

# Anticodon Table of the Chloroplast Genome and Identification of Putative Quadruplet Anticodons in Chloroplast tRNAs

Tapan Kumar Mohanta<sup>1\*</sup>, Yugal Kishore Mohanta<sup>2</sup>, Ahmed Al-Harrasi<sup>1</sup>, Nanaocha Sharma<sup>3</sup>

<sup>1</sup>Natural and Medical Sciences Research Center, University of Nizwa, Nizwa, 616, Oman

<sup>2</sup>Dept. of Applied Biology, School of Biological Sciences, University of Science and Technology Meghalaya, Baridua, 793101, Meghalaya, India

<sup>3</sup>Institute of Bioresources and Sustainable Development, 795001, Imphal, Manipur, India

\*Correspondence: Tapan Kumar Mohanta, [nostoc.tapan@gmail.com](mailto:nostoc.tapan@gmail.com), [tapan.mohanta@unizwa.edu.om](mailto:tapan.mohanta@unizwa.edu.om)

**Running title:** chloroplast tRNA

## Abstract

The chloroplast genome of 5959 species was analyzed to construct the anticodon table of the chloroplast genome. Analysis of the chloroplast transfer ribonucleic acid (tRNA) revealed the presence of a putative quadruplet anticodon containing tRNAs in the chloroplast genome. The tRNAs with putative quadruplet anticodons were UAUG, UGGG, AUAA, GCUA, and GUUA, where the GUUA anticodon putatively encoded tRNA<sup>Asn</sup>. The study also revealed the complete absence of tRNA genes containing ACU, CUG, GCG, CUC, CCC, and CGG anticodons in the chloroplast genome from the species studied so far. The chloroplast genome was also found to encode tRNAs encoding N-formylmethionine (fMet), Ile<sup>2</sup>, selenocysteine, and pyrrolysine. The chloroplast genomes of mycoparasitic and heterotrophic plants have had heavy losses of tRNA genes. Furthermore, the chloroplast genome was also found to encode putative spacer tRNA, tRNA fragments (tRFs), tRNA-derived, stress-induced RNA (tiRNAs), and group I introns. An evolutionary analysis revealed that chloroplast tRNAs had evolved via multiple common ancestors and the GC% had more influence toward encoding the tRNA number in the chloroplast genome compared to the genome size.

**Key words:** Chloroplast, tRNA, Anticodons, Evolution, Quadruplet anticodons

29

30

## 31 **Introduction**

32 The origin of the genetic code and the translation event are considered to be major transition points  
33 in the evolution of biology. The triplet genetic code is hailed as one of the most important and  
34 ultimate evolutionary anchors and an indisputable piece of evidence of life. The triplet genetic code  
35 understands the specific assignment of the amino acids in the translation machinery. It is the  
36 universal manual and a guided dictionary that cells use to translate the corresponding amino acids  
37 into the translating protein. The number of codon combinations on the mRNA can be an astounding  
38 number of feasible protein sequences, from which only a few can be found in nature. It is believed  
39 that the triplet genetic code is universal and degenerate, and accommodates twenty essential amino  
40 acids using sixty-one sense and three stop codons. However, an emerging study has proved that the  
41 “Universal Genetic Code” is no more universal and can be called as canonical [1, 2]. Sometimes nature  
42 enhances the protein functionalities through codon reassignment to incorporate new amino acids.  
43 This has led to the discovery of the role of selenocysteine (Sec) and pyrrolysine (Pyl) amino acid in  
44 the protein, through the assignment of the stop codon as the sense codon. However, sense codon  
45 reassignment requires low frequency codons, and hence, stop codons are used for this purpose. It  
46 has been demonstrated that except for the triple codons, the *Escherichia coli* ribosome can  
47 accommodate codons and anticodons of variable sizes [3]. Taking this opportunity they have  
48 translated four base codon pairs CCCU, AGGA, UAGA, and CUAG using the four base anticodons  
49 [3]. Frame-shifting of the +1 nucleotide is most favorable in the absence of suppressor tRNA in *E. coli*  
50 [3]. The study also reveals that frame maintenance during translation is not absolute and a frame shift  
51 can be promoted by mutant tRNAs and it can occur with high frequencies at the programmed site of  
52 the mRNA [4, 5]. Riddle and Carbon (1973) reported the presence of four base anticodons CCCC in  
53 tRNA<sup>Gly</sup> instead of the wild-type CCC [6]. A study conducted by Mohanta *et al.*, (2020) revealed the  
54 presence of nine nucleotide anticodons instead of seven nucleotides [7]. These features in the tRNAs  
55 certainly explain the presence of extended codons and anticodons. Most possibly these kind of  
56 evolutionary scenarios exist in codons and tRNAs, to meet the novel translational demand.

57 The availability of enormous genome sequencing data is quite valuable to dig deep into the molecular  
58 features of the protein translation machinery. Significant studies have been performed in the field of  
59 codons and tRNAs (anticodons) and yet, a number of things need to be explored. Taking this  
60 opportunity, we have conducted a large-scale study to deduce the anticodon table of the chloroplast  
61 genome, to understand the presence of reduced or extended genetic codes/anticodons in tRNAs.  
62 Furthermore, we have also tried to understand the presence of Sec and Pyl tRNAs, which are part of  
63 the extended genetic code. A Furthermore investigation has also been conducted to understand the  
64 presence of different introns and the presence of a possible spacer tRNA and tRNA fragments.

## 65 **Results**

### 66 *tRNAs with ACU, CUG, GCG, CUC, CCC, and CGG anticodons are absent in the chloroplast genome*

67 Analysis of the chloroplast genome of the 5959 species from Algae (303), Bryophyte (69), Eudicot  
68 (3832), Gymnosperm (153), Magnoliids (182), Monocot (1177), Nymphaeales (34), protist (57),  
69 Pteridophyte (139), and unknown (13) led to the discovery of 215966 tRNA genes. We did not find  
70 any tRNA encode for ACU, CUG, GCG, CUC, CCC, and CGG anticodons from them (Table 1).  
71 Furthermore, we also found several anticodons, which were seemingly very rare in the chloroplast  
72 genome. They were AGU (tRNA<sup>Thr</sup>), AAG (tRNA<sup>Leu</sup>), CGC (tRNA<sup>Ala</sup>), UCA (tRNA<sup>Sup</sup>), AGG  
73 (tRNA<sup>Pro</sup>), AUU (tRNA<sup>Asn</sup>), UAU (tRNA<sup>Ile</sup>), AUA (tRNA<sup>Tyr</sup>), CAG (tRNA<sup>Leu</sup>), CUU (tRNA<sup>Lys</sup>), CCU  
74 (tRNA<sup>Arg</sup>), AAU (tRNA<sup>Ile</sup>), and GAG (tRNA<sup>Leu</sup>) (Table 1, Supplementary File 1). The tRNA with  
75 anticodons AAG, AGU, and CGC was found only once, whereas, the tRNA with anticodon UCA,  
76 AGG, and AUU was found twice for each (Table 1, Supplementary File 1, Supplementary File 2).  
77 However, the percentage of the CAU (5.47%, tRNA<sup>Met</sup>) anticodon was the highest among all the 64  
78 anticodons. The abundance of the CAU anticodon was followed by GUU, UGC, ACG, and others  
79 (Table 1).

### 80 *Chloroplast genome encodes tRNA for N-formylmethionine, Ile2, Selenocysteine, and Pyrrolysine*

81 A study revealed, a chloroplast genome was found to encode tRNAs for tRNA<sup>fMet</sup>, tRNA<sup>Ile2</sup>, tRNA<sup>Sel</sup>,  
82 and tRNA<sup>Pyl</sup> (Table 1). The tRNA<sup>fMet</sup> was encoded by the same CAU anticodon that coded tRNA<sup>Met</sup>.  
83 We found 709 (0.33%) genes that encoded tRNA<sup>fMet</sup> (Table 1). Also, tRNA<sup>Ile</sup> encoded by the CAU

84 anticodon was commonly referred to as tRNA<sup>Ile2</sup> (Table 1). We found at least 10575 (4.93%) tRNA  
85 genes encoding tRNA<sup>Ile2</sup> (Table 1). Selenocysteine amino acid was encoded by a previously known  
86 stop codon UCA. At least, 204 chloroplast genes were found to encode the UCA anticodon for tRNA<sup>Sel</sup>  
87 (Table 1, Supplementary File 3). A chloroplast genome was also found to encode 197 genes for CUA  
88 anticodons that encoded tRNA<sup>Py1</sup> (Table 1). However, we did not find any CUA anticodon that  
89 encoded the suppressor tRNA (Table 1).

### 90 *Chloroplast genome encodes putative duplet and quadruplet anticodons*

91 We have already mentioned that the triplet genetic code is not universal, it is canonical. Therefore, it  
92 is possible that the genome might have suppressed or extended the genetic code, which is yet to be  
93 elucidated, to a greater extent. In our study, we have found that the chloroplast genome encodes the  
94 putative duplet and quadruplet anticodons (Supplementary File 4, Supplementary File 5). The  
95 annotation of tRNA with quadruplet anticodon had been found when chloroplast genomes were  
96 annotated in the GeSeq chlorobox (<https://chlorobox.mpimp-golm.mpg.de/geseq.html>). However,  
97 re-analysis of the tRNA with the quadruplet anticodon in tRNAscan-SE did not result in a tRNA with  
98 a quadruplet anticodon, which might be due to the default setting for identification of a tRNA with  
99 a triplet anticodon. We are the first to report the presence of duplet and quadruplet anticodons in the  
100 chloroplast genome of the plant kingdom. We found that at least 91 species were encoded quadruplet  
101 anticodons (Supplementary File 4). The quadruplet anticodons were UAUG, UGGG, AUAA, GCUA,  
102 and GUUA (Supplementary File 4). The quadruplet anticodon GUUA found in *Gossypium sturtianum*  
103 (NC\_023218.1) putatively encoded tRNA<sup>Asn</sup>. Similarly, at least 13 species were found to encode duplet  
104 (two nucleotides) anticodons in the tRNAs of the chloroplast genome (Supplementary File 5). Among  
105 them, there were at least eight putative unique duplet anticodons namely UG, AG, AU, CA, GA, GG,  
106 GU, and UA (Supplementary File 5). The putative duplet anticodons might have been caused by the  
107 loss of a nucleotide from the anticodon, because, if there were duplet anticodons, the genome could  
108 encode only 16 anticodons in its genome and would not be able to accommodate all the 20 coding  
109 amino acids in the protein. However, there is a high possibility of having quadruplet anticodons in  
110 the tRNAs, because, in a quadruplet anticodon table, there are 256 possibilities to encode different  
111 amino acids into the protein (Supplementary Table 1).

## 112 *Parasitic organisms have lost the tRNA genes in their chloroplast genome*

113 We found that some of the chloroplast genomes had lost the tRNA genes. The species that have been  
114 found to have lost the tRNA genes are *Pilostyles aethiopica* (NC\_029235.1) (Figure 1) and *Pilostyles*  
115 *hamiltonii* (NC\_029236.1) (Supplementary File 6). *Pilostyles aethiopica* and *Pilostyles hamiltonii* are  
116 endoparasitic plants. Furthermore, some other plants have encoded a fewer number of tRNAs in their  
117 chloroplast genome (Supplementary File 6). They are *Asarum minus* (5), *Gastrodia elata* (5), *Sciaphila*  
118 *densiflora* (6), *Epirixanthes elongata* (8), *Burmannia oblonga* (8), *Lecanorchis japonica* (8), *Lecanorchis*  
119 *kiusiana* (9), and *Selaginella tamariscina* (9)(Supplementary File 6). All of the mentioned species  
120 encoded less than 10 tRNA genes in their chloroplast genome. *Gastrodia elata* is a saprophyte,  
121 whereas, *Sciaphila densiflora*, *Epirixanthes elongate*, *Burmannia oblonga*, *Lecanorchis japonica*, and  
122 *Licanorchis kiusiana* are mycoheterotrophs, and *Cystopteris chinensis* is an endangered species.

123 The chloroplast genome of *Asarum minus* encoded UUU (tRNA<sup>Lys</sup>), UUG (tRNA<sup>Gln</sup>), GCU (Trn<sup>Ser</sup>),  
124 UCC (tRNA<sup>Gly</sup>), and UCU (tRNA<sup>Arg</sup>); *Gastrodia elata* encoded UUG (tRNA<sup>Gln</sup>), GCA (tRNA<sup>Cys</sup>),  
125 UUC(tRNA<sup>Glu</sup>), CAU(tRNA<sup>fMet</sup>), and CCA(tRNA<sup>Trp</sup>); *Sciaphila densiflora* encoded UUG (tRNA<sup>Gln</sup>),  
126 CAU (tRNA<sup>Ile</sup>), CCA(tRNA<sup>Trp</sup>), CAU(Trn<sup>fMet</sup>), UUC(tRNA<sup>Glu</sup>), and GCA(tRNA<sup>Cys</sup>); *Epirixanthes*  
127 *elongata* encoded CCA (tRNA<sup>Trp</sup>), CAU(tRNA<sup>fMet</sup>), UUG(tRNA<sup>Gln</sup>), GUC(tRNA<sup>Asp</sup>), GUA(tRNA<sup>Tyr</sup>),  
128 and UUC(tRNA<sup>Glu</sup>); *Burmannia oblonga* encoded UUG (tRNA<sup>Gln</sup>), GCA (tRNA<sup>Cys</sup>), GUA (tRNA<sup>Tyr</sup>),  
129 UCC (tRNA<sup>Glu</sup>), CAU (tRNA<sup>fMet</sup>), GUG (tRNA<sup>His</sup>), and CAU (tRNA<sup>Ile</sup>); *Lecanorchis japonica* encoded  
130 UUG (tRNA<sup>Gln</sup>), GCA (tRNA<sup>Cys</sup>), GUC (tRNA<sup>Asp</sup>), CAU (tRNA<sup>fMet</sup>), GAA (tRNA<sup>Phe</sup>), CAU (tRNA<sup>Ile</sup>),  
131 and GUU (tRNA<sup>Asn</sup>); *Lecanorchis kiusiana* encoded UUG (tRNA<sup>Gln</sup>), GCA (tRNA<sup>Cys</sup>), GUC (tRNA<sup>Asp</sup>),  
132 UUC (tRNA<sup>Glu</sup>), CAU (tRNA<sup>fMet</sup>), GAA (tRNA<sup>Phe</sup>), CAU (tRNA<sup>Ile</sup>), and GUU (tRNA<sup>Asn</sup>); and *Selaginella*  
133 *tamariscina* encoded GUG (tRNA<sup>His</sup>), GUC (tRNA<sup>Asp</sup>), GUA (tRNA<sup>Tyr</sup>), UUC (tRNA<sup>Glu</sup>), GUU  
134 (tRNA<sup>Asn</sup>), and CCA (tRNA<sup>Trp</sup>). These species encoded only 14 anticodons CAU, CCA, GAA, GCA,  
135 GCU, GUA, GUC, GUG, GUU UCC, UCU, UUC, UUG, and UUU.

## 136 *Chloroplast genome encode putative spacer tRNAs*

137 Spacer RNA genes are usually found in the spacer region, between the 16S and 23S rRNAs, in  
138 bacterial genomes. When we focused our study on the spacer RNA in the chloroplast genome, we  
139 found that chloroplast genomes were also encoded in the putative spacer tRNAs between the 16S

140 and 23S rRNA genes. tRNA<sup>Ala</sup>(UGC) and tRNA<sup>Ile</sup> (GAU) were the most predominant spacer tRNAs  
141 found in the chloroplast genome (Figure 2). The percentages of the UCG and GAU anticodons in the  
142 chloroplast genome were 5.13 and 4.98, respectively. This showed that spacer tRNAs were more  
143 common in the chloroplast genome. Sometimes, it contained tRNA<sup>fMet</sup> (CAU) and tRNA<sup>Ser</sup> (GCU) in  
144 the spacer region. All the chloroplast genomes did not encode the spacer tRNAs (Supplementary File  
145 7). None of a mycoparasitic plants was found to encode the putative spacer tRNA in their chloroplast  
146 genome. However, the majority of the species encoded putative spacer tRNAs.

### 147 *The Majority of chloroplast tRNAs encode group I intron*

148 It was found that the majority of chloroplast-encoding tRNAs encode introns. Except for tRNA<sup>Arg</sup>,  
149 tRNA<sup>Asn</sup>, tRNA<sup>Asp</sup>, tRNA<sup>Gln</sup>, tRNA<sup>His</sup>, tRNA<sup>Pro</sup>, tRNA<sup>Trp</sup>, and tRNA<sup>Val</sup> all other tRNA genes were found  
150 to contain group I introns (Table 2). The introns found in tRNA seem to be isotype-specific (Table 2).  
151 The introns are conserved within the tRNA isotype and the conserved nucleotide sequences of the  
152 introns of one isotype do not match with the conserved introns of other isotypes (Table 2). When we  
153 cluster the conserved region of the introns, they form four groups (Supplementary Figure 1). We have  
154 named them group A, B, C, and D. Group A contains tRNA<sup>Leu</sup>, tRNA<sup>Tyr</sup>, and tRNA<sup>Cys</sup>; group B  
155 contains tRNA<sup>Ser</sup>; group C contains tRNA<sup>Lys</sup>, tRNA<sup>Met</sup>, and tRNA<sup>Ala</sup>, and group D contains tRNA<sup>Gly</sup>,  
156 tRNA<sup>Ile</sup>, tRNA<sup>Glu</sup>, and tRNA<sup>Thr</sup> (Supplementary Figure 1). However, the introns of tRNA<sup>Phe</sup> do not  
157 group with any other introns (Supplementary Figure 1).

158

### 159 *Chloroplast genome encode putative novel tRNAs*

160 Although we all are well-acquainted with the fact that tRNA makes a clover leaf-like structure, yet  
161 we found some variations in the tRNA structure. Analysis revealed the presence of a few novel tRNA  
162 structure/tRNA-like molecules (Figure 3 and Figure 4). Some putative novel tRNA-like structures  
163 seemed to lack the anticodon loop, whereas, in some cases they had extra sequences near the  
164 anticodon arm region (Figure 3). A tRNA-like structure contained an extended nucleotide sequence  
165 in the region between the D-arm and anticodon arm (Figure 4). At least 42 species were found to  
166 encode novel tRNA-like structures that contained extended nucleotide sequences between the D-arm  
167 and anticodon arm (Figure 4). Furthermore, a few tRNAs were found to have lost the pseudouridine



168 loop (Figure 5), suggesting the presence of novel tRNAs/tRNA-like structures in the chloroplast  
169 genome.

### 170 *Chloroplast genome encodes putative tRNA Fragments (tRFs)*

171 The tRFs are small 14-32 nucleotides novel class of small, non-coding RNAs, derived from the mature  
172 or precursor tRNAs that are different from the tRNA-derived, stress-induced tRNAs (tiRNAs) [8, 9].  
173 Analysis revealed the presence of at least 55 tRFs in the chloroplast genome. The tRFs found were for  
174 tRNA<sup>Glu</sup>, tRNA<sup>Arg</sup>, tRNA<sup>Gly</sup>, tRNA<sup>His</sup>, tRNA<sup>Val</sup>, tRNA<sup>Ile</sup>, tRNA<sup>Thr</sup>, tRNA<sup>Leu</sup>, tRNA<sup>Lys</sup>, and tRNA<sup>Ala</sup>  
175 (Supplementary File 8). The tRFs of tRNA<sup>Glu</sup> were found to contain conserved nucleotide sequence  
176 GGCCTTATCGTCTAGTGAT, whereas, those of tRNA<sup>Gly</sup> were found to contain conserved  
177 GCGGGTATAGTTTGTGGTAAA nucleotides (Supplementary File 8). As such, we did not find  
178 conserved nucleotide sequences for the other tRFs. The tRFs of tRNA<sup>Ala</sup>, tRNA<sup>Gly</sup>, tRNA<sup>Ile</sup>, tRNA<sup>Lys</sup>,  
179 and tRNA<sup>Leu</sup> were 5'-tRFs, whereas, the tRFs of tRNA<sup>His</sup>, tRNA<sup>Thr</sup>, and tRNA<sup>Val</sup> were 3'-tRFs. The tRFs  
180 of tRNA<sup>Glu</sup> did not match either the 5'- or 3'-end of the tRNA, and hence, might have originated from  
181 the precursor tRNA transcript. Therefore, they can be classified as tRF-1.

### 182 *Chloroplast genome encode putative tiRNAs*

183 The longer tRFs (tRNA fragments) of 30–50 nucleotide-long sequences are called tRNA-derived,  
184 stress-induced RNAs (tiRNAs)[8]. Therefore, we searched for the presence of 30–50 nucleotide tRFs.  
185 We found at least 244 tRNA sequences, which encoded the 30–50 nucleotides (Supplementary File 9).  
186 The tiRFs were part of putative tRNA<sup>Ala</sup> (UGC), tRNA<sup>Phe</sup> (GAA), tRNA<sup>fMet</sup> (CAU), tRNA<sup>Gly</sup> (GCC,  
187 UCC), tRNA<sup>His</sup> (GUG), tRNA<sup>Ile</sup> (CAU, GAU), tRNA<sup>Lys</sup> (UUU), tRNA<sup>Leu</sup> (UAA), tRNA<sup>Asn</sup> (GUU), and  
188 tRNA<sup>Val</sup> (GAC, UAC) (Supplementary File 9). Among them, tiRFs of tRNA<sup>His</sup> (GUG) and tRNA<sup>fMet</sup>  
189 (CAU) were found only once, whereas, tRNA<sup>Lys</sup> (UUU) was the highest (72) encoding tiRF. The tiRFs  
190 of tRNA<sup>Lys</sup> (UUU) was followed by tRNA<sup>Ile</sup> (GAU) and tRNA<sup>Ala</sup> (UGC), which were found to contain  
191 51 and 52 putative tiRFs, respectively (Supplementary File 9).

192

193

194 ***Machine Machine-learning approach showed GC% influences the tRNA number in the chloroplast***  
195 ***genome***

196 We grouped the chloroplast genomes of all the species according to their clade and conducted a  
197 comparative study. The analysis revealed that the average tRNA gene number in monocot (37.80%)  
198 plants is comparatively higher than that in other plants (Supplementary File 6). The protists showed  
199 the lowest (29.5%) average tRNA gene number, followed by algae (30.12%) (Supplementary File 6).  
200 A correlation analysis of the GC% with the tRNA number showed a positive correlation ( $r = 0.362$ )  
201 for the monocot clade (Figure 6). The chloroplast genomes of the species *Isolepis setacea* and *Vitis*  
202 *romanetii* were found to encode the highest number of tRNAs, that is, 52 each (Supplementary File 6).  
203 On an average, the chloroplast genomes were found to encode 36 tRNA genes per genome. A  
204 machine-learning approach was used to understand the role of the GC content and genome size in  
205 the tRNA number in the chloroplast genome. The boosting analysis revealed that the relative  
206 influence of the GC% was more than the genome size (Figure 7). A principal component analysis was  
207 conducted to see their association with different clades.

208 ***Chloroplast tRNAs Evolve from Multiple Common Ancestors***

209 We conducted a phylogenetic analysis by considering the tRNA genes of the chloroplast genome.  
210 The phylogenetic analysis revealed clear and distinct phylogenetic clusters of tRNAs. The  
211 phylogenetic tree showed two major distinct clusters suggesting their origin from multiple common  
212 ancestors (Figure 8). In cluster I, anticodons GCU, GGA, UGA, GCC, UCC, CGU, CGA, GCC, CGA,  
213 CGU, UUC, UCU, CAU, UAA, CAA, GUA, UAG, UAU, UAUG, CAA, GCU, UCG, UCU, GAA, CUA,  
214 UAG, and GAG, grouped together, whereas, in cluster II, anticodons UUG, GUG, GCA, GAA, UUU,  
215 GUU, UGG, GGG, CCA, UGU, GGU, CAU, UAC, GCC, GUC, GAC, GAU, UUC, CGU, ACG, CCG,  
216 ACA, and UGC, grouped together (Figure 8). The anticodons GAA (tRNA<sup>Phe</sup>), CAU (tRNA<sup>Met</sup>), GCC  
217 (tRNA<sup>Gly</sup>), UUC (tRNA<sup>Glu</sup>), and CGU (tRNA<sup>Thr</sup>) were shared in both the clusters. The phylogenetic  
218 analysis of quadruplet anticodons revealed that quadruplet anticodon AUAA shares a phylogenetic  
219 relationship with UAUG anticodons, whereas, the UGGG and GUUA anticodons fall in a distinct  
220 cluster (Figure 9).



221 Genes undergo mutation, which is a common phenomenon. Although it was a common phenomenon  
222 in coding genes, non-coding genes also showed frequent mutation. Therefore, a  
223 transition/transversion bias study was conducted for the chloroplast tRNAs. The analysis revealed  
224 that transition predominates transversion (Supplementary Table 2). The transition/transversion bias  
225 was found to be the highest for tRNA<sup>Asn</sup> (R = 13.71), whereas, tRNA<sup>Ser</sup> (1.22) had the lowest bias  
226 (Supplementary Table 2). The transition/transversion bias of tRNA<sup>Asn</sup> was followed by tRNA<sup>Tyr</sup>  
227 (11.51) and tRNA<sup>Trp</sup> (8.63). Although, tRNA<sup>Arg</sup>, tRNA<sup>Leu</sup>, and tRNA<sup>Ser</sup> encoded six Isoacceptors, their  
228 transition/transversion bias was comparatively lower than that of others (Supplementary Table 2).

## 229 Discussion

230 The chloroplast genome harbors several coding sequences and a few non-coding sequences including  
231 rRNA and tRNA. These genetic elements and their potential to translate codons make them semi-  
232 autonomous organelles of the plant cell. A detailed genomic analysis of the chloroplast tRNA reveals  
233 that it does not encode all the 64 anticodons required for the tRNAs. The tRNAs with anticodons  
234 ACU (tRNA<sup>Ser</sup>), CUG (tRNA<sup>Gln</sup>), GCG (tRNA<sup>Arg</sup>), CUC (tRNA<sup>Glu</sup>), CCC (tRNA<sup>Gly</sup>), and CGG (tRNA<sup>Pro</sup>),  
235 are absent in the chloroplast genome of the studied species. therefore, these anticodons can be  
236 classified as rare anticodons of the chloroplast genome. The ACU anticodon of tRNA<sup>Ser</sup> and the GCG  
237 anticodon of tRNA<sup>Arg</sup> are from the hexa-isoacceptor group, whereas, the CCC anticodon of tRNA<sup>Gly</sup>  
238 and the CGG anticodon of tRNA<sup>Pro</sup> are from the tetra-isoacceptor group. Therefore, a lack of these  
239 anticodons from their isoacceptor group does not make any difference in the genome as other  
240 isoacceptors are available for their use, to encode the codon. However, tRNA<sup>Gln</sup> is encoded only by  
241 CUG and UUG anticodons, whereas, tRNA<sup>Glu</sup> is encoded by the CUC and UUC anticodons. The lack  
242 of the CUG anticodon from tRNA<sup>Gln</sup> and the CUC anticodon from tRNA<sup>Glu</sup> in the chloroplast genome  
243 has left these tRNA isotypes with only one choice of anticodon (Table 1). The lack of the CUG  
244 anticodon in tRNA<sup>Gln</sup> and the CUC anticodon in tRNA<sup>Glu</sup>, in the chloroplast genome, may be due to  
245 a strong selection pressure to establish UUG (tRNA<sup>Gln</sup>) and UUC (tRNA<sup>Glu</sup>) anticodons as the  
246 dominant anticodons. The tRNA anticodons followed by nucleotides CUx (x = any nucleotide) may  
247 have undergone a strong evolutionary pressure, and hence, anticodons CUA, CUU, CUG, and CUC,  
248 encode only 197, 7, 0, and 0 anticodons, respectively, in the chloroplast genome (Table 1). However,

249 the CAU anticodon encoding tRNA<sup>Met</sup> has been seen to have the highest percentage (5.47%) in the  
250 chloroplast genome (Supplementary File 1). The CAU anticodon of tRNA<sup>Met</sup>, of the nuclear encoded  
251 genome has also been found in the highest (5.03%) abundance [10], thus corroborating CAU, as the  
252 most abundant anticodon in the nuclear and chloroplast genomes. The anticodons CAU (tRNA<sup>Met</sup>),  
253 GUU (tRNA<sup>Asn</sup>), UGC (tRNA<sup>Ala</sup>), and ACG (tRNA<sup>Arg</sup>) have been found to encode more than 5% each  
254 of the total anticodons, suggesting the role of positive selection pressure in these anticodons  
255 (Supplementary File 1). However, at the isotype/isodecoder level, tRNA<sup>Leu</sup> (10.27%) has been found  
256 to contain the highest percentage of anticodons followed by tRNA<sup>Ile</sup> (9.93%) and tRNA<sup>Arg</sup> (7.96%)  
257 (Supplementary File 1). A similar level of abundance has been found for tRNA<sup>Leu</sup> (7.80%), for the  
258 nuclear encoded tRNA genes, reflecting a similarity in the anticodon abundance in the nuclear and  
259 chloroplast genomes [10]. However, an abundance of the nuclear-encoded anticodons tRNA<sup>Leu</sup> is  
260 followed by tRNA<sup>Ser</sup> (7.66%), tRNA<sup>Gly</sup> (7.52%), and tRNA<sup>Arg</sup> (7.28%) [10]. Although, tRNA<sup>Leu</sup> is the  
261 highest encoding isotype/isodecoder in nuclear- (7.80%) and chloroplast (10.27%)-encoded genomes,  
262 there is a great difference in their percentage. The chloroplast-encoded CAU anticodon also encodes  
263 tRNA<sup>Ile2</sup> (4.93%). The CAU anticodon for tRNA<sup>fMet</sup> (0.33%) is also quite abundant in the chloroplast  
264 genome. The tRNA<sup>fMet</sup> acts as an initiation anticodon in protein synthesis in mitochondria, bacteria,  
265 and chloroplasts and the presence of tRNA<sup>fMet</sup> in the chloroplast genome is quite justified. However,  
266 only 709 tRNA<sup>fMet</sup> genes were found during the analysis suggesting that tRNA<sup>fMet</sup> is not a universal  
267 tRNA of the chloroplast genome. A majority percentage of the chloroplast genome does not encode  
268 tRNA<sup>fMet</sup>. A few of the chloroplast genomes encode the tRNAs for selenocysteine and pyrrolysine  
269 amino acid (Table 1). However, Zhao *et al.*, (2021) has reported the absence of tRNA<sup>Sec</sup> in gymnosperm  
270 plants [11]. The Sec amino acid specified by the UGA codon, requires the presence of the  
271 selenocysteine insertion sequence (SECIS) element, and the Pyl amino acid encoded by the UAG  
272 codon requires the pyrrolysine insertion sequence (PYLIS) [12]. The presence of tRNA for encoding  
273 Sec and Pyl reflects that the chloroplast genome may have SECIS and PYLIS in it.

274 It was also very peculiar to see the loss of tRNA genes in the chloroplast genome of heterotrophic  
275 and mycoparasitic plants. Our previous study reported the loss of several other genes in the  
276 chloroplast genome in mycoparasitic and heterotrophic plants [13]. Similar is true for the tRNA genes  
277 as well. In the absence of tRNA genes in the chloroplast genome, the cell most probably uses the

278 tRNA genes from the nuclear-encoded genome. However, the loss of tRNA genes in the chloroplast  
279 genome seems independent of the nuclear genome. The parasitic and heterotrophic plants require  
280 less effort to complete their lifecycle, as they are completely dependent on their host. Hence, they do  
281 not need a lot of genes for their function, and hence, may be under constant pressure to eliminate  
282 genes. Therefore, these mycoparasitic and heterotrophic plants contain only 14 (CAU, CCA, GAA,  
283 GCA, GCU, GUA, GUC, GUG, GUU UCC, UCU, UUC, UUG, and UUU) anticodons in their  
284 chloroplast genome.

285 It is well known that the triplet genetic code is canonical and not universal. The genetic code can be  
286 expanded, where specific codons can be re-allocated to encode non-proteogenic amino acids. The  
287 tRNA genes undergo rapid changes to meet the translational demand of the cell [14]. Therefore, it is  
288 highly possible that tRNA can expand its anticodon nucleotide number. Our study helped us to  
289 discover the presence of quadruplet anticodons in the chloroplast genome of at least 91 plant species  
290 (Supplementary File 4). The quadruplet anticodons found in our study were UAUG, UGGG, AUAA,  
291 GCUA, and GUUA. Studies regarding the presence of functional quadruplet anticodons are reported  
292 in a few cases [15–22]. Anderson *et al.*, (2004) reported the role of the quadruplet codon AGGA  
293 through changes in the tRNA anticodon loop to CUUCCUAAA in a suppressor tRNA<sub>CUA</sub> [15]. The  
294 suppression of the amber tRNA led to the encoding of homoglutamine (hGln), using the AGGA  
295 codon [15]. They also reported that quadruplet codons CCCU or CUAG could be used to suppress  
296 the amber tRNA and allow the incorporation of unnatural amino acid into the protein in *Escherichia*  
297 *coli*[15]. Neumann *et al.*, (2010), reported the encoding of unnatural amino acids through the  
298 evolution of the quadruplet anticodon in response to the amber codon tRNA<sub>CUA</sub>[16].  
299 Chloramphenicol resistance was achieved when tRNA<sub>UCCU</sub><sup>Ser2</sup> translated the AAGA codon and  
300 tRNA<sub>UCCU</sub><sup>Ser2</sup> translated the AGGA codon [16]. Niu *et al.*, (2013) replaced tRNA<sup>Pyl</sup><sub>CUA</sub> with the UCCU  
301 anticodon and generated tRNA<sup>Pyl</sup><sub>UCCU</sub>, which recognized and suppressed the quadruplet codon  
302 AGGA [17]. This provided a qualitative notion for the suppression of the quadruplet codon through  
303 tRNA<sub>UCCU</sub> [17]. Most specifically, the presence of the quadruplet anticodon was associated with  
304 suppression of the amber tRNA and incorporation of the unnatural amino acid into the protein chain.  
305 The tRNA<sub>GCUA</sub> contained an additional G nucleotide prior to the tRNA<sub>CUA</sub> anticodon, suggesting its  
306 role in suppression of the amber codon. In the tRNA<sup>Asn</sup><sub>GUUA</sub> anticodon, most probably, nucleotide A

307 was incorporated after the GUU anticodon, as the tRNA with the GUU anticodon was grouped with  
308 the GUUA anticodon in the phylogenetic tree (Figure 9). Similarly, in the UGGG anticodon, the G  
309 nucleotide got incorporated in the UGG anticodon, as they grouped with the UGG anticodon (Figure  
310 9). The GCUA anticodon was grouped with GCU anticodon suggesting that the A nucleotide was  
311 incorporated at the fourth position of the GCU anticodon, which gave rise to the GCUA anticodon  
312 (Figure 9). However, no such clue was found in the case of the UAUG and AUAA anticodons.  
313 Considering, the incorporation of the additional nucleotide at the fourth position, we could speculate  
314 that the G nucleotide was most probably incorporated in the UAU anticodon, and gave rise to the  
315 UAUG anticodon. Similarly, the A nucleotide was incorporated at the 4<sup>th</sup> position of the AUA  
316 anticodon to give rise to the AUAA anticodon. Although, we found only five putative quadruplet  
317 anticodons, the genome could accommodate at least 256 quadruplet anticodons/codons in the cell  
318 (Table 2). We also found the presence of tRNAs, with only duplet anticodon, where one nucleotide  
319 was possibly deleted from the anticodon (Supplementary File 5). At least 13 species resulted that  
320 contained duplet anticodons in the tRNA of the chloroplast genome (Supplementary File 5).

321 The chloroplast encoding tRNAs were also found to encode the group I introns. These group I introns  
322 were conserved in their respective isotype/isodecoder groups (Table 2). From a total of 20 isotypes,  
323 12 of them were found to encode the group I introns (Table 2). However, the group I intron of one  
324 isotype was not conserved with the intron of another isotype, reflecting the isotype-based  
325 conservation of the group I intron, in the tRNA.

326 It is well-reported that group I introns are found in tRNAs, bacteria, lower eukaryotes, and higher  
327 plants [23–25]. Some of the group I intron encode homing endonucleases catalyze intron mobility,  
328 thus facilitating the movement of the intron from one location to another and from one organism to  
329 another [24]. However, the incorporation of the group I intron in the tRNA gene is isotype-specific,  
330 as only 12 isotypes have been found to encode the intron, while eight isotypes do not have any intron  
331 in their tRNAs (Table 2). From the eight isotypes, tRNA<sup>His</sup>, tRNA<sup>Gln</sup>, tRNA<sup>Asp</sup>, tRNA<sup>Asn</sup>, and tRNA<sup>Arg</sup>  
332 belong to the polar group, whereas, tRNA<sup>Trp</sup>, tRNA<sup>Pro</sup>, and tRNA<sup>Val</sup> belong to the non-polar group.  
333 This shows that the presence of the type I intron tends to be more toward the tRNA that encodes  
334 polar amino acids. Furthermore, it is seen that the chloroplast genome also encodes the putative  
335 spacer tRNAs (Figure 2). It is reported that *E. coli* contains a spacer tRNA (tRNA<sup>Ala</sup> and tRNA<sup>Ile</sup>) that

336 is present in the spacer region of the 16S and 23S rRNA [26]. The tRNAs, tRNA<sup>Ala</sup>, and tRNA<sup>Ile</sup>, have  
337 also been found in the spacer region of 16S and 23S rRNA suggesting the presence of a spacer tRNA  
338 in the chloroplast genome. Although, in a majority of cases, tRNA<sup>Ala</sup> and tRNA<sup>Ile</sup> are the predominant  
339 spacer tRNAs; tRNA<sup>Glu</sup> can be the third most possible spacer tRNA of the chloroplast genome.

340 Analysis also revealed the presence of tRNA fragments (tRFs) in the chloroplast genome. We found  
341 at least 55 tRFs that belonged to ten tRNA isotypes (Supplementary File 8). These tRFs were  
342 putatively derived from the tRNA precursors or from the cleavage of mature tRNAs [27]. The tRFs  
343 were reported to control gene expression, translation control, transposon control, ncRNA, and DNA  
344 damage response [8, 27–29]. Although, we found ten different chloroplast-derived tRFs, the majority  
345 of them belonged to tRNA<sup>Glu</sup> and tRNA<sup>Gly</sup> (Supplementary File 8). Among them are the, tRNA<sup>Glu</sup>are  
346 tRF-1 type, tRNA<sup>Gly</sup>are tRF-5'-type, and tRNA<sup>His</sup>, tRNA<sup>Thr</sup>, and tRNA<sup>Val</sup>are tRF-3' type  
347 (Supplementary File 8). Furthermore, we also noted the presence of a few putative tRNA-derived,  
348 stress-induced RNA (tiRNAs) fragments (tiRFs) in the chloroplast genome. The majority of the tiRFs  
349 were from tRNA<sup>Lys</sup> (UUU). For the first time, tiRFs were reported in the human fetus hepatic tissue  
350 and osteosarcoma cells [30, 31]. These tiRFs could be generated in the cell under different stress  
351 conditions via cleavage of mature tRNAs [30]. However, their presence as independent nucleotide  
352 fragments in the annotated genome sequence reflected their independent presence in the genome.  
353 Although, the cleavage of tRNAs to tiRFs was brought about by the enzyme angiogenin (an RNase  
354 superfamily) [31] in the human cell; its counterpart in plants needs to be identified to understand its  
355 detailed functions. The 5'-tiRNA<sup>Ala</sup> and tiRNA<sup>Cys</sup> were reported to inhibit translation in rabbit  
356 reticulocytes [31] suggesting their inhibitory role in protein translation.

357 This study also found the presence of a putative novel tRNA structure encoded by the chloroplast  
358 genome (Figure 4). The tRNA<sup>Gly</sup> (UCC) was found to contain a long nucleotide sequence between the  
359 D-arm and anticodon arm in several species. This long arm could be most probably be an intron that  
360 might have incorporated in between these two arms. The chloroplast tRNAs which had lost the  
361 pseudouridine loop ( $\Psi$ ) seemed to be metazoan mitochondrial-specific (Figure 5). The loss of the  $\Psi$ -  
362 loop in tRNA was first reported in the 1970s [32–34]. Previous studies also reported loss of the  $\Psi$ -  
363 arm and loop in nematode mitochondrial tRNA [34]. However, in the nematode mitochondrial tRNA,  
364 the  $\Psi$ -arm and loop were present in the tRNA<sup>Ser</sup> (GCU), whereas, it had lost the  $\Psi$ -arm and loop

365 in tRNA<sup>Ser</sup> (GCU and GGA) in the chloroplast genome (Figure 5). The elongation factor (EF) Tu  
366 combined with GTP to form a complex that delivered the amino acyl tRNA to the ribosome A site  
367 through binding of the acceptor arm and  $\Psi$ -arm [35]. In the absence of the  $\Psi$ -arm and loop in the  
368 tRNAs, it might be using some alternative binding mode for EF-Tu [36, 37]. In the case of  
369 *Caenorhabditis elegans* mitochondrial EF-Tu, it has around 60 amino acid extensions at the C-  
370 terminal end that might be playing important role in binding tRNAs that lack the  $\Psi$ -arm [38, 39].  
371 This also suggested that the mitochondrial ribosomal protein might have alternate binding sites  
372 for the truncated tRNA. Furthermore, the presence of the metazoan, mitochondria-specific,  
373 truncated tRNA in the chloroplast genome suggested that these tRNA genes might be shared by  
374 sub-cellular organelle chloroplast and mitochondria.

375 Evolutionary analysis revealed, chloroplast tRNAs are derived from multiple common ancestors  
376 (Figure 8). The phylogenetic tree of the chloroplast tRNA shows two distinct clusters, which  
377 reflect their evolution from multiple common ancestors. In cluster I, anticodons GCC, CGU, CGA,  
378 UCU, CAA, and UAG are seen to make more than one group, whereas, none of the anticodons  
379 from cluster II are found to make more than one group (Figure 8). The anticodons GCC, GCU,  
380 UUC, CAU, and GAA are also found in both the clusters (Figure 8). This suggests that tRNAs  
381 with anticodons GCC, CGU, CGA, UCU, CAA, and UAG, of cluster I, may have undergone vivid  
382 duplication and produced more than one anticodon group.

### 383 **Conclusions**

384 Chloroplast is a semiautonomous organelle of the plant and protist kingdom with a great  
385 potential to encode its own genome and protein translation machinery. The important tRNA  
386 molecules require for protein translation process is well documented. Chloroplast genome  
387 encode putative duplet, triplet, and quadruplet anticodons suggesting their role in recognition  
388 of duplet, triplet, and quadruplet codons in the mRNA. Mycoparasitic plants has lost their  
389 chloroplast genome to a large extent thereby losing several chloroplast encoded tRNA genes.  
390 Further, several of the chloroplast encoded tRNA genes were found to encode introns and the  
391 presence of intron in chloroplast genome suggest the presence of introns in the gene of their  
392 prokaryotic ancestor cyanobacteria. Further, the chloroplast genome is very selective and



393 encoded only a few Isoacceptor abundantly while GCG, CUG, CUC, CCC, CGG, and ACU  
394 anticodons were found to be the rarest form of anticodons in the chloroplast genome. It is  
395 important to understand why chloroplast genome do not encode tRNA with such anticodons.

## 396 **Materials and Methods**

397  
398 All the chloroplast genomes were downloaded from the National Center for Biotechnology  
399 Information (NCBI) database. In total, 5959 chloroplast genomes were used in this study. The  
400 downloaded chloroplast genomes were subjected to tRNA annotation. tRNA annotation was  
401 conducted using tRNAscan-SE 2.0, Aragorn and the GeSeq-Annotation of the organellar genomes  
402 [40–42]. The Linux-based approach was used to annotate the chloroplast tRNA for tRNAscan-SE 2.0  
403 and Aragorn. In the GenSeq-annotation of the organellar genome, the chloroplast genome files were  
404 uploaded with the following parameters; sequence source: Plastid; annotation option: Annotate  
405 plastid inverted repeats; blat search: Default; annotate: CDS, tRNA, and rRNA; and third party tRNA  
406 annotator: Aragorn v1.2.38, tRNAscan-SE v2.0.7. All the tRNA sequences generated from these three  
407 annotation pipelines were corroborated and used for further analysis. All the data obtained from  
408 tRNAscan-SE and Aragorn were further processed in an excel worksheet. The Organellar Genome  
409 Draw (OGDRAW) was used to draw the organellar genome map of the chloroplast genome [43]. The  
410 Genbank file was used to draw the chloroplast genome map in OGDRAW [43].

## 411 **Multiple sequence alignment**

412 The intron sequences retrieved from the chloroplast tRNA were aligned to find the possible  
413 conserved structure. Multiple sequence alignment was conducted using the Multalin software  
414 (<http://multalin.toulouse.inra.fr/multalin/>) that uses hierarchical clustering [44]. Default parameters  
415 were used to construct the alignment.

## 416 **Machine Learning Approach and Statistical Analysis**

417 A machine learning approach was used to understand the role of the genome size and GC% content  
418 in the number of tRNA genes in the chloroplast genome. The random forest regression approach was  
419 used for this purpose. The following parameters were used in the random forest analysis: target  
420 tRNA gene number, predictor's genome size, and GC% content; Plots: data split, out-of-bag error,  
421 predictive performance, mean decrease in accuracy, and total increase in node purity; tables:  
422 evaluation matrix; data split preference: sample 20% of all data; training and validation of data: 20%

423 validation data. The training parameters were as follows, training data used per tree: 50%; predictor  
424 per split: auto; and max tree: 100%. The machine-learning approach was studied using the JASP  
425 software version 0.16.1.0 [45]. The correlation plot for GC% content and tRNA was also conducted  
426 using the JASP 0.16.1.0 software. The following parameters were used for the correlation analysis,  
427 sample correlation coefficient: Pearson's  $r$  and confidence interval: 95% ( $p < 0.05$ )[45].

#### 428 **Phylogenetic tree**

429 The tRNA sequences of the chloroplast genomes were taken to construct the phylogenetic tree. The  
430 phylogenetic tree was constructed using the Clustalw program in a Linux-based environment. A  
431 neighbor joining tree was constructed with 100 bootstrap replicates. The resulting file was saved in  
432 nwk file format and later uploaded in the iTOL Interactive Tree of Life, to view the tree [46]. The  
433 phylogenetic tree of the tRNA quadruplet anticodons, with other anticodons, was constructed using  
434 the MEGA software version 7[47]. Prior to the construction of the phylogenetic tree, the tRNA  
435 sequences were subjected to multiple sequence alignments. Multiple sequence alignments were  
436 conducted using the MUSCLE software[48]. The resulting clustal file was converted to the MEGA file  
437 format (aln) using the MEGA 7 software[47]. The converted file was subjected to construct the  
438 phylogenetic tree in the MEGA 7 software, using the maximum-likelihood approach. The Tamura-  
439 Nei model, with a 500-bootstrap replicate, was used for this analysis. The phylogenetic tree of the  
440 tRNA introns was also constructed using the MEGA 7 software with the same statistical parameters  
441 [47].

442

443 **Acknowledgement:** N/A

444 **Statement and Declarations**

445 **Data availability**

446 All the data used during this study was taken from National Center for Biotechnology Information  
447 database and all the data are available in the public domain. Also, the accession numbers are  
448 provided in the supplementary files.

449 **Competing interest**

450 Authors have no competing interest to declare.

451

452 **References**

- 453 1. Lehman N. Molecular evolution: Please release me, genetic code. *Curr Biol.* 2001;11:R63–  
454 6. doi:[https://doi.org/10.1016/S0960-9822\(01\)00016-1](https://doi.org/10.1016/S0960-9822(01)00016-1).
- 455 2. Keeling PJ, Doolittle WF. A non-canonical genetic code in an early diverging eukaryotic  
456 lineage. *EMBO J.* 1996;15:2285–90. doi:<https://doi.org/10.1002/j.1460-2075.1996.tb00581.x>.
- 457 3. Magliery TJ, Anderson JC, Schultz PG. Expanding the genetic code: Selection of efficient  
458 suppressors of four-base codons and identification of “shifty” four-base codons with a  
459 library approach in *Escherichia coli*. *J Mol Biol.* 2001;307:755–69.
- 460 4. Atkins JF, Weiss RB, Thompson S, Gesteland RF. Towards a genetic dissection of the  
461 basis of triplet decoding, and its natural subversion: Programmed Reading Frame Shifts  
462 and Hops. *Annu Rev Genet.* 1991;25:201–28. doi:10.1146/annurev.ge.25.120191.001221.
- 463 5. Kurland CG. Translational accuracy and the fitness of bacteria. *Annu Rev Genet.*  
464 1992;26:29–50. doi:10.1146/annurev.ge.26.120192.000333.
- 465 6. RIDDLE DL, CARBON J. Frameshift Suppression: a Nucleotide Addition in the  
466 Anticodon of a Glycine Transfer RNA. *Nat New Biol.* 1973;242:230–4.  
467 doi:10.1038/newbio242230a0.
- 468 7. Mohanta TK, Yadav D, Khan A, Hashem A, Abd\_Allah EF, Al-Harrasi A. Analysis of  
469 genomic tRNA revealed presence of novel genomic features in cyanobacterial tRNA. *Saudi*  
470 *J Biol Sci.* 2019;27:124–33. doi:<https://doi.org/10.1016/j.sjbs.2019.06.004>.

- 471 8. Yu M, Lu B, Zhang J, Ding J, Liu P, Lu Y. tRNA-derived RNA fragments in cancer:  
472 current status and future perspectives. *J Hematol Oncol.* 2020;13:121. doi:10.1186/s13045-  
473 020-00955-6.
- 474 9. Kumar P, Mudunuri SB, Anaya J, Dutta A. tRFdb: a database for transfer RNA  
475 fragments. *Nucleic Acids Res.* 2015;43 Database issue:D141–5. doi:10.1093/nar/gku1138.
- 476 10. Mohanta TK, Mishra AK, Hashem A, Abd\_Allah EF, Khan AL, Al-Harrasi A.  
477 Construction of anti-codon table of the plant kingdom and evolution of tRNA  
478 selenocysteine (tRNA<sup>Sec</sup>). *BMC Genomics.* 2020;21:804. doi:10.1186/s12864-020-07216-3.
- 479 11. Zhao Y-H, Zhou T, Wang J-X, Li Y, Fang M-F, Liu J-N, et al. Evolution and structural  
480 variations in chloroplast tRNAs in gymnosperms. *BMC Genomics.* 2021;22:750.  
481 doi:10.1186/s12864-021-08058-3.
- 482 12. Zhang Y, Baranov P V, Atkins JF, Gladyshev VN. Pyrrolysine and Selenocysteine Use  
483 Dissimilar Decoding Strategies. *J Biol Chem.* 2005;280:20740–51.  
484 doi:10.1074/jbc.M501458200.
- 485 13. Mohanta TK, Mishra AK, Khan A, Hashem A, Abd\_Allah EF, Al-Harrasi A. Gene Loss  
486 and Evolution of the Plastome. *Genes.* 2020;11:1133.
- 487 14. Yona AH, Bloom-Ackermann Z, Frumkin I, Hanson-Smith V, Charpak-Amikam Y, Feng  
488 Q, et al. tRNA genes rapidly change in evolution to meet novel translational demands.  
489 *Elife.* 2013;2013:1–17.
- 490 15. Anderson JC, Wu N, Santoro SW, Lakshman V, King DS, Schultz PG. An expanded  
491 genetic code with a functional quadruplet codon. *Proc Natl Acad Sci U S A.* 2004;101:7566  
492 LP – 7571. doi:10.1073/pnas.0401517101.
- 493 16. Neumann H, Wang K, Davis L, Garcia-Alai M, Chin JW. Encoding multiple unnatural  
494 amino acids via evolution of a quadruplet-decoding ribosome. *Nature.* 2010;464:441–4.

495 doi:10.1038/nature08817.

496 17. Niu W, Schultz PG, Guo J. An Expanded Genetic Code in Mammalian Cells with a  
497 Functional Quadruplet Codon. *ACS Chem Biol.* 2013;8:1640–5. doi:10.1021/cb4001662.

498 18. Watanabe T, Muranaka N, Hohsaka T. Four-base codon-mediated saturation  
499 mutagenesis in a cell-free translation system. *J Biosci Bioeng.* 2008;105:211–5.  
500 doi:<https://doi.org/10.1263/jbb.105.211>.

501 19. Anderson JC, Schultz PG. Adaptation of an Orthogonal Archaeal Leucyl-tRNA and  
502 Synthetase Pair for Four-base, Amber, and Opal Suppression. *Biochemistry.* 2003;42:9598–  
503 608. doi:10.1021/bi034550w.

504 20. de la Torre D, Chin JW. Reprogramming the genetic code. *Nat Rev Genet.* 2021;22:169–  
505 84. doi:10.1038/s41576-020-00307-7.

506 21. DeBenedictis EA, Carver GD, Chung CZ, Söll D, Badran AH. Multiplex suppression of  
507 four quadruplet codons via tRNA directed evolution. *Nat Commun.* 2021;12:5706.  
508 doi:10.1038/s41467-021-25948-y.

509 22. Wang K, Schmied WH, Chin JW. Reprogramming the Genetic Code: From Triplet to  
510 Quadruplet Codes. *Angew Chemie Int Ed.* 2012;51:2288–97.  
511 doi:<https://doi.org/10.1002/anie.201105016>.

512 23. Hausner G, Hafez M, Edgell DR. Bacterial group I introns: mobile RNA catalysts. *Mob*  
513 *DNA.* 2014;5:8. doi:10.1186/1759-8753-5-8.

514 24. Zhou Y, Lu C, Wu Q-J, Wang Y, Sun Z-T, Deng J-C, et al. GISSD: Group I Intron  
515 Sequence and Structure Database. *Nucleic Acids Res.* 2008;36 Database issue:D31–7.  
516 doi:10.1093/nar/gkm766.

517 25. Mohanta T, Syed A, Ameen F, Bae H. Novel Genomic and Evolutionary Perspective of  
518 Cyanobacterial tRNAs. *Front Genet.* 2017;8:200.

- 519 26. Lund E, Dahlberg JE. Spacer transfer RNAs in ribosomal RNA transcripts of *E. coli*:  
520 Processing of 30S ribosomal RNA in vitro. *Cell*. 1977;11:247–62.  
521 doi:[https://doi.org/10.1016/0092-8674\(77\)90042-3](https://doi.org/10.1016/0092-8674(77)90042-3).
- 522 27. Molla-Herman A, Angelova MT, Ginestet M, Carré C, Antoniewski C, Huynh J-R.  
523 tRNA Fragments Populations Analysis in Mutants Affecting tRNAs Processing and tRNA  
524 Methylation. *Front Genet*. 2020;11. doi:10.3389/fgene.2020.518949.
- 525 28. Goodarzi H, Liu X, Nguyen HCB, Zhang S, Fish L, Tavazoie SF. Endogenous tRNA-  
526 Derived Fragments Suppress Breast Cancer Progression via YBX1 Displacement. *Cell*.  
527 2015;161:790–802. doi:<https://doi.org/10.1016/j.cell.2015.02.053>.
- 528 29. Kuscu C, Kumar P, Kiran M, Su Z, Malik A, Dutta A. tRNA fragments (tRFs) guide Ago  
529 to regulate gene expression post-transcriptionally in a Dicer-independent manner. *Rna*.  
530 2018;24:1093–105.
- 531 30. Fu H, Feng J, Liu Q, Sun F, Tie Y, Zhu J, et al. Stress induces tRNA cleavage by  
532 angiogenin in mammalian cells. *FEBS Lett*. 2009;583:437–42.  
533 doi:<https://doi.org/10.1016/j.febslet.2008.12.043>.
- 534 31. Yamasaki S, Ivanov P, Hu G, Anderson P. Angiogenin cleaves tRNA and promotes  
535 stress-induced translational repression. *J Cell Biol*. 2009;185:35–42.
- 536 32. Baer RJ, Dubin DT. The sequence of a possible 5S RNA-equivalent in hamster  
537 mitochondria. *Nucleic Acids Res*. 1980;8:3603–10. doi:10.1093/nar/8.16.3603.
- 538 33. Dubin DT, Friend DA. Comparison of cytoplasmic and mitochondrial 4 s RNA from  
539 cultured hamster cells: Physical and metabolic properties. *J Mol Biol*. 1972;71:163–75.  
540 doi:[https://doi.org/10.1016/0022-2836\(72\)90344-0](https://doi.org/10.1016/0022-2836(72)90344-0).
- 541 34. Watanabe Y-I, Suematsu T, Ohtsuki T. Losing the stem-loop structure from metazoan  
542 mitochondrial tRNAs and co-evolution of interacting factors. *Front Genet*. 2014;5:109.



543 doi:10.3389/fgene.2014.00109.

544 35. Poul N, Morten K, Søren T, Galina P, Ludmila R, C. CBF, et al. Crystal Structure of the  
545 Ternary Complex of Phe-tRNA<sup>Phe</sup>, EF-Tu, and a GTP Analog. *Science* (80- ).  
546 1995;270:1464–72. doi:10.1126/science.270.5241.1464.

547 36. Ohtsuki T, Sato A, Watanabe Y, Watanabe K. A unique serine-specific elongation factor  
548 Tu found in nematode mitochondria. *Nat Struct Biol.* 2002;9:669–73. doi:10.1038/nsb826.

549 37. Arita M, Suematsu T, Osanai A, Inaba T, Kamiya H, Kita K, et al. An evolutionary  
550 “intermediate state” of mitochondrial translation systems found in *Trichinella* species of  
551 parasitic nematodes: co-evolution of tRNA and EF-Tu. *Nucleic Acids Res.* 2006;34:5291–9.  
552 doi:10.1093/nar/gkl526.

553 38. Ohtsuki T, Watanabe Y, Takemoto C, Kawai G, Ueda T, Kita K, et al. An “Elongated”  
554 Translation Elongation Factor Tu for Truncated tRNAs in Nematode Mitochondria. *J Biol*  
555 *Chem.* 2001;276:21571–7. doi:https://doi.org/10.1074/jbc.M011118200.

556 39. Sakurai M, Watanabe Y, Watanabe K, Ohtsuki T. A protein extension to shorten RNA:  
557 elongated elongation factor-Tu recognizes the D-arm of T-armless tRNAs in nematode  
558 mitochondria. *Biochem J.* 2006;399:249–56. doi:10.1042/BJ20060781.

559 40. Chan PP, Lin BY, Mak AJ, Lowe TM. tRNAscan-SE 2.0: improved detection and  
560 functional classification of transfer RNA genes. *Nucleic Acids Res.* 2021;49:9077–96.  
561 doi:10.1093/nar/gkab688.

562 41. Laslett D, Canback B. ARAGORN, a program to detect tRNA genes and tmRNA genes  
563 in nucleotide sequences. *Nucleic Acids Res.* 2004;32:11–6. doi:10.1093/nar/gkh152.

564 42. Tillich M, Lehwark P, Pellizzer T, Ulbricht-Jones ES, Fischer A, Bock R, et al. GeSeq -  
565 versatile and accurate annotation of organelle genomes. *Nucleic Acids Res.* 2017;45:W6–11.  
566 doi:10.1093/nar/gkx391.

- 567 43. Greiner S, Lehwark P, Bock R. OrganellarGenomeDRAW (OGDRAW) version 1.3.1:  
568 expanded toolkit for the graphical visualization of organellar genomes. *Nucleic Acids Res.*  
569 2019;;doi.org/10.1093/nar/gkz238. doi:10.1101/545509.
- 570 44. Corpet F. Multiple sequence alignment with hierarchical clustering. *Nucleic Acids Res.*  
571 1988;16:10881–90. <https://www.ncbi.nlm.nih.gov/pubmed/2849754>.
- 572 45. Team J. JASP (Version 0.16.1). 2022;;<https://jasp-stats.org/>.
- 573 46. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree  
574 display and annotation. *Nucleic Acids Res.* 2021;49:W293–6. doi:10.1093/nar/gkab301.
- 575 47. Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis  
576 Version 7.0 for Bigger Datasets. *Mol Biol Evol.* 2016;33:1870–4. doi:10.1093/molbev/msw054.
- 577 48. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high  
578 throughput. *Nucleic Acids Res.* 2004;32:1792–7. doi:10.1093/nar/gkh340.
- 579
- 580
- 581
- 582
- 583
- 584
- 585
- 586
- 587
- 588

## 589 **Figure legends**

### 590 **Figure 1**

591 OGDRAW map of *Pilostyles aethiopica* (NC\_029235.1) chloroplast genome. The map shows the loss of  
592 tRNA genes and inverted repeats.

### 593 **Figure 2**

594 OGDRAWM map of (A) *Asparagus officinalis* (NC\_034777.1) and (B) *Populus yunnanensis*  
595 (NC\_037421.1) chloroplast genomes. (A) *Asparagus officinalis* shows the presence of a putative spacer  
596 tRNAs. tRNA<sup>Ala</sup> (UGC) and tRNA<sup>Ile</sup> (GAU) are present between 16S and 23S rRNA in *A. officinalis*.  
597 No putative spacer tRNA is found in the chloroplast genome of *Populus yunnanensis*.

### 598 **Figure 3**

599 Putative novel tRNAs in chloroplast genome. (A) In the tRNA of *Pedicularis ishidoyana* (NC\_029700.1)  
600 there is a long nucleotide sequence present in between the D arm and anticodon arm. (B) In *Entransia*  
601 *fimbriata* (NC\_030313.1) tRNA (tRNA<sup>Lys</sup><sub>UUU</sub>) a long nucleotide sequence is present in the anticodon  
602 loop region that masks the anticodon loop. (C) In *Syntrichia ruralis* (NC\_012052.1) tRNA<sup>Gly</sup><sub>UCC</sub>, a long  
603 nucleotide sequence is found in between the D-arm and anticodon arm.

### 604 **Figure 4**

605 Putative novel tRNA of chloroplast tRNA. The tRNA contains a long nucleotide sequence in between  
606 the D-arm and anticodon arm. At least 42 chloroplast genomes are found to encode a similar tRNA  
607 structure in it. The structure was predicted using the tRNAscan-SE 2.0 program.

### 608 **Figure 5**

609 Figure 5 shows the presence of putative nematode mitochondrial tRNA in the chloroplast genome.  
610 The tRNAs have been seen to lose the  $\Psi$ -arm and  $\Psi$ -loop. The presence of nematode mitochondrial  
611 genome in the chloroplast genome shows that the truncated tRNAs are shared in between the  
612 chloroplast and mitochondria. The structure has been predicted using the Aragorn software.

### 613 **Figure 6**

614 Correlation regression analysis ( $r = 0.362$ ) of GC % and tRNA gene number in the chloroplast genome.  
615 Analysis showed that there was a slight positive correlation between the GC% and tRNA gene  
616 number in the chloroplast genome. The analysis was conducted at a  $p$ -value  $< 0.05$ . Correlation  
617 analysis was conducted using the JASP 0.16.1.0 version software.

#### 618 **Figure 7**

619 Machine-learning analysis of GC % content and genome size in the tRNA gene number in the  
620 chloroplast genome; the random forest approach was used to run the analysis. Analysis revealed that  
621 the GC% content had more influence toward the number of tRNA gene numbers than the genome  
622 size. In the study, from 5959 species, 3814 species were used as training sets, 954 for validation, and  
623 1191 as test sets. All the analysis was conducted at  $p < 0.05$ .

#### 624 **Figure 8**

625 Phylogenetic tree of chloroplast tRNAs. The phylogenetic tree shows two distinct major clusters  
626 named cluster I and cluster II. The phylogenetic tree shows that chloroplast tRNAs have evolved  
627 from multiple common ancestors. In cluster I anticodons GCC, CGU, CGA, UCU, CAA, and UAG,  
628 are found in more than one group, and the anticodons GCC, GCU, UUC, CAU, and GAA are  
629 found in both the clusters, showing their evolution via duplication. The phylogenetic tree has  
630 been constructed using the neighbor-joining method, using the Clustal W program.

#### 631 **Figure 9**

632 Phylogenetic tree of a putative quadruplet anticodon containing tRNAs with triplet codon-  
633 containing tRNAs. The phylogenetic grouping revealed that the quadruplet anticodons had  
634 evolved via addition of a nucleotide preceding the third nucleotide of the triplet anticodons. The  
635 evolutionary history was inferred by using the Maximum Likelihood method based on the  
636 Tamura-Nei model. The tree with the highest log likelihood (-1053.93) is shown. Initial tree(s) for  
637 the heuristic search were obtained automatically by applying the Neighbor-Join and BioNJ  
638 algorithms to a matrix of pair-wise distances, estimated by using the Maximum Composite  
639 Likelihood (MCL) approach, and then selecting the topology with the superior log likelihood  
640 value. The tree was drawn to scale, with branch lengths measured in the number of substitutions

641 per site. The analysis involved 147 nucleotide sequences. All positions with less than 95% site  
642 coverage were eliminated, that is, fewer than 5% alignment gaps, missing data, and ambiguous  
643 bases, were allowed at any position. There were a total of 26 positions in the final dataset.  
644 Evolutionary analyses were conducted using the MEGA 7.

645

## 646 **Supplementary Materials**

647 **Supplementary File 1.** Percentage of anticodons in the chloroplast genome.

648 **Supplementary File 2.** Name of the species with rare anticodons in their tRNA gene.

649 **Supplementary File 3.** Name of the species encoding UCA anticodon for tRNA selenocysteine.

650 **Supplementary File 4.** Name and accession number of the species encoding putative quadruplet  
651 anticodons in the chloroplast genome.

652 **Supplementary File 5.** Accession number of the species encoding putative duplet anticodons in the  
653 chloroplast tRNA.

654 **Supplementary File 6.** Genomic details of chloroplast genome of different species.

655 **Supplementary File 7.** List of species those do not contain the spacer tRNA.

656 **Supplementary file 8.** The list of tRNA fragments found in the chloroplast genome.

657 **Supplementary File 9.** Putative tRNAs of chloroplast genome.

658 **Supplementary Table 1.** Quadruplet anticodon/codon table. There are 256 possibilities to encode  
659 an amino acid via quadruplet anticodon/codon. It can accommodate maximum of the amino acids  
660 available in the proteome to its protein translation machinery.

## 661 **Supplementary Table 2**

662 Transition and transversion bias (MCL) of different tRNA genes of the chloroplast genome. Each  
663 entry shows the probability of substitution (r) from one base (row) to another base (column). For  
664 simplicity, the sum of r-values is made equal to 100. Rates of different transitional substitutions are  
665 shown in bold and those of transversional substitutions are shown in italics. The  
666 transition/transversion rate ratios are k1 indicates purines and k2 indicates pyrimidines. The overall  
667 transition/transversion bias is mentioned as R where  $R = [A * G * k1 + T * C * k2] / [(A + G) * (T + C)]$ . All  
668 positions with less than 95% site coverage were eliminated. That is, fewer than 5% alignment gaps,

669 missing data, and ambiguous bases were allowed at any position. The evolutionary analyses were  
670 conducted in MEGA7.

671 **Supplementary Figure 1.** Grouping of conserved type II introns found in tRNA of the chloroplast  
672 genome

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696



697 Table 1. Anticodon Table of the chloroplast genome. Study from 5959 chloroplast genomes shows several  
 698 rare anticodons.

699

Ala	AGC: 214	GGC: 203	CGC: 1	UGC: 11017		
Arg	ACG: 10927	GCG: 0	CCG: 304	UCG: 103	CCU: 12	UCU: 5729
Asn	AUU: 2	GUU: 11018				
Asp	AUC: 399	GUC: 6122				
Cys	ACA: 224	GCA: 6156				
Gln	CUG: 0	UUG: 6242				
Glu	CUC: 0	UUC: 5925				
Gly	ACC: 204	GCC: 4917	CCC: 0	UCC: 5684		
His	AUG: 405	GUG: 7148				
Ile	AAU: 22	GAU: 10695	CAUIle2: 10575	UAU: 5		
Leu	AAG: 1	GAG: 42	CAG: 6	UAG: 5745	CAA: 10686	UAA: 5546
Lys	CUU: 7	UUU: 5312				
Met	CAU Met: 11741					
	CAUfMet: 709					
Phe	AAA: 207	GAA: 5982				
Pro	AGG: 2	GGG: 817	CGG: 0	UGG: 5883		
Ser	AGA: 203	GGA: 5520	CGA: 173	UGA: 5174	ACU: 0	GCU: 5980
Thr	AGU: 1	GGU: 5703	CGU: 454	UGU: 5669		
Trp	CCA: 5955					
Tyr	AUA: 6	GUA: 5964				
Val	AAC: 399	GAC: 10687	CAC: 200	UAC: 4748		
SeC	UCA: 204					
Pyl	CUA: 197					
Sup	CUA:	UUA: 205	UCA: 2			

700

701

702

703

704

705

706

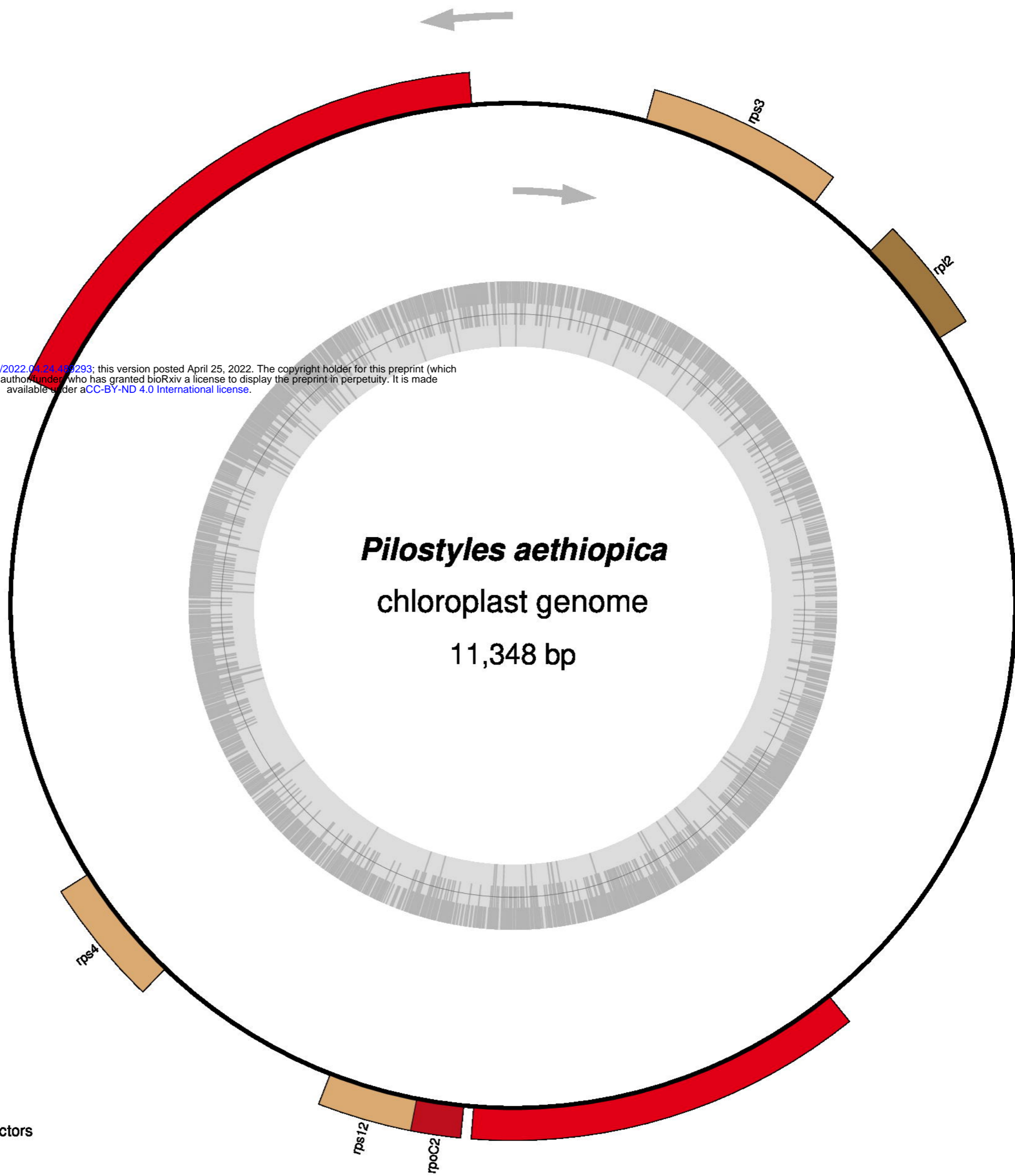
**Table 2.** Conservation of introns in chloroplast tRNAs. From the mentioned tRNA isotypes, at least eight isotypes do not encode any intron in their tRNA genes or its not conserved.

tRNA Isotype	Conserved Consensus sequence of Intron
Alanine	A-U-U-G-G-G-U-C-G-U-U-G-C-G-A-U-U-A-C-G-G-x-G-U-x-U-G-G-A-U-G-U-C-U-A-A-U-U-G
Arginine	Not found
Asparagine	Not found
Aspartic acid	Not found
Cysteine	G-C-G-C-G-C-C-A-A-U-G-U-U-U-U-U-x-C-A-G-x-G-G-A-x-G-U-x-C-A-U-C-A-U-G-x-A-A-U-C-A-A-A-x <sub>3</sub> -U-x-A-U G-x <sub>8</sub> -U-x-U-x <sub>4-5</sub> -C-A-G/A-x <sub>3</sub> -A-x <sub>2</sub> -U-C-x <sub>2</sub> -U-x <sub>5</sub> -A-A-U-x-A-x <sub>2</sub> -A-x <sub>2-3</sub> -U-U-G-A-U-C/U-x <sub>2</sub> -U-U-U-A
Glutamic acid	A-U-U-G-C-G-U-C-G-U-U-G-U-G-C-x-G-G-G-C-U-G-U-G-A-x-G-G-C-U-C-U-C-A U-x <sub>2</sub> -U/C-G-U/C-x-G-x-U-G-x-G-x <sub>8</sub> -C-U
Glutamine	Not found
Glycine	G-x-G-x-C-x <sub>3</sub> -G-C-x <sub>2</sub> -U-U-x <sub>1-5</sub> -C-x <sub>3</sub> -U-A-U-A-x <sub>2</sub> -C
Histidine	Not found
Isoleucine	A/C-G/U-U-G-C-G-x-C-A/G-U-G-U-U/G-U/G-U/C-U-x <sub>1-3</sub> -C-x-G-x <sub>3</sub> -A/G-G-U/G-x <sub>2</sub> -A/C-U-C/U-A-x <sub>2</sub> -U/G-x-C-A-x <sub>5</sub> -A/U-x <sub>4</sub> -U
Leucine	A-A-C-x <sub>5</sub> -A-A-x-U-x <sub>3</sub> -A-G-x-A-x <sub>2</sub> -A-x <sub>2</sub> -A-A
Lysine	A-G-U-G-C-G-x-C-U-x <sub>4</sub> -U-x-U-U-x-A-C-A-C-A-U-U-U-x <sub>2</sub> -A-U-G-A-A
Methionine	U-x-U-G-x-A-x <sub>2</sub> -A-G-A-G-x-U-U-U-x <sub>9-10</sub> -C-G-A-C-U-x <sub>2</sub> -A-A-U-A
Phenylalanine	C-x <sub>2</sub> -G-C-G-C-C-A-A-U-G-x <sub>1-2</sub> -U-U-x-U-C-A-x <sub>2</sub> -G-x-A-G-U-C-x-A-U-x-A-U-G-x-A-A-U-x-A-x-A-A-x-A
Proline	Not found
Serine	A-C-G-U-U-x-A-A-A-x-A-x-U-x <sub>2-7</sub> -G-U-C-G-A-A-C-C-C-C A-x <sub>3</sub> -A-x <sub>2-5</sub> -G-U-C-G-A-A-C-C-C A-x <sub>2-7</sub> -A-U-x <sub>2</sub> -A-C-x <sub>4</sub> -G-x <sub>1-2</sub> -C-x <sub>2</sub> -C C-x <sub>5</sub> -A-A-x <sub>6</sub> -A-x <sub>8</sub> -U-C-x <sub>5</sub> -C-x <sub>1-2</sub> -U-x <sub>2</sub> -A-x <sub>3</sub> -C
Threonine	A-U-U-G-C-G-U-C-G-U-U-G-U-G-C-C-U-G-G-G-C-U-G-U-G-A-G-G-G-C-U-C-U-C-A-G-C-C-A-C-A-U-G-G-A-U-A-G-U-U-C
Tryptophan	Not found
Tyrosine	G-U-U-G-G-G-U-x-U-U/C-C/U-U-x <sub>2</sub> -A-A-C-A-G-U-U-C-A-A-A-U-x-A-U-U-U-G-A-U-A-A-U-A-A-x-A-x-C-U-U-U-G-A-U-C-U-G-U-U-x-U-A G-x-U-U-U-U-x <sub>4</sub> -C-x <sub>5</sub> -A-x <sub>5</sub> -U
Valine	Not found

bioRxiv preprint doi: <https://doi.org/10.1101/2022.04.24.489293>; this version posted April 25, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

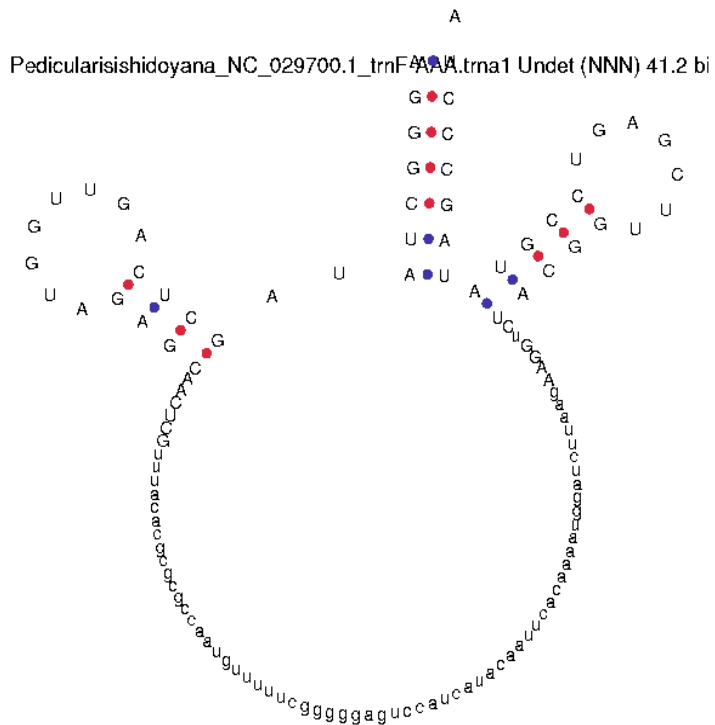
***Pilostyles aethiopica***  
chloroplast genome  
11,348 bp

- photosystem I
- photosystem II
- cytochrome b/f complex
- ATP synthase
- NADH dehydrogenase
- RubisCO large subunit
- photosystem assembly/stability factors
- RNA polymerase
- ribosomal proteins (SSU)
- ribosomal proteins (LSU)
- transfer RNAs
- ribosomal RNAs
- clpP, matK
- other genes
- hypothetical chloroplast reading frames (ycf)
- ORFs
- origin of replication
- polycistronic transcripts
- introns

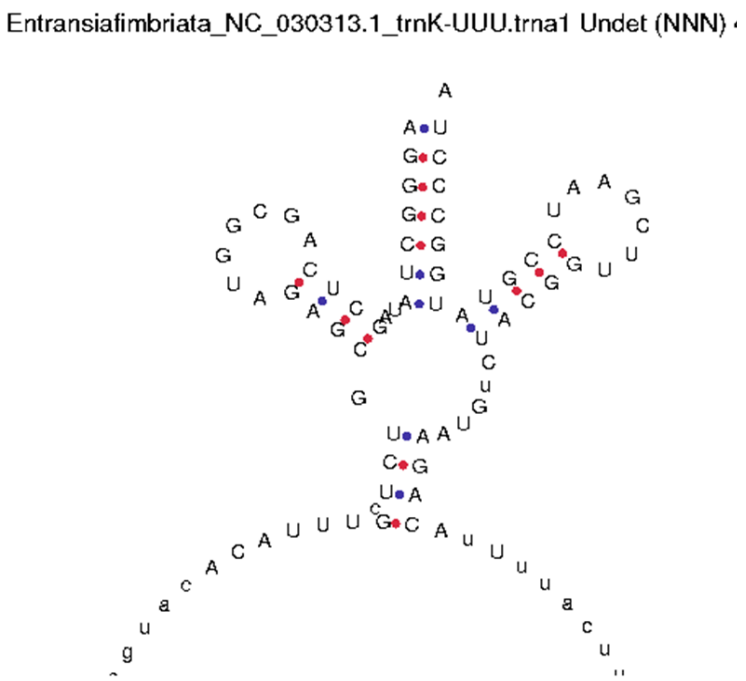




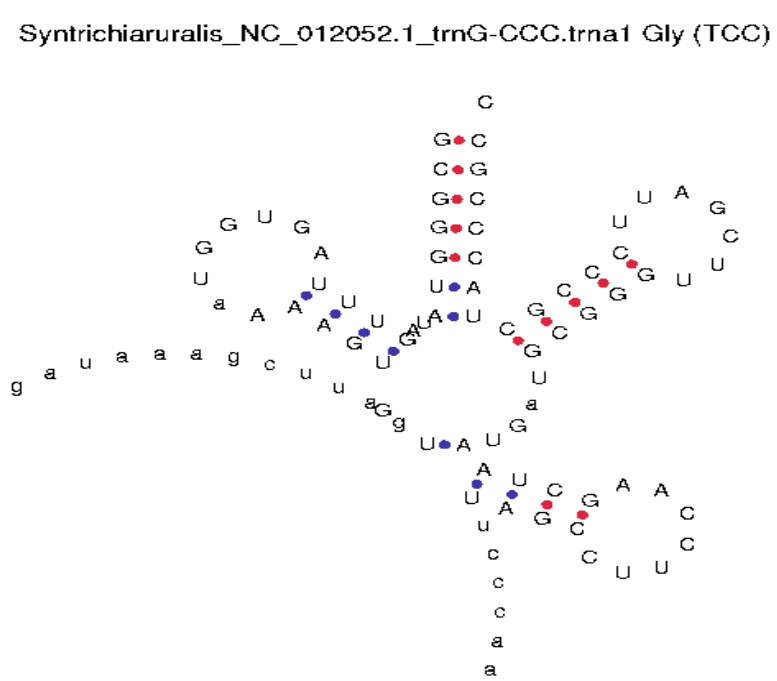
A



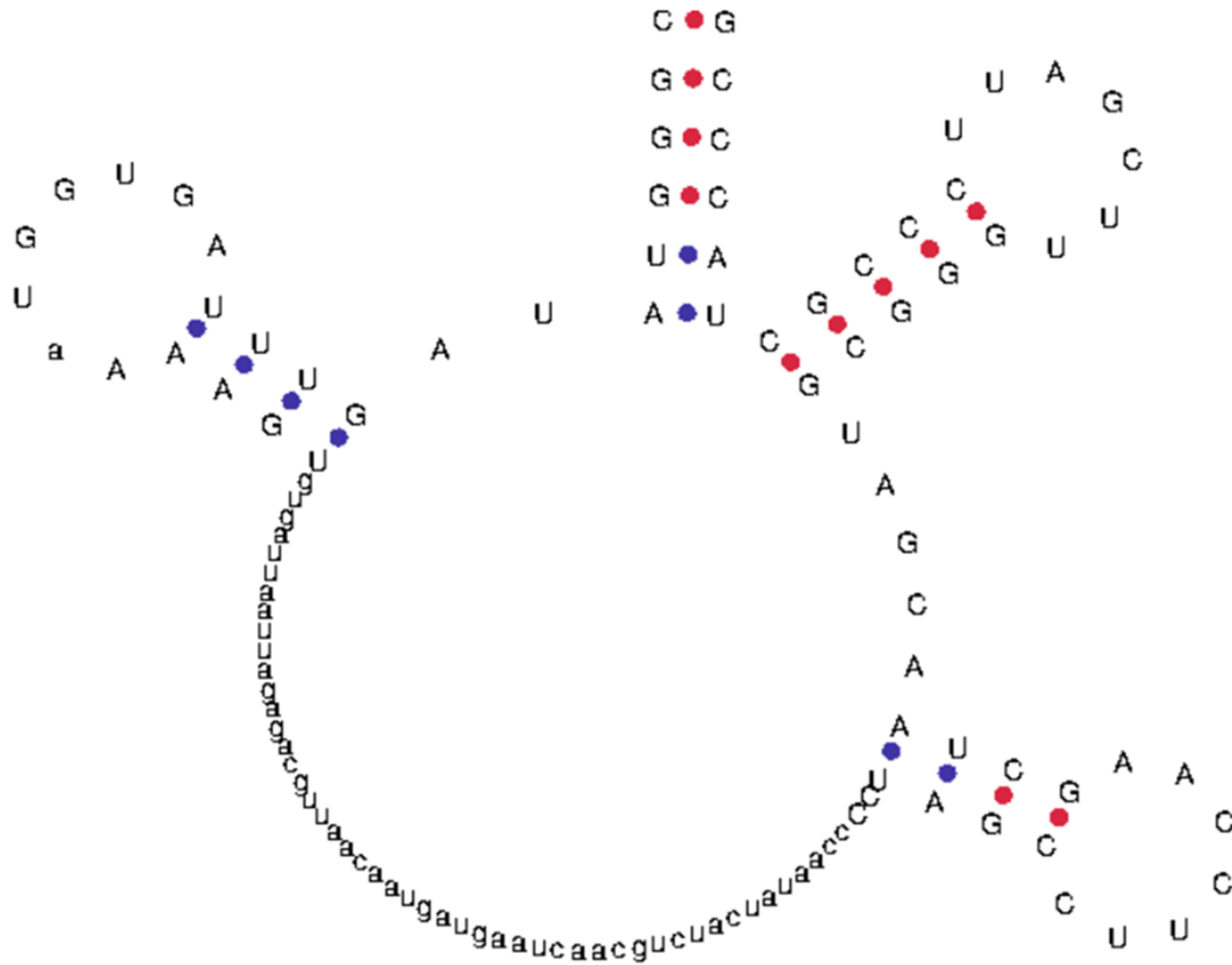
B



C







- Allomaieta villosa* (NC\_031875.1)
- Alseodaphne gracilis* (NC\_037489.1)
- Alseodaphne huanglianshanensis* (NC\_037490.1)
- Alseodaphne semecarpifolia* (NC\_037491.1)
- Arracacia xanthorrhiza* (NC\_032364.1)
- Persea americana* (NC\_031189.1)
- Peucedanum japonicum* (NC\_034644.1)
- Phoebe bournei* (NC\_034926.1)
- Phoebe chekiangensis* (NC\_034925.1)
- Phoebe omeiensis* (NC\_031190.1)
- Phoebe sheareri* (NC\_031191.1)
- Phoebe zhennan* (NC\_036143.1)
- Pterogastra divaricata* (NC\_031885.1)
- Rhexia virginica* (NC\_031886.1)
- Rhynchanthera bracteata* (NC\_031887.1)
- Tibouchina longifolia* (NC\_031889.1)
- Tigridiopalma magnifica* (NC\_036021.1)
- Triolena amazonica* (NC\_031890.1)
- Bertolonia acuminata* (NC\_031876.1)
- Cassytha filiformis* (NC\_036001.1)
- Cinnamomum camphora* (NC\_035882.1)
- Cinnamomum micranthum* (NC\_035802.1)
- Cinnamomum verum* (NC\_035236.1)
- Citrus aurantiifolia* (NC\_024929.1)
- Citrus depressa* (NC\_031894.1)
- Citrus platymamma* (NC\_030194.1)
- Citrus sinensis* (NC\_008334.1)
- Codonopsis minima* (NC\_036311.1)
- Cryptocarya chinensis* (NC\_036002.1)
- Dacrycarpus imbricatus* (NC\_034942.1)
- Floydiella terrestris* (NC\_014346.1)
- Graffenrieda moritziana* (NC\_031879.1)
- Lathyrus sativus* (NC\_014063.1)
- Ledebouriella seseloides* (NC\_034643.1)
- Lindera glauca* (NC\_035953.1)
- Machilus balansae* (NC\_028074.1)
- Machilus thunbergii* (NC\_035319.1)
- Machilus yunnanensis* (NC\_028073.1)
- Melastoma candidum* (NC\_034716.1)
- Merianthera pulchra* (NC\_031881.1)
- Miconia dodecandra* (NC\_031882.1)
- Nepsera aquatica* (NC\_031883.1)

```

      g
    g-c
    g-c
    g.a
    g+t
    a.g
    g-c
    a-ttat
  t       c
tga  g   g
g   gtcg   a
g   :!!!   g
a   aagc   c
cta  g   a
      g+ttgt
      c-g
      g-c
      g-c
      a-t
      t  a
      t  a
      gct

```

NC\_045043.1\_gene\_10\_trnS-GCT

```

      c
    g-c
    g-c
    g-c
    g-c
    a.g
    t-a
    t+g
    g-cttg
  t       g
taa  a   g
t   cttg   c
g   !!!+!   g
g   gagc   t
tca  a   c
      c-gaag
      c-g
      g-c
      c-g
      c-g
      c  a
      t  a
      gtc

```

NC\_045043.1\_gene\_24\_trnD-GTC

```

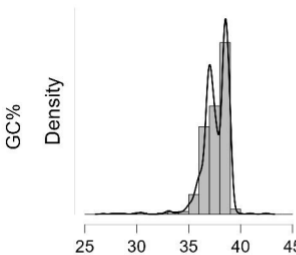
      g
    g-c
    g-c
    a a
    g+t
    a-t
    g+t
    a.gtt
  t       t
tga  g   t
g   gccg   c
g   :!!!   a
t   aggc   g
tca  g   a
      t-atgt
      a-t
      g-c
      c-g
      a-t
      t  c
      t  a
      gga

```

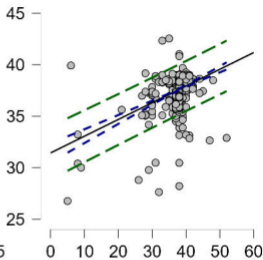
NC\_045043.1\_gene\_38\_trnS-GGA



GC%

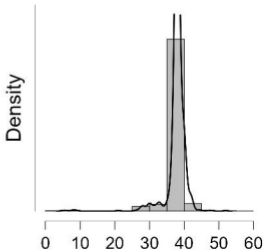


tRNA

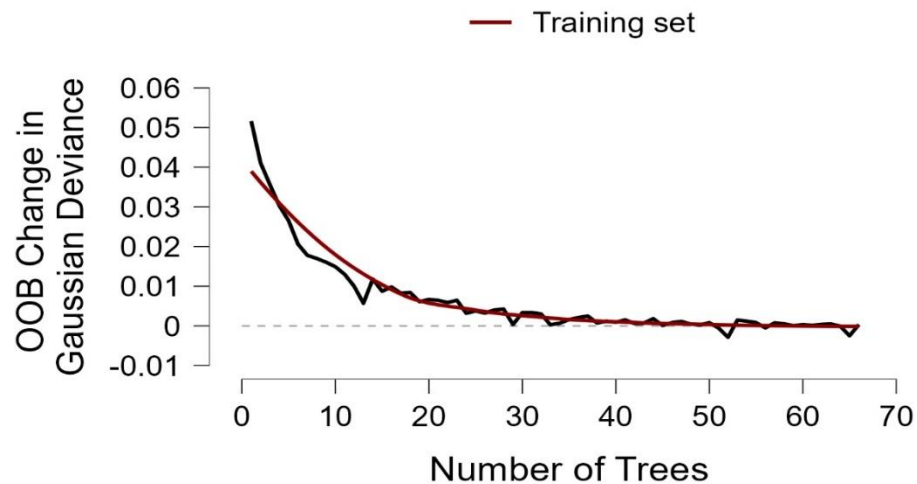


tRNA

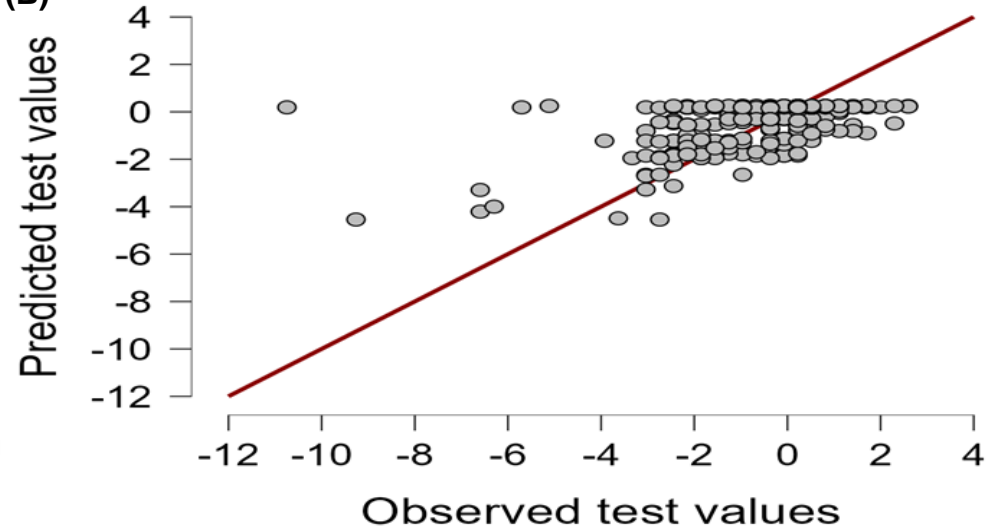
$r = 0.362$   
95% CI: [0.197, 0.504]



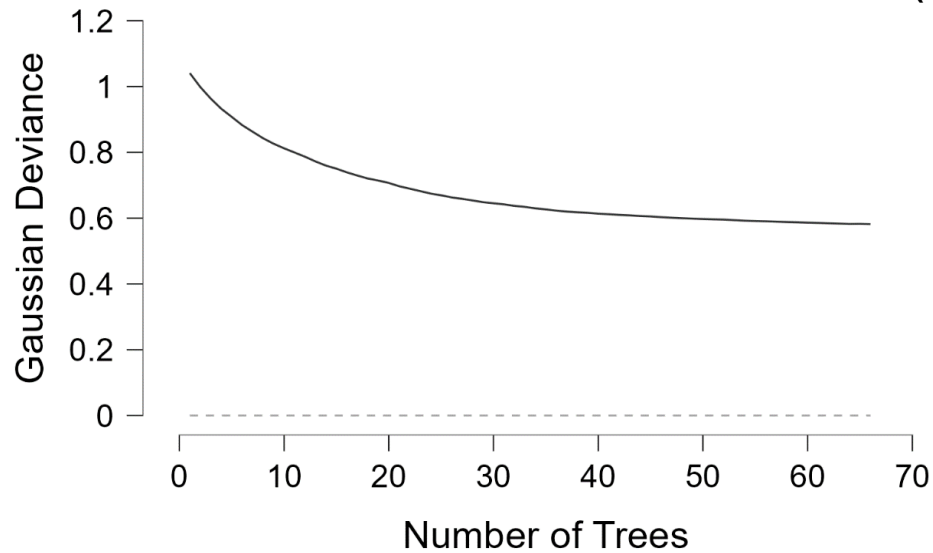
**(A)** Out-of-bag Improvement Plot



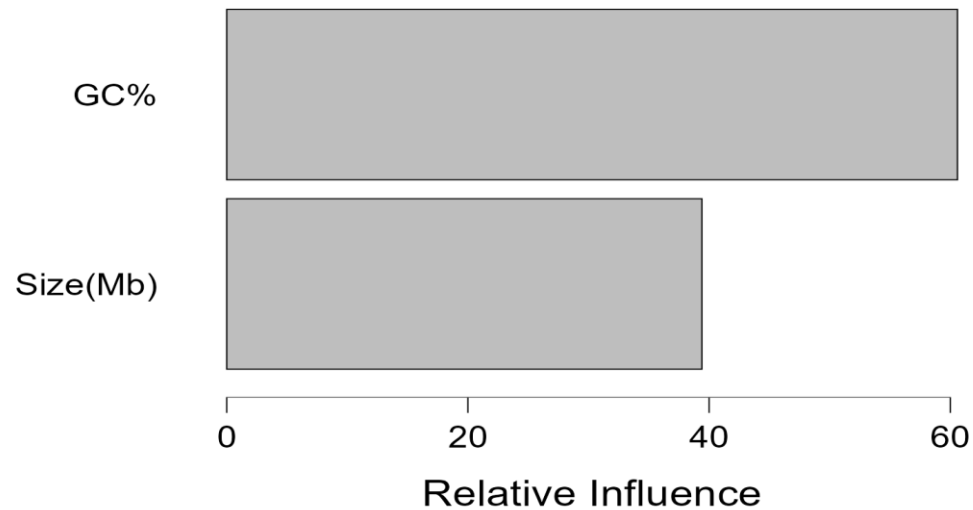
**(B)** Predictive Performance Plot



**(C)** Deviance Plot



**(D)** Relative Influence Plot

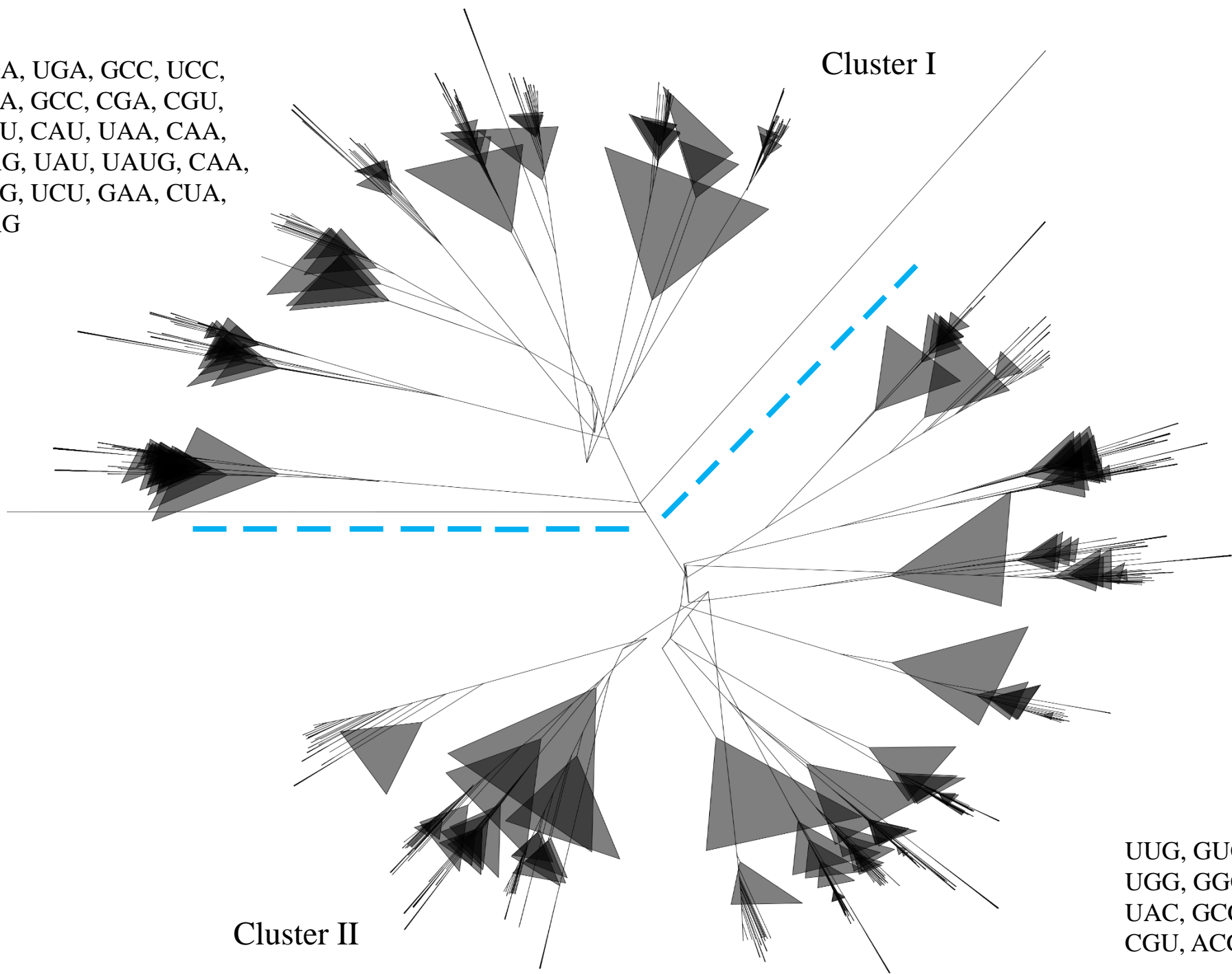


**(E)** Data Split



GCU, GGA, UGA, GCC, UCC,  
CGU, CGA, GCC, CGA, CGU,  
UUC, UCU, CAU, UAA, CAA,  
GUA, UAG, UAU, UAUG, CAA,  
GCU, UCG, UCU, GAA, CUA,  
UAG, GAG

Cluster I



Cluster II

UUG, GUG, GCA, GAA, UUU, GUU,  
UGG, GGG, CCA, UGU, GGU, CAU,  
UAC, GCC, GUC, GAC, GAU, UUC,  
CGU, ACG, CCG, ACA, UGC

