

# Emerging Marker Assisted Selection and Genomic Selection

Zhiwu Zhang





# Zhiwu Zhang Laboratory

for Statistical Genomics

[Home](#)   [People](#)   [Publication](#)   [Research](#)   [Teaching](#)   [Software](#)   [Outreach](#)   [Jobs](#)



## Five ingredients to succeed: CS-VMV

**Culture:** Trying to understand.

**Strategy:** Solve biological problems with analytical and computational challenges.

**Vision:** Genomic and phenomic stream data is stationary water for organisms.

**Mission:** You get data, we help with our analytical methods, tools, and expertise.

**Value:** Every idea makes sense.

[ZZLab.Net/share](http://ZZLab.Net/share)



# Zhiwu Zhang Laboratory

for Statistical Genomics

Home    People    Publication    Research    Teaching    Software    Outreach    Jobs



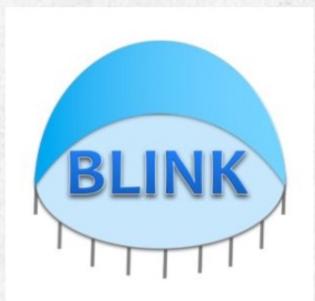
GAPIT



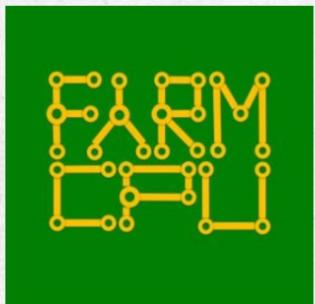
iPat



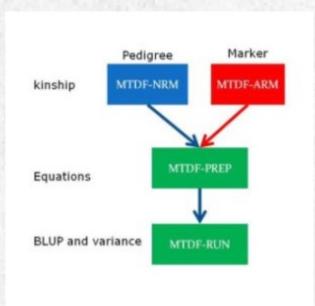
mMAP



Blink



FarmCPU



MTDFREML



GRID



Rooster



Audio4EDU



GridFree



LADDER



AI4EVER



# Genomic study

## ❖ Explanation

- Candidate gene
- Cloning
- Linkage analysis
- GWAS

} Backward  
}

Forward

GWAS  
Assisted GS

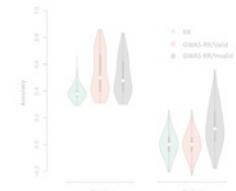
Validation

## ❖ Prediction

- MAS
- GS
- GWAS+GS
- AI

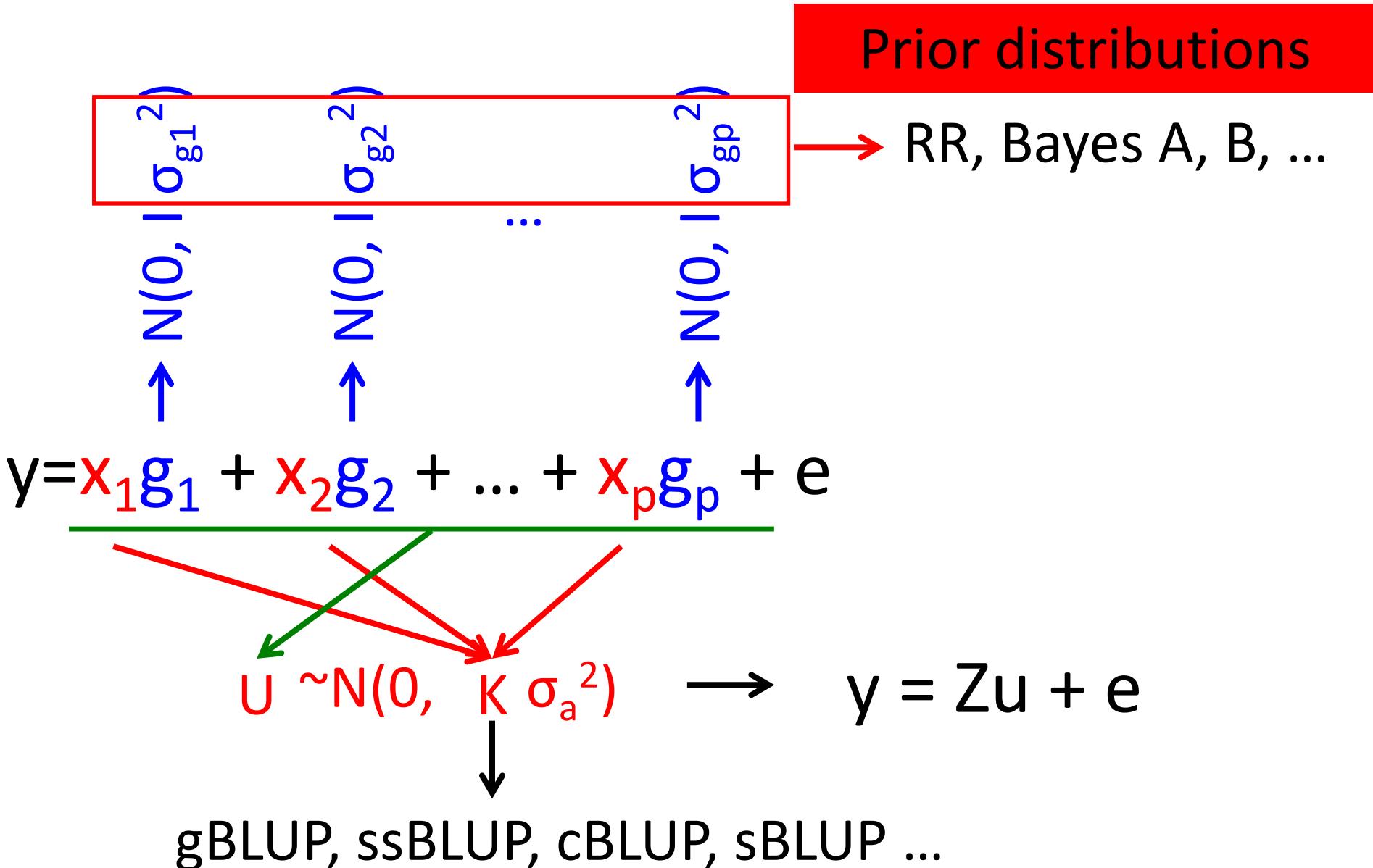
# Outline

- Overview
- Two types of problems: too bad and too good
- Right in the middle



# Genomic Prediction

Bottom up



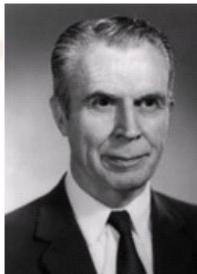
Use of Marker Based Relationships with MTDFREML  
J. Animal Sci., 2007



RFLP kinship  
in maize  
Crop Science  
1994



MAS using BLUP  
GSE, 1989



BLUP  
Biometrics, 1975



Efficient kinship  
J. Dairy Science  
2008

gBLUP



Ridge Regression, Bayes  
A, B, Cpi, ...

ssBLUP

Pedigree & Marker kinship  
single step  
GES, 2011



Super and compression  
Heredity  
2018

Genomic prediction

MAS

Prediction of total genetic value using  
genome wide dense marker maps  
Genetics, 2001



←

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

→

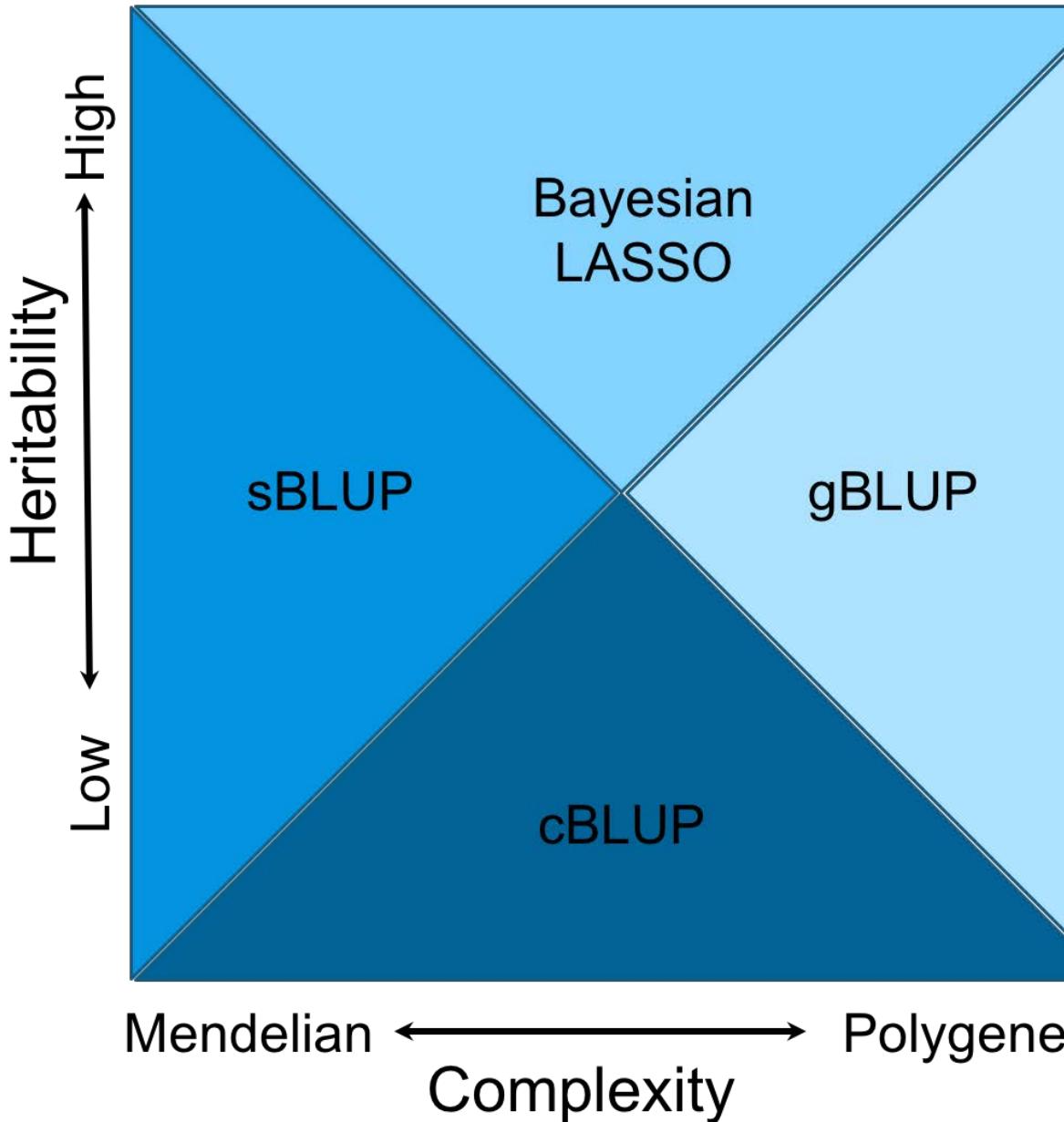
→

→

→

→

# Interaction between data and methods



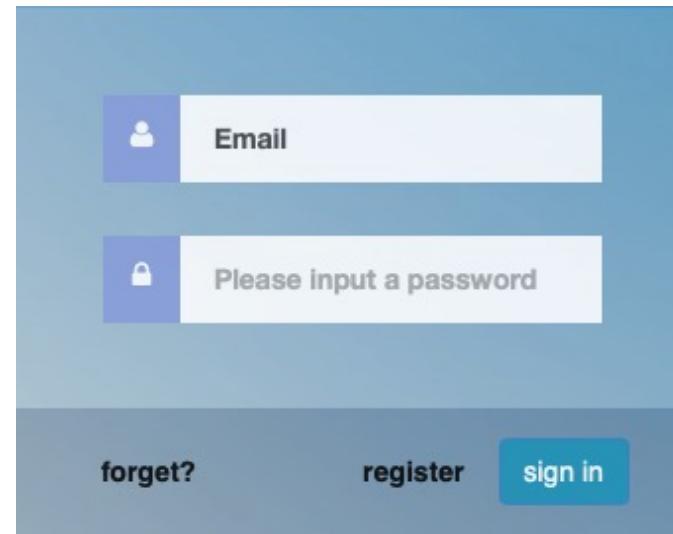
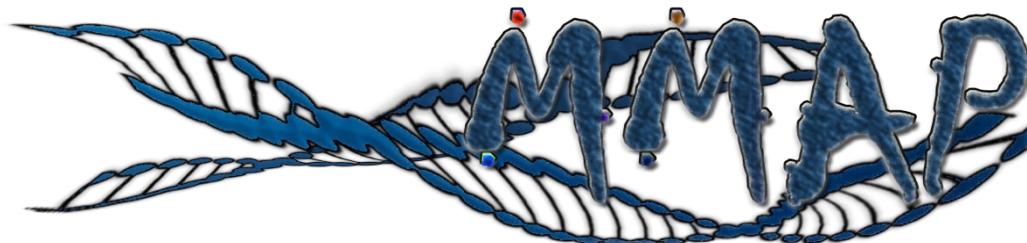
Jiabo Wang

[\*Heredity\*](#) volume 121,  
pages 648–662 (2018)

# mMap: An Online Computing Platform to Transform Genotypes to Phenotypes by Mining the Maximum Accuracy of Prediction



You Tang



Select the best method using machine learning

mMAP website: <http://zzlab.net/mMAP>



We are not luckier than this

**Type I problem: too bad than can be**

# Negative prediction accuracy

**Table 2** Testcross population parameters, parental contributions, testcross heritabilities, and RR-BLUP prediction accuracies within 14 maize biparental populations

Pedigree	Tester	<i>N</i> <sup>a</sup>	<i>N</i> <sup>b</sup> Markers	Parental contribution <sup>c</sup>		Heritability				Prediction accuracy ( <i>r</i> <sub>MG</sub> )	
				Mean	Range	Grain yield	Grain moisture	Root lodging	Stalk lodging	Grain yield	Grain moisture
S1 × S2	N1	214	115	0.47	(0.30, 0.63)	0.19* <sup>d</sup>	0.67*	-0.06	0.11	0.21*	0.23*
S1 × S3	N1	177	214	0.49	(0.19, 0.80)	0.48*	0.7*	0.01	0.19	0.32*	0.38*
S4 × S5	N1	177	231	0.35	(0.12, 0.88)	0.45*	0.75*	-0.03	0.15	-0.30*	0.33*
S4 × S6	N2	185	239	0.48	(0.26, 0.74)	0.59*	0.32*	0.21*	-0.04	0.16*	0.00
S7 × S1	N2	151	203	0.48	(0.17, 0.78)	0.54*	0.86*	0.22*	0.31*	0.14	0.25*
S8 × S9	N2	292	197	0.33	(0.10, 0.90)	0.44*	0.73*	-0.01	-0.02	0.36*	0.39*
S10 × S11	N3	184	232	0.47	(0.14, 0.86)	0.73*	0.63*	-0.08	0.04	0.30*	0.26*
N4 × N5	S10	141	249	0.34	(0.13, 0.52)	0.25*	0.83*	0.12	0.07	0.18	-0.14
N4 × N5	S12	77	249	0.34	(0.16, 0.52)	0.5*	0.56*	0.35*	0.06	0.33*	-0.33*
N6 × N4	S7	171	203	0.43	(0.15, 0.85)	0.39*	0.77*	0.02	-0.05	-0.08	-0.14
N6 × N4	S12	71	203	0.44	(0.14, 0.86)	0.35*	0.28*	0.07	-0.11	-0.12	-0.11
N7 × N3	S4	109	249	0.44	(0.38, 0.62)	0.32*	0.48*	0.11	0.09	0.07	-0.42*
N8 × N6	S13	211	338	0.33	(0.17, 0.48)	0.43*	0.83*	0.002	0.11	0.21*	-0.08
N7 × N9	S4	114	243	0.33	(0.22, 0.78)	0.37*	0.72*	0.11	-0.04	-0.10	-0.24*

\* Significant at *P* = 0.05

<sup>a</sup> Number of individuals in the biparental population

<sup>b</sup> Number of polymorphic SNPs in the biparental population

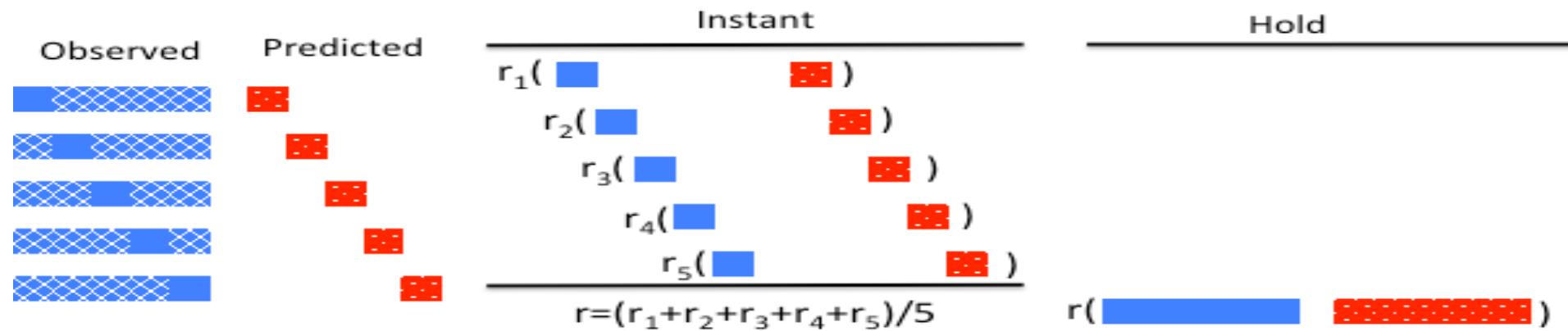
<sup>c</sup> Parental contribution of the less-represented parent

Theor Appl Genet. 2013 Jan;126(1):13-22

Genomewide predictions from maize single-cross data.

Massman JM1, Gordillo A, Lorenzana RE, Bernardo R.

# Two ways of calculating correlation

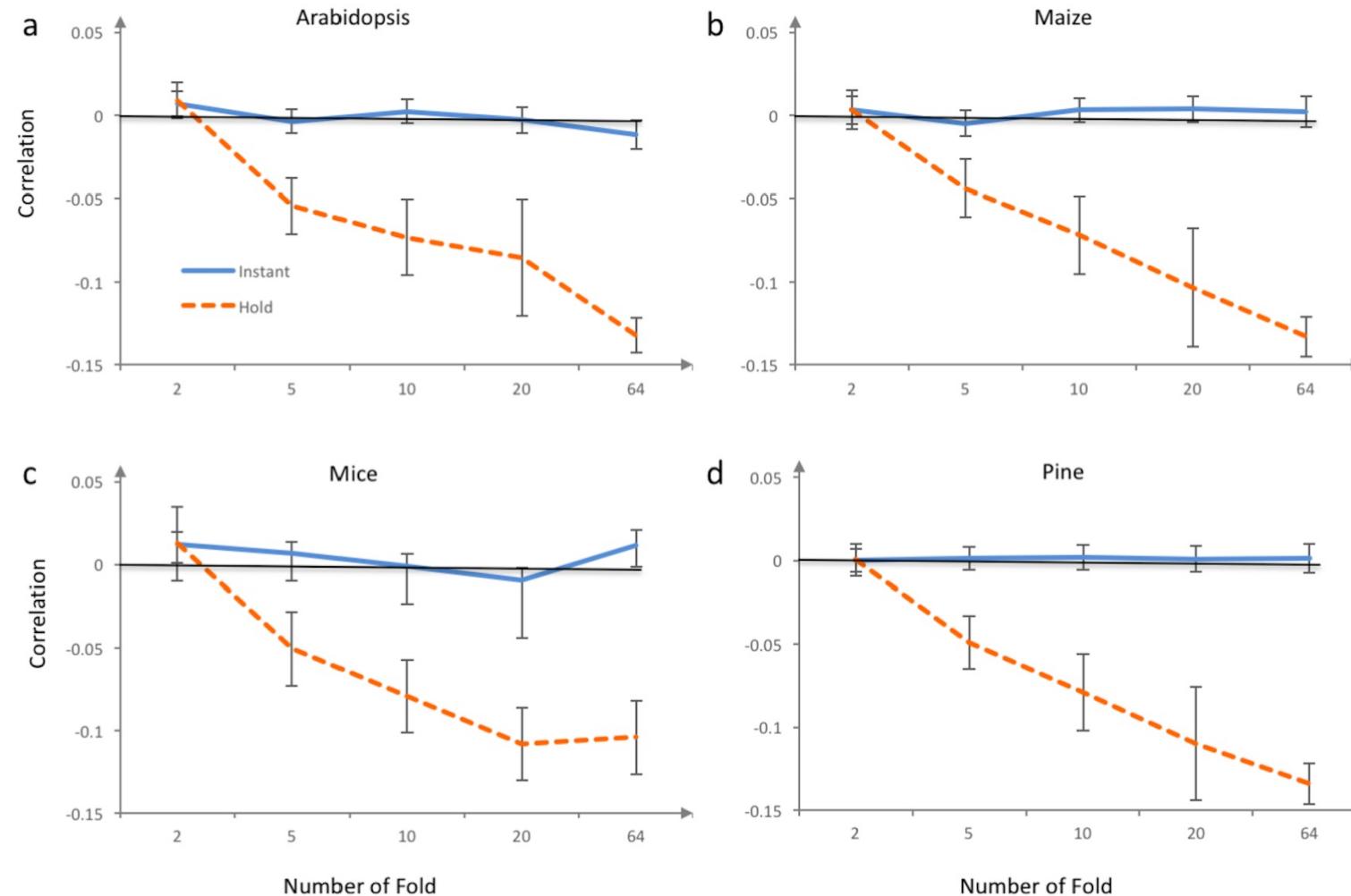


*Yao Zhou and et al., Briefings in Bioinformatics, Volume 18, Issue 5, September 2017, Pages 744–753,  
<https://doi.org/10.1093/bib/bbw064>*

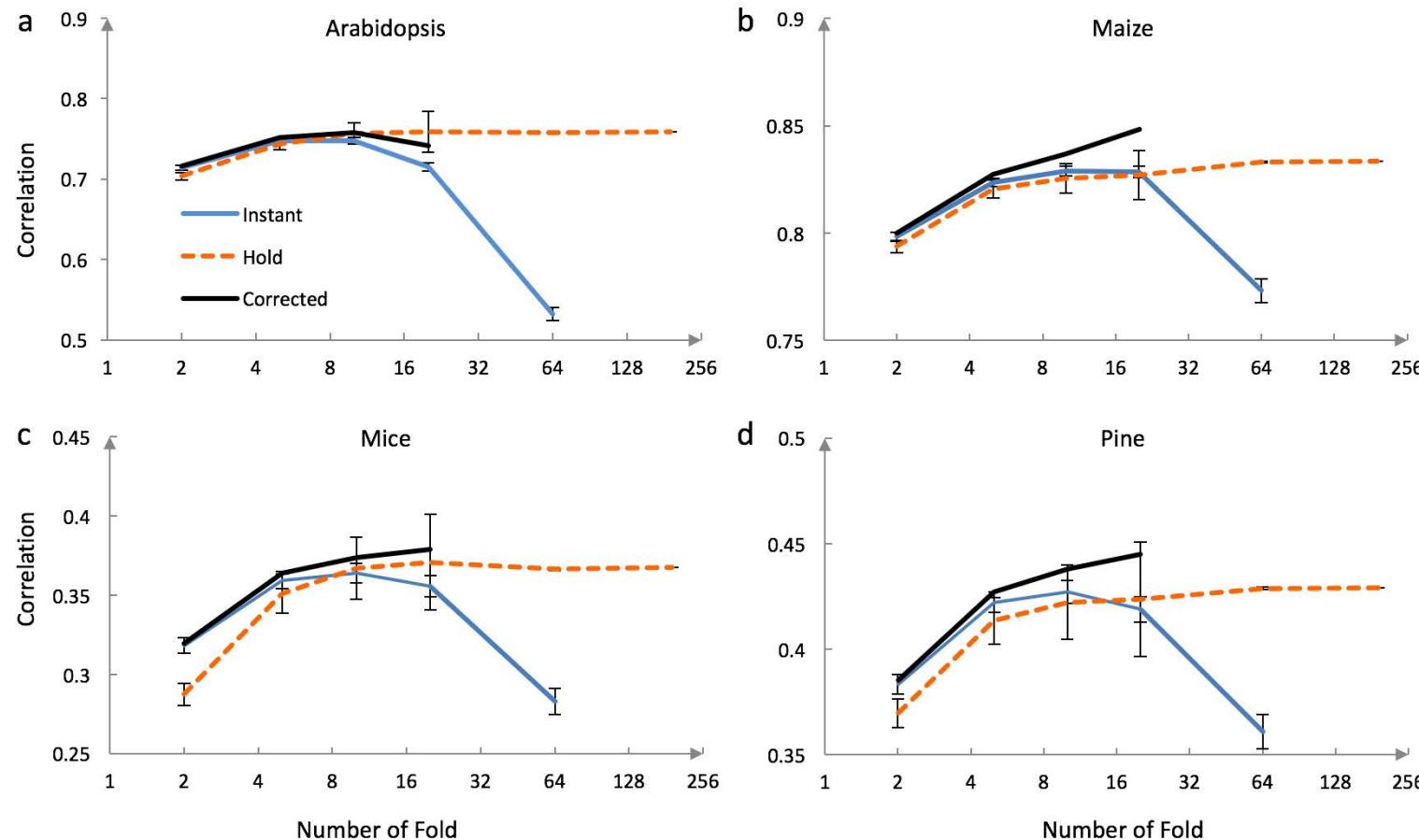
# Underestimates of hold accuracy



Yao Zhou



# Underestimates and correction of instant accuracy



$$\hat{\rho} = r \left[ 1 + \frac{(1 - r^2)}{2(n - 4)} \right]$$

Y Zhou, MI Vales, A  
Wang, Z Zhang  
Briefings in  
bioinformatics 18 (5),  
744-753



Crop Breeding & Genetics

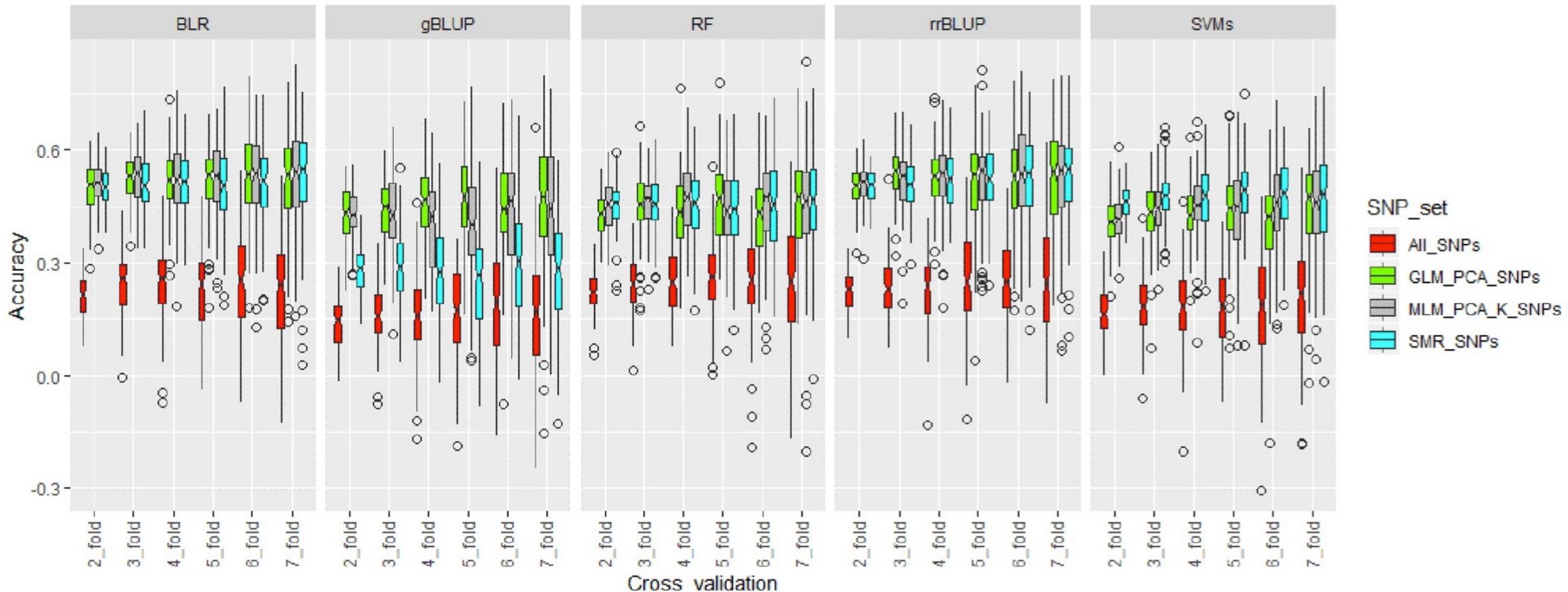
## Genomewide Selection when Major Genes Are Known

Rex Bernardo

First published: 01 January 2014 | <https://doi.org/10.2135/cropsci2013.05.0315> |

Citations: 161

**Type II problem: too good to be true**



Genomic selection accuracy **increased to almost 2-fold** at each level of cross validation when the GWAS-derived SNPs were incorporated into the genomic selection model.

**Table 5.** Comparison of genomic selection (GS) models in 13 phenotypic traits collected in the SolCAP potato diversity panel. Mean and standard deviation of Pearson's correlation obtained by 10-fold cross validation in 10 replicates. SNP weights for WGBLUP were obtained from  $-\log_{10} p$ -values of different GWASpoly models.

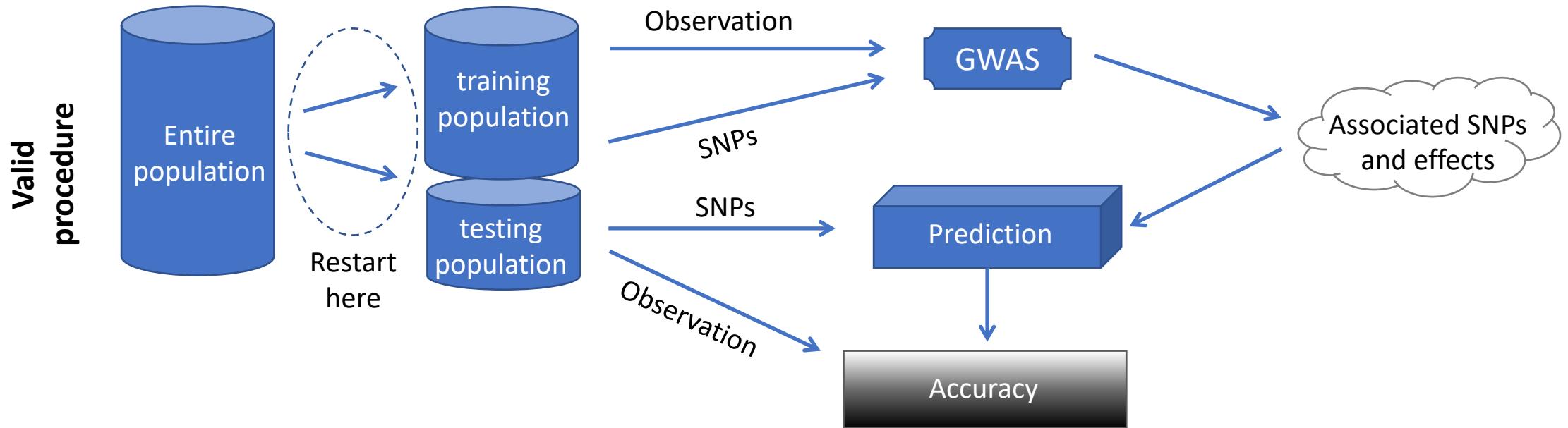
Trait	RRBLUP	GBLUP	WGBLUP							
			1-d-a	1-d-r	2-d-a	2-d-r	General	d-Gen	d-Add	Additive
Chip color	0.723 ( $\pm 0.014$ )	0.721 ( $\pm 0.015$ )	0.826 ( $\pm 0.009$ )	0.798 ( $\pm 0.011$ )	0.859 ( $\pm 0.007$ )	0.850 ( $\pm 0.013$ )	0.867 ( $\pm 0.008$ )	0.849 ( $\pm 0.009$ )	0.855 ( $\pm 0.007$ )	0.896 ( $\pm 0.007$ )
$\log_{10}$ fructose	0.682 ( $\pm 0.024$ )	0.676 ( $\pm 0.025$ )	0.819 ( $\pm 0.014$ )	0.785 ( $\pm 0.017$ )	0.845 ( $\pm 0.007$ )	0.833 ( $\pm 0.011$ )	0.868 ( $\pm 0.011$ )	0.839 ( $\pm 0.015$ )	0.855 ( $\pm 0.003$ )	0.895 ( $\pm 0.008$ )
$\log_{10}$ glucose	0.678 ( $\pm 0.017$ )	0.668 ( $\pm 0.030$ )	0.796 ( $\pm 0.009$ )	0.809 ( $\pm 0.016$ )	0.855 ( $\pm 0.009$ )	0.849 ( $\pm 0.009$ )	0.875 ( $\pm 0.009$ )	0.844 ( $\pm 0.011$ )	0.848 ( $\pm 0.013$ )	0.91 ( $\pm 0.007$ )
Malic acid	0.602 ( $\pm 0.016$ )	0.598 ( $\pm 0.027$ )	0.751 ( $\pm 0.021$ )	0.745 ( $\pm 0.022$ )	0.802 ( $\pm 0.021$ )	0.801 ( $\pm 0.016$ )	0.838 ( $\pm 0.011$ )	0.808 ( $\pm 0.016$ )	0.826 ( $\pm 0.009$ )	0.876 ( $\pm 0.007$ )
Sucrose	0.539 ( $\pm 0.024$ )	0.519 ( $\pm 0.034$ )	0.676 ( $\pm 0.011$ )	0.675 ( $\pm 0.022$ )	0.702 ( $\pm 0.019$ )	0.716 ( $\pm 0.015$ )	0.725 ( $\pm 0.023$ )	0.722 ( $\pm 0.011$ )	0.739 ( $\pm 0.019$ )	0.806 ( $\pm 0.011$ )
Total yield	0.132 ( $\pm 0.023$ )	0.117 ( $\pm 0.041$ )	0.401 ( $\pm 0.026$ )	0.413 ( $\pm 0.030$ )	0.418 ( $\pm 0.031$ )	0.428 ( $\pm 0.017$ )	0.470 ( $\pm 0.029$ )	0.492 ( $\pm 0.030$ )	0.504 ( $\pm 0.030$ )	0.584 ( $\pm 0.028$ )
Tuber eye depth	0.495 ( $\pm 0.026$ )	0.478 ( $\pm 0.019$ )	0.605 ( $\pm 0.029$ )	0.655 ( $\pm 0.016$ )	0.693 ( $\pm 0.025$ )	0.717 ( $\pm 0.014$ )	0.740 ( $\pm 0.020$ )	0.693 ( $\pm 0.020$ )	0.736 ( $\pm 0.018$ )	0.812 ( $\pm 0.007$ )
Tuber length	0.826 ( $\pm 0.012$ )	0.821 ( $\pm 0.014$ )	0.891 ( $\pm 0.006$ )	0.884 ( $\pm 0.009$ )	0.899 ( $\pm 0.006$ )	0.889 ( $\pm 0.012$ )	0.904 ( $\pm 0.008$ )	0.908 ( $\pm 0.008$ )	0.912 ( $\pm 0.005$ )	0.928 ( $\pm 0.009$ )
Tuber shape	0.775 ( $\pm 0.018$ )	0.780 ( $\pm 0.017$ )	0.865 ( $\pm 0.010$ )	0.853 ( $\pm 0.013$ )	0.886 ( $\pm 0.008$ )	0.863 ( $\pm 0.005$ )	0.896 ( $\pm 0.010$ )	0.89 ( $\pm 0.008$ )	0.891 ( $\pm 0.009$ )	0.922 ( $\pm 0.006$ )
Tuber size	0.501 ( $\pm 0.024$ )	0.499 ( $\pm 0.027$ )	0.641 ( $\pm 0.019$ )	0.650 ( $\pm 0.020$ )	0.679 ( $\pm 0.020$ )	0.663 ( $\pm 0.022$ )	0.666 ( $\pm 0.024$ )	0.661 ( $\pm 0.022$ )	0.679 ( $\pm 0.019$ )	0.742 ( $\pm 0.021$ )
Tuber width	0.635 ( $\pm 0.023$ )	0.638 ( $\pm 0.021$ )	0.752 ( $\pm 0.020$ )	0.749 ( $\pm 0.021$ )	0.782 ( $\pm 0.016$ )	0.772 ( $\pm 0.018$ )	0.805 ( $\pm 0.012$ )	0.789 ( $\pm 0.015$ )	0.803 ( $\pm 0.013$ )	0.847 ( $\pm 0.017$ )
Vine maturity 95 days	0.288 ( $\pm 0.035$ )	0.286 ( $\pm 0.042$ )	0.550 ( $\pm 0.028$ )	0.538 ( $\pm 0.020$ )	0.603 ( $\pm 0.022$ )	0.589 ( $\pm 0.028$ )	0.668 ( $\pm 0.022$ )	0.632 ( $\pm 0.019$ )	0.65 ( $\pm 0.025$ )	0.746 ( $\pm 0.017$ )
Vine maturity 120 days	0.321 ( $\pm 0.047$ )	0.323 ( $\pm 0.024$ )	0.495 ( $\pm 0.026$ )	0.569 ( $\pm 0.021$ )	0.636 ( $\pm 0.021$ )	0.633 ( $\pm 0.013$ )	0.669 ( $\pm 0.025$ )	0.616 ( $\pm 0.023$ )	0.666 ( $\pm 0.026$ )	0.755 ( $\pm 0.019$ )

RRBLUP, best linear unbiased prediction using ridge-regression; GBLUP, genomic best linear unbiased prediction using VanRaden G matrix; WGBLUP, weighted GBLUP; 1-d-a and 1-d-r, simplex dominant models; 2-d-a and 2-d-r, duplex dominant models; d-gen, diploidized general; d-add, diploidized additive.

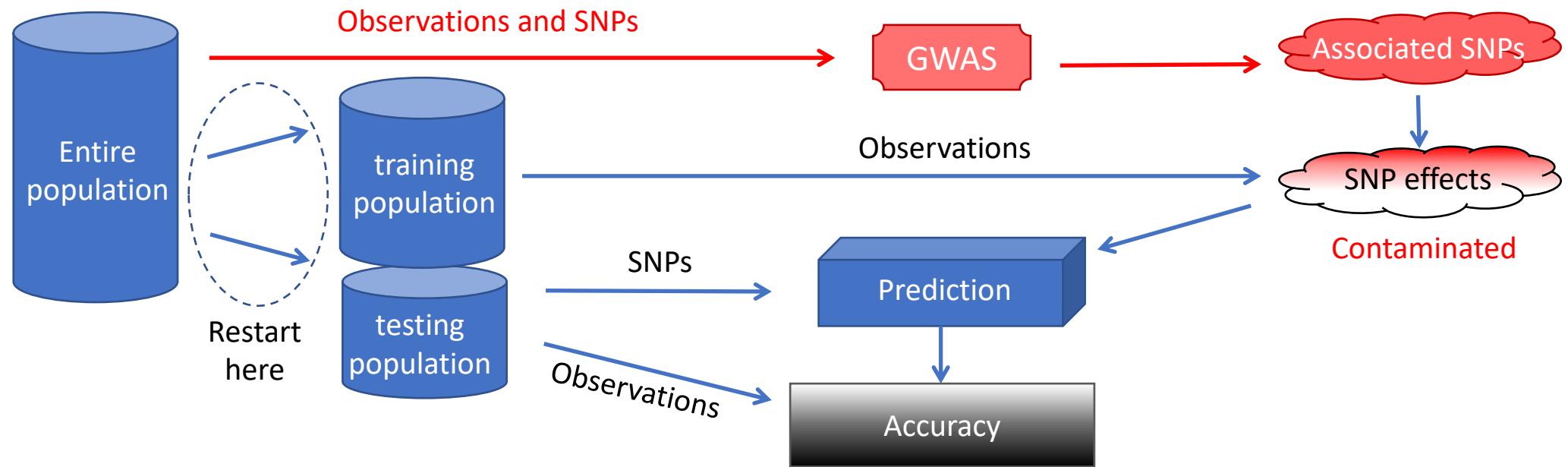
Medina, C. A., Kaur, H., Ray, I. & Yu, L. X. Strategies to increase prediction accuracy in genomic selection of complex traits in alfalfa (*Medicago sativa* l.). Cells vol. 10 Preprint at <https://doi.org/10.3390/cells10123372> (2021).

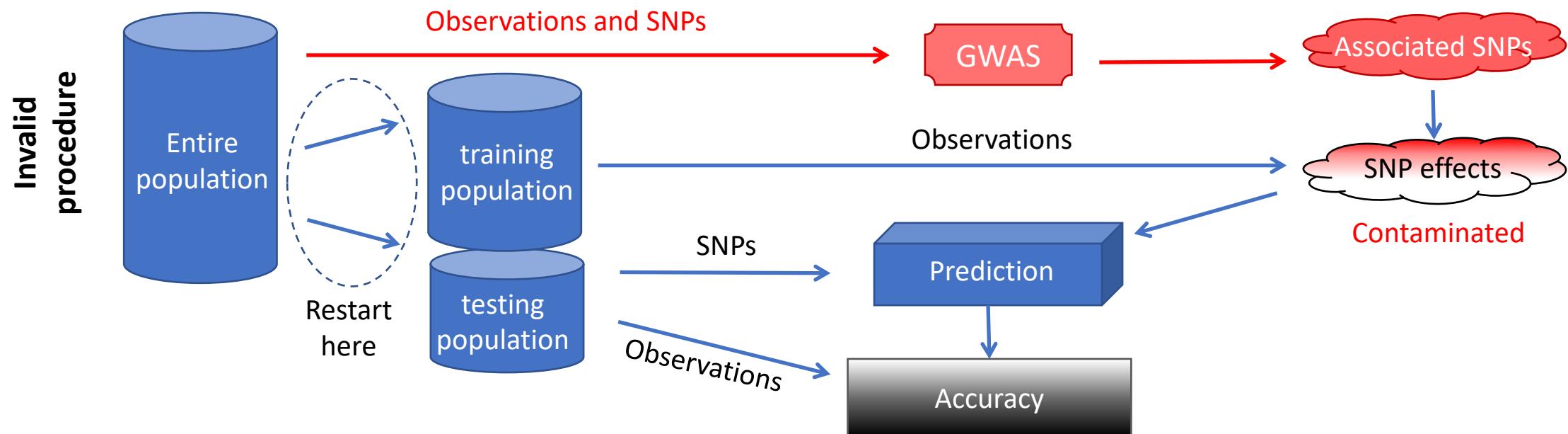
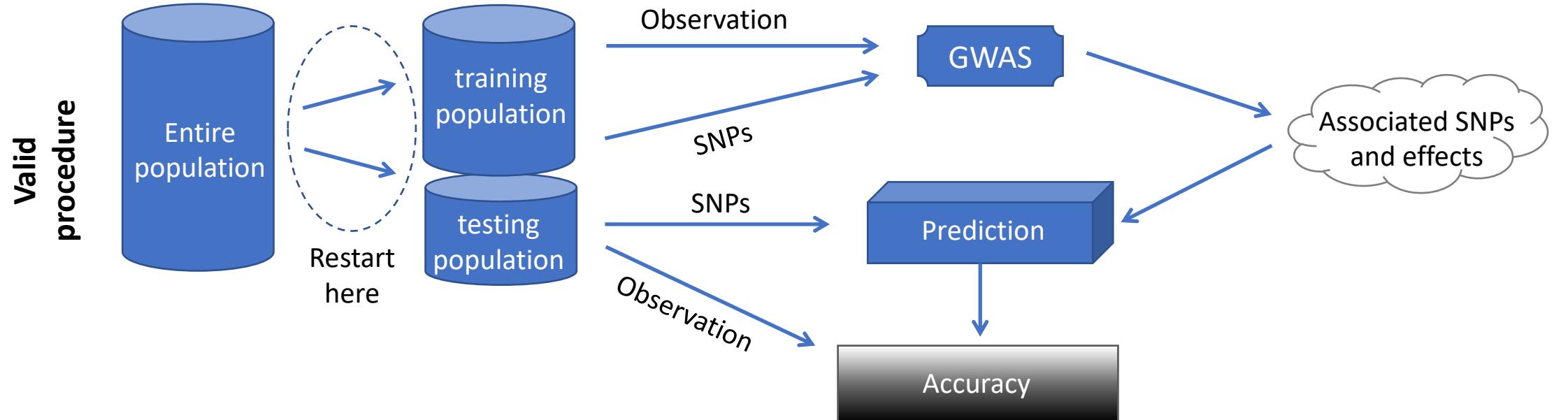
Higher prediction accuracies in all 13 agronomic traits of potato were obtained using the WGBLUP model with SNP  $-\log_{10} p$ -values derived from the additive GWASpoly model. Traits of glucose, tuber length, or tuber shape showed accuracies higher than 0.9. It is important to point out that traits of tuber length or tuber shape had high accuracies (0.82 and 0.78 respectively using RRBLUP and GBLUP models) and the use of the WGBLUP model increased the prediction accuracy up to 0.93. Total yield had low prediction accuracies with RRBLUP or GBLUP models (0.132 and 0.117, respectively), and the use of the WGBLUP model increased prediction accuracy by almost five times (Table 5). These results agree with our previous results in alfalfa (Figure 2).

# Valid procedure



# Invalid procedure





# Invalid procedure create artifact



**Right in the middle**

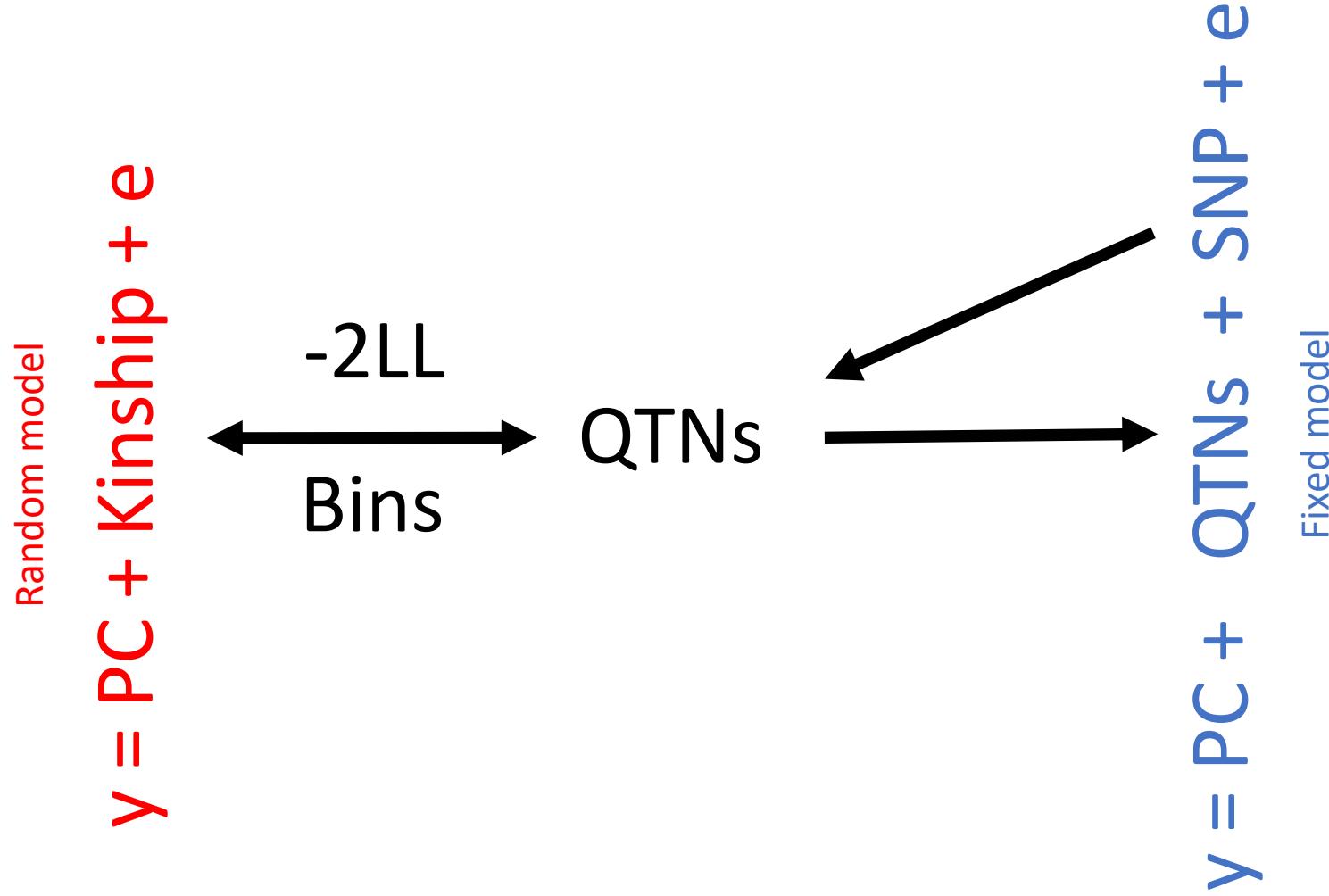
# Evaluation of RR-BLUP Genomic Selection Models that Incorporate Peak Genome-Wide Association Study Signals in Maize and Sorghum

Brian Rice and Alexander E. Lipka\*

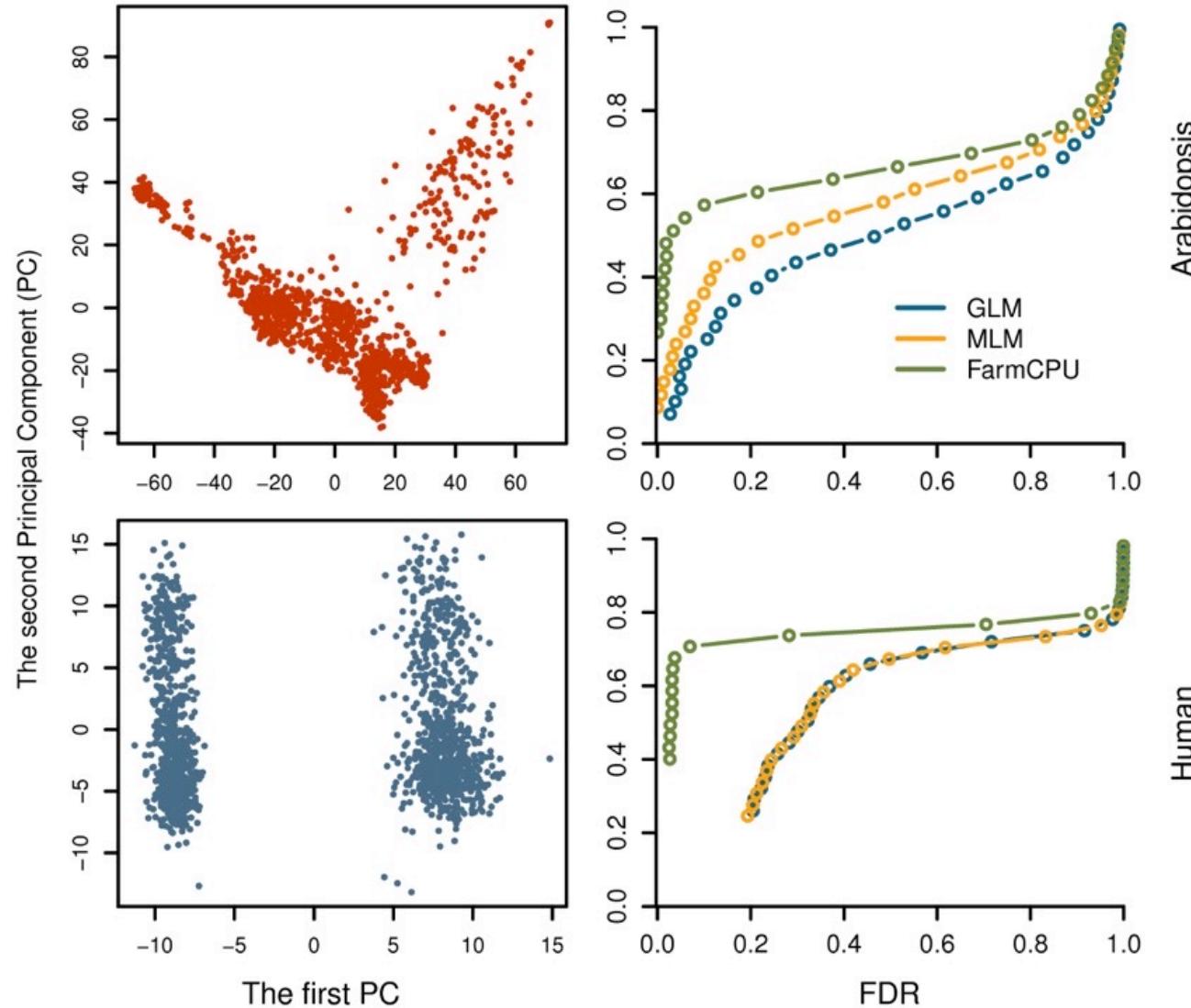
was quantified. Of the 216 genetic architectures that we simulated, we identified 60 where the addition of fixed-effect covariates boosted prediction accuracy. However, for the majority of the simulated data, no increase or a decrease in prediction accuracy was observed. We also noted several instances where the



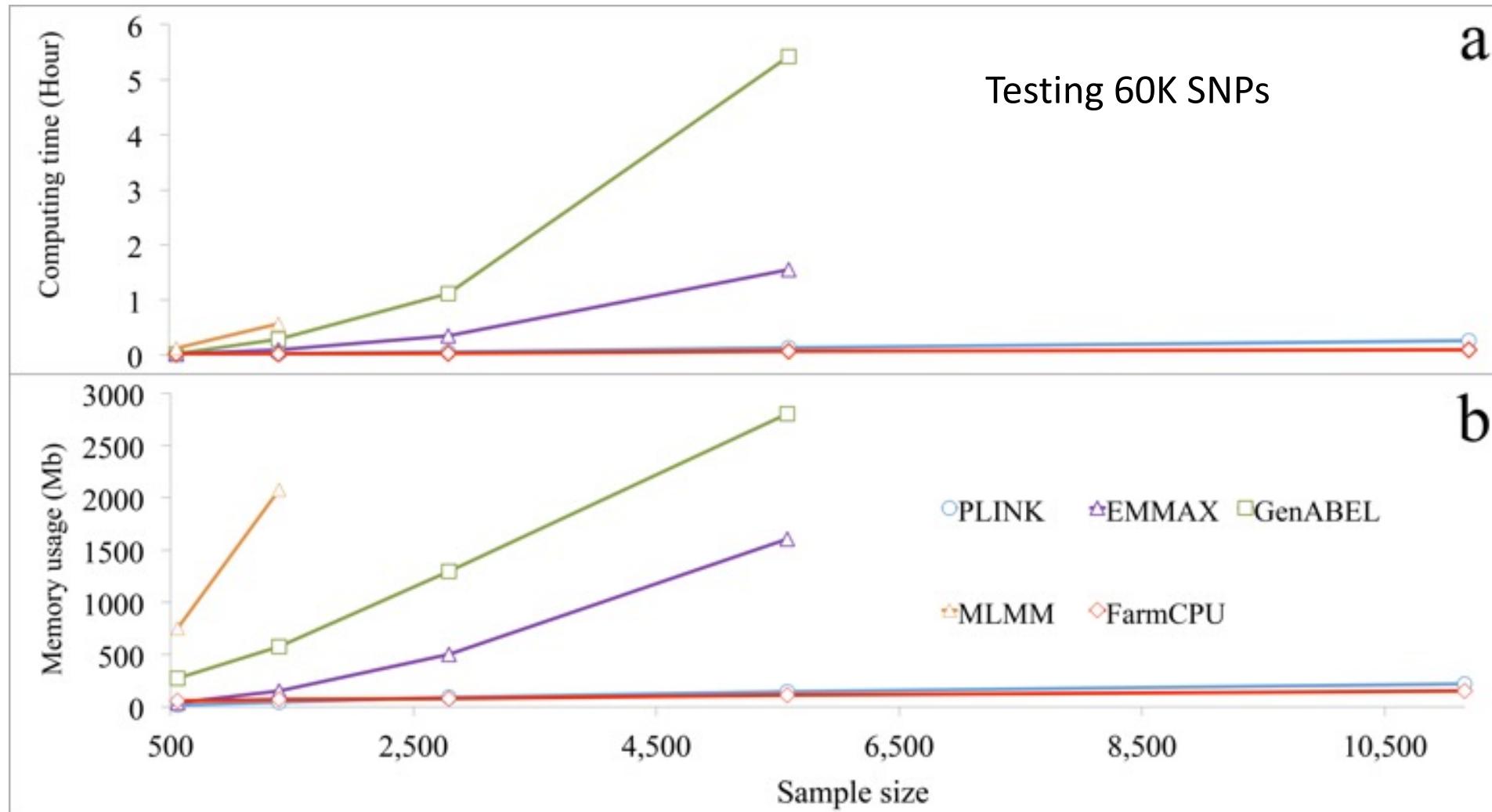
# FarmCPU algorithm



# It is time for human geneticists to move forward



# FarmCPU is computing efficient



ORIGINAL RESEARCH

 Open Access



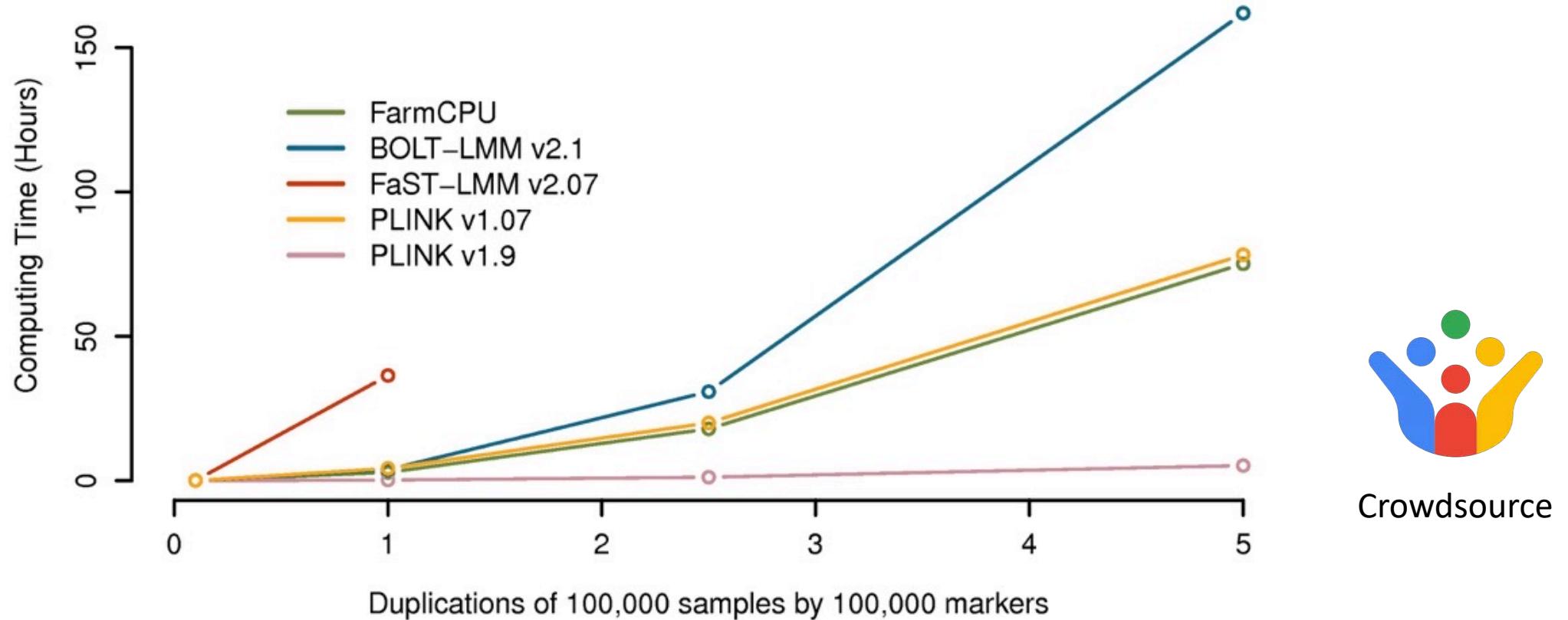
## FarmCPUp: Efficient large-scale genomewide association studies

Aaron Kusmec, Patrick S. Schnable 

First published: 10 April 2018 | <https://doi.org/10.1002/pld3.53> | Citations: 24

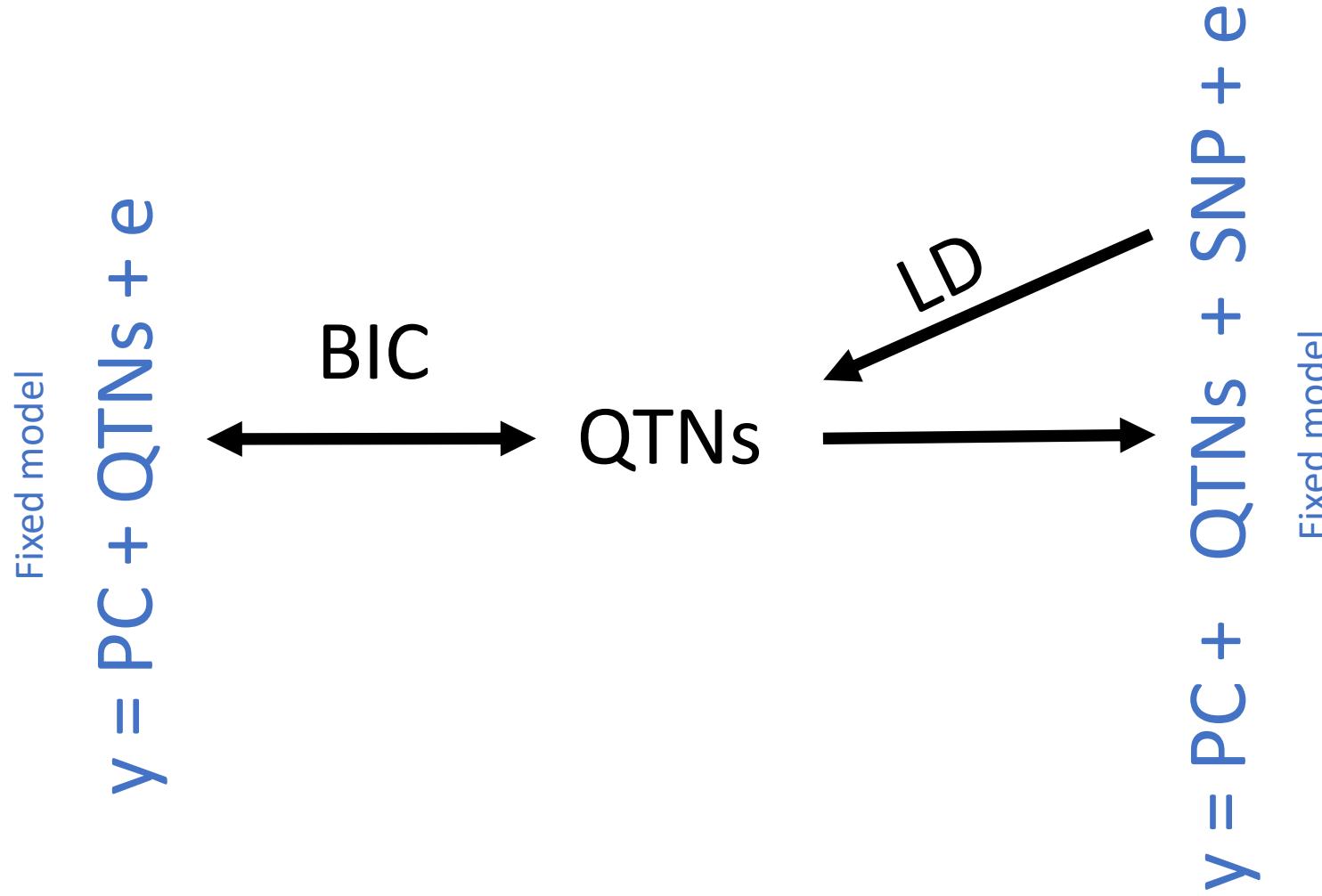
The preprint of this article can be found at  
<https://www.biorxiv.org/content/early/2017/12/24/238832>.

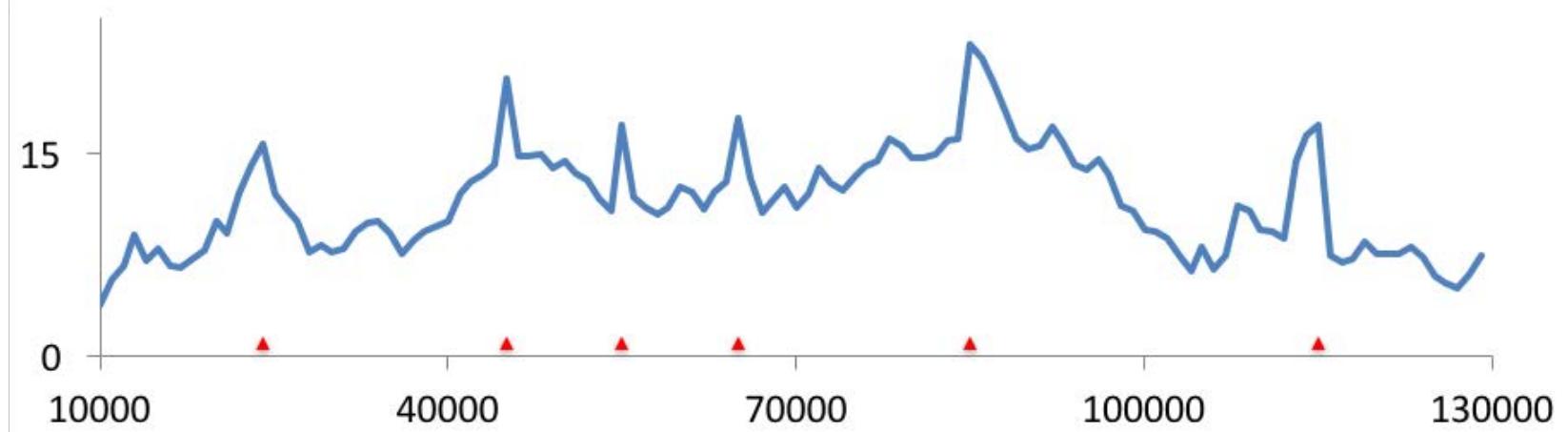
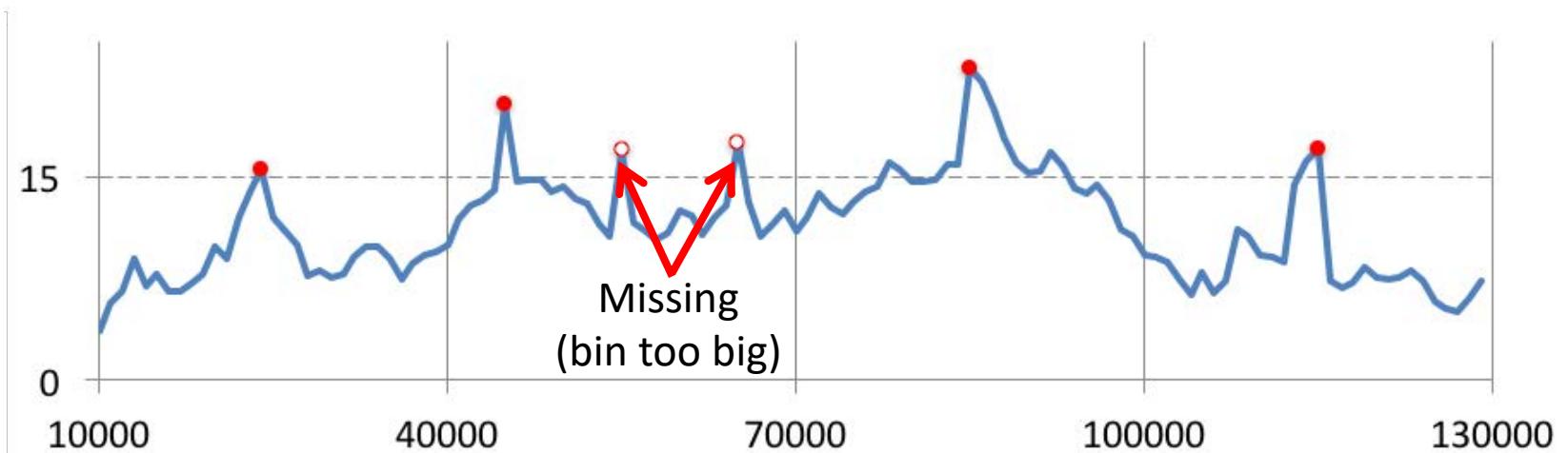
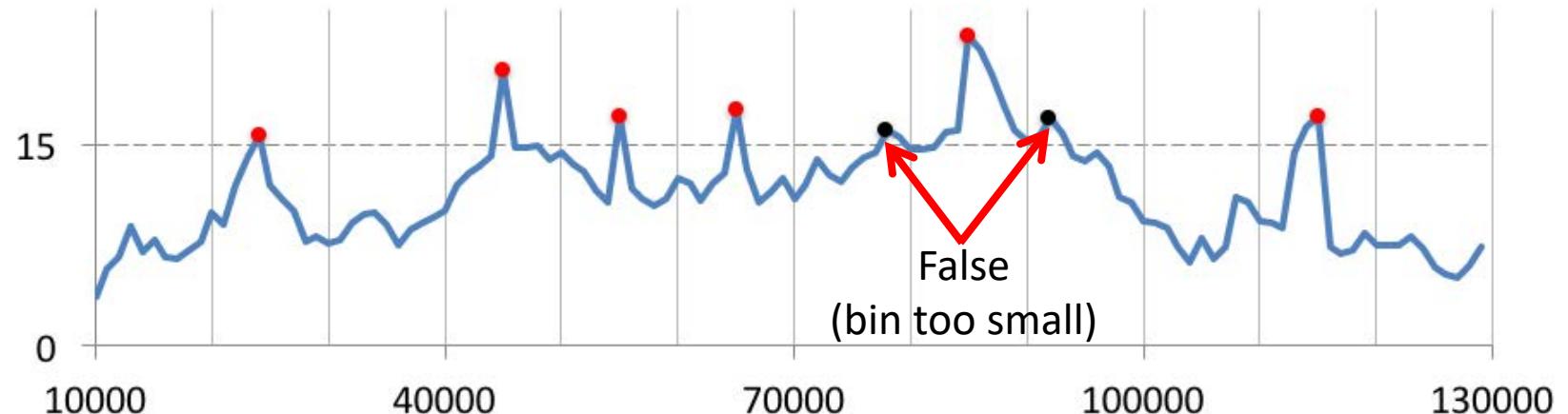
# Half million individuals, half million SNPs: three days



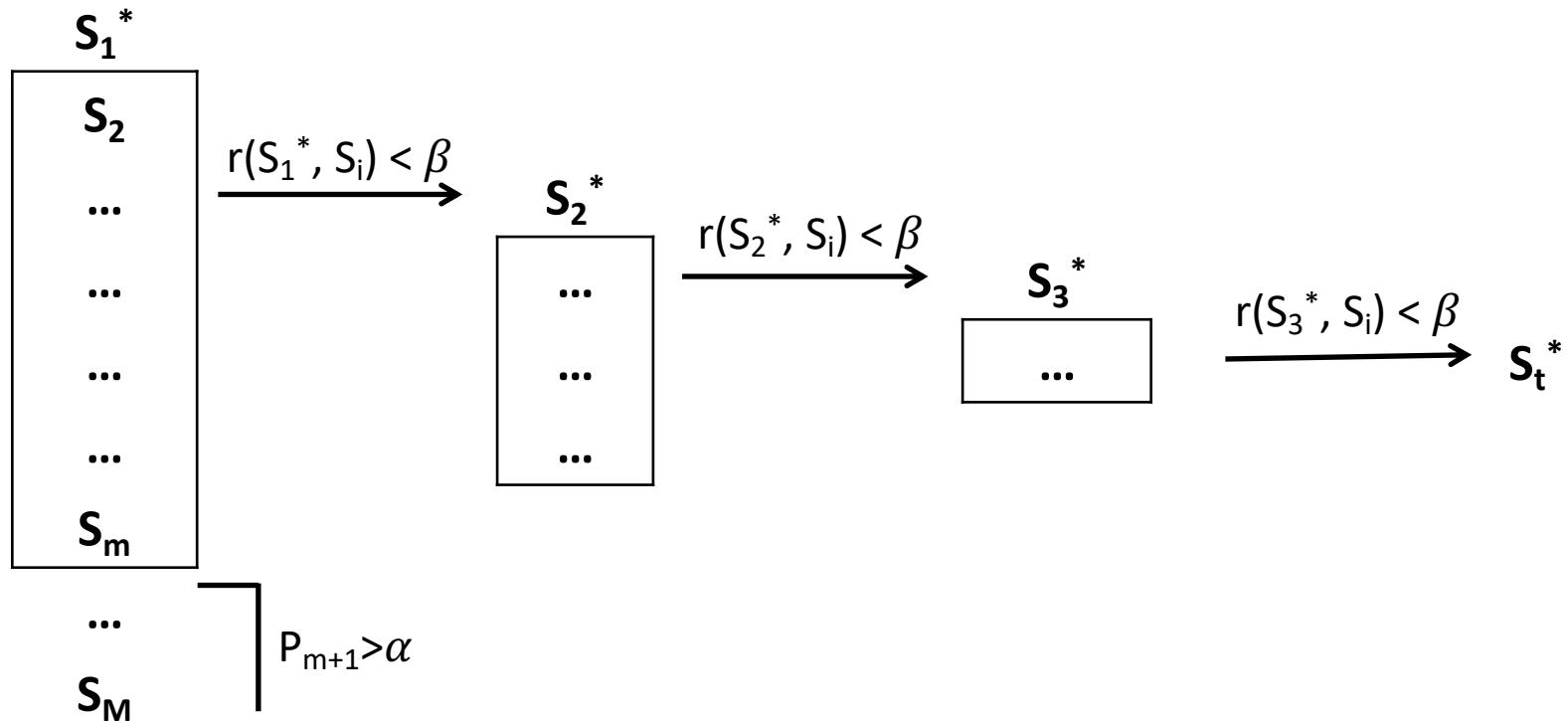
But, PLINK new version is faster

# BLINK algorithm





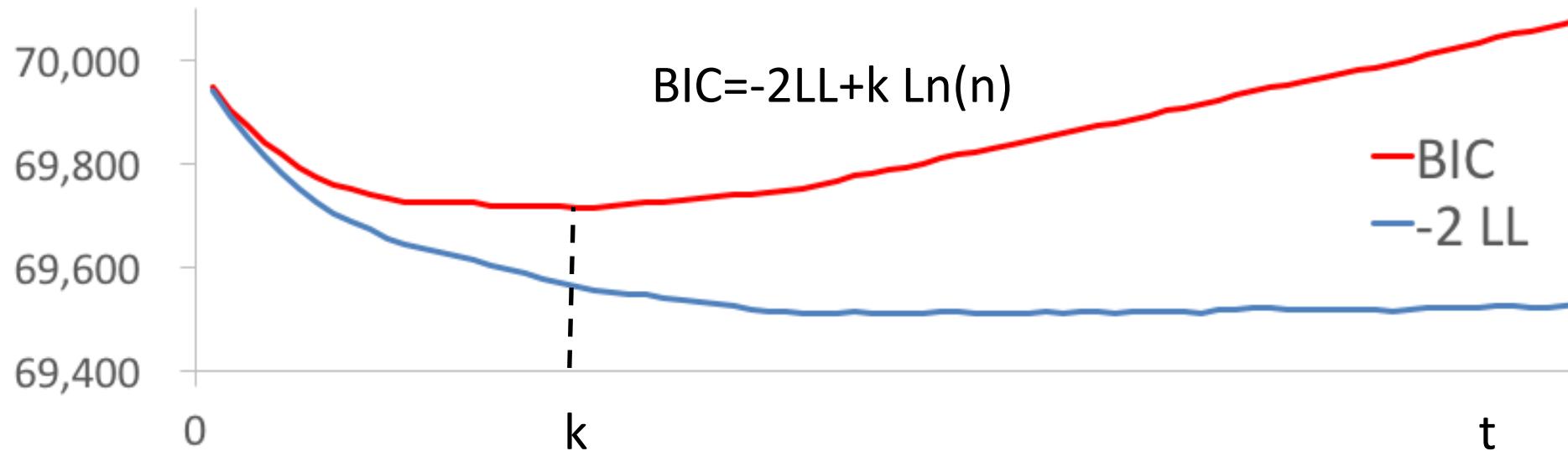
# Elimination of markers with LD

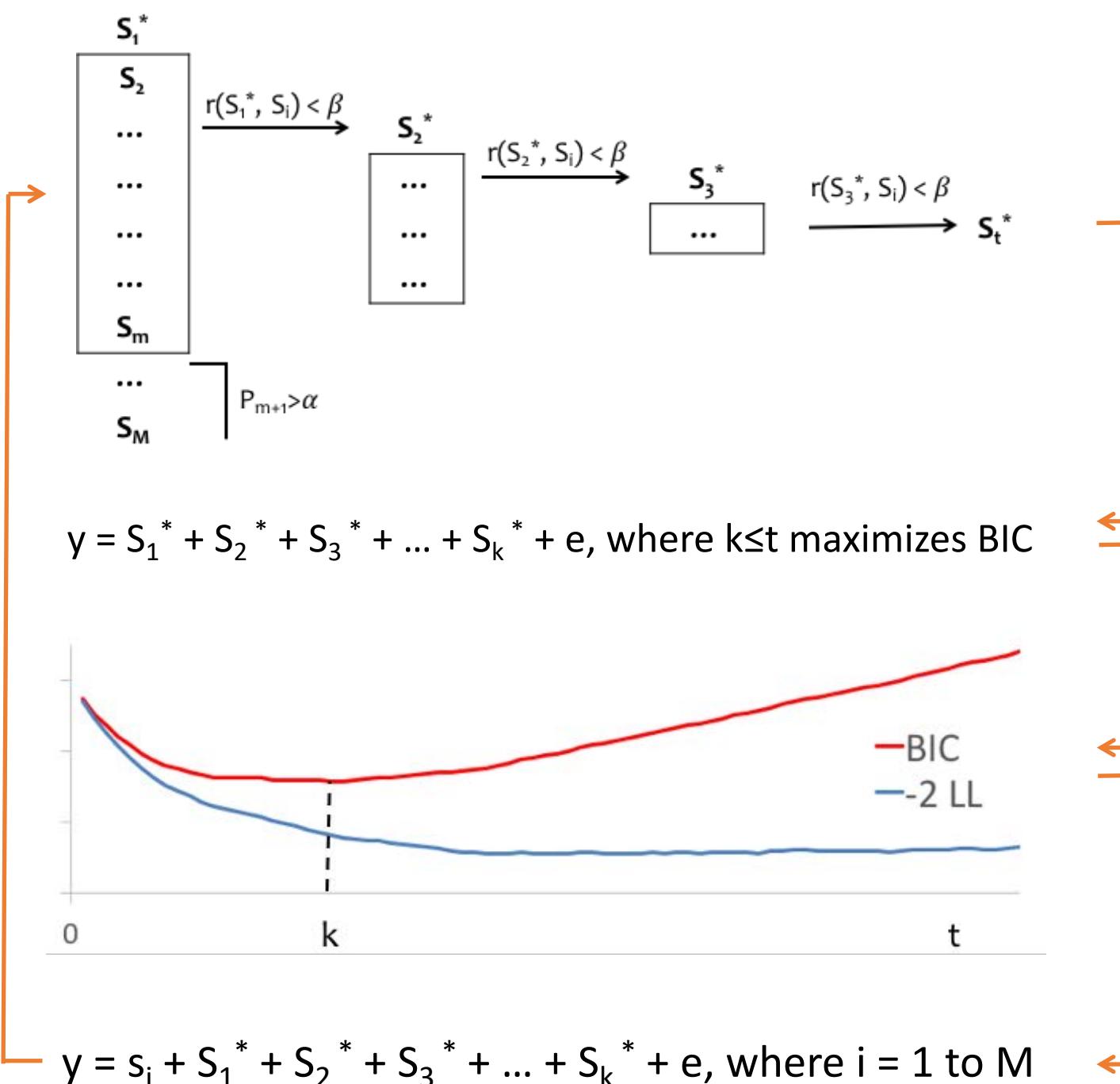


$$y = s_i + S_1^* + S_2^* + S_3^* + \dots + S_k^* + e, \text{ where } i = 1 \text{ to } M$$

# Bayesian information criterion

$$y = S_1^* + S_2^* + S_3^* + \dots + S_t^* + e, \text{ where } k \leq t \text{ maximizes BIC}$$

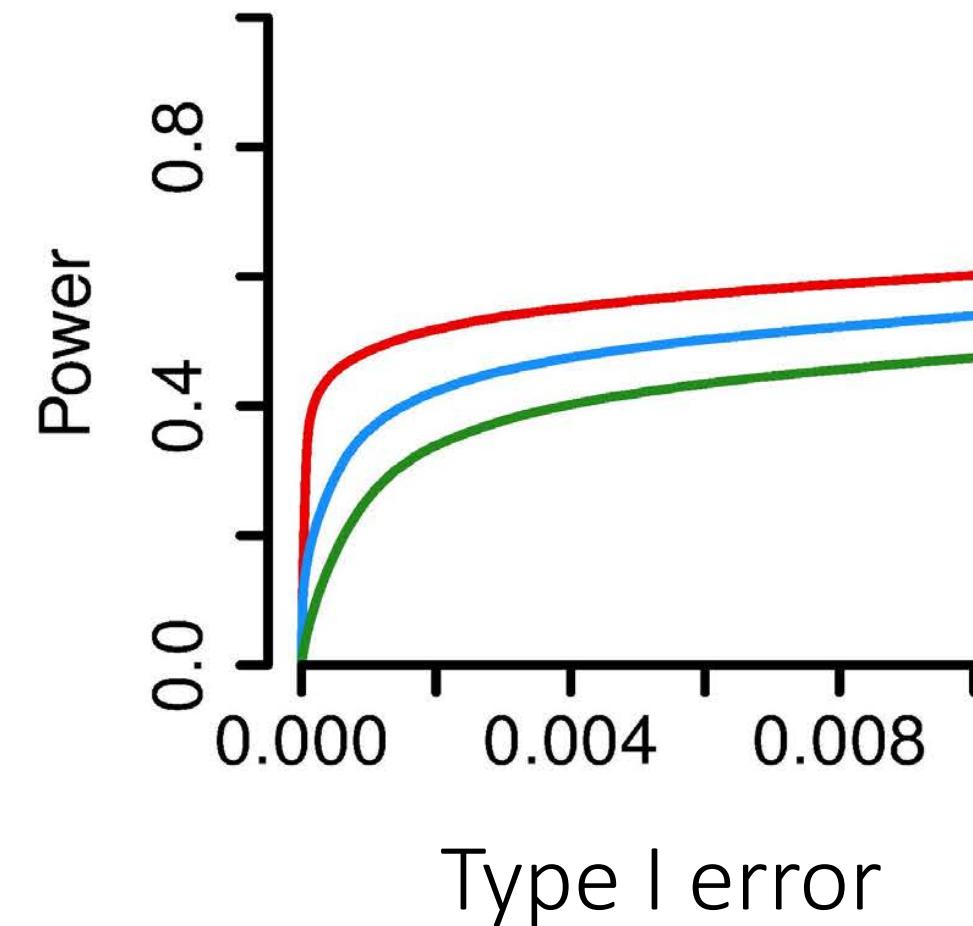
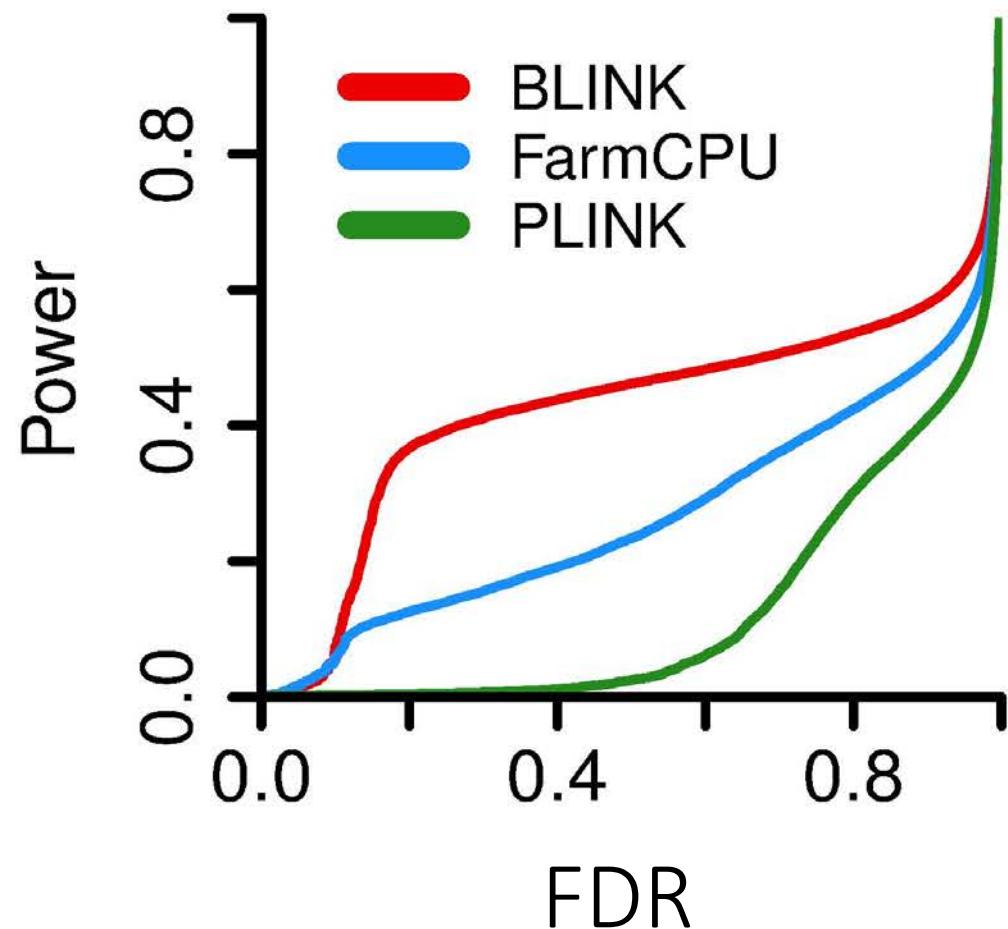




# Simulation study with human data



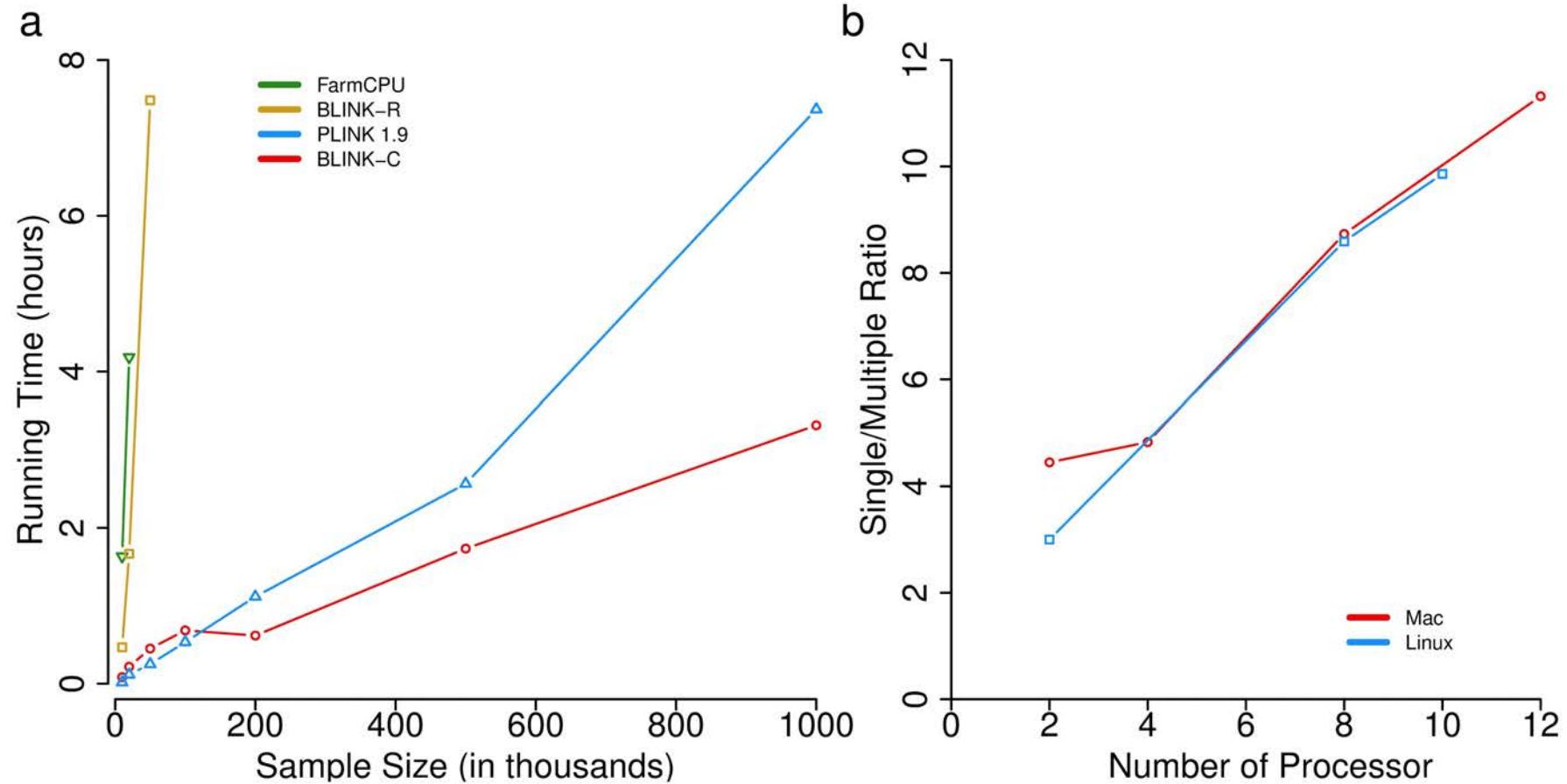
Meng Huang



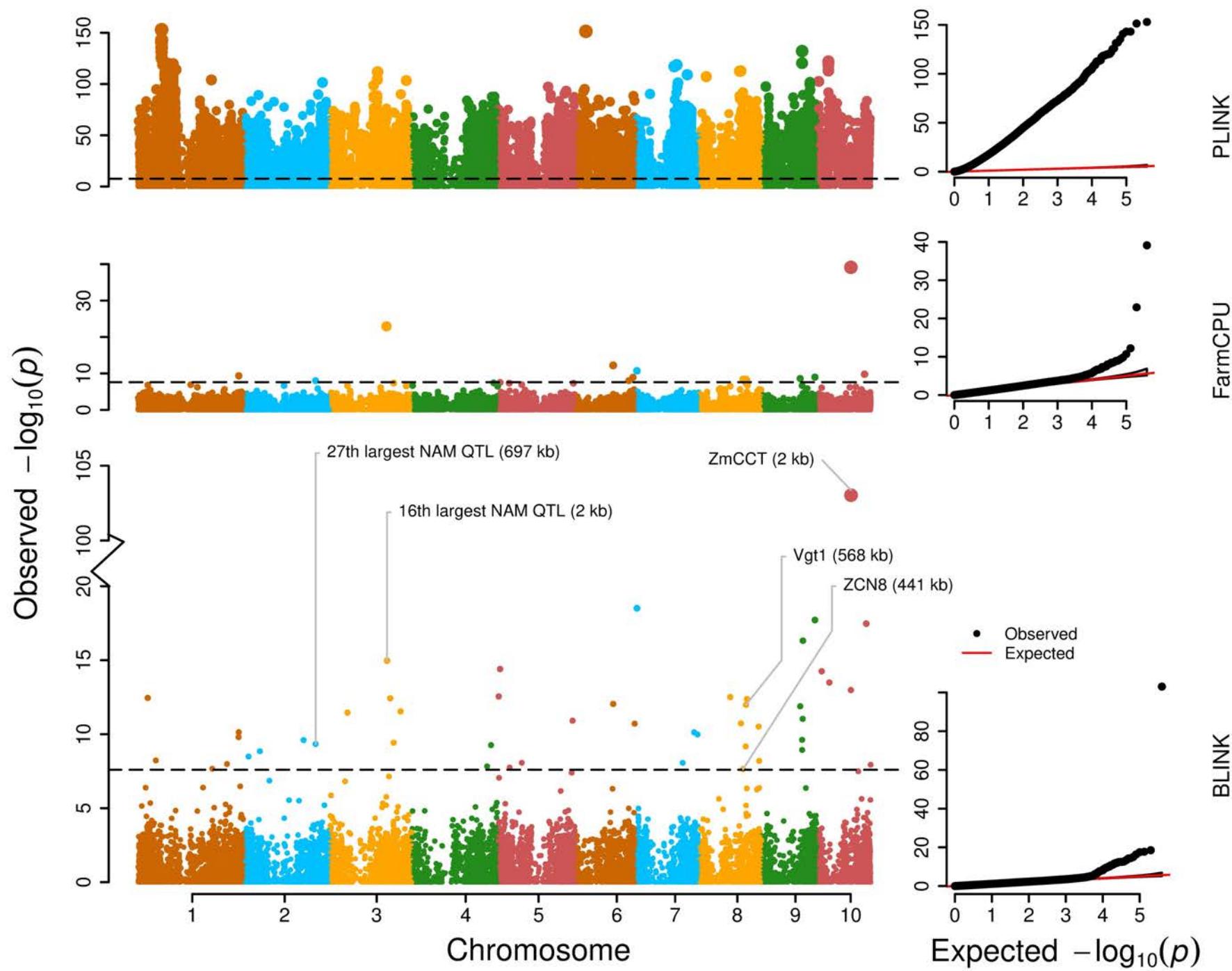
# 1M SNPs + 1M Samples = 3 hours



Meng Huang

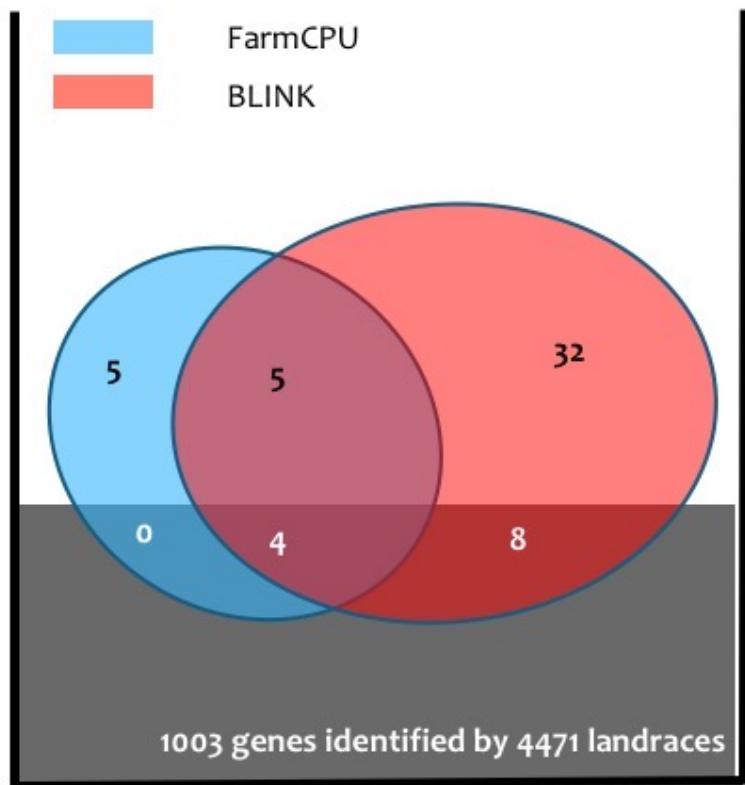


# Application in Maize

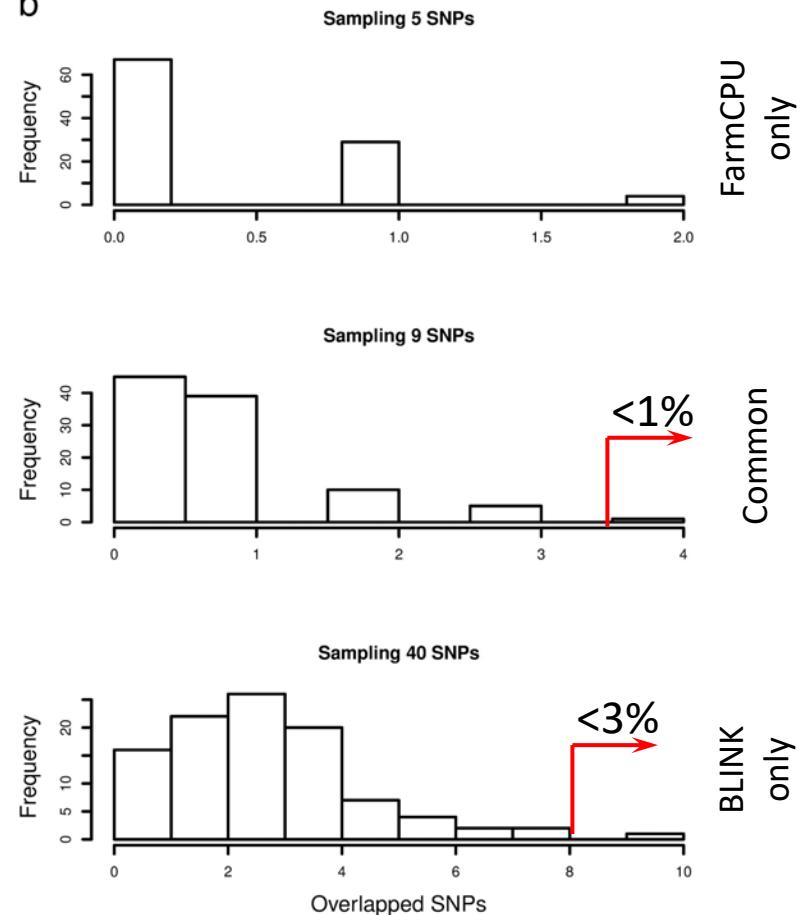


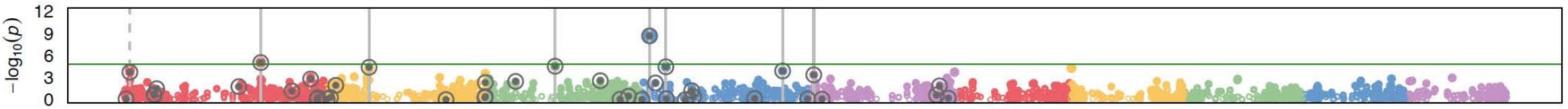
# Enrichment

a



b





All the versions have been implemented to [GAPIT](#) with R

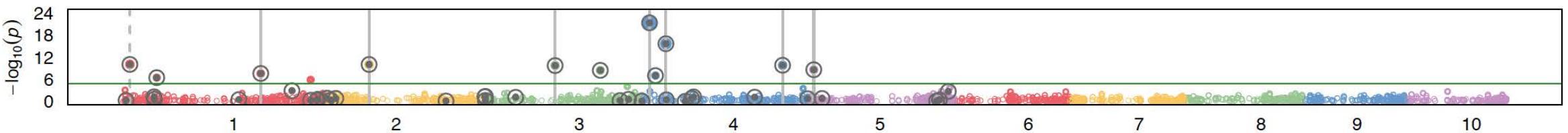
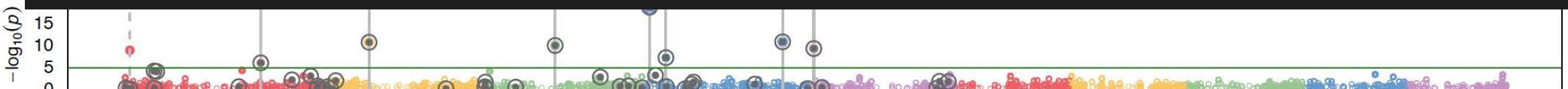
```

Code: source("http://zzlab.net/GAPIT/gapit_functions.txt")
      #Import demo data
      myGD=read.table(file="http://zzlab.net/GAPIT/data/mdp_numeric.txt",head=T)
      myGM=read.table(file="http://zzlab.net/GAPIT/data/mdp_SNP_information.txt",head=T)

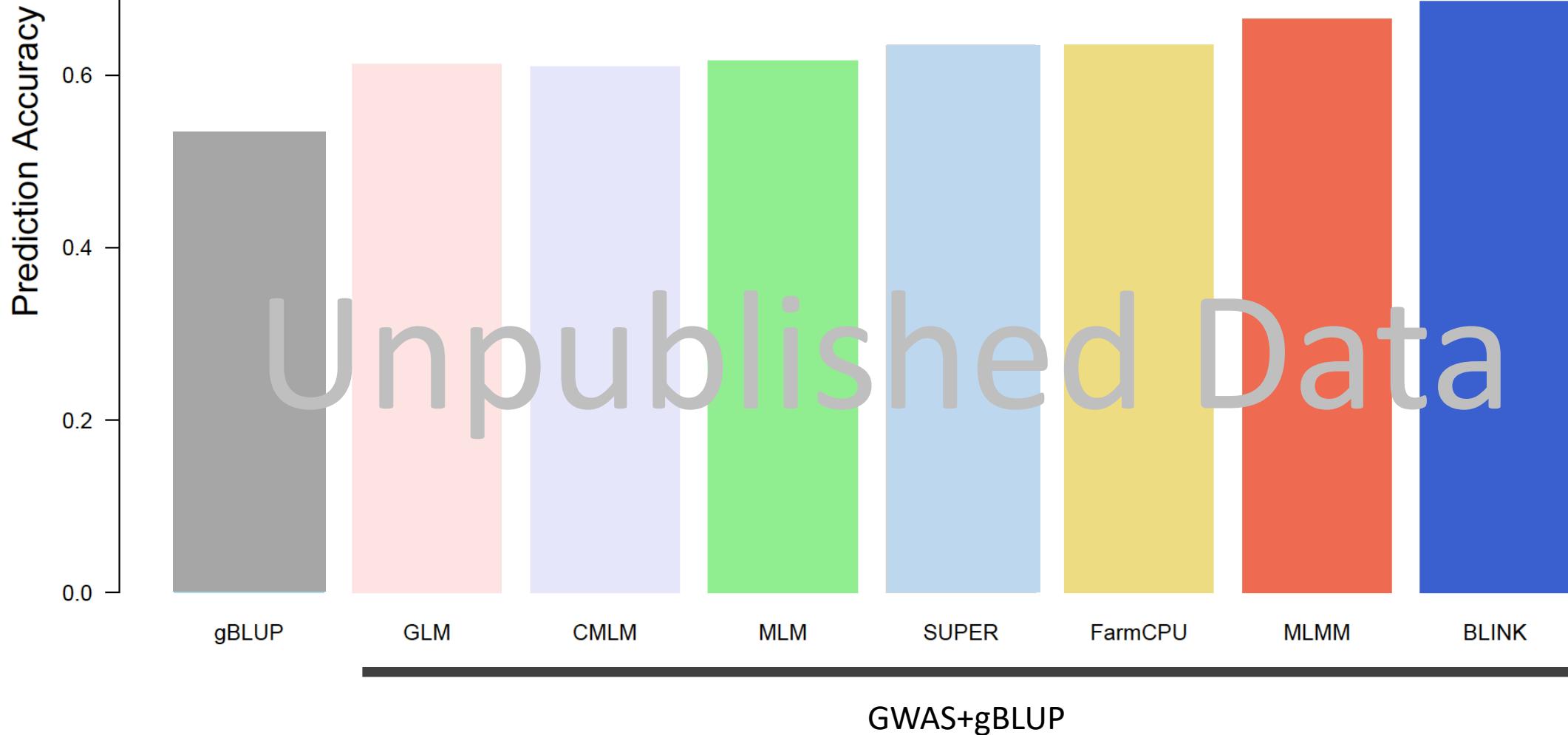
Integrate with GAPIT
believe me, it's not that hard.
with
#Simulate 10 QTN on the first half chromosomes
index1to5=myGM[,2]<6
set.seed(99164)
mySim=GAPIT.Phenotype.Simulation(GD=myGD[,c(TRUE,index1to5)],GM=myGM[index1to5,],h2=.7,NQTN=40, effectunit=.95,QTNDist="normal")

#GWAS with GAPIT
myGAPIT=GAPIT(Y=mySim$Y, GD=myGD, GM=myGM, PCA.total=3,
QTN.position=mySim$QTN.position,
model=c("GLM", "MLM", "CMLM", "SUPER", "MLMM", "FarmCPU", "Blink"))

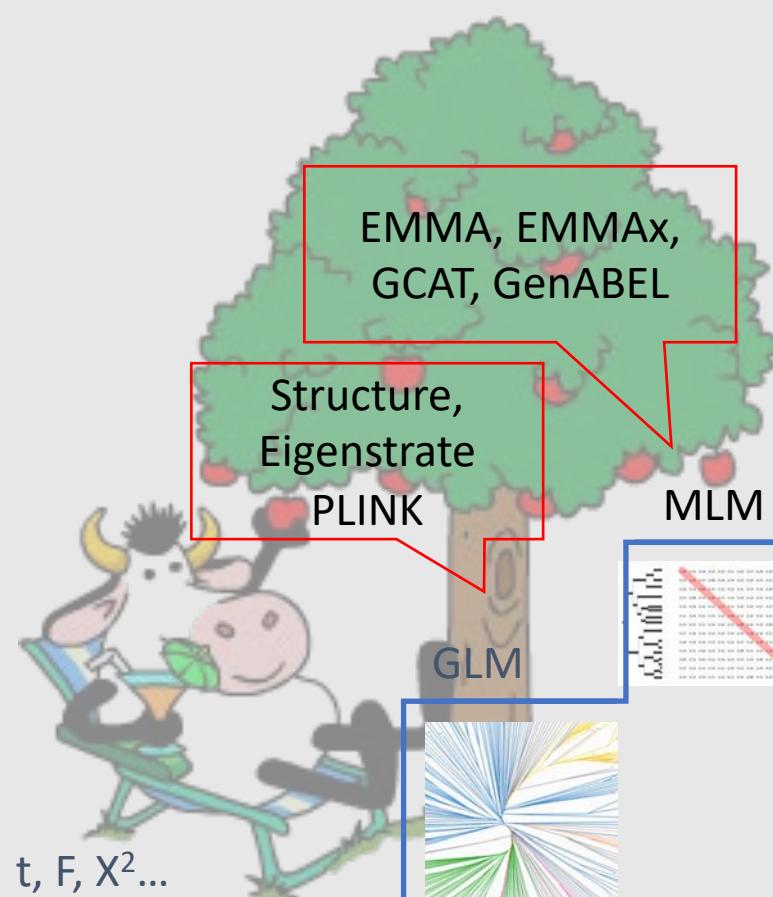
```



# Incorporating GWAS in GS



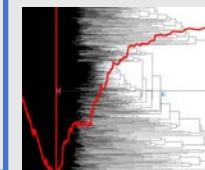
# GAPIT



Uncorrelated or  
equally correlated

TASSEL  
GAPIT

CMLM



MLM

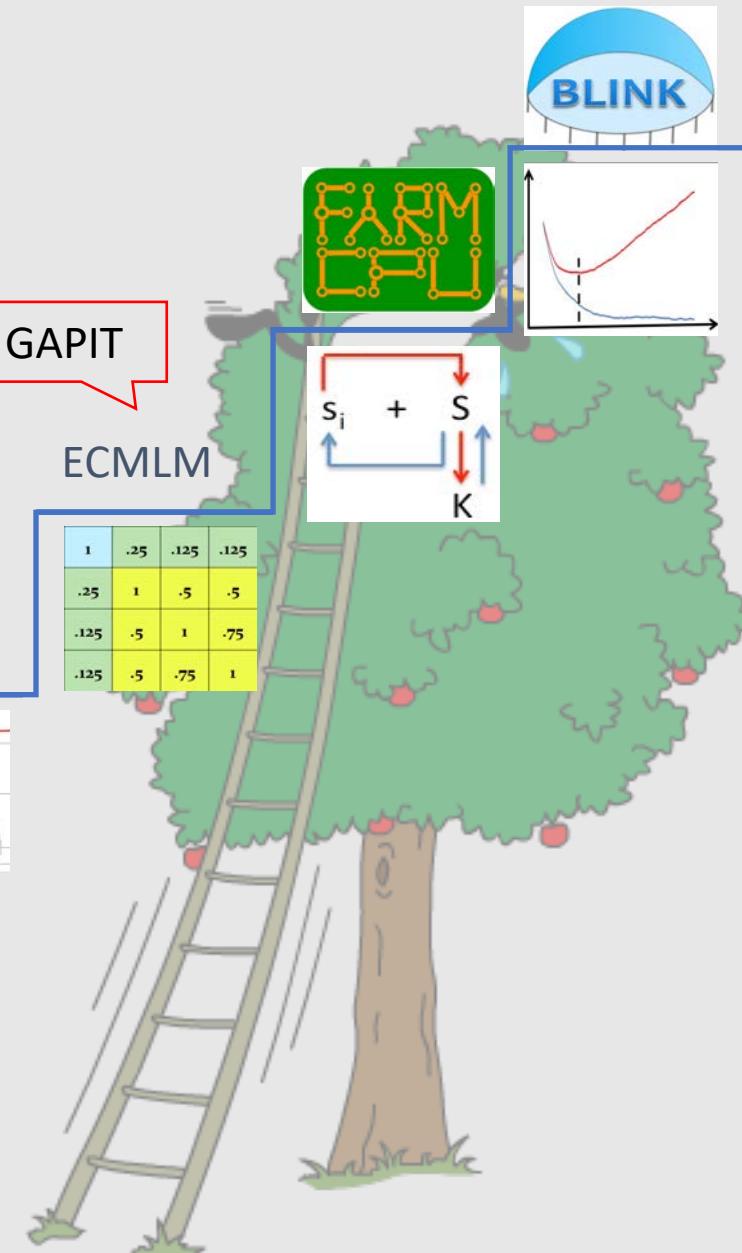
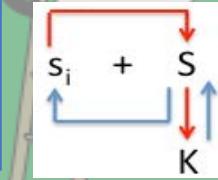


GLM

GAPIT

ECMLM

1	.25	.125	.125
.25	1	-.5	-.5
.125	-.5	1	-.75
.125	-.5	-.75	1



iPat



# 65th ISI World Statistics Congress 2025

The Hague

05 October 2025 - 09 October 2025

## Statistics Concourse Of Machine Learning And Artificial Intelligence

Organiser



Prof. Zhiwu Zhang



Thank you for your attention!

World Class Scenery, Research, Education, & Extension