

Genome wide Association Study

Zhiwu Zhang
Washington State University

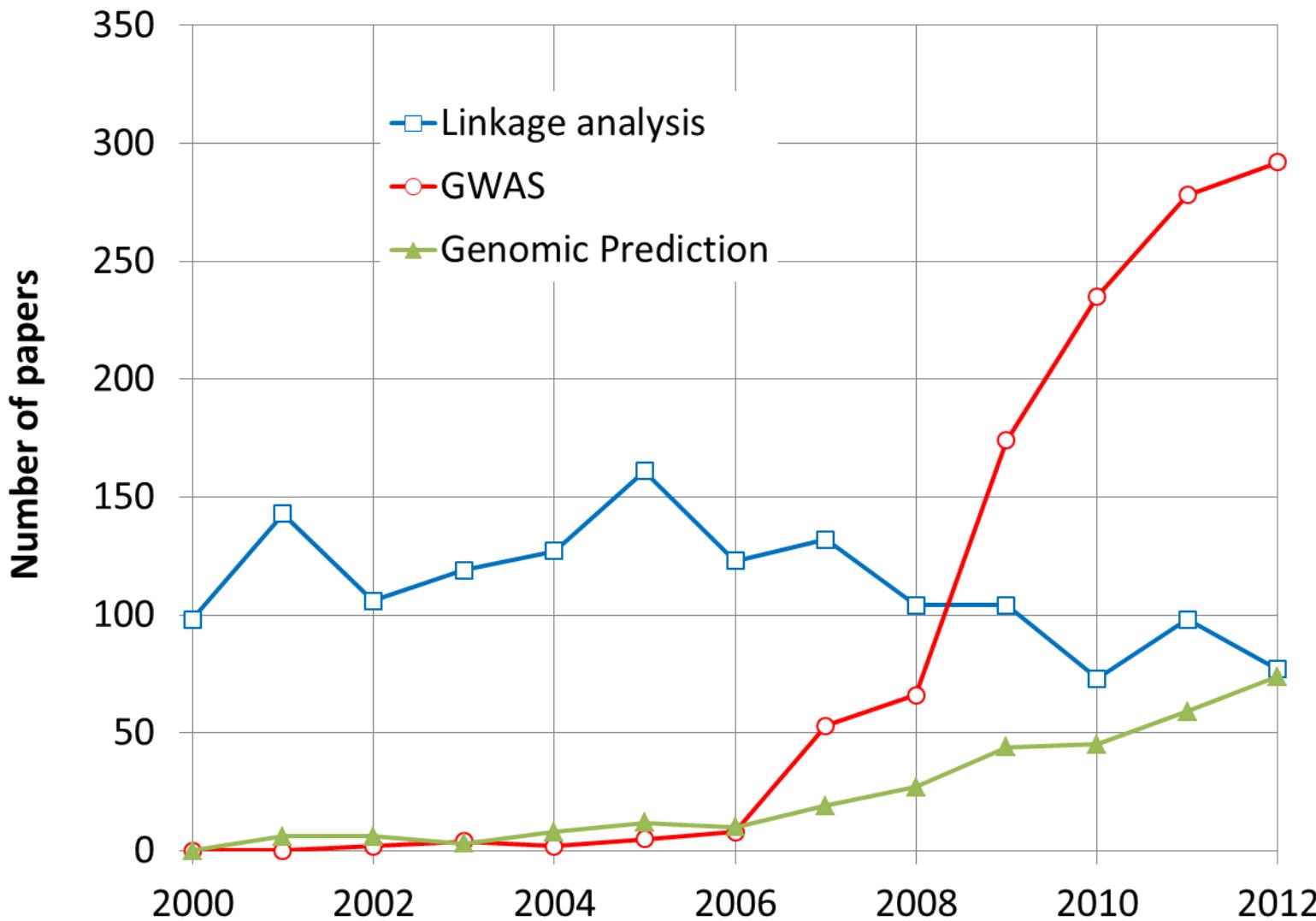


Outline

- **Why GWAS?**
- How does GWAS work?
- How to evaluate GWAS results?
 - Literature
 - Simulation
 - Enrichment analysis
- GWAS Software
- GWAS course at Washington State University

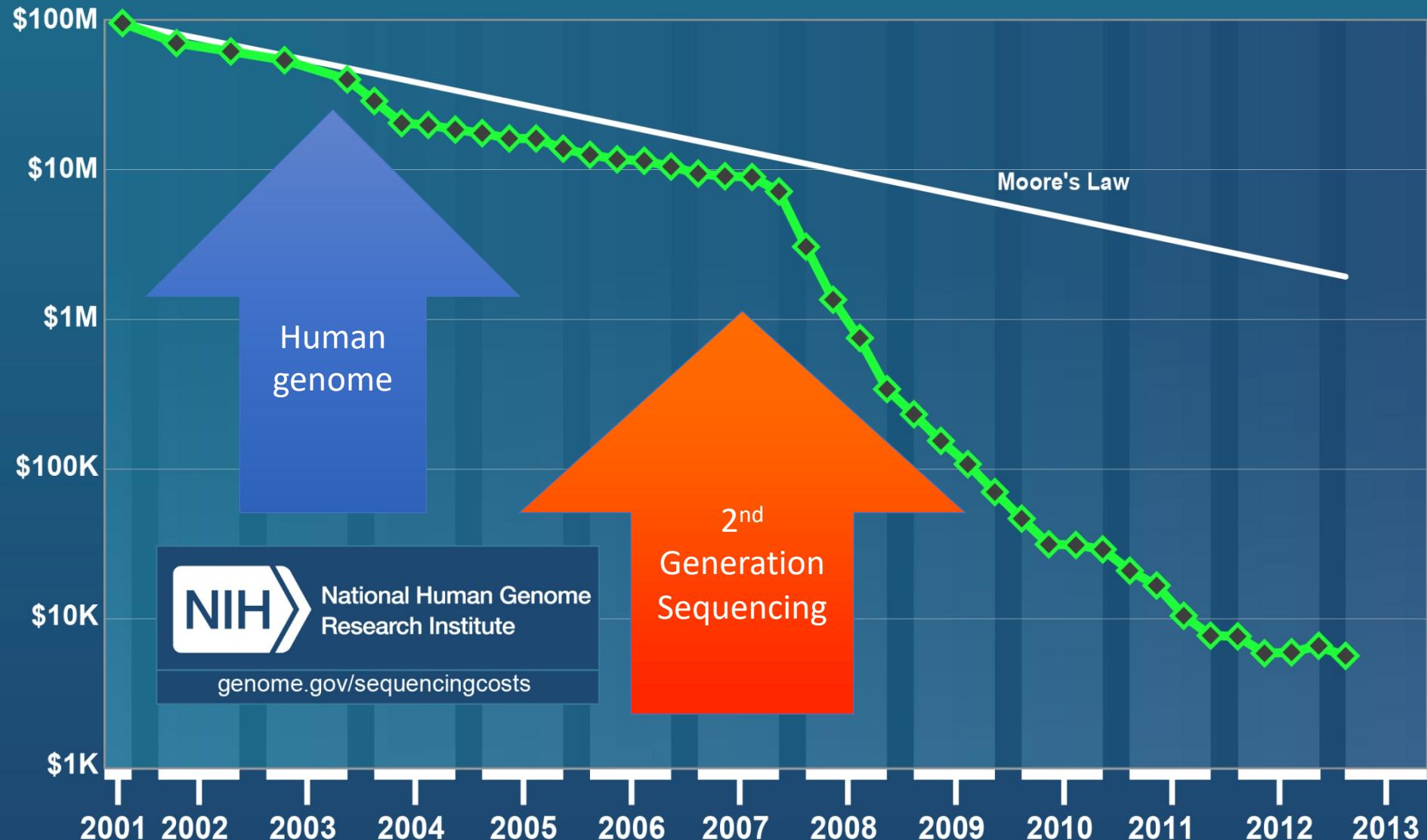


More Research on GWAS and GS



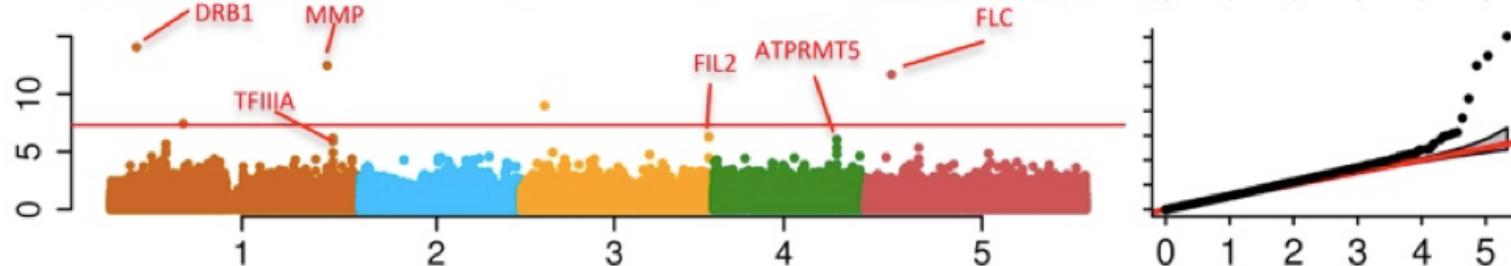
By May 31, 2013

Cost per Genome



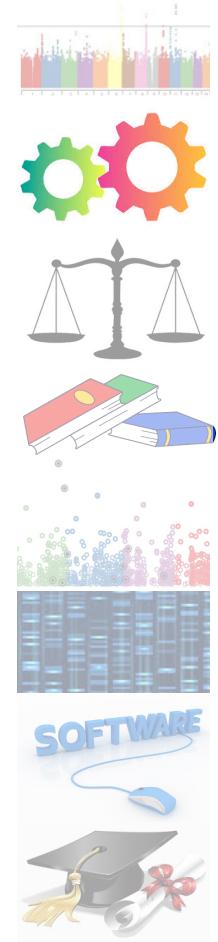
Problems in GWAS

- Computing difficulties: millions of markers, individuals, and traits
- False positives, ex: “Amgen scientists tried to replicate **53** high-profile cancer research findings, but could only replicate **6**”, Nature, 2012, 483: 531
- False negatives

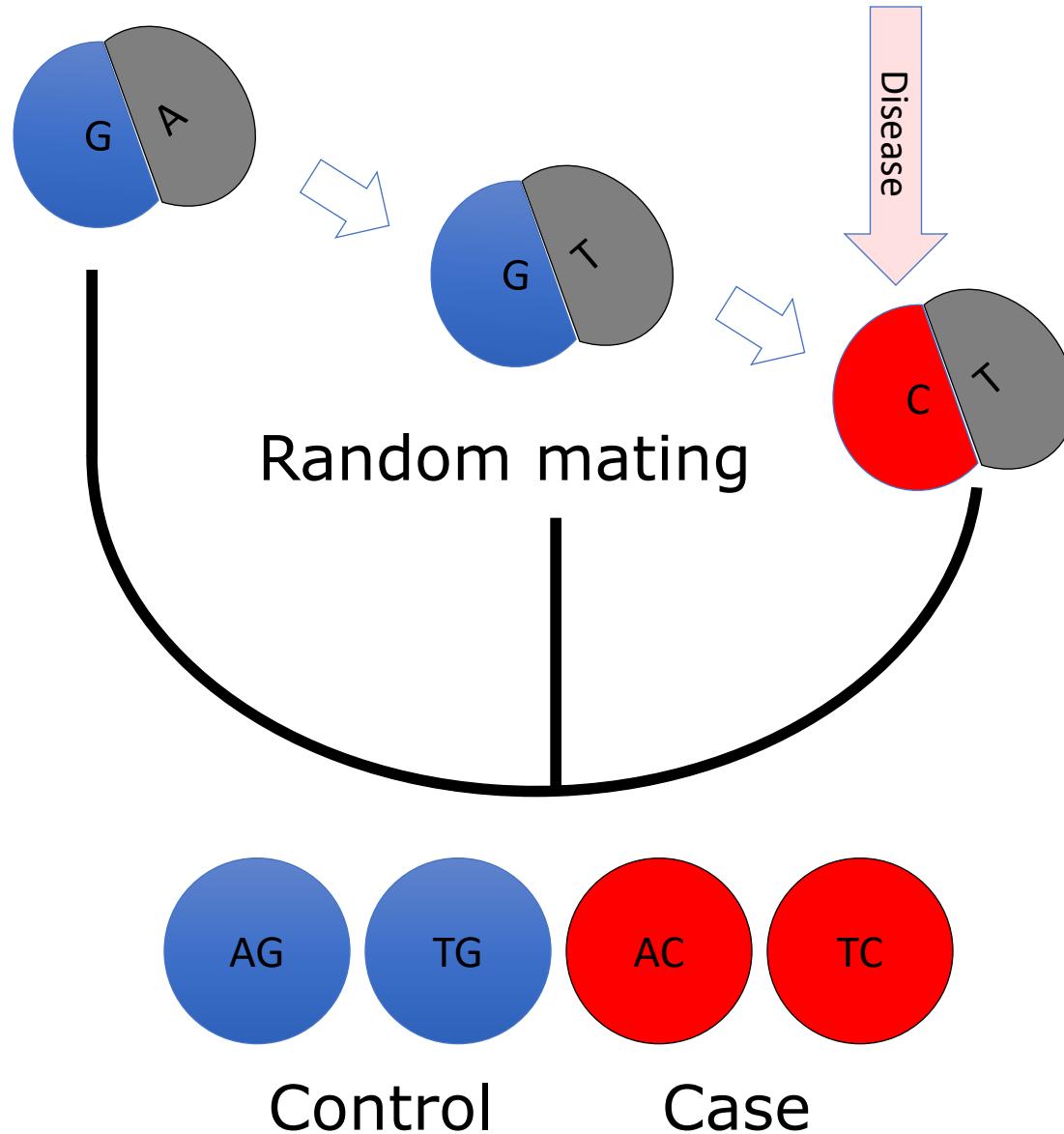


Outline

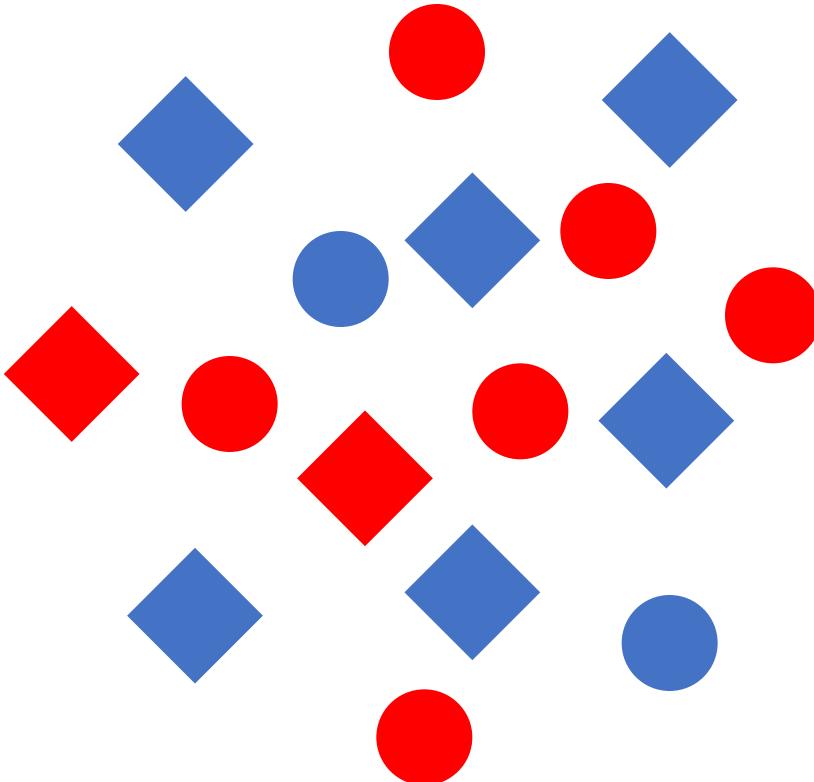
- Why GWAS?
- **How does GWAS work?**
- How to evaluate GWAS results?
 - Literature
 - Simulation
 - Enrichment analysis
- GWAS Software
- GWAS course at Washington State University



Linkage equilibrium



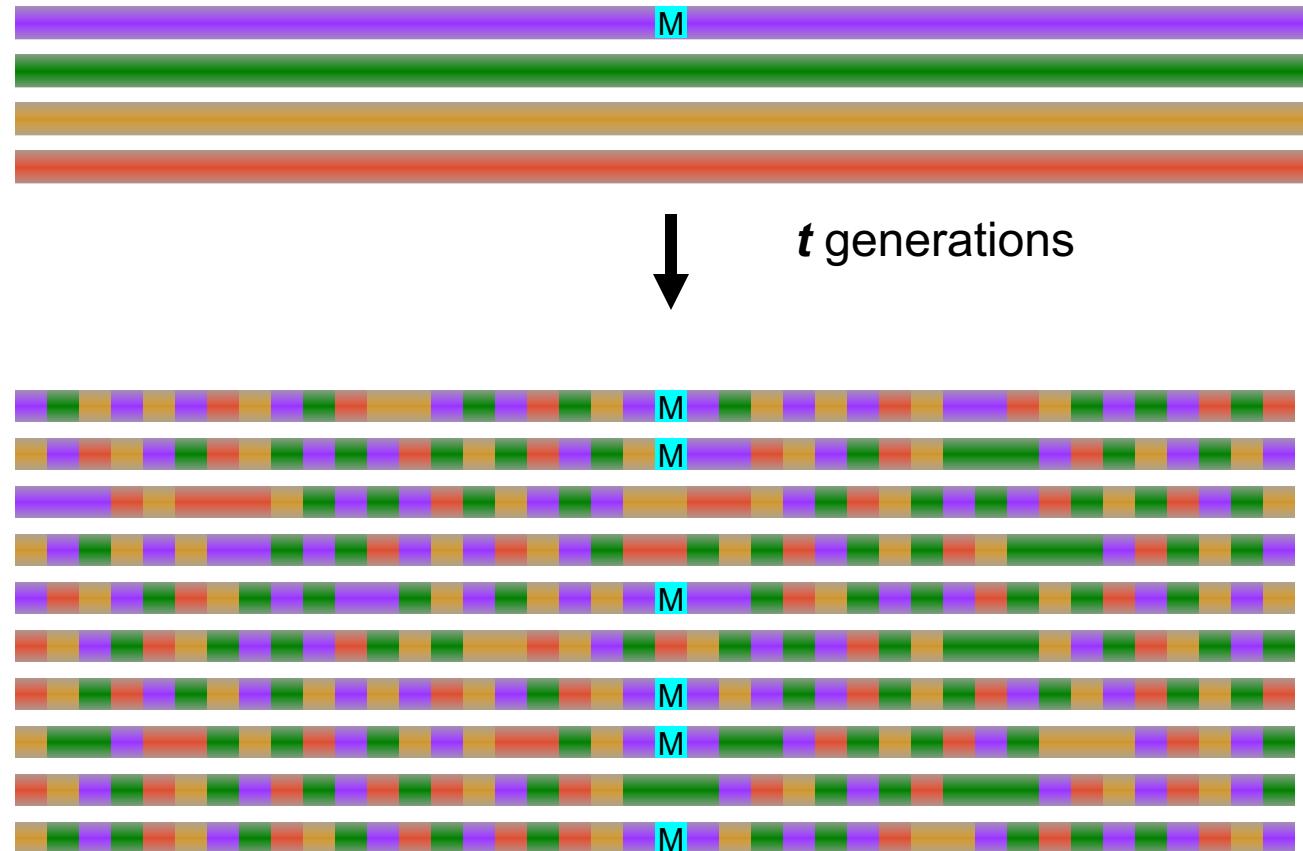
Association study



Marker	Control	Case
	6	2
	2	6

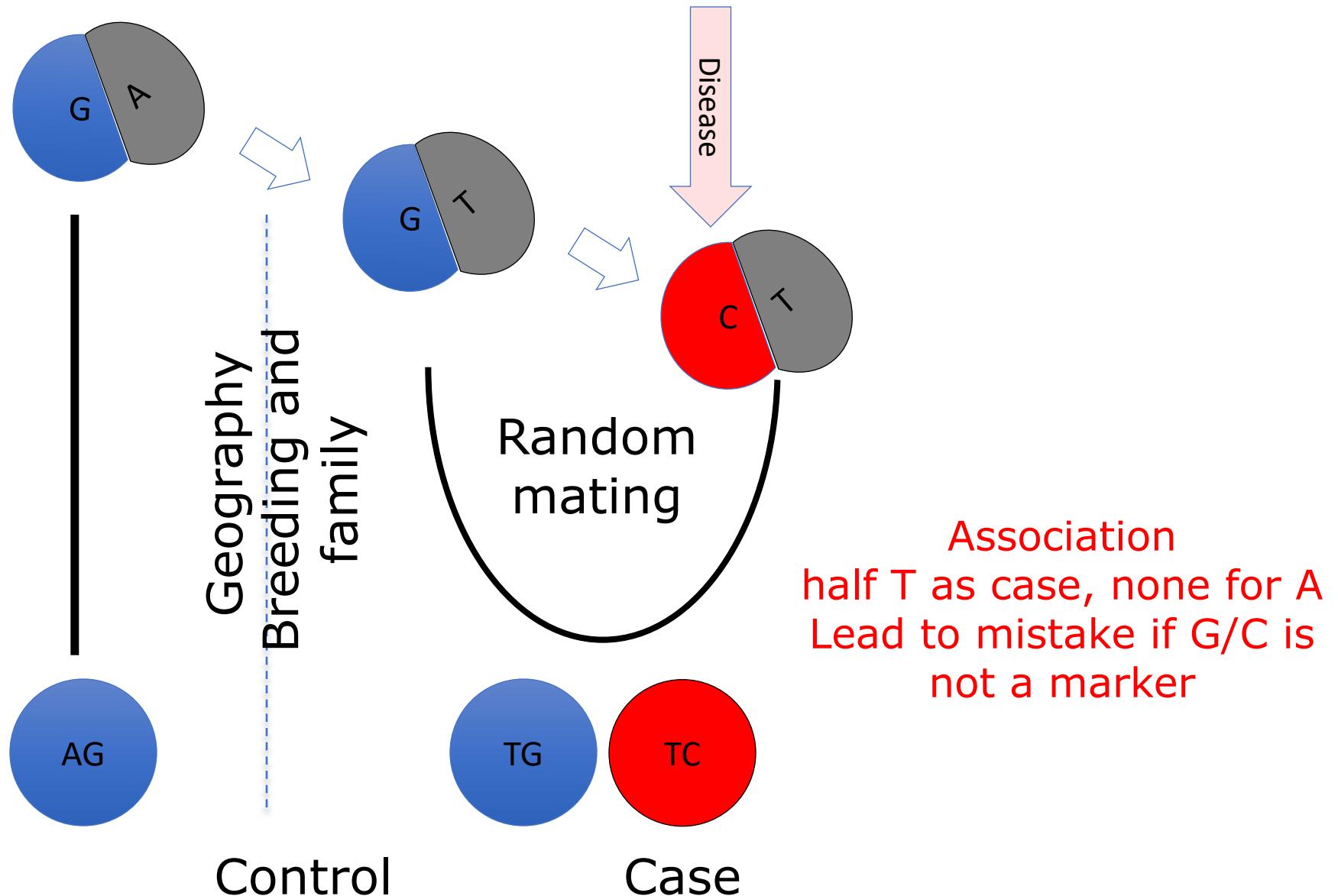
$$\chi^2 = 4(2*2/4) = 4, \text{ df} = 1, \\ P = 4.5\%$$

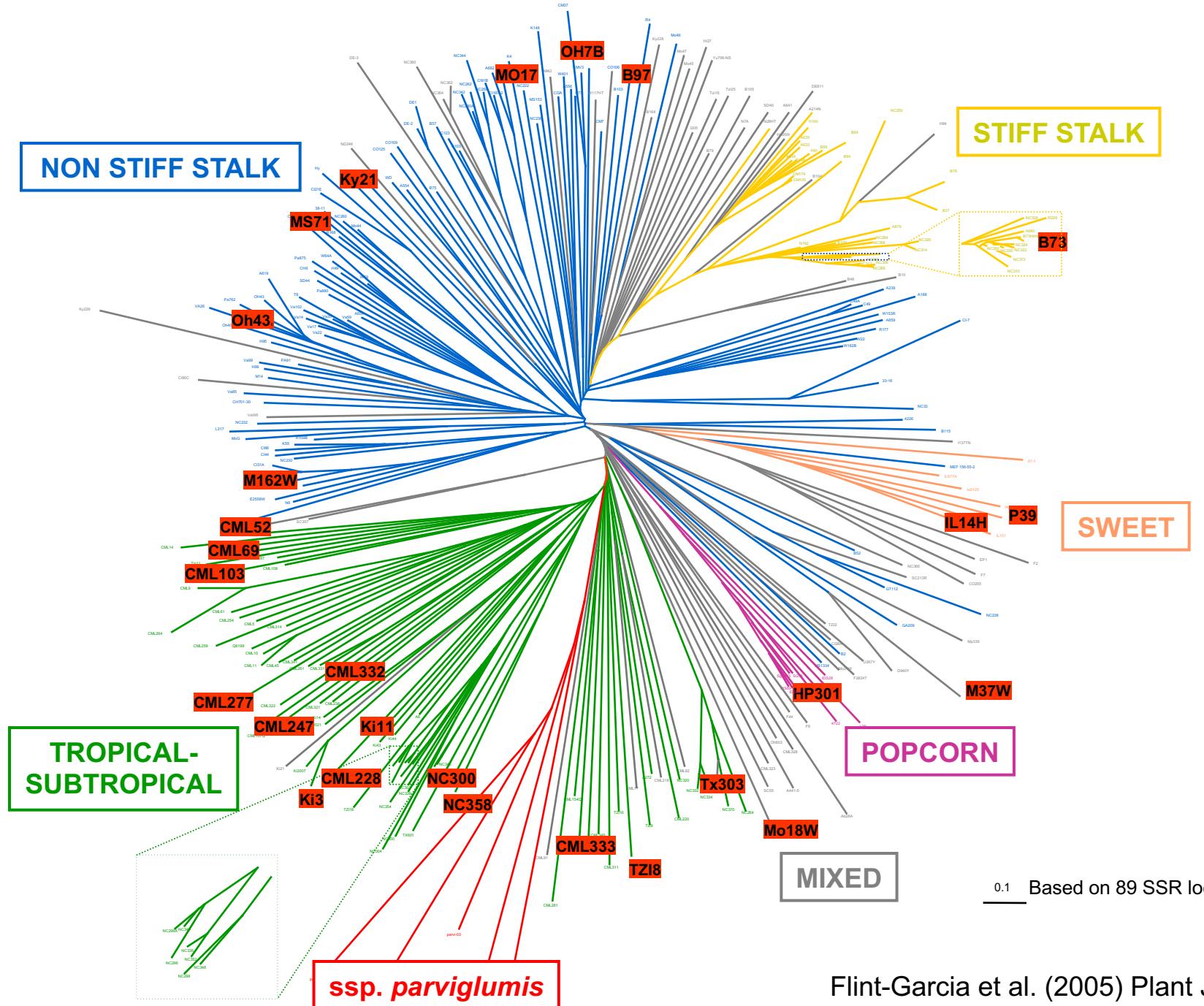
Association study via linkage disequilibrium



Jianming Yu, 2011

Linkage disequilibrium (LD)





Flint-Garcia et al. (2005) Plant J. 44: 1054

Factors Affecting Statistical Power



Number of genes

Gene effect size

Heritability

Population size

Marker density

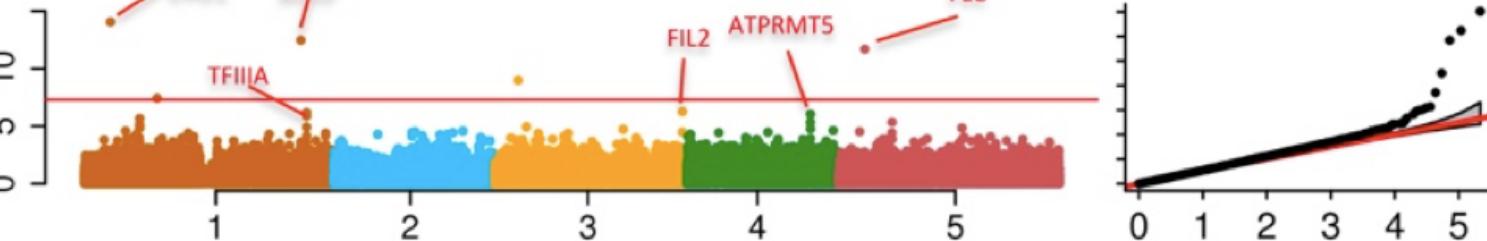
Population structure

Resolution

LD decade

Multiple test correction

Statistical methods



Outline

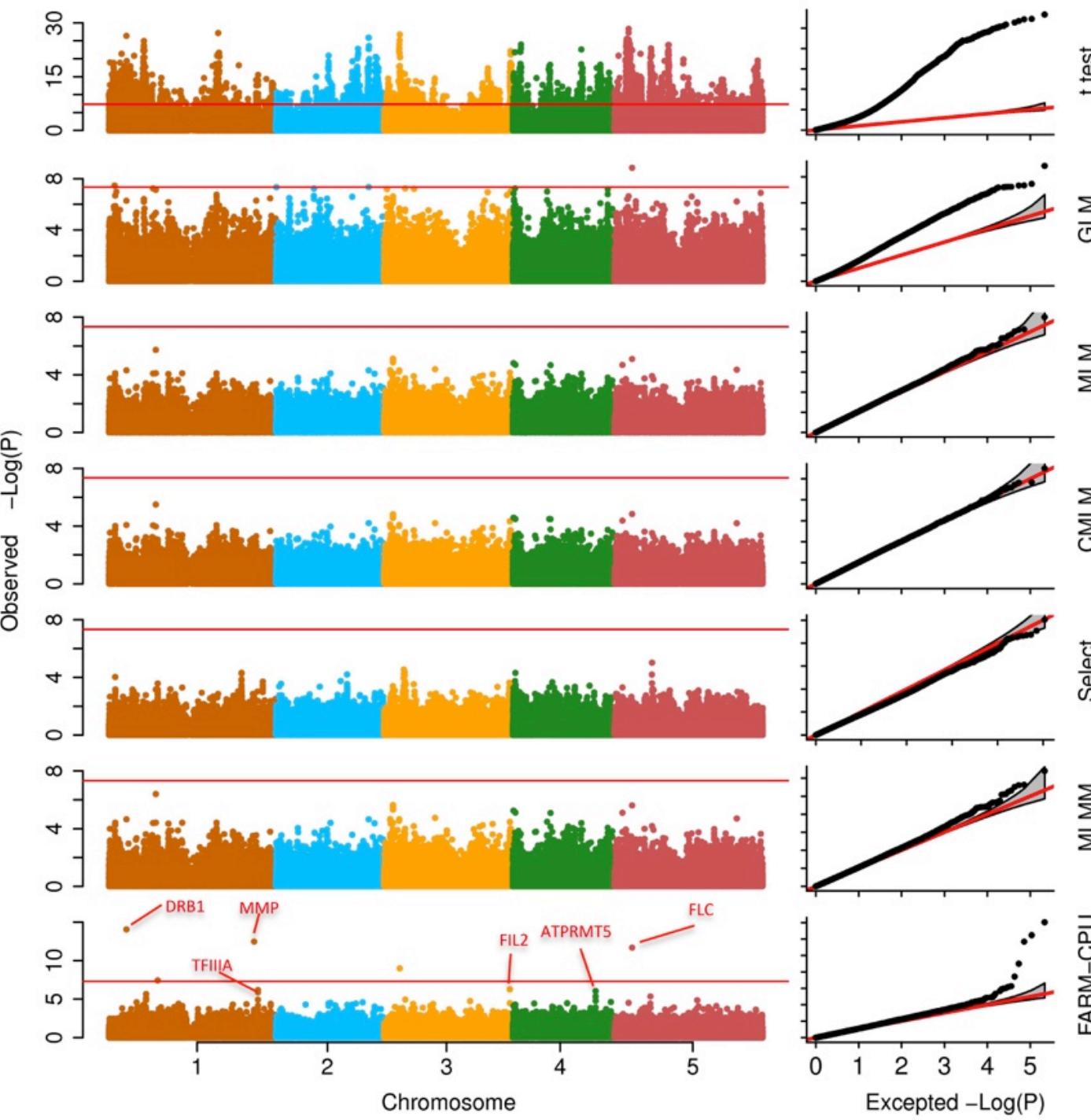
- Why GWAS?
- How does GWAS work?
- **How to evaluate GWAS results?**
 - Literature
 - Simulation
 - Enrichment analysis
- GWAS Software
- GWAS course at Washington State University

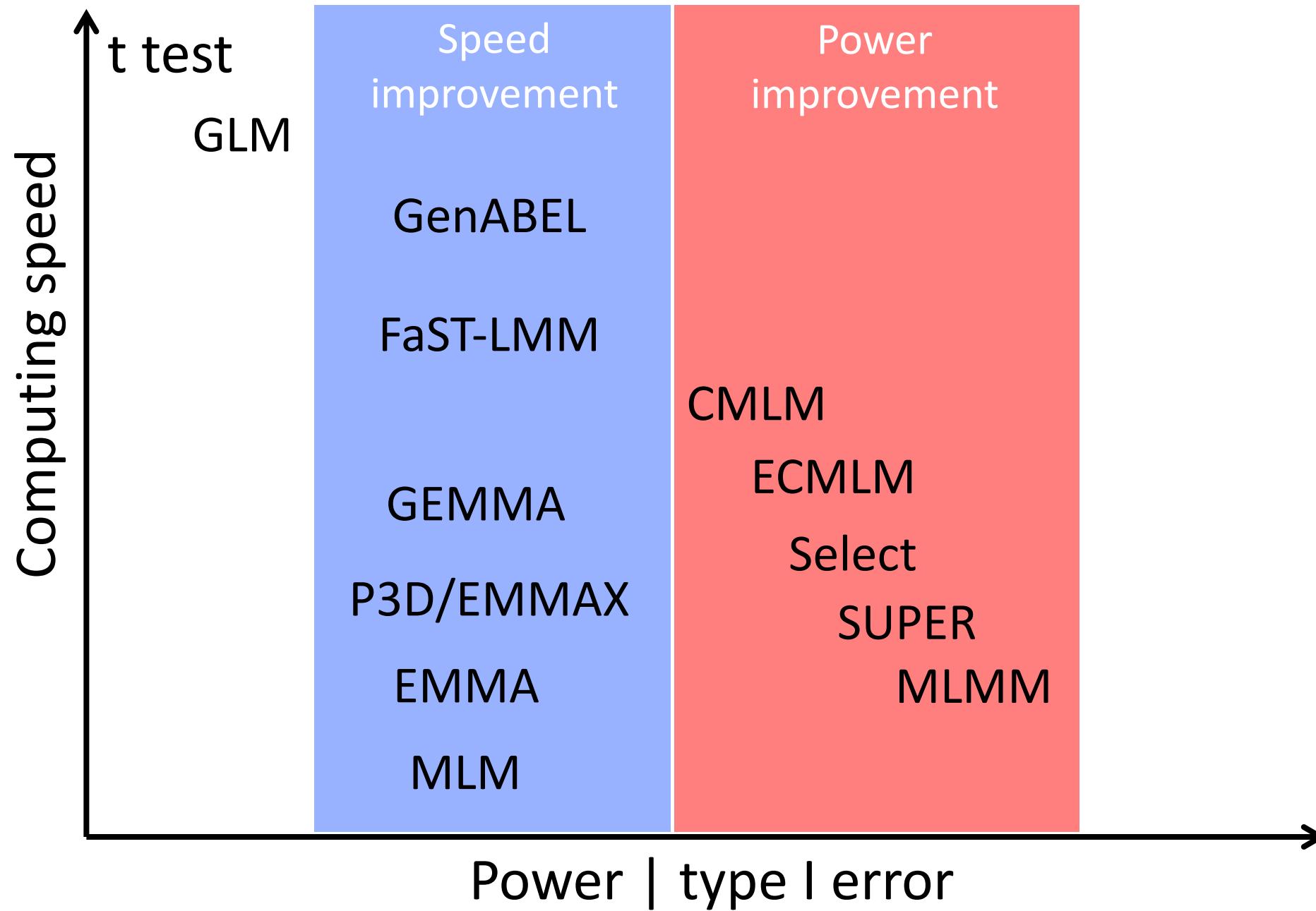


Stream



Associations on flowering time





Model Development

s_i : Testing marker

$$\begin{array}{l} \text{t test} \\ y = s_i + e \end{array}$$

Q: Population structure

$$\begin{array}{l} \text{GLM} \\ y = s_i + Q + e \end{array}$$

→ Adjustment on marker

K: Kinship

$$\begin{array}{l} \text{MLM} \\ y = s_i + Q + K + e \end{array}$$

S: Pseudo QTNs

$$\begin{array}{l} \text{MLMM} \\ y = s_i + S + Q + K + e \end{array}$$

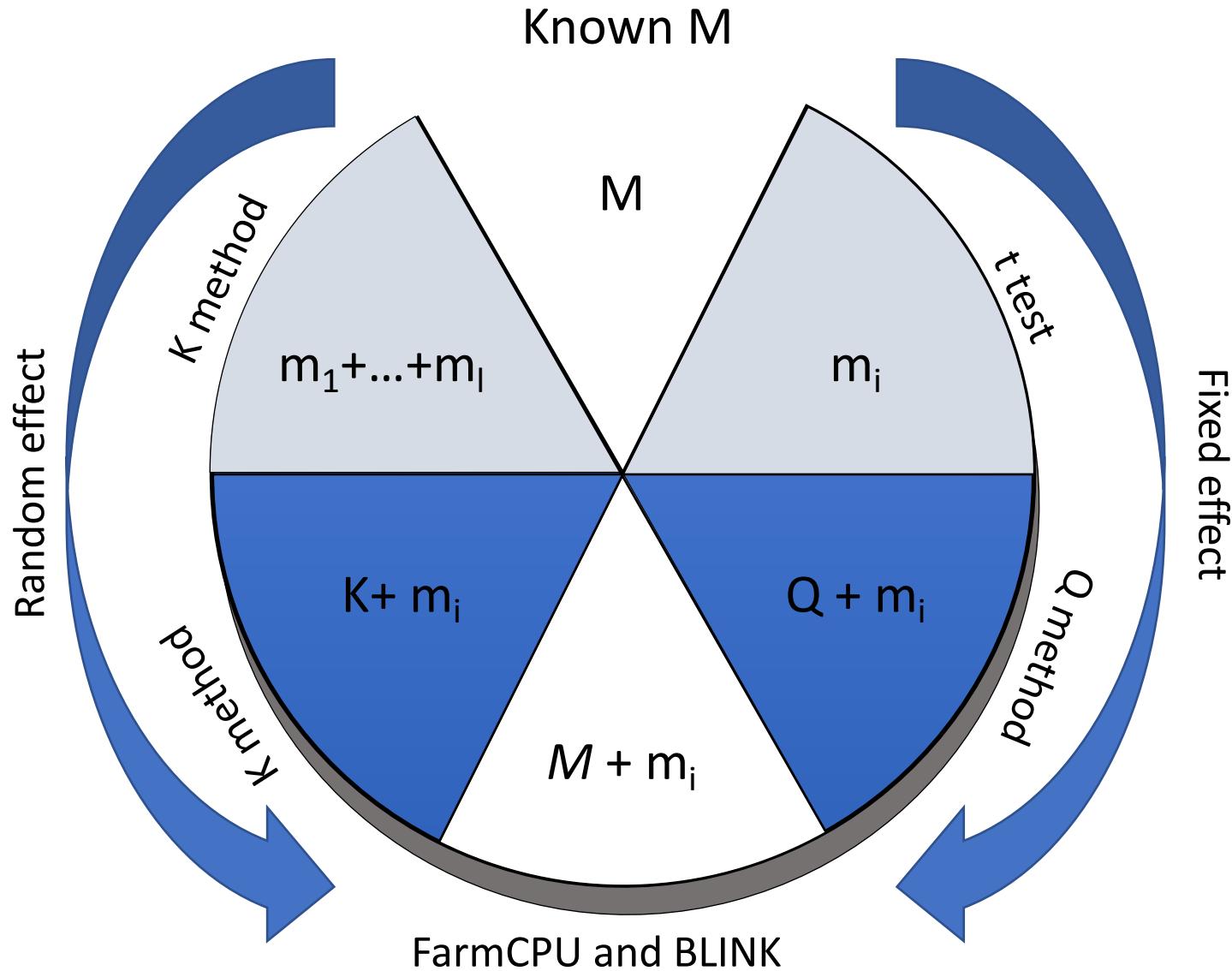
→ Adjustment on covariates

$$\begin{array}{l} \text{SUPER} \\ y = s_i + K + Q + e \end{array}$$

$$\begin{array}{l} \text{FarmCPU} \\ y = s_i + S + e \\ y = K + e \end{array}$$

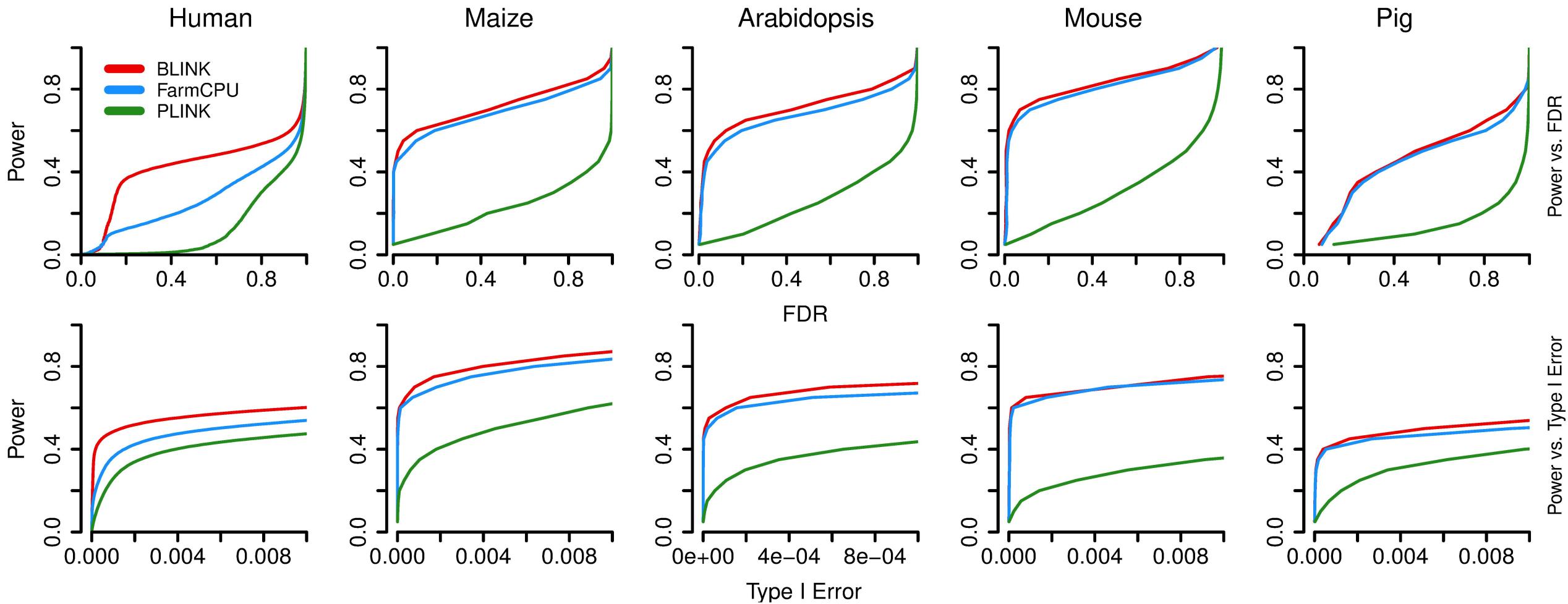
$$\begin{array}{l} \text{BLINK} \\ y = s_i + S + e \\ y = S + e \end{array}$$

GWAS methods



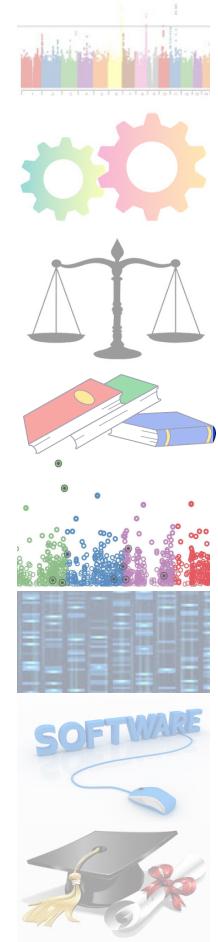
m: marker, M: Mutations, M : Estimated mutations

Same trend across species



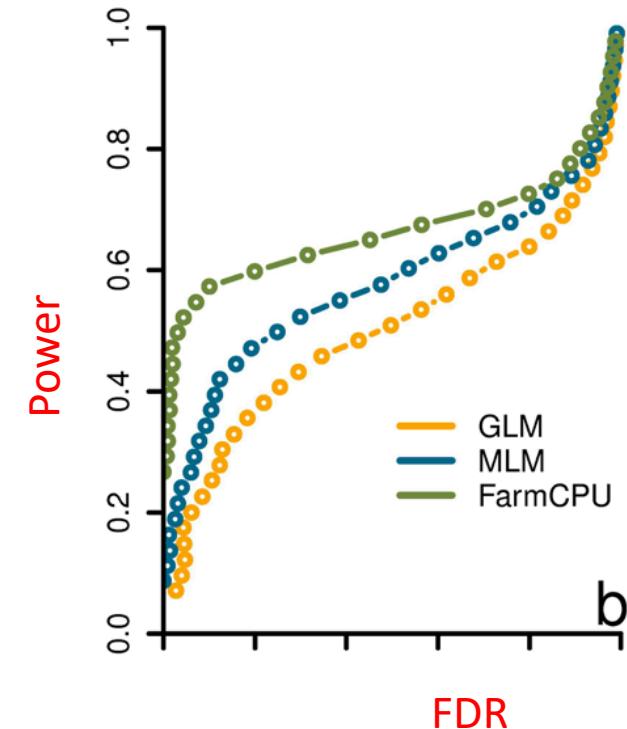
Outline

- Why GWAS?
- How does GWAS work?
- How to evaluate GWAS results?
 - Literature
 - **Simulation**
 - Enrichment analysis
- GWAS Software
- GWAS course at Washington State University



ROC curve

- Receiver Operating Characteristic
- "The curve is created by plotting the true positive rate against the false positive rate at various threshold settings." -Wikipedia



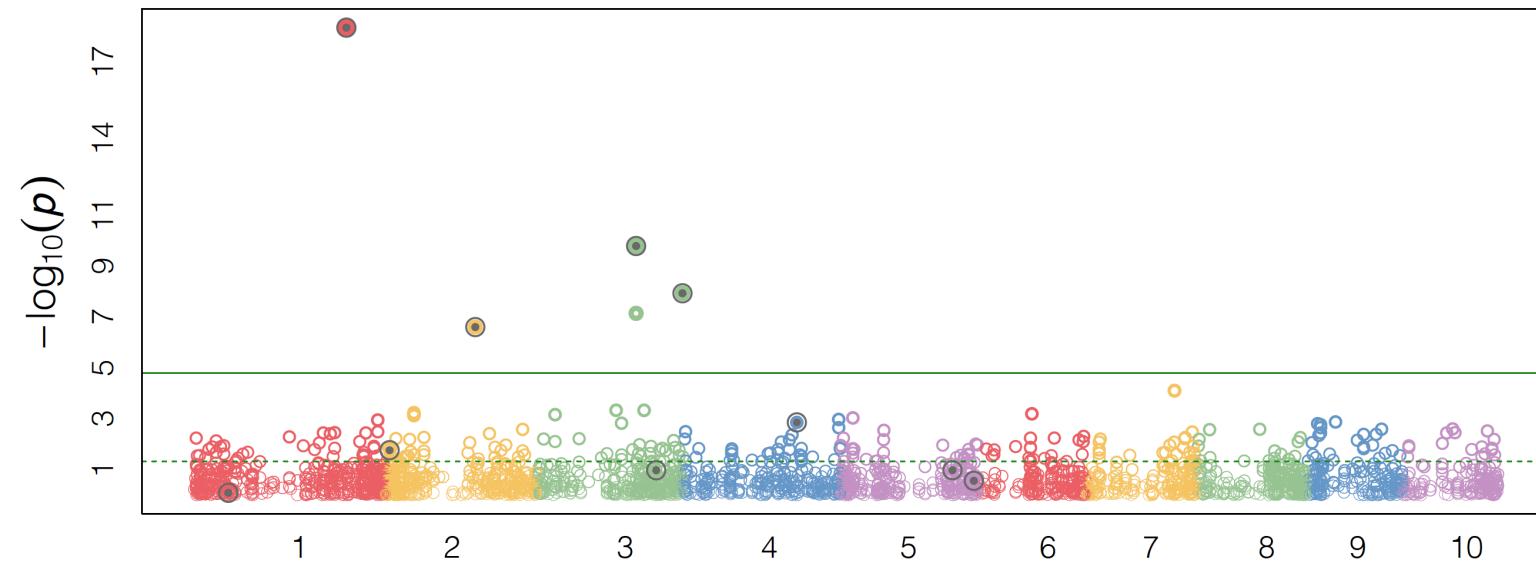
Liu et. al. PLoS Genetics, 2016

Genotypes

taxa	PZB00859.1	PZA01271.1	PZA03613.2	PZA03613.1	PZA03614.2	PZA03614.1	PZA00258.3	SNP	Chromosome	Position
33-16	2	0	0	2	2	2	2	PZB00859.1	1	157104
38-11	2	2	0	2	2	2	0	PZA01271.1	1	1947984
4226	2	0	0	2	2	2	0	PZA03613.2	1	2914066
4722	2	2	0	2	2	2	1	PZA03613.1	1	2914171
A188	0	0	0	2	2	2	0	PZA03614.2	1	2915078
A214N	2	0	2	0	2	0	0	PZA03614.1	1	2915242
A239	0	0	2	2	0	0	0	PZA00258.3	1	2973508
A272	0	0	2	2	0	0	2	PZA02962.13	1	3205252
A441-5	2	0	0	2	2	2	0	PZA02962.14	1	3205262
A554	2	2	2	2	0	2	0	PZA00599.25	1	3206090
A556	2	0	0	2	2	2	1			
A6	0	0	2	2	0	0	0			
A619	2	2	0	2	2	2	0			
A632	2	0	2	0	2	0	0			
A634	2	0	2	0	2	0	0			
A635	2	0	2	0	2	0	0			

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

GWAS by GLM in GAPIT



Import functions and data

Restrict genes on CHR1-5

Simulate phenotype

GWAS with GLM using GAPIT

```
rm(list=ls())
library(compiler) #required for cmpfun
source("http://www.zzlab.net/GAPIT/emma.txt")
source("http://www.zzlab.net/GAPIT/gapit_functions.txt")
#source("http://www.zzlab.net/StaGen/2021/R/G2P.R")
myGD=read.table(file="http://www.zzlab.net/GAPIT/data/mdp_numeric.txt",head=T)
myGM=read.table(file="http://zzlab.net/GAPIT/data/mdp_SNP_information.txt",head=T)

taxa=myGD[,1]
X=myGD[,-1]
index1to5=myGM[,2]<6
X1to5 = X[,index1to5]
GD.candidate=cbind(taxa,X1to5)

set.seed(99164)
setwd("~/Desktop/temp")
#mySim=G2P(X=X1to5,h2=.75,alpha=1,NQTN=10,distribution="normal")
mySim=GAPIT.Phenotype.Simulation(GD=GD.candidate,GM=myGM[index1to5],h2=.75,NQTN=10, effectunit=.95,QTNDist="normal")
myY=mySim$Y

myGAPIT=GAPIT(Y=myY, GD=myGD,
GM=myGM, QTN.position=mySim$QTN.position,
PCA.total=3, group.from=1, group.to=1, group.by=10,
memo="GLM", file.output=TRUE,)
```

Bins (e.g. 100Kb)

```
bigNum=1e9
```

```
resolution=100000
```

```
bin=round((myGM[,2]*bigNum+myGM[,3])/resolution)
```

```
myGWAS=cbind(myGM,myGAPIT$mp,bin)
```

```
head(myGWAS)
```

	SNP	Chromosome	Position	myGAPIT\$mp	bin
1	PZB00859.1	1	157104	0.5928172	10002
2	PZA01271.1	1	1947984	0.2223674	10019
3	PZA03613.2	1	2914066	0.557306	10029
4	PZA03613.1	1	2914171	0.4605953	10029
5	PZA03614.2	1	2915078	0.5398164	10029
6	PZA03614.1	1	2915242	0.3681861	10029

Minimum p value within bin

Bins of QTNs

QTN.bin=myGWAS[mySim\$QTN.position,]

QTN.bin

	SNP	Chromosome	Position	myGAPIT\$mp	bin
1194	PZA00303.6	3	187604478	1.06E-01	31876
1094	PZB01683.3	3	156252478	1.67E-10	31563
1942	PZA01680.3	5	208901002	2.76E-01	52089
1497	PZA01187.1	4	177666738	1.41E-03	41777
1287	PZA00088.3	3	228614270	1.20E-08	32286
150	PZA00617.16	1	53357515	8.19E-01	10534
574	PZA02081.1	2	5923120	1.73E-02	20059
340	PZA00381.3	1	237639087	4.43E-19	12376
761	PZB01487.1	2	139752027	2.54E-07	21398
1849	PZA03718.1	5	175459060	1.06E-01	51755

Sorted bins of QTNs

```
index.qtn.p=order(QTN.bin[,4])
```

```
QTN.bin[index.qtn.p,]
```

	SNP	Chromosome	Position	myGAPIT\$mp	bin
340	PZA00381.3	1	237639087	4.43E-19	12376
1094	PZB01683.3	3	156252478	1.67E-10	31563
1287	PZA00088.3	3	228614270	1.20E-08	32286
761	PZB01487.1	2	139752027	2.54E-07	21398
1497	PZA01187.1	4	177666738	1.41E-03	41777
574	PZA02081.1	2	5923120	1.73E-02	20059
1849	PZA03718.1	5	175459060	1.06E-01	51755
1194	PZA00303.6	3	187604478	1.06E-01	31876
1942	PZA01680.3	5	208901002	2.76E-01	52089
150	PZA00617.16	1	53357515	8.19E-01	10534

FDR and type I error

Total number of bins: **1365** (bin size of 100kb)

SNP	CHR	Position	myGAPIT\$mp
PZA00381.3	1	237639087	4.43E-19
PZB01683.3	3	156252478	1.67E-10
PZA00088.3	3	228614270	1.20E-08
PZB01487.1	2	139752027	2.54E-07
PZA01187.1	4	177666738	1.41E-03
PZA02081.1	2	5923120	1.73E-02
PZA03718.1	5	175459060	1.06E-01
PZA00303.6	3	187604478	1.06E-01
PZA01680.3	5	208901002	2.76E-01
PZA00617.16	1	53357515	8.19E-01

Power	#False bins	FDR	TypeI Error
0.1	0	0	0
0.2	0	0	0
0.3	0	0	0
0.4	0	0	0
0.5	10	0.6666667	0.00732601
0.6	110	0.9482759	0.08058608
0.7	359	0.9808743	0.26300366
0.8	362	0.9783784	0.26520147
0.9	681	0.9869565	0.4989011
1	1141	0.9913119	0.83589744

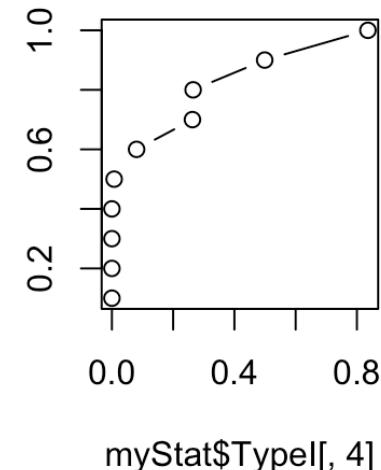
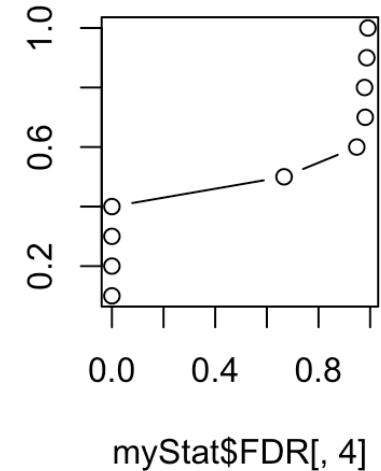
$$0.6666667=10/(10+4)$$

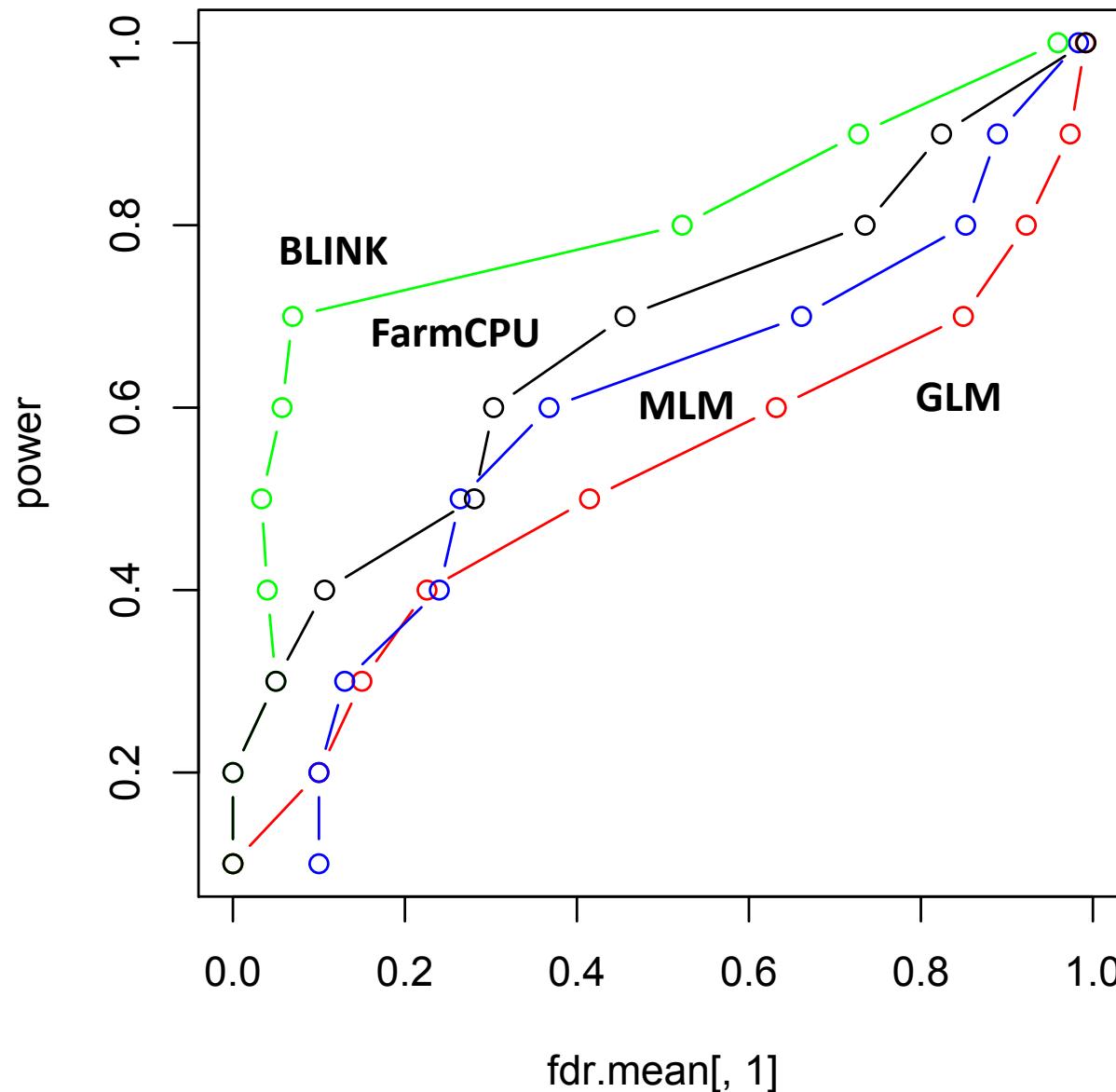
$$0.0073601=10/1365$$

GAPIT.FDR.TypeI Function for Area Under Curve

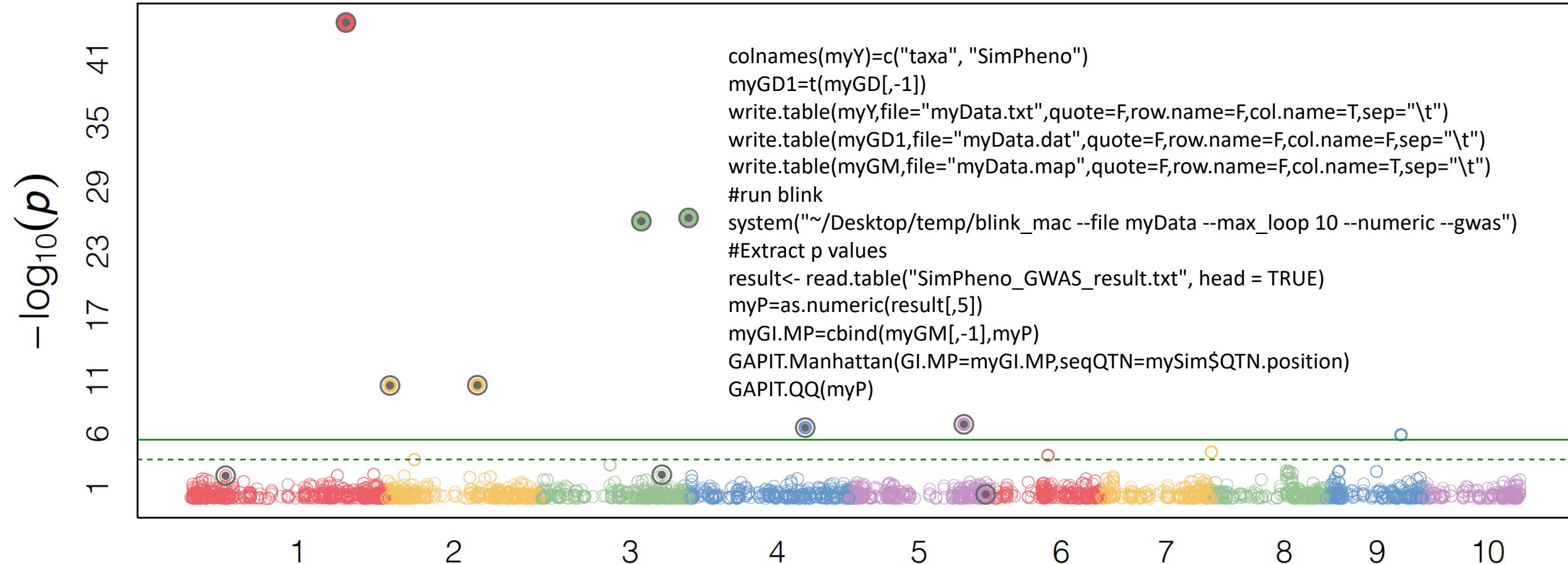
```
library(compiler) #required for cmpfun  
source("http://www.zzlab.net/GAPIT/gapit_functions.  
txt")  
myStat=GAPIT.FDR.TypeI(  
WS=c(1e0,1e3,1e4,1e5), GM=myGM,  
seqQTN=mySim$QTN.position,  
GWAS=myGWAS)
```

```
str(myStat)  
par(mfrow=c(2,1),mar = c(5,2,5,2))  
plot(myStat$FDR[,4],myStat$Power,type="b")  
plot(myStat$TypeI[,4],myStat$Power,type="b")
```



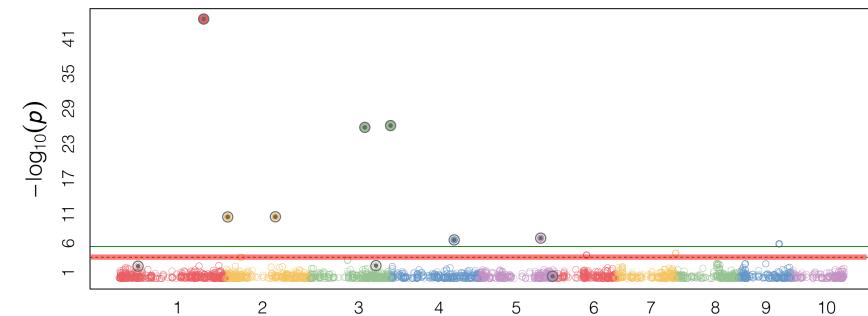
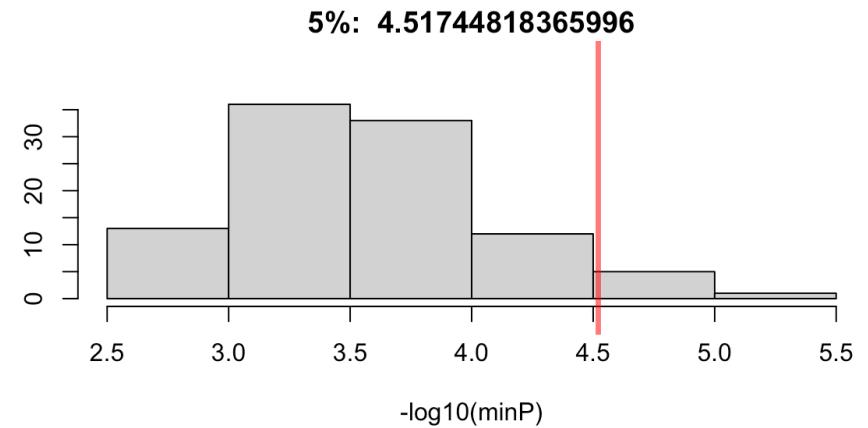


BLINK



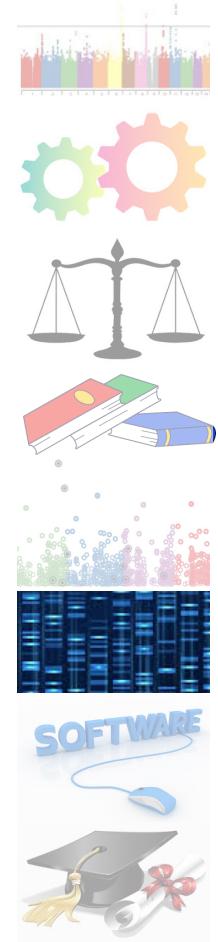
Permutation cutoff

```
nrep=100  
minP=matrix(NA,nrep,1)  
for (i in 1:nrep){  
  randOrder=sample(nrow(myY))  
  write.table(myY[randOrder],file="myData.txt",quote=F,row.  
  name=F,col.name=T,sep="\t")  
  system("~/Desktop/temp/blink_mac --file myData --  
  max_loop 10 --numeric --gwas")  
  result<- read.table("SimPheno_GWAS_result.txt", head =  
  TRUE)  
  minP[i]=min(result[,5])  
}  
hist(-log10(minP),main=paste("5%: ",quantile(-  
  log10(minP),.95)))
```

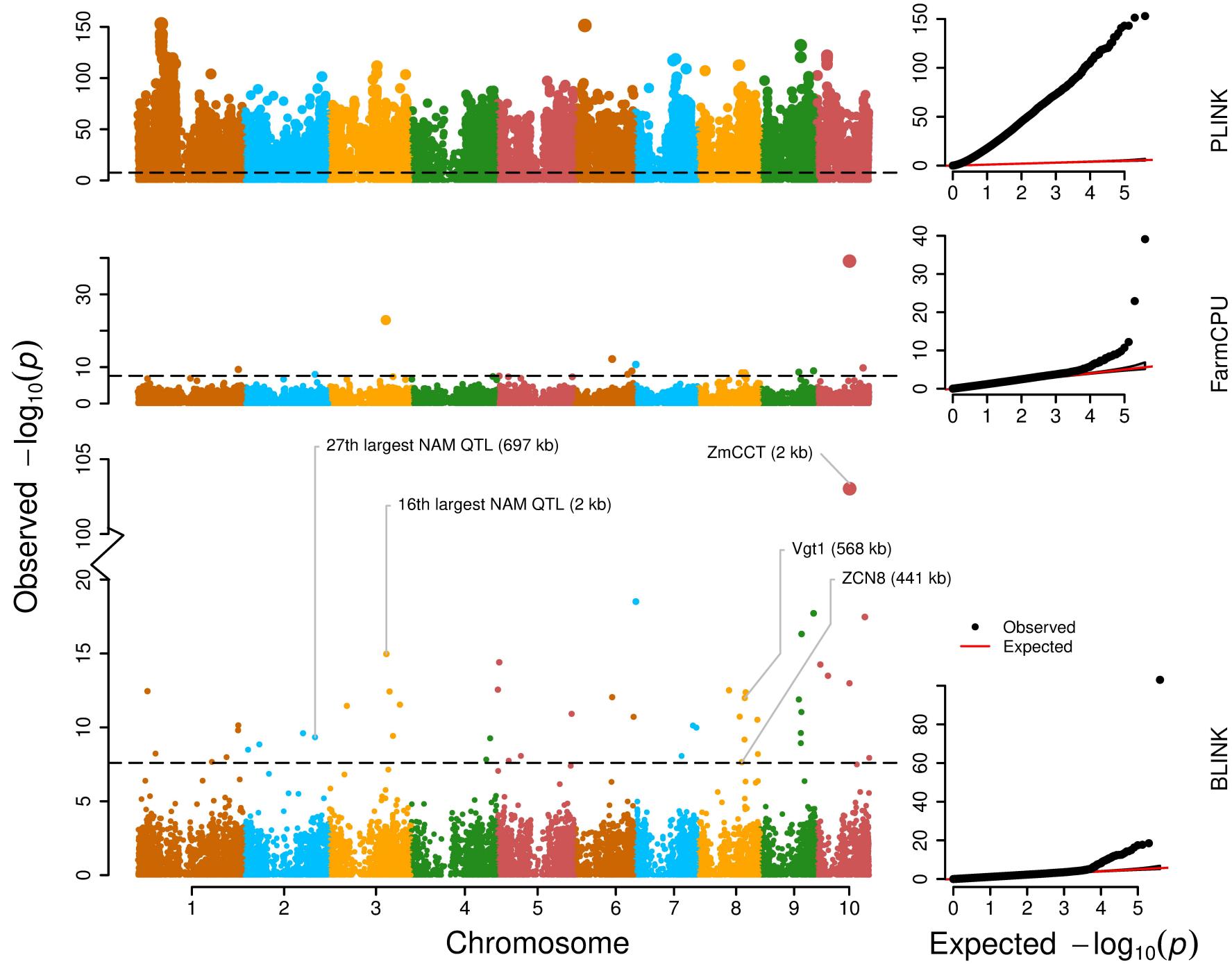


Outline

- Why GWAS?
- How does GWAS work?
- How to evaluate GWAS results?
 - Literature
 - Simulation
 - **Enrichment analysis**
- GWAS Software
- GWAS course at Washington State University

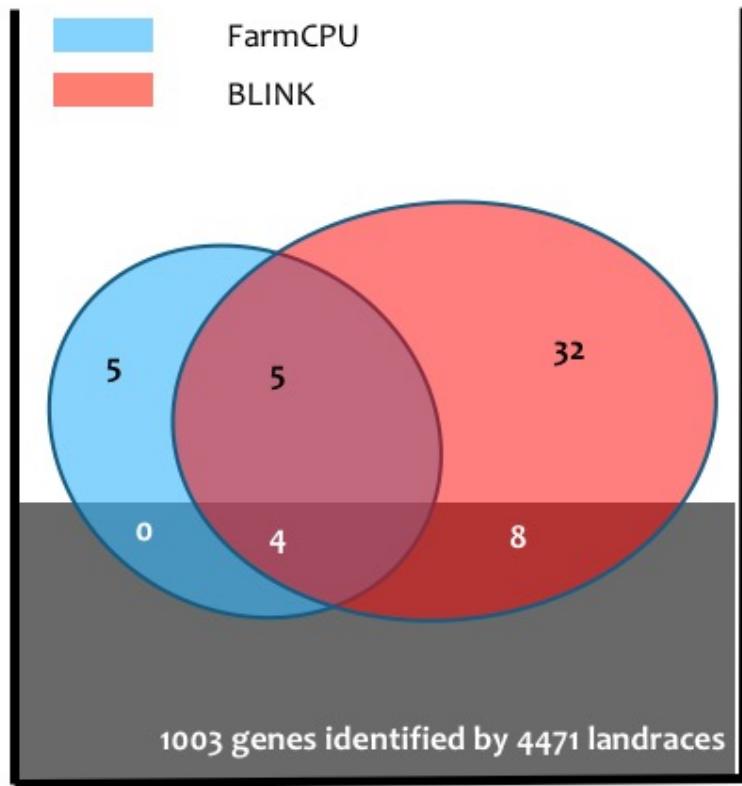


Application in Maize

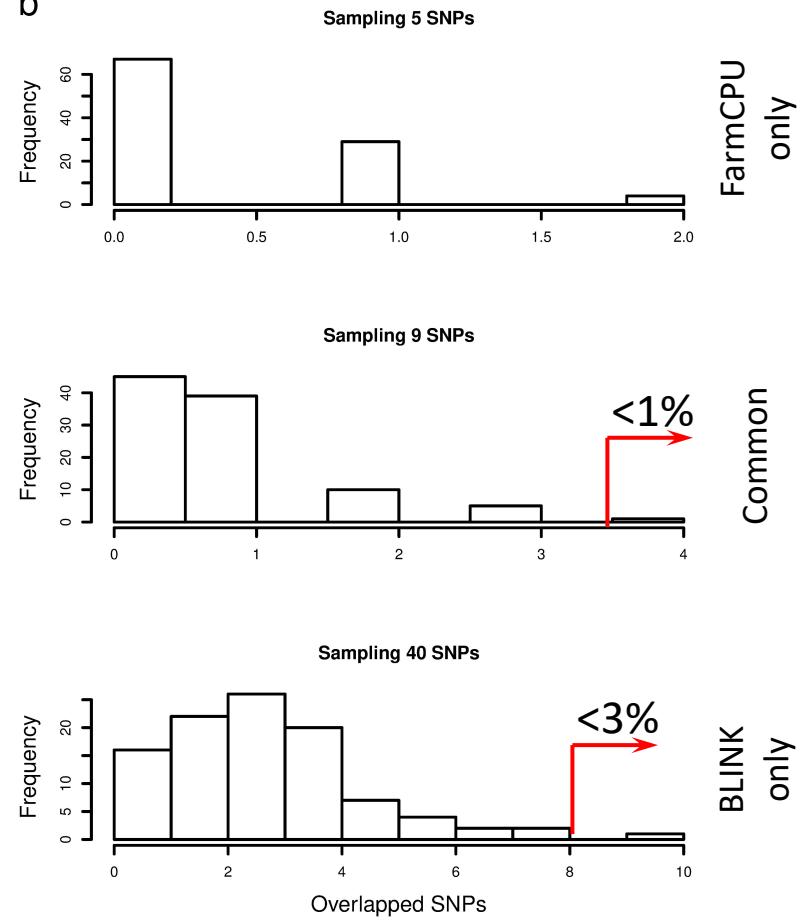


Enrichment

a

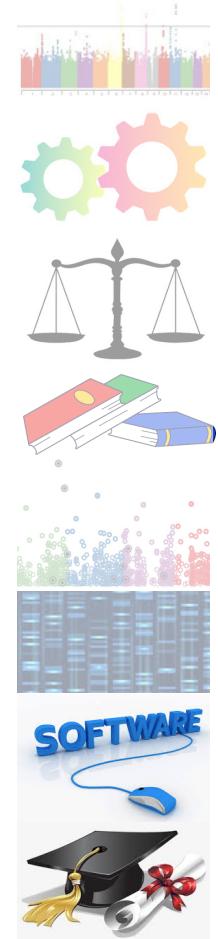


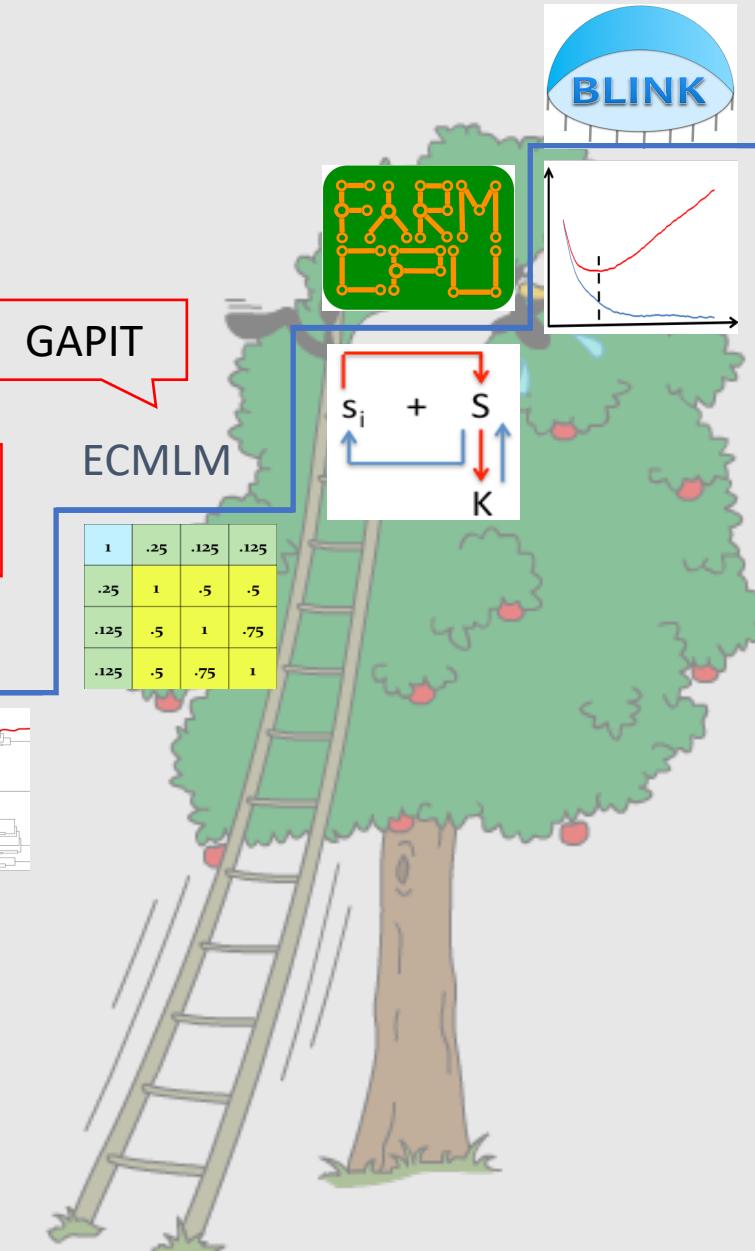
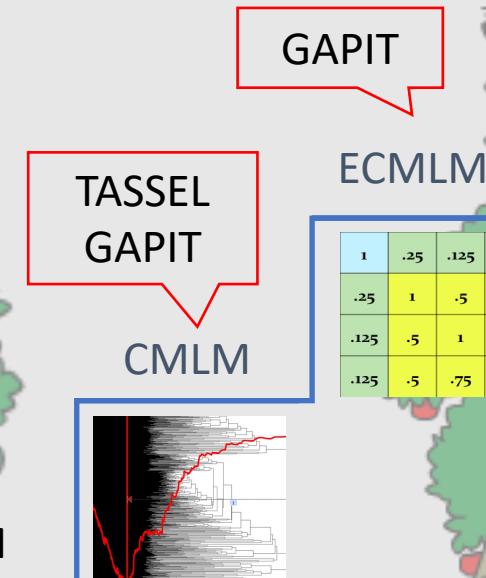
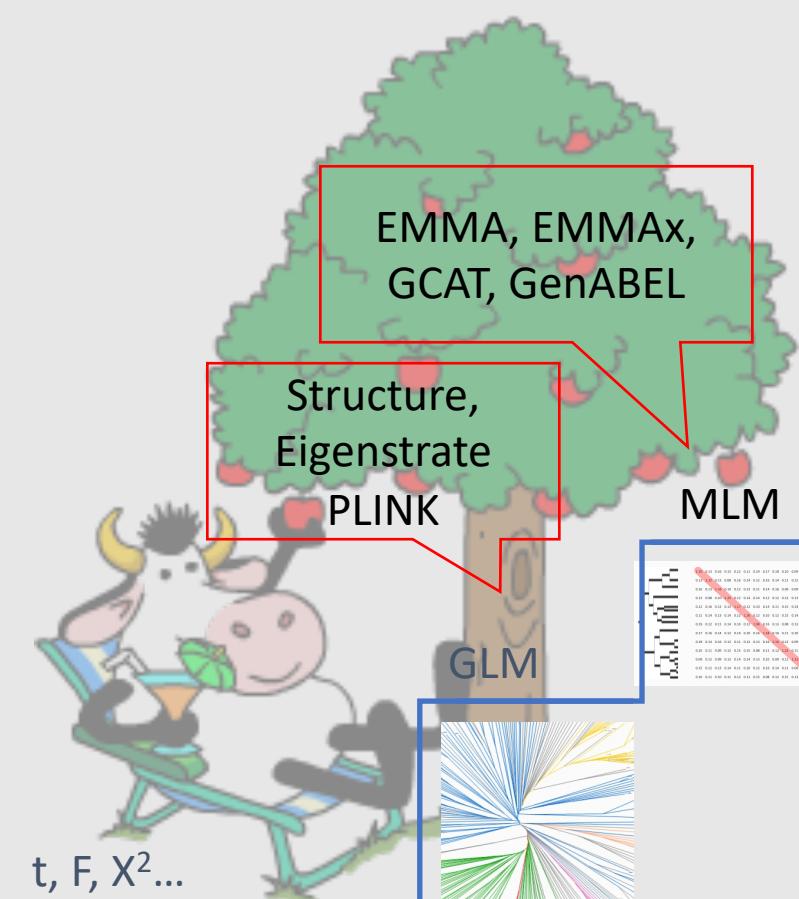
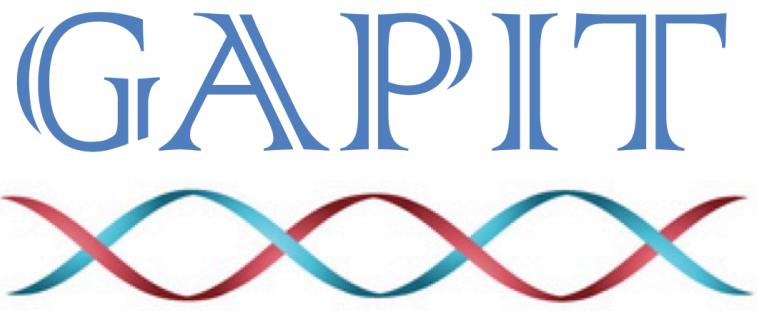
b



Outline

- Why GWAS?
- How does GWAS work?
- How to evaluate GWAS results?
 - Literature
 - Simulation
 - Enrichment analysis
- **GWAS Software**
- **GWAS course at Washington State University**

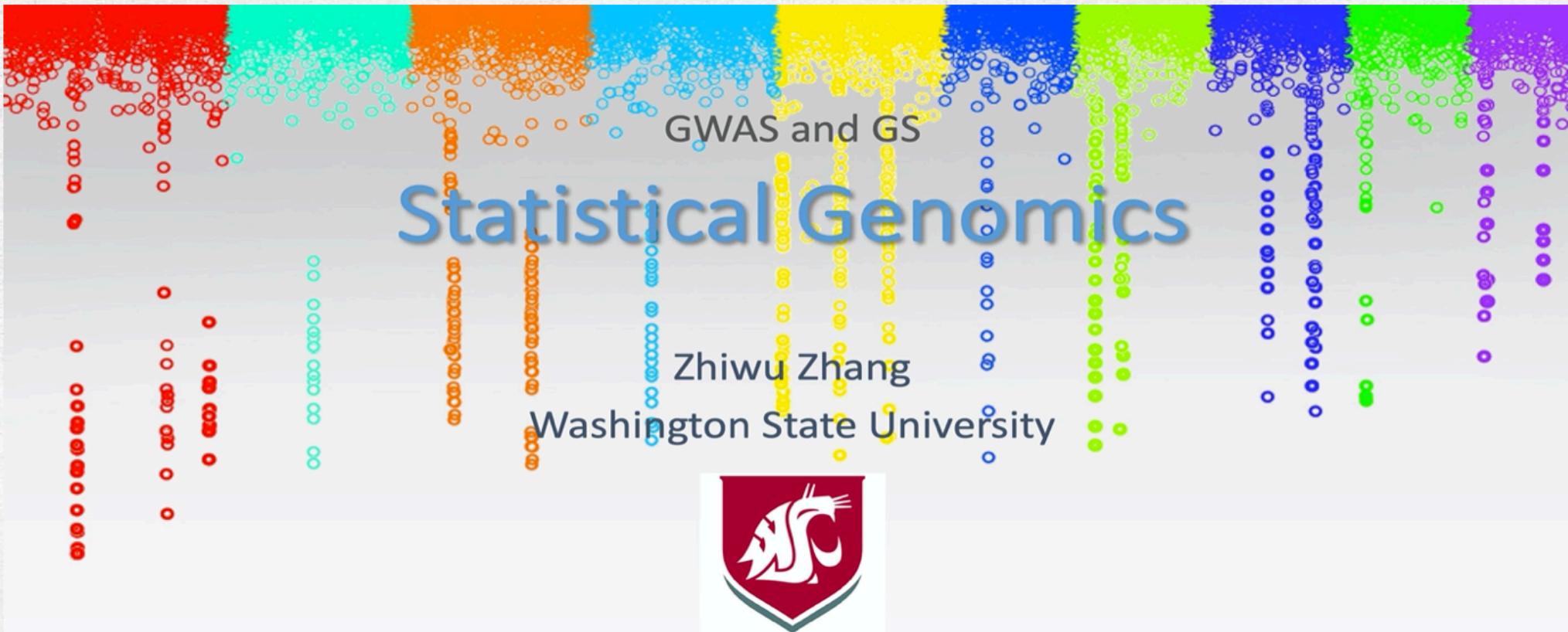




Uncorrelated or
equally correlated



iPat



- Flyer
- Syllabus
- Student evaluation
- Public Genomic Data Resources
- Lecture slides (PPT)
- Source code (R)
- Lab
- Quizzes
- Homework

Offered in spring semesters ([2015](#), [2016](#), [2017](#), [2018](#), [2020](#), and [2021](#)), this graduate course primarily covers Genome Wide Association Study and Genomic Prediction (Selection). The objective is to develop concepts in quantitative genetics and analytical skills in statistics and computation through critical thinking. The course is cross listed by five departments in CAHNRS and College of Art and Science (ANIM_SCI 545, BIOLOGY 545, CROP_SCI 545, HORT 545, and PLP 545). There is no strict prerequisite, however, experience in R programming is strongly recommended. The flyer, syllabus, student evaluation, lecture slides (PPT), and source code (R) are available for all teaching cycles.