



HAL
open science

Detecting directional and non-directional epistasis in bi-parental populations using genomic data

Simon Rio, Alain Charcosset, Laurence Moreau, Tristan Mary-Huard

► **To cite this version:**

Simon Rio, Alain Charcosset, Laurence Moreau, Tristan Mary-Huard. Detecting directional and non-directional epistasis in bi-parental populations using genomic data. *Genetics*, 2023, 224 (3), pp.iyad089. 10.1093/genetics/iyad089 . hal-04215609

HAL Id: hal-04215609

<https://hal.science/hal-04215609v1>

Submitted on 26 Sep 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

Detecting directional and non-directional epistasis in bi-parental populations using genomic data

Simon Rio^{1,2,✉}, Alain Charcosset³, Laurence Moreau³, and Tristan Mary-Huard^{3,4}

¹CIRAD, UMR AGAP, F-34398 Montpellier, France

²AGAP, Université de Montpellier, CIRAD, INRAE, Institut Agro, Montpellier, France

³Université Paris-Saclay, INRAE, CNRS, AgroParisTech, GQE-Le Moulon, Gif-sur-Yvette, France

⁴MIA Paris Saclay, INRAE, AgroParisTech, Université Paris-Saclay, Palaiseau, France

✉simon.rio@cirad.fr

Abstract

Epistasis, commonly defined as interaction effects between alleles of different loci, is an important genetic component of the variation of phenotypic traits in natural and breeding populations. In addition to its impact on variance, epistasis can also affect the expected performance of a population and is then referred to as directional epistasis. Before the advent of genomic data, the existence of epistasis (both directional and non-directional) was investigated based on complex and expensive mating schemes involving several generations evaluated for a trait of interest. In this study, we propose a methodology to detect the presence of epistasis based on simple inbred bi-parental populations, both genotyped and phenotyped, ideally along with their parents. Thanks to genomic data, parental proportions as well as shared parental proportions between inbred individuals can be estimated. They allow the evaluation of the directionality of epistasis through a test of the expected performance and/or the variance of genetic values. This methodology was applied to two large multi-parental populations, i.e., the American maize and soybean nested association mapping populations, evaluated for different traits. Results showed significant epistasis, especially for the test of directional epistasis, e.g., the increase in anthesis to silking interval observed in most maize inbred progenies or the decrease in grain yield observed in several soybean inbred progenies. In general, the effects detected suggested that shuffling allelic associations of both elite parents had a detrimental effect on the performance of their progeny. This methodology is implemented in the EpiTest R-package and can be applied to any inbred bi-/multi-parental population evaluated for a trait of interest.

Author summary

The genetic architecture of complex traits involves the actions of several loci whose allele effects can depend on the presence of specific alleles at other loci. This phenomenon is commonly referred to as epistasis and can affect both the variation and the average performance of a population. We propose a new methodology to detect the presence of epistasis in bi-parental inbred progeny, a very common type of population in plant genetics. It relies on the evaluation of the progeny for a trait of interest as well as on its characterization for genomic variants. Epistasis affecting the phenotypic variation of the trait and/or in the average performance of the bi-parental population can be detected. We applied this methodology to large multi-parental populations of maize and soybean and identified the presence of epistasis in several traits, affecting in particular in the average performance of the progeny. For instance, soybean populations often displayed lower grain yield than expected based on the performance of their elite parents. This methodology will help geneticists to better characterize the role of epistasis in the trait genetic architectures of their species of interest.

Introduction

The term "epistasis" was first introduced by Bateson (1909) in the context of discrete traits with Mendelian segregation to describe the phenomena where the presence of an allele at one locus masks the effect of another locus. The extension of this concept to the context of populations of individuals evaluated for quantitative traits was presented by Fisher (1918) and is commonly referred to as "statistical epistasis". In this context, epistasis characterizes statistical deviations between loci that arise after taking into account the effect of each locus independently (see Phillips (2008) for a review of the different views on epistasis).

Epistasis is prevalent in the genetic architecture of quantitative traits and it arises from the complex transcriptional, metabolic and biochemical networks involved in these traits (Kryazhimskiy, 2021). Despite its prevalence, epistasis has long been considered a nuisance parameter that can be ignored in breeding (Crow, 2010). This can be explained by the limited transmission of the epistatic part of the genotypic value between parents and offspring. However, the contribution of epistasis to the genotypic value of a given commercial cultivar can be substantial, highlighting the importance of characterizing and predicting the epistatic component in plant breeding (Raffo et al., 2022; Varona et al., 2018).

To detect the presence of epistatic interactions, geneticists have designed complex crossing schemes, as summarized in Mather and Jinks (1982). These schemes required the production of several generations of progeny evaluated for the studied trait. The general principle consists of the comparison of the expected performance of the different generations allowing the isolation of epistatic terms whose significance can be tested. Examples include the well-known triple test cross (Jinks et al., 1969; Kearsey and Jinks, 1968) and other designs (Chahal and Jinks, 1978; Hayman, 1958; Melchinger, 1987). In the triple test cross, a sample of parents i are crossed to two inbred testers (with mean progeny performance L_{1i} and L_{2i}) as well as to their F1 hybrid (with mean progeny performance L_{3i}). The ability of such approaches to detect epistasis relies on a set of non-cancelling epistatic genetic effects in the comparison of expected generation performances, leading to so-called directional epistasis (see Rouzic (2014) for a review).

An alternative for epistasis detection is to consider variance rather than expected performance. Cockerham (1954) and Kempthorne (1954) proposed a partitioning of the genetic variance into orthogonal components: additive, dominance, additive \times additive, additive \times dominance, and higher order interactions. Since genetic variances are quadratic functions of genetic effects, the canceling of effects possibly observed in the expected performance is prevented. Before the use of molecular markers, the estimation of epistatic variances in heterozygous populations was very complex and reduced to a handful of sophisticated designs such as double cross hybrids (Rawlings and Cockerham, 1962a) or triallel analysis (Rawlings and Cockerham, 1962b). In the plant community, homozygous inbred individuals can be generated such as double-haploids (DH) or recombinant inbred lines (RILs). As the use of inbred progenies circumvents the difficulty of disentangling dominance from epistasis encountered in most designs, strategies were proposed to estimate additive and epistatic variance components. Choo (1981) proposed a total genetic variance partitioning into across and within F2-derived DH populations, which can be solved to estimate additive and epistatic variance components. Other strategies were proposed based on diallel crosses (Choo, 1980) or random mating populations (Gallais, 1990).

The advent of molecular markers in the late 1980s revolutionized approaches to detect epistatic interactions. The identification of quantitative trait loci (QTL) involved in the genetic architecture of traits opened the way to the identification of epistasis using (i) one-dimensional approaches, i.e. by testing the interaction between a QTL and the genetic background (Blanc et al., 2006; Jannink and Jansen, 2001; Jannink et al., 2009), or (ii) two-dimensional approaches, i.e. by testing the interaction between pairs of QTL (Jannink et al., 2009; Kao et al., 1999). In the context of genome-wide associations studies, similar one-dimensional (Crawford et al., 2017; Jannink, 2007; Rio et al., 2020a) and two-dimensional (Hemani et al., 2011; Prabhu and Pe'er, 2012; Zhang et al., 2010) approaches were proposed. Molecular markers also make it possible to calculate genomic relationship matrices corresponding to the orthogonal partitioning of genetic variance (Álvarez-Castro and Carlborg, 2007;

Vitezica et al., 2017), thus enabling to estimate epistatic variance components without the need for dedicated designs.

Inbred bi-parental populations, including Double Haploid (DH) or recombinant inbred lines (RIL) progenies, are pervasive in the plant genetics community. They are the cornerstone of breeding programs for self-pollinated species like wheat and soybean, as well as cross-pollinated species based on F1 hybrids between pure inbred lines like maize. They are also a fundamental component of genetic mapping studies based on single or multi-parental designs like nested association mapping (NAM) populations that have been generated for a large number of species over the last decade (Gage et al., 2020). With decreasing costs of genotyping, inbred bi-parental population datasets, including both genomic and phenotypic information, are becoming routinely available in most crops.

In this study, we present a framework to test for the existence of epistasis in inbred bi-parental populations genotyped and evaluated for a trait, both through the expectation (i.e., directional epistasis) and the variance (i.e., non-directional epistasis) of genetic values. This framework is implemented in the new R-package "EpiTest" available from the CRAN. Applications are presented to two large NAM populations (the American maize NAM and the soybean NAM) evaluated for agronomy, phenology, morphology and/or quality traits.

Results

Testing for epistasis in inbred bi-parental progeny

The test procedure is based on the following model for phenotypes Y_i of an inbred bi-parental progeny:

$$Y_i = \mu + \beta\pi_i + \delta\pi_i^2 + G_i^S + G_i^{S \times S} + E_i$$

The fixed part of the model includes an intercept μ , a parameter β associated with the linear relationship between the phenotype of an individual Y_i and its proportion of alleles from the alternative parent π_i , and a directional epistatic component δ associated with the quadratic relationship between Y_i and π_i . The random part of the model includes a segregation genetic term G_i^S , an epistatic segregation \times segregation genetic term $G_i^{S \times S}$ whose variance parameter $\sigma_{S \times S}^2$ is non-null in presence of epistasis, and an error term E_i .

Different tests can be proposed whose null hypothesis are:

- $H_0 : \delta = 0$ to test the existence of directional epistasis in the fixed part of the model,
- $H_0 : \sigma_{S \times S}^2 = 0$ to test the contribution of epistasis to the genetic variance,
- $H_0 : \delta = 0$ and $\sigma_{S \times S}^2 = 0$ to test jointly the epistatic terms in the fixed and the variance part of the model.

Using simulations (S1 File), the ability of the model to estimate parameters accurately was confirmed. The proposed test procedures were found to efficiently control the number of False Positive, and exhibited substantial statistical power in particular for the global and the directional epistasis tests. The superiority in statistical power of the latter two tests over the variance test was also confirmed by the total number of significant tests for real data (Table 1), with 108, 99 and 12 significant tests for the global, directional and variance tests, respectively.

Soybean NAM

The method was first applied to the soybean NAM that includes 39 bi-parental progenies evaluated for grain yield, days to maturity, plant height, lodging, grain moisture, and protein/oil/fiber content.

The presence of epistasis was first investigated using the global test (Fig 1B, Table 1, see the S1 Table for details on model parameters estimates and tests). Significant epistasis was detected for all traits, e.g. in 14 out of 39 populations for grain yield. Some bi-parental populations never showed significant epistasis like the IA3023 \times NE3001 population. In contrast, some families showed

Table 1: Number of significant tests based on a Bonferroni threshold of 5% over the number of populations for the global test (Global, δ & $\sigma_{S \times S}^2$), the directional epistasis test (Dir., δ) and the epistatic variance test (Var., $\sigma_{S \times S}^2$) using the model of Eq. (5)

Dataset	Trait	Global (δ & $\sigma_{S \times S}^2$)	Dir. (δ)	Var. ($\sigma_{S \times S}^2$)
Soybean	Plant height	8	3	3
Soybean	Days to maturity	4	0	4
Soybean	Lodging	3	2	1
Soybean	Grain moisture	6	8	0
Soybean	Grain yield	14	15	1
Soybean	Protein content	1	1	0
Soybean	Oil content	4	3	2
Soybean	Fiber content	1	1	0
Maize	Plant height	10	9	0
Maize	Ear height	13	13	1
Maize	Days to silking	13	12	0
Maize	Days to anthesis	16	16	0
Maize	Anthesis to silking interval	15	16	0
Total over datasets and traits		108	99	12

significant epistasis for most traits (e.g., the IA3023×PI574486 population that exhibits significant tests for plant height, maturity, grain moisture, grain yield and oil content).

The presence of directional epistasis was investigated using the directional epistasis test (Table 1, see the S1 Table for details on model parameters estimates and tests). A focus was done on plant height (Fig 2B), grain yield (Fig 2D), and oil content (Fig 2F), results for the other traits are shown in supplementary material (S1B, S2B, S3B, S4B and S5B Fig). Significant directional epistasis was found for 3, 15 and 3 populations out of 39 for plant height, grain yield and oil content, respectively. For grain yield, significant quadratic coefficients δ were all positive, implying that progenies had lower average performances than expected under a model without directional epistasis. Using model parameter estimates presented in the S1 Table, expected performances can be compared between parents and progeny by considering a reference progeny with perfectly balanced parental ancestry proportions of $\pi_i = 0.5$. For instance, for the population IA3023×TN05.3027, IA3023 has an expected performance of 3937.2 kg/ha, TN05.3027 has an expected performance of 3892.3 kg/ha, and the expected performance of their progeny (with perfectly balanced ancestry proportions) is below that of the two parents: 3640.7 kg/ha, indicating the presence of directional epistasis.

The contribution of epistasis to the genetic variance in the progeny was investigated using the test of the variance part of the model (Table 1, see S1 Table for details on model parameters estimates and tests). A focus was done on the same traits as for directional epistasis (see Fig 2A, 2C and 2E), results for the other traits are shown in supplementary material (S1A, S2A, S3A, S4A and S5A Fig). In contrast to the global and directional epistasis tests, only a few epistatic variance tests were significant. Most traits showed a limited contribution of the epistatic variance component ($\sigma_{S \times S}^2$) to the genetic variance, with notable exceptions such as with the IA3023×PI437169B population for grain yield.

American maize NAM

The method was then applied to the American maize NAM that includes 25 bi-parental progenies evaluated for days to anthesis/silking, anthesis to silking interval, and ear/plant height.

Like for the soybean NAM, the presence of epistasis was first investigated using the global test (Fig

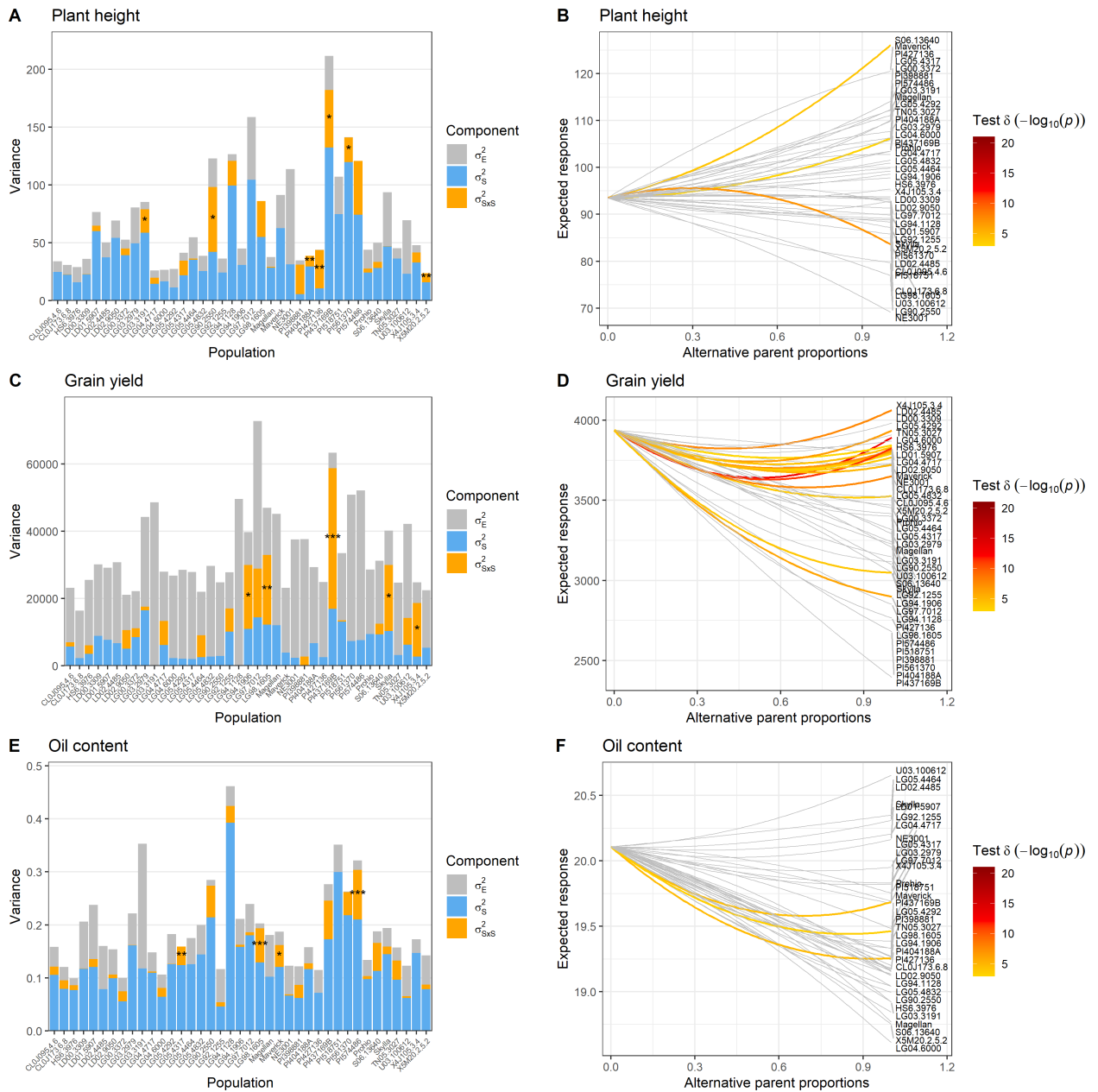


Figure 2: Variance components barplots and directional epistasis plots of the soybean NAM for three traits: plant height (A and B), grain yield (C and D) and oil content (E and F). For the variance components barplots, the range of p-value obtained for the likelihood ratio test of $\sigma^2_{S \times S}$ are indicated as: * for p-values inferior to 0.05, ** for p-values inferior to 0.01, and *** for p-values inferior to 0.001. For the directional epistasis plots, tests with a p-value higher than the Bonferroni threshold at a nominal level of 5% over the number of populations are indicated in grey.

and ear height, respectively. In contrast to the global and directional epistasis tests, only a few epistatic variance tests were significant. Most traits showed a limited to moderate contribution of the epistatic variance component ($\sigma^2_{S \times S}$) to the genetic variance.

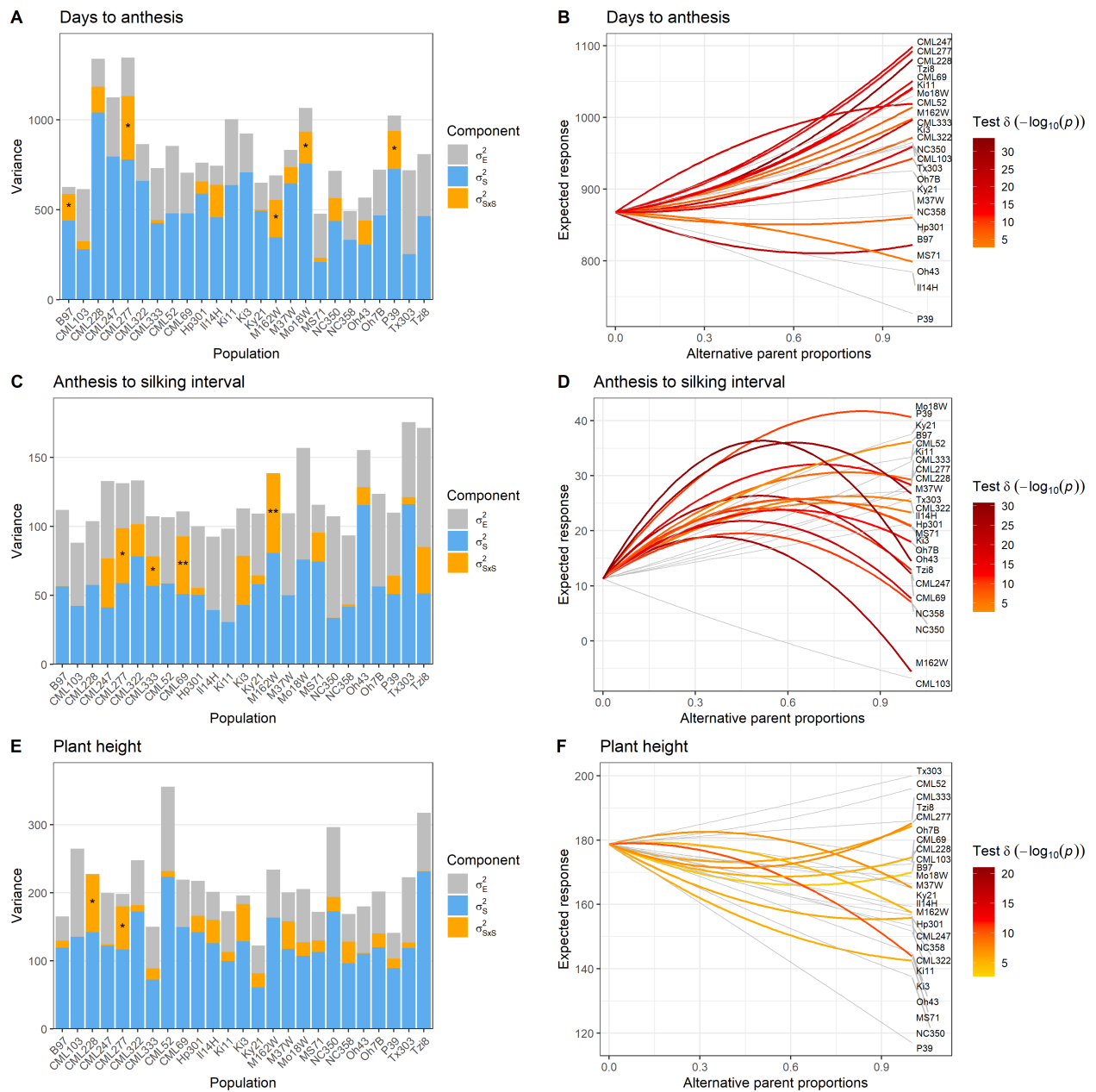


Figure 3: Variance components barplots and directional epistasis plots of the American maize NAM for three traits: days to anthesis (**A** and **B**), anthesis to silking interval (**C** and **D**) and plant height (**E** and **F**). For the variance components barplots, the range of p-value obtained for the likelihood ratio test of $\sigma_{S \times S}^2$ are indicated as: * for p-values inferior to 0.05, ** for p-values inferior to 0.01, and *** for p-values inferior to 0.001. For the directional epistasis plots, tests with a p-value higher than the Bonferroni threshold at a nominal level of 5% over the number of populations are indicated in grey

Discussion

Epistasis tests for inbred bi-parental populations

We developed a procedure to test for the presence of epistasis in inbred bi-parental progeny when genomic information is available. The procedure is based on a Gaussian mixed model where epistasis

is accounted for through a mean component (δ) and a variance component ($\sigma_{S \times S}^2$).

A significant δ for the trait indicates non-cancelling epistatic effects and can be defined as directional epistasis (Rouzic, 2014). In genomic prediction models, marker effects are usually assumed to be centered around zero. While this assumption holds true for additive effects, one expects to observe mostly positive dominance effects especially for traits showing heterosis and thus a directional dominance. Xiang et al. (2016) proposed the use of an inbreeding covariate to account for directional dominance in genomic prediction models. The concept of directional epistasis presented in our study is similar to the latter in that it involves non-cancelling interaction effects over loci pairs. Such directional epistasis is expected not only in bi-parental populations, but also in standard diploid populations. We plan to investigate how to extend our model to account for the directionality of epistasis in genomic prediction models in the near future. The two parameters δ and $\sigma_{S \times S}^2$ only involve epistatic effects and consequently will have value 0 whenever these effects are null. However, δ can be zero in the presence of epistatic interactions, provided that the effects cancel each other out over loci pairs. Also note that the structure of the covariance matrix involves a covariance Δ_{ij} between parental allele ancestries that has already been presented in the context of admixed populations (Aase et al., 2022; Rio et al., 2020b) for the segregation of group-specific allele ancestries. The use of a squared genomic relationship matrix (i.e., Δ_{ij}^2 in our study) to estimate the additive \times additive epistatic variance component has also been suggested in the context of standard diploid populations (Vitezica et al., 2017).

In absence of epistasis, the expected genetic value of an inbred offspring with perfectly balanced parental ancestry proportions does not deviate from the average value of its inbred parents. The genetic variance is then only driven by the effect of segregating QTL parental alleles. When epistasis is present, the expected genetic value can deviate from the average value of its inbred parents, and is also accompanied with an additional source of variability generated by the simultaneous segregation of QTL parental alleles at loci pairs.

The proposed test procedure requires the genomic information to be available in order to calculate the covariance matrix between allele ancestries associated with the variances σ_S^2 and $\sigma_{S \times S}^2$. In absence of genomic information, a test could theoretically be implemented using a linear model that considers expected parental proportions of 0.5 for the whole bi-parental population, provided that both parents are also evaluated. Note that it is not possible to distinguish between the two genetic variances and the model is no longer identifiable whenever the genetic information is not available. As a consequence, the epistatic variance $\sigma_{S \times S}^2$ cannot be tested using either the variance or the global test. Based on simulation, the tests involving the δ parameter (i.e., directional epistasis and global tests) prove to be more powerful than the test of the epistatic variance (S1 File). This observation is consistent with the larger number of significant directional and global tests compared to the test of epistatic variance for real traits in both NAM datasets. However, some traits such as days to maturity in soybean have a higher number of significant tests for $\sigma_{S \times S}^2$ than for δ , which demonstrates the benefit of testing both parameters to detect epistasis. Note that, thanks to genomic information, the test procedure can be applied only based on the progeny data, i.e. in absence of any parental phenotypic information. While this may be useful in some specific cases, we recommend experimental designs that include a thorough evaluation of both parents since the parental information is expected to have a strong (positive) impact on the estimation of model parameters. In practice, parents are often evaluated along with their progeny as replicated checks in bi- and multi-parental experimental designs - as in the maize and soybean NAM datasets - which provides favorable conditions for the application of the proposed procedure.

In this study, only epistatic interactions between pairs of loci were considered. We believe that this constraint is reasonable as the contribution of epistatic interaction to the fitness landscape is expected to decrease with increasing order of interaction (Weinreich et al., 2018). However, this approach could theoretically be generalized to higher order interactions by increasing the power to which the parental proportions in the fixed part of the model and the covariance matrix between allele ancestries are raised. One should keep in mind that this would be done at the cost of increased model complexity.

The strength of the test procedure presented in this study lies in its ability to target epistasis

both through the expectation and the variance of genetic values without requiring a complex design, as is usually the case (Mather and Jinks, 1982). Only simple inbred progenies (e.g., DH or RIL progeny) need to be genotyped and evaluated for a trait of interest. Such datasets are already available in large numbers and can simply be recycled to investigate the presence of epistasis. Our procedure also opens the way to the systematic investigation of epistasis in plant breeding programs and quantitative genetics studies. We highly recommend applying this test procedure prior to any QTL analysis as it helps determine the need to test for pairwise epistatic QTL effects.

Epistasis in breeding

Epistasis has long been considered a nuisance parameter that can be ignored in breeding (Crow, 2010). However, its presence can influence both the average performance of a cross, through the phenomenon of directional epistasis (Rouzig, 2014), and the genetic variance generated by that cross, through an additional epistatic variance component. The average performance and the genetic variance of a cross are the two parameters involved in the usefulness criterion proposed by Schnell and Utz (1975). A good understanding of the genetic determinism of traits subject to selection should allow a better reasoning of breeding mating designs, i.e. choosing which parents to cross along with the minimum progeny size that needs to be generated to maximize the chances of obtaining a superior individual for the trait.

In maize, several experiments have been conducted in which directional epistasis was detected for grain yield, forage yield, plant height, ear height, kernel row number, maturity or flowering traits (Bauman, 1959; Gamble, 1962; Hallauer and Russell, 1962; Melchinger et al., 1986, 1988; Wolf and Hallauer, 1997), or through epistatic variance for grain yield, forage yield and grain dry matter content (Melchinger et al., 1988). Similarly, significant epistasis has also been detected in soybean through variance for oil content (Hanson and Weber, 1961) or through directional epistasis for grain yield (Barona et al., 2012; Uzokwe et al., 2017). The results of these experiments are consistent with the results of our study demonstrating the prevalence of epistasis affecting the mean and variance of breeding traits in maize and soybean.

For some traits like anthesis to silking interval in maize or grain yield in soybean, the significant directional epistasis detected is always accompanied with a deterioration of the genetic values compared to that expected based on parental values under a purely additive model. Regarding maize, anthesis to silking interval increases in most progenies. This desynchronization of male and female flowering is a well-known indicator of stress (Edmeades et al., 2015) and is not a desirable trait in maize breeding. In the case of soybean, a decrease in grain yield is non-desirable as it directly conditions farmers' income. In general, it is reasonable to assume that a good elite inbred line results from both a good additive genetic value and a combination of alleles with favorable epistatic effects. When crossing two unrelated elite inbred lines, favorable allele combinations are likely to be broken, thus leading to deterioration of the progeny mean. A similar perspective is that of less-than-additive epistasis (Eshed and Zamir, 1996), where the effect of QTL becomes low in a genetic background with many favorable alleles. These hypotheses are supported by the tendency for directional epistasis to occur between parents with high grain yield in soybean.

A major challenge remains to understand what other factors determine the apparition of directional epistasis for a given cross. This question is particularly relevant for traits like plant height in maize displaying directional epistasis with opposite sign depending on the bi-parental population. One may hypothesize to observe more directional epistasis when crossing distant lines but our results do not completely support it. For instance, tropical maize inbred lines (e.g., lines with names starting with CML) do not systematically generate the largest directional epistasis when crossed to the distant stiff-stalk parent B73. A better solution to predict the apparition of directional epistasis would probably be a fine characterization of the genetic architecture of traits, including the number and position of QTLs along with their marginal and combined effects.

Material and Methods

Infinitesimal model

To develop a procedure to test for the presence of epistasis, our approach consists in (i) introducing an infinitesimal model of genetic values of inbred bi-parental progeny with digenic epistatic interactions, (ii) deriving the expression of the expected genetic value and the covariance between genetic values, (iii) modeling phenotypic values through a Gaussian mixed model that inherits its design and covariance matrices from the expressions obtained in (ii), and (iv) proposing tests of model parameters involving only epistasis.

Let us consider a population of homozygous inbred individuals derived from a cross between two homozygous inbred parents A and B in absence of selection. Only polymorphic biallelic QTLs are considered with two genotypic states indicating both the genotype and the ancestry of alleles. Let G_i be the genetic value of individual i . One has:

$$G_i = \mu + \sum_{m=1}^M A_{im}\beta_m + \sum_{m=1}^{M-1} \sum_{m'>m}^M A_{im}A_{im'}\delta_{mm'} \quad (1)$$

where μ is the genetic value of parent B, M is the number of loci, A_{im} is the allele ancestry indicator taking value 1 if individual i is homozygous for the parent A allele at locus m and 0 otherwise, β_m is the effect of substituting parent B allele by parent A allele at locus m , and $\delta_{mm'}$ is the interaction effect generated by substituting parent B alleles by parent A alleles at loci m and m' . Note that this genetic model corresponds to the multi-linear model proposed by Hansen and Wagner (2001) to build the genotype-to-phenotype map for the population of interest, considering only digenic interactions between loci.

Let us assume a Bernoulli distribution for allele ancestries: $A_{im} \sim \mathcal{B}(\pi_i)$, where π_i is the proportion of alleles originating from parent A for individual i , and $\text{cov}(A_{im}, A_{jm}) = \Delta_{ij}$ is the covariance between allele ancestries of individuals i and j at locus m . Note that allele ancestries are assumed to be independent between loci, which ignores the physical/genetic linkage between loci located on the same chromosome and constraints on allele ancestries due to a finite number of loci (see S2 File for details). Both π_i and Δ_{ij} are key parameters to describe the genetic content of an individual (or a pair of individuals) in terms of (shared) parent ancestry proportions. These quantities can be used to compute expected genetic values and covariance between genetic values, which will later be used in a statistical model to test for epistasis, as described hereafter.

Expectation and variance of genetic values

Let $E(G_i|\pi_i)$ be the expected genetic value conditional on the proportion of alleles originating from each parent. One shows that:

$$E(G_i|\pi_i) = \mu + \pi_i\beta + \pi_i^2\delta, \quad (2)$$

where $\beta = \sum_{m=1}^M \beta_m$ is the linear "regression" coefficient of the genetic value on the parent A genome proportion and $\delta = \sum_{m=1}^{M-1} \sum_{m'>m}^M \delta_{mm'}$ is the quadratic "regression" coefficient of the genetic value on the parent A genome proportion, which drives directional epistasis.

Similarly, let $\text{Cov}(G_i, G_j|\Delta_{ij})$ be the covariance between genetic values conditional on shared proportion of alleles originating from each parent. One shows that:

$$\text{Cov}(G_i, G_j|\Delta_{ij}) = \Delta_{ij}\sigma_S^2 + \Delta_{ij}^2\sigma_{S \times S}^2 \quad (3)$$

where σ_S^2 is the segregation variance generated by substituting parent B alleles by parent A alleles at loci over the whole genome (see expression in S2 File) and $\sigma_{S \times S}^2 = \sum_{m=1}^{M-1} \sum_{m'>m}^M (\delta_{mm'})^2$ is the segregation \times segregation ($S \times S$) interaction variance generated by the simultaneous substitution of parent B alleles by parent A alleles at loci pairs. For the proofs of Eq. (2) and Eq. (3), see S2 File.

Gaussian mixed model

The following model is assumed for phenotypic values:

$$\mathbf{y} = \mathbf{g} + \mathbf{e} \quad (4)$$

where \mathbf{y} is the vector of reference phenotypes (e.g., best linear unbiased prediction or least-square mean calculated over the whole experimental design), \mathbf{g} is the vector of genetic values defined below, \mathbf{e} is the vector of errors with $\mathbf{e} \sim \mathcal{N}(0, \mathbf{D}\sigma_E^2)$, \mathbf{D} is a diagonal matrix whose elements depend on the number of observations for each individual, σ_E^2 is the error variance, \mathbf{g} and \mathbf{e} are independent.

From the infinitesimal model in Eq. (1), one can derive an approximate Gaussian variance component model that inherits its mean and variance components from Eq. (2) and Eq. (3). The genetic values are then modeled as the sum of a fixed intercept and two random components independent from each other:

$$\mathbf{g} = \mathbf{X}\boldsymbol{\theta} + \mathbf{g}_S + \mathbf{g}_{S \times S} \quad (5)$$

where $\mathbf{X} = (\mathbf{1}, \boldsymbol{\pi}, \boldsymbol{\pi}^2)$ is the design matrix for fixed effects with $\boldsymbol{\pi}$ being the vector of parent A proportions, $\boldsymbol{\theta} = (\mu, \beta, \delta)^T$ being the vector of fixed effects, \mathbf{g}_S is the vector of the segregation component of the genetic value, and $\mathbf{g}_{S \times S}$ is the vector of the $S \times S$ interaction component of the genetic value. Each random genetic effect is modeled as being drawn from a normal distribution with a covariance structure inherited from the infinitesimal model: $\mathbf{g}_S \sim \mathcal{N}(0, \boldsymbol{\Delta}\sigma_S^2)$ and $\mathbf{g}_{S \times S} \sim \mathcal{N}(0, \boldsymbol{\Delta}^2\sigma_{S \times S}^2)$, where $\boldsymbol{\Delta}$ is the matrix of coefficients Δ_{ij} .

Inference

In practice, $\boldsymbol{\pi}$ can be estimated using:

$$(\boldsymbol{\pi})_i = \frac{1}{M} \sum_{m=1}^M A_{im}$$

and $\boldsymbol{\Delta}$ using:

$$(\boldsymbol{\Delta})_{ij} = \frac{1}{M} \sum_{m=1}^M (A_{im}A_{jm} - \pi_i\pi_j)$$

To ease the comparison of variances, all covariance matrices ($\boldsymbol{\Delta}$, $\boldsymbol{\Delta}^2$ and \mathbf{D}) are standardized as in Vitezica et al. (2017) so that the mean of diagonal elements equals 1. The inference of fixed and variance parameters is done using restricted maximum likelihood (ReML). The nullity of the δ parameter (i.e., directional epistasis test) can be tested using a Wald test for which the statistics has a chi-squared distribution. The other tests involving the variance parameter $\sigma_{S \times S}$ are based on a likelihood ratio test for which the distribution of the statistics is specific mixture of chi-square distributions (Self and Liang, 1987; Stram and Lee, 1994). All tests were adjusted for multiple testing using a 5% Bonferroni threshold over the number of populations in each dataset.

Datasets

Two NAM populations were considered in this study, each including multiple bi-parental populations evaluated for several traits of interest.

The soybean NAM is described in Song et al. (2017). It consists of 39 families of 140 recombinant inbred lines generated by crossing 40 diverse inbred lines to the central inbred line "IA3023". All individuals were genotyped using the "SoyNAM6K BeadChip" for which the reference allele (coded 0) corresponds to the homozygous genotype for "IA3023" alleles and the alternative allele (coded 1) corresponds to the homozygous alleles for the other parent. On average, 3,245 SNPs are polymorphic for a given NAM population. As described in Diers et al. (2018) and Xavier et al. (2018), the NAM populations were evaluated along with parental lines in 19 trials for grain yield in kg/ha, days to

maturity, plant height in centimeter, lodging score from 1 to 5, grain moisture in percentage of humidity, protein/oil/fiber in percentage of the grain content.

The American maize NAM is described in Yu et al. (2008) and McMullen et al. (2009). It consists of 25 families of 200 recombinant inbred lines (RILs) generated by crossing 25 temperate and tropical inbred lines to the central inbred line "B73". All individuals were genotyped for 1,106 polymorphic SNPs (Buckler et al., 2009) for which the reference allele (coded 0) corresponds to the homozygous genotype for "B73" alleles and the alternative allele (coded 1) corresponds to the homozygous alleles shared by all parental lines but B73 (i.e., only SNPs involving alleles specific to B73 were considered). The NAM populations were evaluated along with parental lines in ten trials (combinations of five locations and four years) for days to anthesis/silking and anthesis to silking interval in growing degree days, plant height and ear height in centimeter (Peiffer et al., 2014).

The model in Eq. (5) was applied separately to all maize and soybean NAM populations, where parent A refers to the central parent and parent B refers to the alternative parent. For all datasets, a single reference phenotype was considered for each individual after adjusting for experimental design effects. In both the maize and the soybean experiments, parental lines were used as checks and thus had a much larger number of observations, leading to a better precision associated with the resulting reference phenotype. The diagonal elements of \mathbf{D} from Eq. (5) were calculated accordingly and are summarized in Table 2.

	Maize NAM		Soybean NAM	
	N_{obs}	$(\mathbf{D})_{i,i}$	N_{obs}	$(\mathbf{D})_{i,i}$
Central parent	274	$\frac{1}{274}$	200	$\frac{1}{200}$
Alternative parents	10	$\frac{1}{10}$	4	$\frac{1}{4}$
RILs	1	1	1	1

Table 2: Number of observations for parental lines used as checks relative to that of RILs (N_{obs}) and corresponding diagonal elements of the error covariance matrix $(\mathbf{D})_{i,i}$

Data availability statement

The test procedure is implemented in a new R-package "EpiTest", which is based on the mixed model inference R-package "MM4LMM" (Laporte et al., 2022), both available from the CRAN. The American maize NAM dataset is available at: <https://www.panzea.org/>. The soybean NAM dataset is available at <https://www.soybase.org/SoyNAM/>.

References

- Aase, K., Jensen, H., and Muff, S. (2022). Genomic estimation of quantitative genetic parameters in wild admixed populations. *Methods in Ecology and Evolution*, 13(5):1014–1026.
- Álvarez-Castro, J. M. and Carlborg, O. (2007). A unified model for functional and statistical epistasis and its application in quantitative trait loci analysis. *Genetics*, 176(2):1151–1167.
- Barona, M. A. A., Filho, J. M. C., da Silva Santos, V., and Geraldi, I. O. (2012). Epistatic effects on grain yield of soybean [glycine max (l.) merrill]. *Crop Breeding and Applied Biotechnology*, 12(4):231–236.
- Bateson, W. (1909). *Mendel's Principles of Heredity*. Cambridge University Press.
- Bauman, L. F. (1959). Evidence of non-allelic gene interaction in determining yield, ear height, and kernel row number in corn. *Agronomy Journal*, 51(9):531–534.
- Blanc, G., Charcosset, A., Mangin, B., Gallais, A., and Moreau, L. (2006). Connected populations for detecting quantitative trait loci and testing for epistasis: an application in maize. *Theoretical and Applied Genetics*, 113(2):206–224.
- Buckler, E. S., Holland, J. B., Bradbury, P. J., Acharya, C. B., Brown, P. J., Browne, C., Ersoz, E., Flint-Garcia, S., Garcia, A., Glaubitz, J. C., Goodman, M. M., Harjes, C., Guill, K., Kroon, D. E., Larsson, S., Lepak, N. K., Li, H., Mitchell, S. E., Pressoir, G., Peiffer, J. A., Rosas, M. O., Rocheford, T. R., Romay, M. C., Romero, S., Salvo, S., Villeda, H. S., da Silva, H. S., Sun, Q., Tian, F., Upadyayula, N., Ware, D., Yates, H., Yu, J., Zhang, Z., Kresovich, S., and McMullen, M. D. (2009). The genetic architecture of maize flowering time. *Science*, 325(5941):714–718.
- Chahal, G. S. and Jinks, J. L. (1978). A general method of detecting the additive, dominance and epistatic variation that inbred lines can generate using a single tester. *Heredity*, 40(1):117–125.
- Choo, T. M. (1980). Doubled haploids for estimating additive epistatic variances in self-pollinating crops. *Canadian Journal of Genetics and Cytology*, 22(1):125–127.
- Choo, T. M. (1981). Doubled haploids for studying the inheritance of quantitative characters. *Genetics*, 99(3-4):525–540.
- Cockerham, C. C. (1954). an extension of the concept of partitioning hereditary variance of covariances among relatives when epistasis is present. *Genetics*, 39(6):859–882.
- Crawford, L., Zeng, P., Mukherjee, S., and Zhou, X. (2017). Detecting epistasis with the marginal epistasis test in genetic mapping studies of quantitative traits. *PLOS Genetics*, 13(7):e1006869.
- Crow, J. F. (2010). On epistasis: why it is unimportant in polygenic directional selection. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365(1544):1241–1244.
- Diers, B. W., Specht, J., Rainey, K. M., Cregan, P., Song, Q., Ramasubramanian, V., Graef, G., Nelson, R., Schapaugh, W., Wang, D., Shannon, G., McHale, L., Kantartzi, S. K., Xavier, A., Mian, R., Stupar, R. M., Michno, J.-M., An, Y.-Q. C., Goettel, W., Ward, R., Fox, C., Lipka, A. E., Hyten, D., Cary, T., and Beavis, W. D. (2018). Genetic architecture of soybean yield and agronomic traits. *G3 Genes|Genomes|Genetics*, 8(10):3367–3375.

- Edmeades, G. O., Bolaños, J., Elings, A., Ribaut, J.-M., Bänziger, M., and Westgate, M. E. (2015). The role and regulation of the anthesis-silking interval in maize. In *Physiology and Modeling Kernel Set in Maize*, pages 43–73. Crop Science Society of America and American Society of Agronomy.
- Eshed, Y. and Zamir, D. (1996). Less-than-additive epistatic interactions of quantitative trait loci in tomato. *Genetics*, 143(4):1807–1817.
- Fisher, R. A. (1918). XV.—the correlation between relatives on the supposition of mendelian inheritance. *Transactions of the Royal Society of Edinburgh*, 52(2):399–433.
- Gage, J. L., Monier, B., Giri, A., and Buckler, E. S. (2020). Ten years of the maize nested association mapping population: Impact, limitations, and future directions. *The Plant Cell*, 32(7):2083–2093.
- Gallais, A. (1990). Quantitative genetics of doubled haploid populations and application to the theory of line development. *Genetics*, 124(1):199–206.
- Gamble, E. E. (1962). Gene effects in corn (*zea mays* l.): I separation and relative importance of gene effects for yield. *Canadian Journal of Plant Science*, 42(2):339–348.
- Hallauer, A. R. and Russell, W. A. (1962). Estimates of maturity and its inheritance in maize. *Crop Science*, 2(4):289–294.
- Hansen, T. F. and Wagner, G. P. (2001). Modeling genetic architecture: A multilinear theory of gene interaction. *Theoretical Population Biology*, 59(1):61–86.
- Hanson, W. D. and Weber, C. R. (1961). Resolution of genetic variability in self-pollinated species with an application to the soybean. *Genetics*, 46(11):1425–1434.
- Hayman, B. I. (1958). The separation of epistatic from additive and dominance variation in generation means. *Heredity*, 12(3):371–390.
- Hemani, G., Theocharidis, A., Wei, W., and Haley, C. (2011). EpiGPU: exhaustive pairwise epistasis scans parallelized on consumer level graphics cards. *Bioinformatics*, 27(11):1462–1465.
- Jannink, J.-L. (2007). Identifying quantitative trait locus by genetic background interactions in association studies. *Genetics*, 176(1):553–561.
- Jannink, J.-L. and Jansen, R. (2001). Mapping epistatic quantitative trait loci with one-dimensional genome searches. *Genetics*, 157(1):445–454.
- Jannink, J.-L., Moreau, L., Charmet, G., and Charcosset, A. (2009). Overview of QTL detection in plants and tests for synergistic epistatic interactions. *Genetica*, 136(2):225–236.
- Jinks, J. L., Perkins, J. M., and Breese, E. L. (1969). A general method of detecting additive, dominance and epistatic variation for metrical traits II. application to inbred lines. *Heredity*, 24(1):45–57.
- Kao, C.-H., Zeng, Z.-B., and Teasdale, R. D. (1999). Multiple interval mapping for quantitative trait loci. *Genetics*, 152(3):1203–1216.
- Kearsey, M. J. and Jinks, J. L. (1968). A general method of detecting additive, dominance and epistatic variation for metrical traits i. theory. *Heredity*, 23(3):403–409.
- Kempthorne, O. (1954). The correlation between relatives in a random mating population. *Proceedings of the Royal Society of London. Series B - Biological Sciences*, 143(910):103–113.
- Kryazhimskiy, S. (2021). Emergence and propagation of epistasis in metabolic networks. *eLife*, 10.

- Laporte, F., Charcosset, A., and Mary-Huard, T. (2022). Efficient ReML inference in variance component mixed models using a min-max algorithm. *PLOS Computational Biology*, 18(1):e1009659.
- Mather, K. and Jinks, J. L. (1982). *Biometrical genetics*. Springer, New York, NY, 3 edition.
- McMullen, M. D., Kresovich, S., Villeda, H. S., Bradbury, P., Li, H., Sun, Q., Flint-Garcia, S., Thornsberry, J., Acharya, C., Bottoms, C., Brown, P., Browne, C., Eller, M., Guill, K., Harjes, C., Kroon, D., Lepak, N., Mitchell, S. E., Peterson, B., Pressoir, G., Romero, S., Rosas, M. O., Salvo, S., Yates, H., Hanson, M., Jones, E., Smith, S., Glaubitz, J. C., Goodman, M., Ware, D., Holland, J. B., and Buckler, E. S. (2009). Genetic properties of the maize nested association mapping population. *Science*, 325(5941):737–740.
- Melchinger, A. E. (1987). Expectation of means and variances of testcrosses produced from from f2 and backcross individuals and their selfed progenies. *Heredity*, 59(1):105–115.
- Melchinger, A. E., Geiger, H. H., and Schnell, F. W. (1986). Epistasis in maize (*zea mays* l.). *Theoretical and Applied Genetics*, 72(2):231–239.
- Melchinger, A. E., Schmidt, W., and Geiger, H. H. (1988). Comparison of testcrosses produced from f2 and first backcross populations in maize. *Crop Science*, 28(5):743–749.
- Peiffer, J. A., Romay, M. C., Gore, M. A., Flint-Garcia, S. A., Zhang, Z., Millard, M. J., Gardner, C. A. C., McMullen, M. D., Holland, J. B., Bradbury, P. J., and Buckler, E. S. (2014). The genetic architecture of maize height. *Genetics*, 196(4):1337–1356.
- Phillips, P. C. (2008). Epistasis — the essential role of gene interactions in the structure and evolution of genetic systems. *Nature Reviews Genetics*, 9(11):855–867.
- Prabhu, S. and Pe'er, I. (2012). Ultrafast genome-wide scan for SNP–SNP interactions in common complex disease. *Genome Research*, 22(11):2230–2240.
- Raffo, M. A., Sarup, P., Guo, X., Liu, H., Andersen, J. R., Orabi, J., Jahoor, A., and Jensen, J. (2022). Improvement of genomic prediction in advanced wheat breeding lines by including additive-by-additive epistasis. *Theoretical and Applied Genetics*, 135(3):965–978.
- Rawlings, J. O. and Cockerham, C. C. (1962a). Analysis of double cross hybrid populations. *Biometrics*, 18(2):229.
- Rawlings, J. O. and Cockerham, C. C. (1962b). Triallel analysis. *Crop Science*, 2(3):228–231.
- Rio, S., Mary-Huard, T., Moreau, L., Bauland, C., Palaffre, C., Madur, D., Combes, V., and Charcosset, A. (2020a). Disentangling group specific QTL allele effects from genetic background epistasis using admixed individuals in GWAS: An application to maize flowering. *PLOS Genetics*, 16(3):e1008241.
- Rio, S., Moreau, L., Charcosset, A., and Mary-Huard, T. (2020b). Accounting for group-specific allele effects and admixture in genomic predictions: Theory and experimental evaluation in maize. *Genetics*, 216(1):27–41.
- Rouzic, A. L. (2014). Estimating directional epistasis. *Frontiers in Genetics*, 5.
- Schnell, F. and Utz, H. (1975). F1 leistung und elternwahl in der zuchtung von selbstbefruchttern. In *Bericht über die Arbeitstagung der Vereinigung Österreichischer Pflanzenzüchter*, pages 243–248. BAL Gumpenstein, Gumpenstein, Austria.
- Self, S. G. and Liang, K.-Y. (1987). Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *J. Am. Stat. Assoc.*, 82(398):605.

- Song, Q., Yan, L., Quigley, C., Jordan, B. D., Fickus, E., Schroeder, S., Song, B.-H., An, Y.-Q. C., Hyten, D., Nelson, R., Rainey, K., Beavis, W. D., Specht, J., Diers, B., and Cregan, P. (2017). Genetic characterization of the soybean nested association mapping population. *The Plant Genome*, 10(2).
- Stram, D. O. and Lee, J. W. (1994). Variance components testing in the longitudinal mixed effects model. *Biometrics*, 50(4):1171–1177.
- Uzokwe, V. N. E., Asafo-Adjei, B., Fawole, I., Abaidoo, R., Odeh, I. O. A., Ojo, D. K., Dashiell, K., and Sanginga, N. (2017). Generation mean analysis of phosphorus-use efficiency in freely nodulating soybean crosses grown in low-phosphorus soil. *Plant Breeding*, 136(2):139–146.
- Varona, L., Legarra, A., Toro, M. A., and Vitezica, Z. G. (2018). Non-additive effects in genomic selection. *Frontiers in Genetics*, 9.
- Vitezica, Z. G., Legarra, A., Toro, M. A., and Varona, L. (2017). Orthogonal estimates of variances for additive, dominance, and epistatic effects in populations. *Genetics*, 206(3):1297–1307.
- Weinreich, D. M., Lan, Y., Jaffe, J., and Heckendorn, R. B. (2018). The influence of higher-order epistasis on biological fitness landscape topography. *Journal of Statistical Physics*, 172(1):208–225.
- Wolf, D. P. and Hallauer, A. R. (1997). Triple testcross analysis to detect epistasis in maize. *Crop Science*, 37(3):763–770.
- Xavier, A., Jarquin, D., Howard, R., Ramasubramanian, V., Specht, J. E., Graef, G. L., Beavis, W. D., Diers, B. W., Song, Q., Cregan, P. B., Nelson, R., Mian, R., Shannon, J. G., McHale, L., Wang, D., Schapaugh, W., Lorenz, A. J., Xu, S., Muir, W. M., and Rainey, K. M. (2018). Genome-wide analysis of grain yield stability and environmental interactions in a multiparental soybean population. *G3 Genes|Genomes|Genetics*, 8(2):519–529.
- Xiang, T., Christensen, O. F., Vitezica, Z. G., and Legarra, A. (2016). Genomic evaluation by including dominance effects and inbreeding depression for purebred and crossbred performance with an application in pigs. *Genetics Selection Evolution*, 48(1).
- Yu, J., Holland, J. B., McMullen, M. D., and Buckler, E. S. (2008). Genetic design and statistical power of nested association mapping in maize. *Genetics*, 178(1):539–551.
- Zhang, X., Huang, S., Zou, F., and Wang, W. (2010). TEAM: efficient two-locus epistasis tests in human genome-wide association study. *Bioinformatics*, 26(12):i217–i227.