

1 **Meta-GWAS for quantitative trait loci identification in soybean**

2 **Authors:** Johnathon M. Shook<sup>1</sup>, Jiaoping Zhang<sup>1</sup>, Sarah E. Jones<sup>1</sup>, Arti Singh<sup>1</sup>, Brian W. Diers<sup>2</sup>, Asheesh  
3 K. Singh<sup>1,\*</sup>

4 **Affiliation:** <sup>1</sup> Department of Agronomy, Iowa State University, Ames, IA 50011; <sup>2</sup> Department of Crop  
5 Sciences, University of Illinois, Urbana, IL 61801;

6 **\* Corresponding Author:** A.K. Singh ([singhak@iastate.edu](mailto:singhak@iastate.edu))

7 **Keywords:** Meta-analysis, GWAS, QTL, agronomic traits, seed composition traits, disease resistance

8 **ABSTRACT**

9 We report a meta-Genome Wide Association Study involving 73 published studies in soybean  
10 (*Glycine max* L. [Merr.]) covering 17,556 unique accessions, with improved statistical power for robust  
11 detection of loci associated with a broad range of traits. *De novo* GWAS and meta-analysis were conducted  
12 for composition traits including fatty acid and amino acid composition traits, disease resistance traits, and  
13 agronomic traits including seed yield, plant height, stem lodging, seed weight, seed mottling, seed quality,  
14 flowering timing, and pod shattering. To examine differences in detectability and test statistical power  
15 between single- and multi-environment GWAS, comparison of meta-GWAS results to those from the  
16 constituent experiments were performed. Using meta-GWAS analysis and the analysis of individual studies,  
17 we report 483 quantitative trait loci (QTL) at 393 unique loci. Using stringent criteria to detect significant  
18 marker trait associations, 66 candidate genes were identified, including 17 candidate genes for agronomic  
19 traits, 19 for seed related traits, and 33 for disease reaction traits. This study identified potentially valuable  
20 candidate genes that affect multiple traits. The success in narrowing down the genomic region for some loci  
21 through overlapping mapping results of multiple studies is a promising avenue for community-based studies  
22 and plant breeding applications.

23 **INTRODUCTION**

24 Genome-wide association studies (GWAS) analyze the association between a trait of interest and  
25 thousands of genetic variants throughout the genome. The general approach has benefited from the  
26 development of greatly increased numbers of markers due to the advent of next-generation sequencing  
27 approaches (Rico *et al.* 2013), and increased sample size with the formation of biobanks, such as the  
28 100,000 Genomes Project (The 100,000 Genomes Project 2019). Plant scientists now routinely conduct  
29 GWAS in crop species, including soybean [*Glycine max* (L.) Merr.]. Increased marker data availability and  
30 development of new statistic methods provided great opportunities to gain new knowledge from existing  
31 data and address previous lacuna of GWAS experiments (Zeng *et al.* 2017; Chang *et al.* 2016; Chang and

32 Hartman 2017; Bandillo *et al.* 2015; Bandillo *et al.* 2017; Zhou *et al.* 2015; de Azevedo Peixoto *et al.* 2017;  
33 Zhang *et al.* 2015; Zhang *et al.* 2017).

34 Researchers have recognized that while single environment GWAS such as those conducted in the  
35 greenhouse are powerful for genetic studies and candidate gene identification, their extrapolation in field  
36 environment applications require further validation (Zhang *et al.* 2015; de Azevedo Peixoto *et al.* 2017;  
37 Coser *et al.* 2017). When comparing separate studies of the same trait, significant differences in results are  
38 often found. These differences may be caused by allele frequency variation between populations,  
39 inadequate control of population structure, or environmental dependencies (Gibson and Mullen 1996). With  
40 the availability of standardized marker data across the USDA soybean germplasm collection (Song *et al.*  
41 2015), several studies have mapped important major effect quantitative trait loci (QTL) using historical  
42 records and GWAS analysis: for example, insect resistance (Chang and Hartman 2017), disease resistance  
43 (Chang *et al.* 2016), descriptive traits such as flower and pubescence color (Bandillo *et al.* 2017), and seed  
44 oil and protein content (Bandillo *et al.* 2015). However, for many quantitative traits such as seed  
45 composition or plant height, using raw measurements from differing environments introduces bias, which  
46 may erode the power of detection for significant QTL (Chen *et al.* 2010). While results from within the  
47 same environment(s) share a common environmental component, attempting to combine multiple panels  
48 grown in different environments leads to an improper assignment of environmental effects to the differences  
49 between genetics of the panels involved (Zhao *et al.* 2019). Meta-analysis provides an attractive alternative  
50 to address the above-mentioned challenges of individual GWAS, and this analysis can be performed on  
51 results from independent studies using statistical approaches such as those provided by the analysis program  
52 METAL (Willer *et al.* 2010).

53 Quantitative traits, in contrast with qualitative traits, are controlled by many genes and  
54 environmental factors. To fully understand the pathways that determine these traits, interactions between  
55 previously discovered genes and new candidate genes must be added to the existing models. Directly  
56 measured traits often comprise only a portion of the information about a biological pathway, necessitating  
57 the identification of pleiotropic effects (on correlated traits) for an increased biological understanding of  
58 the phenotype. Genes may exhibit pleiotropy either through control of a common pathway such as the  
59 influence of *Dt1* on both plant height and lodging (Diers *et al.* 2018), or through multiple effects of a  
60 chemical as seen in the effect of *T* locus that has a dual role in pigmentation and chilling tolerance through  
61 isoflavones (Takahashi and Asanuma 1996). Identifying genes that control multiple phenotypes of  
62 importance can either suggest candidates for fixation, in cases where both effects are positive, or may  
63 identify possible penalties associated with incorporating particular alleles and improve multi-trait selection  
64 results (Bolormaa *et al.* 2014).

65           Meta-analyses include separately analyzing each individual experiment in order to determine  
66 experiment-specific p-value and allele effect estimates, rather than performing a combined analysis to  
67 leverage extensive data (Bandillo *et al.* 2015). Further genetic insights can be gleaned through an ease in  
68 the identification of pleiotropic effects due to the analysis of a wide range of traits. Moreover, the ability to  
69 compare the results from a combined analysis with those from separate analyses of individual studies allows  
70 for the identification of both environment-dependent associations and for the enrichment and detection of  
71 quantitative traits and rare alleles from more unique but diverse populations. Previous meta-analysis results  
72 have shown the effectiveness of combined panels to identify minor genes that were missed in a single study  
73 (Chang *et al.* 2017). Due to the need for adequate representation of minor alleles in GWAS, rare alleles that  
74 are predominant in a small zone of adaptation may be absent or undetectable within individual studies. The  
75 agronomic screenings for the USDA soybean germplasm collection are arranged based on the influx of new  
76 germplasm into the United States, and therefore serve as a semi-randomized subset of global soybean  
77 variation and spatiotemporal patterns in the origins of new accessions enabling potential detection of rare  
78 variants, which may be enriched in one of these geographical regions (Trotta *et al.* 2016).

79           While combined analyses for disease and insect resistance (Chang *et al.* 2016; Chang and Hartman  
80 2017) and seed composition (Bandillo *et al.* 2015) have previously been reported, we perform a large-scale  
81 meta-analysis utilizing individual studies in soybean. Our study builds on previous studies by integrating  
82 the environmental component that can provide a historical perspective on adaptation, with the inclusion of  
83 quantitative traits of agronomic importance, stress tolerance, and seed composition. Subsequent study of  
84 pleiotropic genes and reporting on gene rich clusters can be useful when attempting to introgress favorable  
85 alleles into breeding lines (Cameron *et al.* 2017), as it improves the understanding of potential  
86 complications of introgression. The multitude of traits examined with our study facilitates the detection of  
87 co-localized peaks indicative of potential pleiotropic effects of genes across a diverse range of phenotypes.  
88 Loci associated with multiple traits identified within this study require additional functional validation, as  
89 GWAS are not designed to definitively differentiate between pleiotropy and (tight) linkage. We included  
90 results from reports published from 1964 to 2009 for a total of 73 individual studies. The design of this  
91 study was intended to identify co-localization of peaks for multiple traits, as well as to identify previously  
92 overlooked genes through meta-analysis approaches. Using meta-GWAS analysis and analysis of  
93 individual studies, we report 393 unique QTL including 66 candidate genes across important traits and  
94 provide confirmation of many previously reported genes. This study provides targets for functional  
95 characterization and introgression of previously untapped diversity for many important traits.

## 96 **MATERIALS AND METHODS**

### 97 **Genotypic data and quality control**

98 Marker data from the testing of 20,087 *Glycine max* and *G. soja* accessions from the USDA  
99 Soybean Germplasm Collection with the SoySNP50K iSelect BeadChip (Song *et al.* 2013) were  
100 downloaded from SoyBase (Song *et al.* 2013). A data imputation pipeline based on Java implementation of  
101 Beagle 5.0 (Browning and Browning 2016) was utilized to impute missing data for the 42,080 SNP markers  
102 that were aligned to the Williams 82 reference genome v2 assembly. Markers aligned to scaffolds but not  
103 assigned to a chromosome were removed prior to processing. Ten burn-in iterations and five phasing  
104 iterations were used to impute missing markers, which accounted for 0.64% of all markers. For each test,  
105 markers remaining after applying cutoffs of minor allele frequency  $\geq 0.05$  for studies involving  $300 \leq n$   
106  $\leq 1000$  accessions, or 0.01 for studies involving  $n \geq 1001$  accessions, were selected for further analysis.

### 107 **Phenotypic data and genetic accessions**

108 Numeric phenotypic data from USDA reports were compiled from the U.S. National Plant  
109 Germplasm System website (<http://npgsweb.ars-grin.gov/gringlobal/descriptors.aspx>) (Descriptors for  
110 Soybean 2019). Subsets of accessions that were part of historical USDA germplasm characterization trials  
111 with a size  $n \geq 300$  were selected for further analysis. Information on the design of the original trials is  
112 available from the technical bulletins in which they were originally published. These technical bulletins are  
113 available online in part at <https://pubs.nal.usda.gov/sites/pubs.nal.usda.gov/files/tb.htm> (Miller 2003).  
114 Alternatively, PDFs of the technical bulletins are available on our GitHub  
115 (<https://github.com/SoylabSingh/META-GWAS>). Additional traits, such as disease resistance and amino  
116 acid composition, were downloaded from the NPGS website.

### 117 **Genome-wide association analysis**

118 Each experiment was analyzed separately with a mixed linear model implemented using GAPIT in  
119 R (Lipka *et al.* 2012) to prevent confounding of environmental effects with marker effects, which would be  
120 expected for several traits (i.e., flowering time, oil, protein, etc.). Population structure was controlled using  
121 the first three PCAs based on the marker data. This resulted in 585 combinations of experiment/trait  
122 analyses. Analysis was subsequently performed for combined panels for each trait. The Bonferroni  
123 threshold (Neyman and Pearson 1928) was employed to minimize the likelihood of false positives in  
124 declaring significance. The significant SNPs were compiled for further analysis (**Supplemental Table 1**).

125 Initial QTL calling was performed trait-by-trait based on marker position. Subsequently, QTL for  
126 related traits (such as flowering date and maturity date) with substantial overlap were merged, resulting in  
127 fewer unique QTL than originally called. Local LD decay analysis was used to further clarify between  
128 separate or overlapping QTL.

129 Markers that were significant for multiple traits and experiments, or were identified during analysis  
130 of the combined trials, were examined for nearby candidate genes. Candidate genes were identified by  
131 examining annotated genes within linkage disequilibrium (LD) of the leading SNP with  $r^2 > 0.7$  for each  
132 experiment and peak (de Azevedo Peixoto *et al.* 2017). Candidate gene identification was performed based  
133 on previously characterized genes, gene family function, and the nearest gene to the peak SNP in cases  
134 where no known function could be identified. For candidate casual genetic variant analysis, we utilized the  
135 SNP dataset from the genome resequencing study of 302 soybean lines (Zhou *et al.* 2015) and searched the  
136 possible causal mutants at the identified candidate genes. We first identified the lead SNP from peaks of  
137 interest in the resequencing dataset, then calculated the pairwise LD  $r^2$  values between the lead SNP and  
138 the SNPs covering the locus of candidate gene. All other analyses here within were aligned to the Glyma2.0  
139 reference genome (<https://soybase.org/gb2/gbrowse/gmax2.0/>). The R package ‘circlize’ was employed to  
140 generate the circular visualizations of significant SNPs for multiple traits throughout the genome. Study  
141 names have been shortened for convenience within the text; a reference file is provided to find the initial  
142 source of phenotypic data used in this work (**Supplemental Table 2**). Trait definitions, as well as the  
143 number of QTL and candidate genes identified for each trait, are provided in **Supplemental Table 3**.

## 144 RESULTS AND DISCUSSION

145 From the individual study GWAS and meta-GWAS 4,919 significant SNPs were detected, of which  
146 787 were reported from the meta-GWAS analysis. Complete listing of the significant SNP identified using  
147 individual study GWAS and meta-GWAS are provided in **Supplemental Table 1**. Among these 787 SNPs  
148 identified using meta-GWAS, 110 were associated with agronomic traits, 106 with seed composition traits,  
149 and 571 with disease resistance traits. Overall, candidate genes were assigned for 66 unique loci; and these  
150 included genes with moderate to large effects. We focus our results on loci that were associated with  
151 multiple traits.

### 152 Agronomic Traits

153 Amongst agronomic traits, we identified 1422 marker-trait associations with traditional GWAS  
154 studies, as well as 110 SNPs associated with agronomic traits when analyzed across studies by meta-  
155 GWAS. In all, 115 QTL across 20 chromosomes were identified, with 17 candidate genes (**Figure 1a**,  
156 **Supplemental Table 1, Supplemental Table 3, Table 1**).

157 In our approach, we used results from individual studies to detect overlapping genomic regions for  
158 the purpose of locating candidate gene for traits, including for genes previously cloned. The locus harboring  
159 *Dt1* (*Glyma.19g194300*) (Liu *et al.* 2010), the major gene conditioning stem termination in soybean, was  
160 significantly associated with oleic acid and linoleic acid content, as well as plant height, stem termination,  
161 and stem lodging (**Supplemental Table 1**). By comparing the mapping results of four studies, we were

162 able to limit the candidate genomic region to a 125 kb fragment harboring previously cloned *Dt1* (from  
163 *ss715635422* to *ss715635460*) (**Supplemental Figure 1**). These results highlight the advantages of meta-  
164 GWAS for finer mapping the candidate gene region. A nonsynonymous SNP (*SNP\_19\_44980087*), in high  
165 LD ( $r^2 = 0.5$ ) with the leading SNP *ss715635424* (also known as *SNP\_19\_45000827*), was found at the  
166 fourth exon of *Dt1* that changes amino acid R (Arg) to W (Trp) (**Supplemental Figure 2**). This SNP is  
167 identical to the R166W mutation previously identified (Liu *et al.* 2010).

168 On chromosome 19, we identified a QTL for stem lodging which was on the opposite end of the  
169 chromosome as *Dt1*. Stem lodging is associated with plant height and this has been reported in multiple  
170 crops (Flint-Garcia *et al.* 2003; Diers *et al.* 2018; Singh *et al.* 2019). As lodging causes significant yield  
171 and quality losses, the development of the shorter statured wheat and rice were promoted which could better  
172 handle high input agriculture. However, this solution is not universally applicable. In soybean, pods are  
173 arranged at nodes on the stem and decreasing the length of stem, and if fewer nodes are present, yield  
174 potential is reduced. Leveraging four studies, we report a peak for tolerance to stem lodging with the  
175 candidate gene *Glyma.19g016400*, an ABC transporter on chromosome 19. This locus was found to affect  
176 lodging tolerance but was not found to be associated with plant height, thereby making it a useful target to  
177 develop lodging resistant soybean cultivars without decreasing stem length and yield potential. While this  
178 is the first genome wide association study identifying this gene, additional evidence towards its validity  
179 comes from several recent patents (US Patents #8697941, 8748695, and 9675071) that relate to molecular  
180 markers in the region of interest and include *Glyma.19g016400* as one of the candidate genes for PPO  
181 inhibitor tolerance in soybean. Significant effects of this region for seed yield, lodging, and plant height  
182 were reported from the SoyNAM project (Diers *et al.* 2018). The results from Hulting *et al.* (2001) on PPO  
183 inhibitor tolerance and our findings on stem lodging susceptibility suggests a tradeoff between PPO  
184 inhibitor tolerance and lodging susceptibility. The soybean accessions highly tolerant to sulfentrazone  
185 contain alleles associated with increased lodging in our study, necessitating further studies to validate these  
186 observations.

187 On chromosome 6, a significant SNP peak was identified that co-located with the *T* gene, a  
188 flavonoid 3' hydroxylase (Toda *et al.* 2002). This region was significant for arginine, cysteine, isoleucine,  
189 and leucine levels, as well as for seed mottling (**Figure 1 a, c**). The cloned *E2* locus (Watanabe *et al.* 2011)  
190 was significantly associated with flowering and maturity date, maturity group, days from flowering to  
191 maturity, plant height, and seed yield (**Figure 1b**). The associations between *E2* and these traits has been  
192 previously reported (Fang *et al.* 2017).

### 193 **Seed Composition Traits**

194 Amongst seed composition traits, we identified 1364 marker-trait associations with traditional  
195 GWAS studies, as well as 106 SNPs associated with compositional traits when analyzed across studies by  
196 meta-GWAS. SNPs associated with composition were found on chromosomes 1-9, 11, 13-15, 17, 19-20,  
197 resulting in 88 QTL with 19 candidate genes (**Figure 1b, Supplemental Table 1, Supplemental Table 3,**  
198 **Table 2**)

199 A cluster of candidate genes for seed composition, including isoleucine, methionine, leucine,  
200 tryptophan, threonine, lysine, and palmitic acid, were located in a region of 30 kb on chromosome 1 between  
201 53.13 – 53.16 Mb, 4 a cysteine desulfurase (*Glyma.01g197100*) and a malate and lactate dehydrogenase  
202 gene (*Glyma.01g197700*) (**Supplemental Figure 1**). Further targeted analysis will be necessary to  
203 determine which gene is influencing each trait, as a single enzyme is unlikely responsible for multiple steps  
204 in the metabolic pathway. We found significant SNPs in high LD ( $r^2 > 0.5$ ) with the detected leading SNP  
205 at the promoter of *Glyma.01g197700*, but not in the coding region of the gene (**Supplemental Figure 2**).

206 A region including the *I* locus on chromosome 8 (Clough *et al.* 2004) was associated with seed  
207 mottling, as well as oil, cysteine, isoleucine, leucine, linoleic acid, lysine, methionine, palmitic acid, stearic  
208 acid, threonine, and valine levels in the seed (**Figure 1b**). The most likely candidate gene for the observed  
209 differences in amino acids levels, *AK-HDSH* (aspartokinase homoserine dehydrogenase,  
210 *Glyma.08g107800*) is a bifunctional enzyme catalyzing the key steps of asparagine phosphatization and the  
211 aspartate-semialdehyde to homoserine conversion by which aspartate family amino acids (lysine, threonine,  
212 methionine, and isoleucine) are synthesized (Zhu-Shimoni and Galili 1998). However, amino acid data  
213 were generated using Near Infrared Reflectance, which may have low precision in estimating amino acid  
214 composition when there is variability in seed coat color (Baianu *et al.* 2011). Therefore, further validation  
215 is needed to establish the association between the *AK-HDSH* or *I* loci and the amino acid profile.

216 *SACPD-C* (*Glyma.14g121400*) was the primary candidate to explain differences in stearic acid  
217 content within seed oil and has been previously functionally validated (Gillman *et al.* 2014). Using the  
218 Wm82.a2 reference genome build, this appeared as three separate peaks; however, a single peak was  
219 observed when using the Wm82.a1 version. We postulate a possible assembly error in the region  
220 surrounding the *SACPD-C* locus in the soybean reference genome Wm82.a2, due to conflicting results  
221 (**Supplemental Table 4**). We attempted to identify false peaks generated due to genome mis-assembly by  
222 fitting the lead SNP as a covariate in the GWAS model, and then observed lower p-values for the remaining  
223 SNPs and detected a weaker signal from surrounding SNPs indicative of a single gene. Presence of stronger  
224 signals in surrounding SNPs would have indicated that two separate genes are in play. Additionally, the  $r^2$   
225 between SNPs in all three regions was greater than 0.7, suggesting physical linkage. The Wm82.a1 results

226 (SNP effects, physical location, LD) provide the most plausible explanation for the presence of a single  
227 gene in this genomic region and suggests that Wm82.a2 has unresolved errors in scaffold positioning.

228 A peak on chromosome 5 associated with palmitic acid content was detected in 3 different studies.  
229 Using data from the ‘2mn81’ study, the locus mapped to a region of over 600 kb. However, two other  
230 studies (2ky81 and ms2000.02) mapped this locus within a smaller region of 130 kb (*ss715592495-*  
231 *ss715592503*) and 182 kb (*ss715592491-ss715592500*), respectively, with an overlap of about 88 kb  
232 (*ss715592495-ss715592500*) (**Supplemental Figure 2**). The candidate gene *FATB1a* (*Glyma.05g012300*)  
233 (Wilson *et al.* 2001) was identified in the overlap. However, no SNP in LD ( $r^2 \geq 0.5$ ) with the leading SNP  
234 of the locus was identified at the coding region or promoter of *FATB1a* based on analysis of resequencing  
235 data (Zhou *et al.* 2015) except the synonymous *SNP\_5\_7995427* (**Supplemental Figure 1**). Causal variants  
236 have been identified in mutagenized breeding material (Thapa *et al.* 2016, Bachleda *et al.* 2016, Goettel *et*  
237 *al.* 2016), but naturally occurring variations are not well characterized.

### 238 **Disease Resistance Traits**

239 Amongst disease traits, we identified 1346 marker-trait associations with traditional GWAS  
240 studies, as well as 571 SNPs associated with disease traits when analyzed across studies by meta-GWAS.  
241 213 QTL mapped to all 20 chromosomes, with 33 candidate genes identified (**Figure 1c, Supplemental**  
242 **Table 1, Supplemental Table 3, Table 3**). Meta-analysis in several instances narrowed the genomic region  
243 for QTL. For example, the association between the *Rps3* region and resistance to race 1 of *Phytophthora*  
244 root rot was mapped to a 144 kb region in the meta-analysis, compared to a 1Mb region in individual studies  
245 (**Supplemental Table 1**). This reduces the search space for causal genes and allows for greater accuracy  
246 when identifying candidate genes.

247 We found a peak that was associated with resistance to races 1, 2, 3, 4, 5, 7, 10, and 17 of  
248 *Phytophthora sojae* that mapped to the position of the *Rps1* locus (Gao and Bhattacharyya 2008). A  
249 previously unreported peak for soybean cyst nematode resistance was identified on chromosome 11 was  
250 mapped to *Glyma.11g234500*, an alpha-soluble N-ethylmaleimide-sensitive factor (NSF) attachment  
251 protein ( $\alpha$ -SNAP). Notably, the candidate genes *GmSNAP11* (*Glyma.11g234500*) and *GmSNAP14*  
252 (*Glyma.14g054900*) (Lakhssassi *et al.* 2017), identified at 7 kb and 84 kb apart from lead SNPs  
253 *ss715610420* and *ss715618859*, respectively, are paralogs and encode a Soluble NSF Attachment Protein  
254 (SNAP). Another soybean SNAP gene on chromosome 18, *GmSNAP18*, has been reported to play a role  
255 in resistance to SCN (Cook *et al.* 2012). On chromosome 1, the locus for seed composition co-localized  
256 with a bacterial pustule resistance QTL. This QTL does not correspond to the previously identified *Rxp*  
257 locus, instead, a candidate gene *Glyma.01g197800* is identified as the potential underlying gene. A peak on  
258 chromosome 3 at 34.24 - 35.18 Mb was found to be significantly associated with iron deficiency chlorosis



259 tolerance and *Pythium irregulare* resistance. This region has previously been investigated as the source of  
260 IDC tolerance in “Isoclark” (Stec et al. 2013). The GWAS analysis identified previously unreported  
261 genomic regions that were associated with resistance to bean pod mottle virus, brown stem rot, frogeye leaf  
262 spot, *Phytophthora* root rot, and soybean cyst nematode (**Figure 1c**). A full list of identified SNPs and  
263 candidate genes for these traits, as well as for all other traits examined in this study using both combined  
264 analyses and analysis of individual experiments are provided in **Supplemental Table 1**.

265 The majority of studies included in this work for disease resistance were germplasm screenings,  
266 where many entries were tested to find new sources of resistance. Such germplasm screening studies were  
267 not originally intended for GWAS; for example, multiple rating systems, ordinal rating scales, and non-  
268 integer ratings used in the studies complicates result comparisons and are not easily amenable to linear  
269 statistical models. Standardization of screening protocols across research groups and inclusion of key data  
270 for comparison of studies such as those suggested by the MIAPPE checklist (Ćwiek-Kupczyńska et al.  
271 2016) will be key for future research into plant disease resistance. In addition, an increased utilization of  
272 image-based phenotyping will play a key role, allowing for digital disease severity ratings on a continuous  
273 scale (Naik et al. 2017; Zhang et al. 2017), minimal inter- and intra-rater variability in measurements  
274 through hyperspectral camera and ML-based analysis (Nagasubramanian et al. 2018; Nagasubramanian et  
275 al. 2019). It will also enable the comparison of results across studies by facilitating reanalysis of previous  
276 experiments with new rating systems or approaches, as long as needed input variables are available.

### 277 **Implications of pleiotropy vs. linked genes**

278 While repeated crossing or careful selection of the donor parent can break linkage drag, negative  
279 pleiotropic effects associated with a gene of interest are more problematic. Candidate gene analysis was  
280 aided by tissue-specific gene expression data available at SoyBase. The use of a blend of individual and  
281 meta-analyses provided improved resolution through examining overlapping peaks and utilizing the  
282 increased power in larger panels in the meta-analysis. When investigating the peak on chromosome 1 for  
283 fatty acid and amino acid composition, a convincing distinction between pleiotropy and linkage could not  
284 be made. This was due to the presence of multiple strong candidate genes. While meta-GWAS approaches  
285 are very beneficial for improving map resolution, they are still limited in their inference in regions with  
286 strong linkage disequilibrium. Meta-GWAS results outputs still require follow-up molecular and functional  
287 validation to confirm the candidate genes as well as to confirm pleiotropy vs. linkage.

288 Pleiotropic effects of major genes significantly alter multiple traits simultaneously, creating a  
289 situation of either rapid improvement across traits, or of tradeoffs, such as is found in most soybean  
290 protein/oil content QTL. Genetic improvement utilizing pleiotropic effects may be limited in applicability  
291 to specific geographic regions if they affect key adaptation genes such as the maturity loci or stem

292 termination. Therefore, it will be necessary for breeders to independently determine whether a gene with  
293 pleiotropic effects is a good fit for their variety development goals. In cases where pleiotropy is associated  
294 with a tradeoff between multiple traits, such as between seed protein and oil content, breeders will need to  
295 weigh the importance of each trait or identify combinations of genes affecting the trait that can provide an  
296 adequate phenotype for each trait considered.

### 297 **Motivation for the use of meta-analysis**

298 For many important row crop species, such as soybean, corn, wheat, and sorghum, it is impractical  
299 or impossible to evaluate the full breadth of the available germplasm at a single location. This is due to  
300 space limitations, availability of labor or funding for phenotyping, or irreconcilable differences between  
301 genotypes preventing them from growing in the same place, such as differences in photoperiod sensitivity  
302 or vernalization requirements. To capture the breadth of the genetic and phenotypic diversity, it is necessary  
303 to test each variety with a similarly adapted cohort. The separate analysis of each environment can increase  
304 the odds of finding alleles which are near fixation in the population or are environmentally dependent (Singh  
305 et al. 2014; Sherman et al. 2019).

306 For simple, qualitative traits such as pubescence color in soybean, there is little benefit in meta-  
307 GWAS due to the consistency with which the gene can be mapped and the lack of environmental  
308 dependence on trait expression. When studying environmentally dependent traits, such as agronomic,  
309 disease resistance and seed composition traits including seed oil or protein content, meta-GWAS provide  
310 advantages particularly in increasing the likelihood of finding small effect genes. When comparing  
311 individual experiments results (**Figure 2a**) with the combined meta-analysis (**Figure 2b**), additional  
312 significant peaks were observed in meta-analysis. For example, the SNP marker *ss715614263* was  
313 previously associated with seed protein using mega-analysis (Bandillo et al. 2015). The same locus was  
314 found to be associated with protein, palmitic, and oleic acid content in an individual panel in the current  
315 study (ms2000.02), but was associated with protein and linoleic acid content in the meta-analysis  
316 (**Supplemental Table 1**). While meta-analysis identified fewer traits in the specific instance of  
317 *ss715614263*, the association with an additional trait (compared to individual analysis) still encourages its  
318 use, as each newly associated trait may provide guidance in identifying putative causal genes. A full listing  
319 of candidate genes detected in each study is provided as **Supplemental Table 5**, which also provides a  
320 reference to candidate genes detected either only in individual studies or only via meta-analysis.  
321 Identification of an association with multiple related traits, although only spanning one to two markers, is  
322 a strong signal that the association may merit additional study to identify a strong candidate gene and further  
323 explore the possible pleiotropic effects this locus is exhibiting, especially when a stringent cut-offs are used  
324 to declare significance.

325 To maximize the effectiveness of soybean breeding programs, we sought to identify as many genes  
326 as possible for numerous traits, ensuring that multiple paths are available for further cultivar improvement.  
327 By maximizing the identified links between markers and phenotypes of interest, meta-GWAS aids efforts  
328 to bridge the gap between genotype and phenotype, allowing for improvements not only in trait prediction  
329 and selections, but also in modelling the interactions between multiple genes in overall trait performance.

### 330 **Future mapping, validation and integration with Phenomics studies**

331 Traditional fine mapping through creating lines sharing homogenous genetic background, such as  
332 near isogenic lines, is a powerful tool to uncover the casual genetic variants. However, it is time consuming  
333 to develop new near-isogenic lines in multiple backgrounds to reduce the potential influence of background-  
334 specific effects. In this study, large variation of LD architecture was observed across populations. This  
335 enables substantially shortening of the candidate chromosomal regions of specific QTL by comparing  
336 mapping results from separate studies using different populations. Considering almost all accessions in the  
337 USDA Soybean Germplasm Collection were genotyped by SoySNP50K BeadChip and are publicly  
338 accessible, mapping populations with a high LD decay rate at specific genomic regions of interest can be  
339 constructed for fine mapping. The consistent identification of major genes, including those affecting  
340 multiple traits of interest, suggests that further improvements in mapping ability would likely require a  
341 model with the major genes treated as covariates. While it is currently possible to account for the effects of  
342 major genes by using SNPs linked to the gene of interest as covariates, this approach is only an  
343 approximation due to incomplete linkage between common SNPs and the underlying gene. Instead, allele-  
344 specific markers should be developed and deployed across both wild-type germplasm and breeding  
345 material.

346 In the future, similar studies will benefit by incorporating weather, soil, or management parameters  
347 in order to explain differences in marker effects between individual studies and in Meta-GWAS (Cook et  
348 al. 2017 ). In this scenario, access to standardized, quality-controlled records will be needed to tease apart  
349 the GxE component and identify the architecture of environmentally mediated expression and decipher  
350 associations between genetics and environmental signals for the traits of interest. The establishment of  
351 standardized tests enabled with advanced sensors and high-throughput phenotyping should improve the  
352 opportunity to identify additional genes influencing traits of interest through the analysis of previously  
353 ignored component traits, such as leaf expansion rate or chlorophyll density in the case of yield, (Dhondt  
354 et al. 2013) which may lead to an increased understanding of the genetic architecture of these traits and  
355 responses to environmental and management conditions (Parmley et al. 2019).

### 356 **CONCLUSION**

357 Combined analysis of all investigated traits found 63 loci that corresponded to previously reported  
358 QTL, characterized genes, and new reported loci backed up with strong candidate genes conditioning the  
359 observed phenotypes. Several of the previously identified loci (for example, *Dt1*, *E2*) were associated with  
360 multiple traits, identifying putative pleiotropic effects of the underlying genes. Differences between results  
361 in individual trials and the combined analyses confirm the importance of multi-environment testing for  
362 identification of key traits, but also provide a strong motivation to create a community database that can be  
363 queried for scientific advancement. Continued publication of raw phenotypic values from screenings will  
364 increase the power for identification of important genes for both mean and plastic responses to reduce the  
365 financial and time burden on any individual program while benefitting future breeders and researchers. For  
366 example, the sharing of phenotypic information across research programs both nationally and globally, as  
367 currently on-going with multi-states and –institutions uniform soybean tests and other cooperatively run  
368 tests in other crops.

369

370 **Data availability statement:** The authors affirm that all data necessary for confirming the conclusions of  
371 the article are present within the article, figures, and tables. Raw data and codes will be available at  
372 <https://github.com/SoylabSingh/Soy-Meta-GWAS>.

373

374 **Author contribution:** AS, AKS conceptualized the study; JS, AS, AKS designed the study; JS conducted  
375 statistical analysis with contributions from AKS and JZ; Figures were prepared by JZ with inputs from JS;  
376 JS interpreted the results with contributions from JZ, SJ, AS, BD, AKS; JS wrote the first draft with AKS;  
377 All authors contributed in writing, reviewing, and approve the manuscript.

378

379 **Acknowledgements:** Authors sincerely thank all researchers past and present who generated data for  
380 individual studies and set up a community resource for advancing soybean research and development. We  
381 thank David Blystone and Dr. David Grant, which greatly helped the manuscript. We thank the Iowa  
382 Soybean Association (to AKS), R F Baker Center for Plant Breeding (to AKS), Monsanto Chair in Soybean  
383 Breeding (to AKS), USDA IOW04314, and National Research Traineeship (to JMS) for the financial  
384 support.

385

386

387

388 **Table 1.** List of candidate genes identified for agronomic traits using GWAS from individual studies and  
 389 Meta-GWAS.

Chromosome	Likely Gene	Meta-GWAS	Individual Studies GWAS	Trait(s)	Studies Source
5	<i>Glyma.05G200100</i>		*	Flower date, Maturity date, Maturity group	4il87, ms1999.01, ms923
6	<i>E1</i>	*	*	Flower date, Maturity date, Maturity group, Stem termination	1il64, 1il66, 2il81.1, 2il81.2, 4il87, 5il90, il0102, il989, meta, mn945
	<i>Glyma.06G068900</i>	*	*	Seed mottling	3mn83.2, meta
	<i>Glyma.06g134400</i>		*	Pod shattering (early), Pod shattering (late)	4il87
	<i>T</i>	*	*	Seed mottling	3il84, meta, ms1999.01, ms923, ms967
7	<i>Glyma.07g049800</i>	*	*	Pod shattering (early), Pod shattering (late)	3il84, meta, ms1999.01, ms923
8	<i>I</i>		*	Seed mottling	1il66, 2ky81, 4il87, il0102, ms923
9	<i>fr1</i>		*	Root fluorescence	fluorjt97
	<i>Glyma.09g090600</i>	*	*	Seed mottling	1il66, 4il87, meta
	<i>Glyma.09g266200</i>		*	Flower date, Maturity group	ms923, ms1999.01
10	<i>E2</i>	*	*	Branching, Flower date, Height, Maturity date, Maturity group, Yield	1il64, 1il66, 2il81.1, 3il83.1, 3il84, il0102, il989, meta, ms1999.01, ms967
11	<i>K1/AGO</i>	*	*	Seed mottling	3mn83.2, il0102, meta, ms923, ms967
13	<i>Rsv1</i>	*	*	Seed mottling	1il66, 2il81.1, 2il81.2, 5il90, meta, ms1999.01, ms2000.02, ms923
14	<i>fan1</i>		*	Seed quality	2ky81
15	<i>Glyma.15g139800</i>	*	*	Pod shattering (early), Pod shattering (late)	1il66, 2il81.2, 2ky81, meta
16	<i>E9</i>	*	*	Flower date, Maturity group	2il81.1, 3il83.1, meta, ms1999.01
	<i>Pdh1</i>	*	*	Pod shattering (early), Pod shattering (late)	1il64, 2il81.1, 4il87, il0102, meta, ms1999.01, ms2000.02, ms923, ms967

18	<i>Dt2</i>	*	*	Stem termination	meta, mn945, ms923
19	<i>ABC, Glyma.19g016400</i>	*	*	Lodging	1il66, 2ky81, ms923, 3il84, meta
	<i>Dt1, Glyma.19g194300</i>	*	*	Height, Lodging, Stem termination	1il64, 1il66, 2il81.1, 2il81.2, 2ky81, 3il83.1, 3il84, 3mn83.2, 4il87, 5il90, il0102, meta, mn945, ms1999.01, ms2000.02, ms923, ms967
	<i>E3</i>		*	Maturity group	2il81.2

390

391

392 **Table 2.** List of candidate genes identified for seed composition traits using GWAS from individual studies  
393 and Meta-GWAS.

Chromosome	Likely Gene	Meta-GWAS	Individual Studies GWAS	Trait(s)	Studies Source
1	<i>BCAT/MDH</i>	*	*	Isoleucine, Leucine, Lysine, Methionine, Palmitic acid, Threonine, Tryptophan	aa op sugar fa 2009, il0102, meta, ms967
3	<i>Glyma.03g173400</i>		*	Methionine	aa op sugar fa 2009
5	<i>fap3</i>	*	*	Iodine number, Palmitic acid, Stearic acid	aa op sugar fa 2009, il164, 2il81.1, 2il81.2, 2ky81, 2mn81, 3il83.1, 3il84, 3il87, il0102, meta, ms1999.01, ms2000.02, ms923, ms967
	<i>MTFL</i>		*	Linoleic acid, Seed oil, Oleic acid, Tryptophan	aa op sugar fa 2009, 2il81.1, il0102, ms1999.01, ms967
6	<i>Glyma.06G214800</i>	*	*	Stearic acid	meta, ms1999.01, ms2000.02
	<i>Glyma.06g275100</i>		*	Cysteine	aa op sugar fa 2009
	<i>T</i>		*	Arginine, Cysteine, Isoleucine, Leucine	aa op sugar fa 2009
8	<i>I/AK-HDSH</i>	*	*	Cysteine, Isoleucine, Leucine, Linoleic acid, Lysine, Methionine, Seed oil, Palmitic acid, Stearic acid, Threonine, Valine	aa op sugar fa 2009, meta, ms967
9	<i>Glyma.09g090600</i>	*	*	Palmitic acid	il0102, meta
	<i>R</i>		*	Tryptophan	aa op sugar fa 2009
13	<i>Glyma.13g149700</i>	*	*	Oleic acid, Palmitic acid, Seed protein	meta, ms2000.02
14	<i>fan1</i>	*	*	Linolenic acid	2mn81, 3il83.1, il0102, meta, mn945, ms967
15	<i>Glyma.15g049200</i> "GmSW EET15"	*	*	Linolenic acid, Seed oil, Seed protein, Threonine	aa op sugar fa 2009, 2ky81, 3il83.1, 3il84, il989, meta, ms1999.01, ms923

<b>19</b>	<i>Dt1,</i> <i>Glyma.1</i> <i>9g19430</i> <i>0</i>		*	Linoleic acid, Oleic acid, Valine	aa op sugar fa 2009, ms1999.01, ms2000.02
<b>20</b>	<i>CHR20</i> <i>OP</i>	*	*	Seed oil, Seed protein	aa op sugar fa 2009, 2il81.1, meta, ms1999.01, ms967
<b>14</b> <b>(3)</b>	<i>SACPD-</i> <i>C</i>	*	*	Stearic acid	1il66, 2il81.1, 2mn81, 3il83.1, 4il87, 5il90, il0102, meta, mn945, ms923

394

395

396



397 **Table 3.** List of candidate genes identified for disease resistance/ stress tolerance traits using GWAS from  
 398 individual studies and Meta-GWAS.

Chromosome	Likely Gene	Meta-GWAS	Individual Studies GWAS	Trait(s)	Studies Source
1	<i>RLK3</i>		*	Bacterial pustule	bp488001
3	<i>Glyma.03g127100</i>		*	Pythium root rot	PYU.11002
	<i>Glyma.03g130600</i>	*	*	Iron deficiency chlorosis	Isslepyeye04, meta
	<i>Glyma.03g262500</i>	*		SCN races: 14	meta
	<i>Rps1</i>	*	*	Phytophthora root rot races: 1, 2, 3, 4, 5, 7, 10, 17	meta, PRR1, PRR1.10001, PRR1.10002, PRR1.10004, PRR1.11002, PRR1.11003, PRR1.461592, PRR1.488001, PRR1.492577, PRR1.492990, PRR10, PRR17, PRR17.491404, PRR17.492448, PRR17.492990, PRR2, PRR3, PRR3.492577, PRR3.492990, PRR4, PRR4.492990, PRR5, PRR5.492990, PRR7, PRR7.491404, PRR7.492448, PRR7.492990, prrd196_1, prrd196_3, prrfs04_17, prrfs04_7, prrrs01_1
	<i>Rps7</i>	*	*	Phytophthora root rot races: 1, 2, 3, 4, 5, 7, 10, 17	meta, PRR1, PRR1.10002, PRR1.10003, PRR1.10004, PRR1.11003, PRR1.488001, PRR1.492577, PRR1.492990, PRR10, PRR17, PRR17.491404, PRR17.492448, PRR17.492990, PRR2, PRR3, PRR3.492990, PRR5, PRR5.492990, PRR7, PRR7.491404, PRR7.492448, PRR7.492990, prrfs04_17, prrfs04_7
4	<i>Glyma.04g190400</i>	*	*	SCN races: 3, 4, 14	meta, SCN14, soyscnyoung94_3
	<i>Glyma.04g227900</i>		*	Brown stem rot	bsrcodeall
5	<i>Glyma.05g137500/800</i>		*	Brown stem rot	bsr97, bsrcode492477

<b>6</b>	<i>Glyma.06g199600,197800</i>	*	*	Frogeye leaf spot, race 2	2ky91, Fe2, meta
<b>7</b>	<i>Glyma.07g192200</i>	*	*	SCN races: 1, 3, 5, 14	meta, SCN14, SCN14.491576, SCN14code.491576, soyscnanand_3, soyscnanand_5, soyscnyoung94_3, soyscnyoung94_5
<b>8</b>	<i>Glyma.08g231100</i>	*	*	SCN races: 3, 5, 14	meta, SCN14, soyscnyoung94_5, soyscnyoung94_14
	<i>Rhg4</i>	*	*	SCN races: 1, 3, 5, 14	meta, SCN1, SCN14, soyscnyoung94_3
<b>10</b>	<i>Glyma.10g273300/276600</i>	*	*	SCN races: 14	meta, SCN14, SCN14.491576, SCN14code.491576, soyscnyoung94_14
<b>11</b>	<i>Glyma.11g233500</i>		*	Phytophthora root rot races: 17	PRR17.492990
	<i>Glyma.11g234500 (SNAP11)</i>	*	*	SCN races 1, 3, 4, 14	meta, SCN14, sojascnarelli00, soyscnanand_5, soyscnyoung88_5, soyscnyoung94_5, soyscnyoung94_14
<b>12</b>	<i>Glyma12g22660</i>		*	SCN races: 1	SCN1
<b>13</b>	<i>Glyma.13g222300</i>	*	*	SCN races: 1, 3, 14	meta, SCN14, sojascnarelli00, soyscnyoung94_14
	<i>Rag2,5</i>		*	Soybean aphid	aphidcm02
	<i>Rps3</i>	*	*	Phytophthora root rot races: 1, 4, 12, 20, 25	PRR1, PRR1.10004, PRR1.11003, PRR1.492990, PRR12, PRR20, PRR25, PRR25.491404, PRR25.492990, PRR4, PRR4.492990, meta
	<i>Rsv1</i>		*	Peanut mottle virus	pmv
<b>14</b>	<i>Glyma.14g098900</i>		*	Brown stem rot	bsr97, bsr492477
	<i>NSC14</i>		*	Northern stem canker	NSC, NSC.491493
<b>15</b>	<i>Glyma.15g052000</i>		*	Phytophthora	PRR2

				root rot races: 2	
<b>16</b>	<i>Glyma.16g096900</i>		*	Phytophthora root rot races: 2	PRR2
	<i>Rag3</i>		*	Soybean aphid	aphidcm02
	<i>Rbs1, Rbs2, Rbs3</i>		*	Brown stem rot	bsr97, bsr491584, bsrall, bsrcodeall
	<i>Rcs3</i>	*	*	Frogeye leaf spot, race 2	2il81.1, Fe2, meta
	<i>Rps2</i>	*	*	Phytophthora root rot races: 2, 25	PRR2, meta
<b>17</b>	<i>Glyma.17g090200</i>		*	Bean pod mottle virus	bpmvall
<b>18</b>	<i>Glyma.18g138700</i>		*	Phytophthora root rot races: 5	PRR5, PRR5.492990
	<i>Rhg1</i>	*	*	SCN races: 3, 4, 5, 14	meta, SCN14, soyscnanand_3, soyscnyoung88_5, soyscnyoung94_3, soyscnyoung94_14
	<i>Rps4</i>	*	*	Phytophthora root rot races: 1, 3, 4, 25	meta, PRR1, PRR1.10001, PRR1.10002, PRR1.10004, PRR1.488001, PRR25, PRR25.491404, PRR4

399

400

401 **FIGURE CAPTIONS**

402 **Figure 1.** Significant SNPs from GWAS from individual studies and meta-GWAS. (a) Peaks for seed  
403 related traits, (b) Peaks for flowering and maturity related traits, (c) Peaks for disease resistance related  
404 traits. Symbol position along the x-axis shows the position (in Mb) along the chromosome, while y-axis  
405 symbol position shows the LOD score of the lead SNP for each QTL. X-axis labels indicate position (in  
406 Mb) of tertile points, while y-axis labels show minimum, maximum, and middle of LOD score range for  
407 the given trait class. Shape and color correspond to unique traits.

408 **Figure 2.** Circle plots of significant SNPs identified with (a) GWAS from individual studies, and (b) meta-  
409 GWAS. The peaks in the innermost ring includes seed composition traits, the middle ring includes disease  
410 resistance traits, and the outermost ring includes agronomic traits. Symbol position along the x-axis shows  
411 the position (in Mb) along the chromosome, while y-axis symbol position shows the LOD score of the lead  
412 SNP for each QTL. X-axis labels indicate position (in Mb) of tertile points, while y-axis labels show  
413 minimum, maximum, and middle of LOD score range for the given trait class. Shape and color correspond  
414 to unique traits

415

416

417 **SUPPLEMENTAL FILES**

418 **Supplemental Figure 1.** Comparison of the chromosomal region of (a) *FATB1a*, (b) *Dt1*, (c) *PMDH1* loci  
419 identified using diverse populations. The x-axis indicates the physical location on each chromosome  
420 referring soybean genome version Glyma2.0. The y-axis indicates the pairwise LD  $r^2$  between the lead SNP  
421 and the rest SNPs in the specific region for each population.

422 **Supplemental Figure 2.** SNP at the region of candidate genes (a) *FATB1a*, (b) *Dt1*, (c) *PMDH1*. SNP were  
423 retrieved from Figshare database ([http://figshare.com/articles/Soybean\\_resequencing\\_project/1176133](http://figshare.com/articles/Soybean_resequencing_project/1176133))  
424 based on the genome resequencing study of the 302 diverse soybean lines. For each panel, the x-axis  
425 indicates the physical location of the specific regions on the chromosome. The y-axis indicates the pairwise  
426 LD  $r^2$  between the SNP(s) in the region and the lead SNP, which was also identified in the resequencing  
427 dataset.

428

429

430 **Supplemental Table 1.** Full list of significant marker-trait associations found in individual GWAS and  
431 meta-GWAS.

432 **Supplemental Table 2.** List of studies, methods, and reference literature used to generate phenotypic  
433 datasets.

434 **Supplemental Table 3.** Significant SNPs for stearic acid levels from 3il83.1. Positions in Wms82.1 and  
435 Wms82.2 provided to show alignment differences between the two reference genome versions.

436 **Supplemental Table 4.** Trait definitions, number of QTL detected, and number of candidate genes  
437 assigned for each trait.

438 **Supplemental Table 5.** Full listing of which candidate genes were detected in which study, as well as  
439 whether the association was detected in only individual studies or only in meta-analysis.

440

441

## 442 **References**

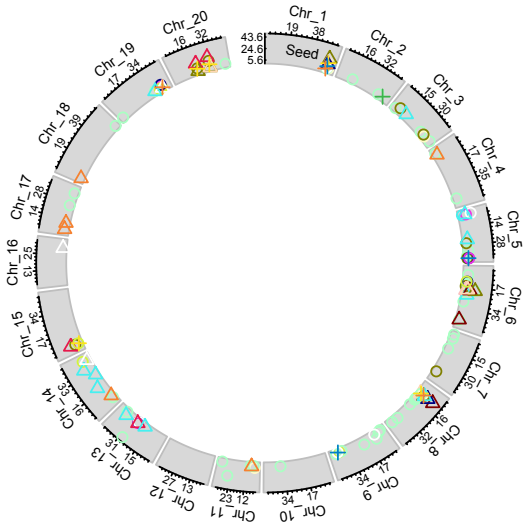
- 443 Assefa, T., J. Zhang, R.V. Chowda-Reddy, A.N. Moran Lauter, A. Singh et al., 2020 Deconstructing the  
444 genetic architecture of iron deficiency chlorosis in soybean using genome-wide approaches. *BMC Plant*  
445 *Biology* 20 (1):42.
- 446 Bachleda, N., Pham, A. & Li, Z. Identifying *FATB1a* deletion that causes reduced palmitic acid content in  
447 soybean N87-2122-4 to develop a functional marker for marker-assisted selection. *Mol Breeding* **36**, 45  
448 (2016). <https://doi.org/10.1007/s11032-016-0468-9>
- 449 Baianu, I., J. Guo, R. Nelson, T. You, and D. Costescu, 2011 NIR Calibrations for Soybean Seeds and Soy  
450 Food Composition Analysis: Total Carbohydrates, Oil, Proteins and Water Contents [v.2] *Nat*  
451 *Proceedings*. <https://doi.org/10.1038/npre.2011.6611.1>
- 452 Bandillo, N., D. Jarquin, Q. Song, R. Nelson, P. Cregan et al., 2015 A Population Structure and Genome-  
453 Wide Association Analysis on the USDA Soybean Germplasm Collection. *The Plant Genome* 8 (3).
- 454 Bandillo, N.B., A.J. Lorenz, G.L. Graef, D. Jarquin, D.L. Hyten et al., 2017 Genome-wide Association  
455 Mapping of Qualitatively Inherited Traits in a Germplasm Collection. *The Plant Genome* 10 (2).
- 456 Bernard, R.L., 1972 Two genes affecting stem termination in soybeans. *Crop Science* 12:235-239.
- 457 Bolormaa, S., J.E. Pryce, A. Reverter, Y. Zhang, W. Barendse et al., 2014 A Multi-Trait, Meta-analysis for  
458 Detecting Pleiotropic Polymorphisms for Stature, Fatness and Reproduction in Beef Cattle. *PLOS*  
459 *Genetics* 10 (3):e1004198.
- 460 Browning, Brian L., and Sharon R. Browning, 2016 Genotype Imputation with Millions of Reference  
461 Samples. *Am J Hum Genet* 98 (1):116-126.
- 462 Cameron, J.N., Y. Han, L. Wang, and W.D. Beavis, 2017 Systematic design for trait introgression projects.  
463 *Theor Appl Genet* 130 (10):1993-2004.
- 464 Chang, D., M.A. Nalls, I.B. Hallgrímsdóttir, J. Hunkapiller, M. van der Brug et al., 2017 A meta-analysis of  
465 genome-wide association studies identifies 17 new Parkinson's disease risk loci. *Nat Genet* 49  
466 (10):1511-1516.
- 467 Chang, H.-X., and G.L. Hartman, 2017 Characterization of Insect Resistance Loci in the USDA Soybean  
468 Germplasm Collection Using Genome-Wide Association Studies. *Frontiers in Plant Science* 8 (670).
- 469 Chang, H.-X., A.E. Lipka, L.L. Domier, and G.L. Hartman, 2016 Characterization of Disease Resistance Loci  
470 in the USDA Soybean Germplasm Collection Using Genome-Wide Association Studies. *Phytopathology*<sup>™</sup>  
471 106 (10):1139-1151.
- 472 Chen, X., F. Zhao, and S. Xu, 2010 Mapping environment-specific quantitative trait loci. *Genetics* 186  
473 (3):1053-1066.
- 474 Clough, S.J., J.H. Tuteja, M. Li, L.F. Marek, R.C. Shoemaker et al., 2004 Features of a 103-kb gene-rich  
475 region in soybean include an inverted perfect repeat cluster of CHS genes comprising the I locus.  
476 *Genome* 47 (5):819-831.

- 477 Cook, D.E., T.G. Lee, X. Guo, S. Melito, K. Wang et al., 2012 Copy number variation of multiple genes at  
478 Rhg1 mediates nematode resistance in soybean. *Science* 338 (6111):1206-1209.
- 479 Cook, J., Mahajan, A. & Morris, A. Guidance for the utility of linear models in meta-analysis of genetic  
480 association studies of binary phenotypes. *Eur J Hum Genet* **25**, 240–245 (2017).  
481 <https://doi.org/10.1038/ejhg.2016.150>
- 482 Coser, S.M., R.V. Chowda Reddy, J. Zhang, D.S. Mueller, A. Mengistu et al., 2017 Genetic Architecture of  
483 Charcoal Rot (*Macrophomina phaseolina*) Resistance in Soybean Revealed Using a Diverse Panel.  
484 *Frontiers in Plant Science* 8 (1626).
- 485 de Azevedo Peixoto, L., T.C. Moellers, J. Zhang, A.J. Lorenz, L.L. Bhering et al., 2017 Leveraging genomic  
486 prediction to scan germplasm collection for crop improvement. *PLOS ONE* 12 (6):e0179191.
- 487 Descriptors for Soybean, 2019. U.S. National Plant Germplasm System.
- 488 Dhondt, S., N. Wuyts, and D. Inzé, 2013 Cell to whole-plant phenotyping: the best is yet to come. *Trends*  
489 *Plant Sci* 18 (8):428-439.
- 490 Diers, B.W., J. Specht, K.M. Rainey, P. Cregan, Q. Song et al., 2018 Genetic Architecture of Soybean Yield  
491 and Agronomic Traits. *G3: Genes|Genomes|Genetics* 8 (10):3367-3375.
- 492 Fang, C., Y. Ma, S. Wu, Z. Liu, Z. Wang et al., 2017 Genome-wide association studies dissect the genetic  
493 networks underlying agronomical traits in soybean. *Genome Biol* 18 (1):161.
- 494 Flint-Garcia, S.A., C. Jampatong, L.L. Darrah, and M.D. McMullen, 2003 Quantitative Trait Locus Analysis  
495 of Stalk Strength in Four Maize Populations Mention of a trademark or proprietary product does not  
496 constitute a guarantee, warranty, or recommendation of the product by the USDA or the University of  
497 Missouri, and does not imply its approval to the exclusion of others that may be more suitable. *Crop*  
498 *Science* 43 (1):13-22.
- 499 Gao, H., and M.K. Bhattacharyya, 2008 The soybean-Phytophthora resistance locus Rps1-k encompasses  
500 coiled coil-nucleotide binding-leucine rich repeat-like genes and repetitive sequences. *BMC Plant Biol*  
501 8:29.
- 502 Gibson, L.R., and R.E. Mullen, 1996 Soybean seed composition under high day and night growth  
503 temperatures. *Int J Mol Sci* 73 (6):733-737.
- 504 Gillman, J.D., M.G. Stacey, Y. Cui, H.R. Berg, and G. Stacey, 2014 Deletions of the SACPD-C locus elevate  
505 seed stearic acid levels but also result in fatty acid and morphological alterations in nitrogen fixing  
506 nodules. *BMC Plant Biol* 14 (1):143.
- 507 Goettel W, Ramirez M, Upchurch RG, An YQ. Identification and characterization of large DNA deletions  
508 affecting oil quality traits in soybean seeds through transcriptome sequencing analysis. *Theor Appl*  
509 *Genet.* 2016;129(8):1577-1593. doi:10.1007/s00122-016-2725-z
- 510 Gu, Z., L. Gu, R. Eils, M. Schlesner, and B. Brors, 2014 circlize implements and enhances circular  
511 visualization in R. *Bioinformatics* 30 (19):2811-2812.
- 512 Hulting, A.G., L.M. Wax, R.L. Nelson, and F.W. Simmons, 2001 Soybean (*Glycine max* (L.) Merr.) cultivar  
513 tolerance to sulfentrazone. *Crop Protection* 20 (8):679-683.

- 514 Lakhssassi, N., S. Liu, S. Bekal, Z. Zhou, V. Colantonio et al., 2017 Characterization of the Soluble NSF  
515 Attachment Protein gene family identifies two members involved in additive resistance to a plant  
516 pathogen. *Sci Rep* 7 (1):45226.
- 517 Lipka, A.E., F. Tian, Q. Wang, J. Peiffer, M. Li et al., 2012 GAPIT: genome association and prediction  
518 integrated tool. *Bioinformatics* 28 (18):2397-2399.
- 519 Liu, B., S. Watanabe, T. Uchiyama, F. Kong, A. Kanazawa et al., 2010 The Soybean Stem Growth Habit  
520 Gene *Dt1* Is an Ortholog of Arabidopsis TERMINAL FLOWER1. *Plant Physiology* 153 (1):198.
- 521 Miller, E.K., 2003 Index to USDA Technical Bulletins, edited by USDA/ARS. National Agricultural Library.
- 522 Nagasubramanian, K., S. Jones, S. Sarkar, A.K. Singh, A. Singh et al., 2018 Hyperspectral band selection  
523 using genetic algorithm and support vector machines for early identification of charcoal rot disease in  
524 soybean stems. *Plant Methods* 14 (1):86.
- 525 Nagasubramanian, K., S. Jones, A.K. Singh, S. Sarkar, A. Singh et al., 2019 Plant disease identification  
526 using explainable 3D deep learning on hyperspectral images. *Plant Methods* 15 (1):98.
- 527 Naik, H.S., J. Zhang, A. Lofquist, T. Assefa, S. Sarkar et al., 2017 A real-time phenotyping framework using  
528 machine learning for plant stress severity rating in soybean. *Plant Methods* 13 (1):23.
- 529 Neyman, J., and E.S. Pearson, 1928 On the use and interpretation of certain test criteria for purposes of  
530 statistical inference, Part I. *Biometrika* 20A (1-2):175-240.
- 531 Parmley, K., K. Nagasubramanian, S. Sarkar, B. Ganapathysubramanian, and A.K. Singh, 2019  
532 Development of Optimized Phenomic Predictors for Efficient Plant Breeding Decisions Using Phenomic-  
533 Assisted Selection in Soybean. *Plant Phenomics* 2019:15.
- 534 Sherman, R.M., J. Forman, V. Antonescu, D. Puiu, M. Daya et al., 2019 Assembly of a pan-genome from  
535 deep sequencing of 910 humans of African descent. *Nat Genet* 51 (1):30-35.
- 536 Singh, A., R.E. Knox, R.M. DePauw, A.K. Singh, R.D. Cuthbert et al., 2014 Stripe rust and leaf rust  
537 resistance QTL mapping, epistatic interactions, and co-localization with stem rust resistance loci in  
538 spring wheat evaluated over three continents. *Theor Appl Genet* 127 (11):2465-2477.
- 539 Singh, D., X. Wang, U. Kumar, L. Gao, M. Noor et al., 2019 High-Throughput Phenotyping Enabled  
540 Genetic Dissection of Crop Lodging in Wheat. *Frontiers in Plant Science* 10 (394).
- 541 Song, Q., D.L. Hyten, G. Jia, C.V. Quigley, E.W. Fickus et al., 2013 Development and Evaluation of  
542 SoySNP50K, a High-Density Genotyping Array for Soybean. *PLOS ONE* 8 (1):e54985.
- 543 Song, Q., D.L. Hyten, G. Jia, C.V. Quigley, E.W. Fickus et al., 2015 Fingerprinting Soybean Germplasm and  
544 Its Utility in Genomic Research. *G3: Genes|Genomes|Genetics* 5 (10):1999.
- 545 Srour, A., A.J. Afzal, L. Blahut-Beatty, N. Hemmati, D.H. Simmonds et al., 2012 The receptor like kinase at  
546 *Rhg1-a/Rfs2* caused pleiotropic resistance to sudden death syndrome and soybean cyst nematode as a  
547 transgene by altering signaling responses. *BMC genomics* 13:368-368.
- 548 Stec, A.O., P.B. Bhaskar, Y.-T. Bolon, R. Nolan, R.C. Shoemaker et al., 2013 Genomic heterogeneity and  
549 structural variation in soybean near isogenic lines. *Frontiers in plant science* 4:104-104.

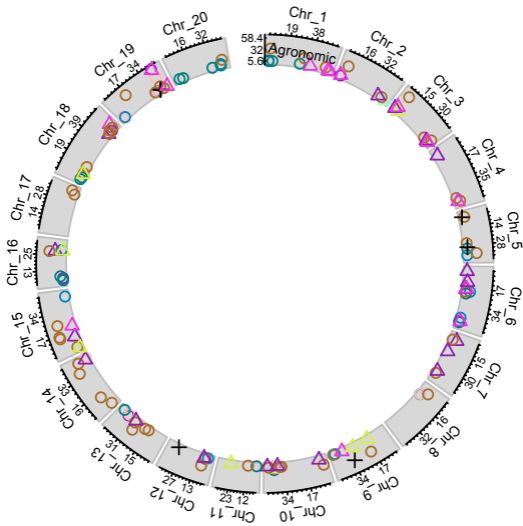


- 550 Takahashi, R., and S. Asanuma, 1996 Association of T gene with chilling tolerance in soybean. *Crop*  
551 *Science* 36:559-562.
- 552 Thapa, R., Carrero-Colón, M. and Hudson, K.A. (2016), New Alleles of *FATB1A* to Reduce Palmitic Acid  
553 Levels in Soybean. *Crop Science*, 56: 1076-1080. doi:[10.2135/cropsci2015.09.0597](https://doi.org/10.2135/cropsci2015.09.0597)
- 554 The 100,000 Genomes Project, 2019. GenomicsEngland.
- 555 Toda, K., D. Yang, N. Yamanaka, S. Watanabe, K. Harada et al., 2002 A single-base deletion in soybean  
556 flavonoid 3'-hydroxylase gene is associated with gray pubescence color. *Plant Mol Biol* 50 (2):187-196.
- 557 Trotta, L., T. Hautala, S. Hämäläinen, J. Syrjänen, H. Viskari et al., 2016 Enrichment of rare variants in  
558 population isolates: single AICDA mutation responsible for hyper-IgM syndrome type 2 in Finland. *Eur J*  
559 *Hum Genet* 24 (10):1473-1478.
- 560 Watanabe, S., Z. Xia, R. Hideshima, Y. Tsubokura, S. Sato et al., 2011 A Map-Based Cloning Strategy  
561 Employing a Residual Heterozygous Line Reveals that the *GIGANTEA* Gene Is Involved in Soybean  
562 Maturity and Flowering. *Genetics* 188 (2):395.
- 563 Willer, C.J., Y. Li, and G.R. Abecasis, 2010 METAL: fast and efficient meta-analysis of genomewide  
564 association scans. *Bioinformatics* 26 (17):2190-2191.
- 565 Zeng, A., P. Chen, K. Korth, F. Hancock, A. Pereira et al., 2017 Genome-wide association study (GWAS) of  
566 salt tolerance in worldwide soybean germplasm lines. *Mol Breeding* 37 (3):30.
- 567 Zhang, J., H.S. Naik, T. Assefa, S. Sarkar, R.V.C. Reddy et al., 2017 Computer vision and machine learning  
568 for robust phenotyping in genome-wide studies. *Sci Rep* 7 (1):44048.
- 569 Zhang, J., A. Singh, D.S. Mueller, and A.K. Singh, 2015 Genome-wide association and epistasis studies  
570 unravel the genetic architecture of sudden death syndrome resistance in soybean. *The Plant Journal* 84  
571 (6):1124-1136.
- 572 Zhang, J., and A.K. Singh, 2020 Genetic Control and Geo-Climatic Adaptation of Pod Dehiscence Provide  
573 Novel Insights into Soybean Domestication. *G3: Genes|Genomes|Genetics* 10 (2):545.
- 574 Zhao, J., C. Sauvage, J. Zhao, F. Bitton, G. Bauchet et al., 2019 Meta-analysis of genome-wide association  
575 studies provides insights into genetic control of tomato flavor. *Nat Comm* 10 (1):1534.
- 576 Zhou, Z., Y. Jiang, Z. Wang, Z. Gou, J. Lyu et al., 2015 Resequencing 302 wild and cultivated accessions  
577 identifies genes related to domestication and improvement in soybean. *Nat Biotech* 33 (4):408-414.
- 578 Zhu-Shimoni, J.X., and G. Galili, 1998 Expression of an Arabidopsis Aspartate Kinase/Homoserine  
579 Dehydrogenase Gene Is Metabolically Regulated by Photosynthesis-Related Signals but Not by  
580 Nitrogenous Compounds. *Plant Physiology* 116 (3):1023-1028.
- 581 Ćwiek-Kupczyńska, H., T. Altmann, D. Arend, E. Arnaud, D. Chen et al., 2016 Measures for  
582 interoperability of phenotypic data: minimum information requirements and formatting. *Plant Methods*  
583 12 (1):44.



## Seed related

 Iodinenum	 Oleic	 Cysteine	 Threonine
 Linoleic	 Palmitic	 Isoleucine	 Tryptophan
 Linolenic	 Protein	 Leucine	 Valine
 Methionine	 Seedweight	 Lysine	
 Mottling	 Stearic	 Quality	
 Oil	 Arginine	 Sucrose	



## Agronomic traits

○ Branching

○ Flowerdate

○ Height

○ Lodging

○ Maturity date

○ Maturity group

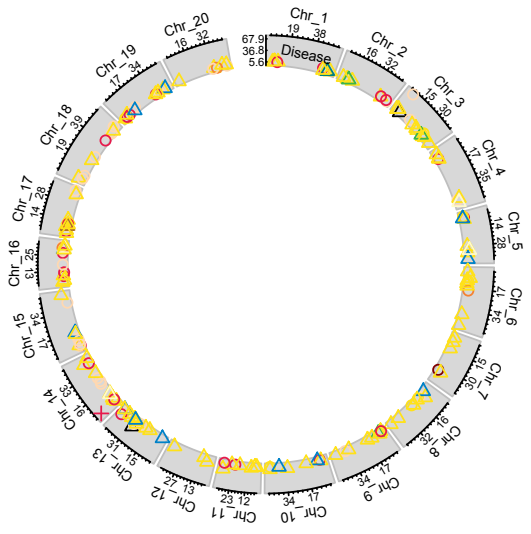
△ Shattering

△ Stemtermination

△ Yield

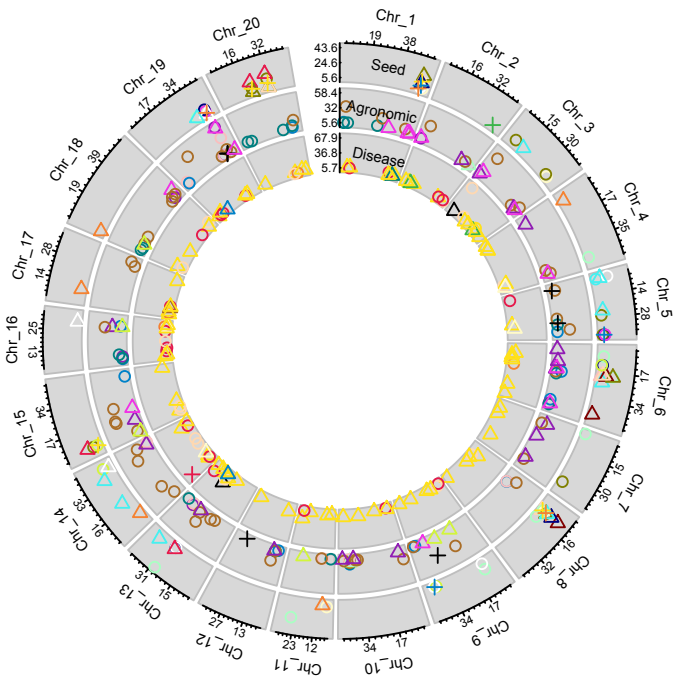
△ Dftm

+ Root fluorescence



## Disease resistance

○ APHID	△ PMV	△ BPMV
○ CHLOROSIS	△ PYTHIUM	△ BSR
○ FROGEYE2	△ SCN	+ STEMCANKER
○ MEXBEANBTL	△ SDS	
○ PRR	△ BP	





## Seed related

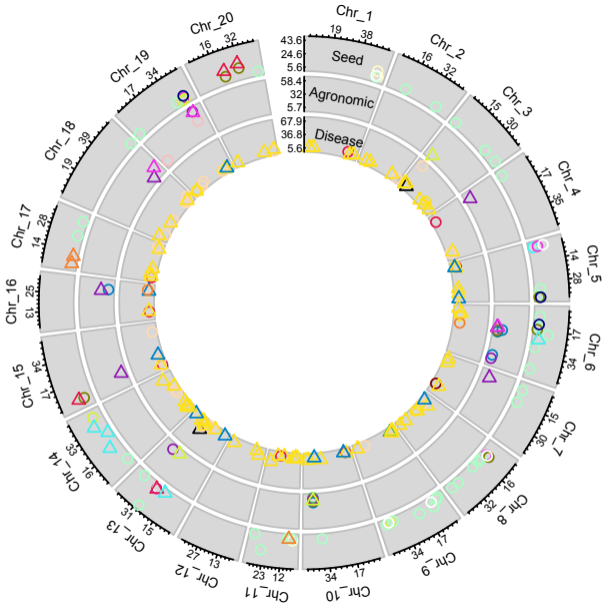
- Iodinenum
- Linoleic
- Linolenic
- Methionine
- Mottling
- Oil
- Oleic
- Palmitic
- △ Protein
- △ Seedweight
- △ Stearic
- △ Arginine
- △ Cysteine
- △ Isoleucine
- △ Leucine
- △ Lysine
- △ Quality
- + Sucrose
- + Threonine
- + Tryptophan
- + Valine

## Agronomic

- Flowerdate
- Height
- Lodging
- Maturity date
- Maturity group
- △ Shattering
- △ Stemtermination
- △ Yield
- △ Dftm
- + Root fluorescence

## Disease resistance

- APHID
- CHLOROSIS
- FROGEYE2
- PRR
- △ PMV
- △ SCN
- △ BP
- △ BPMV
- △ BSR
- △ PYTHIUM
- △ SDS
- + STEMCANKER



## Seed related

- Iodinenum
- Linoleic
- Linolenic
- Methionine
- Mottling
- Oil
- Oleic
- Palmitic
- △ Protein
- △ Seedweight
- △ Stearic

## Agronomic

- Branching
- Flowerdate
- Height
- Lodging
- Maturity date
- Maturity group
- △ Shattering
- △ Stemtermination
- △ Yield

## Disease resistance

- APHID
- CHLOROSIS
- FROGEYE2
- MEXBEANBTL
- PRR
- △ PMV
- △ PYTHIUM
- △ SCN
- △ SDS