# Meta-analysis identifies pleiotropic loci controlling phenotypic trade-offs in sorghum

Ravi V. Mural (ID) ,[1] Marcin Grzybowski (ID) ,[1] Chenyong Miao (ID) ,[1] Alyssa Damke,[2] Sirjan Sapkota,[3,4] Richard E. Boyles (ID) ,[4,5]
Maria G. Salas Fernandez,[6] Patrick S. Schnable (ID) ,[6] Brandi Sigmon,[2] Stephen Kresovich,[4,7] and James C. Schnable (ID) [1,*]

[1]Center for Plant Science Innovation and Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Lincoln, NE 68588, USA,
[2]Department of Plant Pathology, University of Nebraska-Lincoln, Lincoln, NE 68588, USA,
[3]Advanced Plant Technology Program, Clemson University, Clemson, SC 29634, USA,
[4]Department of Plant and Environment Sciences, Clemson University, Clemson, SC 29634, USA,
[5]Pee Dee Research and Education Center, Clemson University, Florence, SC 29532, USA,
[6]Department of Agronomy, Iowa State University, Ames, IA 50010, USA and
[7]Feed the Future Innovation Lab for Crop Improvement, Cornell University, Ithaca, NY 14850, USA

*Corresponding author: Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Beadle Center E207, Lincoln, NE 68583-0660, USA. Email: schnable@unl.edu

## Abstract

Community association populations are composed of phenotypically and genetically diverse accessions. Once these populations are genotyped, the resulting marker data can be reused by different groups investigating the genetic basis of different traits. Because the same genotypes are observed and scored for a wide range of traits in different environments, these populations represent a unique resource to investigate pleiotropy. Here, we assembled a set of 234 separate trait datasets for the Sorghum Association Panel, a group of 406 sorghum genotypes widely employed by the sorghum genetics community. Comparison of genome-wide association studies (GWAS) conducted with two independently generated marker sets for this population demonstrate that existing genetic marker sets do not saturate the genome and likely capture only 35–43% of potentially detectable loci controlling variation for traits scored in this population. While limited evidence for pleiotropy was apparent in cross-GWAS comparisons, a multivariate adaptive shrinkage approach recovered both known pleiotropic effects of existing loci and new pleiotropic effects, particularly significant impacts of known dwarfing genes on root architecture. In addition, we identified new loci with pleiotropic effects consistent with known trade-offs in sorghum development. These results demonstrate the potential for mining existing trait datasets from widely used community association populations to enable new discoveries from existing trait datasets as new, denser genetic marker datasets are generated for existing community association populations.

Keywords: pleiotropy; GWAS; quantitative genetics; sorghum; community genetic resources

## Introduction

The value of common mapping populations in diverse species has been recognized by quantitative genetics for decades. These common populations can be genotyped by a single research group, and genetically identical individuals distributed to the larger research community which could score traits of interest and test for associations between published genetic markers and trait values across the population. Early populations included sets of recombinant inbred lines (RILs) (Bailey 2004). RIL populations were developed, genotyped, and released to the maize (Burr et al. 1988; Lee et al. 2002) and arabidopsis (Lister and Dean 1993) genetics communities. Mapping quantitative trait loci (QTL) in RIL populations provided relatively high power to detect variants with even small numbers of markers as large scale linkage disequilibrium (LD) permits the identification of associations between genetic markers and trait values for markers at quite some distance from the causal variant. Yet the high LD of RIL populations also meant that researchers were unlikely to identify

associations with the causal gene or variant without generating new follow-up populations for fine mapping. Improvements to genotyping technologies ultimately permitted the use of natural populations with more rapid decay of LD (Flint-Garcia et al. 2005; Nordborg et al. 2005; Atwell et al. 2010). Like earlier RIL populations, association populations can be distributed among research groups, permitting association mapping and later genome-wide association studies (GWAS) to be conducted for multiple traits without the need to generate new marker data.

The number of markers needed to saturate the genome of a target population for GWAS is determined by the size of the species' genome, the speed at which LD decays within the target population, and the minimum level of LD between a causal variant and a genotyped marker where statistically significant associations will still be detected, which in turn depends on the proportion of total population variance explained by the causal variant (Bouchet et al. 2012). An early estimate based on 1122 genetic markers suggests that >100,000 genetic markers would be necessary to conduct GWAS in populations spanning global

sorghum diversity when target causal variants explained >10% of total trait variance and >350,000 genetic markers when target causal variants explained 5–10% of total trait variance (Bouchet *et al.* 2012). Estimates of this type are quite sensitive to the rate of LD decay. Estimates of the average distance at which LD decays below an $r^2$ value of 0.1 in sorghum range from 10 kilobases to 350 kilobases depending on the population and genetic marker set employed (Hamblin *et al.* 2005; Bouchet *et al.* 2012; Mace *et al.* 2013; Wang *et al.* 2013; Morris *et al.* 2013b; Wang *et al.* 2020). LD also varies among different portions of the genome, creating an addition challenge to accurate simulation.

Higher marker densities increase the odds of identifying causal variants that explain more modest proportions of the total variance for a target trait. When new sets of genetic marker become available for existing association populations it is possible to reanalyze previously collected trait datasets. In addition, studies based on simulated data suggest multivariate analyses may increase true positive rates relative to trait-by-trait univariate GWAS (Rice *et al.* 2020). Multivariate analysis may provide value in case where individual causal loci have pleiotropic effects on multiple traits which are measured separately. The degree of pleiotropy for loci influencing variation in quantitative traits in plants remains uncertain. A study of maize leaf traits found little incidence of pleiotropy (Tian *et al.* 2011) while detecting modest evidence of pleiotropy between the elongation of leaves tassel and ears (Brown *et al.* 2011). The sorghum QTL Atlas, a meta-analysis of reported QTL map locations from 146 publications of diverse RIL populations identified QTL hotspots on chromosome 2, corresponding to the *brown nucellar layer2* (B2) gene in sorghum, and chromosome 7, corresponding to *dwarf3* (dw3) (Mace *et al.* 2019). However, the relatively large confidence interval of QTL peaks identified via mapping in RIL populations can make it difficult to determine whether QTL hotspots represent a single highly pleiotropic gene or multiple linked genes (Wallace *et al.* 2014). In maize, the pleiotropic effects of a large effect QTL for plant architecture on ear traits was long thought to be explained by a single gene, *teosinte branched1* (tb1) but was later fractionated into multiple partially linked loci for ear-related traits (Studer and Doebley 2011).

The Sorghum Association Panel (SAP) was first assembled in 2008 (Casa *et al.* 2008). After some additions, the population ultimately consisted of 406 (USDA 2010) lines selected to represent the global genetic diversity of sorghum. Because it was intended that this population be grown and phenotyped in the temperate United States, the majority of the lines included in the panel were generated by the Sorghum Conversion program (Casa *et al.* 2008; Boyles *et al.* 2019). Sorghum from many parts of the world fail to flower during the summer growing season in the temperate United States. Sorghum Conversion lines are the result of crossing diverse sorghum germplasm from around the world to a temperate adapted donor parent (BTx406) and then recurrently backcrossing the progeny to the exotic-tropical parent for four generations while selecting for temperate adaptation, including flowering during the summer in temperate latitudes and short stature (Stephens *et al.* 1967). Retrospective genomic analysis of many Sorghum Conversion lines identified three genomic intervals where the haplotype of the donor parent is over represented in the population. These three intervals corresponded to the locations of dwarfing genes *dw1–dw3* (Thurber *et al.* 2013). No comparable independent intervals were identified for loci conferring photoperiod insensitivity; however, *maturity1* (ma1) has a large selection sweep on Chr6 that is linked to *dwarf2* (dw2) (Thurber *et al.* 2013).

Initially, the SAP was genotyped for only a set of 49 simple-sequence repeat markers (Casa *et al.* 2008). Subsequently, the SAP

was employed for a number of genetic association tests using increasing numbers of markers (Brown *et al.* 2008; Casa *et al.* 2008; Sukumaran *et al.* 2012; Wu *et al.* 2012; Hufnagel *et al.* 2014). In 2013, a set of several hundred thousand genetic markers was generated for the population using conventional genotyping by sequencing (Elshire *et al.* 2011; Morris *et al.* 2013a). From 2013 onward, the SAP was widely employed for GWAS by a range of research groups targeting different traits (Table 1). Here, we employ a set of both published and previously unpublished trait datasets from the SAP and multiple genetic marker datasets (Morris *et al.* 2013a; Miao *et al.* 2020a) to empirically evaluate both the degree of saturation achieved by current genetic marker sets and the degree to which detectable loci controlling phenotypic variation in the SAP tend to be pleiotropic or non-pleiotropic using a multi-trait approach based on meta-analysis and adaptive shrinkage (Urbut *et al.* 2019).

## Materials and methods
### Genetic marker datasets

A set of 265,487 single nucleotide polymorphisms (SNPs) generated using conventional genotyping by sequencing and aligned to version 1 of the BTx623 sorghum reference genome were downloaded from http://people.beocat.ksu.edu/gpmorris/sorghum_GBS_data/readme.txt/ (Morris *et al.* 2013a) (the "2013 dataset"). A set of 569,305 SNPs generated using a modified genotyping by sequencing approach and aligned to version 3 of the sorghum genome was obtained from FigShare (https://doi.org/10.6084/m9.figshare.11462469.v5) (Miao *et al.* 2020a) (the "2020 dataset"). In each dataset missing data points were imputed using Beagle (v4.1) with the sliding windows set individually for each chromosome to capture 10% of all markers on that chromosome and overlap windows set to capture 2% of call markers on that chromosome (Browning and Browning 2016). After imputation, marker sets were filtered by removing markers with minor allele frequencies of less than 5% among the set of genotypes employed for a given analysis. This filtering criteria resulted in a set of 107,751 markers scored across 304 lines for the 2013 dataset and a set of 257,882 markers scored across the same 304 lines for the 2020 dataset. Filtering using the same parameters across all 343 lines included in the 2020 dataset produced a set of 256,695 markers.

### Trait datasets

A total of 234 trait datasets scored across all or subsets of the SAP were employed in this study. One hundred ninety of these trait datasets were drawn from published sources as described in Supplementary Table S1. An additional 12 phenotypes were generated using sums or ratios of published trait datasets. The traits and formulas used to generate these 12 phenotypes are provided in Supplementary Table S2. The remaining 32 trait datasets employed in this study were previously unpublished datasets collected at either the University of Nebraska-Lincoln in Nebraska (19 datasets) or Clemson University in South Carolina (13 datasets). All these phenotype data are provided on FigShare (https://doi.org/10.6084/m9.figshare.13143389).

#### *Nebraska trait collection:*

A single replicate of the SAP was grown near Mead, NE in 2016 and 2017. Plant height to inflorescence, plant height to flag leaf, leaf angle (3rd leaf), stem diameter (between the 3rd and 4th leaf), and node number were measured from a representative plant at reproductive maturity in 2016. A ratio of plant height to inflorescence/plant height to flag leaf was also calculated

**Table 1** Papers scoring traits in the sorghum association population

| Reference | Study type | Phenotypes scored | # of SAP accessions evaluated | Genetic associations? | Trait data online? |
|---|---|---|---|---|---|
| Casa *et al.* (2008) | Vegetative | 8 | 377 | AS[a] | No |
| Brown *et al.* (2008) | Height & Inflorescence | 6 | 378 | AS | Yes |
| Vandenbrink *et al.* (2010) | Biomass Composition | 2 | 377 | No | Yes |
| Mutava *et al.* (2011) | Drought Stress | 10 | 300 | No | No |
| Sukumaran *et al.* (2012) | Grain Quality | 10 | 300 | AS | No |
| Wu *et al.* (2012) | Tannins | 1 | 161 | AS | Yes |
| Morris *et al.* (2013a) | Various | 6 | 355 | GWAS | Used public data |
| Morris *et al.* (2013b) | Flavonoids | 2 | 259–387 | GWAS | Yes |
| Hufnagel *et al.* (2014) | Phosphorous Deficiency | 10 | 287 | AS | No |
| Kong *et al.* (2014) | Tillering and inflorescence | 9 | 377 | GWAS | No |
| Perez *et al.* (2014) | Plant Architecture | 6 | 315 | GWAS | Released here |
| Rhodes *et al.* (2014) | Polyphenols | 3 | 308 | GWAS | No |
| Adeyanju *et al.* (2015) | Disease Resistance | 6 | 300 | GWAS | Yes, but IDs ambiguous |
| Lasky *et al.* (2015) | Drought Stress | 28 | 267 | GWAS | Yes |
| Prom *et al.* (2015) | Disease Resistance | 3 | 177 | No | Yes |
| Queiroz *et al.* (2015) | Grain Quality | 12 | 100 | No | Yes |
| Li *et al.* (2015) | Height | 3 | 307 | GWAS | Used public data |
| Zhang *et al.* (2015a) | Height & Inflorescence | 12 | 354 | GWAS | No |
| Zhang *et al.* (2015b) | Seed Size | 6 | 354 | GWAS | Yes |
| Boyles *et al.* (2016) | Various | 13 | 378 | GWAS | Released here |
| Shakoor *et al.* (2016) | Elemental Abundance | 22 | 407 | GWAS | Yes |
| Zhao *et al.* (2016) | Plant Architecture | 9 | 315 | GWAS | Released here |
| Boyles *et al.* (2017) | Grain Quality | 10 | 378 | GWAS | Released here |
| Chen *et al.* (2017) | Heat Stress | 2 | 374 | GWAS | No |
| Chopra *et al.* (2017) | Heat and Cold Stress | 12 | 300 | GWAS | Yes |
| Fernandez *et al.* (2017) | Vegetative (HTP) | 4 | 307 | GWAS | No |
| Ortiz *et al.* (2017) | Photosynthesis/Cold Stress | 24 | 304 | GWAS | No |
| Paiva *et al.* (2017) | Elemental Abundance | 18 | 100 | No | Yes |
| Rhodes *et al.* (2017a) | Polyphenols | 3 | 266 | GWAS | Yes |
| Rhodes *et al.* (2017b) | Grain Quality | 4 | 265 | GWAS | Yes |
| Cuevas *et al.* (2018) | Disease Resistance | 2 | 335 | GWAS | Yes |
| Breitzman *et al.* (2019) | Vegetative (HTP) | 6 | 325 | GWAS | Released here |
| Cuevas *et al.* (2019) | Disease Resistance | 2 | 331 | GWAS | Yes |
| McMaster *et al.* (2019) | Mycotoxin | 2 | 98 | No | Yes, but IDs ambiguous |
| Moghimi *et al.* (2019) | Cold Stress | 13 | 351 | GWAS | Yes |
| Olatoye *et al.* (2019) | Various | 4 | 334 | GWAS | Used public data |
| Prom *et al.* (2019) | Disease Resistance | 1 | 359 | GWAS | Yes |
| Zheng *et al.* (2020) | Root Architecture (HTP) | 12 | 294 | GWAS | Yes |
| Zhou *et al.* (2019) | Inflorescence (HTP) | 8 | 302 | GWAS | Yes |
| Miao *et al.* (2020b) | Height | 1 | 357 | GWAS | Yes |

[a] "AS" means association studies which were not conducted using genome wide sets of markers, and "GWAS" means association studies which did utilize genome wide sets of markers.

(Supplementary Table S2). Inflorescence architecture traits were measured using two representative plants at maturity in 2016, 2017 and included inflorescence length, rachis length, rachis diameter, number of primary branches, length of primary branches at the bottom third of the inflorescence, length of primary branches at the top third of the inflorescence, number of secondary branches on a primary branch, number of third-order branches on a secondary branch, first internode length on a primary branch, prominent awns (binary trait), and prominent glumes (binary trait). During 2017 one additional trait, infertility (scored on a scale from 1 to 4), was also collected (Supplementary Figure S1). Two ratios were also calculated from each year in this dataset: inflorescence length/rachis length, rachis length/rachis diameter (Supplementary Table S2). Best linear unbiased predictors (BLUPs) for each phenotype were calculated by fitting a linear mixed model using R package lme4 (Bates *et al.* 2014) with genotype, year and genotype by year variables fit as random for traits with multiple years of data and only genotypes fit as random variable for the traits with data from only 1 year.

### South Carolina trait collection:

The SAP was grown near Florence, South Carolina in 2013, 2014, and 2017. In each year two replicates per line were grown in a 2× replicated completely randomized design utilizing two row

yield plots. Trait datasets collected at Clemson University included two flowering time-related traits, measured in all 3 years: days to anthesis and grain fill duration (days to maturity—days to anthesis). Two plant height traits were measured: plant height from ground/plant base to panicle apex and flag leaf height. Six reproductive traits were measured: number of grains per primary panicle and grain yield per primary panicle, measured in all 3 years, glume tenacity (0–5 visual rating), primary panicle branch length, panicle length, and exsertion in 2017. Five biochemical traits: magnesium (% dry basis), manganese (ppm), nitrogen (mg), phosphorus (% dry basis), and zinc (ppm) measured from ground grain samples in 2013 and 2014 using near-infrared spectroscopy (Boyles *et al.* 2017). Thirteen seed traits were determined from these experiments: percent moisture, 1000-grain weight, percent of dry mass which was acid detergent fiber and percent of dry mass which was neutral detergent fiber, percent of starch which is amylose, percent of dry grain weight which was oil, protein, or starch, *in vitro* starch digestibility, gross energy per gram (calories/gram), iron (ppm), prolamin as a percentage of dry weight, and seed density (grams per milliliter). All except 1000-grain weight and seed density were measured using near infrared and were evaluated in all 3 years. All phenotypes were measured

in all 3 years except seed density which was measured only in 2017. For each trait, a linear mixed model was fit using the R package lme4 (Bates *et al.* 2014) with genotype, year, genotype by year, and replication nested in a year were fit as random for traits from 2013, 2014, and 2017 combined and the genotype and replication fit as random for the traits collected only during 2017, and the resulting phenotypic BLUPs were employed for further analysis. Of the 28 traits included in this study, 15 were published in part or in totality (Boyles *et al.* 2016, 2017; Sapkota *et al.* 2020) and 13 are previously unpublished data.

*Trait data normalization and heritability calculation:*
With the exception of six binary traits (AnthraClassification_P, PericarpColor_D, Tannins_D, Tannins_F, AwnProminence_U, and GlumProminence_U) each trait dataset was normalized before analysis using the R package, bestNormalize version 1.4.3 (Peterson 2017). The function bestNormalize in the bestNormalize package performs various/suite of normalization transformations, such as Lambert W x F, BoxCox, YeoJohnson, Ordered Quantile, etc. and then select one optimal transformation for each dataset based on minimizing the Pearson P statistic, a test for normality. GWAS analyses for individual traits frequently incorporate a manual examination of trait value distributions to remove extreme outlier values. However, here a total of 468 distinct GWAS analyses (234 traits * 2 sets of individuals) were conducted, making it difficult or impossible to guarantee consistent criteria would be applied to manual removal of outliers across 468 distributions. Hence, a rules-based automated outlier removal strategy was adopted. Subsequent to normalization, values which were more than 1.5 times the interquartile range below the 25th percentile of normalized trait values or more than 1.5 times the interquartile range above the 75th percentile of normalized trait values were converted to missing data. In order to determine the proportion of genetic variation explainable by genetic factors, marker-based estimate of narrow-sense heritability was generated using the R-package sommer (v4.1.1), the reported values for each line for each trait, and the 2020 genetic marker dataset (Covarrubias-Pazaran 2016).

## Genome-wide association analyses

GWAS was performed independently on each of 234 trait datasets. The number and identity of lines evaluated for various traits varied both across and within studies (Supplementary Table S1). For each trait dataset, analyzed with each of the two genetic marker datasets, genetic markers were separately filtered to remove those markers with a minor allele frequency of <5% among individuals with recorded values for the target trait.

The resulting marker sets were employed for GWAS analysis using two single-locus models; generalized linear model (PC), mixed linear model (PC+K) (Price *et al.* 2006), and one multi-locus model; FarmCPU (Liu *et al.* 2016). For all three models, the implementations used were those included in the R package rMVP (v1.0.1) (Yin *et al.* 2020). For the GLM model, the first three principal components (PCs) were fit as covariates in order to control for the population structure. In case of the mixed linear model (MLM), in addition to the first three PCs, a kinship matrix was also integrated in the model for association analysis. The kinship matrix representing the relationship among individuals used in the MLM model was calculated using the first method described by VanRaden (2008) as implemented within the rMVP package, which should be equivalent to the method of Endelman and Jannink (2012) for high-density marker data. The multi-locus mixed linear model; FarmCPU, was run with the first three PCs as

covariates and the kinship matrix calculated internally by the FarmCPU algorithm fitted as random effects. FarmCPU was run using maxLoop = 10, the method for selecting most appropriate bins was run with the MVP option, method.bin = "FaST-LMM" (Lippert *et al.* 2011). The method of variance components analysis, vc.methods was set to "GEMMA" in the association analysis (Zhou 2017).

Bonferroni corrections were applied based on the effective number of independent markers in each genetic marker dataset. Effective SNP numbers were calculated for each dataset using the genetic type I error calculator (GEC (v0.2)) software package with default parameter settings (Li *et al.* 2012). GEC implements an eigenvalue-based method which employs the matrix of correlations in p-values in association testing between SNP markers to estimate the effective number of independent SNPs (Me) which will be lower than or equal to the actual number of SNPs genotyped. Statistical significance thresholds were set to $10^{-5.958}$ (0.05/48,488) and $10^{-6.154}$ (0.05/825,223) for 2013 and 2020 genetic marker datasets, respectively when SNPs were first filtered based on their minor allele frequencies in the set of 304 sorghum varieties shared between the two datasets. The statistical significance threshold for the 2020 genetic marker dataset when markers were filtered based on their frequency in the complete set of 343 lines genotyped in that dataset was $10^{-6.1379}$ (0.05/82,143). Trait values were available for different subsets of lines for different trait datasets. While markers were further filtered to remove low frequency markers in subset populations as described above, the same three P-value cutoffs for statistical significance were employed for all GWAS results to provide consistency across analyses. Manhattan plots were created using the CMplot R package (v3.6.2) (Yin 2020). When multiple genetic markers showed statistically significant associations after correction for multiple testing, markers on the same chromosome and separated by no more than 1 MB were merged into a single peak. A custom python script (CallPeaksBatch.py) was employed to merge nearby SNPs into peaks and identify "summit" SNPs for each peak (Miao 2020).

## Multivariate gene-trait analyses

For multivariate analysis, 176 traits were chosen by excluding 30 traits collected from <100 lines, six binary traits and 22 ionomics trait data (Shakoor *et al.* 2016). Estimated effect size and standard error for each of the markers in each of these 176 traits was extracted from the initial output of rMVP using the results of analysis conducted using the MLM model. First, a subset of 671 strong signals with lfsr <0.1 were chosen by running a condition-by-condition analysis using ash with package ashr/v2.2-47 using the most recent code revision available on github as of 9/8/2020 (Stephens *et al.* 2020). A second control set of estimated effect size and standard error values for a set of 90,000 markers were randomly extracted to aid the mash model in learning the patterns of covariance between SNPs and each phenotype in order to produce improved effect estimates of the SNPs chosen from condition-by-condition analysis. Thus, the control set chosen is an unbiased representation of all the tests considered, including null and non-null tests. The strong signals made up approximately 0.8% of the control set. Furthermore, to overcome the confounding effects caused by the correlated variation among various traits/phenotypes, we estimated the simple correlation matrix V in the 90 K random control set using the MashR function "estimate_nul_correlation_simple" and included the resulting correlation matrix V (V = Vhat) into our analysis. We used simple null correlation to find arbitrary patterns of correlation among

various conditions. Instead of calculating correlation among all the null tests, we estimated simple null correlation among the random subset as it gives a quick approximation of the null correlation matrix. These datasets were analyzed using mashr/v0.2.40 using the most recent code revision available on github as of August 9, 2020 (Urbut *et al.* 2019). Following the recommendations of the MashR documentation, canonical and data-driven covariance matrices were computed. The canonical covariance matrix was calculated using the MashR function "cov_canonical". The data-driven covariance matrix was calculated by using the function "cov_pca". A mash model was fit using both the covariance matrices. The posterior summaries were computed for each SNP in the strong sub set, choosen from condition-by-condition analysis for each phenotype. Furthermore, the Bayes factor extracted from mash output with CDBNgenomics R package (MacQueen *et al.* 2020) was used to determine if a given SNP has significant phenotypic effect. Threshold of local false sign rate (*lfsr*) <0.001 was used to determine the number of traits associated to a given SNP (Stephens 2017). Codes to replicate mashR analysis using the above-mentioned datasets can be found at https://github.com/ravimural/sapmashr. Peaks were called by merging individual significant SNPs which were separated by less than 500 kilobases into peaks and selecting the single SNP within each peak with the largest Bayes factor as representative of that peak. The code used to perform this merging has been deposited on github (Miao 2020). Visualizations of mashr results were generated using the R package CMplot (v3.6.2) (Yin 2020) (Panel A).

### Data availability statement

All genotype data used in this study was drawn from published sources. Combined phenotypic data on both previously published, unpublished phenotypes used in this study, and the GWAS results are provided on FigShare https://doi.org/10.6084/m9.figshare.13143389. Supplemental Material available at figshare: https://doi.org/10.25386/genetics.14721033.

## Results
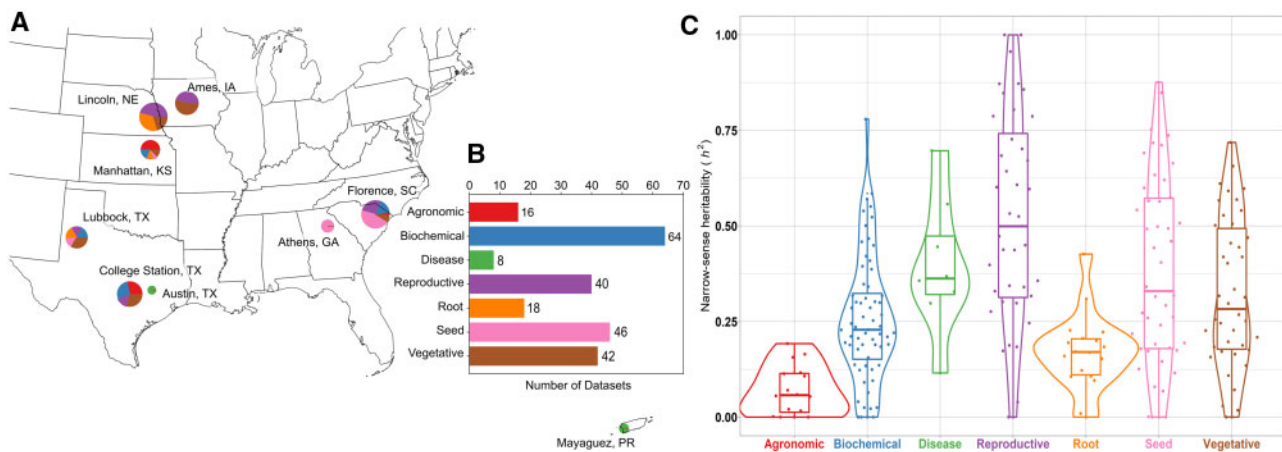
### Properties of SAP trait datasets

A literature review identified 40 papers in which phenotypic data were collected and published from the SAP—or subsets of this population. Of these 40 papers, it was possible to obtain trait values for individual lines in 25 cases. These included 20 papers for which data were provided as Supplementary information and 5 papers for which the data were obtained directly from the authors (Table 1). These 25 papers reported data for a total of 190 traits, although it should be noted that some of these traits are similar or identical measurements conducted in different environments or years. Twelve additional phenotypes were derived based on sums or ratios of trait means (Supplementary Table S1). Previously unpublished data for 19 traits collected in Nebraska and 13 previously unpublished trait datasets collected in South Carolina were also added to the dataset (Supplementary Table S1). Thus, the final dataset consisting of 234 distinct sets of trait data scored for all or subsets of a common sorghum population. Values of all 234 sets of trait data, including those previously unpublished are provided as supplementary data with this paper. A number of studies reported data from plants grown in growth chambers or other controlled environment conditions. However, the majority of published trait datasets came from field trials in six states—Georgia, Iowa, Kansas, Nebraska, South Carolina, and Texas—within the United States with additional trait data collected in Brazil (Figure 1A).

Traits could be broadly classified into seven categories including agronomic phenotypes, biochemical phenotypes, disease-related phenotypes, root phenotypes, above ground vegetative phenotypes (of which 15 plant height or plant height proxies), reproductive phenotypes, and seed phenotypes (Figure 1B; Supplementary Table S1). Most studies provided only trait means or best linear unbiased predictor (BLUP) (Robinson 1991) values (Supplementary Table S1). As a result, it was not possible to estimate broad sense heritabilities for most traits. However, it was possible to estimate narrow sense heritabilities. Estimates of narrow sense heritability had a median of 0.265 (Figure 1C). Traits with high narrow sense heritability tended to be those related to panicle morphology, grain composition and disease, while traits with estimates of narrow sense heritability close to zero tended to be collected from seedlings, a subset of biochemical traits, and measures of some plasticity across environments (Supplementary Table S1). Trait datasets belonging to the same categories, as well as trait datasets collected as part of the same studies tended to be correlated more with each other than with other pairs of traits (Supplementary Figure S2).

### Current sorghum marker sets do not achieve saturation for linkage to causal loci

Previous estimates of the number of markers required to achieve saturation of a sorghum GWAS population have been largely based on simulation studies. The existence of two distinct genetic marker datasets provides an opportunity to empirically estimate the number of markers required to saturate the sorghum genome. In 2013, a set of 265,487 markers identified relative to version 1 of the sorghum genome were generated using conventional genotyping by sequencing for 971 lines including 355 members of the SAP (Elshire *et al.* 2011; Morris *et al.* 2013a) (referred to as the "2013 marker set" below). In 2020, a set of 569,305 markers identified relative to version 3 of the sorghum genome were generated using a modified genotyping by sequencing strategy for 343 members of the SAP (Ott *et al.* 2017; Miao *et al.* 2020a) (referred to as the "2020 marker set" below). A total of 304 members of the SAP were included in both the 2013 and the 2020 genetic marker datasets. Filtering for only the subset of SNPs where the minor allele was present in at least 5% of the 304 shared SAP lines resulted in 107,751 and 257,882 markers in 2013 and 2020 marker sets respectively. The number of markers employed in GWAS were higher than the previous estimates of the number of markers (100,000) required to achieve saturation of a sorghum GWAS population based on a minimum $r^2$ of 0.1 between genotyped markers and causal variants (Bouchet *et al.* 2012).

Different protocols were employed to generate these two datasets and these two protocols sequence different subsets of the sorghum genome. As a result, the sets of specific genetic markers genotyped in each dataset should be largely non-overlapping. Indeed, as the marker sets were generated relative to different versions of the reference genome (v1 and v3), while it is possible to align results at the gene or region level it is not possible to confidently quantify the precise number of markers shared between the two datasets. If current marker datasets are sufficient to saturate the sorghum genome, using different marker datasets would be expected to identify signals from the same regions of the genome. However, if current marker sets are insufficient to saturate the sorghum genome, using different marker datasets to analyze the same trait datasets would be expected to identify only partially overlapping sets of genomic intervals for the same trait datasets.

**Figure 1.** Characteristics of Sorghum Association Panel trait datasets. (A) Geographic distribution of trials where trait datasets were collected. Size of circles indicates number of traits collected at a specific geographic location. Colors of circles indicate types of trait datasets collected at that location. Labels for which colors correspond to which types of traits are given in Panel (B). A set of 30 traits scored in Nova Porteirinha, Minas Gerais, Brazil (Queiroz *et al.* 2015; Paiva *et al.* 2017) are not visible in this panel. (B) Representation of seven broad phenotypic categories among the 234 traits collected here. Category assignments for individual traits are provided in Supplementary Table S1. (C) Distributions of narrow sense heritability values, calculated using the 2020 genetic marker dataset (Miao *et al.* 2020a), across the same seven broad phenotypic categories are shown in panel (B).
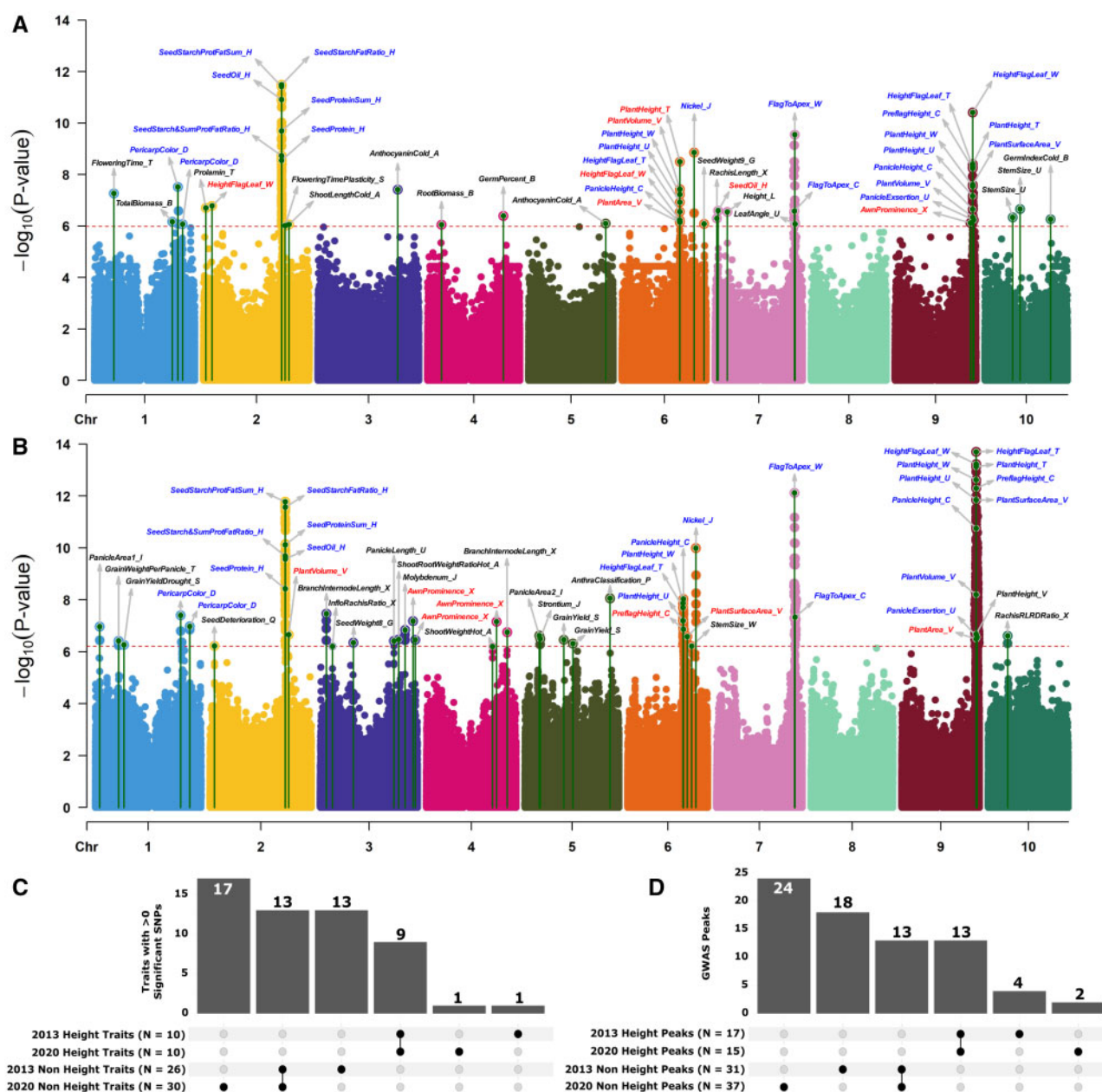
Each genetic marker dataset was employed to conduct GWAS for each of the 234 trait datasets. For each genetic marker dataset, imputation and filtering based on minor allele frequency specifically within the set of 304 shared lines were conducted using a common set of criteria (see Materials and methods). Similarly, a common trait data processing protocol was employed for each of the 234 trait datasets. This protocol incorporated normalization and outlier removal. GWAS for two of the 234 trait datasets produced questionable results, including distributions of observed *P*-values inconsistent with expectations, when analyzed with one of the marker sets (Supplementary Figure S3, A–D). The GWAS results for these two trait datasets were removed from downstream analyses for both marker sets.

Among the remaining 232 trait datasets, 36 trait datasets produced at least one significant marker-trait associations (MTAs) when analyzed using the 2013 marker dataset ($N = 48$ significant peaks). These 48 total significant peaks in individual GWAS with 2013 dataset localized to 26 unique regions of the genome as a result of repeated identification of the same genomic intervals in analyses of different trait datasets (Figure 2A). When the same 232 trait datasets were analyzed using the 2020 genetic marker dataset, 40 trait datasets produced at least one significant peak ($N = 52$ significant peaks). These 52 total significant peaks localized to 32 unique locations on the sorghum genome (Figure 2B). In analyses with each of the two genetic marker datasets, multiple distinct signals were observed in the same region of chromosome six including the canonical signal for *dw2*, as well as a separate peak for nickel abundance (both genetic marker datasets), seed weight (2013 genetic marker dataset) and stem size (2020 genetic marker dataset). The clustering of separate peaks may reflect increased statistical power resulting from elevated minor allele frequencies in this interval from over representation of the BTx406 haplotype among sorghum conversion lines in this region (Thurber *et al.* 2013).

A total of 54 traits exhibited at least one significant trait-associated marker, with 36 traits exhibiting at least one peak in the 2013 dataset and 40 traits exhibiting at least one peak in the 2020 dataset. Among these 54 traits a total of 22 traits exhibited at least one significant peak that was shared between the two marker datasets, while there were 14 unique traits in the 2013

marker set and 18 unique traits in the 2020 marker data set, each of which exhibited at least one unique peak (Supplementary Figure S4, A and B). Of the 22 traits which exhibited at least one significant peak when analyzed with each of the genetic marker datasets, nine were height traits. Among the 32 traits which exhibited at least one significant peak when analyzed with one genetic marker dataset and no significant peaks when analyzed with the other, only two were height-related traits. The non-representative nature of height-related traits may be explained both by the presence of three segregating large effect mutations for height in this population and the large LD blocks which exist around these genes as a result of selection during the temperate adaptation process (Thurber *et al.* 2013). Excluding height-related trait datasets, 13 trait datasets produced at least one significant MTA, when analyzed with either genetic marker dataset and 30 trait datasets produce at least one significant MTA, when analyzed with one and only one of the genetic marker datasets (Supplementary Figure S4C).

However, even when at least one significant MTA was identified when the same trait dataset was analyzed with each of the two genetic marker datasets, these MTAs may not correspond to the same causal loci. Of 74 unique MTAs between a given trait dataset and a given genomic interval identified using the two genetic marker datasets, 26 of the same MTAs were identified using both genetic marker datasets, 22 were identified only using the 2013 dataset and 26 were identified using only the 2020 dataset (Figure 2; Supplementary Figure S4, D and E). Peaks associated with plant height were disproportionately likely to be identified in analyses using both genetic marker datasets. Nineteen total peaks associated with plant height traits were identified of which 13 were identified when the same trait was analyzed with either genetic marker dataset. Excluding plant height-related traits, 55 distinct MTAs were identified between variation in a trait dataset and a given region of the genome. A total of 13 MTAs between non-height traits and a given region of the genome were identified consistently when using each of the two genetic marker datasets. Eighteen MTAs between non-height traits and a given region of the genome were identified only when using the 2013 genetic marker dataset and 24 only when using the 2020 genetic marker dataset (Supplementary Figure S4F). The total non-height

**Figure 2.** Combined Manhattan plots comparing MTAs identified using different marker datasets for the same individuals. (A) Combined Manhattan plot for 36 traits with at least one significant GWAS hit when analyzed using the 2013 genetic marker dataset and considering data from only those 304 sorghum lines genotyped in both the 2013 and 2020 datasets. Green lines topped with circles indicate the physical position and -log$_{10}$ P-value for the single most significant SNP within a GWAS peak identified for a particular trait. Text labels for individual traits employ trait names provided in Supplementary Table S1. Dashed red line indicates the cutoff for statistical significance calculated from the effective SNP number in the 2013 genetic marker dataset. (B) Combined Manhattan plot for 40 traits with at least one significant GWAS hit when analyzed using the 2020 genetic marker dataset and considering data from only those 304 sorghum lines genotyped in both the 2013 and 2020 datasets. Locations and P-values of the most significant SNP within each peak and statistical significance cutoff labeled as above. Blue labels indicate peaks shared between datasets. Red labels indicate traits where at least one significant GWAS peak is identified in both datasets but none of the peaks are shared between datasets. Black labels indicate traits where one or more significant GWAS peaks were identified with one of the marker datasets but no significant GWAS peaks were identified when the other marker dataset was employed. (C) Relationship between the identification of one or more significant GWAS peaks for a given trait dataset in each of the two genetic marker datasets. (D) Number of GWAS peaks which were either identified using both or only one of the two genetic marker datasets tested.

MTAs which would likely be detectable based on allele frequency and effect size in this population of 304 sorghum lines with sufficient numbers of markers was estimated to be approximately 85 (Lincoln Index Method). GWAS with either one of the two genetic marker datasets identified only 35–43% of the estimated total number of MTAs and the combined analysis identified only 63% of the estimated total number of potentially discoverable MTAs

(given sufficient marker density). Hence, increases in the number of genetic markers scored in this population would likely enable the discovery of 50–200% more MTAs when analyzing existing published trait data via single trait GWAS. The Lincoln Index Method, based on the size of two independent samples and the number of overlapping individuals between the two populations is a statistical measure used in several fields to estimate the

number of cases that have not yet been observed based on two independent sets of observed cases (Lincoln 1930). In this case, this approach likely underestimates the true number of detectable but unobserved associations as it assumes complete independence between the two independently observed sets. In reality, extremely large effect loci are more likely to be observed in both datasets even when in only modest LD with a genotyped marker, while smaller effect loci which would be detectable when in high LD with a genotyped marker is more likely to be missed when the LD to the most linked genotyped marker is lower. As a result, this estimate should be treated as a lower bound.

## Limited evidence for pleiotropy from conventional genome-wide association

GWAS was conducted for the complete set of 234 traits and all genetic marker data for the set of 343 accessions present in the 2020 marker dataset with the goal of testing for evidence of pleiotropy for quantitative traits. Analysis was conducted using three distinct statistical approaches to GWAS: GLM, MLM, and FarmCPU. The significant peaks identified by each of these methods are provided on FigShare (see data availability statement) for researchers interested in obtaining lists of loci controlling variation in specific traits. However, here we specifically present the results of MLM-based GWAS. With this larger population of individuals—343 vs the 304 shared between the 2013 and 2020 datasets—a total of 56 significant peaks were identified across 43 traits. After identifying and merging associations between distinct traits within the same genomic intervals, the set of 56 significant peaks collapsed to 31 regions of the genome (Figure 3).

Of these 31 unique genomic regions, 25 were identified in the analysis of only a single trait dataset. Among the remaining six cases where two or more trait datasets identified signals in the same genomic regions, three were identified in the analysis of only two trait datasets. The final three intervals (on chromosomes 2, 6, and 9) were associated with 7, 6, and 13 traits, respectively. The peak on chromosome 2 was identified in GWAS for seed composition traits including oil, protein, and the sums and ratios of seed oil, protein, and starch and likely corresponds to the putative alpha-amylase 3 gene, previously identified in Rhodes *et al.* (2017b). The peaks on chromosomes 6 and 9 correspond to ma1/dw2 (Klein *et al.* 2008; Murphy *et al.* 2011; Hilley *et al.* 2017) and dw1 (Hilley *et al.* 2016; Yamaguchi *et al.* 2016), respectively. Traits associated with these two genomic intervals include measures of both plant height and plant volume/area (Table 2).

One interval on chromosome 7 where a single genomic interval contained MTAs for two and only two trait datasets was the result of measurement of the same trait—distance from the flag leaf to the plant apex—in two different studies conducted in different years in different locations by different research groups. An interval on chromosome 5 was associated with two traits from the same publication, which scored anthracnose resistance in two different ways (Cuevas *et al.* 2018). A third interval on chromosome 4 was associated with both branch length in the inflorescence and acid detergent fiber within the grain. This interval on chromosome 4 may represent genuine pleiotropy or two distinct functional variations in different genes separated by ≤500 kilobases. The relative dearth of evidence for pleiotropy in sorghum is consistent with previous quantitative genetic investigations of pleiotropy in maize for both inflorescence architecture and leaf morphology (Brown *et al.* 2011; Tian *et al.* 2011).

However, given the large number of false negatives expected in any individual GWAS (Korte and Farlow 2013), quantifying
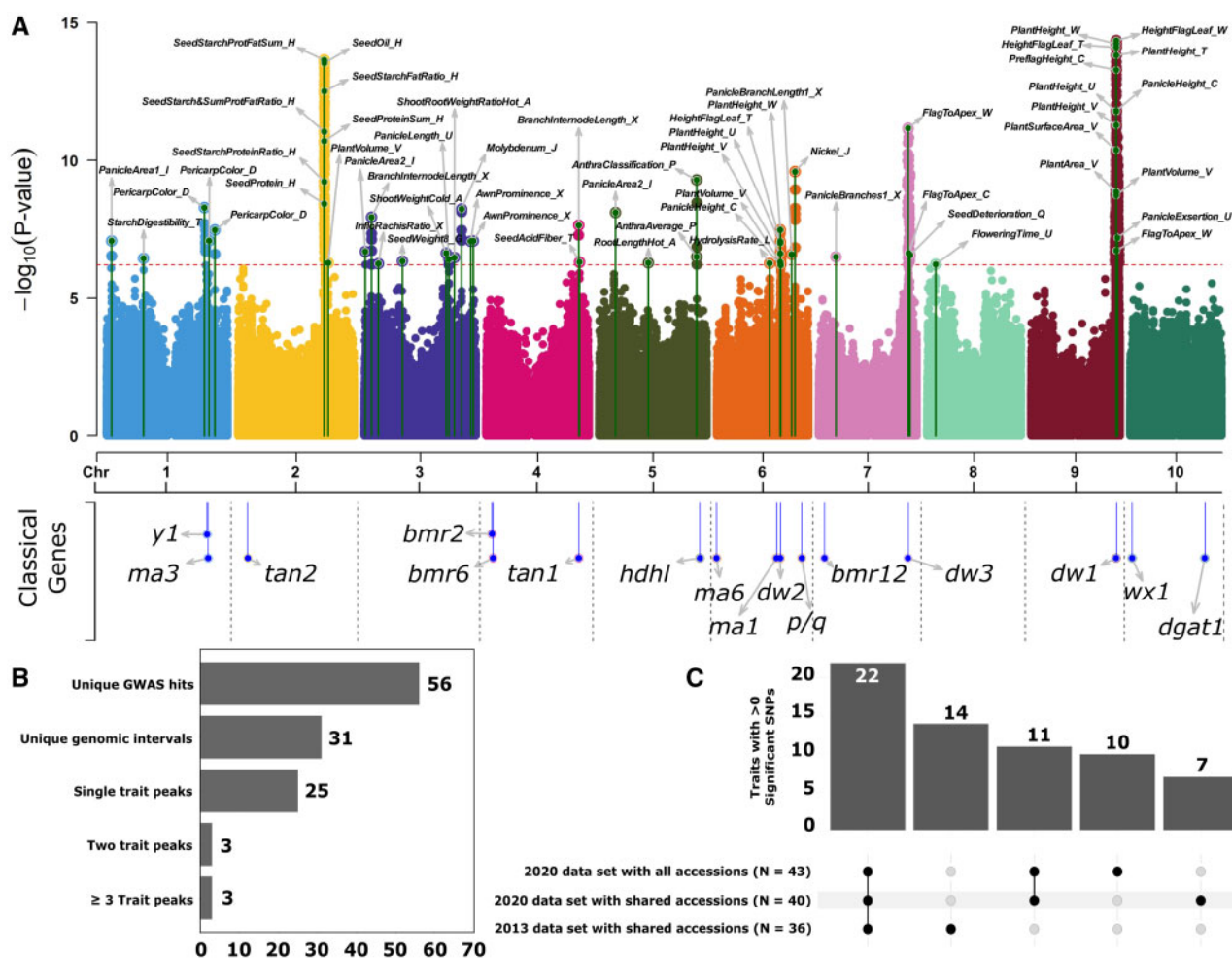
how many intersections exist between independently conducted sets of GWAS in modestly powered populations is likely to underestimate the true extent of pleiotropy (Visscher and Yang 2016). Formal multivariate GWAS approaches face difficulty scaling to large datasets (>3–5 traits) (Zhou and Stephens 2014; Rice *et al.* 2020). Hence, here we employed a multivariate adaptive shrinkage approach to the initial output from MLM based GWAS to estimate the effect of individual markers on separate trait datasets (Urbut *et al.* 2019). This approach provides both a test of which markers are significantly associated with phenotypic variation across the population while also estimating which specific traits a given marker had non-zero effects on, and the directions of those effects.

## Improved interpretability of trait-associated loci under joint analysis

Joint analysis was conducted using MashR for 176 traits, excluding 30 traits scored on no more than 100 individuals, 6 binary traits and 22 ionomics trait data (Shakoor *et al.* 2016). Standard error and effect size from MLM based GWAS were employed for this analysis. While FarmCPU has been shown to exhibit greater power to detect more total causal loci, the inclusion of identified loci as covariates means that only a single marker is identified per locus. If different markers in LD with the same causal locus were identified in FarmCPU based analysis of different traits, the pleiotropic effects of this locus would be undetectable by MashR. A set of 593 markers were identified which both exhibited an association with at least one phenotype with a local false sign rate (*lfsr*) <0.001, and for which the ratio of the likelihood of one or more significant phenotypic effects at an SNP to the likelihood that the SNP had only null effects was estimated to be <$10^4$, which is referred to as the Bayes factor (Urbut *et al.* 2019). An analog of the false discovery rate, *lfsr* requires true discoveries to be not only nonzero but also correctly signed (Stephens 2017). These 593 markers cluster together in 44 unique peaks across the sorghum genome (Figure 4). Within each multi-marker peak, the single marker with the largest Bayes factor was employed to represent the peak. The number of traits upon which markers had significantly nonzero effects ranged from 1 to 141 for a peak on chromosome 6 corresponding to two known large effect loci, dw2 and ma1 (Figure 4, A and B). Figure 4A employs an *lfsr* <0.001 threshold for reporting marker trait associations. A parallel analysis employing an *lfsr* threshold of <0.05 provided roughly equivalent results Supplementary Figure S6. The relationship between the Bayes factor assigned to a peak and the number of trait datasets with which it was associated was not straightforward (Figure 4B; Supplementary Figure S7). The single peak with the largest Bayes factor was associated with only twelve trait datasets, a number of which were independent measures of the same traits in different environments. Other peaks with comparatively modest Bayes factors were associated with modest effects on 90 traits in our datasets (Figure 4A).

Multi-trait analysis was able to recover a number of known pleiotropic features for large effect loci segregating in the population. In addition to plant height, the peak at dw2 was associated with multiple metrics of root size/area, panicle length, plant surface area, and seed weight (Supplementary Figure S8A). The effects of dw2 on panicle length, seed weight, and leaf area have been previously reported (Graham and Lessman 1966; Pereira and Lee 1995) while the reductions in multiple metrics of root size/area associated with the dwarfing allele of dw2 had not. The apparent impact of dw2 on root phenotypes, suggests that the gene may play equivalent roles in determining size of below and

**Figure 3.** Combined Manhattan plot for GWAS using all 343 individuals genotyped in the 2020 SNP set. (A) Combined Manhattan plot for 43 traits with at least one significant GWAS hit when analyzed using the 2020 genetic marker dataset and all 343 sorghum lines genotyped in the 2020 genetic marker dataset. Green lines topped with circles indicate the physical position and -log$_{10}$ P-value for the single most significant SNP within a GWAS peak identified for a given trait. Text labels employ trait names provided in Supplementary Table S1. Dashed red line indicates the cutoff for statistical significance calculated from the effective SNP number in the 2020 genetic marker dataset. Lower panel indicate positions of a set of cloned sorghum mutants, taken from (Boyles *et al.* 2019). Estimates of LD among summit SNPs of each peak are shown in Supplementary Figure S5A. (B) Summary of results from GWAS analysis using all 343 SAP lines included in the 2020 marker dataset. (C) Number of traits where one or more significant GWAS peaks were identified in the 2013 dataset considering only accessions shared with the 2020 dataset, the 2020 dataset considering only accessions shared with the 2013 dataset, and/or all accessions in the 2020 dataset.

**Table 2** Summary of the GWAS results when data from all 343 accessions in the 2020 marker set are employed

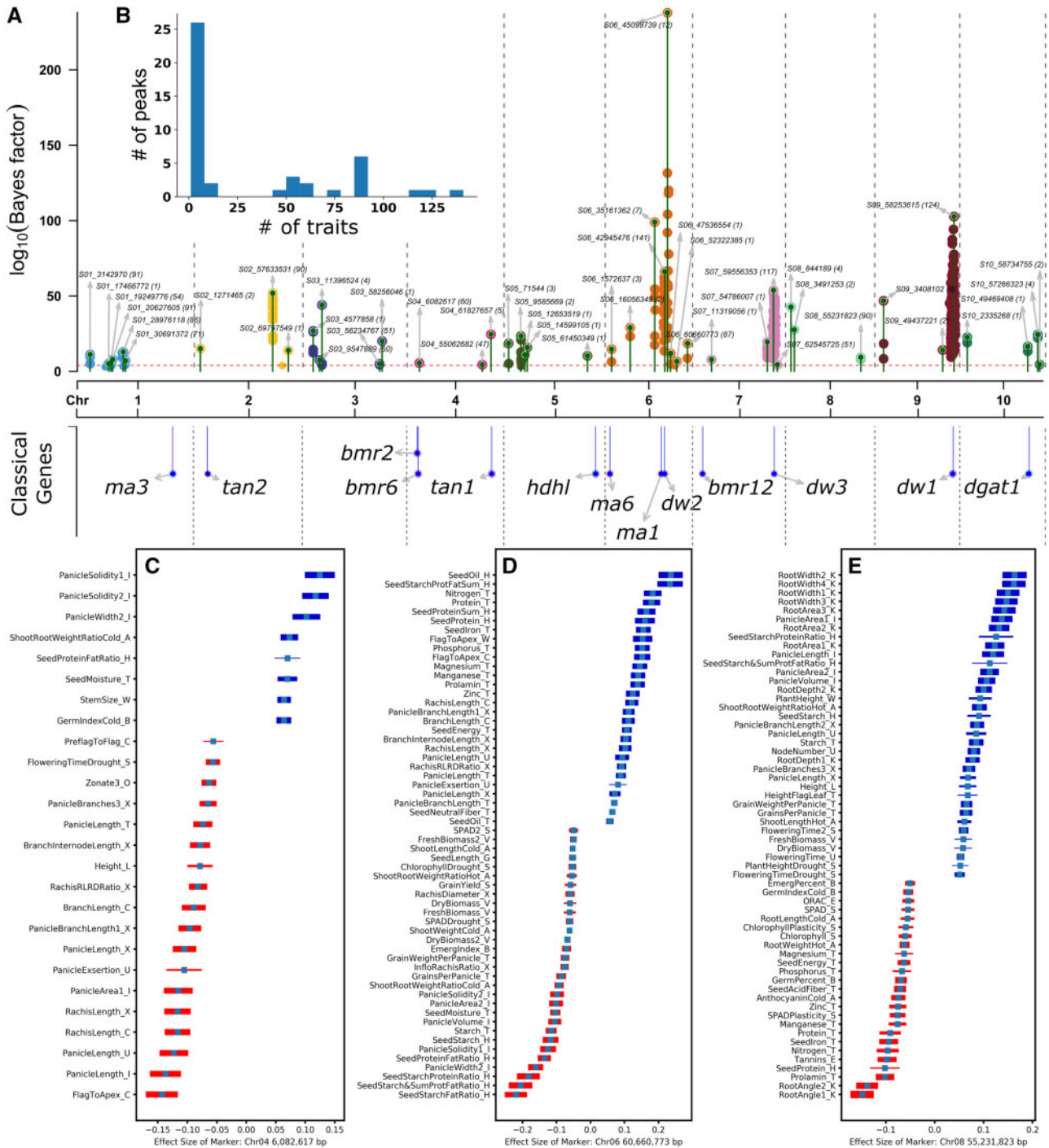| Chr # | GWAS hits | Unique genomic regions[a] | Single trait peaks | Two trait peaks | ≥3 trait peaks |
|---|---|---|---|---|---|
| Chr 1 | 5 | 5 | 5 | 0 | 0 |
| Chr 2 | 8 | 2 | 1 | 0 | 1[b,c] |
| Chr 3 | 10 | 10 | 10 | 0 | 0 |
| Chr 4 | 2 | 1 | 0 | 1 | 0 |
| Chr 5 | 4 | 3 | 2 | 1[c] | 0 |
| Chr 6 | 9 | 4 | 3 | 0 | 1[c,d] |
| Chr 7 | 4 | 3 | 2 | 1[c] | 0 |
| Chr 8 | 1 | 1 | 1 | 0 | 0 |
| Chr 9 | 13 | 2 | 1 | 0 | 1[c,e] |
| Chr 10 | 0 | 0 | 0 | 0 | 0 |
| Total | 56 | 31 | 25 | 3 | 3 |

[a]GWAS hits within 500 kb of each other on the genome were merged into a single interval. Given the low incidence of observed pleiotropy, a conservatively large interval (greater than the 50–350 kb reported range for LD decay in sorghum) was selected to reduce the incidence of false negatives (*i.e.* true cases of pleiotropy effects misclassified as independent signals from distinct loci).
[b]The locus on Chr2 associated with ≥3 traits is associated with seven total traits all associated with seed composition: SeedProteinSum_H, SeedStarch&SumProtFatRatio_H, SeedStarchProtFatSum_H, SeedProtein_H, SeedStarchProteinRatio_H, SeedStarchFatRatio_H, and SeedOil_H.
[c]While these peaks were identified for multiple datasets, the datasets all represent independent measures of similar phenotypes.
[d]The locus on Chr6 associated with ≥3 traits is associated with six total traits all associated with plant height: PlantVolume_V, PlantHeight_U, PlantHeight_V, HeightFlagLeaf_T, PlantHeight_W, and PanicleHeight_C.
[e]The locus on Chr9 associated with ≥3 traits is associated with 13 total traits all associated with plant height: PlantArea_V, PanicleHeight_C, PanicleExsertion_U, HeightFlagLeaf_T, PlantHeight_T, PreflagHeight_C, HeightFlagLeaf_W, PlantHeight_W, FlagToApex_W, PlantVolume_V, PlantHeight_U, PlantHeight_V, and PlantSurfaceArea_V.

**Figure 4.** Pleiotropic analysis of SAP phenotypes. (A) Markers assigned significant Bayes factor values in MashR analysis. Green lines topped with circles indicate the physical position and log$_{10}$ Bayes factor for the most significant SNP within a peak identified for a pleiotropic loci. Text labels indicate the position and name of the most significant marker within each peak. The number of trait datasets significantly associated with a marker at *lfsr* <0.001 is indicated in brackets. It should be noted that trait datasets include both measurements of different traits and the same trait scored across different environments in different studies. Dashed red line indicates the cutoff for statistical significance at log$_{10}$ Bayes factor of 4. Estimates of linkage disequilibrium among the summit SNPs of each distinct peak are shown in Supplementary Figure S5B. (B) Distribution of the number of trait datasets which were significantly associated with each unique peak. (C) Distribution of effect sizes and directions of a subset of the 60 trait datasets for which the genetic marker S04_6082617 has a significant effect (*lfsr* < 0.001). To aid readability, only the subset of trait datasets where the effect size is >0.05 or <0.05 are shown. Bar thickness is proportional to the relative estimated statistical significance of each association with the thickest bars marking the most significantly associated trait for a given marker and the thinnest the least significantly associated trait for a given marker which still passed all filtering criteria. (D) Distribution of effect sizes and directions for a subset of the 87 trait datasets for which the genetic marker S06_60660773 has significant effects. Cutoffs for visualization are the same as applied for panel C. (E) Distribution of effect sizes and directions for a subset of the 90 trait datasets for which the genetic marker (S08_55231823) has significant effects. Cutoffs for visualization are the same as applied for panels C and D.

above ground plant organs. Another plant height-related gene, *dw3* was previously known to have effects on grain yield (Cassady 1965), leaf angle (Truong *et al.* 2015), biomass (George-Jaeggli *et al.* 2011), internode length (Brown *et al.* 2008), stem diameter (Olatoye *et al.* 2020), panicle exsertion (Zhao *et al.* 2016), and panicle architecture (Brown *et al.* 2006). Here, statistically significant links were also observed between *dw3* and variation in plant biomass, leaf angle, stem size/stem diameter, grain yield plasticity, internode length, panicle exsertion, and panicle architecture (panicle length, width, and area) (Supplementary Figure S8B). In addition to plant height, *dw1* is also known to alter biomass and biomass associated traits (Breitzman *et al.* 2019), internode length and lodging resistance in sorghum (Hilley *et al.* 2016; Yamaguchi *et al.* 2016). The set of traits assembled here did not include measurements of either internode length or lodging resistance across the SAP so it was not possible to assess whether these known pleiotropic effects of *dw1* on these traits were recovered. However, the peak associated with *dw1* in chromosome 9 was significantly linked to variation in above ground biomass traits (directly measured biomass, plant surface area), as well as multiple metrics for root size/area (Supplementary Figure S8C). In contrast to *dw3* and similar to *dw2*, *dw1* may play equivalent roles in determining organ size for both below ground and above ground plant organ systems.

Multi-trait analysis also recovered a number of novel signals across the genome. The pairing and direction of effect sizes for these traits enables greater interpretability of the resulting MTAs. In some cases, these are straightforward trade-offs. The trait-associated marker located at 6.08 MB on chromosome 4 one allele is associated with longer panicles, but these panicles are also narrower and less dense. The other allele present at this locus produces shorter, fatter, and denser panicles, resulting in increases in seed moisture levels at harvest (Figure 4C). A number of other multiple trait associations identified were also consistent known trade-offs in plant growth and development. A trait-associated marker located at 60.66 MB on chromosome 6 is associated with increases in seed oil and protein content, and many important micronutrients. However, the same allele is also associated with decrease in panicle volume and solidity (high throughput phenotyping) and decreases in grains per panicle and grain weight per panicle (conventional phenotyping) (Figure 4D). Multiple trait associations can also reveal explanations for potential associations which would otherwise be potential breeding targets. For example, improving root architecture has been proposed as a target for enhancing drought tolerance or nutrient uptake (Paez-Garcia *et al.* 2015). A marker on chromosome 8 located at 55.23 MB had large effects on multiple root traits including greater root width, larger total root area, and a smaller root angle (Figure 4E). In isolation, this might appear a promising target for root-based breeding. However, multiple trait analysis also identified that this allele is associated with delays in flowering time and increases in total node numbers. These results suggest that the observed increases in root extent and root area may be an indirect result of delayed vegetative to reproductive transition.

## Discussion

The SAP has been widely adopted and proven to be a long-lasting resource for the sorghum genetics community. The syntenic conservation of GWAS hits between sorghum and maize means the SAP also provides information on gene function in maize (Zheng *et al.* 2020). In the interval between the start of the analyses in this paper and submission for publication, at least five additional studies employing this population have been posted online including studies of provitamin A (Cruet-Burgos *et al.* 2020), geospatial association with parasitic plants (Bellis *et al.* 2020), resistance to different fungal sources of grain mold (Prom *et al.* 2020), genetic determinants of the root-associated microbiome (Deng *et al.* 2021), and herbicide resistance (Pandian *et al.* 2020). Our results suggest that simply increasing the marker density of the SAP—and similar community association populations—may more than double the number of true positive MTAs detection in future studies with the same population. Care should be taken to record and disseminate accession-specific trait measurements from GWAS in ways that facilitate future reanalyses as additional genetic marker datasets become available. In our study, we identified 40 papers which included the collection of one or more new trait datasets from the SAP. Through data curation and annotation, we have increased the proportion of papers from 50% to 64% where traits have been publicly released with IDs which can be associated back to genetic marker data, facilitating reuse and reanalysis (Table 1). However, adopting both community norms that emphasize the need to store and disseminate trait datasets, as well as developing a central repository for sorghum phenotypic data would likely increase the proportion of trait datasets generated by the sorghum genetics community, which will continue to contribute to new discoveries and understanding in years to come. Similarly, it will be important to maintain and distribute seed of the SAP to lower barriers to entry into sorghum quantitative genetics for new research groups and to avoid the risk of failed or misleading results due to lines that are swapped or duplicated. Seed is currently maintained and distributed by the USDA NPGS; however, resource constraints at NPGS can limit how often the lines of this panel can be increased. Informal lab to lab distribution has acted as a fallback source of seed. However, in the long term, this approach runs the risk of propagating swaps, labeling errors, or pollen contamination, reducing the comparability of data generated by different research groups with different germplasm sources. An analysis of published RNA-seq data labeled as coming from the maize reference inbred B73 found at least three distinct clades of genetically distinct B73 accessions and that relationships between these samples recapitulated advisor–advisee relationships (Liang and Schnable 2016). In *Arabidopsis thaliana*, a widely used commercial source for the reference inbred Col-0 was found to contain substantial introgressions of non Col-0 origin (Shao *et al.* 2016). Storage and dissemination of trait data and associated metadata for future studies will aid both in detecting new associations as genetic marker datasets increase in density, and in expanding our knowledge of pleiotropy, as shown here. Additionally, this community-based strategy will facilitate the development and validation of predictions from empirical crop growth models, and the investigation of the genetic basis of phenotypic plasticity and genotype by environment interactions.

Strong selection will often act on rare, large effect, and pleiotropic loci (Orr 2000). Here, loci identified in a joint analysis of 176 trait datasets tended to fall into one of two categories, either showing associations with large (>40) number of trait datasets or small (<10) trait datasets, with these datasets often representing measures of the same trait in multiple environments or multiple distinct but highly correlated traits (Figure 4B). This pattern does not appear to be an artifact introduced by variation in statistical power as a result of either effect size or allele frequency as both loci associated with many traits or associated with only several traits include examples with both high and low Bayes factor values (Supplementary Figure S7). This stands in contrast to studies

in the related species of maize where little evidence has been found for pleiotropic quantitative genetic loci segregating in populations (Tian *et al.* 2011; Wallace *et al.* 2014), but this difference should be interpreted cautiously until and unless similar wide-scale multivariate analyses are conducted in maize association populations, given the differences in both methodological approaches and patterns of LD. An analysis of historical yield data in common bean employing MashR also identified two genomic intervals associated with pleiotropic effects on different phenotypes (MacQueen *et al.* 2020). If the difference in the prevalence of pleiotropy between maize and sorghum continues to be observed in additional studies, it may reflect the distinct histories of both maize and sorghum in temperate North America, and the distinct histories of widely used association panels in each species. The first reports of sorghum cultivation in the southeastern United States date to either 1838 or 1855, likely as the result of introduction from the Caribbean (Vinall *et al.* 1936). Two temperate adapted strains of sorghum were introduced into the Great Plains approximately 150 years ago followed by rapid selection by farmers for earlier flowering and shorter stature (Quinby 1975). Temperate maize in the United States is much older with adaptation to temperate highlands occurring over an approximate 2000 years period starting 4000 years ago, allowing for more gradual selection and therefore less likely to capture pleiotropic loci (Swarts *et al.* 2017). Similarly, many of the lines in the SAP trace their origin to a conversion process where genes needed for temperate adaptation were introgressed through multiple generations of strong phenotypic selection (Stephens *et al.* 1967; Casa *et al.* 2008), while both the most widely employed maize association panel and the maize nested association panel were assembled from lines already adapted to the temperate United States (Flint-Garcia *et al.* 2005; Gage *et al.* 2020). However, a key limitation of single marker level analyses, including both MashR and conventional GWAS, is that the estimated effects of a given SNP reflect not only the effect of that SNP itself but also all SNPs in LD with the tested marker (Urbut *et al.* 2019). Given the high degree of LD observed in the SAP Morris *et al.* (2013a) and the absence of saturation level genetic marker data, it is not possible to rule out that any given combination of variation in multiple traits associated with a single marker may result from multiple causal polymorphisms in the same or adjacent genes.

As high throughput phenotyping approaches become more widely adopted, direct measurements of pleiotropy may become more feasible. Once sensor datasets are collected (*e.g.* RGB images, LIDAR point clouds, hyperspectral data cubes) and algorithms for numerically quantifying specific traits are developed, the additional cost extracting measurements of non-target traits from the same sensor data is minimal. A better understanding of pleiotropic relationships for specific loci and across groups of plant traits may aid significantly in reducing inadvertent selection and prioritizing candidate loci for introgression into elite germplasm for sorghum and related species. A greater understanding of potential pleiotropic effects may help to prioritize which off target phenotypic effects should be tested for either when evaluating natural variants or when generating gene edits.

## Author contributions

R.V.M., M.G., and J.C.S. conceived of the study. A.D., S.S., R.E.B., M.G.S.F., P.S.S., B.S., and S.K. conducted the experiments and collected the data. R.V.M., M.G., and C.M. analyzed the data. R.V.M. and J.C.S. wrote the manuscript. All authors read and approved the final manuscript.

## Conflicts of interest

J.C.S. and P.S.S. have equity interests in Data2Bio LLC, a company that provides genotyping services using the technology used to generate the 2020 genetic marker set employed in this study. The authors declare no other conflicts of interest.

## Literature cited

Adeyanju A, Little C, Yu J, Tesso T. 2015. Genome-wide association study on resistance to stalk rot diseases in grain sorghum. G3 (Bethesda). 5:1165–1175.

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

Bailey D. 2004. Recombinant inbred strains and bilineal congenic strains. In: HL Foster, JD Small, JG Fox, editors. The Mouse in Biomedical Research. New York: Academic Press. p. 223–239.

Bates D, Mächler M, Bolker B, Walker S. 2015. Fitting linear mixed-effects models using lme4. J. Stat. Softw. 67: 1–48. doi: http://dx.doi.org/10.18637/jss.v067.i01.

Bellis ES, Kelly EA, Lorts CM, Gao H, DeLeo VL, *et al.* 2020. Genomics of sorghum local adaptation to a parasitic plant. Proc Natl Acad Sci USA. 117:4243–4251.

Bouchet S, Pot D, Deu M, Rami J-F, Billot C, *et al.* 2012. Genetic structure, linkage disequilibrium and signature of selection in sorghum: lessons from physically anchored dart markers. PLoS One. 7:e33470.

Boyles RE, Brenton ZW, Kresovich S. 2019. Genetic and genomic resources of sorghum to connect genotype with phenotype in contrasting environments. Plant J. 97:19–39.

Boyles RE, Cooper EA, Myers MT, Brenton Z, Rauh BL, *et al.* 2016. Genome-wide association studies of grain yield components in diverse sorghum germplasm. Plant Genome. 9: 1–17.

Boyles RE, Pfeiffer BK, Cooper EA, Rauh BL, Zielinski KJ, *et al.* 2017. Genetic dissection of sorghum grain quality traits using diverse and segregating populations. Theor Appl Genet. 130:697–716.

Breitzman MW, Bao Y, Tang L, Schnable PS, Salas-Fernandez MG. 2019. Linkage disequilibrium mapping of high-throughput image-derived descriptors of plant architecture traits under field conditions. Field Crops Res. 244:107619.

Brown P, Klein P, Bortiri E, Acharya C, Rooney W, *et al.* 2006. Inheritance of inflorescence architecture in sorghum. Theor Appl Genet. 113:931–942.

Brown PJ, Rooney WL, Franks C, Kresovich S. 2008. Efficient mapping of plant height quantitative trait loci in a sorghum association population with introgressed dwarfing genes. Genetics. 180:629–637.

Brown PJ, Upadyayula N, Mahone GS, Tian F, Bradbury PJ, *et al.* 2011. Distinct genetic architectures for male and female inflorescence traits of maize. PLoS Genet. 7:e1002383.

Browning BL, Browning SR. 2016. Genotype imputation with millions of reference samples. Am J Hum Genet. 98:116–126.

Burr B, Burr FA, Thompson KH, Albertson M, Stuber C. 1988. Gene mapping with recombinant inbreds in maize. Genetics. 118:519–526.

Casa AM, Pressoir G, Brown PJ, Mitchell SE, Rooney WL, *et al.* 2008. Community resources and strategies for association mapping in sorghum. Crop Sci. 48:30–40.

Cassady A. 1965. Effect of a single height (dw) gene of sorghum on grain yield, grain yield components, and test weight 1. Crop Sci. 5:385–388.

Chen J, Chopra R, Hayes C, Morris G, Marla S, *et al.* 2017. Genome-wide association study of developing leaves' heat tolerance during vegetative growth stages in a sorghum association panel. Plant Genome. 10: 1–15.

Chopra R, Burow G, Burke JJ, Gladman N, Xin Z. 2017. Genome-wide association analysis of seedling traits in diverse sorghum germplasm under thermal stress. BMC Plant Biol. 17:12.

Covarrubias-Pazaran G. 2016. Genome-assisted prediction of quantitative traits using the R package sommer. PLoS One. 11:e0156744.

Cruet-Burgos C, Cox S, Ioerger BP, Perumal R, Hu Z, *et al.* 2020. Advancing provitamin a biofortification in sorghum: genome-wide association studies of grain carotenoids in global germplasm. Plant Genome. 13:e20013.

Cuevas HE, Fermin-Pérez RA, Prom LK, Cooper EA, Bean S, *et al.* 2019. Genome-wide association mapping of grain mold resistance in the us sorghum association panel. Plant Genome. 12:180070.

Cuevas HE, Prom LK, Cooper EA, Knoll JE, Ni X. 2018. Genome-wide association mapping of anthracnose (*Colletotrichum sublineolum*) resistance in the U.S. Sorghum Association Panel. Plant Genome. 11:170099.

Deng S, Caddell D, Yang J, Dahlen L, Washington L, *et al.* 2021. Genome wide association study reveals plant loci controlling heritability of the rhizosphere microbiome. The ISME Journal doi: 10.1038/s41396-021-00993-z.

Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, *et al.* 2011. A robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. PLoS One. 6:e19379.

Endelman JB, Jannink J-L. 2012. Shrinkage estimation of the realized relationship matrix. G3 (Bethesda). 2:1405–1413.

Fernandez MGS, Bao Y, Tang L, Schnable PS. 2017. A high-throughput, field-based phenotyping technology for tall biomass crops. Plant Physiol. 174:2008–2022.

Flint-Garcia SA, Thuillet A-C, Yu J, Pressoir G, Romero SM, *et al.* 2005. Maize association population: a high-resolution platform for quantitative trait locus dissection. Plant J. 44:1054–1064.

Gage JL, Monier B, Giri A, Buckler ES. 2020. Ten years of the maize nested association mapping population: Impact, limitations, and future directions. Plant Cell. 32:2083–2093.

George-Jaeggli B, Jordan DR, van Oosterom EJ, Hammer GL. 2011. Decrease in sorghum grain yield due to the *dw*3 dwarfing gene is caused by reduction in shoot biomass. Field Crops Res. 124:231–239.

Graham D, Lessman K. 1966. Effect of height on yield and yield components of two isogenic lines of *Sorghum vulgare* pers. 1. Crop Sci. 6:372–374.

Hamblin MT, Fernandez MGS, Casa AM, Mitchell SE, Paterson AH, *et al.* 2005. Equilibrium processes cannot explain high levels of short-and medium-range linkage disequilibrium in the domesticated grass sorghum bicolor. Genetics. 171:1247–1256.

Hilley J, Truong S, Olson S, Morishige D, Mullet J. 2016. Identification of dw1, a regulator of sorghum stem internode length. PLoS One. 11:e0151271.

Hilley JL, Weers BD, Truong SK, McCormick RF, Mattison AJ, *et al.* 2017. Sorghum dw2 encodes a protein kinase regulator of stem internode length. Sci Rep. 7:13.

Hufnagel B, de Sousa SM, Assis L, Guimaraes CT, Leiser W, *et al.* 2014. Duplicate and conquer: multiple homologs of phosphorus-starvation tolerance1 enhance phosphorus acquisition and sorghum performance on low-phosphorus soils. Plant Physiol. 166:659–677.

Klein RR, Mullet JE, Jordan DR, Miller FR, Rooney WL, *et al.* 2008. The effect of tropical sorghum conversion and inbred development on genome diversity as revealed by high-resolution genotyping. Crop Sci. 48:12.

Kong W, Guo H, Goff VH, Lee T-H, Kim C, *et al.* 2014. Genetic analysis of vegetative branching in sorghum. Theor Appl Genet. 127:2387–2403.

Korte A, Farlow A. 2013. The advantages and limitations of trait analysis with GWAS: a review. Plant Methods. 9:29.

Lasky JR, Upadhyaya HD, Ramu P, Deshpande S, Hash CT, *et al.* 2015. Genome-environment associations in sorghum landraces predict adaptive traits. Sci Adv. 1:e1400218.

Lee M, Sharopova N, Beavis WD, Grant D, Katt M, *et al.* 2002. Expanding the genetic map of maize with the intermated B73× Mo17 (IBM) population. Plant Mol Biol. 48:453–461.

Li M-X, Yeung JM, Cherny SS, Sham PC. 2012. Evaluating the effective numbers of independent tests and significant p-value thresholds in commercial genotyping arrays and public imputation reference datasets. Hum Genet. 131:747–756.

Li X, Li X, Fridman E, Tesso TT, Yu J. 2015. Dissecting repulsion linkage in the dwarfing gene dw3 region for sorghum plant height provides insights into heterosis. Proc Natl Acad Sci USA. 112:11823–11828.

Liang Z, Schnable JC. 2016. RNA-Seq based analysis of population structure within the maize inbred B73. PLoS One. 11:e0157942.

Lincoln FC. 1930. Calculating Waterfowl Abundance on the Basis of Banding Returns. Number 118, US Department of Agriculture.

Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, *et al.* 2011. Fast linear mixed models for genome-wide association studies. Nat Methods. 8:833–835.

Lister C, Dean C. 1993. Recombinant inbred lines for mapping RFLP and phenotypic markers in *Arabidopsis thaliana*. Plant J. 4:745–750.

Liu X, Huang M, Fan B, Buckler ES, Zhang Z. 2016. Iterative usage of fixed and random effect models for powerful and efficient genome-wide association studies. PLoS Genet. 12:e1005767.

Mace E, Innes D, Hunt C, Wang X, Tao Y, *et al.* 2019. The Sorghum QTL atlas: a powerful tool for trait dissection, comparative genomics and crop improvement. Theor Appl Genet. 132:751–766.

Mace ES, Tai S, Gilding EK, Li Y, Prentis PJ, *et al.* 2013. Whole-genome sequencing reveals untapped genetic potential in Africa's indigenous cereal crop sorghum. Nat Commun. 4:9.

MacQueen AH, White JW, Lee R, Osorno JM, Schmutz J, *et al.* 2020. Genetic associations in four decades of multienvironment trials reveal agronomic trait evolution in common bean. Genetics. 215: 267–284.

McMaster N, Acharya B, Harich K, Grothe J, Mehl HL, *et al.* 2019. Quantification of the mycotoxin deoxynivalenol (DON) in sorghum using GC-MS and a stable isotope dilution assay (SIDA). Food Anal Methods. 12:2334–2343.

Miao C, Xu Y, Liu S, Schnable PS, Schnable JC. 2020a. Increased power and accuracy of causal locus identification in time series genome-wide association in sorghum. Plant Physiol. 183: 1898–1909.

Miao C, Xu Z, Rodene E, Yang J, Schnable JC, *et al.* 2020b. Semantic segmentation of sorghum using hyperspectral data identifies genetic associations. Plant Phenomics. 2020:4216373.

Miao C. 2020. schnablelab. https://github.com/freemao/schnablelab.

Moghimi N, Desai JS, Bheemanahalli R, Impa SM, Vennapusa AR, *et al.* 2019. New candidate loci and marker genes on chromosome 7 for improved chilling tolerance in sorghum. J Exp Bot. 70: 3357–3371.

Morris GP, Ramu P, Deshpande SP, Hash CT, Shah T, *et al.* 2013a. Population genomic and genome-wide association studies of agroclimatic traits in sorghum. Proc Natl Acad Sci USA. 110:453–458.

Morris GP, Rhodes DH, Brenton Z, Ramu P, Thayil VM, *et al.* 2013b. Dissecting genome-wide association signals for loss-of-function phenotypes in sorghum flavonoid pigmentation traits. G3 (Bethesda). 3:2085–2094.

Murphy RL, Klein RR, Morishige DT, Brady JA, Rooney WL, *et al.* 2011. Coincident light and clock regulation of pseudoresponse regulator protein 37 (PRR37) controls photoperiodic flowering in sorghum. Proc Natl Acad Sci USA. 108:16469–16474.

Mutava R, Prasad P, Tuinstra M, Kofoid K, Yu J. 2011. Characterization of sorghum genotypes for traits related to drought tolerance. Field Crops Res. 123:10–18.

Nordborg M, Hu TT, Ishino Y, Jhaveri J, Toomajian C, *et al.* 2005. The pattern of polymorphism in *Arabidopsis thaliana*. PLoS Biol. 3:e196.

Olatoye MO, Hu Z, Morris GP. 2020. Genome-wide mapping and prediction of plant architecture in a sorghum nested association mapping population. Plant Genome 13: e20038.

Olatoye MO, Marla SR, Hu Z, Bouchet S, Perumal R, *et al.* 2019. Dissecting adaptive traits with nested association mapping: genetic architecture of inflorescence morphology in sorghum. G3: Genes, Genomes, Genetics, 10, 1785–1796.

Orr HA. 2000. Adaptation and the cost of complexity. Evolution. 54: 13–20.

Ortiz D, Hu J, Salas Fernandez MG. 2017. Genetic architecture of photosynthesis in sorghum bicolor under non-stress and cold stress conditions. J Exp Bot. 68:4545–4557.

Ott A, Liu S, Schnable JC, Yeh C-T, Wang K-S, *et al.* 2017. tGBS® genotyping-by-sequencing enables reliable genotyping of heterozygous loci. Nucleic Acids Res. 45:e178.

Paez-Garcia A, Motes CM, Scheible W-R, Chen R, Blancaflor EB, *et al.* 2015. Root traits and phenotyping strategies for plant improvement. Plants (Basel). 4:334–355.

Paiva CL, Queiroz VAV, Simeone MLF, Schaffert RE, de Oliveira AC, *et al.* 2017. Mineral content of sorghum genotypes and the influence of water stress. Food Chem. 214:400–405.

Pandian BA, Varanasi A, Vennapusa AR, Rajendran S, Lin G, *et al.* 2020. Characterization, Genetic Analyses, and Identification of QTLs Conferring Metabolic Resistance to a 4-Hydroxyphenylpyruvate Dioxygenase Inhibitor in Sorghum (Sorghum bicolor). Front. Plant Sci. 11: 1890.

Pereira M, Lee M. 1995. Identification of genomic regions affecting plant height in sorghum and maize. Theor Appl Genet. 90: 380–388.

Perez MBM, Zhao J, Yin Y, Hu J, Fernandez MGS. 2014. Association mapping of brassinosteroid candidate genes and plant architecture in a diverse panel of sorghum bicolor. Theor Appl Genet. 127: 2645–2662.

Peterson R. 2017. Estimating normalization transformations with bestnormalize. Httpsgithub CompetersonRbestNormalize.

Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, *et al.* 2006. Principal components analysis corrects for stratification in genome-wide association studies. Nat Genet. 38:904–909.

Prom LK, Ahn E, Isakeit T, Magill C. 2019. GWAS analysis of sorghum association panel lines identifies SNPs associated with disease response to Texas isolates of *Colletotrichum sublineola*. Theor Appl Genet. 132:1389–1396.

Prom LK, Cuevas HE, Ahn E, Isakeit T, Rooney WL, *et al.* 2020. Genome-wide association study of grain mold resistance in sorghum association panel as affected by inoculation with *Alternaria alternata* alone and *Alternaria alternata*, *Fusarium thapsinum*, and *Curvularia lunata* combined. Eur J Plant Pathol. 157:783–798.

Prom LK, Isakeit T, Cuevas H, Rooney WL, Perumal R, *et al.* 2015. Reaction of sorghum lines to zonate leaf spot and rough leaf spot. Plant Health Prog. 16:230–234.

Queiroz VAV, da Silva CS, de Menezes CB, Schaffert RE, Guimarães FFM, *et al.* 2015. Nutritional composition of sorghum [sorghum bicolor (l.) moench] genotypes cultivated without and with water stress. J Cereal Sci. 65:103–111.

Quinby J. 1975. The genetics of sorghum improvement. J Heredity. 66:56–62.

Rhodes D, Gadgil P, Perumal R, Tesso T, Herald TJ. 2017a. Natural variation and genome-wide association study of antioxidants in a diverse sorghum collection. Cereal Chem. 94:190–198.

Rhodes DH, Hoffmann L Jr, Rooney WL, Ramu P, Morris GP, *et al.* 2014. Genome-wide association study of grain polyphenol concentrations in global sorghum [Sorghum bicolor (l.) moench] germplasm. J Agric Food Chem. 62:10916–10927.

Rhodes DH, Hoffmann L, Rooney WL, Herald TJ, Bean S, *et al.* 2017b. Genetic architecture of kernel composition in global sorghum germplasm. BMC Genomics. 18:15.

Rice BR, Fernandes SB, Lipka AE. 2020. Multi-trait genome-wide association studies reveal loci associated with maize inflorescence and leaf architecture. Plant Cell Physiol. 61:1427–1437.

Robinson GK. 1991. That BLUP is a good thing: the estimation of random effects. Stat Sci. 6:15–32.

Sapkota S, Boyles R, Cooper E, Brenton Z, Myers M, *et al.* 2020. Impact of sorghum racial structure and diversity on genomic prediction of grain yield components. Crop Sci. 60:132–148.

Shakoor N, Ziegler G, Dilkes BP, Brenton Z, Boyles R, *et al.* 2016. Integration of experiments across diverse environments identifies the genetic determinants of variation in sorghum bicolor seed element composition. Plant Physiol. 170:1989–1998.

Shao M-R, Shedge V, Kundariya H, Lehle FR, Mackenzie SA. 2016. Ws-2 introgression in a proportion of *Arabidopsis thaliana* col-0 stock seed produces specific phenotypes and highlights the importance of routine genetic verification. Plant Cell. 28:603–605.

Stephens J, Miller F, Rosenow D. 1967. Conversion of alien sorghums to early combine genotypes 1. Crop Sci. 7:396–396.

Stephens M, Carbonetto P, Dai C, Gerard D, Lu M, *et al.* 2020. ashr: Methods for adaptive shrinkage, using Empirical Bayes. https://github.com/stephens999/ashr.

Stephens M. 2017. False discovery rates: a new deal. Biostatistics. 18: 275–294.

Studer AJ, Doebley JF. 2011. Do large effect QTL fractionate? A case study at the maize domestication QTL teosinte branched1. Genetics. 188:673–681.

Sukumaran S, Xiang W, Bean SR, Pedersen JF, Kresovich S, *et al.* 2012. Association mapping for grain quality in a diverse sorghum collection. Plant Genome. 5:126–135.

Swarts K, Gutaker RM, Benz B, Blake M, Bukowski R, *et al.* 2017. Genomic estimation of complex traits reveals ancient maize adaptation to temperate North America. Science. 357:512–515.

Thurber CS, Ma JM, Higgins RH, Brown PJ. 2013. Retrospective genomic analysis of sorghum adaptation to temperate-zone grain production. Genome Biol. 14:R68.

Tian F, Bradbury PJ, Brown PJ, Hung H, Sun Q, *et al.* 2011. Genome-wide association study of leaf architecture in the maize nested association mapping population. Nat Genet. 43:159–162.

Truong SK, McCormick RF, Rooney WL, Mullet JE. 2015. Harnessing genetic variation in leaf angle to increase productivity of sorghum bicolor. Genetics. 201:1229–1238.

Urbut SM, Wang G, Carbonetto P, Stephens M. 2019. Flexible statistical methods for estimating and testing effects in genomic studies with multiple conditions. Nat Genet. 51:187–195.

USDA, N. G. R. P., ARS 2010. Germplasm resources information network (GRIN). https://www.ars-grin.gov

Vandenbrink JP, Delgado MP, Frederick JR, Feltus FA. 2010. A sorghum diversity panel biofuel feedstock screen for genotypes with high hydrolysis yield potential. Ind Crops Prod. 31:444–448.

VanRaden PM. 2008. Efficient methods to compute genomic predictions. J Dairy Sci. 91:4414–4423.

Vinall HN, Martin JH, Stephens JC. 1936. Identification, History, and Distribution of Common Sorghum Varieties, Vol. 501. US Dept. of Agriculture.

Visscher PM, Yang J. 2016. A plethora of pleiotropy across complex traits. Nat Genet. 48:707–708.

Wallace J, Larsson S, Buckler E. 2014. Entering the second century of maize quantitative genetics. Heredity (Edinb). 112:30–38.

Wang J, Hu Z, Upadhyaya HD, Morris GP. 2020. Genomic signatures of seed mass adaptation to global precipitation gradients in sorghum. Heredity (Edinb). 124:108–121.

Wang Y-H, Upadhyaya HD, Burrell AM, Sahraeian SME, Klein RR, *et al.* 2013. Genetic structure and linkage disequilibrium in a diverse, representative collection of the c4 model plant, sorghum bicolor. G3 (Bethesda). 3:783–793.

Wu Y, Li X, Xiang W, Zhu C, Lin Z, *et al.* 2012. Presence of tannins in sorghum grains is conditioned by different natural alleles of tannin1. Proc Natl Acad Sci USA. 109:10281–10286.

Yamaguchi M, Fujimoto H, Hirano K, Araki-Nakamura S, Ohmae-Shinohara K, *et al.* 2016. Sorghum dw1, an agronomically important gene for lodging resistance, encodes a novel protein involved in cell proliferation. Sci Rep. 6:28366.

Yin L, Zhang H, Tang Z, Xu J, Yin D, *et al.* 2021. rMVP: A memory-efficient, visualization-enhanced, and parallel-accelerated tool for genome-wide association study. Genom. Proteom. Bioinform.

Yin L. 2020. Cmplot: Circle manhattan plot. https://github.com/YinLiLin/CMplot.

Zhang D, Kong W, Robertson J, Goff VH, Epps E, *et al.* 2015a. Genetic analysis of inflorescence and plant height components in sorghum (Panicoidae) and comparative genetics with rice (Oryzoidae). BMC Plant Biol. 15:107.

Zhang D, Li J, Compton RO, Robertson J, Goff VH, *et al.* 2015b. Comparative genetics of seed size traits in divergent cereal lineages represented by sorghum (panicoidae) and rice (oryzoidae). G3 (Bethesda). 5:1117–1128.

Zhao J, Mantilla Perez MB, Hu J, Salas Fernandez MG. 2016. Genome-wide association study for nine plant architecture traits in sorghum. Plant Genome. 9: 1–14.

Zheng Z, Hey S, Jubery T, Liu H, Yang Y, *et al.* 2020. Shared genetic control of root system architecture between zea mays and sorghum bicolor. Plant Physiol. 182:977–991.

Zhou X, Stephens M. 2014. Efficient multivariate linear mixed model algorithms for genome-wide association studies. Nat Methods. 11:407–409.

Zhou X. 2017. A unified framework for variance component estimation with summary statistics in genome-wide association studies. Ann Appl Stat. 11:2027–2051.

Zhou Y, Srinivasan S, Mirnezami SV, Kusmec A, Fu Q, *et al.* 2019. Semiautomated feature extraction from RGB images for sorghum panicle architecture GWAS. Plant Physiol. 179:24–37.