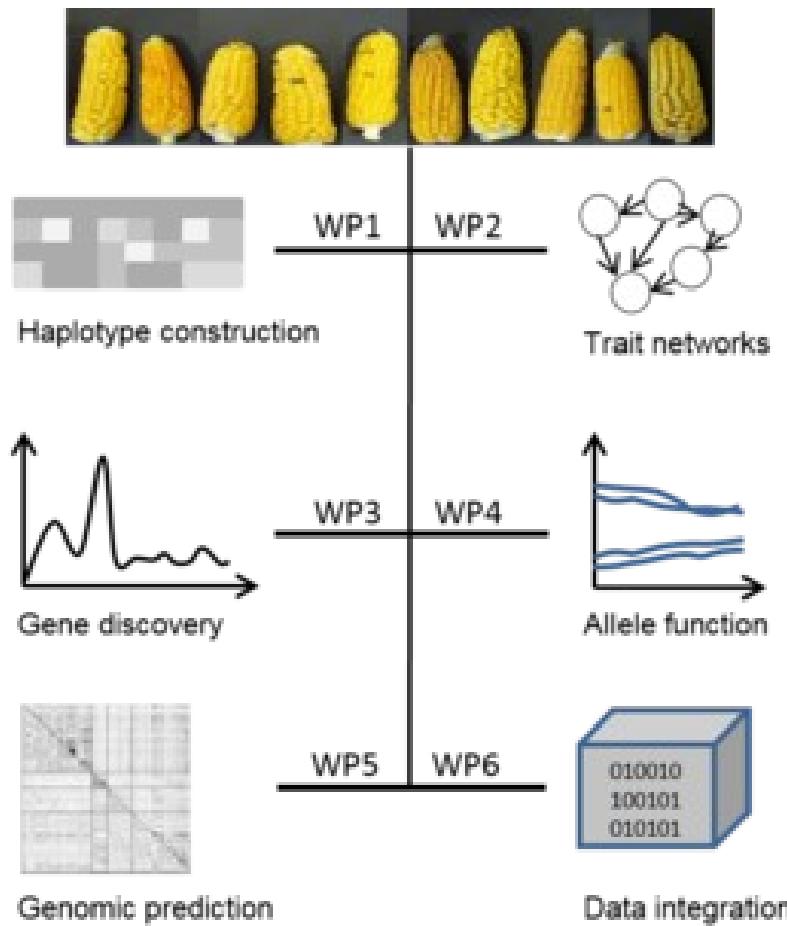
$$\sigma_g = \sigma_A + \sigma_I$$



**perfect for linear models**

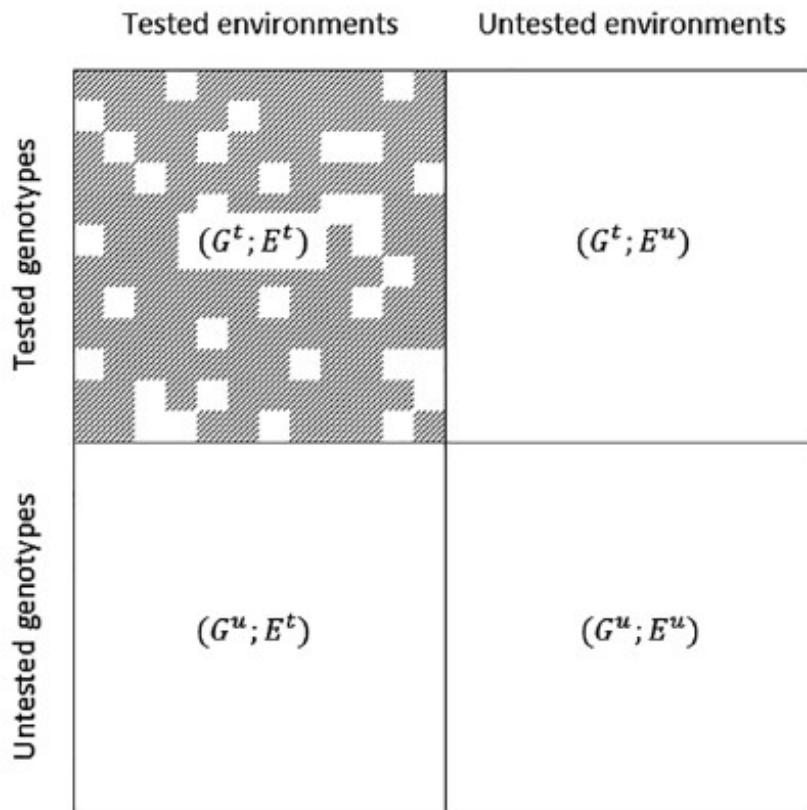
**Non-linear interaction,  
perfect for ML**

# MAZE - Assessing the genomic and functional diversity of maize to improve quantitative traits

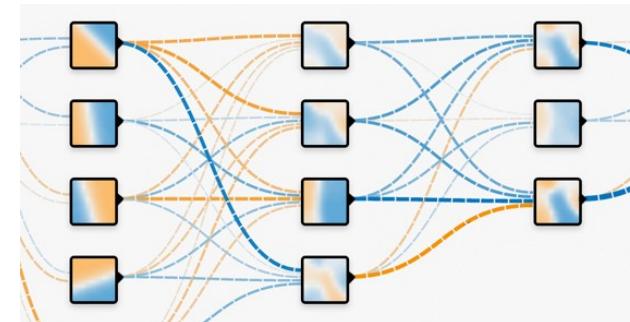


SPONSORED BY THE

# Genomic prediction

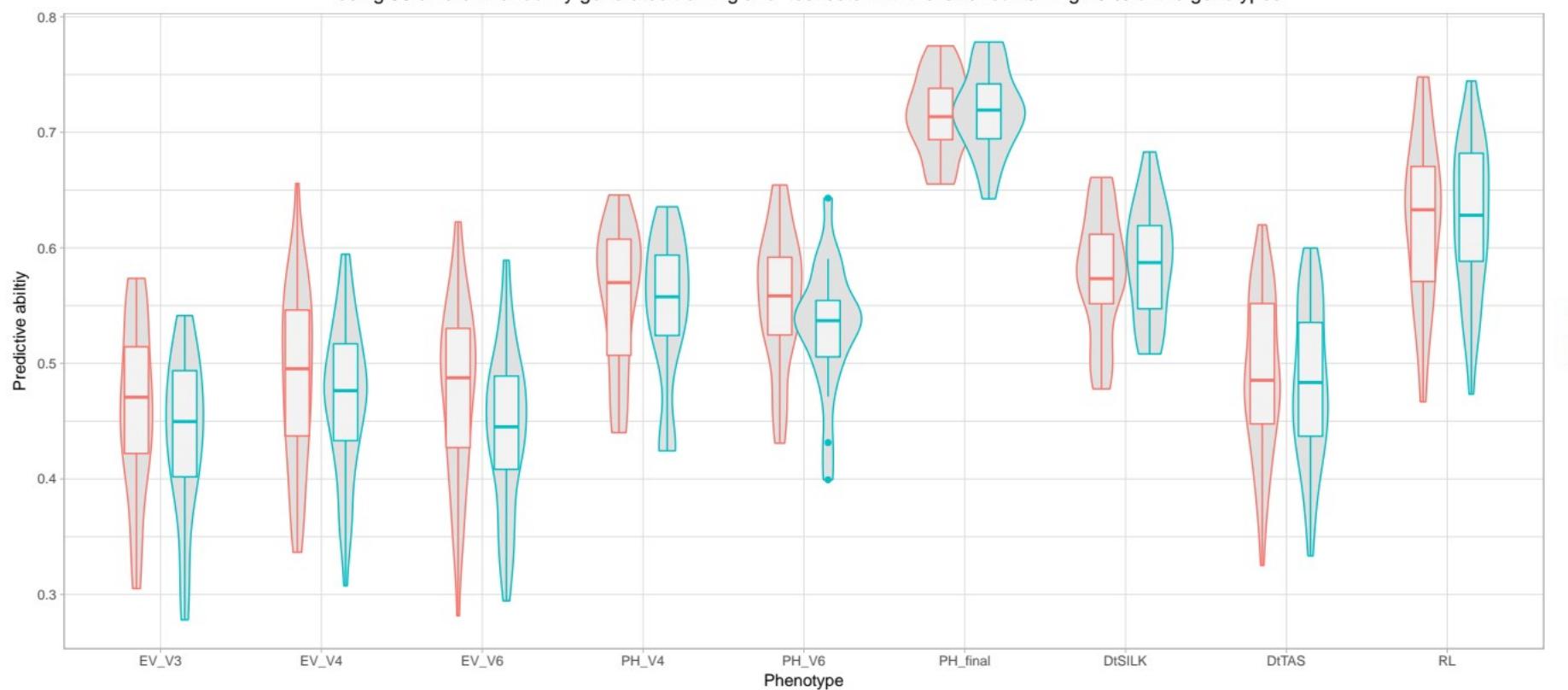


Machine learning



# ML for Genomic Prediction

Predictive abilities of bayesian ridge regression (BRR) and artificial neural networks (ANN)  
in a DH-population derived from Kemater landmaize  
using 50 different randomly generated training and test sets with the latter containing 20 % of the genotypes



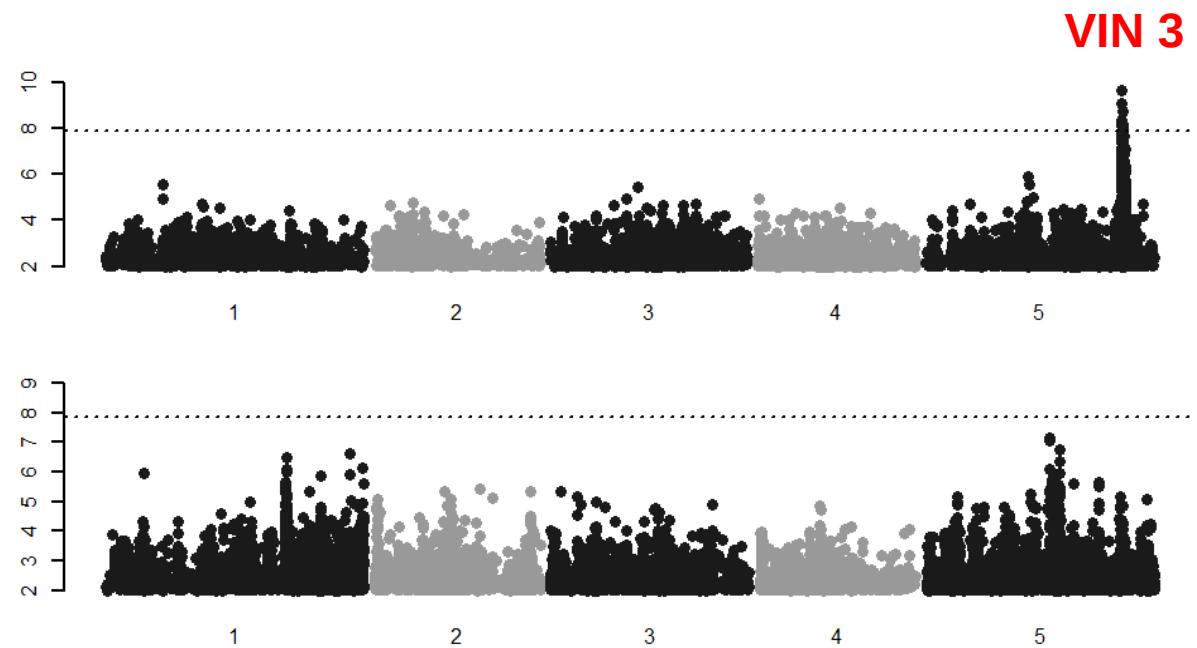
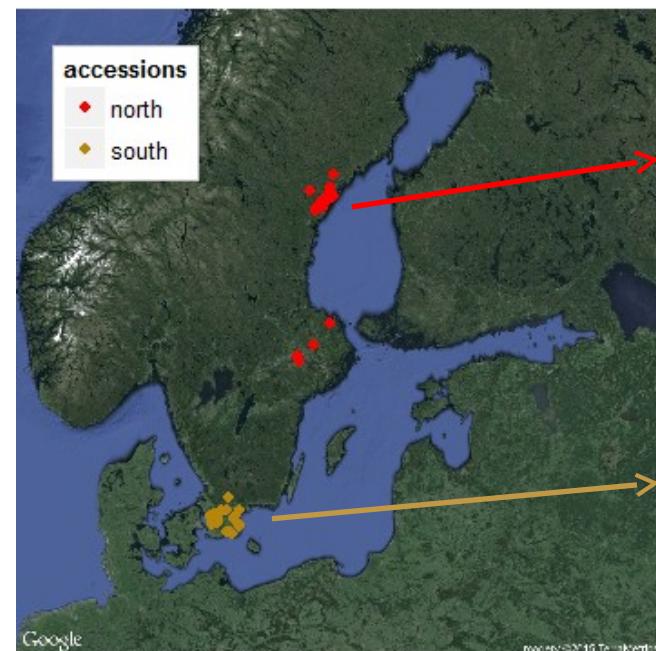
## **Summary GP**

**ML offers some potential for Genomic Prediction**

**but**

**current models do not outperform existing  
additive models**

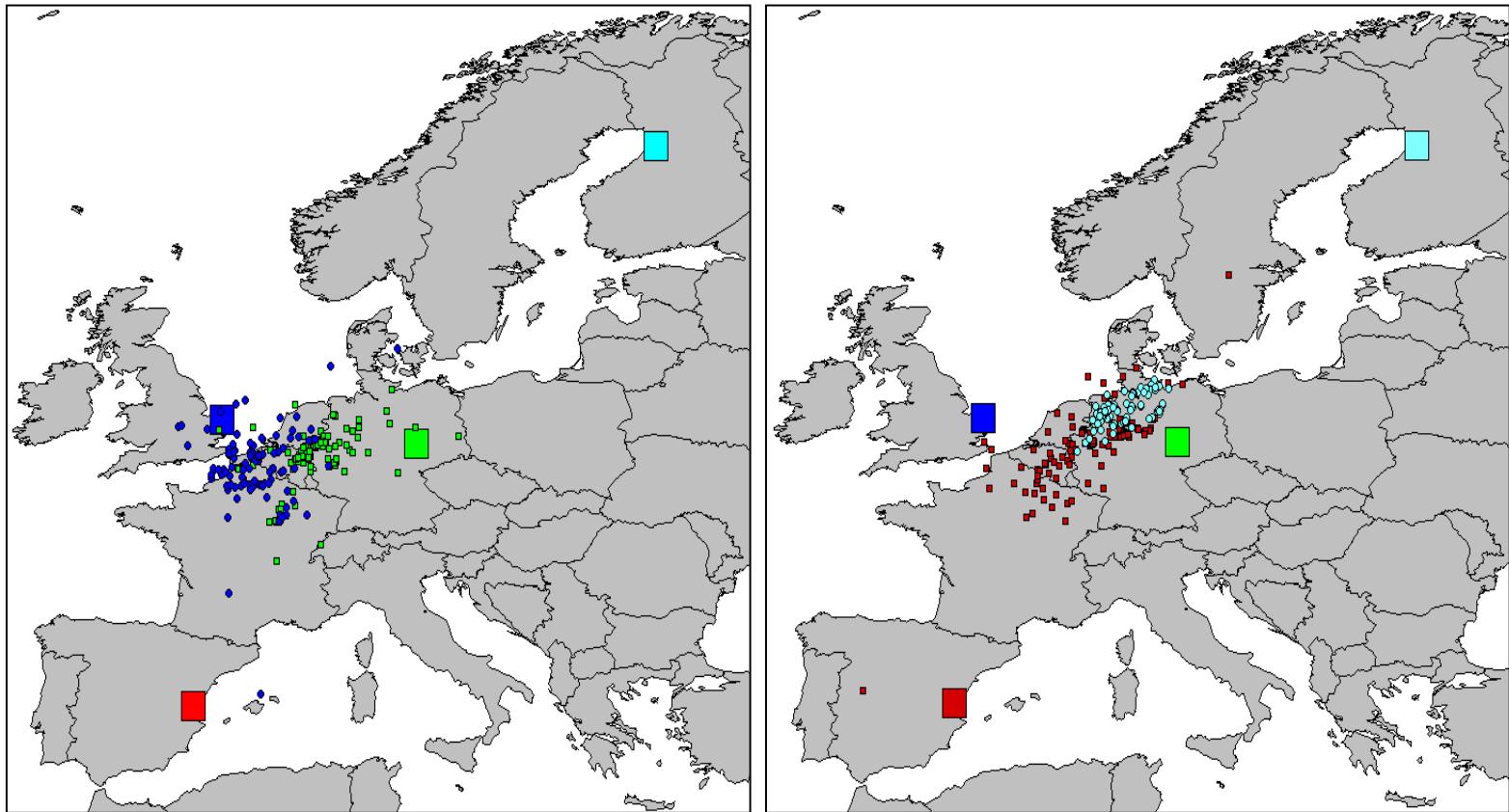
# Back to GWAS



Different results in different subsets

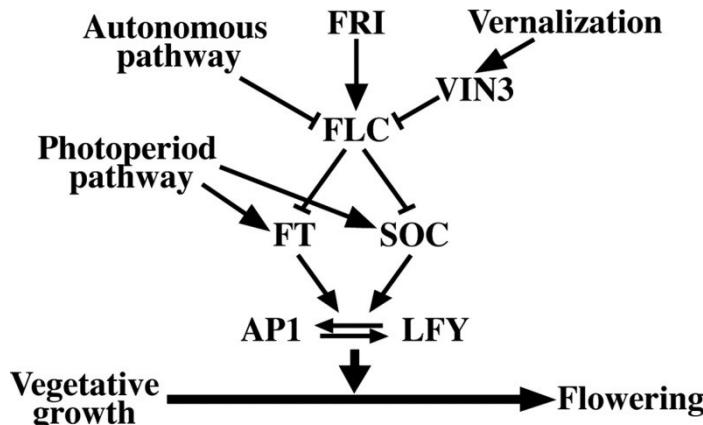
Lopez *et al.* in prep

# Local Adaptation



Genes providing fitness under natural conditions are mostly local

# GWAS in local subsets



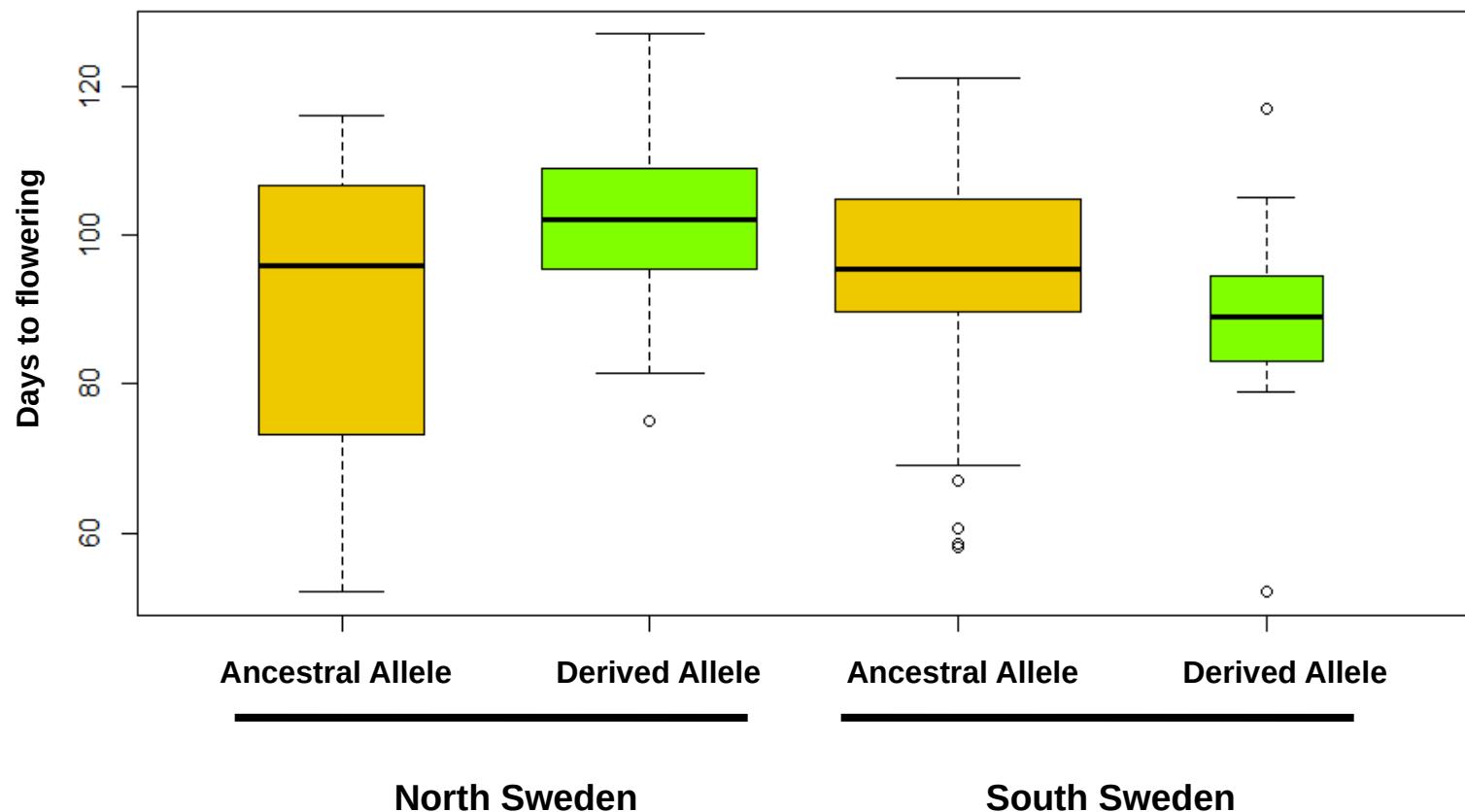
• Epistasis

Caicedo et al. (2004) PNAS

$$\text{GWAS : } Y = \beta_0 + X\beta_1 + u + \epsilon$$

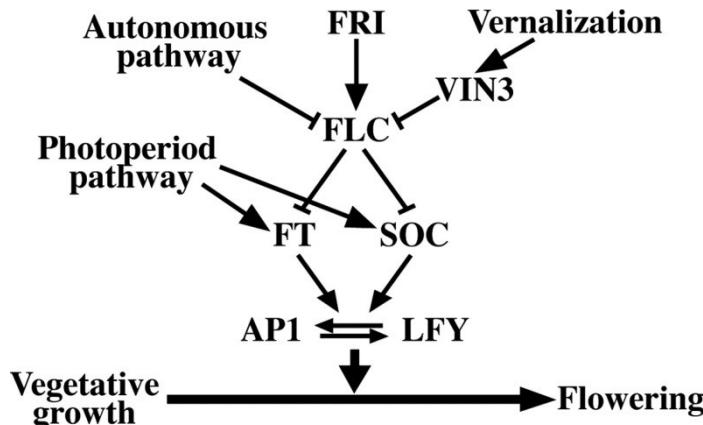
$$\text{GWAS : } Y = \beta_0 + X_1\beta_1 + X_2\beta_2 + \dots + X_1X_2\beta_{12} + u + \epsilon$$

# Effect of the *VIN3* Allele on Flowering time



The effect of the respective allele depends on the genetic background

# GWAS in local subsets



• Epistasis

$$\text{GWAS} : Y = \beta_0 + X\beta_1 + u + \epsilon$$

$$\text{GWAS} : Y = \beta_0 + X_1\beta_1 + X_2\beta_2 + \dots + X_1X_2\beta_{12} + u + \epsilon$$

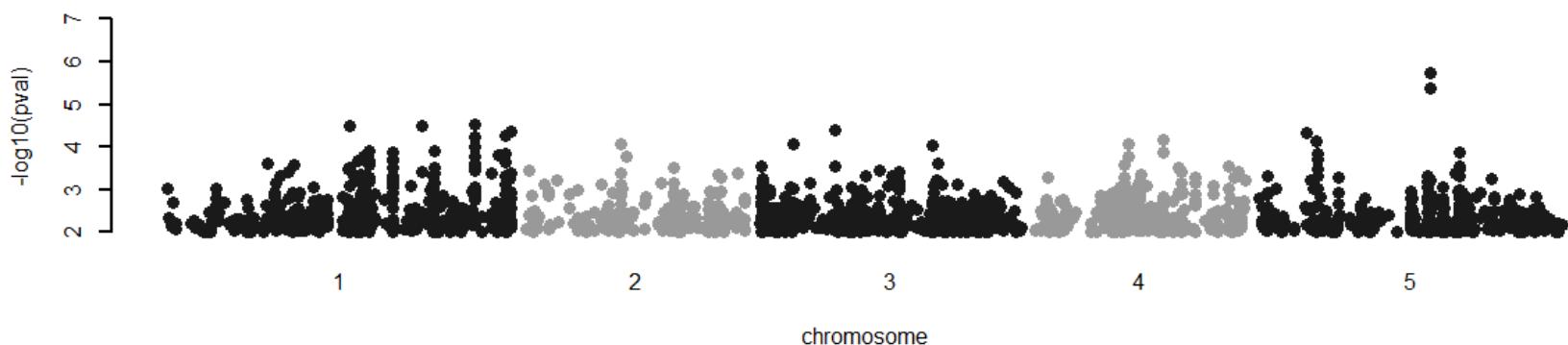
# GWAS models to test for epistasis

**SNP by SNP epistasis**

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

**3 different models with two homozygous SNPs**

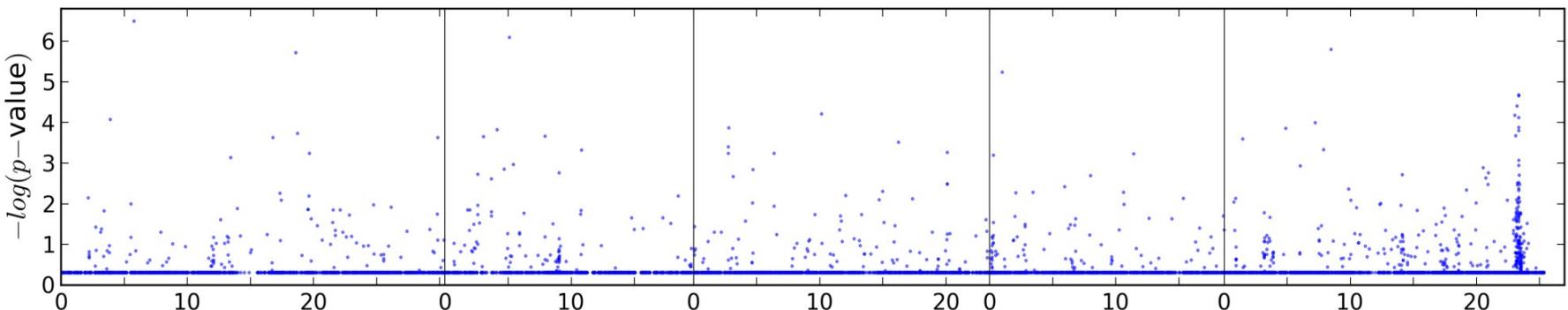
**AND, OR and XOR**



# GWAS models to test for epistasis

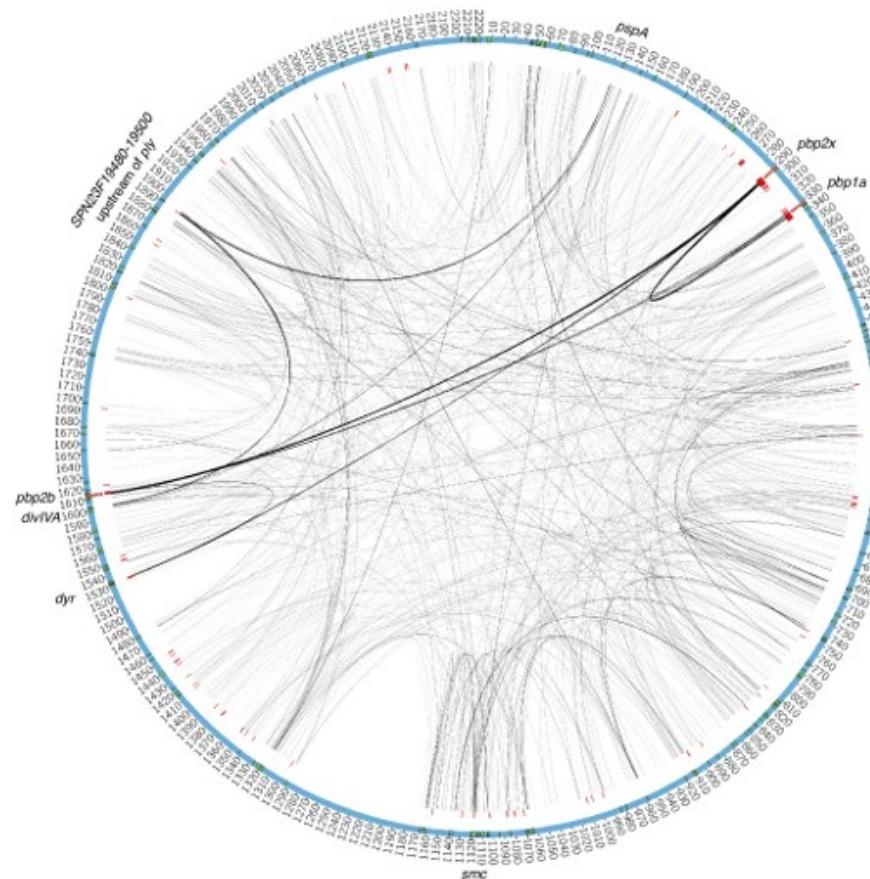
Epistasis with the genomic background in a mixed model setting

$$Y = X\beta + v + Xu + \varepsilon, v \sim N(0, \sigma_g K), \varepsilon \sim N(0, \sigma_e I)$$



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

# Epistatic interactions in a bacterial genome



# Genomic variation in the *A.thaliana* population (1135 accessions)

10,709,466 SNPs segregate in the population

1,854,599 SNPs are located in coding regions

28,148 SNPs lead to a premature STOP codon



Nearly 10,000 genes are knock-out  
in at least 1 accession

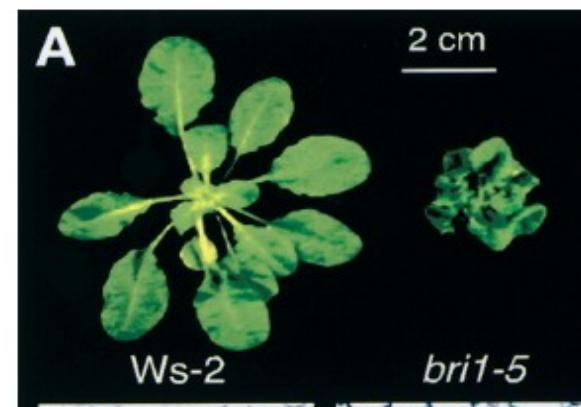
# Natural BRI1 knock out in an accession from Portugal (IP-Alo-0)



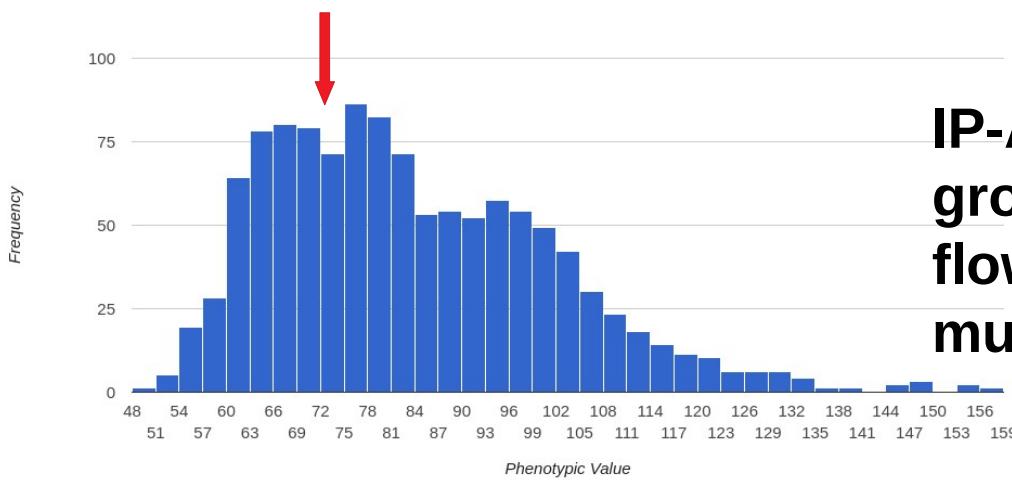
# Natural BRI1 knock out in an accession from Portugal (IP-Alo-0)



# Natural **BRI1** knock out in an accession from Portugal (IP-Alo-0)

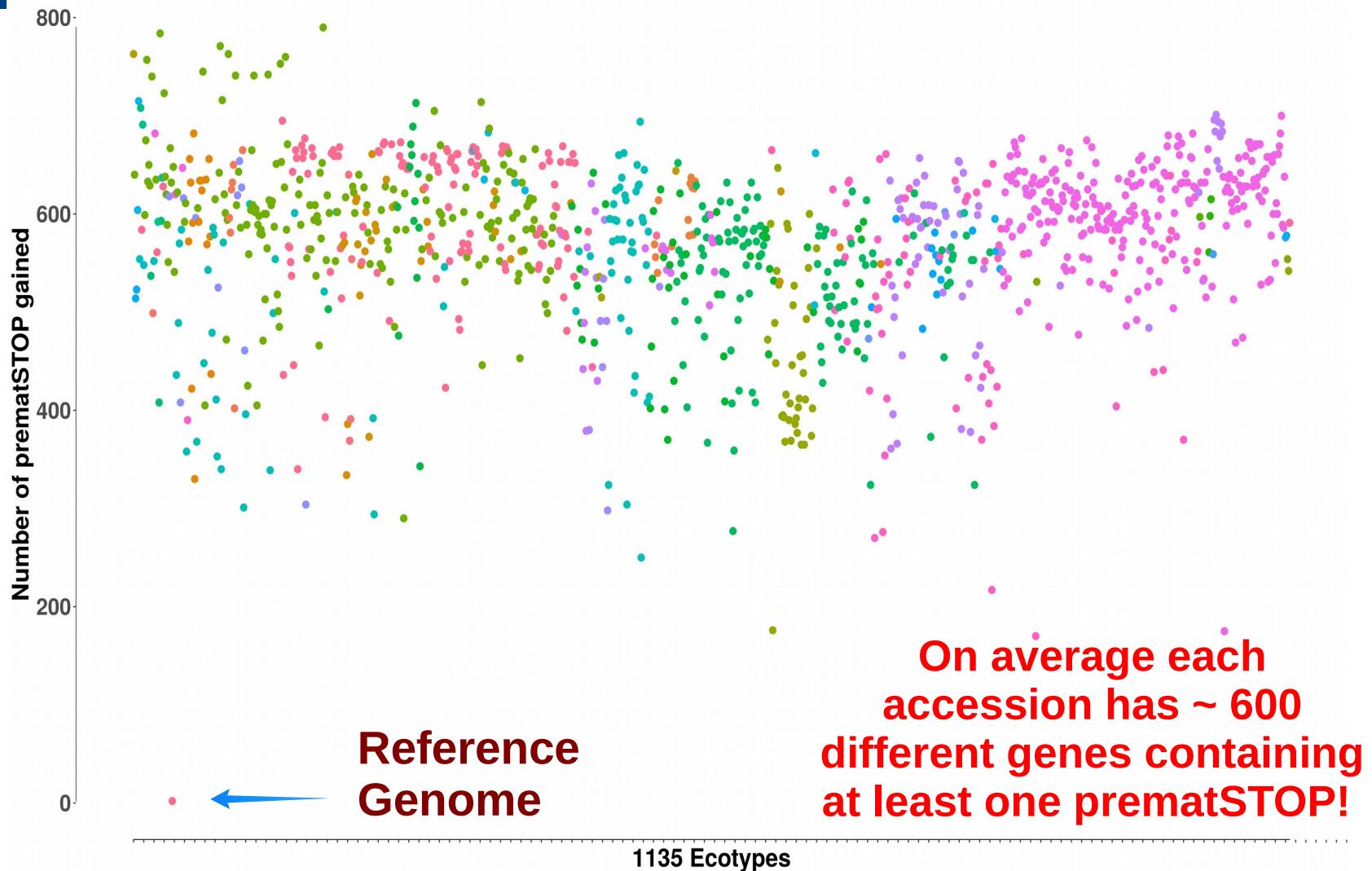


Santiago Mora-García et al. 2004, Genes Dev.



**IP-Alo-0 does neither show the growth defect nor the delay in flowering time observed in *bri1* mutants**

# Frequency of prematSTOP codons by accessions



# Co-occurrence of premature Stop codons

Filter for functional prematSTOP by using available RNAseq data (Kawakatsu *et al.* 2016)



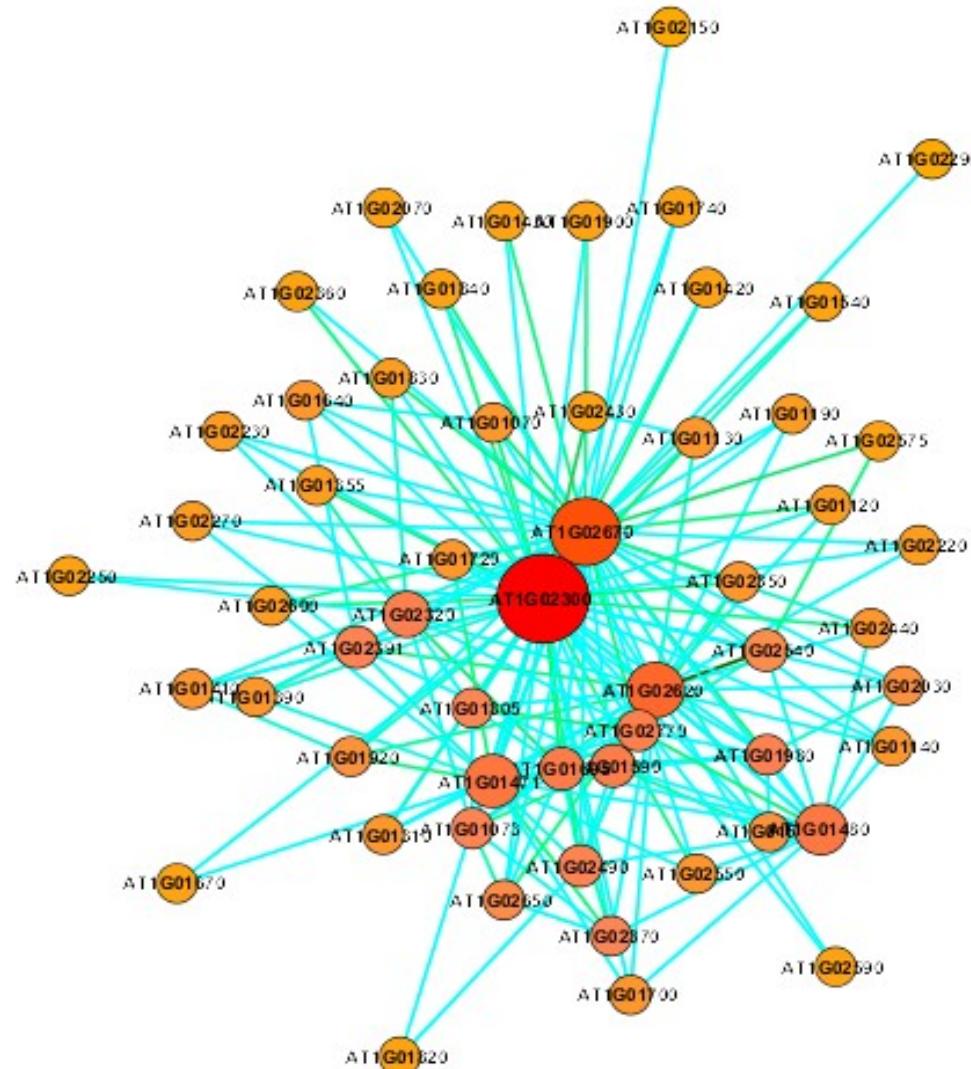
~ 6.000 genes / 16.000 prematSTOP

Gene 1	Gene 2	gene1_count	gene2_count	co-occurrence	P-value over	P-value under
AT2G25850	AT3G45910	27	72	2	0.52	0.75
AT3G55780	AT5G51000	450	435	237	<10 <sup>-20</sup>	1

Many genes are knocked-out together more (or less) often than expected

If we use synonymous SNPs as a control we don't observe this co-occurrence

# Network of epistatic interactions in *A.thaliana*



Use this network architecture as prior knowledge for GWAS

# **MTMM (Multi-trait mixed model)**

**Combine correlated traits  
(same trait in two environments | longitudinal data)**

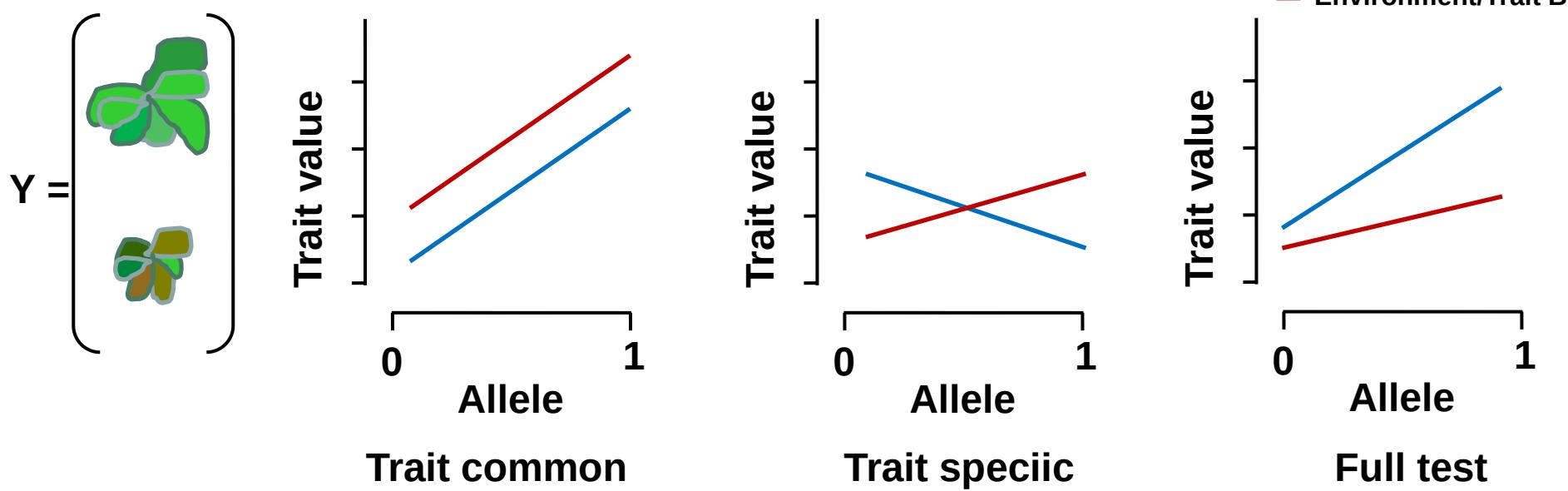
**Flowering time and Yield**

**Leaf growth in normal and water limiting conditions**

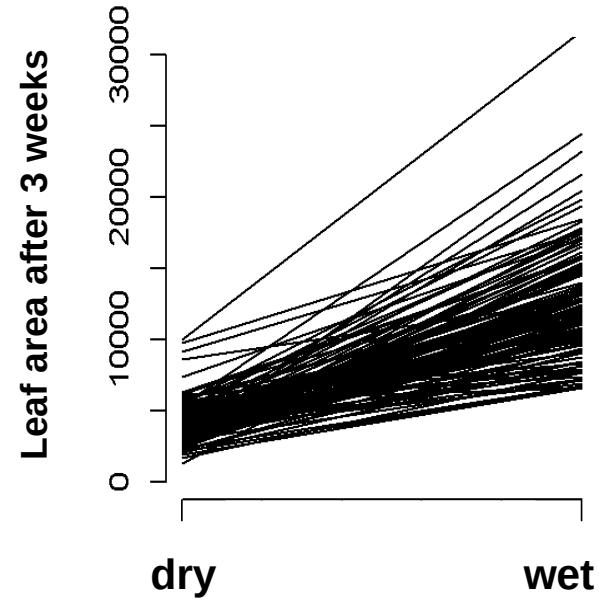
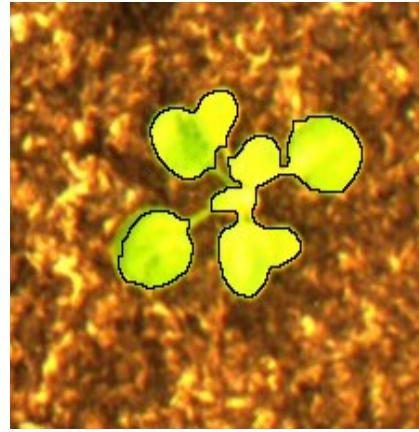
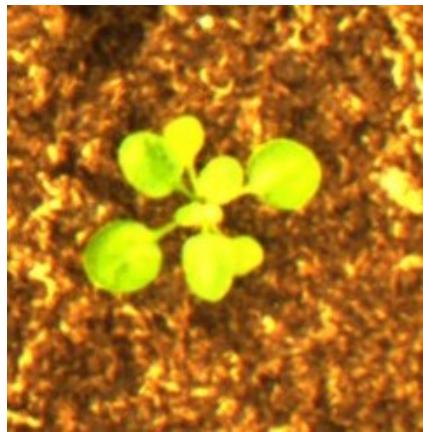
**Plant size on day x as a function of plant size on day x-1**

# MTMM (Multi-trait mixed model)

Combine correlated traits  
(same trait in two environments | longitudinal data)

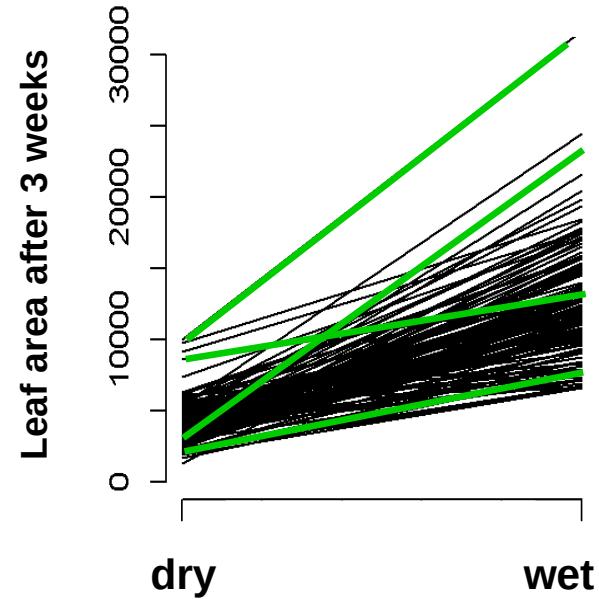
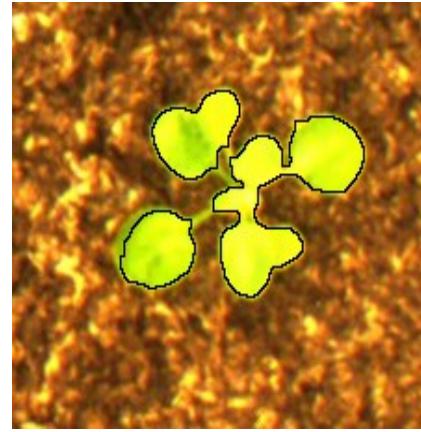
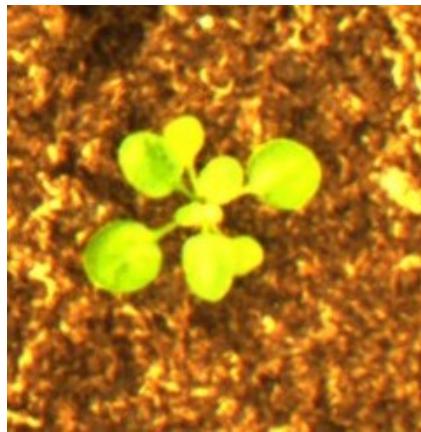


# Drought experiment



**Phenotypic differences in how drought affect seedling growth**

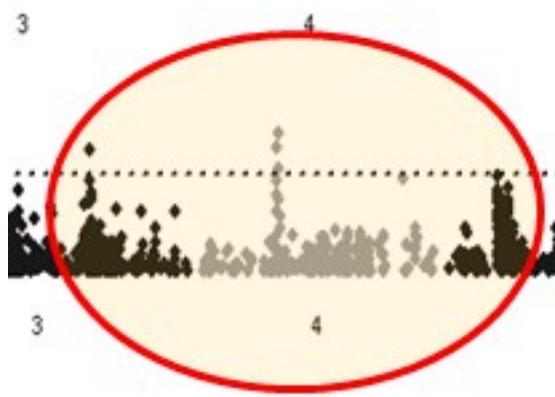
# Drought experiment



**Phenotypic differences in how drought affect seedling growth**

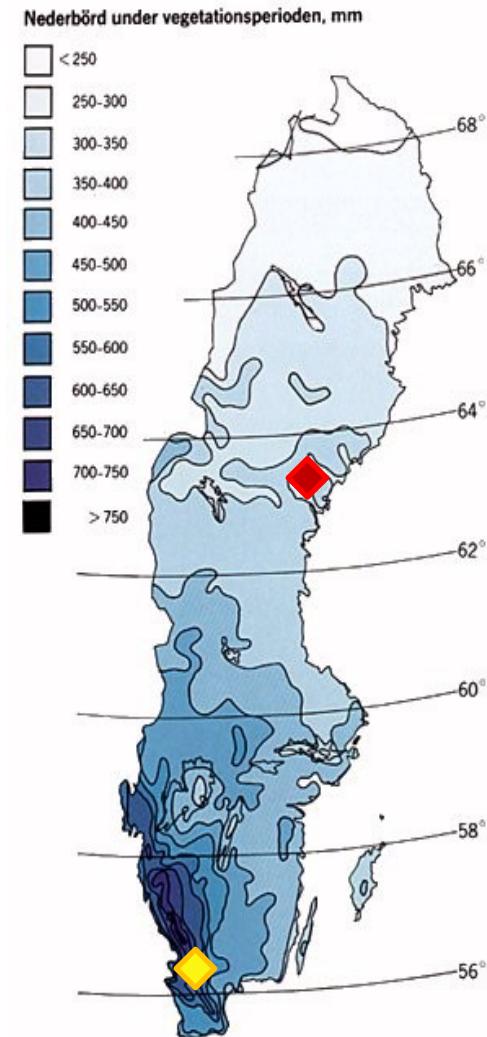
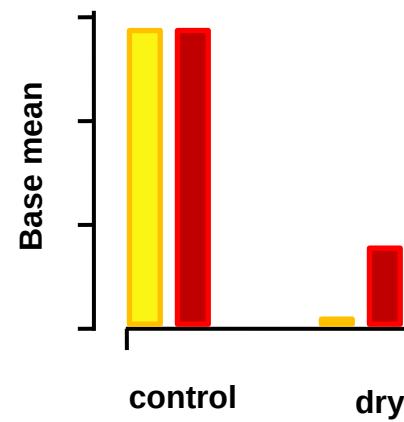
# Example : MTMM on Leaf area day 21

MTMM  
Trait specific test



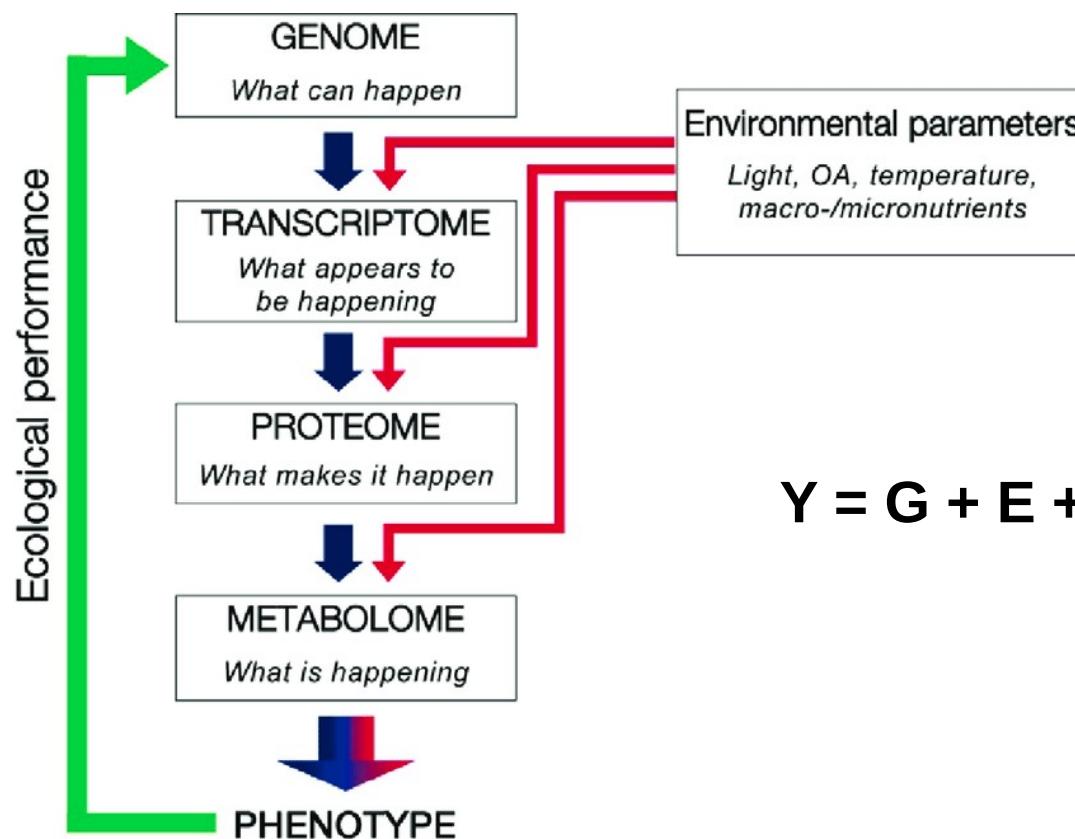
AT4G08950 EXO  
(response to  
brassinosteroids)

RNA expression  
of AT4G08950



# The Omnics cascade

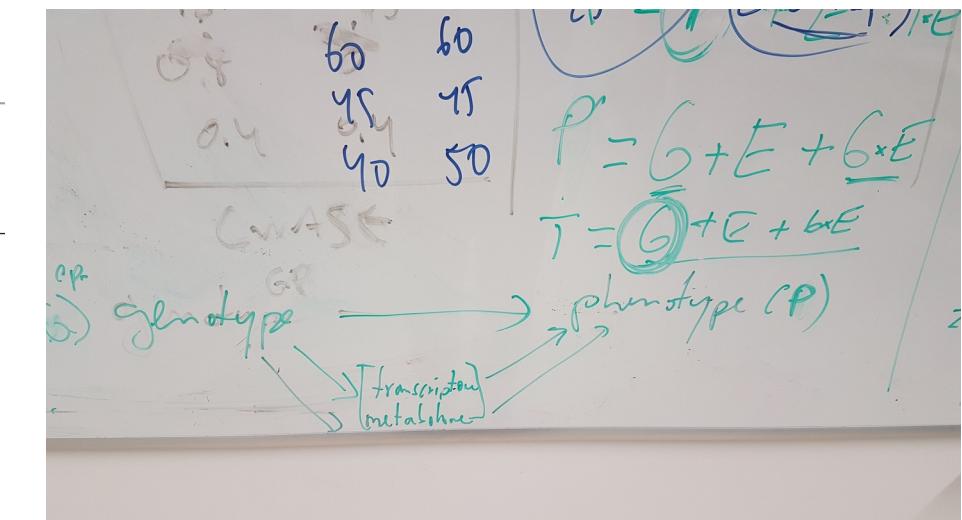
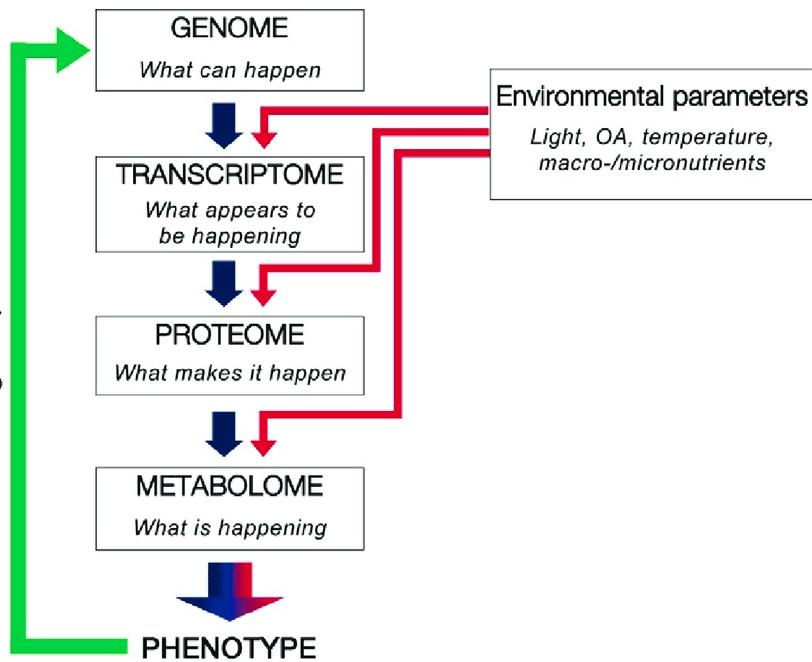
stable  
  
instable



$$Y = G + E + GxE$$

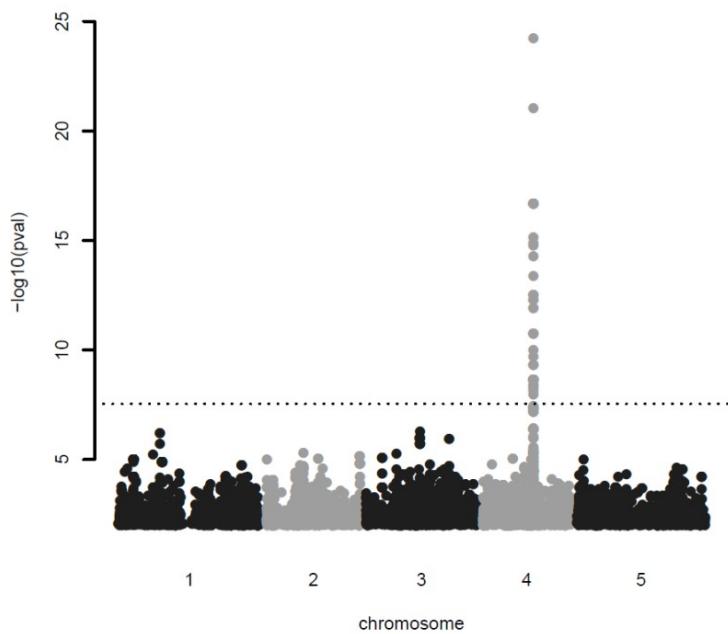
# GWAS / TWAS / eGWAS

Ecological performance



- GWAS : Genome - Phenotype**
- TWAS : Transcriptome - Phenotype**
- eGWAS : Genome - Transcriptome**

# eGWAS

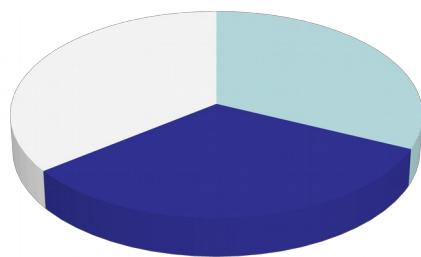


**A SNP on Chromosome 4 affects the expression of the gene AT4G19480  
How does this affect the phenotype ?**

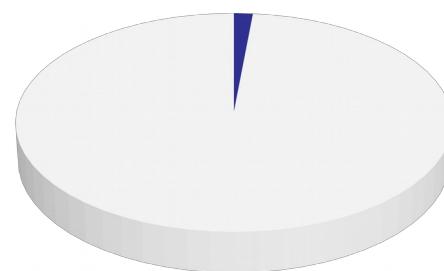
**Integration of different data layers into joint models  
(Structural Equation Models, ML)**

# Summary of MTMM on RNAseq data

Environment independent



G x E

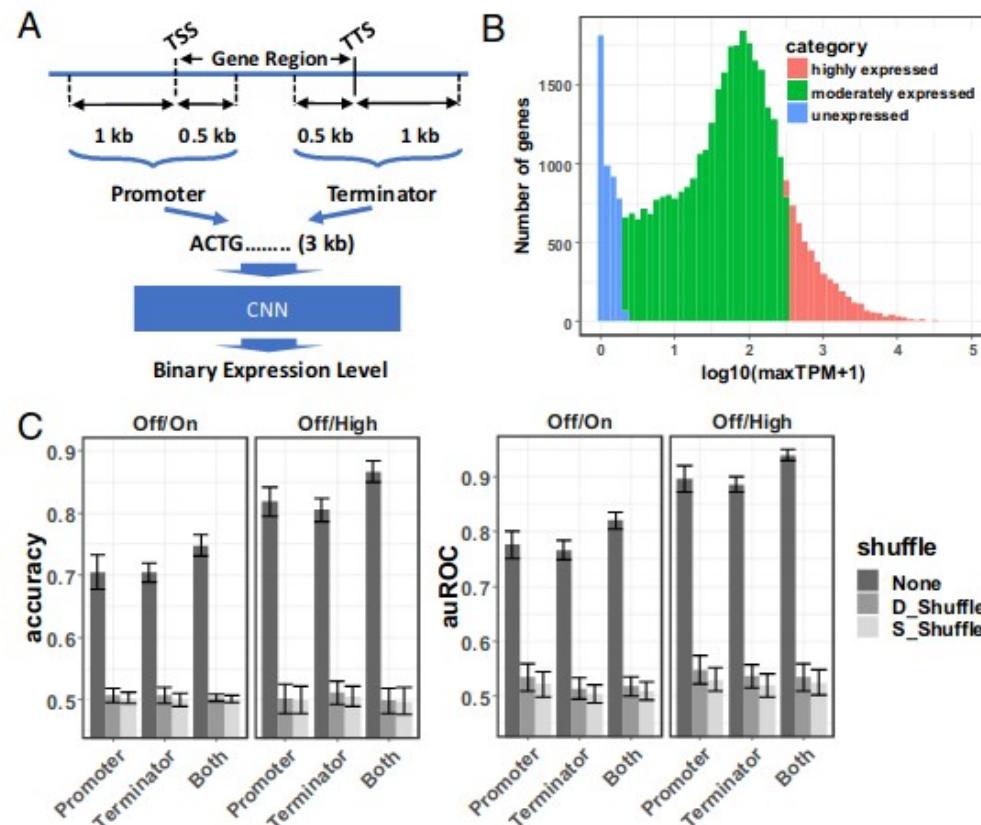


- only cis
- cis and trans
- only trans

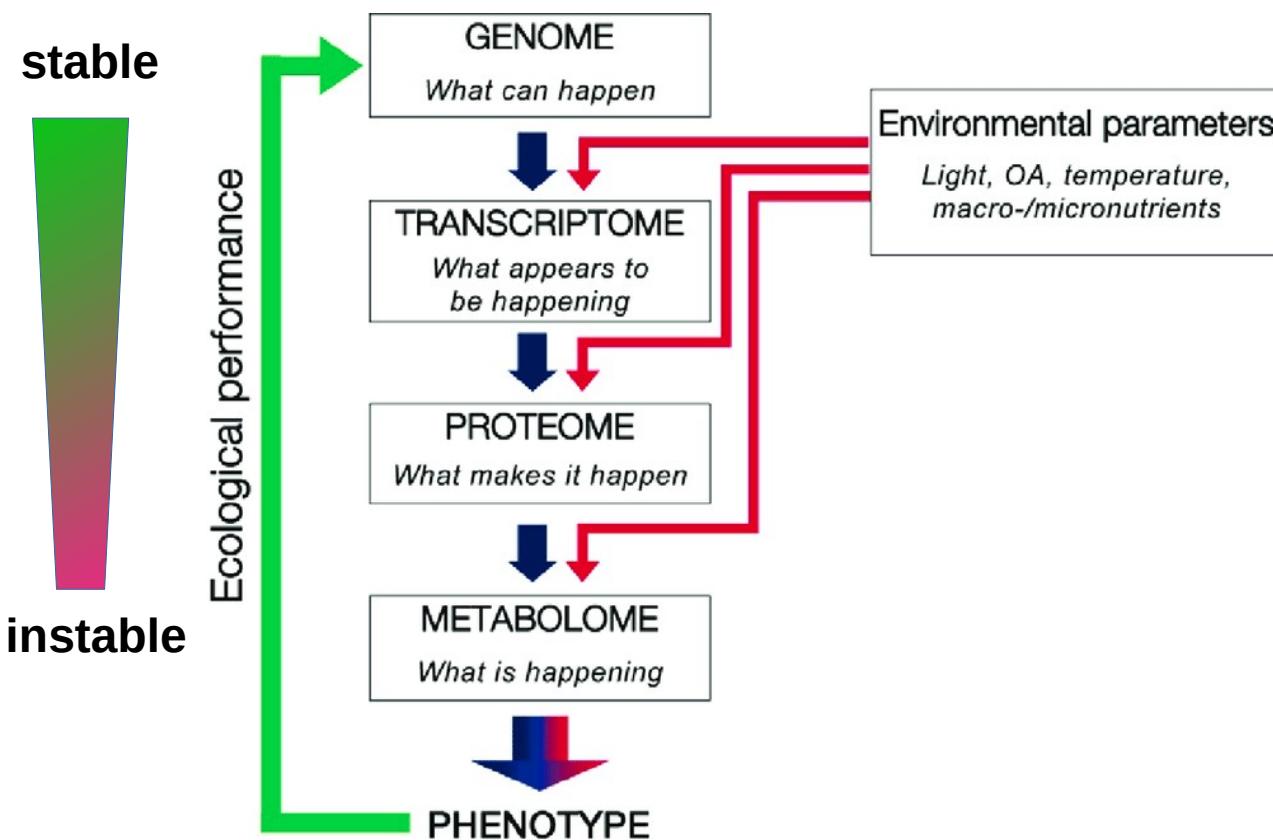
# Evolutionarily informed deep learning methods for predicting relative transcript abundance from DNA sequence

Jacob D. Washburn<sup>a,1</sup>, Maria Katherine Mejia-Guerra<sup>a</sup>, Guillaume Ramstein<sup>a</sup>, Karl A. Kremling<sup>a</sup>, Ravi Valluru<sup>a</sup>, Edward S. Buckler<sup>a,b,2</sup>, and Hai Wang<sup>c,a,1,2</sup>

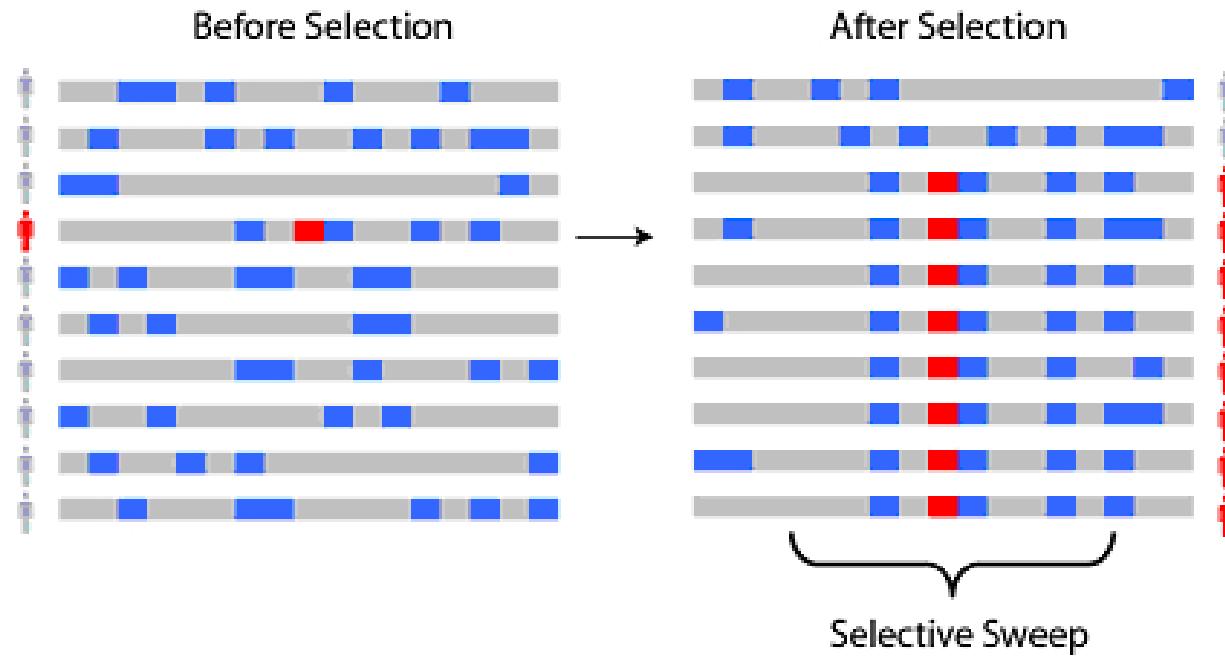
<sup>a</sup>Institute for Genomic Diversity, Cornell University, Ithaca, NY 14853; <sup>b</sup>Agricultural Research Service, United States Department of Agriculture, Ithaca, NY 14850; and <sup>c</sup>Biotechnology Research Institute, Chinese Academy of Agricultural Sciences, 100081 Beijing, China



# Selection shapes the genome



# Selective Sweeps

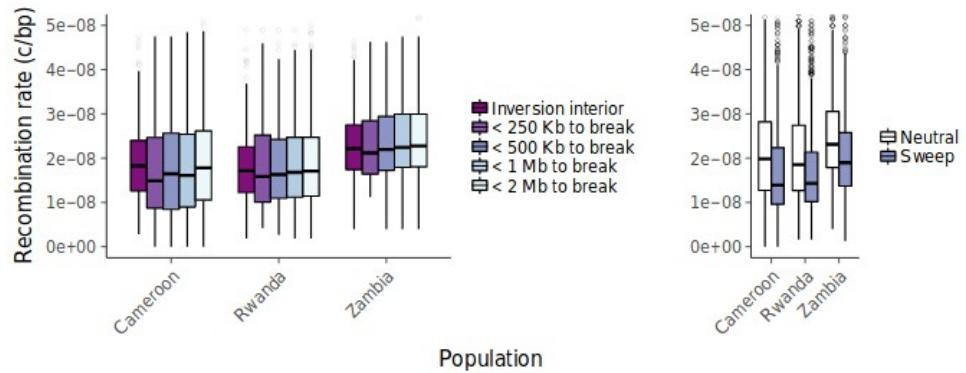
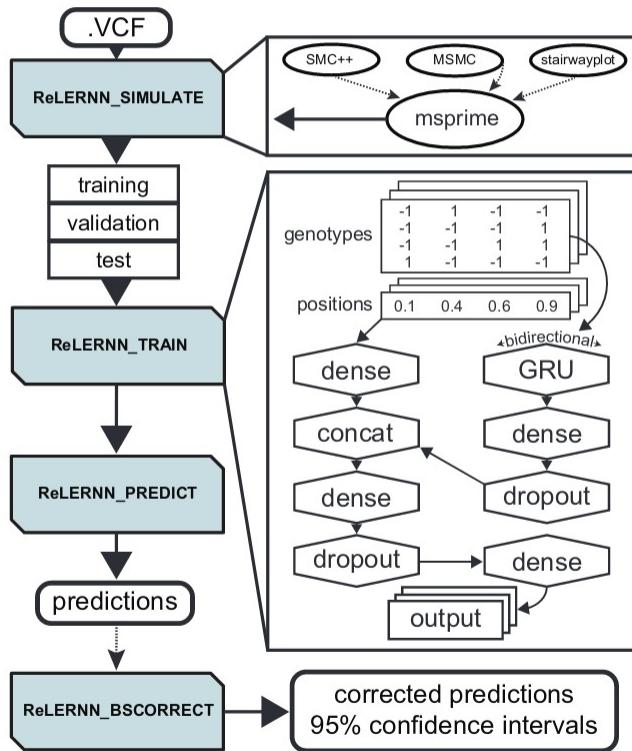


**Beneficial variants increase their frequency**  
**The same holds true for neutral variants in their proximity**

# 1 Inferring the landscape of 2 recombination using recurrent 3 neural networks

4 Jeffrey R. Adrion<sup>1,†</sup>, Jared G. Galloway<sup>1,†</sup>, Andrew D. Kern<sup>1</sup>

5 <sup>1</sup>Institute of Ecology and Evolution, University of Oregon

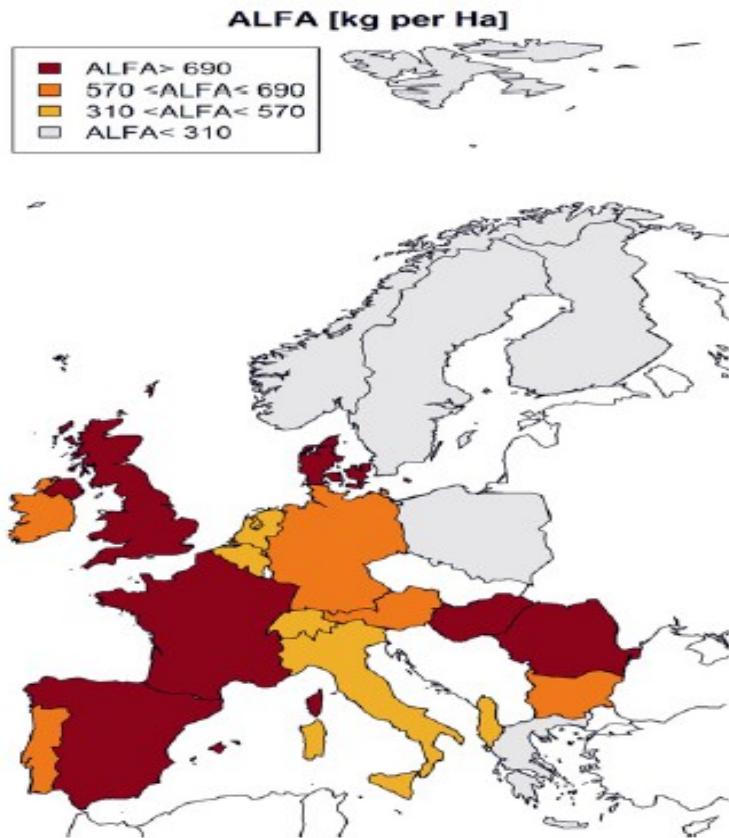


# Summary ML in Genomics

**It still feels a bit like the Wild West  
but we are trying...**

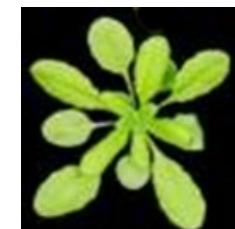
- Genome assembly
  - Imputation
- Detection of structural variants
  - GWAS
  - Genomic Prediction
- Predict gene expression
- Predict recombination rate and selective sweeps
  - joint models including different data layers

# Application to the real world



Naumann et al. Env. Research Letters 2015

Biomass

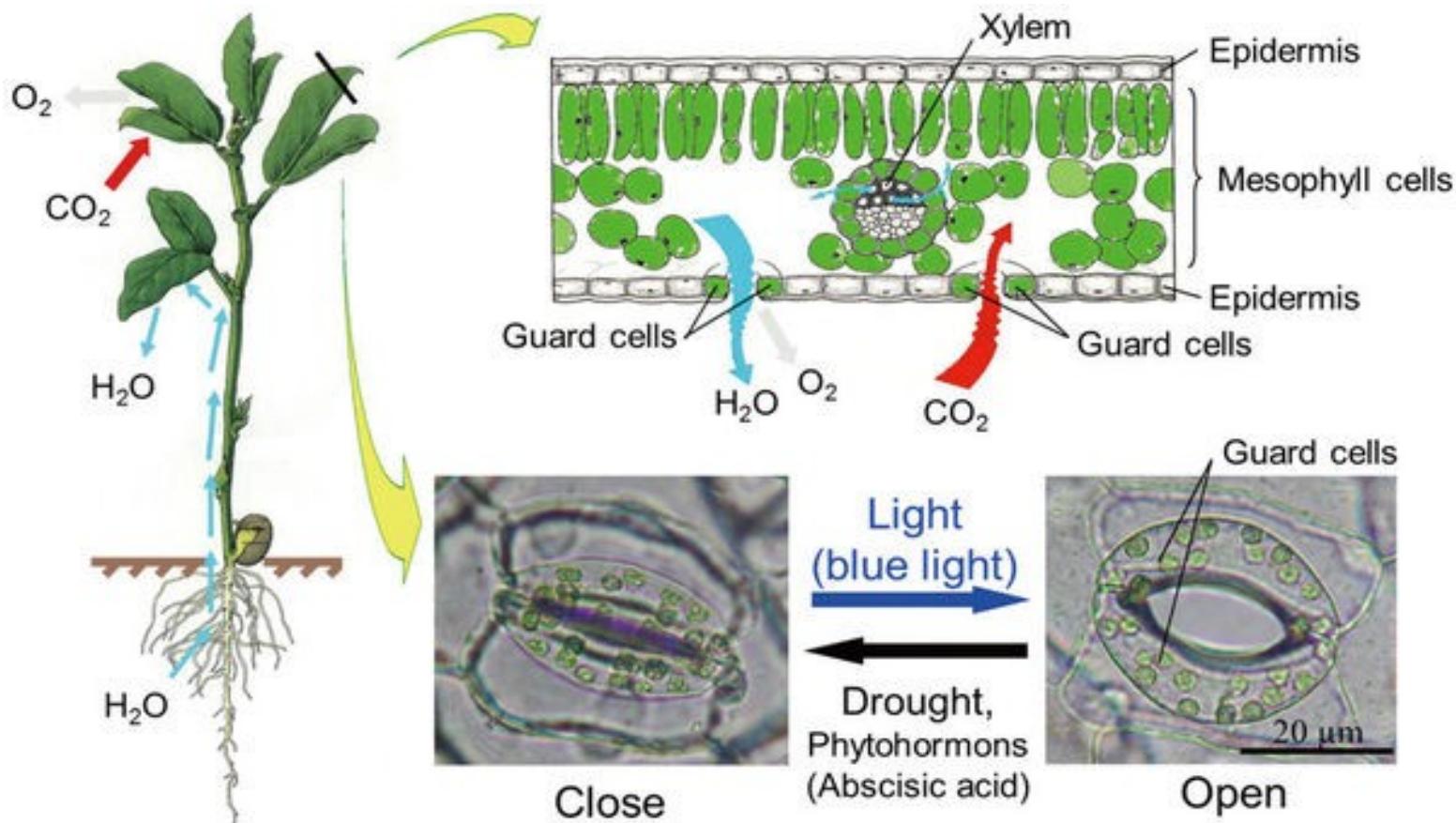


Drought tolerance



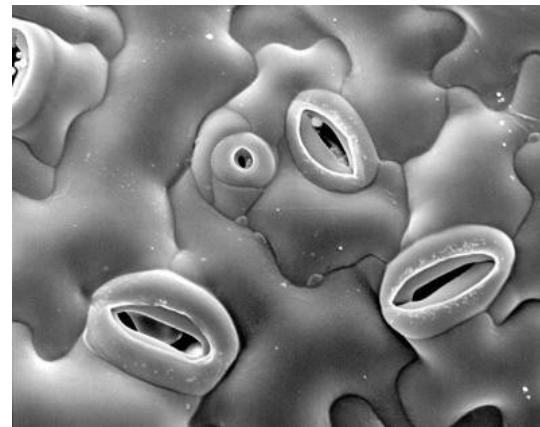
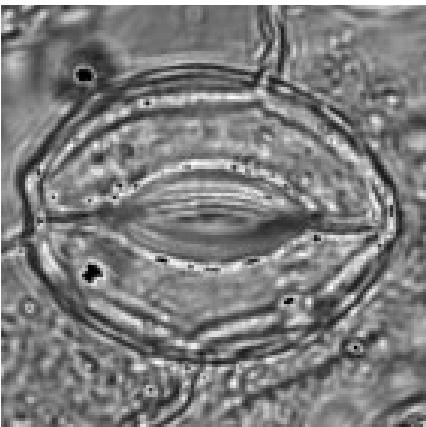
Drought = less biomass production

# Guard Cells regulate plant water status

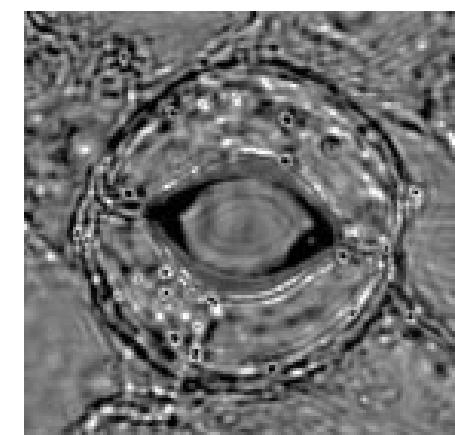


# Natural variation in guard cell signalling

Low humidity  
High CO<sub>2</sub>  
Darkness



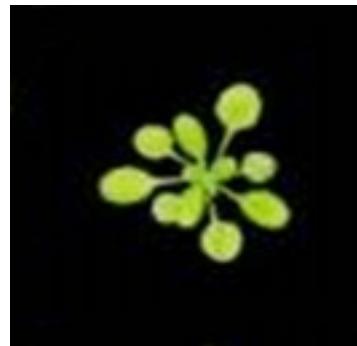
Heat  
Low CO<sub>2</sub>  
Light



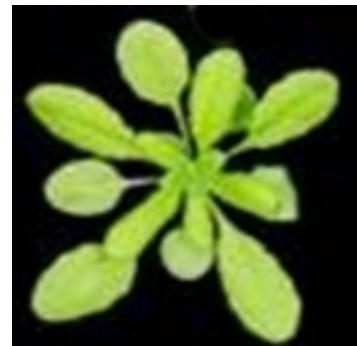
Guard cells directly sense environmental cues and balance water loss and CO<sub>2</sub> uptake

# Natural variation in guard cell signalling

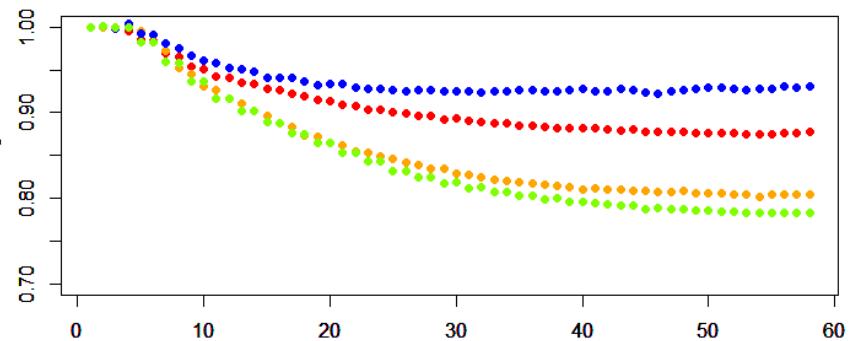
Drought tolerance



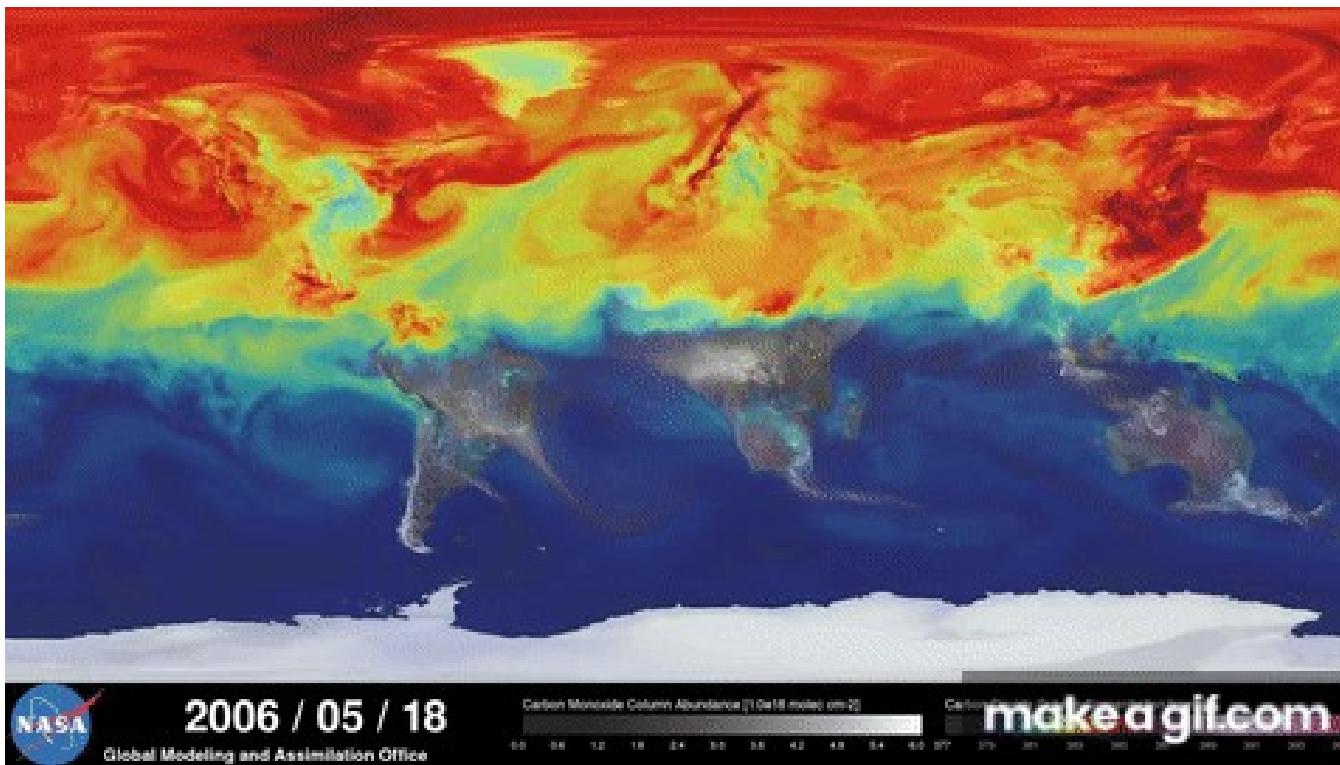
Biomass



Normalized  
transpiration rate  
( $\text{mmolH}_2\text{O/Cm}^2 \text{ leaf area}$ )



# CO<sub>2</sub> distribution in the atmosphere



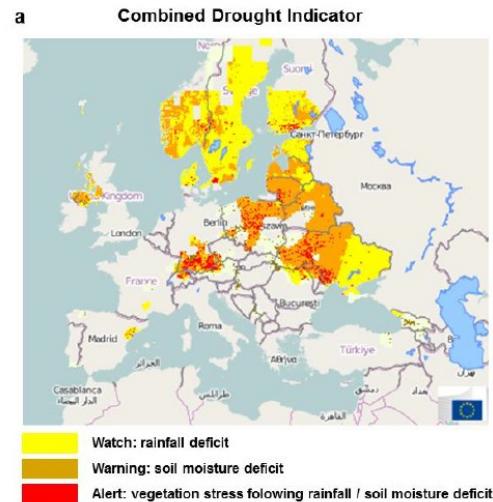
# Understanding genotype – phenotype relationships and how the environment influences these

## Statistical Genetics



$$Y = \beta_0 + \beta_1 X + u + \epsilon$$

## Drought stress



Integration of different data layers into joint models

# Acknowledgments



**Uni Würzburg**  
Jan Freudenthal  
Ammarah Anwar  
Willia Lopez

**GMI Vienna**  
Magnus Nordborg  
Ümit Seren

**VIB Ghent**  
Dirk Inze  
Frederik Coppens  
Pieter Clauw

**Uni Köln**  
Juliette DeMeaux  
Hannes Dittbrenner

