

1 Genes Genomes Genetics G3

2 Article type: Genome Report

3 **A reference genome sequence for the European silver fir (*Abies alba* Mill.): a**
 4 **community-generated genomic resource**

5

6 Elena Mosca ^{*}, Fernando Cruz [†], Jèssica Gómez-Garrido [†], Luca Bianco [‡], Christian Rellstab
 7 [§], Sabine Brodbeck [§], Katalin Csilléry ^{§, ††}, Bruno Fady ^{‡‡}, Matthias Fladung ^{§§}, Barbara Fussi
 8 ^{***}, Dušan Gömöry ^{†††}, Santiago C. González-Martínez ^{†††}, Delphine Grivet ^{§§§}, Marta Gut [†],
 9 ^{****}, Ole Kim Hansen ^{††††}, Katrin Heer ^{††††}, Zeki Kaya ^{§§§§}, Konstantin V. Krutovsky ^{*****},
 10 ^{†††††}, ^{†††††}, Birgit Kersten ^{§§}, Sascha Liepelt ^{††††}, Lars Opgenoorth ^{††††}, Christoph Sperisen [§],
 11 Kristian K. Ullrich ^{§§§§§}, Giovanni G. Vendramin ^{*****}, Marjana Westergren ^{††††††}, Birgit
 12 Ziegenhagen ^{††††}, Tyler Alioto ^{†, ****}, Felix Gugerli [§], Berthold Heinze ^{††††††}, Maria Höhn
 13 ^{§§§§§§}, Michela Troglio [‡], David B. Neale ^{*****}, ²

14 ^{*} C3A - Centro Agricoltura Alimenti Ambiente, University of Trento, via E. Mach 1, 38010 S.
 15 Michele a/Adige (TN); Italy (mosca.helena@gmail.com); [†] CNAG-CRG, Centre for Genomic
 16 Regulation (CRG), The Barcelona Institute of Science and Technology, BaldiriReixac 4,
 17 Barcelona 08028; Spain (fernando.cruz@cnag.crg.eu; jessica.gomez@cnag.crg.eu;
 18 marta.gut@cnag.crg.eu; tyler.alioto@cnag.crg.eu); [‡] Fondazione Edmund Mach, Via Mach 1,
 19 38010 S. Michele a/Adige (TN); Italy (michela.troglio@fmach.it; luca.bianco@fmach.it); [§]
 20 Swiss Federal Research Institute WSL, Zürcherstrasse 111, 8903 Birmensdorf; Switzerland
 21 (felix.gugerli@wsl.ch; christian.rellstab@wsl.ch; christoph.sperisen@wsl.ch;
 22 sabine.brodbeck@wsl.ch; katalin.csillery@wsl.ch); ^{**} Laboratoire d'Ecologie Alpine,
 23 Université Grenoble Alpes (LECA), Université Grenoble Alpes CS 40700; 38058 Grenoble
 24 cedex 9; France; ^{††} University of Zürich, Department of Evolutionary Biology and

25 Environmental Studies, Winterthurerstrasse 190, CH-8057 Zurich; †† Institut National de la
26 Recherche Agronomique (INRA), Unité de Recherche Ecologie des Forêts Méditerranéennes
27 (URFM), Site Agroparc, Domaine Saint Paul, 84914 Avignon; France (Bruno.fady@inra.fr);
28 §§ Thünen-Institute of Forest Genetics, Sieker Landstr. 2, 22927 Grosshansdorf; Germany
29 (matthias.fladung@thuenen.de; birgit.kersten@thuenen.de); *** Bavarian Office for Forest
30 Seeding and Planting (ASP), Applied Forest Genetics, Forstamtsplatz 1, 83317 Teisendorf;
31 Germany (barbara.fussi@asp.bayern.de); ††† Technical University in Zvolen, TG Masaryka
32 24, 96053 Zvolen; Slovakia (gomory@tuzvo.sk); ††† Institut National de la Recherche
33 Agronomique (INRA), UMR1202 Biodiversity, Genes & Communities (BIOGECO),
34 University of Bordeaux, 69, route d'Arcachon, 33610 Cestas; France (santiago.gonzalez-
35 martinez@pierroton.inra.fr); §§§ INIA Forest Research Centre, Carretera de la Coruña km 7.5,
36 28040 Madrid; Spain (dgrivet@inia.es); **** Universitat Pompeu Fabra (UPF), Plaça de la
37 Mercè, 10, 08002 Barcelona, Spain; †††† Department of Geosciences and Natural Resource
38 Management (IGN), University of Copenhagen, Rolighedsvej 23, 1958 Frederiksberg C;
39 Denmark (okh@ign.ku.dk); †††† Philipps-Universität Marburg, Faculty of Biology (PUM),
40 Karl-von-Frisch-Str. 8, 35032 Marburg; Germany (katrin.heer@biologie.uni-marburg.de;
41 liepelt@biologie.uni-marburg.de; lars.opgenoorth@uni-marburg.de;
42 birgit.ziegenhagen@biologie.uni-marburg.de); §§§§ Department of Biological Sciences
43 (METU), Middle East Technical University, 06800 Çankaya/Ankara; Turkey
44 (kayaz@metu.edu.tr); ***** Department of Forest Genetics and Forest Tree Breeding, Georg-
45 August University of Göttingen, Büsgenweg 2, 37077 Göttingen; Germany
46 (konstantin.krutovsky@forst.uni-goettingen.de); ††††† Laboratory of Population Genetics,
47 Vavilov Institute of General Genetics, Russian Academy of Sciences, Gubkina Str. 3,
48 Moscow 119991, Russia; ††††† Laboratory of Forest Genomics, Genome Research and

49 Education Center, Institute of Fundamental Biology and Biotechnology, Siberian Federal
50 University, 50a/2 Akademgorodok, Krasnoyarsk 660036, Russia; §§§§§§ Max Planck Institute
51 for Evolutionary Biology, Department for Evolutionary Genetics (MPI), August Thienemann
52 Str. 2, 24306 Ploen; Germany (ullrich@evolbio.mpg.de); ***** Institute of Biosciences and
53 BioResources, National Research Council, Via Madonna del Piano 10,50019 Sesto Fiorentino
54 (Firenze); Italy (giovanni.vendramin@ibbr.cnr.it); †††††† Slovenian Forestry Institute (SFI),
55 Gozdarskiinštitut Slovenije), Večna pot 2, 1000 Ljubljana; Slovenia
56 (marjana.westergren@gozdis.si); †††††† Federal Research and Training Centre for Forests,
57 Natural Hazards and Landscape (BFW), Seckendorff-Gudent Weg 8, 1130 Wien; Austria
58 (berthold.heinze@bfw.gv.at); §§§§§§ Faculty of Horticultural Science, Department of Botany
59 (SZIU/FHS), Szent Istvan University, 1118 Budapest; Hungary (Hohn.Maria@kertk.szie.hu);
60 ***** Department of Plant Sciences, University of California at Davis (UCD), Davis 95616;
61 USA (dbneale@ucdavis.edu)

62

63 ¹ Eric Bazin deceased in May 2017.

64 ² corresponding author

65 **RUNNING TITLE: SILVER FIR GENOME ABAL 1.1**

66

67 **Keywords:** *Abies alba*, annotation, conifer genome, genome assembly, chloroplast genome

68

69 Corresponding Author:

70 David B. Neale

71 Department of Plant Sciences,

72 University of California at Davis (UCD),

73 Davis 95616; USA

74 Phone: 530-754-8431

75 Email: dbneale@ucdavis.edu

76

77 Word counts 6,257 excluding references.

78

79 **Abstract**

80 Silver fir (*Abies alba* Mill.) is a keystone conifer of European montane forest ecosystems that
81 has experienced large fluctuations in population size during during the Quaternary and, more
82 recently, due to land-use change. To forecast the species' future distribution and survival, it is
83 important to investigate the genetic basis of adaptation to environmental change, notably to
84 extreme events. For this purpose, we here provide a first draft genome assembly and
85 annotation of the silver fir genome, established through a community-based initiative. DNA
86 obtained from haploid megagametophyte and diploid needle tissue was used to construct and
87 sequence Illumina paired-end and mate-pair libraries, respectively, to high depth. The
88 assembled *A. alba* genome sequence accounted for over 37 million scaffolds corresponding to
89 18.16 Gb, with a scaffold N50 of 14,051 bp. Despite the fragmented nature of the assembly, a
90 total of 50,757 full-length genes were functionally annotated in the nuclear genome. The
91 chloroplast genome was also assembled into a single scaffold (120,908 bp) that shows a high
92 collinearity with both the *A. koreana* and *A. sibirica* complete chloroplast genomes. This first
93 genome assembly of silver fir is an important genomic resource that is now publicly available
94 in support of a new generation of research. By genome-enabling this important conifer, this
95 resource will open the gate for new research and more precise genetic monitoring of European
96 silver fir forests.

INTRODUCTION

97

98

99 Conifers represent the dominant trees in some temperate and all boreal ecosystems and have
100 important economic value. They face the effects of climate change, with an increase in
101 temperature and lower precipitation, and increased frequency of extreme events, to which
102 some species may be unable to quickly adapt. Silver fir (*Abies alba* Mill.) is a keystone
103 conifer of European montane forest ecosystems, which is dominant in cool areas of the
104 temperate zone (Ellenberg 2009). Its distribution ranges from the Pyrenees to the Alps and the
105 Carpathians where it reaches its easternmost range edge (Figure S1). Growing interest in
106 silver fir has emerged because of its potential vulnerability to climate change (Cailleret *et al.*
107 2014), which could change conditions for its sustainable use and its economic value. In turn,
108 this species is more drought-resistant than other important species for timber production, such
109 as Norway spruce (Vitali *et al.* 2017), which could turn out to be beneficial under the
110 expected increase in extended future drought periods.

111 Several studies investigated the environmental effect on silver fir genetic diversity across the
112 Italian Alps, showing the association between its genetic diversity and seasonal minimum
113 temperature (Mosca *et al.* 2012) as well as between genetic diversity and both temperature
114 and soil type (Mosca *et al.* 2014). Recent studies confirmed the environmental effect on the
115 species' local adaptation, which was shaped by winter drought in marginal populations
116 (Roschanski *et al.* 2016); while in common gardens, selection on height was driven by
117 thermal stability (Csilléry *et al.* 2018). Another study confirmed the importance of the
118 Apennines as a refugium of genetic diversity (Piotti *et al.* 2017). All these studies were based
119 on a modest number of genetic markers (several hundred of single-nucleotide polymorphisms,

120 SNPs, or tens of simple sequence repeats, SSRs, also called microsatellites) due to the lack of
121 genomic resources.

122 Conifer genomes are very large, ranging from 4 to 35 giga base pairs (Gb) (Bennett and
123 Leitch 2012; Grotkopp *et al.* 2004; Zonneveld 2012), but their gene content is similar to that
124 of other vascular plants (Leitch *et al.* 2005). Conifer genomic resources have grown in recent
125 years due to the application of high throughput sequencing technologies (Reuter *et al.* 2015).
126 To date, a few conifer genomes have been fully sequenced, including: *Picea abies* (L.) Karst
127 (Nystedt *et al.* 2013), *Picea glauca* (Moench) Voss (Warren *et al.* 2015), *Pinus taeda* L.
128 (Neale *et al.* 2014), *Pinus lambertiana* Dougl. (Stevens *et al.* 2016), *Pseudotsuga menziesii*
129 (Mirb.) Franco (Neale *et al.* 2017) and *Larix sibirica* Ledeb (Kuzmin *et al.* 2019).

130 The present research is part of the Silver Fir Genome Project, which is a community
131 effort within the Alpine Forest Genomics Network (AForGeN, IUFRO WP 2.04.11,
132 <https://www.aforgen.org>). This network was established in 2011 with the intent to facilitate
133 information exchange and collaboration among researchers interested in studying adaptation
134 in alpine forest ecosystems to climate change, using landscape genomics approaches (Neale *et*
135 *al.*, 2013). Within this researcher community arose the idea to launch the genome project of
136 an important subalpine conifer species. The genome sequencing was financed by a bottom-up
137 approach among partners (<https://sfgp.faculty.ucdavis.edu/>).

138 The aim of this project was to sequence and assemble the silver fir genome, and to
139 compare this resource with other available conifer genomes. This study provides additional
140 information on the *Abies* chloroplast genome in relation to closely related taxa. A long-term
141 perspective is to identify gene regions involved in drought resistance and late flushing, which
142 are traits found to be important in Mediterranean firs (George *et al.* 2015).

143

MATERIALS AND METHODS

144

145

146 **Reference tree for genome sequencing**

147 Tissue samples for sequencing were obtained from an adult silver fir tree (AA_WSL01)
148 located in a public forest next to the institute of WSL Birmensdorf, Switzerland (47.3624°N,
149 8.4536°E). Seeds were collected directly from the selected tree in November 2016, dried at
150 ambient temperature and stored at -5°C. Fresh needles were harvested shortly after flushing in
151 May 2017. A multilocus SNP analysis across 19 Swiss populations placed the sampled tree
152 mainly in the genetic cluster of other Swiss plateau populations, with some ancestry similar
153 to populations in the Jura Mountains and in the Northern Alps (Figure S2).

154

155 **DNA preparation**

156 *Haploid megagametophyte DNA isolation for paired-end (PE) sequencing*

157 Seeds of the reference tree were incubated in tap water for 24 h at room temperature. Seeds
158 were dissected in a sterile 0.9% sodium-chloride solution under a stereo lens in an
159 environment cleaned with bleach, using micro scissors and forceps. The diploid nucellar and
160 integument tissues were carefully removed. The retained megagametophyte tissue was rinsed
161 with fresh sterile 0.9% sodium-chloride solution, immediately transferred to a 2 mL
162 Eppendorf tube and stored at -80°C. Megagametophyte tissue was lyophilized for 16 h prior
163 to extraction and homogenized for 30 s using a mixer mill (Retsch MM 300, Haan, Germany).
164 DNA extraction was performed with a customized sbeadex kit (LGC Genomics, Berlin,
165 Germany), which included chemicals and reagents as described below. 500 µL LP-PVP, 5 µL
166 Protease, 1 µL RNase and 20 µL debris capture beads were added as lysis buffer to the
167 ground tissue and the mix was incubated at 50°C and 350 rounds per minute (rpm) in a

168 heating block for 30 min. After brief centrifugation, 400 μ L cleared lysate was added to 400
169 μ L binding buffer SB and 10 μ L sbeadex beads. After 15 min binding at room temperature
170 with shaking at 850 rpm, magnetic beads were collected on a magnetic stand for 2 min, and
171 the supernatant was discarded completely. Beads were successively washed with the
172 following buffers: 400 μ L BN1, 400 μ L TN1, 400 μ L TN2, and 400 μ L PN2. Washing time
173 was 7 min for all four steps, with shaking at 850 rpm, followed by a short spin, 2 min of bead
174 collection on a magnetic stand, and careful discarding of wash buffer. DNA was finally eluted
175 in 100 μ L elution buffer AMP at 60°C and 850 rpm on a heating block for 10 min. After a
176 short spin and 3 min of magnetic bead collection on a magnetic stand, DNA was transferred
177 into a new tube, centrifuged at 21,000 x g for 2 min, and transferred without pellet into a new
178 tube.

179 DNA concentration was measured using the QuantiFluor dsDNA System (Promega,
180 Madison, WI, USA). 260/280 and 260/230 ratios were measured using a Nanodrop 1000
181 (Thermo Fisher Scientific, Waltham, MA, USA; Table S1 Supplemental Information), and
182 DNA integrity was visualized by running 5 μ L of DNA on a 1% agarose gel. Nuclear
183 microsatellites were used to test for the contamination of the haploid maternal DNA with
184 diploid DNA deriving from the surrounding tissue and to confirm the presence of only one
185 maternal haplotype (Table S2A). Because different megagametophytes from the same tree
186 represent different haplotypes, only one DNA sample from a single megagametophyte with
187 high DNA quality (260/280 ratio: 1.83, 260/230 ratio: 1.75) and quantity (3.7 μ g at 41 ng/ μ L;
188 Table S1) was used by CNAG-CRG for PE library preparation and sequencing.

189

190 ***Diploid needle DNA isolation for mate-pair (MP) sequencing***

191 Young, bright green needles of the reference tree were collected, frozen at -80 °C and
192 lyophilized for 24 h. For DNA extraction, 25 mg of tissue were ground in a 2 mL Eppendorf
193 tube with two steel balls (d = 3.1 mm) for 1.5 min, using a mixer mill MM300 (Retsch). DNA
194 was extracted with the Dneasy Plant Mini Kit (Qiagen, Hilden, Germany), starting with 600
195 μ L AP1, 1 μ L RNase and 1 μ L DX reagent. Then, DNA extraction was carried out according
196 to the manufacturer's protocol, with an additional washing step with washing buffer AW2.
197 DNA was eluted in 2x 100 μ L nuclease-free water. DNA concentration was measured using
198 QuantiFluor dsDNA System (Promega), 260/280 and 260/230 ratios were measured using a
199 Nanodrop 1000 (ThermoFisher), and DNA integrity was visualized by running 0.6 μ L of
200 DNA on a 1 % agarose gel. DNA samples were checked for contamination again using
201 nuclear microsatellite markers (Table S2A), and one sample (24.5 μ g at 136 ng/ μ L; Table S1)
202 was used for MP sequencing.

203

204 **Sequencing**

205 *Whole-genome sequencing (WGS) library preparation and sequencing*

206 Haploid DNA from the single megagametophyte was used to construct three 300 bp-insert
207 paired-end libraries at the CNAG-CRG Sequencing Unit. The short-insert PE libraries for the
208 whole-genome sequencing were prepared with KAPA HyperPrep kit (Roche-Kapa
209 Biosystems) with some modifications. In short, 1.0 μ g of genomic DNA was sheared on a
210 Covaris™ LE220 (Covaris Woburn, Massachusetts, USA) in order to reach fragment sizes of
211 ~500 bp. The fragmented DNA was further size-selected for fragment sizes of 220-550 bp
212 with AMPure XP beads (Agencourt, Beckman Coulter). The size-selected genomic DNA
213 fragments were end-repaired, adenylated and ligated to Illumina sequencing compatible
214 indexed paired-end adaptors (NEXTflex® DNA Barcodes). The adaptor-modified end library

215 was size-selected and purified with AMPure XP beads to eliminate any non-ligated adapters.
216 The ligation product was split into three samples and in three separate reactions enriched with
217 12 PCR cycles and then validated on an Agilent 2100 Bioanalyzer with the DNA 7500 assay
218 (Agilent) for size and quantity. The resulting libraries had estimated fragment sizes of 304 bp,
219 307 bp and 311 bp. These are referred to as PE300-1, PE300-2, and PE300-3 in Table 1.

220 All three libraries were sequenced in equal proportions on HiSeq 4000 (Illumina, Inc, San
221 Diego, California, USA) in paired-end mode with a read length of 2×151 bp using a HiSeq
222 4000 PE Cluster kit sequencing flow cell, following the manufacturer's protocol. Image
223 analysis, base calling and quality scoring of the run were processed using the manufacturer's
224 software Real Time Analysis (RTA 2.7.6) and followed by generation of FASTQ sequence
225 files by CASAVA.

226

227 *Mate-pair library preparation and sequencing*

228 DNA extracted from diploid needles was used to build three mate-pair (MP) libraries of
229 increasing insert size of 1,500 bp, 3,000 bp and 8,000 bp (MP1500, MP3000, MP8000).
230 Libraries were prepared using the Nextera Mate Pair Library Prep Kit (Illumina) using the
231 gel-plus protocol selecting for three different distribution sizes according to the
232 manufacturer's instructions. After fragmentation, bands of 1.5, 3 and 8 Kb were selected for
233 circularization. The following amounts of size-selected DNA were used for the circularization
234 reaction: 270 ng (1.5 Kb), 285 ng (3 Kb), and 97.4 ng (8 Kb).

235 All three MP libraries were sequenced on HiSeq2000 (Illumina) in paired-end mode with
236 a read length of 2×101 bp using TruSeq SBS Kit v4. Image analysis, base calling and quality
237 scoring of the run were processed using the manufacturer's software Real Time Analysis
238 (RTA 1.18.66.3) and followed by generation of FASTQ sequence files by CASAVA.

239

240 **Assembly**

241 *Genome assembly*

242 Given the nearly equivalent estimated fragment sizes, the reads from the three PE libraries
243 (PE300-1, PE300-2, and PE300-3) were joined into one library for assembly and collectively
244 referred to as PE300. Before assembling the genome, its size and its complexity were
245 evaluated using *k*-mer analyses. Jellyfish v2.2.0 (Marçais and Kingsford 2011) was run on the
246 sequence reads of this PE library to obtain the distribution of 17 *k*-mers. SGA preqc (Simpson
247 and Durbin 2011; Simpson 2014) was then used to estimate the mean fragment size and
248 standard deviation of the PE300 library.

249 First, an initial assembly of the PE300 reads was performed with MaSuRCA v3.2.2
250 (Zimin *et al.* 2013) using default parameters, choosing SOAPdenovo for faster contig and
251 light scaffold assembly. A *k*-mer of 105 was chosen by MaSuRCA for *de Bruijn* graph
252 construction. The initial assembly was run for 33 days on a single 48-core node (4 Intel®
253 Xeon® CPU E7-4830 v3 at 2.10GHz and 2TB of RAM) and with a maximum memory usage
254 of 1.22 TB.

255 Second, the PE300 and the three MP libraries (MP1500, MP3000, MP8000) were used to
256 scaffold the initial assembly with BESSTv2.2.5 (Sahlin *et al.* 2014). It was run with options –
257 *separate_repeats*, *-K=105* *-max_contig_overlap=115* and *-k=466*. Briefly, *-K* specifies the
258 *k*-mer size used in the *de Bruijn* graph for the input assembly to be scaffolded. As 90 % of the
259 input “contigs” were longer than 115 bp, this length was selected, instead of the default value
260 of 200 bp, as the maximum identical overlap to search (*k*). Given the fragmented input
261 assembly, the idea was to avoid using contigs smaller than the original genomic fragment.
262 Therefore, the contig size threshold for scaffolding was set to 466 bp, 10 bp greater than the

263 mean (294) plus two times the standard deviation (81) of the PE300 fragment size as
264 estimated by mapping. The scaffolded genome assembly is referred to as ABAL 1.0.
265 Moreover, an analysis of the spectra copy number (KAT; Mapleson *et al.* 2016) with default
266 $k=27$ was done to compare the PE300 library to the assembly. The KAT program is often
267 used to compare the proportion of k -mers present in the reads that end up in the final
268 assemblies. It shows how much the genome architecture agrees with the final assembly.

269

270 *Chloroplast genome assembly and annotation*

271 All of the 100 bp reads from the MP1500 library (the library with the tightest size distribution
272 and highest complexity) were mapped to the closest complete reference chloroplast sequence
273 available in NCBI, i.e. from *Abies koreana* (NC_026892.1, Yi *et al.* 2015), using BWAmem
274 (Li and Durbin 2010) in paired mode and option $-M$ to discard short split mappings. The
275 mapped reads were then extracted from the alignment using BAM2FASTQ v1.1.0 (Alpha
276 GSLaH). Both the linker sequence and the Nextera adapters present in the MP sequences
277 were removed with Cutadapt (Martin 2011). Finally, they were reversed-complemented in
278 order to obtain an artificial PE library with insert size of $1,387 \pm 327$ bp.

279 The FAST-PLAST pipeline was run producing SPAdes (Bankevich *et al.* 2012)
280 assemblies using a range of k -mers (55, 69, 87). Afterwards, Ragout (Kolmogorov *et al.*
281 2014) was used to obtain a reference-assisted assembly. In this case, *A. sibirica*
282 (NC_035067.1) was used as chloroplast reference to place and orient all the *A. alba* contigs.
283 Finally, Gapfiller (Boetzer and Pirovano 2012) was used to close gaps in the chloroplast
284 genome. The DNA diff module—from MUMMER 3.22 package (Kurtz *et al.* 2004)—was run
285 to compare the intermediate Spases assembly with the *A. koreana* (NC_026892.1) and *A.*

286 *sibirica* (NC_035067.1) complete chloroplast sequences. Finally, the annotation of the
287 chloroplast was carried out with DOGMA (Wyman *et al.* 2004).

288

289 ***Gene completeness***

290 The final nuclear assembly was evaluated for gene completeness using CEGMA v2.5 (Parra
291 *et al.* 2007), which searches for 248 ultra-conserved core eukaryotic genes (CEGs), and
292 BUSCO v3.0.2 (Simão *et al.* 2015), using 956 single-copy orthologues from plants (BUSCO
293 v1 plantae database).

294 To obtain a more comprehensive estimate of genes present in the genome assembly, the
295 STAR software package (Domin and Gingeras 2015) was used to map the genome assembly
296 with the silver fir RNA-seq produced by Roschanski *et al.* (2013) (GenBank accession
297 numbers JV134525–JV157085) as well as 12 transcriptomes originating from Mont Ventoux
298 (France) and the Black Forest (District Oberharmersbach, Germany), as reported in
299 Roschanski *et al.* (2013) and available in the Dryad Digital Repository (Roschanski *et al.*
300 2015; 2016). In addition, the transcripts from *P. taeda* (Zimin *et al.* 2014) were aligned to the
301 genome using GMAP with default options (Wu *et al.* 2016).

302

303 **Annotation**

304 ***Protein-coding gene annotation***

305 Repeats were identified, annotated and masked in the silver fir genome assembly following
306 three sequential steps. First, RepeatMasker v4.0.6 (<http://www.repeatmasker.org>) was run
307 using the Pinaceae-specific repeat library included in the RepeatMasker release. Then, repeats
308 annotated in *P. taeda* and *P. menziesii* were used in a second run of RepeatMasker. Finally,

309 *Abies alba*-specific repeats were detected with RepeatModeler and masked with
310 RepeatMasker.

311 An annotation of the genes present in the assembly was further obtained by combining
312 transcript alignments, protein alignments and *ab initio* gene predictions as follows. The
313 RNAseq reads mentioned above (JV134525– JV157085 in Roschanski et al. 2013; 2015;
314 2016) were aligned to the genome using STAR v2.5.4a (Dobin *et al.* 2013) with default
315 options, and then transcript models were generated from Stringtie (Pertea *et al.* 2015) also
316 with default options. The resulting models were given to PASA v2.2.0 (Haas *et al.* 2008)
317 together with 2,806 *A. alba* Expressed Sequence Tags (ESTs) downloaded from NCBI on
318 January 31st, 2018. Next, the TransDecoder program, which is part of the PASA package,
319 was used to detect coding regions in the PASA assemblies. A BLASTp (Altschul *et al.* 1990)
320 search was performed on the Transdecoder predictions against the Swiss-Prot database (The
321 UniProt Consortium 2017). Sequences with a complete Open Reading Frame (ORF), a
322 BLAST hit against Swiss-Prot (E-value < 1e-9), and not hitting any repeat were considered as
323 potential candidates to train gene predictors. Of this list, the 500 sequences whose length
324 differed the least from the length of their BLAST target were selected as the best candidate
325 genes and used to train the parameters for three gene predictors: GeneID v1.4 (Parra *et al.*
326 2000), Augustus v3.2.3 (Stankeet *et al.* 2006) and Glimmer (Majors *et al.* 2004). These three
327 gene predictors as well as GeneMark v2.3e (Lomsadze *et al.* 2014), which run in a self-
328 trained mode, were then run on the repeat-masked ABAL 1.0 assembly. Finally, an extra run
329 of each GeneID, Augustus and GeneMark was performed using intron data extracted from the
330 RNAseq mappings.

331 The complete Pinaceae protein sets present in PLAZA
332 (<https://bioinformatics.psb.ugent.be/plaza/versions/gymno-plaza/>) in January 2018, were

333 aligned to the repeat-masked genome using exonerate v2.4.7 (Slater and Birney 2005).
334 Moreover, all the data described above were provided as input to Evidence Modeler v1.1.1
335 (Haas *et al.* 2008) and combined into consensus coding sequence (CDS) models. These
336 models were then updated with UTRs and alternative splice isoforms with two rounds of
337 PASA updates.

338 To remove the potential presence of some bacterial genes in the genome annotation, a
339 protein-based bacterial decontamination procedure was performed on the assembly and
340 annotation. This process utilizes a BLASTp search of the annotated proteins against the
341 bacterial non-redundant protein database from NCBI to detect genes likely to belong to
342 bacteria. All the scaffolds containing more than 50% of bacterial genes and without conifer-
343 specific repeats and RNAseq mappings were removed from the assembly, resulting in the
344 final assembly ABAL 1.1.

345 Finally, to check for the presence of the chloroplast genome in the nuclear genome
346 assembly, the chloroplast assembly was mapped to ABAL 1.1 using Minimap2 (Li 2018) with
347 the parameter --asm10. Sixty-six unique mappings longer than 1 Kb were found in the
348 assembly (the longest being 18.49 Kb) but they did not meet the threshold of at least 70%
349 matches. Therefore, these regions were considered as nuclear sequence homologous to
350 chloroplast and were kept in the ABAL 1.1 assembly.

351 The proteins resulting from the structural annotation process described above were
352 functionally annotated using the Blast2GO v4.1 pipeline (Conesa *et al.* 2005) with default
353 parameters. The annotated proteins were first scanned for InterProScan patterns and profiles.
354 Next, a BLASTp search against the NCBI RefSeq database (Uniprot and Swissprot databases)
355 was performed, inheriting the functional annotations of the top-20 BLAST hits with an e-
356 value < 1e-06. Finally, Blast2GO produced a consensus annotation.

357 In addition, the software CateGORize (Zhi-Liang *et al.* 2008) was run to assign all genes
358 to the main Gene Ontology (GO) categories. The software provides the count and percentage
359 of the GO term assigned in each category. Two classification lists (slim2 and myclass2) were
360 used in the analysis. The slim2 list is a subset of gene ontology terms
361 (<http://www.geneontology.org/GO.slims.shtml>). Myclass2 classification list is based on slim2
362 with 50 additional GO term categories (Table S3). The percentages across the two
363 classification lists were visualised using the *geom_col* function of the “ggplot” package in R
364 CRAN.

365

366 ***Comparison with other conifers***

367 The summary statistics on the annotated genes were computed using a custom Python script
368 (Supplementary Material 2). The same script was applied to calculate the length of exons,
369 introns and genes in other conifer assemblies, such as *P. abies* v1.0, *P. glauca* v3.0, *P.*
370 *lambertiana* v1.5, *P. taeda* v2.0 and *P. menziesii* v1.5. The distributions of the exon, intron,
371 gene and transcript lengths across the genome were visualized using the *violinBy* function of
372 the “psych” package in R CRAN (R version 3.3.3; 2017-03-06).

373

374 **Data availability**

375 The silver fir genome assembly ABAL 1.1 is available in the TreeGenes Database under
376 <https://treegenesdb.org/FTP/Genomes/Abal/>. The following data are listed in the
377 supplementary tables: the estimation of DNA concentration (Table S1), the multi-locus
378 microsatellite genotypes of the megagametophyte and needle tissue used for sequencing
379 (Table S2A), the genotype of AA_WSL_01 for the SNP loci (Table S2B), the Gene ontology
380 (GO) term categories used to count the GO terms in *A. alba* (Table S3), the *A. alba* genome

381 annotation statistics (Table S4), the intron and exon statistics for *A. alba* and *Pseudotsuga*
382 *menziesii* (Table S5), and the count and percentage of the GO terms (Table S6). The
383 following supplementary figures are included in the supplementary file: *Abies alba*
384 distribution map (Figure S1), the location of the sampled tree AA_WSL01 along with the
385 location of the other 19 Swiss *A. abies* populations (Figure S2), plot for the comparison
386 between *Abies* chloroplast (Figure S3), boxplots of the gene distribution lengths in *A. alba*
387 (Figure S4) and in other conifers (Figure S5), distribution of the most abundant GO terms.
388 The Python script for the genome summary statistics is presented in Supplementary Material
389 2.

390

391

RESULTS AND DISCUSSION

392

Genome sequencing and genome size estimation

393

394

395

396

397

398

399

400

401

402

403

The sequencing strategy used in this project combined Illumina PE and MP libraries following a protocol similar to that used to sequence other conifer genomes (Neale *et al.* 2017). PE and MP sequencing produced a total of 1,880,827 and 765,104 Mb, respectively (Table 1). The mean fragment size of the PE300 library estimated using *SGA preqc* was 294 bp with a standard deviation of 81 bp.

The estimate of the silver fir genome size, using the distribution of 17-mers (Figure 1) is 17.36 Gb, slightly higher than previous empirical estimates of the haploid C-value of 16.19 Gb (Roth *et al.* 1997). The plot of all 17-mers present in the PE300 aggregated library that were counted and the number of distinct 17-mers (*k*-mer species) for each depth from 1 to 600 shows the existence of a considerable amount of two-, three- and four-copy repeats (17-mers) in this large genome (Figure 1). The main peak at depth 91X corresponds to unique haploid

404 sequences, while the right-most peaks at depths 182, 273, and 364 correspond to considerable
405 fractions of multi-copy repeat sequences (Figure 1).

406

407 **Genome assembly quality**

408 The silver fir genome sequence presented here accounts for 18.17 Gb, with 37 million
409 scaffolds characterized by an N50 of 14.05 Kb (Table 2). The scaffold size ranges between
410 106 bp and 297,427 bp with a mean size of 489.5 bp. The gaps constitute a total of 236.7 Mb
411 and are relatively small on average (29.3 ± 46.8 bp). The assembly size is again slightly
412 higher than the C-value of 16.19 Gb (Roth *et al.* 1997) or the *k*-mer-based estimate of 17.36
413 Gb (Figure 1). A comparison of frequency of spectra of 27-mers from the PE300 reads to the
414 final assembly using KAT (Figure 2) suggests a high level of completeness: most of 27-mers
415 in the reads belong to the haploid or main peak of the genome. Figure 2 also shows that the
416 fractions of the genome corresponding to real 2-copy (violet) and 3-copy (green) repeats were
417 successfully included in the assembly.

418 Genome completeness was estimated with three methods based on the presence of
419 conserved genes. CEGMA estimated 81.5% completeness using 248 conserved eukaryotic
420 genes. Using larger gene sets, BUSCO estimated a completeness of 49%, whereas mapping to
421 the *P. taeda* transcriptome resulted in a completeness estimate of 69%. The contiguity of the
422 silver fir assembly was also compared to those of other available conifer genome assemblies
423 (Tree Gene Database; <https://treegenesdb.org/>). The scaffold N50 (scfN50) of the silver fir
424 assembly was 14.05 Kb, almost double that of the 5.21 Kb scfN50 of the latest *P. abies*
425 assembly (Paab1.0b) and the 6.44 Kb of the *L. sibirica* assembly (Table 3). However, it is still
426 far below those of *P. lambertiana* (2,509.9 Kb), *P. glauca* (110.56 Kb), *P. taeda* (2,108.3 Kb)
427 and *P. menziesii* (372.39 Kb; Table 3).

428 The assembly completeness is estimated to be moderately high with 81.5% of the Core
429 Eukaryotic Genes as estimated by CEGMA, 65% of 956 plant orthologs as estimated by
430 BUSCO and at least 69% *P. taeda* transcripts mapping to the assembly. As each of these
431 methods are also affected by assembly fragmentation, the most likely explanation for less than
432 ideal “completeness” is that the assembly is too fragmented for a good fraction of genes to be
433 detected properly by the programs rather than the genes being truly missing from the
434 assembly. While this first draft of the silver fir genome is highly fragmented, as were earlier
435 conifer genome assemblies, likely due to the presumed density of repetitive sequences typical
436 for plant genomes (Bennetzen 2014), it represents a very valuable reference resource to the
437 community and can be used immediately to facilitate a broad spectrum of genetic and
438 genomic studies in a demographic, evolutionary, and ecological context. Given the size and
439 complexity of the silver fir genome, the low contiguity of the assembly obtained with this
440 sequencing approach was not surprising. However, a comparison of the k -mer spectra of the
441 reads used to assemble contigs (from haploid material) with their copy number in the final
442 assembly shows that we have obtained a fairly complete assembly. In fact, the majority of the
443 k -mers belonging to the main haploid peak are contained in the assembly once and only once,
444 while the peaks of double and triple k -mer depth are almost purely 2-copy and 3-copy repeats.
445 Only minor collapsing of repeats is observed. Given the haploid nature of the sample
446 (megagametophyte), we consider these repeat tails to be a real part of the genome that will
447 mainly contain repeats and, sometimes, partial genes, and might contain repeated genes.
448 Therefore, these regions were included in the assembly for higher reference completeness.

449

450 **Chloroplast assembly**

451 *De novo* chloroplast assembly, using SPADes and the *A. koreana* complete chloroplast
452 sequence as a reference for mapping, gave an assembly totaling 123,546 bp and contig N50 of
453 9,211 bp. The second reference-assisted assembly with Ragout using *A. sibirica* and Gapfiller
454 produced a single scaffold of 120,908 bp, comprised of eleven contigs (Table 2). The
455 estimated contig N50 was 15.8 Kb. Each chloroplast has its own genome (cpDNA) that for
456 most plants is formed by four parts: two large inverted repeats, one large single-copy and one
457 small single-copy region. Pinaceae chloroplast genomes lack the inverted repeats. Moreover,
458 the chloroplast genomes in Pinaceae are characterized by the presence of many small repeats
459 and are known to vary in organization (Hipkins *et al.* 1994). The cpDNA organization in
460 Pinaceae was investigated using the *Cedrus* cpDNA as reference, showing the presence of at
461 least three organization types: one similar to *Cedrus* and also found in *Picea*, another similar
462 to *Pseudotsuga*, and another similar to *Larix* (Wu *et al.* 2011). In addition to *Cedrus/Picea*,
463 *Pseudotsuga* and *Larix* organizations, another form of organization was recognized in *Abies*
464 (Tsumura *et al.* 2000). Using the DNAdiff module for genome alignment, a high collinearity
465 was observed with the *A. koreana* and *A. sibirica* complete chloroplast sequences except for a
466 region of ~45 Kb that aligns in the opposite direction to *A. koreana* due to the presence of
467 inverted repeats (Figure S3).

468 The size of the chloroplast assembly of silver fir was not only close to those of *A. sibirica* and
469 *A. koreana* (Semerikova and Semerikov 2007), as expected, but also to the 124 Kb estimated
470 in *P. abies* (Nystedt *et al.* 2013), the 124.1 Kb in *Picea sitchensis* (Coombe *et al.*, 2017), the
471 121.3 Kb in *Abies nephrolepis* (Yi *et al.* 2015) and 122.6 Kb in *L. sibirica* (Bondar *et al.*
472 2019). By using Dogma 85 protein coding genes, four rRNA genes and 39 tRNA genes have
473 been annotated. With respect to the *A. koreana* and *A. sibirica* chloroplast genomes, the *A.*

474 *alba* chloroplast assembly has four duplicated tRNAs (*trnA*-UGC, *trnI*-GAU, *trnL*-UAA and
475 *trnV*-UAC) and *trnS*-UGA has been replaced by *trnS*-CGA.

476

477 **Annotation**

478 ***Protein-coding gene annotation***

479 According to the repeat annotation performed, 78% (14.25 Gb) of the genome assembly
480 corresponds to repeats. In the non-repetitive fraction, 94,205 genes were annotated whose
481 98,227 transcripts encode 97,750 proteins (Table 4). However, the number of distinct genes is
482 inflated as many partial genes have been annotated due to the large fragmentation of the
483 assembly. Supporting this assessment, the median gene length was 558 bp, half of the genes
484 were mono-exonic and approximately half of the annotated proteins (44,646) contained only
485 partial ORFs (missing a start or stop codon). Of the 97,750 protein sequences, 39,420 (35.8%)
486 were assigned functional labels, while the rest (58,327 proteins) were analyzed with BLAST,
487 but failed to return significant hits against the RefSeq database. In total, 21,612 of the proteins
488 with complete ORFs were functionally annotated successfully.

489 Two types of gene models were used to calculate the genome annotation statistics: the
490 protein-coding genes and the full-length genes, respectively. The coding GC content was
491 46.4% in the protein coding genes and 45.2% in the full-length genes. While the number of
492 exons for the protein-coding genes was 187,740 with a mean length of 327 bp, the number of
493 introns was 89,618 (mean length: 320 bp). The number of full-length genes was 50,757 with a
494 median gene length of 804 bp. The number of exons was 118,168 with mean length of 352 bp,
495 the number of introns was 64,728 (mean length: 330 bp; Table 4, Table S5).

496 The distributions of the transcript, intron and exon lengths across the silver fir genome
497 were similar in the protein coding genes and full-length genes (Figures 3A and S4). The violin

498 plot showed a different length distribution in the low part of the violin between the two gene
499 models, due to the lower number of short genes in the full-length gene model than in all
500 genes.

501

502 ***Comparison with other conifers***

503 The comparison of silver fir genome metrics with other conifer species showed a genome size
504 similar to *P. menziesii* and *P. abies*. Moreover, the gene numbers (94,205) without filtering
505 for quality and completeness were similar to what was found in *P. abies* (70,968), *P.*
506 *lambertiana* (71,117), and *P. glauca* (102,915), but higher than in *P. menziesii* (54,830), *P.*
507 *taeda* (47,602), and *L. sibirica* (49,521). When applying a quality filter, more full-length
508 genes (50,757) were found in silver fir than high-confidence genes in *P. lambertiana*
509 (13,936), *P. glauca* (16,386), *P. abies* (28,354), and *P. menziesii* (20,616). The mean and
510 maximum intron lengths were lower than in the other conifers, while mean exon size was
511 similar to that in *P. taeda*, *P. glauca*, *P. abies* and *L. sibirica* (Table 3).

512 While the distributions of gene length across the genome were similar between silver fir
513 and *P. glauca* (Figure 3B), the mean length in *P. menziesii*, *P. taeda* and *P. lambertiana* was
514 higher than in the other conifers (Table 3). In *P. abies*, the mean gene length was close to that
515 in silver fir, whereas its distribution range was wider (Figure S5A). The density plot using
516 violin visualization confirmed these differences among species. In particular, the shape of this
517 plot showed the distribution of the genes according to their lengths and highlighted the higher
518 number of short genes in *P. abies*, *P. glauca* and silver fir than in the other conifers (Figure
519 3B). This comparison among the distribution of gene lengths estimated in silver fir with the
520 values found in the assemblies of other conifers showed some interesting results. First, the
521 genes of silver fir were on average shorter than in the other conifer species, except for *P.*

522 *glauca* (1,190 bp vs 1,330 bp; Warren *et al.* 2015) and *L. sibirica* (982 bp). However, this
523 might be an effect of the sequencing strategy used and the presence of many short scaffolds in
524 the silver fir assembly, and it will require confirmation with future improvements to the
525 genome sequence.

526 Moreover, the distribution of exon and intron lengths across the silver fir genome was
527 also compared with those found in the other fully sequenced conifers. The exon distribution
528 was similar across species (Figure S5B), with *P. menziesii* and *P. glauca* showing a slightly
529 lower mean value (Table 3). This was due to the short exons in *P. menziesii*, as it is visualized
530 in the density plot (Figure 3C). The comparison of the silver fir exons in the current study
531 with those in the other conifers showed similar values for the number, mean length and
532 maximum length of exons, as well as the total amount of exonic sequence (63.7 Mb versus the
533 mean of 50.8 Mb for all compared annotations). This result confirmed that the number and the
534 length of exons are well conserved across species (Sena *et al.* 2014). The average number of
535 exons per gene was less conserved and the smallest in silver fir (1.92) compared to all other
536 conifers (2.26-8.80). The mean number of exons per gene averaged for all seven species was
537 4.08, which is very close to the value of 3.66 predicted for species such as conifers (Table 2 in
538 Koralewski and Krutovsky 2011). Given that the average amount of exonic sequence in the
539 conifer genomes analyzed here is only 50.8 Mb, the differences in genome size among
540 conifers are presumably due in large part to the large fraction of repetitive sequences they
541 contain (Morse *et al.* 2009; Wegrzyn *et al.* 2013, 2014). Moreover, one of the major
542 components of plant genomes are the transposable elements, which may also affect the
543 evolution of the intron size (Kumar and Bennetzen 1999).

544 Silver fir intron and exon statistics were compared to *P. menziesii*, which was
545 sequenced, assembled and annotated using a similar approach (Table S5). For *P. menziesii*,

546 the genes were classified into two categories that were based on gene quality and
547 completeness (high-quality and high-quality full-length) and the counts were calculated for
548 both categories. While the numbers of exons and their means were similar in the two species
549 (187,740 for the protein-coding gene model in silver fir and 181,475 for the high-quality gene
550 model in *P. menziesii*), a lower number of introns with a lower mean size was found in silver
551 fir than in *P. menziesii* (89,618 and 145,595, respectively).

552 The distribution of intron lengths was similar across all species (Figure 3D), with
553 silver fir showing a narrower distribution range than the other conifer species (Figure S5C).
554 Although intron size has been positively correlated with genome size across eukaryotes
555 (Vinogradov 1999), this trend is not a rule for seed plants (Wan *et al.* 2018). Previous studies
556 have reported larger intron sizes in conifers than in angiosperms (Nystedt *et al.* 2013; Neale *et*
557 *al.* 2014; Guan *et al.* 2016; Sena *et al.* 2014). This difference is probably related to the high
558 percentage of repetitive sequences, which are the major component of all gymnosperm
559 genomes sequenced to date. Across gymnosperms, *Ginkgo biloba* has longer introns (Guan *et*
560 *al.* 2016) than *P. taeda*, but a smaller genome. When comparing the distribution of intron
561 lengths across genomes in several conifers, we found a similar distribution and average
562 between silver fir and *P. glauca* (311 bp vs 511 bp), with the genome size of the latter being
563 almost double (33 Gb) that of silver fir. Moreover, the smallest both mean and maximum
564 intron lengths were observed in *A. alba* and *L. sibirica* that have also the smallest genome
565 sizes, 16.19 Gb (Roth *et al.* 1997) and 12.03 Gb (Ohri and Khoshoo, 1986), respectively.

566 Another aspect related to intron length is the suggestion that the expansion of introns occurred
567 early in conifer evolution (Nystedt *et al.* 2013). This hypothesis was confirmed by the
568 comparison between orthologous introns of *P. taeda* and *G. biloba* that showed a high content
569 of repeats in long introns in both species (Wan *et al.* 2018). In addition, our analysis showed

570 that the maximum intron lengths occur in *P. taeda* and *P. lambertiana*, and their mean intron
571 length was higher than in other conifer species. The geological timescale calculated for the
572 Pinaceae showed that *Pinus* is the oldest genus across the Pinaceae, since its presence was
573 confirmed starting from the Early Cretaceous (Wang *et al.* 2000). The genus *Abies* should be
574 closer to *Pseudotsuga* than to *Picea* and *Pinus* (Wang *et al.* 2000). Nevertheless, likely due to
575 the high fragmentation of the silver fir genome sequence reported here, the estimated
576 maximum intron length in *A. alba* was only half of that estimated for *P. menziesii*.

577 The input file accounted for 462,216 GO terms that were mapped to the slim2
578 classification list categories. The total count (Table S6A) was 27,723 terms corresponding to
579 32,272 genes, of which 12,221 unique terms belonged to at least one of the 110 slim2 classes.
580 The 462,216 GO terms were mapped to the myclass2 classification list categories. The total
581 count (Table S6B) was 31,839 terms corresponding to 32,275 genes, of which 12,361 unique
582 terms belonged to at least one of the 162 myclass2 classes.

583 In both classification lists, the main categories were metabolism (11.1% and 9.7% for
584 slim2 and myclass2, respectively), catalytic activity (7.7%, 6.7%), cell (4.7%, 4.1%) and cell
585 organization (4.3%, 3.7%; Table S5, Figure S6A and Figure S6B).

586

587

588 **Conclusions and Perspectives**

589 Here, we present a draft version of the silver fir genome, which represents a first step
590 towards the full deciphering of this giga-genome in its entire complexity. This research was
591 accomplished by the Alpine Forest Genomics Network (AForGeN). The approach applied in
592 this project could serve as a model for sequencing additional plant and animal genomes. The
593 genome sequencing was financed by a bottom-up approach among partners, which could

594 possibly be a profitable strategy for many (plant) genome-sequencing initiatives in the future
595 (Twyford 2018).

596 Future research projects could utilize the draft silver fir genome as a reference to re-sequence
597 a diverse panel of trees from contrasting environments and to develop a genotyping array with
598 thousands of single-nucleotide polymorphisms (SNP). Such SNP resources will be useful in
599 many types of demographic studies and, along with the gene annotation presented here, will
600 enable genomic studies and experiments aimed at discovering those genes that are relevant for
601 particular traits (e.g. related to growth) and adaptive responses (e.g. drought tolerance).

602 **ACKNOWLEDGEMENT**

603

604 The authors thank Berta Fusté from the CNAG-CRG, Centre for Genomic Regulation for her
605 help in managing this project. Technical support from the WSL nursery staff is highly
606 appreciated. We would also thank Aleksey Zimin, Daniela Puiu and Michael Schatz for their
607 comments and advice on genome and organelle assembly. This work was in part supported by
608 grants of the National Bioinformatics Institute (INB), PRB2-ISCI (PT13/0001/0044 to JG).
609 Authors would like to thank “ELIXIR-ITA HPC@CINECA” for providing the computing
610 resources to complete some bioinformatic tasks within this project. We dedicate this paper to
611 the memory of our colleague Eric Bazin.

612

613 **AUTHOR CONTRIBUTIONS**

614 Project conception: D. Neale, B. Heinze, M. Höhn, and C. Sperisen.

615 Project design: D. Neale, M. Troggio, T. Alioto, M. Gut, F. Gugerli, B. Heinze and M. Höhn.

616 Financial support provided through: E. Bazin, B. Fady, M. Fladung, B. Fussi, D. Gömöry, S.

617 C. González-Martínez, D. Grivet, F. Gugerli, O.K. Hansen, M. Höhn, B. Heinze, K.V.

618 Krutovsky, L. Opgenoorth, G.G. Vendramin, Z. Zaya, B. Ziegenhagen and M. Westergren.

619 Lab work: C. Rellstab, S. Brodbeck, C. Sperisen, M. Gut.

620 Bioinformatic analysis: T. Alioto, L. Bianco, F. Cruz, J. Gómez Garrido, K. Ulrich, E. Mosca

621 Provided background data on the transcriptome: B. Kersten, S. Liepelt, K. Heer, and L.

622 Opgenoorth

623 Clarifying the biogeographic status of the sequenced tree: K. Csilléry and C. Rellstab

624 Manuscript preparation: E. Mosca, D. Neale, L. Bianco, M. Troggio, F. Cruz, J. Gómez-

625 Garrido, T. Alioto, F. Gugerli, C. Rellstab, S. Brodbeck, K. Csilléry, and K. Ullrich.

626 All authors read, commented on, and approved the article.

627

628 **DISCLOSURE DECLARATION**

629 The authors declare no competing interest.

630

631 **LITERATURE CITED**

632 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and D. J. Lipman, 1990 Basic local
633 alignment search tool. *Molecular Biology Journal* **215**: 403-410.

634 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., *et al.*, 2002 SPAdes: a
635 new genome assembly algorithm and its applications to single-cell sequencing. *Journal of*
636 *Computational Biology* **19**: 455-477.

637 Bennett, M. D., and I. J. Leitch, 2011 Nuclear DNA amounts in angiosperms: targets, trends
638 and tomorrow. *Annals of Botany* **107**: 467-590.

639 Bennetzen, J. L., and H. Wang, 2014 The contributions of transposable elements to the 618
640 structure, function, and evolution of plant genomes. *Annual Review of Plant Biology* **65**: 505-
641 530.

642 Boetzer, M., and W. Pirovano, 2012 Toward almost closed genomes with GapFiller. *Genome*
643 *Biology* **13**: R56.

644 Bondar, E. I., Putintseva, Y. A., Oreshkova, N. V., and K. V. Krutovsky, 2019 Siberian larch
645 (*Larix sibirica* Ledeb.) chloroplast genome and development of polymorphic chloroplast
646 markers. *BMC Bioinformatics* **20**: 38.

647 Cailleret, M., Nourtier M., Amm, A., Durand-Gillmann, M., Davi H., 2014 Drought-induced
648 decline and mortality of silver fir differ among three sites in Southern France. *Annals of*
649 *Forest Science* **71**: 643–657.

650

651 Coombe, L., Warren, R.L., Jackman, S.D., Yang, C., Vandervalk, B.P., Moore, R.A., *et al.*
652 2016 Assembly of the complete Sitka spruce chloroplast genome using 10X Genomics'
653 GemCode Sequencing Data. *PLOS ONE* 11: e0163059.

654 Cremer, E., Liepelt, S., Sebastiani, F., Buonamici, A., Michalczyk, I. M., Ziegenhagen, B.,
655 and Vendramin, G. G. 2006 Identification and characterization of nuclear microsatellite
656 loci in *Abies alba* Mill. *Molecular Ecology Notes* 6:374-376.

657 Crepeau, M. W., Langley, C. H., and K. A. Stevens, 2017 From pine cones to read clouds:
658 resc scaffolding the megagenome of sugar pine (*Pinus lambertiana*). *G3* (Bethesda, Md.), 7:
659 1563-1568.

660 Csilléry, K., Sperisen, C., Ovaskainen, O., Widmer, A., and F. Gugerli, 2018 Adaptation to
661 local climate in size, growth and phenology across 19 silver fir (*Abies alba* Mill.)
662 populations from Switzerland. bioRxiv 292540, doi: 10.1101/292540.

663 Conesa, A., Gotz, S., Garcia-Gomez, J. M., Terol, J., Talon, M., and M. Robles, 2005
664 Blast2GO: a universal tool for annotation, visualization and analysis in functional
665 genomics research. *Bioinformatics* 21: 3674-3676.

666 Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., *et al.*, 2013 STAR: ultrafast
667 universal RNA-seq aligner. *Bioinformatics* 29: 15-21.

668 Ellenberg, H. 2009. Coniferous woodland and mixed woods dominated by conifers. Pages
669 191–242. *Vegetation ecology of Central Europe*. Cambridge University Press, Cambridge,
670 UK.

671 Genomic Services Lab of HustonAlpha (2010). BAM2FASTQ version1.1.0.
672 <https://gsl.hudsonalpha.org/information/software/bam2fastq> (18 August 2010).

673 George J.P., Schueler, S., Karanitsch-Ackerl, S., Mayer, K., Klumpp, R. T., and M. Grabner,
674 2015 Inter- and intra-specific variation in drought sensitivity in *Abies* spec. and its relation
675 to wood density and growth traits. *Agricultural and Forest Meteorology* **214-215**: 430-443.

676 Grotkopp, E., Rejmánek, M., Sanderson, M. J., and T. L. Rost, 2004 Evolution of genome
677 size in pines (*Pinus*) and its life history correlates: supertree analyses. *Evolution* **58**: 1705-
678 1729.

679 Guan, R., Zhao, Y., Zhang, H., Fan, G., Liu, X., *et al.*, 2016 Draft genome of the living fossil
680 *Ginkgo biloba*. *GigaScience* **21**: 49.

681 Haas, B. J., Delcher, A. L., Mount, S. M., Wortman, J. R., Smith Jr., R. K., *et al.*, 2003
682 Improving the *Arabidopsis* genome annotation using maximal transcript alignment
683 assemblies. *Nucleic Acids Research* **31**: 5654-5666.

684 Haas, B. J., Salzberg, S. L., Zhu, W., Pertea, M., Allen, J. E., *et al.*, 2008 Automated
685 eukaryotic gene structure annotation using EVIDENCEModeler and the program to assemble
686 spliced alignments. *Genome Biology* **9**: R7.

687 Hansen, O. K., Vendramin, G. G., Sebastiani, F., and Edwards, K. J. 2005 Development of
688 microsatellite markers in *Abies nordmanniana* (Stev.) Spach and cross-species
689 amplification in the *Abies* genus. *Molecular Ecology Notes* **5**:784-787.

690 Hipkins, V. D., Krutovskii, K. V., and S. H. Strauss, 1994 Organelle genomes in conifers:
691 structure, evolution, and diversity. *Forest Genetics* **1**: 179–189.

692 Kolmogorov, M., Raney, B., Paten, B., and S. Pham, 2014 Ragout – a reference-assisted
693 assembly tool for bacterial genomes. *Bioinformatics* **30**: i302-i309.

694 Koralewski, T. E., and K. V. Krutovsky, 2011 Evolution of exon-intron structure and
695 alternative splicing. *PLOS ONE* **6**: e18055.

696 Kumar, A., and J. L. Bennetzen, 1999 Plant retrotransposons. *Annual Review of Genetics*, **33**:
697 479-532. ^[1]_{ISEP}

698 Kurtz, S., Phillippy, A., Delcher, A., Smoot, M., Shumway, M., *et al.*, 2004 Versatile and
699 open software for comparing large genomes. *Genome Biology* **5**: R12

700 Kuzmin, D. A., Feranchuk, S. I., Sharov, V. V., Cybin, A. N., Makolov, S. V., *et al.*, 2019
701 Stepwise large genome assembly approach: a case of Siberian larch (*Larix sibirica*
702 Ledeb.). *BMC Bioinformatics* **20**: 37.

703 Leitch, I. J., Soltis, D.E., Soltis, P.S., and M. D. Bennett, 2005 Evolution of DNA amounts
704 across land plants (Embryophyta). *Annals of Botany* **95**: 207-217.

705 Li, H., 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**:
706 3094-3100.

707 Li, H., and R. Durbin, 2010 Fast and accurate long-read alignment with Burrows-Wheeler
708 Transform. *Bioinformatics* **26**: 589-595.

709 Lomsadze, A., Burns, P.D., and M. Borodovsky, 2014 Integration of mapped RNA-Seq reads
710 into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Research* **42**:
711 e119.

712 Majoros, W.H., Pertea, M., and S. L. Salzberg, 2004 TigrScan and GlimmerHMM: two open
713 source ab initio eukaryotic gene-finders. *Bioinformatics* **20**: 2878-2879.

714 Mapleson, D., Garcia Accinelli, G., Kettleborough, G., Wright, J., and B. J. Clavijo, 2017
715 KAT: a K-mer analysis toolkit to quality control NGS datasets and genome assemblies,
716 *Bioinformatics* **33**: 574-576.

717 Marçais, G., and C. Kingsford, 2011 A fast, lock-free approach for efficient parallel counting
718 of occurrences of k-mers. *Bioinformatics* **27**: 764-770.

719 Martin, M., 2011 Cutadapt removes adapter sequences from high-throughput sequencing
720 reads. *EMBnet.journal* **1**: 17.

721 Morse, A.M., Peterson, D.G., Islam-Faridi, M.N., Smith, K.E., Magbanua, Z., *et al.*, 2009
722 Evolution of genome size and complexity in *Pinus*. *PLOS ONE* **4**: e4332. ^[1]_[SEP]

723 Mosca, E., Eckert, A.J., Di Pierro, E.A., Rocchini, D., La Porta, N., *et al.*, 2012 The
724 geographical and environmental determinants of genetic diversity for four alpine conifers
725 of the European Alps. *Molecular Ecology* **21**: 5530-5545.

726 Mosca, E., Gonzáles-Martínez, S. C., and D. B. Neale, 2014 Environmental versus
727 geographical molecular adaptation in two subalpine conifers. *New Phytologist* **201**: 180-
728 192.

729 Neale, D.B., Mosca, E., and E. A. Di Pierro, 2013 Alpine forest genomics network
730 (AForGeN): a report of the first annual meeting. *Tree Genetics & Genomes* **9**: 879-881.

731 Neale, D.B., Wegrzyn, J.L., Stevens, K. A., Zimin, A. V., Puiu, D., *et al.*, 2014 Decoding the
732 massive genome of loblolly pine using haploid DNA and novel assembly strategies.
733 *Genome Biology* **15**: R59.

734 Neale, D.B., McGuire, P.E., Wheeler, N.C., Stevens, K.A., Crepeau, M.W., *et al.*, 2017 The
735 Douglas-fir genome sequence reveals specialization of the photosynthetic apparatus in
736 Pinaceae. *G3: Genes, Genomes, Genetics* **9**: 3157-3167.

737 Nystedt, B., Street, N. R., Wetterbom, A., Zuccolo, A., Lin, Y.-C., *et al.*, 2013 The Norway
738 spruce genome sequence and conifer genome evolution. *Nature* **497**: 579-584.

739 Ohri, D., and T. N. Khoshoo, 1986 Genome size in gymnosperms. *Plant Systematics and*
740 *Evolution* **153**: 119-132.

741 Parra, G., Bradnam, K., and I. Korf, 2007 CEGMA: a pipeline to accurately annotate core
742 genes in eukaryotic genomes. *Bioinformatics* **23**: 1061-1067.

743 Parra, G., Blanco, E., and R. Guigo, 2000 GeneID in Drosophila. *Genome Resource* **10**: 511-
744 515.

745 Pertea, M., Pertea, G. M., Antonescu, C. M., Chang, T. C., Mendell, J. T., and S. L. Salzberg,
746 2015 StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.
747 *Nature Biotechnology* **33**: 290-295.

748 Piotti, A., Leonarduzzi, C., Postolache, D., Bagnoli, F., Spanu, I., *et al.*, 2017 Unexpected
749 scenarios from Mediterranean refugial areas: disentangling complex demographic
750 dynamics along the Apennine distribution of silver fir. *Journal of Biogeography* **44**: 1547-
751 1558.

752 Reuter, J. A, Spacek, D. V., Snyder, M. P. 2005 High-throughput sequencing technologies.
753 *Molecular Cell* **58**:586-597.

754 Roth, R., Ebert, I., and J. Schmidt, 1997 Trisomy associated with loss of maturation capacity
755 in a long-term embryogenic culture of *Abies alba*. *Theoretical and Applied Genetics* **95**:
756 353-358.

757 Roschanski, A. M., Csilléry, K., Liepelt, S., Oddou-Muratorio, S., Ziegenhagen, B., *et al.*,
758 2015 Data from: Evidence of divergent selection for drought and cold tolerance at
759 landscape and local scales in *Abies alba* Mill. in the French Mediterranean Alps. Dryad
760 Digital Repository. <https://doi.org/10.5061/dryad.t671s>.

761 Roschanski, A. M., Csilléry, K., Liepelt, S., Oddou-Muratorio, S., Ziegenhagen, B., *et al.*,
762 2016 Evidence of divergent selection for drought and cold tolerance at landscape and local
763 scales in *Abies alba* Mill. in the French Mediterranean Alps. *Molecular Ecology* **25**: 776-
764 794.

765 Roschanski, A. M., Fady, B., Ziegenhagen, B., and S. Liepelt, 2013 Annotation and re-
766 sequencing of genes from de novo transcriptome assembly of *Abies alba* (Pinaceae).
767 *Applications in Plant Sciences* **1**: 1-8.

768 Sahlin, K., Vezzi, F., Nystedt, B., Lundeberg, J., and L. Arvestad, 2014 BESST – Efficient
769 scaffolding of large fragmented assemblies. *BMC Bioinformatics* **15**: 281.

770 Semerikova, S. A., and V. L. Semerikov, 2007 The diversity of chloroplast microsatellite loci
771 in Siberian fir (*Abies sibirica* Ledeb.) and two Far East fir species *A. nephrolepis* (Trautv.)
772 Maxim. and *A. sachalinensis* Fr. Schmidt. *Genetika* **43**: 1637-1646.

773 Sena, J. S., Giguère, I., Boyle, B., Rigault, P., Birol, I., *et al.*, 2014 Evolution of gene
774 structure in the conifer *Picea glauca*: a comparative analysis of the impact of intron size.
775 *BMC Plant Biology* **14**: 95

776 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., and E. M. Zdobnov, 2015
777 BUSCO: assessing genome assembly and annotation completeness with single-copy
778 orthologs. *Bioinformatics* **31**: 3210-3212.

779 Simpson, J. T., 2014 Exploring genome characteristics and sequence quality without a
780 reference. *Bioinformatics* **30**: 1228-1235.

781 Simpson, J. T., and R. Durbin, 2012 Efficient de novo assembly of large genomes using
782 compressed data structures. *Genome Research* **22**: 549-556.

783 Slater, G. S., and E. Birney, 2005 Automated generation of heuristics for biological sequence
784 comparison. *BMC Bioinformatics* **6**: 31.

785 Stanke, M., Schoffmann, O., Morgenstern, B., and S. Waack, 2006 Gene prediction in
786 eukaryotes with a generalized hidden Markov model that uses hints from external sources.
787 *BMC Bioinformatics* **7**: 62.

788 Stevens, K. A., Wegrzyn, J. L., Zimin, A., Puiu, D., Crepeau, M., *et al.*, 2016 Sequence of the
789 sugar pine megagenome. *Genetics* **204**: 1613-1626.

790 Tsumura, Y., Suyama, Y., and K. Yoshimura, 2000 Chloroplast DNA inversion
791 polymorphism in populations of *Abies* and *Tsuga*. *Molecular Biology and Evolution* **17**:
792 1302-1312.

793 Twyford, A. D., 2018 The road to 10,000 plant genomes. *Nature Plants* **4**: 312-313.

794 Vinogradov, A. E., 1999 Intron genome size relationship on a large evolutionary scale.
795 *Journal of Molecular Evolution* **49**: 376-384.

796 Vitali, V., Büntgen, U., and J. Bauhus, 2017 Silver fir and Douglas fir are more tolerant to
797 extreme droughts than Norway spruce in south-western Germany. *Global Change Biology*
798 **23**: 5108-5119.

799 Vrška, T., Adam, D., Hort, L., Kolář, T., and D. Janík, 2009 European beech (*Fagus sylvatica*
800 L.) and silver fir (*Abies alba* Mill.) rotation in the Carpathians – a developmental cycle or a
801 linear trend induced by man? *Forest Ecology and Management* **258**: 347-356.

802 Wan, T., Liu, Z. M., Li, L. F., Leitch, A. R., Leitch, I. J., *et al.*, 2018 A genome for
803 gnetophytes and early evolution of seed plants. *Nature Plants* **4**: 82-89.

804 Wang, X.-Q., Tank, D. C., and T. Sang, 2000 Phylogeny and divergence times in Pinaceae:
805 evidence from three genomes. *Molecular Biology and Evolution* **17**: 773-781.

806 Warren, R. L., Keeling, C. I., Yuen, M. M., Raymond, A., Taylor, G. A., *et al.*, 2015
807 Improved white spruce (*Picea glauca*) genome assemblies and annotation of large gene
808 families of conifer terpenoid and phenolic defense metabolism. *Plant Journal* **83**: 189-212.

809 Wegrzyn, J. L., Liechty, J. D., Stevens K. A., Wu, L. S., Loopstra, C. A., *et al.*, 2014 Unique
810 features of the loblolly pine (*Pinus taeda* L.) megagenome revealed through sequence
811 annotation. *Genetics* **196**: 891-909.

812 Wegrzyn, J. L., Lin, B. Y., Zieve, J. J., Dougherty, W. M., Martínez-García, P. J., *et al.*, 2013
813 Insights into the loblolly pine genome: characterization of BAC and fosmid sequences.
814 *PLOS ONE* **8**: e72439.

815 Wegrzyn, J. L., Main, D., Figueroa, B., Choi, M., Yu, J., *et al.*, 2012 Uniform standards for
816 genome databases in forest and fruit trees. *Tree Genetics & Genomes* **8**: 549-557.

817 Wyman, S. K., Jansen, R. K., and J. L. Boore, 2004 Automatic annotation of organellar
818 genomes with DOGMA. *Bioinformatics* **20**: 3252-3255.

819 Wolf, H., 2003 EUFORGEN Technical Guidelines for genetic conservation and use for silver
820 fir (*Abies alba*). Rome: International Plant Genetic Resources Institute; p. 6.

821 Wu, C.-S., Lin, C.-P., Hsu, C.-Y., Wang, R.-J., and S.-M. Chaw, 2011 Comparative
822 chloroplast genomes of Pinaceae: insights into the mechanism of diversified genomic
823 organizations. *Genome Biology and Evolution* **3**: 309-319.

824 Wu, T. D., Reeder, J., Lawrence, M., Becker, G., and M. J. Brauer, 2016 GMAP and GSNAP
825 for genomic sequence alignment: enhancements to speed, accuracy, and functionality.
826 *Methods in Molecular Biology* **1418**: 283-334.

827 Yi, D.-K., Yang, J.C., So, S. K., Joo, M., Kim, D. K., *et al.*, 2015 The complete plastid
828 genome sequence of *Abies koreana* (Pinaceae: Abietoideae). *Mitochondrial DNA* **27**:
829 2351-2353.

830 Zimin, A. V., Marçais, G., Puiu, D., Roberts, M., Salzberg, S. L., and J. A. Yorke, 2013 The
831 MaSuRCA genome assembler. *Bioinformatics* **29**: 2669-2677.

832 Zimin, A., Stevens, K. A., Crepeau, M.W., *et al.* 2014 Sequencing and assembly of the 22-gb
833 loblolly pine genome. *Genetics*;196(3):875–890.

- 834 Zhi-Liang, H., Bao, J., and J. M. Reecy, 2008 CateGORizer: A web-based program to batch
835 analyze Gene Ontology classification categories. *Online Journal of Bioinformatics* **9**: 108-
836 112.
- 837 Zonneveld, B. J. M., 2012 Genome sizes of 172 species, covering 64 out of the 67 genera,
838 range from 8 to 72 picogram. *Nordic Journal of Botany* **30**: 490-502.

839 **Figure captions**

840

841 **FIGURE 1.** Distribution of 17-mers in the whole-genome sequence of *Abies alba* using raw
842 paired-end (PE) 2×151 bp reads generated from the PE300 library with 300 bp long
843 fragment inserts and estimated with Jellyfish 2.2.0 (Marçais and Kingsford 2011). The high
844 peak at very low depths is caused by sequencing errors.

845

846 **FIGURE 2.** Spectra copy number in the *Abies alba* genome ABAL 1.1. Comparison between
847 the k -mer ($k=27$) spectra of paired-end (PE) 300 2×151 bp reads generated from the PE300
848 library with 300 bp long fragment inserts and the ABAL 1.1 assembly. This stacked
849 histogram was produced with KAT (Mapleson et al. 2016) that shows the spectra copy
850 number classes along the assembly.

851

852 **FIGURE 3.** Violin plot of the distribution length of the genes, transcripts, exons and introns
853 across the *Abies alba* (*Abies_al*) high-quality genes and full-length genes (indicated as “full”;
854 **A**). The length was log10 transformed. Violin plot of the distribution lengths of genes (**B**),
855 exons (**C**) and introns (**D**) across the *Abies alba* (*A_alba*) high-quality genes and full-length
856 genes, *Pseudotsuga menziesii* (*Ps_menz*), *Picea abies* (*P_abies*), *Picea glauca* (*P_glauca*),
857 *Pinus taeda* (*P_taeda*), *Pinus lambertiana* (*P_lamb*).

858 **List of supplementary material**

859

860 **TABLE S1.** Estimation of DNA concentration, 260/280, and 260/230 ratios in the two sample
861 types (megagametophyte and needle) used for DNA extraction in *A. alba*.

862

863

864 **TABLE S2 A:** Multi-locus genotypes of the megagametophyte and needle tissue used for
865 sequencing. Samples were genotyped with 11 nuclear microsatellite markers (Hansen *et al.*,
866 2005; Cremer *et al.*, 2006). The genotyping confirms the expected haploid nature of the DNA
867 isolated from the megagametophyte. NA – not available due to PCR failure; but the
868 homozygote genotype in the diploid tissue anyway precludes the option for detecting potential
869 contamination by maternal tissue based on this locus. **B:** The genotype of AA_WSL_01 for
870 the SNP loci that were used in the STRUCTURE analysis in Appendix 2.

871

872 **TABLE S3.** Gene ontology (GO) term categories used to count the GO terms of *A. alba*.
873 GO_slim2 is an option in CateGORize software and myclass2 accounts for 50 additional
874 categories.

875

876 **TABLE S4.** *A. alba* genome annotation statistics considering two types of gene models
877 (protein coding genes and full-length genes).

878

879 **TABLE S5.** Intron (**A**) and exon (**B**) statistics for silver fir (*A. alba*) and Douglas-fir
880 (*Pseudotsuga menziesii*) gene models.

881

882 **TABLE S6.** Count and percentage (fraction) of the GO terms assigned in each category using
883 the two classification lists (**A**: slim2 and **B**: myclass2) to be complemented.

884

885 **FIGURE S1.** Distribution map of *A. alba* natural stands, compiled by the EUFORGEN
886 Network members (EUFORGEN 2009).

887

888 **FIGURE S2.** (**A**) Location of the 19 sampled Swiss populations and tree AA_WSL01.
889 Modified after Csilléry et al. (2018). (**B**) The log-likelihood from Structure runs with $K = 2$ to
890 $K=10$. (**C**) Ancestry proportions of AA_WSL01 and the 19 genotyped Swiss populations for
891 $K=3$ and $K=4$.

892

893 **FIGURE S3.** Plot produced with DNAdiff for the comparison between *A. alba* and *A.*
894 *sibirica* chloroplasts (**A**) and *A. alba* and *A. koreana* chloroplasts (**B**).

895

896 **FIGURE S4.** Boxplots of the distribution lengths of the genes, transcripts, exons and introns
897 across the *A. alba* high-quality genes and full-length genes (indicated as “full”). The
898 distribution is log10 transformed.

899

900 **FIGURE S5.** Boxplots of the distribution lengths of the genes (**A**), exons (**B**), and introns (**C**)
901 across the *Abies alba* (A_alba) high-quality genes and full-length genes (indicated as “full”),
902 *Pseudotsuga menziesii* (Ps_menz), *Picea abies* (P_abies), *Picea glauca* (P_glauca), *Pinus*
903 *taeda* (P_taeda), *Pinus lambertiana* (P_lamb).

904

905 **FIGURE S6.** Distribution of the most abundant Gene Ontology (GO) terms assigned to the *A.*
906 *alba* genome using slim2 categories (**A**) and myclass2 categories (**B**). The percentage
907 (fraction) of the term assigned in each category is represented only for values > 0.2%. All
908 categories are given in Table S3, all count and percentages in Table S6.

909
910
911

TABLE 1 Summary of the raw data for Illumina paired-end (PE) and mate-pair (MP) libraries for whole-genome sequencing of *Abies alba*.

Library	Read length (bp)	Insert size (kb)	Mean fragment size (bp)	Read Pairs (million)	Yield (Mb)	Coverage	Avg. Phix Error R1 (%)	Avg. Phix Error R2 (%)
PE300-1	2 x 151	-	304	3,274	989,029	57.103	0.646	0.908
PE300-2	2 x 151	-	307	1,886	569,617	32.888	0.883	1.126
PE300-3	2 x 151	-	312	1,066	322,181	18.602	0.768	1.081
MP1500	2 x 101	1.5	-	1,255	253,529	14.638	0.214	0.32
MP3000	2 x 101	3	-	1,277	257,985	14.895	0.214	0.32
MP8000	2 x 101	8	-	1,255	253,590	14.641	0.214	0.32
Total PE				6,226	1,880,827	108.593		
Total MP				3,787	765,104	44.175		

912 **TABLE 2** Summary statistics for the *Abies alba* whole-genome assembly version 1.1 (ABAL 1.1) and
 913 chloroplast assembly.
 914

Genome	Feature	
Nuclear	Number of contigs	45,280,944
	Number of scaffolds	37,192,295
	Mean GC%	39.34
	Total length (Mb)	18,167
	Minimum scaffold length (bp)	106
	Maximum scaffold length (bp)	297,427
	Mean scaffold length (bp)	488.50
	Median scaffold length (bp)	115
	Contig N50 (bp)	2,477
	Scaffold N50 (bp)	14,051
Chloroplast	Total length (bp)	120,908
	Number of contigs	11
	Number of scaffolds	1
	Contig N50 (bp)	15,758

915

916 **TABLE 3** Comparison of genome summary metrics from *Abies alba* and other sequenced conifer genomes
 917 (version numbers in parentheses).
 918

Genome summary metric	<i>Abies alba</i> (1.0)	<i>Pseudotsuga menziesii</i> (1.5)	<i>Pinus taeda</i> (2.0)	<i>Pinus lambertiana</i> (1.5)	<i>Picea glauca</i> (3.0)	<i>Picea abies</i> (1.0)	<i>Larix sibirica</i> (1.0) ^a
Total length (Mb)	18,167	15,700	20,613	31,000	32,795	19,600	12,340
N50 scaffold (Kb)	14.05	372.39	2,108.3	2,509.9	110.56	5.21	6.44
N of genes	94,205	54,830	47,602	71,117 ^c	102,915	70,968	49,521
N of full-length genes	50,757	20,616	NA	13,936 ^c	16,386 ^b	28,354 ^d	32,482
N of exons	181,168	181,475	166,465	153,111	232,182	178,049	151,838
N of introns	64,728	145,595	108,809	121,858	124,951	107,313	101,675
Mean gene length (bp)	1,190	10,510	9,066	40,820	1,330	2,427	982
Mean exon length (bp)	352	231	320	241	320	312	324
Mean intron length (bp)	311	2,301	3,004	10,164	511	1,017	353
Maximum exon length (bp)	6,300	8,037	4,946	8,003	9,568	6,068	10,268
Maximum intron length (bp)	36,015	182,831	408,800	805,500	44,116	68,269	10,154
Exons per gene	1.92	8.80	3.50	5.25	2.26	3.78	3.03
Total exonic length	6.4x10 ⁶	4.2x10 ⁶	5.3x10 ⁶	1.8x10 ⁶	7.4x10 ⁶	5.6x10 ⁶	4.9x10 ⁶

919 For the gene annotation and the definition of the “full-length genes” different approaches were used across
 920 species. The scaffold N50 (scfN50) was calculated on the unshuffled assemblies and discarding scaffolds shorter
 921 than 200 bp.
 922

923 ^a Kuzmin *et al.* 2019; K.V. Krutovsky, personal communication

924 ^b high confidence set (Warren *et al.* 2015; PG29 v3) and scaffold N50 calculated using sequences ≥ 500 bp: N50
 925 is 71.5 Kb if considering both clones (WS77111)

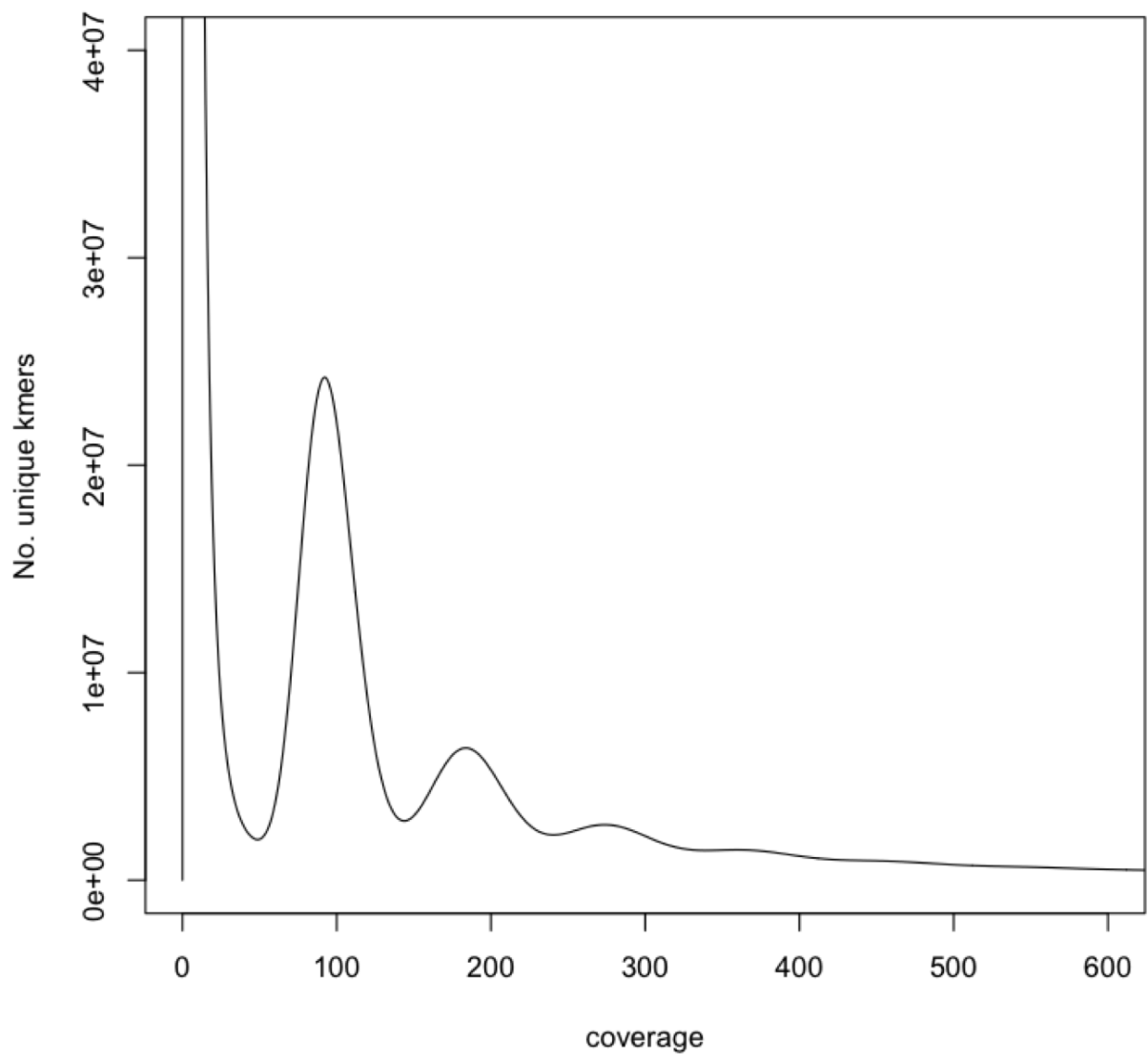
926 ^c low-quality and high-quality gene models from *Pinus lambertiana* v.1 (Stevens *et al.* 2016), the other were
 927 calculated on *Pinus lambertiana* v1.5 (Crepeau *et al.* 2017),

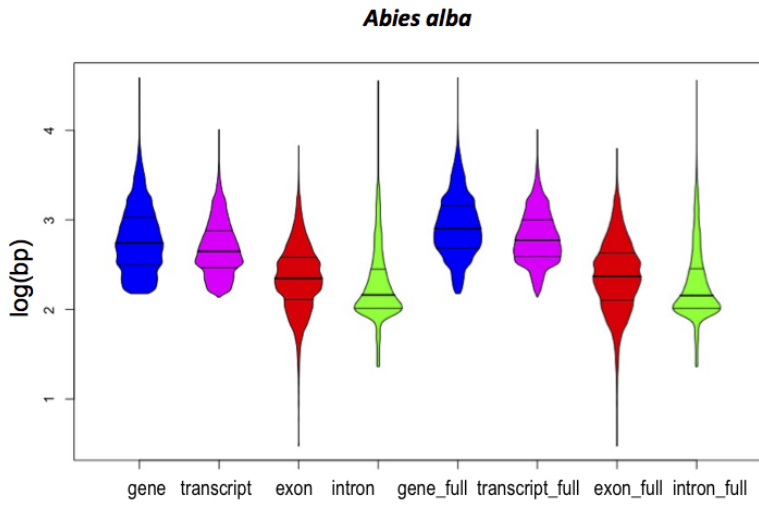
928 ^d high confidence (Nystedt *et al.* 2013)
 929

930 **TABLE 4** Genome annotation statistics for *Abies alba* considering two types of gene models (protein coding
 931 genes and full-length genes). All statistics are given in Table S3.
 932

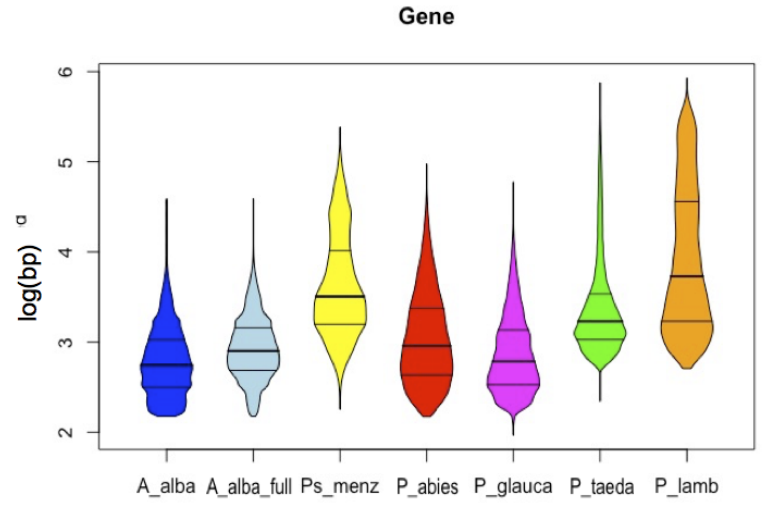
Features	Protein-coding genes	Full-length genes
Number of genes	94,205	50,757
Median gene length (bp)	558	804
Number of transcripts	98,227	53,487
Median transcript length (bp)	445	597
Number of exons	187,740	181,168
Coding GC content	46.4%	45.15%
Median exon length (bp)	224	237
Number of introns	89,618	64,728
Median intron length (bp)	146	145
Exons/transcript	2.00	2.32
Transcripts/gene	1.04	1.05

933

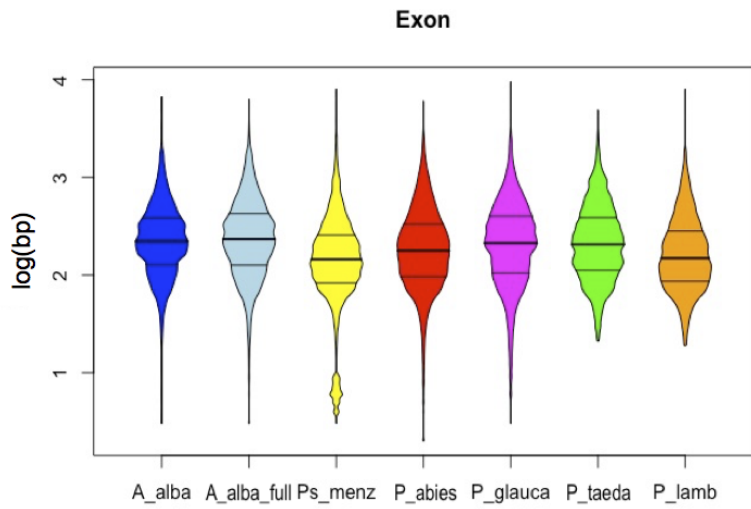




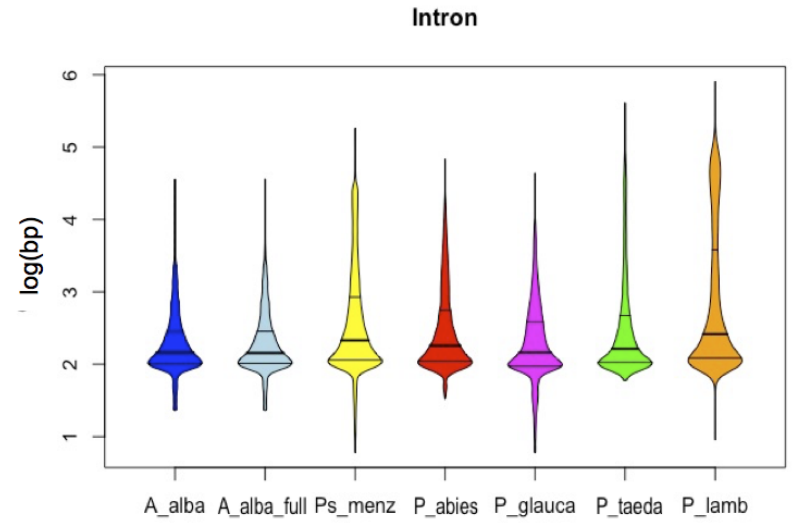
A



B



C



D