# STATISTICAL GENETICS

## Gene Mapping Through Linkage and Association

**Edited by**
**Benjamin M Neale**
**Manuel AR Ferreira**
**Sarah E Medland**
**Danielle Posthuma**

*http://genemapping.org/*

# Epistasis in Association Studies

## David Evans

University of BRISTOL

**e·pis·ta·sis** 🔊)) Audio Help  [i-**pis**-*tuh*-sis] Pronunciation Key – Show IPA Pronunciation

–*noun, plural* –ses 🔊)) Audio Help  [-seez] Pronunciation Key – Show IPA Pronunciation.

1. *Genetics.* a form of interaction between nonallelic genes in which one combination of such genes has a dominant effect over other combinations.

2. *Medicine/Medical.*
   a. the stoppage of a secretion or discharge.
   b. a scum that forms on a urine specimen upon standing.

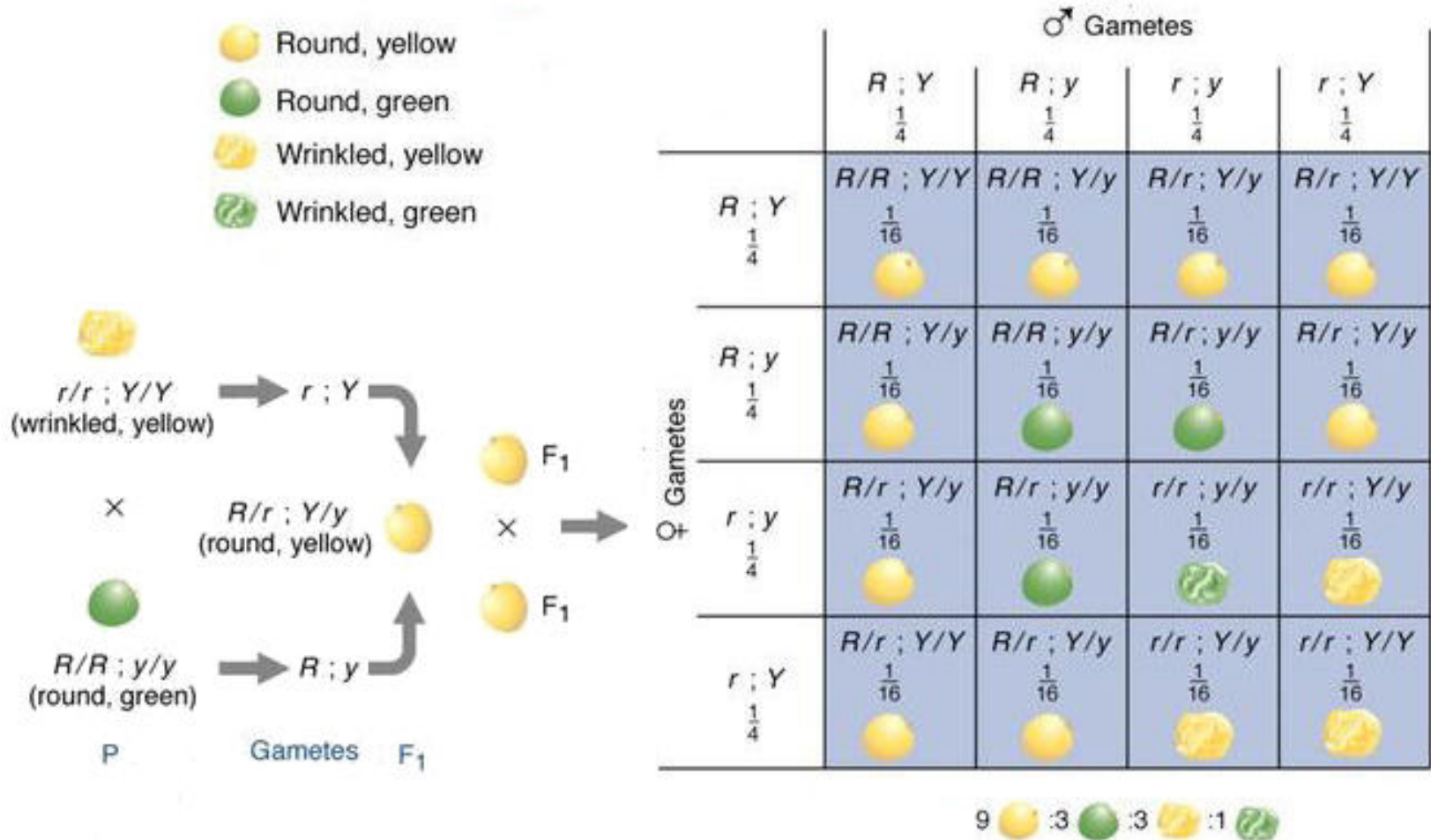[Origin: 1915–20; < Gk *epístasis* stopping, stoppage. See EPI–, STASIS]

—*Related forms*
**ep·i·stat·ic** 🔊)) Audio Help  [ep-*uh*-**stat**-ik] Pronunciation Key – Show IPA Pronunciation, *adjective*
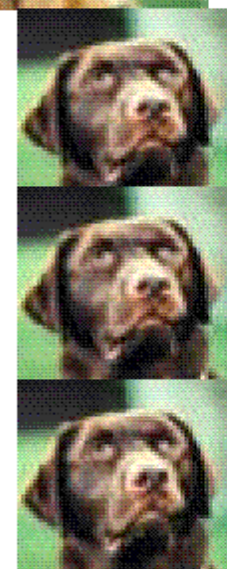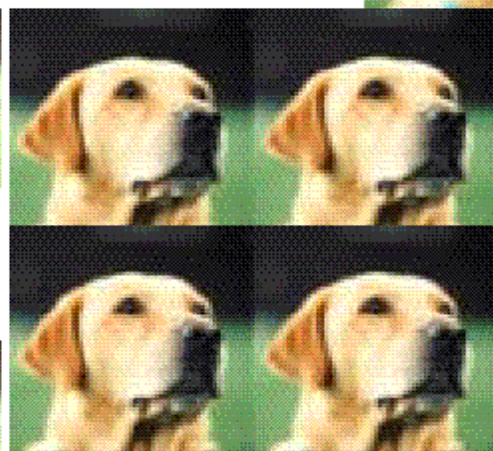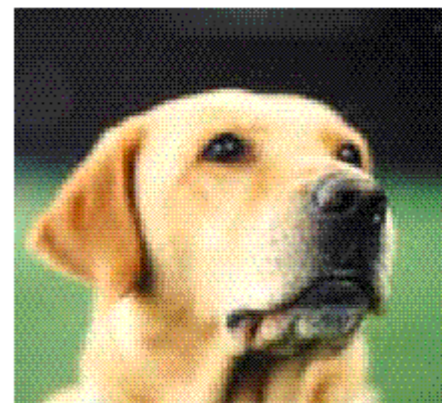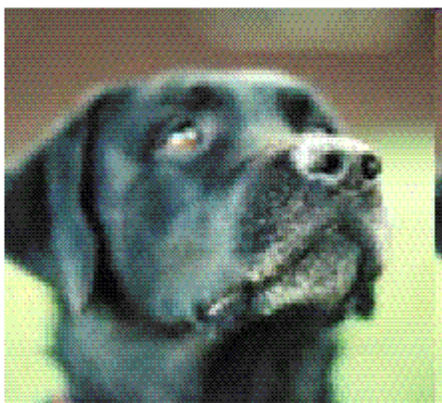
# Law of Independent Assortment

# Biological Epistasis

- Bateson (1909) "a masking effect whereby a variant or allele at one locus prevents the variant at another locus from manifesting its effect…"

- Phenotypic differences among individuals with various genotypes at one locus depend on their genotypes at other loci
  - Does NOT depend on allele frequency

# Epistasis in the labrador retriever dog

# Recessive Epistasis
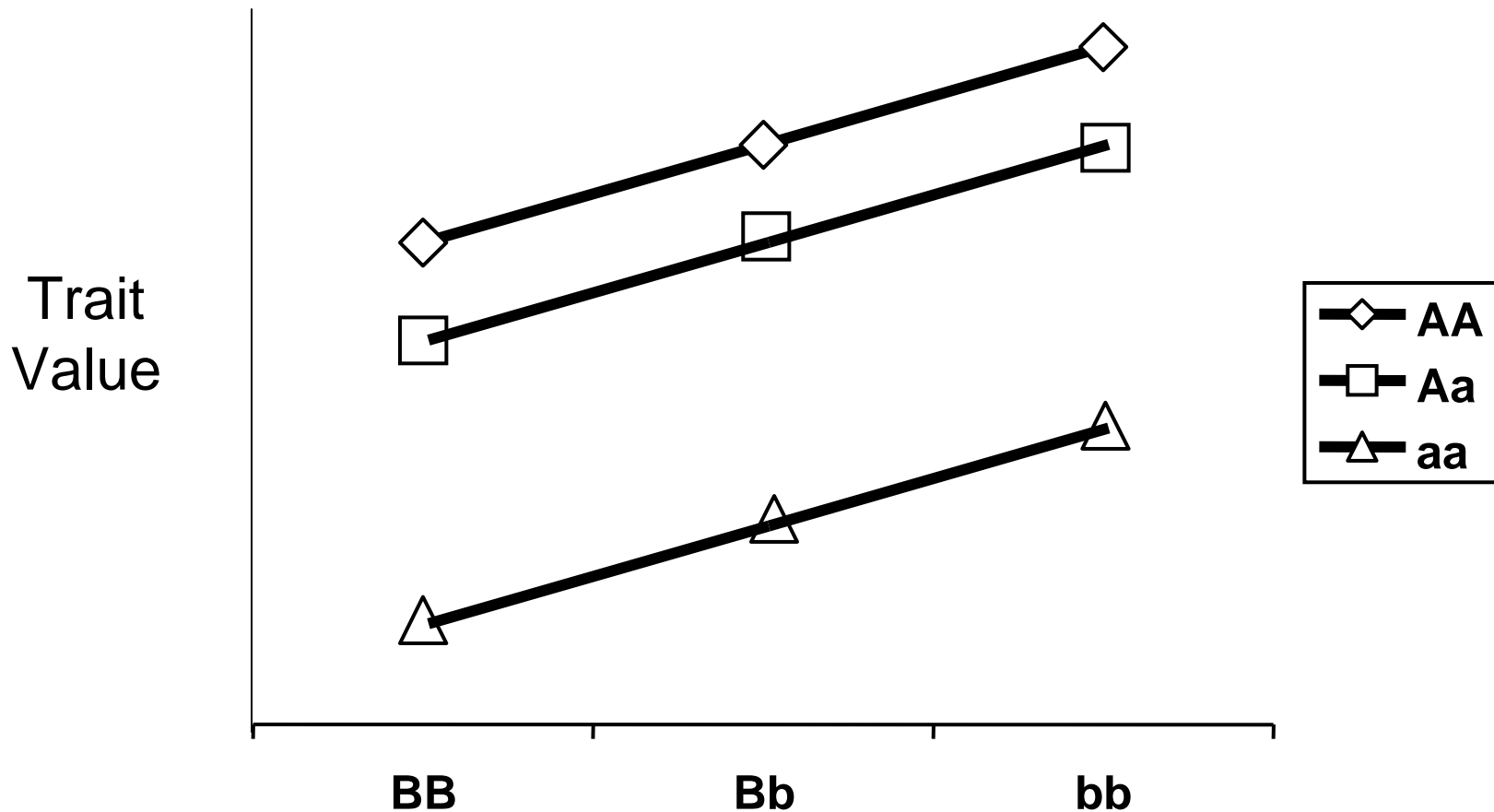
B = Black

b = brown

E/e = gold locus

| ♀\♂ | EB | Eb | eB | eb |
|------|------|------|------|------|
| EB | EEBB | EEBb | EeBB | EeBb |
| Eb | EEBb | EEbb | EeBb | Eebb |
| eB | EeBB | EeBb | eeBb | eeBB |
| eb | EeBb | Eebb | eeBb | eebb |

# Statistical Epistasis

- Deviation of multilocus genotypic values from the additive combination of the single locus components
  - Close to statistical concept of interaction
  - Depends on allele frequencies
  - Population specific
  - May be scale dependent

- Different ways the epistatic values/variance components can be calculated:
  - Hierarchical ANOVA (Sham, 1998)
  - Method of Contrasts (Cockerham, 1954)
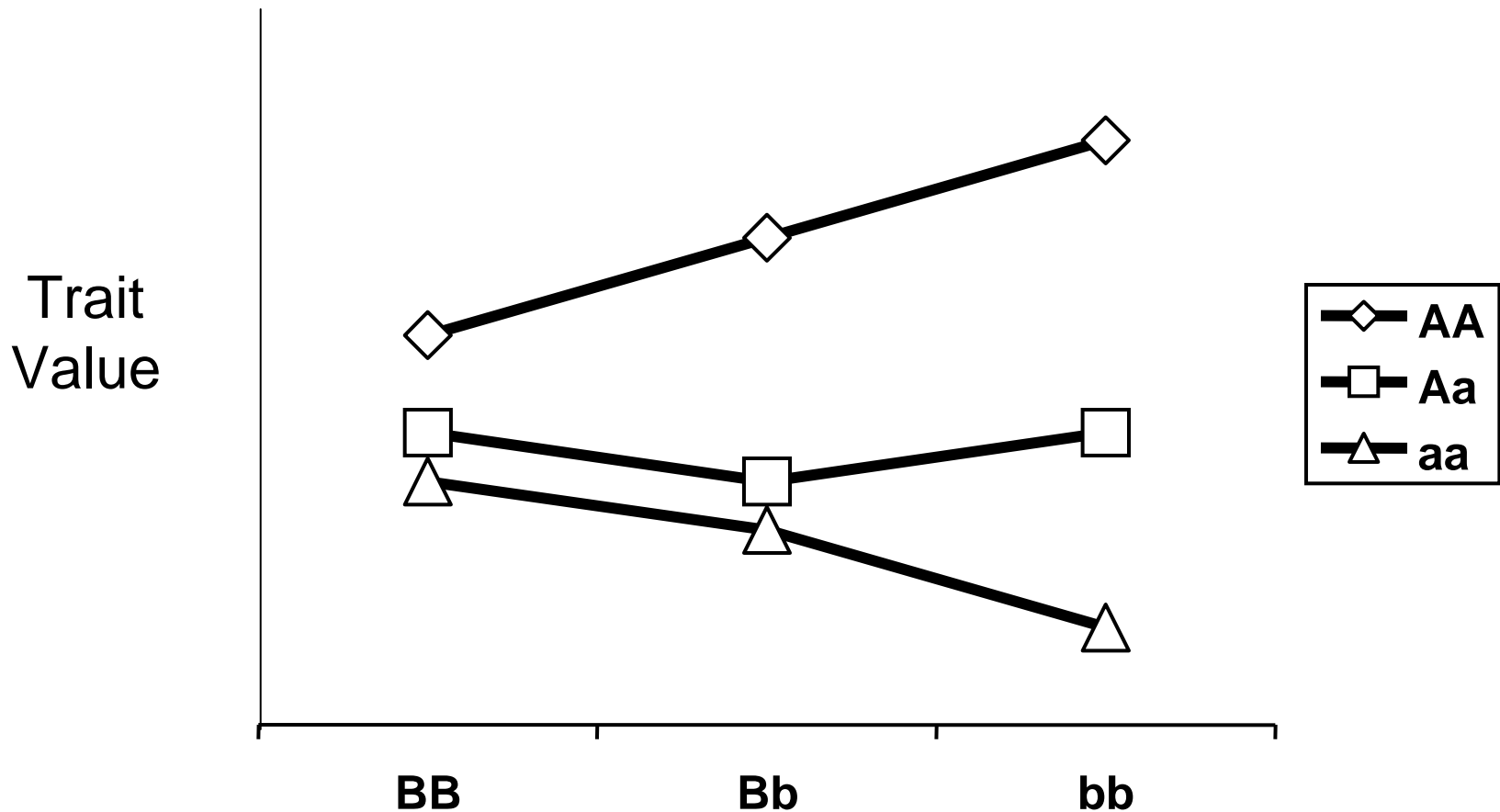  - Using partial derivatives of the population mean (Kojima, 1959; Tiwari & Elston, 1997)

# No Epistasis

▷ Effect at locus B independent of effect at Locus A

# Epistasis Two Loci

▷ Locus A modifies the effect at locus B
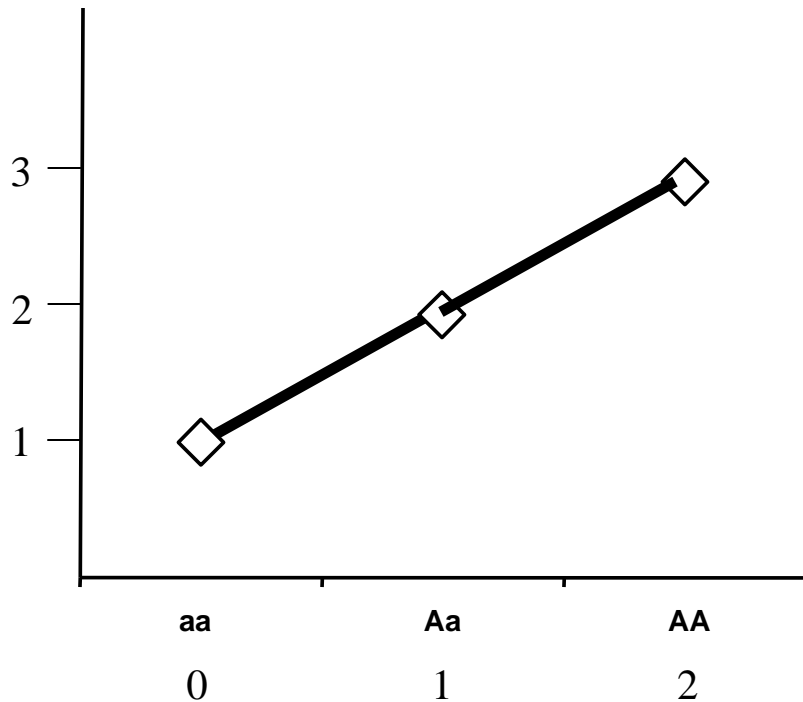
# Genetic Variance – One Locus

|  | AA | Aa | aa |
|---|---|---|---|
| *Genotypic Mean* | $Y_{AA}$ | $Y_{Aa}$ | $Y_{aa}$ |
| *Frequency* | $f_{AA}$ | $f_{Aa}$ | $f_{aa}$ |

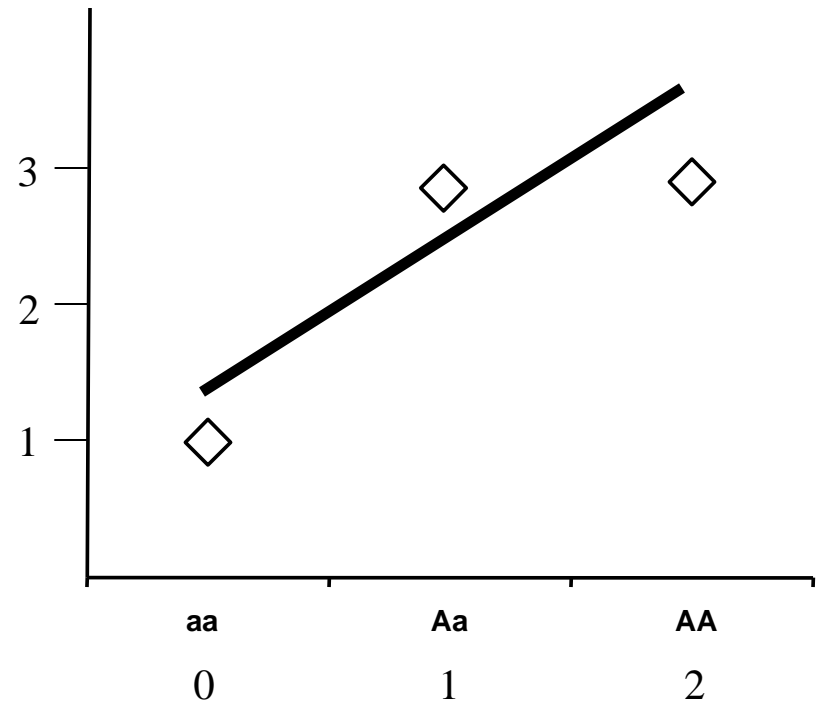$$\mu = \sum_{i=AA,Aa,aa} f_i Y_i \qquad \sigma^2_{A+D} = \sum_{i=AA,Aa,aa} f_i (Y_i - \mu)^2$$

$$\sigma^2_A = 2 f_A f_a [f_A (Y_{AA} - Y_{Aa}) + f_a (Y_{Aa} - Y_{aa})]^2$$

$$\sigma^2_D = f_A^2 f_a^2 (Y_{AA} - 2Y_{Aa} - Y_{aa})^2$$

# Partitioning the Variance One Locus



Additive Model                    Dominant Model

▷ <u>Additive Genetic Variance is variance explained by regression</u>

▷ <u>Dominance variance is residual variance not explained by regression</u>

# Least Squares Regression One Locus

- Represent genotypes of each individual by indicator variables:

| Genotype | Additive Coefficient | Additive and Dom Coefficient | |
|---|---|---|---|
| | $X_1$ | $X_1$ | $Z_1$ |
| aa | -1 | -1 | -½ |
| Aa | 0 | 0 | ½ |
| AA | 1 | 1 | -½ |

$$Y = \mu + a_1 X_1 + d_1 Z_1 + \varepsilon$$

▷ Fit by least squares (or maximum likelihood)

▷ Can provide tests of significance

▷ Partitions data into variance components

# Genetic Variance – Two Loci

|  | AA | Aa | aa |
|---|---|---|---|
| BB | $Y_{AABB}$ $f_{AABB}$ | $Y_{AaBB}$ $f_{AaBB}$ | $Y_{aaBB}$ $f_{aaBB}$ |
| Bb | $Y_{AABb}$ $f_{AABb}$ | $Y_{AaBb}$ $f_{AaBb}$ | $Y_{aaBb}$ $f_{aaBb}$ |
| bb | $Y_{AAbb}$ $f_{AAbb}$ | $Y_{Aabb}$ $f_{Aabb}$ | $Y_{aabb}$ $f_{aabb}$ |

$$A = \{AA, Aa, aa\}$$

$$B = \{BB, Bb, bb\}$$

$$\mu = \sum_{i \in A} \sum_{j \in B} f_{ij} Y_{ij}$$

$$\sigma_G^2 = \sum_{i = A} \sum_{j \in B} f_{ij} (Y_{ij} - \mu)^2$$

$$\sigma_I^2 = \sigma_G^2 - \sigma_{A+D}^2$$

# Components of Variance for a Two Locus Model

Additive genetic variance Locus 1

Additive genetic variance Locus 2
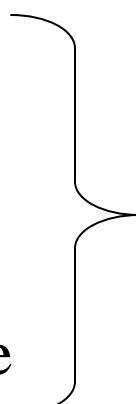
Dominance genetic variance Locus 1

Dominance genetic variance Locus 2

Additive x Additive genetic variance

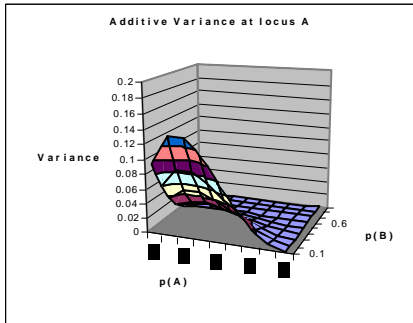Additive x Dominance genetic variance

Dominance x Additive genetic variance
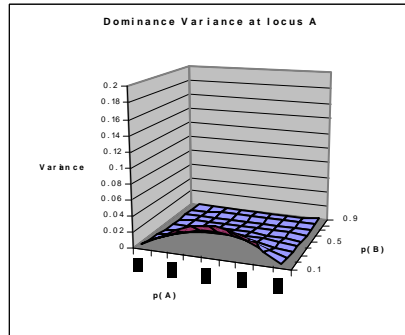
Dominance x Dominance genetic variance
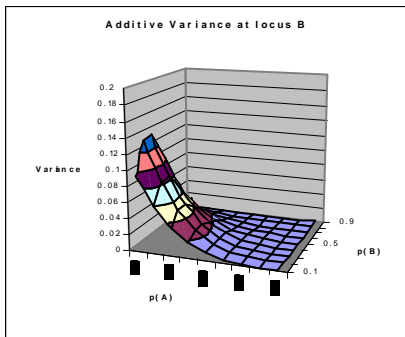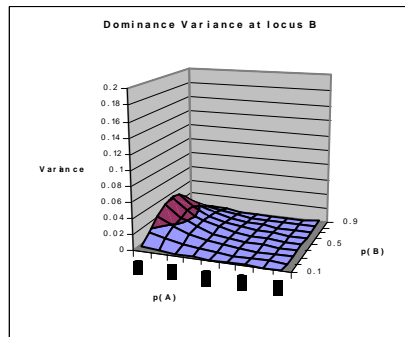
"Epistatic" Variance

## Additive Variance A



Additive Variance at locus A

## Dominance Variance A



Dominance Variance at locus A

|      | BB | Bb | bb |
|------|----|----|----|
| AA   | 0  | 0  | 0  |
| Aa   | 0  | 0  | 0  |
| aa   | 0  | 0  | 1  |

## Additive Variance B



Additive Variance at locus B

## Dominance Variance B



Dominance Variance at locus B

## Additive x Additive



Additive x Additive Variance

## Additive x Dominant



Additive x Dominance Variance

## Dominant x Additive



Additive x Dominance Variance
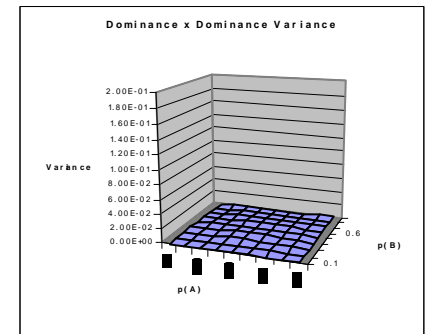
## Dominant x Dominant



Dominance x Dominance Variance

# Least Squares Regression Two Loci

$$Y = \mu + a_1 X_1 + a_2 X_2 + d_1 Z_1 + d_2 Z_2 + i_{aa} w_{aa} + i_{ad} w_{ad} + i_{da} w_{da} + i_{dd} w_{dd} + \varepsilon$$

$$x_1 = \begin{cases} 1 \text{ if } AA \\ 0 \text{ if } Aa \\ -1 \text{ if } aa \end{cases} \qquad x_2 = \begin{cases} 1 \text{ if } BB \\ 0 \text{ if } Bb \\ -1 \text{ if } bb \end{cases}$$

$$z_1 = \begin{cases} \frac{1}{2} \text{ if } Aa \\ -\frac{1}{2} \text{ otherwise} \end{cases} \qquad z_2 = \begin{cases} \frac{1}{2} \text{ if } Bb \\ -\frac{1}{2} \text{ otherwise} \end{cases}$$

$$w_{aa} = x_1 \times x_2 \qquad w_{ad} = x_1 \times z_2 \qquad w_{da} = z_1 \times x_2 \qquad w_{dd} = z_1 \times z_2$$

▷ Can also formulate using logistic regression for dichotomous traits

# Why model Epistasis in GWAS?

▷ <u>Epistasis is important in model organisms</u>

       c.f. Studies in Guinea fowl, yeast, drosophila

▷ <u>Single locus tests will not always detect epistasis !!!</u>

▷ <u>Modeling epistasis may improve power to detect loci (?)</u>

▷ <u>Epistasis is ignored in human studies</u>

       -Main effects hard enough to find!

       -Multiple testing problem: e.g. 100,000 markers gives a cutoff of $p = 1 \times 10^{-11}$ !!!

       -Computational problems

       -Storage problems

# Two-Stage Two-Locus Models in Genome-Wide Association

David M. Evans[1*], Jonathan Marchini[2], Andrew P. Morris[1], Lon R. Cardon[1]

1 Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, United Kingdom, 2 Department of Statistics, University of Oxford, Oxford, United Kingdom

Studies in model organisms suggest that epistasis may play an important role in the etiology of complex diseases and traits in humans. With the era of large-scale genome-wide association studies fast approaching, it is important to quantify whether it will be possible to detect interacting loci using realistic sample sizes in humans and to what extent undetected epistasis will adversely affect power to detect association when single-locus approaches are employed. We therefore investigated the power to detect association for an extensive range of two-locus quantitative trait models that incorporated varying degrees of epistasis. We compared the power to detect association using a single-locus model that ignored interaction effects, a full two-locus model that allowed for interactions, and, most important, two two-stage strategies whereby a subset of loci initially identified using single-locus tests were analyzed using the full two-locus model. Despite the penalty introduced by multiple testing, fitting the full two-locus model performed better than single-locus tests for many of the situations considered, particularly when compared with attempts to detect both individual loci. Using a two-stage strategy reduced the computational burden associated with performing an exhaustive two-locus search across the genome but was not as powerful as the exhaustive search when loci interacted. Two-stage approaches also increased the risk of missing interacting loci that contributed little effect at the margins. Based on our extensive simulations, our results suggest that an exhaustive search involving all pairwise combinations of markers across the genome might provide a useful complement to single-locus scans in identifying interacting loci that contribute to moderate proportions of the phenotypic variance.

# Epistasis in GWAS?

▷ <u>What are the consequences of fitting two locus models when epistasis is absent?</u>

-"If it isn't there, what happens if we go looking for it?"

▷ <u>What are the consequences of fitting two locus models when epistasis is present?</u>

-"If it IS there, what happens if we go looking for it?"

▷ <u>What are the consequences of fitting single locus models when epistasis is present?</u>

-"If it's there, what happens when we ignore it?"

# METHOD

▷ <u>Simulate quantitative variable</u>

-BOTH loci combined are responsible for 10%, 5%, 2% or 1% of total variance

-500, 1000 or 2000 individuals

-Assume 100,000 markers across the genome

-Perfect LD between marker and trait locus

-Comprehensive range of allele frequencies at both loci

-50 different models incorporating varying degrees of epistasis

**MA**

| 1 | ¾ | ½ |
| ¾ | ½ | ¼ |
| ½ | ¼ | 0 |

**M1**

| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 1 |

**M2**

| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 1 | 0 |

**M3**

| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 1 | 1 |

**M5**

| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 1 | 0 | 1 |

**M7**

| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 1 | 1 | 1 |

**M10**

| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 0 | 1 | 0 |

**M11**

| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 0 | 1 | 1 |

**M12**

| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 1 | 0 | 0 |

**M13**

| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 1 | 0 | 1 |

**M14**

| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 1 | 1 | 0 |

**M15**

| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 1 | 1 | 1 |

**M16**

| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 0 |

**M17**

| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |

**M18**

| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 1 | 0 |

**M19**

| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 1 | 1 |

**M21**

| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

**M23**

| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 1 | 1 |

**M26**

| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 0 | 1 | 0 |

**M27**

| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 0 | 1 | 1 |

**M28**

| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 0 |

**M29**

| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |

**M30**

| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

**M40**

| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 0 | 0 |

**M41**

| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 0 | 1 |

**M42**

| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 0 |

**M43**

| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |

**M45**

| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 1 |

**M56**

| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 0 | 0 | 0 |

**M57**

| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 0 | 0 | 1 |

**M58**

| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 0 | 1 | 0 |

**M59**

| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 0 | 1 | 1 |

**M61**

| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 1 | 0 | 1 |

**M68**

| 0 | 0 | 1 |
| 0 | 0 | 0 |
| 1 | 0 | 0 |

**M69**

| 0 | 0 | 1 |
| 0 | 0 | 0 |
| 1 | 0 | 1 |

**M70**

| 0 | 0 | 1 |
| 0 | 0 | 0 |
| 1 | 1 | 0 |

**M78**

| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 1 | 1 | 0 |

**M84**

| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |

**M85**

| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

**M86**

| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 1 | 0 |

**M94**

| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

**M97**

| 0 | 0 | 1 |
| 1 | 0 | 0 |
| 0 | 0 | 1 |

**M98**

| 0 | 0 | 1 |
| 1 | 0 | 0 |
| 0 | 1 | 0 |

**M99**

| 0 | 0 | 1 |
| 1 | 0 | 0 |
| 0 | 1 | 1 |

**M101**

| 0 | 0 | 1 |
| 1 | 0 | 0 |
| 1 | 0 | 1 |

**M106**

| 0 | 0 | 1 |
| 1 | 0 | 1 |
| 0 | 1 | 0 |

**M108**

| 0 | 0 | 1 |
| 1 | 0 | 1 |
| 1 | 0 | 0 |

**M113**

| 0 | 0 | 1 |
| 1 | 1 | 0 |
| 0 | 0 | 1 |

**M114**

| 0 | 0 | 1 |
| 1 | 1 | 0 |
| 0 | 1 | 0 |

**M170**

| 0 | 1 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 0 |

**M186**

| 0 | 1 | 0 |
| 1 | 1 | 1 |
| 0 | 1 | 0 |

Evans et al. (2006)
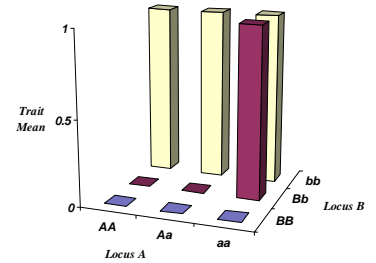*PLOS Genet*

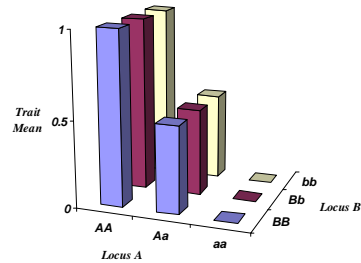**CONTROL**

**LESS FREAKY**

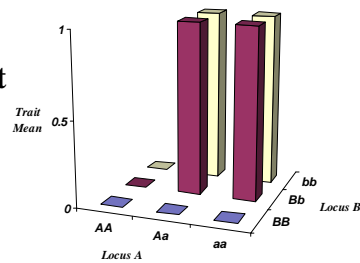**WAY FREAKY**

**Additive**

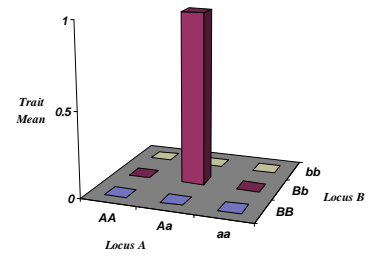Complimentary Gene Action

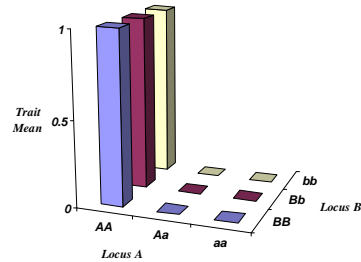**Modifying Effects**

**Additive Locus A**
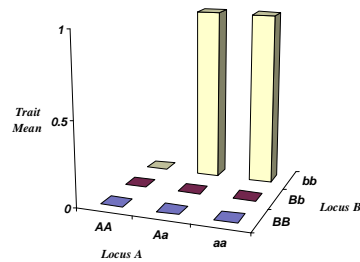
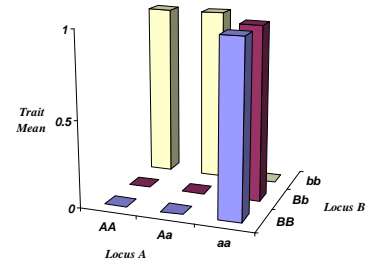Dominant x Dominant Complimentary Gene Action
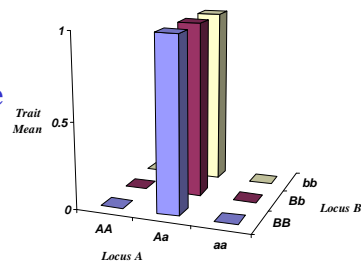
**D x D Epistasis**

**Dominant Locus A**

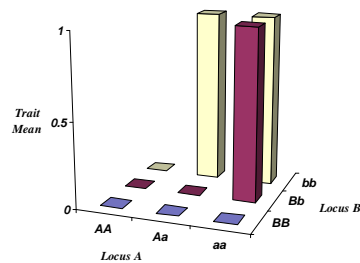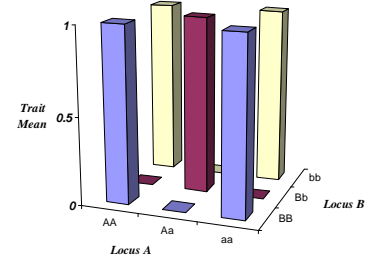Dominant x Recessive Complimentary Gene Action

**XOR**

**Heterozygote Advantage Locus A**

Threshold Model

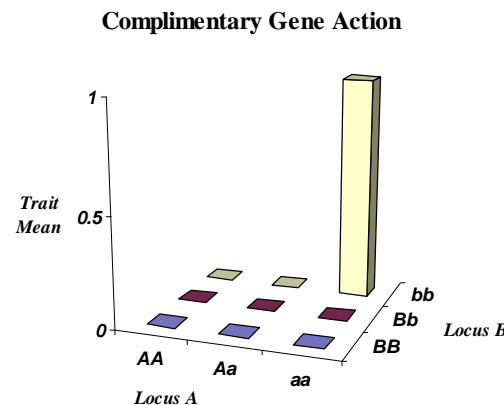**Chequer Board**

# Simulating Genetic Data

▷ <u>Allele frequencies</u>

-Allele frequencies at each locus will determine frequency of genotypes

$p_A = 0.8$

$p_a = 0.2$

→

$p_{AA} = 0.64$

$p_{Aa} = 0.32$

$p_{aa} = 0.04$

▷ <u>Genetic Model</u>

-Need to specify trait means for each genotype combination

**Complimentary Gene Action**



-A random normal deviate can then be placed on these means to simulate the action of the environment

▷ <u>Each combination of parameters will result in a unique variance profile</u>

# Simulating Genetic Data

**Complimentary Gene Action**



*(Variance components shown for complimentary gene action model)*

# METHOD

▷ <u>Quantify power to detect association</u>

    -10,000 simulations for each combination of parameters

▷ <u>Single locus test of association</u>

    -Power to detect BOTH loci

    -Power to detect EITHER locus

▷ <u>Two locus test of association</u>

    -Power to detect BOTH loci

    -Different from power to detect the epistatic variance component explicitly

▷ <u>All models fit via maximum likelihood</u>

    -Significance assessed by minus two log-likelihood chi-square

# Epistasis Isn't There…

**Two Locus Test**

**Power to Detect EITHER Locus
Single Locus Test**

**Power to Detect BOTH Loci
Single Locus Test**



*(Simulations represent 500 individuals, 10% genetic variance, Additive Model)*

▷ The power to detect EITHER locus is greater using the single locus test when epistasis is NOT present

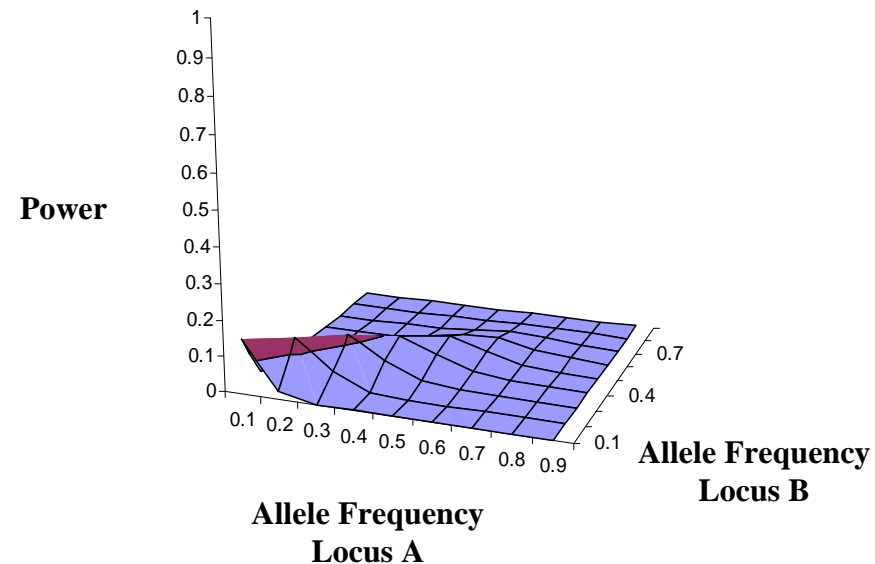▷ The power to detect BOTH loci is actually less using single locus tests even when epistasis is NOT present

# Epistasis IS There…

▷ The power to detect BOTH loci is always better using the two locus test when epistasis is present

**Two Locus Test**

**Power to Detect BOTH Loci
SINGLE Locus Tests**



*(Simulations represent 500 individuals, 10% genetic variance, Dominant x Dominant Complimentary Gene Action Model)*
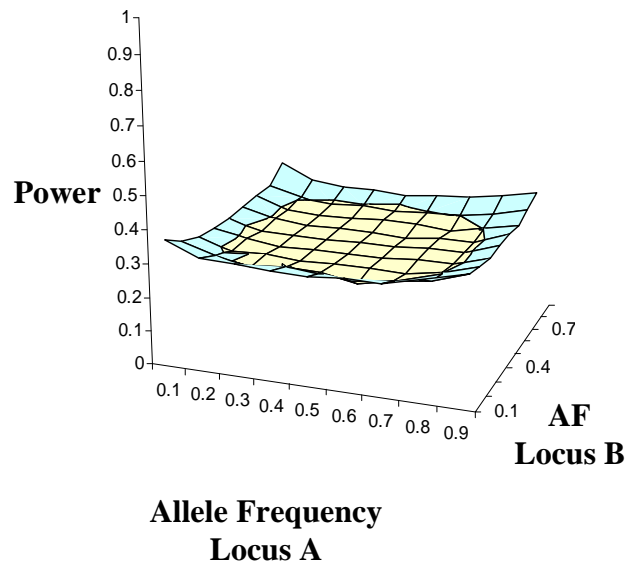
# Power To Detect EITHER Locus

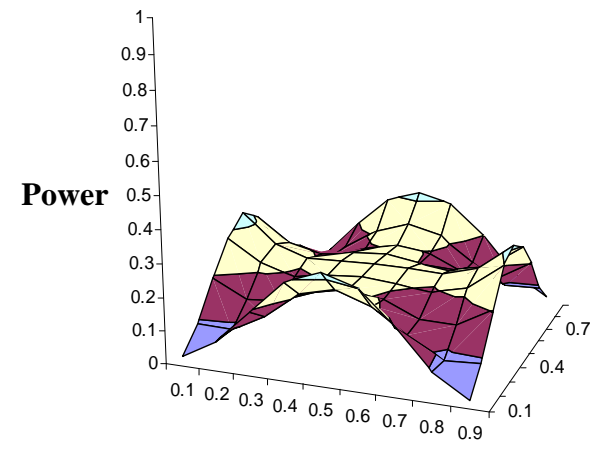▷ When the model is EXTREME, the power of the two locus test is often better than the power to detect EITHER locus using single locus tests

**Model**

**Two Locus Test**

**Power to detect EITHER locus SINGLE Locus Tests**
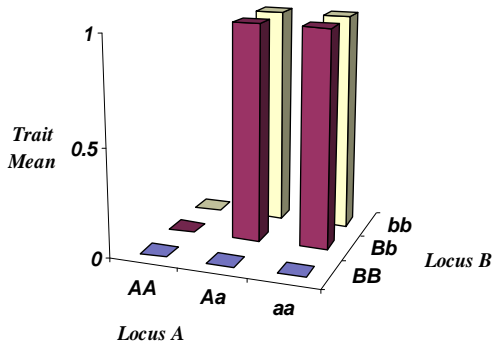
Dominant x Dominant Epistasis

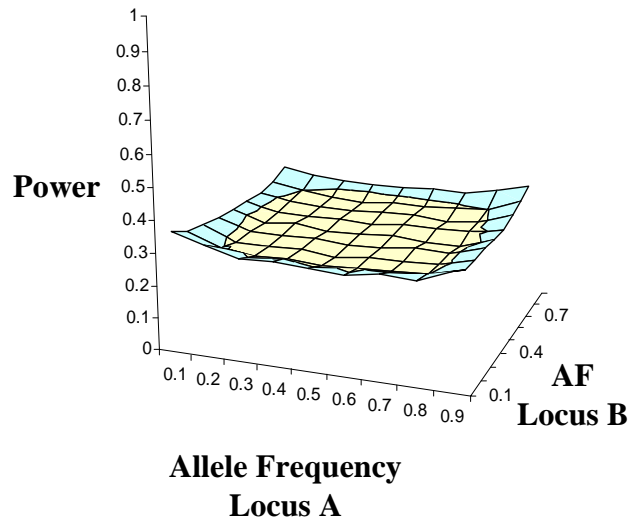*(500 individuals, 10% genetic variance)*

# Power to Detect EITHER Locus

▷ <u>When the model is less extreme, the power to detect EITHER locus using the single locus tests is often better than the two locus test</u>
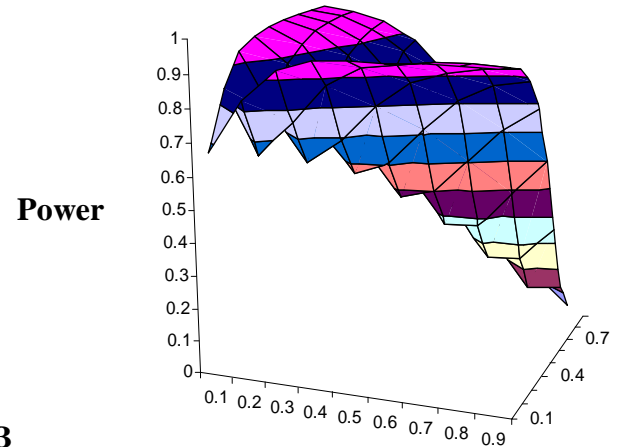
**Model**

**Two Locus Test**

**Power to detect EITHER locus SINGLE Locus Tests**

Dominant x Dominant Complimentary Gene Action Model







*(500 individuals, 10% genetic variance)*

▷ The power to detect BOTH loci is always better fitting the two locus model regardless of the underlying model

▷ There are situations where fitting the full two-locus model will reveal effects which are not identified using single-locus methodology

▷ Multiple testing doesn't kill you as much as you think!!!

# Exhaustive or Two Stage Strategy?

▷ <u>Idea: Two Stage Strategies</u>

　　-Test a subset of $^{100,000}C_2$ comparisons

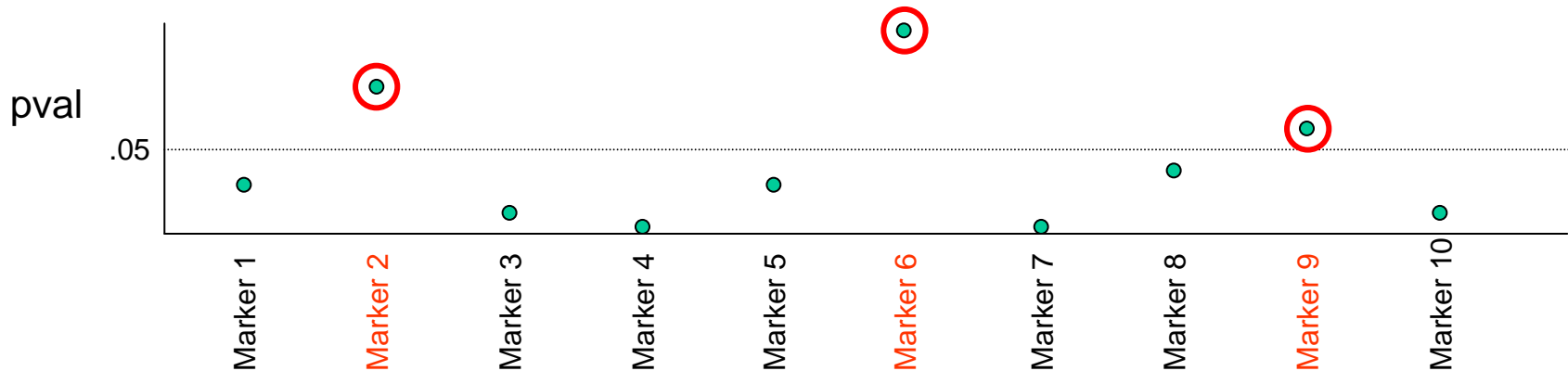　　-Which comparisons are chosen depends on their performance in the single locus
　　 tests

▷ <u>PLUS: Reduce cost due to multiple testing</u>

▷ <u>MINUS: Throw away some comparisons which would be significant
in the two locus test, yet are not significant in the single locus tests</u>
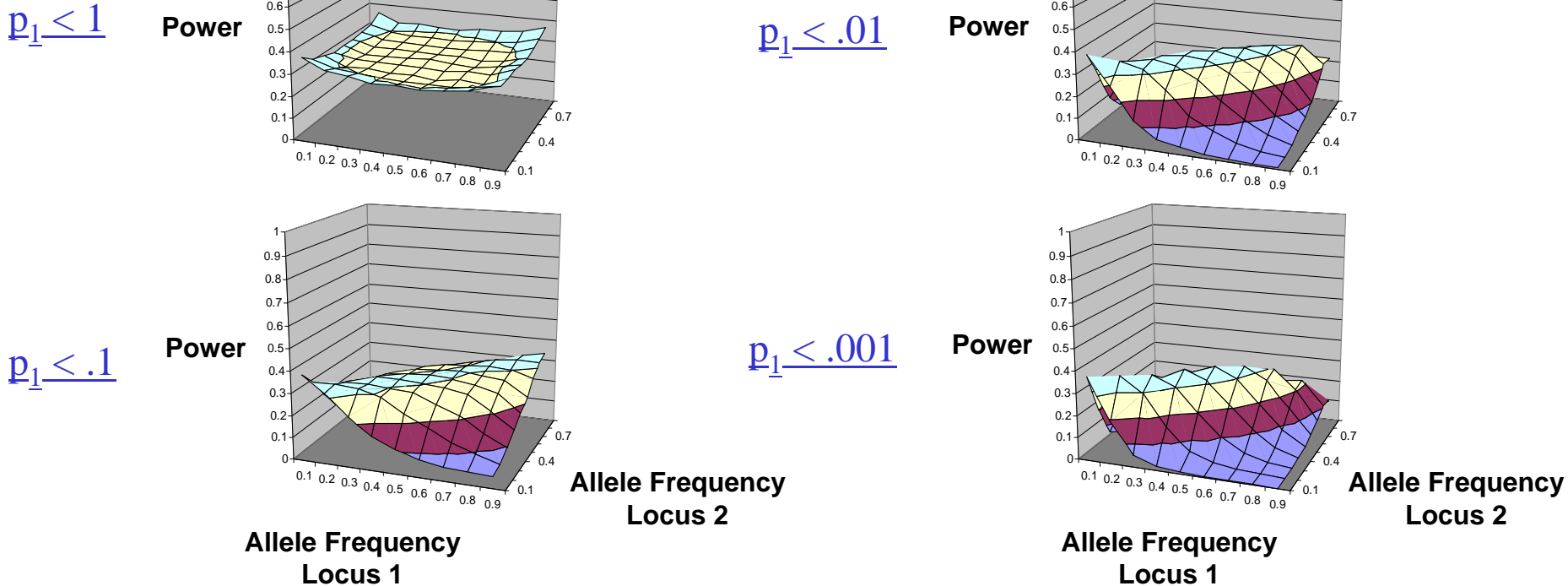
# Two-stage Models

▷ <u>Strategy One</u>

-Only markers which pass first stage threshold in the single locus analysis are tested

-All pair-wise combinations of these markers are tested



-Three comparisons:

    -Marker 2 vs Marker 6

    -Marker 2 vs Marker 9

    -Marker 6 vs Marker 9

# Two Stage Procedure: Strategy One

$p_1 < 1$  **Power**

$p_1 < .01$  **Power**

$p_1 < .1$  **Power**

**Allele Frequency Locus 2**

**Allele Frequency Locus 1**

$p_1 < .001$  **Power**

**Allele Frequency Locus 2**
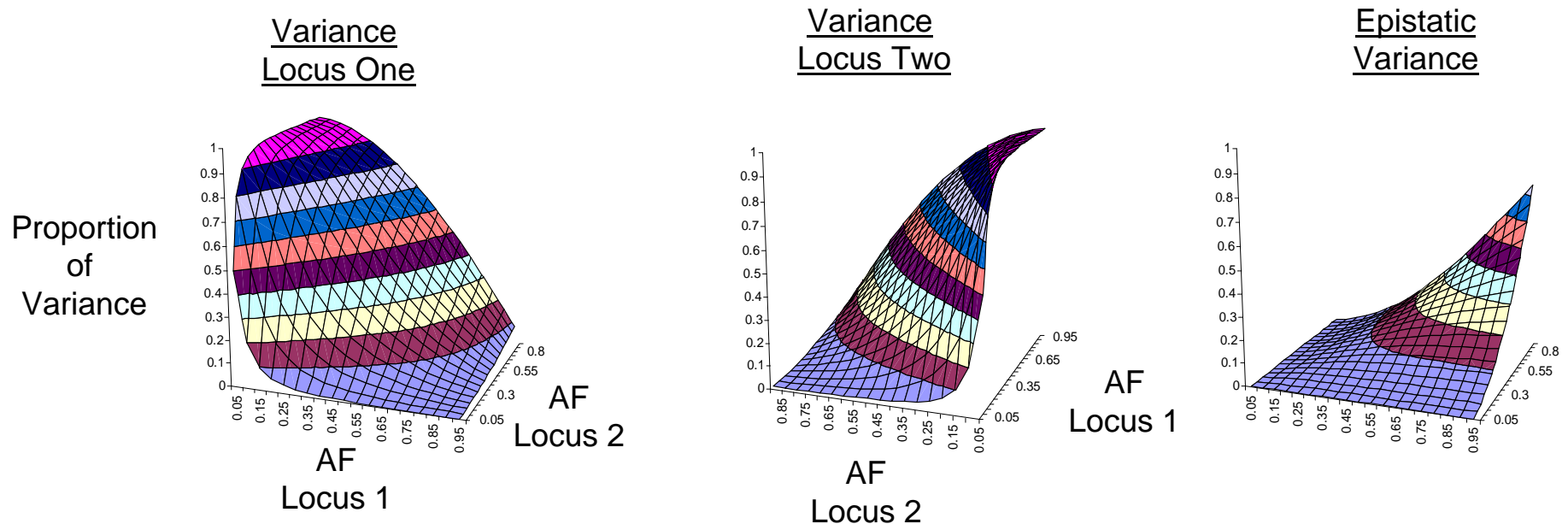
**Allele Frequency Locus 1**

*(Dominant x Dominant Complimentary Gene Action Model, 500 individuals, 10% genetic variance)*

▷ As the first stage threshold becomes more stringent, the power to detect both loci DECREASES for the majority of the parameter space

# Why Does Strategy One Perform Poorly?

▷ For BOTH loci to be included in the second stage, Strategy One requires BOTH single locus tests to meet some threshold

▷ This threshold will not be met when the single locus variance is close to zero

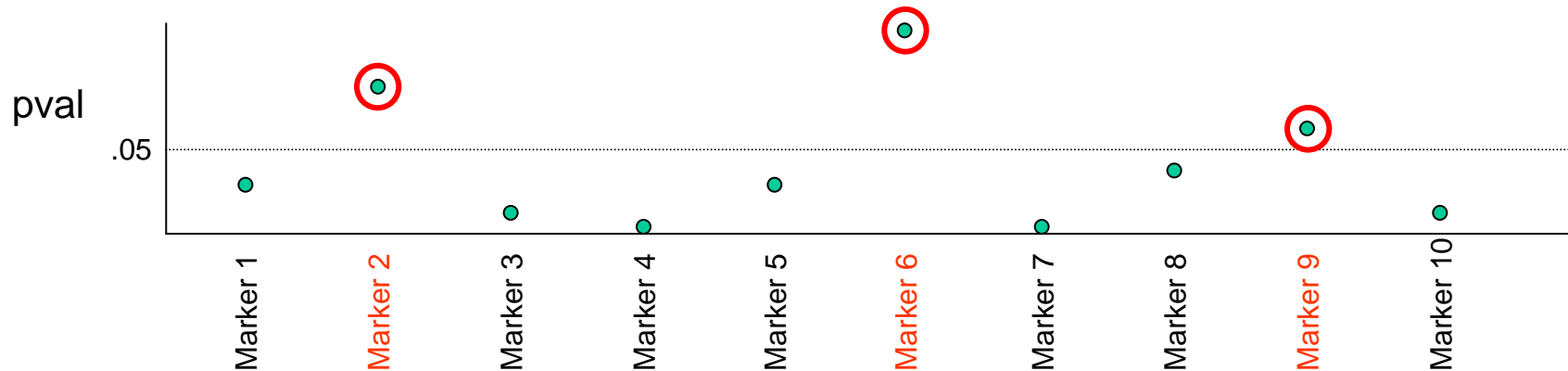▷ Therefore Strategy One will tend to fail whenever EITHER single locus component is close to zero



*(Dominant x Dominant Complimentary Gene Action Model)*

# Two-stage Models

▷ <u>Strategy Two</u>

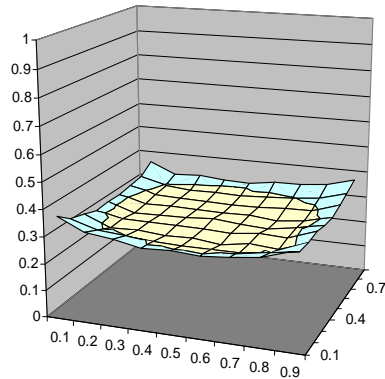-Markers which pass first stage threshold are tested with ALL other markers



-24 comparisons:

    -Marker 2 vs Markers 1, 3, 4, 5, 6, 7, 8, 9, 10

    -Marker 6 vs Marker 1, 3, 4, 5, 6, 7, 8, 9, 10

    -Marker 9 vs Marker 1, 3, 4, 5, 7, 8, 10
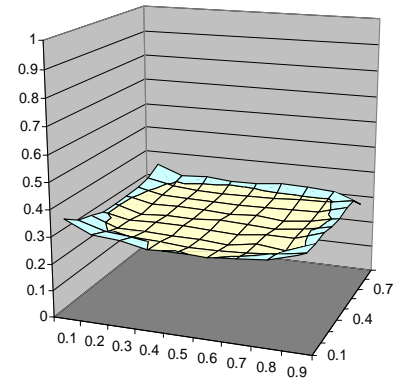
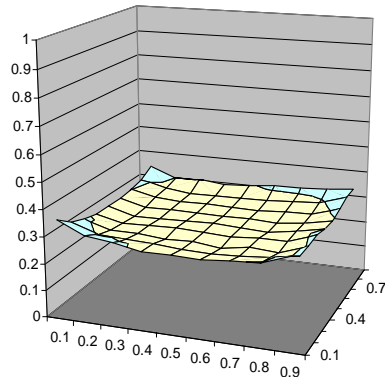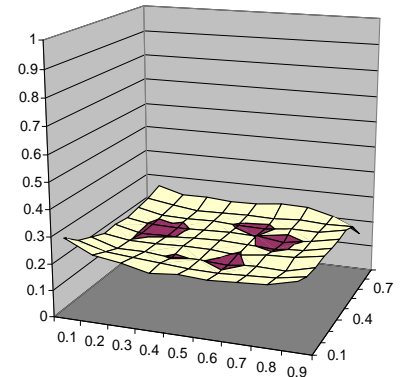# Two-Stage Procedure: Strategy Two

$p_1 < 1$  **Power**

$p_1 < .001$  **Power**

$p_1 < .01$  **Power**

$p_1 < .0001$  **Power**

*(Dominant x Dominant Complimentary Gene Action Model, 500 individuals, 10% genetic variance)*
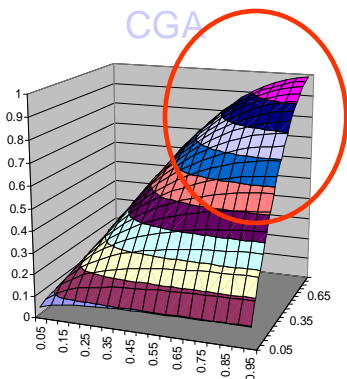
▷ As the first stage threshold becomes more stringent, there is no increase in power

▷ There is a decrease in power at more stringent levels

- Because of the need to condition on the first stage results being significant, there is no increased power in the second stage

- Since loci are included in the second stage if the variance is in EITHER single locus component, the strategy fails when the majority of variance is in the epistatic variance component or the first stage threshold is too severe
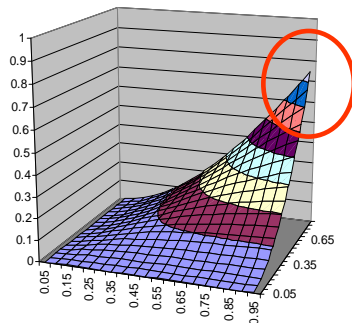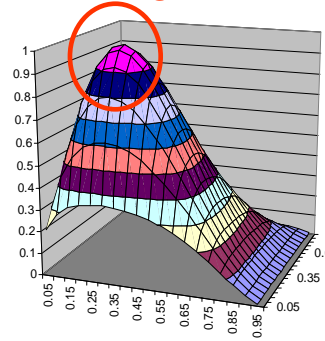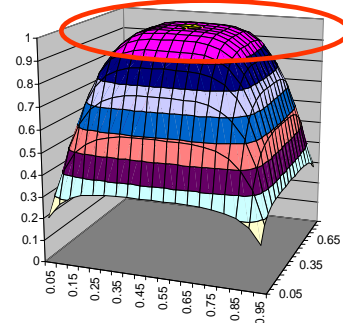
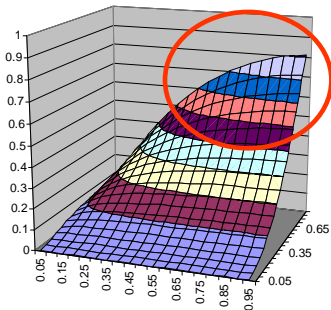Less Extreme Models

Extreme Models



CGA

DxD CGA

XOR
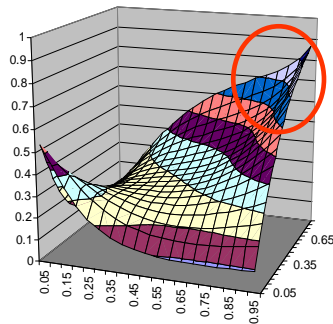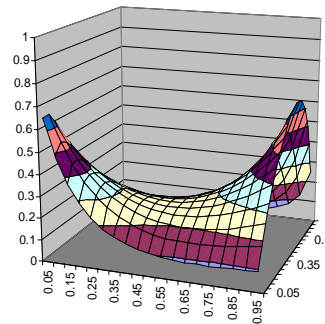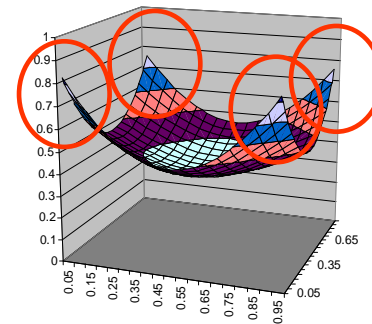
Chequer

DxR CGA

Threshold

ME

DxD

▷ Many simple looking models contain regions of their parameter space where loci would not be able to be identified using single locus analyses or two stage analyses

▷ An exhaustive search involving all pair-wise combinations of markers across the genome is superior to performing a two stage strategy

# Conclusions

▷ An exhaustive search involving all pair-wise combinations of markers across the genome is superior to performing a two stage strategy

▷ Many simple looking models contain sizeable regions of their parameter space where loci would not be able to be identified using single locus analyses

▷ Despite the increased penalty due to multiple testing, it is possible to detect interacting loci which contribute to moderate proportions of the phenotypic variance with realistic sample sizes
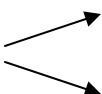
▷ Is it worth incorporating epistasis in GWA?

# Using PLINK to test for Epistasis

▷ Two tests of epistasis are implimented in PLINK

▷ Testing for additive x additive epistasis:

$$Y = \mu + a_1 X_1 + a_2 X_2 + i_{aa} w_{aa} + \varepsilon$$

plink --file *mydata* --epistasis
→ plink.epi.cc
→ plink.epi.cc.summary

▷ "Fast epistasis":

plink --file *mydata* --fast-epistasis

▷ Possible to control output using the flags:

--epi1 0.0001

--epi2 0.0001

# Practical

▷ Copy the files epistasis.ped and epistasis.map from H:/davide/LEUVEN2008

▷ Run single locus tests of association in this dataset

      plink --file *epistasis* --assoc

▷ Run a scan for epistasis using the --fast-epistasis option

      plink --file *epistasis* --fast-epistasis

▷ What are the two top interactions from this analysis?

▷ Are these loci flagged in the single locus analysis?

# Practical

▷ <u>Two locus results:</u>

| CHR1 | SNP1 | CHR2 | SNP2 | STAT | P |
|------|------|------|------|------|---|
| 22 | rs2014410 | 22 | rs9607957 | 15.68 | 7.485e-005 |
| 22 | rs2076672 | 22 | rs2076109 | 15.72 | 7.348e-005 |

▷ <u>Single locus results:</u>

| CHR | SNP | CHISQ | P | OR |
|-----|-----|-------|---|-----|
| 22 | rs2014410 | 0.02691 | 0.8697 | 1.01 |
| 22 | rs9607957 | 3.333 | 0.06791 | 0.7318 |
| 22 | rs2076672 | 3.071 | 0.07969 | 1.122 |
| 22 | rs2076109 | 0.5072 | 0.4763 | 0.9587 |