# Long non-coding RNAs underlie multiple domestication traits and leafhopper resistance

**Jianxin Ma**
*maj@purdue.edu*
Purdue University    https://orcid.org/0000-0002-1474-812X

**Weidong Wang**
*wangwd@cau.edu.cn*
China Agricultural University    https://orcid.org/0000-0002-7110-5630

**Jingbo Duan**
*duan68@purdue.edu*
Purdue University    https://orcid.org/0000-0001-6467-9102

**Xutong Wang**
*wangxt881@gmail.com*
Purdue University

**Xingxing Feng**
*fengxingxing112@163.com*
Chinese Academy of Sciences

**Liyang Chen**
*zjuabcly@gmail.com*
Purdue University

**Chancelor Clark**
*clark367@purdue.edu*
Purdue University    https://orcid.org/0000-0002-2255-2514

**Stephen Swarm**
*stephen.swarm@beckshybrids.com*
University of Illinois

**Jinbin Wang**
*wang5549@purdue.edu*
Purdue University

**Sen Lin**
*lin1595@purdue.edu*
Purdue University

**Randall Nelson**
*rlnelson@illinois.edu*
University of Illinois    https://orcid.org/0000-0001-9482-1763

**Blake Meyers**
*bmeyers@danforthcenter.org*
Donald Danforth Plant Science Center    https://orcid.org/0000-0003-3436-6097

Xianzhong Feng

*fengxianzhong@iga.ac.cn*

Northeast Institute of Geography and Agroecology, CAS   https://orcid.org/0000-0002-7129-3731

Article

Keywords:

**Additional Declarations:** There is **NO** Competing Interest.

**Article:**

# Long non-coding RNAs underlie multiple domestication traits and leafhopper resistance

Weidong Wang[1,6,9], Jingbo Duan[1,9], Xutong Wang[1,7,9], Xingxing Feng[2,9], Liyang Chen[1], Chancelor B. Clark[1], Stephen A. Swarm[3,8], Jinbin Wang[1], Sen Lin[1], Randall L. Nelson[3], Blake C. Meyers[4,5], Xianzhong Feng[2*], Jianxin Ma[1*]

[1]Department of Agronomy, Purdue University; West Lafayette, IN 47907, USA.

[2]Key Laboratory of Soybean Molecular Design Breeding, Northeast Institute of Geography and Agroecology, Chinese Academy of Sciences; Changchun, Jilin 130102, China.

[3]Department of Crop Sciences, University of Illinois at Urbana–Champaign; Urbana, IL 61801, USA.

[4]Donald Danforth Plant Science Center, St. Louis, Missouri, USA.

[5]Division of Plant Science & Technology, University of Missouri-Columbia, Columbia, Missouri, USA.

[6]Present address: College of Agronomy and Biotechnology, China Agricultural University; Beijing 100193, China.

[7]Present address: Hubei Hongshan Laboratory; Wuhan, Hubei 430070, China.

[8]Present address: Beck's Hybrids; Atlanta, IN 46031, USA.

[9]These authors contributed equally to this work.

25

26　*Correspondence: maj@purdue.edu (J.M.), fengxianzhong@iga.ac.cn (X.F.)

**Abstract**

The origination and functionality of long non-coding RNAs (lncRNAs) remain poorly understood. Here, we show that multiple quantitative trait loci modulating distinct domestication traits in soybeans are pleiotropic effects of a locus composed of two tandem lncRNA genes. These lncRNA genes, each containing two inverted repeats (IRs) originated from coding sequences of MYB genes, function by generating clusters of small RNAs in wild soybeans to inhibit the expression of their MYB gene relatives through posttranscriptional regulation. In contrast, the expression of the lncRNA genes in cultivated soybeans is severely repressed, and consequently, the corresponding MYB genes are highly expressed, shaping multiple distinct domestication traits as well as leafhopper resistance. The IRs were formed before the divergence of the Glycine genus from the Phaseolus/Vigna lineage and exhibit strong structure-function constraints. This study exemplifies a new type of targets for selection during plant domestication and uncovers mechanisms of lncRNA formation and action.

## Main

The domestication of a crop from its wild relative is a complex process of artificial selection for a suite of favorable traits, which are generally controlled by different genetic loci[1]. Such a process creates a new form of plants, known as domesticates, to meet human needs. Nevertheless, it also leads to drastic reduction in genetic diversity in domesticates, hindering the sustainability of crop improvement[2]. To better understand the dynamic processes of crop domestication and exploit untapped genetic variation in crop wild relatives for enhancement of elite cultivars, it is important to decipher the genetic and molecular basis underlying domestication-related traits (DRTs).

In the past few decades, tremendous work has been done to identify quantitative trait loci (QTL) underlying DRTs in major crops, such as (cultivated) soybean (*Glycine max*) – an economically important leguminous crop domesticated from wild soybean (*Glycine soja*)[3]. Most wild soybean accessions exhibit a procumbent or climbing growth habit, with long, slender, prolifically branched stems and small leaves that grow appressed pubescence, whereas majority of cultivated soybean varieties display a bush-type upright growth habit, with short, scout primary stems and sparse branches and large leaves with semi-appressed or erect pubescence. Here, we report that multiple QTL underlying different DRTs as well as resistance to leafhoppers in cultivated soybeans are resulted from artificial selection of reduced expression of two tandemly duplicated long non-coding RNA (lncRNA) genes

each carrying MYB gene coding sequence-derived inverted repeats (IRs), which have undergone strong purifying selection in the Glycine genus.

**Results**

**Map-based cloning of multiple DRT QTLs identifies a single locus with pleiotropic effects**

Using a subset of the 2,287 recombinant inbred lines (RILs) derived from a cross between soybean cultivar Williams 82 (Wm82) and *G. soja* accession PI 479752, we initially mapped >100 QTL associated with various DRTs[4]. Remarkably, many of the QTL regions, which underlie different DRTs, physically overlap. One such region, *qDRT12.3* on chromosome 12, was found to harbor five QTL, *qPB-12*, *qMSL-12*, *qLSZ-12*, *qGH-12,* and *qST-12*, which explained 63.3%, 25.0%, 23.0%, 14.8% and 6.4% the phenotypic variation in pubescence form, main stem length, leaf size, growth habit, and stem-twining, respectively (Fig. 1a-1e). To determine whether these QTL are attributed to different genes or pleiotropic effects of the same gene, or both, we first conducted fine mapping of three (*qPB-12*, *qMSL-12,* and *qLSZ-12*) of the five QTL, independently, using the entire RIL population. Two insertion/deletion (InDel) markers, M1 and M10, which initially defined the boundaries of the *qDRT12.3* region, were used to genotype all 2,287 RILs and identified 238 recombinants between the two markers (Fig. 1f). These recombinants were then genotyped with eight additional markers within the *qDRT12.3* region and first examined for pubescence form. Combination of the genotypic and

89    phenotypic data delimited *qPB-12* to a 29-kb region between markers M5 and

90    M7 (Fig. 1f), which was echoed by a genome-wide association study using re-

91    sequencing data from 74 *G. soja* and 594 *G. max* accessions[5] (Extended Data

92    Fig. 1a and 1b). Subsequently, the 238 recombinants were measured for main

93    stem length and leaf size, respectively. Based on the eight markers, these

94    recombinants were divided into 13 haplotypes, and the average phenotypic

95    value of recombinants within each haplotype group was compared to the

96    population mean to calculate the phenotypic scores of individual haplotypes

97    to fine map *qMSL-12* and *qLSZ-12*. Interestingly, these two QTL were also

98    defined to the same 29-kb region (Fig. 1g-1h). According to the Wm82

99    reference genome, this region harbors only two genes, Glyma.12G213800

100   and Glyma.12G213900, both lncRNAs.

101   It has been observed that semi-appressed or erect pubescence is linked to

102   reduced defoliation caused by Cicadellidae insects[6]. To investigate whether

103   *qPB-12* is responsible for such resistance, we conducted a genome-wide

104   association study (GWAS) on leafhopper resistance/susceptibility using the

105   phenotypic and genotypic data from 784 accessions in the USDA soybean

106   germplasm collection[7]. We found that molecular markers within the fine

107   mapped *qPB-12* region were significantly associated with leafhopper

108   resistance (Extended Fig. 1c and 1d) and that erect pubescence indeed

109   contributes to the leafhopper resistance as shown in Supplementary Movies

110   1 and 2. In a set of re-sequenced diverse *G. soja* and *G. max* accessions chosen

111   from the USDA soybean germplasm collection[5], only 13.4% of the *G. soja*

accessions have erect pubescence, whereas 71.3% and 96.7% of the landraces and elite cultivars possess it, respectively (Extended Data Fig. 1e), indicating that erect pubescence and its underlying leafhopper resistance was a target for selection during soybean domestication and improvement. Artificial selection at this locus was also echoed by a selective sweep surrounding it (Extended Fig. 1f), as detected by sequencing data from 103 *G. soja* accessions and 328 landraces[5]. Collectively, these observations suggest that Glyma.12G213800 and Glyma.12G213900 are the candidate genes regulating pubescence form, main stem length, leaf size, as well as leafhopper resistance attributed to erect pubescence.

## The pleiotropic DRT locus is composed of two tandemly duplicated lncRNA genes, *lncRG1* and *lncRG2*

The genes, Glyma.12G213800 and Glyma.12G213900, in Wm82 produce 1,526-nt and 1,565-nt transcripts, which are predicted to encode 37 and 49 amino acids, respectively. Thus, they are defined as long non-coding RNA (lncRNA) genes, referred to as *lncRG1* and *lncRG2*. Both *lncRG1* and *lncRG2* are primarily expressed in stems, leaves, and stem tips of PI 479752 at the vegetative 1 (V1) developmental stage when the first trifoliate leaflets are fully expanded, but are expressed at significantly lower levels in the same tissues of Wm82, as measured by quantitative reverse transcription-PCR (qRT-PCR) (Fig. 2a). RNA-seq data from 45 highly diverse soybean accessions[8] revealed significantly higher expression levels of these two genes

in nine wild soybean accessions than in 36 cultivated soybean accessions (Extended Data Fig. 1g) as well as a pattern of *lncRG1* and *lncRG2* co-expression (Extended Data Fig. 1h). Therefore, the suppressed expression of *lncRG1* and *lncRG2* is most likely to be responsible for the observed phenotypic changes from wild soybeans to cultivated soybeans.

Comparison of *lncRG1* and *lncRG2* with all other soybean genes in the soybean genome reveals that not only the putative coding sequences (CDSs) but the large portions of the non-CDSs of these two lncRNA genes share similarities with typical MYB transcription factor genes (Fig. 2b and 2c). Thus, *lncRG1* and *lncRG2* were derived from MYB genes. Further phylogenetic analysis reveals that *lncRG1* and *lncRG2* were tandemly duplicated before the latest whole genome duplication (WGD) event (Fig. 2b) predicted to have occurred in soybean ~13 million years ago (MYA)[9]. As a result, there were two homologs of *lncRG1* and *lncRG2*, dubbed *lncRG3* and *lncRG4,* respectively (Fig. 2c). Nevertheless, *lncRG3* and *lncRG4* are not associated with any of the domestication QTLs[4]. Interestingly, all four lncRGs in soybean possess IRs, each at ~300-400 bp, corresponding to the third exon of their most closely related MYB genes (Fig. 2c).

**LncRG1 and lncRG2 harbor IRs and produce abundant sRNAs primarily targeting three closely related MYB genes**

Based on prediction, the IRs within the transcripts of *lncRG1* and *lncRG2* may form double-stranded stem loops at 453 bp and 337 bp, respectively (Fig. 2d

158 and 2e), which could be processed to generate small RNAs (sRNAs), such as

159 microRNAs, microRNA (miRNA)-like sRNAs, or short interfering sRNAs

160 (siRNAs). Then, we sequenced sRNAs in the V1-stage stem tips of PI 479752

161 and WM82, respectively. Abundant, overlapping sRNAs, mainly at 21-23

162 nucleotides (nt), across the IRs of both *lncRG1* and *lncRG2* were detected in

163 PI 479752, but their relative abundances vary drastically (Fig. 2f and 2g). The

164 most abundant sRNAs from *lncRG1* are at 23nt, whereas the most abundant

165 sRNAs from *lncRG2* are at 21nt (Fig. 2h and 2i). Nevertheless, much more

166 sRNAs were produced from *lncRG2* than *lncRG1*. Consistent abundances and

167 distribution patterns of the sRNAs produced by *lncRG1* and *lncRG2* were

168 observed in a pair of RILs, RIL186 (*qdrt12.3*) and RIL 334 (*qDRT12.3*)

169 (Extended Data Fig. 2a-2d), suggesting that the abundance of individual

170 sRNAs were tightly regulated and not randomly produced from the IRs.

171 A total of 163 genes were predicted to be targets of 27 distinct sRNAs from

172 *lncRG1* and *lncRG2*, with a relative abundance of >100 copies per million

173 (CPM) sRNA reads (Supplementary Table 1 and 2). Of these putative targets,

174 only Glyma.01G051700, Glyma.02G110000 and Glyma.02G110100 showed

175 significantly reduced levels of expression in PI 479752 compared with Wm82,

176 with at least 2-fold changes in stem tips, stems and leaves as determined by

177 RNA-seq and qRT-PCR (Fig. 2j and Supplementary Table 3). Degradome

178 sequencing revealed that the mRNAs of these three genes were

179 predominantly cleaved at the predicted sRNA target sites in PI 479752 (Fig.

180 2k-2m). Interestingly, all three targets are typical MYB genes that are most

181 closely related to *lncRG1* and *lncRG2* based on the phylogenetic relationships

182 established with the predicted coding sequences (Fig. 2b). Thus, these MYB

183 gene-derived *lncRG1* and *lncRG2* are likely to modulate the DRTs by

184 producing plentiful sRNAs to primarily repress their MYB gene relatives via

185 post-transcriptional regulation.

186

187 **Overproduction of sRNAs in cultivated soybean promotes the wild**

188 **soybean-type phenotypes**

189 To determine whether the sRNAs produced by *lncRG1* and *lncRG2* underlie

190 the DRTs, we first generated Williams 82 transgenic lines that overexpress

191 the "stem-loop" part of each gene by the cauliflower mosaic virus (CaMV) 35S

192 promoter. The transgenic lines displayed elevated abundance of sRNAs from

193 the stem loops (Extended Data Fig. 2e and 2f) and showed expected

194 phenotypic changes including appressed pubescence form, decreased plant

195 height and smaller leaf size in comparison to the Wm82 (Fig. 3a-3c). In

196 addition, we constructed two artificial miRNA precursors (aMIR-sRlncRG1-1

197 and aMIR-sRlncRG2-3) by replacing the miR172a and miR172a* sequences

198 from the soybean miR172a precursor MIR172a with sRlncRG1-1 and its

199 complementary sRlncRG1-1* or with sRlncRG2-3 and its complementary

200 sRlncRG2-3*, respectively. Overexpression of the two artificial miRNA

201 precursors using the 35S promoter in Williams 82 resulted in appressed/semi-

202 appressed pubescence form, reduced plant height, and smaller leaf size

203 compared to the Wm82 (Fig. 3d and 3e). As expected, these transgenic lines

exhibited increased expression levels of the corresponding artificial sRNAs and decreased expression levels of the three MYB genes as determined by stem-loop and regular qRT-PCR, respectively (Fig. 3g and 3h). The mRNAs of the target genes were confirmed to be principally cleaved at the predicted sRlncRG1-1 and sRlncRG2-3 cleavage sites in the transgenic lines, but such cleavages were not detected in the wild-type control using RNA ligase mediated rapid amplification of the 5' cDNA ends (RLM-RACE) technique followed by deep sequencing (Fig. 3i and g). These observations indicate that the specific sRNAs produced from *lncRG1* and *lncRG2* are responsible for forming the DRTs and suggest that these sRNAs use miRNA-like mechanism to repress their targets.

Since *lncRG1* and *lncRG2* are predicted to encode two small peptides, we wonder whether the small peptides also contribute to the DRTs. Then, we generated Williams 82 transgenic lines that overexpress the predicted coding sequence for the small peptide of each gene by the 35S promoter. No phenotypic differences between any of the transgenic lines and the negative controls were observed, suggesting that the predicted coding sequences are unlikely to modulate the DRTs.

**The three MYB genes targeted by the sRNAs exhibit functional redundancy and divergence**

To gain insights into the mechanism by which the three MYB genes regulate the DRTs, we generated Wm82 knockout lines for each MYB gene using

11

CRISPR-Cas9 (Extended Data Fig. 3a-3c). Knocking out any of the three genes resulted in appressed/semi-appressed pubescence, and reduced plant height, and smaller leaf size; however, their effects on each DRT slightly vary (Fig. 4a, 4d and 4e). We then crossed the knockout lines for different MYB genes to generate double mutants, which were further crossed to create triple mutants. Overall, the double and triple mutants exhibited stronger phenotypic changes compared to the single mutants (Fig. 4a-4c and 4f-4i), suggesting an additive effect of the three MYB genes.

Given that protein dimerization often plays a crucial role in transcription factor activity, we wondered whether the three MYB genes enable homo- or hetero-dimerization. Using the yeast two-hybrid (Y2H) system (Fig. 4j and 4k) and the bimolecular fluorescence complementation (BiFC) assay in tobacco leaves (Fig. 4l and 4m), both self- and pairwise-protein-protein interactions were detected among the three MYB genes. As expected, both homo- and hetero-dimers were localized in the nucleus (Fig. 4l and 4m). Furthermore, the three target MYB genes were detected to be able to interact with their more ancestral MYB genes (Fig. 2b), such as Glyma.07G228600, Glyma.20G032900, and Glyma.04G166900; however, the strengths of the interactions involving each of the three target MYB genes vary (Extended Data Fig. 3d). These observations suggest that the three target MYB genes possess both redundant and diverged functions.

**The lncRGs have undergone purifying selection due to the structure-function constraints**

To track the origin and evolutionary variation of the lncRGs, we compared the mapped *lncRG1* and *lncRG2* region and its flanking regions of *G. max* and *G. soja* with the corresponding orthologous regions in seven additional leguminous species belonging to the Phaseolus, Vigna, and Cajanus genera using *Medicago truncatula* as an outgroup. It appears that the tandem duplication event occurred after the divergence of Glycine and Phaseolus/Vigna from a common ancestor ~20 MYA[10,11] (Fig. 5a and 5b). The IRs were also seen in Phaseolus and Vigna, but not seen in Cajanus and *M. truncatula*, suggesting that the IRs were formed before the divergence of Glycine from Phaseolus/Vigna but after its divergence from Cajanus ~20-24MYA[10,11] (Fig. 5a and 5b). The IRs of *lncRG1* and *lncRG2* in *G. soja* and *G. max* exhibited the lowest level of divergence compared with the IRs in the orthologs of *lncRG1*/*lncRG2* in Phaseolus/Vigna (Fig. 5c), indicating that the IRs, as the functional parts of the *lncRG1* and *lncRG2* gene bodies, have experienced strong "purifying selection". It is also noticeable that *lncRG2*, which produced more non-redundant and more abundant sRNAs than *lncRG1* in PI 479752, evolved in a slower pace than *lncRG1*.

**The lncRG-derived sRNAs exhibit diverse distribution patterns at the population level**

13

271    The availability of sRNA sequencing data from 45 *G. soja* and *G. max*

272    accessions[8] allowed us to compare the distribution and relative abundance of

273    sRNAs generated by *lncRG1* and *lncRG2* at the population level (Fig. 5d, 5e,

274    Extended Data Fig. 4a, 4b and Supplementary Table 4). As expected, all the

275    nine *G. soja* accessions and a cultivated soybean accession (Jin Dou No. 23)

276    with appressed pubescence produced abundant sRNAs from *lncRG1* and

277    *lncRG2.* In contrast, few sRNAs were produced from *lncRG1* and *lncRG2* in

278    the remaining 35 cultivated soybean accessions, which possess erect

279    pubescence. Remarkably, the sRNA distribution patterns vary drastically

280    among the 10 accessions with depressed pubescence, and in most cases,

281    different sRNAs are predicted to target the three MYB genes, and up to 41%

282    of the predicted sRNA targets in one accession are not shared by another

283    accession (Supplementary Table 5). As observed in PI 479752 (Fig. 2g),

284    *lncRG2* in each of the 10 accessions produced more non-redundant and more

285    abundant sRNAs than *lncRG1*, with 21-nt and then 22-nt sRNAs as the

286    predominant forms (Extended Data Fig. 4a and 4b).

287

288    **Discussion**

289    LncRNAs are ubiquitously present in eukaryotes and play important roles in

290    regulating gene expression[12]. However, how they are originated and execute

291    their functions remains largely unknown. In this study, we demonstrate that

292    two lncRNA tandem duplicates, *lncRG1* and *lncRG2*, were derived from MYB

293    genes and underwent exonic sequence rearrangement to form IRs. Intragenic

294  IRs are typically lost due to their instability and fitness costs[13]; yet the IRs in

295  *lncRG1* and *lncRG2* have been maintained over the course of 20-24 million

296  years of evolution (Fig. 5), likely due to their crucial role in regulating

297  multiple "wild" adaptive traits in Glycine. IRs can be induced by DNA

298  replication repair or transposable elements[14], which usually leave specific

299  sequence features surrounding the IR junctions. However, as these features

300  do not generally bring any fitness benefits, they would not be preserved over

301  such a long period of evolutionary time. While the processes leading to the

302  formation of the IRs in *lncRG1* and *lncRG2* remain unclear, it is possible that

303  these are newly emerged or incipient miRNA precursors that hasn't yet been

304  selected for production of a single, precisely processed duplex and instead

305  are generating siRNAs from the IR precursor[15]. The IR structures are still

306  detectable across the Phaseolus, Vigna, and Glycine genera, reflecting their

307  functional constraints at variable levels. Given such a great variation in

308  relative abundance and distribution of the *lncRG1*- and *lncRG2*-derived

309  sRNAs among different wild soybean accessions, the functional constraints

310  in soybean may be implemented through purifying selection across the entire

311  IR regions. It would be interesting to explore whether the IRs in other

312  legumes have similar functionality and regulate comparable traits, and

313  whether the IRs were also targeted for selection during domestication of

314  other leguminous crops.

315  A few domestication genes have been shown to exhibit pleiotropic effects on

316  multiple traits[2], such as *TEOSINTE BRANCHED1* in maize, which controls

branching, inflorescence architecture, and plant height[16], and *PROSTRATE GROWTH1* in rice, which controls tiller angle, panicle size, and seed shattering[17]. Compared to these genes, the mechanism by which *lncRG1* and *lncRG2* execute their pleiotropic effects is unique and reflective of evolutionary innovation triggered by varied types of duplications events including exonic duplication, genic duplication, and WGD. In soybean, approximately 75% of the genes existing in multiple copies, which were primarily generated via two rounds of WGD events that occurred 59 and 13 MYA[9]. Consequently, mutations within a single gene can often be "rescued" by its functionally redundant duplicates. In such a case, phenotypic transition of a DRT during soybean domestication would have involved artificial selection of mutations within two or more duplicated genes. As the sRNAs produced by *lncRG1* and *lncRG2* enable simultaneous repression of multiple duplicated MYB genes and most likely additional genes as well, artificial selection of the DRTs regulated by these genes was achieved simply by selecting the reduced expression of *lncRG1* and *lncRG2* within a single locus producing fewer sRNAs.

The causal mutations for the reduced expression of *lncRG1* and *lncRG2* in cultivated soybeans remain unknown. Genome-wide association analysis with the re-sequencing data from 74 *G. soja* and 596 *G. max* accessions[5] revealed numerous polymorphic sites across the entire mapping region that are highly associated with the phenotypic differences in pubescence form (Extended Data Fig. 1a and 1b), but no single polymorphic sites in the putative

promoters of the two genes or other parts of the region could explain the phenotypic differences better than the others. This is not unexpected, given that the entire region has undergone selective sweep (Extended Data Fig. 1f). Because *lncRG1* and *lncRG2* are co-expressed across different tissues and developmental stages, there is a possibility that these two genes are regulated by the same regulatory element(s) within the mapped 29-kb region. Under this caveat, extensive functional assays are needed to pinpoint the causal mutation(s) for reduced *lncRG1* and *lncRG2* expression.

While sRNAs may also repress translation without cleaving mRNAs[18], it is unclear whether the remaining 160 predicted sRNA targets, which show no difference in expression levels between Wm82 and PI 479752 (Supplementary Table 2 and 3), are directly regulated by the sRNAs from *lncRG1* and *lncRG2* through translation inhibition. Given the fact that the three MYB targets also interacts with additional, more diverged copies of MYB genes, that the predominant sizes of sRNAs produced from lncRG1 and lncRG2 are different, and that the sRNAs from *lncRG1* and *lncRG2* and their putative targets are highly variable among different accessions, the pleiotropic effects of *lncRG1* and *lncRG2* and the mechanisms by which they execute their full suite of functions are likely to be more extensive than what has been detected.


**Methods**

**Plant materials**

17

363 The mapping population consisted of 2,287 $F_{6:7}$ recombination inbred lines

364 (RIL) derived from a cross between *G. max* (Wm82) and *G. soja* (PI 479752).

365 The association mapping population for leafhopper resistance were sourced

366 from the USDA soybean germplasm collection (https://www.ars-grin.gov/).

367 Wm82 was used for stable transformation and genome editing; *Nicotiana*

368 *benthamiana* was used for the BiFC assays.

369

**QTL and association mapping**

371 The QTL mapping was performed using composite interval mapping (CIM)

372 method[19] incorporated in the r/qtl package[20]. The phenotypic data for

373 association mapping were downloaded from the USDA National Plant

374 Germplasm System (https://npgsweb.ars-grin.gov/) and the SoySNP50K data

375 were obtained from a previous study[7]. The re-sequencing data were from the

376 soybean pan-genome study[5]. The association mapping was performed using

377 TASSEL 5[21] with a mixed linear model[22].

378

**Recombinants genotyping and phenotyping**

380 All the mapping markers were designed based on the re-sequencing data of

381 PI 479752 from a previous study[23]. The domestication-related traits were

382 examined for all the recombinants in the field at Purdue Agronomy Center for

383 Research & Education in 2018. All the primers used in this study were listed

384 in Supplementary Table 6.

385

**Transgene constructs**

387 For stem loop over-expression, the stem loops of *lncRG1* and *lncRG2* were

388 amplified from genomic DNA of PI 479752. The PCR products and linearized

389 vector were purified using PurLink[TM] Quick Gel Extraction Kit (K210012,

390 ThermoFisher Scientific). The stem loops were inserted into the plasmid

391 vector, linearized by restriction enzymes Nco I and Xba I, using ClonExpress

392 II One Step Cloning Kit (C112, Cellagen Technology).

393 For artificial miRNA over-expression, the soybean MIR172a was used as the

394 backbone following a previous protocol[24]. The miR172a/miR172a* sequences

395 were replaced by sRlncRG1-1 and sRlncRG2-3 and their corresponding

396 reverse complementary sequences. The forward sequence and

397 complementary sequence were annealed to form dimmers and inserted into

398 pPTN1171.

399 For CRISPR-Cas9 editing, 4 sgRNAs were designed for each target gene

400 using CRISPR-P, a web-based guided RNA design tool[25]. The primer pairs

401 were annealed for 5 minutes at 95 °C and then cool down to form dimers. The

402 dimers were inserted into pGEL201, linearized by restriction enzymes Bsa I,

403 vector[26]. During transformation, four agrobacteria with different sgRNA

404 were equally mixed before infection.

405 For yeast two hybrid assays, the full-length coding sequences of the MYB

406 genes were cloned from the cDNA sample of 'Wm82' and then inserted into

407 the vectors pGBKT7 and pGADT7.

408     For bimolecular fluorescence complementation assay, the full-length coding

409     sequences of the three target genes were amplified and cloned into plasmids

410     pCNHP-neYFP-C and pCNHP-ceYFP-C, which express fusion proteins with

411     either N-terminal half of eYFP (neYFP) or C-terminal half of eYFP (ceYFP) at

412     their N-terminus, respectively.

413

414     **Soybean transformation**

415     Mature seeds from soybean cultivar 'Williams 82' were disinfected using

416     chlorine gas for 12 hours. The disinfected seeds were soaked in distilled

417     water for 12 hours at room temperature at dark. Half-seeds were soaked in

418     resuspended agrobacterium liquid co-cultivation medium (OD650 = 0.6,

419     3.21g/L Gamborg B-5 basal medium, 30g/L sucrose, 3.9g/L MES, 0.4g/L L-

420     cystine, 0.1542g/L DTT, 0.25mg/L GA3, 1.67mg/L 6-BA and 0.3924g/L

421     acetosyringone, pH = 5.4) for 30 minutes. After infection, the explants were

422     transferred to solid co-cultivation medium. The plates were sealed with

423     Micropore tape (Catelog #1530-0, 3M, St. Paul, MN) and incubated in the

424     dark at 21 ℃ for 4 days. After co-cultivation, explants were inserted into shoot

425     induction medium plate (3.21g/L Gamborg B-5 basal medium, 30g/L sucrose,

426     0.59g/L MES, 0.25g/L timetin, 0.1 g/L cefradine, 1.67mg 6-BA, 2.5mg/L

427     glufosinate, pH = 5.7, 2g/L gellan gum powder). Shoot induction was carried

428     out at 26 ℃ with a photoperiod of 18 hours and a light intensity of 40-70

429     µE/m2/s. After 4 weeks, the induced shoots were cut from cotyledons and

430     transferred to shoot elongation medium (4.43g/L Murashige & Skoog

modified medium with Gamborg vitamins, 30g/L sucrose, 0.59g/L MES, 0.25g/L timentin, 0.1g/L cefradine, 0.05g/L asparagine, 0.05g/L glutamine , 0.5mg/L GA3, 0.1mg/L IAA, 1mg/L zeatin, 5mg/L glufosinate,  pH = 5.7, 2g/L gellan gum powder) under same temperature and photoperiod. After 2-4 weeks in shoot elongation medium, the glufosinate-resistant shoots were cut and transferred to rooting medium (4.43g/L Murashige & Skoog modified medium with Gamborg vitamins, 30g/L sucrose, 0.59g/L MES, 0.05g/L asparagine, 0.05g/L glutamine, 0.1mg/L IBA, pH = 5.7, 3g/L gellam gum) for further shoot and root elongations. After root grows longer than 1 cm, plants were transferred to moistened Berger BM2 soil (Berger, Saint-Modeste, QC, Canada), and kept enclosed in clear plastic tray in a growth chamber at 26 ℃ with a 16-hour photoperiod at 250 -350 µE/m2/s.

**Genotyping the transgenic and genome editing lines**

Genomic DNA was extracted from $T_0$, $T_1$, and $T_2$ plants. The presence of the transgenes in the transgenic plants was confirmed by PCR with primers specific to the vector and the corresponding transgene. Expression of the transgene were monitored by qRT-PCR or stem loop qRT-PCR for the sRNA. For genome editing lines, the target genes were amplified and sequenced to confirm the presence of frameshift mutation.

**RNA extraction, regular qRT-PCR and stem loop RT-PCR**

Total RNA was extracted using the TRIzol reagent (Cat. # 15596018, Invitrogen). 2 µg DNA-free RNA was used to synthesize cDNA with the

455  Promega M-MLV Reverse Transcriptase (Cat. # M1701, Promega). qRT-PCR

456  was performed using Applied Biosystems™ Power SYBR™ Green PCR Master

457  Mix (Cat. # 4368577, Applied Biosystems) on an Applied Biosystems

458  StepOnePlus$^{TM}$ Real-Time PCR Systerm (Cat. # 4376600, Applied

459  Biosystems). Stem-loop RT-PCR was used to examine the expression levels of

460  miRNAs following a previous protocol[27].

461

**mRNA, small RNA and Degradome sequencing**

463  The cleaned RNA-seq reads were mapped to the soybean reference genome[9]

464  using STAR (v2.5.4b) with only unique mapped reads kept[28]. The expression

465  levels (FPKM) were calculated using the cuffnorm function in cufflinks

466  (v2.2.1)[29]. Degradome libraries were constructed and sequenced at

467  Novogene Corporation Inc. (Sacramento, CA). The potential target genes of

468  miRNA produced by *lncRG1* and *lncRG2* were analyzed using CleaveLand

469  with the following parameters, -r 0.6 and -c 2. (v4.5)[30].

470

**RNA ligase-mediated 5′ rapid amplification of cDNA ends (5' RLM-RACE)**

473  RLM-RACE was performed following the protocol described previously[24]. The

474  mRNAs were then ligated with 5′ RACE oligo adaptors for reverse

475  transcription using the GeneRacer kit (Cat. # L150202, ThermoFisher

476  Scientific) followed by nested PCR. The purified PCR products were

477  sequenced by using the WideSeq method

478  (https://www.purdue.edu/hla/sites/genomics/wideseq-2/).

479

**Phylogenetic analysis and nucleotide diversity calculation**

Sequence alignments and tree construction were performed using the Maximum Likelihood method[31] in MEGA7[32]. Nucleotide diversity was calculated using vcftools (v0.1.16)[33].

484

**RNA secondary structure prediction**

The secondary structures of lncRG1 and *lncRG2* were predicted using the RNAfold server incorporated in the ViennaRNA Web Services (http://rna.tbi.univie.ac.at/).

489

**miRNA target prediction**

Potential targets by the miRNAs from *lncRG1* and *lncRG2* were predicted using the online tool psRNATarget (https://www.zhaolab.org/psRNATarget/, Schema V2 2017 release) with the expectation cutoff set as 2.5[34].

494

**Yeast Two-Hybrid (Y2H) assays**

Y2H assays were performed using the Matchmaker Gold Yeast Two-Hybrid System (Cat. # 630489, Takara Bio USA Inc). Different combinations of the constructs were co-transformed into the yeast strain Y2H Gold. The transformed yeast cells were spread on SD (–Trp/–Leu) medium. The plates were incubated at 30 °C for 3-5 days. 5-10 colonies were picked from each plate and resuspended in 0.9% (w/v) NaCl solution. Then the yeast cells were spotted on SD (–Trp/–Leu/–Ade/–His) selection medium. Plates were

503  incubated at 30°C for 3 days to observe yeast growth. pGADT7-T + pGBKT7-

504  53 was used as positive control; pGADT7-T + pGBKT7-Lam was used as

505  negative control.

506

507  **Bimolecular fluorescence complementation (BiFC)**

508  Different constructs were transformed into Agrobacterium tumefaciens

509  strain EHA105. The agrobacterium suspension was injected into the abaxial

510  surface of 4–6-week-old *Nicotiana Benthamiana* leaves with a needleless

511  syringe. Plasmid expressing mCherry-labeled Pentunia hybrida's histone H1-

512  3 (acted as the nuclear marker) was co-infiltrated with the expression

513  construct of each target gene. 72 h after infiltration, the fluorescent signals

514  in detached leaves were imaged using a Zeiss LSM-880 laser-scanning

515  confocal microscope.

516

517  **Data and code availability**

518  All data are available in the main text, supplemental materials, public

519  databases, or referenced studies. All the raw sequence data generated in this

520  study have been deposited in NCBI database under the BioProject

521  PRJNA876203. This paper does not report original code.

522

**Author contributions**

JM and XZF designed the research. WW, JD, XF, XW LC, CBC, SAS, RLN, SL and JW performed the experiments. WW, XW, BCM and JM analyzed the data. WW and JM wrote the manuscript and BCM edited the manuscript.

**Declaration of interests**

The authors declare no competing interests.

**References**

1.  Olsen, K.M. & Wendel, J.F. A bountiful harvest: genomic insights into crop domestication phenotypes. *Annual review of plant biology* **64**, 47-70 (2013).

2.  Doebley, J.F., Gaut, B.S. & Smith, B.D. The molecular genetics of crop domestication. *Cell* **127**, 1309-1321 (2006).

3.  Sedivy, E.J., Wu, F. & Hanzawa, Y. Soybean domestication: the origin, genetic architecture and molecular bases. *New Phytologist* **214**, 539-553 (2017).

4.  Swarm, S.A. *et al.* Genetic dissection of domestication-related traits in soybean through genotyping-by-sequencing of two interspecific
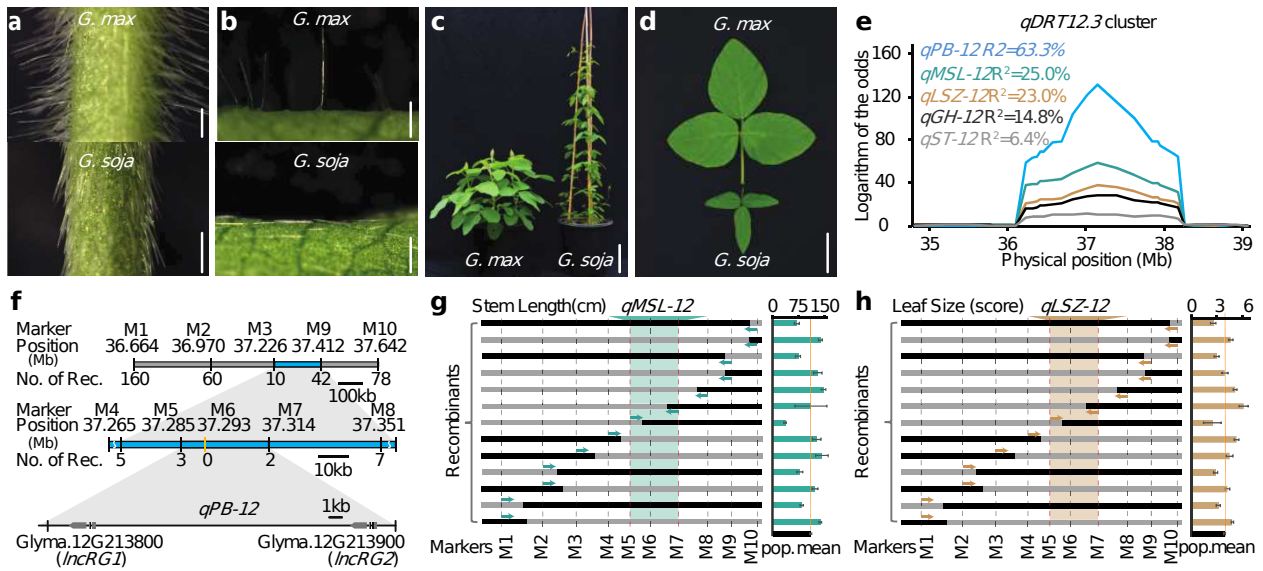
mapping populations. *Theoretical and Applied Genetics* **132**, 1195-1209 (2019).

5. Liu, Y. *et al.* Pan-genome of wild and cultivated soybeans. *Cell* **182**, 162-176. e13 (2020).

6. Broersma, D., Bernard, R. & Luckmann, W. Some effects of soybean pubescence on populations of the potato leafhopper. *Journal of Economic Entomology* **65**, 78-82 (1972).

7. Song, Q. *et al.* Fingerprinting soybean germplasm and its utility in genomic research. *G3: Genes, genomes, genetics* **5**, 1999-2006 (2015).

8. Shen, Y. *et al.* DNA methylation footprints during soybean domestication and improvement. *Genome biology* **19**, 1-14 (2018).

9. Schmutz, J. *et al.* Genome sequence of the palaeopolyploid soybean. *nature* **463**, 178-183 (2010).

10. Choi, H.-K. *et al.* Estimating genome conservation between crop and model legume species. *Proceedings of the National Academy of Sciences* **101**, 15289-15294 (2004).

11. Zheng, F. *et al.* Molecular phylogeny and dynamic evolution of disease resistance genes in the legume family. *BMC genomics* **17**, 1-13 (2016).

12. Statello, L., Guo, C.-J., Chen, L.-L. & Huarte, M. Gene regulation by long non-coding RNAs and its biological functions. *Nature reviews Molecular cell biology* **22**, 96-118 (2021).

574    13.    Parniske, M. *et al.* Novel disease resistance specificities result from

575           sequence exchange between tandemly repeated genes at the Cf-4/9

576           locus of tomato. *Cell* **91**, 821-832 (1997).

577    14.    Reams, A.B. & Roth, J.R. Mechanisms of gene duplication and

578           amplification. *Cold Spring Harbor perspectives in biology* **7**, a016592

579           (2015).

580    15.    Bradley, D. *et al.* Evolution of flower color pattern through selection

581           on regulatory small RNAs. *Science* **358**, 925-928 (2017).

582    16.    Doebley, J., Stec, A. & Hubbard, L. The evolution of apical dominance

583           in maize. *Nature* **386**, 485-488 (1997).

584    17.    Tan, L. *et al.* Control of a key transition from prostrate to erect growth

585           in rice domestication. *Nature genetics* **40**, 1360-1364 (2008).

586    18.    Fabian, M.R. & Sonenberg, N. The mechanics of miRNA-mediated

587           gene silencing: a look under the hood of miRISC. *Nature structural &*

588           *molecular biology* **19**, 586-593 (2012).

589    19.    Zeng, Z.-B. Precision mapping of quantitative trait loci. *Genetics* **136**,

590           1457-1468 (1994).

591    20.    Broman, K.W., Wu, H., Sen, Ś. & Churchill, G.A. R/qtl: QTL mapping in

592           experimental crosses. *Bioinformatics* **19**, 889-890 (2003).

593    21.    Bradbury, P.J. *et al.* TASSEL: software for association mapping of

594           complex traits in diverse samples. *Bioinformatics* **23**, 2633-2635

595           (2007).

596  22.  Yu, J. *et al.* A unified mixed-model method for association mapping

597       that accounts for multiple levels of relatedness. *Nature genetics* **38**,

598       203-208 (2006).

599  23.  Zhou, Z. *et al.* Resequencing 302 wild and cultivated accessions

600       identifies genes related to domestication and improvement in

601       soybean. *Nature biotechnology* **33**, 408-414 (2015).

602  24.  Ren, B., Wang, X., Duan, J. & Ma, J. Rhizobial tRNA-derived small

603       RNAs are signal molecules regulating plant nodulation. *Science* **365**,

604       919-922 (2019).

605  25.  Lei, Y. *et al.* CRISPR-P: a web tool for synthetic single-guide RNA

606       design of CRISPR-system in plants. *Molecular plant* **7**, 1494-1496

607       (2014).

608  26.  Bai, M. *et al.* Generation of a multiplex mutagenesis population via

609       pooled CRISPR⎕Cas9 in soya bean. *Plant Biotechnology Journal* **18**,

610       721-731 (2020).

611  27.  Chen, C. *et al.* Real-time quantification of microRNAs by stem–loop

612       RT–PCR. *Nucleic acids research* **33**, e179-e179 (2005).

613  28.  Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner.

614       *Bioinformatics* **29**, 15-21 (2013).

615  29.  Trapnell, C. *et al.* Differential gene and transcript expression analysis

616       of RNA-seq experiments with TopHat and Cufflinks. *Nature protocols*

617       **7**, 562-578 (2012).

618    30.    Addo-Quaye, C., Miller, W. & Axtell, M.J. CleaveLand: a pipeline for

619            using degradome data to find cleaved small RNA targets.

620            *Bioinformatics* **25**, 130-131 (2009).

621    31.    Tamura, K. & Nei, M. Estimation of the number of nucleotide

622            substitutions in the control region of mitochondrial DNA in humans

623            and chimpanzees. *Molecular biology and evolution* **10**, 512-526

624            (1993).

625    32.    Kumar, S., Stecher, G. & Tamura, K. MEGA7: molecular evolutionary

626            genetics analysis version 7.0 for bigger datasets. *Molecular biology*

627            *and evolution* **33**, 1870-1874 (2016).

628    33.    Danecek, P. *et al.* The variant call format and VCFtools.

629            *Bioinformatics* **27**, 2156-2158 (2011).

630    34.    Dai, X., Zhuang, Z. & Zhao, P.X. psRNATarget: a plant small RNA

631            target analysis server (2017 release). *Nucleic acids research* **46**, W49-

632            W54 (2018).

**Figure legends**



**Fig. 1: Map-based cloning of multiple DRT QTLs identifies a single locus with pleiotropic effects. a-b**, Comparisons of pubescence form on stems (a) and leaves (b) between *G. max* and *G. soja*. Scale bar = 3 mm. **c**, Comparisons of stem height and growth habit between *G. max* and *G. soja*. Scale bar = 10 cm. **d**, Comparison of leaf size between *G. max* and *G. soja*. Scale bar = 5 cm. **e**, Primary mapping region of *qDRT12.3* on chromosome 12. The *y*-axis represents the log10 likelihood ratio and $R^2$ values indicate the phenotypic variations explained by each QTL. **f**, Fine mapping of *qPB-12*. **g-h**, Fin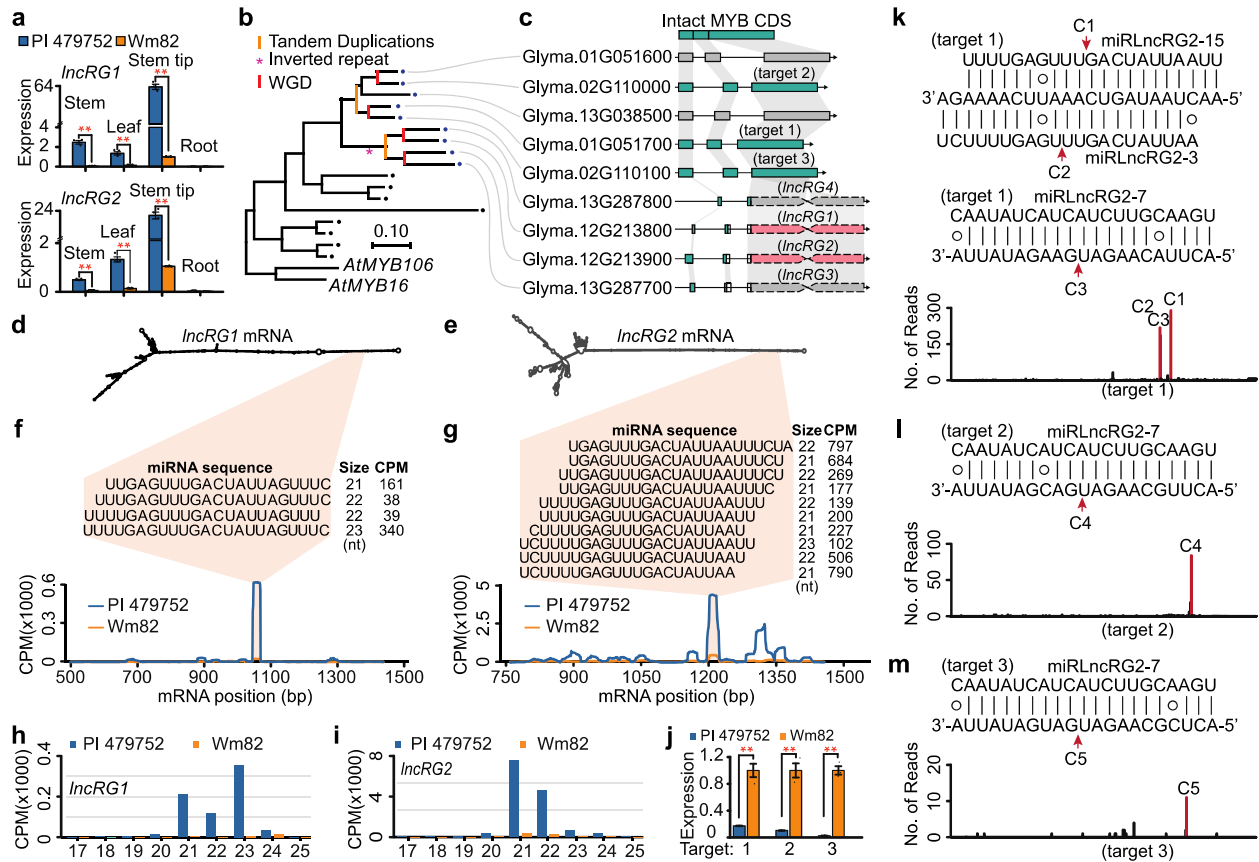e mapping of *qMSL12* (g) and *qLSZ12* (h). Each bar represents the genotype of the recombinants with the same haplotype at all markers. The black color represents the *G. soja* genotype and the grey color represents the *G. max* genotype. Arrows indicate the deduced location of the QTL. Green and brown shades highlight the final mapping interval. Data are represented as mean ± SEM.

**Fig. 2:** ***LncRG1*** **and** ***lncRG2*** **harbor IRs and produce abundant sRNAs primarily targeting three closely related MYB genes. a**, Expression levels of *lncRG1* and *lncRG2* in different tissues as determined by qRT-PCR with Wm82 stem tip set as "1" and the others adjusted accordingly. **b**, Phylogenetic relationships of *lncRG1*, *lncRG2*, and their close MYB relatives. Colored lines indicate duplication events. Red asterisk marks the deduced time when the o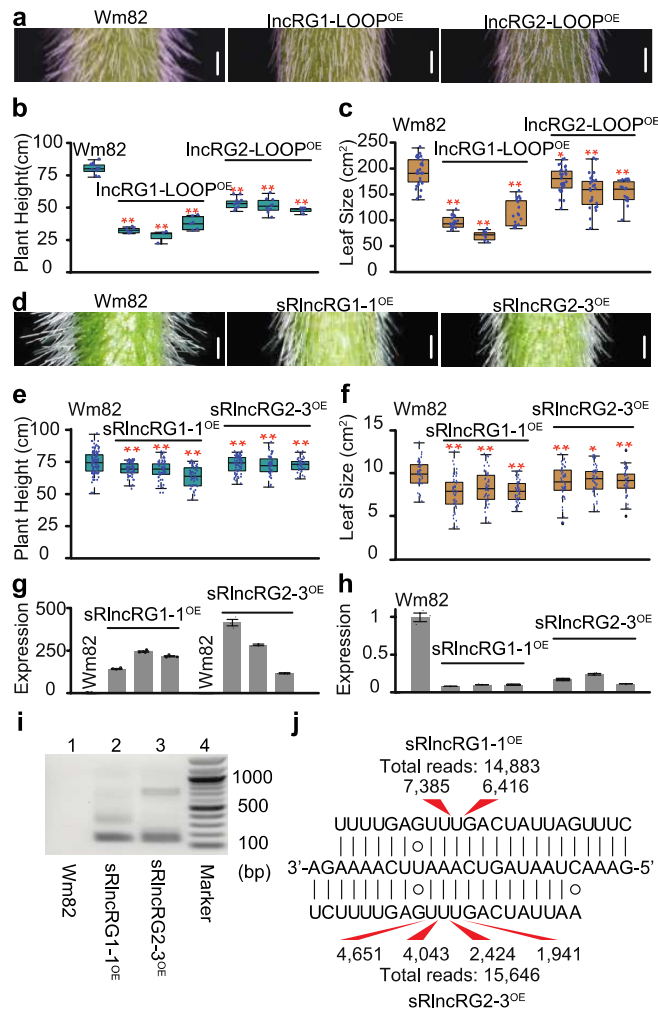riginal IR occurred. **c**, Gene models and alignments of *lncRG1*, *lncRG2*, and their close MYB relatives. Green bars represent coding regions and pink bars represent the inverted repeats (IRs). **d-e**, Predicted secondary structures of *lncRG1* and *lncRG2* transcripts. **f-g**, Distribution, abundance, and the major cluster of sRNAs produced by *lncRG1* and *lncRG2*. **h-i**,

664  Abundance of sRNAs in different sizes produced by *lncRG1* and *lncRG2*. **j**,

665  Expression levels of the target genes, Glyma.01G051700 (target 1),

666  Glyma.02G110000 (target 2), and Glyma.02G110100 (target 3), as

667  determined by qRT-PCR with Wm82 set as "1" and the others adjusted

668  accordingly. **k-m**, The predicted cleavage sites supported by degradome

669  sequencing on the target genes. Letter Cs represents cleavage sites. In (a)

670  and (j), The dots show the values from different biological replicates (n=3),

671  and the red asterisks indicate the significant level at $P < 0.01$ (Student's *t*-

672  test) and data are represented as mean ± SEM.

**Fig. 3: Overproduction of sRNAs in cultivated soybean promotes the wild soybean-type phenotypes. a-c**, The phenotypic changes in pubescence form (a), plant height (b) and leaf size (c) of the stem-loop over-expression lines compared with those of Wm82. Scale bars=3 mm. **d-e**, The phenotypic change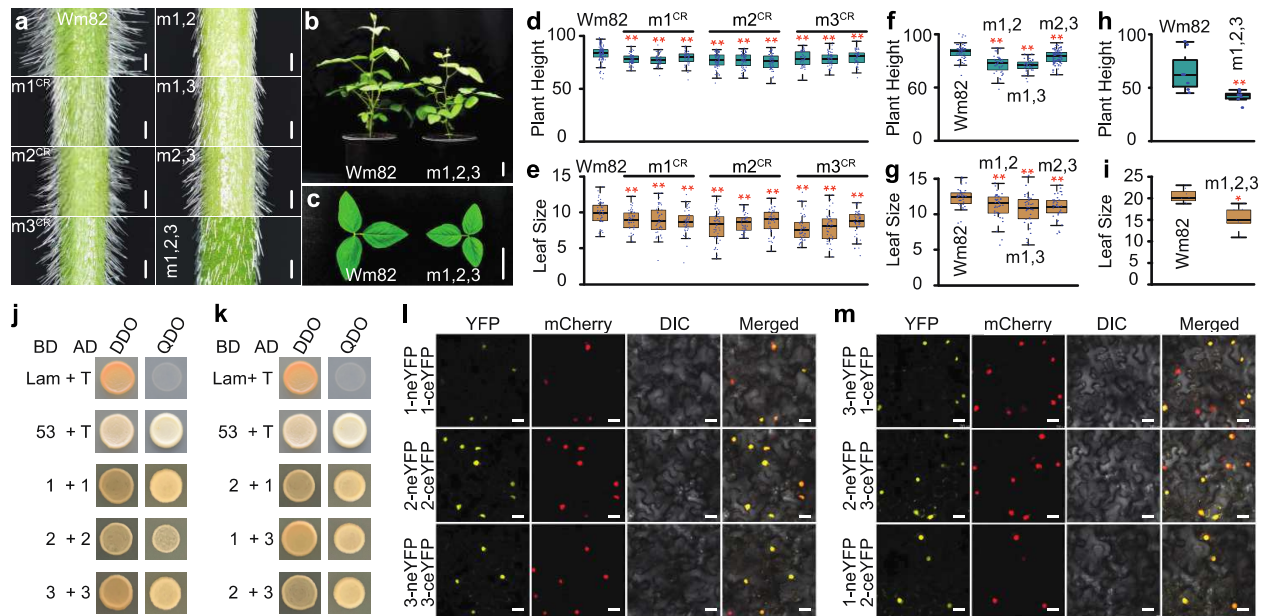s in pubescence form (d), plant height (e) and leaf size (f) of the artificial miRNA over-expression lines compared with those of Wm82. Scale bars=3 mm. **g-h**, The expression levels of the artificial miRNAs (g) and the three target MYB genes (h) in the transgenic lines compared with those inWm82 as determined by stem loop RT-qPCR and regular qRT-PCR,

684    respectively, with Wm82 set as "1" and the others adjusted accordingly. The

685    dots show the values from different biological replicates (n=3). Data are

686    represented as mean ± SEM. **i**, Gel electrophoresis image of RLM-RACE from

687    the transgenic lines and Wm82. **j**, Cleavage frequencies detected by RLM-

688    RACE followed by deep sequencing. Numbers show the total reads number

689    and the read number at each cleavage site. In (b), (c), (e), and (f), the

690    horizontal lines indicate the medians, and the boxes represent the

691    interquartile range (IQR). The whiskers represent the range of 1.5 times IQR

692    and dots beyond the whiskers are outlier values. Red asterisks indicate

693    significant levels at $P<0.01$ or $P<0.05$ (Student's $t$-test).

**Fig. 4: Functional redundancy and divergence of the three MYB genes targeted by the sRNAs. a-c**, Photographic illustration of the phenotypic changes in the pubescence form (a), plant height (b), and leaf size (c) of the gene-edited mutants compare with Wm82. The m1, m2 and m3 are mutants of target 1, target 2 and target 3, respectively. Scale bars=3mm in (a) and 5cm in (b-c). **d-e**, Statistics of the plant height (d) and the leaf size (e) of single mutants and Wm82. **f-g**, Statistics of the plant height (f) and the leaf size (g) of the double mutants and Wm82. **h-i**, Statistics of the plant height (h) and leaf size (i) of the triple mutants and Wm82. **j-k**, Home- (j) and hetero- (k) protein-protein interactions among the three target genes detected by Y2H assays. AD, activation domain; BD, binding domain; DDO, double dropout; QDO, quadruple dropout. **l-m**. Home- (l) and hetero- (m) protein-protein interactions among the three target genes detected by BiFC assay. Scale bars=20μm. In (d-i), horizontal lines indicate the medians, and the boxes

710    represent the interquartile range (IQR). The whiskers represent the range of

711    1.5 times IQR and dots beyond the whiskers are outlier values. Red asterisks

712    indicate significant levels at $P < 0.01$ or $P < 0.05$ (Student's $t$-test).

**Fig. 5: The birth and evolutionary consequences of the lncRGs in legumes. a**, Collinearity analysis of nine legume species at the region harboring the orthologs of *lncRG1* and *lncRG2*. Black boxes present genes and grey shades connect the ortholog genes between species. Red triangles represent the IRs. **b**, Phylogenetic relationships of the nine legume species as determined in previous studies[10,11]. Red lines highlight the genera that carry the IRs and the asterisk indicates the deduced timepoint when the original IRs occurred. **c**, Nucleotide diversity between the forward and reverse repeats in each species. **d-e**, The distribution patterns of the sRNAs

724    produced by *lncRG1* (d) and *lncRG2* (e) in ten diverse soybean accessions as

725    indicated by different colors. Arrows points the position of major sRNA peaks

726    of PI 479752. **f**, The three MYB genes (targets 1, 2, and 3, as shown in Fig.

727    2c) predicted to be targeted by the top 20 sRNAs produced by *lncRG1* and

728    *lncRG2* in each of the ten soybean accessions. Black dots indicate predicted

729    targets, while grey dots indicate they are not predicted to be targets.

**a** qPB12

**b** M5 M7

Glyma.12G213800    Glyma.12G213900

**c** Resistance to Leafhopper    qDRT12.3

**f** Selective Sweep

**d** Pubescense Form

**g** G. soja    G. max

lncRG1    lncRG2

**e** G. soja    Landrace    Elite Cultivar

Appressed    Erect

86.6%    13.4%    28.7%    71.3%    3.3%    96.7%

n=82    n=460    n=182

**h** r = 0.93
P = 9.37 x 10⁻²⁰
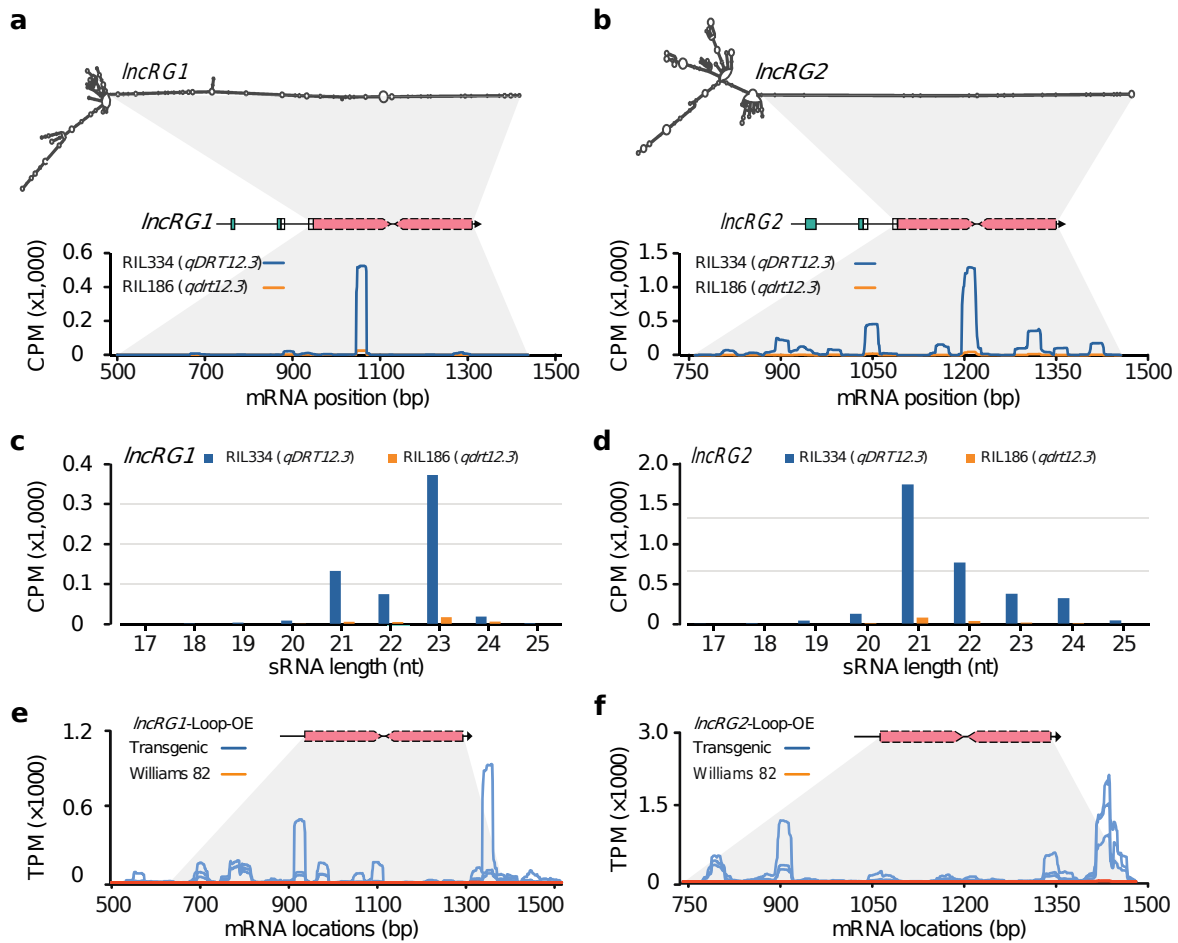
733 **Extended Data Fig. 1 Association studies, selection analyses and**

734 **expression analyses. a-b** Association between genetic variations and

735 expression levels of *lncRG1* and *lncRG2* within the final mapping region.

736 Manhattan plot displays the result of genome-wide association study (GWAS)

737 on pubescence form. The y-axes are the negative log10 of the *P*-values and

738    the red color highlight markers within the final mapping region. **c-d**,

739    Manhattan plots displaying the results of genome-wide association studies

740    (GWAS) on leafhopper resistance (c) and pubescence form (d). The *y*-axes

741    represent the negative log10 of the *P*-values from GWAS. The *x*-axes

742    represent the twenty soybean chromosomes. The rectangle highlights the

743    *qDRT12.3* locus on chromosome 12. The genotypic and phenotypic data of

744    the soybean accessions (n=784) used for the GWAS are from the USDA

745    soybean germplasm collection. **e**, Frequencies of erect and appressed

746    pubescence form in *G. soja*, landrace and elite cultivar sub-populations. n

747    indicates the number of soybean accessions in each sub-population. **f**,

748    Selective sweep surrounding the fine-mapped *qDRT12.3* region. The *y*-axis is

749    the ratio of nucleotide diversity (π) of landraces (n=328) with erect

750    pubescence over *G. soja* (n=103). Each vertical bar represents a 100-kb

751    window (with 10-kb sliding step). The red arrows pinpoint the positions of

752    *lncRG1* and *lncRG2*. The *x*-axis presents the physical positions based on the

753    Zhonghuang 13 (v2) genome assembly. **g**, Comparison of expression levels of

754    *lncRG1* and *lncRG2* between *G. soja* (n=9) and *G. max* (n=36) accessions. The

755    *y*-axis represents the expression level as measured from RNA-seq data and

756    the unit is fragments per kilobase of transcript per million mapped reads

757    (FPKM). The red asterisks indicate the significant level at *P* < 0.01 (Student's

758    *t*-test) and data are represented as mean ± SEM. **h**, Co-expression between

759    *lncRG1* and *lncRG2*. The *x*-axis and y-axis represent the expression levels of

760    *lncRG1* and *lncRG2*, respectively, as measured from RNA-seq data of 45

761   highly diverse soybean accessions. Each dot represents a single soybean

762   accession, with blue dots for *G. soja* haplotype (n=11) and orange dots for *G.*

763   *max* haplotype (n=34). Unit is fragments per kilobase of transcript per million

764   mapped reads (FPKM). Dashed line is the trend line. The Pearson correlation

765   value and the corresponding *P*-value were labeled.

**Extended Data Fig. 2 Abundance and distribution of sRNAs produced by *lncRG1* and *lncRG2* in a pair of RILs and the transgenic lines. a**, Abundance and distribution of sRNAs produced by *lncRG1* in RIL186 (*qdrt12.3*) and RIL334 (*qDRT12.3*). The *x*-axis shows the position on the *lncRG1* transcript, and the *y*-axis is the abundance in copy per million reads (CPM). **b**, Abundance and distribution of sRNAs produced by *lncRG2* in RIL186 (*qdrt12.3*) and RIL334 (*qDRT12.3*). The *x*-axis shows the position on the *lncRG2* transcript, and the *y*-axis is abundance in copy per million reads (CPM). **c**, Frequencies of sRNA from *lncRG1* at different sizes from 17nt to
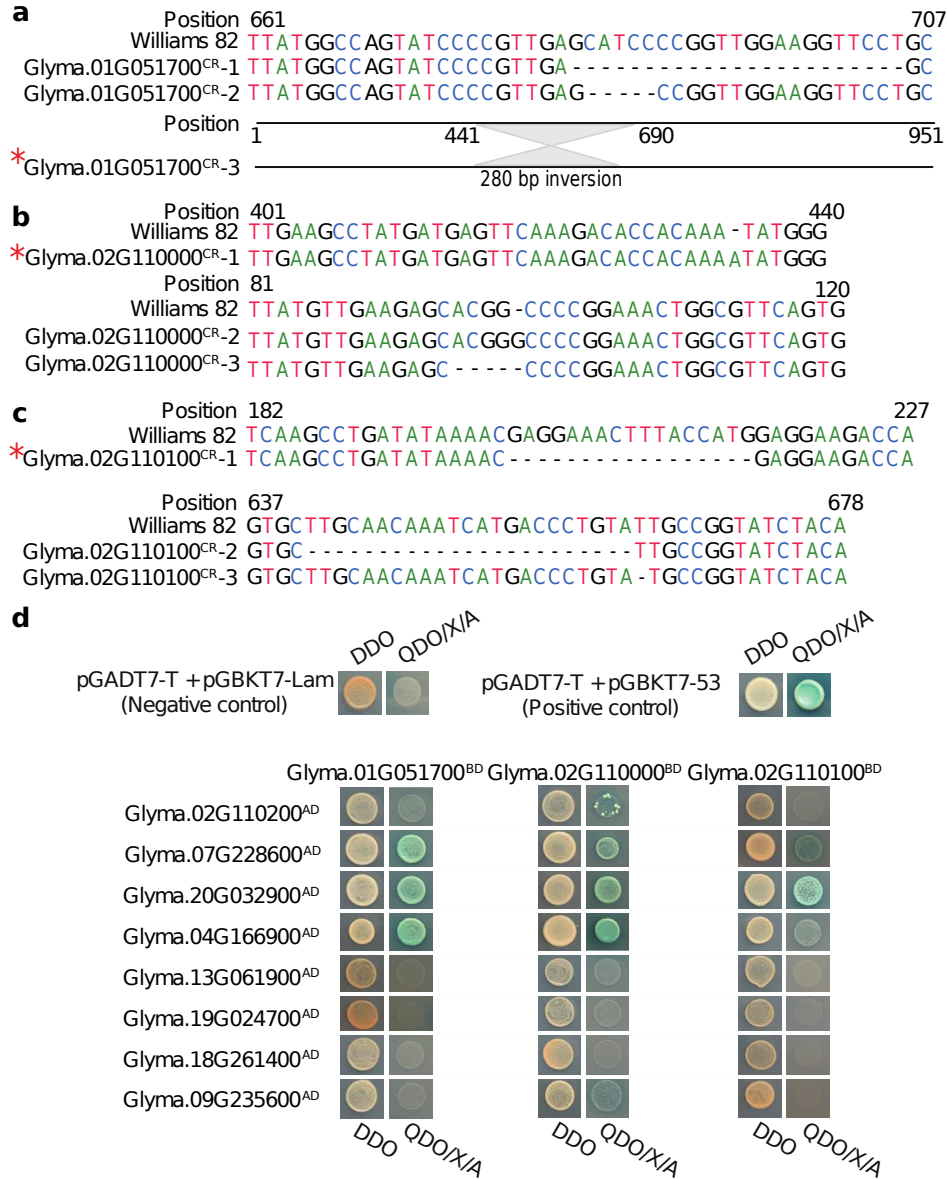
777    25nt in RIL186 (*qdrt12.3*) and RIL334 (*qDRT12.3*). **d**, Frequencies of sRNA

778    from *lncRG2* at different sizes 17nt to 25nt in RIL186 (*qdrt12.3*) and RIL334

779    (*qDRT12.3*). **e**, Abundance and distribution of sRNAs along the transcript of

780    *lncRG1* in the lncRG1-LOOP$^{OE}$ transgenic lines. The *x*-axis shows the position

781    on the *lncRG1* transcript, and the *y*-axis is the abundance in copy per million

782    reads (CPM). **f**, Abundance and distribution of sRNAs along the transcript of

783    *lncRG2* in the lncRG2-LOOP$^{OE}$ transgenic lines. The *x*-axis shows the position

784    on the *lncRG2* transcript, and the *y*-axis is the abundance in copy per million

785    reads (CPM).

**a**

| Position | 661 | 707 |
|---|---|---|
| Williams 82 | TTATGGCCAGTATCCCCGTTGAGCATCCCCGGTTGGAAGGTTCCTGC | |
| Glyma.01G051700^CR-1 | TTATGGCCAGTATCCCCGTTGA- - - - - - - - - - - - - - - - - - - - - - - -GC | |
| Glyma.01G051700^CR-2 | TTATGGCCAGTATCCCCGTTGAG- - - - -CCGGTTGGAAGGTTCCTGC | |

Position 1    441    690    951
*Glyma.01G051700^CR-3
280 bp inversion

**b**

| Position | 401 | 440 |
|---|---|---|
| Williams 82 | TTGAAGCCTATGATGAGTTCAAAGACACCACAAA-TATGGG | |
| *Glyma.02G110000^CR-1 | TTGAAGCCTATGATGAGTTCAAAGACACCACAAAATATGGG | |

| Position | 81 | 120 |
|---|---|---|
| Williams 82 | TTATGTTGAAGAGCACGG-CCCCGGAAACTGGCGTTCAGTG | |
| Glyma.02G110000^CR-2 | TTATGTTGAAGAGCACGGGCCCCGGAAACTGGCGTTCAGTG | |
| Glyma.02G110000^CR-3 | TTATGTTGAAGAGC- - - - -CCCCGGAAACTGGCGTTCAGTG | |

**c**

| Position | 182 | 227 |
|---|---|---|
| Williams 82 | TCAAGCCTGATATAAAACGAGGAAACTTTACCATGGAGGAAGACCA | |
| *Glyma.02G110100^CR-1 | TCAAGCCTGATATAAAAC- - - - - - - - - - - - - - - - - -GAGGAAGACCA | |

| Position | 637 | 678 |
|---|---|---|
| Williams 82 | GTGCTTGCAACAAATCATGACCCTGTATTGCCGGTATCTACA | |
| Glyma.02G110100^CR-2 | GTGC- - - - - - - - - - - - - - - - - - - - - - - -TTGCCGGTATCTACA | |
| Glyma.02G110100^CR-3 | GTGCTTGCAACAAATCATGACCCTGTA-TGCCGGTATCTACA | |

**d**

**Extended Data Fig. 3 Mutations created by CRISPR-Cas9 and protein-protein interaction as detected by Y2H. a-c**, Frameshift mutants created by CRI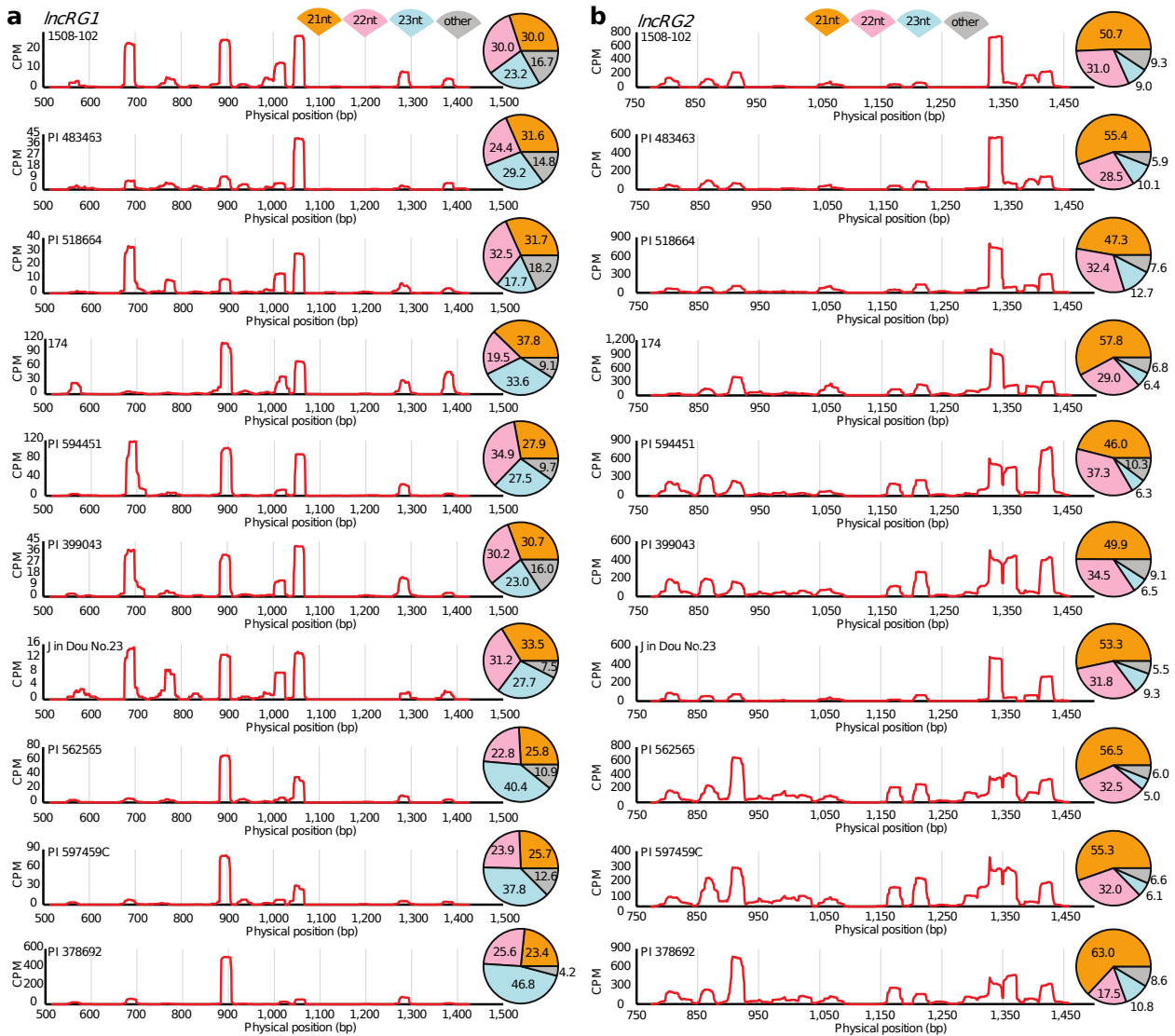SPR-Cas9 for each of the three MYB genes, Glyma.01G051700 (a), Glyma.02G110000 (b) and Glyma.02G110100 (c). The top sequence shows the Wm82 sequence and the position of each base pair in Wm82. - represent deletions in the editing lines. Red asterisk indicates the lines selected for

794    crossing to make double editing lines. **d**, Protein-protein interactions among

795    MYB transcription factors as detected by the yeast two hybrid (Y2H) system.

796    Colonies on DDO plate indicate the successful transformation of the construct

797    in yeast cells. Blue colonies on QDO/X/A plates indicate positive protein-

798    protein interactions. AD, activation domain; BD, binding domain; DDO,

799    double dropout; QDO, quadruple dropout. X, X-alpha-Gal; A, Aureobasidin A.

800

801

**Extended Data Fig. 4 Distribution of the sRNAs produced by *lncRG1* and *lncRG2* in ten diverse soybean accessions.** The *x*-axis shows the position on the *lncRG1* (a) or *lncRG2* (b) transcripts, and the *y*-axis is abundance in copy per million reads (CPM). The relative abundances of sRNAs of different sizes detected in individual accessions are shown in percentage (%) in individual pies.

**Supplemental information**

**Supplementary Table 1**. sRNAs produced by *lncRG1* and *lncRG2* with CPM>10.

**Supplementary Table 2.** List of genes targeted by 27 sRNAs (CPM>100) produced by *lncRG1* and *lincRG2*.

**Supplementary Table 3.** Expression levels (FPKM) of the 163 target genes in shoots, stems, and leaves of Williams 82 and PI 479752.

**Supplementary Table 4.** List of top 20 sRNAs produced by *lncRG1* and *lncRG2* in 10 diverse soybean accessions with *G. soja* haplotype.

**Supplementary Table 5.** List of genes targeted by sRNAs (top 20) produced by *lncRG1* and *lincRG2* in 10 soybean accessions.

**Supplementary Table 6.** List of primers used in this study.

**Supplementary Movie 1.** Erected pubescence confers resistance to leafhopper.

**Supplementary Movie 2.** Appressed pubescence is susceptible to leafhopper.

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- SupplementaryTables.xlsx
- MovieS1.mov
- MovieS2.mov