

# Grouping of genomic markers in populations with family structure

Dörte Wittenburg (✉ [wittenburg@fbn-dummerstorf.de](mailto:wittenburg@fbn-dummerstorf.de))

Leibniz Institute for Farm Animal Biology <https://orcid.org/0000-0002-3639-2574>

Michael Doschoris

Leibniz-Institut für Nutztierbiologie

Jan Klosa

Leibniz-Institut für Nutztierbiologie

---

## Research article

**Keywords:** Single nucleotide polymorphism, Covariance matrix, Clustering, TagSNP, Group lasso, SNP-BLUP

**Posted Date:** January 6th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-54566/v3>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on February 19th, 2021. See the published version at <https://doi.org/10.1186/s12859-021-04010-0>.

## RESEARCH

# Grouping of genomic markers in populations with family structure

Dörte Wittenburg\*, Michael Doschoris and Jan Klosa

\*Correspondence:

wittenburg@fhn-dummerstorf.de

Leibniz Institute for Farm Animal  
Biology, Institute of Genetics and  
Biometry, 18196 Dummerstorf,  
Germany

Full list of author information is  
available at the end of the article

## Abstract

**Background:** Linkage and linkage disequilibrium (LD) between genome regions cause dependencies among genomic markers. Due to family stratification in populations with non-random mating in livestock or crop, the standard measures of population LD such as  $r^2$  may be biased. Grouping of markers according to their interdependence needs to account for the actual population structure in order to allow proper inference in genome-based evaluations.

**Results:** Given a matrix reflecting the strength of association between markers, groups are built successively using a greedy algorithm; largest groups are built first. As an option, a representative marker is selected for each group. We provide an implementation of the grouping approach as a new function to the R package `hscovar`. This package enables the calculation of the theoretical covariance between biallelic markers for half- or full-sib families and the calculation of representative markers. In case studies, we have shown that the number of groups comprising dependent markers was smaller and representative SNPs were spread more uniformly over the investigated chromosome region when the family stratification was respected compared to a population-LD approach. In a simulation study, we observed that sensitivity and specificity of a genome-based association study improved if selection of representative markers took family structure into account.

**Conclusions:** Chromosome segments which frequently recombine in the underlying population can be identified from the matrix of pairwise dependence between markers. Representative markers can be exploited, for instance, for dimension reduction prior to a genome-based association study or the grouping structure itself can be employed in a grouped penalization approach.

**Keywords:** Single nucleotide polymorphism; Covariance matrix; Clustering; TagSNP; Group lasso; SNP-BLUP

## <sup>1</sup>**Background**

<sup>2</sup>Genomic markers are an invaluable source for characterizing genetic variety and to  
<sup>3</sup>elucidate the relationship between genetic and phenotypic variation in breeding pop-  
<sup>4</sup>ulations. Dependencies among genomic markers are caused by linkage and linkage  
<sup>5</sup>disequilibrium (LD) between genome regions. Though this condition complicates  
<sup>6</sup>investigations on which genetic variants are truly associated with trait expression  
<sup>7</sup>[1], dependencies can be advantageous for grouping of markers. For example, clus-  
<sup>8</sup>tering based on a greedy algorithm [2], hierarchical clustering (e.g., [3]) or grouping  
<sup>9</sup>via interval-graph modeling [4] exploit the presence of LD blocks which are regions  
<sup>10</sup>of particularly high correlation. To allow for proper inferences of such approaches, a  
<sup>11</sup>suitable measure for the strength of dependence is needed. For instance, measuring  
<sup>12</sup>LD in terms of  $r^2$  [5] is a natural choice but it is meaningful only for popula-  
<sup>13</sup>tions without stratification. In livestock and crop breeding, however, populations  
<sup>14</sup>are often characterized by strong family stratification due to non-random mating  
<sup>15</sup>of selected individuals. As examples, large paternal half-sib families are typical for  
<sup>16</sup>cattle populations whereas chicken or fish populations consist of full-sib families.  
<sup>17</sup>In plant breeding, maternal half-sib families are often produced in, for instance,  
<sup>18</sup>wheat and clover. Then, linkage between markers within family leads to haplotype  
<sup>19</sup>frequencies among progeny that are not conclusive for estimating  $r^2$ . Hence, there  
<sup>20</sup>is need to promote measures of marker dependence which takes into account the  
<sup>21</sup>particular family structure.

<sup>22</sup>Especially in situations of ultra-dense panels of single nucleotide polymorphisms  
<sup>23</sup>(SNPs), it is often sufficient to investigate representative SNPs ( “tagSNPs” ) out  
<sup>24</sup>of each cluster. This subset can help identifying trait-associated genome regions  
<sup>25</sup>in genome-wide association studies and allows comparing genome characteristics  
<sup>26</sup>between ethnics/species/breeds (e.g., [2]). As the choice of tagSNPs is a consequence  
<sup>27</sup>of grouping, it is also influenced by the underlying population structure.

<sup>28</sup>The objective of this paper is to exploit the family structure of a population for  
<sup>29</sup>specifying groups of associated markers. We generalize the grouping approach of  
<sup>30</sup>Carlson *et al.* [2] in order to allow binning of markers given a correlation matrix or  
<sup>31</sup>any kind of similarity matrix with scaled entries in [0, 1]. We investigate three case  
<sup>32</sup>studies and a simulation study. For each case study, we visually inspect the correla-  
<sup>33</sup>tion matrix and link to the outcome of grouping. Usability for genome-based asso-

ciation studies is shown as one possible field of application. Results were compared<sup>1</sup>  
 to the commonly used population-LD approach which ignores family structure. We<sup>2</sup>  
 provide a new function to the R package `hscovar` (available at CRAN) that enables<sup>3</sup>  
 grouping of markers and selection of representative markers.<sup>4</sup>

## 6 **Methods**

The dependence between pairs of SNPs, each with two alleles A and B, can be<sup>7</sup>  
 expressed in terms of a covariance or correlation matrix. It has already been shown<sup>8</sup>  
 in the literature how to calculate the theoretical covariance between markers in<sup>9</sup>  
 a population consisting of half-sib families [1]. It requires a genetic map, haplo-<sup>10</sup>  
 types of the common parent and LD information (or haplotype frequencies) of the<sup>11</sup>  
 population the individual parent comes from. This approach can be extended to be<sup>12</sup>  
 applicable to full-sib families by adding the paternal and maternal contribution into<sup>13</sup>  
 a single covariance matrix; the derivation is summarized in Additional file 1. Hence,<sup>14</sup>  
 a covariance matrix can be derived for any family structure, and this constitutes<sup>15</sup>  
 the input of the following grouping approach.<sup>16</sup>

### 18 **Grouping of markers**

We generalized the strategy of Carlson *et al.* [2] for binning markers and selecting<sup>19</sup>  
 representatives to be applicable to any symmetric matrix which reflects a measure<sup>20</sup>  
 of dependence between markers and has entries scaled in [0, 1]. In particular, we<sup>21</sup>  
 considered the correlation matrix  $R$ . The idea is that SNPs which are associated to<sup>22</sup>  
 each other are assigned to groups. Groups are built one after the other - the largest<sup>23</sup>  
 group at first. For  $b = 1, 2, \dots$ , the  $b$ -th group is identified by searching for the SNP<sup>24</sup>  
 that has most occurrences of absolute correlation to other SNPs larger than a given<sup>25</sup>  
 threshold  $t$ . More precisely, let  $S_b$  denote the set of SNP indices which have not<sup>26</sup>  
 been binned yet. Then, for each SNP  $k \in S_b$ , the set of highly associated SNPs is<sup>27</sup>  
 determined as (Step A)<sup>28</sup>

$$29 \quad C_k = \{l \mid l \in S_b : |R_{k,l}| > t\}, \quad 29$$

and a set with highest cardinality (operator #) is chosen,<sup>31</sup>

$$32 \quad c \in \arg \max_k \#C_k. \quad 32$$

$$33 \quad c \in \arg \max_k \#C_k. \quad 33$$

<sup>1</sup>Thus,  $C_c$  constitutes the  $b$ -th group, and a tagSNP is selected from this group. SNP<sup>1</sup>  
<sup>2</sup> $c$  has strong correlation with any other SNP in  $C_c$  but it can happen that also other<sup>2</sup>  
<sup>3</sup>SNPs of  $C_c$  fulfill this criterion. Hence, a set of candidates is given by (Step B) <sup>3</sup>

$$T = \{k \mid k \in C_c \wedge \forall l \in C_c : |R_{k,l}| > t\}.$$

<sup>4</sup>  
<sup>5</sup>  
<sup>6</sup>  
<sup>7</sup>If more than one candidate remains, then the  $|T|/2$ -th SNP becomes the repre-<sup>7</sup>  
<sup>8</sup>sentative of group  $b$ . A next round of iteration is started using  $S_{b+1} = S_b \setminus C_c$  until<sup>8,9</sup>  $S_{b+1}$   
<sup>9</sup> $= \emptyset$ . The number of groups only depends on the threshold  $t$ . Similar to [2],<sup>9,10</sup>  $t =$   
<sup>10</sup> $0.8$  is a suitable value. In an extreme case (with  $t$  approaching 1), each SNP<sup>10,11</sup> builds a  
<sup>11</sup>single group, yielding a complexity of this algorithm of  $O(p^2)$ , with  $p$  the<sup>11,12</sup> total  
<sup>12</sup>number of SNPs. This approach is implemented as function tagSNP in the<sup>12</sup> <sup>13</sup>R  
<sup>13</sup>package hscovar. A graphical representation of this algorithm is shown in Figure<sup>13</sup>

<sup>14</sup><sup>1</sup>. <sup>14</sup>

## Evaluation

<sup>15</sup>  
<sup>16</sup>It is an obvious choice to compare the family approach with a population-LD ap-<sup>17</sup>  
<sup>17</sup>proach. Such an approach requires the population frequency of the different hap-<sup>18</sup>  
<sup>18</sup>lotypes ( $f_{A-A}, f_{A-B}, f_{B-A}, f_{B-B}$ ) from which LD between markers is computed  
<sup>19</sup>in terms of  $r^2$  according to [5]. For any marker pair  $k, l$  with allele frequencies  
<sup>20</sup> $f_k = f_{A-A} + f_{A-B}$  and  $f_l = f_{A-A} + f_{B-A}$ , we have  
<sup>21</sup>

$$r_{k,l}^2 = \frac{(f_{A-A}f_{B-B} - f_{A-B}f_{B-A})^2}{f_k(1-f_k)f_l(1-f_l)}.$$

<sup>22</sup>  
<sup>23</sup>  
<sup>24</sup>The LD matrix can be computed from progeny genotypes using the function ld  
<sup>25</sup>from the R package snpStats version 1.38.0 [6]. This function undertakes phasing  
<sup>26</sup>of genotypes using a maximum-likelihood approach [7]. Based on the LD matrix  
<sup>27</sup>containing  $r^2$ , representative markers can be derived as described above.  
<sup>28</sup>  
<sup>29</sup>Family and population-LD approach were compared using the Calinski-Harabasz  
<sup>30</sup>(CH) index [8] which measures the cluster quality with respect to inter- and intra-  
<sup>31</sup>cluster distances. The method with higher CH index performed better. For this, the  
<sup>32</sup>function calinhara from the R package fpc version 2.2-9 was applied to the groups  
<sup>33</sup>obtained; the quality referred to distances based on the genotype matrix centered  
within family and scaled (as described below). Furthermore, we present the pure

<sup>1</sup>number of groups and highlight those groups with group size of at least three. This<sup>1</sup>  
<sup>2</sup>also helped visualizing the location of corresponding representative markers.<sup>2</sup>

<sup>3</sup> Selecting a representative set of markers is a natural tool for dimension reduc-<sup>3</sup>  
<sup>4</sup>tion prior to genomic evaluations in order to reduce the impact of multicollinearity<sup>4</sup>  
<sup>5</sup>among predictor variables. Representative SNPs capture cumulative effects of the<sup>5</sup>  
<sup>6</sup>corresponding LD blocks on trait expression. We investigated a SNP-BLUP ap-<sup>6</sup>  
<sup>7</sup>proach, which is widely used in genomic evaluations (e.g., [9]), and thereby demon-<sup>7</sup>  
<sup>8</sup>strate one possible application of the suggested approach. Representative markers<sup>8</sup>  
<sup>9</sup>selected from the family or from the population-LD approach were employed as<sup>9</sup>  
<sup>10</sup>predictor variables in a regression model<sup>10</sup>

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

<sup>11</sup>with  $\mathbf{y} = (y_1, \dots, y_n)^T$  the phenotype vector,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_\tau)^T$  the vector of genetic<sup>11</sup>  
<sup>12</sup>effects captured by  $\tau$  tagSNPs, the corresponding design matrix  $\mathbf{X}$  with dimensions<sup>12</sup>  
<sup>13</sup> $16n \times \tau$  including the genotype codes in terms of major allele counts. The columns of<sup>13</sup>  
<sup>14</sup> $\mathbf{X}$  and the vector  $\mathbf{y}$  were centered within family and scaled to obtain an empirical<sup>14</sup>  
<sup>15</sup>variance of one. The residual errors were assumed to be independently and normally<sup>15</sup>  
<sup>16</sup>distributed. For convenience, no other effects were assumed. We used the R package<sup>16</sup>  
<sup>17</sup>asreml version 3.0 [10] to estimate the vector of regression coefficients as<sup>17</sup>  
<sup>18</sup>

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}^T\mathbf{y}$$

<sup>19</sup>where the shrinkage parameter  $\lambda$  was estimated via AI-REML. Significance of the<sup>19</sup>  
<sup>20</sup> $k$ -th SNP effect was tested by a  $t$ -like test statistic as in [1],<sup>20</sup>

$$T_k = \frac{\hat{\beta}_k}{SD(\hat{\beta}_k)}.$$

<sup>21</sup>Significance was reported if  $T_k \geq q_{1-\alpha/2}$  or  $T_k < q_{\alpha/2}$  using the  $1 - \alpha/2$  and<sup>21</sup>  
<sup>22</sup> $\alpha/2$  quantile of the standard normal distribution. The SNP-BLUP approach was<sup>22</sup>  
<sup>23</sup>evaluated in terms of sensitivity (i.e., true-positive rate) and specificity (i.e.,  $1 -$ <sup>23</sup>  
<sup>24</sup>false-positive rate) over a range of type-I error  $\alpha$ . We additionally verified the im-<sup>24</sup>  
<sup>25</sup>act of threshold  $t \in \{0.5, 0.6, 0.7, 0.8\}$  on grouping and its consequences on the<sup>25</sup>  
<sup>26</sup>performance of the SNP-BLUP approach.<sup>26</sup>

## <sup>1</sup>**Data**

<sup>2</sup>The data sets used for studying dependencies between SNP markers differed in SNP<sup>2</sup>  
<sup>3</sup>density and family structure. They covered a range of mean inter-marker distances<sup>3</sup>  
<sup>4</sup>from 0.003 cM to 0.23 cM. The case study of mouse data was based on low-density<sup>4</sup>  
<sup>5</sup>genotypes of full-sib families; progeny and parents were genotyped. The case study of<sup>5</sup>  
<sup>6</sup>cattle data comprised medium-density genotypes of half-sib families with genotyped<sup>6</sup>  
<sup>7</sup>progeny only. Furthermore, medium-density SNP data were available for full-sib<sup>7</sup>  
<sup>8</sup>families in maize. High-density genotype data of half-sib families were generated<sup>8</sup>  
<sup>9</sup>in simulations. For evaluation, the SNP-BLUP approach was applied to simulated<sup>9</sup>  
<sup>10</sup>data only. Unless otherwise stated, computations were done using R version 4.0.3<sup>10</sup>  
<sup>11</sup>[11] and  $t = 0.8$ ; all scripts are included as Additional files. 11

12

### <sup>13</sup>**Mouse data**

<sup>14</sup>Genotype data of a heterogeneous stock of mice were available from [https://wp.14  
15cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/](https://wp.14<br/>15cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/). We investigated chro-15  
<sup>16</sup>mosome 17 because it harbors the highly recombining MHC region which affects sev-16  
<sup>17</sup>eral immunological traits [12]. The chromosome data consisted of genotypes at 394<sup>17</sup>  
<sup>18</sup>SNPs of 2002 individuals. After filtering for individual call rate  $\geq 90\%$ , 1998 geno-18  
<sup>19</sup>typed individuals remained comprising 1 759 progeny, 120 fathers and 119 mothers.<sup>19</sup>  
<sup>20</sup>In total, 138 full-sib families (family size ranged from 1 to 47) could be identified.<sup>20</sup>  
<sup>21</sup>The SNP call rate was  $\geq 90\%$ . All genotype data were phased with Beagle version<sup>21</sup>  
<sup>22</sup>5.1 [13] and parental haplotypes were selected to set up the correlation matrix. As-<sup>22</sup>  
<sup>23</sup>suming a 1:1 relationship between physical (Build37 genome assembly) and genetic<sup>23</sup>  
<sup>24</sup>distance of adjacent markers, the genetic length was 91 cM. SNP alleles were coded<sup>24</sup>  
<sup>25</sup>in terms of the major allele in the given sample. The population-LD matrix was<sup>25</sup>  
<sup>26</sup>calculated from progeny genotypes. 26

27

### <sup>27</sup>**Cattle data**

<sup>28</sup>Genotype data of Holstein cattle were available from RADAR [https://dx.doi.28  
29org/10.22000/280](https://dx.doi.28<br/>29org/10.22000/280). The data comprised 50K SNP-chip data of five half-sib fam-  
<sup>30</sup>ilies with  $n = 265$  progeny in total; the family size ranged from 32 to 106. A 30  
<sup>31</sup>chromosome window containing 300 SNPs was selected from BTA1. Based on the 31  
<sup>32</sup>physical ordering of markers according to the genome assembly ARS-UCD1.2, this 32  
<sup>33</sup>region corresponded to 20.59 – 39.44 Mbp. The haplotypes of sires were imputed from 33

<sup>1</sup>progeny genotypes using the R package *hsphase* version 2.0.2 [14]. Maternal LD<sup>1</sup>  
<sup>2</sup>and paternal recombination rates between SNP pairs were estimated according to<sup>2</sup>  
<sup>3</sup>Hampel *et al.* [15]. However, we used a 1:1 relationship between physical (Mbp)<sup>3</sup>  
<sup>4</sup>and genetic positions (cM) for convenience; the genetic length of this window was<sup>4</sup>  
<sup>5</sup>19 cM. Sire haplotypes and maternal LD were also part of the RADAR data set.<sup>5</sup>  
<sup>6</sup>SNP alleles were coded in terms of the major maternal allele among progeny.<sup>6</sup>

### <sup>8</sup>Maize data<sup>8</sup>

<sup>9</sup>Raw marker data were available from NCBI GEO database under Accession Num-<sup>9</sup>  
<sup>10</sup>ber GSE50558, accompanied with physical coordinates corresponding to the genome<sup>10</sup>  
<sup>11</sup>assembly B73. The data set contained two maize panels, Flint and Dent, for which<sup>11</sup>  
<sup>12</sup>about 50K SNPs have been assessed in order to estimate recombination activity in<sup>12</sup>  
<sup>13</sup>different maize populations [16]. We arbitrarily chose the Flint panel and chromo-<sup>13</sup>  
<sup>14</sup>somes 2 for further analysis. In this panel, 13 full-sib families have been obtained by<sup>14</sup>  
<sup>15</sup>crossing an inbred “central” line and several inbred “founder” lines. Double haploid,<sup>15</sup>  
<sup>16</sup>(DH) lines have been derived from the F1 plants. This procedure allowed for study-<sup>16</sup>  
<sup>17</sup>ing maternal meioses only. A cM:Mbp ratio of 0.80 was reported for Flint [16]; the<sup>17</sup>  
<sup>18</sup>genetic length of chromosome 2 was approximately 188 cM. SNP genotypes of DH<sup>18</sup>  
<sup>19</sup>progeny being heterozygous were set to missing value. After filtering the data for<sup>19</sup>  
<sup>20</sup>SNP and individual call rate  $\geq 90\%$ ,  $n = 1\,248$  out of 1 262 DH progeny and 1 447<sup>20</sup>  
<sup>21</sup>out of 2 030 SNPs remained. Rarely missing marker information of DH progeny were<sup>21</sup>  
<sup>22</sup>imputed by sampling the homozygous genotypes according to their frequencies. Af-<sup>22</sup>  
<sup>23</sup>terwards loci with minor allele frequency less than 5% were discarded, yielding 956<sup>23</sup>  
<sup>24</sup>SNPs. As haplotypes of DH progeny were given with certainty, the population-LD<sup>24</sup>  
<sup>25</sup>matrix was set up directly using the squared Spearman correlation between SNPs<sup>25</sup>  
<sup>26</sup>based on haploid data. The family approach solely considered the female part of<sup>26</sup>  
<sup>27</sup>the covariance between SNPs. The haplotypes of F1 individuals were inferred from<sup>27</sup>  
<sup>28</sup>the marker data of inbred lines. SNP alleles were coded according to central-line<sup>28</sup>  
<sup>29</sup>origin.<sup>29</sup>

### <sup>30</sup>Simulated data<sup>30</sup>

<sup>31</sup>The simulation study resembled the population structure of a dairy cattle popula-<sup>31</sup>  
<sup>32</sup>tion. The setup of simulation design is fully described in [1]. Briefly, we considered<sup>32</sup>  
<sup>33</sup> $N = 1, 5, 10$  sires of half siblings. The overall number of progeny was  $n = 1\,000$ <sup>33</sup>



<sup>1</sup>equally partitioned into half-sib families. Quantitative traits were simulated which<sup>1</sup>  
<sup>2</sup>were influenced by 2 and 5 QTLs with equal effect sizes. QTLs contributed 30 %<sup>2</sup>  
<sup>3</sup>to the trait variation (i.e., heritability 0.3). In total, 300 SNPs were simulated on<sup>3</sup>  
<sup>4</sup>a chunk of DNA with 1 cM length. The data were generated using the R package<sup>4</sup>  
<sup>5</sup>AlphaSimR version 0.13.0 [17]. SNP alleles were recoded in terms of the major allele<sup>5</sup>  
<sup>6</sup>in the founder population. The simulation was repeated 100 times. For assessing<sup>6</sup>  
<sup>7</sup>the SNP-BLUP approach, a window of 0.05 cM to both sides of a simulated QTL<sup>7</sup>  
<sup>8</sup>was accepted as true-positive result. 8

## <sup>10</sup>**Results and discussion** 10

### <sup>11</sup>**Case studies** 11

<sup>12</sup>We have shown applicability of the suggested software tool to empirical data. Espe-  
<sup>13</sup>cially for the mouse data consisting of 138 genotyped full-sib families, the correlation  
<sup>14</sup>matrix gave a clear representation of genomic regions with high or low interdepen-  
<sup>15</sup>dence. In contrast to a population-LD approach, which did not account for family  
<sup>16</sup>stratification, it was also possible to identify regions with positive or negative re-  
<sup>17</sup>lationship. The population-LD approach exposed a wide region (13.82 – 21.13 Mbp)  
<sup>18</sup>that had a strong association with the entire chromosome 17 shown as a red band  
<sup>19</sup>in Figure 2b. With the family approach, this region revealed only high positive in-  
<sup>20</sup>terdependence with almost no impact on the remaining chromosome, see Figure 2a.  
<sup>21</sup>Also Carlson *et al.* [2] reported that population stratification may generate artifac-  
<sup>22</sup>tual LD and hence makes an LD-selection algorithm sensitive. Moreover, a highly  
<sup>23</sup>fragmented region appeared in the range of 27.59 – 45.88 Mbp that overlaps MHC re-  
<sup>24</sup>gions 1 and 2 ([https://www.ncbi.nlm.nih.gov/assembly/GCF\\_000001635.18/](https://www.ncbi.nlm.nih.gov/assembly/GCF_000001635.18/)).  
<sup>25</sup>This was caused by high variation of parental haplotypes. Though the total num-  
<sup>26</sup>ber of groups was smaller with the family approach than with the population-LD  
<sup>27</sup>approach (83 vs. 98), the number of large groups, i.e., with group size  $\geq 3$ , was  
<sup>28</sup>almost equal (49 vs. 51), see Table 1. The CH index was about two times higher  
<sup>29</sup>with the family approach. Figure 2 shows that tagSNPs of groups containing at  
<sup>30</sup>least three markers were distributed more uniformly over the chromosome with the  
<sup>31</sup>family approach than with the population-LD approach. The median of distances  
<sup>32</sup>between all tagSNPs was 0.70 cM and 0.52 cM in the family and population-LD  
<sup>33</sup>approach, respectively. 33

<sup>1</sup> For the cattle data consisting of five half-sib families, 239 SNPs were taken into<sup>1</sup>  
<sup>2</sup>account in the family approach. At 61 out of 300 SNPs, all sires were homozygous<sup>2</sup>  
<sup>3</sup>for the major allele, leading to zero entries on the diagonal of the covariance matrix.<sup>3</sup>  
<sup>4</sup>Thus, these loci were discarded when setting up the correlation matrix  $R$ . A clear<sup>4</sup>  
<sup>5</sup>distinction of regions with particularly high interdependence was not possible for<sup>5</sup>  
<sup>6</sup>any of the approaches (Fig. 3). In total, 11 groups with size  $\geq 3$  were found with<sup>6</sup>  
<sup>7</sup>both the family and population-LD approach (Tab. 1) but the representative SNPs<sup>7</sup>  
<sup>8</sup>of these groups were distributed more evenly over the chromosome window based<sup>8</sup>  
<sup>9</sup>on the family approach. The median of distances between all tagSNPs was 0.08 cM<sup>9</sup>  
<sup>10</sup>and 0.07 cM in the family and population-LD approach, respectively. When only<sup>10</sup>  
<sup>11</sup>those SNPs were used in the population-LD approach that were considered in the<sup>11</sup>  
<sup>12</sup>family approach, the outcome of grouping differed: 239 SNPs were binned into 194<sup>12</sup>  
<sup>13</sup>groups (CH index 6.2); 8 out of them had group size of at least three SNPs vs. 300<sup>13</sup>  
<sup>14</sup>SNPs were binned into 210 groups (CH index 2.4); 11 groups with group size of at<sup>14</sup>  
<sup>15</sup>least three SNPs. 15

<sup>16</sup> The correlation matrix corresponding to 13 full-sib families in maize is shown 16  
<sup>17</sup> in Figure 4. In three out of 956 SNPs, maternal SNP alleles of F1 plants were 17  
<sup>18</sup>missing; these loci were discarded in the family approach. In total, 953 SNPs have 18  
<sup>19</sup>been binned into 426 groups based on the correlation matrix but the CH index 19  
<sup>20</sup>was about 50 % lower than with the population-LD approach where 956 SNPs were 20  
<sup>21</sup>grouped into 576 bins (Tab. 1). With both approaches, tagSNPs corresponding to 21  
<sup>22</sup>the bins of at least three SNPs were similarly distributed over the chromosome 22  
<sup>23</sup>(family approach: 93 bins, population-LD approach: 70 bins) except a gap between 23  
<sup>24</sup>SNP index 650 and 750 which was better covered with tagSNPs from the family 24  
<sup>25</sup>approach. The median distance between all representative SNPs was 0.29 cM with 25  
<sup>26</sup>the family and 0.16 cM with the population-LD approach. With both methods, a 26  
<sup>27</sup>block of strong (positive) association among SNPs appeared in the region of 85.04 27  
<sup>28</sup>to 95.31 Mbp which is in the vicinity of the functional centromere [18]. Two F1 28  
<sup>29</sup>plants were the driving factor: they were completely heterozygous in this window. 29  
<sup>30</sup>30

### <sup>31</sup>Simulation study 31

<sup>32</sup> The average number of groups over all repetitions, and groups with at least three 32  
<sup>33</sup>SNPs are listed in Table 1; results are given for  $t = 0.8$  and varying number of 33

<sup>1</sup>half-sib families and QTLs. Grouping of markers appeared rather robust based on<sup>1</sup>  
<sup>2</sup>the family approach, about 60 groups have been built, but the number of groups<sup>2</sup>  
<sup>3</sup>strongly varied with the population-LD approach suggesting a dependence on family<sup>3</sup>  
<sup>4</sup>size. Large families needed more groups. Note that the number of rows in  $\mathbf{X}$  was<sup>4</sup>  
<sup>5</sup>fixed ( $n = 1\ 000$ ). The number of groups containing at least three markers was<sup>5</sup>  
<sup>6</sup>rather constant between methods and for different family sizes and QTLs. The<sup>6</sup>  
<sup>7</sup>CH index was at least two times larger with the family approach than with the<sup>7</sup>  
<sup>8</sup>population-LD approach. The population-LD approach becomes competitive with<sup>8</sup>  
<sup>9</sup>decreasing threshold  $t$  and performed better than the family approach with respect<sup>9</sup>  
<sup>10</sup>to the CH index when  $t \leq 0.6$  (see Additional file 7).<sup>10</sup>

<sup>11</sup> A SNP-BLUP approach was applied to simulated data in order to evaluate sen-<sup>11</sup>  
<sup>12</sup>sitivity and specificity if only tagSNPs were used as predictor variables in linear<sup>12</sup>  
<sup>13</sup>regression. As an example, results based on one half-sib family and two simulated<sup>13</sup>  
<sup>14</sup>QTLs are shown as ROC curve in Figure 5; the number of groups was almost equal<sup>14</sup>  
<sup>15</sup>for this scenario. As expected, considering all SNPs simultaneously yielded highest<sup>15</sup>  
<sup>16</sup>accuracy in terms of true-positive and false-positive rate. However, if the aim was to<sup>16</sup>  
<sup>17</sup>use filtered data in SNP-BLUP, then tagSNPs obtained from the family approach<sup>17</sup>  
<sup>18</sup>was the second best choice. Or in other words, using only one fifth of available<sup>18</sup>  
<sup>19</sup>genotypic information led to almost the same accuracy of genome-based association<sup>19</sup>  
<sup>20</sup>studies as using all genotypic information. The choice of  $t$  had no influence on which<sup>20</sup>  
<sup>21</sup>method performed best, see Figure 5a for  $t = 0.8$  and Figure 5b for  $t = 0.5$ . ROC<sup>21</sup>  
<sup>22</sup>curves looked very similar for all investigated scenarios of simulation though the<sup>22</sup>  
<sup>23</sup>number of groups obtained from the population-LD approach increased with de-<sup>23</sup>  
<sup>24</sup>creasing number of families. Hence, a direct relationship between number of groups<sup>24</sup>  
<sup>25</sup>and sensitivity/specificity seems not to exist.<sup>25</sup>

### <sup>26</sup>Options for statistical-genetics approaches<sup>26</sup>

<sup>27</sup> Instead of selecting representative markers for genome-based evaluations, employ-<sup>27</sup>  
<sup>28</sup>ing the grouping structure itself can be a beneficial option. For instance, the group<sup>28</sup>  
<sup>29</sup>assignment derived from the family approach can directly be considered in a group<sup>29</sup>  
<sup>30</sup>lasso approach [3, 19]. Then the effects of markers in a group of highly dependent<sup>30</sup>  
<sup>31</sup>markers will jointly be shrunk towards zero or enlarged with respect to the rele-<sup>31</sup>  
<sup>32</sup>vance of this group for trait expression. Additional sparsity within group can be<sup>32</sup>  
<sup>33</sup>

<sup>1</sup>achieved with a sparse-group lasso approach [20]. Grouped approaches shall be in-<sup>1</sup>  
<sup>2</sup>vestigated in more detail in future because they hold potential to cope with high<sup>2</sup>  
<sup>3</sup>multicollinearity. Possible benefits will likely depend on characteristics of the sam-<sup>3</sup>  
<sup>4</sup>ple, such as the number of families, SNP density, and population-genetic parameters,<sup>4</sup>  
<sup>5</sup>e.g., heritability and heterozygosity. 5

<sup>6</sup> In future research, the functionality of our package should be extended by grouping<sup>6</sup>  
<sup>7</sup>methods based on LD blocks which can optionally put restrictions to the physical<sup>7</sup>  
<sup>8</sup>distance between SNPs (similar to [21]). Other options for selecting tagSNPs (e.g.,<sup>8</sup>  
<sup>9</sup>depending on allele frequency; [22]) will be verified. 9

10 10

## 11 **Conclusions** 11

12 The extent of dependence among genomic markers is affected by the underlying<sup>12</sup>  
<sup>13</sup>population structure. Representative markers can be selected more efficiently if the<sup>13</sup>  
<sup>14</sup>corresponding matrix of pairwise dependencies takes this structure into account.<sup>14</sup>  
<sup>15</sup>The correlation matrix for half- or full-sib families highlights regions of high depen-<sup>15</sup>  
<sup>16</sup>dence between markers more precisely than the population-LD matrix. Additionally,<sup>16</sup>  
<sup>17</sup>it reveals regions of positive or negative association among markers. We contributed<sup>17</sup>  
<sup>18</sup>a new function tagSNP to the R package hscovar which is suited to samples from<sup>18</sup>  
<sup>19</sup>livestock and crop populations with typical family stratification. The covariance ma-<sup>19</sup>  
<sup>20</sup>trix can be set up in a piecewise manner, either separately for each chromosome or<sup>20</sup>  
<sup>21</sup>based on other meaningful information. The resulting grouping structure can be ex-<sup>21</sup>  
<sup>22</sup>ploited in genome-based evaluations to handle the problem of high multicollinearity<sup>22</sup>  
<sup>23</sup>between markers. 23

## 24 **Abbreviations** 24

25 <b>AI-REML:</b> Average information restricted maximum likelihood	25
26 <b>BLUP:</b> Best linear unbiased prediction	26
27 <b>BTA:</b> Bos taurus autosome	27
28 <b>cM:</b> CentiMorgan	28
29 <b>DNA:</b> Deoxyribonucleic acid	29
30 <b>FPR:</b> False positive rate	30
31 <b>GWAS:</b> Genomewide association study	31
32 <b>LD:</b> Linkage disequilibrium	32
33 <b>Mbp:</b> Mega base pairs	33
<b>MHC:</b> Major histocompatibility complex	33
<b>QTL:</b> Quantitative trait locus	33
<b>ROC:</b> Receiver operating characteristic	33
<b>SNP:</b> Single nucleotide polymorphism	33
<b>TPR:</b> True positive rate	33

1	<b>Declarations</b>	1
2	<b>Ethics approval and consent to participate</b>	2
	Not applicable.	
3		3
4	<b>Consent for publication</b>	4
	Not applicable.	
5		5
	<b>Availability of data and materials</b>	
6	All R scripts used are provided in Additional files 2–6. The R package <code>hscovar</code> version 0.4.0 is available at CRAN.	6
	The cattle data are accessible through RADAR <a href="https://dx.doi.org/10.22000/280">https://dx.doi.org/10.22000/280</a> . The mouse data are	
7	obtainable from <a href="https://wp.cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/">https://wp.cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/</a> . Maize data are available	7
8	at NCBI GEO under Accession Number GSE50558.	8
9	<b>Competing interests</b>	9
	The authors declare that they have no competing interests.	
10		10
	<b>Funding</b>	
11	The project was funded by the German Research Foundation (DFG, WI 4450/2-1). The funder had no role in the	11
12	design of the study and collection, analysis, and interpretation of data and in writing the manuscript.	12
13	<b>Authors' contributions</b>	13
	DW developed the theory, implemented the statistical methods, performed the analysis, and wrote the manuscript.	
14	MD and JK contributed to software development and improved the manuscript. All authors have read and approved	14
15	the final manuscript.	15
16	<b>Acknowledgments</b>	16
	We thank the Reviewers for their helpful comments.	
17		17
	<b>References</b>	
18	1. Wittenburg D, Bonk S, Doschoris M, Reyer H. Design of experiments for fine-mapping quantitative trait loci in	18
19	livestock populations. <i>BMC Genetics</i> . 2020;21:66.	19
20	2. Carlson CS, Eberle MA, Rieder MJ, Yi Q, Kruglyak L, Nickerson DA. Selecting a Maximally Informative Set of	20
	Single-Nucleotide Polymorphisms for Association Analyses Using Linkage Disequilibrium. <i>The American Journal</i>	
21	of Human Genetics. 2004;74(1):106–120.	21
22	3. Dehman A, Ambroise C, Neuvial P. Performance of a Blockwise Approach in Variable Selection Using Linkage	22
	Disequilibrium Information. <i>BMC Bioinformatics</i> . 2015;16:148.	
23	4. Kim SA, Cho CS, Kim SR, Bull SB, Yoo YJ. A New Haplotype Block Detection Method for Dense Genome	23
24	Sequencing Data Based on Interval Graph Modeling of Clusters of Highly Correlated SNPs. <i>Bioinformatics</i> .	24
	2018;34(3):388–397.	
25	5. Hill W, Robertson A. Linkage Disequilibrium in Finite Populations. <i>TAG Theoretical and Applied Genetics</i> .	25
	1968;38(6):226–231.	
26	6. Clayton D. <code>snpStats: SnpMatrix and XSnpmatrix Classes and Methods</code> ; 2017. R package version 1.34.0.	26
	Available from: <a href="http://bioconductor.org/packages/release/bioc/html/snpStats.html">http://bioconductor.org/packages/release/bioc/html/snpStats.html</a> .	
27	7. Clayton D, Leung HT. An R package for analysis of whole-genome association studies. <i>Human heredity</i> .	27
	2007;64(1):45–51.	
28	8. Caliński T, Harabasz J. A dendrite method for cluster analysis. <i>Communications in Statistics-theory and</i>	28
29	<i>Methods</i> . 1974;3(1):1–27.	29
30	9. Koivula M, Strandén I, Su G, Mäntysaari EA. Different Methods to Calculate Genomic	30
	Predictions—Comparisons of BLUP at the Single Nucleotide Polymorphism Level (SNP-BLUP), BLUP at the	
31	Individual Level (G-BLUP), and the One-Step Approach (H-BLUP). <i>Journal of Dairy Science</i> .	31
	2012;95(7):4065–4073.	
32	10. Butler D, Cullis BR, Gilmour A, Gogel B. <i>ASReml-R Reference Manual</i> . The State of Queensland, Department	32
33	of Primary Industries and Fisheries, Brisbane. 2009; Available from: <a href="https://asreml.kb.vsnr.co.uk/wp-content/uploads/sites/3/2018/02/ASReml-R-3-Reference-Manual.pdf">https://asreml.kb.vsnr.co.uk/wp-content/uploads/sites/3/2018/02/ASReml-R-3-Reference-Manual.pdf</a> .	33

- 1 11. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2020. Available 1  
 2 from: <https://www.R-project.org/>. 2  
 3 12. Valdar W, Solberg LC, Gauguier D, Burnett S, Klenerman P, Cookson WO, et al. Genome-Wide Genetic 3  
 4 Association of Complex Traits in Heterogeneous Stock Mice. *Nature Genetics*. 2006 Aug;38(8):879–887. 4  
 5 13. Browning SR, Browning BL. Rapid and Accurate Haplotype Phasing and Missing-Data Inference for 5  
 6 Whole-Genome Association Studies by Use of Localized Haplotype Clustering. *American Journal of Human 6  
 7 Genetics*. 2007 Nov;81(5):1084–1097. 7  
 8 14. Ferdosi M, Kinghorn B, van der Werf J, Lee S, Gondro C. hspbase: An R Package for Pedigree Reconstruction, 8  
 9 Detection of Recombination Events, Phasing and Imputation of Half-Sib Family Groups. *BMC Bioinformatics*. 9  
 10 2014;15(1):172. 10  
 11 15. Hampel A, Teuscher F, Gomez-Raya L, Doschoris M, Wittenburg D. Estimation of Recombination Rate and 11  
 12 Maternal Linkage Disequilibrium in Half-Sibs. *Frontiers in Genetics*. 2018;9:186. 12  
 13 16. Bauer E, Falque M, Walter H, Bauland C, Camisan C, Campo L, et al. Intraspecific variation of recombination 13  
 14 rate in maize. *Genome Biol*. 2013;14:R103. 14  
 15 17. Gaynor RC, Gorjanc G, Hickey JM. AlphaSimR: An R-package for Breeding Program Simulations. *G3: Genes 15  
 16 Genomes Genetics*. 2020;. 16  
 17 18. Schneider KL, Xie Z, Wolfruber TK, Presting GG. Inbreeding drives maize centromere evolution. *Proceedings 17  
 18 of the National Academy of Sciences*. 2016;113(8):E987–E996. 18  
 19 19. Yuan M, Lin Y. Model Selection and Estimation in Regression with Grouped Variables. *Journal of the Royal 19  
 20 Statistical Society B*. 2006;68(1):49–67. 20  
 21 20. Simon N, Friedman J, Hastie T, Tibshirani R. A Sparse-Group Lasso. *Journal of Computational and Graphical 21  
 22 Statistics*. 2013;22(2):231–245. 22  
 23 21. Gabriel SB, Schaffner SF, Nguyen H, Moore JM, Roy J, Blumenstiel B, et al. The Structure of Haplotype 23  
 24 Blocks in the Human Genome. *Science*. 2002;296(5576):2225–2229. 24  
 25 22. Wang S, He S, Yuan F, Zhu X. Tagging SNP-Set Selection with Maximum Information Based on Linkage 25  
 26 Disequilibrium Structure in Genome-Wide Association Studies. *Bioinformatics*. 2017;33(14):2078–2081. 26  
 27 27. 27  
 28 28. 28  
 29 29. 29  
 30 30. 30  
 31 31. 31  
 32 32. 32  
 33 33. 33

## 18 Figures 18

19 19  
 20 **Figure 1** (a) Matrix highlights the SNP pairs with correlation larger than a certain threshold for a 20  
 21 given SNP panel. The SNP involved most is marked by an arrow (Step A). (b) Correlations are 21  
 22 considered within the corresponding SNP subset and a tagSNP is selected (Step B). (c) All SNPs 22  
 23 associated with the tagSNP are removed from the remaining ungrouped SNP panel for the next 23  
 24 round of iteration (Step A). Iterations continue until no ungrouped SNP is left. 24

25 25  
 26 **Figure 2** Correlation (a) and LD matrix (b) for mouse data on chromosome 17. The red dots 26  
 27 highlight representative SNPs of groups with at least three SNPs. In total, 394 SNPs were 27  
 28 considered for the correlation and population-LD matrix. 28

29 29  
 30 **Figure 3** Correlation (a) and LD matrix (b) for cattle data in a target region of chromosome 30  
 31 1:20.6–39.4 Mbp. The red dots highlight representative SNPs of groups with at least three SNPs. 31  
 32 In total, 300 SNPs were considered. The correlation matrix was set up at 239 SNPs and the 32  
 33 population-LD matrix at 300 SNPs; missing values are filled in white color. 33

**Figure 4** Correlation (a) and LD matrix (b) for maize data on chromosome 2. The red dots highlight representative SNPs of groups with at least three SNPs. In total, 956 SNPs were considered. The correlation matrix was set up at 953 SNPs and the population-LD matrix at 956 SNPs; missing values are filled in white color.

**Figure 5** Sensitivity and specificity of testing SNP effects depending on threshold  $t = 0.8$  (a) and  $t = 0.5$  (b). ROC curves are based on 100 repeated simulations of genotypes and phenotypes in  $N = 1$  half-sib family with 1 000 progeny (two QTL signals, heritability 0.3).

**Table 1** Number of groups, number of groups with at least three SNPs and Calinski-Harabasz index (CH). Number of SNPs corresponds to the method applied (family or population LD). Average values of 100 repetitions are presented for simulated data and  $t = 0.8$ . Computing time for grouping markers was up to 0.2s except for the maize data which required 2s.

	Families	QTLs	Family approach				Population-LD approach			
			SNPs	Groups	$\geq 3$	CH	SNPs	Groups	$\geq 3$	CH
Simulation	10	2	283	59	10	80.7	300	21	9	40.0
	10	5	282	61	10	75.0	300	17	9	38.9
	5	2	282	61	10	71.9	300	29	8	32.6
	5	5	281	64	10	85.6	300	24	9	35.0
	1	2	281	59	9	61.4	300	56	6	20.5
	1	5	282	59	9	83.0	300	49	7	25.0
Mouse	138		394	83	49	38.8	394	98	51	20.6
Cattle	5		237	172	11	2.9	300	210	11	2.4
Maize	13		953	426	93	10.4	956	576	70	21.9

**Tables**

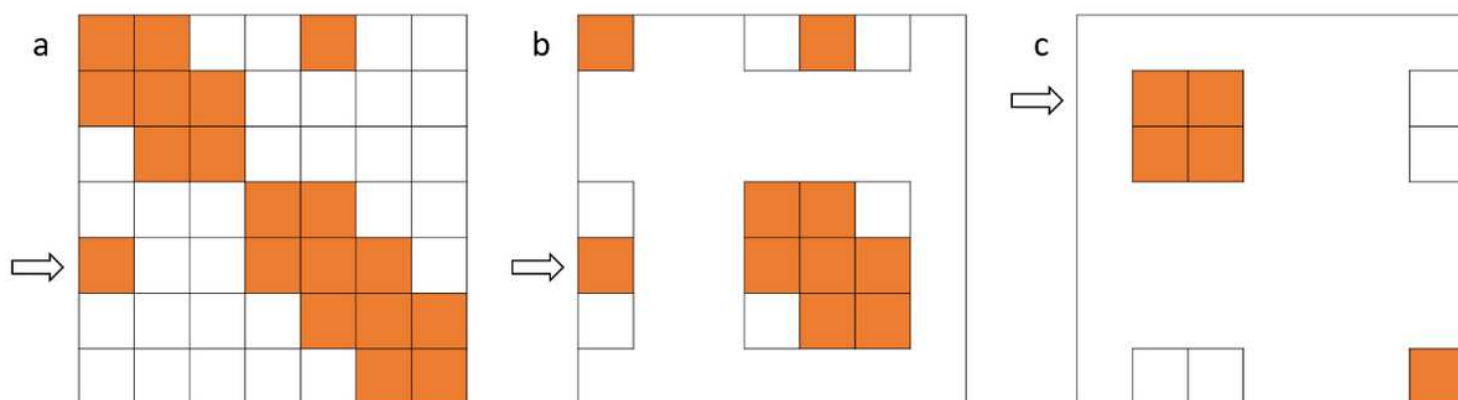
**Additional Files**

- Additional file 1 – PDF file summarizing the derivation of family-based correlation matrices  
Given the population structure, half-sib families or full-sib families, the covariance matrix is analytically retrieved. Its computation using the R package `hscovar` is shown.
- Additional file 2 — TXT file containing R code for mouse data  
Raw mouse data are processed and the matrix of correlation between markers is derived.
- Additional file 3 — TXT file containing R code for cattle data  
Cattle data are processed and the matrix of correlation between markers is derived.
- Additional file 4 — TXT file containing R code for maize data  
Raw maize data are processed and the matrix of correlation between markers is derived.
- Additional file 5 — TXT file containing R code for simulation study  
With this script, genotype and phenotype data of half-sib families are simulated and a genome-based association analysis is carried out.
- Additional file 6 — TXT file containing R code for creating plots  
Plots of correlation matrices and population-LD matrices are produced based on results with additional files 2–5.

<sup>1</sup> Additional file 7 — PDF file with additional tables	1
<sup>2</sup> Number of groups, number of groups with at least three SNPs and Calinski-Harabasz index for different simulation scenarios and varying threshold.	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	23
24	24
25	25
26	26
27	27
28	28
29	29
30	30
31	31
32	32
33	33

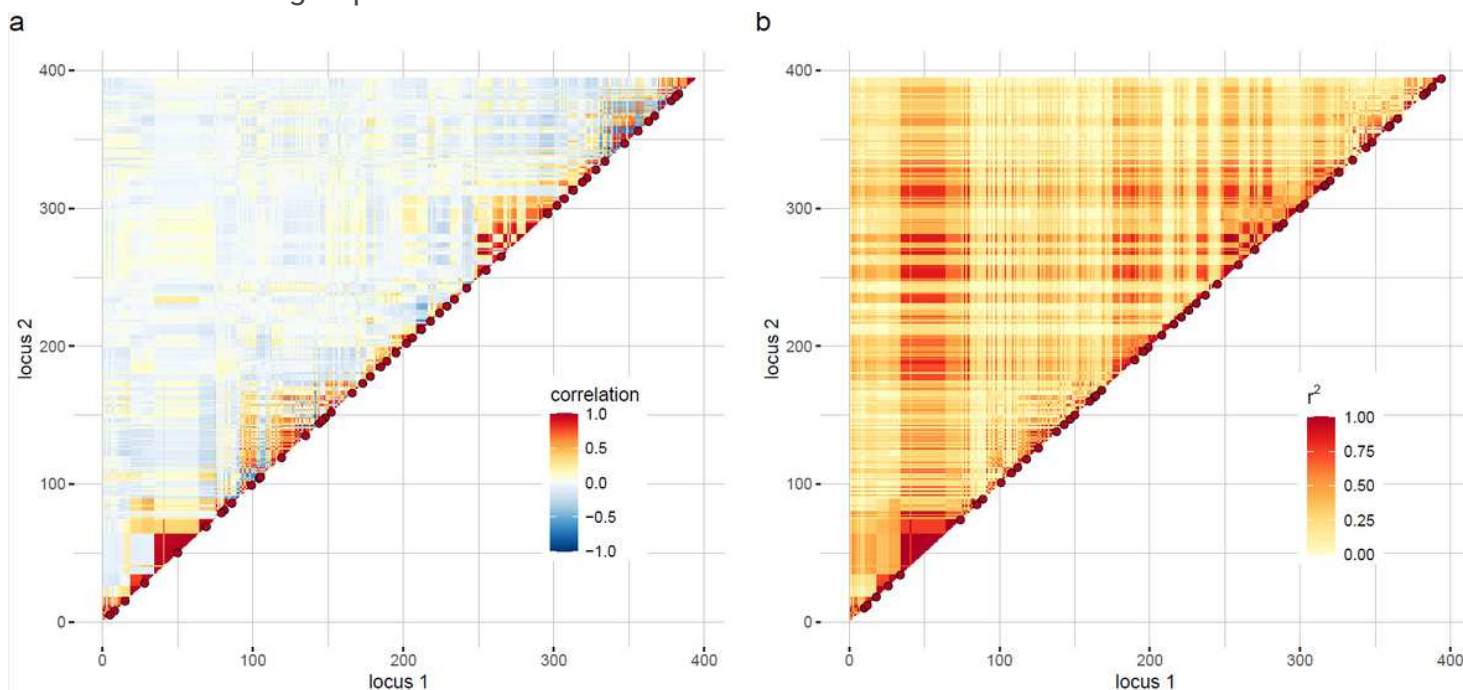


# Figures



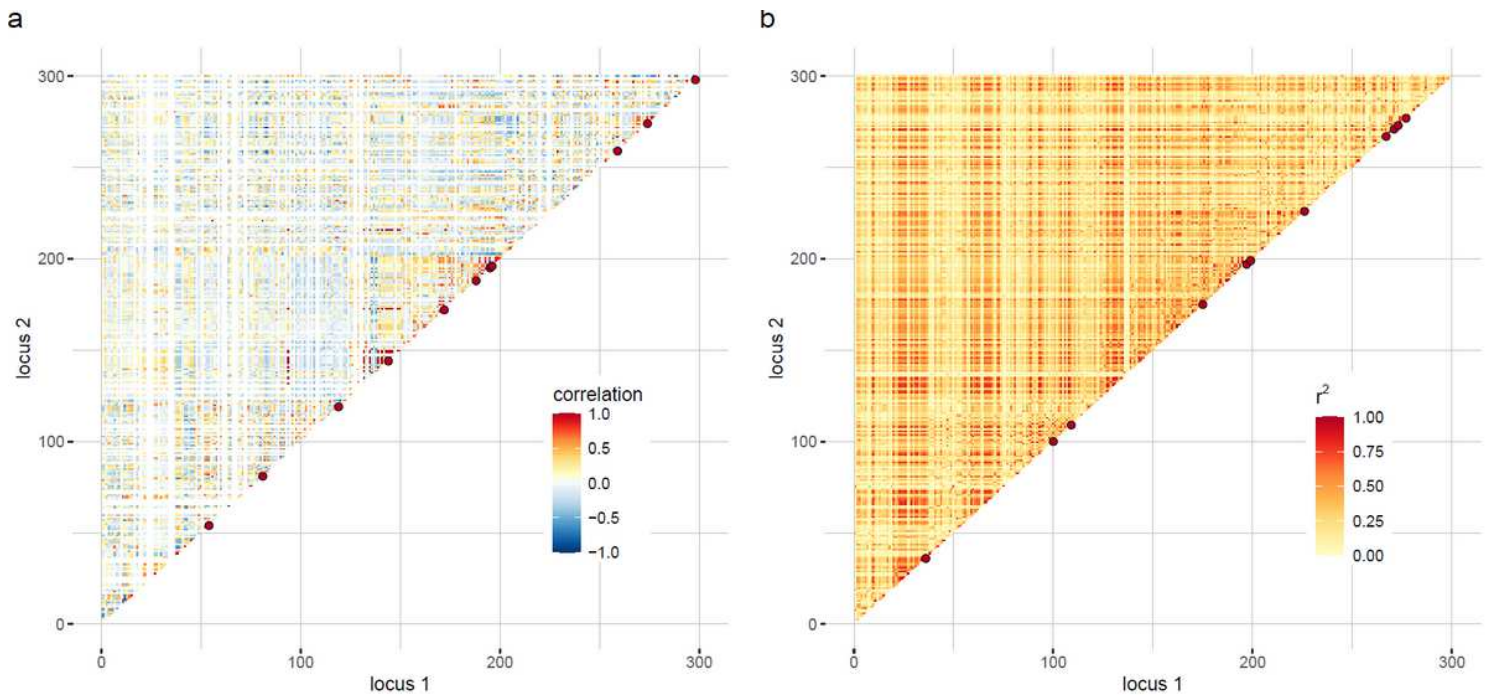
**Figure 1**

(a) Matrix highlights the SNP pairs with correlation larger than a certain threshold for a given SNP panel. The SNP involved most is marked by an arrow (Step A). (b) Correlations are considered within the corresponding SNP subset and a tagSNP is selected (Step B). (c) All SNPs associated with the tagSNP are removed from the remaining ungrouped SNP panel for the next round of iteration (Step A). Iterations continue until no ungrouped SNP is left.



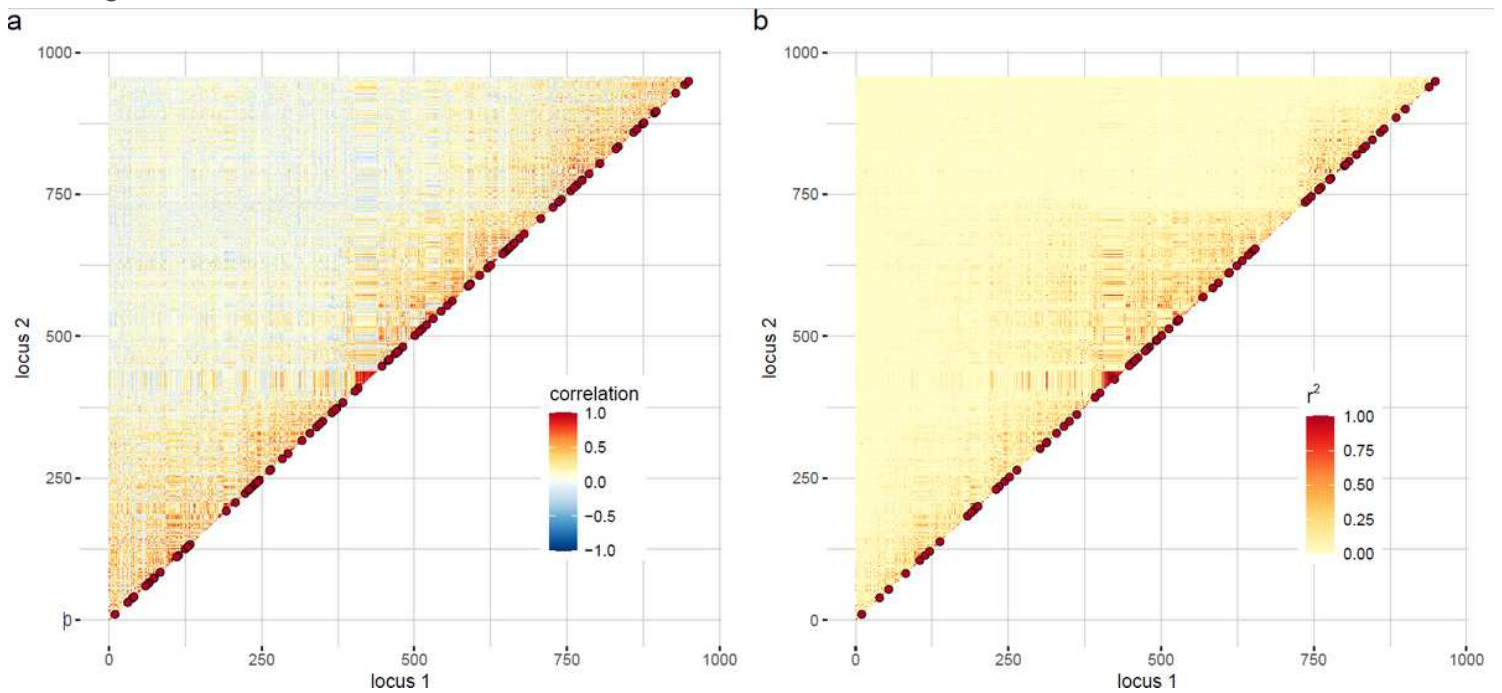
**Figure 2**

Correlation (a) and LD matrix (b) for mouse data on chromosome 17. The red dots highlight representative SNPs of groups with at least three SNPs. In total, 394 SNPs were considered for the correlation and population-LD matrix.



**Figure 3**

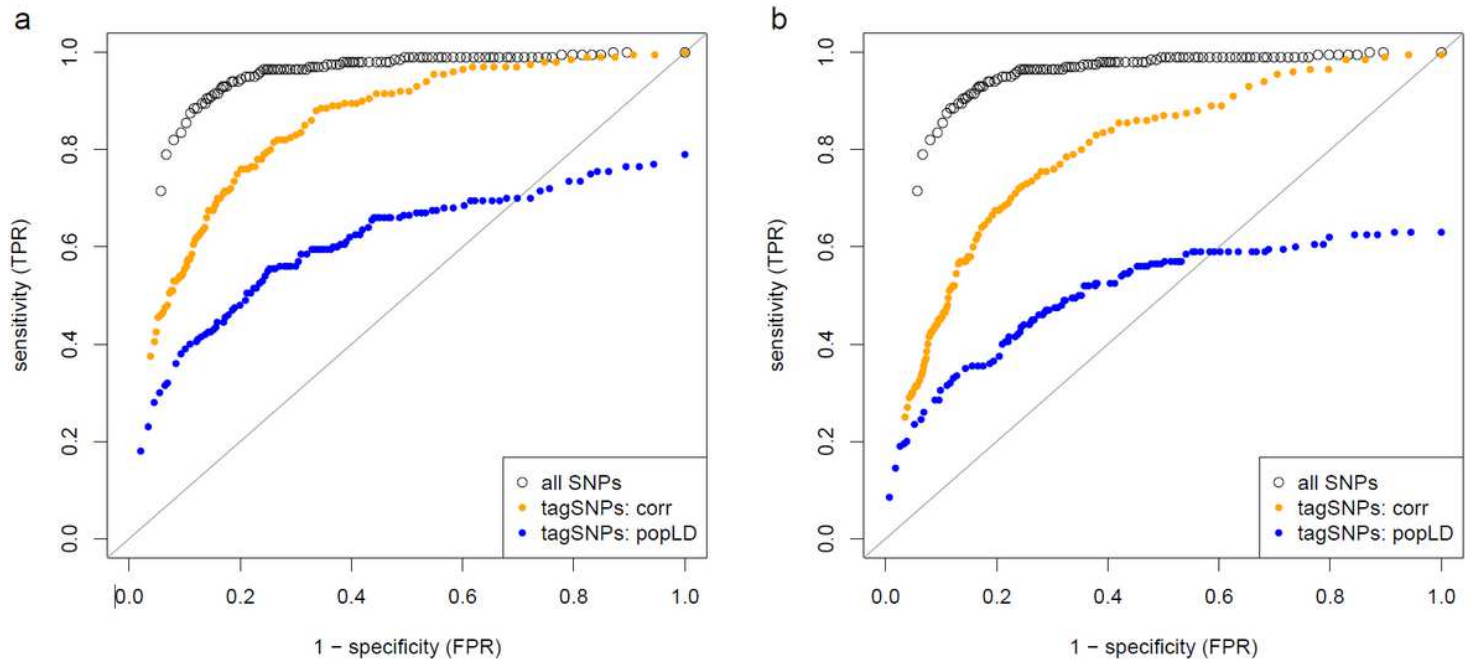
Correlation (a) and LD matrix (b) for cattle data in a target region of chromosome 1:20.6–39.4 Mbp. The red dots highlight representative SNPs of groups with at least three SNPs. In total, 300 SNPs were considered. The correlation matrix was set up at 239 SNPs and the population-LD matrix at 300 SNPs; missing values are filled in white color.



**Figure 4**

Correlation (a) and LD matrix (b) for maize data on chromosome 2. The red dots highlight representative SNPs of groups with at least three SNPs. In total, 956 SNPs were considered. The correlation matrix was

set up at 953 SNPs and the population-LD matrix at 956 SNPs; missing values are filled in white color.



**Figure 5**

Sensitivity and specificity of testing SNP effects depending on threshold  $t = 0.8$  (a) and  $t = 0.5$  (b). ROC curves are based on 100 repeated simulations of genotypes and phenotypes in  $N = 1$  half-sib family with 1 000 progeny (two QTL signals, heritability 0.3).

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [additionalfile1covariance.pdf](#)
- [additionalfile2preparemousedata.R.txt](#)
- [additionalfile3preparecattledata.R.txt](#)
- [additionalfile4preparemaizedata.R.txt](#)
- [additionalfile5simstudygwas.R.txt](#)
- [additionalfile6plotallexamples.R.txt](#)
- [additionalfile7tables.pdf](#)