

## Genetic Diversity Patterns and Domestication Origin of Soybean

Soon-Chun Jeong<sup>1</sup>, Jung-Kyung Moon<sup>2</sup>, Soo-Kwon Park<sup>3</sup>, Myung-Shin Kim<sup>1</sup>, Kwanghee Lee<sup>1</sup>, Soo Rang Lee<sup>1</sup>, Namhee Jeong<sup>3</sup>, Man Soo Choi<sup>3</sup>, Namshin Kim<sup>4</sup>, Sung-Taeg Kang<sup>5</sup>, Euiho Park<sup>6</sup>

<sup>1</sup>Bio-Evaluation Center, Korea Research Institute of Bioscience and Biotechnology, Cheongju, Chungbuk 28116, Korea; <sup>2</sup>Agricultural Genome Center, National Academy of Agricultural Sciences, Rural Development Administration, Jeonju, Jeonbuk 55365, Korea; <sup>3</sup>National Institute of Crop Science, Rural Development Administration, Wanju, Jeonbuk 55365, Korea; <sup>4</sup>Epigenomics Research Center, Genome Institute, Korea Research Institute of Bioscience and Biotechnology, Taejon 34141, Korea; <sup>5</sup>Department of Crop Science and Biotechnology, Dankook University, Cheonan, Chungnam 31116, Korea; <sup>6</sup>School of Biotechnology, Yeungnam University, Gyeongsan, Gyeongbuk 38541, Korea

Corresponding authors:

Soon-Chun Jeong, Bio-Evaluation Center, Korea Research Institute of Bioscience and Biotechnology; (82) 43-240-6540; scjeong@kribb.re.kr; Jung-Kyung Moon, Agricultural Genome Center, National Academy of Agricultural Sciences; (82) 63-238-4763; moonjk2@korea.kr

Keywords: Genetic diversity, soybean, *Glycine max*, *Glycine soja*, domestication

Running title: Genetic Diversity and Origin of Soybean

1 **Abstract**

2 Understanding diversity and evolution of a crop is an essential step to implement a strategy to expand its  
3 germplasm base for crop improvement research. Samples intensively collected from Korea, which is a  
4 small but central region in the distribution geography of soybean, were genotyped to provide sufficient  
5 data to underpin genome-wide population genetic questions. After removing natural hybrids and duplicated  
6 or redundant accessions, we obtained a non-redundant set comprising 1,957 domesticated and 1,079 wild  
7 accessions to perform population structure analyses. Our analysis demonstrates that while wild soybean  
8 germplasm will require additional sampling from diverse indigenous areas to expand the germplasm base,  
9 the current domesticated soybean germplasm is saturated in terms of genetic diversity. We then showed  
10 that our genome-wide polymorphism map enabled us to detect genetic loci underling flower color, seed-  
11 coat color, and domestication syndrome. A representative soybean set consisting of 194 accessions were  
12 divided into one domesticated subpopulation and four wild subpopulations that could be traced back to  
13 their geographic collection areas. Population genomics analyses suggested that the monophyletic group of  
14 domesticated soybeans was originated in eastern Japan. The results were further substantiated by a  
15 phylogenetic tree constructed from domestication-associated single nucleotide polymorphisms identified in  
16 this study.

## 1 **1. Introduction**

2 To fully capitalize on the vast reservoir of favorable alleles that control agronomic traits within wild and  
3 domesticated germplasm, extensive phenotyping and genotyping of germplasm collections are necessary.  
4 Soybean [*Glycine max* (L.) Merr.] is a major crop for dietary protein and oil worldwide. Several hundred  
5 soybean genomes have been resequenced (Lam et al. 2010; Chung et al. 2014; Zhou et al. 2015; Valliyodan  
6 et al. 2016) and three genome-wide high-density SNP arrays have been developed and used to genotype  
7 thousands of soybean accessions (Song et al. 2013; Lee et al. 2015; Wang et al. 2016). These data have  
8 been primarily used to compare the patterns of genetic variation between *G. max* and its wild progenitor (*G.*  
9 *soja* Siebold & Zucc.) to understand the history of soybean domestication and identify selective sweeps  
10 related to the domestication and improvement of soybeans. The data have also been used to identify loci  
11 controlling important agronomic traits, such as protein-and-oil and seed-weight traits (Hwang et al. 2014;  
12 Bandillo et al. 2015; Zhou et al. 2015). However, those studies have been limited to detecting or  
13 confirming the major genetic loci reported in previous genetic mapping studies using biparental  
14 populations. Further efforts will be required to implement genome-wide association studies (GWAS)  
15 (McCarthy et al. 2008) with higher statistical power and mapping resolution in soybean.

16 Soybean was domesticated ~5000 years ago from *G. soja*, its sympatric wild annual progenitor that is  
17 distributed throughout East Asia, including most of China, Korea, Japan, and part of Russia (Hymowitz  
18 2004; Larson et al. 2014). Different regions of China have been proposed as a single center of soybean  
19 domestication on the basis of morphological, cytogenetic, and seed protein variation (Broich and Palmer  
20 1981; Hymowitz and Kaizuma 1981; Hymowitz 2004). Multiple centers of domestication including the  
21 southern areas of Japan and China have also been proposed based on chloroplast sequence variation (Xu et  
22 al. 2002) and archaeological records (Lee et al. 2011). However, recent phylogenetic studies using whole-  
23 genome resequencing data clearly indicated a monophyletic nature of domesticated soybean (Lam et al.  
24 2010; Chung et al. 2014; Zhou et al. 2015). Of late, molecular studies that used hundreds of markers and  
25 accessions have proposed different areas, such as the Yellow River of China (Li et al. 2010) and southern  
26 China (Guo et al. 2010), as a center of soybean domestication. Yet, two other studies suggested the  
27 domestication center as northern and central China using high-density SNP array data (Wang et al. 2016),

1 or central China surrounding the Yellow River using specific-locus amplified fragment sequencing data  
2 (Han et al. 2016).

3 In most of the previous studies, accessions collected from the Korean Peninsula were  
4 underrepresented, although this region is a central region of wild soybean distribution. For example, in the  
5 recent genome-wide analyses reported by Wang et al. (2016) and Han et al. (2016), accessions collected  
6 from China accounted for 91.8% and 100% of the total samples, respectively. Here, we present an analysis  
7 of SNP genotype data from 2,824 domesticated, 1,360 wild, and 50 putative hybrid accessions as part of an  
8 effort to characterize the entire Korean indigenous soybean collection deposited in the country's National  
9 Agrobiodiversity Center. We genotyped soybean accessions using the 180K Axiom<sup>®</sup> SoyaSNP array that  
10 was developed from soybean genome resequencing data (Lee et al. 2015). Our high-density SNP array data  
11 allowed us to evaluate levels of genetic diversity and patterns of population structure. We further  
12 attempted to detect genetic loci underlying soybean domestication and important agronomic traits, as well  
13 as provide a refined model of the evolutionary history of domesticated soybean.

14

## 15 **2. Materials and Methods**

### 16 **2.1. Plant materials and SNP genotyping**

17 The majority of the accessions were from the National Agrobiodiversity Center in Jeonju, Korea, with a  
18 small number of accessions provided by individual laboratories (Table S1). The National Agrobiodiversity  
19 Center collection consists of approximately 12,000 accessions of improved and landrace cultivars (*G. max*)  
20 and wild soybean (*G. soja*). The Korean germplasm collection substantially overlaps with those of other  
21 countries, particularly the United States. Most accessions collected from locations in other countries than  
22 Korea have been donated from the US National Genetic Resources Program. Notable exceptions were the  
23 46 wild accessions from Japan, whose accession codes start with 'B'. These were donated from National  
24 BioResource Project in Japan. As our primary goal was to characterize the indigenous soybean collection  
25 of Korea, we attempted as much as possible to genotype accessions unique to the Korean collection. At the  
26 same time, we analyzed representative sets of landrace accessions from China, North Korea, and Japan and  
27 approximately 400 improved lines, most of which are immediate descendants of ancestral lines of United

1 States soybean cultivars (Gizlice et al. 1994), so that cultivated soybean from Korea could be assessed in  
2 the context of worldwide soybean germplasm pool (Table S2). Representative *G. soja* accessions from  
3 China, Russia, and Japan were also selected, allowing the the geographic distribution of wild soybean in  
4 each of these countries to be sampled. Initially, we planted approximately 5,000 domesticated and 2,400  
5 wild soybean accessions, each of which contains approximately 90% of the accessions collected in Korea.  
6 After pure line selection by single seed descent was performed at least two times, DNA samples from  
7 approximately 4,400 diverse soybean germplasm lines were genotyped. However, because our SNP array  
8 data set ended up with smaller number of *G. soja* accessions from China than those from Korea and Japan,  
9 soybean population structure from the representative set was additionally assessed using genome  
10 resequencing data (downloaded from Figshare database,  
11 [http://figshare.com/articles/Soybean\\_resequencing\\_project/1176133](http://figshare.com/articles/Soybean_resequencing_project/1176133)) from 45 *G. soja* accessions reported  
12 by Zhou et al. (2015).

13 DNA samples from the ~ 4,400 diverse soybean accessions were extracted from a single plant of each  
14 accession and were genotyped with the Axiom<sup>®</sup> SoyaSNP array containing 180,961 SNP sites (Lee et al.  
15 2015). Of the lines genotyped, 4,234 with >97% sample call rate were selected for further analysis. SNPs  
16 were scored following the Axiom<sup>®</sup> Genotyping Solution Data Analysis User Guide  
17 (<http://www.affymetrix.com/>) as described by Lee et al. (2015). Of the 180,961 SNPs, 170,223 were  
18 selected on the basis of the development and validation study. Missing data points in the 170,223 SNPs  
19 were imputed using BEAGLE 4.0 with default settings (Browning and Browning 2007). The 170,223  
20 SNPs were then used to screen out duplicated and redundant accessions, leaving 3,036 non-redundant  
21 accessions. After the initial filtration, SNPs with heterozygous rate > 0.02 and minor allele frequency <  
22 0.02 were discarded from the genotype data of the non-redundant accessions, leaving a total of 117,095  
23 high quality SNPs for the further population analyses.

24 Phenotypic data used for GWAS were obtained primarily from field evaluations in the field at  
25 National Institute of Crop Science, Jeonju, Korea, in 2012 and 2013 (Table S1). The observed phenotype  
26 data were converted into binary data. The flower color phenotypes were divided into absence of color  
27 (white) or presence of colors ranging from light to dark purple. The seed coat color phenotypes were

1 divided into absence of colors (yellow or green) or presence of colors ranging from brown to black.  
2 Domestication phenotypes were divided into presence (*G. max*) or absence (*G. soja*) of domestication.

3

## 4 **2.2. Population structure and genetic diversity pattern analyses**

5 Principal component analysis (PCA) was conducted to summarize the genetic structure and variation  
6 present in the soybean collection using smartpca function in Eigensoft v7.2 (Patterson et al. 2006; Price et  
7 al. 2006). We plotted the first three PCs. NJ trees were constructed by MEGA7 (Kumar et al. 2016) under  
8 the  $p$ -distances model. We used a model-based clustering method implemented in ADMIXTURE v1.23  
9 (Alexander et al. 2009) to investigate the population structure of the soybean accessions. We determined  
10 the optimal  $K$ , the number of clusters based on the smallest cross-validation error calculated from  $v$ -fold  
11 cross-validation procedure. We plotted the membership coefficient using DISTRUCT (Rosenberg 2004).  
12 To investigate the level of genetic diversity maintained in soybean accessions, we calculated the nucleotide  
13 diversity ( $\pi$ ) using VCFtools v 0.1.13 (Danecek et al. 2011). Genetic differentiation (Weir and  
14 Cockerham's  $F_{ST}$  (Weir and Cockerham 1984)) between *G. max* and each of the *G. soja* subpopulations  
15 was calculated using the VCFtools V0.1.13. Hierarchical analyses of molecular variance (AMOVA)  
16 in the whole soybean set and the representative soybean set were performed using  
17 ARLEQUIN v.3.5.2.2 (Excoffier and Lischer 2010). The significance of the values for  $F_{CT}$   
18 (difference among groups),  $F_{SC}$  (difference among populations within groups), and  $F_{ST}$   
19 (difference among populations) was tested by 1023 permutations.

20

## 21 **2.3. Genome-wide association studies**

22 We conducted GWAS using PLINK 1.9 (Purcell et al. 2007). Flower color, seed-coat color, and  
23 domestication phenotype data were converted to binary data to perform a conditional logistic regression  
24 model analysis. Conditional specific SNPs were selected on the basis of minor allele frequencies of SNPs  
25 among the groups defined based on PCA analysis. A Bonferroni correction was used to control for the  
26 multiple testing problem by adjusting the alpha value from  $\alpha = 0.01$  to  $\alpha = (0.01/117,095 \text{ SNPs})$  where

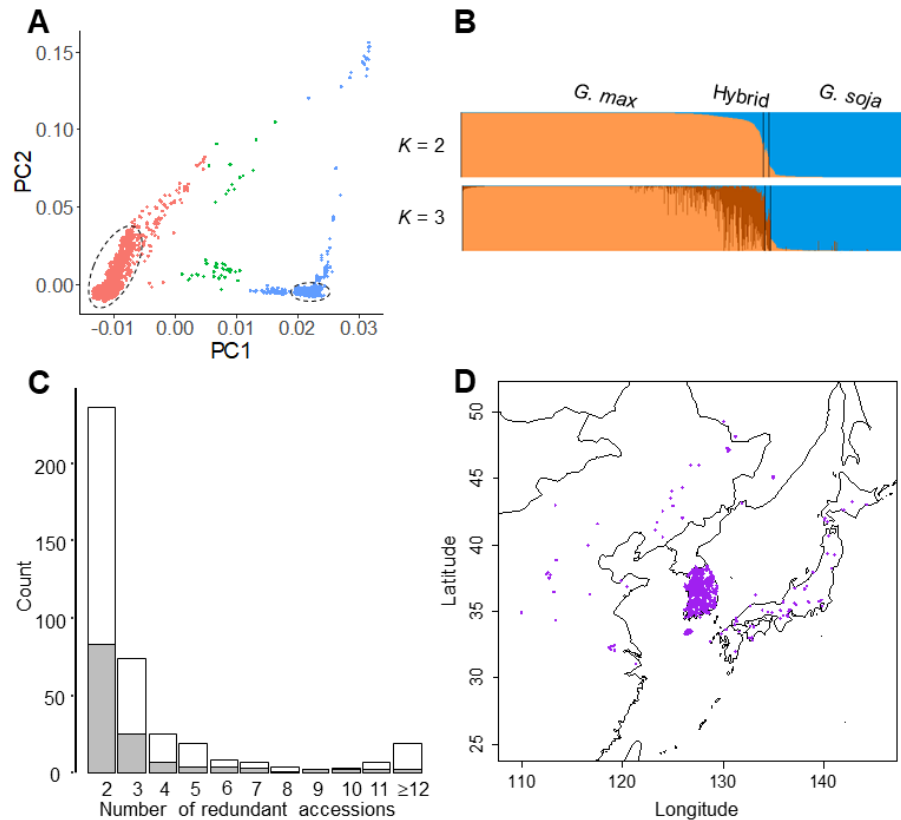
1 117,095 is the number of statistical tests conducted. Therefore, statistical significance of a SNP-trait  
2 association was set at  $8.54e^{-8}$  ( $-\log_{10} P = 7.06$ ). Manhattan plots were produced using the qqman package  
3 (Turner 2014). To define linkage disequilibrium (LD) patterns, correlation coefficient of alleles ( $r^2$ ) were  
4 calculated for SNPs under the peak regions that exhibited significant association using Haploview 4.2  
5 (Barrett et al. 2005). The confidence interval (CI) method of Gabriel et al. (2002) was used to identify LD  
6 blocks.

7

### 8 **3. Results**

#### 9 **3.1. Overall structure of the genotyped soybean germplasm population**

10 Of the approximately 4,400 accessions genotyped, 4,234 exhibited >97% sample call rate. These were used  
11 as the total population for characterizing the Korean soybean germplasm (Table S1). The majority of this  
12 4,234 set contained accessions from Korea (78.7% *G. max* and 91.5% *G. soja*) (Fig. 1A). The rest were *G.*  
13 *max* landrace and *G. soja* accessions from China, Russia, North Korea, and Japan and improved lines  
14 mostly developed in the United States (Table S2). To eliminate potential confounding effects exerted by  
15 hybrids in the comparison of wild and domesticated soybean populations (Vaughan et al. 2008; Wang et al.  
16 2017), we first removed 50 putative hybrid accessions from among the 4,234 accessions (Fig. 1B; Table  
17 S2). In the field evaluation, most of these 50 accessions showed intermediate morphologies between  
18 domesticated soybean (*G. max*) and its wild relative (*G. soja*). Furthermore, principal component analysis  
19 (PCA) using 170,223 high-quality SNPs showed that the accessions were positioned between two large  
20 groups of *G. max* and *G. soja* (Fig. 1A). In further support of their suspected hybrid status, 50 accessions  
21 showed mixed wild or domesticated genome fractions ranging from 30 to 70% in  $K = 2$  or 3 populations in  
22 the ADMIXTURE analysis (Fig. 1B).



1

2 **Figure 1.** Population structure of the genotyped 4,234 soybean accessions. (A) Principal components of  
3 SNP variation. PC1 and PC2 indicate score of principal components 1 and 2, respectively. Each of PC1 and  
4 PC2 explained 15.6% and 2.7% of variance in the data. *Glycine max*, *G. soja*, and hybrids are shown by  
5 red, blue, and green dots, respectively. The majority of Korean accessions cluster together within dashed  
6 ellipses. (B) ADMIXTURE plots. The accessions were divided into three groups: *G. max*, *G. soja*, and  
7 their hybrids. (C) Distribution of number of redundant accession groups that showed <1.25%  
8 inconsistencies between the SNP calls. *G. max* and *G. soja* are shown by white and gray boxes. (D)  
9 Geographic distribution of the collection sites for *G. soja* accessions.

10

11 In our previous development and validation study (Lee et al. 2015), the SNP calls genotyped by the  
12 SoyaSNP array were highly reproducible, with inconsistencies of  $\leq 1.17\%$  observed within pairs of 27  
13 duplicated samples after excluding missing genotypes in either sample. Several sets of near-isogenic



1 isolines were genotyped (Table S3). Single-gene isolines (backcross-derived isolines for single genes)  
2 showed approximately 1.0% inconsistencies after excluding missing genotypes in either sample (e.g., 1.16%  
3 between Harosoy and L67-153 [Harosoy(6) x Higan]). As expected, a slightly higher level of  
4 inconsistency (up to 1.5%) was observed for samples from multiple-gene isolines (e.g. 1.48% between  
5 L62-667 [Harosoy(6) x T204] and OT94-51 [OT89-5/L71-802//OT89-6]). However, we occasionally  
6 observed that some soybean accessions that had no known pedigree relationship showed < 1.50%  
7 inconsistencies (e.g. 0.87% between Williams 82K and KLS85102). Therefore, we used a 1.25%  
8 inconsistency value as the cut-off to remove redundant or highly similar accessions from groups of  
9 duplicates or near isogenic lines. The same cut-off value was applied to filtration of wild soybean  
10 accessions. In each of the genotype duplicate sets, an accession with a sample call rate  $\geq 99\%$  was  
11 preferentially retained. For each of the near-isogenic line groups, the recurrent parent or representative  
12 single-gene isolate in case of the absence of parents was retained. Of the 4,184 accessions genotyped in  
13 this study, 1,148 (867 domesticated and 281 wild) were removed (Fig. 1C). The high rate of redundant and  
14 highly similar accessions has been frequently reported in worldwide germplasm collections (Food and  
15 Agriculture Organization of the United Nations 2010; McCouch et al. 2012). For domesticated soybeans,  
16 the major cause in the National Agrobiodiversity Center in Korea is probably the unknowing submission of  
17 the same accession with different collection sites and designators because there are many accessions with  
18 the same common name but with different collection sites or collectors. For wild soybeans, multiple  
19 accessions were collected from a narrow habitat area.

20 After the filtration, a final set of 3,036 genotyped accessions was available for population structure  
21 analysis (Table S1 and S2). Only a few accessions were removed from countries other than South Korea.  
22 As a result, overall proportion of soybean accessions among countries in this 3,036 set was similar to that  
23 of the 4,234 set (Table S2). In the 3,036 set, 1,957 were *G. max* accessions and 1,079 were *G. soja*  
24 accessions. Representative *G. soja* accessions from China, Russia, and Japan remained to evenly reflect the  
25 geographic distribution of native *G. soja* in each of these country regions (Fig. 1D).

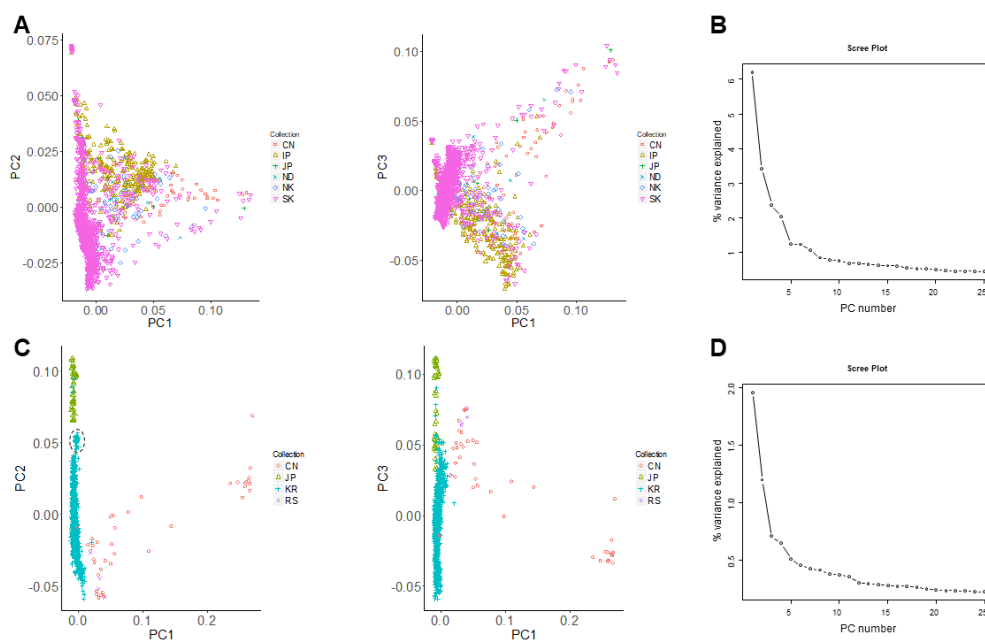
26

### 27 **3.2. Population structure**

1 ADMIXTURE (Alexander et al. 2009) and PCA (Patterson et al. 2006) were used to infer population  
2 structure of the 3,036 non-redundant soybean set using 117,095 SNPs (heterozygous rate < 0.02 and minor  
3 allele frequency > 0.02). As observed in the analysis of our total population of 4,234 accessions, the 3,036  
4 accessions were clearly divided into two large groups, representing *G. max* and *G. soja* (Fig. S1). Both the  
5 estimated cross-validation (CV) error plot from ADMIXTURE and scree plot from the PCA supported the  
6 presence of two large groups (Fig. S1), although the slopes did not level off, which is likely because of  
7 subgroupings within the two large groups. The clear separation of *G. max* and *G. soja* groups might be  
8 expected by the ascertainment bias, which favored selection of *G. max* SNPs (Lee et al. 2015), and the  
9 sampling bias. Unlike the previous observations that the genome diversity level was > 2-fold lower in the  
10 domesticated soybeans relative to that in the wild soybeans, the diversity level of the domesticated  
11 soybeans (mean per-site nucleotide diversity ( $\pi$ ) = 0.189) estimated from the 117,095 SNPs was ~1.58-fold  
12 lower than that of the wild soybeans ( $\pi$  = 0.298) and nearly two times more domesticated soybean  
13 accessions were used for the population structure analyses. The 3,036 set also contained excessive  
14 numbers of accessions collected in South Korea in both the domesticated and wild soybean groups.  
15 Interestingly, in the PCA space constructed with the first two PCs, *G. soja* accessions from Japan and  
16 China that were located at both ends of the *G. soja* cluster were almost equally close to the *G. max* cluster.  
17 However, in the PC1 and PC3 plot, accessions from Japan were closest to the *G. max* cluster and  
18 accessions from China were the most distantly related to the *G. max* cluster (Fig. S1).

19 When we analyzed *G. max* and *G. soja* separately, a somewhat distinct subpopulation structure was  
20 revealed. Within each of *G. max* and *G. soja* populations, CV errors of the ADMIXTURE runs decreased  
21 gradually without a steep drop (Fig. S2), whereas eigenvalues of the PCA runs showed steep decrease up  
22 to  $K = 5$  (Fig. 2), indicating that there were at least four distinct subpopulations in each of the *G. max* and  
23 *G. soja* populations. Both the ADMIXTURE and PCA plots from the 1,957 domesticated soybeans did not  
24 show distinctive grouping (Fig. 2 and S2). The majority of South Korean domesticated accessions (~80%)  
25 formed a dense subpopulation likely because of recent overcollection. This notion is supported by that the  
26 rest of the Korea accessions were well mixed with Chinese, North Korean, and Japanese accessions, which  
27 did not show distinct subgrouping on a geographic basis. Notably, the North Korea accessions, which

1 could be considered true landraces because of their collections during the first half of the 20<sup>th</sup> century  
2 before modern breeding research, were evenly distributed across subpopulations. The improved cultivars  
3 were narrowly clustered in the PCA plot, indicating much lower diversity relative to that of the entire  
4 domesticated soybeans. The 1,079 wild soybean population showed distinctive subpopulations (Fig. 2 and  
5 S2), as shown in the analyses of the entire 4,234 population. The groupings were consistent with  
6 geographic distributions of the collection sites. Korean accessions and Japanese accessions formed a  
7 unique subpopulation, respectively. Chinese formed two subpopulations. Five accessions from the Russian  
8 border clustered together with those from northeast China. A strong relationship between subpopulations  
9 and their geographic distribution was notably exemplified by accessions from Jeju Island located 130 km  
10 off the southern coast of the Korean Peninsula; although they belong to the Korean subpopulation, all  
11 accessions from Jeju Island formed a subgroup that was the closest to the Japanese subpopulation (Fig. 2C).



12  
13 **Figure 2.** Population structures of 1,957 domesticated and 1,079 wild soybean accessions in the 3,036 non-  
14 redundant soybean accession set. (A) Principal components (PC) of SNP variation in the domesticated  
15 population. The plots show the first three principal components. The countries of collection or  
16 improvement status of the soybean accessions in (A) and (C) are represented by two-letter codes—CN,  
17 China; IP, improved breeding line; JP, Japan; ND, not determined; NK, North Korea; RS, Russia; and SK

1 (KR), South Korea. (B) Scree plot of the PC number and their contribution to variance from principal  
2 component analysis of the domesticated accessions. (C) Principal components of SNP variation in the wild  
3 population. The plots show the first three principal components. A cluster of accessions from Jeju Island is  
4 indicated by a dashed ellipse. (D) Scree plot of the PC number and their contribution to variance from  
5 principal component analysis of the wild accessions.

6

### 7 **3.3. Detection of SNPs associated with domestication history**

8 Domestication is a process of continuous artificial selection of a group of traits, collectively called  
9 domestication syndrome. The domestication process has produced selective sweeps with significant  
10 reductions in nucleotide diversity (Doebley et al. 2006; Hufford et al. 2012; Chung et al. 2014) on limited  
11 regions of the genome (approximately 5 ~ 10% of the genome). Numerous recent whole-genome  
12 resequencing studies have effectively detected the selective sweeps, which are associated with  
13 domesticated genes (Meyer and Purugganan 2013), by examining reduction of diversity (ROD) in  
14 windows along chromosomes. However, our SoyaSNP array data are not dense enough to detect the  
15 reduction of diversity. Thus, we attempted to detect SNPs associated with domestication using a case-  
16 control GWA method that analyzed binary domestication phenotypes, which were determined by presence  
17 (*G. max*) or absence (*G. soja*) of domestication.

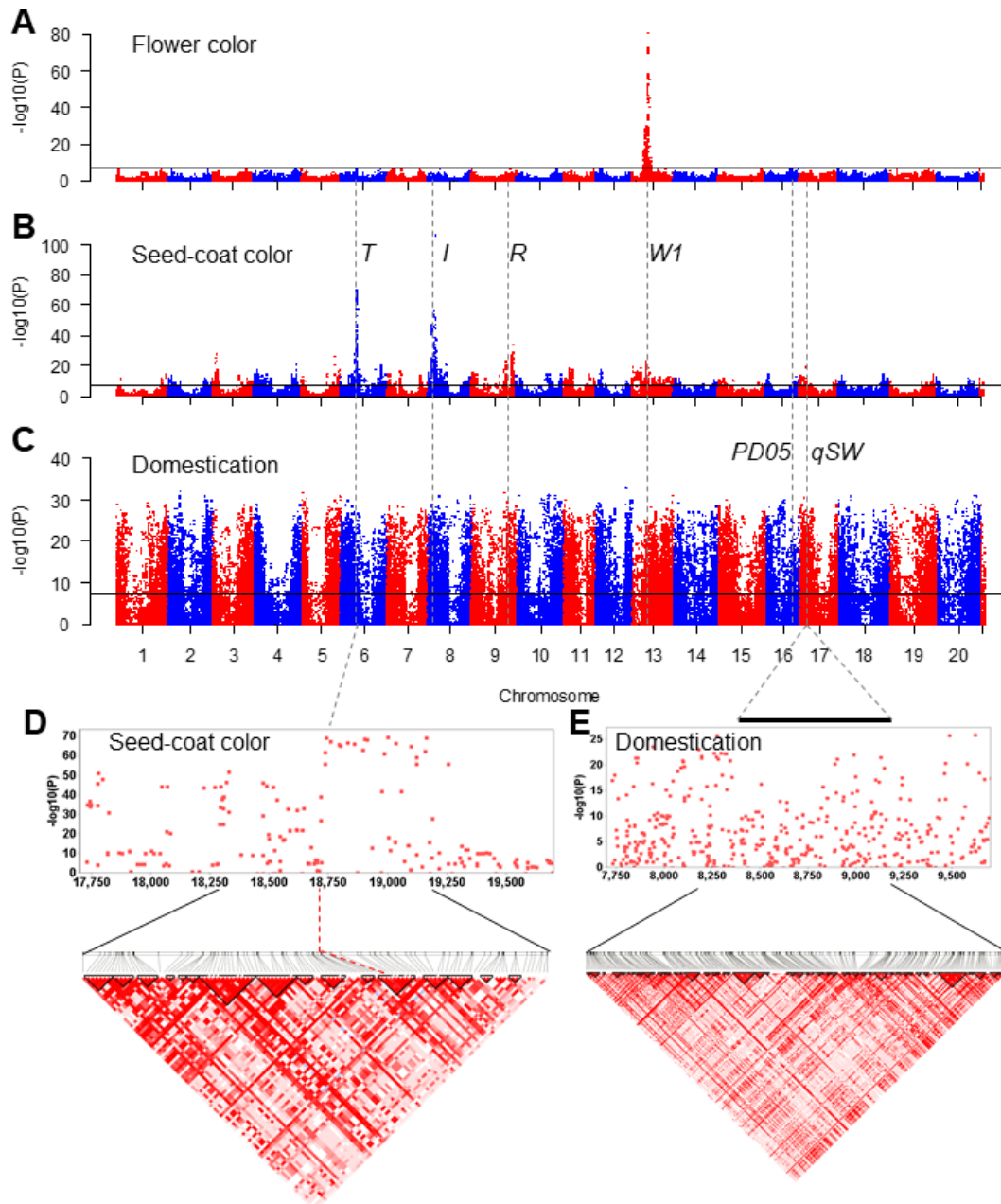
18 To test if our case-control GWA method enabled to find genes or chromosomal regions underlining  
19 binary phenotypes in our 3,036 non-redundant population, we chose two highly studied phenotypes—  
20 flower and seed-coat colors—which are monogenic and multigenic, respectively. Because our population  
21 was highly structured, we performed logistic regression model analysis conditional on a list of  
22 subpopulation-specific SNPs. The selected specific SNPs included one perfect domestication-specific SNP  
23 with each allele being perfectly correlated with *G. max* or *G. soja* membership of soybean accessions, and  
24 ten subpopulation-specific SNP within each of the *G. max* and *G. soja* populations (Table S4). The flower  
25 color phenotypes were divided into absence or presence of anthocyanin deposition colors. The seed-coat  
26 color phenotypes were divided into absence or presence of anthocyanin deposition colors. Using the  
27 conditional logistic regression model, we detected a broad and strong peak for flower color with the most

1 significant SNP (max  $-\log_{10}P = 80.6$ ) located at 17,877,234 on chromosome 13 (Figs. 3A, S3, and S4C).  
2 This peak area contained the *WI* locus, which is the major locus determining flower color (Zabala and  
3 Vodkin 2007). However, the most significant SNPs were located ~ 500 kb off the position of the *flavonoid*  
4 *3'5'-hydroxylase* gene, which is the causal gene of the *WI* locus. To understand this region, we estimated  
5 pairwise LD for SNPs from 16.8 Mb to 18.8 Mb. A strong LD pattern was observed between all the SNPs  
6 under the most significant SNPs, however no clear LD pattern was observed near the *flavonoid 3'5'*-  
7 *hydroxylase* gene. Although the result might be due to a SNP density insufficient for a long LD block, it  
8 has been often observed that a causal gene for a strong peak are not always corresponding with the highest  
9  $-\log_{10} P$  value (Segura et al. 2012; Yano et al. 2016).

10 We detected more than 30 peaks exceeding a significant threshold ( $-\log_{10}P \geq 7.07$ ) for seed-coat color  
11 (Figs. 3B, S3, and S4). The top three peaks were correlated with three known major loci (*I*, *R*, and *T* on  
12 chromosomes 8, 9, and 6, respectively) that control the deposition of various anthocyanin pigments in seed  
13 coat (Yang et al. 2010). The highest peak was generated from a chromosomal region surrounding the  
14 inverted CHS gene repeats, which is the causal region of the *I* locus (Clough et al. 2004). A strong LD  
15 pattern was observed at the *I* locus region. An SNP AX-90432942 on chromosome 6 with the second  
16 highest  $-\log_{10}P$  value = 69.8 was generated from flavonoid 3' hydroxylase, which is the causal gene of the  
17 *T* locus. Finally, the SNPs near the R2R3 MYB transcription factor gene, a strong candidate gene for the *R*  
18 locus, reported by Gillman et al. (2011) were not significantly associated with seed-coat colors. However,  
19 the gene is one of the R2R3 MYB transcription factor genes tandemly repeated in this chromosomal region  
20 and numerous highly significant SNPs were observed 100-kb away from the proposed *R* gene.

21 Interestingly, the peak on chromosome 13 are located on the *WI* locus that influences seed-coat color in a  
22 case of the homozygous recessive *it* genotypes (Palmer et al. 2004). Considering such a wide range of  
23 soybean seed-coat color variations, the detection of numerous minor peaks is not surprising, as reported by  
24 Song et al. (2016), although it is still surprising to detect this large number of significant peaks using  
25 binary phenotyping data. Nevertheless, we think that some of those minor peaks are inevitably false  
26 because the limited number of conditional SNPs could not correct for all inflation of the statistic caused by  
27 population substructure.

1



2

3 **Figure 3.** Genome-wide association scans for 3,036 soybean accessions for flower color, seed-coat color,  
4 and domestication. (A) Manhattan plot for flower color. The solid horizontal line denotes the Bonferroni-  
5 adjusted significance threshold. Chromosomal regions of known genes (*T*, *I*, *R*, *W1*) or loci (*PD05* and  
6 *qSW*) are indicated by dashed vertical lines. (B) Manhattan plot for seed-coat color. (C) Manhattan plot for

1 domestication. (D) Local Manhattan plot (top) and LD heatmap (bottom) surrounding the *T* locus on  
2 chromosome 6. Dashed lines indicate the region of the *T* locus. Physical locations (kb) are indicated under  
3 the Manhattan plot. (E) Local Manhattan plot (top) and LD heatmap (bottom) surrounding the *qSW* locus  
4 on chromosome 17. A bar indicate the region of the *qSW* locus.

5  
6 Since our logistic regression model could readily detect loci associated with flower and seed-coat  
7 colors using binary phenotypes, we performed GWAS for domestication syndrome using the binary  
8 domestication phenotypes. For this analysis, we excluded the perfect domestication-specific SNP in the list  
9 of subpopulation-specific SNPs (Table S4). We detected numerous peaks for domestication syndrome over  
10 the genome, as expected from the previous studies (Hufford et al. 2012; Meyer and Purugganan 2013;  
11 Chung et al. 2014) that showed that domestication features covered approximately ~7% of the crop  
12 genome. To examine if previously detected domestication regions were also detected in the current study,  
13 we compared peak locations from our study with selective signals previously detected for two  
14 domestication traits, pod dehiscence and seed weight (Figs. 3 and S4D), which are two of the most critical  
15 domestication traits among the traits assayed by Zhou et al. (2015). The 190-kb region (*PD05*) responsible  
16 for pod dehiscence was also detected in our study with three highly significant SNPs ( $-\log_{10} P \geq 20$ ) and  
17 the *qSW* locus for seed weight was detected with  $> 30$  highly significant SNPs ( $-\log_{10} P \geq 20$ ). Lengths of  
18 strong LD blocks under the peaks corresponding to the *PD05* and *qSW* loci were not sufficient to define  
19 the selective sweeps that have been known to be  $> 100$  kb, as observed in our LD analyses for the flower  
20 and seed-coat color loci. Flower and seed-coat colors analyzed in this study are considered domestication-  
21 or diversification-related morphological features because nearly all wild soybean accessions have purple  
22 flowers and black seed coats. As expected, their major loci, *W1*, *I*, *R*, and *T* were also detected with highly  
23 significant SNPs ( $-\log_{10} P \geq 20$ ) in this GWAS for domestication syndrome (Fig. 3).

24

### 25 **3.4. Extraction of a representative set**

26 Our population structure and diversity analysis of the 3,036 non-redundant soybean population resulted in  
27 the identification of Japan as a likely center of soybean domestication. Since a fundamental assumption of

1 model-based methods, such as ADMIXTURE and PCA, is that the sample available for analysis is  
2 representative of the entire population distribution, sample sizes of subpopulations can substantially affect  
3 population stratification and ancestral population inference (McVean 2009; Shringarpure and Xing 2014).  
4 To investigate the possibility that excessive numbers of domesticated or Korean soybean accessions might  
5 have caused bias in inference of population structure of wild and domesticated soybeans, we obtained  
6 representative domesticated and wild soybean sets by selecting one from each of tightly distributed  
7 soybean miniclusters in the PCA plots, with a caution that overall distribution patterns are maintained (Fig.  
8 S5). For the representative set of wild soybeans, we filtered the population of the tightly distributed Korean  
9 accessions and selected 50 diverse Korean wild accessions (Table S2). In addition, four wild soybean  
10 anomalies misplaced to subpopulations different from subpopulations predicted by their collection site  
11 records were excluded; three Korean and one Chinese accessions (Fig. S6). For the representative set of  
12 domesticated soybeans, we selected 50 diverse *G. max* accessions that represent diversity of 1,957 *G. max*  
13 accessions. The results of the AMOVA indicated that the overall genetic structure observed in the 3,036  
14 non-redundant soybean population was well represented by the extracted representative set with some  
15 decrease of genetic diversity in the *G. max* population (Table S5). PCA plots from the resultant  
16 representative set of 194 soybean accessions showed distribution patterns similar to those from the 3,036  
17 non-redundant soybean accessions, although relative sizes of *G. max* and *G. soja* distributions in the PCA  
18 spaces constructed with the first and third PCs were reversed. The last drop of eigenvalues from the PCA  
19 runs occurred between  $K = 5$  and  $K = 6$  (Fig. S7), indicating that there were five distinct subpopulations.

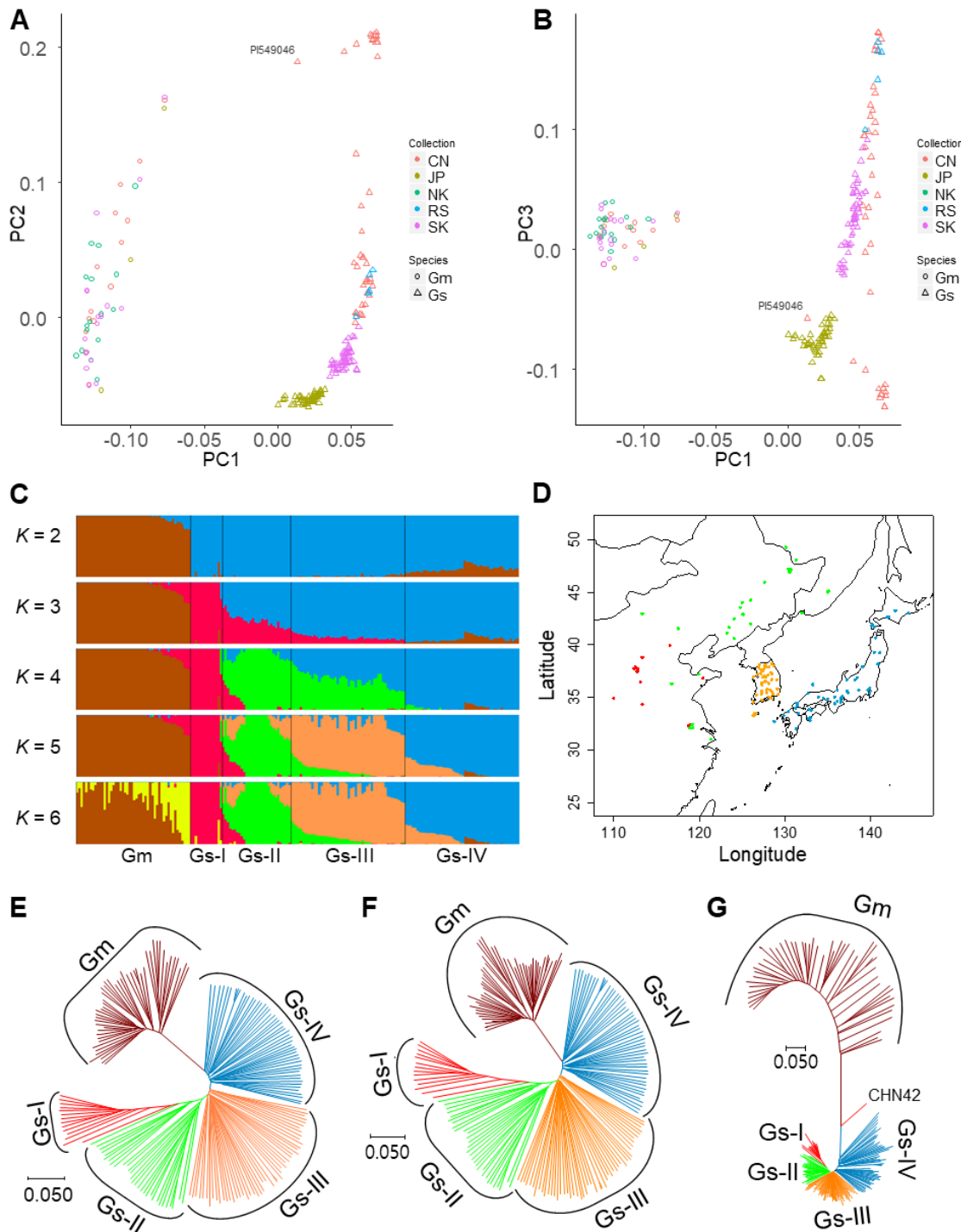
20

### 21 **3.5. Center of soybean domestication**

22 Because *G. soja* can be found *in situ* across most of the East Asia, it is important to establish the  
23 population structure, if any, of a diverse collection of *G. soja* accessions and to associate one or more of  
24 these populations with a collection of domesticated *G. max* varieties. To perform these experiments, we  
25 analyzed the population structure of the representative set of 194 soybean accessions with ADMIXTURE,  
26 and found  $K = 5$  populations based on the estimated CV error plot (Fig. S7). Thus, the soybean accessions  
27 were partitioned into one *G. max* (Gm) and four *G. soja* (Gs-I, Gs-II, Gs-III, and Gs-IV) subgroups (Fig.



1 4A, B). Wild accessions from China were divided into two subgroups, Gs-I and Gs-II. The Gs-I group  
2 showed the least diversity (Table 1) and most of them distributed in the middle region of the Yellow River  
3 basin. The Gs-II group was dispersed across northeast China, south China, and the Russian border of  
4 northeast China. Interestingly, this grouping result is remarkably similar to that obtained by a recent  
5 comprehensive study that showed that, by analyzing a total of 712 *G. soja* individuals from 40 natural  
6 populations in China, Chinese wild soybeans were grouped into two main subgroups, which were one from  
7 the Yellow River basin and the other from northeast China and south China (Guo et al. 2012). The Gs-IV  
8 distributed in Japan. The Gs-III showed the greatest diversity and distributed in South Korea and most of  
9 them appeared to be ancient admixture between Gs-II and Gs-IV. Interestingly, despite clear separation of  
10 the Chinese *G. soja*, diversity level of the combined population of Gs-I and Gs-II was similar to those of  
11 Korean or Japanese *G. soja*. An independent Gm group appeared from  $K = 2$  to  $K = 5$  (Fig. 4C).  
12 Interestingly, major genomic fractions of the Gm subgroup consistently appeared as minor genomic  
13 fractions of the Gs-IV and minor genomic fractions of the Gm group was a genomic fraction of Gs-I  
14 ancestry. The results suggested that after domestication of the Gm subgroup from the Gs-IV subgroup, the  
15 Gm subgroup was substantially diversified by introgression of the Gs-I genomic fractions. One of the Gs-I  
16 accessions, PI 459046, appeared to be a *G. max* x *G. soja* hybrid, although its genomic fraction (~22%)  
17 from *G. max* are lower than our hybrid filtration criteria (30% domesticated ancestry), which was less  
18 stringent than 20% domesticated ancestry used in other admixture studies (e.g. Wang et al. (2017)).  
19



1

2 **Figure 4.** Identification of the domestication center of *G. max*. (A, B) Principal components plots of SNP  
3 variation. PC1, PC2, and PC3 indicate score of principal components 1, 2, and 3, respectively. Each of PC1,

1 PC2, and PC3 explained 12.0%, 5.2%, and 2.6% of variance in the data. Countries of collection of the  
2 soybean accessions and species names are represented by two-letter codes —CN, China; JP, Japan; NK,  
3 North Korea; RS, Russia; SK, South Korea; Gm, *G. max*; and Gs, *G. soja*. A putative hybrid PI 459046 is  
4 labeled. (C) Population structure of 50 *G. max* (Gm) and 144 *G. soja* (Gs-I, Gs-II, Gs-III, and Gs-IV)  
5 accessions inferred using ADMIXTURE. Each color represents one population. PI 459046 showed ~20%  
6 of ancestral genomic fractions from *G. max*. (D) Geographic distribution of the four *G. soja* subgroups.  
7 Gs-I is red, Gs-II green, Gs-III orange, and Gs-IV blue. (E, F, G) Neighbor-joining phylogenetic tree of  
8 194 soybean accessions based on the SNPs genotyped by the 180K AXIOM SoyaSNP array, with  
9 evolutionary distances measured by the  $p$  distance. The taxa used in the neighbor-joining tree and bootstrap  
10 values from 1000 bootstrap replications at branches are described in Fig. S8. (E) Phylogenetic tree based  
11 on 117,095 SNPs. (F) Phylogenetic tree based on 108,897 SNPs, which are weakly or not significantly  
12 associated with domestication traits. (G) Phylogenetic tree based on 8197, which are very significantly  
13 associated with domestication traits. PI 459046 from group Gs-I clusters between Gm and Gs-IV likely  
14 because of contribution of ancestral genomic fraction from Gm.

15

16 We constructed a neighbor-joining (NJ) tree for the representative soybean set (Fig. 4E and S8). The  
17 tree showed that all *G. max* accessions formed a monophyletic cluster. Although *G. max* was artificially  
18 selected recently, terminal branch lengths were similar to those of *G. soja* likely because of ascertainment  
19 bias that more SNPs were selected from *G. max* than from *G. soja* (Lee et al. 2015). The population of the  
20 nearest branches, which were basal to the *G. max* soybean lineage, was *G. soja* subgroup Gs-IV. Within the  
21 Gs-IV that contains wild soybeans from Japan, those from eastern Japan area were closer to the *G. max*  
22 soybean lineage. To measure population differences and similarities, we calculated the fixation index  
23 values ( $F_{ST}$ ) (Holsinger and Weir 2009) between *G. max* and each *G. soja* population (Table 2). The  
24 pairwise  $F_{ST}$  values ranged from 0.201 to 0.334. The value of  $F_{ST}$  for the *G. max* and Gs-IV populations  
25 was the smallest, suggesting that *G. max* was domesticated directly from *G. soja* subpopulation Gs-IV. The  
26 level of population differentiation among *G. soja* subpopulation was much lower than that between *G. soja*  
27 and *G. max*, similar to the case of rice (Huang et al. 2012). However,  $F_{ST}$  values between *G. soja*

1 subpopulations corresponded with their geographic distances. The regions containing those wild  
2 populations that are phylogenetically close with cultivars could be proposed as the domestication region of  
3 crops (Matsuoka et al. 2002; Spooner et al. 2005). Thus, our results suggested that soybeans had been most  
4 likely domesticated only once in eastern Japan.

5

6 Table 1. Mean per-site nucleotide diversity ( $\pi$ ) of *Glycine max* (Gm) and each of *G. soja* groups (Gs-I to  
7 Gs-IV) in the representative soybean set

Group	Representative soybean set		Representative soybean set and 45 wild soybeans from Zhou et al. (2015)	
	Number of accessions	$\pi$	Number of accessions	$\pi$
Gm	50	0.236	50	0.237
Gs-I	14	0.198	29	0.222
Gs-II	30	0.276	45	0.283
Gs-I and Gs-II	44	0.290	74	0.301
Gs-III	50	0.294	57	0.299
Gs-IV	50	0.276	58	0.283
Gs-III and Gs-IV	100	0.296	115	0.302

8

9

10 Because the number of *G. soja* accessions from China in our representative set was smaller than those  
11 from Korea and Japan, the population structure revealed by our representative set was further resolved by  
12 incorporating the SNP data from 62 *G. soja* accession genomes resequenced by Zhou et al. (2015). By  
13 intersecting these SNPs with the set of 117,095 high-quality SNPs selected in this study, we extracted  
14 103,801 SNPs, which were shared between the genome resequencing data and our representative set. Of  
15 the 62 accessions, only 45 were incorporated into the representative set because of the high level (eleven  
16 accessions, > 20%) of heterozygous SNPs, hybrid (one), and overlapping (five) (Fig. S9 and Table S6).  
17 The resultant expanded set contained twelve diverse accessions from Zhejiang, China and one accession  
18 from Taiwan, thus increasing geographic coverage of this study further down to southern China. In results,  
19 the diversity level of *G. soja* accessions from China and its Russian border was similar to those from Korea  
20 or Japan (Table 1). Population structure of the expanded set inferred from ADMIXTURE and PCA was  
21 quite similar to that of our representative set, although the *G. max* accessions appeared to be divided into  
22 two groups likely because of the high level of heterozygous SNPs from the genome resequencing data (Fig.

1 S10, S11, and Table 2). Phylogenetic analysis and estimated  $F_{ST}$  values between subpopulations indicated  
2 that *G. soja* accessions collected from eastern Japan were closest to the *G. max* soybean lineage.

3

4 Table 2.  $F_{ST}$  values between *Glycine max* (Gm) and each of *G. soja* groups (Gs-I to Gs-IV) and between *G.*  
5 *soja* groups

Representative soybean set				
	Gs-I	Gs-II	Gs-III	Gs-IV
Gm	0.334	0.267	0.226	0.201
Gs-I		0.160	0.180	0.214
Gs-II			0.062	0.113
Gs-III				0.058
Representative soybean set and 45 wild soybeans from Zhou et al. (2015)				
Gm	0.323	0.264	0.227	0.204
Gs-I		0.164	0.176	0.209
Gs-II			0.066	0.118
Gs-III				0.055

6

7

8 To evaluate the distribution of SNPs associated with domestication syndrome across soybean  
9 subpopulations, we divided the 117,095 SNPs into 8,197 SNPs highly significantly associated with  
10 domestication traits and 108,898 SNPs weakly or not significantly associated with domestication traits.  
11 The 8,197 SNPs ( $-\log_{10} P > 17$ ) were selected because the previous studies have shown that ~7% of the  
12 crop genome is domestication-related (Hufford et al. 2012; Xu et al. 2012; Chung et al. 2014). Our GWAS  
13 also indicated that, in our GWAS population, strong LD extent under the highly significant SNPs in a peak  
14 tended to be much shorter than chromosomal extent under all the significant SNPs of the peak. The tree  
15 constructed from the 108,898 SNPs (Fig. 4F) was similar to that constructed from the 117,095 SNPs in  
16 their overall grouping and branch patterns, except that the branch length and grouping of the *G. max* clade  
17 were slightly different to each other. Grouping patterns in the tree constructed from the 8,196 SNPs (Fig.  
18 4G) were similar to those in the trees constructed from both 108,898 SNPs and 117,095 SNPs, except that  
19 the putative hybrid PI 459046 moved the closest to the *G. max* clade. Interestingly, the lengths of the basal  
20 and terminal branches for the *G. max* clade and the *G. soja* Gs-IV clade in the tree from the 8,196 SNPs  
21 became distinctively longer than those in the other *G. soja* clades. The results indicated that initial major  
22 artificial selection for soybean domestication was limited to a *G. soja* group from Japan.

1

## 2 **4. Discussion**

3 The present study analyzed genome-wide SNP variations obtained from thousands of soybean accessions,  
4 the majority of which were collected from the Korean peninsula. The results provide insight into the  
5 development of strategies for efficient and directed germplasm use as well as for collection of novel  
6 landraces and wild relatives. Population structure and grouping analyses revealed strong correlations  
7 between genetic distance and geographic distance in *G. soja* (wild soybean) populations and weak  
8 correlations in *G. max* (domesticated soybean) populations. *G. soja* accessions were divided into four  
9 distinct subgroups; Gs-I and Gs-II from China and its Russian border, Gs-III from Korea, and Gs-IV from  
10 Japan. The results suggest that although the Korean territory is much smaller than Chinese and Japanese  
11 territories, the ocean-imposed geographic separation among these countries has been a major contributor to  
12 the evolutionary divergence of *G. soja*. Most of the Gs-III group from Korea appeared to be ancient  
13 admixtures between Gs-II and Gs-IV, suggesting that *G. soja* spread from each of China and Japan might  
14 be mixed in Korea. Thus, Korean wild soybeans are more valuable resources than the other countries' wild  
15 soybeans because they alone provide variations from two large subgroups. Interestingly, accessions from  
16 Jeju Island off the southern coast of the Korean Peninsula are the closest grouped to the Gs-IV group  
17 among members of the Gs-III group, indicating that although our estimated  $F_{ST}$  values between *G. soja*  
18 groups denied appearance of a new distinct group, more extensive sampling from diverse areas will likely  
19 reveal better correlations between geographic and genetic distances among *G. soja* subpopulations. The *G.*  
20 *max* population was divided into four subgroups. However, it was apparent that the subgrouping did not  
21 reflect geographic origins. Particularly, landraces from North Korea that would be considered true  
22 landraces based on their collection time appeared in every subgroup. The majority of South Korean  
23 landraces that had been collected recently were grouped together. The majority of improved accessions  
24 from the United States were clustered closely together, supporting a previous observation (Hyten et al.  
25 2006). Taken together, our results suggest that while *G. soja* germplasm will require additional sampling  
26 from diverse indigenous areas to expand the germplasm base, *G. max* germplasm is saturated in terms of

1 genetic diversity. Thus, extensive genotyping and phenotyping of extant *G. max* germplasm would be the  
2 next step to expand the germplasm base of *G. max*.

3 Our results provided strong support for a single origin of *G. max* from eastern region in Japan,  
4 although pointing to a specific region in Japan likely requires analysis of more extensive wild and landrace  
5 soybean accessions from Japan. Whether a crop species stems from a single domestication event or from  
6 multiple independent domestications has been consistent with whether the domesticated species are  
7 monophyletic or polyphyletic, respectively, in the phylogenetic trees constructed from both the  
8 domesticated and wild progenitor species. Although diversity of chloroplast DNA, which represents  
9 maternal lineage of soybean, revealed multiple lineages of domesticated soybeans, analyses of recent  
10 genome-wide soybean variation data (Guo et al. 2010; Lam et al. 2010; Chung et al. 2014; Zhou et al. 2015)  
11 consistently showed the monophyletic nature of *G. max*, as observed in this study. In other words, recent  
12 soybean phylogenetic studies collectively indicated a single origin of *G. max*. The best examples of  
13 monophyletic grouping are wheat and barley, which appear to have been domesticated once from their wild  
14 ancestors in the Fertile Crescent (Badr et al. 2000; Ozkan et al. 2002). The origin of barley was further  
15 supported by the genome sequences of five 6,000-year-old barley grains (Mascher et al. 2016). In cases of  
16 rice and common bean that showed polyphyletic groupings, single or multiple regions of origin of these  
17 crop species are still contentious (Molina et al. 2011; Bitocchi et al. 2012; Huang et al. 2012). This  
18 controversy may have arisen because most modern wild accessions studied represent descendants of  
19 ancient feralization of admixed accessions that resulted from hybridization events between domesticated  
20 species and wild species populations after domestication (Wang et al. 2017), indicating that one of the  
21 previously thought independent origin regions might be a secondary origin region. The grouping of PI  
22 459046 in this study is a good example that shows how hybrids could mislead inference of relationship  
23 between wild and domesticated crop species.

24 A recent comprehensive study of the archaeological records for soybean from Japan, China, and  
25 Korea indicated that Japan could have been a source of a large-seeded landrace of domesticated soybean  
26 that spread to Korea and subsequently to China (Lee et al. 2011). The archaeological records suggest that  
27 selection of large seed sizes occurred in Japan (Lee et al. 2011; Nakayama 2015) by 5,500 calibrated years

1 (cal) before present (BP) and in Korea (Lee et al. 2011) by 3,500 cal BP. Seed size is clearly a  
2 domestication trait because the seed size of *G. soja* is much smaller than that of *G. max* landraces (Broich  
3 and Palmer 1980). However, the archaeological data were interpreted to suggest the multiple origins  
4 hypothesis of soybean. One particular reason is that the excavated tiny seeds were as old as 9,000–8,600  
5 cal BP in northern China and 7,000 cal BP in Japan. However, the size of the seeds is similar to that of the  
6 seeds of present-day wild soybeans, and so would have been quite different from the landraces already  
7 grown in China by 2,500 BP. Another reason is that the interpretation was greatly influenced by a previous  
8 report that diversity of chloroplast DNA SSRs in wild and domesticated soybeans showed evidence for  
9 multiple origins of domesticated soybean (Xu et al. 2002). However, as mentioned above, numerous recent  
10 genome-wide soybean genome variation studies consistently show a single origin of *G. max*.

11 One of the main reasons that the previous studies pointed different regions in China as the center of  
12 soybean domestication is likely sampling bias. Our results suggested that wild accessions from China had  
13 genetic diversity level almost equal to those from Korea or Japan. However, most previous studies tended  
14 to neglect this fact. In an extreme case (Han et al. 2016), no accession from Korea and Japan was used,  
15 with the conclusion that central China is the initial domestication region. Another confounding factor is the  
16 inclusion of hybrid soybeans from natural mating between *G. soja* and *G. max*. Hybrid soybeans were not  
17 recognized in many previous soybean population studies, although hybrids between wild and domesticated  
18 species have been increasingly regarded as a major problem in studies of crop domestication history  
19 (Bitocchi et al. 2012; Wang et al. 2017). Furthermore, it was often assumed that a region in China is a  
20 center of soybean domestication because hybrid soybeans are frequently found in China (Han et al. 2016).  
21 However, of the 50 hybrids that we removed to avoid their potential confounding effects in this analysis,  
22 the majority (36 of 50) were accessions from Korea. The domestication of domesticated plant species from  
23 their wild ancestors arose from rapid evolutionary changes in the past 13,000 years of Holocene human  
24 history (Diamond 2002; Larson et al. 2014). The list of origins and the list of the most productive areas of  
25 most of major crops in the modern world are almost mutually exclusive. This could be explained by that  
26 the domestication origin of a crop was merely a region to which the most numerous and most valuable  
27 domesticable wild plant species were native. In this respect, our result that shows Japan as the



1 domestication origin of soybean is not totally unexpected one.

2 Expanding on previous studies that reported genome-wide polymorphism data of soybean germplasm  
3 (Lam et al. 2010; Chung et al. 2014; Bandillo et al. 2015; Zhou et al. 2015; Valliyodan et al. 2016; Wang et  
4 al. 2016), our results show that samples intensively collected from Korea, which is a small area of the  
5 entire soybean distribution, provide sufficient amounts of data to underpin genome-wide population  
6 genetic questions that have been neglected or misled in the context of diversity and domestication panels  
7 of extant individuals. Our analysis demonstrates the value of current germplasm collections and how to  
8 expand the germplasm base. Furthermore, the findings show that a single major domestication event had  
9 occurred in a region of Japan. In addition, the high-density SNP array data enabled detection of  
10 domestication-associated SNPs and regions controlling important agronomic traits in a highly accurate  
11 manner. This suggests that our results will likely be useful for marker-assisted selection and genomic  
12 prediction to utilize unexplored genetic diversity in the soybean germplasm.

13

#### 14 **Data accessibility**

15 SNP genotype data are listed in Table S1 and are publicly available at Korean Soya Base  
16 ([http://koreansoyabase.org/Data\\_Resource/](http://koreansoyabase.org/Data_Resource/)).

17

#### 18 **Acknowledgments**

19 We thank Dr. Changyong Lee at Kongju National University for helpful comments in statistical analysis.  
20 We are grateful to Prof. Michael Gore at Cornell University for his critical reading of the revised version  
21 of this paper.

22

#### 23 **Conflict of interest**

24 None declared.

25

#### 26 **Supplementary data**

27 Supplementary data are available at DNARES online.

28

#### 29 **Funding**

30 S.C.J. was supported by the Next-Generation BioGreen 21 Program (PJ01321304), Rural Development  
31 Administration, and partly by the National Research Foundation grant (NRF-2018R1A2A2A05021904)

- 1 funded by the Korea government and by the Korea Research Institute of Bioscience and Biotechnology
- 2 Research Initiative Program. The work at Rural Development Administration was funded in part by Rural
- 3 Development Administration Project No. PJ01155401.

## 1 **References**

- 2 Alexander DH, Novembre J, Lange K. 2009. Fast model-based estimation of ancestry in unrelated  
3 individuals. *Genome Res* **19**: 1655-1664.
- 4 Badr A, Muller K, Schafer-Pregl R, El Rabey H, Effgen S, Ibrahim HH, Pozzi C, Rohde W, Salamini F.  
5 2000. On the origin and domestication history of Barley (*Hordeum vulgare*). *Mol Biol Evol* **17**:  
6 499-510.
- 7 Bandillo N, Jarquin D, Song Q, Nelson RL, Cregan P, Specht J, Lorenz A. 2015. A population  
8 structure and genome-wide association analysis on the USDA soybean germplasm collection.  
9 *The Plant Genome* **8**.
- 10 Barrett JC, Fry B, Maller J, Daly MJ. 2005. Haploview: analysis and visualization of LD and haplotype  
11 maps. *Bioinformatics* **21**: 263-265.
- 12 Bitocchi E, Nanni L, Bellucci E, Rossi M, Giardini A, Zeuli PS, Logozzo G, Stougaard J, McClean P,  
13 Attene G et al. 2012. Mesoamerican origin of the common bean (*Phaseolus vulgaris* L.) is  
14 revealed by sequence data. *Proc Natl Acad Sci U S A* **109**: E788-796.
- 15 Broich SL, Palmer RG. 1980. A cluster analysis of wild and domesticated soybean phenotypes.  
16 *Euphytica* **29**: 23-32.
- 17 Broich SL, Palmer RG. 1981. Evolutionary studies of the soybean: The frequency and distribution of  
18 alleles among collections of *Glycine max* and *G. soja* of various origin. *Euphytica* **30**: 55-64.
- 19 Browning SR, Browning BL. 2007. Rapid and accurate haplotype phasing and missing-data inference  
20 for whole-genome association studies by use of localized haplotype clustering. *Am J Hum*  
21 *Genet* **81**: 1084-1097.
- 22 Chung WH, Jeong N, Kim J, Lee WK, Lee YG, Lee SH, Yoon W, Kim JH, Choi IY, Choi HK et al.  
23 2014. Population structure and domestication revealed by high-depth resequencing of Korean  
24 cultivated and wild soybean genomes. *DNA Res* **21**: 153-167.
- 25 Clough SJ, Tuteja JH, Li M, Marek LF, Shoemaker RC, Vodkin LO. 2004. Features of a 103-kb gene-  
26 rich region in soybean include an inverted perfect repeat cluster of *CHS* genes comprising the  
27 *I* locus. *Genome* **47**: 819-831.
- 28 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth  
29 GT, Sherry ST et al. 2011. The variant call format and VCFtools. *Bioinformatics* **27**: 2156-  
30 2158.
- 31 Diamond J. 2002. Evolution, consequences and future of plant and animal domestication. *Nature* **418**:  
32 700-707.
- 33 Doebley JF, Gaut BS, Smith BD. 2006. The molecular genetics of crop domestication. *Cell* **127**: 1309-  
34 1321.
- 35 Excoffier L, Lischer HE. 2010. Arlequin suite ver 3.5: a new series of programs to perform population  
36 genetics analyses under Linux and Windows. *Mol Ecol Resour* **10**: 564-567.

- 1 Food and Agriculture Organization of the United Nations. 2010. *The second report on the state of the*  
2 *world's plant genetic resources for food and agriculture*. FAO, Rome.
- 3 Gabriel SB, Schaffner SF, Nguyen H, Moore JM, Roy J, Blumenstiel B, Higgins J, DeFelice M,  
4 Lochner A, Faggart M et al. 2002. The structure of haplotype blocks in the human genome.  
5 *Science* **296**: 2225-2229.
- 6 Gillman JD, Tetlow A, Lee JD, Shannon JG, Bilyeu K. 2011. Loss-of-function mutations affecting a  
7 specific *Glycine max* R2R3 MYB transcription factor result in brown hilum and brown seed  
8 coats. *BMC Plant Biol* **11**: 155.
- 9 Gizlice Z, Carter TE, Burton J. 1994. Genetic base for North American public soybean cultivars  
10 released between 1947 and 1988. *Crop Sci* **34**: 1143-1151.
- 11 Guo J, Liu Y, Wang Y, Chen J, Li Y, Huang H, Qiu L, Wang Y. 2012. Population structure of the wild  
12 soybean (*Glycine soja*) in China: implications from microsatellite analyses. *Ann Bot* **110**: 777-  
13 785.
- 14 Guo J, Wang Y, Song C, Zhou J, Qiu L, Huang H, Wang Y. 2010. A single origin and moderate  
15 bottleneck during domestication of soybean (*Glycine max*): implications from microsatellites  
16 and nucleotide sequences. *Ann Bot* **106**: 505-514.
- 17 Han Y, Zhao X, Liu D, Li Y, Lightfoot DA, Yang Z, Zhao L, Zhou G, Wang Z, Huang L et al. 2016.  
18 Domestication footprints anchor genomic regions of agronomic importance in soybeans. *New*  
19 *Phytol* **209**: 871-884.
- 20 Holsinger KE, Weir BS. 2009. Genetics in geographically structured populations: defining, estimating  
21 and interpreting  $F_{ST}$ . *Nat Rev Genet* **10**: 639-650.
- 22 Huang X, Kurata N, Wei X, Wang ZX, Wang A, Zhao Q, Zhao Y, Liu K, Lu H, Li W et al. 2012. A map  
23 of rice genome variation reveals the origin of cultivated rice. *Nature* **490**: 497-501.
- 24 Hufford MB, Xu X, van Heerwaarden J, Pyhajarvi T, Chia JM, Cartwright RA, Elshire RJ, Glaubitz JC,  
25 Guill KE, Kaeppeler SM et al. 2012. Comparative population genomics of maize domestication  
26 and improvement. *Nat Genet* **44**: 808-811.
- 27 Hwang EY, Song Q, Jia G, Specht JE, Hyten DL, Costa J, Cregan PB. 2014. A genome-wide  
28 association study of seed protein and oil content in soybean. *BMC Genomics* **15**: 1.
- 29 Hymowitz T. 2004. Speciation and cytogenetics. In *Soybeans: Improvement, production, and uses* (ed.  
30 HR Boerma, JE Specht), pp. 97-136. American Society of Agronomy, Madison, Wisconsin.
- 31 Hymowitz T, Kaizuma N. 1981. Soybean seed protein electrophoresis profiles from 15 Asian countries  
32 or regions: hypotheses on paths of dissemination of soybeans from China. *Econ Bot* **13**: 10-  
33 23.
- 34 Hyten DL, Song Q, Zhu Y, Choi IY, Nelson RL, Costa JM, Specht JE, Shoemaker RC, Cregan PB.  
35 2006. Impacts of genetic bottlenecks on soybean genome diversity. *Proc Natl Acad Sci U S A*  
36 **103**: 16666-16671.
- 37 Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis Version 7.0

- 1 for Bigger Datasets. *Mol Biol Evol* **33**: 1870-1874.
- 2 Lam HM, Xu X, Liu X, Chen W, Yang G, Wong FL, Li MW, He W, Qin N, Wang B et al. 2010.  
3 Resequencing of 31 wild and cultivated soybean genomes identifies patterns of genetic  
4 diversity and selection. *Nat Genet* **42**: 1053-1059.
- 5 Larson G, Piperno DR, Allaby RG, Purugganan MD, Andersson L, Arroyo-Kalin M, Barton L, Climer  
6 Vigueira C, Denham T, Dobney K et al. 2014. Current perspectives and the future of  
7 domestication studies. *Proc Natl Acad Sci U S A* **111**: 6139-6146.
- 8 Lee GA, Crawford GW, Liu L, Sasaki Y, Chen X. 2011. Archaeological soybean (*Glycine max*) in East  
9 Asia: does size matter? *PLoS One* **6**: e26720.
- 10 Lee YG, Jeong N, Kim JH, Lee K, Kim KH, Pirani A, Ha BK, Kang ST, Park BS, Moon JK et al. 2015.  
11 Development, validation and genetic analysis of a large soybean SNP genotyping array. *Plant*  
12 *J* **81**: 625-636.
- 13 Li YH, Li W, Zhang C, Yang L, Chang RZ, Gaut BS, Qiu LJ. 2010. Genetic diversity in domesticated  
14 soybean (*Glycine max*) and its wild progenitor (*Glycine soja*) for simple sequence repeat and  
15 single-nucleotide polymorphism loci. *New Phytol* **188**: 242-253.
- 16 Mascher M, Schuenemann VJ, Davidovich U, Marom N, Himmelbach A, Hubner S, Korol A, David M,  
17 Reiter E, Riehl S et al. 2016. Genomic analysis of 6,000-year-old cultivated grain illuminates  
18 the domestication history of barley. *Nat Genet* **48**: 1089-1093.
- 19 Matsuoka Y, Vigouroux Y, Goodman MM, Sanchez GJ, Buckler E, Doebley J. 2002. A single  
20 domestication for maize shown by multilocus microsatellite genotyping. *Proc Natl Acad Sci U*  
21 *S A* **99**: 6080-6084.
- 22 McCarthy MI, Abecasis GR, Cardon LR, Goldstein DB, Little J, Ioannidis JP, Hirschhorn JN. 2008.  
23 Genome-wide association studies for complex traits: consensus, uncertainty and challenges.  
24 *Nat Rev Genet* **9**: 356-369.
- 25 McCouch SR, McNally KL, Wang W, Sackville Hamilton R. 2012. Genomics of gene banks: A case  
26 study in rice. *Am J Bot* **99**: 407-423.
- 27 McVean G. 2009. A genealogical interpretation of principal components analysis. *PLoS Genet* **5**:  
28 e1000686.
- 29 Meyer RS, Purugganan MD. 2013. Evolution of crop species: genetics of domestication and  
30 diversification. *Nat Rev Genet* **14**: 840-852.
- 31 Molina J, Sikora M, Garud N, Flowers JM, Rubinstein S, Reynolds A, Huang P, Jackson S, Schaal BA,  
32 Bustamante CD et al. 2011. Molecular evidence for a single evolutionary origin of  
33 domesticated rice. *Proc Natl Acad Sci U S A* **108**: 8351-8356.
- 34 Nakayama S. 2015. Domestication of the soybean (*Glycine max*) and morphological differentiation of  
35 seeds in the Jomon period. *Jpn J Histor Bot* **23**: 33-42.
- 36 Ozkan H, Brandolini A, Schafer-Pregl R, Salamini F. 2002. AFLP analysis of a collection of tetraploid  
37 wheats indicates the origin of emmer and hard wheat domestication in southeast Turkey. *Mol*

- 1 *Biol Evol* **19**: 1797-1801.
- 2 Palmer RG, Pfeiffer TW, Buss GR, Kilen TC. 2004. Qualitative genetics. In *Soybeans: Improvement,*  
3 *production, and uses, 3rd edn.* (ed. HR Boerma, HE Specht), pp. 137-214. ASA, CSSA, and  
4 SSSA, Madison.
- 5 Patterson N, Price AL, Reich D. 2006. Population structure and eigenanalysis. *PLoS Genet* **2**: e190.
- 6 Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. 2006. Principal components  
7 analysis corrects for stratification in genome-wide association studies. *Nat Genet* **38**: 904-909.
- 8 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, de Bakker PI,  
9 Daly MJ et al. 2007. PLINK: a tool set for whole-genome association and population-based  
10 linkage analyses. *Am J Hum Genet* **81**: 559-575.
- 11 Rosenberg NA. 2004. DISTRUCT: a program for the graphical display of population structure. *Mol*  
12 *Ecol Notes* **4**: 137-138.
- 13 Segura V, Vilhjalmsjon BJ, Platt A, Korte A, Seren U, Long Q, Nordborg M. 2012. An efficient multi-  
14 locus mixed-model approach for genome-wide association studies in structured populations.  
15 *Nat Genet* **44**: 825-830.
- 16 Shringarpure S, Xing EP. 2014. Effects of sample selection bias on the accuracy of population  
17 structure and ancestry inference. *G3 (Bethesda)* **4**: 901-911.
- 18 Song J, Liu Z, Hong H, Ma Y, Tian L, Li X, Li YH, Guan R, Guo Y, Qiu LJ. 2016. Identification and  
19 validation of loci governing seed coat color by combining association mapping and bulk  
20 segregation analysis in soybean. *PLoS One* **11**: e0159064.
- 21 Song Q, Hyten DL, Jia G, Quigley CV, Fickus EW, Nelson RL, Cregan PB. 2013. Development and  
22 evaluation of SoySNP50K, a high-density genotyping array for soybean. *PLoS One* **8**: e54985.
- 23 Spooner DM, McLean K, Ramsay G, Waugh R, Bryan GJ. 2005. A single domestication for potato  
24 based on multilocus amplified fragment length polymorphism genotyping. *Proc Natl Acad Sci*  
25 *U S A* **102**: 14694-14699.
- 26 Turner SD. 2014. qqman: an R package for visualizing GWAS results using Q-Q and Manhattan plots.  
27 *bioRxiv* doi:10.1101/005165.
- 28 Valliyodan B, Dan Q, Patil G, Zeng P, Huang J, Dai L, Chen C, Li Y, Joshi T, Song L et al. 2016.  
29 Landscape of genomic diversity and trait discovery in soybean. *Sci Rep* **6**: 23598.
- 30 Vaughan DA, Lu BR, Tomooka N. 2008. The evolving story of rice evolution. *Plant Sci* **174**: 394-408.
- 31 Wang H, Vieira FG, Crawford JE, Chu C, Nielsen R. 2017. Asian wild rice is a hybrid swarm with  
32 extensive gene flow and feralization from domesticated rice. *Genome Res* **27**: 1029-1038.
- 33 Wang J, Chu S, Zhang H, Zhu Y, Cheng H, Yu D. 2016. Development and application of a novel  
34 genome-wide SNP array reveals domestication history in soybean. *Sci Rep* **6**: 20728.
- 35 Weir BS, Cockerham CC. 1984. Estimating F-statistics for the analysis of population structure.  
36 *Evolution* **38**: 1358-1370.
- 37 Xu H, Abe J, Gai Y, Shimamoto Y. 2002. Diversity of chloroplast DNA SSRs in wild and cultivated

- 1 soybeans: evidence for multiple origins of cultivated soybean. *Theor Appl Genet* **105**: 645-653.
- 2 Xu X, Liu X, Ge S, Jensen JD, Hu F, Li X, Dong Y, Gutenkunst RN, Fang L, Huang L et al. 2012.
- 3 Resequencing 50 accessions of cultivated and wild rice yields markers for identifying
- 4 agronomically important genes. *Nat Biotechnol* **30**: 105-111.
- 5 Yang K, Jeong N, Moon JK, Lee YH, Lee SH, Kim HM, Hwang CH, Back K, Palmer RG, Jeong SC.
- 6 2010. Genetic analysis of genes controlling natural variation of seed coat and flower colors in
- 7 soybean. *J Hered* **101**: 757-768.
- 8 Yano K, Yamamoto E, Aya K, Takeuchi H, Lo PC, Hu L, Yamasaki M, Yoshida S, Kitano H, Hirano K
- 9 et al. 2016. Genome-wide association study using whole-genome sequencing rapidly
- 10 identifies new genes influencing agronomic traits in rice. *Nat Genet* **48**: 927-934.
- 11 Zabala G, Vodkin LO. 2007. A rearrangement resulting in small tandem repeats in the *F3'5'H* gene of
- 12 white flower genotypes is associated with the soybean *W1* locus. *Crop Sci* **47**: S113-S124.
- 13 Zhou Z, Jiang Y, Wang Z, Gou Z, Lyu J, Li W, Yu Y, Shu L, Zhao Y, Ma Y et al. 2015. Resequencing
- 14 302 wild and cultivated accessions identifies genes related to domestication and improvement
- 15 in soybean. *Nat Biotechnol* **33**: 408-414.
- 16