

2016

Multi-Population Selective Genotyping to Identify Soybean [*Glycine max* (L.) Merr.] Seed Protein and Oil QTLs

Piyaporn Phansak

University of Nebraska-Lincoln, pphansak@yahoo.com

Watcharin Soonsuwon

University of Nebraska-Lincoln

David L. Hyten

University of Nebraska-Lincoln, david.hyten@unl.edu

Qijian Song


USDA-ARS, qijian.song@ars.usda.gov

Perry B. Cregan

USDA-ARS

See next page for additional authors

Follow this and additional works at: <https://digitalcommons.unl.edu/agronomyfacpub>

 Part of the [Agricultural Science Commons](#), [Agriculture Commons](#), [Agronomy and Crop Sciences Commons](#), [Botany Commons](#), [Horticulture Commons](#), [Other Plant Sciences Commons](#), and the [Plant Biology Commons](#)

Phansak, Piyaporn; Soonsuwon, Watcharin; Hyten, David L.; Song, Qijian; Cregan, Perry B.; Graef, George L.; and Specht, James E., "Multi-Population Selective Genotyping to Identify Soybean [*Glycine max* (L.) Merr.] Seed Protein and Oil QTLs" (2016). *Agronomy & Horticulture -- Faculty Publications*. 894.
<https://digitalcommons.unl.edu/agronomyfacpub/894>

This Article is brought to you for free and open access by the Agronomy and Horticulture Department at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Agronomy & Horticulture -- Faculty Publications by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

Authors

Piyaporn Phansak, Watcharin Soonsuwon, David L. Hyten, Qijian Song, Perry B. Cregan, George L. Graef, and James E. Specht

Multi-Population Selective Genotyping to Identify Soybean [*Glycine max* (L.) Merr.] Seed Protein and Oil QTLs

Piyaporn Phansak,^{*,1} Watcharin Soonswon,^{*,2} David L. Hyten,^{*} Qijian Song,[†] Perry B. Cregan,[†] George L. Graef,^{*} and James E. Specht^{*,3}

^{*}Department of Agronomy and Horticulture, University of Nebraska, Lincoln, Nebraska 68583-0915, and [†]Soybean Genomics and Improvement Laboratory, United States Department of Agriculture - Agricultural Research Service (USDA-ARS), Beltsville, Maryland 20705-2325

ABSTRACT Plant breeders continually generate ever-higher yielding cultivars, but also want to improve seed constituent value, which is mainly protein and oil, in soybean [*Glycine max* (L.) Merr.]. Identification of genetic loci governing those two traits would facilitate that effort. Though genome-wide association offers one such approach, selective genotyping of multiple biparental populations offers a complementary alternative, and was evaluated here, using 48 F_{2:3} populations ($n = \sim 224$ plants) created by mating 48 high protein germplasm accessions to cultivars of similar maturity, but with normal seed protein content. All F_{2:3} progeny were phenotyped for seed protein and oil, but only 22 high and 22 low extreme progeny in each F_{2:3} phenotypic distribution were genotyped with a 1536-SNP chip (ca. 450 bimorphic SNPs detected per mating). A significant quantitative trait locus (QTL) on one or more chromosomes was detected for protein in 35 (73%), and for oil in 25 (52%), of the 48 matings, and these QTL exhibited additive effects of $\geq 4 \text{ g kg}^{-1}$ and R^2 values of 0.07 or more. These results demonstrated that a multiple-population selective genotyping strategy, when focused on matings between parental phenotype extremes, can be used successfully to identify germplasm accessions possessing large-effect QTL alleles. Such accessions would be of interest to breeders to serve as parental donors of those alleles in cultivar development programs, though 17 of the 48 accessions were not unique in terms of SNP genotype, indicating that diversity among high protein accessions in the germplasm collection is less than what might ordinarily be assumed.

KEYWORDS

germplasm
survey tool
QTLs: pleiotropy
or linkage
rare alleles
nonunique SNP
accessions
selection bias

Soybean [*Glycine max* (L.) Merr.], produced mainly in North and South America and Asia, is high in seed protein (40%) and oil (20%). These two seed constituents are consumed worldwide by domestic livestock,

poultry, and fish (*i.e.*, soybean meal), and by humans (*i.e.*, cooking oil and Asian-style soybean food products). Soybean seed protein is inherited quantitatively, though more in an oligenic than a polygenic fashion, and is highly heritable (Burton 1987; Wehrmann *et al.* 1987; Wilcox 1998; Cober and Voldeng 2000). However, highly negative phenotypic and genotypic correlations of seed protein with seed yield and oil content have been routinely detected in biparental breeding populations (Burton 1987). Long-term selection for greater yield has also depressed protein and elevated oil (Rincker *et al.* 2014).

When soybean molecular markers became available in the 1990s (Keim *et al.* 1990), the detection and mapping of soybean quantitative trait loci (QTL) soon began. Diers *et al.* (1992) was the first to detect a major seed protein and oil QTL on soybean chromosome 20. Many seed protein and oil QTL have since been reported, and a listing of these, as well as QTL for other traits, can be found in SoyBase (Grant *et al.* 2010; <http://www.soybase.org>). However, nearly all of the protein and oil QTL reported to date have not been confirmed, and not one has yet been cloned. The additive effect values for these QTL are likely

Copyright © 2016 Phansak *et al.*

doi: 10.1534/g3.116.027656

Manuscript received January 28, 2016; accepted for publication March 25, 2016; published Early Online April 1, 2016.

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

Supplemental material is available online at www.g3journal.org/lookup/suppl/doi:10.1534/g3.116.027656/-/DC1

¹Present address: Division of Biology, Faculty of Science, Nakhon Phanom University, Muang, Nakhon Phanom, 48000 Thailand

²Present address: Department of Plant Science, Faculty of Natural Resources, Prince of Songkla University, Hat Yai, Songkhla, 90112 Thailand

³Corresponding author: 363 Keim Hall, Department of Agronomy and Horticulture, Fair Street, East Campus, University of Nebraska, Lincoln, NE 68583-0915.

E-mail: jspecht1@unl.edu

inflated due to the use of small population sizes in the published reports, because of an intrinsic QTL detection problem known as selection bias (Beavis 1998; Xu 2003; Broman and Sen 2009).

The parental sources of most high protein genes (*i.e.*, QTL alleles) used by soybean breeders are typically the high protein accessions acquired from the USDA Soybean Germplasm Collection. Trait data are documented in the Germplasm Resources Information Network (GRIN) for the 21,728 *G. max* accessions present in the collection as of December 31, 2015 (<http://ars-grin.gov/npgs>). For just the 12,141 *G. max* accessions in maturity groups (MGs) 0–IV, substantive variation clearly exists for each trait (Supplemental Material, Figure S1), though it is also evident in this large set of accessions that seed protein exhibits a negative relationship with seed oil and yield. Knowing the allelic status of the seed protein QTL in these accessions would help breeders select donor parents, and also allow a focus on those QTL that have an allele that exerts a large positive additive effect on seed protein, coupled with a smaller negative pleiotropic effect on seed oil and yield.

The allelic diversity of soybean seed protein QTL in germplasm collections would seem to be best addressed using an association analysis method (Thornsberry *et al.* 2001; Semagn *et al.* 2010; Korte and Farlow 2013). Bandillo *et al.* (2015) recently conducted a genome-wide association study (GWAS) involving all of the *G. max* accessions in the collection with protein and oil phenotypic data (*i.e.*, $n = 12,116$) that had been genotyped with a 50K single nucleotide polymorphism (SNP) chip. Strong signals were detected on chromosomes 20 and 15, plus weaker signals on chromosomes 13, 6, and 5. The authors of the latter study noted that the use of large numbers of accessions for GWAS greatly improved statistical power, and provided exceptional map resolution for the ultimate identification and cloning of the causal genes underpinning the major QTL on chromosomes 20 and 15. However, the rarity of a high protein allele at a QTL (and the coincident rarity of linkage-coupled alleles at QTL-flanking SNPs) can be an issue when selecting accession samples for GWAS. This was evident when Bandillo *et al.* (2015) stratified the 12K accessions into smaller subsets by sorting them into seven different countries of origin, or alternatively, into eight different MG classes. In some country subsets, and some MG subsets, the frequency of the protein-enhancing alleles of the QTL was lower than the GWAS minimum allele frequency (MAF) cutoff value, thereby resulting in no detection of one or two or all of the above-listed QTL. Thus, despite the power offered by GWAS in QTL detection, this “rare allele” problem can result in QTL not being detected in GWAS that were previously identified and confirmed to be present in biparental QTL mapping populations—wherein the frequency of the two parental alleles at any segregating high protein QTL is always expected to be near 0.5.

Selective genotyping (SG) was a term first used by Lander and Botstein (1989) to describe those cases of QTL mapping in which only the most informative individuals—those occupying the lowest and highest tails of a phenotyped trait distribution—were genotyped. A trait-based QTL detection approach had previously been conducted in plants (Stuber *et al.* 1980, 1982). Lebowitz *et al.* (1987) and Darvasi and Soller (1992) subsequently formulated and discussed the statistical issues relevant to SG. When using SG, one must still phenotype the entire population to conduct an unbiased QTL analysis (Darvasi 1997; Darvasi and Soller 1992; Muranty and Goffinet 1997; Sen *et al.* 2005, 2009). Optimal efficiency is usually achieved with SG if one does not genotype more than the upper (and lower) 20–25% of the mapping population for a given trait. Sun *et al.* (2010) noted that the optimum size of the tail proportion of a population was governed by a balance between QTL detection power and total cost, which was reflective of the ratio between genotyping and phenotyping costs.

Soybean breeders typically rely on near-infrared reflection (NIR) instrumentation to estimate the seed protein and oil content of germplasm lines (Hymowitz *et al.* 1974). About 100–200 soybean seed samples can be nondestructively phenotyped per hour of effort. Seed protein and oil phenotyping is relatively inexpensive, though it is labor-intensive. Thus, SG would seem to offer a cost-effective means of conducting a QTL analysis of multiple biparental mapping populations segregating for major-effect high and low alleles at seed protein QTL.

Ayoub and Mather (2002) demonstrated that if SG had been applied to just the lowest 10% and highest 10% of each trait in a North American barley mapping population, the resultant SG-based QTL analyses would have been sufficient to detect *all* of the grain and malt quality QTL that had been identified based on a genotyping of the entire population of about 140–150 lines. This publication triggered our interest in using a multiple mapping population SG approach as a means of surveying a large sample of high protein soybean germplasm accessions for the presence of high protein alleles at known and unknown QTL. The availability of a 1536-SNP marker assay—the Universal Soybean Linkage Panel 1.0 (USLP 1.0) developed by Hyten *et al.* (2010)—in a 96-well genotypic sample format, was another contributing factor leading us to examine the utility of a multi-population SG strategy to identify alleles of QTL that condition high seed protein, but which may have a low frequency in germplasm collection accessions. Rare QTL alleles are difficult to detect in a GWAS, both in theory (Raychaudhuri 2011; Ladouceur *et al.* 2012), and in practice (Bandillo *et al.* 2015). In that regard, we hypothesized that a SG strategy might mitigate the traditional GWAS rare-allele problem (Korte and Farlow 2013).

We thus report here on the use of a SG-based QTL analysis to survey 48 soybean populations, averaging about 224 F₂ plants, derived from the mating of 48 high-seed-protein soybean germplasm accessions in seven MGs (spanning 000 to IV) to one of seven high-yielding lower protein cultivars with a matching MG. The ultimate objective of this study was to discern whether a multiple-population SG approach could be used to identify and map both known and unknown protein QTL in these high protein accessions that might serve as donor parents.

MATERIALS AND METHODS

Parents and population development

To minimize the segregation of major genes controlling date of flowering/maturity in the F₂ generation, the high seed protein accessions of a given MG were mated to a high-yield cultivar of ordinary seed protein content of the same MG. The parents are shown in Table 1, with each M-code-designated male parent listed just below the respective set of female parents to which that male parent was mated (except for MG V P1183, which was reciprocally mated to MG IV P1181M). The phenotypic data in Table 1 (except as footnoted) were extracted from the GRIN website (<https://npgsweb.ars-grin.gov/gringlobal/descriptors.aspx>). The 48 female parents had a GRIN-based seed protein content that ranged from 473 to 529 g kg⁻¹ (*i.e.*, zero seed moisture, dry weight basis), whereas the range for the seven male parents was 382 to 430 g kg⁻¹ (Table 1); the latter range is typical for cultivars currently being grown in the North Central United States soybean production area.

Pollinations for all 48 matings were made in the summer growing season, and were successful in terms of generating putative F₁ seeds that were individually hand-harvested in the fall and packaged by pod. The F₁ to F₂ generation advance was conducted in a greenhouse. To ensure the authenticity of putative F₁ plants, a known parentally polymorphic SSR marker was used to genotype each F₁ to confirm F₁ hybridity in the 48 matings. Marker-confirmed F₁ plants from each mating were

Table 1 The 48 high seed protein accession female parents, and the seven ordinary seed protein cultivar male parents (M-suffixed codes), ordered by soybean maturity group (MG), and then by mating and parent code. The seed protein and oil values listed for the female and male parents are those available in the Germplasm Resources Information Network (GRIN) website (but see footnote for exceptions).

Mating No.	Parent Code ^a	Maturity Group	Seed ^b		Germplasm Accession		Stem Habit	Flower Color	Pubescence		Pod Color	Seed Coat		Hilum Color
			Protein g kg ⁻¹	Oil g kg ⁻¹	Number ^c	Name (if Any)			Origin ^d	Color	Form	Luster	Color	
1	P1001	000	529	151	PI 153296	V-4		D	P	T	E	S	Gn	Bl
2	P1002	000	504	158	PI 189963	Geant Vert		D	P	T	E	D	Gn	Bl
3	P1003	000	522	155	PI 548399	Pando		D	P	T	E	S	Gn	Bl
4	P1004	000	477	156	PI 372423	Ronset 4		D	P	T	E	I	Lgn	Bl
5	P1005	000	512	157	FC 30687	Kosodiguri Ext Early		D	P	T	E	I	Gn	Bl
6	P1006	000	511	158	PI 153293	N-34		D	P	T	E	S	Gn	Bl
7	P1007	000	478	161	PI 372412	Hercumft		D	P	T	E	S	Lgn	Bl
8	P1009	000	522	159	PI 548414	Sioux		D	P	T	E	S	Gn	Bl
—	P1021M	000	430	199	PI 567787	OAC Vision		N	P	T	E	D	Y	Tn
9	P1022	00	507	158	PI 153302	V-16		D	P	T	E	S	Gn	Bl
10	P1023	00	526	157	PI 159764	—		D	P	T	E	S	Gn	Bl
11	P1024	00	485	164	PI 438415	Ronest 4		S	P	T	E	I	Gn	Bl
12	P1025	00	508	147	PI 153301	V-14		D	P	T	E	S	Gn	Bl
13	P1026	00	489	173	PI 189880	Bitterhof		N	P	G	E	S	Y	Y
14	P1027	00	510	148	PI 153297	V-6		D	P	T	E	S	Gn	Bl
15	P2211	00	(-)	(-)	—	HHP		N	Lp	G	?	?	Ib ^f	Ib
16	P2212	00	(486)	(164)	—	AC Proteus		N	P	T	?	D	Y	Br
17	P2213	00	(456)	(181)	—	AC Proteina		N	P	T	?	?	Y	Br
—	P1038M	00	415	185	PI 602897	Jim		N	P	G	E	I	Y	Y
18	P1039	0	480	144	PI 427138	Choseng No. 1		D	W	G	A	D	Y	Bf
19	P1040	0	488	195	PI 261469	Wasedaizu No. 1		N	W	G	A	D	Y	Bf
20	P1041	0	485	177	PI 181571	No. 58		N	W	G	A	D	Y	Bf
21	P1042	0	483	150	PI 424148	Shirome		N	W	G	A	I	Y	Bf
22	P1043	0	473	156	PI 423954	Shirome		D	W	G	Sa	D	Y	Bf
23	P1044	0	494	160	PI 154196	No. 51		D	P	T	E	D	Gn	Bl
—	P1053M	0	403	196	PI 602594	MN0301		N	P	G	E	I	Y	Y
24	P1054	I	484	155	PI 437088A	DV-147		N	P	T	E	D	Y	Br
25	P1055	I	514	144	PI 423949	Saikai 20		D	Lp	G	A	I	Y	Bf
26	P1056	I	495	141	PI 427141	Seuhae No. 20		S	P	T	E	D	Y	Br
27	P1057	I	482	138	PI 437716A	Sjuj-dja-pyn-da-do		S	P	G	Sa	I	Y	Bf
28	P1058	I	489	149	PI 423942	Saikai 1		D	P	G	A	I	Y	Bf
—	P1074M	I	(407)	(195)	PI 602593	MN1301		N	W	G	E	D	Y	Bf
29	P1075	II	499	157	PI 423948A	Saikai 18		N	B	G	E	S	Y	Bf
30	P1076	II	482	154	PI 437112A	Vir 249		N	W	G	E	S	Y	Y
31	P1098	II	484	191	PI 548608	Provar		N	P	T	E	D	Y	Br
—	P1106M	II	382	195	PI 597386	Dwight		N	P	T	E	D	Y	Bl
32	P1107	III	504	132	PI 445845	Szu yueh pa		D	W	G	A	D	Y	Bf
33	P1108	III	494	167	PI 398516	KAERI-GNT 310-1		D	P	Lt	E	D	Y	Y
34	P1109	III	477	170	PI 91725-4	Akazu		D	W	G	Sa	D	Y	Bf
35	P1110	III	493	165	PI 340011	—		D	P	G	E	D	Y	Y
36	P1111	III	478	162	PI 243532	Kariho-takiya		D	W	T	E	S	Y	Br

(continued)

■ **Table 1, continued**

Mating No.	Parent Code ^a	Maturity Group	Seed ^b		Germplasm Accession		Origin ^d	Stem Habit	Flower Color	Pubescence		Pod Color	Seed Coat		Hilum Color
			Protein	Oil	Number ^c	Name (if Any)				Color	Form		Luster	Color	
Descriptor Code ^e															
37	P1113	III	497	168	PI 408138C	KAS 640-7	South Korea	D	P	G	E	Br	D	Y	Y
38	P1121	III	494	177	PI 398672	KAERI-GNT 301-1	South Korea	D	Dp	T	E	Br	S	Rbr ^f	Rbr
39	P1122	III	484	184	PI 360843	Oshimashirome	Japan	N	W	G	E	Br	I	Y	Y
—	P1137M	III	411	194	PI 597387	Pana	Illinois	N	P	G	E	Br	D	Y	Bf
40	P1138	IV	479	157	PI 253666A	No. 17	China	N	W	G	Sa	Br	I	Y	Bf
41	P1139	IV	507	151	PI 407788A	ORD 8113	South Korea	D	P	G	E	Tn	S	Y	Bf
42	P1140	IV	493	155	PI 424286	KAS 239-4	South Korea	D	P	G	E	Tn	D	Y	Bf
43	P1142	IV	488	166	PI 407877B	KAERI 511-11	South Korea	D	P	G	E	Br	D	Y	Bf
44	P1143	IV	488	158	PI 398704	KAS 330-9-1	South Korea	D	P	G	E	Br	I	Y	Bf
45	P1145	IV	491	160	PI 398970	KLS 630-1	South Korea	D	P	G	E	Tn	D	Y	Lbf
46	P1146	IV	493	159	PI 407823	—	South Korea	D	P	G	E	Tn	I	Y	Bf
47	P1152	IV	492	161	PI 407773B	KAS 330-9-2	South Korea	D	W	T	E	Tn	I	Y	Bl
—	P1181M	IV	424	180	PI 606748	Rend	Illinois	N	W	G	E	Br	D	Y	Bf
48	P1183	V	476	195	PI 458256	KAS 578-1	South Korea	D	P	G	Sa	Br	I	Y	Y

The seed protein and oil values listed for the female and male parents are those available in the Germplasm Resources Information Network (GRIN) website (but see footnote for exceptions).

^a Nebraska field nursery parent identification number. The suffix M denotes a male parent (i.e., the seven agronomic cultivars mated to females of the same MG).

^b Seed protein and oil values are not available for these four non-GRIN entries: HHP—Brummer *et al.* (1997) provided details on this high protein accession and its likely source; AC *Proteus* and AC *Proteina*—protein and oil values shown here were reported by Voldeng *et al.* (1996); MN1301—protein and oil values were reported in Hill *et al.* (2008).

^c The solid-line and dashed-lined underscoring identifies two groups of accessions that, within each group, were not unique in terms of their SNP genotype.

^d The non-Asian origin listed for many high protein accessions is, in fact, simply the location of the organization (i.e., mostly European germplasm collection agencies) that donated those accessions to the USDA germplasm collection, but did not provide information as to where in Asia the accession was originally collected.

^e GRIN descriptor codes: D, determinate; IN, indeterminate; S, semi-determinate; Dp, dark purple; Lp, light purple; B, blue; W, white; Bl, black; lb, imperfect black; Y, yellow; Br, brown; Bf, buff; Tn, tan; T, tawny; G, gray; Gn, green; Lgn, light green; E, erect; A, appressed; Sa, semi-appressed; Rb, red-brown; S, shiny; I, intermediate; D, dull.

^f A nonyellow darkly pigmented seed coat color interferes with NIR-based protein and oil measurements. With respect to these two specific female parents, we discarded the homozygous recessive fraction (1/4) of the total F₂ plants that produced F_{2:3} seed progenies that had darkly pigmented (nonyellow) colored seed coats.

individually harvested at maturity to obtain F₂ seeds. Population-specific F_{1,2} seed progeny were planted the following summer into 48 single rows (30 m long; 76.2 cm row spacing). About 300 F₂ seeds of each mating were planted in a row, with a goal of obtaining about 250 F₂ plants bearing F₃ seed. Parental seed and confirmed F₁ seed also were planted in repetitive sections of the same row. All F₂ plants were numerically tagged after emergence (during leaf tissue collection), and surviving tagged plants were gathered at maturity to be individually threshed to obtain F₃ seed.

Phenotypic trait measurement

The F_{2,3} seed progenies, the F_{1,2} seed progenies, and parental seed of a given mating, plus seed of four checks (*i.e.*, breeding lines known to be low or high in seed protein), were evaluated for seed protein, oil, and moisture content using a near-infrared reflectance (NIR) analyzer (Infratec model 1255 NIR Food and Feed Grain Analyzer, Ultra Tec Manufacturing Inc. Santa Ana, CA). The four check samples were used at the beginning and end of each day to confirm that the NIR instrument was operating during the day within its performance standards. Seed protein and oil values were output on a zero per cent seed moisture basis.

One complete replicate of the NIR-measured protein (and oil) data was obtained for all available F_{2,3} progenies in each of the 48 populations. Though each population required about 2 hr of assay time, only two (and on occasion, three) 2-hr assays could be conducted on a given day due to worker availability, instrument warm-up and prep time, etc. Thus, this 48-population NIR assay effort required about five contiguous weeks of workdays to complete. The F_{2,3} seed progeny in each mating were then ranked from lowest to highest based on their measured seed protein value. After completing a second replicate of NIR-assays of all progenies in just two populations (*i.e.*, matings 43 and 44; Table 1), it was determined that the F_{2,3} seed progenies present in highest and lowest 10% fractions of the first and second replicate assays were essentially the same progenies (data not shown). Thus, to reduce the phenotyping effort and time required to identify the F_{2,3} progenies occupying just the lowest and highest 10% fractions, a second replicate of NIR measurement was performed only on the highest and lowest 20% fractions in each of the remaining 46 populations. In each low and high 20% of 2-rep means, those F_{2,3} seed progenies ranking at the extreme ends of those 20% fractions were selected to become the corresponding 10% tail fractions of the seed protein distribution. Leaf tissue samples of the F₂ plant progenitors of just these extreme progenies (*i.e.*, 22 high and 22 low protein) were subsequently used for SG.

SNP marker genotyping

Standard methods for leaf collection and DNA extraction methods were used (for details, see File S1). All steps in the SNP genotyping assays of the parental, F₁, and F₂ DNA samples of the 48 SG populations (*i.e.*, a total of 24 plates) were conducted by personnel at the Soybean Genomics and Improvement Laboratory, USDA-ARS, BARC-West, Beltsville, MD, using the Illumina GoldenGate assay and an Illumina Beadstation 500 (Illumina Inc., San Diego, CA). A soybean-specific USLP 1.0 GoldenGate assay had been developed by Hyten *et al.* (2010) for 1536 SNP markers that were distributed (relatively) uniformly across the 20 chromosomes of the soybean genome. Automatic genotype calling for each SNP locus in each DNA sample in the first 10 two-population plates was conducted using Illumina GeneCall software, but the newer BeadStudio software was used for the 14 remaining two-population plates. All automated genotype call output was manually examined and adjusted as needed. Illumina base-pair allele calls were phase-translated into two-character genotype codes of **AA** for the

high yield (normal protein) elite male parent, **BB** for the high protein accession female parent, and **AB** for the F₁ progenitor of the F₂ population, but were subsequently converted to single character codes of **A H B** – (*i.e.*, dash was assigned to missing genotypes) for use with linkage and analysis software.

Phenotypic data analysis

The distributional statistics of the F_{2,3} phenotypic data collected for seed protein and oil content, and their phenotypic correlation (in each population), were examined using the statistical and graphics R software (<http://cran.r-project.org/>; version 3.1.3; 2015-03-9). A Shapiro-Wilk test of normality (Type I error criterion set to $\alpha=0.01$) was performed on each of the 48 seed protein and oil phenotypic distributions. A Pearson correlation coefficient for protein and oil was also computed for each population.

Individual F₂ plants (and the F₃ seed progeny each produced) were the experimental units in this experiment. Because F₂ plants cannot be naturally replicated to obtain an estimate of environmental variance, NIR assays were performed on the seed progenies harvested from the multiple homozygous female and male parent plants that had been grown in interspersed sections of the same nursery row containing F₂ plants. Parental assay data were used to obtain an indirect estimate of the environmental variance using the following equation:

$$\sigma_e^2 = (1/2)(\sigma_{pFem}^2 + \sigma_{pMal}^2)$$

where σ_{pFem}^2 and σ_{pMal}^2 were the respective phenotypic variances in the seed protein for the seed produced by the high protein female parent, and by the high yield (but ordinary protein) male parent, respectively. The genetic variance component of the F_{2,3} progeny phenotypic variance was then estimated by subtraction, using this formula:

$$\sigma_g^2 = \sigma_p^2 - \sigma_e^2$$

where σ_p^2 was the F_{2,3} progeny phenotypic variance.

A broad sense heritability (H²) estimate was then obtained in the usual manner for each of the 48 populations (Bernardo 2010):

$$H^2 = \frac{\sigma_g^2}{\sigma_p^2} \times 100\%$$

QTL analysis

The R/qtl software package (<http://www.rqtl.org/>) was used in this study. A *.csv file containing phenotypic and genotypic data in a R/qtl csvr format was prepared for each of the 48 populations, and then error-checked prior to the QTL analysis (for details, see File S1). The maximum likelihood method of interval mapping, using the Expectation-Maximization (EM) algorithm, as implemented in R/qtl, was used for QTL detection (Xu and Vogel 2000; Sen *et al.* 2009). Estimates of chromosomal QTL map positions in each of the 48 populations were obtained not only for the SG trait of seed protein, but also for the non-SG trait of seed oil, primarily because of the well-known coinheritance of these two negatively correlated traits. With SG, stratified permutation testing was necessary (Manichaikul *et al.* 2007), and was applied to just the 44 genotyped F₂ progenitors of the selected F_{2,3} progeny (*i.e.*, 22 low/22 high protein phenotypes) to obtain a (QTL peak) LOD score significance criterion for a genome-wide Type I error α of 0.05 +/- SE of 0.005. To attain this degree of precision (see p. 106 in Broman

and Sen 2009), 1900-replicate permutation tests were conducted for each trait in each population. The protein (or oil) additive (*a*) and dominance (*d*) effects conditioned by each marker on each chromosome were first examined graphically, but subsequently, these two effects were numerically estimated for just the putative QTL exhibiting the largest peak LOD score on each chromosome. This estimation used the phenotype means for each of the A, H, and B genotypes of the SNP marker, or a pseudo-marker nearest to the putative QTL. The heritability of each presumptive single QTL on a chromosome is the fraction of the phenotypic variance (*i.e.*, R^2) explained by that QTL, which was estimated with the following equation (p. 77 in Broman and Sen 2009):

$$R^2 = 1 - 10^{(-2/n)\text{LOD}}$$

where *n* is the number of phenotyped $F_{2,3}$ progenies in each population, and LOD is the \log_{10} likelihood ratio (LR) attained by that QTL at its peak map position in the R/qtl scanone output.

The QTL detected in this study were declared statistically significant only if the observed peak LOD score exceeded a population-specific, permutation-generated LOD score computed for a genome-wide Type I error of $\alpha = 0.05$. The chromosomal locations of these QTL were compared with the locations of QTL detected using GWAS in the recent reports, and also the QTL detected in older publications listed in SoyBase. In the latter reports, the authors often used a lower significance threshold for QTL declaration (*i.e.*, $\text{LOD} \geq 3.0$), which in most cases was also a comparison-wise threshold that was not adjusted for multiple testing.

Data availability

Phenotype and genotype data for the 48 F_2 populations and three combined sets of F_2 populations (*.csv files) will be available on SoyBase (www.soybase.com), along with the R/qtl command code (*.txt files). Supplemental files include: File S1 contains additional *Materials and Methods* details; Figure S1 illustrates genetic diversity for seed protein/oil in the Soybean Germplasm Collection; Figure S2 shows the chromosomal map positions of the 1536 SNPs, and the 452 SNPs in the (example) SG mating 1; Figure S3 depicts the chromosomal map positions of SoyBase-listed QTL reported to date; Table S1 documents the original identification codes for the 1536 SNPs aligned with the shorter five-digit Snnnnn names we used to reduce computer memory usage, and to lessen printed table space in this report; Table S2 and Table S3 contain population-specific data for the respective phenotypic and genotypic data after R/qtl error-checking; Table S4 contains the parameter data derived from the population-specific QTL analyses, ordered by either mating number or by chromosome number; Table S5, Table S6, and Table S7 tabulate the QTL analysis information generated in the combined sets of parental matings of MG 000, 00, and 0 in which the high protein accessions were not uniquely different from each other in terms of SNP genotype.

RESULTS

A total of 48 high seed protein soybean accessions were used as female parents in this research (Table 1). Additional high protein accessions have since been added to the germplasm collection, though the 48 used here remain representative of the current group of such accessions (Figure 1). The male parent accessions (*i.e.*, high yielding cultivars of ordinary seed protein content) have a seed oil content that is characteristically higher than that of most of the female parents. Accessions with a maturity greater than MG IV (except for one very early maturing MG V) were not used in this study because the normal fall frost date in Lincoln, NE precludes completion of their normal seed maturation.

Our initial goal was to generate at least 250 $F_{2,3}$ seed progenies in each mating, which was reached in most matings (Table S2), but not in some later MGs, though sufficient F_2 plant numbers were raised per mating. Progeny numbers averaged 224 over the 48 matings, but ranged from 278 to 115. In later MG matings, many F_2 plants produced too few F_3 seed (due to pod shattering) to meet the minimum seed sample requirement of the NIR instrument.

Phenotype data

The $F_{2,3}$ seed protein distributions (Table S2), only three of the 48 seed protein distributions had *P*-values for the Shapiro-Wilkes normality test that were less than the prechosen criterion of $P = 0.01$ (*i.e.*, mating 31, $P = 0.008$; 38, $P = 0.00005$; 47, $P = 0.003$), primarily because of a rightward skew (perceptibly slight in matings 31 and 47, but notably more so in 38). Seven other distributions had *P*-values of less than $P = 0.05$, but these seven were still greater than $P = 0.01$ (*i.e.*, matings 3, 6, 9, 10, 12, 21, and 37).

In the one replicate $F_{2,3}$ seed protein distributions, the minimum and maximum values among the 48 matings ranged from 371 to 402 g kg⁻¹, and from 446 to 497 g kg⁻¹ (Table S2). The $F_{2,3}$ progeny seed protein means in those 48 matings ranged from 411 to 439 g kg⁻¹.

Heritability

The seed protein phenotypic variance in the 48 matings ranged from 13 to 53, with a mean of 27 (Table S2)—typical magnitudes when protein content is NIR-measured using F_3 seed (*i.e.*, F_3 embryos with F_2 seed coats) produced by F_2 plants derived from matings of high protein parents with ordinary protein parents. The F_2 plant phenotypic variance, when divided by the summed parental plant phenotypic variances, led to moderately sized heritability estimates that averaged 66%, but ranged by mating from 30% to 87%. The seed oil phenotypic variance ranged from 12 to 58, with a mean of 26 (Table S2), and the heritability estimates (except for zero in mating 48) averaged 68%, and ranged from 18% to 93% in the other 47 matings. These estimates are based on just one (complete) replicate assay, one location, and one year, and thus do not have the accuracy of multi-environment-based heritability estimates (Visscher *et al.* 2008).

Population SNP genotyping numbers

The SG percentage of the 44 genotyped population individuals was actually a function of the number of phenotyped individuals which, in any given mating, deviated from a 48-mating average of $n = 224$. The SG two-tail percentage averaged 20.5% (Table S2), though that percentage by mating varied from 15.8% (*i.e.*, mating 16) to 38.3% (*i.e.*, mating 45).

A majority of the SNPs (*ca.* 60%) in the 1536-SNP chip developed by Hyten *et al.* (2010) were not bimorphic in each of the 48 matings (Table S3). The 48-mating average for parental SNP bimorphism was 29.3% of the 1536, but, on an individual mating basis, ranged from 16.9% (mating 31) to 36.5% (mating 40). On a chromosome basis, the range was 24% (chromosomes 6 and 7) to 38% (chromosome 16). In a few matings, some chromosomes had fewer than 10 bimorphic SNPs, primarily because of the removal, during error-checking, of several problematic SNPs that, when paired with other nearby SNPs, generated recombination fraction values far above the expected 0.50 maximum.

Version 4.0 of the soybean genetic map spans 2296.4 cM (Hyten *et al.* 2010), but, if restricted to just the 1536 SNPs, the map is shorter (*i.e.*, 2156.2 cM). A 5-cM SNP spacing is considered to be sufficiently dense for optimizing QTL detection power in populations of size 200 (Strange *et al.* 2013), implying that 440–460 evenly spaced SNPs would thus be adequate for a 2150–2300 cM map. The mean number of SNP markers segregating per population in this SG study was, in fact, 450

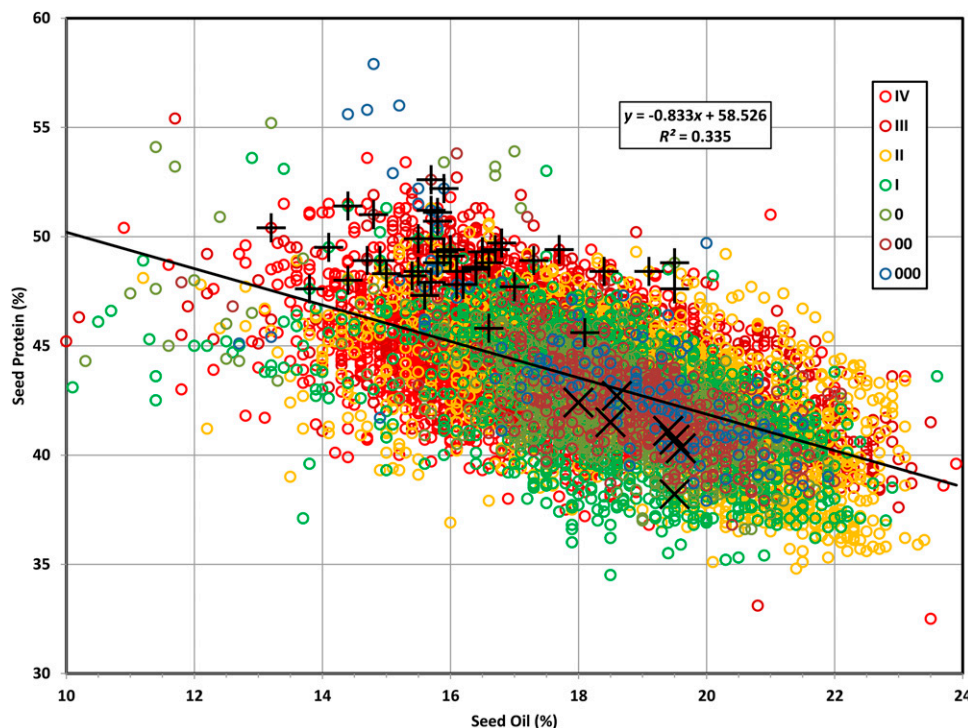


Figure 1 Seed protein values plotted against corresponding seed oil values. These are the GRIN values for 10,762 of the 17,711 *Glycine max* (L.) Merr. accessions in the USDA Soybean Germplasm Collection (as of December 31, 2015) in the seven maturity groups (MGs) of 000 (130), 00 (491), 0 (1179), I (1600), II (1831), III (1731), and IV (3800). Also shown are graph coordinates for 47 of the 48 high-protein female parents (+), and seven agronomic male parents (x) used in this study (see Table 1).

(Table S3), but ranged from a maximum of 560 (mating 40) to a minimum of 259 (mating 31). The number of genotyped SNPs was low in two other cases (317 in mating 7; 305 in mating 8), but 396 or more SNPs did segregate in 40 of the 48 matings, with 348 SNPs or more segregating in five of the remaining eight matings. The 1536-SNP chip was designed to position SNPs as uniformly possible over the chromosomes (Figure S2A), but less than *ca.* one-third of those SNPs segregated in any given mating. An example is mating 1, in which only 452 SNPs were bimorphic (Figure S2B). Marker monomorphism did result in SNP-coverage gaps of 30 cM or more in some matings (in the mating 1 example, chromosomes 1, 6, 7, 12, and 20), but marker gaps are not *a priori* predictable when using a SNP chip for genotyping in a multi-mating SG strategy.

QTL identified for seed protein and oil

The QTL analysis data obtained for soybean seed protein and seed oil in each of the 48 matings (Table S4) were translated into a heat map (Figure 2) to display the QTL peak LOD scores observed for seed protein (Figure 2A) or oil (Figure 2B) on any given soybean chromosome in each mating. The permutation-derived LOD score significance criterion (*i.e.*, genome-wide α of 0.05) for evaluating those observed QTL peak scores varied by population from 3.2 to 4.6 for protein, averaging *ca.* 4.0, and, for oil, varied from 3.2 to 5.6, also averaging *ca.* 4.0 (Table S4). Using the stratified permutation-based significance criterion, a QTL was detected on at least one chromosome for the SG trait of seed protein in 35 (73% of the 48) matings, and detected for the non-SG trait of seed oil in 25 (52% of the 48) matings (Figure 2, A and B, red-center bubbles). In two of the 48 matings (*i.e.*, 22 and 45), LOD score values on all 20 chromosomes were < 3.0 , indicating the absence of any protein or oil QTL.

The LOD score heat map makes evident the near-ubiquitous segregation of the well-known chromosome 20 QTL for protein and/or oil in many SG matings. The protein QTL was significant in 27 (77%) of the above-noted 35 matings (*i.e.*, 56% of all 48) (Figure 2A), with the oil

QTL being significant in 20 (80%) of the above-noted 25 matings (*i.e.*, 42% of all 48) (Figure 2B). Significant protein QTL were also detected on chromosome 10 (Figure 2A) in five matings (*i.e.*, 30 of MG II; 35, 37, 38, and 39 of MG III), but only in one mating (30), was a significant colocalized oil QTL detected (Figure 2B). The QTL region on chromosome 20 is known to be highly homologous with the long arm of chromosome 10 (Schmutz *et al.* 2010). Diers *et al.* (1992) reported that protein-oil QTL existed on chromosome 20 and 15. A protein and oil QTL was SG-detected on 15 in two MG 00 matings (13 and 17), but only for protein in MG IV mating 44 (Figure 2). Other less common SG-detected QTL were on chromosome 6 for protein (matings 18, 30, 33, and 38) and oil (30, 38, and 46), on 7 for both protein and oil (matings 2 and 34) and 18 (33), but just oil on 14 (27 and 46), and 18 (12, 33, and 42). Significant QTL were detected on chromosomes 2, 4, 12, 16, and 18 for protein, and on chromosomes 2, 8, 9 and 13 for oil, but only in single (separate) matings.

With respect to the significant seed protein QTL on chromosomes 20, 10, and 15, plus the protein QTL on chromosome 7 (matings 2 and 34), the QTL allele contributed by the high protein parent *enhanced* protein content, but coordinately *decreased* oil (Figure 2, A and B; +/- additive effects are denoted by a green/orange bubble color). Conversely, for the protein QTL repeatedly detected on chromosome 6 (matings 18, 30, 33, and 38), plus the protein QTL on chromosomes 2 and 18 (matings 31 and 18), the high protein parent allele *decreased* protein but *enhanced* oil.

For those significant protein and oil QTL that had coincident map positions, the protein and oil additive effects were directionally inverted (cf. Figure 2, A and B, and Table S4). Fewer oil QTL than protein QTL were detected, but this was expected, due to a SG focus only on protein, and a protein-oil correlation in the SG matings that, while strong, was clearly not unity, ranging from -0.66 to -0.88 , averaging -0.78 (Table S2). Chung *et al.* (2003) noted the chromosome 20 segment had opposite effects on protein and oil contents, perhaps due to pleiotropy. A single-QTL pleiotropy hypothesis is easily falsifiable upon detection of a

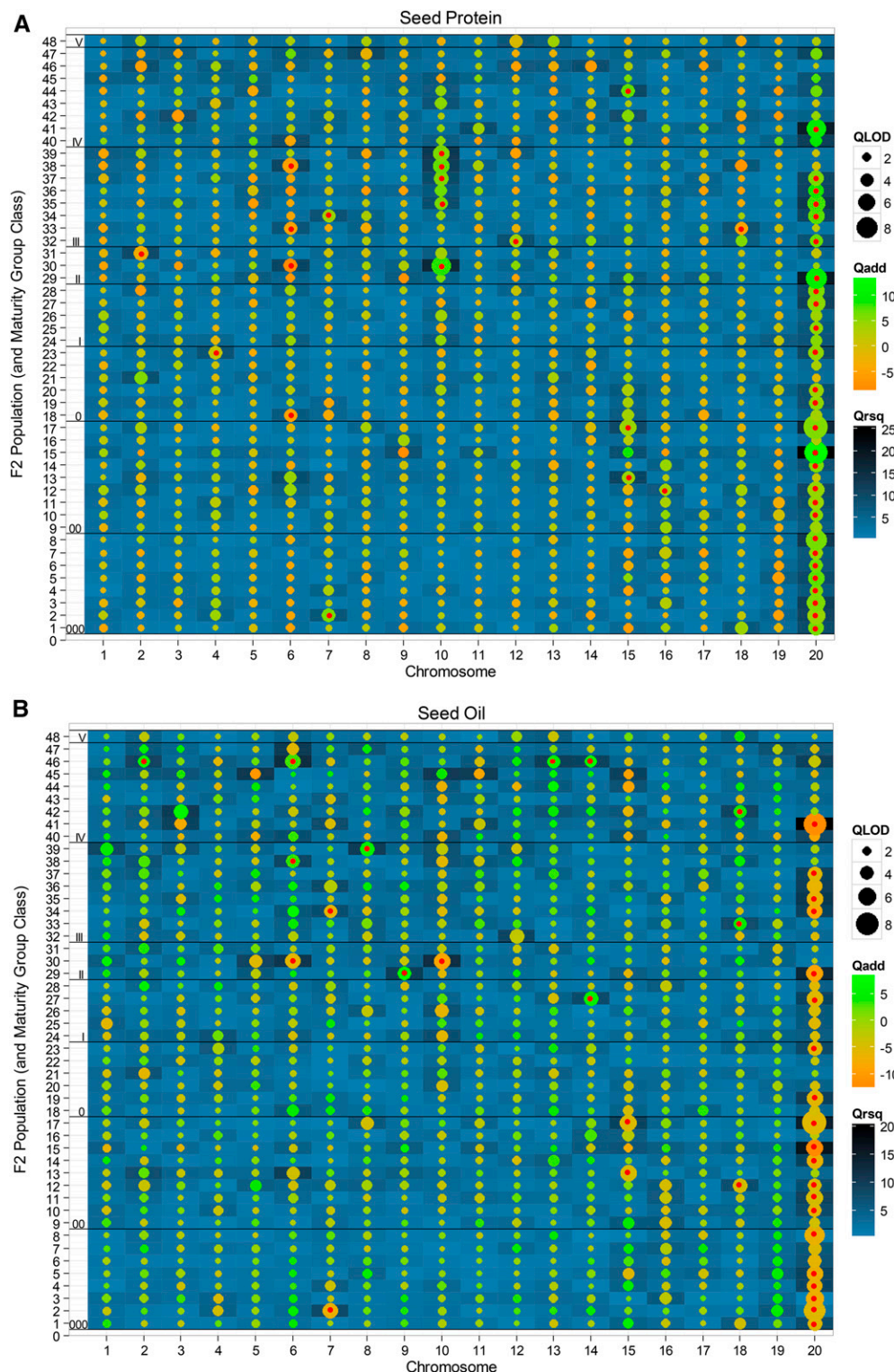


Figure 2 A heat map depicting parameter estimates for the SG-detected QTL for protein (A) and oil (B). The 48 MG-class matings are listed on the left axis, and 20 soybean chromosomes on the horizontal axis. The LOD score peak magnitudes are denoted by bubble size; those exceeding a genome-wide $\alpha = 0.05$ significance threshold derived from trait- and population-specific SG-stratified permutation tests ($n = 1900$) have red dot centers. Additive effect magnitude is denoted by bubble color intensity; green denoting a positive and orange a negative directional effect of the female parent B allele. The magnitude of the R^2 values is denoted by square tile color (light blue to deep black). See Table S4 for numerical values of QTL analysis parameters and permutation values.

recombinant with a coupling- (instead of a repulsion-) phased phenotype, but no recombinant individuals with a *high protein-high oil* seed content were detected in this study.

DISCUSSION

Selective genotyping was a term first defined and used by Lander and Botstein (1989), though the method had been essentially described

earlier by Lebowitz *et al.* (1987) as a “trait-based QTL analysis”, in which genotyping resources could be more efficiently allocated, with minimal loss of information, to just a fraction of progeny in a given mating. Indeed, Navabi *et al.* (2009) used simulation to document that, if 30 to 50 of 200 phenotyped progeny of a mating were genotyped in a bidirectional SG, Type I error would not exceed 0.02. With 20% genotyping, QTL detection power was still nearly 0.8 (*i.e.*, a Type II

error of 0.20), though detection of QTL of moderate to large effect size would require a marker spacing of at least 5 cM. These results led Navabi *et al.* (2009) to conclude that SG would be a very effective tool for screening large numbers of potential donors for large-effect QTL alleles governing a particular trait of interest.

That strategy was evaluated here by genotyping *ca.* 20% of *ca.* 224 phenotypes in each of the 48 F₂ populations created by using 48 high protein donor parents. We calculated, using the R program qtlDesign (Sen *et al.* 2007), that with a 5-cM SNP spacing, a Type I error (α) set to 0.05, and a Type II error (β) set to 0.2 to achieve a power ($1 - \beta$) of 0.8, QTL with an additive effect size of 5 g kg⁻¹ (accounting for *ca.* 15% of the phenotypic variance) could be detected in such populations. In our 48-mating SG study, for which the significance threshold (genome-wise Type I error of 0.05) in each population was obtained by permutation ($n = 1900$), significant QTL with additive effects of ≥ 4 g kg⁻¹, and R^2 values of 0.07 or more, were detected for protein on 10 chromosomes (*i.e.*, 1, 4, 6, 7, 10, 12, 15, 16, 18, and 20; Figure 2A), and for oil on 11 chromosomes (*i.e.*, 2, 6, 7, 8, 9, 10, 13, 14, 15, 18, and 20; Figure 2B, and see Table S4 for QTL summary data), confirming that multiple donor parents can be successfully surveyed for QTL presence using a SG strategy.

Seed protein and oil QTL detected in biparental matings in older publications are summarized in SoyBase (www.soybase.com). The QTL ANOVA F-statistics in old reports are not convertible into LOD scores, but the LOD scores in more recent reports are convertible into an F-statistic (Broman and Sen 2009), so we graphed the ANOVA F-statistic P -value (y -axis) and map position (x -axis) of each SoyBase-reported QTL (Figure S3). The evidence for a SoyBase-reported QTL using these comparison-wise P -values ranged from “merely suggestive” (*i.e.*, $P < 0.01 = 10^{-2}$)—a significance criterion leading to a naïve supposition that a SoyBase-listed QTL exists on every soybean chromosome (except 16 for protein), to “highly likely” (*i.e.*, $P < 0.0001 = 10^{-4}$)—a stringent significance criterion that offsets an intrinsic multiple marker comparison-wise test problem in the older reports. Using the latter criterion, we filtered the SoyBase-reported QTL to just the “most likely” protein QTL on the eight chromosomes of 1, 4, 6, 7, 11, 13, 15, and 20, and the oil QTL on the 10 chromosomes of 2, 5, 6, 8, 11, 14, 15, 16, 19, and 20 (Figure S3), wherein the underscores denote chromosomes in common with those having SG-significant protein or oil QTL (Figure 2). Comparatively, the SG study did reidentify some prior reported QTL; however, none of the 48 high protein SG donor parents were used in any of the 35 to 38 matings listed in SoyBase QTL reports, so this multi-mating SG strategy effectively doubled the number of biparental mapping populations used to date for detection of protein and oil QTL.

Korte and Farlow (2013) noted that GWAS surmounts two key limitations of biparental mapping: a QTL allele in a large-accession GWAS is not restricted to a 0.5 or zero frequency, as might be the case in any given biparental mating, and the QTL mapping resolution is greatly limited by the low number of potential recombination events in a F₂ or RIL population, even if the latter were to be increased to $n > 1000$ individuals to boost the number of recombinant events. Though GWAS does require marker-dense genotyping (*i.e.*, thousands of SNPs) to achieve its signal resolution potential, those SNP numbers are nowadays more easily obtainable in soybean, given the availability of a 50K SNP chip (Song *et al.* 2013), or using genotyping-by-sequencing to generate, *de novo*, several thousands of SNPs (Sonah *et al.* 2014).

The question then is whether a SG strategy is a worthy alternative to just using GWAS. Notably, biparental mapping and GWAS are still considered complementary approaches (Myles *et al.* 2009; Würschum 2012; Sonah *et al.* 2014). In fact, we considered our multi-mating SG

strategy, wherein *ca.* 450 SNPs were used to genotype just the highest 22 (10%) and lowest 22 (10%) protein phenotypes in *ca.* 224 progeny derived from 48 high protein \times low protein parental matings of MG 000 to IV to be contextually analogous to a phenotypic contrast type of GWAS [like the one recently conducted by Song *et al.* (2015) on soybean 100-seed weight]. The GWAS of Hwang *et al.* (2014) involved 31,954 SNP genotypes of 298 accessions of MG II, III, and IV, of which 151 had a high GRIN-based protein values, and 147 had a GRIN-based low protein values [though the contrasting GRIN values were only modestly ($r = 0.6$) correlated their field-based trial estimated values]. They detected significant QTL (using $-\log P \geq 3$) on *ca.* half of the chromosomes (Figure 3, A and B), and some of those QTL had chromosomal map positions coincident with some of our SG-detected significant QTL. Their two groups did include four SG high-protein female parents (32, 34, and 39 of MG III, plus 40 of MG IV), and two SG low-protein male parents (MGII Dwight and MG III Pana) (Table 1). Still, the comparative QTL results demonstrated that a SG survey strategy with *ca.* 48 accessions identified significant QTL with about the same degree of success achievable in a GWAS with *ca.* 300 accessions.

Recently, GWAS was used to detect seed protein and oil QTL (Sonah *et al.* 2014; Vaughn *et al.* 2014; Bandillo *et al.* 2015; Wen *et al.* 2015). The signal strength and map position of the significant QTL detected in our 48-mating SG can be compared to these GWAS-detected QTL (Figure 3). The $-\log P$ significance criterion / MG accession numbers / SNP numbers varied (*i.e.*, Sonah: 4/139 MG 0/17.2K SNPs; Vaughn: 4/619 MG I-II and 977 MG III-IV/~32K SNPs; Bandillo: 5.7/12K MG 000 – X/36.5K SNPs; Wen: 5/1.4K MG I-III/3.75K SNPs), but in all cases, a minor allele frequency (MAF) cut-off of 0.05 was used. Accession number maximization is often sought in GWAS, because doing so increases historical recombinant event numbers, thus enhancing statistical power, and QTL signal resolution. However, if only a few accessions possess an allele of notable effect at a given QTL, nondetection of that QTL will occur in GWAS if those few accessions comprise less than a 0.05 fraction of all of the evaluated accessions. In fact, the routine use of $MAF \geq 0.05$ in GWAS will, *a priori*, remove SNP locus alleles that are in complete linkage disequilibrium with a rare QTL allele that has an *in situ* frequency of < 0.05 . Bandillo *et al.* (2015) documented this by showing that the high protein–low oil allele (of large effect) at the well-known chromosome 20 QTL was present in just over 1% of the 12K accessions they examined. But, when they parsed the 12K accessions into smaller groups, based on seven countries of origin, or on eight MG classes, the high protein allele on chromosome 20 had an $MAF < 0.05$ in all but the Korean accession subset (and in all but the MG V to X subsets). Sonah *et al.* (2014) and Wen *et al.* (2015) did detect the chromosome 20 QTL allele in their respective sets of MG 0 and MG I–III accessions, but Vaughn *et al.* (2014) did not in two large sets of MG I–II or MG III–IV accessions (Figure 3). Myles *et al.* (2009) commented on the ineffectiveness of GWAS relative to the detection of rare alleles, and noted that controlled crosses and family-based mapping would be needed to artificially inflate the infrequency of rare functional alleles to improve the power needed for their detection, and to thus better understand the role that rare alleles play with regard to heritability of a given trait of interest.

Despite its rarity, the chromosome 20 QTL was obviously detected in many of the 48 MG 000 to IV donor parent accessions surveyed in this SG study (Figure 2). Accessions chosen for a SG-based QTL survey are actually quite likely to possess rare QTL alleles of a large-to-moderate effect in heritable traits, given the use of an “extreme” phenotype criterion to select SG donor accessions. If at least one chosen donor parent accession possesses a rare allele, its frequency will obviously be

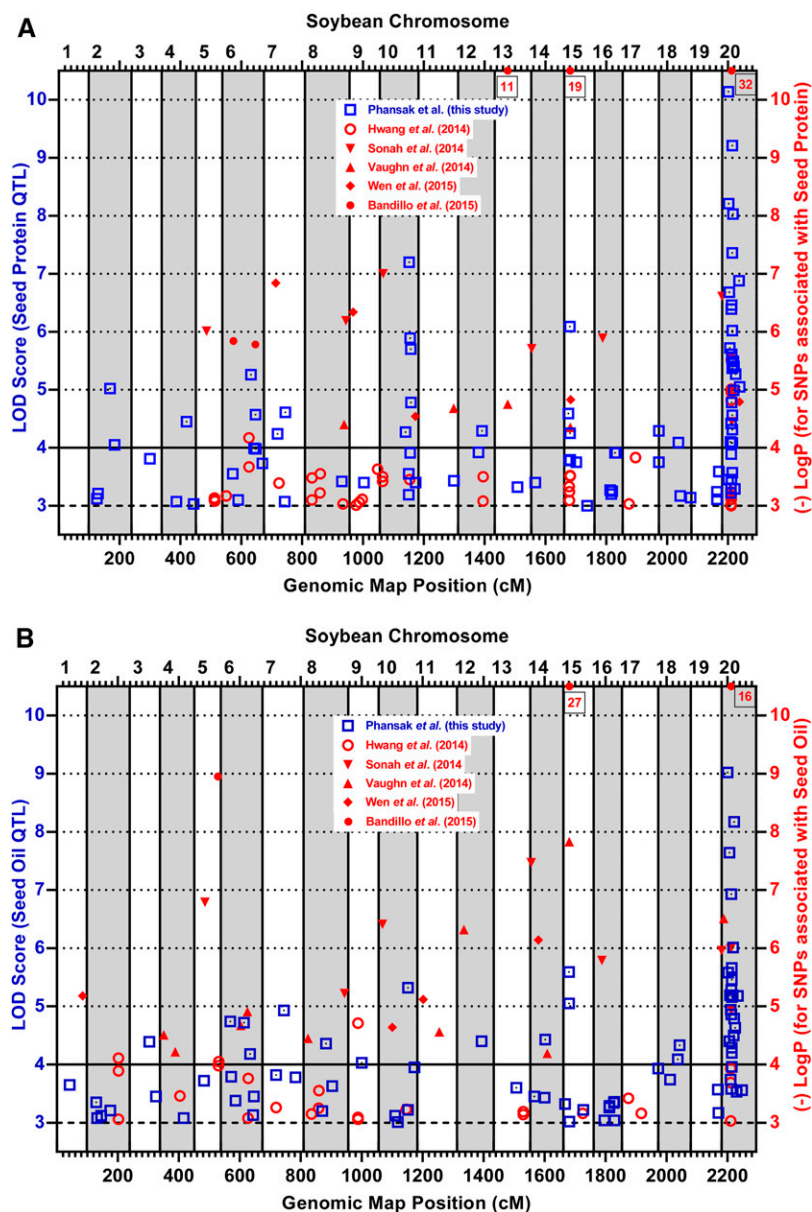


Figure 3 A graph of LOD score magnitudes of SG-detected QTL in 48 F_2 populations for seed protein (A) and oil (B). The bottom axis is scaled in terms of the Version 4.0 cumulative genetic map positions in the 20-chromosome soybean genome. The blue-box symbols with centered blue dots denote SG QTL exceeding a genome-wide $\alpha = 0.05$ significance threshold derived from trait- and population-specific SG-stratified permutation tests ($n = 1900$). Those thresholds varied from 3.6 to 4.6 for protein, and from 3.2 to 4.8 for oil, but averaged ca. 4.0 (horizontal black line). For comparative purposes, QTL detected in five recent GWAS publications are depicted relative to a $-\log P$ scaled right axis, though some Bandillo et al. (2015) values (box-enclosed at graph top) exceeded the scale limit.

0.5 in the progeny of the corresponding SG biparental high \times low mating, thus empowering its detection as noted by Myles et al. (2009).

A multi-mating F_2 population SG strategy can provide multiple estimates of the additive (and dominance) effects for the significant QTL detected in more than one mating. However, our population sizes were ca. 224 in size, and thus the effect estimates are likely overly optimistic, and selectively biased (Beavis 1998; Xu 2003; Broman and Sen 2009). For greater precision and accuracy in effect estimation, a fivefold (or greater) population size is needed, which, along with a more marker-dense SNP chip for genotyping (Strange et al. 2013), potentially mitigates SNP-to-SNP marker linkage map gaps. Despite that problem, our foremost objective in this SG study was evaluating an economical means for *per se* detection of significant protein and oil QTL in a large potential donor accession set. Using GWAS instead of SG offers no panacea for better estimation, given that effect estimates are always specific for the reference population used in either approach, as noted by Würschum (2012). Breeders must obviously conduct follow-up research to precisely estimate the QTL allele effect size in the genetic

backgrounds of their particular high-yielding cultivar sets, and to determine the worthiness of launching any marker-assisted high protein allele introgression program.

Song et al. (2015) found, after conducting a pairwise genetic similarity analysis using the 50K SNP chip, that 9% of the 18,480 accessions in the soybean germplasm collection had SNP genotypes that were not unique. They also reported that, using a 99.9% similarity criterion, 23% could be considered to be not unique. That discovery prompted us to review the 50K SNP genotypes of our 48 accessions. Unfortunately, eight of our MG 000 high protein accession parents (matings 1 to 8 in Table 1), and five of our MG 00 high protein parental accessions (matings 9–12 and 14) were not unique. Four MG 0 parental accessions (matings 18–20 and 22) also were not unique, but these four did differ from the former 13 accessions. Thus, only 33 of our 48 accessions were truly unique. Soybean breeders have used these MG 000, 00, and 0 accessions as a source of high protein alleles (Table 1), generally presuming that their differing GRIN passport data implied source diversity, but the germplasm SNP genotyping data reveals this presumption was mistaken.

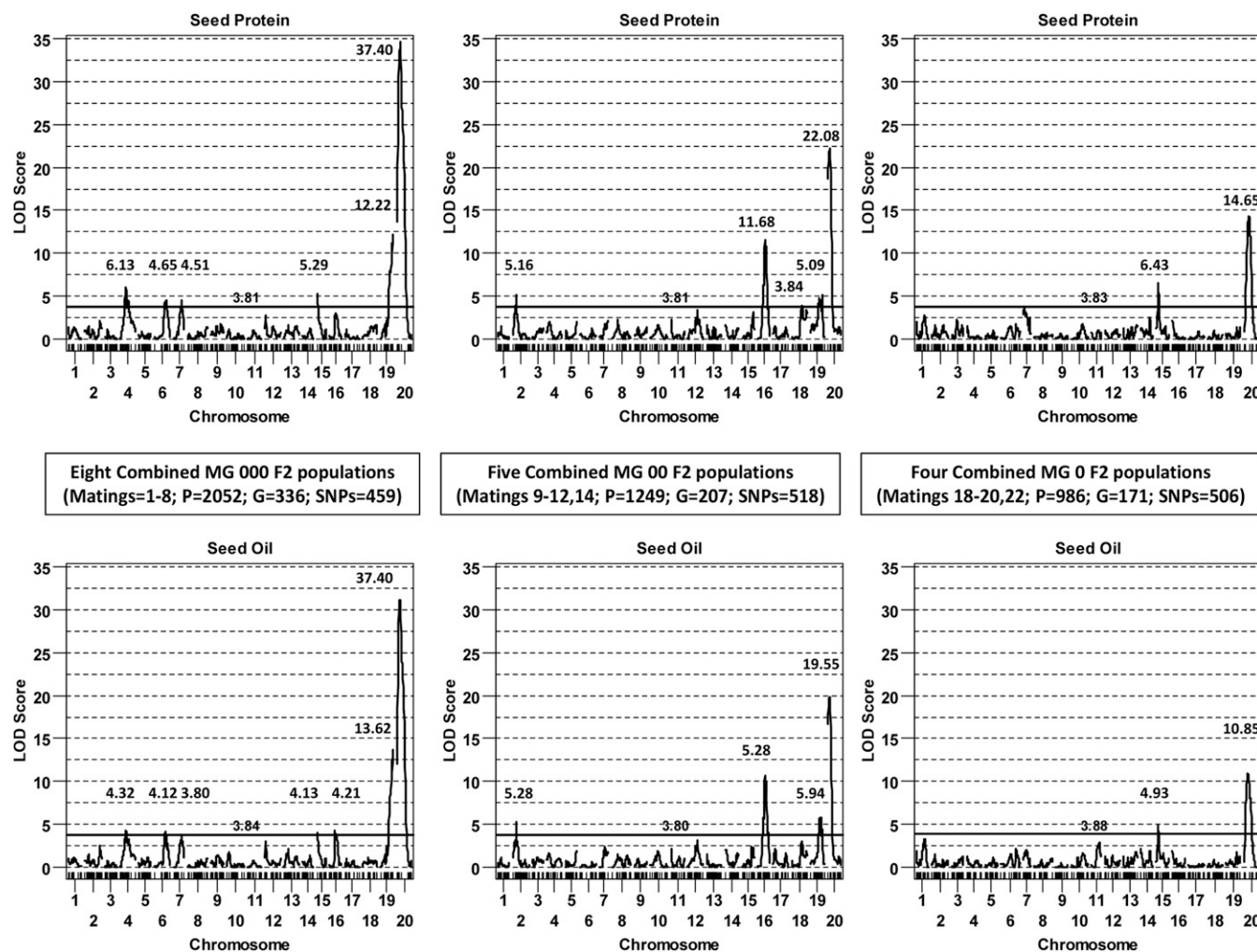


Figure 4 Chromosomal LOD score scans for protein (top panels) and oil (bottom panels). Selectively genotyped F₂ populations derived from parental matings in which the high protein accessions were not unique in terms of SNP genotype were pooled into three MG sets of 000 (left panels), 00 (middle panels), and 0 (right panels). The SG percentages were a respective 16.4, 16.8, and 17.3%, relative to the numbers of phenotypes (P), genotypes (G), and bimorphic SNPs shown for each MG set. Genome-wide $\alpha = 0.05$ significance thresholds, derived from trait- and population-specific SG-stratified permutation tests ($n = 1900$), were nearly-identical (*i.e.*, the 3.80 to 3.88 threshold horizontal lines shown in each panel).

The nonuniqueness of 13 MG 000 and 00 accessions, and the four MG 0 accessions, was a disappointing discovery, but it did offer a serendipitous opportunity to conduct a QTL analysis on three large biparental F₂ populations obtainable by pooling the *ca. n* = 224 F₂ populations of the three parental mating sets (*i.e.*, eight in MG 000, five in MG 00, and four in MG 0) based on the mating of those three sets of female parents to differing MG 000, 00, and 0 male parents. The pooled F₂ phenotype numbers were respectively 2052, 1249, and 986. Soybean populations of this size have not, to our knowledge, been reported for biparental QTL mapping studies, and thus could be used to obtain more precise estimates of QTL peak map positions and, because of the large population sizes, the inflationary impact of selection bias on allele effect estimates would be mitigated (Broman and Sen 2009).

The well-known major QTL located at the proximal end of chromosome 20 was detected in each pooled MG set (Figure 4), and also in the *ca. n* = 224 populations, except mating 9 in MG 00 and mating 18 and 22 in MG 0 (Figure 2; for details see Table S5, Table S6, and Table S7). In contrast, QTL were detected on chromosome 19 in the

MG 000 and 00 sets (Figure 4), but were not detected in any small population comprising those two sets (Figure 2). Similarly, the chromosome 4, 6, 7, and 15 QTL detected in MG 000 were not detected in small populations (except on chromosome 7 in mating 2). Finally, the QTL on chromosomes 2, 16, and 18 that were detected in MG 00 were not detected in small populations (except for mating 12 on 16 and 18 – oil only). These QTL likely had modest allelic effects that did not exceed the QTL detection limit in the small populations (*i.e.*, equivalent to false negatives), but did exceed it in the 5- to 10-fold larger populations.

By using chromosome-specific R/qtl additive and dominance effect scans (Figure 5), one can graphically view the impact of substituting a female parent B allele for the male parent A allele at each successive SNP on a chromosome. Coincident map positions were evident for most of the same-chromosome protein and oil QTL peaks, with such numbers being more concordant with a 1-locus pleiotropy than a 2-locus linkage model (Chung *et al.* 2003). The allele contributed by the high protein parents for the chromosomes 2, 4, 7, 16, and 20 QTL enhanced protein but decreased oil, whereas the allele contributed by

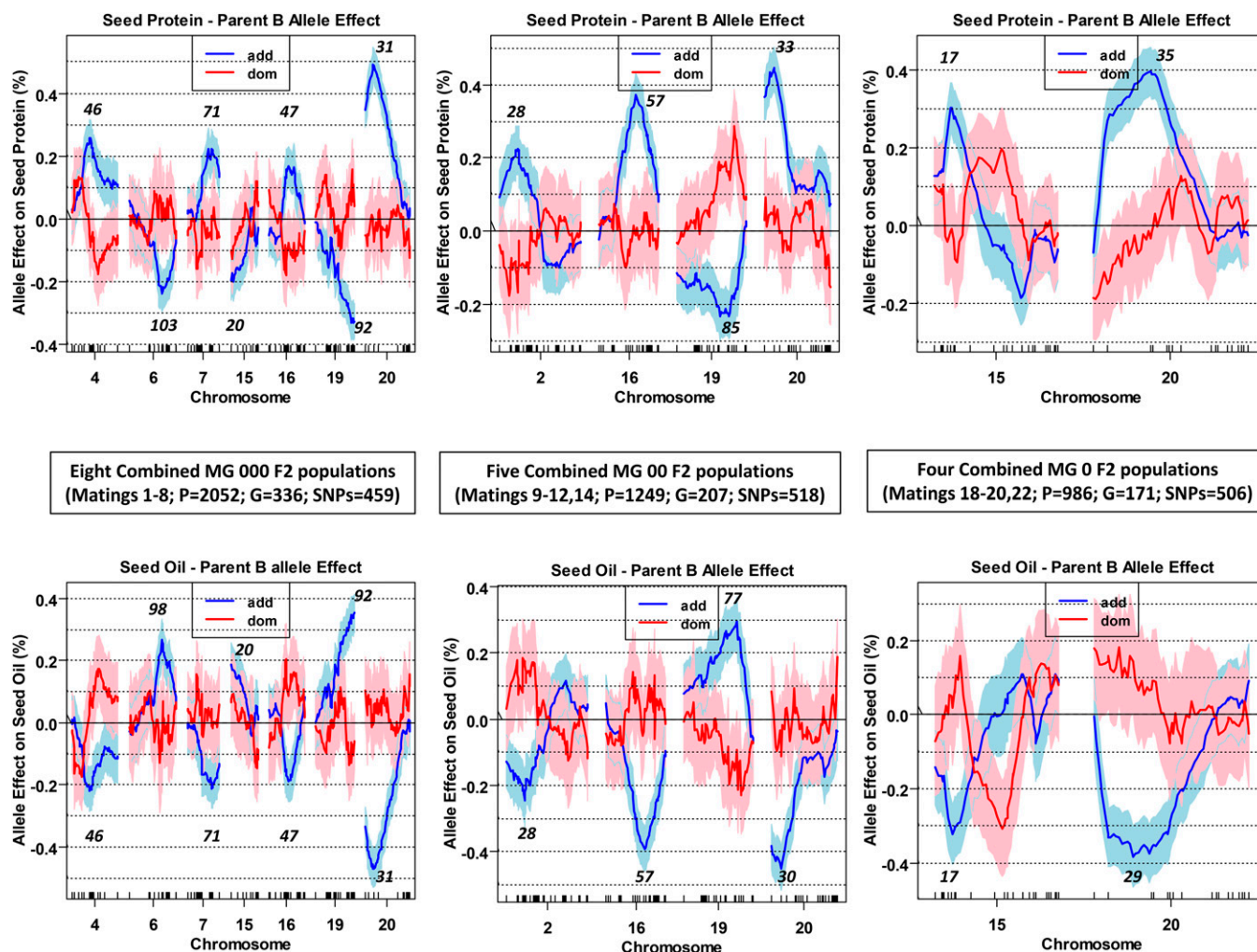


Figure 5 Chromosomal scans of the estimated additive and dominance effects for protein (top panels) and oil (bottom panels). The scans were limited to just the chromosomes exhibiting LOD score peaks shown in Figure 4. Effect magnitude and direction (+/–) reflect the substitution of a female parent B allele for a male parent A allele at any given bimorphic SNP position, with shading denoting the SE of the effect mean at each SNP. The chromosomal cM positions of the positive and negative additive effect maxima are italicized.

the same parents for the chromosome 6 and 19 QTL decreased protein but enhanced seed oil. One peculiarity in these scans was the differential additive effect scan patterns for the chromosome 15 QTL detected in MG 000 vs. MG 0 (Figure 5). The eight high protein female parents in the MG 000 set contributed an allele that decreased seed protein, whereas the four MG 0 high protein female parents contributed an allele that enhanced seed protein. The *trans*-phased phenotypic effect at these two linked QTL (*i.e.*, located at 20 cM in MG 000, but at 17 cM in MG 0) is notable, even if different male parents were used in these MG sets. We are not aware of any soybean linkage mapping study or GWAS documenting a similar *trans*-phased QTL pair.

The tracking of QTL dominance and additive effects in the three MG sets revealed that, at each chromosomal QTL (Figure 5, Table S5, Table S6, and Table S7), when the additive effect was positive, the dominance effect was typically (though not always) negative, and *vice versa*. However the SE boundary for the additive effect was narrower than that bounding the dominance effect—an indication that the latter was less precisely estimated, likely due to heterozygote infrequency in SG phenotypic extremes. Ordinarily, additive, plus additive × additive epistasis, accounts for most of the total trait genetic variance in soybean

(Burton 1987). Only inbred cultivars are used in commercial production, and the creation of F₁ hybrids is not likely anytime soon. Dominance effects would more likely be of breeder interest if made available for yield rather than seed protein and oil.

The discovery, or confirmatory rediscovery, of protein and oil QTL and map positions in this SG-based survey of high protein donor accessions will likely be of relevance to soybean breeders. The SG survey strategy did identify “major” protein and oil QTL in the 48 donor accessions examined here, suggesting that it could be used to detect major QTL alleles (and potentially rare ones) in traits other than seed composition, assuming that such traits can be reliably quantified using individual F₂ plants. A major drawback to the SG strategy is the need to apply it to phenotyping populations much larger than the $n = 224$ size used in this study, if the goal is to detect QTL of more modest additive effect.

ACKNOWLEDGMENTS

The authors thank Dr. Karl W. Broman (Department of Biostatistics and Medical Informatics, School of Medicine and Public Health, University of Wisconsin, Madison, WI) for providing much assistance and advice in the use of Rqtl for the QTL linkage mapping in this

multipopulation SG study. Piyaporn Phansak's Ph.D. program was funded by a scholarship from the Committee Staff Development Project on Higher Education, Royal Thai Government. Watcharin Soonsuwon's Ph.D. program was funded by a scholarship from the Faculty of Natural Resources, Prince of Songkla University, Thailand. Initial funding for the phenotyping phase of this research project was provided by the Nebraska Soybean Board. Subsequent funding for the genotyping phase was provided by the United Soybean Board.

LITERATURE CITED

- Ayoub, M., and D. E. Mather, 2002 Effectiveness of selective genotyping for detection of quantitative trait loci: an analysis of grain and malt quality traits in three barley populations. *Genome* 45: 1116–1124.
- Bandillo, N., D. Jarquin, Q. Song, R. Nelson, P. Cregan *et al.*, 2015 A population structure and genome-wide association analysis on the USDA soybean germplasm collection. *Plant Genome* 8: 1–13.
- Beavis, W. D., 1998 QTL analyses: power, precision, and accuracy, pp. 145–162 in *Molecular Dissection of Complex Traits*, edited by Paterson, A. H., CRC Press LLC, Boca Raton, FL.
- Bernardo, R., 2010 *Breeding for Quantitative Traits in Plants*, Ed. 2. Stemma Press, Woodbury, MN.
- Broman, K. W., and S. Sen, 2009 *A Guide to QTL Mapping with R/qtl*, Springer, New York.
- Brummer, E. C., G. L. Graef, J. Orf, J. R. Wilcox, and R. C. Shoemaker, 1997 Mapping QTL for seed protein and oil content in eight soybean populations. *Crop Sci.* 37: 370–378.
- Burton, J. W., 1987 Quantitative genetics: results relevant to soybean breeding, pp. 211–242 in *Soybeans: Improvement, Production, and Uses*, edited by J. R. Wilcox, ASA, CSSA, and SSSA, Madison, WI.
- Chung, J., H. L. Babka, G. L. Graef, P. E. Staswick, D. J. Lee *et al.*, 2003 The seed protein, oil, and yield QTL on soybean linkage group I. *Crop Sci.* 43: 1053–1067.
- Cober, E. R., and H. D. Voldeng, 2000 Developing high-protein, high-yield soybean populations and lines. *Crop Sci.* 40: 39–42.
- Darvasi, A., 1997 The effect of selective genotyping on QTL mapping accuracy. *Mamm. Genome* 8: 67–68.
- Darvasi, A., and M. Soller, 1992 Selective genotyping for determination of linkage between a marker locus and a quantitative trait locus. *Theor. Appl. Genet.* 85: 353–359.
- Diers, B. W., P. Keim, W. R. Fehr, and R. C. Shoemaker, 1992 RFLP analysis of soybean seed protein and oil content. *Theor. Appl. Genet.* 83: 608–612.
- Grant, D., R. T. Nelson, S. B. Cannon, and R. C. Shoemaker, 2010 Soybean, the USDA-ARS soybean genetics and genomics database. *Nucleic Acids Res.* 38: D843–D846.
- Hill, J. L., E. K. Peregrine, G. L. Sprau, C. R. Cemeens, R. L. Nelson *et al.*, 2008 Evaluation of the USDA soybean germplasm collection: Maturity groups 000–IV (PI 578371–PI 612761). USDA ARS Technical Bulletin 1919.
- Hymowitz, T., J. W. Dudley, F. I. Collins, and C. M. Brown, 1974 Estimations of protein and oil concentration in corn, soybean, and oat seed by near-infrared light reflectance. *Crop Sci.* 14: 713–715.
- Hyten, D. L., I.-Y. Choi, Q. Song, J. E. Specht, T. E. Carter *et al.*, 2010 A high density integrated genetic linkage map of soybean and the development of a 1536 Universal Soy Linkage Panel for quantitative trait locus mapping. *Crop Sci.* 50: 960–968.
- Hwang, E.-Y., Q. Song, G. Jia, J. E. Specht, D. L. Hyten *et al.*, 2014 A genome-wide association study of seed protein and oil content in soybean. *PLoS Genet.* 15: 1.
- Keim, P., B. W. Diers, T. C. Olson, and R. C. Shoemaker, 1990 RFLP mapping in soybean: association between marker loci and variation in quantitative traits. *Genetics* 126: 735–742.
- Korte, A., and A. Farlow, 2013 The advantages and limitations of trait analysis with GWAS: a review. *Plant Methods* 9: 29.
- Ladouceur, M., Z. Dastani, Y. S. Aulchenko, C. M. T. Greenwood, and J. B. Richards, 2012 The empirical power of rare variant association methods: results from sanger sequencing in 1998 individuals. *PLoS Genet.* 8: 1002496.
- Lander, E. S., and D. Botstein, 1989 Mapping Mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics* 121: 185–199.
- Lebowitz, R. J., M. Soller, and J. S. Beckmann, 1987 Trait-based analyses for the detection of linkage between marker loci and quantitative trait loci in crosses between inbred lines. *Theor. Appl. Genet.* 73: 556–562.
- Manichaikul, A., A. A. Abraham, S. Sen, and K. W. Broman, 2007 Significance thresholds for quantitative trait mapping under selective genotyping. *Genetics* 177: 1963–1966.
- Muranty, H. L. N., and B. Goffinet, 1997 Selective genotyping for location and estimation of the effect of a quantitative trait locus. *Biometrics* 53: 629–643.
- Myles, S., J. Peiffer, P. J. Brown, E. S. Ersoz, Z. Zhang *et al.*, 2009 Association mapping: critical considerations shift from genotyping to experimental design. *Plant Cell* 21: 2194–2202.
- Navabi, A., D. E. Mather, J. Bernier, D. M. Spanner, and G. N. Atlin, 2009 QTL detection with bidirectional and unidirectional selective genotyping: marker-based and trait-based analyses. *Theor. Appl. Genet.* 118: 347–358.
- Raychaudhuri, S., 2011 Mapping rare and common causal alleles for complex human diseases. *Cell* 147: 57–69.
- Rincker, K., R. Nelson, J. Specht, D. Sleper, T. Cary *et al.*, 2014 Genetic improvement of U.S. soybean in maturity groups II, III, and IV. *Crop Sci.* 54: 1–14.
- Schmutz, J., S. B. Cannon, J. Schlueter, J. Ma, T. Mitros *et al.*, 2010 Genome sequence of the palaeopolyploid soybean. *Nature* 463: 178–183.
- Semagn, K., A. Bjørnstad, and Y. Yu, 2010 The genetic dissection of quantitative traits in crops. *Electron. J. Biotechnol.* 13: 16–17.
- Sen, S., J. M. Satagopan, and G. A. Churchill, 2005 Quantitative trait locus study design from an information perspective. *Genetics* 170: 447–464.
- Sen, S., J. M. Satagopan, K. W. Broman, and G. A. Churchill, 2007 R/qtlDesign: inbred line cross experimental design. *Mamm. Genome* 18: 87–93.
- Sen, S., F. J. Johannes, and K. W. Broman, 2009 Selective genotyping and phenotyping strategies in a complex trait context. *Genetics* 181: 1613–1626.
- Sonah, S., L. O'Donoghue, E. Cober, I. Rajcan, and F. Belzile, 2014 Identification of loci governing eight agronomic traits using a GBS-GWAS approach and validation by QTL mapping in soya bean. *Plant Biotechnol. J.* 13: 211–221.
- Song, Q., D. L. Hyten, G. Jia, C. V. Quigley, E. W. Fickus *et al.*, 2013 Development and evaluation of SoySNP50K, a high-density genotyping array for soybean. *PLoS One* 8: 1–12.
- Song, Q., D. L. Hyten, G. Jia, C. V. Quigley, E. W. Fickus *et al.*, 2015 Fingerprinting soybean germplasm and its utility in genomic research. *G3 (Bethesda)* 5: 1999–2006.
- Strange, M., H. F. Utz, T. A. Scharag, A. E. Melchinger, and T. Würschum, 2013 High-density genotyping: and overkill for QTL mapping? Lessons learned from a case study in maize and simulations. *Theor. Appl. Genet.* 126: 2563–2574.
- Stuber, C. W., R. H. Moll, M. M. Goodman, H. E. Schaffer, and B. S. Weir, 1980 Allozyme frequency changes associated with selection for increased grain yield in maize (*Zea mays* L.). *Genetics* 95: 225–236.
- Stuber, C. W., M. M. Goodman, and R. H. Moll, 1982 Improvement of yield and ear number resulting from selection at allozyme loci in a maize population. *Crop Sci.* 22: 737–740.
- Sun, Y., J. Wang, J. H. Crouch, and Y. Xu, 2010 Efficiency of selective genotyping for genetic analysis of complex traits and potential application in crop improvement. *Mol. Breed.* 26: 493–511.
- Thornsberry, J. M., M. M. Goodman, J. Doebley, S. Kresovich, D. Nelson *et al.*, 2001 Dwarf8 polymorphisms associate with variation in flowering time. *Nat. Genet.* 28: 286–289.

- Vaughn, J. N., R. L. Nelson, Q. Song, P. B. Cregan, and Z. Li, 2014 The genetic architecture of seed composition in soybean is refined by genome-wide association scans across multiple populations. *G3* (Bethesda) 4: 2283–2294.
- Visscher, P. M., W. G. Wray, and N. R. Wray, 2008 Heritability in the genomics era—concepts and mis-conceptions. *Nat. Rev. Genet.* 9: 255–266.
- Voldeng, H. D., R. J. D. Guillemette, D. A. Leonard, and E. R. Cober, 1996 AC Proteus soybean. *Can. J. Plant Sci.* 76: 153–154.
- Wehrmann, V. K., W. R. Fehr, S. R. Cianzio, and J. F. Cavins, 1987 Transfer of high seed protein to high-yielding soybean cultivars. *Crop Sci.* 27: 927–931.
- Wen, Z., J. F. Boyse, Q. Song, P. B. Cregan, and D. Wang, 2015 Genomic consequences of selection and genome-wide association mapping in soybean. *BMC Genomics* 16: 671.
- Wilcox, J. R., 1998 Increasing seed protein in soybean with eight cycles of recurrent selection. *Crop Sci.* 38: 1536–1540.
- Würschum, T., 2012 Mapping QTL for agronomic traits in breeding populations. *Theor. Appl. Genet.* 125: 201–210.
- Xu, S., 2003 Theoretical basis of the Beavis effect. *Genetics* 165: 2259–2268.
- Xu, S., and C. Vogl, 2000 Maximum likelihood analysis of quantitative loci under selective genotyping. *Heredity* 84: 525–537.

Communicating editor: J. B. Holland

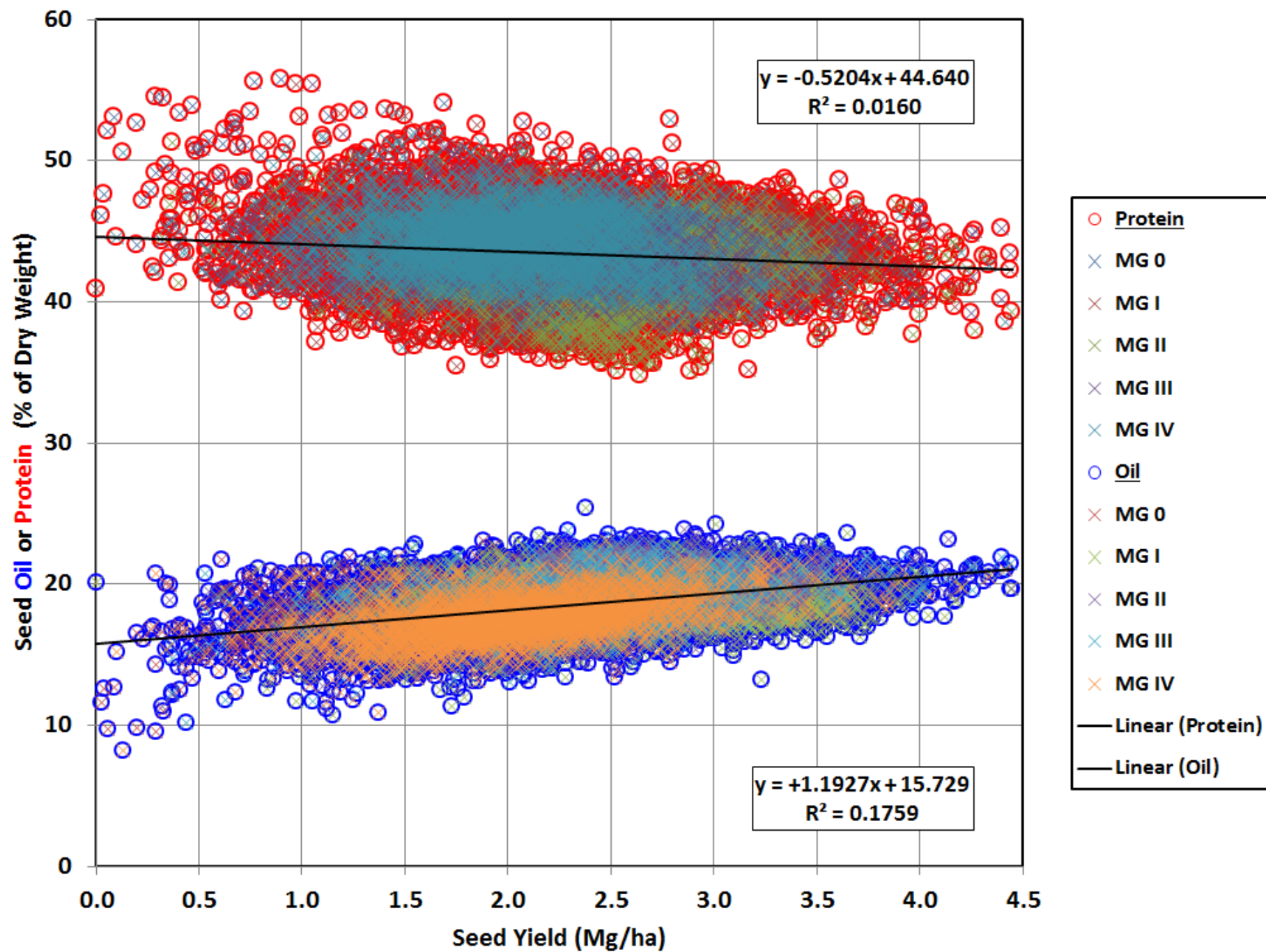
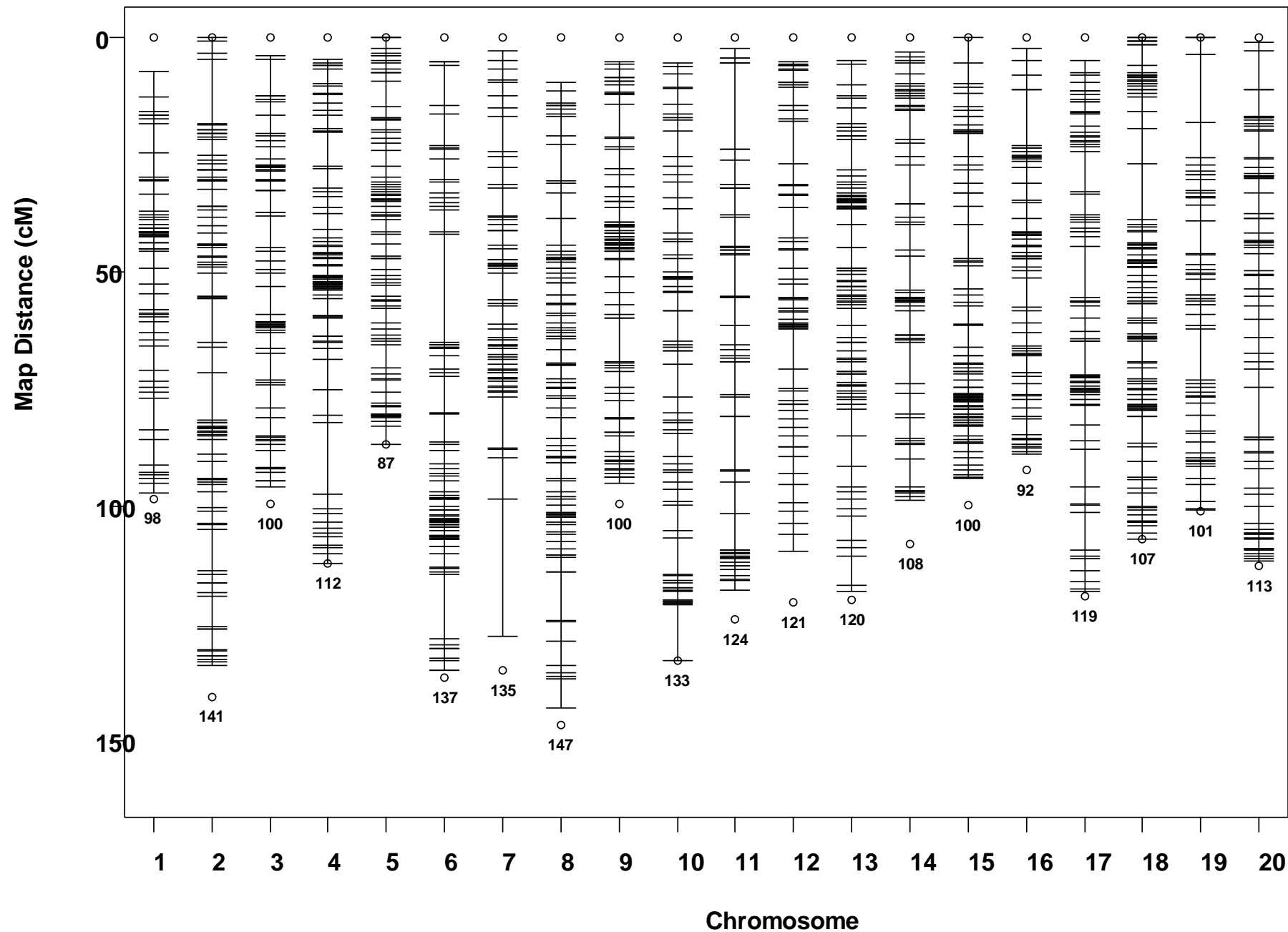


FIGURE S1. Soybean seed protein and oil values plotted against the corresponding seed yield values for maturity group (MG) 0 to IV accessions in the [*Glycine. max* (L.) Merr.] germplasm collection. Only 11,473 of the 12,141 accessions in these MGs have in-common values for all three traits. These data were provided courtesy of the soybean germplasm curator (R.L. Nelson, USDA-ARS, Urbana, IL).

A

B

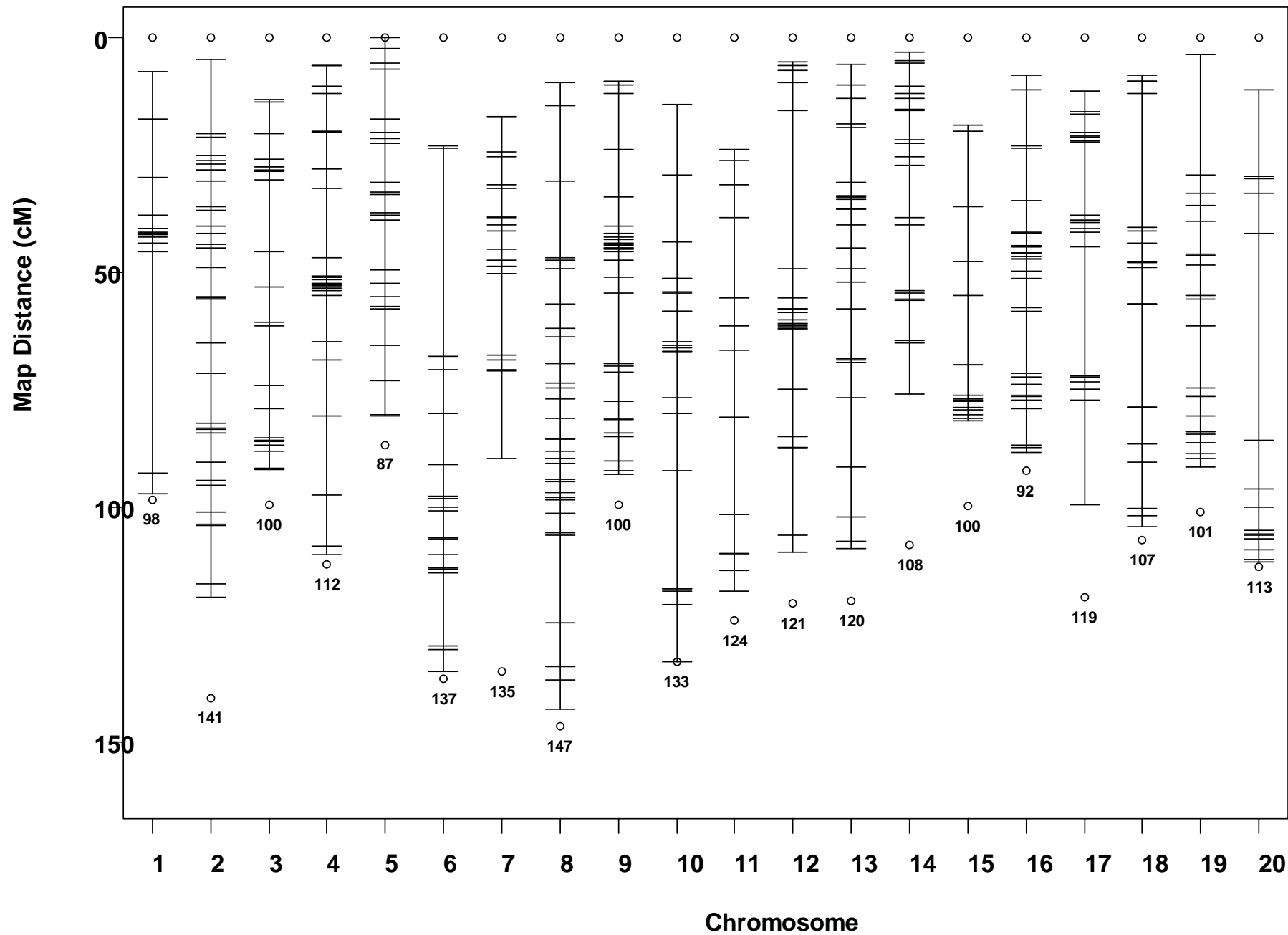
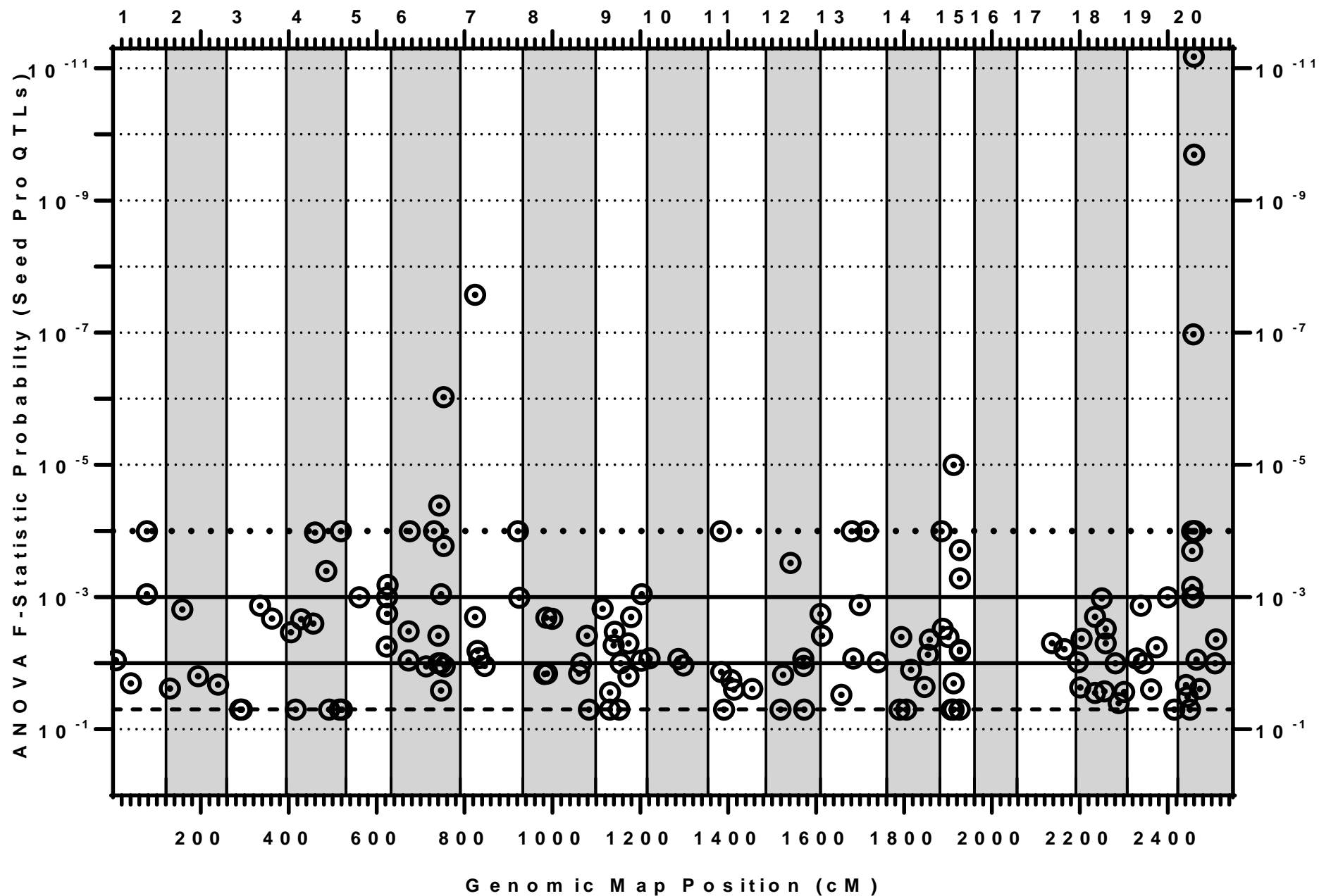


FIGURE S2. SNP marker map position (cross-hairs) in the 20 soybean chromosomes for (A) all of the 1536 SNPs in the chip developed by Hyten *et al.* (2010) and (B) just the 452 SNPs segregating in (the example) mating 1 of the 48 F₂ populations examined in this study. Open circles denote proximal and distal ends of each chromosome in the Version 4.0 genetic map (5500 markers of all types) that has a total genomic (Kosambi) distance of 2296.4 cM, but totals to only 2156.2 cM for a map that includes only the 1536 SNP markers.

A

Soybean Chromosome



B

Soybean Chromosome

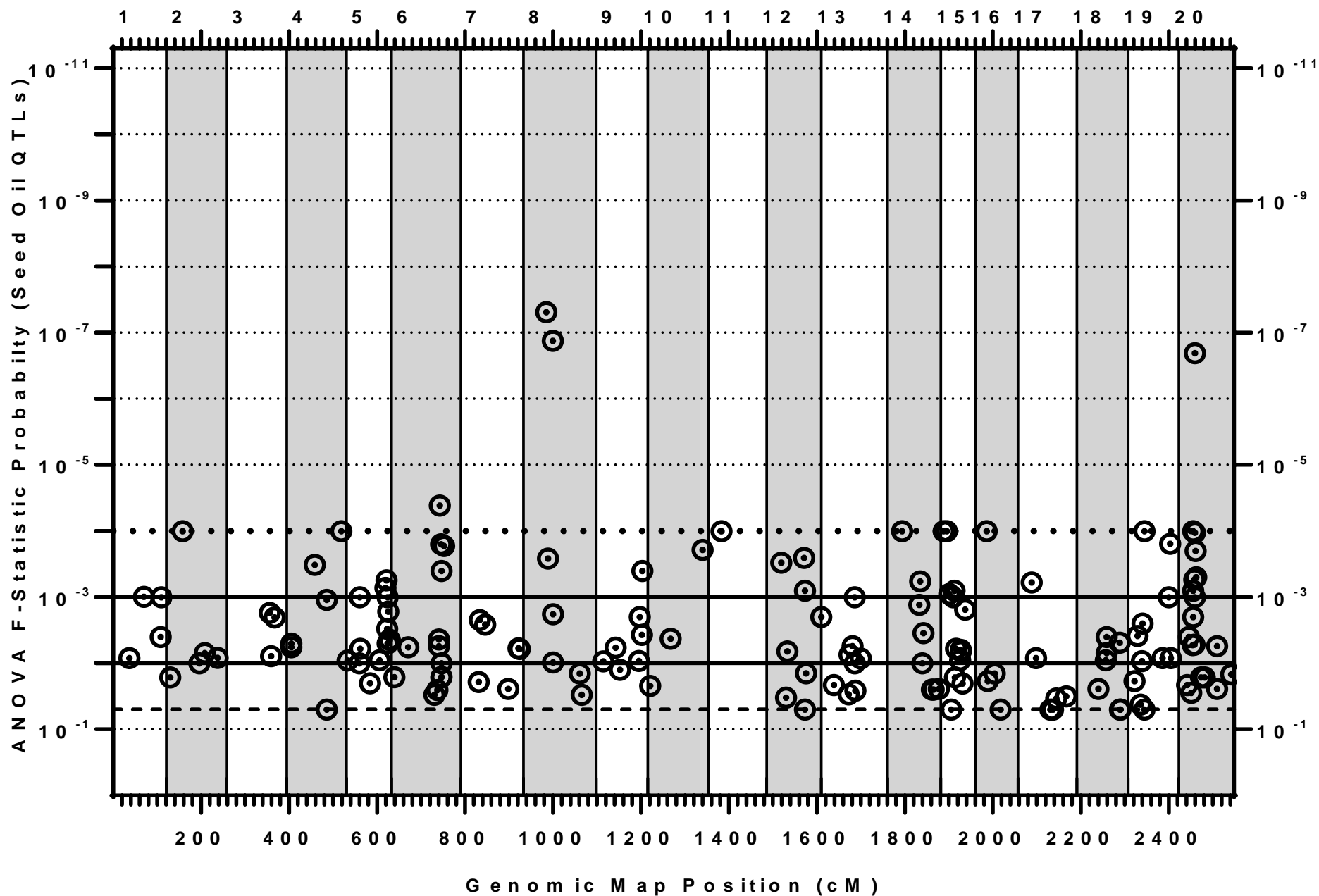


FIGURE S3. A graph of ANOVA F-statistic probabilities for (A) seed protein QTLs and (B) seed oil QTLs listed in SoyBase (Grant *et al.* 2010). SoyBase provides only Version 3.0 genetic map positions. ANOVA-based QTL detection was used in older reports, so for a common-unit comparative purpose, the LOD scores values provided in more recent reports were converted to F-statistic p-values. Because of the comparison-wise nature of F-statistic p-values, horizontal lines were placed at x-axis 10^{-x} values equivalent to p-values of 0.05 (dashed), 0.01 and 0.001 (solid), and 0.0001 (dotted) for reader convenience when comparing this graph with Figure 3 graph.

Multi-Population Selective Genotyping to Identify Soybean

(*Glycine max* (L.) Merr.) Seed Protein and Oil QTLs

Piyaporn Phansak^{*1}, Watcharin Soonsuwon^{*2}, David L. Hyten^{*}, Qijian Song[§],

Perry B. Cregan[§], George L. Graef^{*}, and James E. Specht^{*}

^{*} Department of Agronomy & Horticulture, University of Nebraska, Lincoln, Nebraska, 68583-0915, USA. [§] USDA-ARS, Soybean Genomics and Improvement Laboratory, Beltsville, MD, 20705-2325, USA. 1 Current Address: Division of Biology, Faculty of Science, Nakhon Phanom University, Muang, Nakhon Phanom, 48000 Thailand. 2 Current Address: Department of Plant Science, Faculty of Natural Resources, Prince of Songkla University, Hat Yai, Songkhla, 90112 Thailand.

Corresponding author e-mail: jspecht1@unl.edu

FILE S1

MATERIALS AND METHODS

Leaf Collection and DNA Extraction Procedures

Parents and F₁ Plants: One young trifoliate leaf was collected from each parental plant and from each F₁ plant (to be authenticated as to its hybridity) and placed into a deep 96-well collection plate that was stored at -20°C until subsequent DNA extraction. The DNA was extracted from the leaf tissue samples using a DNA extraction protocol described by (Saghai-Maroo *et al.* 1984), except that the volumes were halved. The DNA concentration of each sample was quantified using a spectrophotometer set to a wavelength of 260 and 280 nm. The stock DNA samples were diluted to 20 ng μL^{-1} for use in subsequent marker analyses.

F₂ Plants: A single newly developed trifoliate leaf with leaflets about 2 cm in length was excised from the main stem apex of each tagged F₂ plant in each of the 48 population-specific nursery rows. The three excised leaflets were gently rolled between the forefingers into a tight bundle that was inserted into a labeled well of a deep-well tissue collection plate. Plates were kept on ice during leaf tissue sampling until all wells were filled with tissue, and then immediately placed into a -20°C freezer. After the completion of phenotyping and the subsequent identification of F_{2:3} progenies in the lowest and highest decile fractions of the seed protein distribution in the fall of 2008, the leaf tissue samples of their F₂ plant progenitors were later retrieved from the collection plates and placed into 96-well extraction plates. Because of the cost of genotyping and the use of a 96-well plate format for genotyping, leaf samples of just the 22 F₂ plant progenitors of the seed progenies residing in the extreme decile seed protein phenotypic tails in each population were selected for genotyping. A 96-well extraction plate thus contained 47 leaf tissue samples of one given population (i.e., its 22 highest and 22 lowest protein F₂ plants, a confirmed F₁ plant, and the high and low protein parent plants), plus 47 leaf samples of the same type from another F₂ population. The two remaining wells were used for a repeat (insurance) sample of the high protein parent of each of the two populations. Plates were shipped on dry ice to a USDA laboratory in Beltsville, MD, for DNA extraction.

Genotypic Error Checking

The R package known as R/qtl (Broman *et al.* 2003; <http://www.rqtl.org/>) was used to error-check the phenotypic and genotypic data in each of the 48 populations, before that same software was used to conduct the QTL analyses. The phenotypic and genotypic data for each given population were organized into an Excel file in a comma-delimited “csvr” format (see Fig. 2.4, Broman and Sen, 2009). The 17-digit BARC number identifier of each of the 1536 SNPs was translated into a more compact S-prefixed 5-digit SNP ID# (Table S1), which minimized tabulation space requirements in the R/qtl output (and in the tables here).

The chromosome number and map position of each segregating SNP marker present in each *.csvr file were those published in the Version 4.0 soybean genetic map (Hyten *et al.* 2010).

A Chi-square goodness-of-fit test was conducted to identify allelic segregation distortion (R/qtl command: *geno.table*). Using a population-specific genome-wise alpha (generated by simply dividing a test-wise Type I alpha of 0.05 by the number of markers in each population), this test identified an average of 4-5 SNPs (range of 0 to 12) that were removed (*drop.markers*) from each population except mating 2 (Table 1), for which 57 such SNPs had to be removed.

A check for duplicate individuals (*comparegeno*) in each population identified a pair of F₂ individuals with 100% identical SNP genotypes in each of two populations (matings 11, 14), and a pair with 95% identity in another population (mating 10). These were likely the result of a harvest, threshing, or labeling error, so only the pair member with the greatest number of SNP marker genotypes was retained (*subset.cross*).

The crossover (XO) count for each genotyped F₂ plant was examined (*countXO*) to check for XO outlier counts that were abnormally high ($n > 175$) or abnormally low ($n < 16$) on a population-specific basis. On average, more than half of the 48 populations had no such outliers, and only 1 or 2 outliers were detected in the remaining populations (though six probable self progeny of the female parent were detected in mating 17). The few outliers were not retained (*subset.cross*).

Closely linked duplicate, triplicate, or quadruplicate SNP loci that did not recombine were identified as haplotype sets (*findDupMakers; exact.only*), and these di-, tri-, and quad sets of loci ranged from a respective 0 to 63, 0 to 18, and 0 to 4 in the 48 populations. This was not a surprising finding, given the limited opportunity for recombination events in these 44-genotype SG populations. In a single population QTL analysis, one ordinarily drops all but one member of a haplotype set. However, in this QTL analysis involving 48 populations, a slight difference (1×10^{-6} cM) in the map positions of the SNPs in each haplotype set was generated (*jittermap*) to retain these SNP markers for the purpose of inter-population comparisons of the markers near any QTL identified at a (near-) identical map position in at least two (or more) of the 48 populations

Allele codes for a given SNP marker can potentially be erroneously recorded in an inverted phase (i.e., A<>B). This condition may arise during the translation of base-pair SNP genotypic code into A-H-B marker genotypic code, usually when genotypic base-pair code is missing (or is heterozygous) for either (or both) inbred parent(s) for a given SNP marker. A range of 0 to 7 SNPs with such inverted allele code errors were detected (*checkAlleles*) in the 48 populations, though one population (mating 20) had 22 such potential errors. If the genotype counts (*geno.crostab*) for a marker pair involving the SNP with the potentially inverted alleles and either its left or right flanking marker provided inferentially sound evidence for doing so, then the erroneously phased allele code for that SNP was purposely re-phased (*switchAlleles*), but if not, then that SNP was dropped (*drop.markers*). The graphing of marker pair recombination fraction (rf) LOD scores against the corresponding rf values (*plot.rf*) also was used (1) to flag those marker pairs with a rf >0.85 (i.e., a criterion chosen based on the potential for an rf deviation of that magnitude from the theoretically expected rf=0.50 in a 44-genotype sample), and (2) to flag any other marker pairs with rf values of just less than 0.85 that had a rf LOD score >4.0. The majority of the marker pairs so flagged typically possessed the same common problematic marker member, and when that problematic marker was removed, the number of flagged pairs substantively decreased or went to zero.

Lincoln and Lander (1992) described an after-the-fact procedure for detecting apparent SNP genotyping errors, and because a variant of that *a posteriori* method is implemented in R/qtl (p. 381, Broman and Sen, 2009), it was used to evaluate the SNP genotype data in each of the 48 populations (*calc.errorlod*). The output list of potential genotyping errors (*toperrorlod; LOD>5*) contained, on average, about eight SNP markers per population (range of 0 to 24) that were present in multiple F₂ individuals (typically three to 20 or more), and these SNP markers were dropped (*drop.markers*). The remaining SNPs in the

outputted error list were not dropped; instead, in the one or two individuals identified with presumed SNP genotype errors, the SNP marker genotypes for those individuals were converted to missing (-) values.

The construction of an intrinsic linkage map for a typical SG population is problematic, not only because chromosomal marker order estimation is much less precise, but also because it is likely to be less accurate, due to the fact that genotype errors and missing SNP genotypes have a greater influence on marker order when the number of genotyped individuals is small (Martinez 1996). In this study, only 44 genotypes (or somewhat fewer, after error-checking) were available for marker order estimation in each population. For a QTL analysis, it is well known that marker order must be as correct as possible to avoid errors in the estimation of QTL map positions (p. 53, Broman and Sen 2009). Chromosomal marker positions can be reliably ascertained using a published genome sequence that has undergone multiple rounds of sequence error correction. Though very tightly linked SNP markers in the Hyten *et al.* (2010) Version 4.0 genetic map may still not be correctly ordered, the soybean breeding and genetics research community considers the linked markers in that map to be otherwise ordered to a sufficient degree of accuracy for mapping purposes. For that reason, we elected to use the Version 4.0 map of 1536 SNP loci as the chromosomal SNP marker order for each of our 48 populations. Population-specific recombination fractions were estimated (*est.rf*), as were the inter-marker map distances (*est.map*, *error.prob=0.001*, *map.function=kosambi*). The choice of 0.001 as an error probability parameter reflects an author conjecture of about one genotype error occurring per 1000 genotypes, thus implying that about 22 such errors would be expected per population amongst the approximate 22,000 genotypes generated when just 44 F₂ individuals per population are genotyped with *ca.* 500 SNPs.

The phenotype numbers and statistics after error-checking are presented in Table S2. A few F₂ progenitors had to be removed in some matings based on phenotypic errors, and, unfortunately, some of those that were dropped were members of the 44 selectively genotyped F₂ progenitors – six in three matings (9, 14, 46) and four in two matings (4, 23). In any event, the removal of a few F₂ plant phenotypes (genotyped or not) only slightly reduced the average SG percentage from 20.5 to 20%.

The genotype numbers and statistics after error-checking are presented in Table S3. The non-missing genotype percentage column in this table was calculated as the total number of SNP genotypes (AA, AB, BB) *minus* any missing ones (--), divided by total number of possible SNP genotypes (i.e., final SNP loci number multiplied by the final F₂ phenotype number). This parameter provided insight as to the degree to which the selectively genotyped individuals in each mating had missing genotypes at individual segregating SNP loci. The average over the 48 matings was 19%, which when compared to the 20% SG, indicates that SG individuals were not missing many SNP genotypes.

Relative to the collective AA: AB: BB genotype ratio (i.e., summed over all SNP loci segregating in a given mating), Chi-square tests (data not shown) revealed a good fit ($P>0.01$) of the observed ratio with the expected 1:2:1 F₂ genotype ratio in most matings (Table S3), though only marginally so in matings 6 ($P=0.038$), 41 ($P=0.029$), and 45 ($P=0.045$). However, a significant lack of fit ($P<0.001$) was evident in five matings, resulting from either an excess of AB heterozygotes (matings 2, 4, 8), or an excess of the male parent AA genotypes (matings 31, 32). Genes giving rise to gametophytic or zygotic differences in sterility, mortality, or vigor are typically the underlying causes of marker segregation distortion (SD). However, in these five populations, the SD was mostly genome-wide, as opposed to being limited to a specific localized chromosomal segment of SD that would arise if the SNPs therein were linked to a SD-causing genetic factor. Although SD can affect QTL detection power if the SD occurs in the closely linked flanking SNPs, Zhang *et al.* (2010) noted that the power can actually increase or decrease and, as long as the SD is not asymmetric to an extreme degree. They also noted that SD does not increase the number of false-positive QTLs, nor does it significantly impact the estimation of QTL position and effect.

For QTL detection with R/qtl, conditional genotypic probabilities were computed on a 2-cM grid basis (first with *calc.genoprob*, *step=2*, *errorprob = 0.001*, and thence with *sim.geno*, *step=2*, *n.draws=128*, *errorprob=0.001*), prior to conducting the interval mapping analysis (*scanone*, *method="em"*, *addcov=ac*, *n.perm=1900*, *perm.strata=strat*). The additive covariate parameter was not used in the analysis of the 48 individual populations, but was used in the analysis of three F₂ progeny sets, wherein each set was comprised of F₂ populations derived from SNP-identical parental matings in MG 000, 00, and 0. The parameter *ac* was a mating number assigned to each of the F₂ populations present in the given set.

Permutation (1900 replicates) was performed on just the selectively genotyped 44 F₂ plant phenotypes and their genotypes (*strat*) in each mating to obtain population-specific LOD score criterion based on a genome-wise $\alpha=0.05$. Permutation was similarly conducted in each of the three combined F₂ sets. The R/qtl output from *scanone* was used to compute a 95% Bayes credible interval for each detected chromosomal QTL in each F₂ population or each combined set (*bayesint*, *prob=0.95*). Numerical data for the protein and oil QTL peak parameters and permutation-based LOD threshold values that were generated with R/qtl are presented in Table S4, ordered by mating code and by chromosome number. Composite interval analysis (CIM) was not conducted with SG data generated in this study, because some experts express caution about using CIM when the missing genotype data is substantial (i.e., ca. 180 individuals were missing genotypes the ca. n=224 population sizes of the 48 matings. For one such expert view, see p. 206 of Broman and Sen (2009), and because SG requires a stratified permutation test (i.e., shuffling of phenotypes within just genotyped individual class separately from the non-genotyped class) for proper computation of a significance threshold LOD score (Manichaikul *et al.* 2007). The same argument applies for not conducting multiple qtl mapping with the SG data sets, given that multiple imputation can potentially result in spurious results when 80% of the genotypes are missing, i.e., p. 312 of Broman and Sen (2009). Attempts to apply the Blanc *et al.* (2009) connected population approach to our SG data sets were not successful, because of the substantial missing data.

LITERATURE CITED

- Blanc, G., A. Charcosset, B. Mangin, A. Gallais, and L. Moreau, 2006 Connected populations for detecting quantitative trait loci and testing for epistasis: and application in maize. *Theor. Appl. Genet.* 113: 206–224.
- Broman, K. W., and S. Sen, 2009 *A Guide to QTL Mapping with R/qtl*, Springer, New York.
- Broman, K. W., H. Wu, S. Sen, and G. A. Churchill, 2003 R/qtl: QTL mapping in experimental crosses. *Bioinformatics* 19: 889–890.
- Hyten, D. L., I.-Y. Choi, Q. Song, J. E. Specht, T. E. Carter et al., 2010 A high density integrated genetic linkage map of soybean and the development of a 1536 Universal Soy Linkage Panel for quantitative trait locus mapping. *Crop Sci.* 50: 960–968.
- Lincoln, S. E., and E. S. Lander, 1992 Systematic detection of errors in genetic linkage data. *Genomics* 14: 604–610.
- Manichaikul, A., A. A. Abraham, S. Sen, and K. W. Broman, 2007 Significance thresholds for quantitative trait mapping under selective genotyping. *Genetics* 177: 1963–1966.
- Martinez, O., 1996 Spurious linkage between markers in QTL mapping. *Theor. Appl. Genet.* 85: 480–488.
- Saghai-Maroo, M. A., K. M. Soliman, R. A. Jorgensen, and R. W. Allard, 1984 Ribosomal DNA spacer-length polymorphisms in barley: Mendelian inheritance, chromosomal location, and population dynamics. *Proc. Natl. Acad. Sci. USA* 81: 8014–8018.
- Zhang, L., S. Wang, H. Li, Q. Deng, A. Zheng et al., 2010 Effects of missing marker and segregation distortion on QTL mapping in F₂ populations. *Theor. Appl. Genet.* 121: 1071–1082.

TABLE S2. F₁, F₂, F_{1:2} and F_{2:3} plants, seeds, and progeny numbers in each of the 48 matings, ordered by MG then by mating code and ID number. Population distributional statistics are shown for the one replicate seed protein (and oil) assay (all progeny) in each mating. Only those progeny in the lowest and highest quintiles of the protein distribution, based on quintile cut-off values, were selected for a two replicate seed protein assay. Using the 2-replicate means, progenies with the 22 lowest and 22 highest seed protein were identified, and their respective 44 F₂ progenitors were selectively genotyped (SG). The final percentage of F₂ plants chosen for selective genotyping is shown in the last column.

			Female					One-Replicate Seed Protein					Seed Oil		Protein-Oil	Two-Replicate Mean Seed Protein Value						F ₂ plants
Mating			Parent	F ₁	F _{1:2}	F ₂	F _{2:3}	F _{2:3} Distributional Statistics					Statistics		Phenotype	Low Pro 22		Quintile Cut		High Pro 22		Selectively
No.	ID	MG	Accession	Plts	Sds	Plts	Prog	Min	MN	Max	Var	H ²	Var	H ²	Correlation	Min	Max	Low	High	Min	Max	Genotyped
				----- number -----				----- g kg ⁻¹ -----					----- g kg ⁻¹ -----								- % -	
1	1001	000	PI 153296	17	397	265	253	388	432	483	32	72	36	70	-0.85	399	414	437	448	465	488	17.4
2	1002	000	PI 189963	19	666	277	269	396	435	487	29	74	34	80	-0.86	398	418	437	442	467	491	16.4
3	1003	000	PI 548399	9	394	260	251	398	437	484	29	74	37	82	-0.85	400	423	440	441	466	489	17.5
4	1004	000	PI 372423	19	659	254	251	400	439	494	30	75	33	80	-0.85	401	422	438	446	468	506	17.5
5	1005	000	FC 30687	19	649	261	255	389	433	483	36	79	38	82	-0.88	392	411	432	440	466	483	17.3
6	1006	000	PI 153293	8	311	267	260	385	438	497	36	79	39	83	-0.88	395	416	443	437	469	505	16.9
7	1007	000	PI 372412	13	389	269	261	401	437	479	28	73	35	81	-0.87	403	419	434	441	465	483	16.9
8	1009	000	PI 548414	17	537	270	268	393	437	489	28	73	36	81	-0.88	402	418	436	444	466	486	16.4
9	1022	00	PI 153302	10	332	265	254	391	422	462	20	82	24	85	-0.76	395	405	419	419	446	462	17.3
10	1023	00	PI 159764	17	358	252	242	397	430	483	31	86	29	88	-0.82	399	415	436	432	456	481	18.2
11	1024	00	PI 438415	16	345	264	259	388	425	457	19	81	22	84	-0.79	399	408	430	419	445	458	17.0
12	1025	00	PI 153301	18	287	265	258	395	423	463	19	80	26	86	-0.79	395	408	428	433	448	467	17.1
13	1026	00	PI 189880	18	184	178	173	395	419	450	14	74	12	72	-0.78	398	408	418	427	435	451	25.4
14	1027	00	PI 153297	16	431	257	250	395	428	481	21	82	31	89	-0.78	400	411	422	423	448	480	17.6
15	2211	00	HHP	8	524	219	147	383	434	482	53	87	57	93	-0.81	387	417	430	450	456	480	29.9
16	2212	00	AC Proteus	9	405	286	278	392	418	450	13	45	15	60	-0.73	393	407	424	414	432	451	15.8
17	2213	00	AC Proteina	16	466	288	275	387	416	446	17	59	20	70	-0.80	390	403	434	414	441	458	16.0
18	1039	0	PI 427138	20	761	266	257	385	419	464	17	61	28	83	-0.75	390	405	424	429	441	462	17.1
19	1040	0	PI 261469	15	691	260	238	375	422	463	23	71	32	85	-0.80	380	407	422	426	448	463	18.5
20	1041	0	PI 181571	24	710	268	248	388	419	454	18	62	30	84	-0.71	394	407	422	424	443	460	17.7
21	1042	0	PI 424148	8	297	219	210	392	424	465	21	68	30	84	-0.76	396	412	427	419	447	462	21.0
22	1043	0	PI 423954	13	379	267	248	383	429	464	26	73	27	82	-0.73	363	408	428	428	445	461	17.7
23	1044	0	PI 154196	18	1117	271	245	384	419	454	21	67	28	83	-0.78	385	405	428	424	443	461	18.0

24	1054	I	PI 437088A	7	789	225	184	402	435	481	24	70	26	20	-0.80	406	423	432	429	454	487	23.9
25	1055	I	PI 423949	12	1014	286	265	386	428	476	24	69	42	50	-0.82	392	410	422	423	448	477	16.6
26	1056	I	PI 427141	11	906	285	258	392	431	471	27	73	39	46	-0.77	397	411	427	431	458	477	17.1
27	1057	I	PI 437716A	6	497	272	246	392	425	460	17	58	27	21	-0.67	396	413	437	433	446	468	17.9
28	1058	I	PI 423942	6	548	279	246	379	432	480	27	73	42	50	-0.76	375	411	431	428	455	482	17.9
29	1075	II	PI 423948A	2	238	208	183	372	430	475	38	72	43	66	-0.75	370	408	423	439	454	474	24.0
30	1076	II	PI 437112A	11	703	217	188	371	431	483	34	69	43	66	-0.76	375	410	425	439	451	479	23.4
31	1098	II	PI 548608	14	685	254	230	389	419	457	15	31	21	29	-0.66	389	405	419	427	435	456	19.1
32	1107	III	PI 445845	14	891	278	222	384	426	471	31	57	36	49	-0.78	385	407	433	421	450	473	19.8
33	1108	III	PI 398516	12	966	280	226	376	423	466	27	50	29	37	-0.74	374	403	419	432	443	464	19.5
34	1109	III	PI 91725-4	15	838	263	225	383	425	480	30	55	50	64	-0.76	387	404	426	439	445	486	19.6
35	1110	III	PI 340011	13	833	254	208	384	425	476	30	55	36	49	-0.81	385	405	421	417	444	474	21.2
36	1111	III	PI 243532	9	489	273	216	386	425	496	42	68	52	65	-0.74	386	405	429	440	449	497	20.4
37	1113	III	PI 408138C	12	1042	259	219	391	430	495	33	59	31	41	-0.78	390	409	429	428	450	498	20.1
38	1121	III	PI 398672	14	670	262	217	377	419	464	31	56	31	42	-0.73	380	400	425	429	443	472	20.3
39	1122	III	PI 360843	18	1544	213	189	373	419	459	19	30	22	18	-0.68	377	402	415	424	435	456	23.3
40	1138	IV	PI 253666A	9	735	242	141	376	431	494	42	76	58	73	-0.80	379	413	434	446	453	493	31.2
41	1139	IV	PI 407788A	17	1399	283	167	391	434	476	32	69	54	71	-0.83	391	412	420	447	455	484	26.3
42	1140	IV	PI 424286	16	554	264	189	389	429	467	27	63	39	60	-0.83	389	411	431	441	448	465	23.3
43	1142	IV	PI 407877B	10	936	278	217	397	429	473	21	52	33	54	-0.80	397	412	423	425	447	482	20.3
44	1143	IV	PI 398704	5	415	279	195	378	431	471	32	69	54	71	-0.83	380	411	422	446	455	475	22.6
45	1145	IV	PI 398970	1	85	165	115	381	428	478	38	74	48	67	-0.81	384	413	416	438	442	473	38.3
46	1146	IV	PI 407823	5	2190	254	171	397	424	467	19	48	40	62	-0.75	396	408	416	432	440	463	25.7
47	1152	IV	PI 407773B	2	500	260	145	394	425	465	26	61	40	62	-0.81	393	410	416	420	439	463	30.3
48	1183	V	PI 458256	9	870	288	231	377	411	451	15	34	15	0	-0.68	376	396	417	405	425	448	19.0

TABLE S3. The number of SNP markers detected as segregating in each soybean chromosome in each of the 48 F₂ populations *versus* the 1536 potentially detectable SNPs on the USLP 1.0 chip (Hyten *et al.* 2010). Some selectively genotyped F₂ individuals and some SNP markers were removed from the populations during the R/qtl error-checking phase of the project. Thus, final F₂ phenotype (Phe) and genotype (Gen) numbers, final selective genotyping (SG) percentage, final total non-missing genotype percentage, and final F₂ AA:AB:BB segregation ratios (summed over SNPs) are presented here for each mating. Relative to the last three F₂ SNP segregation columns, the high protein parent was arbitrarily assigned the BB genotype code at each SNP locus.

						Soybean Chromosome:																						All		SNP Marker		
							1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	Chr		Segregation			
Mating			Final F ₂ Data			Non	Number of SNP markers on the Soybean USLP 1.0 Chip																						Pop	Ratio (%)		
No.	ID	MG	Phe	Gen	SG	Miss	63	90	67	81	80	92	65	96	84	72	54	63	92	61	85	68	78	115	56	74	1536	SNP	AA	AB	BB	
			---- no. ----	%		%	Number of SNP markers detected as parentally polymorphic																						%	----- % -----		
1	1001	000	253	44	17.4	16.1	15	34	24	26	23	20	18	33	32	22	13	24	25	21	16	29	20	20	20	17	452	29.4	23.5	50.9	25.5	
2	1002	000	266	41	15.4	13.8	13	32	24	24	23	17	17	33	30	21	13	22	27	20	14	29	20	15	20	16	430	28.0	21.0	58.1	21.0	
3	1003	000	251	44	17.5	16.0	13	29	16	21	21	20	16	23	28	18	13	21	20	13	13	23	19	18	17	12	374	24.3	23.7	51.7	24.6	
4	1004	000	247	40	16.2	14.5	13	29	18	21	23	20	18	23	27	18	12	21	23	13	13	24	19	16	16	13	380	24.7	20.1	55.9	24.0	
5	1005	000	254	43	16.9	16.5	18	35	26	28	26	22	19	37	33	23	16	22	28	24	17	29	23	21	24	17	488	31.8	24.6	48.2	27.2	
6	1006	000	260	44	16.9	16.0	18	35	26	27	27	22	19	35	31	23	16	22	27	24	17	29	22	21	24	18	483	31.4	21.5	52.4	26.1	
7	1007	000	261	44	16.9	15.5	14	23	13	13	19	13	8	27	24	14	10	15	19	16	12	22	14	17	11	13	317	20.6	25.1	50.8	24.1	
8	1009	000	266	42	15.8	14.3	14	22	11	11	19	13	7	26	23	13	9	14	18	16	12	21	14	17	11	14	305	19.9	19.6	58.3	22.1	
9	1022	00	248	38	15.3	14.1	26	34	27	33	28	21	12	34	21	18	17	19	30	18	30	30	23	26	20	16	483	31.4	22.6	51.9	25.5	
10	1023	00	239	41	17.2	16.2	26	36	25	33	27	21	12	35	21	18	17	19	31	18	29	30	23	26	22	15	484	31.5	22.4	52.1	25.5	
11	1024	00	257	42	16.3	15.8	26	36	27	33	28	21	12	32	21	18	17	19	31	18	29	30	23	24	22	16	483	31.4	23.6	52.9	23.5	
12	1025	00	258	44	17.1	16.9	26	36	27	33	27	21	11	34	21	18	17	19	31	18	30	30	23	26	22	16	486	31.6	23.1	52.0	24.9	
13	1026	00	173	44	25.4	24.6	19	33	27	20	21	22	16	28	22	9	21	18	36	18	30	20	20	19	17	15	431	28.1	24.1	52.0	23.9	
14	1027	00	249	43	17.3	17.0	28	40	31	37	28	23	13	36	20	19	20	22	36	19	30	31	23	26	23	22	527	34.3	24.0	50.5	25.5	
15	2211	00	147	44	29.9	29.5	10	42	19	26	30	26	10	16	20	19	20	20	22	9	29	17	16	46	9	19	425	27.7	25.1	47.6	27.3	
16	2212	00	278	44	15.8	14.8	17	35	15	20	20	26	21	29	15	11	10	25	20	13	28	22	13	34	10	18	402	26.2	24.5	49.7	25.8	
17	2213	00	269	38	14.1	13.8	17	41	20	24	27	24	18	20	16	13	12	22	29	15	27	22	13	23	10	19	412	26.8	27.2	49.0	23.8	
18	1039	0	257	44	17.1	16.5	29	39	24	24	28	25	15	27	27	25	12	30	43	29	28	35	27	31	20	19	537	35.0	24.0	49.1	26.8	
19	1040	0	237	43	18.1	17.6	29	39	24	23	28	25	19	26	29	24	13	30	43	27	28	35	27	33	20	19	541	35.2	26.3	49.2	24.6	
20	1041	0	246	42	17.1	16.2	24	37	21	17	27	21	10	22	24	22	9	25	33	13	18	29	12	22	18	11	415	27.0	25.7	48.7	25.6	
21	1042	0	208	42	20.2	19.4	29	38	28	20	23	24	14	27	27	23	15	28	41	29	31	32	20	30	20	16	515	33.5	22.9	50.4	26.6	
22	1043	0	248	44	17.7	16.8	19	34	16	14	22	15	12	18	23	17	8	24	31	11	14	20	16	17	12	9	352	22.9	24.7	50.1	25.2	
23	1044	0	241	40	16.6	15.7	22	29	27	28	17	20	16	18	27	16	18	20	35	23	31	35	15	35	15	17	464	30.2	24.0	50.5	25.5	

24	1054	I	184	44	23.9	22.3	16	22	14	11	16	29	6	24	14	22	15	17	20	14	25	15	13	24	16	15	348	22.7	23.3	52.5	24.7
25	1055	I	264	43	16.3	15.6	20	28	22	33	20	22	13	24	20	26	17	27	27	21	28	25	20	21	18	18	450	29.3	24.0	52.4	23.6
26	1056	I	257	43	16.7	16.1	21	23	20	29	26	25	14	20	22	24	14	27	27	20	22	22	22	23	15	16	432	28.1	27.3	50.1	22.6
27	1057	I	246	44	17.9	17.1	12	27	19	19	27	18	13	14	27	17	16	11	20	13	26	26	17	30	9	18	379	24.7	24.2	52.5	23.3
28	1058	I	243	41	16.9	16.4	26	23	15	25	20	20	10	22	22	25	18	26	30	26	23	30	19	24	20	20	444	28.9	25.4	48.9	25.7
29	1075	II	181	42	23.2	22.7	24	23	15	21	28	25	14	19	20	18	14	23	19	18	28	28	12	32	21	17	419	27.3	25.3	49.6	25.1
30	1076	II	187	43	23.0	22.6	24	21	25	23	27	43	37	18	34	27	18	24	28	18	33	29	21	34	28	20	532	34.6	24.1	50.9	25.0
31	1098	II	229	43	18.8	18.0	11	19	26	7	18	6	3	5	10	17	7	8	12	13	16	12	16	26	8	19	259	16.9	30.4	46.0	23.6
32	1107	III	221	43	19.5	19.1	20	31	26	26	38	18	19	30	21	21	17	26	32	24	26	23	20	45	16	30	509	33.1	29.6	48.3	22.1
33	1108	III	226	44	19.5	19.1	25	19	22	24	24	27	16	23	20	17	24	28	35	19	35	28	18	42	20	12	478	31.1	25.7	51.2	23.1
34	1109	III	225	44	19.6	19.0	19	37	18	32	29	26	11	28	19	20	17	25	37	19	27	24	23	40	20	24	495	32.2	23.8	50.2	25.9
35	1110	III	208	44	21.2	21.0	18	40	21	27	38	17	22	36	17	25	20	27	36	22	30	27	21	41	18	30	533	34.7	25.8	49.5	24.7
36	1111	III	214	42	19.6	18.3	25	40	20	21	34	22	18	23	18	16	20	24	26	13	23	28	25	29	22	24	471	30.7	24.4	50.2	25.4
37	1113	III	219	44	20.1	19.9	18	37	19	28	28	13	18	34	16	17	18	24	37	28	26	30	20	45	19	28	503	32.7	24.9	51.0	24.1
38	1121	III	217	44	20.3	19.5	18	26	18	22	32	19	12	28	12	26	21	24	21	21	25	34	20	44	20	27	470	30.6	25.2	49.5	25.3
39	1122	III	189	44	23.3	22.5	25	25	15	16	20	14	15	24	19	22	14	14	21	13	24	29	17	26	20	24	397	25.8	23.3	48.7	28.0
40	1138	IV	141	44	31.2	30.7	28	34	19	20	34	31	23	38	38	21	24	28	24	13	40	19	31	40	18	37	560	36.5	26.0	49.9	24.1
41	1139	IV	166	43	25.9	24.7	17	41	24	23	30	22	22	25	24	19	12	20	34	16	21	19	26	47	19	25	486	31.6	24.3	47.1	28.6
42	1140	IV	188	43	22.9	22.8	14	33	24	20	26	30	16	33	27	24	18	30	24	8	22	23	25	39	23	25	484	31.5	25.5	47.3	27.1
43	1142	IV	217	44	20.3	20.0	16	42	23	19	23	30	20	24	16	19	12	26	40	14	29	27	25	46	18	30	499	32.5	25.9	49.8	24.3
44	1143	IV	194	43	22.2	21.3	12	38	22	18	24	28	22	30	20	17	21	23	28	12	25	24	25	46	28	36	499	32.5	25.3	48.2	26.5
45	1145	IV	115	44	38.3	35.8	10	43	20	19	28	22	22	26	21	15	18	25	32	9	19	20	19	33	19	27	447	29.1	24.0	53.8	22.3
46	1146	IV	165	38	23.0	22.1	17	28	17	22	22	24	23	20	13	14	13	20	24	14	21	20	13	36	21	14	396	25.8	26.5	49.8	23.7
47	1152	IV	144	43	29.9	29.7	15	44	15	19	15	30	23	17	14	22	17	23	17	12	23	24	19	39	24	16	428	27.9	24.1	49.4	26.5
48	1183	V	230	43	18.7	18.6	14	44	26	17	26	27	13	25	19	24	21	31	33	10	24	30	30	39	25	34	512	33.3	24.1	49.2	26.7