1  **Organelle genome assembly uncovers the dynamic genome reorganization and cytoplasmic male**

2  **sterility associated genes in tomato.**

3  **Running title:** Organelle genomes of CMS tomato

4

5  Kosuke Kuwabara[1], Issei Harada[1], Yuma Matsuzawa[2], Tohru Ariizumi[1,3*], and Kenta Shirasawa[4*]

6  [1]Graduate School of Life and Environmental Sciences, University of Tsukuba, Tsukuba, Ibaraki 305-8577,

7  Japan

8  [2]TOKITA Seed Co. LTD., Kazo, Saitama 349-1144, Japan

9  [3]Tsukuba Plant Innovation Research Center, Tsukuba, Ibaraki 305-8577, Japan

10  [4]Kazusa DNA Research Institute, Kisarazu, Chiba 292-0818, Japan

11

12  **E-mail addresses of all authors**

13  KK: s1921107@s.tsukuba.ac.jp

14  IH: s2021044@s.tsukuba.ac.jp

15  YM: ymatsuzawa@tokitaseed.co.jp

16  TA: ariizumi.toru.ge@u.tsukuba.ac.jp

17  KS: shirasaw@kazusa.or.jp

18

19  **\*Co-corresponding authors:**

20  Tohru Ariizumi

21  Tel. and Fax: +81-29 853 4710

22  Kenta Shirasawa

23  Tel.: +81-438 52 3935

24  Fax: +81-438 52 3934

25

26  **Abstract**

27  To identify cytoplasmic male sterility (CMS)-associated genes in tomato, we determined the genome

28  sequences of mitochondria and chloroplasts in three CMS tomato lines derived from independent

29  asymmetric cell fusions, their nuclear and cytoplasmic donors, and male fertile weedy cultivated tomato

30  and wild relatives. The structures of the CMS mitochondrial genomes were highly divergent from those

31  of the nuclear and cytoplasmic donors, and genes of the donors were mixed up in these genomes. On the

32  other hand, the structures of CMS chloroplast genomes were moderately conserved across the donors, but

33  CMS chloroplast genes were unexpectedly likely derived from the nuclear donors. Comparative analysis

34  of the structures and contents of organelle genes and transcriptome analysis identified three genes that

35  were uniquely present in the CMS lines, but not in the donor or fertile lines. RNA sequencing analysis

36  indicated that these three genes transcriptionally expressed in anther, two of which were also expressed in

37  pollen. They could be potential candidates for CMS-associated genes. This study suggests that organelle

38  reorganization mechanisms after cell fusion events differ between mitochondria and chloroplasts, and

39  provides insight into the development of new F1 hybrid breeding programs employing the CMS system in

40  tomato.

41

42  **Keywords:** Tomato, cytoplasmic male sterility, organelle genomes, mitochondria, RNA-Seq

43

44 **Introduction**

45 Cytoplasmic male sterility (CMS) is broadly found in the kingdom of Plantae[1]. CMS plants cannot

46 produce seeds by self-pollination due to a lack of male fertility; therefore, pollen from other plants is

47 always required for these plants to produce seeds. CMS is caused by the incompatibility of interactions of

48 genetic information between nuclei and organelles, especially mitochondria[1]. The genes in nuclei and

49 organelles are called *restore of fertility* (*RF*) genes and CMS-associated genes, respectively. Therefore,

50 CMS plants have been used as materials for studies of interactions between nuclear and cytoplasmic

51 genes. Moreover, CMS is used in breeding programs to produce F1 hybrid seeds[1], in which cytoplasmic

52 and pollen donors are employed as maternal and paternal parents, respectively.

53 CMS plants can be artificially generated by recurrent backcrossing or transgenic approaches[2,3],

54 which leads to incompatibility between nuclei and organelles. A tomato CMS line, called CMS-pennellii,

55 which possesses nuclei and cytoplasm from *Solanum pennellii* and *Solanum peruvianum*, respectively,

56 has been developed by recurrent backcrossing[2]. A gene knockdown strategy is also used to develop CMS

57 tomato lines, for which expression of a nuclear gene that regulates mitochondrial substoichiometric

58 shifting has been suppressed[3]. In addition, other types of CMS tomato lines have been generated via

59 asymmetric cell fusion between cultivated tomato lines, namely, *Solanum lycopersicum* as the nuclear

60 donor and a wild potato relative, *Solanum acaule*, as the cytoplasmic donor[4]. Among CMS lines, MSA1

61 has been well-studied to reveal nucleus-organelle incompatibility[4]. A physical map of the mitochondrial

62 genome of MSA1 indicates that this asymmetric cell fusion hybrid has a complex mitochondrial genome

63 structure consisting of the parental genomes[5]. Transcripts of an open reading frame (ORF), *orf206*, of the

64 hybrid mitochondrial genome are heterogeneously edited[6]. However, no candidates of CMS-associated

65 genes have been identified in tomato.

66 Although CMS-associated gene sequences are not conserved across plant species, they have

67 common features[7]. Most CMS-associated gene candidates usually possess transmembrane regions and

68 chimeric structures, so-called fusion genes, of genes involved in respiration. Based on this information,

69 CMS-associated genes have been identified in *Oryza sativa*[8,9], *Helianthus annuus*[10], and *Gossypium*

70 *hirsutum*[11]. RNA-sequencing (RNA-Seq) based on next-generation sequencing technology has been

71 employed to select candidates uniquely expressed in CMS lines of *Brassica juncea*[12]. Further functional

72 studies are required to confirm that these candidates are involved in CMS. Introduction of *RF* genes into

73 CMS lines would be a useful approach because CMS-associated genes can be downregulated in the

74 presence of *RF* genes[7]. Another approach is to introduce CMS-associated genes into fertile lines to induce

75 sterility[13]. More recently, it has become possible to alter or edit gene sequences of mitochondrial genomes

76 with TALEN technology[14]. This technology has been used to disrupt CMS-associated genes in

77 mitochondrial genomes and thereby generate *Arabidopsis thaliana*, *Oryza sativa*, and *Brassica napus*

78 with CMS[14,15].

79 In parallel with MSA1, as shown in Figure 1, two asymmetric cell fusions were developed

3

80    between cultivated tomato lines *S. lycopersicum* ('O' and 'P') as nuclear donors and a wild potato relative,

81    *S. acaule*, as the cytoplasmic donor[16]. The nuclear genome backgrounds of the three cell fusion lines

82    including MSA1 were replaced with the genomes of cultivated tomato lines by a repeated backcrossing

83    strategy. The resultant CMS lines are designated ʻCMS[MSA1]ʼ, ʻCMS[O]ʼ, and ʻCMS[P]ʼ. Therefore,

84    it may be possible to identify CMS-associated genes by comparative analysis of the genomes and

85    transcriptomes of the CMS lines and their nuclear donors. In this study, we determined the sequences of

86    the organelle genomes of the CMS lines and their donors. Subsequently, the genome sequences and gene

87    expression patterns were compared to identify CMS-associated gene candidates. Furthermore, the results

88    of this analysis may provide insights into the cytoplasmic genome features of asymmetric cell fusions.

89

90    **Results**

91    *De novo assembly of chloroplast and mitochondrial genomes*

92    A total of 10.5 Gb reads per sample were obtained from three CMS tomato lines ( ʻCMS[MSA1]ʼ,

93    ʻCMS[O]ʼ, and ʻCMS[P]ʼ), three nuclear donors ('Sekai-ichi', 'O', and 'P'), and one cytoplasmic

94    donor (*S. acaule*). Of them, 374 Mb (3.6%) and 566 Mb (5.4%) of reads per sample were aligned on

95    publicly available sequences of mitochondrial and chloroplast genomes, respectively. The reads mapped

96    on the two sets of reference sequences were separately assembled into contig sequences.

97    Mitochondrial genome sequences were constructed with reads mapped on the mitochondrial

98    reference sequences (Table 1). The mitochondrial genomes of the nuclear donors 'Sekai-ichi', 'O', and 'P'

99    were all constructed from only contigs with assembly sizes of 562.6 kb ($n = 2$, $n$ represents contig

100    numbers), 536.9 kb ($n = 2$), and 553.3 kb ($n = 2$), respectively. In *S. acaule*, 728.4 kb contigs ($n = 7$) for

101    the mitochondrial genome were established. The assembly sizes were longer in the CMS lines than in the

102    nuclear and cytoplasmic donors, specifically, they were 995.2 kb ($n = 7$) in ʻCMS[MSA1]ʼ, 968.4 kb ($n$

103    $= 7$) in ʻCMS[O]ʼ, and 829.3 kb ($n = 5$) in ʻCMS[P]ʼ. For chloroplast genomes, total sequence lengths

104    of 389.2 kb ($n = 2$), 349.5 kb ($n = 2$), and 346.9 kb ($n = 2$) were constructed for 'Sekai-ichi', 'O', and 'P',

105    respectively (Table 1). There were two contig sequences in each of the three nuclear donors. The

106    assembly sizes were shorter in ʻCMS[MSA1]ʼ (296.6 kb, $n = 1$) and ʻCMS[O]ʼ (307.1 kb, $n = 1$) than in

107    the nuclear donors, but longer in ʻCMS[P]ʼ (454.1 kb, $n = 3$).

108    Comparative genome analysis revealed that the mitochondrial genomes of the CMS lines consisted

109    of highly fragmented, repeated, and duplicated sequences derived from both donors throughout the

110    genome (Figure 2). On the other hand, the structures of the chloroplast genomes of the CMS lines were

111    moderately conserved across the nuclear and cytoplasmic donors (Figure 2).

112    In parallel, we determined the mitochondrial and chloroplast genome sequences of *Solanum*

113    *pimpinellifolium* LA1670 and *S. lycopersicum* var. *cerasiforme* LA1673 (Table 1). Sequence reads were

114    obtained from a public DNA database and processed as described above. Assembly sizes of the

115    mitochondrial and chloroplast genomes were 620.6 kb ($n = 3$) and 299.4 kb ($n = 1$) for *S. pimpinellifolium*

4

116    LA1670, respectively, and 569.9 kb (*n* = 2) and 337.7 kb (*n* = 2) for *S. lycopersicum* var. *cerasiforme*

117    LA1673, respectively.

118

119    *Gene prediction from the organelle genomes*

120    ORFs encoding ≥25 amino acids were extracted from the assembled sequences to predict potential genes.

121    The number of potential genes predicted from the chloroplast genome assemblies ranged from 5,130 (*S.*

122    *acaule*) to 8,165 (ʻCMS[P]ʼ) and the number of potential genes predicted from the mitochondrial

123    sequences ranged from 10,326 (ʻOʼ) to 19,170 (ʻCMS[MSA1]ʼ) (Table 1).

124         The ORFs were clustered to identify genes unique to and shared among the CMS lines, nuclear

125    donors, and cytoplasmic donor (Figure 3). The ORFs in the CMS mitochondrial genomes consisted of

126    four types of genes, namely, those unique to the CMS lines (Type 1: 9.4–11.9%), those shared with the

127    nuclear donors only (Type 2: 14.1–17.0%), those shared with the cytoplasmic donor only (Type 3:

128    8.9–13.2%), and those shared with both the nuclear and cytoplasmic donors (Type 4: 61.8–64.1%). By

129    contrast, the ORFs in the CMS chloroplast genomes mostly consisted of three types of genes, namely,

130    those unique to the CMS lines (Type 1: 1.2–5.9%), those shared with the nuclear donors only (Type 2:

131    31.2–33.1%), and those shared with both the nuclear and cytoplasmic donors (Type 4: 62.9–65.7%). Few

132    genes shared with the cytoplasmic donor only were found (Type 3: up to 0.1%).

133         The genome positions of the genes differed according to the gene type and organelle (Figure 4).

134    Type 1 genes in mitochondria were distributed across the genome with some gaps. The positions of Type

135    2 genes were basically the same as those of Type 1 genes, while Type 3 genes were located in the gaps

136    between Type 1 genes. Type 4 genes were also located in the gaps and at the ends of contig sequences. On

137    the other hand, in chloroplast genomes, the positions of Type 1 and 2 genes overlapped and Type 4 genes

138    were located at the ends of contigs.

139

140    *Screening of CMS-associated gene candidates*

141    To identify candidates of CMS-associated genes in the mitochondrial genomes, we set the following four

142    criteria: 1) amino acid length ≥70, 2) absent from male fertile lines, 3) present in all three CMS lines, and

143    4) expressed in anthers of the CMS lines. Among the predicted genes in the ʻCMS[P]ʼ, ʻCMS[MSA1]ʼ,

144    and ʻCMS[O]ʼ mitochondrial genomes, 831, 1,025, and 969 genes encoded ≥70 amino acids,

145    respectively. The gene sequences from the CMS lines were compared with the mitochondrial genomes of

146    the nuclear donors (ʻSekai-ichiʼ, ʻPʼ, and ʻOʼ) and *S. pimpinellifolium* LA1670, *S. lycopersicum* var.

147    *cerasiforme* LA1673), *S. pennellii*, and *Nicotiana tabacum*. In total, 183, 272, and 140 genes were

148    selected because they were absent from the nuclear donors and Solanaceae relatives, all of which possess

149    male fertility. Furthermore, we selected 36, 41, and 33 genes commonly present in the CMS lines. The

150    copy numbers of the genes varied. Finally, RNA-Seq reads were mapped on the mitochondrial genomes

151    of the CMS lines. This analysis limited the number of CMS-associated gene candidates to four, including

5

152    two identical sequences. The three genes were named *orf137* (two copies in the genome of each CMS

153    line: CMS-PMt002g07240 and CMS-PMt005g13392), *orf193* (one copy: CMS-PMt002g06465), and

154    *orf265* (one copy: CMS-PMt010g15739).

155        *De novo* transcriptome assembly was performed in parallel. RNA-Seq data were obtained from the

156    anthers of 'P' and ﹃CMS[P]', and assembled into 62 and 43 transcript sequences, respectively, of which

157    37 'P' and 18 ﹃CMS[P]' transcripts were predicted to have transmembrane domains. Of these sequences,

158    eight were uniquely detected in ﹃CMS[P]'. Two genes (STRG.32.1.p1 and STRG.39.1.p1) were identical

159    to *orf137* and *orf265*.

160        Because two genes were commonly identified in both analyses, a total of nine genes were finally

161    selected as candidates of CMS-associated genes (Table 2). Sequence similarity searches with the

162    mitochondrial and chloroplast genomes indicated that two copies of the STRG.32.1.p1 (*orf137*) sequence

163    (CMS-PMt002g07240 and CMS-PMt005g13392) were present in the mitochondrial genomes of the three

164    CMS lines. A single copy sequence of *orf193* (CMS-PMt002g06465) and a single copy sequence of

165    STRG.39.1.p1 (*orf265*, CMS-PMt010g15739) were found in the mitochondrial genomes of the three

166    CMS lines in addition to that of *S. acaule*. The presence of the three genes in the CMS lines was validated

167    by a PCR assay with the three CMS lines and six fertile lines. The remaining six genes were found in

168    both the CMS and fertile lines. We selected three genes, *orf137*, *orf193*, and *orf265*, as highly potential

169    candidates for CMS-associated genes due to their presence specifically in the CMS mitochondrial

170    genomes and their expression in anthers.

171

172    *Sequence similarity analysis of the candidate genes*

173    The sequence similarity of the candidate genes including their flanking genome regions in the

174    mitochondrial genome of ﹃CMS[P]' was investigated. A 3,045 bp genome sequence around *orf193*

175    showed high sequence similarity to a 4,682 bp region of the tomato chloroplast genome sequences. The

176    3,045 bp sequence was split into three sequences containing 1,590, 488, and 1,007 bp (Figure 5A) with

177    highly conserved boundary sequences (Figure 5B). In the 1,590 bp chloroplast genome sequence, a gene

178    encoding *cytochrome f* was encoded; however, the corresponding sequence in the mitochondrial genome

179    had a single base insertion causing a frame-shift mutation (Figure 5C). This mutation broke the ORF of

180    the *cytochrome f* gene and generated two small ORFs, *orf116* and *orf193*.

181        A portion of *orf265* and its upstream sequences (177 bp in total) showed high similarity to the *ATP*

182    *synthase subunit 8* (*atp8*) gene encoded in the tomato mitochondrial genome (Figure 5D). The remaining

183    sequences of *orf265* lacked similarity to reported sequences. *orf265* was located upstream of the *nad3* and

184    *rps12* genes in the mitochondrial genome. No sequence similarity was observed for *orf137* and the

185    flanking sequence.

186

187    *Expression analysis of the candidate genes*

6

188    The expression patterns of the candidate genes, *orf137*, *orf193*, and *orf265*, were investigated by RT-PCR.

189    First, we validated the results of the transcriptome analysis by detecting expression of the three genes in

190    anthers of 'CMS[P]' and 'CMS[MSA1]' (Figure 6A). *orf265* was tandemly arrayed with *nad3* and

191    *rps12*; therefore, we assumed that these three genes were co-transcribed as an operon. As expected,

192    transcripts spanning the three genes were also detected (Figure 6A). Next, we analyzed gene expression

193    in leaves, stems, roots, ovaries, and pollen in addition to anthers of Dwarf 'CMS[P]' which was a BC3

194    generation of 'CMS[P]' backcrossed with a tomato dwarf cultivar 'Micro-Tom'. Expression of *orf137*

195    and *orf265* was detected in all tested tissues, while that of *orf193* was observed in leaves, stems, roots,

196    ovaries, and anthers (Figure 6B).

197

198    **Discussion**

199    We determined the mitochondrial and chloroplast genome sequences of CMS lines derived from

200    asymmetric cell fusions and those of their nuclear and cytoplasmic donors (Table 1). Comparative

201    analysis of the structures unexpectedly revealed that the cytoplasmic genome structures of the fusions

202    were rearranged and divergent from those of the cytoplasmic donor (*S. acaule*) and nuclear donors (*S.*

203    *lycopersicum*) (Figure 2). CMS-PMt003g09846 and CMS-PMt003g11185 were encoded in both the

204    mitochondrial and chloroplast genomes (Table 2), suggesting that mitochondria and chloroplasts from the

205    two donors were fused with each other and reorganized even though the cytoplasm of the nuclear donors

206    was chemically inactivated to generate asymmetric cell fusions. Interestingly, the mitochondrial genomes

207    of the CMS lines were larger than those of the donors, while the size of chloroplast genomes among the

208    CMS lines were equivalent (Table 1). In addition, the structures of the mitochondrial genomes were

209    divergent, while those of the chloroplast genomes were rather conserved (Figure 2). Gene clustering

210    analysis suggested that both the cytoplasmic and nuclear donors contributed to form mitochondria in the

211    CMS lines (Figure 3). Furthermore, the structures of the CMS mitochondrial genomes contained patches

212    of the two genomes of the donors (Figure 4). These results suggest that the mitochondrial genomes of

213    both donors were highly fragmented at the time of asymmetric cell fusion and reorganized to form a new

214    mitochondrial genome[5]. This is completely different from our expectation that genomes of the

215    cytoplasmic donors should be present in CMS lines derived from asymmetric cell fusions. More

216    interestingly, chloroplasts of the CMS lines consisted only of genes from nuclear donors, not from the

217    cytoplasmic donor. This unexpected finding has been frequently made in tomato[17], tobacco[18], and

218    *Brassica*[19]. We speculate that interactions of genetic information between nuclei and organelles might be

219    strict with chloroplasts rather than with mitochondria. The genome and/or organelle reorganization

220    mechanisms after cell fusions might differ between mitochondria and chloroplasts.

221    Based on genome and transcriptome analyses, nine genes encoded in the mitochondrial genome of

222    'CMS[P]' were selected as candidate CMS-associated genes (Table 2). Among them, three genes

223    (*orf193*, STRG.32.1.p1 = *orf137* and STRG.39.1.p1 = *orf265*) were uniquely present in the genomes of

224    the CMS lines and expressed in their anthers (Table 2 and Figure 6). STRG.32.1.p1 (*orf137*) showed

225    sequence similarity with the CMS-associated protein encoding *cytochrome c subunit 1* (Figure 5).

226    STRG.39.1.p1 (*orf265*) was similar to *ATP synthase subunit 8* at the N-terminus, but lacked similarity in

227    the remaining regions (Figure 5). CMS-associated genes are generally involved in cellular respiration

228    producing energy to generate pollen[20], and this is true of both these genes. In many cases, fusion genes

229    have been reported to be CMS-associated genes, e.g., *orf307* in *Oryza sativa*[21] and *orf72* in *Brassica*

230    *oleracea*[22], and to produce cytotoxic proteins, which lead to male sterility. Knockout mutagenesis with

231    mitoTALENs[14,15] targeting the candidate genes would be useful to identify CMS-associated genes and to

232    generate CMS lines from normal tomato cultivars.

233    CMS lines are powerful tools to produce F1 hybrid seeds in breeding programs[1]. However, in

234    cereals and fruits including tomato, the *RF* genes are essential for F1 plants to set seeds and bear fruits.

235    Restorer genes for CMS lines have been identified in wild tomato relatives, e.g., *S. pimpinellifolium*

236    LA1670 and *S. lycopersicum* var. *cerasiforme* LA1673[16]. Recently, we published the genome sequence

237    data of these two wild relatives[23]. We expect *RF* genes for CMS lines to be discovered soon based on this

238    information, although no candidate genes or genetic loci have been reported. Once CMS-associated genes

239    and restorer genes are identified, tomato F1 hybrid seeds can be produced by employing insect pollinators

240    instead of the currently used hand-pollination systems. We propose that CMS-based F1 hybrid breeding

241    programs with insect pollinators can be implemented in tomato breeding programs to reduce the costs of

242    F1 seed production in the future.

243

244    **Materials and methods**

245    *Plant materials*

246    Three tomato CMS lines (ʻCMS[MSA1]', ʻCMS[O]', and ʻCMS[P]'), three cultivated tomato lines (*S.*

247    *lycopersicum* 'Sekai-ichi', 'O', and 'P'), and one potato wild relative (*S. acaule*) were used (Figure 1).

248    ʻCMS[MSA1]' was developed by repeated backcrossing using 'O' as a recurrent parent and a

249    male-sterile tomato, MSA1, as a cytoplasmic donor. MSA1 is an asymmetric cell fusion between the

250    tomato cultivar Sekai-ichi (as the nuclear donor) and the potato wild relative *S. acaule* (as the cytoplasmic

251    donor)[4]. ʻCMS[O]' was a progeny in repeated backcrossing using 'O' as the paternal parent and an

252    asymmetric cell fusion between 'O' (as the nuclear donor) and *S. acaule* (as the cytoplasmic donor).

253    ʻCMS[P]' was also a progeny in backcrossing using 'P' as the paternal parent and an asymmetric cell

254    fusion between 'P' (as the nuclear donor) and *S. acaule* (as the cytoplasmic donor). Dwarf ʻCMS[P]' was

255    developed from ʻCMS[P]' by backcrossing with *S. lycopersicum* ʻMicro-Tom' (TOMJPF0001), which

256    is a miniature dwarf cultivar[24]. The putative nuclear and cytoplasmic genomes of the materials are shown

257    in Figure 1.

258

259    *Genome sequence analysis*

8

260    Total genomic DNA was extracted from young leaves of the six tomato lines ( ˋCMS[MSA1]',
261    ˋCMS[O]', ˋCMS[P]', 'Sekai-ichi', 'O', and 'P') and *S. acaule* with a Maxwell 16 Instrument and
262    Maxwell 16 Tissue DNA Purification Kits (Promega, Madison, WI, USA). SMRT sequence libraries were
263    constructed with an SMRTbell Express Template Prep Kit (PacBio, Menlo Park, CA, USA) and used for
264    sequencing on a PacBio Sequel system (PacBio). Genome sequence data for S. *pimpinellifolium* LA1670
265    and *S. lycopersicum* var. *cerasiforme* LA1673 were obtained from a public DNA database (DRA
266    accession numbers DRX231405 and DRX231409)[23].

267

268    *Genome assembly and gene prediction*
269    Sequence reads were mapped on reference genome sequences for mitochondria (GenBank accession
270    numbers MF034192, MF034193, NC_035964, and MF98995–MF989957) or chloroplasts (NC_007898)
271    with Organelle_PBA[25]. Reads mapped on the reference sequences were assembled into contig sequences
272    with Canu[26]. Potential sequence errors in the contig sequences were corrected twice with the sequence
273    reads by Arrow (PacBio). The corrected contig sequences were aligned back to the reference sequences
274    with Nucmer[27] to select highly confident organelle genomes. ORFs ($\geq$75 bases) in the organelle genomes
275    were selected as potential genes with ORFfinder (https://www.ncbi.nlm.nih.gov/orffinder). The ORF
276    sequences were clustered with CD-HIT[28]. Transmembrane domains in the gene sequences were predicted
277    by TMHMM[29]. Sequence similarity searches with the mitochondrial genomes of *S. pennellii*
278    (NC_035964) and *N. tabacum* (NC_006581) were performed by BLAST[30] with a threshold E-value of
279    1e-50.

280

281    *RNA expression analysis*
282    Total RNA was extracted from the anthers of P and ˋCMS[P]' with an RNeasy Plant Mini Kit (QIAGEN,
283    Hilden, Germany). RNA was treated with RNase-free DNase (QIAGEN) and used for sequence library
284    preparation with a TruSeq Stranded mRNA Library Prep Kit (Illumina, San Diego, CA, USA). The
285    resultant libraries were sequenced on NextSeq500 (Illumina) in paired-end, 151 bp mode. After trimming
286    adaptors and low-quality reads by Trim_galore (https://github.com/FelixKrueger/TrimGalore) with option
287    -q 30 --length 100 followed by fastp[31] with option -l 100, transcriptomes were *de novo* assembled by the
288    HiSat2-Stringtie pipeline[32] and putative ORFs were searched for annotation by BLASTP[30] against the
289    SWISS-PROT database[33].
290         To validate the RNA-Seq results, RT-PCR was performed. In total, 800 ng of total RNA isolated
291    from anthers of the CMS lines or seedlings of *S. acaule* was converted into cDNA with ReverTra Ace
292    (TOYOBO, Osaka, Japan) using a random primer (TAKARA BIO, Kusatsu, Japan). cDNA diluted
293    10-fold with water was used as a template for PCR. The PCR mixture (10 µL) contained 0.5 µL cDNA,
294    0.3 µM primers (Table 3), 2× PCR buffer (TOYOBO), 400 µM dNTPs, and 1 U DNA polymerase (KOD
295    FX Neo, TOYOBO). The thermal cycling conditions were as follows: initial denaturation at 94°C for 3

296     min; 35 cycles of denaturation at 98℃ for 15 s, annealing at 60☐ for 30 s, and extension at 68℃ for 60 s;

297     and a final extension at 68°C for 3 min. PCR products were separated by electrophoresis in a 1% agarose

298     gel with TAE buffer. Gels were stained with Midori Green Advance (NIPPON Genetics, Tokyo, Japan) to

299     detect DNA bands under ultraviolet illumination.

300

301     **Data availability**

302     The DDBJ accession numbers of the assembled sequences are LC613090-LC613141. Genome

303     information is available at KaTomicsDB (http://www.kazusa.or.jp/tomato).

304

305     **Acknowledgments**

311

312     **Conflicts of interest**

313     YM is an employee of TOKITA Seed Co. LTD. All other authors declare no competing interests.

314

315     **Contributions**

316     TA and KS conceived and coordinated the project. YM established the plant materials. KK, IH, and KS

317     collected the data. KK, IH, TA, and KS analyzed and interpreted the data. KK and KS wrote the

318     manuscript with contributions from TA. All authors read and approved the final manuscript.

319

320     **References**

321     1.      Bohra, A., Jha, U.C., Adhimoolam, P., Bisht, D. & Singh, N.P. Cytoplasmic male sterility (CMS)

322             in    hybrid    breeding    in    field    crops.    *Plant    Cell    Rep*    **35**,    967-93

323             http://dx.doi.org/10.1007/s00299-016-1949-3 (2016).

324     2.      Petrova, M. *et al.* Characterisation of a cytoplasmic male-sterile hybrid line between

325             Lycopersicon peruvianum Mill. ×Lycopersicon pennellii Corr. and its crosses with cultivated

326             tomato. *Theoretical and Applied Genetics* **98**, 825-830 http://dx.doi.org/10.1007/s001220051139

327             (1999).

328     3.      Sandhu, A.P., Abdelnoor, R.V. & Mackenzie, S.A. Transgenic induction of mitochondrial

329             rearrangements for cytoplasmic male sterility in crop plants. *Proc Natl Acad Sci U S A* **104**,

330             1766-70 http://dx.doi.org/10.1073/pnas.0609344104 (2007).

331     4.      Melchers, G., Mohri, Y., Watanabe, K., Wakabayashi, S. & Harada, K. One-step generation of

332  cytoplasmic male sterility by fusion of mitochondrial-inactivated tomato protoplasts with
333  nuclear-inactivated Solanum protoplasts. *Proc Natl Acad Sci U S A* **89**, 6832-6
334  http://dx.doi.org/10.1073/pnas.89.15.6832 (1992).

335  5.  Shikanai, T., Kaneko, H., Nakata, S., Harada, K. & Watanabe, K. Mitochondrial genome
336  structure of a cytoplasmic hybrid between tomato and wild potato. *Plant Cell Rep* **17**, 832-836
337  http://dx.doi.org/10.1007/s002990050493 (1998).

338  6.  Shikanai, T., Nakata, S., Harada, K. & Watanabe, K. Analysis of the heterogeneous transcripts of
339  the highly edited orf206 in tomato mitochondria. *Plant Cell Physiol* **37**, 692-6
340  http://dx.doi.org/10.1093/oxfordjournals.pcp.a029000 (1996).

341  7.  Chen, L. & Liu, Y.G. Male sterility and fertility restoration in crops. *Annu Rev Plant Biol* **65**,
342  579-606 http://dx.doi.org/10.1146/annurev-arplant-050213-040119 (2014).

343  8.  Igarashi, K., Kazama, T., Motomura, K. & Toriyama, K. Whole genomic sequencing of RT98
344  mitochondria derived from Oryza rufipogon and northern blot analysis to uncover a cytoplasmic
345  male sterility-associated gene. *Plant Cell Physiol* **54**, 237-43
346  http://dx.doi.org/10.1093/pcp/pcs177 (2013).

347  9.  Okazaki, M., Kazama, T., Murata, H., Motomura, K. & Toriyama, K. Whole mitochondrial
348  genome sequencing and transcriptional analysis to uncover an RT102-type cytoplasmic male
349  sterility-associated candidate Gene Derived from Oryza rufipogon. *Plant Cell Physiol* **54**,
350  1560-8 http://dx.doi.org/10.1093/pcp/pct102 (2013).

351  10.  Makarenko, M.S. *et al.* Characterization of the mitochondrial genome of the MAX1 type of
352  cytoplasmic male-sterile sunflower. *BMC Plant Biol* **19**, 51
353  http://dx.doi.org/10.1186/s12870-019-1637-x (2019).

354  11.  Li, S. *et al.* The comparison of four mitochondrial genomes reveals cytoplasmic male sterility
355  candidate genes in cotton. *BMC Genomics* **19**, 775 http://dx.doi.org/10.1186/s12864-018-5122-y
356  (2018).

357  12.  Wu, Z. *et al.* Mitochondrial genome and transcriptome analysis of five alloplasmic male-sterile
358  lines in Brassica juncea. *BMC Genomics* **20**, 348 http://dx.doi.org/10.1186/s12864-019-5721-2
359  (2019).

360  13.  Yang, J., Liu, X., Yang, X. & Zhang, M. Mitochondrially-targeted expression of a cytoplasmic
361  male sterility-associated orf220 gene causes male sterility in Brassica juncea. *BMC Plant Biol* **10**,
362  231 http://dx.doi.org/10.1186/1471-2229-10-231 (2010).

363  14.  Kazama, T. *et al.* Curing cytoplasmic male sterility via TALEN-mediated mitochondrial genome
364  editing. *Nat Plants* **5**, 722-730 http://dx.doi.org/10.1038/s41477-019-0459-z (2019).

365  15.  Arimura, S.I. *et al.* Targeted gene disruption of ATP synthases 6-1 and 6-2 in the mitochondrial
366  genome of Arabidopsis thaliana by mitoTALENs. *Plant J* **104**, 1459-1471
367  http://dx.doi.org/10.1111/tpj.15041 (2020).

368    16.    Harada, K., Kondo, K. & Watabe, K. Creation of f1 tomato plant recovered in fertility. (1995).

369    17.    Bonnema, A.B., Melzer, J.M. & O'Connell, M.A. Tomato cybrids with mitochondrial DNA from
370           Lycopersicon pennelli. *Theor Appl Genet* **81**, 339-48  http://dx.doi.org/10.1007/BF00228674
371           (1991).

372    18.    Sidorov, V.A., Menczel, L., Nagy, F. & Maliga, P. Chloroplast transfer in Nicotiana based on
373           metabolic complementation between irradiated and iodoacetate treated protoplasts. *Planta* **152**,
374           341-5 http://dx.doi.org/10.1007/BF00388259 (1981).

375    19.    Morgan, A. & Maliga, P. Rapid chloroplast segregation and recombination of mitochondrial
376           DNA in Brassica cybrids. *Mol Gen Genet* **209**, 240-6  http://dx.doi.org/10.1007/BF00329649
377           (1987).

378    20.    Touzet, P. & Meyer, E.H. Cytoplasmic male sterility and mitochondrial metabolism in plants.
379           *Mitochondrion* **19 Pt B**, 166-71 http://dx.doi.org/10.1016/j.mito.2014.04.009 (2014).

380    21.    Fujii, S., Kazama, T., Yamada, M. & Toriyama, K. Discovery of global genomic re-organization
381           based on comparison of two newly sequenced rice mitochondrial genomes with cytoplasmic
382           male sterility-related genes. *BMC Genomics* **11**, 209 http://dx.doi.org/10.1186/1471-2164-11-209
383           (2010).

384    22.    Shinada, T., Kikuchi, Y., Fujimoto, R. & Kishitani, S. An alloplasmic male-sterile line of
385           Brassica oleracea harboring the mitochondria from Diplotaxis muralis expresses a novel
386           chimeric    open    reading    frame,    orf72.    *Plant    Cell    Physiol*    **47**,    549-53
387           http://dx.doi.org/10.1093/pcp/pcj014 (2006).

388    23.    Takei, H. *et al.* De novo genome assembly of two tomato ancestors, solanum pimpinellifolium
389           and    S.    lycopersicum    var.    cerasiforme,    by    long-read    sequencing.    *DNA    Res*
390           http://dx.doi.org/10.1093/dnares/dsaa029 (2021).

391    24.    Scott, J.W. & Harbaugh, B.K. Micro-tom: A Miniature Dwarf Tomato (Agricultural Experiment
392           Station, Institute of Food and Agricultural Sciences, University of Florida, 1989).

393    25.    Soorni, A., Haak, D., Zaitlin, D. & Bombarely, A. Organelle_PBA, a pipeline for assembling
394           chloroplast and mitochondrial genomes from PacBio DNA sequencing data. *BMC Genomics* **18**,
395           49 http://dx.doi.org/10.1186/s12864-016-3412-9 (2017).

396    26.    Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting
397           and repeat separation. *Genome Res* **27**, 722-736 http://dx.doi.org/10.1101/gr.215087.116 (2017).

398    27.    Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol* **5**, R12
399           http://dx.doi.org/10.1186/gb-2004-5-2-r12 (2004).

400    28.    Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation
401           sequencing    data.    *Bioinformatics*    **28**,    3150-2    http://dx.doi.org/10.1093/bioinformatics/bts565
402           (2012).

403    29.    Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E.L. Predicting transmembrane protein

404    topology with a hidden Markov model: application to complete genomes. *J Mol Biol* **305**,
405    567-80 http://dx.doi.org/10.1006/jmbi.2000.4315 (2001).

406  30.  Altschul, S.F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database
407    search programs. *Nucleic Acids Res* **25**, 3389-402  http://dx.doi.org/10.1093/nar/25.17.3389
408    (1997).

409  31.  Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
410    *Bioinformatics* **34**, i884-i890 http://dx.doi.org/10.1093/bioinformatics/bty560 (2018).

411  32.  Pertea, M., Kim, D., Pertea, G.M., Leek, J.T. & Salzberg, S.L. Transcript-level expression
412    analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat Protoc* **11**, 1650-67
413    http://dx.doi.org/10.1038/nprot.2016.095 (2016).

414  33.  Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its supplement
415    TrEMBL in 2000. *Nucleic Acids Res* **28**, 45-8 http://dx.doi.org/10.1093/nar/28.1.45 (2000).

416
417

418 **Figure legends**

419 **Figure 1** Pedigree of the CMS tomato lines.

420 Squares and circles indicate cytoplasm and nuclei, respectively. Arrows with dashed lines and doubled

421 lines indicate cell fusions and crossings, respectively.

422 **Figure 2** Comparative maps of the organelle genomes of the CMS tomato lines.

423 Mitochondrial (A) and chloroplast (B) genomes of the three CMS lines, nuclear donors, and cytoplasmic

424 donor. Dots indicate sequence similarity between the genome sequences.

425 **Figure 3** Organelle genes in the CMS lines, nuclear donors, and cytoplasmic donor.

426 Numbers of genes unique to the CMS lines, nuclear donors, and cytoplasmic donor are indicated in bold,

427 standard, and italic fonts, respectively. Percentages of genes are shown in parentheses.

428 **Figure 4** Distributions of CMS tomato genes across the organelle genomes.

429 Dots indicate gene positions on contig sequences of the organelle genomes. Genes are grouped into the

430 following four types: Type 1, genes unique to the CMS lines; Type 2, genes shared with the nuclear

431 donors; Type 3, genes shared with the cytoplasmic donor; and Type 4, genes shared with both the nuclear

432 and cytoplasmic donors.

433 **Figure 5** Structures of mitochondrial genes in 'CMS[P]'.

434 A. Genome structure of the *orf193* region. Homologous sequences between the two genomes are

435 indicated by gray boxes. Highly conserved sequences at the borders are shown in red and blue. B.

436 Sequence alignments of the borders. C. Details of the genome structure of the *orf193* region. A single

437 nucleotide insertion causing a frame-shift mutation is indicated with a red arrow. D. Genome structure of

438 the *orf265* region.

439 **Figure 6** RT-PCR analysis of the CMS-associated gene candidates.

440 Gene expression patterns in anthers of two CMS lines (A) and in seven samples of Dwarf 'CMS[P]' (B).

441 *cox2* is a positive control.

442

14

443 **Table 1** Assembly data of the organelle genomes.

| Organelle | Male sterile | | | Male fertile | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | CMS lines | | | Nuclear donors | | | Cytoplasmic donor | Tomato wild relative and weedy tomato | |
| | CMS[MSA1] | CMS[O] | CMS[P] | Sekai-ichi | O | P | *S. acaule* | LA1670 | LA1673 |
| ***Mitochondrion*** | | | | | | | | | |
| Number of sequences | 7 | 7 | 5 | 2 | 2 | 2 | 7 | 3 | 2 |
| Total length (bp) | 995,217 | 968,425 | 829,310 | 562,630 | 536,932 | 553,289 | 728,387 | 620,567 | 569,852 |
| Number of genes | 19,170 | 18,623 | 15,912 | 10,782 | 10,326 | 10,653 | 13,898 | 11,920 | 10,965 |
| ***Chloroplast*** | | | | | | | | | |
| Number of sequences | 1 | 1 | 3 | 2 | 2 | 2 | 1 | 1 | 2 |
| Total length (bp) | 296,583 | 307,105 | 454,083 | 389,209 | 349,506 | 346,936 | 284,040 | 299,393 | 337,655 |
| Number of genes | 5,279 | 5,456 | 8,165 | 6,971 | 6,267 | 6,192 | 5,130 | 5,346 | 5,995 |

444

445    **Table 2** Copy numbers of CMS-associated gene candidates in the organelle genomes.

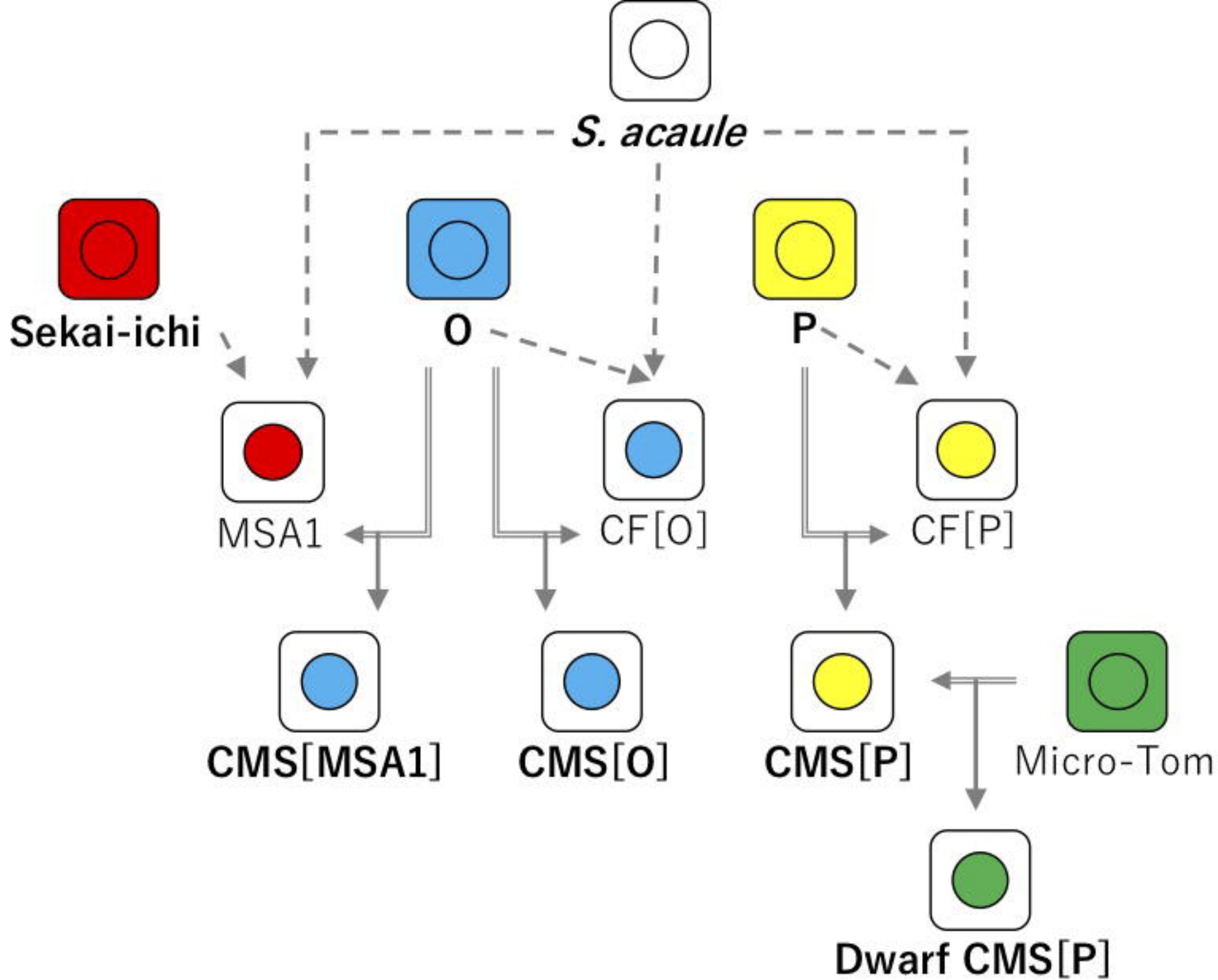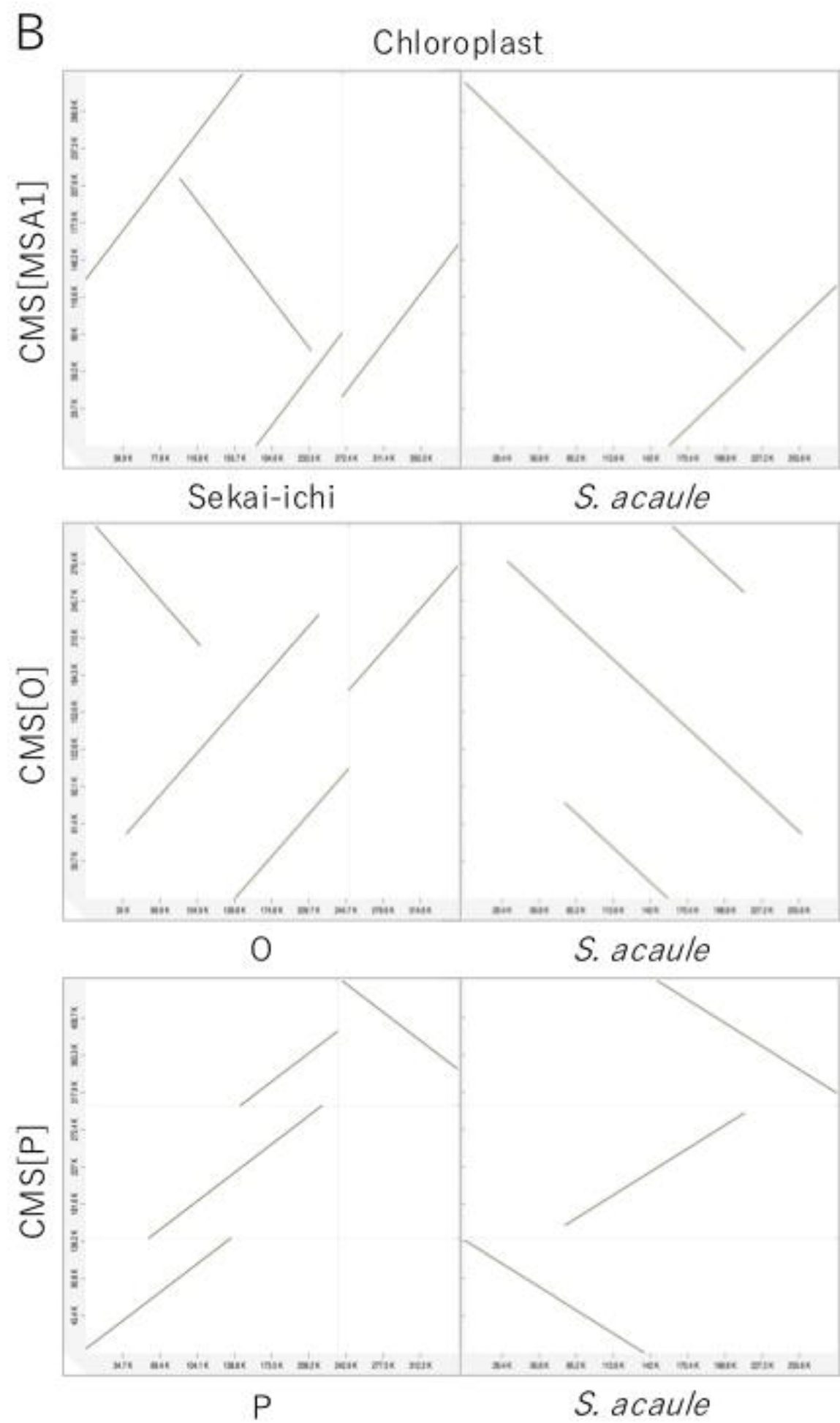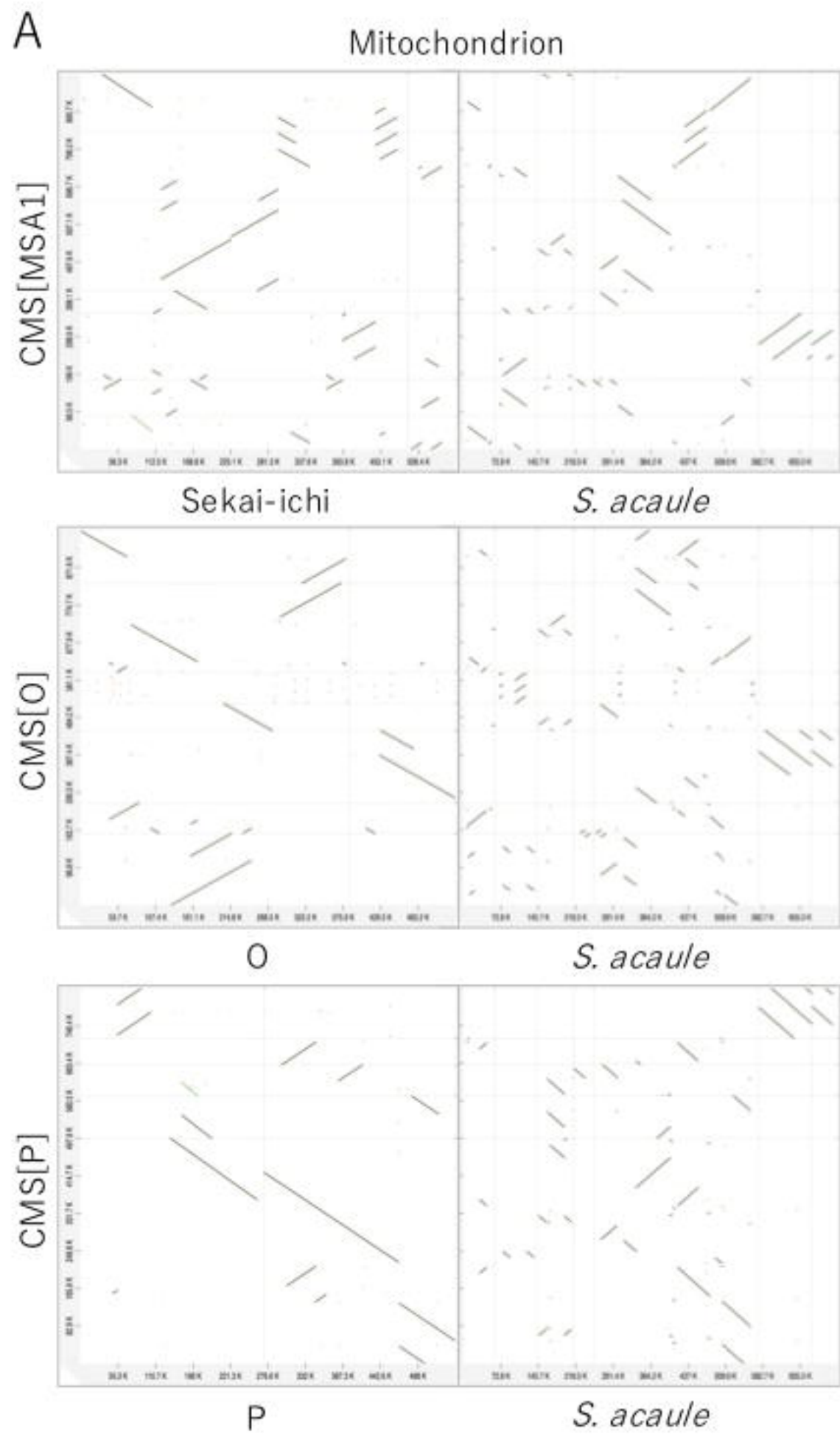| Gene ID of CMS[PF] | Candidates from genome analysis | Candidates from transcriptome analysis | Copy number in mitochondrial genome | | | | | | | Copy number in chloroplast genome | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | CMS[MSA1]' | CMS[O] | CMS[P] | Sekai-ichi | O | P | S. acaule | CMS[MSA1]' | CMS[O] | CMS[P] | Sekai-ichi | O | P | S. acaule |
| CMS-PMt002g06465 | *orf193* | | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CMS-PMt002g07240 and CMS-PMt005g13392 | *orf137* | STRG.32.1.p1 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CMS-PMt002g07993, CMS-PMt003g11130, and CMS-PMt004g12510 | | STRG.22.1.p1 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CMS-PMt003g09515 | | STRG.5.1.p1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CMS-PMt003g09846 | | STRG.8.1.p1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 3 | 2 | 1 | 2 | 1 |
| CMS-PMt003g11185 | | STRG.18.1.p1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 2 | 2 | 3 | 1 | 2 | 2 | 2 |
| CMS-PMt010g15327 and CMS-PMt010g15548 | | STRG.31.1.p1 | 2 | 2 | 2 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CMS-PMt010g15739 | *orf265* | STRG.39.1.p1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CMS-PMt010g15740 | | STRG.39.1.p3 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

446
447

16

448    **Table 3** Oligonucleotide sequences of PCR primers

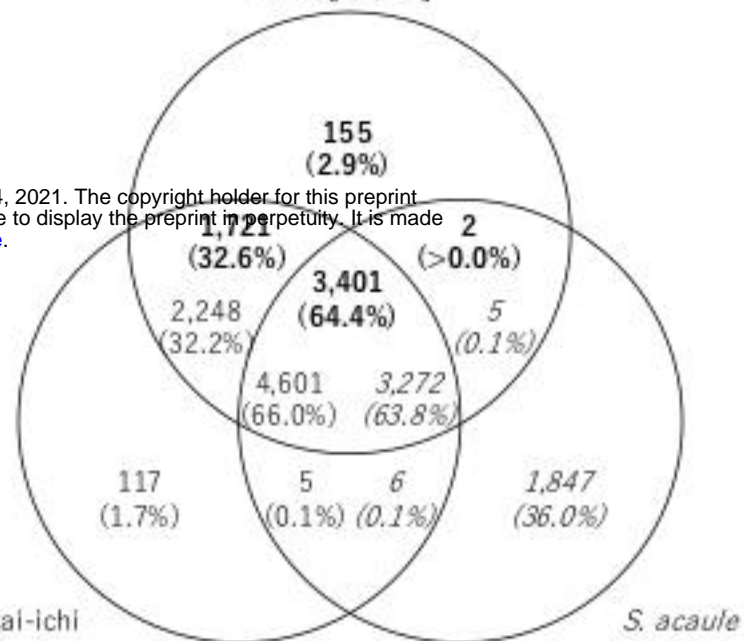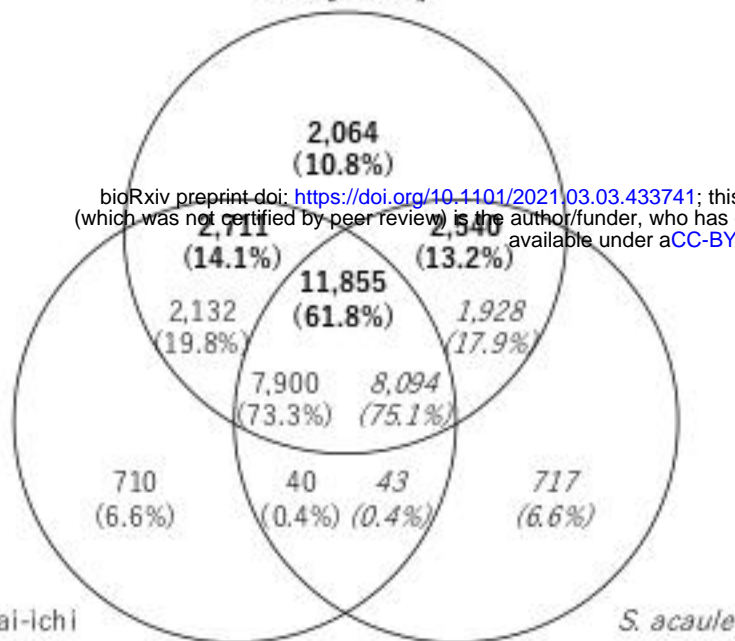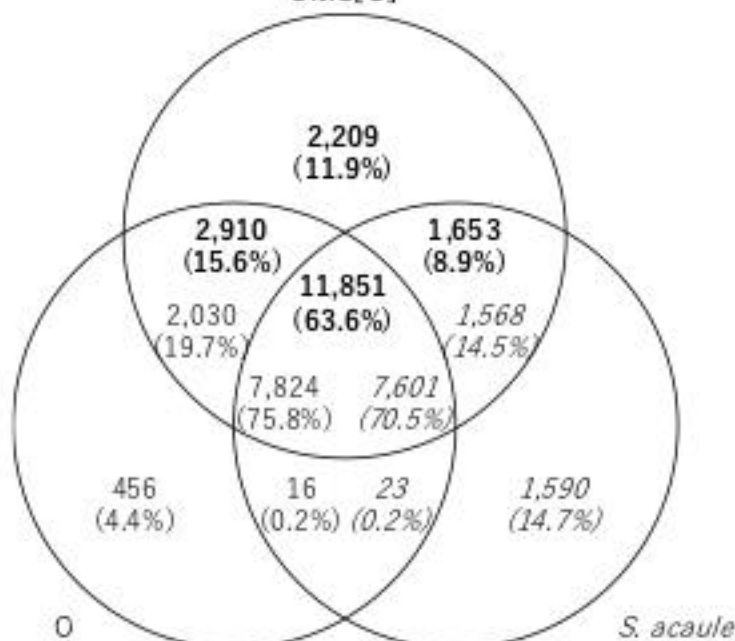| Target gene | Forward primer (5' - 3') | Reverse primer (5' - 3') |
| --- | --- | --- |
| *orf137* | CGATTGAGAAAGCGGCAGGC | GTTATTTTCGCTGCAACGGCG |
| *orf193* | GGGGAATCGGCCTTCTTTAGTC | GGGGAGGGTTTAATAAAGGAGCTG |
| *orf265* | CGGAGTGAAGCTGTATTGAGGG | GAGGAGAGGAACGAAGAACGAAAC |
| *orf265-rps12* | CGGAGTGAAGCTGTATTGAGGG | GATCCGGAATTCCCAGCAAATCC |
| *cox2* | CCCGCAAAGGATTGTTCATGG | CGTATAGGGCTCTTTGCTGGTAG |

449

*S. acaule*

Sekai-ichi

O

P

MSA1

CF[O]

CF[P]

**CMS[MSA1]**

**CMS[O]**

**CMS[P]**

Micro-Tom

**Dwarf CMS[P]**

A                    Mitochondrion                    B                    Chloroplast

CMS[MSA1]    Sekai-ichi          S. acaule              CMS[MSA1]    Sekai-ichi          S. acaule

CMS[O]       O                   S. acaule             CMS[O]       O                   S. acaule

CMS[P]       P                   S. acaule             CMS[P]       P                   S. acaule

**CMS[MSA1]**

2,064
(10.8%)

2,711
(14.1%)

2,546
(13.2%)

11,855
(61.8%)

2,132
(19.8%)

1,928
(17.9%)

7,900
(73.3%)

8,094
(75.1%)

710
(6.6%)

40
(0.4%)

43
(0.4%)

717
(6.6%)

Sekai-ichi

S. acaule

**CMS[MSA1]**

155
(2.9%)

1,721
(32.6%)

2
(>0.0%)

3,401
(64.4%)

2,248
(32.2%)

5
(0.1%)

4,601
(66.0%)

3,272
(63.8%)

117
(1.7%)

5
(0.1%)

6
(0.1%)

1,847
(36.0%)

Sekai-ichi

S. acaule

**CMS[O]**

2,209
(11.9%)

2,910
(15.6%)

1,653
(8.9%)

11,851
(63.6%)

2,030
(19.7%)

1,568
(14.5%)

7,824
(75.8%)

7,601
(70.5%)

456
(4.4%)

16
(0.2%)

23
(0.2%)

1,590
(14.7%)

O

S. acaule

**CMS[O]**

63
(1.2%)

1,807
(33.1%)

1
(>0.0%)

3,585
(65.7%)

1,982
(31.6%)

3
(0.1%)

3,986
(63.6%)

3,275
(63.8%)

269
(4.7%)

3
(>0.0%)

7
(0.1%)

1,845
(36.0%)

O

S. acaule

**CMS[P]**

1,501
(9.4%)

2,702
(17.0%)

1,516
(9.5%)

10,193
(64.1%)

1,969
(18.5%)

1,314
(12.2%)

7,942
(74.6%)

7,592
(70.4%)

714
(6.7%)

28
(0.3%)

30
(0.3%)

1,846
(17.1%)

P

S. acaule

**CMS[P]**

478
(5.9%)

2,547
(31.2%)

8
(0.1%)

5,132
(62.9%)

1,993
(32.2%)

27
(0.5%)

4,028
(65.1%)

3,275
(63.8%)

167
(2.7%)

4
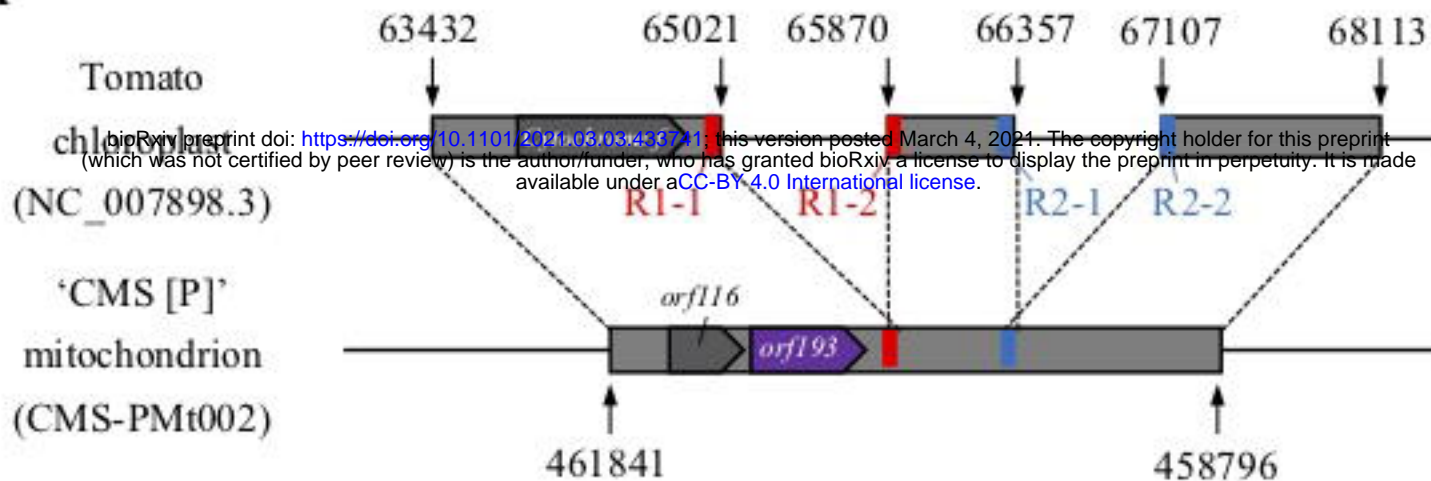(0.1%)

12
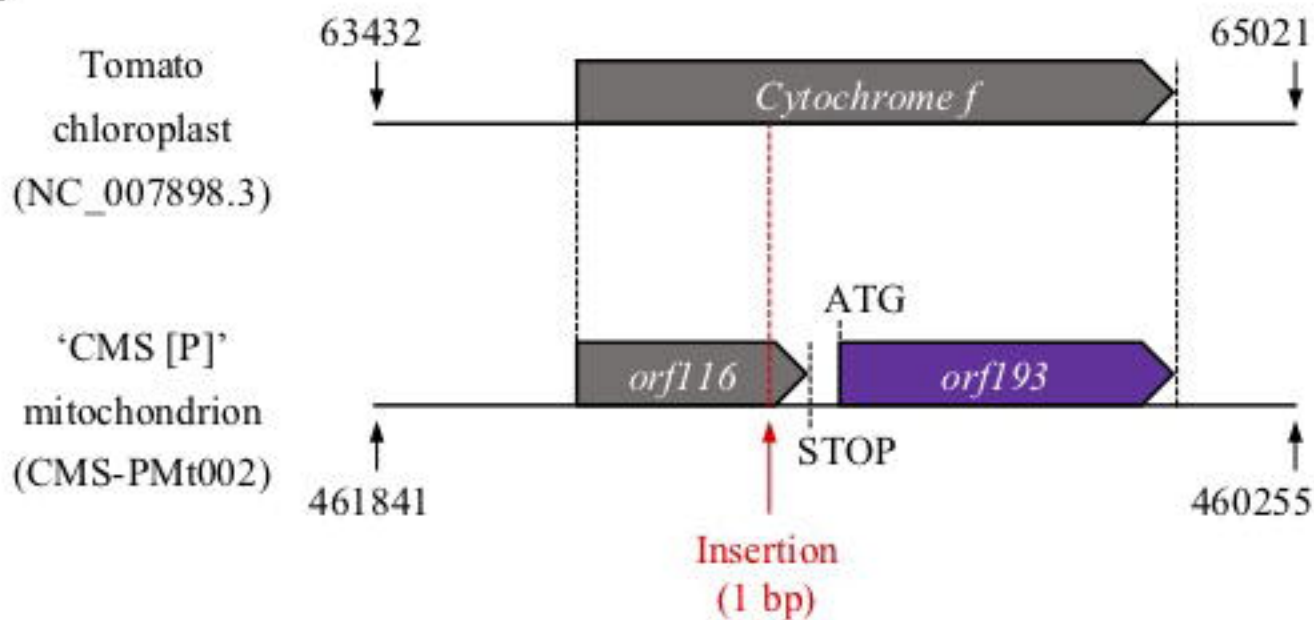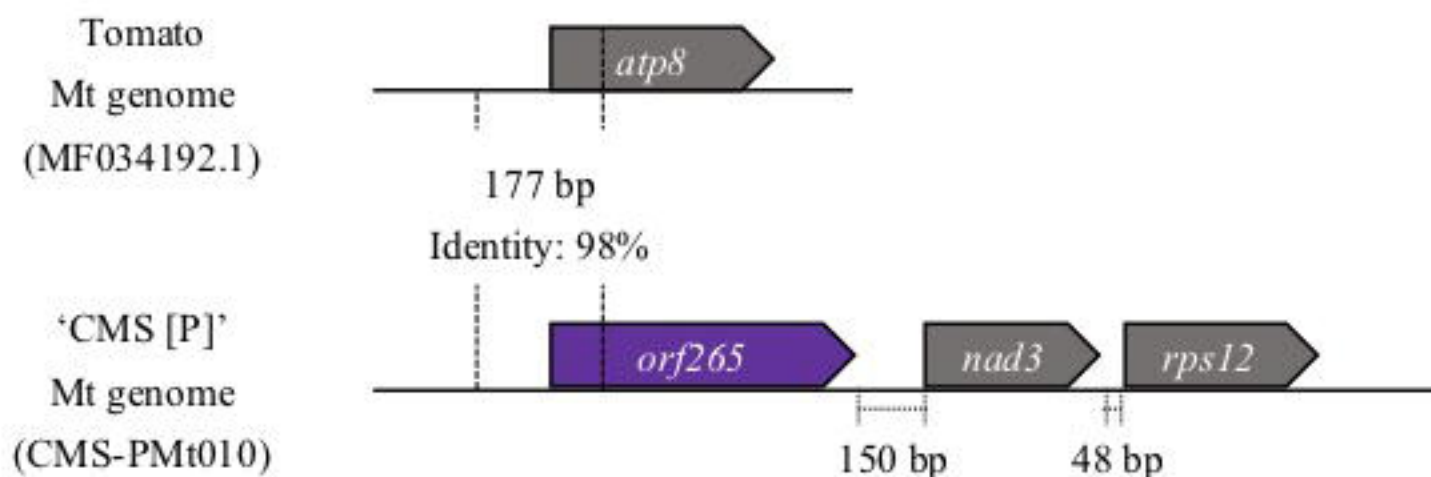(0.2%)

1,816
(35.4%)

P

S. acaule

# Mitochondrion

# Chloroplast

**CMS[MSA1] Mt genome**

**CMS[MSA1] Cp genome**

**CMS[O] Mt genome**

**CMS[O] Cp genome**

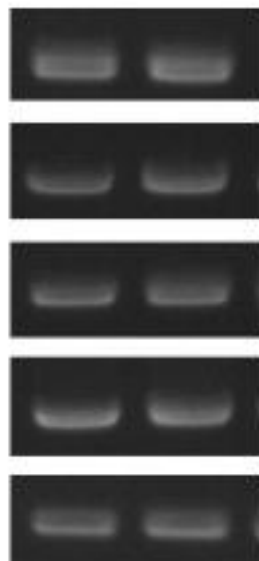**CMS[P] Mt genome**

**CMS[P] Cp genome**

**A**

Tomato chloroplast (NC_007898.3)

63432    65021    65870    66357    67107    68113

R1-1    R1-2    R2-1    R2-2

'CMS [P]' mitochondrion (CMS-PMt002)

orf116    orf193

461841    458796

**B**

(R1-1) -------TAACGGGATTCCC
(R1-2) CAA - GGGATTCCC-------

(R2-1) -------TCTTTTTTTTTG
(R2-2) TCTTTTTTTTTG-------

**C**

63432    65021

Tomato chloroplast (NC_007898.3)

Cytochrome f

'CMS [P]' mitochondrion (CMS-PMt002)

orf116    ATG    orf193    STOP

461841    460255

Insertion (1 bp)

**D**

Tomato Mt genome (MF034192.1)

atp8

177 bp

Identity: 98%

'CMS [P]' Mt genome (CMS-PMt010)

orf265    nad3    rps12

150 bp    48 bp

**A**

'CMS[P]' anther
'CMS[MSA1]' anther

**B**

'Dwarf CMS[P]'

Pollen (Dried)
Pollen (Incubated)
Anther
Leaf
Stem
Root
Ovary

*orf137*

*orf193*

*orf265*

*orf265 - rps12*

*cox2*