1  **On the relationship between high-order linkage disequilibrium and**
2  **epistasis**

3

4  Yanjun Zan[1*], Simon K. G. Forsberg[1*+] and Örjan Carlborg[1§]

5  [1]Department of Medical Biochemistry and Microbiology, Uppsala University, SE-751
6  23 Uppsala, Sweden

7

8  [*]Authors contributed equally
9  [+]Present address: Ecology and Evolutionary Biology Department; Lewis Sigler
10  Institute for Integrative Genomics, Princeton University, Princeton, New Jersey,
11  08540 and Department of Neuroscience, Functional Pharmacology, Uppsala
12  University, BMC, Box 593, 751 24 Uppsala, Sweden
13  [§]Corresponding author: orjan.carlborg@imbim.uu.se
14  Örjan Carlborg
15  Uppsala University
16  Medical Biochemistry and Microbiology
17  BMC Box 582, SE-751 23, Uppsala, Sweden.
18  Phone: +46 18 4714592

19

22

23  Data availability:
24  i) Genome wide re-sequencing data are available as part of *Arabidopsis thaliana* 1001
25  genomes project http://1001genomes.org/data-center.html, ii) 250 K SNP chip data
26  were available in the *Arabidopsis thaliana* Regmap panel
27  (http://bergelson.uchicago.edu/?page_id=790), iii) The Molybdenum level for 340
28  *Arabidopsis thaliana* are available in
29  https://doi.org/10.1371/journal.pgen.1005648.s005. iv) Corresponding genotypes are
30  extracted from a subset of the Regmap panel
31  (http://bergelson.uchicago.edu/?page_id=790).

32

33  Running title: High-order LD and statistical epistasis

1

## ABSTRACT

34

35  A plausible explanation for statistical epistasis revealed in genome wide association

36  analyses is the presence of high order linkage disequilibrium (LD) between the

37  genotyped markers tested for interactions and unobserved functional polymorphisms.

38  Based on findings in experimental data, it has been suggested that high order LD

39  might be a common explanation for statistical epistasis inferred between local

40  polymorphisms in the same genomic region. Here, we empirically evaluate how

41  prevalent high order LD is between local, as well as distal, polymorphisms in the

42  genome. This could provide insights into whether we should account for this when

43  interpreting results from genome wide scans for statistical epistasis. An extensive and

44  strong genome wide high order LD was revealed between pairs of markers on the high

45  density 250k SNP-chip and individual markers revealed by whole genome sequencing

46  in the *A. thaliana* 1001-genomes collection. The high order LD was found to be more

47  prevalent in smaller populations, but present also in samples including several

48  hundred individuals. An empirical example illustrates that high order LD might be an

49  even greater challenge in cases when the genetic architecture is more complex than

50  the common assumption of bi-allelic loci. The example shows how significant

51  statistical epistasis is detected for a pair of markers in high order LD with a complex

52  multi allelic locus. Overall, our study illustrates the importance of considering also

53  other explanations than functional genetic interactions when genome wide statistical

54  epistasis is detected, in particular when the results are obtained in small populations of

55  inbred individuals.

56

## INTRODUCTION

The genetic architecture of most biological traits is complex and involves multiple genes, whose effects are often influenced by interactions with other genes and environmental factors. To study the relative contributions by genes, environmental factors and their interactions in segregating populations, statistical genetic approaches are commonly used to partition the genetic variance to additive and dominance variance of individual loci and epistatic interaction variance between them (Lynch and Walsh 1998). In principle, the variance partitioning is performed by associating the phenotypic variation for a trait in a population with linear combinations of the genotypes within and/or across loci. How the genotypes are combined (parameterized) in the model is determined by the genetic model used in the analysis. The classic quantitative genetics models are parameterized to capture the genetic variance in a hierarchical manner. First, a main additive allele-substitution is defined. Then, if accounted for, dominance is modeled as a single-locus deviation from additivity and genetic interactions as multi-locus deviations from single locus additivity and dominance (Nelson *et al.* 2013). As a consequence of this, the genetic contributions of individual and combinations of loci described as additive, dominance and epistatic variances are unlikely to reflect the underlying biological mechanisms (Carlborg *et al.* 2006; Phillips 2008; Huang *et al.* 2012; Sackton and Hartl 2016; Forsberg *et al.* 2017).

Although the ultimate aim of a genetic association study is generally to detect functional polymorphisms, most often genotypes are only scored for a reduced set of polymorphisms (genetic markers). These reduced marker sets are selected with the aim to tag as many of the unobserved functional polymorphisms as possible. The statistical inferences of the underlying genetic architecture made from such reduced sets of markers can, however, be problematic in some cases. For example, multiple unobserved functional polymorphisms can lead to associations to individual markers that do not properly represent the causal variants (Platt *et al.* 2010), and high order linkage disequilibrium (LD) to single functional polymorphism can lead to indirect statistical epistatic associations to pairs of markers (Wood *et al.* 2014). Here, we focus on high order linkage disequilibrium defined as when two genotyped markers tag an un-genotyped polymorphism (see Materials and Methods section). It is still unknown how prevalent and strong such high order LD is in the genome, making it difficult to

3

90    estimate how many reported pairwise statistical epistatic interactions are due to such
91    LD. However, the study by *Wood et al* (Wood *et al.* 2014) presented results
92    suggesting that many of the significant statistical epistatic interactions detected
93    between pairs of local markers by Hemani *et al.* (Hemani *et al.* 2014) might be due to
94    high-order LD to unobserved, linked sequence polymorphisms in the same genomic
95    region. Many past and current studies of genetic interactions in, for example,
96    Drosophila, plant, animal and human populations (Shimomura *et al.* 2001; Anholt *et*
97    *al.* 2003; Caicedo *et al.* 2004; Segrè *et al.* 2004; Carlborg *et al.* 2006; Hemani *et al.*
98    2014) rely on genome-wide statistical analyses of pairwise interactions between
99    selected sets of markers as in (Hemani *et al.* 2014). With the increasing interest in,
100   and availability of, sufficiently large datasets for epistatic association analyses it is
101   therefore important to also evaluate the risk of making false inferences about loci
102   being involved in functional genetic interactions from findings of statistical epistasis,
103   when they instead are due to high order LD.

104

105   Here, we empirically explore the prevalence and strength of high order LD within and
106   between chromosomes in publically available high-density SNP and whole-genome
107   re-sequencing data from the model plant *Arabidopsis thaliana*. Two locus LDs are
108   calculated between the markers selected for the 250k *A. thaliana* SNP chip that have
109   been the basis for many GWAS analyses in the past, and the additional SNPs revealed
110   by whole genome sequencing using data from the 1001 genomes project (Atwell *et al.*
111   2010; Cao *et al.* 2011; Horton *et al.* 2012; Schmitz *et al.* 2013; Alonso-Blanco *et al.*
112   2016). Strong high order LD was found to be common both within and across
113   chromosomes between pairs of markers from the SNP-chip and the sequencing
114   polymorphisms and often the combined genotype of the marker pair tagged the
115   genotype of the sequencing markers better than any single marker on the SNP chip.
116   The risk of falsely inferring genetic interactions between markers on different
117   chromosomes in a two-locus interaction analysis might increase in situations when the
118   underlying genetic architecture is more complex, for example when a single locus
119   contains multiple functional alleles. This is illustrated using an empirical example
120   from a second public *A. thaliana* dataset (Forsberg *et al.* 2015). Overall, this study
121   provides new insights that deepen our understanding about the link between high
122   order LD and statistical epistasis to guide researchers when interpreting results

4

123    obtained from epistatic genetic association analyses.

124

125                    **MATERIALS AND METHODS**

126

127    **Methods**

128    When an individual marker is in complete linkage disequilibrium ($r^2 = 1$) with a

129    functional polymorphism affecting a studied trait, a single-locus association test

130    between the marker and the trait will capture all the phenotypic variance contributed

131    by the functional polymorphism. A basic assumption in genetic association studies is

132    that at least one genotyped marker will be in sufficiently high LD with each functional

133    polymorphism to detect it in this way. In reality, however, not all functional

134    polymorphism will be in such perfect LD with a genotyped marker, and then there is a

135    risk that the joint genotype of two (or more) markers tags the genotype of the

136    functional polymorphism better than any single marker (high-order LD > single-

137    marker LD). This will, as discussed below, influence the significances of the trait-

138    marker associations detected in a genetic association analysis and the inferences made

139    about the genetic architecture of the trait.

140

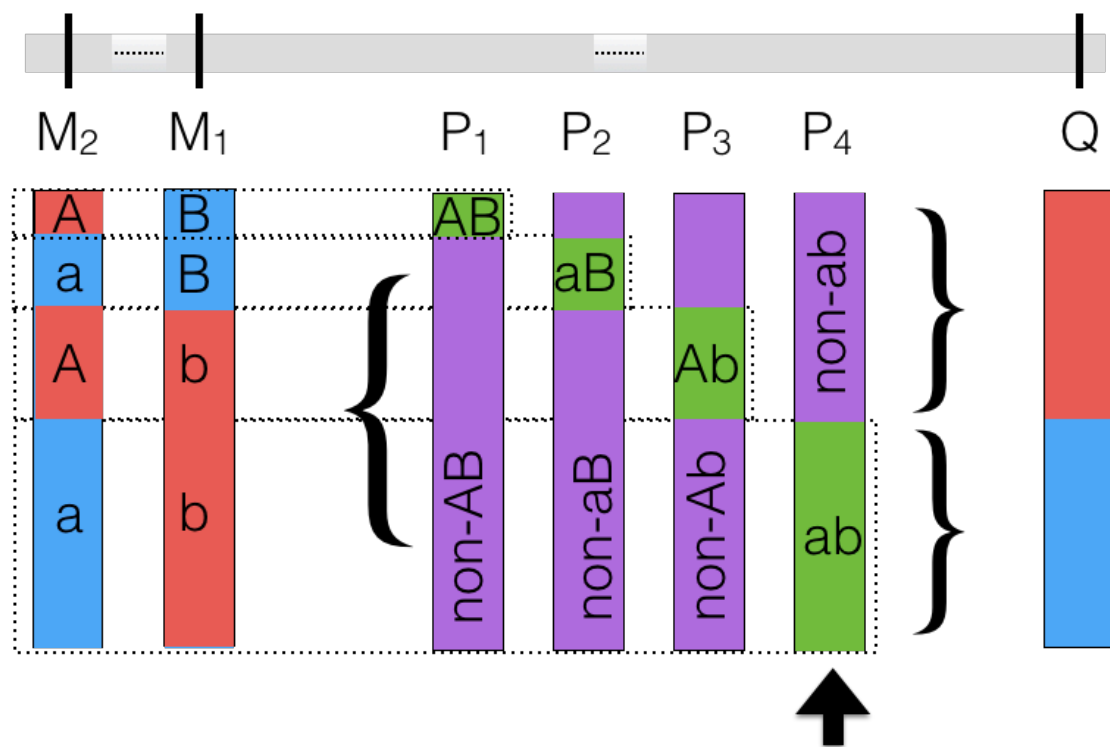141    *Quantifying high order linkage disequilibrium*

142    We calculate the high order LD between pairs of predictors (here genotyped SNP

143    markers) and single targets (here un-genotyped SNP polymorphisms) following (Hao

144    *et al.* 2007).

145

146    Consider a pair of bi-allelic predictor SNPs ($M_1$ and $M_2$; Figure 1). These markers can

147    together form four two-locus genotypes: AB, Ab, aB and ab (Figure 1). We now want

148    to know whether any two-locus predictor could tag the single locus target genotype

149    better than any of the individual predictor genotypes (i.e. evaluate whether

150    max(second-order LD) > max(single order LD)). To calculate the high order LD

151    between the two predictors ($M_1$ and $M_2$) and the single target (Q), the two-locus

152    $M_1M_2$ genotype is used to create a multi-allelic pseudo marker (P) with four alleles

153    (Figure 1). In this way, a second-order LD ($r^2$) can be calculated for each of the

154    possible ways that $M_1$ and $M_2$ together can tag the genotype at Q (Figure 1).

155

        5

156    The calculation of the second order LD therefore first involves creating the four

157    possible bi-allelic pseudomarkers ($P_1$, $P_2$, $P_3$ & $P_4$; Figure 1) from the two locus

158    $M_1/M_2$ genotypes. These are  assigned the genotypes $P_1${AB, non-AB}, $P_2${Ab, non-

159    Ab}, $P_3${aB, non-aB} and $P_4${ab, non-ab}, respectively. The LD-$r^2$ is then computed

160    between the target (Q) and the four bi-allelic pseudomarkers ($P_1$, $P_2$, $P_3$ & $P_4$). For

161    each pair of predictors, the second order LD is then defined as the LD-$r^2$ for the

162    pseudomarker with the highest LD-$r^2$ to the target. Pseudomarkers with higher LD-$r^2$

163    to the target (Q) than 0.3 are kept for further analyses. The LD-$r^2$ values were

164    computed using the software *LdCompare* (Hao *et al.* 2007).



165

166    ***Figure 1.*** *Illustration of how the pseudomarkers ($P_1$, $P_2$, $P_3$, $P_4$) used in the estimation of the second*

167    *order linkage disequilibrium between a pair of linked or unlinked markers (predictors; $M_1$ and $M_2$,)*

168    *and a third linked or unlinked functional polymorphism (target; Q) are created. The pseudomarkers*

169    *together represent the possible bi-allelic formulations of the two-locus $M_1M_2$ genotypes. The maximum*

170    *pairwise LD-$r^2$ between the target and the four pseudomarkers ($P_4$) defines the second order LD*

171    *between the predictors ($M_1$, $M_2$) and the target (Q).*

172

173

174    *Statistical epistasis emerging from high order linkage disequilibrium*

175    In a genetic association study in an inbred or haploid population, two-locus epistasis

6

176     is typically modelled as:

177

178     $Y = a_1\beta_1 + a_2\beta_2 + a_1a_2\beta_{12} + e$                    [1]

179

180     Here, $a_1$ and $a_2$ are indicator variables for the genotypes at two genotyped markers, $M_1$

181     and $M_2$, taking values 1/-1 for the two alternative homozygous genotypes AA vs aa

182     and BB vs bb, respectively. $a_1a_2$ is an indicator variable for the interaction between $a_1$

183     and $a_2$ taking value 1 for the two-locus genotypes AABB and aabb and -1 for AAbb

184     and aaBB. $\beta_1$, $\beta_2$ and $\beta_{12}$ are the corresponding estimates for the marginal (additive)

185     effects and the additive-by-additive interaction between the loci.

186

187     The aim of a statistical epistatic analysis is to include an interaction term in the model

188     [1] to estimate the deviations of the two-locus genotype-values (AABB, AAbb, aaBB

189     and aabb) from the predictions obtained by the marginal (additive) effects (Alvarez-

190     Castro and Carlborg 2007). However, a non-zero estimate of the interaction term in

191     model [1] does not, as noted e.g. by Wood et al. (Wood *et al.* 2014) necessarily have

192     to result from a genetic interaction. It could, for example, instead emerge from a

193     second-order LD between two markers and a single functional polymorphism. Here,

194     refer back to Figure 1. Now assume that a trait is determined by a single functional

195     locus (Q). Two markers, $M_1$ and $M_2$, are genotyped but neither of these markers

196     individually tag the causal genotype (blue) at Q well. However, the causal (blue)

197     allele at Q is,tagged perfectly by one of the two-locus $M_1M_2$ genotypes (ab; Figure 1),

198     while the other three $M_1M_2$ two-locus genotypes (aB, Ab and AA; Figure 1) are only

199     present together with the no-effect (red) allele at locus Q. When fitting model [1] to

200     the genotypes of marker $M_1$ and $M_2$, the estimate for the interaction term ($\beta_{12}$) will be

201     non-zero, illustrating how statistical epistasis can emerge from the second-order LD

202     between $M_1$ and $M_2$ and Q. This example illustrates a scenario similar to what was

203     empirically observed in (Wood *et al.* 2014), where physically linked markers in low

204     LD with each other tagged   haplotypes that were in high order LD with a

205     polymorphism that was unobserved in the original study.

206

207     *Classifying identified high order linkage disequilibrium triplets depending on the*

208     *distance between the loci*

        7

209    Here, we evaluate the prevalence and strength of high order LD between pairs of

210    markers selected for genotyping on a 250k SNP chip (predictors) and a third locus

211    revealed by whole genome sequencing (targets) using publicly available datasets in *A.*

212    *thaliana* (Cao *et al.* 2011; Alonso-Blanco *et al.* 2016). Three types of high order LD

213    are defined based on the locations of the predictors relative to the target. If both

214    predictors are located within 1Mb of the target it is classified as cis-cis. If only one

215    predictor is closer than 1Mb it is classified as cis-trans. If none is closer than 1Mb it is

216    classified as trans-trans. The choice of a 1Mb threshold to define cis vs trans

217    predictors is arbitrary, but we consider it useful for evaluating how common high

218    order LD is between predictors near (local/cis) and far (global/trans) from the target.

219

**Material**

221

222    *The genome wide prevalence of high order linkage disequilibrium in publically*

223    *available* Arabidopsis thaliana *datasets*

224    The *A. thaliana* 1001-genomes project has released complete genome sequences for

225    hundreds of wild collected accessions (http://www.1001genomes.org). Here, we used

226    whole-genome SNP data on 728 accessions scored by whole genome re-sequencing

227    (Cao *et al.* 2011; Alonso-Blanco *et al.* 2016). The predictors used in our analysis was

228    a subset of the SNPs selected for the 250k *A. thaliana* SNP chip (Horton *et al.* 2012)

229    (n = 200,352 in total; MAF > 0.05) and the targets a subset of the SNPs revealed

230    using whole-genome re-sequencing (n = 1,641,240 in total; MAF > 0.05) (Table 1).

231    Although the results from the analyses of this data will be specific to this species and

232    dataset, it is assumed that the relationships between targets and predictors will be a

233    realistic representation of what to be expected also in other populations. This is

234    because the selection of markers for the high-density 250k SNP chip, was done for the

235    purpose of genetic association studies following similar procedures as used also in

236    other species and populations.

237

238    The reason for only studying a subset of the possible targets and predictors is that it

239    was not computationally feasible to exhaustively evaluate the high order LD between

240    all possible pairs of predictors selected for the 250k SNP chip and all the targets

241    revealed by genome sequencing. Instead, the second order LD was exhaustively

8

242 calculated for all targets and predictors i) within a randomly selected 6 Mb window on

243 chromosome 2 as well as ii) between three randomly selected windows from different

244 chromosomes (Table 1). Computations were performed for the entire population (n =

245 728 individuals) and two smaller random samples of n = 100 and n = 50 individuals.

246 The results for the populations with n = 100 and n = 728 are reported in the main

247 manuscript and the results for n = 50 is reported in the Supplementary material.

248

249 **Table 1.** *Regions and SNPs selected for evaluation of second order LD.*

|  | Window 1 | Window 2 | Targets[1] | Predictors[2] | Filtered targets[3] |
|---|---|---|---|---|---|
| Region 1 | Chr2: 8-14Mb | - | 70,712 | 6,053 | - |
| Regionpair 1 | Chr1: 10-12Mb | Chr3: 10-12Mb | 29,133 | 6,245 | 20,239 |
| Regionpair 2 | Chr2: 10-12Mb | Chr4: 10-12Mb | 23,751 | 5,302 | 15,887 |
| Regionpair 3 | Chr2: 10-12Mb | Chr3: 10-12Mb | 23,751 | 5,212 | 15,884 |
| Genome |  |  | 1,486,942 | 154,298 | 1,229,012 |

250 [1]*Total number of polymorphic SNPs in the evaluated windows/genome in the population revealed via*

251 *whole-genome re-sequencing (Alonso-Blanco et al. 2016).* [2]*Total number of polymorphic SNPs in the*

252 *two windows/genome included on the 250k AT SNP-chip (Horton et al. 2012);* [3]*Number of target SNPs*

253 *in the two windows/genome with LD-$r^2$ < 0.6 to any individual predictor.*

254

255 The predictor pairs in the evaluated windows in the genome with high order LD-$r^2$ >

256 0.6 to a target were classified as cis-cis/cis-trans/trans-trans. To extrapolate these

257 findings to the genome level, the proportions of all evaluated predictor pairs that

258 displayed these patterns were calculated and then multiplied with the total number of

259 possible cis-cis/cis-trans/trans-trans pairs in the genome (Table S1).

260

261 *Analyzing a public* A. thaliana *dataset for two locus statistical epistasis*

262 A publicly available dataset including 340 *Arabidopsis thaliana* accessions were used

263 for a genome wide association analysis. In short, the plants were grown in a controlled

264 environment with 6 biological replicate plants per accession. Analyses by Inductively

265 Coupled Mass Spectroscopy (ICP-MS) provided estimates of leaf molybdenum

266 concentration as described in (Baxter *et al.* 2010; Forsberg *et al.* 2015). The

267 accessions were genotyped for 141,385 SNP markers with MAF > 0.15 (Atwell *et al.*

268 2010; Baxter *et al.* 2010; Shen *et al.* 2012; Forsberg *et al.* 2015). A more thorough

269 description of the dataset can be found in (Baxter *et al.* 2010; Forsberg *et al.* 2015). In

270 an earlier study of this dataset (Forsberg *et al.* 2015), it was revealed that a large

9

271    fraction of the genetic variance for this trait was explained by a single linkage block

272    containing several low-frequency, large effect structural variants that were poorly

273    tagged by the genotyped SNPs. This linkage block was originally identified due to its

274    large marginal, variance heterogeneity effect in the population (Shen *et al.* 2012). It is

275    known that statistical epistasis and genetic variance heterogeneity can emerge from

276    similar genetic architectures (Forsberg *et al.* 2015),  and this population was therefore

277    selected for further evaluations of whether high order LD between the genotyped

278    SNPs and these hidden polymorphisms could lead to statistical epistasis in a two locus

279    association analysis. We performed an exhaustive, two-dimensional genome scan for

280    pairwise statistical epistasis between the genotyped markers and the level of

281    molybdenum in the leaf using the software *plink* (Purcell *et al.* 2007) without control

282    for population structure. Thereafter, each pair of loci that passed the genome wide

283    significance threshold in the initial scan was fitted in a two-locus epistatic genetic

284    model [1] using *hglm* function in *hglm* package (Rönnegård *et al.* 2010) to correct for

285    the possible effects of population structure via the genomic kinship matrix as in

286    (Forsberg *et al.* 2015). The significance threshold used to infer significant interacting

287    pairs ($p < 3.2 \times 10^{-10}$) was defined as a Bonferroni corrected nominal 5% significance

288    threshold. The correction was done for an estimated number of independent

289    association tests assumed to equal the number of independent LD blocks in the *A.*

290    *thaliana* genome as described in (Lachowiec *et al.* 2015).

291

292    *Data availability*

293    Genome wide re-sequencing data are available as part of the *Arabidopsis thaliana*

294    1001 genomes project http://1001genomes.org/data-center.html. The 250 K SNP chip

295    data are available as part of the genotype data for the *Arabidopsis thaliana* Regmap

296    panel (http://bergelson.uchicago.edu/?page_id=790). The Molybdenum levels for the

297    340 *Arabidopsis thaliana* accessions are available in

298    https://doi.org/10.1371/journal.pgen.1005648.s005
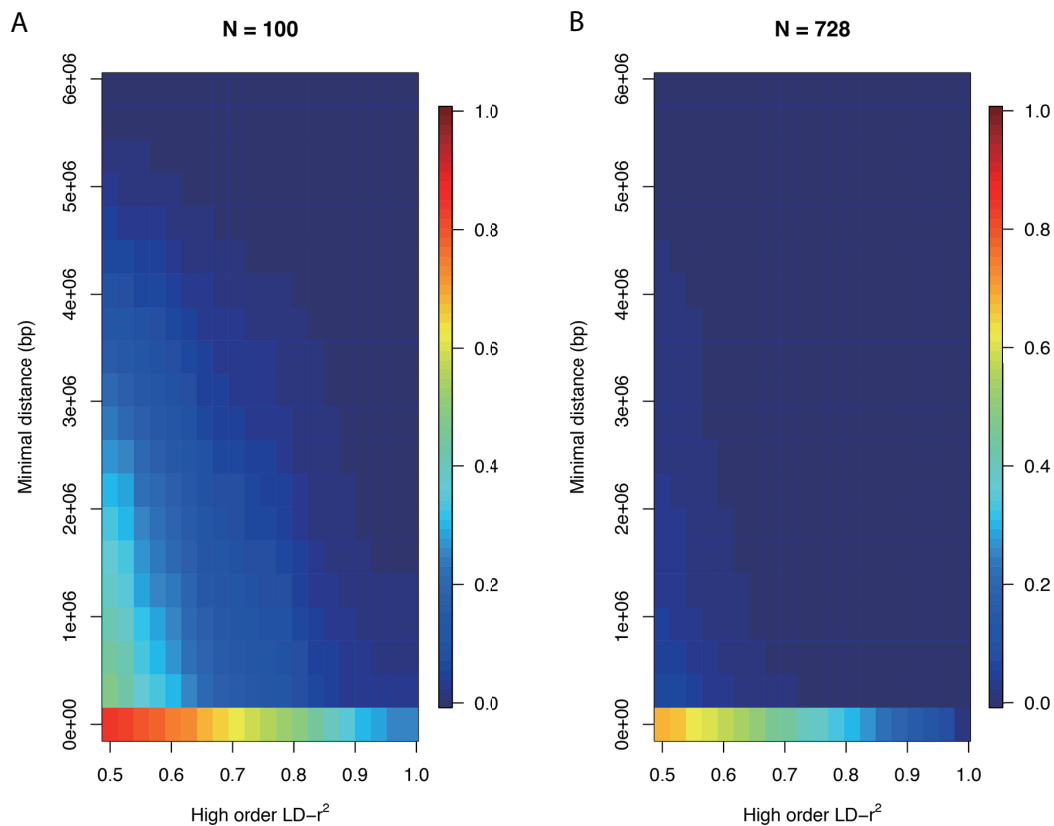
299

300

301                           **RESULTS**

302    This study aims to answer the following questions by analyzing two public *A.*

303    *thaliana* datasets: How common can we expect high order LD to be between pairs of

10

304     SNPs selected for genotyping and hidden sequence variants in the genome? Is high

305     order LD primarily observed between predictors tightly linked to a target functional

306     polymorphism (in cis) as in (Wood *et al.* 2014), or is it also observed for predictors

307     unlinked to the target (in trans)? How dependent is the prevalence of high order LD

308     and cis vs trans predictors on the population size? We also present an empirical

309     example where high order LD exists between a cis-trans predictor pair with

310     significant statistical epistasis and a locus displaying a strong genetic variance

311     heterogeneity due to independent contributions by multiple linked polymorphisms

312     (Forsberg *et al.* 2015). This illustrates how complex inheritance patterns of individual

313     loci, something usually not explored in GWAS data, further complicates the

314     interpretation of detected statistical epistatic signals.

315

316     *The population size affects the prevalence and location of predictors in high order LD*

317     The high order LD-$r^2$ values for all pairs of predictors and individual targets in a 6Mb

318     window on Chromosome 2 (Table 1) is shown for populations with n = 100 and n =

319     728 individuals in Figure 2. The strongest second order LD-$r^2$ was observed where at

320     least one predictor is located near the target (y-axis). When the sample size was

321     smaller (n = 100; Figure 2A), strong second order LD-$r^2$ was rather common also

322     when both predictors were located far from the target. For example, 20% of the

323     targets had a high order LD-$r^2$ > 0.65 with a predictor pair where at least one of the

324     predictors was located more than 1Mb away from it. Even though the prevalence of

325     strong high order LD-$r^2$ decreases when the sample size increases, it is still common

326     in the large population (n = 728; Figure 2B), with the highest prevalence when at least

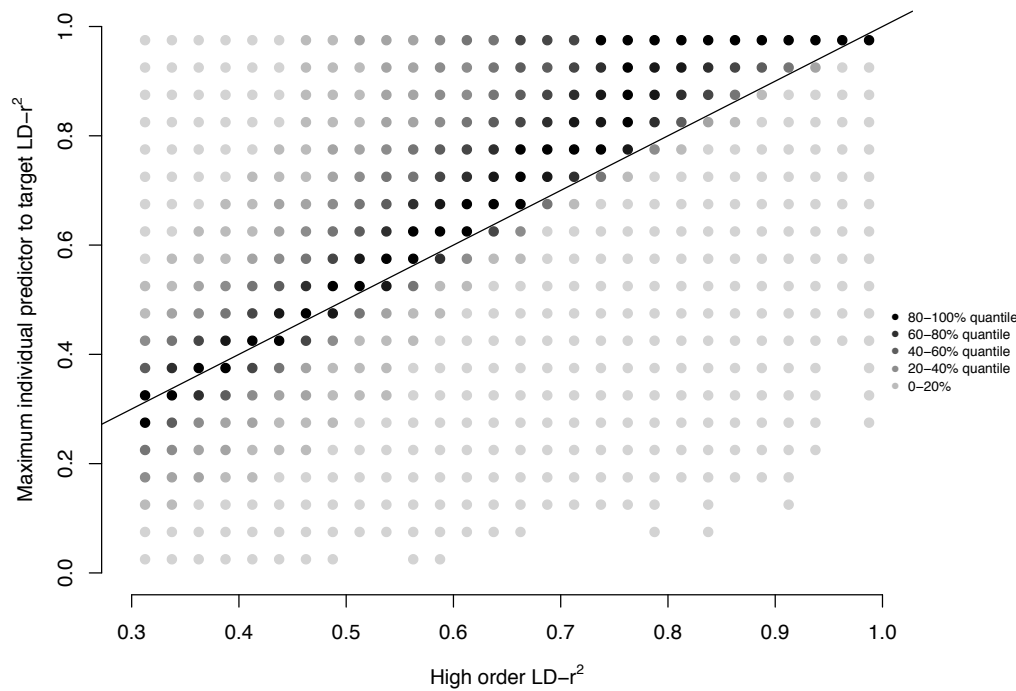327     one of the predictors is located close to the target.

11

**Figure 2.** *Illustration of how the prevalence of high order LD-$r^2$ to the targets in a 6Mb window on* A. thaliana *chromosome 2 (8 – 14Mb) depends on distance of the predictors from the target. The color gradient illustrates the proportion of predictor pairs that reach a particular LD-$r^2$ (x-axis) depending on the distance between the nearest predictor and the target (y-axis). Results are presented for populations with n = 100 (**A**) and n = 728 (**B**) individuals.*

Strong high-order LD-$r^2$ between a predictor pair and a target is mostly observed when at least one of the predictors is in strong individual LD-$r^2$ with the target. However, as illustrated in Figure 3, many cases also exist where the high order LD-$r^2$ is strong while the LD-$r^2$ to the individual predictors is weak.

12

**Figure 3.** *Strong second order LD-$r^2$ exists also when the individual predictor to target LD-$r^2$ is weak. The intensity of each dot illustrates the number of cases with a particular high order LD-$r^2$ / maximum individual predictor to target LD-$r^2$ combination. Dots below the line are cases where the high order LD-$r^2$ stronger than any individual predictor to target LD-$r^2$ (n=728).*

*Estimating the genome wide prevalence of strong high order linkage disequilibrium*

Figure 2 illustrates that high-order LD-$r^2$ exists where one or both predictors are located close to the target as well as when one or both predictors are located further away in the evaluated 6Mb window. The genome-wide prevalence of high order LD-$r^2$ for the three different classes of predictor pairs, cis-cis/cis-trans/trans-trans (as defined above) were next explored in three pairs of distant 2Mb windows in the genome (Table 1) to provide data to estimate their genome-wide prevalence. Here, only cases when individual predictors in the windows had lower individual LD-$r^2$ than 0.6 to the targets were considered.
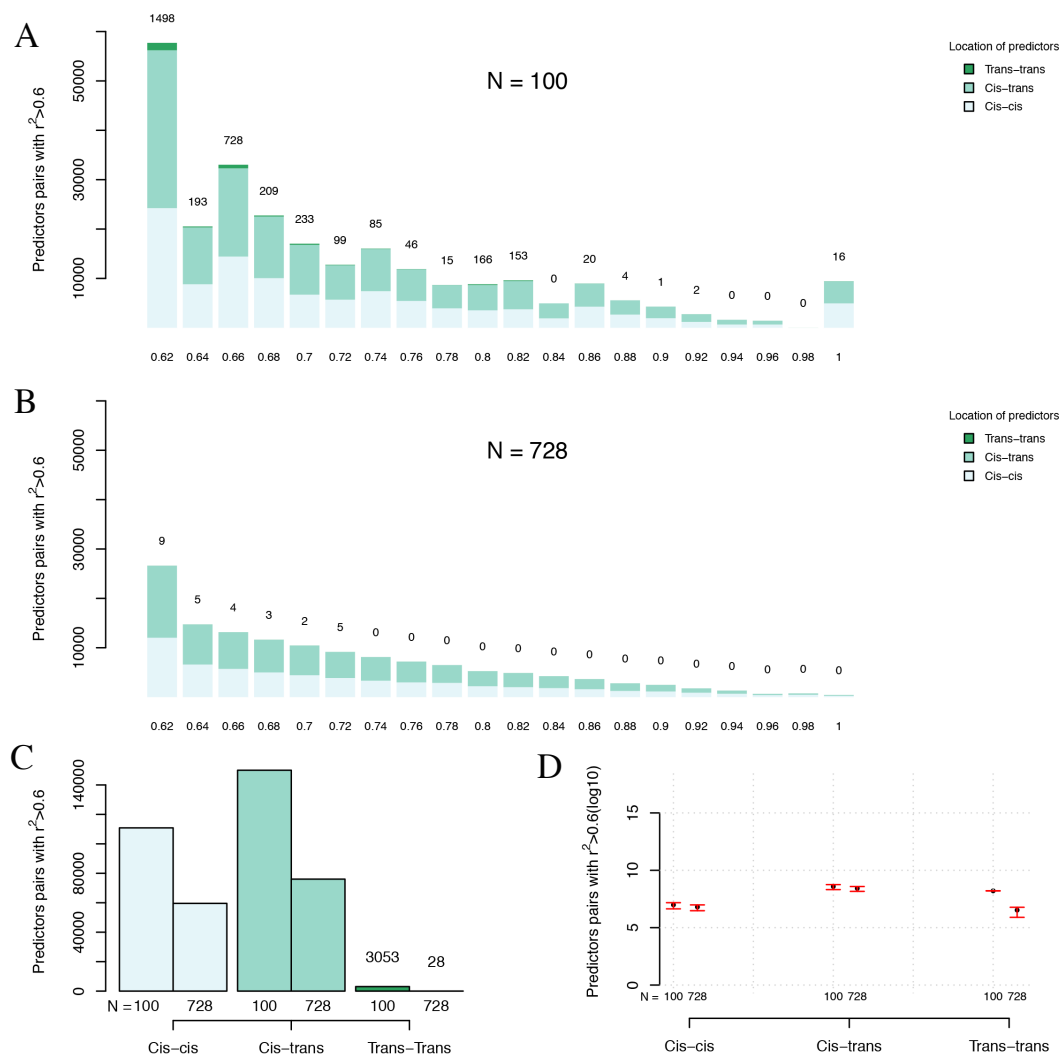
Overall, the fraction of predictor pairs that display higher second-order LD (LD-$r^2$ > 0.6) is low. In the smaller population (n = 100), less than 1 out of $10^6$ evaluated predictor pairs and in the larger population (n = 768) less than 1 out of $10^7$ (Table S1). However, since the total number of evaluated pairs was very large (around $10^{11}$),

13

359  many cases were still detected. Regardless of population size, cis-cis and cis-trans

360  pairs dominated (42/44% for n = 100, and 56/58% for n = 728; Figure 4A-C; Table

361  S1). Trans-trans pairs existed, but were much less common (~1% for n = 100, <0.01%

362  for n = 728, respectively, Figure 4A-C; Table S1). When extrapolating these results to

363  a genome wide scale, this picture, however, changes dramatically (Figure 4D). Trans-

364  trans and cis-trans predictor pairs are now much more common than cis-cis pairs due

365  to their much higher genome-wide prevalence (35/18-fold for n = 100 and 35/0.3 for n

366  = 728 more common; Figure 4D, Table S1). This result illustrates that it is a

367  considerable risk to disregard high-order LD as a possible explanation for statistical

368  epistatic interactions even at larger sample-sizes.

369



371  ***Figure 4.*** *Number of predictor pairs of different classes in strong high order LD-r² (>0.6) to targets*

372  *detected in the evaluated windows and estimated genome wide. The distribution of LD-r² values > 0.6*

14

373     *for the cis-cis, cis-trans,trans-trans predictor pairs for (**A**; n = 100) and (**B**; n = 728) The total number*

374     *of predictor pairs with high order LD-$r^2$ above 0.6 in the three classes are summarized in (**C**) and used*

375     *to estimate the total expected number of predictor pairs in the entire genome (**D**; error bars show the*

376     *estimation error estimated from the results obtained for the three window (Materials and Methods).*

377

378     *Linking high order LD and statistical epistasis in a two locus epistatic association*

379     *analysis in* A. thaliana

380     A publicly available dataset including 340 *Arabidopsis thaliana* accessions were used

381     for a genome wide association analysis for leaf molybdenum concentration This

382     dataset was earlier used by (Forsberg *et al.* 2015) to dissect a locus with a highly

383     significant variance heterogeneity association for leaf molybdenum concentration

384     (Shen *et al.* 2012) to the contributions of four independent associations in an extended

385     LD block on chromosome 2. Several of these associations were found to structural

386     variants that were poorly tagged by the SNP markers (Forsberg *et al.* 2015). Our

387     pairwise genome wide scans for pairs of epistatic loci identified 396 significant SNP

388     pairs. For 290 pairs both markers were located in the narrow region on chromosome 2

389     that was earlier dissected in detail (Forsberg *et al.* 2015). All these are examples of

390     cis-cis predictor pairs. The remaining 106 pairs contained one predictor in the

391     chromosome 2 region and another one elsewhere in the genome, being examples of

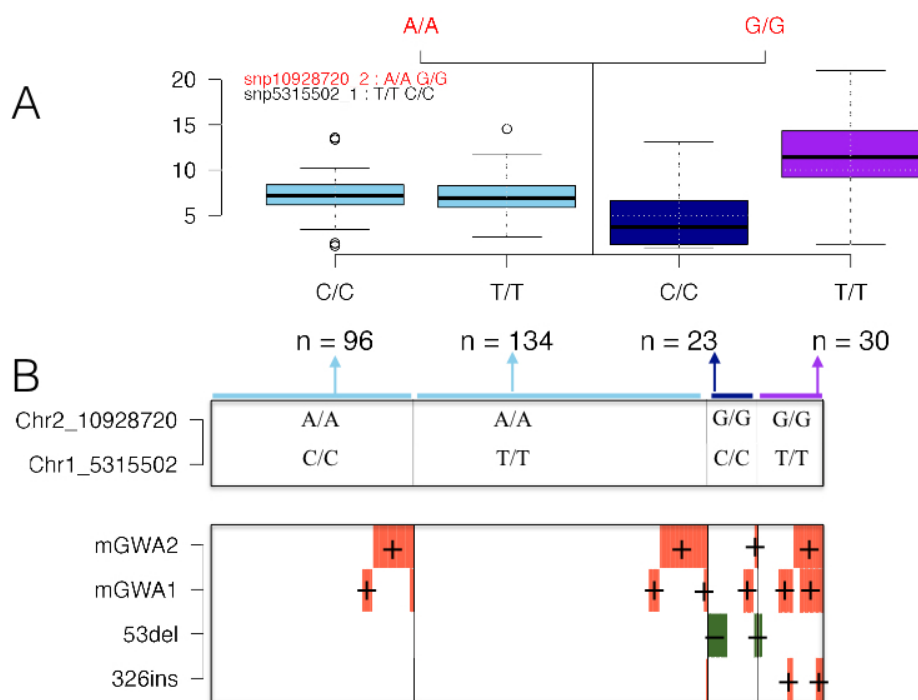392     cis-trans predictor pairs.

393

394     The strongest pairwise epistasis was detected for a cis-trans predictor pair (Figure 5A).

395     The accessions with the AA genotype at the predictor located in trans to the

396     chromosome 2 region (chromosome 1:5,315,502 bp) all have an intermediate

397     molybdenum level in the leaf (Figure 5A). The accessions with the GG allele at the

398     trans predictor have different levels of molybdenum in their leaves depending on

399     whether they carry the CC or TT genotype at the cis predictor in on chromosome 2

400     (10,928,720 bp). These differences explain the significant statistical epistasis detected

401     when fitting the two-locus epistatic model [1] to this data.

402

403     This statistical interaction could be due to a true genetic interaction. An alternative

404     explanation is however presented in Figure 5 There, the overlap between the two

405     locus genotypes for the cis-trans predictor pair (Figure 5A) and the alleles at the four

15

406    loci earlier reported to be associated with leaf molybdenum levels in this region
407    (Forsberg *et al.* 2015) are illustrated. The multi-locus genotypes of the predictor pair
408    tags different combinations of minor alleles at the four loci that were found to either
409    increase (mGWA1, mGWA2, 326ins) or decrease (53del) leaf molybdenum levels in
410    the accessions (Forsberg *et al.* 2015). The statistical epistatic interaction was detected
411    due to the difference in molybdenum levels between accessions carrying the GGCC
412    genotype (low molybdenum) and GGTT (high molybdenum). Figure 5B shows that
413    the accessions with the GGCC genotype have the lowest frequency of the
414    molybdenum increasing allele mGWA2 and the highest frequency of the molybdenum
415    decreasing allele 53del. The accessions with the GGTT genotype instead have the
416    highest frequencies of the molybdenum increasing alleles at mGWA2, mGWA1 and
417    326ins. The genotypes AACC and AATT, with intermediate molybdenum levels, both
418    have intermediate frequencies of the mGWA1 and mGWA1 increasing alleles and
419    lack the 53del and 326ins alleles. A more parsimonious interpretation of these results
420    is thus that the statistical epistasis at the predictor pair is due to the high order LD
421    between them and the genotypes at the four loci located in the region on chromosome
422    2.



423

**Figure 5.** *An illustration of how the high order LD between four polymorphisms affecting the level of molybdenum in the* A. thaliana *leaf (Forsberg et al. 2015), likely explains the significant statistical*

16

426    *epistasis detected for a cis-trans predictor pair. (**A**) Boxplots illustrating the phenotypic distribution in*

427    *the four genotype classes defined by the cis-trans predictor pair with the strongest significant epistatic*

428    *interaction to the level of molybdenum in the* A. thaliana *leaf. (**B**) Illustration of the connection*

429    *between the two-locus genotypes of the predictor pair and the minor alleles at the four linked loci*

430    *associated with this trait on chromosome 2 (Forsberg et al. 2015). The top box in (**B**) illustrates the*

431    *two-locus genotype for the predictor pair, with the width of each sub-box indicating the number of*

432    *individuals in each genotype class in the population. In the bottom box in (**B**), each individual is*

433    *represented as a column, where green (molybdenum decreasing) and orange (molybdenum increasing)*

434    *colors indicates that the individual carry the minor alleles at the four loci identified in (Forsberg et al.*

435    *2015). mGWA1 and mGWA2 are SNP markers associated with the trait and 53del and 326 are*

436    *structural polymorphisms (Forsberg et al. 2015).*

437

438                                        **DISCUSSION**

439    High order linkage disequilibrium between combinations of genotyped markers, and

440    unobserved functional polymorphisms, can result in significant statistical epistasis in

441    genome wide association analyses. This was earlier illustrated empirically for linked

442    pairs of genotyped predictor SNPs and ungenotyped target polymorphisms in humans

443    by Wood *et al.(Wood et al.* 2014). Here, we present a new example from *A. thaliana*

444    where significant statistical epistasis between pairs of predictors is due to the effects

445    at a single loci and that only one of the statistically interacting loci was located near

446    the target. By exploring the prevalence of second order LD in the genome of the

447    1001-genomes *A. thaliana* collection, we find that although the total amount of high

448    order linkage disequilibrium decreases with increasing population sizes, it is still

449    highly prevalent both within and across chromosomes even in relatively large

450    populations (n = 728). It is was found to be most common when one predictor is in

451    high LD to (and located physically near) the target, but many cases exist where the

452    LD to the individual predictors is very weak but the high order LD is strong. The

453    choice of target and predictor SNPs in this study is arbitrary and therefore it it is

454    difficult to assess how representative they are for the prevalence of high order LD in

455    other populations. However, they do suggest that strong high order LD can be

456    prevalent also in larger populations, indicating that statistical epistasis observed in

457    studies based on reduced representation genotyping (such as SNP-chips) need to be

458    interpreted with caution.

17

459

460    The most prevalent type of high order LD on a genome wide basis is that of cis-trans

461    predictor pairs, but also cis-cis pairs are common regardless of population size. The

462    prevalence of trans-trans pairs is high in smaller populations but decreases rapidly as

463    the population size increases. A possible biological explanation for the observation

464    that cis-cis and cis-trans high order LD pairs is relatively prevalent also at larger

465    population sizes would be that the number of, and variation in, the trans located

466    predictors is sufficiently large on a genome-wide basis to complement any

467    imperfection in the tagging of the functional polymorphism by the cis located

468    predictor. Whereas trans-trans high order LD will always result in falsely associated

469    loci, cis-trans and cis-cis high order LD presents an opportunity to identify true

470    functional loci for the trait. The problem in a real data analysis is that statistical

471    epistasis between a pair of predictors can emerge from true interactions or high order

472    LD within and across chromosomes. However, as the sample sizes increase the risk of

473    detecting pairs of predictors where none is located close to the true functional

474    polymorphism decreases. Before concluding that the detected association is due to

475    two interacting loci, further analyses of the associated pair are however recommended.

476

477    Whole-genome sequencing provides unprecedented opportunities to genotype most

478    segregating single nucleotide polymorphisms in the genome. Despite this, it is

479    unlikely that these will be able to tag all functional polymorphisms, such as larger

480    structural variants or multi-allelic functional loci due to tandem repeats. Hence, even

481    though the scenario of reduced representation genotyping with SNP-chips or similar

482    will soon be a technology of the past, association analyses will still be challenged by

483    the need to tag hidden polymorphisms with imperfect markers as illustrated in our

484    analyses of the complex locus affecting molybdenum levels in the *A. thaliana* leaf. In

485    fact, it is not unlikely that the problem with high order LD between SNP predictors

486    and hidden, complex functional loci will remain a major challenge in the future as the

487    increased number of markers generated by sequencing also increases the chance of

488    finding combinations of cis-cis or cis-trans predictors that tag these functional

489    polymorphisms better than any single marker. To evaluate the extent of this problem

490    one will, however, need a more comprehensive dataset than the one studied here

491    including a more complete scoring of all types of non-SNP polymorphisms in the

18

492    genome with potential effect on traits of interest.

493

494    The prevalence of high order LD is likely to be more of a concern in populations of

495    inbred or haploid individuals. These include, for example, inbred lines derived from

496    bi- and multi-parental crosses of plants and animals, as well as populations of wild

497    collected inbred plants (Churchill *et al.* 2004; Valdar *et al.* 2006; Kover *et al.* 2009;

498    Cao *et al.* 2011; Mackay *et al.* 2012). As heterozygotes are not present in these

499    populations, the number of multi locus genotype classes is smaller than in outbred

500    populations, making them attractive for studies of genetic interactions. As a common

501    approach to detect interactions in such populations is to identify pairs of loci

502    displaying significant statistical epistasis, such results need to be interpreted with

503    caution, as the analyzed populations are generally small. If one, or more, of the

504    functional polymorphisms in the genome are unknown and poorly tagged by the

505    genotyped markers, there is a risk that statistical interactions arise from high-order LD

506    between the genotyped markers and the hidden functional polymorphisms. Hence,

507    even though these populations increase the power to map loci displaying statistical

508    epistasis, there is also a risk of falsely concluding that the underlying genetic

509    architecture involves genetic interactions.

510

511                                    **CONCLUSIONS**

512    Statistical epistasis detected in genome wide association analyses can result from high

513    order LD between genotyped markers and unobserved functional polymorphisms.

514    This study revealed extensive and strong genome wide high order LD between pairs

515    of markers on a high density 250k SNP-chip and individual markers revealed by

516    whole genome sequencing in the *A. thaliana* 1001-genomes collection. The high

517    prevalence of strong high order LD in this dataset suggests that epistatic variance

518    detected between pairs of markers in association analyses, especially in small inbred

519    populations genotyped for reduced representation sets of markers, need to be

520    interpreted with caution. An empirical example is presented where a pair of markers

521    with significant statistical epistasis in a genome wide association analysis is in high

522    order LD with a complex multi allelic locus with large effects on the analyzed trait.

523    As complex functional loci such as this are unlikely to be captured by individual bi-

524    allelic SNP markers, even if millions of them are scored by whole genome sequencing,

       19

525     it is important to evaluate also other explanations of statistical epistasis than
526     underlying genetic interactions in particular when small populations of inbred
527     individuals are studied.

528

529                                   **ACKNOWLEDGEMENTS**

532

533                                 **AUTHOR CONTRIBUTIONS**

534     ÖC and SF initiated the study. ÖC, SF and YZ designed the project and the statistical
535     analyses; SF and YZ wrote analysis scripts and performed the data analyses. ÖC and
536     YZ summarized the results and wrote the initial version of the manuscript. All authors
537     contributed to the writing of the final version of the manuscript.

538

539                                 **DISCLOSURE DECLARATION**

540     The authors declare no competing interest.

541

542                                 **SUPPLEMENTARY MATERIAL**

543     Supplementary material is provided in Supplementary Figure 1 and Supplementary
544     Table 1.

545

546                                 **REFERENCES:**

547     1001 Genomes Consortium., 2016 1,135 Genomes Reveal the Global Pattern of
548        Polymorphism in Arabidopsis thaliana. Cell **166**: 481–491.

549     Anholt R. R. H., Dilda C. L., Chang S., Fanara J.-J., Kulkarni N. H., et al., 2003 The
550        genetic architecture of odor-guided behavior in Drosophila: epistasis and the
551        transcriptome. Nat Genet **35**: 180–184.

552     Atwell S., Huang Y. S., Vilhjálmsson B. J., Willems G., Horton M., et al., 2010
553        Genome-wide association study of 107 phenotypes in Arabidopsis thaliana inbred
554        lines. Nature **465**: 627–631.

555     Aulchenko Y. S., Ripke S., Isaacs A., van Duijn C. M., 2007 GenABEL: an R

20

556     package for genome-wide association analysis. Bioinformatics **23**: 1294–1296.

557  Baxter I., Brazelton J. N., Yu D., Huang Y. S., Lahner B., et al., 2010 A coastal cline
558     in sodium accumulation in Arabidopsis thaliana is driven by natural variation of
559     the sodium transporter AtHKT1;1. PLoS Genet **6**: e1001193.

560  Caicedo A., Stinchcombe J., Olsen K., Schmitt J., Purugganan M., 2004 Epistatic
561     interaction between Arabidopsis FRI and FLC flowering time genes generates a
562     latitudinal cline in a life history trait. Proceedings of the National Academy of
563     Sciences of the United States of America **101**: 15670.

564  Cao J., Schneeberger K., Ossowski S., Günther T., Bender S., et al., 2011 Whole-
565     genome sequencing of multiple Arabidopsis thaliana populations. Nat Genet **43**:
566     956–963.

567  Carlborg Ö., Jacobsson L., Åhgren P., Siegel P., Andersson L., 2006 Epistasis and the
568     release of genetic variation during long-term selection. Nat Genet **38**: 418–420.

569  Churchill G. A., Airey D. C., Allayee H., Angel J. M., Attie A. D. et al., 2004 The
570     Collaborative Cross, a community resource for the genetic analysis of complex
571     traits. Nat Genet **36**: 1133–1137.

572  Forsberg S. K. G., Carlborg Ö., 2017 On the relationship between epistasis and
573     genetic variance heterogeneity. Journal of Experimental Biology 68:5341-5438.

574  Forsberg S. K. G., Andreatta M. E., Huang X.-Y., Danku J., Salt D. E. et al., 2015
575     The Multi-allelic Genetic Architecture of a Variance-Heterogeneity Locus for
576     Molybdenum Concentration in Leaves Acts as a Source of Unexplained Additive
577     Genetic Variance. (GP Copenhaver, Ed.). PLoS Genet **11**: e1005648.

578  Forsberg S. K. G., Bloom J. S., Sadhu M. J., Kruglyak L., Carlborg Ö., 2017
579     Accounting for genetic interactions improves modeling of individual quantitative
580     trait phenotypes in yeast. Nat Genet **49**: 497–503.

581  Hao K., Di X., Cawley S., 2007 LdCompare: rapid computation of single- and
582     multiple-marker r2 and genetic coverage. Bioinformatics **23**: 252–254.

21

583  Hemani G., Shakhbazov K., Westra H.-J., Esko T., Henders A. K. et al., 2014
584      Detection and replication of epistasis influencing transcription in humans. Nature
585      **508**: 249–253.

586  Horton M. W., Hancock A. M., Huang Y. S., Toomajian C., Atwell S. et al, 2012
587      Genome-wide patterns of genetic variation in worldwide Arabidopsis thaliana
588      accessions from the RegMap panel. Nat Genet **44**: 212–216.

589  Huang W., Mackay T. F. C., 2016 The Genetic Architecture of Quantitative Traits
590      Cannot Be Inferred from Variance Component Analysis (X Zhu, Ed.). PLoS
591      Genet **12**: e1006421.

592  Kover P. X., Valdar W., Trakalo J., Scarcelli N., Ehrenreich I. M. et al, 2009 A
593      Multiparent Advanced Generation Inter-Cross to fine-map quantitative traits in
594      Arabidopsis thaliana. PLoS Genet **5**: e1000551.

595  Lachowiec J., Shen X., Queitsch C., Carlborg Ö., 2015 A Genome-Wide Association
596      Analysis Reveals Epistatic Cancellation of Additive Genetic Variance for Root
597      Length in Arabidopsis thaliana. PLoS Genet **11**: e1005541.

598  Lynch M., Walsh B., 1997 *Genetics and Analysis of Quantitative Traits*. Sinauer
599      Assoc.

600  Mackay T. F. C., Richards S., Stone E. A., Barbadilla A., Ayroles J. F. et al., 2012
601      The Drosophila melanogaster Genetic Reference Panel. Nature **482**: 173–178.

602  Nelson R. M., Pettersson M. E., Carlborg Ö., 2013 A century after Fisher: time for a
603      new paradigm in quantitative genetics. Trends Genet **29**: 669–676.

604  Phillips P. C., 2008 Epistasis--the essential role of gene interactions in the structure
605      and evolution of genetic systems. Nat Rev Genet **9**: 855–867.

606  Platt A., Vilhjálmsson B. J., Nordborg M., 2010 Conditions under which genome-
607      wide association studies will be positively misleading. Genetics **186**: 1045–1052.

608  Sackton T. B., Hartl D. L., 2016 Genotypic Context and Epistasis in Individuals and
609      Populations. Cell **166**: 279–287.

610 Segrè D., Deluna A., Church G. M., Kishony R., 2004 Modular epistasis in yeast
611     metabolism. Nat Genet: 1–7.

612 Shen X., Pettersson M., Rönnegård L., Carlborg Ö., 2012 Inheritance beyond plain
613     heritability: variance-controlling genes in Arabidopsis thaliana. PLoS Genet **8**:
614     e1002839.

615 Shimomura K., Low-Zeddies S., King D., Steeves T., Whiteley A. et al., 2001
616     Genome-wide epistatic interaction analysis reveals complex genetic determinants
617     of circadian behavior in mice. Genome Research **11**: 959.

618 Valdar W., Solberg L. C., Gauguier D., Burnett S., Klenerman P. et al., 2006
619     Genome-wide genetic association of complex traits in heterogeneous stock mice.
620     Nat Genet **38**: 879–887.

621 Wood A. R., Tuke M. A., Nalls M. A., Hernandez D. G., Bandinelli S. et al., 2014
622     Another explanation for apparent epistasis. Nature **514**: E3–5.

623