

1 **Siberian larch (*Larix sibirica* Ledeb.) chloroplast genome and**  
2 **development of polymorphic chloroplast markers**

3 Eugeniya I. Bondar<sup>1</sup>, Yuliya A. Putintseva<sup>1</sup>, Nataliya V. Oreshkova<sup>1,2</sup>,

4 Konstantin V. Krutovsky<sup>1,3,4,5</sup>

5 <sup>1</sup>Laboratory of Forest Genomics, Genome Research and Education Center, Siberian Federal  
6 University, Krasnoyarsk, Russian Federation

7 <sup>2</sup>Laboratory of Forest Genetics and Selection, V.N. Sukachev Institute of Forest, Siberian  
8 Branch of Russian Academy of Sciences, Krasnoyarsk, Russian Federation

9 <sup>3</sup>Department of Forest Genetics and Forest Tree Breeding, Georg-August University of  
10 Göttingen, Göttingen, Germany

11 <sup>4</sup>Laboratory of Population Genetics, Vavilov Institute of General Genetics, Russian Academy  
12 of Sciences, Moscow, Russia

13 <sup>5</sup>Department of Ecosystem Science and Management, Texas A&M University, College  
14 Station, Texas, USA

15 **Abstract**

16 **Background:** The main objectives of this study were sequencing, assembling and annotation  
17 of chloroplast genome of one of the main Siberian boreal forest tree conifer species Siberian  
18 larch (*Larix sibirica* Ledeb.) and detection of polymorphic genetic markers – microsatellite loci  
19 or simple sequence repeats (SSRs) and single nucleotide polymorphisms (SNPs).

20 **Results:** We used data of the whole genome sequencing of three Siberian larch trees from  
21 different regions - Urals, Krasnoyarsk, and Khakassia, respectively. Sequence reads were  
22 obtained using the Illumina HiSeq2000 in the Laboratory of Forest Genomics at the Genome  
23 Research and Education Center in the Siberian Federal University. The assembling was done  
24 using the Bowtie2 mapping program and the SPAdes genomic assembler. The genome  
25 annotation was performed using the RAST service. We used the SciRoKo program for the SSRs  
26 search, and the Bowtie2 and UGENE programs for the SNPs detection. Length of the assembled  
27 chloroplast genome was 122,561 bp, which is similar to 122,474 bp in the closely related  
28 European larch (*Larix decidua* Mill.). As a result of annotation and comparison of the data with

29 existing data available only for three larch species - *L. decidua*, *L. potaninii var. chinensis*  
30 (complete genome 122,492 bp) and *L. occidentalis* (partial genome of 119,680 bp), we  
31 identified 110 genes, 34 of which represented tRNA, 4 rRNA and 72 protein-coding genes. In  
32 total, 13 SNPs were detected; two of them were in the *tRNA-Arg* and *Cell division protein FtsH*  
33 genes, respectively.

34 **Conclusions:** The complete chloroplast genome sequence was obtained for Siberian larch for  
35 the first time. The reference complete chloroplast genomes, such as one described here would  
36 greatly help in the chloroplast resequencing and search for additional genetic markers using  
37 population samples. The results of this research will be useful for further phylogenetic and gene  
38 flow studies in conifers.

39 **Keywords:** Chloroplast genome, *Larix sibirica*, Sequencing, Siberian larch, SNPs, SSRs

40 **Correspondence:** Konstantin V. Krutovsky, Department of Forest Genetics and Forest Tree  
41 Breeding, Georg-August University of Göttingen, Büsgenweg 2, D-37077 Göttingen, Germany;  
42 Fax:+49 (551)398367; E-mail: [konstantin.krutovsky@forst.uni-goettingen.de](mailto:konstantin.krutovsky@forst.uni-goettingen.de); ORCID:  
43 <http://orcid.org/0000-0002-8819-7084>

44

## 45 **Background**

46 Chloroplast genome in conifers, including larch species [1] has a unique, strictly paternal  
47 inheritance via pollen, unlike angiosperms, where it has a maternal inheritance via seeds [2]. It  
48 allows tracing paternal gene flow and lineages separately from maternal (mitochondrial genes)  
49 and bi-parental (nuclear genes) ones. Therefore, chloroplast DNA sequences are the most  
50 important source of genetic markers to study distribution of paternal genes and paternally based  
51 molecular phylogenetic relationships in conifers.

52 Larch species, as well as many other conifer species are the main boreal forest tree species,  
53 which comprise ~30% of the world's forested lands [3]. Boreal forests play very important  
54 ecological role, but are also affected by the global climate change. On one hand, they suffer

55 now from more frequent and drastic droughts, but on the other hand their area is expanding in  
56 the northern regions, and their tree line is moving towards north creating an ecotone, a highly  
57 dynamic transition area [4]. It is important to know how much of paternal associated gene flow  
58 by pollen contributes into establishing this zone compared to the maternal and bi-parental  
59 contributions by seeds. Such studies require chloroplast markers. Next generation sequencing  
60 (NGS) technique allows whole chloroplast genome sequencing in multiple individuals and  
61 makes a search for the molecular genetic markers more efficient. For instance, Parks *et al.* [5]  
62 nearly completely sequenced chloroplast genomes in 37 pine species using NGS. They found  
63 significant amount of variation (especially in two loci *ycf1* and *ycf2*) that provided them with  
64 additional data for inferring intrageneric phylogeny of genus *Pinus*.

65 Whole chloroplast genome comparison across different species and genera allows also  
66 studying organelle evolution and how it is associated with speciation and dispersal. Complete  
67 chloroplast genome sequences are available in NCBI Genbank for multiple plant species,  
68 including conifers. However, most of them represent the *Pinus* genus, and only three chloroplast  
69 genomes are available for the *Larix* genus: complete for European (*Larix decidua* Mill.;  
70 AB501189.1) and Chinese (*L. potaninii* var. *chinensis* Beissn.; KX808508) larch and partial for  
71 Western larch (*L. occidentalis* Nutt.; FJ899578.1).

72 Variation in the chloroplast genome is effectively used in phylogenetics at different levels.  
73 It allowed discriminating different subgenera and genera. For instance, Cronn *et al.* [6]  
74 compared chloroplast genome sequences of seven pine and one spruce species and found three  
75 regions that have deletions corresponded to the subgenera specific deletions in three genes:  
76 *ycf12* (78 bp at the nucleotide starting position 51051), *psaM* (93 bp at position 51442) and  
77 *ndhI* (371 bp at position 101988), respectively. These are common deletions in the chloroplast  
78 genome in the pine species of the subgenus *Strobus* (i.e., *P. gerardiana*, *P. krempfii*, *P.*

79 *lambertiana*, *P. longaeva*, *P. monophylla*, *P. nelsonii*, *P. koraiensis*); the corresponding genes  
80 were present in the subgenus *Pinus* (*P. contorta*, *P. ponderosa*, *P. thunbergii*) and in spruce  
81 *Picea sitchensis* [6].

82 Variation in the chloroplast genome can be also effectively used in discriminating different  
83 populations of the same species. For instance, Whittall *et al.* [7] demonstrated a strong  
84 differentiation between mainland and island populations of Torrey pine (*Pinus torreyana*) based  
85 on 5 SNPs found in the entire chloroplast genome of 120 Kbp.

## 86 **Methods**

87 We used data of the whole genome sequencing of three Siberian larch trees generated by  
88 Illumina HiSeq2000 [8]. DNA samples were isolated from needles and haploid callus of three  
89 Siberian larch trees, representing different regions in Russia – Ural Mountains, Krasnoyarsk  
90 Region and Khakassia Republic, respectively. *Larix decidua* Mill. [9] and *L. occidentalis* Nutt.  
91 [5] chloroplast genomes were used as reference (NCBI Genbank accession numbers  
92 AB501189.1 and FJ899578.1, respectively). We did not use the chloroplast genome of *L.*  
93 *potaninii* [10] as a reference, because it was assembled by using the chloroplast genome of *L.*  
94 *decidua* (NC\_016058; [9]) as a reference, but we used it in the comparative analysis. The paired-  
95 end (PE) and mate-pair (MP) libraries with fragment sizes of 400-500 bp (Ural and Krasnoyarsk  
96 trees) and 300-400 bp (Khakassia tree), respectively, were used for sequencing via  $2 \times 100$   
97 cycles by Illumina HiSeq2000.

98 The sequence reads were mapped to the reference chloroplast genomes using the Bowtie2  
99 software [11], which is good for mapping short sequence reads to medium-sized and large  
100 genomes. This software implements an algorithm to derive FM-index based on Burrows-  
101 Wheeler Transform. The SPAdes genome assembler has been used to assemble the larch

102 genome, which implements the De Bruijn graph approach [12]. The Rapid Annotation service  
103 with Subsystem Technology (RAST) has been used for annotation [13].

104 The first step in our assembly procedure consisted of mapping short reads to the available  
105 chloroplast genome references of *L. decidua* and *L. occidentalis* using the Bowtie2 software.  
106 Then, the aligned reads were assembled by SPAdes. Obtained contigs were aligned again on  
107 the reference of *L. decidua* using BLAST. At the third step, the selected contigs were verified  
108 to get the “trusted” status. Then, the assembly was carried out using SPAdes. The final step of  
109 the assembly was the scaffolding, which was done using the generated contigs and MP reads  
110 using the SSPACE program [14].

111 Considering a well-known fact that chloroplast organelle originated from cyanobacteria, and  
112 that, therefore, chloroplast genes are still very similar to the bacterial ones, the RAST service,  
113 which was designed for annotation of bacterial and archaeal genomes, was used for the larch  
114 genome annotation. The annotation obtained by the RAST contained both confirmed known  
115 genes and predicted genes, potentially coding hypothetical proteins. In order to clarify the roles  
116 of these hypothetical coding regions our annotation was compared with annotations of two  
117 closely related species *L. decidua* and *L. occidentalis*, respectively. In addition, some fragments  
118 of the genome have been also selectively aligned with BLAST. Sites of hypothetical proteins  
119 confirmed by BLAST were identified and recorded.

120 SNPs were search using the Bowtie2 and UGENE [15] software (option *Call Variants with*  
121 *SAMtools*). The search was done across the three above mentioned trees. First, reads of Urals  
122 and Khakassian trees were mapped to the finally assembled genome of the Krasnoyarsk tree.  
123 The resulting *sam*-file together with the assembled genome was used by the UGENE program  
124 to search for SNPs.

## 125 **Results**

126 The total length of the final Siberian larch chloroplast genome assembly was 122,560 bp, which  
127 is very close to 122,474 bp in closely related European larch (*Larix decidua*). The annotation  
128 through the comparison with available data for *L. decidua* and *L. occidentalis* identified 110  
129 genes, from which 34 represented tRNA genes, 4 rRNA and 72 protein-coding genes. In three  
130 trees 13 SNPs were detected. Two of them were found in the coding regions of the *tRNA-Arg*  
131 and *Cell division protein ycf2* genes.

132 We used available software, such as Bowtie2, BLAST and SPAdes to assemble chloroplast  
133 genome using reads generated in the whole genome sequencing of Siberian larch project. We  
134 used SSPACE for scaffolding and the RAST service for annotation of obtained chloroplast  
135 genome. We developed a procedure that allowed us to successfully extract chloroplast genome  
136 specific reads and then assemble and annotate the resulting sequences. We identified and  
137 verified 110 coding regions representing 38 RNA and 72 protein genes, which is equal to the  
138 number of genes in chloroplast sequences of *L. decidua* and *L. potaninii* and close to 105 genes  
139 in a partial chloroplast genome sequence of *L. occidentalis*. A gene map of the genome was  
140 generated using OGDRAW [16] and presented in Fig. 1. Search for SNPs using UGENE  
141 revealed a relatively small number of SNPs (Fig. 2; Additional file 1), but it is only preliminary  
142 data based on a limited sample size.

## 143 **Discussion**

144 The chloroplast genome variation in most plants is often limited due to a relatively low  
145 frequency of mutations in this organelle. For example, the mutation rate of the chloroplast  
146 genome in pines is approximately  $0.2-0.4 \times 10^{-9}$  synonymous substitutions per nucleotide per  
147 year [17, 18]. However, with an average length of 120-160 Kbp and 130 genes chloroplast

148 genomes are sufficiently large and complex and include structural and point mutations that  
149 reflect population differentiation and evolutionary divergence [6].

150 Unlike angiosperms, conifer chloroplast DNA (cpDNA) lacks large inverted repeats (IR),  
151 but contains dispersed repetitive DNA that is associated with structural rearrangements. In  
152 addition to large dispersed repeated sequences, conifer cpDNA also possess a number of small  
153 repeats. It contains variable numbers of tandem repeats of 124 to 150 bp in size, which are  
154 associated with polymorphic rearranged region near *trnK-psbA*, where the *psbA* gene has been  
155 duplicated [19].

156 Most variation in the chloroplast genome is associated with microsatellite loci [20, 21].  
157 However, these markers have a too high mutation rate that can lead to the incorrect phylogenetic  
158 inferences [22-24]. SNPs could be better markers for phylogenetic inferences, and comparative  
159 complete chloroplast genome studies are needed to discover these markers. The reference  
160 complete chloroplast genomes, such as one described here would greatly help in the chloroplast  
161 resequencing and search for SNPs using population samples.

## 162 **Conclusions**

163 The complete chloroplast genome sequence was obtained for Siberian larch for the first time.  
164 Annotation and comparison of the obtained data with data available only for two other larch  
165 species helped us identify and verify 110 coding regions representing 38 RNA and 72 protein  
166 genes. Total 13 SNPs were detected; two of them were in the coding regions of the genome.  
167 The results of this research will be useful for further phylogenetic and gene flow studies in  
168 conifers.

## 169 **Additional files**

170 **Additional file 1:** excel file representing the Siberian larch chloroplast genetic variant data.

171 **Acknowledgements**

172 We thank Dr. M. Sadovsky for discussing results and his valuable comments, Dr. I. N.  
173 Tretyakova and M. E. Pak for providing with callus culture, Dr. V. L. Semerikov and K. Deich  
174 for help with sampling, Dr. A. Ibe – with laboratory analysis, and Dr. D. Kuzmin and S.  
175 Makolov - with computing. This study was supported by a research grant No. 14.Y26.31.0004  
176 from the Government of the Russian Federation.

177 **Funding**

178 The presented study was a part of the project "Genomic studies of major boreal coniferous  
179 forest tree species and their most dangerous pathogens in the Russian Federation" funded by  
180 the Government of the Russian Federation (grant № 14.Y26.31.0004).

181 **Availability of data and materials**

182 The annotated chloroplast genome of *L. sibirica* has been deposited in the NCBI GenBank with  
183 the accession number MF795085.

184 **Authors' contributions**

185 EIB and YAP assembled and annotated the chloroplast genome, analysed the data and wrote the  
186 draft paper, NVO prepared and sequenced DNA libraries, KVK designed and coordinated  
187 research and wrote the paper.

188 **Ethics approval and consent to participate**

189 Not applicable.

190 **Consent for publication**

191 Not applicable.

192 **Competing interests**

193 The authors declare that they have no competing interests.

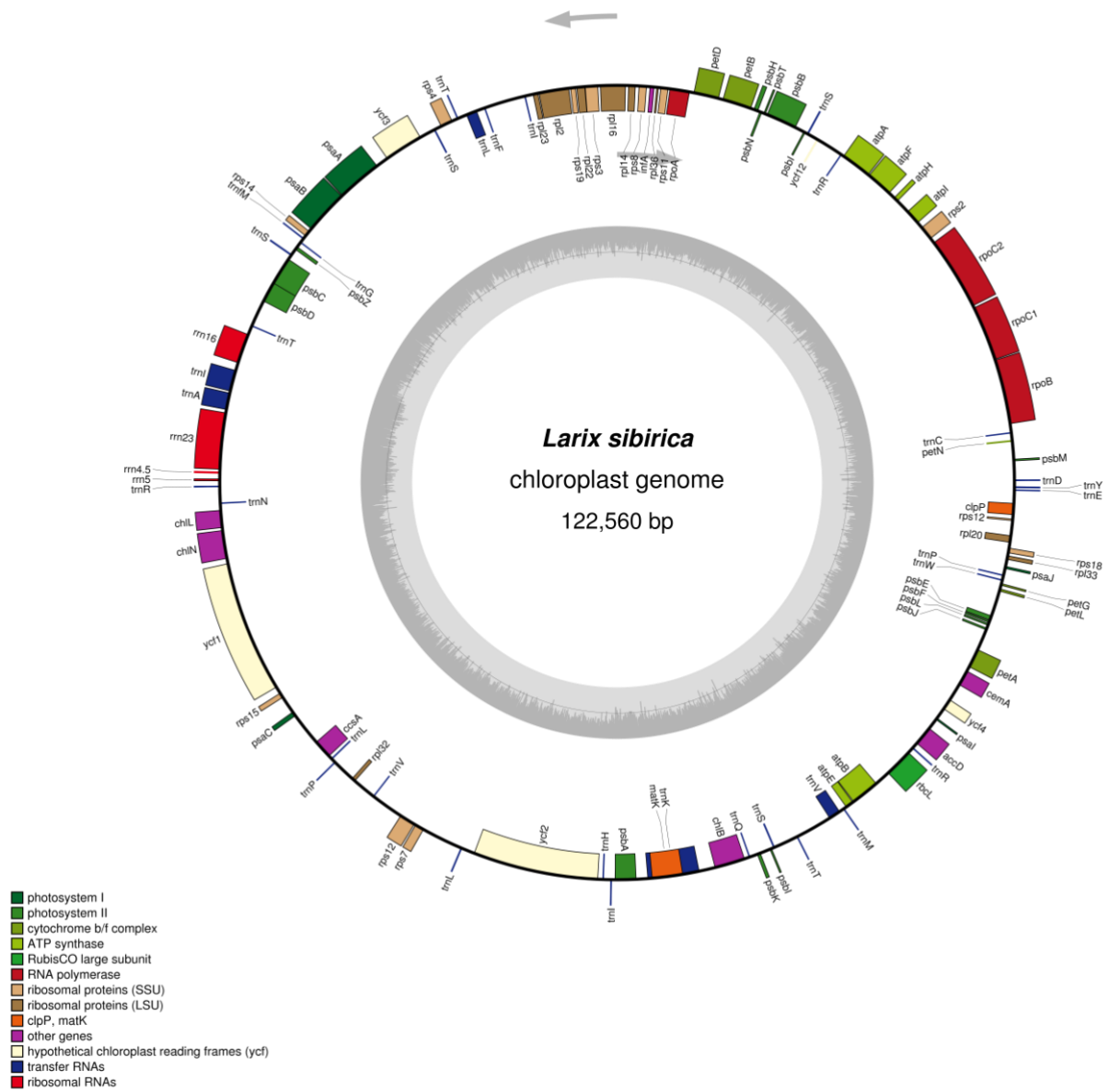


194 **References**

- 195 1. Szmidt AE, Aldén T, Hällgren J-E. Paternal inheritance of chloroplast DNA in *Larix*.  
196 Plant Mol Biol. 1987;9:59-64. <https://doi.org/10.1007/BF00017987>
- 197 2. Hipkins VD, Krutovskii KV, Strauss SH. Organelle genomes in conifers: structure,  
198 evolution, and diversity. Forest Genet. 1995;1:179-89.
- 199 3. Brandt JP, Flannigan MD, Maynard DG, Thompson ID, Volney WJA. An introduction to  
200 Canada's boreal zone: ecosystem processes, health, sustainability, and environmental  
201 issues. Environm Rev. 2013;21:207–26. <https://doi.org/10.1139/er-2013-0040>
- 202 4. Johnson JS, Gaddis KD, Cairns DM, Konganti K, Krutovsky KV. Landscape genomic  
203 insights into the historic migration of mountain hemlock in response to Holocene climate  
204 change. Am J Bot. 2017;104(3):439-50. <https://doi.org/10.3732/ajb.1600262>
- 205 5. Parks M, Cronn R, Liston A. Increasing phylogenetic resolution at low taxonomic levels  
206 using massively parallel sequencing of chloroplast genomes. BMC Biol. 2009;7:84.  
207 <https://doi.org/10.1186/1741-7007-7-84>
- 208 6. Cronn R, Liston A, Parks M, Gernandt DS, Shen R, Mockler T. Multiplex sequencing of  
209 plant chloroplast genomes using Solexa sequencing-by-synthesis technology. Nucleic  
210 Acids Res. 2008;36(19):e122. <https://doi.org/10.1093/nar/gkn502>
- 211 7. Whittall JB, Syring J, Parks M, Buenrostro J, Dick C, Liston A, Cronn R. Finding a (pine)  
212 needle in a haystack: chloroplast genome sequence divergence in rare and widespread  
213 pines. Mol Ecol. 2010;19:100–14. <https://doi.org/10.1111/j.1365-294X.2009.04474.x>
- 214 8. Krutovsky KV, Oreshkova NV, Putintseva YuA, Ibe AA, Deich KO, Shilkina EA.  
215 Preliminary results of *de novo* whole genome sequencing of Siberian larch (*Larix sibirica*  
216 Ledeb.) and Siberian stone pine (*Pinus sibirica* Du Tour.). Siberian J Forest Sci.  
217 2014;1:79-83.
- 218 9. Wu CS, Lin CP, Hsu CY, Wang RJ, Chaw SM. Comparative chloroplast genomes of  
219 Pinaceae: insights into the mechanism of diversified genomic organizations. Genome Biol  
220 Evol. 2011; 3:309-19. <https://doi.org/10.1093/gbe/evr026>
- 221 10. Han K, Li J, Zeng S, Liu Z. Complete chloroplast genome sequence of Chinese larch  
222 (*Larix potaninii* var. *chinensis*), an endangered conifer endemic to China. Conservation  
223 Genet Res. 2017;9:111–3. <https://doi.org/10.1007/s12686-016-0633-9>
- 224 11. Langmead B, Salzberg S. Fast gapped-read alignment with Bowtie 2. Nature Methods.  
225 2012;9:357-9. <https://doi.org/10.1038/nmeth.1923>

- 226 12. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes:  
227 A new genome assembly algorithm and its applications to single-cell sequencing. *J*  
228 *Comput Biol.* 2012;19:455-77. <https://doi.org/10.1089/cmb.2012.0021>
- 229 13. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, et al. The SEED and the  
230 Rapid Annotation of microbial genomes using Subsystems Technology (RAST). *Nucleic*  
231 *Acids Res.* 2014;42(Database issue):D206-D214. <https://doi.org/10.1093/nar/gkt1226>
- 232 14. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-assembled  
233 contigs using SSPACE. *Bioinformatics.* 2011;27:578-9.  
234 <https://doi.org/10.1093/bioinformatics/btq683>
- 235 15. Okonechnikov K, Golosova O, Fursov M, the UGENE team. Unipro UGENE: a unified  
236 bioinformatics toolkit. *Bioinformatics.* 2012;28:1166-7.  
237 <https://doi.org/10.1093/bioinformatics/bts091>.
- 238 16. Lohse M, Drechsel O, Kahlau S, Bock R. OrganellarGenomeDRAW – a suite of tools  
239 for generating physical maps of plastid and mitochondrial genomes and visualizing  
240 expression data sets. *Nucleic Acids Res.* 2013;41:575–81.  
241 <https://doi.org/10.1093/nar/gkt289>
- 242 17. Willyard A, Syring J, Gernandt DS, Liston A, Cronn R. Fossil calibration of molecular  
243 divergence infers a moderate mutation rate and recent radiations for *Pinus*. *Mol Biol*  
244 *Evol.* 2007;24:90-101. <https://doi.org/10.1093/molbev/msl131>
- 245 18. Gernandt DS, Magallon S, Lopez GG, Flores OZ, Willyard A, Liston A. Use of  
246 simultaneous analyses to guide fossil-based calibrations of Pinaceae phylogeny. *Int J*  
247 *Plant Sci.* 2008;169:1086–99. <https://doi.org/10.1086/590472>
- 248 19. Lidholm J, Gustafsson P. A three-step model for the rearrangement of the chloroplast  
249 *trnK-psbA* region of the gymnosperm *Pinus contorta*. *Nucleic Acids Res.* 1991. 19:2881–  
250 87.
- 251 20. Provan J, Soranzo N, Wilson NJ, Goldstein D B, Powell W A. A low mutation rate for  
252 chloroplast microsatellites. *Genetics.* 1999;153:943–47.
- 253 21. Ebert D, Peakall R. Chloroplast simple sequence repeats (cpSSRs): technical resources  
254 and recommendations for expanding cpSSR discovery and applications to a wide array of  
255 plant species. *Mol Ecol Res.* 2009;9:673–90. [https://doi.org/10.1111/j.1755-](https://doi.org/10.1111/j.1755-0998.2008.02319.x)  
256 [0998.2008.02319.x](https://doi.org/10.1111/j.1755-0998.2008.02319.x)

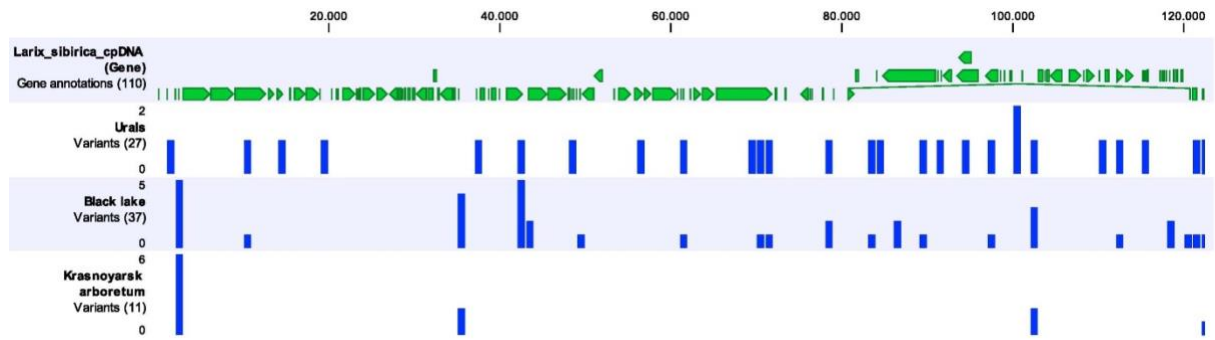
- 257 22. Afzal-Rafii Z, Dodd A. Chloroplast DNA supports a hypothesis of glacial refugia over  
258 postglacial recolonization in disjunct populations of black pine (*Pinus nigra*) in Western  
259 Europe. *Mol Ecol.* 2007;16:723–36. <https://doi.org/10.1111/j.1365-294X.2006.03183.x>
- 260 23. Höhn M, Gugerli F, Abran P, Bisztray G, Buonamici A, Cseke K., et al. Variation in the  
261 chloroplast DNA of Swiss stone pine (*Pinus cembra* L.) reflects contrasting post-glacial  
262 history of populations from the Carpathians and the Alps. *J Biogeography.* 2009;36:1798–  
263 1806. <https://doi.org/10.1111/j.1365-2699.2009.02122.x>
- 264 24. Moreno-Letelier A, Piñero D. Phylogeographic structure of *Pinus strobiformis* Engelm.  
265 across the Chihuahuan Desert filter-barrier. *J Biogeography.* 2009;36:121–31.  
266 <https://doi.org/10.1111/j.1365-2699.2008.02001.x>
- 267



269

270 **Fig. 1** Gene map of the *Larix sibirica* chloroplast genome. Genes belonging to different  
 271 functional groups are color-coded. The dark and light grey in the inner circle represents the  
 272 GC and AT content, respectively

273



274

275 **Fig. 2** Variation detected in the *Larix sibirica* chloroplast genome (see also Additional file 1)