# Identification of DNA molecular markers by comparison of *Pinus densiflora* and *Pinus sylvestris* chloroplast genomes

**Sang-Chul Kim** [Corresp., 1] , **Jei-Wan Lee** [Corresp., 1] , **Seung-Hoon Baek** [1] , **Ji-Young Ahn** [1] , **Kyung-Nak Hong** [1]

[1] National Institute of Forest Science, Division of Forest Genetic Resources, Suwon, Republic of Korea

Corresponding Authors: Sang-Chul Kim, Jei-Wan Lee
Email address: majin01@korea.kr, leejeiwan@korea.kr

**Background:** Identifying and characterizing genetic variation can clarify the molecular basis of biological phenomena in plants. In particular, related or morphologically similar species can be distinguished by molecular markers. *Pinus densiflora* Siebold & Zucc. is a species that is distributed in the Korean peninsula, the Japanese archipelago, and China's Shandong and Manchu Provinces and has long been harvested for timber. However, it is difficult to distinguish *P. densiflora* from *Pinus sylvestris* L. both morphologically and phylogenetically. The complete chloroplast genome of *P. densiflora* has not yet been reported. In this study, we sequenced the *P. densiflora* chloroplast genome in order to identify the molecular markers that can be used to distinguish this species from *P. sylvestris.*

**Methods:** Genomic DNA was extracted from *P. densiflora* samples obtained from the clone bank of the National Forest Seed Variety Center and was sequenced on an Ion Torrent platform. Filtered sequences were assembled with *P. sylvestris* sequences used as a reference and gene annotation was performed. The chloroplast genome sequences of the two species were aligned and the number and location of forward, reverse, complement and palindromic matches were determined. Single nucleotide polymorphisms (SNPs) and insertion/deletion mutations (Indels) were identified and analyzed by PCR.

**Results:** The *P. densiflora* chloroplast genome consisted of circular double-stranded DNA with 119,835 bp compared to 119,758 bp for *P. sylvestris*. Between the two *Pinus* chloroplast genomes, we identified 73 SNPs and 171 Indels; two gene regions with amplification products ≤ 300 bp (*rpoC1* and *trnM-trnV*) were validated as molecular markers.

**Discussion:** PCR restriction fragment length polymorphism analysis revealed differences between *P. sylvestris* and *P. densiflora* at the molecular level. These differences can be used to distinguish between these two species, which is not possible by microscopy-based morphological examination.

## 1 Introduction

2   Pinaceae, comprising 11 genera and more than 200 species, is the largest extant family of

3   gymnosperms. Many species of the pine family constitute the major forest elements in the

4   northern temperate region (Wang, Tank & Sang. 2000). *Pinus* L. is one of 11 genera in

5   Pinaceae, a monophyletic family among gymnosperms (Farjon, 2010). Approximately 110

6   species comprise 50% of Pinaceae, making it the largest genus of existing gymnosperms

7   (Syring *et al*., 2005), most of which are distributed in the temperate zone of the Northern

8   Hemisphere. *Pinus* is divided into two subgenera *Strobus* and *Pinus*, according to the number of

9   fibrovascular bundles in the needle (Geada López, Kamiya & Harada, 2002). Approximately 20

10  species are native to or are cultivated in Korea (Korea National Arboretum and The Plant

11  Taxonomic Society of Korea, 2007; Hong *et al*., 2014). *Pinus densiflora* Siebold & Zucc. is

12  distributed throughout Korea and is one of the most economically important species sustaining

13  forest ecosystems and is harvested for wood and fuel (Lee *et al*., 2004; Kim, Kim & Lim, 2017).

14  *Pinus sylvestris* L. is the most abundant species in Europe and is found from Scotland and Spain

15  to Siberia and northern Asia.

16   Chloroplasts a type of plastid in plants and algae, are intracellular organelles that carry out

17  photosynthesis (Howe *et al*., 2003). They are presumed to have originated from an

18  endosymbiotic event between cyanobacteria and non-photosynthetic host cells (Dyall, Brown &

19  Johnson, 2004). Plastid genomes are stable in terms of structure, gene content, and gene order

20  across land plants (Jansen *et al*., 2005). The chloroplast genome of higher plants consists of a

21  circular double strand ranging from 120 to 210 kb that usually contains two inverted repeat (IR)

22  regions (IRA and IRB) separated by large and small single-copy regions (LSC and SSC,

23  respectively) (Ravi *et al*., 2008). Most plant genomes have 66–82 protein-coding genes, 29–32

24  genes encoding tRNAs, and four genes encoding rRNAs, With the exception of non-

25  photosynthetic parasitic plants,  gene composition, sequence, content, and orientation are highly

26  conserved among seed plants (Jansen & Ruhlman, 2012). However, structural modifications

27  such as loss of IR domains or entire genes and gene rearrangement have been reported in

28  gymnosperms such as conifers (Lin *et al*., 2010; Wu *et al*., 2011; Wu & Chaw, 2014; Yi *et al*.,

29  2016). The first complete sequences of *Pinus* cpDNA were reported in *Pinus thunbergii* Parl.,

30  with 4 rRNA genes and 32 tRNA genes, and the most striking feature is the loss of all 11

31  functional genes (*ndh* genes) for in subunits of a putative NADH dehydrogenase that are found

32  in the chloroplast genomes of angliosperms and a bryophyte (Wakasugi *et al*., 1994). There are

33  currently; 2,245 complete chloroplast genomes of seed plants in the National Center for

34  Biotechnology Information (NCBI) Organelle Genome Resources database

35    (http://www.ncbi.nlm.nih.gov/genomes/).

36    Plastid genome sequences are widely used for DNA barcoding, species conservation, genomic

37    evolution, and molecular phylogenetic studies (Moore *et al*., 2007). Identifying and

38    characterizing genetic variation can clarify the molecular basis the of biological phenomena in

39    plants (Agarwal, Shrivastava & Padh, 2008) and provide insight into the mechanisms of

40    evolution and natural selection. In particular, species that are difficult to differentiate

41    morphologically can be distinguished using molecular markers. The complete chloroplast

42    genome can be rapidly sequenced at a relatively low cost (Yi *et al*., 2016). Also, PCR restriction

43    fragment length polymorphism (RFLP) analysis, also known as cleaved amplified polymorphic

44    sequence (CAPS), is widely used to detect intra- and interspecies variation (Rasmussen, 2012).

45    *P. sylvestris* and *P. densiflora* belong to the subgenus *Pinus*, section *Pinus*, subsection *Pinus*.

46    These trees are characterized by the shedding of bud-scales along with the leaves and by two

47    cross-sectional vascular bundles in the leaves (Lee, 2003). Futher, molecular phylogenetic

48    studies show that *P. sylvestris* and *P. densiflora* formseparate strongly supported groups, with

49    common morphological features, including irregular cracking of 2-year-old bark (Wang *et al*.,

50    1999; Gernandt, 2005; Hong *et al*., 2014). Thus, *P. sylvestris* and *P. densiflora* are difficult to

51    distinguish morphologically and phylogenetically; as such, the timber of the two species is often

52    combined or illegally substituted. A complete chloroplast genome is available for *P. sylvestris* but

53    not for *P. densiflora.* In this study, we sequenced the chloroplast genome of *P. densiflora* and

54    compared it with that of *P. sylvestris* in order to identify polymorphisms that can serve as

55    molecular markers to distinguish between the two species by PCR-RFLP analysis.


56    **Materials and Methods**


57    *Sample collection, DNA extraction, and sequencing*

58    *P. densiflora* samples were obtained from the clone bank of the National Forest Seed Variety

59    Center (Anmyeondo, Korea; elite tree: Gyeongbuk No. 4) and genomic DNA was isolated from

60    fresh leaves using the Plasmid SV mini kit (GeneAll, Seoul, Korea). DNA samples from plants

61    used in this study are now stored in the DNA Bank of the Forest Genetic Resources Department

62    of the National Institute of Forest Science. Total genomic DNA was extracted from 10 g of fresh

63    leaves using a Plasmid SV mini kit. Whole genome sequencing was performed on the Ion

64    Torrent platform (Life Technologies, Carlsbad, CA, USA). Libraries were sequenced on Ion

65    Proton using the Ion PI Chip kit v3 deposited at full density according to the protocol for 200 bp

66    sequencing supplied by the manufacturer.


67    *Chloroplast genome assembly and annotation*

68    Filtered sequences were assembled using Bowtie2 v. 2.2.3 software (http://bowtie-

69    bio.sourceforge.net/bowtie2/index.shtml; Langmead & Salzberg, 2012) with *P. sylvestris*

70    sequence (GenBank: NC035069) as a reference. In total, 1,449,103 reads were mapped to the

71    reference sequence with an average coverage of 851.1X. Finally, the contigs were assembled

72    using Geneious 10.2.3 (Biomatters, Auckland, New Zealand; Kearse *et al.*, 2012). Gene

73    annotation was performed using the Basic Local Alignment Search Tool (BLAST and BLASTX)

74    available on the NCBI website. All tRNA sequences were confirmed using the web-based tool,

75    tRNAScan-SE (Schattner, Brooks & Lowe, 2005) with default settings to corroborate tRNA

76    boundaries identified by Geneious. Genome maps were generated using

77    OrganellarGenomeDRAW (Lohse, Drechsel & Bock, 2007), followed by manual modification.

78    *Comparison of Pinus chloroplast genome sequences*

79    Simple sequence repeats (SSRs) were analyzed using Phobos v. 3.3.12 (Mayer, 2010), with

80    thresholds of eight repeat units for mononucleotide SSRs, four for di- and trinucleotide SSRs,

81    three  for tetra- and pentanucleotide SSRs, and two for hexanucleotide SSRs. All detected

82    repeats were manually verified, and redundant results were removed. We aligned the plastid

83    genome sequences of the two *Pinus* species using MAFFT (Katoh *et al.*, 2002). For long repeat

84    sequences, the REPuter program was used to assess the number and location of forward,

85    reverse, complement and palindromic matches (Kurtz *et al.*, 2001). Repeat identity and size

86    were limited to > 90% and ≥ 25 bp, respectively.

87    *Identification of molecular markers to distinguish between Pinus species*

88    Single nucleotide polymorphisms (SNPs) and insertion/deletion mutations (Indels) were

89    identified using Geneious 10.2.3 and analyzed by PCR. Primer3 software was used to design

90    primers ranging in size from 18 to 22 mer. The temperature ranged from 57 °C to 63 °C. There

91    was one GC clamp, and primer amplification products ranged from 250 to 350 bp. Each thirty

92    individuals of both *Pinus* species were tested for the species-specific DNA markers. The

93    individuals of *P. densiflora* were sampled from thirteen populations in South Korea (Ahn *et al.*,

94    2015), and the individuals of *P. sylvestris* were sampled from the 22 provenance in Sweden,

95    which were introduced into Korea in order to select superior provenances that are well adapted

96    to Korean environment (Ryu *et al.*, 2013). It is also stored in the DNA Bank of Forest Genetic

97    Resources Department (NIFS_122059323 to 122059343).  PCR reaction mixtures contained 10

98    pmol of each primers pairs, 25 ng total DNA, 0.5 µl of 10 mM dNTPs, 2.5 µl of 10× reaction

99    buffer (2.5 mM $MgCl_2$, 20 mM Tris-HCl [pH 8.4], 50 mM KCl), and 1 U Taq DNA polymerase

100   (BioFACT, Daejeon, Korea). Reactions were performed on a GeneAmp PCR System 9700

101   thermal cycler (Applied Biosystems, Foster City, CA, USA) under the following conditions:  94 °C

102    for 5 min; 45 cycles of 94 °C for 30 s, 60 °C for 30 s, and 72 °C for 30 s; and 72 °C for 10 min.

103    PCR products were confirmed on a 2 % agarose gel, and were digested using 5 U of *Hinf*I,

104    *BsaW*I and *Sph*I followed by electrophoresis on a 2 % acrylamide gel. DNA fragment sizes were

105    estimated by comparison with 100 bp Plus Ladder (Thermo Fisher Scientific, Waltham, MA,

106    USA).


107    **Results**


108    *Chloroplast genome assembly and features*

109     The complete chloroplast genome of *P. densiflora* was determined to be a circular double-

110    stranded DNA sequence of 119,835 bp (GenBank accession number: MF990371). The genome

111    showed a typical quadripartite structure including LSC (65,896 bp) and SSC (53,219 bp) regions

112    and IRs (360 bp) (Fig. 1). The genome had a similar GC content to that of *P. sylvestris*. The GC

113    content was the highest in the SSC region (39.4%), moderate in the LSC region (37.8%), and

114    the lowest in the IR region (34.7%). We identified 113 genes including 74 encoding proteins and

115    36 and four encoding tRNA and rRNA, respectively. Twelve genes (six protein-coding and six

116    tRNA genes) contained one intron, while two of the protein-coding genes (*ycf3* and *rps12*) had

117    two introns. We also confirmed that *trnS*-GCU and *psaM* were duplicated in two chloroplasts

118    (Table 1).


119    *Analyses of repetitive sequences*

120     We detected SSRs > 8 bp in *P. densiflora* and *P. sylvestris* chloroplast genomes according to a

121    previously published method (Qian *et al*., 2013). We set the threshold based on the fact that

122    SSRs > 8 bp are prone to strand slippage and mispairing, which is thought to be the primary

123    mechanism underlying the high rate of polymorphism. In our analysis, there were 103 and 106

124    SSRs accounting for 1,236 bp in *P. densiflora* and 1,254 bp in *P. sylvestris*, respectively. These

125    included 18 mono-, five di-, one tri-, four tetra-, and 75 hexa-nucleotide repeats for *P. densiflora*

126    and 21 mono-, seven di-, one tri-, four tetra-, and 73 hexanucleotide repeats for *P. sylvestris*.

127    Hexanucleotide repeats accounted for 72.8% and 68.9% of total SSRs in *P. densiflora* and *P.*

128    *sylvestris*, respectively. The majority of mononucleotide SSRs were thymine and adenine. The

129    majority of the identified repeats were located in the non-coding regions (intergenic spacers and

130    introns) except 17 protein-coding genes (*matK, atpH, atpI, rpoC2, rpoB, petL, psbL, petA, rbcL,*

131    *atpB, rpl22, psbC, rrn16, rrn23, ycf1, rpl32,* and *ycf2*). In total, 35 long repeat sequences > 25 bp

132    were identified in the *P. densiflora* chloroplast genome, including 21 forward and 14 palindromic

133    matches. Thirty long repeat sequences were identified in *P. sylvestris*, including 16 forward, one

134    reverse, and 13 palindromic matches (Fig. 2).

135    *Comparison of indels and SNPs in Pinus species*

136    In total, 171 indels were found to differ between *P. densiflora* and *P. sylvestris*, most of which

137    were located in intergenic spacer regions (91.2%), with 69.6% and 30.4% in the LSC and SSC,

138    respectively. The average Indel length was 6 bp, and the longest was located in *cemA-ycf4* and

139    *trnE-clpP*. The frequency of 1 bp Indels was 36.7%, while 30% was > 8 bp. In addition, Indels

140    were detected in two coding genes of both species (*psaM* and *ycf2*; Table 2). Seventy-three

141    SNPs differed between *P. densiflora* and *P. sylvestris*, of which 46 were transversions (63%). In

142    total, 34 (46.6%) SNPs were located in coding regions, whereas 39 (53.4%) were in intergenic

143    spacer regions or introns (Table 3).

144    The *rpoC1* and *trnM-trnV* gene regions were amplified by PCR (using the primers shown in

145    Table 4) in order to validate their capacity to distinguish between the two *Pinus* species. The

146    amplification product obtained using the Pdest-cp1 primer was digested with *Hinf*I and visualized

147    by agarose gel electrophoresis. Approximately 200 bp fragment was observed in *P. densiflora*

148    but not in *P. sylvestris* (Fig. 3). On the other hand, digestion of the Pidest-cp2 amplification

149    product with *BsaW*I yielded approximately 200 bp fragment that was observed in *P. sylvestris* but

150    not in *P. densiflora* (Fig. 4). In addition, digesting the Pidest-cp2 amplification product with *Sph*I

151    produced a fragment of approximately 200 bp that was observed in *P. densiflora* but not in *P.

152    sylvestris* (Fig. 5).


153    **Discussion**


154    *Comparison of the complete plastid genomes of P. densiflora and P. sylvestris*

155    IRs are known to stabilize the plastid genome because of its low base exchange rate and high

156    copy-correcting activity. Thus, the loss of IRs can result in the shortening of intergenic spaces,

157    gene loss, and structural variations in plastids. Reductions in Irs have been observed in most

158    gymnosperm and in some legumes. In *P. densiflora*, the total chloroplast genome was 77 bp

159    longer than that of *P. sylvestris,* while the gene content, order, and orientation were similar to

160    those of ther *Pinus* chloroplast genomes (Wakasugi *et al*.,1994; Duan *et al*., 2016; Fang *et al*.,

161    2016; Celiński *et al*., 2017; Ni *et al*., 2017).

162    We also found that SSRs of 1–6 bp per unit—known as, microsatellites— were distributed

163    throughout the *P. densiflora* chloroplast genome. SSRs are important molecular markers of

164    genomic variation within species or populations due to their high polymorphism, and have been

165    extensively used to analyze plant population structure, diversity, differentiation and fertility (Kim

166    & Kim, 2016). The SSRs detected in the present study will provide basic information for future

167    analyses of genetic diversity in Pinaceae.


168    *Identification molecular markers for distinguish between Pinus species based on SNPs*

169     *P. sylvestris* and *P. densiflora* are difficult to distinguish because they form separate strongly

170    supported groups in molecular phylogenetic studies (Wang *et al*., 1999; Gernandt, 2005; Hong

171    *et al*., 2014) and also have very similar morphology (Lee, 2003; Hong *et al*., 2014).

172     Complete chloroplast genome sequences (plastomes) have been very useful for understanding

173    phylogenetic relationships in angiosperms at the family level and above and have been used to

174    resolve previously recalcitrant nodes (Barrett *et al.*, 2016).

175     In order to identify molecular markers that can be used to distinguish the two *Pinus* species, we

176    compared their chloroplast genomes and found a total of 171 indels and 73 SNPs. We amplified

177    the *rpoC1* and *trnM-trnV* gene regions by PCR-RFLP and found that SNPs in these two regions

178    could be clearly distinguished by restriction enzyme digestion.

179     The differences between individuals can be clearly detected by separating differently sized

180    fragments based on the single nucleotide and Indel polymorphisms of the restriction enzyme

181    sites rather than by using dominant markers (Lee *et al*., 2012). In this context, the co-dominant

182    CAPS markers developed in this study provide a means of unambiguously identifying *P.*

183    *densiflora* and *P. sylvestris*.

184     It is likely that the SNP loci found in this study exist in other *Pinus* species. These markers have

185    many research and commercial applications, including studies of genetic diversity, breeding, and

186    species identification for the timber market.


187    **Conclusions**

188

189     This study provides the complete chloroplast sequences of *P. densiflora*. These sequences

190    revealed significant similarity in the structural organization of the chloroplast genomes in

191    Pinaceae, such as loss of IR and reduction of *ndh* genes, i.e., all *ndh* genes were transferred to

192    the nucleus or that NADH dehydrogenase is not essential in pine chloroplasts. In addition,

193    molecular markers that can distinguish between the phylogenetically and morphologically similar

194    *P. sylvestris* and *P. densiflora* were identified, providing a more objective and reliable method of

195    identification than that by conventional visual identification methods. Further, the data generated

196    here can be used to develop additional molecular markers for comparing with other *Pinus*

197    species. These markers can also have research and commercial applications such as in genetic

198    diversity studies, breeding, and identification of species for the timber market.

199     Further research is necessary to determine whether the restriction site differences occur across

200    the entire geographic range of both species, given that additional mutations may have caused

201    one or more of the restriction sites to disappear in some populations.

## References

Agarwal M, Shrivastava N, Padh H. 2008. Advances in molecular marker techniques and their applications in plant sciences. *Plant Cell Reports*, 27: 617–631. DOI: 10.1007/s00299-008-0507-z

Ahn, J. Y., Hong, K. N., Lee, J. W., Hong, Y. P., & Kang, H. 2015. Genetic Variation of *Pinus densiflora* Populations in South Korea Based on ESTP Markers. *Korean Journal of Plant Resources*, 28(2), 279-289.

Celiński K, Kijak H, Barylski J, Grabsztunowicz M, Wojnicka-Półtorak A, Chudzińska E. 2017. Characterization of the complete chloroplast genome of *Pinus uliginosa* (Neumann) from the *Pinus mugo* complex. *Conservation Genetics Resources*, 9: 209–212. DOI: 10.1007/s12686-016-0652-6.

Duan RY, Yang LM, Lv T, Wu GL, Huang MY. 2016. The complete chloroplast genome sequence of *Pinus dabeshanensis*. *Conservation Genetics Resources*, 8: 395–397. DOI: 10.1007/s12686-016-0567-2

Dyall SD, Brown MT, Johnson PJ. 2004. Ancient invasions: from endosymbionts to organelles. *Science*, 304: 253–257. DOI: 10.1126/science.1094884

Farjon, A. (2010). *A Handbook of the World's Conifers (2 vols.)*(Vol. 1). Brill.

Fang MF, Wang YJ, Zu YM, Dong WL, Wang RN, Deng TT, Li ZH. 2016. The complete chloroplast genome of the Taiwan red pine *Pinus taiwanensis* (Pinaceae). *Mitochondrial DNA Part A*, 27: 2732–2733. DOI: 10.3109/19401736.2015.1046169

Geada López G, Kamiya K, Harada K. 2002. Phylogenetic relationships of Diploxylon pines (subgenus *Pinus*) based on plastid sequence data. *International Journal of Plant Sciences*, 163: 737–747. DOI: 10.1086/342213

Gernandt, D. S., G. Lopez, S. O. Garcia and A. Liston. 2005. Phylogeny and classification of *Pinus*. *Taxon* 54: 29-42.

Hong JK, Yang JC, Lee YM, Kim JH. 2014. Molecular phylogenetic study of *Pinus* in Korea based on chloroplast DNA *psbA-trnH* and *atpF-H* sequences data. *Korean Journal of Plant Taxonomy*, 44: 111–118. DOI: 10.11110/kjpt.2014.44.2.111

Howe CJ, Barbrook AC, Koumandou VL, Nisbet RER, Symington HA, Wightman TF. 2003. Evolution of the chloroplast genome. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 358: 99-107. DOI: 10.1098/rstb.2002.1176

Jansen RK, Ruhlman TA. 2012. Plastid genomes of seed plants. In: Bock R, Knoop V, eds. *Genomics of chloroplasts and mitochondria*. Dordrecht: Springer Netherlands. 103–126. DOI: 10.1007/978-94-007-2920-9_5

Jansen RK, Raubeson LA, Boore JL, Chumley TW, Haberle RC, Wyman SK, Kuehl JV. 2005.

237      Methods for obtaining and analyzing whole chloroplast genome sequences. *Methods in*

238      *Enzymology* 395: 348–384. DOI: 10.1016/S0076-6879(05)95020-9

239   Katoh K, Misawa K, Kuma KI, Miyata T. 2002. MAFFT: a novel method for rapid multiple

240      sequence alignment based on fast Fourier transform. *Nucleic Acids Research* 30: 3059–

241      3066. DOI: 10.1093/nar/gkf436

242   Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Thierer T. 2012. Geneious

243      Basic: an integrated and extendable desktop software platform for the organization and

244      analysis of sequence data. *Bioinformatics* 28: 1647–1649. DOI:

245      10.1093/bioinformatics/bts199

246   Kim JB, Kim ES, Lim JH. 2017. Topographic and meteorological characteristics of *Pinus*

247      *densiflora* dieback areas in Sogwang-Ri, Uljin. *Korean Journal of Agricultural and Forest*

248      *Meteorology* 19: 10–18. DOI: 10.5532/KJAFM.2017.19.1.10

249   Kim SC, Kim JS, Kim JH. 2016. Insight into infrageneric circumscription through complete

250      chloroplast genome sequences of two *Trillium* species. *AoB Plants* 8. DOI:

251      10.1093/aobpla/plw015

252   Korea National Arboretum and The Plant Taxonomic Society of Korea. 2007. A Synonymic List of

253      Vascular Plants in Korea. Korea National Arboretum. Pocheon. (in Korean)

254   Kurtz S, Choudhuri JV, Ohlebusch E, Schleiermacher C, Stoye J, Giegerich R. 2001. REPuter:

255      the manifold applications of repeat analysis on a genomic scale. *Nucleic Acids Research* 29:

256      4633–4642. DOI: 10.1093/nar/29.22.4633

257   Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods* 9:

258      357–359. DOI: 10.1038/nmeth.1923

259   Lee CS, Kim JH, Yi H, You YH. 2004. Seedling establishment and regeneration of Korean red

260      pine (*Pinus densiflora* S. et Z.) forests in Korea in relation to soil moisture. *Forest Ecology*

261      *and Management* 199: 423–432. DOI: 10.1016/j.foreco.2004.05.053

262   Lee JW, Bang KH, Kim YC, Seo AY, Jo IH, Lee JH, Cho JH. 2012. CAPS markers using

263      mitochondrial consensus primers for molecular identification of *Panax* species and Korean

264      ginseng cultivars (*Panax ginseng* CA Meyer). *Molecular Biology Reports* 39: 729–736. DOI:

265      10.1007/s11033-011-0792-4

266   Lee TB. 2003. Coloured Flora of Korea. Hyangmunsa, Seoul (2003) (in Korean)

267   Lin CP, Huang JP, Wu CS, Hsu CY, Chaw SM. 2010. Comparative chloroplast genomics reveals

268      the evolution of Pinaceae genera and subfamilies. *Genome Biology and Evolution* 2: 504–

269      517. DOI: 10.1093/gbe/evq036

270   Lohse M, Drechsel O, Bock R. 2007. OrganellarGenomeDRAW (OGDRAW): a tool for the easy

271      generation of high-quality custom graphical maps of plastid and mitochondrial genomes.

272      *Current Genetics* 52: 267–274. DOI: 10.1007/s00294-007-0161-y

273   Mayer C. 2010. Phobos Version 3.3.12. A tandem repeat search program, 20.

274   Moore MJ, Bell CD, Soltis PS, Soltis DE. 2007. Using plastid genome-scale data to resolve

275       enigmatic relationships among basal angiosperms. *Proceedings of the National Academy of*

276       *Sciences of the United States of America* 104: 19363–19368. DOI:

277       10.1073/pnas.0708072104

278   Ni Z, Ye Y, Bai T, Xu M, Xu LA. 2017. Complete chloroplast genome of *Pinus massoniana*

279       (Pinaceae): Gene rearrangements, loss of ndh genes, and short inverted repeats contraction,

280       expansion. *Molecules* 22: 1528. DOI: 10.3390/molecules22091528

281   Qian J, Song J, Gao H, Zhu Y, Xu J, Pang X, Liu J. 2013. The complete chloroplast genome

282       sequence of the medicinal plant *Salvia miltiorrhiza. PLoS One* 8: e57607. DOI:

283       10.1371/journal.pone.0057607

284   Rasmussen HB. 2012. Restriction fragment length polymorphism analysis of PCR-amplified

285       fragments (PCR-RFLP) and gel electrophoresis-valuable tool for genotyping and genetic

286       fingerprinting. In: Magdeldin S, ed. *Gel Electrophoresis – Principles and Basics.* InTech. DOI:

287       10.5772/37724

288   Ravi V, Khurana JP, Tyagi AK, Khurana P. 2008. An update on chloroplast genomes. *Plant*

289       *Systematics and Evolution* 271: 101–122. DOI: 10.1007/s00606-007-0608-0

290   Ryu, K. O., Han, M. S., Kim, I. S., Lee, J. H., & Lee, J. C. 2013. Adaptation Test of Scotch Pine

291       (*Pinus sylvestris* L.) in Korea-Thirty-six-year-old Growth Performance of Twenty-two

292       Provenances. *Korean Journal of Plant Resources*, *26*(1), 26-35.

293   Schattner P, Brooks AN, Lowe TM. 2005. The tRNAscan-SE, snoscan and snoGPS web servers

294       for the detection of tRNAs and snoRNAs. *Nucleic Acids Research* 33: W686–W689. DOI:

295       10.1093/nar/gki366

296   Syring J, Willyard A, Cronn R, Liston A. 2005. Evolutionary relationships among *Pinus*

297       (Pinaceae) subsections inferred from multiple low-copy nuclear loci. *American Journal of*

298       *Botany* 92: 2086–2100. DOI: 10.3732/ajb.92.12.2086

299   Wakasugi, T., Tsudzuki, J., Ito, S., Nakashima, K., Tsudzuki, T., & Sugiura, M. 1994 Loss of all

300       ndh genes as determined by sequencing the entire chloroplast genome of the black pine

301       *Pinus thunbergii. Proceedings of the National Academy of Sciences*, *91*(21), 9794-9798.

302   Wang, X. Q., Tank, D. C., & Sang, T. (2000). Phylogeny and divergence times in Pinaceae:

303       evidence from three genomes. *Molecular Biology and Evolution*, 17(5), 773-781.

304   Wang, X. R., Y. Tsumura, H. Yoshimaru, K. Nagasaka and A. E. Szmidt. 1999. Phylogenetic

305       relationships of Eurasian pines (*Pinus*, Pinaceae) based on chloroplast *rbcL*, *matK*, *rpl20-*

306       *rps18* spacer, and *trnV* intron sequences. *American Journal of Botany* 86: 1742-1753.

307   Wu CS, Chaw SM. 2014. Highly rearranged and size□variable chloroplast genomes in conifers

308       II clade (cupressophytes): evolution towards shorter intergenic spacers. *Plant Biotechnology*

309       *Journal* 12: 344–353. DOI: 10.1111/pbi.12141

310    Wu CS, Lin CP, Hsu CY, Wang RJ, Chaw SM. 2011. Comparative chloroplast genomes of

311       Pinaceae: insights into the mechanism of diversified genomic organizations. *Genome Biology*

312       *and Evolution* 3: 309–319. DOI: 10.1093/gbe/evr026

313    Yi DK, Choi K, Joo M, Yang JC, Mustafina FU, Han JS, Lee YM. 2016. The complete chloroplast

314       genome sequence of *Abies nephrolepis* (Pinaceae: Abietoideae). *Journal of Asia-Pacific*

315       *Biodiversity* 9: 245–249. DOI: 10.1016/j.japb.2016.03.014

# Figure 1

Gene maps and summary of the *Pinus densiflora* S. *et* Z. chloroplast genome.

Genes lying outside the circle are transcribed in a clockwise direction, whereas genes inside are transcribed in a counterclockwise direction. Different colors denote known functional groups. The GC and AT contents of the genome are denoted by dashed darker and lighter gray in the inner circle. LSC, SSC, and IR indicate large single-copy, small single-copy, and inverted repeat regions, respectively.

# Figure 2

Distribution of repeats present in *Pinus* chloroplast genomes. A) Distribution of SSRs present in *Pinus* chloroplast genomes. B) Distribution of long repeat sequences present in *Pinus* chloroplast genomes.
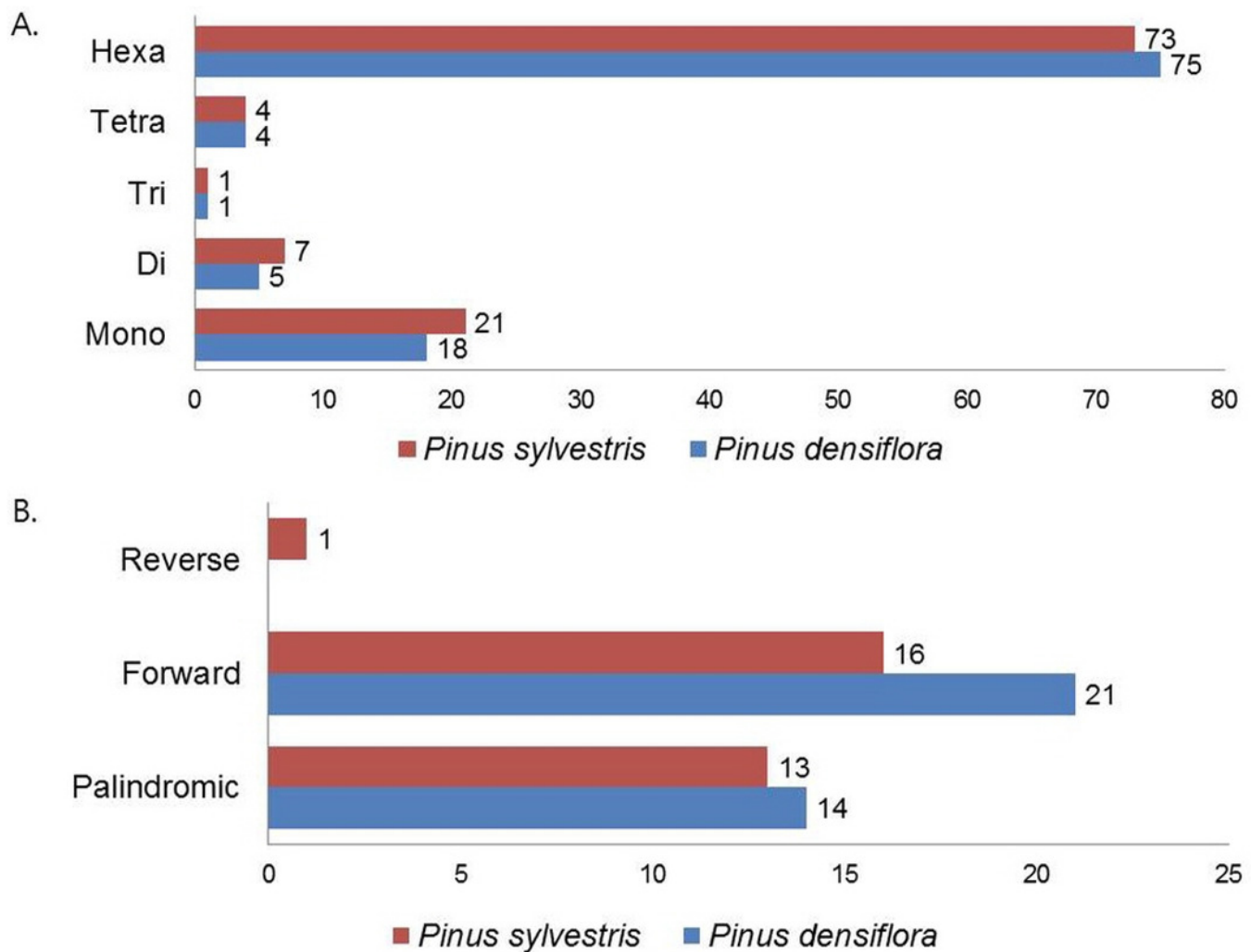
# Figure 3

Relevant part of the SNP multiple sequence alignment in the *trnM-trnV* gene regions.

The recognition site of *Hinf*I restriction enzyme (G/ANTC) is altered by one SNP at position 84 (A/G transition). A fragment degraded to 198 bp was observed In *P. densiflora* (1-9), but not in *P. sylvestris* (10-21).
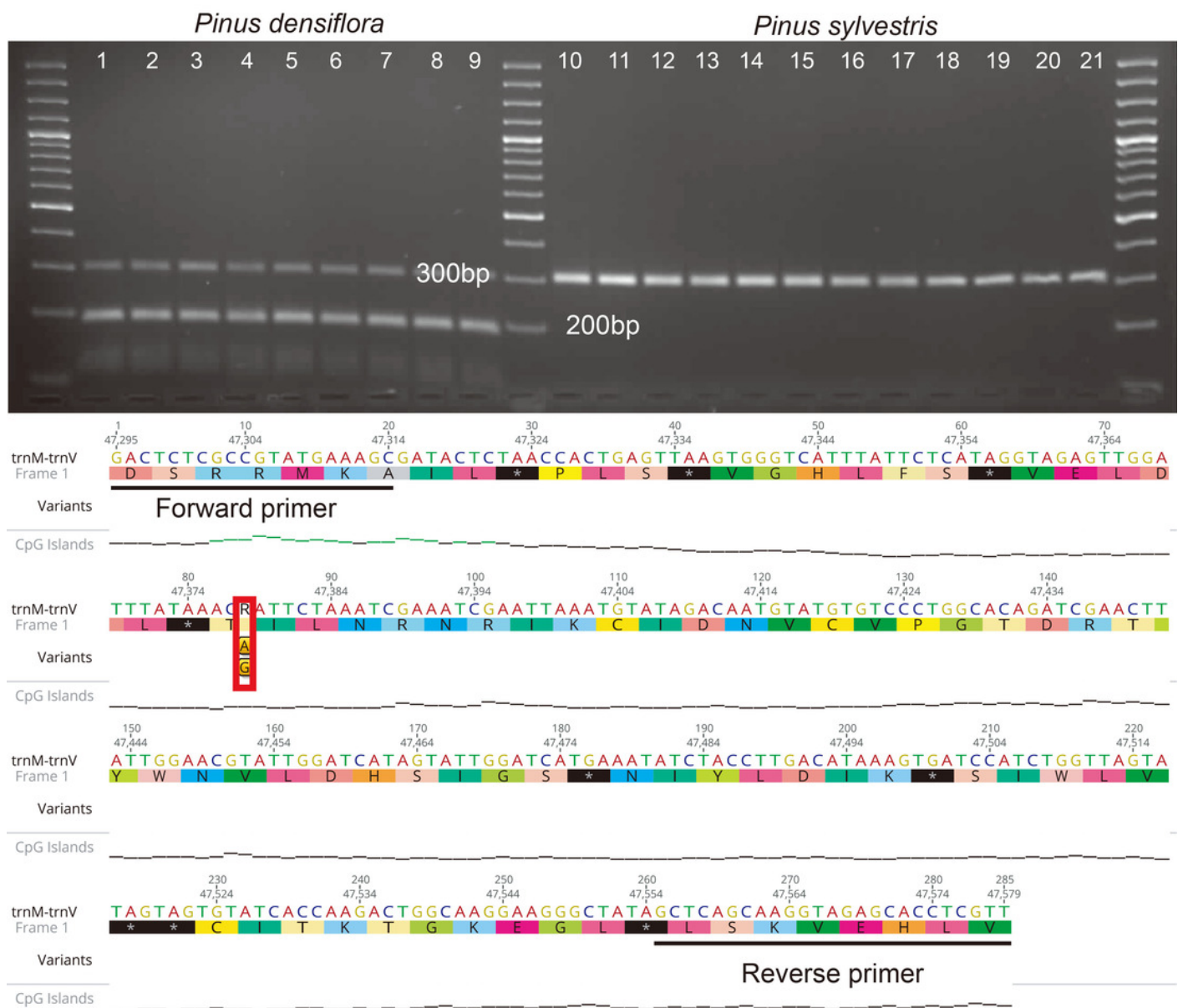
# Figure 4

Relevant part of the SNP multiple sequence alignment in the *rpoC1* gene.

The recognition site of *BsaW*I restriction enzyme (W/CCGGW) is altered by one SNP at position 184 (A/G transition). A fragment degraded to 179 bp was observed In *P. sylvestris* (10-21), but not in *P. densiflora* (1-9).
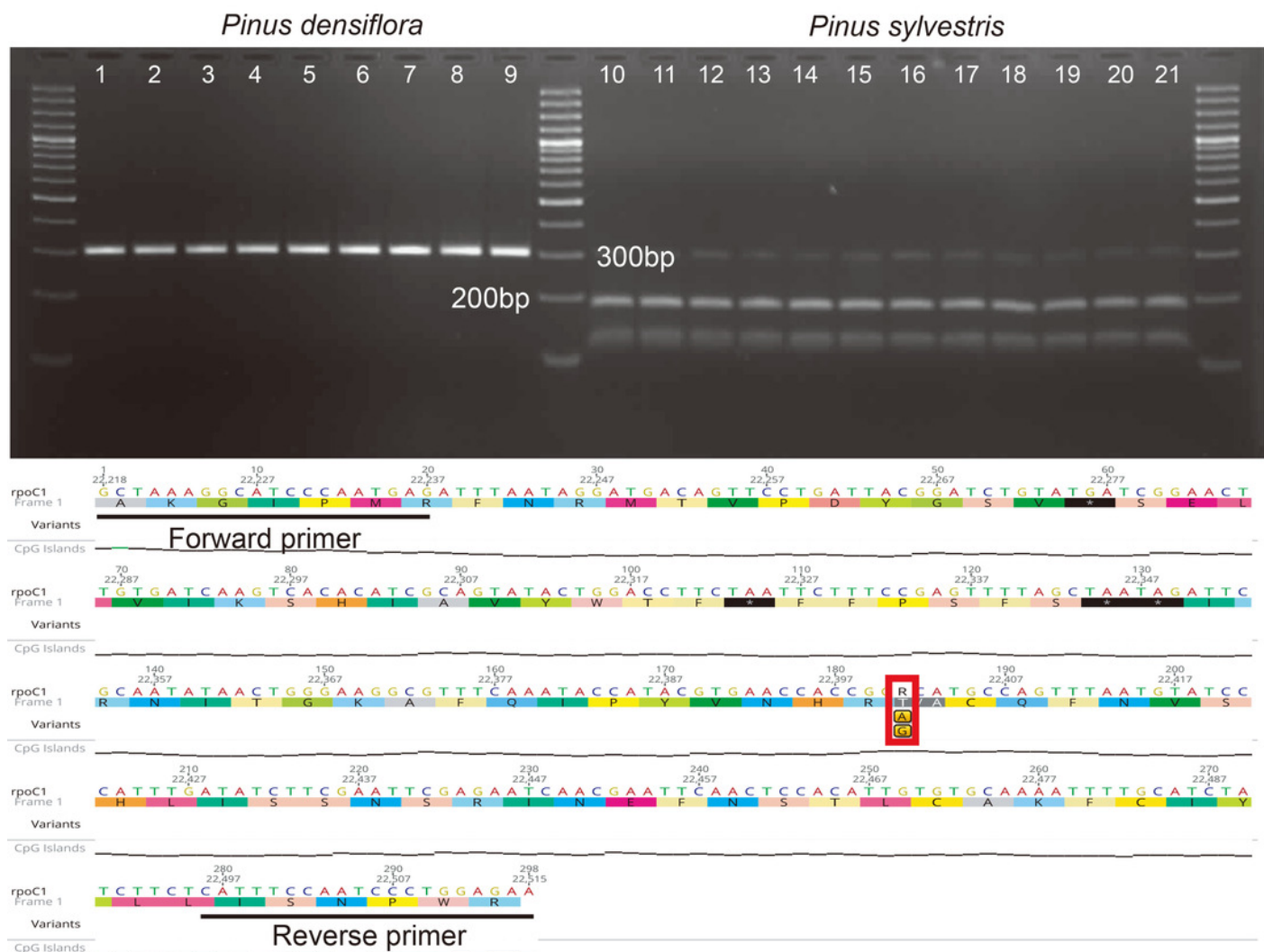
# Figure 5

Relevant part of the SNP multiple sequence alignment in the *rpoC1* gene.

The recognition site of *Sph*I restriction enzyme (GCATG/C) is altered by one SNP at position 184 (A/G transition). A fragment degraded to 184 bp was observed In *P. densiflora* (1-9), but not in *P. sylvestris* (10-21).
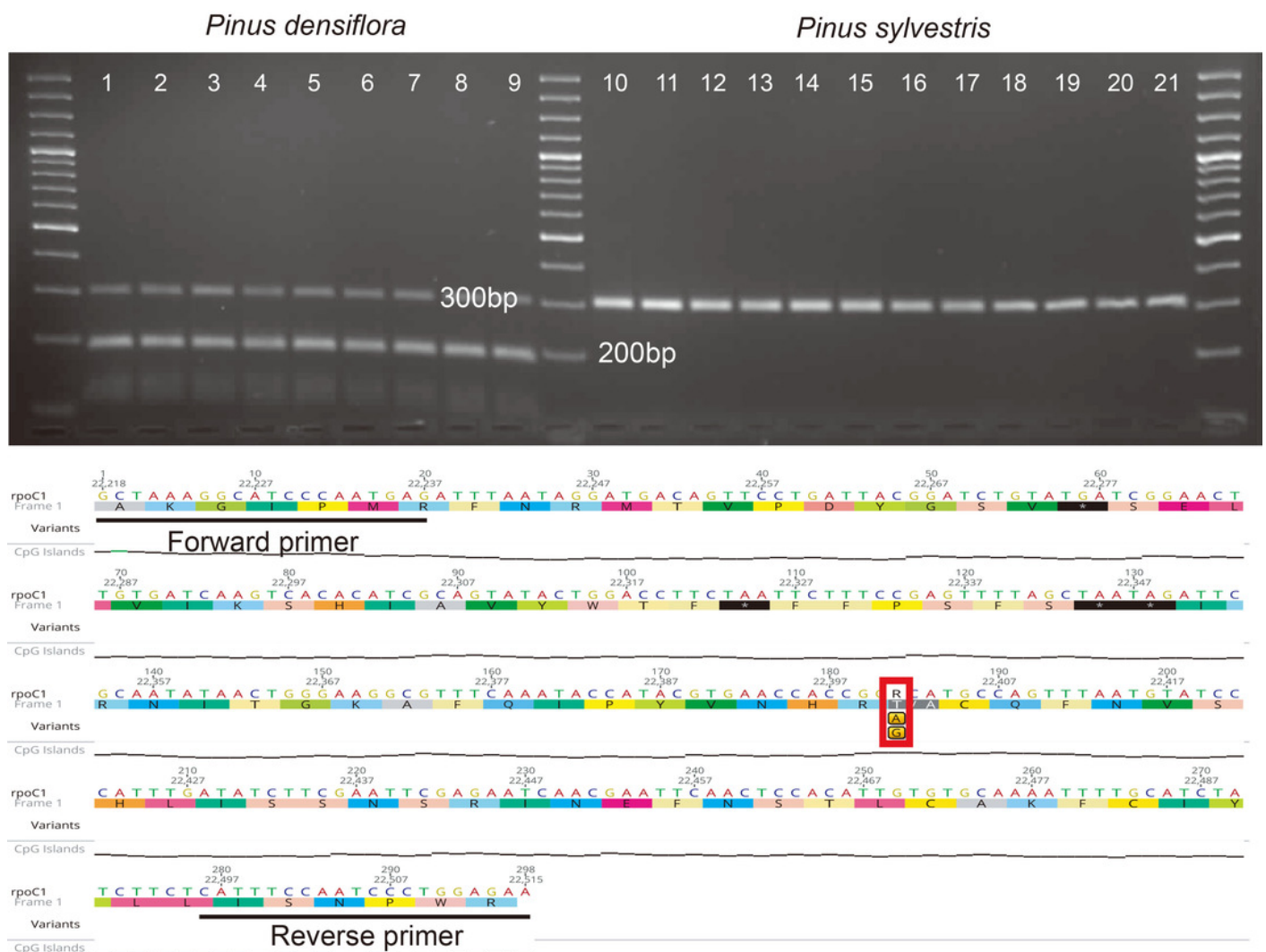
# Table 1(on next page)

List of genes encoded by the *P. densiflora* chloroplast genome.

Table 1. List of genes encoded by the *P. densiflora* chloroplast genome.

| Gene types | Gene products | |
|---|---|---|
| Ribosomal RNAs | *rrn4.5, rrn5, rrn16, rrn23* | 4 |
| Transfer RNAs | *trnA*-UGC[a], *trnC*-GCA, *trnD*-GUC, *trnE*-UUC, *trnF*-GAA, *trnfM*-CAU, *trnG*-GCC, *trnG*-UCC[a], *trnH*-GUG, *trnH*-GUG, *trnI*-CAU(x2), *trnI*-GAU[a], *trnK*-UUU[a], *trnL*-CAA, *trnL*-UAA[a], *trnL*-UAG, *trnM*-CAU, *trnN*-GUU, *trnP*-GGG, *trnP*-UGG, *trnQ*-UUG, *trnR*-ACG, *trnR*-CCG, *trnR*-UCU, *trnS*-GCU, *trnS*-GCU, *trnS*-GGA, *trnS*-UGA, *trnT*-GGU, *trnT*-GGU, *trnT*-UGU, *trnV*-GAC, *trnV*-UAC[a], *trnW*-CCA, *trnY*-GUA | 36 |
| Photosystem I | *psaA, psaB, psaC, psaI, psaJ, psaM*(x2) | 7 |
| Photosystem II | *psbA, psbB, psbC, psbD, psbE, psbF, psbH, psbI, psbJ, psbK, psbL, psbM, psbN, psbT, psbZ* | 15 |
| Cytochrome b/f complex | *petA, petB[a], petD[a], petG, petL, petN* | 6 |
| ATP synthase | *atpA, atpB, atpE, atpF[a], atpH, atpI* | 6 |
| Large subunit of rubisco | *rbcL* | 1 |
| Chloroplast envelope membrane protein | *cemA* | 1 |
| Large subunit ribosomal proteins | *rpl2[a], rpl14, rpl16[a], rpl20, rpl22, rpl23, rpl32, rpl33, rpl36* | 9 |
| Small subunit ribosomal proteins | *rps2, rps3, rps4, rps7, rps8, rps11, rps12[b], rps14, rps15, rps18, rps19* | 11 |
| RNA polymerase | *rpoA, rpoB, rpoC1[a], rpoC2* | 4 |
| Translational initiation factor | *infA* | 1 |
| Subunit of acetyl-CoA-carboxylase | *accD* | 1 |
| C-type cytochrome synthesis gene | *ccsA* | 1 |
| Maturase | *matK* | 1 |
| Chlorophyll biosynthesis | *chlB, chlL, chlN* | 3 |
| ATP-dependent protease | *clpP* | 1 |
| Conserved open reading frames | *ycf1, ycf2, ycf3[b], ycf4, ycf12, ycf68* | 6 |
| Total | | 114 |

[a]: Gene containing a single intron.
[b]: Gene containing two introns.

3

Peer Preprints

**Table 2**(on next page)

Distribution of Indels in *Pinus* chloroplast genomes

| P. densiflora | P. sylvestris | Minimum | Maximum | Length | Loci |
|---|---|---|---|---|---|
| | T | 1376 | 1376 | 1 | psbA-trnK |
| CAG | | 8212 | 8214 | 3 | psaM |
| | GGG | 9850 | 9852 | 3 | prnG-trnR |
| | CC | 15202 | 15203 | 2 | atpl-rps2 |
| | TCTA | 15236 | 15239 | 4 | atpl-rps2 |
| CTATTTCTCAAGA | | 16150 | 16162 | 13 | rps2-rpoC2 |
| A | | 26128 | 26128 | 1 | rpoB-trnC |
| ATTTAAATAATTTTGATAATTTTAATT | | 29221 | 29247 | 27 | trnE-clpP |
| | G | 30329 | 30329 | 1 | clpP-rps12 |
| TATTTTCTTC | | 36493 | 36502 | 10 | psbJ-petA |
| AATTTCAATAAATATTTCATTGTATGAAAATGG | | 39594 | 39626 | 33 | cemA-ycf4 |
| T | | 41146 | 41146 | 1 | psaI-accD |
| | TTTTTTTATTT | 45123 | 45133 | 11 | rbcL-atpB |
| CTG | | 51621 | 51623 | 3 | psaM |
| T | | 51963 | 51963 | 1 | trnS-psbB |
| | T | 63817 | 63817 | 1 | rps19-rpl2 |
| | GC | 65871 | 65872 | 2 | psbA-trnI |
| | AA | 68122 | 68123 | 2 | trnF-trnL |
| CTCCCCTTCT | | 68911 | 68920 | 10 | trnL-trnT |
| TTTTTTTT | | 72059 | 72066 | 8 | ycf3 intron |
| AT | | 73572 | 73572 | 2 | ycf3-psaA |
| | C | 82837 | 82837 | 1 | psbD-trnT |
| CAATTTGTTGT | | 93896 | 93906 | 11 | chlL-chlN |
| | T | 101253 | 101253 | 1 | ycf1-rps15 |
| | T | 102793 | 102793 | 1 | ndhI-ndhE |
| A | | 108576 | 108576 | 1 | trnV-rps12 |
| | TCATA | 109477 | 109481 | 5 | trnV-rps12 |
| | A | 109734 | 109734 | 1 | trnV-rps12 |
| | AA | 110195 | 110196 | 2 | rps12 intron |
| | AGAAAAAAA | 115341 | 115349 | 9 | ycf2 |

# Table 3 (on next page)

Distribution of SNPs in *Pinus* chloroplast genomes

| P. densiflora | P. sylvestris | Seqeunce No. | Mutation | Loci | P. densiflora | P. sylvestris | Seqeunce No. | Mutation | Loci |
|---|---|---|---|---|---|---|---|---|---|
| A | G | 4190 | transition | trnK-chlB | A | T | 48380 | transversion | trnV-trnH |
| A | C | 5708 | transversion | chlB | G | T | 59337 | transversion | rps11 |
| G | T | 7608 | transversion | psbK-psbI | C | T | 62762 | transition | rps3 |
| A | T | 7613 | transversion | psbK-psbI | G | A | 69894 | transition | rps4 |
| A | G | 8182 | transition | trnS-psaM | G | T | 79253 | transversion | trnfM-psbZ |
| T | A | 8183 | transversion | trnS-psaM | G | T | 80551 | transversion | psbC |
| C | T | 8184 | transition | trnS-psaM | G | A | 82941 | transition | psbD-trnT |
| A | G | 9855 | transition | trnG-trnR | G | C | 83659 | transversion | psbD-trnT |
| A | G | 9857 | transition | trnG-trnR | A | T | 93066 | transversion | chlL |
| A | G | 10538 | transition | atpA | G | A | 95123 | transition | chlN |
| T | G | 14210 | transversion | atpH-atpI | C | A | 96483 | transversion | ycf1 |
| A | C | 15041 | transversion | atpI-rps2 | C | A | 96893 | transversion | ycf1 |
| C | A | 22063 | transversion | rpoC1 intron | T | G | 97302 | transversion | ycf1 |
| G | A | 22401 | transition | rpoC1 | T | G | 97763 | transversion | ycf1 |
| C | A | 22743 | transversion | rpoC1-rpoB | C | T | 99931 | transition | ycf1 |
| T | A | 29216 | transversion | trnE-clpP | G | T | 100184 | transversion | ycf1 |
| G | A | 29217 | transition | trnE-clpP | G | A | 100412 | transition | ycf1 |
| A | T | 29274 | transversion | trnE-clpP | A | G | 100435 | transition | ycf1 |
| A | T | 29283 | transversion | trnE-clpP | T | A | 100541 | transversion | ycf1 |
| A | T | 29288 | transversion | trnE-clpP | A | C | 100552 | transversion | ycf1 |
| A | G | 29289 | transition | trnE-clpP | C | A | 101591 | transversion | rps15-ndhH |
| T | A | 30313 | transversion | clpP-rps12 | C | A | 103342 | transversion | ndhE-psaC |
| C | A | 30314 | transversion | clpP-rps12 | T | G | 108063 | transversion | rpl32-trnV |
| T | C | 31522 | transition | rpl20 | C | A | 108688 | transversion | trnV-rps12 |
| T | G | 33070 | transversion | psaJ-trnP | G | T | 110180 | transversion | rps12 intron |
| G | T | 34025 | transversion | petL | A | C | 111793 | transversion | ndhB |
| A | C | 34211 | transversion | petL-psbE | T | G | 111948 | transversion | ndhB-trnL |
| A | G | 42181 | transition | accD | C | A | 112340 | transversion | trnL-ycf2 |
| C | T | 42668 | transition | accD-trnR | C | A | 113412 | transversion | ycf2 |
| A | G | 42992 | transition | trnR-rbcL | C | A | 115338 | transversion | ycf2 |
| G | A | 43272 | transition | rbcL | T | C | 115351 | transition | ycf2 |
| G | A | 43740 | transition | rbcL | A | T | 115354 | transversion | ycf2 |
| C | T | 43762 | transition | rbcL | A | G | 115791 | transition | ycf2 |
| G | C | 44583 | transversion | rbcL-atpB | A | C | 116833 | transversion | ycf2 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| C | T | 46241 | transition | *atpB* | G | T | 118062 | transversion | *ycf2* |
| G | A | 47378 | transition | *trnM-trnV* | T | G | 118560 | transversion | *ycf2* |
| T | A | 48379 | transversion | *trnV-trnH* | | | | | |

**Table 4**(on next page)

Primers and restriction enzymes used for identification of *Pinus* species

|  | forward primer | reverse primer | enzyme | Loci |
|---|---|---|---|---|
| Pdest-cp1 | GACTCTCGCCGTATGAAAGC | GCAAGGTAGAGCACCTCGTT | *Hinf*I | *rpoC1* |
| Pdest-cp2 | GCTAAAGGCATCCCAATGAG | TTCTCCAGGGATTGGAAATG | *BsaW*I, *Sph*I | *trnM-trnV* |