**BMC Genomics**

# Comparative chloroplast genomes: insights into the evolution of the chloroplast genome of *Camellia sinensis* and the phylogeny of *Camellia*

Li Li[1*†] , Yunfei Hu[1†], Min He[1†], Bo Zhang[1], Wei Wu[2], Pumo Cai[1], Da Huo[1] and Yongcong Hong[1*]

## Abstract

**Background:** Chloroplast genome resources can provide useful information for the evolution of plant species. Tea plant (*Camellia sinensis*) is among the most economically valuable member of *Camellia*. Here, we determined the chloroplast genome of the first natural triploid Chinary type tea ('Wuyi narcissus' cultivar of *Camellia sinensis var. sinensis*, CWN) and conducted the genome comparison with the diploid Chinary type tea (*Camellia sinensis var. sinensis*, CSS) and two types of diploid Assamica type teas (*Camellia sinensis var. assamica*: Chinese Assamica type tea, CSA and Indian Assamica type tea, CIA). Further, the evolutionary mechanism of the chloroplast genome of *Camellia sinensis* and the relationships of *Camellia* species based on chloroplast genome were discussed.

**Results:** Comparative analysis showed the evolutionary dynamics of chloroplast genome of *Camellia sinensis* were the repeats and insertion-deletions (indels), and distribution of the repeats, indels and substitutions were significantly correlated. Chinese tea and Indian tea had significant differences in the structural characteristic and the codon usage of the chloroplast genome. Analysis of sequence characterized amplified region (SCAR) using sequences of the intergenic spacers (trnE/trnT) showed none of 292 different *Camellia sinensis* cultivars had similar sequence characteristic to triploid CWN, but the other four *Camellia* species did. Estimations of the divergence time showed that CIA diverged from the common ancestor of two Assamica type teas about 6.2 Mya (CI: 4.4–8.1 Mya). CSS and CSA diverged to each other about 0.8 Mya (CI: 0.4–1.5 Mya). Moreover, phylogenetic clustering was not exactly consistent with the current taxonomy of *Camellia*.

**Conclusions:** The repeat-induced and indel-induced mutations were two important dynamics contributed to the diversification of the chloroplast genome in *Camellia sinensis*, which were not mutually exclusive. Chinese tea and Indian tea might have undergone different selection pressures. Chloroplast transfer occurred during the polyploid evolution in *Camellia sinensis*. In addition, our results supported the three different domestication origins of Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea. And, the current classification of some *Camellia* species might need to be further discussed.

**Keywords:** *Camellia sinensis*, *Camellia*, Chloroplast genome, Evolutionary dynamics, Chloroplast transfer, Divergence time, Taxonomy

* Correspondence: zizheng2006@163.com; 10817788@qq.com
Li Li, Yunfei Hu and Min He are first authors.
Li Li, Yunfei Hu and Min He contribute equally to this work.
[1]College of Tea and Food Science, Wuyi University, 358# Baihua Road, Wuyishan 354300, China
Full list of author information is available at the end of the article

## Background

Because of frequent hybridization and polyploidization, the mechanisms operating in the evolution of *Camellia* has always been focus of botanical and ecological research [1–3]. Tea plant (*Camellia sinensis*) is a member of the Theaceae family of angiosperms, and is highly regarded as the oldest and most popular nonalcoholic beverage with huge economic values in the world [4]. Cultivated tea plants have been divided into three distinct groups: *Camellia sinensis* var. *sinensis* (L.) O. Kuntze (Chinary type), *Camellia sinensis* var. *assamica* (Masters) Chang (Assamica type) and *C. sinensis* var. *assamica* subssp. *Lasiocalyx* Planch (Cambodia type). Of which, the most obvious distinction is between *C. sinensis* var. *sinensis* and *C. sinensis* var. *assamica.* In brief, *C. sinensis* var. *sinensis* has small leaves and is major cultivated in China and some Southeast Asian countries, while *C. sinensis* var. *assamica* has large leaves and widely grown in India and some hot countries except for southern China [5–7]. It has long been suggested that *C. sinensis* var. *sinensis* and *C. sinensis* var. *assamica* may have distinct origins, but the idea that *C. sinensis* var. *assamica* consists of two distinct lineages (Chinese Assamica type and Indian Assamica type) that were domesticated separately is more controversial [8].

Chloroplast (cp) genomes are highly conserved in sequence and structure due to their non-recombinant, haploid, and uniparentally inherited nature [9]. Nonetheless, the gene losses and/or additions, rearrangements and repeats within cp genomes had been revealed in many angiosperm lineages [10–13]. Additionally, gene transfer between plastome, chondrome and nucleus had also been found in plants [14, 15]. Therefore, cp genome structural variations are accompanied by speciation over time, which can provide a wealth of evolutionary information [16]. In previous studies, cp genomes had been found to be particularly useful for phylogenetic and phylogeographic studies in the contexts of reticulate evolution (i.e. hybridization) and polyploidization that characterize the history of most plant lineages [17–20]. Some studies also had found that the cp genome resources could provide useful data for eliciting the evolutionary relationships of tea plants, thus reflecting important evidence for a well-supported hypothesis of classification [21]. Up to now, more than 30 complete cp genomes of *Camellia* species had been sequenced [22]. These massive data, helped from their conserved evolution, promotes the use of cp sequences as an effective tool for *Camellia* species phylogenomic analyses.

In addition to interspecific hybridization, polyploidization is another important factor in the diversification of angiosperm plants [23, 24]. cpDNA variation could provide valuable genetic markers for the analysis of polyploids. Non-recombination and uniparental inheritance
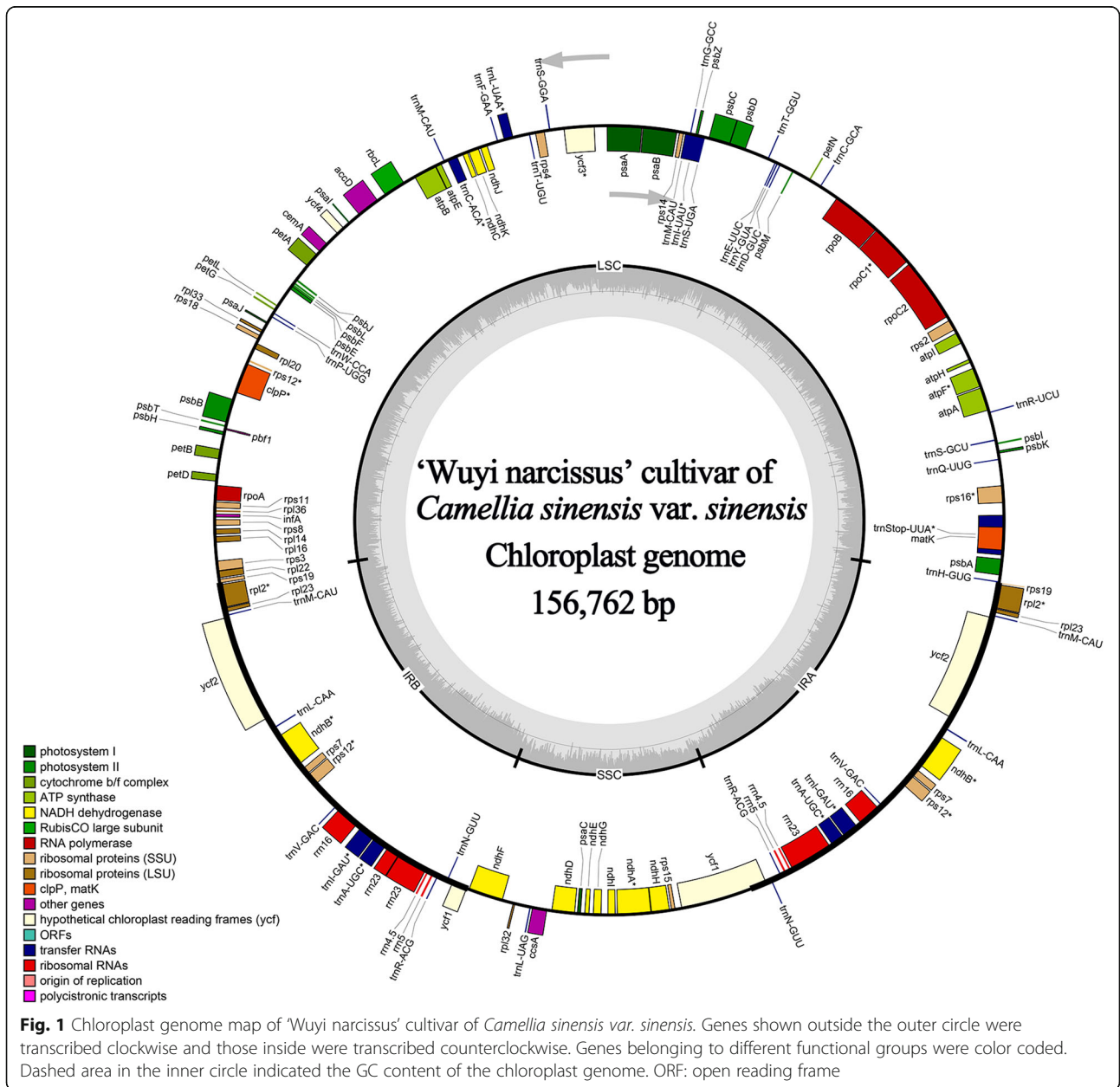
had made cpDNA marker a good indicator of maternal ancestry which could be easily identified in putative hybrid progeny in the absence of parental information, regardless of how many generations had past [25–28]. Using cpDNA marker as sequence characterized amplified region (SCAR) to screen for cp differences between species had proven to be utility in analysis of maternal ancestry of polyploid [29]. In a previous study on the evolution of allotetraploid *Brassicas*, cpDNA data revealed not only the maternal origin of three allotetraploids, but also specific populations of diploids that contributed the cytoplasm to each allotetraploid, and proposed the possibility of introgressive hybridization (chloroplast transfer) [30]. So far, the cp genome of the polyploid tea plant has not been reported, and the possible effects of polyploidization on the cp genome of tea plant need to be further explored.

In this study, we generated the complete cp genome of the first natural triploid tea plant ('Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis*) which belong to asexual propagation cultivar and was recognized as one of the national quality tea varieties by China National Crop Variety Examination Committee in 1985 (GS13009–1985) [31]. Then, we presented the detailed sequence and structural variations of the cp genome among the four representative tea plants, including 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (CWN, a natural triploid Chinary type tea), a diploid *C. sinensis* var. *sinensis* (CSS, Chinary type tea) and two diploid *C. sinensis* var. *assamica* (CSA, Chinese Assamica type tea and CIA, Indian Assamica type tea). Through comparative analysis, we explored the evolutionary dynamics of cp genome and the effects of polyploidization in *C. sinensis.* Furthermore, the phylogenetic analysis and the divergence time estimation based on complete cp genomes were conducted to explore the evolutionary relationship between Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea, and to further improve our understanding of the taxonomic classification of *Camellia*.

## Results

### Chloroplast genome sequencing and assembly

The cp genome of 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* was constructed by PacBio long-reads with Illumina paired-ends data support. In total, 46,941,086 Illumina reads (7.04 Gb, Average read length 145 bp) and 364,638 PacBio reads (10,383 reads > 5000 bp, Average read length 1139 bp) were mapped to the complete genome, respectively. The average organelle coverage reached 43,419× and 2650× sequencing depth, respectively. The de novo assembly using error-corrected PacBio reads resulted in a circular genome of 156,762 bp length (Fig. 1). Raw reads, assembled cp genome sequences and accompanying gene annotations had been

**Fig. 1** Chloroplast genome map of 'Wuyi narcissus' cultivar of *Camellia sinensis var. sinensis*. Genes shown outside the outer circle were transcribed clockwise and those inside were transcribed counterclockwise. Genes belonging to different functional groups were color coded. Dashed area in the inner circle indicated the GC content of the chloroplast genome. ORF: open reading frame

deposited in the NCBI GenBank (SRA: SRR12002624, Accession numbers: MT612435).

**Chloroplast genome structure and characteristics analyses**

All four complete cp genomes displayed the typical quadripartite structure of most angiosperms, including the large single copy (LSC), the small single copy (SSC) and a pair of inverted repeats (IRa and IRb). Among these cp genomes, genome size ranged from 156,762 bp to 157,353 bp due to expansion and contraction of cp genomes. The length varied from 86,301 bp to 87,214 bp in the LSC region, from 18,079 bp to 18,285 bp in the

SSC region, and from 26,030 bp to 26,090 bp in IR region (Table 1).

Each cp genome contained a total of 137 genes, including 92 protein-coding genes, 37 transfer RNA (tRNA) genes and 8 ribosomal RNA (rRNA) (Supplementary Tab. S1). Of them, 60 protein-coding and 22 tRNA genes were located within LSC, 16 protein-coding genes, 14 tRNA coding genes and eight rRNA coding genes were located within IRs and 11 protein-coding and one tRNA gene were located within SSC. The rps12 gene was a divided gene with the 5′ end exon located in the LSC region while two copies of 3′ end exon and intron were located in the IRs. The ycf1 was located in the

Li *et al. BMC Genomics*     (2021) 22:138

Page 4 of 22

**Table 1** Summary of four chloroplast genome features

| Genome Features | CWN (MT612435) | CSS (KJ806281) | CSA (MH019307) | CIA (MH460639) |
|---|---|---|---|---|
| Location of sample | Fujian, China | Yunnan, China | Yunnan, China | Assam, India |
| Longitude | 118.004001 | 102.714601 | 102.714601 | 94.228661 |
| Latitude | 27.72846 | 25.04915 | 25.04915 | 26.73057 |
| Genome size (bp) | 156,762 | 157,117 | 157,100 | 157,353 |
| LSC length (bp) | 86,301 | 86,662 | 86,649 | 87,214 |
| SSC length (bp) | 18,281 | 18,275 | 18,285 | 18,079 |
| IR length (bp) | 26,090 | 26,090 | 26,083 | 26,030 |
| Number of genes | 137 | 137 | 137 | 137 |
| Number of Protein-coding genes | 92 | 92 | 92 | 92 |
| Number of tRNA genes | 37 | 37 | 37 | 37 |
| Number of rRNA genes | 8 | 8 | 8 | 8 |
| GC content of LSC (%) | 35.32 | 35.31 | 35.31 | 35.38 |
| GC content of SSC (%) | 30.55 | 30.56 | 30.51 | 30.59 |
| GC content of IR (%) | 42.94 | 42.95 | 42.95 | 42.96 |
| Overall GC content (%) | 37.3 | 37.3 | 37.29 | 37.34 |

*CWN* 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea), *CSS C. sinensis* var. *sinensis* (diploid Chinary type tea), *CSA C. sinensis* var. *assamica* (diploid Chinese Assamica type tea), *CIA C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

boundary regions between IR/SSC, leading to incomplete duplication of the gene within IRs. There were 18 genes containing introns, including 6 tRNA genes and 12 protein-coding genes. Except for two introns in the ycf3 and clpP genes, all other genes contained only one intron. MatK gene was located within the intron of trnK-UUU with the largest intron (2489 bp). Overlaps of adjacent genes were found in the complete genome, rps3-rpl22, atpB-atpE, and psbD-psbC had a 16 bp, 4 bp, and 53 bp overlapping region, respectively. Unusual initiator codons were observed in rps19 with GTG and orf42 with ATC in four cp genomes. The initiation codon of ndhD in *CIA* was ATG, while that of other three cp genomes was GTG.

## Sequence variation analyses

The differences and evolutionary divergences among four cp genomes were compared using nucleotide substitutions and sequence distance. Across all four species, the value of nucleotide differences was 70–185, and the p-distance was 0.00045–0.00118. The value of nucleotide difference (70) and the p-distance (0.00045) between triploid *CWN* and diploid *CSS* was smallest (Table 2).

To identify the potential genome rearrangements and inversions, the cp genome sequences of four species were plotted to check their identity using the program mVISTA. No gene rearrangement and inversion events were detected (Fig. 2). Sequence divergence analyses showed four regions (including rp12/trnH-UGU, psaA/ycf3, atpB/rbcL and psbT/psbH) had relatively higher divergence values (Pi > 0.006) (Fig. 3). Mutations of the
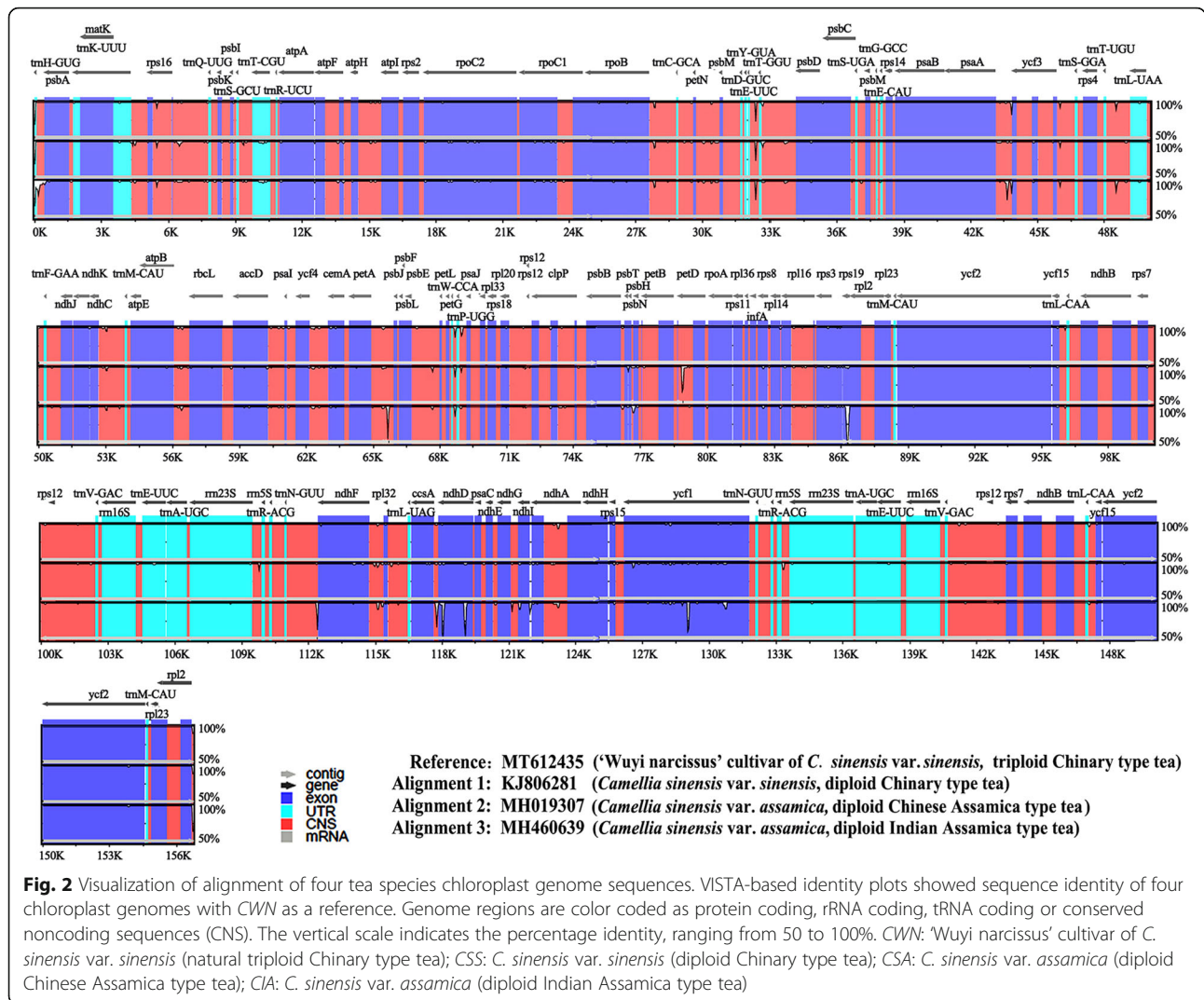
base replacement or deletion may cause changes in the length of the coding gene sequence, leading to changes in the coding and non-coding regions. Therefore, the variable characters in coding and non-coding regions among four cp genomes were further analyzed. The results showed that the proportion of variability in non-coding regions was with a mean value of 1.82%, while in the coding regions was 1.15%. Five coding genes had over 4% variability proportion, such as rps19, ndhF, ndhD, ndhI and ycf1. Five non-coding regions had over 10% variability proportions, such as rpl2/trnH-GUG, trnE-UUC/trnT-GGU, ndhD/psaC, ndhI/ndhA and rps15/ycf1 (Fig. 4).

To further observe the potential contraction and expansion of IR regions, the gene variation at the IR/SSC and IR/LSC boundary regions of the four plastomes was compared (Fig. 5). The genes rps19, ycf1–5'end/ndhF, ycf1 and rp12/trnH-GUG were located in the junctions

**Table 2** Numbers of nucleotide substitutions and sequence distance in four complete cp genomes

|  | CWN | CSS | CSA | CIA |
|---|---|---|---|---|
| CWN |  | 0.00045 | 0.00118 | 0.00115 |
| CSS | 70 |  | 0.00115 | 0.00105 |
| CSA | 185 | 180 |  | 0.00100 |
| CIA | 180 | 164 | 157 |  |

The lower triangle shows the number of nucleotide substitutions and the upper triangle indicates the number of sequence distance in complete cp genomes. *CWN* 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea), *CSS C. sinensis* var. *sinensis* (diploid Chinary type tea), *CSA C. sinensis* var. *assamica* (diploid Chinese Assamica type tea), *CIA C. sinensis* var. *assamica* (diploid Indian Assamica type tea)
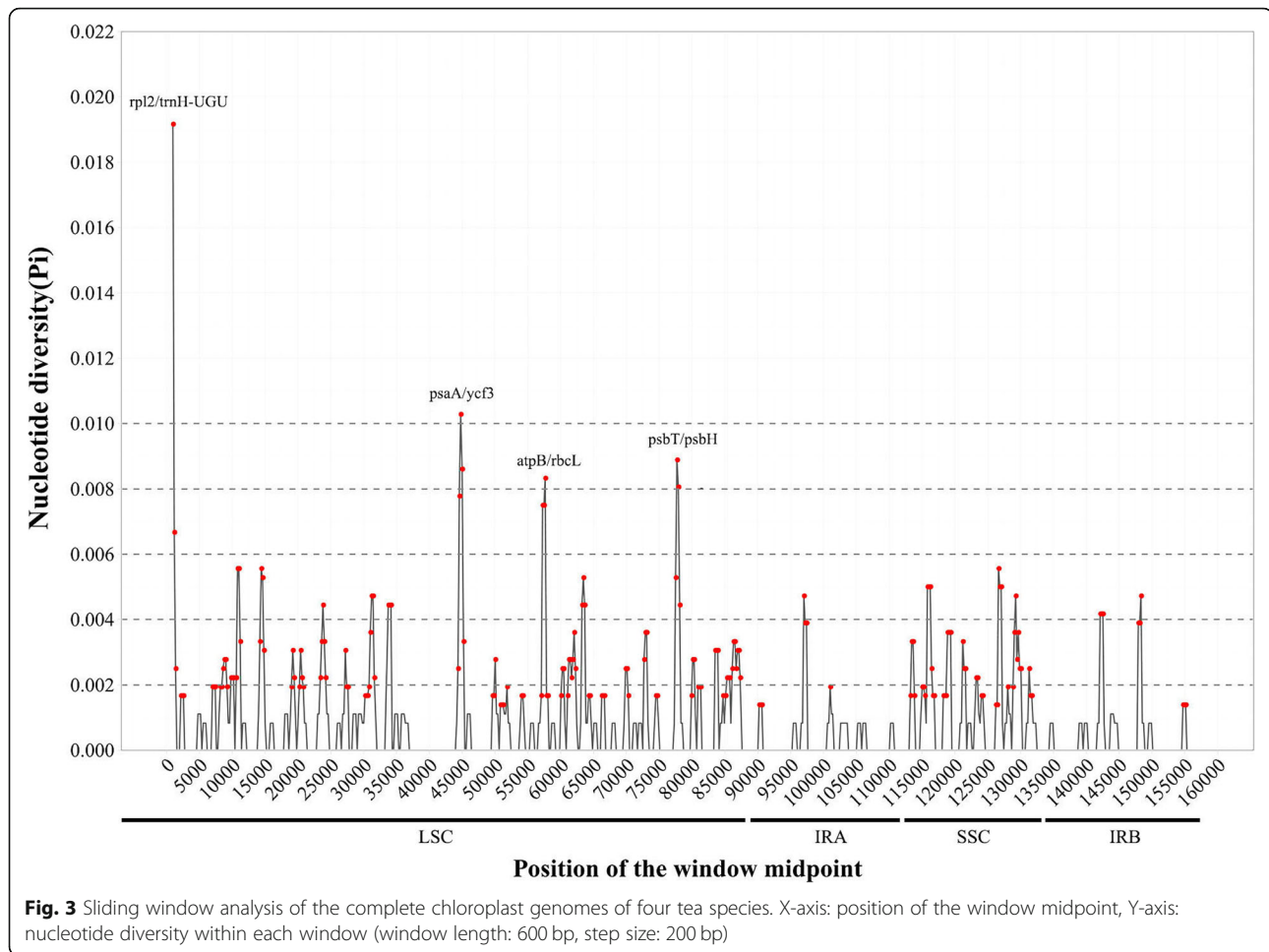
**Fig. 2** Visualization of alignment of four tea species chloroplast genome sequences. VISTA-based identity plots showed sequence identity of four chloroplast genomes with *CWN* as a reference. Genome regions are color coded as protein coding, rRNA coding, tRNA coding or conserved noncoding sequences (CNS). The vertical scale indicates the percentage identity, ranging from 50 to 100%. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

of LSC/IR and SSC/IR regions. The rps19 gene in *CSS*, *CSA*, and *CWN* was 279 bp, and crossed the LSC/IR region by 46 bp while the rps19 gene in *CIA* was just 150 bp, and all located in the LSC region, 1 bp away from the IR region. The ycf1–5'end gene in *CSS*, *CSA*, and *CWN* was 1071 bp, and crossed the IR/SSC region by 2 bp while in *CIA* was 1065 bp, and crossed the IR/SSC region by 33 bp. The ndhF gene in all four cp genomes was located in the SSC region. The ndhF gene in *CSA*, *CIA*, and *CWN* was 2247 bp while in *CSS* was 2139. The ndhF gene in *CSS* was 165 bp away from the IR region, in *CSA* or *CWN* was 57 bp away from the IR region while in *CIA* was 88 bp away from the IR region. The ycf1 gene in *CSS* or *CWN* was 5622 bp, in *CSA* was 5628 bp while in *CIA* was only 1038 bp. The ycf1 genes in all four cp genomes crossed the IR/SSC region. The ycf1 gene in *CSS* or *CWN* was with 4553 bp located in the SSC region and 1069 bp in IR region, in *CSA* was with 4559 bp located in the SSC region and 1069 bp in

IR region while in *CIA* was with only 6 bp located in the SSC region and 1032 bp in IR region. The rpl2 gene in *CSS*, *CSA* or *CWN* was 107 bp away from the LSC region while in *CIA* was 82 bp away from the LSC region. The trnH-GUG gene in *CSS*, *CSA* or *CWN* was 2 bp away from the IR region while in *CIA* was 637 bp away from the IR region.

## Repeat and indel sequence analyses

Simple sequence repeats (SSRs) are small repeating units of cpDNA, a total of 671 SSRs were identified in four cp genomes (Fig. 6a), of which 57% were in IGS, 34% were in CDS, and 9% were in Intron (Fig. 6b). 74.0% of these SSRs were monomers, 19.3% of dimers, 0.5% of trimers, 5.3% of tetramers, 0.9% of hexamers and no pentamers found. Comparing the four genomes, except for 167 SSRs of *CIA*, the other three were all 168. A total of 128 SSRs were identical among four cp genomes (Fig. 6c). There were 47 loci

**Fig. 3** Sliding window analysis of the complete chloroplast genomes of four tea species. X-axis: position of the window midpoint, Y-axis: nucleotide diversity within each window (window length: 600 bp, step size: 200 bp)

with different SSR types, most of which existed in the LSC region. Among them, *CSS* had 7 unique types, *CSA* had 18 unique types, *CIA* had 9 unique types, and *CWN* had 14 unique types (Fig. 6c, Supplementary Tab. S2).

A total of 270 long repeats were detected in four plastomes, including three categories of long repeats: tandem, forward and palindromic. The number of the three repeated types was consistent in *CSS* and *CWN*, as follows: 23, 20, 23. However, it was 19, 20, 23 in *CSA* and 21, 23, 32 in *CIA*. The sizes of repeats ranged from 11 to 82 bp (Fig. 7a, c). The four cp genomes have a total 57 identical long repeat sequences. In addition, *CSS* had 1 unique long repeat, *CIA* had 1 unique long repeat, *CWN* had 2 unique long repeats, while *CSA* had no unique long repeat (Fig. 7b). These unique repeats were found mainly in the intergenic psaA/ycf3, atpB/rbcL, trnW-CCA/ trnP-UGG, rps19/rpl2, psbT/psbN, rpl2/trnH-GUG and gene rpl2, ycf1, ycf2. Only one repeat was in the intron regions (ndhA) (Supplementary Tab. S3).
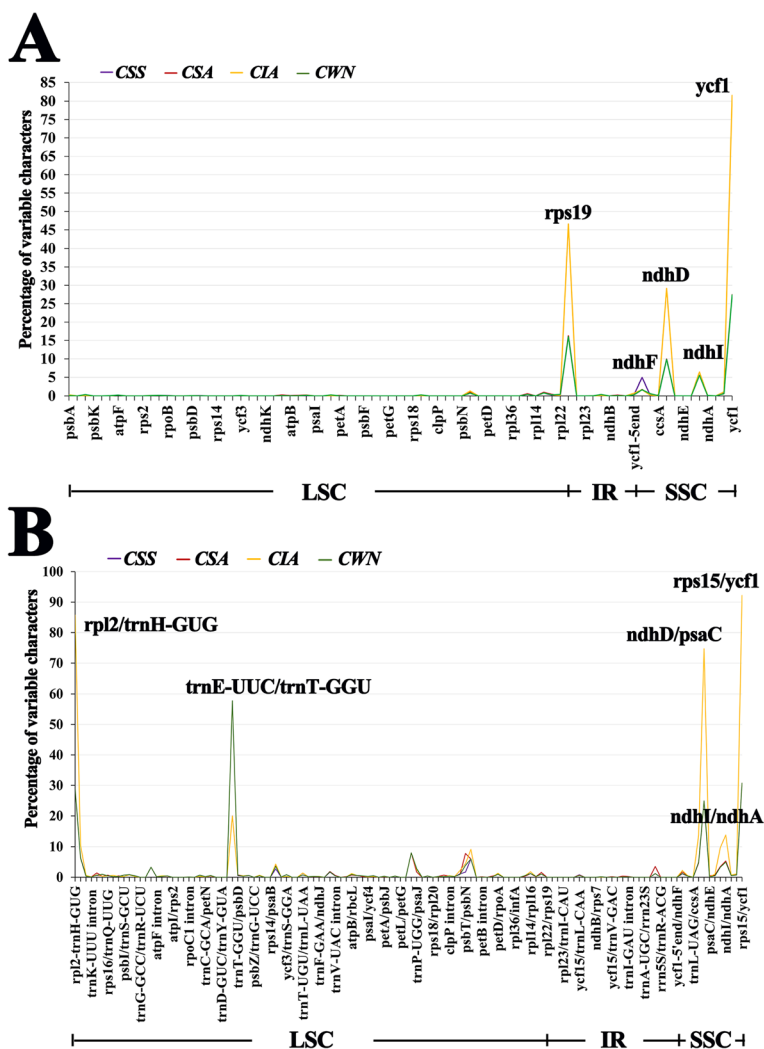
A total of 100 indels were found, and indels ranged in size from 1 to 637 bp (Fig. 8a). Most of the indels events occurred in IGS regions (70%), with 23% in CDS regions

and only 7% in Intron regions (Fig. 8b). As expected, single-nucleotide indels (1 bp) were the most common, but some long indels also were found. The longest one was an insertion of 637 bp in *CIA* (intergenic rp12/trnH-GUG), followed by a 335 bp deletion in *CWN* (intergenic trnE-UUC/trnT-GGU) and a 107 bp deletion in *CIA* (gene rps19). Paired comparison showed that the *CIA* had the most indels compared to the other three species (Fig. 8c). In addition, *CIA* also possessed the most species-specific indels, with 49, followed by *CSA* with 16, *CWN* with 11 and *CSS* with 5 (Fig. 8d, Supplementary Tab. S4).
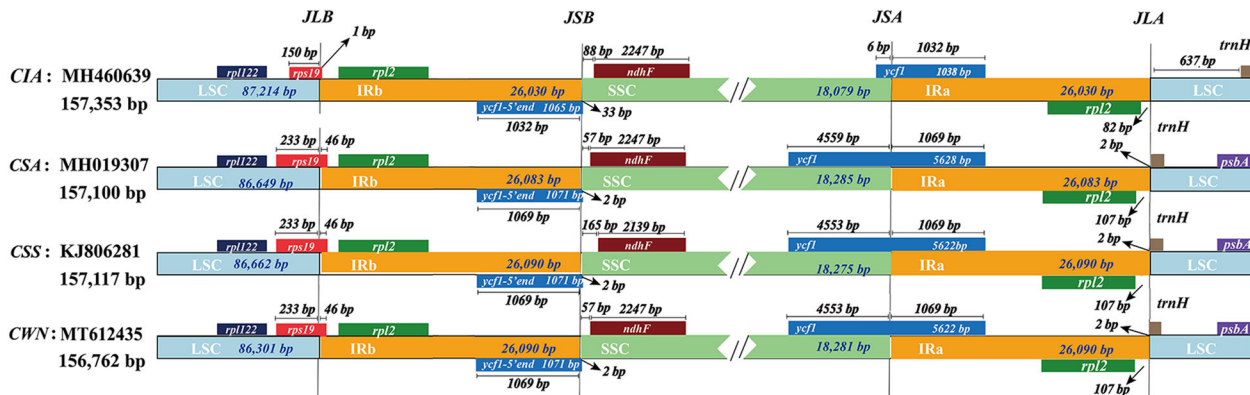
The regions with relatively high divergence values (rp12/trnH-UGU, psaA/ycf3, atpB/rbcL and psbT/psbH, Pi > 0.006) (Fig. 3) all were associated with the repeat and the indel sequences. For example, the repeat sequences could be found within the region of rp12/trnH-UGU, atpB/rbcL and psbT/psbH. The indel sequences could be found within the region of rp12/trnH-UGU, psaA/ycf3 and psbN/psbH.

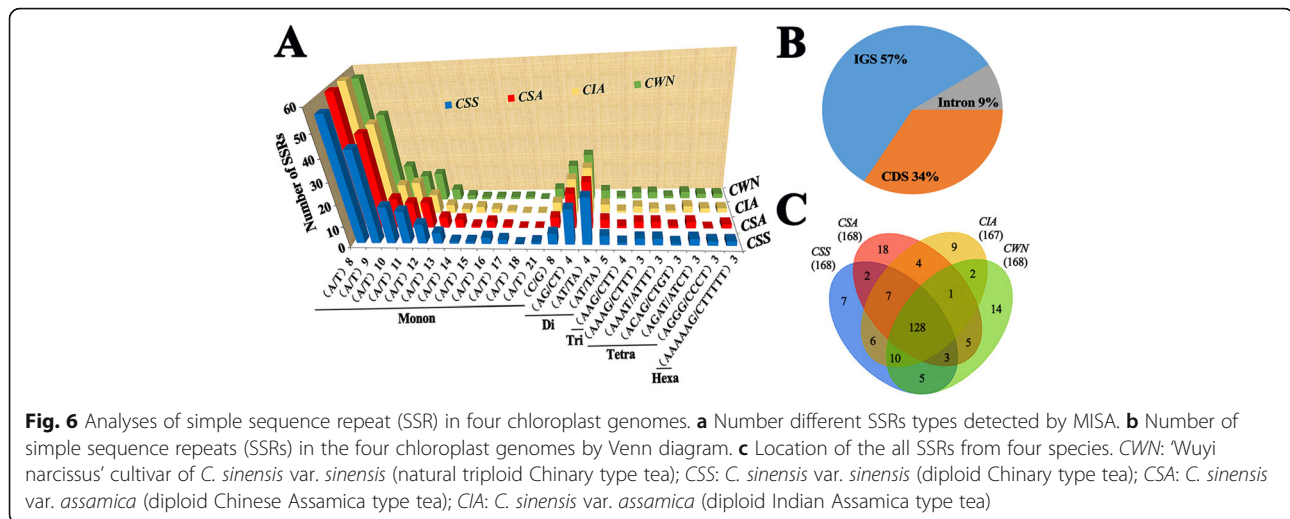## Correlation analysis of three types of mutation

Correlations were highly significant in the pairwise comparisons between the three types of mutations: "repeats

**Fig. 4** Percentages of variable characters in homologous regions across the four chloroplast genomes. **a** Coding regions. **b** Non-coding regions. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)



**Fig. 5** The comparison of the LSC, IR and SSC border regions among the four chloroplast genomes. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

**Fig. 6** Analyses of simple sequence repeat (SSR) in four chloroplast genomes. **a** Number different SSRs types detected by MISA. **b** Number of simple sequence repeats (SSRs) in the four chloroplast genomes by Venn diagram. **c** Location of the all SSRs from four species. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

and substitutions", "indels and substitutions" and "repeats and indels". The strength of correlations was greatest for "indels and substitutions" (r: 0.165–0.435) followed by "repeats and indels" (r: 0.090–0.120) and then "repeats and substitutions" (r: 0.028–0.049), and "indels and substitutions" had relatively higher significance value (t: 0.144–0.195) than "repeats and substitutions" (t: 0.103–0.145) (Table 3).

### Codon usage analyses

ENc plots analysis showed only a few points lie near the curve, however, most of the genes with lower ENc values than expected values lay below the curve (Fig. 9), suggesting the codon usage bias of the cp genome was slightly affected by the mutation pressure, but selection and other factors play an important role. To further investigate the extent of influence between mutation pressure and natural selection on the codon usage patterns, Neutrality plot (GC12 vs. GC3) was performed. The correlation between GC1 and GC2 was strong (*CSS*: r = 0.445; *CSA*: r = 0.453; *CIA*: r = 0.445; *CWN*: r = 0.464, $p < 0.01$). However, no significant correlation was found for GC1 with GC3 (*CSS*: r = 0.141; *CSA*: r = 0.139; *CIA*: r = 0.078; *CWN*: r = 0.141) or GC2 with GC3 (*CSS*: r = 0.146; *CSA*: r = 0.143; *CIA*: r = 0.078; *CWN*: r = 0.152), which suggested mutation pressure had a minor effect on the codon usage bias. The slope of Neutrality plot showed that mutation pressure accounts for only 0.52–8.42% on the codon usage patterns in four cp genomes, while natural selection accounts for 91.58–99.48% (Fig. 10).

The distributions of codon usage in four cp genomes showed that RSCU values of the 37 codons (37/64, 57.81%) were identical in the three Chinese teas, but different from those in Indian tea (Table 4).

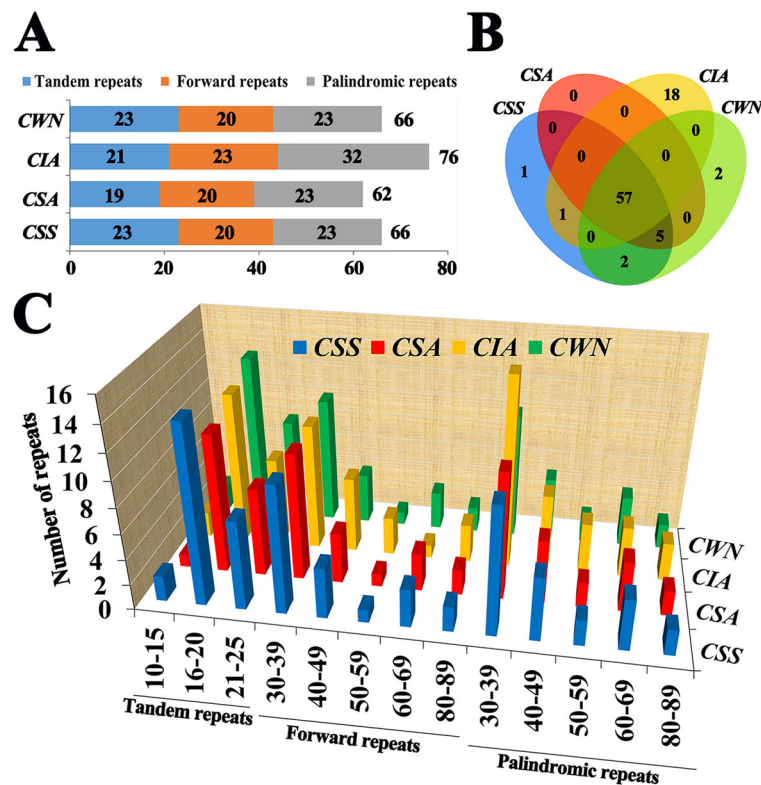### Analysis of cp sequence characterized amplified region (SCAR)

By comparing with the cp genomes of three representative diploid *C. sinensis* species, a 335 bp long deletion in the trnE/trnT intergenic spacer was found in triploid *CWN* (Fig. 11a). We used this marker for SCAR analysis in 292 individuals covering the majority of *C. sinensis* cultivars in China. No cultivar with similar sequence deletion characteristics to triploid *CWN* was detected (Fig. 11b, Supplementary Fig. S1, Supplementary Tab. S5). However, we could find similar sequence deletion in *C. cuspidate* (Accession numbers: NC022459), *C. renshanxiangiae* (Accession numbers: NC041672), *C. elongata* (Accession numbers: NC035652) and *C. gymnogyna* (Accession numbers: NC039626) by comparing the cp genome sequences (Fig. 11a).

### Phylogenetic analysis and the divergence time estimation of three tea plants

Phylogenetic trees were generated by ML and BI analysis based on 44 complete cp genomes showed the same topology. Cultivated tea plants were clustered into a single clade, within which Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea were in separate lineages with high support, respectively (Figs. 12 and 13, Supplementary Tab. S6).

Excluding seven non-*Camellia* species, the sequence variation of the 37 *Camellia* species associated with the six datasets (Complete cp genome, LSC, SSC, IR, PCGs, and non-PCGs) showed different percentage variation (Supplementary Tab. S7). SSC had the highest percentage variation at 2.32%, followed by non-PCGs at 1.65%. The IR regions were least variable at 0.5%. The cp genome, LSC, and PCGs, were 1.3, 1.54 and 1.21%, respectively. Phylogenetic trees based on six different data sets

**Fig. 7** Analyses of repeated sequences in four chloroplast genomes. **a** Number of the three repeat types. **b** Number of repeat sequences in the four chloroplast genomes by Venn diagram. **c** Number of the repeats by different length. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

showed mostly similar topologies. A few individual species were retrieved incongruently among different clades across the six data partitions, but all *Camellia* species remained grouped separately, except IR regions that were shown to be mixed with *Polyspora* species of *Theaceae*. The support values of nodes increased significantly with the increasing of the sequence length in the different data partitions. In terms of interspecific relationships of three tea plants (Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea), the results showed the same topology across all six datasets (Figs. 12 and 13, Supplementary Fig. S2, S3, S4, S5, S6).
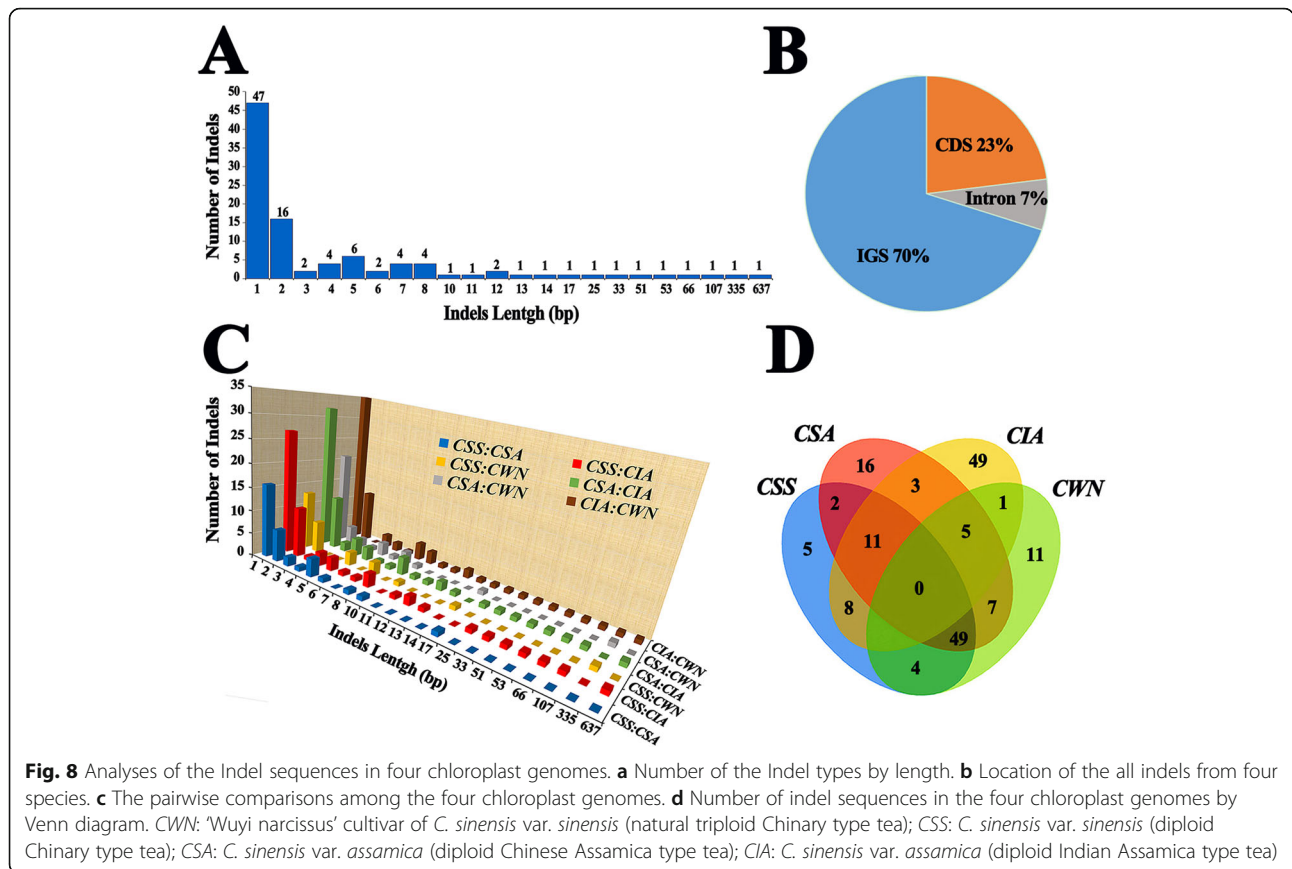
Estimated divergence time showed the three types of tea plant were diverged to each other during 0.8–6.2 million years ago (Mya) (CI: 0.3–8.1 Mya). Indian Assamica type tea diverged from the ancestor of Indian Assamica type tea and Chinese Assamica type tea about 6.2 Mya (CI: 4.4–8.1 Mya, Miocene), Chinese Assamica type tea diverged separately about 0.8 Mya (CI: 0.3–1.6 Mya, Quaternary), and Chinary type tea diverged separately from the ancestor of Indian Assamica type tea and Chinary type tea about 0.8 Mya (CI: 0.4–1.5 Mya, Quaternary) (Fig. 12).

## Discussion

### Genetic variation and mutational dynamics of the chloroplast genome in tea plant

The four cp genomes of the tea plants showed a high degree of conservation in genome structure, gene content, gene order, intron number, and also GC content. To better understand the sequence variation in tea plant, the three important types of genetic variation in cp genome, inducing nucleotide substitutions, repeats and indels [33–36], were identified. In addition to nucleotide substitutions, 671 SSRs (simple repeat) were identified (another 32, 31, 31 and 30 SSRs occurred in compound formations for *CSS*, *CSA*, *CIA* and *CWN*, respectively). The number of SSRs was consistent with a previous study [37]. In addition, a total of 270 long repeats and 100 indels also were identified. The repeats and indels identified here might provide information for markers development to further species identification and population genetic studies [38, 39].
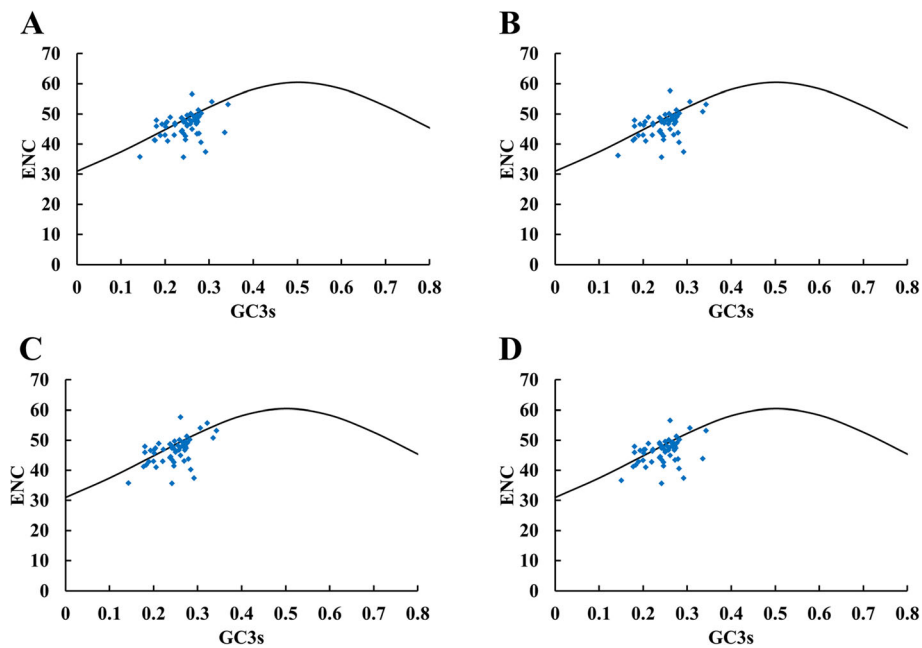
A characteristic feature of eukaryote and prokaryote genomes is the co-occurrence of nucleotide substitution and insertion/deletion (indel) mutations [40]. We also found that the divergent regions of cp genomes were
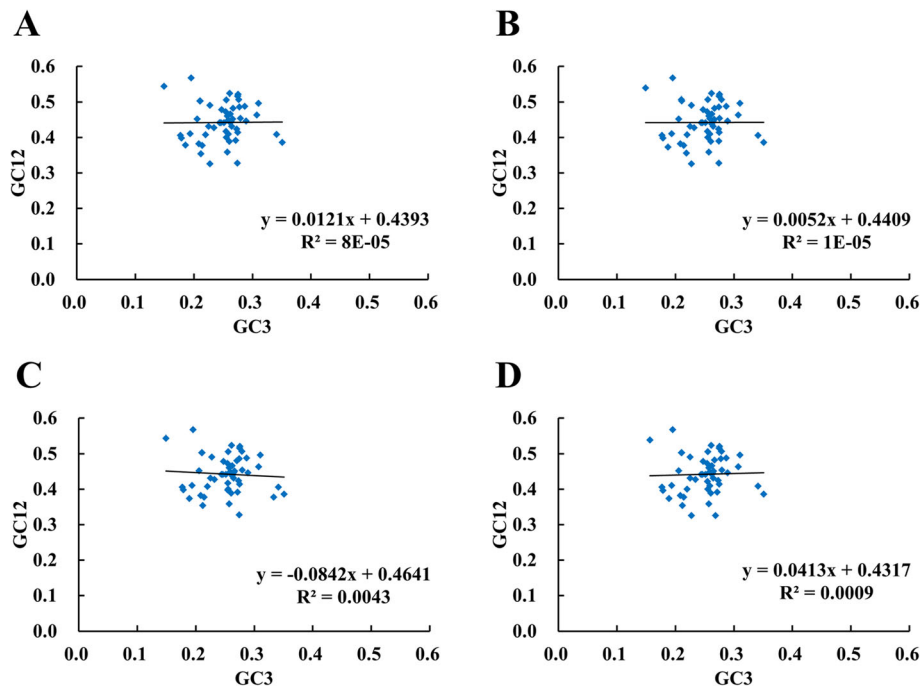
Li *et al. BMC Genomics*    (2021) 22:138

Page 10 of 22



**Fig. 8** Analyses of the Indel sequences in four chloroplast genomes. **a** Number of the Indel types by length. **b** Location of the all indels from four species. **c** The pairwise comparisons among the four chloroplast genomes. **d** Number of indel sequences in the four chloroplast genomes by Venn diagram. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

**Table 3** Correlation analysis of three types of mutation

| Comparison | CSA | CIA | CWN |
|---|---|---|---|
| **Repeats and Substitutions** | | | |
| Correlation between repeats and substitutions (r) | 0.033 | 0.049 | 0.028 |
| Significance of correlation (t) | 0.103** | 0.103** | 0.145** |
| Coefficient of determination ($r^2$) | 0.0011 | 0.0024 | 0.0008 |
| **Indels and Substitutions** | | | |
| Correlation between indels and substitutions (r) | 0.207 | 0.435 | 0.165 |
| Significance of correlation (t) | 0.158** | 0.195** | 0.144** |
| Coefficient of determination ($r^2$) | 0.043 | 0.189 | 0.0273 |
| **Repeats and Indels** | | | |
| Correlation between repeats and indels (r) | 0.090 | 0.099 | 0.120 |
| Significance of correlation (t) | 0.195** | 0.221** | 0.268** |
| Coefficient of determination ($r^2$) | 0.0081 | 0.0098 | 0.0145 |

Comparisons among the pairwise alignments (*CSS* taken as a Reference) to calculate the correlations between Repeats and Substitutions, Insertion-Deletions (Indels) and Substitutions, and Repeats and Indels. The alignments were partitioned into 630 nonoverlapping bins of 250 bp size each to calculate these correlations. ** indicated high significance. *CWN* 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea), *CSS* *C. sinensis* var. *sinensis* (diploid Chinary type tea), *CSA* *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea), *CIA* *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

**Fig. 9** ENc-plot of chloroplast genomes of four tea species. **a** *C. sinensis* var. *sinensis* (diploid Chinary type tea); **b** *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); **c** *C. sinensis* var. *assamica* (diploid Indian Assamica type tea); **d** 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea)



**Fig. 10** Neutrality plot of chloroplast genomes of four tea species. **a** *C. sinensis* var. *sinensis* (diploid Chinary type tea); **b** *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); **c** *C. sinensis* var. *assamica* (diploid Indian Assamica type tea); **d** 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea)

**Table 4** The relative synonymous codon usage (RSCU) values of four chloroplast genomes

| AA | Condon | Species | | | | AA | Condo n | Species | | | | AA | Condon | Species | | | |
|----|--------|-----|-----|-----|-----|----|---------|-----|-----|-----|-----|----|--------|-----|-----|-----|-----|
| | | CSS | CSA | CIA | CWN | | | CSS | CSA | CIA | CWN | | | CSS | CSA | CIA | CWN |
| Phe | UUU | 1.32 | 1.32 | 1.3 | 1.32 | Pro | CCU | 1.65 | 1.66 | 1.69 | 1.65 | Lys | AAA | 1.53 | 1.53 | 1.51 | 1.53 |
| | UUC | 0.68 | 0.68 | 0.7 | 0.68 | | CCC | 0.71 | 0.71 | 0.7 | 0.71 | | AAG | 0.47 | 0.47 | 0.49 | 0.47 |
| Leu | UUA | 1.97 | 1.97 | 1.96 | 1.98 | | CCA | 1.17 | 1.16 | 1.14 | 1.16 | Asp | GAU | 1.63 | 1.63 | 1.62 | 1.63 |
| | UUG | 1.24 | 1.23 | 1.25 | 1.24 | | CCG | 0.47 | 0.47 | 0.47 | 0.47 | | GAC | 0.37 | 0.37 | 0.38 | 0.37 |
| | CUU | 1.25 | 1.25 | 1.24 | 1.25 | Thr | ACU | 1.66 | 1.66 | 1.68 | 1.66 | Glu | GAA | 1.53 | 1.54 | 1.51 | 1.54 |
| | CUC | 0.39 | 0.39 | 0.4 | 0.39 | | ACC | 0.74 | 0.74 | 0.74 | 0.74 | | GAG | 0.47 | 0.46 | 0.49 | 0.46 |
| | CUA | 0.78 | 0.78 | 0.8 | 0.77 | | ACA | 1.22 | 1.22 | 1.2 | 1.22 | Cys | UGU | 1.52 | 1.52 | 1.51 | 1.52 |
| | CUG | 0.37 | 0.37 | 0.35 | 0.37 | | ACG | 0.38 | 0.38 | 0.38 | 0.38 | | UGC | 0.48 | 0.48 | 0.49 | 0.48 |
| Ile | AUU | 1.46 | 1.46 | 1.48 | 1.46 | Ala | GCU | 1.84 | 1.84 | 1.85 | 1.84 | Trp | UGG | 1 | 1 | 1 | 1 |
| | AUC | 0.58 | 0.58 | 0.59 | 0.58 | | GCC | 0.62 | 0.62 | 0.63 | 0.63 | Arg | CGU | 1.39 | 1.39 | 1.44 | 1.39 |
| | AUA | 0.96 | 0.96 | 0.93 | 0.96 | | GCA | 1.16 | 1.15 | 1.14 | 1.16 | | CGC | 0.33 | 0.34 | 0.35 | 0.33 |
| Met | AUG | 1 | 1 | 1 | 1 | | GCG | 0.38 | 0.38 | 0.39 | 0.38 | | CGA | 1.42 | 1.42 | 1.44 | 1.43 |
| Val | GUU | 1.49 | 1.48 | 1.48 | 1.49 | Tyr | UAU | 1.62 | 1.62 | 1.63 | 1.62 | | CGG | 0.39 | 0.39 | 0.4 | 0.39 |
| | GUC | 0.43 | 0.43 | 0.44 | 0.43 | | UAC | 0.38 | 0.38 | 0.37 | 0.38 | Ser | AGU | 1.31 | 1.31 | 1.33 | 1.31 |
| | GUA | 1.53 | 1.54 | 1.53 | 1.53 | TER | UAA | 1.56 | 1.5 | 1.5 | 1.5 | | AGC | 0.32 | 0.32 | 0.33 | 0.32 |
| | GUG | 0.55 | 0.55 | 0.56 | 0.55 | | UAG | 0.63 | 0.63 | 0.63 | 0.63 | Arg | AGA | 1.86 | 1.85 | 1.75 | 1.85 |
| Ser | UCU | 1.78 | 1.78 | 1.8 | 1.79 | | UGA | 0.81 | 0.87 | 0.87 | 0.87 | | AGG | 0.6 | 0.61 | 0.62 | 0.6 |
| | UCC | 0.91 | 0.91 | 0.92 | 0.91 | His | CAU | 1.57 | 1.57 | 1.56 | 1.57 | Gly | GGU | 1.29 | 1.29 | 1.31 | 1.29 |
| | UCA | 1.16 | 1.16 | 1.11 | 1.16 | | CAC | 0.43 | 0.43 | 0.44 | 0.43 | | GGC | 0.45 | 0.45 | 0.45 | 0.45 |
| | UCG | 0.51 | 0.51 | 0.52 | 0.51 | Gln | CAA | 1.53 | 1.53 | 1.51 | 1.53 | | GGA | 1.59 | 1.59 | 1.57 | 1.58 |
| Asn | AAU | 1.57 | 1.57 | 1.55 | 1.57 | | CAG | 0.47 | 0.47 | 0.49 | 0.47 | | GGG | 0.67 | 0.67 | 0.67 | 0.67 |
| | AAC | 0.43 | 0.43 | 0.45 | 0.43 | | | | | | | | | | | | |

RSCU values of the 37 codons (37/64, 57.81%) were identical in all three Chinese teas (Blue background), but different in Indian tea (Yellow background). *CWN* 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea), *CSS C. sinensis* var. *sinensis* (diploid Chinary type tea), *CSA C. sinensis* var. *assamica* (diploid Chinese Assamica type tea), *CIA C. sinensis* var. *assamica* (diploid Indian Assamica type tea)

almost all associated with repeat sequences and indel sequences. Therefore, the genome-wide association between repeat, indel and nucleotide substitution was further analyzed. Correlations were significant in the pairwise comparisons between the three types of mutations: substitutions and indels, repeats and indels, and repeats and substitutions, which provided further support for the repeat-and indel-induced mutation hypothesis. The indel-induced mutation hypothesis assumes that the changes are induced by indels that trigger the DNA repair process, in which error-prone DNA polymerases are recruited [41, 42]. Instead, the repeat-induced mutation hypothesis assumes that mutations are produced because the presence of repeat sequences rather than indels per se, which promote replication fork arrest that led to the recruitment of error-prone polymerases and cause indels and nucleotide substitution [43]. Our results showed that while the "repeats and indels" model was significantly correlated, the "indels and substitutions" model exhibited the highest strength and relative statistical significance for correlation (Table 3). Therefore, we considered these two hypotheses as occurring not mutually exclusive in the cp genome evolution of tea plants. It also implied that if the distribution of the repeat or indel sequences had been identified, the mutational hotspot regions and the appropriate sequences for genetic analysis could be predicted based on their location.

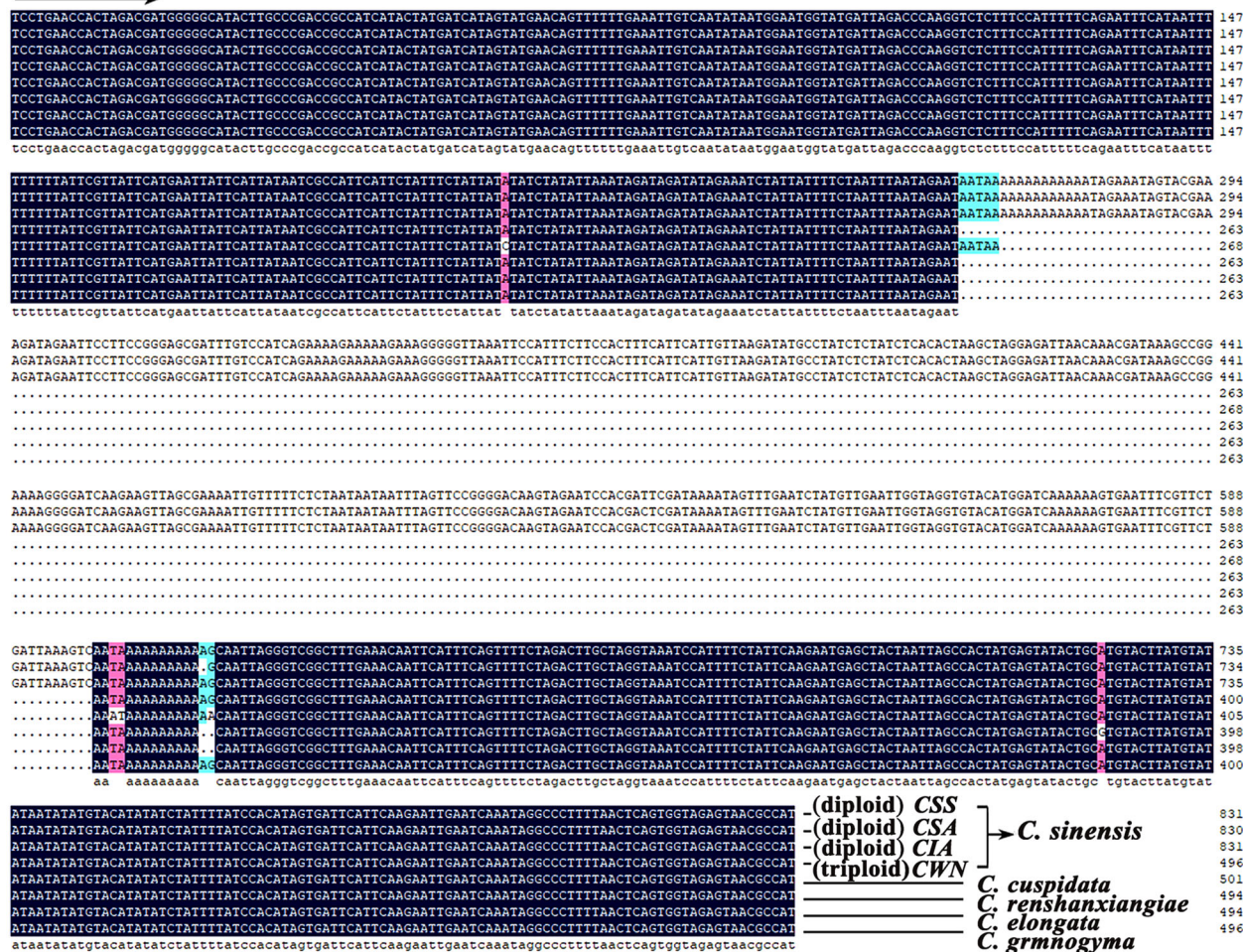The expansion and contraction of cp genome is a common evolutionary phenomenon in plants [44]. Of the four cp genomes, the IR regions of *CIA* were the smallest, and the more sequence deletions were found in SC/IR regions of *CIA*, suggesting that deletions might cause the contraction of the IR region. On the other hand, the LSC region of *CIA* had the longest length, and we found more repeat sequences in LSC/IR region, suggesting the repeats contributed to the length of cp genome. Moreover, a species-specific 335 bp long deletion was found in *CWN* with a length close to the difference length of cp genome between *CWN* and two Chinese tea plants (355 and 338 bp for *CSS* and *CSA*, respectively). Likewise, the difference length of cp genome between *CWN* and *CIA* was 913 bp, which was close to the sum of a species-specific deletion length of *CWN* (335 bp) plus a species-specific long insertion length of *CIA* (637 bp). Therefore, the repeats and indels might have critical effects on the structure of cp genome in tea plant, and the indels might play more important role, which was also well supported with the indel-induced mutation having the highest strength of correlation in the pairwise comparisons between the three types of mutations (Table 3).
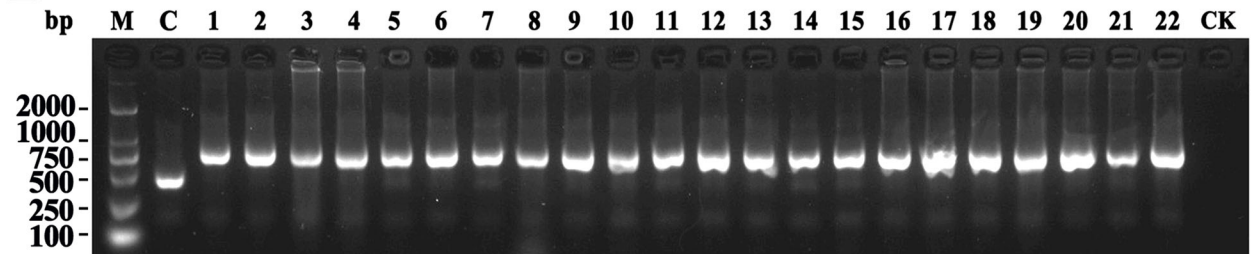
### Difference between Chinese tea and Indian tea
In tea plant, the most obvious division is Chinary type tea (*C. sinensis* var. *sinensis*) and Assamica type tea (*C. sinensis* var. *assamica*), which have markedly different appearance. Further, we compared the cp genome
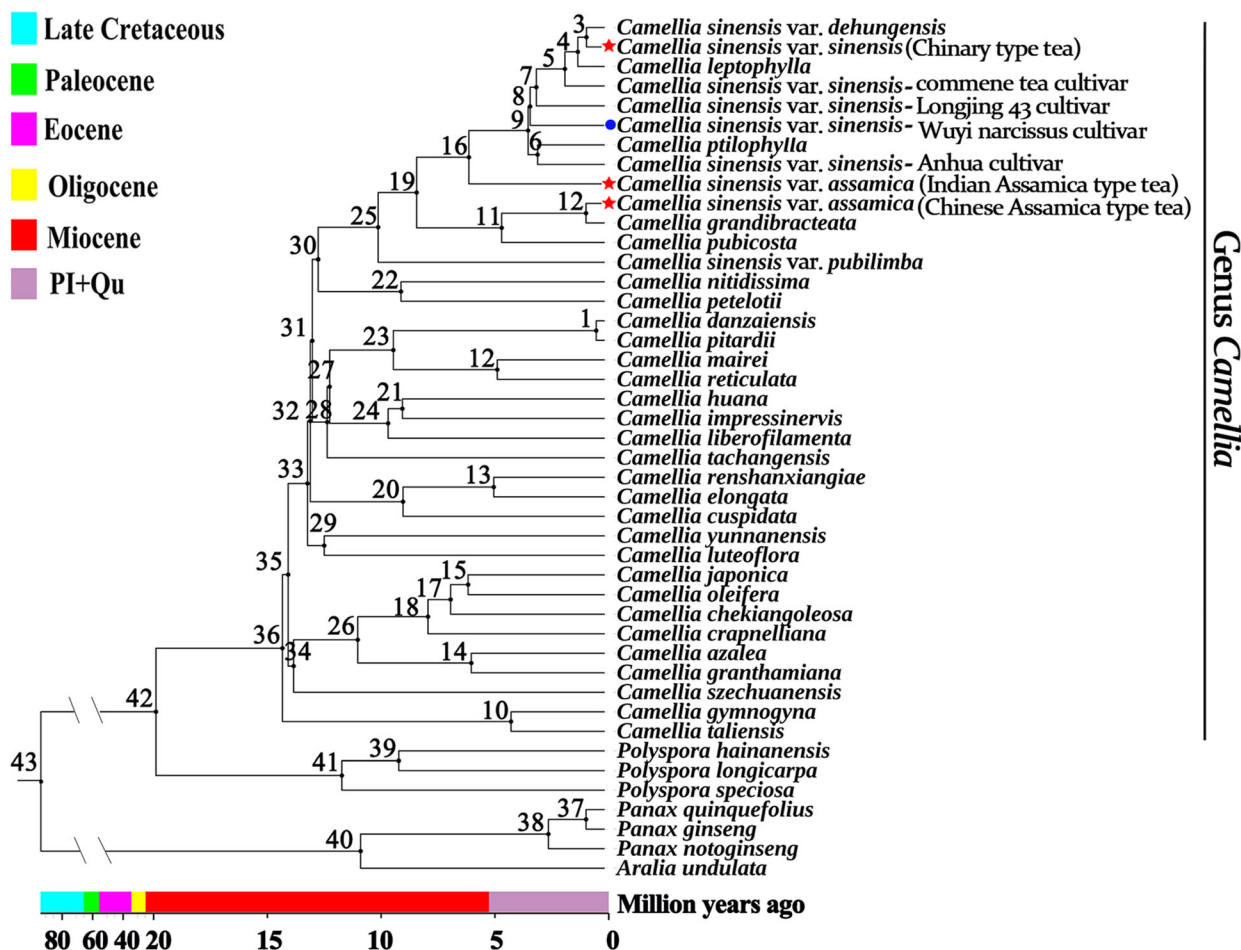
**Fig. 11** Analysis of cp sequence characterized amplified region (SCAR) using a 335 bp deletion of the intergenic spacers (trnE/trnT). **a** By comparing with the cp genomes of three representative diploid *C. sinensis* species, a 335 bp long deletion was observed in triploid *CWN*, and the similar long deletion could be found in four other *Camellia* species, including: *C. cuspidate*, *C. renshanxiangiae*, *C. elongata* and *C. gymnogyna*. The shadow Blue mean homology = 100%, pink mean homology≥70%, turquoise mean homology≥50% and the dot indicated missing. The number on the right showed the sequence length. The position of the primers was indicated by the arrow. **b** Only PCR products of *CWN* had a 335 bp long sequence deletion, while those of other 292 *C. sinensis* did not. M: D2000 DNA molecular marker; **c** PCR products of *CWN*; Lane 1–22: PCR products of 22 examples of randomly selected cultivars from 292 different cultivars covering the majority of *C. sinensis* cultivars in China. CK: Control. All 292 PCR products were shown in Supplementary Fig. S1. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea)
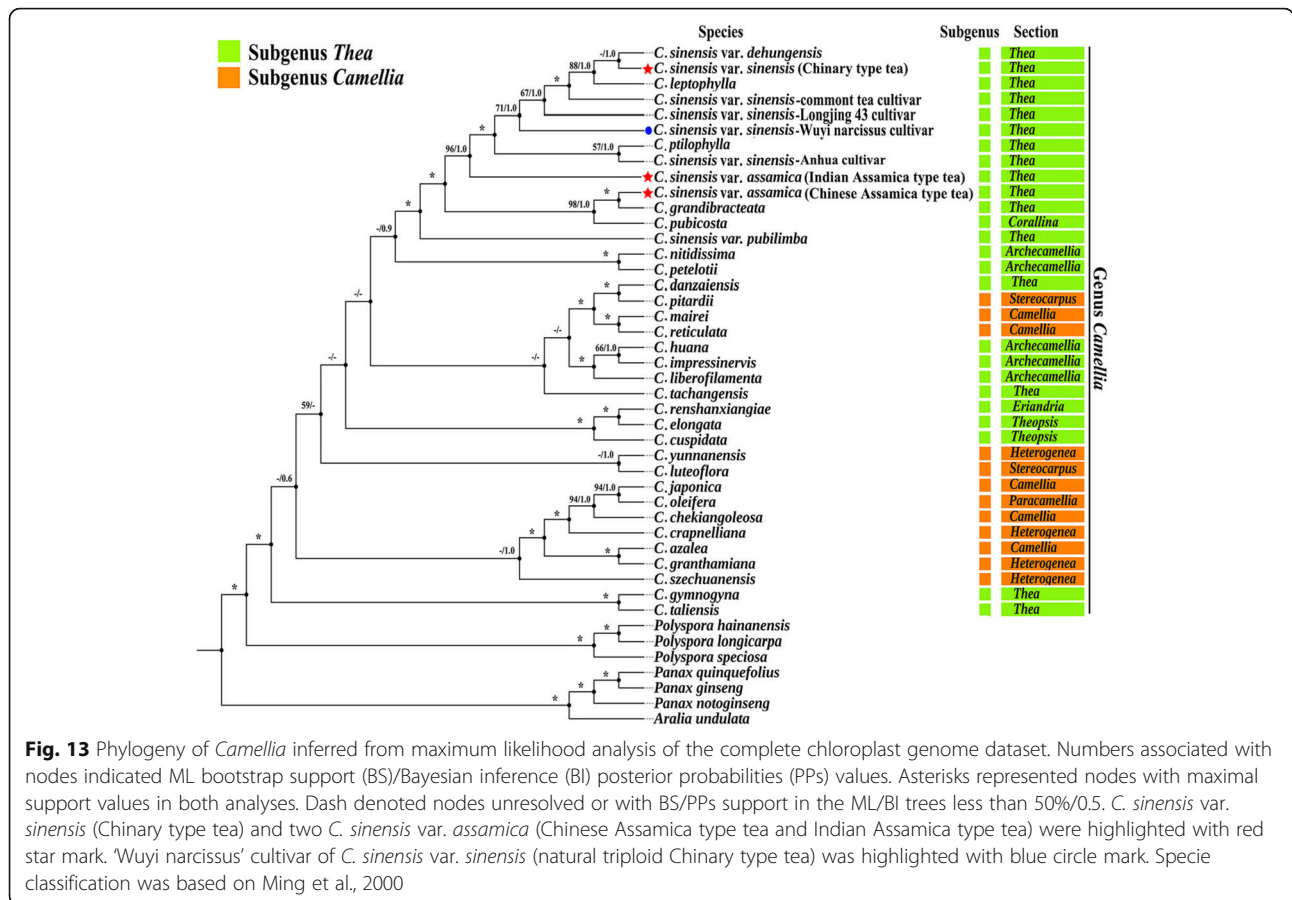
**Divergence time (Million year ago) of *Camellia* species inferred from the complete chloroplast genome data.**

| *Camellia* species | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.4 (0.1-1.1) | 2 | 0.8 (0.3-1.6) | 3 | 0.8 (0.4-1.5) | 4 | 1.2 (0.7-1.9) | 5 | 1.8 (1.1-2.7) | 6 | 3.0 (1.8-4.4) |
| 7 | 3.1 (2.1-4.3) | 8 | 3.4 (2.3-4.7) | 9 | 3.5 (2.4-4.8) | 10 | 4.2 (1.8-7.1) | 11 | 4.7 (2.8-7.1) | 12 | 4.9 (2.5-7.5) |
| 13 | 5.0 (2.7-7.5) | 14 | 6.1 (1.9-10.3) | 15 | 6.2 (4.4-8.1) | 16 | 6.2 (4.6-8.0) | 17 | 7.0 (5.3-8.8) | 18 | 8.0 (6.1-10.0) |
| 19 | 8.5 (6.6-10-6) | 20 | 9.2 (6.6-11.8) | 21 | 9.2 (6.9-11.5) | 22 | 9.3 (5.5-12.3) | 23 | 9.6 (7.3-11.9) | 24 | 9.8 (7.6-12.1) |
| 25 | 10.3 (8.1-12.5) | 26 | 11.2 (9.2-13.2) | 27 | 12.5 (10.2-14.6) | 28 | 12.6 (10.3-14.7) | 29 | 12.8 (10.4-14.9) | 30 | 13.0 (10.7-15.1) |
| 31 | 13.3 (10.9-15.4) | 32 | 13.4 (11.0-15.5) | 33 | 13.5 (11.1-15.6) | 34 | 14.2 (11.7-16.3) | 35 | 14.4 (11.9-16.5) | 36 | 14.7 (12.2-16.8) |
| **Other species** | | | | | | | | | | | |
| 37 | 0.8 (10.9-15.4) | 38 | 2.5 (1.3-4.6) | 39 | 9.4 (6.1-12.6) | 40 | 11.1 (4.5-19.6) | 41 | 12.0 (8.7-15.5) | 42 | 20.4 (17.7-22.3) | 43 | 93.5 (73.6-107.0) |

**Estimated divergence time correspond to the nodes was in order.**

**Fig. 12** Divergence times and topologies of maximum likelihood trees based on the complete chloroplast genome. Length of each branch was not shown in the two phylogenetic trees. *C. sinensis* var. *sinensis* (Chinary type tea) and two *C. sinensis* var. *assamica* (Chinese Assamica type tea and Indian Assamica type tea) were highlighted with red star mark. 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea) was highlighted with blue circle mark. The ages of stratigraphic boundaries were obtained from the International Chronostratigraphic Chart (Pl, Pliocene; Qu, Quaternary) [32], with a scale as million years ago (Mya). The divergence time of each node among *Camellia* species was showed in the table below

**Fig. 13** Phylogeny of *Camellia* inferred from maximum likelihood analysis of the complete chloroplast genome dataset. Numbers associated with nodes indicated ML bootstrap support (BS)/Bayesian inference (BI) posterior probabilities (PPs) values. Asterisks represented nodes with maximal support values in both analyses. Dash denoted nodes unresolved or with BS/PPs support in the ML/BI trees less than 50%/0.5. *C. sinensis* var. *sinensis* (Chinary type tea) and two *C. sinensis* var. *assamica* (Chinese Assamica type tea and Indian Assamica type tea) were highlighted with red star mark. 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea) was highlighted with blue circle mark. Specie classification was based on Ming et al., 2000

sequences between them in detail. We found some significant differences between Chinese tea and Indian tea by comparing the four cp genomes (the triploid Chinary type tea, the diploid Chinary type tea, the diploid Chinese Assamica type tea, and the diploid Indian Assamica type tea): (1) The three Chinese teas were more similar in IR/LSC boundary pattern, which was obvious difference from Indian tea (Fig. 5). (2) Chinese teas and Indian tea exhibited the distinct repeat patterns, Indian tea possessed high level of species-specific long repeats, with 18, while three Chinese teas only possessed 0–2, and most of the long repeats in three Chinese teas were shared (Fig. 7b). Similarly, the indel patterns also exhibited obviously different between Chinese teas and Indian tea, Indian tea possessed 49 species-specific indels patterns, accounting for 49% of the total number of indel loci, while only 5–16 indels were specific to each of the three Chinese teas (Fig. 8d). (3) In particular, the analysis showed that selection played an important role in the codon usage of cp genome in all four species (Figs. 9 and 10). However, 37 codons of the three Chinese teas had the identical RSCU value, accounting for 57.8% of the codons, all of which were different from that of Indian tea (Table 4). Codon use bias had long been used

to reflect the origin, evolution, and mutation mode of species or genes [45]. Therefore, it suggested that Chinese tea and Indian tea might have undergone different selections, as did Chinese Assamica type tea and Indian Assamica type tea.

## Analysis of possible effects of polyploid evolution on cp genome

CpDNA data could be used not only to identify the maternal origin of polyploidy, but also to explore the possible effects of polyploid evolution, such as introgressive hybridization [26]. Our phylogenetic analysis based on complete cp genome showed the triploid *CWN* was clustered with *C. sinensis* species together (Fig. 13), and the p-distance between *CWN* and *CSS* was only 0.00045 (Table 2), which demonstrated the maternal ancestor of triploid came from *C. sinensis* species. However, SCAR analysis using a 335 bp deletion of the intergenic spacers (trnE/trnT) in 292 cultivars of *C. sinensis* showed that none of the cultivars had this sequence deletion like triploid *CWN* (Fig. 11b, Supplementary Fig. S1). This could be explained by limited sample sizes, possible extinction of its progenitors and/or the occurrence of chloroplast transfer through hybridization with other *Camellia*

species. Chloroplast transfer (or 'chloroplast capture'), the introgression of a chloroplast from one species into another (introgressive hybridization) [46], has been reported [47–49]. In our study, none of the 292 individuals covering the majority of *C. sinensis* cultivars in China had a 335 bp long sequence deletion similar to triploid *CWN*, but it could be found in four another *Camellia* species, including: *C. cuspidate, C. renshanxiangiae, C. elongata* and *C. gymnogyna* (Fig. 11a). Therefore, the triploid *CWN* should have originated through allopolyploidization involving two parental species. The further field surveys and molecular studies about parental origin remain to be performed. Nevertheless, our study suggested chloroplast transfer occurred during the polyploidization in *C. sinensis*, that was, the progenitor of the triploid *CWN* (*C. sinensis*) itself captured the cp genome of another *Camellia* species when hybridization occurred where the ranges of the species overlap.

## Divergence time estimation and indication for the three different domestication

The controversy over the domestication origin of the tea plant has existed for a long time. Previous studies using nSSRs, SNPs and RADseq data found that there were three different gene pools of tea plant (Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea), considering the breeding history of Assamica type tea was not long enough to produce a new lineage, so suggesting that Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea were likely the result of three independent domestication events [3, 50, 51]. Another view, however, was that the three different gene pools was a far cry from three different domestications. It was also possible that Assamica type tea was domesticated only once and then brought from one region to another, allowing it to evolve separately in both regions [8]. To further understand the domestication origin of these three types of tea plant, in this study, we first used the complete cp genome to estimate the divergence time of tea plant (Fig. 12). The result showed the three types of tea plant (Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea) were in different branches in the phylogenetic tree, and diverged to each other during 0.8–6.2 Mya (CI: 0.3–8.1 Mya). Indian Assamica type tea diverged from the ancestor of Indian Assamica type tea and Chinese Assamica type tea about 6.2 Mya (CI: 4.4–8.1 Mya). Thereafter, Chinary type tea and Chinese Assamica type tea diverged to each other about 0.8 Mya (CI: 0.4–1.5 Mya), which was in good agreement with that previously estimated by Wei et al. using the collinear nuclear genes between Chinary type tea and Chinese Assamica type tea genomes, about 0.38–1.54 Mya [4]. These results suggested that the three types of tea plant diverged from each other much

earlier than the known domestication time (2737 BC) [52], thus supporting the hypothesis of the three different domestication origins [50, 51].

## Phylogenetic relationships within *Camellia*

*Camellia* is taxonomically and phylogenetically ranked as one of the most challengingly difficult taxa in plants, and the taxonomic classification of *Camellia* based on morphology has been controversial. In terms of morphological classification, Chang et al. classified the genus *Camellia* into 4 subgenera with 22 sections [53], while Ming et al. revised the classification of Chang and classified the genus *Camellia* into 2 subgenera with a total of 14 sections [54, 55]. Our phylogenetic analysis (Fig. 13) of *Camellia* did not agree with any of the current traditional classification methods used recently in *Camellia* taxonomy, such as: *C. danzaiensis* was clustered with *C. pitardii, C. mairei* and *C. reticulata* together ($BS_{ML}$ = 100%, $PP_{BI}$ = 1.00), suggesting *C. danzaiensis* might belong to subgenus *Camellia*, rather than subgenus *Thea*. Similarly, *C. pubicosta* was sister to *C. sinensis* var. *assamica* and *C. grandibracteata* ($BS_{ML}$ = 98%, $PP_{BI}$ = 1.00), suggesting this species probably should be classified into sect. *Thea* by Chang et al. [53], instead of being classified into sect. *Corallina*. These results were all consistent with two previous studies of *Camellia* [21, 56].

In addition, our phylogenetic tree constructed based on complete cp genome in this study showed some incongruence with the previous phylogeny study of *Camellia*. A previous phylogenetic analysis using RAPDs showed that *Camellia* species could be divided into 5 ovary groups and 3 ovary groups according to the number of their ovaries, while our phylogenetic relationships of *Camellia* species did not well follow the number of locule ovary. For example, *C. grandibracteata* with 5 ovaries and *C. sinensis* var. *assamica* with 3 ovaries were well supported as monophyletic ($BS_{ML}$ = 100%, $PP_{BI}$ = 1.00). Phylogenetic analysis was not consistent with the taxonomy according to ovarian number, which was also found in Huang et al. 's study [21]. However, in Huang et al. 's study, *C. ptilophylla* with 3 ovaries was clustered with some species with five ovaries (*C. tachangensis, C. kwangsiensis* and *C. crassicolumna* var. *crassicolumna*), which did not occur in our results. This might be because the cp gene sequences used in their phylogenetic analysis were incomplete cp genomes (83,585 to 83,835 bp). A recent study had emphasized that the extensive heterogeneity of nucleotide substitution rate among different plastid genes and among different functional groups of genes were likely contributing to phylogenetic ambiguity [57]. In our phylogenetic analysis of the data partitions of cp genome, the sequence variations in the six datasets showed different percentage variation (Supplementary Tab. S7). Although the phylogenetic trees from the six

datasets generated mostly similar topological structures, there were still some individual species showing different locations in the six trees, and the support values of nodes increased significantly with the increasing of the sequence length (Fig. 13, Supplementary Fig. S2, S3, S4, S5, S6). It suggested that attention should also be paid to the effects of such heterogeneity when functional genes or plastid fragments were used to study phylogenetic evolution of cp genome of *Camellia*. In addition, phylogenomic analysis also tend to suffer from the poor sampling [58], and the number of *Camellia* species sampled in Huang et al. 's study was 18, while ours was more, at 37. Therefore, further taxon sampling and more complete cp genomes of *Camellia* were needed to resolve the controversial taxonomy of *Camellia* in future study.

While the analyses of the complete cp genomes provided a feasible way to clarify relationships [59], it might still be insufficient to fully resolve all phylogenetic relationships [60, 61]. Our results suggested that this type of complete genome phylogenomic analyses would resolve many controversies and guide the way for phylogeny in *Camellia*. Since plastome was regarded as a linked single locus due to its uniparental inheritance and lack of sexual recombination, in future work researchers should also attach importance to the use of multilocus approaches (including nuclear genes and mitochondrial genes), so as to provide abundant and detailed molecular data for the systematic classification and evolutionary study of *Camellia*.

## Conclusion
In this study, we had found that the repeats and indels were two most important evolutionary dynamics contributed to the diversification of the cp genome, which were not mutually exclusive. Chinese tea and Indian tea exhibited significantly differences in the structural characteristic and the codon usage of the cp genome, suggesting they might have undergone different selection pressures. In addition, our result demonstrated that the chloroplast transfer occurred during the polyploidization. Further, the phylogenomic analysis combined with divergence time estimation implied that Chinary type tea, Chinese Assamica type tea and Indian Assamica type tea might have three different domestication origins, and the current classification of some *Camellia* species might need to be further discussed. Our data would not only provide insights into the chloroplast genome evolution of *C. sinensis* but also offer valuable information for taxonomic classification of *Camellia*.

## Methods
### Plant material and DNA extraction
Fresh leaves of 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* used in this study were obtained from the tea tree germplasm garden in Wuyi University. The voucher specimen (No. 20200317) was authenticated by Prof. Yongcong Hong (Wuyi University) and deposited in the Wuyi University herbarium.

Genomic DNA was extracted from leaves using CTAB extraction method. Final DNA quality was assessed by a NanoDrop spectrophotometer (Thermo Scientific, Carlsbad, CA, USA), and their integrity was examined by electrophoresis on a 0.8% agarose gel. DNA samples were preserved at − 80 °C at the Key Laboratory of Tea germplasm Genetic Resources of Wuyi University.

### Chloroplast genome sequencing, assembly and annotation
The previous studies had shown that cp genomes of closely related species share highly similarity and these sequences could be used as reference genomes to obtain the order of contigs for new cp genome assembly. This efficient strategy did not require organelle DNA isolation and had been well accepted by the scientific community [62–66]. Here, the three published cp genomes of tea plants (Accession number: KJ806281, MH019307 and MH460639) [21, 22, 62] were selected as reference genomes. A combined approach of PacBio sequencing complemented with Illumina sequencing was performed by Biozeron Biotechnology Co., Ltd. (Shanghai, China).

High quality total genomic DNA was applied to 500 bp paired-end library construction using the NEBNext Ultra DNA Library Prep Kit for Illumina sequencing. Sequencing was carried out on the Illumina NovaSeq 6000 platform. For PacBio sequencing, more than 5 μg of sheared and concentrated DNA was applied to size selection by the Blue Pippin (Sage Science, Beverly MA, USA). Approximately 20 kb SMRTbell libraries were prepared according to the manufacturer's instructions (PacBio, Menlo Park, CA, USA). The samples were sequenced on PacBio Sequel instrument using the P6 polymerase C4 chemistry combination. All raw reads were processed with Trimmomatic v0.39 software [67] to remove adapter sequences, short reads (length < 75 bp), and to trim low- quality bases (Q-value < 20). The Pac-Bio raw reads (polymerase reads) were filtered by discarding low-quality polymerase reads (Q-value < 0.80), short reads (length < 100 bp), short sub reads (length < 500 bp), and adapters. Both of Illumina reads and PacBio reads were mapped to the cp reference genomes of three closely related species to extract the cp reads. Short reads were aligned to the reference set using NCBI-BLAST-2.2.30 [68], and long reads were aligned to the reference set using Minimap2 [69]. Prior to the assembly, all filtered PacBio reads were error corrected with Illumina data using Pilon v1.22 [70]. All error-corrected PacBio reads were assembled into a single contig using Canu v2.0 [71]. To estimate the quality and coverage of the assembled genome [72], we mapped all Illumina and

PacBio reads to the assembly using BWA v0.7.12 [73], respectively.

The cp genes were annotated using CpGAVAS [74] and verified the sequence coordinates of each of the annotated genes using BLAST search against ref. cp genes. Annotation errors were manually corrected. The four junction regions (SC/IR) and the six randomly selected cp genes were furtherly validated with PCR-based conventional Sanger sequencing. The final annotated cp genome sequence was subjected to OGDRAW software [75] to generate the circular cp genome map and deposited to NCBI GenBank. The four junction regions (SC/IR) and 6 randomly selected cp genes were validated with PCR-based conventional Sanger sequencing (Supplementary Tab. S8).

### Comparative analysis of four cp genomes

Four cp genome sequences of *Camellia* species were aligned using MAFFT Version 7.017 [76]. Full alignments with annotation were visualized using the mVISTA software [77]. Number of nucleotide substitutions was calculated by MEGA 6.0 [78]. A sliding window analysis was conducted to compare π among the complete cp genomes, using DnaSP v5.0 [79]. The window length was 600 bp with a 200 bp step size. The percentage of variable characters for coding and noncoding regions in the genome was calculated as described previously [21]. The proportion of mutation events = [(NS + ID)/L] × 100, where NS = the number of nucleotide substitutions, ID = the number of indels, L = the aligned sequence length.

Repeat sequences were searched by REPuter [80] with a minimal size of 30 bp and > 90% identity (Hamming distance equal to 3) between the two repeats. Gap size between the repeats was restricted to a maximal length of 3 kb. Tandem repeats were identified by Tandem Repeats Finder (http://tandem.bu.edu/trf/trf.html) [81] with default settings. Simple sequence repeats (SSRs) were predicted using MISA [82] with the parameters: monomer (one nucleotide, $n \geq 8$), dimer (two nucleotides, $n \geq 4$), trimer (three nucleotides, n ≥ 4), tetramer (four nucleotides, $n \geq 3$), pentamer (five nucleotides, n ≥ 3), hexamer (six nucleotides, n ≥ 3). A combination of SSRs separated by the maximum distance of 100 bp was considered as an SSR (i.e. an SSR of the compound formation).

According to Ibrar et al.'s method [40], the correlations were analyzed in the pairwise comparisons between the three types of mutations: substitutions and indels, repeats and indels, and repeats and substitutions. Setting cp genome of *CSS* (Accession number: KJ806281) as a reference, indels and substitutions were counted in the comparison of 630 nonoverlapping bins each with a size of 250 bp.

### Codon usage analyses

In order to avoid sampling bias, each protein-coding gene (CDS) in cp genome were checked for being full-length and for the presence of proper start and stop codons. CDS < 300 bp were excluded in codon usage calculations [83].

ENc (Effective number of codons) and GC3s (GC content at the third synonymously variable coding position excluding Met and Trp) were calculated using CodonW v1.4.4 [84]. ENc value is a measure of general non-uniformity of usage within synonymous groups of codons, ranging from 20 (extreme bias where only one codon is used in each amino acid) to 61 (random codon usage) [85]. ENc plot analysis (ENc vs GC3s) was used to examine whether the codon usages were affected only by mutation or other factors. If codon usage is constrained only by mutation pressure, ENc value lie on or slightly below the expected curve, and if codon usage is subject to natural selection, ENc value will lie considerably below the expected curve [86].

Neutrality plot (GC12 vs. GC3) was used to investigate the effects of mutation pressure and natural selection on codon use patterns. GC12 (the average value of GC contents at the first and second positions of codon) and GC3 (the GC content at the third position) were calculated by Perl script. GC3 was calculated excluding the three termination codons (TAA, TAG and TGA) and the three codons for Ile (ATT, ATC and ATA). Meanwhile, two single codons for Met (ATG) and Trp (TGG) were also excluded in all three patterns [87]. The slope of the plot regression is zero which indicates that there is no effect on directional mutation pressure (complete selection constraint). Slope 1 indicates that the codon usage bias is completely affected by the directional mutation pressure, and represents complete neutrality [88].

RSCU (Relative synonymous codon usage) for CDS was calculated using CodonW v1.4.4 [84]. RSCU value for a particular codon refers to the ratio of its actual usage frequency to expected frequency when it is used without bias. The preferred codons with RSCU > 1.0 occur when they are used with higher frequencies than random, and the rare codons with RSCU < 1.0 means the opposite [89].

### Analysis of cp sequence characterized amplified region (SCAR)

A sequence part of the chloroplast trnE/trnT intergenic spacer, where a 335 bp long deletion was observed in triploid *CWN* by comparing with another three tea species, was used to SCAR analysis. Screening of the distribution of the similar deletion in the trnE/trnT intergenic spacer was conducted by PCR amplification of the respective chloroplast region in 292 individuals covering the majority of *C. sinensis* cultivars in China (Supplementary Tab. S5).

All individual samples were collected from the tea germplasm resource garden of Wuyi University and the DNA extraction method was described above. PCR products were observed on 1.5% agarose gels against a D2000 DNA molecular marker. A set of the primers (trnE: 5′-TCCTGAACCACTAGACGATG-3′; trnT: 5′-ATGGCG TTACTCTACCACTG-3′) were designed in conserved regions on either side of the 335 bp long deletion region by comparing cp genome sequences.

### Phylogenetic analysis and divergence time estimation

The complete cp genome sequences of 37 *Camellia* species, 3 *Polyspora* species, 3 *Panax* species and the outgroup (*Aralia undulata*) (Supplementary Tab. S6) were aligned with the program MAFFT version [76]. The ambiguously aligned loci (e.g., 'N' or 'K') were excluded from the analyses and the poorly aligned regions were removed from the complete plastome dataset using Gblocks v0.91b [90]. Maximum likelihood (ML) analyses were implemented in RAxML version 7.2.6 [91]. Nonparametric bootstrapping test was implemented in the "fast bootstrap" algorithm of RAxML with 1000 replicates. Bayesian inference (BI) analyses were performed using the program MrBayes version 3.1.2 [92]. The best-fitting models were determined by the Akaike Information Criterion [93] as implemented in the program Modeltest 3.7 [94]. The bootstrap value above 70% ($BS_{ML} > 70\%$) for Maximum Likelihood and the posterior probability above 0.95 ($PP_{BI} > 0.95$) for Bayesian Inference were used to determine a well-supported clade [95]. ML and BI analyses were performed based on the following six datasets: (1) the complete cp DNA sequences, (2) the large single copy region (LSC), (3) the small single copy region (SSC), (4) the inverted repeat region (IR), (5) a set of the common protein coding genes (PCGs) and (6) a set of the common non-coding genes (Non-PCGs). The best-fit models for each data set were showed in Supplementary Tab. S7.

The divergence times of *Camellia* species were calculated using mcmctree of PAML [96]. The empirical divergence times of *P. ginseng*/*P. quinquefolius* (0.8–1.2 Mya), *P. ginseng*/*P. notoginseng* (3.5–5.2 Mya) [97–99] were assigned to constrain the age of the *Panax*. Yu et al. [100] have demonstrated the influence of fossil calibration and divergence time estimation in *Theaceae*, and tested safety and risk scenarios of 9 scenarios based on a complete cp genome phylogenetic framework. Based on the results of Yu et al., the divergence times of *Polyspora* species/*Camellia* species (6.9–21.2 Mya) were assigned to constrain the stem of *Camellia* family. A Birth-Death prior on branching rates was employed and three independent analyses were run for 10,000 generations.

## Supplementary Information

The online version contains supplementary material available at https://doi.org/10.1186/s12864-021-07427-2.

**Additional file 1: Supplementary Fig. S1.** SCAR analysis was conducted by PCR amplification of the respective chloroplast region in 292 different cultivars of *Camellia sinensis*. PCR amplification to screen the distribution of a 335 bp deletion of the intergenic spacers (trnE/trnT) in 292 different cultivars of *Camellia sinensis*. None of 292 cultivars had this sequence deletion like triploid *CWN*. M: D2000 DNA molecular marker; C: PCR products of *CWN*; Lane 1–292: PCR products of 292 different cultivars covering the majority of *C. sinensis* cultivars in China. CK: Control. The corresponding 292 cultivars were shown in Supplementary Tab. S5.

**Additional file 2: Supplementary Fig. S2.** Phylogeny of *Camellia* inferred from maximum likelihood (ML) analysis of LSC. Numbers associated with nodes indicated ML bootstrap support (BS)/Bayesian inference (BI) posterior probabilities (PPs) values. Asterisks represented nodes with maximal support values in both analyses. Dash denoted nodes unresolved or with BS/PPs support in the ML/BI trees less than 50%/0.5. *C. sinensis* var. *sinensis* (Chinary type tea) and two *C. sinensis* var. *assamica* (Chinese Assamica type tea and Indian Assamica type tea) were highlighted with star mark.

**Additional file 3: Supplementary Fig. S3.** Phylogeny of *Camellia* inferred from maximum likelihood (ML) analysis of SSC. The meaning of the figures and symbols was consistent with Fig. S2.

**Additional file 4: Supplementary Fig. S4.** Phylogeny of *Camellia* inferred from maximum likelihood (ML) analysis of IR. The meaning of the figures and symbols was consistent with Fig. S2.

**Additional file 5: Supplementary Fig. S5.** Phylogeny of *Camellia* inferred from maximum likelihood (ML) analysis of PCGs. The meaning of the figures and symbols was consistent with Fig. S2.

**Additional file 6: Supplementary Fig. S6.** Phylogeny of *Camellia* inferred from maximum likelihood (ML) analysis of Non-PCGs. The meaning of the figures and symbols was consistent with Fig. S2.

**Additional file 7: Supplementary Tab. S1.** List of genes in the cp genome. * Genes containing a single introns; ** Genes containing two introns.

**Additional file 8: Supplementary Tab. S2.** Distribution of the different and the shared SSRs in four *Camellia* chloroplast genomes. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea). Symbol "/" indicated missing.

**Additional file 9: Supplementary Tab. S3.** A list of repeated sequences and their locations identified in four *Camellia* chloroplast genomes. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea).

**Additional file 10: Supplementary Tab. S4.** Distribution of indel in four chloroplast genomes. *CWN*: 'Wuyi narcissus' cultivar of *C. sinensis* var. *sinensis* (natural triploid Chinary type tea); *CSS*: *C. sinensis* var. *sinensis* (diploid Chinary type tea); *CSA*: *C. sinensis* var. *assamica* (diploid Chinese Assamica type tea); *CIA*: *C. sinensis* var. *assamica* (diploid Indian Assamica type tea).

**Additional file 11: Supplementary Tab. S5.** Two hundred ninety two different cultivars of *Camellia sinensis* in China adopted for SCAR analysis in this study. Two hundred ninety two different cultivars represented the majority of *Camellia sinensis* varieties from 14 provinces in China.

**Additional file 12: Supplementary Tab. S6.** Chloroplast genomes of species adopted in this study. No. 1–37: 37 *Camellia* species; No. 38–40: 3 *Polyspora* species; No. 41–43: 3 *Panax* species; No. 44: Outgroup (*Aralia undulata*); No. 3, 13 and 17: 3 *Camellia sinensis* chloroplast reference genomes.

**Additional file 13: Supplementary Tab. S7.** Variable sites in 37 *Camellia* chloroplast genomes and the best-fitting models for phylogenetic analysis.

**Additional file 14: Supplementary Tab. S8.** Primers used for assembly and junction verification of the chloroplast genome.

Li *et al. BMC Genomics*        (2021) 22:138

Page 20 of 22

## Author details
[1]College of Tea and Food Science, Wuyi University, 358# Baihua Road, Wuyishan 354300, China. [2]College of Mathematics and Computer Science, Wuyi University, 358# Baihua Road, Wuyishan 354300, China.

## References
1. Lu H, Jiang W, Ghiassi M, Lee S, Nitin M. Classification of *Camellia* (*Theaceae*) species using leaf architecture variations and pattern recognition techniques. PLoS One. 2012;7:e29704. https://doi.org/10.1371/journal.pone.0029704.
2. Meegahakumbura MK. Genetic assessment of Asian tea Germplasm and the domestication history of the tea plant (*Camellia sinensis*). PhD dissertation. Beijing: University of Chinese Academy of Sciences; 2016.
3. Yang H, Wei CL, Liu HW, Wu JL, et al. Genetic divergence between *Camellia sinensis* and its wild relatives revealed via genome wide SNPs from RAD sequencing. PLoS One. 2016;11:e0151424.
4. Wei C, Yang H, Wang S, et al. Draft genome sequence of *Camellia sinensis var. sinensis* provides insights into the evolution of the tea genome and tea quality [J]. Proc Natl Acad, U S A. 2018;115(18):201719622.
5. Mondal TK, Bhattacharya A, Laxmikumaran M, Ahuja PS. Recent advances in tea (*Camellia sinensis*) biotechnology. Plant Cell Tissue Organ Cult. 2004;76: 195–254.
6. Mondal TK. Breeding and biotechnology of tea and its wild species. Springer Science & Business Media. 2014:1–167. https://doi.org/10.1007/978-81-322-1704-6.
7. Prathibhani C, Kumarihami H, Eun U, et al. Comparative study on cross-compatibility between *Camellia sinensis* var. *sinensis* (China type) and *C. sinensis* var. *assamica* (Assam type) tea. Afr J Agric Res. 2016;11:1092–101. https://doi.org/10.5897/AJAR2015.9951.
8. Drew L. The growth of tea. Genetic studies of today's tea trees are providing clues to how the plant was first domesticated. Nature Outlook. 2019;566:S2–4. https://doi.org/10.1038/d41586-019-00395-4.
9. Wicke S, Schneeweiss GM, de Pamphilis CW, Muller KF, Quandt D. The evolution of the plastid chromosome in land plants: gene content, gene order, gene function. Plant Mol Biol. 2011;76:273–97.
10. Barrett CF, Freudenstein JV, Li J, et al. Investigating the path of plastid genome degradation in an early-transitional clade of heterotrophic orchids, and implications for heterotrophic angiosperms. Mol Biol Evol. 2014;31:3095–112.
11. Molina J, Hazzouri KM, Nickrent D, et al. Possible loss of the chloroplast genome in the parasitic flowering plant *Rafflesia lagascae* (*Rafflesiaceae*). Mol Biol Evol. 2014;31:793–803.
12. Li Y, Zhou JG, Chen XL, et al. Gene losses and partial deletion of small single copy regions of the chloroplast genomes of two *hemiparasitic Taxillus* species. Sci Rep. 2017;7:12834. https://doi.org/10.1038/s41598-017-13401-4.
13. Guisinger MM, Kuehl JV, Boore JL, Jansen RK. Extreme reconfiguration of plastid genomes in the angiosperm family *Geraniaceae*: rearrangements, repeats, and codon usage. Mol Biol Evol. 2011;28:583–600. https://doi.org/10.1093/molbev/msq229.
14. Jansen RK, Saski C, Lee S-B, Hansen AK, Daniell H. Complete plastid genome sequences of three rosids (*Castanea*, *Prunus*, *Theobroma*): evidence for at least two independent transfers of rpl22 to the nucleus. Mol Biol Evol. 2011; 28:835–47.
15. Ma PF, Zhang YX, Guo ZH, Li DZ. Evidence for horizontal transfer of mitochondrial DNA to the plastid genome in a bamboo genus. Sci Rep. 2015;5:11608–16.
16. Cosner ME, Raubeson LA, Jansen RK. Chloroplast DNA rearrangements in *Campanulaceae*: phylogenetic utility of highly rearranged genomes. BMC Evol Biol. 2004;4:27–43.
17. Kugita M, et al. The complete nucleotide sequence of the hornwort (*Anthoceros formosae*) chloroplast genome: insight into the earliest land plants. Nucleic Acids Res. 2003;31:716–21.
18. McKinnon G. Reticulate evolution in higher plants. In: Henry R, editor. Plant diversity and evolution. Wallingford: CABI publishing; 2004. p. 81–96.
19. Wendel JF, Doyle JJ. Polyploidy and evolution in plants. In: Henry RJ, editor. Plant diversity and evolution; 2005. p. 97–117.
20. Yamane K, Yasui Y, Ohnishi O. Intraspecific cpDNA variations of diploid and tetraploid perennial buckwheat, *Fagopyrum cymosum* (*Polygonaceae*). Am J Bot. 2003;90:339–46.
21. Huang H, Shi C, Liu Y, et al. Thirteen *Camellia* chloroplast genome sequences determined by high-throughput sequencing: genome structure and phylogenetic relationships [J]. BMC Evol Biol. 2014;14(1):151.
22. Zhang F, Li W, Gao C, et al. Deciphering tea tree chloroplast and mitochondrial genomes of *Camellia sinensis* var *assamica*. Sci Data. 2019;6: 209. https://doi.org/10.1038/s41597-019-0201-8.
23. Fawcett JA, Yves VDP. Angiosperm polyploids and their road to evolutionary success. Trends Evol Biol 2010; 2:16–21.
24. Soltis DE, Albert VA, Leebens-Mack J, et al. Polyploidy and angiosperm diversification. Am J Bot. 2009;96:336–48.
25. Pillay M, Hilu KW. Chloroplast DNA variation in diploid and polyploid species of Bromus (*Poaceae*) subgenera *Festucaria* and *Ceratochloa* [J]. Theor Appl Genet. 1990;80(3):326–32.

26. Soltis DE, Soltis PS. 7 - chloroplast DNA and nuclear rDNA variation: insights into autopolyploid and allopolyploid evolution. In: Shoichi K, editor. Biological approaches and evolutionary trends in plants: Academic; 1990. p. 97–117. ISBN 9780124029606. https://doi.org/10.1016/B978-0-12-402960-6.50012-3.

27. Rousseau-Gueutin M, Bellot S, Martin G, et al. The chloroplast genome of the hexaploid *Spartina maritima* (*Poaceae*, *Chloridoideae*): comparative analyses and molecular dating. Mol Phylogenet Evol. 2015;93:5–16.

28. Jiang LY, Qian ZQ, Guo ZG, et al. Polyploid origins in *Gynostemma pentaphyllum* (*Cucurbitaceae*) inferred from multiple gene sequences [J]. Molecular Phylogenetics & Evolution. 2009;52(1):183–91.

29. Nemati Z, Dörte H, Ciftci A, et al. Saffron (*Crocus sativus*) is an autotriploid that evolved in Attica (Greece) from wild Crocus cartwrightianus. Mol Phylogenet Evol. 136. https://doi.org/10.1016/j.ympev.2019.03.022.

30. Palmer JD, Shields CR, Cohen DB, Orton TJ. Chloroplast DNA evolution and the origin of amphidiploid *Brassica*. Theor Appl Genet. 1983;65:181–9.

31. Liang YR, Liu ZS. Study on chromosome number and karyotype of five tea Clonal varieties. J Tea Sci (In Chinese). 1988;2.

32. Cohen K, Finney S, Gibbard P, Fan JX. The ICS international chronostratigraphic chart. Episodes. 2013;36:199–204.

33. Gaudeul M, Giraud T, Kiss L, Shykoff JA. Nuclear and chloroplast microsatellites show multiple introductions in the worldwide invasion history of common ragweed, *Ambrosia artemisiifolia*. PLoS One. 2011;6(3):e17658.

34. Cavalier-Smith T. Chloroplast evolution: secondary symbiogenesis and multiple losses. Curr Biol. 2002;12:R62–4.

35. Timme RE, Kuehl JV, Boore JL, Jansen RK. A comparative analysis of the *Lactuca* and *Helianthus* (*Asteraceae*) plastid genomes: identification of divergent regions and categorization of shared repeats. Am J Bot. 2007;94: 302–12.

36. Britten RJ, Rowen L, Williams J, Cameron RA. Majority of divergence between closely related DNA samples is due to indels. Proc Natl Acad Sci U S A. 2003;100:4661–5.

37. Dubey H, Rawal H, Rohilla M, et al. TeaMiD: a comprehensive database of simple sequence repeat markers of tea. Database (*Oxford*). 2020;2020: baaa013. https://doi.org/10.1093/database/baaa013.

38. Pyo YJ, Kwon KC, Kim A, Cho MH. Seedling Lethal1, a pentatricopeptide repeat protein lacking an E/EC or DYW domain in Arabidopsis, is involved in plastid gene expression and early chloroplast development. Plant Physiol. 2013;163:1844–58. https://doi.org/10.1104/pp.113.227199.

39. Bapteste E, Philippe H. The potential value of indels as phylogenetic markers: position of Trichomonads as a case study. Mol Biol Evol. 2002;19:972–7.

40. Ibrar A, Biggs PJ, Matthews PJ, et al. Mutational dynamics of aroid chloroplast genomes. Genome Biol Evol. 2012;4(12):1316–23. https://doi.org/10.1093/gbe/evs110.

41. Tian D, Qiang W, Zhang P, et al. Single-nucleotide mutation rate increases close to insertions/deletions in eukaryotes. Nature. 2008;455:105–8.

42. Zhu L, Wang Q, Tang P, Araki H, Tian D. Genome-wide association between insertions/deletions and the nucleotide diversity in bacteria. Mol Biol Evol. 2009;26:2353–61.

43. McDonald MJ, WangW-C HH-D, Leu J-Y. Clusters of nucleotide substitutions and insertion/deletion mutations are associated with repeat sequences. PLoS Biol. 2011;9:e1000622.

44. Plunkett GM, Downie SR. Expansion and contraction of the chloroplast inverted repeat in Apiaceae subfamily Apioideae [J]. Syst Bot. 2000;25(4):648–67.

45. Shackelton LA, Parrish CR, Holmes EC. Evolutionary basis of codon usage and nucleotide composition bias in vertebrate DNA viruses. J Mol Evol. 2006;62:551–63.

46. Weitzman JB. Chloroplast transfer [J]. Genome Biol. 2003;4:spotlight-20030206-01. https://doi.org/10.1186/gb-spotlight-20030206-01.

47. Soltis DE, Soltis PS, Collier TG, et al. Chloroplast DNA variation within and among genera of the Heuchera group (*Saxifragaceae*): evidence for chloroplast transfer and Paraphyly [J]. Am J Bot. 1991;78(8):1091–112.

48. Sidorov VA, Menczel L, Nagy F, et al. Chloroplast transfer in Nicotiana based on metabolic complementation between irradiated and iodoacetate treated protoplasts [J]. Planta. 1981;152(4):341–5.

49. Naciri Y, Manen JF. Potential DNA transfer from the chloroplast to the nucleus in Eryngium alpinum [J]. Mol Ecol Resour. 2010;10(4):728–31.

50. Meegahakumbura MK, Wambulwa MC, Thapa KK, et al. Indications for three independent domestication events for the tea plant (*Camellia sinensis* (L.) O. Kuntze) and new insights into the origin of tea Germplasm in China and India revealed by nuclear microsatellites. PLoS One. 2016;11(5):e0155369. https://doi.org/10.1371/journal.pone.0155369.

51. Meegahakumbura M, Wambulwa M, Li DZ, Gao LM. Preliminary investigations on the genetic relationships and origin of domestication of the tea plant (*Camellia sinensis* (L.)) using genotyping by sequencing. Tropical Agricultural Research. 2018;29(3):229–40.

52. Hara Y, Liu SJ, Wickermasinghe RL, et al. Special issue on tea. Food Rev Int (USA). 1995;11:371–545.

53. Chang HD, Ren SX. Flora of China. Science Press Tomus. 1998;49(3):1–251.

54. Ming TL. Monograph of the genus camellia. Kunming: Yunnan Science and Technology Press; 2000.

55. Ming TL, Bruce B. Flora of China. Beijing: Science Press; 2010.

56. Yang JB, Yang SX, Li HT, et al. Comparative chloroplast genomes of *Camellia* species. PLoS One. 2013;8:e73053. https://doi.org/10.1371/journal.pone.0073053.

57. Zhang X, Sun YX, Landis J, Lv ZY, et al. Plastome phylogenomic study of *Gentianeae* (*Gentianaceae*): widespread gene tree discordance and its association with evolutionary rate heterogeneity of plastid genes. BMC Plant Biol. 2020;20. https://doi.org/10.1186/s12870-020-02518-w.

58. Peng ZH, Lu TT, Li LB, et al. Genome-wide characterization of the biggest grass, bamboo, based on 10,608 putative full-length cDNA sequences. BMC Plant Biol. 2010;10:116–29.

59. Li X, Yang Y, Henry RJ, Rossetto M, Wang Y, Chen S. Plant DNA barcoding: from gene to genome. Biol Rev Camb Philos Soc. 2015;90(1):157–66. https://doi.org/10.1111/brv.12104.

60. Wortley AH, Rudall PJ, Harris DJ, Scotland RW. How much data are needed to resolve a difficult phylogeny? case study in Lamiales. Syst Biol. 2005;54: 697–709. https://doi.org/10.1080/10635150500221028.

61. Petersen G, Aagesen L, Seberg O, Larsen IH. When is enough, enough in phylogenetics? A case in point from Hordeum (Poaceae). Cladistics. 2011;27: 428–46. https://doi.org/10.1111/j.1096-0031.2011.00347.x.

62. Rawal HC, Kumar PM, Bera B, Singh NK, Mondal TK. Decoding and analysis of organelle genomes of Indian tea (Camellia assamica) for phylogenetic confirmation. Genomics. 2020;112(1):659–68. https://doi.org/10.1016/j.ygeno.2019.04.018.

63. Cronn R, Liston A, Parks M, et al. Multiplex sequencing of plant chloroplast genomes using Solexa sequencing-by-synthesis technology. Nucleic Acids Res. 2008;36:e122. https://doi.org/10.1093/nar/gkn502.

64. Moore MJ, Dhingra A, Soltis PS, Shaw R, et al. Rapid and accurate pyrosequencing of angiosperm plastid genomes. BMC Plant Biol. 2006;6:17. https://doi.org/10.1186/1471-2229-6-17.

65. Tangphatsornruang S, Sangsrakru D, Chanprasert J, et al. The chloroplast genome sequence of mungbean (Vigna radiata) determined by high-throughput pyrosequencing: structural organization and phylogenetic relationships. DNA Res. 2009;17:11–22. https://doi.org/10.1093/dnares/dsp025.

66. Wang X, Cheng F, Rohlsen D, et al. Organellar genome assembly methods and comparative analysis of horticultural plants. Hortic Res. 2018;5:3. https://doi.org/10.1038/s41438-017-0002-1.

67. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics. 2014;30(15):2114–20.

68. Acland A, Agarwala R, Barrett T, et al. Database resources of the National Center for Biotechnology Information Nucleic Acids Res. Nucleic Acids Res. 2014;42:D7-17.

69. Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics. 2018;34(18):3094–100.

70. Walker BJ, Abeel T, Shea T, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. PLoS ONE. 2014;9(11):e112963.

71. Koren S, Walenz BP, Berlin K, et al. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome Res. 2017;27(5):722–36. https://doi.org/10.1101/gr.215087.116.

72. Twyford AD, Ness RW. Strategies for complete plastid genome sequencing[J]. Mol Ecol Resour. 2017;17:858–68.

73. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler Transform. Bioinformatics. 2009;25:1754–60.

74. Liu C, Shi LC, Zhu YJ, et al. CpGAVAS, an integrated web server for the annotation, visualization, analysis, and GenBank submission of completely sequenced chloroplast genome sequences[J]. BMC Genomics. 2012;13:715.

75. Lohse M, Drechsel O, Bock R. Organellar Genome DRAW (OGDRAW): a tool for the easy generation of high quality custom graphical maps of plastid and mitochondrial genomes [J]. Curr Genet. 2007;52(5–6):267–74.

76. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. Mol Biol Evol. 2013;30:772–80.

77.  Frazer KA, Pachter L, Poliakov A, Rubin EM, Dubchak I. VISTA: computational tools for comparative genomics. Nucleic Acids Res. 2004;32(Suppl. 2):W273–9.

78.  Tamura K, Stecher G, Peterson D, Filipski A, Kumar S. MEGA6: molecular evolutionary genetics analysis version 6.0. Mol Biol Evol. 2013;30:2725–9.

79.  Librado P, Rozas J. DnaSP v5: a software for comprehensive analysis of DNA polymorphism data. Bioinformatics. 2009;25(11):1451–2.

80.  Kurtz S, Choudhuri JV, Ohlebusch E, Schleiermacher C, Stoye J, Giegerich R. REPuter: the manifold applications of repeat analysis on a genomic scale. Nucleic Acids Res. 2001;29:4633–42.

81.  Benson G. Tandem repeats finder: A program to analyze DNA sequences. Nucleic Acids Res. 1999;27:573–80.

82.  Beier S, Thiel T, Münch T, Scholz U, Mascher M. MISA-web: a web server for microsatellite prediction. Bioinformatics. 2017;33:2583–5.

83.  Rosenberg MS, Subramanian S, Kumar S. Patterns of transitional mutation biases within and among mammalian genomes. Mol Biol Evol. 2003;20:988–93.

84.  Peden JF. Analysis of codon usage. Ph. D. Thesis. Nottingham: University of Nottingham.;1999.

85.  Wu Y, Li Z, Zhao D, Tao J. Comparative analysis of flower-meristem-identity gene APETALA2 (AP2) codon in different plant species. J Integr Agric. 2018; 17:867–77. https://doi.org/10.1016/S2095-3119(17)61732-5.

86.  Wright F. The 'effective number of codons' used in a gene. Gene. 1990;87: 23–9.

87.  Sueoka N. Directional mutation pressure and neutral molecular evolution. Proc Natl Acad Sci U SA. 1988;85:2653–7. https://doi.org/10.1073/pnas.85.8.2653.

88.  Wen Y, Zou Z, Li H, Xiang Z, He N. Analysis of codon usage patterns in Morus notabilis based on genome and transcriptome data. Genome. 2017; 60:473–84. https://doi.org/10.1139/gen-2016-0129.

89.  Gupta SK, Bhattacharyya TK, Ghosh TC. Synonymous codon usage in Lactococcus lactis: mutational bias versus translational selection. J Biomol Struct Dyn. 2004;21:527–36.

90.  Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood [J]. Syst Biol. 2003;52(5):696–704.

91.  Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. Bioinformatics. 2006;22:2688–90.

92.  Ronquist F, Huelsenbeck JP. MrBayes 3: Bayesian phylogenetic inference under mixed models. Bioinformatics. 2003;19:1572–4.

93.  Posada D, Buckley TR. Model Selection and model averaging in phylogenetics: advantages of Akaike Information Criterion and Bayesian approaches over likelihood ratio tests. Syst Biol. 2004;53:793–808.

94.  Posada D, Crandall KA. Modeltest: testing the model of DNA substitution. Bioinformatics. 1998;14:817–8.

95.  Alfaro ME, Holder MT. The posterior and the prior in Bayesian phylogenetics. Annu Rev Ecol Evol Syst. 2006;37:19–42.

96.  Xu B, Yang Z. PAMLX: a graphical user interface for PAML. Mol Biol Evol. 2013;30:2723–4.

97.  Choi HI, Kim NH, Lee J, et al. Evolutionary relationship of Panax ginseng and P. quinquefolius inferred from sequencing and comparative analysis of expressed sequence tags. Genet Resour Crop Evol. 2013;60:1377–87.

98.  Choi HI, Kim NH, Kim JH, et al. Development of reproducible EST-derived SSR markers and assessment of genetic diversity in Panax ginseng cultivars and related species. J Ginseng Res. 2011;35:399–412.

99.  Shi F, Li M, Li Y, et al. The impacts of polyploidy, geographic and ecological isolations on the diversification of Panax (Araliaceae). BMC Plant Biol. 2015; 15:297. https://doi.org/10.1186/s12870-015-0669-0.

100. Yu XQ, Gao LM, Soltis DE, et al. Insights into the historical assembly of East Asian subtropical evergreen broadleaved forests revealed by the temporal history of the tea family [J]. New Phytologist. 2017;215(3):1235.

## Publisher's Note