frontiers

# A chromosome-level genome assembly of the European Beech (*Fagus sylvatica*) reveals anomalies for organelle DNA integration, repeat content and distribution of SNPs

1

2  **Bagdevi Mishra[1,2], Bartosz Ulaszewski[3], Joanna Meger[3], Jean-Marc Aury[4], Catherine
3  Bodénès[5], Isabelle Lesur-Kupin[5,6,7], Markus Pfenninger[1], Corinne Da Silva[4], Deepak K
4  Gupta[1,2,8], Erwan Guichoux[5], Katrin Heer[10], Céline Lalanne[5], Karine Labadie[4], Lars
5  Opgenoorth[7], Sebastian Ploch[1], Grégoire Le Provost[5], Jérôme Salse[9], Ivan Scotti[10], Stefan
6  Wötzel[1,2], , Christophe Plomion[5], Jaroslaw Burczyk[3], Marco Thines[1,2,8***

7  [1] Senckenberg Biodiversity and Climate Research Centre (BiK-F), Senckenberg Gesellschaft für
8  Naturforschung, Senckenberganlage 25, D-60325 Frankfurt am Main, Germany

9  [2] Goethe University, Department for Biological Sciences, Institute of Ecology, Evolution and
10  Diversity, Max-von-Laue-Str. 9, D-60438 Frankfurt am Main, Germany

11  [3] Kazimierz Wielki University, Department of Genetics, ul. Chodkiewicza 30, 85-064 Bydgoszcz,
12  Poland

13  [4] Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université
14  Paris-Saclay, F-91057, Evry, France

15  [5] INRAE, Univ. Bordeaux, BIOGECO, F-33610 Cestas, France

16  [6] HelixVenture, F-33700, Mérignac, France

17  [7] Philipps University Marburg, Faculty of Biology, Plant Ecology and Geobotany, 35043, Marburg,
18  Germany

19  [8] LOEWE Centre for Translational Biodiversity Genomics (TBG), Georg-Voigt-Str. 14-16, D-60325
20  Frankfurt am Main (Germany)

21  [9] INRAE, UCA, GDEC, F-63100 Clermont-Ferrand, France

22  [10] INRAE, URFM, F-84914, Avignon, France

23  Catherine Bodénès[5], Isabelle Lesur-Kupin[5,6,10] , Jean-Marc Aury[9], ,

24

25  **\* Correspondence:**
26  Marco Thines
27  m.thines@thines-lab.eu

**Keywords: Chromosomes, *Fagaceae*, genome architecture, genomics, Hi-C, repeat elements, SNPs**

## Abstract

The European Beech is the dominant climax tree in most regions of Central Europe and valued for its ecological versatility and hardwood timber. Even though a draft genome has been published recently, higher resolution is required for studying aspects of genome architecture and recombination. Here we present a chromosome-level assembly of the more than 300 year-old reference individual, Bhaga, from the Kellerwald-Edersee National Park (Germany). Its nuclear genome of 541 Mb was resolved into 12 chromosomes varying in length between 28 Mb and 73 Mb. Multiple nuclear insertions of parts of the chloroplast genome were observed, with one region on chromosome 11 spanning more than 2 Mb of the genome in which fragments up to 54,784 bp long and covering the whole chloroplast genome were inserted randomly. Unlike in *Arabidopsis thaliana*, ribosomal cistrons are present in *Fagus sylvatica* only in four major regions, in line with FISH studies. On most assembled chromosomes, telomeric repeats were found at both ends, while centromeric repeats were found to be scattered throughout the genome apart from their main occurrence per chromosome. The genome-wide distribution of SNPs was evaluated using a second individual from Jamy Nature Reserve (Poland). SNPs, repeat elements and duplicated genes were unevenly distributed in the genomes, with one major anomaly on chromosome 4. The genome presented here adds to the available highly resolved plant genomes and we hope it will serve as a valuable basis for future research on genome architecture and for understanding the past and future of European Beech populations in a changing climate.

## 1. Introduction

Many lowland and mountainous forests in Central Europe are dominated by the European Beech (*Fagus sylvatica*) (Durrant et al., 2016). This tree is a shade-tolerant hardwood tree that can survive as a sapling in the understorey for decades until enough light becomes available for rapid growth and maturation (Wagner et al., 2010; Ligot et al., 2013). Beech trees reach ages of 200-300 years, but older individuals are known e.g. from suboptimal habitats, especially close to the tree line (Di Filippo et al., 2012). Under optimal water availability, European Beech is able to outcompete most other tree species, forming monospecific stands (Leuschner et al., 2006), but both stagnant soil water and drought restrict its presence in natural habitats (Jump at al., 2006; Geßler at al., 2007). Particularly, dry summers, which have recently been observed in Central Europe and that are predicted to increase as a result of climate change (Coumou and Rahmstorf, 2012; Spinoni at al., 2015), will intensify climatic stress as already now severe damage has been observed (Geßler at al., 2007; Reif at al., 2017). In order to cope with this, human intervention in facilitating regeneration of beech forests with more drought-resistant genotypes might be a useful strategy (Rose et al., 2009; Bolte and Degen, 2010). However, for the selection of drought-resistant genotypes, whole genome sequences of trees that thrive in comparatively dry conditions and the comparison with trees that are declining in drier conditions are necessary to identify genes associated with tolerating these adverse conditions (Pfenninger et al., 2020). Such genome-wide association studies rely on well-assembled reference genomes onto which genome data from large-scale resequencing projects can be mapped (e.g. (Atwell et al., 2010)).

Due to advances in library construction and sequencing, chromosome-level assemblies have been achieved for a variety of genomes from various kingdoms of live, including animals (Michael and VanBuren, 2020; Priest at al., 2020; Rhie at al., 2020). While the combination of short- and long-read sequencing has brought about a significant improvement in the assembly of the gene space and regions with moderate repeat-element presence, chromosome conformation information libraries, such as Hi-C (Lieberman-Aiden et al., 2009), have enabled associating scaffolds across highly repetitive regions, enabling the construction of super-scaffolds of chromosomal scale (e.g. (Yin et al., 2020)). Recently, the first chromosome-level assemblies have been published for tree and shrub species, e.g. the tea tree (*Camellia sinensis* (Chen et al., 2020)), loquat (*Eriobotrya japonica* (Jiang et al., 2020)), walnut (*Juglans regia* (Marrano et al., 2020)), Chinese tupelo (*Nyssa sinensis* (Yang et al., 2019)), fragrant rosewood (*Dalbergia odorifera* (Hong et al., 2020)), wheel tree (*Trochodendron aralioides* (Strijk at. Al., 2019)), azalea (*Rhododendron simsii* (Yang et al., 2020)), agarwood tree (*Aquilaria sinensis* (Nong et al., 2020)), and tea olive (*Osmanthus fragrans* (Yang et al., 2018)). However, such resources are currently lacking for species of the *Fagaceae*, which includes the economically and ecologically important genera *Castanea*, *Fagus*, and *Quercus* (Kermer at al., 2012). For this family, various draft assemblies have been published (Sork et al., 2016; Martínez-García et al., 2016; Plomion et al., 2016), including European Beech (Mishra et al., 2018), but none is so far resolved on a chromosome scale. To achieve this, we have sequenced the genome of the more than 300 year-old beech individual, Bhaga, from the Kellerwald-Edersee National Park (Germany), and compared it to an individual from the Jamy Nature Reserve (Poland), to get first insights into the genome architecture and variability of *Fagus sylvatica*.

## 2. Materials and Methods

2.1. Sampling and processing

2.1.1 Reference genome

The more than 300 year-old beech individual Bhaga (Fig. 1) lives on a rocky outcrop on the edge of a cliff in the Kellerwald-Edersee National Park in Hesse, Germany (51°10'09"N 8°57'47"E). Dormant buds were previously collected for the extraction of high molecular weight DNA and obtaining the sequence data described in Mishra et al. (2018). The same tree was sampled again in February 2018 for obtaining bud samples for constructing Hi-C libraries. Hi-C libraries construction and sequencing was done by a commercial sequencing provider (BGI, Hong Kong, China). For an initial assessment of genome variability and to obtain its genome sequence, Illumina reads derived from the Polish individual, Jamy, reported in Mishra et al. (Mishra et al., 2021a), were used.

2.1.2 Progeny trial and linkage map construction

For a progeny trial establishment seeds were sampled from a single mother tree (accession MSSB). About 1,000 beechnuts were collected during two successive campaigns in the fall 2013 and 2016 using a net under the mother tree located in the southern range of the species in the south-west of France (Saint- Symphorien 44° 25' 41.138" N 0° 29' 23.125" W). Seeds were germinated and raised the following springs at the National Forest Office nursery in Guémené-Penfao (47° 37' 59.99" N -1° 49' 59.99" W) and then planted at the Nouzilly (47° 32′ 36″ N  0° 45′ 0″ E) experimental unit PAO of INRAE in February 2017 (537 saplings corresponding to the 1st campaign, used for the paternity

114 reconstruction) and at the National Forest Office nursery in Guémené-Penfao in January 2019 (429
115 saplings corresponding to the 2$^{nd}$ campaign, used for linkage mapping). for relatedness assessment
116 among the half-sib progeny of MSSB, young leaves after bud burst were sampled from saplings in
117 the nursery in spring 2014 (1$^{st}$ campaign) and 2017 (2$^{nd}$ campaign), immediately frozen in dry ice and
118 then stored at -80°C before subsequent genetic analyses. Likewise, leaves were sampled on the
119 mother tree and 19 surrounding adult trees (expected fathers). Nuclear DNA was extracted
120 individually from 10 mg of tissue using the DNeasy Plant Mini Kit (QIAGEN, DE) following the
121 manufacturer's instructions. DNA concentration was measured on a ND-8000 NanoDrop
122 spectrophotometer (Thermo Scientific, Wilmington, USA). For additional transcriptome construction
123 a total of six different organs were sampled on the MSSB accession, including: two types of buds
124 (quiescent buds and swelling buds just before bud break) during dormancy release the 15th of March
125 2017, male flowers and female flowers collected the 3rd of May 2017, leaves and xylem collected the
126 28th of June 2017. Each organ was immediately flash-frozen in liquid nitrogen and stored at -80°C
127 before RNA extraction. For short read sequencing (Illumina), total RNA was extracted from these six
128 samples following the procedure described in Le Provost et al. (2007). Residual genomic DNA was
129 removed before purification using DNase RQ1 (Promega, Madisson, WI, USA) according to the
130 manufacturer's instructions. The quantity and the quality of each extract was determined using an
131 Agilent 2100 Bioanalyser (Agilent Technologies, Inc., Santa Clara, CA, USA). For long read
132 sequencing (Oxford Nanopore Technologies) total RNA was extracted as described above and
133 depleted using the Ribo-Zero rRNA Removal Kit Plant Leaves (Illumina, San Diego, CA, USA).
134 RNA was then purified and concentrated on a RNA Clean Concentrator™-5 column (Zymo
135 Research, Irvine, CA, USA).
136 For the linkage mapping, vegetative buds from the individuals from the first and second
137 campaign were sampled on the 28$^{th}$ of February 2018 in Nouzilly at the ecodormancy stage from 200
138 genotypes (i.e. 200 half-sibs that constitute the mapping population) and were frozen on dry ice and
139 then stored at -80°C. RNA was extracted from bud scale-free leaves following the procedure
140 described above. These 200 genotypes included two relatively large full-sib families comprising 49
141 full-sibs (family MSSBxSSP12) and 36 full-sibs (family MSSBxMSSH) (see results section).
142

143 2.2. Chromosomal pseudo-molecules and their annotation

144 2.2.1 Building of chromosomal pseudo-molecules using Hi-C reads

145 The previous scaffold-level assembly was constructed with Illumina shotgun short reads and PacBio
146 long reads (Mishra et al., 2018). For a chromosome-level assembly, intermediate results from the
147 previous assembly were used as the starting material. Sequence homology of the 6699 scaffolds
148 generated from the DBG2OLC hybrid assembler (Ye et al., 2016), to the separately assembled
149 chloroplast and mitochondria of beech, were inferred using blast v2.10.1 (Altschul et al., 1990). All
150 scaffolds that match in full length to any of the organelle with identity > 99 % and gaps and/or
151 mismatches ≤ 3 were discarded. The remaining 6657 scaffolds along with Hi-C data (116 Mb) were
152 used in ALLHiC (Zhang et al., 2019) for building the initial chromosome-level assembly. The
153 cleaned Illumina reads were aligned to the initial assembly using Bowtie2 software (Langmead
154 and Salzberg, 2012) and then, sorted and indexed bam files of the concordantly aligned read pairs for
155 all the sequences were used in Pilon (Walker et al., 2014) to improve the correctness of the assembly.
156 The final assemblies for Bhaga and Jamy were deposited under the accession numbers PRJEB43845
157 and PRJNA450822, respectively.

158  The completeness of the assembly was evaluated with plant-specific (viridiplantae_odb10.2019-11-
159  20) and eudicot-specific (eudicots_odb10.2019-11-20) Benchmarking Universal Single-Copy
160  Orthologs (BUSCO v4.1.4) (Seppey et al., 2019).

161      2.2.2. Gene prediction

162  Cleaned transcriptomic Illumina reads (minimum read length: 70; average read quality: 25 and read
163  pairs containing no N) were aligned to the assembly using Hisat (Kim et al., 2015) in order to
164  generate splice-aware alignments. The sorted and indexed bam file (samtools, v1.9 (Li et al., 2009))
165  of the splice alignments was used in "Eukaryotic gene finding" pipeline of OmicsBox (Accessed
166  March 3, 2020) which uses Augustus (Stanke and Morgenstern, 2005) for gene prediction. For
167  prediction, few parameters were changed from the default values. Minimum intron length was set to
168  20 and minimum exon length was set to 200 and complete genes (with start and stop codon) of a
169  minimum of 180 bp length were predicted, by choosing *Arabidopsis thaliana* as the closest organism.

170      2.2.3. Assessment of the gene space

171  The protein sequences of the PLAZA genes for *A. thaliana*, *Vitis vinifera*, and *Eucalyptus grandis*
172  were downloaded from plaza v4.5 dicots (Accessed October 21, 2020) dataset and were used along
173  with the predicted proteins from our assembly to make protein clusters using cd-hit v.4.8.1 (Li and
174  Godzik, 2006; Fu et al., 2012). The number of exons per genes was assessed and compared to the
175  complete coding genes from *A. thaliana*, *Populus trichocarpa*, and *Castanea mollissima*, in line with
176  the comparison made in the scaffold level assembly (Mishra et al., 2018).

177      2.2.4. Functional annotation of genes

178  The predicted genes were translated into proteins using transeq (EMBOSS:6.6.0.0 (Rice et al., 2000))
179  and were queried against the non-redundant database from NCBI (downloaded on 2020-06-24) using
180  diamond (v0.9.30) software (Buchfink and Xie, 2015) to find homology of the predicted proteins to
181  sequences of known functions. For prediction of protein family membership and the presence of
182  functional domains and sites in the predicted proteins, Interproscan v5.39.77 (Jones et al., 2014) was
183  used. Result files from both diamond and Interproscan (in Xml format) were used in the blast2go
184  (Götz et al., 2008) module of OmicsBox and taking both homology and functional domains into
185  consideration, the final functional annotations were assigned to the genes. The density of coding
186  space for each 100 kb region stretch was calculated for all the chromosomes.

187      2.2.5. Repeat prediction and analysis

188  A repeat element database was generated using RepeatScout (v1.0.5) (Price et al., 2005), which was
189  used in RepeatMasker (v4.0.5) (Smit and Hubley, 2007) to predict repeat elements. The predicted
190  repeat elements were further filtered on the basis of their copy numbers. Those repeats represented
191  with at least 10 copies in the genome were retained as the final set of repeat elements of the genome.
192  Repeat fractions per 100 kb region for each of the chromosomes were calculated for accessing
193  patterns of repeat distribution over the genome.

194  In a separate analysis, repeat elements present in *Fagus sylvatica* were identified by a combination of
195  homology-based and de novo approaches using RepeatModeler 2.0 (Flynn et al., 2020) and
196  RepeatMasker v. 4.1.1 (Tarailo-Graovac and Chen, 2009). First, we identified and classified
197  repetitive elements de novo and generated a library of consensus sequences using RepeatModeler 2.0

198 (Flynn et al., 2020). We then annotated repeats in the assembly with RepeatMasker 4.1.1 (Tarailo-
199 Graovac and Chen, 2009) using the custom repeat library generated in the previous step.

### 2.2.6. Telomeric and Centromeric repeat identification

201 Tandem repeat finder (TRF version 4.0.9) (Benson, 1999) was used with parameters 2, 7, 7, 80, 10,
202 50 and 500 for Match, Mismatch, Delta, PM, PI, Minscore and MaxPeriod, respectively (Marrano et
203 al., 2020), and all tandem repeats with monomer length up to 500 bp were predicted. Repeat
204 frequencies of all the monomers were plotted against the length of the monomers to identify all high-
205 frequency repeats. As the repeats were fetched by TRF program with different start and end positions
206 and the identical repeats were falsely identified as different ones, the program MARS (Ayad and
207 Pissis, 2017) was used to align the monomers of the different predicted repeats, and the repeat
208 frequencies were adjusted accordingly. The chromosomal locations of telomeric and centromeric
209 repeats were identified by blasting the repeats to the chromosomes. For confirmation of centromeric
210 locations, pericentromeres of *A. thaliana* were blasted against the chromosomes of Bhaga.

### 2.2.7. Organelle integration

212 Separately assembled chloroplast (Mishra et al. 2021a) and mitochondrial (Mishra et al. 2021b)
213 genomes were aligned to the genomic assembly using blastn with an e-value cut-off of 10e-10.
214 Information for different match lengths and different identity cut-offs were tabulated and analysed.
215 Locations of integration into the nuclear genome were inferred at different length cut-offs for
216 sequence homology (identity) equal to or more than 95%. The number of insertions per non-
217 overlapping window of 100 kb was calculated separately for both organelles.

### 2.2.8. SNP identification and assessment

219 The DNA isolated from the Polish individual Jamy was shipped to Macrogen Inc. (Seoul, Rep. of
220 Korea) for library preparation with 350 bp targeted insert size using TruSeq DNA PCR Free
221 preparation kit (Illumina, USA) and sequencing on HiSeq X device (Illumina, USA) using PE-150
222 mode. The generated 366,127,860 raw read pairs (55.3 Gb) were processed with AfterQC v 0.9.1
223 (Chen et al., 2017) for quality control, filtering, trimming and error removal with default parameters
224 resulting in 54.12 Gbp of high-quality data. Illumina shotgun genomic data from Jamy was mapped
225 to the chromosome-level assembly using stringent parameters (--very-sensitive mode of mapping) in
226 bowtie2 (Li, 2011). The sam formatted output of Bowtie2 was converted to binary format and sorted
227 according to the coordinates using samtools version 1.9 (Li et al., 2009). SNPs were called from the
228 sorted mapped data using bcftools (version: 1.10.2) (Li, 2011) call function. SNPs were called for
229 only those genomic locations with sequencing depth $\geq$ 10 bases. All locations 3 bp upstream and
230 downstream of gaps were excluded. For determining heterozygous and homozygous states in Bhaga,
231 sites with more than one base called and a ratio between the alternate and the reference allele of $\geq$
232 0.25 and < 0.75 in were considered as heterozygous SNP. Where the ratio was $\geq$ 0.75, the position
233 was considered homozygous. In addition, homozygous SNPs were called by comparison to Jamy,
234 where the consensus base in Jamy was different than in Bhaga and Bhaga was homozygous at that
235 position. SNP density was calculated for each chromosome in 100 kb intervals.

### 2.2.9. Genome browser

237 A genome browser was set up using JBrowse v.1.16.10 (Buels et al., 2016). Tracks for the predicted
238 gene model, annotated repeat elements were added using the gff files. Separate tracks for the SNP

239 locations and the locations of telomere and centromere were added as bed files. A track depicting the
240 GC content was also added. The genome browser can be accessed from http://beechgenome.net.

241        2.3. Pedigree reconstruction

242        2.3.1. SNP assay design and genotyping for relatedness assessment among half-sibs

243 We used a multiplexed assay using the MassARRAY® MALDI-TOF platform (iPLEX MassArray,
244 Agena BioScience, USA) to genotype the mother tree (MSSB), its half-sib progeny from the 1[st]
245 campaign and 19 putative fathers. PCR and extension primers were designed from flanking
246 sequences (60pb of either side) of 40 loci (Supplementary file 5) available from Lalagüe et al. (2014)
247 and Ouayjan and Hampe (2018). Data analysis was performed with Typer Analyzer 4.0.26.75 (Agena
248 BioScience). We filtered out all monomorphic SNPs, as well as loci with a weak or ambiguous signal
249 (i.e., displaying more than three clusters of genotypes or unclear cluster delimitation). Thirty-six
250 SNPs were finally retained for the paternity analysis.

251        2.3.2. Sibship assignment

252 Paternity analysis was carried out using Cervus 3.0 software (Kalinowski et al., 2007, Marshall et al.
253 1998) to check the identity of the maternal parent and identify the paternal parent among 19
254 candidate fathers growing in the neighbourhood of mother tree MSSB. Cervus was run assuming a
255 0.1% genotyping error rate. The pollen donor of each offspring was assigned by likelihood ratios
256 assuming the strict confidence criterion (95%). We performed simulations with the following
257 parameters: number of offspring genotypes = 100 000, number of candidate fathers = 19, mistyping
258 rate = 0.01 and proportion of loci typed = 0.9755. Zero mismatch was allowed for each offspring and
259 the supposed father. The Cervus selfing option was used because self-pollination may occur.

260        2. 4. Unigene set construction

261        2.4.1. Library construction and sequencing

262 Six Illumina RNA-Seq libraries (one for each organ) were constructed from 500ng total RNA using
263 the TruSeq Stranded mRNA kit (Illumina, San Diego, CA, USA), which allows for mRNA strand
264 orientation (the orientation of sequences relative to the antisense strand is recorded). Each library was
265 sequenced using 151 bp paired end reads chemistry on a HS4000 Illumina sequencer.

266 One Nanopore cDNA library was also prepared from entire female flowers RNA. The cDNA library
267 was obtained from 50 ng RNA according to the Oxford Nanopore Technologies (Oxford Nanopore
268 Technologies Ltd, Oxford, UK) protocol "cDNA-PCR Sequencing (SQK-PCS108)" with a 14 cycles
269 PCR (6 minutes for elongation time). ONT adapters were ligated to 190 ng of cDNA. The Nanopore
270 library was sequenced using a MinION Mk1b with R9.4.1 flowcells.

271        2.4.2. Bioinformatic analysis

272 Short-read RNA-Seq data (Illumina) from the six tissues were assembled using Velvet (Zerbino et
273 al., 2010) 1.2.07 and Oases (Schulz et al., 2012) 0.2.08, using a k-mer size of 63 bp. Reads were
274 mapped back to the contigs with BWA-mem (Li et al., 2009) and the consistent paired-end reads
275 were selected. Chimeric contigs were identified and splitted (uncovered regions) based on coverage
276 information from consistent paired-end reads. Moreover, open reading frames (ORF) and domains
277 were searched using respectively TransDecoder (Haas et al., 2013) and CDDsearch (Marchler-Bauer

7

278    et al., 2011). We only allowed breaks outside ORF and domains. Finally, the read strand information
279    was used to correctly orient the RNA-seq contigs.

280    Long-read RNA-Seq data (Oxford Nanopore Technologies) from female flowers were corrected
281    using NaS (Madoui et al., 2015) with default parameters.

282    Contigs obtained from short reads as well as corrected long reads were then aligned on a draft version
283    of MSSB genome assembly (unpublished) using BLAT (Kent, 2002). The best matches (based on
284    BLAT score) for each contig were selected. Then, Est2genome (Mott, 1997) was used to refine the
285    alignments and we kept alignments with an identity percent and a coverage at least of 95% and 80%,
286    respectively. Finally, for each genomic cluster, the sequence with the best match against *Quercus*
287    *robur* or *Castanea mollissima* proteins was kept. This procedure yielded 34,987 unigenes (below
288    referred to as the 35K unigene set).

289        2.5. Genotyping-by-sequencing of the mapping population

290        2.5.1. RNAseq libraries construction

291    The 200 RNA samples were prepared as described above (Unigene set construction section), using
292    the TruSeq Stranded mRNA kit (Illumina, San Diego, CA, USA), from 500 ng total RNA. Libraries
293    were multiplexed onto Illumina Novaseq 6000 using S4 chemistry (2x150 read length), targeting
294    approximately 30 million reads per sample.

295        2.5.2. RNAseq reads processing for the MSSB accession

296    We first identified SNPs in the MSSB reference unigene. To this end, a trimming procedure was
297    applied to the MSSB sequences to remove adapters, primers, ribosomal reads and nucleotides with
298    quality value lower than 20 from both ends of the reads and reads shorter than 30 nucleotides as
299    described previously (Alberti et al., 2017). Trimmed reads were aligned onto the 35K unigene set
300    using bwa mem 0.7.17. Biallelic SNPs were identified using two methods: samtools 1.8 / bcftools 1.9
301    (Danecek et al. 2021) and GATK 3.8 (van der Auwera et al. 2020) with java 1.8.0_72. We kept SNPs
302    identified by both methods.

303        2.5.3. Identification of SNPs from RNAseq data and offspring genotype inference

304    We called SNPs and bioinformatically genotyped the mapping population at each MSSB
305    polymorphic site, based on the paired-end Illumina sequencing of 200 RNAseq libraries. The 200
306    raw-read datasets were trimmed following the same procedure used for MSSB. Reads were aligned to
307    the 35K unigene set using bwa mem 0.7.17. Genotypes were recovered from the 200 libraries at the
308    507,905 polymorphic positions, identified in MSSB, using GATK 3.8.

309    We then applied the following four-step filtering procedure: i/ for each SNP of a given half-sib,
310    polymorphic genotypes were set to monomorphic if the sequencing depth for this individual at this
311    position was lower than 20X; ii/ we kept SNPs only if at least 50% of the mapping population (i.e.
312    100 half-sibs) were heterozygous at this site; iii/ we kept only polymorphic sites consistent with a
313    1:1 heterozygote:homozygote genotype ratio, according to a Chi-square test with a 90% confidence
314    interval (Chi-square < 6.635, 1 d.f.), corresponding to heterozygous loci in the mother tree and
315    monomorph in all possible fathers; iv/ finally, for each contig, we retained only the SNP with fewest
316    missing data in the mapping population.

317        2.6. Linkage map construction

318 Half-sibs presenting too many missing data were discarded. As a result, 182 individuals (out of 200
319 selected from the first and second campaign) with valid genotypes for at least 4,127 loci were kept
320 for further analyses. A preliminary analysis was then performed using R-qtl package to group linked
321 SNP markers into robust linkage groups (LG) (LOD = 8) (Supplementary file 6). Given the large
322 number of markers per LG, marker ordering was performed within each LG using JoinMap 4.1
323 (Kyazma, Wageningen, NL). To this end, linkage groups of the maternal parent (MSSB) were
324 constructed using a four-step procedure: i) The maximum likelihood (ML) algorithm of JoinMap was
325 first used with a minimum linkage LOD score of 5 to calculate the number of crossing-overs (CO)
326 for each individual and to estimate the position of all mapped SNPs, ii/ then, the regression algorithm
327 (with a minimum LOD of 5 and default parameters: recombination frequency of 0.4 and maximum
328 threshold value of 1 for the jump) was used for a subset of evenly spaced SNPs (referred to below as
329 set #1 SNPs) along each LG, iii) the maternal linkage maps of the two full-sib families, identified
330 from the paternity test, were constructed using this subset of markers and individuals, providing two
331 genetic maps (referred to below as set #2 and set #3 SNPs) with higher confidence in genetic distance
332 estimates and marker ordering, both parents being known; iv) finally, from these two SNP datasets,
333 we created a final dataset (set #4) combining sets #2 and #3. For these 3 marker sets (#2, #3 and #4),
334 a first map was constructed using the ML algorithm to calculate the number of CO and a second map
335 was established using the regression algorithm excluding SNPs with high conflict of positions and
336 reducing the number of CO.

337        2.7. Genomic scaffold anchoring

338 Sequences of the unigenes encompassing SNP markers included in the linkage map, were aligned on
339 the genome assembly using BLAT with default parameters, except "-minScore=80". Unigenes
340 presenting more than one alignment were filtered out. In other words, when a second best match
341 having a score equal to or greater than 90% of the best score the marker was tagged as ambiguous.
342 For all the remaining alignments we kept only the alignment with the best score.

343

344 **3. Results**
345

346    3.1. General genome features

347    3.1.1. Genomic composition and completeness

348 The final assembly of the Bhaga genome was based on hybrid assembly of PacBio and Illumina reads
349 as well as scaffolding using a Hi-C library. It was resolved into 12 chromosomes, spanning 535.4 Mb
350 of the genome and 155 unassigned contigs of 4.9 Mb, which to 79% consisted of unplaced repeat
351 regions that precluded their unequivocal placement. It revealed a high level of BUSCO gene
352 detection (97.4%), surpassing that of the previous assembly and other genome assemblies available
353 for members of the *Fagaceae* (Table 1). Of the complete assembly, 57.12% were annotated as
354 interspersed repeat regions and 1.97% consisted of simple sequence repeats (see Supplementary File
355 1 for details regarding the repeat types and abundances).

356 The gene prediction pipeline yielded 63,736 complete genes with start and stop codons and a
357 minimum length of 180 bp. Out of these, 2,472 genes had alternate splice variants. For 86.8% of all

358 genes, a functional annotation could be assigned. Gene density varied widely in the genome, ranging
359 from zero per 100 kb window to 49.7%, with an average and median of 18.2% and 17.6%,
360 respectively. Gene lengths ranged from 180 to 54,183 bp, with an average and median gene length of
361 3,919 and 3,082 bp, respectively. In *Fagus sylvatica* 4.9 exons per gene were found on average,
362 corresponding well to other high-quality plant genome drafts. The distribution of exons and introns in
363 comparison to *J. regia* and *A. thaliana* are presented in Table 2. An analysis of PLAZA genes
364 identified 28,326 such genes in *F. sylvatica*, out of which 1,776 genes were present in three other
365 species used for comparison (Supplementary File 2).
366

367        3.1.2.  Telomere and centromere predictions

368 The results given above indicate a high quality of the genome assembly and the gene annotations. To
369 ascertain that the chromosomes were fully resolved, telomeric and centromeric regions were
370 predicted in the genome. The tandem repeat element TTTAGGG was the most abundant repeat in the
371 genome and was the building block of the telomeric repeats. Out of 12 chromosomes, 8 have
372 stretches of telomeric repeats towards both ends of the chromosomes and the other 4 chromosomes
373 have telomeric repeats towards only one end of chromosomes (Fig. 2). One unplaced scaffold of
374 110,653 bp which is composed of 12,051 bp of telomeric repeats at one end, probably represents one
375 of the missing chromosome-ends.
376 Two different types of potential centromeric repeats were observed, consisting of 79 bp and 80 bp
377 monomer units (Supplementary File 3). Centromeric repeats were also observed in higher numbers
378 outside the main centromeric region on several chromosomes (Supplementary File 3). However,
379 except for chromosome 10, there was a clear clustering of centromeric repeats within each of the
380 chromosomes, likely corresponding to the actual centromere of the respective chromosomes, and
381 supported also by complementary evidence, such as similarities to centromeric regions of *A. thaliana*,
382 high gypsy element content and low GC content (Supplementary File 3).
383

384        3.1.3. Integration of organelle DNA in the nuclear genome

385 As it has previously been shown that organelle DNA insertions can be uneven across the genome and
386 associated with chromatin structure (Wang and Timmis 2013), their distribution in the genome of
387 Bhaga was analysed. For both chloroplast (Mishra et al. 2021a) and mitochondria (Mishra et al.
388 2021b), multiple integrations of fragments of variable length of their genomic DNA were observed in
389 all chromosomes (Figs. 3, 4). These fragments varied in length from the minimum size threshold
390 (100 bp) to 54,784 bp for the chloroplast and 26,510 bp for the mitochondrial DNA. The identity of
391 the integrated organelle DNA with the corresponding stretches in the organelle genome ranged from
392 the minimum threshold tested of 95% to 100%. Nuclear-integrated fragments of organelle DNA
393 exceeding 10 kbp were found on six chromosomes for the chloroplast, but only on one chromosome
394 for the mitochondrial genome (Figs. 3, 4).
395 Nuclear insertions with sequence identity > 99% were about ten times more frequent for chloroplast
396 than for mitochondrial DNA with 173 vs. 16 for fragments > 1 kb and 115 vs. 11 for fragments > 5
397 kb, respectively. Eight of these matches of mitochondria were located on unplaced contigs. Overall,
398 mitochondrial insertions tended to be smaller and show a slightly higher sequence similarity
399 (Supplementary File 4), suggesting that they might be purged from the nuclear genome quicker than
400 the chloroplast genome insertions.
401 The integration of organelle DNA into the nuclear genome was mostly even, but tandem-like
402 integrations of chloroplast DNA on chromosome 2 were observed (Fig. 3). In addition, insertions of
403 both organelles were found close to the ends in 4 of the 24 chromosome ends (4, 6, 7, and 8). For the

404    insertions further than 500 kb away from the chromosome ends the integration sites of mitochondrion
405    DNA were sometimes found within the same 100 kb windows where the chloroplast DNA insertion
406    was found. If some regions of the genome are more amenable for the integration of organelle DNA
407    than others needs to be clarified in future studies. A major anomaly was found on chromosome 11,
408    where in a stretch of about 2 Mb (from about Mb 16-18 on that chromosome) consisting mainly of
409    multiple insertions of both chloroplast and mitochondrial DNA was observed. In this region, an
410    insertion of more than 20 kb of mitochondrial DNA was flanked by multiple very long integrations of
411    parts of the chloroplast genome on both sides (Figs. 3, 4). Thus, these integrations appeared almost
412    repeat-like at this particular location.

413

414          3.1.4. Repeat elements and gene space

415    The most abundant repeat elements were LTR elements and LINEs, covering 11.49% and 3.66% of
416    the genome, respectively. A detailed list of the element types found, their abundance and proportional
417    coverage of the genome is given in Supplementary File 1. Repeat elements presence was variable
418    across the chromosomes (Fig. 5). While the repeat content per 100 kb window exceeded 50 % over
419    more than 88% of chromosome 1, this was the case for only 37.5% of chromosome 9. Chromosomes
420    showed an accumulation of repeat elements towards their ends, except for chromosome 10, where
421    only a moderate increase was observed on one of the ends, and chromosome 1, where repeat
422    elements were more evenly distributed. Repeat content was unevenly distributed, with a patchy
423    distribution of repeat-rich and repeat-poor regions of variable length.
424    A conspicuous anomaly was noticed in chromosome 4, where at one end a large region of about 10
425    Mb was found in which 97% of the 100 kb windows had a repeat content greater than 70%. This
426    region also contained a high proportion of duplicated or multiplicated genes (Fig. 5). Additional
427    regions containing more than 20% of duplicated genes within a window of at least 1 Mb were
428    identified on chromosomes 4, 10, and 11. On chromosome 11, two clusters were detected, one of
429    which corresponded to the site of organelle DNA insertions described above.
430    The ribosomal cistrons were reported to be located at the telomeres of four different chromosomes in
431    *F. sylvatica* (Ribeiro et al., 2011). Due to the highly repetitive nature of the ribosomal repeats and
432    their placement near the telomeres, they could not be assigned with certainty to specific
433    chromosomes and thus remained in four unplaced contigs. However, the 5S unit, which is separate
434    from the other ribosomal units in F. sylvatica, could be placed near the centromeric locations of
435    chromosomes 1 and 2, in line with the locations inferred by fluorescence microscopy (Ribeiro et al.,
436    2011).
437    Coding space was more evenly distributed over the chromosomes, with the exception of the regions
438    with high levels of duplicated or multiplied genes. Apart from this, a randomly fluctuating proportion
439    of coding space was observed, with only few regions that seemed to be slightly enriched or depleted
440    in terms of coding space, e.g. in the central part of chromosome 8.

441

442          3.1.5. Distribution of single nucleotide polymorphisms

443    To study, if the distribution of single nucleotide polymorphisms (SNPs) correlates with the feature
444    reported above, they were identified on the basis of the comparison of the two individuals
445    investigated in this study, Bhaga and Jamy. A total of 2,787,807 SNPs were identified out of which
446    1,271,410 SNPs were homozygous (i.e. an alternating base on both chromosomes between Bhaga
447    and Jamy) and 1,582,804 were heterozygous (representing two alleles within Bhaga). A total of
448    269,756 SNPs fell inside coding regions out of which 119,946 were homozygous.

449 Heterozygous SNPs were very unequally distributed over the chromosomes (Fig. 6). Several regions,
450 the longest of which comprised more than 30 Mb on chromosome 6, contained only very low
451 amounts of heterozygous SNPs. Apart from the chromosome ends, where generally few heterozygous
452 positions were observed, all chromosomes contained at least one window of 1 Mb where only very
453 few heterozygous SNPs were present. On chromosomes 2, 3, 4, 6, and 9 such areas extended beyond
454 5 Mb. On chromosome 4 this region corresponded to the repeat region anomaly reported in the
455 previous paragraph, but for the region poor in heterozygous SNPs on chromosome 9, no association
456 with a repeat-rich region could be observed.
457 Homozygous SNPs differentiating Bhaga and Jamy, often followed a different pattern. All regions
458 with low heterozygous SNP frequency longer than 5 Mb had an above-average homozygous SNP
459 frequency, with the exception of the anomalous repeat-rich region on chromosome 4, which had very
460 low frequencies for both homozygous and heterozygous SNPs. However, there were also two regions
461 of more than 1 Mb length on chromosome 11 that also showed low frequencies of both SNP
462 categories (Fig. 6).
463 Generally, the frequency of overall and intergenic SNPs per 100 kb window corresponded well for
464 both heterozygous and homozygous SNPs, suggesting neutral evolution. However, there were some
465 regions in which genic and intergenic SNP frequencies were uncoupled. For example, on
466 chromosome 1 a high overall heterozygous SNP frequency was observed at 37.7, 48.2 and 56 Mb,
467 but genic heterozygous SNP frequency was low despite normal gene density, suggesting the presence
468 of highly conserved genes. In line with this, also the frequency of homozygous genic SNPs was
469 equally low in the corresponding areas. Similary, homozygous SNP frequencies were also decoupled
470 on chromosome 1, where a low frequency was observed at 4.2, 7.1, 38.2, 62.1, and 64.8 Mb, but a
471 high genic SNP frequency was observed. This suggests the presence of diversifying genes in the
472 corresponding 100 kb windows, such as genes involved in coping with biotic or abiotic stress.

473 In line with the different distribution over the chromosomes, with large areas poor in heterozygous
474 SNPs, there were much more windows with low numbers of heterozygous SNPs than windows with
475 homozygous SNPs (Fig. 7). Notably, at intermediate SNP frequencies, homozygous SNPs were
476 found in more 100 kb windows, while at very high SNP frequencies, heterozygous SNPs were more
477 commonly found. This pattern is consistent with predominant local pollination, but occasional
478 introgression of highly distinct genotypes.

479       3.1.6. Genome browser

480 A genome browser for the genome of Bhaga, with the various genomic features outlined above
481 annotated, is available at beechgenome.net. Predicted genes, annotated repeat elements and
482 homozygous and heterozygous SNPs are available in "B. Annotations". The telomeric and
483 centromeric locations, as well as the GC content details are available in "C. Other Details".

484       3.2. Validation of chromosomal-scale pseudomolecules

485       3.2.1. Pedigree reconstruction

486 The analysis of the 36 SNPs using Cervus allowed the identification of candidate fathers and
487 reconstruct full-sib families. For 317 of the 537 offspring a likely father was identified. The 19
488 candidate fathers were represented in the progeny, although their contributions were variable (0.8 to
489 21%). For the other offpring, no father could be assigned, i.e. the pollen donor is not present among
490 the surrounding trees (corresponding to 210 genotypes, i.e. 39.1% of the samples when 0 mismatch is
491 allowed, and 22 % when 1 mismatch is allowed). The two largest families comprised 68

492 (MSSBxMSSH) and 86 (MSSBxSSP12) full-sibs. Few years after plantations, 36 genotypes for the
493 former and 49 for the latter survived (Table 3).

### 3.2.2. A new unigene set for European beech

495 Our study provides a new reference unigene set for *Fagus sylvatica* based on short and long NGS
496 reads obtained from cDNA libraries constructed from six different tissues. The first unigene set for
497 this species was established back in 2015 using a combination of Sanger and Roche-454 reads (Lesur
498 et al. 2015). The sequences were assembled into 21000 contigs. A second step was achieved by
499 Müller et al. (2017) using NGS data (Illumina) resulting in 44000 contigs. Tis third transcript catalog
500 contains a total of 34,987 items. When compared to the oak proteome (to date the best annotated
501 among Fagaceae species), this new reference provides the most complete transcript catalog (Table 4).

### 3.2.3. Identification of RNAseq-based SNP markers for linkage mapping

503 Sequencing of the six tissues (collected on the MSSB accession) using an RNA-Seq approach, led to
504 408,111,505 Illumina paired-end reads. A total of 383,149,091 trimmed sequences were used to
505 identify putative segregating SNPs in MSSB.

506 On average, 82.67% of the reads were properly aligned on the reference unigene, ranging from
507 72.94% for the male flowers to 86.46% for leaves. We identified 613,885 and 507,905 SNPs using
508 Samtools/bcftools and GATK, respectively. A total of 507,905 SNPs in MSSB were finally identified
509 by both methods.

510 Sequencing of the 200 siblings, followed by trimming of the raw data, led to a total of 9,155,925,565
511 reads. On average, 78.64% of the reads were properly aligned on the reference unigene (min. 72.6% -
512 max. 83.04%). We found 267,361 polymorphic sites in at least one out of the 200 half-sibs. Our four-
513 step filtering process yielded a final set of 6,385 SNPs spread over 6,385 contigs, with at least 20X
514 coverage.

### 3.2.4. Linkage map construction

516 Beech is a diploid species with 2n=2x=24. The 12 expected linkage groups (LG) were retrieved using
517 SNPs from set #1 using the R-qtl package. The number of SNP markers per LG ranged from 231 to
518 412. However, the detailed linkage analysis, carried out with JoinMap for each LG, revealed an
519 unexpectedly high number of crossing-overs and oversized LGs compared to previous linkage
520 mapping analyses performed in beech (Scalfi et al., 2004) or oak (Bodénès et al., 2016), probably
521 owing to genotyping errors among the 182 hal-sibs. Because of this, we established genetic linkage
522 maps based on the two largest full-sib families identified from the paternity analysis, and only used
523 the corresponding two sets of mapped SNPs (sets #2 and #3) to create a combined genetic linkage
524 map based on the analysis of 182 half-sibs. A total of 768 SNPs were available for the combined
525 maternal linkage map, 368 of which were unambiguously mapped on the 13 longest LGs. The size of
526 LGs varied from 64 to 279 cM and comprised 8 to 56 SNPs (Table 5). High colinearity was observed
527 between the homologous linkage groups obtained from the three different maps (Fig. 8).

### 3.2.5. Alignment of Bhaga genomic scaffolds to the SNP-based linkage map of beech

529 The 368 mapped markers were aligned on the 12 genomic scaffolds (Bagha_1 to Bagha_12) of the
530 Fagus sylvatica genome assembly. The alignments were filtered and congruence between scaffolds
531 and linkage groups were checked. Most of the markers from a given LG mapped on a single scaffold

13

532 (Table 6) providing a genetic validation of the physical assembly obtained for the Bagha genome
533 sequence. Notable exceptions were: (i) LG11 and LG12, which corresponded to Bagha_#8; these two
534 chromosomal arms could not be merged into a single LG, and (ii) LG13 and scaffold #11, which
535 presented too few markers for unambiguous assignment to one or more scaffolds and LGs,
536 respectively.
537

## 4. Discussion

539 4.1. General genome features

540 The genome assembled and analysed in this study compares well with previously published
541 *Fagaceae* genomes, both in terms of size and gene space. We here confirm the base chromosome
542 number of 12, as was previously reported based on chromosome counts (Ribeiro et al., 2011). The
543 number of exons per gene is moderately higher than in the previously published genome of the same
544 individual (Mishra et al., 2018), reflecting the higher contiguity of the presented chromosome-level
545 assembly. Despite the lower chromosome number of the beech genome, it is structurally similar to
546 the available genomes of genus *Juglans*, which is the most closely related genus for which
547 chromosome-level assemblies are available, with continuous sequences from telomere to telomere (*J. 
548 regia* (Marrano et al., 2020); *J. sigillata* (Ning et al., 2020); *J. regia × J. microcarpa* (Zhu et al.,
549 2019)).

550 4.2. Telomere and centromere predictions

551 Telomeres are inherently difficult to resolve because of long stretches of GC-rich repeats that can
552 cause artefacts during library preparation (Aird et al., 2011) and can lead to biased mapping (Dohm
553 et al., 2008). However, using long-read sequencing and Hi-C scaffolding, we could identify telomeric
554 repeats on all chromosomes. It seems likely that several of the unplaced contigs of 4.9 Mb, which
555 included telomeric sequences, were not correctly anchored in the assembly due to ambiguous Hi-C
556 association data resulting from the high sequence similarity of telomeric repeats, because of which
557 for four chromosomes we could identify telomeric repeats only on one of the ends. This might also
558 be due to the presence of ribosomal cistrons on four chromosome ends, which might have interfered
559 with the Hi-C linkage due to their length and very high sequence similarity. On the outermost regions
560 of the chromosomes, no longer telomeric repeat stretches were present most likely due to their
561 ambiguous placement in the assembly, because of very high sequence similarity.

562 Centromere repeats were identified by screening the genome for repeats of intermediate sizes, and
563 were found to be present predominantly within a single location per chromosome. However, lower
564 amounts of centromeric repeat units were also observed to be scattered throughout the genome. The
565 function of the centromeric repeats outside of the centromere remains largely enigmatic but could be
566 associated with chromosome structuring (Alves et al., 2012) or centromere repositioning
567 (Mandáková et al., 2020; Klein and O'Neill, 2018). Interestingly, we could find two major groups of
568 potential centromeric repeat units of different lengths, which did not always coincide. The location of
569 the main occurrence of the centromere-defining repeat unit agreed well with the location previously
570 inferred using chromosome preparations and fluorescence microscopy (Ribeiro et al., 2011).

571 4.3. Integration of organelle DNA in the nuclear genome

572 Organelle DNA integration has been frequently found in all kingdoms of life for which high-
573 resolution genomes are available (Zhang et al., 2020; Guo et al., 2008; Stegemann et al., 2003). It can

14

574     be assumed that this transfer of organelle DNA to the nucleus is the seed of transfer of chloroplast
575     genes to the nuclear genome (Huang et al., 2003). However, apart from a few hints (Yang et al.,
576     2017) it is unclear, which factors stabilise the chloroplast genome so that its content in non-parasitic
577     plants stays relatively stable over long evolutionary timescales (Xiong et al., 2009; Wang et al.,
578     2007). In the present study, it has been found that the insertion of organelle DNA insertions are
579     located mainly in repeat-rich regions of the beech genome. However, their presence in regions
580     without pronounced repeat density might suggest that repeats are not the only factor associated with
581     the insertion of organelle DNA. Nevertheless, it appears that some regions are generally amenable to
582     the integration of organelle DNA, as in several cases chloroplast and mitochondrion insertions were
583     observed in close proximity. The reason for this is unclear, but is known that open chromatin is more
584     likely to accumulate insertions (Wang and Timmis 2013). The potential presence of areas in the
585     genome that are less protected from the insertion of foreign DNA could open up potential molecular
586     biology applications for creating stable transformants.

587     An anomaly regarding organelle DNA insertion was observed on chromosome 11. Around a central
588     insertion of mitochondrion DNA, multiple insertions of chloroplast DNA were found. The whole
589     region spans more than 2 Mb, which is significantly longer than the organelle integration hotspots
590     reported in other species (Zhang et al., 2020). The evolutionary origin of this large chromosome
591     region is unclear, but given its repetitive nature it is conceivable that it resulted from a combination
592     of an integration of long fragments and repeat element activity. The presence of multiple copies at the
593     location implies an unusual genome structure in this area, but further analyses, ideally including
594     multiple additional individuals, will be necessary to elucidate the basis for this.

595         4.4. Distribution of single nucleotide polymorphisms (SNPs)

596     SNP content was found to vary across all chromosomes leading to a mosaic pattern. While most of
597     the areas of high or low SNP density were rather short and not correlated to any other patterns, there
598     were several regions > 1 Mbp that exhibited a similar polymorphism type, suggesting non-neutral
599     evolution.

600     The longest of those stretches poor in both heterozygous and homozygous positions was found on
601     chromosome 4, and corresponded to a region rich in both genes and repeat elements. This is
602     remarkable and probably due to a recent proliferation, as repeat-rich regions are usually less stable
603     and more prone to accumulate mutations (Wang et al., 2020; Flynn et al., 2018; Ho et al., 2020).

604     Most regions with lower abundance of heterozygous SNPs than on average were found to be
605     particularly high in homozygous SNPs. The longest of such stretches was found on chromosome 6,
606     comprising about two thirds of the entire chromosome. Three more such regions longer than 5 Mbp
607     were found on other chromosomes. The evolutionary significance of this is unclear, but it is
608     conceivable that these areas contain locale specific variants for which no alternative alleles are shared
609     within the same stand. For confirmation of this hypothesis, it would be important to evaluate genetic
610     markers from additional individuals of the same stand. Locally adaptive alleles could be fixed
611     relatively easy by local inbreeding (Ceballos et al., 2018), considering the low seed dispersal kernel
612     of European Beech (Martínez and González-Taboada, 2009). The presence of genes involved in local
613     adaptation could explain the rather high amount of homozygous SNPs in the same location, as the
614     stands from which the two studied individuals came from differ in soil, water availability,
615     continentality, and light availability. However, more individuals from geographically separated
616     similar stands need to be investigated to disentangle the effects of inbreeding and local adaptation.

617  In summary, homozygous and heterozygous SNPs were rather uniformly distributed throughout the
618  major part of the genome, suggesting neutral evolution or balancing selection.

619

620  ### 5. Conclusions

621  The chromosome-level assembly of the ultra-centennial individual Bhaga from the Kellerwald-
622  Edersee National Park in Germany and its comparison with the individual Jamy from the Jamy
623  Nature Reserve in Poland has revealed several notable genomic features. The prediction of the
624  telomeres and centromeres as well as ribosomal DNA corresponded well with data gained from
625  chromosome imaging (Ribeiro et al., 2011), suggesting state-of-the-art accuracy of the assembly.
626  Interestingly, several anomalies were observed in the genome, corresponding to regions with
627  abundant integrations of organelle DNA, low frequency of both heterozygous and homozygous
628  SNPs, and long chromosome stretches almost homozygous but with a high frequency of SNPs
629  differentiating the individuals.

630  Taken together, the data presented here suggest a strongly partitioned genome architecture and
631  potentially divergent selection regimes in the stands of the two individuals investigated here. Future
632  comparisons of additional genomes to the reference will help understanding the significance of
633  variant sites identified in this study and shed light on the fundamental processes involved in local
634  adaptation of a long-lived tree species exposed to a changing climate.

635  ### 6. Data availability

636  The data sets supporting the results of this article are available in the GenBank repository, under the
637  accession number PRJEB24056 for the *Fagus sylvatica* reference individual Bhaga, PRJNA450822
638  for the individual Jamy, PRJEB46583 for sequencing of a new unigene set and PRJEB46593 for
639  RNA-seq-based genetic mapping in European beech.

640

641  ### 7. Conflict of Interest

642  The authors declare that they have no competing interest.

643

644  ### 8. Author Contributions

645  M.T. conceived the study, wih contributions from C.P. and J.B. B.U., C.B., J.B., J.M., J.M.A, L.O.,
646  M.T., and S.P. provided materials. All authors conducted laboratory experiments or analysed the
647  data. All authors were involved in data interpretation. B.M. and M.T. wrote the manuscript with
648  contributions from the other authors. All authors read and approved the final manuscript.

649

650  ### 9. Funding and Acknowledgment

16

**10. References**

Aird, D., Ross, M.G., Chen, W.S., Danielsson, M., Fennell, T., Russ, C., et al. (2011). Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. Genome Biol. 12(R18), 1–14.

Alberti, A., Poulain, J., Engelen, S., Labadie, K., Romac, S., Ferrera, T., et al. (2017). Viral to metazoan marine plankton nucleotide sequences from the Tara Oceans expedition. Sci Data 4, 17009.

Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J. (1990). Basic local alignment search tool. J Mol Biol. 215, 403–410.

Atwell, S., Huang, Y.S., Vilhjálmsson, B.J., Willems, G., Horton, M., Li, Y., et al. (2010). Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. Nature. 465, 627–631.

Alves, S., Ribeiro, T., Inácio, V., Rocheta, M., Morais-Cecílio, L. (2012). Genomic organization and dynamics of repetitive DNA sequences in representatives of three *Fagaceae* genera. Genome. 55, 348–359.

Ayad, L.A. and Pissis, S.P. (2017). MARS: improving multiple circular sequence alignment using refined sequences. BMC Genomics. 18(86), 1–10.

Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids Res. 27, 573–580.

17

701  Bodénès, C., Chancerel, E., Ehrenmann, F., Kremer, A., Plomion, C. (2016). High-density linkage
702  mapping and distribution of segregation distortion regions in the oak genome. DNA Research. 23,
703  115-124.
704
705  Bolte, A., Degen, B. (2010). Forest adaptation to climate change - options and limitations.
706  Landbauforsch Volk. 60, 111–117.
707
708  Buchfink, B. and Xie, C. (2015). Huson DH. Fast and sensitive protein alignment using DIAMOND.
709  Nat Meth. 12, 59–60.
710
711  Buels, R., Yao, E., Diesh, C.M., Hayes, R.D., Munoz-Torres, M., Helt, G., et al. (2016). JBrowse: a
712  dynamic web platform for genome visualization and analysis. Genome Biol. 17(66), 1–12.
713
714  Caudullo, G., Durrant, T.H., Mauri, A. (Luxembourg: Publication Office of the European Union),
715  94–95.
716
717  Chen, J.D., Zheng, C., Ma, J.Q., Jiang, C.K., Ercisli, S., Yao, M.Z., et al. (2020). The chromosome-
718  scale genome reveals the evolution and diversification after the recent tetraploidization event in tea
719  plant. Hortic Res. 7(63), 1–11.
720
721  Ceballos, F.C., Joshi, P.K., Clark, D.W., Ramsay, M., Wilson, J.F. (2018). Runs of homozygosity:
722  windows into population history and trait architecture. Nat Rev Gen.19, 220.
723
724  Chen, S., Huang, T., Zhou, Y., Han, Y., Xu, M., Gu, J. (2017). AfterQC: automatic filtering,
725  trimming, error removing and quality control for fastq data. BMC Bioinf. 18, 80.
726  doi:10.1186/s12859-017-1469-3
727
728  Coumou, D., Rahmstorf, S. (2012). A decade of weather extremes. Nat Clim Change. 2, 491–496.
729
730  Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., et. al. (2021). Twelve
731  years of SAMtools and BCFtools. GigaScience. 10(2)Di Filippo, A., Biondi, F., Maugeri, M.,
732  Schirone, B., Piovesan, G. (2012). Bioclimate and growth history affect beech lifespan in the Italian
733  Alps and Apennines in Glob Change Biol. 18, 960–972.
734  Dohm, J.C., Lottaz, C., Borodina, T., Himmelbauer, H. (2008). Substantial biases in ultra-short read
735  data sets from high-throughput DNA sequencing. Nucleic Acids Res. 36, e105.
736
737  Durrant, T.H., De Rigo, D., Caudullo, G. (2016). "*Fagus sylvatica* in Europe: distribution, habitat,
738  usage and threats." in European atlas of forest tree species, ed. San-Miguel-Ayanz J, de Rigo D,
739
740  Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., Smit, A.F. (2020).
741  RepeatModeler2 for automated genomic discovery of transposable element families. PNAS. 117,
742  9451–9457.
743
744  Flynn, J.M., Lower, S.E., Barbash, D.A., Clark, A.G. (2018). Rates and patterns of mutation in
745  tandem repetitive DNA in six independent lineages of *Chlamydomonas reinhardtii*. Genome Biol
746  Evol. 10, 1673–1686.
747
748  Fu, L., Niu, B., Zhu, Z., Wu, S., Li, W. (2012), CD-HIT: accelerated for clustering the next-
749  generation sequencing data. Bioinformatics. 28, 150–152.

750

751 Geßler, A., Keitel, C., Kreuzwieser, J., Matyssek, R., Seiler, W., Rennenberg, H. (2007). Potential
752 risks for European beech (*Fagus sylvatica* L.) in a changing climate. Trees. 21, 1–11.

753

754 Guo, X., Ruan, S., Hu, W., Cai, D., Fan, L. (2008). Chloroplast DNA insertions into the nuclear
755 genome of rice: the genes, sites and ages of insertion involved. Funct Integr Genomic. 8, 101–108.

756

757 Götz, S., García-Gómez, J.M., Terol, J., Williams, T.D., Nagaraj, S.H., Nueda, M.J., et al. (2008).
758 High-throughput functional annotation and data mining with the Blast2GO suite. Nucleic Acids Res.
759 36, 3420–3435.

760

761 Haas, B.J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P.D., Bowden, J., et al. (2013). De
762 novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference
763 generation and analysis. Nat. Protoc. 8, 1494–1512.

764

765 Ho, E.K., Bellis, E.S., Calkins, J., Adrion, J.R., Latta, L.C., Schaack, S. (2020). Engines of change:
766 Transposable element mutation rates are high and vary widely among genotypes and populations of
767 *Daphnia magna*. bioRxiv. doi: 10.1101/2020.09.21.307181

768

769 Hong, Z., Li, J., Liu, X., Lian, J., Zhang, N., Yang, Z., et al. (2020). The chromosome-level draft
770 genome of *Dalbergia odorifera*. GigaScience. 9, giaa084.

771

772 Huang, C.Y., Ayliffe, M.A., Timmis, J.N. (2003). Direct measurement of the transfer rate of
773 chloroplast DNA into the nucleus. Nature. 422, 72–76.

774

775 Jiang, S., An, H., Xu, F., Zhang, X. (2020). Chromosome-level genome assembly and annotation of
776 the loquat (*Eriobotrya japonica*) genome. GigaScience. 9, giaa015.

777

778 Jones, P., Binns, D., Chang, H., Fraser, M., Li, W., Mc Anulla, C., et al. (2014). InterProScan 5:
779 genome-scale protein function classification. Bioinf. 30, 1236–1240.

780

781 Jump, A.S., Hunt, J.M., Penuelas, J. (2006). Rapid climate change-related growth decline at the
782 southern range edge of *Fagus sylvatica*. Glob Change Biol. 12, 2163–2174.

783

784 Kalinowski, S.T., Taper, M.L., Marshall, T.C. (2007). Revising how the computer program CERVUS
785 accommodates genotyping error increases success in paternity assignment. Molecular Ecology. 16,
786 1099–1106.

787

788 Kent, W. J. (2002). BLAT—The BLAST-Like Alignment Tool. Genome Res. 12, 656–664.

789

790 Kim, D., Langmead, B., Salzberg, S.L. (2015). HISAT: a fast spliced aligner with low memory
791 requirements. Nat Methods. 12, 357–360.

792

793 Klein, S.J. and O'Neill, R.J. (2018). Transposable elements: genome innovation, chromosome
794 diversity, and centromere conflict. Chromosome Res. 26, 5–23.

795

796 Kremer, A., Abbott, A.G., Carlson, J.E., Manos, P.S., Plomion, C., Sisco, P., et al. (2012). Genomics
797 of *Fagaceae*. Tree Genet Genomes. 8, 583–610.

798

799  Lalagüe, H., Csilléry, K., Oddou-Muratorio, S., Safrana, J., de Quattro, C., Fady, B., et al. (2014).
800  Nucleotide diversity and linkage disequilibrium at 58 stress response and phenology candidate genes
801  in a European beech (Fagus sylvatica L.) population from southeastern France. Tree Genetics and
802  Genomes. 10, 15-26.
803
804  Langmead, B. and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. Nat Methods. 9,
805  357–359.
806
807  Le Provost, G., Herrera, R., Paiva, J.A.P., Chaumeil P, Salin F, Plomion C (2007). A micromethod
808  for high throughput RNA extraction in forest trees. Biological Research. 40, 291–297.
809
810  Lesur, I., Bechade, A., Lalanne, C., Klopp, C., Noirot, C., Leplé, J.C., et al. (2015). A unigene set for
811  European beech (Fagus sylvatica L.) and its use to decipher the molecular mechanisms involved in
812  dormancy regulation. Mol Ecol Res. 15, 1192-1204.
813
814  Leuschner, C., Meier, I.C., Hertel, D. (2006) On the niche breadth of *Fagus sylvatica*: soil nutrient
815  status in 50 Central European beech stands on a broad range of bedrock types. Ann For Sci. 63, 355–
816  368.
817
818  Lieberman-Aiden, E., Van Berkum, N.L., Williams, L., Imakaev, M., Ragoczy, T., Telling, A., et al.
819  (2009). Comprehensive mapping of long-range interactions reveals folding principles of the human
820  genome. Science. 326, 289–293.
821
822  Ligot, G., Balandier, P., Fayolle, A., Lejeune, P. (2013). Claessens H. Height competition between
823  *Quercus petraea* and *Fagus sylvatica* natural regeneration in mixed and uneven-aged stands. Forest
824  Ecol Manag. 304, 391–398.
825
826  Li H. (2011). A statistical framework for SNP calling, mutation discovery, association mapping and
827  population genetical parameter estimation from sequencing data. Bioinf. 27, 2987–2993.
828
829  Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler
830  transform. Bioinf. 25, 1754–1760.
831
832  Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009). The sequence
833  alignment/map format and SAMtools. Bioinf. 25, 2078–2079.
834
835  Li, W. and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of
836  protein or nucleotide sequences. Bioinf. 22, 1658–1659.
837
838  Madoui, M. A., Engelen, S., Cruaud, C., Belser, C., Bertrand, L., Alberti, A., et al. (2015). Genome
839  assembly using Nanopore-guided long and error-free DNA reads. BMC Genomics. 16, 327.
840
841  Mandáková,T., Hloušková, P., Koch, M.A., Lysak, M.A. (2020). Genome evolution in Arabideae
842  was marked by frequent centromere repositioning. Plant Cell. 32, 650–665.
843
844  Marchler-Bauer, A., Lu, S., Anderson, J.B., Chitsaz, F., Derbyshire, M.K., DeWeese-Scott, C., et al.
845  (2011). CDD: a Conserved Domain Database for the functional annotation of proteins. Nucleic Acids
846  Res. 39, D225–D229.
847

This is a provisional file, not the final typeset article

848  Marrano, A., Britton, M., Zaini, P.A., Zimin, A.V., Workman, R.E., Puiu, D., et al. (2020). High-
849  quality chromosome-scale assembly of the walnut (Juglans regia L.) reference genome. GigaScience.
850  9, giaa050.
851
852  Marshall, T.C., Slate, J., Kruuk, L.E.B., Pemberton, J.M. (1998). Statistical confidence for
853  likelihood-based paternity inference in natural populations. Molecular Ecology. 7, 639–655.
854
855  Martínez-García, P.J., Crepeau, M.W., Puiu, D., Gonzalez-Ibeas, D., Whalen, J., Stevens, K.A., et al.
856  (2016). The walnut (Juglans regia) genome sequence reveals diversity in genes coding for the
857  biosynthesis of non-structural polyphenols. Plant J. 87, 507–32.
858
859  Martínez, I. and González-Taboada, F. (2009). Seed dispersal patterns in a temperate forest during a
860  mast event: performance of alternative dispersal kernels. Oecologia. 159, 389–400.
861
862  Michael, T.P., VanBuren, R. (2020). Building near-complete plant genomes. Curr Opin Plant Biol.
863  54, 26–33.
864
865
866
867
868  Mishra, B., Gupta, D.K., Pfenninger, M., Hickler, T., Langer, E., Nam, B., et al. (2018). A reference
869  genome of the European beech (*Fagus sylvatica L.*). GigaScience. 7, giy063.
870
871  Mishra, B., Ulaszewski, B., Ploch, S., Burczyk, J., Thines, M. (2021a). A circular chloroplast
872  genome of *Fagus sylvatica* reveals high conservation between two individuals from Germany and
873  one individual from Poland and an alternate direction of the small single-copy region. Forests. 12,
874  180.
875
876  Mishra, B., Ulaszewski, B., Meger, J., Ploch, S., Burczyk, J., Thines, M. (2021b). A comparison of
877  three circular mitochondrial genomes of *Fagus sylvatica* from Germany and Poland reveals low
878  variation and complete identity of the gene space. Forests. 12, 571.
879
880  Mott, R. (1997). EST_GENOME: a program to align spliced DNA sequences to unspliced genomic
881  DNA. Comput. Appl. Biosci. CABIOS 13, 477–478.
882
883  Müller, M., Seifert, S., Lübbe, T., Leuschner, C., and Finkeldey, R. (2017). De novo transcriptome
884  assembly and analysis of differential gene expression in response to drought in European beech. PloS
885  one. 12(9), e0184167. https://doi.org/10.1371/journal.pone.0184167
886
887  NCBI nr database. https://ftp.ncbi.nlm.nih.gov/blast/db/ [accessed June 24, 2020].
888
889  Ning, D.L., Wu, T., Xiao, L.J., Ma, T., Fang. W.L., Dong, R.Q., Cao, F.L. (2020). Chromosomal-
890  level assembly of *Juglans sigillata* genome using Nanopore, BioNano, and Hi-C analysis.
891  GigaScience. 9, giaa006.
892
893  Nong, W., Law, S.T., Wong, A.Y., Baril, T., Swale, T., Chu, L.M., et al. (2020). Chromosomal-level
894  reference genome of the incense tree *Aquilaria sinensis*. Mol Ecol Resour. 20, 971.
895

896  OmicsBox - Bioinformatics Made Easy, BioBam Bioinformatics. (2020).
897  https://www.biobam.com/omicsbox [Accessed March 3, 2020]
898
899  Ouayjan, A. and Hampe, A. (2018). Extensive sib-mating in a refugial population of beech (Fagus
900  sylvatica) growing along a lowland river. Forest Ecology and Management. 407, 66–74.
901
902  Pfenninger, M., Reuss, F., Kiebler, A., Schönnenbeck, P., Caliendo, C., Gerber, S., et al. (2020).
903  Genomic basis of drought resistance in *Fagus sylvatica*. bioRxiv. doi: 10.1101/2020.12.04.411264
904
905  PLAZA 4.0: an integrative resource for functional, evolutionary and comparative plant genomics
906  Nucleic Acids Res (online access). Accessed October 21, 2020.
907
908  Plomion, C., Aury, J.M., Amselem, J., Alaeitabar, T., Barbe, V., Belser, C., et al. (2016). Decoding
909  the oak genome: public release of sequence data, assembly, annotation and publication strategies.
910  Mol Ecol Resour. 16, 254–65.
911
912  Plomion C, Aury JM, Amselem J, Leroy T, Murat F, Duplessis S., et al. (2018). Oak genome reveals
913  facets of long lifespan. Nature Plants. 4, 440-452.
914
915  Price, A.L., Jones, N.C., Pevzner, P.A. (2005). De novo identification of repeat families in large
916  genomes. Bioinf. 21(suppl_1), i351–358.
917
918  Priest, S.J., Yadav, V., Heitman, J. (2020). Advances in understanding the evolution of fungal
919  genome architecture. F1000Research. 9(Faculty Rev), 776, doi:10.12688/f1000research.25424.1
920
921  Reif, A., Xystrakis, F., Gaertner, S., Sayer, U. (2017). Floristic change at the drought limit of
922  European beech (*Fagus sylvatica L.*) to downy oak (*Quercus pubescens*) forest in the temperate
923  climate of central Europe. Not Bot Horti Agrobo. 45, 646–54.
924
925  Rhie, A., Mc Carthy, S.A., Fedrigo, O., Damas, J., Formenti, G., Koren, S., et al. (2020). Towards
926  complete and error-free genome assemblies of all vertebrate species. BioRxiv. doi:
927  10.1101/2020.05.22.110833.
928
929  Ribeiro, T., Loureiro, J., Santos, C., Morais-Cecílio, L. (2011). Evolution of rDNA FISH patterns in
930  the Fagaceae. Tree Genet & Genomes. 7, 1113-1122.
931
932  Rice, P., Longden, I., Bleasby, A. (2000). EMBOSS: the European molecular biology open software
933  suite. Trends Genet. 16, 276–277.
934
935  Rose, L., Leuschner, C., Köckemann, B., Buschmann, H. (2009). Are marginal beech (*Fagus
936  sylvatica L.*) provenances a source for drought tolerant ecotypes? Eur J For Res. 128, 335–343.
937
938  Scalfi, M., Troggio, M., Piovani, P., Leonardi, S., Magnaschi, G., Vendramin, G.G., et al. (2004). A
939  RAPD, AFLP and SSR linkage map, and QTL analysis in european beech (Fagus sylvatica L.).
940  Theoretical and Applied Genetics. 108(3), 433-441.
941
942  Schulz, M.H., Zerbino, D.R., Vingron, M. and Birney, E. (2012). Oases: robust de novo RNA-seq
943  assembly across the dynamic range of expression levels. Bioinformatics. 28, 1086–1092.
944

945 Seppey, M., Manni, M., Zdobnov, E.M. (2019). "BUSCO: assessing genome assembly and
946 annotation completeness." in Gene Prediction. Methods in Molecular Biology, vol 1962, ed. Kollmar,
947 M. (New York: Humana), 227–245.

949 Smit, A.F.A., Hubley, R. (2007). RepeatMasker Open-4.0.5. 2007–2014.
950 http://www.repeatmasker.org. [Accessed Nov 16, 2020].

952 Spinoni, J., Naumann, G., Vogt, J., Barbosa, P. (2015). European drought climatologies and trends
953 based on a multi-indicator approach. Global Planet Change. 127, 50–57.

955 Sork, V.L., Squire, K., Gugger, P.F., Steele, S.E., Levy, E.D., Eckert, A.J. (2016). Landscape
956 genomic analysis of candidate genes for climate adaptation in a California endemic oak, *Quercus*
957 *lobata*. American J Bot. 103, 33–46.

959 Stanke, M. and Morgenstern, B. (2005). AUGUSTUS: a web server for gene prediction in eukaryotes
960 that allows user-defined constraints. Nucleic Acids Res. 33(suppl_2), 465–467.

962 Stegemann, S., Hartmann, S., Ruf, S., Bock, R. (2003). High-frequency gene transfer from the
963 chloroplast genome to the nucleus. PNAS. 100, 8828–8833.

965 Strijk, J.S., Hinsinger, D.D., Zhang, F., Cao, K. (2019), Trochodendron aralioides, the first
966 chromosome-level draft genome in *Trochodendrales* and a valuable resource for basal eudicot
967 research. GigaScience. 8, giz136.

969 Tarailo-Graovac, M., Chen, N. (2009). Using RepeatMasker to identify repetitive elements in
970 genomic sequences. Curr Prot Bioinf. 25, 4.10.1–4.10.14.

972 Van der Auwera, G.A. and O'Connor, B.D. (2020). Genomics in the Cloud: Using Docker, GATK,
973 and WDL in Terra (1st Edition). O'Reilly Media.

975 Wagner, S., Collet, C., Madsen, P., Nakashizuka, T., Nyland, R.D., Sagheb-Talebi, K. (2010). Beech
976 regeneration research: from ecological to silvicultural aspects. Forest Ecol Manag. 259, 2172–2182.

978 Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., et al. (2014). Pilon: an
979 integrated tool for comprehensive microbial variant detection and genome assembly improvement.
980 PloS ONE. 9, e112963.

982 Wang, D. and Timmis, J.N. (2013). Cytoplasmic organelle DNA preferentially inserts into open
983 chromatin. Genome Biol Evol. 5, 1060–1064.

985 Wang, D., Wu, Y.W., Shih, A.C.C., Wu, C.S., Wang, Y.N., Chaw, S.M. (2007). Transfer of
986 chloroplast genomic DNA to mitochondrial genome occurred at least 300 MYA. Mol Biol Evol. 24,
987 2040–2048.

989 Wang, J., Tian, S., Sun, X., Cheng, X., Duan, N., Tao, J., Shen, G. (2020). Construction of
990 Pseudomolecules for the Chinese Chestnut (*Castanea mollissima*) Genome. G3. 10, 3565–3574.

992 Wang, L., Sun, Y., Sun, X., Yu, L., Xue, L., He, Z., et al. (2020). Repeat-induced point mutation in
993 *Neurospora crassa* causes the highest known mutation rate and mutational burden of any cellular
994 life. Genome Biol. 21(142), 1–23.

996 Xiong, A.S., Peng, R.H., Zhuang, J., Gao, F., Zhu, B., Fu, X.Y., et al. (2009). Gene duplication,
997 transfer, and evolution in the chloroplast genome. Biotechnol Adv. 27, 340–347.

999 Yang, F.S., Nie, S., Liu, H., Shi, T.L., Tian, X.C., Zhou, S.S., et al. (2020). Chromosome-level
1000 genome assembly of a parent species of widely cultivated azaleas. Nat Commun. 11(1), 1–13.

1002 Yang, X., Kang, M., Yang, Y., Xiong, H., Wang, M., Zhang, Z., et al. (2019). A chromosome-level
1003 genome assembly of the Chinese tupelo *Nyssa sinensis*. Sci Data. 6(282), 1–7.

1005 Yang, X., Yue, Y., Li, H., Ding, W., Chen, G., Shi, T., et al. (2018). The chromosome-level quality
1006 genome provides insights into the evolution of the biosynthesis genes for aroma compounds of
1007 *Osmanthus fragrans*. Hortic Res. 5(72), 1–13.

1009 Yang, Z., Hou, Q., Cheng, L., Xu, W., Hong, Y., Li, S., et al. (2017). RNase H1 cooperates with
1010 DNA gyrases to restrict R-loops and maintain genome integrity in *Arabidopsis* chloroplasts. Plant
1011 Cell. 29, 2478–2497.

1013 Ye, C., Hill, C.M., Wu, S., Ruan, J., Ma, Z.S. (2016). DBG2OLC: efficient assembly of large
1014 genomes using long erroneous reads of the third generation sequencing technologies. Sci Rep. 6, 1–9.

1016 Yin, X., Arias-Pérez, A., Kitapci, T.H., Hedgecock, D. (2020). High-Density Linkage Maps Based on
1017 Genotyping-by-Sequencing (GBS) Confirm a Chromosome-Level Genome Assembly and Reveal
1018 Variation in Recombination Rate for the Pacific Oyster *Crassostrea gigas*. G3 - Genes Genom Genet.
1019 10, 4691–4705.

1021 Zerbino, D.R. and Birney, E. (2008). Velvet: algorithms for de novo short read assembly using de
1022 Bruijn graphs. Genome research. 18, 821-829.

1024 Zhang, G.J., Dong, R., Lan, L.N., Li, S.F., Gao, W.J., Niu, H.X. (2020). Nuclear integrants of
1025 organellar DNA contribute to genome structure and evolution in plants. Int J Mol Sci. 21, 707.

1027 Zhang, X., Zhang, S., Zhao, Q., Ming, R., Tang, H. (2019). Assembly of allele-aware, chromosomal-
1028 scale autopolyploid genomes based on Hi-C data. Nat Plants. 5, 833–845.

1030 Zhu, T., Wang, L., You, F.M., Rodriguez, J.C., Deal, K.R., Chen, L., et al. (2019). Sequencing a
1031 *Juglans regia* × *J. microcarpa* hybrid yields high-quality genome assemblies of parental species.
1032 Hortic Res. 6(55), 1–16.
1033

1034

1035 **11. Tables**

1036

1037 **Table 1**. Comparison of BUSCO completeness in Fagaceae genomes available and in the present

1038 study (*Fagus sylvatica* V2).

| Species | Complete genes | Single genes | Duplicated genes | Fragmented genes | Missing genes |
|---|---|---|---|---|---|
| *Fagus sylvatica* V2 | 97.4% | 90.3% | 7.1% | 1.3% | 1.3% |
| *Fagus sylvatica* V1 (Mishra et al., 2018) | 96.6% | 85.6% | 11% | 1.8% | 1.6% |
| *Castanea mollissima* (Wang et al., 2020) | 92.4% | 88.8% | 3.7% | 1.5% | 6.1% |
| *Quercus lobata* v3 (Sork et al., 2016) | 93.5% | 87.6% | 5.9% | 1.0% | 5.5% |

1039

1040

1041

1042 **Table 2**. Distribution of exons in *Fagus sylvatica* in comparison to *Juglans regia* and *Arabidopsis*

1043 *thaliana.*

| Species | Minimum exons / gene | First quartile | Mean exons / gene | Median exons / gene | Third quartile | Maximum exons / gene |
|---|---|---|---|---|---|---|
| *Fagus sylvatica* V2 | 1 | 2 | 4.916 | 4 | 7 | 70 |
| *Juglans regia* (Martínez-García et al., 2016) | 1 | 2 | 5.301 | 4 | 7 | 70 |
| *Arabidopsis thaliana* (GCA_000001735) | 1 | 1 | 5.299 | 4 | 7 | 79 |

1044
1045
1046
1047 **Table 3.** Size of the full-sib families identified from pedigree reconstruction.

| candidate father | size of the full-sib family |
|---|---|
| MSSB | 47 |
| MSSH | 68 |
| SSP01 | 24 |
| SSP02 | 27 |
| SSP03 | 4 |

| | |
|---|---|
| SSP04 | 10 |
| SSP05 | 16 |
| SSP06 | 13 |
| SSP07 | 9 |
| SSP08 | 17 |
| SSP09 | 12 |
| SSP10 | 9 |
| SSP11 | 17 |
| SSP12 | 86 |
| SSP13 | 15 |
| SSP14 | 10 |
| SSP15 | 2 |
| SSP16 | 13 |
| SSP17 | 3 |
| SSP18 | 8 |
| sum | 410 |

1048

1049 **Table 4.** Summary statistics for three Fagus sylvatica unigene sets. The last column gives the number
1050 of homologous proteins (blastX E10-5) against the most complete fagaceae proteome (25,808
1051 proteins) to date, that of Quercus robur (Plomion et al. 2018).
1052

27

| | Technologies | assembler | # contigs in the unigene | Identified oak proteins | # contigs with identified proteins |
|---|---|---|---|---|---|
| Lesur et al. 2015 | Sanger 454 Roche | MIRA | 21,057 | 22,684 | 16,512 |
| Muller et al. 2017 | Illumina | CLCBio | 44,335 | 24,804 | 24,480 |
| This study* | Illumina ONT | Velvet Oases | 34,987 | 24,826 | 22,347 |
| | | | 33,013** (≥200bp) | 24,811 | 21,886 |

1053 *In addition to Illumina and ONT RNAseq, contigs obtained from Lesur et al. 2015 were also
1054 included in the analysis. This first unigene provided a total of 609 transcripts to the new reference
1055 unigene. * *Transcripts longer than 200bp are available online (ENA accession HBVZ01000000).
1056 Smaller contigs are available upon request.

1057

1058 **Table 5.** Characteristics of the combined maternal linkage map in terms of genetic size (cM) and
1059 number of SNP markers for each linkage group (LG).

| LG | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Size (cM) | 279 | 152 | 224 | 137 | 168 | 192 | 146 | 172 | 182 | 171 | 186 | 64 | 140 | 2213 |
| # of SNPs | 37 | 30 | 56 | 36 | 49 | 24 | 24 | 22 | 29 | 15 | 22 | 16 | 8 | 368 |

1060

1061 **Table 6.** Number of SNP markers of a given linkage group (LG) aligned to a specified scaffold
1062 (Bhaga_i) of the *Fagus sylvatica* assembly.

| | Bhaga_1 | Bhaga_2 | Bhaga_3 | Bhaga_4 | Bhaga_5 | Bhaga_6 | Bhaga_7 | Bhaga_8 | Bhaga_9 | Bhaga_10 | Bhaga_11 | Bhaga_12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LG1 | 2 | | | | | 26 | | | | | | |

|      |    |    |    |    |    |    |    |    |    |    |    |    |    |
|------|----|----|----|----|----|----|----|----|----|----|----|----|----|
| LG2  | 1  |    |    |    | 1  |    | 23 |    |    |    |    |    |    |
| LG3  |    | 42 |    |    |    |    |    |    |    |    |    |    |    |
| LG4  |    | 1  | 1  |    | 22 | 1  | 1  |    |    |    | 1  |    |    |
| LG5  |    |    |    |    |    |    |    |    |    |    |    |    | 42 |
| LG6  | 1  |    | 16 | 1  |    |    |    |    |    |    |    |    |    |
| LG7  | 16 |    | 1  |    |    |    |    | 1  |    | 2  |    |    |    |
| LG8  |    |    |    |    | 1  |    |    |    |    | 15 |    |    |    |
| LG9  |    |    |    |    |    | 1  |    |    | 25 |    |    |    |    |
| LG10 | 1  |    |    | 12 |    |    |    |    |    | 1  |    |    |    |
| LG11 |    |    |    |    |    |    |    | 20 |    |    |    |    |    |
| LG12 | 1  |    |    |    | 1  |    | 1  | 10 |    |    |    |    |    |
| LG13 |    | 1  |    |    | 1  | 1  |    |    |    |    | 2  |    |    |

1063

1064

1065
1066    **12. Figure captions**
1067

1068    **Fig. 1**. The more than 300 year-old *Fagus sylvatica* reference individual Bhaga on a cliff over the

1069    Edersee in the Kellerwald Edersee National Park (Germany)


1070    **Fig. 2**. Locations of probable centromeric repeats on the chromosomes presented as red lines and

1071    telomeric locations as blue line on the chromosomes.

1072    **Fig. 3**. Chloroplast genome insertions within 100 kb windows on the chromosomes. Each

1073    chromosome is represented as three rows, the first with insertions more than 100 bp long, the second

1074    row with more than 1 kb and the third with more than 10 kb.


1075    **Fig. 4**. Mitochondrion genome insertions within 100 kb windows on the chromosomes. Each

1076    chromosome is represented as three rows, the first with insertions more than 100 bp long, the second

1077    row with more than 1 kb and the third with more than 10 kb.


1078    **Fig. 5**. Repeat regions, coding regions, and regions coding for genes present within 100 kb windows

1079    on the chromosomes.


1080    **Fig. 6**. Homozygous and heterozygous SNPs in *Fagus sylvatica* present within 100 kb windows on

1081    the chromosomes.


1082    **Fig. 7**. Distribution of homozygous and heterozygous SNPS in non-overlapping 100 kb windows.


1083    **Fig. 8**. Example of the high collinearity between homologous maternal (MSSB) linkage group #4

1084    obtained from the analysis of three sets of offspring: xMSSH and xSSP12 correspond to the two

1085    largest full-sib families and x182 correspond to the cosegregation analysis of their mapped markers

1086    in the 182 half-sibs.


1087


1088


1089


1090    **Supplementary Files**


1091    **Supplementary file 1**. Details of annotated repeat elements in Fagus sylvatica.

1092    **Supplementary file 2**. Venn diagram showing shared PLAZA proteins of *Arabidopsis thaliana*

1093    (27615), *Eucalyptus grandis* (36331), and *Vitis vinifera* (26346) with those of *Fagus sylvatica*

1094    (28326).

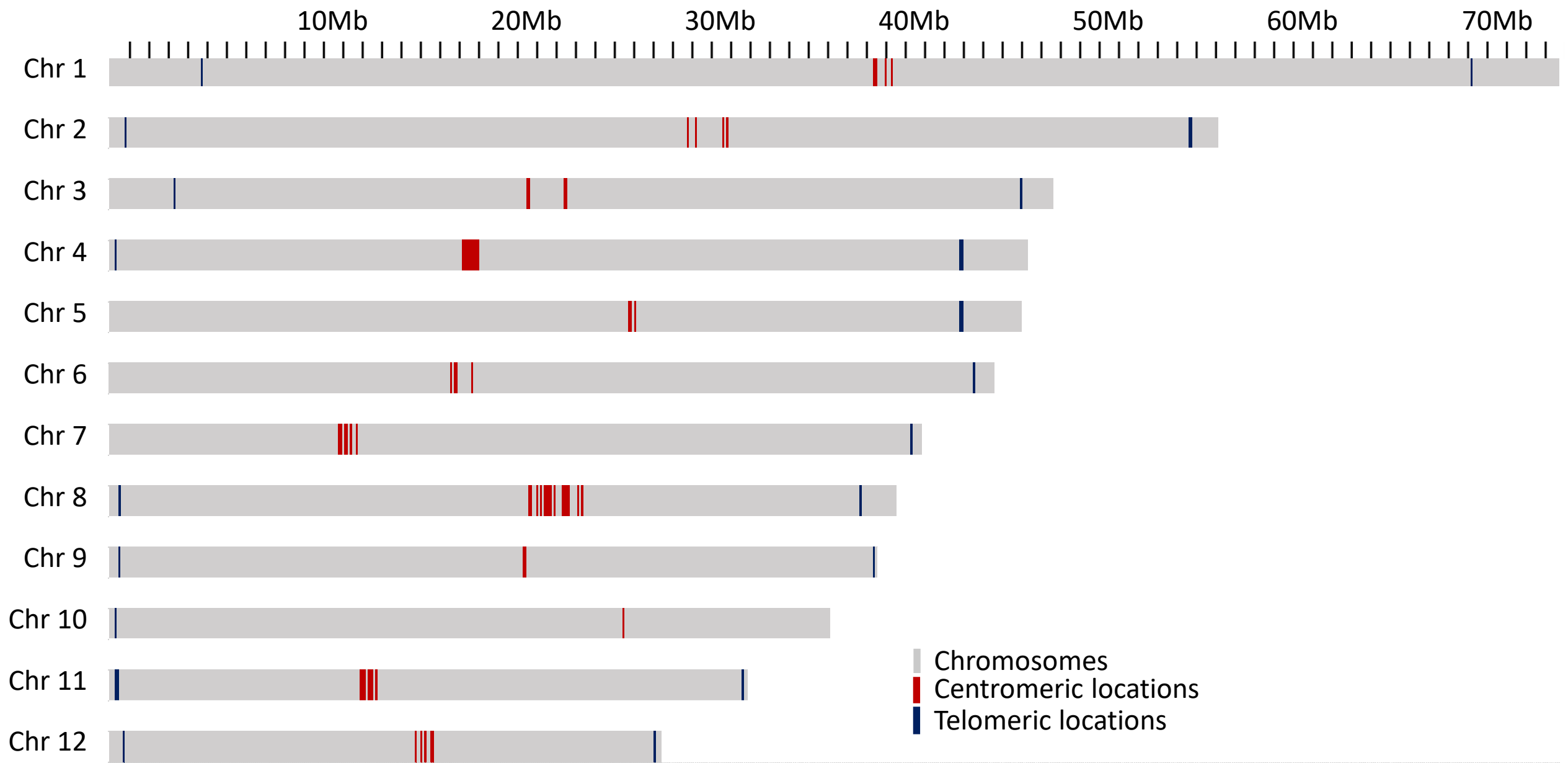1095    **Supplementary file 3**. Centromeric feature annotation.

1096    **Supplementary file 4**. Details of the conservation of organelle DNA insertions in the nuclear

1097    genome.

1098    **Supplementary file 5.** Multiplexed SNP assay. SNP_ID (a) following Ouayjan et al. (2018). SNPs

1099    discarded from the analyses are highlighted in yellow; Locus_Name_pos_SNP (b) corresponds to the

1100    locus name given by Lalagüe et al. (2014); seq SNP corresponds to sequences of the SNP flanking

1101    regions. The targeted SNP is indicated in brackets [ / ].

1102    **Supplementary file 6.** Genotyping data. List of 4127 SNPs and their associated linkage groups based

1103    on R_qtl (second raw) and JoinMap (third raw) analyses.

1104

31

Row 1: ▮ chloroplast insertion for >= 100 bases with 95% identity.

Row 2: ▮ chloroplast insertion for >= 1000 bases with 95% identity.

Row 3: ▮ chloroplast insertion for >= 10 Kb bases with 95% identity.

Row 1: ▮ mitochondrial insertion for >= 100 bases with 95% identity.

Row 2: ▮ mitochondrial insertion for >= 1000 bases with 95% identity.

Row 3: ▮ mitochondrial insertion for >= 10 Kb bases with 95% identity.

Repeat region per 100 Kb (2694 - 99857)

Coding region per 100Kb (0 - 49668)

Coding region of duplicated genes per 100Kb (0 – 47227)

Heterozygous SNPs per 100Kb (0 - 1294)

Heterozygous genic SNPs per 100Kb (0 - 331)

Homozygous SNPs per 100Kb (0 - 1532)

Homozygous genic SNPs per 100Kb (0 - 310)

**LG4_x182_Reg**

**LG_xMSSH**

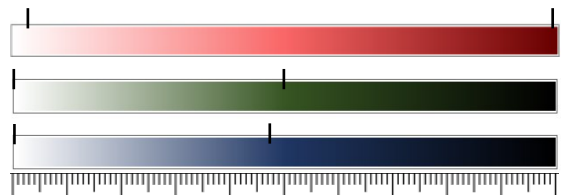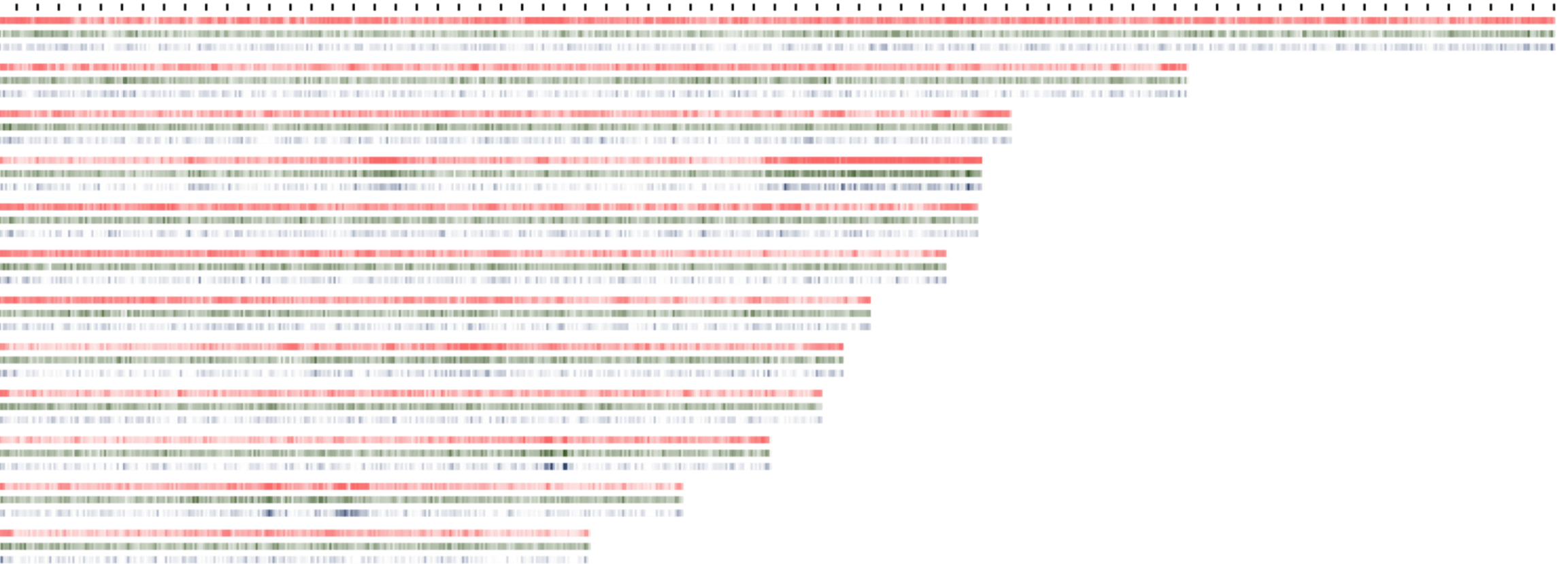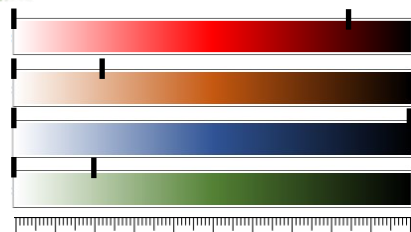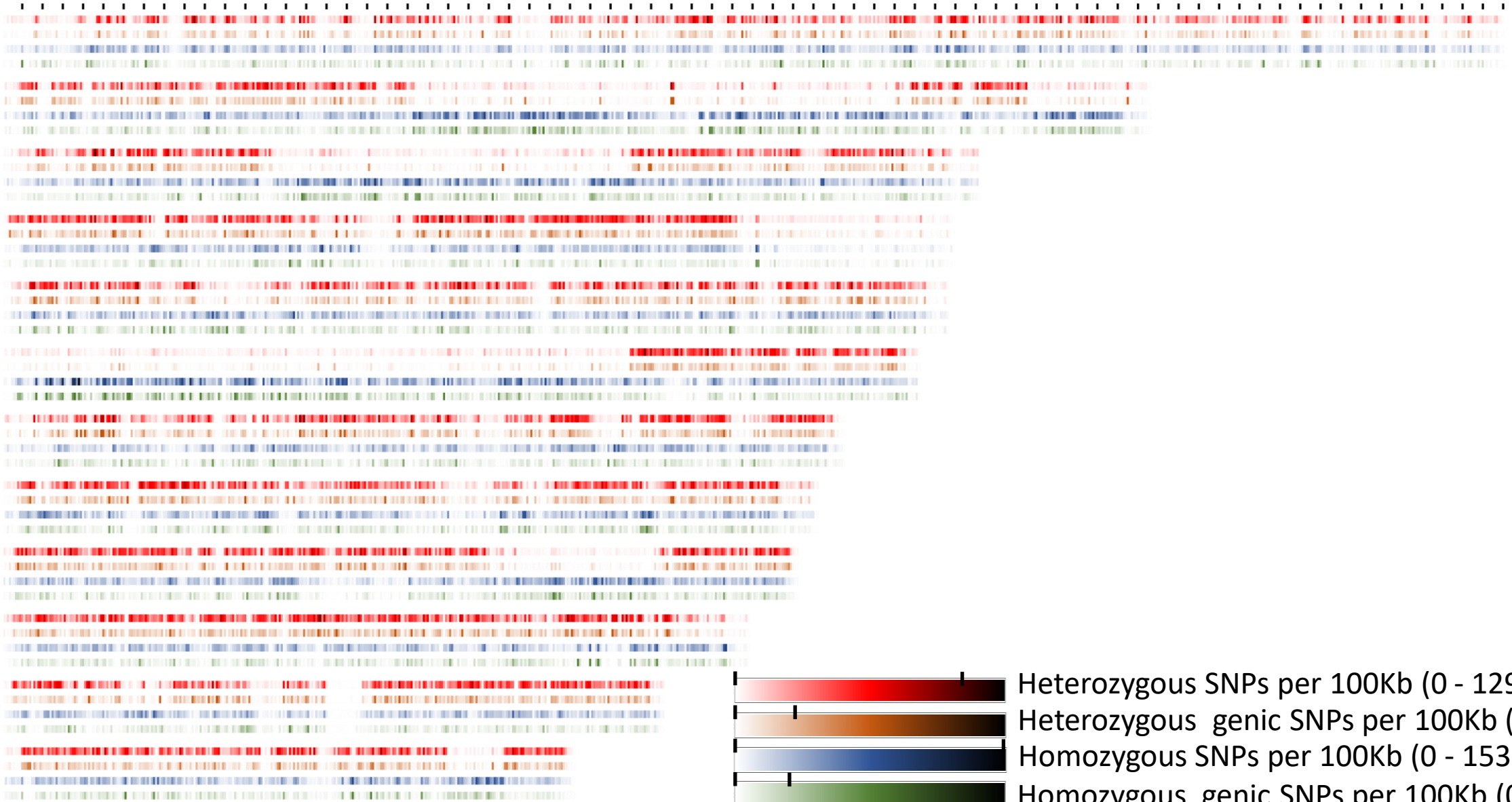| | |
|---|---|
| 0,0 | snp_B4428 |
| 1,2 | snp_B6053 |
| 5,2 | snp_B5603 |
| 5,4 | snp_B0038 |
| 5,5 | snp_B2562  snp_B0482 |
| 5,6 | snp_B5240 |
| 5,8 | snp_B1388 |
| 6,3 | snp_B1833 |
| 6,5 | snp_B0759 |
| 9,4 | snp_B3803 |
| 9,8 | snp_B1616 |
| 12,2 | snp_B3397 |
| 12,4 | snp_B4570 |
| 13,4 | snp_B0893 |
| 18,2 | snp_B3108 |
| 18,8 | snp_B0485 |
| 20,8 | snp_B5713 |
| 26,4 | snp_B1498 |
| 28,6 | snp_B2665 |

**LG4_x182**

| | |
|---|---|
| 0,0 | snp_B3397 |
| 45,3 | snp_B1498 |
| 53,0 | snp_B1604 |
| 57,6 | snp_B1882 |
| 64,2 | snp_B1627 |
| 69,4 | snp_B4947 |
| 70,9 | snp_B0759 |
| 73,4 | snp_B3852 |
| 76,7 | snp_B4334 |
| 77,4 | snp_B4156 |
| 78,5 | snp_B2915 |
| 79,0 | snp_B3333 |
| 79,8 | snp_B4174 |
| 81,9 | snp_B1705 |
| 86,4 | snp_B2522 |
| 88,4 | snp_B3430 |
| 90,9 | snp_B3803 |
| 92,7 | snp_B4371 |
| 104,1 | snp_B5713 |
| 119,3 | snp_B1051 |
| 131,7 | snp_B1050 |
| 146,3 | snp_B1747 |
| 148,1 | snp_B4263 |
| 157,4 | snp_B2267 |
| 160,8 | snp_B5184 |
| 165,8 | snp_B1258 |

**LG_xSSP12**

| | |
|---|---|
| 0,0 | snp_B6333 |
| 4,2 | snp_B1498 |
| 11,0 | snp_B1882 |
| 13,0 | snp_B1604 |
| 18,5 | snp_B2976 |
| 25,9 | snp_B5240 |
| 27,7 | snp_B3852  snp_B4156 |
| | snp_B4947  snp_B1627 |
| 28,7 | snp_B893 |
| 29,9 | snp_B2915  snp_B3333 |
| | snp_B4334  snp_B4174 |
| 32,4 | snp_B1705 |
| 36,4 | snp_B2517 |
| 38,1 | snp_B246 |
| 44,5 | snp_B3803 |
| 45,4 | snp_B2522 |
| 46,0 | snp_B3430 |
| 46,3 | snp_B1616  snp_B5603 |
| 48,5 | snp_B3108 |
| 49,5 | snp_B1177 |
| 50,3 | snp_B4371 |
| 54,8 | snp_B4158 |
| 56,9 | snp_B5713 |
| 61,9 | snp_B1903 |
| 63,4 | snp_B3390 |
| 67,7 | snp_B3891 |
| 70,5 | snp_B1747 |
| | snp_B5184  snp_B4263 |
| 73,5 | snp_B2267 |
| 74,1 | snp_B5035 |
| 74,3 | snp_B1050 |
| 75,5 | snp_B1258 |
| 76,7 | snp_B4694 |
| 77,9 | snp_B3568 |
| 81,8 | snp_B1051 |
| 83,2 | snp_B2573 |
| 86,2 | snp_B3154 |
| 88,1 | snp_B4361 |