

# Comparison of Chloroplast Genomes and Phylogenetic Analysis of Four Species in *Quercus* section *Cyclobalanopsis*

Xiaoli Chen

China West Normal University

Xuemei Zhang (✉ [zhangmei103127@sina.com](mailto:zhangmei103127@sina.com))

China West Normal University


---

## Article

**Keywords:**

**Posted Date:** July 17th, 2023

**DOI:** <https://doi.org/10.21203/rs.3.rs-3009025/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

**Additional Declarations:** No competing interests reported.

---

**Version of Record:** A version of this preprint was published at Scientific Reports on October 31st, 2023. See the published version at <https://doi.org/10.1038/s41598-023-45421-8>.

## Abstract

The identification in *Quercus* L. species was considered to be difficult all the time. The fundamental phylogenies of *Quercus* have already been discussed by morphological and molecular means. However, the morphological characteristics of some *Quercus* groups may not be consistent with the molecular results (such as the group *Helferiana*), which may lead to blurring of species relationships and prevent further evolutionary researches. To understand the interspecific relationships and phylogenetic positions, we sequenced and assembled the CPGs (160715 bp ~ 160842 bp) of four *Quercus* section *Cyclobalanopsis* species by Illumina pair-end sequencing. The genomic structure, GC content and IR/SC boundaries exhibited significant conservatism. Six highly variable hotspots were detected in comparison analysis, among which *rpoC1*, *clpP* and *ycf1* could be used as molecular markers. Besides, two genes (*petA*, *ycf2*) were detected to be under positive selection pressure. The phylogenetic analysis showed: *Trigonobalanus* genus and *Fagus* genus located at the base of the phylogeny tree; the *Quercus* genus were distinguished to two clades, including five sections. All CTB species clustered into a single branch, which was in accordance with the results of the morphological studies. But neither of group *Gilva* nor group *Helferiana* had formed a monophyly. Six CTB species gathered together in pairs to form one branch respectively (*Quercus kerrii* and *Quercus chungii*; *Quercus austrocochinchinensis* with *Quercus gilva*; *Quercus helferiana* and *Quercus rex*). Due to a low support rate (0.338) in the phylogeny tree, the interspecific relationship between the two branches differentiated by this node remained unclear. We believe that *Q. helferiana* and *Q. kerrii* can exist as independent species due to their distance in the phylogeny tree. Our study provided genetic information in *Quercus* genus, which could be applied to further studies in taxonomy and phylogenetics.

## Introduction

*Quercus* L. is the most diverse genus in Fagaceae, with 430 species worldwide, which is one of the most widely distributed woody genera in Northern Hemisphere. Based on the morphology, molecular, and evolutionary history researches, *Quercus* genus was separated into 2 subgenera, namely *Quercus* and *Cerris*, including 8 sections<sup>1; 2</sup>. China is the second center of oak diversity and has identified and utilized the Fagaceae plants for the first time. *Quercus* section *Cyclobalanopsis* (ca. 150 species) mainly distributed in tropical and subtropical regions in Asia, which was divided into 6 groups by morphological features<sup>3</sup>, and *Quercus austrocochinchinensis*, *Quercus kerrii*, and *Quercus rex* were considered to belong to the group *Helferiana* inside. The three species were clustered to a branch based on leaf epidermal features, but when using RAD-Seq data, they were dispersive and did not represent as monophyletic<sup>4; 5</sup>, suggesting that the phylogeny location of these three species remained doubtful. In addition, *Quercus helferiana* showed high similarity with *Q. kerrii* in morphology. Some scholars believed that these two species should be classified as the same species, but subsequent research found that the similarity of these two species is inconsistent in different populations, the kinship between these two species therefore remains to be studied.

While it is a consensus for a long time that characters of the abaxial leaf surface and pollen provide valuable information for the species delimitation at infrageneric level<sup>5-7</sup>. However, when molecular sequence data were used to recognize (sub)sections/series, the result do not always conform to groups identified by means of traditional morphological classification within oaks<sup>8</sup>. For example, the research based on ITS sequences indicate that the species of compound trichome base (CTB) group of *Quercus* section *Cyclobalanopsis* converge into the same branch and with *Quercus* section *Cerris*, which is greatly different from the traditional classification of morphology<sup>3</sup>. Due to the similarities of leaf characteristics and gene introgression among different groups, despite lots of studies on morphological characteristics of *Quercus* section *Cyclobalanopsis*, more molecular evidence is needed for interspecific relationship and infrageneric phylogenetic status within *Quercus* genus.

Plastid exhibits key functions in plant growth and photosynthesis, and had independent genetic material, manifesting a tetrad structure<sup>9-11</sup>. Due to the maternal inheritance of the chloroplast genome (CPG), which was smaller in size, lower in nucleotide substitution rate, and more stable in structure compared to the nuclear genome, it exhibited conservative genetic variation<sup>12-15</sup>. CPGs can significantly enhance resolution at lower taxonomic levels and facilitate recovery of monophyletic lineages<sup>16</sup>, and are therefore considered ideal material in phylogenetics and population genetics<sup>13; 17-19</sup>. In recent years, DNA sequencing technology has shifted to high-throughput, and CPGs of a large number of plants have been sequenced and published<sup>20</sup>, which was in turn used to identification and classification of plant<sup>21-23</sup>, lineage geography<sup>24</sup> and phylogenetic relationship researches<sup>25-27</sup>. Due to the existence of overlap and mosaicism in important taxonomic morphological traits among the species of *Quercus* section *Cyclobalanopsis*<sup>5</sup>, molecular means such as chloroplast genomes can be used to explore intragroup interspecific relationships, identify species, and inform the implementation of species conservation strategies.

Currently, 50 CPGs for *Quercus* spp. could be queried in the National Center for Biotechnology Information (NCBI) database, 14 of which are from the sect. *Cyclobalanopsis*<sup>28</sup>. Here, we newly present CPGs sequences of 4 *Quercus* section *Cyclobalanopsis* species, including: *Quercus austrocochinchinensis*, *Quercus kerrii*, *Quercus helferiana*, and *Quercus rex*. Using these CPGs, we performed: (1) Structure and gene annotation. (2) Comparative genomics analysis. (3) Selection Pressure and Phylogenetic Analysis. This study aims to investigate: Characteristics and differences of CPGs among the four species; Hypervariable regions for the CPGs studied; Phylogenetic status of *Quercus* genus species. Our study will enrich the molecular data for the phylogenetic study and conservation of endangered species in the *Quercus* section *Cyclobalanopsis*.

## Materials And Methods

**Plant Samples, DNA Extraction and Sequencing.** Tender, unwounded leaf of 4 *Quercus* section *Cyclobalanopsis* species were harvested from 3 provinces in China: Yunnan, Hainan and Guizhou. Silica gel was used to dry the materials collected. Voucher specimens were saved in China West Normal University (CWNU) and sample information was listed in Table1. The improved CTAB protocol was used to extract and purify total genomic DNA from leaf tissues<sup>29</sup>. We used the high-quality genomic DNA to constructed a 400 bp Illumina Nova Seq library according to the manufacturer's protocol. Then the sequencing was performed on the Illumina Nova Seq PE150 platform, using pair-end strategies. Quality control on the raw data used FastQC<sup>30</sup>. Use Adapter Removal(v. 2)<sup>31</sup> to leach the joint contamination at the 3'end; quality filtration by sliding window method. Sequencing information was provided in Table 1.

**Chloroplast Genome Assembly, Annotation and Visualization.** CPGs were assembled by following steps: Firstly, clean reads were assembled by GetOrganelle<sup>32</sup>, with the iterative k-mers sizes setting to 21,45,65,85 and 105. Secondly, the assembled results were edited into circular sequences using Bandage<sup>33</sup>. Thirdly, the Geneious v.9.0.2<sup>34</sup> were using to adjust the initiations and find inverted repeat region. Assembled CPGs were annotated by Online website CPGAVAS2<sup>35</sup> and the complete plastome sequence of *Quercus ningangensis* (NC\_061582) as a reference. The intron/exon boundaries of annotation sequence were checked by Geneious. The CPG sequences and annotations were put in NCBI database. CPGs map were drawn on OGDRAW<sup>36</sup>.

**Table 1.** Basic information of 4 *Quercus* section *Cyclobalanopsis* species

| Taxa                            | Voucher  | Clean bases (G) | Average coverage (x) | NCBI accession number |
|---------------------------------|--|-----------------|----------------------|-----------------------|
| <i>Q. austrocochinchinensis</i> | CHINA. Hainan, 19°7'21.648"N, 109°9'36.828"E, 613m   | 2.11            | 92.7                 | OQ998918              |
| <i>Quercus. helferiana</i>      | CHINA. Guizhou, 25°3'47.559"N, 106°23'1.311"E, 932m  | 2.17            | 83.3                 | OQ998919              |
| <i>Quercus. kerrii</i>          | CHINA. Hainan, 19°7'24.708"N, 109°9'39.672"E, 606m   | 2.24            | 72.0                 | OQ998920              |
| <i>Quercus. rex</i>             | CHINA. Yunnan, 22°36'46.742"N, 101°6'13.042"E, 1595m | 2.29            | 59.9                 | OQ998921              |

**Genome Structure and Codon Usage Analyses.** In order to understand the framework of whole chloroplast genomes, Geneious was used to identify the size, genes and GC content in CPGs. Then confirmed and visualized the boundaries between IRs/ SCs by IRscope<sup>37</sup>. The totality of codons and RSCU (relative synonymous codon usage) values were calculated by CodonW v.1.4.2 with default parameters<sup>38</sup>.

**Sequence Divergence and Comparative Analysis.** The types of long sequence repeats (LSRs) were predicted by REPuter<sup>39</sup>, including type forward (F), type reverse (R), type complementary (C) and type palindromic (P), with parameters setting to: 30 bp for minimum repeat sequence, 3 for Hamming distance. In addition, MISA v1.0<sup>40</sup> with parameters setting of  $\geq 10$  for type mononucleotides,  $\geq 5$  for type dinucleotides,  $\geq 4$  for type trinucleotides, and  $\geq 3$  for type tetranucleotides, pentanucleotides and hexa-nucleotides were applied to predicted SSRs quantity and types. Multiple sequence alignment of CPGs were performed in mVISTA<sup>41</sup>, selecting Shufe-LAGAN mode when analyzing with *Quercus gilva* (MG678009) as a reference. After alignment the sequence by MAFFT v7.0<sup>42</sup> with default parameters, nucleotide diversity (Pi) values of CPGs evaluating were performed using DnaSP v6.0<sup>43</sup>.

**Selection Pressure and Phylogenetic Analysis.** KaKs\_Calculator v2.0<sup>44</sup> was adopted to calculated the rate of nonsynonymous mutation (Ka), synonymous mutation (Ks) in protein-coding genes. So that the results of Ka/Ks could be used to assess the role of selection for

each gene in CPGs of 11 *Quercus* species, 7 species of which were downloaded from NCBI (Supplementary Table S7).

For the purpose of acquainting with the phylogenetic relationships, phylogenetic tree of *Quercus* genus were implemented using Bayesian (BI) analysis methods, based on the CPG data. The CPG sequences required for the phylogenetic analysis are shown in the tableS7, including 27 Fagaceae species downloaded from NCBI. Apply all selected CPG sequences to MAFFT<sup>42</sup> to align. Later MrBayes v.3.2.<sup>45</sup> was utilized to carry out the BI tree analysis on the basis of following processes: infer the best-fit nucleotide substitution model (GTR+F+I+G4) by Modeltest v.3.7<sup>46</sup> and PAUP v.4.0<sup>47</sup>; Run 6,000,000 generations in Markov chain Monte Carlo (MCMC) analysis; Sample the trees each 1,000 generations and ignore the initial 0.25 as burnin fraction.

## Results

Characteristics of the CPGs. The length of 4 CPGs assembled scoped from 160715 bp in *Q. kerrii* to 160842 bp in *Q. rex*. All the structures manifest same circular quadripartite tetrad, comprising of 2 single-copy areas (LSC, SSC) and a couple of inverted repeats (IRs). The length of each region was shown in Table 2. The GC content of general sequences was 36.9% for all samples. Besides, the GC content in IRs lead to 42.8%, which was greater than that in LSC and SSC areas (34.8% and 31.1%). Additionally, all the 4 CPGs encoded 131 genes, namely 86 CDS, 37 tRNA and 8 rRNA, and it should be noted that eighteen (7 CDS, 7 tRNA and 4 rRNA) of these were iterant in the IRs. Among all of the genes, 15 have an intron and 3 genes (*rps12*, *clpP*, *ycf3*) with two. The specific distribution and function of the genes were shown in Figure 1, Supplementary Table S1.

**Table 2.** A summary of the statistics for the CPGs of 4 *Quercus* sect. *Cyclobalanopsis* species

| Species               | <i>Quercus kerrii</i> | <i>Quercus austrocochinchinensis</i> | <i>Quercus helferiana</i> | <i>Quercus rex</i> |
|-----------------------|-----------------------|--------------------------------------|---------------------------|--------------------|
| Genome size (bp)      | 160715                | 160768                               | 160801                    | 160842             |
| Length of LSC (bp)    | 90135                 | 90231                                | 90216                     | 90281              |
| Length of IRs (bp)    | 25841                 | 25835                                | 25840                     | 25839              |
| Length of SSC (bp)    | 18898                 | 18867                                | 18905                     | 18883              |
| Number of genes       | 131                   | 131                                  | 131                       | 131                |
| protein-coding genes  | 86                    | 86                                   | 86                        | 86                 |
| tRNA genes            | 37                    | 37                                   | 37                        | 37                 |
| rRNA genes            | 8                     | 8                                    | 8                         | 8                  |
| GC content (%)        | 36.9                  | 36.9                                 | 36.9                      | 36.9               |
| GC content of LSC (%) | 34.8                  | 34.8                                 | 34.8                      | 34.8               |
| GC content of IRs (%) | 42.8                  | 42.8                                 | 42.8                      | 42.8               |
| GC content of SSC (%) | 31.1                  | 31.1                                 | 31.1                      | 31.1               |

Figure 2 gave the results of CPGs boundary comparison in 6 *Quercus* section *Cyclobalanopsis* species, which could show the borderlines of the IRs and SCs regions in CPGs. The junction of LSC and IRb (*JLB*) laid in IGS (intergenic spacer) of *rps19* and *rpm2* gene. Most samples had 11 bp shift away from the borderline for *rps19* gene in JLB, except *Q. helferiana* and *Q. neglecta*, which had 13 bp and 4 bp shift respectively. Moreover, the demarcation of LSC and IRa was situated in the IGS of *rpm2* and *trnH* gene, with the *trnH* gene shifting 15 or 16 bp from JLA. IRa/SSC boundary (*JSA*) was reposed within gene *ycf1*. What should be noted was that the 5' end of gene *ycf1* stood in the IRa area but the 3' end stood in SSC area, therefore created a 5' end pseudogene (*ycf1Ψ*) in the IRb in all CPGs compared, resulting in all IRb/SSC (*JSB*) boundaries lying within the pseudogene *ycf1Ψ*.

The codon usage analysis summarized in Table 3. According to the results, sequence sizes range of extracted protein-coding genes were 64359 - 64377 bp in 4 *Quercus* section *Cyclobalanopsis* species; 21453 - 21459 codons were encoded. The ENC (effective number of codons) value was between 49.93 and 49.97. The FOP (Frequency of optimal codons) value was 0.353 in *Q. kerrii* and 0.354 in other three samples. The GC content was between 37.93% and 37.95%. The codon preference indexes of the four species varied slightly, indicating that they had similar codon usage. The GC3 of four species ranged between 29.85% and 29.88%, indicating that they prefer codons ending with A/U.

The CDSs of 17 CPGs (4 newly sequenced and 13 species of Fagaceae released in NCBI) were extracted using Geneious. Subsequently, based on the extracted sequences, the ratio of RSCU in all samples were calculated and clustered. The results showed in Figure 3, Supplementary Table S2. We found that:

(1) Leucine (Leu) encoded with the maximum number of codons, arranging from 2,044 to 2,268, with the number of isoleucine(Ile) (1699 ~1892) following. The minimum number of codons (213 to241) presented in Cysteine(Cys). (2) The (RSCU) values varied marginally among the CDSs of 17 species. 31 codons were frequently manipulated since RSCU ratio was greater than 1 and the remaining codons were less frequently used as their RSCU ratios were less than 1. (3) The frequently used codons include: UUA, AGA UAA(\*) GCU UCU GAU ACU, and the codon usage frequency of UAC CUC CGC CUG AGC GAC is on the low side. Thereinto the UUA codon showed a bias in 17 CPGs due to its highest usage. No usage frequency bias (RSCU=1) showed in the starting codons of AUG and UGG, which encoded methionine (Met) and tryptophan (Try).

**Table 3.** Codon preference index of four species of *Quercus* section *Cyclobalanopsis*

| Index                       | <i>Quercus. kerrii</i> | <i>Quercus. austrocochinchinensis</i> | <i>Quercus. helferiana</i> | <i>Quercus. rex</i> |
|-----------------------------|------------------------|---------------------------------------|----------------------------|---------------------|
| Length (bp)                 | 64359                  | 64359                                 | 64359                      | 64377               |
| Codon number                | 21453                  | 21453                                 | 21453                      | 21459               |
| Effective number of codons  | 49.94                  | 49.97                                 | 49.96                      | 49.93               |
| Codon adaptation index      | 0.166                  | 0.166                                 | 0.166                      | 0.166               |
| Codon bias index            | -0.100                 | -0.100                                | -0.100                     | -0.100              |
| Frequency of optimal codons | 0.353                  | 0.354                                 | 0.354                      | 0.354               |
| GC content (%)              | 37.93                  | 37.95                                 | 37.94                      | 37.93               |
| GC1 content (%)             | 46.05                  | 46.05                                 | 46.06                      | 46.03               |
| GC2 content (%)             | 37.82                  | 37.84                                 | 37.82                      | 37.83               |
| GC3 content (%)             | 29.92                  | 29.95                                 | 29.93                      | 29.92               |

Repeated sequences analysis. A total of 163 LSRs were identified among the 4 CPGs examined. As a whole, the amount of LSRs identified in every CPG was scoping from 37 in *Q. rex* to 44 LSRs in *Q. helferiana*. Thereinto, 14~18 were type F, 20~22 were palindromic repeats, and the number of type R was 2 in *Q. rex* when other three species were 3 (Figure 4A, Supplementary Table S3). Just one complement repeat was filtrated from 4 species. Among these repeats, the longest repeat was 56 bp in every species, and the most common length of repeats are 30 bp. 44.5% LSRs located in IGS, and 23.5% were found in the intron region. About half repeat sequences (46.8%) were distinguished in the IR areas (Supplementary Table S3).

The total quantity of SSRs identified in the CPGs of four *Quercus* section *Cyclobalanopsis* species was 453, ranging from 110 in *Q. austrocochinchinensis* to 115 in *Q. rex*, among which 74-81 were type mono-, 15 were type di-, 7 were type tri-, 9-11 were type tetra-, and 3 were penta- (Figure 4B). The most universal unit of SSRs was A/T (mono-), whose amount ranged from 69 to 76, far higher than in the other types. 68% of SSRs were type mononucleotide made up of unit A/T and C/G. What's more, most of the SSRs (70.8%) were located in the IGS (Supplementary Table S4). All the type din- comprised multiple copies of unit AT/TA and AG/CT (Figure 4C). The type of hexanucleotide was not detected in all species. Taken as a whole, no significant distinction in the number of SSR units among the 4 species, except the slight differences in unit of mono- and penta-.

Sequence Divergence, Hotspots and Selection Pressure Estimation. CPGs Comparative analysis could be seen in Figure 5, revealing that high sequence similarity among the four sect. *Cyclobalanopsis* species. Sequences in noncoding areas were more variant than in coding areas generally. Besides, the level of sequence divergence in SCs areas were evidently higher than that in IR areas. We found 8 intergenic regions were in a high degree of variation, of which 7 were located at LSC areas as follows: *psbA/trnH*, *rps16/trnK*, *trnQ/rps16*, *trnE/trnT*, *rbcl/accD*, *psbE/petL*, *ndhF/rpl32*. One located at SSC areas, namely *ndhI/ndhG*. Other than aforementioned areas, the intron area of *rpoC1* showed high level of sequence divergence too.

Window length setting to 600-bp, we calculated the nucleotide diversity values to elucidate levels of diversity for all CPGs assembled in this study. The Pi values were recorded in Supplementary Table S5, ranging from 0 to 0.01083, with 0.00041 on average. When the amount of

polymorphic loci outweighed the sum of mean and twofold standard deviation, the region is defined as a highly variable region<sup>48</sup>. Ultimately, six hotspots ( $P_i > 0.0022$ ) were discovered, coding and noncoding regions each accounting for half. The greatest  $P_i$  value (0.01083) appeared in the region between gene *trnK-UUU* and *rps16*. The distribution of highly variable regions was shown in Figure 6A. In general, these regions were not located at the IR areas but all at the SC areas, which reflected an identical pattern of CPG structural variation.

To estimate the role of selection of the *Quercus* section *Cyclobalanopsis* species,  $K_a$  and  $K_s$  values of 79 unique CDS were calculated in 11 CPGs using *Quercus chenii* as a reference. The  $K_a/K_s$  values were simply calculated and recorded in the Supplementary Table S6, ranging from 0 to 1.471. Among which 40 protein-coding genes showed significance (Figure 6B) in 11 species. Based on the calculation results, we speculated that the purifying selection may affect on most protein coding genes, as their  $K_a/K_s$  values were less than 1. At the same time, when  $K_a/K_s$  values  $>1$  demonstrated that the positive selection was working on the genes. Therefore we identified two genes were under the positive selection, namely *petA* gene in *Q. aliena*, and *ycf2* gene in *Q. austrocochinchinensis*, *Q. rex*, *Q. kerrii*, *Q. sichouensis*, and *Q. neglecta*.

Phylogenetic relationships. Resorting to approaches BI, the phylogenetic relationships were reconstructed among the members of the four CPGs sequenced in this study and closely related species in *Quercus* genus, according to the whole chloroplast genome data. The *Trigonobalanus doichangensis* (NC\_023959) was used as the outgroup. A total of 31 taxa were used and reconstructed phylogeny tree was shown in Figure 7, and most branches obtained high support bootstrap values. Genus *Trigonobalanus* and *Fagus* located at the base of the phylogeny tree. Two distinct clades were recognized among all *Quercus* species analysed: the first clade consisted of two sections (four species in *Quercus* and three species in *Lobatae*). Another clade was divided into two nodes, including three sections, namely *Cyclobalanopsis*, *Cerris* and *Ilex*. In section *Cyclobalanopsis*, the species were divided to STB (Single-celled trichome base) and CTB (compound trichome base). All the CTB species were clustered into a single branch including the four species we studied.

## Discussion

CPG Architectures in 4 *Quercus* Section *Cyclobalanopsis* species. Four species CPGs were successfully assembled of *Quercus* Section *Cyclobalanopsis* in the present paper. The size of four CPGs (ca. 160 kb) conformed to the photosynthetic land plant plastid chromosomes, whose size varied from 120 kb to 160 kb<sup>49</sup>. The same quadripartite circular structure were found in the four assembled CPGs and other *Quercus* species<sup>50-52</sup>. Overall GC content had no distinction within the four species. After CPGs comparison, it was found that the totality, order, and function of genes were highly conservative in genus *Quercus*, which were also in accordance with most Fagaceae species<sup>24, 53, 54</sup>, evidencing a highly conservative CPG construction in *Quercus* Section *Cyclobalanopsis*.

Due to the duplicative nature of the IR reduced the substitution rate within this region, therefore it was of great significance to analyze the contraction and expansion of IRs in evolutionary researches<sup>55</sup>. In addition, the IR regions were vital in stabilizing the structure of the CPGs, which were also the main factor affecting the total length<sup>56, 57</sup>. The results showed that boundaries of four areas in the CPGs were conserved in 6 *Quercus* section *Cyclobalanopsis* species. The IRs/SCs boundary of all species compared in this study were located within similar positions except for slight difference in *JLB*, whose displacements from *rps19* presented subtle variations in different species. Most of the compared species found no significant expansion or contraction in the IR regions except the *Quercus neglecta*, which had a only 4 bp displacement between the *JLB* and *rps19*. The conservatism of *Quercus* section *Cyclobalanopsis* was demonstrated by the relatively constant length of CPGs and the minor variations in their region borders, as the same conditions with other *Quercus* species<sup>24, 51</sup>.

Codon usage bias was a natural phenomenon caused by mutation, natural selection, genome composition, etc<sup>58-60</sup>. In the CPGs of four cp genomes, total 64 codons were detected, encoding 20 amino acids. We could tell from the values of RSCU and content of GC3 that the bias in codon usage towards A / U at the third position, a phenomenon that is widespread in angiosperms<sup>61-64</sup>.

Large Repeats and Simple Sequence Repeats. Dispersed in CPGs, long repeat sequences played an significant role the genomic inheritance, variation and the evolution of species<sup>49, 56, 65</sup>. Our study identified a total of 163 LSRs with palindromic being the most common type. The variations observed in CPGs could partially attributed to the differences in the number and types of LSRs<sup>66</sup>. Therefore, due to their genetic variations, LSRs can potentially provide valuable information for researches of phylogenetic relationship and population genetics. After analysis, it was found that the repeat sequences of this study were in accord with the general pattern above: about half (43.2 - 46.3%) of LSRs were identified in IGS, following by the coding regions and introns. Current studies had suggested that most repeats in CPG were situated in the IGS, comparing to the coding regions<sup>14, 67</sup>. SSR, had been extensively studied as a kind of effective molecular marker in various fields such as discrimination, breeding, conservation and phylogenetic research at the species and population levels<sup>68-70</sup>. A strong A/T bias, SC regions concentration (90.9 - 91.3%), and IGS concentration (69.6 - 72.7%) were detected in SSRs of four *Quercus* section *Cyclobalanopsis* species, similar to other *Quercus* genus species<sup>28, 71</sup>. The numbers and types of SSRs varied slightly in

*Quercus* genus but extensively in other families<sup>72-74</sup>. The numbers of SSRs were almost identical between *Quercus* section *Cyclobalanopsis* and section *Cerris*<sup>73</sup>, so we speculated that such case might imply that the two sections were phylogenetically more closely related.

Conservatism, Highly Variable Regions and Selection Pressure Estimation. We compared the whole sequences of CPG in four species with *Quercus gilva* as the reference. The results indicated that there were differences in the degree of variation between regions of CPGs, with the single-copy (SC) regions having higher variation than IR regions, simultaneously the IGS regions having higher variation than coding regions. Same phenomena were found in other *Quercus* species<sup>50; 51; 75; 76</sup>. The copy-dependent repair mechanism of CPGs could guarantee the stability of IRs construction and thereby advance the steadiness and conservation of genomes, which possibly explain the different degree of variation between IRs and SCs. In addition, due to natural selection, the coding areas tend to exhibit higher conservation than the noncoding areas<sup>77-79</sup>. The gene regions of high variability we found (namely *rpoC1*, *clpP* and *ycf1*) in both sequence divergence analysis and nucleotide variability (pi) assessment could be used to develop DNA barcodes, conduct species identification and systematic classification<sup>80</sup>. Out of the highly variable regions identified, the *ycf1* gene<sup>81</sup> and two IGS regions: *trnH-psbA*, *trnK-rps16* had already been selected as practicable barcode for plants<sup>82-84</sup>.

In our study, most of the Ka/Ks values were < 1 or not available, suggesting that the emergence frequency of synonymous nucleotide substitution was more than that of non-synonymous nucleotide substitution due to the purify selection process<sup>85; 86</sup>. We conjectured that positive selection was operating only in two genes: *petA* in *Q. aliena* and the *ycf1* in multiple *Quercus* taxa. The *ycf1* was indicated to contain multiple SSRs in many taxa and it was claimed that these SSRs were undoubted in detecting population-level polymorphisms and could also be used to compare phylogenetic relationships at the genus level or higher taxonomic levels<sup>18; 72</sup>. Whether these divergence hotspots found in the above analysis could be utilized for DNA barcodes or estimating taxonomic evolution in genus *Quercus* needs more further researches.

Inference of Phylogenetic Relationship. Due to the complex evolutionary issues such as convergent evolution, extensive infiltration, hybridization, and serious hybridization introgression in the *Quercus* genus, great challenges remain in the phylogenetic relationship research of Oak trees<sup>1; 87; 88</sup>. CPGs have been demonstrated considerable utility in addressing the phylogeny relationships of angiosperms<sup>89</sup>. The phylogenetic trees we reconstructed based on CPGs indicated two major clades corresponding to geographic distribution: sections of *Quercus* and *Lobatae* constituted a "New World Clade" (subgenus *Quercus*), while the sections of *Cyclobalanopsis*, *Cerris* and *Ilex* forming an "Old World Clade" (subgenus *Cerris*)<sup>8; 88; 90</sup>. The section *Ilex* was paraphyletic, nested into the lineage formed by section *Cerris*, which was similar to the results based on plastid genome but differed from the phylogenetic relationships inferred from RAD-seq data<sup>28; 87; 91</sup>.

The morphological studies found that the four species we studied possessed compound trichome base (CTB) so that clustered into a single branch with other CTB species, distincting to the group STB (single-celled trichome base)<sup>5; 92</sup>, similar to the results of our phylogenetic study based on the CPGs. In the CTB group, *Q. kerrii* and *Q. chungii* clustered into a clade that diverging before the other four species, which had simple uniseriate thin-walled trichomes, distinct from other CTB species<sup>5</sup>. *Q. austrocochinchinensis* then clustered with *Q. gilva* into sister groups, which differed from the clustering results of RAD-seq data<sup>4; 93</sup>. *Q. helferiana* and *Q. rex* gathered together, they all possessed Fasciculation trichomes<sup>5</sup>. Deng M.<sup>3</sup> divided the CTB species into two groups based on their comprehensive morphological characteristics, namely group *Gilva* (containing *Q. chungii* and *Q. Gilva*) and group *Helferiana*, including the four species we studied. From our results, we can see that neither of these two groups had formed a monophyly, and there were mosaics between these species. *Q. helferiana* and *Q. kerrii* were far apart in the phylogeny tree, so we believe that they can exist as independent species. Nevertheless the interspecific relationship within the four species remained some controversies: for instance, the *Q. kerrii* and *Q. austrocochinchinensis* gathered for a monophyletic sister branch in multiple studies, different from our BI tree. The *Q. rex* was thought to be the base of 4 species, but in our analysis the *Q. kerrii* differentiated firstly<sup>4; 5</sup>. Due to the presence of one node with a low support rate (0.338) in the phylogeny tree, the interspecies relationship between the two branches differentiated by this node was still unclear. The continued advancements of sequencing techniques will allow for the inclusion of more taxa and samples in future studies, facilitating further exploration of the interspecific relationships and phylogenomics of the *Quercus* section *Cyclobalanopsis*.

## Conclusions

In this study, we successfully completed the CPG basic analysis of four species in *Quercus* section *Cyclobalanopsis*. Despite the overall conservation of CPG structure and gene content were obviously found, distinct sequence divergences were uncovered in alternating regions of these genomes among the studied species. The regions of divergence hotspots and multitudinous SSRs detected in the CPGs had

potential for development as molecular markers, which could aid in further population genetic studies and facilitate the establishment of appropriate protection policies for vulnerable species. The phylogenetic analysis based on CPG data suggested: genus *Trigonobalanus* and *Fagus* differentiated earlier; All *Quercus* species were divided into two categories (STB and CTB). The phylogenetic relationships between the various sections in *Quercus* genus were strongly supported. In a word, the findings obtained will facilitate further investigations into the taxonomy, phylogenetic evolution and preservation of *Quercus* genus.

## Abbreviations

RAD-Seq Restriction-site associated DNA sequencing

ITS Internally Transcribed Spacer

CPG Chloroplast Genome

CTAB Cetyltrimethylammonium Bromide

LSC Large single-copy

SSC Small single-copy

IR Inverted repeat

RSCU Relative Synonymous Codon Usage

LSRs Long Sequence Repeats

SSRs Simple Sequence Repeat

BI Bayesian inference

CDS Coding sequence

IGS Intergenic Spacer

## Declarations

### Data Availability Statement

The datasets generated and/or analysed during the current study are available in the [National Center for Biotechnology Information] repository, [Accession Number: OQ998918, OQ998919, OQ998920, OQ998921]

### Author contributions

X.L.C. conducted plant collection, data analysis, and paper writing. X.M.Z. conducted plant identification, experimental guidance, and paper revisions. All authors read and approved the final manuscript.

### Funding

This work was supported by grants from National Specimen Platform Teaching Standard Subplatform(<http://mnh.scu.edu.cn/>) (2005DKA21403-JK), Research and Innovation Team of China West Normal University (KCXTD2022-4).

### Competing interests

The authors declare no competing interests.

### Guideline Statement

The plant materials directly used in this study do not contain any rare or endangered plants. The collection of plant materials has been approved by the National Park of Hainan Tropical Rainforest, Yunnan Taiyanghe Provincial Nature Reserve and Guizhou Wangmo *Cycas* Nature Reserve.



Correspondence and requests for materials should be addressed to X.M.Z.

## References

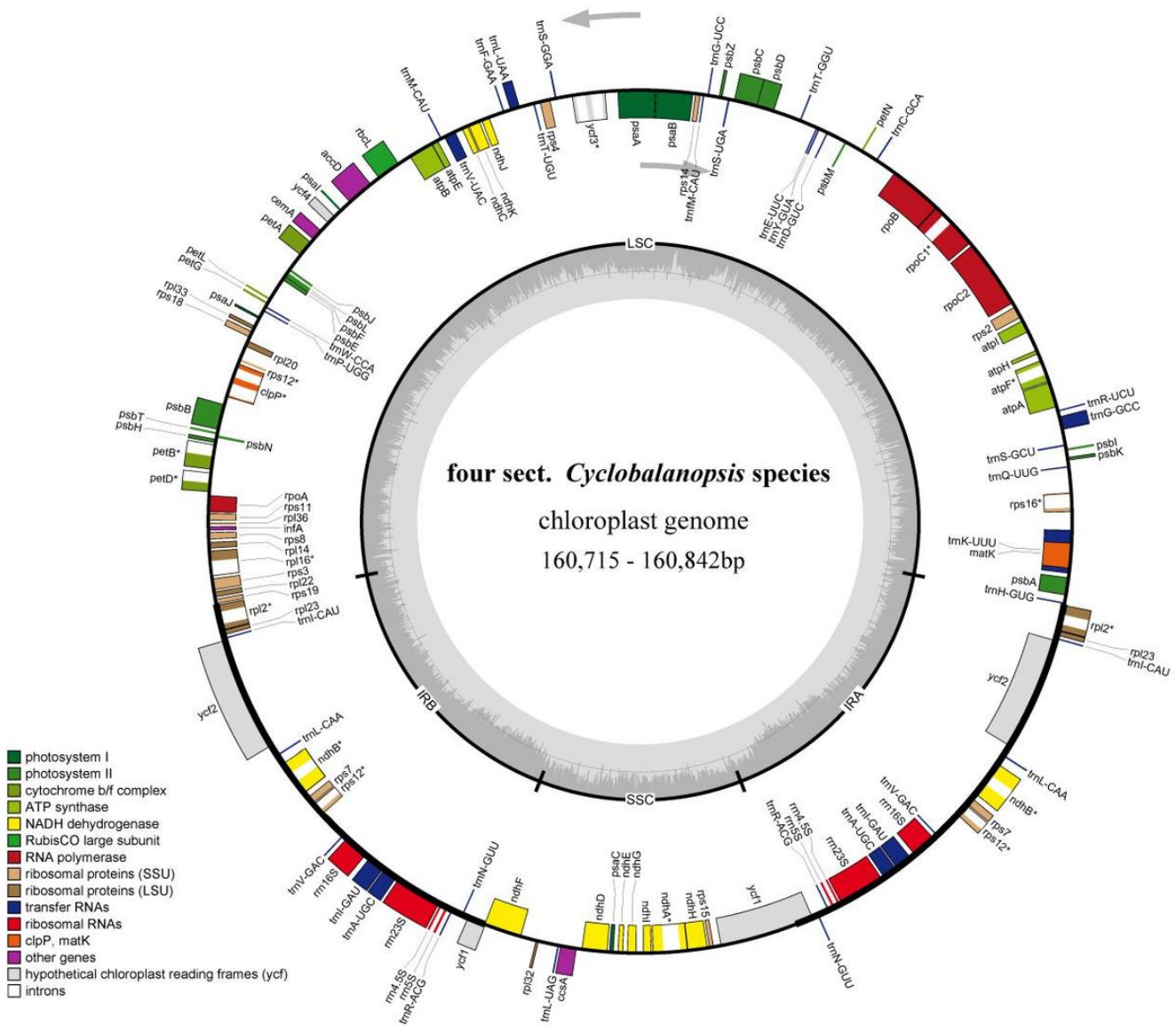
1. Manos, P. S., Doyle, J. J. & Nixon, K. C. Phylogeny, Biogeography, and Processes of Molecular Differentiation in *Quercus* Subgenus *Quercus* (Fagaceae). *Mol. Phylogenet. Evol.* **12**, 333-349 (1999).
2. Denk, T., Grimm, G. W., Manos, P. S., Deng, M. & Hipp, A. L. An Updated Infrageneric Classification of the Oaks: Review of Previous Taxonomic Schemes and Synthesis of Evolutionary Patterns. In: Gil-Pelegri n E, Peguero-Pina JJ, Sancho-Knapik D, eds. *Oaks Physiological Ecology. Exploring the Functional Diversity of Genus Quercus L.* Cham: Springer International Publishing, 2017:13-38.
3. Deng, M. *Anatomy, Taxonomy, Distribution and Phylogeny of Quercus Subg. Cyclobalanopsis (Oersted) Schneid. (Fagaceae)*. Chinese Academy of Sciences, 2007.
4. Deng, M., Jiang, X., Hipp, A. L., Manos, P. S. & Hahn, M. Phylogeny and Biogeography of East Asian Evergreen Oaks (*Quercus* Section *Cyclobalanopsis*; Fagaceae): Insights Into the Cenozoic History of Evergreen Broad-Leaved Forests in Subtropical Asia. *Mol. Phylogenet. Evol.* **119**, 170-181 (2018).
5. Deng, M. et al. Leaf Epidermal Features of *Quercus* Subgenus *Cyclobalanopsis* (Fagaceae) and their Systematic Significance. *Bot. J. Linnean Soc.* **176**, 224-259 (2014).
6. TSCHAN, G. F. & DENK, T. Trichome Types, Foliar Indumentum and Epicuticular Wax in the Mediterranean Gall Oaks, *Quercus* Subsection *Galliferae* (Fagaceae): Implications for Taxonomy, Ecology and Evolution. *Bot. J. Linnean Soc.* **169**, 611-644 (2012).
7. Sauquet, H. & Cantrill, D. J. Pollen Diversity and Evolution in Proteoideae (Proteales: Proteaceae). *Syst. Bot.* **32**, 271-316 (2007).
8. Denk, T., Grimm, G. W., Manos, P. S., Min, D. & Hipp, A. L. An Updated Infrageneric Classification of the Oaks: Review of Previous Taxonomic Schemes and Synthesis of Evolutionary Patterns. *Tree Physiology: Springer, Cham.* **7**, 13-38 (2017).
9. Daniell, H., Lin, C., Yu, M. & Chang, W. Chloroplast Genomes: Diversity, Evolution, and Applications in Genetic Engineering. *Genome Biol.* **17**, 134 (2016).
10. K, S. et al. The Complete Nucleotide Sequence of the Tobacco Chloroplast Genome: Its Gene Organization and Expression. *The Embo Journal.* **5**, 2043-2049 (1986).
11. Bobik, K. & Burch-Smith, T. M. Chloroplast Signaling within, Between and Beyond Cells. *Front. Plant Sci.* **6**, (2015).
12. Yang, J. B., Yang, S. X., Li, H. T., Yang, J. & Li, D. Z. Comparative Chloroplast Genomes of *Camellia* Species. *Plos One.* **8**, e73053 (2013).
13. Li, X. et al. Plant Dna Barcoding: From Gene to Genome. *Biol. Rev.* **90**, 157-166 (2015).
14. Hong, Z. et al. Comparative Analyses of Five Complete Chloroplast Genomes From the Genus *Pterocarpus* (Fabaceae). *International Journal of Molecular Sciences.* **21**, 3758 (2020).
15. Korpelainen, H. The Evolutionary Processes of Mitochondrial and Chloroplast Genomes Differ From those of Nuclear Genomes. *Sci. Nat.* **91**, 505-518 (2004).
16. Parks, M., Cronn, R. & Liston, A. Increasing Phylogenetic Resolution at Low Taxonomic Levels Using Massively Parallel Sequencing of Chloroplast Genomes. *Bmc Biology.* **7**, 84-100 (2009).
17. Wei, W. et al. Pcr-Rflp Analysis of Cpdna and Mtdna in the Genus *Houttuynia* in some Areas of China. *Hereditas.* **142**, 24-32 (2005).
18. Huang, H., Shi, C., Liu, Y., Mao, S. Y. & Gao, L. Z. Thirteen *Camellia* Chloroplast Genome Sequences Determined by High-Throughput Sequencing: Genome Structure and Phylogenetic Relationships. *Bmc Evol. Biol.* **14**, 151 (2014).
19. Xue, S. et al. Comparative Analysis of the Complete Chloroplast Genome Among *Prunus Mume*, *P. Armeniaca*, and *P. Salicina*. *Hortic. Res.-England.* **6**, 89 (2019).
20. Shendure, J. & Ji, H. Next-Generation Dna Sequencing. *Nat. Biotechnol.* **26**, 1135-1145 (2008).
21. Itoh, M. O. K. B. Possibility of Grouping of *Cyclobalanopsis* Species (Fagaceae) Grown in Japan Based On an Analysis of Several Regions of Chloroplast Dna. *The Japan Wood Research Society.* 498-501 (1999).
22. Catherine J Nock et al. Chloroplast Genome Sequences From Total Dna for Plant Identification. *Plant Biotechnol. J.* **9**, 328-333 (2011).
23. Gaixia Zhang et al. Identification of the Original Plants of Cultivated *Bupleuri Radix* Based On Dna Barcoding and Chloroplast Genome Analysis. *Peerj.* **10**, e13208 (2022).
24. Xu, J. et al. Phylogeography of *Quercus Glauca* (Fagaceae), a Dominant Tree of East Asian Subtropical Evergreen Forests, Based On Three Chloroplast Dna Interspace Sequences. *Tree Genet. Genomes.* **11**, 805 (2014).
25. Kamiya, K., Harada, K., OGINO, K., CLYDE, M. & LATIFF, A. Phylogeny and Genetic Variation of Fagaceae in Tropical Montane Forests. *Tropics.* **13**, 119-125 (2003).

26. Asaf, S. et al. Comparative Analysis of Complete Plastid Genomes From Wild Soybean (*Glycine Soja*) and Nine Other *Glycine* Species. *Plos One*. **12**, e182281 (2017).
27. Ruihong, Y., Runfang, G., Yuguang, L., Ziqian, K. & Baosheng, S. Identification and Phylogenetic Analysis of the Genus *Syringa* Based On Chloroplast Genomic Dna Barcoding. *Plos One*. **17**, e271633 (2022).
28. Li, Y. et al. Complete Chloroplast Genome of an Endangered Species *Quercus Litseoides*, and its Comparative, Evolutionary, and Phylogenetic Study with Other *Quercus* Section *Cyclobalanopsis* Species. *Genes*. **13**, 1184 (2022).
29. Allen, G. C., Flores-Vergara, M. A., Krasynanski, S., Kumar, S. & Thompson, W. F. A Modified Protocol for Rapid Dna Isolation From Plant Tissues Using Cetyltrimethylammonium Bromide. *Nat. Protoc.* **1**, 2320-2325 (2006).
30. Andrews, S. Fastqc a Quality Control Tool for High Throughput Sequence Data., 2014.
31. Mikkel et al. Adapterremoval V2: Rapid Adapter Trimming, Identification, and Read Merging. *Bmc Res. Notes*. **9**, 88 (2016).
32. Jin, J. et al. Getorganelle: A Fast and Versatile Toolkit for Accurate De Novo Assembly of Organelle Genomes. *Genome Biol.* **21**, 241 (2020).
33. Wick, R. R., Schultz, M. B., Justin, Z. & Holt, K. E. Bandage: Interactive Visualization of De Novo Genome Assemblies. *Bioinformatics*. 3350-3352 (2015).
34. Matthew Kearse, R. M. A. W. et al. Geneious Basic: An Integrated and Extendable Desktop Software Platform for the Organization and Analysis of Sequence Data. *Bioinformatics*. **12**, 1647-1649 (2012).
35. Shi, L. et al. Cpgavas2, an Integrated Plastome Sequence Annotator and Analyzer. *Nucleic Acids Research*. **47**, W65-W73 (2019).
36. Lohse, M., Drechsel, O., Kahlau, S. & Bock, R. Organellargenomedraw—A Suite of Tools for Generating Physical Maps of Plastid and Mitochondrial Genomes and Visualizing Expression Data Sets. *Nucleic Acids Res.* **41**, W575-W581 (2013).
37. Amiryousefi, A., Hyvönen, J. & Poczai, P. Irscope: An Online Program to Visualize the Junction Sites of Chloroplast Genomes. *Bioinformatics*. **34**, 3030-3031 (2018).
38. Sharp, Paul, M., Li & Wen-Hsiung. The Codon Adaptation Index—a Measure of Directional Synonymous Codon Usage Bias, and its Potential Applications. *Nucl Acids Res.* **15**, 1281-1295 (1987).
39. Kurtz, S. et al. Reputer: The Manifold Applications of Repeat Analysis On a Genomic Scale. *Nucleic Acids Res.* **29**, 4633-4642 (2001).
40. Sebastian et al. Misa-Web: A Web Server for Microsatellite Prediction. *Bioinformatics*. **33**, 2583-2585 (2017).
41. Frazer, K. A., Pachter, L., Poliakov, A., Rubin, E. M. & Dubchak, I. Vista: Computational Tools for Comparative Genomics. *Nucleic Acids Res.* **32**, W273-W279 (2004).
42. Katoh, K. & Standley, D. M. Mafft Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. *Mol. Biol. Evol.* **30**, 772-780 (2013).
43. Rozas, J. et al. Dnasp 6: Dna Sequence Polymorphism Analysis of Large Data Sets. *Mol. Biol. Evol.* **34**, 3299-3302 (2017).
44. Wang, D., Zhang, Y., Zhang, Z., Zhu, J. & Yu, J. Kaks\_Calculator 2.0: A Toolkit Incorporating Gamma-Series Methods and Sliding Window Strategies. *Genom. Proteomics Bioinformatics*. **8**, 77-80 (2010).
45. Fredrik Ronquist et al. Mrbayes 3.2: Efficient Bayesian Phylogenetic Inference and Model Choice Across a Large Model Space. *Syst. Biol.* **61**, 539-542 (2012).
46. Posada, D. & Crandall, K. A. Modeltest: Testing the Model of Dna Substitution. *Bioinformatics*. **14**, 817-818 (1998).
47. Matthews, L. J. & Rosenberger, A. L. Taxon Combinations, Parsimony Analysis (Paup\*), and the Taxonomy of the Yellow-Tailed Woolly Monkey, *Lagothrix Flavicauda*. *Wiley Subscription Services, Inc., A Wiley Company*. **137**, 245-255 (2008).
48. Wei Wang et al. Comparative and Phylogenetic Analyses of the Complete Chloroplast Genomes of Six Almond Species (*Prunus* Spp. L.). *Sci Rep.* **10**, 10137 (2020).
49. Wicke, S., Schneeweiss, G. M., DePamphilis, C. W., Müller, K. F. & Quandt, D. The Evolution of the Plastid Chromosome in Land Plants: Gene Content, Gene Order, Gene Function. *Plant Mol.Biol.* **76**, 273-297 (2011).
50. Li, X., Li, Y., Zang, M., Li, M. & Fang, Y. Complete Chloroplast Genome Sequence and Phylogenetic Analysis of *Quercus Acutissima*. *International Journal of Molecular Sciences*. **19**, 2443 (2018).
51. Wang, T., Wang, Z., Song, Y. & Kozłowski, G. The Complete Chloroplast Genome Sequence of *Quercus Ningangensis* and its Phylogenetic Implication. *Plant and Fungal Systematics*. **66**, 155-165 (2021).
52. Chen, S. et al. The Complete Chloroplast Genome Sequence of *Quercus Sessilifolia* Blume (Fagaceae). *Mitochondrial Dna. Part B. Resources*. **7**, 182-184 (2022).

53. Liang, D., Wang, H., Zhang, J., Zhao, Y. & Wu, F. Complete Chloroplast Genome Sequence of *Fagus Longipetiolata* Seemen (Fagaceae): Genome Structure, Adaptive Evolution, and Phylogenetic Relationships. *Life*. **12**, 92 (2022).
54. Yang, X., Yin, Y., Feng, L., Tang, H. & Wang, F. The First Complete Chloroplast Genome of *Quercus Coccinea* (Scarlet Oak) and its Phylogenetic Position within Fagaceae. *Mitochondrial Dna. Part B, Resources*. **4**, 3634-3635 (2019).
55. Cai, Z. et al. Complete Plastid Genome Sequences of Drimys, Liriodendron, and Piper: Implications for the Phylogenetic Relationships of Magnoliids. *Bmc Evol. Biol.* **6**, 77 (2006).
56. Maréchal, A. & Brisson, N. Recombination and the Maintenance of Plant Organelle Genome Stability. *New Phytol.* **186**, 299-317 (2010).
57. Chumley, T. W. et al. The Complete Chloroplast Genome Sequence of Pelargonium × Hortorum: Organization and Evolution of the Largest and Most Highly Rearranged Chloroplast Genome of Land Plants. *Mol. Biol. Evol.* **23**, 2175-2190 (2006).
58. Xu, C. et al. Factors Affecting Synonymous Codon Usage Bias in Chloroplast Genome of *Oncidium Gower Ramsey*. *Evol. Bioinform.* **7**, 271-278 (2011).
59. Ikemura, T. Codon Usage and Trna Content in Unicellular and Multicellular Organisms. *Mol. Biol. Evol.* **2**, 13-34 (1985).
60. Bernardi, G. & Bernardi, G. Compositional Constraints and Genome Evolution. *J. Mol. Evol.* **24**, 1-11 (1986).
61. Chi, X., Zhang, F., Dong, Q. & Chen, S. Insights Into Comparative Genomics, Codon Usage Bias, and Phylogenetic Relationship of Species From Biebersteiniaceae and Nitrariaceae Based On Complete Chloroplast Genomes. *Plants*. **9**, 1605 (2020).
62. Ren, T. et al. Plastomes of Eight *Ligusticum* Species: Characterization, Genome Evolution, and Phylogenetic Relationships. *Bmc Plant Biol.* **20**, 519 (2020).
63. Delannoy, E., Fujii, S., Colas Des Francs-Small, C., Brundrett, M. & Small, I. Rampant Gene Loss in the Underground Orchid *Rhizanthella Gardneri* Highlights Evolutionary Constraints On Plastid Genomes. *Mol. Biol. Evol.* **28**, 2077-2086 (2011).
64. Tangphatsornruang, S. et al. The Chloroplast Genome Sequence of Mungbean (*Vigna Radiata*) Determined by High-Throughput Pyrosequencing: Structural Organization and Phylogenetic Relationships. *Dna Research : An International Journal for Rapid Publication of Reports On Genes and Genomes*. **17**, 11-22 (2009).
65. Chen, Y., Hu, N. & Wu, H. Analyzing and Characterizing the Chloroplast Genome of *Salix Wilsonii*. *Biomed Res. Int.* **2019**, 1-14 (2019).
66. Yang, F. et al. Complete Chloroplast Genome Sequence of Poisonous and Medicinal Plant *Datura Stramonium*: Organizations and Implications for Genetic Engineering. *Plos One*. **9**, e110656 (2014).
67. Deng, Y. Complete Chloroplast Genome of *Michelia Shiluensis* and a Comparative Analysis with Four Magnoliaceae Species. *Forests*. **11**, 267 (2020).
68. Yan, X. et al. Chloroplast Genomes and Comparative Analyses Among Thirteen Taxa within Myrsinaceae S.Str. Clade (Myrsinoideae, Primulaceae). *International Journal of Molecular Sciences*. **20**, 4534 (2019).
69. Yamamoto, T. Dna Markers and Molecular Breeding in Pear and Other Rosaceae Fruit Trees. *The Horticulture Journal*. **90**, 1-13 (2021).
70. Mohammad-Panah, N., Shabaniyan, N., Khadivi, A., Rahmani, M. & Emami, A. Genetic Structure of Gall Oak (*Quercus Infectoria*) Characterized by Nuclear and Chloroplast Ssr Markers. *Tree Genet. Genomes*. **13**, 70 (2017).
71. Zhang, R. et al. A High Level of Chloroplast Genome Sequence Variability in the Sawtooth Oak *Quercus Acutissima*. *Int. J. Biol. Macromol.* **152**, 340-348 (2020).
72. Liu, X., Chang, E., Liu, J. & Jiang, Z. Comparative Analysis of the Complete Chloroplast Genomes of Six White Oaks with High Ecological Amplitude in China. *J. For. Res.* **32**, 2203-2218 (2021).
73. Yang, Y., Hu, Y., Ren, T., Sun, J. & Zhao, G. Remarkably Conserved Plastid Genomes of Quercus Group *Cerris* in China: Comparative and Phylogenetic Analyses. *Nord. J. Bot.* **36**, e1921 (2018).
74. Li, Y. et al. The Complete Plastid Genome of *Magnolia Zenii* and Genetic Comparison to Magnoliaceae Species. *Molecules*. **24**, 261 (2019).
75. Liu, X. et al. Complete Chloroplast Genome Sequence and Phylogenetic Analysis of *Quercus Bawanglingensis* Huang, Li Et Xing, a Vulnerable Oak Tree in China. *Forests*. **10**, 587 (2019).
76. Yang, Y. et al. Comparative Analysis of the Complete Chloroplast Genomes of Five *Quercus* Species. *Front. Plant Sci.* **7**, (2016).
77. Shaw, J., Lickey, E. B., Schilling, E. E. & Small, R. L. Comparison of Whole Chloroplast Genome Sequences to Choose Noncoding Regions for Phylogenetic Studies in Angiosperms: The Tortoise and the Hare III. *Am. J. Bot.* **94**, 275-288 (2007).
78. Khakhlova, O. & Bock, R. Elimination of Deleterious Mutations in Plastid Genomes by Gene Conversion. *Plant J.* **46**, 85-94 (2006).
79. Perry, A. S. & Wolfe, K. H. Nucleotide Substitution Rates in Legume Chloroplast Dna Depend On the Presence of the Inverted Repeat. *J. Mol. Evol.* **55**, 501-508 (2002).

80. Dong, W., Liu, J., Yu, J., Wang, L. & Zhou, S. Highly Variable Chloroplast Markers for Evaluating Plant Phylogeny at Low Taxonomic Levels and for Dna Barcoding. *Plos One*. **7**, e35071 (2012).
81. Dong, W. et al. *Ycf1*, the Most Promising Plastid Dna Barcode of Land Plants. *Sci Rep*. **5**, 8348 (2015).
82. Group, C. P. W. et al. A Dna Barcode for Land Plants. *Proceedings of the National Academy of Sciences*. **106**, 12794-12797 (2009).
83. Yang, J. et al. Development of Chloroplast and Nuclear Dna Markers for Chinese Oaks (*Quercus* Subgenus *Quercus*) and Assessment of their Utility as Dna Barcodes. *Front. Plant Sci*. **8**, 816 (2017).
84. Zecca, G. et al. The Timing and the Mode of Evolution of Wild Grapes (*Vitis*). *Mol. Phylogenet. Evol.* **62**, 736-747 (2012).
85. Castle, J. Snps Occur in Regions with Less Genomic Sequence Conservation. *Plos One*. **6**, e20660 (2011).
86. Matsuoka, Y., Yamazaki, Y., Ogiwara, Y. & Tsunewaki, K. Whole Chloroplast Genome Comparison of Rice, Maize, and Wheat: Implications for Chloroplast Gene Diversification and Phylogeny of Cereals. *Mol. Biol. Evol.* **19**, 2084-2091 (2002).
87. Yang, Y., Zhou, T., Qian, Z. & Zhao, G. Phylogenetic Relationships in Chinese Oaks (Fagaceae, *Quercus*): Evidence From Plastid Genome Using Low-Coverage Whole Genome Sequencing. *Genomics*. **113**, 1438-1447 (2021).
88. Curtu, A. L., Gailing, O. & Finkeldey, R. Evidence for Hybridization and Introgression within a Species-Rich Oak (*Quercus* Spp.) Community. *Bmc Evol. Biol.* **7**, 218 (2007).
89. Li, H. et al. Plastid Phylogenomic Insights Into Relationships of All Flowering Plant Families. *Bmc Biol.* **19**, 232 (2021).
90. Grímsson, F. et al. Fagaceae Pollen From the Early Cenozoic of West Greenland: Revisiting Engler's and Chaney's Arcto-Tertiary Hypotheses. *Plant Syst. Evol.* **301**, 809-832 (2015).
91. Hipp, A. L. et al. Genomic Landscape of the Global Oak Phylogeny. *New Phytol.* **226**, 1198-1212 (2020).
92. Deng, M., Zhou, Z. K. & Li, Q. S. Taxonomy and Systematics of *Quercus* Subgenus *Cyclobalanopsis*. *Int Oaks*. **24**, 48-60 (2013).
93. Xiaolong, J. *Phylogenetic Relationship and Population Genetic Structure of Quercus Chungii and Q. Championii*. Changsha: Central South University of Forestry and Technology, 2020.

## Figures



**Figure 1**

CPGs Gene map of four *Quercus* section *Cyclobalanopsis* species. Different colors refer to different functions of genes.

## Inverted Repeats



Figure 2

CPGs boundary comparison in 6 *Quercus* section *Cyclobalanopsis* species. Genes shown above the lines were transferred in reverse and those displayed below were transferred forward.

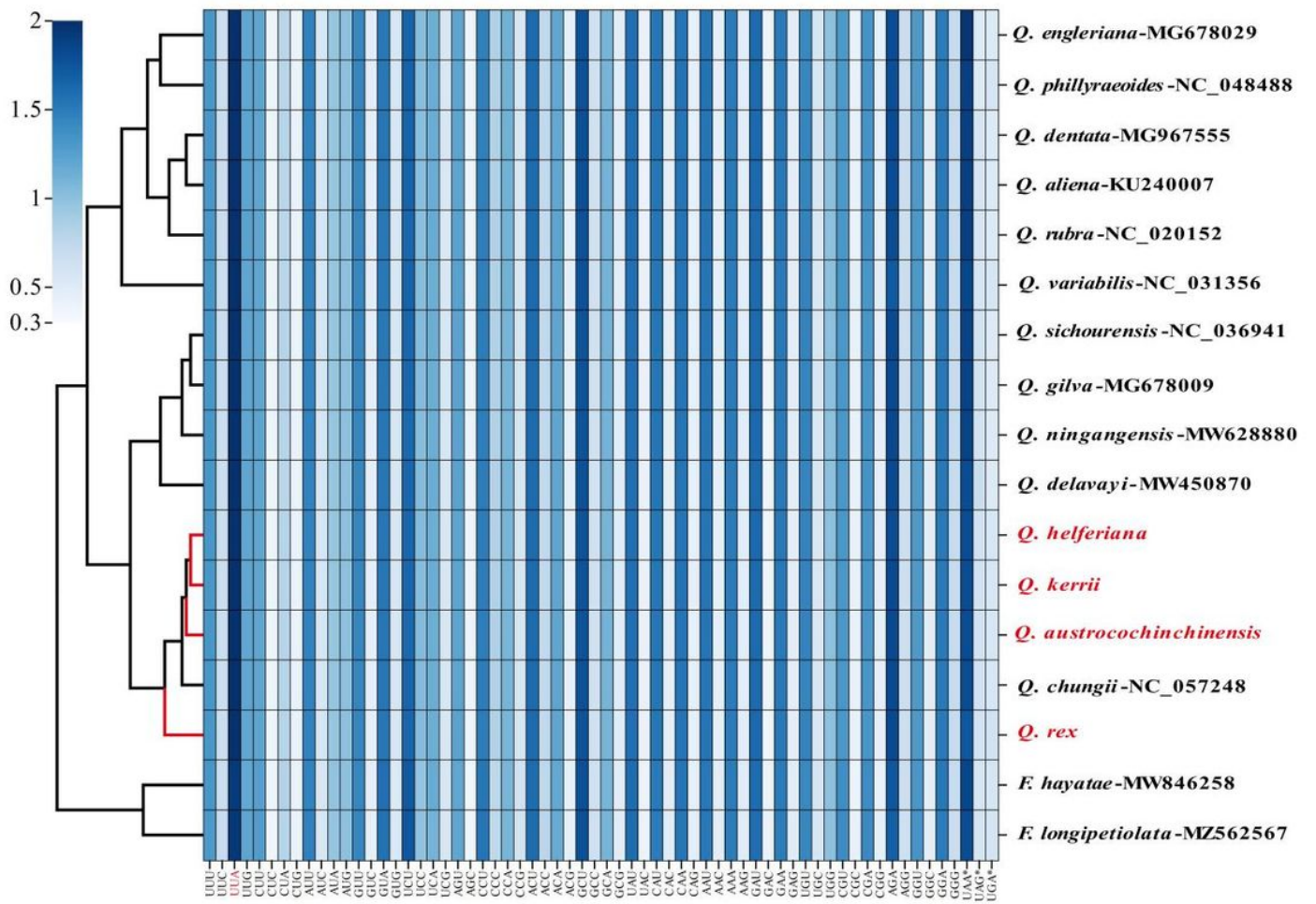
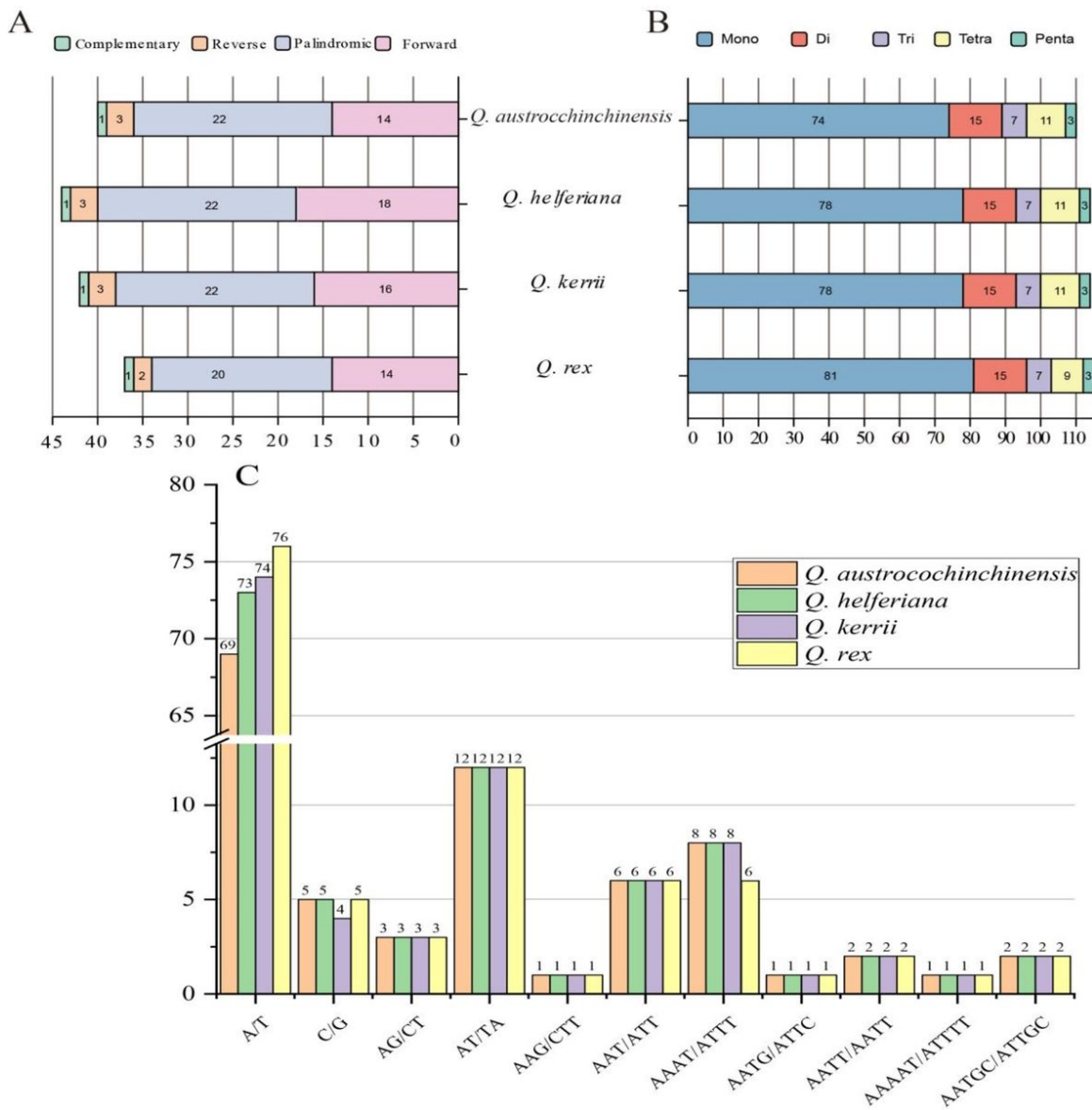


Figure 3

RSCU ratios of CDS genes for CPGs in four species and its sibling species of Fagaceae. (\*) indicated the stop codon



**Figure 4**

Distribution of repeats in 4 *Quercus* section *Cyclobalanopsis* samples. **A** Type and number of LSRs. **B** Distribution of SSRs types. **C** The number of SSR units detected.



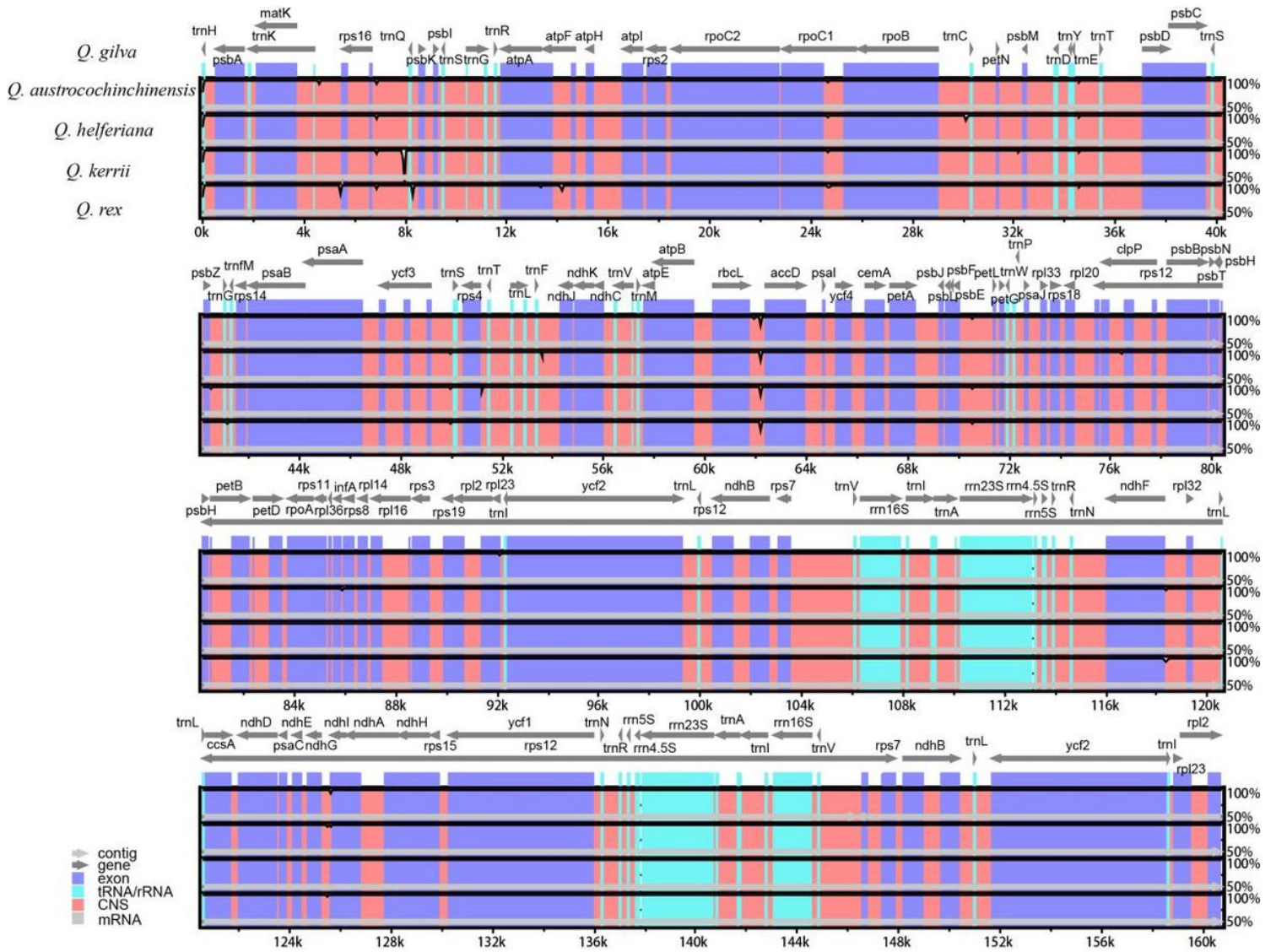
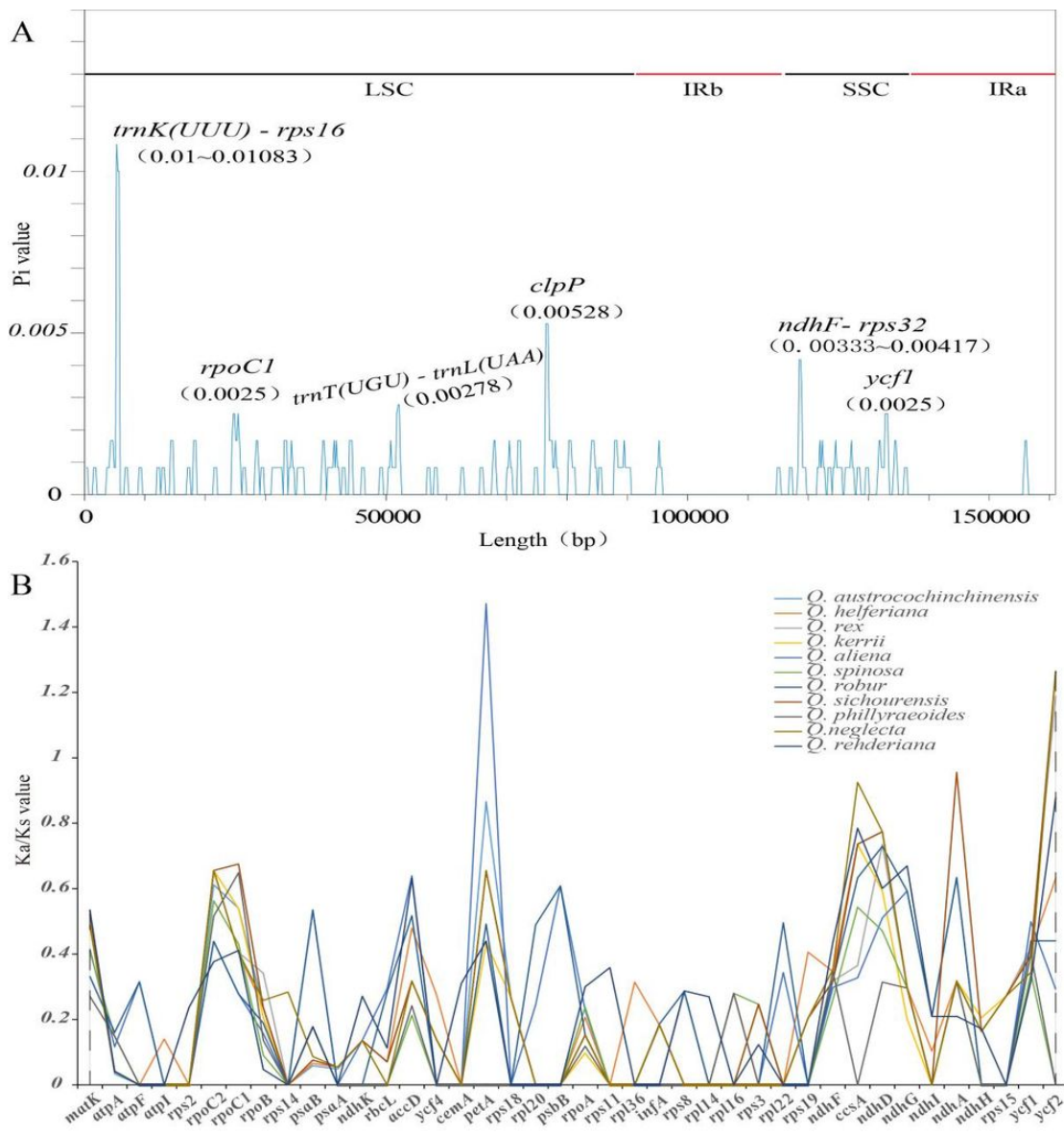


Figure 5

Sequence alignment of the CPGs of four *Quercus* section *Cyclobalanopsis* species. The *Quercus. gilva*(MG678009) was used as reference.



**Figure 6**

**A** Pi values in the multiple alignments of 4 CPGs, details in Supplementary Table S5; **B** Ka/Ks of 40 protein-coding genes (details in Supplementary Table S6) in 11 *Quercus* CPGs

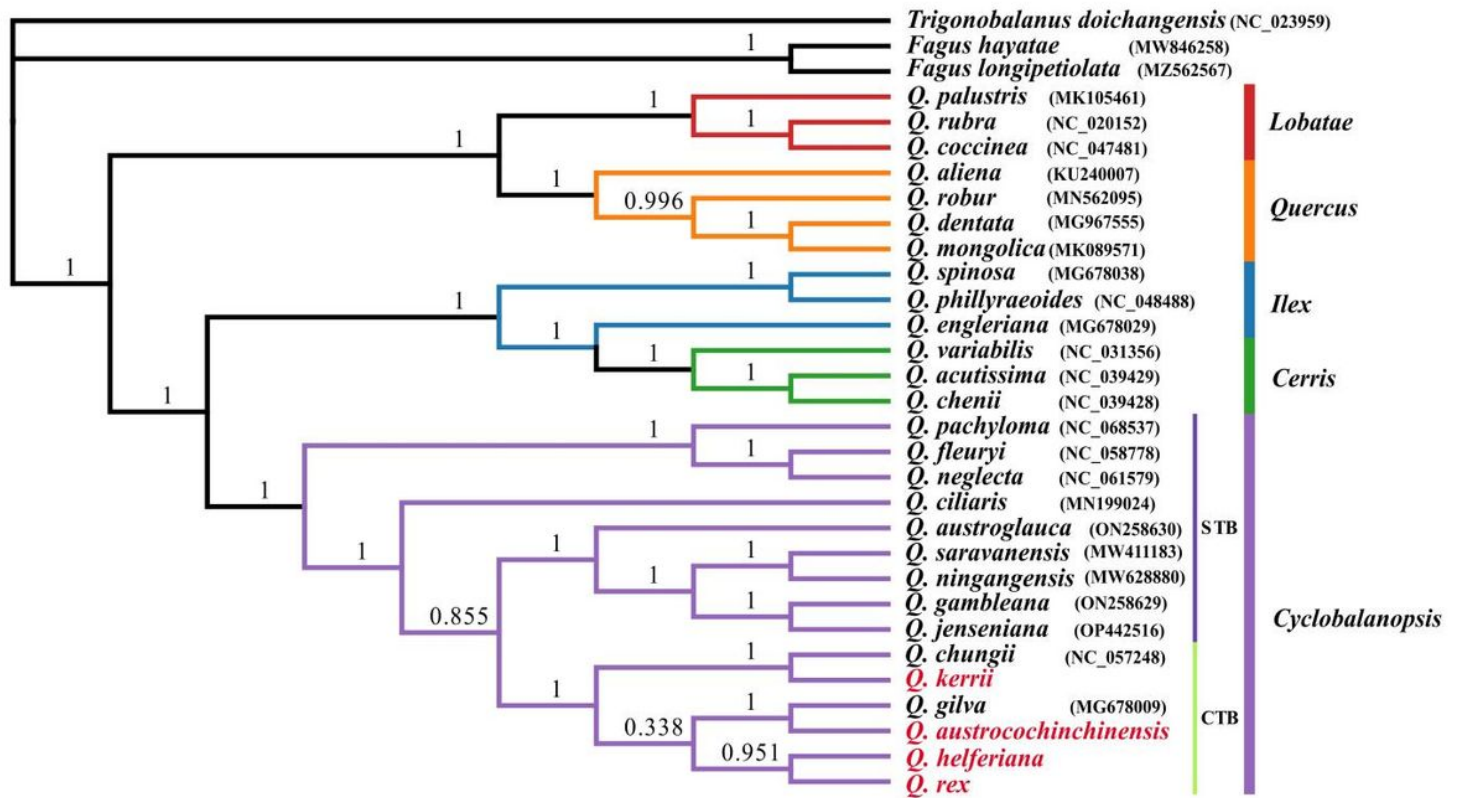


Figure 7

Bayesian(BI) analysis phylogenetic tree among 31 CPGs of Fagaceae species. Values above the branch represented bootstrap support

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryTableS1GenelistintheCPGsof4Quercussect.Cyclobalanopsispecies.xlsx](#)
- [SupplementaryTableS2CodonusageandrelativesynonymouscodonusageRSCUvaluesofproteinencodinggenesofthe16Fagaceae.CPGs.xlsx](#)
- [SupplementaryTableS3Comparisonofdispersedrepeatsamong4QuercussectionCyclobalanopsisCPGs.xlsx](#)
- [SupplementaryTableS4ComparisonofSSRsamong4QuercussectionCyclobalanopsisCPGs.xlsx](#)
- [SupplementaryTableS5ThenucleotidevariabilityPi of4QuercussectionCyclobalanopsisCPGs.xlsx](#)
- [SupplementaryTableS6KaKsvaluesforproteinencodinggenesusingQuercuschenii as a reference.xlsx](#)
- [SupplementaryTableS7Informationonthechloroplastgenomesusedinthisstudy.xlsx](#)