

Prediction of Subgenome Additive and Interaction Effects in Allohexaploid Wheat

Nicholas Santantonio,^{*1} Jean-Luc Jannink,^{*†} and Mark Sorrells^{*}

^{*}Cornell University, Plant Breeding and Genetics Section, School of Integrated Plant Sciences, College of Agriculture and Life Sciences, 240 Emerson Hall, Ithaca, NY 14853 and [†]USDA ARS, Robert W. Holley Center for Agriculture & Health, Ithaca, NY 14853

ORCID IDs: 0000-0002-4351-4023 (N.S.); 0000-0003-4849-628X (J.-L.J.); 0000-0002-7367-2663 (M.S.)

ABSTRACT Whole genome duplications have played an important role in the evolution of angiosperms. These events often occur through hybridization between closely related species, resulting in an allopolyploid with multiple subgenomes. With the availability of affordable genotyping and a reference genome to locate markers, breeders of allopolyploids now have the opportunity to manipulate subgenomes independently. This also presents a unique opportunity to investigate epistatic interactions between homeologous orthologs across subgenomes. We present a statistical framework for partitioning genetic variance to the subgenomes of an allopolyploid, predicting breeding values for each subgenome, and determining the importance of inter-genomic epistasis. We demonstrate using an allohexaploid wheat breeding population evaluated in Ithaca, NY and an important wheat dataset from CIMMYT previously shown to demonstrate non-additive genetic variance. Subgenome covariance matrices were constructed and used to calculate subgenome interaction covariance matrices for variance component estimation and genomic prediction. We propose a method to extract population structure from all subgenomes at once before covariances are calculated to reduce collinearity between subgenome estimates. Variance parameter estimation was shown to be reliable for additive subgenome effects, but was less reliable for subgenome interaction components. Predictive ability was equivalent to current genomic prediction methods. Including only inter-genomic interactions resulted in the same increase in accuracy as modeling all pairwise marker interactions. Thus, we provide a new tool for breeders of allopolyploid crops to characterize the genetic architecture of existing populations, determine breeding goals, and develop new strategies for selection of additive effects and fixation of inter-genomic epistasis.

KEYWORDS

Allopolyploidy
Epistasis
Heterosis
Genomic prediction

Gene duplication is known to be a primary driver of evolution by providing the raw genetic material for gene diversification through sub- and neofunctionalization (Haldane 1933; Ohno 1970). Whole genome duplication events, in which an entire set of genes is duplicated, occurs either through duplication of the same genome (autopolyploidy) or the union of two closely related genomes (allopolyploidy). Both types of

polyploids can exhibit non-additive genetic variation from the presence of multiple alleles (Segovia-Lerma *et al.*, 2004; Birchler *et al.*, 2010; Jiang *et al.*, 2017), although how these non-additive effects are classified needs clarification.

Statistical deviations from additivity (*i.e.*, interactions) are important contributors to genetic variation. Homologous gene interactions, also known as dominance, are deviations from an additive expectation due to different allele combinations at a single locus. Non-homologous gene interactions, commonly referred to as epistasis, are deviations from an additive expectation due to different allele combinations at two or more loci (Fisher 1919). When epistasis occurs between non-homologous loci with similar function, such as across orthologs or paralogs, these interactions are comparable to dominance effects. If interactions occur between homeologous orthologs on separate subgenomes of an allopolyploid, should we call this epistasis or dominance?

In classical hybrid variety production, divergent sets of alleles are intentionally isolated into heterotic groups and then brought back

Copyright © 2019 Santantonio *et al.*

doi: <https://doi.org/10.1534/g3.118.200613>

Manuscript received July 26, 2018; accepted for publication November 14, 2018; published Early Online November 19, 2018.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Supplemental material available at Figshare: <https://doi.org/10.25387/g3.6870110>.

¹Corresponding Author: Cornell University, Plant Breeding and Genetics Section, School of Integrated Plant Sciences, College of Agriculture and Life Sciences, 240 Emerson Hall, Ithaca, NY 14853, Email: ns722@cornell.edu

together to form a hybrid. This establishes heterozygosity (by descent) at all loci to form a homogeneous population. The union of two divergent suites of genes during the formation of an allopolyploid also results in a homogeneous population, but heterozygosity is established across homeologs rather than homologs. Diploid hybrids lose heterozygosity through segregation in following filial generations, but heterozygosity across homeologous genes is subsequently preserved through selfing in the allopolyploid (Mac Key 1970; Ozkan *et al.* 2001; Abel *et al.* 2005). Allelic interactions contribute to dominance variance in the diploid hybrid, whereas homeoallelic interactions will be present as part of the additive by additive epistatic variance in an inbred allopolyploid population. As such, allopolyploids may be thought of as an immortalized hybrid (Ellstrand and Schierenbeck 2000; Feldman *et al.* 2012), although it is not yet clear that these exhibit a true heterotic response as traditional hybrids have demonstrated.

Birchler *et al.* (2010) note that newly synthesized allopolyploids often outperform their sub-genome progenitors, and that the heterotic response appears to be exaggerated in wider inter-specific crosses. This seems to hold true even within species, where autopolyploids tend to exhibit higher vigor from wider crosses (Bingham *et al.* 1994; Segovia-Lerma *et al.* 2004). Complementation of deleterious recessive alleles (or pseudo-dominance) has long been the primary explanation of the heterotic response (Stuber *et al.* 1992; Cockerham and Zeng 1996). However, Birchler *et al.* (2010) indicate evidence against this, where purging deleterious alleles has increased the additive value of inbred parents but has not reduced the heterotic response observed in the hybrid (Duvick 1999). Complementation also seems an unlikely driver of a heterotic response in allopolyploids, as the inbred subgenome progenitors would supposedly need functional copies of these genes to survive.

The availability of affordable genome-wide markers has sparked a revolution in selection on additive variation through the use of genomic prediction models. The additive genetic merit of an individual can be estimated as the sum of its additive marker effects to produce a genomic estimated breeding value (GEBV) (Meuwissen *et al.* 2001). When the number of markers is large, marker effects are typically considered random and normally distributed such that only one parameter need be estimated. Alternatively, the additive genetic covariance between individuals can be estimated from the same genome-wide markers and used to predict additive genetic values of individuals based on relatedness (Nejati-Javaremi *et al.* 1997; Van Raden 2008). These models are equivalent for prediction under the same set of assumptions (Garrick 2007; Van Raden 2008; Strandén and Garrick 2009). Genomic prediction models have since become popular for their ability to predict the performance of genotyped individuals with no phenotypic observations. Selections on unobserved individuals allows for reduction in the cost of phenotyping and breeding cycle time, increasing the rate of genetic gain (Goddard and Hayes 2007; Heffner *et al.* 2009; Jannink *et al.* 2010; Heslot *et al.* 2015).

The potential utility of genome-wide markers has also drawn renewed interest in non-additive genetic variation in recent years (Vitezica *et al.* 2013; Martini *et al.* 2016; Jiang and Reif 2015; Huang and Mackay 2016; Jiang *et al.* 2017). Genomic prediction models that use genome-wide markers can incorporate non-additive genetic components to obtain better estimates of individual performance than based on additivity alone (Technow *et al.* 2012; Vitezica *et al.* 2013; Jiang and Reif 2015; Akdemir and Jannink 2015; Akdemir *et al.* 2017; Wolfe *et al.* 2016). In outcrossing species such as maize, prediction of dominance effects is key to harnessing heterosis in unobserved hybrids (Technow *et al.* 2012). In inbred species, additive by additive epistatic effects have been shown to significantly increase genomic prediction

accuracy (Cossa *et al.* 2010; Martini *et al.* 2016). Epistatic effects can be added to the prediction model by extending the method of expected epistatic covariance estimation Henderson (1985) to marker based covariance estimation (Jiang and Reif 2015; Martini *et al.* 2016).

The use of genome-wide markers has allowed for the partitioning of genetic variance to specific units of chromatin, previously infeasible with phenotypes alone (Visscher *et al.* 2007; Yang *et al.* 2011; Bernardo and Thompson 2016; Gage *et al.* 2017). Allopolyploids have been traditionally treated as diploids because they undergo disomic inheritance (Mac Key 1970), such that the contribution of each subgenome to the genetic variance is ignored. By assigning markers to each subgenome, an additive genetic covariance based on each subgenome can be calculated. Using these covariances in a genomic prediction model, the genetic merit of an allopolyploid individual can be assigned to each of its subgenomes. These subgenomic estimated breeding values (SGBV) can then be used to identify parents with complementary subgenome effects for crossing.

Under Hardy Weinberg equilibrium, subgenomes segregate independently, and realized estimates of additive covariance of individuals based on each subgenome will be independent. However, this does not generally hold true in breeding programs, where population structure from non-random mating is inherent. As a consequence, the estimates of additive covariance between individuals based on different subgenomes will not be independent, potentially leading to confounding of effects from each subgenome and problems partitioning variance reliably. In an attempt to circumvent this obstacle, we present an approach for removing the largest sources of genetic variance (*i.e.*, population structure) using singular value decomposition of the matrix of marker scores.

Common wheat (*Triticum aestivum*) is a staple allopolyploid crop which has undergone two allopolyploid events, resulting in three genomes, denoted A, B and D. The A genome ancestor, *Triticum uratu*, still exists today and was an early domesticate from the fertile crescent important in the neolithic revolution (Dvořák *et al.* 1993). The B genome ancestor (an *Aegilops spp.*) is believed to have since gone extinct (Blake *et al.* 1999), but the tetraploid formed by these two genomes, *Triticum turgidum*, is still cultivated today primarily as emmer and durum wheat. The D genome comes from a goat grass, *Aegilops tauschii*, which may have been incorporated in a single hybridization event as recently as 10,000 years ago (Salamini *et al.* 2002). However, recent evidence based on sequence divergence of the D genome from the A and B genome has suggested a much earlier D genome incorporation around 400,000 years ago (Marcussen *et al.* 2014). Other evidence shows that limited gene flow into the D genome may have occurred after the polyploidization event, but appears to be from a single lineage (Wang *et al.* 2013). As a result, the D genome has significantly lower genetic variation than either the A or B genome.

We demonstrate methodology to partition subgenome additive variance to estimate SGBVs as well as subgenome interactions using two allohexaploid wheat data sets, the Cornell small grains breeding program soft winter wheat breeding population dataset presented here (CNLM) and the W-GY wheat data from Cossa *et al.* (2010).

MATERIALS AND METHODS

Empirical data sets

CNLM population: The CNLM dataset consists of 8,692 phenotypic records of 1,447 soft winter wheat inbred lines evaluated at four locations near Ithaca, NY from 2007 to 2016, representing 26 environments. These phenotypic evaluations serve primarily as a first round of evaluation for grain yield and other agronomic traits before relatively few are selected for replicated regional trials around New York State. Lines are

introduced and then removed after they are deemed either fit for advanced field trials or to be discarded or recycled in the breeding program. As such, this dataset is unbalanced in nature. Most lines were not replicated within a given trial (*i.e.*, year and location), but various check varieties were used throughout these years and are typically replicated several times within a given trial.

Data were recorded for four agronomic traits, grain yield (GY), plant height (PH), test weight (TW) and heading date (HD). GY and TW have no missing data, but 842 records are missing PH and 246 records are missing HD. To facilitate comparison across traits, all phenotypes were standardized by subtracting the mean and dividing by the standard deviation of the phenotype vector (Table A1). Data were not standardized within environment to preserve all relationships within the data.

The population was genotyped with genotyping by sequencing (GBS; Elshire *et al.* 2011) markers aligned to the International Wheat Genome Sequencing Consortium (IWGSC) RefSeq v1.0 wheat genome sequence of ‘Chinese Spring’ (IWGSC 2018). Markers were filtered for minor allele frequency of at least 0.01 (Figure A1), no more than 30% missing scores, and no more than 10% heterozygous calls. Missing marker scores were imputed using categorical random forest imputation by chromosome, and all heterozygous calls (< 2% of all calls) were subsequently replaced with the population mode (*i.e.*, homozygous major allele). Marker scores are presented as boolean indicators of the minor allele. Further details of the CNLM dataset can be found in Appendix 1.

W-GY population: The W-GY wheat dataset of 599 historical wheat lines from the CIMMYT Global Wheat Breeding program (Crossa *et al.* 2010) was used in this study due to its importance in genomic prediction of an inbred population with non-additive variation (Crossa *et al.* 2010; Martini *et al.* 2016). The W-GY dataset consists of genotypic values of all lines for grain yield in each of four environments. The genetic correlations between these environments ranged from -0.19 to 0.66 and can be found in Martini *et al.* (2016). As performance between these environments is not highly correlated, we refer to grain yield performance in each environment as a trait. The dataset was used in its entirety with one exception. Of the 1,279 available DArT markers, only the 1,188 with known chromosomal positions as indicated by (Crossa *et al.* 2010) were utilized in this study. This information was required to know which markers belonged to which subgenome, such that subgenome specific relationship matrices could be calculated.

Statistical framework

Subgenome additive effects: To illustrate, we begin with a linear mixed model depicting environments (*i.e.*, trials) as fixed effects and genotypes as random

$$y_{ijk} = \mu + E_i + G_j + \varepsilon_{ijk} \quad (1)$$

where μ is the population mean, E_i and G_j are the fixed environmental and random genetic effects, respectively, of the j^{th} genotype evaluated in the i^{th} environment, and ε is the error associated with the k^{th} observation. Using matrix notation, model 1 (denoted G) can be rewritten

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{g}_G + \boldsymbol{\varepsilon} \quad (2)$$

where $\mathbf{1}$ is a vector of ones, \mathbf{X} is the design matrix of dummy variables for each trial, and $\boldsymbol{\beta}$ is the vector of fixed environmental effects. \mathbf{Z} is the incidence matrix linking observations in the vector \mathbf{y} to their respective genotype effects, in the vector \mathbf{g}_G . Normality was assumed for genotype effects and the residuals, where $\mathbf{g}_G \sim N(0, \sigma_G^2 \mathbf{K}_G)$ and

$\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I})$. The genetic covariance, \mathbf{K}_G , can be derived from the expectation (or coefficient) of co-ancestry between individuals from a pedigree (Henderson 1985), or by an empirical estimation of the realized genetic relationship calculated with genome-wide markers (Van Raden 2008). When genome-wide markers are used to estimate \mathbf{K}_G , the genomic prediction model initially suggested by Nejati-Javaremi *et al.* (1997) and Meuwissen *et al.* (2001) is obtained.

Given an $n \times m$ matrix, \mathbf{M} , of m markers scored as reference allele counts (*i.e.*, $\{0, 1, 2\}$) for n individuals, method I of Van Raden (2008) finds the genetic relationship \mathbf{K} as

$$\mathbf{K} = c^{-1}(\mathbf{M} - \mathbf{P})(\mathbf{M} - \mathbf{P})^T + 0.01\mathbf{I} \quad (3)$$

where $\mathbf{P} = \mathbf{1}_n \otimes 2\mathbf{p}^T$, $c = 2\mathbf{p}^T(\mathbf{1}_m - \mathbf{p})$ and \mathbf{p} is the vector of allele frequencies. The small coefficient of 0.01 was added to the diagonal to recover full rank after centering the matrix, such that \mathbf{K}_G is invertible.

We use allohexaploid wheat to illustrate, but this method is easily truncated to allotetraploids, or extended to higher level allopolyploids. If we allow the total genetic effect, G_j , to be decomposed into individual additive effects for each subgenome, such that $G_j = A_j + B_j + D_j$, the following model (ABD) is obtained

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{g}_A + \mathbf{Z}\mathbf{g}_B + \mathbf{Z}\mathbf{g}_D + \boldsymbol{\varepsilon} \quad (4)$$

In model 4, each subgenome is allowed to have its own additive genetic variance and covariance between individuals, such that $\mathbf{g}_A \sim N(0, \sigma_A^2 \mathbf{K}_A)$, $\mathbf{g}_B \sim N(0, \sigma_B^2 \mathbf{K}_B)$ and $\mathbf{g}_D \sim N(0, \sigma_D^2 \mathbf{K}_D)$. The realized additive genetic covariances for each subgenome, \mathbf{K}_A , \mathbf{K}_B and \mathbf{K}_D , are estimated using only markers corresponding to the respective subgenome, and calculated as described above.

Subgenome epistatic interactions: Following Henderson (1985), the epistatic covariance of individuals can be calculated as the Hadamard product of the component covariance matrices. Jiang and Reif (2015) and Martini *et al.* (2016) provide proofs of Henderson’s method using genome-wide markers to estimate the additive by additive covariance matrix, \mathbf{H} . An additional linear kernel can then be added for an additive by additive epistatic interaction term, I_j , once the additive covariance is estimated to obtain the following model (G \times G)

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{g}_G + \mathbf{Z}\mathbf{g}_I + \boldsymbol{\varepsilon} \quad (5)$$

where $\mathbf{g}_I \sim N(0, \sigma_{g_I}^2 \mathbf{H})$. As shown by Jiang and Reif (2015) and Martini *et al.* (2016), \mathbf{H} is calculated from the marker data as

$$\mathbf{H} = \mathbf{K} \odot \mathbf{K} - c^{-2}(\mathbf{W} \odot \mathbf{W})(\mathbf{W} \odot \mathbf{W})^T \quad (6)$$

where $\mathbf{W} = \mathbf{M} - \mathbf{P}$. Jiang and Reif (2015, Appendix A1) prove that \mathbf{H} is asymptotically equivalent to $\mathbf{K} \odot \mathbf{K}$ when the number of markers is large.

The additive by additive epistatic interaction term, \mathbf{g}_I , can also be decomposed into across subgenome interactions and within subgenome epistatic interactions such that $I_j = AB_j + AD_j + BD_j + I_j^-$, where AB_j , AD_j and BD_j are the subgenome interaction effects and I_j^- is the remaining epistatic effects due to within subgenome epistasis. Since no markers are shared across subgenomes, subgenome interaction covariances can be estimated by the Hadamard product of their component covariance matrices (Martini *et al.* 2016). These interactions can then be incorporated in the following model (ABD \times ABD)

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{g}_A + \mathbf{Z}\mathbf{g}_B + \mathbf{Z}\mathbf{g}_D + \mathbf{Z}\mathbf{g}_{AB} + \mathbf{Z}\mathbf{g}_{AD} + \mathbf{Z}\mathbf{g}_{BD} + \mathbf{Z}\mathbf{g}_{ABD} + \boldsymbol{\varepsilon} \quad (7)$$

where $\mathbf{g}_{AB} \sim \mathcal{N}(0, \sigma_{g_{AB}}^2 (\mathbf{K}_A \odot \mathbf{K}_B))$, $\mathbf{g}_{AD} \sim \mathcal{N}(0, \sigma_{g_{AD}}^2 (\mathbf{K}_A \odot \mathbf{K}_D))$, $\mathbf{g}_{BD} \sim \mathcal{N}(0, \sigma_{g_{BD}}^2 (\mathbf{K}_B \odot \mathbf{K}_D))$ and $\mathbf{g}_{ABD} \sim \mathcal{N}(0, \sigma_{g_{ABD}}^2 (\mathbf{K}_A \odot \mathbf{K}_B \odot \mathbf{K}_D))$. The three way interaction is included here for biological completeness, but was found to be estimated on the boundary (*i.e.*, zero) for all traits, and was therefore dropped from further analyses.

Accounting for population structure

Under Hardy Weinberg equilibrium, subgenomes segregate independently, such that for subgenome effects, $\text{Cov}(A, B) = \text{Cov}(A, D) = \text{Cov}(B, D) = 0$ and $\text{Var}(G) = \text{Var}(A) + \text{Var}(B) + \text{Var}(D)$. A breeding program, however, intentionally violates this assumption, and therefore may contain significant population structure. Price *et al.* (2006) demonstrated that the first k largest principal components (PCs) of the kinship matrix can be used to control for population structure in genome-wide association studies, and its use has since become wide spread. Because most realized estimates of additive covariance are proportional to $\mathbf{M}\mathbf{M}^T$, singular value decomposition of \mathbf{M} , instead of $\mathbf{M}\mathbf{M}^T$, can be used to separate the population structure as the first few principal components from the entire matrix of marker scores before it is divided into its subgenome components. This is accomplished by first extracting the first k principal components in the $n \times k$ matrix \mathbf{Q} . The marker matrix can then be reconstructed by setting the first k singular values of the diagonal matrix to zero and multiplying to produce a matrix $\tilde{\mathbf{M}}$ with the population structure removed from each subgenome simultaneously (Appendix 2)

Additive covariance matrices with reduced collinearity can then be calculated for each subgenome from $\tilde{\mathbf{M}}$ and incorporated into the model as previously described. \mathbf{Q} can then be added to the model as a set of fixed covariates, with slopes $\boldsymbol{\gamma}$, such that the model will now be of the form

$$\mathbf{y} = 1\boldsymbol{\mu} + \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{Q}\boldsymbol{\gamma} + \sum_l \mathbf{Z}\mathbf{g}_l + \boldsymbol{\varepsilon} \quad (8)$$

for all l genetic terms in the model. Genomic estimated breeding values are then predicted by summing the centered population structure and genetic effects. For this study, a population structure of dimension $k = 5$ was chosen for both the CNLM and W-GY datasets, and used to compare to the $k = 0$ models that do not correct for population structure. The number of PCs, $k = 5$, was chosen based on Supplementary Figure S1, where the correlation of additive covariance estimates between subgenomes appeared to level off in both populations.

Genomic prediction

To determine the predictability of genetic effects and the variability of variance component estimates, five-fold cross-validation was performed with 10 replications. For each replicate, the set of individuals was randomly split into five groups, with 4 groups of 289 and one of 291. For each fold, records of individuals in the fold were removed (*i.e.*, masked) from the dataset. Each model was subsequently fit with the remaining lines and used to predict the whole genome effect of the masked lines in the fold. Predictions for all five folds were gathered and correlated to the “true” genetic values once for each replicate. In this way, prediction results are directly comparable between the different models, and not subject to differences in the individuals sampled. The whole genome values were calculated as the sum of the genotypic additive and epistatic effects in the model as previously described. Due to the unbalanced nature of the CNLM dataset, “true” genetic values were calculated as in equation 2 but were considered independent with a covariance $\mathbf{K}_G = \mathbf{I}$.

Software

Models were fit using restricted maximum likelihood (REML) with the software ASReml (Gilmour 1997) implemented in R (Butler 2009). The Tassel 5.0 GBS pipeline v2 (Glaubitz *et al.* 2014) along with the ‘bwa’ alignment tool (Li and Durbin 2009) were used for aligning GBS markers to the reference genome. All additional computation, analyses and figures were made using base R (R Core Team 2015) implemented in the Microsoft Open R environment 3.3.2 (Microsoft 2017).

Data availability

Phenotypes for the CNLM population are included in the file ‘pheno.txt’. Marker information and imputed marker scores for the CNLM population are included in files ‘snpInfo.txt’ and ‘snpMatrix.txt’, respectively. Best Linear Unbiased Predictors (BLUPs) for whole and subgenome additive effects (GEBVs and SGEBVs, respectively), as well as non-additive whole and subgenome interaction effects can be found in the ‘effectTable.txt’ file. Genotype and phenotype data for the W-GY population can be found in the ‘BGLR’ package of R (de los Campos and Pérez Rodriguez 2015), and marker chromosome information can be found in Crossa *et al.* (2010). Supplemental material available at Figshare: <https://doi.org/10.25387/g3.6870110>.

RESULTS

Model fit and variance components

Model fit was assessed using Akaike’s Information Criterion (AIC). Whole genome models tended to have the lowest AIC values, with the exception of the PH and HD traits for the epistatic ABD×ABD models in the CNLM population. When whole genome models had lower AIC values, the comparable subgenome models had only marginally higher AIC values (Supplementary Tables S1 and S2). Whole genome predictions between comparable whole genome and subgenome models were correlated at $\rho > 0.999$ or $\rho > 0.993$ for traits within the CNLM and W-GY populations, respectively. This indicates that little, if any, genetic information was lost by splitting the whole genome into biologically relevant subgenome effects. The lack of perfect correlation is at least partially due to floating point rounding errors during model fitting and summation of genotype effects.

Subgenome additive variance parameter estimates were positive for all models, but subgenome interaction variance parameter estimates were often estimated on the boundary (*i.e.*, near zero). Variance parameters estimated on the boundary were thus considered to be exactly zero. Shifts in variance component importance were seen when the epistatic terms were added in the model. For example, for the TW and E1 traits in the CNLM and W-GY populations, respectively, the A genome component was the largest in the additive only model, but was reduced to less than that of the B genome component in the epistatic model. Additive variance components were generally reduced in epistatic models compared to additive only models, but this reduction in additive variance was accompanied by non-zero subgenome interaction components. The B genome contributed the greatest amount of additive variance in the epistatic ABD×ABD models for all traits except HD. While the D genome variance component was far smaller than the A subgenome component for GY and TW in the CNLM population, it was comparable to the A subgenome component for all traits in the W-GY population.

The A×B component was particularly important for the W-GY traits, E1, E3 and E4, as well as the HD and TW in the CNLM population. The A×D component also featured prominently for the PH and TW traits in the CNLM population. The B×D component appeared to be less important, having the largest effect for PH. No epistatic terms

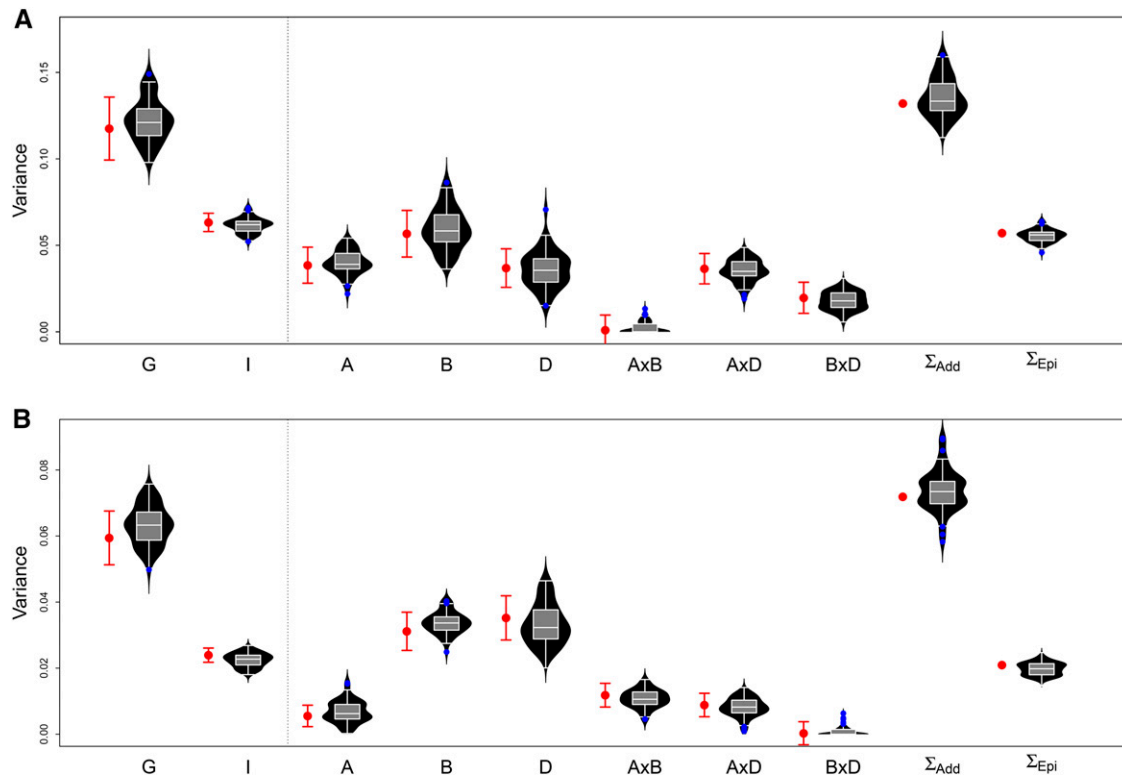


Figure 1 Estimates and standard errors of variance components from the full model (red) compared to the sampling distribution of variance component estimates from the cross-validation scheme (black violins). Two traits from the CNLM population, (A) PH and (B) HD, with contrasting genetic architectures are shown. GxG (5) and ABDxABD (7) models are shown to the left and right of the dotted line, respectively. The sum of the additive (Σ_{Add}) and epistatic (Σ_{Epi}) variance components is also shown for the ABDxABD model.

were significantly greater than zero for the E2 trait in the W-GY population. Addition of epistatic interactions resulted in a significant likelihood ratio test at $p < 10^{-6}$ for all traits except GY, which was significant at $p < 10^{-2}$. Despite the significant addition of epistatic terms, additive GEBVs were highly correlated with whole genome predictions from the epistatic models, at $\rho \geq 0.988$ for the CNLM population and $\rho \geq 0.869$ for the W-GY population. A model containing the three-way subgenome epistatic term was fit for all traits, but estimates of the three-way interaction variance parameter were zero for all traits.

The distributions of variance component estimates from repeated sub-sampling of the data during k -fold cross-validation were centered near the point estimate from the full model fit. These distributions were either as wide (≈ 2 standard errors from the center) or tighter than expected based on the standard error from the full model fit (Figure 1, Supplementary Figures S2 and S3). Standard errors were generally larger for epistatic variance components relative to their magnitude than additive variance components. Standard errors relative to their respective parameter estimates tended to be larger for all terms in models with more estimated variance parameters (Supplementary Tables S1 and S2).

Subgenome additive effects

Subgenome estimated breeding values (SGBEVs) were moderately correlated with the whole genome effect, but weakly correlated with one another (Supplementary Tables S3 and S4). The individuals with the highest SGBEV for one subgenome never had the highest SGBEV for the other two subgenomes, and were often not in the top 95% quantile of the population based on the other two subgenomes

(Figure 2, Supplementary Figures S4 and S5). For example, the individuals with the highest A, B and D SGBEV for GY in the CNLM population ranked 43rd, 39th and 60th for the whole genome effect, respectively. In contrast, the individual with the best A SGBEV for GY ranked 1067th, 952nd for the B and D genome, respectively. The individual with the highest B genome breeding value for GY ranked 221th and 1393rd for the A and D genomes, respectively. The individual with the highest D genome breeding value for GY ranked 347th and 123rd for the A and B subgenome, respectively. The individual with the highest whole genome GEBV for GY ranked 6th, 22nd and 519th for the A, B and D SGBEV, respectively. In several cases, the top individual based on a SGBEV was not in the top 95% quantile based on their whole genome effect, particularly in the W-GY population.

Prediction accuracy

Including epistasis kernels significantly improved genomic prediction accuracy for all traits except GY and E2 (Table 1). Subgenome models had either comparable or slightly lower mean prediction accuracy than whole genome models for all traits except HD, for which subgenome models had superior accuracy. The variability in the prediction accuracy based on the individuals sampled was either the same (GY and TW) or lower (PH and HD) for the epistatic models compared to the additive models in the CNLM population, but was similar in the W-GY population (Table 1). The variability in prediction accuracy was increased for the subgenome models compared to the whole genome models in the W-GY population for some traits (E2 and E3), but was either the same or decreased in the the CNLM population.

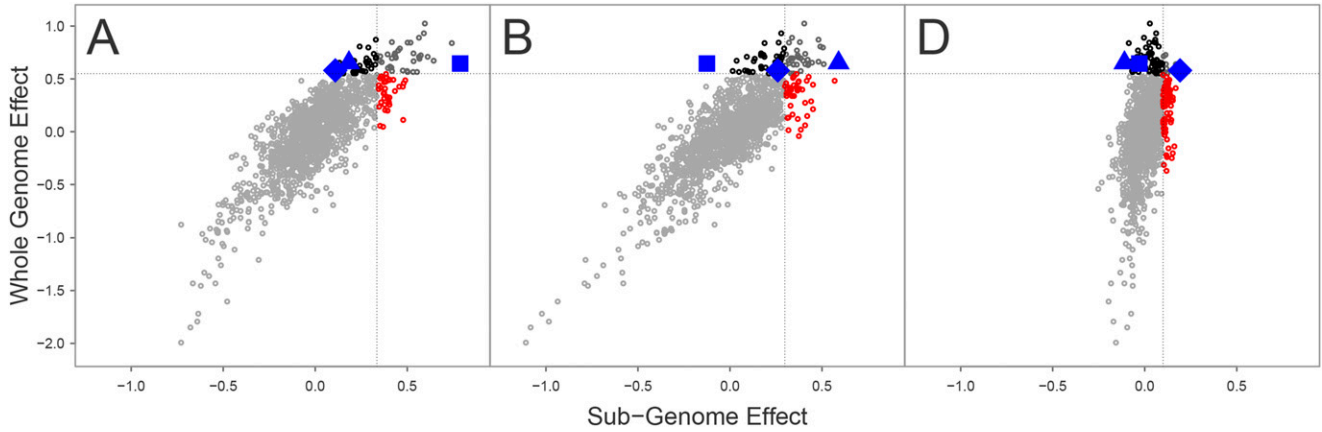


Figure 2 Plot of whole genome additive effects (GEBV) by subgenome additive effects (SGEBV) for GY in the CNLM populations. The dotted line indicates the 95% quantiles for whole or subgenome effects. Blue squares, triangles and diamonds indicate the line with the highest SGEBV for each of the A, B and D subgenomes, respectively.

Adjustment for population structure

The first two principal components explained 17% and 19% of the variance of \mathbf{M} in the CNLM and W-GY populations, respectively, indicating that some population structure exists in both populations (Supplementary Figure S6). The correlation of additive genetic covariance estimates between individuals based on the three subgenomes declined as PCs were removed from \mathbf{M} , but appeared to level out between 5 to 10 PCs (Supplementary Figure S1). Correlation of whole genome effects between additive models, G and ABD, for $k = 0$ and $k = 5$ was ≥ 0.999 and ≥ 0.996 for the CNLM and W-GY populations respectively. Whole genome effect correlations were lower between epistatic models $G \times G$ and $ABD \times ABD$, with coefficients of ≥ 0.998 in the CNLM population and ≥ 0.980 in the W-GY population.

Removing population structure with $k = 5$ reduced most of the SGEBV effect correlation coefficients by up to 0.06 in the CNLM population, but there was one instance in which one correlation coefficient increased from 0.15 to 0.19 between A and B SGEBVs for PH (Supplementary Tables S3 and S4). This was not the case for the W-GY population, where many of the SGEBV effect correlations increased by up to 0.21.

Additive variance generally decreased as k was increased from 0 to 10 (Figure 3, Supplementary Figure S7). Ranks of additive variance components relative to one another were stable for most traits, while epistatic variance components were more sensitive to changes in k . Significant epistatic variance component rank changes occurred for the PH, TW and E4 traits. For PH, the $A \times D$ term was comparable in magnitude with the additive variance components for A and D when $k = 0$, but declined as k increased. The reduction in $A \times D$ variance for PH was accompanied by an increase in both the $A \times B$ and $B \times D$ terms. Similarly, a decline in $A \times B$ variance was followed by an increase in $B \times D$ for TW and $A \times D$ for E4 as k was increased.

Correlations of variance component estimates were calculated from the average information matrix for models $k \in \{0, 1, \dots, 10\}$ (Supplementary Figures S8 and S9). Correlations between subgenome additive variance estimates were generally low (0.2-0.4), while correlations of subgenome interaction variance estimates were high (0.8-0.95), and correlations between the two were moderate (0.4-0.6). Despite a small reduction in the correlation of SGEBVs as k was increased from 0 to 5, little reduction in variance component estimate correlations was observed as k was increased from 0. Generally, correlations of additive

variance parameter estimates were slightly reduced while correlations between interaction variance parameter estimates increased slightly.

DISCUSSION

Model fit and variance components

While whole genome models tended to be the most parsimonious, subgenome models are worth consideration because they provide insight into the biology of the allopolyploid organism. Given the stability of variance component estimation and that no genetic information appears to be lost by partitioning the whole genome into its individual subgenome additive effects, such a partition is informative.

The method presented here could be used for any set of independent loci, such as estimating a variance component and breeding value for each chromosome (Yang *et al.* 2011). However, this will become computationally burdensome as the number of variance components to be estimated increases. If the number of variance parameters to estimate is high and the data set is small this may become infeasible. It is also unclear if the estimates from larger numbers of additive kernels would be reliable.

Bernardo and Thompson (2016) assigned a breeding value for each of the 10 maize chromosomes by fitting a single ridge regression model to estimate marker effects. They subsequently summed marker effects by chromosome to produce a breeding value for each chromosome. However, this method does not allow for direct estimation of variance components for each unit of chromatin. By fitting each unit simultaneously, variance attributable to sets of loci will be split, and the sum of the variance estimates should not exceed the total genetic variance (Yang *et al.* 2011). It is unclear what effect linkage disequilibrium across chromosomes has on the variance parameters estimated.

Here we assumed that the subgenome effects are independent, but this is clearly not the case. Generally, we can express the genetic variance due to the three subgenomes as

$$\text{Var} \begin{pmatrix} \mathbf{g}_A \\ \mathbf{g}_B \\ \mathbf{g}_D \end{pmatrix} = \mathbf{S} \otimes \mathbf{J}_n \odot \mathbf{K} \quad (9)$$

where \mathbf{S} is the subgenome covariance matrix, \mathbf{J} is an $n \times n$ matrix of ones for n genotypes, and \mathbf{K} is the additive relationship matrix for within and across subgenomes. In this report, we have assumed that \mathbf{S} is diagonal with $\mathbf{S}_{ii} = \sigma_i^2$ for the i^{th} subgenome, and \mathbf{K} is a block

■ **Table 1** Table of genomic prediction accuracies for eight traits in the CNLM (GY, PH, TW and HD) or W-GY (E1, E2, E3, E4) populations with $k = 0$ and $k = 5$. k is the number of principal components removed from the marker matrix prior to calculating subgenome covariance matrices. The first 1 to k principal components were included as fixed effects in the model fit for $k > 0$

CNLM	k	GY	PH	TW	HD
G (2) ^a	0	0.601 ^b (0.008) ^c	0.559 (0.007)	0.515 (0.010)	0.664 (0.009)
ABD (4)		0.600 (0.008)	0.557 (0.008)	0.514 (0.011)	0.679 (0.007)
G×G (5)		0.604 (0.008)	0.637 (0.004)	0.576 (0.010)	0.712 (0.008)
ABD×ABD (7)		0.603 (0.008)	0.638 (0.005)	0.569 (0.011)	0.720 (0.006)
G	5	0.600 (0.009)	0.558 (0.007)	0.514 (0.011)	0.663 (0.010)
ABD		0.600 (0.009)	0.556 (0.008)	0.513 (0.011)	0.678 (0.008)
G×G		0.602 (0.008)	0.624 (0.005)	0.560 (0.010)	0.701 (0.008)
ABD×ABD		0.602 (0.007)	0.618 (0.005)	0.557 (0.010)	0.708 (0.006)
W-GY	k	E1	E2	E3	E4
G	0	0.501 (0.010)	0.493 (0.016)	0.356 (0.008)	0.457 (0.010)
ABD		0.492 (0.012)	0.481 (0.023)	0.346 (0.010)	0.449 (0.011)
G×G		0.568 (0.010)	0.494 (0.017)	0.396 (0.013)	0.520 (0.010)
ABD×ABD		0.549 (0.011)	0.484 (0.023)	0.393 (0.015)	0.509 (0.013)
G	5	0.502 (0.010)	0.491 (0.017)	0.354 (0.007)	0.458 (0.010)
ABD		0.495 (0.011)	0.475 (0.024)	0.345 (0.010)	0.453 (0.011)
G×G		0.526 (0.010)	0.491 (0.017)	0.381 (0.007)	0.493 (0.011)
ABD×ABD		0.520 (0.012)	0.475 (0.023)	0.373 (0.013)	0.486 (0.012)

^aEquation.

^bMean genomic prediction accuracy across ten replicates of five fold cross validation.

^cStandard deviation of prediction accuracy across ten replicates are shown in parentheses.

diagonal with the i^{th} diagonal block represented by the subgenome additive covariance matrix.

$$S = \begin{bmatrix} \sigma_A^2 & 0 & 0 \\ 0 & \sigma_B^2 & 0 \\ 0 & 0 & \sigma_D^2 \end{bmatrix} \text{ and } K = \begin{bmatrix} K_A & 0 & 0 \\ 0 & K_B & 0 \\ 0 & 0 & K_D \end{bmatrix} \quad (10)$$

An unstructured covariance matrix, S , could be estimated, with correlation coefficients between subgenomes. The subgenome effects would be allowed to have a correlation such that

$$S = \begin{bmatrix} \sigma_A^2 & \sigma_{AB} & \sigma_{AD} \\ \sigma_{AB} & \sigma_B^2 & \sigma_{BD} \\ \sigma_{AD} & \sigma_{BD} & \sigma_D^2 \end{bmatrix} \text{ and } K = \begin{bmatrix} K_A & K_{AB} & K_{AD} \\ K_{AB} & K_B & K_{BD} \\ K_{AD} & K_{BD} & K_D \end{bmatrix} \quad (11)$$

However, it is unclear what the covariance structure should be between subgenomes (e.g., K_{AB}). If consensus haplotypes from uniquely identifiable sequences could be determined with two or more alleles segregating in at least two subgenomes, a covariance across the subgenomes could be constructed. Polymorphisms that predate speciation would be used to identify the consensus haplotypes, while post speciation polymorphisms would be used to identify the subgenome origin. Individuals would then receive a score based on the number of consensus haplotypes they have in common between two subgenomes. This could prove to be a formidable challenge given the evolutionary time between the subgenome ancestors. The Hadamard product of the two additive covariance matrices is a tempting candidate for these off diagonal blocks, however, this would substitute a covariance between additive effects in place of an epistatic variance. It is unclear to these authors if epistasis variance can be thought of and modeled as a covariance between additive effects.

Genetic architecture

The genetic architecture of grain yield (GY, E1, E2, E3, E4) in the two populations investigated here are markedly similar, despite the divergent genetic backgrounds of the two populations. The CNLM population is

primarily comprised of breeding lines and varieties derived from germplasm historically grown in the North East, in contrast to the W-GY population which has a broader pedigree.

The D genome is known to have low genetic diversity due to limited gene flow from a single *Ae. tauschii* lineage after the most recent allopolyploidization event (Wang *et al.* 2013), estimated to have taken place as recently as 10,000 years ago (Salamini *et al.* 2002; Marcussen *et al.* 2014). The International Maize and Wheat Improvement Center (Centro Internacional de Mejoramiento de Maíz y Trigo, CIMMYT) has introgressed some genetic material from the D genome ancestor, *Ae. tauschii*, through the use of synthetically produced hexaploid wheat to increase the genetic diversity of the historically bottle-necked D genome. The higher proportion of D genome variance in the W-GY population may be due to the increased use of wild *Ae. tauschii* in their breeding program, highlighting the merit of the strategy.

Many of the subgenome epistatic variance parameters were estimated at zero, possibly due to a lack of power to detect them. Greater genetic diversity, larger numbers of individuals, and higher allele frequencies would allow for increased power to detect true interactions. Hill *et al.* (2008) emphasized the effect of low allele frequencies on epistatic interactions, proving that as allele frequencies (and therefore joint frequencies of alleles at two loci) approach zero or one, most of the epistatic variance becomes additive. For example, suppose two loci have a large interaction, such that one pair of alleles is selected. Once one locus becomes fixed, all remaining variance is due to segregation at the other locus, and becomes strictly additive. The low joint frequency is magnified in the three way interactions, likely causing the inability to detect any three way epistatic interaction signal between the three subgenomes.

This is also apparent in the reduction of additive variance components upon the addition of epistatic terms to the model. These epistatic components were often estimated to be rather large compared to the additive components, but did not change the final whole genome value drastically. This suggests that the additive terms absorb much of the epistatic variance in the absence of the epistatic kernels.

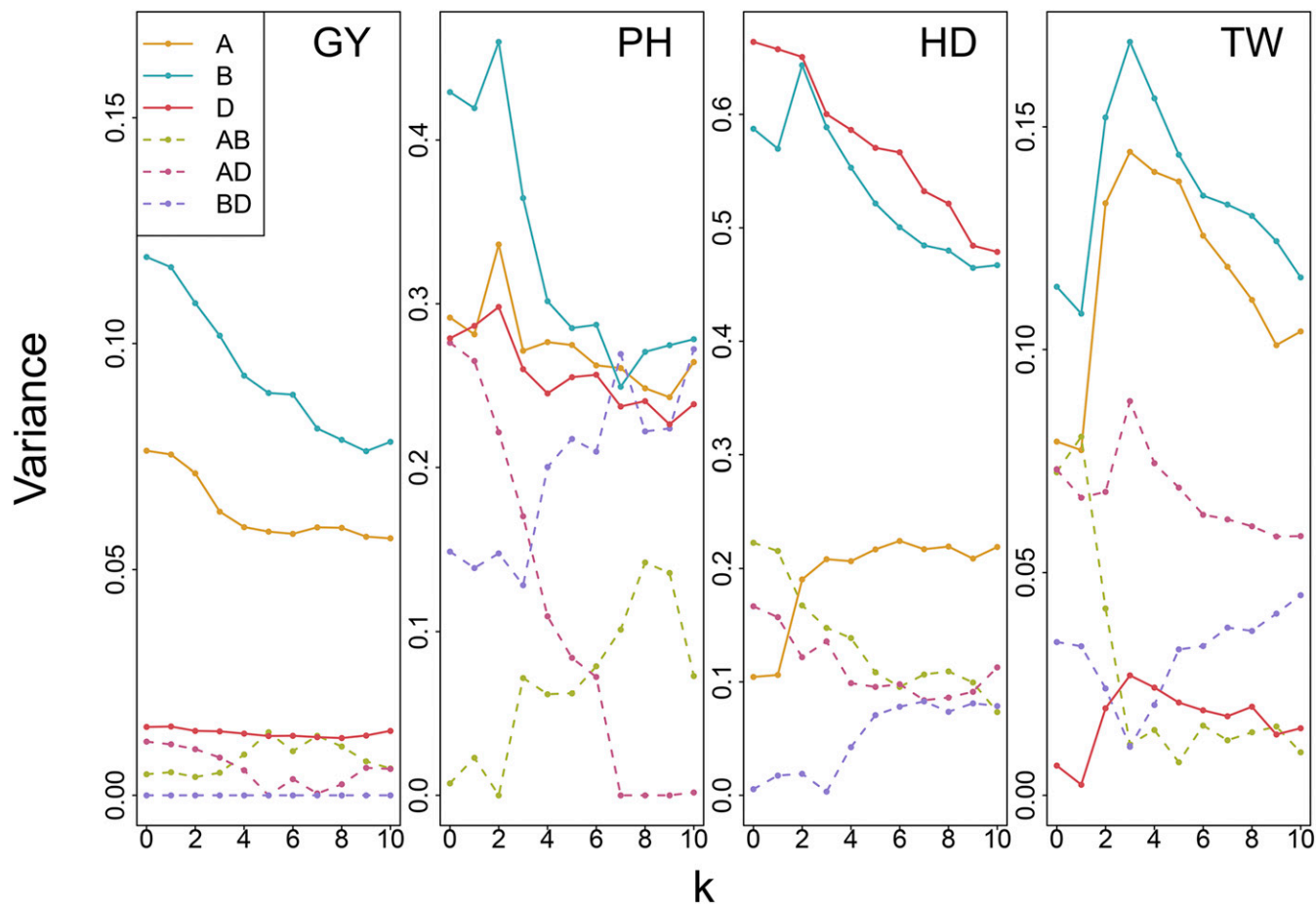


Figure 3 Subgenome additive and interaction variance parameter estimates from the ABD×ABD (7) model correcting for population structure with $k \in \{0, 1, \dots, 10\}$ principal components as fixed effects. Models were fit with four traits for the CNLM population.

The A×B epistatic terms were the most important for many of the traits, reflecting the greater genetic variation of these two subgenomes. Subgenome interaction terms including the D genome were notably more important for traits known to have important loci on the D genome. PH is partially governed by two dwarfing genes, *Rht-1D* and *Rht-1B* on 4B and 4D, respectively. These two genes have been shown to exhibit a less-than-additive (Eshed and Zamir 1996) epistatic interaction, where the double wildtype is less tall than expected based on the additive effects of the two semi dwarfs from the double dwarf (Santantonio *et al.*, 2018b). The B×D term was large for PH, particularly after correction for population structure. Population structure is common for these genes, as many breeding programs primarily utilize one or the other dwarfing gene to avoid producing double dwarfs during crossing, which are often agronomically undesirable.

Selection on SGEbVs

Partitioning genetic variance to the subgenomes of an allopolyploid provides a method for identifying individuals with complementary subgenomes as potential parents for crossing. If we consider the upper 95% quantiles as candidates for parental selection, many of the top candidates based on subgenome breeding values would not be considered candidates based on their whole genome breeding value. When they would be considered, they were typically not the top candidates. The low correlation between SGEbVs highlights the opportunity to identify individuals with complementary subgenomes for crossing.

These individuals may or may not be among the top performing selection candidates, demonstrating that the optimum set of crosses are not always between the top performing individuals (Akdemir and Sánchez 2016).

While it is difficult to evaluate the predictability of SGEbVs *per se*, the sum of their values appears highly predictive. The low correlation of SGEbVs also suggests that individual subgenomes may be directly manipulated as never before. Prior to the discovery and use of genetic markers to track genomic regions, the phenotype (or some summary statistic thereof) was the only indicator of the genetic structure of a genetically distinct individual. Variety releases still demonstrate this legacy, with phenotypic descriptors that define a new variety as genetically distinct from other similar varieties. One breeding strategy will be selecting parents for crossing that have complementary SGEbVs to increase the potential of transgressive segregation in the resulting offspring. We envision other breeding strategies beyond simply choosing parents with complementary subgenomes, and see an opportunity to weight SGEbVs according to some breeding goal.

For example, a newly formed population could undergo several rounds of genomic selection only on the D genome SGEbVs (*i.e.*, weights of 0, 0 and 1 for the A, B and D subgenomes respectively) before phenotypic or whole genome selection. Because the D genome contributes the least to the total genetic variance, phenotypic selection on D genome loci is challenging. Selection will act on the largest sources of genetic variance first, potentially leading to fixation of small effect loci

in the D genome by drift, while selection acts on the large effect loci on the A and B genomes first. By selecting on D genome SGBEVs, gains can be made to the D genome directly with little to no selection on the A and B genomes, a feat previously impossible with phenotypic selection. Signatures of selection on the D genome under this scheme may also help establish the accuracy of SGBEV prediction.

Subgenome interactions

Genomic prediction of GY and E2 did not appear to benefit from including epistatic interactions as it did for the other six traits. This may be due in part to the highly polygenic nature of grain yield, which is the culmination of essentially all functional genetic variants subjected to stress throughout the growth cycle. The E2 trait in the W-GY population has previously been shown to be invariant to the addition of various epistatic terms (Crossa *et al.* 2010; Martini *et al.* 2016), and it is unclear why this population does not exhibit non-additive variation in this environment. It may be that important epistatic interactions of GY in the CNLM population are too small to detect or are involved with differing performance across years or locations, such that they are lacking in a model that does not include genotype by environment interactions.

Subgenome epistatic terms increased genomic prediction accuracy comparable to modeling all pairwise interactions across the subgenomes, suggesting that the most important interacting loci are on different subgenomes. This result is consistent with the observation that newly formed allopolyploids undergo considerable changes in gene expression, known as genome shock (McClintock 1984). This shock has been suggested to be caused by incompatibilities of genetic pathways across the subgenomes (Comai *et al.* 2003). Residual subgenome incompatibility may still be affecting the germplasm pool, even thousands of years after the last polyploidization event. Decay of negative duplicated gene interactions may take hundreds or thousands of generations before all interacting genes are lost or silenced (Lynch and Conery 2000, 2003).

It is unclear what proportion of this non-additive signal is due to homeoallelic interactions. The proposed method models all pairwise interactions across subgenomes, of which homeoallelic gene interactions are a small minority in number. Smaller homeoallelic regions (Santantonio *et al.* 2018a) or homeoallele specific marker sets (Santantonio *et al.* 2018b) have been constructed to determine the relative importance of these interactions relative to other gene interactions across the subgenomes. The usefulness of the epistatic subgenome interactions is currently unclear and warrants further investigation.

Regardless of the source of the epistasis, we suggest that a breeding scheme should be designed to take advantage of beneficial subgenome interactions. If a suitable training set related to the breeding material can be established, subgenome interactions can be predicted in new, genotyped breeding materials. We suggest that a series of small bi-parental populations be constructed from important contributors to the breeding program, and be used in the development of a training population to balance high genetic diversity and high allele frequencies. This training population would be used to predict SGBEVs and subgenome interactions in individuals formed from new crosses. Individuals that contain favorable interactions could then be selected such that they are fixed in early filial generations. After fixation in a given line, phenotypic, whole genome, or subgenome selection could be used for further line development until complete homozygosity is reached.

Adjustment for population structure

The efficacy of the proposed method to handle population structure may need to be improved, or a different approach may need to be taken. While this method reduced the correlation of additive genetic subgenome

covariance estimates across the three genomes, variance parameter estimate correlations were not drastically reduced. The correlation of subgenome interaction variance parameter estimates tended to increase slightly when accounting for higher levels of population structure, counter to the assumption that removing this structure should result in better estimates of subgenome interactions.

The lower correlation between epistatic models that correct and do not correct for population structure is likely due to removing Q from the marker matrix. Correcting for population structure also had a small, but negative effect on genomic prediction accuracy for epistatic models. The population structure fixed effect predictors are strictly additive and the loss of accuracy may be due to epistasis variance associated with these PCs (*i.e.*, population structure epistasis). Epistatic variance related to these PCs may be recovered by using the squares of the PC scores, although this was not done in this study. At least for the additive models, it appears little to no genetic information is lost using the population structure adjustment proposed here.

Determining the best value for k will be at the crux for implementing this methodology for various traits and populations. The same population may need different values of k for different traits, depending on the covariance of the trait and population structure. Traits such as PH or HD may have lower covariance with population structure than traits such as TW or GY, due to different marker effect distributions and the history of the breeding population. Several methods might be used to determine k empirically from the marker matrix (Patterson *et al.* 2006, *e.g.*), however, these methods do not account for the covariance of trait and structure. We used the first one to k PCs in this study, but there is no reason why we must include all PCs up to some value k . There may be certain PCs that are important for a given trait, and could be tested as fixed effects for inclusion or exclusion.

This method may have better performance in populations with greater degrees of population structure than in the populations presented here. Use of this method for partitioning genetic variance to biologically important sets of chromatin and estimating epistatic interactions will need further testing and validation before widespread use.

CONCLUSION

To our knowledge, we provide the first attempt to assign a breeding value to each subgenome of an allopolyploid crop. With estimates of subgenome additive effects, parents with complementary subgenomes can be selected for crossing. Weighted selection of subgenomes using genomic prediction could be key to increasing the diversity of the D genome in wheat germplasm. Direct selection on the D genome may allow targeted introgression from *Ae. tauschii* while mitigating the effects of introducing unimproved alleles. Subgenome additive genetic variances appear to be estimated well, and no genetic information is lost partitioning the genome into its subgenome components. This demonstrates that partitioning genetic variance to the subgenomes of an allopolyploid can provide useful information for genomics assisted breeding efforts.

Subgenome interactions increase prediction accuracy, but it is unclear how well the epistatic variance is partitioned to the three interaction terms and what proportion of that variance is due to homeologous gene interactions. Because the homeologous interactions make up relatively few of the possible interactions across subgenomes, they may only explain a small portion of the observed epistatic variance. Yet, seeing as how homeologous genes likely operate in the same or similar physiological pathway, the likelihood for interactions between homeologous loci is high. Further research is needed to investigate the efficacy of modeling subgenome interaction terms, and to what degree this is explained by interactions between homeologous orthologs.

Allopolyploids have traditionally been treated as diploids in breeding programs because they undergo disomic inheritance. With modern DNA marker technology and ever increasing computational power, breeders of allopolyploids can further exploit the genetic complexity of their crops.

ACKNOWLEDGMENTS

Funding of this research was provided by the USDA National Needs Fellowship for N. Santantonio, in partial fulfillment of the requirements for a Ph.D in Plant Breeding and Genetics at Cornell University. Additionally, the field trials comprising the phenotypic data for the CNLM population were funded in part by the Hatch Project # 149-447. Genotyping was funded by the Wheat Coordinated Agricultural Project (WheatCAP). We are also grateful to Jesse Poland's research group at Kansas State University for their contribution to genotyping of CNLM materials. The authors thank the International Wheat Genome Sequencing Consortium for pre-publication access to IWGSC RefSeq v1.0. Finally, we would like to acknowledge the Cornell small grains staff, particularly David Benschler and James Tanaka, who were vital in implementing, collecting and processing the materials used to build the CNLM dataset.

Note added in proof: See Santantonio *et al.* 2019 (pp. 675–684) in this issue and (pp. 1105–1122) in the Genetics March 2019 issue, for related works.

LITERATURE CITED

- Abel, S., C. Möllers, and H. Becker, 2005 Development of synthetic brassica napus lines for the analysis of “fixed heterosis” in allopolyploid plants. *Euphytica* 146: 157–163. <https://doi.org/10.1007/s10681-005-3364-7>
- Akdemir, D., and J.-L. Jannink, 2015 Locally epistatic genomic relationship matrices for genomic association and prediction. *Genetics* 199: 857–871. <https://doi.org/10.1534/genetics.114.173658>
- Akdemir, D., J.-L. Jannink, and J. Isidro-Sánchez, 2017 Locally epistatic models for genome-wide prediction and association by importance sampling. *Genet. Sel. Evol.* 49: 74. <https://doi.org/10.1186/s12711-017-0348-8>
- Akdemir, D., and J. I. Sánchez, 2016 Efficient breeding by genomic mating. *Front. Genet.* 7: 210. <https://doi.org/10.3389/fgene.2016.00210>
- Bernardo, R., and A. M. Thompson, 2016 Germplasm architecture revealed through chromosomal effects for quantitative traits in maize. *Plant Genome* 9: 1–11. <https://doi.org/10.3835/plantgenome2016.03.0028>
- Bingham, E., R. Groose, D. Woodfield, and K. Kidwell, 1994 Complementary gene interactions in alfalfa are greater in auto-tetraploids than diploids. *Crop Sci.* 34: 823–829. <https://doi.org/10.2135/cropsci1994.0011183X003400040001x>
- Birchler, J. A., H. Yao, S. Chudalayandi, D. Vaiman, and R. A. Veitia, 2010 Heterosis. *Plant Cell* 22: 2105–2112. <https://doi.org/10.1105/tpc.110.076133>
- Blake, N. K., B. R. Lehfeldt, M. Lavin, and L. E. Talbert, 1999 Phylogenetic reconstruction based on low copy dna sequence data in an allopolyploid: the b genome of wheat. *Genome* 42: 351–360. <https://doi.org/10.1139/g98-136>
- Butler, D., 2009 *asreml: asreml() fits the linear mixed model*. R package version 3.0.
- Cockerham, C. C., and Z.-B. Zeng, 1996 Design iii with marker loci. *Genetics* 143: 1437–1456.
- Comai, L., A. Madlung, C. Josefsson, and A. Tyagi, 2003 Do the different parental ‘heteromes’ cause genomic shock in newly formed allopolyploids? *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 358: 1149–1155. <https://doi.org/10.1098/rstb.2003.1305>
- Crossa, J., G. de Los Campos, P. Pérez, D. Gianola, J. Burgueño *et al.*, 2010 Prediction of genetic values of quantitative traits in plant breeding using pedigree and molecular markers. *Genetics* 186: 713–724. <https://doi.org/10.1534/genetics.110.118521>
- de los Campos, G. and P. Pérez Rodríguez, 2015 *BGLR: Bayesian Generalized Linear Regression*. R package version 1.0.4.
- Duvick, D. N., 1999 Heterosis: feeding people and protecting natural resources. *Genetics and Exploitation of Heterosis in Crops*, Coors J.G., Pandey S., eds Madison, WI: American Society of Agronomy and Crop Science Society of America, pp. 19–29.
- Dvořák, J., P. Terlizzi, H.-B. Zhang, and P. Resta, 1993 The evolution of polyploid wheats: identification of the a genome donor species. *Genome* 36: 21–31. <https://doi.org/10.1139/g93-004>
- Eckart, C., and G. Young, 1936 The approximation of one matrix by another of lower rank. *Psychometrika* 1: 211–218. <https://doi.org/10.1007/BF02288367>
- Ellstrand, N. C., and K. A. Schierenbeck, 2000 Hybridization as a stimulus for the evolution of invasiveness in plants? *Proc. Natl. Acad. Sci. USA* 97: 7043–7050. <https://doi.org/10.1073/pnas.97.13.7043>
- Elshire, R. J., J. C. Glaubitz, Q. Sun, J. A. Poland, K. Kawamoto *et al.*, 2011 A robust, simple genotyping-by-sequencing (gbs) approach for high diversity species. *PLoS One* 6: e19379. <https://doi.org/10.1371/journal.pone.0019379>
- Eshed, Y., and D. Zamir, 1996 Less-than-additive epistatic interactions of quantitative trait loci in tomato. *Genetics* 143: 1807–1817.
- Feldman, M., A. A. Levy, T. Fahima, and A. Korol, 2012 Genomic asymmetry in allopolyploid plants: wheat as a model. *J. Exp. Bot.* 63: 5045–5059. <https://doi.org/10.1093/jxb/ers192>
- Fisher, R. A., 1919 Xv—the correlation between relatives on the supposition of mendelian inheritance. *Earth and Environmental Science Transactions of the Royal Society of Edinburgh* 52: 399–433. <https://doi.org/10.1017/S0080456800012163>
- Gage, J. L., D. Jarquin, C. Romay, A. Lorenz, E. S. Buckler *et al.*, 2017 The effect of artificial selection on phenotypic plasticity in maize. *Nat. Commun.* 8: 1348. <https://doi.org/10.1038/s41467-017-01450-2>
- Garrick, D., 2007 Derivation of equivalent computing algorithms for genomic predictions and reliabilities of animal merit. *J. Dairy Sci.* 90: 376.
- Gilmour, A., 1997 Asreml for testing fixed effects and estimating multiple trait variance components. *Proceedings of the Association for the Advancement of Animal Breeding and Genetics* 12: 386–390.
- Glaubitz, J. C., T. M. Casstevens, F. Lu, J. Harriman, R. J. Elshire *et al.*, 2014 Tassel-gbs: a high capacity genotyping by sequencing analysis pipeline. *PLoS One* 9: e90346. <https://doi.org/10.1371/journal.pone.0090346>
- Goddard, M., and B. Hayes, 2007 Genomic selection. *J. Anim. Breed. Genet.* 124: 323–330. <https://doi.org/10.1111/j.1439-0388.2007.00702.x>
- Haldane, J., 1933 The part played by recurrent mutation in evolution. *Am. Nat.* 67: 5–19. <https://doi.org/10.1086/280465>
- Heffner, E. L., M. E. Sorrells, and J.-L. Jannink, 2009 Genomic selection for crop improvement. *Crop Sci.* 49: 1–12. <https://doi.org/10.2135/cropsci2008.08.0512>
- Henderson, C., 1985 Best linear unbiased prediction of nonadditive genetic merits in noninbred populations. *J. Anim. Sci.* 60: 111–117. <https://doi.org/10.2527/jas1985.601111x>
- Heslot, N., J.-L. Jannink, and M. E. Sorrells, 2015 Perspectives for genomic selection applications and research in plants. *Crop Sci.* 55: 1–12. <https://doi.org/10.2135/cropsci2014.03.0249>
- Hill, W. G., M. E. Goddard, and P. M. Visscher, 2008 Data and theory point to mainly additive genetic variance for complex traits. *PLoS Genet.* 4: e1000008. <https://doi.org/10.1371/journal.pgen.1000008>
- Huang, W., and T. F. Mackay, 2016 The genetic architecture of quantitative traits cannot be inferred from variance component analysis. *PLoS Genet.* 12: e1006421. <https://doi.org/10.1371/journal.pgen.1006421>
- IWGSC, I. W. G. S. C., 2018 Shifting the limits in wheat research and breeding using a fully annotated reference genome by the international wheat genome sequencing consortium (iwgsc). *Science* 361: eaar7191.
- Jannink, J.-L., A. J. Lorenz, and H. Iwata, 2010 Genomic selection in plant breeding: from theory to practice. *Brief. Funct. Genomics* 9: 166–177. <https://doi.org/10.1093/bfpg/elq001>
- Jiang, Y., and J. C. Reif, 2015 Modeling epistasis in genomic selection. *Genetics* 201: 759–768. <https://doi.org/10.1534/genetics.115.177907>
- Jiang, Y., R. H. Schmidt, Y. Zhao, and J. C. Reif, 2017 A quantitative genetic framework highlights the role of epistatic effects for grain-yield heterosis in bread wheat. *Nat. Genet.* 49: 1741–1746. <https://doi.org/10.1038/ng.3974>

- Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics* 25: 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>
- Liaw, A., and M. Wiener, 2002 Classification and regression by random-forest. *R News* 2: 18–22.
- Lynch, M., and J. S. Conery, 2000 The evolutionary fate and consequences of duplicate genes. *Science* 290: 1151–1155. <https://doi.org/10.1126/science.290.5494.1151>
- Lynch, M. and J. S. Conery, 2003 The origins of genome complexity. *Science* 302: 1401–1404.
- Mac Key, J., 1970 Significance of mating systems for chromosomes and gametes in polyploids. *Hereditas* 66: 165–176.
- Marcussen, T., S. R. Sandve, L. Heier, M. Spannagl, M. Pfeifer *et al.*, 2014 Ancient hybridizations among the ancestral genomes of bread wheat. *Science* 345: 1250092. <https://doi.org/10.1126/science.1250092>
- Martini, J. W., V. Wimmer, M. Erbe, and H. Simianer, 2016 Epistasis and covariance: how gene interaction translates into genomic relationship. *Theor. Appl. Genet.* 129: 963–976. <https://doi.org/10.1007/s00122-016-2675-5>
- McClintock, B., 1984 The significance of responses of the genome to challenge. *Science* 266: 792–801. <https://doi.org/10.1126/science.15739260>
- Meuwissen, T., B. J. Hayes, and M. Goddard, 2001 Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157: 1819–1829.
- Microsoft and R Core Team, Microsoft R Open 2017 Microsoft Corporation. Redmond, Washington. <https://mran.microsoft.com>
- Nejati-Javaremi, A., C. Smith, and J. Gibson, 1997 Effect of total allelic relationship on accuracy of evaluation and response to selection. *J. Anim. Sci.* 75: 1738–1745. <https://doi.org/10.2527/1997.7571738x>
- Ohno, S., 1970 *Evolution by Gene Duplication*, Springer, New York. <https://doi.org/10.1007/978-3-642-86659-3>
- Ozkan, H., A. A. Levy, and M. Feldman, 2001 Allopolyploidy-induced rapid genome evolution in the wheat (aegilops-triticum) group. *Plant Cell* 13: 1735–1747. <https://doi.org/10.1105/tpc.13.8.1735>
- Patterson, N., A. L. Price, and D. Reich, 2006 Population structure and eigenanalysis. *PLoS Genet.* 2: e190. <https://doi.org/10.1371/journal.pgen.0020190>
- Poland, J. A., P. J. Brown, M. E. Sorrells, and J.-L. Jannink, 2012 Development of high-density genetic maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing approach. *PLoS One* 7: e32253. <https://doi.org/10.1371/journal.pone.0032253>
- Price, A. L., N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick *et al.*, 2006 Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* 38: 904–909. <https://doi.org/10.1038/ng1847>
- R Core Team, 2015 *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria.
- Rutkoski, J. E., J. Poland, J.-L. Jannink, and M. E. Sorrells, 2013 Imputation of unordered markers and the impact on genomic selection accuracy. *G3 (Bethesda)* 3: 427–439. <https://doi.org/10.1534/g3.112.005363>
- Salamini, F., H. Özkan, A. Brandolini, R. Schäfer-Pregl, and W. Martin, 2002 Genetics and geography of wild cereal domestication in the near east. *Nat. Rev. Genet.* 3: 429–441. <https://doi.org/10.1038/nrg817>
- Santantonio, N., J.-L. Jannink, and M. E. Sorrells, 2018a A low resolution epistasis mapping approach to identify chromosome arm interactions in allohexaploid wheat. *G3 (Bethesda)* 9: 675–684. <https://doi.org/10.1534/g3.118.200646>
- Santantonio, N., J.-L. Jannink, and M. E. Sorrells, 2018b A subfunctionalization epistasis model to evaluate homeologous gene interactions in allopolyploid wheat. *bioRxiv* 376731. <https://doi.org/10.1101/376731>
- Segovia-Lerma, A., L. Murray, M. Townsend, and I. Ray, 2004 Population-based diallel analyses among nine historically recognized alfalfa germ-plasms. *Theor. Appl. Genet.* 109: 1568–1575. <https://doi.org/10.1007/s00122-004-1784-8>
- Stekhoven, D. J., and P. Bühlmann, 2011 Missforest—non-parametric missing value imputation for mixed-type data. *Bioinformatics* 28: 112–118. <https://doi.org/10.1093/bioinformatics/btr597>
- Strandén, I., and D. Garrick, 2009 Derivation of equivalent computing algorithms for genomic predictions and reliabilities of animal merit. *J. Dairy Sci.* 92: 2971–2975. <https://doi.org/10.3168/jds.2008-1929>
- Stuber, C. W., S. E. Lincoln, D. Wolff, T. Helentjaris, and E. Lander, 1992 Identification of genetic factors contributing to heterosis in a hybrid from two elite maize inbred lines using molecular markers. *Genetics* 132: 823–839.
- Technow, F., C. Riedelsheimer, T. A. Schrag, and A. E. Melchinger, 2012 Genomic prediction of hybrid performance in maize with models incorporating dominance and population specific marker effects. *Theor. Appl. Genet.* 125: 1181–1194. <https://doi.org/10.1007/s00122-012-1905-8>
- Van Raden, P., 2008 Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91: 4414–4423. <https://doi.org/10.3168/jds.2007-0980>
- Visscher, P. M., S. Macgregor, B. Benyamin, G. Zhu, S. Gordon *et al.*, 2007 Genome partitioning of genetic variation for height from 11,214 sibling pairs. *Am. J. Hum. Genet.* 81: 1104–1110. <https://doi.org/10.1086/522934>
- Vitezica, Z. G., L. Varona, and A. Legarra, 2013 On the additive and dominant variance and covariance of individuals within the genomic selection scope. *Genetics* 195: 1223–1230. <https://doi.org/10.1534/genetics.113.155176>
- Wang, J., M.-C. Luo, Z. Chen, F. M. You, Y. Wei *et al.*, 2013 Aegilops tauschii single nucleotide polymorphisms shed light on the origins of wheat d-genome genetic diversity and pinpoint the geographic origin of hexaploid wheat. *New Phytol.* 198: 925–937. <https://doi.org/10.1111/nph.12164>
- Wolfe, M. D., P. Kulakow, I. Y. Rabbi, and J.-L. Jannink, 2016 Marker-based estimates reveal significant non-additive effects in clonally propagated cassava (*manihot esculenta*): implications for the prediction of total genetic value and the selection of varieties. *G3 (Bethesda)* 6: 3497–3506. <https://doi.org/10.1534/g3.116.033332>
- Yang, J., T. A. Manolio, L. R. Pasquale, E. Boerwinkle, N. Caporaso *et al.*, 2011 Genome partitioning of genetic variation for complex traits using common snps. *Nat. Genet.* 43: 519–525. <https://doi.org/10.1038/ng.823>

Communicating editor: F. van Eeuwijk

APPENDIX 1 CNLM DATASET

Field plots 1.5 m by 3 m in size were planted with 100 g of seed in September or October of the year prior to the harvest year. Data were recorded for four agronomic traits: grain yield (GY), plant height (PH), test weight (TW) and heading date (HD). Plots were harvested for grain with a plot combine after physiological maturity and oven dried to a grain moisture of approximately 12%. Dried grain was cleaned, weighed and measured for moisture content using a grain moisture analyzer (GAC 2100, Dickey-John). GY was standardized to a uniform grain moisture of 12%. PH was measured as the distance from the ground to the top of the grain head at full extension. TW is used as a measure of grain quality and was measured as the mass of a volume of grain (gL^{-1}) using the grain moisture analyzer (GAC 2100 Dickey-John) which corrects TW for moisture content. HD is a proxy for flowering time and was defined as the number of Julian days until 50% of the primary grain heads have extended out of the boot.

The data set initially consisted of 1,552 lines evaluated in 10,069 1.5 m by 3 m plots planted in September to October of the year prior to the harvest year. Thirty one lines from 2007 were not harvested for GY, nor were they genotyped, and were dropped from the data set. Because GY is of primary interest for breeding, plots that were not harvested or had missing values for GY were dropped, resulting in 9,090 plots with GY measurements. This also caused two additional lines with missing GY measurements to be dropped from the dataset.

Due to the reasonable size of the dataset, small physical area of most trials, lack of replication within environment for most lines and the availability of genetic markers, raw plot observations were used and no attempt was made to correct plot level data for various spatial effects or otherwise. Preliminary results also indicated relatively high genomic prediction accuracy, suggesting that spatial correction, such as an $AR1 \times AR1$ row column autocorrelation structure, would be unlikely to reduce error variance drastically. Instead, 59 plots that included breeder comments about bad seed or significant damage to the plot, via animal or otherwise, were removed from the dataset. Observations outside of a four standard deviation interval from the grand mean of uncorrected GY phenotypic observations were also removed to account for any significant undocumented damage, grain spillage or other undocumented mistakes. This included two observations that were deemed too high, and 20 observations that were deemed too low.

Observations of 11 lines lacking at least one phenotypic record in at least two separate trials, 60 lines that were missing at least 50% of their genotypic calls and one line that had greater than 20% heterozygous genotype calls were also removed from the dataset. This resulted in 8,692 phenotypic observations of 1,447 lines across 26 environments, representing 96.6% of the plots with grain yield measurements. HD was not recorded for the 246 observations from 2007, and PH was not recorded for the 840 observations from 2009. Two additional plots were missing PH measurements from the Ketola location, one that was recorded at 2 meters in 2008 which was likely a recording mistake, and one in 2010 which was simply missing a record.

While most of the genotypes were directly from the Cornell small grains breeding program, a few varieties and breeding lines from other breeding programs that had been genotyped and evaluated were not excluded from the dataset as long as they met the previous criteria. This included 75 lines from The Ohio State University wheat breeding program and 93 lines from the Michigan State University wheat breeding programs that were part of the Allele Based Breeding initiative, among other lines from various breeding programs.

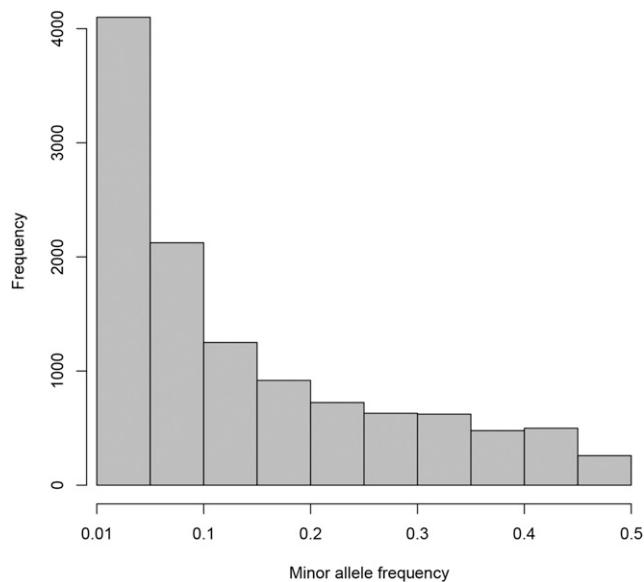


Figure A1 Distribution of minor allele frequencies for 11,604 GBS markers in the CNLM population.

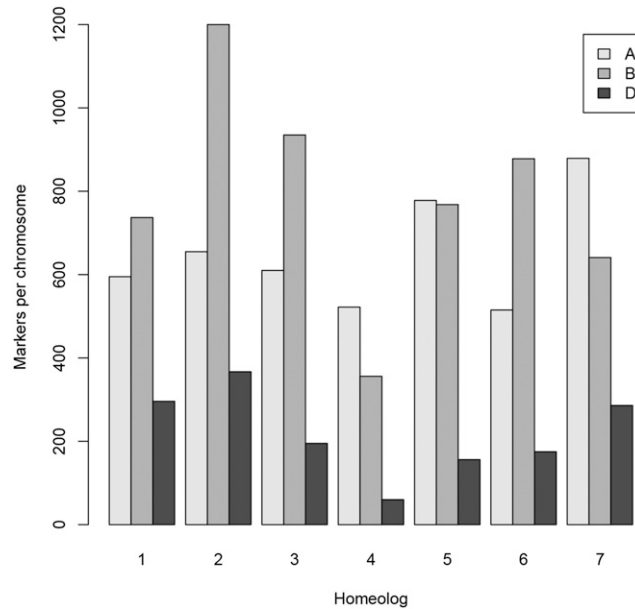


Figure A2 Distribution of 11,604 GBS markers on the 21 wheat chromosomes comprised of 7 homeologs of three subgenomes, A, B and D, for the CNLM population.

Genotyping by sequencing (GBS) libraries (Elshire *et al.* 2011) of 1,521 CNLM were developed using the protocol described by Poland *et al.* (2012) at Kansas State University, and subsequently sequenced at the Genomic Diversity Facility at Cornell University. Genotyping calls were accomplished using standard parameters of the Tassel 5.0 GBS v2 Pipeline (Glaubitz *et al.* 2014) and were aligned to the International Wheat Genome Sequencing Consortium (IWGSC) RefSeq v1.0 wheat genome sequence of ‘Chinese Spring’ (IWGSC 2018). Following Poland *et al.* (Poland *et al.* 2012), 64 bp sequence tags containing no more than three Single Nucleotide Polymorphisms (SNPs) per tag were included to increase the likelihood of obtaining subgenome specific markers. Only markers with a minor allele frequency of at least 0.01 (Figure A1), no more than 30% missing scores, and no more than 10% heterozygous calls were kept for the following analyses. Then individuals with greater than 20% heterozygous calls and individuals with more than 50% missing genotype calls were excluded from the dataset. The process was repeated iteratively, starting by filtering on markers until the number of markers and genotypes converged. This resulted in 11,604 available GBS markers distributed throughout the three subgenomes (Figure A2). Of the 61 lines removed, 60 were removed due to missing marker information and 1 was removed due to high heterozygosity in the final iteration. The 11 lines with a single observation and the two lines without grain yield observations were subsequently removed to produce genotypic information for the 1,447 lines.

Marker scores were coded using $\{-1, 0, 1\}$ for homozygous major allele, heterozygous and homozygous minor allele, respectively. Categorical marker imputation was done independently for each chromosome using random forest imputation via the R package ‘missForest’ (Stekhoven and Bühlmann 2011) which relies on the R package ‘randomForest’ (Liaw and Wiener 2002). Random forest has been shown to be effective for genotype imputation in wheat (Rutkoski *et al.* 2013). To allow all individuals to be considered completely inbred, the remaining heterozygous calls ($< 2\%$ of all marker scores) were conservatively replaced with the population mode for that marker (*i.e.* the homozygous major allele, -1). Marker scores were then converted to $\{0, 1\}$ coding for presence of the minor allele.

Genetic correlations of traits were estimated in a multivariate model fit (Table A2). This was accomplished by treating genotypes as independent, or having a realized additive covariance structure calculated from genetic markers.

Table A1 Means (μ) and standard deviations (σ) of four traits in the CNLM population

	units	μ	σ
GY	kg ha	5315.20	1015.76
PH	cm	90.84	11.99
HD	Julian days	151.64	3.87
TW	g L	74.95	3.09

■ **Table A2** Estimated genetic correlation of traits with additive (below diagonal) and independent genetic relationships (above diagonal). Genetic standard deviations of scaled traits estimated with a realized additive covariance between individuals and assuming independence are shown in parentheses on the diagonal, respectively

	GY	PH	TW	HD
GY	(0.29, 0.36)	-0.39	-0.24	0.16
PH	-0.44	(0.72, 0.65)	0.31	0.05
HD	-0.05	0.11	(0.44, 0.44)	-0.28
TW	-0.04	0.3	-0.22	(0.5, 0.49)

APPENDIX 2 REMOVAL OF POPULATION STRUCTURE FROM \mathbf{M}

Let \mathbf{M} be the $n \times m$ matrix of m marker scores for n genotypes. Markers can be sorted into their respective genome, such that

$$\mathbf{M} = [\mathbf{M}_A \mathbf{M}_B \mathbf{M}_D] \quad (12)$$

\mathbf{M} can be factored using singular value decomposition as follows:

$$\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^T \quad (13)$$

where \mathbf{U} , and \mathbf{V} are unitary matrices of left and right singular vectors, and \mathbf{D} is a diagonal matrix of singular values.

The first k principal components of the marker matrix can be extracted by selecting the first k columns of \mathbf{U} and the first k rows and columns of \mathbf{D} and multiplying.

$$\text{Let } \mathbf{Q} = \mathbf{U}_{n \times k} \mathbf{D}_{k \times k} \quad (14)$$

In a manner similar to Eckart and Young (Eckart and Young 1936), an approximation, $\tilde{\mathbf{M}}$, of the marker matrix, \mathbf{M} with the first q principal components removed can be reconstructed by setting the first k singular values in \mathbf{D} to zero (denoted $\tilde{\mathbf{D}}$).

$$\tilde{\mathbf{M}} = \mathbf{U}\tilde{\mathbf{D}}\mathbf{V}^T = [\tilde{\mathbf{M}}_A \tilde{\mathbf{M}}_B \tilde{\mathbf{M}}_D] \quad (15)$$