

1 **Short Title: Comparison of *Nicotiana* plastid genomes**

2 **Plastid genomics of *Nicotiana* (Solanaceae): insights into molecular evolution, positive**  
3 **selection and the origin of the maternal genome of Aztec tobacco (*Nicotiana rustica*)**

4 Furrukh Mehmood<sup>1</sup>, Abdullah<sup>1</sup>, Zartasha Ubaid<sup>1</sup>, Iram Shahzadi<sup>1</sup>, Ibrar Ahmed<sup>2</sup>, Mohammad  
5 Tahir Waheed<sup>1</sup>, Péter Poczai\*<sup>3</sup>, Bushra Mirza\*<sup>1</sup>

6 <sup>1</sup>Department of Biochemistry, Quaid-i-Azam University, Islamabad, Pakistan

7 <sup>2</sup>Alpha Genomics Private Limited, Islamabad, Pakistan

8 <sup>3</sup>Finnish Museum of Natural History (Botany Unit), University of Helsinki, Helsinki, Finland

9 \*Corresponding authors: **Bushra Mirza (bushramirza@qau.edu.pk)**

10 **Péter Poczai (peter.poczai@helsinki.fi)**

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

## 27 Abstract

28 The genus *Nicotiana* of the family Solanaceae, commonly referred to as tobacco plants, are a  
29 group cultivated as garden ornamentals. Besides their use in the worldwide production of  
30 tobacco leaves, they are also used as evolutionary model systems due to their complex  
31 development history, which is tangled by polyploidy and hybridization. Here, we assembled the  
32 plastid genomes of five tobacco species, namely *N. knightiana*, *N. rustica*, *N. paniculata*, *N.*  
33 *obtusifolia* and *N. glauca*. *De novo* assembled tobacco plastid genomes showed typical  
34 quadripartite structure, consisting of a pair of inverted repeats (IR) regions (25,323–25,369 bp  
35 each) separated by a large single copy (LSC) region (86,510–86,716 bp) and a small single copy  
36 (SSC) region (18,441–18,555 bp). Comparative analyses of *Nicotiana* plastid genomes showed  
37 similar GC content, gene content, codon usage, simple sequence repeats, oligonucleotide repeats,  
38 RNA editing sites and substitutions with currently available Solanaceae genomes sequences. We  
39 identified twenty highly polymorphic regions mostly belonging to intergenic spacer regions  
40 (IGS), which could be appropriate for the development of robust and cost-effective markers to  
41 infer the phylogeny of genus *Nicotiana* and family Solanaceae. Our comparative plastid genome  
42 analysis revealed that the maternal parent of the tetraploid *N. rustica* was the common ancestor  
43 of *N. paniculata* and *N. knightiana*, and the later species is more closely related to *N. rustica*.  
44 The relaxed molecular clock analyses estimated that the speciation event between *N. rustica* and  
45 *N. knightiana* appeared 0.56 Ma (HPD 0.65–0.46). The biogeographical analysis showed a  
46 south-to-north range expansion and diversification for *N. rustica* and related species, where *N.*  
47 *undulata* and *N. paniculata* evolved in North/Central Peru, while *N. rustica* developed in  
48 Southern Peru and separated from *N. knightiana*, which adapted to the Southern coastal climatic  
49 regimes. We further inspected selective pressure on protein-coding genes among tobacco species  
50 to determine if this adaptation process affected the evolution of plastid genes. These analyses  
51 indicated that four genes involved in different plastid functions, such as DNA replication (*rpoA*)  
52 and photosynthesis (*atpB*, *ndhD* and *ndhF*), came under positive selective pressure as a result of  
53 specific environmental conditions. Genetic mutations of the following genes might have  
54 contributed to the survival and better adaptation during the evolutionary history of tobacco  
55 species.

56 **Key words:** *Nicotiana*, Chloroplast genome, Substitution and InDels, Mutational hotspots,  
57 substitutions, positive selection.

58

## 59 **1. Introduction**

60 The plant family Solanaceae consists of 98 genera and ~ 2700 species (Olmstead et al., 2008;  
61 Olmstead & Bohs, 2007). This megadiverse family consists of herbaceous annual species to  
62 perennial trees with a natural distribution ranging from deserts to rainforests (Knapp et al.,  
63 2004). *Nicotiana* L. is the fifth largest genus in the family, comprising 76 species, which were  
64 subdivided into three subgenera and fourteen sections by Goodspeed (1954). The subgenera of  
65 *Nicotiana*, as proposed by Goodspeed (1954), were not monophyletic (Aoki & Ito, 2000; Chase  
66 et al., 2003), but most of Goodspeed's sections were natural groups. The formal classification of  
67 the genus has been refined to reflect the growing body of evidence on *Nicotiana*, consisting of  
68 thirteen sections (Knapp, Chase & Clarkson, 2004). Most *Nicotiana* species are diploid ( $2n = 2x$   
69  $= 24$ ) while allopolyploid species are also reported (Leitch et al., 2008). Phylogenetic studies  
70 have shown that these allopolyploids were formed 0.2 million (*N. rustica* L. and *N. tabacum* L.)  
71 to more than 10 million years ago (species of sect. *Suaveolentes*) (Clarkson et al., 2004; Leitch et  
72 al., 2008). Cultivated tobacco (*N. tabacum* L.), commonly grown for its leaves and an important  
73 economic and agricultural crop around the world (Occhialini et al., 2016), is a natural  
74 amphiploidy derived from two progenitors (Smith 1974). *Nicotiana* species, especially *N.*  
75 *tabacum*, are also used as model organisms in plant sciences and genetics (Zhang et al., 2011).  
76 The first complete chloroplast genome sequence was also published for this species (Shinozaki et  
77 al., 1986). Since the publication of this sequence, the structure and composition of chloroplast  
78 genomes has become widely utilized in identifying unique genetic changes and the evolutionary  
79 relationships of various groups of plants, while plastid genes have also been linked with  
80 important crop traits such as yield and resistance to various pest and pathogens (Jin & Daniell,  
81 2015). Chloroplasts (cp) are large double membrane organelles with a genome size of 75-250 kb  
82 (Palmer, 1985). Proteins are used not only for photosynthesis but also for the synthesis of fatty  
83 acids and amino acids (Cooper, 2000). Angiosperm plastomes commonly contain ~130 genes  
84 with up to 80 protein-coding, 30 transfer RNA (tRNA), and four ribosomal RNA (rRNA) genes  
85 (Daniell et al., 2016). The plastid genome exists in circular and linear forms (Oldenburg &  
86 Bendich, 2015) and the percentage of each form varies within plant cells (Oldenburg & Bendich,  
87 2016). Circular formed plastomes have a typical quadripartite structure, with two inverted repeat  
88 regions (IRa and IRb), separated by one large single-copy (LSC) and one small single-copy

89 (SSC) region (Palmer, 1985; Amiryousefi, Hyvönen & Poczai, 2018a; Abdullah et al., 2019b).  
90 Numerous mutational events occur in plastid genomes: variations in tandem repeats, insertion  
91 and deletions (indels), and point mutations, but inversions and translocations are also common  
92 (Jheng et al., 2012; Xu et al., 2015; Abdullah et al., 2019a). The plastid genome of angiosperms  
93 have a uniparental maternal inheritance (Daniell, 2007) but paternal inheritance has been  
94 recorded in a few gymnosperm species (Neale & Sederoff, 1989). The conserved organization of  
95 the plastid genome makes it extremely useful in exploring the phylogenetic relationships at  
96 various taxonomic levels (Ravi et al., 2008). Polymorphism in the chloroplast genome has been  
97 exploited to solve taxonomic issues, infer phylogeny and to investigate species adaptation to  
98 their natural habitats (Daniell et al., 2016). Genes in the plastid genome encode proteins and  
99 several types of RNA molecules, which play a vital role in functional plant metabolism, and can  
100 consequently undergo selective pressures. Most plastid protein-coding genes are under purifying  
101 selection to maintain their function, while positive selection might act on some genes in response  
102 to environmental changes. Complete plastid genome sequences are also useful tools in  
103 population genetics (Ahmad, 2014), species barcoding (Nguyen et al., 2017) transplastomic  
104 (Waheed et al., 2011, 2015) and conservation of endangered species (Wambugu et al., 2015).

105 Here, we assembled the plastid genome of five *Nicotiana* species and compared their sequences  
106 to gain insight into the chloroplast genome structure of the genus *Nicotiana*. We also inferred the  
107 phylogenetic relationship of genus *Nicotiana* and investigated the selection pressures acting on  
108 protein-coding genes, then identified mutational hotspots in the *Nicotiana* plastid that might be  
109 used for the development of robust and cost-effective markers in crop breeding or taxonomy.

## 110 **2. Materials and Methods**

### 111 **2.1. Chloroplast genomes assembly and annotation**

112 Illumina sequence data of *Nicotiana knightiana* L. (13.1 Gb, accession number SRR8169719),  
113 *N. rustica* (15.5 Gb, SRR8173839), *N. paniculata* (35.1 Gb, SRR8173256), *N. obtusifolia* (23  
114 Gb, SRR3592445) and *N. glauca* (12.5 Gb, SRR6320052) were downloaded from the Sequence  
115 Read Archive (SRA). The chloroplast genome sequence contigs were selected by performing the  
116 BWA alignment with default settings (Li & Durbin, 2009) using *Nicotiana tabacum* (GenBank  
117 accession number: NC\_001879) as a reference. Geneious R8.1 *de novo* assembler (Kearse et al.,  
118 2012) was used to order the selected contigs for final assembly. The genome sequence was

119 annotated using GeSeq (Tillich et al., 2017) and CPGAVAS2 (Shi et al., 2019). Following *de*  
120 *novo* annotation, start/stop codons and the position of introns were manually inspected and  
121 curated. The tRNA genes were verified by tRNAscan-SE version 2.0 with default settings (Lowe  
122 & Chan, 2016) and Aragorn version 1.2.38 (Laslett & Canback, 2004). Circular genome maps  
123 were drawn with OGDRAW v1.3.1 (Greiner, Lehwerk & Bock, 2019). The average coverage  
124 depth of *Nicotiana* species plastid genomes was determined by mapping all reads to *de novo*  
125 assembled plastid genomes with BWA (Li & Durbin, 2009) visualized with Tablet (Milne et al.,  
126 2009). Novel *Nicotiana* plastid genomes were deposited in NCBI and the assigned accession  
127 numbers are shown in Table 1.

## 128 **2.2. Comparative genome analysis and RNA editing prediction**

129 Novel plastid genome sequences were compared through multiple alignments using MAFFT v7  
130 (Kato & Standley, 2013). Every part of the genome, such as intergenic spacer regions (IGS),  
131 introns, protein-coding genes, and ribosomal RNAs and tRNAs, was considered for comparison.  
132 Each part was extracted and used to determine nucleotide diversity in DnaSP v6 (Rozas et al.,  
133 2017). Substitution, transition and transversion rates were also calculated compared to the *N.*  
134 *tabacum* reference using Geneious R8.1 (Kearse et al., 2012). Structural units of the plastid  
135 genome (LSC, SSC and IR) were individually aligned to determine the rate of substitutions and  
136 to further search for indels using DnaSP v6. The expansion and contraction of inverted repeats  
137 and their border positions were compared for ten selected *Nicotiana* species using IRscope  
138 (Amiryousefi, Hyvönen and Poczai, 2018b). The online software PREP-cp (Putative RNA  
139 Editing Predictor of Chloroplast) was used with default settings to determine putative RNA  
140 editing sites (Mower, 2009) and the codon usage and amino acids frequencies were determined  
141 by Geneious R8.1 software (Kearse et al., 2012).

## 142 **2.3. Repeats analyses**

143 Microsatellites repeats within the plastid genomes of five *Nicotiana* species were detected using  
144 MISA (Beier et al., 2017) with the minimal repeat number of 7 for mononucleotide repeats, 4 for  
145 di- and 3 for tri-, tetra-, penta- and hexanucleotide SSRs. We also used REPuter software (Kurtz,  
146 2002) with the following parameters: minimal repeats size was set to 30 bp, Hamming distance  
147 to 3, minimum similarity percentage of two repeats copies up to 90%, maximum computed  
148 repeats numbers to 500 bp for scanning and visualizing forward (F), reverse (R), palindromic (P)

149 and complementary (C) repeats. Tandem repeats were found with the tandem repeats finder  
150 using default parameters (Benson, 1999).

#### 151 **2.4. Synonymous ( $K_s$ ) and non-synonymous ( $K_a$ ) substitution rate analysis**

152 The synonymous ( $K_s$ ) and non-synonymous ( $K_a$ ) substitution were analyzed using the chloroplast  
153 genome of *Nicotiana tabacum* as reference for all *de novo* assembled *Nicotiana* plastid genomes.  
154 For this purpose, protein-coding genes were extracted from *Nicotiana* plastomes, then aligned  
155 with the corresponding genes of *S. dulcamara* as a reference using MAFFT (Kato & Standley,  
156 2013) and analyzed using DnaSP software (Rozas et al., 2017). We further assessed the impact  
157 of positive selection using additional codon models to estimate the rates of synonymous and  
158 nonsynonymous substitution. The signs of positive selection were further assessed using fast  
159 unconstrained Bayesian approximation (FUBAR) (Murrell et al., 2013) and the mixed effects  
160 model of evolution (MEME) (Murrell et al., 2012) as implemented in the DATAMONKEY web  
161 server (Delport et al., 2010). Sites with cut-off values of PP < 0.9 in FUBAR were considered as  
162 candidates to have evolved under positive selection. Out of all analyses performed in  
163 DATAMONKEY, the most suited model of evolution for each data set, directly estimated on this  
164 web server, was used. In addition, the mixed effects model of evolution (MEME), a branch-site  
165 method incorporated in the DATAMONKEY server, was used to test for both pervasive and  
166 episodic diversifying selection. MEME applies models variable  $\omega$  across lineages at individual  
167 sites, restricting  $\omega$  to be  $\leq 1$  in a proportion p of branches and unrestricted at a proportion (1 - p)  
168 of branches per site. Positive selection was inferred with this method for a P value < 0.05.

#### 169 **2.5. Phylogenomic analyses**

170 Plastid genome sequences from the genus *Nicotiana* were selected from Organelle Genome  
171 Resources of NCBI, accessed on 21.2.2019. We included all available plastid genome sequences  
172 of tobacco species in our analysis and added *de novo* assembled sequences while *S. dulcamara*  
173 was used as an outgroup. For the species included in our analysis, coding alignments were  
174 constructed from the excised plastid genes using MACSE (Ranwez et al., 2011), including the  
175 following seventy-five protein coding genes: *atpA*, B, E, F, H, I; *ccsA*; *cemA*; *clpP*; *matK*;  
176 *ndhA*, B, C, D, E, F, G, H, I, J, K; *petA*, B, D, G, J, L, N; *psaA*, C, I, J; *psbA*, B, C, D, E, F, I, L,  
177 M, N, T, Z; *rbcL*; *rpl2*, 7, 14, 16, 19, 20, 22, 23, 32, 33, 36; *rpoA*, B, C1; *rps2*, 3, 4, 7, 8, 11, 14,  
178 15, 16, 18, 19; *ycf2*, 3, 4. For phylogenetic analysis we used a matrix of protein-coding genes of



179 twelve species with a concatenated matrix length of 75,449 bp. The best fitting model  
180 (GY+F+I+G4) was determined by ModelFinder (Kalyaanamoorthy et al., 2017) as implemented  
181 in IQ-TREE according to the Akaike information criterion (AIC), and Bayesian information  
182 criterion (BIC). Maximum likelihood (ML) analyses were performed with IQ-TREE (Nguyen et  
183 al., 2015) using the ultrafast bootstrap approximation (UFBoot; Hoang et al., 2018) with 1,000  
184 replicates and the SH-like approximate likelihood ratio test (SH-aLRT), also with 1,000  
185 bootstrap replicates, and TreeDyn was used for further enhancement of phylogenetic tree  
186 analysis (Dereeper et al., 2008; Lemoine et al., 2019).

187 Relative divergence times were estimated for the species *N. rustica* and putative parental species  
188 using BEAST v.1.8.4 (Drummond et al., 2012), applying GTR + I + G rate substitution to the  
189 protein-coding plastid gene matrix. A Yule speciation tree prior and a relaxed uncorrelated  
190 clock- model that allows rates to vary independently along branches (Drummond et al., 2006)  
191 were used, with all other parameters set to default. The median time split between the *S.*  
192 *dulcamara* and *N. undulata* (mean = 25 Myr; standard deviation = 0.5) was used as a temporal  
193 constraint to calibrate the BEAST analyses derived from the Solanaceae-wide phylogeny of  
194 Särkinen et al. (2013) and the Time Tree of Life (Kumar et al. 2017). Uncertainty regarding  
195 these dates was incorporated by assigning normal prior distributions to the two calibration points  
196 (Couvreur et al., 2008; Evans et al., 2014). Four independent BEAST runs were conducted, each  
197 with 10 million generations, sampling every 10,000 generations. Convergence of all parameters  
198 was assessed in Tracer 1.5 (Rambaut et al., 2014) and 10% of each chain was removed as burn-  
199 in. The Markov chains were combined in LogCombiner 1.7.2. (Drummond et al., 2012) to  
200 calculate the maximum clade credibility tree.

201 We defined six biogeographical areas based on Köppen-Geiger climatic and further  
202 biogeographic evidence and distributions: (A) Colombian/Ecuadorian mountain range mixed  
203 equatorial (Af), monsoon (Am) and temperate oceanic climate (Cfb), (B) Northern Peruvian  
204 mountain range with tropical savanna climate (Aw), (C) Central Peru with equatorial climate  
205 (Af), (D) Coastal Peru with cold semi-arid and desert climate (*Bsk*, *BWk*), (E) Peruvian  
206 Mountain range with humid subtropical/oceanic highland climate (*Cwb*), (F) Bolivian/Chilean  
207 alpine/mountain range with mixed semi-arid cold (*Bsk*, *BWk*) and humid subtropical climate  
208 (*Cwa*). These areas were used in the Bayesian Binary Method (BBM) model implemented in  
209 RASP (Yu, Blair & He, 2019) to investigate the biogeographic history of the selected four

210 *Nicotiana* species. BBM infers ancestral area using a full hierarchical Bayesian approach and  
211 hypothesizes a special “null distribution”, meaning that an ancestral range contains none of the  
212 unit areas (Ronquist 2004). The analysis was performed on the BEAST maximum clade  
213 credibility tree using default settings, i.e. fixed JC + G (Jukes-Cantor + Gamma) with null root  
214 distribution. Ancestral area reconstruction for each node was manually plotted on the BEAST  
215 tree using pie charts. Species distributions were determined from data stored in the Solanaceae  
216 Source Database (<http://solanaceaesource.org/>) and Global Biodiversity Information Facility  
217 (GBIF) (<https://www.gbif.org/>).

### 218 **3. RESULTS**

#### 219 **3.1. Characteristics of *Nicotiana* plastid genomes**

220 Five *Nicotiana* species chloroplast genomes were assembled and the lengths of these plastid  
221 genomes were: *Nicotiana knightiana* (155,968 bp), *Nicotiana rustica* (155,849 bp), *Nicotiana*  
222 *paniculata* (155,689 bp), *Nicotiana obtusifolia* (156,022 bp) and *Nicotiana glauca* (155,917 bp).  
223 Further details of the characteristics of the assembled plastid genomes are summarized in Table  
224 S1. The coverage of assembled plastid genomes was 811× *Nicotiana knightiana*, 1,951×  
225 *Nicotiana rustica*, 1,032× *Nicotiana paniculata*, 1,412× *Nicotiana obtusifolia* and 327×  
226 *Nicotiana glauca*. The GC content of IR regions were highest (43.2%) followed by LSC (35.9%)  
227 and SSC (32.1%) (Table S1). The high GC content of IR was due to high GC content of the  
228 tRNAs (52.9%) and rRNAs (55.4%) genes.

229 *De novo* assembled *Nicotiana* plastid genomes had 134 unique genes, whereas eighteen genes  
230 were duplicated in the IR region (Table S2, Fig.1). Out of 134 genes, 86 were protein-coding  
231 genes, 37 were tRNA genes and 8 were rRNA genes. Among 18 duplicated genes in IR region, 7  
232 were protein-coding, 7 were tRNA genes and 4 were rRNA genes. 18 intron-containing genes  
233 were present in the plastome of *Nicotiana* species. The *rps12* gene is a trans-spliced gene, its 1<sup>st</sup>  
234 exon existing in the LSC region while the 2<sup>nd</sup> and 3<sup>rd</sup> exons are in the IR region.

#### 235 **3.2. Comparative analyses, codon usage and RNA editing sites**

236 The nucleotide composition of *Nicotiana* species was compared, and all genomes had similar  
237 nucleotide composition indicating high synteny in the LSC, SSC, IR and CDSs but also in non-  
238 coding regions. Detailed comparison of the base composition is shown in Table S3. A high  
239 percentage of hydrophobic amino acids were encoded in *Nicotiana* plastid genomes, while acidic



240 amino acids were present at lower rates. The amino acids are AT rich sequences as compared to  
241 GC (Fig. 2A). Relative synonyms codon usage (RSCU) and frequency of amino acid revealed  
242 that leucine is the most abundant and cysteine was the least encoded amino acid in these  
243 genomes (Fig S1). The codon usage revealed a high frequency of codons with A/T at 3<sup>rd</sup> codon  
244 position as compared to C/G at 3<sup>rd</sup> codon position (Table S4).

245 The number of predicted RNA editing sites using PREP-cp varied between 34 and 37, distributed  
246 among fifteen genes (see Table S3). Among these genes, *ndhB* (9) possessed the most of these  
247 sites, followed by *ndhD* (6-8) and *rpoB* (4). The *ndhD* gene revealed a fraction of variation  
248 among species: *N. knightiana*, *N. rustica* and *N. paniculata* having six RNA editing sites whereas  
249 seven were observed in *N. obtusifolia* and eight in *N. glauca*. Most of the RNA editing sites were  
250 C to U edits on the first and second base of the codons, but the frequency of second base codon  
251 edits was much higher. The conversions from serine to leucine were the most frequent and these  
252 changes helped in the formation of hydrophobic amino acids, i.e. valine, leucine and  
253 phenylalanine (Table S5).

### 254 **3.3. IR contraction and expansion**

255 The LSC/IR and IR/SSC border positions of *Nicotiana* plastid genomes were compared (Fig 3)  
256 using IRscope. The length of the IR regions was similar, ranging from 25,331bp to 25,436bp  
257 showing some expansion. The endpoint of the Solanaceae JLA (IRa/SSC) is characteristically  
258 located upstream of the *rps19* and downstream of the *trnH-GUG*, which was confirmed in  
259 *Nicotiana*. In *N. tomentosiformis*, the IR expanded to partially include *rps19*, creating a truncated  
260  $\psi$ *rps19* copy at JLA, which was thought to be missing from the entire *Nicotiana* clade  
261 (Amiryousefi, Hyvönen & Poczai, 2018a). In this species the IR region has expanded to include  
262 60 bp of *rps19*. The extent of the IR expansion to *rps19* varied from 2 to 60 bp and the end point  
263 seems to be conserved to the following intergenic spacer region. Furthermore, *infA*, *ycf15*, and a  
264 copy of *ycf1* located on the JSB were detected as pseudogenes. The position of *ycf1* in the  
265 IRb/SSC region varied. It left a 36 bp pseudogene in *N. knightiana*, *N. rustica* and *N. glauca*, 33  
266 bp pseudogene in *N. obtusifolia* and a 72 bp one in *N. paniculata*.

### 267 **3.4. Non-synonymous ( $K_a$ ) and synonymous ( $K_s$ ) substitution rate analysis**

268 Synonymous/non-synonymous substitutions ratio is widely used as an indicator of adaptive  
269 evolution or positive selection (Kimura, 1979). We have calculated the  $K_s$ ,  $K_a$  and  $K_a/K_s$  ratio for

270 77 protein-coding genes for five selected *Nicotiana* species using *S. dulcamara* as a reference  
271 (Table S6). Among the analyzed genes, 31 had  $K_s=0$ , 19 had  $K_a=0$ , and 39 genes had both  $K_s$  and  
272  $K_a=0$  values. Of the investigated genes, 21 showed a  $K_a/K_s$  ratio of more than 0.5. Eight of these  
273 genes (*atpF*, *psaA*, *ycf4*, *psbB*, *infA*, *ndhB*, *rpl32* and *ccsA*) had a  $K_a/K_s$  ratio greater than 0.5 for  
274 one species, *ycf1* had  $K_a/K_s$  greater than 0.5 for two species, while *atpA*, *rps2*, *rpoB*, *rps12*, *ycf2*,  
275 *ndhG* had  $K_a/K_s$  greater than 0.5 for three species whereas genes *rpoC1*, *atpB*, *rpoA*, *ndhD* had  
276  $K_a/K_s$  ratio more than 4 species and *rpoC2* and *ndhF* had  $K_a/K_s$  ratio for all species. We selected  
277 the genes *atpB*, *rpoA*, *ndhD*, *ndhF*, *rpoC1* and *C2* for further analysis using FUBAR and  
278 MEME. FUBAR estimates the number of nonsynonymous and synonymous substitutions at each  
279 codon given a phylogeny, and provides the posterior probability of every codon belonging to a  
280 set of classes of  $\omega$  (including  $\omega = 1$ ,  $\omega < 1$  or  $\omega > 1$ ) (Murrell et al., 2013). MEME estimates the  
281 probability for a codon to have undergone episodes of positive evolution, allowing the  $\omega$  ratio  
282 distribution to vary across codons and branches in the phylogeny. This last attribute allows  
283 identification of the proportion of codons that may have been evolving neutrally or under  
284 purifying selection, while the remaining codons can also evolve under positive selection (Murrell  
285 et al., 2012). The two models indicated positive selection on the codons only found in *atpB*,  
286 *rpoA*, *ndhF* and *rpoA* (Table 1). Thus, the methods described suggested six amino acid  
287 replacements altogether as candidates for positive selection, of which three were fixed in all  
288 *Nicotiana*, and three were restricted to diverse groups of species (see Table 1).

### 289 **3.5. Repetitive sequences in novel *Nicotiana* plastid genomes**

290 Repeat analysis performed with MISA revealed high similarity in chloroplast microsatellites  
291 (cpSSRs) ranging from 368 to 384 among tobacco species. The majority of the SSRs in these  
292 plastid genomes were mononucleotide rather than trinucleotide or dinucleotide. The most  
293 dominant of the SSRs were A/T motifs mononucleotides, and in dinucleotides AT/TA motifs  
294 were the second most predominant. Mononucleotide SSRs varied from 7-17 units repeats;  
295 dinucleotide SSRs from 4-5-unit repeats while other SSRs types were present mainly in 3-unit  
296 repeats. Mostly the SSRs existed in LSC, in comparison to IR and SSC (Fig 4) (Table S7).  
297 REPuter software was used to identify and locate forward (F), reverse (R), palindromic (P), and  
298 complementary (C) repeats in all the species of *Nicotiana*. In the plastomes of five *Nicotiana*  
299 species, we found 117 oligonucleotide repeats: 25 in *N. knightiana*, 23 in *N. rustica*, 21 in *N.*  
300 *paniculata*, 23 in *N. obtusifolia*, 25 in *N. glauca*. Forward (F) and palindromic repeats were

301 present in large numbers as compared to others in all species: 11 (44%) (F) and 14 (56%) (P) in  
302 *N. knightiana*, 14 (60%) (F) and 9 (39%) (P) in *N. rustica*, and 12 (57%) (F) and 9 (42%) (P) in  
303 *N. paniculata*, 14 (56%) (F) and 11 (44%) (P) in *N. obtusifolia*, 9 (39%) (F) and 11 (52%) (P) in  
304 *N. glauca*. The size of oligonucleotide repeats varied from 30-65 bp, and many of these repeats  
305 were 30-35 bp in length. The LSC region held most of the identified oligonucleotide repeats as  
306 compared to SSC and IR. The LSC region contained 13 in *N. knightiana*, 11 in *N. rustica*, 14 in  
307 *N. paniculata*, 15 in *N. obtusifolia* and 17 in *N. glauca*. In plastid genome regions, the repeats  
308 existed mostly in IGS, followed by CDS and intronic regions (Fig. 5) (Table S8). The number of  
309 tandem repeats varies from 24-27 between these *Nicotiana* species. The IGS region contains the  
310 most tandem repeats followed by the CDS region. The size of these repeats varied between 20 to  
311 88 among *Nicotiana* (Fig. 6).

### 312 **3.6. Single nucleotide polymorphism and insertion/deletion analyses in *Nicotiana***

313 We investigated substitution types in the five plastomes of *Nicotiana* species (one IR removed),  
314 using *Nicotiana tabacum* as a reference. *Nicotiana knightiana* (786), *Nicotiana rustica* (775),  
315 *Nicotiana paniculata* (861), *Nicotiana obtusifolia* (847) and *Nicotiana glauca* (509) substitutions  
316 were seen in the whole plastid genome. The types of substitutions exhibited among *Nicotiana*  
317 species were similar. Most of the conversions were A/G and C/T in comparison to other SNPs  
318 (single nucleotide polymorphism) (Table 2). Ts/Tv ratio were as follows: *Nicotiana knightiana*  
319 LSC (1.5), SSC (0.968) and IR (1.047), *Nicotiana rustica* LSC (1.496), SSC (0.978) and IR (1),  
320 *Nicotiana paniculata* LSC (1.461), SSC (0.886) and IR (0.833), *Nicotiana obtusifolia* LSC  
321 (1.097), SSC (1.020) and IR (1.194), *Nicotiana glauca* LSC (0.924), SSC (0.819) and IR (0.783)  
322 (Table S9). The substitutions in different regions of these genomes are *Nicotiana knightiana*  
323 contains 560 (LSC), 43 (IR) and 183 (SSC) SNPs, *Nicotiana rustica* contains 599 (LSC), 32 (IR)  
324 and 183 (SSC) substitutions, *Nicotiana paniculata* has 630 (LSC), 33 (IR) and 198 (SSC)  
325 substitutions, *Nicotiana obtusifolia* consists of 671 (LSC), 68 (IR) and 210 (SSC) substitutions  
326 while *Nicotiana glauca* has 327 (LSC), 82 (IR) and 100 (SSC). Insertions and Deletions (indels)  
327 were also examined using DnaSP in all regions of the chloroplast genome. In total, *Nicotiana*  
328 *knightiana* (110), *Nicotiana rustica* (107), *Nicotiana paniculata* (116), *Nicotiana obtusifolia*  
329 (143) and *Nicotiana glauca* (113) indels were found. The LSC region held the majority of the  
330 indels, followed by SSC, whereas IR contained minimum indels (Table 3).

### 331 **3.7. Divergence hotspot regions in *Nicotiana***

332 The CDS, intron and IGS regions of the whole plastid genome of five *Nicotiana* species were  
333 compared to discover polymorphic regions (mutational hotspots). High polymorphism was found  
334 in intronic regions (average  $\pi=0.167$ ) in comparison to IGS ( $\pi=0.031$ ) and CDS regions (average  
335  $\pi=0.002$ ). Among *Nicotiana* species, the nucleotide diversity values varied from 0 (*ycf3*) to 0.306  
336 (*rps12 intron* region) (Fig. 7). Here, 20 highly polymorphic regions were determined that might  
337 be used as potential makers to reconstruct the phylogeny for identifying *Nicotiana* species (Table  
338 4).

### 339 **3.8. Phylogenomic analyses**

340 Phylogenetic analysis within *Nicotiana* plastid genomes were reconstructed with the maximum  
341 likelihood method, based on selected and concatenated protein-coding genes. Our phylogenetic  
342 analyses resulted in a highly resolved tree (Fig 8), with almost all clades recovered having  
343 maximum branch support values. After the elimination of indels, the tree was reconstructed  
344 based on alignment size of 75,449 bp with the best fitting model GY+F+I+G4 (Fig 8). We  
345 further concentrated on the species phylogeny of *N. rustica* and putative parental species where  
346 relative divergence times were estimated using a relaxed uncorrelated clock implemented in  
347 BEAST. This analysis found that the divergence of *N. undulata* appeared 5.36 (highest posterior  
348 density, HPD 6.38–4.43) million years ago (Ma), while *N. paniculata* diverged 1.17 Ma (HPD  
349 2.18–0.63) followed by the most recent split of *N. rustica* and *N. knightiana* 0.56 Ma (HPD  
350 0.65–0.46). This analysis showed that the *Nicotiana* species included in the analysis are not older  
351 than the end of the Pliocene and that most subsequent evolution must have occurred in the  
352 Pleistocene. The timing of these lineage splits, in addition to the current distributions of four  
353 closely related species, were used to infer the progression of migratory steps in RASP (Fig 9).  
354 The most recent common ancestor (MRCA) area illustrated a dispersal event for *N. paniculata* in  
355 Northern (B) and Southern Peru (E) and the vicariance of *N. knightiana* in Coastal Peru (D). The  
356 overall dispersal pattern of the examined species showed a south-to-north expansion pattern from  
357 Central Peru to Colombia and Ecuador (*N. rustica*) to Bolivia (*N. undulata*).

## 358 **4. Discussion**

### 359 **4.1. Molecular evolution of *Nicotiana* plastid genomes**

360 We compared five chloroplast genomes of *Nicotiana* species, which revealed similar genomic  
361 features. These comparative analyses produced an insight into the phylogeny and evolution of  
362 *Nicotiana* species. The GC content of the *Nicotiana* species referred to above were similar to  
363 those of other *Nicotiana* species (Sugiyama et al., 2005; Yukawa, Tsudzuki & Sugiura, 2006) i.e.  
364 the GC content in the IR is high, which might be a result of the existence of ribosomal RNA  
365 (Qian et al., 2013; Cheng et al., 2017; Zhao et al., 2018). The genome organization, gene order  
366 and content and of these *Nicotiana* species were also similar for *N. slyvestris* and *N. tabacum*  
367 (Sugiyama et al., 2005; Yukawa, Tsudzuki & Sugiura, 2006). The intron plays an important role  
368 in the regulation of gene expression (Xu et al., 2003). The *trnK* intron is important because it  
369 expresses an unusual form of a group II intron derived from a mobile group of mitochondrial-  
370 like intron open reading frames (ORFs) (Hausner et al., 2006). As in the plastid genomes of  
371 many land plants, the abundance of A/T content at 3<sup>rd</sup> base of codons was reported due to high  
372 concentration of A/T nucleotides in the whole plastid genome (Menezes et al., 2018).

373 The plastomes of land plants have conserved structure but diversity prevails at the border  
374 position of LSC/SSC/IR of the genome. The size range of LSC, SSC and IR varies between the  
375 plastid genomes of species that advances to alterations in several genes and leads to the deletion  
376 of one copy of a gene or duplication of functional and non-functional genes of different sizes  
377 (Menezes et al., 2018; Saina et al., 2018). In the current study, all these ten *Nicotiana* species  
378 showed similarities with some variation as compared to the *Nicotiana tabacum*: in all these  
379 plants except *Nicotiana tomentosiformis* which have 60 bp in IRb region, the *rps19* gene is  
380 present entirely in the LSC region but in *Nicotiana tabacum* *rps19* gene extended 5bp in the IR  
381 region. The fluctuations at the border positions of various regions of the plastid genome might be  
382 helpful in determining the evolution of species (Menezes et al., 2018). Liu et al., (2018) reported  
383 that the similarities at the junction regions may be useful in explaining the relationship between  
384 the species and that those plants which have a high level of relatedness show minimal  
385 fluctuations at the junctions of the chloroplast genome. The resemblance at junctions reveals a  
386 close relationship between the *Nicotiana* species.

387 Repeats in the chloroplast genome are useful in evolutionary studies and play a vital role in  
388 genome arrangement (Zhang et al., 2016). Here, we detected that the mononucleotide repeats  
389 (A/T), and trinucleotide SSRs (ATT/TAA) were present in large amounts in all the species of  
390 *Nicotiana*, which may be a result of the A/T rich proportion of chloroplast genome. A similar

391 result was also reported in *Nicotiana otophora* (Asaf et al., 2016). In all the species of  
392 *Nicotiana*, the LSC region contained a greater amount of SSRs in comparison to SSC and IR,  
393 which has also been demonstrated in other studies of angiosperm plastomes (Shahzadi et al.,  
394 2019; Mehmood et al., 2019). When less genomic resources are available for revealing  
395 divergence hotspot regions, the oligonucleotide repeats might be utilized as a substitute for  
396 identifying polymorphic regions (Ahmed et al., 2012; Ahmad, 2014). The current study results of  
397 oligonucleotide repeats are similar to the previously reported results of *Nicotiana* species and  
398 other angiosperm plastome studies (Asaf et al., 2016; Yang et al., 2019). Thus, the presence of  
399 both the high divergence regions in IGS and oligonucleotide repeats suggest that these regions  
400 are suitable for the development of markers to demonstrate phylogenetics relationships.

401 To understand the molecular evolution, it is important to know about the nucleotide substitution  
402 rates (Muse & Gaut, 1994). LSC and SSC regions are more prone to substitutions and indels  
403 whereas the IR regions are more conserved in the chloroplast genome (Ahmed et al., 2012;  
404 Abdullah et al., 2019b). Our results also showed similar results in that the IR region is mostly  
405 conserved, and most of the substitutions occurs in the LSC and SSC regions. Thus, the ratio  
406 (Ts/Tv) was equal to or more than 1. Similar results were shown in the chloroplast genome of  
407 *Dioscorea polystachya* (Yam) (Cao et al., 2018).

408 Divergence hotspot regions of the plastid genome could be used to develop accurate, robust and  
409 cost-effective molecular markers for population genetics, species barcoding and evolutionary  
410 based studies. (Ahmed et al., 2013; Ahmad, 2014; Nguyen et al., 2017). Previously, in several  
411 studies, polymorphic loci were identified based on comparisons of chloroplast genome to  
412 provide information about suitable loci for the development of molecular markers (Choi, Chung  
413 & Park, 2016; Li et al., 2018; Menezes et al., 2018). We found 20 polymorphic regions such as  
414 *infA*, *rps12 intron*, *rps16-trnQ-UUG* which have 0.25942, 0.15275, 0.08451 nucleotide diversity  
415 respectively that were more polymorphic than frequently used markers such as *rbcL*, and *matK*.  
416 These regions could be suitable markers for population genetics and phylogenetic analyses,  
417 especially in the genus *Nicotiana*.

#### 418 **4.2. Positive selection on *Nicotiana* plastid genes**

419 Plants have evolved complex physiological and biochemical adaptations to adjust and adapt to a  
420 variety of environmental stresses. *Nicotiana*, originating in South America, has spread to many



421 regions of the world and members of the genus have successfully adapted to harsh environmental  
422 conditions to survive. This great variation in their distributional range induced distinctive habits  
423 and morphology in the inflorescence and flowers, indicative of the physiological specialization  
424 to the area where they evolved. Desert ephemeral *Nicotiana* species are short while subtropical  
425 perennials have tall and robust habits with variable inflorescences ranging from pleiochasial  
426 cymes to solitary flowers and diffuse paniculate-cymose mixtures. For example, members of  
427 *Nicotiana* section *Suaveolentes* Goodsp. evolving in isolation faced several cycles of harsh  
428 climate change. In Australia, the native range of the species, a predominantly warm and wet  
429 environment went through intensive aridification (Poczai, Hyvönen & Symon, 2011).  
430 Throughout this climate change and increasing central aridification, many species either retreated  
431 to the wetter coastline or adapted to and still survive in this hostile inland environment (Bally et  
432 al., 2018). Tobacco plants also developed specialized biosynthetic pathways and metabolites,  
433 such as nicotine, which serve complex functions for ecological adaptations to biotic and abiotic  
434 stresses, most importantly serving as a defense mechanism against herbivores (Xu et al. 2017).  
435 Therefore, *Nicotiana* is a rich reservoir of genetic resources for evolutionary biological research,  
436 since several members of the genus went through changing climatic events and adopted to  
437 environmental fluctuations.

438 The patterns of synonymous ( $K_s$ ) and non-synonymous ( $K_a$ ) substitution of nucleotides are  
439 essential markers in evolutionary genetics defining slow and fast evolving genes (Kimura, 1979).  
440  $K_a/K_s$  values  $>1$ ,  $=1$ , and  $<1$  indicate positive selection, natural evolution and purifying selection,  
441 respectively (Lawrie et al., 2013), while a minimal ratio of  $K_a/K_s$  ( $<0.5$ ) in many genes  
442 represents purifying selection working on them. Many proteins and RNA molecules encoded by  
443 the chloroplast genomes are under purifying selection since they are involved in important  
444 functions of plant metabolism, self-replication and photosynthesis and therefore play a pivotal  
445 role in plant survival (Piot et al., 2018). Departure from the main purifying selection in case of  
446 plastid genes might happen in response to certain environmental changes when advantageous  
447 genetic mutations might contribute to survival and better adaptation. The  $K_a/K_s$  ratios in our  
448 analysis for *Nicotiana* species indicated changes in selective pressures. The genes *atpB*, *ndhD*,  
449 *ndhF*, *rpoA*, *rps2* and *rps12* had greater  $K_a/K_s$  value ( $> 1$ ), possibly due to positive selective  
450 pressure as a result of specific environmental conditions. This has been conclusively supported  
451 by an integrative analysis using Fast Unconstrained Bayesian AppRoximation (FUBAR) and

452 Mixed Effects Model of Evolution (MEME) methods, which identified the set of positively  
453 selected codons in case of *atpB*, *ndhD*, *ndhF* and *rpoA* (Table 5), but provided no further  
454 evidence for *rps2* and *rps12*.

455 These genes are involved in different plastid functions, such as DNA replication (*rpoA*) and  
456 photosynthesis (*atpB*, *ndhD* and *ndhF*). The *rpoA* gene encodes the alpha subunit of PEP, which  
457 is believed to predominantly transcribe photosynthesis genes (Hajdukiewicz, Allison & Maliga,  
458 1997). The transcripts of plastid genes encoding the PEP core subunits are transiently  
459 accumulated during leaf development (Kusumi et al., 2011), thus the entire *rpoA* polycistron is  
460 essential for chloroplast gene expression and plant development (Zhang et al., 2018). The  
461 housekeeping gene *atpB* encodes the  $\beta$ -subunit of the ATP synthase complex, which has a highly  
462 conserved structure that couples proton translocation across membranes with the synthesis of  
463 ATP (Gatenby, Rothstein & Nomura, 1989), which is the main source of energy for the  
464 functioning of plant cells. In chloroplasts, linear electron transport mediated by PSII and PSI  
465 produces both ATP and NADPH, whereas PSI cyclic electron transport preferentially contributes  
466 to ATP synthesis without the accumulation of NADPH (Peng & Shikanai, 2011). Chloroplast  
467 NDH monomers are sensitive to high light stress, suggesting that the *ndh* genes encoding the  
468 NAD(P)H dehydrogenase (NDH) may also be involved in stress acclimation through the  
469 optimization of photosynthesis (Casano, Martín & Sabater, 2001; Martin et al., 2002; Rumeau,  
470 Peltier & Cournac, 2007). During acclimation to growth light environments, many plants change  
471 biochemical composition and morphology (Terashima et al., 2005). The highly responsive  
472 regulatory system controlled by cyclic electron transport around PSI could optimize  
473 photosynthesis and plant growth under naturally fluctuating light (Yamori, 2016). When the  
474 demand for ATP is higher than that for NADPH (e.g., during photosynthetic induction, at high or  
475 low temperature, at low CO<sub>2</sub> concentration, or under drought), cyclic electron transport around  
476 PSI is likely to be activated (Yamori, 2016; Yamori & Shikanai, 2016). Thus, positive selection  
477 acting on ATP synthase and NAD(P)H dehydrogenase encoding genes is probably evidence for  
478 adaptation to novel ecological conditions in *Nicotiana*.

479 These findings might be also supported by our observation that RNA editing sites occurred  
480 frequently in *Nicotiana ndh* genes (Table S3). It has been shown that *ndhB* mutants under lower  
481 air humidity conditions or following exposure to ABA present a reduction in the photosynthetic  
482 level, likely mediated through stomatal closure triggered under these conditions (Horvath et al.,

483 2000). Therefore, a protein structure modification resulting from a loss or decrease in RNA  
484 editing events could affect adaptations to stress conditions or cause other unknown changes  
485 (Rodrigues et al., 2017). Previous studies have demonstrated that abiotic stress influences the  
486 editing process and consequently plastid physiology (Nakajima & Mulligan, 2001). Alterations  
487 in editing site patterns resulting from abiotic stress could be associated with susceptibility to  
488 photo-oxidative damage (Rodrigues et al., 2017) and indicate that *Nicotiana* species experienced  
489 abiotic stresses during their evolution, which resulted in positive selection of some of the plastid  
490 genes. Up to this point, positive selection has rarely been detected in chloroplast genes except for  
491 *clpP1* (Erixon & Oxelman, 2008), *ndhF* (Peng et al., 2011), *matK* (Hao, Chen & Xiao, 2010) and  
492 *rbcL* (Kapralov et al., 2011). However, a recent study by Piot et al. (2018) showed that one-third  
493 of the plastid genes in 113 species of grasses (Poaceae) evolved under positive selection. This  
494 might indicate that positive selection might be overlooked among diverse groups of plant taxa.

#### 495 **4.3. Phylogenetic relationships and the origin of tetraploid *Nicotiana rustica***

496 Our comparative plastid genome analysis revealed that the maternal parent of the tetraploid *N.*  
497 *rustica* was the common ancestor of *N. paniculata* and *N. knightiana*, and the later species is  
498 more closely related to *N. rustica*. The relaxed molecular clock analyses estimated that the  
499 speciation event between *N. rustica* and *N. knightiana* appeared ~0.56 Ma (HPD 0.65–0.46) in  
500 line with previous findings (Sierro et al., 2018). Comparative analysis of the genomes of four  
501 related *Nicotiana* species revealed that *N. rustica* inherited about 41% of its nuclear genome  
502 from its paternal progenitor, *N. undulata*, the rest from its maternal progenitor, the common  
503 ancestor of *N. paniculata* and *N. knightiana* (Sierro et al., 2018), which has also been confirmed  
504 by our study. It has been shown that *N. rustica* and in fact all *Nicotiana* tetraploids, except  
505 species included in section *Suaveolentes*, originated from a doubling of the diploid chromosome  
506 for the genus. Thus, they should be regarded as natural allopolyploids (Leitch et al., 2008). We  
507 also revealed that *N. knightiana* is more closely related to *N. rustica* than *N. paniculata*, which  
508 can be further corroborated by the distribution of indels highlighted in the present study. The  
509 biogeographical analysis carried out suggests that *N. undulata* and *N. paniculata* evolved in  
510 North/Central Peru, while *N. rustica* developed in Southern Peru and separated from *N.*  
511 *knightiana*, which adapted to the Southern coastal climatic regimes. Positively selected plastid  
512 genes with functions such as DNA replication (*rpoA*) and photosynthesis (*atpB*, *ndhD* and *ndhF*)  
513 might have been associated with successful adaptation to, for example, a coastal environment.

514 However, our results should be regarded as tentative, as our survey excludes several broad  
515 ecological variables from testing, including variation in salinity, island versus mainland, and East  
516 versus West of the Andes. We aim to highlight that many potential environmental variables  
517 might be highly correlated with speciation processes, as has been demonstrated in the same  
518 region for another Solanaceae group in the tomato clade (*Solanum* sect. *Lycopersicon*), where  
519 amino acid differences in genes associated with seasonal climate variation and intensity of  
520 photosynthetically active radiation were correlated with speciation processes (Pease et al., 2016).  
521 Another example of rapid adaptive radiation from the family is the genus *Nolana* L.f., where  
522 several clades gained competitive advantages in water-dependent environments by succeeding  
523 and diverging in Peru and Northern Chile (Dillon et al., 2009). In the case of *N. rustica* and  
524 related species we assume that diversification was driven by the ecologically variable  
525 environments of the Andes. Our molecular clock analysis provides recent species diversification  
526 in the Pleistocene and Pliocene while substantial climatic transitions in Peru predate these events.  
527 For example, the uplift of the central region of the Andes and the formation of the Peruvian  
528 coastal desert ended (~14 – 150 Mya; Hoorn et al., 2010; Gerreaud et al., 2010) before the  
529 geographical and ecological expansion of *N. rustica* and related parental species.

530 The dispersal of *N. rustica* and related species shows a south-to-north range expansion and  
531 diversification which has been suggested by phylogenies of other plant and animal groups in the  
532 Central Andes (Picard, Sempere & Plantard, 2008; Lueber and Weigend, 2014). Based on the  
533 south-to-north progression scenario, habitats located at high altitudes were first available for  
534 colonization in the south, recently continuing to northward. Erosion and orogenic progression  
535 caused dispersal barriers of species colonizing these high habitats to diversify in a south-to-north  
536 pattern, frequently following allopatric speciation. Thus, for taxonomic groups currently residing  
537 throughout a large portion of the high Andes, a south-to-north speciation pattern is expected  
538 (Doan, 2003). In this case the most basal species (*N. undulata*) has more southern geographic  
539 ranges, and the most derived species (*N. rustica*) has more northern geographic ranges except for  
540 *N. knightiana*, which presumably colonized the coastal range of Peru. Although the four  
541 *Nicotiana* species examined show overlaps in their distribution, it is probable that speciation was  
542 caused by fragmentation of populations during the glacial period (see Simpson, 1975). Utilizing  
543 fewer chloroplast loci for phylogenetic analyses of plant species may limit the solution of  
544 phylogenetic relationships, specifically at low taxonomic levels (Hilu & Alice, 2001; Majure et

545 al., 2012). Previously, genus *Nicotiana* was subdivided into 13 sections using multiple  
546 chloroplast markers, i.e. *trnL* intron and *trnL-F* spacer, *trnS-G* spacer and two genes, *ndhF* and  
547 *matK* (Clarkson et al., 2004). Recently, inference of phylogeny based on complete chloroplast  
548 genomes has provided deep insight into the phylogeny of certain families and genera (Henriquez  
549 et al., 2014; Amiryousefi, Hyvönen & Poczai, 2018a; Abdullah et al., 2019a). Here, we  
550 reconstructed a phylogenetic tree among eleven species of genus *Nicotiana* that belong to nine  
551 sections (Clarkson et al., 2004) based on 75 protein-coding genes by using *S. dulcamara* as an  
552 outgroup which attests the previous classification of genus *Nicotiana* with high bootstrapping  
553 values. Species of each section are well resolved whereas the *N. tabacum* of section *Nicotiana*  
554 and *N. sylvestris* of section *Sylvestres* show close resemblance. The *N. paniculata* and *N.*  
555 *knightiana* belong to section *Paniculatae* but here did not appear on the same node. This revealed  
556 that further data is required to elucidate the phylogenetic relationship among these two species.  
557 Overall, our phylogenetic analyses support the previous classification of genus *Nicotiana*, but  
558 enrichment of chloroplast genomic resources can provide further insight into the phylogeny of  
559 the genus *Nicotiana*.

## 560 **5. Conclusion**

561 In the present study, we assembled, annotated and analyzed the whole cp genome sequence of  
562 five *Nicotiana* species. The genomic structure and organization of their chloroplast genome was  
563 like those of previously reported Solanaceae plastomes. Divergences of LSC, SSC and IR region  
564 sequences were identified, as well as the distribution and location of repeat sequences. The  
565 identified mutational hotspots sequences could be utilized as potential molecular markers to  
566 investigate phylogenetic relationships in the genus. As we demonstrated in our study to elucidate  
567 the maternal genome origins of *N. rustica*, our results could provide further help in  
568 understanding the evolutionary history of tobaccos.

## 569 **Acknowledgements**

570 We thank Kenneth Quek for editing the manuscript.

## 571 **Competing interest**

572 The authors declare that they have no conflict of interest.

## 573 **Authors contributions**

574 Furrukh Mehmood: Conceptualization, Genome assembly and annotation, Data analysis, Data  
575 interpretation, prepared figures and tables, Manuscript drafting and editing.

576 Abdullah: Genome annotation, Data analysis, Data interpretation, Manuscript drafting.

577 Zartasha Ubaid: Data analysis, Data interpretation, Manuscript drafting.

578 Iram Shehzadi: Data analysis, Data interpretation, Manuscript drafting.

579 Ibrar Ahmed: Conceptualization, Manuscript editing.

580 Mohammad Tahir Waheed: Conceptualization, Manuscript editing.

581 Péter Poczai: Supervision, carried out selection tests and phylogenetic analysis, prepared figures  
582 and tables, authored and reviewed drafts of the paper, approved the final draft.

583 Bushra Mirza: Supervision, authored or reviewed drafts of the paper, approved the final draft.

## 584 **References**

585 **Abdullah, Mehmood F, Shahzadi I, Waseem S, Mirza B, Ahmed I, Waheed MT. 2019a.**  
586 Chloroplast genome of *Hibiscus rosa-sinensis* (Malvaceae): Comparative analyses and  
587 identification of mutational hotspots. *Genomics* DOI: 10.1016/j.ygeno.2019.04.010.

588 **Abdullah, Shahzadi I, Mehmood F, Ali Z, Malik MS, Waseem S, Mirza B, Ahmed I,**  
589 **Waheed MT. 2019b.** Comparative analyses of chloroplast genomes among three *Firmiana*  
590 species: Identification of mutational hotspots and phylogenetic relationship with other  
591 species of Malvaceae. *Plant Gene*:100199. DOI: 10.1016/j.plgene.2019.100199.

592 **Ahmad I. 2014.** Evolutionary dynamics in taro. Massey University, Palmerston North, New  
593 Zealand.

594 **Ahmed I, Biggs PJ, Matthews PJ, Collins LJ, Hendy MD, Lockhart PJ. 2012.** Mutational  
595 dynamics of aroid chloroplast genomes. *Genome Biology and Evolution* **4**:1316–1323. DOI:  
596 10.1093/gbe/evs110.

597 **Ahmed I, Matthews PJ, Biggs PJ, Naeem M, Mclenachan PA, Lockhart PJ. 2013.**  
598 Identification of chloroplast genome loci suitable for high-resolution phylogeographic  
599 studies of *Colocasia esculenta* (L.) Schott (Araceae) and closely related taxa. *Molecular*  
600 *Ecology Resources* **13**:929–937. DOI: 10.1111/1755-0998.12128.



- 601 **Amiryousefi A, Hyvönen J, Poczai P. 2018a.** The chloroplast genome sequence of bittersweet  
602 (*Solanum dulcamara*): Plastid genome structure evolution in Solanaceae. *PLoS ONE* **13**: 1–  
603 23. DOI: 10.1371/journal.pone.0196069.
- 604 **Amiryousefi A, Hyvönen J, Poczai P. 2018b.** IRscope: an online program to visualize the  
605 junction sites of chloroplast genomes. *Bioinformatics* **34**: 3030–3031 DOI:  
606 10.1093/bioinformatics/bty220.
- 607 **Aoki S, Ito M. 2000.** Molecular phylogeny of *Nicotiana* (Solanaceae) based on the nucleotide  
608 sequence of the *matK* gene. *Plant Biology* **2**: 316–324. DOI: 10.1055/s-2000-3710.
- 609 **Asaf S, Khan AL, Khan AR, Waqas M, Kang S-M, Khan MA, Lee S-M, Lee I-J. 2016.**  
610 Complete Chloroplast Genome of *Nicotiana otophora* and its Comparison with Related  
611 Species. *Frontiers in Plant Science* **7**: 1–12. DOI: 10.3389/fpls.2016.00843.
- 612 **Beier S, Thiel T, Münch T, Scholz U, Mascher M. 2017.** MISA-web: a web server for  
613 microsatellite prediction. *Bioinformatics* **33**: 2583–2585. DOI:  
614 10.1093/bioinformatics/btx198.
- 615 **Benson G. 1999.** Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids*  
616 *Research* **27**: 573–580. DOI: 10.1093/nar/27.2.573.
- 617 **Cao J, Jiang D, Zhao Z, Yuan S, Zhang Y, Zhang T, Zhong W, Yuan Q, Huang L. 2018.**  
618 Development of Chloroplast Genomic Resources in Chinese Yam (*Dioscorea polystachya*).  
619 *BioMed Research International*. DOI: 10.1155/2018/6293847.
- 620 **Casano LM, Martín M, Sabater B. 2001.** Hydrogen peroxide mediates the induction of  
621 chloroplastic *ndh* complex under photooxidative stress in Barley. *Plant Physiology* **125**:  
622 1450–1458 DOI: 10.1104/pp.125.3.1450.
- 623 **Chase MW, Knapp S, Cox A V., Clarkson JJ, Butsko Y, Joseph J, Savolainen V,**  
624 **Parokony AS. 2003.** Molecular systematics, GISH and the origin of hybrid taxa in  
625 *Nicotiana* (Solanaceae). *Annals of Botany* **92**: 107–127 DOI: 10.1093/aob/mcg087.
- 626 **Chen X, Cui Y, Nie L, Hu H, Xu Z, Sun W, Gao T, Song J, Yao H. 2019.** Identification and  
627 Phylogenetic Analysis of the Complete Chloroplast Genomes of Three *Ephedra* Herbs  
628 Containing Ephedrine. *BioMed Research International* **2019**: 1–10. DOI:  
629 10.1155/2019/5921725.

- 630 **Cheng H, Li J, Zhang H, Cai B, Gao Z, Qiao Y, Mi L. 2017.** The complete chloroplast  
631 genome sequence of strawberry (*Fragaria × ananassa* Duch.) and comparison with related  
632 species of Rosaceae. *PeerJ* **5**: e3919 DOI: 10.7717/peerj.3919.
- 633 **Choi KS, Chung MG, Park S. 2016.** The complete chloroplast genome sequences of three  
634 Veroniceae species (Plantaginaceae): comparative analysis and highly divergent regions.  
635 *Frontiers in Plant Science* **7**:1–8. DOI: 10.3389/fpls.2016.00355.
- 636 **Cooper G. 2000.** *Chloroplasts and other plastids in the cell: A molecular approach. 2nd edition.*
- 637 **Couvreur TLP, Chatrou LW, Sosef MSM, Richardson JE. 2008.** Molecular phylogenetics  
638 reveal multiple tertiary vicariance origins of the African rain forest trees. *BMC Biology* **6**:  
639 54.
- 640 **Daniell H. 2007.** Transgene containment by maternal inheritance: Effective or elusive?  
641 *Proceedings of the National Academy of Sciences of the USA* **104**: 6879–6880 DOI:  
642 10.1073/pnas.0702219104.
- 643 **Daniell H, Lin C-S, Yu M, Chang W-J. 2016.** Chloroplast genomes: diversity, evolution, and  
644 applications in genetic engineering. *Genome Biology* **17**:134. DOI: 10.1186/s13059-016-  
645 1004-2.
- 646 **Delpont W, Poon AFY, Frost SDW, Kosakovsky Pond SL. 2010.** Datamonkey 2010: A suite  
647 of phylogenetic analysis tools for evolutionary biology. *Bioinformatics* **26**: 2455–2457 DOI:  
648 10.1093/bioinformatics/btq429.
- 649 **Dereeper A, Guignon V, Blanc G, Audic S, Buffet S, Chevenet F, Dufayard J-F, Guindon S,  
650 Lefort V, Lescot M, Claverie J-M, Gascuel O. 2008.** Phylogeny.fr: robust phylogenetic  
651 analysis for the non-specialist. *Nucleic Acids Research* **36**: W465–W469. DOI:  
652 10.1093/nar/gkn180.
- 653 **Dillon MO, Tu T, Xie L, Quipuscia Silvestre V. 2009.** Biogeographic diversification in *Nolana*  
654 (Solanaceae), a ubiquitous member of the Atacama and Peruvian Deserts along the western  
655 coast of South America. *Journal of Systematics and Evolution* **47**: 457–476
- 656 **Doan TM. 2003.** A south-to-north biogeographic hypothesis for Andean speciation: evidence  
657 from the lizard genus *Proctoporus* (Reptilia, Gymnophthalmidae). *Journal of Biogeography*  
658 **30**: 361–374

- 659 **Drummond AJ, Ho SYW, Phillips MJ, Rambaut A. 2006.** Relaxed phylogenetics and dating  
660 with confidence. *PLoS Biology* **4**: e88
- 661 **Drummond AJ, Suchard M, Xie D, Rambaut A. 2012.** Bayesian phylogenetics with BEAUti  
662 and the BEAST 1.7. *Molecular Biology and Evolution* **29**:1969 – 1973
- 663 **Erixon P, Oxelman B. 2008.** Whole-gene positive selection, elevated synonymous substitution  
664 rates, duplication, and indel evolution of the chloroplast *clpP1* gene. *PLoS ONE* **3**: e1386  
665 DOI: 10.1371/journal.pone.0001386.
- 666 **Evans M, Aubriot X, Hearn D, Lanciaux M, Lavergne S, Cruaud C, Lowry II PP,  
667 Haevermans R. 2014.** The evolution of succulence: insights from a remarkable radiation in  
668 Madagascar. *Systematic Biology* **63**: 698 – 711.
- 669 **Garreaud RD, Molina A, Farias M. 2010.** Aean uplift ocean cooling and Atacama  
670 hyperaridity: a climate modeling perspective. *Earth and Planetary Science Letters* **292**: 39–  
671 50
- 672 **Gatenby AA, Rothstein SJ, Nomura M. 1989.** Translational coupling of the maize chloroplast  
673 *atpB* and *atpE* genes. *Proceedings of the National Acadademy of Sciences of the USA* **86**:  
674 4066–4070
- 675 **Goodspeed TH. 1954.** The Genus *Nicotiana*. Chronica Botanica, New York, USA  
676 **Greiner S, Lehwark P, Bock R. 2019.** OrganellarGenomeDRAW (OGDRAW) version 1.3.1:  
677 expanded toolkit for the graphical visualization of organellar genomes. *Nucleic Acids  
678 Research* **47**: W59–W64 DOI: 10.1093/nar/gkz238.
- 679 **Hajdukiewicz PTJ, Allison LA, Maliga P. 1997.** The two RNA polymerases encoded by the  
680 nuclear and the plastid compartments transcribe distinct groups of genes in tobacco plastids.  
681 *EMBO Journal* **16**: 4041–4048 DOI: 10.1093/emboj/16.13.4041.
- 682 **Hao DC, Chen SL, Xiao PG. 2010.** Molecular evolution and positive Darwinian selection of the  
683 chloroplast maturase *matK*. *Journal of Plant Research* **123**: 241–247 DOI: 10.1007/s10265-  
684 009-0261-5.
- 685 **Hausner G, Olson R, Simon D, Johnson I, Sanders ER, Karol KG, McCourt RM,  
686 Zimmerly S. 2006.** Origin and evolution of the chloroplast *trnK (matK)* intron: A model for  
687 evolution of group II intron RNA structures. *Molecular Biology and Evolution* **23**: 380–391

- 688 DOI: 10.1093/molbev/msj047.
- 689 **Henriquez CL, Arias T, Pires JC, Croat TB, Schaal BA. 2014.** Phylogenomics of the plant  
690 family Araceae. *Molecular Phylogenetics and Evolution* **75**: 91–102. DOI:  
691 10.1016/j.ympev.2014.02.017.
- 692 **Hilu KW, Alice LA. 2001.** A phylogeny of Chloridoideae (Poaceae) based on *matK* sequences.  
693 *Systematic Botany* **26**: 386–405
- 694 **Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS. 2018.** UFBoot2: improving  
695 the ultrafast bootstrap approximation. *Molecular Biology and Evolution* **35**: 518–522
- 696 **Hoorn C, Wesselingh FP, ter Steege H, Bermudez MA, Mora A, Sevink J, et al. 2010.**  
697 Amazonia through tie: Andean uplift, climate change, landscape evolution, and biodiversity.  
698 *Science* **330**: 927–931
- 699 **Horvath EM, Peter SO, Joet T, Rumeau D, Cournac L, Horvath G V., Kavanagh TA,**  
700 **Schafer C, Peltier G, Medgyesy P. 2000.** Targeted inactivation of the plastid *ndhB* gene in  
701 tobacco results in an enhanced sensitivity of photosynthesis to moderate stomatal closure.  
702 *Plant Physiology* **123**: 1337–1350
- 703 **Jheng C-F, Chen T-C, Lin J-Y, Chen T-C, Wu W-L, Chang C-C. 2012.** The comparative  
704 chloroplast genomic analysis of photosynthetic orchids and developing DNA markers to  
705 distinguish *Phalaenopsis* orchids. *Plant Science* **190**: 62–73. DOI:  
706 10.1016/j.plantsci.2012.04.001.
- 707 **Jin S, Daniell H. 2015.** Engineered chloroplast genome just got smarter. *Trends in Plant Science*  
708 **20**:622-640 DOI: 10.1016/j.tplants.2015.07.004.
- 709 **Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermin LS. 2017.**  
710 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature Methods* **14**:  
711 587–589. DOI: 10.1038/nmeth.4285.
- 712 **Kapralov M V., Kubien DS, Andersson I, Filatov DA. 2011.** Changes in Rubisco kinetics  
713 during the evolution of C4 Photosynthesis in *Flaveria* (Asteraceae) are associated with  
714 positive selection on genes encoding the enzyme. *Molecular Biology and Evolution* **28**:  
715 1491–1503. DOI: 10.1093/molbev/msq335.

- 716 **Katoh K, Kuma KI, Toh H, Miyata T. 2005.** MAFFT version 5: Improvement in accuracy of  
717 multiple sequence alignment. *Nucleic Acids Research* **33**: 511–518. DOI:  
718 10.1093/nar/gki198.
- 719 **Katoh K, Standley DM. 2013.** MAFFT multiple sequence alignment software version 7:  
720 Improvements in performance and usability. *Molecular Biology and Evolution* **30**: 772–780.  
721 DOI: 10.1093/molbev/mst010.
- 722 **Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A,  
723 Markowitz S, Duran C, Thierer T, Ashton B, Meintjes P, Drummond A. 2012.**  
724 Geneious Basic: an integrated and extendable desktop software platform for the  
725 organization and analysis of sequence data. *Bioinformatics* **28**: 1647–1649. DOI:  
726 10.1093/bioinformatics/bts199.
- 727 **Kimura M. 1979.** Model of effectively neutral mutations in which selective constraint is  
728 incorporated. *Proceedings of the National Academy of Sciences of the USA* **76**: 3440–3444  
729 DOI: 10.1073/pnas.76.7.3440.
- 730 **Knapp S, Bohs L, Nee M, Spooner DM. 2004.** Solanaceae - A model for linking genomics with  
731 biodiversity. *Comparative and Functional Genomics* **5**: 285–291. DOI: 10.1002/cfg.393.
- 732 **Knapp S, Chase MW, Clarkson JJ. 2004.** Nomenclatural changes and a new sectional  
733 classification in *Nicotiana* (Solanaceae). *Taxon* **53**: 73-82. DOI: 10.2307/4135490.
- 734 **Kumar S, Strecher G, Suleski M, Hedges SB. 2017.** TimeTree: a resource for timelines,  
735 timetrees, and divergence times. *Molecular Biology and Evolution* **34**: 1812–1819.
- 736 **Kurtz S. 2002.** REPuter: the manifold applications of repeat analysis on a genomic scale.  
737 *Nucleic Acids Research* **29**: 4633–4642. DOI: 10.1093/nar/29.22.4633.
- 738 **Kusumi K, Sakata C, Nakamura T, Kawasaki S, Yoshimura A, Iba K. 2011.** A plastid  
739 protein NUS1 is essential for build-up of the genetic system for early chloroplast  
740 development under cold stress conditions. *Plant Journal* **68**: 1039–1050. DOI:  
741 10.1111/j.1365-313X.2011.04755.x.
- 742 **Laslett D, Canback B. 2004.** ARAGORN, a program to detect tRNA genes and tmRNA genes  
743 in nucleotide sequences. *Nucleic Acids Research* **32**: 11–16. DOI: 10.1093/nar/gkh152.

- 744 **Lawrie DS, Messer PW, Hershberg R, Petrov DA. 2013.** Strong Purifying Selection at  
745 Synonymous Sites in *D. melanogaster*. *PLoS Genetics* **9**: e1003527. DOI:  
746 10.1371/journal.pgen.1003527.
- 747 **Leitch IJ, Hanson L, Lim KY, Kovarik A, Chase MW, Clarkson JJ, Leitch AR. 2008.** The  
748 ups and downs of genome size evolution in polyploid species of *Nicotiana* (Solanaceae).  
749 *Annals of Botany* **101**: 805–814. DOI: 10.1093/aob/mcm326.
- 750 **Lemoine F, Correia D, Lefort V, Doppelt-Azeroual O, Mareuil F, Cohen-Boulakia S,**  
751 **Gascuel O. 2019.** NGPhylogeny.fr: new generation phylogenetic services for non-  
752 specialists. *Nucleic Acids Research* **47**: W260–W265. DOI: 10.1093/nar/gkz303.
- 753 **Li H, Durbin R. 2009.** Fast and accurate short read alignment with Burrows-Wheeler transform.  
754 *Bioinformatics* **25**:1754–1760. DOI: 10.1093/bioinformatics/btp324.
- 755 **Li Y, Zhang Z, Yang J, Lv G. 2018.** Complete chloroplast genome of seven *Fritillaria* species,  
756 variable DNA markers identification and phylogenetic relationships within the genus. *PLoS*  
757 *ONE* **13**: e0194613 DOI: 10.1371/journal.pone.0194613.
- 758 **Liu L, Wang Y, He P, Li P, Lee J, Soltis DE, Fu C. 2018.** Chloroplast genome analyses and  
759 genomic resource development for epilithic sister genera *Oresitrophe* and *Mukdenia*  
760 (Saxifragaceae), using genome skimming data. *BMC Genomics* **19**: 1–17. DOI:  
761 10.1186/s12864-018-4633-x.
- 762 **Lowe TM, Chan PP. 2016.** tRNAscan-SE On-line: integrating search and context for analysis of  
763 transfer RNA genes. *Nucleic Acids Research* **44**: W54–W57 DOI: 10.1093/nar/gkw413.
- 764 **Luebert F, Weigend M. 2014.** Phylogenetic insights into Andean plant diversification.  
765 *Frontiers in Ecology and Evolution* **2**: 27
- 766 **Majure LC, Puente R, Patrick Griffith M, Judd WS, Soltis PS, Soltis DE. 2012.** Phylogeny  
767 of *Opuntia* s.s. (Cactaceae): Clade delineation, geographic origins, reticulate evolution.  
768 *American Journal of Botany* **99**: 847–864 DOI: 10.3732/ajb.1100375.
- 769 **Martin W, Rujan T, Richly E, Hansen A, Cornelsen S, Lins T, Leister D, Stoebe B,**  
770 **Hasegawa M, Penny D. 2002.** Evolutionary analysis of *Arabidopsis*, cyanobacterial, and  
771 chloroplast genomes reveals plastid phylogeny and thousands of cyanobacterial genes in the  
772 nucleus. *Proceedings of the National Academy of Sciences of the USA* **99**: 12246–12251.



- 773 DOI: 10.1073/pnas.182432999.
- 774 **Mehmood F, Abdullah, Shahzadi I, Ahmed I, Waheed MT, Mirza B. 2019.** Characterization  
775 of *Withania somnifera* chloroplast genome and its comparison with other selected species of  
776 Solanaceae. *Genomics*. DOI: 10.1016/J.YGENO.2019.08.024.
- 777 **Menezes APA, Resende-Moreira LC, Buzatti RSO, Nazareno AG, Carlsen M, Lobo FP,**  
778 **Kalapothakis E, Lovato MB. 2018.** Chloroplast genomes of *Byrsonima* species  
779 (Malpighiaceae): Comparative analysis and screening of high divergence sequences.  
780 *Scientific Reports* **8**: 1–12. DOI: 10.1038/s41598-018-20189-4.
- 781 **Milne I, Bayer M, Cardle L, Shaw P, Stephen G, Wright F, Marshall D. 2009.** Tablet-next  
782 generation sequence assembly visualization. *Bioinformatics* **26**: 401–402. DOI:  
783 10.1093/bioinformatics/btp666.
- 784 **Mower JP. 2009.** The PREP suite: Predictive RNA editors for plant mitochondrial genes,  
785 chloroplast genes and user-defined alignments. *Nucleic Acids Research* **37**: W253–W259.  
786 DOI: 10.1093/nar/gkp337.
- 787 **Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Kosakovsky Pond SL, Scheffler K.**  
788 **2013.** FUBAR: A fast, unconstrained bayesian AppRoximation for inferring selection.  
789 *Molecular Biology and Evolution* **30**: 1196-1205. DOI: 10.1093/molbev/mst030.
- 790 **Murrell B, Wertheim JO, Moola S, Weighill T, Scheffler K, Kosakovsky Pond SL. 2012.**  
791 Detecting individual sites subject to episodic diversifying selection. *PLoS Genetics* **8**:  
792 e1002764. DOI: 10.1371/journal.pgen.1002764.
- 793 **Muse S V, Gaut BS. 1994.** A likelihood approach for comparing synonymous and  
794 nonsynonymous nucleotide substitution rates, with application to the chloroplast genome.  
795 *Molecular Biology and Evolution* **11**: 715–724. DOI:  
796 10.1093/oxfordjournals.molbev.a040152.
- 797 **Nakajima Y, Mulligan RM. 2001.** Heat stress results in incomplete C-to-U editing of maize  
798 chloroplast mRNAs and correlates with changes in chloroplast transcription rate. *Current*  
799 *Genetics* **40**: 209–213. DOI: 10.1007/s002940100249.
- 800 **Neale DB, Sederoff RR. 1989.** Paternal inheritance of chloroplast DNA and maternal  
801 inheritance of mitochondrial DNA in loblolly pine. *Theoretical and Applied Genetics* **77**:

- 802 212–216. DOI: 10.1007/BF00266189.
- 803 **Nguyen VB, Park H-S, Lee S-C, Lee J, Park JY, Yang T-J. 2017.** Authentication markers for  
804 five major *Panax* species developed via comparative analysis of complete chloroplast  
805 genome sequences. *Journal of Agricultural and Food Chemistry* **65**: 6298–6306 DOI:  
806 10.1021/acs.jafc.7b00925.
- 807 **Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015.** IQ-TREE: A fast and effective  
808 stochastic algorithm for estimating Maximum-likelihood phylogenies. *Molecular Biology*  
809 *and Evolution* **32**: 268–274. DOI: 10.1093/molbev/msu300.
- 810 **Occhialini A, Lin MT, Andralojc PJ, Hanson MR, Parry MAJ. 2016.** Transgenic tobacco  
811 plants with improved cyanobacterial Rubisco expression but no extra assembly factors grow  
812 at near wild-type rates if provided with elevated CO<sub>2</sub>. *Plant Journal* **85**: 148–160. DOI:  
813 10.1111/tpj.13098.
- 814 **Oldenburg DJ, Bendich AJ. 2015.** DNA maintenance in plastids and mitochondria of plants.  
815 *Frontiers in Plant Science* **6**: 883. DOI: 10.3389/fpls.2015.00883.
- 816 **Oldenburg DJ, Bendich AJ. 2016.** The linear plastid chromosomes of maize: terminal  
817 sequences, structures, and implications for DNA replication. *Current Genetics* **62**: 431–442.  
818 DOI: 10.1007/s00294-015-0548-0.
- 819 **Olmstead RG, Bohs L. 2007.** A summary of molecular systematic research in solanaceae: 1982-  
820 2006. *Acta Horticulturae* **745**: 255–268. DOI: 10.17660/ActaHortic.2007.745.11.
- 821 **Olmstead RG, Bohs L, Migid HA, Santiago-valentin E. 2008.** A molecular phylogeny of the  
822 Solanaceae. *Taxon* **57**: 1159–1181
- 823 **Palmer JD. 1985.** Comparative organization of chloroplast genomes. *Annual Review of Genetics*  
824 **19**: 325–354. DOI: 10.1146/annurev.ge.19.120185.001545.
- 825 **Pease JB, Haak DC, Hahn MW, Moyle LC. 2016.** Phylogenomics reveals three sources of  
826 adaptive variation during a rapid radiation. *PLoS Biology* **14**: e1002379
- 827 **Peng L, Shikanai T. 2011.** Supercomplex formation with photosystem I is required for the  
828 stabilization of the chloroplast NADH dehydrogenase-like complex in *Arabidopsis*. *Plant*  
829 *Physiology* **155**: 1629–1639. DOI: 10.1104/pp.110.171264.

- 830 **Picard D, Sempere T, Plantard O. 2008.** Direction and timing of uplift propagation in the  
831 Peruvian Andes deduced from molecular phylogenetics of highland biotaxa. *Earth and*  
832 *Planetary Science Letters* **271**: 326–336
- 833 **Piot A, Hackel J, Christin PA, Besnard G. 2018.** One-third of the plastid genes evolved under  
834 positive selection in PACMAD grasses. *Planta* **247**: 255–266. DOI: 10.1007/s00425-017-  
835 2781-x.
- 836 **Poczai P, Hyvönen J, Symon DE. 2011.** Phylogeny of kangaroo apples (*Solanum* subg.  
837 *Archaeosolanum*, Solanaceae). *Molecular Biology Reports* **38**: 5243–5259. DOI:  
838 10.1007/s11033-011-0675-8.
- 839 **Qian J, Song J, Gao H, Zhu Y, Xu J, Pang X, Yao H, Sun C, Li X, Li C, Liu J, Xu H, Chen**  
840 **S. 2013.** The complete chloroplast genome sequence of the medicinal plant *Salvia*  
841 *miltiorrhiza*. *PLoS ONE* **8**: e57607. DOI: 10.1371/journal.pone.0057607.
- 842 **Rambaut A, Suchard MA, Xie D, Drummond AJ. 2014.** Tracer v1.6. Computer program and  
843 documentation distributed by the author, website: <http://beast.community/tracer> . [accessed  
844 3 December 2019].
- 845 **Ranwez V, Harispe S, Delsuc F, Douzery EJP. 2011.** MACSE: multiple alignment of coding  
846 SEquences Accounting for frameshifts and stop codons. *PLoS One* **6**: e22594
- 847 **Ravi V, Khurana JP, Tyagi AK, Khurana P. 2008.** An update on chloroplast genomes. *Plant*  
848 *Systematics and Evolution* **271**: 101–122. DOI: 10.1007/s00606-007-0608-0.
- 849 **Rodrigues NF, Christoff AP, da Fonseca GC, Kulcheski FR, Margis R. 2017.** Unveiling  
850 chloroplast RNA editing events using next generation small RNA sequencing data.  
851 *Frontiers in Plant Science* **8**: 1686. DOI: 10.3389/fpls.2017.01686.
- 852 **Ronquist F. 2004.** Bayesian inference of character evolution. *Trends in Ecology and Evolution*  
853 **19**: 475–481
- 854 **Rozas J, Ferrer-Mata A, Sanchez-DelBarrio JC, Guirao-Rico S, Librado P, Ramos-Onsins**  
855 **SE, Sanchez-Gracia A. 2017.** DnaSP 6: DNA sequence polymorphism analysis of large  
856 data sets. *Molecular Biology and Evolution* **34**: 3299–3302. DOI: 10.1093/molbev/msx248.
- 857 **Rumeau D, Peltier G, Cournac L. 2007.** Chlororespiration and cyclic electron flow around PSI

- 858 during photosynthesis and plant stress response. *Plant, Cell and Environment* **30**: 1041–  
859 1051. DOI: 10.1111/j.1365-3040.2007.01675.x.
- 860 **Saina JK, Li ZZ, Gichira AW, Liao YY. 2018.** The complete chloroplast genome sequence of  
861 tree of heaven (*Ailanthus altissima* (mill.) (sapindales: Simaroubaceae), an important  
862 pantropical tree. *International Journal of Molecular Sciences* **19**: E929 DOI:  
863 10.3390/ijms19040929.
- 864 **Shahzadi I, Abdullah, Mehmood F, Ali Z, Ahmed I, Mirza B. 2019.** Chloroplast genome  
865 sequences of *Artemisia maritima* and *Artemisia absinthium*: Comparative analyses,  
866 mutational hotspots in genus *Artemisia* and phylogeny in family Asteraceae. *Genomics*.  
867 DOI: 10.1016/J.YGENO.2019.08.016.
- 868 **Shi L, Chen H, Jiang M, Wang L, Wu X, Huang L, Liu C. 2019.** CPGAVAS2, an integrated  
869 plastome sequence annotator and analyzer. *Nucleic Acids Research* **47**: W65–W73. DOI:  
870 10.1093/nar/gkz345.
- 871 **Shinozaki K, Ohme M, Tanaka M, Wakasugi T, Hayashida N, Matsubayashi T, Zaita N,  
872 Chunwongse J, Obokata J, Yamaguchi-Shinozaki K, Ohto C, Torazawa K, Meng BY,  
873 Sugita M, Deno H, Kamogashira T, Yamada K, Kusuda J, Takaiwa F, Kato A,  
874 Tohdoh N, Shimada H, Sugiura M. 1986.** The complete nucleotide sequence of the  
875 tobacco chloroplast genome: its gene organization and expression. *The EMBO Journal* **5**:  
876 2043–2049. DOI: 10.1002/j.1460-2075.1986.tb04464.x.
- 877 **Sierro N, Battey JND, Bovet L, Liedschulte V, Ouadi S, Thomas J, Broye H, Laparra H,  
878 Vuarnoz A, Lang G, Goepfert S, Peitsch MC, Ivanov N V. 2018.** The impact of genome  
879 evolution on the allotetraploid *Nicotiana rustica* - An intriguing story of enhanced alkaloid  
880 production. *BMC Genomics* **19**: 855 DOI: 10.1186/s12864-018-5241-5.
- 881 **Simpson, BB. 1975.** Pleistocene changes in the Flora of the high tropical Andes. *Paleobiology* **1**:  
882 273–294
- 883 **Smith, H.H. 1974.** *Nicotiana*. In: King RC (eds) Handbook of Genetics 2. Plenum Press, New  
884 York, pp 281-314
- 885 **Sugiyama Y, Watase Y, Nagase M, Makita N, Yagura S, Hirai A, Sugiura M. 2005.** The  
886 complete nucleotide sequence and multipartite organization of the tobacco mitochondrial

- 887 genome: Comparative analysis of mitochondrial genomes in higher plants. *Molecular*  
888 *Genetics and Genomics* **272**: 603–615. DOI: 10.1007/s00438-004-1075-8.
- 889 **Terashima I, Araya T, Miyazawa SI, Sone K, Yano S. 2005.** Construction and maintenance of  
890 the optimal photosynthetic systems of the leaf, herbaceous plant and tree: An eco-  
891 developmental treatise. *Annals of Botany* **95**: 507–519. DOI: 10.1093/aob/mci049.
- 892 **Tillich M, Lehwark P, Pellizzer T, Ulbricht-Jones ES, Fischer A, Bock R, Greiner S. 2017.**  
893 GeSeq – versatile and accurate annotation of organelle genomes. *Nucleic Acids Research*  
894 **45**: W6–W11. DOI: 10.1093/nar/gkx391.
- 895 **Waheed MT, Ismail H, Gottschamel J, Mirza B, Lössl AG. 2015.** Plastids: the green frontiers  
896 for vaccine production. *Frontiers in Plant Science* **6**: 1005. DOI: 10.3389/fpls.2015.01005.
- 897 **Waheed MT, Thönes N, Müller M, Hassan SW, Razavi NM, Lössl E, Kaul HP, Lössl AG.**  
898 2011. Transplastomic expression of a modified human papillomavirus L1 protein leading to  
899 the assembly of capsomeres in tobacco: A step towards cost-effective second-generation  
900 vaccines. *Transgenic Research* **20**: 271–281. DOI: 10.1007/s11248-010-9415-4.
- 901 **Wambugu PW, Brozynska M, Furtado A, Waters DL, Henry RJ. 2015.** Relationships of  
902 wild and domesticated rices (*Oryza* AA genome species) based upon whole chloroplast  
903 genome sequences. *Scientific Reports* **5**:13957. DOI: 10.1038/srep13957.
- 904 **Xu J, Feng D, Song G, Wei X, Chen L, Wu X, Li X, Zhu Z. 2003.** The first intron of rice  
905 EPSP synthase enhances expression of foreign gene. *Science in China. Series C, Life*  
906 *sciences / Chinese Academy of Sciences* **46**: 561–569. DOI: 10.1360/02yc0120.
- 907 **Xu J-H, Liu Q, Hu W, Wang T, Xue Q, Messing J. 2015.** Dynamics of chloroplast genomes in  
908 green plants. *Genomics* **106**: 221–231. DOI: 10.1016/J.YGENO.2015.07.004.
- 909 **Yamori W. 2016.** Photosynthetic response to fluctuating environments and photoprotective  
910 strategies under abiotic stress. *Journal of Plant Research* **129**: 379–395. DOI:  
911 10.1007/s10265-016-0816-1.
- 912 **Yamori W, Shikanai T. 2016.** Physiological functions of cyclic electron transport around  
913 photosystem I in sustaining photosynthesis and plant growth. *Annual Review of Plant*  
914 *Biology* **67**: 81–106. DOI: 10.1146/annurev-arplant-043015-112002.

- 915 **Yang Z, Wang G, Ma Q, Ma W, Liang L, Zhao T. 2019.** The complete chloroplast genomes  
916 of three Betulaceae species: Implications for molecular phylogeny and historical  
917 biogeography. *PeerJ* **7**: e6320. DOI: 10.7717/peerj.6320.
- 918 **Yu Y, Blair C, He XJ. 2019.** RASP (Reconstruct Ancestral State in Phylogenies): a tool for  
919 historical biogeography. *Molecular Biology and Evolution* doi: 10.1093/molbev/msz257
- 920 **Yukawa M, Tsudzuki T, Sugiura M. 2006.** The chloroplast genome of *Nicotiana sylvestris* and  
921 *Nicotiana tomentosiformis*: Complete sequencing confirms that the *Nicotiana sylvestris*  
922 progenitor is the maternal genome donor of *Nicotiana tabacum*. *Molecular Genetics and*  
923 *Genomics* **275**: 367–373. DOI: 10.1007/s00438-005-0092-6.
- 924 **Zhang Y, Cui YL, Zhang XL, Yu QB, Wang X, Yuan XB, Qin XM, He XF, Huang C, Yang**  
925 **ZN. 2018.** A nuclear-encoded protein, mTERF6, mediates transcription termination of *rpoA*  
926 polycistron for plastid-encoded RNA polymerase-dependent chloroplast gene expression  
927 and chloroplast development. *Scientific Reports* **8**: 11929. DOI: 10.1038/s41598-018-  
928 30166-6.
- 929 **Zhang Y, Du L, Liu A, Chen J, Wu L, Hu W, Zhang W, Kim K, Lee S-C, Yang T-J, Wang**  
930 **Y. 2016.** The Complete chloroplast genome sequences of five *Epimedium* species: lights  
931 into phylogenetic and taxonomic analyses. *Frontiers in Plant Science* **7**: 306. DOI:  
932 10.3389/fpls.2016.00306.
- 933 **Zhang J, Zhang Y, Du Y, Chen S, Tang H. 2011.** Dynamic metabonomic responses of tobacco  
934 (*Nicotiana tabacum*) plants to salt stress. *Journal of Proteome Research* **10**: 1904–1914.  
935 DOI: 10.1021/pr101140n.
- 936 **Zhao Z, Wang X, Yu Y, Yuan S, Jiang D, Zhang Y, Zhang T, Zhong W, Yuan Q, Huang L.**  
937 **2018.** Complete chloroplast genome sequences of *Dioscorea*: Characterization, genomic  
938 resources, and phylogenetic analyses. *PeerJ* **6**: e6032. DOI: 10.7717/peerj.6032.

939

940 **Figure 1.** Chloroplast genome map of *Nicotiana knightiana*, *Nicotiana rustica*, *Nicotiana*  
941 *paniculata*, *Nicotiana obtusifolia* and *Nicotiana glauca*.

942 Genes that lie outside the circle are transcribed clockwise while the genes that transcribed  
943 counterclockwise are inside the circle. Different colors indicate the genes belonging to various



944 functional groups. GC and AT content of genome are plotted light grey and dark, respectively, in  
945 the inner circle. Large single copy (LSC), inverted repeat A (IRa), inverted repeat B (IRb) and  
946 small single copy (SSC) are shown in the circular diagram. Inverted repeat regions are  
947 highlighted with *cinderella* color.

948 **Figure 2.** (A) Comparison of amino acid groups in *Nicotiana knightiana*, *Nicotiana rustica*,  
949 *Nicotiana paniculata*, *Nicotiana obtusifolia*, *Nicotiana glauca*. (B) Comparison of amino acid  
950 frequency in *Nicotiana knightiana*, *Nicotiana rustica*, *Nicotiana paniculata*, *Nicotiana*  
951 *obtusifolia*, *Nicotiana glauca*.

952 **Figure 3** Comparison of the border positions of LSC, SSC and IR among the five *Nicotiana*  
953 chloroplast genomes.

954 Positive strand transcribed genes are indicated under the line while the genes that are transcribed  
955 by negative strands are indicated above the line. Gene names are expressed in boxes, and the  
956 lengths of relative regions are showed above the boxes. The number of bp (base pairs) that are  
957 written with genes reveal the part of the genes that exists in the region of chloroplast or away  
958 from region of chloroplast i.e. bp written with *ycf1* indicate that sequences exist in that region of  
959 the plastid genome.

960 **Figure 4.** Comparison of microsatellite repeats among *Nicotiana knightiana*, *Nicotiana rustica*,  
961 *Nicotiana paniculata*, *Nicotiana obtusifolia*, *Nicotiana glauca*. (A) Indicate numbers of various  
962 types of microsatellites present in the plastid genome of *Nicotiana* species. (B) Distribution of  
963 SSRs in different regions of the plastid genome of *Nicotiana* species. (C) SSRs motifs  
964 distribution in different regions of the plastid genome of *Nicotiana* species.

965 **Figure 5.** (A). Indication of various kinds of oligonucleotide repeats exist in all *Nicotiana*  
966 *species* (B). Indicate repeats that exist range of size i.e. 30–35 indicate numbers of repeats within  
967 the size vary from 30 and 35. (C). Indicate number of repeats exist in separate areas of plastid  
968 genome. LSC: Large single copy, SSC: small single copy, IR: inverted repeat region, LSC/SSC:  
969 one copy of LSC and another in SSC, LSC/IR: one copy of LSC and another in SSC, IR/SSC:  
970 one copy of IR and another in SSC, LSC/SSC/IR: one copy of LSC, one in SSC and another in  
971 IR. (D). Indicate number of repeats in different regions of plastid genome. IGS: Intergenic spacer  
972 region, CDS: coding DNA sequences, Intron: intronic regions, IGS/Intron: one copy of  
973 intergenic spacer region and another in intronic regions. Intron/CDS: one copy intron region and

974 another in CDS regions. IGS/CDS: intergenic spacer region copy of repeat and one more in  
975 coding regions.

976 **Figure 6.** Comparison of tandem repeats among *Nicotiana knightiana*, *Nicotiana rustica*,  
977 *Nicotiana paniculata*, *Nicotiana obtusifolia*, *Nicotiana glauca*. (A) Number of tandem repeats in  
978 the chloroplast genome of *Nicotiana knightiana*, *Nicotiana rustica*, *Nicotiana paniculata*,  
979 *Nicotiana obtusifolia*, *Nicotiana glauca*. (B) Location and number of tandem repeats in the  
980 plastid genome of *Nicotiana knightiana*, *Nicotiana rustica*, *Nicotiana paniculata*, *Nicotiana*  
981 *obtusifolia*, *Nicotiana glauca*. (C) Number, size, distribution of tandem repeats across the plastid  
982 genome of *Nicotiana knightiana*, *Nicotiana rustica*, *Nicotiana paniculata*, *Nicotiana obtusifolia*,  
983 *Nicotiana glauca*

984 **Figure 7.** Nucleotide diversity of various regions of the chloroplast genome among *Nicotiana*  
985 species. The X-axis indicate the chloroplast regions and Y-axis indicate the nucleotide diversity.  
986

987 **Figure 8.** Maximum likelihood (ML) tree was reconstructed based on seventy-five protein  
988 coding plastid genes of eleven *Nicotiana* species and *Solanum dulcamara* as an outgroup.  
989 Bootstrap support values are shown above or below the nodes.

990  
991 **Figure 9.** Plastome phylogeny and biogeography of the tetraploid *Nicotiana rustica* and related  
992 species. A) Map showing the six biogeographic areas used to infer the biogeographic history of  
993 the *Nicotiana rustica* in South America. Arrows illustrate the dispersal events inferred from the  
994 biogeographic analysis. Geographical distribution for each terminal is indicated using the  
995 biogeographic regions subdivision. The most probable ancestral area is figured at each node of  
996 the phylogenetic tree. Pie-charts represent relative probabilities of ancestral states at each node.  
997 B) Node-calibrated Bayesian maximum clade credibility tree with 95% highest posterior density  
998 (HPD) interval for node ages presented as horizontal bars and mean values are displayed above  
999 each node. All nodes have PP  $\geq$  0.97 and BS  $\geq$  87%. Trace plot of the combined chains showing  
1000 the sampled joint probability and the convergence of the chains.

1001  
1002 **Table 1.** List of amino acid replacements and results of positive selection tests on codons  
1003 underlying these replacements.

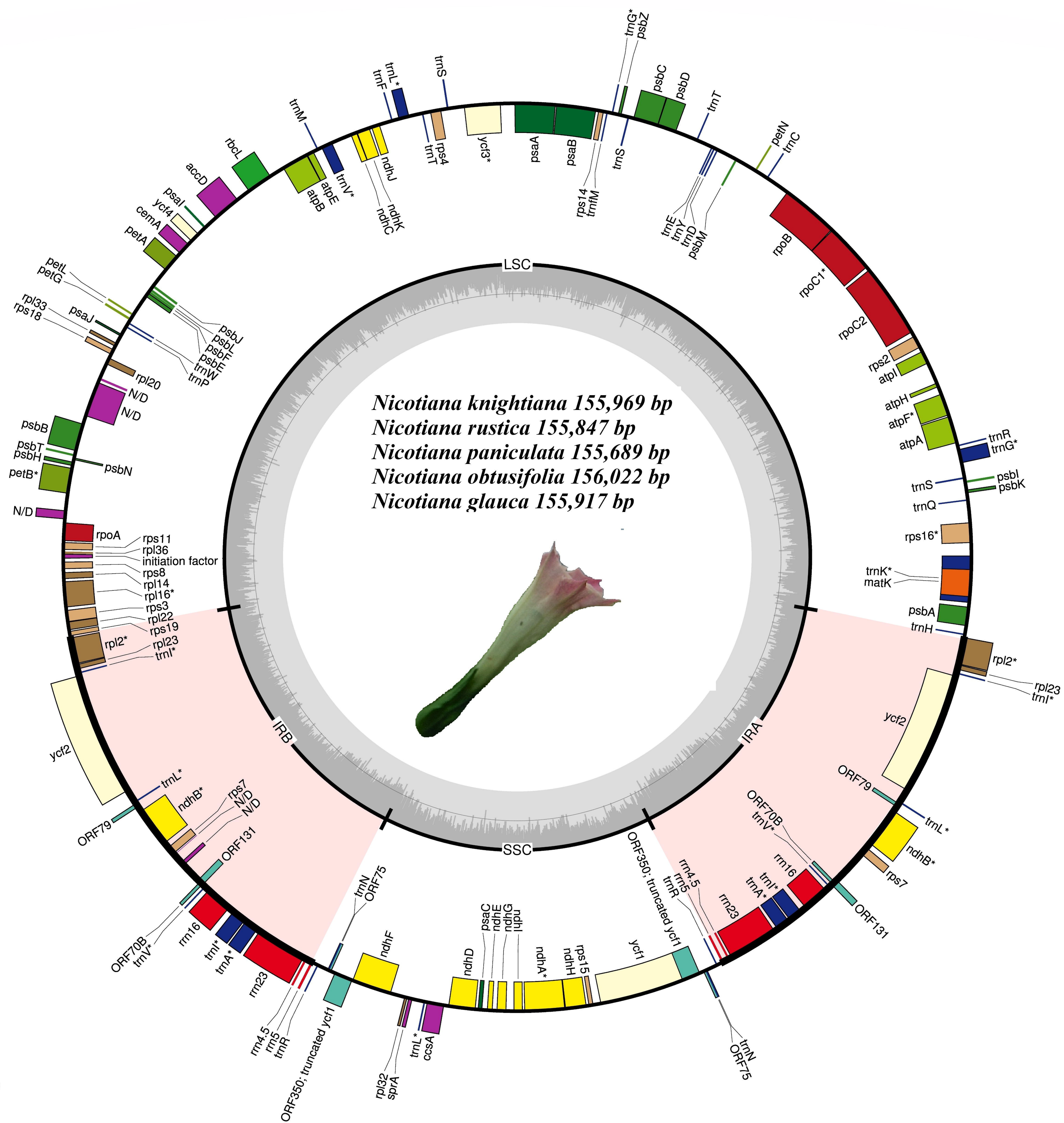
1004 **Table 2.** Comparison of substitution in *Nicotiana* species

1005 **Table 3.** Distribution of indels in *Nicotiana* chloroplast genome

1006 **Table 4.** Mutational hotspots among *Nicotiana* species

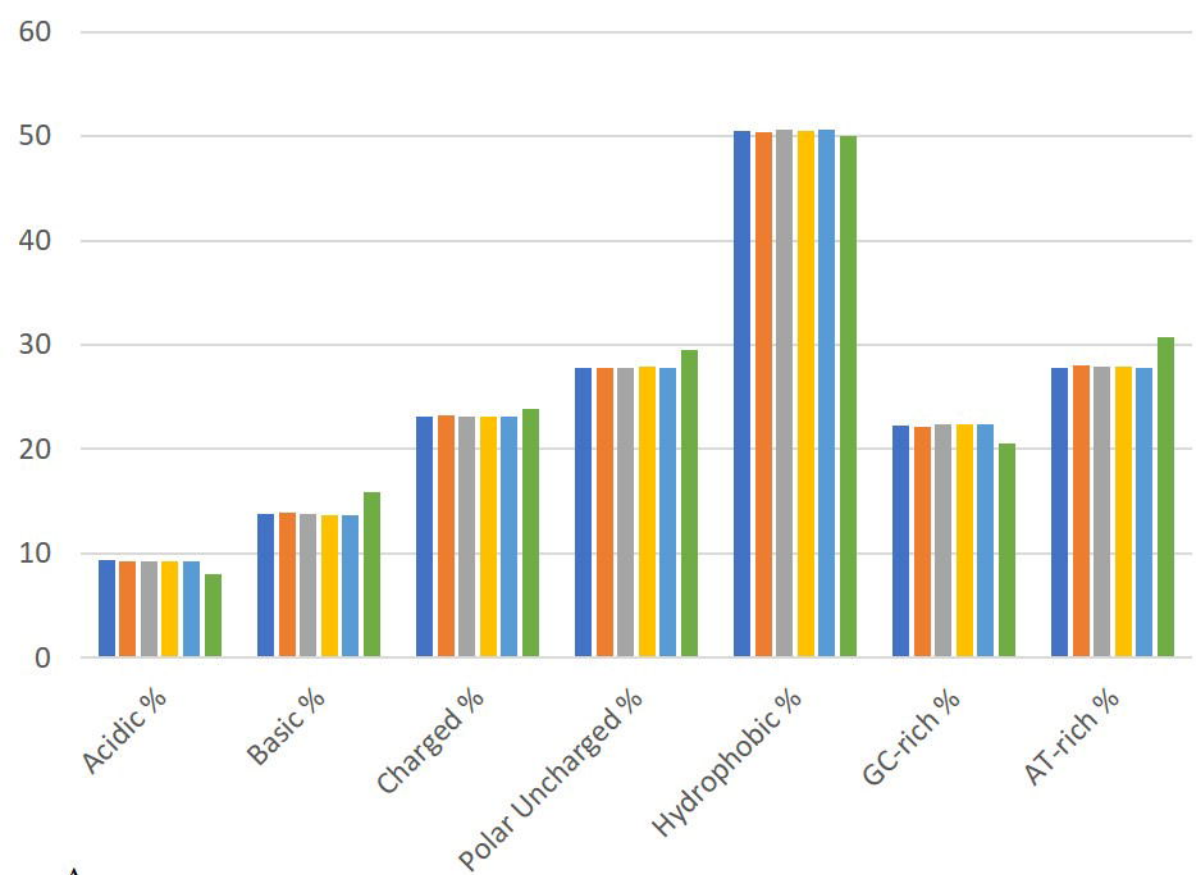
1007





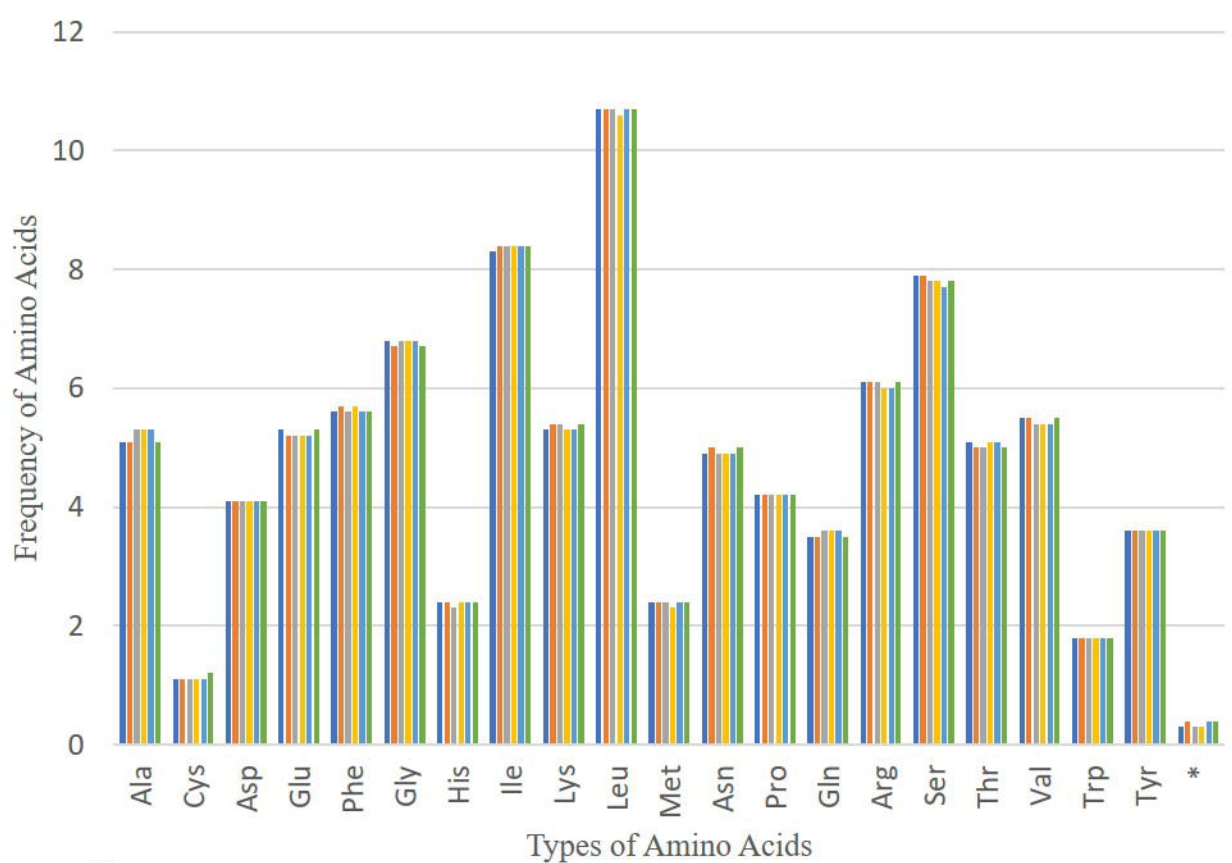
- photosystem I
- photosystem II
- cytochrome b/f complex
- ATP synthase
- NADH dehydrogenase
- RubisCO large subunit
- RNA polymerase
- ribosomal proteins (SSU)
- ribosomal proteins (LSU)
- clpP, matK
- other genes
- hypothetical chloroplast reading frames (ycf)
- ORFs
- transfer RNAs
- ribosomal RNAs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





**A**

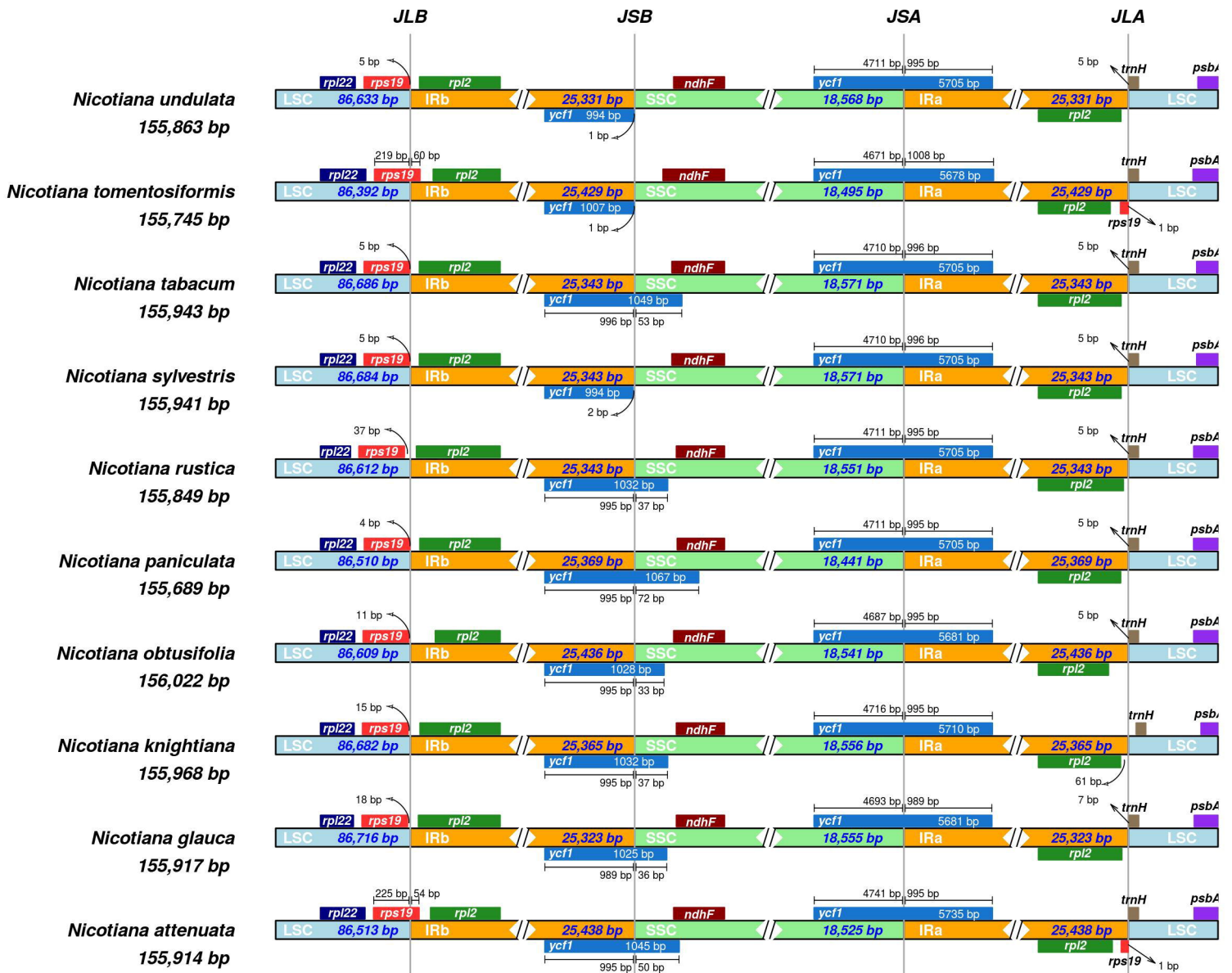
■ *Nicotiana knightiana*   
 ■ *Nicotiana rustica*   
 ■ *Nicotiana paniculata*  
■ *Nicotiana obtusifolia*   
 ■ *Nicotiana glauca*   
 ■ *Nicotiana tabacum*



**B**

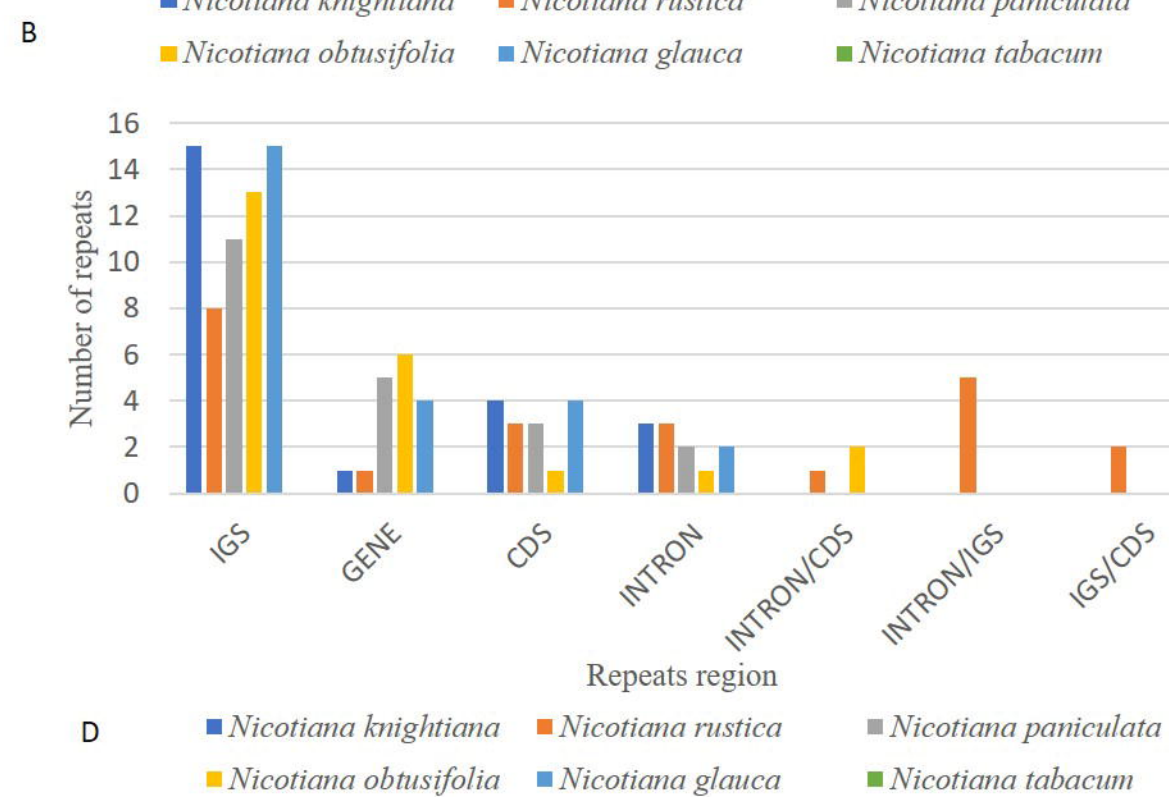
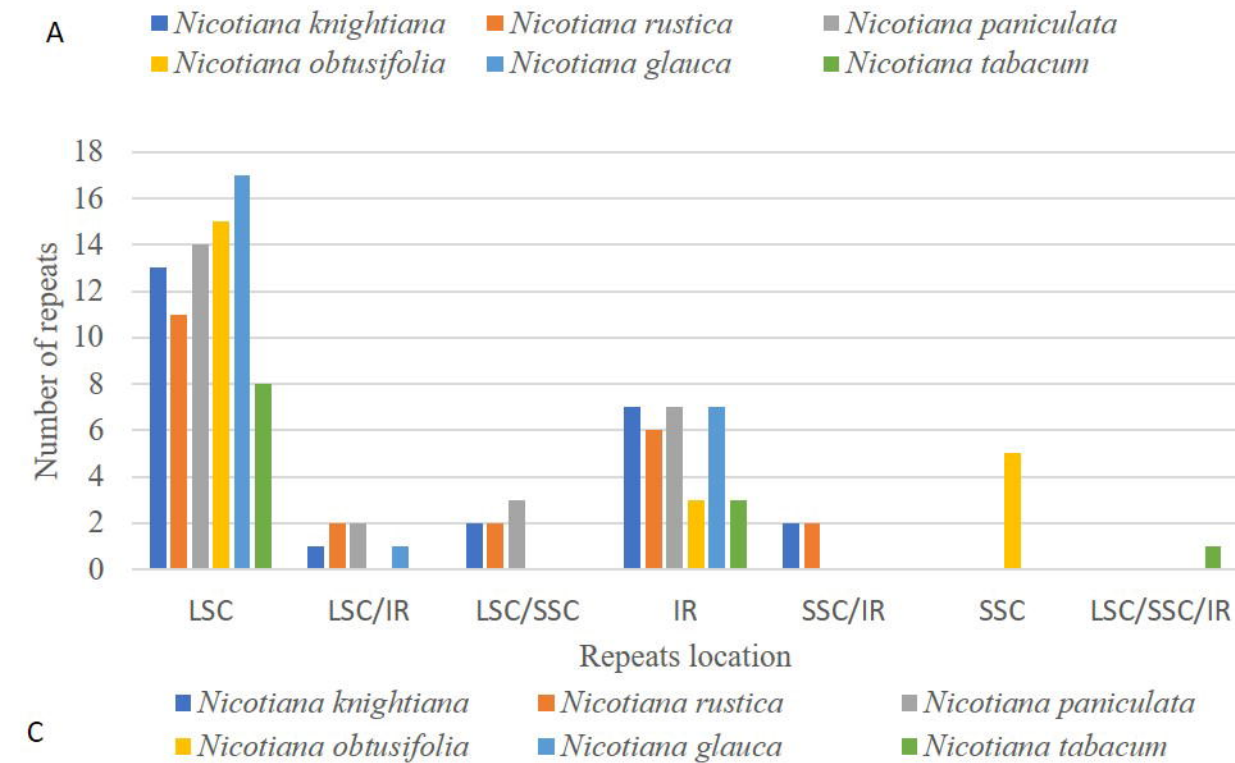
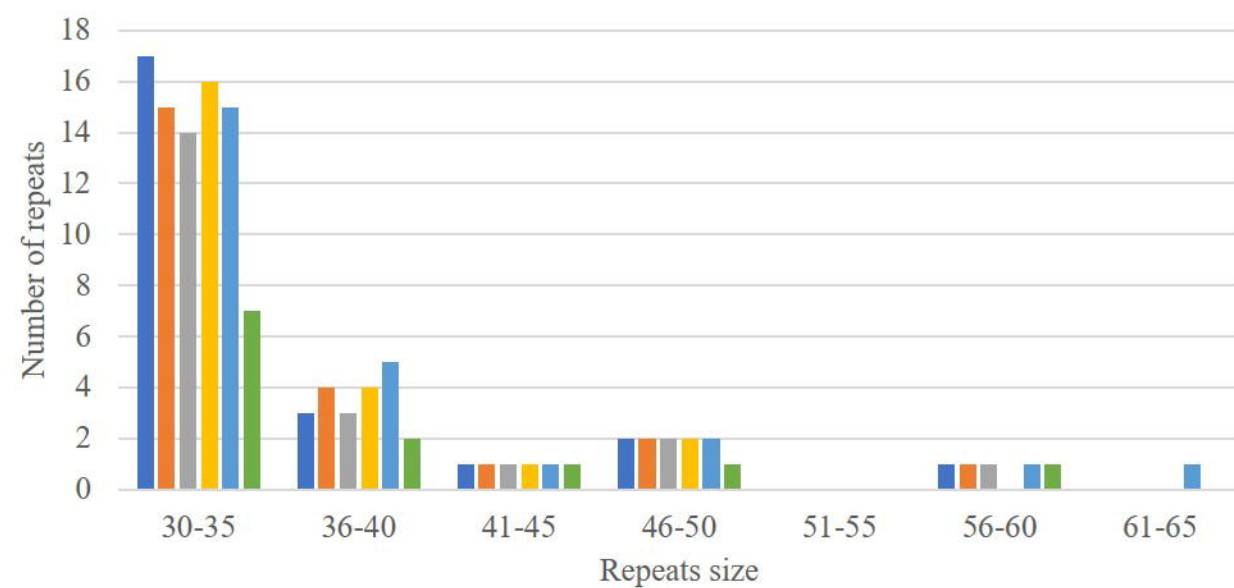
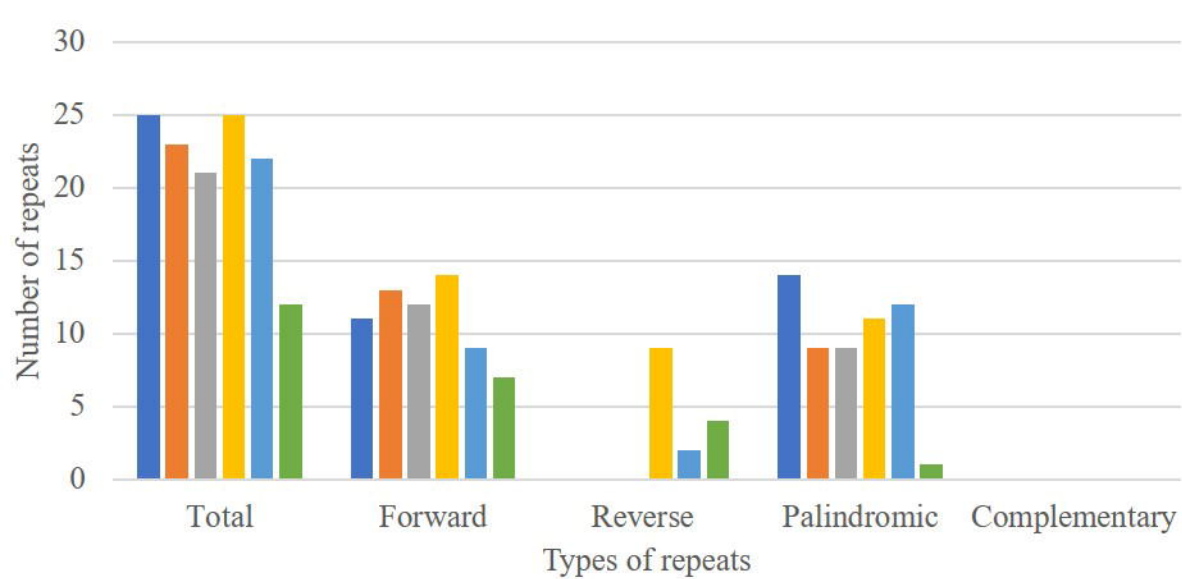
■ *Nicotiana knightiana*   
 ■ *Nicotiana rustica*   
 ■ *Nicotiana paniculata*  
■ *Nicotiana obtusifolia*   
 ■ *Nicotiana glauca*   
 ■ *Nicotiana tobaccum*

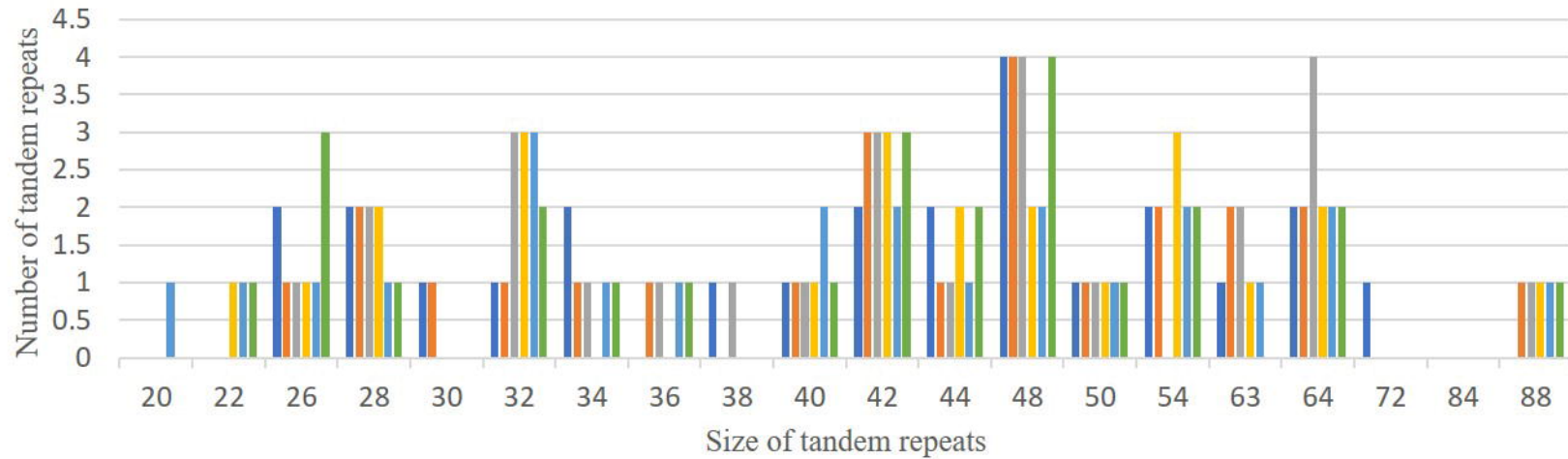
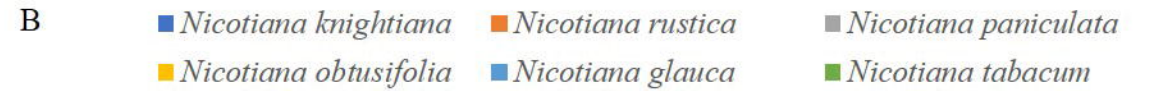
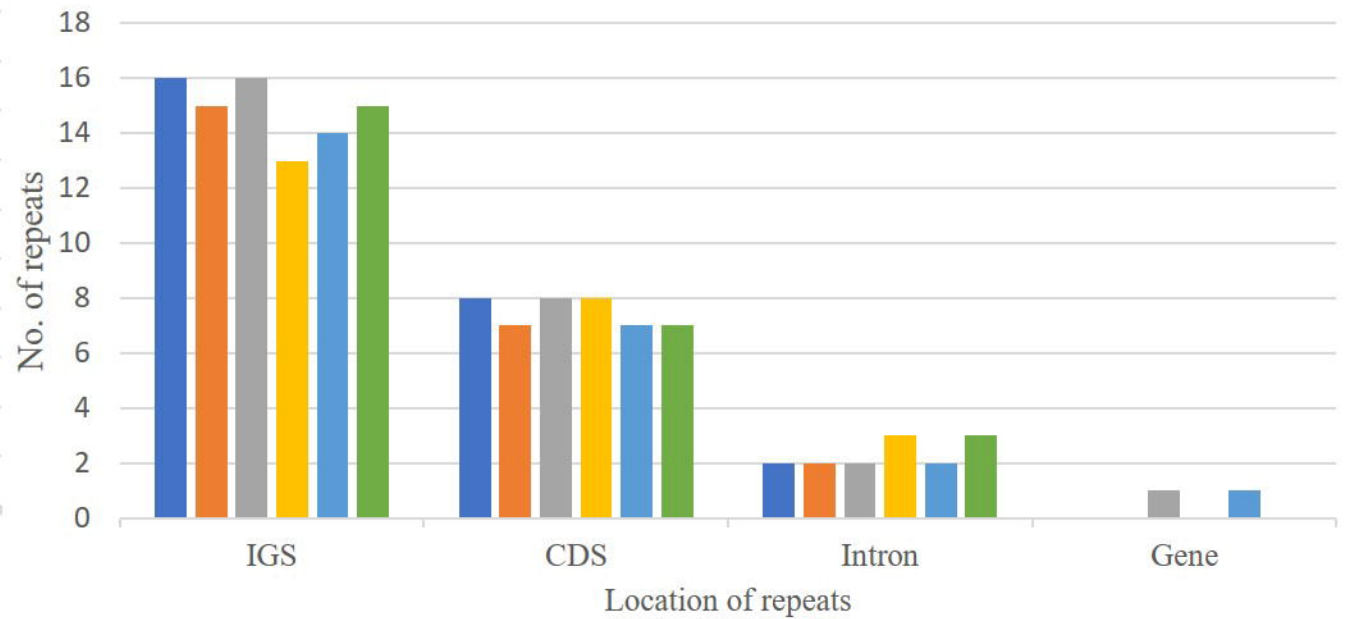
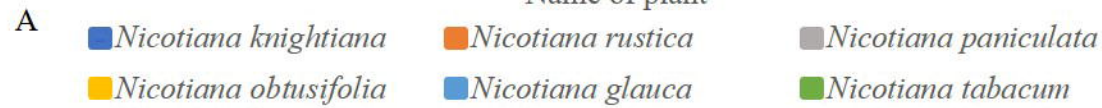
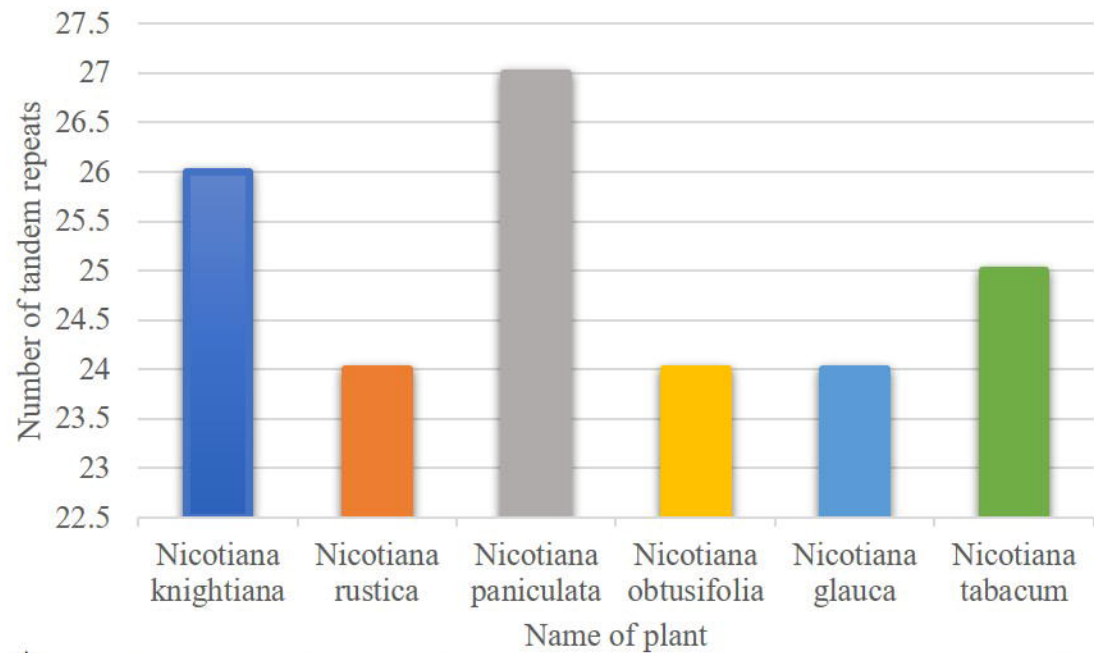
# Inverted Repeats



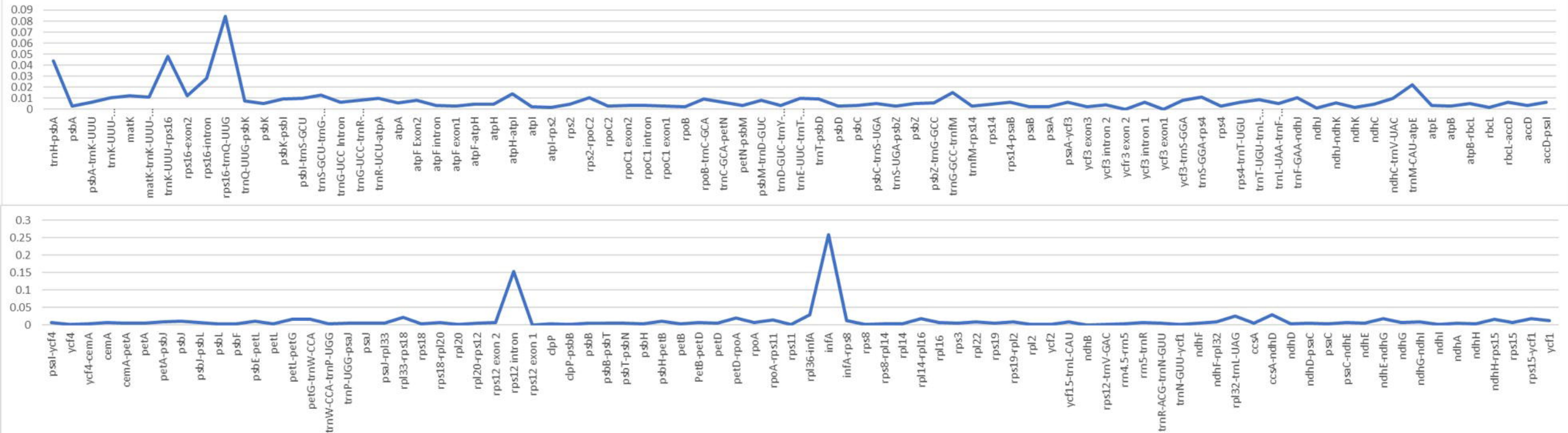


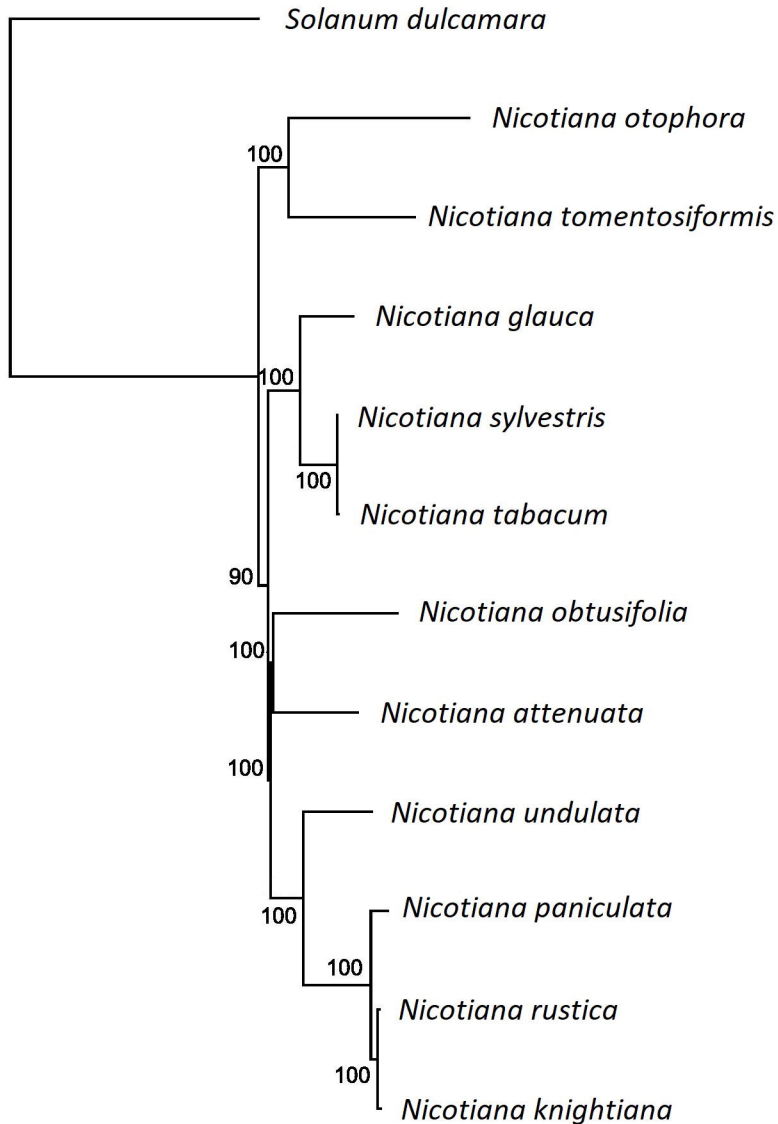






## Nucleotide Diversity





0.004

