1    **Genome sequence analysis of a giant-rooted 'Sakurajima daikon' radish (*Raphanus sativus*)**

2    Running title: Genome of 'Sakurajima daikon' radish

3

4    Kenta Shirasawa[1*], Hideki Hirakawa[1], Nobuko Fukino[2†], Hiroyasu Kitashiba[3], and Sachiko Isobe[1]

5    [1]Kazusa DNA Research Institute, Kisarazu, Chiba 292-0818, Japan, [2]Institute of Vegetable and

6    Floriculture Science, NARO, Tsu, Mie 514-2392, Japan, and [3]Graduate School of Agricultural

7    Science, Tohoku University, Sendai, Miyagi 980-0845, Japan

8    †Current address: Strategic Planning Headquarters, NARO, Tsukuba, Ibaraki 305-8517, Japan

9    *Correspondence: shirasaw@kazusa.or.jp

10

11    **Abstract**

12    Daikon radish (*Raphanus sativus*) roots vary in size and shape between cultivars. This study reports

13    the genome sequence assembly of a giant-rooted 'Sakurajima daikon' radish variety,

14    'Okute-Sakurajima', which produces extremely large round roots. Radish genome assembly is

15    hampered by the repetitive and complex nature of the genome. To address this, single-molecule

16    real-time technology was used to obtain long-read sequences at 60× genome coverage. *De novo*

17    assembly of the long reads generated 504.5 Mb contig sequences consisting of 1,437 sequences with

18    contig N50 length of 1.2 Mb, including 94.1% of the core eukaryotic genes. Nine pseudomolecule

19    sequences, comprising 69.3% of the assembled contig length, were generated with high-density SNP

20    genetic maps. The chromosome-level sequences revealed structure variations and rearrangements

21    among Brassicaceae genomes. In total, 89,915 genes were predicted in the 'Okute-Sakurajima'

22    genome, 30,033 of which were unique to the assembly in this study. The improved genome

23    information generated in this study will not only form a new baseline resource for radish genomics,

24    but will also provide insights into the molecular mechanisms underlying formation of giant radish

25    roots.

26    **Footnote**

27    References for data analysis tools used in this study, which are indicated with back quotes in the text,

28    are listed in Supplementary Table S1.

29

30    **1.   Introduction**

31    Daikon radish (*Raphanus sativus*) is a member of the Brassicaceae family of flowering plants.

32  Daikon radish roots of different cultivars vary considerably in their size and shape.[1] For example,

33  among daikon cultivars, 'Sakurajima daikon' radishes exhibit the largest roots. 'Sakurajima daikon'

34  is the name of a group of daikon radish varieties that are mainly cultivated in the Kagoshima

35  prefecture of Japan. Soil in this region contains volcanic ash from Mt. Sakurajima and is thought to

36  be particularly suitable for cultivation of 'Sakurajima daikon'. One well-known 'Sakurajima daikon'

37  line is 'Okute-Sakurajima', which has large round roots that can exceed 20 kg.[1] The size and shape

38  of 'Okute-Sakurajima' roots are desirable traits for plant breeding, but the molecular mechanisms

39  underpinning these characteristics remain unknown.

40      Four genome assemblies based on next-generation sequencing technologies have been reported

41  for radish,[2-5] but these do not include 'Okute-Sakurajima' or 'Sakurajima daikon'. The sequence

42  contiguities of the genome assemblies are relatively short, and the entire radish genome is not

43  covered.[6] One reason for the short assembly size may be the complexity of the radish genome

44  because *Raphanus*, like other species in the *Brassica* genus, underwent genome triplication

45  sometime after divergence from Arabidopsis.[7] Furthermore, radish, in common with other *Brassica*,

46  is self-incompatible and allogamous, resulting in a highly heterozygous genome.[8]

47      Recent advances in long-read sequence technology have allowed highly heterozygous genomes

48  of several plant species to be successfully sequenced.[9] In this study, the 'Okute-Sakurajima' genome

49  was sequenced using long-read technology. Sequences were aligned to radish chromosomes,

50  establishing pseudomolecules and allowing gaps in previous radish genome assemblies to be

51  resolved. This enhanced radish genome will provide insights into radish evolutionary development

52  and, specifically, into the molecular underpinnings of giant daikon root formation.

53

54  **2.  Materials and methods**

55  **2.1.  'Okute-Sakurajima'** *de novo* **genome assembly**

56  Total DNA was extracted from young leaves of the 'Okute-Sakurajima' daikon radish cultivar

57  (NARO GeneBank accession number: JP27228) using a Genomic-tip kit (Qiagen, Hilden, Germany).

58  Short-read sequence data were obtained using a MIGSEQ-2000 DNA sequencer (also known as a

59  DNBSEQ-G400; MGI Tech, Shenzhen, China) and were used to estimate the size of the

60  'Okute-Sakurajima' genome, with *k*-mer distribution analysis performed using `Jellyfish`. To gain

61  long-read sequence data, an SMRT sequence library was constructed with an SMRTbell Express

62  Template Prep Kit (PacBio, Menlo Park, CA, USA) and sequenced on a PacBio Sequel system

63 (PacBio). The long reads from the Sequel system were assembled, and the haplotypes were phased

64 with `Falcon-unzip`. The assembly was polished twice with `Arrow` and designated as RSAskr_r1.0.

65 Assembly completeness was evaluated with `BUSCO`.

66 **2.2. Construction of map-based pseudomolecule sequences**

67 An F2 mapping population (n=115), termed SNF2, derived from a cross between an inbred line via

68 self-pollination of radish cultivar 'Shogoin Daikon' and a line of *R. sativus* var. *raphanistroides*

69 'Nohara 1', collected from Nohara, Maizuru, Kyoto, Japan, was used to establish genetic maps

70 according to the methods of Shirasawa and Kitashiba.[6] In brief, DNA was extracted from leaves of

71 each line and used for ddRAD-Seq library construction. The library was sequenced on a HiSeq4000

72 sequencer (Illumina, San Diego, CA, USA). Data analysis was also performed as described by

73 Shirasawa and Kitashiba.[6] After trimming low-quality sequences and adapter sequences using

74 `FASTX-Toolkit` and `PRINSEQ`, the remaining high-quality reads were mapped onto the

75 RSAskr_r1.0 assembly, using `Bowtie2` to call SNPs using the mpileup command in `SAMtools`

76 followed by filtering out the low-quality data with `VCFtools`. In parallel, ddRAD-Seq data (DRA

77 accession number: DRA005069) from another F2 population (n=95), namely ASF2,[2] derived from a

78 cross between 'Aokubi *S-h*' and 'Sayatori 26704', was also analyzed as above. SNP data were used

79 for genetic map construction with `Lep-Map3`. On the basis of genetic maps constructed with

80 `ALLMAPS`, the RSAskr_r1.0 sequence assembly was assigned to radish chromosomes to produce

81 pseudomolecule sequences, termed RSAskr_r1.0.pmol. The genome structure of the

82 'Okute-Sakurajima' genome was compared with those of *R. sativus*,[4] *Brassica rapa*,[10] and

83 *Arabidopsis thaliana*[11] using `D-Genies`.

84 **2.3. Gene identification in the genome sequences**

85 The gene models predicted in the RSA_r1.0 assembly[2] were mapped onto the RSAskr_r1.0.pmol

86 pseudomolecule sequences using `Minimap2`. In addition, *ab initio* gene prediction was performed

87 on the RSAskr_r1.0.pmol sequences using `Augustus` as described in Kitashiba et al.[2]

88

89 **3. Results**

90 **3.1. *De novo* assembly of the 'Okute-Sakurajima' radish genome**

91 *K*-mer distribution analysis of the 147.3 Gb short-read data indicated that the 'Okute-Sakurajima'

92 radish genome was highly heterozygous, and that the estimated haploid genome size was 592.4 Mb

93 (Supplementary Figure S1). Subsequent long-read sequencing produced 36.0 Gb data (60.7×

3

94 coverage of the estimated genome size) with 2.3 million reads with N50 length of 29.1 kb. After two

95 rounds of polishing, the long-read assembly consisted of 504.5 Mb primary contigs (including 1,437

96 sequences with N50 length of 1.2 Mb) and 263.5 Mb haplotig sequences (including 2,373 sequences

97 with N50 length of 154.6 kb) (Table 1). A BUSCO analysis of the primary contigs indicated that

98 71.0% and 23.1% of sequences were single-copy complete BUSCOs and duplicated complete

99 BUSCOs, respectively (Table 1), suggesting that most of the gene regions were represented in the

100 primary contigs.

101 **3.2. Construction of pseudomolecule sequences based on genetic maps**

102 In total, 5,872 and 2,830 high-quality SNPs were obtained from the SNF2 and ASF2 mapping

103 populations, respectively, and employed for linkage analysis. The resultant genetic map for SNF2

104 consisted of nine linkage groups with 5,570 SNPs covering 867.2 cM, and the map for ASF2

105 comprised nine linkage groups with 2,680 SNPs covering 895.3 cM (Supplementary Tables S2).

106 Contig sequences of RSAskr_r1.0 were assigned to the radish chromosomes in accordance with the

107 two genetic maps. Nine pseudomolecule sequences, termed RSAskr_r1.0.pmol, spanning 349.8 Mb

108 (69.3%) were established with 293 contigs (Table 2), of which 95 sequences (218.5 Mb, 43.3%)

109 were oriented. The nine resulting sequences were named using the nomenclature (R1–R9) proposed

110 by Shirasawa and Kitashiba[6]. The remaining unassigned sequences (n=1,144, 154.7 Mb, 30.7%)

111 were designated as R0. The structure of the 'Okute-Sakurajima' genome was conserved in *R. sativus*

112 but partially disrupted in *B. rapa* and *A. thaliana*, as indicated previously[2].

113 **3.3. Genes on the 'Okute-Sakurajima' genome sequence**

114 In total, 89,915 gene models were predicted in the RSAskr_r1.0 assembly using an *ab initio* gene

115 prediction method (Table 2). To assess the availability of gene spaces in RSAskr_r1.0,[2] the 80,521

116 predicted gene models were aligned to the RSAskr_r1.0.pmol pseudomolecules. Of the predicted

117 gene models, 78,645 (97.6%) were aligned, suggesting that most of the genes were represented in

118 the RSAskr_r1.0.pmol assembly. The genome positions of 59,882 predicted genes of RSAskr_r1.0

119 and 77,496 mapped genes of RSA_r1.0 overlapped. The remaining 30,033 genes (=89,915-59,882)

120 were unique to the assembled sequences of 'Okute-Sakurajima'.

121

122 **4. Discussion**

123 In this study, we report the genome sequence assembly of 'Okute-Sakurajima' radish, a variety of

124 'Sakurajima-diakon', based on long-read sequence technology. The total assembly size of 504.5 Mb

125  is the longest reported for any radish genome to date,[2-5] suggesting that the long reads might span the

126  repetitive sequences found throughout the radish genome. However, since the assembly size was

127  approximately 90 Mb shorter than the estimated genome size, the long sequencing technology was

128  not sufficient to fully resolve the difficulties in assembling the complex genome of this

129  mesopolyploid species.[7]

130     Map-based pseudomolecule sequences comprising 69.3% of assembled sequences were

131  produced. Unexpectedly, 30.7% of assembled sequence remained unassigned to radish chromosomes.

132  As genetic mapping is reliant upon SNP availability in the genome, sequences cannot be assigned

133  where SNPs are not present. To overcome this genetic limitation, optical mapping and Hi-C

134  technologies, both of which are based on physical mapping strategies, have been developed.[9] These

135  technologies, alongside traditional genetic mapping, would allow the genome coverage of the

136  assembly and the completeness of pseudomolecules to be further improved.

137     In this study, 30,033 genes were discovered that were unique to the current assembly,

138  suggesting that these genes may not have been identified in previous studies.[2] The expanded

139  genomic information obtained in this study is expected to provide new insights into radish growth

140  and development. For example, this study is the first genome report for a giant radish cultivar, and

141  some of the genes unique to the 'Okute-Sakurajima' genome may be involved in giant root

142  formation. Further comparative analysis with radish cultivars with divergent root shapes and sizes

143  would provide insights into the genetic mechanisms contributing to giant root formation.

144

156   **Conflict of interest:** None declared.

157   **Supplementary Data:**

158   **Supplementary Table S1** Program tools used for genome assembly and gene prediction.

159   **Supplementary Table S2** Genetic map length and number of SNPs for F2 radish populations.

160   **Supplementary Figure S1** Genome size estimation for 'Okute-Sakurajima' with the distribution of

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

162

163   **References**

164   1.   Yamagishi, H. 2017, Speciation and Diversification of Radish. In: Nishio, T. and Kitashiba,
165        H. (eds), *The Radish Genome*, Springer, Cham, pp. 11-30.

166   2.   Kitashiba, H., Li, F., Hirakawa, H., et al. 2014, Draft sequences of the radish (*Raphanus*
167        *sativus* L.) genome. *DNA Res*, **21**, 481-490.

168   3.   Mitsui, Y., Shimomura, M., Komatsu, K., et al. 2015, The radish genome and
169        comprehensive gene expression profile of tuberous root formation and development. *Sci*
170        *Rep*, **5**, 10835.

171   4.   Jeong, Y. M., Kim, N., Ahn, B. O., et al. 2016, Elucidating the triplicated ancestral genome
172        structure of radish based on chromosome-level comparison with the *Brassica* genomes.
173        *Theor Appl Genet*, **129**, 1357-1372.

174   5.   Zhang, X., Yue, Z., Mei, S., et al. 2015, A *de novo* genome of a Chinese radish cultivar. *Hort*
175        *Plant J*, **1**, 155-164.

176   6.   Shirasawa, K. and Kitashiba, H. 2017, Genetic Maps and Whole Genome Sequences of
177        Radish. In: Nishio, T. and Kitashiba, H. (eds), *The Radish Genome*, Springer, Cham, pp.
178        31-42.

179   7.   Moghe, G. D., Hufnagel, D. E., Tang, H., et al. 2014, Consequences of Whole-Genome
180        Triplication as Revealed by Comparative Genomic Analyses of the Wild Radish Raphanus
181        raphanistrum and Three Other Brassicaceae Species. *Plant Cell*, **26**, 1925-1937.

182   8.   Nishio, T. and Sakamoto, K. 2017, Polymorphism of Self-Incompatibility Genes. In: Nishio,
183        T. and Kitashiba, H. (eds), *The Radish Genome*, Springer, Cham, pp. 177-188.

184   9.   Michael, T. P. and VanBuren, R. 2020, Building near-complete plant genomes. *Curr Opin*
185        *Plant Biol*, **54**, 26-33.

186   10.  Wang, X., Wang, H., Wang, J., et al. 2011, The genome of the mesopolyploid crop species
187        *Brassica rapa*. *Nat Genet*, **43**, 1035-1039.

188   11.  Arabidopsis Genome, I. 2000, Analysis of the genome sequence of the flowering plant
189        *Arabidopsis thaliana*. *Nature*, **408**, 796-815.

190

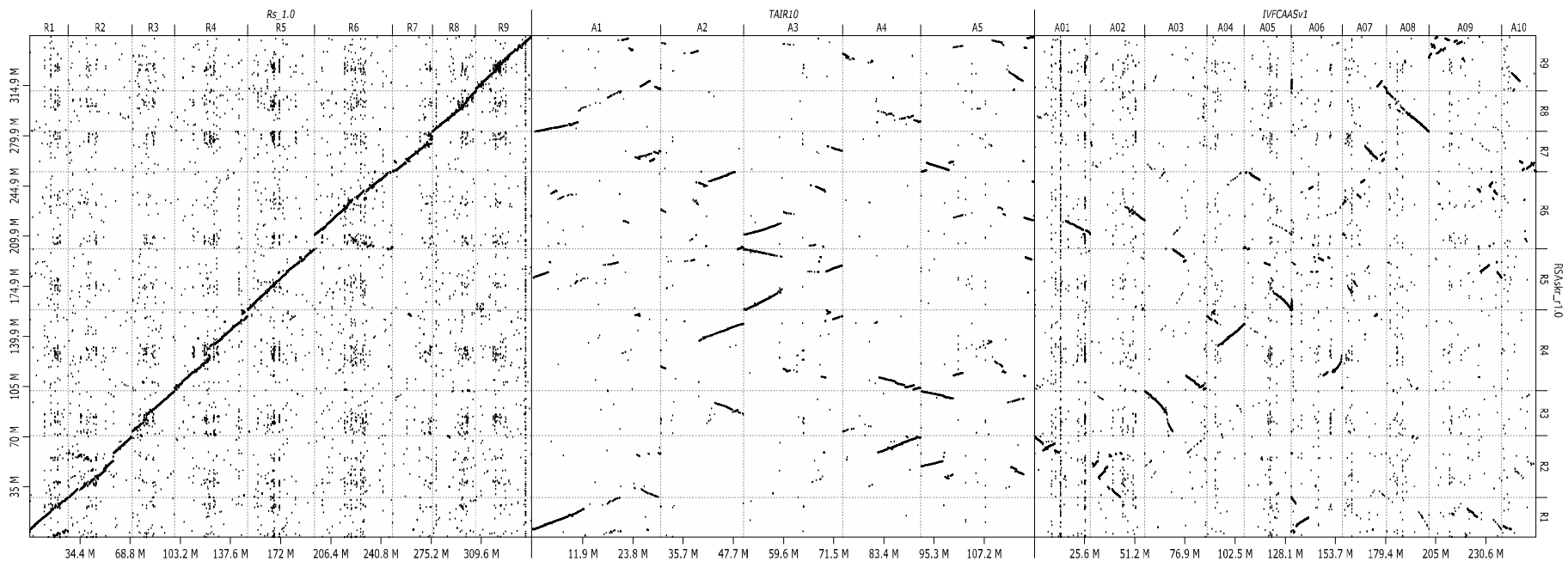191    **Table 1** Statistics of the primary contig sequences of 'Okute-Sakurajima'

|  | RSAskr_r1.0 |
| --- | --- |
| Total contig size (bases) | 504,534,164 |
| Number of contigs | 1,437 |
| N50 length (bases) | 1,247,688 |
| Longest contig size (bases) | 8,317,732 |
| Gap (%) | 0.0 |
| Complete BUSCOs | 94.1 |
| (Single-copy BUSCOs | 71.0) |
| (Duplicated BUSCOs | 23.1) |
| Fragmented BUSCOs | 2.0 |
| Missing BUSCOs | 3.9 |
| #Genes (*ab initio*) | 89,915 |
| #Genes (mapping) | 78,645 |

192

193 **Table 2** Statistics of the 'Okute-Sakurajima' pseudomolecule sequences, RSAskr_r1.0.pmol.

| Chr | #Contigs | (%) | Contig size (bp) | (%) | #Genes | (%) |
|---|---|---|---|---|---|---|
| R1 | 25 | 1.7 | 27,719,058 | 5.5 | 5,426 | 6.0 |
| R2 | 40 | 2.8 | 42,944,316 | 8.5 | 8,295 | 9.2 |
| R3 | 30 | 2.1 | 31,410,669 | 6.2 | 5,923 | 6.6 |
| R4 | 35 | 2.4 | 56,498,296 | 11.2 | 10,843 | 12.1 |
| R5 | 34 | 2.4 | 42,357,306 | 8.4 | 8,477 | 9.4 |
| R6 | 41 | 2.9 | 53,940,652 | 10.7 | 10,403 | 11.6 |
| R7 | 17 | 1.2 | 28,108,325 | 5.6 | 5,545 | 6.2 |
| R8 | 30 | 2.1 | 28,319,830 | 5.6 | 5,474 | 6.1 |
| R9 | 41 | 2.9 | 38,520,653 | 7.6 | 7,143 | 7.9 |
| **R1-R9** | **293** | **20.4** | **349,819,105** | **69.3** | **67,529** | **75.1** |
| R0 | 1,144 | 79.6 | 154,715,059 | 30.7 | 22,386 | 24.9 |
| **Total** | **1,437** | **100.0** | **504,534,164** | **100.0** | **89,915** | **100.0** |

**Figure 1** Comparative maps of the 'Okute-Sakurajima' genome.

Dots indicate sequence similarities of the 'Okute-Sakurajima' genome (RSAskr_r1.0) on the vertical axis with those of *R. sativus* (Rs_1.0), *A. thaliana* (TAIR10), and *B. rapa* (IVFCAASv1) on the horizontal axis.

9