1 **Genome sequence of *Hydrangea macrophylla* and its application in analysis of the double**

2 **flower phenotype**

3

4 **Authors**

5 Nashima K[*1], Shirasawa K[*2], Ghelfi A[2], Hirakawa H[2], Isobe S[2], Suyama T[3], Wada T[3], Kurokura T[4],

6 Uemachi T[5], Azuma M[1], Akutsu M[6], Kodama M[6], Nakazawa Y[6], Namai K[6]

7

8 1. College of Bioresource Sciences, Nihon University, Kameino 1866, Fujisawa, Kanagawa, 252-

9  0880 Japan

10 2. Kazusa DNA Research Institute, Kazusa-Kamatari 2-6-7, Kisarazu, Chiba, 292-0813 Japan

11 3. Fukuoka Agriculture and Forestry Research Center, Yoshiki 587, Chikushino, Fukuoka, 818-8549

12  Japan

13 4. Faculty of Agriculture, Utsunomiya University, Mine 350, Utsunomiya, Tochigi, 321-8505 Japan

14 5. School of Environmental Science, University of Shiga Prefecture, Hassakacho 2500, Hikone,

15  Shiga, 522-0057 Japan

16 6. Tochigi Prefectural Agricultural Experimental Station, Kawarayacho 1080, Utsunomiya, Tochigi,

17  320-0002 Japan

18

19 *equally contributed as first author

20 Corresponding author: Nashima K

21 College of Bioresource Sciences, Nihon University, Kameino 1866, Fujisawa, Kanagawa, 252-0880

22 Japan

23 Tel: +81-466-84-3507

24 Mail: nashima.kenji@nihon-u.ac.jp

**Abstract**

Owing to its high ornamental value, the double flower phenotype of hydrangea (*Hydrangea macrophylla*) is one of its most important traits. In this study, genome sequence information was obtained to explore effective DNA markers and the causative genes for double flower production in hydrangea. Single molecule real-time sequencing data followed by a HiC analysis was employed. The resultant haplotype-phased sequences consisted of 3,779 sequences (2.256 Gb in length and N50 of 1.5 Mb), and 18 pseudomolecules comprising 1.08 Gb scaffold sequences along with a high-density SNP genetic linkage map. Using the genome sequence data obtained from two breeding populations, the SNPs linked to double flower loci ($D_{jo}$ and $D_{su}$), were discovered for each breeding population. DNA markers J01 linked to $D_{jo}$ and S01 linked to $D_{su}$ were developed, and these could be used successfully to distinguish the recessive double flower allele for each locus respectively. The *LEAFY* gene was suggested as the causative gene for $D_{su,}$ since frameshift was specifically observed in double flower accession with $d_{su}$. The genome information obtained in this study will facilitate a wide range of genomic studies on hydrangea in the future.

Keywords:

Hydrangea, double flower, de novo genome sequencing, DNA marker

49    **1. Introduction**

50       *Hydrangea macrophylla* (Thunb.) Ser., commonly known as hydrangea, originated in Japan,

51    and since it is the place of origin, there are rich genetic resources for this plant in Japan. Wild

52    hydrangea accessions with superior characteristics have been bred to create attractive cultivars, and it

53    has a long history of use as an ornamental garden plant in temperate regions. There are both decorative

54    and non-decorative flowers in an inflorescence. Decorative flowers have large ornamental sepals that

55    attract pollinators, whereas non-decorative flowers have inconspicuous perianths that instead play a

56    major role in seed production[1-3]. In hydrangea, there are two types of decorative flower phenotype:

57    single flower and double flower. Single flowers generally have four petaloid sepals per decorative

58    flower, while this number in double flowers is approximately fourteen. Double flowers do not have

59    stamens or petals[4]. Therefore, petals and stamens would be converted to petaloid sepals since number

60    of petaloid sepals are increased and stamens and petals are lost. Because of their high ornamental value,

61    producing double flower is an important breeding target in hydrangea cultivation.

62       To obtain double flower progenies, the double flower cultivars 'Sumidanohanabi' (Figure

63    1A) and 'Jogasaki' (Figure 1B) were crossbred in Japan[4]. Previous studies have suggested that double

64    flower phenotype is a recessive characteristic controlled by a single major gene[4,5]. Suyama et al.[4]

65    found that crosses between the progeny of 'Sumidanohanabi' and the progeny of 'Jogasaki' produced

66    only single flower descendants. Thus, it was also suggested that genes controlling the double flower

67    phenotype are different[4]. While Suyama et al.[4] suggested that a single locus with different double

68    flower alleles controls the phenotype, Waki et al.[5] speculated that two different loci control double

69    flower production individually. Therefore, it is not clear whether a single locus or two loci control the

70    phenotype. We term the double flower locus $D_{su}$ as the locus controlling the double flower phenotype

71    of 'Sumidanohanabi' and the double flower locus $D_{jo}$ as the locus controlling the double flower

72    phenotype of 'Jogasaki.' Waki et al.[5] identified $D_{su}$ on the genetic linkage map. They also found that

73    the DNA marker STAB045 was the nearest marker to $D_{su}$, and that STAB045 could help in

74    distinguishing flower phenotype with a 98.6% fitting ratio[5]. Contrarily, $D_{jo}$ has not been identified,

75    and the DNA marker linked to $D_{jo}$ has not been developed. It is still not known whether $D_{jo}$ and $D_{su}$

76    are at the same loci.

77         The mechanisms and genes controlling double flower phenotype in hydrangea have not been

78    clarified. Waki et al.[5] hypothesized that the mutation of C-class genes could be associated with the

79    double flower phenotype of 'Sumidanohanabi', since the C-class gene mutant of *Arabidopsis thaliana*

80    and C-class gene-repressed petunias produce double flowers[6]. However, the double flower phenotype

81    of hydrangea is morphologically different from that of *A. thaliana* and petunia—petals and stamens

82    would be converted to petaloid sepals, while stamens converted to petals in *A. thaiana* and petunia.

83    This suggests that the genes controlling double flower production in hydrangea are different from

84    corresponding genes in other plant species. Identification of the genes controlling double flower

85    production in hydrangea could reveal novel regulatory mechanisms of flower development.

86         Genomic information is essential for DNA marker development and identification of genes

87    controlling specific phenotypes. However, no reference genome sequence is publicly available for

88    hydrangea so far. Although a genome assembly of hydrangea (1.6 Gb) using only short-read data has

89    been reported[7], the resultant assembly is so fragmented that it comprises 1,519,429 contigs with an

90    N50 size of 2,447 bp and has not been disclosed. Improved, advanced long-read technologies and

91    bioinformatics methods would make it possible to determine the sequences of complex genomes. An

92    assembly strategy for single molecule real-time sequencing data followed by a HiC analysis has been

93    developed to generate haplotype-phased sequences in heterozygous regions of diploid genomes[8].

94    Genome sequences at the chromosome level could be obtained with a HiC clustering analysis[9] as well

95    as with a genetic linkage analysis[10]. Such genomic sequence will provide basic information to identify

96    genes and DNA markers of interest, and to discover allelic sequence variations. In this study, we

97    constructed the genomic DNA sequence, obtained SNPs information, and performed gene prediction.

98    We also developed DNA markers linked to $D_{jo}$ using SNP information obtained by double digest

99    restriction site associated DNA sequence (ddRAD-Seq) analysis of breeding population 12GM1,

100    which segregated double flower phenotypes of $D_{jo}$. In addition, we attempted to identify the causative

101    genes for $D_{jo}$ and $D_{su}$.

102

103    **2. Materials and Methods**

104    **2.1. De novo assembly of the hydrangea genome**

105    For genomic DNA sequencing, *H. macrophylla* 'Aogashima-1,' collected from Aogashima

106    island of the Izu Islands in Tokyo Prefecture, Japan, was used. Genomic DNA was extracted from the

107    young leaves with Genomic-Tip (Qiagen, Hilden, Germany). First, we constructed a sequencing

108    library (insert size of 500 bp) with TruSeq DNA PCR-Free Library Prep Kit (Illumina, San Diego, CA,

109    USA) to sequence on HiSeqX (Illumina). The size of the 'Aogashima-1' genome was estimated using

110    Jellyfish v2.1.4[11]. After removing adapter sequences and trimming low-quality reads, high-quality

111    reads were assembled using Platanus[12]. The resultant sequences were designated HMA_r0.1.

112    Completeness of the assembly was assessed with sets of BUSCO v.1.1b[13].

113    Next, a SMRT library was constructed with SMRTbell Express Template Prep Kit 2.0

114    (PacBio, Menlo Park, CA, USA) in accordance with the manufacture's protocol and sequenced with

115    SMRT Cell v2.1 on a Sequel System (PacBio). The sequence reads were assembled using FALCON

116    v.1.8.8[14] to generate primary contig sequences and to associate contigs representing alternative alleles.

117    Haplotype‐resolved assemblies (i.e. haplotigs) were generated using FALCON-Unzip v.1.8.8[14].

118    Potential sequence errors in the contigs were corrected twice with ARROW v.2.2.1 implemented in

119    SMRT Link v.5.0 (PacBio) followed by one polishing with Pilon[15]. Subsequently, a HiC library was

120    constructed with Proximo Hi-C (Plant) Kit (Phase Genomics, Seattle, WA, USA) and sequenced on

121  HiSeqX (Illumina). After removing adapter sequences and trimming low-quality reads, high-quality

122  HiC reads were used to generate two haplotype-phased sequences from the primary contigs and

123  haplotig sequences with FALCON-Phase[8].

124       To validate the accuracy of the sequences, we developed a genetic map based on SNPs,

125  which were from a ddRAD-Seq analysis on an F2 mapping population (n = 147), namely 12GM1,

126  maintained at the Fukuoka Agriculture and Forestry Research Center, Japan. The 12GM1 population

127  was generated from a cross between 'Posy Bouquet Grace' (Figure 1C) and 'Blue Picotee Manaslu'

128  (Figure 1D). Genomic DNA was extracted from the leaves with DNeasy Plant Mini Kit (Qiagen). A

129  ddRAD-Seq library was constructed as described in Shirasawa et al.[16] and sequenced with HiSeq4000.

130  Sequence reads were processed as described by Shirasawa et al.[16] and mapped on the HMA_r1.2 as a

131  reference. From the mapping alignment, high-confidence biallelic SNPs were obtained with the

132  following filtering options: --minDP 5 --minQ 10 --max-missing 0.5. The genetic map was constructed

133  with Lep-Map3[17].

134       Potential mis-jointed points in the phase 0 and 1 sequences of HMA_r1.2 were cut and re-

135  joined, based on the marker order in the genetic map, for which we employed ALLMAPS[18]. The

136  resultant sequences were named HMA_r1.3.pmol, as two haplotype-phased pseudomolecule

137  sequences of the 'Aogashima-1' genome. Sequences that were unassigned to the genetic map were

138  connected and termed chromosome 0.

139

140  **2.2 Gene prediction**

141       For gene prediction, we performed Iso-Seq analysis. Total RNA was extracted from 12

142  samples of 'Aogashima-1': flower buds (2 stages); decorative flowers (2 stages); colored and colorless

143  non-decorative flowers; fruits; shoots; roots; buds, and one-day light-intercepted leaves and buds. In

144  addition, the 29 samples listed in Supplementary Table S1 were included. Iso-Seq libraries were

145   prepared with the manufacture's Iso-Seq Express Template Preparation protocol, and sequenced on a

146   Sequel System (PacBio). The raw reads obtained were treated with ISO-Seq3 pipeline, implemented

147   in SMRT Link v.5.0 (PacBio) to generate full-length, high-quality consensus isoforms. In parallel,

148   RNA-Seq data was also obtained from the 16 samples listed in Supplementary Table S1. Total RNA

149   extracted from the samples was converted into cDNA and sequenced on HiSeq2000, Hiseq2500

150   (Illumina), and NovaSeq6000 (Illumina). The Iso-Seq isoform sequences and the RNA-Seq short-

151   reads were employed for gene prediction.

152           To identify putative protein-encoding genes in the genome assemblies, ab-initio-, evidence-,

153   and homology-based gene prediction methods were used. For this prediction, unigene sets generated

154   from 1) the Iso-Seq isoforms; 2) de novo assembly of the RNA-Seq short-reads with Trinity-v2.4.0[19];

155   3) peptide sequences predicted from the genomes of *Arabidopsis thaliana*, *Arachis hypogaea*,

156   *Cannabis sativa*, *Capsicum annuum*, *Cucumis sativus*, *Populus trichocarpa*, and *Quercus lobata*; and

157   4) *ab-initio* genes, were predicted with Augustus-v3.3.1[20]. The unigene sequences were aligned onto

158   the genome assembly with BLAT[21] and genome positions of the genes were listed in general feature

159   format version 3 with blat2gff.pl (https://github.com/vikas0633/perl/blob/master/blat2gff.pl). Gene

160   annotation was performed with Hayai-annotation Plants[22]. Completeness of the gene prediction was

161   assessed with sets of BUSCO v4.0.6[13].

162

163   **2.3 Detection of SNPs linked to the double flower phenotype**

164           For identification of SNPs linked to double flower loci $D_{jo}$ and $D_{su}$, ddRAD-Seq data

165   analysis was performed. ddRAD-Seq data of the 12GM1 population described above was used to

166   identify $D_{jo}$. For identification of SNPs linked to double flower locus $D_{su}$, KF population[5]—93 F2

167   specimens of 'Kirakiraboshi' (Figure 1E) and 'Frau Yoshimi' (Figure 1F)—were used for ddRAD-Seq

168   analysis. The KF population was maintained at Tochigi Prefectural Agricultural Experimental Station,

169    Japan. ddRAD-Seq analysis of the KF population was performed using the same method used for the

170    12GM1 population.

171        ddRAD-Seq data of the 12GM1 and KF populations were processed as follows: Low-quality

172    sequences were removed and adapters were trimmed using Trimmomatic-0.36[23] (LEADING:10,

173    TRAILING:10, SLIDINGWINDOW:4:15, MINLEN:51). BWA-MEM (version 0.7.15-r1140) was

174    used for mapping onto genome sequence. The resultant sequence alignment/map format (SAM) files

175    were converted to binary sequence alignment/map format files and subjected to SNP calling using the

176    mpileup option of SAMtools[24] (version 1.4.1) and the view option of BCFtools (parameter -vcg). If

177    the DP of called SNP in individuals was under 5%, the genotype was treated as missing. SNPs with

178    5% or more of missing genotype were filtered out. Each SNP was evaluated, fitting ratios with the

179    flower phenotype.

180

181    **2.4 DNA marker development and analysis for $D_{jo}$**

182        A CAPS marker was designed based on SNP (Scaffold:0008F-2, position: 780104) that was

183    completely linked to the double flower locus $D_{jo}$. Primers were designed using Primer3[25] under

184    conditions with product size ranging from 150 to 350 bp, primer size from 18 to 27 bp, and primer

185    TM from 57 to 63°C. Primer sequences of the designed CAPS marker named J01 were: Forward: 5′-

186    CTGGCAGATTCCTCCTGAC-3′ and Reverse: 5′-TATTTCCTTGGGGAGGCTCT-3′. PCR assays

187    were done in a total volume of 10 $\mu$ L, containing 5 μL of GoTaq Master Mix (Promega, Mdison, WI,

188    USA), 1 mM each of forward and reverse primer, and 5 ng of template DNA. The PCR conditions

189    were 94°C for 2 min, 35 cycles of denaturation at 94°C for 1 min, annealing at 55°C for 1 min, and

190    extension at 72°C for 1 min; and a final extension step at 72°C for 3 min. Then, restriction enzyme

191    assay was done in a total volume of 10 μL, containing 5 μL of PCR product, ten units of restriction

192    enzyme TaqI (New England Biolabs, Ipswich, MA, USA), and 1 μL of cut smart buffer. Restriction

193  enzyme assay was performed at 65°C for 3 h. The restriction assay product was stained with 1x

194  GRRED (Biocraft, Tokyo, Japan) and separated in 1.5% (w/v) agarose gel in TAE buffer. Designed

195  CAPS marker J01 was applied to the 12GM1 population, 14GT77 population (64 F2 specimens of

196  'Posy Bouquet Grace' × 'Chibori') and the 15IJP1 population (98 F1 specimens of 'Izunohana' ×

197  03JP1) that segregate the double flower locus $D_{jo}$.

198

199  **2.5 Resequencing and comparison of LEAFY gene sequence and DNA marker development**

200  To compare sequences, resequencing of genomic DNA was performed for accessions of

201  'Kirakiraboshi,' 'Frau Yoshimi,' 'Posy Bouquet Grace,' and 'Blue Picotee Manaslu.' Sequencing

202  libraries (insert size of 500 bp) for the four lines were constructed with TruSeq DNA PCR-Free Library

203  Prep Kit (Illumina) to sequence on a HiSeqX (Illumina). From the sequence reads obtained, low-

204  quality bases were deleted with PRINSEQ v.0.20.4[26] and adaptor sequences were trimmed with fastx

205  clipper        (parameter,        -a        AGATCGGAAGAGC)        in        FASTX-Toolkit        v.0.0.13

206  (http://hannonlab.cshl.edu/fastx_toolkit). High-quality reads were aligned on the HMA_r1.2 with

207  Bowtie2[27] v.2.2.3 to detect sequence variant candidates by with the mpileup command in SAMtools

208  v.0.1.19[24]. High-confidence variants were selected using VCFtools[29] v.0.1.12b with parameters of --

209  minDP 10, --maxDP 100, --minQ 999, --max-missing 1.

210  For comparison of *LEAFY (LFY)* sequence in 'Kirakiraboshi,' 'Frau Yoshimi,' 'Posy

211  Bouquet Grace,' and 'Blue Picotee Manaslu,' BLAST analysis using genomic sequence of *LFY*

212  (Scaffold 0577F, position 678200-684639) as query, and genomic DNA sequence of each cultivar as

213  database, was performed to confirm detected sequence variants. These data analyses were performed

214  using CLC main workbench (Qiagen). INDEL marker S01 that amplifies the second intron of *LFY*,

215  was designed by visual inspection (Forward: 5′-CATCATTAATAGTGGTGACAG-3′, Reverse: 5′-

216  CACACATGAATTAGTAGCTC-3′). The PCR conditions were 94°C for 2 min, 35 cycles of

217    denaturation at 94°C for 1 min, annealing at 55°C for 1 min, extension at 72°C for 1 min; and a final

218    extension step at 72°C for 3 min. The PCR product was stained with 1x GRRED (Biocraft) and

219    separated in 2.5% (w/v) agarose gel in TAE buffer.

220

221    **2.6 Cloning and sequence determination of LFY gene of 'Kirakiraboshi' and 'Frau Yoshimi'**

222    Total RNA was isolated from the flower buds of 'Kirakiraboshi,' and 'Frau Yoshimi' using RNAiso

223    Plus (TaKaRa, Japan), and reverse transcribed using PrimeScript Ⅱ 1st strand cDNA Synthesis Kit

224    (TaKaRa, Japan). The sequence of the *LFY* gene was amplified by PCR in 50-µL reaction mixture by

225    using TaKaRa Ex Taq Hot Start Version (TaKaRa Bio, Shiga, Japan) and the *LFY* specific primer

226    (Forward: 5′-ATGGCTCCACTACCTCCACC-3′ and Reverse: 5′-CTAACACCCTCTAAAAGCAG-

227    3′). These PCR products were purified, and inserted into a pMD20-T vector using the Mighty TA-

228    cloning kit (TaKaRa Bio). The sequence of *LFY* coding sequence (CDS) in pMD20-T vector was

229    analyzed by 3130xl DNA sequencer (Applied Biosystems, Foster City, CA, USA). Sequence

230    alignments were obtained by using CLC main workbench (Qiagen).

231

232    **2.7 DNA marker assessment across hydrangea accessions**

233    For assessment of DNA markers for the double flower phenotype, 35 *H. macrophylla*

234    accessions were used. Genotyping for J01 was performed as described above. Genotyping for S01 was

235    performed by fragment analysis as follows. PCR amplification was performed in a 10-µL reaction

236    mixture containing 5 µL of GoTaq Master Mix (Promega), 5 pmol FAM-labeled universal primer (5′

237    - FAM-gctacggactgacctcggac -3′), 2.5 pmol forward primer with universal adapter sequence (5′ -

238    gctacggactgacctcggacCATCATTAATAGTGGTGACAG -3′), 5 pmol reverse primer, and 5 ng of

239    template DNA. DNA was amplified in 35 cycles of 94°C for 1 min, 55°C for 1 min, and 72°C for 2

240    min; and a final extension of 5 min at 72°C. The amplified PCR products were separated and detected

241    in a PRISM 3130xl DNA sequencer (Applied Biosystems, USA). The sizes of the amplified bands

242    were scored against internal-standard DNA (400HD-ROX, Applied Biosystems, USA) by

243    GeneMapper software (Applied Biosystems, USA).

244

245    **3. Results and Discussion**

246    **3.1 Draft genome assembly with long-read and HiC technologies**

247    The size of the hydrangea genome was estimated by k-mer-distribution analysis with the short-read of

248    132.3 Gb data. The resultant distribution pattern indicated two peaks, representing homozygous (left

249    peak) and heterozygous (right peak) genomes, respectively (Figure 2). The haploid genome of

250    hydrangea was estimated to be 2.2 Gb in size. The short reads were assembled into 612,846 scaffold

251    sequences. The total length of the resultant scaffolds, i.e. HMA_r0.1, was 1.7 Gb with an N50 length

252    of 9.1 kb (Supplementary Table S2). Only 72.2% of complete single copy orthologues in plant

253    genomes were identified in a BUSCO analysis (Supplementary Table S2).

254        Next, we employed long sequence technology to extend the sequence contiguity and to

255    improve the genome coverage. A total of 106.9 Gb of reads (49.4×) with an N50 read length of 28.8

256    kb was obtained from 14 SMRT Cells. The long-reads were assembled, followed by sequence error

257    corrections into 15,791 contigs consisting of 3,779 primary contigs (2.178 Gb in length and N50 of

258    1.4 Mb), and 12,012 haplotig sequences (1.436 Gb in length and N50 of 184 kb). To obtain two

259    haplotype-phased complete-length sequences, 697 M reads of HiC data (105.3 Gb) were obtained and

260    subjected to FALCON-Phase. The resultant haplotype-phased sequences consisted of 3,779 sequences

261    (2.256 Gb in length and N50 of 1.5 Mb) for "phase 0," and 3,779 sequences (2.227 Gb in length, and

262    N50 of 1.4 Mb) for "phase 1."

263

**3.2 Pseudomolecule sequences based on genetic mapping**

To detect potential errors in the assembly and to assign the contig sequences onto the hydrangea chromosomes, we established an F2 genetic map based on SNPs derived from a ddRAD-Seq technology. Approximately 1.8 million high-quality ddRAD-Seq reads per sample were obtained from the mapping population and mapped to either of the two phased sequences with alignment rates of 88.4% and 88.7%, respectively. A set of SNPs detected from the alignments were classified into 18 groups and ordered to construct two genetic maps for the two phased sequences (2,849.3 cM in length with 3,980 SNPs, and 2,944.5 cM in length with 4,071 SNPs). The nomenclature of the linkage groups was named in accordance with the previous genetic map based on SSRs[5]. The phased sequences were aligned on each genetic map to establish haplotype-phased, chromosome-level pseudomolecule sequences. During this process, one contig was cut due to possible mis-assembly. The resultant sequences for phase 0 had 730 contigs with a total length of 1,078 Mb and the other for phase 1 had 743 contigs spanning 1,076 Mb.

**3.3. Transcriptome analysis followed by gene prediction**

In the Iso-Seq analysis, Circular Consensus Sequence (CCS) reads were generated from the raw sequence reads. The CCS reads were classified in full-length and non-full length reads and the full-length reads were clustered to produce consensus isoforms. In total, 116,634 high-quality isoforms were used for gene prediction. In the RNA-Seq analysis, on the contrary, a total of 80.7 Gb reads were obtained and assembled into 12,265 unigenes. The high-quality isoforms and unigenes together with gene sequences predicted from the *Arabidopsis thaliana*, *Arachis hypogaea*, *Cannabis sativa*, *Capsicum annuum*, *Cucumis sativus*, *Populus trichocarpa*, and *Quercus lobate* genomes were aligned onto the assembly sequence of the hydrangea genome. By adding ab-initio on genes, 32,205 and 32,222 putative protein-encoding genes were predicted from the phase 0 and phase 1 sequences,

288    respectively. This gene set included 91.4% complete BUSCOs. Out of the 10,108 genes, 16,725, and

289    21,985 were assigned to Gene Ontology slim terms in the biological process, cellular component, and

290    molecular function categories, respectively. Furthermore, 4,271 genes had assigned enzyme

291    commission numbers.

292

293    **3.4 Identification of SNPs tightly linked to double flower phenotype**

294    To identify SNPs tightly linked to the double flower phenotype of 'Jogasaki,' ddRAD-Seq

295    analysis was performed on the 12GM1 population, which segregates the double flower phenotype of

296    'Jogasaki.' As a result, 14,006 of SNPs were called by ddRAD-Seq analysis of the 12GM1 population.

297    In this population, the double flower phenotype was expected when the plant was homozygous for the

298    'Posy Bouquet Grace' genotype, and the single flower phenotype was expected when the plant was

299    homozygous for 'Blue Picotee Manaslu' or was heterozygous. Each SNP was tested for its fitting rate

300    to this model. As a result, ten SNPs were found to have more than a 95% fitting rate, and six SNPs

301    were completely co-segregated with flower phenotype (Table 1).

302    CAPS marker J01 was developed based on SNP at scaffold 0008F-2_780104. J01 CAPS

303    marker amplified 167 bp of fragment by PCR, and digestion with Taq I restriction enzyme generated

304    50 bp and 117 bp fragments in the double flower allele (Figure 3). J01 marker was fitted with flower

305    phenotype at 99.3% in the 15IJP1 and 14GT77 populations, which segregated the double flower

306    phenotype of 'Jogasaki' (Supplementary Table S3, S4). This indicated that J01 marker was tightly

307    linked to the $D_{jo}$ locus. Thus, $D_{jo}$ is suggested to be located adjacent to J01, which is located at position

308    46,326,384 in CHR17, (Figure 4).

309    For identification of SNPs linked to the double flower phenotype of 'Sumidanohanabi,' the

310    KF population that segregates the double flower phenotype derived from 'Sumidanohanabi' were used.

311    First, we tried to find co-segregated scaffolds with the double flower phenotype by ddRAD-Seq

312    analysis of the KF population. As a result of ddRAD-Seq analysis, 15,102 of SNPs were called. In this

313    population, the double flower phenotype was expected when the plant was homozygous for the

314    'Kirakiraboshi' genotype, and the single flower phenotype was expected when the plant was

315    homozygous for 'Frau Yoshimi' or was heterozygous. Each SNP was tested for its fitting rate to this

316    model. As a result, five SNPs on three scaffolds were found to have more than a 95% fitting rate with

317    the model (Table 2). Since SNPs on scaffold 3145F all had the same genotype across the KF population,

318    three loci—on scaffold 0577F, 3145F, 0109F—were detected. According to genotypes of the KF

319    population, these three loci were tightly linked within 5 cM; 0109F (0 cM) - 3145F (3.9 cM) - 0577F

320    (5.0 cM). Since the SNP at position 868569 in 0109F was found at the position 57,436,162 in CHR04,

321    locus $D_{su}$, which controls the double flower phenotype of 'Sumidanohanabi,' was suggested to be

322    located on terminal of CHR04 (Figure 4).

323

324    **3.5 Prediction of genes controlling double flower**

325        To find the gene controlling $D_{su}$ and $D_{jo}$, we searched the homeotic genes on scaffolds shown

326    in Table 1 and Table 2. We did not find any notable homeotic gene controlling flower phenotype for

327    $D_{jo}$. For $D_{su}$, the g182220 gene, which encoded a homeotic gene *LFY*, was found on scaffold 0577F.

328    To investigate the possibility that it was the causative gene for $D_{su}$, sequence variants on *LFY* genomic

329    sequence were searched to identify 'Kirakiraboshi' specific mutation, using resequencing data of

330    'Kirakiraboshi,' 'Frau Yoshimi,' 'Posy Bouquet Grace,' and 'Blue Picotee Manaslu.' As a result, five

331    INDELs and six sequence variants were found as 'Kirakiraboshi' specific mutations (Figure 5).

332        Cloning and sequencing of *LFY* CDS was performed on 'Kirakiraboshi' and 'Frau Yoshimi.'

333    From 'Frau Yoshimi,' a single CDS comprising three exons was obtained. From 'Kirakiraboshi,' two

334    CDSs with splice variants were obtained. While splicing 1 CDS resulted in three exons, splicing 2

335    CDS resulted in only two exons, corresponding to the first and third splice products of splicing 1 CDS

336    (Supplementary Figure S1). The deduced amino acid sequences were aligned using CDSs of 'Frau

337    Yoshimi' and 'Kirakiraboshi,' g182220 sequence, protein LFY of *Arabidopsis thaliana,* and protein

338    FLO of *Antirrhium majos*. While the deduced amino acid sequences of 'Frau Yoshimi' and g182220

339    showed sequence similarity in the entire region, frameshift occurred in the two CDSs obtained from

340    'Kirakiraboshi' and the resulting products had no sequence similarity across the latter half (Figure 6).

341    Frameshift observed in splicing 1 CDS was due to one bp of DNA insertion in the second exon, at

342    position 1,931 (Figure N3A). On the contrary, frameshift observed in splicing 2 CDS was due to the

343    complete loss of the second exon (Figure 6).

344            To develop a DNA marker for distinguishing the $d_{su}$ allele from the $D_{su}$ alleles in the *LFY*

345    genomic sequence, we focused and designed a DNA marker on 'Kirakiraboshi' specific 14 bp deletion

346    at position 3,617 from initiation codon (Figure 5). We developed INDEL S01 marker amplified 236

347    bp fragment for the double flower allele of 'Kirakiraboshi,' and 250 bp and 280 bp fragments for the

348    single flower allele of 'Frau Yoshimi' (Figure 7A). Three types of alleles resulted from the presence

349    or absence of a 30 bp deletion at position 3,482 in addition to the 14 bp INDEL. These were both 30

350    bp and 14 bp deletions on the 236 bp allele, 30 bp deletion on the 250 bp allele, and no deletion on the

351    280 bp allele (Figure 7B).

352

353    **3.6 Genotyping of hydrangea accessions using J01 and S01 markers**

354            Since the J01 marker could distinguish $D_{jo}/d_{jo}$ alleles and the S01 marker could distinguish

355    $D_{su}/d_{su}$ alleles, a combined use of J01 and S01 DNA markers was expected to reveal the origin of the

356    double flower phenotype, $d_{jo}$ or $d_{su}$, in various accessions. Therefore, DNA marker genotyping on *H.*

357    *macrophylla* accessions were performed using two DNA markers, J01 and S01. All tested double

358    flower accessions showed homozygous genotypes of J01 or S01; ten of the double flower accessions

359    were homozygous of 117_50 in J01, and four were homozygous of 236 in S01 (Table 3). Contrarily,

360    all single flower accessions showed other genotypes.

361         Previously, the double flower phenotype has been revealed to be controlled by a single locus

362    with the inheritance of single flower dominant and double flower recessive genes[4,5]. It was also

363    suggested that genes controlling the double flower phenotype were different between 'Jogasaki' and

364    'Sumidanohanabi' based on confirmation of the segregation ratio of crossed progenies[4]. Our study

365    revealed that the double flower phenotype of 'Jogasaki' was controlled by a single $D_{jo}$ locus on CHR17,

366    and the double flower phenotype of 'Sumidanohanabi' was controlled by a single $D_{su}$ locus on CHR04.

367    In addition, all double flower accessions showed homozygosity for the double flower allele at one

368    locus, $D_{jo}$ or $D_{su}$. Contrarily, all single flowers have dominant single flower alleles on both $D_{jo}$ and $D_{su}$

369    loci. This indicated that each locus independently controls flower phenotype.

370         Developed DNA markers J01 and S01 could successfully identify recessive double flower

371    alleles for $D_{jo}$ and $D_{su}$, respectively. Both markers showed high fitting ratio with phenotype and were

372    applicable to the examined *H. macrophylla* accessions. The S01 marker is superior to the DNA marker

373    STAB045 linked to $D_{su}$ and which was discovered by Waki et al.[5] because the former has a wide range

374    of applicability. While the S01 marker genotype completely fitted with the phenotype in all tested

375    accessions, STAB045 did not (data not shown). Because both J01 and S01 showed a wide range of

376    applicability, it is advantageous to use them in combination to reveal the existence of the double flower

377    allele in *H. macrophylla* accessions. This information will help in selection of candidate parents with

378    heterozygous recessive double flower alleles to obtain double flower progenies. In addition, these

379    DNA markers should be useful in marker assisted selection (MAS) of double flower progenies. To

380    obtain double flower progenies, at least the paternal parent should be of the single flower phenotype

381    because very few or none at all pollen grains are produced in double flower individuals. In addition,

382    it requires approximately 2 years to confirm the flower phenotype from the time of crossing.

383    Identification of flower phenotype at the seedling stage by MAS would enable the discarding of single

384    flower individuals and allow the growth of double flower individuals. The developed DNA markers

385    should accelerate the breeding of double flower phenotypes.

386         In the genomic sequence of 'Kirakiraboshi,' an insertion was detected in the second exon of

387    the *LFY* gene. This insertion actually resulted in frameshift of cloned mRNA in both splice variants.

388    Therefore, it was speculated that the function of the *LFY* gene was suppressed or lost in 'Kirakiraboshi'.

389    The *LFY* gene and its homologue *FLO* have been identified in many plants, such as *Arabidopsis*

390    *thaliana* and *Antirrhinum majus*, and are known as transcription factors for major flowering signals[29-

391    31]. Additionally, many types of phenotypes in *Arabidopsis lfy* mutants have been reported[32,33]. In the

392    *lfy* strong phenotype, most organs are sepal-like, or mosaic sepal/carpels organs, and the sepal-like

393    organs are characteristic of wild-type cauline leaves[33]. Therefore, the flowers of the *lfy* mutant

394    appeared to be double flowers that are formed from leaves or sepals. Additionally, a similar phenotype

395    has been reported in *LFY* homologue mutants or transgenic plants such as the *flo* mutant of

396    *Antirrhinum majus*[34], *uni* mutant of pea[35], and co-suppressed *NFL* transgenic plant of tobacco[37].

397    Therefore, generally, when the *LFY* gene function is lost, petal, stamens, and a carpel are likely to be

398    replaced by sepal-like organs. In decorative flowers of hydrangea, sepals show petaloid characteristics

399    including pigmentation and enlarged organ size. It is possible that sepal-like organs in decorative

400    flowers show petaloid characteristics and form double flowers. Therefore, we assumed that *LFY* is a

401    causative gene of the double flower phenotype of 'Sumidanohanabi'.

402         However, there remain several unexplained observations in this study. The double flower of

403    'Kirakiraboshi' did not exhibit the exact same phenotype of the *lfy* mutant. Generally, the flowers of

404    *lfy* or its orthologous gene mutants have only leaf-like or sepal-like organs that have chlorophyll,

405    stomata, and trichome, and these organs have almost no petal identity[33,34]. When flowering signals in

406    *lfy* mutant were lost completely, floral organs were not fully formed[33-35]. In the double flowers of

407    'Kirakiraboshi', the floral organs keep their petal identity, have papilla cells, and are pink or blue.

408    These phenotypes of 'Kirakiraboshi' might reflect partial remaining of LFY function. Additionally, it

409    has been reported that *lfy* mutants with an intermediate or weak phenotype sometimes develop petaloid

410    organs[33]. According to the genomic sequence of *H. macrophylla*, no other *LFY* gene was observed. It

411    could be considered that the double flowers of 'Kirakiraboshi' were induced via partial repression of

412    the LFY function.

413    On the contrary, we could not find any candidate gene that controls the double flower

414    phenotype for the $D_{jo}$ locus. One possible reason was that SNPs were not called in scaffold with

415    causative gene. In pseudomolecules, about half of the total scaffolds length was not included since

416    relevant SNPs were not called. Improvement of SNP density would be effective for discovering

417    additional scaffolds that are tightly linked to $D_{jo}$. Although candidate gene for $D_{jo}$ could not be

418    identified from the linkage information, we predicted several candidate genes. In hydrangea, stamens

419    and petals were absent from decorative flowers of the double flower plant, and there was an increased

420    number of sepals[4]. Since causative genes should explain the changes in formation, the B-class genes

421    of the ABC model, *PI* and *AP3*, were predicted as candidate genes. In *A. thaliana*, the B-class gene *pi*

422    or *ap3* mutants showed an increase in the number of sepals converted from petals[37]. If these genes

423    were mutated in hydrangea, an increase in sepals would be expected. In hydrangea, *HmPI*, *HmAP3*,

424    and *HmTM6* were identified as B-class genes[38,39]. As *HmAP3* was located on CHR13, it was not

425    considered as a causative gene for $D_{jo}$. In this study, *HmPI* and *HmTM6* were not included in the

426    pseudomolecule. Ascertaining the loci of these genes might reveal the causative gene for $D_{jo}$.

427    In this study, we report DNA markers and possible causative genes for the double flower

428    phenotype observed in two hydrangea cultivars. For this analysis, we established a reference sequence

429    for the hydrangea genome using advanced sequencing technologies including the long-read

430    technology (PacBio) and the HiC method[9], bioinformatics techniques for the diploid genome

431    assembly[14], and haplotype phasing[8]. To the best of our knowledge, this is the first report on the

432    chromosome-level haplotype-phased sequences in hydrangea at the level of the species (*H.*

433    *macrophylla*), genus (*Hydrangea*), family (Hydrangeaceae), and order (Cornales). The genomic

434    information from this study based on NGS technology is a significant contribution to the genetics and

435    breeding of hydrangea and its relatives. It will serve to accelerate the knowledge base of the evolution

436    of floral characteristics in Hydrangeaceae.

437

441

444

445    **Supporting information:**

446    **Supplementary Table S1.** RNA samples used for Iso-Seq and RNA-Seq

447    **Supplementary Table S2.** Statistics of the genome sequences of *Hydrangea macrophylla*

448    'Aogashima-1'

449    **Supplementary Table S3.** J01 marker genotypes and double flower phenotypes of 15IJP1 population.

450    **Supplementary Table S4.** J01 marker genotypes and double flower phenotypes of 14GT77

451    population.

452    **Supplementary Figure S1.** Alignment of *LFY* genomic sequence and CDS.

453

454    **Data availability:**

455    The sequence reads are available from the DNA Data Bank of Japan (DDBJ) Sequence Read Archive

456    (DRA) under the accession numbers DRA010300, DRA010301, and DRA010302. The assembled

457    sequences are available from the BioProject accession number PRJDB10054. The genome information

458    is available at Plant GARDEN (https://plantgarden.jp).

459

460    **References**

461    1. Uemachi, T., Kato, Y., and Nishio, T. 2004, Comparison of decorative and non-decorative flowers

462    in *Hydrangea macrophylla* (Thunb.) Ser., *Sci. Hortic*., 102, 325–334

463

464    2. Uemachi, T., Kurokawa, M., and Nishio, T. 2006, Comparison of inflorescence composition and

465    development in the lacecap and its sport, hortensia *Hydrangea macrophylla* (Thunb.) Ser., *J. Japan.*

466    *Soc. Hort. Sci.,* 75, 154–160.

467

468    3. Uemachi, T. and Okumura, A. 2012, The inheritance of inflorescence types in *Hydrangea*

469    *macrophylla*, *J. Japan. Soc. Hort. Sci*., 81, 263–268.

470

471    4. Suyama, T., Tanigawa, T., Yamada, A. et al. 2015, Inheritance of the double-flowered trait in

472    decorative hydrangea flowers, *Hortic. J*., 84, 253-260.

473

474    5. Waki, T., Kodama, M., Akutsu, M. et al. 2018, Development of DNA markers linked to double-

475    flower and hortensia traits in *Hydrangea macrophylla* (Thunb.) Ser., *Hortic J*., 87, 264-273.

476

477    6. Heijmans, K., Ament, K., Rijpkema, A.S. et al. 2012, Redefining C and D in the petunia ABC, *Plant*

478    *Cell*, 24, 2305-2317.

479

480    7. Tränkner, C., Krüger, J., Wanke, S., Naumann, J., Wenke, T. and Engel, F. 2019, Rapid identification

481    of inflorescence type markers by genotyping-by-sequencing of diploid and triploid F1 plants of

482    *Hydrangea macrophylla*, *BMC Genet.*, 20, 60.

483

484    8. Kronenberg, Z. N., Hall, R. J., Hiendleder, S., et al. 2018, FALCON-Phase: Integrating PacBio

485    and Hi-C data for phased diploid genomes, *BioRxiv*, 327064.

486

487    9. Dudchenko, O., Batra, S. S., Omer, A. D., et al. 2017, De novo assembly of the *Aedes aegypti*

488    genome using Hi-C yields chromosome-length scaffolds. *Science*, 356, 92-95.

489

490    10. Mascher, M. and Stein, N. 2014, Genetic anchoring of whole-genome shotgun assemblies,

491    *Front Genet*, 5, 208.

492

493    11. Marcais, G. and Kingsford, C. 2011, A fast, lock-free approach for efficient parallel counting

494    of occurrences of k-mers, *Bioinformatics*, 27, 764-770.

495

496    12. Kajitani, R., Toshimoto, K., Noguchi, H., et al. 2014, Efficient de novo assembly of highly

497    heterozygous genomes from whole-genome shotgun short reads, *Genome Res*, 24, 1384-1395.

498

499    13. Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., and Zdobnov, E. M. 2015,

500    BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs,

501    *Bioinformatics*, 31, 3210-3212.

502

503    14. Chin, C. S., Peluso, P., Sedlazeck, F. J., et al. 2016, Phased diploid genome assembly with

504    single-molecule real-time sequencing, *Nat Methods*, 13, 1050-1054.

505

506    15. Walker, B. J., Abeel, T., Shea, T., et al. 2014, Pilon: an integrated tool for comprehensive

507    microbial variant detection and genome assembly improvement, *PLoS One*, 9, e112963.

508

509    16. Shirasawa, K., Hirakawa, H., and Isobe, S. 2016, Analytical workflow of double-digest

510    restriction site-associated DNA sequencing based on empirical and in silico optimization in

511    tomato, *DNA Res*, 23, 145-153.

512

513    17. Rastas, P. 2017, Lep-MAP3: robust linkage mapping even for low-coverage whole genome

514    sequencing data, *Bioinformatics*, 33, 3726-3732.

515

516    18. Tang, H., Zhang, X., Miao, C., et al. 2015, ALLMAPS: robust scaffold ordering based on

517    multiple maps, *Genome Biol*, 16, 3.

518

519    19. Grabherr, M. G., Haas, B. J., Yassour, M., et al. 2011, Full-length transcriptome assembly

520    from RNA-Seq data without a reference genome, *Nat Biotechnol*, 29, 644-652.

521

522    20. Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S., and Morgenstern, B. 2006,

523    AUGUSTUS: ab initio prediction of alternative transcripts, *Nucleic Acids Res*, 34, W435-439.

524

525    21. Kent, W. J., 2002, BLAT - the BLAST-like alignment tool, Genome Res, 12, 656-664.

526

527    22. Ghelfi, A., Shirasawa, K., Hirakawa, H., and Isobe, S. 2019, Hayai-Annotation Plants: an

528    ultra-fast and comprehensive functional gene annotation system in plants, *Bioinformatics*, 35,

529    4427-4429.

530

531    23. Bolger, A.M., Lohse, M., and Usadel, B. 2014, Trimmomatic: a flexible trimmer for Illumina

532    sequence data, *Bioinformatics*, 30, 2114-2120.

533

534    24. Li, H., Handsaker, B., Wysoker, A., et al. 2009, The Sequence Alignment/Map format and

535    SAMtools, *Bioinformatics*, 25, 2078-2079.

536

537    25. Untergasser, A., Cutcutache, I., Koressaar, T. et al. 2012, Primer3--new capabilities and interfaces.

538    *Nucleic Acids Res*., 40, e115.

539

540    26. Schmieder, R. and Edwards, R. 2011, Quality control and preprocessing of metagenomic

541    datasets, *Bioinformatics*, **27**, 863-864.

542

543    27. Langmead, B. and Salzberg, S. L. 2012, Fast gapped-read alignment with Bowtie 2, *Nat*

544    *Methods*, 9, 357-359.

545

546    28. Danecek, P., Auton, A., Abecasis, G., et al. 2011, The variant call format and VCFtools,

547    *Bioinformatics*, 27, 2156-2158.

548

549    29. Jaeger, K.E., Pullen, N., Lamzin, S., Morris, R.J., and Wigge, P.A. 2013, Interlocking feedback

550    loops govern the dynamic behavior of the floral transition in Arabidopsis, *Plant Cell*, 25, 820–

551    833.

552

553    30. Krizek, B.A. and Fletcher, J.C. 2005, Molecular mechanisms of flower development: an

554    armchair guide, *Nat. Rev. Genet.*, 6, 688.

555

556    31. William, D.A., Su, Y., Smith, M.R., Lu, M., Baldwin, D.A., and Wagner, D. 2004, Genomic

557    identification of direct target genes of LEAFY, *Proc. Nat. Acad. Sci.*, 101, 1775–1780.

558

559    32. Okamuro, J.K., Den Boer, B.G., and Jofuku, K.D. 1993, Regulation of Arabidopsis flower

560    development, *Plant Cell*, 5, 1183-1193.

561

562    33. Weigel, D., Alvarez, J., Smyth, D.R., Yanofsky, M.F., and Meyerowitz, E.M. 1992, LEAFY

563    controls floral meristem identity in Arabidopsis, *Cell*, 69, 843-859.

564

565    34. Carpenter, R. and Coen, E.S. 1990, Floral homeotic mutations produced by transposon-

566    mutagenesis in *Antirrhinum majus*, *Gene. Dev.*, 4, 1483–1493.

567

568    35. Hofer, J., Turner, L., Hellens, R. et al. 1997, UNIFOLIATA regulates leaf and flower

569    morphogenesis in pea, *Curr. Biol.*, 7, 581–587.

570

571    36. Ahearn, K.P., Johnson, H.A., Weigel, D., and Wagner, D.R. 2001, NFL1, a *Nicotiana tabacum*

572    LEAFY-like gene, controls meristem initiation and floral structure, *Plant Cell Physiol.*, 42, 1130–

573    1139.

574

575    37. Bowman, J.L., Smyth, D.R., and Meyerowitz, E.M. 1989, Genes directing flower development in

576    Arabidopsis, *Plant Cell*, 1, 37-52.

577

578    38. Kitamura, Y., Hosokawa, M., Uemachi, T., and Yazawa, S. 2009, Selection of ABC genes for

579    candidate genes of morphological changes in hydrangea floral organs induced by phytoplasma

580    infection, *Sci. Hort.*, 122, 603-609.

581

582    39. Kramer, E.M. and Irish, V.F. 2000. Evolution of the petal and stamen development programs:

583    Evidence from comparative studies of the lower eudicots and basal angiosperms, *Int. J. Plant Sci.*, 161,

584    s29-s40

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

Table 1. SNPs correlated (fitting rate more than 95%) with double flower phenotype in 12GM1 population

| Scaffold | Position at Phase 0 | Sequence variant | | Fitting rate (%) | Frequency of double flower phenotype (double flower/all) | | |
|---|---|---|---|---|---|---|---|
| | | Posy Bouquet Grace | Blue Picotee manasulu | | Homozygous of 'Posy Bouquet Grace' | Heterozygous | Homozygous of 'Blue Picotee Manasulu' |
| 0008F-2 | 3250598 | A | G | 100 | 37/37 | 0/61 | 0/47 |
| 0008F-2 | 3250523 | A | C | 100 | 37/37 | 0/61 | 0/47 |
| 0008F-2 | 780104 | C | A | 100 | 37/37 | 0/60 | 0/48 |
| 0259F | 404610 | T | A | 100 | 37/37 | 0/60 | 0/48 |
| 1207F | 365533 | C | T | 100 | 38/38 | 0/61 | 0/48 |
| 1207F | 372121 | C | A | 100 | 38/38 | 0/61 | 0/47 |
| 0012F | 1318350 | T | C | 97.9 | 37/39 | 1/59 | 0/48 |
| 0437F | 170787 | G | A | 97.9 | 36/37 | 1/60 | 1/49 |
| 0437F | 180821 | A | G | 97.9 | 36/37 | 1/60 | 1/49 |
| 0994F | 216439 | C | T | 97.9 | 36/37 | 1/60 | 1/49 |

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

Table 2. SNPs correlated (fitting rate more than 95%) with double flower phenotype in KF population

| Scaffold | Position at Phase 0 | Sequence variant | | Fitting rate (%) | Frequency of double flower phenotype (double flower/all) | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | Kirakiraboshi | Frau Yoshimi | | Homozygous of 'Kirakiraboshi' | Heterozygous | Homozygous of 'Frau Yoshimi' |
| 0577F | 1204837 | AG | AAACATG | 98.9 | 22/22 | 0/51 | 1/20 |
| 3145F | 55089 | TA | TAA | 98.9 | 22/22 | 0/51 | 1/20 |
| 3145F | 55109 | G | A | 98.9 | 22/22 | 0/51 | 1/20 |
| 3145F | 55446 | G | A | 98.9 | 22/22 | 0/51 | 1/20 |
| 0109F | 868569 | C | G | 95.7 | 22/25 | 0/44 | 1/24 |

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

Table 3. Genotypes of DNA marker J01 and S01 in *H. macrophylla* varieties

| Accession name | Phenotype | Genotype | |
|---|---|---|---|
| | | J01 | S01 |
| Jogasaki | Double | 117_50/117_50 | 250/280 |
| Posy Bouquet Grace | Double | 117_50/117_50 | 280/280 |
| Izunohana | Double | 117_50/117_50 | 250/280 |
| Chikushinokaze | Double | 117_50/117_50 | 250/280 |
| Chikushinomai | Double | 117_50/117_50 | 280/280 |
| Chikushiruby | Double | 117_50/117_50 | 280/280 |
| Corsage | Double | 117_50/117_50 | 280/280 |
| Dance Party | Double | 117_50/117_50 | 280/280 |
| Fairy Eye | Double | 117_50/117_50 | 250/280 |
| Posy Bouquet Casey | Double | 117_50/117_50 | 250/280 |
| Sumidanohanabi | Double | 167/167 | 236/236 |
| Kirakiraboshi | Double | 167/167 | 236/236 |
| HK01 | Double | 167/167 | 236/236 |
| HK02 | Double | 167/167 | 236/236 |
| 03JP1 | Single | 117_50/167 | 280/280 |
| Amethyst | Single | 167/167 | 250/280 |
| Blue Picotee Manaslu | Single | 167/167 | 280/280 |
| Blue Sky | Single | 167/167 | 280/280 |
| Bodensee | Single | 167/167 | 250/250 |
| Chibori | Single | 167/167 | 280/280 |
| Furau Mariko | Single | 167/167 | 250/250 |
| Furau Yoshiko | Single | 167/167 | 280/280 |
| Furau Yoshimi | Single | 167/167 | 250/280 |
| Green Shadow | Single | 167/167 | 280/280 |
| Kanuma Blue | Single | 167/167 | 250/280 |
| Mrs. Kumiko | Single | 167/167 | 280/280 |
| Paris | Single | 167/167 | 280/280 |
| Peach Hime | Single | 167/167 | 280/280 |
| Picotee | Single | 167/167 | 282/282 |
| Ruby Red | Single | 167/167 | 280/280 |
| Shinkai | Single | 167/167 | 280/280 |
| Tokimeki | Single | 167/167 | 280/282 |
| Uzuajisai | Single | 167/167 | 250/280 |

635    Genotypes shown as gray indicate homozygous of double flower allele.
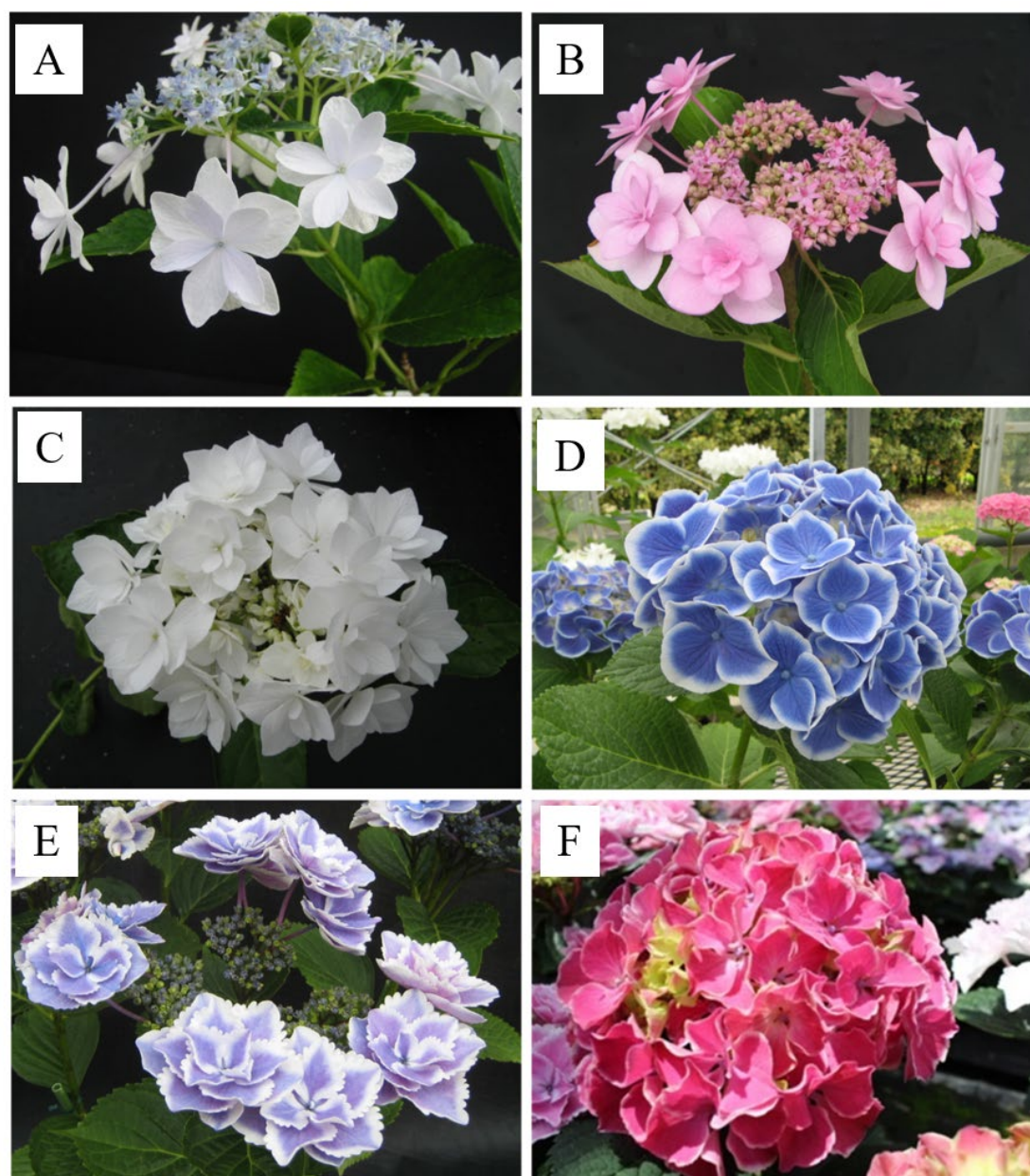
636

637    Figure 1. Flower phenotypes of hydrangea accessions

638    A: 'Sumidanohanabi' (double flower). B: 'Jogasaki' (double flower). C: 'Posy Bouquet Grace' (double

639    flower). D: 'Blue Picotee Manaslu' (single flower). E: 'Kirakiraboshi' (double flower). F: 'Frau
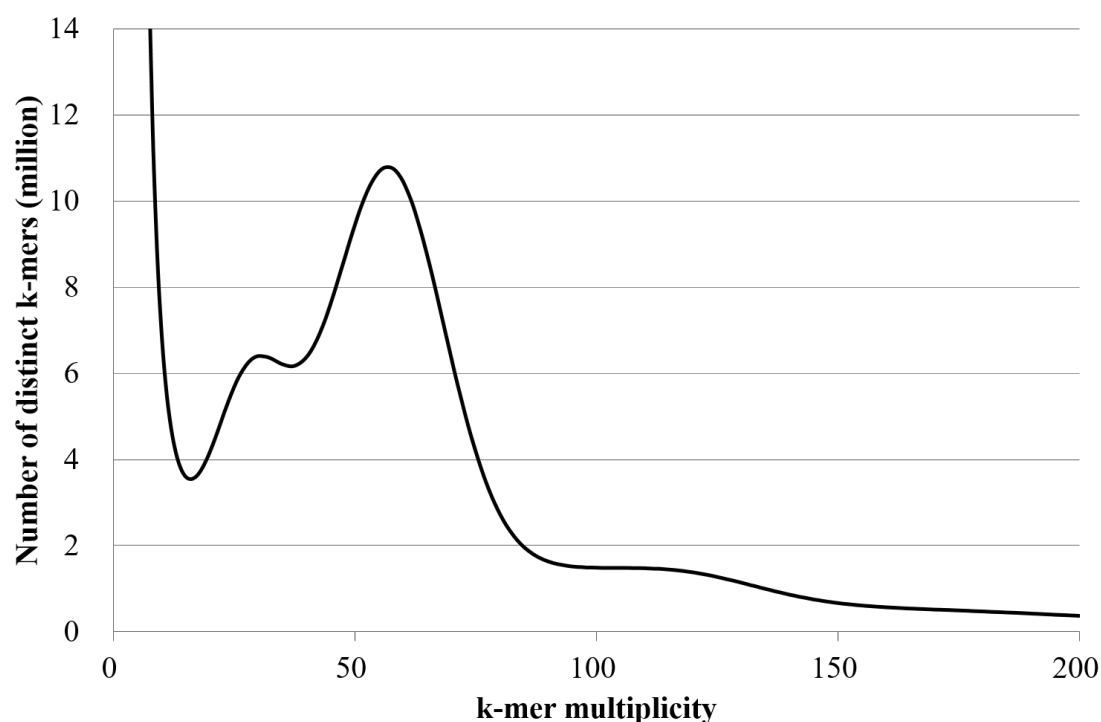
640    Yoshimi' (single flower).

641

642

643    Figure 2. Genome size estimation for the hydrangea line 'Aogashima-1' with the distribution of the

644    number of distinct $k$-mers ($k$=17), with the given multiplicity values.

645

646

647

648

649

650

651

652

653

654

655

L   P1   P2   F1

F2
Double flower

F2
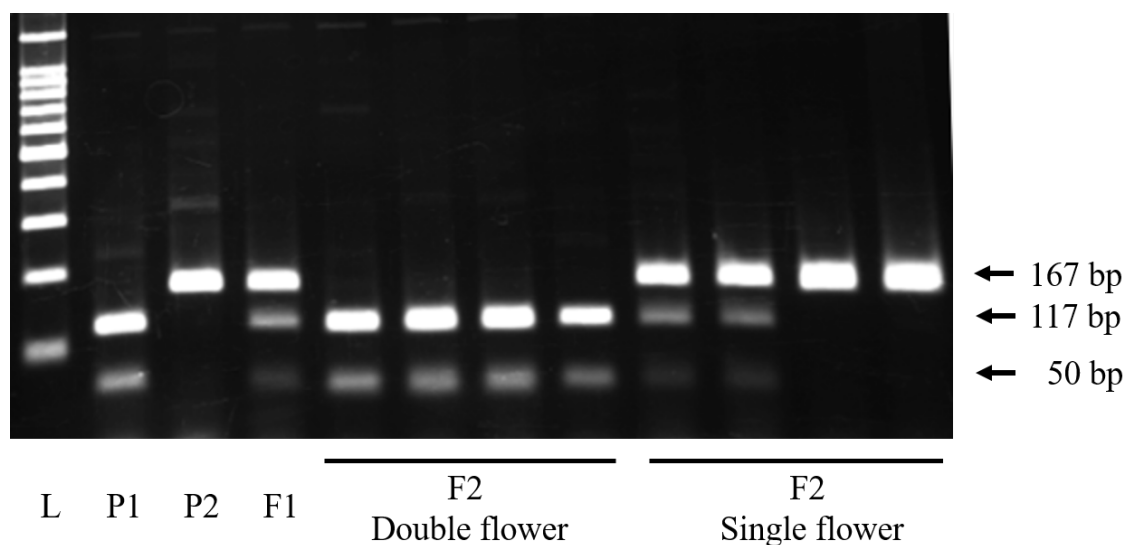Single flower

← 167 bp
← 117 bp
← 50 bp

656

657    Figure 3. Fragment pattern of J01 DNA marker

658    Dominant single flower allele is shown as undigested 167 bp fragment. Recessive double flower allele

659    is shown as digested 117 and 50 bp fragments. L: 100 bp ladder, P1: 'Posy Bouquet Grace'

660    (117_50/117_50), P2: 'Blue Picotee Manaslu' (167/167).

661

662

663

664

665

666

667
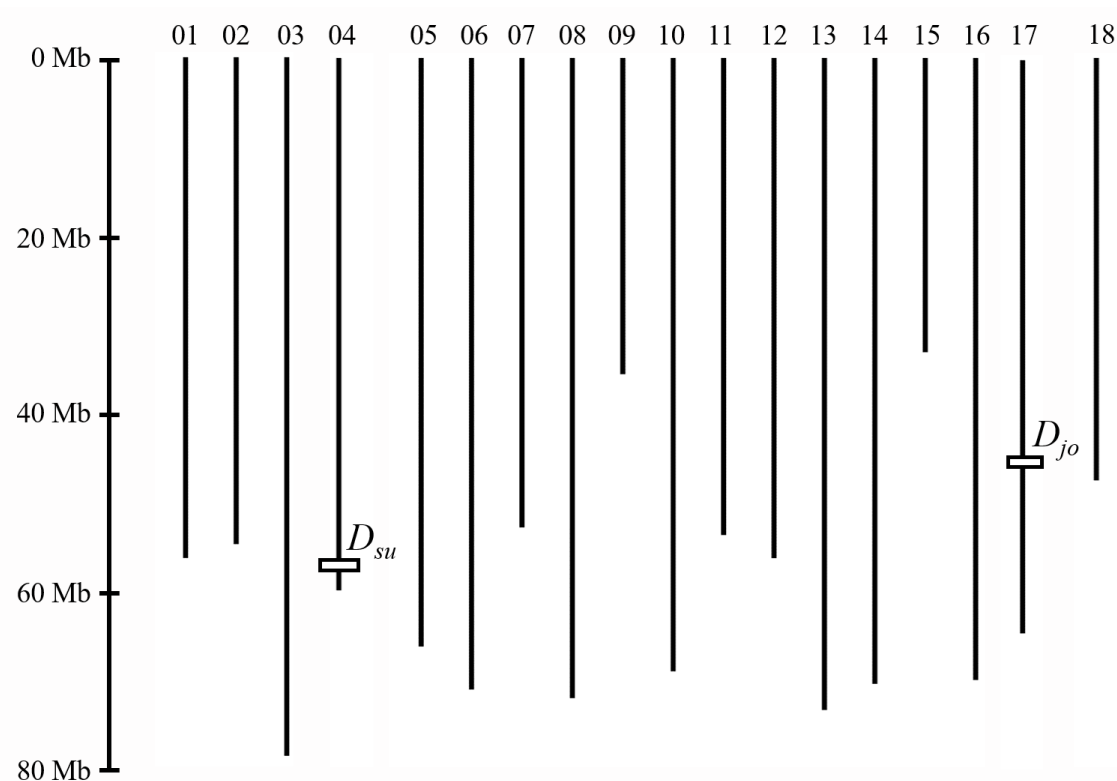
668

669

670

671

672

673    Figure 4. Schematic model of pseudomolecules

674    Double flower phenotype controlling loci $D_{su}$ and $D_{jo}$ are shown. $D_{jo}$ is shown as J01 marker position

675    46,326,384 in CHR17. $D_{su}$ is shown as tightly linked SNP at 0109F_868569, since the S01 marker
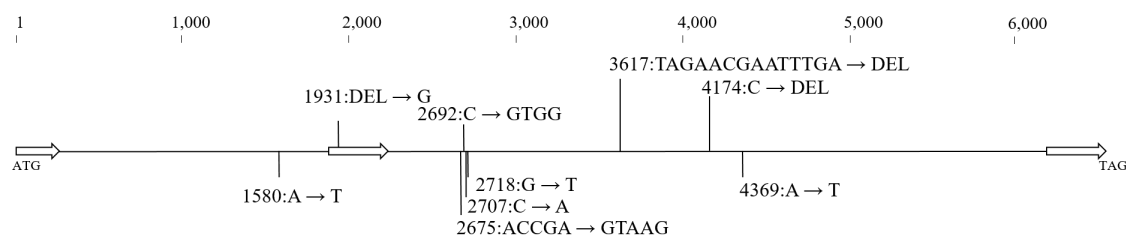
676    sequence was not on the pseudomolecule.

677

678

679

680

681

682

683

684

685

686

687    Figure 5. DNA polymorphisms in *LFY* genomic sequence

688    *LFY* sequence polymorphisms observed specifically in 'Kirakiraboshi' genomic sequence

689    The sequence is started from the initiation codon (ATG) at 678,200 to the termination signal (TAG) at

690    684,639 in phase 1 sequence of 0577F of HMA_r1.2. White arrows indicate coding sequences, CDS1:

691    1 to 454 bp, CDS2: 1,888 to 2,255 bp, CDS3: 6,078 to 6,440 bp. Genetic variants are shown as from

692    Hma1.2 sequence to 'Kirakiraboshi'.

693
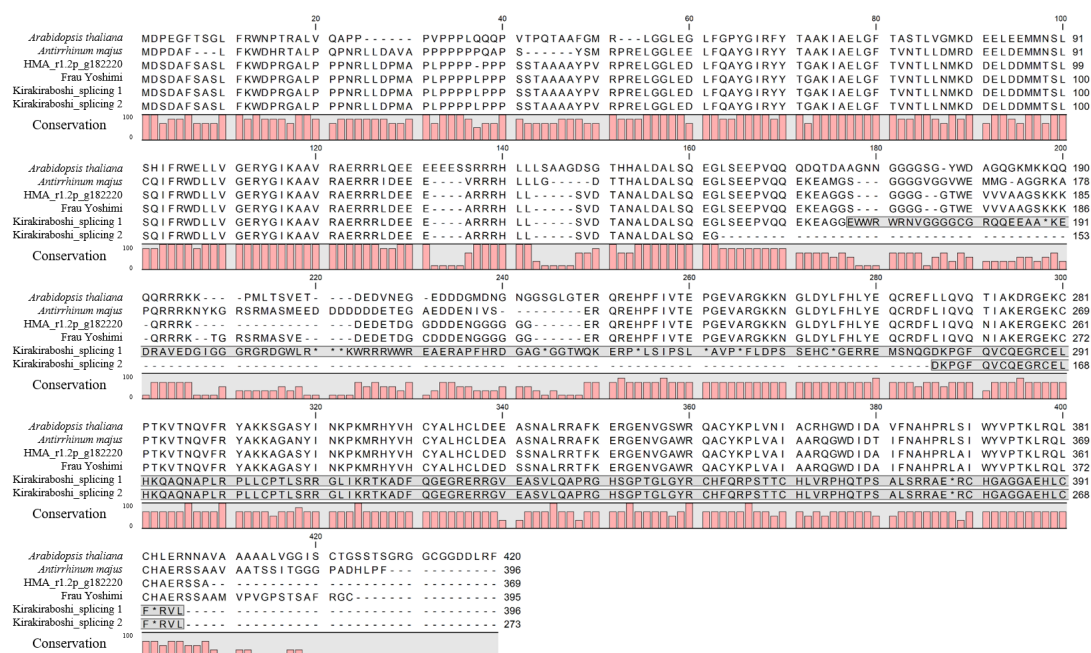
694

695

696

697

698

699

700

701

702

703

704

705 Figure 6. Alignment of LFY protein sequences

706 Amino acids with gray background show frameshifted regions. Splicing variant was observed, and

707 both sequences showed frameshift in 'Kirakiraboshi'. *Arabidopsis thaliana:* ABE66271.1 *Antirrhium*

708 *majus*: AAA62574.1.

709

710

711

712

713

714

715

716

717

A



L    P1    P2        Double flower        Single flower

B



Forward primer                                    Reverse primer

280 bp ($D_{su}$)

250 bp ($D_{su}$)

236 bp ($d_{su}$)

3464                                              3743

3617:TAGAACGAATTTGA → DEL

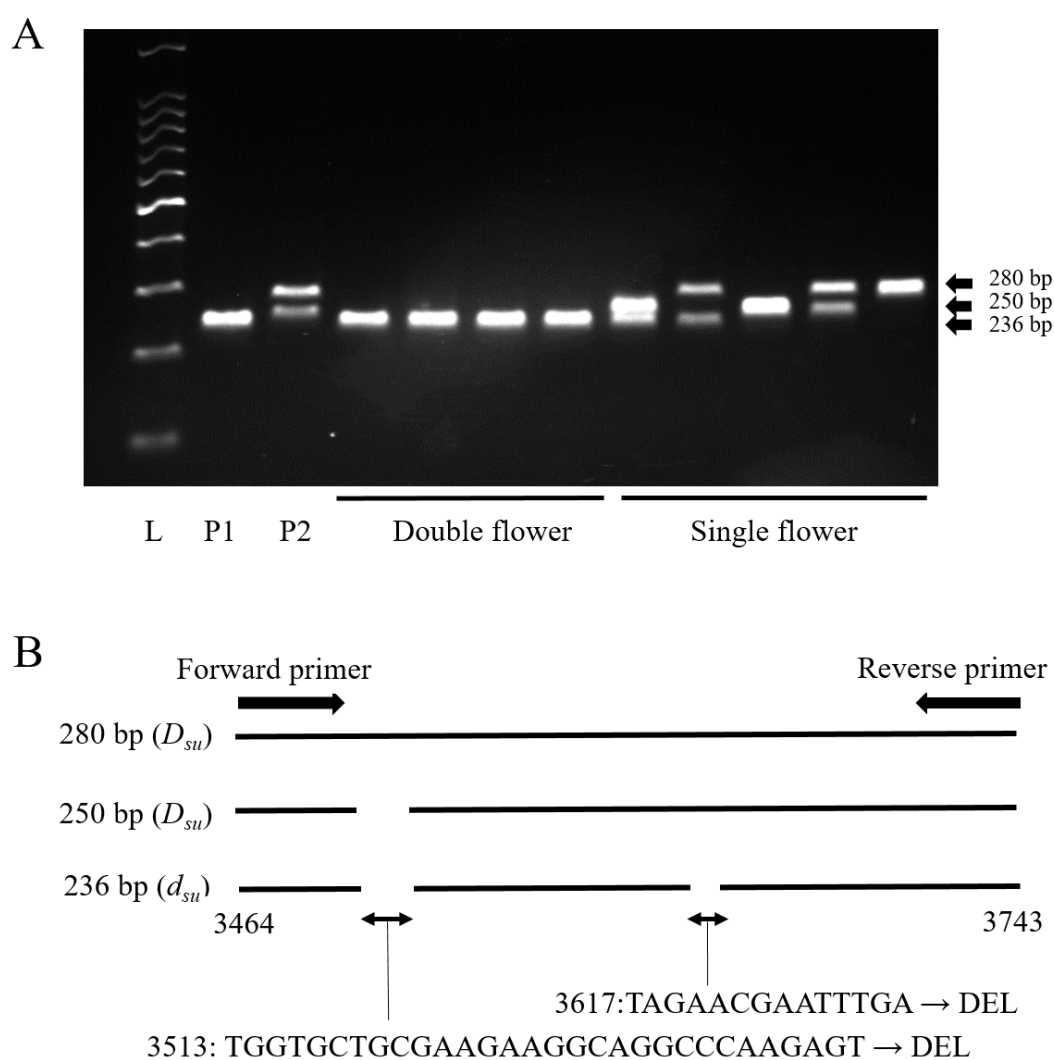3513: TGGTGCTGCGAAGAAGGCAGGCCCAAGAGT → DEL

718

719    Figure 7. Fragment pattern of S01 DNA marker

720    A. Fragment pattern of S01 DNA marker. Dominant single flower alleles are shown as 250 bp and 280

721    bp fragments. Recessive double flower allele is shown as 236 bp fragments. L: 100 bp ladder, P1:

722    'Kirakiraboshi' (236/236), P2: 'Frau Yoshimi' (250/280).

723    B. INDEL polymorphisms in alleles of DNA marker S01 amplified sequences. Position on schematic

724    models were the same as in Figure 5.