

ARTICLE

Open Access

# Organelle genome assembly uncovers the dynamic genome reorganization and cytoplasmic male sterility associated genes in tomato

Kosuke Kuwabara<sup>1</sup>, Issei Harada<sup>1</sup>, Yuma Matsuzawa<sup>2</sup>, Tohru Ariizumi<sup>1,3</sup>✉ and Kenta Shirasawa<sup>1,4</sup>✉

## Abstract

To identify cytoplasmic male sterility (CMS)-associated genes in tomato, we determined the genome sequences of mitochondria and chloroplasts in three CMS tomato lines derived from independent asymmetric cell fusions, their nuclear and cytoplasmic donors, and male fertile weedy cultivated tomato and wild relatives. The structures of the CMS mitochondrial genomes were highly divergent from those of the nuclear and cytoplasmic donors, and genes of the donors were mixed up in these genomes. On the other hand, the structures of CMS chloroplast genomes were moderately conserved across the donors, but CMS chloroplast genes were unexpectedly likely derived from the nuclear donors. Comparative analysis of the structures and contents of organelle genes and transcriptome analysis identified three genes that were uniquely present in the CMS lines, but not in the donor or fertile lines. RNA-sequencing analysis indicated that these three genes transcriptionally expressed in anther, and identified different RNA editing levels in one gene, *orf265*, that was partially similar to *ATP synthase subunit 8*, between fertile and sterile lines. The *orf265* was a highly potential candidate for CMS-associated gene. This study suggests that organelle reorganization mechanisms after cell fusion events differ between mitochondria and chloroplasts, and provides insight into the development of new F1 hybrid breeding programs employing the CMS system in tomato.

## Introduction

Cytoplasmic male sterility (CMS) is broadly found in the kingdom of Plantae<sup>1</sup>. CMS plants cannot produce seeds by self-pollination due to a lack of male fertility; therefore, pollen from other plants is always required for these plants to produce seeds. CMS is caused by the incompatibility of interactions of genetic information between nuclei and organelles, especially mitochondria<sup>1</sup>. The genes in nuclei and organelles are called *restorer of fertility (RF)* genes and CMS-associated genes, respectively. Therefore, CMS plants have been used as materials for studies of interactions between nuclear and cytoplasmic genes. Moreover, CMS is used in breeding programs to produce F1 hybrid

seeds<sup>1</sup>, in which cytoplasmic and pollen donors are employed as maternal and paternal parents, respectively.

CMS plants can be artificially generated by recurrent backcrossing or transgenic approaches<sup>2,3</sup>, which leads to incompatibility between nuclei and organelles. A tomato CMS line, called CMS-pennellii, which possesses nuclei and cytoplasm from *Solanum pennellii* and *Solanum peruvianum*, respectively, has been developed by recurrent backcrossing<sup>2</sup>. A gene knockdown strategy is also used to develop CMS tomato lines, for which expression of a nuclear gene that regulates mitochondrial substoichiometric shifting has been suppressed<sup>3</sup>. In addition, other types of CMS tomato lines have been generated via asymmetric cell fusion between cultivated tomato lines, namely, *Solanum lycopersicum* as the nuclear donor and a wild potato relative, *Solanum acaule*, as the cytoplasmic donor<sup>4</sup>. Among CMS lines, MSA1 has been well-studied to reveal nucleus-organelle incompatibility<sup>4</sup>. A physical map of the mitochondrial genome of MSA1 indicates that this

Correspondence: Tohru Ariizumi (ariizumi.toru.ge@u.tsukuba.ac.jp) or Kenta Shirasawa (shirasawa@kazusa.or.jp)

<sup>1</sup>Graduate School of Life and Environmental Sciences, University of Tsukuba, Tsukuba, Ibaraki 305-8577, Japan

<sup>2</sup>TOKITA Seed Co. LTD., Kazo, Saitama 349-1144, Japan

Full list of author information is available at the end of the article

© The Author(s) 2021



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

asymmetric cell fusion hybrid has a complex mitochondrial genome structure consisting of the parental genomes<sup>5</sup>. Transcripts of an open reading frame (ORF), *orf206*, of the hybrid mitochondrial genome are heterogeneously edited<sup>6</sup>. However, no candidates of CMS-associated genes have been identified in tomato.

Although CMS-associated gene sequences are not conserved across plant species, they have common features<sup>7</sup>. Most CMS-associated gene candidates usually possess transmembrane regions and chimeric structures, so-called fusion genes, of genes involved in respiration. Based on this information, CMS-associated genes have been identified in *Oryza sativa*<sup>8,9</sup>, *Helianthus annuus*<sup>10</sup>, and *Gossypium hirsutum*<sup>11</sup>. RNA-sequencing (RNA-Seq) based on next-generation sequencing technology has been employed to select candidates uniquely expressed in CMS lines of *Brassica juncea*<sup>12</sup>. The RNA-Seq technique is useful to detect base substitution, mainly C to U, in mitochondrial RNA, the so-called RNA editing, which is related to CMS<sup>13</sup>. Further functional studies are required to confirm that these candidates are involved in CMS. The introduction of *RF* genes into CMS lines would be a useful approach because CMS-associated genes can be down-regulated in the presence of *RF* genes<sup>7</sup>. Another approach is to introduce CMS-associated genes into fertile lines to induce sterility<sup>14</sup>. More recently, it has become possible to alter or edit gene sequences of mitochondrial genomes with TALEN technology<sup>15</sup>. This technology has been used to disrupt CMS-associated genes in mitochondrial genomes and thereby generate *Arabidopsis thaliana*, *Oryza sativa*, and *Brassica napus* with CMS<sup>15,16</sup>.

In parallel with MSA1, as shown in Fig. 1, two asymmetric cell fusions were developed between cultivated tomato lines *S. lycopersicum* ('O' and 'P') as nuclear donors and a wild potato relative, *S. acaule*, as the cytoplasmic donor<sup>17</sup>. The nuclear genome backgrounds of the three cell fusion lines including MSA1 were replaced with

the genomes of cultivated tomato lines by a repeated backcrossing strategy. The resultant CMS lines are designated 'CMS[MSA1]', 'CMS[O]', and 'CMS[P]'. Therefore, it may be possible to identify CMS-associated genes by comparative analysis of the genomes and transcriptomes of the CMS lines and their nuclear donors. In this study, we determined the sequences of the organelle genomes of the CMS lines and their donors. Subsequently, the genome sequences and gene expression patterns were compared to identify CMS-associated gene candidates. Furthermore, the results of this analysis may provide insights into the cytoplasmic genome features of asymmetric cell fusions.

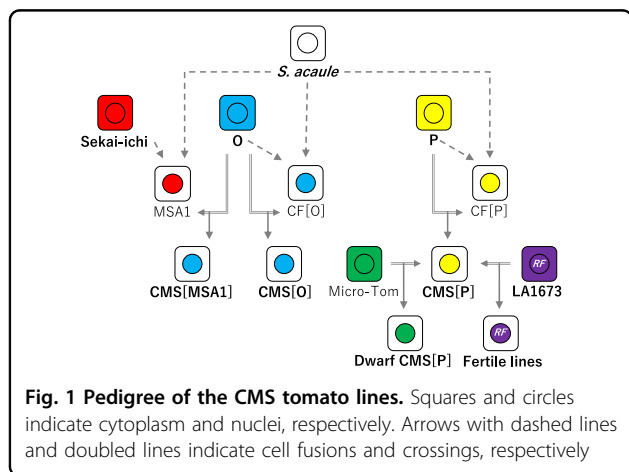
## Results

### De novo assembly of chloroplast and mitochondrial genomes

A total of 10.5 Gb reads per sample were obtained from three CMS tomato lines ('CMS[MSA1]', 'CMS[O]', and 'CMS[P]'), three nuclear donors ('Sekai-ichi', 'O', and 'P'), and one cytoplasmic donor (*S. acaule*). Of them, 374 Mb (3.6%) and 566 Mb (5.4%) of reads per sample were aligned on publicly available sequences of mitochondrial and chloroplast genomes, respectively. The reads mapped on the two sets of reference sequences were separately assembled into contig sequences.

Mitochondrial genome sequences were constructed with reads mapped on the mitochondrial reference sequences (Table 1). The mitochondrial genomes of the nuclear donors 'Sekai-ichi', 'O', and 'P' were all constructed from only contigs with assembly sizes of 562.6 kb ( $n = 2$ ,  $n$  represents contig numbers), 536.9 kb ( $n = 2$ ), and 553.3 kb ( $n = 2$ ), respectively. In *S. acaule*, 728.4 kb contigs ( $n = 7$ ) for the mitochondrial genome were established. The assembly sizes were longer in the CMS lines than in the nuclear and cytoplasmic donors, specifically, they were 995.2 kb ( $n = 7$ ) in 'CMS[MSA1]', 968.4 kb ( $n = 7$ ) in 'CMS[O]', and 829.3 kb ( $n = 5$ ) in 'CMS[P]'. For chloroplast genomes, total sequence lengths 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', respectively (Table 1). There were two contig sequences in each of the three nuclear donors. The assembly sizes were shorter in 'CMS[MSA1]' (296.6 kb,  $n = 1$ ) and 'CMS[O]' (307.1 kb,  $n = 1$ ) than in the nuclear donors, but longer in 'CMS[P]' (454.1 kb,  $n = 3$ ).

Comparative genome analysis revealed that the mitochondrial genomes of the CMS lines consisted of highly fragmented, repeated, and duplicated sequences derived from both donors throughout the genome (Fig. 2). On the other hand, the structures of the chloroplast genomes of the CMS lines were moderately conserved across the nuclear and cytoplasmic donors (Fig. 2).



**Fig. 1** Pedigree of the CMS tomato lines. Squares and circles indicate cytoplasm and nuclei, respectively. Arrows with dashed lines and doubled lines indicate cell fusions and crossings, respectively

**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	<i>S. acaule</i>	LA1670	LA1673
<i>Mitochondrion</i>									
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
<i>Chloroplast</i>									
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	5279	5456	8165	6971	6267	6192	5130	5346	5995

In parallel, we determined the mitochondrial and chloroplast genome sequences of *Solanum pimpinellifolium* LA1670 and *S. lycopersicum* var. *cerasiforme* LA1673 (Table 1). Sequence reads were obtained from a public DNA database and processed as described above. Assembly sizes of the mitochondrial and chloroplast genomes were 620.6 kb ( $n = 3$ ) and 299.4 kb ( $n = 1$ ) for *S. pimpinellifolium* LA1670, respectively, and 569.9 kb ( $n = 2$ ) and 337.7 kb ( $n = 2$ ) for *S. lycopersicum* var. *cerasiforme* LA1673, respectively.

**Gene prediction from the organelle genomes**

ORFs encoding  $\geq 25$  amino acids were extracted from the assembled sequences to predict potential genes. The number of potential genes predicted from the chloroplast genome assemblies ranged from 5,130 (*S. acaule*) to 8,165 ('CMS[P]') and the number of potential genes predicted from the mitochondrial sequences ranged from 10,326 ('O') to 19,170 ('CMS[MSA1]') (Table 1).

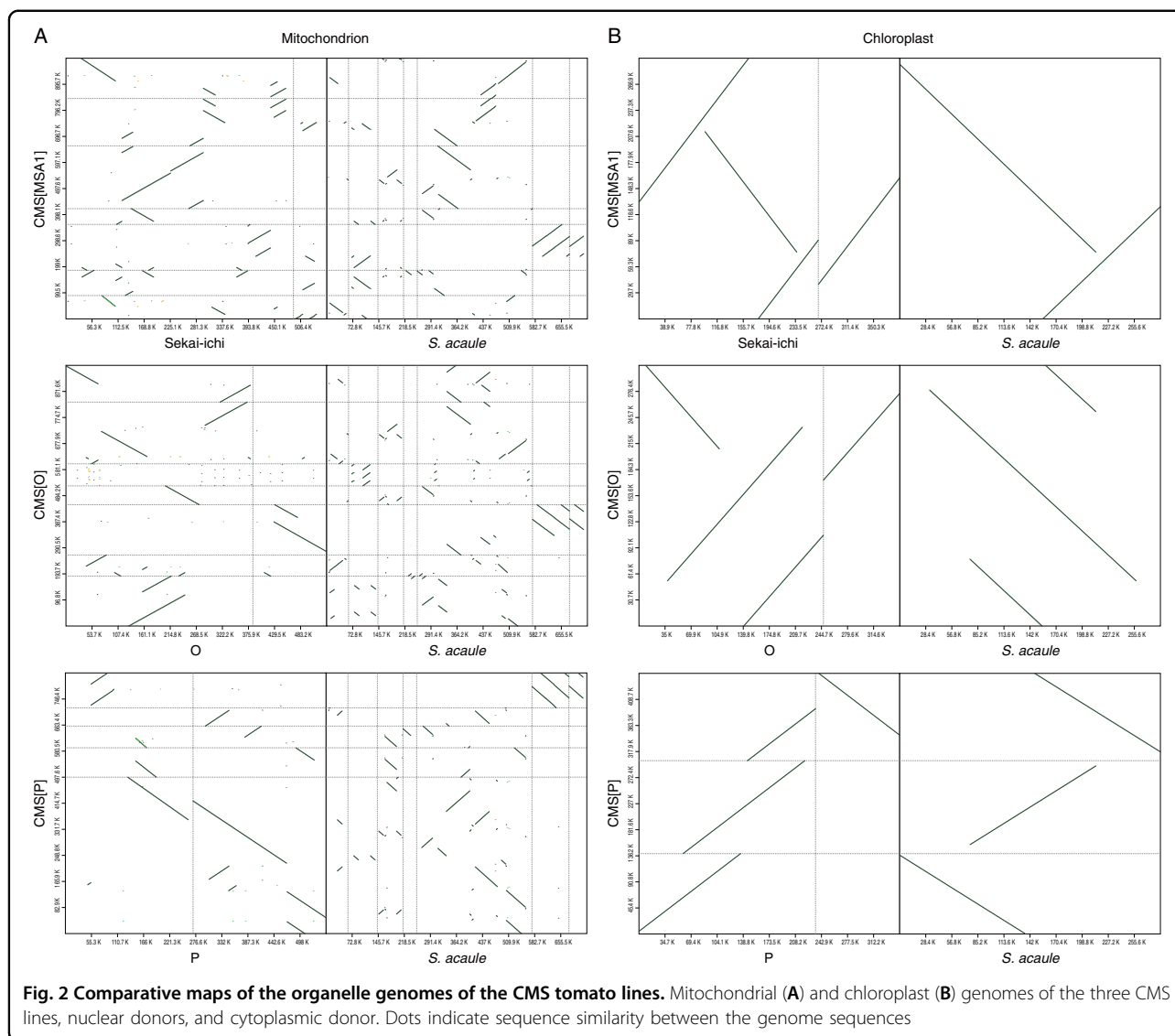
The ORFs were clustered to identify genes unique to and shared among the CMS lines, nuclear donors, and cytoplasmic donor (Fig. 3). The ORFs in the CMS mitochondrial genomes consisted of four types of genes, namely, those unique to the CMS lines (Type 1: 9.4–11.9%), those shared with the nuclear donors only (Type 2: 14.1–17.0%), those shared with the cytoplasmic donor only (Type 3: 8.9–13.2%), and those shared with both the nuclear and cytoplasmic donors (Type 4: 61.8–64.1%). By contrast, the ORFs in the CMS chloroplast genomes mostly consisted of three types of genes, namely, those unique to the CMS lines (Type 1: 1.2–5.9%), those shared with the nuclear donors only (Type 2: 31.2–33.1%), and those shared with both the

nuclear and cytoplasmic donors (Type 4: 62.9–65.7%). Few genes shared with the cytoplasmic donor only were found (Type 3: up to 0.1%).

The genome positions of the genes differed according to the gene type and organelle (Fig. 4). Type 4 genes in mitochondria were distributed across the genome with some gaps. The positions of Type 2 genes were basically the same as those of Type 4 genes, while Type 3 genes were located in the gaps between Type 4 genes. Type 1 genes were also located in the gaps and at the ends of contig sequences. On the other hand, in chloroplast genomes, the positions of Type 2 and 4 genes overlapped and Type 1 genes were located at the ends of contigs.

**Screening of CMS-associated gene candidates**

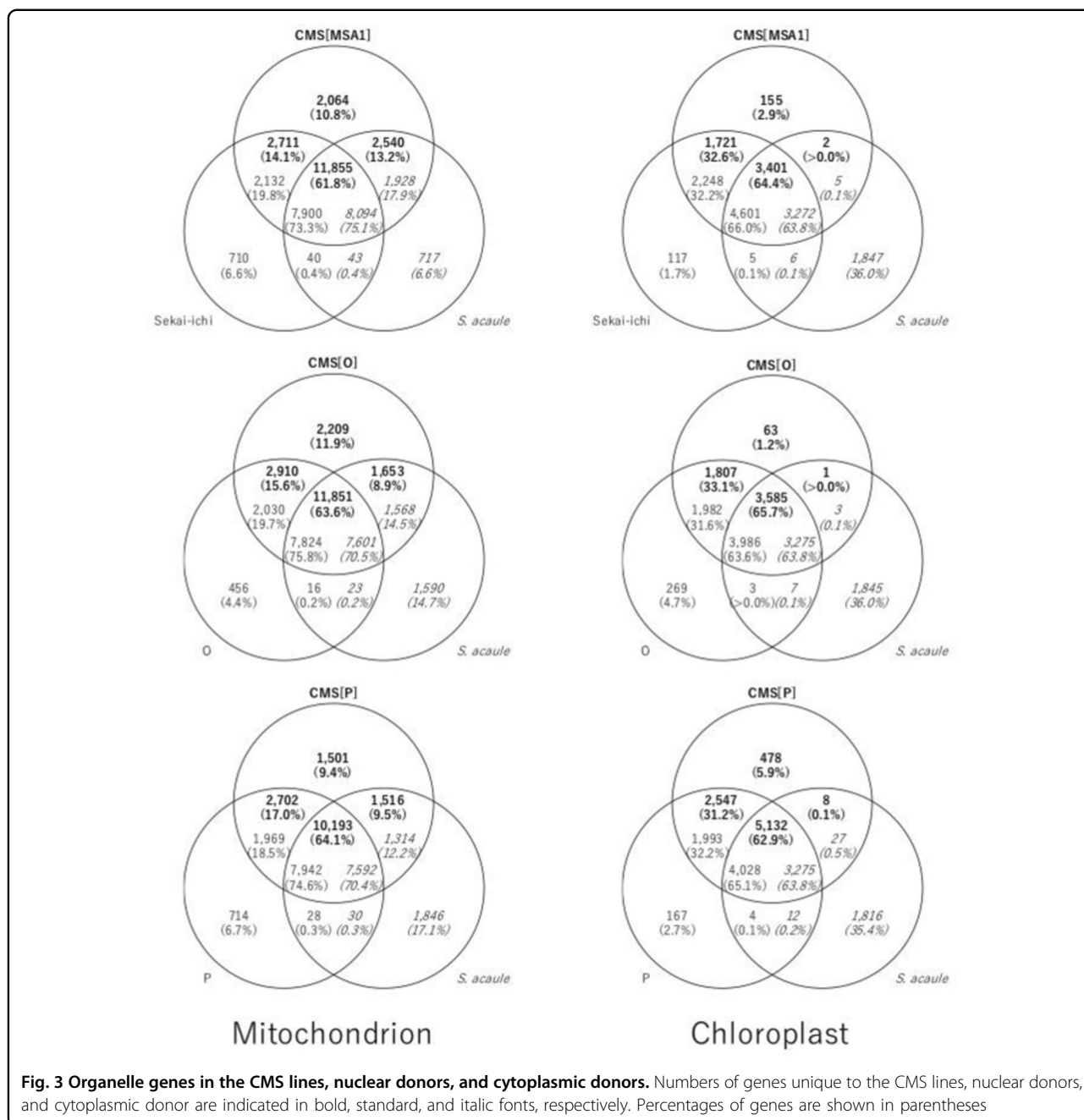
To identify candidates of CMS-associated genes in the mitochondrial genomes, we set the following four criteria: (1) amino acid length  $\geq 70$ , (2) absent from male fertile lines, (3) present in all three CMS lines, and (4) expressed in anthers of the CMS lines. Among the predicted genes in the 'CMS[P]', 'CMS[MSA1]', and 'CMS[O]' mitochondrial genomes, 831, 1025, and 969 genes encoded  $\geq 70$  amino acids, respectively. The gene sequences from the CMS lines were compared with the mitochondrial genomes of the nuclear donors ('Sekai-ichi', 'P', and 'O') and *S. pimpinellifolium* LA1670, *S. lycopersicum* var. *cerasiforme* LA1673), *S. pennellii*, and *Nicotiana tabacum*. In total, 183, 272, and 140 genes were selected because they were absent from the nuclear donors and Solanaceae relatives, all of which possess male fertility. Furthermore, we selected 36, 41, and 33 genes commonly present in the CMS lines. The copy numbers of the genes varied.



Finally, RNA-Seq reads were mapped on the mitochondrial genomes of the CMS lines. This analysis limited the number of CMS-associated gene candidates to four, including two identical sequences. The three genes were named *orf137* (two copies in the genome of each CMS line: CMS-PMt002g07240 and CMS-PMt005g13392), *orf193* (one copy: CMS-PMt002g06465), and *orf265* (one copy: CMS-PMt010g15739). Among these genes, four RNA editing sites, where C was substituted with U, were found in only *orf265* at 60,019th, 60,030th, 60,038th, and 60,047th positions of the contig of CMS-PMt010, which were corresponded to the positions of 58th, 47th, 39th, and 30th positions from the initial codon of ATG, respectively. While two substitutions at 39th and 30th positions were silent mutations, those at 58th and 47th induced non-synonymous substitutions, leucine to phenylalanine (L20F) and serine to leucine (S16L).

De novo transcriptome assembly was performed in parallel. RNA-Seq data were obtained from the anthers of 'P' and 'CMS[P]', and assembled into 62 and 43 transcript sequences, respectively, of which 37 'P' and 18 'CMS[P]' transcripts were predicted to have transmembrane domains. Of these sequences, eight were uniquely detected in 'CMS[P]'. Two genes (STRG.32.1.p1 and STRG.39.1.p1) were identical to *orf137* and *orf265*.

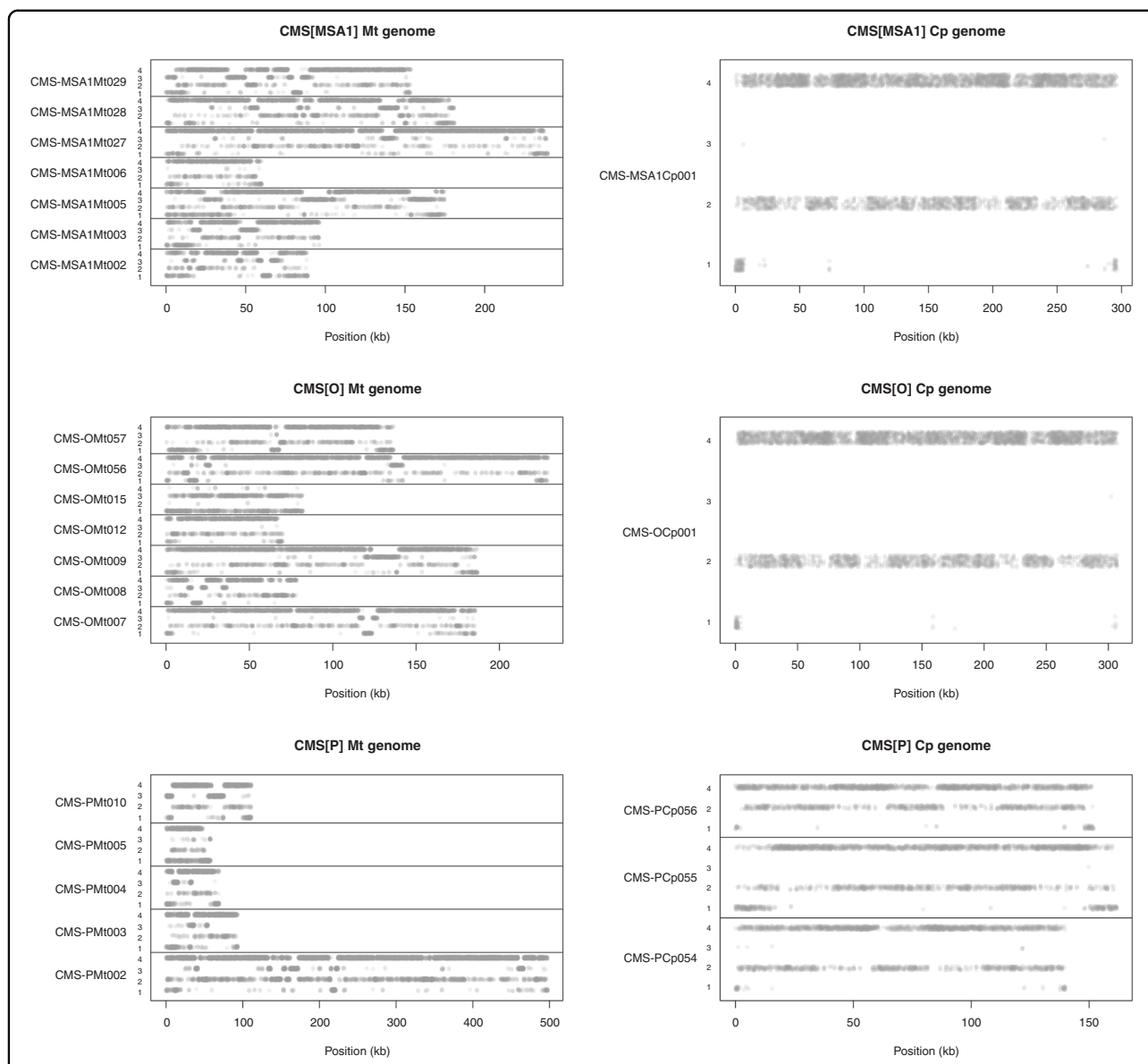
Because two genes were commonly identified in both analyses, a total of nine genes were finally selected as candidates of CMS-associated genes (Table 2). Sequence similarity searches with the mitochondrial and chloroplast genomes indicated that two copies of the STRG.32.1.p1 (*orf137*) sequence (CMS-PMt002g07240 and CMS-PMt005g13392) were present in the mitochondrial genomes of the three CMS lines. A single copy sequence of *orf193* (CMS-PMt002g06465) and a single copy sequence



of STRG.39.1.p1 (*orf265*, CMS-PMt010g15739) were found in the mitochondrial genomes of the three CMS lines in addition to that of *S. acaule*. The presence of the three genes in the CMS lines was validated by a PCR assay with the three CMS lines and six fertile lines. The remaining six genes were found in both the CMS and fertile lines. We selected three genes, *orf137*, *orf193*, and *orf265*, as highly potential candidates for CMS-associated genes due to their presence specifically in the CMS mitochondrial genomes and their expression in anthers.

**Sequence similarity analysis of the candidate genes**

The sequence similarity of the candidate genes including their flanking genome regions in the mitochondrial genome of ‘CMS[P]’ was investigated. A 3,045 bp genome sequence around *orf193* showed high sequence similarity to a 4,682 bp region of the tomato chloroplast genome sequences. The 3,045 bp sequence was split into three sequences containing 1,590, 488, and 1,007 bp (Fig. 5A) with highly conserved boundary sequences (Fig. 5B). In the 1,590 bp chloroplast genome sequence, a gene



**Fig. 4 Distributions of CMS tomato genes across the organelle genomes.** Dots indicate gene positions on contig sequences of the organelle genomes. Genes are grouped into the following four types: Type 1, genes unique to the CMS lines; Type 2, genes shared with the nuclear donors; Type 3, genes shared with the cytoplasmic donor; and Type 4, genes shared with both the nuclear and cytoplasmic donors

encoding *cytochrome f* was encoded; however, the corresponding sequence in the mitochondrial genome had a single base insertion causing a frame-shift mutation (Fig. 5C). This mutation broke the ORF of the *cytochrome f* gene and generated two small ORFs, *orf116* and *orf193*.

A portion of *orf265* and its upstream sequences (177 bp in total) showed high similarity to the *ATP synthase subunit 8 (atp8)* gene encoded in the tomato mitochondrial genome (Fig. 5D). The remaining sequences of *orf265* lacked similarity to reported sequences. *orf265* was located upstream of the *nad3* and *rps12* genes in the

mitochondrial genome. No sequence similarity was observed for *orf137* and the flanking sequence.

**Expression analysis of the candidate genes**

The expression patterns of the candidate genes, *orf137*, *orf193*, and *orf265*, were investigated by RT-PCR. First, we validated the results of the transcriptome analysis by detecting the expression of the three genes in anthers of ‘CMS[P]’ and ‘CMS[MSA1]’ (Fig. 6A). *orf265* was tandemly arrayed with *nad3* and *rps12*; therefore, we assumed that these three genes were co-transcribed as an

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

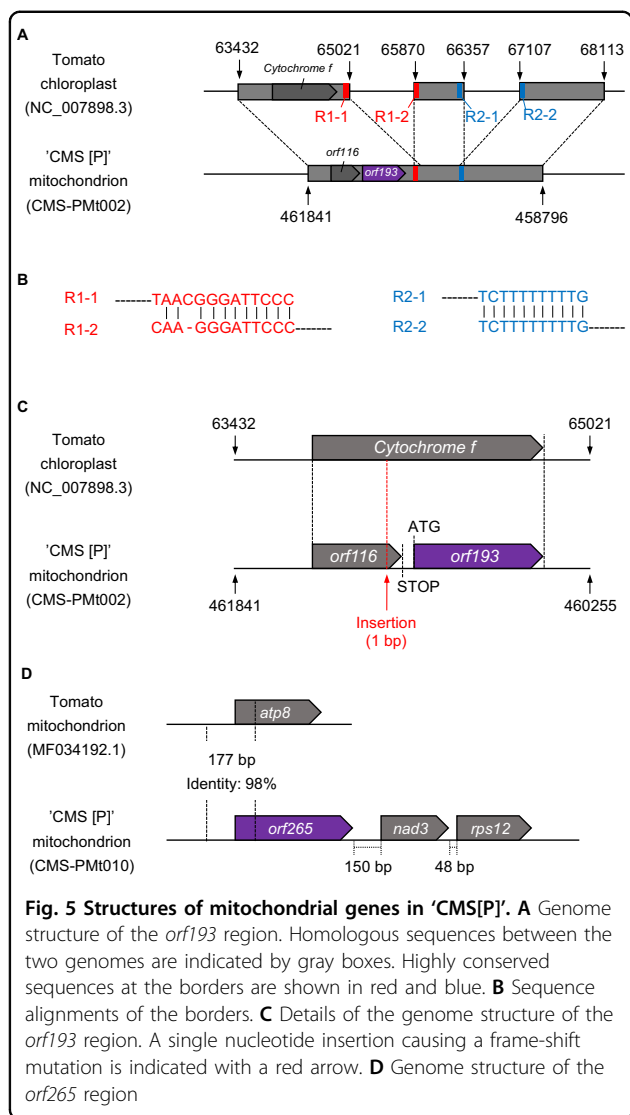
Gene ID of CMS[PF]	Candidates from genome analysis	Candidates from transcriptome analysis	Copy number in mitochondrial genome						Copy number in chloroplast genome							
			CMS [MSA1]			S. <i>acaule</i>			CMS [P]			S. <i>acaule</i>				
			CMS [MSA1]	CMS [O]	CMS [P]	O	P	S. <i>acaule</i>	CMS [P]	CMS [O]	CMS [P]	O	P	S. <i>acaule</i>		
CMS-PMt002g06465	<i>orf193</i>		1	1	1	0	0	1	0	0	0	0	0	0	0	0
CMS-PMt002g07240 and CMS-PMt005g13392	<i>orf137</i>	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	1	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	0	1	3	2	1	2	1
CMS-PMt003g11185		STRG.18.1.p1	0	0	1	0	0	0	0	0	2	3	1	2	2	2
CMS-PMt010g15327 and CMS-PMt010g15548		STRG.31.1.p1	2	2	2	0	0	1	1	1	0	0	0	0	0	0
CMS-PMt010g15739	<i>orf265</i>	STRG.39.1.p1	1	1	1	0	0	0	0	0	0	0	0	0	0	0
CMS-PMt010g15740		STRG.39.1.p3	1	1	1	1	1	1	1	1	0	0	0	0	0	0

operon. As expected, transcripts spanning the three genes were also detected (Fig. 6A). Next, we analyzed gene expression in leaves, stems, roots, ovaries, and pollen in addition to anthers of Dwarf ‘CMS[P]’ which was a BC3 generation of ‘CMS[P]’ backcrossed with a tomato dwarf cultivar ‘Micro-Tom’. Expression of *orf137* and *orf265* was detected in all tested tissues, while that of *orf193* was observed in leaves, stems, roots, ovaries, and anthers (Fig. 6B).

In parallel, to quantify the RNA edited rate found in *orf265*, we sequenced the cDNA of three sterile plants of ‘CMS[P]’ and three fertile lines of F4 progenies obtained from a cross between ‘CMS[P]’ and a fertility-restoring line, *S. lycopersicum* var. *cerasiforme* LA1673. Each data point was covered with 32,850 RNA reads in average. RNA edits were observed at the positions of 60,019 (L20F), 60,030 (S16L), 60,038 (F13F), and 60,047 (F10F) in the cDNA samples of fertile and sterile lines in comparison with genome DNA as a negative control. The RNA edits at three positions 60,019 (L20F), 60,030 (S16L), and 60,047 (F10F) were significantly higher in fertile lines (59.8%, 65.4%, and 64.7%) than sterile lines (39.5%, 50.7%, and 43.9%) (Fig. 7A). The amino acid sequence translated from the edited RNA was different from the conserved sequenced in *Solanum* species including the unedited sequence (Fig. 7B).

**Discussion**

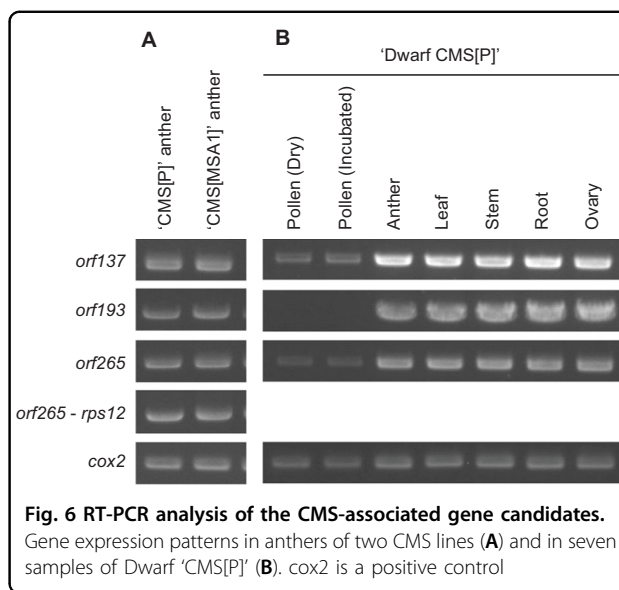
We determined the mitochondrial and chloroplast genome sequences of CMS lines derived from asymmetric cell fusions and those of their nuclear and cytoplasmic donors (Table 1). Comparative analysis of the structures unexpectedly revealed that the cytoplasmic genome structures of the fusions were rearranged and divergent from those of the cytoplasmic donor (*S. acaule*) and nuclear donors (*S. lycopersicum*) (Fig. 2). CMS-PMt003g09846 and CMS-PMt003g11185 were encoded in both the mitochondrial and chloroplast genomes (Table 2), suggesting that mitochondria and chloroplasts from the two donors were fused with each other and reorganized even though the cytoplasm of the nuclear donors was chemically inactivated to generate asymmetric cell fusions. Interestingly, the mitochondrial genomes of the CMS lines were larger than those of the donors, while the size of chloroplast genomes among the CMS lines was equivalent (Table 1). In addition, the structures of the mitochondrial genomes were divergent, while those of the chloroplast genomes were rather conserved (Fig. 2). Gene clustering analysis suggested that both the cytoplasmic and nuclear donors contributed to form mitochondria in the CMS lines (Fig. 3). Furthermore, the structures of the CMS mitochondrial genomes contained patches of the two genomes of the donors (Fig. 4). These results suggest that the mitochondrial genomes of both



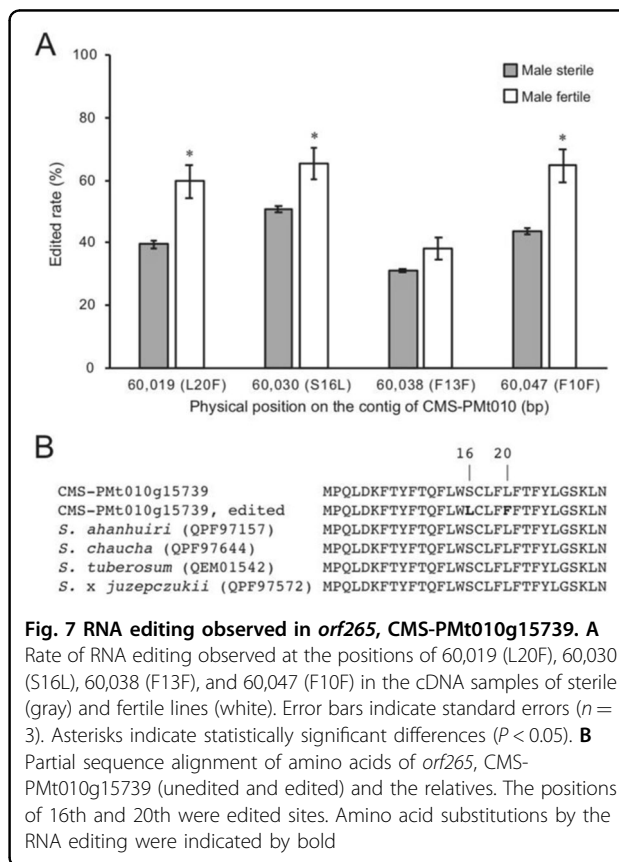
**Fig. 5 Structures of mitochondrial genes in 'CMS[P]'. A** Genome structure of the *orf193* region. Homologous sequences between the two genomes are indicated by gray boxes. Highly conserved sequences at the borders are shown in red and blue. **B** Sequence alignments of the borders. **C** Details of the genome structure of the *orf193* region. A single nucleotide insertion causing a frame-shift mutation is indicated with a red arrow. **D** Genome structure of the *orf265* region

donors were highly fragmented at the time of asymmetric cell fusion and reorganized to form a new mitochondrial genome<sup>5</sup>. This is completely different from our expectation that genomes of the cytoplasmic donors should be present in CMS lines derived from asymmetric cell fusions. More interestingly, chloroplasts of the CMS lines consisted only of genes from nuclear donors, not from the cytoplasmic donor. This unexpected finding has been frequently made in tomato<sup>18</sup>, tobacco<sup>19</sup>, and *Brassica*<sup>20</sup>. We speculate that interactions of genetic information between nuclei and organelles might be strict with chloroplasts rather than with mitochondria. The genome and/or organelle reorganization mechanisms after cell fusions might differ between mitochondria and chloroplasts.

Based on genome and transcriptome analyses, nine genes encoded in the mitochondrial genome of 'CMS [P]' were selected as candidate CMS-associated genes



**Fig. 6 RT-PCR analysis of the CMS-associated gene candidates.** Gene expression patterns in anthers of two CMS lines (A) and in seven samples of Dwarf 'CMS[P]' (B). *cox2* is a positive control



**Fig. 7 RNA editing observed in *orf265*, CMS-PMt010g15739. A** Rate of RNA editing observed at the positions of 60,019 (L20F), 60,030 (S16L), 60,038 (F13F), and 60,047 (F10F) in the cDNA samples of sterile (gray) and fertile lines (white). Error bars indicate standard errors ( $n = 3$ ). Asterisks indicate statistically significant differences ( $P < 0.05$ ). **B** Partial sequence alignment of amino acids of *orf265*, CMS-PMt010g15739 (unedited and edited) and the relatives. The positions of 16th and 20th were edited sites. Amino acid substitutions by the RNA editing were indicated by bold

(Table 2). Among them, three genes (*orf193*, STRG.32.1.p1 = *orf137* and STRG.39.1.p1 = *orf265*) were uniquely present in the genomes of the CMS lines and expressed in their anthers (Table 2 and Fig. 6). STRG.32.1.p1 (*orf137*) showed sequence similarity with the CMS-associated protein-encoding *cytochrome c subunit 1*. STRG.39.1.p1



(*orf265*) was similar to *ATP synthase subunit 8* at the N-terminus, but lacked similarity in the remaining regions (Fig. 5). CMS-associated genes are generally involved in cellular respiration producing energy to generate pollen<sup>21</sup>, and this is true of both these genes. In many cases, fusion genes have been reported to be CMS-associated genes, e.g., *orf307* in *Oryza sativa*<sup>22</sup> and *orf72* in *Brassica oleracea*<sup>23</sup>, and to produce cytotoxic proteins, which lead to male sterility. Knockout mutagenesis with mitoTALENs<sup>15,16</sup> targeting the candidate genes would be useful to identify CMS-associated genes and to generate CMS lines from normal tomato cultivars. However, these technologies are still limited in *Arabidopsis*, rice, and *Brassica*<sup>15,16</sup>. We have struggled to develop a conventional mitoTALEN method for the CMS tomato lines (data not shown).

Among the three candidates, *orf265* had a remarkable difference in the RNA edited rate between fertile and sterile lines (Fig. 7), in which pentatricopeptide repeat (PPR) proteins are well known to involve as *RF* gene for CMS<sup>24</sup>. Since the RNA editing was observed in both fertile and sterile lines, we hypothesized that the *RF* genes for the tomato CMS lines used in this study might be encoded in the genomes of CMS lines as well as *RF* lines, but the levels of expressions, enzyme activities, and/or affinity to the mitochondrial RNA might be different. The different levels of the edited RNA might control the fertility of the CMS tomato lines, even though RNA-processing and translation inhibition by PPR proteins as the *RF* gene should be still considered. Restorer genes for CMS lines have been identified in wild tomato relatives, e.g., *S. pimpinellifolium* LA1670 and *S. lycopersicum* var. *cerasiforme* LA1673<sup>17</sup>. Recently, we published the genome sequence data of these two wild relatives<sup>25</sup>. We expect *RF* genes for CMS lines to be discovered soon based on this information, although no candidate genes or genetic loci have been reported.

CMS lines are powerful tools to produce F1 hybrid seeds in breeding programs<sup>1</sup>. However, in cereals and fruits including tomato, the *RF* genes are essential for F1 plants to set seeds and bear fruits. Once CMS-associated genes and *RF* genes are identified, tomato F1 hybrid seeds can be produced by employing insect pollinators instead of the currently used hand-pollination systems. We propose that CMS-based F1 hybrid breeding programs with insect pollinators can be implemented in tomato breeding programs to reduce the costs of F1 seed production in the future.

## Materials and methods

### Plant materials

Three tomato CMS lines ('CMS[MSA1]', 'CMS[O]', and 'CMS[P]'), three cultivated tomato lines (*S. lycopersicum* 'Sekai-ichi', 'O', and 'P'), and one potato wild relative (*S. acaule*) were used (Fig. 1). 'CMS[MSA1]' was developed

by repeated backcrossing using 'O' as a recurrent parent and a male-sterile tomato, MSA1, as a cytoplasmic donor. MSA1 is an asymmetric cell fusion between the tomato cultivar Sekai-ichi (as the nuclear donor) and the potato wild relative *S. acaule* (as the cytoplasmic donor)<sup>4</sup>. 'CMS [O]' was a progeny in repeated backcrossing using 'O' as the paternal parent and an asymmetric cell fusion between 'O' (as the nuclear donor) and *S. acaule* (as the cytoplasmic donor). 'CMS[P]' was also a progeny in backcrossing using 'P' as the paternal parent and an asymmetric cell fusion between 'P' (as the nuclear donor) and *S. acaule* (as the cytoplasmic donor). Dwarf 'CMS[P]' was developed from 'CMS[P]' by backcrossing with *S. lycopersicum* 'Micro-Tom' (TOMJPF0001), which is a miniature dwarf cultivar<sup>26</sup>. From a cross between 'CMS [P]' and a fertility-restoring line, *S. lycopersicum* var. *cerasiforme* LA1673, we selected fertility restored lines of F4 progenies, which possessed CMS-associated genes and *RF* genes in cytoplasm and nuclei, respectively. The putative nuclear and cytoplasmic genomes of the materials are shown in Fig. 1.

### Genome sequence analysis

Total genomic DNA was extracted from young leaves of the six tomato lines ('CMS[MSA1]', 'CMS[O]', 'CMS[P]', 'Sekai-ichi', 'O', and 'P') and *S. acaule* with a Maxwell 16 Instrument and Maxwell 16 Tissue DNA Purification Kits (Promega, Madison, WI, USA). SMRT sequence libraries were constructed with an SMRTbell Express Template Prep Kit (PacBio, Menlo Park, CA, USA) and used for sequencing on a PacBio Sequel system (PacBio). Genome sequence data for *S. pimpinellifolium* LA1670 and *S. lycopersicum* var. *cerasiforme* LA1673 were obtained from a public DNA database (DRA accession numbers DRX231405 and DRX231409)<sup>25</sup>.

### Genome assembly and gene prediction

Sequence reads were mapped on reference genome sequences for mitochondria (GenBank accession numbers MF034192, MF034193, NC\_035964, and MF98995–MF989957) or chloroplasts (NC\_007898) with Organelle\_PBA<sup>27</sup>. Reads mapped on the reference sequences were assembled into contig sequences with Canu<sup>28</sup>. Potential sequence errors in the contig sequences were corrected twice with the sequence reads by Arrow (PacBio). The corrected contig sequences were aligned back to the reference sequences with Nucmer<sup>29</sup> to select highly confident organelle genomes. ORFs ( $\geq 75$  bases) in the organelle genomes were selected as potential genes with ORFfinder (<https://www.ncbi.nlm.nih.gov/orffinder>). The ORF sequences were clustered with CD-HIT<sup>30</sup>. Transmembrane domains in the gene sequences were predicted by TMHMM<sup>31</sup>. Sequence similarity searches with the

**Table 3** Oligonucleotide sequences of PCR primers

Target gene	Forward primer (5' - 3')	Reverse primer (5' - 3')
<i>orf137</i>	CGATTGAGAAAGCGGCAGGC	GTTATTTTCGCTGCAACGGCG
<i>orf193</i>	GGGGAATCGGCCTTCTTAGTC	GGGGAGGGTTAATAAAGGAGCTG
<i>orf265</i>	CGGAGTGAAGCTGTATTGAGGG	GAGGAGAGGAACGAAGAACGAAAC
<i>orf265-rps12</i>	CGGAGTGAAGCTGTATTGAGGG	GATCCGGAATCCAGCAAATCC
<i>cox2</i>	CCCGCAAAGGATTGTCATGG	CGTATAGGGCTCTTTGCTGGTAG

mitochondrial genomes of *S. pennellii* (NC\_035964) and *N. tabacum* (NC\_006581) were performed by BLAST<sup>32</sup> with a threshold E-value of 1e-50.

### RNA expression analysis

Total RNA was extracted from the anthers of 'P' and 'CMS[P]' with an RNeasy Plant Mini Kit (QIAGEN, Hilden, Germany). RNA was treated with RNase-free DNase (QIAGEN) and used for sequence library preparation with a TruSeq Stranded mRNA Library Prep Kit (Illumina, San Diego, CA, USA). The resultant libraries were sequenced on NextSeq500 (Illumina) in paired-end, 151 bp mode. After trimming adaptors and low-quality reads by Trim\_galore (<https://github.com/FelixKrueger/TrimGalore>) with option -q 30 -length 100 followed by fastp<sup>33</sup> with option -l 100, transcriptomes were de novo assembled by the HiSat2-Stringtie pipeline<sup>34</sup> and putative ORFs were searched for annotation by BLASTP<sup>32</sup> against the SWISS-PROT database<sup>35</sup>.

To validate the RNA-Seq results, RT-PCR was performed. In total, 800 ng of total RNA isolated from anthers of the CMS lines was converted into cDNA with ReverTra Ace (TOYOBO, Osaka, Japan) using a random primer (TAKARA BIO, Kusatsu, Japan). cDNA diluted 10-fold with water was used as a template for PCR. The PCR mixture (10 µL) contained 0.5 µL cDNA, 0.3 µM primers (Tables 3), 2× PCR buffer (TOYOBO), 400 µM dNTPs, and 1 U DNA polymerase (KOD FX Neo, TOYOBO). The thermal cycling conditions were as follows: initial denaturation at 94 °C for 3 min; 35 cycles of denaturation at 98 °C for 15 s, annealing at 60 °C for 30 s, and extension at 68 °C for 60 s; and a final extension at 68 °C for 3 min. PCR products were separated by electrophoresis in a 1% agarose gel with TAE buffer. Gels were stained with Midori Green Advance (NIPPON Genetics, Tokyo, Japan) to detect DNA bands under ultraviolet illumination.

The RT-PCR was applied to amplify *orf265* from anther RNAs of 'CMS[P]' ( $n = 3$ ) and the fertile F4 lines ( $n = 3$ ) obtained from a cross between 'CMS[P]' and a fertility-restoring line, *S. lycopersicum* var. *cerasiforme* LA1673. DNA amplicon from the genome DNA of 'CMS[P]' was used as a negative control. The PCR products were ligated

and indexed with TruSeq DNA CD Indexes (Illumina) and sequenced on MiSeq (Illumina) in paired-end, 151 bp mode. RNA reads were mapped on the mitochondrial genome of 'CMS[P]' with Bowtie2 (version 2.3.5.1)<sup>36</sup> and sequence variants due to RNA editing was counted with BCFtools (version 1.9)<sup>37</sup>. The statistical analysis was performed by t test.

### Acknowledgements

The authors are grateful to Prof. Kazuo Watanabe (University of Tsukuba, Japan) for providing the DNA material of *S. acaule*. We also thank all technical and administrative members in T-PIRC center in the University of Tsukuba and Kazusa DNA Research Institute. This work was supported by the Kazusa DNA Research Institute Foundation to K.S., the Project of the NARO Bio-oriented Technology Research Advancement Institution (Research Program on Development of Innovative Technology, Grant numbers: 30010A and ID21448196) to K.S., Y.M., and T.A., JSPS KAKENHI (Grant Numbers: 17H03761 and 21H02181) to T.A., and JSPS Research Fellowships for Young Researchers (Grant Number: 21J20479) to K.K. Seeds of Micro-Tom (TOMJPF0001) was provided from National BioResource Project Tomato (NBRP tomato).

### Author details

<sup>1</sup>Graduate School of Life and Environmental Sciences, University of Tsukuba, Tsukuba, Ibaraki 305-8577, Japan. <sup>2</sup>TOKITA Seed Co. LTD., Kazo, Saitama 349-1144, Japan. <sup>3</sup>Tsukuba Plant Innovation Research Center, Tsukuba, Ibaraki 305-8577, Japan. <sup>4</sup>Kazusa DNA Research Institute, Kisarazu, Chiba 292-0818, Japan

### Author contributions

T.A. and K.S. conceived and coordinated the project. Y.M. established the plant materials. K.K., I.H., and K.S. collected the data. K.K., I.H., T.A., and K.S. analyzed and interpreted the data. K.K. and K.S. wrote the manuscript with contributions from T.A. All authors read and approved the final manuscript.

### Data availability

The DDBJ accession numbers of the assembled sequences are LC613090-LC613141. Genome information is available at KaTomicsDB (<http://www.kazusa.or.jp/tomato>).

### Conflict of interest

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

Received: 3 March 2021 Revised: 21 June 2021 Accepted: 2 August 2021  
Published online: 01 December 2021

### References

- Bohra, A., Jha, U. C., Adhimoalam, P., Bisht, D. & Singh, N. P. Cytoplasmic male sterility (CMS) in hybrid breeding in field crops. *Plant Cell Rep.* **35**, 967–93 (2016).

2. Petrova, M. et al. Characterisation of a cytoplasmic male-sterile hybrid line between *Lycopersicon peruvianum* Mill. × *Lycopersicon pennellii* Corr. and its crosses with cultivated tomato. *Theor. Appl. Genet.* **98**, 825–830 (1999).
3. Sandhu, A. P., Abdelnoor, R. V. & Mackenzie, S. A. Transgenic induction of mitochondrial rearrangements for cytoplasmic male sterility in crop plants. *Proc. Natl Acad. Sci. USA* **104**, 1766–70 (2007).
4. Melchers, G., Mohri, Y., Watanabe, K., Wakabayashi, S. & Harada, K. One-step generation of cytoplasmic male sterility by fusion of mitochondrial-inactivated tomato protoplasts with nuclear-inactivated *Solanum* protoplasts. *Proc. Natl Acad. Sci. USA* **89**, 6832–6 (1992).
5. Shikanai, T., Kaneko, H., Nakata, S., Harada, K. & Watanabe, K. Mitochondrial genome structure of a cytoplasmic hybrid between tomato and wild potato. *Plant Cell Rep.* **17**, 832–836 (1998).
6. Shikanai, T., Nakata, S., Harada, K. & Watanabe, K. Analysis of the heterogeneous transcripts of the highly edited orf206 in tomato mitochondria. *Plant Cell Physiol.* **37**, 692–6 (1996).
7. Chen, L. & Liu, Y. G. Male sterility and fertility restoration in crops. *Annu Rev. Plant Biol.* **65**, 579–606 (2014).
8. Igarashi, K., Kazama, T., Motomura, K. & Toriyama, K. Whole genomic sequencing of RT98 mitochondria derived from *Oryza rufipogon* and northern blot analysis to uncover a cytoplasmic male sterility-associated gene. *Plant Cell Physiol.* **54**, 237–43 (2013).
9. Okazaki, M., Kazama, T., Murata, H., Motomura, K. & Toriyama, K. Whole mitochondrial genome sequencing and transcriptional analysis to uncover an RT102-type cytoplasmic male sterility-associated candidate Gene Derived from *Oryza rufipogon*. *Plant Cell Physiol.* **54**, 1560–8 (2013).
10. Makarenko, M. S. et al. Characterization of the mitochondrial genome of the MAX1 type of cytoplasmic male-sterile sunflower. *BMC Plant Biol.* **19**, 51 (2019).
11. Li, S. et al. The comparison of four mitochondrial genomes reveals cytoplasmic male sterility candidate genes in cotton. *BMC Genomics* **19**, 775 (2018).
12. Wu, Z. et al. Mitochondrial genome and transcriptome analysis of five alloplasmic male-sterile lines in *Brassica juncea*. *BMC Genomics* **20**, 348 (2019).
13. Wang, Z. et al. Cytoplasmic male sterility of rice with boro II cytoplasm is caused by a cytotoxic peptide and is restored by two related PPR motif genes via distinct modes of mRNA silencing. *Plant Cell* **18**, 676–87 (2006).
14. Yang, J., Liu, X., Yang, X. & Zhang, M. Mitochondrially-targeted expression of a cytoplasmic male sterility-associated orf220 gene causes male sterility in *Brassica juncea*. *BMC Plant Biol.* **10**, 231 (2010).
15. Kazama, T. et al. Curing cytoplasmic male sterility via TALEN-mediated mitochondrial genome editing. *Nat. Plants* **5**, 722–730 (2019).
16. Arimura, S. I. et al. Targeted gene disruption of ATP synthases 6-1 and 6-2 in the mitochondrial genome of *Arabidopsis thaliana* by mitoTALENs. *Plant J.* **104**, 1459–1471 (2020).
17. Harada, K., Watabe, K. & Kondo, K. A method for producing a tomato plant with restored fertility. Japanese Patent No. 3386292. (2003).
18. Bonnema, A. B., Melzer, J. M. & O'Connell, M. A. Tomato cybrids with mitochondrial DNA from *Lycopersicon pennellii*. *Theor. Appl. Genet.* **81**, 339–48 (1991).
19. Sidorov, V. A., Menczel, L., Nagy, F. & Maliga, P. Chloroplast transfer in *Nicotiana* based on metabolic complementation between irradiated and iodoacetate treated protoplasts. *Planta* **152**, 341–5, <https://doi.org/10.1007/BF00388259> (1981).
20. Morgan, A. & Maliga, P. Rapid chloroplast segregation and recombination of mitochondrial DNA in *Brassica* cybrids. *Mol. Gen. Genet.* **209**, 240–6 (1987).
21. Touzet, P. & Meyer, E. H. Cytoplasmic male sterility and mitochondrial metabolism in plants. *Mitochondrion* **19 Pt B**, 166–71 (2014).
22. Fujii, S., Kazama, T., Yamada, M. & Toriyama, K. Discovery of global genomic reorganization based on comparison of two newly sequenced rice mitochondrial genomes with cytoplasmic male sterility-related genes. *BMC Genomics* **11**, 209 (2010).
23. Shinada, T., Kikuchi, Y., Fujimoto, R. & Kishitani, S. An alloplasmic male-sterile line of *Brassica oleracea* harboring the mitochondria from *Diplotaxis muralis* expresses a novel chimeric open reading frame, orf72. *Plant Cell Physiol.* **47**, 549–53 (2006).
24. Fujii, S. & Small, I. The evolution of RNA editing and pentatricopeptide repeat genes. *N. Phytol.* **191**, 37–47 (2011).
25. Takei, H. et al. De novo genome assembly of two tomato ancestors, *Solanum pimpinellifolium* and *S. lycopersicum* var. *cerasiforme*, by long-read sequencing. *DNA Res.* <https://doi.org/10.1093/dnares/dsaa029> (2021).
26. Scott, J.W. & Harbaugh, B.K. *Micro-tom: A Miniature Dwarf Tomato* (Agricultural Experiment Station, Institute of Food and Agricultural Sciences, University of Florida, 1989).
27. Soorni, A., Haak, D., Zaitlin, D. & Bombarely, A. Organelle\_PBA, a pipeline for assembling chloroplast and mitochondrial genomes from PacBio DNA sequencing data. *BMC Genomics* **18**, 49 (2017).
28. Koren, S. et al. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
29. Kurtz, S. et al. Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004).
30. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* **28**, 3150–2 (2012).
31. Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E. L. Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J. Mol. Biol.* **305**, 567–80 (2001).
32. Altschul, S. F. et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–402 (1997).
33. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ pre-processor. *Bioinformatics* **34**, i884–i890 (2018).
34. Pertea, M., Kim, D., Pertea, G. M., Leek, J. T. & Salzberg, S. L. Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.* **11**, 1650–67 (2016).
35. Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res.* **28**, 45–8 (2000).
36. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–9 (2012).
37. Danecek, P. et al. Twelve years of SAMtools and BCFtools. *Gigascience* **10**, <https://doi.org/10.1093/gigascience/giab008> (2021).