

Gene Hunting, Bread Making, and GAPIT Engineering

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



Shiwu Zhang Laboratory for Statistical Genomics

Home

People

Publication

Research

Teaching

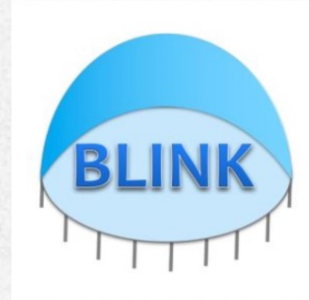
Software

Outreach

Jobs



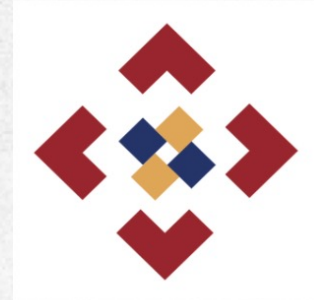
GAPIT



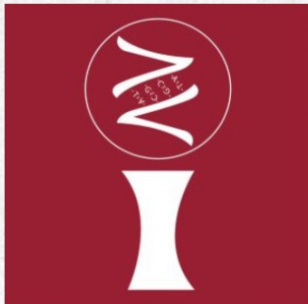
Blink



GRID



GridFree



iPat



FarmCPU



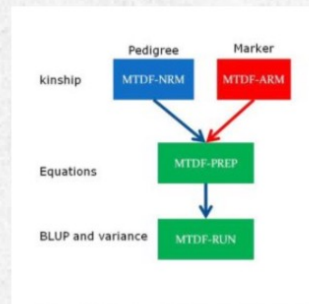
Rooster



LADDER



mMAP



MTDFREML



Audio4EDU

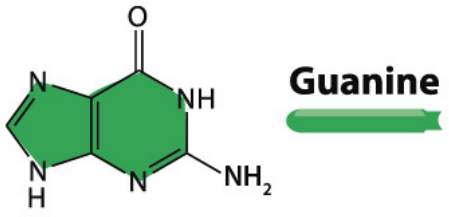
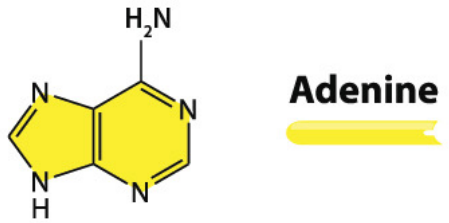
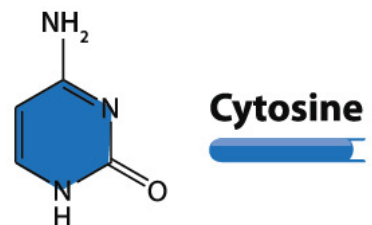
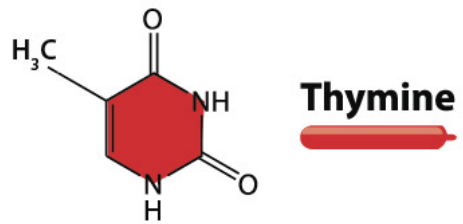


AI4EVER

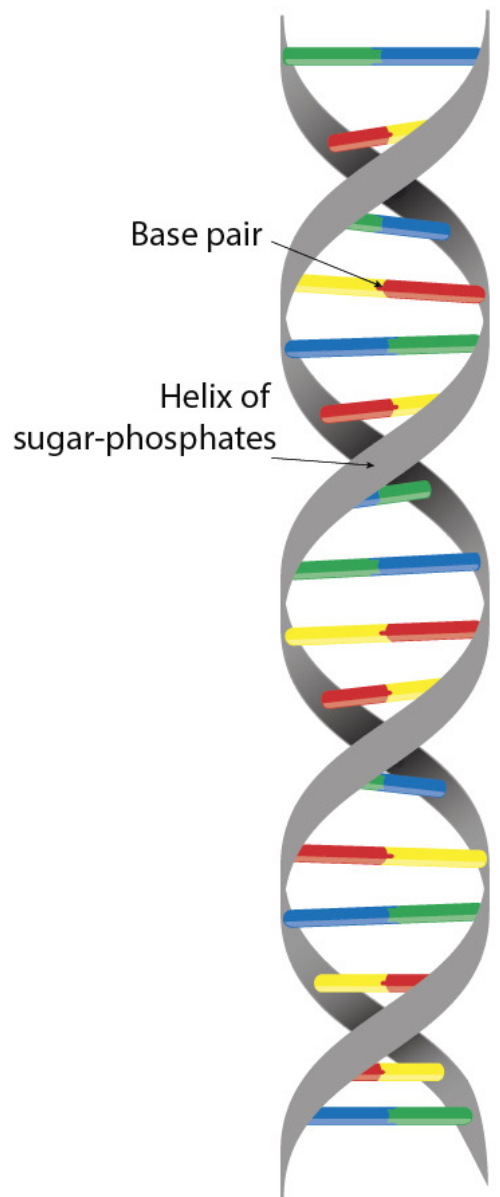




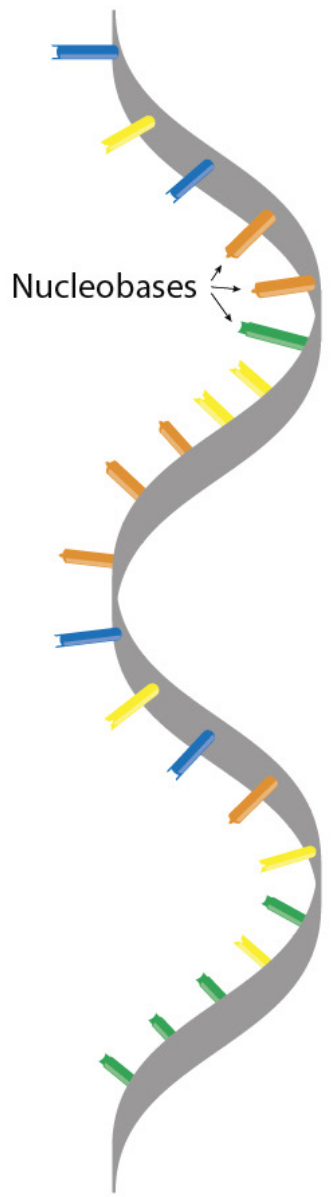
**James Watson and Francis Crick DNA Model
(Nobel Prize, 1962)**



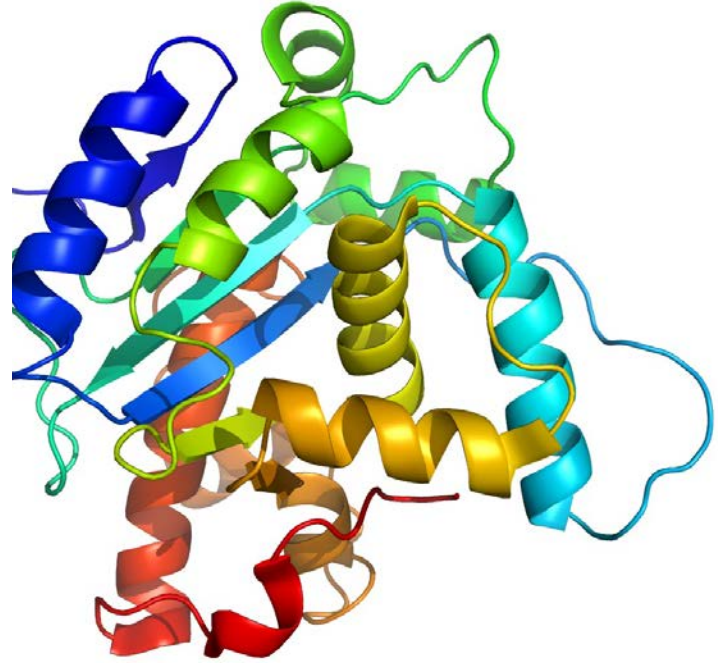
Nucleobases
of DNA



DNA
Deoxyribonucleic acid





RNA
Ribonucleic Acid



Genotypes

AA: 0

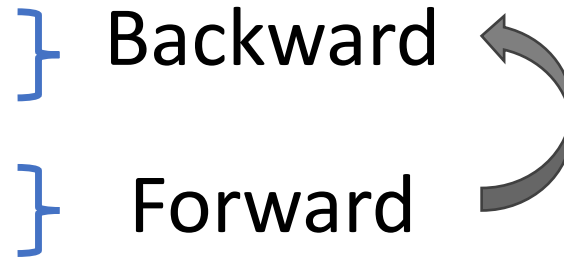
AT: 1 

TT: 2  

Genomic study

❖ Explanation

- Candidate gene
- Cloning
- Linkage analysis
- GWAS



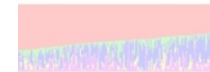
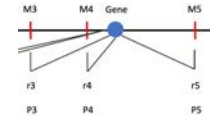
❖ Prediction

- MAS
- GS
- GWAS+GS
- AI

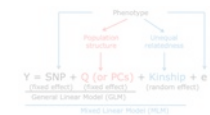


Outline

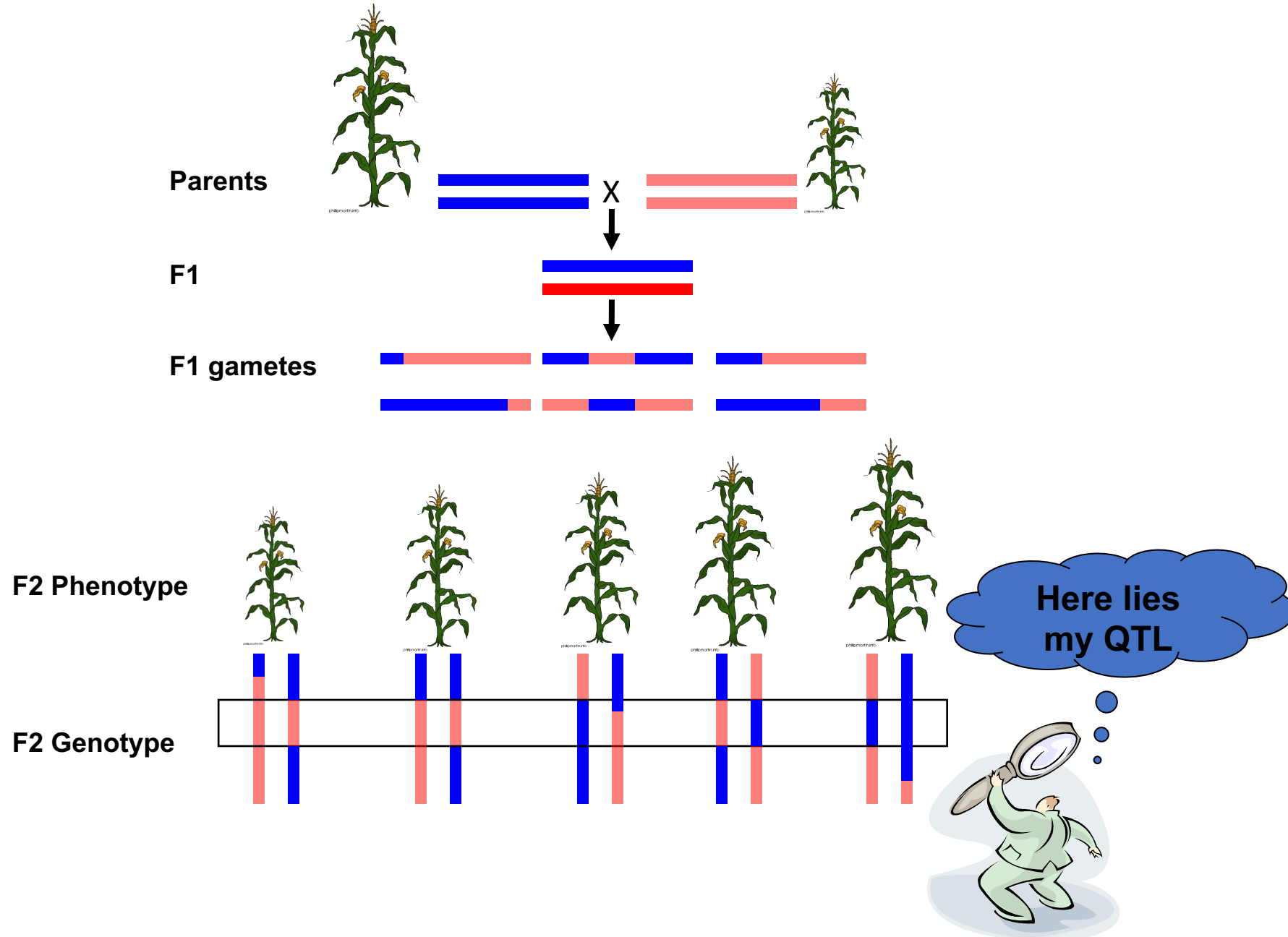
- Linkage analysis
- Association study
- Population structure and GLM
- Kinship and MLM
- BLINK



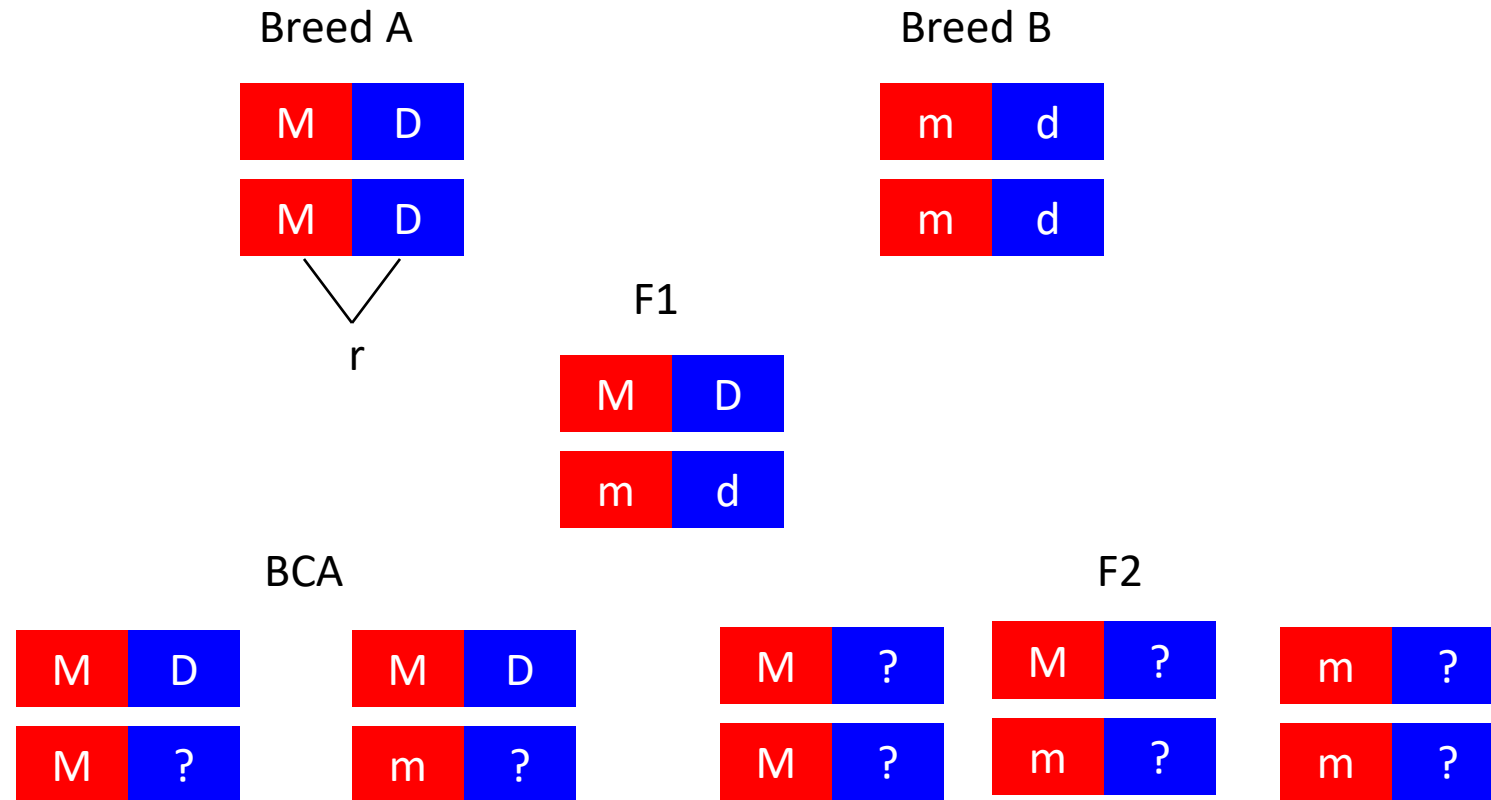
Q1	Q2	Q3
0.014	0.972	0.014
0.003	0.993	0.004
0.071	0.917	0.012
0.035	0.854	0.111
0.013	0.982	0.005



Linkage analysis

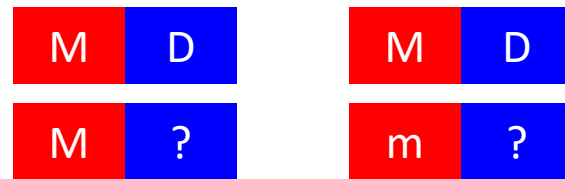


Crosses



Probability

BCA



$$P(?=D \mid MM)=1-r$$

$$P(?=d \mid MM)=r$$

$$P(?=D \mid Mm)=r$$

$$P(?=d \mid Mm)=1-r$$

	DD	Dd
MM	n1	n2
Mm	n3	n4

Recombine

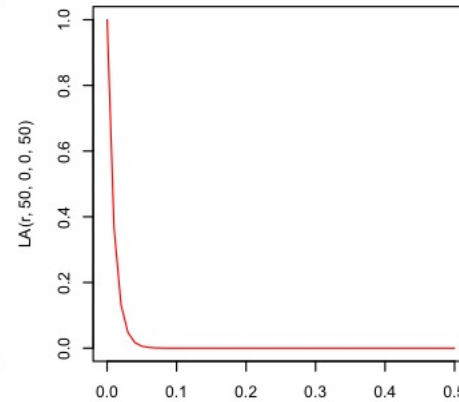
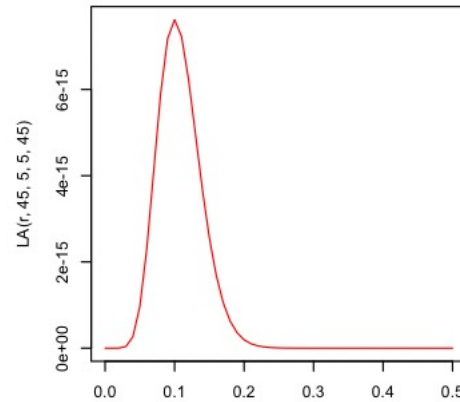
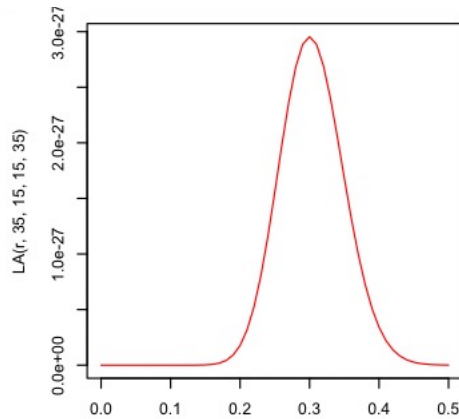
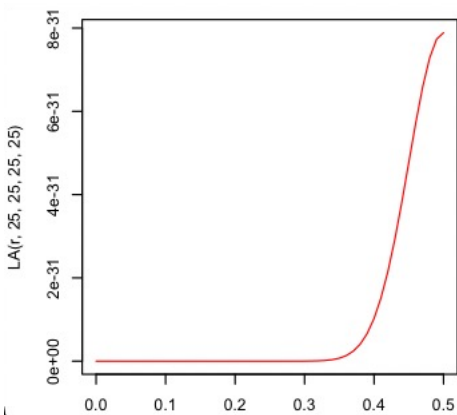
Unrecombine

$$P = r^{(n2+n3)} (1-r)^{(n1+n4)}$$

Mapping: vary r to maximize P

$$P = r^{(n2+n3)} (1-r)^{(n1+n4)}$$

	D	d		D	d		D	d		D	d
MM	25	25	MM	35	15	MM	45	5	MM	50	0
Mm	25	25	Mm	15	35	Mm	5	45	Mm	0	50



```
r=seq(0, .5, .01)
```

```
LA=function(r, n1, n2, n3, n4) {return(r^(n2+n3) * (1-r)^(n1+n4)) }
```

```
par(mfrow=c(1,4),mar = c(3,4,1,1))
```

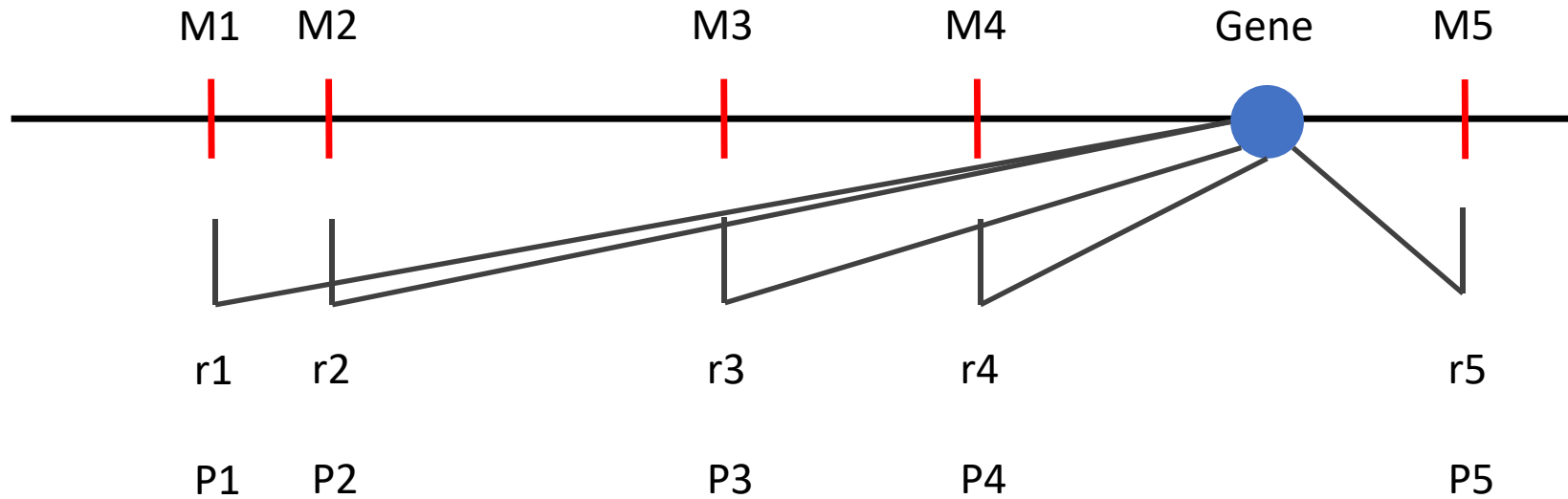
```
plot(r, LA(r, 25, 25, 25, 25), type="l", col="red")
```

```
plot(r, LA(r, 35, 15, 15, 35), type="l", col="red")
```

```
plot(r, LA(r, 45, 5, 5, 45), type="l", col="red")
```

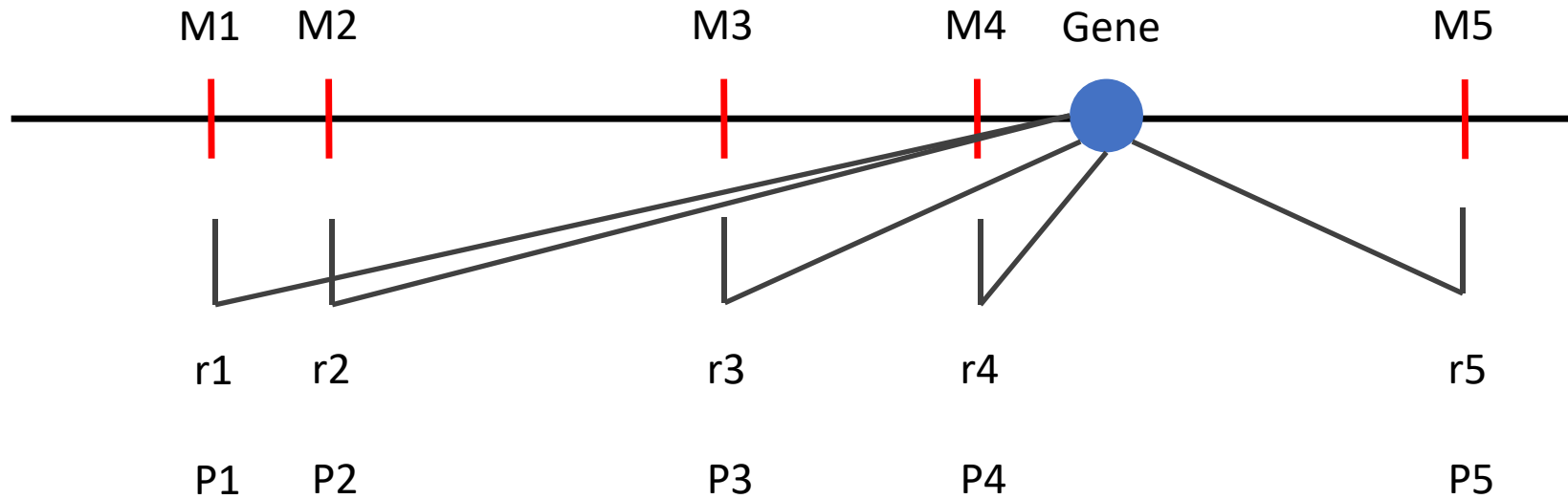
```
plot(r, LA(r, 50, 0, 0, 50), type="l", col="red")
```

Multiple markers



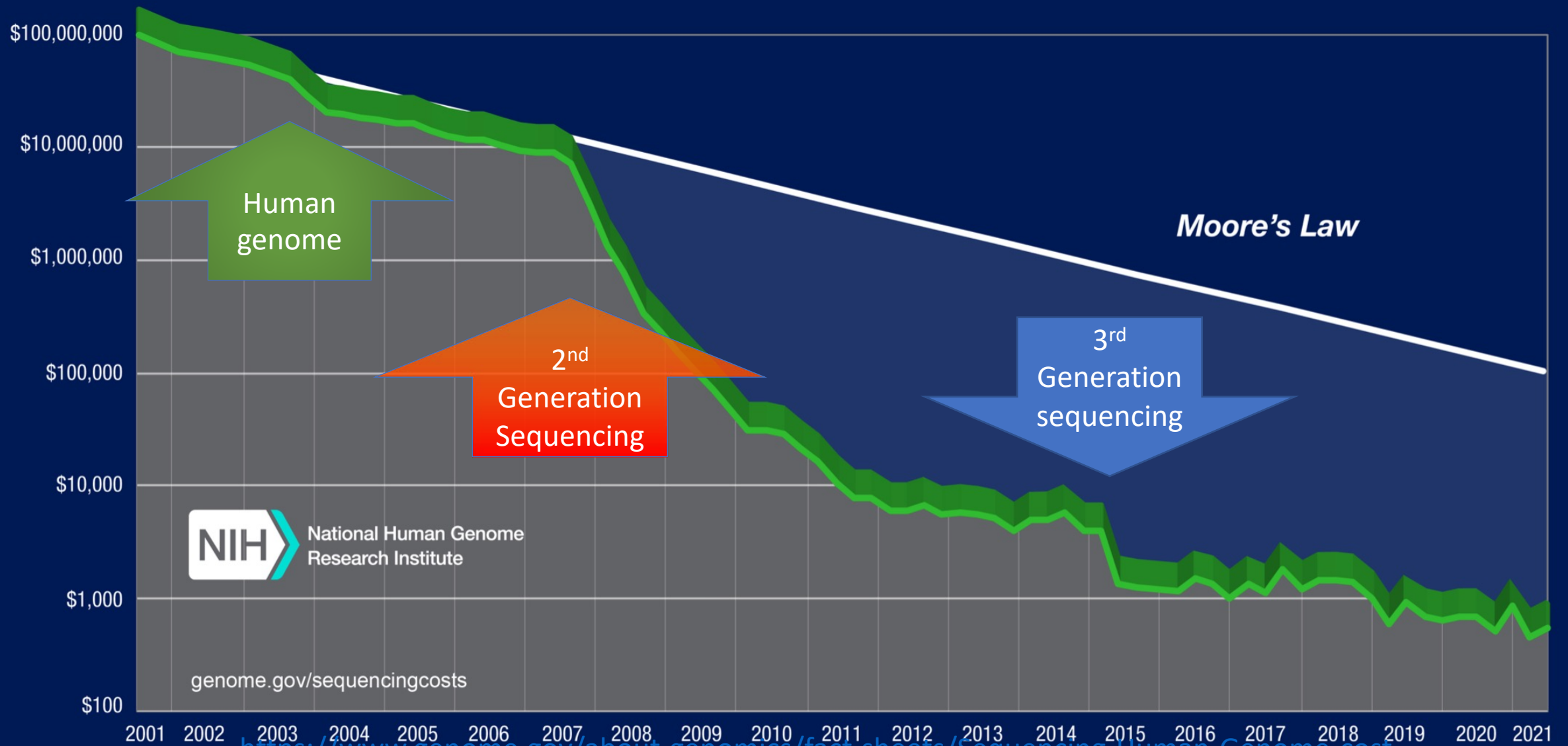
$$P = P1 * P2 * P3 * P4 * P5$$

Multiple markers



$$P = P1 * P2 * P3 * P4 * P5$$

Cost per Human Genome

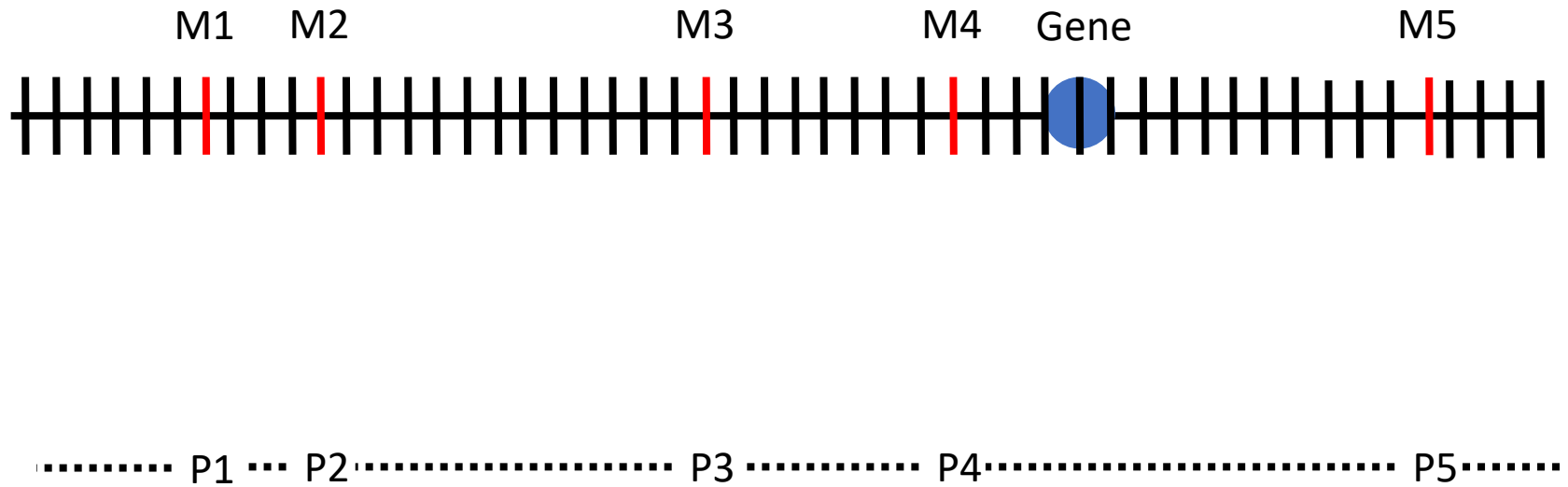


NIH National Human Genome Research Institute

genome.gov/sequencingcosts

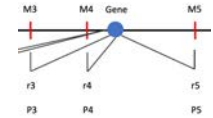
<https://www.genome.gov/about-genomics/fact-sheets/Sequencing-Human-Genome-cost>

Dense markers (GWAS)

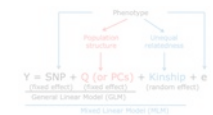


Outline

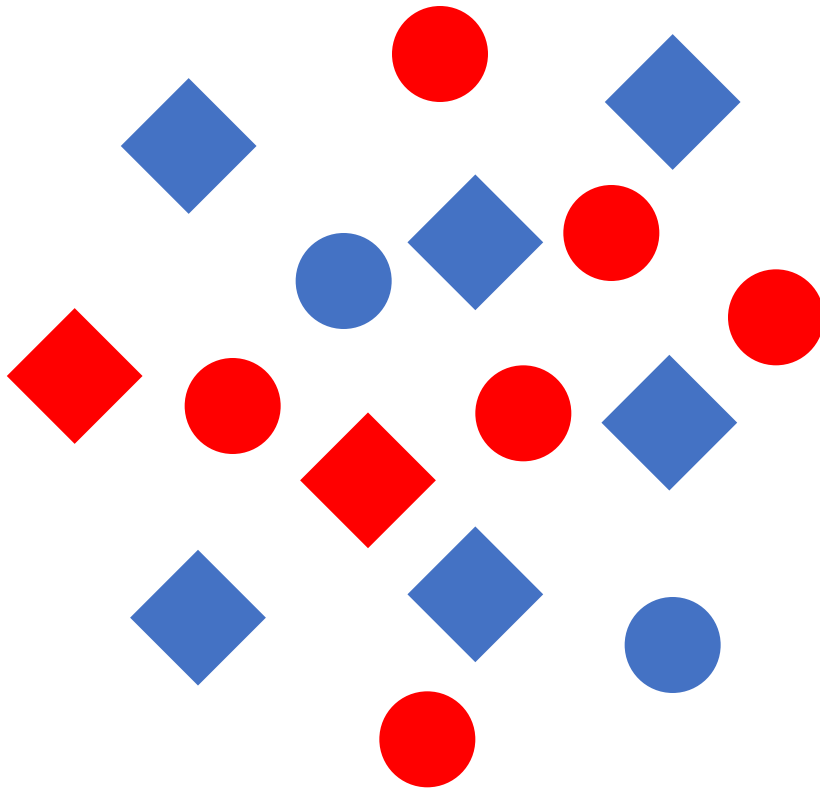
- Linkage analysis
- Association study
- Population structure and GLM
- Kinship and MLM
- BLINK



Q1	Q2	Q3
0.014	0.972	0.014
0.003	0.993	0.004
0.071	0.917	0.012
0.035	0.854	0.111
0.013	0.982	0.005



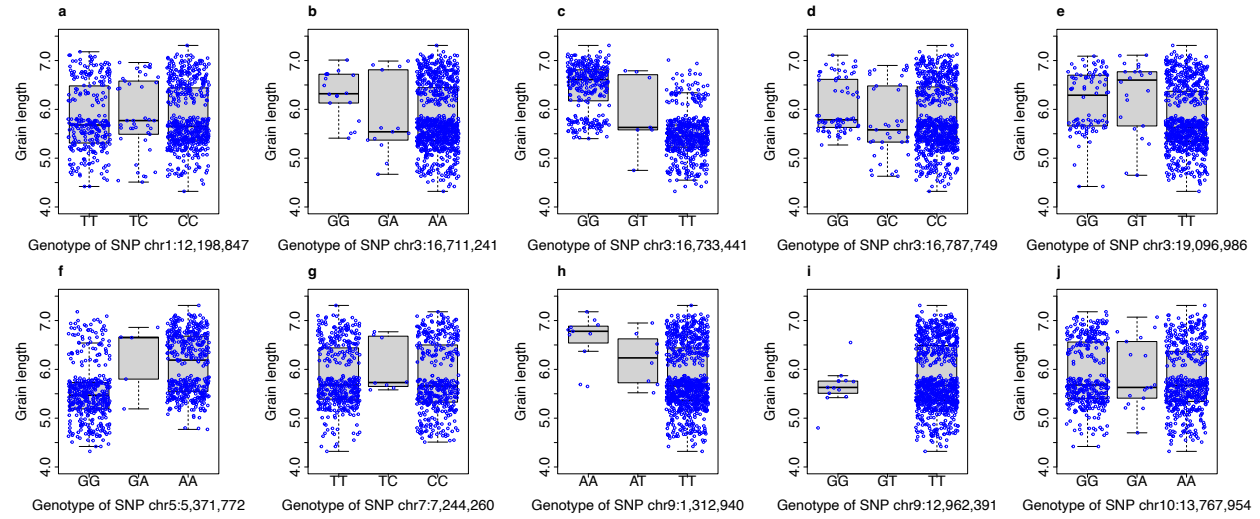
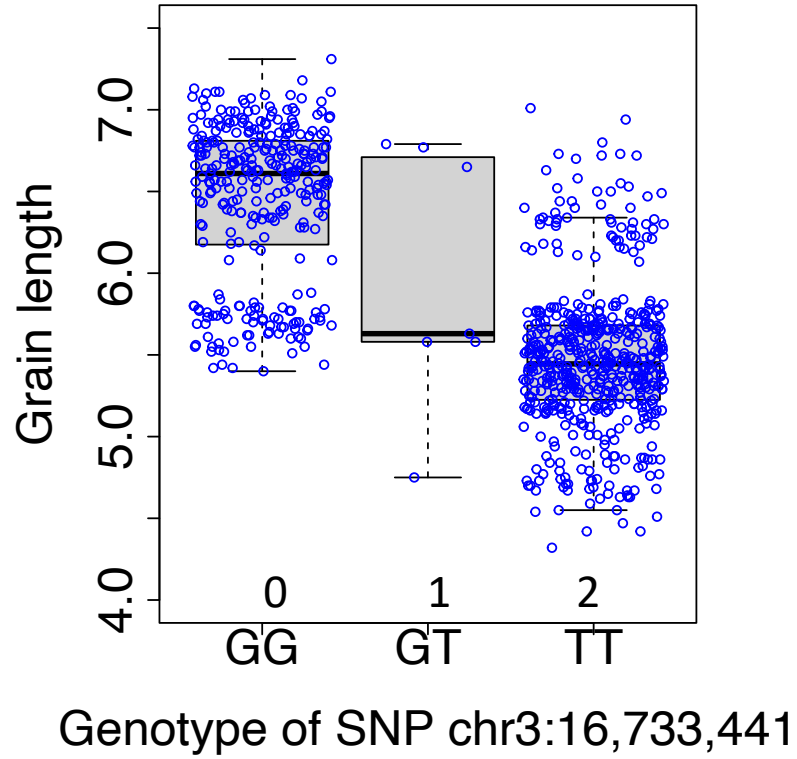
Association study






Marker	Control	Case
	6	2
	2	6

$$X^2 = 4(2 * 2 / 4) = 4, \text{ df} = 1, \\ P = 4.5\%$$

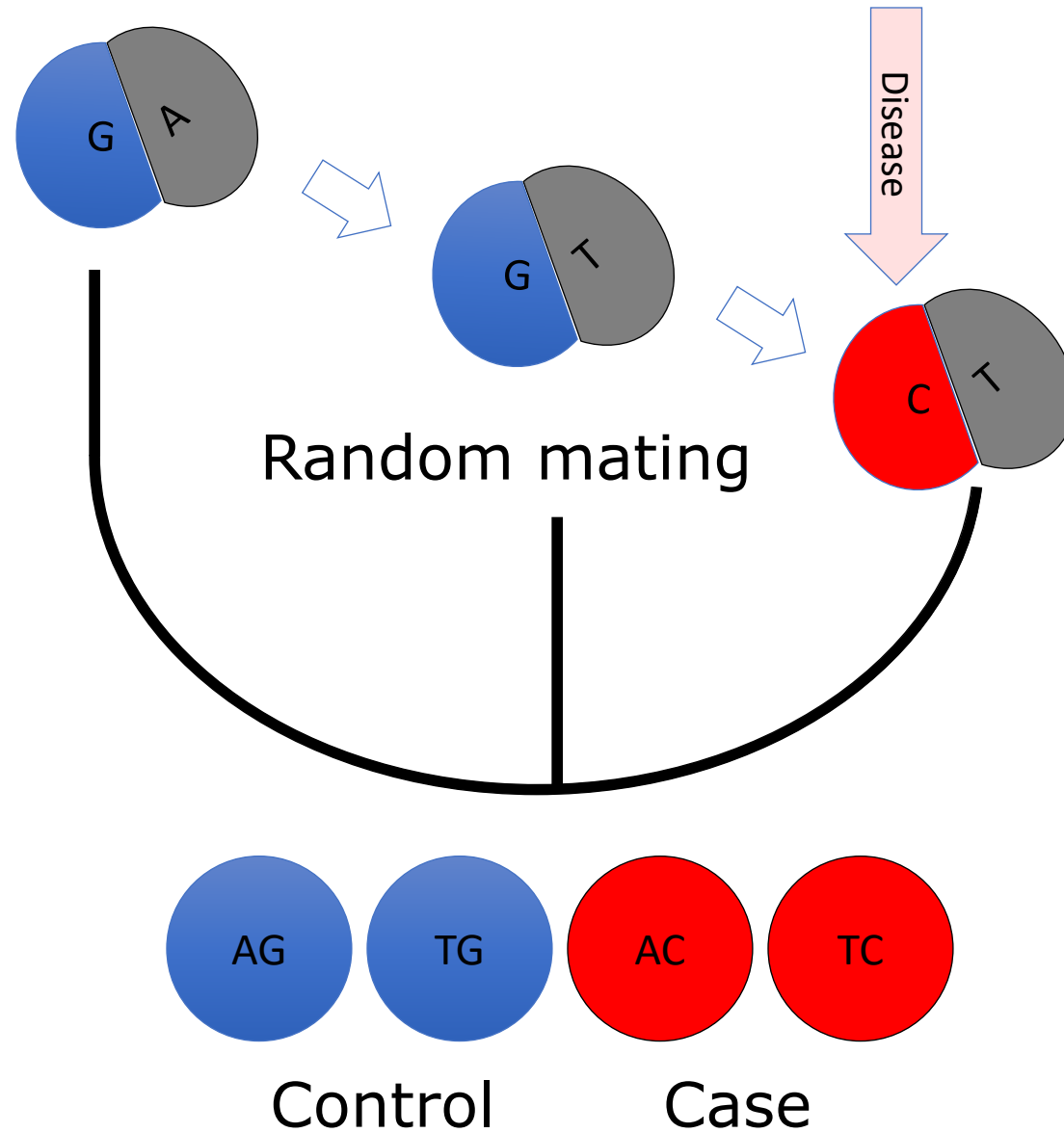
Correlation



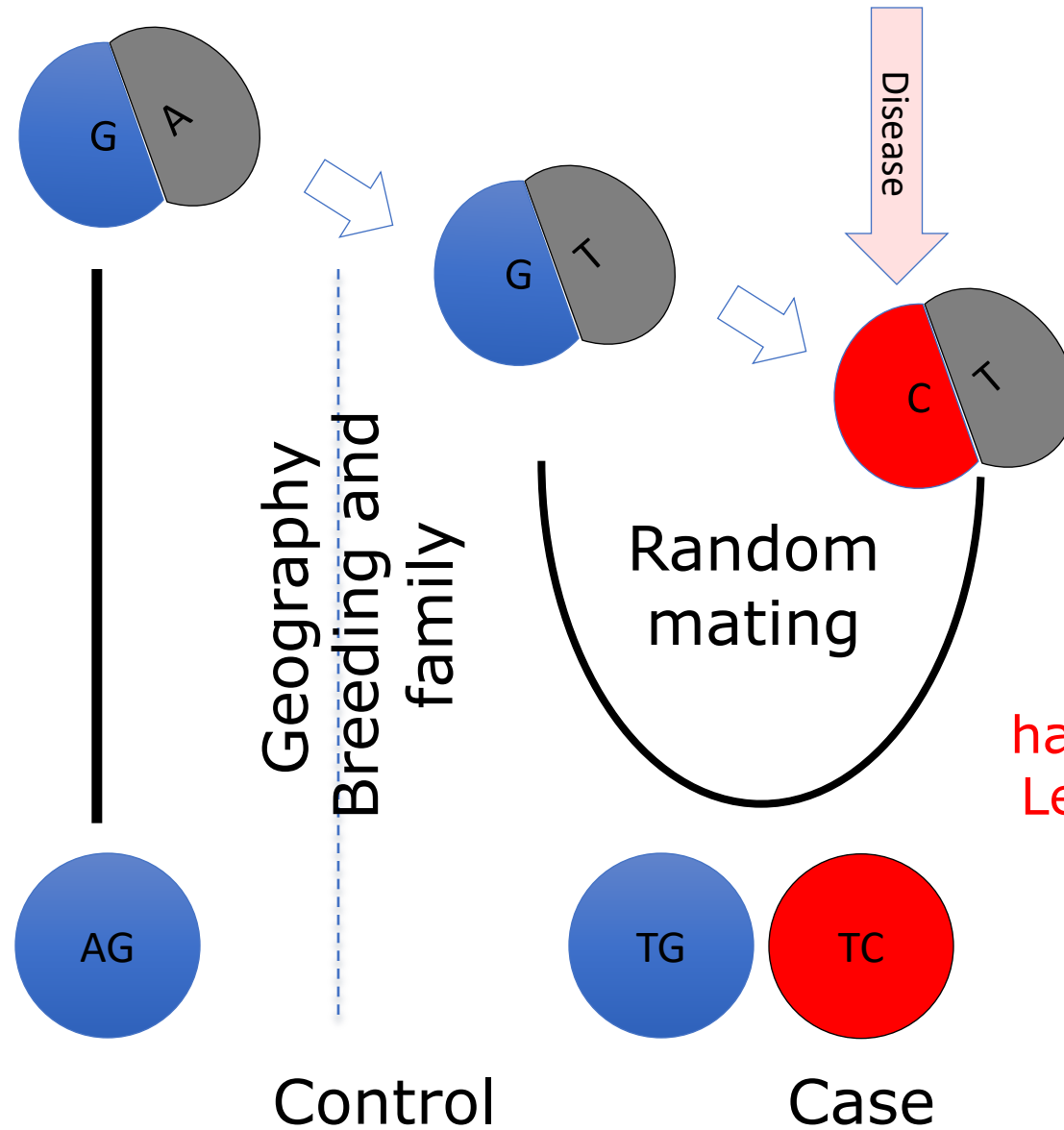
Comparison between linkage analysis and GWAS

Property	Linkage analysis	GWAS
Resolution		
Generation		
Genetic base		

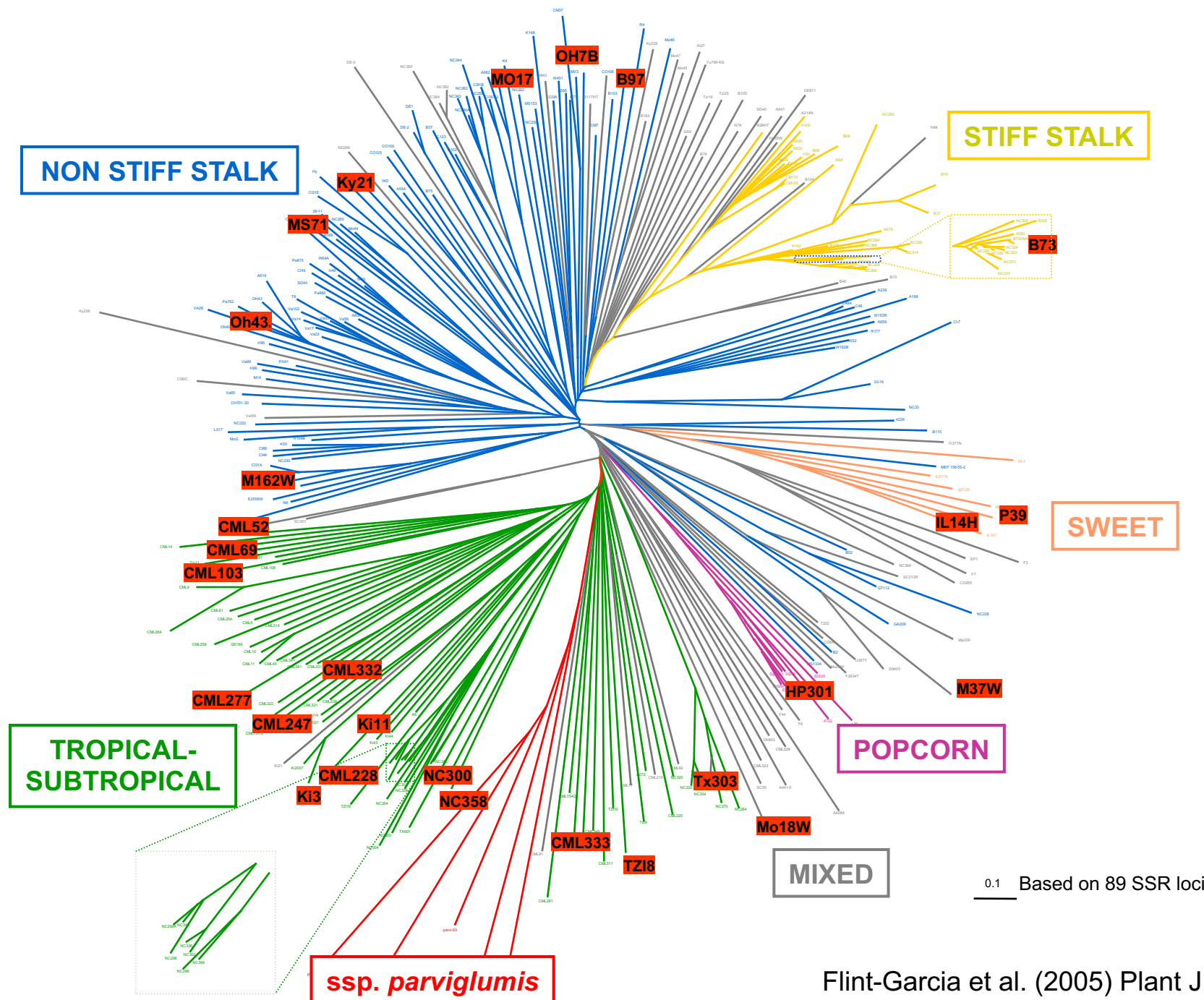
Linkage equilibrium



Linkage disequilibrium (LD)

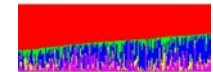
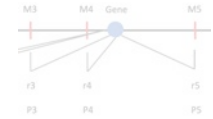


Association
half T as case, none for A
Lead to mistake if G/C is
not a marker

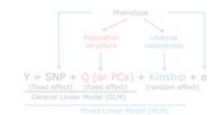


Outline

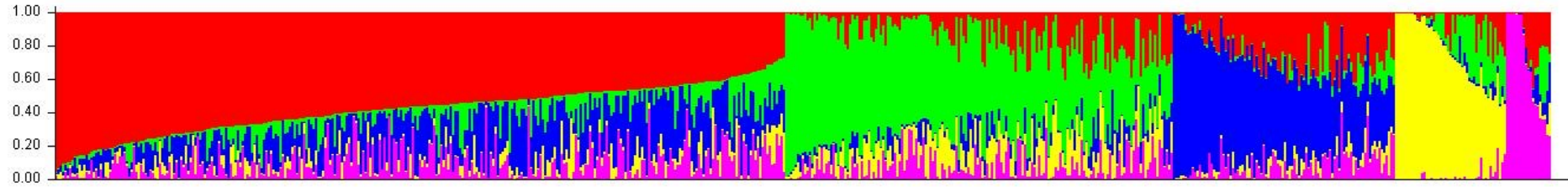
- Linkage analysis
- Association study
- Population structure and GLM
- Kinship and MLM
- BLINK



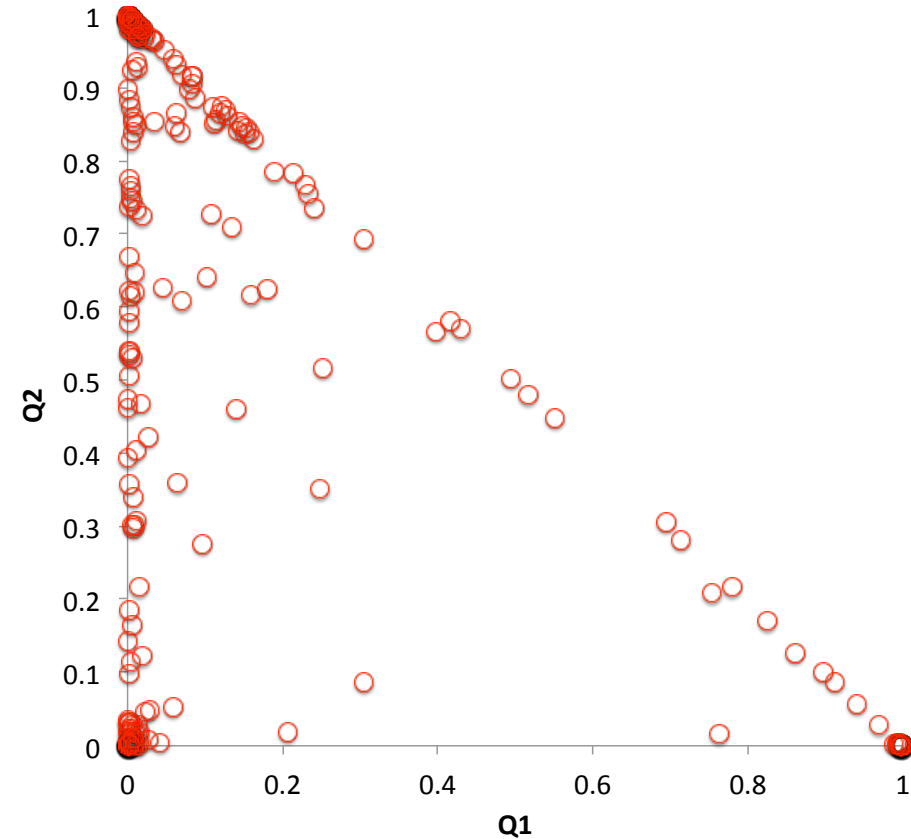
Q1	Q2	Q3
0.014	0.972	0.014
0.003	0.993	0.004
0.071	0.917	0.012
0.035	0.854	0.111
0.013	0.982	0.005



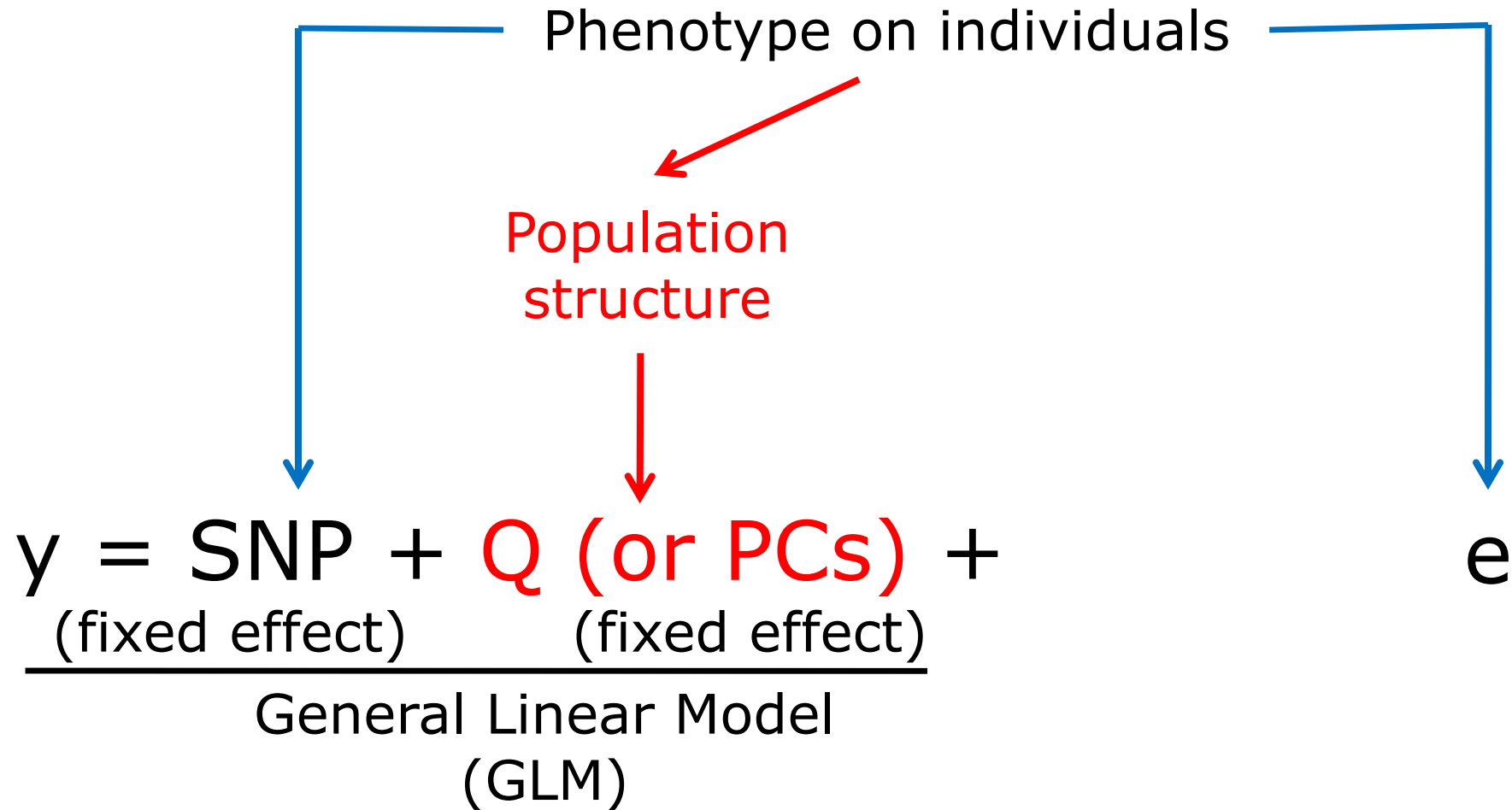
Population structure of maize



Taxa	Q1	Q2	Q3
33-16	0.014	0.972	0.014
38-11	0.003	0.993	0.004
4226	0.071	0.917	0.012
4722	0.035	0.854	0.111
A188	0.013	0.982	0.005
B73	0.999	0.001	1.10E-16
B73HTRHM	0.999	0.001	1.10E-16
B75	0.005	0.993	0.002
WD	0.014	0.97	0.016
WF9	0.005	0.994	0.001
YU796NS	0.189	0.785	0.026



GLM (Conceptual)



GLM on individuals

observation	mean	PC2	SNP	Ind1	Ind2	...	Ind9	Ind10		
	$b = [$	b_0	b_1	$b_2]$	$u = [$	u_1	u_2	$...$	u_9	$u_{10}]$
-4.709379	1	8.040247	2	1	0	...	0	0		
-5.103188	1	4.824156	2	0	1	...	0	0		
-2.782490	1	4.750749	2							
-3.835722	1	-5.773005	2							
-9.195871	1	-14.023364	2							
-3.283042	1	-7.073483	2							
-5.659523	1	8.636867	2							
-4.264048	1	-15.491325	2							
-10.486154	1	13.363734	0	0	0	...	1	0		
-3.057630	1	-2.142841	2	0	0	...	0	1		

$$y \quad [\quad 1 \quad \quad x_1 \quad \quad x_2 \quad] = X$$

Z

$$y = Xb + Zu + e$$

General linear model

$$y = b_0 + x_1 b_1 + x_2 b_2 + \dots + x_p b_p + e$$

y: observation, dependent variable

x: Explanatory/independent variables

e: Residuals/errors

$$\Delta = e_1^2 + e_2^2 + \dots + e_n^2$$

$$= e'e$$

$$= (y - Xb)'(y - Xb)$$

Optimization to minimize residual

$$\begin{aligned}\Delta &= e'e \\ &= e^2 = (y - Xb)^2\end{aligned}$$

$$\begin{aligned}\partial\Delta/\partial b &= 2X'(y - Xb) \\ &= 2X'y - 2X'Xb = 0\end{aligned}$$

$$X'Xb = X'y$$

$$b = [X'X]^{-1}[X'Y]$$

Statistical test

$$\hat{y} = X' \hat{b}$$

$$\sigma_e^2 = (y - \hat{y})'(y - \hat{y})/n$$

$$\text{Var}(\hat{b}) = [X'X]^{-1} \sigma_e^2$$

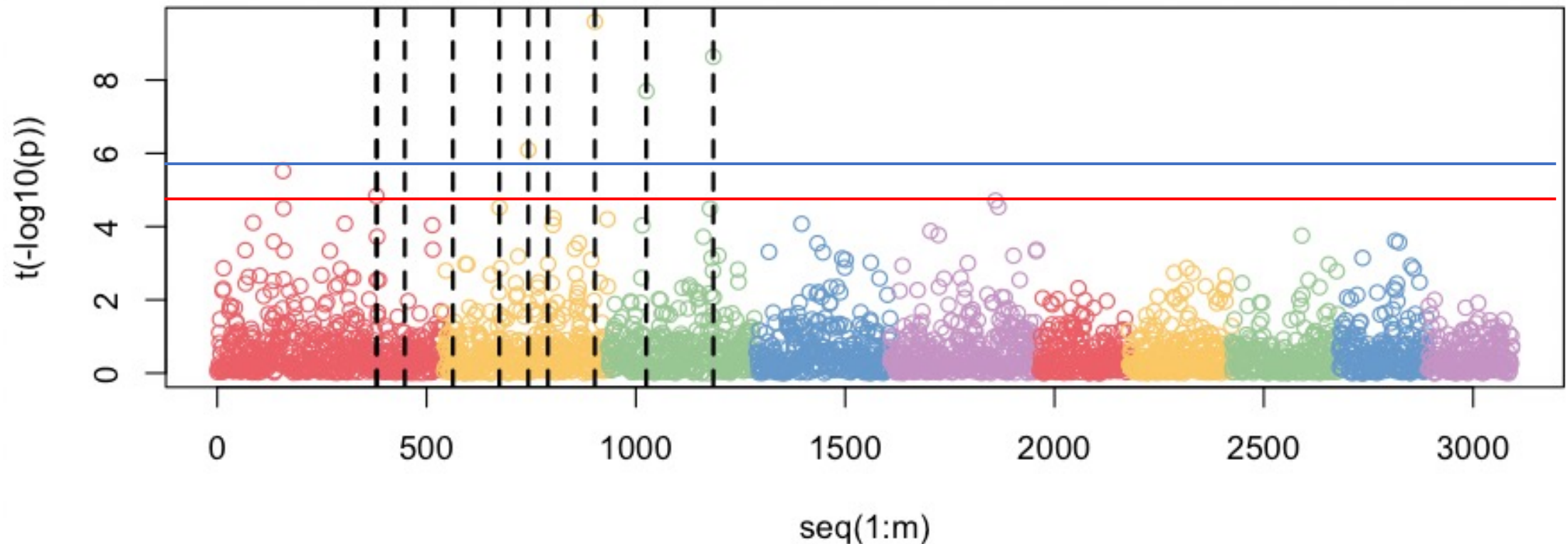
$$t = \hat{b} / \sqrt{\text{Var}(\hat{b})} \quad \sim t(n - 1)$$

QTNs On CHR 1-5, leave 6-10 empty

```
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)
source("http://zzlab.net/StaGen/2020/R/G2P.R")
source("http://zzlab.net/StaGen/2020/R/GWASbyCor.R")
X=myGD[,-1]
index1to5=myGM[,2]<6
X1to5 = X[,index1to5]
set.seed(99164)
mySim=G2P(X= X1to5,h2=.75,alpha=1,NQTN=10,distribution="norm")
p= GWASbyCor(X=X,y=mySim$y)
```

False positives

```
color.vector <- rep(c('#EC5f67', '#FAC863', '#99C794', '#6699CC', '#C594C5'),10)
m=nrow(myGM)
plot(t(-log10(p))~seq(1:m),col=color.vector[myGM[,2]])
abline(v=mySim$QTN.position, lty = 2, lwd=2, col = "black")
```



QQ plot

```
p.obs=p[!index1to5]
```

```
m2=length(p.obs)
```

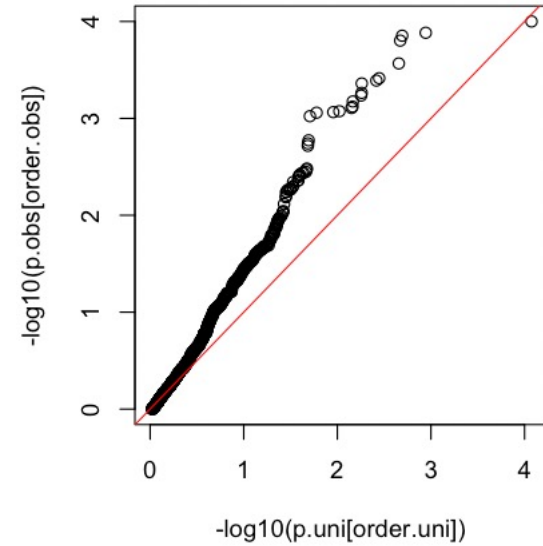
```
p.uni=runif(m2,0,1)
```

```
order.obs=order(p.obs)
```

```
order.uni=order(p.uni)
```

```
plot(-log10(p.uni[order.uni]),-log10(p.obs[order.obs]))
```

```
abline(a = 0, b = 1, col = "red")
```

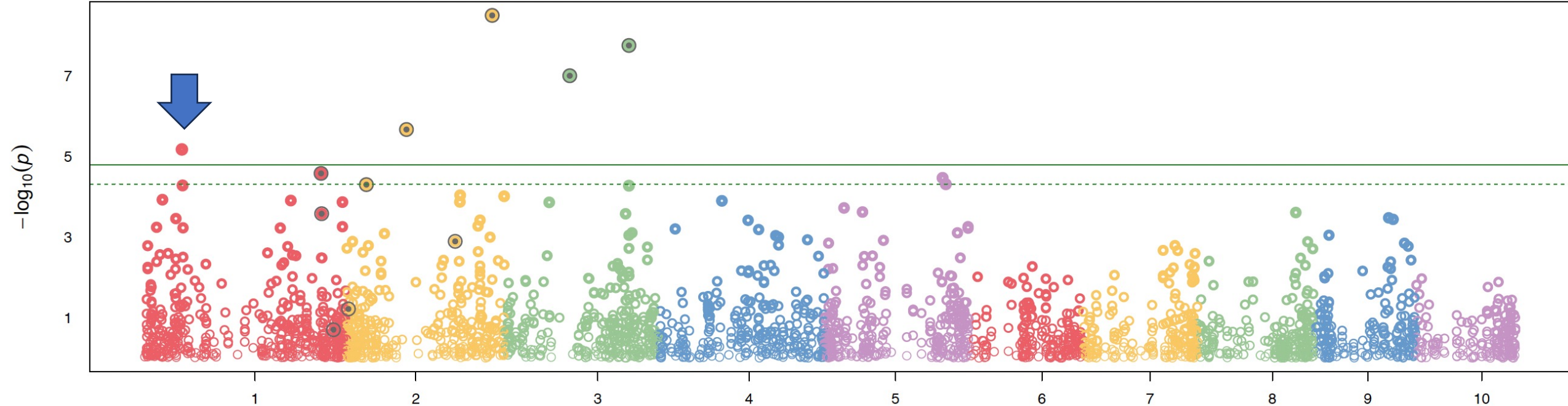


T Test



```
setwd("~/Desktop/temp")  
myY=(cbind(myGD[,1], as.data.frame(mySim$y)))  
source("http://zzlab.net/GAPIT/gapit_functions.txt")
```

```
#GWAS by GAPIT  
myGAPIT=GAPIT(  
  Y=myY,  
  GD=myGD,  
  GM=myGM,  
  QTN.position=mySim$QTN.position,  
  PCA.total=0,  
  model="GLM",  
  memo="tTest")
```



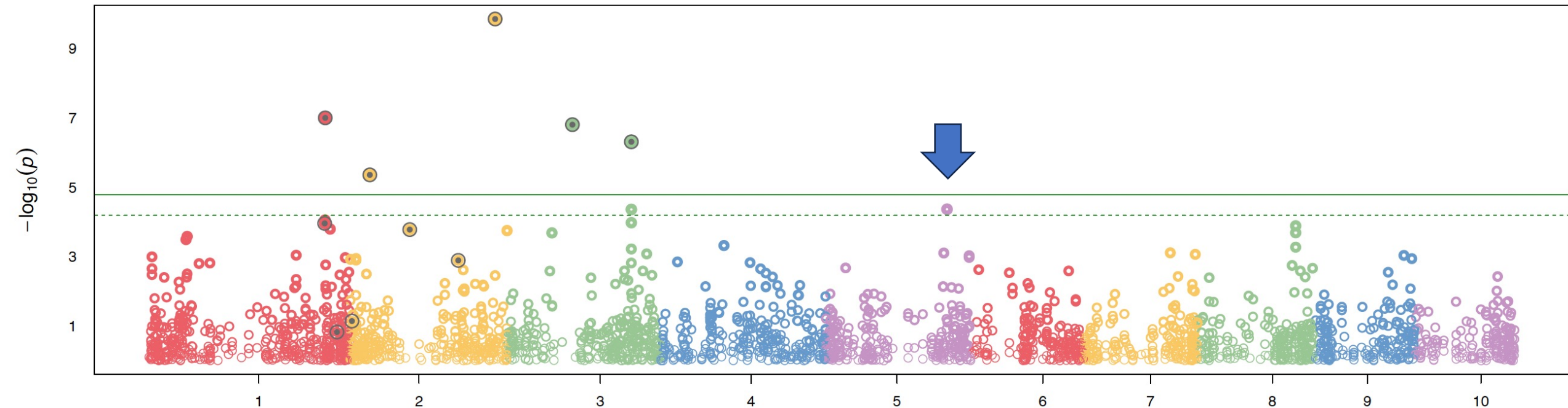
```
setwd("~/Desktop/temp")
```

```
myY=(cbind(myGD[,1], as.data.frame(mySim$y)))
```

```
source("http://zzlab.net/GAPIT/gapit_functions.txt")
```

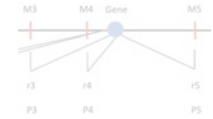
GLM

```
#GWAS by GAPIT  
myGAPIT=GAPIT(  
  Y=myY,  
  GD=myGD,  
  GM=myGM,  
  QTN.position=mySim$QTN.position,  
  PCA.total=3,  
  model="GLM",  
  memo="GLM")
```

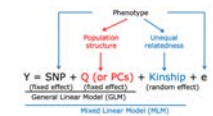


Outline

- Linkage analysis
- Association study
- Population structure and GLM
- Kinship and MLM
- BLINK

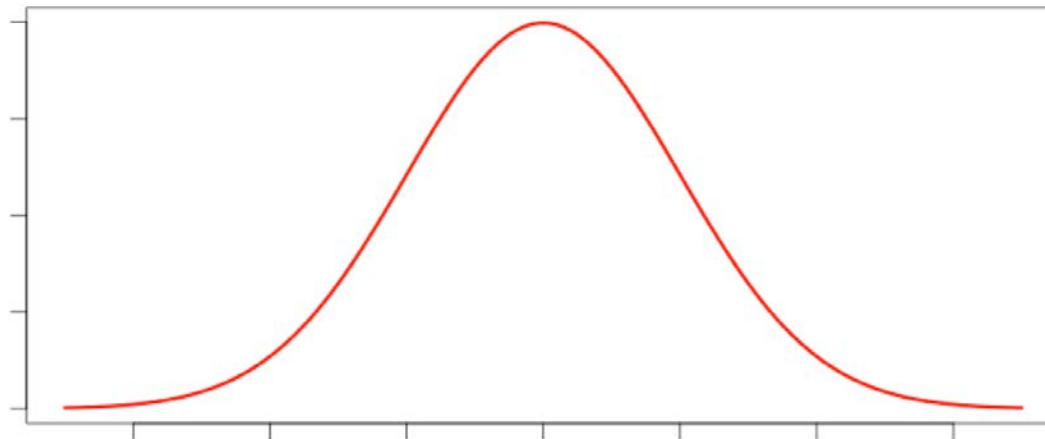


Q1	Q2	Q3
0.014	0.972	0.014
0.003	0.993	0.004
0.071	0.917	0.012
0.035	0.854	0.111
0.013	0.982	0.005



Minimize residual does not work for adding individuals' effects

- More parameters than observations
- Residuals are always zero due to over model fitting
- Needs new rules

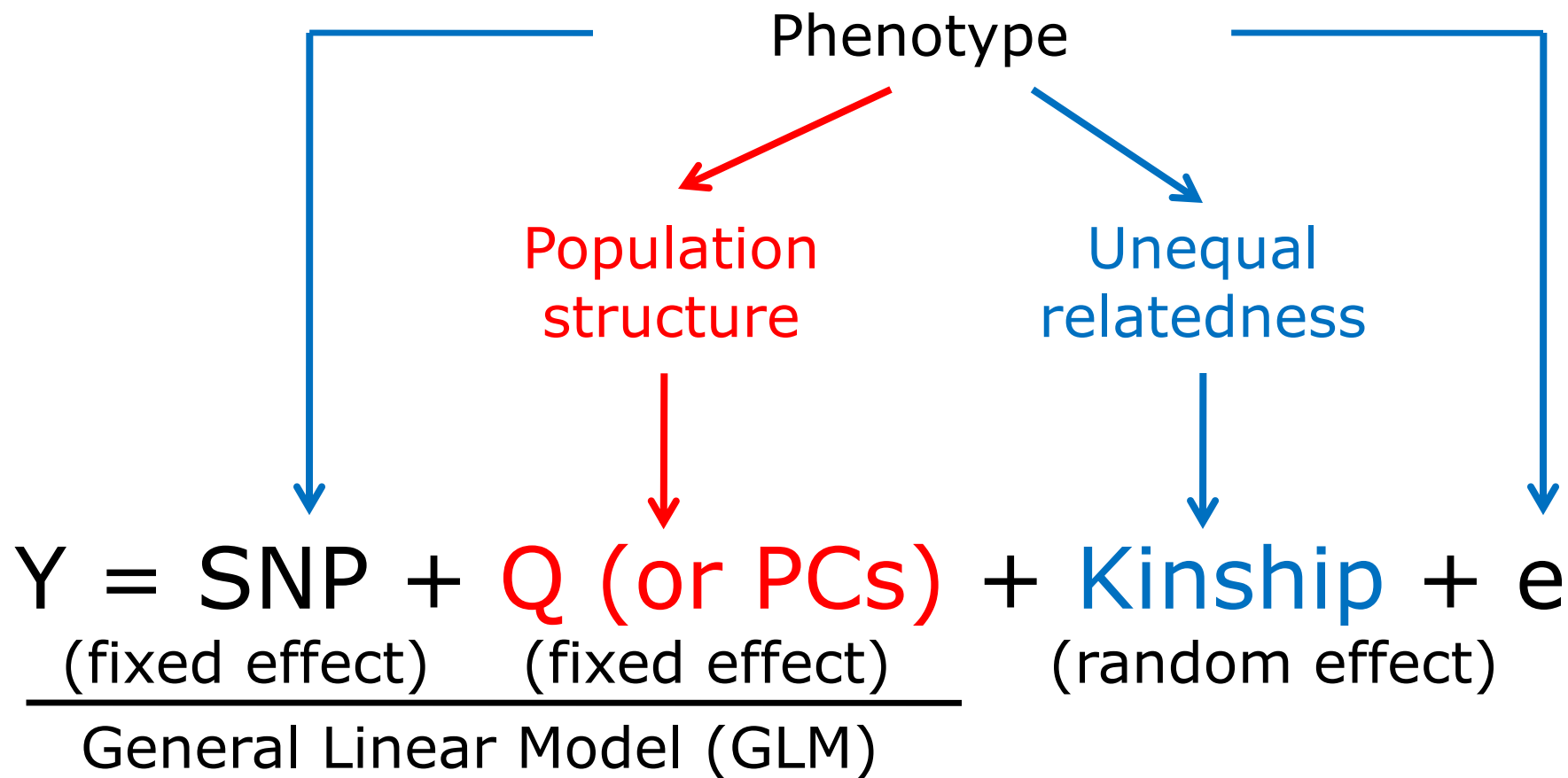


X

New rules

- Free individuals effects
- Only regulate the features for the population where they from: means (0) and variances
- Optimize variances according to kinship and observation to maximize the likelihood

MLM for GWAS



Mixed Linear Model (MLM)

(Yu et al. 2005, Nature Genetics)

Mixed Linear Model (MLM)

$$y = Xb + Zu + e$$

$$\text{Var}(y) = V = \text{Var}(u) + \text{Var}(e)$$

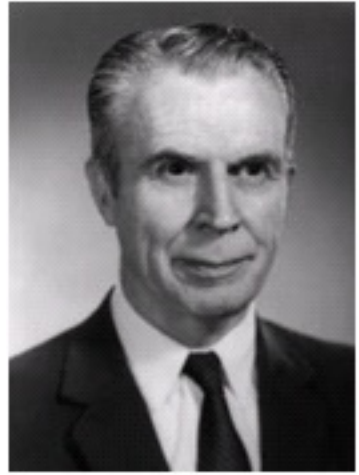
$$\text{Var}(u) = G = 2K\sigma_a^2$$

$$\text{Var}(e) = R = I\sigma_e^2$$

u prediction: Best Linear Unbiased Prediction, BLUP)

b prediction: Best Linear Unbiased Estimate, BLUE)

Mixed Model Equation



C.R. Henderson

$$y = Xb + Zu + e$$

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + \frac{\sigma_e^2}{\sigma_a^2} A^{-1} \end{bmatrix} \begin{bmatrix} b \\ u \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

$$\begin{bmatrix} b \\ u \end{bmatrix} = \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + \frac{\sigma_e^2}{\sigma_a^2} A^{-1} \end{bmatrix}^{-1} \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

$$\text{Var}\left(\begin{bmatrix} b \\ u \end{bmatrix}\right) = \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + \frac{\sigma_e^2}{\sigma_a^2} A^{-1} \end{bmatrix}^{-1} \sigma_e^2$$

SPAGeDi

Hardy OJ, Vekemans X (2002) SPAGeDi: a versatile computer program to analyse spatial genetic structure at the individual or population levels. *Molecular Ecology Notes* 2: 618-620.

- Kinship coefficient
 - Loiselle et al. (1995)
 - Ritland (1996)
- Relationship coefficient
 - Queller & Goodnight (1989)
 - Hardy & Vekemans (1999)
 - Lynch & Ritland (1999)
 - Wang (2002);
- Genetic distance: Rousset (2000)



Identical by status

	AA	AT	TT
AA	1	.5	0
AT	.5	.5	.5
TT	0	.5	1

Proportion of shared alleles

	-1	0	1
-1	1	0	-1
0	0	0	0
1	-1	0	1

Genotype coding

Identical by status

Efficient algorithm

- M: n individual by m SNPs
- M: -1, 0 and 1
- p_i : frequency of 2nd allele for SNP i
- P: Column of i is $2(p_i - 0.5)$
- $Z = M - P$

$$G = \frac{ZZ'}{2 \sum p_i (1 - p_i)}$$

J. Dairy Sci. 2008. 91 (11) 4414-4423. Efficient Methods to Compute Genomic Predictions P. M. VanRaden



Paul VanRaden: Image Number K7168-6

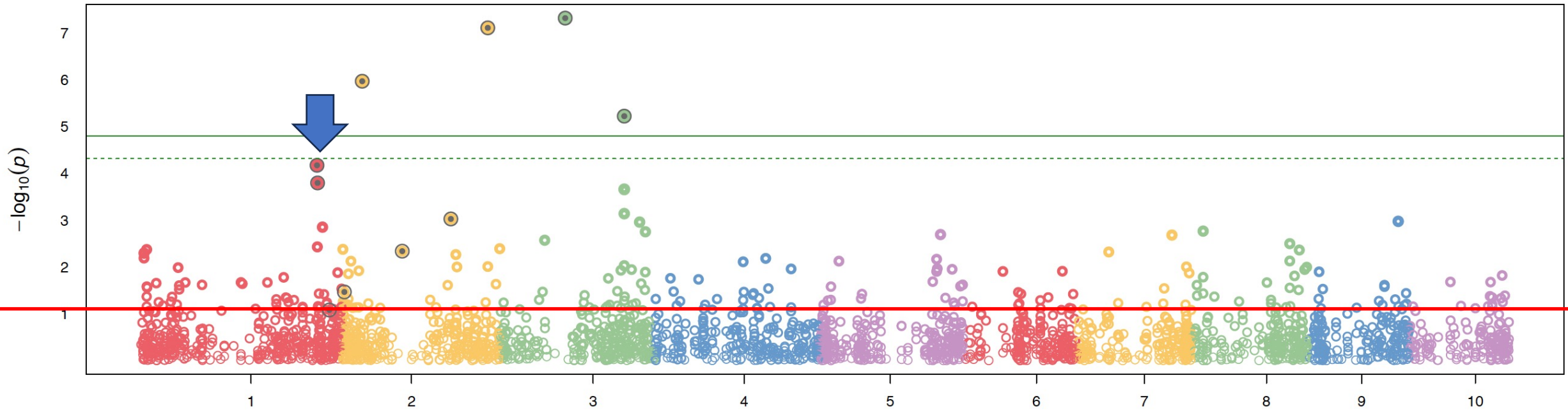
```
setwd("~/Desktop/temp")
```

```
myY=(cbind(myGD[,1], as.data.frame(mySim$y)))
```

```
source("http://zzlab.net/GAPIT/gapit_functions.txt")
```

MLM

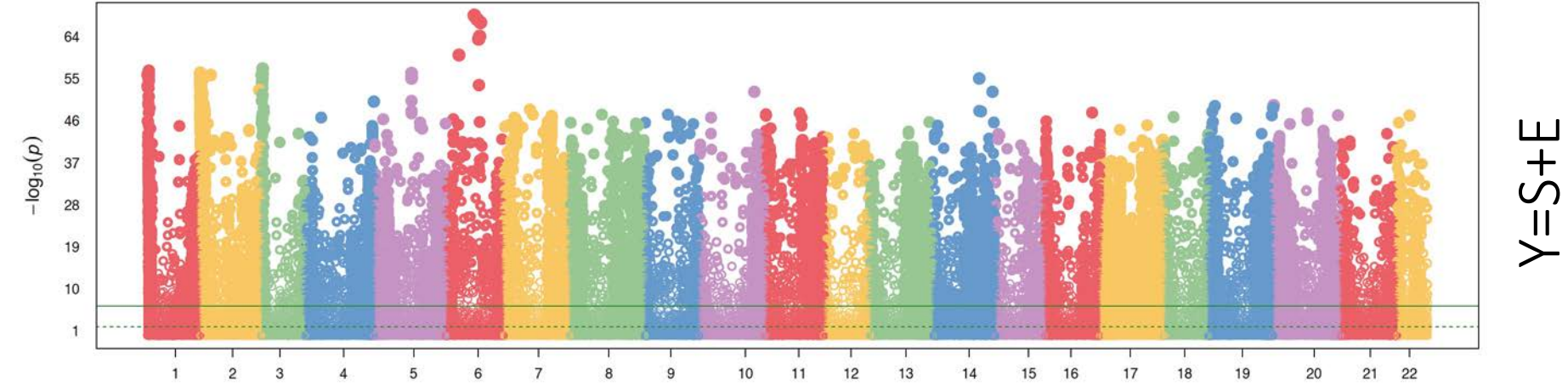
```
#GWAS by GAPIT  
myGAPIT=GAPIT(  
  Y=myY,  
  GD=myGD,  
  GM=myGM,  
  QTN.position=mySim$QTN.position,  
  PCA.total=3,  
  model="MLM",  
  memo="MLM_3PC")
```



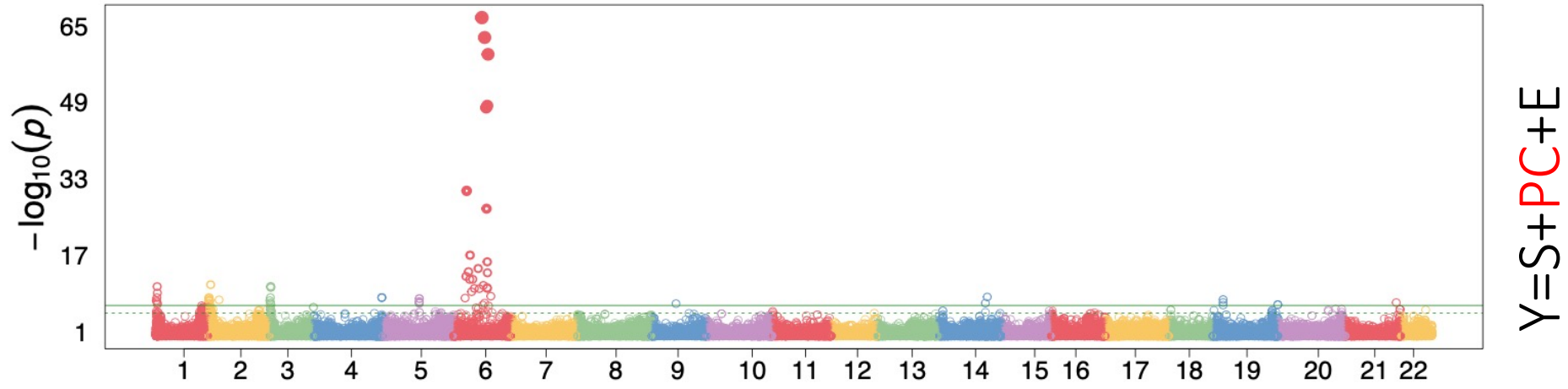
Physical mapping of the wheat genes in low-recombination regions:
radiation hybrid mapping of the C-locus, Kajla and et al., TAG, 2023

Club wheat (CW) with compact spike architecture

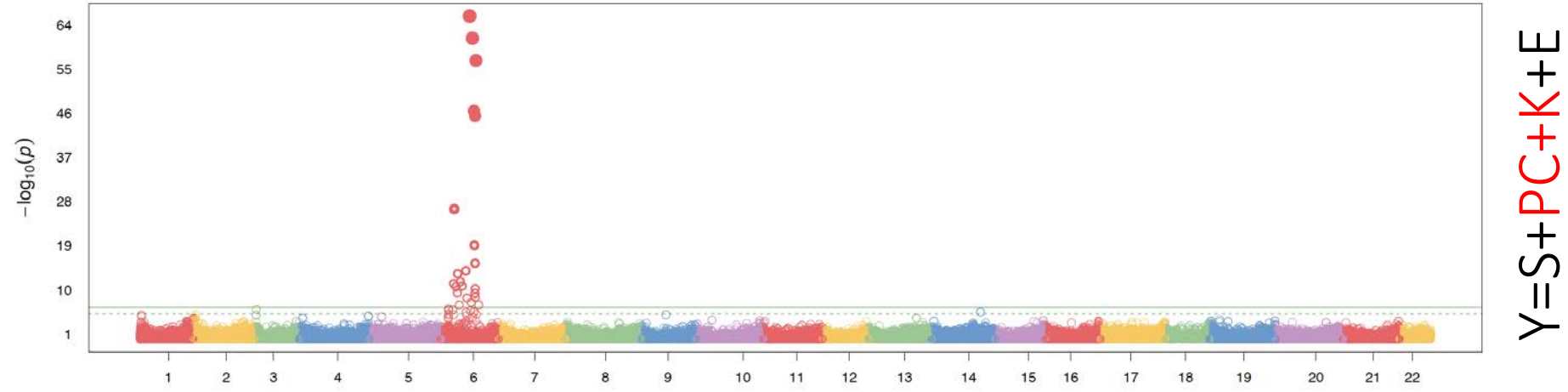




$Y=S+E$



$Y=S+PC+E$



$Y=S+PC+K+E$

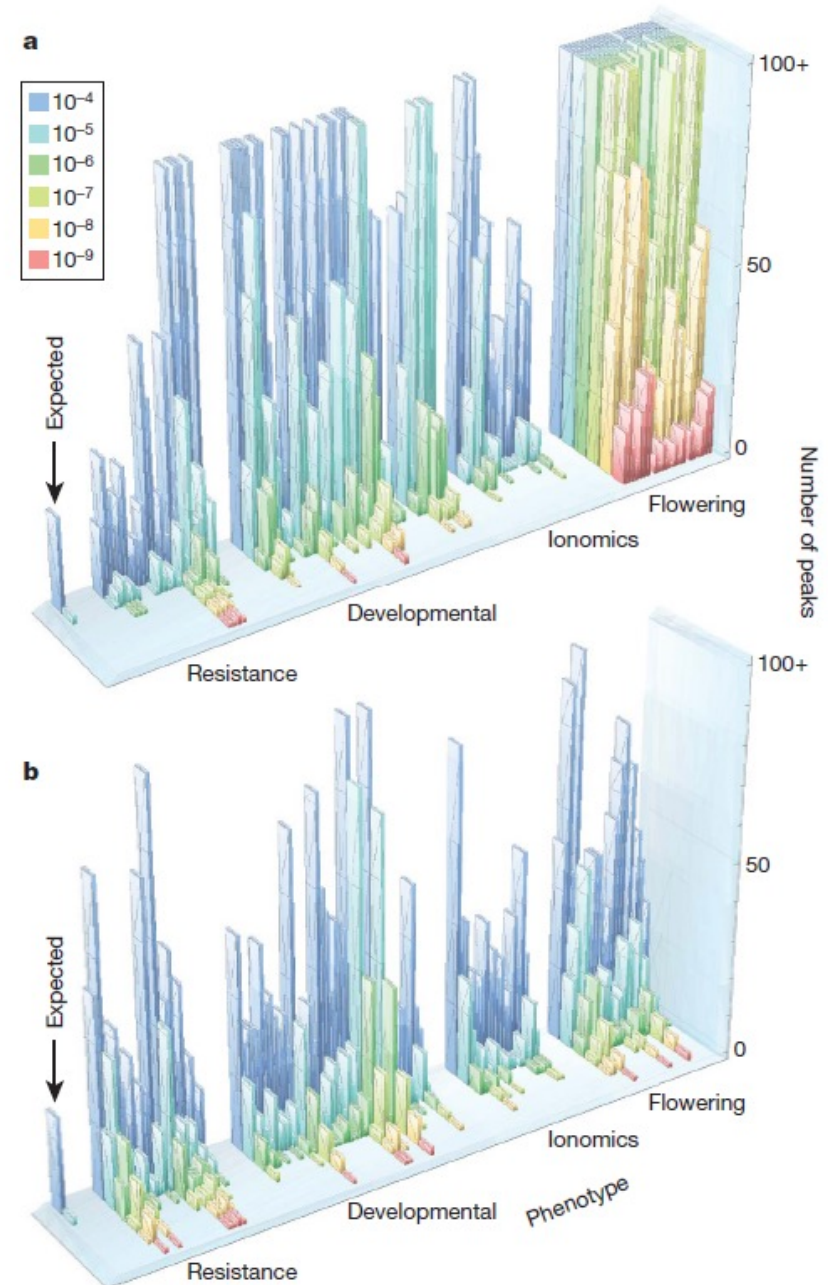
Atwell et al Nature 2010

a, No correction test

b, Correction with MLM



Magnus Nordborg



GWAS does not work for traits associated with structure

Queen + King





Flour>Water>Yeast>Salt



EMMAX



EMMA



PC+K



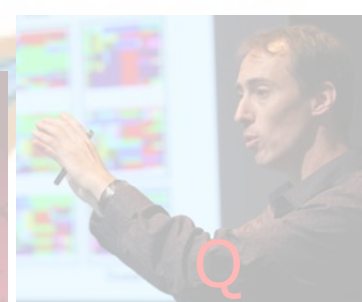
PC



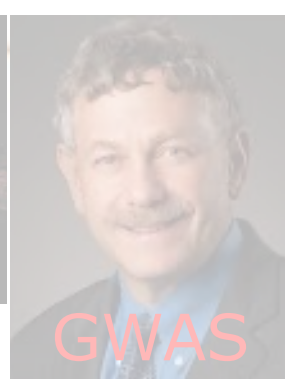
CMLM



Q+K



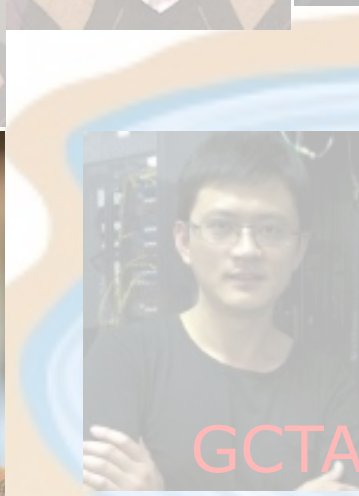
Q



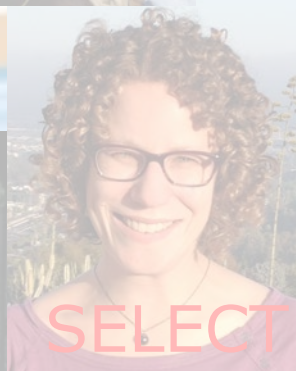
GWAS



P3D



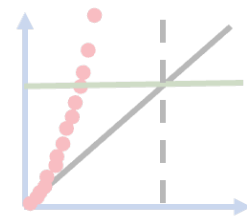
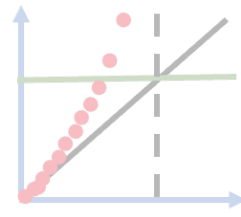
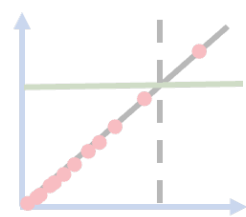
GCTA



SELECT



MLMM



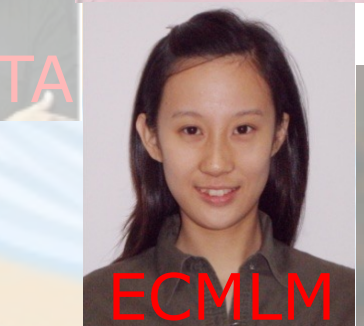
GWAS Stream



FST-LMM



GEMMA



ECMLM



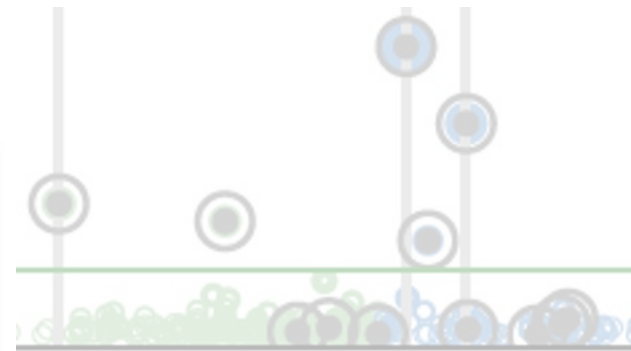
SUPER



FarmCPU

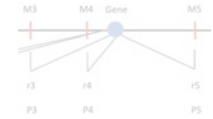


BLINK

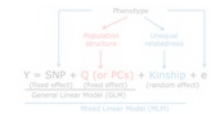


Outline

- Linkage analysis
- Association study
- Population structure and GLM
- Kinship and MLM
- **BLINK**

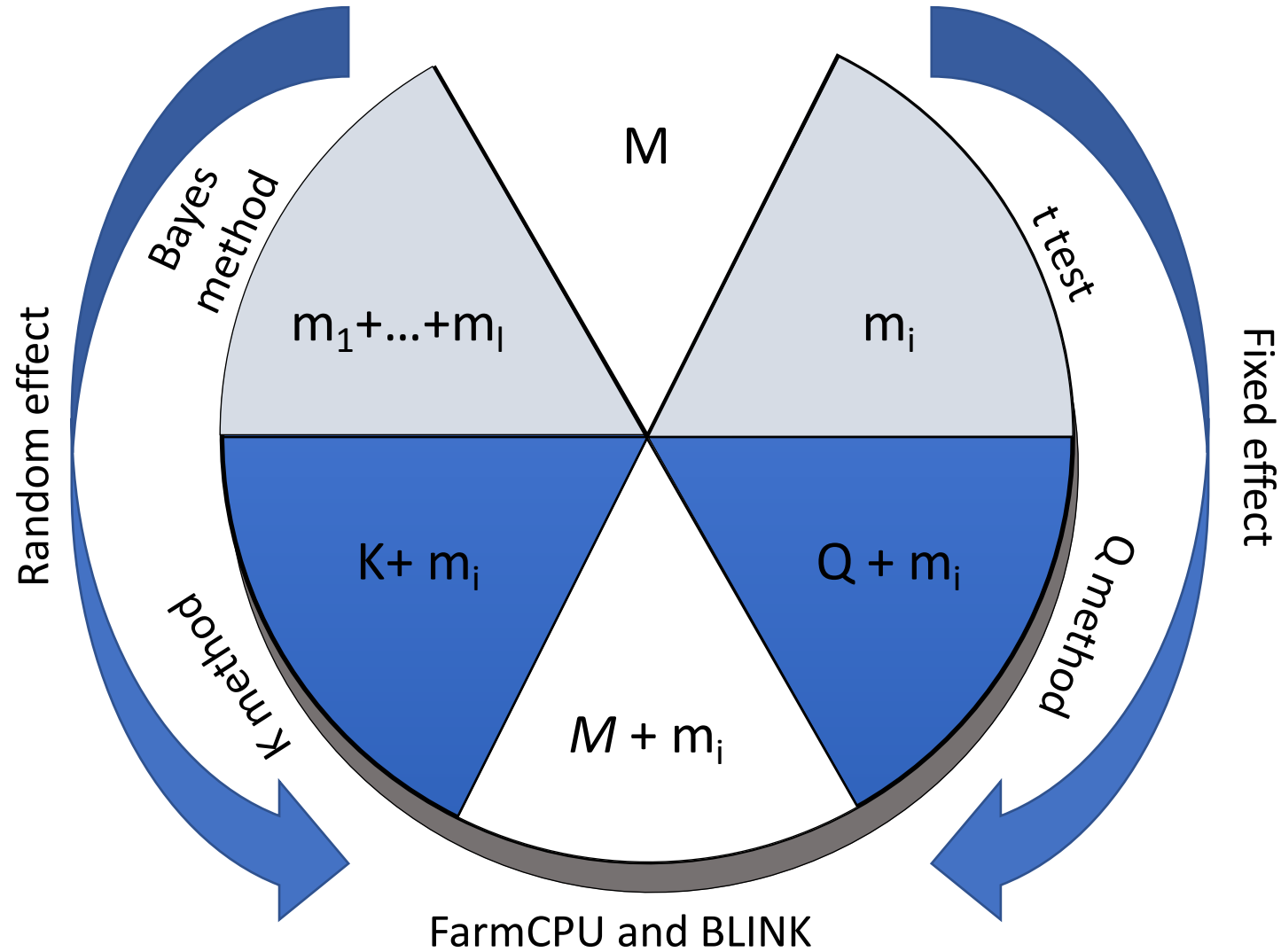


Q1	Q2	Q3
0.014	0.972	0.014
0.003	0.993	0.004
0.071	0.917	0.012
0.035	0.854	0.111
0.013	0.982	0.005



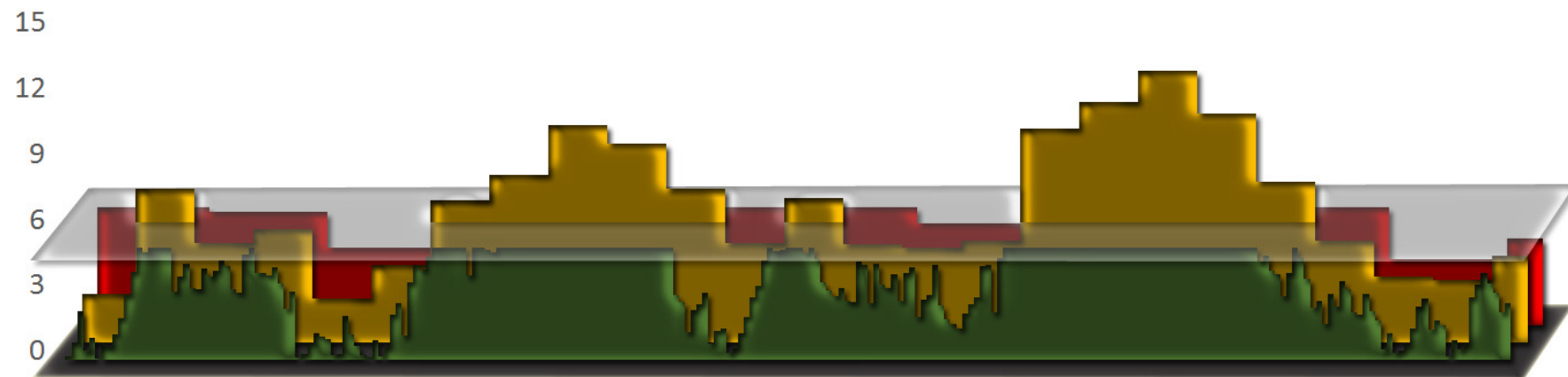
GWAS methods

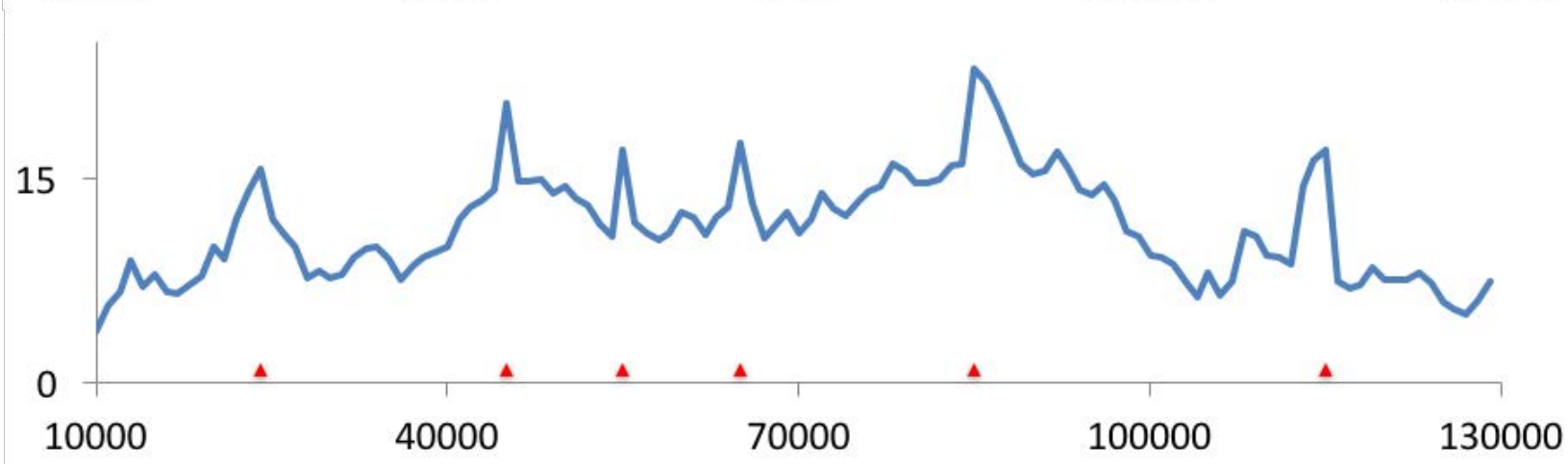
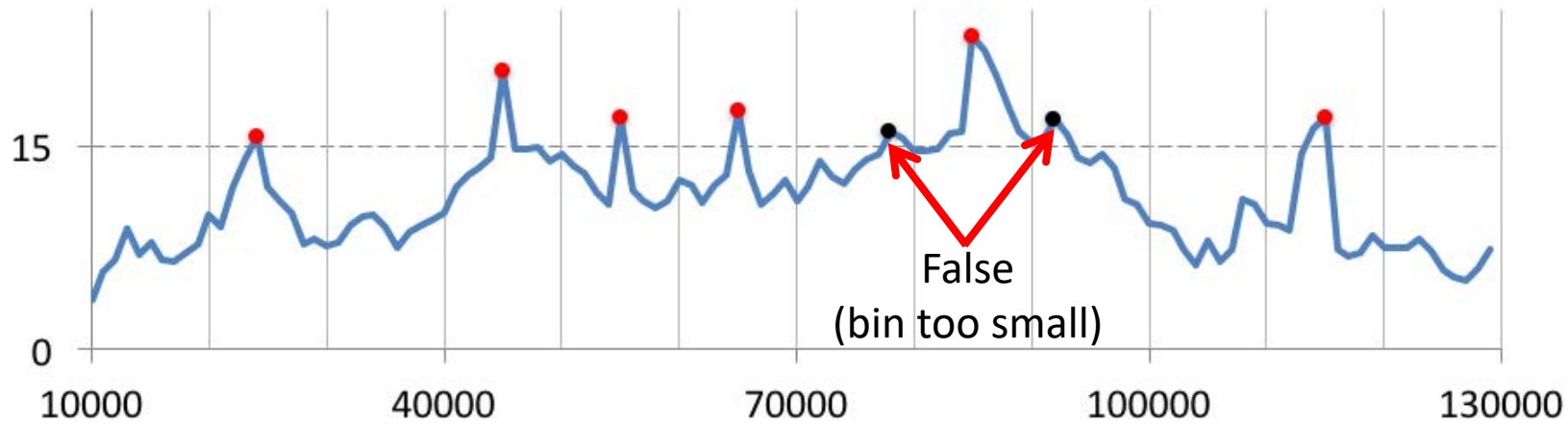
Known M



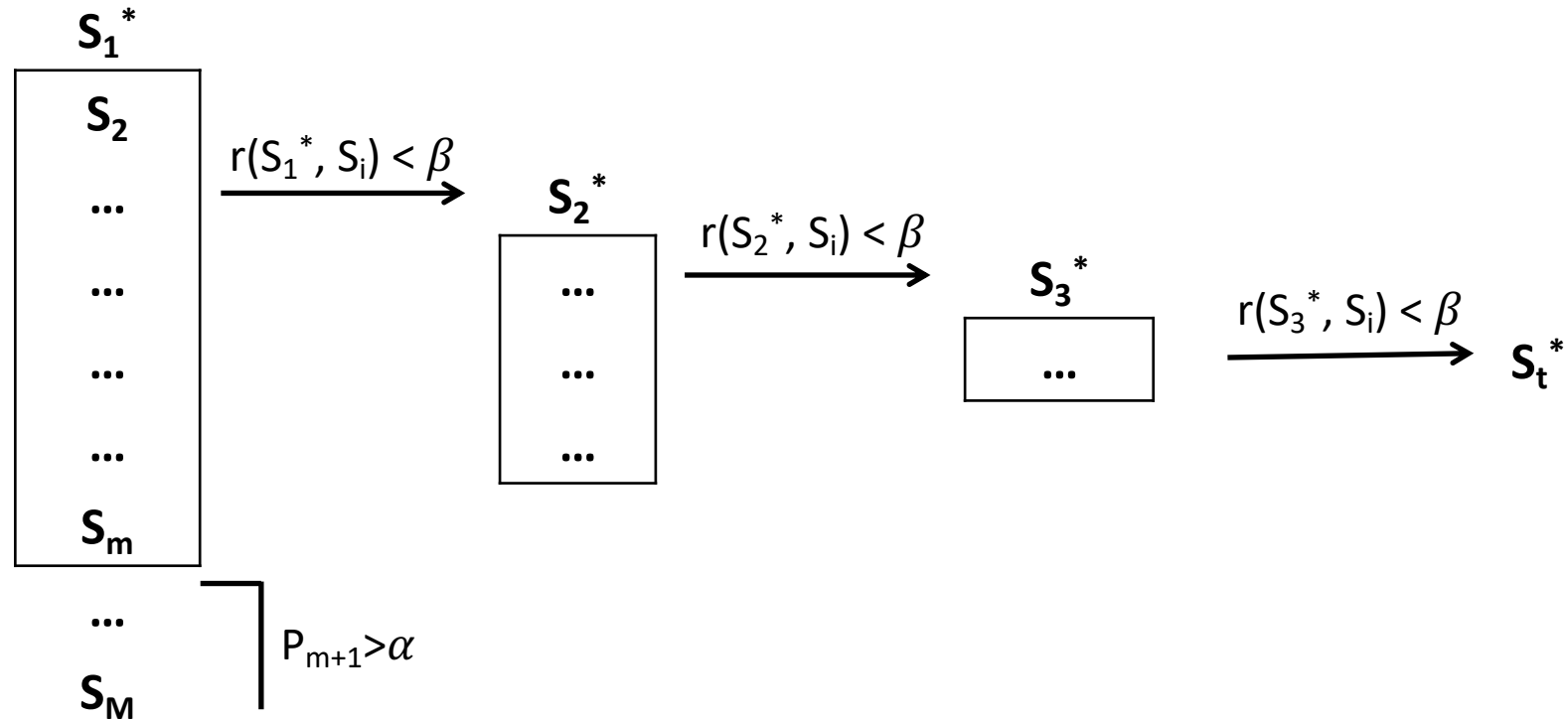
m: marker, M: Mutations, M : Estimated mutations

Bin approach





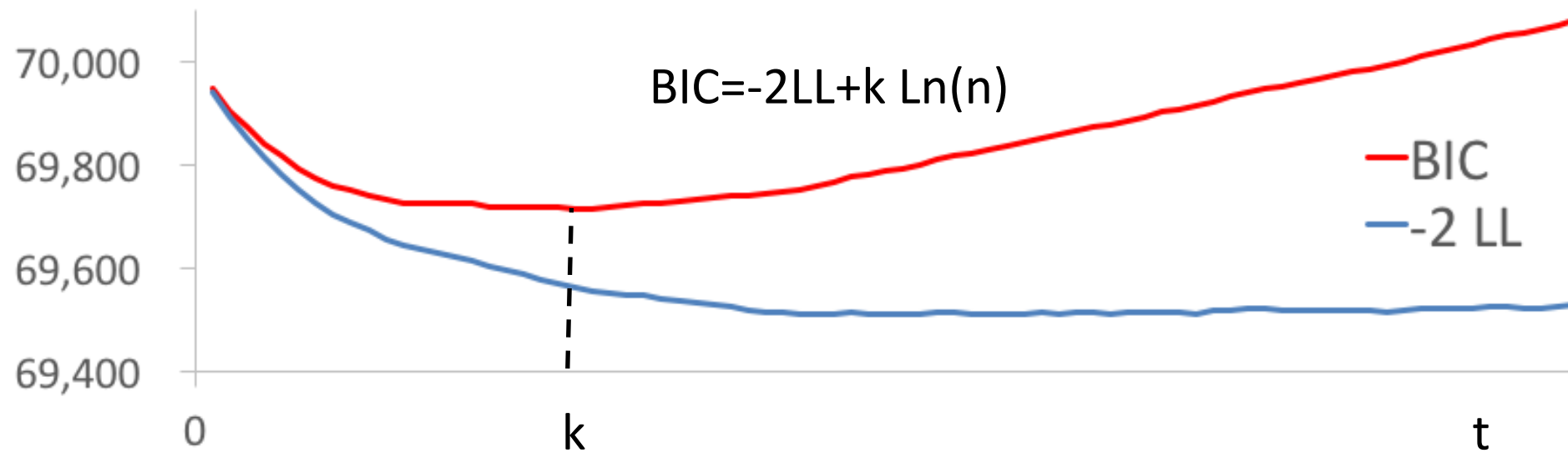
Elimination of markers with LD

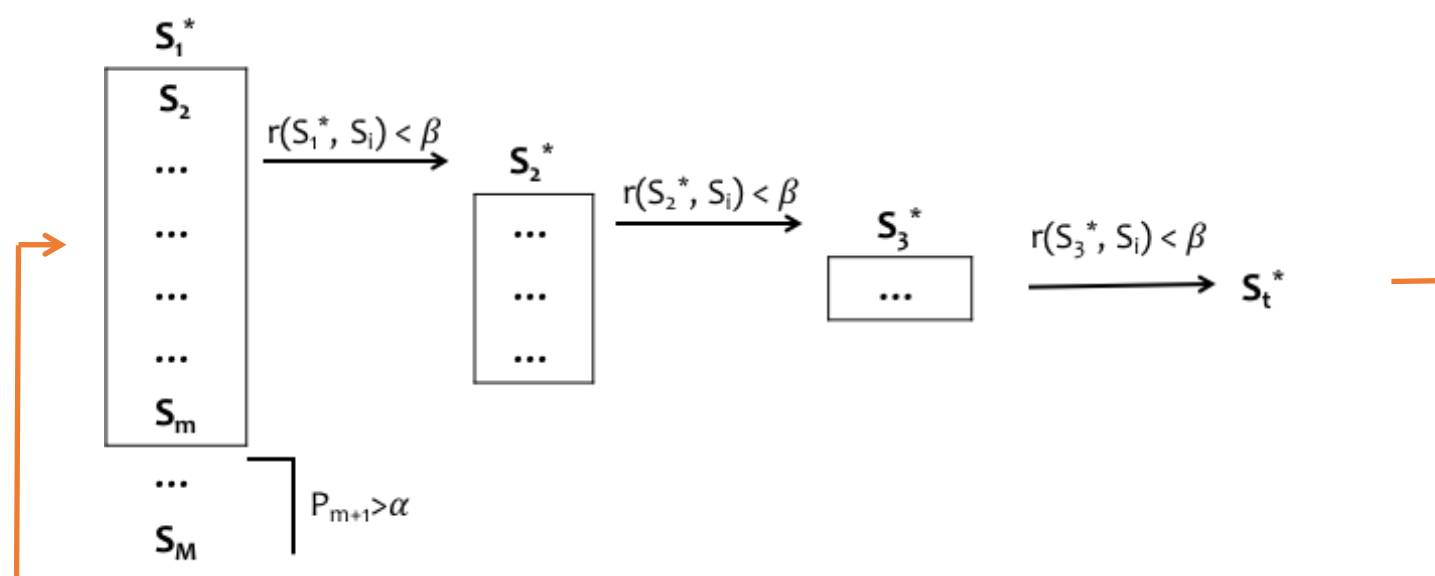


$$y = s_i + S_1^* + S_2^* + S_3^* + \dots + S_k^* + e, \text{ where } i = 1 \text{ to } M$$

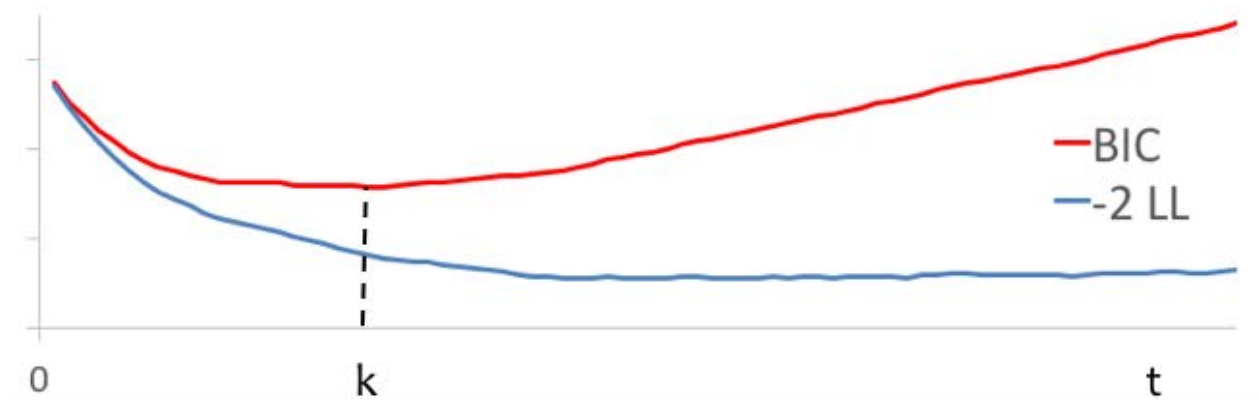
Bayesian information criterion

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





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

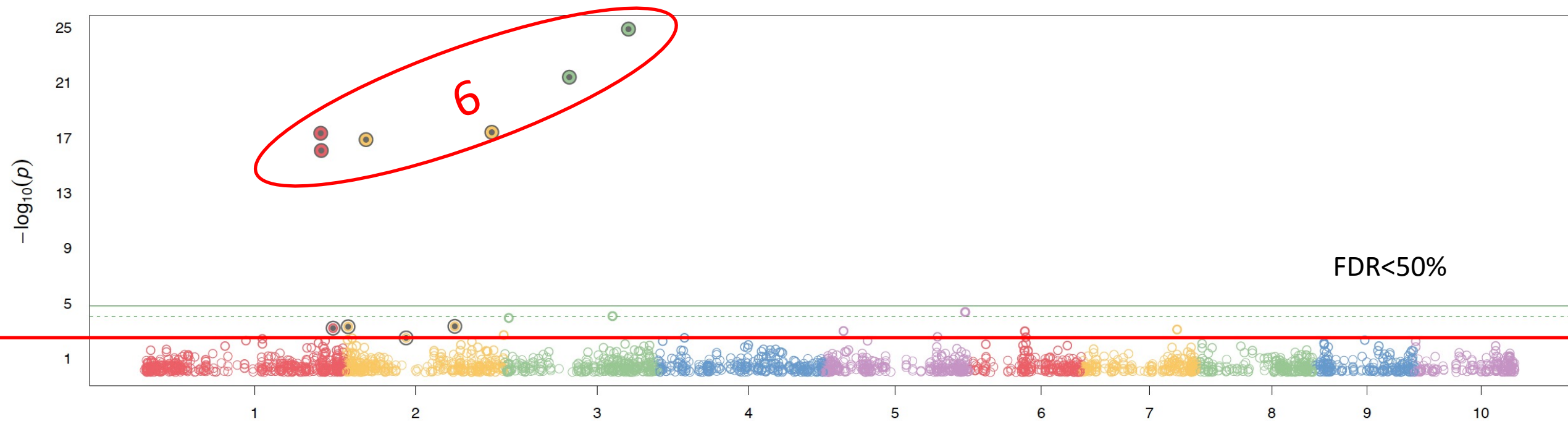
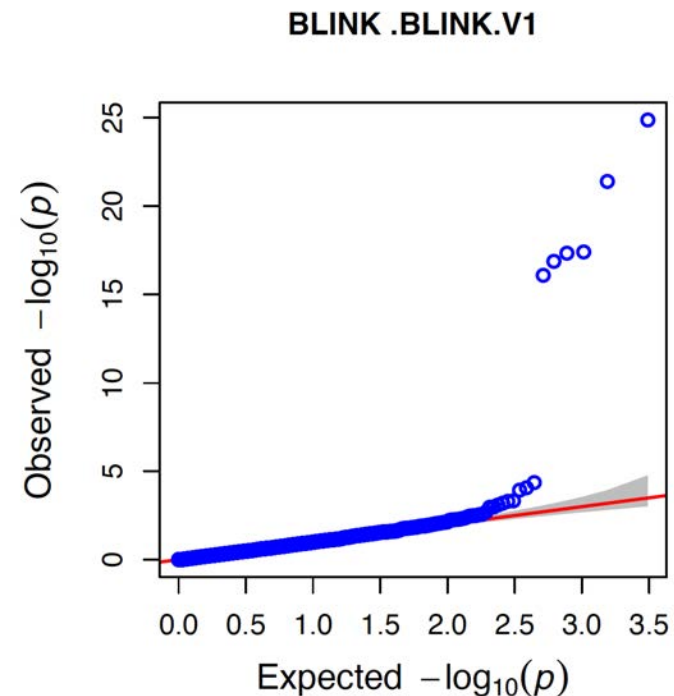


$$y = s_i + S_1^* + S_2^* + S_3^* + \dots + S_k^* + e, \text{ where } i = 1 \text{ to } M$$

BLINK (R in GAPIT)

```
myGAPIT=GAPIT(  
  Y=mySim$Y,  
  GD=myGD,  
  GM=myGM,  
  QTN.position=mySim$QTN.position,  
  PCA.total=3,  
  model="Blink")
```

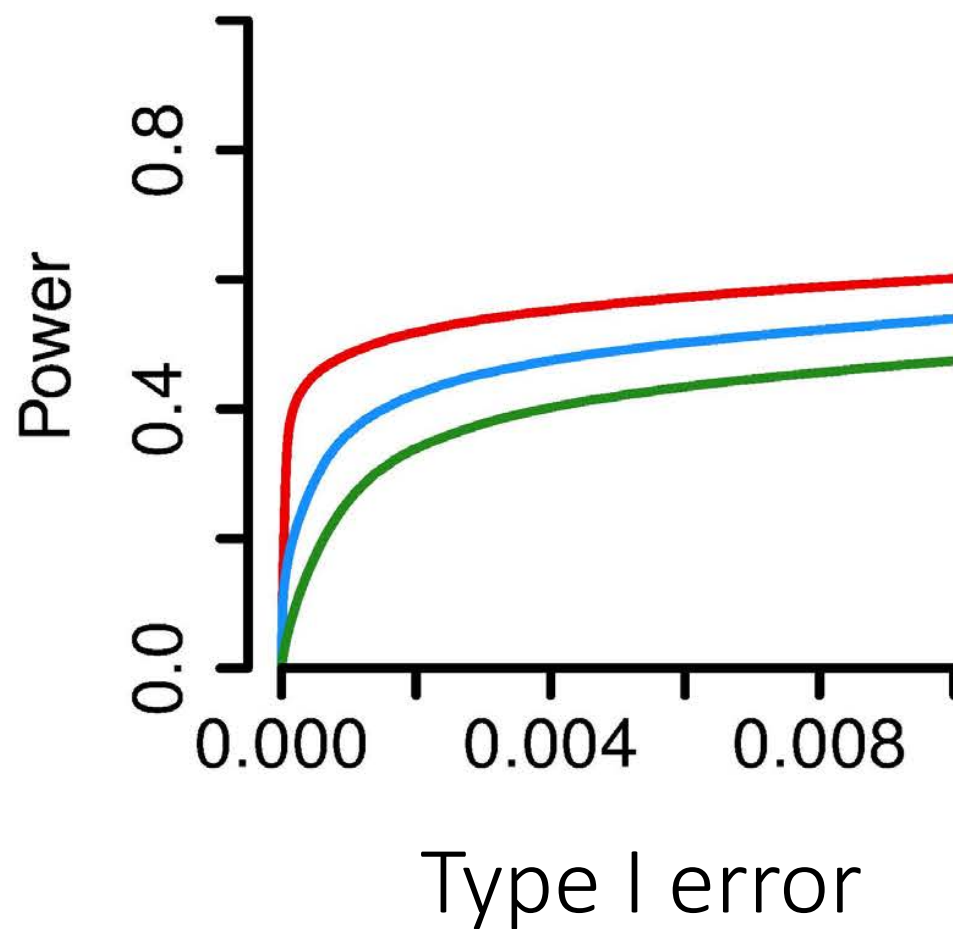
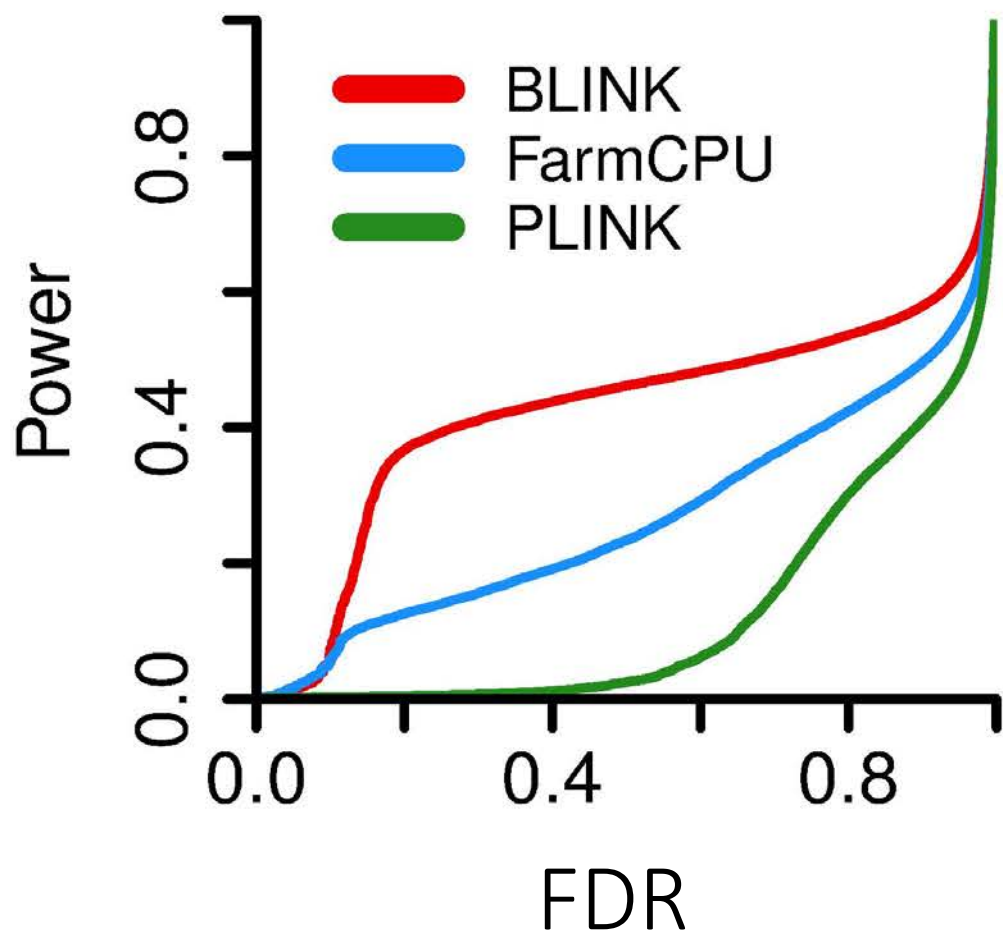
BLINK .BLINK.V1



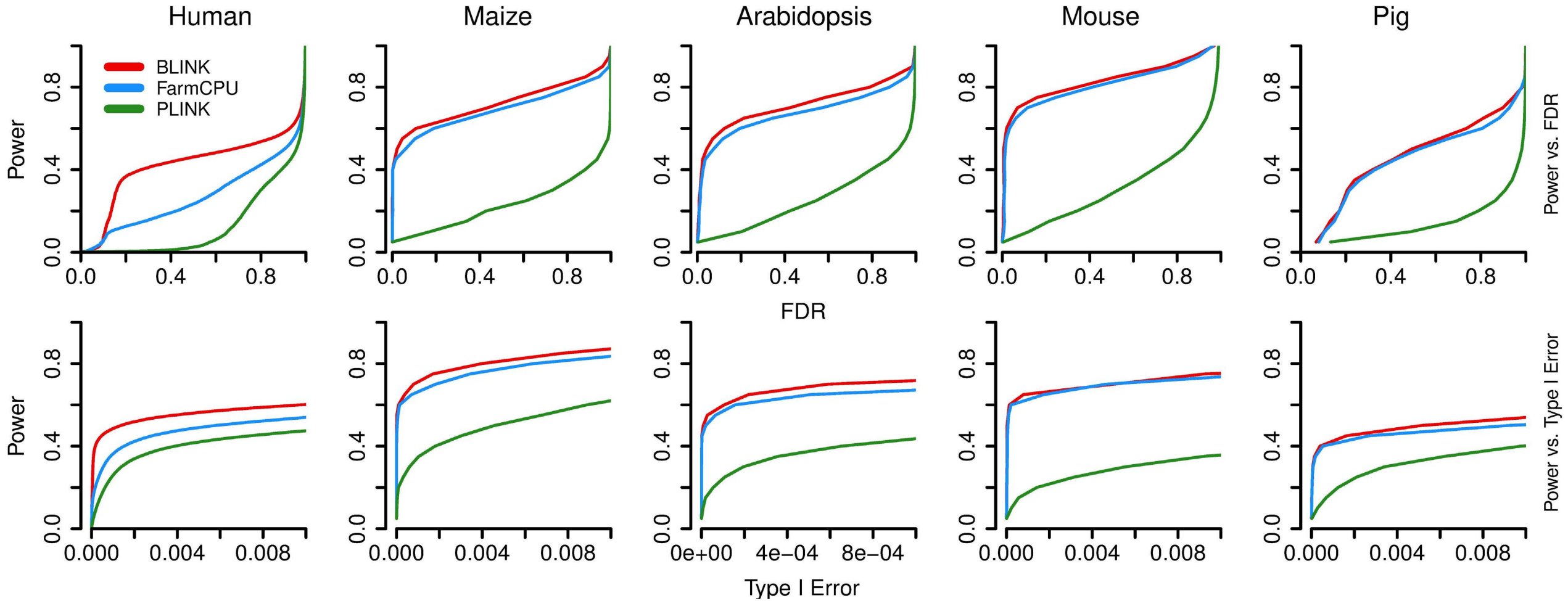
Simulation study with human data



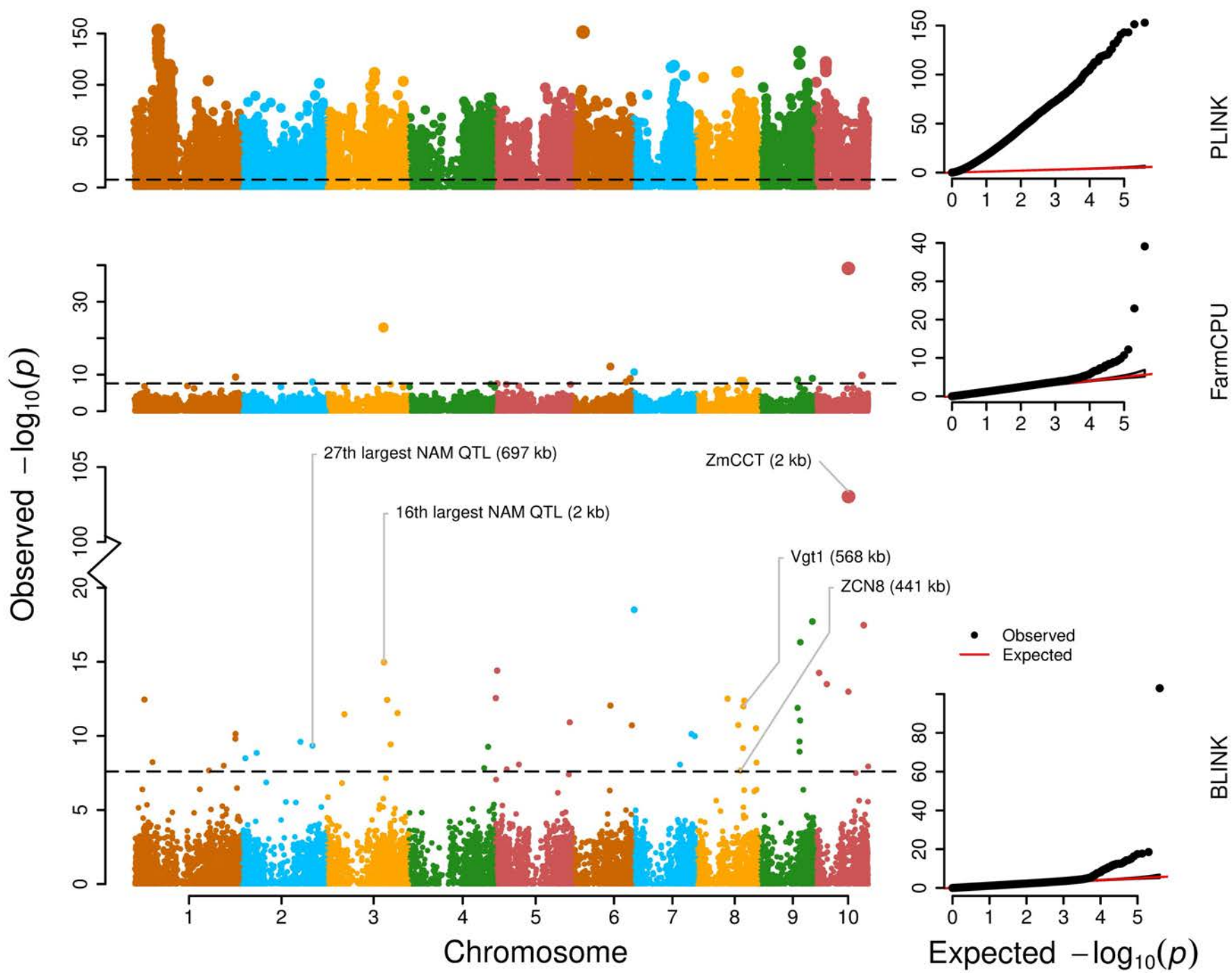
Meng Huang



Same trend across species

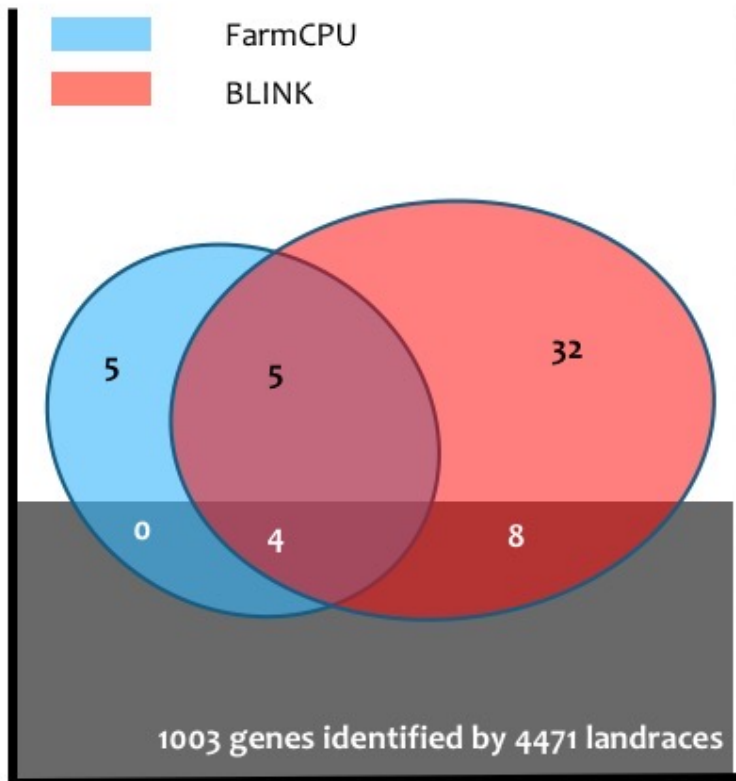


Application in Maize

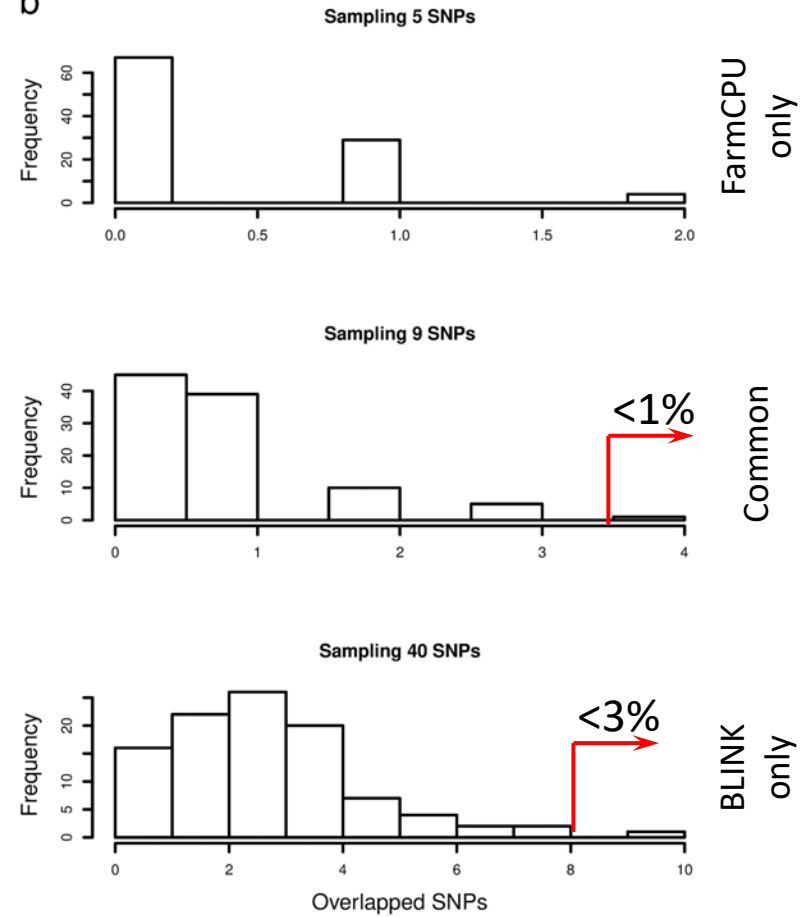


Enrichment

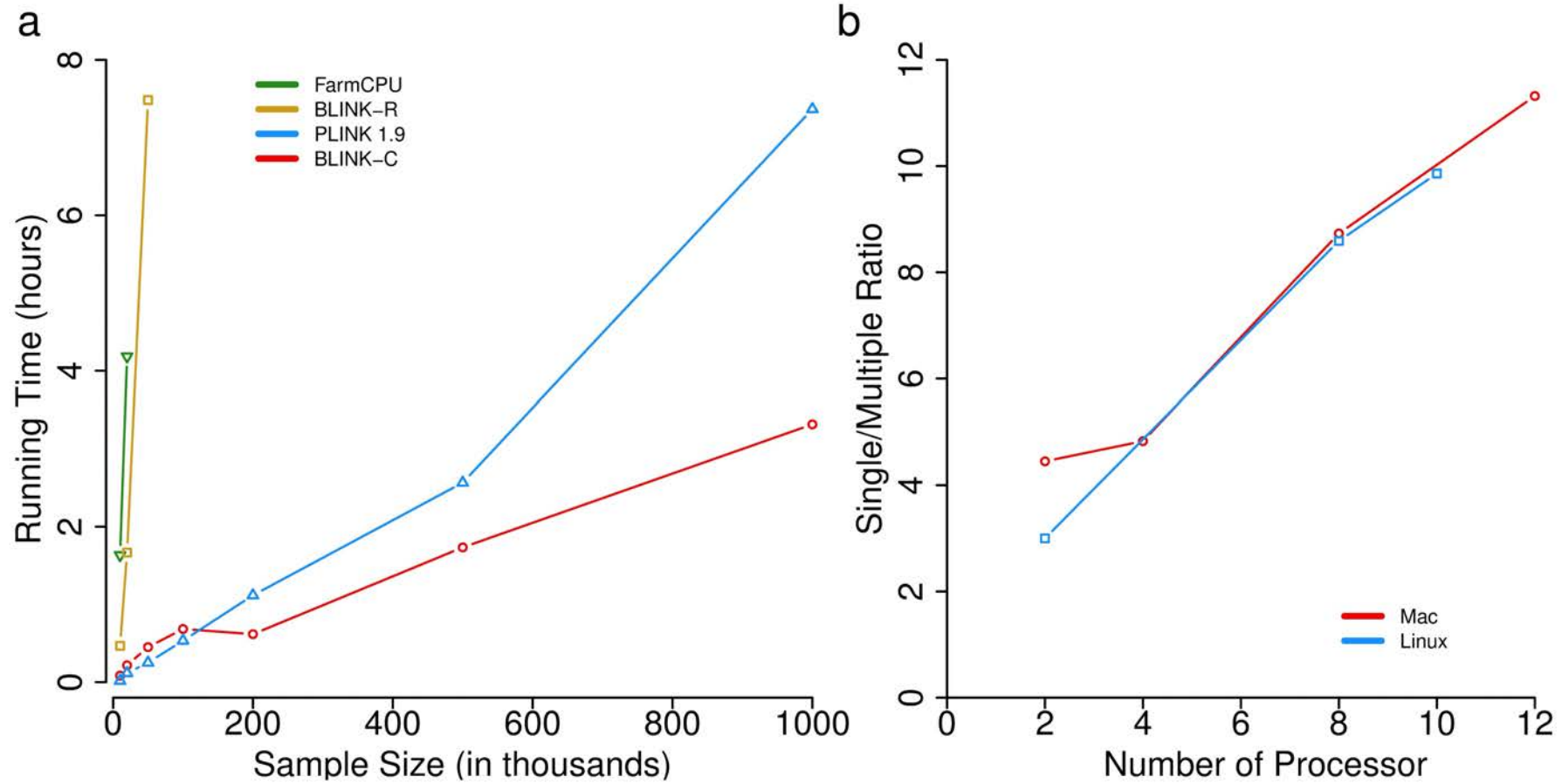
a

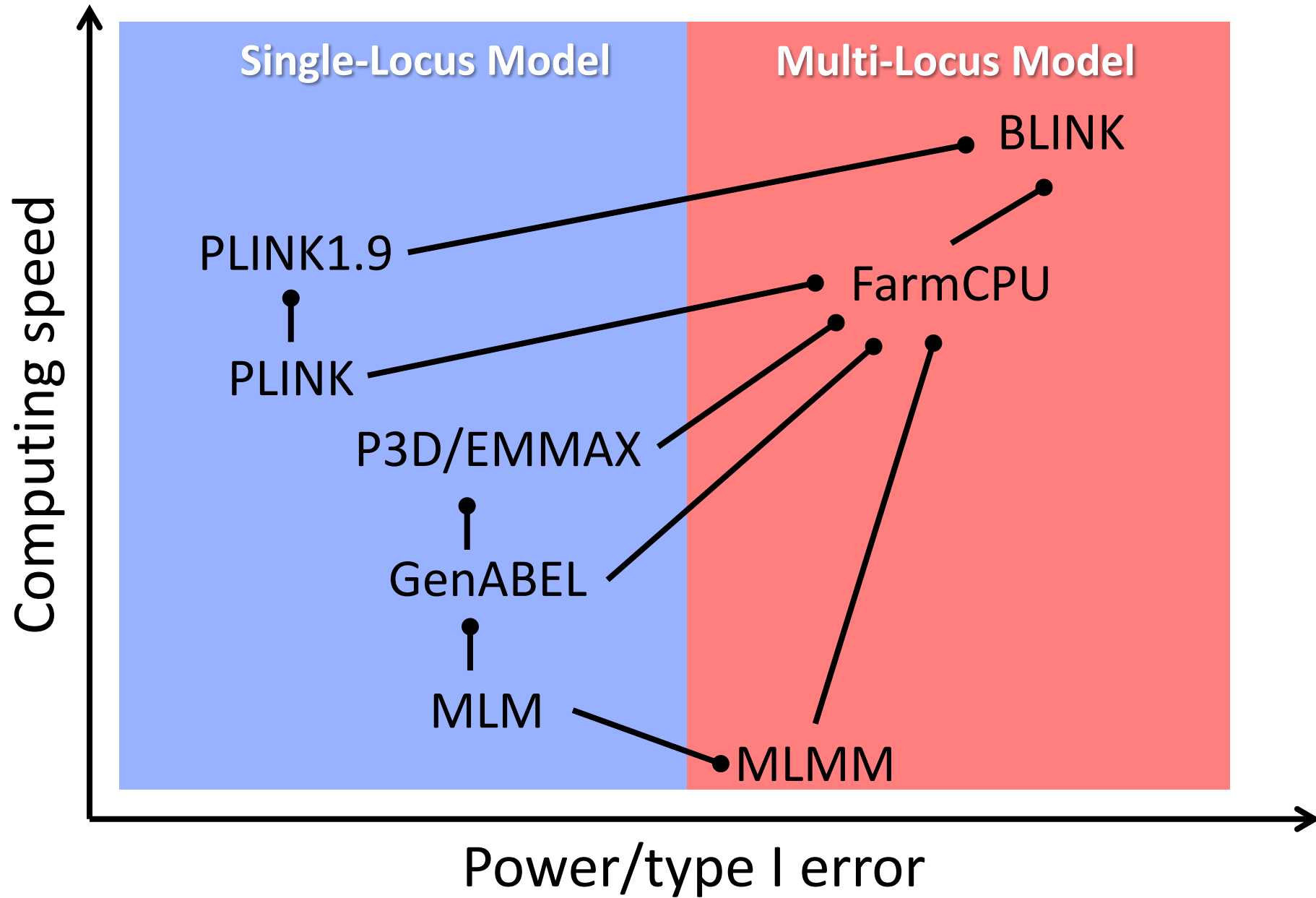


b


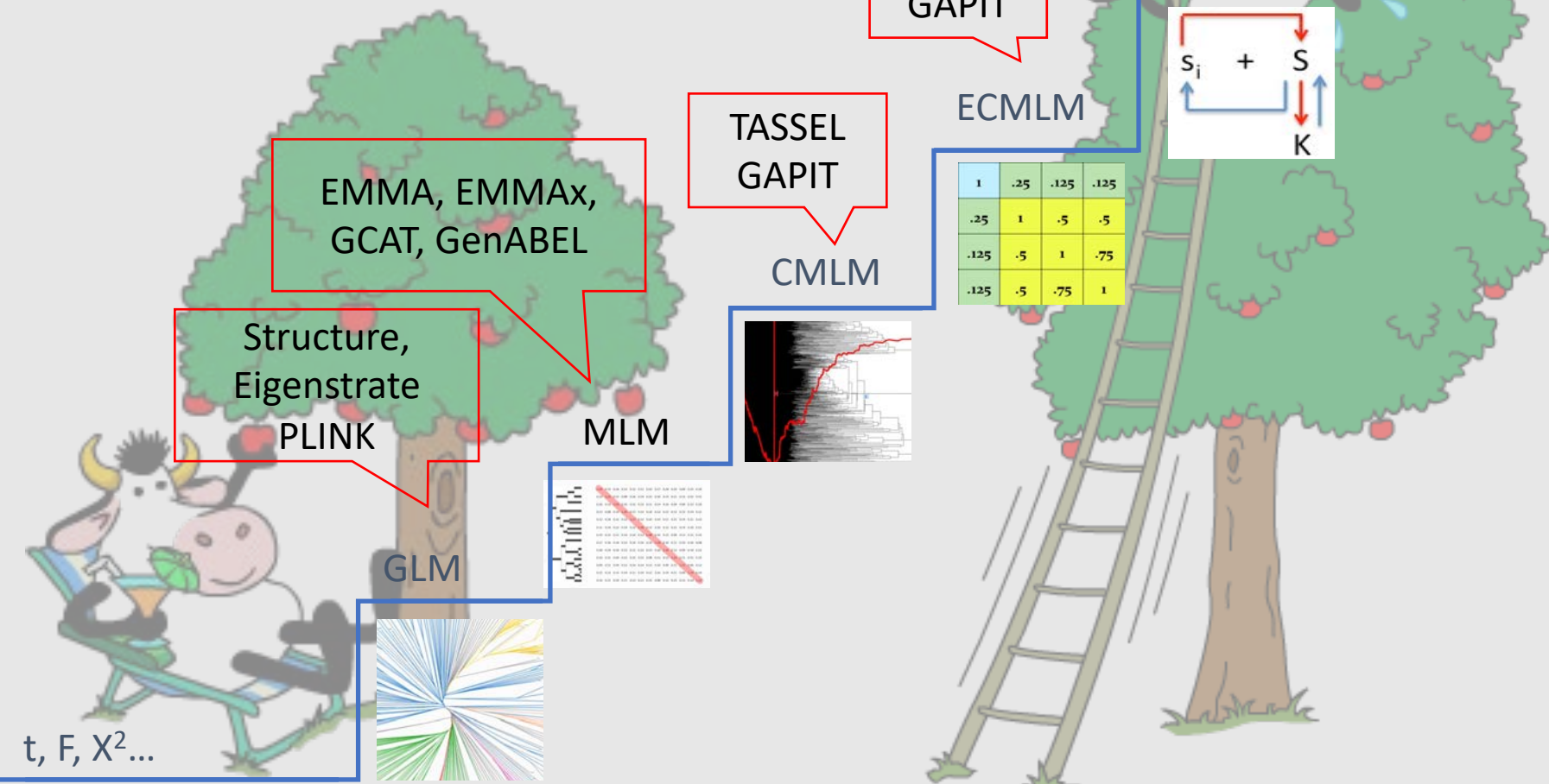


Computation efficiency



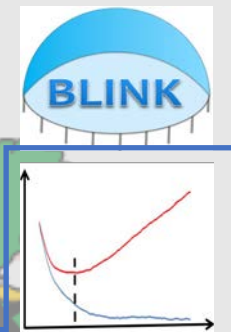
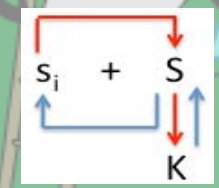


GAPIT

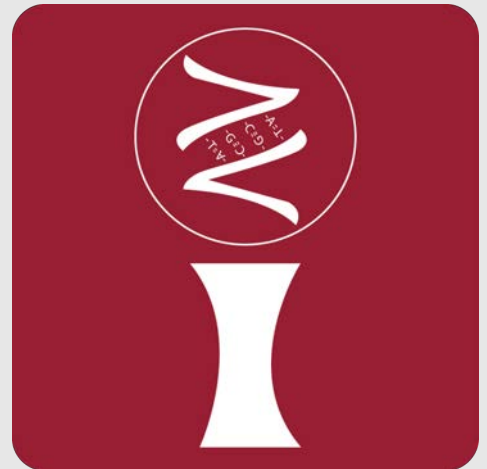



ECMLM

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



Uncorrelated or equally correlated



iPat

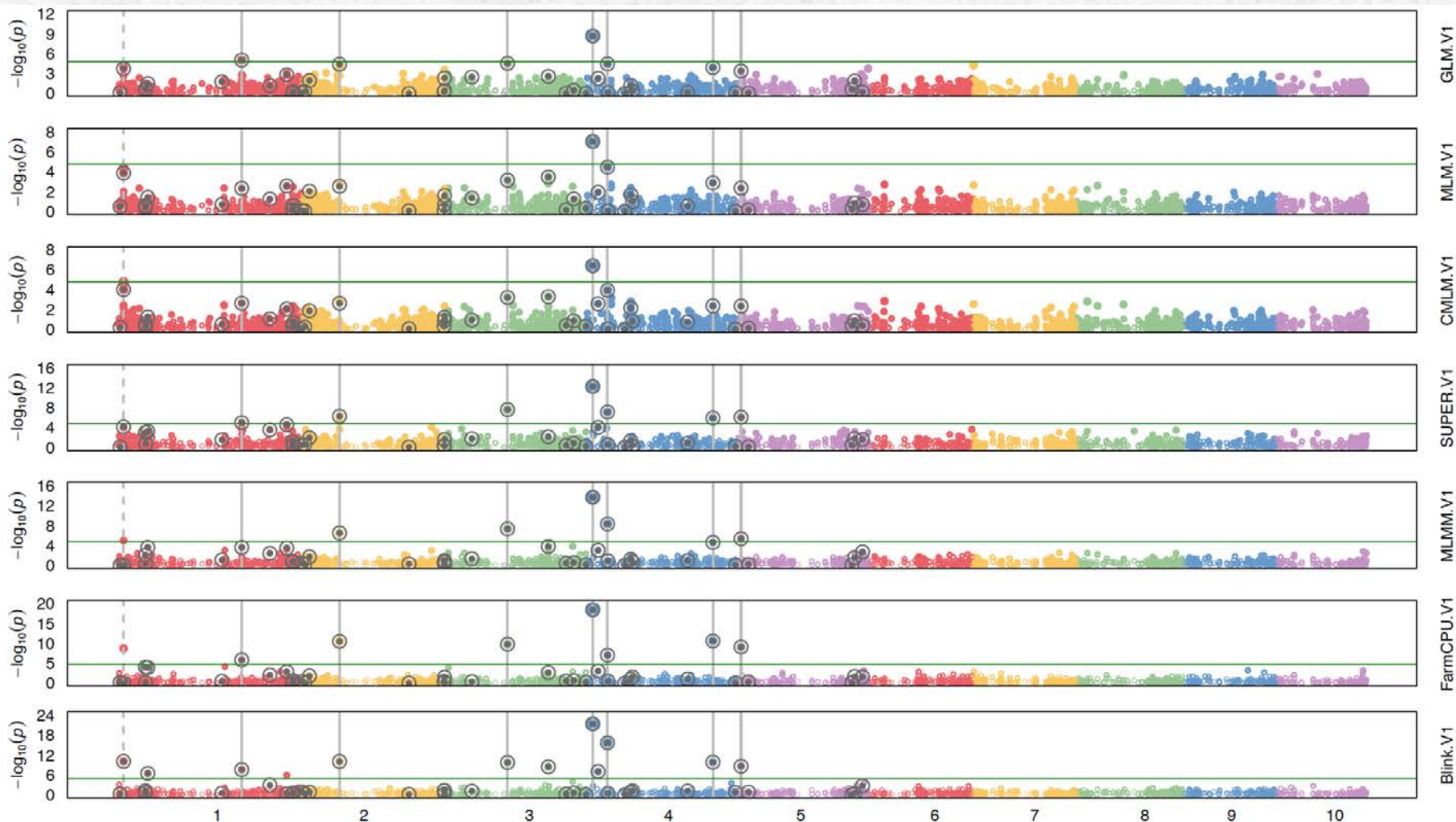
```
source("http://zzlab.net/GAPIT/gapit_functions.txt") #Import demo data
myGD=read.table(file="http://zzlab.net/GAPIT/data/mdp_numeric.txt",head=T)
myGM=read.table(file="http://zzlab.net/GAPIT/data/mdp_SNP_information.txt",head=T)
```

#Simultate 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"))
```



Collaborators and funding



Arron Carter



Mike Pumphrey



Karen Sanguinet



Kawamu Tanaka



Sindhuja Sankaran



Longxi Yu



Jack Brown



Ananth Kalyanaraman



Kim Campbell



Deven See



Camille Steber

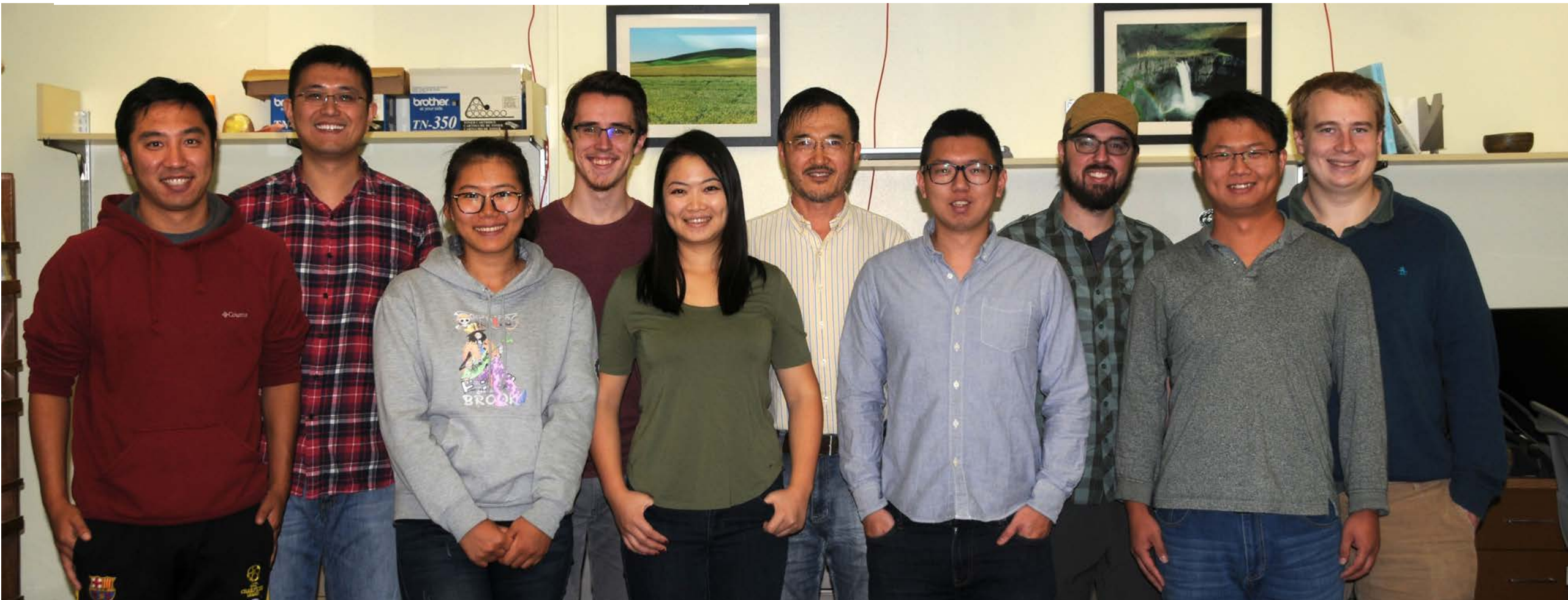


Mike Peel



Shiwu Zhang Laboratory *for Statistical Genomics*

WASHINGTON STATE
UNIVERSITY





Thank you for your attention!