**Whole genome sequencing of a wild yam species *Dioscorea tokoro* reveals a genomic region associated with sex**

Satoshi Natsume[1], Hiroki Yaegashi[1], Yu Sugihara[2], Akira Abe[1], Motoki Shimizu[1], Kaori Oikawa[1],

Benjamen White[3], Aoi Kudoh[2], Ryohei Terauchi[1,2*]

[1] Iwate Biotechnology Research Center, Kitakami, Iwate, 024-0003, Japan

[2] Laboratory of Crop Evolution, Kyoto University, Mozume, Muko, Kyoto, 617-0001, Japan

[3] Earlham Institute, Norwich NR4 7UZ, United Kingdom

*Correspondence to  Ryohei Terauchi  (terauchi@ibrc.or.jp)

**Abstract**

*Dioscorea tokoro* is a wild species distributed in East Asia including Japan. Typical of the genus *Dioscorea*, *D. tokoro* is dioecious with male and female flowers borne on separate individuals. To understand its sex determination system and to serve as a model species for population genomics of obligate outcrossing wild species, we set out to determine the whole genome sequence of the species. Here we show 443 Mb genome sequence of *D. tokoro* distributed over 2,931 contigs that were anchored on 10 linkage groups. Linkage analysis of sex in a segregating F1 family revealed a sex determination locus residing on Pseudochromosome 3 with XY-type male heterogametic sex determination system.

**key words**

(Keywords. *Dioscorea*, yam, dioecy, genome, sex determination.)

**Introduction**

The genus *Dioscorea* belongs to the monocotyledons and has 450 - 600 species distributed mainly in tropical and subtropical area of the world (Coursey, 1972; Sugihara et al. 2021). Cultivated species of *Dioscorea* are collectively called yam, which includes Guinea yam (*D. rotundata*) of West Africa accounting for more than 90% of the world yam production (FAOSTAT, 2018). The entire genus of *Dioscorea* is characterized by dioecy, with male and female flowers borne on separate individuals. Consequently, the species of *Dioscorea* have obligate outcrossing, resulting in a higher level of heterozygosity and frequent inter-species hybridization (Terauchi et al. 1992; Chaïr et al. 2010, 2016;

38  Girma et al. 2014; Scarcelli et al. 2006, 2017; Siadjeu et al. 2018; Sugihara et al. 2020, 2021; Bredeson

39  et al. 2022). Previously we reported the whole genome sequence of *D. rotundata* with 570 Mb in size

40  (Tamiru et al. 2017), which served as a reference to study population genomics of *D. rotundata* and

41  its wild relatives to reveal the origin of Guinea yam (Scarcelli et al. 2019; Sugihara et al. 2020).

42  *Dioscorea tokoro* is a wild species belonging to the section Stenophora. It is widely

43  distributed in East Asia including Japan. *D. tokor*o is a diploid species with a chromosome number 2n

44  = 2x = 20. It is a perennial species with rhizomes. In spring, shoots emerge from rhizomes and develop

45  to vines that twine around nearby trees in an anti-clockwise direction which expand alternate leaves

46  (**Fig 1**). The species commonly occurs along the fringes of forests in Japan. Crossing experiment is

47  easy and the generation time is relatively short (1-2 years), so that the species has been serving as a

48  model species of *Dioscorea*. The species was subjected to studies of population genetics (Terauchi

49  1990 : Terauchi and Konuma 1994; Terauchi et al.1997), linkage mapping and elucidation of sex

50  determination mechanisms (Terauchi and Kahl, 2004).

51  To serve as a platform for future genomics study of the species, we here report the whole

52  genome sequence of *D. tokoro*. We combined Oxford Nanopore long read sequences and Illumina

53  short read sequences for de novo assembly to generate contigs. The contigs were further anchored on

54  linkage maps to generate pseudochromosomes. RNA-seq data were used for gene prediction. A

55  putative locus involved in *D. tokoro* sex determination was identified on Pseudochromosome 3.

56

57  **Results**

58

59  **Estimation of size of *D. tokoro* genome by flow cytometry**

60

61  We used a *D. tokoro* individual Kita1 collected at Kitakami, Iwate, Japan as well as *D. rotundata*

62  accession TDR96-F1 with known genome size (~570 Mb, Tamiru et al. 2017), as the material for flow

63  cytometry (FCM) analysis using nuclei prepared from fresh leaf samples. DNA of isolated nuclei were

64  stained with propidium iodide and analyzed by a flow cytometer. The value of G1 peak mean of *D.*

65  *tokoro* was 206.5, whereas that of *D. rotundata* was 303.6. The ratio between the two species was 0.68

66  (206.5/303.6). From these values the genome size of *D. tokoro* was estimated to be ~388 Mb (570 Mb

67  × 0.68) (**Fig. S1**).

68

69  **Reference assembly using Oxford Nanopore Technology**

70

71  Genomic DNA was extracted from fresh leaves of *D. tokoro* Kita1 and subjected to Oxford Nanopore

72  Technologies (ONT) sequencing. As a result, we obtained a total of 2,515,235 reads amounting 27.4

73  Gb in size (**Table S1**). We also performed Illumina sequencing of 35 - 251 bp read-length (total 24.6

74    Gb; obtained by MiSeq) as well as 150 bp read-length (total 37.8 Gb; obtained by HiSeq4000) (**Table**

75    **S2**). We assembled ONT reads and Illumina sequence reads using a hybrid assembler MaSuRCA

76    v3.3.4 (Zimin et al. 2013) with Flye assembler v2.6 (Kolmogorov et al. 2019) running internally,

77    which generated *D. tokoro* draft genome sequence consisting of 2,931 contigs amounting 443.5 Mb

78    with $N_{50}$ being 586,368 bp (**Table 1**). The estimated genome size by *k*-mer analysis of the reads was

79    438.7 Mb. These estimated genome size based on DNA sequencing were larger than 388 Mb as

80    estimated by the FCM analysis.

81

82    **Anchoring of contigs on *D. tokoro* linkage maps**

83    To generate *D. tokoto* pseudochromosomes, we mapped the contigs on ten linkage groups. For this

84    purpose, we crossed a female individual Waka1 (P1) with a male individual Kita1 (P2) to obtain F1

85    progeny comprising 186 individuals (**Fig. S2**; **Table S3**). These plants were genotyped by RAD

86    markers (**Fig. S3**; Baird et al, 2008). We identified 946 SNPs and 180 presence/absence

87    polymorphisms (PAs) that are heterozygous in P1 and homozygous in P2 parents, and 724 SNPs and

88    880 PAs that are homozygous in P1 and heterozygous in P2 parents (**Table S4**). These DNA markers

89    were used for construction of linkage map using pseudo-testcross approach (Grattapaglia and Sederoff,

90    1994). We obtained two linkage maps, one for DNA markers heterozygous in P1 parent, and the other

91    for markers heterozygous in P2 parent (**Fig. S4**). Since each RAD marker has ~75 bp sequence, this

92    information was used to associate RAD marker to contigs generated by the de novo assembly (**Fig.**

93    **S5**). This method allowed us to anchor contigs amounting 378.8 Mb (85.4% of the total genome size)

94    to the linkage maps (**Table S5**) and to combine the two linkage groups and generate

95    pseudochromosomes 1-10 with sizes ranging from 31.5 Mb (Pseudochromosome 5) to 54.6 Mb

96    (Pseudochromosome 1) (**Fig. 2**).

97        BUSCO analysis (Mosè et al. 2020) showed that complete BUSCO value of 98%, indicating

98    that our *D. tokoro* genome sequence is of a sufficient quality as the reference (**Table 1**).

99

100    **Gene prediction**

101

102    We performed RNA-seq of 18 samples representing different organs and developmental stages of *D.*

103    *tokoro* (**Table S6**). The total size of RNA-seq reads amounted 31.17 Gb. These RNA-seq reads were

104    mapped to the contigs, revealing a total of 29,084 genes, among which 25,447 genes were assigned to

105    pseudochromosomes (**Table 2**).

106

107    **Sex determination in *D. tokoro***

108

109 The 186 F1 progeny derived from a cross between Waka1 female (P1) and Kita1 male (P2) parents
110 segregated in 38 female, 89 male, and 59 non-flowering in 2011 (**Table S3**). We attempted to identify
111 genomic region that shows association with sex of the F1 individuals. As a result of Fisher's exact test
112 based on the sex and genotype contingency table of each progeny, we found a significant association
113 of the middle position of Pseudochromosome 3 with sex when the DNA markers heterozygous in the
114 male parent (P2) were used. By contrast, there was no association detected if we use the markers
115 heterozygous in the female parent (P1; **Fig. 3**). This result indicates a male heterogametic sex
116 determination (XY) system in *D. tokoro*, and supports our previous analysis with AFLP markers
117 (Terauchi and Kahl, 1999).

118

119 **Discussion**

120

121 Here we report *D. tokoro* draft genome sequence of 443.5 Mb in size. For the species of the genus
122 *Dioscorea,* whole genome sequences are available for *D. rotundata* (Tamiru, 2017; Sugihara et al.
123 2020), *D. dumetorum* (Siadjeu et al. 2020) and *D. alata* (Bredeson et al. 2022). Genome sizes of these
124 species were 570 Mb (*D. rotundata*), 485 Mb (*D. dumetorum*) and 480 Mb (*D. alata*). The genome of
125 *D. tokoro* is slightly smaller than the genomes of *Dioscorea* species so far reported.

126 Basic chromosome number of *Dioscorea* is suggested to be ten. *D. tokoro* revealed to have ten
127 linkage groups in this study, which is in line with the report of linkage group obtained by AFLP
128 analysis (Terauchi and Kahl, 1999). It is contrasting to *D. rotundata* ($2n = 2x = 40$) and *D. alata* ($2n
129 = 2x = 40$), both belonging to the section Enantiophyllum. It is likely that during the evolution of
130 *Dioscorea*, chromosome duplication occurred. However, no signature of genome duplication observed
131 in *D. rotundata* genome (Tamiru et al. 2017), suggesting that genome duplication happened in a
132 remote past.

133 Sex determination of *D. tokoro* was confirmed to be XY type male heterogametic system. It is
134 similar to the XY system in *D. alata* (Cormier et al. 2019), but in contrast to ZW female heterogametic
135 system in *D. rotundata* (Tamiru et al. 2017). It is likely that sex determination locus has shifted
136 multiple times in the genus. Future study will identify the genes involved in *D. tokoro* sex
137 determination.

138 In summary, we determined a draft genome sequence of a wild yam species *D. tokoro*. This
139 chromosome level sequence information will serve as a platform to understand population genomics
140 of this obligate outcrossing species and to elucidate the mechanism and evolution of its sex
141 determination system.

142

**Materials and Methods**

**Plant Materials**

A female plant, Waka1 (original code: DT49) (**Fig. S2A**; **Fig. S2D**), was collected from Tahara, Wakayama Pref. in Central Japan. A male plant, Kita1 (original code: 110628-5), was collected from Waga-Sennin of Iwate Pref. in Northern Japan (**Fig. S2B**; **Fig. S2D**). To construct a linkage map, we obtained F1 seeds derived from a cross between Waka1 (P1) and Kita1 (P2) in 2011. We started growing 206 F1 individuals in 2012 and obtained sex data for 186 F1 individuals in 2014 and 2015 (**Fig. S2C**).

**Flow cytometry**

Flow cytometry (FCM) analysis was carried out using nuclei prepared from fresh leaf samples. Nuclei were isolated and stained with propidium iodide (PI) and analyzed using a Cell Lab QuantaTM SC Flow Cytometer (Beckman Coulter, USA) following the manufacturer's protocol.

**Whole genome sequencing and *de novo* assembly**

To generate *Dioscorea tokoro* reference genome sequence, we sequenced the male plant Kita1 using the PromethION sequencer (Oxford Nanopore Technologies). First, Kita1 DNA was extracted from fresh leaves as described in our previous report (Tamiru et al., 2017). The extracted DNA was subjected to size selection and purification with a gel extraction kit (Large Fragment DNA Recovery Kit; Zymo Research). Finally, the purified DNA was sequenced by PromethION at GeneBay company, Yokohama, Japan (http://genebay.co.jp). As the first step for genome assembly, we removed the lambda phage genome from Nanopore fastq using NanoLyse 1.1.0 (De Coster et al. 2018) and filtered out reads with an average read quality score less than seven and those shorter than 1,000 bases with Nanofilt v2.2 (De Coster et al. 2018). We also performed two types of Illumina sequencing, 251 bp paired-end sequencing using MiSeq and 150 bp paired-end sequencing using HiSeq4000. Next, we assembled the filtered long DNA sequence reads with the hybrid assembler MaSuRCA v3.3.4 (Zimin et al. 2013), run internally by Flye assembler v2.6 (Kolmogorov et al. 2019).

**Assessing genome completeness**

To evaluate the completeness of the gene set in the assembled genome, we applied BUSCO analysis (Bench-Marking Universal Single Copy) v5.1.2 (Mosè et al. 2020). We used the default gene search

179  method 'metaeuk gene search' instead of the traditional gene search method using AUGUSTUS (Hoff
180  and Stanke, 2013) and TBLASTN (Camacho et al. 2009). We set "genome" as the assessment mode
181  and used embryophyte_odb10 as the lineage datasets.

182

183  **Gene prediction and annotation**

184

185  For gene prediction, we used RNA-seq data from 18 samples of *D. tokoro*, representing seven organs
186  of Kita1 individual (leaves, stems, root apex, rhizome bud, rhizome root, rhizome stem, and rhizome
187  storage) and 11 different flowering stages in female and male *D. tokoro* plants from the wild (**Table**
188  **S6**). First, according to the manufacturer's instructions, total RNAs were used to construct cDNA
189  libraries using a TruSeq RNA Sample Prep Kit V2 (Illumina, USA). Then, the bulked cDNA library
190  was sequenced using the Illumina NextSeq500 platform for 75 bp single-end reads. In the fastq quality
191  control step, we first remove adapters, poly(A), and the reads shorter than 50 bp using FaQCs (Lo and
192  Chain, 2014). Subsequently, we removed low-quality bases from the read end (window size = 5, base
193  quality average = 20) and low-quality reads with an average read quality below 20 using PRINSEQ
194  lite 0.20.4 (Schmieder and Edwards, 2011). Quality trimmed reads were aligned to the assembled
195  genome with HISAT2 v2.1 (Kim et al. 2019) with the options "--max-intronlen 15000 --dta". Next,
196  transcript alignments were assembled with StringTie v1.3.6 (Pertea et al. 2015) separately for each
197  BAM file. Finally, these GFF files were integrated with TACO v0.7.3 (Niknafs et al. 2017) with the
198  option "--filter-min-length 90", generating 24,148 gene models within the assembled genome (**Table**
199  **2**). Additionally, 34,539 peptide sequences that were predicted in *D. rotundata* genome (Tamiru et al.,
200  2017) [ENSEMBL (http://plants.ensembl.org/Dioscorea_rotundata/Info/Index).] were aligned to
201  assembled genome with Spaln2 v2.3.3 (Iwata and Gotoh, 2012). Consequently, 1,900 CDSs that did
202  not overlap with the new gene models were added to the new gene models (**Table 2**). In addition, the
203  3,036 transcripts that were assembled in the StringTie program but rejected in the TACO program
204  were added manually. Finally, gene models shorter than 75 bases were removed, and InterProScan
205  v5.36 (N) was used to predict ORFs (open reading frames) and strand information for each gene model.
206  We predicted 29,084 genes, including 54,847 transcript variants (**Table 2**). For gene annotation, the
207  predicted gene models were searched in the Pfam protein family database using InterProScan (Blum
208  et al. 2021) and with the blastx command in BLAST+ (Camacho et al. 2009) with the option "-evalue
209  1e-10", using the Viridiplantae database from UniProt as the target database. The resulting gene
210  models and annotations were uploaded to https://genome-e.ibrc.or.jp/resource/dioscorea-tokoro/.

211

212  **Identification of parental line-specific heterozygous markers**

213

214  *RAD sequencing*

6

215 We performed RAD-seq to develop the linkage map as previously described (Tamiru et al., 2017).
216 Genomic DNA was extracted from fresh leaves of Waka1, Kita1, and 186 F1 individuals and digested
217 with the restriction enzymes PacI and NlaIII to prepare libraries used to generate 75-bp paired-end
218 reads by Illumina NextSeq500. We remove adapters and the unpaired reads using FaQCs and
219 PRINSEQ lite as previously described. The filtered RAD-seq reads were used as RAD-tags (**Fig. S3**).
220

221 *SNP-type heterozygous markers*
222 RAD-tags were aligned to the assembled genome of *D. tokoro* in this study using BWA (ver. 0.7.12).
223 SNP-based genotypes for P1, P2, and F1 individuals was obtained as a variant call format (VCF) file.
224 The VCF file was generated from BAM files of P1, P2, and F1 individuals using SAMtools (ver 1.5),
225 and the VCF variants were called and filtered using BCFtools (ver 1.5). As a result, 5,894 P1- or P2-
226 heterozygous SNP markers were selected (shown as "All RAD markers" in **Table S4**). Next, to
227 increase the accuracy of the selected markers, their segregation (1:1 ratio) was confirmed in F1
228 individuals obtained from a cross between P1 and P2. If the segregation ratio was out of the confidence
229 interval ($P < 0.001$) hypothesized by the binomial distribution, $B$ (n = number of individuals, $P = 0.5$),
230 the markers were excluded from further analysis. Finally, 3,057 P1-heterozygous SNP markers and
231 1,559 P2-heterozygous SNP markers were selected (shown as "Confirmed segregation ratio" in **Table**
232 **S4**). Additional details are provided in the **Supplementary Method**.
233

234 *Presence/absence-type heterozygous markers*
235 The presence/absence-type markers were defined based on the alignment depth of parental line RAD-
236 tags. The presence/absence-type markers were called by the following method: First, the VCF file was
237 generated from BAM files of P1 and P2 and selected the region where either P1 or P2 had sufficient
238 read depth ($\geq 8$) and that the other parental line had no read depth in that region. Next, BEDtools (ver
239 2.26) converted continuous positions in the VCF file to a feature, and only sufficiently wide features
240 (width $\geq 50$ bp) were retained as the BED file. For these regions in the BED file, the F1 individual's
241 genotypes were classified into three categories (depth $\geq 3$, depth $= 0$, others) and three genotypes
242 ("presence," "absence," "NA"). As a result, 5,071 PA markers were selected (shown as "All RAD
243 markers" in **Table S4**). We then applied the same binomial test as for the SNP-type heterozygous
244 markers. Finally, 480 P1-heterozygous PA markers and 1,682 P2-heterozygous PA markers were
245 selected (shown as "Confirmed segregation ratio" in **Table S4**). Additional details are provided in the
246 **Supplementary Method**.
247

248 *Integration of SNP-type and presence/absence-type heterozygous markers.*
249

7

250    We integrated SNP-type and PA-type heterozygous markers to develop parental line-specific linkage

251    maps. Two types of markers were defined: P1-heterozygous markers and P2-heterozygous markers.

252    If an SNP-type marker was heterozygous in P1 but homozygous in P2 or if a PA-type marker was

253    present in P1 and absent in P2, it was classified as a P1-heterozygous marker set. Conversely, if a

254    SNP-type marker was homozygous and heterozygous in P1 and P2, respectively, or if a PA-type

255    marker was absent in P1 but present in P2, it was classified as a P2-heterozygous marker set.

256

257    **Linkage maps construction**
258

259    *Pruning and flanking markers by Spearman's correlation coefficients*

260    Pairwise matrix of Spearman's correlation coefficients ($\rho$) were calculated for every maker pair in

261    each contig in each marker set (P1-heterozygous marker set and P2-heterozygous marker set).

262    According to the histogram of absolute $\rho$ calculated from each contig, most markers on the same

263    contigs were correlated with each other. Therefore, we pruned correlated flanking markers to remove

264    redundant markers. Finally, we obtained 2,818 markers for linkage mapping (shown as "Pruning and

265    flanking" in **Table S4**).

266

267    *Linkage mapping*

268    We converted the flanking markers obtained as described in the previous section into the genotype-

269    formatted data for constructing genetic linkage maps using MSTmap (Wu, 2008) with following

270    parameter sets: "populationtype DH; distancefunction kosambi; cutoffpvalue 0.000000000001;

271    nomapdist 15.0; nomapsize 0; missingthreshold 25.0; estimationbeforeclustering no; detectbaddata

272    no; objective_function ML" for P1-heterozygous marker set and P2-heterozygous marker set. After

273    trimming the orphan linkage groups, we solved the complemented-phased duplex linkage groups

274    caused by coupling-type and repulsion-type markers in the pseudo-testcross method. Finally, two

275    parental-specific linkage maps were constructed. These two linkage maps were designated as P1-map

276    (constructed using P1-heterozygous marker set) and P2-map (constructed using P2-heterozygous

277    marker set) (**Fig. S4A; Fig. S4C**). The order and names of each linkage group were organized

278    according to the P2-map (**Fig. S4 and Fig. S5**). The linkage groups were visualized by R/qtl (Broman

279    et al., 2003).

280

281    **Generation of pseudochromosomes**

282

283    Based on a matrix derived from the contigs shared between the P1- and P2-maps, i.e., linkage groups

284    (**Fig. S5**), the contigs were anchored and linearly ordered as pseudochromosomes. First, we identified

285    contigs whose markers were allocated to different linkage groups during the anchoring and ordering

8

286 process. Such contigs were further divided into sub-contigs to ensure that they were not allocated to

287 wrong pseudochromosomes. Next, we divided the contigs at the proper positions as described

288 previously (Tamiru et al. 2017). Finally, we followed the described method (Tamiru et al. 2017) to

289 generate ten pseudochromosomes.

290

291 **Identification of sex associated region**

292

293 To identify the sex-associated genomic region, we performed Fisher's exact test using the genotype

294 of 127 F1 individuals based on 2,730 markers located on the pseudochromosomes (**Table S4**) and

295 their sex phenotype (**Table S3**). Fisher's exact test was performed using the fisherexact function in

296 the python SciPy package. A significance threshold of 5% false discovery rate (FDR) was calculated

297 using the multipletests function in the python statsmodels package with the option method=" fdr_bh"

298 (Benjamini / Hochberg procedure).

299

300

301 **Data Availability**

302

303 All sequencing read data generated for this work have been deposited at the DNA Databank of Japan

304 (DDBJ) database under BioProject PRJDB12945; see **Table S1** and **S2** for individual sample

305 accession numbers. The genomic sequence file (fasta), gene annotation file (gff3), and gene/protein

306 sequences file (fasta) are available at the following URL: https://genome-

307 e.ibrc.or.jp/resource/dioscorea-tokoro/

308

309

310 **Funding**

311

312 This work was supported by Iwate Biotechnology Research Center.

313

314

315 **Disclosures**

316

317 The authors have no conflicts of interest to declare.

318

319

320 **Acknowledgments**

321

325

326

**References**

328

329    Beck, H. E., Zimmermann, N. E., McVicar, T. R., Vergopolan, N., Berg, A., and Wood, E. F. (2018)

330        Present and future Köppen-Geiger climate classification maps at 1-km resolution. Scientific

331        data, 5.1: 1-12.

332

333    Blum, M., Chang, H. Y., Chuguransky, S., Grego, T., Kandasaamy, S., Mitchell, A., Gift N., et al.

334        (2021) The InterPro protein families and domains database: 20 years on. Nucleic acids

335        research, 49.D1: D344-D354.

336

337    Bredeson, J. V., Lyons, J. B., Oniyinde, I. O., Okereke, N. R., Kolade, O., Nnabue, I., Nwadili, C. O.,

338        et al. (2022) Chromosome evolution and the genetic basis of agronomically important traits

339        in greater yam. Nature communications, 13.1: 1-16.

340

341    Broman, K. W., Wu, H., Sen, Ś., and Churchill, G. A. (2003) R/qtl: QTL mapping in experimental

342        crosses. bioinformatics, 19.7: 889-890.

343

344    Burkill, I. H. (1960) The organography and the evolution of Dioscoreaceae, the family of the yams.

345        Journal of the Linnean Society, 56: 319-412.

346

347    Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T. L.

348        (2009) BLAST+: architecture and applications. BMC bioinformatics, 10.1: 1-9.

349

350    Chaïr, H., Cornet, D., Deu, M., Baco, M. N., Agbangla, A., Duval, M. F., and Noyer, J. L. (2010)

351        Impact of farmer selection on yam genetic diversity. Conservation Genetics. 11.6: 2255-2265.

352

353    Chaïr, H., Sardos, J., Supply, A., Mournet, P., Malapa, R., and Lebot, V. (2016) Plastid phylogenetics

354        of Oceania yams (*Dioscorea* spp., Dioscoreaceae) reveals natural interspecific hybridization

355        of the greater yam (*D. alata*). Botanical Journal of the Linnean Society. 180.3: 319-333.

356

357  Cormier, F., Lawac, F., Maledon, E., Gravillon, M.-C., Nudol, E., Mournet, P., Vignes H, et al. (2019)
358      A reference high-density genetic map of greater yam (*Dioscorea alata* L.). Theor Appl Genet.
359      132: 1733–1744.
360

361  Danecek, P., Bonfield, J. K., Liddle, J., Marshall, J., Ohan, V., Pollard, M. O., Andrew W., et al.
362      (2021) Twelve years of SAMtools and BCFtools. Gigascience, 10.2: giab008.
363

364  De Coster, W., D'Hert, S., Schultz, D. T., Cruts, M., and Van Broeckhoven, C. (2018) NanoPack:
365      visualizing and processing long-read sequencing data. Bioinformatics, 34.15: 2666-2669.
366

367  FAOSTAT (2018) Food and Agriculture Organization. http://www.fao.org/statistics.
368

369  Girma, G., Hyma, K. E., Asiedu, R., Mitchell, S. E., Gedil, M., and Spillane, C. (2014) Next-
370      generation sequencing based genotyping, cytometry and phenotyping for understanding
371      diversity and evolution of guinea yams. Theoretical and Applied Genetics. 127.8: 1783-1794.
372

373  Grattapaglia, D., and Sederoff, R. (1994) Genetic linkage maps of *Eucalyptus grandis* and *Eucalyptus*
374      *urophylla* using a pseudo-testcross: mapping strategy and RAPD markers. Genetics, 137.4:
375      1121-1137.
376

377  Hoff, K. J., and Stanke, M. (2013). WebAUGUSTUS—a web service for training AUGUSTUS and
378      predicting genes in eukaryotes. Nucleic acids research, 41.W1: W123-W128.
379

380  Iwata, H., and Gotoh, O. (2012) Benchmarking spliced alignment programs including Spaln2, an
381      extended version of Spaln that incorporates additional species-specific features. Nucleic acids
382      research, 40.20: e161-e161.
383

384  Kim, D., Paggi, J. M., Park, C., Bennett, C., and Salzberg, S. L. (2019) Graph-based genome alignment
385      and genotyping with HISAT2 and HISAT-genotype. Nature biotechnology, 37.8: 907-915.
386

387  Kolmogorov, M., Yuan, J., Lin, Y., and Pevzner, P. A. (2019) Assembly of long, error-prone reads
388      using repeat graphs. Nature biotechnology, 37.5: 540-546.
389

390  Li, H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv
391      preprint arXiv:1303.3997.
392

393    Lo, C. C., and Chain, P. S. G. (2014) Rapid evaluation and quality control of next generation
394            sequencing data with FaQCs. BMC Bioinformatics. 15: 366.

395

396    Mosè M., Matthew R. B., Mathieu S., Felipe A. S., and Evgeny M. Z., (2021) BUSCO Update: Novel
397            and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for
398            Scoring of Eukaryotic, Prokaryotic, and Viral Genomes, Molecular Biology and Evolution,
399            Volume 38, Issue 10, Pages 4647-4654.

400

401    Niknafs, Y. S., Pandian, B., Iyer, H. K., Chinnaiyan, A. M., and Iyer, M. K. (2017) TACO produces
402            robust multisample transcriptome assemblies from RNA-seq. Nature methods, 14.1: 68-70.

403

404    Okagami, N. and Kawai, M. (1982) Dormancy in *Dioscorea*: Differences of temperature responses in
405            seed germination among six Japanese species. The botanical magazine= Shokubutsu-gaku-
406            zasshi. Tokyo. 95.2: 155-166.

407

408    Oyama, M., Tokiwano, T., Kawaii, S., Yoshida, Y., Mizuno, K., Oh, K., et al. (2017) Protodioscin,
409            Isolated from the Rhizome of *Dioscorea tokoro* Collected in Northern Japan is the Major
410            Antiproliferative Compound to HL-60 Leukemic Cells. Current bioactive compounds. 13.2:
411            170-174.

412

413    Pertea, M., Pertea, G. M., Antonescu, C. M., Chang, T. C., Mendell, J. T., and Salzberg, S. L. (2015)
414            StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. Nature
415            biotechnology, 33.3: 290-295.

416

417    Quinlan, A. R., and Hall, I. M. (2010) BEDTools: a flexible suite of utilities for comparing genomic
418            features. Bioinformatics, 26.6: 841-842.

419

420    Scarcelli, N., Tostain, S., Vigouroux, Y., Agbangla, C., Daïnou, O., and Pham, J. L. (2006) Farmers'
421            use of wild relative and sexual reproduction in a vegetatively propagated crop. The case of
422            yam in Benin. Molecular Ecology. 15.9: 2421-2431.

423

424    Scarcelli, N., Chaïr, H., Causse, S., Vesta, R., Couvreur, T. L. P., and Vigouroux, Y. (2017) Crop wild
425            relative conservation: Wild yams are not that wild. Biological Conservation. 210: 325-333.

426

427   Scarcelli, N., Cubry, P., Akakpo, R., Thuillet, A. C., Obidiegwu, J., Baco, M. N., Otoo, E., et al. (2019)
428         Yam genomics supports West Africa as a major cradle of crop domestication. Science
429         advances. 5.5: eaaw1947.
430

431   Schmieder, R., and Edwards, R. (2011) Quality control and preprocessing of metagenomic datasets.
432         Bioinformatics. 27.6: 863-864.
433

434   Siadjeu, C., Mayland-Quellhorst, E., and Albach, D. C. (2018) Genetic diversity and population
435         structure of trifoliate yam (*Dioscorea dumetorum* Kunth) in Cameroon revealed by
436         genotyping-by-sequencing (GBS). BMC Plant Biology. 18.1: 1-14.
437

438   Siadjeu, C., Pucker, B., Viehöver, P., Albach, D.C., and Weisshaar, B. (2020) High Contiguity de
439         novo Genome Sequence Assembly of Trifoliate Yam (*Dioscorea dumetorum*) Using Long
440         Read Sequencing. Genes. 11: 274.
441

442   Sugihara, Y., Darkwa, K., Yaegashi, H., Natsume, S., Shimizu, M., Abe, A., Hirabuchi, A., et al.
443         (2020) Genome analyses reveal the hybrid origin of the staple crop white Guinea yam
444         (*Dioscorea rotundata*). Proceedings of the National Academy of Sciences. 117.50: 31987-
445         31992.
446

447   Sugihara, Y., Kudoh, A., Oli, M. T., Takagi, H., Natsume, S., Shimizu, M., Abe, A., et al. (2021)
448         Population Genomics of Yams: Evolution and Domestication of *Dioscorea* Species. In:
449         Population Genomics. pp. 1-28. Springer, Cham.
450

451   Tamiru, M., Natsume, S., Takagi, H., White, B., Yaegashi, H., Shimizu, M., et al. (2017) Genome
452         sequencing of the staple food crop white Guinea yam enables the development of a molecular
453         marker for sex determination. BMC biology. 15.1: 1-20.
454

455   Terauchi, R. (1990) Genetic diversity and population structure of *Dioscorea tokoro* Makino, a
456         dioecious climber. Plant Species Biology, 5.2: 243-253.
457

458   Terauchi, R., Chikaleke, V. A., Thottappilly, G., and Hahn, S. K. (1992) Origin and phylogeny of
459         Guinea yams as revealed by RFLP analysis of chloroplast DNA and nuclear ribosomal DNA.
460         Theoretical and Applied Genetics. 83.6: 743-751.
461

13

462    Terauchi, R. and Konuma, A. (1994) Microsatellite polymorphism in *Dioscorea tokoro*, a wild yam
463        species. Genome, 37: 794-801.

464

465    Terauchi, R., Terachi, T., and Miyashita, N. T. (1997) DNA polymorphism at the Pgi locus of a wild
466        yam, *Dioscorea tokoro*. Genetics, 147.4: 1899-1914.

467

468    Terauchi, R. and Kahl, G. (1999) Mapping of the *Dioscorea tokoro* genome: AFLP markers linked to
469        sex. Genome 42: 752-762.

470

471    Terauchi, R. and Kahl, G. (2004) Sex determination in *Dioscorea tokoro*, a wild yam species. In Sex
472        determination in plants (pp. 165-174). Garland Science.

473

474    Wessel, P., Luis, J. F., Uieda, L., Scharroo, R., Wobbe, F., Smith, W. H. F., and Tian, D. (2019) The
475        Generic Mapping Tools version 6. Geochemistry, Geophysics, Geosystems, 20, 5556-5564.
476        https://doi.org/10.1029/2019GC008515

477

478    Wu, Y., Bhat, P. R., Close, T. J., and Lonardi, S. (2008) Efficient and accurate construction of genetic
479        linkage maps from the minimum spanning tree of a graph. PLoS genetics, 4.10: e1000212.

480

481    Zimin, A. V., Marçais, G., Puiu, D., Roberts, M., Salzberg, S. L., and Yorke, J. A. (2013) The
482        MaSuRCA genome assembler. Bioinformatics, 29.21: 2669-2677.

483

484

485


486    **Figures Legends**

487

488    **Fig. 1 Botanical characteristics of *Dioscorea tokoro*.** (A) *D. tokoro* is a herbaceous climber species.

489    Aerial stems twine around tree trunks. (B) Stem twines in an anti-clockwise direction (left-handed;

490    sinistrorse). Leaves alternate. (C) Leaf shape is usually heart-shaped. Leaf blades are typically 5-12

491    cm long and 5-12 cm wide. (D) Female pendulous inflorescences. (E) Close-up view of a female

492    flower. Three-locular ovary are below the petal. Three-lobed pistil and six degenerated stamens around

493    the pistil are seen. Petal apex is round and curled inward. (F) Male upright inflorescences. (G) Close-

494    up view of a male flower. Pedicel branches from the base and has a few flowers. Six stamens, and

495    degenerated pistil in the center. Petal apex is round and curled inward. The scale is same as (E). (H)

496    Female inflorescence with immature obovate-elliptic capsules. Capsules reflex and dehisce at maturity.

497    (I) Male inflorescence. The scale is same as (H). (J) Mature fruit has three capsules, with winged two

498    seeds placed alternately overlapped near its base. Seed's wing is biased wider toward capsule apex.

499    (K) Underground rhizome of *D. tokoto*. (L) A side view of the rhizome from a different angle. The

500    direction of the white arrow corresponds to (K).

501

502    **Fig. 2 An integrated linkage and physical map of *D. tokoro*.** Approximately 85.4% of the *D. tokoro*

503    contig sequences were anchored using a RAD-based genetic map generated with 186 F1 individuals

504    obtained from a cross between Waka1 (P1: female) and Kita1 (P2: male). The 10 pseudochromosomes

505    are numbered from chrom_01 to chrom_10. Markers are located according to genetic distance (cM).

506    The black frame in the center of each group represents the reconstructed pseudochromosome and

507    orange and green bars indicate P1-map and P2-map, respectively. Thin grey lines connecting linkage

508    map and pseudochromosome indicate the positions of markers. The blue dots indicate the positions of

509    PA markers.

510

511    **Fig. 3 Genome-wide association mapping of sex in the F1 progeny derived from a cross between**

512    **Waka1 (P1: female) and Kita1 (P2: male) in *D. tokoro*.** Manhattan plot of markers associated with

513    sex phenotype as determined by Fisher's exact test with (A) P1-heterozygous marker set and (B) with

514    P2-heterozygous marker set. Orange and blue dots indicate SNP and presence/absence markers,

515    respectively, showing significant association with sex based on a 5 % false discovery rate ($q < 0.05$).
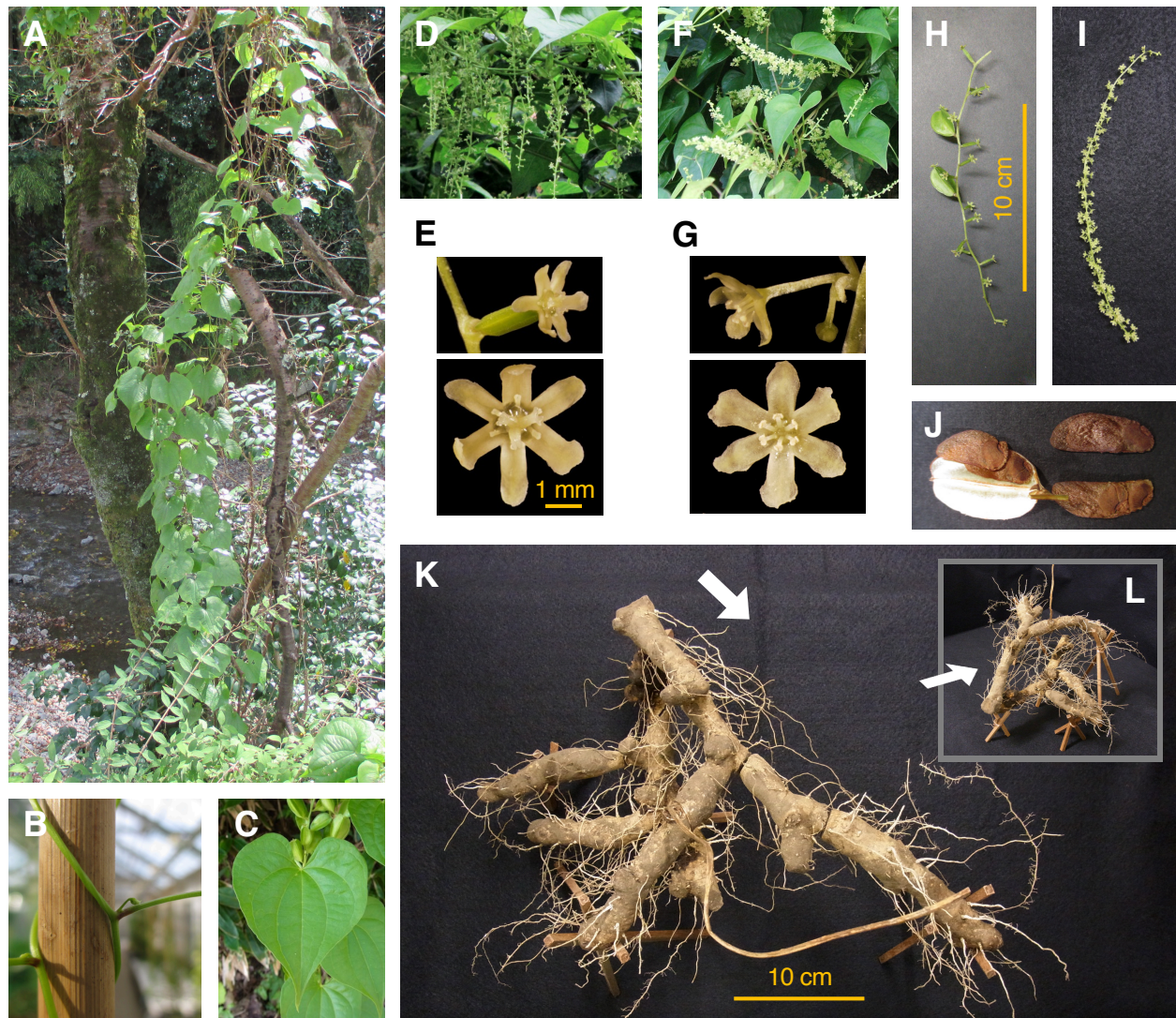
516

517

15

**Fig. 1 Botanical characteristics of *Dioscorea tokoro*.** (A) *D. tokoro* is a herbaceous climber species. Aerial stems twine around tree trunks. (B) Stem twines in an anti-clockwise direction (left-handed; sinistrorse). Leaves alternate. (C) Leaf shape is usually heart-shaped. Leaf blades are typically 5-12 cm long and 5-12 cm wide. (D) Female pendulous inflorescences. (E) Close-up view of a female flower. Three-locular ovary are below the petal. Three-lobed pistil and six degenerated stamens around the pistil are seen. Petal apex is round and curled inward. (F) Male upright inflorescences. (G) Close-up view of a male flower. Pedicel branches from the base and has a few flowers. Six stamens, and degenerated pistil in the center. Petal apex is round and curled inward. The scale is same as (E). (H) Female inflorescence with immature obovate-elliptic capsules. Capsules reflex and dehisce at maturity. (I) Male inflorescence. The scale is same as (H). (J) Mature fruit has three capsules, with winged two seeds placed alternately overlapped near its base. Seed's wing is biased wider toward capsule apex. (K) Underground rhizome of *D. tokoto*. (L) A side view of the rhizome from a different angle. The direction of the white arrow corresponds to (K).

**Table 1.** Summary of a reference genome of *D. tokoro* (Kita1).

| | |
|---|---:|
| Total number of contigs | 2,931 |
| Total base-pairs (bp) | 443,501,147 |
| Estimated genome size from k-mer (bp) | 438,704,233 |
| Average contig size (bp) | 151,313 |
| Longest contig (bp) | 6,172,819 |
| Shortest contig (bp) | 502 |
| N50 (bp) | 586,368 |
| | |
| Total BUSCO groups searched | 1,614 |
| Complete BUSCOs (%) | 98.0 |
| Complete and single-copy BUSCOs (%) | 92.9 |
| Complete and duplicated BUSCOs (%) | 5.1 |
| Fragmented BUSCOs (%) | 1.1 |
| Missing BUSCOs (%) | 0.9 |

**Table 2.** Summary of predicted genes in the *D. tokoro* genome

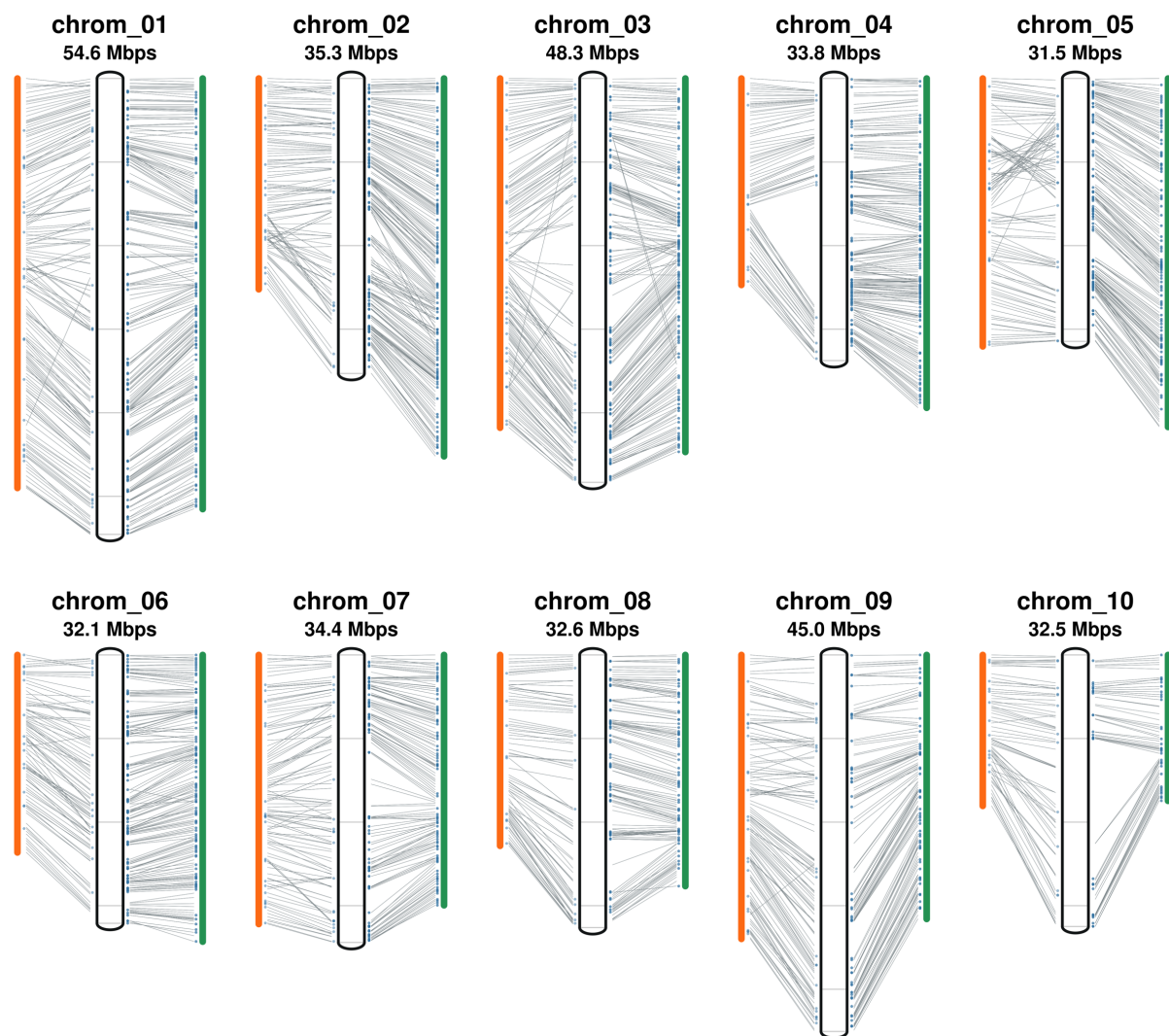| | Contigs (2,931) | Pseuodochrom. (01-10) |
|---|---:|---:|
| No. genes | 29,084 | 25,447 |
| (total transcript variant) | (54,847) | (48,271) |
| ORF status | | |
| Complete | 21,610 | 19,069 |
| 5' partial | 195 | 166 |
| 3' partial | 4,752 | 4,091 |
| Internal | 157 | 127 |
| No ORF | 2,370 | 1,994 |
| Prediction software | | |
| TACO | 24,148 | 21,239 |
| Spaln2 | 1,900 | 1,581 |
| StringTie | 3,036 | 2,627 |

**Fig. 2 An integrated linkage and physical map of *D. tokoro*.** Approximately 85.4% of the *D. tokoro* contig sequences were anchored using a RAD-based genetic map generated with 186 F1 individuals obtained from a cross between Waka1 (P1: female) and Kita1 (P2: male). The 10 pseudochromosomes are numbered from chrom_01 to chrom_10. Markers are located according to genetic distance (cM). The black frame in the center of each group represents the reconstructed pseudochromosome and orange and green bars indicate P1-map and P2-map, respectively. Thin grey lines connecting linkage map and pseudochromosome indicate the positions of markers. The blue dots indicate the positions of PA markers.
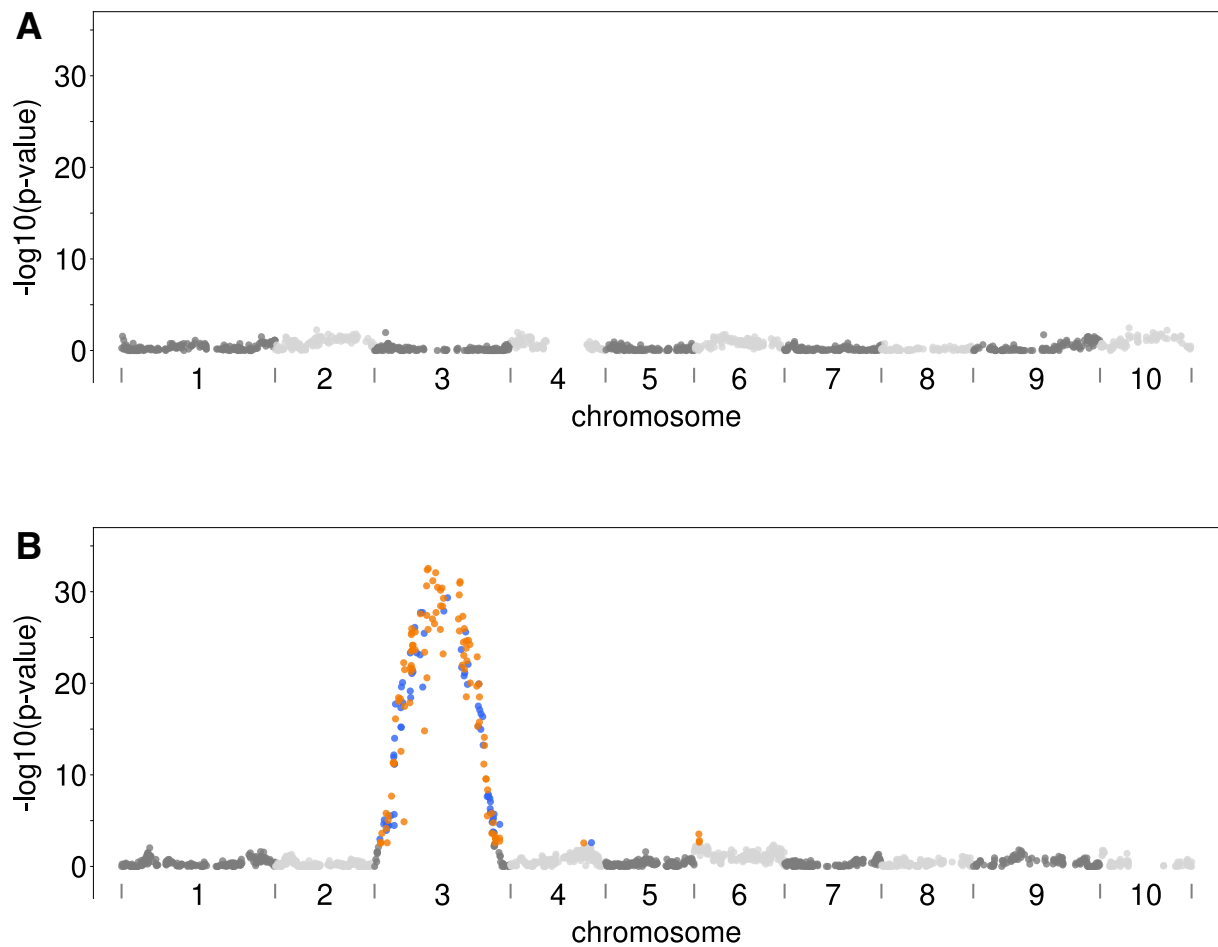
**Fig. 3 Genome-wide association mapping of sex in the F1 progeny derived from a cross between Waka1 (P1: female) and Kita1 (P2: male) in *D. tokoro*.** Manhattan plot of markers associated with sex phenotype as determined by Fisher's exact test with (A) P1-heterozygous marker set and (B) with P2-heterozygous marker set. Orange and blue dots indicate SNP and presence/absence markers, respectively, showing significant association with sex based on a 5 % false discovery rate (q < 0.05).
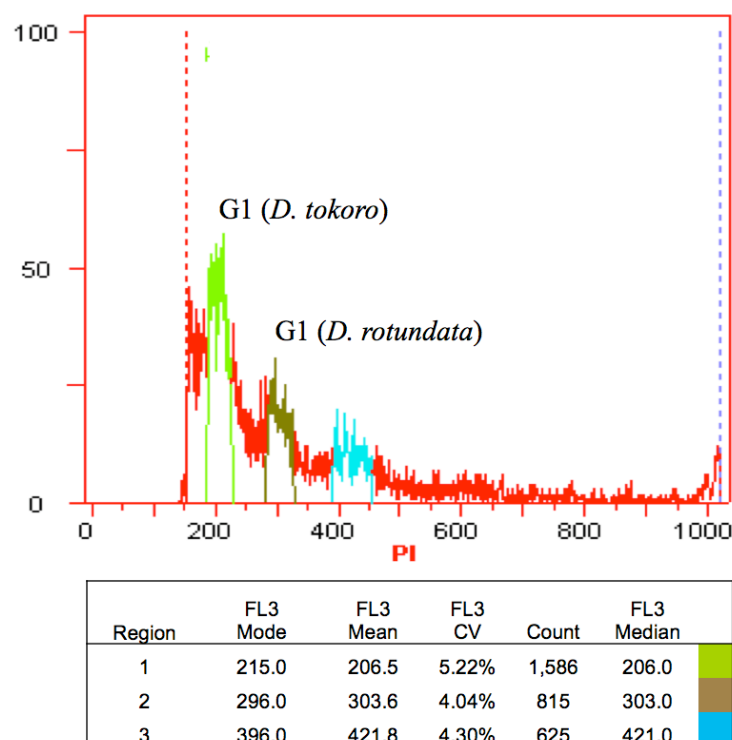
| Region | FL3 Mode | FL3 Mean | FL3 CV | Count | FL3 Median | |
|--------|----------|----------|--------|-------|------------|---|
| 1 | 215.0 | 206.5 | 5.22% | 1,586 | 206.0 | |
| 2 | 296.0 | 303.6 | 4.04% | 815 | 303.0 | |
| 3 | 396.0 | 421.8 | 4.30% | 625 | 421.0 | |

**Fig. S1 Size estimation of *D. tokoro* genome by flow cytometry.** Flow cytometry analysis was carried out using nuclei prepared from fresh leaf samples of a wild plant of *D. tokoro* collected in Kitakami, Iwate, Japan and a plant of *D. rotundata* maintained in a greenhouse at Iwate Biotechnology Research Center (IBRC). *D. rotundata* (570 Mb) was served as an internal reference standard of known genome size (Tamiru et al. 2017). Nuclei were isolated and stained with propidium iodide (PI) and analyzed using a Cell Lab QuantaTM SC Flow Cytometer (Beckman Coulter, USA) following the manufacturer's protocol. The ratio of G1 peak mean [*D. tokoro* (206.5): *D. rotundata* (303.6) = 0.680] was used to estimate the genome size of *D. tokoro* to be 388 Mb (570 Mb $\times$ 0.68).

**Table S1.** Summary of filtered ONT reads.

| | |
|---|---|
| Number of reads | 2,515,235 |
| (before filtering) | (3,126,676) |
| | |
| Total base-pairs (Gbp) | 27.4 |
| (before filtering) | (32.7) |
| | |
| Genome coverage* (x) | 70.6 |
| Mean read length (kbp) | 10.9 |
| Longest fragment (kbp) | 11.0 |
| Shortest fragment (bp) | 1 |
| | |
| Accession No. | DRR344532 |

*Genome coverage is estimated from the expected genome size of *D. tokoro* (388 Mb).

**Table S2.** Summary of non-filtered Illumina short reads.

| Sequence platform | Read length (bp) | Total base pairs(Gbp) | Genome coverage | Accession No. |
|---|---|---|---|---|
| Illumina MiSeq | 35-251 | 24.6 | 63.4x | DRR344531 |
| Illumina HiSeq 4000 | 150 | 37.8 | 97.4x | DRR347075 |
| | Total | 62.4 | 160.8x | |

**Table S3.** Number of male, female and non-flowering progeny derived from a Waka1 (P1) x Kita1 (P2) cross.

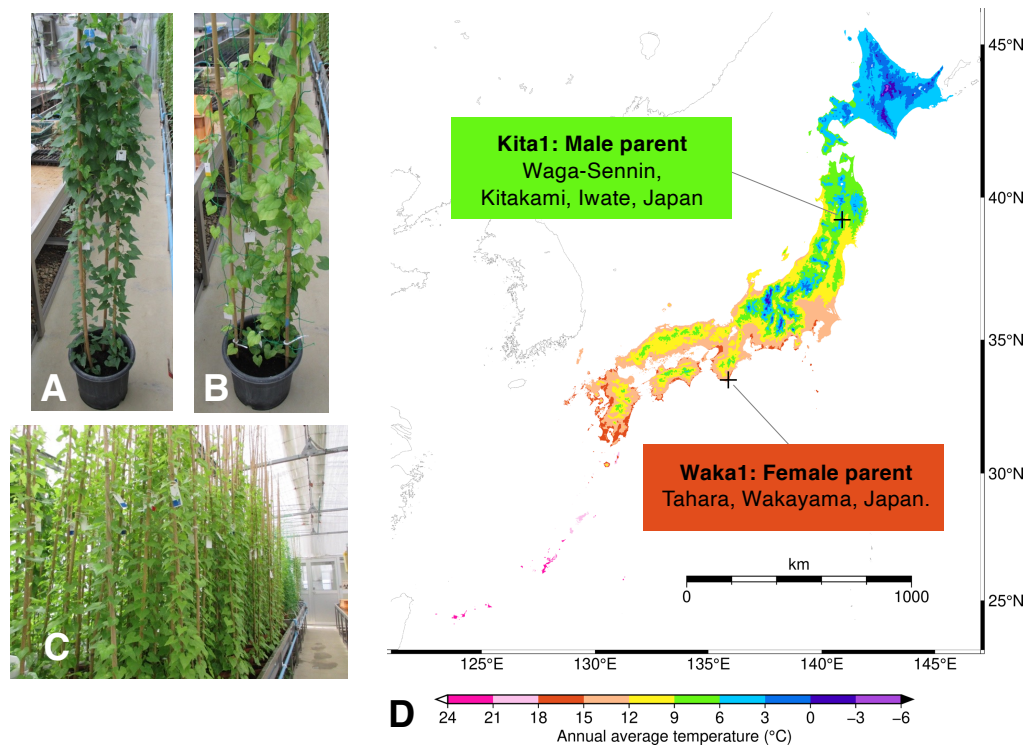| | | Female | Male | Not flowered |
|---|---|---|---|---|
| Parent | Waka1 (P1) | 1 | | |
| | Kita1 (P2) | | 1 | |
| Progeny | 186 individuals | 38 | 89 | 59 |

**Fig. S2  *D. tokoro* plants used for a genetic cross and information of sites of origin.** (A) A female individual Waka1. (B) A male individual Kita1. (C) Sideview of 186 F1 individuals obtained from a cross between Waka1 and Kita1. (D) Waka1 was collected from Tahara, Wakayama Pref. in the Kinki district of Japan. This place is close to the coast. The latitude and longitude are 33.538, 135.860, respectively, and the average annual temperature is 15-18°C. Kita1 was collected from Waga-Sennin in Kitakami, Iwate Pref. in the Tohoku district of Japan. This place is a mountainous. The latitude and longitude are 39.295, 140.896, respectively, and the average annual temperature is 6-9 °C. This map was created with GMT, Version 6.1.1 (Wessel et al. 2019). The annual average temperature on the map of Japan is drawn using "Annual average (climate) mesh data", downloaded from the National Land Numerical Information Download Service (JPGIS2.1) (https://nlftp.mlit.go.jp/ksj/index.html) published by the Ministry of Land, Infrastructure, Transport and Tourism.
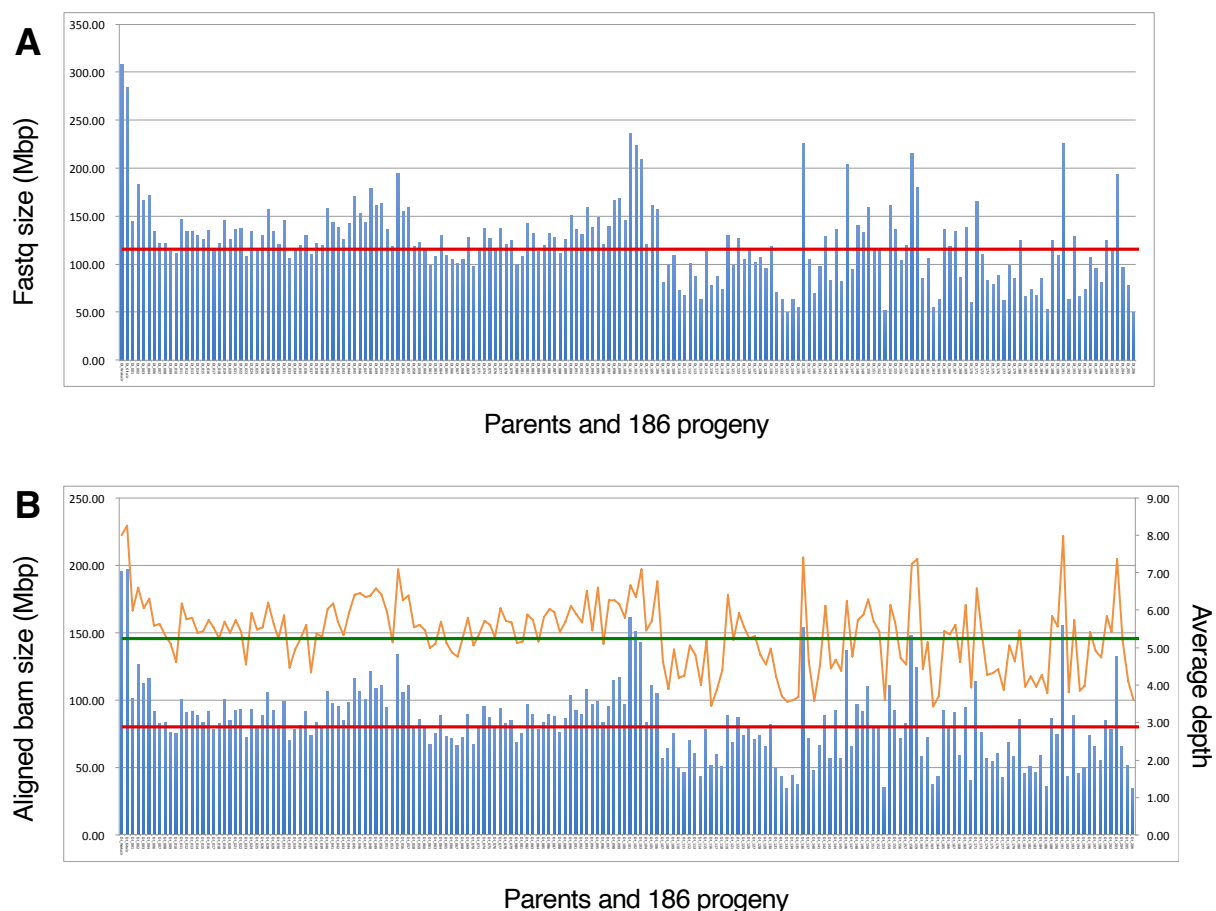
**Fig S3 Summary of RAD tags generated for 186 F1 individuals derived from a cross between Waka1 (P1: female) and Kita1 (P2: male).** In all graphs, the two bars on the left end indicate the parents Waka1 and Kita1, and the other indicate 186 F1 individuals. (A) The total size of filtered fastq of each individual (blue bars). The horizontal red line indicates the average fastq size of the progeny (120.7 Mbp). (B) Aligned bam size (blue bars) and average read depth at genomic regions in the reference genome aligned by the RAD-tags (orange line). The horizontal red line indicates the average aligned bam size of the progeny (82.3 Mbp), and the horizontal green line indicates the average read depth of the progeny (5.34 Mbp).

4

**Table S4.** Number of RAD markers used for anchoring the contigs after filtering.

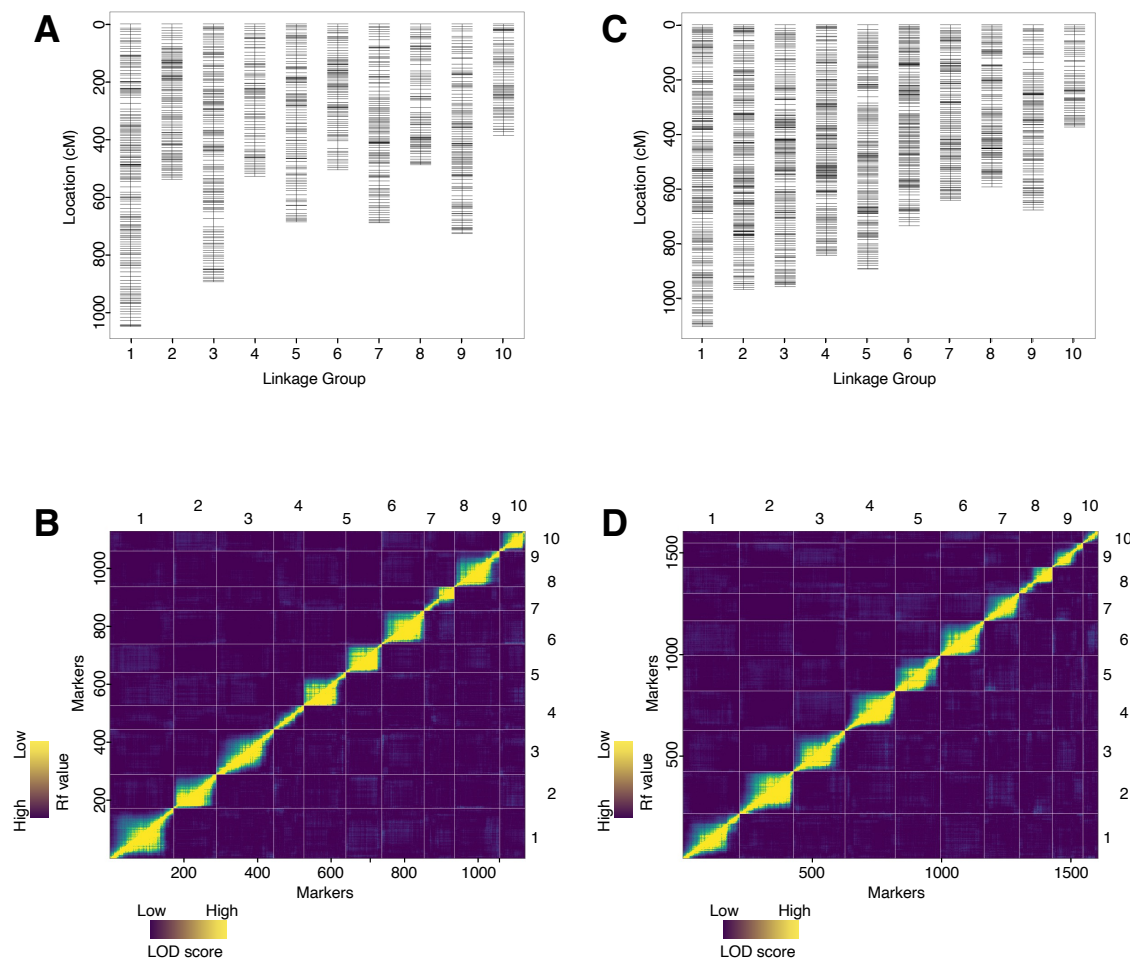| Filtering step | Type | LG | SNP | PA | total |
|---|---|---|---|---|---|
| All RAD markers | testcross | | **5,894** | **5,071** | **10,965** |
| Confirmed segregation ratio | P1-hetero | | 3,057 | 480 | 3,537 |
| | P2-hetero | | 1,559 | 1,682 | 3,241 |
| | **(total)** | | **4,616** | **2,162** | **6,778** |
| Pruning and flanking | P1-hetero | 49 | 988 | 181 | 1,169 |
| | P2-hetero | 55 | 768 | 881 | 1,649 |
| | **(total)** | | **1,756** | **1,062** | **2,818** |
| Anchoring contigs | P1-hetero | 10 | 946 | 180 | 1,126 |
| | P2-hetero | 10 | 724 | 880 | 1,604 |
| | **(total)** | | **1,670** | **1,060** | **2,730** |

**Fig. S4 RAD-seq-based linkage map of *D. tokoro* generated by the pseudo-testcross method using 186 F1 individuals.** (A) P1-map generated using P1-heterozygous marker set. (B) Plots of estimated recombination fractions (upper-left triangle) and LOD score (lower-right triangle) for P1-Map. (C) P2-map generated using P2-heterozygous marker set. (D) Plots of estimated recombination fractions (upper-left triangle) and LOD score (lower-right triangle) for P2-Map. Yellow indicates linked (large LOD score or small recombination fraction) and blue indicates not linked (small LOD score or large recombination fraction).

|  |  | P2-map Linkage group | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| P1-map Linkage group | 1 | 64 | 2 | 2 | 1 | 0 | 0 | 2 | 1 | 1 | 0 |
|  | 2 | 2 | 52 | 1 | 0 | 0 | 2 | 2 | 0 | 0 | 1 |
|  | 3 | 1 | 3 | 38 | 1 | 3 | 1 | 1 | 0 | 0 | 0 |
|  | 4 | 1 | 1 | 0 | 51 | 1 | 0 | 0 | 1 | 0 | 0 |
|  | 5 | 2 | 2 | 0 | 2 | 41 | 1 | 0 | 3 | 0 | 1 |
|  | 6 | 3 | 1 | 2 | 0 | 1 | 38 | 0 | 1 | 0 | 0 |
|  | 7 | 0 | 1 | 1 | 2 | 0 | 0 | 39 | 1 | 1 | 0 |
|  | 8 | 2 | 0 | 1 | 1 | 0 | 1 | 0 | 30 | 0 | 0 |
|  | 9 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 27 | 0 |
|  | 10 | 0 | 3 | 0 | 1 | 1 | 0 | 2 | 0 | 2 | 17 |

**Fig. S5. A matrix of the number of shared contigs between the P1-map and P2-map.**
Both P1-Map and P2-Map contained 10 LGs.

**Table S5.** Summary of contigs anchored by RAD markers of the 10 linkage groups.

|  | P1 hetero | P2 hetero | P1 P2 total |
|---|---|---|---|
| Markers in 10 LGs | 1,126 | 1,604 | 2,730 |
| contig information anchored by marker: |  |  |  |
| number | 566 | 963 | 1,123 |
| contig number % | 19.3 | 32.9 | 38.3 |
| contig total bps | 303,270,609 | 320,556,813 | 378,798,395 |
| contig total % | 68.4 | 72.3 | 85.4 |

*The reference fasta size is 443,501,147 bp.

**Table S6.** Details of 18 RNA-seq samples.

| No. | Organ | Phase | Sex | Stage* | Collected material | Fastq size |
|---|---|---|---|---|---|---|
| 1 | leaves | | male | | Kita1 | 1.71 |
| 2 | stems | | male | | Kita1 | 1.35 |
| 3 | root apex | | male | | Kita1 | 1.66 |
| 4 | rhizome bud | | male | | Kita1 | 1.86 |
| 5 | rhizome root | | male | | Kita1 | 1.61 |
| 6 | rhizome stem | | male | | Kita1 | 1.79 |
| 7 | rhizome storage | | male | | Kita1 | 1.78 |
| 8 | shoot apex | vegetative | unknown | | multiple wild | 1.62 |
| 9 | shoot apex | reproductive | female | 0 | multiple wild | 1.87 |
| 11 | shoot apex | reproductive | female | 1 | multiple wild | 1.70 |
| 13 | shoot apex | reproductive | female | 2 | multiple wild | 1.81 |
| 10 | shoot apex | reproductive | male | 0 | multiple wild | 1.71 |
| 12 | shoot apex | reproductive | male | 1 | multiple wild | 1.74 |
| 14 | shoot apex | reproductive | male | 2 | multiple wild | 1.91 |
| 15 | mature bud | | female | | multiple wild | 1.85 |
| 16 | mature bud | | male | | multiple wild | 1.69 |
| 17 | flower | | female | | multiple wild | 1.83 |
| 18 | flower | | male | | multiple wild | 1.83 |
| | | | | | total | 31.17 |

*The reproductive shoot stage are indicated as follows. 0: shoot apical meristem (SAM), 1: inflorescence with unopened bracts, 2: inflorescences below 10 mm with unseparated bottom buds.

8

**Supplementary Method**

Identification of parental line-specific heterozygous RAD markers

Heterozygous SNP markers

SNP genotypes for P1, P2, and F1 progenies were obtained as a VCF file. The VCF file was generated as follows: (i) SAMtools v1.5 mpileup command with the option "-t DP,AD,SP -B -Q 18 -C 50"; (ii) BCFtools v1.5 call command with the option "-P 0 -v -m -f GQ,GP"; (iii) BCFtools view command with the options "-i 'INFO/MQ≥40, INFO/MQ0F≤0.1, and AVG(GQ)≥10"; and (iv) BCFtools norm command with the option "-m+any." We rejected the variants with low read depth (< 10) or low genotype quality scores (< 10) in the two parents. Likewise, we regarded variants with low read depth (< 8) or low genotype quality scores (< 5) in F1 progenies as missing and only retained the variants with low missing rates (< 0.3). Subsequently, only bi-allelic SNPs were selected by the BCFtools view command with the option "-m 2 -M 2 -v snps". Referring to the genotypes in the VCF file, heterozygous genotypes called by unbalanced allele frequency (out of 0.1-0.9 in F1 progenies) were regarded as missing, and filtering for missing rate (< 0.1) was applied again. As a result, 5,894 P1- or P2-heterozygous SNP markers were selected (shown as "All RAD markers" in Table S5). Next, a binomial test was performed to reject SNPs affected by segregating distortion in the F1 progenies. This binomial test assumes that the probability of success rate is 0.5 based on the two-side hypothesis, and we regarded variants having p-value less than 0.001 as segregation distortion. Finally, 3,057 P1-heterozygous SNP markers and 1,559 P2-heterozygous SNP markers were selected (shown as "Confirmed segregation ratio" in Table S5).

Heterozygous presence/absence RAD markers

The presence/absence markers were defined based on the alignment depth of parental line RAD-tags. A VCF file was generated to search for positions with contrasting read depth between the two parental plants P1 and P2 using the following commands: (i) SAMtools mpileup command with the option "-B -Q 18 -C 50"; (ii) BCFtools call command with the option "-A"; and (iii) BCFtools view command with the options "-i 'MAX(FMT/DP)≥8 and MIN(FMT/DP)≤0' -g miss -V indels". This means that one of the parents (P1 or P2) has enough read depth (≥ 8) and another parent has no reads aligned on that region. Subsequently, we converted continuous positions in the VCF file to a feature that provides a region's start and end coordinate information using the BEDTools v.2.26 merge command with the option "-d 10 -c 1 -o count". We only retained sufficiently wide features (≥ 50 bp) in the BED file. Using the depth value in each feature given in

the BED file, presence/absence (PA) -based genotypes for parental plants P1 and P2 and F1 progenies were determined. For P1 and P2, we regarded genotypes having depth $\geq 4$ as present genotypes, meaning the heterozygosity of presence and absence, while those having depth $= 0$ were classified as absent genotypes, meaning the homozygosity of absence. For F1 progenies, we classified markers with depth $> 2$ and $= 0$ as present and absent markers, respectively. Referring to the genotypes, heterozygous genotypes called by unbalanced allele frequency (out of 0.1-0.9 in F1 progenies) were rejected. As a result, 5,071 PA markers were selected (shown as "All RAD markers" in Table S5). Next, we applied the same binomial test for PA heterozygous markers as SNP-type heterozygous markers. Finally, 480 P1-heterozygous PA markers and 1,682 P2-heterozygous PA markers were selected (shown as "Confirmed segregation ratio" in Table S5).