1    # Imperfect Linkage Disequilibrium Generates

2    # Phantom Epistasis (& Perils of Big Data)

3

4

5    by

6

7

8    G. de los Campos[1*], D. Sorensen[2] & M. A. Toro[3],

9

10

11

12

13

14    1: Epidemiology & Biostatistics, Statistics & Probability departments, IQ-Institute for Quantitative

15       Health Science and Engineering, Michigan State University;

16    2: Department of Molecular Biology and Genetics, Faculty of Science and Technology, Aarhus

17       University.

18    3: Producción Animal, Universidad Politécnica de Madrid.

19

20    *: Corresponding Author. 775 Woodlot Dr. (1311), East Lansing, MI, 48824
21    (gustavoc@msu.edu )

22    **ABSTRACT.** The genetic architecture of complex human traits and diseases is affected by large

23    number of possibly interacting genes, but detecting epistatic interactions can be challenging. In

24    the last decade, several studies have alluded to problems that linkage disequilibrium can create

25    when testing for epistatic interactions between DNA markers. However, these problems have

26    not been formalized nor have their consequences been quantified in a precise manner. Here we

27    use a conceptually simple three locus model involving a causal locus and two markers to show

28    that imperfect LD can generate the illusion of epistasis, even when the underlying genetic

29    architecture is purely additive. We describe necessary conditions for such "*phantom epistasis*" to

30    emerge and quantify its relevance using simulations. Our empirical results demonstrate that

31    phantom epistasis can be a very serious problem in GWAS studies (with rejection rates against

32    the additive model greater than 0.2 for nominal p-values of 0.05, even when the model is purely

33    additive). Some studies have sought to avoid this problem by only testing interactions between

34    SNPs with R-sq. <0.1. We show that this threshold is not appropriate and demonstrate that the

35    magnitude of the problem is even greater with large sample size. We conclude that caution must

36    be exercised when interpreting GWAS results derived from very large data sets showing strong

37    evidence in support of epistatic interactions between markers.

38

39    **Keywords**: epistasis, apparent epistasis, phantom epistasis, GWAS, linkage disequilibrium,

40    imperfect LD, missing heritability, Big Data.

## Introduction

42    A big challenge in genetics is to understand how variation at the DNA sequences translates into

43    phenotypic variation. Genome-wide-association (GWA) studies address part of this challenge by

44    testing for the association between phenotype (or a disease indicator) with genotype, one locus

45    at a time. In the last decade, many GWA studies were conducted; these studies have reported

46    thousands of SNP's (single nucleotide polymorphism) associated to complex traits and diseases

47    (http://www.ebi.ac.uk/gwas).

48        Recently, several studies in model organisms (e.g., Mackay 2014), humans (Strange, Ask, and

49    Nielsen 2013) and agricultural species (e.g., Huang, Xu, and Cai 2014), have used genotype data

50   linked to phenotypes to investigate the presence of epistatic interactions between loci. Cordell

51   (2002, 2009) and Wei, Hermani, and Haley (2014) provide comprehensive reviews of the methods

52   commonly used to detect epistatic interactions.

53   There are several issues associated with studies aimed at detecting interactions, including

54   matters of scale, the importance of the contribution of epistasis at the level of the genotype

55   effects or at the level of the genotypic variance (e.g., Hill, Goddard, and Visscher 2008) and how

56   an interaction detected in a linear statistical model may be associated to biological pathways that

57   underlie a complex trait (e.g., Wang, Elston, and Zhu 2010; Aschard 2016). The latter becomes

58   particularly problematic when the markers used to assess associations between SNPs and

59   phenotypes (or a disease indicator) are in imperfect linkage disequilibrium (LD) with the alleles

60   at the causal loci (i.e., those responsible for inter-individual genetic differences in a trait or

61   disease phenotype). Under those conditions, evidence supporting the existence of a non-null

62   interaction between markers do not necessarily provide definite evidence of epistasis at causal

63   loci. Indeed, when the SNPs used in association analyses are in imperfect LD with the alleles at

64   causal loci, linear regression on SNPs may lead to unaccounted variance, or *missing heritability*

65   (e.g., Manolio et al. 2009; de los Campos, Sorensen, and Gianola 2015). Furthermore, the un-

66   accounted additive signal may be correlated with interaction contrasts, thus creating the

67   "illusion" of epistasis even for traits that are purely additive.

68   Several authors have expressed concerns about the role that LD can have on the detection of

69   epistasis (e.g., Wei, Hermani, and Haley 2014).  However, these problems have not been

70   quantified nor have they been given a precise mathematical treatment. In this study, we present

71   a simple three locus model involving a causal (unobserved) locus and two markers that makes

72   explicit how *phantom epistasis* may emerge even in systems that are strictly additive. We use

73   this model to derive a set of conditions that are necessary for the occurrence of phantom

74   epistasis, and quantify the magnitude of the problem using simulations based on real human

75   genotypes from the UK-Biobank. Our results suggest that imperfect LD can lead to seriously

76   inflated type-I error rates. We also show that the rate of detection of phantom epistatic

77   interactions increases with sample size; this should be considered when testing for epistatic

78   interactions using big data sets such as the ones that are becoming available.

3

## Materials and Methods

79

80  To study what factors may induce phantom epistasis we consider a simple model with three bi-

81  allelic loci. One of them, denoted as $z_i$, represents a causal locus (also referred as to the

82  'quantitative trait locus', QTL) and has a direct effect on the expression of a phenotype $y_i$. The

83  other two loci, denoted as $x_{1i}$ and $x_{2i}$, are markers that are possibly in LD with the QTL but have

84  no causal effect on $y_i$. For SNPs, a standard practice is to code genotypes $(z_i, x_{1i}, x_{2i})$ by counting

85  at each of the loci the number of copies of a reference allele carried by the $i^{th}$ individual. Here, to

86  facilitate the presentation we assume that genotypic codes and phenotypes are expressed as

87  deviations from their corresponding means; therefore $E(z_i) = E(x_{1i}) = E(x_{2i}) = E(y_i) = 0$. In

88  this setting, a single-locus strictly additive model takes the form

89  $$y_i = z_i b + \delta_i, \qquad [1]$$

90  where $b$ is the additive effect of an allele substitution at locus $z$, and $\delta_i$ is an error term. Evidently,

91  with only one causal locus there is no epistasis. We assume that [1] represents the causal model.

92  Next, suppose that an instrumental regression of the form

93  $$y_i = x_{1i}\beta_1 + x_{2i}\beta_2 + x_{1i}x_{2i}\beta_{12} + \varepsilon_i \qquad [2]$$

94  is used to investigate the presence of epistasis. Here, the $\beta's$ are regression coefficients that are

95  functions of the QTL effect ($b$) and of the (multi-locus) LD involving the two markers and the QTL

96  genotypes. In the population, given the centered genotype codes, the regression coefficients

97  entering in the right-hand-side of [2] are

98  $$\begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_{12} \end{bmatrix} = \begin{bmatrix} E(x_{1i}^2) & E(x_{1i}x_{2i}) & E(x_{1i}^2 x_{2i}) \\ E(x_{1i}x_{2i}) & E(x_{2i}^2) & E(x_{1i}x_{2i}^2) \\ E(x_{1i}^2 x_{2i}) & E(x_{1i}x_{2i}^2) & E(x_{1i}^2 x_{2i}^2) \end{bmatrix}^{-1} \begin{bmatrix} E(y_i x_{1i}) \\ E(y_i x_{2i}) \\ E(y_i x_{1i}x_{2i}) \end{bmatrix}.$$

99      If the random residual $\delta_i$ in expression [1] is orthogonal to the genotypes, then $E(y_i x_{1i}) =$

100  $E(z_i x_{1i})b$, $E(y_i x_{2i}) = E(z_i x_{2i})b$ and $E(y_i x_{1i}x_{2i}) = E(z_i x_{1i}x_{2i})b$. Thus, the population

101  regression coefficients are defined by

102  $$\begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_{12} \end{bmatrix} = \begin{bmatrix} E(x_{1i}^2) & E(x_{1i}x_{2i}) & E(x_{1i}^2 x_{2i}) \\ E(x_{1i}x_{2i}) & E(x_{2i}^2) & E(x_{1i}x_{2i}^2) \\ E(x_{1i}^2 x_{2i}) & E(x_{1i}x_{2i}^2) & E(x_{1i}^2 x_{2i}^2) \end{bmatrix}^{-1} \begin{bmatrix} E(z_i x_{1i}) \\ E(z_i x_{2i}) \\ E(z_i x_{1i}x_{2i}) \end{bmatrix} b. \qquad [3]$$

4

103    This indicates that the regression coefficients of the instrumental model [2] are not only

104    functions of the QTL effect ($b$) and of pair-wise (1$^{st}$ order) LD but also of higher order LD, e.g.,

105    joint disequilibrium at three loci, $E(z_i x_{1i} x_{2i})$. The moments involved in the right hand-side of [3]

106    are diploid genotypic measurements of LD. Under random mating these genotypic measures of

107    LD are equal to twice the standard haploid measures of LD (the $D$-coefficients for two and tree

108    loci linkage disequilibrium; see Section 1 of the Supplementary Methods for further details).

109    In the population, the interaction effect $\beta_{12}$ is given by a linear combination of two-loci LD

110    between each of the markers and the QTL and by three-loci LD involving the two markers and

111    the QTL:  $\beta_{12} = [t_{31}E(z_i x_{1i}) + t_{32}E(z_i x_{2i}) + t_{33}E(z_i x_{1i} x_{2i})]b$. Here, the $t's$ are the entries of

112    the third row of the inverse of the coefficient matrix

113
$$T^{-1} = \begin{bmatrix} E(x_{1i}^2) & E(x_{1i} x_{2i}) & E(x_{1i}^2 x_{2i}) \\ E(x_{1i} x_{2i}) & E(x_{2i}^2) & E(x_{1i} x_{2i}^2) \\ E(x_{1i}^2 x_{2i}) & E(x_{1i} x_{2i}^2) & E(x_{1i}^2 x_{2i}^2) \end{bmatrix}^{-1} .$$

114    expressions to study the conditions that lead to a null interaction between markers.

**Conditions that lead to phantom epistasis**

116    Next, we describe sufficient conditions for $\beta_{12} = 0$. These sufficient conditions also imply

117    necessary conditions for phantom epistasis, $\beta_{12} \neq 0$, to emerge.

118    *Complete Linkage Equilibrium.*  If the QTL is in LE with the two markers, then  $(z_i, x_{1i}, x_{2i}) =$

119    $p(z_i)p(x_{1i}, x_{2i})$. Consequently, $E(z_i x_{1i}) = E(x_{1i})E(z_i) = 0$, $E(z_i x_{2i}) = E(x_{2i})E(z_i) = 0$, and

120    $E(z_i x_{1i} x_{2i}) = E(x_{1i} x_{2i})E(z_i) = 0$. Therefore, all elements of the right-hand-side of [3] are

121    equal to zero and, thus $\beta_1 = \beta_2 = \beta_{12} = 0$. Therefore, *a first necessary condition for phantom*

122    *epistasis to emerge is that the QTL must be in LD with at least one of the SNPs.*

123    *Perfect Linkage Disequilibrium*. On the other extreme, if there is perfect LD between the QTL

124    and the marker pair $(x_{1i} x_{2i})$, then the QTL genotype can be expressed as a linear function of the

125    two marker genotypes $z_i = x_{1i}\beta_1 + x_{2i}\beta_2$. In this case, a linear regression on the two markers

126    captures fully the QTL variance and therefore the interaction term will be equal to zero. (A

127    derivation of this intuitive result is presented section 2 of the Supplementary Methods.)

128    Therefore, perfect LD is a sufficient condition for $\beta_{12} = 0$. Consequently, *a second necessary*

129    *condition for phantom epistasis to emerge is imperfect LD between the QTL and the marker pair*.

130     This guarantees that some fraction of the QTL variance is not captured by linear regression on

131     the two marker genotypes. Furthermore, if the left-out QTL signal is not orthogonal to the

132     interaction contrast $x_{1i}x_{2i}$, then $\beta_{12} \neq 0$.

133     ***Independence of one of the markers prevents phantom epistasis***. Consider now an

134     intermediate case where one of the markers (say $x_{2i}$) is independent of the pair formed by the

135     QTL and the other marker $(z_i, x_{1i})$. This implies that $p(z_i, x_{1i}, x_{2i}) = p(z_i, x_{1i})p(x_{2i})$. Under this

136     condition, because the two markers are in LE, the coefficient matrix and its inverse $(T^{-1})$ is

137     diagonal; therefore, $\beta_{12} = \frac{E(z_i x_{1i} x_{2i})}{E(x_{1i}^2 x_{2i}^2)} b$. Moreover, $E(z_i x_{1i} x_{2i}) = E(z_i x_{1i})E(x_{2i}) = 0$, implying

138     that $\beta_{12} = 0$. Therefore, a _third necessary condition_ for phantom epistasis to emerge is that the

139     _three loci must be jointly in LD_.

140     In summary, _phantom epistasis can emerge if the three loci are in mutual but imperfect LD_.

## Simulation

141 **Simulation**

142 The analytical results presented in the previous section indicates that multi-locus LD plays an

143     important role in determining whether phantom epistasis may emerge. To shed light on the

144     nature and the magnitude of the problem we conducted Monte Carlo simulations to assess how

145     LD among the three genotypes $(z_i, x_{1i}, x_{2i})$ affects the rates at which $H_0: \beta_{12} = 0$ is rejected.

146     Data were generated according to an additive model with a single causal locus that had strictly

147     additive gene action (as in expression [1]) and then analyzed using an instrumental model such

148     as the one in [2]. In this setting, rejection of $H_0: \beta_{12} = 0$ is indicative of phantom epistasis.

149     Simulations were based on real human genotypes of distantly related white Caucasian

150     individuals from the ***UK-Biobank***, a cohort study consisting of about half a million participants

151     aged between 40-69 years who were recruited in 2006-2010. The National Research Ethics

152     Committee approved the study and informed consent was obtained from all participants. Study

153     details are described elsewhere (Sudlow et al. 2015).

154     To avoid confounding due to population structure and long-range LD due to family

155     relationships we focused on distantly related white Caucasian individuals. Therefore, we only

156     considered subjects whose self-reported ethnicity was Caucasian and confirmed their genetic

157     race/ethnicity using SNP-derived principal components. From these individuals, we identified

6

158    ~270,000 subjects that have pairwise genomic relationships, $G_{ij} = p^{-1} \sum_{k=1}^{p} \frac{(x_{ik}-2\theta_k)(x_{jk}-2\theta_k)}{2\theta_k(1-\theta_k)}$,

159    smaller than 0.03. Here, $x_{ik}$ and $x_{jk}$ are genotypes (coded as 0, 1, 2) at the $k^{th}$ SNP of the $i^{th}$

160    and $j^{th}$ individual, respectively, and $\theta_k$ is the frequency of the allele counted at the $k^{th}$ loci.

161    Genomic relationships were computed using the *getG()* function of the BGData R-package

162    (Grueneberg and de los Campos 2017).

163    ***Genotypes*** where from the Affymetrix UK BiLEVE Axiom and Affymetrix UK Biobank

164    Axiom® arrays.  Only SNPs with minor-allele-frequency greater than 0.1% and those with

165    missing calling rate smaller than 3% were used for simulations. Furthermore, since we focused

166    on a single locus model, we used only SNPs mapped to chromosome 1. There were 66,331 SNPs

167    mapped to chromosome 1, of those, 45,866 passed our minor-allele frequency and calling rate

168    inclusion criteria.

169    ***Marker-QTL pairs***. The position of the QTL genotype $z_i$ was determined by randomly

170    choosing a marker position on chromosome 1. In a first simulation scenario, the two chosen

171    markers were those flanking the QTL (i.e., those immediately adjacent to it). In subsequent

172    simulation scenarios, the marker locus "to the right" ($x_{2i}$) of the QTL $z_i$ was placed at increasing

173    base-pair lags from the QTL, whereas the marker locus to the left ($x_{1i}$) of $z_i$ remained always the

174    most proximal marker "to the left of $z_i$". In this manner, the LD between one of the markers and

175    the QTL was approximately constant whereas the LD between the distal marker, $x_{2i}$, and the

176    marker-QTL pair ($x_{1i}$, $z_i$) decreased as base-pair distance between the two markers increased.

177    For each simulation scenario, we conducted 10,000 Monte Carlo replicates with random

178    assignment of the QTL position within chromosome 1.

179    ***Phenotypes*** were generated according to a single-locus additive model (expression [1]). with

180    the QTL explaining one-half-of-one percent (0.005) of the phenotypic variance.

181    ***Inferences*** were based on a linear model such as that of expression [2] extended with

182    inclusion of an intercept and the top 5 SNP-derived PCs, that is

183    $$y_i = \mu + \sum_{j=1}^{5} PC_{ji}\gamma_j + x_{1i}\beta_1 + x_{2i}\beta_2 + x_{1i}x_{2i}\beta_{12} + \varepsilon_i \qquad [2b]$$

184    Principal components were included to avoid any confounding that may emerge from any

185    remaining substructure that may have been present. The PCs used in [2b] were derived using 50K

7

186    SNPs evenly distributed in the entire genome. The model of expression [2b] was fitted via least

187    squares using the `ls.fit()` function of R (R Development Core Team 2012). Then for each

188    scenario and MC replicate we saved the p-value associated to the interaction term and counted

189    the proportion of times that $H0: \beta_{12} = 0$ was rejected when using a significance level of 0.05.

190    ***Genotypic measures of LD*** between pairs of loci, $R^2(x_1, x_2)$, $R^2(x_1, z)$ and $R^2(x_2, z)$, were

191    computed using the squared correlation between genotypes at the two loci. This information was

192    stored for each MC replicate of each simulated scenario. The proportion of variance of the QTL

193    genotype explained by linear regression on the two markers, $R^2(z \sim x_1 + x_2)$, was computed by

194    Analysis of Variance, of a linear model where the QTL genotype was regressed, via least squares,

195    on the two markers using a main effects additive model of the form: $z_i = \mu + x_{1i}a_1 + x_{2i}a_2 +$

196    $\varepsilon_i$. The R-squared from this model was also saved for each MC replicate of each scenario and

197    then used to analyze the relationship between this LD measure and the rate of rejection of

198    $H0: \beta_{12} = 0$.

199    ***Effect of sample size***. The power to detect a non-null interaction effect depends on two main

200    factors: the proportion of variance explained by that interaction and sample size. The first factor

201    is controlled in our simulation by controlling the distance between the QTL and the distal marker;

202    this affects LD among the three loci and thus the size of the marker-interaction (see Methods).

203    To assess the effect of sample size on inferences we carried out simulations using four different

204    sample sizes: n=10K (K=1,000, this is representative of the size of a standard GWAS cohort) and

205    n=50K, 100K and 250K (these sample sizes are more representative of modern large biomedical

206    data sets).

207

**Data availability**

209    The genotypes used in the simulation were from the UK Biobank. Data was acquired under

210    project identification number 15326. The data are available for all bonafide researchers and can

211    be obtained by applying at http://www.ukbiobank.ac.uk/register-apply/. The Institutional

212    Review Board (IRB) of Michigan State University has approved this research with the IRB number

213    15-745.

8

## Results

214

**Figure 1** shows measures of linkage disequilibrium between the three loci $(z_i, x_{1i}, x_{2i})$ involved

215

216    in the system. The average (across Monte Carlo replicates) proportion of variance of the QTL $(z_i)$

217    explained by the most adjacent marker $(x_{1i})$ averaged was about 0.085 (**Figure 1**); however, the

218    distribution of this statistic is highly skewed.  When $x_{1i}$ and $x_{2i}$ were the two flanking markers of

219    the QTL, on average they jointly explained on average 15% of the QTL variance. Therefore, on

220    average there was a sizable fraction of imperfect LD between the QTL and the markers. This leads

221    to a sizable rate of "missing" heritability. The LD between $x_{2i}$ and the pair $(x_{1i}, z_i)$ decreased as

222    the distance between $x_{2i}$ and the pair $(x_{1i}, z_i)$ increased. The R-sq. between $x_{2i}$ and either the

223    other marker or the QTL, falls very quickly for lags between 0-0.5Mb and reached near zero values
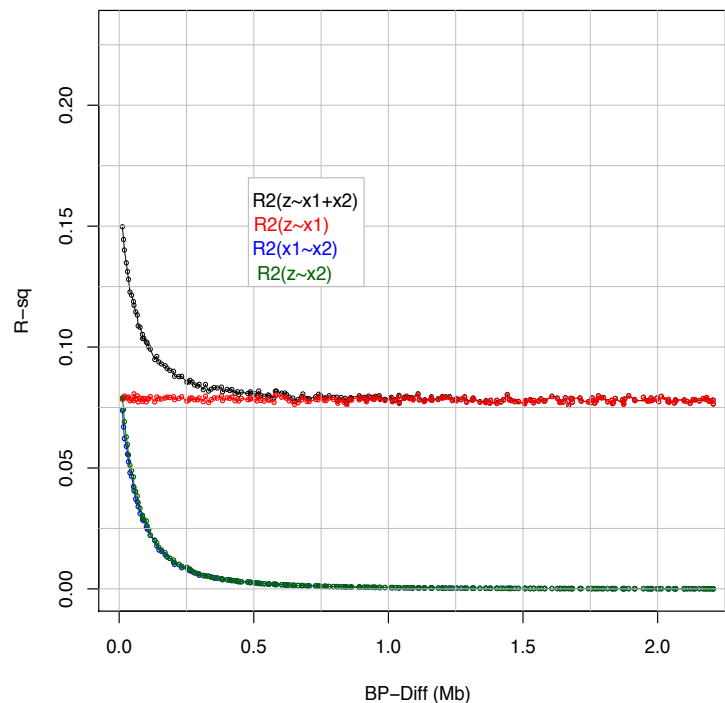
224    at approximately 1 Mb (Figure 1).

225



**Figure 1**. Average R-squared between pairs of loci and proportion of variance of the QTL genotype explained by the two markers, $R^2(z_i \sim x_{1i} + x_{2i})$, versus distance between the QTL $(z_i)$ and the distal marker $(x_{2i})$. Marker $x_{1i}$ was always adjacent to the QTL.

9

226

227    In our simulation rejection of $H_0: \beta_{12} = 0$ was performed using a significance level of 0.05. **Figure**

228    **2** displays the estimated rates of rejection by BP-distance between the QTL and the distal marker

229    $(x_{2i})$ and sample size. For the largest sample size, the curve relating empirical rejection rates with

230    BP distance was clearly above 0.05 for distances of up to 2MB. The highest rejection rate was

231    observed for $n$=250,000 when $x_{2i}$ and the QTL were at a distance of about 0.15 MB; here the

232    empirical rejection rate was ~0.13–this is more than twice the value expected under the absence

233    of phantom epistasis. The curves relating empirical rejection rates with physical distance reach

234    the nominal rejection rate of 0.05 at ~1Mb for n=10,000; however, for larger sample size the

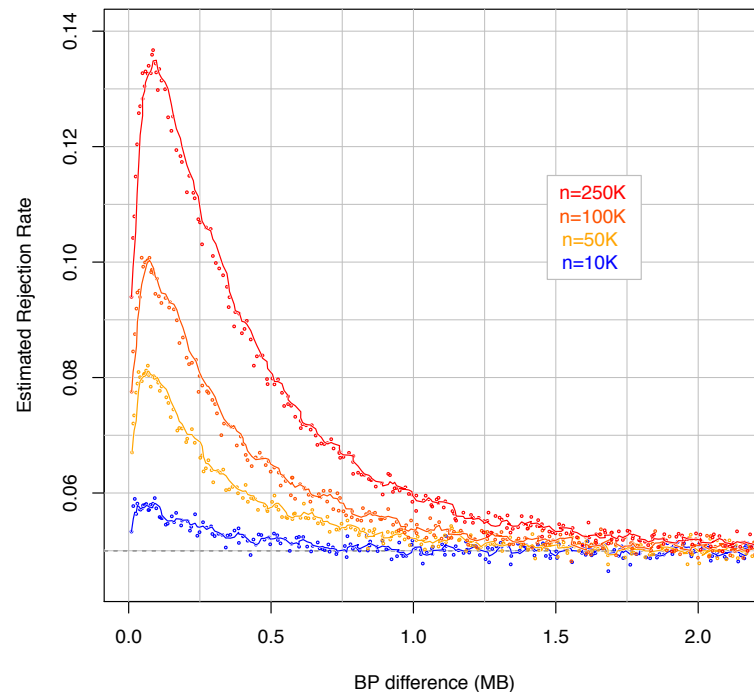235    curves stayed above 0.05 even for distances longer than 1Mb.



**Figure 2**. Estimated rejection rates versus distance between the QTL and the distal marker, by sample size. In the simulations one of the markers $(x_{1i})$ was adjacent to the QTL $(z_i)$ and the other marker $(x_{2i})$ was placed at increasing distance from the pair $(x_{1i}, z_i)$.

236

237

10

238    The extent of LD varies substantially along the genome; therefore, for a given BP distance some

239    regions may have very weak LD while others may have, at the same distance, SNPs in moderate

240    or high LD. **Figure 3** shows another way of viewing the simulation results displayed in **Figure 2**

241    where the average rejection rate is calculated within bins of R-sq. between the two markers.

242    When the two markers were un-correlated, rejection rates were very close to 0.05 indicating

243    absence of phantom epistasis. However very small LD between the two markers generates

244    considerably higher rejection rates: an $R^2(x_{1i}, x_{2i}) \sim 0.1$ leads to rejection rates as high as 0.21

245    with the largest sample size. The maximum rejection rates occur when the R-sq. between

246    markers is between 0.1 to 0.2. Beyond this value in the range (0.2-0.9) rejection rates shows a
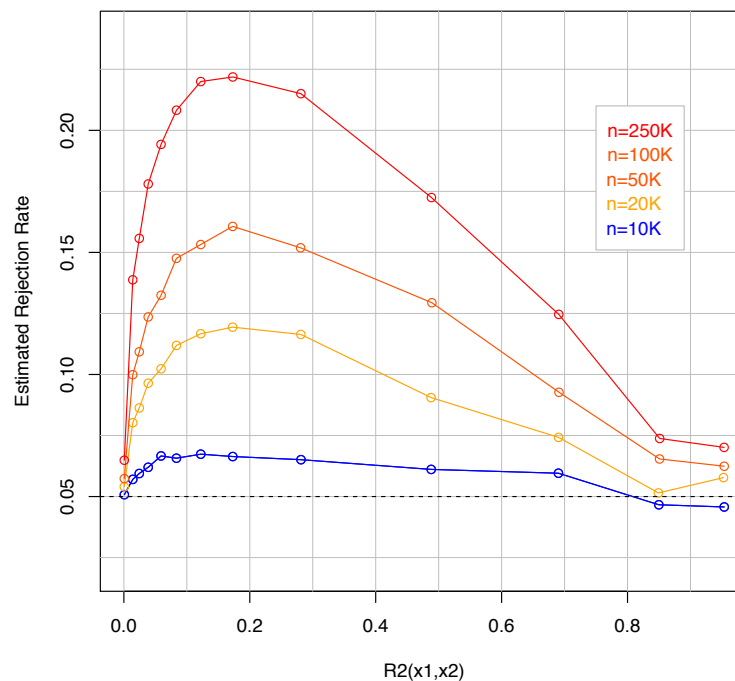
247    linear decline.



**Figure 3**. Estimated rejection rates by R-sq. between the two markers and sample size.

In the simulations one of the markers $(x_{1i})$ was adjacent to the QTL $(z_i)$) and the other

marker $(x_{2i})$ was placed at increasing distance from the pair $(z_i, x_{1i})$.

248

249

11

250   The results in Figures 2 and 3 are in line with the conceptual model described in the previous

251   section. Analytically, the conditions needed for phantom epistasis to emerge include

252   simultaneous but imperfect LD between the three loci. When the distal marker becomes

253   independent of the QTL-proximal-marker pair, there is no phantom epistasis. This happens at

254   about 2MB (Figure 2) and requires the R-sq. between the two markers to be very close to zero

255   (Figure 3).  When the LD between the QTL and the marker pair is very high but imperfect (e.g.

256   $0.9 < R^2(z \sim x_1 + x_2) < 1$), some phantom epistasis is generated. However, for those R-sq.

257   values the amount of signal that is not captured by the linear regression on the two markers and

258   that can be recaptured by an interaction term involving both is small.  Therefore, a very large

259   sample size is required to detect the phantom epistasis (compare the empirical rejection rates in

260   Figure 3 for R-sq. in the range 0.9-1).

## Discussion

262   There is a substantial amount of literature reporting the presence of epistasis affecting complex

263   traits but results, when scrutinized, have been controversial. Sometimes the controversy spawns

264   from the suspicion that epistatic interactions may be capturing additive signals that were missed

265   by the baseline additive model used to test interactions. For instance, Hermani et al. (2014)

266   identified 30 pairs of SNPs that interact influencing gene expression and that were replicated

267   across two independent studies. In a subsequent study (Wood et al. 2014) replicated many of the

268   interactions reported by Hermani et al.; however, in each case, using sequence data, a single

269   third variant could explain all the apparent epistasis. This happened even after removal of all

270   pairs of SNPs with $r^2 < 0.1$ which was suggested by Wei, Hermani, and Haley (2014) to minimize

271   confounding due to "haplotype effects".

272   However, the problem of why and under what conditions additive effects may generate

273   "epistatic signals" has not be formalized. In this work, we use a simple three locus model to reveal

274   the conditions that lead to phantom epistasis. We show that phantom epistasis emerges in the

275   presence of simultaneous but imperfect mutual LD between the three loci (the QTL and the two

276   markers involved in the interaction). This conceptually simple three loci model can be extended

277   to more complex settings (e.g., multiple QTL-marker trios) without affecting the underlying

12

278   source of the principle: if additive QTL variance is imperfectly captured by linear regression on

279   markers and the unexplained variation is not orthogonal to interaction contrasts, then phantom

280   epistasis emerges.

281   ***Testing interactions among weakly correlated SNPs only*** (e.g., considering only SNP-pairs

282   with $r^2 < 0.1$) ***is not a solution***. Our simulations demonstrate that phantom epistasis can

283   emerge even when the two markers involved in the interaction are very weakly correlated. R-

284   squared values greater than 0.05 or even smaller generate strong evidence of phantom epistasis

285   particularly when sample size is large (Figure 3).

286   ***Inferences under imperfect LD***. In a series of recent studies, we (de Los Campos et al. 2013;

287   de los Campos, Sorensen, and Gianola 2015b; Gianola et al. 2015) and others (e.g., M E Goddard

288   2009) have studied the role of imperfect LD on related inferential problems, including missing

289   heritability (i.e., in generating a gap between the trait heritability and the amount of variance

290   that can be captured by a SNP set) and whether imperfect LD can lead to estimates of genomic

291   correlations between traits that are different than the underlying genetic correlations (Gianola

292   et al. 2015). In all these cases, imperfect LD generates inferential difficulties; therefore, phantom

293   epistasis should be seen as one of many issues arising when the markers used for inferences are

294   in imperfect LD with causal variants.

295   ***Perils of Big Data***. The power to detect a small non-null interaction between markers

296   emerging from phantom epistasis increases with sample size. Our simulation results demonstrate

297   this clearly: for pairwise R-sq. between markers of 0.1 there are clear signs of phantom epistasis;

298   however, rejection rates are not highly elevated over the significance level when sample size was

299   moderate (n=10k) because at that R-sq. the size of the interaction effect is small and therefore

300   the power to detect such small interaction effect with moderate sample size is low. Big Data is a

301   blessing for genomic analysis of complex traits; however, some problems cannot be addressed

302   with larger sample size. Moreover, in some cases, large sample size can make an inferential

303   problem even more problematic.

304   ***Dominance can also contribute to phantom epistasis.*** The conceptual and empirical

305   model used in the simulation was based on a purely additive genetic architecture. In the

306   presence of dominance, the true single-locus model becomes $y_i = az_i + dz_i^2 + \delta_i$ where $a$ and

307     $d$ are additive and dominance values, respectively. If the empirical model of expression [2] is used

308     to test for epistatic interactions then the left-hand-side of expression [3] remains unchanged, but

309     the right-hand-side becomes

310
$$\left[\begin{pmatrix} E(x_{1i}z_i) \\ E(x_{2i}z_i) \\ E(x_{1i}x_{2i}z_i) \end{pmatrix} a + \begin{pmatrix} E(x_{1i}z_i^2) \\ E(x_{2i}z_i^2) \\ E(x_{1i}x_{2i}z_i^2) \end{pmatrix} d \right]$$

311     indicating that both dominance and additive effects can contribute to phantom epistasis. The

312     conditions needed for phantom epistasis to emerge are similar to those under the pure additive

313     model. These include, first, imperfect LD between $(z_i, z_i^2)$ and the marker pair $(x_{1i}, x_{2i})$ such that

314     neither $z_i$ nor $z_i^2$ can be fully explained by a linear combination of the two markers. Secondly,

315     phantom epistasis requires mutual LD at the three loci. If one of the markers (say $x_{2i}$) is

316     independent of the other-marker-QTL pair, then, $E(x_{1i}x_{2i}z_i)a + E(x_{1i}x_{2i}z_i^2)d =$

317     $E(x_{2i})[E(x_{1i}x_{2i}z_i)a + E(x_{1i}z_i^2)d] = 0$.

318     ***Local epistasis?*** Several studies have reported results highlighting the importance of 'local'

319     epistatic interactions (e.g., Wei, Hermani, and Haley 2014; He et al. 2017). From a biological

320     perspective it is plausible that multiple mutations in a gene may have collectively a larger impact

321     than the simple sum of the effects of each mutation individually. And this could manifest as

322     "haplotype effects" (e.g., Haig 2011). However, phantom epistasis provides an alternative

323     explanation of why most of the epistatic interactions detected in GWAS occur between loci that

324     are physically close. Indeed, we show analytically and empirically that LD between SNPs is

325     required for phantom epistasis to appear, thus, phantom epistasis is expected to be

326     predominantly a 'local' phenomena.

327     ***The additive-non-additive conundrum***. Quantitative genetics studies properties of

328     complex traits using regression analysis. In the field a careful distinction is made between

329     observable and causal features of complex traits. For instance, it is well established that the

330     linear regression of a phenotype on allele content yields estimates of the average effect of

331     allele substitution and that both truly additive as well as dominance and epistatic effects can

332     contribute to allele substitution effects. Furthermore, theoretical and empirical research has

333     demonstrated that highly non-linear systems can generate signals that can often be

334   explained almost completely with a linear model (Hill, Goddard, and Visscher 2008). For this
335   reason, in general, one cannot make causal statements about gene action from observational
336   variance component analyses (e.g., W. Huang and T. F. C. Mackay, 2016). Complicating
337   matters even further we show in this study that the opposite can happen: under a purely
338   additive model, imperfect LD can generate non-additive signals!

339   The recognition that phantom epistasis may be an important phenomenon does not
340   negate the relevance of gene-gene interactions at the causal level. It simply stresses the
341   difficulties that one faces when trying to learn about causal features of a system using
342   observational data and inputs (markers) which are proxies for the underlying variants that
343   may have causal effects on traits.

344   ***Phantom epistasis: an opportunity to improve predictive performance?*** In this work
345   we have stressed that imperfect LD can limit the possibility to learn about causal effects.
346   However, linear and non-linear genomic regressions can be very powerful predictive machines,
347   and it is well-established that the model that is best for inferences is not necessarily the best
348   predictive tool. Phantom epistasis creates inferential problems but also opens opportunities for
349   improving prediction models. Indeed, by capturing signals that are missed by an additive model,
350   non-linear models using interactions between markers may increase the amount of genetic
351   variance captured and improve prediction accuracy. This may explain, for instance why some
352   non-linear models such as kernel regressions have shown better predictive performance than
353   additive models, especially in breeding populations with long-span LD and low marker density
354   (de los Campos et al. 2010).

355

363 at the 2017 EAAP (European Federation of Animal Science) meetings. We are grateful for the

364 comments received and would like to particularly thank the feedback offered by Johanes

365 Martini.

## Literature Cited

367

368 Aschard, H. 2016. "A Perspective on Interaction Effects in Genetic Association Studies." *Genetic*

369 *Epidemiology*. https://doi.org/10.1002/gepi.21989.

370 Bennett, J H. 1954. "On the Theory of Random Mating." *Annals of Eugenics* 18 (4):311–17.

371 http://www.ncbi.nlm.nih.gov/pubmed/13148997.

372 Cordell, H J. 2002. "Epistasis: What It Means, What It Doesn't Mean, and Statistical Methods to

373 Detect It in Humans." *Human Molecular Genetics* 11:2463–68.

374 ———. 2009. "Detecting Gene-Gene Interactions That Underlie Human Diseases." *Nature*

375 *Reviews Genetics* 10:392–404.

376 Gianola, D., G. de los Campos, M. A. Toro, H. Naya, C.-C. Schon, and D. Sorensen. 2015. "Do

377 Molecular Markers Inform About Pleiotropy?" *Genetics* 201 (1):23–29.

378 https://doi.org/10.1534/genetics.115.179978.

379 Goddard, M E. 2009. "Genomic Selection: Prediction of Accuracy and Maximisation of Long

380 Term Response." *Genetica* 136:245–52.

381 Grueneberg, Alexander, and Gustavo de los Campos. 2017. "BGData: A Suite of R-Packages for

382 Analysis of Big Genomic Data [v 1.0.0]." Comprehensive R Archive Network (CRAN).

383 https://cran.r-project.org/web/packages/BGData/index.html.

384 Haig, D. 2011. "Does Heritability Hide in Epistasis between Linked SNPs?" *European Journal of*

385 *Human Genetics* 19:123.

386 He, Sang, Jochen C. Reif, Viktor Korzun, Reiner Bothe, Erhard Ebmeyer, and Yong Jiang. 2017.

387 "Genome-Wide Mapping and Prediction Suggests Presence of Local Epistasis in a Vast Elite

388 Winter Wheat Populations Adapted to Central Europe." *Theoretical and Applied Genetics*

389 130 (4). Springer Berlin Heidelberg:635–47. https://doi.org/10.1007/s00122-016-2840-x.

390 Hermani, G, K Shakhbazov, H J Westra, T Esko, A K Henders, A F McRae, J Yang, et al. 2014.

391      "Detection and Replication of Epistasis Influencing Transcription in Humans." *Nature* 508.

392 Hill, W G, M E Goddard, and P M Visscher. 2008. "Data and Theory Point to Mainly Additive

393      Genetic Variance for Complex Traits." In *PLos Genetics*.

394 Huang, A, S Xu, and X Cai. 2014. "Whole-Genome Quantitative Trait Locus Mapping Reveals

395      Major Role of Epistasis on Yield of Rice." *Plos One*.

396 los Campos, Gustavo de, Daniel Gianola, Guilherme J. M. Rosa, Kent A. Weigel, and José Crossa.

397      2010. "Semi-Parametric Genomic-Enabled Prediction of Genetic Values Using Reproducing

398      Kernel Hilbert Spaces Methods." *Genetics Research* 92 (04). Cambridge University

399      Press:295–308. https://doi.org/10.1017/S0016672310000285.

400 los Campos, Gustavo de, Daniel Sorensen, and Daniel Gianola. 2015a. "Genomic Heritability:

401      What Is It?" Edited by Gregory S. Barsh. *PLOS Genetics* 11 (5):e1005048.

402      https://doi.org/10.1371/journal.pgen.1005048.

403 ———. 2015b. "Genomic Heritability: What Is It?" Edited by Gregory S. Barsh. *PLOS Genetics* 11

404      (5):e1005048. https://doi.org/10.1371/journal.pgen.1005048.

405 Los Campos, Gustavo de, Ana I Vazquez, Rohan Fernando, Yann C Klimentidis, and Daniel

406      Sorensen. 2013. "Prediction of Complex Human Traits Using the Genomic Best Linear

407      Unbiased Predictor." Edited by Michael E. Goddard. *PLoS Genetics* 9 (7):e1003608.

408      https://doi.org/10.1371/journal.pgen.1003608.

409 Mackay, Trudy F C. 2014. "Epistasis and Quantitative Traits: Using Model Organisms to Study

410      Gene-Gene Interactions." *Nature Reviews. Genetics* 15 (1):22–33.

411      https://doi.org/10.1038/nrg3627.

412 Manolio, T A, F S Collins, N J Cox, D B Goldstein, L A Hindorff, D J Hunter, M I McCarthy, E M

413      Ramos, L R Cardon, and \textit{et. al}. 2009. "Finding the Missing Heritability of Complex

414      Diseases." *Nature* 461:747–53.

415 R Development Core Team. 2012. "R: A Language and Environment for Statistical Computing."

416      Vienna, Austria. http://www.r-project.org/.

417 Strange, T, B Ask, and B Nielsen. 2013. "Genetic Parameters of the Piglet Mortality Traits

418      Stillborn, Weak at Birth, Starvation, Crushing, and Miscellaneous in Crossbred Pigs."

17

419        *Journal of Animal Science* 91:1562–69.

420    Sudlow, Cathie, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul

421        Downey, et al. 2015. "UK Biobank: An Open Access Resource for Identifying the Causes of a

422        Wide Range of Complex Diseases of Middle and Old Age." *PLoS Medicine* 12 (3). Public

423        Library of Science:e1001779. https://doi.org/10.1371/journal.pmed.1001779.

424    Wang, X, R C Elston, and X Zhu. 2010. "The Meaning of Interaction." *Human Heredity* 70:269–

425        77.

426    Wei, W.-H., G Hermani, and C S Haley. 2014. "Detecting Epistasis in Human Complex Traits."

427        *Nature Reviews Genetics* 15:722–33.

428    Wood, R W, M A Tuke, M A Nalls, D G Hernandez, S Bandinelli, A B Singleton, D Melzer, L

429        Ferrucci, T M Frayling, and M N Weedon. 2014. "Another Explanation for Apparent

430        Epistasis." *Nature* 508.

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

## Supplementary Methods

446

447

448 **1. Equivalence between genotype and haplotype measures of LD**

449 In this section we present the standard haplotype two- and three-loci measures of LD and

450 establish the connection between these measures and the genotype moments involved in

451 expression [3].

452 ***Two-loci haplotype measure of LD.*** Consider a pair of bi-allelic loci (A and B, with alleles $A_1/A_2$

453 and $B_1/B_2$, respectively). The haplotype linkage disequilibrium parameter is $D_{AB} = p(A_1, B_1) -$

454 $p(A_1)p(B_1)$. Let $X = 1$ when allele $A_1$ is present and $X = 0$ when allele $A_2$ is present. Likewise,

455 let $Y = 1$ when allele $B_1$ is present and $Y = 0$ when allele $B_2$ is present. Then $E(X) = p(A_1)$,

456 $E(Y) = p(B_1)$, $E(XY) = p(A_1, B_1)$ and the covariance between $X$ and $Y$ is $Cov(X, Y) =$

457 $E(XY) - E(X)E(Y)$ which reduces to $D_{AB}$, thus

458 $$D_{AB} = Cov(X, Y) = E(XY) - E(X)E(Y)$$

459 If the two genotypes are centered, then $E(X) = E(Y) = 0$ and

460 $$D_{AB} = Cov(X, Y) = E(XY) \qquad\qquad [4]$$

461 This is a haplotype analog of the 1$^{st}$ order measures of LD entering in expression [3].

462

463 ***Three-loci haplotype measure of LD***. For a system involving three bi-allelic loci ($A$, $B$ and $C$, with

464 alleles $A_1/A_2$, $B_1/B_2$, and $C_1/C_2$, respectively) a three-loci haplotype measure of LD can be defined

465 as (Bennett 1954)

466 $$D_{ABC} = p(A_1, B_1, C_1) - p(A_1)D_{BC} - p(B_1)D_{AC} - p(C_1)D_{AB} - p(A_1)p(B_1)p(C_1).$$

467 Extending the two-loci system by introduction of three binary random variables $X$, $Y$, and $Z$, that

468 take values 1 when the allelic forms $A_1$, $B_1$ and $C_1$ are present, respectively, and take values 0

469 othewise, yields

470 $$D_{ABC} = E(XYZ) - E(X)[E(YZ) - E(Y)E(Z)]$$

471 $$-E(Y)[E(XZ) - E(X)E(Z)]$$

472 $$-E(Z)[E(XY) - E(X)E(Y)]$$

473 $$-E(X)E(Y)E(Z)$$

19

474    The three terms in square brackets represent pairwise disequilibria. When the three random

475    variables are centered their marginal expectations are zero and the expression reduces to

476 $$D_{ABC} = E(XYZ). \qquad [5]$$

477    ***Relationship with genotype measures of LD***. The disequilibria measures described above involve

478    associations between alleles within gametes, whereas in the body of the paper, the expectations

479    involve different genotypes. Assuming random mating, the expectations involving genotypes

480    result in twice those involving gametes.

481

482    **2. Perfect LD between markers and QTL prevents phantom epistasis**

483    We demonstrate (the very intuitive result) that if a response $(z_i)$ can be fully captured by

484    regression on a set of predictors $(x_i)$, then the regression of $z_i$ on $x_i$ plus $w_i$,

485 $$z_i = x_i'b_X + w_i'b_W + \psi, \qquad [6]$$

486    yields $b_W = 0$ in the population.

487    ***Demonstration:*** In the population, the regression coefficients of [6], are defined by the following

488    system

489 $$\begin{bmatrix} \Sigma_X & \Sigma_{XW} \\ \Sigma_{WX} & \Sigma_W \end{bmatrix} \begin{bmatrix} b_X \\ b_W \end{bmatrix} = \begin{bmatrix} \Sigma_{Xz} \\ \Sigma_{Wz} \end{bmatrix}$$

490    Where the $\Sigma$.'s represent covariance matrices: $\Sigma_X = Cov(x_i, x_i')$, $\Sigma_W = Cov(w_i, w_i')$ and

491    $\Sigma_{XW} = Cov(x_i, w_i') = \Sigma_{WX}$. It follows that

492 $$\Sigma_X b_X + \Sigma_{XW} b_W = \Sigma_{Xz} \qquad [7]$$

493    and

494 $$\Sigma_{WX} b_X + \Sigma_W b_W = \Sigma_{Wz}. \qquad [8]$$

495    Solving [7] for $b_X$ yields $b_X = \Sigma_X^{-1}(\Sigma_{Xz} - \Sigma_{XW} b_W)$. Plugging this into [8] yields,

496    $\Sigma_{WX}\Sigma_X^{-1}(\Sigma_{Xy} - \Sigma_{XW} b_W) + \Sigma_W b_W = \Sigma_{Wz}$. And solving for $b_W$ gives

497 $$b_W = [\Sigma_W - \Sigma_{WX}\Sigma_X^{-1}\Sigma_{XW}]^{-1}[\Sigma_{Wz} - \Sigma_{WX}\Sigma_X^{-1}\Sigma_{Xz}]. \qquad [9]$$

20

498    Now if $z_i$ can be fully explained by regression on $\boldsymbol{x}_i$, that is if $z_i = \boldsymbol{x}_i'\boldsymbol{a}$, with $\boldsymbol{a} = \boldsymbol{\Sigma}_X^{-1}\boldsymbol{\Sigma}_{Xz}$,

499    then, $\boldsymbol{\Sigma}_{Wz} = Cov(\boldsymbol{w}_i, \boldsymbol{x}_i'\boldsymbol{a}) = \boldsymbol{\Sigma}_{WX}\boldsymbol{a} = \boldsymbol{\Sigma}_{WX}\boldsymbol{\Sigma}_X^{-1}\boldsymbol{\Sigma}_{Xz}$, thus, $[\boldsymbol{\Sigma}_{Wz} - \boldsymbol{\Sigma}_{WX}\boldsymbol{\Sigma}_X^{-1}\boldsymbol{\Sigma}_{Xz}] = \boldsymbol{0}$ and

500    therefore, $\boldsymbol{b}_W = \boldsymbol{0}$. QED.

501    **Implication.** Let $z_i$ be the QTL genotype, $\boldsymbol{x}_i = (x_{1i}, x_{2i})'$ be a vector containing the two marker

502    genotypes and $\boldsymbol{w}_i = x_{1i}x_{2i}$ be the two-marker interaction contrast. Under perfect LD between

503    the QTL and the markers the QTL genotype can be fully explained by linear regression on the two

504    markers, that is $z_i = \boldsymbol{x}_i'\boldsymbol{a}$. Therefore, the above result ($\boldsymbol{b}_W = \boldsymbol{0}$) implies $\beta_{12} = 0$, i.e., absence of

505    phantom epistasis.

506