

# UC Davis

## UC Davis Previously Published Works

### Title

Average semivariance yields accurate estimates of the fraction of marker-associated genetic variance and heritability in complex trait analyses.

### Permalink

<https://escholarship.org/uc/item/00t8d6h6>

### Journal

PLoS genetics, 17(8)

### ISSN

1553-7390

### Authors

Feldmann, Mitchell J  
Piepho, Hans-Peter  
Bridges, William C  
[et al.](#)

### Publication Date

2021-08-01

### DOI

10.1371/journal.pgen.1009762

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

## RESEARCH ARTICLE

## Average semivariance yields accurate estimates of the fraction of marker-associated genetic variance and heritability in complex trait analyses

Mitchell J. Feldmann<sup>1</sup>, Hans-Peter Piepho<sup>2</sup>, William C. Bridges<sup>3</sup>, Steven J. Knapp<sup>1\*</sup>

**1** Department of Plant Sciences, University of California, Davis, California, United States of America, **2** Biostatistics Unit, Institute of Crop Science, University of Hohenheim, Stuttgart, Germany, **3** Department of Mathematical Sciences, Clemson University, Clemson, South Carolina, United States of America

\* [sjknapp@ucdavis.edu](mailto:sjknapp@ucdavis.edu)**OPEN ACCESS**

**Citation:** Feldmann MJ, Piepho H-P, Bridges WC, Knapp SJ (2021) Average semivariance yields accurate estimates of the fraction of marker-associated genetic variance and heritability in complex trait analyses. *PLoS Genet* 17(8): e1009762. <https://doi.org/10.1371/journal.pgen.1009762>

**Editor:** Trudy F. C. Mackay, Clemson University, UNITED STATES

**Received:** March 31, 2021

**Accepted:** August 9, 2021

**Published:** August 26, 2021

**Copyright:** © 2021 Feldmann et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The custom R scripts for reproducing our simulations have been deposited in a public GitHub repository <https://github.com/mjfeldmann/VarCompSim>. The simulated data shown in Figs 1, 2, S1 and S2 have been deposited in a public Zenodo repository <https://doi.org/10.5281/zenodo.3742421>.

**Funding:** This research was supported by grants to SJK from the United States Department of Agriculture (<http://dx.doi.org/10.13039/>

**Abstract**

The development of genome-informed methods for identifying quantitative trait loci (QTL) and studying the genetic basis of quantitative variation in natural and experimental populations has been driven by advances in high-throughput genotyping. For many complex traits, the underlying genetic variation is caused by the segregation of one or more ‘large-effect’ loci, in addition to an unknown number of loci with effects below the threshold of statistical detection. The large-effect loci segregating in populations are often necessary but not sufficient for predicting quantitative phenotypes. They are, nevertheless, important enough to warrant deeper study and direct modelling in genomic prediction problems. We explored the accuracy of statistical methods for estimating the fraction of marker-associated genetic variance ( $p$ ) and heritability ( $H_M^2$ ) for large-effect loci underlying complex phenotypes. We found that commonly used statistical methods overestimate  $p$  and  $H_M^2$ . The source of the upward bias was traced to inequalities between the expected values of variance components in the numerators and denominators of these parameters. Algebraic solutions for bias-correcting estimates of  $p$  and  $H_M^2$  were found that only depend on the degrees of freedom and are constant for a given study design. We discovered that average semivariance methods, which have heretofore not been used in complex trait analyses, yielded unbiased estimates of  $p$  and  $H_M^2$ , in addition to best linear unbiased predictors of the additive and dominance effects of the underlying loci. The cryptic bias problem described here is unrelated to selection bias, although both cause the overestimation of  $p$  and  $H_M^2$ . The solutions we described are predicted to more accurately describe the contributions of large-effect loci to the genetic variation underlying complex traits of medical, biological, and agricultural importance.

**Author summary**

The contributions of individual genes to the phenotypic variation observed for genetically complex traits has been an ongoing and important challenge in biology, medicine, and

100000199) National Institute of Food and Agriculture (NIFA) Specialty Crops Research Initiative (# 2017-51181-26833) and California Strawberry Commission (<http://dx.doi.org/10.13039/100006760>), in addition to funding from the University of California, Davis (<http://dx.doi.org/10.13039/100007707>). HPP was supported by the German Research Foundation (DFG) grant PI 377/18-1. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

agriculture. While many genes have statistically undetectable effects, those with large effects often warrant in-depth study and can be important predictors of complex phenotypes such as disease risk in humans or disease resistance in domesticated plants and animals. The genes identified through associations with genetic markers in complex trait analyses typically account for a fraction of the heritable variation, a genetic parameter we called ‘marker heritability’. We discovered that textbook statistical methods systematically overestimate marker heritability and thus overestimate the contributions of specific genes to the phenotypic variation observed for complex traits in natural and experimental populations. We describe the source of the upward bias, validate our findings through computer simulation, describe methods for bias-correcting estimates of marker heritability, and illustrate their application through empirical examples. The statistical methods we describe supply investigators with more accurate estimates of the contributions of specific genes or networks of interacting genes to the heritable variation observed in complex trait studies.

## Introduction

The genetic variation observed in nature is frequently caused by genes with quantitative effects [1–7]. Their discovery and characterization has been a dominant feature of quantitative genetic studies in biology, evolution, agriculture, and medicine since the introduction of methods for genotyping DNA variants genome-wide [8–11], and the parallel development of statistical methods for finding associations between DNA variants and the underlying genes or quantitative trait loci (QTL) [2, 4, 5, 7, 12–16]. A significant breakthrough was achieved when Lander and Botstein [12] introduced ‘interval mapping’ and showed that genomes could be systematically searched to identify QTL in populations genotyped with a genome-wide framework of genetically mapped DNA markers. As genotyping technologies advanced and marker densities increased, genome-wide association study (GWAS) methods emerged to search genomes for genotype-to-phenotype associations by exploiting the historical recombination in populations [14, 15, 17–19]. The concept of genomic prediction emerged as a counterpart to GWAS, initially for estimating genomic-estimated breeding values (GEBVs) in domesticated plants and animals and later for estimating polygenic risk scores (PRSs) in humans and model organisms [20–23]. These technical advances precipitated a consequential shift in the study of quantitative traits from analyses of phenotypic variation limited and informed by pedigree or family data to genome-wide analyses of genotype-to-phenotype associations and genomic prediction informed by genotypic data [6, 7, 13, 16, 20, 24–31].

The phenotypic variation observed in a population is customarily partitioned into genetic and non-genetic components to estimate heritability, repeatability, and reliability of the quantitative traits under study [24, 25, 30, 32]. The genetic component can be caused by any number of genes with quantitative effects, even a single gene, but more often by multiple genes with a range of effects [31, 33–43]. For most quantitative traits, that number is unknown but presumed to be large and undiscoverable [3, 6, 7, 22, 32, 34]. Because genes with small effects are challenging to identify and validate, the ‘many genes with small effects’ hypothesis has been difficult to conclusively falsify [21, 22, 32]. Despite the uncertainty surrounding the identity, number, effects, and interactions of genes in the undiscovered fraction [6], three decades of complex trait analyses in humans, domesticated plants and animals, *Drosophila*,

Arabidopsis, yeast, mice, zebrafish, and other organisms have shown that the ‘discovered’ genes are typically small in number, large in effect, and collectively only explain a fraction of the genetic variance ( $\sigma_G^2$ ) [13, 16, 28, 32, 34–36, 44, 45]. The unexplained fraction has been called ‘missing heritability’ [46–48].

The discovered genes in polygenic systems of genes are often necessary but not sufficient for predicting quantitative phenotypes, e.g., disease risk in humans or yield in domesticated plants and animals [3, 21, 34, 42, 44, 49]. There is a large body of evidence that the QTL effects for many quantitative traits are gamma family distributed, where the discovered genes are found in the upper or thin tail of the distribution above the threshold of statistical significance [34]. The presumption is that the lower or heavy tail of the gamma family distribution is caused by many genes with small effects, the chief tenet of the infinitesimal model of quantitative genetics [6, 26, 32, 50]. Genes with large effects often dominate the ‘non-missing heritability’, mask or obscure the effects of other quantitatively acting genes, and pleiotropically affect multiple quantitative phenotypes [16, 35, 39, 51], e.g., mutations in the *BRCA2* gene can have large effects, are incompletely penetrant, interact with other genes, and are necessary but not sufficient for predicting breast, ovarian, and other cancer risks in women [52]. The large-effect QTL BTA19 pleiotropically affects milk yield, protein yield, and productive life in Guernsey cattle (*Bos taurus*) [43], and branching and pigment genes (*BR*, *PHY*, and *HYP*) have large effects, interact, and pleiotropically affect several genetically correlated seed biomass traits in sunflower (*Helianthus annuus*) [53]. Despite decades of directional selection, loci with large effects often segregate (have not been fixed) in domesticated plant and animal populations [33, 34, 37, 38, 40, 54, 55]. The fractions of the genetic variances explained by *BRCA2*, BTA19, *BR*, *PHY*, and *HYP* were not reported in those studies. What fraction of the heritability for breast cancer risk, for example, can be explained by the known mutations in *BRCA2*? Our study explored the accuracy of methods for estimating that parameter.

Our surveys and others substantiate that the missing and non-missing fractions of the genetic variance are commonly either not estimated or inaccurately estimated in GWAS and other gene finding studies, e.g., the statistical significance of individual marker loci from sequential regression analyses are typically reported without correcting for the effects of other discovered marker loci through multilocus partial regression analyses or Type III ANOVA [17, 19, 22, 34, 56]. Such analyses are necessary for accurately assessing the statistical importance of the underlying gene and gene-gene interaction effects in a multilocus system, e.g., when multiple loci are identified by GWAS (sequential analyses of individual loci), their effects are more accurately estimated by simultaneous analysis using partial regression analysis approaches and even then can be upwardly biased [51]. The estimation problem we studied is intertwined with the broader problem of accurately describing multilocus systems of genes with large effects. We show that the discovered fraction of the genetic variance can be grossly overestimated and that the cause of the problem is a mathematical artifact in the expected values of variance components and their ratios. We revisited the problem of estimating the non-missing and missing fractions of heritability in candidate gene and other complex trait analyses, in part because of the systematic upward bias we discovered, in addition to inconsistencies in the methods commonly applied to the problem. The solutions to the problem presented here are straightforward and primarily applicable to the study of genes with large effects, especially those affecting the accuracy of genomic predictions for disease risk or breeding value [21, 43]. The optimum approaches for weighting or correcting for loci with large effects in genomic prediction are not completely clear; however, in artificial selection settings where the favorable alleles for discovered loci are unequivocally known, those alleles can be directly

selected via marker-assisted selection (MAS) with genomic selection exerting pressure on unknown loci underlying the additive genetic variance not explained by the segregation of known large effect loci [54, 57–61].

Lande and Thompson [62] proposed the parameter  $p = \sigma_M^2 / \sigma_G^2$  to estimate the discovered or non-missing fraction of the genetic variance, where  $\sigma_M^2$  is the fraction of the genetic variance associated with statistically significant markers in linkage disequilibrium (LD) with genes or QTL affecting the trait under study (here QTL refers to a chromosome segment predicted to harbor a gene or genes affecting a quantitative trait). Similarly, marker heritability ( $H_M^2 = \sigma_M^2 / \sigma_P^2$ ) estimates the non-missing fraction of the phenotypic variance ( $\sigma_P^2$ ) associated with statistically significant markers in LD with causal genes or QTL. Here a distinction needs to be made between  $H_M^2$  and genomic heritability, a parameter estimated by summing the effects of a dense genome-wide sample of markers, only some of which are predicted to be in LD with the underlying causal genes or QTL [27, 30, 63]. We are not proposing marker heritability as a replacement or substitute for genomic heritability but as a parameter for parsing out the non-missing fraction of heritability associated with discovered loci, especially loci like *BRCA2* and *BTA19* [43, 52]. The genetic variance component ( $\sigma_G^2$ ) in these ratios can be estimated from pedigree or family information (as shown in our examples) or genomic information (as reviewed by [30] and [63]). For either,  $\sigma_M^2$  is simply the variance explained by marker loci with effects large enough to be statistically detected and important enough to be specifically studied and modeled, perhaps as fixed effects [22, 39, 40, 51, 61]. Despite a direct and logical connection to heritability, estimates of  $p$  and  $H_M^2$  are seldom reported in complex trait studies, whereas genomic heritability estimates are commonly reported in genomic prediction studies [30, 34, 62].

Here we show that  $p$  and  $H_M^2$  are often overestimated in complex trait analyses. The problem we discovered is unrelated to selection bias, the phenomena where the effects of discovered QTL are inflated by biased sampling from truncated distributions with small sample sizes [64–69], and unrelated to the upward biases known to arise in GWAS [70]. While selection bias is a well known and widely cited problem in complex trait analyses, we describe a previously unreported and cryptic source of bias in estimates of  $p$  and  $H_M^2$ . To identify the source of the bias and explore the problem in greater depth, we compared the accuracy of average marginal variance (AMV) [71, 72] and average semivariance (ASV) [73] methods for estimating  $p$  and  $H_M^2$ . AMV is the acronym applied throughout this paper for the ANOVA and REML variance component estimation methods commonly described in textbooks and implemented in statistical software for the analysis of generalized linear mixed models (GLMMs), e.g., the ‘lme4’ R package and the SAS packages ‘GLM’ and ‘GLIMMIX’ [24, 25, 72, 74–79]. We introduced the average marginal variance terminology here to facilitate comparisons of the differences between AMV and ASV methods for estimating variance component *ratios*. The ASV methods we applied to the problem are extensions of those described by Piepho [73] for estimating the total variance and coefficient of determination ( $R^2$ ) in GLMM analyses. For the AMV and ASV analyses shown throughout this paper, REML was used to estimate the variance components [56, 72, 75, 79]. The source of the bias was discovered, however, through algebraic analyses of the expected mean squares (EMSs) from ANOVA. We describe that source and approaches for bias-correcting ANOVA or REML estimates of  $p$  and  $H_M^2$  from the commonly applied AMV methods. We show that ASV methods directly yield unbiased estimates of  $p$  and  $H_M^2$  that are identical to bias-corrected AMV estimates. Finally, we discuss the connection of these random effects methods to the fixed effect methods commonly applied in QTL mapping and genome-wide association studies [51, 80, 81].

## Results and discussion

### Overestimation of the genetic variance explained by markers in linkage disequilibrium with causative genes or QTL

The overestimation problem described here was originally discovered in a reanalysis of data from genetic studies in plants where REML estimates of  $H_M^2$  exceeded REML estimates of broad-sense heritability ( $H^2$ ) and REML estimates of  $p$  and  $H_M^2$  exceeded 1.0, the theoretical upper limit for these parameters (Table 1). We initially suspected that selection bias might be the culprit [68, 69, 82–84] but concluded that selection bias alone could not explain  $\hat{p} > 1.0$  or  $H_M^2 > 1.0$ . Although proof was lacking and the bias was non-obvious, we hypothesized that many estimates in the theoretical range ( $0.0 \leq p \leq 1.0$ ) must also be upwardly biased. The proof was found through algebraic analyses of the ANOVA estimators of  $\sigma_M^2$ ,  $\sigma_G^2$ , and  $\sigma_p^2$  for balanced and unbalanced data (S1, S2 and S4 Texts). Although variance components are commonly estimated using REML, as was done in the analyses shown throughout this paper, algebraic analyses of ANOVA expected mean squares (EMSs) identified the source of the bias and yielded explicit algebraic solutions for bias correcting ANOVA and REML estimates of  $p$  and  $H_M^2$ .

The source of the bias was identified by expressing the estimator of  $p$  as a function of the ANOVA estimators of  $\sigma_M^2$  and  $\sigma_G^2$  for balanced data and algebraically simplifying the equations. The linear mixed models (LMMs) and ANOVA estimators of the variance components needed to show this are described here. We start with the analysis of a single marker locus in an experiment where entries (e.g., individuals, families, or strains) are replicated,  $\sigma_G^2$  can be estimated, and the data for entries and markers are balanced. Extensions for one to three marker loci with unbalanced data are shown in S1, S2 and S3 Texts. Two LMMs are needed for estimating  $\sigma_M^2$ ,  $\sigma_G^2$ ,  $p$ , and  $H_M^2$ . Consider a study where  $n_G$  entries are phenotyped for a normally distributed quantitative trait using a balanced completely randomized study design with  $r_G$  replications/entry,  $n_M$  marker genotypes/locus, and  $r_M$  replications/marker genotype. The LMM needed for estimating  $\sigma_G^2$  (the between entry variance component) is:

$$y_{jk} = \mu + G_j + \epsilon_{jk} \tag{1}$$

where  $y_{jk}$  is the  $jk^{th}$  phenotypic observation,  $\mu$  is the population mean,  $G_j$  is the random effect of the  $j^{th}$  entry,  $\epsilon_{jk}$  is the random effect of the  $jk^{th}$  residual,  $G_j \sim N(0, \sigma_G^2)$ ,  $\epsilon_{jk} \sim N(0, \sigma_\epsilon^2)$ ,  $j = 1, 2, \dots, n_G$ , and  $k = 1, 2, \dots, r_G$ . Suppose entries are genotyped for a single marker locus ( $M$ ) in linkage disequilibrium with a gene or QTL affecting the quantitative phenotype ( $y_{jk}$ ). The between entry source of variation from LMM (1) can be partitioned into marker ( $M$ ) and entry nested in marker ( $G : M$ ) sources of variation (this is the residual genetic variation among entries not explained by markers in the model). The LMM for estimating  $\sigma_M^2$  and  $\sigma_{G:M}^2$  is:

$$y_{ijk} = \mu + M_i + G : M_{i(j)} + \epsilon_{ijk} \tag{2}$$

where  $y_{ijk}$  is the  $ijk^{th}$  phenotypic observation,  $M_i$  is the random effect of the  $i^{th}$  marker genotype at locus  $M$ ,  $G : M_{i(j)}$  is the random effect of the  $j^{th}$  entry nested in the  $i^{th}$  marker genotype,  $\epsilon_{ijk}$  is the random effect of the  $ijk^{th}$  residual,  $i = 1, 2, 3$ ,  $M_i \sim N(0, \sigma_M^2)$ ,  $G : M_{i(j)} \sim N(0, \sigma_{G:M}^2)$ , and  $\epsilon_{ijk} \sim N(0, \sigma_\epsilon^2)$ .

The ANOVA estimator of the between-entry variance component ( $\sigma_G^2$ ) from LMM (1) with balanced data is:

$$\hat{\sigma}_G^2 = \frac{MS_G - MS_\epsilon}{r_G} = \frac{SS_G/df_G - SS_\epsilon/df_\epsilon}{r_G} = \frac{1}{df_G r_G} SS_G - \frac{1}{df_\epsilon r_G} SS_\epsilon \tag{3}$$

**Table 1. REML estimates of marker-associated variance ( $\sigma_M^2$ ), the fraction of the genetic variance explained by markers ( $p = \sigma_M^2/\sigma_G^2$ ), and marker heritability ( $H_M^2 = \sigma_M^2/\sigma_p^2$ ) from random marker effects analyses and coefficients of determination ( $R^2$ ) from Type II and Type III fixed marker effects analyses for large effect loci identified in cattle, sunflower, and strawberry studies.**

Study	Source	df	$k_M$	Variance Component	Uncorrected			Bias Corrected			Type II	Type III
					$\hat{\sigma}^2$	$\hat{p}$	$\hat{H}^2$	$\hat{\sigma}_+^2$	$\hat{p}_+$	$\hat{H}_+^2$	$\hat{R}^{2d}$	$\hat{R}^{2e}$
Cattle White Spotting <sup>a</sup>	M	25	—	$\sigma_{rs10}^2 + \dots + \sigma_{rs10 \times rs45 \times rs20}^2$	7.92	—	0.76	3.88	—	0.37	—	—
	rs10	2	0.35	$\sigma_{rs10}^2$	0.62	—	0.06	0.21	—	0.02	0.04	0.00
	rs45	2	0.41	$\sigma_{rs45}^2$	2.91	—	0.28	1.20	—	0.11	0.21	0.08
	rs20	2	0.54	$\sigma_{rs20}^2$	3.81	—	0.37	2.04	—	0.20	0.23	0.10
	rs10 × rs45	4	0.58	$\sigma_{rs10 \times rs45}^2$	0.00	—	0.00	0.00	—	0.00	0.00	0.00
	rs10 × rs20	4	0.67	$\sigma_{rs10 \times rs20}^2$	0.00	—	0.00	0.00	—	0.00	0.01	0.01
	rs45 × rs20	4	0.70	$\sigma_{rs45 \times rs20}^2$	0.37	—	0.04	0.26	—	0.02	0.01	0.01
	rs10 × rs45 × rs20	7	0.77	$\sigma_{rs10 \times rs45 \times rs20}^2$	0.22	—	0.02	0.17	—	0.02	0.01	0.01
	G : rs10 × rs45 × rs20	2,935	—	$\sigma_{G:rs10 \times rs45 \times rs20}^2$	5.26	—	—	5.26	—	—	—	—
Sunflower Oil Content <sup>b</sup>	Entry (G)	145	—	$\sigma_G^2$	21.61	—	0.95	21.61	—	0.95	—	—
	M + G : M	145	—	$\sigma_B^2 + \dots + \sigma_{G:M}^2$	30.76	1.42	1.35	22.15	1.02	0.98	—	—
	M	7	—	$\sigma_B^2 + \dots + \sigma_{B \times P \times HYP}^2$	17.85	0.83	0.79	9.24	0.43	0.41	—	—
	BR	1	0.48	$\sigma_B^2$	11.57	0.54	0.51	5.59	0.26	0.25	0.21	0.26
	PHY	1	0.47	$\sigma_{PHY}^2$	1.26	0.06	0.06	0.60	0.03	0.03	0.02	0.04
	HYP	1	0.49	$\sigma_{HYP}^2$	2.9	0.13	0.13	1.41	0.07	0.06	0.05	0.10
	BR × PHY	1	0.77	$\sigma_{B \times P}^2$	0.21	0.01	0.01	0.17	0.01	0.01	0.01	0.01
	BR × HYP	1	0.78	$\sigma_{B \times HYP}^2$	1.89	0.09	0.08	1.46	0.07	0.06	0.03	0.04
	PHY × HYP	1	0.77	$\sigma_{PHY \times HYP}^2$	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	BR × PHY × HYP	1	0.88	$\sigma_{B \times PHY \times HYP}^2$	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.01
	G : BR × PHY × HYP	138	—	$\sigma_{G:B \times PHY \times HYP}^2$	12.91	0.60	0.57	12.91	0.60	0.57	—	—
	Residual (ε)	144	—	$\sigma_\epsilon^2$	2.07	—	—	2.07	—	—	—	—
	Strawberry Fusarium Wilt <sup>c</sup>	Entry (G)	557	—	$\sigma_G^2$	3.26	—	0.98	3.26	—	0.98	—
M + G : M		557	—	$\sigma_{AX396}^2 + \sigma_{G:AX396}^2$	4.77	1.46	1.44	2.39	0.73	0.72	—	—
AX396 (M)		2	0.47	$\sigma_{AX396}^2$	4.48	1.37	1.35	2.09	0.64	0.63	0.84	0.84
G : AX396		555	—	$\sigma_{G:AX396}^2$	0.30	0.09	0.09	0.30	0.09	0.09	—	—
Residual (ε)		1,631	—	$\sigma_\epsilon^2$	0.23	—	—	0.23	—	—	—	—
Strawberry Fusarium Wilt <sup>c</sup>	Entry (G)	540	—	$\sigma_G^2$	3.30	—	0.98	3.30	—	0.98	—	—
	M + G : M	540	—	$\sigma_{AX493}^2 + \sigma_{G:AX493}^2$	4.01	1.21	1.20	3.45	1.05	1.03	—	—
	AX493 (M)	2	0.62	$\sigma_{AX493}^2$	1.48	0.45	0.44	0.93	0.28	0.28	0.22	0.22
	G : AX493	538	—	$\sigma_{G:AX493}^2$	2.53	0.77	0.75	2.53	0.77	0.75	—	—
	Residual (ε)	1,584	—	$\sigma_\epsilon^2$	0.23	—	—	0.23	—	—	—	—

<sup>a</sup>Statistics are shown for three marker loci (*rs10*, *rs45*, and *rs20*) associated with genetic variation for white spotting (%) in a cattle population ( $n_G = 2, 973$ ) with a single phenotypic observation per individual and highly unbalanced marker data [85]. The marker loci were identified by GWAS. The linear mixed model for the cattle analysis was identical to that for the sunflower analysis without replications ( $r_G = 1$ ).  $k_M$  coefficient equations for three loci with unbalanced data are shown in S3 Text.

<sup>b</sup>Statistics are shown for three marker loci (*BR*, *PHY*, and *HYP*) associated with genetic variation for seed oil content (%) in a sunflower recombinant inbred line (RIL) population ( $n_G = 146$ ) with nearly balanced marker data and multiple phenotypic observations (replications) per RIL [53]. The marker loci were identified by QTL mapping. Variance components were estimated from LMM (27) for the AMV method and LMM (S13) for the ASV method.

<sup>c</sup>Statistics are shown for two SNP markers (*AX396* and *AX493*) associated with genetic variation for resistance to Fusarium wilt in a strawberry population ( $n_G = 565$ ) with unbalanced SNP marker data and multiple phenotypic observations per individual [86]. *AX396* and *AX493* are tightly linked and both were in LD with a dominant gene (*FW1*) conferring resistance to Fusarium wilt but had significantly different genotypic ratios among individuals in the population. Variance components were estimated from LMM (2) for the AMV method and LMM (15) for the ASV method. The  $k_M$  coefficient a single locus with unbalanced data are shown in S1 Text.

<sup>d</sup>Type II  $R^2$  is the coefficient of partial determination estimated from a Type II ANOVA, where the main and interactions effects of markers are fixed. For the cattle example, the reduction in sums of squares for main effects were estimated with the other main effects in the genetic model without interactions, e.g., the reduction in SS for *rs10* was  $R(rs10|rs45, rs20)$ . Similarly, the reduction in SS for each two-locus interaction was estimated without main or three-way interaction effects in the genetic model, e.g., the Type II reduction in sum of squares for the *rs10* × *rs45* interaction was  $R(rs10 \times rs45|rs45, rs20, rs10 \times rs20, rs45 \times rs20)$  and so on for the other two-locus interactions. Finally, the reduction in SS for the three-locus interaction was  $R(rs10 \times rs45 \times rs20|rs10, rs45, rs20, rs10 \times rs45, rs10 \times rs20, rs45 \times rs20)$ .

<sup>e</sup>Type III  $R^2$  is the coefficient of partial determination estimated from a Type III ANOVA, where the main and interactions effects of markers are fixed, e.g., the reduction in sums of squares for *rs10* in the cattle example was estimated by fitting *rs10* with all other factors in the model:  $R(rs10|rs45, rs20, rs10 \times rs45, rs10 \times rs20, rs45 \times rs20, rs10 \times rs45 \times rs20)$ .

<https://doi.org/10.1371/journal.pgen.1009762.t001>

where  $MS_G = SS_G/df_G$  is the between entry mean square,  $SS_G$  is the between entry sum of squares,  $df_G = n_G - 1$  is the between entry degrees of freedom,  $MS_\epsilon = SS_\epsilon/df_\epsilon = \sigma_\epsilon^2$  is the residual mean square,  $SS_\epsilon$  is the residual sum of squares,  $df_\epsilon = n_G(r_G - 1) - 1$  is the residual degrees of freedom,  $\sigma_\epsilon^2$  is the residual variance component, and  $r_G$  is the number of replications per entry [74]. The between-entry variance component has a theoretical genetic interpretation when entries are progeny with genetic relationships known from pedigrees, e.g., monozygotic twins, full-sib families, or recombinant inbred lines [24, 25, 30]. ANOVA estimators of the marker locus  $M$  and entry nested in  $M$  variance components from LMM (2) with balanced data are:

$$\hat{\sigma}_M^2 = \frac{MS_M - MS_{G:M}}{r_G n_{G:M}} \tag{4}$$

and

$$\hat{\sigma}_{G:M}^2 = \frac{MS_{G:M} - MS_\epsilon}{r_G} \tag{5}$$

respectively, where  $n_{G:M}$  is the number of entries nested in each marker genotype,  $E(\hat{\sigma}_M^2) = \sigma_M^2$ ,  $E(\hat{\sigma}_{G:M}^2) = \sigma_{G:M}^2$ ,  $MS_{G:M}$  is the entry nested in  $M$  mean square, and  $MS_M$  is the mean square for marker locus  $M$ . The residuals in LMMs (1) and (2) are identical when the data are balanced ( $\hat{\sigma}_\epsilon^2 = MS_\epsilon$ ). Hence, for a single marker locus with balanced data, the ANOVA estimator of  $p$  is:

$$\hat{p} = \frac{\hat{\sigma}_M^2}{\hat{\sigma}_G^2} \tag{6}$$

and the ANOVA estimator of broad-sense marker heritability on an entry-mean basis is:

$$\hat{H}_M^2 = \frac{\hat{\sigma}_M^2}{\hat{\sigma}_G^2 + \hat{\sigma}_\epsilon^2/r_G} \tag{7}$$

where  $\hat{\sigma}_p^2 = \hat{\sigma}_G^2 + \hat{\sigma}_\epsilon^2/r_G$  is the phenotypic variance on an entry-mean basis [25, 76].

The overestimation of  $p$  and  $H_M^2$  was not obvious from inspection of ANOVA estimators (6) and (7). The source of the bias was discovered by substituting  $SS_M + SS_{G:M}$  for  $SS_G$  in the ANOVA estimator of  $\sigma_G^2$  from (3) and simplifying:

$$\begin{aligned} \hat{\sigma}_G^2 &= \frac{1}{df_G r_G} (SS_M + SS_{G:M}) - \frac{1}{df_\epsilon r_G} SS_\epsilon \\ &= \frac{1}{df_G r_G} [df_M (\hat{\sigma}_\epsilon^2 + r_{G:M} \hat{\sigma}_{G:M}^2 + r_M \hat{\sigma}_M^2) + df_{G:M} (\hat{\sigma}_\epsilon^2 + r_{G:M} \hat{\sigma}_{G:M}^2)] - \frac{1}{df_\epsilon r_G} df_\epsilon \hat{\sigma}_\epsilon^2 \\ &= \frac{df_M r_M}{df_G r_G} \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 = \frac{df_M n_{G:M}}{df_G} \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 = k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 \end{aligned} \tag{8}$$

where the fraction  $k_M$  is source of the bias,  $0 < k_M < 1$ ,  $r_M$  is the number of replications per marker genotype,  $n_{G:M}$  is the number of entries nested in marker loci,  $SS_M$  is the marker sum of squares,  $df_M$  is the marker degrees of freedom,  $r_M$  is the number of replicates of each marker genotype,  $SS_{G:M}$  is the entry nested in marker sum of squares, and  $df_{G:M}$  is the entry nested in marker degrees of freedom. The term  $k_M$  in (8) depends on degrees of freedom and  $n_{G:M}$  and is hereafter referred to as the  $k_M$  bias coefficient, where the subscript  $M$  indexes the intralocus and interlocus effects of marker loci.



Eq (8) shows that the sum of ANOVA estimates of  $\sigma_M^2$  and  $\sigma_{G:M}^2$  from LMM (1) are greater than the ANOVA estimate of  $\sigma_G^2$  from LMM (2):

$$\hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 > \hat{\sigma}_G^2 = k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 \tag{9}$$

Although the SS for sources of variation in these LMMs are additive ( $SS_M + SS_{G:M} = SS_G$ ), the mean squares are not ( $MS_M + MS_{G:M} \neq MS_G$ ). Because  $\hat{\sigma}_G^2 = k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2$ , the sum  $\hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2$  from LMM (2) overestimates  $\hat{\sigma}_G^2$  by a factor of  $(1 - k_M) \hat{\sigma}_M^2$ . The ANOVA estimators of  $p$  and  $H_M^2$  from analyses of LMMs (1) and (2) are upwardly biased because  $\hat{\sigma}_M^2$  is multiplied by the fraction  $k_M$  in their denominators, and not the numerators:

$$\hat{p} = \frac{\hat{\sigma}_M^2}{\hat{\sigma}_G^2} = \frac{\hat{\sigma}_M^2}{k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2} \tag{10}$$

and

$$\hat{H}_M^2 = \frac{\hat{\sigma}_M^2}{\hat{\sigma}_p^2} = \frac{\hat{\sigma}_M^2}{k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 + \hat{\sigma}_\epsilon^2 / r_G} \tag{11}$$

Substituting  $\hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2$  for  $\hat{\sigma}_G^2$  in the denominators of  $p$  and  $H_M^2$  decreases but does not eliminate the bias because  $\hat{\sigma}_M^2$  is multiplied by  $k_M$  in the denominator (S1 Fig). For a single marker with balanced data, we found that:

$$k_M = \frac{df_M r_M}{df_G r_G} = \frac{df_M n_{G:M}}{df_G} \tag{12}$$

and

$$\left( \frac{\hat{\sigma}_M^2}{\hat{\sigma}_G^2} = \frac{\hat{\sigma}_M^2}{k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2} \right) > \frac{\hat{\sigma}_M^2}{\hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2} > \left( \frac{k_M \hat{\sigma}_M^2}{k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2} = \frac{k_M \hat{\sigma}_M^2}{\hat{\sigma}_G^2} \right) \tag{13}$$

where  $0 < k_M < 1$ . Hence, the bias is caused by the  $k_M$  multiplier in the expected values of the ANOVA estimators of  $p$  and  $H_M^2$ . As shown later, simulation analyses confirmed that (9) and (13) accurately predict the upward bias caused by  $k_M$ . Moreover, we concluded that the bias could be corrected by multiplying ANOVA or REML estimates of  $\sigma_M^2$  by  $k_M$  in the numerators of  $p$  and  $H_M^2$  estimates.

### Genetic models with unbalanced genotypic data

We started with the special case of balanced data, which seldom arises in practice, but develop results here for the general case of unbalanced data. Following the same approach as that shown above for a single locus with balanced data, we found  $k_M$  coefficients for bias-correcting ANOVA and REML estimates of  $p$  and  $H_M^2$  for analyses of one to three marker loci with unbalanced genotypic data (S1, S2 and S3 Texts). For a single marker locus with unbalanced genotypic data, we found:

$$k_M = \frac{n_G - n_G^{-1} \sum_h n_{G:M_h}^2}{df_G} \tag{14}$$

where  $n_G$  is the number of entries,  $df_G$  are the degrees of freedom for entries, and  $n_{G:M_h}$  is the number of entries nested in the  $h^{th}$  marker genotype (S1 Text). This simplifies to (12) for a single marker locus with balanced genotypic data.

The  $k_M$  coefficients become slightly more complicated as the number of marker loci increases but nevertheless follow a predictable algebraic pattern, e.g., for a two locus genetic

model, see equations (S10)-(S12) in [S2 Text](#). Similarly, for a three locus genetic model, see equations (S19)-(S25) in [S3 Text](#).  $k_M$  is greater ( $k_M$  bias is proportionally smaller) for interaction than main effects, e.g., for two marker loci,  $k_{M1} < k_{M1 \times M2} < 1$  and  $k_{M2} < k_{M1 \times M2} < 1$ , where  $k_{M1}$  is the coefficient for  $M_1$ ,  $k_{M2}$  is the coefficient for  $M_2$ , and  $k_{M1 \times M2}$  is the coefficient for the  $M_1 \times M_2$  epistatic interaction ([S2 Text](#)).  $k_M$  for the two-locus interaction ( $k_{M1 \times M2}$ ) is larger than  $k_M$  for the individual marker loci ( $k_{M1}$  and  $k_{M2}$ ) because the denominator ( $df_G r_G$ ) is constant, whereas the numerators increase and approach the denominator as the degrees of freedom for marker effects increase. Therefore, the upward bias is proportionally smaller for the  $M1 \times M2$  variance component than the  $M_1$  or  $M_2$  variance components for a two locus genetic model. Similarly, for a three locus genetic model, the upward bias is proportionally smaller for the  $M_1 \times M_2 \times M_3$  interaction variance component than the two-way interaction variance components ( $M_1 \times M_2$ ,  $M_1 \times M_3$ , and  $M_2 \times M_3$ ). These results naturally extend to genetic models with more than three loci. Algebraic results are only shown for three marker loci because we found that the  $k_M$  bias problem can be directly solved using average semivariance estimation methods when analyzing more complex genetic models (see below). Although certainly not limited to three marker loci, the methods described herein are primarily designed to study the effects of one to a few genes with large effects, e.g., *BRCA2* [52], *BTA19* [43], and the examples shown in [Tables 1](#) and [2](#), and not to replace GWAS or QTL mapping.

### Study designs without replications or repeated measures of individuals or families

LMMs (1) and (2) arise in study designs where entries (individuals, families, or strains) are replicated, e.g., in studies with domesticated plants, biological replicates of half-sib or full-sib families, doubled haploid or recombinant inbred lines, or testcross hybrids are commonly phenotyped [24, 25, 31, 76, 87, 88] (see the sunflower example in [Table 1](#)). These same LMMs apply to study designs for monozygotic twins in humans and other mammals and clonally replicated individuals in asexually propagated plants, e.g., cassava (*Manihot esculenta*), strawberry (*Fragaria × ananassa*), and apple (*Malus × domestica*) (see the strawberry examples in [Table 1](#)). The extension of the proposed  $k_M$  bias correction solutions to LMMs with repeated measures is straightforward and should have applications in studies where large effect loci are important determinants of the genetic variation underlying quantitative traits in both replicable and unreplicable organisms or populations [88–94].

When entries are unreplicated, the random error or residual source of variation in LMM (2) disappears ( $\sigma_{G:M}^2$  becomes the residual) and  $\sigma_G^2$ ,  $\sigma_e^2$ , and  $p$  cannot be estimated; however, the marker heritability can be estimated using the phenotypic variance among unreplicated individuals ( $k_M \sigma_M^2 + \sigma_{G:M}^2$ ). As before, this variance component ratio is upwardly biased by the factor  $k_M$  (see the cattle example in [Table 1](#)). Without the insights gained from the algebra shown in equations (10), (S3), (S9), and (S18), and [S1](#), [S2](#) and [S3](#) Texts, the bias would not be obvious unless one or more estimates of marker heritability exceeded 1, which only happens when the loci under study have very large effects. That was exactly how we originally discovered the bias problem in the first place ([Table 1](#)). The bias is systematic and ubiquitous but not immediately obvious when estimates fall within the expected range ( $0 < \hat{H}_M^2 < 1$ ). The same bias correction solutions we proposed for study designs with replications of entries can be applied in study designs where entries are unreplicated. When unreplicated entries are genotyped with a dense genome-wide of markers,  $\sigma_G^2$  be estimated using a genomic or pedigree relationship matrix [92, 95–97], which yields an estimate of  $p$ .

**Table 2. Type I, II, and III sums of squares for fixed effect analyses of markers associated with QTL identified in GWAS and QTL mapping experiments in cattle and sunflower.**

Study	Source	Type I SS <sup>a</sup>						Type II SS	Type III SS
		ABC	ACB	BAC	BCA	CAB	CBA		
Cattle White Spotting <sup>b</sup>	<i>rs10</i>	3,552.3	3,552.3	1,707.2	591.4	1,208.7	591.4	542.5	22.1
	<i>rs45</i>	6,539.7	4,259.6	8,384.8	8,384.8	4,259.6	4,876.8	4,282.7	1,394.9
	<i>rs20</i>	4,880.5	7,160.7	4,880.5	5,996.4	9,504.4	9,504.4	4,834.3	1,788.4
	<i>rs10</i> × <i>rs45</i>	12.7	12.7	12.7	12.7	12.7	12.7	14.3	47.4
	<i>rs10</i> × <i>rs20</i>	132.7	132.7	132.7	132.7	132.7	132.7	107.4	234.0
	<i>rs45</i> × <i>rs20</i>	193.1	193.1	193.1	193.1	193.1	193.1	193.1	91.5
	<i>rs10</i> × <i>rs45</i> × <i>rs20</i>	143.5	143.5	143.5	143.5	143.5	143.5	143.5	143.5
	<i>G</i> : <i>rs10</i> × <i>rs45</i> × <i>rs20</i>	15,512.9	15,512.9	15,512.9	15,512.9	15,512.9	15,512.9	15,512.9	15,512.9
Sunflower Oil Content <sup>c</sup>	<i>BR</i>	1,624.0	1,624.0	1,708.8	1,904.0	1,829.7	1,904.0	1,881.4	1,711.2
	<i>PHY</i>	298.2	254.2	213.4	213.4	254.3	180.0	220.2	208.3
	<i>HYP</i>	537.1	581.0	537.1	342.0	375.4	375.4	507.0	511.6
	<i>BR</i> × <i>PHY</i>	57.9	57.9	57.9	57.9	57.9	57.9	49.7	50.0
	<i>BR</i> × <i>HYP</i>	168.0	168.0	168.0	168.0	168.0	168.0	172.1	195.5
	<i>PHY</i> × <i>HYP</i>	11.1	11.1	11.1	11.1	11.1	11.1	11.1	7.6
	<i>BR</i> × <i>PHY</i> × <i>HYP</i>	36.6	36.6	36.6	36.6	36.6	36.6	36.6	36.6
	<i>G</i> : <i>BR</i> × <i>PHY</i> × <i>HYP</i>	4,113.4	4,113.4	4,113.4	4,113.4	4,113.4	4,113.4	4,113.4	4,113.4
	Residual	553.8	553.8	553.8	553.8	553.8	553.8	553.8	553.8

<sup>a</sup>For each Type I ANOVA, the six possible orders of the three main effects (marker loci *A*, *B*, and *C*) were tested in the genetic model, where *A* = *rs10*, *B* = *rs45*, and *C* = *rs20* for the cattle example and *A* = *BR*, *B* = *PHY*, and *C* = *HYP* for the sunflower example. The interactions were added to the genetic model in a single sequence: *A* × *B*, *A* × *C*, *B* × *C*, and *A* × *B* × *C*. The three letters indicate the sequence with which markers loci entered the genetic model, e.g., for the ABC order, the sums of squares for the three main effects were  $SS(A|\mu)$ ,  $SS(B|A, \mu)$ , and  $SS(C|A, B, \mu)$ , where  $\mu$  is the population mean and factors to the right of the vertical bar were included in the model. Similarly, for the CBA order, the sums of squares for the three main effects were  $SS(C|\mu)$ ,  $SS(B|C, \mu)$ , and  $SS(A|B, C, \mu)$ . The sequences with which interactions were added to the genetic model were identical in the six Type I analyses, e.g., the sums of squares for the *A* × *B* interaction was  $SS(A \times B|A, B, C, \mu)$  and for the three-way interaction was  $SS(A \times B \times C|A, B, C, A \times B, A \times C, B \times C, \mu)$ .

<sup>b</sup>Statistics are shown for three marker loci (*rs10*, *rs45*, and *rs20*) associated with genetic variation for white spotting (%) in a cattle population ( $n_G = 2,973$ ) with a single phenotypic observation per individual and highly unbalanced marker data [85]. The markers were identified by GWAS. The linear model for the cattle analysis was identical to the linear model for the sunflower analysis without replications ( $r_G = 1$ ); hence, the residual in the cattle analysis was the entry nested in marker source of variation.  $k_M$  coefficients for three loci with unbalanced data are shown in S3 Text.

<sup>c</sup>Statistics are shown for three marker loci (*BR*, *PHY*, and *HYP*) associated with genetic variation for seed oil content (%) in a sunflower recombinant inbred line (RIL) population ( $n_G = 146$ ) with nearly balanced marker data and multiple phenotypic observations (replications) per RIL [53].  $k_M$  coefficients for three loci with unbalanced data are shown in S3 Text.

<https://doi.org/10.1371/journal.pgen.1009762.t002>

## Average semivariance estimation directly solves the bias problem

The AMV methods proposed above for bias correcting ANOVA or REML estimates of  $p$  and  $H_M^2$  are straightforward to apply in practice because they are the methods widely described in textbooks and implemented in popular statistical software packages, e.g., the R package ‘lme4’ and SAS package ‘GLIMMIX’ [78, 98]. Here we show that the bias problem can be directly solved by applying average semivariance (ASV) estimation methods [73]. As before, we start by showing results for a single marker locus with balanced genotypic data. AMV notation and estimators are reformulated in matrix notation here to build the foundation for describing ASV notation and estimators. The input for both are the adjusted entry-level means ( $\bar{y}_{ij\bullet}$ ) from LMM (1) stored in an  $n_G$ -element vector. These are the best linear unbiased estimates (BLUEs) for entries [73, 99]. The LMM equivalent to (2) for the entry-level means analysis of the effect

of a single marker locus ( $M$ ) is:

$$\bar{y}_{ij\bullet} = \mu + M_i + G : M_{i(j)} + \bar{\epsilon}_{ij\bullet} \tag{15}$$

where  $\bar{y}_{ij\bullet}$  is the phenotypic mean for the  $ij^{th}$  entry,  $\mu$  is the population mean,  $M_i$  is the random effect of the  $i^{th}$  marker genotype,  $var(M_i) = \sigma_M^2$ ,  $G : M_{i(j)}$  is the random effect of entries nested in  $M$ ,  $var(G : M_{i(j)}) = \sigma_{G:M}^2$ ,  $\bar{\epsilon}_{ij\bullet}$  is the residual error, and  $var(\bar{\epsilon}_{ij\bullet}) = r_G^{-1}\sigma_\epsilon^2$ . The residual variance-covariance matrix ( $R$ ) is estimated in the first stage of a two-stage analysis [99–101]. The between-entry variance can be partitioned into  $\sigma_M^2$  and  $\sigma_{G:M}^2$  with individual variance-covariance matrices  $G_c$  defined by the genetic model, e.g., different main and interaction effects among marker loci.

The AMV estimator of the phenotypic (total) variance among observations for LMM (15) is:

$$\hat{\theta}_P^{AMV} = n_G^{-1}tr(V) = \sum_c \hat{\theta}_{g_c}^{AMV} + \hat{\theta}_\epsilon^{AMV} \tag{16}$$

where  $V$  is the variance-covariance matrix of the phenotypic observations,  $n_G$  is the number of entries,  $tr(V)$  is the trace of  $V$ ,  $\theta_{g_c}^{AMV} = n_G^{-1}tr(Z_c G_c Z_c^T)$  is the marginal variance explained by the  $c^{th}$  genetic factor in the model (e.g.,  $M$  and  $G : M$ ),  $Z_c$  are design matrices for the  $c$  genetic factors,  $\theta_\epsilon^{AMV} = n_G^{-1}tr(R)$  is the AMV estimator of the residual variance, and  $R$  is the residual variance-covariance matrix. The AMV estimator of the genetic variance among entries ( $G$ ) is:

$$\hat{\theta}_G^{AMV} = (n_G)^{-1} \hat{\sigma}_G^2 tr(Z_{u_G} Z_{u_G}^T) = \hat{\sigma}_G^2 \tag{17}$$

where  $Z_{u_G}$  is a  $n_G$  identity matrix. From LMM (15), the AMV estimator of the variance associated with a single marker locus with balanced data is:

$$\hat{\theta}_M^{AMV} = (n_G)^{-1} \hat{\sigma}_M^2 tr(Z_{u_M} Z_{u_M}^T) = \frac{(n_M)n_{G:M}}{n_G} \hat{\sigma}_M^2 = \hat{\sigma}_M^2 \tag{18}$$

where  $E(\hat{\sigma}_M^2) = \sigma_M^2 = \theta_M^{AMV}$ ,  $Z_{u_M} = I_{n_M} \otimes 1_{n_{G:M}}$ ,  $I_{n_M}$  is a  $n_M$  identity matrix,  $1_{n_{G:M}}$  is an  $n_{G:M}$ -element unit vector, and  $u_M$  is a vector of random effects for  $M$ . The AMV estimator of the variance associated with the residual genetic variation among entries nested in  $M$  is:

$$\hat{\theta}_{G:M}^{AMV} = (n_G)^{-1} \hat{\sigma}_{G:M}^2 tr(Z_{u_{G:M}} Z_{u_{G:M}}^T) = \frac{n_G}{n_G} \hat{\sigma}_{G:M}^2 = \hat{\sigma}_{G:M}^2 \tag{19}$$

where  $u_{G:M}$  is a vector of random entry nested in  $M$  effects and  $Z_{u_{G:M}}$  is a  $n_G$  identity matrix. Hence, the AMV estimators of  $\sigma_M^2$  and  $\sigma_G^2$  are identical to ANOVA estimators (4) and (5), respectively, with entry means as input for the former and original observations as input for the latter.

ASV, or the average variance of differences among observations, leads to a definition of the total variance that provides a natural way to account for the heterogeneity of variance and covariance among observations [73, 102]. ASV can be defined for any variance-covariance structure in a generalized LMM and allows for missing and unbalanced data [73]. The ASV estimator of total variance is half the average variance of pairwise differences among entries and can be partitioned into independent sources of variance, e.g., genetic and non-genetic or residual:

$$\hat{\theta}_P^{ASV} = (n_G - 1)^{-1} tr(VD_{n_G}) = \sum_c \hat{\theta}_{g_c}^{ASV} + \hat{\theta}_\epsilon^{ASV} \tag{20}$$

where  $D_{n_G} = I_{n_G} - n_G^{-1}J_{n_G}$  is the idempotent matrix used for column-wise mean-centering,  $I_{n_G}$  is an  $n_G \times n_G$  identity matrix, and  $J_{n_G}$  is an  $n_G \times n_G$  unit matrix [73].  $\theta_p^{ASV}$  accounts for the variance and covariance of the phenotypic observations. From (20),  $\theta_{g_c}^{ASV} = (n_G - 1)^{-1}tr(Z_c G_c Z_c^T D_{n_G})$  is the variance explained by the  $c^{th}$  genetic factor ( $u_c$ ), where  $c$  indexes genetic factors, the genetic factors are marker locus effects and entries nested in marker locus effects, and  $\theta_{\epsilon}^{ASV} = (n_G - 1)^{-1}tr(RD_{n_G})$  is the residual variance. The variance explained by the  $c^{th}$  genetic factor is  $\theta_{g_c}^{ASV} = E(s_{g_c}^2)$ , e.g., for a single marker locus  $M$ ,  $\theta_M^{ASV} = E(s_M^2)$ ,  $E(s_{g_c}^2)$ ,  $E(s_M^2)$ , and the biases of these ASV estimators are defined in S4 Text.

The ASV estimator of the genetic variance among entries ( $G$ ) is:

$$\hat{\theta}_G^{ASV} = (n_G - 1)^{-1} \hat{\sigma}_G^2 tr(Z_{u_G} Z_{u_G}^T D_{n_G}) = \hat{\sigma}_G^2 \tag{21}$$

where  $Z_{u_G}$  is a  $n_G$  identity matrix. Hence, from Eqs (8), (17) and (21), AMV and ASV estimators of the between-entry variance component ( $\sigma_G^2$ ) are equivalent ( $\hat{\sigma}_G^2 = \hat{\theta}_G^{AMV} = \hat{\theta}_G^{ASV}$ ). The ASV estimator of the variance associated with  $M$  is:

$$\hat{\theta}_M^{ASV} = (n_G - 1)^{-1} \hat{\sigma}_M^2 tr(Z_{u_M} Z_{u_M}^T D_{n_G}) = \frac{(n_M - 1)n_{G:M}}{n_G - 1} \hat{\sigma}_M^2 = \frac{df_M n_{G:M}}{df_G} \hat{\sigma}_M^2 = k_M \hat{\sigma}_M^2 \tag{22}$$

where  $k_M = df_M n_{G:M} / df_G$  is the bias correction coefficient,  $Z_{u_M} = I_{u_M} \otimes 1_{n_{G:M}}$ ,  $df_G = n_G - 1$ ,  $df_M = n_M - 1$ , and  $df_{G:M} = df_G - df_M$ . This definition of the  $k_M$ -bias coefficient is identical to the earlier definition with  $r_G$  factored out (see Eq 12). Eq (22) shows that the ASV estimator of  $\sigma_M^2$  is corrected by the fraction  $k_M$ , which correctly scales the estimate of  $\sigma_M^2$  to the genetic variance and yields unbiased estimates of  $p$  and  $H_M^2$ . From Eqs (9) and (22), we found that  $\hat{\theta}_M^{ASV} < \hat{\theta}_M^{AMV}$  by the factor  $k_M$ . The ASV estimator of the variance associated with  $G : M$  is:

$$\hat{\theta}_{G:M}^{ASV} = (n_G - 1)^{-1} \hat{\sigma}_{G:M}^2 tr(Z_{u_{G:M}} Z_{u_{G:M}}^T D_{n_G}) = \frac{n_G - 1}{n_G - 1} \hat{\sigma}_{G:M}^2 = \hat{\sigma}_{G:M}^2 \tag{23}$$

The ASV estimator of  $p$  for a single marker locus ( $M$ ) is:

$$\hat{p}_* = \frac{\hat{\theta}_M^{ASV}}{\hat{\theta}_G^{ASV}} = \frac{k_M \hat{\sigma}_M^2}{\hat{\sigma}_G^2} = \frac{k_M \hat{\sigma}_M^2}{k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2} \tag{24}$$

Similarly, the ASV estimator of  $H_M^2$  for a single marker locus is:

$$\hat{H}_{M*}^2 = \frac{\hat{\theta}_M^{ASV}}{\hat{\theta}_p^{ASV}} = \frac{k_M \hat{\sigma}_M^2}{\hat{\sigma}_G^2 + r_G^{-1} \hat{\sigma}_{\epsilon}^2} = \frac{k_M \hat{\sigma}_M^2}{k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 + r_G^{-1} \hat{\sigma}_{\epsilon}^2} \tag{25}$$

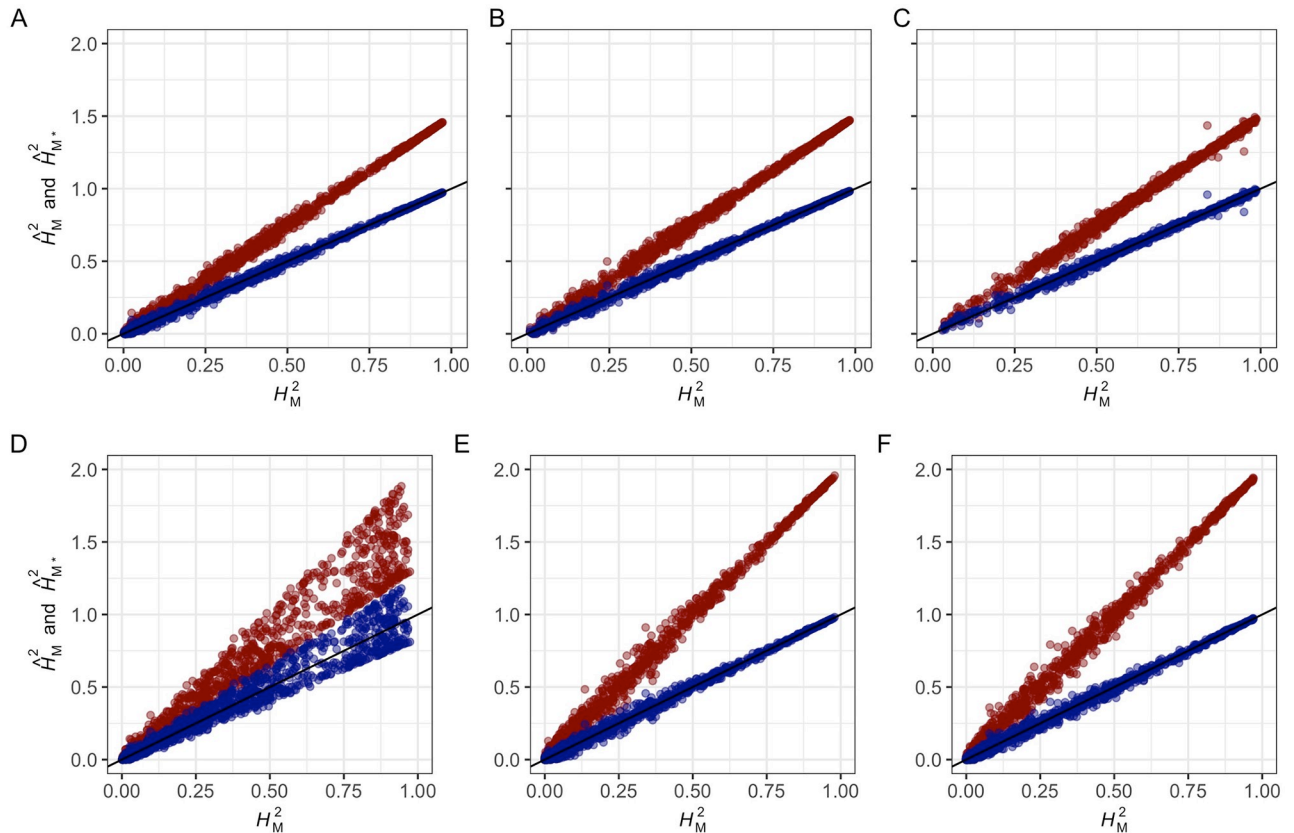
where  $\hat{\sigma}_G^2 + \hat{\sigma}_{\epsilon}^2 / r_G = \hat{\sigma}_p^2$  is the phenotypic variance on an entry-mean basis [25]. From these results, we found that:

$$\hat{\theta}_M^{ASV} + \hat{\theta}_{G:M}^{ASV} = \hat{\theta}_G^{ASV} = \hat{\theta}_G^{AMV} = \hat{\sigma}_G^2 < \hat{\theta}_M^{AMV} + \hat{\theta}_{G:M}^{AMV} \tag{26}$$

and showed that ASV estimators of  $p$  and  $H_M^2$  are unbiased (automatically corrected for  $k_M$ ).

### Computer simulations confirmed that ASV-REML estimates of $p$ and $H_M^2$ are unbiased

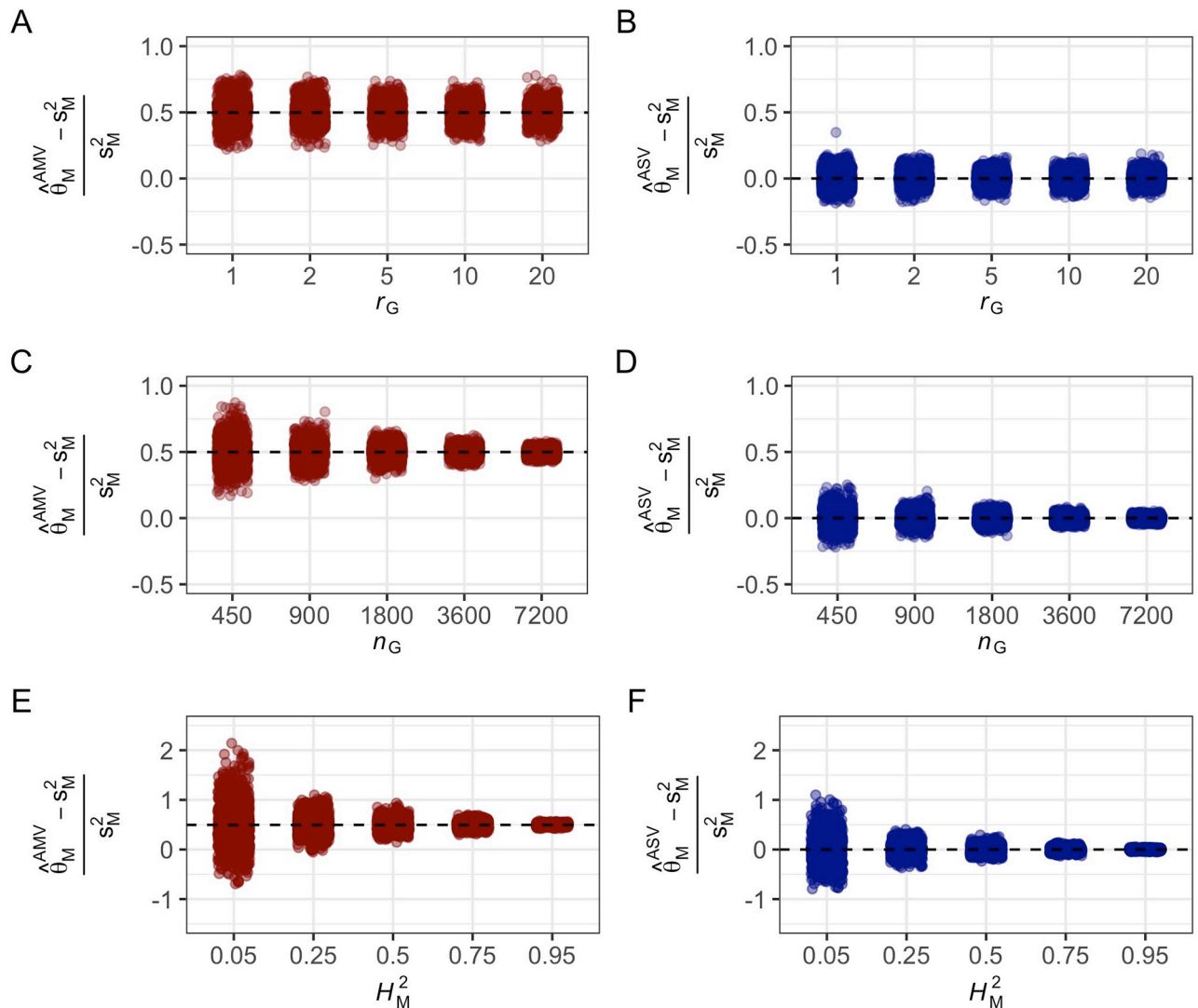
Computer simulations confirmed that AMV-REML estimates of  $p$  (6) and  $H_M^2$  (7) are upwardly biased by the factor  $k_M$  and that ASV-REML estimates of these parameters form (24) and (25)



**Fig 1. Accuracy of AMV and ASV estimators of marker heritability.** AMV and ASV estimates of  $H_M^2$  are shown for 1,000 segregating populations simulated for different numbers of entries ( $n_G$  individuals, families, or strains), five replications/entry ( $r_G = 5$ ), true marker heritability ( $H_M^2$ ) ranging from 0 to 1, and one to three marker loci with three genotypes/marker locus ( $n_{M1} = 3$ ). AMV estimates of marker heritability ( $\hat{H}_M^2$ ; red highlighted observations) and ASV estimates of marker heritability ( $\hat{H}_M^2$ ; blue highlighted observations) are shown for: (A) one locus with balanced data for  $n_G = 540$  entries (study design 1); (B) two marker loci with interaction ( $M1$ ,  $M2$ , and  $M1 \times M2$ ) and balanced data for  $n_G = 540$  (study design 2); (C) three marker loci with interactions ( $M1$ ,  $M2$ ,  $M3$ ,  $M1 \times M2$ ,  $M1 \times M3$ ,  $M2 \times M3$ , and  $M1 \times M2 \times M3$ ) and balanced data for  $n_G = 540$  (study design 3); (D) an population segregating 1:2:1 for one marker locus with  $r_{G:M} = 135$  entries for both homozygotes and  $r_{G:M} = 270$  heterozygous entries, and  $n_G = 540$  (study design 4); (E) one locus with 10% randomly missing data among 540 entries (study design 5); and (F) one locus with 33% randomly missing data among 540 entries (study design 6). Study design details are shown in [S1 Table](#).

<https://doi.org/10.1371/journal.pgen.1009762.g001>

are unbiased (Figs 1 and 2). The mean of AMV-REML estimates of  $p$  and  $H_M^2$  from 21 different simulation study designs ([S1 Table](#)) were identical to those predicted by the  $k_M$  coefficients shown in [S1](#), [S2](#) and [S3](#) Texts. Several insights arose from the simulation analyses. First, the bias caused by  $k_M$  increased as  $H_M^2$  increased but was proportionally constant for different  $H_M^2$  (Fig 1). These results show that the overestimation of  $p$  and  $H_M^2$  is greatest for genes and gene-gene interactions with large effects (Fig 1). Their effects could be inflated by selection bias over and above  $k_M$  bias [[67](#), [68](#), [82](#), [83](#), [103](#)]; hence, we concluded that  $k_M$ -bias and selection bias could operate in combination to inflate estimates of the contribution of a locus to the heritable variation in a population ([S1](#), [S2](#) and [S3](#) Texts). Moreover, because the bias increases as the effect of the locus increases, we concluded that the overestimation problem is worst for large-effect QTL (Fig 1). Second,  $k_M$  bias was greater for unbalanced than balanced data (Fig 1D and 1E). The effect of unbalanced data was more extreme for the  $F_2$  simulation (Fig 1D) where the expected genotypic ratio was 1 AA: 2 Aa: 1 aa than for simulations where 10 or 33% of the observations were randomly missing for markers with roughly equal numbers of replicates/



**Fig 2. Effect of  $r_G$ ,  $n_G$ , and  $H_M^2$  on the relative bias of AMV and ASV estimators of  $\sigma_M^2$ .** (A and B) Phenotypic observations were simulated for 1,000 populations segregating for a single marker locus with three genotypes ( $n_M = 3$ ),  $n_G = 900$  progeny, and  $r_G = 1, 2, 5, 10$ , or 20 (study designs 7–11). The marker locus was assumed to be in complete linkage disequilibrium with a single QTL that explains 50% of the phenotypic variance ( $H_M^2 = 0.50$ ). (A) Distribution of the relative biases of AMV estimates of  $\sigma_M^2$  for different  $r_G$ . The relative bias  $RB[\theta_M^{AMV}] = 0.498$  was identical for different  $r_G$ . (B) Distribution of the relative biases of ASV estimates of  $\sigma_M^2$  for different  $r_G$ . The relative bias  $RB[\theta_M^{ASV}] = 0.00$  was identical for different  $r_G$ . (C and D) Phenotypic observations were simulated for 1,000 populations segregating for a single marker locus with three genotypes ( $n_M = 3$ ), five replications/entry ( $r_G = 5$ ), and  $n_G = 450, 900, 1,800, 3,600$ , or 7,200 entries/population (study designs 12–16). The marker locus was assumed to be in complete linkage disequilibrium with a single QTL that explains 50% of the phenotypic variance ( $H_M^2 = 0.50$ ). (C) Distribution of the relative biases of AMV estimates of  $\sigma_M^2$  for different  $n_G$ . The relative bias  $RB[\theta_M^{AMV}] = 0.499$  was identical across the variables tested. (D) Distribution of the relative biases of ASV estimates of  $\sigma_M^2$  for different  $n_G$ . The relative bias ( $RB[\theta_M^{ASV}] = 0.00$ ) was identical across the variables tested. (E and F) Phenotypic observations were simulated for 1,000 populations segregating for a single marker locus with three genotypes ( $n_M = 3$ ), five replications/entry ( $r_G = 5$ ), and  $n_G = 450$  entries/population. The marker locus was assumed to be in complete linkage disequilibrium with a single QTL that explains 5–95% of the phenotypic variance ( $H_M^2 = 0.05$  to 0.95 (study designs 17–21)). (E) Distribution of the relative biases of AMV estimates of  $\sigma_M^2$  for different  $H_M^2$ . The relative bias  $RB[\theta_M^{AMV}] = 0.496$  was identical across the variables tested. (F) Distribution of the relative biases of ASV estimates of  $\sigma_M^2$  for different  $H_M^2$ . The relative bias  $RB[\theta_M^{ASV}] = 0.0$  was identical across the variables tested.

<https://doi.org/10.1371/journal.pgen.1009762.g002>

marker genotype (Fig 1E and 1F). Third, the F<sub>2</sub> and missing data simulations further showed that the precision of estimates of these parameters decreased as the genotypic data imbalance increased. Even though bias-corrected AMV and ASV estimates of these parameters are unbiased, the sampling variances among the simulated F<sub>2</sub> samples were larger than observed for

the 10 and 33% missing data samples and yielded a small percentage of  $H_M^2$  estimates slightly greater than 1.0 (Fig 1D). For the other simulation study designs (Fig 1), none of the ASV estimates exceeded 1.0. The sample variances of  $p$  and  $H_M^2$  can be estimated using data resampling methods, e.g., bootstrapping [104], or the estimators we developed using the Delta method (S5 Text) [25, 105, 106]. Equations (S44) and (S45) in S5 Text show that ASV estimates are more precise than AMV estimates by a factor of  $k_M^2$ . These predictions perfectly aligned with the empirical bootstrap estimates. Fourth, the relative biases were not affected by the number of replications of entries or the number of entries, although the precision of  $\sigma_M^2$  estimates increased as  $n_G$  and  $H_M^2$  increased (Fig 2). Predictably, the number of entries ( $n_G$ ) dramatically affected the precision of estimates of  $H_M^2$  (Fig 2C and 2D). The relative biases were not affected by  $r_G$  or  $H_M^2$ ; however, the sampling variances were strongly affected by  $H_M^2$  and decreased as  $H_M^2$  increased (Fig 2E and 2F and S2 Fig).

### GWAS example: A single marker locus with highly unbalanced genotypic data

The bias-correction methods described above are illustrated here for highly unbalanced genotypic data from a GWAS experiment. Variance components were estimated for two SNP markers (AX493 and AX396) in LD with a gene (*FW1*) conferring resistance to Fusarium wilt in a strawberry (*Fragaria* × *ananassa*) GWAS population ( $n_G = 564$ ) genotyped with a genome-wide framework of SNP markers [86]. Both SNP markers had highly significant GWAS effects with  $-\log_{10}(p) = 6.61 \times 10^{-31}$  for AX493 and  $2.95 \times 10^{-222}$  for AX396. Genotype frequencies were highly unbalanced for both markers with a scarcity of AA homozygotes (2.8%) for AX396 (16AA : 177Aa : 371aa) and a 1 : 2 : 1 ratio for AX493 (141AA : 282Aa : 141aa). For both loci, the minor allele frequency was  $>0.05$ . The  $k_M$  for these data ( $k_{AX493} = 0.62$  and  $k_{AX396} = 0.47$ ) were calculated as shown in S1 Text. The AMV-REML estimate of  $H_M^2$  for AX396 exceeded 1.0, a telltale sign of  $k_M$ -bias (Table 1). AMV-REML estimates of  $\sigma_M^2$  and  $H_M^2$  for both SNP markers were double or nearly double their bias-corrected ASV-REML estimates (Table 1). The bias-corrected estimate of marker heritability for AX396 was 0.62, versus 1.33 for the uncorrected estimate. Even with bias-correction, the sum of ASV-REML estimates of  $\sigma_M^2$  and  $\sigma_{G:M}^2$  for AX493 was slightly greater than the ASV-REML estimate of  $\sigma_G^2$ . This result was consistent with findings for highly unbalanced marker genotypic data in our simulation studies where a certain fraction of bias-corrected estimates exceeded the theoretical limit for heritability because of decreased precision (Fig 1). The  $k_M$ -bias problem would not necessarily have been detected in the analysis of AX396 because the  $p$  and  $H_M^2$  estimates fell within the expected range, e.g.,  $\hat{\theta}_{AX396}^{ASV} + \hat{\theta}_{G:AX396}^{ASV} / \hat{\theta}_P^{ASV} = 0.71$  (Table 1). Although both SNP markers were closely associated with *FW1*, they accounted for dramatically different fractions of genetic variance because of historic recombination and because neither are causal DNA variants or in complete LD with causal DNA variants [17, 19, 86, 107].

### QTL mapping example: Three marker loci with slightly unbalanced genotypic data

Statistics are shown here for an analysis of three marker loci (*BR*, *PHY*, and *HYP*) affecting seed oil content in a sunflower (*Helianthus annuus*) RIL population using LMM (27) [53]. The genotypic data were only slightly unbalanced and the three marker loci were identified by QTL mapping. The  $k_M$  needed for bias-correcting AMV-REML estimates of  $p$  and  $H_M^2$  are shown in S3 Text (Table 1). The AMV-REML estimates of  $p$  and  $H_M^2$  were nearly double the bias-corrected ASV-REML estimates, e.g., the AMV-REML estimate of  $H_M^2$  for the three-locus genetic



model (0.79) was nearly two-fold greater than the ASV-REML estimate (0.41) (Table 1). Similarly, the AMV-REML estimate of  $p$  for the *BR* locus (0.54) was slightly more than double the bias-corrected (ASV-REML) estimate (0.26). Hence, the uncorrected REML estimates of  $p$  and  $H_M^2$  grossly inflated the predicted contributions of the three marker loci to genetic variation for seed oil content (Table 1).

### GWAS example: Three marker loci with unbalanced genotypic data and unreplicated entries

The application of bias-correction is illustrated here for a genetic model with three marker loci, highly unbalanced genotypic data, and a single phenotypic observation per individual— $\sigma_G^2$  and  $p$  could not be estimated for this example because individuals were unreplicated. Variance components were estimated for three SNP markers (*rs10*, *rs45*, and *rs20*) on chromosomes 2, 6, and 22, respectively, affecting white spotting (%) in a Holstein–Friesian cattle (*Bos taurus*) population ( $n_G = 2,973$ ) [85]. These SNP markers had the largest effects among those predicted to be in LD with genes affecting white spotting. The genotypic frequencies were 50AA : 586Aa : 2, 337aa for *rs10*, 78AA : 736Aa : 2, 159aa for *rs45*, and 237AA : 976Aa : 1, 760aa for *rs20*. The  $k_M$  for these data ( $k_{rs10} = 0.35$ ,  $k_{rs45} = 0.41$ , and  $k_{rs20} = 0.54$ ) were calculated as shown in S3 Text. The uncorrected AMV-REML estimate of  $H_M^2$  for the three-locus genetic model (0.76) was substantially greater than the bias-corrected ASV-REML estimate (0.37) (Table 1). Similar differences were observed for the three marker loci.

### Candidate gene analysis: Fixed or random, BLUE or BLUP?

Our study was partly motivated by inconsistencies in the statistical approaches applied in candidate gene and other complex trait analyses when testing hypotheses and fitting genetic models for multiple large-effect loci. With the high densities of genome-wide markers commonly assayed in gene finding studies, investigators often identify markers tightly linked to candidate or known causal genes, as exemplified by diverse real world examples [17, 19, 33, 34, 37, 38, 40, 42, 43, 52, 54, 108]. The candidate marker loci are nearly always initially identified by genome-wide searches using sequential (marker-by-marker) approaches [56, 72, 75, 79, 109, 110]. Complicated and often misunderstood problems arise in the estimation and interpretation of statistics from sequential fixed effect analyses when the data are unbalanced [79, 111, 112]. Most importantly, there are multiple model fitting and analysis options (Type I, II, and III ANOVA) and the reduction in error sums of squares (SSE), test statistics, and parameter estimates differ among them, a problem that disappears when the data are balanced or when single large effect loci are discovered [79, 111–113]. Our review of the literature uncovered substantial variation and inconsistencies in the statistical approaches applied to the problem of fitting multilocus genetic models, testing multilocus genetic hypotheses, and calculating best linear unbiased estimates (BLUEs) from a fixed effects analysis of marker loci.

The problems that arise in fixed effect analyses of unbalanced data profoundly affect parameter estimates and statistical inferences but have not been universally recognized or addressed in complex trait analyses [79, 112]. We reanalyzed the cattle and sunflower examples with markers as fixed effects (Table 2) to show this, illustrate the challenges and nuances of fixed effects analyses of unbalanced data, and facilitate comparisons between random and fixed effects analyses of marker loci [56, 56, 75, 79, 109–112]. Following the discovery of statistically significant marker-trait associations from a marker-by-marker genome-wide scan, the natural progression would be to analyze multilocus genetic models where the effects of the discovered loci are simultaneously corrected for the effects of other discovered loci [79, 112], as shown in

our multilocus analysis examples (Tables 1 and 2). This is straightforward when the genotypic data are balanced or nearly balanced (as in the sunflower example) but more complicated and convoluted when the genotypic data are unbalanced (as in the cattle example) [75, 79, 111, 112]. Although methods for fixed effect analyses of factorial treatment designs (multilocus genetic models) with unbalanced data are well known [56, 79, 109, 110, 112], there are several model fitting and parameter estimation variations that can lead to dramatically different parameter estimates and statistical inferences. This is perfectly illustrated by the cattle example where the coefficients of determination (analogous but not identical to  $H_M^2$ ) from Type I, II, and III analyses were substantially different from each other and from  $H_M^2$  estimates from the random effects analysis (Tables 1 and 2). The differences and ambiguities among the different fixed effects approaches disappear when the random effects approach is applied to the problem.

The analysis of markers as random effects in multilocus analyses of known or candidate genes with large effects with ASV, although historically uncommon, simultaneously yields unbiased estimates of the variance component ratios investigated in the present study ( $p$  and  $H_M^2$ ) and best linear unbiased predictors (BLUPs) of the additive and dominance effects of the causative loci identified by marker associations, in addition to solving the often ambiguous problems that arise in fixed effects analyses of unbalanced data [32, 75, 77, 79, 112, 113]. As discussed in depth below and illustrated through a reanalysis of the cattle and sunflower examples (Table 2), the random effects approach we described (ASV with REML estimation of the variance components) yields accurate estimates of the underlying genetic parameters (variance component ratios and BLUPs of marker effects) from a *single* unambiguous generalized linear mixed model analysis, whereas wildly different parameter estimates can arise among the multitude of fixed effects analyses that investigators might elect to apply in practice when the underlying genotypic and phenotypic data are unbalanced (Tables 1 and 2).

As substantiated by our simulation analyses (Figs 1 and 2), ASV with REML estimation of the underlying variance components yields accurate estimates of  $p$  and  $H_M^2$  for marker loci and interactions between marker loci, both individually and collectively, and BLUPs of the additive and dominance effects of marker loci [76, 113–115]. When the genotypic data are unbalanced, the order with which marker and marker  $\times$  marker effects enter the genetic model profoundly affects parameter estimates and statistical inferences in fixed effect analyses [56, 72, 74, 116]. To illustrate this, the main effects of marker loci A, B, and C were estimated for the six possible Type I ANOVA orders of the three loci (ABC, ACB, BAC, BCA, CAB, and CBA) (Table 2). Predictably, the reduction in the error sums of squares for a particular locus differed for each Type I order in the cattle example: the Type I SS ranged from 591.4 to 3,552.3 for  $rs10$ , 4,880.5 to 9,504.4 for  $rs20$ , and 4,259.6 to 8,384.8 for  $rs45$ . The  $R^2$ , or PVE, estimates for marker loci were radically different among the six Type I ANOVA and Type II and III analyses. The Type I SS were, in addition, significantly greater than the Type III SS for nearly every factor. Although Type III statistics are commonly estimated and reported in analyses of factorial treatment designs with unbalanced data, there are compelling arguments for estimating Type II statistics [109, 110]; nevertheless, as we have argued, the fixed effects approach is unnecessary.

Broadly speaking, the large effect loci segregating in a population are typically necessary but not sufficient for predicting genetic merit or disease risks but are often important enough to warrant deeper study and, in animal and plant breeding, direct selection via MAS or direct modelling in genome selection applications [21, 32, 57]. The BLUP (random marker effects) approach we applied was designed to align the study of loci with large and highly predictive effects with the BLUP approaches commonly applied to genomic prediction problems that are

agnostic or indifferent to the effects of individual loci, the so-called “black box” of genomic prediction [6, 7, 20, 21, 88, 117–121]. The predictive markers associated with large effect marker loci can be integrated into the genome-wide framework of marker loci applied in genomic prediction or incorporated as fixed effects when estimating GEBVs or PRSs [21, 54, 57–61]. One of the greatest strengths of the random effects (BLUP) approach is that the genetic parameters can be estimated from a single REML analysis free of the challenges and uncertainty associated with the fixed effects model building process [79, 109, 110, 112]. Finally, if our conclusions are correct, the complex trait analysis literature is riddled with overestimates of the genotypic and phenotypic variances explained by specific genes or QTL.

## Materials and methods

### Simulation studies

We used computer simulation to estimate the bias and assess the accuracy of uncorrected and bias-corrected REML estimates of  $p$  and  $H_M^2$  for 21 study designs (S1 Table and S4 Text). Phenotypic observations ( $y_{ijk}$ ) for LMMs (1) and (2) were simulated for  $n_M = 3$  genotypes/marker locus and 21 combinations of study design variables ( $n_G$ ,  $r_G$ ,  $r_M$ , and  $H^2$ ) with balanced or unbalanced data (S1 Table). Simulations were performed to assess the accuracy of REML estimates of  $p$  and  $H_M^2$  for 21 study designs with 1,000 replicates per study design (S1 Table). The phenotypic observations for each sample were obtained by generating random normal variables for entries, markers, and residuals using the R function *rnorm()* with known means and variances [122] as described by [123, 124]. The simulated random effects of entries, markers, and replications in LMMs (1) and (2) were summed to obtain  $n = n_G r_G$  phenotypic observations for each study design. Variance components for the random effects in LMMs (1) and (2) were estimated using the REML function implemented in and assess the accuracy of AMV and ASV estimators of  $p$  and  $H_M^2$ . For study designs 1–6, the true marker heritability randomly varied from 0 to 1. Study designs 1–6 demonstrate how different numbers of marker loci ( $m$ ) and unbalanced data affect estimates of  $p$  and  $H_M^2$  (Fig 1; S1 Table). For study designs 5 and 6, we randomly deleted 10 and 33% of the phenotypic observations, respectively, to create unbalanced data. For study designs 7–21, the true variances of the independent variables were fixed for all samples, which allowed us to estimate the bias and relative bias associated with the different estimators (the biases are shown in S4 Text). Study designs 7–21 illustrate how  $r_G$ ,  $n_G$ , and  $H_M^2$  affected the biases and relative biases of  $p$  and  $H_M^2$  (Fig 2; S1 Table). The variance components were estimated using REML in the *lme4::lmer()* v1.1–21 [78] package in R v4.0.2 [122]. We estimated the sample variances of AMV and ASV estimates of  $p$  for each study design (S1 Table). Finally, we developed estimators of the sampling variances of  $p$  and  $H_M^2$  using the delta method [25, 106], as shown in S5 Text.

### Estimation examples

To illustrate the application of bias-correction methods and the differences between AMV and bias-corrected AMV estimates of  $p$  and  $H_M^2$ , we reanalyzed data from a GWAS study in cattle (*Bos taurus*), a QTL mapping study in oilseed sunflower (*Helianthus annuus* L.) [53], and a GWAS study of Fusarium wilt resistance in strawberry (*Fragaria × ananassa* Duchesne ex Rozier) [86]. For the sunflower study, two replications ( $r_G = 2$ ) of  $n_G = 146$  recombinant inbred lines (RILs) were phenotyped for seed oil concentration (g/kg) and genotyped for three marker loci (*BR*, *PHY*, and *HYP*) with two homozygous marker genotypes/locus [53]. For the cattle study, unreplicated entries ( $r_G = 1$ ;  $n_G = 2$ , 973) were phenotyped for white spotting (%) and genotyped for three marker loci (*rs10*, *rs45*, *rs20*) with three marker genotypes per locus [85].

LMM (2) expanded to three marker loci with all possible interactions among marker loci is:

$$y_{hijkl} = \mu + BR_h + PHY_i + HYP_j + BR \times PHY_{hi} + BR \times HYP_{hj} + PHY \times HYP_{ij} + BR \times PHY \times HYP_{hij} + G : (BR \times PHY \times HYP)_{hij(k)} + \epsilon_{hijkl} \tag{27}$$

where  $BR_h$  is the  $h^{th}$  effect of the *BR* locus,  $PHY_i$  is the  $i^{th}$  effect of the *PHY* locus,  $HYP_j$  is the  $j^{th}$  effect of the *HYP* locus,  $G : (BR \times PHY \times HYP)_{hij(k)}$  is the  $k^{th}$  effect of entries nested in the  $hij^{th}$  *BR* × *PHY* × *HYP* interaction, and  $\epsilon_{hijkl}$  is the  $hijkl^{th}$  residual effect. The data for RILs were balanced, whereas the data for marker genotypes were slightly unbalanced. Each of the eight *BR* × *PHY* × *HYP* homozygotes were observed in the RIL population; however, the number of entries nested in each marker genotype ( $n_{G:M}$ ) varied from  $n_{G:BR} = 81 : 65$ ,  $n_{G:PHY} = 60 : 86$ , and  $n_{G:HYP} = 70 : 76$ . Variance components for LMMs (1) and (27) were estimated using the REML method in *lme4::lmer()* [78]. The marker-associated genetic variances for individual marker loci and two- and three-way interactions among marker loci were bias-corrected using the formula described in S1, S2 and S3 Texts.

For the strawberry study, four replications ( $r_G = 4$ ) of 565 entries ( $n_G = 565$ ) from a genome-wide association study (GWAS) were phenotyped for resistance to *Fusarium wilt* and genotyped for single nucleotide polymorphism (SNP) markers in LD with *FW1*, a dominant gene conferring resistance to *Fusarium oxysporum* f.sp. *fragariae*, the causal pathogen [86]. The replications were asexually propagated clones of individuals; hence, the expected causal variance among individuals was equal to the total genetic variation in the population, analogous to monozygotic twins [25]. Genetic parameters were estimated for two SNP markers (AX493 and AX396) that were tightly linked to *FW1* [86]. The genotypic data for both markers were highly unbalanced. Genotype numbers were 141 AA : 282 Aa : 141 aa for AX493 and 16 AA : 177 Aa : 371 aa for AX396, where A and a are alternate SNP alleles. The variance components were estimated for LMMs (1) and (2) using REML method implemented in the R package *lme4::lmer()* [78]. REML estimates of the marker-associated genetic variances for both marker loci were bias-corrected using the approach described in S1 Text.

For the cattle study, we used a model similar to (27) for the analysis. However, because entries are unreplicated in this experiment, we cannot include the entries nested in the three-way marker interaction ( $G : M$ ) term because it has the same levels as the residual. The LMM for this case study is:

$$y_{hijk} = \mu + rs10_h + rs45_i + rs20_j + rs10 \times rs45_{hi} + rs10 \times rs20_{hj} + rs45 \times rs20_{ij} + rs10 \times rs45 \times rs20_{hij} + G : (rs10 \times rs45 \times rs20)_{hij(k)} \tag{28}$$

where  $rs10_h$  is the  $h^{th}$  effect of the peak SNP (*rs109979909*) on chromosome 2,  $rs45_i$  is the  $i^{th}$  effect of the peak SNP on chromosome 6 (*rs451683615*),  $rs20_j$  is the  $j^{th}$  effect of the peak SNP on chromosome 22 (*rs209784468*), and  $G : (rs10 \times rs45 \times rs20)_{hij(k)}$  is the  $hij(k)^{th}$  residual effect comprising residual genetic effects  $G : M$  and residual error. In this experiment, there were  $k$  entries and  $k$  observations, and because of this we cannot fit LMM (1) without incorporating pedigree or genomic relatedness. In this single case, we estimate  $\hat{s}_p^2 = 10.42$  from the log transformed data to use in the denominator of  $\hat{H}_M^2$ .

We used the SAS package PROC GLM [77] for Type I and III analyses of the sunflower and cattle data with marker loci as fixed effects. Type I analyses were done for the six possible orders of main effects (ABC, ACB, BAC, BCA, CAB, and CBA) and a single order for marker × marker interactions (A × B, A × C, B × C, and A × B × C), where A, B, and C are the three marker loci (factors). For the ABC order, the reduction SS for the main effects were R(A |  $\mu$ ), R(B |  $\mu, A$ ), and R(C |  $\mu, A, B$ ), where  $\mu$  is the population mean. Similarly, for the ACB

order, the reduction SS for the main effects were  $R(A | \mu)$ ,  $R(C | \mu, A)$ , and  $R(B | \mu, A, C)$ , and so on for the other four orders (BAC, BCA, CAB, and CBA).  $C, A \times B, A \times C, B \times C, A \times B \times C$ ), the reduction SS for the main effect of B was  $R(B | B, C, A \times B, A \times C, B \times C, A \times B \times C)$ . For comparison, the Type III reduction SS for the main effects were  $R(A | \mu, B, C, A \times B, A \times C, B \times C, \text{ and } A \times B \times C)$ ,  $R(B | \mu, A, C, A \times B, A \times C, B \times C, \text{ and } A \times B \times C)$ , and  $R(C | \mu, A, B, A \times B, A \times C, B \times C, \text{ and } A \times B \times C)$ .

**Supporting information**

**S1 Table. Simulation study designs and variables.** Normally distributed phenotypic observations were simulated for 21 study designs and associated linear mixed models by varying the number of observations ( $n = n_G \times r_G$ ), the number of entries ( $n_G$ ), the number of replications/entry ( $r_G$ ), the number of marker loci ( $m$ ),  $n_M = 3$  genotypes/marker locus, the number of entries/marker genotype ( $n_{G:M}$ ), and marker heritability ( $H_M^2$ ). One thousand samples of size  $n$  were simulated for each study design. The segregation of a single marker locus in an  $F_2$  population was simulated in study design 4. The number of entries nested in marker genotypes for study design 4 was equivalent to the expected number for the segregation of a co-dominant DNA marker in a population segregating 1 AA : 2 Aa : 1 aa for a single marker locus. In this example, there are 135 entries nested in AA, 270 entries nested in Aa, and 135 entries nested in aa and each are replicated 5 times. simulates the segregation of a single locus in an  $F_2$  population. The number of entries/genotype for study design 4. (PDF)

**S1 Fig. Accuracy of AMV and ASV estimators of marker heritability when the phenotypic variance is estimated by pooling marker and residual genetic sources of variation** ( $\sigma_M^2 + \sigma_{G:M}^2$ ). AMV and ASV estimates of  $H_M^2$  when  $\sigma_G^2$  from LMM (1) is replaced with  $\hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2$  for AMV from LMM (2) or  $k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2$  for ASV. Estimates are shown for 1,000 segregating populations simulated for different numbers of entries ( $n_G$  individuals, families, or strains), five replications/entry ( $r_G = 5$ ), true marker heritability ( $H_M^2$ ) ranging from 0 to 1, and one to three marker loci with three genotypes/marker locus ( $n_{M1} = 3$ ). The AMV estimates (shown in red) equal  $\hat{\sigma}_M^2 / (\hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 + r_G^{-1} \hat{\sigma}_\epsilon^2)$ , whereas the ASV estimates (shown in blue) equal  $k_M \hat{\sigma}_M^2 / (k_M \hat{\sigma}_M^2 + \hat{\sigma}_{G:M}^2 + r_G^{-1} \hat{\sigma}_\epsilon^2)$ . AMV estimates of marker heritability ( $\hat{H}_{M^*}^2$ ; red highlighted observations) and ASV estimates of marker heritability ( $\hat{H}_{M^*}^2$ ; blue highlighted observations) are shown for: (A) one locus with balanced data for  $n_G = 540$  entries (study design 1); (B) two marker loci with interaction ( $M1, M2$ , and  $M1 \times M2$ ) and balanced data for  $n_G = 540$  (study design 2); (C) three marker loci with interactions ( $M1, M2, M3, M1 \times M2, M1 \times M3, M2 \times M3$ , and  $M1 \times M2 \times M3$ ) and balanced data for  $n_G = 540$  (study design 3); (D) a population segregating 1:2:1 for a single marker locus with  $r_{G:M} = 135$  entries for both homozygotes and  $r_{G:M} = 270$  heterozygous entries, and  $n_G = 540$  (study design 4); (E) one locus with 10% randomly missing data among 540 entries (study design 5); and (F) one locus with 33% randomly missing data among 540 entries (study design 6). Study design details are shown in S1 Table. (TIFF)

**S2 Fig. Relative bias of AMV and ASV estimators of marker heritability.** Relative biases of AMV and ASV estimates of  $H_M^2$  are shown for 1,000 segregating populations simulated for different numbers of entries ( $n_G$  individuals, families, or strains), five replications/entry ( $r_G = 5$ ), true marker heritability ( $H_M^2$ ) ranging from 0 to 1, and one to three marker loci with three genotypes/marker locus ( $n_{M1} = 3$ ). AMV estimates of marker heritability ( $\hat{H}_{M^*}^2$ ; red highlighted observations) and ASV estimates of marker heritability ( $\hat{H}_{M^*}^2$ ; blue highlighted observations)

are shown for: (A) one locus with balanced data for  $n_G = 540$  entries (study design 1); (B) two marker loci with interaction ( $M1$ ,  $M2$ , and  $M1 \times M2$ ) and balanced data for  $n_G = 540$  (study design 2); (C) three marker loci with interactions ( $M1$ ,  $M2$ ,  $M3$ ,  $M1 \times M2$ ,  $M1 \times M3$ ,  $M2 \times M3$ , and  $M1 \times M2 \times M3$ ) and balanced data for  $n_G = 540$  (study design 3); (D) a population segregating 1:2:1 for one marker locus with  $r_{G:M} = 135$  entries for both homozygotes and  $r_{G:M} = 270$  heterozygous entries, and  $n_G = 540$  (study design 4); (E) one locus with 10% randomly missing data among 540 entries (study design 5); and (F) one locus with 33% randomly missing data among 540 entries (study design 6). Study design details are shown in [S1 Table](#). (TIFF)

**S1 Text. ASV estimator of the fraction of the genetic variance associated with a single marker locus for unbalanced data.**

(PDF)

**S2 Text. ASV estimator of the fraction of the genetic variance associated with two marker loci for unbalanced data.**

(PDF)

**S3 Text. ASV estimator of the fraction of the genetic variance associated with three marker loci for unbalanced data.**

(PDF)

**S4 Text. Biases of AMV and ASV estimators of marker-associated variance.**

(PDF)

**S5 Text. Sample variances for AMV and ASV estimators of  $p$  and  $H_M^2$ .**

(PDF)

## Author Contributions

**Conceptualization:** Mitchell J. Feldmann, Hans-Peter Piepho, William C. Bridges, Steven J. Knapp.

**Data curation:** Mitchell J. Feldmann.

**Formal analysis:** Mitchell J. Feldmann, Hans-Peter Piepho, William C. Bridges, Steven J. Knapp.

**Funding acquisition:** Hans-Peter Piepho, Steven J. Knapp.

**Investigation:** Mitchell J. Feldmann, Hans-Peter Piepho, Steven J. Knapp.

**Methodology:** Mitchell J. Feldmann, Hans-Peter Piepho, William C. Bridges, Steven J. Knapp.

**Project administration:** Steven J. Knapp.

**Resources:** Mitchell J. Feldmann, Steven J. Knapp.

**Software:** Mitchell J. Feldmann.

**Supervision:** Mitchell J. Feldmann, Steven J. Knapp.

**Validation:** Mitchell J. Feldmann, Hans-Peter Piepho, William C. Bridges, Steven J. Knapp.

**Visualization:** Mitchell J. Feldmann, Hans-Peter Piepho, Steven J. Knapp.

**Writing – original draft:** Mitchell J. Feldmann, Hans-Peter Piepho, Steven J. Knapp.

**Writing – review & editing:** Mitchell J. Feldmann, Hans-Peter Piepho, William C. Bridges, Steven J. Knapp.

## References

1. Lander ES, Schork NJ. Genetic dissection of complex traits. *Science*. 1994; 265(5181):2037–2048. <https://doi.org/10.1126/science.8091226> PMID: 8091226
2. Risch NJ. Searching for genetic determinants in the new millennium. *Nature*. 2000; 405(6788):847–856. <https://doi.org/10.1038/35015718> PMID: 10866211
3. Glazier AM, Nadeau JH, Aitman TJ. Finding genes that underlie complex traits. *science*. 2002; 298(5602):2345–2349. <https://doi.org/10.1126/science.1076641> PMID: 12493905
4. Consortium CT, et al. The nature and identification of quantitative trait loci: a community's view. *Nature reviews Genetics*. 2003; 4(11):911. <https://doi.org/10.1038/nrg1206>
5. Hirschhorn JN, Daly MJ. Genome-wide association studies for common diseases and complex traits. *Nature reviews genetics*. 2005; 6(2):95–108. <https://doi.org/10.1038/nrg1521> PMID: 15716906
6. Hill WG. Understanding and using quantitative genetic variation. *Philos Trans R Soc London, Ser B*. 2010; 365(1537):73–85. <https://doi.org/10.1098/rstb.2009.0203> PMID: 20008387
7. Hill WG. Quantitative genetics in the genomics era. *Curr Genomics*. 2012; 13(3):196–206. <https://doi.org/10.2174/138920212800543110> PMID: 23115521
8. Botstein D, White RL, Skolnick M, Davis RW. Construction of a genetic linkage map in man using restriction fragment length polymorphisms. *American journal of human genetics*. 1980; 32(3):314. PMID: 6247908
9. Sachidanandam R, Weissman D, Schmidt SC, Kakol JM, Stein LD, Marth G, et al. A map of human genome sequence variation containing 1.42 million single nucleotide polymorphisms. *Nature*. 2001; 409(6822):928–934. <https://doi.org/10.1038/35057149> PMID: 11237013
10. Huang X, Feng Q, Qian Q, Zhao Q, Wang L, Wang A, et al. High-throughput genotyping by whole-genome resequencing. *Genome research*. 2009; 19(6):1068–1076. <https://doi.org/10.1101/gr.089516.108> PMID: 19420380
11. Li R, Li Y, Fang X, Yang H, Wang J, Kristiansen K, et al. SNP detection for massively parallel whole-genome resequencing. *Genome research*. 2009; 19(6):1124–1132. <https://doi.org/10.1101/gr.088013.108> PMID: 19420381
12. Lander ES, Botstein D. Mapping mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics*. 1989; 121(1):185–199. <https://doi.org/10.1093/genetics/121.1.185> PMID: 2563713
13. Mackay TFC. The genetic architecture of quantitative traits. *Annu Rev Genet*. 2001; 35(1):303–339. <https://doi.org/10.1146/annurev.genet.35.102401.090633> PMID: 11700286
14. Klein RJ, Zeiss C, Chew EY, Tsai JY, Sackler RS, Haynes C, et al. Complement factor H polymorphism in age-related macular degeneration. *Science*. 2005; 308(5720):385–389. <https://doi.org/10.1126/science.1109557> PMID: 15761122
15. Consortium WTCC, et al. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*. 2007; 447(7145):661. <https://doi.org/10.1038/nature05911>
16. Mackay TFC, Stone EA, Ayroles JF. The genetics of quantitative traits: challenges and prospects. *Nat Rev Genet*. 2009; 10(8):565. <https://doi.org/10.1038/nrg2612> PMID: 19584810
17. Visscher PM, Brown MA, McCarthy MI, Yang J. Five years of GWAS discovery. *Am J Hum Genet*. 2012; 90(1):7–24. <https://doi.org/10.1016/j.ajhg.2011.11.029> PMID: 22243964
18. Pasaniuc B, Price AL. Dissecting the genetics of complex traits using summary association statistics. *Nat Rev Genet*. 2017; 18(2):117. <https://doi.org/10.1038/nrg.2016.142> PMID: 27840428
19. Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 years of GWAS discovery: biology, function, and translation. *Am J Hum Genet*. 2017; 101(1):5–22. <https://doi.org/10.1016/j.ajhg.2017.06.005> PMID: 28686856
20. Meuwissen T, Hayes B, Goddard M. Prediction of total genetic value using genome-wide dense marker maps. *Genetics*. 2001; 157(4):1819–1829. <https://doi.org/10.1093/genetics/157.4.1819> PMID: 11290733
21. Wray NR, Kempner KE, Hayes BJ, Goddard ME, Visscher PM. Complex trait prediction from genome data: contrasting EBV in livestock to PRS in humans: genomic prediction. *Genetics*. 2019; 211(4):1131–1141. <https://doi.org/10.1534/genetics.119.301859> PMID: 30967442

22. Crouch DJ, Bodmer WF. Polygenic inheritance, GWAS, polygenic risk scores, and the search for functional variants. *Proceedings of the National Academy of Sciences*. 2020; 117(32):18924–18933. <https://doi.org/10.1073/pnas.2005634117> PMID: 32753378
23. Wray NR, Lin T, Austin J, McGrath JJ, Hickie IB, Murray GK, et al. From basic science to clinical application of polygenic risk scores: a primer. *JAMA psychiatry*. 2021; 78(1):101–109. <https://doi.org/10.1001/jamapsychiatry.2020.3049> PMID: 32997097
24. Falconer D, Mackay T. *Introduction to Quantitative Genetics*. Harlow, Essex, UK. Longmans Green; 1996.
25. Lynch M, Walsh B. *Genetics and analysis of quantitative traits*. vol. 1. Sinauer Sunderland, MA; 1998.
26. Walsh B. Quantitative genetics in the age of genomics. *Theoretical Population Biology*. 2001; 59(3):175–184. <https://doi.org/10.1006/tpbi.2001.1512> PMID: 11444958
27. Visscher PM, Medland SE, Ferreira MA, Morley KI, Zhu G, Cornes BK, et al. Assumption-free estimation of heritability from genome-wide identity-by-descent sharing between full siblings. *PLoS Genet*. 2006; 2(3). <https://doi.org/10.1371/journal.pgen.0020041> PMID: 16565746
28. Roff DA. A centennial celebration for quantitative genetics. *Evolution*. 2007; 61(5):1017–1032. <https://doi.org/10.1111/j.1558-5646.2007.00100.x> PMID: 17492957
29. Hill WG, Goddard ME, Visscher PM. Data and theory point to mainly additive genetic variance for complex traits. *PLoS Genet*. 2008; 4(2):e1000008. <https://doi.org/10.1371/journal.pgen.1000008>
30. Visscher PM, Hill WG, Wray NR. Heritability in the genomics era—concepts and misconceptions. *Nat Rev Genet*. 2008; 9(4):255–266. <https://doi.org/10.1038/nrg2322> PMID: 18319743
31. Roff DA. *Evolutionary quantitative genetics*. Springer Science & Business Media; 2012.
32. Bernardo R. Reinventing quantitative genetics for plant breeding: something old, something new, something borrowed, something BLUE. *Heredity*. 2020; 125(6):375–385. <https://doi.org/10.1038/s41437-020-0312-1> PMID: 32296132
33. Andersson L. Genetic dissection of phenotypic diversity in farm animals. *Nature Reviews Genetics*. 2001; 2(2):130–138. <https://doi.org/10.1038/35052563> PMID: 11253052
34. Hayes B, Goddard ME. The distribution of the effects of genes affecting quantitative traits in livestock. *Genetics Selection Evolution*. 2001; 33(3):1–21. <https://doi.org/10.1186/1297-9686-33-3-209> PMID: 11403745
35. Mackay TF. Quantitative trait loci in *Drosophila*. *Nature reviews genetics*. 2001; 2(1):11–20. <https://doi.org/10.1038/35047544> PMID: 11253063
36. Andersson L, Georges M. Domestic-animal genomics: deciphering the genetics of complex traits. *Nature Reviews Genetics*. 2004; 5(3):202–212. <https://doi.org/10.1038/nrg1294> PMID: 14970822
37. Anderson JA, Chao S, Liu S. Molecular breeding using a major QTL for Fusarium head blight resistance in wheat. *Crop Science*. 2007; 47:S–112. <https://doi.org/10.2135/cropsci2007.04.0006IPBS>
38. Septiningsih EM, Pamplona AM, Sanchez DL, Neeraja CN, Vergara GV, Heuer S, et al. Development of submergence-tolerant rice cultivars: the Sub1 locus and beyond. *Annals of Botany*. 2009; 103(2):151–160. <https://doi.org/10.1093/aob/mcn206> PMID: 18974101
39. Lorenz K, Cohen BA. Small-and large-effect quantitative trait locus interactions underlie variation in yeast sporulation efficiency. *Genetics*. 2012; 192(3):1123–1132. <https://doi.org/10.1534/genetics.112.143107> PMID: 22942125
40. Saatchi M, Schnabel RD, Taylor JF, Garrick DJ. Large-effect pleiotropic or closely linked QTL segregate within and across ten US cattle breeds. *BMC Genomics*. 2014; 15(1):1–17. <https://doi.org/10.1186/1471-2164-15-442>
41. Huang H, Cao J, Hanif Q, Wang Y, Yu Y, Zhang S, et al. Genome-wide association study identifies energy metabolism genes for resistance to ketosis in Chinese Holstein cattle. *Animal genetics*. 2019; 50(4):376–380. <https://doi.org/10.1111/age.12802> PMID: 31179571
42. Freebern E, Santos DJ, Fang L, Jiang J, Gaddis KLP, Liu GE, et al. GWAS and fine-mapping of livability and six disease traits in Holstein cattle. *BMC genomics*. 2020; 21(1):1–11. <https://doi.org/10.1186/s12864-020-6461-z> PMID: 31931710
43. Li B, VanRaden P, Null D, O'Connell J, Cole J. Major quantitative trait loci influencing milk production and conformation traits in Guernsey dairy cattle detected on *Bos taurus* autosome 19. *Journal of Dairy Science*. 2021; 104(1):550–560. <https://doi.org/10.3168/jds.2020-18766> PMID: 33189290
44. Bernardo R. Molecular markers and selection for complex traits in plants: learning from the last 20 years. *Crop Sci*. 2008; 48(5):1649–1664. <https://doi.org/10.2135/cropsci2008.03.0131>
45. Bernardo R. Bandwagons I, too, have known. *Theor Appl Genet*. 2016; 129(12):2323–2332. <https://doi.org/10.1007/s00122-016-2772-5> PMID: 27681088



46. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, et al. Finding the missing heritability of complex diseases. *Nature*. 2009; 461(7265):747. <https://doi.org/10.1038/nature08494> PMID: 19812666
47. Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, et al. Missing heritability and strategies for finding the underlying causes of complex disease. *Nat Rev Genet*. 2010; 11(6):446–450. <https://doi.org/10.1038/nrg2809> PMID: 20479774
48. Young AL. Solving the missing heritability problem. *PLoS Genet*. 2019; 15(6):e1008222. <https://doi.org/10.1371/journal.pgen.1008222> PMID: 31233496
49. Bernardo R. What if we knew all the genes for a quantitative trait in hybrid crops? *Crop Sci*. 2001; 41(1):1–4.
50. Hill WG. Applications of population genetics to animal breeding, from Wright, Fisher and Lush to genomic prediction. *Genetics*. 2014; 196(1):1–16. <https://doi.org/10.1534/genetics.112.147850> PMID: 24395822
51. De Villemereuil P, Morrissey MB, Nakagawa S, Schielzeth H. Fixed-effect variance and the estimation of repeatabilities and heritabilities: issues and solutions. *Journal of Evolutionary Biology*. 2018; 31(4):621–632. <https://doi.org/10.1111/jeb.13232> PMID: 29285829
52. Gaudet MM, Kirchoff T, Green T, Vijai J, Korn JM, Guiducci C, et al. Common genetic variants and modification of penetrance of BRCA2-associated breast cancer. *PLoS Genet*. 2010; 6(10):e1001183. <https://doi.org/10.1371/journal.pgen.1001183> PMID: 21060860
53. Tang S, Leon A, Bridges WC, Knapp SJ. Quantitative trait loci for genetically correlated seed traits are tightly linked to branching and pericarp pigment loci in sunflower. *Crop Sci*. 2006; 46(2):721–734. <https://doi.org/10.2135/cropsci2005.0006-7>
54. Hayes BJ, Pryce J, Chamberlain AJ, Bowman PJ, Goddard ME. Genetic architecture of complex traits and accuracy of genomic prediction: coat colour, milk-fat percentage, and type in Holstein cattle as contrasting model traits. *PLoS genetics*. 2010; 6(9):e1001139. <https://doi.org/10.1371/journal.pgen.1001139> PMID: 20927186
55. Seabury CM, Oldeschulte DL, Saatchi M, Beever JE, Decker JE, Halley YA, et al. Genome-wide association study for feed efficiency and growth traits in US beef cattle. *BMC genomics*. 2017; 18(1):1–25. <https://doi.org/10.1186/s12864-017-3754-y> PMID: 28521758
56. Littell RC. Analysis of unbalanced mixed model data: a case study comparison of ANOVA versus REML/GLS. *Journal of Agricultural, Biological, and Environmental Statistics*. 2002; 7(4):472. <https://doi.org/10.1198/108571102816>
57. Daetwyler HD, Calus MP, Pong-Wong R, de Los Campos G, Hickey JM. Genomic prediction in animals and plants: simulation of data, validation, reporting, and benchmarking. *Genetics*. 2013; 193(2):347–365. <https://doi.org/10.1534/genetics.112.147983> PMID: 23222650
58. Gianola D. Priors in whole-genome regression: the Bayesian alphabet returns. *Genetics*. 2013; 194(3):573–596. <https://doi.org/10.1534/genetics.113.151753> PMID: 23636739
59. Moore JK, Manmathan HK, Anderson VA, Poland JA, Morris CF, Haley SD. Improving Genomic Prediction for Pre-Harvest Sprouting Tolerance in Wheat by Weighting Large-Effect Quantitative Trait Loci. *Crop Science*. 2017; 57(3):1315–1324. <https://doi.org/10.2135/cropsci2016.06.0453>
60. Rice B, Lipka AE. Evaluation of RR-BLUP Genomic Selection Models that Incorporate Peak Genome-Wide Association Study Signals in Maize and Sorghum. *The Plant Genome*. 2019; 12(1). <https://doi.org/10.3835/plantgenome2018.07.0052> PMID: 30951091
61. Spindel J, Begum H, Akdemir D, Collard B, Redoña E, Jannink J, et al. Genome-wide prediction models that incorporate de novo GWAS are a powerful new tool for tropical rice improvement. *Heredity*. 2016; 116(4):395–408. <https://doi.org/10.1038/hdy.2015.113> PMID: 26860200
62. Lande R, Thompson R. Efficiency of marker-assisted selection in the improvement of quantitative traits. *Genetics*. 1990; 124(3):743–756. <https://doi.org/10.1093/genetics/124.3.743> PMID: 1968875
63. de los Campos G, Sorensen D, Gianola D. Genomic heritability: what is it? *PLoS Genet*. 2015; 11(5):e1005048. <https://doi.org/10.1371/journal.pgen.1005048> PMID: 25942577
64. Beavis WD. QTL analyses: power, precision, and accuracy. In Patterson AH (ed) *Molecular Dissection of Complex Traits*. 1998; p. 145–162.
65. Melchinger AE, Utz HF, Schön CC. Quantitative trait locus (QTL) mapping using different testers and independent population samples in maize reveals low power of QTL detection and large bias in estimates of QTL effects. *Genetics*. 1998; 149(1):383–403. <https://doi.org/10.1093/genetics/149.1.383> PMID: 9584111
66. Utz HF, Melchinger AE, Schön CC. Bias and sampling error of the estimated proportion of genotypic variance explained by quantitative trait loci determined from experimental data in maize using cross

- validation and validation with independent samples. *Genetics*. 2000; 154(4):1839–1849. <https://doi.org/10.1093/genetics/154.4.1839> PMID: 10866652
67. Allison DB, Fernandez JR, Heo M, Zhu S, Etzel C, Beasley TM, et al. Bias in estimates of quantitative-trait–locus effect in genome scans: demonstration of the phenomenon and a method-of-moments procedure for reducing bias. *Am J Hum Genet*. 2002; 70(3):575–585. <https://doi.org/10.1086/339273> PMID: 11836648
  68. Xu S. Theoretical basis of the Beavis effect. *Genetics*. 2003; 165(4):2259–2268. <https://doi.org/10.1093/genetics/165.4.2259> PMID: 14704201
  69. Bernardo R. What proportion of declared QTL in plants are false? *Theor Appl Genet*. 2004; 109(2):419–424. <https://doi.org/10.1007/s00122-004-1639-3> PMID: 15085262
  70. Göring HH, Terwilliger JD, Blangero J. Large upward bias in estimation of locus-specific effects from genomewide scans. *Am J Hum Genet*. 2001; 69(6):1357–1369. <https://doi.org/10.1086/324471> PMID: 11593451
  71. Henderson CR. Estimation of variance and covariance components. *Biometrics*. 1953; 9(2):226–252. <https://doi.org/10.2307/3001853>
  72. Searle SR, Gruber MH. *Linear models*. Wiley Online Library; 1971.
  73. Piepho HP. A coefficient of determination ( $F^2$ ) for generalized linear mixed models. *Biom J*. 2019; 61:860–872. PMID: 30957911
  74. Searle SR. *Linear models for unbalanced data*. Wiley; 1987.
  75. Littell RC, Milliken GA, Stroup WW, Wolfinger RD, Schabenberger O. *SAS system for mixed models*. vol. 633. SAS institute Cary, NC; 1996.
  76. Bernardo R. *Breeding for quantitative traits in plants*. vol. 1. Stemma press Woodbury, MN; 2002.
  77. Inc SI. *SAS/STAT 13.1 User's Guide: Chapter 43—The GLIMMIX Procedure*. Author Cary, NC; 2013. Available from: <https://support.sas.com/documentation/onlinedoc/stat/131/glimmix.pdf>.
  78. Bates D, Mächler M, Bolker B, Walker S. Fitting Linear Mixed-Effects Models Using lme4. *J Stat Softw*. 2015; 67(1):1–48. <https://doi.org/10.18637/jss.v067.i01>
  79. Gbur EE, Stroup WW, McCarter KS, Durham S, Young LJ, Christman M, et al. *Analysis of generalized linear mixed models in the agricultural and natural resources sciences*. vol. 156. John Wiley & Sons; 2020.
  80. Broman KW, Sen S. *A Guide to QTL Mapping with R/QTL*. vol. 46. Springer; 2009.
  81. Cantor RM, Lange K, Sinsheimer JS. Prioritizing GWAS results: a review of statistical methods and recommendations for their application. *The American Journal of Human Genetics*. 2010; 86(1):6–22. <https://doi.org/10.1016/j.ajhg.2009.11.017> PMID: 20074509
  82. Beavis W, Smith O, Grant D, Fincher R. Identification of quantitative trait loci using a small sample of topcrossed and F4 progeny from maize. *Crop Sci*. 1994; 34(4):882–896. <https://doi.org/10.2135/cropsci1994.0011183X003400040010x>
  83. Luo L, Mao Y, Xu S. Correcting the bias in estimation of genetic variances contributed by individual QTL. *Genetica*. 2003; 119(2):107–114. <https://doi.org/10.1023/A:1026028928003> PMID: 14620950
  84. Zhang J, Yue C, Zhang Y. Bias correction for estimated QTL effects using the penalized maximum likelihood method. *Heredity*. 2012; 108(4):396–402. <https://doi.org/10.1038/hdy.2011.86> PMID: 21934700
  85. Jivanji S, Worth G, Lopdell TJ, Yeates A, Couldrey C, Reynolds E, et al. Genome-wide association analysis reveals QTL and candidate mutations involved in white spotting in cattle. *Genet Sel Evol*. 2019; 51(1):62. <https://doi.org/10.1186/s12711-019-0506-2> PMID: 31703548
  86. Pincot DD, Poorten TJ, Hardigan MA, Harshman JM, Acharya CB, Cole GS, et al. Genome-wide association mapping uncovers Fw1, a dominant gene conferring resistance to Fusarium wilt in strawberry. *G3: Genes, Genomes, Genet*. 2018; 8(5):1817–1828. <https://doi.org/10.1534/g3.118.200129> PMID: 29602808
  87. Conner JK, Hartl DL, et al. *A primer of ecological genetics*. Sinauer Associates Incorporated; 2004.
  88. Mrode RA. *Linear models for the prediction of animal breeding values*. Cabi; 2014.
  89. Choy Y, Brinks J, Bourdon R. Repeated-measure animal models to estimate genetic components of mature weight, hip height, and body condition score. *Journal of animal science*. 2002; 80(8):2071–2077. <https://doi.org/10.2527/2002.8082071x> PMID: 12211374
  90. Dekkers JC. Commercial application of marker-and gene-assisted selection in livestock: strategies and lessons. *J Anim Sci*. 2004; 82(suppl\_13):E313–E328. PMID: 15471812

91. Dekkers J. Prediction of response to marker-assisted and genomic selection using selection index theory. *J Anim Breed Genet.* 2007; 124(6):331–341. <https://doi.org/10.1111/j.1439-0388.2007.00701.x> PMID: 18076470
92. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci.* 2008; 91(11):4414–4423. <https://doi.org/10.3168/jds.2007-0980> PMID: 18946147
93. Sun X, Habier D, Fernando RL, Garrick DJ, Dekkers JC. Genomic breeding value prediction and QTL mapping of QTLMAS2010 data using Bayesian methods. In: *BMC proceedings.* BioMed Central; 2011. p. 1–8.
94. Isik F, Holland J, Maltecca C. *Genetic data analysis for plant and animal breeding.* Springer; 2017.
95. Astle W, Balding DJ. Population structure and cryptic relatedness in genetic association studies. *Statistical Science.* 2009; 24(4):451–471. <https://doi.org/10.1214/09-STS307>
96. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. Common SNPs explain a large proportion of the heritability for human height. *Nat Genet.* 2010; 42(7):565. <https://doi.org/10.1038/ng.608> PMID: 20562875
97. Endelman JB, Jannink JL. Shrinkage estimation of the realized relationship matrix. *G3: Genes, Genomes, Genet.* 2012; 2(11):1405–1413. <https://doi.org/10.1534/g3.112.004259> PMID: 23173092
98. Cary N. *SAS/STAT 13.1 User's Guide;* 2013.
99. Piepho HP, Moehring J, Schulz-Streeck T, Ogutu JO. A stage-wise approach for the analysis of multi-environment trials. *Biom J.* 2012; 54(6):844–860. <https://doi.org/10.1002/bimj.201100219>
100. Damesa TM, Möhring J, Worku M, Piepho HP. One step at a time: stage-wise analysis of a series of experiments. *Agron J.* 2017; 109(3):845–857. <https://doi.org/10.2134/agnonj2016.07.0395>
101. Damesa TM, Hartung J, Gowda M, Beyene Y, Das B, Semagn K, et al. Comparison of Weighted and Unweighted Stage-Wise Analysis for Genome-Wide Association Studies and Genomic Selection. *Crop Sci.* 2019; 59:2572–2584. <https://doi.org/10.2135/cropsci2019.04.0209>
102. Schmidt P, Hartung J, Bennewitz J, Piepho HP. Heritability in Plant Breeding on a Genotype-Difference Basis. *Genetics.* 2019; 212(4):991–1008. <https://doi.org/10.1534/genetics.119.302134> PMID: 31248886
103. Estaghevrou SBO, Ogutu JO, Schulz-Streeck T, Knaak C, Ouzunova M, Gordillo A, et al. Evaluation of approaches for estimating the accuracy of genomic prediction in plant breeding. *BMC Genomics.* 2013; 14(1):860. <https://doi.org/10.1186/1471-2164-14-860>
104. Efron B, Tibshirani R. The bootstrap method for assessing statistical accuracy. *Behaviormetrika.* 1985; 12(17):1–35. [https://doi.org/10.2333/bhmk.12.17\\_1](https://doi.org/10.2333/bhmk.12.17_1)
105. Oehlert GW. A note on the delta method. *The American Statistician.* 1992; 46(1):27–29. <https://doi.org/10.1080/00031305.1992.10475842>
106. Johnson NL, Kemp AW, Kotz S. *Univariate discrete distributions.* vol. 444. John Wiley & Sons; 2005.
107. Korte A, Farlow A. The advantages and limitations of trait analysis with GWAS: a review. *Plant Methods.* 2013; 9(1):29. <https://doi.org/10.1186/1746-4811-9-29> PMID: 23876160
108. Jensen J, Su G, Madsen P. Partitioning additive genetic variance into genomic and remaining polygenic components for complex traits in Dairy cattle. *BMC Genet.* 2012; 13(1):44. <https://doi.org/10.1186/1471-2156-13-44> PMID: 22694746
109. Herr DG. On the history of ANOVA in unbalanced, factorial designs: The first 30 years. *The American Statistician.* 1986; 40(4):265–270. <https://doi.org/10.2307/2684597>
110. Langsrud Ø. ANOVA for unbalanced data: Use Type II instead of Type III sums of squares. *Statistics and Computing.* 2003; 13(2):163–167. <https://doi.org/10.1023/A:1023260610025>
111. Stroup W, Littell R. Impact of variance component estimates on fixed effect inference in unbalanced linear mixed models. *Proceedings of the Kansas State University Conference on Applied Statistics in Agriculture.* 2002; 14:32–48.
112. Stroup WW, Milliken GA, Claassen EA, Wolfinger RD. *SAS for mixed models: introduction and basic applications.* SAS Institute; 2018.
113. Piepho H, Möhring J, Melchinger A, Büchse A. BLUP for phenotypic selection in plant breeding and variety testing. *Euphytica.* 2008; 161(1):209–228. <https://doi.org/10.1007/s10681-007-9449-8>
114. Searle SR, Casella G, McCulloch C. *Variance components.* John Wiley & Sons; 1992.
115. Molenaar H, Boehm R, Piepho HP. Phenotypic selection in ornamental breeding: it's better to have the BLUPs than to have the BLUEs. *Frontiers in plant science.* 2018; 9:1511. <https://doi.org/10.3389/fpls.2018.01511> PMID: 30455707
116. Hector A, Von Felten S, Schmid B. Analysis of variance with unbalanced data: an update for ecology & evolution. *Journal of animal ecology.* 2010; 79(2):308–316. <https://doi.org/10.1111/j.1365-2656.2009.01634.x> PMID: 20002862

117. Gianola D, Perez-Enciso M, Toro MA. On marker-assisted prediction of genetic value: beyond the ridge. *Genetics*. 2003; 163(1):347–365. <https://doi.org/10.1534/genetics.109.103952> PMID: [12586721](https://pubmed.ncbi.nlm.nih.gov/12586721/)
118. Goddard M, Hayes B. Genomic selection. *J Anim Breed Genet*. 2007; 124(6):323–330. <https://doi.org/10.1111/j.1439-0388.2007.00702.x> PMID: [18076469](https://pubmed.ncbi.nlm.nih.gov/18076469/)
119. Habier D, Fernando RL, Garrick DJ. Genomic BLUP decoded: a look into the black box of genomic prediction. *Genetics*. 2013; 194(3):597–607. <https://doi.org/10.1534/genetics.113.152207> PMID: [23640517](https://pubmed.ncbi.nlm.nih.gov/23640517/)
120. Meuwissen T, Hayes B, Goddard M. Genomic selection: A paradigm shift in animal breeding. *Animal Frontiers*. 2016; 6(1):6–14. <https://doi.org/10.2527/af.2016-0002>
121. Crossa J, Pérez-Rodríguez P, Cuevas J, Montesinos-López O, Jarquín D, de los Campos G, et al. Genomic selection in plant breeding: methods, models, and perspectives. *Trends Plant Sci*. 2017; 22(11):961–975. <https://doi.org/10.1016/j.tplants.2017.08.011> PMID: [28965742](https://pubmed.ncbi.nlm.nih.gov/28965742/)
122. R Core Team. R: A Language and Environment for Statistical Computing; 2020. Available from: <https://www.R-project.org/>.
123. Burton A, Altman DG, Royston P, Holder RL. The design of simulation studies in medical statistics. *Stat Med*. 2006; 25(24):4279–4292. <https://doi.org/10.1002/sim.2673> PMID: [16947139](https://pubmed.ncbi.nlm.nih.gov/16947139/)
124. Morris TP, White IR, Crowther MJ. Using simulation studies to evaluate statistical methods. *Stat Med*. 2019; 38(11):2074–2102. <https://doi.org/10.1002/sim.8086> PMID: [30652356](https://pubmed.ncbi.nlm.nih.gov/30652356/)